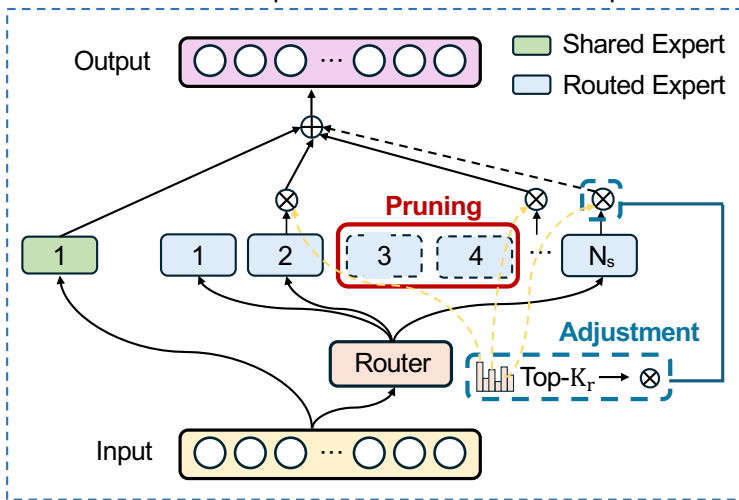
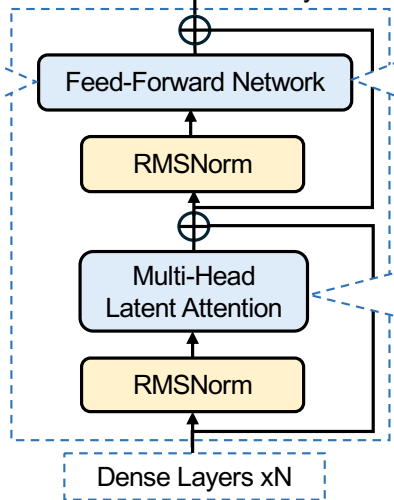


Pre-Quantization Optimization for Maximum Compression



MoE Layers xM



Dynamic Mixed Precision Quantization

