

GPU Compute

MoE LLM

Dual Positional
Encoding
§ 3.2.2

Dense Transformer-based
Expert Predictor
§ 3.2.1

GPU Memory

Non-Expert
Model

LLM K-V
Cach

Pred Model
§ 3.2.1

Cached Experts
§ 3.4

§ 3.3.2
Unique Experts

CPU Memory



Expert Selection

§ 3.3.1
Expert Deduplication

§ 3.3.2

Unique Experts

CPU Compute