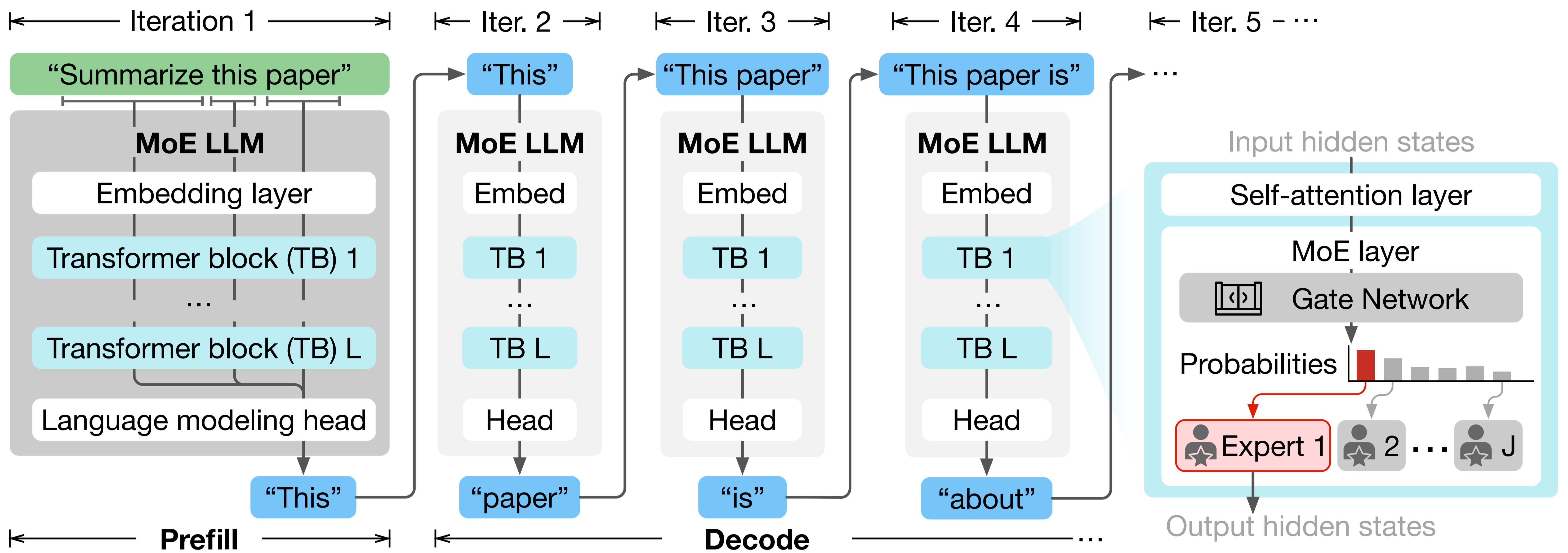


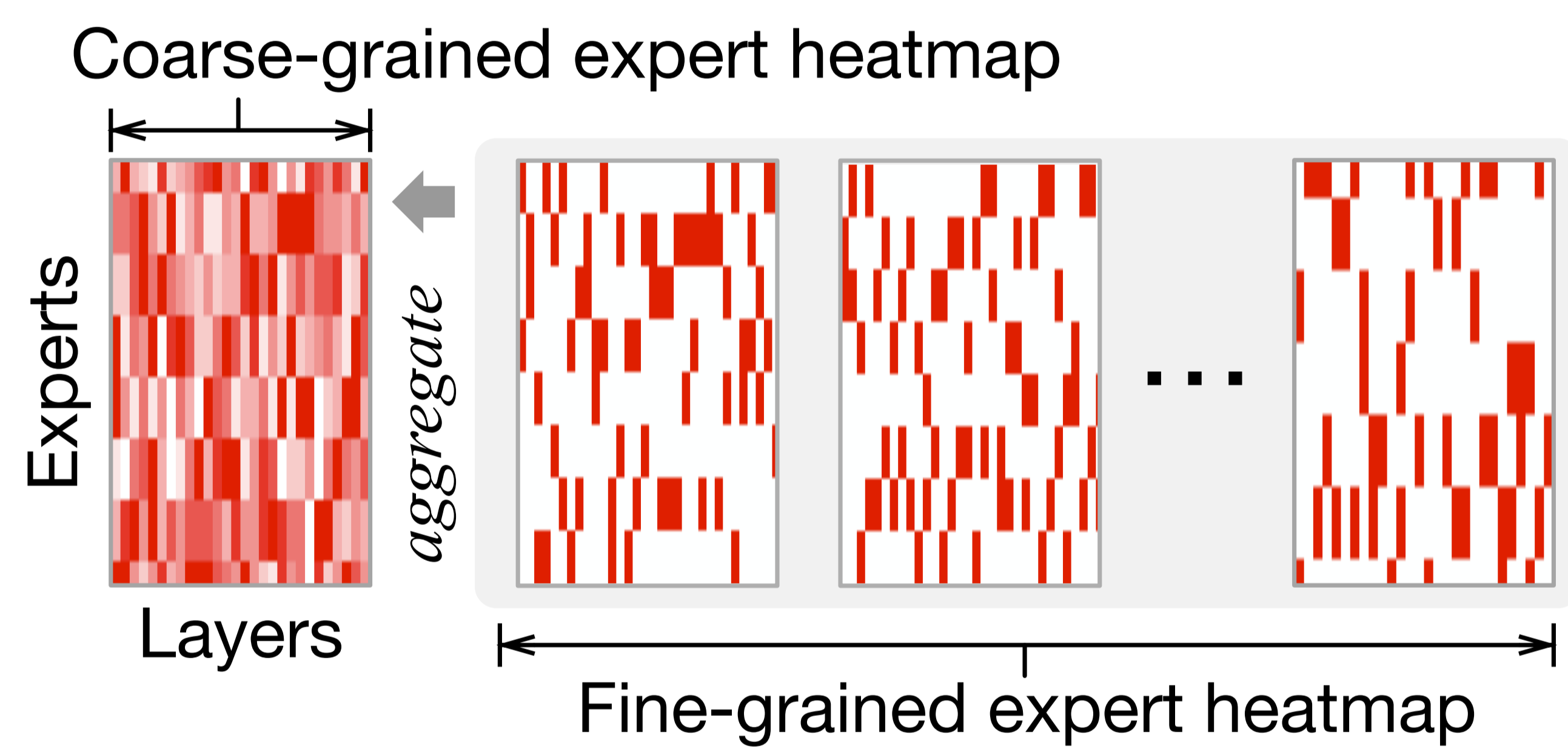
Taming Latency-Memory Trade-Off in MoE-Based LLM Serving via Fine-Grained Expert Offloading

Hanfei Yu¹, Xingqi Cui², Hong Zhang³, Hao Wang⁴, Hao Wang¹ Stevens Institute of Technology¹, Rice University², University of Waterloo³, Rutgers University⁴

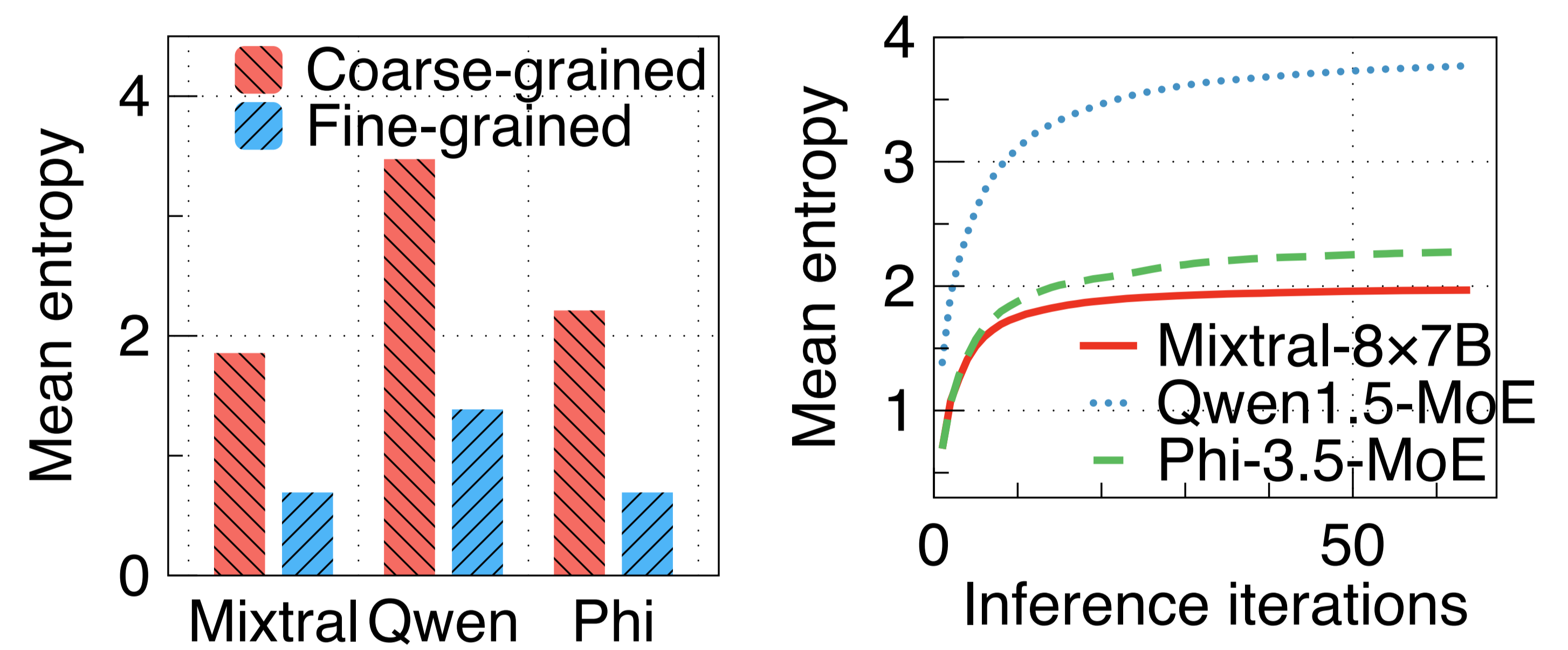
0. Mixture-of-Experts (MoE) Based Large Language Model (LLM) Serving



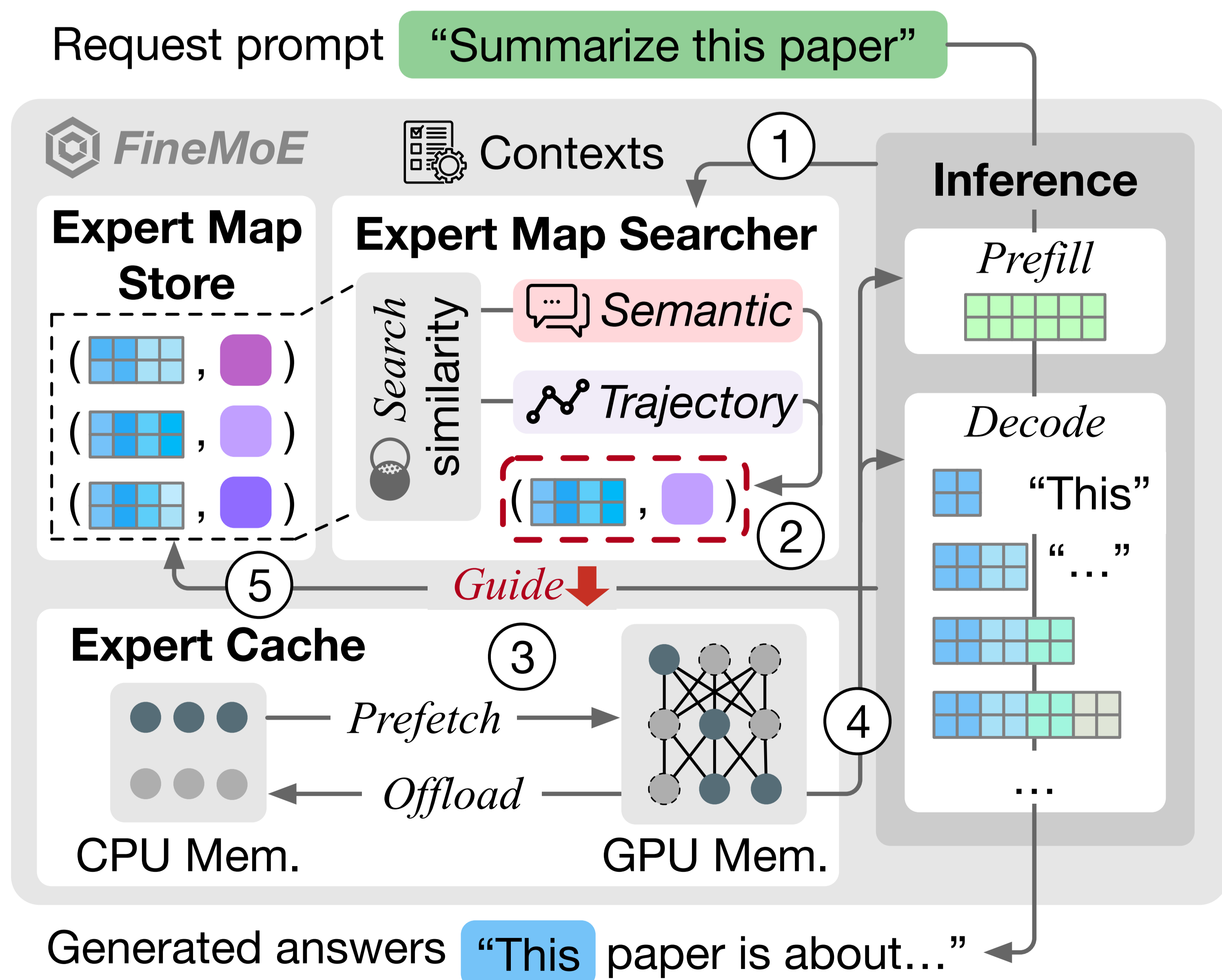
1. Coarse- vs. Fine-grained Expert Offloading



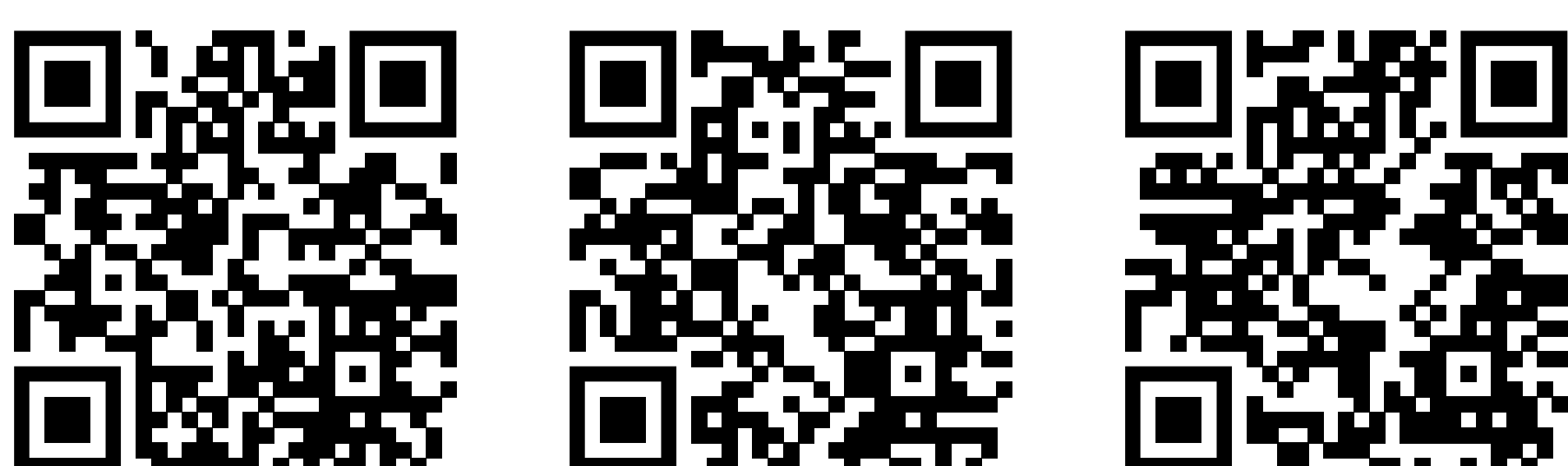
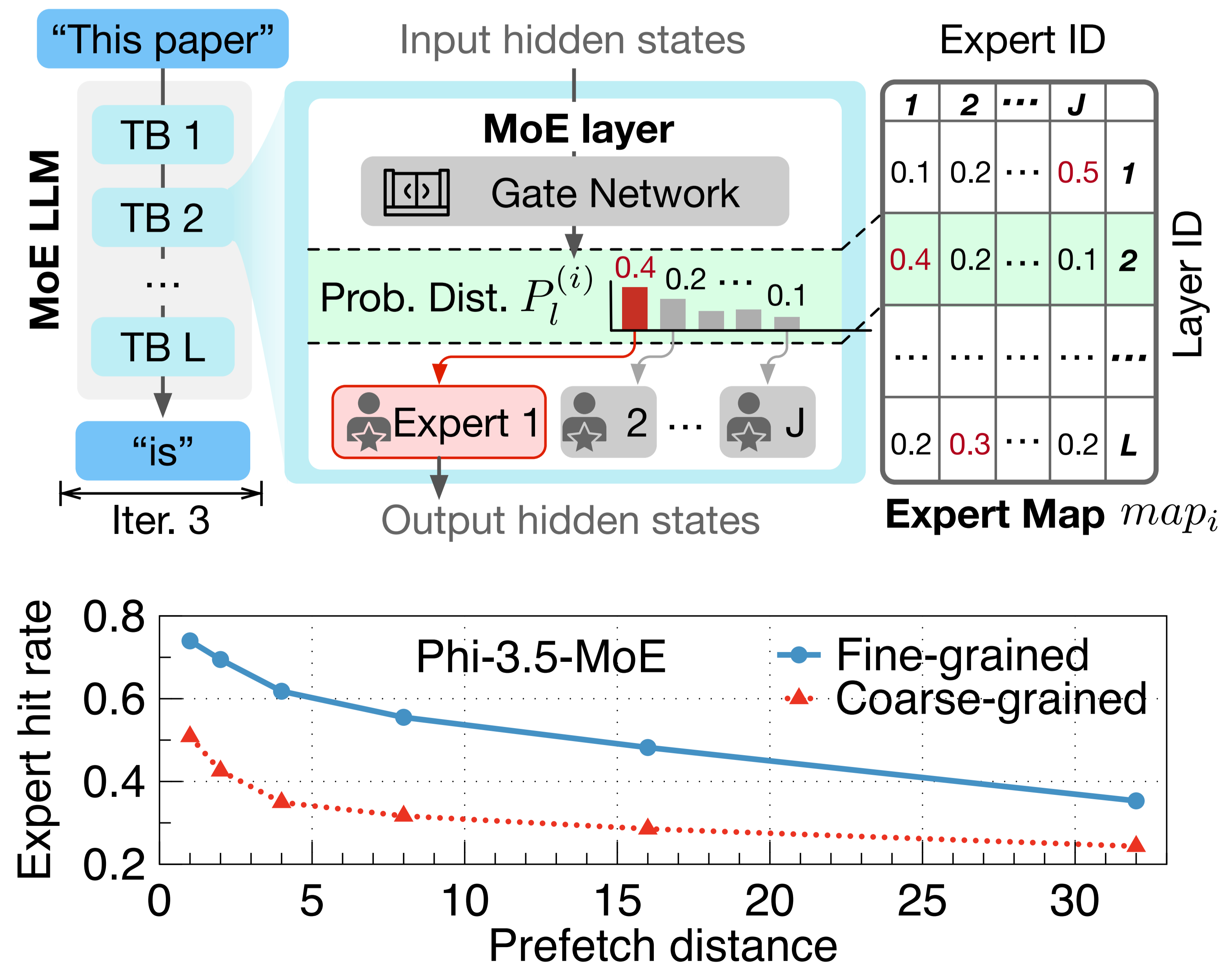
2. Entropy of Expert Activation Patterns



3. Workflow of FineMoE



4. Design of Expert Probability Map



We thank anonymous reviewers and our shepherd, Dr. Yaniv David, for their valuable feedback. This project is supported by U.S. National Science Foundation and AWS Cloud Credit for Research.

