**Random (1500 ⇒ 250) (100 Request Rate)**

E2E Latency + 100ms Network Latency

Legend: No CC, CC

LLama 3.3 70B Int4 — 120.2s (No CC), 118.4s (CC)
GPT OSS 120B — 11.8s (No CC), 12.5s (CC)
Mistral 3.1 24B — 28.9s (No CC), 29.7s (CC)
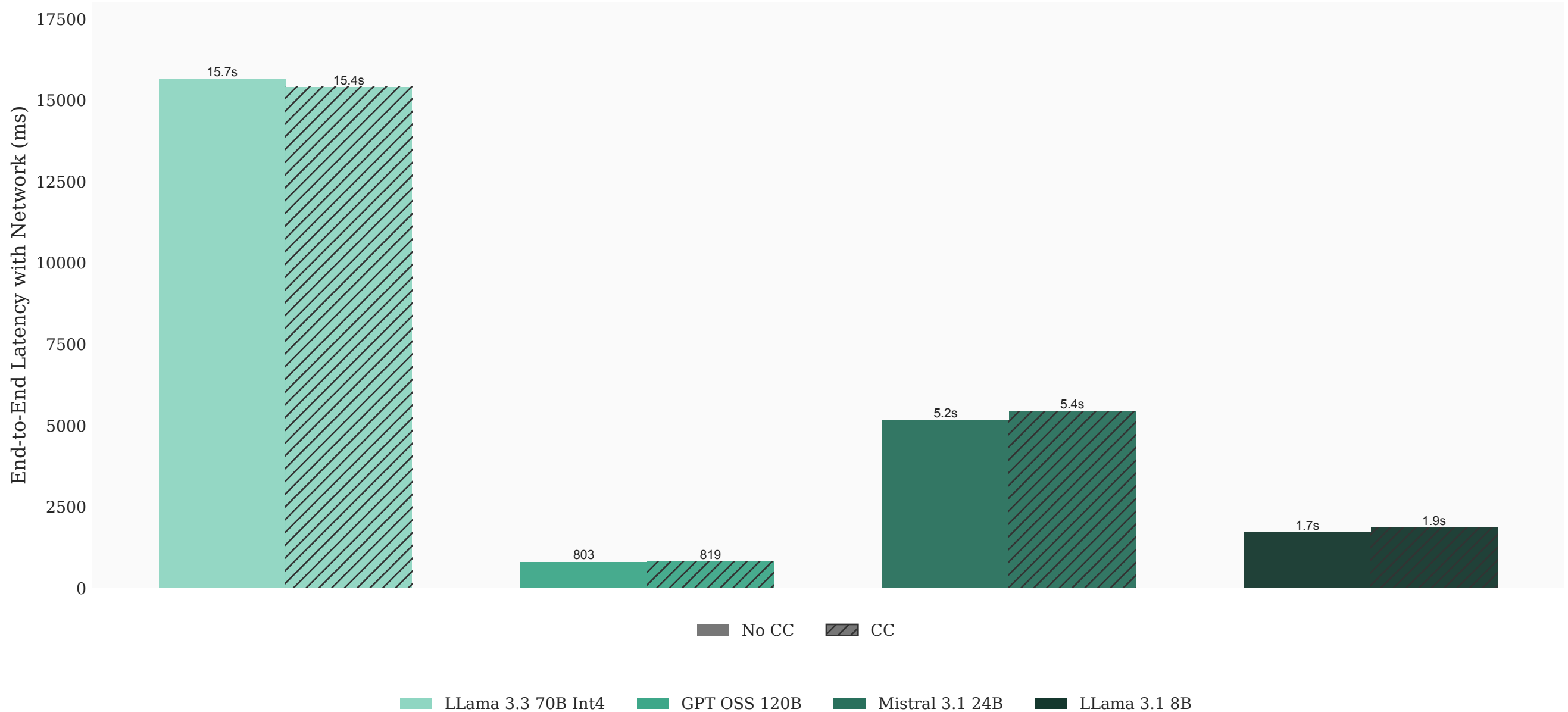LLama 3.1 8B — 11.1s (No CC), 11.8s (CC)

Y-axis: End-to-End Latency with Network (ms)

**Random (1500 ⇒ 250) (50 Request Rate)**

E2E Latency + 100ms Network Latency

**Random (1500 ⇒ 250) (Single Request)**

E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

- LLama 3.3 70B Int4: 15.7s (No CC), 15.4s (CC)
- GPT OSS 120B: 803 (No CC), 819 (CC)
- Mistral 3.1 24B: 5.2s (No CC), 5.4s (CC)
- LLama 3.1 8B: 1.7s (No CC), 1.9s (CC)

No CC  CC

LLama 3.3 70B Int4    GPT OSS 120B    Mistral 3.1 24B    LLama 3.1 8B

# Random (4000 ⇒ 1000) (100 Request Rate)

## E2E Latency + 100ms Network Latency



Bar chart titled "Random (4000 ⇒ 1000) (100 Request Rate)" with subtitle "E2E Latency + 100ms Network Latency". Y-axis: End-to-End Latency with Network (ms), ranging 0 to 300000.

Values (No CC / CC):
- LLama 3.3 70B Int4: 280.4s / 279.7s
- GPT OSS 120B: 33.9s / 34.3s
- Mistral 3.1 24B: 122.0s / 126.2s
- LLama 3.1 8B: 49.4s / 52.1s

Legend: No CC, CC (hatched)

**Random (4000 ⇒ 1000) (50 Request Rate)**

E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

- 279.9s
- 277.5s
- 34.0s
- 36.1s
- 121.0s
- 125.6s
- 50.2s
- 52.9s

No CC   CC

LLama 3.3 70B Int4   GPT OSS 120B   Mistral 3.1 24B   LLama 3.1 8B

# Random (4000 ⇒ 1000) (Single Request)
## E2E Latency + 100ms Network Latency

**End-to-End Latency with Network (ms)**

- 182.5s
- 180.2s
- 3.0s
- 3.7s
- 35.7s
- 40.0s
- 9.0s
- 10.0s

Legend: No CC | CC

LLama 3.3 70B Int4 | GPT OSS 120B | Mistral 3.1 24B | LLama 3.1 8B

# Random (1000 ⇒ 1000) (100 Request Rate)
## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

- LLama 3.3 70B Int4: 110.0s (No CC), 109.1s (CC)
- GPT OSS 120B: 13.0s (No CC), 12.3s (CC)
- Mistral 3.1 24B: 54.9s (No CC), 57.9s (CC)
- LLama 3.1 8B: 22.8s (No CC), 25.5s (CC)

Legend: No CC, CC

LLama 3.3 70B Int4    GPT OSS 120B    Mistral 3.1 24B    LLama 3.1 8B

# Random (1000 ⇒ 1000) (50 Request Rate)

## E2E Latency + 100ms Network Latency



**End-to-End Latency with Network (ms)**

- LLama 3.3 70B Int4: 109.0s (No CC), 108.8s (CC)
- GPT OSS 120B: 12.2s (No CC), 12.5s (CC)
- Mistral 3.1 24B: 54.3s (No CC), 57.1s (CC)
- LLama 3.1 8B: 19.5s (No CC), 21.9s (CC)

Legend: No CC, CC

LLama 3.3 70B Int4   GPT OSS 120B   Mistral 3.1 24B   LLama 3.1 8B

# Random (1000 ⇒ 1000) (Single Request)

## E2E Latency + 100ms Network Latency



Bar chart showing End-to-End Latency with Network (ms) for four models, comparing No CC and CC configurations.

- LLama 3.3 70B Int4: No CC 18.8s, CC 19.5s
- GPT OSS 120B: No CC 2.7s, CC 2.7s
- Mistral 3.1 24B: No CC 19.9s, CC 21.4s
- LLama 3.1 8B: No CC 7.1s, CC 7.9s

Legend: No CC, CC (hatched)

**ShareGPT (100 Request Rate)**

**E2E Latency + 100ms Network Latency**

End-to-End Latency with Network (ms)

- 29.3s / 29.2s — LLama 3.3 70B Int4
- 7.5s / 7.9s — GPT OSS 120B
- 6.2s / 6.6s — Mistral 3.1 24B
- 2.6s / 2.9s — LLama 3.1 8B

Legend: No CC / CC

LLama 3.3 70B Int4 | GPT OSS 120B | Mistral 3.1 24B | LLama 3.1 8B

**ShareGPT (50 Request Rate)**

**E2E Latency + 100ms Network Latency**

End-to-End Latency with Network (ms)

- 11.8s
- 11.9s
- 6.2s
- 7.2s
- 4.3s
- 4.7s
- 2.1s
- 2.4s

No CC    CC

LLama 3.3 70B Int4    GPT OSS 120B    Mistral 3.1 24B    LLama 3.1 8B

# ShareGPT (Single Request)

## E2E Latency + 100ms Network Latency

**End-to-End Latency with Network (ms)**

| | LLama 3.3 70B Int4 | GPT OSS 120B | Mistral 3.1 24B | LLama 3.1 8B |
|---|---|---|---|---|
| No CC | 4.4s | 1.7s | 3.7s | 1.6s |
| CC | 4.6s | 1.8s | 3.8s | 1.7s |

Legend: No CC, CC

LLama 3.3 70B Int4  GPT OSS 120B  Mistral 3.1 24B  LLama 3.1 8B

# Edit 10K Characters (100 Request Rate)

## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

- 445.2s
- 447.1s
- 192.8s
- 203.7s
- 199.1s
- 209.7s
- 83.2s
- 90.8s

Legend: No CC | CC

LLama 3.3 70B Int4 | GPT OSS 120B | Mistral 3.1 24B | LLama 3.1 8B

# Edit 10K Characters (50 Request Rate)

## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

- 439.5s
- 441.6s
- 191.7s
- 203.4s
- 197.2s
- 207.9s
- 78.3s
- 85.8s

No CC    CC

LLama 3.3 70B Int4    GPT OSS 120B    Mistral 3.1 24B    LLama 3.1 8B

**Edit 10K Characters (Single Request)**

E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

LLama 3.3 70B Int4: 360.4s (No CC), 368.0s (CC)
GPT OSS 120B: 126.4s (No CC), 139.1s (CC)
Mistral 3.1 24B: 133.4s (No CC), 144.6s (CC)
LLama 3.1 8B: 32.4s (No CC), 35.9s (CC)

No CC    CC

LLama 3.3 70B Int4    GPT OSS 120B    Mistral 3.1 24B    LLama 3.1 8B

**Numina Math (100 Request Rate)**

E2E Latency + 100ms Network Latency

Legend: ■ No CC  ▨ CC

■ LLama 3.3 70B Int4   ■ GPT OSS 120B   ■ Mistral 3.1 24B   ■ LLama 3.1 8B

Y-axis: End-to-End Latency with Network (ms)

Values:
- LLama 3.3 70B Int4: No CC 34.2s, CC 34.4s
- GPT OSS 120B: No CC 13.3s, CC 14.2s
- Mistral 3.1 24B: No CC 7.9s, CC 8.5s
- LLama 3.1 8B: No CC 4.5s, CC 5.3s

**Numina Math (50 Request Rate)**

E2E Latency + 100ms Network Latency

# Numina Math (Single Request)

## E2E Latency + 100ms Network Latency



**End-to-End Latency with Network (ms)**

- 7.5s
- 7.8s
- 3.4s
- 3.6s
- 5.8s
- 6.2s
- 2.8s
- 3.0s

Legend: No CC, CC

LLama 3.3 70B Int4 · GPT OSS 120B · Mistral 3.1 24B · LLama 3.1 8B