

Thesis: Analytical Approximation of Statistical Rounding

The **Christian Doppler Laboratory Embedded Machine Learning** at the *Institute of Computer Technology* does research on Deep Neural Networks (DNN) in resource constrained embedded devices. It studies how energy consumption and resource usage can be minimized while keeping high accuracy. The solution space is characterized by architecture parameters, DNN optimization and transformations, implementation platform configurations, and mapping options. This design space is huge, poorly understood, and it is rapidly evolving.

In this work, statistical rounding[1][2] is the main focus. The question is can we approximate statistical rounding, by a modified non-linearity and if that is possible what are the benefits and downsides. This thesis project consists of the following steps:

- Get a good understanding of statistics and handling distributions, and statistical rounding.
- Calculate the resulting distribution of a linear/convolution layer under the assumption, that each weight has a binomial distribution.
- Fold this distribution into the non-linearity.
- Calculate the overhead of the parameters of this new non-linearity, as they should be calculated each training step.
- Calculate the back-propagation.
- Implement and test it against classical statistical rounding.

This thesis offers you an excellent opportunity to get into the hot topic of deep learning. It allows you to become an expert in configuring neural networks. Moreover, you acquire critical skills in optimizing neural networks for efficient embedded inference.

For details, please consult:

- Matthias Wess (matthias.wess@tuwien.ac.at)
- Daniel Schnöll (daniel.schnoell@tuwien.ac.at)

References

- [1] Matteo Croci et al. "Stochastic rounding: implementation, error analysis and applications". In: *Royal Society Open Science* 9.3 (2022), p. 211631. DOI: 10.1098/rsos.211631.
- [2] Nick Higham. *What is stochastic rounding?* July 2020. URL: <https://nhigham.com/2020/07/07/what-is-stochastic-rounding/>.