Adaptive Layer Compression and Storage with QoS-Aware Loading for LLM Serving

Meriem Bouzouad¹, Vincent Lannurien¹, Yuan-Hao Chang², Jalil Boukhobza¹

¹Lab-STICC, CNRS, UMR 6285, ENSTA, Institut Polytechnique de Paris, Brest ²Institute of Information Science, Academia Sinica, Taipei, Taiwan



Presentation Outline



Context



Problem justification



Contribution

Context

Generative Ai has infiltrated the world due to its ability to approximate or even surpass human intelligence.



"Who doesn't use generative AI? Exactly."



Problem justification

"Behind the success of generative AI"



- storage.
- Optimizations include: Quantization, pruning, distillation..etc [2-3]

Rule of thumb: Take a model that can fit in memory. However edge workload vary and static model deployment doesnt can underestimate or overestimate system's needs.

Problem statement: How can we design an edge-optimized solution that accelerates LLM inference to meet real-time workload and Quality of Service (QoS) constraints?

They are Large: they consume memory, compute and

 \Rightarrow Tradeoff between latency and generation quality [2].

02

Contribution

- Dynamic layer compression depending on its importance.
- Leverages a pool of models with different sizes to adapt to QoS dynamically such as latency, energy, and accuracy.
- Efficient storage to switch between models on demand.

Intelligent compression + Adaptability.



03

References

[1] Zhenyan Lu et al. "Small Language Models: Survey, Measurements, and Insights". In: arXiv preprint arXiv:2409.15790 (2024).

[2] Zhihang Yuan et al. "LLM Inference Unveiled: Survey and Roofline Model Insights". In: arXiv preprint arXiv:2402.16363 (2024)

[3] Razvan-Gabriel Dumitru et al. "Layer-Wise Quantization: A Pragmatic and Effective Method for Quantizing LLMs Beyond Integer Bit-Levels". In: arXiv preprint arXiv:2406.17415 (2024).



Thank you

