

Prefill Stage
Optimization Goal

max Cache Reuse s.t. TTFT SLO, Minimum MFU, KVCache < DRAM

Decoding Stage Optimization Goal

max Throughput s.t. TBT SLO,

TBT SLO, KVCache < VRAM