# COMPSCI 690AB SYSTEMS FOR DEEP LEARNING
## Reading Assignment 6

## The case for 4-bit precision: k-bit Inference Scaling Laws – Summary

This paper mainly addresses the challenges of quantizing large language models for efficient inference during maintenance of the model's performance. Due to their large size and computational demands, there are many problems during deployment because of resource constrained environments. The paper suggests that by reducing the number of bits through quantization can reduce latency, and thus aiming to get an optimal performance with zero-shot accuracy. The author then proceeds to explain the relationship between latency and the total model bits, explaining the importance of lower precision numbers in case of reducing latency. The paper works on the hypothesis of proxy quantization from developing outlier-dependent quantization. Here, weights are aimed to be quantized to a higher precision for the outlier feature dimensions. This is done by using standard deviation method where each layer's hidden unit weights act as a proxy in order to identify outlier-feature dimensions. This will aim in providing a model which has a constant memory footprint across all the tasks. One limitation of the paper proposal is that there isa dependency on constant memory footprint across all the models and tasks. Even though this approach is simpler and more efficient, it does not take into account the requirements and complexities of all different LLMs and its architecture. In future studies, there could be some exploration done on dynamic memory allocation which will adjust the memory footprint according to the specific characteristics of the model.

Srimathi Mahalingam