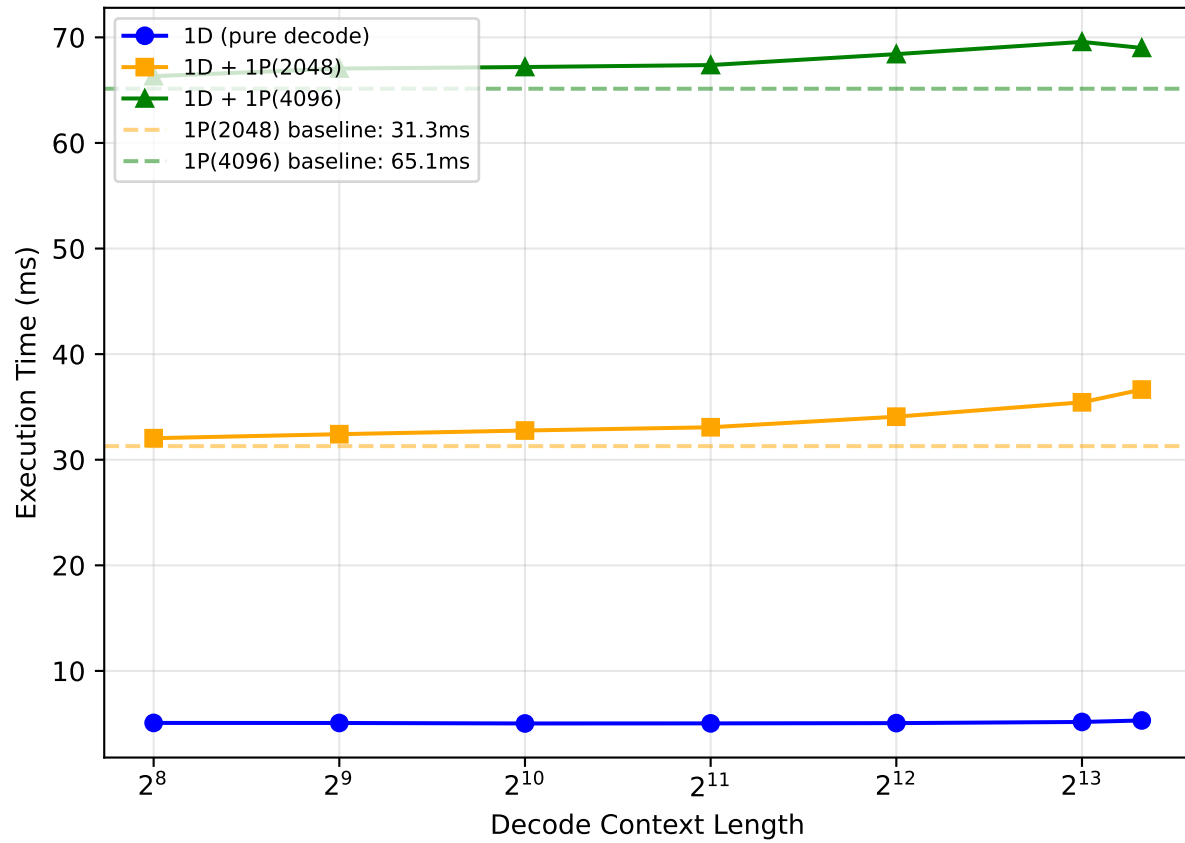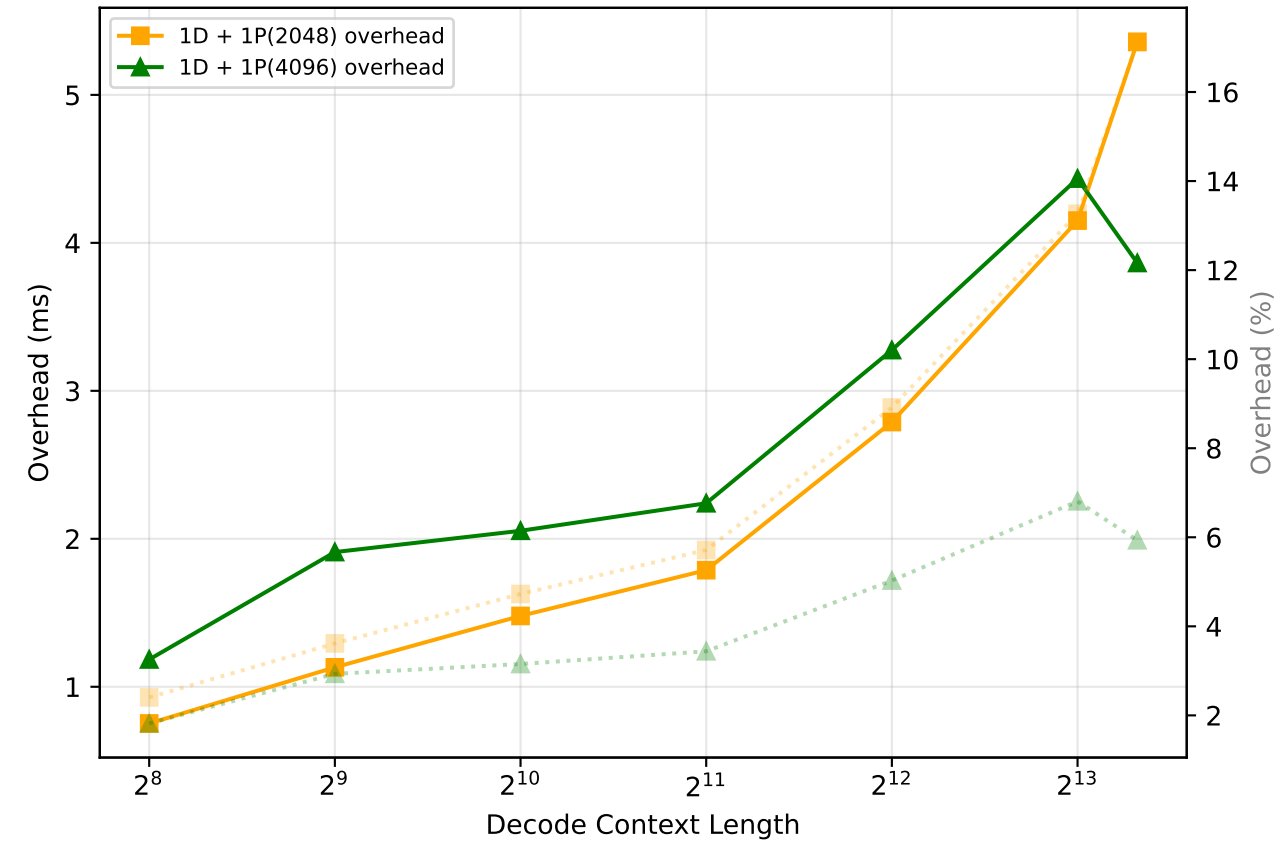vLLM Benchmark Results - Qwen/Qwen3-4B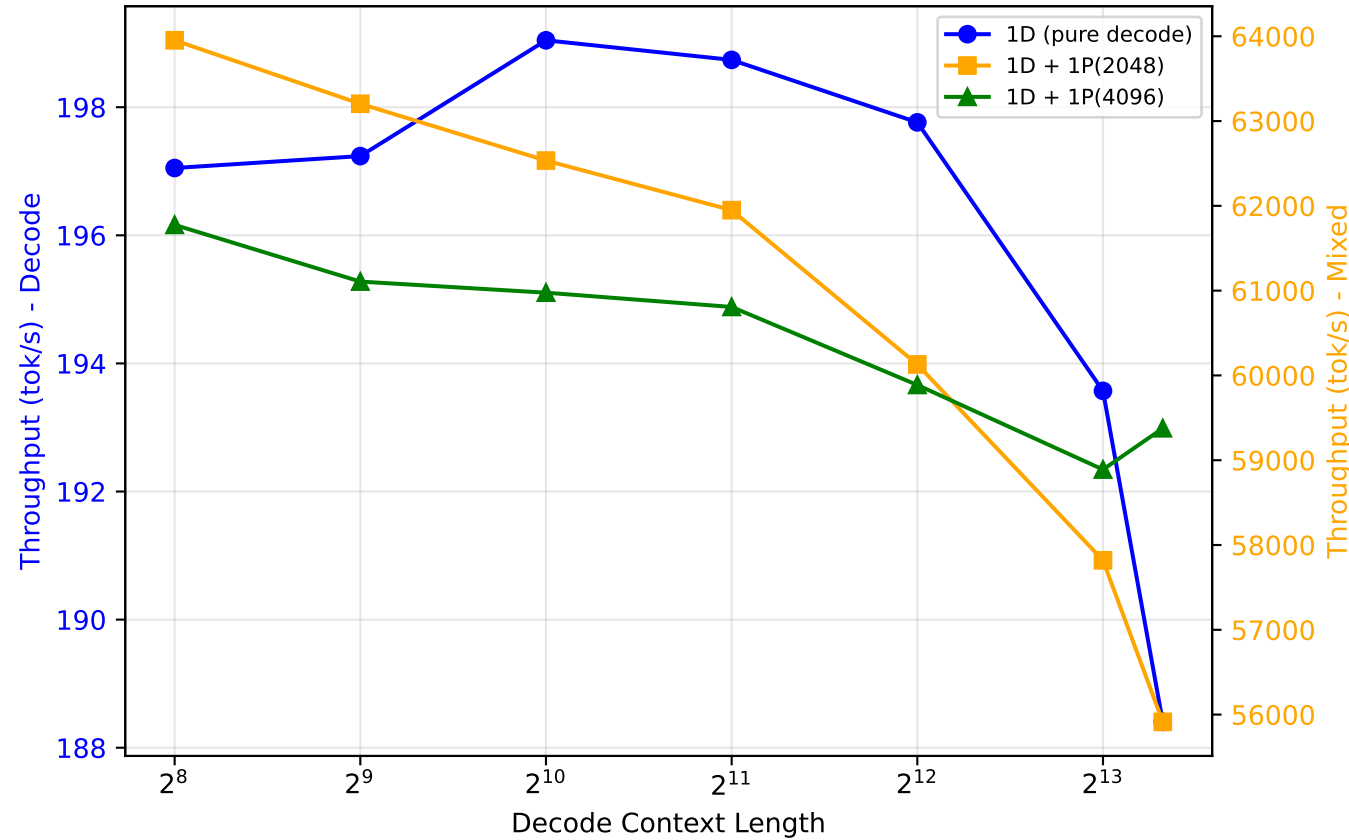