# D I S S E R T A T I O N

# Numerical Analysis of the 1D Satellite Beam Equation

ausgeführt zum Zwecke der Erlangung des akademischen Grades eines Doktors der technischen Wissenschaften unter der Leitung von

Prof. Dr.rer.nat. Dipl.-Math. Dipl.-Ing. Carsten Carstensen
Institut für Mathematik
Humboldt-Universität zu Berlin

und

Ao. Univ. Prof. Dr.techn. Dipl.-Math. Dirk Praetorius
Institut für Analysis und Scientific Computing
Technische Universität Wien

eingereicht an der Technischen Universität Wien
Fakultät für Mathematik und Geoinformation

von

## Dipl.-Math. Dipl.-Ing. Jelena Bojanić
Matr. Nummer 0125467
Hernalser Hauptstrasse 104/8
A-1170 Wien

Wien am 14. Juni 2005

# Kurzfassung

Das Verständnis der Dynamik verkabelter Satellitensysteme (engl. *tethered satellite systems*) ist von großer Bedeutung für die Raumfahrt. Besonders interessant sind Systeme von zwei Satelliten, bei denen ein kleiner Subsatellit von einem massiven Space-Shuttle an einem bis zu hundert Kilometer langen visko-elastischen Kabel ausgesetzt wird. Um die Bewegungsgleichungen solcher Systeme genauer zu untersuchen, werden diese durch mechanische Modelle ähnlicher Struktur beschrieben. Darüber hinaus lassen sich numerische und analytische Aspekte der Modellierung mittels des Fadenpendel-Modell genauer untersuchen.

In der vorliegenden Arbeit wird das Fadenpendel-Modell in vertikaler Lage diskutiert, wobei die Herleitung der entsprechenden Bewegungsgleichung auf die partiellen Differentialgleichung

$$\ddot{y} = \Delta y + \varepsilon \Delta \dot{y} + f, \tag{1}$$

mit $\varepsilon \geq 0$ führt. Diese partielle Differentialgleichung wird in der Literatur als stark gedämpfte Wellengleichung (engl. *strongly damped wave equation*) bezeichnet. Sie ist für Dämpfungsparameter $\varepsilon > 0$ von hyperbolisch-parabolischen Charakter.

In der Literatur ist bisher keine numerische Analysis für den Fall $\varepsilon > 0$ zu finden, vorhandene Resultate beschränken sich auf die Wellengleichung, d.h. $\varepsilon = 0$. Die vorliegende Arbeit schließt diese Lücke. Zur numerischen Lösung wird (1) dabei als System erster Ordnung in der Variablen $u = (y, \dot{y})$ formuliert,

$$\dot{u} + \mathcal{A}u = F. \tag{2}$$

Aufgrund der analytischen Eigenschaften des Operators $\mathcal{A}$ lässt sich die eindeutige Lösbarkeit von (2) in einem Hilbert-Raum $\mathcal{H}$ beweisen. Die Diskretisierung von (2) erfolgt mit Galerkin-Verfahren in Ort und Zeit für beide Komponenten von $u$. Die Orts-Diskretisierung basiert auf stückweise affinen bzw. stückweise kubischen Ansatzfunktionen. Für die Zeitdiskretisierung verwenden wir stetige und unstetige Galerkin-Verfahren $cG(1)$ bzw. $dG(q)$ mit $q = 0, 1$.

In Rahmen dieser Arbeit werden die vorgeschlagenen Diskretisierungsverfahren im Zeit-Orts-Raum analytisch untersucht. In Kapitel 2 beweisen wir zunächst die eindeutige Existenz diskreter Lösungen für alle vorgeschlagenen Diskretisierungsverfahren. Der Schwerpunkt liegt hiernach auf der a priori und a posteriori Analysis des Diskretisierungsfehlers in der Energienorm.

Unter einer a priori Abschätzung des Fehlers versteht man dabei im wesentlichen eine Abschätzung, bei der die rechte Seite von der exakten Lösung $u$, aber nicht von der diskreten Lösung $U$ abhängt, sodass man (unter Regularitätsannahmen an $u$) die Konvergenzordnung des Verfahrens erhält. Im Gegensatz dazu hängt die rechte Seite einer a posteriori Abschätzung von $U$ aber nicht von $u$ ab, d.h. die Fehlerschranke ist explizit berechenbar. A posteriori Fehlerabschätzungen werden daher zur effizienten Netzgenerierung in adaptiven Algorithmen verwendet.

Unsere Fehleranalysis basiert im wesentlichen auf drei verschiedenen Beweistechniken: der Energie-Methode in Kapitel 3, der dualen Methode in Kapitel 4 und der zielorientierten Methode in Kapitel 5. Die Vorteile und Schwierigkeiten der drei Techniken werden herausgearbeitet und im Einzelnen diskutiert. Die analytischen Resultate verallgemeinern Arbeiten von ERIKSSON-JOHNSON für den Fall $\varepsilon = 0$ auf den Fall $\varepsilon \geq 0$. Teilweise werden dabei suboptimale a priori Abschätzungen verbessert, und wir erhalten optimale Konvergenzordnungen in den Netzschrittweiten $h$ für die Ortsdiskretisierung und $k$ für die Zeitdiskretisierung. Abschließende numerische Experimente in Kapitel 6 belegen die bewiesene Analysis.

# Contents

# List of Tables

# List of Figures

# Introduction

## Motivation

During the last decade a new space technology, called tethered satellite systems (TSS), which are a system of two or more satellites (rigid bodies) in orbit around the Earth connected by thin flexible visco-elastic cables, has been of great interest. For instance, there were space flights where one of the satellites was a space-shuttle and the other one a satellite of significantly less mass. The length of the connecting tether can be up to 100 km, see KRUPA et al. [48]. An overview concerning the practically important applications of TSS, can be found in [34, 76, 77].

The mathematical model of a TSS involves continuous and discrete bodies and hence the equations of motion are a set of coupled nonlinear partial and ordinary differential equations. The nonlinearities follow from finite geometry of the large displacement. The equations are stiff in the sense that motions are present, i.e. axial and transversal motions, which evolve on different time scales. An efficient numerical simulation of such systems and a reliable approximation requires an appropriate implicit time integration. The spatial discretisation is successfully conducted by the finite element method [63, 70, 74]. In order to express the property that the equations are stiff already in their formulation, instead of the usual displacement coordinates of the deformation of the tether, the orientation of the tangent vector to the deformed tether as slow variable and the axial strain of the tether as fast variable are used. Then the resulting string equations separate into a fast and a slow part and are named after Minakov. It is referred to [12, 36, 58] for an integration of the resulting tether equations without any systematic analysis of efficiency. Especially problems in the numerical simulation arise if, for large amplitude motions, the tether force becomes zero and hence the tether is not strained but slack POTH et al. [58]. Hence one important question is how is it possible to achieve high accuracy in simulations of such motions which include the case of slack string configuration.

For a mathematical investigation all essential properties and difficulties of the equations of a TSS are already given by the model of a string pendulum, [47, 62]. There, one point mass is connected by a massive visco-elastic string and with this model in KUHN et al. [49] the numerical calculations in the displacement formulation are validated. If the restriction to the vertical position of the string is made, the equation of motion for the string pendulum in the space-fixed non-rotating frame reduces to the strongly damped scalar wave equation [49]. This equation is of the hyperbolic-parabolic character due to the presence of the damping parameter $\varepsilon \geq 0$. This thesis emphasises the analysis of the strongly damped wave equation in which $\varepsilon = 0$, i.e., the analysis of the wave equation, is included.

# The strongly damped wave equation

Given some arbitrary external force $f$, the strongly damped wave equation reads

$$\frac{\partial^2 y}{\partial t^2}(t,x) - \frac{\partial^2 y}{\partial x^2}(t,x) - \varepsilon \frac{\partial^3 y}{\partial x^2 \partial t}(t,x) = f(t,x) \quad \text{in} \quad [0,T] \times \Omega. \tag{3}$$

A more general overview on the spatial boundary and initial conditions in time will be given in Section 1.3 with Dirichlet conditions in $x = 0,1$ or mixed Dirichlet-Neumann boundary conditions.

The initial condition in time describes the initial length and the initial velocity.

Apart from the choice of the boundary conditions, the initial problem (3) allows the conversion into an equivalent vector formulation, typical for the hyperbolic-type problem, see HUGHES-HULBERT [38]. Indeed, for $u = (y, \dot{y})$, the problem (3) is equivalent to

$$u + \mathcal{A}u = F. \tag{4}$$

Therein, $\mathcal{A}$ is some maximal monotone operator. This vector setting takes place in the Hilbert space $\mathcal{H} = H_D^1(\Omega) \times L^2(\Omega)$, cf. NEVES [54]. The problem (4) can be also defined with respect to $L^2(\Omega) \times L^2(\Omega)$ scalar product, cf. BANGERTH [9]. Here, however we rely on the first formulation with $\mathcal{H}$. This allows us also to define and analyse the energy behaviour of the continuous as well as of the discrete solution. Namely, the energy, given as a sum of the potential and kinetic energy of the solution is defined through the $\mathcal{H}$ norm, namely,

$$\mathcal{E}(u) = \frac{1}{2}\|u\|_{\mathcal{H}}^2.$$

Furthermore, use of $\mathcal{H}$ serves as a better basis for the analysis of the solvability, i.e. the uniqueness and the existence of the solution. We prove that according to the properties of the operator $\mathcal{A}$, there exists a unique solution of the problem (4).

# Discretisation in space and time

Within this work we present a study of the certain approximation techniques, applied to the vector form of the strongly damped wave equation (4) in a way that the analysis of the error becomes much more efficient and an optimal or nearly optimal a posteriori and a priori error estimates can be obtained.

In particular, we consider the numerical approximation of both components of the vector solution separately but using the same time-space discrete ansatz. Time approximation is done by applying the discontinuous Galerkin scheme $dG(q)$ of order $q = 0,1$ as well as the continuous Galerkin $cG(1)$ method. We also apply the semi-discretisation in space, so-called Method of Lines. The use of finite elements to discretise the temporal as well as the spatial domain was first proposed in [4, 31, 56].

Most time-dependent problems use the semidiscrete methods for designing the algorithms for the computation of the approximative solution, cf. SÜLI-WILKINS [67], for the analysis of the damped wave equation and BABUŠKA et al. [6] where the class of evolution problems has been solved. In these references, the convergence results concerning the a priori and a posteriori error estimates obtained with aid of the dual technique are provided. The method of lines is favourable when approximating the solutions which are smooth. It is also easy to perform

since the solvability of the discrete solution relies onto the solvability of the system of the ordinary differential equations. It retains the energy as the $cG(1)$ method, when comparing the continuous and discrete solution. However, their performance is less satisfactory when solving problems with discontinuities or sharp gradients in the solution. One other disadvantage of the method of lines is that the corresponding space-time discretisation is structured, cf. HUGHES-HULBERT [38]. This disables the adaptive mesh refinement.

The method which seems to cover these difficulties is the discontinuous Galerkin method. First introduced for 1D hyperbolic problems, cf. [45, 51, 59], this method yields an A-stable, higher-order accurate method, cf. [51, 44]. Furthermore, this time discretisation framework seems to be conducive to the establishment of the rigorous convergence proofs and error estimates for other evolution equations, cf. [19, 25, 37, 43, 66, 69]. For the strongly damped wave equation as well as the wave equation, the analysis of the convergence of the error is given in e.g. JOHNSON [42] and LARSSON-THOMÉE-WAHLBIN [50]. In JOHNSON [42], the wave equation with Dirichlet boundary conditions is discussed. The methods of proofs for the establishment of the a posteriori and a priori error bounds rely on the dual method technique. If the mesh does not change from one time slab to another, then the convergence result can be found very satisfactory. In LARSSON-THOMÉE-WAHLBIN [50], the strongly damped wave equation has been discussed whereby the derivation of the error bounds for several discretisation methods has been conducted with aid of the semi-group theory.

Another approach towards space-time discretisation is the $cG(1)$ method in which the unknown quantities are assumed to be globally continuous, piecewise affine in time. This method allows the use of unstructured space-time grid, but the problems may be effective discretisation of the solution which are less smooth in time. See [5, 9, 8, 10, 11, 30, 29] for some examples of continuous time space finite element discretisation. In BANGERTH [9], the acoustic wave equation has been discussed, whereby the a posteriori error analysis has been conducted with aid of the goal-oriented arguments. The analysis presented by FRENCH-JENSEN [29] refers to the strongly damped wave equation for which the long time behaviour has been analysed. In FRENCH-PETERSON [30], the a priori error estimates have been proved for the wave equation by using the energy arguments. The a priori error estimates in the negative as well as in the $\mathcal{H}$ norm have been proved for the wave equation in BALES-LASIECKA [8].

For the hyperbolic problems in general, the main difficulty occurs in the interpretation of the time derivative in case of non-smooth problems. For first order hyperbolic problems, this can be solved by modifying some finite element methods, in order to improve the convergence, we refer to JOHNSON [43]. For the wave equation as well as for the strongly damped wave equation, some authors propose so-called shock-capturing artificial viscosity, c.f. [42, 71], but this wont be discussed within this work.

The approximation of the spatial variable follows by means of the conforming $\mathcal{P}_1$ elements (linear splines) as well as of the Hermite cubic splines ($\mathcal{C}^1$ elements), respectively. The $\mathcal{C}^1$ elements provide a better basis for construction of the a priori and a posteriori error bounds due to the continuity in the first derivative when the discrete function is concerned. This is obvious when the derivation of the a posteriori error estimates by using the energy techniques is discussed.

# Outline of the thesis and some results

An outline of the thesis reads as follows:

- In Chapter 1, we first introduce the mechanical model of the string pendulum. In its vertical position this model can be described by the strongly damped wave equation where the damping parameter $\varepsilon \geq 0$ characterises the visco-elastic property of the string. The subsequent Section 1.2 provides some backup facilities needed for the understanding of the mathematical theory for observed model. This assumes the introduction of the Sobolev spaces, distributions etc. In Section 1.3, we analyse the strongly damped wave equation from the mathematical point of view. Furthermore, we introduce the vector formulation in the Hilbert space $\mathcal{H}$ which enable us to define an energy and energy scalar product which will be frequently used in the error analysis. Finally, we recapitulate the analytical theory by proving the existence and the uniqueness of the solution of the strongly damped wave equation. The stability of the continuous solution is also showed within this section.

- In Chapter 2 we introduce the discrete model for each time discretisation ansatz in particular and derive the stability estimates for the same. We also provide the proof of the unique solvability for each time-space discrete ansatz.

- In Chapter 3 we first introduce the necessary techniques needed for the error analysis in general. This implies the definition of the interpolation and projection operators and their approximation properties, see Section 3.1. In the following we analyse a derivation of the a priori and a posteriori error bounds for the full discrete problems, see Section 3.2 and 3.3, respectively. Techniques used here are new. In the a priori error analysis, see Table 1, the estimates for $\|e\|_{L^\infty(\mathcal{H})}$ when $cG(1) \otimes \mathcal{C}^1$ and $MoL \otimes \mathcal{C}^1$ are optimal when

| $\otimes$ | $\mathcal{P}_1$ | $\mathcal{C}^1$ |
|---|---|---|
| $dG(0)$ | $\mathcal{O}(h+k^{-1/2}h+k)$ | $\mathcal{O}(h^3+k^{-1/2}h^3+k)$ |
| $cG(1)$ | – | $\mathcal{O}(h^3+k^2),\ \varepsilon=0$ <br> $h-$quasi uniform spatial mesh |
| $MoL$ | – | $\mathcal{O}(h^3),\ \varepsilon=0$ or $(DD)$ <br> $h-$quasi uniform spatial mesh |

Table 1: Proven a priori error estimates; energy method.

restricted to the wave equation, i.e. of order $\mathcal{O}(h^3 + k^2)$ and $\mathcal{O}(h^3)$, respectively. For $dG(0)$ method we prove a convergence order $\mathcal{O}(h^p + k^{-1/2}h^p + k)$ where $p=1$ for $\mathcal{P}^1$ and $p=3$ for $\mathcal{C}^1$ elements. These estimates are nearly optimal, due to the presence of factor $\mathcal{O}(k^{-1/2}h^p)$.

In case of the residual-based a posteriori error analysis, see Table 2 for $dG(0) \otimes \mathcal{C}^1$ and $MoL \otimes \mathcal{C}^1$ time-space ansatz, the proven error bound retain the optimal convergence of the exact error, i.e. $\mathcal{O}(h^3 + k)$ and $\mathcal{O}(h^3)$, respectively. For $dG(0)$ we additionally assume that the spatial meshes are hierarchical in time. Also note that the a posteriori error analysis for $dG(0)$ in time, relies on the use of the affine approximation $\widetilde{U}$, bilinear form $\widetilde{\mathcal{B}}$ and residual $\widetilde{Res}$. All a posteriori error bounds proved within this Chapter are

| $\otimes$ | $\mathcal{C}^1$ |
|---|---|
| **dG(0)** | $\mathcal{O}(h^3+k)$ $\scriptstyle \mathcal{S}^{j-1}\subseteq\mathcal{S}^j$ <br> $\mathcal{O}(h^3+k^{-1}h^3+k)$ $\scriptstyle \text{otherwise}$ |
| **dG(1)** | $\mathcal{O}(h^3+k)$ $\scriptstyle \mathcal{S}^{j-1}\subseteq\mathcal{S}^j$ <br> $\mathcal{O}(h^3+k^{-1}h^3+k)$ $\scriptstyle \text{otherwise}$ |
| **cG(1)** | $\mathcal{O}(h^3+k)$ |
| **MoL** | $\mathcal{O}(h^3)$ |

Table 2: Proven a posteriori error estimates; energy method.

computable, i.e. the constants are known. This analysis can be also easily extended to the 2D case.

- The previous chapter motivates the further analysis on the error by use of the duality technique, cf. Chapter 4. We first introduce the adjoint problem and derive the strong stability estimates, cf. Section 4.1. Then we analyse a derivation of the a priori and a posteriori error bounds for each time-space ansatz in particular, see Section 4.2 and 4.3, respectively. The analysis in case of $dG(q), q=0,1$ method relies on the arguments used in JOHNSON[42], where the $dG(1)\otimes\mathcal{P}_1$ finite element discretisation of the $2D$ wave equation has been studied. Here however, we proved an error bounds when $dG(0)$ and $cG(1)$ method in time additionally.
The proved a a priori error bounds for the strongly damped wave equation, see Table 3,

| | $\varepsilon \geq 0$ | $\varepsilon = 0$ |
|---|---|---|
| **dG(0)** | $\mathcal{O}(h^p+k)$ $\scriptstyle \mathcal{S}^{j-1}=\mathcal{S}^j$ <br> $\mathcal{O}(h^p+k^{-1/2}h^{p+1}+k)$ $\scriptstyle \text{otherwise}$ | |
| **dG(1)** | $\mathcal{O}(h^p+k^2)$ $\scriptstyle \mathcal{S}^{j-1}=\mathcal{S}^j$ <br> $\mathcal{O}(h^p+k^{-1/2}h^p+k^2)$ $\scriptstyle \text{otherwise}$ | $\mathcal{O}(h^p+k^3)$ $\scriptstyle \mathcal{S}^{j-1}=\mathcal{S}^j$ <br> $\mathcal{O}(h^p+k^{-1/2}h^p+k^3)$ $\scriptstyle \text{otherwise}$ |
| **cG(1)** | $\mathcal{O}(h^p+k)$ | $\mathcal{O}(h^p+k^2)$ |

Table 3: Proven a priori error estimates; dual method.

are of the optimal order when $dG(0)$ method in time and $\mathcal{P}^1$ or $\mathcal{C}^1$ discretization in space. Namely, for $\mathcal{P}^1$ $(p=1)$ and $\mathcal{C}^1$ $(p=3)$ we have $\|e^{N-}\|_{\mathcal{H}}=\mathcal{O}(h^p+k)$. This assumes

also that the spatial mesh does not change in time. When wave equation is discretized by $dG(1)$ resp. $cG(1)$ method we proved the optimal convergence order $\mathcal{O}(h^p+k^3)$ and $\mathcal{O}(h^p+k^2)$, respectively. Our estimates when $dG(1)$ in time generalize and improve those of JOHNSON[42] of order $\mathcal{O}(h^{1/2})$ when $k=h^3$, when restricted to the 1D case. The result for $cG(1)$ in time is new.

As far as the a posteriori error bound is concerned, see Table 4, we analyse a derivation

| $\otimes$ <br> if $s=1$ then $\varepsilon=0$ or (DD) | | $\mathcal{P}^1$ | $\mathcal{C}^1$ |
|---|---|---|---|
| **$dG(0)$** | $s=0$ | – | $\mathcal{O}(h^3+k)$, (DD*), $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, <br> $\mathcal{O}(h^3+k^{-1}h^3+k)$ otherwise |
| | $s=1$ | $\mathcal{O}(h+k)$ $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, <br> $\mathcal{O}(h+k^{-1}h^2+k)$ otherwise | $\mathcal{O}(h^4+k)$, $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, <br> $\mathcal{O}(h^4+k^{-1}h^4+k)$ otherwise |
| **$dG(1)$** | $s=0$ | – | $\mathcal{O}(h^3+k^3)$, $\varepsilon=0$, (DD*), $\mathcal{S}^{j-1}=\mathcal{S}^j$ <br> $\mathcal{O}(h^3+k^{-1}h^3+k^3)$ otherwise |
| | $s=1$ | $\mathcal{O}(h+k^3)$, $\varepsilon=0$, $\mathcal{S}^{j-1}=\mathcal{S}^j$, <br> $\mathcal{O}(h+k^{-1}h+k^3)$ otherwise | $\mathcal{O}(h^4+k^3)$, (DD*), $\mathcal{S}^{j-1}=\mathcal{S}^j$ <br> $\mathcal{O}(h^4+k^{-1}h^4+k^3)$ otherwise |
| **$cG(1)$** | $s=0$ | – | $\mathcal{O}(h^{s+3}+k^2)$, (DD*) |
| | $s=1$ | $\mathcal{O}(h+k^2)$ | |

Table 4: Proven a posteriori error estimates; dual method.

of the a posteriori error bounds for $\|e^{N-}\|_{\mathcal{H}}$ ($s=0$) and $\|e^{N-}\|_{\widehat{\mathcal{H}}}$ ($s=1$). The bounds of optimal order are proved for all three time discretisation methods and $\mathcal{C}^1$ functions in space with the restriction to the (DD) boundary conditions ($\Gamma(\Omega)=\Gamma_D(\Omega)$). For $dG(1)$ method in time, the optimal order $\mathcal{O}(h^3+k^3)$ is proved only for the wave equation, whereby in case of the $dG(0)$ and $cG(1)$ method in time, we proved the optimal order for the case $\varepsilon>0$ also. Case $\mathcal{P}^1$ elements in space does not allow us to obtain an optimal error estimates since in $1D$ we can not make necessary requirements on the mesh as in 2D, cf. [42]. There, the error bound of order $\mathcal{O}(h+k^3)$ was proved for the $dG(1)\otimes\mathcal{P}^1$ discrete model of the 2D wave equation with the Dirichlet boundary conditions.

- Chapter 5 is devoted to the goal oriented a posteriori error analysis. We first introduce the general idea of the goal-oriented method and describe the possible ansatz when the strongly damped wave equation is considered. This implies the definition of the target functional and its linearised version, as well as the proposed approximation method when dual solution is considered. In the following, we deduce the a posteriori error estimates for $dG(q), q=0,1$ and $cG(1)$ ansatz in time combined with $\mathcal{P}_1$ or $\mathcal{C}^1$ ansatz in space.

This theory is generalisation of the one presented in [9, 10, 11], where the $cG(1)\otimes\mathcal{P}^1$ discretisation of the acoustic wave equation is treated. However, a detailed numerical discussion concerning the strongly damped wave equation is not included here.

- Some numerical experiments are discussed in Chapter 6. Here we also provide a full algorithm for the computation of the finite element solution using the Galerkin methods in time and space, cf. Section 6.1. Section 6.4 is more theoretically inclined. We also discuss some possible ways of solving the equation by relying on some of the adaptivity techniques. As a basis we take the a posteriori error estimates obtained by using the energy techniques where adaptivity applies only to the temporal mesh.

The presented study due to simplicity of the observed mathematical model, denotes the platform for the future investigations.

# Acknowledgments

I would like to express my deep gratitude to all people who have supported me throughout my research. First of all, I offer my sincerest gratitude to my supervisor, Carsten Carstensen. He introduced me to research while giving me the freedom and the possibility to pursue my own ideas. Many thanks go to Dirk Praetorius, my second advisor, for his patience and helpful comments while reviewing the thesis. His strict and extensive comments and many discussions and the interaction with him had a significant impact on the final form and quality of this thesis. I also want to thank Prof. Hans Troger from the Mechanical Department at Vienna University of Technology who together with Prof. Carstensen initialised the topic of my research. His valuable comments and suggestions helped me to understand the mechanical aspect of the problem. Special thanks goes to all my colleagues from the research group in Berlin as well as to some former members, in particular my dear friend Jan Bolte.

I would also like to all members from the Institute for Analysis and Scientific Computing at Vienna University of Technology for helping me to finalise my thesis in Vienna. Among them I especially want to thank Mrs. Kovalj, Mrs. Frank and Mrs. Schweigler for their efforts on solving my administrative problems and Dieter Kvasnicka for his friendly advises.

I am greatly indebted to my beloved husband Nikola for his patience and support during the time of writing this dissertation. I would also like to express my gratitude to my dear parents and my brother who encouraged and supported me all over my life.

# Chapter 1

# Satellite Beam Equation

We consider the sting pendulum as a simple version of the mechanical model of TSSs, which is believed to convey all the difficulties in the mathematical and numerical treatment. We derive the equations of motion for the string pendulum in the vertical position where the length of the string is set to some constant value. This results in a partial differential equation which is often referred to as the strongly damped wave equation.
The second part of the chapter deals with the mathematical aspect of the strongly damped wave equation, cf. Section 1.3. Here we introduce some necessary definitions and formulations and derive the main result which provides the uniqueness and existence of the continuous solution.

## 1.1   Physical model

The string pendulum, see Figure 1.1, consists of a visco-elastic massive string with an end mass $m_S$, modeled as a mass point, moving in the constant gravitational field with $\boldsymbol{g}$ being the gravitational acceleration. Notice that the deformations of the string pendulum, described by



Figure 1.1: Mechanical model of the string pendulum.

the position vector $\mathbf{r}(\bar{s}, \tau)$ are given with respect to the inertial frame $(\boldsymbol{e}_x, \boldsymbol{e}_y)$.

Here $\tau$ denotes time and $\bar{s}$ the unstrained arclength which is a material coordinate along the string and $0 \leq \bar{s} \leq \ell$ for $\ell$ the total length of the longitudinally unstrained string. The unit tangential vector is denoted by $\boldsymbol{t}$ whereby $\epsilon$ stands for the strain. If $d\bar{s}$ and $ds$ are the unstrained and strained element lengths, then

$$\epsilon := \frac{ds - d\bar{s}}{d\bar{s}}.$$

The large amplitude motion of the string pendulum system are described by a field equation

$$\bar{\mu}\frac{\partial^2 \boldsymbol{r}}{\partial \tau^2} = \frac{\partial \boldsymbol{n}}{\partial \bar{s}} - \bar{\mu}\boldsymbol{g}, \tag{1.1}$$

which is a nonlinear hyperbolic partial differential equation. Here $\bar{\mu}$ denotes the mass of the unstrained string per unit arc length and $\boldsymbol{n}$ is the axial string force given by

$$\boldsymbol{n} = N\boldsymbol{t} \quad \text{where} \quad N = EA\Big(\epsilon + \alpha\frac{\partial \epsilon}{\partial \tau}\Big), \tag{1.2}$$

where $E$ stands for Young's modulus, $A$ is the area of the cross-section of the string, $N$ the axial tension and $\alpha$ is a damping constant which describes the viscosity.

The equation (1.1) arises from the application of the principle of conservation of linear momentum to a particular string element, whereby the constitutive law gives the relationship (1.2).

The boundary condition at the suspension point $\bar{s}=0$ is given by

$$\boldsymbol{r}(0, \tau) = \boldsymbol{0}. \tag{1.3a}$$

The boundary condition at $\bar{s}=\ell$ follows from the equation of motion of the end mass $m_S$, i.e.

$$m_S\frac{\partial^2 \boldsymbol{r_0}}{\partial \tau^2} = \boldsymbol{n}(\ell) + m(\tau)\boldsymbol{g}. \tag{1.3b}$$

from which the cable force $\boldsymbol{n}(\ell) = -N(\ell)(\partial \boldsymbol{r}(\ell)/\partial \bar{s})/\|\frac{\partial \boldsymbol{r}}{\partial \bar{s}}\|$ can be obtained. For a more detailed description, we refer to KUHN et al. [49].

In the mechanical formulation above we considered the string equation with respect to unstrained coordinate $\bar{s}$ (*Lagrange's description*). The advantage of the use of the coordinate $\bar{s}$ is that the boundary conditions can be taken at $\bar{s}=0$ and $\bar{s}=\ell$.

### 1.1.1   Problem statement for the vertical state of the string

Subsequently, we consider the simple state of our mechanical model, i.e. the vertical position of the string pendulum which yields a less complicated formulation of the equation of motion. If the string is in vertical position, then it is obvious that the position vector $\vec{r}(\bar{s}, \tau)$ to the characteristic point of the string can be decomposed in the following form

$$\vec{r} = 0 \cdot \boldsymbol{e}_x + (\bar{s} + u(\bar{s}, \tau))(-\boldsymbol{e}_y). \tag{1.4}$$

Here $u(\bar{s}, \tau)$ denotes the string displacement.

Using the scalar form of the vector $\boldsymbol{r} = \bar{s} + u(\bar{s}, \tau)$ it is easy to see that

$$\frac{\partial^2 \boldsymbol{r}}{\partial \tau^2} = \frac{\partial^2 u}{\partial \tau^2}. \tag{1.5}$$

Moreover, an asymptotic expansion of the function $u(\bar{s}, t)$ for $d\bar{s} \to 0$ yields

$$
\begin{aligned}
\epsilon &= \lim_{d\bar{s} \to 0} \frac{ds - d\bar{s}}{d\bar{s}} \\
&= \lim_{d\bar{s} \to 0} \frac{d\bar{s} + u(\bar{s} + d\bar{s}, \tau) - u(\bar{s}, \tau) - d\bar{s}}{d\bar{s}} \\
&= \lim_{d\bar{s} \to 0} \frac{u(\bar{s}, \tau) + d\bar{s}\frac{\partial u}{\partial \bar{s}} + \ldots - u(\bar{s}, \tau)}{d\bar{s}} \\
&= \lim_{d\bar{s} \to 0} \left( \frac{\partial u}{\partial \bar{s}} + \mathcal{O}(d\bar{s}) \right) = \frac{\partial u}{\partial \bar{s}}.
\end{aligned}
\tag{1.6}
$$

For the string in the vertical position the axial string force and the tangential vector have the same direction, i.e. the scalar part of axial force and axial tension coincide

$$
\boldsymbol{n} = N.
\tag{1.7}
$$

From (1.2), (1.5), (1.6) and (1.7), the equation (1.1) reads

$$
\begin{aligned}
\bar{\mu}\frac{\partial^2 u}{\partial \tau^2} &= \frac{\partial \boldsymbol{n}}{\partial \bar{s}} - \bar{\mu}\boldsymbol{g} \\
&= EA\frac{\partial}{\partial \bar{s}}\left( \epsilon + \alpha\frac{\partial \epsilon}{\partial \tau} \right) - \bar{\mu}\boldsymbol{g} \\
&= EA\frac{\partial^2 u}{\partial \bar{s}^2} + \alpha EA\frac{\partial^3 u}{\partial \bar{s}^2 \partial \tau} - \bar{\mu}\boldsymbol{g}.
\end{aligned}
\tag{1.8}
$$

Moreover, if we substitute a pair of new variables $(x, t) := (\bar{s}(EA)^{-1/2}, \tau\bar{\mu}^{-1/2})$ in (1.8), for $\varepsilon := \alpha\bar{\mu}^{-1/2}$ and $f(x, t) := -\bar{\mu}\boldsymbol{g}$, the equation simplifies to

$$
\frac{\partial^2 u}{\partial t^2}(x, t) = \frac{\partial^2 u}{\partial x^2}(x, t) + \varepsilon\frac{\partial^3 u}{\partial x^2 \partial t}(x, t) + f(x, t).
\tag{1.9}
$$

From the boundary condition (1.3a), we have for $\bar{s} = 0$

$$
0 = \boldsymbol{r}(0, \tau) = u(0, t).
\tag{1.10}
$$

Similarly, if we assume that the intensity of the axial force at the end mass suspension equals zero (free end), the second boundary condition (1.3b) is equivalent to

$$
0 = \boldsymbol{n}(\ell, \tau) = N(\ell, \tau).
$$

For $\ell = \sqrt{EA}$ according to (1.2), the last equation simplifies to

$$
0 = EA\left( \epsilon(\ell, \tau) + \alpha\frac{\partial \epsilon}{\partial \tau}(\ell, \tau) \right) = \sqrt{EA}\left( \frac{\partial u}{\partial x}(1, t) + \varepsilon\frac{\partial^2 u}{\partial x \partial t}(1, t) \right).
\tag{1.11}
$$

Finally, from (1.9), (1.10) and (1.11) it is obvious that the displacement of the string $u(x, t)$ satisfies the following system of equations

$$
\frac{\partial^2 u}{\partial t^2}(x, t) = \frac{\partial^2 u}{\partial x^2}(x, t) + \varepsilon\frac{\partial^3 u}{\partial x^2 \partial t}(x, t) + f(x, t)
\tag{1.12a}
$$

$$
u(0, t) = 0, \frac{\partial u}{\partial t}(0, t) = 0 \quad \text{and} \quad \frac{\partial u}{\partial x}(\ell, t) + \varepsilon\frac{\partial^2 u}{\partial x \partial t}(\ell, t) = 0,
\tag{1.12b}
$$

$$
u(x, 0) = u_0(x) \quad \text{and} \quad \frac{\partial u}{\partial t}u(x, 0) = u_1(x).
\tag{1.12c}
$$

Here function $u_0(x)$ is set to be a function which represents the initial length of the strain and accordingly $u_1(x)$ is the initial velocity. Moreover, the equation (1.12a) is of the hyperbolic-parabolic type due to the presence of $\varepsilon \geq 0$.

## 1.2 Mathematical foundations

In this section, we recall the definitions of some familiar function spaces, including the theory of distributions. For proofs and further details, we refer to ALBERTY [2], GRIFFEL [33], EVANS [27]. Let in the following $\Omega$ be some open set in $\mathbb{R}^n$.

### 1.2.1 Spaces of continuous functions

**Definition 1.2.1.1 (Multi-index).** Let $\mathbb{N}_0$ denote the set of all natural numbers. An $n$-tuple $\alpha = (\alpha_1, \ldots, \alpha_n) \in \mathbb{N}^n$ is called *multi-index* of length $|\alpha| := |\alpha_1| + \cdots + |\alpha_n|$ such that

$$D^\alpha = \partial_1^{\alpha_1} \ldots \partial_n^{\alpha_n}, \quad \text{where} \quad \partial_j = \partial/\partial x_j, \quad \text{for all } j = 1, \ldots, n. \qquad \square$$

**Definition 1.2.1.2 ($C^k$ space).** For $k \in \mathbb{N}$, we denote by $C^k(\Omega)$ the set of all continuous real-valued functions $u$, defined on $\Omega$, such that $D^\alpha u$ is continuous on $\Omega$ for every multi-index $\alpha$, $|\alpha| \leq k$. In case $k = \infty$, $C^\infty(\Omega)$ is defined as intersection $\bigcap_{k \geq 0} C^k(\Omega)$.
By $C^k(\bar{\Omega})$, we also denote the space of all $u \in C^k(\Omega)$ such that $D^\alpha u$ can be continuously extended from $\Omega$ to $\bar{\Omega}$ for every multi-index $\alpha$, $|\alpha| \leq k$. $C^\infty(\bar{\Omega})$ and $C^0(\bar{\Omega})$ are defined analogously.
In case of $\Omega$ being a bounded open set in $\mathbb{R}^n$, the linear space $C^k(\bar{\Omega})$, $k \in \mathbb{N}$ is a Banach space equipped with a norm

$$\|u\|_{C^k(\bar{\Omega})} = \max_{|\alpha| \leq k} \sup_{x \in \Omega} |\partial^\alpha u(x)|. \qquad \square$$

**Definition 1.2.1.3 (Support, $C_0^k$ space).** The support of a continuous function $u \in C^k(\Omega)$ is defined by

$$\text{supp } u := \overline{\{x \in \Omega \mid u(x) \neq 0\}}.$$

For all $k \in \mathbb{N}$, $C_0^k(\Omega)$ is defined as a set of all $u \in C^k(\Omega)$, whose support is a compact subset of $\Omega$. $\qquad \square$

### 1.2.2 Spaces of integrable functions

**Definition 1.2.2.1 ($L^p$ space).** $(L^p(\Omega), \|\cdot\|_{L^p(\omega)})$ on a Lebesgue measurable subset $\Omega \in \mathbb{R}^n$ is a Banach space of Lebesgue integrable functions equipped with a norm

$$\|u\|_{L^p(\Omega)} := \begin{cases} \left( \int_\Omega |u|^p d\Omega \right)^{1/p}, & p < \infty, \\ \text{supess}_{x \in \Omega} |u(x)|, & p = \infty. \end{cases}$$

In particular, for $p = 2$, $(L^2(\Omega), (\cdot; \cdot))$ is a Hilbert space with the inner product

$$(u; v) := \int_\Omega uv \, d\Omega. \qquad \square$$

**Definition 1.2.2.2 (Locally integrable).** Let $0 \leq p \leq \infty$. We denote by

$$L_{loc}^p(\Omega) := \{f | f \in L^p(K) \text{ for all compact } K \subset int \, \Omega\},$$

the space of all locally integrable $L^p$ functions. $\qquad \square$

**Lemma 1.2.2.1 (Hölder inequality).** Let $u \in L^p(\Omega)$ and $v \in L^q(\Omega)$, where $1/p + 1/q = 1$, $1 \leq p, q \leq \infty$. Then $uv \in L^1(\Omega)$ and

$$\left| \int_\Omega uv d\Omega \right| \leq \|u\|_{L^p(\Omega)} \|v\|_{L^q(\Omega)}.$$

For $p = q = 2$, this is referred to as the *Cauchy-Schwarz inequality*.

**Proof.** See EVANS [27, Appendix B.2]. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

### 1.2.3   Theory of Distributions

Let $(H, (\cdot\,; \cdot))$ be some Hilbert space with corresponding scalar product.

**Definition 1.2.3.1 (Piecewise).** Given the partition $\bigcup_{j=1}^N I_j = \mathcal{I}$ of some arbitrary domain $\mathcal{I}$ we say that function $f$ has some property *piecewise* on $\mathcal{I}$ if for each $I_j$, the restriction $f|_{I_j}$ has the same property. Note that the function $f$ does not have to satisfy the property on the whole interval $\mathcal{I}$. $\qquad\qquad$ $\square$

**Definition 1.2.3.2 (Space of test functions).** Let $C_0^\infty(\Omega; H) =: \mathcal{D}(\Omega; H)$ denote the space of infinitely differentiable functions $\phi : H \to \mathbb{R}$ with compact support in $\Omega$. We say that the function belonging to $\mathcal{D}(\Omega; H)$ is a *test function* and $\mathcal{D}(\Omega; H)$ is a *test space*. $\qquad$ $\square$

**Definition 1.2.3.3 (Distribution).** A linear, sequentially continuous functional

$$f : \mathcal{D}(\Omega; H) \to \mathbb{R}$$

is called *distribution* on $\mathcal{D}(\Omega; H)$.
For the definition of the convergence of sequence in the locally convex vector space, see MCLEAN [53, Chapter 3]. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

**Remark 1.2.3.1.** In one-dimensional case, $\Omega \subseteq \mathbb{R}$, we use the following notation

$$\mathcal{D}(\Omega; H) =: \mathcal{D}(\Omega) \quad \text{and accordingly} \quad \mathcal{D}^*(\Omega; H) =: \mathcal{D}(\Omega),$$

where $H$ denotes the usual $L^2(\Omega)$ Hilbert space with a corresponding scalar product. $\qquad$ $\square$

**Example 1.2.1 (Dirac delta distribution $\delta$).** The Dirac delta distribution $\delta \in \mathcal{D}^*(\mathbb{R})$ is defined by

$$(\delta; \phi) := \phi(0), \quad \text{for all} \quad \phi \in \mathcal{D}(\mathbb{R}).$$

**Example 1.2.2 (Shifted Dirac distribution).** For some $a \in \mathbb{R}$ and all $\phi \in \mathcal{D}(\mathbb{R})$ test functions, distribution $\delta_a$ with a property $(\delta_a; \phi) = \phi(a)$ for all $\phi \in \mathcal{D}(\mathbb{R})$ is called *shifted Dirac distribution*.

**Definition 1.2.3.4 (Distributional derivative, weak derivative).** Suppose $u \in L^1_{loc}(\Omega; H)$ and $\alpha$ is a multi index of order $|\alpha| := \alpha_1 + \cdots + \alpha_n$. If there exists some distribution $v \in \mathcal{D}^*(\Omega; H)$ such that

$$\int_\Omega (v; \phi) d\Omega = (-1)^{|\alpha|} \int_\Omega (u; D^\alpha \phi) d\Omega \quad \text{for all} \quad \phi \in \mathcal{D}(\Omega; H),$$

we say that $v$ is a distributional derivative of $u$ with respect to $\alpha$, where

$$D^\alpha \phi := \left( \frac{\partial^\alpha}{\partial_1^{\alpha_1} \dots \partial_n^{\alpha_n}} \right) \phi$$

with $D^\alpha := v$ $\alpha$th *weak partial derivative* of $u$. $\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

**Remark 1.2.3.2.** If $u \in L^1_{loc}$ is $m$ times continuously differentiable, then we may replace $v$ from the previous definition with $D^\alpha u$ in case when $|\alpha| \leq m$. $\qquad\square$

**Example 1.2.3 (Distributional derivative of Heaviside function).** For the Heaviside function $\mathbb{H} : \mathbb{R} \to \mathbb{R}$, defined as

$$\mathbb{H}(t) := \left\{ \begin{array}{ll} 0, & t \leq 0 \\ 1, & t > 1 \end{array} \right.$$

and Dirac distribution $\delta \in \mathcal{D}^*(\mathbb{R})$, the following result is valid

$$\mathbb{H}' = \delta \quad \text{in the distributional sense.}$$

The proof is by integration by parts, if we take some test function $\phi \in \mathcal{D}(\mathbb{R})$ such that $\text{supp}(\phi) = (-n, n)$, then

$$-\int_{\mathbb{R}} \mathbb{H}(t) \frac{\partial \phi}{\partial t} dt = -\int_0^n \frac{\partial \phi}{\partial t} dt = -\phi(n) + \phi(0) = \phi(0) = (\delta; \phi).$$

**Definition 1.2.3.5 (Derivative of distribution).** The derivative of some distribution $u \in \mathcal{D}^*(\Omega; H)$ is also distribution $v \in \mathcal{D}^*(\Omega; H)$ such that

$$(v; \phi) = (-1)^{|\alpha|}(f; D^\alpha \phi) \quad \text{for all} \quad \phi \in \mathcal{D}(\Omega; H). \qquad\square$$

### 1.2.4   Sobolev Spaces

Let $k$ be some nonnegative integer and $1 \leq p \leq \infty$.

**Definition 1.2.4.1 (Sobolev space).** We define Sobolev space $W^{k,p}(\Omega)$ as

$$W^{k,p}(\Omega) := \left\{ u \in L^1_{loc}(\Omega) \mid \text{for all} \ |\alpha| \leq k : D^\alpha \text{ exists in a weak sense and } \|D^\alpha u\|_{L^p(\Omega)} < \infty \right\}.$$
$\qquad\square$

**Definition 1.2.4.2 (Sobolev norm).** If $u \in W^k_p(\Omega)$, we define its norm to be

$$\|u\|_{W^{k,p}(\Omega)} := \left\{ \begin{array}{ll} \left( \sum_{|\alpha| \leq k} \|D^\alpha u\|^p_{L^p(\Omega)} \right)^{1/p}, & p < \infty, \\ \max_{|\alpha| \leq k} \|D^\alpha u\|_{L^\infty(\Omega)}, & p = \infty. \end{array} \right. \qquad\square$$

**Remark 1.2.4.1.** For all $1 \leq p \leq \infty$, $(W^{k,p}; \|\cdot\|_{W^{k,p}(\Omega)})$ is a Banach space, whereas in case of $p = 2$, $(W^{k,2}, (\cdot; \cdot)_{k,\Omega})$ is a Hilbert space equipped with a scalar product defined as

$$(u; v)_{k,\Omega} := \sum_{|\alpha| \leq k} \int_\Omega D^\alpha u D^\alpha v \, d\Omega.$$

We will henceforth write $H^k(\Omega) := W^{k,2}(\Omega)$ and accordingly for the corresponding norm. In particular, $H^0(\Omega) := L^2(\Omega)$ and $\|\cdot\|_{H^0(\Omega)} := \|\cdot\|_{L^2(\Omega)}$. $\qquad\square$

## 1.3 Mathematical model

Suppose that $\Omega$ is a one dimensional spatial domain, $\Omega := (0,1)$ and $T > 0$. In time-space domain $Q:=(0,T)\times\Omega$ for $\Delta := D^2$, $D^m := \partial^m/\partial x^m$ and $\dot{()} := \partial/\partial t$, $\ddot{()} := \partial^2/\partial t^2$ and given data

$$y_0 \in H^1(\Omega), \; y_1 \in L^2(\Omega), \; \text{and} \; f \in L^2(Q)$$

with a damping parameter $\varepsilon \geq 0$, the **continuous problem** reads: Find $y := Q \to \mathbb{R}$ such that

$$\ddot{y} - \Delta y - \varepsilon \Delta \dot{y} = f \; \text{on} \quad Q, \tag{1.13a}$$

subject to the initial conditions

$$y(0,x) = y_0(x), \; \dot{y}(0,x) = y_1(x) \text{ on } \Omega, \tag{1.13b}$$

and homogenous mixed Dirichlet-Neumann boundary conditions

$$y(t,0) = 0, \; Dy(t,1)+\varepsilon D\dot{y}(t,1) = 0, \quad \text{on} \quad [0,T] \qquad \text{(DN)}, \tag{1.13c}$$

or Dirichlet boundary conditions

$$y(t,0) = 0, \; y(t,1) = 0, \quad \text{on} \quad [0,T] \qquad \qquad \text{(DD)}. \tag{1.13d}$$

We refer to continuous problem as Problem (DN) when (1.13a)–(1.13b)–(1.13c) hold and Problem (DD) provided (1.13a)–(1.13b)–(1.13d). Here D designates *Dirichlet* and N *Neumann.*

**Remark 1.3.0.2.** For the visco-elastic parameter $\varepsilon = 0$, (1.13) is the wave equation. For $\varepsilon > 0$, (1.13) is called strongly damped wave equation. $\qquad\square$

**Remark 1.3.0.3.** Note that the consistency condition $(y_0, y_1) \in \mathscr{D}(A)$ in Definition 1.3.0.7 below requires further conditions on the data $y_0$ and $y_1$. $\qquad\square$

In the following, we introduce the steady state equation corresponding to the strongly damped equation (1.13). Namely, let $z$ be the unique solution of

$$-\Delta z = f \quad \text{on} \quad \Omega, \tag{1.14a}$$

subject to the homogenous boundary conditions

$$z|_{\Gamma_D} = 0, \; Dz|_{\Gamma_N} = 0. \tag{1.14b}$$

Here the boundary $\Gamma := \partial\Omega = \{0,1\}$ and $\Gamma_D = \Gamma \setminus \Gamma_N$ such that for (DN) we have $\Gamma_D = \{0\}, \Gamma_N = \{1\}$ and in case of (DD) $\Gamma_D = \{0,1\} = \Gamma$, $\Gamma_D = \emptyset$.

**Remark 1.3.0.4.** Given $f \in L^2(\Omega)$, the existence of a unique solution $z \in H^2(\Omega)$ of (1.14) is proven, e.g. in AFEM [16, Chapter 1, Theorem 1.5]. $\qquad\square$

Let here and in the following

$$a(u; v) := \int_\Omega DuDv \, dx \quad \text{and} \quad (u; v) := \int_\Omega uv \, dx$$

for all $u, v$ from $H^1(\Omega)$ and $L^2(\Omega)$, respectively.

We define $H_D^1(\Omega) := \{w \in H^1(\Omega) \mid w|_{\Gamma_D} = 0\}$. Because of $\Gamma_D \neq \emptyset$, $\|u\|_{H_D^1(\Omega)} := (a(u, u))^{1/2}$ defines a Hilbert norm on $H_D^1(\Omega)$ and we write $\|\cdot\|_{H^1(\Omega)}$ instead of $\|\cdot\|_{H_D^1(\Omega)}$.
Moreover, let $H^{-1}(\Omega)$ be the dual space of $H_D^1(\Omega)$ equipped with the usual norm

$$\|u\|_{H^{-1}} := \sup_{\|v\|_{H^1(\Omega)} \leq 1} (u; v) = \sup_{v \neq 0} \frac{(u; v)}{\|v\|_{H^1(\Omega)}} \quad \text{for} \quad u \in H^{-1}(\Omega), \tag{1.15}$$

Here the duality brackets $(\cdot, \cdot)$ extend the $L^2$ scalar product. In the following we introduce some definitions needed for the analysis of the long time behavior of the initial equation (1.13).

**Definition 1.3.0.3 (Energy scalar product, Energy norm, Hilbert space $\mathcal{H}$).** Let $\mathcal{H}$ be the Hilbert space

$$\mathcal{H} := H_D^1(\Omega) \times L^2(\Omega) \tag{1.16}$$

equipped with the canonical inner product defined for all $u = (u_1, u_2)$, $v = (v_1, v_2)$ in $\mathcal{H}$ by

$$\langle u; v \rangle_\mathcal{H} := a(u_1; v_1) + (u_2; v_2). \tag{1.17}$$

The scalar product $\langle \cdot; \cdot \rangle_\mathcal{H}$ is often referred to as the *energy scalar product* and the corresponding *energy norm* is accordingly defined for all $u = (u_1, u_2) \in \mathcal{H}$ by

$$\|u\|_\mathcal{H}^2 := \langle u; u \rangle_\mathcal{H} = \|Du_1\|_{L^2(\Omega)}^2 + \|u_2\|_{L^2(\Omega)}^2. \tag{1.18}$$

$\square$

**Definition 1.3.0.4 (Energy).** The energy of the continuous solution of Problem (1.13) at a time point $t \in [0, T]$ is the halved sum of potential and kinetic energy, i.e.

$$\mathcal{E}(y(t)) := \frac{1}{2} \int_\Omega |Dy(t, x)|^2 \, dx + \frac{1}{2} \int_\Omega |\dot{y}(t, x)|^2 \, dx. \qquad \square$$

**Remark 1.3.0.5.** From the definition of the energy norm, cf. (1.18), it is obvious that $\|(y, \dot{y})(t)\|_\mathcal{H}^2 = 2\mathcal{E}(y)$.

$\square$

**Lemma 1.3.0.1 (Conservation of Energy).** The energy of the continuous solution $y$ of the homogenous strongly damped wave equation (1.13) for $\varepsilon > 0$ dissipates in time

$$\mathcal{E}(y(t)) = \mathcal{E}(y(0)) - \varepsilon \int_0^t \|D\dot{y}(\tau)\|_{L^2(\Omega)} d\tau \quad \text{for all} \quad t \in [0, T]. \tag{1.19}$$

**Proof.** To prove this, we multiply the equation (1.13a) by $\dot{y}$ with respect to the inner product in $L^2(\Omega)$. Together with the boundary conditions, an integration by parts by proves

$$(\ddot{y}(t); \dot{y}(t)) + a(y(t); \dot{y}(t)) + \varepsilon a(\dot{y}(t); \dot{y}(t)) = 0.$$

An integration over $(0, T)$ leads to

$$\frac{1}{2} \int_0^T \frac{\partial}{\partial t} \|\dot{y}(t)\|^2 dt + \frac{1}{2} \int_0^T \frac{\partial}{\partial t} \|y(t)\|^2_{H^1(\Omega)} + \varepsilon \int_0^T \|\dot{y}(t)\|^2_{H^1(\Omega)} dt = 0, \qquad (1.20)$$

and the main theorem of calculus concludes the proof. $\qquad\square$

**Definition 1.3.0.5 (Operator $\mathcal{K}$).** For given $f \in H^{-1}(\Omega)$, there is a unique weak solution $u \in H^1_D(\Omega)$ of the problem (1.14), i.e. $a(u, v) = (f, v)$ for all $v \in H^1_D(\Omega)$, as follows from the fundamental Riesz theorem. Hence we may define the operator $\mathcal{K} : H^{-1}(\Omega) \to H^1_D(\Omega)$ which maps $f$ onto the solution $\mathcal{K}f = u$ of (1.14). $\qquad\square$

**Lemma 1.3.0.2.** Since $H^1_D(\Omega) \subset L^2(\Omega) \subset H^{-1}(\Omega)$, we may consider the restricted operator $\mathcal{K} : L^2(\Omega) \to L^2(\Omega)$. This restriction satisfies

(a) $\mathcal{K}$ is self-adjoint, positive semi-definite and compact from $L^2(\Omega)$ to $L^2(\Omega)$,

(b) Therefore, $\mathcal{K}^{1/2} : L^2(\Omega) \to L^2(\Omega)$ is well-defined, self-adjoint, compact and positive semi-definite.

(c) $\mathcal{K}^{1/2}$ can be extended to an operator $\mathcal{K}^{1/2} : H^{-1}(\Omega) \to L^2(\Omega)$ and there holds

$$\|\mathcal{K}^{1/2} f\|_{L^2(\Omega)} = \|f\|_{H^{-1}(\Omega)} \quad \text{for all} \quad f \in H^{-1}(\Omega), \qquad (1.21)$$

i.e. $\mathcal{K}^{1/2}$ is an isometry.

**Proof.** (a) We first show that $\mathcal{K}$ is self-adjoint: For $f \in L^2(\Omega)$ there holds $a(\mathcal{K}f; v) = (f, v)$ for all $v \in H^1_D(\Omega)$ according to the definition of $\mathcal{K}$. Thus, $f, g \in L^2(\Omega)$ implies

$$(f, \mathcal{K}g) = a(\mathcal{K}f; \mathcal{K}g) = a(\mathcal{K}g; \mathcal{K}f) = (g; \mathcal{K}f). \qquad (1.22)$$

The same argument implies

$$(\mathcal{K}f; f) = a(\mathcal{K}f; \mathcal{K}f) = \|\mathcal{K}f\|^2_{H^1(\Omega)} \geq 0.$$

i.e. $\mathcal{K}$ is positive semi-definite.
It remains to prove that $\mathcal{K}$ is continuous from $L^2(\Omega)$ to $H^1_D(\Omega)$ which would yield the compactness according to the Rellich theorem, cf. ZEIDLER [75, Band II, Theorem 19.25].
Due to $\mathcal{K}f \in H^1_D(\Omega)$ a Friedrich's inequality yields

$$\|\mathcal{K}f\|^2_{H^1(\Omega)} = a(\mathcal{K}f; \mathcal{K}f) = (f, \mathcal{K}f) \leq \|f\|_{L^2(\Omega)} \|\mathcal{K}f\|_{L^2(\Omega)} \leq C \|f\|_{L^2(\Omega)} \|\mathcal{K}f\|_{H^1(\Omega)}. \qquad (1.23)$$

This concludes the proof for (a).

(b) The existence of a positive, self-adjoint and compact square root $\mathcal{K}^{1/2}$ is proven e.g. in WERNER [73, Theorem VI.3.4].

(c) For the proof of (1.21), recall the definition of $\|\cdot\|_{H^{-1}(\Omega)}$ from (1.15).

First, let $f \in L^2(\Omega)$, $f \neq 0$. Obviously, there holds $\|\mathcal{K}f\|_{H^1(\Omega)} = a(\mathcal{K}f; \mathcal{K}f)^{1/2} = (f; \mathcal{K}f)^{1/2} = \|\mathcal{K}^{1/2}f\|_{L^2(\Omega)}$. Thus, the Cauchy inequality yields

$$\|f\|_{H^{-1}(\Omega)} = \sup_{v \in H^1(\Omega)} \frac{(f;v)}{\|v\|_{H^1(\Omega)}} = \sup_{v \in H^1(\Omega)} \frac{a(\mathcal{K}f;v)}{\|v\|_{H^1(\Omega)}} \leq \sup_{v \in H^1_D(\Omega)} \frac{\|\mathcal{K}f\|_{H^1(\Omega)}\|v\|_{H^1(\Omega)}}{\|v\|_{H^1(\Omega)}} = \|\mathcal{K}^{1/2}f\|_{L^2(\Omega)}.$$

In particular $\mathcal{K}^{1/2}f \neq 0$. Therefore, we may also conclude that

$$\|f\|_{H^{-1}(\Omega)} = \sup_{v \neq 0}(f;v) \geq \frac{(f;\mathcal{K}f)}{\|\mathcal{K}f\|_{H^1(\Omega)}} = \frac{(f;\mathcal{K}f)}{(f;\mathcal{K}f)^{1/2}} = \|\mathcal{K}^{1/2}f\|_{L^2(\Omega)}.$$

From the latter two inequalities we obtain $\|f\|_{H^{-1}(\Omega)} = \|\mathcal{K}^{1/2}f\|_{L^2(\Omega)}$ for $f \in L^2(\Omega)$. In particular, $\mathcal{K}^{1/2}$ is an isometric operator from $(L^2(\Omega), \|\cdot\|_{H^{-1}(\Omega)})$ to $L^2(\Omega)$. Continuity allows to extend $\mathcal{K}^{1/2}$ to the whole of $H^{-1}(\Omega)$, since $L^2(\Omega) \hookrightarrow H^{-1}(\Omega)$, see ODEN-REDDY [57, Section 4.5]. This concludes the proof.  □

**Definition 1.3.0.6 (Hilbert space $\widehat{\mathcal{H}}$).** Let $\widehat{\mathcal{H}}$ be the Hilbert space

$$\widehat{\mathcal{H}} := L^2(\Omega) \times H^{-1}(\Omega) \tag{1.24}$$

equipped with the inner products defined for all $u = (u_1, u_2)$, $v = (v_1, v_2)$ by

$$\langle u;v\rangle_{\widehat{\mathcal{H}}} := (u_1;v_1) + (\mathcal{K}^{1/2}u_2;v_2). \tag{1.25}$$

Here $\mathcal{K}$ is an operator from Definition 1.3.0.5.
The corresponding norm is defined for all $u = (u_1, u_2) \in \widehat{\mathcal{H}}$ through

$$\|u\|_{\widehat{\mathcal{H}}}^2 := \langle u;u\rangle_{\widehat{\mathcal{H}}} = \|u_1\|_{L^2(\Omega)}^2 + \|\mathcal{K}^{1/2}u_2\|_{L^2(\Omega)}^2. \tag{1.26}$$
□

In the following we introduce the vector notation, in order to present initial equation as the evolution equation.

**Definition 1.3.0.7 (Vector formulation, Operator $\mathcal{A}$).** For the solution of the initial problem (1.13) let

$$u(t,x) := (y(t,x), \dot{y}(t,x)), \ u_0(x) = (y_0(x), y_1(x)), \ F(t,x) := (0, f(t,x)),$$

and define the $2 \times 2$ operator matrix $\mathcal{A}:\mathcal{H}\to\mathcal{H}$ by

$$\mathcal{A} := \begin{bmatrix} 0 & -1 \\ -\Delta & -\varepsilon\Delta \end{bmatrix}. \tag{1.27}$$

The problem (1.13) is equivalent to: Find $u \in H^1(0,T;\mathcal{H})$ such that

$$\dot{u} + \mathcal{A}u = F, \quad \text{on} \quad Q, \tag{1.28a}$$
$$u(0,x) = u_0(x) \quad \text{on} \quad \Omega, \tag{1.28b}$$

for $\mathcal{A} : \mathscr{D}(\mathcal{A}) \subset \mathcal{H} \to \mathcal{H}$ where

$$\mathscr{D}(\mathcal{A}) := \{(u_1,u_2) \subset H^1_D(\Omega) \times H^1_D(\Omega) \mid u_1 + \varepsilon u_2 \in H^2(\Omega) \text{ and } D(u_1 + \varepsilon u_2)|_{\Gamma_N} = 0\}. \quad □$$

**Lemma 1.3.0.3.** The operator $\mathcal{A}$ is maximal monotone.

**Proof.** First we need to prove the monotonicity. Let $u = (u_1, u_2)$ and $v = (v_1, v_2)$ belong to $\mathscr{D}(\mathcal{A})$. An integration by parts in space using the boundary conditions yields to

$$
\begin{aligned}
\langle u - v \, ; \mathcal{A}u - \mathcal{A}v \rangle_{\mathcal{H}} &= a(u_1 - v_1; -u_2 + v_2) - (u_2 - v_2; \Delta(u_1 - v_1) + \varepsilon\Delta(u_2 - v_2)) \\
&= -a(u_1 - v_1; u_2 - v_2) + a(u_2 - v_2; u_1 - v_1) + \varepsilon a(u_2 - v_2; u_2 - v_2) \\
&\quad - [(u_2 - v_2)(D(u_1 - v_1 + \varepsilon(u_2 - v_2)))]_0^1 = \varepsilon \int_\Omega (D(u_2 - v_2))^2 dx \geq 0.
\end{aligned}
$$

This shows that the operator $\mathcal{A}$ is a monotone operator.
By definition it remains to prove that the range of $\mathcal{A} + I$ satisfies $R(\mathcal{A} + I) = \mathcal{H}$ with $I$ the identity operator on $\mathcal{H}$.
For arbitrary $(e, g) \in \mathcal{H}$, we consider the following ODE problem

$$z - (1 + \varepsilon)\Delta z = e + (1 + \varepsilon)g, \tag{1.29a}$$

$$z|_{\Gamma_D} = 0, \quad Dz|_{\Gamma_N} = 0. \tag{1.29b}$$

According to AFEM [16, Chapter 1, Theorem 1.5], there exists a unique solution $z \in H^2(\Omega) \cap H_D^1(\Omega)$ of (1.29).
Accordingly $e \in H_D^1(\Omega)$ and the functions

$$m := \frac{1}{1 + \varepsilon}(z + \varepsilon e) \quad \text{and} \quad n := \frac{1}{1 + \varepsilon}(z - e)$$

satisfy $m, n \in H_D^1(\Omega)$, $m + \varepsilon n = z \in H^2(\Omega)$, $D(m + \varepsilon n)|_{\Gamma_N} = 0$, and

$$m - n = e, \tag{1.30a}$$

$$n - \Delta(m + \varepsilon n) = g. \tag{1.30b}$$

The system (1.30a) is equivalent to

$$[I + \mathcal{A}] \begin{bmatrix} m \\ n \end{bmatrix} = \begin{bmatrix} e \\ g \end{bmatrix}. \tag{1.31}$$

Hence $\mathcal{H} \subseteq R(\mathcal{A} + I)$ and from $R(\mathcal{A} + I) \subseteq \mathcal{H}$, which is obvious due to the definition of $\mathscr{D}(\mathcal{A})$, we may conclude the proof of lemma. $\qquad \square$

**Lemma 1.3.0.4.** $\mathcal{H}$ is separable Hilbert space as a product of two separable Hilbert spaces $H_D^1(\Omega)$ and $L_2(\Omega)$.

**Proof.** Follows directly from the definition of separable spaces. $\qquad \square$

**Theorem 1.3.0.1 (ZEIDLER [75, Band II, Section 31.1, Theorem 31.A]).** If $\mathcal{A}: \mathscr{D}(\mathcal{A}) \subseteq \mathcal{H} \to \mathcal{H}$ is a maximal monotone operator with respect to domain $\mathscr{D}(\mathcal{A})$ and $F \in H^1(0, T; \mathcal{H})$, then problem (1.28) has exactly one solution $u \in H^1(0, T; \mathcal{H})$.

**Proof.** The proof follows since $\mathcal{H}$ is a real separable space and $\mathcal{A}$ maximal monotone operator. $\qquad \square$

**Remark 1.3.0.6.** The motivation to use the above introduced vector formulation (Definition 1.3.0.7 and 1.3.0.3) comes from the analytical study of the strongly damped wave equation with parameter $\varepsilon = 1$, cf. NEVES [54]. $\qquad\square$

**Lemma 1.3.0.5 ($H^2$ stability of continuous solution).** For the solution $y$ of (1.13) with $f \equiv 0$, we have

$$\|\ddot{y}(t)\|^2_{L^2(\Omega)} + \|\dot{y}(t)\|^2_{H^1(\Omega)} + \|\Delta(y + \varepsilon\dot{y})(t)\|^2_{L^2(\Omega)} + 2\varepsilon\int_0^t \|\ddot{y}(\tau)\|^2_{H^1(\Omega)}d\tau \leq 2\big(\|y_1\|^2_{H^1(\Omega)} + \|\Delta(y_0 + \varepsilon y_1)\|^2_{L^2(\Omega)}\big).$$

**Proof.** According to Theorem 1.3.0.1, there is a unique solution $u$ of Problem (1.28) for $u_0 = (y_0, y_1) \in \mathscr{D}(\mathcal{A})$. This solution is the unique solution of the problem (1.13) too, meaning $u = (y, \dot{y})$ for $y$ the solution of (1.13).
The proof follows by multiplying the initial equation (1.13a) with $-\Delta(\dot{y} + \varepsilon\ddot{y})$ with respect to $L^2$ scalar product in space and further integrating over time interval $[0, t]$ where $t$ is the arbitrary point in time. Namely

$$0 = -\int_0^t (\ddot{y}; \Delta(\dot{y} + \varepsilon\ddot{y}))d\tau + \int_0^t (\Delta(y + \varepsilon\dot{y}); \Delta(\dot{y} + \varepsilon\ddot{y}))d\tau. \tag{1.32}$$

The partial integration in space in case of the first term on the LHS of (1.32) shows

$$0 = \int_0^t a(\ddot{y}; \dot{y} + \varepsilon\ddot{y})d\tau + \int_0^t (\Delta(y + \varepsilon\dot{y}); \Delta(\dot{y} + \varepsilon\ddot{y}))d\tau$$
$$= \frac{1}{2}\int_0^t \frac{\partial}{\partial t}\|\dot{y}(\tau)\|^2_{H^1(\Omega)}d\tau + \varepsilon\int_0^t \|\ddot{y}(\tau)\|^2_{H^1(\Omega)}d\tau + \frac{1}{2}\int_0^t \frac{\partial}{\partial t}\|\Delta(y + \varepsilon\dot{y})(\tau)\|^2_{L^2(\Omega)}d\tau \tag{1.33}$$

Here we used the boundary conditions (1.13c) as well as (1.13d). Moreover, owing to the main theorem in calculus, the last expression is equivalent to

$$\|\dot{y}(t)\|^2_{H^1(\Omega)} + \|\Delta(y + \varepsilon\dot{y})(t)\|^2_{L^2(\Omega)} + 2\varepsilon\int_0^t \|\ddot{y}(\tau)\|^2_{H^1(\Omega)}d\tau = \|y_1\|^2_{H^1(\Omega)} + \|\Delta(y_0 + \varepsilon y_1)\|^2_{L^2(\Omega)}. \tag{1.34}$$

From (1.34) and the fact that $\ddot{y} = \Delta(y + \varepsilon\dot{y})$ we additionally have

$$\|\ddot{y}(t)\|^2_{L^2(\Omega)} \leq \|y_1\|^2_{H^1(\Omega)} + \|\Delta(y_0 + \varepsilon y_1)\|^2_{L^2(\Omega)}. \tag{1.35}$$

The proof follows by summing (1.34) and (1.35). $\qquad\square$

In case of the norms and normed spaces defined on the whole domain $Q$, we shorten the notation and write e.g. $\|\cdot\|_{L^1(L^2)}$ instead of $\|\cdot\|_{L^1(0,T;L^2(\Omega))}$ and proceed similarly in case of Sobolev norms.

# Chapter 2

# Discrete model

For the numerical simulation of the strongly damped wave equation and a proof of a sharp a priori and a posteriori error estimates, we developed several approximation schemes which result in a fully discrete analogon of initial problem (1.13). The proposed discretisation methods are applied to the vectorised form of the initial equation (1.28), such that both components of the solution vector belong to the same discrete space. The usage and approximation of the vectorised form instead of the initial form of equation (1.13) is motivated by the fact that across from very few known methods for the discretisation of the second time derivative, the variety of numerical methods for solving differential equations of first order is widely used and fully discussed. This also yields to the better approximation of the second variable $\dot{y}$ (velocity) and is therefore favourable. As an example for both approaches, we refer to HUGHES-HULBERT [38, 41].

In case of spatial discretisation, due to the presence of a second order differential operator in space, we employ globally continuous finite elements, e.g. linear splines ($\mathcal{P}_1$ conform finite elements) and cubic splines ($\mathcal{C}^1$ finite elements), respectively, cf. Section 2.1. A time discretisation is carried out with aid of the discontinuous Galerkin methods of order $0, 1$ as well as the continuous Galerkin method of order 1, see Section 2.2. As a successful tool for the analysis of the long time behaviour of the initial equation and its discrete counterpart, we also concern the spatial semi-discrete version of the strongly damped wave equation. Each of these methods are explained in the following, along with a detailed overview concerning the implementation, cf. Section 6.1. We will also show that each discrete model provides the unique existence of a discrete solution, see Section 2.4.

First, let us introduce the time-space partition of the domain $Q = [0, T] \times \Omega$.

**Definition 2.0.0.8 (Partition in time and space).** Let $\mathscr{T}$ be a partition of the time interval $[0, T]$,

$$\mathscr{T} : 0 = t_0 < t_1 = t_0 + k_1 < t_2 = t_1 + k_2 < \cdots < t_N = t_{N-1} + k_N = T,$$

where $I_j := (t_{j-1}, t_j]$, $k_j := |I_j|$ and $k \in L^\infty(0, T)$ is defined by $k|_{I_j} := k_j$ for each $j = 1, \ldots, N$. With each time interval $I_j$, we associate a triangulation $\mathcal{T}_j$ of the spatial domain $\Omega = [0, 1]$

$$\mathcal{T}_j : x_0 = 0 < x_1 < x_2 < \ldots < x_{n-1} < x_n = 1.$$

As above we write $T_k = (x_{k-1}, x_k]$, $h_k := |T_k|$ and define $h \in L^\infty(\Omega)$ by $h|_{T_k} := h_k$ for each $k = 1, \ldots, n$. $\qquad\square$

The space discretisation may vary in both space and time, but the time steps are only variable in time, not in space so that the corresponding time-space mesh is not fully optimal. This variation of spatial domain arises when coercing/refinement procedure in space is applied.

For different ways of discretisation where the triangulation of temporal domain is not fixed as a reference for space discretisation we refer to RICHTER [61].

Moreover, when the mesh varies in space, we need to find a way of transferring information from one time level to the next, so that we can approximate the numerical solution at a point of the previous spatial mesh, which may not have been a node. This can be done in a several ways: Some authors, e.g. BANGERTH [9] propose the usage of hierarchical meshes (important for a continuous Galerkin ansatz in time) and some introduce the idea of a global $L^2$ projection, see SÜLI-WILKINS [67]. In this work, we deal with the one-dimensional spatial domain and accordingly the straight-forward interpolation can be used.

**Remark 2.0.0.7 (Additional time step control, CFL condition).** Depending on the time discretisation, an additional time-space-step control may be needed for the stability of the discrete method. This is the case when, e.g. explicit methods are used for the discretisation in time, cf. COURANT et al. [21]. However, in this work only the implicit methods in time are considered and therefore the additional time-space step control is not necessary.         □

## 2.1   Discretisation in space

The spatial discretisation follows by means of linear splines and Hermite cubic splines, respectively.

With respect to the partition of the whole domain $Q = [0, T] \times \Omega$ from Definition 2.0.0.8, for each $j$ and the corresponding time slab $I_j \times \Omega$ we may define a spatial conforming FE space $\mathcal{S}^j \subset H_D^1(\Omega)$ such that in case of linear splines

$$\mathcal{S}^j := \mathcal{S}_D^1(\mathcal{T}_j) = \operatorname{span}\{\phi_1, \ldots, \phi_n\} \subset \mathcal{C}_D(\Omega), \tag{2.1a}$$

and when the discretisation in space is by Hermite cubic finite elements

$$\mathcal{S}^j := \mathcal{S}_D^3(\mathcal{T}_j), \quad \text{such that} \quad \mathcal{S}_D^3(\mathcal{T}_j)|_{T_k} = \operatorname{span}\{\phi_1, \ldots, \phi_4\} \subset \mathcal{C}_D^1(\Omega), \ T_k \in \mathcal{T}_j, \ 1 \leq k \leq n. \tag{2.1b}$$

The number $n + 1$ is the number of nodes in triangulation of the spatial domain $\Omega$.

In case of (2.1b), $\{\phi_\ell\}_{\ell=1}^4$ are chosen to be piecewise Hermite cubic finite elements and they are defined on each space interval separately, cf. Subsection 2.1.2. Both ansatz spaces consist of globally continuous functions and therefore conforming elements.

We choose the linear finite elements because the theory derived for the proof of a posteriori error bound can easily be extended to the case of the strongly damped wave equation with two-dimensional space variable. They also decrease the computational cost compared to the use of the cubic splines. On the other hand, we rather use cubic elements because their improved accuracy yields a much sharper a posteriori error bound. This will be shown in Chapter 3 and Chapter 4.

Note that the structure of the FE spaces $\mathcal{S}^j$ with respect to the number of elements in the corresponding triangulation $\mathcal{T}^j$ may vary between two neighbouring time slabs.

### 2.1.1 Linear Splines ($\mathcal{P}_1$)

In case of piecewise affine ($\mathcal{P}_1$), globally continuous functions which are often referred to as *linear splines*, we consider the following basis for the discrete space $\mathcal{S}^j$, namely for $k=1,\ldots,n$

$$\phi_k(x) := \begin{cases} (x - x_{k-1})/h_k, & \text{for } x \in T_k, \\ (x_{k+1} - x)/h_{k+1}, & \text{for } x \in T_{k+1}, \\ 0, & \text{elsewhere.} \end{cases} \tag{2.2}$$

These functions are also called *hat functions*, see Figure 2.1.



Figure 2.1: Basis function $\phi_1$, linear splines.

### 2.1.2 Piecewise Hermite cubic finite elements ($\mathcal{C}^1$)

Globally continuous $\mathcal{C}^1$ cubic polynomials, so-called *cubic splines* are defined on each space interval separately. Namely, if $\zeta \in [-1, 1]$ denotes the local coordinate relative to the interval $T_k = (x_{k-1}, x_k]$, where

$$x = \frac{1}{2}(1 - \zeta)x_{k-1} + \frac{1}{2}(1 + \zeta)x_k,$$

then the basis functions read

$$\phi_1(\zeta) := \frac{1}{4}(1 - \zeta)^2(2 + \zeta), \quad \phi_2(\zeta) := \frac{1}{4}(1 - \zeta)^2(1 + \zeta),$$

$$\phi_3(\zeta) := \frac{1}{4}(1 + \zeta)^2(2 - \zeta), \quad \phi_4(\zeta) := -\frac{1}{4}(1 + \zeta)^2(1 - \zeta). \tag{2.3}$$

Accordingly, the discrete solution on the $k$th element, $U \in \mathcal{S}^j|_{T_k}$ can be represented as

$$U(t, x(\zeta)) := U(t, x_{k-1})\phi_1(\zeta) + \frac{h_k}{2}DU(t, x_{k-1})\phi_2(\zeta) + U(t, x_k)\phi_3(\zeta) + \frac{h_k}{2}DU(t, x_k)\phi_4(\zeta). \tag{2.4}$$

**Remark 2.1.2.1.** The Hermite cubic finite elements are continuous in first derivative and therefore admissible for the incorporation of Neumann boundary condition (1.13c) in the computation of the discrete solution in case of Problem (DN). For details we refer to Subsection 6.1.4. $\qquad\square$

Figure 2.2: Basis functions $\phi_1, \phi_2, \phi_3, \phi_4$, cubic splines.

## 2.2   Time discretisation

In case of time discretisation we apply the following methods, namely

- Discontinuous Galerkin method, abbreviated $dG(q)$, for $q = 0, 1$ polynomial degree,

- Continuous Galerkin method, abbreviated by $cG(q)$ for $q = 1$,

- Method of lines, abbreviated by $MoL$.

The discontinuous Galerkin method generates an implicit A-stable time-stepping schemes of order $q + 1$ in time, where the super convergence of order $2q + 1$ occurs at the nodal points $t_j$ from $\mathscr{T}$. The approximate solution is sought on each time level separately. Although the $dG$ methods are strongly dissipative, due to the smoothing effects which cause the additional energy loss, see Remark 2.3.3.4, they can be understood as a general method for solving time-dependent problems adaptively. Namely, the jumps of $dG$ functions can be regarded as a proper refinement or coarsening indicators. One of the first applications of $dG$ methods for the second order hyperbolic problems was by Hulbert and Hughes, see [38, 41]. Note that in the case $q = 0$, the $dG(0)$ method coincides with the backward Euler scheme up to the computation of the right-hand side (RHS).

Another very useful method for solving linear as well as nonlinear wave problems is the continuous Galerkin method where the discretisation is carried out with aid of globally continuous functions. These methods where first introduced in context of ordinary differential equations in Hulme [39, 40]. There was also shown that $cG$ methods are closely related to Runge-Kutta schemes based on Gauss-Legendre quadrature rules with a convergence rate $k^{2q}$, where $k$ is the step size and $q$ is the polynomial degree. With this approach, methods of any desired order of accuracy can be easily formulated. The advantage of using this particular method for the time discretisation is that it conserves the energy in the same way as the continuous problem, see Remark 2.3.4.2. We will only consider the case $cG(1)$, i.e. time discretisation with globally continuous piecewise linear functions. Note that in case of a homogenous problem and equidistant time-stepping, the $cG(1)$ method coincides with an A-stable implicit mid-point scheme of convergence order 2. The subsequent analysis from Section 2.3.4 is based on

some standard ideas from BANGERTH [9], FRENCH-JENSEN [29] and FRENCH-PETERSON [30].

A very important method for the analysis of the space discretisation as well as the long time behavior of initial equation is the vertical method of lines. This technique converts the initial hyperbolic problem (1.28) to a system of ordinary differential equations which can be solved by very efficient solvers for large systems of ordinary differential equations. These solvers, e.g. Matlab solvers ode23,ode45,ode113,... (for details, we refer to [60]), possess a high accuracy and therefore can be considered as relatively reliable. The disadvantage of using this method is that this approach compared to $dG$ method does not enable us to detect the discontinuities of sharp gradients in the discrete solution. Therefore an effective time adaptive scheme can be developed here. The related analysis is based on some general ideas from BABUŠKA et al. [6] and SÜLI-WILKINS [67].

Note that the stability of all methods mentioned above, is independent of the choice of the time and space mesh increment, recall Remark 2.0.0.7 and HARTMANN [37] as well as the references therein.

In order to illustrate the respective time approximation properties, in Figure 2.3 we present the energy distribution of the discrete function in time, obtained for the application of four different methods in time and cubic splines in space in case of Example 6.2.2. The Figure 2.4 displays the same values only with different scaling where only the the method of lines, $cG(1)$ and $dG(1)$ method in time are represented.
In Figure (2.3), the dissipation of $dG(0)$ method is apparent, whereas the dissipation of $dG(1)$ method is much more visible in the Figure (2.4).



Figure 2.3: Energy of discrete solution in time $\|U\|_{\mathcal{H}}$; $MoL$, $cG(1)$, $dG(0)$, $dG(1)$ in time and $\mathcal{C}^1$ in space; Example 6.2.2, $\varepsilon = 0$, (DN), $T = 1$.

Figure 2.4: Energy of discrete solution $\|U\|_{\mathcal{H}}$ on time interval $[0,4]$; $MoL$, $cG(1)$, $dG(1)$ in time and $\mathcal{C}^1$ in space; Example 6.2.2, $\varepsilon = 0$, (DN), $T = 1$.

## 2.3   Fully discrete model

In the subsequent section we introduce the fully discrete model obtained by approximating the initial problem (1.28) in time and space. This will be done first by introducing the general formulation of the weak form and some related terms such as the residual, cf. Subsection 2.3.1. In Subsection 2.3.2 the error of the time-space discretisation is defined. The remainder of the section is devoted to the formulation and the stability analysis of the fully discrete problem for each ansatz in time in particular. This includes the formulation of the weak form for the corresponding time approximation method and analysis of the appropriate discrete space.

### 2.3.1   Weak form

Due to the choice of the time-space discrete scheme, let $\mathcal{Q}$ and $\mathcal{W}$ be the corresponding fully discrete spaces with corresponding indices for each method in time, whereby $\mathcal{Q}$ denotes the space of trial functions and $\mathcal{W}$ is referred to as the space of test functions.
Henceforth, we seek to find some approximation $U := (U_1, U_2) \in \mathcal{Q}$ of the continuous problem (1.28) such that

$$\mathcal{B}(U, V) = \mathscr{L}(V) \quad \text{for all} \quad V = (V_1, V_2) \in \mathcal{W}, \tag{2.5}$$

with an initial condition noted by $U(0)$ or $U^{0-}$ in dependence on the time approximation method. Moreover, the initial solution denotes the discrete variant of the continuous initial solution $u_0$, cf. (1.28b), i.e. $U(0) = \Pi u_0$. Here $\Pi$ denotes a spatial multi projection which will be introduced in the following.
We suppose that $\mathcal{B}$ is a bilinear form

$$\mathcal{B}(\cdot, \cdot) : H^1(\mathscr{T}; \mathcal{H}) \times L^2(0, T; H_D^1(\Omega) \times H_D^1(\Omega)) \to \mathbb{R},$$

where $\mathscr{T}$ stands for some arbitrary partition of the time interval introduced in Definition 2.0.0.8. $H^1(\mathscr{T}; \mathcal{H})$ is the space of all $\mathscr{T}$-piecewise $H^1$ functions on $[0, T]$ which are not necessary in $H^1(0, T)$.
The linear functional $\mathscr{L}$ is defined such that

$$\mathscr{L} : L^2(0, T; H_D^1(\Omega) \times H_D^1(\Omega)) \to \mathbb{R}.$$

The definition of the bilinear form $\mathcal{B}$ and the functional $\mathcal{L}$ will differ for different fully discrete spaces due to the different choice of time discretisation. This will be explicitly explained in Subsections 2.3.3–2.3.5.

**Remark 2.3.1.1 (Local discrete form, $dG(q)$, $cG(1)$ time discretisation).** In case of the Galerkin time discretisation, we define by $\mathcal{B}^j := \mathcal{B}|_{I_j}$, $\mathcal{L}^j := \mathcal{L}|_{I_j}$ the $j$th contributions of $\mathcal{B}$ and $\mathcal{L}$ respectively, which are related to the time interval $I_j$. Then, the weak form (2.5) reads: Find $U \in \mathcal{Q}|_{I_j}$ such that for given $U|_{I_{j-1}}$

$$\mathcal{B}^j(U, V) = \mathcal{L}^j(V) \quad \text{for all} \quad V \in \mathcal{W}|_{I_j} \quad \text{and} \quad j = 1, \dots, N. \tag{2.6}$$

The formulation (2.6) is equivalent to (2.5) and can be obtained through decoupling of time steps, using the fact that for each time step $I_j$ we may choose a test function $V \in \mathcal{W}$ such that $V|_{I_k} \equiv 0$ for all $k \neq j$ and since the test functions in Galerkin time discretisation method, include the characteristic functions. $\qquad\square$

**Definition 2.3.1.1 (Consistency).** The solution $u$ of initial vector problem (1.28) is called *strong solution*. A solution $U$ of weak problem (2.5) is called *weak solution*.
We also say that the bilinear form $\mathcal{B}$ is consistent with the strong formulation (1.28) if strong and weak solution coincide. $\qquad\square$

**Definition 2.3.1.2 (Residual).** Let for any $v = (v_1, v_2) \in L^2(\mathcal{T}; \mathcal{H})$ the residual be defined as

$$Res(v) := \mathcal{L}(v) - \mathcal{B}(U, v), \tag{2.7}$$

where $U$ denotes the discrete solution from (2.5). $\qquad\square$

The residual has the orthogonality property,

$$Res(V) = 0, \quad \text{for all} \quad V \in \mathcal{Q}. \tag{2.8}$$

## 2.3.2 Error

**Definition 2.3.2.1 (Error).** Given the exact solution $u$ of (1.28) and the discrete solution $U \in \mathcal{Q}$ of (2.5), we define the error

$$e := (e_1, e_2) = u - U \in H^1(\mathcal{T}; \mathcal{H}). \qquad\square$$

In view of (2.5) and Remark 2.3.1.1, $e$ satisfies so-called *Galerkin orthogonality*

$$\mathcal{B}(e, V) = 0 \quad \text{for all} \quad V \in \mathcal{Q}. \tag{2.9}$$

**Remark 2.3.2.1.** In case of a consistent discretisation scheme, see Definitions 2.3.1.1, 2.3.1.2, we may find that

$$\mathcal{B}(e, v) = Res(v) \quad \text{for all} \quad v \in L^2(0, T; \mathcal{H}). \tag{2.10}$$
$$\square$$

**Definition 2.3.2.2.** Let $E(t) := \mathcal{E}(e(t))$ denote the energy of the error $e = (e_1, e_2)$ at time $t$, i.e.

$$E(t) = \frac{1}{2} \int_\Omega |De_1(t)|^2 \, dx + \frac{1}{2} \int_\Omega |e_2(t)|^2 \, dx.$$

Note that the energy of the error is defined according to the Definition 1.3.0.4 with $e$ instead of the exact solution $u$. $\qquad\square$

### 2.3.3   Discontinuous Galerkin time discretisation

On the basis of the main properties of discontinuous Galerkin methods, we design a class of discrete functions as set of piecewise polynomials such that on each time interval $I_j \in \mathscr{T}$ a discrete function belongs to $\mathcal{P}_q(I_j)$ with the coefficients in spatial discrete space $\mathcal{S}^j$. Moreover, the discrete functions are not globally continuous, i.e. on the whole time interval $[0, T]$, and trial and test spaces coincide ($\mathcal{W} \equiv \mathcal{Q}$) which refers to *Ritz-Galerkin methods*.

The concept of distributional time derivatives of discontinuous $dG(q)$ function arises in the formulation of the weak form.

**Definition 2.3.3.1 ($dG(q)$ space, $q = 0, 1$, full discretisation).** For each $I_j \in \mathscr{T}$ from Definition 2.0.0.8, let $\mathcal{Q}_q^j$ be defined by

$$\mathcal{Q}_q^j := \left\{ U : I_j \times \Omega \to \mathbb{R}^2 \mid U(t, \cdot) \in \mathcal{S}^j \times \mathcal{S}^j, \ U(\cdot, x) \in \mathcal{P}_q(I_j) \right\} \subseteq H^1(I_j; \mathcal{H}). \qquad (2.11a)$$

Accordingly, we may define

$$\mathcal{Q}_q := \left\{ U : Q \to \mathbb{R}^2 \mid U|_{I_j \times \Omega} \in \mathcal{Q}_q^j, \text{ for all } I_j \in \mathscr{T} \right\} \subseteq H^1(\mathscr{T}; \mathcal{H}), \qquad (2.11b)$$

as a set of polynomials in time $U$, such that for each interval $I_j \in \mathscr{T}$, $U|_{I_j} := U^j$ is a polynomial of degree $q$ with coefficients in $\mathcal{S}^j \times \mathcal{S}^j$. The space $\mathcal{S}^j$ takes a form due to the choice of spatial discrete functions.

We say that $\mathcal{Q}_q$ defines the $dG(q)$ space with full discretisation in space. The functions in $\mathcal{Q}_q$ are called $dG(q)$ functions, $q = 0, 1$. $\qquad \square$

**Definition 2.3.3.2.** Let $\mathscr{T}$ be a partition of the domain $[0, T]$ introduced in Definition 2.0.0.8 and $U \in \mathcal{Q}_q$ a piecewise smooth function on each interval $I_j$, $j = 1, \ldots, N$. For given $U^{0-}$, (which stands for the initial condition), we may define the extension of function $U$ onto the entire space $\mathbb{R}$ as

$$U(t) := \begin{cases} U^{0-}, & t \in (-\infty, 0) \\ U(t), & t \in [0, T] \\ U(T), & t \in (T, +\infty). \end{cases} \qquad (2.12)$$

On account of this, the one-sided limits exist for all $j = 0, \ldots, N$ and are defined by

$$U^{j\pm} := \lim_{t \to t_j^{\pm}} U(t). \qquad (2.13)$$

Accordingly, we introduce the jumps

$$[U]^j := U^{j+} - U^{j-}. \qquad (2.14)$$

In particular, our definition implies $[U]_0 = U^{0+} - U^{0-}$ and $[U]^N = 0$. $\qquad \square$

**Lemma 2.3.3.1 (ALBERTY [2, Lemma 3.13]).** Let $U \in \mathcal{Q}_q$ be an arbitrary $dG$ function and $\mathbb{H}$, Heaviside function from Example 1.2.3 and $\delta_{t_j}$, the distributional derivative of Heaviside function, cf. Example 1.2.2. Owing to the definition of one-sided limits, cf. Definition 2.3.3.2, we have for $t \in I_j$

$$\frac{\partial}{\partial t}\left(U^j(t)\mathbb{H}(t - t_{j-1})\right) = U^{j-1+}\delta_{t_{j-1}} + \dot{U}^j(t)\mathbb{H}(t - t_{j-1}) \qquad (2.15a)$$

$$\frac{\partial}{\partial t}\left(U(t)\mathbb{H}(t_j - t)\right) = -U^{j-}\delta_{t_j} + \dot{U}^j(t)\mathbb{H}(t_j - t). \qquad (2.15b)$$

| $dG(q)$ functions with space discretisation by linear splines | |
|---|---|
| $q=0$ | $U(t,x)\|_{I_j\times\Omega}:=U^j(x)=\sum_{k=1}^m U_k^j\phi_k(x),\ \ U_k^j=(U_{k,1}^j,U_{k,2}^j)\in\mathbb{R}^2$ |
| $q=1$ | $U(t,x)\|_{I_j\times\Omega}:=U^j(t,x)=\sum_{k=1}^m\big(U_k^{j,0}+{}^{(t-t_{j-1})}\!/\!_{k_j}U_k^{j,1}\big)\phi_k(x),U_k^{j,0\|1}=(U_{k,1}^{j,0\|1},U_{k,2}^{j,0\|1})\in\mathbb{R}^2$ |

| $dG(q)$ functions with space discretisation by cubic splines | |
|---|---|
| $q=0$ | $U(t,x(\zeta))\|_{I_j\times T_k}:=U_k^j(x(\zeta))=U_{k-1}^j\phi_1(\zeta)+{}^{h_k}\!/\!_2 DU_{k-1}^j\phi_2(\zeta)+U_k^j\phi_3(\zeta)+{}^{h_k}\!/\!_2 DU_k^j\phi_4(\zeta),$ <br> $U_{k-1\|k}^j:=(U_{k-1\|k,1}^j,U_{k-1\|k,2}^j),\ \ DU_{k-1\|k}^j:=(DU_{k-1\|k,1}^j,DU_{k-1\|k,2}^j)\in\mathbb{R}^2$ |
| $q=1$ | $U(t,x(\zeta))\|_{I_j\times T_k}:=U_k^j(t,x(\zeta))={}^{(t-t_{j-1})}\!/\!_{k_j}\big(U_{k-1}^{j,0}\phi_1(\zeta)+{}^{h_k}\!/\!_2 DU_{k-1}^{j,0}\phi_2(\zeta)+U_k^{j,0}\phi_3(\zeta)$ <br> $+{}^{h_k}\!/\!_2 DU_k^{j,0}\phi_4(\zeta)\big)+{}^{(t_j-t)}\!/\!_{k_j}\big(U_{k-1}^{j,1}\phi_1(\zeta)+{}^{h_k}\!/\!_2 DU_{k-1}^{j,1}\phi_2(\zeta)+U_k^{j,1}\phi_3(\zeta)+{}^{h_k}\!/\!_2 D_k^{j,1}\phi_4(\zeta)\big),$ <br><br> $U_{k-1\|k}^{j,0\|1}:=(U_{k-1\|k,1}^{j,0\|1},U_{k-1\|k,2}^{j,0\|1}),\ \ DU_{k-1\|k}^{j,0\|1}:=(DU_{k-1\|k,1}^{j,0\|1},DU_{k-1\|k,2}^{j,0\|1})\in\mathbb{R}^2$ |

Table 2.1: $dG(q)$ functions for different ansatz in space, $U_k^j=U^j(x_k)$, $m=n-1,n$. For the definition of basis functions $\phi_\ell$, see Subsections 2.1.1 and 2.1.2.

**Lemma 2.3.3.2 (Distributional time derivative of $dG(q)$ function).** Due to the properties of discrete space $\mathcal{Q}_q$, if the piecewise time derivative is defined by $U_\tau|_{I_j}:=\dot{U}^j$ then the distributional time derivative may be seen as

$$\dot{U}=U_\tau+\sum_{j=1}^N[U]^{j-1}\delta_{t_{j-1}},\tag{2.16}$$

where $\delta_{t_j}$ denotes the delta distribution supported in $t_j$.

**Proof.** The proof follows by using the following representation of $dG(q)$ function $U\in\mathcal{Q}^q$,

$$U(t)=\sum_{j=1}^{N-1}U^j(t)\mathbb{H}(t-t_{j-1})\mathbb{H}(t_j-t)+U^N(t)\mathbb{H}(t-t_{N-1})-U^{0-}\mathbb{H}(t_0-t).$$

According to the definition of the distributional derivative, cf. Definition 1.2.3.4 and results

from Lemma 2.3.3.1, we have

$$\dot{U}(t) = \sum_{j=1}^{N-1} \dot{U}^j(t)\mathbb{H}(t - t_{j-1})\mathbb{H}(t_j - t) + \dot{U}^N(t)\mathbb{H}(t - t_{N-1})$$

$$+ \sum_{j=1}^{N-1} \left( U^{j-1+}\delta_{t_{j-1}}\mathbb{H}(t_j - t) - U^{j-}\mathbb{H}(t - t_{j-1})\delta_{t_j} \right) + U^{N-1+}\delta_{t_{N-1}} - U^{0-}\delta_{t_0}$$

$$= U_\tau + \sum_{j=1}^{N-1} \left( U^{j-1+}\delta_{t_{j-1}} - U^{j-}\delta_{t_j} \right) + U^{N-1+}\delta_{t_{N-1}} - U^{0-}\delta_{t_0}$$

$$= U_\tau + \sum_{j=0}^{N-1} [U]^j \delta_{t_j}.$$

This concludes the proof of lemma.                                              $\square$

To obtain the weak formulation in terms of weak bilinear form $\mathcal{B}$, we multiply the initial equation (1.28) with test function $V \in \mathcal{W}_q$ with respect to $\mathcal{H}$ scalar product from Definition 1.3.0.3. An integration with respect to each interval $I_j$ and sum over all $j = 1, \ldots, N$ yield

$$\sum_{j=1}^N \int_{I_j} \langle \dot{u} ; V \rangle_{\mathcal{H}} dt + \sum_{j=1}^N \int_{I_j} \langle \mathcal{A}u ; V \rangle_{\mathcal{H}} dt = \sum_{j=1}^N \int_{I_j} \langle F ; V \rangle_{\mathcal{H}} dt. \tag{2.17}$$

For discrete function $U \in \mathcal{Q}_q$ instead of smooth $u$ in (2.17), the time derivative $\dot{U}$ is to be interpreted in distributional sense, see (2.16). By extending the equation (2.17) in terms of $\mathcal{H}$ scalar product, we may formulate the discrete problem: Find $U = (U_1, U_2) \in \mathcal{Q}_q$ such that $U$ solves (2.5) where

$$\mathcal{B}(U, V) := \sum_{j=1}^N \int_{I_j} a(U_{1,\tau}; V_1)dt + \sum_{j=1}^N \int_{I_j} (U_{2,\tau}; V_2)dt - \sum_{j=1}^N \int_{I_j} a(U_2; V_1)dt + \sum_{j=1}^N \int_{I_j} a(U_1; V_2)dt$$

$$+ \varepsilon \sum_{j=1}^N \int_{I_j} a(U_2; V_2)dt + \sum_{j=1}^N a([U_1]^{j-1}; V_1^{j-1+}) + \sum_{j=1}^N ([U_2]^{j-1}; V_2^{j-1+}), \tag{2.18}$$

and the functional $\mathscr{L}$ is defined by

$$\mathscr{L}(V) := \sum_{j=1}^N \int_{I_j} (f; V_2)dt, \tag{2.19}$$

for all $V := (V_1, V_2) \in \mathcal{Q}_q$ and given initial solution $U^{0-} = \Pi u_0$.

**Definition 2.3.3.3 (Affine approximation).** Given $U \in \mathcal{Q}_q$, we define a globally continuous and piecewise affine function $\widetilde{U}$ with respect to triangulation $\mathscr{T}$ of time interval $[0, T]$ such that

$$\widetilde{U}(t, x)|_{I_j} := \frac{t - t_{j-1}}{k_j} U^j(t_j, x) + \frac{t_j - t}{k_j} U^{j-1}(t_{j-1}, x)$$

for all $(t, x) \in I_j \times \Omega$ and $j = 1, \ldots, n$.                                              $\square$

On the basis of Definition 2.3.3.3, we notice that in case $q=0$, for all $U, V \in \mathcal{Q}_0$, the bilinear form $\mathcal{B}$ (2.18) is equivalent to

$$\widetilde{\mathcal{B}}(U, V) := \sum_{j=1}^{N} \int_{I_j} a(\dot{\tilde{U}}_1; V_1)dt + \sum_{j=1}^{N} \int_{I_j} (\dot{\tilde{U}}_2; V_2)dt - \sum_{j=1}^{N} \int_{I_j} a(U_2; V_1)dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} a(U_1; V_2)dt + \varepsilon \sum_{j=1}^{N} \int_{I_j} a(U_2; V_2)dt. \tag{2.20}$$

This follows from the fact that $\left\langle [U]^j ; V^{j-1+} \right\rangle_{\mathcal{H}} dt = k_j \left\langle \dot{\tilde{U}} ; V^j \right\rangle_{\mathcal{H}} = \int_{I_j} \left\langle \dot{\tilde{U}} ; V \right\rangle_{\mathcal{H}} dt$ if $U, V \in \mathcal{Q}_0^j$.

**Remark 2.3.3.1.** The bilinear form $\mathcal{B}$, cf. (2.18), coincides with the bilinear form $\tilde{\mathcal{B}}$ (2.18) only on a space of test functions (piecewise constant in time). The bilinear form $\mathcal{B}$ is also consistent with a strong formulation unlike $\tilde{\mathcal{B}}$, cf. Definition 2.3.1.1. $\qquad\square$

Consequently, we define for all $v \in L^2(0, T; \mathcal{H})$

$$\widetilde{Res}(v) := \mathscr{L}(v) - \widetilde{\mathcal{B}}(U, v). \tag{2.21}$$

where the linear functional $\mathscr{L}$ is defined as in (2.19). Furthermore, the identity (2.21) can also be written as

$$\widetilde{Res}(v) = \sum_{j=1}^{N} \int_{I_j} a(\dot{e}_1; v_1)dt + \sum_{j=1}^{N} \int_{I_j} (\dot{e}_2; v_2)dt - \sum_{j=1}^{N} \int_{I_j} a(e_2; v_1)dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} a(e_1; v_2)dt + \varepsilon \sum_{j=1}^{N} \int_{I_j} a(e_2; v_2)dt. \tag{2.22}$$

Notice that this residual retains the orthogonality property, namely for $V \in \mathcal{Q}_0$

$$\widetilde{Res}(V) = \mathscr{L}(V) - \widetilde{\mathcal{B}}(V) = \mathscr{L}(V) - \mathcal{B}(V) = Res(V) = 0. \tag{2.23}$$

**Lemma 2.3.3.3.** For $U \in \mathcal{Q}_q$, the following two identities are valid

$$\int_{I_j} \left\langle U_\tau ; U \right\rangle_{\mathcal{H}} dt = \frac{1}{2} \|U^{j-}\|_{\mathcal{H}}^2 - \frac{1}{2} \|U^{j-1+}\|_{\mathcal{H}}^2 \quad \text{and}$$

$$\left\langle [U]^{j-1} ; U^{j-1+} \right\rangle_{\mathcal{H}} = \frac{1}{2} \|U^{j-1+}\|_{\mathcal{H}}^2 + \frac{1}{2} \|[U]^{j-1}\|_{\mathcal{H}}^2 - \frac{1}{2} \|U^{j-1-}\|_{\mathcal{H}}^2.$$

**Proof.** For the first identity we have

$$\int_{I_j} \left\langle U_\tau ; U \right\rangle_{\mathcal{H}} dt = \int_{I_j} a(U_{1,\tau}; U_1)dt + \int_{I_j} (U_{2,\tau}; U_2)dt = \frac{1}{2} \int_{I_j} \frac{\partial}{\partial t} \|U_1\|_{H^1(\Omega)}^2 dt + \int_{I_j} \frac{\partial}{\partial t} \|U_2\|_{L^2(\Omega)}^2 dt$$

$$= \|U_1^{j-}\|_{H^1(\Omega)}^2 - \|U_1^{j-1+}\|_{H^1(\Omega)}^2 + \|U_2^{j-}\|_{L^2(\Omega)}^2 - \|U_2^{j-1+}\|_{L^2(\Omega)}^2.$$

In case of the second identity,

$$\left\langle [U]^{j-1} ; U^{j-1+} \right\rangle_{\mathcal{H}} = \frac{1}{2} \left\langle U^{j-1+} - U^{j-1-} ; U^{j-1+} \right\rangle_{\mathcal{H}} + \frac{1}{2} \left\langle U^{j-1+} - U^{j-1+} ; U^{j-1+} - U^{j-1-} \right\rangle_{\mathcal{H}}$$

$$+ \frac{1}{2} \left\langle U^{j-1+} - U^{j-1+} ; U^{j-1-} \right\rangle_{\mathcal{H}}$$

$$= \frac{1}{2} \|U^{j-1+}\|_{\mathcal{H}}^2 - \frac{1}{2} \left\langle U^{j-1-} ; U^{j-1+} \right\rangle_{\mathcal{H}} + \frac{1}{2} \|[U]^{j-1}\|_{\mathcal{H}}^2$$

$$+ \frac{1}{2} \left\langle U^{j-1+} ; U^{j-1-} \right\rangle_{\mathcal{H}} - \frac{1}{2} \|U^{j-1-}\|_{\mathcal{H}}^2.$$

Owing to the symmetry of the $\mathcal{H}$ scalar product, we conclude the proof of Lemma. $\qquad\square$

Before we establish the stability estimate for $U$ in terms of initial data $U^{0-}$, where $U$ is the solution of (2.5), we introduce the discrete operator $\mathcal{K}_h$, i.e. a discrete variant of operator $\mathcal{K}$ from Definition 1.3.0.5.

**Definition 2.3.3.4 (Discrete operator $\mathcal{K}_h$).** For given $f \in H^{-1}(\Omega)$, there is a unique solution $U \in \mathcal{S}$ of the weak problem

$$a(U,V) = (f;V) \quad \text{for all} \quad V \in \mathcal{S}, \tag{2.25}$$

where $\mathcal{S}$ denotes a finite dimensional subspace of $H_D^1(\Omega)$. Hence we may define the operator $\mathcal{K}_h : H^{-1}(\Omega) \to \mathcal{S}$ which maps $f \in H^{-1}(\Omega)$ onto the unique solution $U = \mathcal{K}_h f \in \mathcal{S}$ of (2.25). $\square$

**Lemma 2.3.3.4.** Since $\mathcal{S} \subset H_D^1 \subset L^2(\Omega) \subset H^{-1}(\Omega)$, we may consider the restricted operator $\mathcal{K}_h : L^2(\Omega) \to \mathcal{S}$ from Definition 2.3.3.4. This restriction satisfies

  a) $\mathcal{K}_h$ is self-adjoint, positive semi-definite and compact operator from $L^2(\Omega)$ onto $L^2(\Omega)$,

  b) $\mathcal{K}_h^{1/2} : L^2(\Omega) \to L^2(\Omega)$ is well-defined, positive, self-adjoint, compact operator.

  c) $\mathcal{K}_h^{1/2}$ can be extended to an operator $\mathcal{K}_h^{1/2} : H^{-1}(\Omega) \to L^2(\Omega)$. There holds

$$\|\mathcal{K}_h^{1/2} f\|_{L^2(\Omega)} = \|f\|_{H^{-1}(\Omega)} \quad \text{for all} \quad f \in H^{-1}(\Omega), \tag{2.26}$$

   i.e. $\mathcal{K}_h^{1/2}$ is an isometry.

  d) The restriction $\mathcal{K}_h : \mathcal{S} \to \mathcal{S}$ is injective and hence invertible. In particular, there exists $\mathcal{K}_h^{-1} : \mathcal{S} \to \mathcal{S}$ with

$$(\mathcal{K}_h^{-1} U; V) := a(U; V) \quad \text{for all} \quad V \in \mathcal{S}. \tag{2.27}$$

**Proof.** The proof for a),b),c) follows analogously as the proof of Lemma 1.3.0.2 with $\mathcal{K}_h$ instead of $\mathcal{K}$.
To see that $\mathcal{K}_h : \mathcal{S} \to \mathcal{S}$ is injective, let $U \in \mathcal{S}$ satisfy $\mathcal{K}_h U = 0$. By definition (2.25), there holds $0 = a(\mathcal{K}_h U; V) = (U; V)$ for all $V \in \mathcal{S}$. If $V = U$ we have $\|U\|_{L^2(\Omega)}^2 = 0$ whence $U = 0$. $\square$

**Remark 2.3.3.2.** The operator $\mathcal{K}_h^{-1}$ is often referred to as discrete Laplacian, c.f. ERIKSSON-JOHNSON [23]. $\square$

**Remark 2.3.3.3.** The operator $\mathcal{K}_h$ can be seen as Galerkin projection of the operator $\mathcal{K}$ from Definition 1.3.0.5, i.e. $\mathcal{K}_h = \mathcal{G}\mathcal{K}$ where $\mathcal{G}$ is defined in Definition 3.1.0.3. This results from the following identity. Namely, if $U$ is the discrete solution of problem (2.25) and $u$ a continuous solution of (1.14), then for each $V \in \mathcal{S}$

$$a(\mathcal{K}_h f; V) = a(U; V) = (f; V) = a(u; V) = a(\mathcal{K} f; V) = a(\mathcal{G}\mathcal{K} f; V). \tag{2.28}$$

Meaning, the discrete and continuous operator $\mathcal{K}_h$ and $\mathcal{K}$, respectively, coincide on the space of discrete functions. $\square$

**Lemma 2.3.3.5 (Stability of $dG$ solution).** Let $U$ be the $dG(q), q = 0, 1$ solution of the homogeneous problem (2.5), i.e. $f \equiv 0$ with $\mathcal{B}$ defined in (2.18). Then there holds for all $\varepsilon \geq 0$ and $1 \leq n \leq N$

$$\|U^{n-}\|_{\mathcal{H}}^2 + 2\varepsilon \sum_{j=1}^{n} \int_{I_j} \|U_2(t)\|_{H^1(\Omega)}^2 dt + \sum_{j=1}^{n} \|[U]^j\|_{\mathcal{H}}^2 = \|U^{0-}\|_{\mathcal{H}}^2. \tag{2.29a}$$

Moreover, in case when $q = 1$

$$\|U_{1,\tau}\|_{L^\infty(L^2)} \leq \|U^{0-}\|_{\mathcal{H}}, \tag{2.29b}$$

where additionally for $\varepsilon = 0$ there holds

$$\|U_{2,\tau}\|_{L^\infty(H^{-1})} \leq \|U^{0-}\|_{\mathcal{H}}. \tag{2.29c}$$

**Proof.** In order to prove the basic stability estimate (2.29a), we chose $V = U$ in the weak form (2.5), where $\mathscr{L}(U) = 0$ to obtain

$$0 = \mathcal{B}(U, U) = \sum_{j=1}^{N} \int_{I_j} a(U_{1,\tau}; U_1) dt + \sum_{j=1}^{N} \int_{I_j} (U_{2,\tau}; U_2) dt + \varepsilon \sum_{j=1}^{N} \int_{I_j} a(U_2; U_2) dt$$

$$+ \sum_{j=1}^{N} a([U_1]^{j-1}; U_1^{j-1+}) + \sum_{j=1}^{N} ([U_2]^{j-1}; U_2^{j-1+})$$

$$= \sum_{j=1}^{N} \int_{I_j} \langle U_\tau; U \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \langle [U]^{j-1}; U^{j-1+} \rangle_{\mathcal{H}} + \varepsilon \sum_{j=1}^{N} \int_{I_j} \|U_2(t)\|_{H^1(\Omega)} dt. \tag{2.30}$$

According to the results of Lemma 2.3.3.3, this is equivalent to

$$0 = \frac{1}{2} \sum_{j=1}^{N} \|U^{j-}\|_{\mathcal{H}}^2 - \frac{1}{2} \sum_{j=1}^{N} \|U_1^{j-1-}\|_{\mathcal{H}}^2 + \frac{1}{2} \sum_{j=1}^{N} \|[U_2]^{j-1}\|_{\mathcal{H}}^2 + \varepsilon \sum_{j=1}^{N} \int_{I_j} \|U_2(t)\|_{H^1(\Omega)}^2 dt.$$

Hence, we conclude the proof of (2.29a).

In case of (2.29b), let $V \in \mathcal{Q}_1$ be chosen such that

$$V(t) := \begin{cases} ((t - t_{j-1}) \mathcal{K}_h U_{1,\tau}^j, 0), & t \in I_j, \\ \mathbf{0}, & t \notin I_j \end{cases}$$

for $\mathcal{K}_h$, the discrete variant of operator $\mathcal{K}$ defined in Definition 2.3.3.4.

A substitution of $V$ into the homogeneous weak form (2.5) with $\mathcal{B}$ from (2.18) yields

$$0 = \int_{I_j} a(U_{1,\tau}; (t - t_{j-1}) \mathcal{K}_h U_{1,\tau}) dt - \int_{I_j} a(U_2; (t - t_{j-1}) \mathcal{K}_h U_{1,\tau}^j) dt$$

$$= \int_{I_j} (t - t_{j-1})(U_{1,\tau}; U_{1,\tau}) dt - \int_{I_j} (U_2(t - t_{j-1}); U_{1,\tau}) dt \tag{2.31}$$

In the second equality above we used the identity (2.25), cf. Definition 2.3.3.4.
The equation (2.31) is now equivalent to

$$\frac{k_j^2}{2}\|U_{1,\tau}^j\|_{L^2(\Omega)}^2 = \int_{I_j} (U_2(t-t_{j-1}); U_{1,\tau}) dt$$

$$\leq \|U_{1,\tau}^j\|_{L^2(\Omega)} \int_{I_j} (t-t_{j-1}) \|U_2(t)\|_{L^2(\Omega)} dt \leq \frac{k_j^2}{2}\|U_{1,\tau}^j\|_{L^2(\Omega)} \|U_2\|_{L^\infty(L^2)}.$$

The discrete function $U_1$ is piecewise affine in time and therefore

$$\|U_{1,\tau}\|_{L^\infty(L^2)}^2 \leq \|U_2\|_{L^\infty(L^2)}^2 \leq \|U^{0-}\|_{\mathcal{H}}, \tag{2.32}$$

where the first stability estimate (2.29a) was used for the second inequality.

To prove the third stability estimate (2.29c), we need to bound $\|U_{2,\tau}\|_{L^\infty(H^{-1})}$. Let $V \in \mathcal{Q}_1$ be defined as

$$V(t) := \begin{cases} (\varepsilon \mathcal{K}_h U_{2,\tau}(t-t_{j-1}), \mathcal{K}_h U_{2,\tau}(t-t_{j-1})), & t \in I_j, \\ \mathbf{0}, & t \notin I_j. \end{cases}$$

Hence, the homogeneous weak problem $\mathcal{B}(U, V) = 0$ reads

$$0 = \varepsilon \int_{I_j} (U_{1,\tau}; U_{2,\tau}(t-t_{j-1})) dt + \int_{I_j} (U_{2,\tau}; \mathcal{K}_h U_{2,\tau}(t-t_{j-1})) dt - \varepsilon \int_{I_j} (U_2; U_{2,\tau}(t-t_{j-1})) dt$$

$$+ \int_{I_j} (U_1; U_{2,\tau}(t-t_{j-1})) dt + \varepsilon \int_{I_j} (U_2; U_{2,\tau}(t-t_{j-1})) dt$$

$$= \int_{I_j} (U_{2,\tau}; \mathcal{K}_h U_{2,\tau}(t-t_{j-1})) dt - \int_{I_j} (U_1; U_{2,\tau}(t-t_{j-1})) dt - \varepsilon \int_{I_j} (U_{1,\tau}; U_{2,\tau}(t-t_{j-1})) dt.$$

On account of the definition of the identity of the dual norm $\|\cdot\|_{H^{-1}(\Omega)}$ and $\|\mathcal{K}_h^{1/2}(\cdot)\|_{L^2(\Omega)}$, cf. (2.26), by applying the Hölder inequality in space the last equation can be recast to

$$\frac{k_j^2}{2}\|U_{2,\tau}^j\|_{H^{-1}(\Omega)}^2 = \frac{k_j^2}{2}\|\mathcal{K}_h^{1/2} U_{2,\tau}^j\|_{L^2(\Omega)}^2 = \int_{I_j} (U_{2,\tau}; \mathcal{K}_h U_{2,\tau}(t-t_{j-1})) dt$$

$$= -\int_{I_j} (U_1; U_{2,\tau}(t-t_{j-1})) dt - \varepsilon \int_{I_j} (U_{1,\tau}; U_{2,\tau}(t-t_{j-1})) dt$$

$$\leq \frac{k_j^2}{2}\|U_1\|_{L^\infty(H^1)}\|U_{2,\tau}^j\|_{H^{-1}(\Omega)} + \frac{k_j^2}{2}\varepsilon\|U_{1,\tau}\|_{L^\infty(H^1)}\|U_{2,\tau}^j\|_{H^{-1}(\Omega)}.$$

From the fact that $U_{2,\tau}$ is piecewise constant in time and the stability estimate (2.29a), there holds

$$\|U_{2,\tau}\|_{L^\infty(H^{-1})} \leq \|U_1\|_{L^\infty(H^1)} + \varepsilon\|U_{1,\tau}\|_{L^\infty(H^1)} \leq \|U^{0-}\|_{\mathcal{H}} + \varepsilon\|U_{1,\tau}\|_{L^\infty(H^1)}. \tag{2.33}$$

Obviously, we need to estimate term $\varepsilon\|U_{1,\tau}\|_{L^\infty(H^1)}$ such that the RHS of the estimate (2.33) depends only on $\|U^{0-}\|_{\mathcal{H}}$. We will show that this is in general possible only in case when $\varepsilon = 0$.

First, let $V \in \mathcal{Q}_1$ be a test function

$$V(t) := \begin{cases} (\varepsilon U_{1,\tau}(t-t_{j-1}), U_{1,\tau}(t-t_{j-1})), & t \in I_j \\ \mathbf{0}, & t \notin I_j, \end{cases}$$

such that the weak form $\mathcal{B}(U,V) = 0$ reads

$$0 = \varepsilon \int_{I_j} a(U_{1,\tau}; U_{1,\tau}(t-t_{j-1}))dt + \int_{I_j} (U_{2,\tau}; U_{1,\tau}(t-t_{j-1}))dt - \varepsilon \int_{I_j} a(U_2; U_{1,\tau}(t-t_{j-1}))dt$$
$$+ \int_{I_j} a(U_1; U_{1,\tau}(t-t_{j-1}))dt + \varepsilon \int_{I_j} a(U_2; U_{1,\tau}(t-t_{j-1}))dt.$$

This simplifies to

$$\varepsilon \int_{I_j} \|U_{1,\tau}\|^2_{H^1(\Omega)}(t-t_{j-1})dt = -\int_{I_j} (U_{1,\tau}; U_{2,\tau})(t-t_{j-1})dt - \int_{I_j} a(U_1; U_{1,\tau})(t-t_{j-1})dt. \qquad (2.34)$$

The discrete function $U$ is piecewise affine in time and therefore the RHS of the equality (2.34) can be estimated such that

$$\varepsilon \|U_{1,\tau}\|^2_{L^\infty(H^1)} \leq \left( \|U_{2,\tau}\|_{L^\infty(H^{-1})} + \|U_1\|_{L^\infty(H^1)} \right) \|U_{1,\tau}\|_{L^\infty(H^1)} \qquad (2.35)$$

Hence, by means of (2.29a)

$$\varepsilon \|U_{1,\tau}\|_{L^\infty(H^1)} \leq \|U_{2,\tau}\|_{L^\infty(H^{-1})} + \|U^{0-}\|_{\mathcal{H}}. \qquad (2.36)$$

Obviously, from (2.36) and (2.32) we can not obtain the estimate for $\|U_{2,\tau}\|_{L^\infty(H^{-1})}$ when $\varepsilon > 0$. On that basis, we conclude from (2.32) that for $\varepsilon = 0$ the following estimate is valid

$$\|U_{2,\tau}\|^2_{L^\infty(H^{-1})} \leq \|U^{0-}\|_{\mathcal{H}}. \qquad (2.37)$$

This concludes the proof of the estimate (2.29c) and lemma also. $\qquad \square$

**Remark 2.3.3.4 (Strong energy dissipation by $dG(q)$ method).** Due to the definition of the energy, see Definition 1.3.0.4, it is obvious that in case of $dG(q)$, $q = 0, 1$ time approximation, the energy of the discrete solution dissipates more than the energy of the continuous solution. For the continuous model, there holds

$$\|u(t)\|^2_{\mathcal{H}} = \|u_0\|^2_{\mathcal{H}} - 2\varepsilon \int_0^T \|u_2(t)\|^2_{H^1(\Omega)}dt$$

according to Lemma 1.3.0.1 and Remark 1.3.0.5.
Contrary to that, the discrete solution satisfies

$$\|U^{j-}\|^2_{\mathcal{H}} = \|U^{j-1-}\|^2_{\mathcal{H}} - \|[U]^{j-1}\|^2_{\mathcal{H}} - 2\varepsilon \int_{I_j} \|U_2(t)\|^2_{H^1(\Omega)}dt. \qquad (2.38)$$

For the proof of (2.38), we proceed similarly as in case of the Lemma 2.3.3.5, but substitute $V = U$ in the $j$th contribution of the homogeneous weak form, cf. Remark 2.3.1.1. The equation (2.38) clearly indicates that $\|U^{j-}\|_{\mathcal{H}}$ decreases compared to $\|U^{j-1-}\|_{\mathcal{H}}$. $\qquad \square$

### 2.3.4  Continuous Galerkin time discretisation

In the following we consider another concept of time discretisation by use of affine and globally continuous functions, abb. $cG(1)$. The discretisation in space remains as in the case of the $dG$ method, i.e. by using the linear or cubic splines.

Note that in case of a continuous Galerkin approximation in time, the test functions are one degree lower in time to account for the fact that the discrete solution is fixed a priori by continuity at each time node, i.e. for $t = t_j$ where $j = 0, \dots, N$. Such a method, where test and trial functions do not belong to the same discrete space, is called *Petrov-Galerkin* method. Petrov-Galerkin methods also comprise the inverse choice of trial and test functions, namely when $U \in \mathcal{W}$ is one degree lower than the test function $V \in \mathcal{Q}$. For more details we refer to HARTMANN [37], even if this inverse choice will not be discussed within this work.

On the basis of the definition of the continuous Galerkin method, where the approximative solution is globally continuous, and the fact that the spatial refining method assumes arbitrary refinement on each time step, we restrict here to so called hierarchical meshes in space. These meshes allow the continuity condition to be satisfied, i.e. $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j = 1, \dots, N$.

**Definition 2.3.4.1 ($cG(1)$ space, full discretisation).** Let $\mathcal{Q}_c$ be denoted in the following as $cG(1)$ space, i.e. space of piecewise affine globally continuous functions in time with coefficients in $\mathcal{S}^j$ on each time interval $I_j \in \mathcal{T}$, cf. Definition 2.0.0.8. Then,

$$\mathcal{Q}_c := \left\{ U : Q \to \mathbb{R}^2 \mid U(t, \cdot)|_{I_j} \in \mathcal{S}^j \times \mathcal{S}^j, \ U(\cdot, x)|_{I_j} \in \mathcal{P}_1(I_j), \ I_j \in \mathcal{T} \right\} \subseteq \mathcal{C}(0, T; \mathcal{H}). \quad (2.39)$$

Accordingly, let

$$\mathcal{W}_c := \left\{ V : Q \to \mathbb{R}^2 \mid V(t, \cdot)|_{I_j} \in \mathcal{S}^j \times \mathcal{S}^j, \ V(\cdot, x)|_{I_j} \in \mathcal{P}_0(I_j), \ I_j \in \mathcal{T} \right\} \in H^1(\mathcal{T}; \mathcal{H}) \quad (2.40)$$

denote the space of piecewise constant functions on $[0, T]$ with respect to $\mathcal{T}$. This space will be further referred to as the test space. To abbreviate the notation, we write in the following $\mathcal{Q}_c|_{I_j} =: \mathcal{Q}_c^j$ etc.

The form of piecewise continuous discrete function in time, abb. $cG(1)$ function, $U := (U_1, U_2)$, is presented in Table 2.2 below.                                                                                        □

Hereby, from the Definition 2.3.4.1, the discrete problem reads: Find $U = (U_1, U_2) \in \mathcal{Q}_c$ such that $U$ solves (2.5) where

$$\mathcal{B}(U, V) := \int_0^T a(\dot{U}_1; V_1) dt + \int_0^T (\dot{U}_2; V_2) dt - \int_0^T a(U_2; V_1) dt$$
$$+ \int_0^T a(U_1; V_2) dt + \varepsilon \int_0^T a(U_2; V_2) dt, \quad (2.41)$$

and the linear functional $\mathscr{L}$ is defined by

$$\mathscr{L}(V) := \int_0^T (f, V_2) dt, \quad (2.42)$$

for all $V = (V_1, V_2) \in \mathcal{W}_c$ and given initial solution $U(0) = \Pi u_0$.

**Remark 2.3.4.1.** Notice that the bilinear form $\mathcal{B}$ (2.41) is consistent with the strong formulation, cf. Remark 2.3.1.1.                                                                          □

---

**$cG(1)$ functions with space discretisation by linear splines**

$$U(t,x)|_{I_j}:=U^j(t,x)=\sum_{k=1}^{m}\big((t-t_{j-1})/k_j U_k^j+(t_j-t)/k_j U_k^{j-1}\big)\phi_k(x),\ \ U_k^{j|j-1}=(U_{k,1}^{j|j-1},U_{k,2}^{j|j-1})\in\mathbb{R}^2$$

---

**$cG(1)$ functions with space discretisation by cubic splines**

$$U(t,x(\zeta))|_{I_j\times T_k}:=U_k^j(t,x)=(t-t_{j-1})/k_j\big(U_{k-1}^j\phi_1(\zeta)+h_k/2 DU_{k-1}^j\phi_2(\zeta)+U_k^j\phi_3(\zeta)+h_k/2 DU_k^j\phi_4(\zeta)\big)$$
$$+(t_j-t)/k_j\big(U_{k-1}^{j-1}\phi_1(\zeta)+h_k/2 DU_{k-1}^{j-1}\phi_2(\zeta)+U_k^{j-1}\phi_3(\zeta)+h_k/2 DU_k^{j-1}(t)\phi_4(\zeta)\big),$$

$$U_{k|k-1}^{j|j-1}=(U_{k|k-1,1}^{j|j-1},U_{k|k-1,2}^{j|j-1}),\ \ DU_{k|k-1}^{j|j-1}=(DU_{k|k-1,1}^{j|j-1},DU_{k|k-1,2}^{j|j-1})\in\mathbb{R}^2$$

---

Table 2.2: $cG(1)$ functions for different ansatz in space, $\mathcal{S}^{j-1}\subseteq S^j$, $U_k^j=U(t_j,x_k)$, $m=n-1,n$. For the definition of basis functions $\phi_\ell$, see Subsection 2.1.1 and Subsection 2.1.2.

**Lemma 2.3.4.1 (Stability of $cG(1)$ solution).** Let $U$ be the $cG(1)$ solution of the homogeneous problem (2.5), i.e. $f\equiv0$ with $\mathcal{B}$ defined in (2.41). Then, there holds for all $\varepsilon\geq0$ and $1\leq n\leq N$

$$\|U(t_n)\|_{\mathcal{H}}^2+2\varepsilon\int_0^{t_n}\|\dot{U}_1(t)\|_{H^1(\Omega)}^2 dt=\|U(0)\|_{\mathcal{H}}^2, \tag{2.43a}$$

$$\|\dot{U}_1\|_{L^\infty(L^2)}\leq\|U(0)\|_{\mathcal{H}}. \tag{2.43b}$$

Moreover, if $\varepsilon=0$, then

$$\|\dot{U}_2\|_{L^\infty(H^{-1})}\leq\|U(0)\|_{\mathcal{H}}. \tag{2.43c}$$

**Proof.** We start with the proof of the first stability estimate (2.43a).
Given the bilinear form $\mathcal{B}$ defined in (2.41) and a test function $V=(\mathcal{K}_h\dot{U}_2-\varepsilon\dot{U}_1,-\dot{U}_1)$ for $\mathcal{K}_h$, discrete operator from Definition 2.3.3.4, then the homogeneous weak problem (2.5) reads

$$0=\mathcal{B}(U,V)=\int_0^T a(\dot{U}_1;\mathcal{K}_h\dot{U}_2)dt-\varepsilon\int_0^T a(\dot{U}_1;\dot{U}_1)dt-\int_0^T(\dot{U}_2;\dot{U}_1)dt-\int_0^T a(U_2;\mathcal{K}_h\dot{U}_2)dt$$
$$+\varepsilon\int_0^T a(U_2;\dot{U}_1)dt-\int_0^T a(U_1;\dot{U}_1)dt-\varepsilon\int_0^T a(U_2;\dot{U}_1)dt$$
$$=\int_0^T(\dot{U}_1;\dot{U}_2)dt-\varepsilon\int_0^T a(\dot{U}_1;\dot{U}_1)dt-\int_0^T(\dot{U}_2;\dot{U}_1)dt-\int_0^T(U_2;\dot{U}_2)dt$$
$$+\varepsilon\int_0^T a(U_2;\dot{U}_1)dt-\int_0^T a(U_1;\dot{U}_1)dt-\varepsilon\int_0^T a(U_2;\dot{U}_1)dt.$$

This simplifies to

$$\frac{1}{2}\int_0^T\frac{\partial}{\partial t}\|U_1\|_{H^1(\Omega)}^2 dt+\frac{1}{2}\int_0^T\frac{\partial}{\partial t}\|U_2\|_{L^2(\Omega)}^2 dt+\varepsilon\int_0^T\|\dot{U}_1\|_{H^1(\Omega)}^2 dt=0.$$

By means of the main theorem in calculus we conclude the proof of (2.43a).

In case of (2.43b), we choose a test function $V = (\mathcal{K}_h \dot{U}_1, 0)$ in (2.5) in order to bound $\dot{U}_1$ in terms of $U_2$. This leads to

$$0 = \mathcal{B}(U, V) = \int_0^T a(\dot{U}_1; \mathcal{K}_h \dot{U}_1)dt - \int_0^T a(U_2; \mathcal{K}_h \dot{U}_1)dt = \int_0^T (\dot{U}_1; \dot{U}_1)dt - \int_0^T (U_2; \dot{U}_1)dt.$$

Furthermore, the discrete function $U_1$ is piecewise affine in time. This implies

$$\|\dot{U}_1\|_{L^\infty(L^2)} \leq \|U_2\|_{L^\infty(L^2)} \leq \|U(0)\|_{\mathcal{H}} \tag{2.44}$$

where we used the same arguments as in the proof of (2.32) and stability estimate (2.43a).

We continue with the proof of (2.43c). To bound $\|\dot{U}_2\|_{L^\infty(L^2)}$ we choose $V = (\varepsilon\mathcal{K}_h\dot{U}_2, \mathcal{K}_h\dot{U}_2)$. For this choice of test function, the homogeneous discrete problem (2.5) reads

$$0 = \varepsilon\int_0^T a(\dot{U}_1, \mathcal{K}_h\dot{U}_2)dt + \int_0^T (\dot{U}_2, \mathcal{K}_h\dot{U}_2)dt - \varepsilon\int_0^T a(U_2; \mathcal{K}_h\dot{U}_2)dt + \int_0^T a(U_1; \mathcal{K}_h\dot{U}_2)dt + \varepsilon\int_0^T a(U_2; \mathcal{K}_h\dot{U}_2)dt$$

$$= \varepsilon\int_0^T (\dot{U}_1, \dot{U}_2)dt + \int_0^T (\dot{U}_2, \mathcal{K}_h\dot{U}_2)dt - \varepsilon\int_0^T (U_2; \dot{U}_2)dt + \int_0^T (U_1; \dot{U}_2)dt + \varepsilon\int_0^T (U_2; \dot{U}_2)dt.$$

Moreover, from (2.26), cf. Lemma 2.3.3.4 we have

$$\int_0^T \|\dot{U}_2\|_{H^{-1}(\Omega)}^2 dt = \int_0^T \|\mathcal{K}_h^{1/2}\dot{U}_2\|_{L^2(\Omega)}^2 = \varepsilon\int_0^T (\dot{U}_1; \dot{U}_2)dt + \int_0^T (U_1; \dot{U}_2)dt.$$

The Hölder inequality in time and space leads to

$$\|\dot{U}_2\|_{L^\infty(H^{-1})} \leq \varepsilon\|\dot{U}_1\|_{L^\infty(H^1)}dt + \|U_1\|_{L^\infty(H^1)}. \tag{2.45}$$

Obviously, to apply the stability estimate (2.43a) to the RHS of (2.45) we need to assume that $\varepsilon = 0$ similar as in (2.29c). Then

$$\|\dot{U}_2\|_{L^\infty(H^{-1})} \leq \|U_1\|_{L^\infty(H^1)} \leq \|U(0)\|_{\mathcal{H}} \tag{2.46}$$

This concludes the proof.                                                                              $\square$

**Remark 2.3.4.2 (Energy conservation by $cG(1)$ method).** Due to the definition of the energy, see Definition 1.3.0.4 and the stability estimate (2.43a) it is obvious that in case $\varepsilon = 0$, $cG(1)$ time approximation conserves the energy.                                                   $\square$

## 2.3.5   Method of Lines

The method of lines approach applies to the semi-discretisation in space only, where time remains continuous. The discretisation in space follows by using the already known approximation methods such as linear or cubic splines. Here we also distinguish between the trial and test space. On account to the fact that time is continuous, we pass on the usual notation valid in case of Galerkin methods which comprises the indexing $j$ related to triangulation of time interval.

**Definition 2.3.5.1 (*MoL* space, semi-discretisation).** Let $\mathcal{S}$ be a finite dimensional subspace of $H_D^1(\Omega)$ which consists of linear or cubic splines, see (2.1). We denote by $\mathcal{Q}_s$ a *MoL* space, i.e. a space of functions globally continuous in time with coefficients in $\mathcal{S}$, i.e.

$$\mathcal{Q}_s := \left\{ U : Q \to \mathbb{R}^2 \mid U(t, \cdot) \in \mathcal{S} \times \mathcal{S}, \ U(\cdot, x) \in \mathcal{C}^1(0, T) \right\} \subseteq H^1(0, T; \mathcal{H}). \tag{2.47}$$

Let also $W_s$ be the appropriate test space

$$\mathcal{W}_s := \left\{ V : Q \to \mathbb{R}^2 \mid V(t, \cdot) \in \mathcal{S} \times \mathcal{S}, \ V(\cdot, x) \in L^2(0, T) \right\} \subseteq L^\infty(0, T; \mathcal{H}). \tag{2.48}$$

Discrete functions belonging to the discrete space $\mathcal{Q}_s$ are called the semi-discrete or *MoL* functions and their structure is presented in Table 2.3 below. □

| ***MoL* functions with space discretisation by linear splines, $t \in [0, T]$** |
|:---:|
| $U(t, x) := \sum_{k=1}^m U_k(t)\phi_k(x), \ U_k(t) = (U_{k,1}(t), U_{k,2}(t)) \in \mathbb{R}^2$ |

| ***MoL* functions with space discretisation by cubic splines, $t \in [0, T]$** |
|:---:|
| $U(t, x(\zeta))\|_{T_k} := U_k(t, x(\zeta)) = U_{k-1}(t)\phi_1(\zeta) + {}^{h_k}\!/{}_2 DU_{k-1}(t)\phi_2(\zeta) + U_k(t)\phi_3(\zeta) + {}^{h_k}\!/{}_2 DU_k(t)\phi_4(\zeta),$ $U_{k-1\|k} = (U_{k-1\|k,1}, U_{k-1\|k,2}), \ DU_{k-1\|k} = (DU_{k-1\|k,1}, DU_{k-1\|k,2}) \in \mathbb{R}^2$ |

Table 2.3: Semi-discrete functions for different ansatz in space, $U_k(t) = U(t, x_k)$, $m = n-1, n$. For the definition of basis functions $\phi_\ell$, see Subsection 2.1.1 and Subsection 2.1.2.

The discrete problem reads: Find $U = (U_1, U_2) \in \mathcal{Q}_s$ such that $U$ solves the discrete problem (2.5) where

$$\mathcal{B}(U, V) := a(\dot{U}_1; V_1) + (\dot{U}_2; V_2) - a(U_2; V_1) + a(U_1 + \varepsilon U_2; V_2) \tag{2.49}$$

and the linear functional $\mathscr{L}$ is defined by

$$\mathscr{L}(V) := (f(t), V_2), \tag{2.50}$$

for all $V = (V_1, V_2) \in \mathcal{W}_s$ and $U := U(t)$ and $V := V(t)$ for each $t \in [0, T]$.
Additionally, $U$ is imposed to satisfy the initial condition $U(0, x) = \Pi u_0(x)$ for all $x \in \Omega$.

**Remark 2.3.5.1.** The bilinear form $\mathcal{B}$ from (2.49) is consistent with a strong formulation, cf. Remark 2.3.1.1. □

## 2.4 Existence and uniqueness of the discrete solution

Within this section we present a proof for the uniqueness of the discrete solution.
In case of the Galerkin time approximation, the idea is to show that for each time interval $I_j$, the reduced homogeneous equation (2.6) has only trivial solution $U^j \equiv 0$ provided that the

solution from the previous time interval is zero, i.e. $U^{j-1} \equiv 0$. Existence of the discrete solution follows then directly from the uniqueness owing to the fact that the weak problem (2.5) is linear, discrete and therefore finite-dimensional problem. In case of the method of lines, the existence and uniqueness follow directly from the matrix formulation. Here the discrete solution is the solution of the system of ordinary differential equations (ODE system).

The point of departure is the discretisation in space, namely, the proofs will be derived by considering the spatial discretisation first, and thereon combining the same with time discretisation. For the sake of clarity, we adopt the notation i.e. the abbreviations for various fully discrete problems from Table 2.4.

|  | linear splines in space | cubic splines in space |
|---|---|---|
| discontinuous Galerkin in time | $dG(q) \otimes \mathcal{P}_1$ | $dG(q) \otimes \mathcal{C}^1$ |
| continuous Galerkin in time | $cG(1) \otimes \mathcal{P}_1$ | $cG(1) \otimes \mathcal{C}^1$ |
| method of lines | $MoL \otimes \mathcal{P}_1$ | $MoL \otimes \mathcal{C}^1$ |

Table 2.4: Abbreviations for different time-space discretisation.

## 2.4.1   Linear splines ($\mathcal{P}_1$)

Recall the definition of the discrete space $\mathcal{S}$ with respect to $\Omega$ given in (2.1) and definition of basis functions $\{\phi_\ell\}_{\ell=0}^n$ from Subsection 2.1.1. Thereafter, one may introduce the stiffness matrix

$$S, \text{ a } (n+1) \times (n+1) \text{ matrix with the entries } S_{k,\ell} := a(\phi_k; \phi_\ell), \qquad (2.51)$$

and the mass matrix

$$M, \text{ a } (n+1) \times (n+1) \text{ matrix with the entries } M_{k,\ell} := (\phi_k; \phi_\ell). \qquad (2.52)$$

Note that these matrices are symmetric, tridiagonal and positive definite and therefore invertible.

In Subsection 2.1.1 we restricted the number of the basic functions $\phi_\ell(x)$ with respect to spatial boundary conditions i.e. Dirichlet boundary condition in $x = 0$. Therefore, we left out the definition of $\phi_0$. Here however, we consider rather the general formulation where the Dirichlet boundary conditions are not embedded in the definition and dimension of the vector formulations of the discrete functions as well as of the stiffness and mass matrices. The incorporation of the Dirichlet boundary conditions follows afterwards and will be explained in detail in Subsection 6.1.3.

In the following, we introduce a vector notation needed for the definition of the vector variant of the discrete function. Its structure depends on the time discretisation. Recall the definitions of discrete functions from Tables 2.1– 2.3, case linear splines in space.

For each $j = 1, \ldots, N$ or $t \in [0, T]$ we may define $\mathbb{U}^j$ and $\mathbb{U}(t)$, respectively, such that their structure for particular time discretisation method corresponds to the one from Table 2.5.

| $dG(0), cG(1)$ | $\mathbb{U}^j := \begin{bmatrix} \mathbb{U}_1^j \\ \mathbb{U}_2^j \end{bmatrix}$ | $\mathbb{U}_{1\mid2}^j := (U_{0,1\mid2}^j, \ldots, U_{n,1\mid2}^j)^T$ |
|:---:|:---:|:---:|
| $dG(1)$ | $\mathbb{U}^j := \begin{bmatrix} \mathbb{U}_1^{j,0} \\ \mathbb{U}_1^{j,1} \\ \mathbb{U}_2^{j,0} \\ \mathbb{U}_2^{j,1} \end{bmatrix}$ | $\mathbb{U}_{1\mid2}^{j,0\mid1} := (U_{0,1}^{j,0}, \ldots, U_{n,1}^{j,0})^T$ |
| $MoL$ | $\mathbb{U}(t) := \begin{bmatrix} \mathbb{U}_1(t) \\ \mathbb{U}_2(t) \end{bmatrix}$ | $\mathbb{U}_{1\mid2}(t) := (U_{0,1\mid2}(t), \ldots, U_{n,1\mid2}(t))^T$ |

Table 2.5: Vector notation with respect to $\mathcal{P}_1$ space discretisation and different time approximation. Here $U_{k,1}^j = U_1^j(x_k)$ and $U_{k,1\mid2}^{j,0\mid1} = U_{1\mid2}^{j,0\mid1}(x_k)$ for each $k = 0, \ldots, n$. Moreover, for $cG(1)$ in time, $U_{k,1\mid2}^j = U_{1\mid2}(t_j, x_k)$.

In the following we apply the particular time discretisation and deduce the resulting theorem for each case.

### 2.4.1.1 Existence and uniquence, $dG(q) \otimes \mathcal{P}_1$, $q = 0, 1$

For the understanding of the analysis below, recall the notation and definitions from Subsection 2.3.3 and Subsection 2.1.1 where $dG(q)$ method and linear splines are introduced.

**Theorem 2.4.1.1 (Existence of $dG(q) \otimes \mathcal{P}_1$ discrete solution, $q = 0, 1$).** There exists a unique function $U \in \mathcal{Q}_q$ such that $U$ solves the weak problem (2.5), where $\mathcal{B}, \mathscr{L}$ take the form of (2.18) and (2.19), respectively.

**Proof.** In the following, we provide only the proof for the case $q = 1$. The case $q = 0$ can be shown easily by using similar arguments.
From (2.18) and (2.19), for $U^j \in \mathcal{Q}_1^j$, a discrete solution of the reduced weak form (2.6) and

arbitrary test function $V^j \in \mathcal{Q}_1^j$, a weak problem (2.6) reads

$$
\frac{1}{k_j}\int_{I_j}\sum_{k,\ell=0}^{n} U_{k,1}^{j,1}(V_{\ell,1}^{j,0}+\frac{t-t_{j-1}}{k_j}V_{\ell,1}^{j,1})a(\phi_k;\phi_\ell)dt+\frac{1}{k_j}\int_{I_j}\sum_{k,\ell=0}^{n} U_{k,2}^{j,1}(V_{\ell,2}^{j,0}+\frac{t-t_{j-1}}{k_j}V_{\ell,2}^{j,1})(\phi_k;\phi_\ell)dt
$$

$$
-\int_{I_j}\sum_{k,\ell=0}^{n}(U_{k,2}^{j,0}+\frac{t-t_{j-1}}{k_j}U_{k,2}^{j,1})(V_{\ell,1}^{j,0}+\frac{t-t_{j-1}}{k_j}V_{\ell,1}^{j,1})a(\phi_k;\phi_\ell)dt
$$

$$
+\int_{I_j}\sum_{k,\ell=0}^{n}(U_{k,1}^{j,0}+\frac{t-t_{j-1}}{k_j}U_{k,1}^{j,1})(V_{\ell,2}^{j,0}+\frac{t-t_{j-1}}{k_j}V_{\ell,2}^{j,1})a(\phi_k;\phi_\ell)dt
$$

$$
+\varepsilon\int_{I_j}\sum_{k,\ell=0}^{n}(U_{k,2}^{j,0}+\frac{t-t_{j-1}}{k_j}U_{k,2}^{j,1})(V_{\ell,2}^{j,0}+\frac{t-t_{j-1}}{k_j}V_{\ell,2}^{j,1})a(\phi_k;\phi_\ell)dt
$$

$$
+\sum_{k,\ell=0}^{n} U_{k,1}^{j,0}V_{\ell,1}^{j,0}a(\phi_k;\phi_\ell)+\sum_{k,\ell=0}^{n} U_{k,2}^{j,0}V_{\ell,2}^{j,0}(\phi_k;\phi_\ell)
$$

$$
=\int_{I_j}\sum_{\ell=0}^{n}(V_{\ell,2}^{j,0}+\frac{t-t_{j-1}}{k_j}V_{\ell,2}^{j,1})(f;\phi_\ell)dt
$$

$$
+\sum_{k,\ell=0}^{n}\left((U_{k,1}^{j-1,0}+U_{k1}^{j-1,1})V_{\ell,1}^{j,0}a(\phi_k;\phi_\ell)+(U_{k,2}^{j-1,0}+U_{k,2}^{j-1,1})V_{\ell,2}^{j,0}(\phi_k;\phi_\ell)\right). \qquad (2.53)
$$

Note that for $j=1$, $U^{j-1,0}+U^{j-1,1}=\Pi u_0$. Here we added the zero components, i.e. components of the discrete functions $U,V$ whose indices correspond to the indices of the Dirichlet nodes.

In order to prove uniqueness let $f\equiv 0$, $U^{j-1}\equiv 0$ and $V^j\equiv U^j$. Recall the vector notation from Table 2.5. Hence, the equation (2.53) simplifies to

$$
0=\frac{1}{k_j}\int_{I_j}(\mathbb{U}_1^{j,1})^T S(\mathbb{U}_1^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_1^{j,1})dt+\frac{1}{k_j}\int_{I_j}(\mathbb{U}_2^{j,1})^T M(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1})dt
$$

$$
-\int_{I_j}(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1})^T S(\mathbb{U}_1^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_1^{j,1})dt+\int_{I_j}(\mathbb{U}_1^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_1^{j,1})^T S(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1})dt
$$

$$
+\varepsilon\int_{I_j}(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1})^T S(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1})dt+(\mathbb{U}_1^{j,0})^T S\mathbb{U}_1^{j,0}+(\mathbb{U}_2^{j,0})^T M\mathbb{U}_2^{j,0}=:\sum_{\ell=1}^{7}E_\ell.
$$

Since $S$ is symmetric, we have $E_3+E_4=0$. Also $E_5\geq 0$ because of $\varepsilon\geq 0$.
It is left to estimate $E_1, E_2, E_6, E_7$. Let

$$
g(t):=\frac{1}{2}(\mathbb{U}_1^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_1^{j,1})^T S(\mathbb{U}_1^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_1^{j,1}),
$$

$$
h(t):=\frac{1}{2}(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1})^T M(\mathbb{U}_2^{j,0}+\frac{t-t_{j-1}}{k_j}\mathbb{U}_2^{j,1}).
$$

From the main theorem of calculus we have

$$
E_1=\int_{I_j}\dot{g}(t)dt=g(t_j)-g(t_{j-1})
$$

and consequently

$$E_1 + E_6 = g(t_j) - g(t_{j-1}) + 2g(t_{j-1}) = g(t_j) + g(t_{j-1}) \geq 0.$$

The same holds for $E_2 + E_7$ by use of the function $h(t)$ instead of $g(t)$. This implies $0 = \sum_{\ell=1}^{7} E_\ell \geq 0$ which is valid only if $\mathbb{U}^j \equiv \mathbf{0}$. This concludes the proof. $\square$

### 2.4.1.2 Existence and uniquence, $cG(1) \otimes \mathcal{P}_1$

In case of the $cG(1)$ time discretisation and linear splines in space, recall the notation and definitions from Subsection 2.3.4 and Subsection 2.1.1 where $cG(1)$ method and linear splines are introduced. There holds the following theorem.

**Theorem 2.4.1.2 (Existence of $cG(1) \otimes \mathcal{P}_1$ discrete solution).** There is a unique function $U \in \mathcal{Q}_c$ which solves the weak problem (2.5), where $\mathcal{B}, \mathcal{L}$ take the form as in (2.41) and (2.42), respectively.

**Proof.** Due to the definition of the test space $\mathcal{W}_c$, we choose $V \in \mathcal{W}_c$ to be constant on each time interval $I_j$, i.e. $V|_{I_j}(t,x) = V^j(x) = \sum_{\ell=0}^{n} V_\ell^j \phi_\ell(x)$, where $V_\ell^j = (V_{\ell,1}^j, V_{\ell,2}^j) \in \mathbb{R}^2$. Note that here as in the Subsection 2.4.1.1, we add the zero components which correspond to the Dirichlet boundary conditions. For $U \in \mathcal{Q}_c$ and $V$ as above, the weak problem (2.6) simplifies to

$$\frac{1}{k_j} \int_{I_j} \sum_{k,\ell=0}^{n} U_{k,1}^j V_{\ell,1}^j a(\phi_k; \phi_\ell) + \frac{1}{k_j} \int_{I_j} \sum_{k,\ell=0}^{n} U_{k,2}^j V_{\ell,2}^j(\phi_k; \phi_\ell) - \int_{I_j} \sum_{k,\ell=0}^{n} \frac{t-t_{j-1}}{k_j} U_{k,2}^j V_{\ell,1}^j a(\phi_k; \phi_\ell) dt$$

$$+ \int_{I_j} \sum_{k,\ell=0}^{n} \frac{t-t_{j-1}}{k_j} U_{k,1}^j V_{\ell,2}^j a(\phi_k; \phi_\ell) dt + \varepsilon \int_{I_j} \sum_{k,\ell=0}^{n} \frac{t-t_{j-1}}{k_j} U_{k,2}^j V_{\ell,2}^j a(\phi_k; \phi_\ell) dt$$

$$= \int_{I_j} \sum_{\ell=0}^{n} V_{\ell,2}^j(f; \phi_\ell) dt + \frac{1}{k_j} \int_{I_j} \sum_{k,\ell=0}^{n} U_{k,1}^{j-1} V_{\ell,1}^j a(\phi_k; \phi_\ell) + \frac{1}{k_j} \int_{I_j} \sum_{k,\ell=0}^{n} U_{k,2}^{j-1} V_{\ell,2}^j(\phi_k; \phi_\ell)$$

$$+ \int_{I_j} \sum_{k,\ell=0}^{n} \frac{t_j-t}{k_j} U_{k,2}^{j-1} V_{\ell,1}^j a(\phi_k; \phi_\ell) dt - \int_{I_j} \sum_{k,\ell=0}^{n} \frac{t_j-t}{k_j} U_{k,1}^{j-1} V_{\ell,2}^j a(\phi_k; \phi_\ell) dt$$

$$- \varepsilon \int_{I_j} \sum_{k,\ell=0}^{n} \frac{t_j-t}{k_j} U_{k,2}^{j-1} V_{\ell,2}^j a(\phi_k; \phi_\ell) dt. \tag{2.54}$$

Note that the initial condition reads $U(0) = \Pi u_0$. It suffices to show that for each $j = 1, \ldots, N$, the weak solution $U^j(t,x)$ of the homogeneous equation (2.54) ($f \equiv 0$ and $U^{j-1} \equiv 0$), equals zero. Let also $V^j \equiv U^j$. Hence, owing to the vector notation from Table 2.5, the homogeneous equation (2.54) reads

$$0 = (\mathbb{U}_1^j)^T S \mathbb{U}_1^j + (\mathbb{U}_2^j)^T M \mathbb{U}_2^j - \frac{k_j}{2} (\mathbb{U}_2^j)^T S \mathbb{U}_1^j + \frac{k_j}{2} (\mathbb{U}_1^j)^T S \mathbb{U}_2^j + \varepsilon \frac{k_j}{2} (\mathbb{U}_2^j)^T S \mathbb{U}_2^j =: \sum_{\ell=1}^{5} E_\ell.$$

Since $S, M$ are symmetric, we have $E_1, E_2 \geq 0$ and $E_3 + E_4 = 0$. Because of $\varepsilon \geq 0$, we also have $E_5 \geq 0$.
This yields $0 = \sum_{\ell=1}^{5} E_\ell \geq 0$ which implies $\mathbb{U}^j \equiv 0$. This concludes the proof. $\square$

### 2.4.1.3 Existence and uniquence, $MoL \otimes \mathcal{P}_1$

Before we start with a proof of the main theorem, recall the notation and definitions from Subsection 2.3.5 and Subsection 2.1.1 where the method of lines and linear splines are introduced. Thereafter, we derive the following theorem.

**Theorem 2.4.1.3 (Existence of $MoL \otimes \mathcal{P}_1$ semi-discrete solution).** There exists a unique function $U \in \mathcal{Q}_s$ which solves the weak problem (2.5), where $\mathcal{B}, \mathscr{L}$ take the form as in (2.49) and (2.50), respectively. This function also satisfies the following initial conditions

$$\mathbb{U}(0) = \left[ \begin{array}{c} \mathbb{U}_1(0) \\ \mathbb{U}_2(0) \end{array} \right] = \left[ \begin{array}{c} M^{-1}((y_0; \phi_0), \ldots, (y_0; \phi_n))^T \\ M^{-1}((y_1; \phi_0), \ldots, (y_1; \phi_n))^T \end{array} \right], \tag{2.55}$$

defined owing to (1.13b) and the notation from Table 2.5.

**Proof.** Given the weak problem (2.5), let $V := (\phi_\ell, 0)$ and $V := (0, \phi_\ell)$ for $\ell \in \{0, \ldots, n\}$, respectively. Then the weak problem (2.5) simplifies to the following two equations

$$\sum_{k=0}^{n} \dot{U}_{k,1}(t) a(\phi_k; \phi_\ell) - \sum_{k=0}^{n} U_{k,2}(t) a(\phi_k; \phi_\ell) = 0, \tag{2.56a}$$

$$\sum_{k=0}^{n} \dot{U}_{k,2}(t)(\phi_k; \phi_\ell) + \sum_{k=0}^{n} a(\phi_k; \phi_\ell) + \varepsilon \sum_{k=0}^{n} U_{k,2}(t) a(\phi_k; \phi_\ell) = (f(t); \phi_\ell), \tag{2.56b}$$

where $\ell = 1, \ldots, n$ and $t \in [0, T]$. With

$$F(t) := \big((f(t); \phi_0), \ldots, (f(t); \phi_n)\big)^T,$$

and notation from Table 2.5, (2.56) becomes the linear system of ODE

$$\left[ \begin{array}{cc} S & 0 \\ 0 & M \end{array} \right] \left[ \begin{array}{c} \dot{\mathbb{U}}_1(t) \\ \dot{\mathbb{U}}_2(t) \end{array} \right] + \left[ \begin{array}{cc} 0 & -S \\ S & \varepsilon S \end{array} \right] \left[ \begin{array}{c} \mathbb{U}_1(t) \\ \mathbb{U}_2(t) \end{array} \right] = \left[ \begin{array}{c} 0 \\ F(t) \end{array} \right] \tag{2.57}$$

which satisfies the initial conditions (2.4.1.3). According to the theory of ordinary differential equations, there exists a unique vector function $\mathbb{U}(t)$ defined as in Table 2.5 which solves the initial value problem stated above. This concludes the proof.                     $\square$

**Remark 2.4.1.1.** The proof of Theorem 2.4.1.3 imposes also the existence of the discrete solution. In addition to that, it comprises the algorithm for the calculation of the discrete solution in each time point $t \in (0, T]$ which is obvious from (2.57).                     $\square$

## 2.4.2   Hermite cubic splines ($\mathcal{C}^1$)

We proceed in the following by discretising in space by means of the Hermite cubic finite elements.
Let $S_k$ and $M_k$ denote the $4 \times 4$ element defined stiffness and mass matrix respectively, defined for each $k = 1, \ldots, n$, where $k$ stands for the number of elements in space.
If $\phi_1, \ldots, \phi_4$ are the basic functions from (2.3) related to each space interval $T_k \in \mathcal{T}_j$, we have

$$S_{k,p,q} := \left\{ \begin{array}{ll} 2/h_k \int_{-1}^{1} D\phi_p D\phi_q d\zeta & \text{for } p, g \text{ odd,} \\ \int_{-1}^{1} D\phi_p D\phi_q d\zeta & \text{for } p + g \text{ odd,} \\ h_k/2 \int_{-1}^{1} D\phi_p D\phi_q d\zeta & \text{for } p, g \text{ even,} \end{array} \right. \tag{2.58}$$

and

$$M_{k,p,q} := \begin{cases} h_k/2 \ \int_{-1}^{1} \phi_p \phi_q d\zeta & \text{for } p, g \text{ odd}, \\ h_k^2/4 \ \int_{-1}^{1} \phi_p \phi_q d\zeta & \text{for } p + g \text{ odd}, \\ h_k^3/8 \ \int_{-1}^{1} \phi_p \phi_q d\zeta & \text{for } p, g \text{ even}. \end{cases} \tag{2.59}$$

Besides, we will also use the full form of $2(n+1) \times 2(n+1)$ block mass and stiffness matrix, noted as $M$ and $S$. The structure of both element matrices $S_k$ and $M_k$ is explained in detail in Section 6.1 where the algorithm for the calculation of FE solution is introduced, cf. (6.6) and (6.7).

In the following, we introduce the vector notation which depends on time discretisation and corresponds to the discrete functions from Table 2.1–2.3, case cubic splines in space. On account to the definition of the cubic splines we first define the local vectors for each $T_k$ space interval and thus the global (full) vectors related to the whole space domain $\Omega$.

For each time-space slab $I_j \times T_k$ where $(k,j) \in \{1, \ldots, n\} \times \{1, \ldots, N\}$, or each $t \in [0, T]$, these vectors take form as in Table 2.6.

| $dG(0), cG(1)$ | $\mathbb{U}^j := \begin{bmatrix} \mathbb{U}_1^j \\ \mathbb{U}_2^j \end{bmatrix}$ | $\mathbb{U}_{k,1\vert 2}^j := (U_{k-1,1\vert 2}^j, DU_{k-1,1\vert 2}^j, U_{k,1\vert 2}^j, DU_{k,1\vert 2}^j)^T,$ $\mathbb{U}_{1\vert 2}^j := (U_{0,1\vert 2}^j, DU_{0,1\vert 2}^j, \ldots, DU_{n,1\vert 2}^j)^T$ |
|---|---|---|
| $dG(1)$ | $\mathbb{U}^j := \begin{bmatrix} \mathbb{U}_1^{j,0} \\ \mathbb{U}_1^{j,1} \\ \mathbb{U}_2^{j,0} \\ \mathbb{U}_2^{j,1} \end{bmatrix}$ | $\mathbb{U}_{k,1\vert 2}^{j,0\vert 1} := (U_{k-1,1\vert 2}^j, DU_{k-1,1\vert 2}^j, U_{k,1\vert 2}^j, DU_{k,1\vert 2}^j)^T,$ $\mathbb{U}_{1\vert 2}^{j,0\vert 1} := (U_{0,1\vert 2}^{j,0\vert 1}, DU_{0,1\vert 2}^{j,0\vert 1}, \ldots, DU_{n,1\vert 2}^{j,0\vert 1})^T$ |
| $MoL$ | $\mathbb{U}(t) := \begin{bmatrix} \mathbb{U}_1(t) \\ \mathbb{U}_2(t) \end{bmatrix}$ | $\mathbb{U}_{k,1\vert 2}(t) := (U_{k-1,1\vert 2}(t), DU_{k-1,1\vert 2}(t), U_{k,1\vert 2}(t), DU_{k,1\vert 2}(t))^T,$ $\mathbb{U}_{1\vert 2}(t) := (U_{0,1\vert 2}(t), DU_{0,1\vert 2}(t), \ldots, DU_{n,1\vert 2}(t))^T,$ |

Table 2.6: Full and local vector notation with respect to $\mathcal{C}^1$ space discretisation and different time approximation. Here $U_{k,1\vert 2}^{j,0\vert 1} = U_{1\vert 2}^{j,0\vert 1}(x_k)$ and $DU_{k,1\vert 2}^{j,0\vert 1} = DU_{1\vert 2}^{j,0\vert 1}(x_k)$ for each $k = 1, \ldots, n$. Moreover, if $cG(1)$ in time, $U_{k,1\vert 2}^j := U_{1\vert 2}(t_j, x_k)$ and $DU_{k,1\vert 2}^j = DU_{1\vert 2}(t_j, x_k)$.

We continue by applying the corresponding time discretisation methods to the local $k$-contribution of the weak problem, cf. Remark 2.3.1.1. If we denote by $\mathcal{B}^j|_{T_k} =: \mathcal{B}_k^j$ and $\mathcal{L}^j|_{T_k} =: \mathcal{L}_k^j$, then there holds

$$\sum_{k=1}^{n} \mathcal{B}_k^j(U, V) = \mathcal{B}^j(U, V) = \mathcal{L}^j(U, V) = \sum_{k=1}^{n} \mathcal{L}_k^j(U, V) \quad \text{for all} \quad j = 1, \ldots, N. \tag{2.60}$$

### 2.4.2.1 Existence and uniquence, $dG(q) \otimes \mathcal{C}^1$, $q = 0, 1$

The related definitions and notation are taken over from Subsection 2.1.2, where cubic splines are introduced and Subsection 2.3.3 where disountinous Galerkin method for time discretisation is described.

**Theorem 2.4.2.1 (Existence of $dG(q) \otimes \mathcal{C}^1$ discrete solution, $q = 0, 1$).** There exists a unique function $U \in \mathcal{Q}_q$ such that $U$ solves the weak problem (2.5), where where $\mathcal{B}, \mathcal{L}$ take the form as in (2.18) and (2.19), respectively.

**Proof.** In the proof we restrict to the case $q = 0$. The case $q = 1$ can be shown by following the ideas from Subsection 2.4.1.1 where the proof for uniqueness of $dG(1) \otimes \mathcal{P}_1$ discrete solution is derived.

The approximated weak problem (2.6) on $k$th element $T_k$ of the triangulation $\mathcal{T}_j$ can be rewritten in form of contributions $\mathcal{B}_k^j$, $\mathcal{L}_k^j$ such that

$$-\int_{I_j}\int_{T_k} DU_2^j DV_1^j \, dxdt + \int_{I_j}\int_{T_k} DU_1^j DV_2^j \, dxdt + \varepsilon \int_{I_j}\int_{T_k} DU_2^j DV_2^j \, dxdt$$
$$+ \int_{T_k} DU_1^j DV_1^j \, dx + \int_{T_k} U_2^j V_2^j \, dx$$
$$= \int_{I_j}\int_{T_k} f V_2^j \, dxdt + \int_{T_k} DU_1^{j-1} DV_1^j \, dx + \int_{T_k} U_2^{j-1} V_2^j \, dx. \quad (2.61)$$

Note that for $j = 1$, $U^{j-1} = \Pi u_0$. Given the Galerkin form (2.4), by means of the notation from Table 2.6 and definitions of $S_k$ and $M_k$, cf. (2.58) and (2.59), respectively, the weak form above reads

$$(\mathbb{V}_{k,1}^j)^T (S_k \mathbb{U}_{k,1}^j - k_j S_k \mathbb{U}_{k,2}^j) + (\mathbb{V}_{k,2}^j)^T (k_j S_k \mathbb{U}_{k,1}^j + (M_k + \varepsilon k_j S_k) \mathbb{U}_{k,2}^j)$$
$$= F_k^j + (\mathbb{V}_{k,1}^j)^T G_{k,1}^j + (\mathbb{V}_{k,2}^j)^T G_{k,2}^j. \quad (2.62)$$

The structure of the $4 \times 1$ element forcing vector $F_k^j$ and the element mass vectors $G_{k,1}^j$, $G_{k,2}^j$ is explained in detail in Subsection 6.1.1.2. The corresponding full formulations are also given there.
From (2.60) and (2.62) by using the full vector and matrix form, we obtain

$$(\mathbb{V}_1^j)^T (S\mathbb{U}_1^j - k_j S\mathbb{U}_2^j) + (\mathbb{V}_2^j)^T (k_j S\mathbb{U}_1^j + (M + \varepsilon k_j S)\mathbb{U}_2^j) = F^j + (\mathbb{V}_1^j)^T G_1^j + (\mathbb{V}_2^j)^T G_2^j. \quad (2.63)$$

In order to prove the uniqueness of the discrete solution in each time step, it is sufficient to show that the weak solution $\mathbb{U}^j$ of the homogeneous problem (2.63) ($F^j \equiv 0$, $\mathbb{U}^{j-1} \equiv 0$, and therefore $G_1^j$, $G_2^j \equiv 0$) is $\mathbb{U}^j \equiv 0$.
To verify this, let $\mathbb{V}^j \equiv \mathbb{U}^j$ and since $M$ and $S$ are also symmetric matrices, the equation (2.63) simplifies to

$$(\mathbb{U}_1^j)^T S\mathbb{U}_1^j + (\mathbb{U}_2^j)^T M\mathbb{U}_2^j + \varepsilon k_j (\mathbb{U}_2^j)^T S\mathbb{U}_2^j = 0. \quad (2.64)$$

From $\varepsilon \geq 0$ and the positive definiteness of $S, M$, we conclude that the LHS of (2.64) must be positive. Then the equality (2.64) holds only if $\mathbb{U}^j \equiv 0$. This concludes the proof. $\qquad\square$

## 2.4.2.2 Existence and uniquence, $cG(1) \otimes \mathcal{C}^1$

In case of $cG(1)$ time approximation, we refer to the notation and definitions from Subsection 2.3.4. Furthermore, the definition of $\mathcal{C}^1$ elements, i.e. Hermite cubic splines, can be found in Subsection 2.1.2. Thereafter, we derive the following theorem.

**Theorem 2.4.2.2 (Existence of $cG(1) \otimes \mathcal{C}^1$ discrete solution).** There exists a unique function $U \in \mathcal{Q}_c$ which solves the weak problem (2.5), where $\mathcal{B}, \mathscr{L}$ take the form as in (2.41) and (2.42), respectively.

**Proof.** On account to the definition of the test space $\mathcal{W}_c$, we choose a test function $V$ to be constant in time. Then, on each time-space slab $I_j \times T_k$, where $T_k$ denotes arbitrary element of the space triangulation $\mathcal{T}_j$, there holds

$$\int_{I_j} \int_{T_k} \frac{1}{k_j} DU_{k,1}^j DV_1^j \, dx + \int_{I_j} \int_{T_k} \frac{1}{k_j} U_2^j V_2^j \, dx - \int_{I_j} \int_{T_k} \frac{t-t_{j-1}}{k_j} DU_2^j DV_1^j \, dxdt$$

$$+ \int_{I_j} \int_{T_k} \frac{t-t_{j-1}}{k_j} DU_1^j DV_2^j \, dxdt + \varepsilon \int_{I_j} \int_{T_k} \frac{t-t_{j-1}}{k_j} DU_2^j DV_2^j \, dxdt$$

$$= \int_{I_j} \int_{T_k} f V_2^j \, dxdt + \int_{I_j} \int_{T_k} \frac{1}{k_j} DU_1^{j-1} DV_1^j \, dx + \int_{I_j} \int_{T_k} \frac{1}{k_j} U_2^{j-1} V_2^j \, dx$$

$$+ \int_{I_j} \int_{T_k} \frac{t_j-t}{k_j} DU_2^{j-1} DV_1^j \, dxdt - \int_{I_j} \int_{T_k} \frac{t_j-t}{k_j} DU_1^{j-1} DV_2^j \, dxdt$$

$$- \varepsilon \int_{I_j} \int_{T_k} \frac{t_j-t}{k_j} DU_2^{j-1} DV_2^j \, dxdt. \tag{2.65}$$

Note that for $j = 1$, $U^{j-1} = \Pi u_0$. By expanding discrete function $U^j, V^j$ in terms of basic functions, see (2.4), and using the notation from Table 2.6, (2.65) can be rewritten such that

$$(\mathbb{V}_{k,1}^j)^T (S_k \mathbb{U}_{k,1}^j - \frac{k_j}{2} S_k \mathbb{U}_{k,2}^j) + (\mathbb{V}_{k,2}^j)^T (\frac{k_j}{2} S_k \mathbb{U}_{k,1}^j + (M_k + \varepsilon \frac{k_j}{2} S_k) \mathbb{U}_{k,2}^j)$$

$$= F_k^j + (\mathbb{V}_{k,1}^j)^T G_{k,1}^j + (\mathbb{V}_{k,2}^j)^T G_{k,2}^j, \tag{2.66}$$

The structure of the $4 \times 1$ element forcing vector $F_k^j$ and mass vectors $G_{k,1|2}^j$, which depends on $\mathbb{U}^{j-1}$, is explained in detail in the Subsection 6.1.2.2. The corresponding full formulations are also given there.

From (2.60) and (2.66), using the full vector and matrix form we obtain

$$(\mathbb{V}_1^j)^T (S \mathbb{U}_1^j - \frac{k_j}{2} S \mathbb{U}_2^j) + (\mathbb{V}_2^j)^T (\frac{k_j}{2} S \mathbb{U}_1^j + (M + \varepsilon \frac{k_j}{2} S) \mathbb{U}_2^j) = F^j + (\mathbb{V}_1^j)^T G_1^j + (\mathbb{V}_2^j)^T G_2^j. \tag{2.67}$$

In order to prove the uniqueness of the discrete solution it is sufficient to show that the solution $\mathbb{U}^j$ of the homogeneous problem (2.67) where $F^j \equiv 0$, $\mathbb{U}^{j-1} \equiv 0$, and therefore $G_1^j$, $G_2^j \equiv 0$, is $\mathbb{U}^j \equiv 0$.

To verify this, let $\mathbb{V}^j = \mathbb{U}^j$ and since $M$ and $S$ are also symmetric matrices, the equation (2.67) simplifies to

$$(\mathbb{U}_1^j)^T S \mathbb{U}_1^j + (\mathbb{U}_2^j)^T M \mathbb{U}_2^j + \varepsilon \frac{k_j}{2} (\mathbb{U}_2^j)^T S \mathbb{U}_2^j = 0. \tag{2.68}$$

From $\varepsilon \geq 0$ and positive definiteness of $S, M$, we conclude that the LHS of (2.68) must be positive. This holds only if $\mathbb{U}^j \equiv 0$ and the proof of theorem follows. $\qquad \square$

### 2.4.2.3 Existence and unique, $MoL \otimes \mathcal{C}^1$

Notation and definitions used below are taken over from the Subsection 2.3.5 and Subsection 2.1.2 where the method of lines and the Hermite cubic splines, respectively, are introduced.

**Theorem 2.4.2.3 (Existence of $MoL \otimes C^1$ semi-discrete solution).** There exists a unique function $U \in \mathcal{Q}_s$ which solves the weak problem (2.5), where $\mathcal{B}, \mathcal{L}$ take the form of (2.49) and (2.50), respectively. This function also satisfies the following initial conditions,

$$\mathbb{U}_{k,1}(0) = M_k^{-1} \begin{bmatrix} h_k/2 \ \int_{-1}^{1} y_0 \phi_1 d\zeta \\ h_k^2/4 \ \int_{-1}^{1} y_0 \phi_2 d\zeta \\ h_k/2 \ \int_{-1}^{1} y_0 \phi_3 d\zeta \\ h_k^2/4 \ \int_{-1}^{1} y_0 \phi_4 d\zeta \end{bmatrix}, \quad \mathbb{U}_{k,2}(0) = M_k^{-1} \begin{bmatrix} h_k/2 \ \int_{-1}^{1} y_1 \phi_1 d\zeta \\ h_k^2/4 \ \int_{-1}^{1} y_1 \phi_2 d\zeta \\ h_k/2 \ \int_{-1}^{1} y_1 \phi_3 d\zeta \\ h_k^2/4 \ \int_{-1}^{1} y_1 \phi_4 d\zeta \end{bmatrix}, \tag{2.69}$$

defined owing to (1.13d) and the notation from Table 2.6,

**Proof.** From (2.49) the approximated weak form (2.5) on $k$th element $T_k$ of the triangulation $\mathcal{T}$ can be rewritten in form of contributions such that

$$\int_{T_k} D\dot{U}_1(t) DV_1 dx + \int_{T_k} \dot{U}_2(t) V_2 dx - \int_{T_k} DU_2(t) DV_1 dx$$

$$+ \int_{T_k} DU_1(t) DV_2 dx + \varepsilon \int_{I_j} \int_{T_k} DU_2(t) DV_2 dx$$

$$= \int_{T_k} f(t) V_2 dx \quad \text{for all} \quad t \in [0, T].$$

If we set in the equation above a constant function in time $V \in \mathcal{S} \times \mathcal{S}$, by use of (2.4) and notation from Table 2.6, we obtain for all $t \in [0, T]$

$$(\mathbb{V}_{k,1})^T (S_k \dot{\mathbb{U}}_{k,1}(t) - S_k \mathbb{U}_{k,2}(t)) + (\mathbb{V}_{k,2})^T (S_k \mathbb{U}_{k,1}(t) + \varepsilon S_k \mathbb{U}_{k,2}(t) + M_k \dot{\mathbb{U}}_{k,2}(t)) = F_k(t). \tag{2.70}$$

$F_k(t)$ is the $4 \times 1$ element force function vector defined by

$$F_k(t) := \begin{bmatrix} h_k/2 \ \int_{-1}^{1} f(t)\phi_1 \ d\zeta \\ h_k^2/4 \ \int_{-1}^{1} f(t)\phi_2 \ d\zeta \\ h_k/2 \ \int_{-1}^{1} f(t)\phi_3 \ d\zeta \\ h_k^2/4 \int_{-1}^{1} f(t)\phi_4 \ d\zeta \end{bmatrix}. \tag{2.71}$$

If we choose a discrete function $V$ such that

$$V = (V_1, V_2) \in \{(\phi_p, 0), (0, \phi_p)\} \quad \text{with} \ 1 \le p \le 4,$$

then for all $k = 1, \ldots, n$, $\mathbb{V}_{k,1|2} \in \{e_1, e_2, e_3, e_4\} \in \mathbb{R}^4$. By summing the contributions (2.70) over $k = 1, \ldots, n$ we obtain two systems of equations. Namely,

1. for $V = (\phi_p, 0)$, $1 \le p \le 4$

$$\sum_{k=1}^{n} S_k \dot{\mathbb{U}}_{k,1}(t) - S_k \mathbb{U}_{k,2}(t) = 0,$$

2. for $V = (0, \phi_p)$, $1 \le p \le 4$

$$\sum_{k=1}^{n} S_k \mathbb{U}_{k,1}(t) + \varepsilon S_k \mathbb{U}_{k,2}(t) + M_k \dot{\mathbb{U}}_{k,2} = \sum_{k=1}^{n} F_k(t).$$

Given the global mass and stiffness matrices $M, S$ and similarly defined $2(n+1) \times 1$ block vector $F(t)$, cf. (2.71), two systems above can be written in form of the system of ODEs such that

$$\begin{bmatrix} S & 0 \\ 0 & M \end{bmatrix} \begin{bmatrix} \dot{\mathbb{U}}_1(t) \\ \dot{\mathbb{U}}_2(t) \end{bmatrix} + \begin{bmatrix} 0 & -S \\ S & \varepsilon S \end{bmatrix} \begin{bmatrix} \mathbb{U}_1(t) \\ \mathbb{U}_2(t) \end{bmatrix} = \begin{bmatrix} 0 \\ F(t) \end{bmatrix} \tag{2.72a}$$

with the initial condition

$$\mathbb{U}(0) = \begin{bmatrix} \mathbb{U}_1(0) \\ \mathbb{U}_2(0) \end{bmatrix}. \tag{2.72b}$$

From the theory of ordinary differential equations, there exists an unique vector function $\mathbb{U}(t)$ who solves the initial value problem stated above. $\qquad\square$

**Remark 2.4.2.1.** The proof of Theorem 2.4.2.3 imposes also the existence of the discrete solution. Above that, it comprises the algorithm for the calculation of the discrete solution in each time point $t \in (0, T]$. This is obvious from (2.72).

# Chapter 3

# Energy Method

In the following chapter the main focus will be posed on derivation and analysis of *a priori* and *a posteriori* error estimates by using the energy techniques.

The proof of both a priori and a posteriori error bound has the similar structure. Namely, it is based on the representation of the error in terms of the residual (a posteriori error analysis) or in terms of the exact solution (a priori error analysis) combined with the Galerkin orthogonality. It also assumes the usage of some interpolation and projection operators and appropriate estimates. The following analysis will be conducted for a priori and a posteriori error bounds separately. We start first by introducing the time discretisation, i.e. Galerkin discretisation methods ($dG(q)$, $q = 0, 1$, $cG(1)$) and method of lines and then by coupling it with two space ansatz, linear and cubic splines. Thereafter we analyse the error bounds with respect to the corresponding fully discrete form.

Note that the realization and construction of adaptive refinement strategy is based on a posteriori error estimates. This will be further emphasised in Chapter 6.4.

| | | |
|---|---|---|
| $u$ | exact solution | solution of (1.28) |
| $U$ | discrete solution | solution of (2.5) |
| $\widetilde{U}$ | affine interpolant of the discrete solution $U$ | Definition 2.3.3.3 |
| $e$ | error, $e = u - U$ | Definition 2.3.2.1 |
| $\tilde{e}$ | error, $\tilde{e} = u - \widetilde{U}$ | definition (3.89) |
| $\tilde{\tilde{e}}$ | error, $\tilde{\tilde{e}} = \tilde{u} - \widetilde{U}$ | definition (3.51) |
| $\mathcal{B}$ | weak bilinear form | Subsection 2.3.1 |
| $\widetilde{\mathcal{B}}$ | weak bilinear form, $dG(0)$ in time | definition (2.20) |
| $\mathscr{L}$ | RHS in the variational formulation | Subsection 2.3.1 |
| $Res(v)$ | residual, $Res(v) = \mathcal{B}(u, v) - \mathcal{B}(U, v)$ | Definition 2.3.1.2 |
| $\widetilde{Res}(v)$ | residual, $\widetilde{Res}(v) = \mathcal{B}(u, v) - \widetilde{\mathcal{B}}(U, v)$ | definition (2.21) |
| $\Pi$ | spatial multi projection | Definition 3.1.0.5 |
| $\mathcal{I}$ | nodal or cubic Hermite interpolation operator | Definition 3.1.0.1, 3.1.0.2 |
| $\mathcal{G}$ | spatial $H^1$ projection, Galerkin projection | Definition 3.1.0.3 |
| $\mathcal{L}$ | spatial $L^2$ projection | Definition 3.1.0.4 |
| $\mathcal{J}$ | temporal projection | Definition 3.1.0.9 |
| $D_{h,m}$ | discrete variant of Laplace operator | Definition 3.1.0.7 |
| $\langle \cdot ; \cdot \rangle_{\mathcal{H}}, \| \cdot \|_{\mathcal{H}}$ | energy scalar product, energy norm | Definition 1.3.0.3 |

Table 3.1: Notation used in Chapter 3.

# 3.1   Preliminaries

Within this section some general remarks concerning the finite element method, interpolation and projection operators together with appropriate bounds are introduced. We also provide some basic theory results such are trace theorem and Friedrics and Cauchy inequality. These theoretical results denote the main tool for the subsequent error analysis and therefore need to be explained in detail.

We do this by following the definitions of cf. BRENNER-SCOT [15], CIARLET [20], EVANS [27] as well as by developing some of the ideas which will be referred in particular while introducing the corellative theorem, lemma or definition.

**Definition 3.1.0.1 (Nodal interpolant).** Let $\mathcal{S}^j$ be the spatial discrete space related to the time interval $I_j \in \mathcal{T}$, cf. Definition 2.1 which consists of piecewise linear globally continuous functions. Then for any given function $u \in \mathcal{C}(\Omega)$, we define nodal interpolant by

$$\mathcal{I}u(x) := \sum_{j=0}^{n} u(x_j)\phi_j(x), \quad \text{for all} \quad x \in \Omega,$$

where $(\phi_j)_0^n$ are the basis functions introduced in Subsection 2.1.                               $\square$

**Definition 3.1.0.2 (Hermite cubic interpolant).** The Hermite cubic interpolant to $u \in H^2(\Omega)$ is the unique polynomial $\mathcal{I}u$ defined such that for each interval $T_k, k = 1, \ldots, n$ from the spatial triangulation $\mathcal{T}_j, j = 1, \ldots, N$, see Definition 2.0.0.8, $\mathcal{I}_k := \mathcal{I}|_{T_k}$ satisfies

$$(u - \mathcal{I}_k u)(x(\zeta)) = 0, \text{ and } D(u - \mathcal{I}_k u)(x(\zeta)) = 0 \text{ for each } x \in T_k \text{ and } \zeta \in \{-1, 1\}. \quad \square$$

**Lemma 3.1.0.1 (Friedrichs inequality in 1D, AFEM [16, Section 5]).** On the bounded interval $\Omega = (a, b)$ for all functions $u \in H_0^1(a, b)$ there holds

$$\|u\|_{L^2(a,b)} \leq \frac{b-a}{\pi} \|u\|_{H^1(a,b)}.$$

Moreover, if $u \in H^1(a, b)$ vanishes somewhere on the compact interval $[a, b]$, then

$$\|u\|_{L^2(a,b)} \leq \frac{2(b-a)}{\pi} \|u\|_{H^1(a,b)}.$$

**Lemma 3.1.0.2 (Approximation properties of Nodal Interpolation).** Given $u \in H^1(\Omega)$ and its nodal interpolant $\mathcal{I}u$, there holds

$$\|u - \mathcal{I}u\|_{L^2(\Omega)} \leq \frac{1}{\pi} \|hDu\|_{L^2(\Omega)}. \tag{3.1a}$$

Moreover, if $u \in H^2(\Omega)$ then additionally holds

$$\|u - \mathcal{I}u\|_{L^2(\Omega)} \leq \frac{1}{4} \|h^2 \Delta u\|_{L^2(\Omega)}, \tag{3.1b}$$

$$\|D(u - \mathcal{I}u)\|_{L^2(\Omega)} \leq \frac{1}{\sqrt{2}} \|h\Delta u\|_{L^2(\Omega)}. \tag{3.1c}$$

**Proof.** We provide only the proof for the first estimate (3.1a) is provided. For the other two estimates we refer to AFEM [16, Theorem 2.4].

Let us assume that there exists a constant $C$ such that

$$\|u - \mathcal{I}u\|_{L^2(\Omega)} \leq C\|hu\|_{H^1(\Omega)}. \tag{3.2}$$

The idea is to find an optimal one for the case $\Omega = (0,1)$. The Friedrichs inequality, restricted to the interval $T_k$ yields for some $u \in H^1_D(\Omega)$

$$\|u - \mathcal{I}u\|_{L_2(T_k)} \leq \frac{h_k}{\pi}\|u - \mathcal{I}u\|_{H^1(T_k)}. \tag{3.3}$$

Here we used the fact that $u$ and $\mathcal{I}$ coincide in each $h_k, k = 1, \ldots, n$.
Owing to the properties of discrete functions we have that $a(u - \mathcal{I}u; \mathcal{I}u) = 0$. This implies

$$\|u - \mathcal{I}u\|^2_{H^1(T_k)} = \|u\|^2_{H^1(T_k)} - \|\mathcal{I}u\|^2_{H^1(T_k)}. \tag{3.4}$$

From (3.3) and (3.4) there holds for each $T_k$

$$\|u - \mathcal{I}u\|_{L_2(T_k)} \leq \frac{h_k}{\pi}\|u\|_{H^1(T_k)}.$$

The last inequality is equivalent to

$$\|u - \mathcal{I}u\|_{L_2(\Omega)} \leq \frac{1}{\pi}\|hu\|_{H^1(\Omega)}.$$

According to (3.2), we have that $C \leq 1/\pi$.
Notice that the equality is satisfied with the function $u(x) = sin(x\pi/h)$ where $h$ is the uniform space step. This proves that the constant $C = 1/\pi$ is optimal.  $\square$

**Lemma 3.1.0.3 (Trace Identity in 1D, AFEM [16, Chapter II, Claim 3.1.1]).** For some function $u \in H^1(a,b)$, the following estimate holds

$$\max_{a \leq x \leq b} |u(x)| \leq (b-a)^{-1/2}\|u\|_{L^2(a,b)} + (b-a)^{1/2}\|Du\|_{L^2(a,b)}.$$

**Lemma 3.1.0.4 (Inverse estimate, CIARLET [20, pp. 142]).** Let $\mathcal{S} \subset H^1_D(\Omega)$ be a finite element space and $h$ the step size with respect to triangulation of the domain $\Omega$. Then there exists a positive constant $C_{inv}$, such that for all $U \in \mathcal{S}$,

$$\|U\|_{H^1(\Omega)} \leq C_{inv}\|h^{-1}U\|_{L^2(\Omega)}. \tag{3.5}$$

**Lemma 3.1.0.5 (ODEN-REDDY [57, Theorem 6.8, Section 8.5]).** Given $u \in H^4(\Omega)$ on some bounded interval $\Omega$ and its Hermite cubic interpolation operator $\mathcal{I}$ from Definition 3.1.0.2, there exists a constant $C$ independent of $u$ and $h$ such that for all $0 \leq m \leq 4$ there holds

$$\|u - \mathcal{I}u\|_{H^m(\Omega)} \leq C\|h^{4-m}u\|_{H^4(\Omega)}. \tag{3.6a}$$

Moreover, if $u$ is not smooth enough e.g. $u \in H^r\Omega$ for $r < 4$, then

$$\|u - \mathcal{I}u\|_{H^m(\Omega)} \leq C\|h^{r-m}u\|_{H^r(\Omega)}. \tag{3.6b}$$

where $0 \leq m \leq r$.

**Remark 3.1.0.2.** For interpolation operator $\mathcal{I}$, there hold the local interpolation estimates. Namely, for each element $T_k$ from the arbitrary triangulation of the spatial domain $\Omega$ there exists a constant $C$ such that for all functions $u \in H^{p+1}(T_k)$

$$\|u - \mathcal{I}u\|_{H^m(T_k)} \leq C h_k^{p+1-m} \|D^{p+1-m}u\|_{L^2(T_k)}, \quad 0 \leq m \leq p+1,$$

where $p$ is the polynomial order of the discrete functions, i.e. $p=1$ for linear splines and $p=3$ for cubic Hermite splines. For details, see CIARLET [20, THEOREM 3.1.6]. $\qquad\square$

**Definition 3.1.0.3 (Galerkin projection).** We consider $\mathcal{G}$ to be a Galerkin (elliptic) projection operator onto the space of discrete functions in space. Then for each $I_j$ element of triangulation $\mathscr{T}$ from Definition 2.0.0.8 and corresponding discrete space $\mathcal{S}^j$, see (2.1), $\mathcal{G}|_{I_j} : H_D^1(\Omega) \to \mathcal{S}^j$ is defined by

$$a(\mathcal{G}u; v) = a(u; v) \quad \text{for all} \quad v \in \mathcal{S}^j \text{ and } j = 1 \ldots, N. \qquad\square$$

**Lemma 3.1.0.6 (HACKBUSCH [35, Theorem 8.5.1, Remark 8.5.2]).** There exists a constant $C$ independent of $u$ and $h$ such for $\mathcal{G}$, Galerkin projection from Definition 3.1.0.3 and all $u \in H_D^1(\Omega) \cap H^s(\Omega)$, $0 \leq r \leq 1 \leq s \leq p+1$ there holds

$$\|u - \mathcal{G}u\|_{H^r(\Omega)} \leq C \|h^{s-r} D^s u\|_{L^2(\Omega)},$$

where $p=1$ in case of $\mathcal{P}^1$ elements and $p=3$ for $\mathcal{C}^1$ elements in space.

**Remark 3.1.0.3.** Let $\mathcal{T}_h$ be an arbitrary triangulation of $\Omega = [0,1] \subset \mathbb{R}$ with a step size $h$. In $1D$, the Galerkin projection $\mathcal{G}$ and the nodal interpolation operator $\mathcal{I}$ coincide, since there holds the Galerkin orthogonality for $\mathcal{I}$, i.e.

$$a(u - \mathcal{I}u; v) = \int_\Omega D(u - \mathcal{I}u)Dv\,dx = Dv \int_\Omega D(u - \mathcal{I}u)dx = 0 \quad \text{for all} \quad v \in \mathcal{S} \subset \mathcal{P}_1(\mathcal{T}_h) \cap \mathcal{C}(\Omega),$$

i.e. $\mathcal{G} = \mathcal{I}$ by uniqueness of the orthogonal projection.
The idea is to show that the nodal interpolation operator is not continuous with respect to the $L^2$ norm, i.e.

$$\forall C > 0, \ \exists u \in H_D^1(\Omega) \qquad \|\mathcal{I}u\|_{L^2(\Omega)} \not\lesssim \|u\|_{L^2(\Omega)}. \tag{3.7}$$

We argue by contradiction and assume

$$\exists C > 0 \quad \text{such that } \forall u \in H_D^1(\Omega) \ \|h^{-1}\mathcal{I}u\|_{L_2(\Omega)} \leq C\|h^{-1}u\|_{L_2(\Omega)}. \tag{3.8}$$

In the following we introduce the counter example which contradicts (3.8).
Let $\mathcal{T}_h$ be a uniform triangulation of interval $\Omega$ with the step size $h$ and let the function $u(x)$ be defined by

$$u(x) := \left(\frac{2}{h}\left(x - \frac{(2k-1)h}{2}\right)\right)^{2n} \quad \text{for all} \quad (k-1)h \leq x \leq kh, \quad k = 1, \ldots, n, \quad n = 1/h.$$

It is clear that $u(x) \in H_D^1(\Omega)$ with $\Gamma_D = \{0\}$. Note that the case $\Gamma_D = \Gamma$ can be analysed using the similar arguments and is therefore omitted in the following.
From the definition of the nodal interpolation we have

$$\|\mathcal{I}u\|_{L_2(\Omega)} = \Big( \int_0^h \Big(\frac{x}{h}\Big)^2 dx + \int_h^1 1 dx \Big)^{1/2} = \Big(1 - \frac{2h}{3}\Big)^{1/2} = \Big(1 - \frac{2}{3n}\Big)^{1/2} \xrightarrow{n\to\infty} 1. \qquad (3.9)$$

On the other hand, for $\|u\|_{L_2(\Omega)}$ we have

$$\|u\|_{L_2(\Omega)} = \Big( \sum_{k=1}^n \int_{T_k} \Big(\frac{2}{h}\Big(x - \frac{(2k-1)h}{2}\Big)\Big)^{4n} dx \Big)^{1/2} = \Big(\frac{2}{h}\Big)^{2n} \Big(\frac{1}{4n+1} \sum_{k=1}^n \Big[\Big(\frac{h}{2}\Big)^{4n+1} + \Big(\frac{h}{2}\Big)^{4n+1}\Big]\Big)^{1/2}$$

$$= \Big(\frac{2}{h}\Big)^{2n} \Big(\frac{2}{4n+1} \sum_{k=1}^n \Big(\frac{h}{2}\Big)^{4n+1}\Big)^{1/2} = \Big(\frac{1}{4n+1}\Big)^{1/2} \xrightarrow{n\to\infty} 0. \qquad (3.10)$$

This shows that the assertion (3.7) is valid. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Definition 3.1.0.4 ($L^2$ spatial projection).** Let $\mathcal{L}$ be the $L^2$ projection onto the space of the spatial discrete functions with respect to domain $\Omega$. Then for each $I_j \in \mathscr{T}$, cf. Definition 2.0.0.8 and corresponding discrete space $\mathcal{S}^j$, cf. (2.1), the restriction of the projection $\mathcal{L}|_{I_j} : L^2(\Omega) \to \mathcal{S}^j$ is defined for each $u \in L^2(\Omega)$ by the relation

$$(\mathcal{L}u; v) = (u; v) \quad \text{for all} \quad v \in \mathcal{S}^j.$$

Moreover, for $\mathcal{S}^j$ consisting of $\mathcal{P}_1$ functions, the projection $\mathcal{L}$ of an arbitrary function $u \in L^2(\Omega)$ can be seen as

$$\mathcal{L}u(x) := \sum_{k=0}^n u_k \phi_k(x) \quad \text{where} \quad (u_0, \dots, u_n) := ((u; \phi_0), \dots, (u; \phi_n)) M^{-1} \in \mathbb{R}^{n+1}.$$

For $\mathcal{C}^1$ functions, we define the projection only interval wise, i.e. for each $T_k \in \mathcal{T}_j$

$$\mathcal{L}u(x(\zeta))|_{T_k} := u_{k,1}\phi_1(\zeta) + \frac{h_k}{2}u_{k,2}\phi_2(\zeta) + u_{k,3}\phi_3(\zeta) + \frac{h_k}{2}u_{k,4}\phi_4(\zeta) \quad \text{such that}$$

$$u_k := \frac{h_k}{2}\Big( \int_{-1}^1 u\phi_1 d\zeta, \frac{h_k}{2}\int_{-1}^1 u\phi_2 d\zeta, \int_{-1}^1 u\phi_3 d\zeta, \frac{h_k}{2}\int_{-1}^1 u\phi_4 d\zeta \Big) M_k^{-1} \in \mathbb{R}^4.$$

Note that $M$ and $M_k$ are mass matrices for linear and cubic space ansatz, defined in (2.52) and (2.59) respectively. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Lemma 3.1.0.7 (Approximation properties of $L^2$ spatial projection).** Given $u \in H_D^1(\Omega) \cap H^s(\Omega)$ and projection $\mathcal{L}$, cf. Definition 3.1.0.4, there exists a constant $C$ independent of $u$ and $h$ such that for $0 \le s \le p+1$ there holds

$$\|u - \mathcal{L}u\|_{L^2(\Omega)} \le C\|h^s D^s u\|_{L^2(\Omega)},$$

where $p=1$ in case of $\mathcal{P}^1$ elements and $p=3$ for $\mathcal{C}^1$ elements in space.

**Proof.** The inequality is obvious due to the properties of $L^2$ projection. Namely,

$$\|u - \mathcal{L}u\|_{L^2(\Omega)} = \inf_{v \in \mathcal{S}^j} \|u - v\|_{L^2(\Omega)} \leq \|u - \mathcal{I}u\|_{L^2(\Omega)}$$

for $\mathcal{I}$ nodal or Hermite cubic interpolant, respectively. Then, if we recall the Lemma 3.1.0.2 and Lemma 3.1.0.5, we may conclude the proof. $\qquad\square$

**Lemma 3.1.0.8 ($H^1$ stability of projection $\mathcal{L}$).** Let $\mathcal{L}$ be the spatial projection from Definition 3.1.0.4. Then, there exists a constant $C$ independent of $u$ and $h$ such that for all $u \in H_D^1(\Omega)$ there holds

$$\|\mathcal{L}u\|_{H^1(\Omega)} \leq C\|u\|_{H^1(\Omega)}.$$

**Proof.** We refer to CROUZEIX-THOMÉE [22, Theorem 2], where the case of (DD) boundary conditions has been considered and note that the similar result can be obtained for (DN) boundary conditions. See also CARSTENSEN [17]. $\qquad\square$

**Definition 3.1.0.5 (Spatial multi projection).** Let $\Pi$ be a spatial multi projection onto the product of two spatial discrete spaces, such that for each interval $I_j \in \mathcal{T}$, cf. Definition 2.0.0.8, and correponding spatial discrete space $\mathcal{S}^j$, see (2.1) there holds

$$\Pi|_{I_j} : \mathcal{H} \to \mathcal{S}^j \times \mathcal{S}^j,$$

where $\mathcal{H}$ is the Hilbert space introduced in (1.16). $\qquad\square$

**Lemma 3.1.0.9.** If $\mathcal{K}_h$ is the discrete operator from Definition 2.3.3.4 and $\mathcal{G}, \mathcal{L}$ are the spatial Galerkin projection and $L^2$ projection defined in Definitions 3.1.0.3 and 3.1.0.7, respectively, then there holds for all $u \in \mathscr{D}(\mathcal{K}_h^{-1})$,

$$-(\mathcal{L}\Delta u; v) = (\mathcal{K}_h^{-1}\mathcal{G}u; v) \quad \text{for all} \quad v \in \mathcal{S}, \tag{3.11a}$$

where $\mathcal{S}$ is the spatial discrete space.
Moreover, there also holds

$$-(\mathcal{L}\Delta u; v) = (\mathcal{K}_h^{-1}u; v) \quad \text{for all} \quad v \in \mathcal{S}. \tag{3.11b}$$

**Proof.** Recall the orthogonality properties of projections $\mathcal{G}, \mathcal{L}$, see Lemma 3.1.0.3 and Lemma 3.1.0.4, respectively. Thereafter we have for all $v \in \mathcal{S}$

$$-(\mathcal{L}\Delta u; v) = -(\Delta u; v) = a(u; v) = a(\mathcal{G}u; v) = (\mathcal{K}_h^{-1}\mathcal{G}u; v).$$

The proof of (3.11b) is then obvious. $\qquad\square$

**Remark 3.1.0.4.** In case of the $dG(q), q = 0, 1$ and $cG(1)$ time discretisation we may choose to work with the spatial meshes which change between neighbouring time slabs. Therefore, we have to take into account that the definition of the certain operators which map is contained in the spatial discrete space $\mathcal{S}$, must be somehow restricted to the time slab. Namely by the definition of the discrete operator $\mathcal{K}_h$, we have to differentiate between $\mathcal{K}_h^j$ and $\mathcal{K}_h^{j-1}$ in case when $\mathcal{S}^{j-1} \neq \mathcal{S}^j$. Same holds for spatial projections $\mathcal{L}, \mathcal{G}$ and interpolant $\mathcal{I}$. $\qquad\square$

**Definition 3.1.0.6 (Jumps in space).** Let $U$ be an arbitrary piecewise affine, globally continuous function with respect to the spatial domain $\Omega$ and triangulation $\mathcal{T}$, cf. Definition 2.0.0.8. Then $[DU]_k$ can be seen as the jump in the first derivative with respect to node $x_k$ such that

$$[DU]_k := DU(x_{k+1}) - DU(x_k) \quad \text{for all} \quad k = 1, \dots m, \tag{3.12}$$

where $m = n-1$ for (DD) boundary conditions. Additionally, if the initial problem satisfies the (DN) boundary conditions, then $m = n$ and $[DU]_n := -DU(1)$. $\qquad\square$

**Definition 3.1.0.7.** For all $v \in \mathcal{S}$, where $\mathcal{S}$ is discrete space with respect to domain $\Omega$ with a basis which consists of $\mathcal{P}_1$ conform elements and the step size $h$, cf. Definition 2.0.0.8, let $D_{h,\ell}$ denote a discrete counterpart of $\|h^\ell \Delta v\|_{L^2(\Omega)}$ such that

$$D_{h,\ell}(v) := \Big( \sum_{k=1}^m (h_k)^\ell ([Dv]_k)^2 \Big)^{1/2} \quad \text{for} \quad \ell = 1, 2,$$

where $m = n$ for (DN) boundary conditions and $m = n-1$ for (DD) boundary conditions. The jump terms $([Dv]_k)_{k=1}^m$ are defined according to Definition 3.1.0.6. $\qquad\square$

**Lemma 3.1.0.10.** For $D_{h,\ell}$ from Definition 3.1.0.7 and $u \in H_D^1(\Omega) \cap H^\ell(\Omega), \ell = 1, 2$, there holds

$$a(v; u - \mathcal{L}u) \le C D_{h,2\ell-1}(v) \|D^\ell u\|_{L^2(\Omega)}, \qquad \text{for all} \quad v \in \mathcal{S}. \tag{3.13}$$

**Proof.** An integration by parts in space and the fact that $u - \mathcal{L}u \in H_D^1(\Omega)$ yields

$$a(v; u - \mathcal{L}u) = -(\Delta v; u - \mathcal{L}u) - \sum_{k=1}^m [Dv]_k (u - \mathcal{L}u)(x_k), \tag{3.14}$$

where $m = n$ for (DN) and $m = n-1$ for (DD) boundary conditions. Obviously, from the fact that $v$ is $\mathcal{P}^1$ function, we have $\Delta v = 0$. Then, the RHS of the equation (3.14) can be estimated by means of the trace inequality such that

$$a(v; u - \mathcal{L}u) \le -\sum_{k=1}^m [Dv]_k \big( h_k^{-1/2} \|u - \mathcal{L}u\|_{L^2(T_k)} + h_k^{1/2} \|u - \mathcal{L}u\|_{H^1(T_k)} \big).$$

Furthermore, owing to the approximation properties of the projection $\mathcal{L}$, see Lemma 3.1.0.7 and the $H^1$ stability of the same, see Lemma 3.1.0.8, we have by applying the discrete Cauchy inequality

$$a(v; u - \mathcal{L}u) \le C \sum_{k=1}^m [Dv]_k h_k^{(2\ell-1)/2} \|D^\ell u\|_{L^2(T_k)} \le C \Big( \sum_{k=1}^m (h_k)^{2\ell-1} ([Dv]_k)^2 \Big)^{1/2} \|D^\ell u\|_{L^2(\Omega)}.$$

From the Definition of the discrete Laplace operator $D_{h,2\ell-1}$, cf. Definition 3.1.0.7, we conclude the proof. $\qquad\square$

**Definition 3.1.0.8 (Integral mean).** For some $u \in L_2(\mathscr{T})$, where $\mathscr{T}$ is some arbitrary tri-angulation of the interval $[0,T]$, we define the piecewise constant integral mean of the function $u$ by $\bar{u}$ such that

$$\bar{u}|_{I_j} := \fint_{I_j} u(t)dt \quad \text{for each} \quad I_j \in \mathscr{T}. \qquad \square$$

**Definition 3.1.0.9 (Temporal projection).** Let $\mathcal{J}$ be the projection operator on the space of the piecewise polynomials $\mathcal{P}_q(\mathscr{T})$, $0 \leq q \leq 1$ with respect to the arbitrary triangulation $\mathscr{T}$ of interval $[0,T]$, cf. Definition 2.0.0.8, defined for each time interval $I_j \in \mathscr{T}$ by

$$\int_{I_j} (\mathcal{J}u - u)V\, dt = 0 \quad \text{for all} \quad V \in \mathcal{P}_0(I_j). \qquad (3.15)$$

We define $\mathcal{J}$ only for when the time is approximated by the Galerkin discretisation method.

| $cG(1),\ dG(0)$ | $\mathcal{J}\|_{I_j} : L^2(I_j) \to \mathcal{P}_0(I_j)$ | $\mathcal{J}\|_{I_j} u := \fint_{I_j} u(t)dt,$ |
|---|---|---|
| $dG(1)a$ | | $\mathcal{J}u(t) := u(t_{j-1}) + 2(t - t_{j-1})/k_j^2 \int_{I_j} u(s) - u(t_{j-1})ds$ |
| | $\mathcal{J} : L^2(I_j) \to \mathcal{P}_1(I_j),$ | |
| $dG(1)b$ | | $\mathcal{J}u(t) := u(t_j) + 2(t_j - t)/k_j^2 \int_{I_j} u(s) - u(t_j)ds.$ |

Obviously, in case of the $dG(1)$ time approximation, for $dG(1)a$ variant of the operator $\mathcal{J}$ we have $(\mathcal{J}u)^{j-1+} = u^{j-1+}$. Similarly, in case of $dG(1)b$ there holds $(\mathcal{J}u)^{j-} = u^{j-}$. $\qquad \square$

**Remark 3.1.0.5.** The projection operator $\mathcal{J}$ in case of $dG(0)$ time approximation is orthogonal $L^2$ projection. $\qquad \square$

**Lemma 3.1.0.11 (Approximation properties of temporal projection $\mathcal{J}$).** For projection $\mathcal{J}$ from the Definition 3.1.0.9 and all $u \in H^\ell(I_j)$ where $\ell = 1$ in case of $cG(1)$ and $dG(0)$ and $\ell = 2$ in case of $dG(1)$ approximation in time, there is a constant $C$ independent of the triangulation $\mathscr{T}$ such that there holds

$$\|u - \mathcal{J}u\|_{L^\infty(0,T)} \leq C \|k^s \frac{\partial^s}{\partial t} u\|_{L^\infty(0,T)}, \quad 0 \leq s \leq \ell. \qquad (3.16a)$$

Moreover, there also holds

$$\|u - \mathcal{J}u\|_{L^1(0,T)}^2 \leq C \|k^s \frac{\partial^s}{\partial t} u\|_{L^1(0,T)}^2, \quad 1 \leq s \leq \ell, \qquad (3.16b)$$

$$\|u - \mathcal{J}u\|_{L^2(0,T)}^2 \leq C \|k^s \frac{\partial^s}{\partial t} u\|_{L^2(0,T)}^2, \quad 1 \leq s \leq \ell. \qquad (3.16c)$$

**Proof (sketch).** Follows from the Taylor expansions of the function $u$, where for each $I_j \in \mathscr{T}$ we have

$$u(t) = u(t_{j-1}) + \int_{I_j} \dot{u}(\tau)d\tau.$$

Then for each $t \in I_j$,

$$u(t) - \mathcal{J}u(t) = \int_{I_j} \left( \frac{2}{k_j^2}(t_j - t)(s - t_{j-1}) - \chi_{[t,t_j]} \right) \dot{u}(s) ds \qquad (3.17)$$

where $\chi_{[t,t_n]}$ is the characteristic function of interval $[t, t_n]$. The representation (3.17) is used for the estimate when $s = 1$. For $s = 2$, we make use of the equivalent representation

$$u(t) - \mathcal{J}u(t) = \int_{I_j} \left( (s - t)\chi_{[t,t_j]} - \frac{1}{k_j^2}(t_j - t)(s - t_{j-1})^2 \right) \ddot{u}(s) ds. \qquad (3.18)$$
$\square$

**Remark 3.1.0.6.** Note that the projection $\mathcal{J}$ from Definition 3.1.0.9 is continuous only with respect to $L^\infty$ norm. This follows from the estimate (3.16c). $\square$

**Definition 3.1.0.10 ($H^1$ time projection).** Let $\mathcal{J}_1$ denote the operator of the temporal $H^1$ projection on the space of the piecewise polynomials $\mathcal{P}_q(\mathcal{T})$, $q \geq 1$ with respect to triangulation $\mathcal{T}$ defined by

$$\int_0^T \frac{\partial}{\partial t}\left( \mathcal{J}_1 u - u \right) \frac{\partial}{\partial t} V \, dt = 0 \quad \text{for all} \quad V \in \mathcal{P}_q(\mathcal{T}), \qquad (3.19)$$

with the initial condition $\mathcal{J}_1 u(0) = u(0)$.
Note that by taking $V = t$ in (3.19) for $0 \leq t \leq t_j$ and $V = t_j$ for $t > t_j$, we conclude that

$$\mathcal{J}_1 u(t_j) = u(t_j) \quad \text{for all} \quad 1 \leq j \leq N. \qquad \square$$

**Lemma 3.1.0.12 (Approximation properties of $H^1$ temporal projection).** Let $\mathcal{J}_1$ be the projection from the Definition 3.1.0.10. For all $u \in H^s(0,T)$ there is a constant $C$ independent of the time interval such that there holds for all $0 \leq r \leq 1 \leq s \leq q + 1$,

$$\| u - \mathcal{J}_1 u \|_{H^r(0,T)} \leq C \| k^{s-r} u \|_{H^s(0,T)}.$$

**Proof.** For the proof we recall the proof of Lemma 3.1.0.6 where the approximation properties of Galerkin, i.e. spatial $H^1$ projection are proven. In particular, if $q = 1$, then $\mathcal{J}_1$ can be seen as the nodal interpolant in time defined as in the Definition 3.1.0.1 for which the approximation properties from Lemma 3.1.0.2 hold. $\square$

## 3.2 A priori energy error estimate

In the following we present and analyse the methods for the determination of the a priori error bound $\eta_a$ for the error of the fully (time-space) discrete problem. The error bound $\eta_a := \eta_a(u, h, k)$ is a quantity depending on the mesh data $h, k$ and exact solution $u$ together with its derivatives.
An overview concerning the accuracy order of developed estimates with respect to time and space discretisation is presented in Table 3.2.
To ensure that $\eta_a$ can be taken as an optimal or quasi-optimal bound, we also present the results of some simple numerical experiments, see Figure 3.1–3.3. We plot the error in the energy norm for different choice of $h$ and $k$. Obviously, the $dG(0)$ method provides the

| $\otimes$ | $\mathcal{P_1}$ | $\mathcal{C}^1$ |
|---|---|---|
| **$dG(0)$**<br><br>Subsection 3.2.1.1 | $\mathcal{O}(h+k^{-1/2}h+k)$ | $\mathcal{O}(h^3+k^{-1/2}h^3+k)$ |
| **$dG(1)$**<br><br>Subsection 3.2.1.2 | — | — |
| **$cG(1)$**<br><br>Subsection 3.2.2 | — | $\mathcal{O}(h^3+k^2),\ \varepsilon=0$<br><br>$h-$quasi uniform spatial mesh |
| **$MoL$**<br><br>Subsection 3.2.3 | — | $\mathcal{O}(h^3),\ \varepsilon=0$ or $(DD)$<br><br>$h-$quasi uniform spatial mesh |

Table 3.2: Proven a priori error estimates for $\|e\|_{L^\infty(\mathcal{H})}$ in case of $dG(0)$ and $MoL$ and $\max_{t_j\in\mathcal{T}}\|e(t_j)\|_\mathcal{H}$ in case of $cG(1)$; energy method.

optimal result when expected convergence rate in time is considered, i.e. $\mathcal{O}(k)$. For the space discretisation we have the nearly optimal result $\mathcal{O}(h^p+k^{-1/2}h^p)$ where $p=1,3$ for $\mathcal{P_1},\mathcal{C}^1$ elements respectively, see Figure 3.1. The expected order of accuracy in space is $\mathcal{O}(h^p)$. In case of the $dG(1)$ method in time we did not succeed to derive any a priori estimate by using the avaliable techniques. This is still an open question and discussion is provided in Subsection 3.2.1.2. Also for $cG(1)$ in time and $\mathcal{P_1}$ in space, a lack of continuity in the first derivative caused that the derivation of the a priori error estimate can not be completed. On the other hand, under certain requirements concerning the choice of parameter $\varepsilon$, for $\mathcal{C}^1$ elements in



Figure 3.1: Convergence of $\|e\|_{L^\infty(\mathcal{H})}$ for $dG(0)\otimes\mathcal{P_1}$ (left) and $dG(0)\otimes\mathcal{C}^1$ (right) with respect to the number of elements in space; The a priori error bound $\eta_a=\mathcal{O}(h^p+k^{-1/2}h^p+k)$ is suboptimal for both $p=1,3$. Namely for $p=1$ and $h=k$, the exact error is of order $\mathcal{O}(h)$ whereby $\eta_a=\mathcal{O}(h^{1/2})$. For $\mathcal{C}^1$ elements and $k=h^2$ we have the best convergence order $\eta_a=\mathcal{O}(h^2)$, but still $\|e\|_{L^\infty(\mathcal{H})}=\mathcal{O}(h^3)$. On the other hand, a duality approach, cf. Subsection 4.2.1.1, provides a bound with an optimal convergence order, i.e. $\mathcal{O}(h^p+k)$; Example 6.2.5, $\varepsilon=0$, (DN), $T=1$.

space, the optimal results in time and space are obtained, see Figure 3.2. When the method



Figure 3.2: Convergence of $\max_{t \in \mathcal{T}} \|e(t_j)\|_{\mathcal{H}}$ for $cG(1) \otimes \mathcal{C}^1$ with respect to the number of elements in space; The optimal convergence order of the a priori error estimate $\eta_a = \mathcal{O}(h^3 + k^2)$ with respect to the step size $h$, i.e. $\mathcal{O}(h^3)$ is achieved already by $k = h^{3/2}$; Example 6.2.5, $\varepsilon = 0$, (DN), $T = 1$.

of lines is considered, the same problems occur when $\mathcal{P}_1$ discretisation in space is applied, as it was the case in $cG(1) \otimes \mathcal{P}_1$ discretisation. For $\mathcal{C}^1$ elements, the optimal order of accuracy is obtained under weaker requirements than the $cG(1)$ time approximation allowed. For an example, see Figure 3.3.
Note that improved results are obtained by a different analytical technique, namely by a duality approach, cf. Section 4.2.



Figure 3.3: Convergence of $\|e\|_{L^\infty(\mathcal{H})}$ for $MoL \otimes \mathcal{C}^1$ with respect to the number of elements in space; The approximated error bound $\eta_a = \mathcal{O}(h^3)$ is optimal for arbitrary choice of $k$; Example 6.2.5, $\varepsilon = 0$, (DN), $T = 1$.

### 3.2.1  $dG(q)$ time approximation, $q=0,1$

Recall the notation and definitions related to the discontinuous Galerkin method in time, introduced in Subsection 2.3.3 as well as Section 2.1, where the particular space discretisation methods are introduced, namely linear $(\mathcal{P}_1)$ and cubic $(\mathcal{C}^1)$ elements.

**Lemma 3.2.1.1 (Error representation, $dG(q)$ time approximation).** In case of $dG(q)$ time discretisation, for every $V \in \mathcal{Q}_q$ the following identity is valid

$$\frac{1}{2}\|e^{N-}\|_{\mathcal{H}}^2 + \varepsilon\|e_2\|_{L^2(H^1)}^2 \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \frac{1}{2}\sum_{j=1}^{N}\|(e-V)^{j-1+}\|_{\mathcal{H}}^2 + \sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; e-V\rangle_{\mathcal{H}} dt$$
$$-\sum_{j=1}^{N}\int_{I_j}a(e_2; e_1-V_1)dt + \sum_{j=1}^{N}\int_{I_j}a(e_1; e_2-V_2)dt$$
$$+\varepsilon\sum_{j=1}^{N}\int_{I_j}a(e_2; e_2-V_2)dt. \tag{3.20}$$

**Remark 3.2.1.1.** The error representation (3.20) is also valid for each $1 \leq n \leq N$ with respect to triangulation $\mathscr{T}$, due to the local representation of weak problem, see Remark 2.3.1.1.  $\square$

**Proof.** Since $e$ is discontinuous as a function of time, we derive from Lemma 2.3.3.3 that

$$\sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; e\rangle_{\mathcal{H}} dt = \frac{1}{2}\sum_{j=1}^{N}\|e^{j-}\|_{\mathcal{H}}^2 - \frac{1}{2}\sum_{j=1}^{N}\|e^{j-1+}\|_{\mathcal{H}}^2. \tag{3.21}$$

With the definition (2.18) of the bilinear form $\mathcal{B}$ and the Galerkin orthogonality (2.9), there holds, for all $V \in \mathcal{Q}_q$,

$$\sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; e\rangle_{\mathcal{H}} dt = \sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; e-V\rangle_{\mathcal{H}} dt + \sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; V\rangle_{\mathcal{H}} dt$$
$$= \sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; e-V\rangle_{\mathcal{H}} dt + \sum_{j=1}^{N}\int_{I_j}a(e_2; V_1)dt$$
$$-\sum_{j=1}^{N}\int_{I_j}a(e_2; V_1)dt - \varepsilon\sum_{j=1}^{N}\int_{I_j}a(e_2; V_2)dt - \sum_{j=1}^{N}\langle[e]^{j-1} ; V^{j-1+}\rangle_{\mathcal{H}}.$$

The combination of the last equality and identity (3.21) shows

$$\frac{1}{2}\sum_{j=1}^{N}\|e^{j-}\|_{\mathcal{H}}^2 - \frac{1}{2}\sum_{j=1}^{N}\|e^{j-1+}\|_{\mathcal{H}}^2 = \sum_{j=1}^{N}\int_{I_j}\langle e_\tau ; e-V\rangle_{\mathcal{H}} dt - \sum_{j=1}^{N}\int_{I_j}a(e_2; e_1-V_1)dt$$
$$+\sum_{j=1}^{N}\int_{I_j}a(e_1+\varepsilon e_2; e_2-V_2)dt + \sum_{j=1}^{N}\langle[e]^{j-1} ; (e-V_1)^{j-1+}\rangle_{\mathcal{H}}$$
$$-\varepsilon\sum_{j=1}^{N}\int_{I_j}a(e_2; e_2)dt - \sum_{j=1}^{N}\langle[e]^{j-1} ; e^{j-1+}\rangle_{\mathcal{H}}. \tag{3.22}$$

With Lemma 2.3.3.3 the last jump contribution from (3.22) reads

$$-\sum_{j=1}^{N}\langle [e]^{j-1}\,;e^{j-1+}\rangle_{\mathcal{H}}=-\frac{1}{2}\sum_{j=1}^{N}\||e^{j-1+}\||_{\mathcal{H}}^{2}-\frac{1}{2}\||[e]^{j-1}\||_{\mathcal{H}}^{2}+\frac{1}{2}\sum_{j=1}^{N}\||e^{j-1-}\||_{\mathcal{H}}^{2}.\qquad(3.23)$$

For the first jump contribution in (3.22), the Cauchy inequality yields

$$\sum_{j=1}^{N}\langle [e]^{j-1}\,;(e-V)^{j-1+}\rangle_{\mathcal{H}}\leq\sum_{j=1}^{N}\||[e]^{j-1}\||_{\mathcal{H}}\||(e-V)^{j-1+}\||_{\mathcal{H}}$$

$$\leq\frac{1}{2}\sum_{j=1}^{N}\||[e]^{j-1}\||_{\mathcal{H}}^{2}+\frac{1}{2}\sum_{j=1}^{N}\||(e-V)^{j-1+}\||_{\mathcal{H}}^{2}.\qquad(3.24)$$

With (3.23)-(3.24), the equation (3.22) is equivalent to

$$\frac{1}{2}\||e^{N-}\||_{\mathcal{H}}^{2}+\varepsilon\|e_{2}\|_{L^{2}(H^{1})}^{2}\leq\frac{1}{2}\||e^{0-}\||_{\mathcal{H}}^{2}+\frac{1}{2}\sum_{j=1}^{N}\||(e-V)^{j-1+}\||_{\mathcal{H}}^{2}+\sum_{j=1}^{N}\int_{I_{j}}\langle e_{\tau}\,;e-V\rangle_{\mathcal{H}}dt$$

$$-\sum_{j=1}^{N}\int_{I_{j}}a(e_{2};e_{1}-V_{1})dt+\sum_{j=1}^{N}\int_{I_{j}}a(e_{1}+\varepsilon e_{2};e_{2}-V_{2})dt.\qquad(3.25)$$

This completes the proof of lemma. □

In the following the time approximation methods, in this case $dG(0)$ and $dG(1)$ method, have to be analysed separately. The space approximation assumes the use of linear $(\mathcal{P}_{1})$ and Hermite cubic splines $(\mathcal{C}^{1})$.

### 3.2.1.1 $dG(0)$ time approximation

**Theorem 3.2.1.1 (A priori energy error estimate, $dG(0)$ time approximation).** There exists a constant $C$, which is independent of $u$ and its $dG(0)$ approximant, such that

$$\|e\|_{L^{\infty}(\mathcal{H})}\leq C\Big\{\|k\dot{u}_{1}\|_{L^{\infty}(H^{1})}+\|k\dot{u}_{2}\|_{L^{\infty}(L^{2})}+\|h^{p}D^{p+1}y_{0}\|_{L^{2}(\Omega)}+\|h^{p+1}D^{p+1}y_{1}\|_{L^{2}(\Omega)}$$

$$+\Big(\sum_{j=1}^{N}\|h^{p}D^{p+1}u_{1}(t_{j})\|_{L^{2}(\Omega)}^{2}\Big)^{1/2}+\Big(\sum_{j=1}^{N}\|h^{p+1}D^{p+1}u_{2}(t_{j})\|_{L^{2}(\Omega)}^{2}\Big)^{1/2}$$

$$+\|k\Delta\dot{u}_{1}\|_{L^{1}(L^{2})}+\|h^{p}D^{p+1}u_{2}\|_{L^{1}(L^{2})}+\|k\dot{u}_{2}\|_{L^{1}(H^{1})}+\sqrt{\varepsilon}\|k\dot{u}_{2}\|_{L^{2}(H^{1})}$$

$$+\sqrt{\varepsilon}\|h^{p}D^{p+1}u_{2}\|_{L^{2}(L^{2})}+\sqrt{\varepsilon}\|k\Delta\dot{u}_{2}\|_{L^{2}(L^{2})}\Big\},\qquad(3.26)$$

provided that $u\in H^{p+1}(\Omega)^{2}$. Here $p=1$ when $\mathcal{P}_{1}$ and $p=3$ when $\mathcal{C}^{1}$ ansatz in space is applied. The last term on the RHS of (3.26) does not appear if the initial problem satisfies the (DD) boundary conditions.

**Remark 3.2.1.2.** Note that the estimate of Theorem 3.2.1.1 is of order $\mathcal{O}(h^{p}+k^{-1/2}h^{p}+k)$. □

**Proof.** In case of $dG(0)$ time approximation the third term on the RHS of the error representation (3.20) can be rewritten as

$$\sum_{j=1}^{N}\int_{I_j}\langle e_\tau\,;e-V\rangle_{\mathcal{H}}dt=\sum_{j=1}^{N}\int_{I_j}\langle (e-V)_\tau\,;e-V\rangle_{\mathcal{H}}dt$$

$$=\frac{1}{2}\sum_{j=1}^{N}|\!|\!|(e-V)^{j-}|\!|\!|_{\mathcal{H}}^2-\frac{1}{2}\sum_{j=1}^{N}|\!|\!|(e-V)^{j-1+}|\!|\!|_{\mathcal{H}}^2. \qquad (3.27)$$

Inserting (3.27) into (3.20), we estimate the energy error and dissipative term as follows

$$\frac{1}{2}\|e^{N-}\|_{\mathcal{H}}^2+\varepsilon\|e_2\|_{L^2(H^1)}^2\leq\frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2+\frac{1}{2}\sum_{j=1}^{N}|\!|\!|(e-V)^{j-}|\!|\!|_{\mathcal{H}}^2-\sum_{j=1}^{N}\int_{I_j}a(e_2;e_1-V_1)dt$$

$$+\sum_{j=1}^{N}\int_{I_j}a(e_1;e_2-V_2)dt+\varepsilon\sum_{j=1}^{N}\int_{I_j}a(e_2;e_2-V_2)dt=:\sum_{\ell=1}^{5}E_\ell. \qquad (3.28)$$

The idea in the following is to estimate the contributions $E_1,\ldots,E_5$ on the RHS of (3.28) with respect to $\|e_1\|_{L^\infty(H^1)},\|e_2\|_{L^\infty(L^2)},\ \sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}$ and some terms which are priori known. This is done in several steps formulated as Lemmas. Before we skip to the proof of the same, we define a test function $V$ by

$$V|_{I_j}:=\Pi e^{j-}\in\mathcal{Q}_0^j,\quad \Pi=(\mathcal{G},\mathcal{G}). \qquad (3.29)$$

This choice of $V$ applies to the rest of the proof. We also use $p=1$ for $\mathcal{P}^1$ and $p=3$ for $\mathcal{C}^1$ elements in space in the following analysis.

**Lemma 3.2.1.2.** Let the discrete initial data be defined such that

$$U^{0-}:=(\mathcal{I}y_0,\mathcal{I}y_1) \qquad (3.30)$$

for the interpolation operator $\mathcal{I}$ corresponding to the space discretisation ansatz and $u_0=(y_0,y_1)$ the initial solution, then there exists a constant $C$ such that

$$E_1\leq C\big\{\|h^pD^{p+1}y_0\|_{L^2(\Omega)}^2+\|h^{p+1}D^{p+1}y_1\|_{L^2(\Omega)}^2\big\}.$$

**Proof.** The approximation properties of the interpolant $\mathcal{I}$ given in Lemma 3.1.0.2 and 3.1.0.5 for $\mathcal{P}_1$ and $\mathcal{C}^1$ elements in space, respectively, imply

$$E_1=\frac{1}{2}\|y_0-\mathcal{I}y_0\|_{H^1(\Omega)}^2+\frac{1}{2}\|y_1-\mathcal{I}y_1\|_{L^2(\Omega)}^2$$

$$\leq C\big\{\|h^pD^{p+1}y_0\|_{L^2(\Omega)}^2+\|h^{p+1}D^{p+1}y_1\|_{L^2(\Omega)}^2\big\}. \qquad \square$$

**Lemma 3.2.1.3.** There is a constant $C$ such that

$$E_2\leq C\Big\{\sum_{j=1}^{N}\|h^pD^{p+1}u_1(t_j)\|_{L^2(\Omega)}^2+\sum_{j=1}^{N}\|h^{p+1}D^{p+1}u_2(t_j)\|_{L^2(\Omega)}^2\Big\}.$$

**Proof.** Owing to the approximation properties from Lemmas 3.1.0.6 and 3.1.0.7 there holds

$$E_2 = \frac{1}{2} \sum_{j=1}^{N} \|(e - \Pi e)^{j-}\|_{\mathcal{H}}^2$$

$$= \frac{1}{2} \sum_{j=1}^{N} \|(e_1 - \mathcal{G}e_1)^{j-}\|_{H^1(\Omega)}^2 + \frac{1}{2} \sum_{j=1}^{N} \|(e_2 - \mathcal{G}e_2)^{j-}\|_{L^2(\Omega)}^2$$

$$\leq C \left\{ \sum_{j=1}^{N} \|h^p D^{p+1} u_1(t_j)\|_{L^2(\Omega)}^2 + \sum_{j=1}^{N} \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)}^2 \right\}.$$

This concludes the proof. □

**Lemma 3.2.1.4.** For $E_3$ there is a constant $C$ such that

$$E_3 \leq \|k\Delta \dot{u}_1\|_{L^1(L^2)} \|e_2\|_{L^\infty(L^2)} + \|k\dot{u}_2\|_{L^2(H^1)} \varepsilon \|e_2(t)\|_{L^2(H^1)}$$
$$+ \|k\Delta \dot{u}_2\|_{L^2(L^2)} \varepsilon \|e_2(t)\|_{L^2(H^1)} + C \|h^p D^{p+1} u_2\|_{L^1(L^2)} \|e_1\|_{L^\infty(H^1)}. \tag{3.31}$$

The second and the third term on the RHS of (3.31) does not appear in case of the (DD) boundary conditions.

**Proof.** In order to simplify the estimation, we decompose $E_3$ such that

$$E_3 = -\sum_{j=1}^{N} \int_{I_j} a(e_2; e_1 - e_1^{j-}) dt - \sum_{j=1}^{N} \int_{I_j} a(e_2; e_1^{j-} - \mathcal{G}e_1^{j-}) dt. \tag{3.32}$$

The first term on the RHS of (3.32) is independent of the choice of the multi projection $\Pi$. Since for all $t \in I_j$, $e_1(t) - e_1^{j-} = u(t) - u(t_j)$, an integration by parts in space yields

$$-\sum_{j=1}^{N} \int_{I_j} a(e_2; e_1 - e_1^{j-}) dt = \sum_{j=1}^{N} \int_{I_j} (e_2; \Delta(u_1 - u_1(t_j))) dt - \sum_{j=1}^{N} \int_{I_j} e_2(t, 1) D(u_1(t, 1) - u_1(t_j, 1)) dt, \tag{3.33}$$

where the second term on the RHS of the equation above equals zero in case of the (DD) boundary conditions.

On account of the properties of $dG(0)$ functions, the following identity is valid

$$e(s) = e^{j-} + \int_s^{t_j} e_\tau(\tau) d\tau = e^{j-} + \int_s^{t_j} \dot{u}(t) dt \quad \text{for all} \quad s \in I_j. \tag{3.34}$$

Thereafter, we may easily conclude that

$$\|\Delta(e_1 - e_1^{j-})\|_{L^1(I_j; L^2(\Omega))} \leq \int_{I_j} \int_s^{t_j} \|\Delta \dot{u}_1(t)\| dt \leq k_j \|\Delta \dot{u}_1\|_{L^1(I_j; L^2(\Omega))}. \tag{3.35}$$

From (3.35) the first term on the RHS of (3.33) can be estimated by

$$\sum_{j=1}^{N} \int_{I_j} (e_2; \Delta(e_1 - e_1^{j-})) dt \leq \sum_{j=1}^{N} \int_{I_j} \|e_2\|_{L^2(\Omega)} \|\Delta(e_1 - e_1^{j-})\|_{L^2(\Omega)}$$

$$\leq \|e_2\|_{L^\infty(L^2)} \sum_{j=1}^{N} \int_{I_j} \|\Delta(e_1 - e_1^{j-})\|_{L^2(\Omega)}$$

$$\leq \|e_2\|_{L^\infty(L^2)} \|k\Delta \dot{u}_1\|_{L^1(L^2)}. \tag{3.36}$$

In case of the of (DN) boundary conditions, the equation (1.13c) allows us to reformulate the second term on the RHS of (3.33) such that

$$-\sum_{j=1}^{N}\int_{I_j}e_2(t,1)D(u_1(t,1)-u_1(t_j,1))dt=\sum_{j=1}^{N}\int_{I_j}\varepsilon e_2(t,1)D(u_2(t,1)-u_2(t_j,1))dt \qquad (3.37)$$

where from $e_2(t,0)=0$

$$e_2(t,1)=e_2(t,0)+\int_{\Omega}De_2(t,x)dx=\|e_2(t)\|_{H^1(\Omega)}. \qquad (3.38)$$

An application of the trace inequality with respect to the whole interval $\Omega=(0,1)$, see Lemma 3.1.0.3 and identity (3.34), which is also valid for $u_2$ instead of $e$, leads to

$$D(u_2(t,1)-u_2(t_j,1))\leq \|u_2(t)-u_2(t_j)\|_{H^1(\Omega)}+\|\Delta(u_2(t)-u_2(t_j))\|_{L^2(\Omega)}$$
$$\leq k_j^{1/2}\|\dot u_2(t)\|_{L^2(I_j;H^1(\Omega))}+k_j^{1/2}\|\Delta\dot u_2(t)\|_{L^2(I_j;L^2(\Omega))}. \qquad (3.39)$$

With (3.38) and (3.39), the RHS of (3.37) can be estimated by

$$-\sum_{j=1}^{N}\int_{I_j}e_2(t,1)D(u_1(t,1)-u_1(t_j,1))dt\leq\varepsilon\sum_{j=1}^{N}k_j^{1/2}\big(\|\dot u_2(t)\|_{L^2(I_j;H^1(\Omega))}+\|\Delta\dot u_2(t)\|_{L^2(I_j;L^2(\Omega))}\big)\int_{I_j}\|e_2(t)\|_{H^1(\Omega)}$$
$$\leq\varepsilon\sum_{j=1}^{N}k_j\big(\|\dot u_2(t)\|_{L^2(I_j;H^1(\Omega))}+\|\Delta\dot u_2(t)\|_{L^2(I_j;L^2(\Omega))}\big)\|e_2\|_{L^2(I_j;H^1(\Omega))}.$$
$$\leq\big(\|k\dot u_2\|_{L^2(H^1)}+\|k\Delta\dot u_2\|_{L^2(L^2)}\big)\varepsilon\|e_2(t)\|_{L^2(H^1)} \qquad (3.40)$$

where we have used a discrete Cauchy inequality in the last estimate.
From (3.36) and (3.40) we obtain

$$-\sum_{j=1}^{N}\int_{I_j}a(e_2;e_1-e_1^{j-})dt\leq\|k\Delta\dot u_1\|_{L^1(L^2)}\|e_2\|_{L^\infty(L^2)}$$
$$+\big(\|k\dot u_2\|_{L^2(H^1)}+\|k\Delta\dot u_2\|_{L^2(L^2)}\big)\varepsilon\|e_2(t)\|_{L^2(H^1)}. \qquad (3.41)$$

This completes the estimation of the first term on the RHS of (3.32).
On account of the properties of the Galerkin projection, cf. Lemma 3.1.0.6, the second term from (3.32) can be estimated by

$$-\sum_{j=1}^{N}\int_{I_j}a(e_2;e_1^{j-}-\mathcal{G}e_1^{j-})dt=-\sum_{j=1}^{N}\int_{I_j}a(e_2-\mathcal{G}e_2;e_1^{j-})dt\leq C\|h^pD^{p+1}u_2\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)}. \quad (3.42)$$

The combination of (3.41) and (3.42) concludes the proof. $\qquad\qquad\Box$

**Lemma 3.2.1.5.** There exists a constant $C$ such that

$$E_4\leq C\|h^pD^{p+1}u_2\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)}+\|k\dot u_2\|_{L^1(H^1)}\|e_1\|_{L^\infty(H^1)}.$$

**Proof.** We start from the following decomposition

$$E_4 = \sum_{j=1}^{N} \int_{I_j} a(e_1; e_2 - \mathcal{G}e_2)dt + \sum_{j=1}^{N} \int_{I_j} a(e_1; \mathcal{G}e_2 - \mathcal{G}e_2^{j-})dt. \tag{3.43}$$

The first term on the RHS of (3.43) can be estimated such that

$$\sum_{j=1}^{N} \int_{I_j} a(e_1; e_2 - \mathcal{G}e_2)dt \le C\|h^p D^{p+1} u_2\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)}. \tag{3.44}$$

In case of the second term, from the identity (3.34) applied for $u_2$, there holds

$$\sum_{j=1}^{N} \int_{I_j} a(e_1; \mathcal{G}e_2 - \mathcal{G}e_2^{j-})dt \le \|e_1\|_{L^\infty(H^1)} \sum_{j=1}^{N} \|u_2 - u_2(t_j)\|_{L^1(I_j; H^1(\Omega))}$$

$$\le \|k\dot{u}_2\|_{L^1(H^1)}\|e_1\|_{L^\infty(H^1)}. \tag{3.45}$$

The latter estimations combined with (3.43) complete the proof. □

**Lemma 3.2.1.6.** For $E_5$ there holds

$$E_5 \le C\varepsilon\|h^p D^{p+1} u_2\|_{L^2(L^2)}\|e_2\|_{L^2(H^1)} + \varepsilon\|k\dot{u}_2\|_{L^2(H^1)}\|e_2\|_{L^2(H^1)},$$

where $C$ is some positive constant.

**Proof.** We start from

$$E_5 = \varepsilon \sum_{j=1}^{N} \int_{I_j} a(e_2; e_2 - \mathcal{G}e_2)dt + \varepsilon \sum_{j=1}^{N} \int_{I_j} a(e_2; \mathcal{G}e_2 - \mathcal{G}e_2^{j-})dt. \tag{3.46}$$

Arguing as in $E_4$ for the first term on the RHS of (3.46), we obtain

$$\varepsilon \sum_{j=1}^{N} \int_{I_j} a(e_2; e_2 - \mathcal{G}e_2)dt \le C\varepsilon\|h^p D^{p+1} u_2\|_{L^2(L^2)}\|e_2\|_{L^2(H^1)}.$$

By use of (3.34), there holds

$$\varepsilon \sum_{j=1}^{N} \int_{I_j} a(e_2; \mathcal{G}e_2 - \mathcal{G}e_2^{j-})dt \le \varepsilon\|k\dot{u}_2\|_{L^2(H^1)}\|e_2\|_{L^2(H^1)}.$$

A substitution of the last two estimates in (3.46) yields the proof. □

We continue with the proof of Theorem 3.2.1.1.

The combination of Lemma 3.2.1.2–Lemma 3.2.1.6, the error representation (3.28) and the Young inequality proves the following result

$$\|e^{N-}\|_{\mathcal{H}}^2 + \varepsilon \|e_2\|_{L^2(H^1)}^2 \leq C\Big\{ \|h^p D^{p+1} y_0\|_{L^2(\Omega)}^2 + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)}^2$$
$$+ \sum_{j=1}^N \|h^p D^{p+1} u_1(t_j)\|_{L^2(\Omega)}^2 + \sum_{j=1}^N \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)}^2$$
$$+ \|k \Delta \dot{u}_1\|_{L^1(L^2)} \|e_2\|_{L^\infty(L^2)} + \|h^p D^{p+1} u_2\|_{L^1(L^2)} \|e_1\|_{L^\infty(H^1)}$$
$$+ \|k \dot{u}_2\|_{L^1(H^1)} \|e_1\|_{L^\infty(H^1)} + \varepsilon \|k \dot{u}_2\|_{L^2(H^1)}^2$$
$$+ \varepsilon \|h^p D^{p+1} u_2\|_{L^2(L^2)} + \varepsilon \|k \Delta \dot{u}_2\|_{L^2(L^2)}^2 \Big\}. \tag{3.47}$$

According to Lemma 3.2.1.3, the last term on the RHS of (3.47) does not appear if the initial problem satisfies the (DD) boundary conditions.

Finally, we introduce the lemma which provides the relation between $\|e\|_{L^\infty(\mathcal{H})}$ and $\|e^{N-}\|_{\mathcal{H}}$. The motivation for the derivation of such a result stems from the fact that the RHS of the estimate (3.47) depends on $\|e\|_{L^\infty(\mathcal{H})}$.

**Lemma 3.2.1.7.** Since $e$ is piecewise continuous function in time with respect to triangulation $\mathcal{T}$, there is a $t \in [0, T]$ with $\|e(t)\|_{\mathcal{H}} = \|e\|_{L^\infty(\mathcal{H})}$. Let $I_j$ be a time interval in $\mathcal{T}$ with $t \in I_j$. Then there holds

$$\|e\|_{L^\infty(\mathcal{H})} \leq \|e^{j-}\|_{\mathcal{H}} + \|k \dot{u}\|_{L^\infty(\mathcal{H})}. \tag{3.48}$$

**Proof.** For $t \in I_j$, the representation (3.34) implies

$$\|e(t)\|_{\mathcal{H}} \leq \|e^{j-}\|_{\mathcal{H}} + \int_{I_j} \|\dot{u}(\tau)\|_{\mathcal{H}} d\tau \leq \|e^{j-}\|_{\mathcal{H}} + k_j \|\dot{u}\|_{L^\infty(I_j;\mathcal{H})} \leq \|e^{j-}\|_{\mathcal{H}} + \|k \dot{u}\|_{L^\infty(\mathcal{H})}.$$

This completes the proof. □

To conclude the proof of the Theorem 3.2.1.1, we recall that the estimate (3.47) can be derived for each $t_j$ where $1 \leq j \leq N$ on the basis of Remark 3.2.1.1 which also applies to the error representation (3.28). Then, the RHS of (3.47) is a global bound for each $j = 1, \ldots, N$ instead for only $j = N$. From Lemma 3.2.1.7 we have

$$\|e\|_{L^\infty(\mathcal{H})}^2 \leq C\Big\{ \|k \dot{u}\|_{L^\infty(\mathcal{H})}^2 + \|h^p D^{p+1} y_0\|_{L^2(\Omega)}^2 + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)}^2 + \sum_{j=1}^N \|h^p D^{p+1} u_1(t_j)\|_{L^2(\Omega)}^2$$
$$+ \sum_{j=1}^N \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)}^2 + \|k \Delta \dot{u}_1\|_{L^1(L^2)} \|e_2\|_{L^\infty(L^2)}$$
$$+ \|h^p D^{p+1} u_2\|_{L^1(L^2)} \|e_1\|_{L^\infty(H^1)} + \|k \dot{u}_2\|_{L^1(H^1)} \|e_1\|_{L^\infty(H^1)} + \varepsilon \|k \dot{u}_2\|_{L^2(H^1)}^2$$
$$+ \varepsilon \|h^p D^{p+1} u_2\|_{L^2(L^2)}^2 + \varepsilon \|k \Delta \dot{u}_2\|_{L^2(L^2)}^2 \Big\}. \tag{3.49}$$

With the help of the Young inequality we can absorb the terms $\|e_1\|_{L^\infty(H^1)}$ and $\|e_2\|_{L^\infty(L^2)}$ on the RHS of (3.49) through the LHS. i.e. $\|e\|_{L^\infty(\mathcal{H})}$. This completes the proof of theorem. □

### 3.2.1.2 *dG*(1) time approximation

Our intention within this section was to derive an a priori error estimate for fully discrete problem in case of the $dG(1)$ time approximation and $\mathcal{P}_1$ or $\mathcal{C}^1$ discretisation in space, by using the similar arguments as for the $dG(0)$ method in time, see Subsection 3.2.1.1. This needs bounds for each $\|e^{j-}\|_{\mathcal{H}}$, $j = 1, \ldots, N$ in terms of $\|e\|_{L^\infty(\mathcal{H})}$ and some a priori known terms, see (3.47). Apart from this estimate, which is very difficult to obtain for $dG(1)$ functions, we would need to establish a relation between $\|e\|_{L^\infty(\mathcal{H})}$ and $\|e^{j-}\|_{\mathcal{H}}$ in a form

$$\|e\|_{L^\infty(\mathcal{H})} \leq \|e^{j-}\|_{\mathcal{H}} + g(u), \tag{3.50}$$

where $g$ is a function depending on the exact solution and its derivatives. This relation was the reference argument in the a priori analysis for the $dG(0)$ approximation with $g(u) := \|k\dot{u}\|_{L^\infty(\mathcal{H})}$, see Lemma 3.2.1.7 in the proof of Theorem 3.2.1.1.
However, we did not succeed to derive the estimate of type (3.50) in this setting.

Another possibility would be to find some error function which provides the relation (3.50) directly owing to its properties. Obviously, if we define $\tilde{\tilde{e}}$ by

$$\tilde{\tilde{e}} := \widetilde{U} - \widetilde{U} \tag{3.51}$$

then $\tilde{\tilde{e}}$ is a globally continuous and piecewise affine function in time and therefore there exists some $j$, $1 \leq j \leq N$ such that

$$\|\tilde{\tilde{e}}\|_{L^\infty(\mathcal{H})} = \|\tilde{\tilde{e}}(t_j)\|_{\mathcal{H}} = \|e^{j-}\|_{\mathcal{H}}. \tag{3.52}$$

However, the application of the Galerkin orthogonality yields an additional jump term on the RHS of the error representation. We did neither succeed to absorb this additional term nor could we transfer it into proper a priori term.

## 3.2.2  *cG*(1) time approximation

Within this section we prove an a priori error estimate for the $cG(1)$ time discretisation and $\mathcal{C}^1$ discretisation in space when $\varepsilon = 0$ and the spatial mesh is $h-$quasi uniform.

The main argument in the following analysis is the use of the error function $\tilde{\tilde{e}}$ defined by

$$\tilde{\tilde{e}} := \tilde{u} - U \tag{3.53}$$

which is an affine, globally continuous function in time.
The idea to use $\tilde{\tilde{e}}$ instead of $e$ arises from the fact that $\tilde{\tilde{e}}$ reaches its maximum in some node of the triangulation $\mathcal{T}$, whereby for $e$ we can not determine whether the maximum point is an interior point of some interval from $\mathcal{T}$ or not.
For notation and definitions used in the following analysis, we refer to Chapter 2, where the particular space discretisation methods are introduced, cf. Section 2.1 and Subsection 2.3.4 where $cG(1)$ method is analysed.

**Lemma 3.2.2.1 (Error representation, $cG(1)$ time approximation).** In case of $cG(1)$ discretisation in time, for every $V \in \mathcal{W}_c$, the following identity is valid

$$
\begin{aligned}
\frac{1}{2}\|e(T)\|_{\mathcal{H}}^2 + \varepsilon\|\tilde{\tilde{e}}_2\|_{L^2(H^1)}^2 = {}& \frac{1}{2}\|e(0)\|_{\mathcal{H}}^2 + \int_0^T \langle \dot{\tilde{\tilde{e}}} - \dot{e}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt + \int_0^T \langle \dot{e}\,;\tilde{\tilde{e}}-V\rangle_{\mathcal{H}}\,dt \\
& - \int_0^T a(e_2;\tilde{\tilde{e}}_1 - V_1)\,dt + \int_0^T a(e_1;\tilde{\tilde{e}}_2 - V_2)\,dt + \varepsilon\int_0^T a(\tilde{\tilde{e}}_2;\tilde{\tilde{e}}_2 - V_2)\,dt \\
& + \int_0^T a(e_2 - \tilde{\tilde{e}}_2;\tilde{\tilde{e}}_1)\,dt + \int_0^T a(\tilde{\tilde{e}}_1 - e_1;\tilde{\tilde{e}}_2)\,dt \\
& + \varepsilon\int_0^T a(\tilde{\tilde{e}}_2 - e_2;\tilde{\tilde{e}}_2)\,dt.
\end{aligned}
\tag{3.54}
$$

**Proof.** From the fundamental theorem of calculus in time and the fact that the errors $e$ and $\tilde{\tilde{e}}$ coincide in each node $t_j, j=1,\dots,N$ we have

$$
\frac{1}{2}\|e(T)\|_{\mathcal{H}}^2 - \frac{1}{2}\|e(0)\|_{\mathcal{H}}^2 = \int_0^T \langle \dot{\tilde{\tilde{e}}}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt.
\tag{3.55}
$$

From definition of the bilinear form $\mathcal{B}$, (2.41), and the Galerkin orthogonality (2.9), we infer for all $V \in \mathcal{W}_c$

$$
\begin{aligned}
\int_0^T \langle \dot{\tilde{\tilde{e}}}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt = {}& \int_0^T \langle \dot{\tilde{\tilde{e}}} - \dot{e}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt + \int_0^T \langle \dot{e}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt \\
= {}& \int_0^T \langle \dot{\tilde{\tilde{e}}} - \dot{e}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt + \int_0^T \langle \dot{e}\,;\tilde{\tilde{e}}-V\rangle_{\mathcal{H}}\,dt + \int_0^T a(e_2;V_1)\,dt - \int_0^T a(e_1 + \varepsilon e_2;V_2)\,dt \\
= {}& \int_0^T \langle \dot{\tilde{\tilde{e}}} - \dot{e}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt + \int_0^T \langle \dot{e}\,;\tilde{\tilde{e}}-V\rangle_{\mathcal{H}}\,dt - \int_0^T a(e_2;\tilde{\tilde{e}}_1 - V_1)\,dt \\
& + \int_0^T a(e_1 + \varepsilon e_2\,\tilde{\tilde{e}}_2 - V_2)\,dt + \int_0^T a(e_2;\tilde{\tilde{e}}_1)\,dt - \int_0^T a(e_1 + \varepsilon e_2;\tilde{\tilde{e}}_2)\,dt \\
= {}& \int_0^T \langle \dot{\tilde{\tilde{e}}} - \dot{e}\,;\tilde{\tilde{e}}\rangle_{\mathcal{H}}\,dt + \int_0^T \langle \dot{e}\,;\tilde{\tilde{e}}-V\rangle_{\mathcal{H}}\,dt - \int_0^T a(e_2;\tilde{\tilde{e}}_1 - V_1)\,dt + \int_0^T a(e_1 + \varepsilon e_2;\tilde{\tilde{e}}_2 - V_2)\,dt \\
& + \int_0^T a(e_2 - \tilde{\tilde{e}}_2;\tilde{\tilde{e}}_1)\,dt + \int_0^T a(\tilde{\tilde{e}}_1 - e_1;\tilde{\tilde{e}}_2)\,dt + \varepsilon\int_0^T a(\tilde{\tilde{e}}_2 - e_2;\tilde{\tilde{e}}_2)\,dt - \varepsilon\int_0^T a(\tilde{\tilde{e}}_2;\tilde{\tilde{e}}_2)\,dt.
\end{aligned}
$$

With the identity (3.55), we complete the proof. $\qquad\square$

**Theorem 3.2.2.1 (A priori energy error estimate, $cG(1)$ time approximation).** There exists a constant $C$, which is independent of $u$ and its $cG(1) \otimes \mathcal{C}^1$ discrete counterpart, such that for $\varepsilon = 0$, $p = 3$ and the $h-$quasi uniform spatial mesh, there holds

$$
\begin{aligned}
\|\tilde{\tilde{e}}\|_{L^\infty(\mathcal{H})} \leq C\Big\{ & \|h^p D^{p+1} y_0\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)} + \|h^p D^{p+1}\dot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1}\dot{u}_2\|_{L^1(L^2)} \\
& + \|k^2\ddot{u}_1\|_{L^1(H^1)} + \|h^p D^{p+1} u_2\|_{L^2(L^2)} + \|k^2\Delta\ddot{u}_1\|_{L^1(H^1)} + \|h^p D^{p+1} u_1\|_{L^2(L^2)} \\
& + \|h^p D^{p+1}\tilde{u}_2\|_{L^2(L^2)} + \|k^2\ddot{u}_2\|_{L^1(H^1)} + \|k^2\Delta\ddot{u}_1\|_{L^1(L^2)} \Big\},
\end{aligned}
\tag{3.56}
$$

provided that $u \in H^{p+1}(\Omega)^2$.

**Remark 3.2.2.1.** The estimate of Theorem 3.2.2.1 is of order $\mathcal{O}(h^p+k^2)$ where $p=3$. ☐

**Proof.** The globally continuous and piecewise affine function $\tilde{\tilde{e}}$ reaches its maximum in some of the nodes $t_j$ from $\mathscr{T}$, i.e.

$$\max_{t\in[0,T]}\|\tilde{\tilde{e}}(t)\|_{\mathcal{H}}=\|\tilde{\tilde{e}}\|_{L^\infty(\mathcal{H})}=\|\tilde{\tilde{e}}(t_j)\|_{\mathcal{H}}. \tag{3.57}$$

Owing to the fact that the representation (3.54) is also valid for this particular choice of $T=t_j$, (3.54) can be rewritten such that

$$\|\tilde{\tilde{e}}\|^2_{L^\infty(\mathcal{H})}+\varepsilon\int_0^{t_j}\|\tilde{\tilde{e}}_2(t)\|^2_{H^1(\Omega)}dt=:\sum_{\ell=1}^9 E_\ell. \tag{3.58}$$

The main idea in the following is to estimate $E_1,\dots,E_9$ such that the final estimate consists of the exact solution and its derivatives or the error contributions $\|\tilde{\tilde{e}}_1\|_{L^\infty(H^1)}, \|\tilde{\tilde{e}}_2\|_{L^\infty(L^2)}$ and $\varepsilon\|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))}$. The error terms can be then absorbed through the LHS of (3.58).
We assume in the following that $p=1$ for $\mathcal{P}_1$ and $p=3$ for $\mathcal{C}^1$ elements in space.

First we choose a test function $V$ such that

$$V:=\mathcal{J}\Pi\tilde{\tilde{e}}\in\mathcal{W}_c, \quad\text{where}\quad \Pi=(\mathcal{G},\mathcal{L}). \tag{3.59}$$

Here $\mathcal{J}$ denote the temporal projection, $\mathcal{J}\tilde{\tilde{e}}|_{I_j}:=\fint_{I_j}\tilde{\tilde{e}}(t)dt$ for each $I_j\in\mathscr{T}$. The idea is to use the time projection $\mathcal{J}$ in order to apply the orthogonal properties with respect to $L^2$ norm in time.

The term $E_1$ is estimated with aid of Lemma 3.2.1.2.
We continue by estimating each of $E_2+E_3$, $E_4+E_5$, $E_6$, $E_7$, $E_8$, $E_9$.

**Lemma 3.2.2.2.** For $E_2+E_3$ there exists some constant $C$ such that

$$E_2+E_3\leq C\big\{\|h^pD^{p+1}\dot{u}_1\|_{L^1(L^2)}\|\tilde{\tilde{e}}_1\|_{L^\infty(H^1)}+\|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)}\|\tilde{\tilde{e}}_2\|_{L^\infty(L^2)}\big\}. \tag{3.60}$$

**Proof.** We start from

$$E_3=\int_0^{t_j}\langle\dot{e};\tilde{\tilde{e}}-\mathcal{J}\Pi\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt=\int_0^{t_j}\langle\dot{e};\tilde{\tilde{e}}-\mathcal{J}\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt+\int_0^{t_j}\langle\dot{e};\mathcal{J}\tilde{\tilde{e}}-\mathcal{J}\Pi\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt. \tag{3.61}$$

Moreover, with the definition of $E_2$ and the fact that $\int_{I_j}\dot{\tilde{\tilde{e}}}(t)dt=\int_{I_j}\dot{e}(t)dt$ for each $I_j\in\mathscr{T}$, there holds

$$\int_0^{t_j}\langle\dot{e};\tilde{\tilde{e}}-\mathcal{J}\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt=\int_0^{t_j}\langle\dot{e};\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt-\int_0^{t_j}\langle\dot{e};\mathcal{J}\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt$$
$$=\int_0^{t_j}\langle\dot{e};\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt-\int_0^{t_j}\langle\dot{\tilde{\tilde{e}}};\mathcal{J}\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt$$
$$=\int_0^{t_j}\langle\dot{e};\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt-\int_0^{t_j}\langle\dot{\tilde{\tilde{e}}};\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt=\int_0^{t_j}\langle\dot{e}-\dot{\tilde{\tilde{e}}};\tilde{\tilde{e}}\rangle_{\mathcal{H}}dt=-E_2.$$

Note that $\mathcal{J}$ is here a $L^2$ projection on piecewise constant functions in time.

For the second term on the RHS of (3.61), we have according to the properties of the projections $\mathcal{G}$ and $\mathcal{L}$, see Lemma 3.1.0.6 and Lemma 3.1.0.7, respectively,

$$
\begin{aligned}
\int_0^{t_j} \langle \dot{e}\,; \mathcal{J}\tilde{e} - \mathcal{J}\Pi\tilde{e} \rangle_{\mathcal{H}} dt &= \int_0^{t_j} a(\dot{e}_1; \mathcal{J}\tilde{e}_1 - \mathcal{G}\mathcal{J}e_1) dt + \int_0^{t_j} (\dot{e}_2; \mathcal{J}\tilde{e}_2 - \mathcal{L}\mathcal{J}e_2) dt \\
&= \int_0^{t_j} a(\dot{e}_1 - \mathcal{G}\dot{e}_1; \mathcal{J}\tilde{e}_1) dt + \int_0^{t_j} (\dot{e}_2 - \mathcal{L}\dot{e}_2; \mathcal{J}\tilde{e}_2) dt \\
&\leq C\big( \|h^p D^{p+1}\dot{u}_1\|_{L^1(L^2)} \|\tilde{e}_1\|_{L^\infty(H^1)} + \|h^{p+1} D^{p+1}\dot{u}_2\|_{L^1(L^2)} \|\tilde{e}_2\|_{L^\infty(L^2)} \big).
\end{aligned}
$$

This completes the proof of Lemma.                                                     $\square$

In the following lemma why we estimate $E_4 + E_5$ only under certain requirements concerning the parameter $\varepsilon$, the choice of the boundary conditions and of the spatial mesh where the spatial discretisation is carried out with aid of $\mathcal{C}^1$ elements in space, only.

**Lemma 3.2.2.3.** For $E_4 + E_5$, provided that either $\varepsilon = 0$ or the initial problem satisfies the (DD) boundary conditions, for $p=3$ and $h-$quasi uniform spatial mesh, there holds

$$
\begin{aligned}
E_4 + E_5 \leq C\Big\{ &\|k^2\ddot{u}_1\|_{L^1(H^1)} \|\tilde{e}_1\|_{L^\infty(H^1)} + \|h^p D^{p+1}u_2\|_{L^1(L^2)} \|\tilde{e}_1\|_{L^\infty(H^1)} \\
&+ \|k^2\Delta\ddot{u}_1\|_{L^1(H^1)} \|\tilde{u}_2\|_{L^\infty(L^2)} + \|h^p D^{p+1}u_1\|_{L^2(L^2)} \|h^p D^{p+1}\tilde{u}_2\|_{L^2(L^2)} \Big\},
\end{aligned}
\tag{3.62}
$$

where $C$ is some numerical constant.

**Proof.** We start from $E_4$

$$
E_4 = -\int_0^{t_j} a(e_2; \tilde{e}_1 - \mathcal{J}\mathcal{G}\tilde{e}_1) dt = -\int_0^{t_j} a(e_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt - \int_0^{t_j} a(e_2; \mathcal{J}\tilde{e}_1 - \mathcal{J}\mathcal{G}\tilde{e}_1) dt.
\tag{3.63}
$$

With aid of the Lemma 3.1.0.9, the second term on the RHS of (3.63) can be estimated such that

$$
\int_0^{t_j} a(e_2; \mathcal{J}\tilde{e}_1 - \mathcal{J}\mathcal{G}\tilde{e}_1) dt = -\int_0^{t_j} a(e_2 - \mathcal{G}e_2; \mathcal{J}\tilde{e}_1) dt \leq C\|h^p D^{p+1}u_2\|_{L^1(L^2)} \|\tilde{e}_1\|_{L^\infty(H^1)}.
$$

Furthermore, the first term on the RHS of (3.63) can be rewritten as

$$
\begin{aligned}
-\int_0^{t_j} a(e_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt &= -\int_0^{t_j} a(e_2 - \tilde{e}_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt - \int_0^{t_j} a(\tilde{e}_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt \\
&= \int_0^{t_j} a(\tilde{u}_2 - u_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt - \int_0^{t_j} a(\tilde{e}_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt
\end{aligned}
\tag{3.64}
$$

Moreover, the fact that $u(t_j) = \tilde{u}(t_j)$ for each $t_j$ from $\mathscr{T}$ and the fundamental theorem of calculus prove

$$
u(t) - \tilde{u}(t) \leq Ck_j^2 \|\ddot{u}\|_{L^1(I_j)} \quad \text{for all} \quad t \in I_j, \ I_j \in \mathscr{T},
\tag{3.65}
$$

and some numerical constant $C > 0$. Then, there holds

$$\int_0^{t_j} a(\tilde{u}_2 - u_2; \tilde{e}_1 - \mathcal{J}\tilde{e}_1) dt \leq C \|k^2 \ddot{u}_1\|_{L^1(H^1)} \|\tilde{e}_1\|_{L^\infty(H^1)}. \tag{3.66}$$

Similarly, if $E_5$ is decomposed such that

$$E_5 = \int_0^{t_j} a(e_1; \tilde{e}_2 - \mathcal{J}\mathcal{L}\tilde{e}_2) dt = \int_0^{t_j} a(e_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt + \int_0^{t_j} a(e_1; \mathcal{J}\tilde{e}_2 - \mathcal{J}\mathcal{L}\tilde{e}_2) dt. \tag{3.67}$$

then the first term can be rewritten as

$$\int_0^{t_j} a(e_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt = \int_0^{t_j} a(e_1 - \tilde{e}_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt + \int_0^{t_j} a(\tilde{e}_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt. \tag{3.68}$$

An integration by parts in space shows

$$\int_0^{t_j} a(e_1 - \tilde{e}_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt = \int_0^{t_j} a(u_1 - \tilde{u}_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt$$

$$= \int_0^{t_j} (\Delta(\tilde{u}_1 - u_1); \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt + \int_0^{t_j} D(u_1 - \tilde{u}_1)(t, 1)(\tilde{e}_2 - \mathcal{J}\tilde{e}_2)(t, 1) dt$$

$$= \int_0^{t_j} (\Delta(\tilde{u}_1 - u_1); \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt - \varepsilon \int_0^{t_j} D(u_2 - \tilde{u}_2)(t, 1)(\tilde{e}_2 - \mathcal{J}\tilde{e}_2)(t, 1) dt \tag{3.69}$$

where the second term in the equality above equals zero when the initial problem satisfies the (DD) boundary conditions.
From (3.65)

$$\int_0^{t_j} (\Delta(\tilde{u}_1 - u_1); \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt \leq C \|k^2 \Delta \ddot{u}_1\|_{L^1(H^1)} \|\tilde{e}_2\|_{L^\infty(L^2)}. \tag{3.70}$$

The trace and Friedrichs inequality provide the following estimate for the second term

$$\varepsilon \int_0^{t_j} D(u_2 - \tilde{u}_2)(t, 1)(\tilde{e}_2 - \mathcal{J}\tilde{e}_2)(t, 1) dt \leq \left( \|k^2 \ddot{u}_2\|_{L^2(H^1)} + \|k^2 \Delta \ddot{u}_2\|_{L^2(L^2)} \right) \varepsilon \|\tilde{e}_2\|_{L^2([0,t_j]; H^1(\Omega))}.$$

In order to neglect the last term on the RHS of (3.68), we recall the second term from (3.64). Because of the symmetry properties of $\mathcal{J}$, there holds

$$\int_0^{t_j} a(\tilde{e}_1; \tilde{e}_2 - \mathcal{J}\tilde{e}_2) dt - \int_0^{t_j} a(\tilde{e}_1 - \mathcal{J}\tilde{e}_1; \tilde{e}_2) dt = 0. \tag{3.71}$$

It remains to estimate the last term on the RHS of (3.67). An integration by parts in space leads to

$$\int_0^{t_j} a(e_1; \mathcal{J}\tilde{e}_2 - \mathcal{J}\mathcal{L}\tilde{e}_2) dt = - \int_0^{t_j} (\Delta e_1; \mathcal{J}\tilde{e}_2 - \mathcal{J}\mathcal{L}\tilde{e}_2) dt$$

$$- \int_0^{t_j} \sum_{k=1}^m [De_1(t)]_k (\mathcal{J}\tilde{e}_2 - \mathcal{J}\mathcal{L}\tilde{e}_2)(t, x_k) dt, \tag{3.72}$$

where $m=n$ in case of Problem (DN) and $m=n-1$ for Problem (DD).
The first term can be estimated by using the properties of the projection $\mathcal{L}$ such that

$$-\int_0^{t_j}(\Delta e_1; \mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt=-\int_0^{t_j}(\Delta e_1-\mathcal{L}\Delta e_1; \mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt$$

$$\leq C\int_0^{t_j}\|h^{p-1}D^{p+1}u_1(t)\|_{L^2(\Omega)}\|h^{p+1}D^{p+1}\tilde{u}_2(t)\|_{L^2(\Omega)}dt$$

$$\leq C\|h^pD^{p+1}u_1\|_{L^2(L^2)}\|h^pD^{p+1}\tilde{u}_2\|_{L^2(L^2)}. \tag{3.73}$$

In the last inequality we used the fact that the spatial mesh is $h-$quasi uniform.
As far as the jump term on the RHS of (3.72) and (DN) boundary conditions are concerned, we can not derive an estimate which depends either on a priori terms or on the error terms $\varepsilon\|\tilde{\tilde{e}}_2\|_{L^2(0,t_j;H^1)}$ and $\|\tilde{\tilde{e}}\|_{L^\infty(\mathcal{H})}$. In case of $cG(1)\otimes\mathcal{P}_1$ discretisation, the problem is to estimate the jump terms $[De]_k$ which do not vanish due to the fact that $De$ is piecewise constant function in space. Therefore we assume in the following that the discretisation in space is done with $\mathcal{C}^1$ functions. Then, the jump term simplifies to

$$-\int_0^{t_j}\sum_{k=1}^m[De_1(t)]_k(\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)(t,x_k)dt=\int_0^{t_j}De_1(t,1)(\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)(t,1)dt, \tag{3.74}$$

since $(\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)(t,0)=0$. Obviously if $\varepsilon=0$, then $De_1(t,1)=0$ and the last term vanishes. Also, if the the initial problem satisfies the (DD) boundary conditions, then $(\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)(t,1)=0$.
We did not succeed to estimate the latter terms in other cases then this two cases mentioned. Therefore, we assume that either one of these two conditions is satisfied, namely that either $\varepsilon=0$ or (DD) boundary conditions. Then for $p=3$ and $h-$quasi uniform mesh, there holds

$$\int_0^{t_j}a(e_1,\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt\leq C\|h^pD^{p+1}u_1\|_{L^2(L^2)}\|h^pD^{p+1}\tilde{u}_2\|_{L^2(L^2)}. \tag{3.75}$$

This concludes the proof. $\qquad\square$

**Lemma 3.2.2.4.** For $E_6$ there exists a constant $C$ such that if the initial problem satisfies the (DD) boundary conditions for $p=3$ and $h-$quasi uniform mesh, there holds

$$E_6\leq C\left\{\varepsilon\|h^pD^{p+1}u_2\|_{L^2(L^2)}\|h^pD^{p+1}\tilde{u}_2\|_{L^2(L^2)}+\|k^2\ddot{u}_2\|_{L^2(H^1)}\varepsilon\|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))}\right\}+\varepsilon\|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))}^2.$$

**Proof.** We start from the following decomposition of $E_6$

$$E_6=\varepsilon\int_0^{t_j}a(e_2;\tilde{\tilde{e}}_2-\mathcal{J}\tilde{\tilde{e}}_2)dt+\varepsilon\int_0^{t_j}a(e_2;\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt. \tag{3.76}$$

An integration by parts in the second integral on the RHS of (3.76) yields

$$\varepsilon\int_0^{t_j}a(e_2;\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt=-\varepsilon\int_0^{t_j}(\Delta e_2;\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt-\varepsilon\int_0^{t_j}\sum_{k=1}^m[De_2(t)]_k(\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)(t,x_k).$$

For the first term we have similar as in (3.73) for the $h-$quasi uniform spatial mesh

$$-\varepsilon\int_0^{t_j}(\Delta e_2;\mathcal{J}\tilde{\tilde{e}}_2-\mathcal{J}\mathcal{L}\tilde{\tilde{e}}_2)dt\leq C\varepsilon\|h^pD^{p+1}u_2\|_{L^2(L^2)}\|h^pD^{p+1}\tilde{u}_2\|_{L^2(L^2)}. \tag{3.77}$$

In order to estimate the jump term, we recall the estimation of (3.74). Obviously, we need to assume that the initial problem satisfies the (DD) boundary conditions and $p = 3$. Then the jump term vanishes.

Moreover, for the first term on the RHS of (3.76) we have

$$\varepsilon \int_0^{t_j} a(e_2; \tilde{\tilde{e}}_2 - \mathcal{J}\tilde{\tilde{e}}_2) dt = \varepsilon \int_0^{t_j} a(e_2 - \tilde{\tilde{e}}_2; \tilde{\tilde{e}}_2 - \mathcal{J}\tilde{\tilde{e}}_2) dt + \varepsilon \int_0^{t_j} a(\tilde{\tilde{e}}_2; \tilde{\tilde{e}}_2 - \mathcal{J}\tilde{\tilde{e}}_2) dt$$

$$\leq C \|k^2 \ddot{u}_2\|_{L^2(H^1)} \varepsilon \|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))} + \varepsilon \|\tilde{\tilde{e}}_2\|^2_{L^2([0,t_j];H^1(\Omega))}.$$

This completes the proof of lemma. □

**Remark 3.2.2.2.** Apperantly, the estimate from Lemma 3.2.2.4 is not optimal due to the presence of $\varepsilon \|e_2\|^2_{L^2(I_j;H^1(\Omega))}$, which is not the only $\varepsilon \|\tilde{\tilde{e}}_2\|_{L^2(I_j;H^1(\Omega))}$ contribution. However, we do not abandon this result because it shows why the estimation for $\varepsilon > 0$ is difficult and does not yield a result, which can be further used. □

**Lemma 3.2.2.5.** For $E_7$ there holds

$$E_7 \leq C \|k^2 \ddot{u}_2\|_{L^1(H^1)} \|\tilde{\tilde{e}}_1\|_{L^\infty(H^1)}, \tag{3.78}$$

with $C > 0$ some numerical constant.

**Proof.** From the definition of $\tilde{\tilde{e}}$ and the identity (3.65) we have that

$$E_7 = \int_0^{t_j} a(e_2 - \tilde{\tilde{e}}_2; \tilde{\tilde{e}}_1) dt = \int_0^{t_j} a(u_2 - \tilde{u}_2; \tilde{\tilde{e}}_1) dt \leq C \|k^2 \ddot{u}_2\|_{L^1(H^1)} \|\tilde{\tilde{e}}_1\|_{L^\infty(H^1)}.$$

This completes the proof. □

**Lemma 3.2.2.6.** For $E_8$ there is a constant $C$ such that

$$E_8 \leq C \left\{ \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} \|\tilde{\tilde{e}}_2\|_{L^\infty(L^2)} + \left( \|k^2 \ddot{u}_2\|_{L^2(H^1)} + \|k^2 \Delta \ddot{u}_2\|_{L^2(L^2)} \right) \varepsilon \|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))} \right\}, \tag{3.79}$$

where the second summand does not appear in case of the (DD) boundary conditions.

**Proof.** An integration by parts in space and the fact that $u_1 - \tilde{u}_1$ is the continuous functions in space yield

$$E_8 = \int_0^{t_j} a(\tilde{\tilde{e}}_1 - e_1; \tilde{\tilde{e}}_2) dt = \int_0^{t_j} (\Delta(u_1 - \tilde{u}_1); \tilde{\tilde{e}}_2) dt + \int_0^{t_j} D(\tilde{u}_1 - u_1)(t, 1) \tilde{\tilde{e}}_2(t, 1) dt$$

$$= \int_0^{t_j} (\Delta(u_1 - \tilde{u}_1); \tilde{\tilde{e}}_2) dt - \varepsilon \int_0^{t_j} D(\tilde{u}_2 - u_2)(t, 1) \tilde{\tilde{e}}_2(t, 1) dt.$$

We assumed here that the initial problem $u = (y, \dot{y})$ satisfies the (DN) boundary condition (1.13c). On the other hand for (DD) boundary conditions (1.13d), the second term on the RHS of either equalities above vanishes.

Finally, the identity (3.65) and the trace and Friedrichs inequality lead to

$$E_8 \leq C \left\{ \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} \|\tilde{\tilde{e}}_2\|_{L^\infty(L^2)} + \left( \|k^2 \ddot{u}_2\|_{L^2(H^1)} + \|k^2 \Delta \ddot{u}_2\|_{L^2(L^2)} \right) \varepsilon \|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))} \right\},$$

which proves the Lemma. □

**Lemma 3.2.2.7.** For $E_9$ the following is valid

$$E_9 \leq C \|k^2 \ddot{u}_2\|_{L^2(H^1)} \varepsilon \|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))}, \tag{3.80}$$

with $C$ some numerical constant.

**Proof.** The Hölder inequality in time and space and estimate (3.65) yield the proof. Namely,

$$E_9 = \varepsilon \int_0^{t_j} a(\tilde{\tilde{e}}_2 - e_2; \tilde{\tilde{e}}_2) dt \leq \|\tilde{u}_2 - u_2\|_{L^2(H^1)} \varepsilon \|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))}$$

$$\leq C \|k^2 \ddot{u}_2\|_{L^2(H^1)} \varepsilon \|\tilde{\tilde{e}}_2\|_{L^2([0,t_j];H^1(\Omega))}. \qquad \square$$

We now substitute the results of Lemma 3.2.1.2 and Lemma 3.2.2.2 -Lemma 3.2.2.7 into the error representation (3.58) and assume $\varepsilon = 0$, $p = 3$ and $h-$quasi uniformity of the spatial mesh. Then, the application of the Young inequality allows us to move the terms $\|\tilde{\tilde{e}}_1\|_{L^\infty(H^1)}$, $\|\tilde{\tilde{e}}_2\|_{L^\infty(L^2)}$ onto the LHS of (3.58) so that we finally obtain

$$\|\tilde{\tilde{e}}\|_{L^\infty(\mathcal{H})} \leq C \Big\{ \|h^p D^{p+1} y_0\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)} + \|h^p D^{p+1} \dot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}$$

$$+ \|k^2 \ddot{u}_1\|_{L^1(H^1)} + \|h^p D^{p+1} u_2\|_{L^2(L^2)} + \|k^2 \Delta \ddot{u}_1\|_{L^1(H^1)} + \|h^p D^{p+1} u_1\|_{L^2(L^2)}$$

$$+ \|h^p D^{p+1} \tilde{u}_2\|_{L^2(L^2)} + \|k^2 \ddot{u}_2\|_{L^1(H^1)} + \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} \Big\}.$$

This completes the proof of theorem.                                                                                $\square$

### 3.2.3   Method of lines

Within this subsection we prove an a priori error estimate for the error of the semi-discrete approximation where the discretisation in space follows by use of $\mathcal{C}^1$ function and the (DD) boundary condition and $h-$quasi uniformity of the spatial mesh are also assumed. The error function for which the analysis applies,

$$e := u - U, \tag{3.81}$$

is a globally continuous function in time.

For notation and definitions used in the following, we refer to Chapter 2, where the particular space discretisation methods, cf. Section 2.1 and the method of lines, cf. Subsection 2.3.5 are introduced.

**Lemma 3.2.3.1 (Error representation, *MoL*).** In case of the semi-discretisation with respect to time, there holds for each $t \in [0, T]$ and all $V \in \mathcal{W}_s$

$$\frac{1}{2} \|e(t)\|_{\mathcal{H}}^2 + \varepsilon \int_0^t \|e_2(\tau)\|_{H^1(\Omega)}^2 d\tau = \frac{1}{2} \|e(0)\|_{\mathcal{H}}^2 + \int_0^t \langle \dot{e} \, ; e - V \rangle_{\mathcal{H}} d\tau - \int_0^t a(e_2; e_1 - V_1) d\tau$$

$$+ \int_0^t a(e_1; e_2 - V_2) d\tau + \varepsilon \int_0^t a(e_2; e_2 - V_2) d\tau. \tag{3.82}$$

**Proof.** From the fundamental theorem of calculus in time and the fact that the error $e$ is a continuous function in time, we have

$$\frac{1}{2} \frac{\partial}{\partial t} \|e(t)\|_{\mathcal{H}}^2 = \langle \dot{e}(t) \, ; e(t) \rangle_{\mathcal{H}}, \quad t \in [0, T]. \tag{3.83}$$

With the definition of the bilinear form $\mathcal{B}$ (2.49) and the Galerkin orthogonality (2.9), we have for all $V \in \mathcal{W}_s$ and continously in time

$$
\begin{aligned}
\langle \dot{e} ; e \rangle_{\mathcal{H}} &= \langle \dot{e} ; e - V \rangle_{\mathcal{H}} + \langle \dot{e} ; V \rangle_{\mathcal{H}} \\
&= \langle \dot{e} ; e - V \rangle_{\mathcal{H}} + a(e_2; V_1) - a(e_1 + \varepsilon e_2; V_2) \\
&= \langle \dot{e} ; e - V \rangle_{\mathcal{H}} - a(e_2; e_1 - V_1) + a(e_1 + \varepsilon e_2; e_2 - V_2) - \varepsilon a(e_2, e_2),
\end{aligned}
$$

where $e = e(t), \dot{e} = \dot{e}(t)$ and also $V = V(t)$.

If we integrate the last equation with respect to time interval $[0, t]$ and apply the identity (3.83), we may conclude the proof of the Lemma.      $\square$

**Theorem 3.2.3.1 (A priori energy error estimate, *MoL*).** There is a constant $C$ such that for $u \in H^{p+1}(\Omega)^2$, $p = 3$ and $h-$quasi uniform spatial mesh, there holds

(a) in case of the (DD) boundary conditions

$$
\begin{aligned}
\|e\|_{L^\infty(\mathcal{H})} \leq C \Big\{ &\|h^p D^{p+1} y_0\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)} + \|h^p D^{p+1} \dot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} \\
&+ \|h^p D^{p+1} u_2\|_{L^1(L^2)} + \|h^p D^{p+1} u_1\|_{L^2(L^2)} + \|h^p D^{p+1} u_2\|_{L^2(L^2)} \\
&+ \sqrt{\varepsilon} \|h^p D^{p+1} u_2\|_{L^2(L^2)} \Big\},
\end{aligned}
\tag{3.84}
$$

(b) if $\varepsilon = 0$

$$
\begin{aligned}
\|e\|_{L^\infty(\mathcal{H})} \leq C \Big\{ &\|h^p D^{p+1} y_0\|_{L^2(\Omega)}^2 + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)}^2 + \|h^p D^{p+1} \dot{u}_1\|_{L^1(L^2)}^2 + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}^2 \\
&+ \|h^p D^{p+1} u_2\|_{L^1(L^2)}^2 + \|h^p D^{p+1} u_1\|_{L^2(L^2)}^2 + \|h^p D^{p+1} u_2\|_{L^2(L^2)}^2 \Big\}.
\end{aligned}
\tag{3.85}
$$

**Remark 3.2.3.1.** Note that both estimates of Theorem 3.2.3.1 are of order $\mathcal{O}(h^p)$ i.e. $\mathcal{O}(h^3)$. $\square$

**Proof.** Fix some $t \in [0, T]$ such that

$$
\|e(t)\|_{\mathcal{H}} = \|e\|_{L^\infty(\mathcal{H})}.
\tag{3.86}
$$

Then, from the representation (3.82), we have

$$
\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon \int_0^t \|e_2(\tau)\|_{H^1(\Omega)}^2 d\tau =: \sum_{\ell=1}^5 E_\ell.
\tag{3.87}
$$

The main idea in the following is to estimate the terms $E_1, \ldots, E_5$ from (3.87) such that the final estimate contains the contribution of the exact solution only or the error terms $\|e\|_{L^\infty(\mathcal{H})}$ or $\varepsilon \|e_2\|_{L^2(0,t;H^1(\Omega))}^2$ which can then be absorbed through the LHS of (3.87), by an application of the Young inequality.

As before, we write $p = 1$ for $\mathcal{P}_1$ elements and $p = 3$ for $\mathcal{C}^1$ elements in space. The proof will be conducted in a similar way as in case of the $cG(1)$ and $dG(0)$ time discretisation, i.e. the estimation of terms $E_1, \ldots, E_5$ will be given through the several lemmas.

Let us first assume that the test function $V$ is defined by

$$V := \Pi e \in \mathcal{W}_s, \ \Pi = (\mathcal{G}, \mathcal{L}). \tag{3.88}$$

For fix $t \in [0, T]$, we have $V(t) \in \mathcal{S} \times \mathcal{S}$. We also recall the approximation properties of the projections $\mathcal{G}, \mathcal{L}$ given in Lemma 3.1.0.6 and Lemma 3.1.0.7, respectively.

$E_1$ is estimated by means of Lemma 3.2.1.2. We continue with the estimation of $E_2, E_3, E_4, E_5$.

**Lemma 3.2.3.2.** There is a constant $C$ such that

$$E_2 \leq C\big(\|h^p D^{p+1}\dot{u}_1\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)} + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)}\|e_2\|_{L^\infty(L^2)}\big).$$

**Proof.** The approximation properties of the projections $\mathcal{G}, \mathcal{L}$ imply

$$\begin{aligned}
E_2 &= \int_0^t a(\dot{e}_1; e_1 - \mathcal{G}e_1)d\tau + \int_0^t (\dot{e}_2; e_2 - \mathcal{L}e_2)d\tau \\
&= \int_0^t a(\dot{e}_1 - \mathcal{G}\dot{e}_1; e_1)d\tau + \int_0^t (\dot{e}_2 - \mathcal{L}\dot{e}_2; e_2)d\tau \\
&\leq C\big(\|h^p D^{p+1}\dot{u}_1\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)} + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)}\|e_2\|_{L^\infty(L^2)}\big). \qquad \square
\end{aligned}$$

**Lemma 3.2.3.3.** For $E_3$ there exists a constant $C$ such that

$$E_3 \leq C\|h^p D^{p+1}u_2\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)}.$$

**Proof.** By using the symmetry of Galerkin projection $\mathcal{G}$ with respect to $H^1$ scalar product, we obtain

$$E_3 = -\int_0^t a(e_2; e_1 - \mathcal{G}e_1)d\tau = -\int_0^t a(e_2 - \mathcal{G}e_2; e_1)d\tau \leq C\|h^p D^{p+1}u_2\|_{L^1(L^2)}\|e_1\|_{L^\infty(H^1)}. \qquad \square$$

In the following lemma we show why the estimation of $E_4$ can be accomplished only for $p = 3$ under certain restrictions concerning the parameter $\varepsilon$ and the choice of the boundary conditions and the spatial mesh.

**Lemma 3.2.3.4.** For $E_4$ and $p = 3$ there exists a constant $C$ such that if $\varepsilon = 0$ or the initial problem satisfies the (DD) boundary conditions, there holds

$$E_4 \leq C\|h^p D^{p+1}u_1\|_{L^2(L^2)}\|h^p D^{p+1}u_2\|_{L^2(L^2)}.$$

**Proof.** A partial integration in space yields

$$E_4 = \int_0^t a(e_1; e_2 - \mathcal{L}e_2)d\tau = -\int_0^t (\Delta e_1; e_2 - \mathcal{L}e_2)d\tau - \int_0^t \sum_{k=1}^m [De_1(\tau)]_k(e_2 - \mathcal{L}e_2))(\tau, x_k)d\tau.$$

As in the proof of Lemma 3.2.2.3, for $p = 3$ if $\varepsilon = 0$ or the initial problem satisfies the (DD) boundary conditions the second integral vanishes. If we assume additional that the spatial mesh is $h-$quasi uniform, then

$$E_4 = \int_0^t a(e_1; e_2 - \mathcal{L}e_2)d\tau = -\int_0^t (\Delta e_1; e_2 - \mathcal{L}e_2)d\tau \leq C\|h^p D^{p+1}u_1\|_{L^2(L^2)}\|h^p D^{p+1}u_2\|_{L^2(L^2)}.$$

This concludes the proof. $\qquad \square$

**Lemma 3.2.3.5.** For $E_5$ there exists a constant $C$ such that for $p = 3$, (DD) boundary conditions and $h-$quasi uniform spatial mesh, there holds

$$E_5 \leq C\varepsilon \|h^p D^{p+1} u_2\|_{L^2(L^2)}^2.$$

**Proof.** As in the proof of Lemma 3.2.3.4, we may conclude

$$E_5 = \varepsilon \int_0^t a(e_2; e_2 - \mathcal{L}e_2) d\tau \leq C\varepsilon \|h^p D^{p+1} u_2\|_{L^2(L^2)}^2. \qquad \square$$

From an application of Young inequality, provided (DD) boundary conditions and $h-$quasi uniform spatial mesh, there holds for $p = 3$

$$\|e\|_{L^\infty(\mathcal{H})} \leq C \Big\{ \|h^p D^{p+1} y_0\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)} + \|h^p D^{p+1} \dot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}$$

$$+ \|h^p D^{p+1} u_2\|_{L^1(L^2)} + \|h^p D^{p+1} u_1\|_{L^2(L^2)} + \|h^p D^{p+1} u_2\|_{L^2(L^2)} + \varepsilon \|h^p D^{p+1} u_2\|_{L^2(L^2)} \Big\}.$$

Moreover, if $\varepsilon = 0$ and (DD) or (DN) boundary conditions, there holds additionally

$$\|e\|_{L^\infty(\mathcal{H})} \leq C \Big\{ \|h^p D^{p+1} y_0\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)} + \|h^p D^{p+1} \dot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}$$

$$+ \|h^p D^{p+1} u_2\|_{L^1(L^2)} + \|h^p D^{p+1} u_1\|_{L^2(L^2)} + \|h^p D^{p+1} u_2\|_{L^2(L^2)} \Big\}.$$

This yields the proof of the theorem. $\qquad \square$

## 3.3 A posteriori energy error estimate

Within this section we derive a posteriori error bounds for the fully discrete model with the discontinuous and continuous Galerkin method in time, i.e. $dG(q), q = 0, 1$, $cG(1)$, respectively, and cubic and linear splines in space, abb. $\mathcal{C}^1$ and $\mathcal{P}_1$ finite elements. We also provide a short analysis for the semi-discrete model obtained from the discretisation in space, where time remains continuous.

The analysis follows via residual arguments introduced in Definition 2.3.1.2. The a posteriori error bound $\eta$ for the error of the fully (time-space) discrete problem and error damping term is a quantity depending on the mesh data $h, k$, the discrete solution $U$ and the initial data $u_0, f$. An overview concerning the accuracy order of developed estimates, with respect to time and space discretisation is presented in Table 3.3.

The spatial approximation by $\mathcal{P}_1$ functions was not an adequate choice for the derivation of a posteriori error bounds by use of the energy techniques. This holds independently of the time approximation method due to the lack of continuity in the first derivative of the discrete solution. For a further discussion see Subsection 3.3.1.1.

On the other hand, the $\mathcal{C}^1$ elements in space enable the derivation of an a posteriori error bound and provide the expected optimal accuracy of order $\mathcal{O}(h^3)$ in space.

Concerning the convergence order in time, only the $dG(0)$ method provides the optimal result, i.e. $\mathcal{O}(k)$, see Figure 3.4. The $dG(1)$ and $cG(1)$ method in time yield estimates with a suboptimal order of convergence $\mathcal{O}(k)$, where $\mathcal{O}(k^3)$ and $\mathcal{O}(k^2)$, respectively, are expected. A detailed discussion follows in Subsections 3.3.2.2 and 3.3.3.2, respectively. For an example see Figure 3.6.

| $\otimes$ | $\mathcal{P}_1$ | $\mathcal{C}^1$ |
|:---:|:---:|:---:|
| **$dG(0)$** <br><br> Subsection 3.3.1 | – | $\mathcal{O}(h^3+k)$ $_{\mathcal{S}^{j-1}\subseteq\mathcal{S}^j}$ <br> $\mathcal{O}(h^3+k^{-1}h^3+k)$ $_{\text{otherwise}}$ |
| **$dG(1)$** <br><br> Subsection 3.3.2 | – | $\mathcal{O}(h^3+k)$ $_{\mathcal{S}^{j-1}\subseteq\mathcal{S}^j}$ <br> $\mathcal{O}(h^3+k^{-1}h^3+k)$ $_{\text{otherwise}}$ |
| **$cG(1)$** <br><br> Subsection 3.3.3 | – | $\mathcal{O}(h^3+k)$ |
| **$MoL$** <br><br> Subsection 3.3.4 | – | $\mathcal{O}(h^3)$ |

Table 3.3: Proven a posteriori error estimates for $\|\tilde{e}\|_{L^\infty(\mathcal{H})}+\sqrt{\varepsilon}\|\tilde{e}_2\|_{L^2(H^1)}$ in case of $dG(0)$ and $\|e\|_{L^\infty(\mathcal{H})}+\sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}$ in case of $dG(1), cG(1)$ and $MoL$; energy method.

However, the estimates proved within this section are computable, i.e. the constants which occur in the error bound are exact. Improved results in case of the $dG(1)$ and $cG(1)$ method in time with an optimal convergence order in time are obtained by use of the duality technique, see Section 4.3.



Figure 3.4: Convergence of $\|\tilde{e}\|_{L^\infty(\mathcal{H})}+\sqrt{\varepsilon}\|\tilde{e}_2\|_{L^2(H^1)}$ and $\eta$ for $dG(0)\otimes\mathcal{C}^1$ with respect to the number of elements in space (time) for $h=k$; Proven a posteriori error bound $\eta=\mathcal{O}(h^3+k)=$ $\mathcal{O}(k)$ and $\|\tilde{e}\|_{L^\infty(\mathcal{H})}+\sqrt{\varepsilon}\|\tilde{e}_2\|_{L^2(H^1)}=\mathcal{O}(h^3+k)=\mathcal{O}(k)$; Example 6.2.4, $\varepsilon=0.1$, (DN), $T=1$; energy method.

Figure 3.5: Convergence of $\|e\|_{L^\infty(\mathcal{H})} + \sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}$ and $\eta$ for $dG(1)\otimes\mathcal{C}^1$ with respect to the number of elements in space (time) for $h=k$; Proven a posteriori error bound $\eta=\mathcal{O}(h^3+k)=\mathcal{O}(h^3)=\mathcal{O}(k)$ is suboptimal compared with $\|e\|_{L^\infty(\mathcal{H})}+\sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}=\mathcal{O}(h^3+k^2)=\mathcal{O}(h^2)=\mathcal{O}(k^2)$ and $\|e\|_{L^\infty(\mathcal{H})}=\mathcal{O}(h^3+k^3)=\mathcal{O}(h^3)=\mathcal{O}(k^3)$; Left figure: Example 6.2.3, $\varepsilon=0$, (DD), $T=1$; Right figure: Example 6.2.4, $\varepsilon=0.1$, (DN), $T=1$; energy method.



Figure 3.6: Convergence of $\|e\|_{L^\infty(\mathcal{H})} + \sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}$ and $\eta$ for $cG(1)\otimes\mathcal{C}^1$ with respect to the number of elements in space (left) and time (right) for $h=k^{2/3}$; Proven a posteriori error bound $\eta = \mathcal{O}(h^3 + k) = \mathcal{O}(h^{3/2}) = \mathcal{O}(k)$ is suboptimal compared with $\|e\|_{L^\infty(\mathcal{H})} + \sqrt{\varepsilon}\|e_2\| = \mathcal{O}(h^3+k^2) = \mathcal{O}(h^3) = O(k^2)$; Left figure: Example 6.2.5, $\varepsilon = 0$, (DN), $T = 1$; Right figure: Example 6.2.4, $\varepsilon=0.1$, (DN), $T=1$; energy method.

### 3.3.1 $dG(0)$ time approximation

The main tool in the following a posteriori error analysis is the "modified" bilinear form $\widetilde{\mathcal{B}}$ introduced in (2.20), the linear functional $\mathscr{L}$ from (2.19) and the residual $\widetilde{Res}$ from (2.21). Through these definitions, the concept of affine approximation concerning the discrete solution $\widetilde{U}$ is introduced, see Definition 2.3.3.3. This allows the temporal jump terms to be included in the definition of the residual as a function of time. Thereafter we may introduce the error

$$\tilde{e}:=u-\widetilde{U}\in H^1(0,T;\mathcal{H}),\tag{3.89}$$

which is a globally continuous function in time. The term of interests in the following a posteriori error analysis will be precisely this error function, i.e.

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 + \varepsilon \|\tilde{e}_2\|_{L^2(H^1)}^2,$$

for which we seek to derive the computable bound $\eta$.

**Lemma 3.3.1.1 (Error representation, $dG(0)$ time approximation).** In case of $dG(0)$ time discretisation, there holds for every $V \in \mathcal{Q}_0$

$$\frac{1}{2}\|e^{N-}\|_{\mathcal{H}}^2 + \varepsilon \sum_{j=1}^{N} \int_{I_j} \|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 dt = \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \widetilde{Res}(\tilde{e}-V) + \sum_{j=1}^{N} \int_{I_j} a(e_2 - \tilde{e}_2; \tilde{e}_1)dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} a(\tilde{e}_1 - e_1; \tilde{e}_2)dt + \varepsilon \sum_{j=1}^{N} \int_{I_j} a(\tilde{e}_2 - e_2; \tilde{e}_2)dt. \quad (3.90)$$

**Proof.** The fundamental theorem of calculus and the fact that $\tilde{e}(t_j) = e(t_j)$, for each $t_j \in \mathscr{T}$, lead to

$$\frac{1}{2}\|e^{N-}\|_{\mathcal{H}}^2 - \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 = \sum_{j=1}^{N} \int_{I_j} \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} dt. \quad (3.91)$$

With the residual representation (2.22) and orthogonality of the residual from (2.23), we show that for arbitrary $V \in \mathcal{Q}_0$

$$\sum_{j=1}^{N} \int_{I_j} \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} dt = \sum_{j=1}^{N} \int_{I_j} \langle \dot{\tilde{e}}; \tilde{e}-V \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \int_{I_j} \langle \dot{\tilde{e}}; V \rangle_{\mathcal{H}} dt$$

$$= \sum_{j=1}^{N} \int_{I_j} \langle \dot{\tilde{e}}; \tilde{e}-V \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \int_{I_j} a(e_2; V_1)dt - \sum_{j=1}^{N} \int_{I_j} (e_1; V_2)dt - \varepsilon \sum_{j=1}^{N} \int_{I_j} (e_2; V_2)dt$$

$$= \widetilde{Res}(\tilde{e}-V) + \sum_{j=1}^{N} \int_{I_j} a(e_2; \tilde{e}_1)dt - \sum_{j=1}^{N} \int_{I_j} a(e_1; \tilde{e}_2)dt - \varepsilon \sum_{j=1}^{N} \int_{I_j} a(e_2; \tilde{e}_2)dt$$

$$= \widetilde{Res}(\tilde{e}-V) + \sum_{j=1}^{N} \int_{I_j} a(e_2 - \tilde{e}_2; \tilde{e}_1)dt + \sum_{j=1}^{N} \int_{I_j} a(\tilde{e}_1 - e_1; \tilde{e}_2)dt$$

$$+ \varepsilon \sum_{j=1}^{N} \int_{I_j} a(\tilde{e}_2 - e_2; \tilde{e}_2)dt - \varepsilon \sum_{j=1}^{N} \int_{I_j} a(\tilde{e}_2; \tilde{e}_2)dt. \quad (3.92)$$

The combination of the last equality and identity (3.91) yields the proof. $\qquad \square$

**Lemma 3.3.1.2.** In case of $dG(0)$ time approximation and $\mathcal{P}_1$, i.e. $\mathcal{C}^1$ approximation in space, the residual (2.21) may be expressed in terms of the local residuals and local jumps such that for all $v \in L^2(H_D^1 \times H_D^1)$,

$$\widetilde{Res}(v) = \sum_{j=1}^{N} \int_{I_j} \langle \widetilde{R}_j; v \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \int_{I_j} \widetilde{J}_j(U, v)dt, \quad (3.93a)$$

where for all $(t,x) \in I_j \times \Omega$,

$$\widetilde{R}_j(t,x) := F(t,x) - \dot{\tilde{U}}(t,x) - \mathcal{A}U(t,x), \tag{3.93b}$$

$$\widetilde{J}_j(U,v)(t) := \begin{cases} \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)]_k v_2(t,x_k) & \text{for } \mathcal{P}_1 \text{ elements,} \\ 0 & \text{for } \mathcal{C}^1 \text{ elements.} \end{cases} \tag{3.93c}$$

Here $m = n$ in case of (DN) and $m = n-1$ in case of (DD) boundary conditions. Furthermore, $[D(U_1 + \varepsilon U_2)]_k$ are defined as jumps of the piecewise constant functions $DU_1$, $DU_2$ with respect to the space coordinate, see Definition 3.1.0.6.

**Proof.** Given the residual definition from (2.21) with the bilinear form $\widetilde{\mathcal{B}}$ and functional $\mathscr{L}$ defined in (2.20) and (2.19), respectively, an integration by parts in space yields the following result for all $v \in L^2(H_D^1 \times H_D^1)$

$$\mathscr{L}(v) - \widetilde{\mathcal{B}}(U,v) = \sum_{j=1}^{N} \int_{I_j} (f; v_2) dt - \sum_{j=1}^{N} \int_{I_j} a(\dot{\tilde{U}}_1; v_1) - \sum_{j=1}^{N} (\dot{\tilde{U}}_2; v_2) + \sum_{j=1}^{N} \int_{I_j} a(U_2; v_1) dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} (\Delta(U_1 + \varepsilon U_2); v_2) dt + \sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)]_k v_2(x_k) dt.$$

In case of the spatial approximation by Hermite cubic splines, the jump terms equal zero due to the fact that discrete functions are globally continuous in the first derivative and they also satisfy the (DN) boundary conditions too, cf. Subsections 6.1.3 and 6.1.4. With the notation from Definition 1.3.0.7 we conclude the proof of the lemma. $\square$

**Remark 3.3.1.1.** If $\mathscr{T}$ is an arbitrary discretisation of the time interval $[0,T]$, then from the definition of the local weak problem (2.6) and residual representation in Lemma 3.3.1.2, there holds for each $t_n \in \mathscr{T}$

$$\frac{1}{2}\|e^{n-}\|_{\mathcal{H}}^2 + \varepsilon \sum_{j=1}^{n} \int_{I_j} \|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 dt = \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \sum_{j=1}^{n} \int_{I_j} \langle \widetilde{R}_j ; \tilde{e} - V \rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \int_{I_j} \widetilde{J}_j(U; \tilde{e} - V) dt$$

$$+ \sum_{j=1}^{n} \int_{I_j} a(\widetilde{U}_2 - U_2; \tilde{e}_1) dt + \sum_{j=1}^{n} \int_{I_j} a(U_1 - \widetilde{U}_1; \tilde{e}_2) dt$$

$$+ \varepsilon \sum_{j=1}^{n} \int_{I_j} a(U_2 - \widetilde{U}_2; \tilde{e}_2) dt. \tag{3.94}$$

Here $V$ is an arbitrary test function, $V|_{I_j} \in \mathcal{Q}_0^j$ for all $j = 1, \ldots, n$ for which the residual orthogonality (2.8) applies. $\square$

The following technical lemma will be used latter on for the different space discretisation separately. Then the abstract contributions $M_1$, $M_2$, $M_{3,n}$, $M_{4,n}$ will be defined accordingly.

**Lemma 3.3.1.3.** Let $\widetilde{R}$ and $\tilde{J}$ be defined by (3.93b) and (3.93c), respectively. Let $\tilde{e}$ be the error of the affine approximation of the $dG(0)$ solution with respect to $\mathcal{P}_1$ or $\mathcal{C}^1$ space discretisation. Let $\Pi := (\mathcal{G}, \mathcal{L})$ be a spatial orthogonal projection with respect to the $\mathcal{H}$-scalar

product. Finally, we assume that there are some positive $M_1, M_2, M_{3,n}, M_{4,n}$ such that there holds for each $n = 1, \ldots, N$ and both $\alpha \in \{2, 4\}$,

$$\sum_{j=1}^{n} \int_{I_j} \langle \widetilde{R}_j ; \tilde{e} - \Pi \bar{\tilde{e}} \rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \int_{I_j} \widetilde{J}_j(U; \tilde{e} - \Pi \bar{\tilde{e}}) dt \leq \alpha M_1 + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}^2, \tag{3.95a}$$

$$\sum_{j=1}^{n} \int_{I_j} a(\widetilde{U}_2 - U_2; \tilde{e}_1) dt + \sum_{j=1}^{n} \int_{I_j} a(U_1 - \widetilde{U}_1 + \varepsilon(U_2 - \widetilde{U}_2); \tilde{e}_2) dt \leq \alpha M_2 + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}^2, \tag{3.95b}$$

$$\int_{I_n} \langle \widetilde{R}_n ; \tilde{e} \rangle_{\mathcal{H}} dt + \int_{I_n} \widetilde{J}_n(U; \tilde{e}) dt \leq \alpha M_{3,n} + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}^2, \tag{3.95c}$$

$$\int_{I_n} a(\widetilde{U}_2 - U_2; \tilde{e}_1) dt + \int_{I_n} a(U_1 - \widetilde{U}_1 + \varepsilon(U_2 - \widetilde{U}_2); \tilde{e}_2) dt \leq \alpha M_{4,n} + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}^2. \tag{3.95d}$$

Then the following estimate is valid for all $\varepsilon \geq 0$

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 + \varepsilon \|\tilde{e}_2\|_{L^2(H^1)}^2 \leq 3\|e^{0-}\|_{\mathcal{H}}^2 + 22(M_1 + M_2) + 20 \max_{1 \leq j \leq N} (M_{3,j} + M_{4,j}). \tag{3.96}$$

Moreover if $\varepsilon = 0$, then there holds additionally

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 \leq 2\|e^{0-}\|_{\mathcal{H}}^2 + 16(M_1 + M_2) + 16 \max_{1 \leq j \leq N} (M_{3,j} + M_{4,j}). \tag{3.97}$$

**Proof.** Let $t \in [0, T]$ be some time point where $\|\tilde{e}(t)\|_{\mathcal{H}} = \|\tilde{e}\|_{L^\infty(\mathcal{H})}$. Due to the fact that $\tilde{e}$ is not an affine function in time, we have that either

  1. $t = t_\ell$ for some $\ell = 1, \ldots, N$, or

  2. $t \in (t_{\ell-1}, t_\ell)$ for some $\ell = 1, \ldots, N$, i.e. $t = t_{\ell-1} + \delta$, $\delta < k_\ell$.

In the following we discuss both cases separately.

  1. In the first case given representation from Remark 3.3.1.1 where $V := \Pi \bar{\tilde{e}}$ we may see that there holds for each $n = 1, \ldots, N$

$$\frac{1}{2}\|e^{n-}\|_{\mathcal{H}}^2 + \varepsilon \sum_{j=1}^{n} \int_{I_j} \|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 dt \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \sum_{j=1}^{n} \int_{I_j} \langle \widetilde{R}_j ; \tilde{e} - \Pi \bar{\tilde{e}} \rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \int_{I_j} \widetilde{J}_j(U; \tilde{e} - \Pi \bar{\tilde{e}}) dt$$

$$+ \sum_{j=1}^{n} \int_{I_j} a(\widetilde{U}_2 - U_2; \tilde{e}_1) dt + \sum_{j=1}^{n} \int_{I_j} a(U_1 - \widetilde{U}_1; \tilde{e}_2) dt$$

$$+ \varepsilon \sum_{j=1}^{n} \int_{I_j} a(U_2 - \widetilde{U}_2; \tilde{e}_2) dt. \tag{3.98}$$

Furthermore from (3.95) with $\alpha = 2$ and the fact that (3.98) is also valid for $n = \ell$, where $\|\tilde{e}^{\ell-}\|_{\mathcal{H}} = \|\tilde{e}\|_{L^\infty(\mathcal{H})}$, we obtain

$$\frac{1}{2}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + 2M_1 + 2M_2 + \frac{1}{4}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2. \tag{3.99}$$

Moving the last term onto the LHS of (3.99) proves

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 \leq 2\|e^{0-}\|_{\mathcal{H}}^2 + 8M_1 + 8M_2. \tag{3.100}$$

If $\varepsilon = 0$, the estimate (3.100) holds. We continue by estimating the damping term in case of $\varepsilon > 0$. This also hold for the case $\varepsilon = 0$.

For $n = N$ in (3.98), the estimation with (3.95) for $\alpha = 2$ and with (3.100), yields

$$\varepsilon \sum_{j=1}^{N} \int_{I_j} \|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + 2M_1 + 2M_2 + \frac{1}{4}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 \leq \|e^{0-}\|_{\mathcal{H}}^2 + 4M_1 + 4M_2. \quad (3.101)$$

Summing the estimates (3.100) and (3.101), we finally obtain for $\varepsilon \geq 0$

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})} + \varepsilon\|\tilde{e}_2\|_{L^2(H^1)}^2 \leq 3\|e^{0-}\|_{\mathcal{H}}^2 + 12M_1 + 12M_2. \quad (3.102)$$

2. In case when $t$ is not a node from the triangulation $\mathcal{T}$, the fundamental theorem of calculus yields

$$\frac{1}{2}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 - \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 = \int_0^t \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} dt = \int_0^{t_{\ell-1}} \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} dt + \int_{t_{\ell-1}}^t \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} d\tau. \quad (3.103)$$

From the error representation formula (3.94) where $V = \Pi\bar{\tilde{e}}$, and an application of the fundamental theorem of calculus, we have for all $n = 1, \ldots, N$

$$\int_0^{t_n} \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} dt = \sum_{j=1}^{n} \int_{I_j} \langle \widetilde{R}_j; \tilde{e} - \Pi\bar{\tilde{e}} \rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \int_{I_j} \widetilde{J}_j(U; \tilde{e} - \Pi\bar{\tilde{e}}) dt + \sum_{j=1}^{n} \int_{I_j} a(\widetilde{U}_2 - U_2; \tilde{e}_1) dt$$

$$+ \sum_{j=1}^{n} \int_{I_j} a(U_1 - \widetilde{U}_1; \tilde{e}_2) dt + \varepsilon \sum_{j=1}^{n} \int_{I_j} a(U_2 - \widetilde{U}_2; \tilde{e}_2) dt$$

$$- \varepsilon \sum_{j=1}^{n} \int_{I_j} \|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 dt. \quad (3.104)$$

Then, if we assume that $n = \ell - 1$ in (3.104), the first term on the RHS of (3.103) can be estimated by (3.95) with $\alpha = 4$ such that

$$\int_0^{t_{\ell-1}} \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} dt \leq 4M_1 + 4M_2 + \frac{1}{8}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2. \quad (3.105)$$

With the definition of the residual where the jumps with respect to time variable are included in the definition of the volume term $\widetilde{R} := F - \dot{\widetilde{U}} - \mathcal{A}U$, see Lemma 3.3.1.2, the second term on the RHS of (3.103) is equivalent to

$$\int_{t_{\ell-1}}^t \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} d\tau = \int_{t_{\ell-1}}^t \langle \widetilde{R}_\ell; \tilde{e} \rangle_{\mathcal{H}} d\tau + \int_{t_{\ell-1}}^t \widetilde{J}_\ell(U; \tilde{e}) d\tau + \int_{t_{\ell-1}}^t a(\widetilde{U}_2 - U_2; \tilde{e}_1) d\tau$$

$$+ \int_{t_{n-1}}^t a(U_1 - \widetilde{U}_1; \tilde{e}_2) d\tau + \varepsilon \int_{t_{\ell-1}}^t a(U_2 - \widetilde{U}_2; \tilde{e}_2) d\tau - \varepsilon \int_{t_{\ell-1}}^t \|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 d\tau.$$

From (3.95) with $\alpha = 4$ and the fact that $\varepsilon \geq 0$, we obtain

$$\int_{t_{\ell-1}}^t \langle \dot{\tilde{e}}; \tilde{e} \rangle_{\mathcal{H}} d\tau \leq 4M_{3,\ell} + 4M_{4,\ell} + \frac{1}{8}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2. \quad (3.106)$$

The sum of (3.105) and (3.106) combined with the representation (3.103) yield

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 \leq 2\|e^{0-}\|_{\mathcal{H}}^2 + 16(M_1 + M_2) + 16(M_{3,n} + M_{4,n})$$
$$\leq 2\|e^{0-}\|_{\mathcal{H}}^2 + 16(M_1 + M_2) + 16 \max_{1 \leq j \leq N}(M_{3,j} + M_{4,j}). \tag{3.107}$$

For $\varepsilon = 0$, this concludes the proof. For $\varepsilon > 0$, we proceed as follows. Given (3.104) with $n = N$, the fundamental theorem of calculus combined with (3.95) for $\alpha = 2$, yields

$$\varepsilon\sum_{j=1}^{N}\int_{I_j}\|\tilde{e}_2(t)\|_{H^1(\Omega)}^2 dt \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + 2(M_1 + M_2) + \frac{1}{4}\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2$$
$$\leq \|e^{0-}\|_{\mathcal{H}}^2 + 6(M_1 + M_2) + 4\max_{1 \leq j \leq N}(M_{3,n} + M_{4,n}), \tag{3.108}$$

where we used the estimate (3.107) in the second inequality above.
Combination of (3.107) and (3.108) proves

$$\|\tilde{e}\|_{L^\infty(\mathcal{H})}^2 + \varepsilon\|\tilde{e}_2\|_{L^2(H^1)}^2 \leq 3\|e^{0-}\|_{\mathcal{H}}^2 + 22(M_1 + M_2) + 20\max_{1 \leq j \leq N}(M_{3,j} + M_{4,j}). \tag{3.109}$$

In general, it is not easy to determine whether $t$ is mesh point or not. Therefore, we have to compare the error estimates (3.102) and (3.109), for $\varepsilon > 0$ and estimates (3.100) and (3.107) for $\varepsilon = 0$. Obviously, the RHS of (3.102) is smaller then the RHS in (3.109) and the analogous result holds in case $\varepsilon = 0$. Therefore, (3.107) and (3.109) are the estimates valid in both cases, i.e. independent of whether $t$ is mesh point or not. This concludes the proof of the Lemma.□

### 3.3.1.1 A posteriori energy error analysis, $dG(0) \otimes \mathcal{P}_1$

In the following we consider whether Lemma 3.3.1.3 can be applied in case of the $dG(0) \otimes \mathcal{P}_1$ approximation which would yield an a posteriori error estimate. To check the assumptions (3.95), we should estimate

$$E_1 := \sum_{j=1}^{N}\int_{I_j}\langle\widetilde{R}_j\,;\tilde{e} - \Pi\bar{\tilde{e}}\rangle_{\mathcal{H}}dt + \sum_{j=1}^{N}\int_{I_j}\widetilde{J}_j(U, \tilde{e} - \Pi\tilde{e})dt,$$

$$E_2 := \int_{I_j}a(\widetilde{U}_2 - U_2; \tilde{e}_1)dt + \int_{I_j}a(U_1 - \widetilde{U}_1 + \varepsilon(U_2 - \widetilde{U}_2); \tilde{e}_2)dt,$$

$$E_3 := \int_{I_j}\langle\widetilde{R}_j\,;\tilde{e}\rangle_{\mathcal{H}}dt + \int_{I_j}\widetilde{J}_j(U, \tilde{e})dt,$$

$$E_4 := \int_{I_j}a(\widetilde{U}_2 - U_2; \tilde{e}_1)dt + \int_{I_j}a(U_1 - \widetilde{U}_1 + \varepsilon(U_2 - \widetilde{U}_2); \tilde{e}_2)dt,$$

by $\|\tilde{e}\|_{L^\infty(\mathcal{H})}$. Unfortunately, we did not succeed to derive an upper bound for $E_1, E_2, E_3, E_4$ in terms of $\|\tilde{e}\|_{L^\infty(\mathcal{H})}$. For instance, we could prove that

$$E_1 \leq \left(\sum_{j=1}^{N}\|(I-\mathcal{G})U_1^{j-1}\|_{H^1(\Omega)}\right)\|\tilde{e}_1\|_{L^\infty(H^1)} + \|(I-\mathcal{L})(f + \frac{1}{k_j}U_2^{j-1})\|_{L^1(L^2)}\|\tilde{e}_2\|_{L^\infty(L^2)}$$
$$+ C\|D_{h,1}(U_1 + \varepsilon U_2)\|_{L^2(0,T)}\|\tilde{e}_2\|_{L^2(H^1)} + \|f - \bar{f}\|_{L^1(L^2)}\|\tilde{e}_2\|_{L^\infty(L^2)}, \tag{3.111}$$

with a numerical constant $C \geq 0$. However, it is clear that the term $\|\tilde{e}_2\|_{L^2(H^1)}$ can not be dominated by $\|\tilde{e}_2\|_{L^\infty(L^2)}$. The problem is the estimation of the jump term in space. Here we used the techniques as before, e.g. trace and Friedrichs inequality, inverse estimates and we could not obtain a bound in terms of $\|\tilde{e}_2\|_{L^2(\Omega)}$. Another possibility would be to use the nodal interpolation operator $\mathcal{I}$ instead of $L^2$ projection operator $\mathcal{L}$, but according to Remark 3.1.0.3, the nodal interpolation operator $\mathcal{I}$ is not $L^2$ stable. We therefore skip the verification of (3.111). The same difficulty arises if we aim to estimate $E_2, E_3$ and $E_4$ instead of $E_1$.

Obviously, the choice of the discrete space may be relevant for the validity of the error estimate. Namely, by means of the discretisation in the higher degree spaces (concerning the spatial variable), we may be able to avoid the jump terms. This idea will be further developed in Subsection 3.3.1.2.

**Remark 3.3.1.2.** The idea applied in Lemma 3.3.1.3 was to estimate residual and other discrete solution contributions with respect to $\|\tilde{e}\|_{L^\infty(\mathcal{H})}$ in order to absorb the error terms through the $\|\tilde{e}\|_{L^\infty(\mathcal{H})}$. Note that $\|\tilde{e}_2\|_{L^2(H^1)}$ could also be absorbed through the viscosity term $\varepsilon\|\tilde{e}_2(t)\|^2_{L^2(H^1)}$ from the LHS only account to the loss in the final estimate. Namely in the final estimate we would then have some residual contributions in corresponding norm multiplied with $1/\sqrt{\varepsilon}$. For $\varepsilon \to 0$ factor $1/\sqrt{\varepsilon} \to \infty$ and this yields the estimate which is not sharp enough to be further used. $\qquad\square$

### 3.3.1.2 A posteriori energy error analysis, $dG(0)\otimes\mathcal{C}^1$

Within this subsection we analyse whether the result of Lemma 3.3.1.3 can be applied in case of the $dG(0)\otimes\mathcal{C}_1$ approximation.

**Theorem 3.3.1.1 (A posteriori energy error estimate, $dG(0)\otimes\mathcal{C}^1$).** For $u$ and its discrete $dG(0)\otimes\mathcal{C}^1$ counterpart $U$, there holds for all $\varepsilon \geq 0$

$$
\|\tilde{e}\|^2_{L^\infty(\mathcal{H})} + \varepsilon\|\tilde{e}\|^2_{L^2(H^1)} \leq 3\|y_0 - \mathcal{I}y_0\|^2_{H^1(\Omega)} + 3\|y_1 - \mathcal{I}y_1\|^2_{L^2(\Omega)} + 22\Big(\sum_{j=1}^N \|(I-\mathcal{G})U_1^{j-1}\|_{H^1(\Omega)}\Big)^2
$$

$$
+ 44\|f - \bar{f}\|^2_{L^1(L^2)} + 44\|(I-\mathcal{L})(\bar{f} + \frac{1}{k_j}U_2^{j-1} + \Delta(U_1 + \varepsilon U_2))\|^2_{L^1(L^2)}
$$

$$
+ 22\|\widetilde{U}_2 - U_2\|^2_{L^1(H^1)} + 22\|\Delta(\widetilde{U}_1 - U_1 + \varepsilon(\widetilde{U}_2 - U_2))\|^2_{L^1(L^2)}
$$

$$
+ 20 \max_{1 \leq j \leq N} \Big\{ \|U_2 - \dot{\widetilde{U}}_1\|^2_{L^1(I_j;H^1(\Omega))} + \|\widetilde{U}_2 - U_2\|^2_{L^1(I_j;H^1(\Omega))}
$$

$$
+ \|f - \dot{\widetilde{U}}_2 + \Delta(U_1 + \varepsilon U_2)\|^2_{L^1(I_j;L^2(\Omega))}
$$

$$
+ \|\Delta(\widetilde{U}_1 - U_1 + \varepsilon(\widetilde{U}_2 - U_2))\|^2_{L^1(I_j;L^2(\Omega))} \Big\}. \tag{3.112}
$$

Moreover if $\varepsilon = 0$, then there holds additionally

$$\|\tilde{e}\|^2_{L^\infty(\mathcal{H})} \le 2\|y_0 - \mathcal{I}y_0\|^2_{H^1(\Omega)} + 2\|y_1 - \mathcal{I}y_1\|^2_{L^2(\Omega)} + 16\Big(\sum_{j=1}^{N} \|(I-\mathcal{G})U_1^{j-1}\|_{H^1(\Omega)}\Big)^2$$

$$+ 32\|f - \bar{f}\|^2_{L^1(L^2)} + 32\|(I-\mathcal{L})(\bar{f} + \frac{1}{k_j}U_2^{j-1} + \Delta U_1)\|^2_{L^1(L^2)} + 16\|\widetilde{U}_2 - U_2\|^2_{L^1(H^1)}$$

$$+ 16\|\Delta(\widetilde{U}_1 - U_1)\|^2_{L^1(L^2)} + 16 \max_{1 \le j \le N} \Big\{ \|U_2 - \dot{\widetilde{U}}_1\|^2_{L^1(I_j;H^1(\Omega))} + \|f - \dot{\widetilde{U}}_2 + \Delta U_1\|^2_{L^1(I_j;L^2(\Omega))}$$

$$+ \|\widetilde{U}_2 - U_2\|^2_{L^1(I_j;H^1(\Omega))} + \|\Delta(\widetilde{U}_1 - U_1)\|^2_{L^1(I_j;L^2(\Omega))} \Big\}. \tag{3.113}$$

Note that if $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$, then $(I-\Pi)U^{j-1} = 0$ on the RHS of (3.112) and (3.113). Moreover, the RHS of both estimates is denoted by $\eta^2$.

**Remark 3.3.1.3.** For sufficiently smooth initial data $f, y_0, y_1$, the estimate of Theorem 3.3.1.1 is of order $\mathcal{O}(h^3 + k)$ when $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j = 1, \ldots, N$ and $\mathcal{O}(h^3 + k^{-1}h^3 + k)$ otherwise.  $\square$

**Proof.** The proof follows from an application of Lemma 3.3.1.3. We therefore check the assumptions of the same in the following.

**Lemma 3.3.1.4.** With

$$M_1 := 2\|f - \bar{f}\|^2_{L^1(L^2)} + 2\|(I-\mathcal{L})(\bar{f} + \frac{1}{k_j}U_2^{j-1} + \Delta(U_1 + \varepsilon U_2))\|^2_{L^1(L^2)} + \Big(\sum_{j=1}^{N} \|(I-\mathcal{G})U_1^{j-1}\|_{H^1(\Omega)}\Big)^2,$$

there holds for $n = 1, \ldots, N$ and both $\alpha \in \{2, 4\}$,

$$E_1 := \sum_{j=1}^{n} \int_{I_j} \langle \widetilde{R}_j ; \tilde{e} - \Pi\bar{\tilde{e}} \rangle_{\mathcal{H}} dt \le \alpha M_1 + \frac{1}{4\alpha}\|\tilde{e}\|_{L^\infty(\mathcal{H})}. \tag{3.114}$$

Note that $E_1$ is the LHS of (3.95a) since the jump term vanishes for $\mathcal{C}^1$ functions.

**Proof.** Due to the orthogonal properties of the projections $\mathcal{G}, \mathcal{L}$ and Lemma 3.3.1.2, we have

$$E_1 = \sum_{j=1}^{N} \int_{I_j} \langle \widetilde{R}_j - \overline{\widetilde{R}}_j ; \tilde{e} \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \int_{I_j} \langle \overline{\widetilde{R}}_j - \Pi\overline{\widetilde{R}}_j ; \tilde{e} \rangle_{\mathcal{H}} dt$$

$$= \sum_{j=1}^{N} \int_{I_j} \langle F - \overline{F} ; \tilde{e} \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \int_{I_j} \langle \overline{F} - \Pi\overline{F} - \dot{\widetilde{U}} + \Pi\dot{\widetilde{U}} - \mathcal{A}U + \Pi\mathcal{A}U ; \tilde{e} \rangle_{\mathcal{H}} dt.$$

Note that $(\Pi\dot{\widetilde{U}} - \dot{\widetilde{U}})|_{I_j} = {}^1\!/_{k_j}(U^{j-1} - \Pi|_{I_j}U^{j-1})$ and this term equals zero only in case of hierarchical grids, i.e. when $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$.

Expanding the last equation in terms of the $\mathcal{H}$ scalar product, we obtain

$$E_1 = \sum_{j=1}^{N} \int_{I_j} (f - \bar{f}; \tilde{e}_2) dt + \sum_{j=1}^{N} \int_{I_j} a((I-\mathcal{G})\frac{1}{k_j}U_1^{j-1}; \tilde{e}_1) dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} ((I-\mathcal{L})(\bar{f} + \frac{1}{k_j}U_2^{j-1} + \Delta(U_1 + \varepsilon U_2)); \tilde{e}_2) dt.$$

The Hölder inequality in time and in space yields

$$E_1 \leq \|f - \bar{f}\|_{L^1(L^2)} \|\tilde{e}_2\|_{L^\infty(L^2)} + \Big( \sum_{j=1}^N \|(I - \mathcal{G})U_1^{j-1}\|_{H^1(\Omega)} \Big) \|\tilde{e}_1\|_{L^\infty(H^1)}$$

$$+ \|(I - \mathcal{L})(\bar{f} + \frac{1}{k_j}U_2^{j-1} + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)} \|\tilde{e}_2\|_{L^\infty(L^2)}.$$

With the aid of the Young inequality, we conclude for $\alpha \in \{2, 4\}$

$$E_1 \leq \alpha \Big( \sum_{j=1}^N \|(I - \mathcal{G})U_1^{j-1}\|_{H^1(\Omega)} \Big)^2 + 2\alpha \|f - \bar{f}\|_{L^1(L^2)}^2$$

$$+ 2\alpha \|(I - \mathcal{L})(\bar{f} + \frac{1}{k_j}U_2^{j-1} + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}^2 + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}^2. \qquad \square$$

**Lemma 3.3.1.5.** For all $n = 1, \ldots, N$ and $M_2$ and $M_{4,n}$ defined by

$$M_2 := \|\widetilde{U}_2 - U_2\|_{L^1(H^1)}^2 + \|\Delta(\widetilde{U}_1 - U_1 + \varepsilon(\widetilde{U}_2 - U_2))\|_{L^1(L^2)}^2,$$

$$M_{4,n} := \|\widetilde{U}_2 - U_2\|_{L^1(I_j; H^1(\Omega))}^2 + \|\Delta(\widetilde{U}_1 - U_1 + \varepsilon(\widetilde{U}_2 - U_2))\|_{L^1(I_j; L^2(\Omega))}^2.$$

there holds for both $\alpha \in \{2, 4\}$

$$E_2 := \sum_{j=1}^n \int_{I_j} a(\widetilde{U}_2 - U_2; \tilde{e}_1)dt + \sum_{j=1}^n \int_{I_j} a(U_1 - \widetilde{U}_1 + \varepsilon(U_2 - \widetilde{U}_2); \tilde{e}_2)dt \leq \alpha M_2 + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}, \quad (3.116a)$$

and also

$$E_4 := \int_{I_n} a(\widetilde{U}_2 - U_2; \tilde{e}_1)dt + \int_{I_n} a(U_1 - \widetilde{U}_1 + \varepsilon(U_2 - \widetilde{U}_2); \tilde{e}_2)dt \leq \alpha M_{4,n} + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}. \qquad (3.116b)$$

**Proof.** It is obvious that if we prove (3.116b), the estimate (3.116a) follows automatically.

The Hermite cubic polynomials are globally $\mathcal{C}^1$ functions in space which satisfy the (DD) or (DN) boundary conditions. Then, the partial integration in space yields

$$E_4 = \int_{I_n} a(\widetilde{U}_2 - U_2; \tilde{e}_1)dt + \int_{I_n} (\Delta(\widetilde{U}_1 - U_1 + \varepsilon(\widetilde{U}_2 - U_2)); \tilde{e}_2)dt.$$

With Hölder inequality we obtain

$$E_4 \leq \|\widetilde{U}_2 - U_2\|_{L^1(I_n; H^1(\Omega))} \|\tilde{e}_1\|_{L^\infty(H^1)} + \|\Delta(\widetilde{U}_1 - U_1 + \varepsilon(\widetilde{U}_2 - U_2))\|_{L^1(I_n; L^2(\Omega))} \|\tilde{e}_2\|_{L^\infty(L^2)}. \quad (3.117)$$

With the aid of the Young inequality we prove (3.116b). The estimate (3.116a) can be proven by taking a sum over all $n = 1, \ldots, N$ in (3.117). $\qquad \square$

**Lemma 3.3.1.6.** For $n = 1, \ldots, N$ and $M_{3,n}$ defined as

$$M_{3,n} := \|U_2 - \dot{\widetilde{U}}_1\|_{L^1(I_n; H^1(\Omega))}^2 + \|f - \dot{\widetilde{U}}_2 - \Delta(U_1 + \varepsilon U_2)\|_{L^1(I_n; L^2(\Omega))}^2,$$

there holds for both $\alpha \in \{2, 4\}$,

$$E_3 := \int_{I_n} \langle \widetilde{R}_n ; \tilde{e} \rangle_{\mathcal{H}} dt \leq \alpha M_{3,n} + \frac{1}{4\alpha} \|\tilde{e}\|_{L^\infty(\mathcal{H})}. \qquad (3.118)$$

Note that $E_3$ is the LHS of (3.95c) since the jump term vanishes for $\mathcal{C}^1$ functions.

**Proof.** Given $E_3$ from (3.118), the residual representation from Lemma 3.3.1.2 implies

$$E_3 = \int_{I_n} \left\langle F - \dot{\tilde{U}} - \mathcal{A}U ; \tilde{e} \right\rangle_{\mathcal{H}} dt$$
$$= \int_{I_n} a(U_2 - \dot{\tilde{U}}_1 ; \tilde{e}_1) dt + \int_{I_n} (f - \dot{\tilde{U}}_2 - \Delta(U_1 + \varepsilon U_2) ; \tilde{e}_2) dt.$$

An application of the Hölder inequality in time and in space shows

$$E_3 \leq \|U_2 - \dot{\tilde{U}}_1\|_{L^1(I_n ; H^1(\Omega))} \|\tilde{e}_1\|_{L^\infty(H^1)} + \|f - \dot{\tilde{U}}_2 - \Delta(U_1 + \varepsilon U_2)\|_{L^1(I_n ; L^2(\Omega))} \|\tilde{e}_2\|_{L^\infty(L^2)}.$$

The Young inequality for $\alpha \in \{2, 4\}$ yields the proof.                                                    $\square$

Latter lemmas have shown that Lemma 3.3.1.3 can be applied in case of the $dG(0) \otimes \mathcal{C}^1$ approximation. However, if we recall the estimates (3.96) and (3.97), it remains to estimate term $\|e^{0-}\|_{\mathcal{H}}$ in order to complete the proof of theorem.

**Lemma 3.3.1.7.** Let the discrete variant of initial solution be defined as

$$U^{0-} := (\mathcal{I}y_0, \mathcal{I}y_1)$$

where $u_0 = (y_0, y_1)$ is the initial solution and $\mathcal{I}$, the Hermite cubic interpolant. Then there holds,

$$\|e^{0-}\|_{\mathcal{H}}^2 = \|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)}^2 + \|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)}^2. \tag{3.119}$$

**Proof.** Proof follows directly from

$$\|e^{0-}\|_{\mathcal{H}}^2 = \|u_0 - U^{0-}\|_{\mathcal{H}}^2 = \|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)}^2 + \|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)}^2. \tag{3.120}$$
$$\square$$

This concludes the proof of theorem.                                                    $\square$

## 3.3.2   $dG(1)$ time approximation

The main argument in the following error analysis is the use of the residual $Res$ from (2.7), where $\mathcal{B}$ stands for the bilinear form (2.18) and $\mathscr{L}$ is defined as in (2.19). We derive an error bound $\eta$ for

$$\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon \|e_2\|_{L^2(H^1)}^2.$$

Contrariwise to the analysis developed for the $dG(0)$ time approximation presented in Subsection 3.3.1, the error function $\tilde{e}$ and the affine interpolant of the discrete solution $\widetilde{U}$ are not involved in the error analysis when $dG(1)$ time approximation is considered. This is due to the fact that here the time projection $\mathcal{J}$, cf. Definition 3.1.0.9, case $dG(1)a$, enables the optimal estimation of the temporal jump terms. Namely, the projection $\mathcal{J}u$ and $u$ coincide in each $t_{j-1}$ on $I_j$.

**Lemma 3.3.2.1 (Error representation, $dG(1)$ time approximation).** In case of $dG(1)$ time discretisation, there holds for every $V \in \mathcal{Q}_1$

$$\frac{1}{2}\|e^{N-}\|_{\mathcal{H}}^2 + \frac{1}{2}\sum_{j=1}^{N}\|[e]^{j-1}\|_{\mathcal{H}}^2 + \varepsilon\sum_{j=1}^{N}\int_{I_j}\|e_2(t)\|_{H^1(\Omega)}^2 = \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + Res(e - V). \tag{3.121}$$

**Proof.** The fundamental theorem of calculus leads to

$$\frac{1}{2}\sum_{j=1}^{N}\|e^{j-}\|_{\mathcal{H}}^2-\frac{1}{2}\sum_{j=1}^{N}\|e^{j-1+}\|_{\mathcal{H}}^2=\sum_{j=1}^{N}\int_{I_j}\langle e_\tau\,;e\rangle_{\mathcal{H}}dt. \tag{3.122}$$

Given the residual representation (2.10) and the orthogonality (2.9), for arbitrary test function $V\in\mathcal{Q}_1$ we have

$$\sum_{j=1}^{N}\int_{I_j}\langle e_\tau\,;e\rangle_{\mathcal{H}}dt=\sum_{j=1}^{N}\int_{I_j}\langle e_\tau\,;e-V\rangle_{\mathcal{H}}dt+\sum_{j=1}^{N}\int_{I_j}a(e_2;V_1)dt-\sum_{j=1}^{N}\int_{I_j}(e_1;V_2)dt$$

$$-\varepsilon\sum_{j=1}^{N}\int_{I_j}(e_2;V_2)dt-\sum_{j=1}^{N}\langle[e]^{j-1}\,;V^{j-1+}\rangle_{\mathcal{H}}. \tag{3.123}$$

From (3.122) and the second identity from Lemma 2.3.3.3, the last equation (3.123) simplifies to

$$\frac{1}{2}\sum_{j=1}^{N}\|e^{j-}\|_{\mathcal{H}}^2-\frac{1}{2}\sum_{j=1}^{N}\|e^{j-1+}\|_{\mathcal{H}}^2=Res(e-V)-\varepsilon\sum_{j=1}^{N}\int_{I_j}a(e_2;e_2)dt-\sum_{j=1}^{N}\langle[e_1]^{j-1}\,;e^{j-1+}\rangle_{\mathcal{H}}$$

$$=Res(e-V)-\varepsilon\sum_{j=1}^{N}\int_{I_j}a(e_2;e_2)dt$$

$$+\frac{1}{2}\sum_{j=1}^{N}\|e^{j-1-}\|_{\mathcal{H}}^2-\frac{1}{2}\sum_{j=1}^{N}\|[e]^{j-1}\|_{\mathcal{H}}^2-\frac{1}{2}\sum_{j=1}^{N}\|e^{j-1+}\|_{\mathcal{H}}^2.$$

Moving the last four terms on the RHS of latter equation to the LHS of the same, we may prove the lemma. $\qquad\square$

**Lemma 3.3.2.2.** In case of $dG(q),q=0,1$ time approximation and $\mathcal{P}_1$, i.e. $\mathcal{C}_1$ approximation in space, the residual (2.7) may be expressed in terms of local residuals and local jumps such that for all $v\in L^2(H_D^1\times H_D^1)$

$$Res(v)=\sum_{j=1}^{N}\int_{I_j}\langle R_j\,;v\rangle_{\mathcal{H}}dt+\sum_{j=1}^{N}\int_{I_j}J_j(U,v)dt, \tag{3.124a}$$

where for all $t\in[t_{j-1},t_{j-1}+\delta)$, $0<\delta\le k_j$

$$\langle R_j\,;v\rangle_{\mathcal{H}}(t):=\langle F\,;v\rangle_{\mathcal{H}}(t)-\langle U_\tau\,;v\rangle_{\mathcal{H}}(t)-\langle\mathcal{A}U\,;v\rangle_{\mathcal{H}}(t)-\langle[U]^{j-1}\,;v^{j-1+}\rangle_{\mathcal{H}}, \tag{3.124b}$$

$$J_j(U,v)(t):=\begin{cases}\sum_{k=1}^{m}[D(U_1+\varepsilon U_2)(t)]_kv_2(t,x_k), & \text{for }\mathcal{P}_1\text{ elements,}\\ 0, & \text{for }\mathcal{C}^1\text{ elements.}\end{cases} \tag{3.124c}$$

Here $m=n,n-1$ in case of Problem (DN), (DD) respectively. The jump terms $[D(U_1+\varepsilon U_2)]_k$ are defined in Definition 3.1.0.6.

**Proof.** An integration by parts in space in (2.21) with the bilinear form $\mathcal{B}$ as in (2.18) and $\mathscr{L}$ from (2.19), yields for all $v \in L^2(H_D^1 \times H_D^1)$

$$Res(v) = \sum_{j=1}^{N} \int_{I_j} (f; v_2)dt - \sum_{j=1}^{N} \int_{I_j} a(U_{1,\tau}; v_1)dt - \sum_{j=1}^{N} \int_{I_j} (U_{2,\tau}; v_2)dt + \sum_{j=1}^{N} \int_{I_j} a(U_2; v_1)dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} (\Delta U_1; v_2)dt + \varepsilon \sum_{j=1}^{N} \int_{I_j} (\Delta U_2; v_2)dt + \sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)]_k v_2(x_k)dt$$

$$- \sum_{j=1}^{N} a([U_1]^{j-1}; v_1^{j-1+}) - \sum_{j=1}^{N} ([U_2]^{j-1}; v_2^{j-1+}).$$

The jump terms in space $[D(U_1 + \varepsilon U_2)]_k$ equal zero for all $k = 1, \ldots, m$ in case of spatial approximation by $\mathcal{C}^1$ functions due to the continuity of discrete function in first derivative and the fact that discrete solution satisfies the (DD) as well as (DN) boundary conditions, see Subsections 6.1.3, 6.1.4. With the notation from Definition 1.3.0.7, we conclude the proof of the lemma.                                                                                                  $\square$

**Remark 3.3.2.1.** If $\mathscr{T}$ is an arbitrary discretisation of the time interval $[0, T]$, then from the definition of the local weak problem (2.6) and the residual representation in Lemma 3.3.2.2, there holds for each $t_n \in \mathscr{T}$

$$\frac{1}{2}\|e^{n-}\|_{\mathcal{H}}^2 + \frac{1}{2}\sum_{j=1}^{n} \|[e]^{j-1}\|_{\mathcal{H}}^2 + \varepsilon \sum_{j=1}^{n} \int_{I_j} \|e_2(t)\|_{H^1(\Omega)}^2 dt = \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \sum_{j=1}^{n} \int_{I_j} \langle R_j; e - V\rangle_{\mathcal{H}} dt$$

$$+ \sum_{j=1}^{n} \int_{I_j} J_j(U; e - V)dt. \qquad (3.125)$$

Here $V$ is an arbitrary test function, $V|_{I_j} \in \mathcal{Q}_1^j$ for all $j = 1, \ldots, n$ for which the residual orthogonality (2.8) applies.                                                                                                  $\square$

The following lemma will be used latter on for the $\mathcal{P}_1$ and $\mathcal{C}^1$ approximation in space separately.

**Lemma 3.3.2.3.** Let $R_j$ and $J_j$ be defined by (3.124b) and (3.124c), respectively. Let $e$ be the error of $dG(1)$ time approximation and $\mathcal{P}_1$ or $\mathcal{C}^1$ discretisation in space. Let $\Pi := (\mathcal{G}, \mathcal{L})$ be a spatial orthogonal projection with respect to the $\mathcal{H}$ scalar product and let also $\mathcal{J}$ be a temporal $L^2$ projection from Definition 3.1.0.9, case $dG(1)a$. If there exists some positive $M_1, M_{2,n}$ such that there holds for all $n = 1, \ldots, N$ and both $\alpha = \{1, 2\}$

$$\sum_{j=1}^{n} \int_{I_j} \langle R_j; e - \mathcal{J}\Pi e\rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \int_{I_j} J_j(U; e - \mathcal{J}\Pi e)dt \le \alpha M_1 + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2, \qquad (3.126a)$$

$$\int_{I_n} \langle R_n; e\rangle_{\mathcal{H}} dt + \int_{I_n} J_n(U; e)dt \le \alpha M_{2,n} + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2. \qquad (3.126b)$$

Then the following estimate is valid for all $\varepsilon \ge 0$

$$\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon\|e_2\|_{L^2(H^1)}^2 \le 3\|e^{0-}\|_{\mathcal{H}}^2 + 11M_1 + 10 \max_{1 \le j \le N} M_{2,j}. \qquad (3.127)$$

Moreover for $\varepsilon = 0$ there holds the sharper estimate

$$\|e\|_{L^\infty(\mathcal{H})}^2 \le 2\|e^{0-}\|_{\mathcal{H}}^2 + 8M_1 + 8 \max_{1 \le j \le N} M_{2,j}. \qquad (3.128)$$

**Proof.** Similar as in the proof of Lemma 3.3.1.3, let $t \in [0,T]$ be some time point where $\|e(t)\|_{\mathcal{H}} = \|e\|_{L^\infty(\mathcal{H})}$. Due to the fact that $e$ is not an affine, globally continuous function in time, it is difficult to determine whether $t$ is mesh point or not. Analogously to the proof of Lemma 3.3.1.3, where we proved that the estimate in case $t \notin \mathscr{T}$ also holds if $t \in \mathscr{T}$, we may proceed here by using the same arguments. Therefore, we treat only the case $t \notin \mathscr{T}$ in the following.

Let us assume that $t \in (t_{\ell-1}, t_\ell)$. Given the error representation formula (3.125) where $V = \mathcal{J}\Pi e$, from (3.126) there holds for all $n = 1, \ldots, N$ and $\alpha = \{1,2\}$

$$\frac{1}{2}\|e^{n-}\|_{\mathcal{H}}^2 + \varepsilon\sum_{j=1}^{n}\int_{I_j}\|e_2(t)\|_{H^1(\Omega)}^2 dt \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \sum_{j=1}^{n}\int_{I_j}\langle R_j\,;e-\Pi\bar{e}\rangle_{\mathcal{H}}dt + \sum_{j=1}^{n}\int_{I_j}J_j(U;e-\Pi\bar{e})dt$$

$$\leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + \alpha M_1 + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2. \tag{3.129}$$

The fundamental theorem of calculus, Lemma 3.3.2.2 and Lemma 2.3.3.3 imply

$$\frac{1}{2}\|e\|_{L^\infty(\mathcal{H})}^2 - \frac{1}{2}\|e^{\ell-1+}\|_{\mathcal{H}}^2 = \int_{t_{\ell-1}}^{t}\langle \dot{e}\,;e\rangle_{\mathcal{H}}d\tau$$

$$= \int_{t_{\ell-1}}^{t}\langle R_\ell\,;e\rangle_{\mathcal{H}}dt + \int_{t_{\ell-1}}^{t}J_\ell(U;e)d\tau - \varepsilon\int_{t_{\ell-1}}^{t}\|e_2(\tau)\|_{H^1(\Omega)}^2 d\tau - \langle [e]^{\ell-1}\,;e^{\ell-1+}\rangle$$

$$\leq \int_{t_{\ell-1}}^{t}\langle R_\ell\,;e\rangle_{\mathcal{H}}dt + \int_{t_{\ell-1}}^{t}J_\ell(U;e)d\tau + \frac{1}{2}\|e^{\ell-1-}\|_{\mathcal{H}}^2 - \frac{1}{2}\|[e]^{\ell-1}\|_{\mathcal{H}}^2 - \frac{1}{2}\|e^{\ell-1+}\|_{\mathcal{H}}^2.$$

According to (3.126), the last inequality simplifies to

$$\frac{1}{2}\|e\|_{L^\infty(\mathcal{H})}^2 \leq \frac{1}{2}\|e^{\ell-1-}\|_{\mathcal{H}}^2 + \alpha M_{2,\ell} + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2 \leq \frac{1}{2}\|e^{\ell-1-}\|_{\mathcal{H}}^2 + \alpha\max_{1\leq j\leq N}M_{2,j} + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2. \tag{3.130}$$

With (3.129) where $n = \ell-1$ and (3.130), owing to $\varepsilon \geq 0$ we obtain for $\alpha = 2$

$$\|e\|_{L^\infty(\mathcal{H})}^2 \leq 2\|e^{0-}\|_{\mathcal{H}}^2 + 8M_1 + 8\max_{1\leq j\leq N}M_{2,j}. \tag{3.131}$$

For $\varepsilon = 0$, this concludes the proof. For $\varepsilon \geq 0$, from (3.129) with $n = N$ and $\alpha = 1$ we have

$$\varepsilon\sum_{j=1}^{N}\int_{I_j}\|e_2(t)\|_{H^1(\Omega)}^2 dt \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + M_1 + \frac{1}{4}\|e\|_{L^\infty(\mathcal{H})}^2 \leq \|e^{0-}\|_{\mathcal{H}}^2 + 3M_1 + 2\max_{1\leq j\leq N}M_{2,j}. \tag{3.132}$$

Note that in the second inequality above, we used the estimate (3.131).

Finally, the combination of (3.131) and (3.132) proves

$$\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon\|e_2\|_{L^2(H^1)}^2 \leq 3\|e^{0-}\|_{\mathcal{H}}^2 + 11M_1 + 10\max_{1\leq j\leq N}M_{2,j}. \tag{3.133}$$

This concludes the proof of the Lemma. $\qquad\square$

### 3.3.2.1 A posteriori energy error analysis, $dG(1)\otimes\mathcal{P}_1$

Owing to the same complexity that obviated the derivation of an a posteriori error estimate for $dG(0)\otimes\mathcal{P}_1$ discretisation, cf. Subsection 3.3.1.1, we were not able to derive an a posteriori bound in case of $dG(1)\otimes\mathcal{P}_1$ discrete problem. This yields the conclusion that in case of $\mathcal{P}_1$ space discretisation by estimating according to the techniques of energy method, the error analysis and derivation of the a posteriori error bound can not be completed at least by using the result of Lemma 3.3.2.3.

### 3.3.2.2  A posteriori energy error analysis, $dG(1)\otimes\mathcal{C}^1$

Contrary to the $dG(1)\otimes\mathcal{P}^1$ According to results and conclusions from the Subsection 3.3.2.1, spatially $\mathcal{C}^1$ approximation is required in order to prove a sharp residual-based a posteriori error bound in the energy norm.

**Theorem 3.3.2.1 (A posteriori energy error estimate, $dG(1)\otimes\mathcal{C}^1$).** For $u$ and its discrete $dG(1)\otimes\mathcal{C}^1$ counterpart $U$, there holds for all $\varepsilon\geq 0$

$$\|e\|^2_{L^\infty(\mathcal{H})}+\varepsilon\|e_2\|^2_{L^2(H^1)}\leq 3\|y_0-\mathcal{I}y_0\|^2_{H^1(\Omega)}+3\|y_1-\mathcal{I}y_1\|^2_{L^2(\Omega)}+33\|(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(L^2)}$$
$$+33\|f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2))\|^2_{L^1(L^2)}$$
$$+33\Big(\sum_{j=1}^{N}\|(\mathcal{L}-I)U_2^{j-1-}\|_{L^2(\Omega)}\Big)^2+22\|U_2-\overline{U}_2\|^2_{L^1(H^1)}$$
$$+22\Big(\sum_{j=1}^{N}\|(\mathcal{G}-I)U_1^{j-1-}\|_{H^1(\Omega)}\Big)^2$$
$$+20\max_{1\leq j\leq N}\Big\{\|U_2-U_{1,\tau}\|^2_{L^1(I_j;H^1(\Omega))}+\|[U]^{j-1}\|^2_{\mathcal{H}}$$
$$+\|f-U_{2,\tau}+\Delta(U_1+\varepsilon U_2)\|^2_{L^1(I_j;L^2(\Omega))}\Big\}. \tag{3.134}$$

Moreover, for $\varepsilon=0$, then there holds the sharper estimate

$$\|e\|^2_{L^\infty(\mathcal{H})}\leq 2\|y_0-\mathcal{I}y_0\|^2_{H^1(\Omega)}+2\|y_1-\mathcal{I}y_1\|^2_{L^2(\Omega)}+24\|(I-\mathcal{L})(f+\Delta U_1)\|_{L^1(L^2)}$$
$$+24\|f-\bar{f}+\Delta(U_1-\overline{U}_1)\|^2_{L^1(L^2)}+24\Big(\sum_{j=1}^{N}\|(\mathcal{L}-I)U_2^{j-1-}\|_{L^2(\Omega)}\Big)^2+16\|U_2-\overline{U}_2\|^2_{L^1(H^1)}$$
$$+16\Big(\sum_{j=1}^{N}\|(\mathcal{G}-I)U_1^{j-1-}\|_{H^1(\Omega)}\Big)^2+16\max_{1\leq j\leq N}\Big\{\|U_2-U_{1,\tau}\|^2_{L^1(I_j;H^1(\Omega))}$$
$$+\|[U]^{j-1}\|^2_{\mathcal{H}}+\|f-U_{2,\tau}+\Delta U_1\|^2_{L^1(I_j;L^2(\Omega))}\Big\}. \tag{3.135}$$

Note that if $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, then $(I-\Pi)U^{j-1}=0$ on the RHS of (3.134) and (3.135). Moreover, the RHS of both estimates is denoted by $\eta^2$.

**Remark 3.3.2.2.** The estimate of Theorem 3.3.2.1 is of order $\mathcal{O}(h^3+k)$ when $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$ and $\mathcal{O}(h^3+k^{-1}h^3+k)$ otherwise. This holds for sufficiently smooth initial data $y_0,y_1,f$.  $\square$

**Proof.** The proof follows from an application of Lemma 3.3.2.3. We prove the validity of the assumptions (3.126) in the following two lemmas.

**Lemma 3.3.2.4.** With

$$M_1:=3\|(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|^2_{L^1(L^2)}+3\|f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2))\|^2_{L^1(L^2)}$$
$$+3\Big(\sum_{j=1}^{N}\|(\mathcal{L}-I)U_2^{j-1-}\|_{L^2(\Omega)}\Big)^2+2\|U_2-\overline{U}_2\|^2_{L^1(H^1)}+2\Big(\sum_{j=1}^{N}\|(\mathcal{G}-I)U_1^{j-1-}\|_{H^1(\Omega)}\Big)^2.$$

there holds for $n = 1, \ldots, N$ and both $\alpha = \{1, 2\}$

$$E_1 := \sum_{j=1}^{n} \int_{I_j} \langle R_j ; e - \mathcal{J}\Pi e \rangle_{\mathcal{H}} dt \leq \alpha M_1 + \frac{1}{4\alpha} \|e\|_{L^\infty(\mathcal{H})}^2. \tag{3.136}$$

Note that $E_1$ is the LHS of (3.126a) since the jump term vanishes for $\mathcal{C}^1$ functions.

**Proof.** Due to the orthogonal properties of the projections $\mathcal{G}, \mathcal{L}, \mathcal{J}$ and Lemma 3.3.2.2, we have

$$\begin{aligned}
E_1 &= \sum_{j=1}^{n} \int_{I_j} \langle R_j - \Pi R_j ; e - \Pi e \rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \int_{I_j} \langle R_j - \bar{R}_j ; \Pi e - \Pi \bar{e} \rangle_{\mathcal{H}} dt \\
&= \sum_{j=1}^{n} \int_{I_j} \langle F - \Pi F - \mathcal{A}U + \Pi \mathcal{A}U ; e - \Pi e \rangle_{\mathcal{H}} dt + \sum_{j=1}^{n} \langle \Pi U^{j-1-} - U^{j-1-} ; (e - \Pi e)^{j-1+} \rangle_{\mathcal{H}} \\
&\quad + \sum_{j=1}^{n} \int_{I_j} \langle F - \bar{F} - \mathcal{A}(U - \bar{U}) ; \Pi e - \Pi \bar{e} \rangle_{\mathcal{H}} dt. \tag{3.137}
\end{aligned}$$

Note that the term $\Pi U^{j-1} - U^{j-1}$ from the last equality above equals zero for $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$. Expanding the equation above due to the definition of the $\mathcal{H}$ scalar product, we obtain

$$\begin{aligned}
E_1 &= \sum_{j=1}^{n} \int_{I_j} ((I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2)); e_2) dt + \sum_{j=1}^{n} a((\mathcal{G} - I)U_1^{j-1-}; e_1^{j-1+}) \\
&\quad + \sum_{j=1}^{n} ((\mathcal{L} - I)U_2^{j-1-}; e_2^{j-1+}) + \sum_{j=1}^{n} \int_{I_j} a(U_2 - \bar{U}_2; \mathcal{G}e_1) dt \\
&\quad + \sum_{j=1}^{n} \int_{I_j} (f - \bar{f} + \Delta(U_1 - \bar{U}_1 + \varepsilon(U_2 - \bar{U}_2)); \mathcal{L}e_2) dt.
\end{aligned}$$

The Hölder inequality yields

$$\begin{aligned}
E_1 &\leq \|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)} \|e_2\|_{L^\infty(L^2)} + \left( \sum_{j=1}^{N} \|(\mathcal{G} - I)U_1^{j-1-}\|_{H^1(\Omega)} \right) \|e_1\|_{L^\infty(H^1)} \\
&\quad + \left( \sum_{j=1}^{N} \|(\mathcal{L} - I)U_2^{j-1-}\|_{L^2(\Omega)} \right) \|e_2\|_{L^\infty(L^2)} + \|U_2 - \bar{U}_2\|_{L^1(H^1)} \|e_1\|_{L^\infty(H^1)} \\
&\quad + \|f - \bar{f} + \Delta(U_1 - \bar{U}_1 + \varepsilon(U_2 - \bar{U}_2))\|_{L^1(L^2)} \|e_2\|_{L^\infty(L^2)}.
\end{aligned}$$

With the aid of the Young inequality we obtain for $\alpha = \{1, 2\}$

$$\begin{aligned}
E_1 &\leq 3\alpha \|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}^2 + 3\alpha \|f - \bar{f} + \Delta(U_1 - \bar{U}_1 + \varepsilon(U_2 - \bar{U}_2))\|_{L^1(L^2)}^2 \\
&\quad + 3\alpha \left( \sum_{j=1}^{N} \|(\mathcal{L} - I)U_2^{j-1-}\|_{L^2(\Omega)} \right)^2 + 2\alpha \left( \sum_{j=1}^{N} \|(\mathcal{G} - I)U_1^{j-1-}\|_{H^1(\Omega)} \right)^2 \\
&\quad + 2\alpha \|U_2 - \bar{U}_2\|_{L^1(H^1)}^2 + \frac{1}{4\alpha} \|e\|_{L^\infty(\mathcal{H})}^2.
\end{aligned}$$

This concludes the proof of the lemma. $\qquad \square$

**Lemma 3.3.2.5.** For all $n = 1, \ldots, N$ and $M_{2,n}$ defined by

$$M_{2,n} := 2\|U_2 - U_{1,\tau}\|^2_{L^1(I_n;H^1(\Omega))} + 2\|[U_1]^{n-1}\|^2_{H^1(\Omega)}$$
$$+ 2\|f - U_{2,\tau} + \Delta(U_1 + \varepsilon U_2))\|^2_{L^1(I_n;L^2(\Omega))} + 2\|[U_2]^{n-1}\|^2_{L^2(\Omega)},$$

there holds for both $\alpha = \{1, 2\}$

$$E_2 := \int_{I_n} \langle R_n ; e \rangle_{\mathcal{H}} dt \leq \alpha M_{2,n} + \frac{1}{4\alpha}\|e\|^2_{L^\infty(\mathcal{H})}. \tag{3.138}$$

Note that $E_2$ is the LHS of (3.126b) since the jump term vanishes for $\mathcal{C}^1$ functions.

**Proof.** In case of $E_2$ we have

$$E_2 = \int_{I_n} \langle F - U_\tau - \mathcal{A}U ; e \rangle_{\mathcal{H}} dt + \langle [U]^{n-1} ; e^{j-1+} \rangle_{\mathcal{H}}$$
$$= \int_{I_n} a(U_2 - U_{1,\tau} ; e_1) dt + \int_{I_n} (f - U_{2,\tau} + \Delta(U_1 + \varepsilon U_2); e_2) dt$$
$$+ a([U_1]^{n-1}; e_1^{n-1+}) + ([U_2]^{n-1}; e_2^{n-1+}).$$

Applying the Hölder and then the Young inequality inequality for $\alpha = \{1, 2\}$, we conclude

$$E_2 \leq 2\alpha\|U_2 - U_{1,\tau}\|^2_{L^1(I_n;H^1(\Omega))} + 2\alpha\|f - U_{2,\tau} + \Delta(U_1 + \varepsilon U_2))\|^2_{L^1(I_n;L^2(\Omega))}$$
$$+ 2\alpha\|[U_1]^{n-1}\|^2_{H^1(\Omega)} + 2\alpha\|[U_2]^{n-1}\|^2_{L^2(\Omega))} + \frac{1}{4\alpha}\|e\|^2_{L^\infty(\mathcal{H})}. \qquad \square$$

If we assume that the discrete initial solution $U^{0-}$ is defined as in Lemma 3.3.1.7, then the results of the latter lemmas and the estimate (3.120), combined with (3.127) and (3.128) yield the proof of theorem. $\qquad \square$

**Remark 3.3.2.3.** In Subsection 3.3.1.2, we used the bilinear form $\widetilde{\mathcal{B}}$ for the discrete model to derive an a posteriori error estimate for $dG(0) \otimes \mathcal{C}^1$. On the other hand, we may also perform the same strategy for the derivation of a posteriori error estimates for the $dG(0) \otimes \mathcal{C}^1$ approximation as for $dG(1) \otimes \mathcal{C}^1$. In this case, the resulting estimate reads for $\varepsilon \geq 0$

- if the temporal operator $\mathcal{J}$ is defined by $\mathcal{J}|_{I_j} e := \fint_{I_j} e$,

$$\|e\|^2_{L^\infty(\mathcal{H})} + \varepsilon\|e_2\|^2_{L^2(H^1)} dt \leq 3\|y_0 - \mathcal{I}y_0\|^2_{H^1(\Omega)} + 3\|y_1 - \mathcal{I}y_1\|^2_{L^2(\Omega)} + 44\|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|^2_{L^1(L^2)}$$
$$+ 22\Big(\sum_{j=1}^{N} \|(\mathcal{G} - I)U_1^{j-1-}\|_{H^1(\Omega)}\Big)^2 + 22\Big(\sum_{j=1}^{N} \|(\mathcal{L} - I)U_2^{j-1-}\|_{L^2(\Omega)}\Big)^2$$
$$+ 44\|f - \bar{f}\|^2_{L^1(L^2)} + 22\Big(\sum_{j=1}^{N} \|[U]_1^{j-1}\|_{L^1(H^1)}\Big)^2$$
$$+ 44\Big(\sum_{j=1}^{N} \|[U_2]^{j-1}\|_{L^1(L^2)}\Big)^2 + 20 \max_{1 \leq j \leq N} \Big\{ \|U_2\|^2_{L^1(I_j;H^1(\Omega))}$$
$$+ \|[U]^{j-1}\|^2_{\mathcal{H}} + \|f + \Delta(U_1 + \varepsilon U_2)\|^2_{L^1(I_j;L^2(\Omega))} \Big\}. \tag{3.139}$$

- if the temporal operator $\mathcal{J}$ is defined by $\mathcal{J}|_{I_j}e := e^{j-1+}$

$$
\begin{aligned}
\|e\|^2_{L^\infty(\mathcal{H})} + \varepsilon\|e_2\|^2_{L^2(H^1)} \leq{}& 3\|y_0 - \mathcal{I}y_0\|^2_{H^1(\Omega)} + 3\|y_1 - \mathcal{I}y_1\|^2_{L^2(\Omega)} + 22\|U_2\|^2_{L^1(H^1)} \\
&+ 33\|(I-\mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|^2_{L^1(L^2)} + 33\|f + \Delta(U_1 + \varepsilon U_2)\|^2_{L^1(L^2)} \\
&+ 22\Big(\sum_{j=1}^N \|(\mathcal{G}-I)U_1^{j-1-}\|_{H^1(\Omega)}\Big)^2 + 33\Big(\sum_{j=1}^N \|(\mathcal{L}-I)U_2^{j-1-}\|_{L^2(\Omega)}\Big)^2 \\
&+ 20 \max_{1\leq j \leq N} \Big\{ \|U_2\|^2_{L^1(I_j;H^1(\Omega))} + \|[U]^{j-1}\|^2_{\mathcal{H}} \\
&\qquad\qquad + \|f + \Delta(U_1 + \varepsilon U_2)\|^2_{L^1(I_j;L^2(\Omega))} \Big\}.
\end{aligned}
\tag{3.140}
$$

The estimate in case $\varepsilon = 0$ can be derived similarly. The comparison of estimate (3.112) with (3.139) and (3.140) shows that the estimates (3.139) and (3.140) are of order $\mathcal{O}(1)$ in time, whereby the estimate (3.112) shows the convergence order $\mathcal{O}(k)$ which makes it better. □

### 3.3.3  $cG(1)$ time approximation

Within this section we provide an a posteriori error analysis for $cG(1)\otimes \mathcal{P}_1$ and $cG(1)\otimes \mathcal{C}^1$ discretisation. The presented analysis follows via residual arguments introduced in Definition 2.3.1.2 with $\mathcal{B}, \mathscr{L}$ from (2.41) and (2.42), respectively. For the notation and definitions, we refer to Section 2.1 and Subsection 2.3.4. Note that the analysis of this section is completely analogous to the analysis for the $dG(1)$ case. The only difference is that the jump terms in time of the discrete solution do not occur here.

Recall that in case of the $cG(1)$ time discretisation, we require that $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j = 1, \ldots, N$ in order to provide the continuity of the discrete solution with respect to the neighbouring time levels.

**Lemma 3.3.3.1 (Error representation, $cG(1)$ time approximation).** In case of $cG(1)$ time discretisation there holds for every $V \in \mathcal{W}_c$,

$$
\frac{1}{2}\|e(T)\|^2_{\mathcal{H}} + \varepsilon\int_0^T \|e_2(t)\|^2_{H^1(\Omega)}\,dt = \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} + Res(e - V).
\tag{3.141}
$$

**Proof.** The fundamental theorem of calculus and the fact that the error $e$ is a continuous function in time, leads to

$$
\frac{1}{2}\|e(T)\|^2_{\mathcal{H}} - \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} = \int_0^T \langle \dot{e}\,;e\rangle_{\mathcal{H}}\,dt.
\tag{3.142}
$$

With the residual representation (2.7) and the orthogonality of the residual from (2.8), we have for arbitrary $V \in W_c$

$$
\begin{aligned}
\int_0^T \langle \dot{e}\,;e\rangle_{\mathcal{H}}\,dt &= \int_0^T \langle \dot{e}\,;e - V\rangle_{\mathcal{H}}\,dt + \int_0^T a(e_2;V_1)\,dt - \int_0^T a(e_1 + \varepsilon e_2;V_2)\,dt \\
&= Res(e - V) - \varepsilon\int_0^T a(e_2;e_2)\,dt.
\end{aligned}
$$

The combination of the last equality and the identity (3.142) yields the proof of lemma. □

**Lemma 3.3.3.2.** In case of $cG(1)$ time approximation and $\mathcal{P}_1$ or $\mathcal{C}^1$ approximation in space, the residual (2.7) may be expressed in terms of local residuals and local jumps such that for any $v \in L^2(H_D^1 \times H_D^1)$

$$Res(v) = \int_0^T \langle R\,;v\rangle_{\mathcal{H}}dt + \int_0^T J(U,v)dt \qquad (3.143a)$$

where for all $(t,x) \in [0,T] \times \Omega$

$$R(t,x):= F(t,x) - \dot{U}(t,x) - \mathcal{A}U(t,x) \ \ \text{in} \ \ Q, \qquad (3.143b)$$

$$J(U,v)(t):= \begin{cases} \sum_{k=1}^m ([D(U_1+\varepsilon U_2)]_k)(t)v_2(x_{t,k}), & \text{for } \ \mathcal{P}_1 \ \text{ elements,} \\ 0, & \text{for } \ \mathcal{C}^1 \ \text{ elements.} \end{cases} \qquad (3.143c)$$

Here $m = n, n{-}1$ for Problem (DN), (DD) respectively. Furthermore, $[D(U_1+\varepsilon U_2)]_k$ are defined as jumps of the piecewise constant functions $DU_1, DU_2$ with respect to the space coordinate, see Definition 3.1.0.6.

**Proof.** Given (2.7) with $\mathcal{B}, \mathscr{L}$ as in (2.41), (2.42), respectively, an integration by parts in space yields for each $v \in L^2(H_D^1 \times H_D^1)$

$$Res(v) = \int_0^T (f; v_2)dt - \int_0^T a(\dot{U}_1; v_1)dt - \int_0^T (\dot{U}_2; v_2)dt + \int_0^T a(U_2; v_1)dt + \int_0^T (\Delta U_1; v_2)dt$$

$$+ \varepsilon \int_0^T (\Delta U_2; v_2)dt + \int_0^T \sum_{k=1}^m [D(U_1+\varepsilon U_2)]_k v_2(x_k)dt.$$

In case of the spatial approximation by Hermite cubic splines, the jump terms equal zero due to the fact that the discrete functions satisfy the boundary conditions, see Subsection 2.1.2. With the notation from Definition 1.3.0.7, we conclude the proof of the lemma. $\qquad\square$

**Remark 3.3.3.1.** If $\mathscr{T}$ is an arbitrary discretisation of time interval $[0,T]$, then from the definition of the local weak problem (2.6) and the residual representation from Lemma 3.3.3.2, there holds for each $t_n \in \mathscr{T}$

$$\frac{1}{2}\|e(t_n)\|_{\mathcal{H}}^2 + \varepsilon \int_0^{t_n}\|e_2(t)\|_{H^1(\Omega)}^2 dt = \frac{1}{2}\|e(0)\|_{\mathcal{H}}^2 + \int_0^{t_n}\langle R\,;e-V\rangle_{\mathcal{H}}dt + \int_0^{t_n} J(U\,;e-V)dt. \quad (3.144)$$

Here $V$ is an arbitrary test function, $V|_{I_j} \in \mathcal{Q}_0^j$ for all $j = 1, \ldots, n$ for which the residual orthogonality (2.8) holds. $\qquad\square$

**Lemma 3.3.3.3.** Let $R$ and $J$ be defined by (3.143b) and (3.143c), respectively. Let $e$ be the error of $cG(1)$ solution with respect to the $\mathcal{P}_1$ i.e. $\mathcal{C}^1$ discretisation in space. Let $\Pi := (\mathcal{G}, \mathcal{L})$ be a spatial orthogonal projection with respect to the $\mathcal{H}$ scalar product. We assume that there are some positive $M_1, M_{2,n}$ such that there holds for all $n = 1, \ldots, N$ and both $\alpha = \{1,2\}$

$$\int_0^{t_n}\langle R\,;e-\Pi\bar{e}\rangle_{\mathcal{H}}dt + \int_0^{t_n} J(U\,;e-\Pi\bar{e})dt \le \alpha M_1 + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2, \qquad (3.145a)$$

$$\int_{I_n}\langle R\,;e\rangle_{\mathcal{H}}dt + \int_{I_n} J(U\,;e)dt \le \alpha M_{2,n} + \frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2, \qquad (3.145b)$$

then the following estimate is valid for all $\varepsilon \geq 0$

$$\|e\|^2_{L^\infty(\mathcal{H})} + \varepsilon \|e_2\|^2_{L^2(H^1)} \leq 3\|e(0)\|^2_{\mathcal{H}} + 11 M_1 + 10 \max_{1 \leq j \leq N} M_{2,j}. \qquad (3.146)$$

Moreover if $\varepsilon = 0$, then there holds additionally

$$\|e\|^2_{L^\infty(\mathcal{H})} \leq 2\|e(0)\|^2_{\mathcal{H}} + 8 M_1 + 8 \max_{1 \leq j \leq N} M_{2,j}. \qquad (3.147)$$

**Proof.** The proof is similar to the proof of Lemma 3.3.2.3. Here we also discuss only the case $t \notin \mathcal{T}$ where $\|e(t)\|_{\mathcal{H}} = \|e\|_{L^\infty(\mathcal{H})}$ and $\mathcal{T}$ is the arbitrary triangulation of the time interval $[0, T]$.

Let us assume that $t \in [t_{\ell-1}, t_\ell)$ for some $1 \leq \ell \leq N$.
Given the error representation formula (3.144) with $V = \Pi \bar{e}$, from (3.145) we have for all $n = 1, \ldots, N$ and $\alpha \in \{1, 2\}$

$$\frac{1}{2}\|e(t_n)\|^2_{\mathcal{H}} + \varepsilon \int_0^{t_n} \|e_2(t)\|^2_{H^1(\Omega)} dt = \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} + \int_0^{t_n} \langle R\,;\, e - \Pi\bar{e} \rangle_{\mathcal{H}} dt + \int_0^{t_n} J(U\,;\, e - \Pi\bar{e}) dt$$

$$\leq \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} + \alpha M_1 + \frac{1}{4\alpha}\|e\|^2_{L^\infty(\mathcal{H})}. \qquad (3.148)$$

The fundamental theorem of calculus, Lemma 3.3.3.2 and assumption (3.145) with $\alpha = 2$, imply

$$\frac{1}{2}\|e\|^2_{L^\infty(\mathcal{H})} - \frac{1}{2}\|e(t_{\ell-1})\|^2_{\mathcal{H}} = \int_{t_{\ell-1}}^t \langle e_\tau\,;\, e \rangle_{\mathcal{H}} d\tau$$

$$= \int_{t_{n-1}}^t \langle R\,;\, e \rangle_{\mathcal{H}} d\tau + \int_{t_{n-1}}^t J(U\,;\, e) d\tau - \varepsilon \int_{t_{n-1}}^t \|e_2(\tau)\|^2_{H^1(\Omega)} d\tau$$

$$\leq 2 M_{2,\ell} + \frac{1}{8}\|e\|^2_{L^\infty(\mathcal{H})} \leq 2 \max_{1 \leq j \leq N} M_{2,j} + \frac{1}{8}\|e\|^2_{L^\infty(\mathcal{H})}. \qquad (3.149)$$

With (3.148) where $n = \ell - 1$ and $\alpha = 2$, we obtain from (3.149) and the fact that $\varepsilon \geq 0$

$$\|e\|^2_{L^\infty(\mathcal{H})} \leq 2\|e(0)\|^2_{\mathcal{H}} + 8 M_1 + 8 \max_{1 \leq j \leq N} M_{2,j}. \qquad (3.150)$$

For $\varepsilon = 0$ this concludes the proof. For all $\varepsilon \geq 0$, from (3.148) with $n = N$ and $\alpha = 1$ we have

$$\varepsilon \int_0^T \|e_2(t)\|^2_{H^1(\Omega)} dt \leq \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} + M_1 + \frac{1}{4}\|e\|^2_{L^\infty(\mathcal{H})} \leq \|e(0)\|^2_{\mathcal{H}} + 3 M_1 + 2 \max_{1 \leq j \leq N} M_{2,j}. \qquad (3.151)$$

Note that in the second inequality above we used the estimate (3.150).
Finally, for $\varepsilon \geq 0$, the combination of (3.150) and (3.151) proves

$$\|e\|^2_{L^\infty(\mathcal{H})} + \varepsilon \|e_2\|^2_{L^2(H^1)} \leq 3\|e(0)\|^2_{\mathcal{H}} + 11 M_1 + 10 \max_{1 \leq j \leq N} M_{2,j}. \qquad (3.152)$$

$$\square$$

### 3.3.3.1 A posteriori energy error analysis, $cG(1)\otimes\mathcal{P}_1$

From Lemma 3.3.3.3, an a posteriori error estimate can be derived for the case of $cG(1)\otimes\mathcal{P}_1$ approximation if the following terms are estimated with respect to $\|e\|_{L^\infty(\mathcal{H})}$, namely

$$E_1 := \int_0^{t_n} \langle R \,; e - \Pi\bar{e}\rangle_{\mathcal{H}} dt + \int_0^{t_n} J(U, e - \Pi\bar{e}) dt, \tag{3.153a}$$

$$E_2 := \int_{I_n} \langle R \,; e\rangle_{\mathcal{H}} dt + \int_{I_n} J(U, e) dt. \tag{3.153b}$$

Here we deal with the same complicity in the derivation of the a posteriori error estimate as in the case of $dG(0)\otimes\mathcal{P}_1$ and $dG(1)\otimes\mathcal{P}_1$ approximation, see Subsection 3.3.1.1 and 3.3.2.1, respectively. Namely, we can not accomplish the estimation of $(e_2 - \mathcal{L}e_2)(x_k)$ with respect to $\|e_2\|_{L^\infty(L^2)}$. However, the term in question can be estimated with respect to $\|e_2\|_{H^1(\Omega)}$, but then we can not apply Lemma 3.3.3.3. For further details, we refer to the Subsection 3.3.1.1.

### 3.3.3.2 A posteriori energy error analysis, $cG(1)\otimes\mathcal{C}^1$

According to results and conclusions from the previous section a spatially $\mathcal{C}^1$ approximation is required in order to prove a sharp residual-based a posteriori error bound in the energy norm.

**Theorem 3.3.3.1 (A posteriori energy error estimate, $cG(1)\otimes\mathcal{C}^1$).** For $u$ and its discrete $cG(1)\otimes\mathcal{C}^1$ counterpart $U$, there holds for all $\varepsilon \geq 0$

$$\begin{aligned}
\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon\|e_2\|_{L^2(H^1)}^2 \leq{}& 3\|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)}^2 + 3\|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)}^2 + 11\|U_2 - \overline{U}_2\|_{L^1(H^1)}^2 \\
&+ 22\|f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2))\|_{L^1(L^2)}^2 \\
&+ 22\|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}^2 + 10 \max_{1 \leq j \leq N} \Big\{ \|U_2 - \dot{U}_1\|_{L^1(I_j;H^1)}^2 \\
&+ \|f - \dot{U}_2 + \Delta(U_1 + \varepsilon U_2)\|_{L^1(I_j;\Omega)}^2 \Big\}.
\end{aligned} \tag{3.154}$$

Moreover if $\varepsilon = 0$, there holds additionally

$$\begin{aligned}
\|e\|_{L^\infty(\mathcal{H})}^2 \leq{}& 2\|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)}^2 + 2\|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)}^2 + 8\|U_2 - \overline{U}_2\|_{L^1(H^1)}^2 \\
&+ 16\|f - \bar{f} + \Delta(U_1 - \overline{U}_1)\|_{L^1(L^2)}^2 + 16\|(I - \mathcal{L})(f + \Delta U_1)\|_{L^1(L^2)}^2 \\
&+ 8 \max_{1 \leq j \leq N} \Big\{ \|U_2 - \dot{U}_1\|_{L^1(I_j;H^1)}^2 + \|f - \dot{U}_2 + \Delta U_1\|_{L^1(I_j;\Omega)}^2 \Big\}.
\end{aligned} \tag{3.155}$$

Note that the RHS of both estimates is denoted by $\eta^2$.

**Remark 3.3.3.2.** For sufficiently smooth initial data, the error estimate of Theorem 3.3.3.1 is of order $\mathcal{O}(h^2 + k)$. $\qquad\square$

**Proof.** The proof follows from an application of Lemma 3.3.3.3. We prove the validity of the assumptions (3.145) in the following two lemmas.

**Lemma 3.3.3.4.** With

$$\begin{aligned}
M_1 :={}& 2\|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}^2 + \|U_2 - \overline{U}_2\|_{L^1(H^1)}^2 \\
&+ 2\|f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2))\|_{L^1(L^2)}^2,
\end{aligned}$$

there holds for $n=1,\ldots,N$ and both $\alpha\in\{1,2\}$

$$E_1:=\int_0^{t_n}\langle R\,;e-\Pi\bar{e}\rangle_{\mathcal{H}}dt\leq\alpha M_1+\frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2.\qquad(3.156)$$

Note that $E_1$ is the LHS of (3.145a), since the jump terms vanishes for $\mathcal{C}^1$ functions.

**Proof.** Due to the orthogonal properties of projections $\mathcal{G},\mathcal{L}$ and the fact that $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$ for all $j=1,\ldots,N$, we have

$$E_1=\int_0^{t_n}\langle R-\Pi R\,;e-\Pi e\rangle_{\mathcal{H}}dt+\int_0^{t_n}\langle R-\bar{R}\,;\Pi e-\Pi\bar{e}\rangle_{\mathcal{H}}dt$$
$$=\int_0^{t_n}\langle F-\Pi F-\mathcal{A}U+\Pi\mathcal{A}U\,;e\rangle_{\mathcal{H}}dt+\int_0^{t_n}\langle F-\overline{F}-\mathcal{A}(U-\overline{U})\,;\Pi e\rangle_{\mathcal{H}}dt.$$

Expanding the last equation in terms of $\mathcal{H}$ scalar product we show

$$E_1=\int_0^{t_n}((I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2));e_2)dt+\int_0^{t_n}a(U_2-\overline{U}_2;\mathcal{G}e_1)dt$$
$$+\int_0^{t_n}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2));\mathcal{L}e_2)dt.$$

The main properties of projections $\mathcal{G},\mathcal{L}$ and the Cauchy inequality yield

$$E_1\leq\|(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(L^2)}\|e_2\|_{L^\infty(L^2)}+\|U_2-\overline{U}_2\|_{L^1(H^1)}\|e_1\|_{L^\infty(H^1)}$$
$$+\|f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2))\|_{L^1(L^2)}\|e_2\|_{L^\infty(L^2)}.$$

With the aid of the Young inequality, we conclude for $\alpha\in\{1,2\}$

$$E_1\leq 2\alpha\|(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(I_j;L^2)}^2+\alpha\|U_2-\overline{U}_2\|_{L^1(H^1)}^2$$
$$+2\alpha\|f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2))\|_{L^1(L^2)}^2+\frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2.\qquad\square$$

**Lemma 3.3.3.5.** For all $n=1,\ldots,N$ and $M_{2,n}$ defined by

$$M_{2,n}:=\|U_2-\dot{U}_1\|_{L^1(I_n;H^1(\Omega))}^2+\|f-\dot{U}_2+\Delta(U_1+\varepsilon U_2)\|_{L^1(I_n;L^2(\Omega))}^2.$$

there holds for both $\alpha\in\{1,2\}$

$$E_2:=\int_{I_n}\langle R\,;e\rangle_{\mathcal{H}}dt\leq\alpha\max_{1\leq j\leq N}M_{2,j}+\frac{1}{4\alpha}\|e\|_{L^\infty(\mathcal{H})}^2.$$

Note that $E_2$ is the LHS of (3.145b), since the jump terms vanishes for $\mathcal{C}^1$ functions.

**Proof.** Similar as in the proof of of lemma above, we have

$$E_2=\int_0^{t_n}\langle F-\dot{U}-\mathcal{A}U+\Pi\mathcal{A}U\,;e\rangle_{\mathcal{H}}dt$$
$$=\int_0^{t_n}a(U_2-\dot{U}_1;e_1)dt+\int_0^{t_n}(f-\dot{U}_2+\Delta(U_1+\varepsilon U_2);e_2)dt$$
$$\leq\|U_2-\dot{U}_1\|_{L^1(I_n;H^1(\Omega))}\|e_1\|_{L^\infty(H^1)}+\|f-\dot{U}_2+\Delta(U_1+\varepsilon U_2)\|_{L^1(I_n;L^2(\Omega))}\|e_2\|_{L^\infty(L^2)}.$$

An application of the Young inequality for both $\alpha\in\{1,2\}$ yields the proof of lemma. $\quad\square$

If we assume that the discrete initial solution $U(0)$ is defined as $U^{0-}$ in Lemma 3.3.1.7, then the latter Lemmas and the estimate (3.120) combined with (3.146) and (3.147) yield the proof of theorem. $\quad\square$

### 3.3.4  Method of lines

In this section we provide a full error analysis concerning the derivation of a posteriori error bound in case of the method of lines and $\mathcal{P}_1$ or $\mathcal{C}^1$ spatial discretisation. Namely, we seek to find an error bound $\eta$, such that

$$\|e\|^2_{L^\infty(\mathcal{H})} + \varepsilon \|e_2\|^2_{L^2(H^1)} \le \eta.$$

The error analysis follows via residual argument introduced in Definition 2.3.1.2, with $\mathcal{B}, \mathcal{L}$ from (2.49) and (2.50), respectively. For the notations and definitions, we refer to Section 2.1 and Subsection 2.3.5.

**Lemma 3.3.4.1 (Error representation, *MoL*).** In case of the *MoL* discretisation, there holds for each $t \in [0, T]$

$$\frac{1}{2}\|e(t)\|^2_{\mathcal{H}} + \varepsilon \int_0^t \|e_2(\tau)\|^2_{H^1(\Omega)} d\tau = \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} + \int_0^t Res(e-V)d\tau \quad \text{for all} \quad V \in \mathcal{W}_s. \qquad (3.157)$$

**Proof.** From the fundamental theorem of calculus, we may deduce for each $t \in [0, T]$

$$\frac{1}{2}\|e(t)\|^2_{\mathcal{H}} - \frac{1}{2}\|e(0)\|^2_{\mathcal{H}} = \int_0^t \langle \dot{e} \, ; e \rangle_{\mathcal{H}} d\tau.$$

According to the definition of the residual $Res$, cf. (2.7), and the orthogonality (2.8), we have for arbitrary $V \in \mathcal{W}_s$

$$\int_0^t \langle \dot{e} \, ; e \rangle_{\mathcal{H}} d\tau = \int_0^t \langle \dot{e} \, ; e-V \rangle_{\mathcal{H}} d\tau + \int_0^t a(e_2; V_1)d\tau - \int_0^t a(e_1 + \varepsilon e_2; V_2)d\tau$$

$$= Res(e-V) - \varepsilon \int_0^t a(e_2; e_2)d\tau.$$

A combination of the latter two inequalities yields the proof of the lemma.                    $\square$

**Lemma 3.3.4.2.** In case of the method of lines and $\mathcal{P}_1$ or $\mathcal{C}^1$ approximation in space, the residual (2.7) may be expressed in terms of local residuals and local jumps such that for any $v \in C(0, T; H^1_D(\Omega)^2)$ and $t \in [0, T]$,

$$Res(v)(t) = \langle R \, ; v \rangle_{\mathcal{H}}(t) + J(U, v)(t), \qquad (3.158a)$$

where in particular

$$R(t, x) := F(t, x) - \dot{U}(t, x) - \mathcal{A}U(t, x) \quad \text{in} \quad Q, \qquad (3.158b)$$

$$J(U, v)(t) := \begin{cases} \sum_{k=1}^m [D(U_1 + \varepsilon U_2)(t)]_k v_2(t, x_k), & \text{for} \ \ \mathcal{P}_1 \ \ \text{elements}, \\ 0, & \text{for} \ \ \mathcal{C}^1 \ \ \text{elements}. \end{cases} \qquad (3.158c)$$

Here $m = n, n-1$ for Problem (DN), (DD) respectively. The jumps terms $[DU_1]_k$, $[DU_2]_k$ are defined as in Definition 3.1.0.6.

**Proof.** Given (2.7) with $\mathcal{B}, \mathcal{L}$ as in (2.41) and (2.42), respectively, an integration by parts in space yields for each $v \in C(0, T; H_D^1(\Omega)^2)$ and $t \in [0, T]$,

$$Res(v)(t) = (f(t); v_2(t)) - a(\dot{U}_1(t); v_1(t)) - (\dot{U}_2(t); v_2(t)) + a(U_2(t); v_1(t)) + (\Delta U_1(t); v_2(t))$$

$$+ \varepsilon(\Delta U_2(t); v_2(t)) + \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)(t)]_k v_2(t, x_k).$$

In case of the spatial approximation by $\mathcal{C}^1$ elements, the jump terms equal zero due to the fact that the discrete functions satisfy the (DD) or (DN) boundary conditions, see Subsection 2.1.2. With the notation from Definition 1.3.0.7, we conclude the proof of the lemma. $\quad\square$

**Remark 3.3.4.1.** From the error representation (3.157) and the residual representation from Lemma 3.3.4.2, there holds for each $t \in [0, T]$ and $V \in \mathcal{W}_s$

$$\frac{1}{2}\|e(t)\|_{\mathcal{H}}^2 + \varepsilon \int_0^t \|e_2(\tau)\|_{H^1(\Omega)}^2 d\tau = \frac{1}{2}\|e(0)\|_{\mathcal{H}}^2 + \int_0^t \langle R; e - V \rangle_{\mathcal{H}} d\tau + \int_0^t J(U, e - V) d\tau. \quad (3.159)$$
$$\square$$

**Lemma 3.3.4.3.** Let $R$ and $J$ be defined by (3.158b) and (3.158c), respectively. Let $e$ be the the error of semi-discrete solution with respect to the $\mathcal{P}_1$ i.e. $\mathcal{C}^1$ discretisation in space. Let $\Pi := (\mathcal{G}, \mathcal{L})$ be a spatial orthogonal projection with respect to the $\mathcal{H}$ scalar product. We assume that there is some positive $M$ such that there holds for each $t \in [0, T]$

$$\int_0^t \langle R; e - \Pi e \rangle_{\mathcal{H}} d\tau + \int_0^t J(U; e - \Pi e) d\tau \leq M + \frac{1}{4}\|e\|_{L^\infty(\mathcal{H})}^2, \quad (3.160)$$

then the following estimate is valid for all $\varepsilon \geq 0$

$$\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon\|e_2\|_{L^2(H^1)}^2 \leq 3\|e(0)\|_{\mathcal{H}}^2 + 6M. \quad (3.161)$$

Moreover, if $\varepsilon = 0$, then there holds a sharper estimate

$$\|e\|_{L^\infty(\mathcal{H})}^2 \leq 2\|e(0)\|_{\mathcal{H}}^2 + 4M. \quad (3.162)$$

**Proof.** Given the error representation formula (3.159) where $V = \Pi e$, from (3.160) there holds for all $t \in [0, T]$

$$\frac{1}{2}\|e(t)\|_{\mathcal{H}}^2 + \varepsilon \int_0^t \|e_2(\tau)\|_{H^1(\Omega)}^2 d\tau = \frac{1}{2}\|e(0)\|_{\mathcal{H}}^2 + \int_0^t \langle R; e - \Pi e \rangle_{\mathcal{H}} d\tau + \int_0^t J(U; e - \Pi e) d\tau$$

$$\leq \frac{1}{2}\|e(0)\|_{\mathcal{H}}^2 + M + \frac{1}{4}\|e\|_{L^\infty(\mathcal{H})}^2. \quad (3.163)$$

The function $e$ is a monotone function in time. Let $t_{max} \in [0, T]$ be some point in time such that $\|e\|_{L^\infty(\mathcal{H})} = \|e(t_{max})\|_{\mathcal{H}}$.

If we choose $t = t_{max}$ and $\alpha = 1$ in (3.163), then for $\varepsilon = 0$ we may conclude

$$\|e\|_{L^\infty(\mathcal{H})}^2 \leq 2\|e^{0-}\|_{\mathcal{H}}^2 + 4M. \quad (3.164)$$

From (3.163) with $t = T$ and $\alpha = 1$, we have for all $\varepsilon \geq 0$

$$\varepsilon \int_0^T \|e_2(\tau)\|_{H^1(\Omega)}^2 d\tau \leq \frac{1}{2}\|e^{0-}\|_{\mathcal{H}}^2 + M + \frac{1}{4}\|e\|_{L^\infty(\mathcal{H})}^2 \leq \|e^{0-}\|_{\mathcal{H}}^2 + 2M. \quad (3.165)$$

Note that in the second inequality above we used the estimate (3.164).

Finally, for $\varepsilon \geq 0$, the combination of (3.164) and (3.165) proves

$$\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon\|e_2\|_{L^2(H^1)}^2 \leq 3\|e(0)\|_{\mathcal{H}}^2 + 6M. \quad (3.166)$$
$$\square$$

### 3.3.4.1 A posteriori energy error analysis, $MoL \otimes \mathcal{P}_1$

From Lemma 3.3.4.3, we can derive an a posteriori error estimate if the assumption (3.160) hold. Unfortunately, in case of $MoL \otimes \mathcal{P}_1$ discrete problem, we did not succeed to derive an upper bound for

$$E_1 := \int_0^t \langle R\,;e - \Pi e\rangle_{\mathcal{H}} d\tau + \int_0^t J(U, e - \Pi e)d\tau$$

in terms of $\|e\|_{L^\infty(\mathcal{H})}$. Namely, we can estimate the second term on the RHS of the equality above only with respect to $\|e_2\|_{H^1(\Omega)}$ and thereafter the adequate a posteriori error bound can not be derived. We refer to Subsection 3.3.1.1 for further explanation.

### 3.3.4.2 A posteriori energy error analysis, $MoL \otimes \mathcal{C}^1$

From the Subsection 3.3.4.1, a spatially $\mathcal{C}^1$ approximation is required in order to prove a sharp residual-based a posteriori error bound in the energy norm.

**Theorem 3.3.4.1 (A posteriori energy error estimate, $MoL \otimes \mathcal{C}^1$).** For $u$ and its semi-discrete $MoL \otimes \mathcal{C}^1$ counterpart $U$, there holds for all $\varepsilon \geq 0$

$$\|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon \|e_2\|_{L^2(H^1)}^2 dt \leq 3\|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)}^2 + 3\|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)}^2$$
$$+ 6\|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}^2. \tag{3.167}$$

Moreover if $\varepsilon = 0$, then there holds additionally

$$\|e\|_{L^\infty(\mathcal{H})}^2 \leq 2\|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)}^2 + 2\|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)}^2 + 4\|(I - \mathcal{L})(f + \Delta U_1)\|_{L^1(L^2)}^2. \tag{3.168}$$

Note that the RHS of both estimates is denoted by $\eta^2$.

**Remark 3.3.4.2.** For sufficiently smooth initial data $f, y_0, y_1$, the estimates (3.167) and (3.168) are of order $\mathcal{O}(h^3)$. □

**Proof.** The proof follows from an application of Lemma 3.3.4.3. Therefore, we need to to prove the validity of the assumption (3.160). This is given in the following lemma.

**Lemma 3.3.4.4.** With

$$M := \|(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}^2,$$

there holds for $t \in [0, T]$

$$E_1 := \int_0^t \langle R\,;e - \Pi e\rangle_{\mathcal{H}} d\tau \leq M + \frac{1}{4}\|e\|_{L^\infty(\mathcal{H})}^2. \tag{3.169}$$

**Proof.** Owing to the orthogonal properties of projection $\Pi = (\mathcal{G}, \mathcal{L})$ with respect to $\mathcal{H}$ scalar product, the LHS of (3.169) is equivalent to

$$E_1 = \int_0^t \langle R - \Pi R\,;e\rangle_{\mathcal{H}} d\tau$$
$$= \int_0^t \langle F - \Pi F - \dot{U} + \Pi\dot{U} - \mathcal{A}U + \Pi\mathcal{A}U\,;e\rangle_{\mathcal{H}} d\tau$$
$$= \int_0^t (f - \mathcal{L}f + \Delta(U_1 + \varepsilon U_2) + \mathcal{L}\Delta(U_1 + \varepsilon U_2); e_2)d\tau.$$

Note that in the second inequality above we used $\dot{U} = \Pi\dot{U}$. An application of the Hölder and Young inequality yield the proof of lemma. □

In order to complete the proof, its left to estimate the error of initial solution $\|e(0)\|_{\mathcal{H}}$. If we assume that the discrete initial solution $U(0)$ is defined as $U^{0-}$ in Lemma 3.3.1.7, then the latter Lemma and the estimate (3.120) combined with (3.167) and (3.168) yield the proof of theorem. $\qquad\square$

# Chapter 4

# Dual method

Within the following chapter, we give an introduction to the dual method as one of the easily applicable methods for the derivation and analysis of a priori and a posteriori error bounds for finite elements discrete methods.

In determining the error bounds, the dual method approach assumes the usage of the so-called *backward problem* which is also referred to as *dual* or *adjoint problem*. The initial problem is then according to this notation often referred to as *forward problem*. The backward problem is introduced through its continuous (strong) and discrete (weak) formulation. For the strong formulation we refer to Section 4.1. The weak formulation is introduced for each ansatz in time separately. However, the derivation of the a priori resp. a posteriori error bound mainly relies on the so-called weak resp. strong stability estimates which are derived from the corresponding adjoint problem. Moreover, the a posteriori error analysis also relies on the residual arguments similar as in Chapter 3.

A priori and a posteriori error analysis are discussed in Section 4.2 and 4.3, respectively. We first introduce the time approximation, i.e. discontinuous and continuous Galerkin methods as well as the method of lines and then combine with different spatial ansatzes, i.e. $\mathcal{P}_1$ and $\mathcal{C}^1$ elements in space.

The notation used in the following is the one from Table 3.1 and also the additional notation presented in Table 4.1. The analytical techniques in case of $dG(q)$, $q=0,1$ time discretisation

| | | |
|---|---|---|
| $\Phi$ | exact dual solution | solution of (4.3) |
| $\Psi$ | discrete dual solution | solution of (4.9), (4.35) |
| $\mathcal{B}^*$ | weak dual bilinear form | Definition 4.1.0.2 |
| $\mathcal{K}$ | continuous operator $\mathcal{K}$ | Definition 1.3.0.5 |
| $\mathcal{K}_h$ | discrete $\mathcal{K}$ operator | Definition 2.3.3.4 |
| $\mathcal{J}_1$ | temporal $H^1$ projection | Definition 3.1.0.10 |
| $\langle \cdot\,;\cdot \rangle_{\widehat{\mathcal{H}}}, \| \cdot \|_{\widehat{\mathcal{H}}}$ | weak scalar product and corresponding norm | Definition 1.3.0.6 |

Table 4.1: Additional notation used in Chapter 4.

rely on the duality approach used in JOHNSON[42] for the $dG(1)\otimes\mathcal{P}^1$ discretisation of the wave equation. In case of continuous Galerkin approximation in time and space our analysis is similar to the one presented in AZIZ-MONK [5] and LUSKIN-RANNACHER[?], where the parabolic equation has been discretised. Both papers consider the wave equation , i.e strongly damped wave equation with $\varepsilon = 0$. In our analysis we extend these methods to the more general case $\varepsilon \geq 0$ with certain restrictions for $\varepsilon=0$ and $\varepsilon>0$, respectively.

## 4.1    Adjoint problem

In the following we introduce the adjoint problem related to the initial problem (1.13) and derive strong stability estimates for which. This will allow us to construct the discrete scheme for particular discrete ansatzes in space and time and also to develop the a priori and a posteriori error bounds.

First, we introduce the adjoint problem, whereby we recall the notations and definitions from Section 1.3. The adjoint problem reads: For given data $(\phi_1^{N-}, \phi_2^{N-}) \in H_D^1 \times L^2(\Omega)$, for $s = 0$, resp. $(\phi_1^{N-}, \phi_2^{N-}) \in (H^2 \cap H_D^1) \times H_D^1(\Omega)$, for $s = 1$, find $\phi : Q \to \mathbb{R}$ such that

$$\ddot{\phi} - \Delta\phi + \varepsilon\Delta\dot{\phi} = 0 \text{ on } Q \tag{4.1a}$$

subject to the initial conditions

$$\phi(T, x) = \phi_1^{N-}(x), \ \dot{\phi}(T, x) = \phi_2^{N-}(x) \text{ on } \Omega \tag{4.1b}$$

and boundary conditions of either

$$\phi(t, 0) = 0, \ D\phi(t, 1) - \varepsilon D\dot{\phi}(t, 1) = 0, \quad \text{on} \quad [0, T] \qquad (\text{DN}^*), \tag{4.1c}$$

or

$$\phi(t, 0) = 0, \ \dot{\phi}(t, 0) = 0, \quad \text{on} \quad [0, T] \qquad (\text{DD}^*). \tag{4.1d}$$

We refer to the dual continuous problem as Problem (DN*) when (4.1a)–(4.1b)–(4.1c) hold and Problem (DD*) provided (4.1a)–(4.1b)–(4.1d).
Similar to the notation from Definition 1.3.0.7 we may also rewrite the dual problem (4.1) in the vector form.

**Definition 4.1.0.1 (Adjoint operator $\mathcal{A}^*$).** Let $\Phi = (\phi, \dot{\phi})$ the vector form of the dual solution of problem (4.1) and define the dual operator matrix $\mathcal{A}^* : \mathcal{H} \to \mathcal{H}$ by

$$\mathcal{A}^* := \begin{bmatrix} 0 & 1 \\ \Delta & -\varepsilon\Delta \end{bmatrix}. \tag{4.2}$$

The problem (4.1) is equivalent to: Find $\Phi \in H^1(0, T; \mathcal{H})$ such that

$$-\dot{\Phi}(t, x) + \mathcal{A}^*\Phi(t, x) = 0, \quad \text{on} \quad Q, \tag{4.3a}$$
$$\Phi(T, x) = \Phi^{N-}(x) \quad \text{on} \quad \Omega, \tag{4.3b}$$

for $\mathcal{A}^* : \mathscr{D}(A^*) \subset \mathcal{H} \to \mathcal{H}$ where

$$\mathscr{D}(\mathcal{A}^*) := \{(u_1, u_2) \subset H_D^1(\Omega) \times H_D^1(\Omega) \mid u_1 - \varepsilon u_2 \in H^2(\Omega) \text{ and } D(u_1 - \varepsilon u_2)|_{\Gamma_N} = 0\}$$

and given $\Phi^{N-} \in \mathscr{D}(\mathcal{A}^*)$.                                                                              $\square$

**Remark 4.1.0.3.** Note that for $u \in \mathscr{D}(\mathcal{A})$ and $v \in \mathscr{D}(\mathcal{A}^*)$ there holds $\langle \mathcal{A}u ; v \rangle_{\mathcal{H}} = \langle u ; \mathcal{A}^* v \rangle_{\mathcal{H}}$ since

$$\begin{aligned}
\langle \mathcal{A}u ; v \rangle_{\mathcal{H}} &= \langle (-u_2, -\Delta(u_1 + \varepsilon u_2)) ; (v_1, v_2) \rangle_{\mathcal{H}} \\
&= -a(u_2; v_1) - (\Delta(u_1 + \varepsilon u_2); v_2) \\
&= a(u_1; v_2) - a(u_2; v_1) + \varepsilon a(u_2; v_2) \\
&= a(u_1; v_2) + (u_2; \Delta(v_1 - \varepsilon v_2)) = \langle u ; \mathcal{A}^* v \rangle_{\mathcal{H}}.
\end{aligned}$$                                         $\square$

In the following we derive the strong stability estimates.

**Lemma 4.1.0.5.** Let $\Phi = (\phi, \dot{\phi})$ be the solution of the vector dual problem (4.3), then the following stability estimates for all $t \in [0, T]$ hold

$$\|\Phi(t)\|_{\mathcal{H}}^2 + 2\varepsilon \int_t^T \|\dot{\phi}(\tau)\|_{H^1(\Omega)}^2 d\tau = \|\Phi^{N-}\|_{\mathcal{H}}^2, \tag{4.4a}$$

$$\|\Phi(t)\|_{\widehat{\mathcal{H}}}^2 + 2\varepsilon \int_t^T \|\dot{\phi}(\tau)\|_{L^2(\Omega)}^2 d\tau = \|\Phi^{N-}\|_{\widehat{\mathcal{H}}}^2. \tag{4.4b}$$

Moreover in case of (DD*) boundary conditions there holds

$$\|\Delta\phi\|_{L^\infty(L^2)} + \|\dot{\phi}\|_{L^\infty(H^1)}^2 + \varepsilon \int_0^T \|\Delta\dot{\phi}(\tau)\|_{L^2(\Omega)}^2 d\tau \leq C\{\|\phi_1^{N-}\|_{L^2(\Omega)}^2 + \|\phi_2^{N-}\|_{H^1(\Omega)}^2\}. \tag{4.4c}$$

If $\varepsilon = 0$ then

$$\|\ddot{\phi}\|_{L^\infty(L^2)} + \|\Delta\phi\|_{L^\infty(L^2)} + \|\dot{\phi}\|_{L^\infty(H^1)}^2 = C\{\|\phi_1^{N-}\|_{L^2(\Omega)}^2 + \|\phi_2^{N-}\|_{H^1(\Omega)}^2\}, \tag{4.4d}$$

**Proof.** Scalar multiplication of (4.1a) with $\dot{\phi}$ with respect to $L^2(\Omega)$ gives

$$(\ddot{\phi}; \dot{\phi}) - (\Delta\phi; \dot{\phi}) + \varepsilon(\Delta\dot{\phi}; \dot{\phi}) = 0.$$

An integration by parts with respect to the spatial variable $x$ and the boundary conditions (DN*) or (DD*) from (4.1) lead to

$$\frac{1}{2}\frac{\partial}{\partial\tau}\|\dot{\phi}(\tau)\|_{L^2(\Omega)}^2 + \frac{1}{2}\frac{\partial}{\partial\tau}\|\phi(\tau)\|_{H^1(\Omega)}^2 d\tau - \varepsilon\|\dot{\phi}(\tau)\|_{H^1(\Omega)}^2 = 0.$$

An integration over the time interval $[t, T]$ and the main theorem of calculus in time, yield

$$\frac{1}{2}\|\dot{\phi}(T)\|_{L^2(\Omega)}^2 - \frac{1}{2}\|\dot{\phi}(t)\|_{L^2(\Omega)}^2 + \frac{1}{2}\|\phi(T)\|_{H^1(\Omega)}^2 - \frac{1}{2}\|\phi(t)\|_{H^1(\Omega)}^2 - \varepsilon\int_t^T \|\dot{\phi}(\tau)\|_{H^1(\Omega)}^2 d\tau = 0.$$

Finally, the use of the condition (4.1b) proves the statement (4.4a).

For the proof of the second stability assertion, we proceed similarly by multiplying (4.1a) with $\mathcal{K}\dot{\phi}$ with respect to the $L^2(\Omega)$ scalar product. This leads to

$$(\ddot{\phi}; \mathcal{K}\dot{\phi}) - (\Delta(\phi - \varepsilon\dot{\phi}); \mathcal{K}\dot{\phi}) = 0.$$

An integration by parts in space in the second term shows

$$(\ddot{\phi}; \mathcal{K}\dot{\phi}) + a(\phi - \varepsilon\dot{\phi}; \mathcal{K}\dot{\phi}) - (D\phi - \varepsilon D\dot{\phi})(1)\mathcal{K}\dot{\phi}(1) + (D\phi - \varepsilon D\dot{\phi})(0)\mathcal{K}\dot{\phi}(0) = 0.$$

From $\Phi \in \mathscr{D}(\mathcal{A}^*)$ and the fact that $\mathcal{K}\dot{\phi}$ is the solution of the steady-state problem (1.14), the boundary terms in the equation above vanish for either (DD*) or (DN*) boundary conditions. Moreover, for $\dot{\phi} \in L^2(\Omega)$ there holds $a(\mathcal{K}\dot{\phi}; v) = (\dot{\phi}; v)$ for all $v \in H_D^1(\Omega)$ with the definition of the operator $\mathcal{K}$. Therefore, the last equality simplifies to

$$(\mathcal{K}^{1/2}\ddot{\phi}; \mathcal{K}^{1/2}\dot{\phi}) + (\phi; \dot{\phi}) - \varepsilon(\dot{\phi}; \dot{\phi}) = 0.$$

Furthermore, an integration with respect to the time interval $[t, T]$ and the main theorem of calculus yield

$$\frac{1}{2}\|\mathcal{K}^{1/2}\dot{\phi}(T)\|^2_{L^2(\Omega)} - \frac{1}{2}\|\mathcal{K}^{1/2}\dot{\phi}(t)\|^2_{L^2(\Omega)} + \frac{1}{2}\|\phi(T)\|^2_{L^2(\Omega)} - \frac{1}{2}\|\phi(t)\|^2_{L^2(\Omega)} - \varepsilon\int_t^T \|\dot{\phi}(\tau)\|^2_{L^2(\Omega)}d\tau = 0.$$

This concludes the proof of the second stability assertion.

For the proof of the third stability estimate (4.4c), we proceed similarly as in case of the first two estimates. Namely, by multiplying (4.1a) with $-\Delta\dot{\phi}$ with respect to $L^2(\Omega)$ scalar product and integrating in time, we obtain

$$-\int_t^T (\ddot{\phi}; \Delta\dot{\phi})d\tau + \int_t^T (\Delta\phi; \Delta\dot{\phi})d\tau - \varepsilon\int_t^T (\Delta\dot{\phi}; \Delta\dot{\phi})d\tau = 0.$$

An integration by parts in space yields

$$\int_t^T a(\ddot{\phi}; \dot{\phi})d\tau - \int_t^T \ddot{\phi}(t, 1)D\dot{\phi}(t, 1)dt + \int_t^T (\Delta\phi; \Delta\dot{\phi})d\tau - \varepsilon\int_t^T (\Delta\dot{\phi}; \Delta\dot{\phi})d\tau = 0.$$

Obviously, the second term on the LHS in the equation above equals zero if $(DD^*)$ boundary conditions or $\varepsilon = 0$ hold. The main theorem of calculus applied to the last equation proves the third stability assertion, i.e.

$$0 = \int_t^T \frac{1}{2}\frac{\partial}{\partial\tau}\|D\dot{\phi}\|^2_{L^2(\Omega)}d\tau + \int_t^T \frac{1}{2}\frac{\partial}{\partial\tau}\|\Delta\phi\|^2_{L^2(\Omega)}d\tau - \varepsilon\int_t^T \|\Delta\dot{\phi}(\tau)\|^2_{L^2(\Omega)}d\tau$$

$$= \frac{1}{2}\|\phi_2^{N-}\|^2_{H^1(\Omega)} - \frac{1}{2}\|\phi_2(t)\|^2_{H^1(\Omega)} + \frac{1}{2}\|\Delta\phi_1^{N-}\|^2_{L^2(\Omega)} - \frac{1}{2}\|\Delta\phi_1(t)\|^2_{L^2(\Omega)} - \varepsilon\int_t^T \|\Delta\dot{\phi}(\tau)\|^2_{L^2\Omega}d\tau.$$

Furthermore, in order to prove (4.4d), recall that if $\varepsilon = 0$, $\ddot{\phi} = \Delta\phi$. If we replace $\|\Delta\phi\|_{L^\infty(L^2)}$ in (4.4c) by $\|\ddot{\phi}\|_{L^\infty(L^2)}$ and combine the resulting estimate with (4.4c), we complete the proof of the fourth stability estimate.                                                                                              $\square$

**Definition 4.1.0.2 (Dual bilinear form $\mathcal{B}^*$).** Let $\mathcal{B}^*$ be the dual bilinear form with respect to the bilinear form $\mathcal{B}$ defined in Subsection 2.3.1 defined such that

$$\mathcal{B}^*(v, u) := \mathcal{B}(u, v) \quad \text{for all} \quad u, v \in H^1(\mathscr{T}; \mathcal{H}). \tag{4.5}$$

$\square$

An explicit formel for $\mathcal{B}^*$ is derived later on with integration by parts in time for each discretisation in time separately.

## 4.2   A priori dual error estimate

In the following we try to determine the a priori error bound $\eta_a$ such that

$$\|e^{N-}\|_{\mathcal{H}} \le \eta_a.$$

Here $e^{N-}$ is the error related to the final time point $T$. The bound $\eta_a$ depends on space and time mesh-size $h, k$, respectively, as well as on the initial solution $u$.

| | $\varepsilon \geq 0$ | $\varepsilon = 0$ |
|---|---|---|
| **dG(0)** <br> Subsection 4.2.1.1 | $\mathcal{O}(h^p+k)$ $_{\mathcal{S}^{j-1}=\mathcal{S}^j}$ <br> $\mathcal{O}(h^p+k^{-1/2}h^{p+1}+k)$ $_{\text{otherwise}}$ | |
| **dG(1)** <br> Subsection 4.2.1.2 | $\mathcal{O}(h^p+k^2)$ $_{\mathcal{S}^{j-1}=\mathcal{S}^j}$ <br> $\mathcal{O}(h^p+k^{-1/2}h^p+k^2)$ $_{\text{otherwise}}$ | $\mathcal{O}(h^p+k^3)$ $_{\mathcal{S}^{j-1}=\mathcal{S}^j}$ <br> $\mathcal{O}(h^p+k^{-1/2}h^p+k^3)$ $_{\text{otherwise}}$ |
| **cG(1)** <br> Subsection 4.2.2 | $\mathcal{O}(h^p+k)$ | $\mathcal{O}(h^p+k^2)$ |

Table 4.2: Proven a priori error estimates for $\|e^{N^-}\|_{\mathcal{H}}$; $p=1$ for $\mathcal{P}_1$ elements and $p=3$ for $\mathcal{C}^1$ elements in space; dual method.

The main arguments in the subsequent dual-based a priori error analysis are the weak stability estimates. They depend on the properties of the weak dual problem, in particular on the time discretisation method.

The proven convergence order of a priori error estimates for different time discretisation methods and $\mathcal{P}_1$ and $\mathcal{C}^1$ ansatz in space are given in Table 4.2.

For comparison only, note that the results obtained by the dual method show better convergence rates than the one obtained by the energy method, cf. Table 4.2 and 3.2, respectively. In particular, by use of the energy method, we did not succeed to derive an a priori error estimate for $dG(1)$ time discretisation as well as for the $cG(1)\otimes\mathcal{P}_1$ case. Apart from that, the dual method also provides the better estimate for the $dG(0)$ time discretisation, see Figure 4.1. Namely for $h=k$ and $dG(0)\otimes\mathcal{P}_1$, the duality method proves convergence of order $\mathcal{O}(h)$.



Figure 4.1: Convergence of $\|e\|_{L^\infty(\mathcal{H})}$ for $dG(0)\otimes\mathcal{C}^1$ with respect to the number of elements in space; The a priori error bound $\eta_a=\mathcal{O}(h^p+k)$ is optimal for both $p=1,3$. For $p=1$, see Figure 3.1. For $p=3$ and $k=h^3$, the exact error is of order $\mathcal{O}(h^3)$; Example 6.2.1, $\varepsilon=0$, (DN), $T=1$.

The energy metod provides an estimate of order $\mathcal{O}(h^{1/2})$ for the same choice of time and space steps. For the proof of the optiamlity of the proven a priori error bound when $dG(1)$ in time, we refer to the Figure 3.5 where besides the proven a posteriori energy error estimate, the

error in the energy norm and the dissipative term have been ploted. We also refer to Figure 3.6, where the bahavior of the exact error terms when $cG(1)$ in time is showed.

The approach and techniques used in the following adopt the one presented in JOHNSON [42] for the problem $dG(1) \otimes \mathcal{P}_1$ with $\varepsilon = 0$ whereby our estimate show better convergence rates. Namely for each $k = h^\alpha$ where $\alpha > 2$, we have an estimate of order $\mathcal{O}(h)$ whereby in JOHNSON [42] proven a priori estimate is weaker, i.e. $\mathcal{O}(h^{2-\alpha/2})$.

### 4.2.1   $dG(q)$ time approximation, $q=0,1$

For the notation and definitions used within this section we refer to Section 2.1 and Subsection 2.3.3 where space discretisation and discontinuous Galerkin approximation are introduced, respectively. The main tool in the following analysis is the bilinear form $\mathcal{B}$ from (2.18) and its dual counterpart $\mathcal{B}^*$.

We start with an explicit representation of the dual bilinear form $\mathcal{B}^*$ for $dG$ ansatz in time.

**Lemma 4.2.1.1 (Dual bilinear form $\mathcal{B}^*$, $dG$ time approximation).** If discontinuous Galerkin time discretisation is performed, for sufficiently (piecewise) smooth functions in time $U, V$, the dual bilinear form $\mathcal{B}^*$ of the bilinear form $\mathcal{B}$ from (2.18) reads

$$\mathcal{B}^*(V,U) = -\sum_{j=1}^{N} \int_{I_j} \langle V_\tau ; U \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \int_{I_j} a(V_2 ; U_1) dt - \sum_{j=1}^{N} \int_{I_j} a(V_1 - \varepsilon V_2 ; U_2) dt$$

$$- \sum_{j=2}^{N} \langle [V]^{j-1} ; U^{j-1-} \rangle_{\mathcal{H}} + \langle V^{N-} ; U^{N-} \rangle_{\mathcal{H}} - \langle V^{0+} ; U^{0-} \rangle_{\mathcal{H}}. \tag{4.6}$$

**Proof.** The representation of the bilinear dual form $\mathcal{B}^*$ (4.6) arises from the definition of $\mathcal{B}$ (2.18) if we integrate by parts with respect to time variable. Namely,

$$\mathcal{B}(U,V) = -\sum_{j=1}^{N} \int_{I_j} \langle V_\tau ; U \rangle_{\mathcal{H}} dt + \sum_{j=1}^{N} \langle V^{j-} ; U^{j-} \rangle_{\mathcal{H}} - \sum_{j=1}^{N} \langle V^{j-1+} ; U^{j-1+} \rangle_{\mathcal{H}}$$

$$+ \sum_{j=1}^{N} \int_{I_j} a(V_2 ; U_1) dt - \sum_{j=1}^{N} \int_{I_j} a(V_1 - \varepsilon V_2 ; U_2) dt + \sum_{j=1}^{N} \langle V^{j-1+} ; U^{j-1+} - U^{j-1-} \rangle_{\mathcal{H}}. \tag{4.7}$$

Furthermore, the sum of the jump terms and the boundary terms equals

$$\sum_{j=1}^{N} \langle V^{j-} ; U^{j-} \rangle_{\mathcal{H}} - \sum_{j=1}^{N} \langle V^{j-1+} ; U^{j-1+} \rangle_{\mathcal{H}} + \sum_{j=1}^{N} \langle V^{j-1+} ; U^{j-1+} \rangle_{\mathcal{H}} - \sum_{j=1}^{N} \langle V^{j-1+} ; U^{j-1-} \rangle_{\mathcal{H}}$$

$$= \sum_{j=2}^{N+1} \langle V^{j-1-} ; U^{j-1-} \rangle_{\mathcal{H}} - \sum_{j=1}^{N} \langle V^{j-1+} ; U^{j-1-} \rangle_{\mathcal{H}}$$

$$= -\sum_{j=2}^{N} \langle [V]^{j-1} ; U^{j-1-} \rangle_{\mathcal{H}} + \langle V^{N-} ; U^{N-} \rangle_{\mathcal{H}} + \langle V^{0+} ; U^{0-} \rangle_{\mathcal{H}}. \tag{4.8}$$

A substitution of (4.8) into (4.7) and the definition (4.5) yield the proof.                           $\square$

Having introduced the dual bilinear form, we can formulate the discrete weak dual problem: Given $\Psi^{N-} \in (\mathcal{S}^N)^2 \subseteq \mathcal{H}$, find $\Psi \in \mathcal{Q}_q$ such that

$$\mathcal{B}^*(\Psi, V) + \left\langle \Psi^{0+} ; V^{0-} \right\rangle_{\mathcal{H}} = \left\langle \Psi^{N-} ; V^{N-} \right\rangle_{\mathcal{H}} \quad \text{for all} \quad V \in \mathcal{Q}_q. \tag{4.9}$$

$\Psi$ is called the discrete weak solution of the dual problem (4.3).

**Remark 4.2.1.1.** Analogue to the analysis presented in Subsections 2.4.1.1 and 2.4.2.1, we may show that the problem (4.9) has an unique solution. $\qquad\square$

**Remark 4.2.1.2.** Note that for every $V \in L^2(\mathscr{T}, H_D^1 \times H_D^1)$ the continuous dual solution $\Phi$ of (4.3) is the weak solution of (4.9). $\qquad\square$

**Lemma 4.2.1.2 (Stability of the discrete weak dual $dG$ solution).** The solution $\Psi = (\Psi_1, \Psi_2)$ of (4.9) satisfies the following stability estimate

$$\|\Psi^{n-}\|_{\mathcal{H}}^2 + 2\varepsilon \sum_{j=n}^N \int_{I_j} \|\Psi_2(t)\|_{H^1(\Omega)}^2 dt + \sum_{j=n}^N \|[\Psi]^j\|_{\mathcal{H}}^2 = \|\Psi^{N-}\|_{\mathcal{H}}^2, \tag{4.10a}$$

where $1 \leq n \leq N$. Moreover, in case of $\Psi \in \mathcal{Q}_1$

$$\|\Psi_{1,\tau}\|_{L^\infty(L^2)}^2 \leq \|\Psi^{N-}\|_{\mathcal{H}}^2, \tag{4.10b}$$

where additionally for $\varepsilon = 0$

$$\|\Psi_{2,\tau}\|_{L^\infty(H^{-1})}^2 \leq \|\Psi^{N-}\|_{\mathcal{H}}^2. \tag{4.10c}$$

**Proof.** The proof follows analogously as in case of weak "forward problem", see Lemma 2.3.3.5. $\qquad\square$

To derive an a priori error bound for $\|e^{N-}\|_{\mathcal{H}}$, the error $e := u - U$ will be decomposed in two parts, namely

$$e = \rho - \theta, \quad \text{where} \quad \rho := u - \mathcal{J}\Pi u \quad \text{and} \quad \theta := U - \mathcal{J}\Pi u. \tag{4.11}$$

Here $\Pi$ denotes the spatial multi projection. $\mathcal{J}$ is the mapping onto the space of discrete functions in time which will be introduced for $dG(0)$ or $dG(1)$ method separately.

The idea is to dominate $\|e^{N-}\|_{\mathcal{H}}$ by estimating $\|\rho^{N-}\|_{\mathcal{H}}$ and $\|\theta^{N-}\|_{\mathcal{H}}$. The estimation of $\|\rho^{N-}\|_{\mathcal{H}}$ is done for each ansatz in time owing to the approximation properties of the corresponding projections. To bound $\|\theta^{N-}\|_{\mathcal{H}}$, we introduce the following lemma.

**Lemma 4.2.1.3.** For $\theta, \rho$ from (4.11) there holds

$$\|\theta^{N-}\|_{\mathcal{H}}^2 = \left\langle \Psi^{0+} ; \theta^{0-} \right\rangle_{\mathcal{H}} + \sum_{j=1}^N \int_{I_j} a(\Psi_2 - \Psi_{1\tau}; \rho_1) dt - \sum_{j=1}^N \int_{I_j} (\Psi_{2,\tau}; \rho_2) dt - \sum_{j=1}^N \int_{I_j} a(\Psi_1; \rho_2) dt$$

$$+ \varepsilon \sum_{j=1}^N \int_{I_j} a(\Psi_2; \rho_2) dt - \sum_{j=2}^N a([\Psi_1]^{j-1}; \rho_1^{j-1-}) - \sum_{j=2}^N ([\Psi_2]^{j-1}; \rho_2^{j-1-})$$

$$+ a(\Psi_1^{N-}; \rho_1^{N-}) + (\Psi_2^{N-}; \rho_2^{N-}) =: \sum_{\ell=1}^9 E_\ell, \tag{4.12}$$

when $\Psi$ solves (4.9) with $\Psi^{N-} = \theta^{N-}$.

**Proof.** Owing to the definition of $\theta$, we have $\theta \in \mathcal{Q}_q$. Given (4.9) with $\Psi^{N-} = \theta^{N-}$, let $V = \theta$. This yields

$$\|\theta^{N-}\|_{\mathcal{H}}^2 = \left\langle \Psi^{0+} ; \theta^{0-} \right\rangle + \mathcal{B}^*(\Psi, \theta).$$

From Definition 4.1.0.2 and the Galerkin orthogonality (2.9), we further have

$$\mathcal{B}^*(\Psi, \theta) = \mathcal{B}(\theta, \Psi) = \mathcal{B}(U - \mathcal{J}\Pi u, \Psi) = \mathcal{B}(u - \mathcal{J}\Pi u, \Psi) = \mathcal{B}(\rho, \Psi) = \mathcal{B}^*(\Psi, \rho).$$

The combination of the last two equations and the representation of $\mathcal{B}^*$ in (4.6) yield the proof. □

With the representation from Lemma 4.2.1.3, we derive in the following an a priori error estimate for each ansatz in time separately.

### 4.2.1.1 A priori dual error analysis, $dG(0)$ time approximation

**Theorem 4.2.1.1 (A priori dual error estimate, $dG(0)$ time approximation).** There is a constant $C > 0$ independent of $u$, its discrete $dG(0)$ counterpart $U$ and $h, k$ such that

1. if $\mathcal{S}^{j-1} = \mathcal{S}^j$ for all $j = 1, \ldots, N$

$$\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C\Big\{ &\|k\Delta\dot{u}_1\|_{L^1(L^2)} + \sqrt{\varepsilon}\big(\|k\dot{u}_2\|_{L^2(H^1)} + \|k\Delta\dot{u}_2\|_{L^2(L^2)}\big) + \|k\dot{u}_2\|_{L^1(H^1)} \\
&+ \sqrt{\varepsilon}\|k\dot{u}_2\|_{L^2(H^1)} + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)}^2 + \|h^{p+1}D^{p+1}y_1\|_{L^1(L^2)}^2 \\
&+ \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)} + \|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)} \Big\},
\end{aligned} \tag{4.13}$$

2. if $\mathcal{S}^{j-1} \neq \mathcal{S}^j$ for all at least one $j = 1, \ldots, N$

$$\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C\Big\{ &\|k\Delta\dot{u}_1\|_{L^1(L^2)} + \sqrt{\varepsilon}\big(\|k\dot{u}_2\|_{L^2(H^1)} + \|k\Delta\dot{u}_2\|_{L^2(L^2)}\big) + \|k\dot{u}_2\|_{L^1(H^1)} + \sqrt{\varepsilon}\|k\dot{u}_2\|_{L^2(H^1)} \\
&+ \Big(\sum_{j=2}^{N} \|h^p D^{p+1}u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} + \Big(\sum_{j=2}^{N} \|h^{p+1}D^{p+1}u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \\
&+ \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)} + \|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)} \Big\},
\end{aligned} \tag{4.14}$$

provided that $u \in (H^{p+1}(\Omega))^2$ where $p = 1$ for linear splines in space and $p = 3$ for cubic splines in space. The 2th summand on the RHS of (4.13) and (4.14) does not appear in case of the $(DD^*)$ problem.

**Remark 4.2.1.3.** The estimate (4.13) is of order $\mathcal{O}(h^p + k)$. The estimate (4.14) of order $\mathcal{O}(h^p + k^{-1/2}h^p + k)$. □

The remaining part of this subsection is devoted to the proof of the Theorem 4.2.1.1.

Given the error decomposition (4.11), we may choose $\Pi = (\mathcal{G}, \mathcal{G})$. Let the mapping $\mathcal{J}$ be defined such that

$$\mathcal{J}u|_{I_j} := u^{j-} \quad \text{for all} \quad j = 1, \ldots, N, \quad \text{and} \quad u \in H^1(\mathscr{T}) \tag{4.15a}$$

and

$$(\mathcal{J}u)^{0-} := \begin{cases} u^{0-} & \text{if} \quad u \in H^1((-\infty, 0] \cup \mathscr{T}), \\ u(0) & \text{if} \quad u \in \mathcal{C}[0, T]. \end{cases} \tag{4.15b}$$

Obviously $\mathcal{J}u$ belongs to the space of constant functions in time but $\mathcal{J}$ is not $L^2$ orthogonal. We also make use of the approximation properties of the projection $\mathcal{G}$. These are given in Lemma 3.1.0.6.

**Lemma 4.2.1.4.** If the discrete variant of the initial solution $u_0$ is defined by

$$U^{0-} := \Pi u_0, \tag{4.16}$$

then $E_1 = 0$.

**Proof.** After the definition of the $\mathcal{H}$ scalar product and projection $\mathcal{J}$ we have

$$E_1 = \left\langle \Psi^{0+} ; U^{0-} - (\mathcal{J}\Pi u)^{0-} \right\rangle_{\mathcal{H}} = \left\langle \Psi^{0+} ; U^{0-} - \Pi u_0 \right\rangle_{\mathcal{H}} = 0. \qquad \square$$

**Lemma 4.2.1.5.** There holds for $E_2$,

$$E_2 \leq \|k\Delta \dot{u}_1\|_{L^1(L^2)} \|\Psi_2\|_{L^\infty(L^2)} + \left(\|k\dot{u}_2\|_{L^2(H^1)} + \|k\Delta \dot{u}_2\|_{L^2(L^2)}\right)\varepsilon\|\Psi_2\|_{L^2(H^1)},$$

where the second summand does not appear in case of the (DD*) boundary conditions.

**Proof.** Note that $\Phi_{1,\tau} = 0$ according to the $dG(0)$ time discretisation. With the orthogonality properties of the Galerkin projection and integration by parts in space, we obtain

$$E_2 = \sum_{j=1}^{N} \int_{I_j} a(\Psi_2; u_1 - \mathcal{G}u_1)dt + \sum_{j=1}^{N} \int_{I_j} a(\Psi_2; \mathcal{G}u_1 - \mathcal{G}u_1^{j-})dt$$

$$= -\sum_{j=1}^{N} \int_{I_j} (\Psi_2; \Delta(u_1 - u_1^{j-}))dt + \sum_{j=1}^{N} \int_{I_j} \Psi_2(t, 1)D(u_1 - u_1^{j-})(t, 1)dt. \tag{4.17}$$

The second term on the RHS of the second equality above vanishes for (DD*) boundary conditions. If the initial problems satisfies the (DN*) boundary conditions, the main theorem of calculus with respect to space interval $\Omega = [0, 1]$ implies

$$\Psi_2(t, 1) = \Psi_2(t, 0) + \int_{\Omega} D\Psi_2(x)dx \leq \|\Psi_2(t)\|_{H^1(\Omega)}. \tag{4.18}$$

Furthermore, from the boundary condition $D(u_1(t, 1) + \varepsilon u_2(t, 1)) = 0$ and a trace inequality we have for all $t \in I_j$

$$D(u_1(t, 1) - u_1(t_j, 1)) = -\varepsilon D(u_2(t, 1) - u_2(t_j, 1))$$
$$\leq \varepsilon \|D(u_2(t) - u_2(t_j))\|_{L^2(\Omega)} + \varepsilon\|\Delta(u_2(t) - u_2(t_j))\|_{L^2(\Omega)}.$$

Combining the last two estimates and using

$$u(t) - u(t_j) = \int_{t_j}^{t} \dot{u}(\tau)d\tau, \tag{4.19}$$

we obtain

$$\sum_{j=1}^{N}\int_{I_j}\Psi_2(t,1)D(u_1(t,1)-u_1(t_j,1))dt\leq\varepsilon\|\Psi_2\|_{L^2(H^1)}(\|k\dot{u}_2\|_{L^2(H^1)}+\|k\Delta\dot{u}_2\|_{L^2(L^2)}).\qquad(4.20)$$

With the same technique we prove

$$-\sum_{j=1}^{N}\int_{I_j}(\Psi_2;\Delta(u_1-u_1^{j-}))dt\leq\|k\Delta\dot{u}_1\|_{L^1(L^2(\Omega))}\|\Psi_2\|_{L^\infty(L^2)}.$$

From (4.20) and the last inequality we conclude the proof of the Lemma.                          □

**Lemma 4.2.1.6.** $E_3=0.$

**Proof.** Follows from the fact that in case of the $dG(0)$ time approximation method $\Phi_{2,\tau}=0.$□

**Lemma 4.2.1.7.** For $E_4$ there holds

$$E_4\leq\|k\dot{u}_2\|_{L^1(H^1)}\|\Psi_1\|_{L^\infty(H^1)}.$$

**Proof.** On account to the orthogonality properties of the projection $\mathcal{G}$ we have

$$E_4=-\sum_{j=1}^{N}\int_{I_j}a(\Psi_1;u_2-\mathcal{G}u_2)dt-\sum_{j=1}^{N}\int_{I_j}a(\Psi_1;\mathcal{G}u_2-\mathcal{G}u_2^{j-})dt=-\sum_{j=1}^{N}\int_{I_j}a(\Psi_1;u_2-u_2^{j-})dt.$$

Then, similar as in the proof of $E_3$, we deduce

$$E_4\leq\|k\dot{u}_2\|_{L^1(H^1)}\|\Psi_1\|_{L^\infty(H^1)}.$$                                         □

**Lemma 4.2.1.8.** For $E_5$ there holds

$$E_5\leq\|k\dot{u}_2\|_{L^2(H^1)}\varepsilon\|\Psi_2\|_{L^2(H^1)}.$$

**Proof.** With the same arguments as for the Lemma 4.2.1.7, there holds

$$E_5=\varepsilon\sum_{j=1}^{N}\int_{I_j}a(\Psi_2;u_2-u_2^{j-})dt\leq\|k\dot{u}_2\|_{L^2(H^1)}\varepsilon\|\Psi_2\|_{L^2(H^1)}.$$         □

**Lemma 4.2.1.9.** If $\mathcal{S}^{j-1}\neq\mathcal{S}^{j}$ for at least one $j$ there is a constant $C$ such that

$$E_6\leq C\Big(\sum_{j=2}^{N}\|h^pD^{p+1}u_1(t_{j-1})\|^2_{L^2(\Omega)}\Big)^{1/2}\Big(\sum_{j=2}^{N}\|[\Psi_1]^{j-1}\|^2_{H^1(\Omega)}\Big)^{1/2}.$$

Otherwise $E_6=0.$

**Proof.** If for all $j=1,\ldots,N$ there holds $\mathcal{S}^{j-1}=\mathcal{S}^j$, then owing to the orthogonal property of the projection $\mathcal{G}$ and the definition of the projection $\mathcal{J}$, we have

$$E_6 = -\sum_{j=2}^{N} a(\Psi_1^j - \Psi_1^{j-1}; (u_1 - \mathcal{G}u_1)^{j-1-}) = 0.$$

Otherwise, the application of Hölder, a discrete Cauchy inequality, and the approximation property of the Galerkin projection $\mathcal{G}$ yield

$$E_6 \leq C\Big(\sum_{j=2}^{N} \|h^p D^{p+1} u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \Big(\sum_{j=2}^{N} \|[\Psi_1]^{j-1}\|_{H^1(\Omega)}\Big)^{1/2},$$

for some constant $C > 0$. $\qquad\square$

**Lemma 4.2.1.10.** There exists a constant $C$ such that

$$E_7 + E_8 + E_9 \leq C\big(\|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} y_1\|_{L^1(L^2)}\big)\|\Psi_2\|_{L^\infty(L^2)}.$$

if $\mathcal{S}^{j-1}=\mathcal{S}^j$ for all $j=1,\ldots,N$. Otherwise

$$E_7 + E_8 + E_9 \leq C\Big(\sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \Big(\sum_{j=2}^{N} \|[\Psi_2]^{j-1}\|_{L^2(\Omega)}^2\Big)^{1/2}$$
$$+ C\|h^{p+1} D^{p+1} u_2(T)\|_{L^1(L^2)}\|\theta_2^{N-}\|_{L^2(\Omega)}.$$

**Proof.** Owing to the orthogonality property of $\mathcal{J}, \mathcal{G}$, we have

$$E_7 + E_8 + E_9 = -\sum_{j=2}^{N} ([\Psi_2]^{j-1}; (u_2 - \mathcal{G}u_2)^{j-1-}) + a(\theta_1^{N-}; (u_1 - \mathcal{G}u_1)^{N-}) + (\theta_2^{N-}; (u_2 - \mathcal{G}u_2)^{N-})$$
$$= -\sum_{j=2}^{N} ([\Psi_2]^{j-1}; (u_2 - \mathcal{G}u_2)^{j-1-}) + (\theta_2^{N-}; (u_2 - \mathcal{G}u_2)^{N-}). \qquad (4.21)$$

If we assume that space mesh does not change in time, then $(\mathcal{G}u_2)^{j-1+} = (Gu_2)^{j-1-} = \mathcal{G}u_2(t_{j-1})$. An integration by parts in time yields

$$E_7 + E_8 + E_9 = \sum_{j=1}^{N} (\Psi_2; \dot{u}_2 - \mathcal{G}\dot{u}_2) + (\Psi^{0+}; y_1 - \mathcal{G}y_1).$$

Hölder inequality and approximation properties of $\mathcal{G}$ yield the following estimate

$$E_7 + E_8 + E_9 \leq C\big(\|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} y_1\|_{L^1(L^2)}\big)\|\Psi_2\|_{L^\infty(L^2)}.$$

We continue by estimating $E_7 + E_8 + E_9$ under assumption that there exists at least one $j$ such that $\mathcal{S}^{j-1} \neq \mathcal{S}^j$. Using the same arguments as before, we have from (4.21)

$$E_7 + E_8 + E_9 \leq C\Big(\sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \Big(\sum_{j=2}^{N} \|[\Psi_2]^{j-1}\|_{L^2(\Omega)}^2\Big)^{1/2}$$
$$+ C\|h^{p+1} D^{p+1} u_2(T)\|_{L^1(L^2)}\|\theta_2^{N-}\|_{L^2(\Omega)}.$$

This concludes the proof. $\qquad\square$

Incorporating all previous estimates into the representation (4.12) and recalling the stability Lemma 4.2.1.2 we have

1. if $\mathcal{S}^{j-1}=\mathcal{S}^j$ for all $j=1,\ldots,N$

$$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{\|k\Delta\dot{u}_1\|_{L^1(L^2)}+\sqrt{\varepsilon}\big(\|k\dot{u}_2\|_{L^2(H^1)}+\|k\Delta\dot{u}_2\|_{L^2(L^2)}\big)+\|k\dot{u}_2\|_{L^1(H^1)}$$
$$+\sqrt{\varepsilon}\|k\dot{u}_2\|_{L^2(H^1)}+\|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)}+\|h^{p+1}D^{p+1}y_1\|_{L^1(L^2)}\Big\}. \qquad (4.22)$$

2. if $\mathcal{S}^{j-1}\neq\mathcal{S}^j$ for at least one $j=1,\ldots,N$

$$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{\|k\Delta\dot{u}_1\|_{L^1(L^2)}+\sqrt{\varepsilon}\big(\|k\dot{u}_2\|_{L^2(H^1)}+\|k\Delta\dot{u}_2\|_{L^2(L^2)}\big)+\|k\dot{u}_2\|_{L^1(H^1)}+\sqrt{\varepsilon}\|k\dot{u}_2\|_{L^2(H^1)}$$
$$+\Big(\sum_{j=2}^{N}\|h^p D^{p+1}u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2}+\Big(\sum_{j=2}^{N}\|h^{p+1}D^{p+1}u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2}$$
$$+\|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)}\Big\}. \qquad (4.23)$$

In order to complete the proof of Theorem 4.2.1.1 we need to estimate $\|\rho^{N-}\|_{\mathcal{H}}$. This is done in the following lemma.

**Lemma 4.2.1.11.** For $\rho=u-\mathcal{J}\Pi u$ where $\Pi=(\mathcal{G},\mathcal{G})$ and the temporal projection $\mathcal{J}$ where $(\mathcal{J}u)^{j-}=u^{j-}$ there is a constant $C$ such that

$$\|\rho^{N-}\|_{\mathcal{H}} \leq C(\|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)}+\|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)}). \qquad (4.24)$$

**Proof.** Due to the approximation properties of projection $\Pi$ and mapping $\mathcal{J}$ we have

$$\|\rho^{N-}\|_{\mathcal{H}}^2 = \|(u-\Pi u)^{N-}\|_{\mathcal{H}}^2 = \|(u_1-\mathcal{G}u_1)(T)\|_{H^1(\Omega)}^2+\|(u_2-\mathcal{G}u_2)(T)\|_{L^2(\Omega)}^2$$
$$\leq C\big(\|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)}^2+\|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)}^2\big). \qquad \square$$

According to the decomposition (4.11), the sum of (4.23) and (4.24) yields the proof of Theorem 4.2.1.1.

### 4.2.1.2 A priori dual error analysis, $dG(1)$ time approximation

**Theorem 4.2.1.2 (A priori dual error estimate, $dG(1)$ time approximation, $\mathcal{S}^{j-1}=\mathcal{S}^j$).** If $\mathcal{S}^{j-1}=\mathcal{S}^j$ for all $j=1,\ldots,N$ then there is a constant $C$ independent of $u$, its discrete $dG(1)$ counterpart $U$ and mesh size $h,k$ such that

1. $\varepsilon\geq0$, (DD)

$$\|e^{N-}\|_{\mathcal{H}} \leq C\Big\{\|k^2\Delta\ddot{u}_1\|_{L^1(L^2)}+\|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)}+\|h^{p+1}D^{p+1}y_1\|_{L^1(L^2)}$$
$$+\|k^3\Delta\ddot{u}_2\|_{L^1(L^2)}+\sqrt{\varepsilon}\|k^2\ddot{u}_2\|_{L^2(H^1)}+\|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)}$$
$$+\|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)}\Big\}, \qquad (4.25a)$$

2. if $\varepsilon \geq 0$ and (DN)

$$
\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C \Big\{ &\|k^2\Delta\dddot{u}_1\|_{L^1(L^2)} + \sqrt{\varepsilon}\big(\|k^2\ddot{u}_2\|_{L^2(H^1)} + \|k^2\Delta\ddot{u}_2\|_{L^2(L^2)}\big) + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)} \\
&+ \|h^{p+1}D^{p+1}y_1\|_{L^1(L^2)} + \|k^2\ddot{u}_2\|_{L^1(H^1)} \\
&+ \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)} + \|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)} \Big\},
\end{aligned}
\tag{4.25b}
$$

3. if $\varepsilon = 0$ and (DD) or (DN)

$$
\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C \Big\{ &\|k^3\Delta\dddot{u}_1\|_{L^1(H^1)} + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)} + \|h^{p+1}D^{p+1}y_1\|_{L^1(L^2)} \\
&+ \|k^3\Delta\dddot{u}_2\|_{L^1(L^2)} + \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)} + \|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)} \Big\},
\end{aligned}
\tag{4.25c}
$$

provided $u \in (H^{p+1}(\Omega))^2$ where $p=1$ for linear and $p=3$ for cubic splines in space. Additionally, we demand also $u \in (H^{p+2}(\Omega))^2$ for $p=1$ in the third estimate.

**Remark 4.2.1.4.** The estimates (4.25a) and (4.25b) are of order $\mathcal{O}(h^p + k^2)$. The estimate (4.25c) is of order $\mathcal{O}(h^p + k^3)$. $\qquad\square$

**Theorem 4.2.1.3 (A priori dual error estimate, $dG(1)$ time approximation, $\mathcal{S}^{j-1} \neq \mathcal{S}^j$).**
If $\mathcal{S}^{j-1} \neq \mathcal{S}^j$ for at least one $j = 1, \ldots, N$ then there is a constant $C$ independent of $u$, its discrete $dG(1)$ counterpart $U$ and mesh size $h, k$ such that

1. if $\varepsilon \geq 0$, (DD)

$$
\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C \Big\{ &\|k^2\Delta\dddot{u}_1\|_{L^1(L^2)} + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)} + \sum_{j=0}^{N}\|h^{p+1}D^{p+1}u_2(t_j)\|_{L^2(\Omega)} \\
&+ \Big(\sum_{j=2}^{N}\|h^{p+1}D^{p+1}u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} + \|k^3\Delta\dddot{u}_2\|_{L^1(L^2)} + \sqrt{\varepsilon}\|k^2\ddot{u}_2\|_{L^2(H^1)} \\
&+ \Big(\sum_{j=1}^{N}\|h^p D^{p+1}u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} + \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)} \Big\},
\end{aligned}
\tag{4.26a}
$$

2. if $\varepsilon \geq 0$ and (DN)

$$
\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C \Big\{ &\|k^2\Delta\dddot{u}_1\|_{L^1(L^2)} + \sqrt{\varepsilon}\big(\|k^2\ddot{u}_2\|_{L^2(H^1)} + \|k^2\Delta\ddot{u}_2\|_{L^2(L^2)}\big) + \|h^{p+1}D^{p+1}\dot{u}_2\|_{L^1(L^2)} \\
&+ \sum_{j=0}^{N}\|h^{p+1}D^{p+1}u_2(t_j)\|_{L^2(\Omega)} + \Big(\sum_{j=2}^{N}\|h^{p+1}D^{p+1}u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \\
&+ \|k^2\ddot{u}_2\|_{L^1(H^1)} + \Big(\sum_{j=2}^{N}\|h^p D^{p+1}u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} + \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)} \Big\},
\end{aligned}
\tag{4.26b}
$$

3. $\varepsilon = 0$ and (DD) or (DN)

$$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{\|k^3\Delta\ddot{u}_1\|_{L^1(H^1)} + \|h^p D^{p+1}u_2\|_{L^1(L^2)} + \Big(\sum_{j=2}^{N}\|h^{p+1}D^{p+1}u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2}$$

$$+ \|h^{p+1}D^{p+1}u_2(T)\|_{L^2(\Omega)} + \|k^3\Delta\ddot{u}_2\|_{L^1(L^2)}$$

$$+ \Big(\sum_{j=1}^{N}\|h^p D^{p+1}u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} + \|h^p D^{p+1}u_1(T)\|_{L^2(\Omega)}\Big\}, \qquad (4.26\text{c})$$

provided $u \in (H^{p+1}(\Omega))^2$ where $p=1$ for linear and $p=3$ for cubic splines in space. Additionally, we demand also $u \in (H^{p+2}(\Omega))^2$ for $p=1$ in the third estimate.

**Remark 4.2.1.5.** The estimates (4.26) are of order $\mathcal{O}(h^p + k^{-1/2}h^p + k^2)$. $\qquad\square$

We prove both theorems simultaneously.

**Proof (Theorem 4.2.1.2, Theorem 4.2.1.3).** Given the error decomposition (4.11), let $\Pi = (\mathcal{G}, \mathcal{G})$ and $\mathcal{J}$ be the time projection operator as in Definition 3.1.0.9, case $dG(1)b$. In the following we make use of the corresponding approximation lemmas, i.e. Lemma 3.1.0.6 and Lemma 3.1.0.11, respectively.

The idea is to estimate $E_1 - E_9$ such that the final estimates contains a priori known terms and discrete dual solution contributions which can be further estimated by means of Lemma 4.2.1.2. To estimate $E_1$, recall that $(\mathcal{J}\Pi u)^{0-} = \Pi u_0$ according to the definition of the temporal projection $\mathcal{J}$. Then there holds the same estimate for $E_1$ as the one already proven in case of $dG(0)$ method in time, see Lemma 4.2.1.4. The estimates for the remaining terms are given throughout the following lemmas.

**Lemma 4.2.1.12.** There exists a constant $C > 0$ such that for $\varepsilon > 0$ there holds

$$E_2 \leq C\big\{\|k^2\Delta\ddot{u}_1\|_{L^1(L^2)}\|\Psi_2\|_{L^\infty(L^2)} + (\|k^2\ddot{u}_2\|_{L^2(H^1)} + \|k^2\Delta\ddot{u}_2\|_{L^2(L^2)})\varepsilon\|\Psi_2\|_{L^2(H^1)}\big\},$$

where the second summand does not appear for $(DD^*)$ boundary conditions.
Moreover, if $\varepsilon = 0$, then

$$E_2 \leq C\|k^3\Delta\ddot{u}_1\|_{L^1(H^1)}\|\Psi_{2,\tau}\|_{L^\infty(H^{-1})}.$$

**Proof.** Owing to the properties of $\mathcal{G}$ and $\mathcal{J}$, we have

$$E_2 = \sum_{j=1}^{N}\int_{I_j} a(\Psi_2 - \Psi_{1,\tau}; u_1 - \mathcal{J}u_1)dt = \sum_{j=1}^{N}\int_{I_j} a(\Psi_2; u_1 - \mathcal{J}u_1)dt,$$

since $u_1 - \mathcal{J}u_1$ is orthogonal to the piecewise constant functions in time.
An integration by parts in space yields

$$E_2 = -\sum_{j=1}^{N}\int_{I_j}(\Psi_2; \Delta(u_1 - \mathcal{J}u_1))dt + \sum_{j=1}^{N}\int_{I_j}\Psi_2(t,1)D(u_1 - \mathcal{J}u_1)(t,1)dt. \qquad (4.27)$$

The second term on the RHS of the equation above equals zero when $(DD^*)$ boundary conditions hold or $\varepsilon = 0$. On the other hand, if $\varepsilon > 0$ and $(DN^*)$, i.e. $(DN)$ boundary conditions hold, the second term can be estimated by using the same arguments as in the proof of Lemma 4.2.1.5. Namely,

$$\sum_{j=1}^{N} \int_{I_j} \Psi_2(t,1) D(u_1 - \mathcal{J}u_1)(t,1) dt = -\varepsilon \sum_{j=1}^{N} \int_{I_j} \Psi_2(t,1) D(u_2 - \mathcal{J}u_2)(t,1) dt$$
$$\leq C\big(\|k^2 \ddot{u}_2\|_{L^2(H^1)} + \|k^2 \Delta \ddot{u}_2\|_{L^2(L^2)}\big)\varepsilon \|\Psi_2\|_{L^2(H^1)}.$$

We continue by estimating the first term on the RHS of (4.27). This completes the estimation for $E_2$. Namely, if $\varepsilon > 0$ then

$$-\sum_{j=1}^{N} \int_{I_j} (\Psi_2; \Delta(u_1 - \mathcal{J}u_1)) dt \leq C\|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} \|\Psi_2\|_{L^\infty(L^2)}.$$

For $\varepsilon = 0$ we conclude

$$E_2 = -\sum_{j=1}^{N} \int_{I_j} (\Psi_2 - \overline{\Psi}_2; \Delta(u_1 - \mathcal{J}u_1)) dt \leq C\|k^3 \Delta \ddot{u}_1\|_{L^1(H^1)} \|\Psi_{2,\tau}\|_{L^\infty(H^{-1})}. \qquad \square$$

In the following lemma we provide the estimate for $E_3 + E_7 + E_9$.

**Lemma 4.2.1.13.** If $\mathcal{S}^{j-1} = \mathcal{S}^j$ for all $j = 1, \dots, N$ and $\varepsilon \geq 0$ there is a constant $C$ such that there holds

$$E_3 + E_7 + E_9 \leq C\big(\|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)}\big) \|\Psi_2\|_{L^\infty(L^2)}. \qquad (4.28\text{a})$$

Otherwise, if there exists at least one $j$ such that $\mathcal{S}^{j-1} \neq \mathcal{S}^j$, then for all $\varepsilon \geq 0$

$$E_3 + E_7 + E_9 \leq C\Big\{ \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} \|\Psi_2\|_{L^\infty(L^2)} + \Big( \sum_{j=0}^{N} \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)} \Big) \|\Psi_2\|_{L^\infty(L^2)}$$
$$+ \Big( \sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2 \Big)^{1/2} \Big( \sum_{j=2}^{N} \|[\Psi_2]^{j-1}\|_{L^2(\Omega)}^2 \Big)^{1/2}$$
$$+ \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} \|\theta_2^{N-}\|_{L^2(\Omega)} \Big\} \qquad (4.28\text{b})$$

where for $\varepsilon = 0$ additionally holds

$$E_3 + E_7 + E_9 \leq C\Big\{ \|h^p D^{p+1} u_2\|_{L^1(L^2)} \|\Psi_{2,\tau}\|_{L^\infty(H^{-1})}$$
$$+ \Big( \sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2 \Big)^{1/2} \Big( \sum_{j=2}^{N} \|[\Psi_2]^{j-1}\|_{L^2(\Omega)}^2 \Big)^{1/2}$$
$$+ \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} \|\theta_2^{N-}\|_{L^2(\Omega)} \Big\}, \qquad (4.28\text{c})$$

**Proof.** Let us assume that $\mathcal{S}^{j-1} = \mathcal{S}^j$ for all $j = 1, \ldots, N$. Owing to the properties of $\mathcal{J}$ we have

$$E_3 + E_7 + E_9 = -\sum_{j=1}^{N} \int_{I_j} (\Psi_{2,\tau}; u_2 - \mathcal{G}u_2) dt - \sum_{j=2}^{N} ([\Psi_2]^{j-1}; (u_2 - \mathcal{G}u_2)^{j-1-}) + (\theta_2^{N-}; (u_2 - \mathcal{G}u_2)^{N-}).$$

If we integrate by parts in time in the first term, we obtain

$$E_3 + E_7 + E_9 = \sum_{j=1}^{N} \int_{I_j} (\Psi_2; \dot{u}_2 - \mathcal{G}\dot{u}_2) dt - \sum_{j=1}^{N} \left\{ (\Psi_2^{j-}; (u_2 - \mathcal{G}u_2)^{j-}) - (\Psi_2^{j-1+}; (u_2 - \mathcal{G}u_2)^{j-1+}) \right\}$$

$$- \sum_{j=2}^{N} ([\Psi_2]^{j-1}; (u_2 - \mathcal{G}u_2)^{j-1-}) + (\theta_2^{N-}; (u_2 - \mathcal{G}u_2)^{N-}).$$

Obviously, from $\mathcal{S}^{j-1} = \mathcal{S}^j$ we have $(\mathcal{G}u_2)^{j-1-} = (\mathcal{G}u_2)^{j-1+} = \mathcal{G}u_2(t_j)$. Then

$$E_3 + E_7 + E_9 = \sum_{j=1}^{N} \int_{I_j} (\Psi_2; \dot{u}_2 - \mathcal{G}\dot{u}_2) dt - (\theta_2^{N-}; (u_2 - \mathcal{G}u_2)(T)) - (\Psi_2^{0+}; y_1 - \mathcal{G}y_1)$$

$$+ \sum_{j=2}^{N} ([\Psi_2]^{j-1}; (u_2 - \mathcal{G}u_2)(t_{j-1})) - \sum_{j=2}^{N} ([\Psi_2]^{j-1}; (u_2 - \mathcal{G}u_2)(t_{j-1}))$$

$$+ (\theta_2^{N-}; (u_2 - \mathcal{G}u_2)(T)).$$

Hölder inequality and the approximation properties of $\mathcal{G}$ imply

$$E_3 + E_7 + E_9 \leq C \left( \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} y_1\|_{L^2(\Omega)} \right) \|\Psi_2\|_{L^\infty(L^2)}.$$

This concludes the proof for the case $\varepsilon \geq 0$ when the space mesh does not change in time, i.e. $\mathcal{S}^{j-1} = \mathcal{S}^j$ for all $j = 1, \ldots, N$.

We continue by assuming that there exists at least one $j$ such that $\mathcal{S}^{j-1} \neq \mathcal{S}^j$. Then

$$E_7 + E_9 \leq C \left( \sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2 \right)^{1/2} \left( \sum_{j=2}^{N} \|[\Psi_2]^{j-1}\|_{L^2(\Omega)}^2 \right)^{1/2}$$

$$+ C \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} \|\theta_2^{N-}\|_{L^2(\Omega)}.$$

In order to estimate $E_3$, we differ between two cases, case $\varepsilon = 0$ and case $\varepsilon > 0$. Let us assume that $\varepsilon = 0$. Then we my deduce

$$E_3 \leq C \|h^p D^{p+1} u_2\|_{L^1(L^2)} \|\Psi_{2,\tau}\|_{L^\infty(H^{-1})}. \tag{4.29}$$

For $\varepsilon > 0$, an integration by parts in time leads to

$$E_3 = \sum_{j=1}^{N} \int_{I_j} (\Psi_2; \dot{u}_2 - \mathcal{G}\dot{u}_2) dt - \sum_{j=1}^{N} (\Psi_2^{j-}; (u_2 - \mathcal{G}u_2)^{j-}) - (\Psi_2^{j-1+}; (u_2 - \mathcal{G}u_2)^{j-1+})$$

$$\leq C \left\{ \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \sum_{j=0}^{N} \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)} \right\} \|\Psi_2\|_{L^\infty(L^2)}. \tag{4.30}$$

Note that the estimate (4.30) also holds for $\varepsilon = 0$ when the space mesh differs from one time slab to another. However, we choose (4.29) because it shows the better convergence rates then (4.30). This concludes the proof. $\qquad \square$

**Lemma 4.2.1.14.** There exists a constant $C$ such that for $\varepsilon \geq 0$ and (DN*)

$$E_4 \leq C \| k^2 \ddot{u}_2 \|_{L^1(H^1)} \| \Psi_1 \|_{L^\infty(H^1)}.$$

Furthermore, if (DD*) or $\varepsilon = 0$, then

$$E_4 \leq C \| k^3 \Delta \ddot{u}_2 \|_{L^1(L^2)} \| \Psi_{1,\tau} \|_{L^\infty(L^2)}.$$

**Proof.** The orthogonality property of $\mathcal{G}$ implies

$$E_4 = -\sum_{j=1}^{N} \int_{I_j} a(\Psi_1; u_2 - \mathcal{J} u_2) dt. \tag{4.31}$$

If the (DD*) boundary conditions hold or $\varepsilon = 0$, an integration by parts in space and the approximation properties of $\mathcal{J}$ lead to

$$E_4 = \sum_{j=1}^{N} \int_{I_j} (\Psi_1 - \overline{\Psi}_1; \Delta(u_2 - \mathcal{J} u_2)) dt \leq C \| k^3 \Delta \ddot{u}_2 \|_{L^1(L^2)} \| \Psi_{1,\tau} \|_{L^\infty(L^2)}.$$

Otherwise, if (DN), i.e. (DN*) and $\varepsilon > 0$, we deduce from (4.31)

$$E_4 \leq C \| k^2 \ddot{u}_2 \|_{L^1(H^1)} \| \Psi_1 \|_{L^\infty(H^1)}. \qquad \square$$

**Lemma 4.2.1.15.** There is a constant $C$ such that

$$E_5 \leq C \| k^2 \ddot{u}_2 \|_{L^2(H^1)} \varepsilon \| \Psi_2 \|_{L^2(H^1)}.$$

**Proof.** From an application of the Hölder inequality, we may deduce

$$E_5 = -\varepsilon \sum_{j=1}^{N} \int_{I_j} a(\Psi_2; u_2 - \mathcal{J} u_2) dt \leq \varepsilon \| k^2 \ddot{u}_2 \|_{L^2(H^1)} \| \Psi_2 \|_{L^2(H^1)}. \qquad \square$$

**Lemma 4.2.1.16.** For $E_6$ there is a constant $C$ such that

$$E_6 \leq C \Big( \sum_{j=2}^{N} \| h^p D^{p+1} u_1(t_{j-1}) \|_{L^2(\Omega)}^2 \Big)^{1/2} \Big( \sum_{j=2}^{N} \| [\Psi_1]^{j-1} \|_{H^1(\Omega)}^2 \Big)^{1/2},$$

where the both sums go only over such $j$ where $\mathcal{S}^{j-1} \neq \mathcal{S}^j$.

**Proof.** If we assume that for each $j = 1, \ldots, N$ there holds $\mathcal{S}^{j-1} = \mathcal{S}^j$, then

$$E_6 = \sum_{j=2}^{N} a(\Psi_1^{j-1+} - \Psi_1^{j-1-}; (\mathcal{G} u_1 - u_1)^{j-1-}) = 0.$$

Otherwise, the approximation properties of $\mathcal{G}$ and the discrete Cauchy inequality yield for some $C > 0$,

$$E_6 \leq C \Big( \sum_{j=2}^{N} \| h^p D^{p+1} u_1(t_{j-1}) \|_{L^2(\Omega)}^2 \Big)^{1/2} \Big( \sum_{j=2}^{N} \| [\Psi_1]^{j-1} \|_{H^1(\Omega)}^2 \Big)^{1/2}. \qquad \square$$

**Lemma 4.2.1.17.** $E_8 = 0$.

**Proof.** On account to the properties of projection $\mathcal{J}$ and $\mathcal{G}$, we have

$$E_8 = \sum_{j=1}^{N} a(\theta_1^{N-}; (u_1 - \mathcal{G}u_1)^{N-}) = 0. \tag{4.32}$$

$\square$

Incorporating estimates for $E_1 - E_9$ into the representation (4.12) and recalling the stability estimate from Lemma 4.2.1.2 we have

1. if $\mathcal{S}^{j-1} = \mathcal{S}^j$ for all $j = 1, \ldots, N$ and

   a) $\varepsilon \geq 0$, (DD)

   $$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{ \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} y_1\|_{L^1(L^2)}$$
   $$+ \|k^3 \Delta \ddot{u}_2\|_{L^1(L^2)} + \sqrt{\varepsilon}\|k^2 \ddot{u}_2\|_{L^2(H^1)} \Big\}, \tag{4.33a}$$

   b) if $\varepsilon \geq 0$ and (DN)

   $$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{ \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} + \sqrt{\varepsilon}\big(\|k^2 \ddot{u}_2\|_{L^2(H^1)} + \|k^2 \Delta \ddot{u}_2\|_{L^2(L^2)}\big) + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}$$
   $$+ \|h^{p+1} D^{p+1} y_1\|_{L^1(L^2)} + \|k^2 \ddot{u}_2\|_{L^1(H^1)} \Big\}, \tag{4.33b}$$

   c) if $\varepsilon = 0$ and (DD) or (DN)

   $$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{ \|k^3 \Delta \ddot{u}_1\|_{L^1(H^1)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}$$
   $$+ \|h^{p+1} D^{p+1} y_1\|_{L^1(L^2)} + \|k^3 \Delta \ddot{u}_2\|_{L^1(L^2)} \Big\}, \tag{4.33c}$$

2. if $\mathcal{S}^{j-1} \neq \mathcal{S}^j$ for at least one $j = 1, \ldots, N$

   a) $\varepsilon \geq 0$, (DD)

   $$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{ \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)} + \sum_{j=0}^{N} \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)}$$
   $$+ \Big(\sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} + \|k^3 \Delta \ddot{u}_2\|_{L^1(L^2)}$$
   $$+ \sqrt{\varepsilon}\|k^2 \ddot{u}_2\|_{L^2(H^1)} + \Big(\sum_{j=1}^{N} \|h^p D^{p+1} u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \Big\}, \tag{4.33d}$$

   b) $\varepsilon \geq 0$ and (DN)

   $$\|\theta^{N-}\|_{\mathcal{H}} \leq C\Big\{ \|k^2 \Delta \ddot{u}_1\|_{L^1(L^2)} + \sqrt{\varepsilon}\big(\|k^2 \ddot{u}_2\|_{L^2(H^1)} + \|k^2 \Delta \ddot{u}_2\|_{L^2(L^2)}\big) + \|h^{p+1} D^{p+1} \dot{u}_2\|_{L^1(L^2)}$$
   $$+ \sum_{j=0}^{N} \|h^{p+1} D^{p+1} u_2(t_j)\|_{L^2(\Omega)} + \Big(\sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_2(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2}$$
   $$+ \|k^2 \ddot{u}_2\|_{L^1(H^1)} + \Big(\sum_{j=2}^{N} \|h^p D^{p+1} u_1(t_{j-1})\|_{L^2(\Omega)}^2\Big)^{1/2} \Big\}, \tag{4.33e}$$

c) $\varepsilon = 0$ and (DD) or (DN)

$$\|\theta^{N-}\|_{\mathcal{H}} \leq C \Big\{ \|k^3 \Delta \ddot{u}_1\|_{L^1(H^1)} + \|h^p D^{p+1} u_2\|_{L^1(L^2)} + \Big( \sum_{j=2}^{N} \|h^{p+1} D^{p+1} u_1(t_{j-1})\|_{L^2(\Omega)}^2 \Big)^{1/2}$$

$$+ \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} + \|k^3 \Delta \ddot{u}_2\|_{L^1(L^2)}$$

$$+ \Big( \sum_{j=1}^{N} \|h^p D^{p+1} u_1(t_{j-1})\|_{L^2(\Omega)}^2 \Big)^{1/2} \Big\}, \tag{4.33f}$$

The estimation of $\|\rho^{N-}\|_{\mathcal{H}}$ follows as Lemma 4.2.1.11. According to the decomposition (4.11), the combination of the estimates above with an estimate (4.24) yields an upper bound for $\|e^{N-}\|_{\mathcal{H}}$. This concludes the proof of Theorem 4.2.1.2 and Theorem 4.2.1.3. □

### 4.2.2 $cG(1)$ time approximation

For the notation and definitions used in the following, we refer to Section 2.3.4 for $cG(1)$ time approximation and to Section 2.1 for the spatial Galerkin discretisation.

**Lemma 4.2.2.1 (Dual bilinear form $\mathcal{B}^*$, $cG(1)$ time approximation).** In case of the continuous Galerkin discretisation in time, for (piecewise) smooth functions, globally continuous in time $u, v$ the dual bilinear form of the bilinear form (2.41) reads

$$\mathcal{B}^*(v, u) = -\int_0^T \langle \dot{v} ; u \rangle_{\mathcal{H}} dt + \int_0^T a(v_2; u_1) dt - \int_0^T a(v_1 - \varepsilon v_2; u_2) dt$$

$$+ \langle v(T) ; u(T) \rangle_{\mathcal{H}} - \langle v(0) ; u(0) \rangle_{\mathcal{H}}. \tag{4.34}$$

**Proof.** Given the definition of $\mathcal{B}^*$ in (4.5) and the definition of $\mathcal{B}$ in (2.41), an integration by parts in time yields

$$\mathcal{B}^*(v, u) = \mathcal{B}(u, v) = -\int_0^T \langle \dot{v} ; u \rangle_{\mathcal{H}} dt + \langle v(T) ; u(T) \rangle_{\mathcal{H}} - \langle v(0) ; u(0) \rangle_{\mathcal{H}}$$

$$+ \int_0^T a(v_2; u_1) dt - \int_0^T a(v_1 - \varepsilon v_2; u_2) dt. \qquad \square$$

The discrete weak dual problem reads: Given $\Psi(T) \in \mathcal{S} \times \mathcal{S} \subset \mathcal{H}$, find $\Psi \in \mathcal{Q}_c$ such that

$$\mathcal{B}^*(\Psi, V) + \langle \Psi(0) ; V(0) \rangle_{\mathcal{H}} = \langle \Psi(T) ; V(T) \rangle_{\mathcal{H}} \quad \text{for all} \quad V \in \mathcal{W}_c. \tag{4.35}$$

**Remark 4.2.2.1.** Note that the strong dual solution $\Phi$ of problem (4.3) is also a weak solution of problem (4.35). □

**Lemma 4.2.2.2 (Stability of the discrete weak dual $cG(1)$ solution).** The solution $\Psi = (\Psi_1, \Psi_2)$ of the problem (4.35) satisfies the following stability estimates

$$\|\Psi(t_n)\|_{\mathcal{H}}^2 + 2\varepsilon \int_{t_n}^T \|\dot{\Psi}_1(t)\|_{H^1(\Omega)}^2 dt = \|\Psi(T)\|_{\mathcal{H}}^2, \tag{4.36a}$$

where $1 \leq n \leq N$ and

$$\|\dot{\Psi}_1\|_{L^\infty(L^2)}^2 \leq \|\Psi(T)\|_{\mathcal{H}}^2. \tag{4.36b}$$

Furthermore, for $\varepsilon = 0$, there holds

$$\|\dot{\Psi}_2\|_{L^\infty(H^{-1})}^2 \leq \|\Psi(T)\|_{\mathcal{H}}^2. \tag{4.36c}$$

**Proof.** The proof follows analogously as in case of weak "forward problem", see Lemma 2.3.4.1. $\qquad\square$

**Theorem 4.2.2.1 (A priori dual error estimate, $cG(1)$ time approximation).** There is a constant $C$, independent of $u$ and its discrete $cG(1)$ counterpart $U$, such that for $\varepsilon \geq 0$

$$\|e(T)\|_{\mathcal{H}} \leq C \Big\{ \|h^{p+1} D^{p+1} \ddot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} \dot{y}_1\|_{L^2(\Omega)}$$
$$+ \|k^2 \Delta \ddot{u}_2\|_{L^1(L^2)} + \big( \|k\ddot{u}_2\|_{L^1(H^1)} \|k\Delta \ddot{u}_2\|_{L^1(L^2)} \big) + \|k\Delta \ddot{u}_1\|_{L^1(L^2)}$$
$$+ \|k\Delta(\ddot{u}_1 + \varepsilon u_2)\|_{L^1(L^2)} + \|h^p D^{p+1} u_1(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} \Big\}, \tag{4.37a}$$

provided $u \in (H^{p+1}(\Omega))^2$ where $p = 1$ for linear and $p = 3$ for cubic elements in space. Here the 5th does not appear if the (DD*) boundary conditions hold. Moreover, if $\varepsilon = 0$, there holds additionally

$$\|e(T)\|_{\mathcal{H}} \leq C \Big\{ \|h^{p+1} D^{p+1} \ddot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} \dot{y}_1\|_{L^2(\Omega)}$$
$$+ \|k^2 \Delta \ddot{u}_2\|_{L^1(L^2)} + \|k^2 \Delta(\ddot{u}_1 + \varepsilon \ddot{u}_2)\|_{L^1(H^1)}$$
$$+ \|h^p D^{p+1} u_1(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} \Big\}, \tag{4.37b}$$

provided $u \in (H^{p+2}(\Omega))^2$ for $p = 1$ and $u \in (H^{p+2}(\Omega))^2$ for $p = 3$.

**Remark 4.2.2.2.** The estimate (4.37a) is of order $\mathcal{O}(h^p + k)$ whereby the estimate (4.37b) attains the expected order of convergence $\mathcal{O}(h^p + k^2)$. $\qquad\square$

**Proof.** Let the error $e = u - U$ be decomposed as

$$e = \rho - \theta, \quad \rho := u - \mathcal{J}_1 \Pi u, \quad \theta := U - \mathcal{J}_1 \Pi u, \tag{4.38}$$

where $\Pi := (\mathcal{G}, \mathcal{G})$ and $\mathcal{J}_1$ is the $H^1$ temporal projection onto the space of $cG(1)$ functions from Definition 3.1.0.10. For the properties of the projections $\mathcal{J}_1$ and $\mathcal{G}$ we recall Lemma 3.1.0.12 and 3.1.0.6, respectively. To bound $\|\theta(T)\|_{\mathcal{H}}$, we introduce the following Lemma.

**Lemma 4.2.2.3.** For $\theta, \rho$ as in (4.38) there holds

$$\|\theta(T)\|_{\mathcal{H}}^2 = a(\theta_1(0); \Psi_2(0)) - a(\theta_2(0); \Psi_1(0)) + \varepsilon a(\theta_2(0); \Psi_2(0)) + \int_0^T a(\dot{\rho}_1; \dot{\Psi}_1) dt + \int_0^T (\dot{\rho}_2; \dot{\Psi}_2) dt$$

$$- \int_0^T a(\rho_2; \dot{\Psi}_1) dt + \int_0^T a(\rho_1; \dot{\Psi}_2) dt + \varepsilon \int_0^T a(\rho_2; \dot{\Psi}_2) dt =: \sum_{\ell=1}^8 E_\ell, \tag{4.39}$$

where $\Psi(T)$ is chosen such that $\Psi(T) = (\varepsilon \theta_1(T) - \mathcal{K}_h \theta_2(T), \theta_1(T))$.

**Proof.** Given (4.35), let the test function $V \in \mathcal{W}_c$ be chosen such that $V = \dot{\theta}$. This yields

$$\langle \Psi(T) ; \dot{\theta}(T) \rangle_{\mathcal{H}} - \langle \Psi(0) ; \dot{\theta}(0) \rangle_{\mathcal{H}} = \mathcal{B}^*(\Psi, \dot{\theta})$$
$$= \int_0^T a(\Psi_2 - \dot{\Psi}_1; \dot{\theta}_1) dt - \int_0^T (\dot{\Psi}_2; \dot{\theta}_2) dt - \int_0^T a(\Psi_1; \dot{\theta}_2) dt$$
$$+ \varepsilon \int_0^T a(\Psi_2; \dot{\theta}_2) dt + \langle \Psi(T) ; \dot{\theta}(T) \rangle_{\mathcal{H}} - \langle \Psi(0) ; \dot{\theta}(0) \rangle_{\mathcal{H}}.$$

The last equation is equivalent to

$$0 = - \int_0^T a(\dot{\theta}_1; \dot{\Psi}_1) dt - \int_0^T (\dot{\theta}_2; \dot{\Psi}_2) dt - \int_0^T a(\dot{\theta}_2; \Psi_1) dt + \int_0^T a(\dot{\theta}_1; \Psi_2) dt + \varepsilon \int_0^T a(\dot{\theta}_2; \Psi_2) dt.$$

Furthermore, an integration by parts in time yields

$$0 = - \int_0^T a(\dot{\theta}_1; \dot{\Psi}_1) dt - \int_0^T (\dot{\theta}_2; \dot{\Psi}_2) dt + \int_0^T a(\theta_2; \dot{\Psi}_1) dt - a(\theta_2(T); \Psi_1(T)) + a(\theta_2(0); \Psi_1(0))$$
$$- \int_0^T a(\theta_1; \dot{\Psi}_2) dt + a(\theta_1(T); \Psi_2(T)) - a(\theta_1(0); \Psi_2(0))$$
$$- \varepsilon \int_0^T a(\theta_2; \dot{\Psi}_2) dt + \varepsilon a(\theta_2(T); \Psi_2(T)) - \varepsilon a(\theta_2(0); \Psi_2(0)).$$

From the definition of the bilinear form $\mathcal{B}$, see (2.41), the last equation simplifies to

$$a(\theta_1(T); \Psi_2(T)) - a(\theta_2(T); (\Psi_1 - \varepsilon \Psi_2)(T)) = \mathcal{B}(\theta, \dot{\Psi}) + a(\theta_1(0); \Psi_2(0))$$
$$- a(\theta_2(0); (\Psi_1 - \varepsilon \Psi_2)(0)). \qquad (4.40)$$

The Galerkin orthogonality (2.9) provides

$$\mathcal{B}(\theta, \dot{\Psi}) = \mathcal{B}(U - \mathcal{J}_1 \Pi u, \dot{\Psi}) = \mathcal{B}(u - \mathcal{J}_1 \Pi u, \dot{\Psi}) = \mathcal{B}(\rho, \dot{\Psi}). \qquad (4.41)$$

Furthermore, due to the choice of $\Psi$, we obtain

$$a(\theta_1(T); \Psi_2(T)) - a(\theta_2(T); (\Psi_1 - \varepsilon \Psi_2)(T)) = a(\theta_1(T); \theta_1(T)) + a(\theta_2(T); \mathcal{K}_h \theta_2(T))$$
$$= \|\theta_1(T)\|_{H^1(\Omega)}^2 + \|\theta_2(T)\|_{L^2(\Omega)}^2$$
$$= \|\theta(T)\|_{\mathcal{H}}^2. \qquad (4.42)$$

If we substitute (4.41) and (4.42) into (4.40) we obtain

$$\|\theta(T)\|_{\mathcal{H}}^2 = \mathcal{B}(\rho, \dot{\Psi}) + a(\theta_1(0); \Psi_2(0)) - a(\theta_2(0); (\Psi_1 - \varepsilon \Psi_2)(0)).$$

From the definition of the bilinear form $\mathcal{B}$, cf. (2.41), we may conclude the proof. $\qquad \square$

In order to estimate $\|\theta(T)\|_{\mathcal{H}}^2$, we have to estimate each of $E_\ell$, $\ell = 1, \ldots, 8$.

**Lemma 4.2.2.4.** Let $U(0) := \Pi u_0$ be a discrete variant of the initial solution $u^0$, then

$$E_1 + E_2 + E_3 = 0.$$

**Proof.** Owing to the properties of the projection $\mathcal{G}, \mathcal{J}_1$ we have

$$\theta(0) = (U - \mathcal{J}_1 \Pi u_0)(0) = U(0) - \Pi u_0 = \Pi u_0 - \Pi u_0 = 0. \tag{4.43}$$

$\square$

**Remark 4.2.2.3.** If we choose $U(0) = (\mathcal{I}y_0, \mathcal{I}y_1)$, we may obtain the estimate in terms of $\|\Psi_2(0)\|_{H^1(\Omega)}$. This is not optimal due to the results from the stability Lemma 4.2.2.2. $\square$

**Lemma 4.2.2.5.** $E_4 = 0$.

**Proof.** According to the orthogonal properties of projection $\mathcal{G}$ and $\mathcal{J}_1$

$$E_4 = \int_0^T a(\dot{u}_1 - \mathcal{G}\dot{u}_1; \dot{\Psi}_1)dt + \int_0^T a(\frac{\partial}{\partial t}(\mathcal{G}u_1 - \mathcal{J}_1\mathcal{G}u_1); \dot{\Psi}_1)dt = 0. \qquad \square$$

**Lemma 4.2.2.6.** There exists a constant $C$ such that for $\varepsilon \geq 0$

$$\begin{aligned}
E_5 \leq C\big\{ &\|h^{p+1}D^{p+1}\ddot{u}_2\|_{L^1(L^2)}\|\Psi_2\|_{L^\infty(L^2)} + \|h^{p+1}D^{p+1}\dot{u}_2(T)\|_{L^2(\Omega)}\|\theta_2(T)\|_{L^2(\Omega)} \\
&+ \|h^{p+1}D^{p+1}\dot{y}_1\|_{L^2(\Omega)}\|\Psi_2\|_{L^\infty(L^2)} \big\}.
\end{aligned} \tag{4.44}$$

**Proof.** Similar as in the proof of $E_4$, the orthogonality properties of $\mathcal{J}_1$ imply

$$E_5 = \int_0^T (\dot{u}_2 - \mathcal{G}\dot{u}_2; \dot{\Psi}_2)dt + \int_0^T (\frac{\partial}{\partial t}(\mathcal{G}u_2 - \mathcal{J}_1\mathcal{G}u_2); \dot{\Psi}_2)dt = \int_0^T (\dot{u}_2 - \mathcal{G}\dot{u}_2; \dot{\Psi}_2)dt.$$

An integration by parts in time and the approximation properties of the Galerkin projections yield

$$\begin{aligned}
E_5 = &-\int_0^T (\ddot{u}_2 - \mathcal{G}\ddot{u}_2; \Psi_2)dt + ((\dot{u}_2 - \mathcal{G}\dot{u}_2)(T); \Psi_2(T)) - (\dot{y}_1 - \mathcal{G}\dot{y}_1; \Psi_2(0)) \\
\leq &\|h^{p+1}D^{p+1}\ddot{u}_2\|_{L^1(L^2)}\|\Psi_2\|_{L^\infty(L^2)} + \|h^{p+1}D^{p+1}\dot{u}_2(T)\|_{L^2(\Omega)}\|\theta_2(T)\|_{L^2(\Omega)} \\
&+ \|h^{p+1}D^{p+1}\dot{y}_1\|_{L^2(\Omega)}\|\Psi_2\|_{L^\infty(L^2)}. \qquad\qquad\qquad\qquad\qquad \square
\end{aligned}$$

**Lemma 4.2.2.7.** There exists a constant $C$ such that

$$E_6 \leq C\Big( \|k^2\Delta\ddot{u}_2\|_{L^1(L^2)}\|\dot{\Psi}_1\|_{L^\infty(L^2)} + \big(\|k\ddot{u}_2\|_{L^1(H^1)} + \|k\Delta\ddot{u}_2\|_{L^1(L^2)}\big)\|\phi_1\|_{L^\infty(H^1)} \Big) \tag{4.45}$$

where the second summand does not appear if the (DD*) boundary conditions or $\varepsilon = 0$.

**Proof.** The orthogonality properties of projection $\mathcal{G}$ lead to

$$E_6 = -\int_0^T a(u_2 - \mathcal{G}u_2; \dot{\Psi}_1)dt - \int_0^T a(\mathcal{G}u_2 - \mathcal{J}_1\mathcal{G}u_2; \dot{\Psi}_1)dt = -\int_0^T a(u_2 - \mathcal{J}_1 u_2; \dot{\Psi}_1)dt.$$

Furthermore, an integration by parts in space yields

$$E_6 = \int_0^T (\Delta(u_2 - \mathcal{J}_1 u_2); \dot{\Psi}_1)dt - \int_0^T D(u_2 - \mathcal{J}_1 u_2)(t, 1)\dot{\Psi}_1(t, 1)dt. \tag{4.46}$$

For the first term on the RHS we have

$$\int_0^T (\Delta(u_2 - \mathcal{J}_1 u_2); \dot{\Psi}_1) dt \leq C \|k^2 \Delta \ddot{u}_2\|_{L^1(L^2)} \|\dot{\Psi}_1\|_{L^\infty(L^2)}. \tag{4.47}$$

The second term equals zero if (DD) or $\varepsilon = 0$ and (DN) hold. Otherwise, if we integrate by parts in time and make use of the fact that $\mathcal{J} u_2$ and $u_2$ coincide in each time-point $t_j$, then

$$-\int_0^T D(u_2 - \mathcal{J}_1 u_2)(t, 1) \dot{\Psi}_1(t, 1) dt = \int_0^T \frac{\partial}{\partial t}(D(u_2 - \mathcal{J}_1 u_2))(t, 1) \Psi_1(t, 1) dt.$$

From

$$\Psi_1(t, 1) \leq \|\Psi_1(t)\|_{H^1(\Omega)},$$
$$\frac{\partial}{\partial t}(D(u_2 - \mathcal{J}_1 u_2))(t, 1) \leq \|\frac{\partial}{\partial t}(u_2 - \mathcal{J}_1 u_2)(t)\|_{H^1(\Omega)} + \|\frac{\partial}{\partial t}\Delta(u_2 - \mathcal{J}_1 u_2)(t)\|_{L^2(\Omega)},$$

we obtain for some numerical constant $C$,

$$-\int_0^T D(u_2 - \mathcal{J}_1 u_2)(t, 1) \dot{\Psi}_1(t, 1) dt \leq C\left(\|k\ddot{u}_2\|_{L^1(H^1)} + \|k\Delta\ddot{u}_2\|_{L^1(L^2)}\right) \|\Psi_1\|_{L^\infty(H^1)}. \tag{4.49}$$

A substitution of (4.47) and (4.49) into (4.46) yields the proof of theorem. $\square$

**Lemma 4.2.2.8.** There exists a constant $C$ such that for $\varepsilon \geq 0$

$$E_7 + E_8 \leq C \|k\Delta(\ddot{u}_1 + \varepsilon \ddot{u}_2)\|_{L^1(L^2)} \|\Psi_2\|_{L^\infty(L^2)},$$

Furthermore, for $\varepsilon = 0$ there holds additionally

$$E_7 + E_8 \leq C \|k^2 \Delta(\ddot{u}_1 + \varepsilon \ddot{u}_2)\|_{L^1(H^1)} \|\dot{\Psi}_2\|_{L^\infty(H^{-1})}.$$

**Proof.** The orthogonality property of $\mathcal{G}$ and integration by parts in space yield

$$E_7 + E_8 = \int_0^T a((I - \mathcal{J})(u_1 + \varepsilon u_2); \dot{\Psi}_2) dt = -\int_0^T ((I - \mathcal{J})\Delta(u_1 + \varepsilon u_2); \dot{\Psi}_2) dt. \tag{4.50}$$

Let us assume that $\varepsilon = 0$. Then

$$E_7 + E_8 = -\int_0^T ((I - \mathcal{J})\Delta(u_1 + \varepsilon u_2); \dot{\Psi}_2) dt \leq C \|k^2 \Delta(\ddot{u}_1 + \varepsilon \ddot{u}_2)\|_{L^1(H^1)} \|\dot{\Psi}_2\|_{L^\infty(H^{-1})}.$$

If $\varepsilon > 0$, an integration by parts in time in the RHS of (4.50) and the Hölder inequality yield

$$E_7 + E_8 = -\int_0^T ((I - \mathcal{J})\Delta(u_1 + \varepsilon u_2); \dot{\Psi}_2) dt = \int_0^T (\frac{\partial}{\partial t}(I - \mathcal{J})\Delta(u_1 + \varepsilon u_2); \Psi_2) dt$$
$$\leq C \|k\Delta(\ddot{u}_1 + \varepsilon \ddot{u}_2)\|_{L^1(L^2)} \|\Psi_2\|_{L^\infty(L^2)}. \tag{4.51}$$

This concludes the proof. $\square$

Given the stability Lemma 4.2.2.2, if we substitute the estimates for $E_1, \ldots, E_8$ derived above into the representation (4.39), we have for $\varepsilon \geq 0$,

$$
\begin{aligned}
\|\theta(T)\|_{\mathcal{H}} \leq C \Big\{ &\|h^{p+1} D^{p+1} \ddot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} \dot{y}_1\|_{L^2(\Omega)} \\
&+ \|k^2 \Delta \ddot{u}_2\|_{L^1(L^2)} + \big( \|k \ddot{u}_2\|_{L^1(H^1)} \|k \Delta \ddot{u}_2\|_{L^1(L^2)} \big) + \|k \Delta \ddot{u}_1\|_{L^1(L^2)} \\
&+ \|k \Delta (\ddot{u}_1 + \varepsilon \ddot{u}_1)\|_{L^1(L^2)} \Big\}.
\end{aligned}
\tag{4.52a}
$$

If $\varepsilon = 0$, then

$$
\begin{aligned}
\|\theta(T)\|_{\mathcal{H}} \leq C \Big\{ &\|h^{p+1} D^{p+1} \ddot{u}_2\|_{L^1(L^2)} + \|h^{p+1} D^{p+1} \dot{u}_2(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} \dot{y}_1\|_{L^2(\Omega)} \\
&+ \|k^2 \Delta \ddot{u}_2\|_{L^1(L^2)} + \|k^2 \Delta (\ddot{u}_1 + \varepsilon \ddot{u}_2)\|_{L^1(H^1)} \Big\}.
\end{aligned}
\tag{4.52b}
$$

Its left to estimate $\|\rho(T)\|_{\mathcal{H}}$. We recall the proof of Lemma 4.2.1.11. Then there holds

$$
\|\rho(T)\|_{\mathcal{H}} \leq C \big( \|h^p D^{p+1} u_1(T)\|_{L^2(\Omega)} + \|h^{p+1} D^{p+1} u_2(T)\|_{L^2(\Omega)} \big),
\tag{4.53}
$$

for some numerical constant $C > 0$.

According to the decomposition (4.38), estimate (4.52), and (4.53) yield the proof.  $\square$

## 4.3   A posteriori dual error estimate

The a posteriori error estimates accomplish two main goals. First they provide a computable error bound for the given finite element computation. Secondly, they are used to perform the adaptive mesh refinement. Within this section we analyse and derive these bounds by using the dual method techniques which rely on the stability estimates of the strong dual solution introduced in Lemma 4.1.0.5. Their use in adaptive refinement process will be emphasised in Chapter 6.4.

We start by first considering the time discretisation methods, and then combine them with the two different space ansatz, i.e. $\mathcal{P}_1$ and $\mathcal{C}^1$ elements. The proven convergence order of a posteriori error estimates for different time discretisation methods and $\mathcal{P}_1$ and $\mathcal{C}^1$ ansatz in space are given in Table 4.3.

Note that the error convergence rates proved by the dual method under certain restrictions show the better convergence behaviour than the ones obtained by the energy method. This is obvious when compared $dG(1) \otimes \mathcal{C}^1$ and $cG(1) \otimes \mathcal{C}^1$ discretisation. However, we still did not succeed to prove the optimal convergence order for $\mathcal{P}^1$ space ansatz. The problem is that in 1D we can not make any additional requirements on the spatial mesh.

Here we also proved the a posteriori error estimates in the negative norm.

| $\|e_1^{N-}\|_{H^{1-s}}+\|e_2^{N-}\|_{H^{-s}}$ if $s=1$ then $\varepsilon=0$ or (DD) | | $\mathcal{P}^1$ | $\mathcal{C}^1$ |
|---|---|---|---|
| **$dG(0)$** Subsection 4.3.1.1 | $s=0$ | – | $\mathcal{O}(h^3+k)$, (DD*), $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, $\mathcal{O}(h^3+k^{-1}h^3+k)$ otherwise |
| | $s=1$ | $\mathcal{O}(h+k)$ $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, $\mathcal{O}(h+k^{-1}h^2+k)$ otherwise | $\mathcal{O}(h^4+k)$, $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$, $\mathcal{O}(h^4+k^{-1}h^4+k)$ otherwise |
| **$dG(1)$** Subsection 4.3.1.2 | $s=0$ | – | $\mathcal{O}(h^3+k^3)$, $\varepsilon=0$, (DD*), $\mathcal{S}^{j-1}=\mathcal{S}^j$ $\mathcal{O}(h^3+k^{-1}h^3+k^3)$ otherwise |
| | $s=1$ | $\mathcal{O}(h+k^3)$, $\varepsilon=0$, $\mathcal{S}^{j-1}=\mathcal{S}^j$, $\mathcal{O}(h+k^{-1}h+k^3)$ otherwise | $\mathcal{O}(h^4+k^3)$, (DD*), $\mathcal{S}^{j-1}=\mathcal{S}^j$ $\mathcal{O}(h^4+k^{-1}h^4+k^3)$ otherwise |
| **$cG(1)$** Subsection 4.3.2 | $s=0$ | – | $\mathcal{O}(h^{s+3}+k^2)$, (DD*) |
| | $s=1$ | $\mathcal{O}(h+k^2)$ | |

Table 4.3: Proven a posteriori error estimates for $\|e^{N-}\|_{\mathcal{H}}$ ($s=0$) and $\|e^{N-}\|_{\widehat{\mathcal{H}}}$ ($s=1$); dual method.

## 4.3.1  $dG(q)$ time approximation, $q=0,1$

In the following, the error representation lemma for the case $dG(q)$ in time will be introduced. The notations and definitions are adopted from Subsection 2.3.3 and Section 2.1. Note also that the analysis presented below employs the residual $Res$ from Definition 2.3.1.2.

**Lemma 4.3.1.1 (Dual error representation, $dG(q)$ time approximation).** If $u$ is a solution of (1.28), $U$ its discrete $dG(q)$ variant and $\Phi$ a strong solution of the corresponding dual problem (4.3), there holds

$$\left\langle\Phi^{N-};e^{N-}\right\rangle_{\mathcal{H}}=\left\langle\Phi^{0+};e^{0-}\right\rangle_{\mathcal{H}}+Res(e,\Phi-V)\quad\text{for all}\quad V\in\mathcal{Q}_q. \tag{4.54}$$

**Proof.** Since $\Phi$ is also a solution of the discrete weak dual problem (4.9), cf. Remark 4.2.1.2, we have

$$\left\langle\Phi^{N-};e^{N-}\right\rangle_{\mathcal{H}}-\left\langle\Phi^{0+};e^{0-}\right\rangle_{\mathcal{H}}=\mathcal{B}^*(\Phi,e). \tag{4.55}$$

From the definition of the dual bilinear form, see Definition 4.1.0.2, the Galerkin orthogonality (2.9) and the definition of the residual $Res$ in (2.10), we have for all $V\in\mathcal{Q}_q$

$$\mathcal{B}^*(\Phi,e)=\mathcal{B}(e,\Phi)=\mathcal{B}(e,\Phi-V)=Res(\Phi-V).$$

A substitution of the last equation into (4.55) completes the proof. $\square$

### 4.3.1.1 Dual a posteriori error analysis, $dG(0) \otimes \mathcal{P}_1$

The following analysis implies the estimation of the error in the energy norm $\|e^{N-}\|_{\mathcal{H}}$ for Problem (1.13) with an arbitrary boundary conditions and of the error $\|e^{N-}\|_{\widehat{\mathcal{H}}}$ in case of (DD*) boundary conditions only.

**Theorem 4.3.1.1 (A posteriori dual error estimate, $dG(0) \otimes \mathcal{P}_1$).** There exists a constant $C$ such that the error of the $dG(0) \otimes \mathcal{P}_1$ approximation satisfies the following a posteriori error bound for $s=1$

$$\|e^{N-}\|_{\widehat{\mathcal{H}}} \leq C \Big\{ \|h(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} + \|h(y_1 - \mathcal{I}y_1)\|_{H^1(\Omega)} + \|kU_2\|_{L^1(H^1)} + \|h(f - \mathcal{L}f)\|_{L^1(L^2)}$$

$$+ (T^{1/2} + \varepsilon^{1/2}) \|kf\|_{L^1(L^2)} + \sum_{j=1}^{N} \|h(I - \mathcal{G})U_1^{j-1-}\|_{H^1(\Omega)} + \sum_{j=1}^{N} \|h(I - \mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}$$

$$+ (T^{1/2} + \varepsilon^{1/2}) \|k\mathcal{K}_h^{-1}(U_1 + \varepsilon U_2)\|_{L^1(L^2)} \Big\}. \tag{4.56}$$

The sums $\sum_{j=1}^{N} \|(I - \mathcal{G})U_1^{j-1-}\|_{H^1(\Omega)} + \sum_{j=1}^{N} \|(I - \mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}$ in equation above vanish provided $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j$.

**Remark 4.3.1.1.** The estimate (4.56) is of order $\mathcal{O}(h+k)$ when $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j = 1, \ldots, N$. Otherwise, the estimate is of order $\mathcal{O}(h + h^2 k^{-1} + k)$. $\qquad\square$

**Proof.** Given the residual representation from Lemma 3.3.2.2, case $\mathcal{P}^1$, the error represenation (4.54) is equivalent to

$$\big\langle e^{N-} ; \Phi^{N-} \big\rangle_{\mathcal{H}} = \big\langle e^{0-} ; \Phi^{0+} \big\rangle_{\mathcal{H}} + \sum_{j=1}^{N} \int_{I_j} a(U_2; \phi - V_1) dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} (f; \dot{\phi} - V_2) dt - \sum_{j=1}^{N} \big\langle [U]^{j-1} ; (\Phi - V)^{j-1+} \big\rangle_{\mathcal{H}}$$

$$- \sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)(t)]_k (\dot{\phi} - V_2)(t, x_k) dt =: \sum_{\ell=1}^{5} E_\ell. \tag{4.57}$$

Here $m = n-1$ for (DD*) or $m = n$ for (DN*) boundary conditions.

In order to determine the upper bound for $\|e^{N-}\|_{\widehat{\mathcal{H}}}$, we need to estimate each of the terms $E_1, \ldots, E_5$ from (4.57). Therefore, we choose a test function $V \in \mathcal{Q}_0$ such that

$$V = \mathcal{J}\Pi\Phi, \quad \Pi = (\mathcal{G}, \mathcal{L}). \tag{4.58}$$

Here, $\Phi$ denotes the continuous solution of the dual problem (4.1) and $\mathcal{J}$ is the mapping onto the space of constant functions in time defined by

$$\mathcal{J}u|_{I_j} := u^{j-1+} \quad \text{for all} \quad j = 1, \ldots, N \quad \text{and} \quad u \in H^1(\mathscr{T}). \tag{4.59}$$

Obviously, $\mathcal{J}$ is not a $L^2$ projection in time.

**Lemma 4.3.1.2.** If a discrete variant $U^{0-}$ of the initial solution $u_0 = (y_0, y_1)$ is defined by

$$U^{0-} := (\mathcal{I}y_0, \mathcal{I}y_1),$$

where $\mathcal{I}$ is the nodal interpolation operator, then

$$E_1 \leq C \Big\{ \|h(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} \|\Delta\phi\|_{L^\infty(L^2)} + \|h(y_1 - \mathcal{I}y_1)\|_{H^1(\Omega)} \|\dot{\phi}\|_{L^\infty(H^1)} \Big\}.$$

**Proof.** We start from

$$E_1 = a(y_0 - \mathcal{I}y_0; \phi^{0+}) + (y_1 - \mathcal{I}y_1; \dot{\phi}^{0+}).$$

The Friedrichs inequality and the fact that the nodal interpolation operator and the Galerkin projection $\mathcal{G}$ coincide in $1D$, yield the following estimate

$$E_1 \leq C \Big\{ \|h(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} \|\Delta\phi\|_{L^\infty(L^2)} + \|h(y_1 - \mathcal{I}y_1)\|_{H^1(\Omega)} \|\dot{\phi}\|_{L^\infty(H^1)} \Big\}. \qquad \square$$

**Lemma 4.3.1.3.** There holds

$$E_2 \leq \|kU_2\|_{L^1(H^1)} \|\dot{\phi}\|_{L^\infty(H^1)}.$$

**Proof.** The properties of $\mathcal{G}$ imply

$$E_2 = \sum_{j=1}^{N} \int_{I_j} a(U_2; \phi - \mathcal{G}\phi) dt + \sum_{j=1}^{N} \int_{I_j} a(U_2; \mathcal{G}\phi - \mathcal{G}\phi^{j-1+}) dt = \sum_{j=1}^{N} \int_{I_j} a(U_2; \phi - \phi^{j-1+}) dt.$$

Since $\phi(t) - \phi(t_{j-1}) = \int_{t_{j-1}}^{t} \dot{\phi}(\tau) d\tau$ we may deduce

$$E_2 \leq \sum_{j=1}^{N} \int_{I_j} \|U_2\|_{H^1(\Omega)} \|\phi - \phi^{j-1+}\|_{H^1(\Omega)} dt \leq \|kU_2\|_{L^1(H^1)} \|\dot{\phi}\|_{L^\infty(H^1)}.$$

Note that the inequality above can be further used only if the (DD*) boundary conditions or $\varepsilon = 0$ hold. $\qquad \square$

**Lemma 4.3.1.4.** There exists a constant $C > 0$ such that there holds

$$E_3 \leq C \big\{ \|h(f - \mathcal{L}f)\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^1)} + \|kf\|_{L^2(L^2)} \big( T^{1/2} \|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon \|\Delta\dot{\phi}\|_{L^2(L^2)} \big) \big\}.$$

**Proof.** According to the symmetry properties of $\mathcal{L}$, we have

$$E_3 = \sum_{j=1}^{N} \int_{I_j} (f - \mathcal{L}f; \dot{\phi} - \mathcal{L}\dot{\phi}) dt + \sum_{j=1}^{N} \int_{I_j} (f; \mathcal{L}\dot{\phi} - \mathcal{L}\dot{\phi}^{j-1+}) dt.$$

For the first term there holds

$$\sum_{j=1}^{N} \int_{I_j} (f - \mathcal{L}f; \dot{\phi} - \mathcal{L}\dot{\phi}) dt \leq C \|h(f - \mathcal{L}f)\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^1)}.$$

Owing to $\ddot{\phi} = \Delta(\phi - \varepsilon\dot{\phi})$, we deduce the following estimate for the second term

$$\sum_{j=1}^{N} \int_{I_j} (f; \mathcal{L}(\dot{\phi} - \dot{\phi}^{j-1+})) dt = \sum_{j=1}^{N} \int_{I_j} (\mathcal{L}f; \dot{\phi} - \dot{\phi}^{j-1+}) dt$$

$$\leq \big( \|k\mathcal{L}f\|_{L^2(0,T)}; T^{1/2} \|\Delta\phi\|_{L^\infty(0,T)} + \varepsilon \|\Delta\dot{\phi}\|_{L^2(0,T)} \big). \qquad (4.60)$$

This follows since

$$\sum_{j=1}^{N}\int_{I_j}(\mathcal{L}f;\dot\phi-\dot\phi^{j-1+})dt=\int_{\Omega}\sum_{j=1}^{N}\int_{I_j}(\mathcal{L}f)(\dot\phi-\dot\phi^{j-1+})dtdx$$

$$\leq\int_{\Omega}\sum_{j=1}^{N}\|k_j\mathcal{L}f\|_{L^2(I_j)}\|\Delta(\phi-\varepsilon\dot\phi)\|_{L^2(I_j)}dx$$

$$\leq\int_{\Omega}\sum_{j=1}^{N}\|k_j\mathcal{L}f\|_{L^2(I_j)}\big(\|\Delta\phi\|_{L^2(I_j)}+\varepsilon\|\Delta\dot\phi\|_{L^2(I_j)}\big)dx$$

$$\leq\int_{\Omega}\sum_{j=1}^{N}\|k_j\mathcal{L}f\|_{L^2(I_j)}\big(k_j^{1/2}\|\Delta\phi\|_{L^\infty(I_j)}+\varepsilon\|\Delta\dot\phi\|_{L^2(I_j)}\big)dx$$

$$\leq\int_{\Omega}\Big\{\|k\mathcal{L}f\|^2_{L^2(0,T)}\Big(\Big(\sum_{j=1}^{N}k_j\Big)^{1/2}\|\Delta\phi\|_{L^\infty(0,T)}+\varepsilon\|\Delta\dot\phi\|_{L^2(0,T)}\Big\}dx.$$

By use of the Hölder inequality, the RHS of the equation above can be dominated by

$$\sum_{j=1}^{N}\int_{I_j}(f;\mathcal{L}(\dot\phi-\dot\phi^{j-1+}))dt\leq\|kf\|_{L^2(L^2)}\big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)}+\varepsilon\|\Delta\dot\phi\|_{L^2(L^2)}\big).\qquad\square$$

**Lemma 4.3.1.5.** There exists a constant $C$ such that

$$E_4\leq C\Big(\sum_{j=1}^{N}\|h(I-\mathcal{G})U_1^{j-1-}\|_{H^1(\Omega)}\Big)\|\Delta\phi\|_{L^\infty(L^2)}+C\Big(\sum_{j=1}^{N}\|h(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big)\|\dot\phi\|_{L^\infty(H^1)},$$

where both sums go only over such $j$ where $\mathcal{S}^{j-1}\not\subseteq\mathcal{S}^{j}$.

**Proof.** Owing to the properties of the mapping $\mathcal{J}$, we have

$$E_4=-\sum_{j=1}^{N}a([U_1]^{j-1};(\phi-\mathcal{G}\phi)^{j-1+})dt-\sum_{j=1}^{N}([U_2]^{j-1};(\dot\phi-\mathcal{L}\dot\phi)^{j-1+}).$$

If we assume that for all $j=1,\ldots,N$, $\mathcal{S}^{j-1}\subseteq\mathcal{S}^{j}$, then $E_4=0$ owing to the orthogonality properties of $\mathcal{G}$ and $\mathcal{L}$. Otherwise,

$$E_4=\sum_{j=1}^{N}a(U_1^{j-1-}-\mathcal{G}U_1^{j-1-};(\phi-\mathcal{G}\phi)^{j-1+})dt+\sum_{j=1}^{N}(U_2^{j-1-}-\mathcal{L}U_2^{j-1-};(\dot\phi-\mathcal{L}\dot\phi)^{j-1+})$$

$$\leq C\Big(\sum_{j=1}^{N}\|h(I-\mathcal{G})U_1^{j-1-}\|_{H^1(\Omega)}\Big)\|\Delta\phi\|_{L^\infty(L^2)}+C\Big(\sum_{j=1}^{N}\|h(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big)\|\dot\phi\|_{L^\infty(H^1)}.\ \square$$

**Lemma 4.3.1.6.** There exists a constant $C$ such that there holds

$$E_5\leq C\Big(\sum_{j=1}^{N}\|h(U_1+\varepsilon U_2)^j\|_{H^1(\Omega)}\Big)\|\Delta\phi\|_{L^\infty(L^2)}$$

$$+C\|k\mathcal{K}_h^{-1}(U_1+\varepsilon U_2)\|_{L^1(L^2)}\big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)}+\varepsilon\|\Delta\dot\phi\|_{L^2(L^2)}\big).$$

**Proof.** We start from

$$E_5 = -\sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1+\varepsilon U_2)]_k (\dot\phi - \mathcal{L}\dot\phi)(x_k) dt$$

$$-\sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1+\varepsilon U_2)]_k (\mathcal{L}(\dot\phi - \dot\phi^{j-1+}))(x_k) dt. \tag{4.61}$$

For the first term on the RHS of the equation above we have

$$-\sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1+\varepsilon U_2)]_k (\dot\phi - \mathcal{L}\dot\phi)(x_k) dt = \sum_{j=1}^{N} \int_{I_j} a(U_1+\varepsilon U_2; \dot\phi - \mathcal{L}\dot\phi) dt. \tag{4.62}$$

An integration by parts in time yields

$$\sum_{j=1}^{N} \int_{I_j} a(U_1+\varepsilon U_2; \dot\phi - \mathcal{L}\dot\phi) dt = \sum_{j=1}^{N} a((U_1+\varepsilon U_2)^{j-}; (\phi - \mathcal{L}\phi)^{j-})$$

$$- \sum_{j=1}^{N} a((U_1+\varepsilon U_2)^{j-1+}; (\phi - \mathcal{L}\phi)^{j-1+}).$$

The Hölder inequality and the approximation properties of $\mathcal{L}$ yield the following estimate for the first term

$$\sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1+\varepsilon U_2)]_k (\dot\phi - \mathcal{L}\dot\phi)(x_k) dt \leq C \Big( \sum_{j=1}^{N} \|h(U_1+\varepsilon U_2)^j\|_{H^1(\Omega)} \Big) \|\Delta\phi\|_{L^\infty(L^2)}.$$

Arguing in the similar way, the second term can be rewritten such that

$$-\sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1+\varepsilon U_2)]_k (\mathcal{L}(\dot\phi - \dot\phi^{j-1+}))(x_k) dt = \sum_{j=1}^{N} \int_{I_j} a(U_1+\varepsilon U_2; \mathcal{L}(\dot\phi - \dot\phi^{j-1+})) dt. \tag{4.63}$$

Since $\ddot\phi = \Delta(\phi - \varepsilon\dot\phi)$, we may deduce

$$\sum_{j=1}^{N} \int_{I_j} a(U_1+\varepsilon U_2; \mathcal{L}(\dot\phi - \dot\phi^{j-1+})) dt \leq a\big(\|k(U_1+\varepsilon U_2)\|_{L^1(0,T)}; T^{1/2}\|\mathcal{L}\Delta\phi\|_{L^\infty(0,T)} + \varepsilon\|\mathcal{L}\Delta\dot\phi\|_{L^2(0,T)}\big).$$

From the definition of the operator $\mathcal{K}_h^{-1}$, we have

$$a\big(\|k(U_1+\varepsilon U_2)\|_{L^1(0,T)}; T^{1/2}\|\mathcal{L}\Delta\phi\|_{L^\infty(0,T)} + \varepsilon\|\mathcal{L}\Delta\dot\phi\|_{L^2(0,T)}\big)$$
$$= \big(\mathcal{K}_h^{-1}\|k(U_1+\varepsilon U_2)\|_{L^1(0,T)}; T^{1/2}\|\mathcal{L}\Delta\phi\|_{L^\infty(0,T)} + \varepsilon\|\mathcal{L}\Delta\dot\phi\|_{L^2(0,T)}\big). \tag{4.64}$$

Finally, we may deduce

$$-\sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1+\varepsilon U_2)]_k (\mathcal{L}(\dot\phi - \dot\phi^{j-1+}))(x_k) dt \leq \|k\mathcal{K}_h^{-1}U_1\|_{L^1(L^2)} \big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon\|\Delta\dot\phi\|_{L^2(L^2)}\big).$$

If $s=0$, an integration by parts in space in the RHS of (4.64) yields This concludes the proof. $\square$

Owing to the results of the latter lemmas, if we recall the stability Lemma 4.1.0.5 and the error representation (4.57) we may conclude the proof of theorem. $\square$

**Remark 4.3.1.2.** We may also use $V = (\mathcal{J}\mathcal{G}\phi, \mathcal{J}\mathcal{I}\dot\phi)$ in order to neutralise the jumps terms in $E_5$. However, since $\mathcal{I}$ is not $L^2$ orthogonal, we can not obtain the optimal estimate in e.g. $E_3$. $\square$

**4.3.1.2 Dual a posteriori error analysis, *dG*(0)⊗*C*₁**

For the derivation of the a posteriori error estimate in case of $dG(0) \otimes \mathcal{C}^1$ approximation, recall the definition of the bilinear form $\mathcal{B}$ (2.18), its dual form $\mathcal{B}^*$ (4.6) and the error representation formula from Lemma 4.3.1.1.

**Theorem 4.3.1.2 (A posteriori dual error estimate, *dG*(0)⊗*C*¹).** There exists a constant $C > 0$ such that for $s = 0, 1$, the error of the $dG(0) \otimes \mathcal{C}^1$ approximation satisfies the following a posteriori error bound

$$
\begin{aligned}
\|e_1^{N-}\|_{H^{1-s}(\Omega)} + \|e_2^{N-}\|_{H^{-s}(\Omega)} \leq C \Big\{ & \|h^s(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} + \|h^s(y_1 - \mathcal{I}y_1)\|_{H^s(\Omega)} + \|kD^{2-s}U_2\|_{L^1(L^2)} \\
& + \big(\|U_2\|_{H^1(\Omega)} + \|\Delta U_2\|_{L^2(\Omega)}\big) + \|h^s(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2)\|_{L^1(L^2)} \\
& + (T^{1/2} + \varepsilon^{1/2})\|k\mathcal{L}^{1-s}(f + \Delta(U_1 + \varepsilon U_2))\|_{L^2(H^{1-s})} \\
& + \sum_{j=1}^{N} \|h^s(I - \mathcal{G})U_1^{j-1-}\|_{H^1(\Omega)} + \sum_{j=1}^{N} \|h^s(I - \mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}^2 \Big\}. (4.65)
\end{aligned}
$$

Moreover, for $s = 1$ we assume that either $\varepsilon = 0$ or (DD*) hold. Furthermore, the 4th summand on the RHS of the equation above appears only for $s = 0$ and (DN)* boundary conditions. The sums over $j$ vanish for all $j$ where $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$.

**Remark 4.3.1.3.** Assume $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j = 1, \ldots, N$. If $s = 0$ and (DD)*, the estimate (4.65) is of order $\mathcal{O}(h^3 + k)$. If $s = 1$ and either $\varepsilon = 0$ or (DD*) hold, proven estimate has the order of convergence $\mathcal{O}(h^4 + k)$. Otherwise, if we assume that $\mathcal{S}^{j-1} \nsubseteq \mathcal{S}^j$ for all $j = 1, \ldots, N$, then the estimate is of order $\mathcal{O}(h^3 + h^3 k^{-1} + k)$ for $s = 0$ and (DD*) and of order $\mathcal{O}(h^4 + h^4 k^{-1} + k)$ for $s = 1$ where either $\varepsilon = 0$ or (DD)* hold.                                    □

**Proof.** According to the definition of the residual $Res$ from Lemma 3.3.1.2, case $\mathcal{C}^1$, the error representation (4.54) reads

$$
\begin{aligned}
\langle e^{N-} ; \Phi^{N-} \rangle_{\mathcal{H}} = \langle e^{0-} ; \Phi^{0+} \rangle_{\mathcal{H}} & + \sum_{j=1}^{N} \int_{I_j} a(U_2; \phi - V_1)dt + \sum_{j=1}^{N} \int_{I_j} (f + \Delta(U_1 + \varepsilon U_2); \dot{\phi} - V_2)dt \\
& - \sum_{j=1}^{N} \langle [U]^{j-1} ; (\Phi - V)^{j-1+} \rangle_{\mathcal{H}} =: \sum_{\ell=1}^{4} E_\ell. \quad (4.66)
\end{aligned}
$$

In order to derive an a posteriori error bound for $\|e^{N-}\|_{\mathcal{H}}$ and $\|e^{N-}\|_{\hat{\mathcal{H}}}$, we need to estimate $E_1, \ldots, E_4$ such that the final bound consists of some computable terms and terms which can be dominated according to the stability Lemma 4.1.0.5. If we chose a test function $V$ as in (4.72) with $\mathcal{J}$ as in (4.59), then for $E_4$ we may a bounds as in Lemma 4.3.1.5, where $\mathcal{G}$ and $\mathcal{L}$ are the projections onto the space of cubic polynomials in space. It remains to estimate $E_1, E_2$ and $E_3$. This is emphasised in the following lemmas.

**Lemma 4.3.1.7.** If a discrete variant $U^{0-}$ of the initial solution $u_0 = (y_0, y_1)$ is defined by

$$
U^{0-} := (\mathcal{I}y_0, \mathcal{I}y_1),
$$

where $\mathcal{I}$ is the Hermite cubic interpolation operator, then

$$
E_1 \leq C^s \Big\{ \|h^s(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} \|D^{s+1}\phi\|_{L^\infty(L^2)} + \|h^s(y_1 - \mathcal{I}y_1)\|_{H^s(\Omega)} \|\dot{\phi}\|_{L^\infty(H^s)} \Big\}.
$$

**Proof.** We start from the definition

$$E_1 = a(y_0 - \mathcal{I}y_0; \phi^{0+}) + (y_1 - \mathcal{I}y_1; \dot{\phi}^{0+}).$$

For $s = 0$ we have

$$E_1 \leq \|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)} \|\phi\|_{L^\infty(H^1)} + \|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)} \|\dot{\phi}\|_{L^\infty(L^2)}.$$

For $s = 1$ an integration by parts in space leads to

$$E_1 = -(y_0 - \mathcal{I}y_0; \Delta\phi^{0+}) + (y_1 - \mathcal{I}y_1; \dot{\phi}^{0+}).$$

An application of the Friedrichs inequality yields the following estimate

$$E_1 \leq C \left\{ \|h(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} \|\Delta\phi\|_{L^\infty(L^2)} + \|h(y_1 - \mathcal{I}y_1)\|_{H^1(\Omega)} \|\dot{\phi}\|_{L^\infty(H^1)} \right\}.$$

Note that for $s = 1$ we applied the following inequalities in the estimation of the second term

$$\|\phi\|_{L^2(\Omega)} \leq \|\phi\|_{H^1(\Omega)} \quad \text{and} \quad \|y_0 - \mathcal{I}y_0\|_{L^2(\Omega)} \leq \frac{1}{\pi} \|h(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)}.$$

This concludes the proof. $\qquad\square$

**Lemma 4.3.1.8.** There holds the following estimate

$$E_2 \leq \|kD^{2-s}U_2\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^s)} + \left(2 + \frac{2}{\pi}\right) \left(\|U_2\|_{L^1(H^1)} + \|\Delta U_2\|_{L^1(L^2)}\right) \|\phi\|_{L^\infty(H^1)},$$

where the second term appears only for $s = 0$, $\varepsilon > 0$ and $(DN)^*$ boundary conditions.

**Proof.** The properties of $\mathcal{G}$ imply

$$E_2 = \sum_{j=1}^{N} \int_{I_j} a(U_2; \phi - \phi^{j-1+}) dt. \qquad (4.67)$$

For $s = 0$, an integration by parts in space yields

$$E_2 = -\sum_{j=1}^{N} \int_{I_j} (\Delta U_2; \phi - \phi^{j-1+}) dt + \sum_{j=1}^{N} \int_{I_j} DU_2(t, 1)(\phi - \phi^{j-1+})(t, 1).$$

The second term on the RHS of the equation above vanishes for $(DD)^*$. Otherwise an application of the trace, Friedrichs and then of the Hölder inequality leads to

$$\sum_{j=1}^{N} \int_{I_j} DU_2(t, 1)(\phi - \phi^{j-1+})(t, 1) \leq \left(2 + \frac{2}{\pi}\right)\left(\|U_2\|_{L^1(H^1)} + \|\Delta U_2\|_{L^1(L^2)}\right) \|\phi\|_{L^\infty(H^1)}.$$

Hence,

$$E_2 \leq \|k\Delta U_2\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(L^2)} + \left(2 + \frac{2}{\pi}\right)\left(\|U_2\|_{L^1(H^1)} + \|\Delta U_2\|_{L^1(L^2)}\right) \|\phi\|_{L^\infty(H^1)}.$$

For $s = 1$, since $\phi(t) - \phi(t_{j-1}) = \int_{t_{j-1}}^{t} \dot{\phi}(\tau) d\tau$ we may deduce from (4.67)

$$E_2 \leq \sum_{j=1}^{N} \int_{I_j} \|U_2\|_{H^1(\Omega)} \|\phi - \phi^{j-1+}\|_{H^1(\Omega)} dt \leq \|kU_2\|_{L^1(H^1)} \|\dot{\phi}\|_{L^\infty(H^1)}. \qquad\square$

**Lemma 4.3.1.9.** For $s=0,1$ there is a constant $C>0$ such that there holds

$$E_3 \le C^s \Big\{ \|h^s(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(L^2)} \|\dot\phi\|_{L^\infty(H^s)}$$
$$+ \|k\mathcal{L}^{1-s}(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(H^{1-s})}\big(T^{1/2}\|D^{s+1}\phi\|_{L^\infty(L^2)}+\varepsilon\|D^{s+1}\dot\phi\|_{L^2(L^2)}\big) \Big\}.$$

**Proof.** We rewrite $E_3$ owing to the symmetry properties of $\mathcal{L}$ such that

$$E_3 = \sum_{j=1}^N \int_{I_j} (f-\mathcal{L}f+\Delta(U_1-\mathcal{L}U_1)+\varepsilon\Delta(U_2-\mathcal{L}U_2); \dot\phi-\mathcal{L}\dot\phi)dt$$
$$+ \sum_{j=1}^N \int_{I_j} (f+\Delta(U_1+\varepsilon U_2); \mathcal{L}(\dot\phi-\dot\phi^{j-1+}))dt. \tag{4.68}$$

For the first term on the RHS of the equation above we may deduce

$$\sum_{j=1}^N \int_{I_j} (f-\mathcal{L}f+\Delta(U_1-\mathcal{L}U_1)+\varepsilon\Delta(U_2-\mathcal{L}U_2); \dot\phi-\mathcal{L}\dot\phi)dt$$
$$\le C^s \|h^s(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(L^2)}\|\dot\phi\|_{L^\infty(H^s)}.$$

For the second term on the RHS of (4.68) we have since $\ddot\phi=\Delta(\phi-\varepsilon\dot\phi)$

$$\sum_{j=1}^N \int_{I_j} (f+\Delta(U_1+\varepsilon U_2); \mathcal{L}(\dot\phi-\dot\phi^{j-1+}))dt$$
$$\le \big(\|k(f+\Delta(U_1+\varepsilon U_2))\|_{L^2(0,T)}; \|\mathcal{L}\ddot\phi\|_{L^2(0,T)}\big)$$
$$\le \big(\|k(f+\Delta(U_1+\varepsilon U_2))\|_{L^2(0,T)}; \|\mathcal{L}\Delta(\phi-\varepsilon\dot\phi)\|_{L^2(0,T)}\big). \tag{4.69}$$

For $s=1$ a Hölder inequality in space proves

$$\sum_{j=1}^N \int_{I_j} (f+\Delta(U_1+\varepsilon U_2); \mathcal{L}(\dot\phi-\dot\phi^{j-1+}))dt \le \|k(f+\Delta(U_1+\varepsilon U_2))\|_{L^2(L^2)}\big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)}+\varepsilon\|\Delta\dot\phi\|_{L^2(L^2)}\big).$$

For $s=0$, we continue by using the same arguments as in the proof of Lemma 4.3.1.4. From (4.69) using the symmetry property of $\mathcal{L}$ and integrating by parts in space we obtain

$$\sum_{j=1}^N \int_{I_j} (f+\Delta(U_1+\varepsilon U_2); \mathcal{L}(\dot\phi-\dot\phi^{j-1+}))dt \le \big(\mathcal{L}\|k(f+\Delta(U_1+\varepsilon U_2))\|_{L^2(0,T)}; \|\Delta(\phi-\varepsilon\dot\phi)\|_{L^2(0,T)}\big)$$
$$\le a\big(\|\mathcal{L}k(f+\Delta(U_1+\varepsilon U_2))\|_{L^2(0,T)}; \|\phi-\varepsilon\dot\phi\|_{L^2(0,T)}\big)$$
$$\le \|k\mathcal{L}(f+\Delta(U_1+\varepsilon U_2))\|_{L^2(H^1)}\big(T^{1/2}\|\phi\|_{L^\infty(H^1)}+\varepsilon\|\dot\phi\|_{L^2(H^1)}\big).$$

This concludes the proof of the lemma. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

If we substitute the results of the lemmas above into the error representation (4.66) and recall the stability results from Lemma 4.1.0.5, we may conclude the proof of theorem. Note that in case $s=1$, the stability result (4.4c) and (4.4d) can be applied if (DD*) or $\varepsilon=0$, respectively.□

### 4.3.1.3 Dual a posteriori error analysis, $dG(1)\otimes\mathcal{P}_1$

**Theorem 4.3.1.3 (A posteriori dual error estimate, $dG(1)\otimes\mathcal{P}_1$).** There is a constant C such that the error of the $dG(1)\otimes\mathcal{P}_1$ approximation satisfies the following a posteriori error bound for $\varepsilon=0$

$$\|e^{N-}\|_{\widehat{\mathcal{H}}}\leq C\Big\{\|h(y_0-\mathcal{I}y_0)\|_{H^1(\Omega)}+\|h(y_1-\mathcal{I}y_1)\|_{H^1(\Omega)}+\|D_{h,3}(U_2)\|_{L^1(0,T)}+T^{1/2}\|k^2(U_2-\overline{U}_2)\|_{L^2(L^2)}$$

$$+\|h(f-\mathcal{L}f)\|_{L^1(L^2)}+T^{1/2}\|k^2\mathcal{L}(f-\bar{f})\|_{L^1(L^2)}+\sum_{j=1}^{N}\|h(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}$$

$$+D_{h,3}(U_1^{N-})+D_{h,3}(U_1^{0-})+\sum_{j=2}^{N}D_{h,3}(U_1^{j-1-})$$

$$+T^{1/2}\|k^2\mathcal{L}\mathcal{K}_h^{-1}(U_1-\overline{U}_1)\|_{L^2(H^1)}\Big\}. \tag{4.70}$$

The sums on the RHS of the estimate above go only over such $j=1,\dots,N$ where $\mathcal{S}^{j-1}\neq\mathcal{S}^j$.

**Remark 4.3.1.4.** The estimate (4.3.1.3) is of order $\mathcal{O}(h+k^3)$ if $\mathcal{S}^{j-1}=\mathcal{S}^j$ and of order $\mathcal{O}(h+hk^{-1}+k^3)$ otherwise. $\square$

**Proof.** Given the residual representation from Lemma 3.3.2.2, the error representation (4.54) can be rewritten such that

$$\big\langle e^{N-}\,;\,\Phi^{N-}\big\rangle_{\mathcal{H}}=\big\langle e^{0-}\,;\,\Phi^{0+}\big\rangle_{\mathcal{H}}+\sum_{j=1}^{N}\int_{I_j}a(U_2;\phi-V_1)dt-\sum_{j=1}^{N}\int_{I_j}a(U_{1\tau};\phi-V_1)dt$$

$$+\sum_{j=1}^{N}\int_{I_j}(f-U_{2\tau};\dot{\phi}-V_2)dt-\sum_{j=1}^{N}a([U_1]^{j-1};(\phi-V_1)^{j-1+})-\sum_{j=1}^{N}([U_2]^{j-1};(\dot{\phi}-V_2)^{j-1+})$$

$$+\sum_{j=1}^{N}\int_{I_j}\sum_{k=1}^{m}[D(U_1+\varepsilon U_2)]_k(\dot{\phi}-V_2)(x_k)dt=:\sum_{\ell=1}^{7}E_\ell. \tag{4.71}$$

Note that here $m=n$ if we deal with (DN*) boundary conditions and $m=n-1$ for (DD*) boundary conditions.
The idea is to estimate each of $E_1,\dots,E_7$ such that the dual solution contribution in the final estimate can be dominated by means of the stability results from Lemma 4.1.0.5.
Therefore, we choose a test function $V\in\mathcal{Q}_1$ such that

$$V=\mathcal{J}\Pi\Phi,\quad \Pi=(\mathcal{L},\mathcal{L}). \tag{4.72}$$

Here, $\Phi$ denotes the continuous solution of the dual problem (4.1) and $\mathcal{J}$ is the mapping onto the space of $dG(1)$ functions from Definition 3.1.0.9, case $dG(1)a$. Like in the Subsection 4.3.1.1, see Lemma 4.3.1.7, if $U^{0-}$ is chosen such that $U^{0-}=(\mathcal{I}y_0,\mathcal{I}y_1)$ where $\mathcal{I}$ denotes the nodal interpolant, we have

$$E_1\leq C\Big\{\|h(y_0-\mathcal{I}y_0)\|_{H^1(\Omega)}\|\Delta\phi\|_{L^\infty(L^2)}+\|h(y_1-\mathcal{I}y_1)\|_{H^1(\Omega)}\|\dot{\phi}\|_{L^\infty(H^1)}\Big\}.$$

The estimation of the remaining terms $E_2,\dots,E_7$ is given in the Lemmas below.

**Lemma 4.3.1.10.** There exists a constant $C$ such that for $\varepsilon=0$

$$E_2 \leq C\Big(\|D_{h,3}(U_2)\|_{L^1(0,T)}+T^{1/2}\|k^2(U_2-\overline{U}_2)\|_{L^2(L^2)}\Big)\|\Delta\phi\|_{L^\infty(L^2)}.$$

**Proof.** Due to the orthogonality property of $\mathcal{J}$, we have

$$E_2 = \sum_{j=1}^N \int_{I_j} a(U_2;\phi-\mathcal{L}\phi)dt+\sum_{j=1}^N \int_{I_j} a(U_2-\overline{U}_2;\mathcal{L}(\phi-\mathcal{J}\phi))dt.$$

The approximation properties of $\mathcal{L}$ given in Lemma 3.1.0.10 prove

$$\sum_{j=1}^N \int_{I_j} a(U_2;\phi-\mathcal{L}\phi)dt \leq C\|D_{h,3}(U_2)\|_{L^1(0,T)}\|\Delta\phi\|_{L^\infty(L^2)}.$$

For the second term the approximation properties of $\mathcal{J}$ imply

$$\sum_{j=1}^N \int_{I_j} a(U_2-\overline{U}_2;\mathcal{L}(\phi-\mathcal{J}\phi))dt \leq C\ a\big(\|k^2(U_2-\overline{U}_2)\|_{L^2(0,T)};\|\mathcal{L}\ddot{\phi}\|_{L^2(0,T)}\big).$$

Since $\ddot{\phi}=\Delta\phi$ for $\varepsilon=0$, the last inequality can be recast to

$$\sum_{j=1}^N \int_{I_j} a(U_2-\overline{U}_2;\mathcal{L}(\phi-\mathcal{J}\phi))dt \leq C\ a\big(\|k^2(U_2-\overline{U}_2)\|_{L^2(0,T)};T^{1/2}\|\mathcal{L}\Delta\phi\|_{L^\infty(0,T)}\big). \tag{4.73}$$

From the definition of the operator $\mathcal{K}_h^{-1}$, the last inequality is equivalent to

$$\sum_{j=1}^N \int_{I_j} a(U_2-\overline{U}_2;\mathcal{L}(\phi-\mathcal{J}\phi))dt \leq C\ \big(\|k^2\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|_{L^2(0,T)};T^{1/2}\|\mathcal{L}\Delta\phi\|_{L^\infty(0,T)}\big)$$

$$\leq CT^{1/2}\|k^2\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|_{L^2(L^2)}\|\Delta\phi\|_{L^\infty(L^2)}.$$

This concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Lemma 4.3.1.11.** There exists a constant $C$ such that for $\varepsilon=0$,

$$E_4 \leq C\big(\|h(f-\mathcal{L}f)\|_{L^1(L^2)}+T^{1/2}\|k^2\mathcal{L}(f-\bar{f})\|_{L^2(H^1)}\big)\|\dot{\phi}\|_{L^\infty(H^1)}.$$

**Proof.** Using the symmetry and the approximation properties of $\mathcal{L}$ and $\mathcal{J}$, we may rewrite $E_4$ such that

$$E_4 = \sum_{j=1}^N \int_{I_j} (f-\mathcal{L}f;\dot{\phi}-\mathcal{L}\dot{\phi})dt+\sum_{j=1}^N \int_{I_j} (f-\bar{f};\mathcal{L}\dot{\phi}-\mathcal{J}\mathcal{L}\dot{\phi})dt. \tag{4.74}$$

For the first term on the RHS of the equation above we may deduce,

$$\sum_{j=1}^N \int_{I_j} (f-\mathcal{L}f;\dot{\phi}-\mathcal{L}\dot{\phi})dt \leq C\|h(f-\mathcal{L}f)\|_{L^1(L^2)}\|\dot{\phi}\|_{L^\infty(H^1)}. \tag{4.75}$$

The second term can be estimated such that

$$\sum_{j=1}^{N}\int_{I_j}(f-\bar{f};\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt\leq C\big(\|k^2(f-\bar{f})\|_{L^2(0,T)};\|\mathcal{L}\frac{\partial^3}{\partial t}\phi\|_{L^2(0,T)}\big).$$

From the dual equation we have $\dddot{\phi}=\Delta(\dot{\phi}-\varepsilon\ddot{\phi})$. Then, owing to the symmetry properties of $\mathcal{L}$, we may deduce

$$\sum_{j=1}^{N}\int_{I_j}(f-\bar{f};\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt\leq C\big(\|k^2\mathcal{L}(f-\bar{f})\|_{L^2(0,T)};\|\Delta(\dot{\phi}-\varepsilon\ddot{\phi})\|_{L^2(0,T)}\big). \tag{4.76}$$

After integrating by parts in space, the term on the RHS of the equation (4.76) is equivalent to

$$\big(\|k^2\mathcal{L}(f-\bar{f})\|_{L^2(0,T)};\Delta\|\dot{\phi}-\varepsilon\ddot{\phi}\|_{L^2(0,T)}\big)=-a\big(\|k^2\mathcal{L}(f-\bar{f})\|_{L^2(0,T)};\|\dot{\phi}-\varepsilon\ddot{\phi}\|_{L^2(0,T)}\big). \tag{4.77}$$

If $\varepsilon=0$ we have

$$a\big(\mathcal{L}\|k^2(f-\bar{f})\|_{L^2(0,T)};\|\dot{\phi}-\varepsilon\ddot{\phi}\|_{L^2(0,T)}\big)\leq CT^{1/2}\|k^2\mathcal{L}(f-\bar{f})\|_{L^2(H^1)}\|\dot{\phi}\|_{L^\infty(H^1)}. \tag{4.78}$$

A combination of (4.78), (4.77) and (4.76) yields an estimate for the second term from the RHS of (4.74). If we recall an estimate (4.75), we may complete the proof. $\qquad\square$

**Lemma 4.3.1.12.** There exists a constant $C$ such that

$$E_6\leq C\Big(\sum_{j=1}^{N}\|h(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big)\|\dot{\phi}\|_{L^\infty(H^1)},$$

where the sum goes only over such $j$ where $\mathcal{S}^{j-1}\not\subseteq\mathcal{S}^{j}$.

**Proof.** Owing to the properties of $\mathcal{L},\mathcal{J}$, we have

$$E_6=\sum_{j=1}^{N}((I-\mathcal{L})U_2^{j-1-};(\dot{\phi}-\mathcal{L}\dot{\phi})^{j-1+}).$$

The approximation properties of $\mathcal{L}$ lead to

$$E_6\leq C\Big(\sum_{j=1}^{N}\|h(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big)\|\dot{\phi}\|_{L^\infty(H^1)}.$$

This concludes the proof. $\qquad\square$

**Lemma 4.3.1.13.** There exists a constant $C>0$ such that for $(DD^*)$ and $\varepsilon=0$

$$E_3+E_5+E_7\leq C\Big\{D_{h,3}(U_1^{N-})\|\Delta\phi^{N-}\|_{L^2(\Omega)}+D_{h,3}(U_1^{0-})\|\Delta\phi\|_{L^\infty(L^2)}$$

$$+\Big(\sum_{j=2}^{N}D_{h,3}(U^{j-1-})\Big)\|\Delta\phi\|_{L^\infty(L^2)}+T^{1/2}\|k^2\mathcal{L}\mathcal{K}_h^{-1}(U_1-\overline{U}_1)\|_{L^2(H^1)}\|\dot{\phi}\|_{L^\infty(H^1)}\Big\},$$

where the third summand (sum) on the RHS does not appear if for all $j=1,\ldots,N$, $\mathcal{S}^{j-1}=\mathcal{S}^{j}$.

**Proof.** Let $E_7$ be rewritten such that

$$E_7 = \sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)]_k (\dot{\phi} - \mathcal{L}\dot{\phi})(x_k) dt + \sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} [D(U_1 + \varepsilon U_2)]_k (\mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi}))(x_k) dt$$

$$= -\sum_{j=1}^{N} \int_{I_j} a(U_1 + \varepsilon U_2; \dot{\phi} - \mathcal{L}\dot{\phi}) dt - \sum_{j=1}^{N} \int_{I_j} a(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2); \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi})) dt. \qquad (4.79)$$

In the following we assume $\varepsilon = 0$. An integration by parts in time yields for the first term

$$-\sum_{j=1}^{N} \int_{I_j} a(U_1; \dot{\phi} - \mathcal{L}\dot{\phi}) dt = \sum_{j=1}^{N} \int_{I_j} a(U_{1,\tau}; \phi - \mathcal{L}\phi) dt - \sum_{j=1}^{N} \int_{I_j} a(U_1^{j-}; (\phi - \mathcal{L}\phi)^{j-}) dt$$

$$+ \sum_{j=1}^{N} \int_{I_j} a(U_1^{j-1+}; (\phi - \mathcal{L}\phi)^{j-1+}) dt.$$

Then

$$E_3 + E_5 - \sum_{j=1}^{N} \int_{I_j} a(U_1; \dot{\phi} - \mathcal{L}\dot{\phi}) dt = -a(U_1^{N-}; (\phi - \mathcal{L}\phi)^{N-}) dt$$

$$+ a(U_1^{0-}; (\phi - \mathcal{L}\phi)^{0+}) + \sum_{j=2}^{N} a([\phi - \mathcal{L}\phi]^{j-1}; U^{j-1-})$$

$$\leq C \Big\{ D_{h,3}(U_1^{N-}) \|\Delta\phi^{N-}\|_{L^2(\Omega)} + D_{h,3}(U_1^{0-}) \|\Delta\phi\|_{L^\infty(L^2)}$$

$$+ \Big( \sum_{j=2}^{N} D_{h,3}(U^{j-1-}) \Big) \|\Delta\phi\|_{L^\infty(L^2)} \Big\}. \qquad (4.80)$$

Furthermore, on account to the projection properties of $\mathcal{J}$, we have for the second term on the RHS of (4.79)

$$-\sum_{j=1}^{N} \int_{I_j} a(U_1 - \overline{U}_1; \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi})) dt \leq Ca\big(\|k^2(U_1 - \overline{U}_1)\|_{L^2(0,T)}; \|\mathcal{L}\Delta\dot{\phi}\|_{L^2(0,T)}\big)$$

$$\leq C\big(\|k^2 \mathcal{L}\mathcal{K}_h^{-1}(U_1 - \overline{U}_1)\|_{L^2(0,T)}; \|\Delta\dot{\phi}\|_{L^2(0,T)}\big)$$

$$\leq Ca\big(\|k^2 \mathcal{L}\mathcal{K}_h^{-1}(U_1 - \overline{U}_1)\|_{L^2(0,T)}; \|\dot{\phi}\|_{L^2(0,T)}\big)$$

$$\leq CT^{1/2} \|k^2 \mathcal{L}\mathcal{K}_h^{-1}(U_1 - \overline{U}_1)\|_{L^2(H^1)} \|\dot{\phi}\|_{L^\infty(H^1)}.$$

Here we used the symmetry properties of $\mathcal{L}$, properties of operator $\mathcal{K}_h^{-1}$ and the fact that $\ddot{\phi} = \Delta\phi$ for $\varepsilon = 0$. $\qquad \square$

If we substitute the estimates for $E_1, \ldots, E_7$ given below into the error representation (4.71) and apply the stability results of Lemma 4.1.0.5, we may conclude the proof. $\qquad \square$

### 4.3.1.4 Dual a posteriori error analysis, $dG(1) \otimes C^1$

In the following we provide the estimate for $\|e^{N-}\|_{\widehat{\mathcal{H}}}$ when (DD*), i.e. (DD) hold.

**Theorem 4.3.1.4 (Dual a posteriori error estimate, $dG(1) \otimes C^1$).** There exists a constant C such that the error of $dG(1) \otimes C^1$ finite element approximation satisfies the following a posteriori error bound if (DD*) holds

1. For $s=0$ and $\varepsilon=0$

$$
\begin{aligned}
\|e^{N-}\|_{\mathcal{H}} \leq C\Big\{ &\|y_0 - \mathcal{I}y_0\|_{H^1(\Omega)} + \|y_1 - \mathcal{I}y_1\|_{L^2(\Omega)} + \|h(I-\mathcal{L})U_2\|_{L^1(L^2)} \\
&+ T^{1/2}\|k^2\mathcal{L}\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|^2_{L^1(H^1)} + \|(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2)\|^2_{L^1(L^2)} \\
&+ T^{1/2}\|k^2(f-\bar{f}+\Delta(U_1-\overline{U}_1)\|^2_{L^2(L^2)} + \|(I-\mathcal{L})\Delta U_1\|_{L^1(L^2)} \\
&+ \sum_{j=1}^{N+1}\|h(I-\mathcal{L})\Delta U_1^{j-1-}\|_{L^2(\Omega)} + \sum_{j=1}^{N}\|(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big\}. \quad (4.81\text{a})
\end{aligned}
$$

2. For $s=1$ and $\varepsilon \geq 0$

$$
\begin{aligned}
\|e^{N-}\|_{\widehat{\mathcal{H}}} \leq C\Big\{ &\|h(y_0-\mathcal{I}y_0)\|_{H^1(\Omega)} + \|h(y_1-\mathcal{I}y_1)\|_{H^1(\Omega)} + \|h^2(I-\mathcal{L})U_2\|_{L^1(L^2)} \\
&+ (T^{1/2}+\varepsilon^{1/2})\|k^2\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|_{L^2(L^2)} + \|h(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2)\|_{L^1(L^2)} \\
&+ (T^{1/2}+\varepsilon T^{1/2}+\varepsilon^{3/2})\|k^2\Delta\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2))\|_{L^2(L^2)} \\
&+ \|h(I-\mathcal{L})\Delta U_1\|_{L^1(L^2)} + \sum_{j=1}^{N+1}\|h^2(I-\mathcal{L})\Delta U_1^{j-1}\|_{L^2(\Omega)} \\
&+ \sum_{j=1}^{N}\|h(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big\}. \quad (4.81\text{b})
\end{aligned}
$$

The last two summands on the RHS of the both estimates above do not appear if for all $j=1,\dots,N$, $\mathcal{S}^{j-1}=\mathcal{S}^j$. We also assume that the mesh in space is quasi-uniform.

**Remark 4.3.1.5.** The proven estimates are of order $\mathcal{O}(h^{s+3}+k^3)$ if $\mathcal{S}^{j-1} = \mathcal{S}^j$ for all $j = 1,\dots,N$ and of order $\mathcal{O}(h^{s+3}+h^{s+3}k^{-1}+k^3)$ otherwise. $\square$

**Proof.** From residual representation, cf. Lemma 3.3.2.2, the error representation (4.54) is equivalent to

$$
\begin{aligned}
\langle e^{N-} ; \Phi^{N-} \rangle_{\mathcal{H}} = &\langle e^{0-} ; \Phi^{0+} \rangle_{\mathcal{H}} + \sum_{j=1}^{N}\int_{I_j} a(U_2;\phi-V_1)dt - \sum_{j=1}^{N}\int_{I_j} a(U_{1\tau};\phi-V_1)dt \\
&+ \sum_{j=1}^{N}\int_{I_j}(f-U_{2\tau}+\Delta(U_1+\varepsilon U_2);\dot\phi-V_2)dt \\
&- \sum_{j=1}^{N}\langle [U]^{j-1} ; (\Phi-V)^{j-1+} \rangle_{\mathcal{H}} =: \sum_{\ell=1}^{5} E_\ell. \quad (4.82)
\end{aligned}
$$

We choose a test function $V \in \mathcal{Q}_1$ such that

$$V := \mathcal{J} \Pi \Phi \quad \Pi = (\mathcal{L}, \mathcal{L}),$$

where $\Phi$ is strong solution and $\mathcal{J}$ is the temporal $L^2$ projection introduced in Definition 3.1.0.9, case $dG(1)a$.

The idea is to estimate each of the terms $E_1, \ldots, E_4$ so that the stability Lemma 4.1.0.5 can be applied.

If $U^{0-}$ is chosen as $U^{0-} = (\mathcal{I}y_0, \mathcal{I}y_1)$ where $\mathcal{I}$ denotes the cubic Hermite interpolant, then there holds for $s = 0, 1$

$$E_1 \leq C \Big\{ \|h^s(y_0 - \mathcal{I}y_0)\|_{L^2(\Omega)} \|D^{s+1}\phi\|_{L^\infty(L^2)} + \|h^s(y_1 - \mathcal{I}y_1)\|_{H^s(\Omega)} \|\dot{\phi}\|_{L^\infty(H^s)} \Big\},$$

cf. Lemma 4.3.1.7.

**Lemma 4.3.1.14.** If $(DD^*)$ then for $s = 0, 1$ there holds

$$E_2 \leq C \Big\{ \|h^{s+1}(I - \mathcal{L})U_2\|_{L^1(L^2)} \|\Delta\phi\|_{L^\infty(L^2)}$$
$$+ \|k^2 \mathcal{L}^{1-s} \mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(H^{1-s})} \big( T^{1/2} \|D^{s+1}\phi\|_{L^\infty(L^2)} + \varepsilon \|D^{s+1}\dot{\phi}\|_{L^2(L^2)} \big) \Big\}.$$

**Proof.** We start from the following decomposition

$$E_2 = \sum_{j=1}^N \int_{I_j} a(U_2; \phi - \mathcal{L}\phi)dt + \sum_{j=1}^N \int_{I_j} a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt, \qquad (4.83)$$

since $\mathcal{J}$ is orthogonal to the functions constant in time. If we assume that $(DD)$ ie.e $(DD^*)$ hold, an integration by parts in space yields

$$\sum_{j=1}^N \int_{I_j} a(U_2; \phi - \mathcal{L}\phi)dt = -\sum_{j=1}^N \int_{I_j} (\Delta U_2; \phi - \mathcal{L}\phi)dt.$$

Then the first term on the RHS of (4.84) can be estimated for $s = 0, 1$ such that

$$\sum_{j=1}^N \int_{I_j} a(U_2; \phi - \mathcal{L}\phi)dt \leq C \|h^{s+1}(I - \mathcal{L})U_2\|_{L^1(L^2)} \|D^{s+1}\phi\|_{L^\infty(L^2)}.$$

From $\ddot{\phi} = \Delta(\phi - \varepsilon\dot{\phi})$, the definition of the operator $\mathcal{K}_h^{-1}$, and the symmetry properties of $\mathcal{L}$, the second term reads

$$\sum_{j=1}^N \int_{I_j} a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt \leq C\, a\big( \|k^2(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\mathcal{L}\ddot{\phi}\|_{L^2(0,T)} \big)$$

$$= C\, a\big( \|k^2(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\mathcal{L}\Delta(\phi - \varepsilon\dot{\phi})\|_{L^2(0,T)} \big)$$

$$= C\, \big( \|k^2 \mathcal{L}\mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\Delta(\phi - \varepsilon\dot{\phi})\|_{L^2(0,T)} \big). \qquad (4.84)$$

For $s = 1$, the RHS of the inequality above can be further estimated such that

$$\sum_{j=1}^N \int_{I_j} a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt \leq C\, \|k^2 \mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(L^2)} \big( T^{1/2} \|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon \|\Delta\dot{\phi}\|_{L^2(L^2)} \big).$$

If $s=0$, an integration by parts in space in the RHS of (4.84) yields

$$\sum_{j=1}^{N}\int_{I_j} a(U_2-\overline{U}_2;\mathcal{L}(\phi-\mathcal{J}\phi))dt \leq C\; a\big(\|k^2\mathcal{L}\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|_{L^2(0,T)}; \|\phi-\varepsilon\dot{\phi}\|_{L^2(0,T)}\big).$$

Hence,

$$\sum_{j=1}^{N}\int_{I_j} a(U_2-\overline{U}_2;\mathcal{L}(\phi-\mathcal{J}\phi))dt \leq C\;\|k^2\mathcal{L}\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|_{L^2(H^1)}\big(T^{1/2}\|\phi\|_{L^\infty(H^1)}+\varepsilon\|\dot{\phi}\|_{L^2(H^1)}\big).$$

This concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

**Lemma 4.3.1.15.** For (DD*) there exists a constant $C>0$ such that there holds

1. for $s=0$ and $\varepsilon=0$

$$E_4 \leq C\big\{\|(I-\mathcal{L})(f+\Delta U_1)\|_{L^1(L^2)}\|\dot{\phi}\|_{L^\infty(L^2)}+T^{1/2}\|k^2\Delta\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1))\|_{L^2(L^2)}\|\dot{\phi}\|_{L^\infty(L^2)}\big\},$$

2. for $s=1$

$$\begin{aligned}
E_4 \leq C\big\{&\|h(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(L^2)}\|\dot{\phi}\|_{L^\infty(H^1)}\\
&+\|k^2\Delta\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)))\|_{L^2(L^2)}\big(T^{1/2}\|\dot{\phi}\|_{L^\infty(L^2)}\\
&+\varepsilon T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)}+\varepsilon^2\|\Delta\dot{\phi}\|_{L^2(L^2)}\big)\big\}.
\end{aligned}$$

**Proof.** According to the properties of $\mathcal{L}$ and $\mathcal{J}$, we have in $E_4$

$$\begin{aligned}
E_4 = &\sum_{j=1}^{N}\int_{I_j}((I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2)); \dot{\phi}-\mathcal{L}\dot{\phi})dt\\
&+\sum_{j=1}^{N}\int_{I_j}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)));\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt
\end{aligned}$$

For the first term on the RHS of the equation above we have by using the approximation properties of $\mathcal{L}$ for both $s=0,1$

$$\sum_{j=1}^{N}\int_{I_j}((I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2)); \dot{\phi}-\mathcal{L}\dot{\phi})dt \leq C^s\|h^s(I-\mathcal{L})(f+\Delta(U_1+\varepsilon U_2))\|_{L^1(L^2)}\|\dot{\phi}\|_{L^\infty(H^s)}.$$

Furthermore, the second term can be estimated such that

$$\sum_{j=1}^{N}\int_{I_j}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)));\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt$$

$$\leq C\;\big(\|k^2(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)))\|_{L^2(0,T)}; \|\mathcal{L}\frac{\partial^3}{\partial t}\phi\|_{L^2(0,T)}\big). \qquad (4.85)$$

Since $\frac{\partial^3}{\partial t}\phi = \Delta(\dot{\phi} - \varepsilon\ddot{\phi})$, an integration by parts in space yields

$$\left(\|k^2(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(0,T)}; \|\mathcal{L}\frac{\partial^3}{\partial t}\phi\|_{L^2(0,T)}\right)$$

$$\leq \left(\|k^2\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(0,T)}; \|\Delta(\dot{\phi} - \varepsilon\ddot{\phi})\|_{L^2(0,T)}\right)$$

$$\leq a\left(\|k^2\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(0,T)}; \|\dot{\phi} - \varepsilon\ddot{\phi}\|_{L^2(0,T)}\right) \qquad (4.86)$$

If $(DD^*)$, we have by integrating by parts in the last term on the RHS of the inequality above

$$a\left(\|k^2\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(0,T)}; \|\dot{\phi} - \varepsilon\ddot{\phi}\|_{L^2(0,T)}\right)$$

$$\leq \left(\|k^2\Delta\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(0,T)}; \|\dot{\phi} - \varepsilon\ddot{\phi}\|_{L^2(0,T)}\right) \qquad (4.87)$$

For $s = 0$ and $\varepsilon = 0$ we may then conclude

$$\sum_{j=1}^{N}\int_{I_j}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2))); \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi}))dt$$

$$\leq CT^{1/2}\|k^2\Delta\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(L^2)}\|\dot{\phi}\|_{L^\infty(L^2)}.$$

For $s = 1$, from (4.85), (4.86) and (4.87) we have

$$\sum_{j=1}^{N}\int_{I_j}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2))); \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi}))dt$$

$$\leq C\ \|k^2\Delta\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(L^2)}\big(T^{1/2}\|\dot{\phi}\|_{L^\infty(H^1)}$$

$$+ \varepsilon T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon^2\|\Delta\dot{\phi}\|_{L^2(L^2)}\big)$$

This concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Lemma 4.3.1.16.** There exists a constant $C > 0$ such that

$$E_3 + E_5 \leq C\Big\{\|h^s(I - \mathcal{L})\Delta U_1\|_{L^1(L^2)}\|\dot{\phi}\|_{L^\infty(H^s)} + \|h^{s+1}(I - \mathcal{L})\Delta U_1^{N-}\|_{L^2(\Omega)}\|D^{s+1}\phi^{N-}\|_{L^2(\Omega)}$$

$$+ \|h^{s+1}(I - \mathcal{L})\Delta U_1^{0-}\|_{L^2(\Omega)}\|D^{s+1}\phi\|_{L^\infty(L^2)} + \sum_{j=2}^{N}\|(h^{s+1}(I - \mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\|D^{s+1}\dot{\phi}\|_{L^\infty(L^2)}$$

$$+ \Big(\sum_{j=1}^{N}\|h^s(I - \mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big)\|\dot{\phi}\|_{L^\infty(H^s)}\Big\},$$

where the last two sums on the RHS of the equation above go only over such $j = 1, \ldots, N$ where $\mathcal{S}^{j-1} \neq \mathcal{S}^j$.

**Proof.** Owing to he properties of $\mathcal{L}$ and $\mathcal{J}$ we have

$$E_3 + E_5 = -\sum_{j=1}^{N}\int_{I_j}a(U_{1,\tau}; \phi - \mathcal{L}\phi)dt - \sum_{j=1}^{N}a([U_1]^{j-1}; (\phi - \mathcal{L}\phi)^{j-1+})$$

$$+ \sum_{j=1}^{N}((I - \mathcal{L})U_2^{j-1-}; (\dot{\phi} - \mathcal{L}\dot{\phi})^{j-1+}).$$

If we integrate by parts in time in the first term, we obtain

$$E_3+E_5=\sum_{j=1}^{N}\int_{I_j}a(U_1;\dot\phi-\mathcal{L}\dot\phi)dt-a(U_1^{N-};(\phi-\mathcal{L}\phi)^{N-})+a(U_1^{0-};(\phi-\mathcal{L}\phi)^{0+})$$

$$+\sum_{j=2}^{N}a(U_1^{j-1-};[\phi-\mathcal{L}\phi]^{j-1})+\sum_{j=1}^{N}((I-\mathcal{L})U_2^{j-1-};(\dot\phi-\mathcal{L}\dot\phi)^{j-1+}).$$

If we assume that the spatial mesh does not change i.e. $\mathcal{S}^{j-1}=\mathcal{S}^j$ for all $j=1,\ldots,N$, then the 4th and the 5th summand on the RHS of the equation above equals zero.
Anyhow, if we integrate by parts in space and make use of the (DD) boundary conditions, then

$$E_3+E_5= -\sum_{j=1}^{N}\int_{I_j}((I-\mathcal{L})\Delta U_1;\dot\phi-\mathcal{L}\dot\phi)dt+((I-\mathcal{L})U_1^{N-};(\phi-\mathcal{L}\phi)^{N-})$$

$$-((I-\mathcal{L})\Delta U_1^{0-};(\phi-\mathcal{L}\phi)^{0+})+\sum_{j=2}^{N}((I-\mathcal{L})\Delta U_1^{j-1-};[\phi-\mathcal{L}\phi]^{j-1})dt$$

$$+\sum_{j=1}^{N}((I-\mathcal{L})U_2^{j-1-};(\dot\phi-\mathcal{L}\dot\phi)^{j-1+}).$$

Finally, the approximation properties of $\mathcal{L}$ imply

$$E_3+E_5\le C\Big\{\|h^s(I-\mathcal{L})\Delta U_1\|_{L^1(L^2)}\|\dot\phi\|_{L^\infty(H^s)}+\|h^{s+1}(I-\mathcal{L})\Delta U_1^{N-}\|_{L^2(\Omega)}\|D^{s+1}\phi^{N-}\|_{L^2(\Omega)}$$

$$+\|h^{s+1}(I-\mathcal{L})\Delta U_1^{0-}\|_{L^2(\Omega)}\|D^{s+1}\phi\|_{L^\infty(L^2)}+\sum_{j=2}^{N}\|(h^{s+1}(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\|D^{s+1}\dot\phi\|_{L^\infty(L^2)}$$

$$+\Big(\sum_{j=1}^{N}\|h^s(I-\mathcal{L})U_2^{j-1-}\|_{L^2(\Omega)}\Big)\|\dot\phi\|_{L^\infty(H^s)}\Big\}.$$

This concludes the proof.  □

Recalling Lemma 4.1.0.5 and the error representation from (4.82) we may conclude the proof. □

**Remark 4.3.1.6.** In [42], the author proved the estimate for $\|e^{N-}\|_{\mathcal{H}}$ by deriving the estimates in terms of $\|\cdot\|_{L^2(\Omega)^*}$ norm which is equivalent to the $L^2$ norm on the space of finite element functions. Here we did not used this equivalent norm and therefore only the estimate for $\|e^{N-}\|_{\widehat{\mathcal{H}}}$ is proven.  □

### 4.3.2  $cG(1)$ time approximation

Within the following subsection we analyse an a posteriori error bound and its derivation for the case of $cG(1)$ time approximation and $\mathcal{P}_1(\mathcal{C}^1)$ approximation in space. The notations and definitions are adopted from Subsections 2.1.1 and 2.3.4. The presented analysis employs the residual $Res$ defined in Definition 2.3.1.2 with bilinear form $\mathcal{B}$ defined in (2.41) and its dual form $\mathcal{B}^*$ from (4.34). Note that in case of $cG(1)$ method in time, $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$ for all $j=1,\ldots,N$.

**Lemma 4.3.2.1 (Dual error representation, $cG(1)$ time approximation).** If $u$ is a solution of (1.28), $U$ its discrete $cG(1)$ variant and $\Phi$, a strong solution of the corresponding dual problem (4.3) there holds

$$\big\langle \Phi(T)\,;V(T)\big\rangle_{\mathcal{H}} = Res(e,\Phi-V)+\big\langle \Phi(0)\,;V(0)\big\rangle_{\mathcal{H}} \quad \text{for all} \quad V\in\mathcal{W}_c. \tag{4.88}$$

**Proof.** Since $\Phi$ is solves the weak discrete problem (4.35), see Remark 4.2.2.1, and both $\Phi,e$ belong to the space of the continuous functions, from definition of the dual bilinear form, see Definition 4.1.0.2, we have

$$\big\langle \Phi(T)\,;e(T)\big\rangle_{\mathcal{H}} - \big\langle \Phi(0)\,;e(0)\big\rangle_{\mathcal{H}} = \mathcal{B}^*(\Phi,e) = \mathcal{B}(e,\Phi) = Res(\Phi) = Res(\Phi-V),$$

for all $V\in\mathcal{W}_c$.                                                                           $\square$

### 4.3.2.1 Dual a posteriori error analysis, $cG(1)\otimes\mathcal{P}_1$

**Theorem 4.3.2.1 (A posteriori dual error estimate, $cG(1)\otimes\mathcal{P}_1$).** There exists a constant $C>0$ such that the error of the $cG(1)\otimes\mathcal{P}_1$ finite element approximation satisfies the following a posteriori error bound if (DD$^*$) or $\varepsilon=0$,

$$\begin{aligned}
\|e(T)\|_{\widehat{\mathcal{H}}} \leq C\Big\{ &\|h(y_0-\mathcal{I}y_0)\|_{H^1(\Omega)}+\|h(y_1-\mathcal{I}y_1)\|_{H^1(\Omega)}+\|k(U_2-\overline{U}_2)\|_{L^1(H^1)} \\
&+\|h(f-\mathcal{L}f)\|_{L^1(L^2)}+(T^{1/2}+\varepsilon^{1/2})\|k(f-\bar{f})\|_{L^1(L^2)} \\
&+\|D_{h,3}(\dot{U}_1+\varepsilon\dot{U}_2)\|_{L^1(\Omega)}+D_{h,3}((U_1+\varepsilon U_2)(T))+D_{h,3}((U_1+\varepsilon U_2)(0)) \\
&+(T^{1/2}+\varepsilon^{1/2})\|k\mathcal{K}_h^{-1}(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)\|_{L^2(L^2)} \Big\}.
\end{aligned} \tag{4.89}$$

**Remark 4.3.2.1.** The estimate (4.3.2.1) is of order $\mathcal{O}(h+k^2)$.                    $\square$

**Proof.** A substitution of into error representation formula (4.88) and use of the residual definition from Lemma 3.3.3.2 yields

$$\begin{aligned}
\big\langle e(T)\,;\Phi(T)\big\rangle_{\mathcal{H}} = \big\langle e(0)\,;\Phi(0)\big\rangle_{\mathcal{H}} &+ \int_0^T a(U_2-\dot{U}_1;\Phi-V_1)dt+\int_0^T (f-\dot{U}_2;\dot{\Phi}-V_2)dt \\
&+\int_0^T\sum_{k=1}^m [D(U_1+\varepsilon U_2)]_k(\dot{\Phi}-V_2)(x_k)dt =: \sum_{\ell=1}^4 E_\ell.
\end{aligned} \tag{4.90}$$

We may choose a test function $V\in\mathcal{W}_c$ such that

$$V:=\mathcal{J}\Pi\Phi,\quad \Pi=(\mathcal{G},\mathcal{L}) \tag{4.91}$$

where $\Phi$ is the continuous solution of the dual problem and $\mathcal{J}$ is the temporal projection on the space of constant functions in time i.e. integral mean, see Definition 3.1.0.9, case $cG(1)$. In order to derive an estimate for $\|e(T)\|_{\widehat{\mathcal{H}}}$, we need to estimate each of $E_1,\dots,E_4$ such that the stability Lemma 4.1.0.5 can be applied. This is done in the following Lemmas.

**Lemma 4.3.2.2.** If a discrete variant $U(0)$ of the continuous initial solution $u_0=(y_0,y_1)$ is chosen such that

$$U(0) = (\mathcal{I}y_0,\mathcal{I}y_1)$$

for $\mathcal{I}$ a nodal interpolant, then there holds

$$E_1\leq C\Big\{\|h(y_0-\mathcal{I}y_0)\|_{H^1(\Omega)}\|\Delta\phi\|_{L^\infty(L^2)}+\|h(y_1-\mathcal{I}y_1)\|_{H^1(\Omega)}\|\dot{\phi}\|_{L^\infty(H^1)}\Big\}.$$

**Proof.** We start from

$$E_1 = a(y_0 - \mathcal{I}y_0; \phi(0)) + (y_1 - \mathcal{I}y_1; \dot{\phi}(0)).$$

Since the nodal interpolation $\mathcal{I}$ and the Galerkin projection $\mathcal{G}$ coincide in $1D$, we may derive

$$E_1 \leq C\Big\{ \|h(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} \|\Delta\phi\|_{L^\infty(L^2)} + \|h(y_1 - \mathcal{I}y_1)\|_{H^1(\Omega)} \|\dot{\phi}\|_{L^\infty(H^1)} \Big\}.$$

This completes the proof. $\qquad\square$

**Lemma 4.3.2.3.** There holds

$$E_2 \leq \|k(U_2 - \overline{U}_2)\|_{L^1(H^1)} \|\dot{\phi}\|_{L^\infty(H^1)}.$$

**Proof.** Owing to the orthogonality properties of $\mathcal{G}$ and $\mathcal{J}$ we have

$$E_2 = \int_0^T a(U_2 - \overline{U}_2; \phi - \bar{\phi}) dt.$$

Since $(\phi - \bar{\phi})|_{I_j} \leq \int_{I_j} \dot{\phi}(t) dt$, an application of the Hölder inequality yields the proof

$$E_2 \leq \|k(U_2 - \overline{U}_2)\|_{L^1(H^1)} \|\dot{\phi}\|_{L^\infty(H^1)}.\qquad\square$$

**Lemma 4.3.2.4.** There exists a constant $C$ such that

$$E_3 \leq \Big\{ \|h(f - \mathcal{L}f)\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^1)} + \|k(f - \bar{f})\|_{L^2(L^2)} \big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon\|\Delta\dot{\phi}\|_{L^2(L^2)}\big) \Big\}.$$

**Proof.** Using the orthogonality properties of $\mathcal{L}$ and $\mathcal{J}$, we have

$$E_3 = \int_0^T (f - \mathcal{L}f; \dot{\phi} - \mathcal{L}\dot{\phi}) dt + \int_0^T (f - \bar{f}; \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi})) dt.$$

For the first term we have

$$\int_0^T (f - \mathcal{L}f; \dot{\phi} - \mathcal{L}\dot{\phi}) dt \leq C\|h(f - \mathcal{L}f)\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^1)}$$

In case of the second term we may deduce by using the fact that $\ddot{\phi} = \Delta(\phi - \varepsilon\dot{\phi})$

$$\int_0^T (f - \bar{f}; \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi})) dt \leq \big(\|k(f - \bar{f})\|_{L^2(0,T)}; \|\mathcal{L}\ddot{\phi}\|_{L^2(0,T)}\big)$$

$$\leq \big(\|k\mathcal{L}(f - \bar{f})\|_{L^2(0,T)}; \|\Delta(\phi - \varepsilon\dot{\phi})\|_{L^2(0,T)}\big).$$

Furthermore, an application of the Hölder inequality in time and in space yields

$$\big(\|k\mathcal{L}(f - \bar{f})\|_{L^2(0,T)}; \|\Delta(\phi - \varepsilon\dot{\phi})\|_{L^2(0,T)}\big) \leq \|k(f - \bar{f})\|_{L^2(L^2)}\big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon\|\Delta\dot{\phi}\|_{L^2(L^2)}\big).$$

From the last two inequalities we derive an estimate for the second term and this concludes the proof of the Lemma. $\qquad\square$

**Lemma 4.3.2.5.** There exists a constant $C > 0$ such that

$$E_4 \leq C\Big\{\|D_{h,3}(\dot{U}_1 + \varepsilon\dot{U}_2)\|_{L^1(\Omega)}\|\Delta\phi\|_{L^\infty(L^2)} + D_{h,3}((U_1 + \varepsilon U_2)(T))\|\Delta\phi(T)\|_{L^2(\Omega)}$$
$$+ D_{h,3}((U_1 + \varepsilon U_2)(0))\|\Delta\phi\|_{L^\infty(L^2)}$$
$$+ \|k\mathcal{K}_h^{-1}(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)\|_{L^2(L^2)}\big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon\|\Delta\dot{\phi}\|_{L^2(L^2)}\big)\Big\}.$$

**Proof.** We may rewrite $E_4$ such that

$$E_4 = \int_0^T \sum_{k=1}^m [D(U_1 + \varepsilon U_2)]_k(\dot{\phi} - \mathcal{L}\dot{\phi}) - \int_0^T \sum_{k=1}^m [D(U_1 + \varepsilon U_2)]_k\mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi})(x_k)dt$$
$$= -\int_0^T a(U_1 + \varepsilon U_2; \dot{\phi} - \mathcal{L}\dot{\phi})dt - \int_0^T a(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2); \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi}))dt. \qquad (4.92)$$

An integration by parts in time in the first term yields

$$\int_0^T a(U_1 + \varepsilon U_2; \dot{\phi} - \mathcal{L}\dot{\phi})dt = -\int_0^T a(\dot{U}_1 + \varepsilon\dot{U}_2; \phi - \mathcal{L}\phi)dt + a((U_1 + \varepsilon U_2)(T); (\phi - \mathcal{L}\phi)(T))$$
$$- a((U_1 + \varepsilon U_2)(0); (\phi - \mathcal{L}\phi)(0).$$

According to Lemma 3.1.0.10, we may conclude

$$\int_0^T a(U_1 + \varepsilon U_2; \dot{\phi} - \mathcal{L}\dot{\phi})dt \leq \|D_{h,3}(\dot{U}_1 + \varepsilon\dot{U}_2)\|_{L^1(\Omega)}\|\Delta\phi\|_{L^\infty(L^2)} + D_{h,3}((U_1 + \varepsilon U_2)(T))\|\Delta\phi(T)\|_{L^2(\Omega)}$$
$$+ D_{h,3}((U_1 + \varepsilon U_2)(0))\|\Delta\phi\|_{L^\infty(L^2)}.$$

For the second term on the RHS of (4.92) we have

$$\int_0^T a(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2); \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi}))dt$$
$$\leq a\big(\|k(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\mathcal{L}\ddot{\phi}\|_{L^2(0,T)}\big)$$
$$= \big(\|k\mathcal{K}_h^{-1}(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\mathcal{L}\ddot{\phi}\|_{L^2(0,T)}\big)$$
$$\leq \|k\mathcal{K}_h^{-1}(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)\|_{L^2(L^2)}\big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon\|\Delta\dot{\phi}\|_{L^2(L^2)}\big).$$

This concludes the proof of Lemma.                                                                                  □

Recalling Lemma 4.1.0.5 and error representation (4.90) we may conclude the proof.                  □

### 4.3.2.2 Dual a posteriori error estimate, $cG(1) \otimes \mathcal{C}^1$

**Theorem 4.3.2.2 (A posteriori dual error estimate, $cG(1) \otimes \mathcal{C}^1$).** There is a constant $C > 0$ such that for (DD*), the error of the $cG(1) \otimes \mathcal{C}^1$ approximation satisfies the following a posteriori error bound

$$\|e_1(T)\|_{H^{1-s}} + \|e_2(T)\|_{H^{-s}} \leq C\Big\{\|h^s(y_0 - \mathcal{I}y_0)\|_{H^1(\Omega)} + \|h^s(y_1 - \mathcal{I}y_1)\|_{H^s(\Omega)} + \|h^{s+1}\Delta(U_2 - \overline{U}_2)\|_{L^1(H^1)}$$
$$+ (T^{1/2} + \varepsilon)\|k\mathcal{L}\mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^1(H^{1-s})}$$
$$+ \|h^s(I - \mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)}$$
$$+ (T^{1/2} + \varepsilon)\|k\mathcal{L}(f - \overline{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2))\|_{L^1(L^2)}\Big\}. (4.93)$$

**Remark 4.3.2.2.** The estimate (4.93) is of convergence order $\mathcal{O}(h^{s+3}+k^2)$. $\qquad\square$

**Proof.** Given the error representation formula (4.88) where residual takes a form as in Lemma 3.3.3.2, case $cG(1)$, we have

$$\langle e(T)\,;\Phi(T)\rangle_{\mathcal{H}}=\langle e(0)\,;\Phi(0)\rangle_{\mathcal{H}}+\int_0^T a(U_2-\dot{U}_1;\phi-V_1)dt$$

$$+\int_0^T (f-\dot{U}_2+\Delta(U_1+\varepsilon U_2);\dot{\phi}-V_2)dt=:\sum_{\ell=1}^3 E_\ell. \qquad (4.94)$$

In the following we choose a test function $V\in\mathcal{W}_c$ such that

$$V:=\mathcal{J}\Pi\Phi,\quad \Pi=(\mathcal{L},\mathcal{L}),$$

where $\mathcal{J}$ is the temporal $L^2$ projection orthogonal to the constant functions in time, see Definition 3.1.0.9, case $cG(1)$.
The idea is to estimate $E_1$, $E_2$ and $E_3$ such that the final estimate consists of the dual solution contributions which can be further estimated by means of stability Lemma 4.1.0.5. This is proved in the following three Lemmas.

**Lemma 4.3.2.6.** If $U(0)$ is a discrete variant of the initial solution $u_0=(y_0,y_1)$ defined as

$$U(0):=(\mathcal{I}y_0,\mathcal{I}y_1),$$

where $\mathcal{I}$ is the Hermite cubic interpolant, then

$$E_1\leq C^s\Big\{\|h^s(y_0-\mathcal{I}y_0)\|_{H^1(\Omega)}\|D^{s+1}\phi\|_{L^\infty(L^2)}+\|h^s(y_1-\mathcal{I}y_1)\|_{H^s(\Omega)}\|\dot{\phi}\|_{L^\infty(H^s)}\Big\}.$$

**Proof.** We start from the following representation

$$E_1=a(y_0-\mathcal{I}y_0;\phi)+(y_1-\mathcal{I}y_1;\dot{\phi}).$$

If $s=0$, then

$$E_1\leq\|y_0-\mathcal{I}y_0\|_{H^1(\Omega)}\|\phi\|_{H^1(\Omega)}+\|y_1-\mathcal{I}y_1\|_{L^2(\Omega)}\|\dot{\phi}\|_{L^2(\Omega)}.$$

For $s=1$, an integration by parts in the first term, Hölder and the Friedrichs inequality yield

$$E_1=-(y_0-\mathcal{I}y_0;\Delta\phi)+(y_1-\mathcal{I}y_1;\dot{\phi})$$
$$\leq C\Big\{\|h(y_0-\mathcal{I}y_0\|_{H^1(\Omega)}\|\Delta\phi\|_{L^\infty(L^2)}+\|h(y_1-\mathcal{I}y_1)\|_{H^1(\Omega)}\|\dot{\phi}\|_{L^\infty(H^1)}\Big\}. \qquad\square$$

**Lemma 4.3.2.7.** There exists a constant $C$ such that for $s=0,1$

$$E_2\leq C\Big\{\|h^{s+1}\Delta(U_2-\dot{U}_1)\|_{L^1(L^2)}\|D^{s+1}\phi\|_{L^\infty(L^2)}$$

$$+\|k\mathcal{L}\mathcal{K}_h^{-1}(U_2-\overline{U}_2)\|_{L^2(H^{1-s})}(T^{1/2}\|D^{s+1}\phi\|_{L^\infty(L^2)}+\varepsilon\|D^{s+1}\dot{\phi}\|_{L^2(L^2)})\Big\}.$$

**Proof.** In order to simplify the estimation, we may rewrite $E_2$ owing to the properties of $\mathcal{J}$ such that

$$E_2 = \int_0^T a(U_2 - \dot{U}_1; \phi - \mathcal{L}\phi)dt + \int_0^T a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt.$$

An integration by parts in space in the first term and the use of (DD) boundary conditions implies

$$\int_0^T a(U_2 - \dot{U}_1; \phi - \mathcal{L}\phi)dt = -\int_0^T ((I-\mathcal{L})\Delta(U_2 - \dot{U}_1); \phi - \mathcal{L}\phi)dt$$
$$\leq C \|h^{s+1}\Delta(U_2 - \dot{U}_1)\|_{L^1(L^2)} \|D^{s+1}\phi\|_{L^\infty(L^2)}.$$

For the second term we may conclude

$$\int_0^T a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt \leq a\big(\|k(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\mathcal{L}\ddot{\phi}\|_{L^2(0,T)}\big)$$
$$\leq a\big(\|k(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\mathcal{L}\Delta(\phi - \varepsilon\dot{\phi})\|_{L^2(0,T)}\big)$$
$$\leq \big(\|k\mathcal{L}\mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\Delta(\phi - \varepsilon\dot{\phi})\|_{L^2(0,T)}\big). \qquad (4.95)$$

For $s=1$ we may conclude from the last inequality

$$\int_0^T a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt \leq \|k\mathcal{L}\mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(L^2)} \big(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)} + \varepsilon\|\Delta\dot{\phi}\|_{L^2(L^2)}\big).$$

For $s=0$, we proceed by integrating by parts in space in the last inequality in (4.95), i.e.

$$\int_0^T a(U_2 - \overline{U}_2; \mathcal{L}(\phi - \mathcal{J}\phi))dt \leq a\big(\|k\mathcal{L}\mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(0,T)}; \|\phi - \varepsilon\dot{\phi}\|_{L^2(0,T)}\big)$$
$$\leq \|k\mathcal{L}\mathcal{K}_h^{-1}(U_2 - \overline{U}_2)\|_{L^2(H^1)} \big(T^{1/2}\|\phi\|_{L^\infty(H^1)} + \varepsilon\|\dot{\phi}\|_{L^2(H^1)}\big).$$

This concludes the proof of lemma. □

**Lemma 4.3.2.8.** There exists a constant $C > 0$ such that

$$E_3 \leq C^s \Big\{ \|h^s(I-\mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^s)}$$
$$+ \|k\mathcal{L}(f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)))\|_{L^2(H^{1-s})} \big(T^{1/2}\|D^{s+1}\phi\|_{L^\infty(L^2)} + \varepsilon\|D^{s+1}\dot{\phi}\|_{L^2(L^2)}\big) \Big\}.$$

**Proof.** By using the projection properties of $\mathcal{L}$ and $\mathcal{J}$ we may write

$$E_3 = \int_0^T ((I-\mathcal{L})(f + \Delta(U_1 + \varepsilon U_2)); \dot{\phi} - \mathcal{L}\dot{\phi})dt$$
$$+ \int_0^T (f - \bar{f} + \Delta(U_1 - \overline{U}_1 + \varepsilon(U_2 - \overline{U}_2)); \mathcal{L}(\dot{\phi} - \mathcal{J}\dot{\phi}))dt.$$

The first term can be estimated such that

$$\int_0^T ((I-\mathcal{L})(f + \Delta(U_1 + \varepsilon U_2)); \dot{\phi} - \mathcal{L}\dot{\phi})dt \leq C^s \|h^s(I-\mathcal{L})(f + \Delta(U_1 + \varepsilon U_2))\|_{L^1(L^2)} \|\dot{\phi}\|_{L^\infty(H^s)}.$$

For the second term we have

$$\int_0^T (f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)));\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt$$
$$\leq \left(\|k\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)))\|_{L^2(0,T)}; \|\Delta(\phi-\varepsilon\dot{\phi})\|_{L^2(0,T)}\right).$$

If $s=1$, we may further estimate such that

$$\int_0^T (f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)));\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt$$
$$\leq \|k\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)))\|_{L^2(L^2)}\left(T^{1/2}\|\Delta\phi\|_{L^\infty(L^2)}+\varepsilon\|\Delta\dot{\phi}\|_{L^2(L^2)}\right).$$

If $s=0$, then

$$\int_0^T (f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)));\mathcal{L}(\dot{\phi}-\mathcal{J}\dot{\phi}))dt$$
$$\leq a\left(\|k\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)))\|_{L^2(0,T)}; \|\phi-\varepsilon\dot{\phi}\|_{L^2(0,T)}\right)$$
$$\leq \|k\mathcal{L}(f-\bar{f}+\Delta(U_1-\overline{U}_1+\varepsilon(U_2-\overline{U}_2)))\|_{L^2(H^1)}\left(T^{1/2}\|\phi\|_{L^\infty(L^2)}+\varepsilon\|\dot{\phi}\|_{L^2(H^1)}\right).$$

This completes the proof of Lemma.                                      □

Recalling Lemma 4.1.0.5 and error representation (4.94) we may conclude the proof of theorem.                                                                    □

# Chapter 5

# Goal-oriented error analysis

In the following, we provide the basic theory and ideas underlying the goal-oriented approach in the error analysis and mesh adaptivity process. This includes only the definitions and derivation of the a posteriori error bounds, whereby the verification in form of numerical results in not provided. The notation used within this chapter includes the one used in Chapter 3, see Table 3.1 and the additional given in Table 5.1 below.

| | | |
|---|---|---|
| $\Phi$ | exact dual solution | solution of (5.6) |
| $\Psi$ | discrete dual solution | solution of (5.7) |
| $\mathscr{J}$ | target functional | definition (5.1) |
| $\mathcal{B}_g^*$ | weak dual bilinear form | |

Table 5.1: Additional notation used in Chapter 5.

The goal is the efficient and accurate computation of certain (locally) defined quantities, so-called "target quantities" which arise from the physical formulation of the problem. They are quantified through the linear output functional

$$\mathscr{J} : \mathcal{H} \to \mathbb{R}. \tag{5.1}$$

This functional can be chosen differently, e.g. it can be the energy flux over some curve of interest, the energy of the whole system at some point in time, etc. $\mathscr{J}$, i.e. $\mathscr{J}(e) = \mathscr{J}(u) - \mathscr{J}(U)$ pretends to be the quantity of interest for the error control with $e = u - U$ and $U$, the Galerkin approximation.

From the dual formulation it can easily be seen that the computation of the target functional is closely related to the computation of the unknown continuous dual solution. In particular, the formulation of the dual problem involves the target functional as the right hand side.

The dual and energy approach in the a posteriori error analysis, see Section 3.3 and 4.3, respectively, make use of the projection and interpolation estimates when the computation of the error bound requires some exact terms, which are not necessarily a priori known. On the other hand, the goal-oriented method, instead of estimating these exact terms, replaces it by some suitable numerical approximation and then calculates an a posteriori error bound directly. To attain the optimal order of convergence, we commonly choose some approximation method based on higher order elements in space and the same order in time. It also makes sense to use higher order approximations in time. However, this approach apart from being more expensive, has to deal with a difficulty of adjusting the data between neighboring time slabs. The analysis employed here is closely related to that used in BANGERTH [9], BANGERTH-RANNACHER [10].

## 5.1    Application to the strongly damped wave equation

In case of the strongly damped wave equation, we are mainly interested in the control of the energy term at the final time point $T$ as well as in the control of the dissipative term which arises only in damped case ($\varepsilon > 0$). Hence, for some $u \in L^2(\mathcal{T}; \mathcal{H})$, we define

$$\widehat{\mathscr{J}}(u) := \frac{1}{2}\|u(T)\|_{\mathcal{H}}^2 + \varepsilon \sum_{j=1}^N \int_{I_j} \|u(t)\|_{H^1(\Omega)}^2 dt. \tag{5.2}$$

Here $\mathcal{T}$ is some arbitrary triangulation of the time domain $[0, T]$. Obviously, $\widehat{\mathscr{J}}$ is not linear in $u$. We may encounter this difficulty, by introducing the linearised target functional $\mathscr{J}$, such that

$$\mathscr{J}(v) := \langle u(T) ; v(T) \rangle_{\mathcal{H}} + 2\varepsilon \sum_{j=1}^N \int_{I_j} a(u_2; v_2) dt \quad \text{for all} \quad v \in H^1(\Omega).$$

Due to the definition of the energy norm, there holds

$$\mathscr{J}(e) = \frac{1}{2}\|u(T)\|_{\mathcal{H}}^2 - \frac{1}{2}\|U(T)\|_{\mathcal{H}}^2 + \frac{1}{2}\|e(T)\|_{\mathcal{H}}^2$$

$$+ \varepsilon \sum_{j=1}^N \int_{I_j} \|u(t)\|_{H^1(\Omega)}^2 dt - \varepsilon \sum_{j=1}^N \int_{I_j} \|U(t)\|_{H^1(\Omega)}^2 dt + \varepsilon \sum_{j=1}^N \int_{I_j} \|e(t)\|_{H^1(\Omega)}^2 dt, \tag{5.3}$$

and thus

$$\mathscr{J}(e) = \mathscr{J}(u) - \mathscr{J}(U), \text{ and} \quad \mathscr{J}(e) \to \mathscr{J}(u) - \mathscr{J}(U) \quad \text{when} \quad e \to 0 \quad \text{in} \quad \mathcal{H}. \tag{5.4}$$

From the Riesz-representation theorem, owing to the fact that $\mathscr{J}$ is some linear and bounded functional, there exists some density function $j = (j_1, j_2) \in L^2(\mathcal{T}; \mathcal{H})$, such that

$$\mathscr{J}(v) := \int_0^T \langle j ; v \rangle_{\mathcal{H}} dt \quad \text{for all} \quad v \in L^2(\mathcal{T}; \mathcal{H}). \tag{5.5}$$

Following the general concept of the goal-oriented method, we may now introduce the continuous dual problem, similar to the one defined in Section 4.1.

Then the vector form of the goal-oriented dual problem reads: Find $\Phi \in H^1(0, T; \mathcal{H})$ such that

$$-\dot{\Phi}(t, x) - \mathcal{A}^* \Phi(t, x) = j(t, x) \quad \text{on} \quad Q, \tag{5.6a}$$

$$\Phi(T, x) = 0 \qquad \text{on} \quad \Omega, \tag{5.6b}$$

where $\mathcal{A}^*$ takes the form of (4.2).

Note that the goal-oriented dual problem (5.6) has the same structure as (4.1) except for the RHS which is in this case inhomogeneous.

## 5.2    $dG(q)$ time approximation, $q = 0, 1$

For the notation and definitions used in the following, we refer to Section 2.1 and Subsection 2.3.3 where space discretisation methods and discontinuous Galerkin time approximation are

introduced.

The discrete weak dual problem reads: Find $\Psi \in \mathcal{Q}_q \subset H^1(\mathscr{T};\mathcal{H})$ such that

$$\mathcal{B}_g^*(\Psi, V) = \mathscr{J}(V) \quad \text{for all} \quad V \in \mathcal{Q}_q. \tag{5.7}$$

Here, $\mathcal{B}_g^*$ denotes the weak dual form, especially adapted for the goal oriented analysis. Namely for all $v = (v_1, v_2)$, $u = (u_1, u_2)$ sufficiently (piecewise) smooth functions in time, we define

$$\mathcal{B}_g^*(v, u) := \sum_{j=1}^N \int_{I_j} a(v_2 - v_{1\tau}; u_1) dt - \sum_{j=1}^N \int_{I_j} (v_{2,\tau}; u_2) dt - \sum_{j=1}^N \int_{I_j} a(v_1; u_2) dt$$

$$+ \varepsilon \sum_{j=1}^N \int_{I_j} a(v_2; u_2) dt - \sum_{j=2}^N a([v_1]^{j-1}; u_1^{j-1-}) - \sum_{j=2}^N ([v_2]^{j-1}; u_2^{j-1-}). \tag{5.8}$$

**Remark 5.2.0.3.** If we recall the definition of the dual bilinear form $\mathcal{B}^*$ from (4.6) in Chapter 4, we may notice that the goal-oriented dual form $\mathcal{B}_g^*$ (5.8) satisfies

$$\mathcal{B}^*(v, u) = \mathcal{B}_g^*(v, u) - \left\langle v^{N-}; u^{N-} \right\rangle_{\mathcal{H}} + \left\langle v^{0+}; u^{0-} \right\rangle_{\mathcal{H}}.$$

In Section 4.2.1, the additional terms allowed to control the error on $\|e^{N-}\|_{\mathcal{H}}$. In case of the goal-oriented dual method all terms of interest (error in the energy norm etc.) are covered by the functional $\mathcal{J}$. □

**Remark 5.2.0.4.** The solution $\Phi = (\Phi_1, \Phi_2)$ of (5.6) is also a solution of (5.7). Scalar multiplication of (5.6) with respect to $\mathcal{H}$ scalar product, by some function $V \in \mathcal{Q}$, and an integration by parts in space where the boundary conditions are imposed through the definition of the dual operator $\mathcal{A}^*$ provide (5.7). Notice that the jump terms vanish owing to the continuity of the dual solution $\Phi$ in time. □

**Lemma 5.2.0.9 (Goal-oriented error representation, *dG(q)* time approximation).**
For $u$, its discrete $dG(q)$ variant $U$, $q=0,1$ and the solution $\Phi$ of the dual problem (5.6), we have the following error representation

$$\mathscr{J}(e) = \left\langle e^{0-}; \Phi^{0+} \right\rangle_{\mathcal{H}} + Res(\Phi - V) \quad \text{for all} \quad V \in \mathcal{W}_q. \tag{5.9}$$

$Res(\Phi - V) = \mathscr{L}(\Phi - V) - \mathcal{B}(U, \Phi - V)$ is the residual with $\mathcal{B}$ as in (2.18) and $\mathscr{L}$ from (2.19).

**Proof.** The proof follows due to the properties of the discrete functions. Given (5.7), an integration by parts in time yields

$$\mathscr{J}(e) = \mathcal{B}_g^*(\Phi, e) = \sum_{j=1}^N \int_{I_j} \left\langle \Phi; e_\tau \right\rangle_{\mathcal{H}} dt - \sum_{j=1}^N \left\langle \Phi^{j-}; e^{j-} \right\rangle_{\mathcal{H}} + \sum_{j=1}^N \left\langle \Phi^{j-1+}; e^{j-1+} \right\rangle_{\mathcal{H}}$$

$$- \sum_{j=1}^N \int_{I_j} a(\Phi_1; e_2) dt + \sum_{j=1}^N \int_{I_j} a(\Phi_2; e_1) dt + \varepsilon \sum_{j=1}^N \int_{I_j} a(\Phi_2; e_2) dt$$

$$- \sum_{j=2}^N \left\langle [\Phi]^{j-1}; e^{j-1-} \right\rangle_{\mathcal{H}}. \tag{5.10}$$

The sum of the jump terms and additional $j-$ and $j-1+$ contributions can be rewritten such that

$$-\sum_{j=1}^{N}\left\{\left\langle\Phi^{j-};e^{j-}\right\rangle_{\mathcal{H}}+\left\langle\Phi^{j-1+};e^{j-1+}\right\rangle_{\mathcal{H}}\right\}-\sum_{j=2}^{N}\left\langle[\Phi]^{j-1};e^{j-1-}\right\rangle_{\mathcal{H}}$$

$$=-\sum_{j=1}^{N}\left\{\left\langle\Phi^{j-};e^{j-}\right\rangle_{\mathcal{H}}+\left\langle\Phi^{j-1+};e^{j-1+}\right\rangle_{\mathcal{H}}\right\}-\sum_{j=2}^{N}\left\{\left\langle\Phi^{j-1+};e^{j-1-}\right\rangle_{\mathcal{H}}+\left\langle\Phi^{j-1-};e^{j-1-}\right\rangle_{\mathcal{H}}\right\}$$

$$=-\sum_{j=1}^{N}\left\langle\Phi^{j-};e^{j-}\right\rangle_{\mathcal{H}}+\sum_{j=1}^{N}\left\langle[e]^{j-1};\Phi^{j-1+}\right\rangle_{\mathcal{H}}+\left\langle e^{0-};\Phi^{0+}\right\rangle_{\mathcal{H}}+\sum_{j=1}^{N-1}\left\langle\Phi^{j-};e^{j-}\right\rangle_{\mathcal{H}}$$

$$=\sum_{j=1}^{N}\left\langle[e]^{j-1};\Phi^{j-1+}\right\rangle_{\mathcal{H}}+\left\langle e^{0-};\Phi^{0+}\right\rangle_{\mathcal{H}}-\left\langle\Phi^{N-};e^{N-}\right\rangle_{\mathcal{H}}. \tag{5.11}$$

With (5.11) and the definition of the bilinear form $\mathcal{B}$, cf. (2.18), equation (5.10) simplifies to

$$\mathscr{J}(e)=\left\langle e^{0-};\Phi^{0+}\right\rangle_{\mathcal{H}}+\mathcal{B}(e,\Phi).$$

Here we used the condition (5.6b), i.e. $\left\langle\Phi^{N-};e^{N-}\right\rangle_{\mathcal{H}}=0$.
Moreover, from (2.10) and Galerkin orthogonality (2.9), we may conclude for each $V\in\mathcal{Q}$

$$\mathscr{J}(e)=\left\langle e^{0-};\Phi^{0+}\right\rangle_{\mathcal{H}}+Res(\Phi)=\left\langle e^{0-};\Phi^{0+}\right\rangle_{\mathcal{H}}+Res(\Phi-V). \qquad \square$$

Having derived the error representation (5.9), it is our aim to compute the RHS of (5.9) numerically. Since the strong dual solution $\Phi$ is unknown, the residual $Res$ is not a computable quantity. Therefore, we replace $\Phi$ by an appropriate approximation. This is emphasised in the following lemma.

**Lemma 5.2.0.10.** Let $\Phi$ be a sufficiently smooth solution of the dual problem (5.6). If we denote by $\check{\Phi}$ a discrete variant of $\Phi$ such that $\check{\Phi}$ is a polynomial of order $p'>p$ in space and $dG(q)$ function in time, then

$$\mathscr{J}(e)=\left\langle e^{0-};\check{\Phi}^{0+}\right\rangle_{\mathcal{H}}+Res(\check{\Phi}-V)+\mathcal{O}(h^{p'+p}+k^{2q+1}).$$

**Proof.** We start from (5.9). Then,

$$\mathscr{J}(e)=\left\langle e^{0-};\check{\Phi}^{0+}\right\rangle_{\mathcal{H}}+\left\langle e^{0-};(\Phi-\check{\Phi})^{0+}\right\rangle_{\mathcal{H}}+Res(\check{\Phi}-V)+Res(\Phi-\check{\Phi}). \tag{5.12}$$

It remains to prove that

$$\left\langle e^{0-};(\Phi-\check{\Phi})^{0+}\right\rangle_{\mathcal{H}}+Res(\Phi-\check{\Phi})=\mathcal{O}(h^{p'+p}+k^{2q+1}).$$

Owing to the approximation properties, we have

$$\left\langle e^{0-};(\Phi-\check{\Phi})^{0+}\right\rangle_{\mathcal{H}}=\mathcal{O}(h^{p+p'}). \tag{5.13}$$

If we assume that the mesh is quasi-uniform in space, then

$$Res(\Phi-\check{\Phi})=\mathcal{B}(e,\Phi-\check{\Phi})$$

$$=\sum_{j=1}^{N}\int_{I_j}\left\langle\dot{e};\Phi-\check{\Phi}\right\rangle_{\mathcal{H}}dt-\sum_{j=1}^{N}\int_{I_j}\left\langle\mathcal{A}e;\Phi-\check{\Phi}\right\rangle_{\mathcal{H}}dt+\sum_{j=1}^{N}\left\langle[U]^{j-1};(\Phi-\check{\Phi})^{j-1+}\right\rangle_{\mathcal{H}}$$

$$=\mathcal{O}(h^{p'+p}+k^{2q+1})+\mathcal{O}(h^{p+p'}+k^{2q+2})+\mathcal{O}(h^{p'+p}+k^{2q+1}). \tag{5.14}$$

A substitution of (5.13) and (5.14) into (5.12) yields the proof. $\square$

There are also some other ways of approximating $\Phi$ in space, e.g. by use of the biquadratic, patch-wise interpolation, see BANGERTH-RANNACHER [10, Section 4.1].

When the continuous dual solution is approximated, it is left to find a proper choice for test function $V$. This will be done for each time-space ansatz separately.

## 5.2.1   A posteriori goal-oriented error analysis, $dG(q) \otimes \mathcal{P}_1, q=0,1$

**Theorem 5.2.1.1 (A posteriori goal error estimate, $dG(q) \otimes \mathcal{P}_1$).** For $u$, its discrete $dG(q) \otimes \mathcal{P}_1$ counterpart $U$ and linearised error functional $\mathscr{J}$ there holds

1. if $q = 0$

$$
\mathscr{J}(e) = a(y_0 - \mathcal{I}y_0; \check{\Phi}_1^1) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2^1) + \sum_{j=1}^{N} \int_{I_j} \left( f; (I - \mathcal{I})\check{\Phi}_2^j \right) dt
$$

$$
+ \sum_{j=1}^{N} k_j a\left(U_2^j; (I - \mathcal{I})\check{\Phi}_1^j\right) + \sum_{j=1}^{N} \sum_{k=1}^{m} k_j [D(U_1^j + \varepsilon U_2^j)]_k (I - \mathcal{I})\check{\Phi}_2^j(x_k)
$$

$$
+ \sum_{j=1}^{N} a\left(U_1^{j-1} - U_1^j; (I - \mathcal{I})\check{\Phi}_1^j\right) dt + \sum_{j=1}^{N} \int_{I_j} (U_2^{j-1} - U_2^j; (I - \mathcal{I})\check{\Phi}_2^j) dt
$$

$$
+ \mathcal{O}(h^{p'+1} + k), \tag{5.15}
$$

2. if $q = 1$

$$
\mathscr{J}(e) = a(y_0 - \mathcal{I}y_0; \check{\Phi}_1^{1,0}) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2^{1,0}) + \sum_{j=1}^{N} \int_{I_j} \left( f; (I - \mathcal{I})(\check{\Phi}_2^{j,0} + \frac{t - t_{j-1}}{k_j}\check{\Phi}_2^{j,1}) \right) dt
$$

$$
+ \sum_{j=1}^{N} a(U_1^{j,1}; (\mathcal{I} - I)(\check{\Phi}_1^{j,0} + \frac{1}{2}\Phi_1^{j,1})) + \sum_{j=1}^{N} (U_2^{j,1}; (\mathcal{I} - I)(\check{\Phi}_2^{j,0} + \frac{1}{2}\Phi_2^{j,1}))
$$

$$
+ \sum_{j=1}^{N} k_j a\left(U_2^{j,0}; (I - \mathcal{I})(\check{\Phi}_1^{j,0} + \frac{1}{2}\check{\Phi}_1^{j,1})\right) + \sum_{j=1}^{N} k_j a\left(U_2^{j,1}; (I - \mathcal{I})(\frac{1}{2}\check{\Phi}_1^{j,0} + \frac{1}{3}\check{\Phi}_1^{j,1})\right)
$$

$$
+ \sum_{j=1}^{N} \sum_{k=1}^{m} k_j [D(U_1^{j,0} + \varepsilon U_2^{j,0})]_k \left((I - \mathcal{I})(\check{\Phi}_2^{j,0} + \frac{1}{2}\check{\Phi}_2^{j,1})\right)(x_k)
$$

$$
+ \sum_{j=1}^{N} \sum_{k=1}^{m} k_j [D(U_1^{j,1} + \varepsilon U_2^{j,1})]_k \left((I - \mathcal{I})(\frac{1}{2}\check{\Phi}_1^{j,0} + \frac{1}{3}\check{\Phi}_1^{j,1})\right)(x_k)
$$

$$
+ \sum_{j=1}^{N} a\left(U_1^{j-1,0} + U_1^{j-1,1} - U_1^{j,0}; (I - \mathcal{I})\check{\Phi}_1^{j,0}\right)
$$

$$
+ \sum_{j=1}^{N} \int_{I_j} (U_2^{j-1,0} + U_2^{j-1,1} - U_2^{j,0}; (I - \mathcal{I})\check{\Phi}_2^{j,0}) + \mathcal{O}(h^{p'+1} + k^3). \tag{5.16}
$$

**Proof.** In the following we provide only the proof for the $dG(1) \otimes \mathcal{P}_1$ case. An estimate in case of the $dG(0) \otimes \mathcal{P}_1$ discrete problem can be derived analogously and will be given only in its final form.

Given the error representation from Lemma 5.2.0.10 with the residual defined as in Lemma 3.3.2.2, case $\mathcal{P}^1$ in space, we may deduce

$$
\mathscr{J}(e) = \left\langle e^{0-}; \Phi^{0+}\right\rangle_{\mathcal{H}} + \sum_{j=1}^{N}\int_{I_j}(f, \check{\Phi}_2 - V_2)dt - \sum_{j=1}^{N}\int_{I_j}a(\dot{U}_1; \check{\Phi}_1 - V_1)dt - \sum_{j=1}^{N}\int_{I_j}(\dot{U}_2; \check{\Phi}_2 - V_2)dt
$$

$$
+ \sum_{j=1}^{N}\int_{I_j}a(U_2; \check{\Phi}_1 - V_1)dt + \sum_{j=1}^{N}\int_{I_j}\sum_{k=1}^{m}([D(U_1 + \varepsilon U_2)]_k)(\check{\Phi}_2 - V_2)(x_k)dt
$$

$$
- \sum_{j=1}^{N}a([U_1]^{j-1}; (\check{\Phi}_1 - V_1)^{j-1+}) - \sum_{j=1}^{N}([U_2]^{j-1}; (\check{\Phi}_2 - V_2)^{j-1+})
$$

$$
+ \mathcal{O}(h^{p'+1} + k^3) =: \sum_{\ell=1}^{9}E_\ell. \tag{5.17}
$$

Our aim is to compute $E_1, \ldots, E_8$.

Before we start, recall that the continuous dual solution $\Phi$ is replaced by $\check{\Phi}$ according to Lemma 5.2.0.10 where $p' > 1$. Then fix $V$ to be a nodal interpolant of $\check{\Phi}$ at the midpoint of each time interval (piecewise constant in time), i.e.

$$
V|_{I_j} := \mathcal{I}\left(\check{\Phi}^{j,0} + \frac{t - t_{j-1}}{k_j}\check{\Phi}^{j,1}\right) \quad \text{for all} \quad I_j \in \mathscr{T}. \tag{5.18}
$$

If we assume that the discrete solution $U^{0-} = (\mathcal{I}y_0, \mathcal{I}y_1)$ for the initial solution $u_0$ from (1.28b), we may deduce for $E_1$

$$
E_1 = \left\langle e^{0-}; \check{\Phi}^{0+}\right\rangle_{\mathcal{H}} = a(y_0 - \mathcal{I}y_0; \check{\Phi}_1^{1,0}) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2^{1,0}). \tag{5.19}
$$

Using (5.18),

$$
E_2 = \sum_{j=1}^{N}\int_{I_j}(f; \check{\Phi}_2 - V_2)dt = \sum_{j=1}^{N}\int_{I_j}\left(f; (I - \mathcal{I})(\check{\Phi}_2^{j,0} + \frac{t - t_{j-1}}{k_j}\check{\Phi}_2^{j,1})\right)dt, \tag{5.20}
$$

$$
E_3 = -\sum_{j=1}^{N}\int_{I_j}a(\dot{U}_1; \check{\Phi}_1 - V_1)dt = \sum_{j=1}^{N}\int_{I_j}\frac{1}{k_j}a\left(U_1^{j,1}; (\mathcal{I} - I)(\check{\Phi}_1^{j,0} + \frac{t - t_{j-1}}{k_j}\check{\Phi}_1^{j,1})\right)dt
$$

$$
= \sum_{j=1}^{N}a(U_1^{j,1}; (\mathcal{I} - I)(\check{\Phi}_1^{j,0} + \frac{1}{2}\check{\Phi}_1^{j,1})). \tag{5.21}
$$

Arguing as in case of $E_3$, we have

$$
E_4 = -\sum_{j=1}^{N}\int_{I_j}(\dot{U}_2; \check{\Phi}_2 - V_2)dt = \sum_{j=1}^{N}(U_2^{j,1}; (\mathcal{I} - I)(\check{\Phi}_2^{j,0} + \frac{1}{2}\check{\Phi}_2^{j,1})). \tag{5.22}
$$

Furthermore,

$$E_5 = \sum_{j=1}^{N} \int_{I_j} a(U_2; \check{\Phi}_1 - V_1) dt$$

$$= \sum_{j=1}^{N} \int_{I_j} a\big(U_2^{j,0} + \frac{t-t_{j-1}}{k_j} U_2^{j,1}; (I-\mathcal{I})(\check{\Phi}_1^{j,0} + \frac{t-t_{j-1}}{k_j}\check{\Phi}_1^{j,1})\big) dt$$

$$= \sum_{j=1}^{N} k_j a\big(U_2^{j,0}; (I-\mathcal{I})\big(\check{\Phi}_1^{j,0} + \frac{1}{2}\check{\Phi}_1^{j,1}\big)\big) + \sum_{j=1}^{N} k_j a\big(U_2^{j,1}; (I-\mathcal{I})(\frac{1}{2}\check{\Phi}_1^{j,0} + \frac{1}{3}\check{\Phi}_1^{j,1})\big). \quad (5.23)$$

Similarly,

$$E_6 = \sum_{j=1}^{N} \int_{I_j} \sum_{k=1}^{m} ([D(U_1+\varepsilon U_2)]_k)(\check{\Phi}_2 - V_2)(x_k) dt$$

$$= \sum_{j=1}^{N} \sum_{k=1}^{m} k_j [D(U_1^{j,0}+\varepsilon U_2^{j,0})]_k \big((I-\mathcal{I})(\check{\Phi}_2^{j,0} + \frac{1}{2}\check{\Phi}_2^{j,1})\big)(x_k)$$

$$+ \sum_{j=1}^{N} \sum_{k=1}^{m} k_j [D(U_1^{j,1}+\varepsilon U_2^{j,1})]_k \big((I-\mathcal{I})(\frac{1}{2}\check{\Phi}_2^{j,0} + \frac{1}{3}\check{\Phi}_2^{j,1})\big)(x_k). \quad (5.24)$$

We continue with the jump terms $E_7$ and $E_8$.

$$E_7 = -\sum_{j=1}^{N} a([U_1]^{j-1}; (\check{\Phi}_1 - V_1)^{j-1+}) = \sum_{j=1}^{N} a\big(U_1^{j-1,0} + U_1^{j-1,1} - U_1^{j,0}; (I-\mathcal{I})\check{\Phi}_1^{j,0}\big), \quad (5.25)$$

$$E_8 = -\sum_{j=1}^{N} ([U_2]^{j-1}; (\check{\Phi}_2 - V_2)^{j-1+}) = \sum_{j=1}^{N} (U_2^{j-1,0} + U_2^{j-1,1} - U_2^{j,0}; (I-\mathcal{I})\check{\Phi}_2^{j,0}). \quad (5.26)$$

After a discrete solution $U$ and the approximation of the strong dual solution $\check{\Phi}$ have been calculated, it is easy to compute the terms $E_1, \ldots, E_8$. This yields the proof of theorem. $\square$

## 5.2.2 A posteriori goal-oriented error analysis, $dG(q) \otimes \mathcal{C}^1, q=0,1$

**Theorem 5.2.2.1 (A posteriori goal error estimate, $dG(q) \otimes \mathcal{C}^1$).** For $u$, its discrete $dG(q) \otimes \mathcal{C}^1$ counterpart $U$ and linearised error functional $\mathscr{J}$, there holds

1. if $q=0$

$$\mathscr{J}(e) = a(y_0 - \mathcal{I}y_0; \check{\Phi}_1^1) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2^1) + \sum_{j=1}^{N} \int_{I_j} \big(f; (I-\mathcal{I})\check{\Phi}_2^j\big) dt$$

$$+ \sum_{j=1}^{N} k_j a\big(U_2^j; (I-\mathcal{I})\check{\Phi}_1^j\big) + \sum_{j=1}^{N} k_j \big(\Delta(U_1^j + \varepsilon U_2^j); (I-\mathcal{I})\check{\Phi}_2^j\big)$$

$$+ \sum_{j=1}^{N} a\big(U_1^{j-1} - U_1^j; (I-\mathcal{I})\check{\Phi}_1^j\big) dt + \sum_{j=1}^{N} \int_{I_j} (U_2^{j-1} - U_2^j; (I-\mathcal{I})\check{\Phi}_2^j) dt$$

$$+ \mathcal{O}(h^{3+p'} + k), \quad (5.27)$$

2. if $q = 1$

$$
\mathscr{J}(e) = a(y_0 - \mathcal{I}y_0; \check{\Phi}_1^{1,0}) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2^{1,0}) + \sum_{j=1}^{N} \int_{I_j} \left( f; (I - \mathcal{I})(\check{\Phi}_2^{j,0} + \frac{t - t_{j-1}}{k_j} \check{\Phi}_2^{j,1}) \right) dt
$$

$$
+ \sum_{j=1}^{N} a(U_1^{j,1}; (\mathcal{I} - I)(\check{\Phi}_1^{j,0} + \frac{1}{2}\check{\Phi}_1^{j,1})) + \sum_{j=1}^{N} (U_2^{j,1}; (\mathcal{I} - I)(\check{\Phi}_2^{j,0} + \frac{1}{2}\check{\Phi}_2^{j,1}))
$$

$$
+ \sum_{j=1}^{N} k_j a\left(U_2^{j,0}; (I - \mathcal{I})(\check{\Phi}_1^{j,0} + \frac{1}{2}\check{\Phi}_1^{j,1})\right) + \sum_{j=1}^{N} k_j a\left(U_2^{j,1}; (I - \mathcal{I})(\frac{1}{2}\check{\Phi}_1^{j,0} + \frac{1}{3}\check{\Phi}_1^{j,1})\right)
$$

$$
+ \sum_{j=1}^{N} k_j \left(\Delta(U_1^{j,0} + \varepsilon U_2^{j,0}); (I - \mathcal{I})(\check{\Phi}_2^{j,0} + \frac{1}{2}\check{\Phi}_2^{j,1})\right)
$$

$$
+ \sum_{j=1}^{N} k_j \left(\Delta(U_1^{j,1} + \varepsilon U_2^{j,1}); (I - \mathcal{I})(\frac{1}{2}\check{\Phi}_2^{j,0} + \frac{1}{3}\check{\Phi}_2^{j,1})\right)
$$

$$
+ \sum_{j=1}^{N} a\left(U_1^{j-1,0} + U_1^{j-1,1} - U_1^{j,0}; (I - \mathcal{I})\check{\Phi}_1^{j,0}\right)
$$

$$
+ \sum_{j=1}^{N} \int_{I_j} (U_2^{j-1,0} + U_2^{j-1,1} - U_2^{j,0}; (I - \mathcal{I})\check{\Phi}_2^{j,0}) + \mathcal{O}(h^{3+p'} + k^3). \quad (5.28)
$$

**Proof.** Same as in the proof of Theorem 5.2.2.1, we derive only the proof for the case $dG(1)$ in time.

Given the error representation from Lemma 5.2.0.10 with the residual as in Lemma 3.3.2.2, case $\mathcal{C}^1$ in space, we have

$$
\mathscr{J}(e) = \left\langle e^{0-}; \check{\Phi}^{0+} \right\rangle_{\mathcal{H}} + \sum_{j=1}^{N} \int_{I_j} (f, \check{\Phi}_2 - V_2) dt - \sum_{j=1}^{N} \int_{I_j} a(\dot{U}_1; \check{\Phi}_1 - V_1) dt - \sum_{j=1}^{N} \int_{I_j} (\dot{U}_2; \check{\Phi}_2 - V_2) dt
$$

$$
+ \sum_{j=1}^{N} \int_{I_j} a(U_2; \check{\Phi}_1 - V_1) dt + \sum_{j=1}^{N} \int_{I_j} (\Delta(U_1 + \varepsilon U_2); \check{\Phi}_2 - V_2) dt
$$

$$
- \sum_{j=1}^{N} a([U_1]^{j-1}; (\check{\Phi}_1 - V_1)^{j-1+}) - \sum_{j=1}^{N} ([U_2]^{j-1}; (\check{\Phi}_2 - V_2)^{j-1+})
$$

$$
+ \mathcal{O}(h^{p'+3} + k^3) =: \sum_{\ell=1}^{9} E_\ell. \quad (5.29)
$$

The idea in the following is to compute $E_1, \dots, E_8$. First, let a discrete function $V$ be defined as in (5.18) where $\mathcal{I}$ stands for the cubic Hermite interpolation operator in space, i.e. $p = 3$ and then $p' > 4$. Then, we may see that for $E_1 - E_5$ and $E_7 - E_8$ there hold the analogous estimates to the one derived for $\mathcal{P}_1$ ansatz in space, see (5.19)–(5.23) and (5.25)–(5.26), respectively.

With the following approximation for $E_6$, see (5.29), we conclude the proof of theorem.

$$E_6 = -\sum_{j=1}^{N} \int_{I_j} (\Delta(U_1 + \varepsilon U_2); \check{\Phi}_2 - V_2) dt = \sum_{j=1}^{N} k_j \Big( \Delta(U_1^{j,0} + \varepsilon U_2^{j,0}); (\mathcal{I} - I)(\check{\Phi}_1^{j,0} + \frac{1}{2}\check{\Phi}_1^{j,1}) \Big)$$

$$+ \sum_{j=1}^{N} k_j \Big( \Delta(U_1^{j,1} + \varepsilon U_2^{j,1}); (\mathcal{I} - I)(\frac{1}{2}\check{\Phi}_1^{j,0} + \frac{1}{3}\check{\Phi}_1^{j,1}) \Big). \quad \square$$

## 5.3  *cG*(1) time approximation

For the notation and definitions used in the following, we refer to Section 2.1 and Subsection 2.3.4 where space discretisation methods and continuous Galerkin time approximation are introduced.

The discrete weak dual problem reads: Find $\Psi \in \mathcal{Q}_c$ such that

$$\mathcal{B}_g^*(\Phi, V) = \mathcal{J}(V) \quad \text{for all} \quad V \in W_c \subset L^2(\mathcal{T}; \mathcal{H}). \tag{5.30}$$

Here, $\mathcal{B}_g^*$ denotes the weak dual form adapted for the goal-oriented analysis. Namely, for all $v = (v_1, v_2), u = (u_1, u_2)$ sufficiently (piecewise) smooth functions in time, $\mathcal{B}_g^*$ is defined as

$$\mathcal{B}_g^*(v, u) := \int_0^T a(v_2 - \dot{v}_1; u_1) dt - \int_0^T (\dot{v}_2; u_2) dt - \int_0^T a(v_1; u_2) dt + \varepsilon \int_0^T a(v_2; u_2) dt. \tag{5.31}$$

**Remark 5.3.0.1.** Dual bilinear form $\mathcal{B}^*$ from (4.34) and the goal-oriented dual form $\mathcal{B}_g^*$ defined in (5.31) satisfy

$$\mathcal{B}^*(v, u) = \mathcal{B}_g^*(v, u) + \big\langle v(T); u(T) \big\rangle_{\mathcal{H}} - \big\langle v(0); u(0) \big\rangle_{\mathcal{H}}. \qquad \square$$

**Remark 5.3.0.2.** The solution $\Phi$ of (5.6) solves (5.30). $\qquad \square$

**Lemma 5.3.0.1 (Goal-oriented error representation, *cG*(1) time approximation).**
For $u$, its discrete *cG*(1) variant $U$ and the solution $\Phi$ of the dual problem (5.6) we have the following error representation

$$\mathcal{J}(e) = \big\langle e(0); \Phi(0) \big\rangle_{\mathcal{H}} + Res(\Phi - V) \quad \text{for all} \quad V \in \mathcal{W}_c. \tag{5.32}$$

Here, the residual $Res(\Phi - V) = \mathcal{L}(\Phi - V) - \mathcal{B}(U, \Phi - V)$ with $\mathcal{B}$ from (2.41) and $\mathcal{L}$ from (2.42).

**Proof.** The error $e$ is globally continuous in time. Given (5.30) and (5.31) we have by integrating by parts in time

$$\mathcal{J}(e) = \mathcal{B}^*(\Phi, e) = \int_0^T a(\dot{e}_1; \Phi_1) dt - a(e_1(T); \Phi_1(T)) + a(e_1^{0-}; \Phi_1(0))$$

$$+ \int_0^T (\dot{e}_2; \Phi_2) dt - (e_2(T); \Phi_2(T)) + (e_2^{0-}; \Phi_2(0))$$

$$- \int_0^T a(e_2; \Phi_1) dt + \int_0^T a(e_1; \Phi_2) dt + \varepsilon \int_0^T a(e_2; \Phi_2) dt.$$

Since $\Phi(T)=0$ if we recall the definition of the bilinear form $\mathcal{B}$, see (2.41), we may deduce

$$\mathscr{J}(e)=\langle e(0)\,;\Phi(0)\rangle_{\mathcal{H}}+\mathcal{B}(e,\Phi). \tag{5.33}$$

Finally, from the residual representation (2.10) and the Galerkin orthogonality (2.9) we have for all $V\in W_c$

$$\mathscr{J}(e)=\langle e(0)\,;\Phi(0)\rangle_{\mathcal{H}}+Res(\Phi)=\langle e(0)\,;\Phi(0)\rangle_{\mathcal{H}}+Res(\Phi-V). \qquad\qquad \square$$

Having derived the error representation in terms of the functional of interest (5.32), the idea is to compute the RHS of the same. Since $\Phi$ is unknown we approximate it by some discrete function. This is emaphasized in the following lemma.

**Lemma 5.3.0.2.** Let $\Phi$ be a a sufficiently smooth solution of problem (5.6). Let $\check{\Phi}$ be a discrete variant of $\Phi$ such that $\check{\Phi}$ is a polynomial of order $p'>p$ in space and $cG(1)$ function in time, then there holds

$$\mathscr{J}(e)=\langle e(0)\,;\check{\Phi}(0)\rangle_{\mathcal{H}}+Res(\check{\Phi}-V)+\mathcal{O}(h^{p'+p}+k^3)$$

**Proof.** We start from (5.32). Then,

$$\mathscr{J}(e)=\langle e(0)\,;\check{\Phi}(0)\rangle_{\mathcal{H}}+\langle e(0)\,;(\Phi-\check{\Phi})(0)\rangle_{\mathcal{H}}+Res(\check{\Phi}-V)+Res(\Phi-\check{\Phi}). \tag{5.34}$$

We need to prove that

$$\langle e(0)\,;(\Phi-\check{\Phi})(0)\rangle_{\mathcal{H}}+Res(\Phi-\check{\Phi})=\mathcal{O}(h^{p+p'}+k^3).$$

Owing to the approximation properties of both time and space discretisation method we have

$$\langle e(0)\,;(\Phi-\check{\Phi})(0)\rangle_{\mathcal{H}}=\mathcal{O}(h^{p+p'}). \tag{5.35}$$

Using the same arguments as above

$$Res(\Phi-\check{\Phi})=\int_0^T\langle\dot{e}\,;\Phi-\check{\Phi}\rangle_{\mathcal{H}}dt-\int_0^T\langle\mathcal{A}e\,;\Phi-\check{\Phi}\rangle_{\mathcal{H}}dt$$
$$=\mathcal{O}(h^{p+p'}+k^3)+\mathcal{O}(h^{p+p'}+k^4). \tag{5.36}$$

This completes the proof. $\qquad\qquad \square$

Furthermore, let a test function $V$ be chosen as the interpolation of $\check{\Phi}$ at the midpoint of each time interval (piecewise constant in time), i.e.

$$V|_{I_j}:=\frac{1}{2}\mathcal{I}\left(\check{\Phi}(t_j)+\check{\Phi}(t_{j-1})\right)\quad\text{for all}\quad I_j\in\mathscr{T}. \tag{5.37}$$

Note that the interpolation operator $\mathcal{I}$ need to be applied interval-wise in time according to the general idea of the discretisation where the grid may vary in space, from one time slab to the another. However, in case of the $cG(1)$ method in time, we need to assume that $\mathcal{S}^{j-1}\subseteq\mathcal{S}^j$.

### 5.3.1 A posteriori goal-oriented error analysis, $cG(1)\otimes\mathcal{P}_1$

**Theorem 5.3.1.1 (A posteriori goal error estimate, $cG(1)\otimes\mathcal{P}^1$).** For $u$, its discrete $cG(1)\otimes\mathcal{P}_1$ counterpart $U$ and linearised error functional $\mathscr{J}$ there holds

$$
\begin{aligned}
\mathscr{J}(e) = {}& \mathcal{O}(h^{p'+1}+k^3) + a(y_0 - \mathcal{I}y_0; \check{\Phi}_1(0)) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2(0)) \\
& + \sum_{j=1}^{N} \int_{I_j} \left( f; \frac{1}{k_j}\left((t-t_{j-1})\check{\Phi}_2(t_j) + (t_j-t)\check{\Phi}_2(t_{j-1})\right) - \frac{1}{2}\mathcal{I}\left(\check{\Phi}_2(t_j)+\check{\Phi}_2(t_{j-1})\right) \right) dt \\
& + \sum_{j=1}^{N} \frac{1}{2} a\left( U_1(t_j) - U_1(t_{j-1}); (\mathcal{I}-I)\left(\check{\Phi}_1(t_j)+\check{\Phi}_1(t_{j-1})\right) \right) \\
& + \sum_{j=1}^{N} \frac{1}{2}\left( U_2(t_j) - U_2(t_{j-1}); (\mathcal{I}-I)\left(\check{\Phi}_2(t_j)+\check{\Phi}_2(t_{j-1})\right) \right) \\
& + \sum_{j=1}^{N} \frac{k_j}{4} a\left( U_2(t_j) + U_2(t_{j-1}); (I-\mathcal{I})(\check{\Phi}_1(t_j)+\check{\Phi}_1(t_{j-1})) \right) \\
& + \sum_{j=1}^{N} \frac{k_j}{12} a\left( U_2(t_j) - U_2(t_{j-1}); \check{\Phi}_1(t_j)-\check{\Phi}_1(t_{j-1}) \right) \\
& + \sum_{j=1}^{N}\sum_{k=1}^{m} \frac{k_j}{4}[D(U_1(t_j)+U_1(t_{j-1})+\varepsilon(U_2(t_j)+U_2(t_{j-1})))]_k (I-\mathcal{I})\left(\check{\Phi}_1(t_j)+\check{\Phi}_1(t_{j-1})\right)(x_k) \\
& + \sum_{j=1}^{N}\sum_{k=1}^{m} \frac{k_j}{12}[D(U_1(t_j)-U_1(t_{j-1})+\varepsilon(U_2(t_j)-U_2(t_{j-1})))]_k (\check{\Phi}_2(t_j)-\check{\Phi}_2(t_{j-1}))(x_k). \quad (5.38)
\end{aligned}
$$

**Proof.** Given the error representation from Lemma 5.3.0.2 with residual as in Lemma 3.3.3.2, case $\mathcal{P}_1$ in space, we may deduce

$$
\begin{aligned}
\mathscr{J}(e) = {}& \left\langle e(0); \check{\Phi}(0)\right\rangle_{\mathcal{H}} + \int_0^T (f, \Phi_2 - V_2)dt - \int_0^T a(\dot{U}_1; \check{\Phi}_1 - V_1)dt - \int_0^T (\dot{U}_2; \check{\Phi}_2 - V_2)dt \\
& + \int_0^T a(U_2; \check{\Phi}_1 - V_1)dt - \int_0^T \sum_{k=1}^m [D(U_1+\varepsilon U_2)]_k (\check{\Phi}_2 - V_2)(x_k)dt \\
& + \mathcal{O}(h^{1+p'}+k^3) =: \sum_{\ell=1}^{7} E_\ell. \quad (5.39)
\end{aligned}
$$

The idea is to compute $E_1 - E_6$. If we assume that the discrete variant of the initial solution $u_0$ reads $U(0) = (\mathcal{I}y_0, \mathcal{I}y_1)$, see (1.28b), then we may conclude

$$
E_1 = \left\langle e(0); \check{\Phi}(0)\right\rangle_{\mathcal{H}} = \left\langle u_0 - U(0); \check{\Phi}(0)\right\rangle_{\mathcal{H}} = a(y_0 - \mathcal{I}y_0; \check{\Phi}_1(0)) + a(y_1 - \mathcal{I}y_2; \check{\Phi}_2(0)). \quad (5.40)
$$

Notice that within this subsection $\mathcal{I}$ stands for the nodal interpolation operator in space, cf. Definition 3.1.0.1.

Furthermore,

$$E_2 = \int_0^T (f; \check{\Phi}_2 - V_2) dt$$

$$= \sum_{j=1}^N \int_{I_j} \left( f; \frac{t-t_{j-1}}{k_j} \check{\Phi}_2(t_j) + \frac{t_j-t}{k_j} \check{\Phi}_2(t_{j-1}) - \frac{1}{2} \mathcal{I}(\check{\Phi}_2(t_j) + \check{\Phi}_2(t_{j-1})) \right) dt. \qquad (5.41)$$

Using the same arguments,

$$E_3 = -\int_0^T a(\dot{U}_1; \check{\Phi}_1 - V_1) dt$$

$$= \sum_{j=1}^N \int_{I_j} a\left( \frac{1}{k_j}(U_1(t_j) - U_1(t_{j-1})); \frac{1}{2}\mathcal{I}(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) - \frac{t-t_{j-1}}{k_j}\check{\Phi}_1(t_j) - \frac{t_j-t}{k_j}\check{\Phi}_1(t_{j-1}) \right) dt$$

$$= \sum_{j=1}^N \frac{1}{2}a\left( U_1(t_j) - U_1(t_{j-1}); (\mathcal{I}-I)(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right). \qquad (5.42)$$

Similarly,

$$E_4 = -\int_0^T (\dot{U}_2; \check{\Phi}_2 - V_2) dt = \sum_{j=1}^N \frac{1}{2}\left( U_2(t_j) - U_2(t_{j-1}); (\mathcal{I}-I)(\check{\Phi}_2(t_j) + \check{\Phi}_2(t_{j-1})) \right). \qquad (5.43)$$

In case of $E_5$ we have

$$E_5 = \int_0^T a(U_2; \Phi_1 - V_1) dt$$

$$= \sum_{j=1}^N \int_{I_j} \left( \frac{t-t_{j-1}}{k_j}U_2(t_j) + \frac{t_j-t}{k_j}U_2(t_{j-1}); \frac{t-t_{j-1}}{k_j}\check{\Phi}_1(t_j) + \frac{t_j-t}{k_j}\check{\Phi}_1(t_{j-1}) \right) dt$$

$$- \sum_{j=1}^N \int_{I_j} \frac{1}{2}a\left( \frac{t-t_{j-1}}{k_j}U_2(t_j) + \frac{t_j-t}{k_j}U_2(t_{j-1}); \mathcal{I}(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right) dt$$

$$= \sum_{j=1}^N \frac{k_j}{4}a\left( U_2(t_j) + U_2(t_{j-1}); (I-\mathcal{I})(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right)$$

$$+ \sum_{j=1}^N \frac{k_j}{12}a\left( U_2(t_j) - U_2(t_{j-1}); \check{\Phi}_1(t_j) - \check{\Phi}_1(t_{j-1}) \right). \qquad (5.44)$$

Likewise,

$$E_6 = \int_0^T \sum_{k=1}^m [D(U_1 + \varepsilon U_2)]_k(\check{\Phi}_2 - V_2)(x_k) dt$$

$$= \sum_{j=1}^N \sum_{k=1}^m \frac{k_j}{4}[D(U_1(t_j) + U_1(t_{j-1})) + \varepsilon D(U_2(t_j) + U_2(t_{j-1})))]_k(I-\mathcal{I})(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1}))(x_k)$$

$$+ \sum_{j=1}^N \sum_{k=1}^m \frac{k_j}{12}[D(U_1(t_j) - U_1(t_{j-1})) + \varepsilon D(U_2(t_j) - U_2(t_{j-1})))]_k(\check{\Phi}_1(t_j) - \check{\Phi}_1(t_{j-1}))(x_k). \qquad (5.45)$$

A substitution of the approximations for $E_1 - E_6$ into (5.39) yields the proof of theorem. $\quad\square$

## 5.3.2 A posteriori goal-oriented error analysis, $cG(1) \otimes \mathcal{C}^1$

**Theorem 5.3.2.1 (A posteriori goal error estimate, $cG(1) \otimes \mathcal{C}^1$).** For $u$, its discrete $cG(1) \otimes$ $\mathcal{C}_1$ counterpart $U$ and linearised error functional $\mathscr{J}$ there holds

$$\mathscr{J}(e) = \mathcal{O}(h^{p'+3} + k^3) + a(y_0 - \mathcal{I}y_0; \check{\Phi}_1(0)) + (y_1 - \mathcal{I}y_1; \check{\Phi}_2(0))$$

$$+ \sum_{j=1}^{N} \int_{I_j} \left( f; \frac{1}{k_j} \left( (t - t_{j-1}) \check{\Phi}_2(t_j) + (t_j - t) \check{\Phi}_2(t_{j-1}) \right) - \frac{1}{2} \mathcal{I} \left( \check{\Phi}_2(t_j) + \check{\Phi}_2(t_{j-1}) \right) \right) dt$$

$$+ \sum_{j=1}^{N} \frac{1}{2} a \left( U_1(t_j) - U_1(t_{j-1}); (\mathcal{I} - I)(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right)$$

$$+ \sum_{j=1}^{N} \frac{1}{2} \left( U_2(t_j) - U_2(t_{j-1}); (\mathcal{I} - I)(\check{\Phi}_2(t_j) + \check{\Phi}_2(t_{j-1})) \right)$$

$$+ \sum_{j=1}^{N} \frac{k_j}{4} a \left( U_2(t_j) + U_2(t_{j-1}); (I - \mathcal{I})(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right)$$

$$+ \sum_{j=1}^{N} \frac{k_j}{12} a \left( U_2(t_j) - U_2(t_{j-1}); \check{\Phi}_1(t_j) - \check{\Phi}_1(t_{j-1}) \right)$$

$$+ \sum_{j=1}^{N} \frac{k_j}{4} \left( \Delta(U_1(t_j) + U_1(t_{j-1}) + \varepsilon(U_2(t_j) + U_2(t_{j-1}))); (I - \mathcal{I})(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right)$$

$$+ \sum_{j=1}^{N} \frac{k_j}{12} \left( \Delta(U(t_j) - U_1(t_{j-1}) + \varepsilon(U_2(t_j) - U_2(t_{j-1}))); \check{\Phi}_1(t_j) - \check{\Phi}_1(t_{j-1}) \right). \qquad (5.46)$$

**Proof.** Given the error representation Lemma 5.3.0.2 with the residual as in Lemma 3.3.3.2, case $\mathcal{C}^1$ elements, we may deduce

$$\mathscr{J}(e) = \left\langle e(0); \check{\Phi}(0) \right\rangle_{\mathcal{H}} + \int_0^T (f, \check{\Phi}_2 - V_2) dt - \int_0^T a(\dot{U}_1; \check{\Phi}_1 - V_1) dt - \int_0^T (\dot{U}_2; \check{\Phi}_2 - V_2) dt$$

$$+ \int_0^T a(U_2; \Phi_1 - V_1) dt + \int_0^T (\Delta(U_1 + \varepsilon U_2); \check{\Phi}_2 - V_2) dt + \mathcal{O}(h^{p'+3} + k^3) =: \sum_{\ell=1}^{7} E_\ell. \quad (5.47)$$

The idea is to find a suitable approximations for $E_1 - E_6$. Obviously, if we choose a test function $V$ as in (5.37) then we may derive an analogous estimates for $E_1 - E_5$ as in (5.40)-(5.44) where $\mathcal{I}$ stands for the cubic Hermite interpolation operator in space. To complete the proof of theorem, its left to estimate $E_6$.

$$E_6 = \int_0^T (\Delta(U_1 + \varepsilon U_2); \check{\Phi}_2 - V_2) dt$$

$$= \sum_{j=1}^{N} \frac{k_j}{4} \left( \Delta(U_1(t_j) + U_1(t_{j-1}) + \varepsilon(U_2(t_j) + U_2(t_{j-1}))); (I - \mathcal{I})(\check{\Phi}_1(t_j) + \check{\Phi}_1(t_{j-1})) \right)$$

$$+ \sum_{j=1}^{N} \frac{k_j}{12} \left( \Delta(U(t_j) - U_1(t_{j-1}) + \varepsilon(U_2(t_j) - U_2(t_{j-1}))); \check{\Phi}_1(t_j) - \check{\Phi}_1(t_{j-1}) \right). \qquad (5.48)$$

A substition of approximation for $E_1 - E_6$ into (5.47), yields the proof of theorem. $\qquad \square$

# Chapter 6

# Numerical experiments

## 6.1 FE solution - Algorithmen

Our intention within this section is to determine the algorithm for the computation of the discrete solution on each time interval $I_j$. This algorithm is based on the localised weak form (2.6). We consider only the Galerkin methods for time discretisation introduced in Subsections 2.3.3 and 2.3.4, and spatial discretisation by $\mathcal{P}_1$ and $\mathcal{C}^1$ elements, see Section 2.1. The algorithm for the computation of the semi-discrete solution is already presented in Subsection 2.4.1.3 and Subsection 2.4.2.3, for particular space ansatz, respectively.

The analysis follows by considering first the different time approximations and thereupon the spatial approximation.

### 6.1.1 FE solution, $dG(q)$ time approximation, $q=0,1$

Below we only consider the case $q = 1$, whereas for $q = 0$ we provide a final result without going into details.

#### 6.1.1.1 FE Solution, $dG(q) \otimes \mathcal{P}_1$

Given the equation (2.53), a substitution of four different variants for the test function $V^j \in \mathcal{Q}_1^j$ yields the following systems,

1. for $V^j = (V_1^j, 0) = (\phi_m, 0)$, $m = 0, \dots, n$

$$\sum_{k=0}^{n}(U_{k,1}^{j,0}+U_{k,1}^{j,1})S_{m,k} - \sum_{k=0}^{n}(k_j U_{k,2}^{j,0}+\frac{k_j}{2}U_{k,2}^{j,1})S_{m,k} = G_{1,m}^j,$$

2. for $V^j = (V_1^j, 0) = ((t-t_{j-1})/k_j\phi_m, 0)$, $m = 0, \dots, n$

$$\sum_{k=0}^{n}\frac{1}{2}U_{k,1}^{j,1}S_{m,k} - \sum_{k=0}^{n}(\frac{k_j}{2}U_{k,2}^{j,0}+\frac{k_j}{3}U_{k,2}^{j,1})S_{m,k} = 0,$$

3. for $V^j = (0, V_2^j) = (0, \phi_m)$, $m = 0, \ldots, n$

$$\sum_{k=0}^{n}(k_j U_{k,1}^{j,0} + \frac{k_j}{2} U_{k,1}^{j,1}) S_{m,k} + \varepsilon \sum_{k=0}^{n}(k_j U_{k,2}^{j,0} + \frac{k_j}{2} U_{k,2}^{j,1}) S_{m,k} + \sum_{k=0}^{n}(U_{k,2}^{j,0} + U_{k,2}^{j,1}) M_{m,k}$$

$$= \int_{I_j}(f; \phi_m) dt + G_{m,2}^{j},$$

4. for $V^j = (0, V_2^j) = (0, (t - t_{j-1})/k_j \phi_m)$, $m = 0, \ldots, n$

$$\sum_{k=0}^{n}(\frac{k_j}{2} U_{k,1}^{j,0} + \frac{k_j}{3} U_{k,1}^{j,1}) S_{m,k} + \varepsilon \sum_{k=1}^{n}(\frac{k_j}{2} U_{k,2}^{j,0} + \frac{k_j}{3} U_{k,2}^{j,1}) S_{m,k} + \sum_{k=0}^{n}\frac{1}{2} U_{k,2}^{j,1} M_{m,k}$$

$$= \int_{I_j}(f; \frac{t - t_{j-1}}{k_j}\phi_m) dt.$$

Here, for all $m = 0, \ldots, n$,

$$G_{m,1}^{j} := \begin{cases} \sum_{k=0}^{n}(U_{k,1}^{j-1,0} + U_{k,1}^{j-1,1}) S_{m,k}, & \text{for } j > 1, \\ (y_0; \phi_m), & \text{for } j = 1. \end{cases} \tag{6.1a}$$

$$G_{m,2}^{j} := \begin{cases} \sum_{k=0}^{n}(U_{k,2}^{j-1,0} + U_{k,2}^{j-1,1}) M_{m,k} & \text{for } j > 1, \\ (y_1; \phi_m) & \text{for } j = 1. \end{cases} \tag{6.1b}$$

If we denote by

$$F_0 := \begin{bmatrix} \int_{I_j}(f, \phi_0) dt \\ \vdots \\ \int_{I_j}(f, \phi_n) dt \end{bmatrix} \quad \text{and} \quad F_1 := \begin{bmatrix} \int_{I_j}(f, {}^{t-t_{j-1}}/k_j \phi_0) dt \\ \vdots \\ \int_{I_j}(f, {}^{t-t_{j-1}}/k_j \phi_n) dt \end{bmatrix}, \tag{6.2}$$

owing to the vector notation from Table 2.5, the systems above read in the equivalent vector form

$$\mathbb{L}\mathbb{U}^j = \mathbb{F}, \tag{6.3}$$

where $\mathbb{L}$ is a $4(n+1) \times 4(n+1)$ block matrix and $\mathbb{F}$ is a $4(n+1) \times 1$ block vector of the form

$$\mathbb{L} := \begin{bmatrix} S & S & -k_j S & -k_j/2 S \\ 0 & 1/2 S & -k_j/2 S & -k_j/3 S \\ k_j S & k_j/2 S & \varepsilon k_j S + M & \varepsilon k_j/2 S + M \\ k_j/2 S & k_j/3 S & \varepsilon k_j/2 S & \varepsilon k_j/3 S + 1/2 M \end{bmatrix}, \quad \mathbb{F} := \begin{bmatrix} G_1^j \\ \mathbf{0} \\ F_0 + G_2^j \\ F_1 \end{bmatrix}. \tag{6.4}$$

Here, $G_1^j, G_2^j$ stand for the global matrices (6.1).

In case of $q = 0$, we solve a system (6.3) with a solution vector $\mathbb{U}^j := \begin{bmatrix} \mathbb{U}_1^j \\ \mathbb{U}_2^j \end{bmatrix}$ where $\mathbb{U}_1^j, \mathbb{U}_2^j \in \mathbb{R}^n$

coincides with $\mathbb{U}_1^{j,0}, \mathbb{U}_2^{j,0}$ from above, and $2(n+1) \times 2(n+1)$ matrix $\mathbb{L}$ and $2(n+1) \times 1$ dimensional vector $\mathbb{F}$ where

$$\mathbb{L} := \begin{bmatrix} S & -k_j S \\ k_j S & \varepsilon k_j S + M \end{bmatrix}, \quad \mathbb{F} := \begin{bmatrix} G_1^j \\ F_0 + G_2^j \end{bmatrix}, \tag{6.5}$$

and $G_1^j := S\mathbb{U}_1^{j-1}, G_2^j := M\mathbb{U}_2^{j-1}$.

For both $q = 0, 1$, $\mathbb{U}^0$ denotes the discrete variant of the initial solution $u_0$.

### 6.1.1.2 FE solution, $dG(q) \otimes \mathcal{C}^1$

In the following we only consider the case $q = 0$. Given the weak equation (2.61), for different choice of test function $V \in \mathcal{Q}_1^j$ a sum over $k = 1, \ldots, n$ of the corresponding vector form (2.62) yields four systems. Namely,

1. for $V^j = (V_1^j, 0) = (\phi_p, 0)$, $1 \leq p \leq 4$

$$\sum_{k=1}^n S_k(\mathbb{U}_{k,1}^{j,0} + \mathbb{U}_{k,1}^{j,1}) - \sum_{k=1}^n S_k(k_j \mathbb{U}_{k,2}^{j,0} + \frac{k_j}{2}\mathbb{U}_{k,2}^{j,1}) = \sum_{k=0}^n G_{k,1}^j,$$

2. for $V^j = (V_1^j, 0) = ((t - t_{j-1})/k_j \phi_p, 0)$, $1 \leq p \leq 4$

$$\sum_{k=1}^n \frac{1}{2}S_k\mathbb{U}_{k,1}^{j,1} - \sum_{k=1}^n S_k(\frac{k_j}{2}\mathbb{U}_{k,2}^{j,0} + \frac{k_j}{3}\mathbb{U}_{k,2}^{j,1}) = 0,$$

3. for $V^j = (0, V_2^j) = (0, \phi_p)$, $1 \leq p \leq 4$

$$\sum_{k=1}^n S_k(k_j \mathbb{U}_{k,1}^{j,0} + \frac{k_j}{2}\mathbb{U}_{k,1}^{j,1}) + \varepsilon \sum_{k=1}^n S_k(k_j \mathbb{U}_{k,2}^{j,0} + \frac{k_j}{2}\mathbb{U}_{k,2}^{j,1}) + \sum_{k=1}^n M_k(\mathbb{U}_{k,2}^{j,0} + \mathbb{U}_{k,2}^{j,1}) = \sum_{k=1}^n F_k^{j,0} + G_{k,2}^j,$$

4. for $V^j = (0, V_2^j) = (0, (t - t_{j-1})/k_j \phi_p)$, $1 \leq p \leq 4$

$$\sum_{k=1}^n S_k(\frac{k_j}{2}\mathbb{U}_{k,1}^{j,0} + \frac{k_j}{3}\mathbb{U}_{k,1}^{j,1}) + \varepsilon \sum_{k=1}^n S_k(\frac{k_j}{2}\mathbb{U}_{k,2}^{j,0} + \frac{k_j}{3}\mathbb{U}_{k,2}^{j,1}) + \sum_{k=1}^n \frac{1}{2}M_k\mathbb{U}_{k,2}^{j,1} = \sum_{k=1}^n F_k^{j,1}.$$

The coefficient matrices $\mathbb{U}_{k,1|2}^{j,0|1}$ have the same form as in Table 2.6. Note also that for these choices of test functions, coefficient matrices $\mathbb{V}_{k,1|2}^{j,0|1} = e_p \in \mathbb{R}^4$. $S_k, M_k$ are the $4 \times 4$ element stiffness and mass matrix defined in (2.58), (2.59). In particular they read for $k = 1, \ldots, n$

$$S_k = \frac{1}{30h_k} \begin{bmatrix} 36 & 3h_k & -36 & 3h_k \\ 3h_k & 4(h_k)^2 & -3h_k & -(h_k)^2 \\ -36 & -3h_k & 36 & -3h_k \\ 3h_k & -(h_k)^2 & -3h_k & 4(h_k)^2 \end{bmatrix}, \tag{6.6}$$

$$M_k = \frac{h_k}{420} \begin{bmatrix} 156 & 22h_k & 54 & -13h_k \\ 22h_k & 4(h_k)^2 & 13h_k & -3(h_k)^2 \\ 54 & 13h_k & 156 & -22h_k \\ -13h_k & -3(h_k)^2 & -22h_k & 4(h_k)^2 \end{bmatrix}. \tag{6.7}$$

Furthermore, $F_k^{j,0}, F_k^{j,1}$ are the $4 \times 1$ element force vectors defined as

$$F_k^{j,0} := \begin{bmatrix} h_k/2 \int_{I_j} \int_{-1}^1 f\phi_1 \, dt d\zeta \\ h_k^2/4 \int_{I_j} \int_{-1}^1 f\phi_2 \, dt d\zeta \\ h_k/2 \int_{I_j} \int_{-1}^1 f\phi_3 \, dt d\zeta \\ h_k^2/4 \int_{I_j} \int_{-1}^1 f\phi_4 \, dt d\zeta \end{bmatrix}, F_k^{j,1} := \begin{bmatrix} h_k/2 \int_{I_j} \int_{-1}^1 {}^{t-t_{j-1}}/k_j f\phi_1 \, dt d\zeta \\ h_k^2/4 \int_{I_j} \int_{-1}^1 {}^{t-t_{j-1}}/k_j f\phi_2 \, dt d\zeta \\ h_k/2 \int_{I_j} \int_{-1}^1 {}^{t-t_{j-1}}/k_j f\phi_3 \, dt d\zeta \\ h_k^2/4 \int_{I_j} \int_{-1}^1 {}^{t-t_{j-1}}/k_j f\phi_4 \, dt d\zeta \end{bmatrix}. \tag{6.8}$$

The $4 \times 1$ element mass vectors $G_{k,1}^j, G_{k,2}^j$ are related to the terms which include the vector form of the discrete solution from the previous time step, i.e. $\mathbb{U}_{k,1}^{j-1}$ and $\mathbb{U}_{k,2}^{j-1}$ respectively, such that

$$G_{k,1}^j := S_k(\mathbb{U}_{k,1}^{j-1,0} + \mathbb{U}_{k,1}^{j-1,1}), \tag{6.9a}$$

$$G_{k,2}^j := M_k(\mathbb{U}_{k,2}^{j-1,0} + \mathbb{U}_{k,2}^{j-1,1}). \tag{6.9b}$$

From the fact that they include some terms from the previous time step, $G_{k,1|2}^j$ are dependent on the spatial refining method used on the interval $I_{j-1}$ and therefore their implementation can be of special interest in case of different (adaptive) spatial refinement of neighbouring time slabs.

Our intention is to determine the solution vector $\mathbb{U}^j$ defined as in Table 2.6 of the following global system

$$\mathbb{L}\mathbb{U}^j = \mathbb{F}, \tag{6.10}$$

where

$$\mathbb{L} := \begin{bmatrix} S & S & -k_j S & -k_j/2 S \\ k_j S & k_j/2 S & \varepsilon k_j S + M & \varepsilon k_j/2 S + M \\ 0 & 1/2 S & -k_j/2 S & -\frac{k_j}{3} S \\ k_j/2 S & k_j/3 S & \varepsilon k_j/2 S & \varepsilon k_j/3 S + 1/2 M \end{bmatrix}, \quad \mathbb{F} := \begin{bmatrix} G_1^j \\ F^{j,0} + G_2^j \\ \mathbf{0} \\ F^{j,1} \end{bmatrix}. \tag{6.11}$$

Here, $S, M$ are denoted as global stiffness and mass $2(n+1) \times 2(n+1)$ block matrices and $F^{j,0}, F^{j,1}, G_1^j, G_2^j$ are force and mass vectors of dimension $2(n+1) \times 1$.

In case of $dG(0)$ discretisation in time, the system has the similar structure to the one from Subsection 6.1.1.1. We solve a problem (6.10) where matrices $\mathbb{L}, \mathbb{F}$ read

$$\mathbb{L} := \begin{bmatrix} S & -k_j S \\ k_j S & \varepsilon k_j S + M \end{bmatrix}, \quad \mathbb{F} := \begin{bmatrix} G_1^j \\ F^{j,0} + G_2^j \end{bmatrix}, \tag{6.12}$$

with $G_1^j := S\mathbb{U}_1^{j-1}$, $G_2^j := M\mathbb{U}_2^{j-1}$. The structure of $S$ and $M$ is determined due to cubic spline ansatz in space. The solution vector $\mathbb{U}^j$ has the same structure as in Table 2.6, case $dG(0)$.

For both $q = 0, 1$, $\mathbb{U}^0$ is the coefficient matrix with respect to $U^{0-} = \Pi u_0$.

## 6.1.2   FE solution, $cG(1)$ time approximation

The derivation of the $I_j$-interval based algorithm for the computation of the discrete solution vector $\mathbb{U}^j$ relies on the notation from Table 2.5 and Table 2.6. We also assume that $\mathcal{S}^{j-1} \subseteq \mathcal{S}^j$ for all $j = 1, \ldots, N$, where $\mathcal{S}^j$ denote the space of spatial discrete functions related to time slab $I_j \times \Omega$.

### 6.1.2.1 FE solution, $cG(1) \otimes \mathcal{P}_1$

Given the weak problem (2.54) related to time interval $I_j$, by applying two different variants of constant test function in time $V \in \mathcal{W}_c$, we obtain the following two systems, namely

1. for $V^j = (V_1^j, 0) = (\phi_m, 0)$, $m = 0, \ldots, n$

$$\sum_{k=0}^{n}(U_{k,1}^j - \frac{k_j}{2}U_{k,2}^j)S_{m,k} = \sum_{k=0}^{n}(U_{k,1}^{j-1} + \frac{k_j}{2}U_{k,2}^{j-1})S_{m,k},$$

2. for $V^j = (0, V_2^j) = (0, \phi_m)$, $m = 0, \ldots, n$

$$\sum_{k=0}^{n}\frac{k_j}{2}(U_{k,1}^j + \varepsilon U_{k,2}^j)S_{m,k} + \sum_{k=0}^{n}U_{k,2}^j M_{m,k} = \int_{I_j}(f, \phi_m)dt - \sum_{k=0}^{n}\frac{k_j}{2}(U_{k,1}^{j-1} + \varepsilon U_{k,2}^{j-1})S_{m,k}$$
$$+ \sum_{k=0}^{n}U_{k,2}^{j-1}M_{m,k}.$$

Here $U^{j-1}$ denotes the solution from the previous time slab which is a priori known due to the fact that in case when $j = 1$, $U^0 = \Pi u_0$, i.e. $U^0$ represents the discrete variant of the initial solution from (1.28b).

We rewrite the systems above such that the global solution system reads

$$\mathbb{L}\mathbb{U}^j = \mathbb{F}, \tag{6.13}$$

where $\mathbb{U}^j$ is a solution vector related to the time interval $I_j$ and defined as in Table 2.5. $\mathbb{L}$ is $4(n+1) \times 4(n+1)$ block matrix and $\mathbb{F}$ is a $4(n+1) \times 1$ block matrix such that

$$\mathbb{L} := \begin{bmatrix} S & -k_j/2S \\ k_j/2S & \varepsilon k_j/2S + M \end{bmatrix}, \quad \mathbb{F} := \begin{bmatrix} G_1^j \\ F + G_2^j \end{bmatrix}, \tag{6.14}$$

with force and mass vectors defined such that

$$F := \left(\int_{I_j}(f, \phi_0)dt, \ldots, \int_{I_j}(f, \phi_n)dt\right)^T, \tag{6.15a}$$

$$G_1^j := S\mathbb{U}_1^{j-1} + \frac{k_j}{2}S\mathbb{U}_2^{j-1}, \tag{6.15b}$$

$$G_2^j := -\frac{k_j}{2}S\mathbb{U}_1^{j-1} + (M - \varepsilon\frac{k_j}{2}S)\mathbb{U}_2^{j-1}. \tag{6.15c}$$

Note that in case $j = 0$, $\mathbb{U}^0$ is a coefficient matrix of $\Pi u_0$.

### 6.1.2.2 FE solution, $cG(1) \otimes \mathcal{C}^1$

Given (2.65), for two different choices of test function $V \in \mathcal{W}_c$, a sum of $jk$-contributions (2.66) over $k = 1, \ldots, n$ yields two systems, namely

1. for $V^j = (V_1^j, 0) = (\phi_p, 0)$, $1 \leq p \leq 4$

$$\sum_{k=1}^{n}S_k(\mathbb{U}_{k,1}^j - k_j/2\mathbb{U}_{k,2}^j) = \sum_{k=1}^{n}G_{k,1}^j$$

2. for $V^j = (0, V_2^j) = (0, \phi_p)$, $1 \le p \le 4$

$$\sum_{k=1}^{n} \frac{k_j}{2} S_k \mathbb{U}_{k,1}^j + \sum_{k=1}^{n} (M_k + \varepsilon \frac{k_j}{2} S_k) \mathbb{U}_{k,2}^j = \sum_{k=1}^{n} F_k^j + G_{k,2}^j.$$

Here we used the vector notation for $\mathbb{U}_{k,1|2}^j$ from Table 2.6. Note also that for test functions $V$ above, we have $\mathbb{V}_{k,1,2}^j = e_p \in \mathbb{R}^4$ for $1 \le p \le 4$. $S_k, M_k$ are $4 \times 4$ element mass and stiffness matrices from (6.6), (6.7). $F_k^j$ is the $4 \times 1$ element forcing vector which has the same structure as vector $F_k^{j,0}$ from (6.8).

The $4 \times 1$ element mass vectors $G_{k,1}^j, G_{k,2}^j$ are related to the terms which include the discrete solution from the previous time step $\mathbb{U}_k^{j-1}$ respectively, i.e.

$$G_{k,1}^j := S_k(\mathbb{U}_{k,1}^{j-1} + \frac{k_j}{2} \mathbb{U}_{k,2}^{j-1}), \tag{6.16a}$$

$$G_{k,2}^j := -\frac{k_j}{2} S_k \mathbb{U}_{k,1}^j + (M_k - \varepsilon \frac{k_j}{2} S_k) \mathbb{U}_{k,2}^{j-1}. \tag{6.16b}$$

Global matrices $G_1^j, G_2^j$ depend on the spatial refining method on the previous time slab $I_{j-1} \times \Omega$. This need to be considered in case of different (adaptive) spatial refinement of two neighbouring time slabs.

Finally, we seek to find a $2(n+1) \times 1$ dimensional global solution vector $\mathbb{U}^j$ defined as in Table 2.6, of the system

$$\mathbb{L}\mathbb{U}^j = \mathbb{F}, \tag{6.17}$$

where

$$\mathbb{L} := \begin{bmatrix} S & -k_j/2 S \\ k_j/2 S & \varepsilon k_j/2 S + M \end{bmatrix}, \quad \mathbb{F} := \begin{bmatrix} G_1^j \\ F^j + G_2^j \end{bmatrix}.$$

$S, M$ are denoted as global stiffness and mass $2(n+1) \times 2(n+1)$ block matrices and $F^j, G_1^j, G_2^j$ are $2(n+1) \times 1$ dimensional force and mass vectors.

### 6.1.3   Incorporation of Dirichlet boundary conditions

In order to include the homogeneous Dirichlet boundary conditions into the global system related to time interval $I_j$, the first components concerning both vectors $\mathbb{U}_1^j$, $\mathbb{U}_2^j$ from solution vector $\mathbb{U}^j = (\mathbb{U}_1^j, \mathbb{U}_2^j)$ need to be set to zero.

Additional in case of Problem (DD), the ultimate (linear splines in space) and penultimate (cubic splines) component in both $\mathbb{U}_1^j$ and $\mathbb{U}_2^j$ must be also set to zero. A detailed review

| | $dG(0), cG(1)$ | $dG(1)$ |
|---|---|---|
| **DN** | $U_{0,1|2}^j = 0$ | $U_{0,1|2}^{j,0} = U_{0,1|2}^{j,1} = 0$ |
| **DD** | $U_{0,1|2}^j = U_{n,1|2}^j = 0$ | $U_{0,1|2}^{j,0} = U_{0,1|2}^{j,1} = U_{n,1|2}^{j,0} = U_{n,1|2}^{j,1} = 0$ |

Table 6.1: Zero components in solution vector $\mathbb{U}^j$ for homogeneous Dirichlet boundary data.

concerning the different discretisation in time can be found in Table 6.1.

If the global solution system reads

$$\mathbb{L}\mathbb{U}^j = \mathbb{F}$$
$$\mathbb{U}^j(dirichlet) = \mathbf{0}$$

where *dirichlet* denotes the set of indices at Dirichlet nodes and $\mathbf{0}$ is the zero vector, by proper indexing where $FN$ denotes the set of all remaining nodes, the system can be rewritten as

$$\begin{bmatrix} \mathbb{L}(FN, FN) & \mathbb{L}(FN, dirichlet) \\ \mathbb{L}(dirichlet, FN) & \mathbb{L}(dirichlet, dirichlet) \end{bmatrix} \begin{bmatrix} \mathbb{U}^j(FN) \\ \mathbb{U}^j(dirichlet) \end{bmatrix} = \begin{bmatrix} \mathbb{F}(FN) \\ \mathbb{F}(dirichlet) \end{bmatrix}.$$

Taking into account that we deal with homogeneous Dirichlet boundary conditions, the first block of equations can be rewritten as

$$\mathbb{L}(FN, FN)\mathbb{U}^j(FN) = \mathbb{F}(FN).$$

The second block of equations is not of interest and is omitted in the following which simplifies the implementation.

## 6.1.4   Incorporation of Neumann boundary conditions

In case of cubic splines in space, we may additionally require that the Neumann boundary condition

$$DU_1(t, 1) + \varepsilon DU_2(t, 1) = 0 \tag{6.18}$$

is also incorporated in the computation of the discrete solution. This can be done owing to the fact that cubic splines are continuous in first space derivative, whereas in case of linear splines this is not possible.

If $\mathbb{U}^j$ is the solution vector related to the time interval $I_j$ which can be obtained from the system

$$\mathbb{L}\mathbb{U}^j = \mathbb{F},$$

than the extended system with incorporated Neumann boundary condition reads

$$\begin{bmatrix} \mathbb{L} & B_1^T & \cdots & B_\ell^T \\ B_1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ B_\ell & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \mathbb{U}^j \\ \lambda_1 \\ \vdots \\ \lambda_\ell \end{bmatrix} = \begin{bmatrix} \mathbb{F} \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Here $\lambda_1, \ldots, \lambda_\ell$ stand for the Lagrange multipliers. In particular, $\ell \in \{1, 2\}$ such that $\ell = 1$ in case of $dG(0)$ and $cG(1)$ time discretisation and $\ell = 2$ for $dG(1)$ method in time.
Furthermore, $B_1, \ldots, B_\ell$ are a $(n + 1) \times 1$ matrices which arise from a substitution of the discrete form (2.4) for $U_1$ and $U_2$ into (6.18).

## 6.2   Examples

Within this section, we present five examples of the (strongly damped) wave equation used in the numerical experiments. The exact solutions of all presented examples are a priori known. The difference between them is in the choice of the boundary conditions, parameter $\varepsilon$ and homogeneous and non homogeneous right hand side $f$.

**Experiment 6.2.1.** We observe the inhomogeneous equation (1.13) without damping, i.e. $\varepsilon = 0$, on the domain $Q = [0, T] \times [0, 1]$ with forcing vector

$$f(t, x) = -(5 + t)(x - 2)e^{-x} \tag{6.19a}$$

and initial conditions

$$u_0(x) = 5xe^{-x}, \tag{6.19b}$$
$$u_1(x) = xe^{-x}. \tag{6.19c}$$

The advantage of using this example is that we know the analytical solution,

$$u(t, x) = (5 + t)xe^{-x} \tag{6.19d}$$

so that we can compare our numerical solution to this and calculate the error in energy norm. It should be also noted that, as the exact solution is very smooth, this is very well-behaved problem.

**Experiment 6.2.2.** We first consider the solution of the homogeneous problem (1.13)-(DN), i.e. $f = 0$, without damping ($\varepsilon = 0$) where the initial conditions read

$$y_0(x) = 0 \tag{6.20a}$$
$$y_1(x) = \frac{\pi}{2} \sin(\frac{\pi x}{2}). \tag{6.20b}$$

This problem is defined on the domain $Q = [0, T] \times [0, 1]$, $T > 0$.
The analytical solution is known and reads

$$y(t, x) = \sin(\frac{t\pi}{2}) \sin(\frac{x\pi}{2}). \tag{6.20c}$$

Owing to the fact that the exact solution is $C^\infty$ with respect to time and space variable, this is very well-behaved problem which does not give rise to oscillating discrete solution in space.

**Experiment 6.2.3.** Next we have on the same time-space domain $Q = [0, T] \times [0, 1]$, a homogenous wave equation (1.13)-(DD), i.e. $f = 0$ and $\varepsilon = 0$ which solution satisfies

$$y(0, x) = y_0(x) = \sin(x\pi) \tag{6.21a}$$
$$\dot{y}(0, x) = y_1(x) = 0 \tag{6.21b}$$

The exact solution is smooth and reads

$$y(t, x) = \cos(-t\pi) \sin(x\pi). \tag{6.21c}$$

**Experiment 6.2.4.** This problem is non homogeneous with damping $\varepsilon > 0$. The solution of this damped wave equation (1.13)-(DN) with the forcing vector

$$f(t,x) = e^{-t/\varepsilon} \sin(x\pi/2), \tag{6.22a}$$

and the initial conditions

$$y_0(x) = \varepsilon^2 \sin(x\pi/2), \tag{6.22b}$$
$$y_1(x) = -\varepsilon \sin(x\pi/2), \tag{6.22c}$$

is known and reads

$$y(t,x) = \varepsilon^2 e^{-t/\varepsilon} \sin(x\pi/2). \tag{6.22d}$$

The domain is $Q = [0,T] \times [0,1]$.

**Experiment 6.2.5.** Finally, we consider the solution of the wave equation (1.13)-(DN) i.e. with $\varepsilon = 0$, where the forcing vector reads

$$f(t,x) = -(5+t)(x-2)e^{-x}. \tag{6.23a}$$

The initial conditions are defined by

$$y_0(x) = 5e^{-x}x \tag{6.23b}$$
$$y_1(x) = 0 \tag{6.23c}$$

on the domain $Q = [0,T] \times [0,1]$. The analytical solution is known and reads

$$y(t,x) = (5+t^3)e^{-x}x. \tag{6.23d}$$

## 6.3 Validity of error and error bound with respect to change of parameters

Within this section we provide some numerical results which were conducted to examine the algorithm as well as the a posteriori error estimator obtained by energy method, cf. Chapter **??**. For all examples the domain was $Q = [0,T] \times [0,1]$ where $T = 1$, only in Subsection **??** the final time point changes owing to the purpose of experiments which assumes the study of the long time behaviour.

In Figure 6.1, we study the effects of saturation with respect to convergence in space for different choices of time step-size $k$. Similar behaviour is observed when step site $h$ changes and the convergence with respect to $k$ is plotted, see Figure 6.2. Notice that we plotted only the convergence behaviour of the exact error when $\mathcal{P}_1$ elements in space, but the same effects are obvious when $\mathcal{C}^1$ in space. The approximation by $\mathcal{P}_1$ elements in space is not so exact as the one by $\mathcal{C}^1$ elements, and therefore the effect of saturation becomes more clear.

In the following we study the behaviour of the reliability constant

$$C_r := \left( \|e\|_{L^\infty(\mathcal{H})}^2 + \varepsilon \|e_2\|_{L^2(H^1)} \right)^{1/2} \eta^{-1}, \tag{6.24}$$

in dependence of the parameter $\varepsilon$ and $T$. In Figure 6.3 when $cG(1), dG(0)$ and $dG(1)$ time approximation combined with $\mathcal{C}^1$ elements in space were applied, we may observe the nearly asymptotical behaviour of $C_r$. The approximated error is obtained with help of the energy arguments. This can be treated as the effect of the error accumulation. In Figure 6.4, we plot $C_r$ in case of the $cG(1), dG(0)$ and $dG(1)$ discretisation in time combined with $\mathcal{C}^1$ elements in space for different choice of terminal point in time $T$. The error bound is obtained with help of the energy arguments. The instability of the observed reliability constant may come from the fact that time step-size $k$ also changes in time as $T$, i.e. $k = T/2$.

## 6.4 Adaptivity

The solutions of the (strongly damped) wave equation, posses localised features, like e.g. wave fronts. Employing the adaptivity in the choice of the computational grids, may help the establishing of the efficient Algorithmen, which may resolve this features at certain portion and improve the convergence of the discrete solution toward the exact one.
Thereafter, the main task is to find the most optimal time-space mesh and then calculate the approximate solution $U$ on that mesh with a minimal cost concerning the time of computation and efficiency of the algorithm in a sense that the algorithm can be used as a platform for development of solution strategies for much more complicated problems.
In this work, we discussed this problem only partly. Namely, from Section 3.3, where the derivation of a posteriori error estimates by use of the energy techniques were discussed, we have

$$\text{quantity of interest} \leq \eta \text{ (value of the error estimator)},$$

where the quantities of interest represent the energy of the error and dissipative term. We may also notice, that in the formulation of the error estimator $\eta$, some terms arise form the global and some only form the local discretisation, i.e.

$$\eta = \eta_{gl} + \max_{1 \leq j \leq N} \eta_j \qquad (6.25)$$

where indices $j$ clearly indicate that in case of $\eta_j$ we refer to the local quantity computed only on time interval $I_j$. $\eta_{gl}$ is the quantity which is computed over he whole time interval $[0, T]$. The main idea is to use $\eta_j$ as the indicator for the mesh refinement in time. Therefore, the adaptive method may be formulated as follows

**Algorithm 1.**

Input: Spatial mesh $\mathcal{S}$ which is fix for each time slab; Discrete variant of $u_0$ with respect to $\mathcal{S}$; An initial coarse time partition $\mathcal{T}^0$ of the time interval $[0, T]$, $\mathcal{T}^0 = \bigcup_{\ell=1}^{N^i} I_\ell$, $i = 0$; $j = 1$.

1. Compute a Galerkin solution $U^j$ on the time-space slab $I_j \times \Omega$ where $I_j \in \mathcal{T}^i$.

2. Compute $\eta_j$. If $j < N^i$ increase $j$ and goto 1.
   Else goto 3.

3. For all $\eta_j$ where $\eta_j \geq \frac{1}{2} \max\{\eta_1, \ldots, \eta_{N^i}\}$, halve $I_j \in \mathcal{T}^i$ such that $I'_j \cup I''_j = I_j$, define $\mathcal{T}^{i+1} = \{\mathcal{T}^i \setminus I_j\} \cup \{I'_j, I''_j\}$, set $j = 1$, and goto 1.
   Else end.

It's obvious that $i$ in the algorithm above can be iterated a number of times, but we van also set an indicator for $i$ in order to provide an finiteness.

From the simple numerical computations, see Figures 6.5, 6.6 , it is obvious that this algorithm do not improve the convergence rate of the estimator. It just makes the exact error and the estimator smaller.  Possible improvements of the algorithm above which would include the adaptive mesh refinement as in time as in space, and also for which the efficiency can be easily proved is a matter of future investigation which are not included in this work.  Here we give just an overview which gives rise to the future discussions.   We also note that this strategy can not be applied for the estimates obtained by using the duality arguments due to the different form of the estimator then (6.25), cf. Section 4.3, where the local quantities do not occur separately in the error estimate. For some possible ways of designing an adaptive algorithms, we refer to the following literature sources [66, 64, 13].

As far as the adaptivity in goal oriented method is concerned, we refer to [9, 10, 11] for further discussion and references. This has not also been treated within this work.
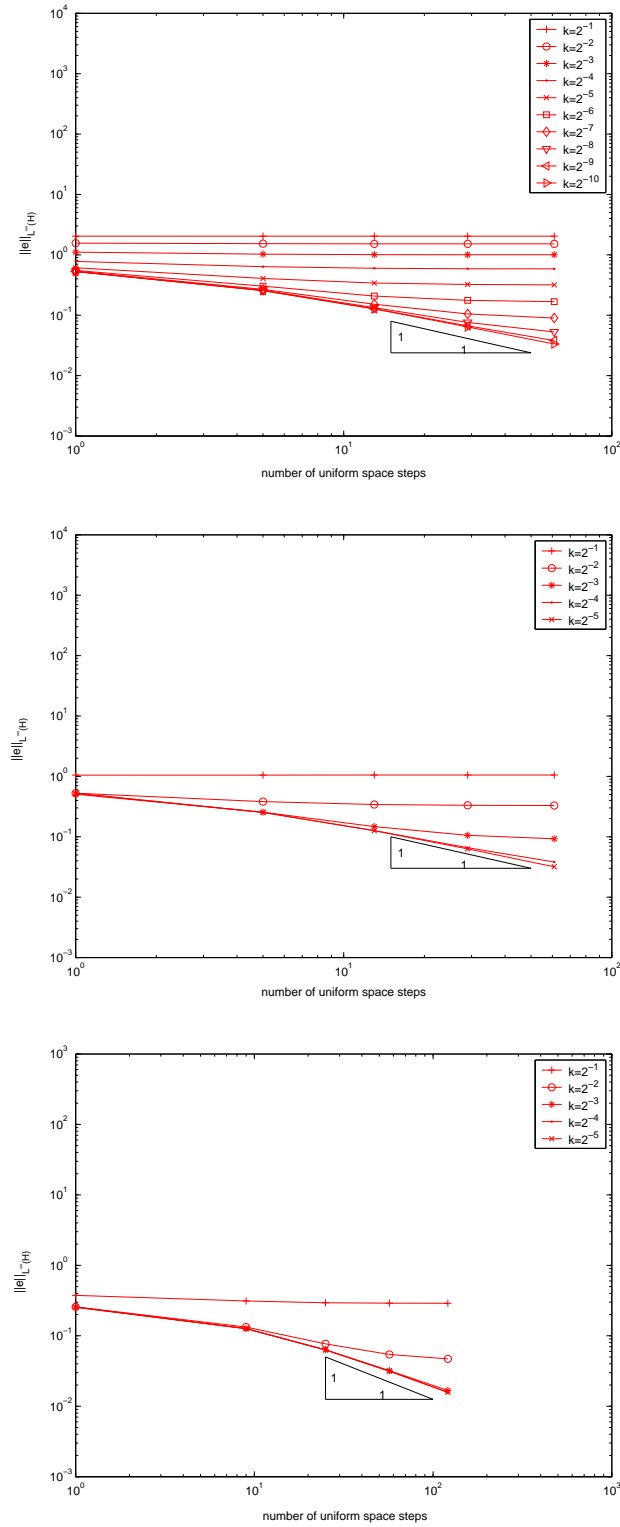
Figure 6.1: The exact energy of the error $\|e\|_{L^\infty(\mathcal{H})}$ with respect to the number of elements in space for the different time-step sizes. We observe the effect of saturation i.e. as soon as the step size $k$ becomes small, the convergence behaviour with respect to $h$, i.e. $\mathcal{O}(h)$ is obvious; First row: $dG(0)\otimes\mathcal{P}_1$, Second row: $cG(1)\otimes\mathcal{P}_1$, Third row: $dG(1)\otimes\mathcal{P}_1$; Example 6.2.3, $\varepsilon=0$, (DD), $T=1$.
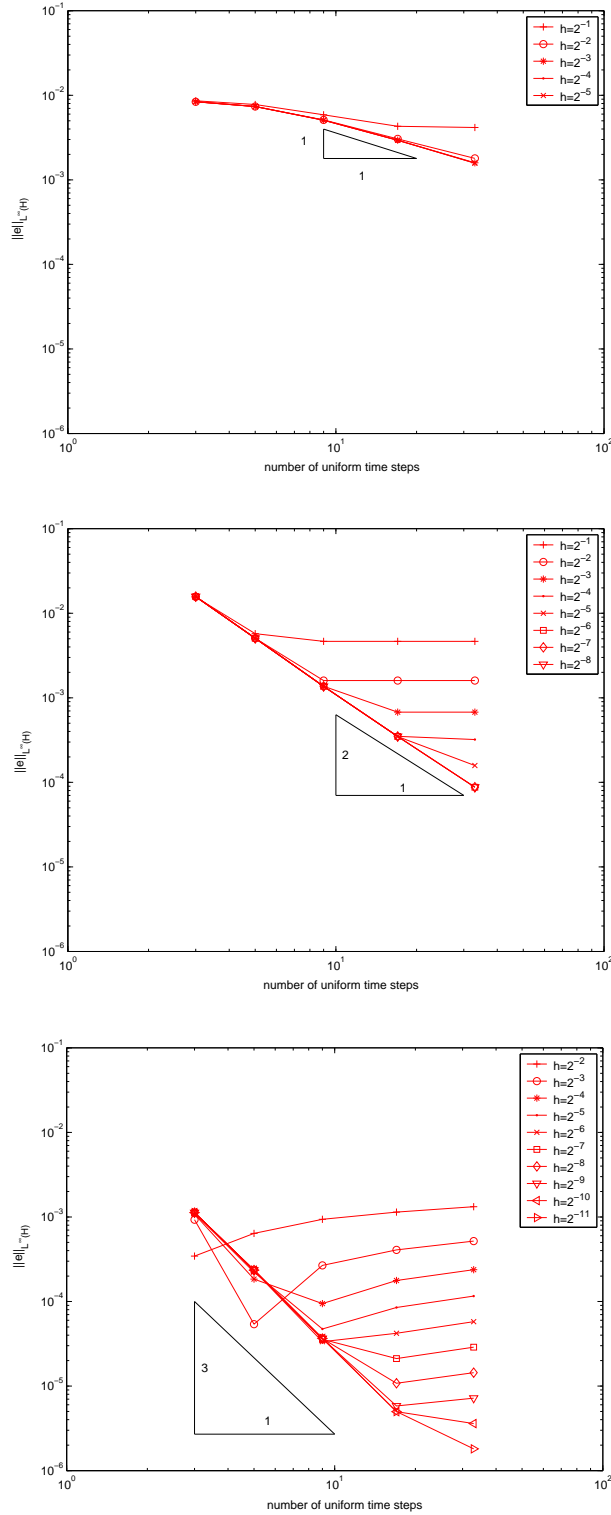
Figure 6.2: The exact energy of the error $\|e\|_{L^\infty(\mathcal{H})}$ with respect to the number of elements in time for the different time-step sizes. We observe the effect of saturation i.e. as soon as the step size $h$ becomes small, the convergence behaviour with respect to $k$ is obvious; The expected convergence order is $\mathcal{O}(k)$ for $dG(0)$, $\mathcal{O}(k^2)$ for $cG(1)$ and $\mathcal{O}(k^3)$ for $dG(1)$. First row: $dG(0)\otimes\mathcal{P}_1$, Second row: $cG(1)\otimes\mathcal{P}_1$, Third row: $dG(1)\otimes\mathcal{P}_1$; Example 6.2.4, $\varepsilon=0.1$, (DD), $T=1$.
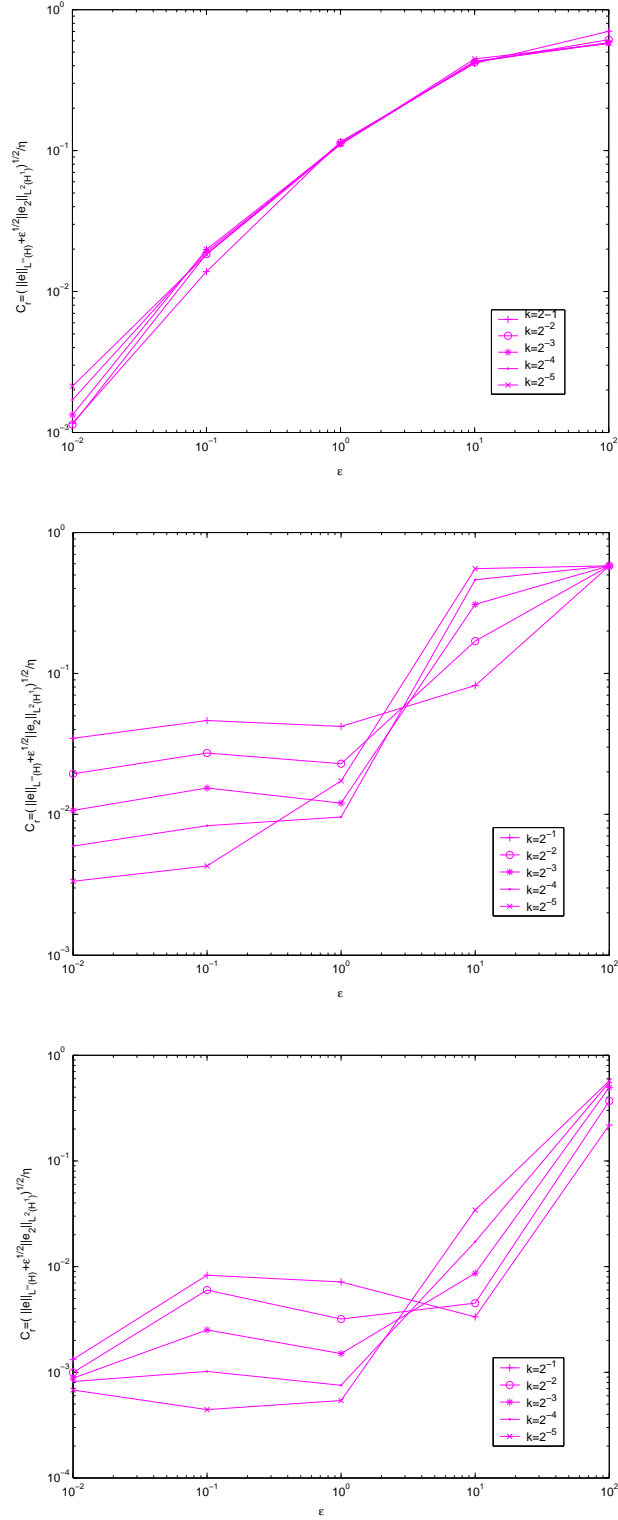
Figure 6.3: Reliability constant of $\eta$ in dependence of $\varepsilon$. $h$ is fix and $k \rightarrow 0$; First row: $dG(0) \otimes \mathcal{C}^1$, Second row: $cG(1) \otimes \mathcal{C}^1$, Third row: $dG(1) \otimes \mathcal{C}^1$. Obviously, the tendency towards the asymptotical behaviour of $\mathcal{C}_r$ is obvious when $dG(0)$ in time. If $cG(1)$ and $dG(1)$ this tendency is not so clear, but can be assumed. The growth of $C_r$ can be treated as the effect of the accumulation of the true error; Example 6.2.4, $\varepsilon > 0$, (DN), $T = 1$; energy method.
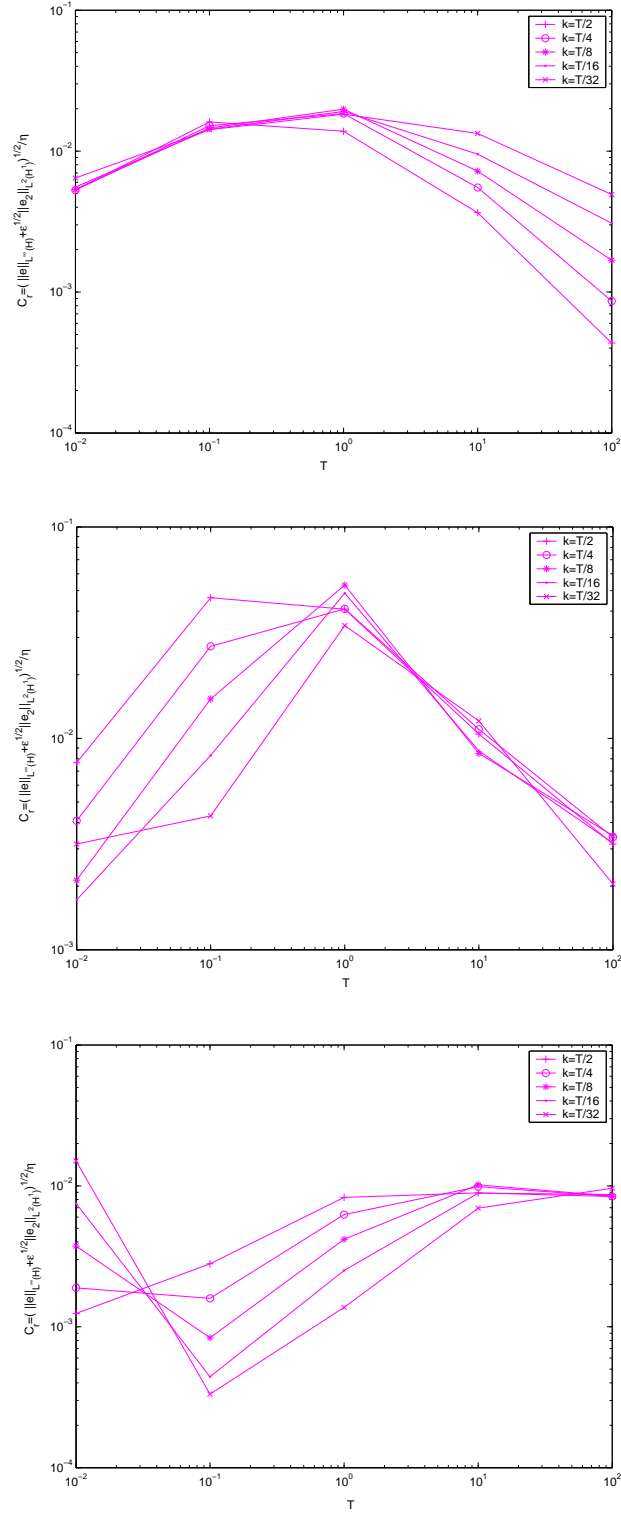
Figure 6.4: Reliability constant of $\eta$ in dependence of final time point $T$; First row: $dG(0) \otimes \mathcal{C}^1$, Second row: $cG(1) \otimes \mathcal{C}^1$, Third row: $dG(1) \otimes \mathcal{C}^1$; Example 6.2.4, $\varepsilon = 0.1$, (DN), $T = 1$; energy method.

Figure 6.5: Convergence behaviour of $\|e\|_{L^\infty(\mathcal{H})} + \sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}$ and $\eta$ for $dG(0)\otimes\mathcal{C}^1$ with respect to the number of elements in time; uniform and adaptive refinement (time); The adaptive refinement minimises the error and its bound; Example 6.2.4, $\varepsilon = 0.1$, (DN), $T = 1$; energy method.
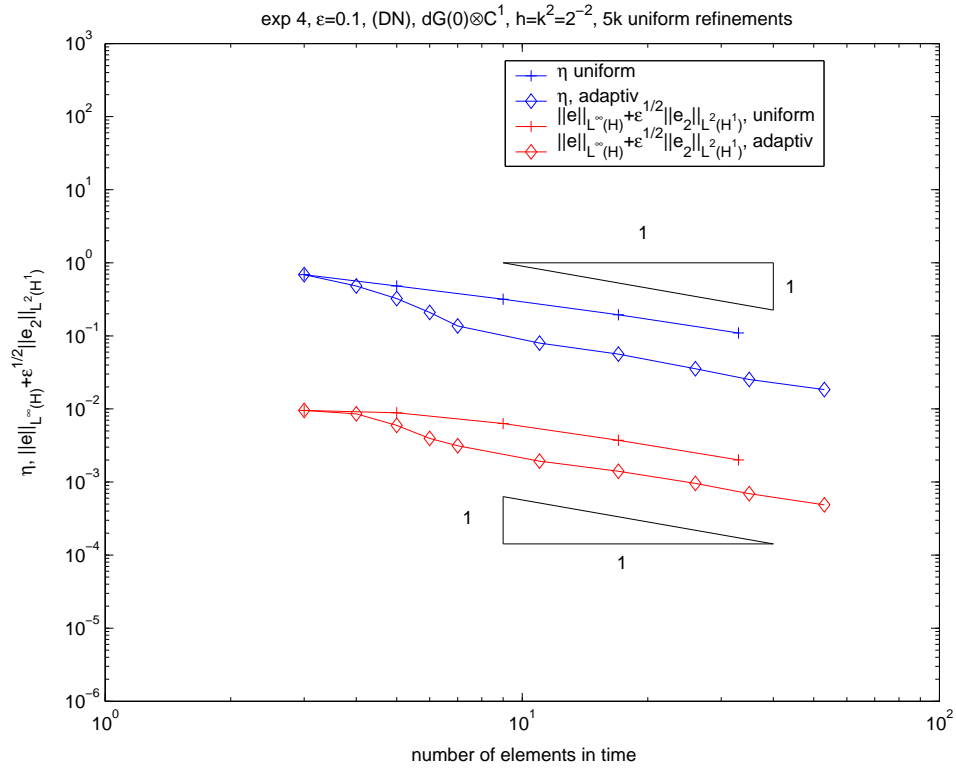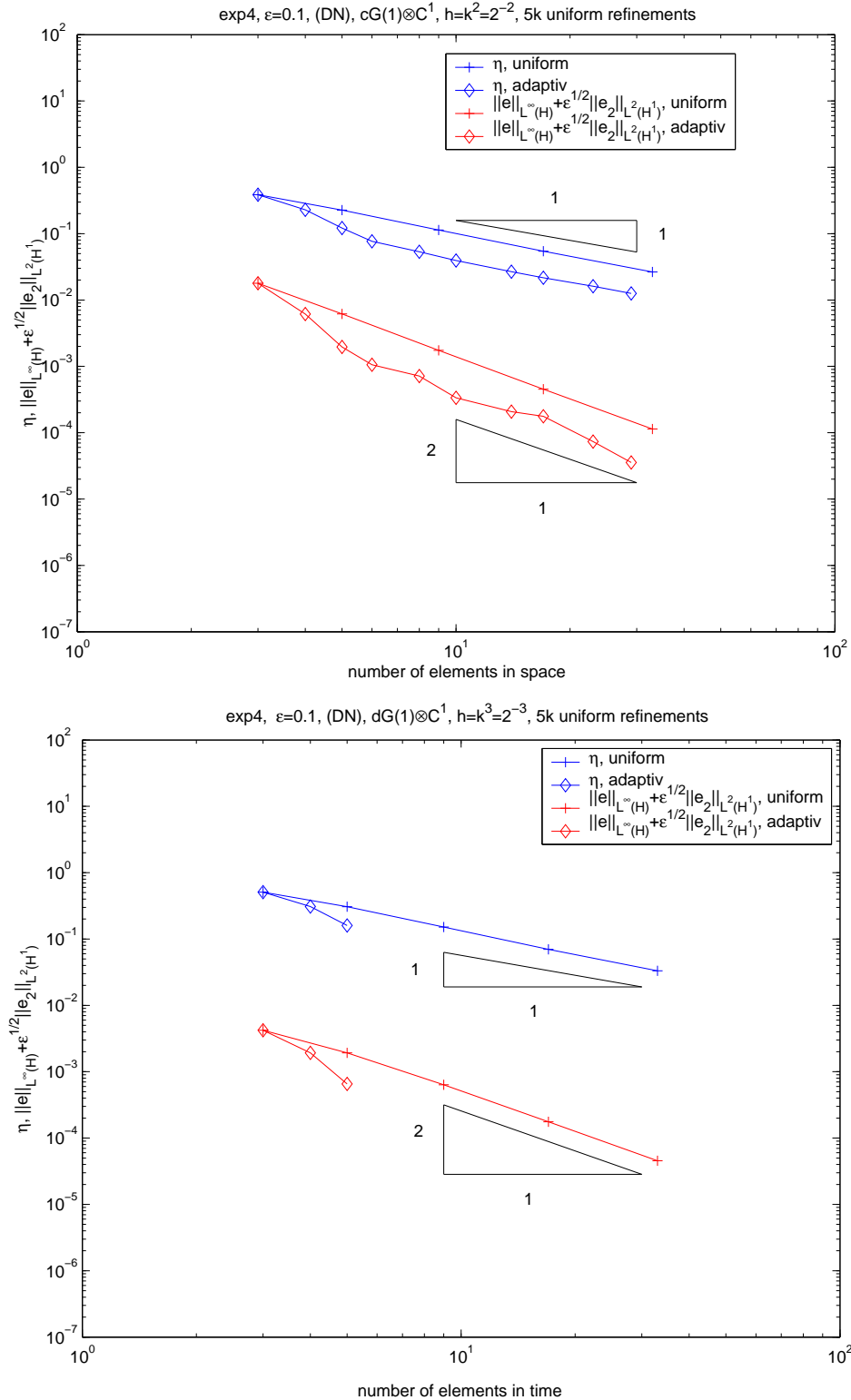
Figure 6.6: Convergence behaviour of $\|e\|_{L^\infty(\mathcal{H})} + \sqrt{\varepsilon}\|e_2\|_{L^2(H^1)}$ and $\eta$ for $cG(1)\otimes\mathcal{C}^1$ (first row) and $dG(1) \otimes \mathcal{C}^1$ (second row) with respect to the number of elements in time; uniform and adaptive refinement; The adaptive refinement in time does not improve the convergence order, it only minimises the error and its bound: Note that in case of the $dG(1)$ method (bottom figure), the efficiency of the adaptive refinement strategy is obvious only at the beginning of the discretisation owing to the influence of the initial solution when approximated on the coarse grid; Example 6.2.4, $\varepsilon = 0.1$, (DN), $T = 1$; energy method.

# Bibliography

[1] R.P. AGARWAL, P.J.Y. WONG. *Error inequalities in Polynomial Interpolation and Their Applications*, Dodrecht: Kluwer Academic publishers, 1993.

[2] J. ALBERTY. *Zeitdiskretisierungsverfahren für elastoplastische Probleme der Kontinuumsmechanik*, Ph.D. Thesis, University of Kiel, FRG, Tenea Verlag, 2001.

[3] J. ALBERTY, C. CARSTENSEN, S. FUNKEN. *Remarks around 50 lines of Matlab: short finite element implementation*, Numer. Algorithms 20, no. 2-3, pp. 117–137.

[4] J.H. ARGYRIS, D.W. SCHARP. *Finite elements in time and space*, Nucl. Engrg. Des., no. 10, pp. 456–464, 1969.

[5] A.K. AZIZ, P. MONK. *Continuous Finite Elements in Space and Time for the Heat Equation*, Math. Comp., vol. 52, no. 186, pp. 255–274, 1989.

[6] I. BABUŠKA, M. FEISTAUER, P. ŠOLIN. *On one approach to a posteriori error estimates for evolution problems solved by the method of lines*, Numer. Math. 89, pp. 225–246, 2001.

[7] G.A. BAKER. *Error Estimates for Finite Element Methods for Second Order Hyperbolic Equations*, SIAM. J. Numer. Anal., vol. 13, no. 4, 1976.

[8] L. BALES, I. LASIECKA. *Continuous Finite Elements in Space and Time for the Non homogeneous Wave Equation*, Comput. Math. Appl., vol. 27, no. 3, pp. 91–102, 1994.

[9] W. BANGERTH. *Adaptive Finite-Elemente-Methoden zur Lösung der Wellengleichung mit Anwendung in der Physik der Sonne*, Thesis, Universität Heidelberg, 1998.

[10] W. BANGERTH, R. RANNACHER. *Finite Element Approximation of the Acoustic Wave Equation: Error Control and Mesh Adaptivity*, East-West J. Numer. Math., vol. 7, no. 4, pp. 263–282, 1999.

[11] W. BANGERTH, R. RANNACHER. *Adaptive Finite Element for Differential Equations*, Lecture in Mathematics, ETH Zürich.

[12] V.V. BELETSKY, E.M. LEVIN. *Dynamics of Space Tether Systems*, Advances in the Astronautical Sciences, no. 83, 1993.

[13] C. BERNARDI, E. SÜLI. *Time and space adaptivity for the second-order wave equation*, Math. Models Methods Appl. Sci., vol. 3, no. 15, 2005.

[14] V. BONFIM, A.F. NEVES. *A One Dimensional Heat Equation with Mixed Boundary Conditions*, J. Differential Equations, no. 139, pp. 319–338, 1997.

[15] S.C. Brenner, L.R. Scott. *The Mathematical Theory of Finite Element Methods*, Springer-Verlag, 2002.

[16] C. Carstensen. *Adaptive Finite Element Methods, Vienna Lectures*, 2001.

[17] C. Carstensen. *An adaptive Mesh-Refining Algorithm Allowing for an $H^1$ Stable $L^2$ Projection onto Courant Finite Element Spaces*, Constr. Approx. 20, pp. 549–546, 2004.

[18] C. Carstensen, J. Alberty. *Discontinuous Galerkin time discretisation in elastoplasticity: motivation, numerical algorithms and application*, Comput. Methods Appl. Mech. Engrg., no. 191, pp. 4949–4968, 2002.

[19] F. Chen, B. Guo, P. Wang. *Long Time Behavior of Strongly Damped Nonlinear Wave Equations*, J. Differential Equations, no. 147, pp. 231–241, 1998.

[20] P.G. Ciarlet. *The Finite Element Method for Elliptic Problems*, North Holland, Amsterdam, 1978.

[21] R. Courant, K. Friedrichs, H. Lewy. *Über die partiellen Differenzengleichungen der mathematischen Physik*, Math. Ann., no. 100, pp. 32–74, 1928.

[22] M. Crouzeix, V. Thomée. *The stability in $L^p$ and $W^{1,p}$ of the $L^2$-Projection onto Finite Element Function Spaces*, Math. Comp., no. 48, pp. 521–532, 1987.

[23] K. Eriksson, C. Johnson. *Adaptive Finite Element Methods for Parabolic Problems I: Linear Model Problem*, SIAM J. Numer. Anal., vol. 28, no. 1, pp. 43–77, 1991.

[24] K. Eriksson, C. Johnson. *Adaptive Finite Element Methods for Parabolic Problems II: Optimal Error Estimates in $L_\infty L_2$ and $L_\infty L_\infty$*, SIAM J. Numer. Anal., vol. 32, no. 3, pp. 706–740, 1995.

[25] K. Eriksson, C. Johnson. *Adaptive Finite Element Methods for Parabolic Problems V: Long-time Integration*, SIAM J. Numer. Anal., vol. 32, no. 6, pp. 1750–1763, 1995.

[26] K. Eriksson, C. Johnson, Stig Larsson. *Adaptive Finite Element Methods for Parabolic Problems V: Analytical Semigroups*, SIAM J. Numer. Anal., vol. 35, no. 4, pp. 1315–1325, 1998.

[27] L.C. Evans. *Partial Differential Equations*, Graduate Studies in Mathematics, vol. 19, American Mathematical Society, 1999.

[28] D.A. French. *A space-time finite element method for the wave equation*, Comput. Methods Appl. Mech. Engrg., no. 107, pp. 145–157, 1993.

[29] D.A. French, S. Jensen. *Long-time behavior of arbitrary order continuous time Galerkin schemes for some one-dimensional phase transition problems*, IMA J. Numer. Anal., no. 14, pp. 421–442, 1994.

[30] D.A. French, T. Peterson. *A continuous space-time finite element method for the wave equation*, Math. Comp., vol. 65, no. 214, pp. 491–506, 1996.

[31] I. Fried. *Finite element analysis of time-dependent phenomena*, AIAA J., no. 7, pp. 1170–1173, 1969.

[32] V. GIRLAUT, P.A. RAVIART. *Finite Element Methods for Navier-Stokes Equations - Theory and algorithms*, Springer-Verlag, 1980.

[33] D.H. GRIFFEL. *Applied functional analysis*, Mathematics and its Applications, Ellis Horwood Ltd, 1981.

[34] L. GUERRIERO, I. BAKEY. *Space Tethers for Science in the Space Station Ero*, Societa Italiana di Fisica, Conference Proceedings, vol. 14, Bologna, 1988.

[35] W. HACKBUSCH. *Elliptic Differential Equations*, Springer-Verlag, 1992.

[36] E. HAIRER, G. WANNER. *Solving Ordinary Differential Equations, II. Stiff and Differential-Algebraic Problems*, Springer Verlag, 1991.

[37] R. HARTMANN. *A-posteriori Fehlerschätzung und adaptive Schrittweiten- und Ortsgittersteuerung bei Galerkin-Verfahren für die Wärmeleitungsgleichung*, Thesis, Fakultät für Mathematik, Universität Heidelberg, 1998.

[38] G.M. HULBERT, T.J.R. HUGHES. *Space-time finite element methods for second order hyperbolic equations*, Comput. Methods Appl. Mech. Engrg., no. 84, pp. 327–348, 1990.

[39] B.L. HULME. *One-step piecewise polynomial Galerkin methods for initial value problems*, Math. Comput., no. 26, pp. 415–426, 1972.

[40] B.L. HULME. *Discrete Galerkin and related one-step methods for ordinary differential equations*, Math. Comput., no. 26, pp. 881–891.

[41] T.J.R. HUGHES, G.M. HULBERT. *Space-time finite element methods for elastodynamic: Formulations and error estimates*, Comput. Methods Appl. Mech. Engrg., no. 66, pp. 339–363, 1988.

[42] C. JOHNSON. *Discontinuous Galerkin finite element methods for second order hyperbolic problems*, Comput. Methods Appl. Mech. Engrg., no. 107, pp. 117–129, 1993.

[43] C. JOHNSON. *Numerical solutions of partial differential equations by the finite element method*, Cambridge University Press, 1987.

[44] C. JOHNSON. *Error estimates and automatic time step control for numerical methods for stiff ordinary differential equations*, Technical report 1984-27, Department of Mathematics, Chalmers University of Technology, 1998.

[45] C. JOHNSON, J. PITKÄRNTA. *Finite Element Methods For Linear Hyperbolic Problems*, Math. Comp., vol. 46, no. 173, pp. 1–26, 1986.

[46] C. JOHNSON, U. NÄVERT, J. PITKÄRNTA. *An Analysis of the Discontinuous Galerkin Method for a Scalar Hyperbolic Equation*, Comput. Methods Appl. Mech. Engrg., no. 45, pp. 285–312, 1984.

[47] P. KOHLER, W. MAAG, R. WEHRLI, R. WEBER, H. BRAUCHLI. *Dynamics of system of two Satellites Connected by a Deployable and Extensible Tether of Finite Mass*, Societa Italiana di Fisica, Conference Proceedings, vol. 14, Bologna, 1988.

[48] M. KRUPA, M. SCHAGERL, A. STEINDL, P. SZMOLYAN, H. TROGER. *Relative equilibria of tethered satellite systems and their stability for very stiff tethers*, Dynamical systems, no. 16, pp. 253–278, 2001.

[49] A. KUHN, W. STEINER, J.ZEMANN, D. DINEVSKI, H. TROGER. *A Comparison of Various Mathematical Formulations and Numerical Solution Methods for the Large Amplitude Oscillations of a String Pendulum*, Appl. Math. Comput., no. 67, pp. 227–264, 1995.

[50] S. LARSSON, V. THOMÉE, L. WAHLBIN. *Finite-element methiods for a strongly damped wave equation*, IMA J. Numer. Anal., vol. 11, no. 1, pp. 115–142, 1991.

[51] P. LESAINT, P.A. RAVIART. *On a finite element method for solving the neutron transport equation*, Mathematical Aspects of Finite Elements in Partial Differential Equations, Academic Press New York, pp. 89–123, 1974.

[52] M. LUSKIN, R. RANNACHER. *On the Smoothing Property of the Galerkin Method for Parabolic Equations*, SIAM J. Numer. Anal., vol. 19, no. 1, pp. 99–113, 1982.

[53] W. MCLEAN. *Strongly elliptic systems and boundary integral equations*, Cambridge University Press, 2000.

[54] A.F. NEVES. *On the strongly damped wave equation and the heat equation with mixed boundary conditions*, Abstr. Appl. Anal., vol.5, no. 3, pp. 175–189, 2000.

[55] R.H. NOCHETTO, G. SAVARÉ, C. VERDI. *A posteriori error estimates for variable time-step discretizations of nonlinear evolution equations* Comm. Pure Appl. Math. 53, no. 5, pp. 525–589, 2000.

[56] J.T. ODEN. *A general theory of finite elements II. Applications*, Internat. J. Numer. Methods Engrg., no. 1, pp. 247–259, 1969.

[57] J.T. ODEN, J.N. REDDY. *An introduction to the mathematical theory of finite elements*, John Wiley & Sons, 1976.

[58] W. POTH, M. MATZL, W. AUZINGER, A. STEINDL, H. TROGER. *Comparison of displacement versus natural variables for the numerical simulation of a string pendulum*, Nonlinear Dynamics, Chaos, Control and Their Applications to the Engineering Sciences, vol. 5, pp. 127–136, 2002.

[59] W.H. REED, T.R. HILL. *Triangular mesh methods for the neutron transport equation*, Report LA-UR-73-479, Los Alamos Scientific Laboratory, 1973.

[60] D. REDFERN, C. CAMPBELL. *The Matlab 5 Handbook*, Springer-Verlag, 1997.

[61] G.R. RICHTER. *An explicit finite element method for the wave equation,* Appl. Numer. Math., vol. 16, pp. 51–64, 1994.

[62] W. STEINER. *Numerische Untersuchungen an verkabellten Satellitensystemen*, Thesis, Technische Universität Wien, 1992.

[63] W. Steiner, A. Steindl, H. Troger. *Dynamics of a Space Tethered Satellite system with Two Rigid Endbodies*, Proceedings of the Fourth International Conference on Tethers in Space, Smithsonian Institution, Washington DC, pp. 1367–1379, 1995.

[64] E. Süli. *A posteriori error analysis and adaptivity for finite element approximations of hyperbolic problems*, Technical Report no. 97/21, Oxford University Computing Laboratory.

[65] E. Süli, P. Houston, C. Schwab. *hp-Finite Element Methods for Hyperbolic Problems*, Research Report no. $99 - 14$, ETH Zürich.

[66] E. Süli, N. Jackson. *Adaptive Finite Element Solution of $1D$ European Option Pricing Problems*, Technical Report no. 97/05, Oxford University Computing Laboratory.

[67] E. Süli, C. Wilkins. *Adaptive Finite Element Methods for the Damped Wave equation*, Technical Report no. 96/23, Oxford University Computing Laboratory.

[68] E. Süli, C. Wilkins. *A Priori Analysis for the Semi-Discrete Approximation to the Nonlinear Damped Wave Equation*, Technical Report no. 99/09, Oxford University Computing Laboratory.

[69] V. Thomée. *Galerkin Element Methods for Parabolic Problems*, Springer-Verlag, 1984.

[70] L. Vu-Quoc, J.C. Simo. *Dynamics of Earth-Orbiting Flexible Satellites with Multibody Components*, J. Guidance Control Dynam., vol. 10, pp. 549–558, 1987.

[71] L. Wahlbin. *A Modified Galerkin Procedure with Hermite Cubics for Hyperbolic Problems*, Math. Comp., vol. 29, no. 132, pp. 978–984, 1975.

[72] S. Waldron. *$L_p$-error bounds for Hermite interpolation and associated Wirtinger inequalities*, Constr. Approx., vol. 13, no. 4, pp. 461–479, 1997.

[73] D. Werner. *Funktionalanalysis*, Springer-Verlag, 1995.

[74] G. Wiedermann, M. Schagerl, A. Steindl, H. Troger. *Computation of Force Controlled Deployment and Retrieval of a Tethered Satellite by the Finite Element Method*, Proceedings of ECCM'99, (W.Wunderlich Ed.), pp. 410–429, München, 1999.

[75] E. Zeidler. *Nonlinear Functional Analysis and its Applications*, Band I-IV, Springer-Verlag, 1990.

[76] *Tethers in Space Handbook*, NASA Report NASW-3921, 1986.

[77] *Tethers in Space*, Advances in the Astronautical Sciences, vol. 62, 1986.

# Curriculum Vitae

# JELENA BOJANIĆ

Institut für Analysis und Scientific Computing
Technische Universität Wien
Wiedner Hauptstrasse 8 - 10
A - 1040 Wien
Tel: ++435880110163
e-mail: jelena@aurora.anum.tuwien.ac.at

- **Persönliche Daten**

| | |
|---|---|
| Name | Jelena Bojanić, geboren Petričković |
| Geburtstag | 17. 04. 1979 |
| Geburtsort | Zadar, Kroatien |
| Staatsbürgerschaft | Serbien |
| Famillienstand | verheiratet seit Mai 2005 |

- **Akademische Abschlüsse**

| | |
|---|---|
| 08. 02. 2003 | Sponsion zur Diplom-Ingenieurin, Technische Universität Wien |
| 28. 08. 2001 | Sponsion zur Diplommathematikerin, Universität Belgrad |

- **Ausbildung**

| | |
|---|---|
| seit 01. 01. 2003 | wissenschaftliche Mitarbeiterin am Institut für Analysis und Scientific Computing, Technische Universität Wien |
| 02. 2002-02. 2003 | Studium der Technischen Mathematik an der Technischen Universität Wien Studienzweig: Mathematik in den Naturwissenschaften |
| 09. 2001-01. 2002 | Deutschkurs, Vorstudienlehrgang, Wien |
| 1997 - 2001 | Mathematikstudium an der Fakultät für Mathematik, Universität Belgrad Studienzweig: Numerische Mathematik und Optimierung |
| 06. 1993 | Reifeprüfung |
| 1993 - 1997 | Gymnasium, "Zemunska Gimnazija", Belgrad, Serbien, Naturwissenschaftlich-mathematische Fachrichtung, |
| 1993 - 1997 | Mittlere Musikschule mit dem Schwerpunkt Klavier |
| 1985 - 1993 | Grundschule, Grundschule für Musik |

- **Forschung**

| | |
|---|---|
| seit 01. 2005 | Mitglied in der Forschungsgruppe "Numerics and Simulation of Differential Equations", Technische Universität Wien |
| 01. 2004 - 05. 2005 | Doktorand in der FWF Projekt P15274-N04 *Numerische Algorithmen in Computational Micromagnetics*, Projektleiter: Prof. Carsten Carstensen |
| 06. 2003 - 01. 2004 | Doktorand in der FWF Projekt P16461-N12 *Adaptive Raumdiskretisierungen in 4 Beispielen* , Projektleiter: Prof. Carsten Carstensen |

Wien, den 14. 06. 2005