**Baseline Performance Comparison (10K tokens)**
**Larger Models → Better Parallelism**

**Throughput**

IOPS (×1000)

- Llama2-7B: 91K
- Llama2-13B: 97K
- Llama2-70B: 150K

**Response Time**

Average Latency (ms)

- Llama2-7B: 47.6ms
- Llama2-13B: 42.7ms
- Llama2-70B: 30.0ms

**Reliability**

ECC Failures

- Llama2-7B: 18,117
- Llama2-13B: 16,929
- Llama2-70B: 10,822