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

E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

- No CC
- CC

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

Values shown:
- 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)

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

**E2E Latency + 100ms Network Latency**

End-to-End Latency with Network (ms)

119.2s — 117.1s — 11.0s — 10.8s — 27.9s — 28.6s — 4.1s — 4.8s

No CC — CC

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

# Random (1500 ⇒ 250) (Request Rate 1)
## E2E Latency + 100ms Network Latency

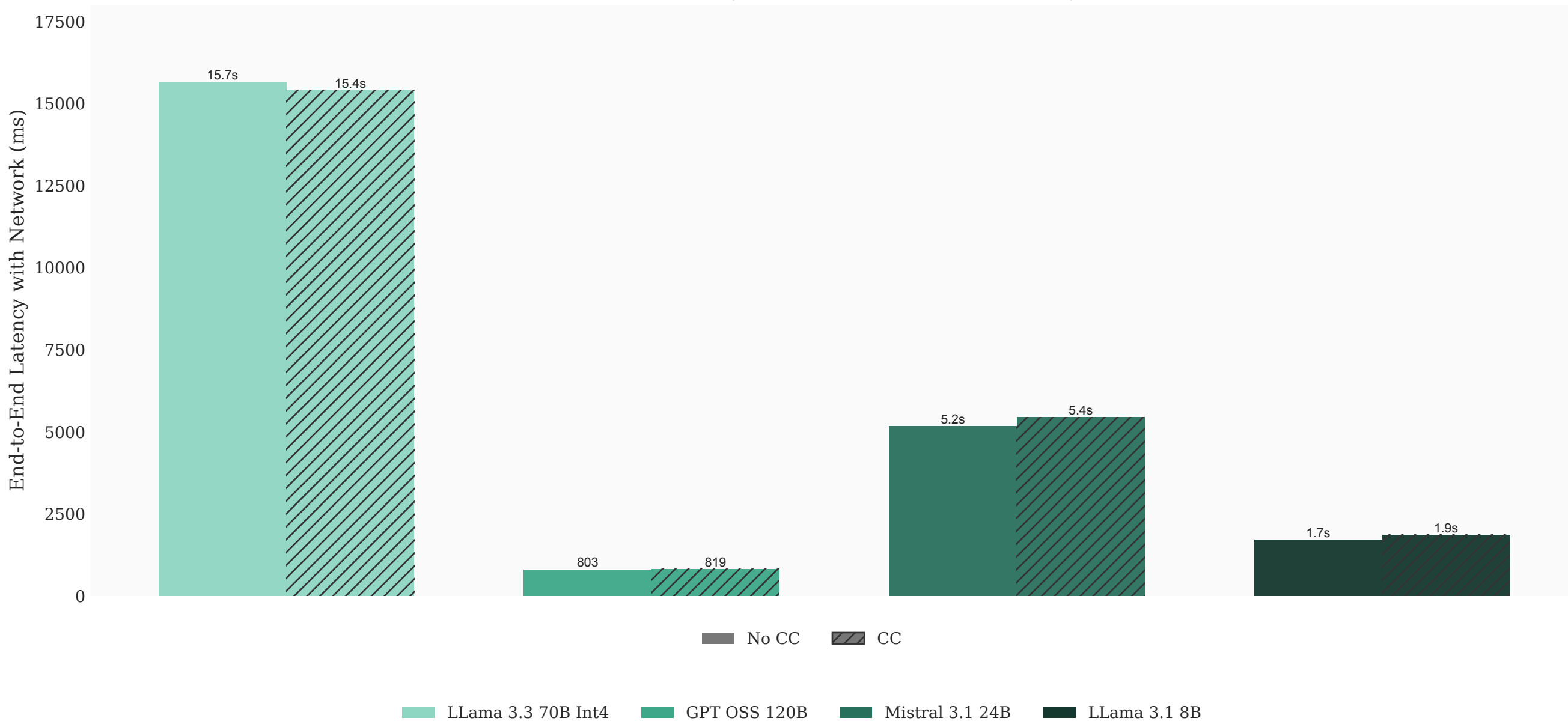End-to-End Latency with Network (ms)

15.7s 15.4s 803 819 5.2s 5.4s 1.7s 1.9s

No CC    CC

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

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

## E2E Latency + 100ms Network Latency

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

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

Legend: No CC, CC

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

Random (4000 ⇒ 1000) (Request Rate 50)

E2E Latency + 100ms Network Latency

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) (Request Rate 1)**

**E2E Latency + 100ms Network Latency**

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

- LLama 3.3 70B Int4: No CC 182.5s, CC 180.2s
- GPT OSS 120B: No CC 3.0s, CC 3.7s
- Mistral 3.1 24B: No CC 35.7s, CC 40.0s
- LLama 3.1 8B: No CC 9.0s, CC 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) (Request Rate 100)

## E2E Latency + 100ms Network Latency



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

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

Legend: No CC, CC

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

# Random (1000 ⇒ 1000) (Request Rate 50)
## E2E Latency + 100ms Network Latency



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

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

Legend: No CC, CC

# Random (1000 ⇒ 1000) (Request Rate 1)
## E2E Latency + 100ms Network Latency



Bar chart. Y-axis: End-to-End Latency with Network (ms), ranging from 0 to over 20000.

Legend (patterns): No CC (solid), CC (hatched)

Legend (colors): LLama 3.3 70B Int4, GPT OSS 120B, Mistral 3.1 24B, LLama 3.1 8B

Values:
- 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

**ShareGPT (Request Rate 100)**

**E2E Latency + 100ms Network Latency**

End-to-End Latency with Network (ms)

| | No CC | CC |
|---|---|---|
| LLama 3.3 70B Int4 | 29.3s | 29.2s |
| GPT OSS 120B | 7.5s | 7.9s |
| Mistral 3.1 24B | 6.2s | 6.6s |
| LLama 3.1 8B | 2.6s | 2.9s |

No CC    CC

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

**ShareGPT (Request Rate 50)**

**E2E Latency + 100ms Network Latency**

End-to-End Latency with Network (ms)

- 11.8s / 11.9s — LLama 3.3 70B Int4
- 6.2s / 7.2s — GPT OSS 120B
- 4.3s / 4.7s — Mistral 3.1 24B
- 2.1s / 2.4s — 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 (Request Rate 1)

## E2E Latency + 100ms Network Latency



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

- 4.4s
- 4.6s
- 1.7s
- 1.8s
- 3.7s
- 3.8s
- 1.6s
- 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 (Request Rate 100)

## 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 (Request Rate 50)

## 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

Legend: No CC, CC

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

**Edit 10K Characters (Request Rate 1)**

E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

| Model | No CC | CC |
|---|---|---|
| LLama 3.3 70B Int4 | 360.4s | 368.0s |
| GPT OSS 120B | 126.4s | 139.1s |
| Mistral 3.1 24B | 133.4s | 144.6s |
| LLama 3.1 8B | 32.4s | 35.9s |

No CC    CC

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

# Numina Math (Request Rate 100)
## E2E Latency + 100ms Network Latency



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

- 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

Legend: No CC, CC

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

# Numina Math (Request Rate 50)
## E2E Latency + 100ms Network Latency



Bar chart. Y-axis: End-to-End Latency with Network (ms), ranging from 0 to 30000.

- LLama 3.3 70B Int4: No CC 27.9s, CC 27.9s
- GPT OSS 120B: No CC 12.8s, CC 13.9s
- Mistral 3.1 24B: No CC 7.7s, CC 8.0s
- LLama 3.1 8B: No CC 4.0s, CC 4.8s

Legend: No CC, CC (hatched)

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

# Numina Math (Request Rate 1)

## 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