



RESEARCH INTERNSHIP IN TRANSFORMER MODEL QUANTIZATION & INFERENCE

fortiss is the research institute of the Free State of Bavaria for the development of software-intensive systems with headquarters in Munich. The scientists at the institute cooperate in research, development and transfer projects with universities and technology companies in Bavaria, Germany and Europe. The focus is on research into state-of-the-art methods, techniques and tools for the development of software and AI-based technologies for dependable, secure cyber-physical systems such as the Internet of Things (IoT). fortiss is organized in the legal form of a non-profit limited liability company. Shareholders are the Free State of Bavaria (majority shareholder) and the Fraunhofer-Gesellschaft zur Förderung der angewandten Forschung e.V. www.fortiss.org

We are looking for a Master student who would like to complete their research internship in the field of autonomous driving where we aim at deploying a Transformer-based multi-modal sensor fusion model for 3D object detection to our autonomous driving system, which is based on open-source software stack Baidu [Apollo](https://github.com/Baidu-Apollo). This internship offers an exciting opportunity to work on challenging problems at the intersection of deep learning, AI inference, and sensor fusion.

Your tasks:

- Explore and implement state-of-the-art techniques for optimizing model inference performance, with a focus on TensorRT integration and model quantization.
- Conduct experiments and performance evaluations to assess the effectiveness and efficiency of different optimization strategies and quantization techniques.
- Demonstration in Apollo with nuScenes dataset.

Your profile:

- Master student currently enrolled in Electrical and Computer Engineering or a related field.
- Practical experience in programming languages such as Python, C++.
- Good background knowledge in deep learning and model quantization.
- Excellent communication skills in English.
- Prior experience with NVIDIA GPU programming and TensorRT is a plus.

Our offer:

- An international and dynamic work environment with highly qualified colleagues.
- Increased experience with deep learning and machine learning engineering.
- Flexible working conditions, e.g., home office, flexible working hours.
- Possibility to pursue your Master's thesis on further autonomous driving topics.

Please submit your application with a detailed CV and a current transcript.

Contact for details or direct application: Xiangzhong Liu, xliu@fortiss.org

Published on 05.03.2024.