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

## E2E Latency + 100ms Network Latency



Legend: CC, No CC

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

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

Values: 11.8s, 11.1s, 29.7s, 28.9s, 12.5s, 11.8s, 118.4s, 120.2s

# Random (1500 ⇒ 250) (50 Request Rate)
## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

- 4.8s / 4.1s — LLama 3.1 8B
- 28.6s / 27.9s — Mistral 3.1 24B
- 10.8s / 11.0s — GPT OSS 120B
- 117.1s / 119.2s — LLama 3.3 70B Int4

Legend: CC | No CC

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

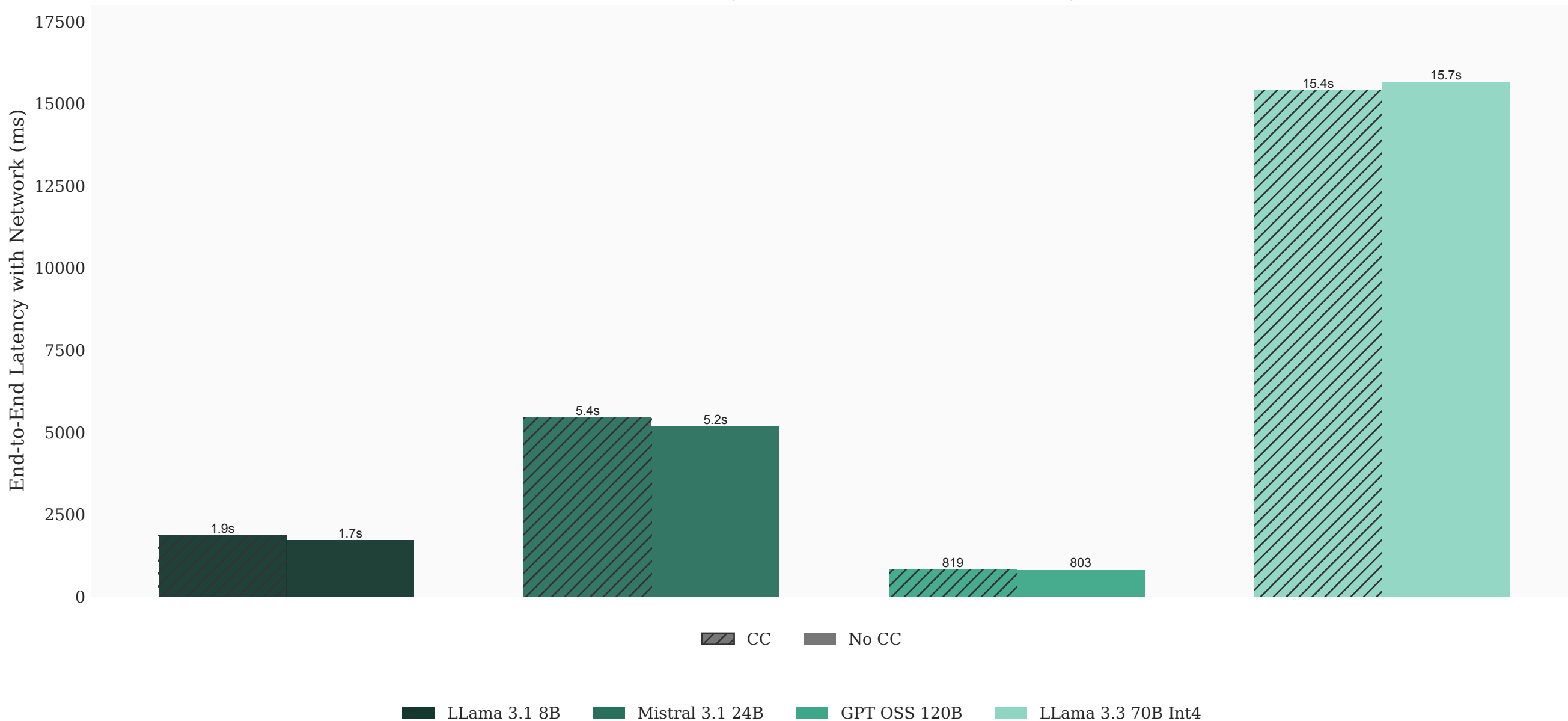# Random (1500 ⇒ 250) (1 Request Rate)

## E2E Latency + 100ms Network Latency



Legend: CC, No CC

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

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

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

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

## E2E Latency + 100ms Network Latency



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

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

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

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

## E2E Latency + 100ms Network Latency



Bar chart showing End-to-End Latency with Network (ms) for four models with CC and No CC variants:
- LLama 3.1 8B: CC 52.9s, No CC 50.2s
- Mistral 3.1 24B: CC 125.6s, No CC 121.0s
- GPT OSS 120B: CC 36.1s, No CC 34.0s
- LLama 3.3 70B Int4: CC 277.5s, No CC 279.9s

Legend: CC, No CC

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

# Random (4000 ⇒ 1000) (1 Request Rate)

## E2E Latency + 100ms Network Latency



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

- LLama 3.1 8B: CC 10.0s, No CC 9.0s
- Mistral 3.1 24B: CC 40.0s, No CC 35.7s
- GPT OSS 120B: CC 3.7s, No CC 3.0s
- LLama 3.3 70B Int4: CC 180.2s, No CC 182.5s

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

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

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

## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

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

Legend: CC | No CC

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

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

E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

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

CC    No CC

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

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



Bar chart showing End-to-End Latency with Network (ms). CC (hatched) vs No CC (solid) for four models:
- LLama 3.1 8B: CC 7.9s, No CC 7.1s
- Mistral 3.1 24B: CC 21.4s, No CC 19.9s
- GPT OSS 120B: CC 2.7s, No CC 2.7s
- LLama 3.3 70B Int4: CC 19.5s, No CC 18.8s

Legend: CC, No CC; LLama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, LLama 3.3 70B Int4

# ShareGPT (100 Request Rate)
## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

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

Legend:
- CC
- No CC
- LLama 3.1 8B
- Mistral 3.1 24B
- GPT OSS 120B
- LLama 3.3 70B Int4

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



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

- LLama 3.1 8B: CC 2.4s, No CC 2.1s
- Mistral 3.1 24B: CC 4.7s, No CC 4.3s
- GPT OSS 120B: CC 7.2s, No CC 6.2s
- LLama 3.3 70B Int4: CC 11.9s, No CC 11.8s

Legend: CC (hatched), No CC (solid); LLama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, LLama 3.3 70B Int4

**ShareGPT (1 Request Rate)**

E2E Latency + 100ms Network Latency

Legend: CC, No CC

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

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

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

# Edit 10K Characters (100 Request Rate)
## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

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

Legend: CC | No CC

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

# Edit 10K Characters (50 Request Rate)

## E2E Latency + 100ms Network Latency



Bar chart showing End-to-End Latency with Network (ms) for four models with CC and No CC:
- LLama 3.1 8B: CC 85.8s, No CC 78.3s
- Mistral 3.1 24B: CC 207.9s, No CC 197.2s
- GPT OSS 120B: CC 203.4s, No CC 191.7s
- LLama 3.3 70B Int4: CC 441.6s, No CC 439.5s

Legend: CC, No CC

# Edit 10K Characters (1 Request Rate)
## E2E Latency + 100ms Network Latency

End-to-End Latency with Network (ms)

- 35.9s
- 32.4s
- 144.6s
- 133.4s
- 139.1s
- 126.4s
- 368.0s
- 360.4s

CC | No CC

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

# Numina Math (100 Request Rate)

## E2E Latency + 100ms Network Latency



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

Legend: CC, No CC

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

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

# Numina Math (50 Request Rate)

## E2E Latency + 100ms Network Latency



End-to-End Latency with Network (ms)

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

Legend: CC, No CC

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

# Numina Math (1 Request Rate)

## E2E Latency + 100ms Network Latency



Bar chart showing End-to-End Latency with Network (ms) for four models (LLama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, LLama 3.3 70B Int4), comparing CC and No CC configurations:

- LLama 3.1 8B: CC 3.0s, No CC 2.8s
- Mistral 3.1 24B: CC 6.2s, No CC 5.8s
- GPT OSS 120B: CC 3.6s, No CC 3.4s
- LLama 3.3 70B Int4: CC 7.8s, No CC 7.5s

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

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