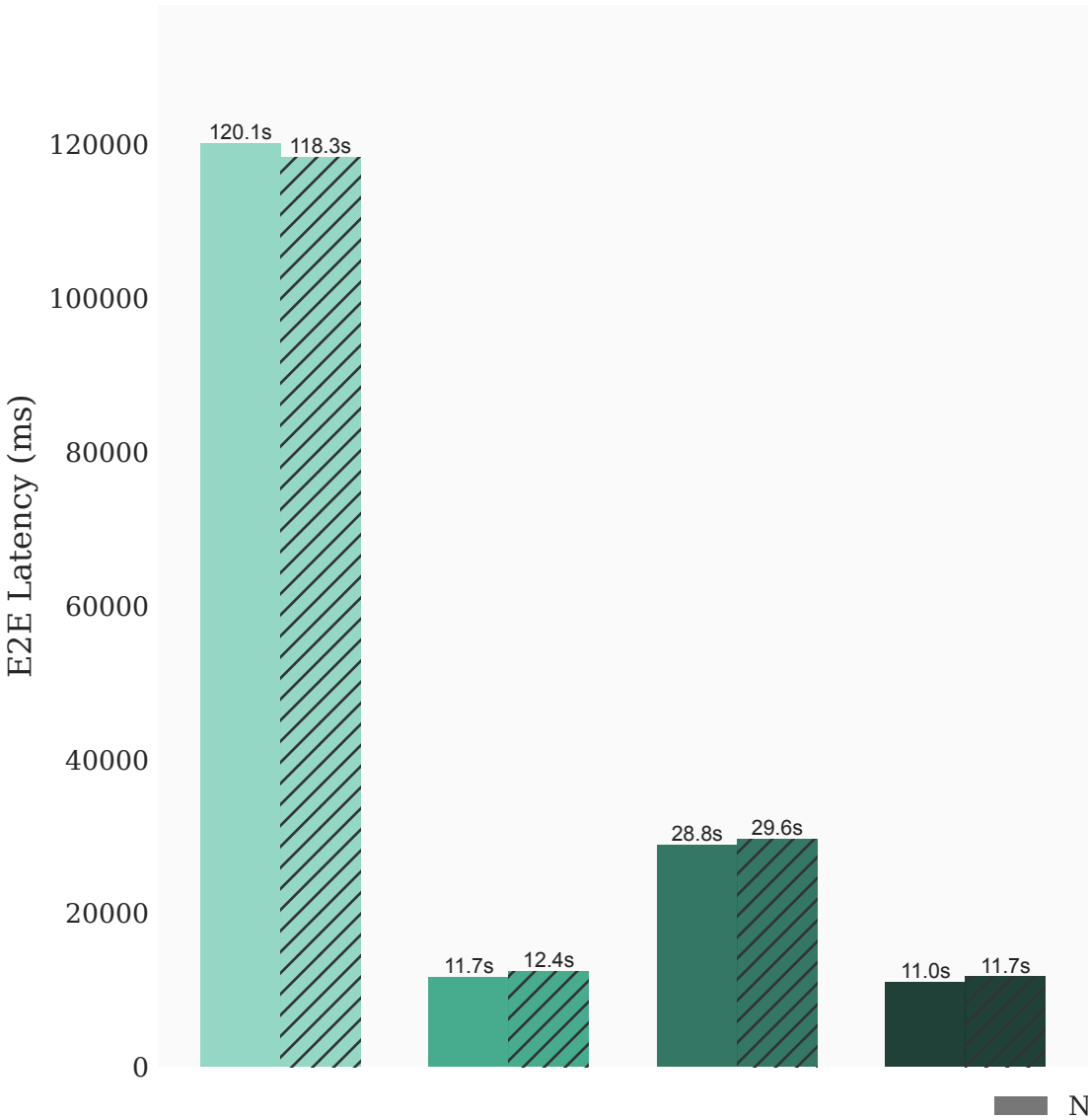# Random (1500 ⇒ 250) (100 Request Rate)

## End-to-End Latency (Mean)



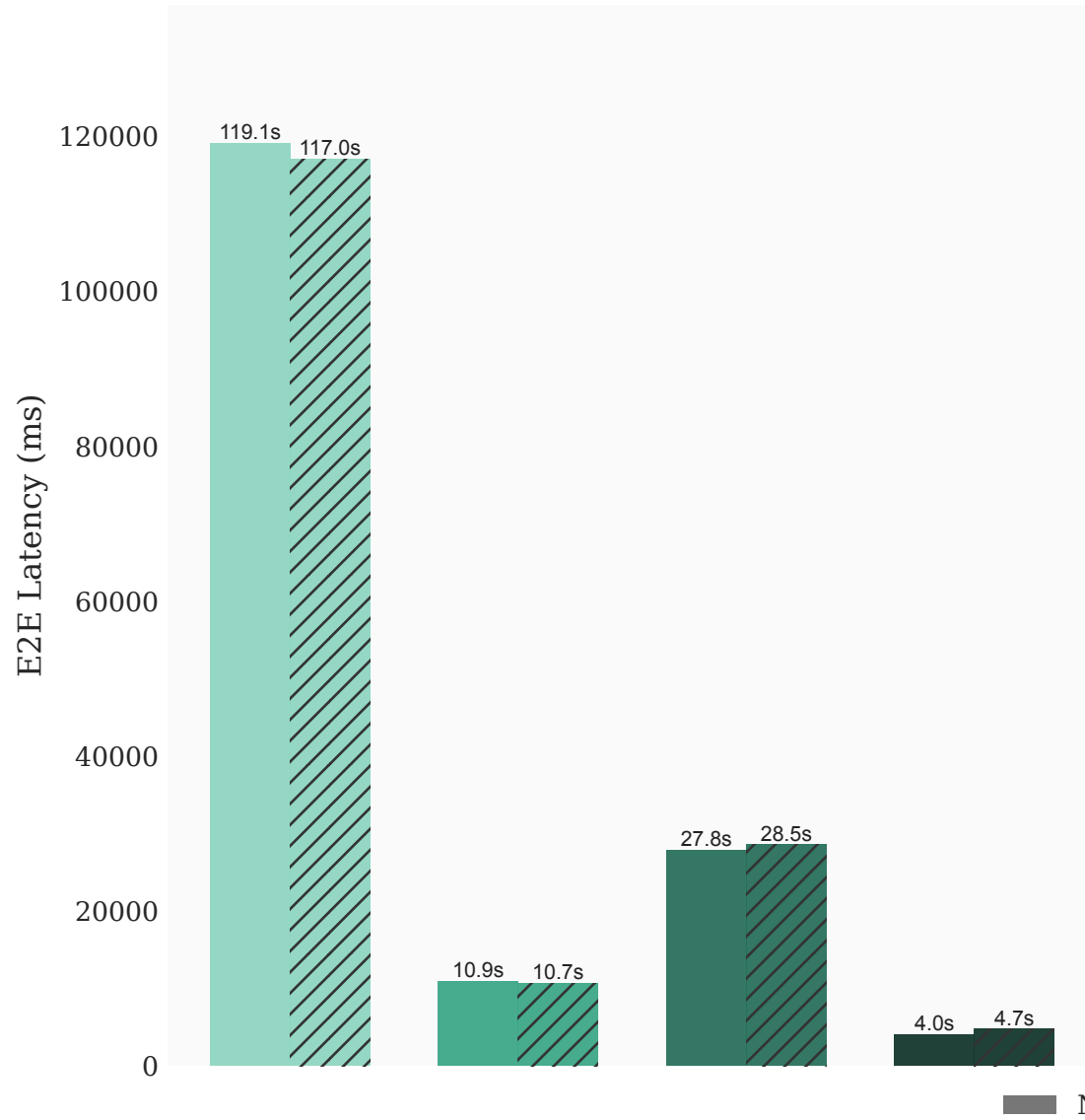## End-to-End Latency (P99)

Legend:
- No CC
- CC

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

Mean values:
- LLama 3.3 70B Int4: 120.1s (No CC), 118.3s (CC)
- GPT OSS 120B: 11.7s (No CC), 12.4s (CC)
- Mistral 3.1 24B: 28.8s (No CC), 29.6s (CC)
- LLama 3.1 8B: 11.0s (No CC), 11.7s (CC)

P99 values:
- LLama 3.3 70B Int4: 164.5s (No CC), 162.6s (CC)
- GPT OSS 120B: 17.5s (No CC), 18.4s (CC)
- Mistral 3.1 24B: 34.8s (No CC), 35.9s (CC)
- LLama 3.1 8B: 12.3s (No CC), 13.0s (CC)

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

## End-to-End Latency (Mean)

E2E Latency (ms)

- 119.1s / 117.0s (LLama 3.3 70B Int4)
- 10.9s / 10.7s (GPT OSS 120B)
- 27.8s / 28.5s (Mistral 3.1 24B)
- 4.0s / 4.7s (LLama 3.1 8B)

## End-to-End Latency (P99)

E2E Latency (ms)

- 162.2s / 160.2s (LLama 3.3 70B Int4)
- 17.6s / 17.7s (GPT OSS 120B)
- 33.6s / 34.6s (Mistral 3.1 24B)
- 4.8s / 5.6s (LLama 3.1 8B)

No CC    CC

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

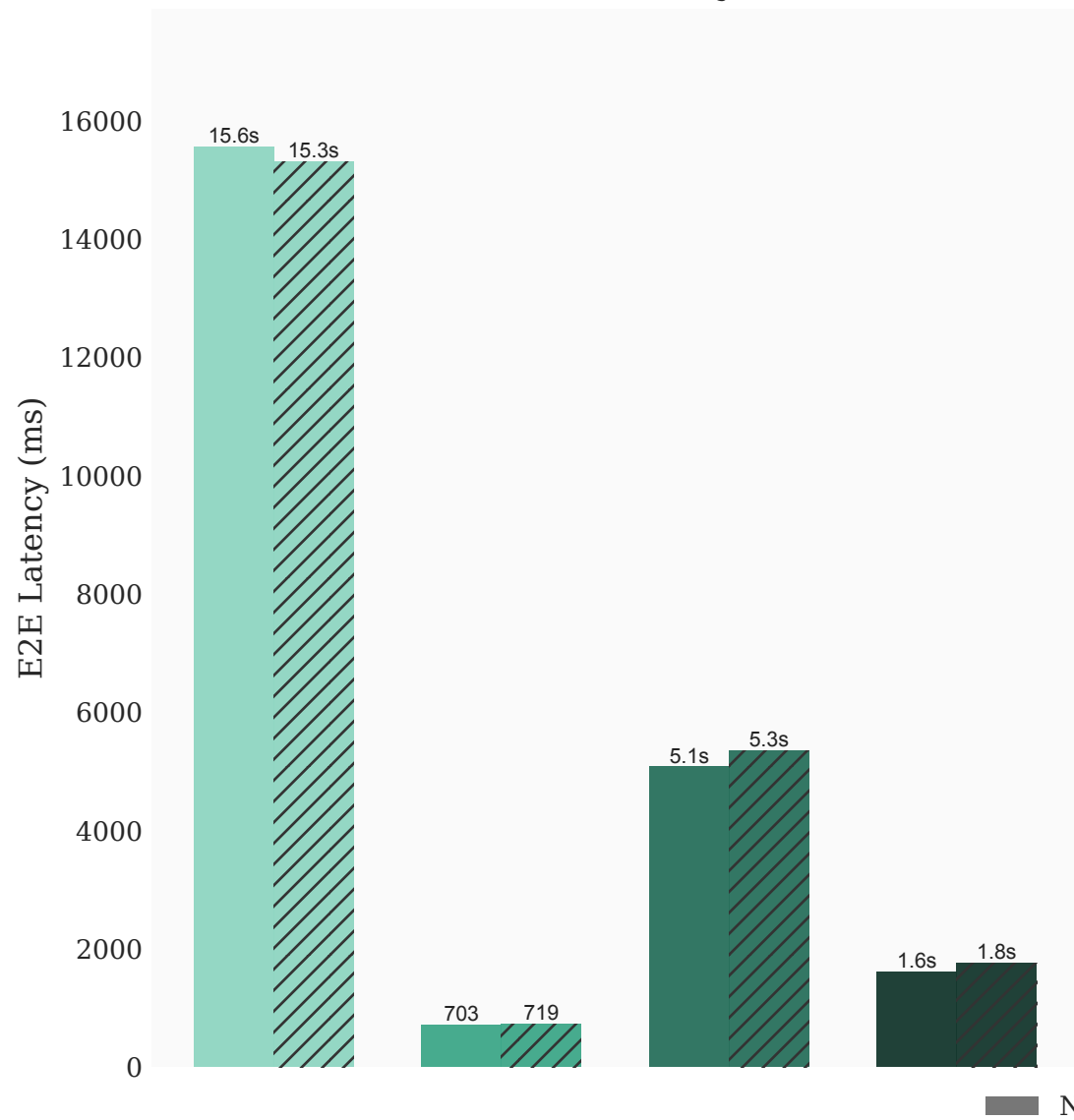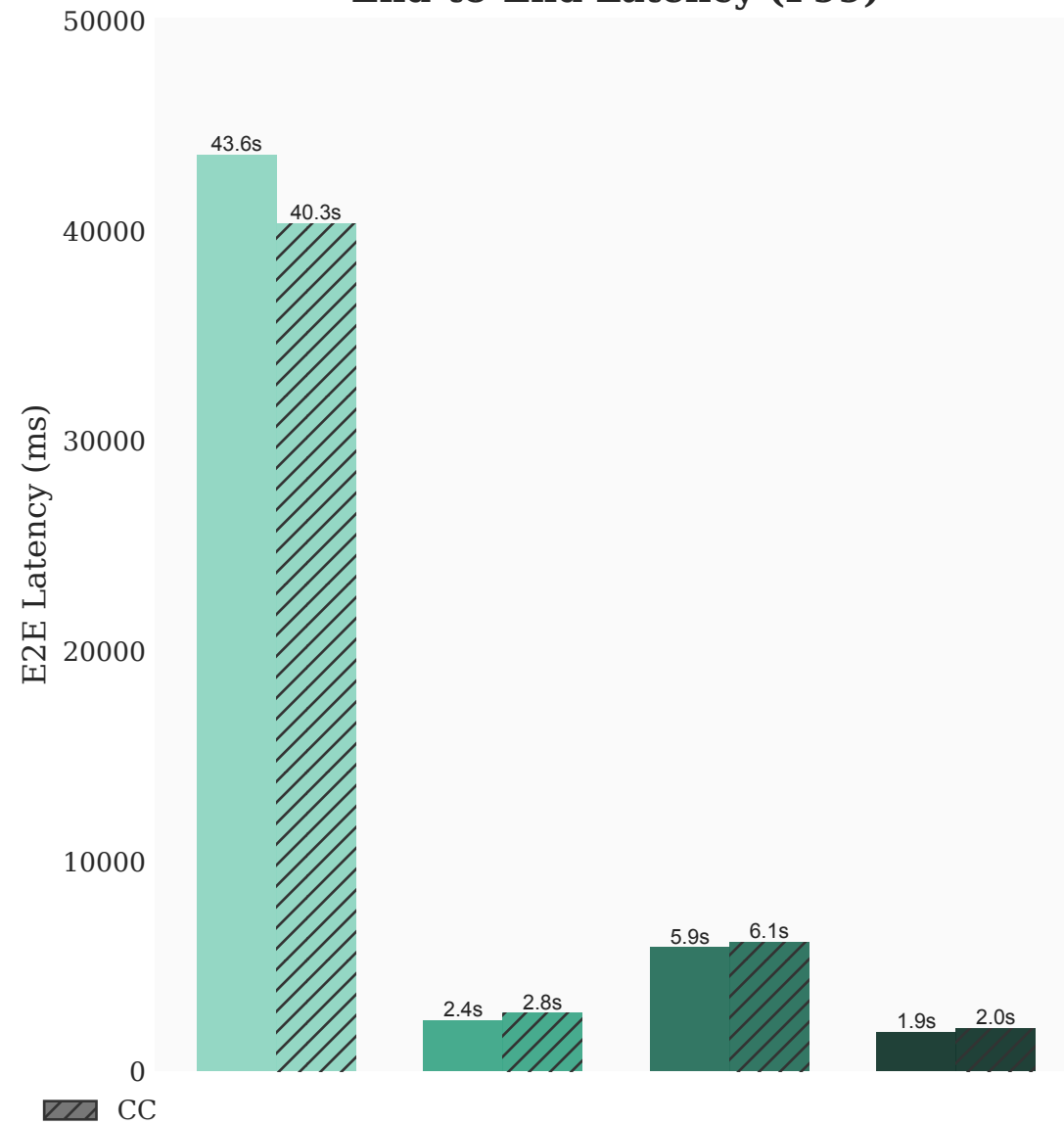# Random (1500 ⇒ 250) (Single Request)

## End-to-End Latency (Mean)



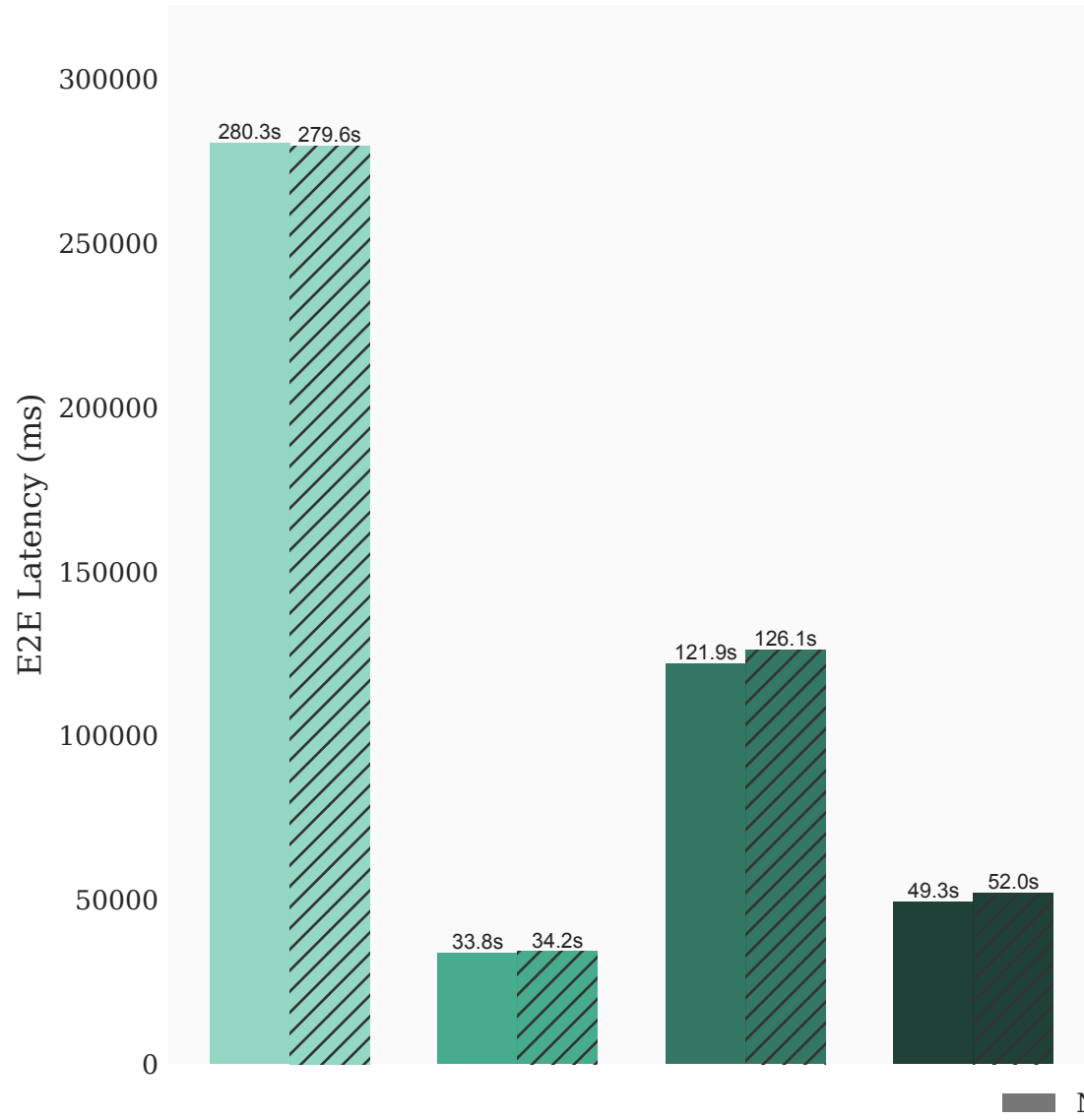E2E Latency (ms)

15.6s  15.3s
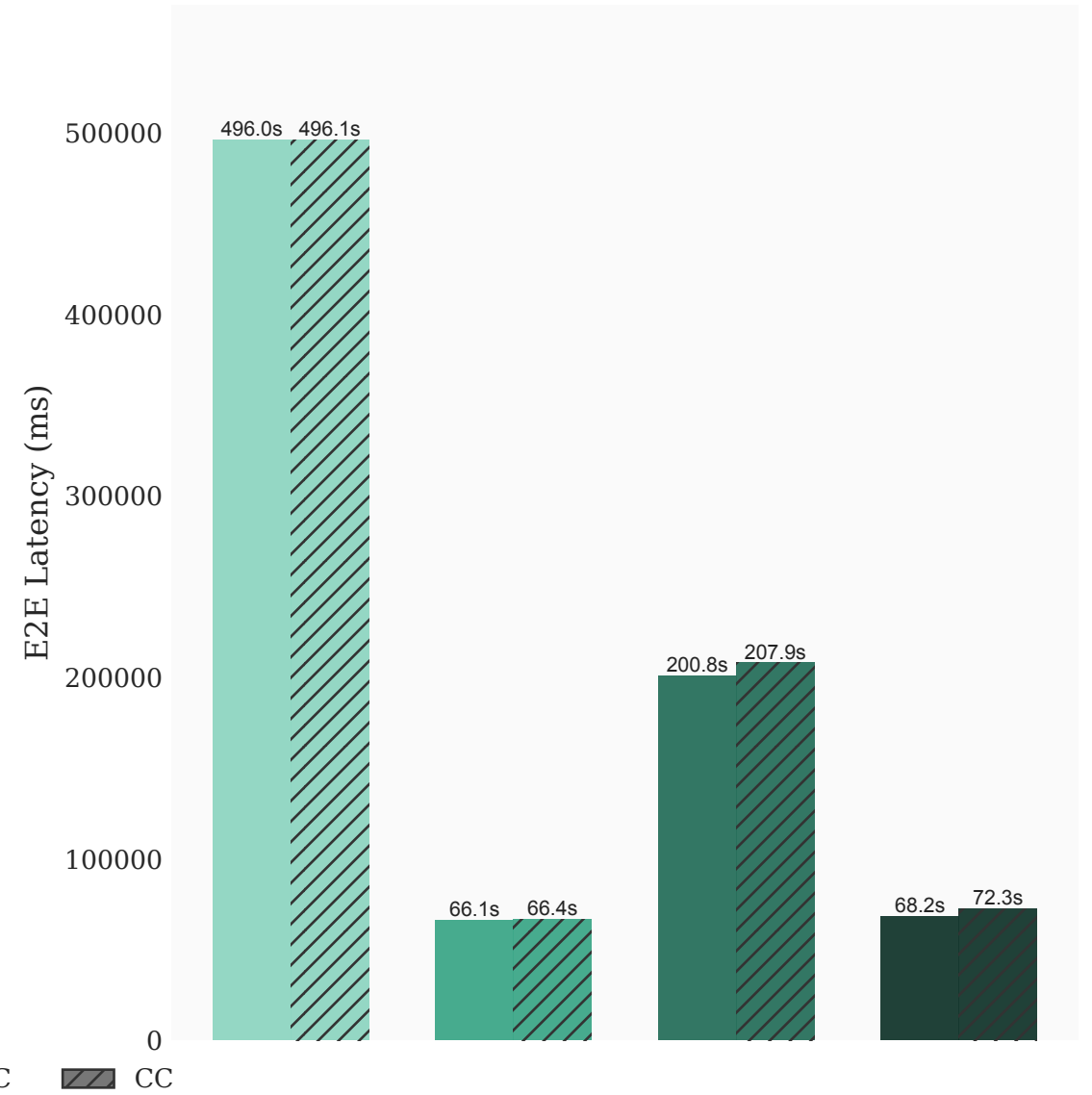
5.1s  5.3s

703  719

1.6s  1.8s

## End-to-End Latency (P99)

E2E Latency (ms)

43.6s  40.3s

2.4s  2.8s

5.9s  6.1s

1.9s  2.0s

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)**

**End-to-End Latency (Mean)**

E2E Latency (ms)

- 280.3s / 279.6s (LLama 3.3 70B Int4)
- 33.8s / 34.2s (GPT OSS 120B)
- 121.9s / 126.1s (Mistral 3.1 24B)
- 49.3s / 52.0s (LLama 3.1 8B)

**End-to-End Latency (P99)**

E2E Latency (ms)

- 496.0s / 496.1s (LLama 3.3 70B Int4)
- 66.1s / 66.4s (GPT OSS 120B)
- 200.8s / 207.9s (Mistral 3.1 24B)
- 68.2s / 72.3s (LLama 3.1 8B)

No CC    CC

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

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

## End-to-End Latency (Mean)

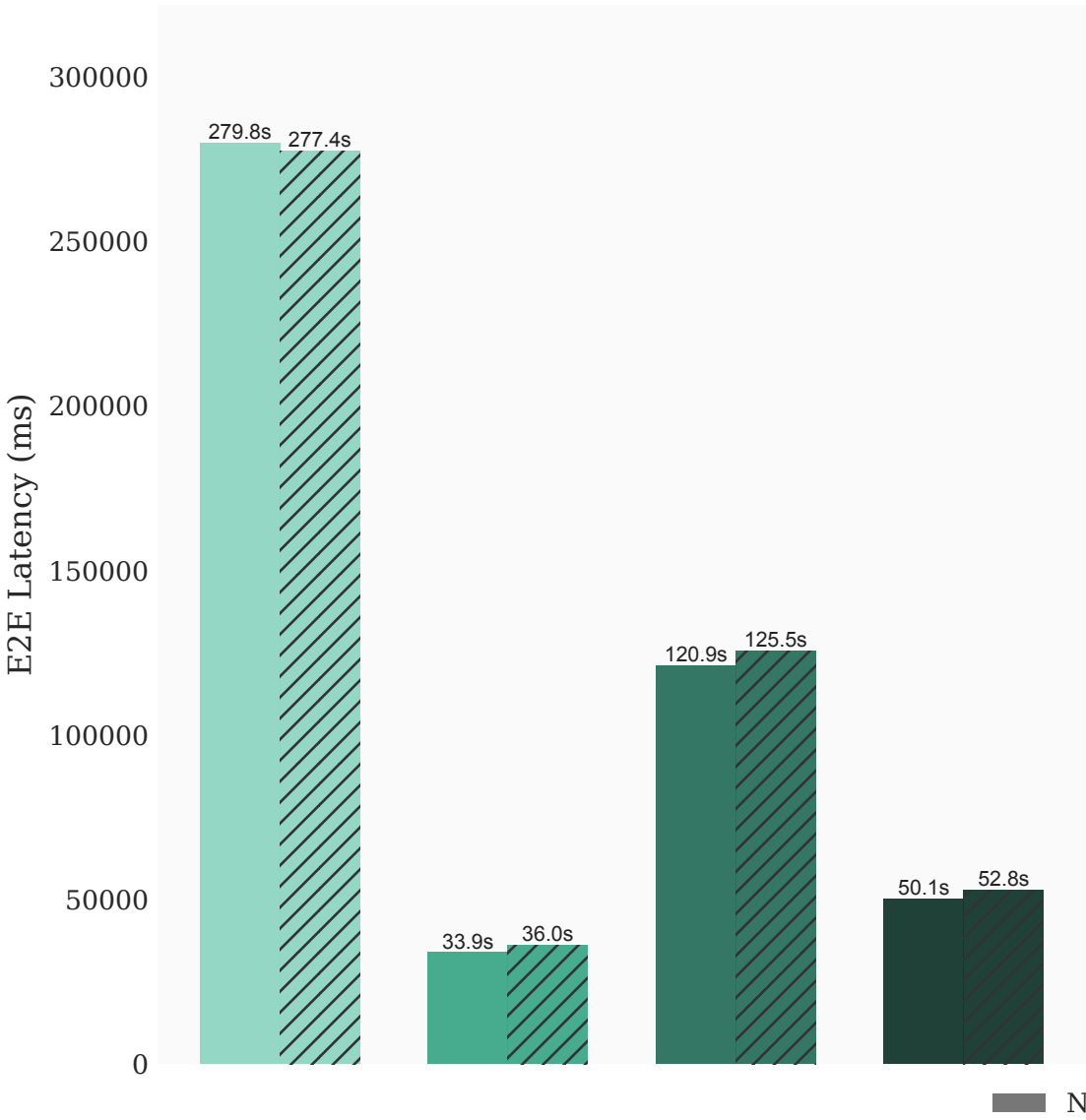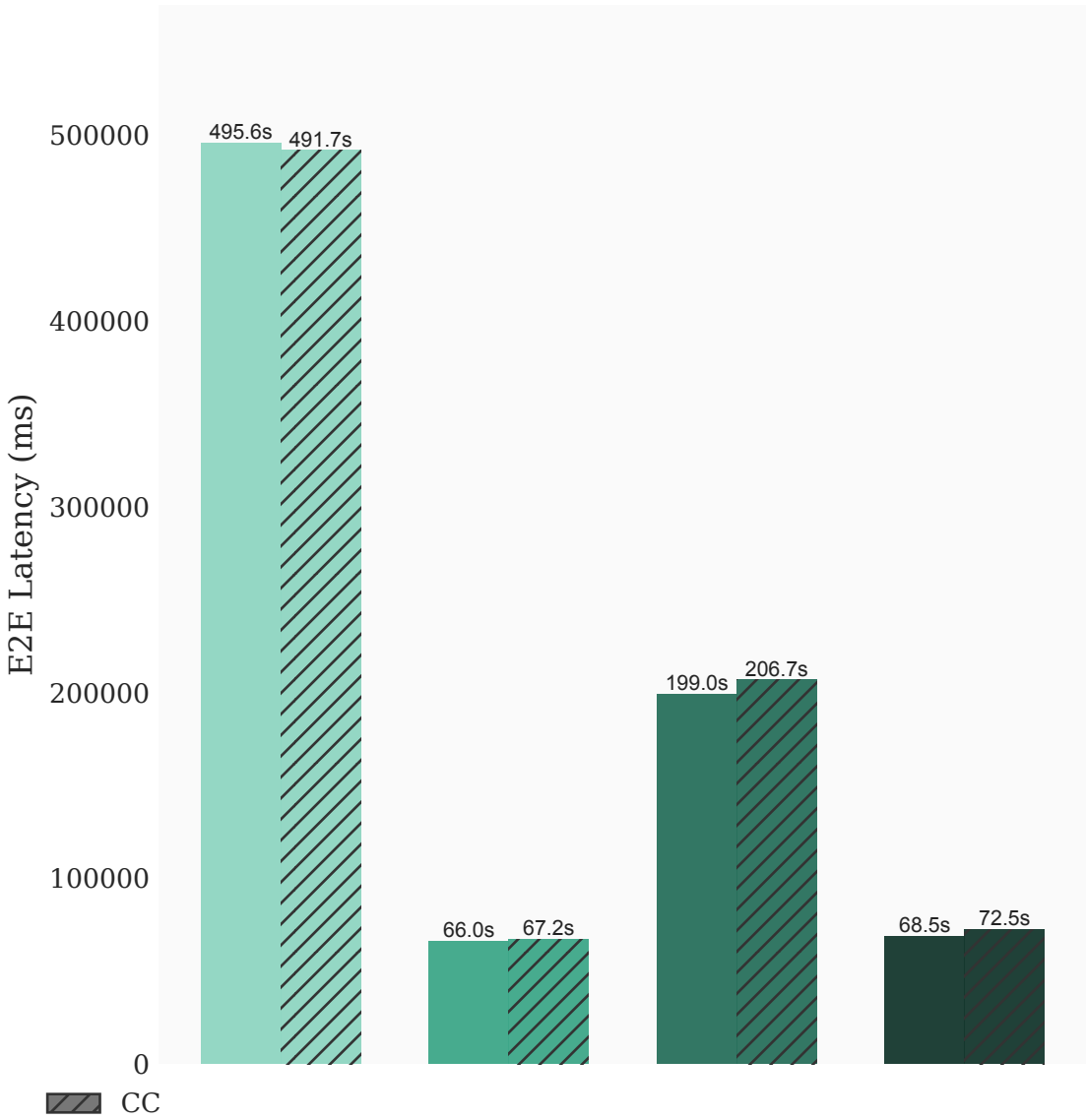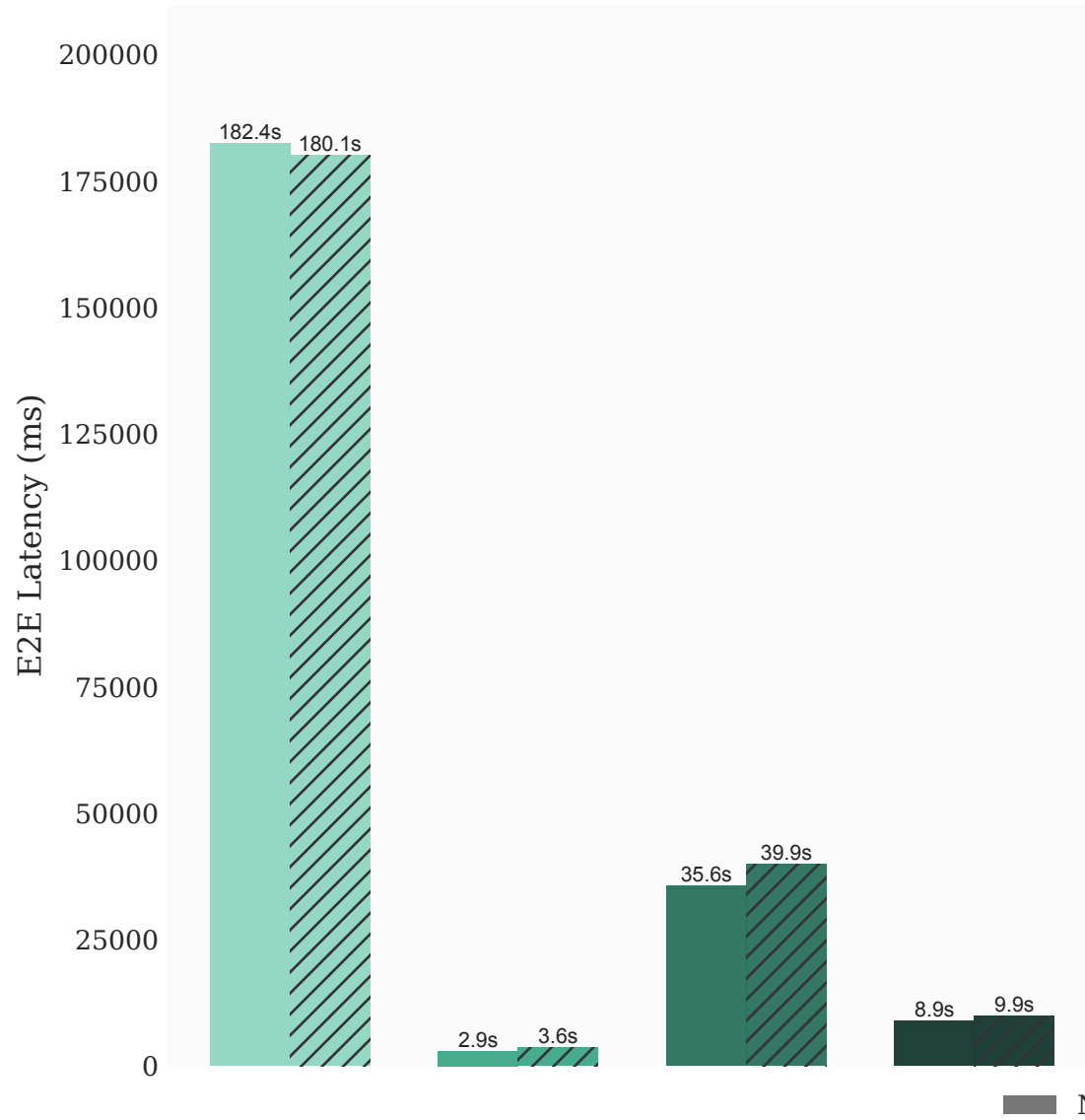E2E Latency (ms)

- 279.8s / 277.4s
- 33.9s / 36.0s
- 120.9s / 125.5s
- 50.1s / 52.8s

## End-to-End Latency (P99)

E2E Latency (ms)

- 495.6s / 491.7s
- 66.0s / 67.2s
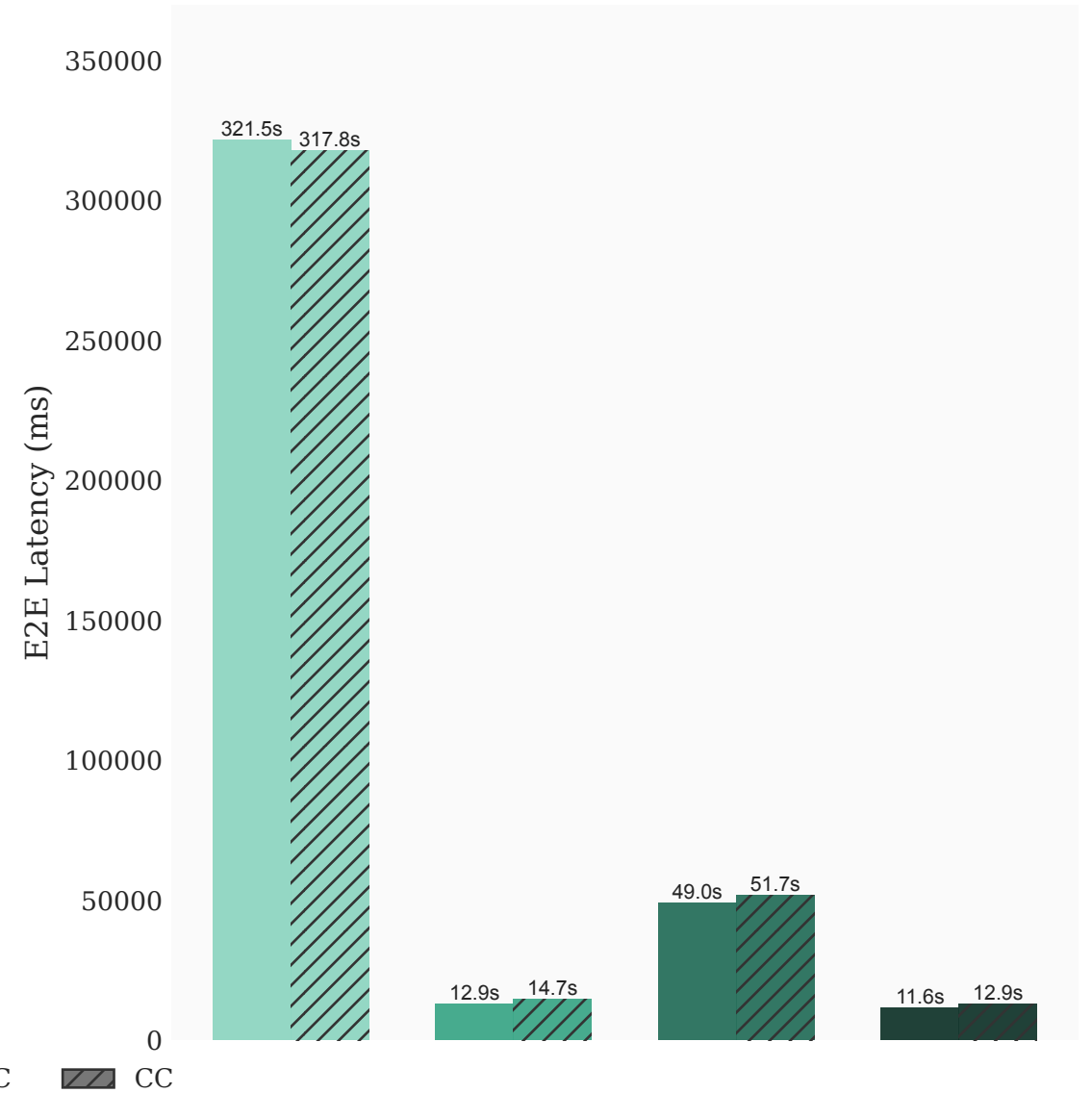- 199.0s / 206.7s
- 68.5s / 72.5s

No CC    CC

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

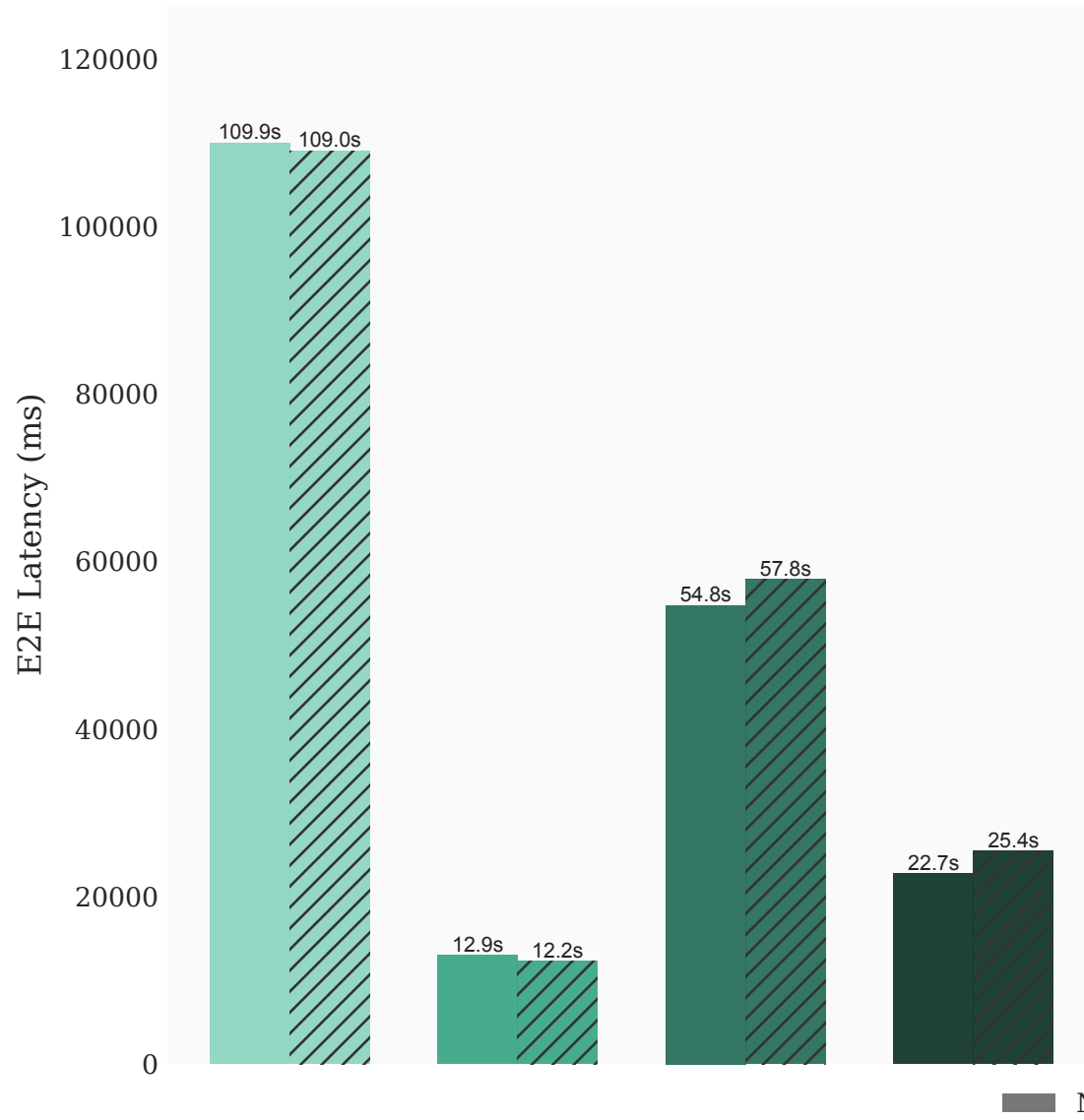# Random (4000 ⇒ 1000) (Single Request)

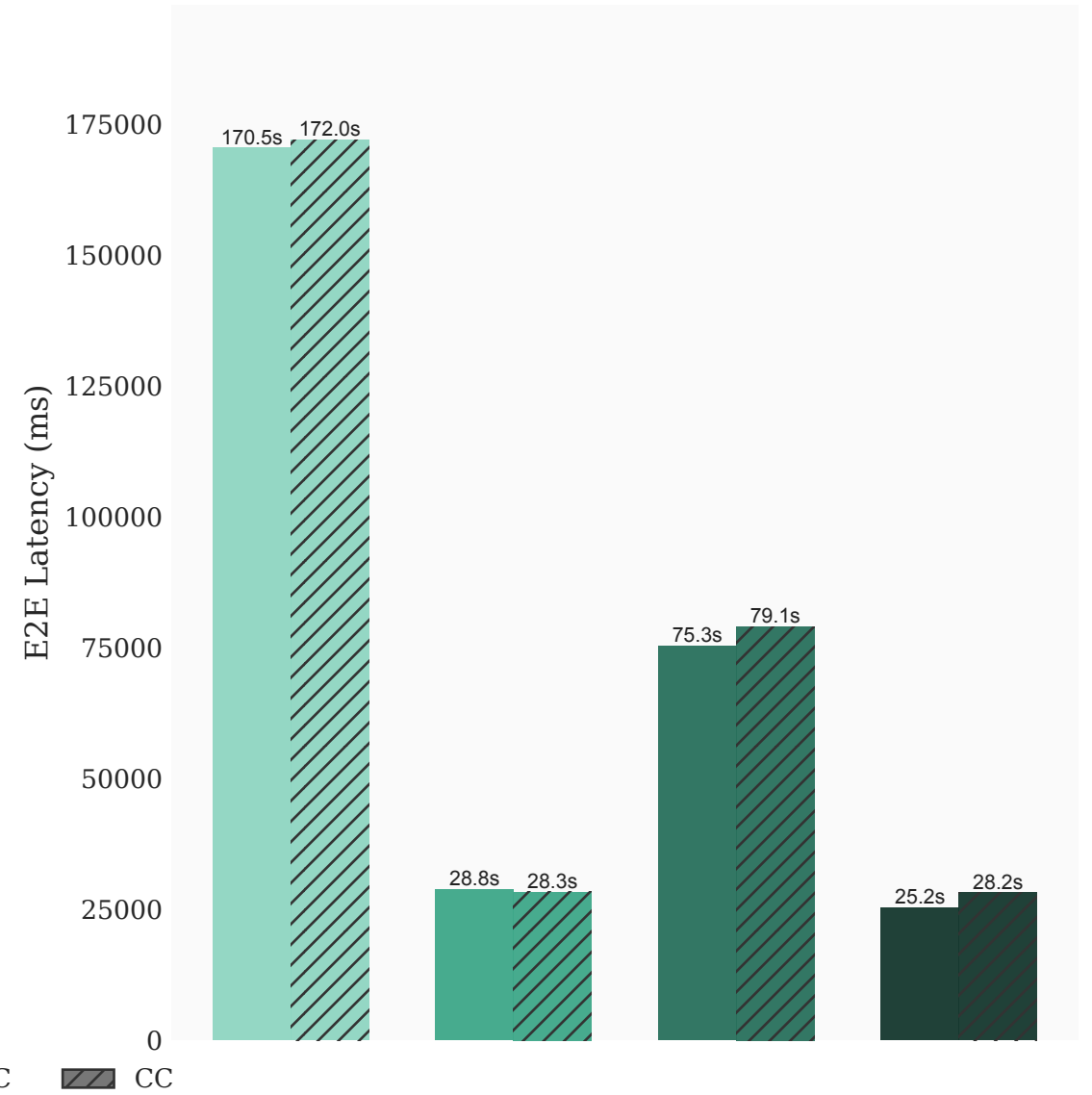## End-to-End Latency (Mean)



## End-to-End Latency (P99)

Legend: No CC, CC

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

Mean values:
- LLama 3.3 70B Int4: 182.4s (No CC), 180.1s (CC)
- GPT OSS 120B: 2.9s (No CC), 3.6s (CC)
- Mistral 3.1 24B: 35.6s (No CC), 39.9s (CC)
- LLama 3.1 8B: 8.9s (No CC), 9.9s (CC)

P99 values:
- LLama 3.3 70B Int4: 321.5s (No CC), 317.8s (CC)
- GPT OSS 120B: 12.9s (No CC), 14.7s (CC)
- Mistral 3.1 24B: 49.0s (No CC), 51.7s (CC)
- LLama 3.1 8B: 11.6s (No CC), 12.9s (CC)

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

## End-to-End Latency (Mean)

E2E Latency (ms)

- 109.9s / 109.0s (LLama 3.3 70B Int4)
- 12.9s / 12.2s (GPT OSS 120B)
- 54.8s / 57.8s (Mistral 3.1 24B)
- 22.7s / 25.4s (LLama 3.1 8B)

## End-to-End Latency (P99)

E2E Latency (ms)

- 170.5s / 172.0s (LLama 3.3 70B Int4)
- 28.8s / 28.3s (GPT OSS 120B)
- 75.3s / 79.1s (Mistral 3.1 24B)
- 25.2s / 28.2s (LLama 3.1 8B)

No CC    CC

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

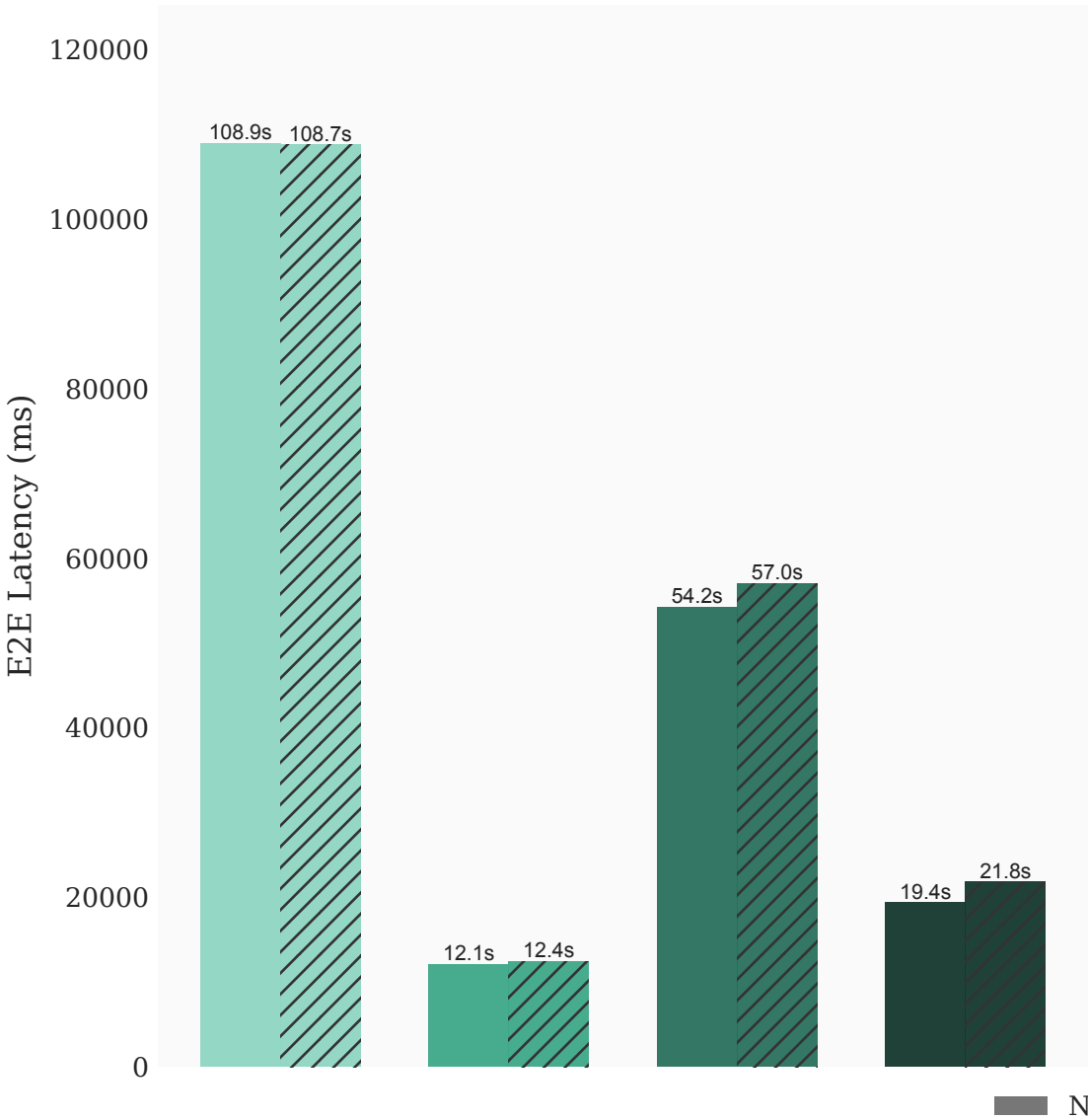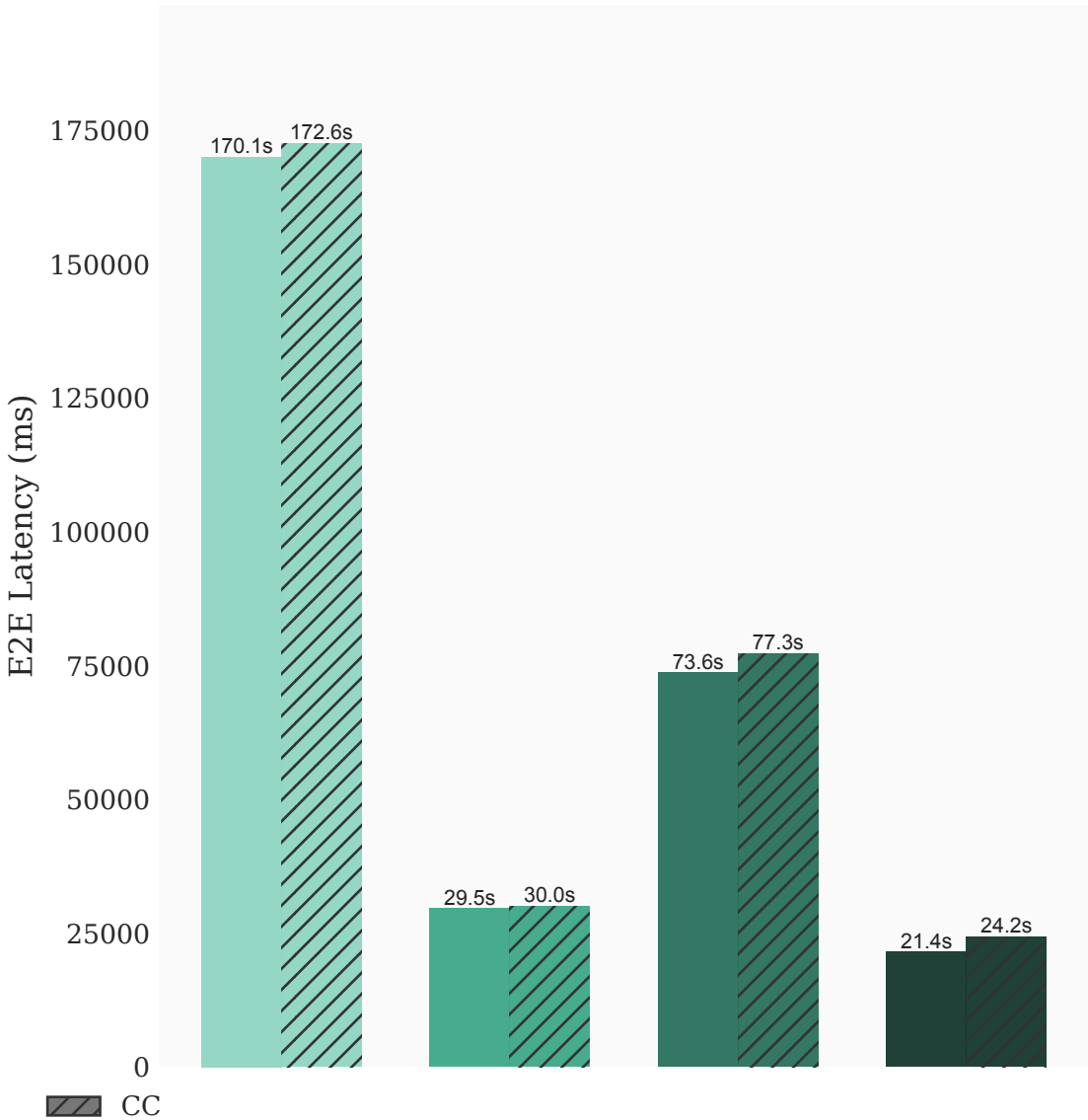# Random (1000 ⇒ 1000) (50 Request Rate)

## End-to-End Latency (Mean)

E2E Latency (ms)

- LLama 3.3 70B Int4: No CC 108.9s, CC 108.7s
- GPT OSS 120B: No CC 12.1s, CC 12.4s
- Mistral 3.1 24B: No CC 54.2s, CC 57.0s
- LLama 3.1 8B: No CC 19.4s, CC 21.8s

## End-to-End Latency (P99)

E2E Latency (ms)

- LLama 3.3 70B Int4: No CC 170.1s, CC 172.6s
- GPT OSS 120B: No CC 29.5s, CC 30.0s
- Mistral 3.1 24B: No CC 73.6s, CC 77.3s
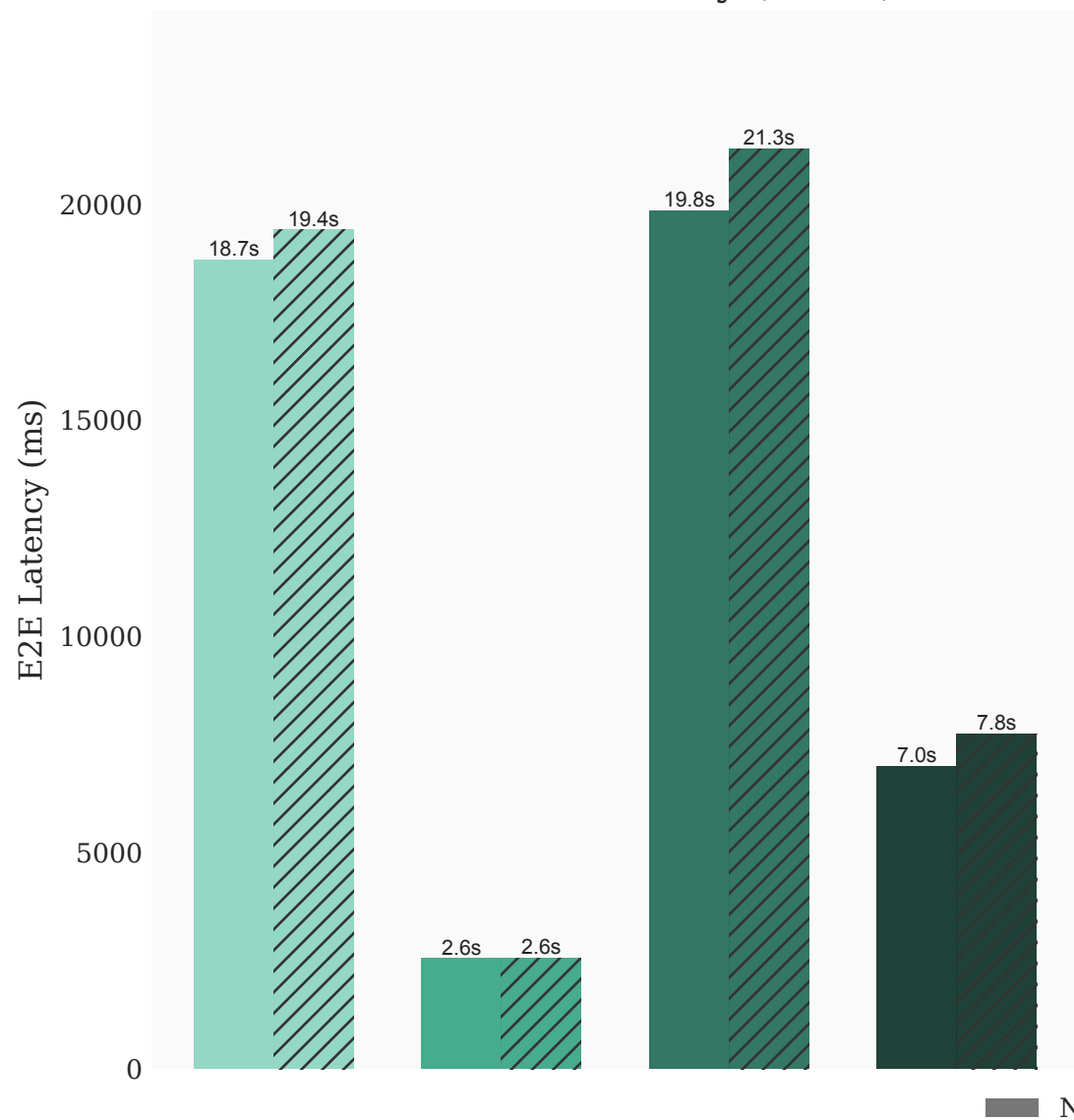- LLama 3.1 8B: No CC 21.4s, CC 24.2s

Legend: No CC, CC

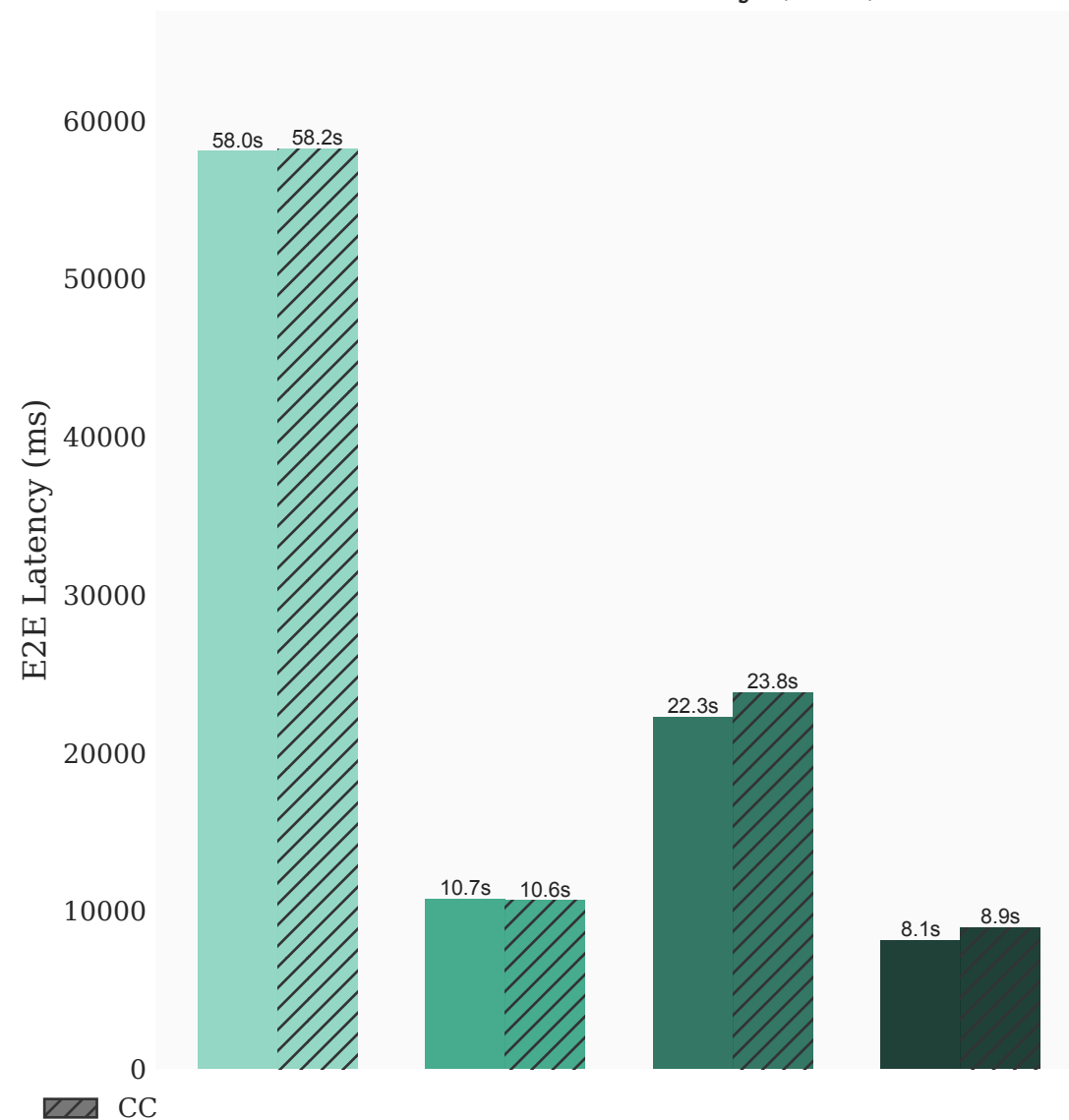LLama 3.3 70B Int4　　GPT OSS 120B　　Mistral 3.1 24B　　LLama 3.1 8B

# Random (1000 ⇒ 1000) (Single Request)

## End-to-End Latency (Mean)



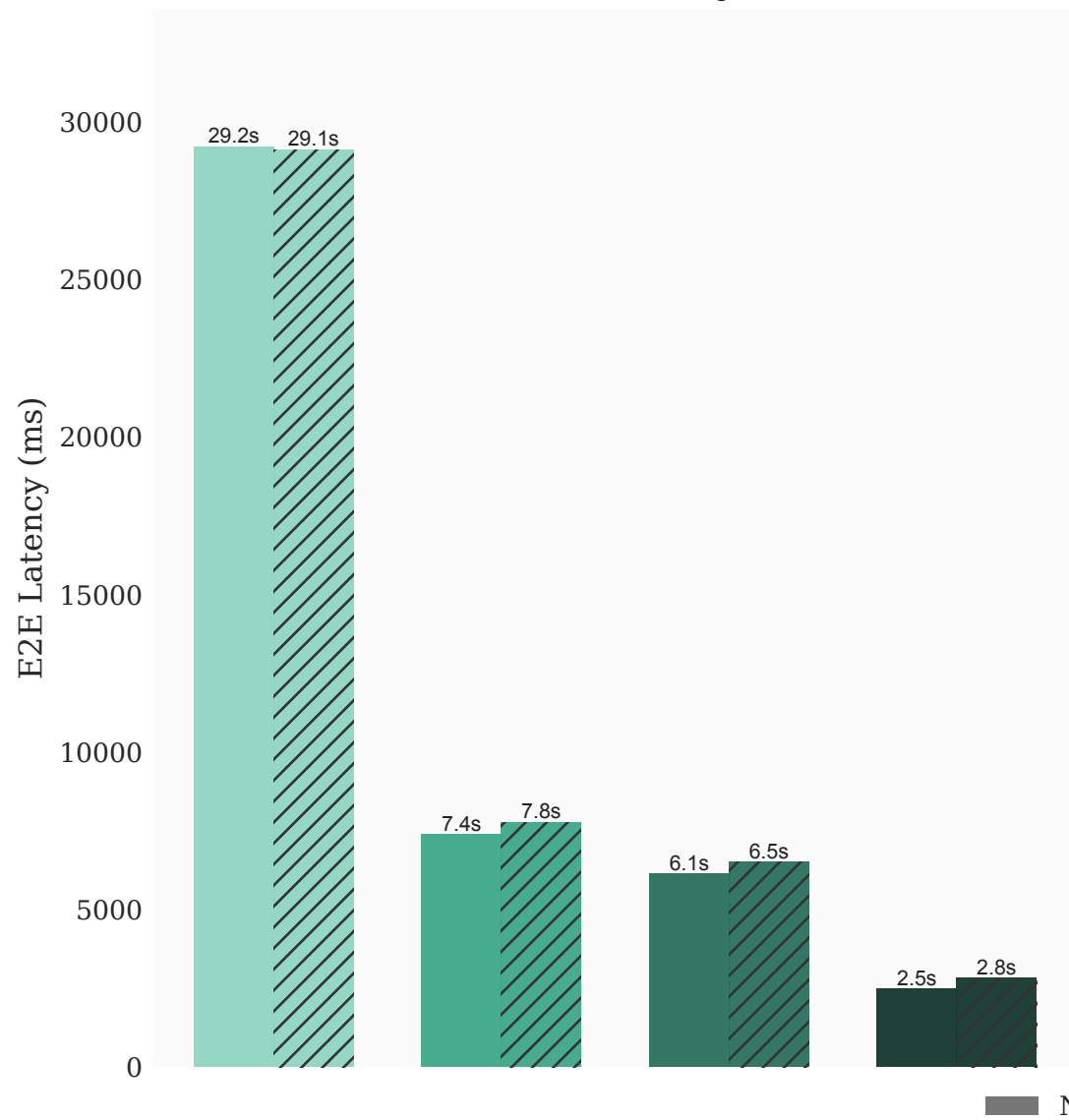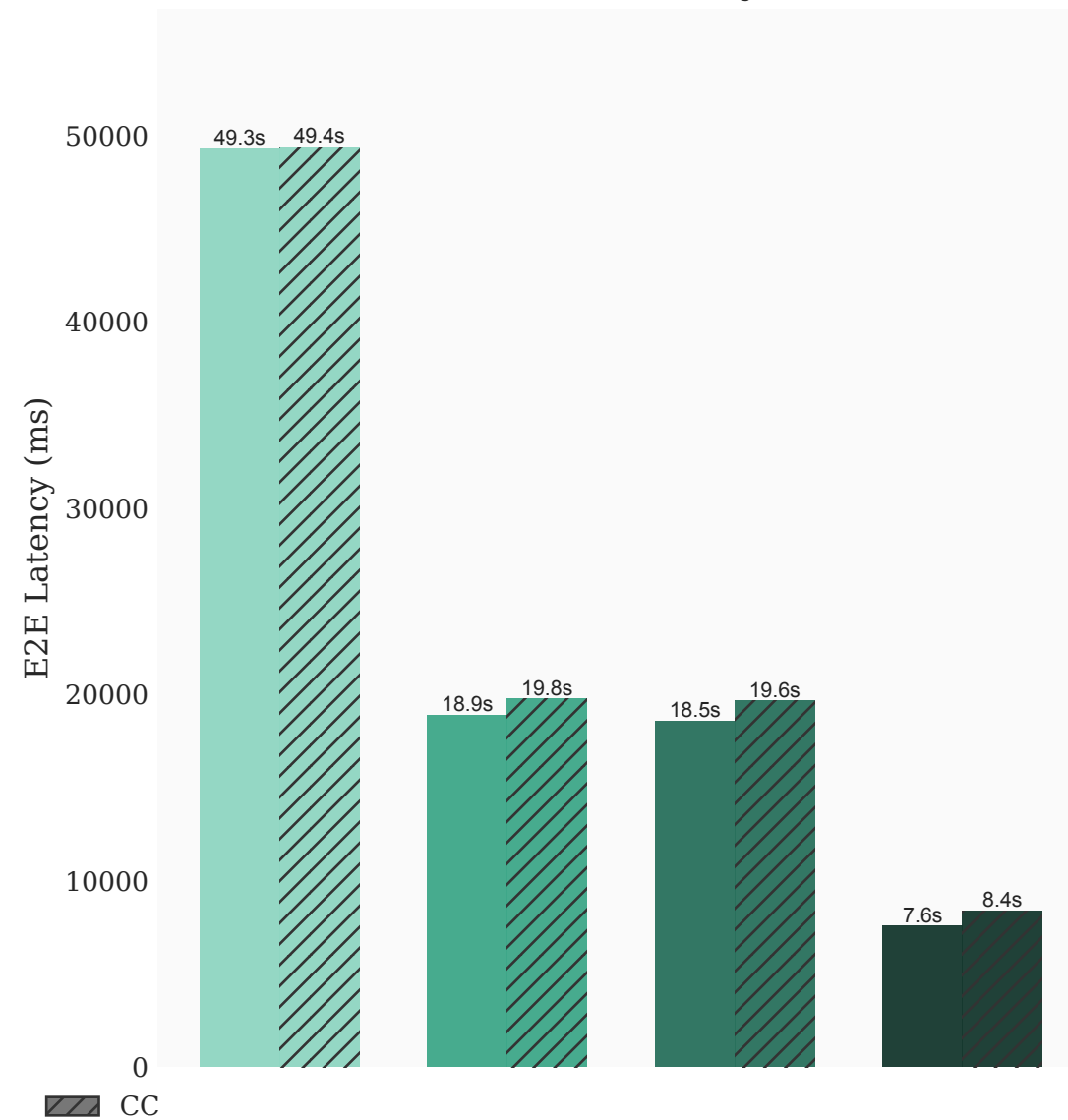## End-to-End Latency (P99)

No CC  ▨ CC

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

# ShareGPT (100 Request Rate)

## End-to-End Latency (Mean)



## End-to-End Latency (P99)

Legend: ■ No CC  ▨ CC

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

Mean values:
- LLama 3.3 70B Int4: 29.2s (No CC), 29.1s (CC)
- GPT OSS 120B: 7.4s (No CC), 7.8s (CC)
- Mistral 3.1 24B: 6.1s (No CC), 6.5s (CC)
- LLama 3.1 8B: 2.5s (No CC), 2.8s (CC)

P99 values:
- LLama 3.3 70B Int4: 49.3s (No CC), 49.4s (CC)
- GPT OSS 120B: 18.9s (No CC), 19.8s (CC)
- Mistral 3.1 24B: 18.5s (No CC), 19.6s (CC)
- LLama 3.1 8B: 7.6s (No CC), 8.4s (CC)

# ShareGPT (50 Request Rate)

## End-to-End Latency (Mean)



## End-to-End Latency (P99)

No CC    CC

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

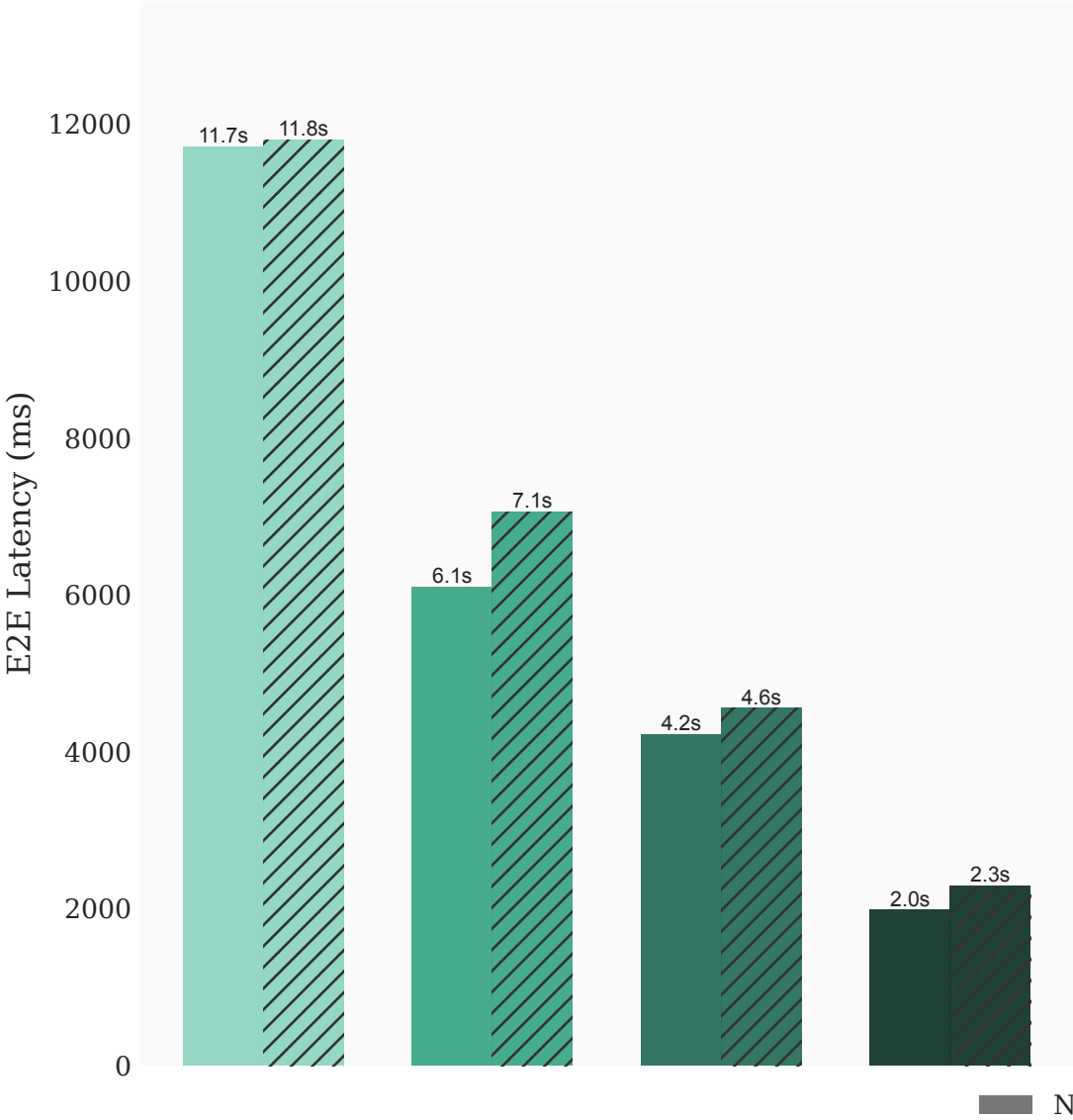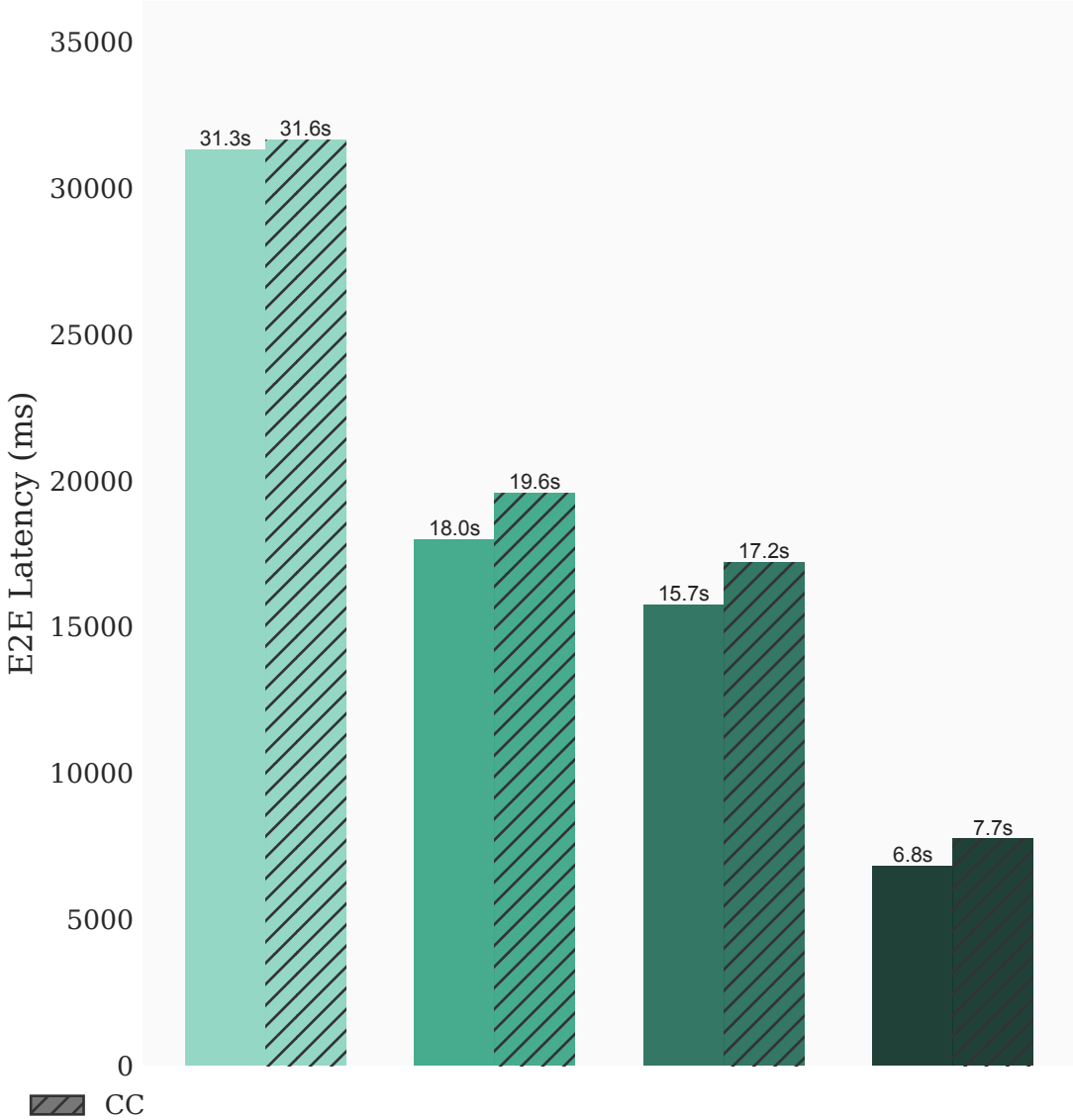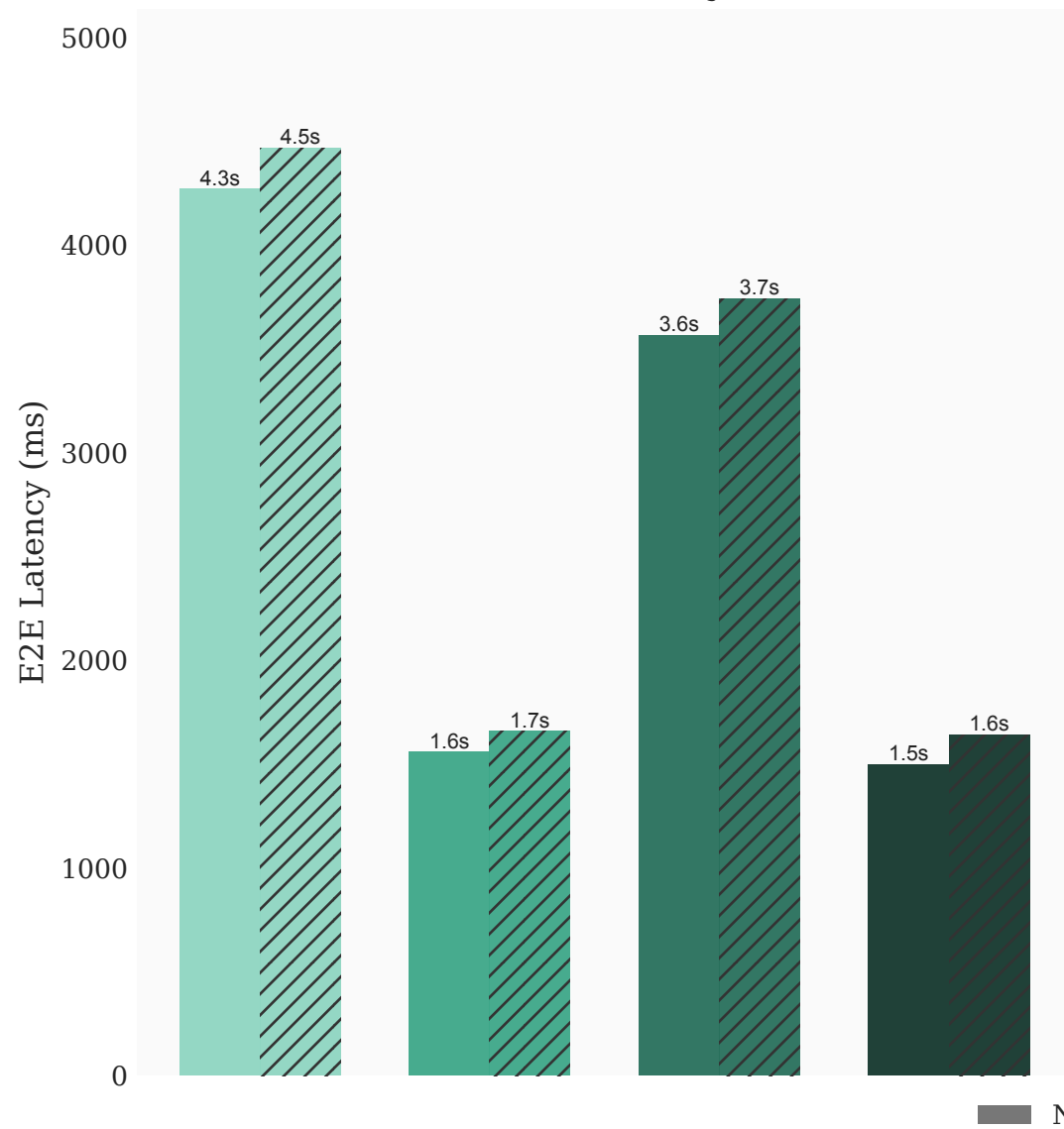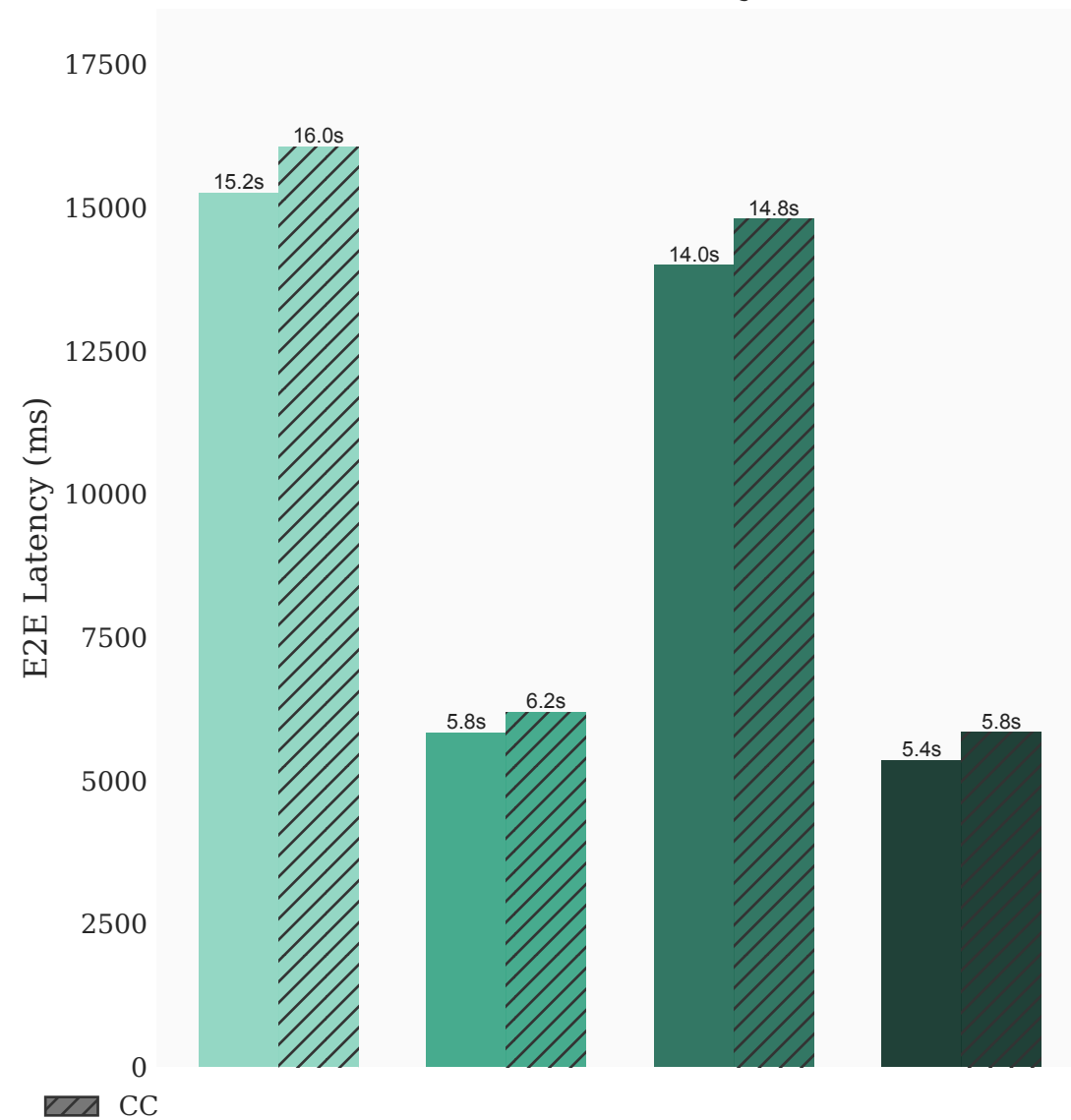# ShareGPT (Single Request)

## End-to-End Latency (Mean)



E2E Latency (ms)

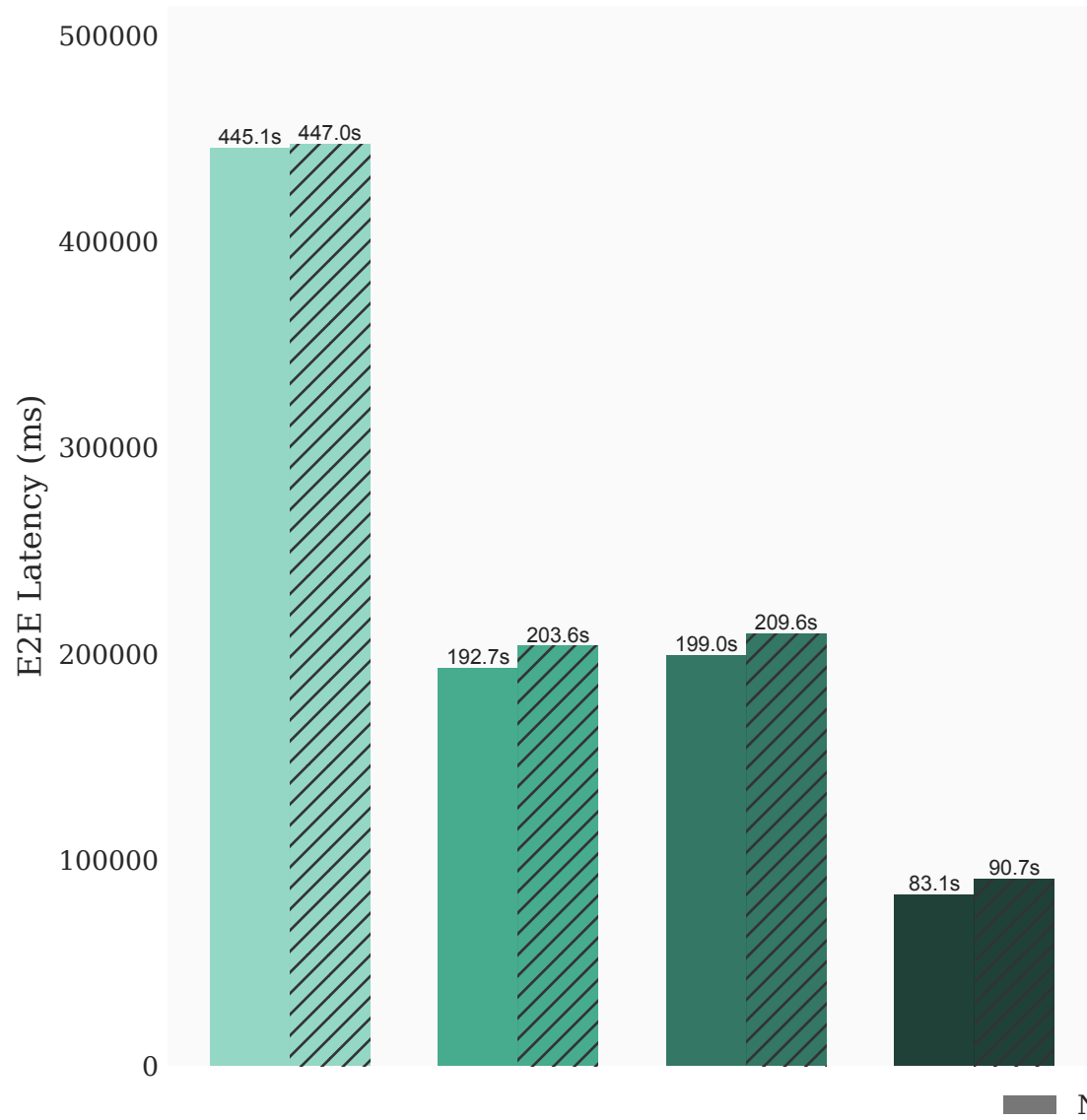- 4.3s
- 4.5s
- 1.6s
- 1.7s
- 3.6s
- 3.7s
- 1.5s
- 1.6s

## End-to-End Latency (P99)

E2E Latency (ms)

- 15.2s
- 16.0s
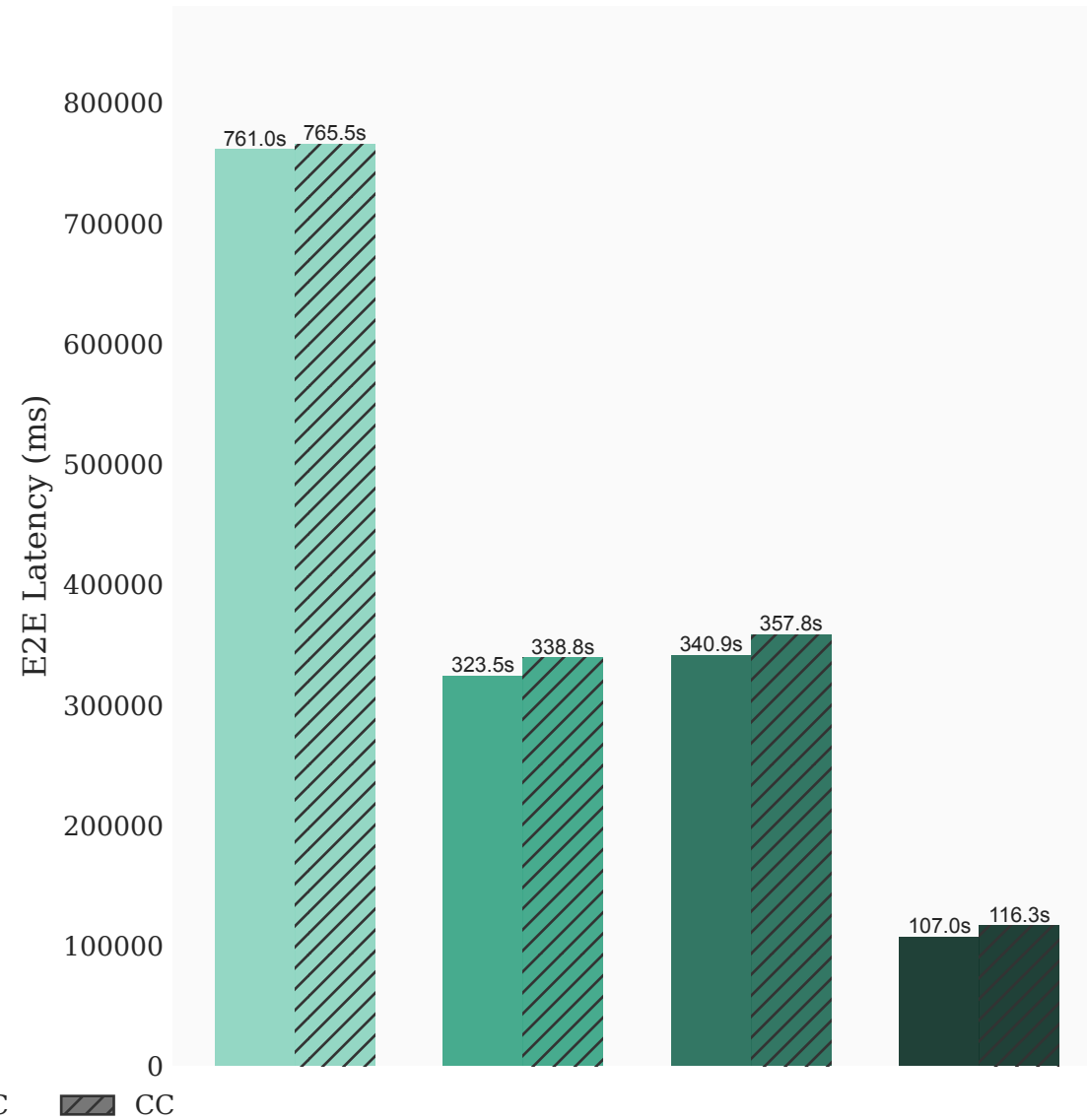- 5.8s
- 6.2s
- 14.0s
- 14.8s
- 5.4s
- 5.8s

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)

## End-to-End Latency (Mean)



## End-to-End Latency (P99)



Legend: No CC, CC

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

Mean values: 445.1s, 447.0s, 192.7s, 203.6s, 199.0s, 209.6s, 83.1s, 90.7s

P99 values: 761.0s, 765.5s, 323.5s, 338.8s, 340.9s, 357.8s, 107.0s, 116.3s

# Edit 10K Characters (50 Request Rate)

## End-to-End Latency (Mean)

E2E Latency (ms)

- LLama 3.3 70B Int4: No CC 439.4s, CC 441.5s
- GPT OSS 120B: No CC 191.6s, CC 203.3s
- Mistral 3.1 24B: No CC 197.1s, CC 207.8s
- LLama 3.1 8B: No CC 78.2s, CC 85.7s

## End-to-End Latency (P99)

E2E Latency (ms)

- LLama 3.3 70B Int4: No CC 753.6s, CC 765.8s
- GPT OSS 120B: No CC 320.7s, CC 338.0s
- Mistral 3.1 24B: No CC 336.1s, CC 355.5s
- LLama 3.1 8B: No CC 101.3s, CC 112.0s

Legend: No CC, CC

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

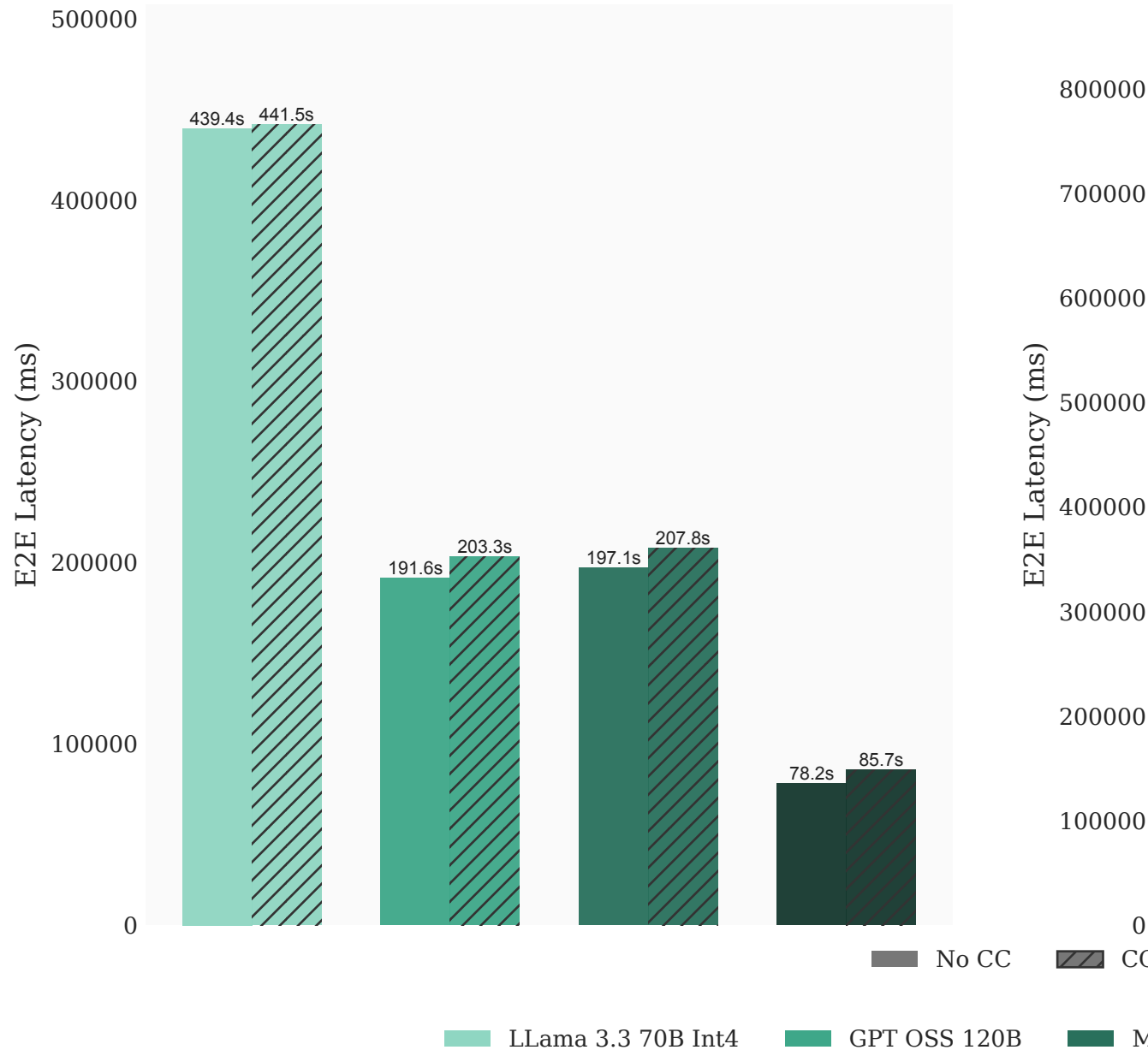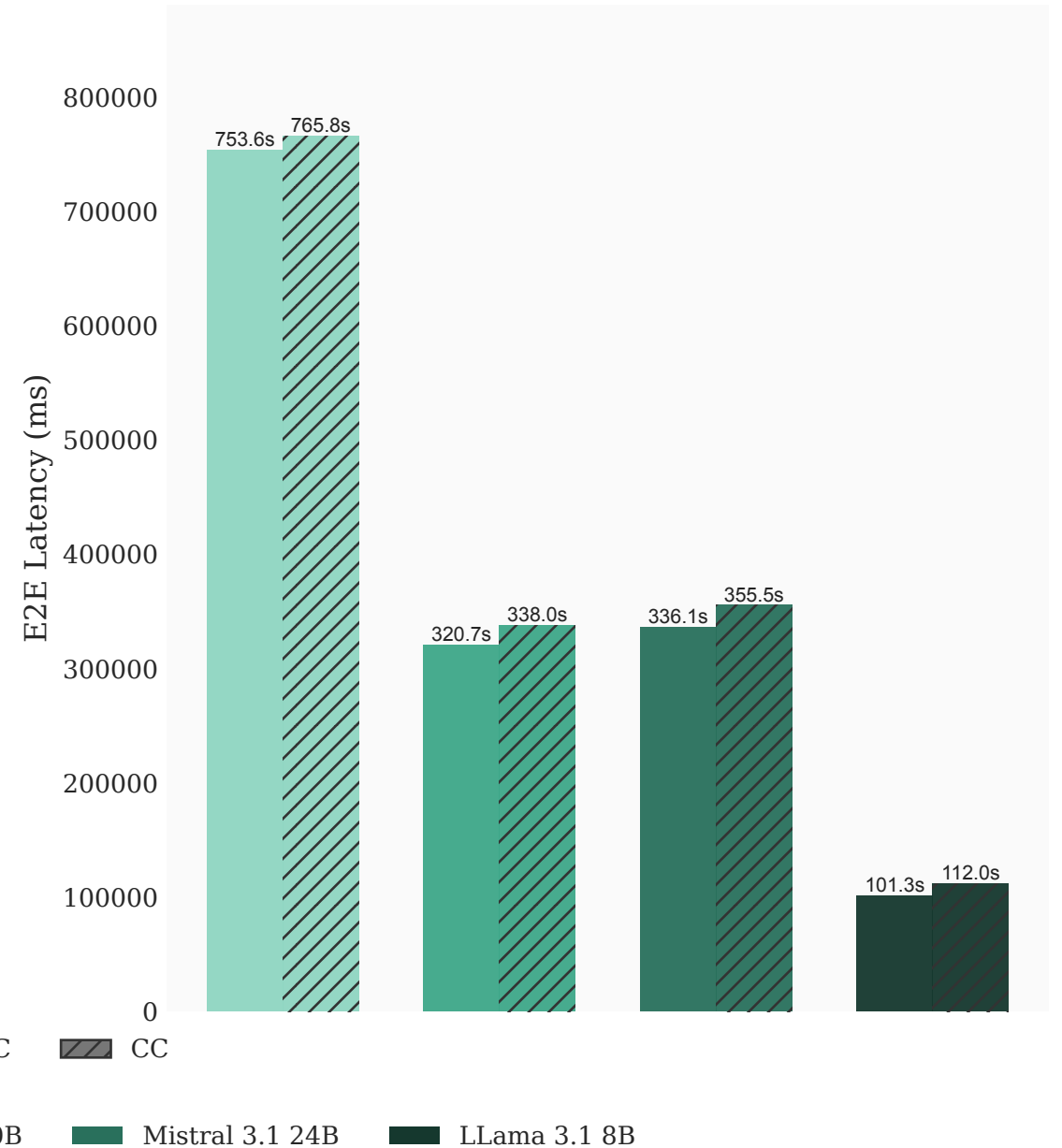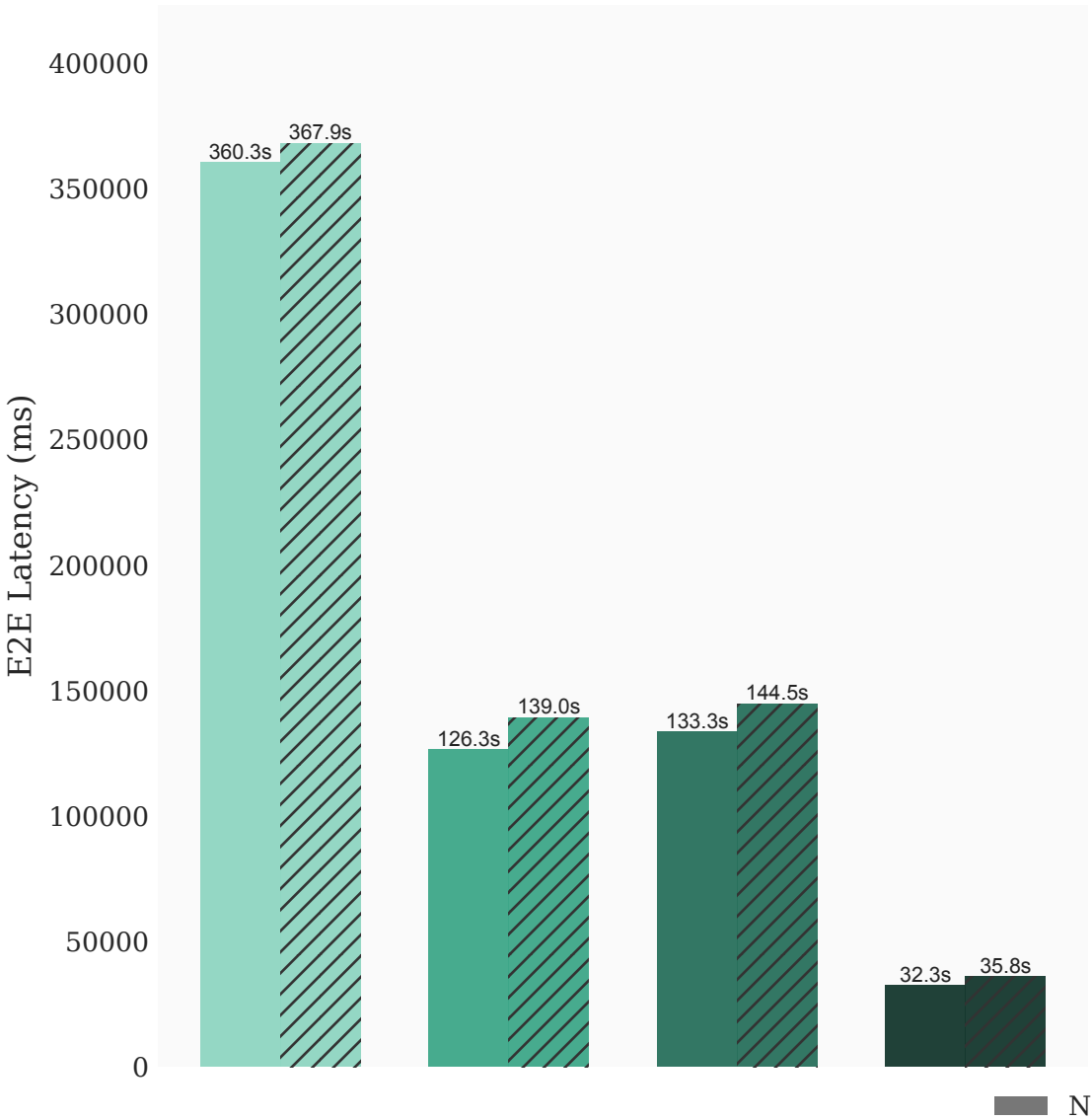# Edit 10K Characters (Single Request)

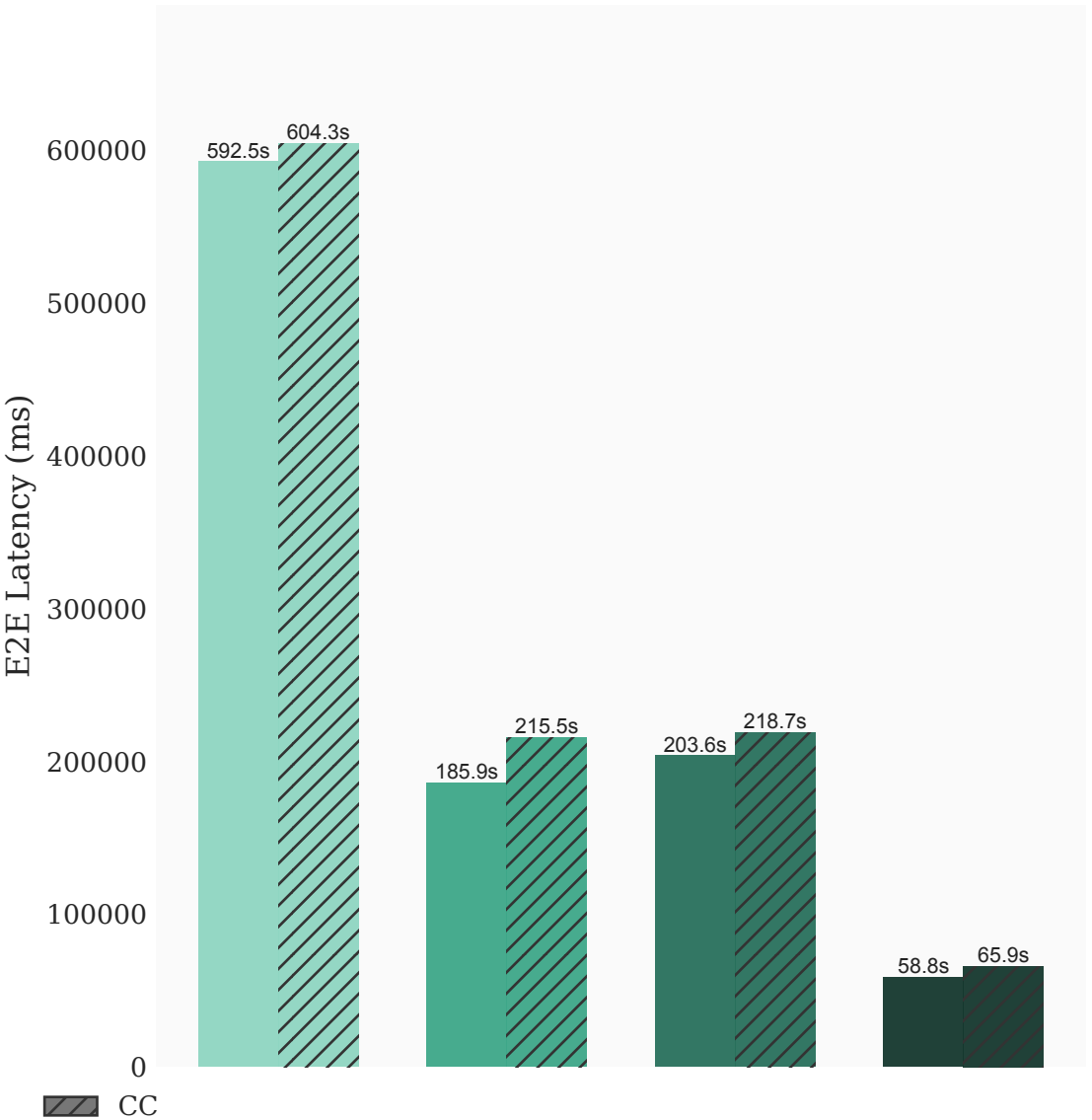## End-to-End Latency (Mean)



## End-to-End Latency (P99)



Legend: No CC | CC

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

Mean values:
- LLama 3.3 70B Int4: 360.3s (No CC), 367.9s (CC)
- GPT OSS 120B: 126.3s (No CC), 139.0s (CC)
- Mistral 3.1 24B: 133.3s (No CC), 144.5s (CC)
- LLama 3.1 8B: 32.3s (No CC), 35.8s (CC)

P99 values:
- LLama 3.3 70B Int4: 592.5s (No CC), 604.3s (CC)
- GPT OSS 120B: 185.9s (No CC), 215.5s (CC)
- Mistral 3.1 24B: 203.6s (No CC), 218.7s (CC)
- LLama 3.1 8B: 58.8s (No CC), 65.9s (CC)

# Numina Math (100 Request Rate)

## End-to-End Latency (Mean)

E2E Latency (ms)

- LLama 3.3 70B Int4: 34.1s (No CC), 34.3s (CC)
- GPT OSS 120B: 13.2s (No CC), 14.1s (CC)
- Mistral 3.1 24B: 7.8s (No CC), 8.4s (CC)
- LLama 3.1 8B: 4.4s (No CC), 5.2s (CC)

## End-to-End Latency (P99)

E2E Latency (ms)

- LLama 3.3 70B Int4: 52.8s (No CC), 54.2s (CC)
- GPT OSS 120B: 24.3s (No CC), 25.5s (CC)
- Mistral 3.1 24B: 19.9s (No CC), 21.1s (CC)
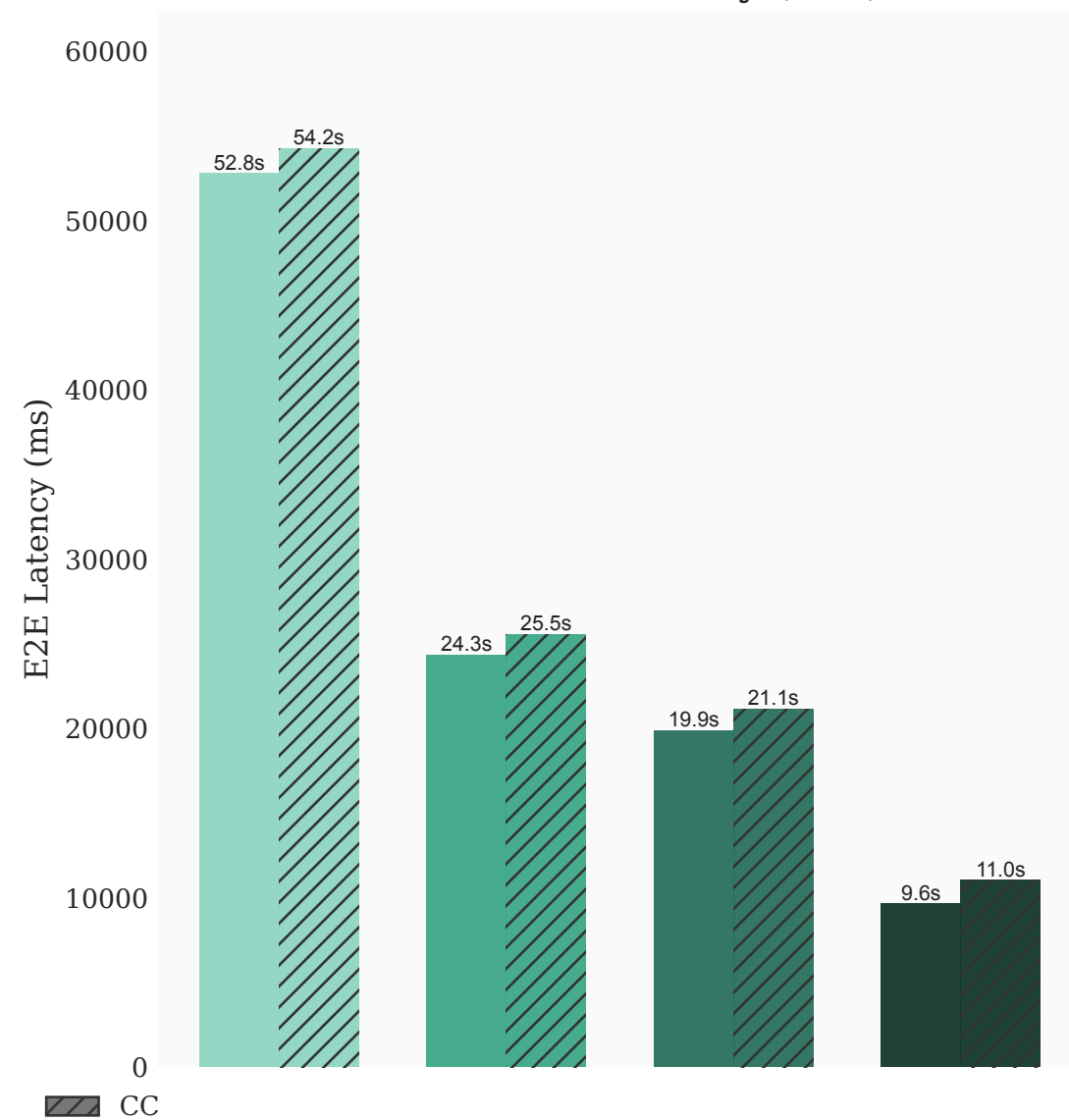- LLama 3.1 8B: 9.6s (No CC), 11.0s (CC)
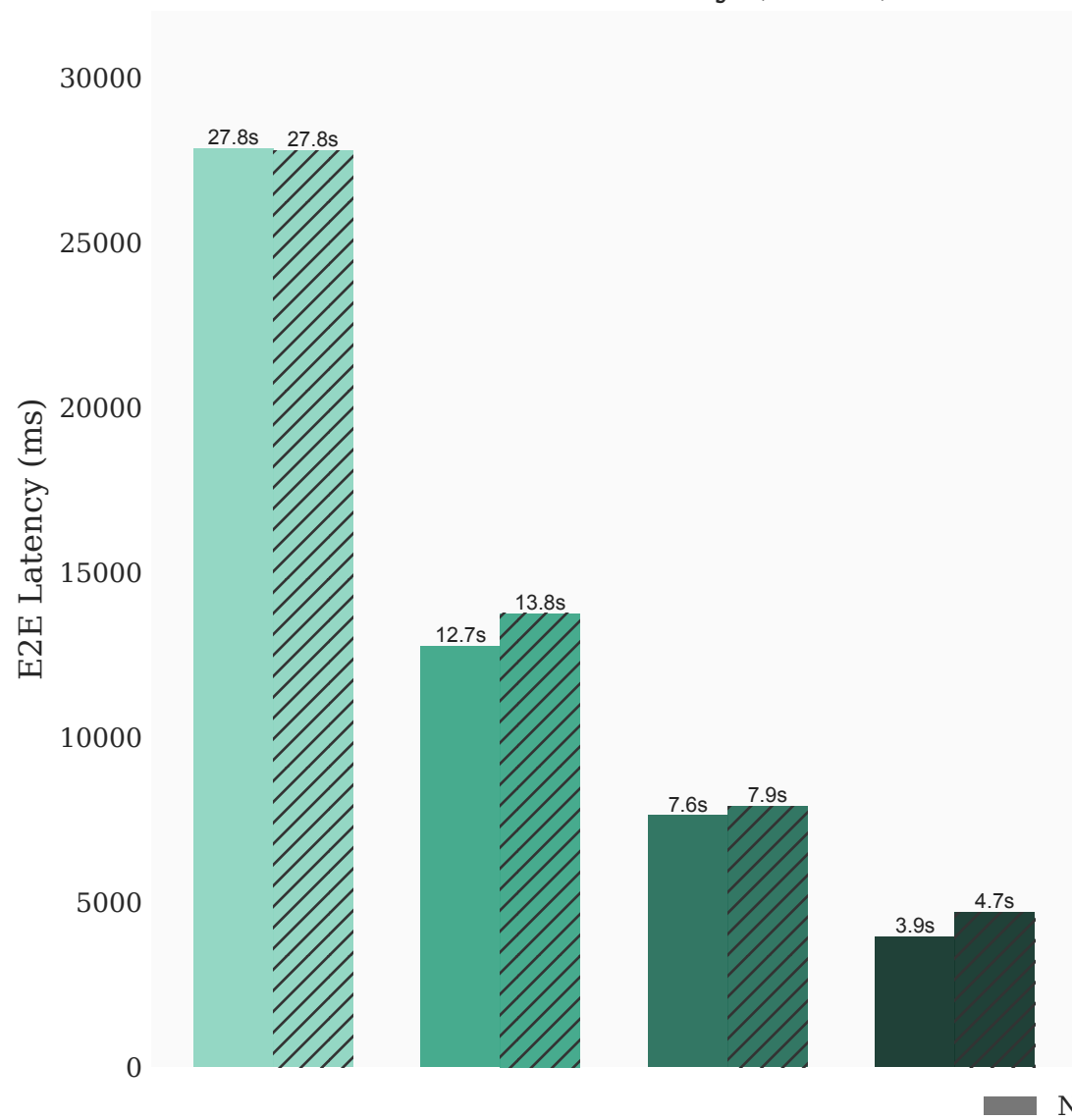
■ No CC  ▨ CC

■ LLama 3.3 70B Int4  ■ GPT OSS 120B  ■ Mistral 3.1 24B  ■ LLama 3.1 8B
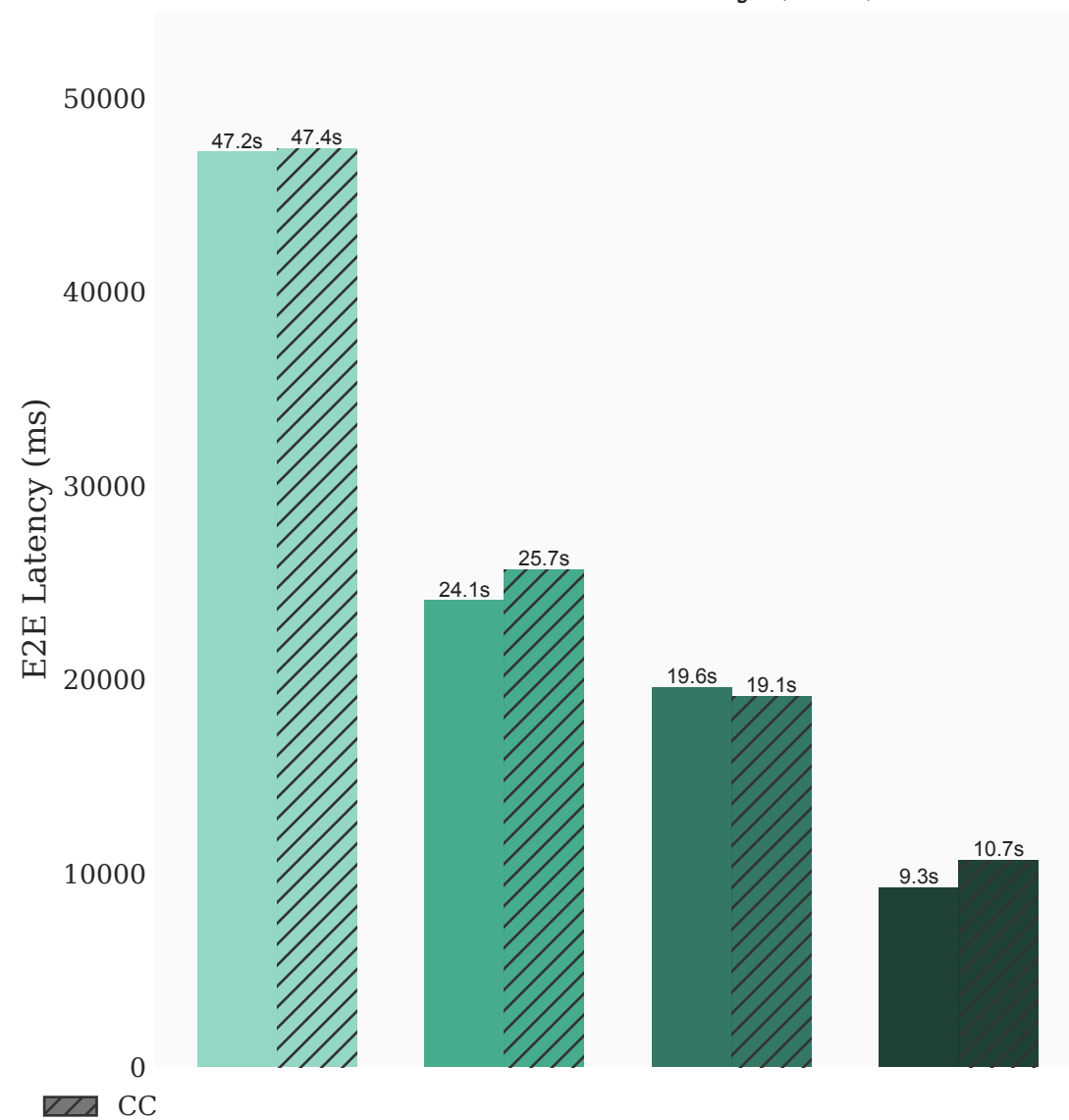
# Numina Math (50 Request Rate)

## End-to-End Latency (Mean)
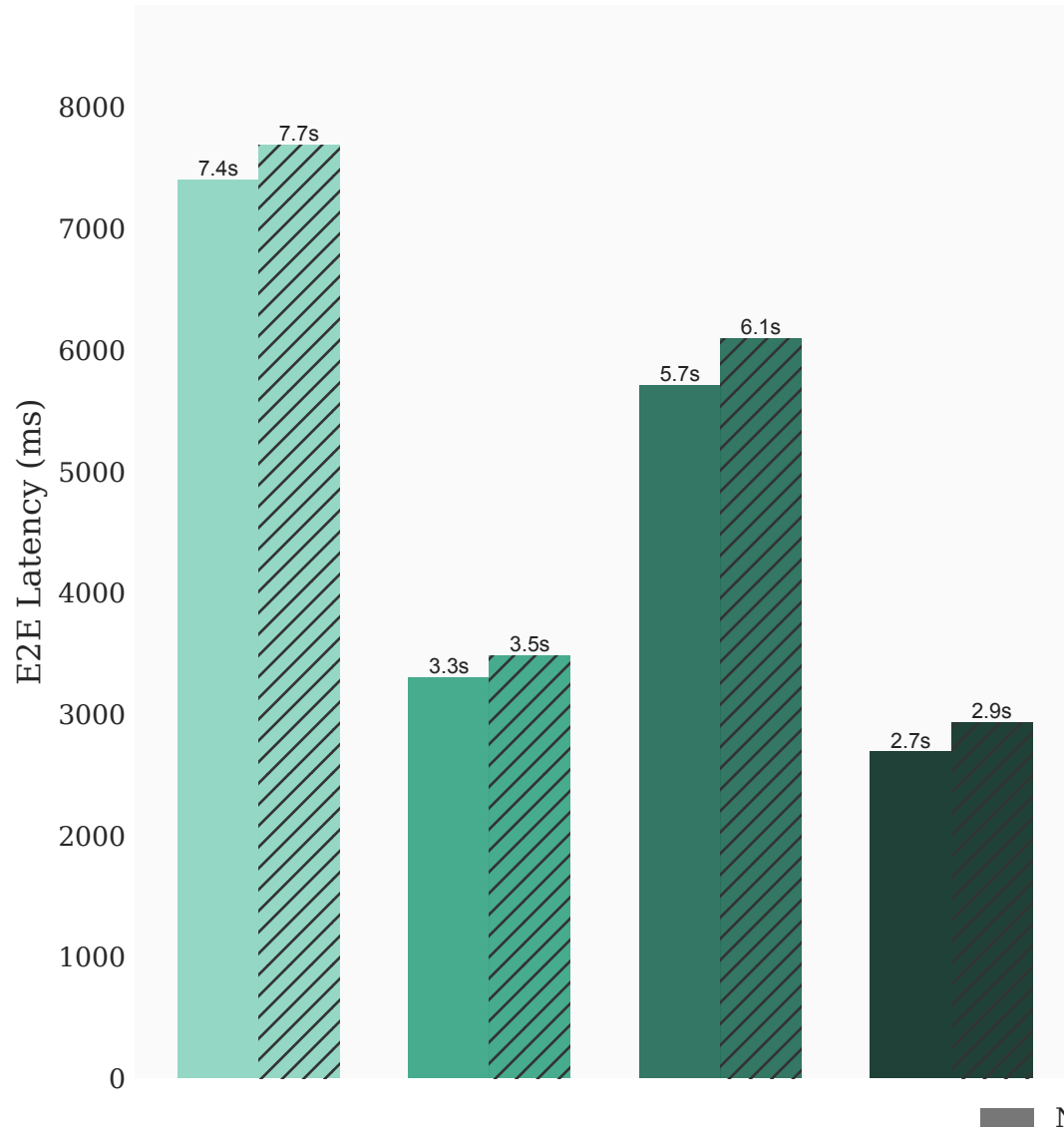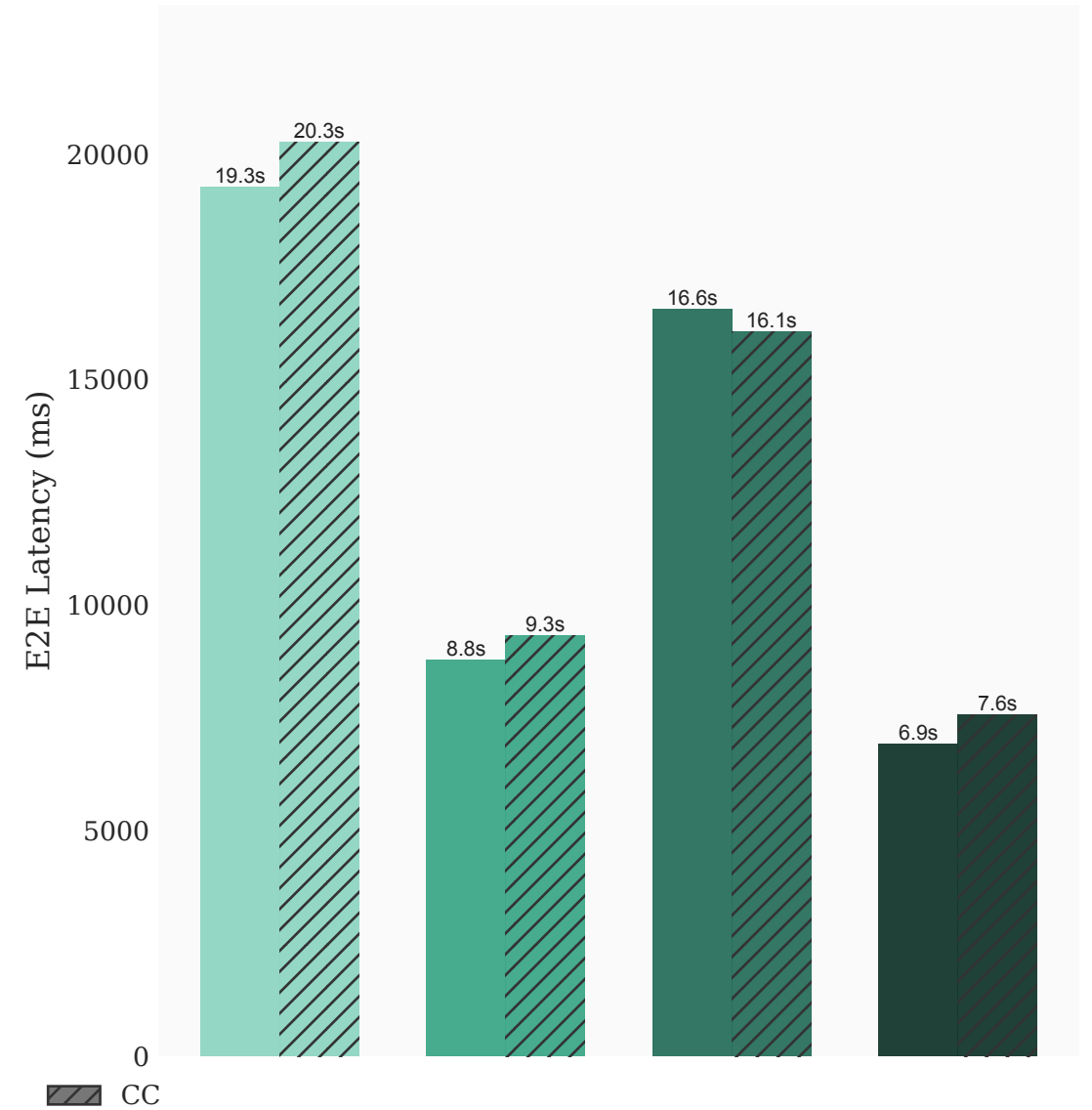


## End-to-End Latency (P99)

No CC     CC

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

# Numina Math (Single Request)

## End-to-End Latency (Mean)



E2E Latency (ms)

- LLama 3.3 70B Int4: No CC 7.4s, CC 7.7s
- GPT OSS 120B: No CC 3.3s, CC 3.5s
- Mistral 3.1 24B: No CC 5.7s, CC 6.1s
- LLama 3.1 8B: No CC 2.7s, CC 2.9s

## End-to-End Latency (P99)

E2E Latency (ms)

- LLama 3.3 70B Int4: No CC 19.3s, CC 20.3s
- GPT OSS 120B: No CC 8.8s, CC 9.3s
- Mistral 3.1 24B: No CC 16.6s, CC 16.1s
- LLama 3.1 8B: No CC 6.9s, CC 7.6s

No CC ▨ CC

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