# Random (1500 ⇒ 250) (Rate 100)

## Time to First Token (Mean)

TTFT (ms)

- 2.8s / 2.8s
- 13.5s / 13.1s
- 5.7s / 5.5s
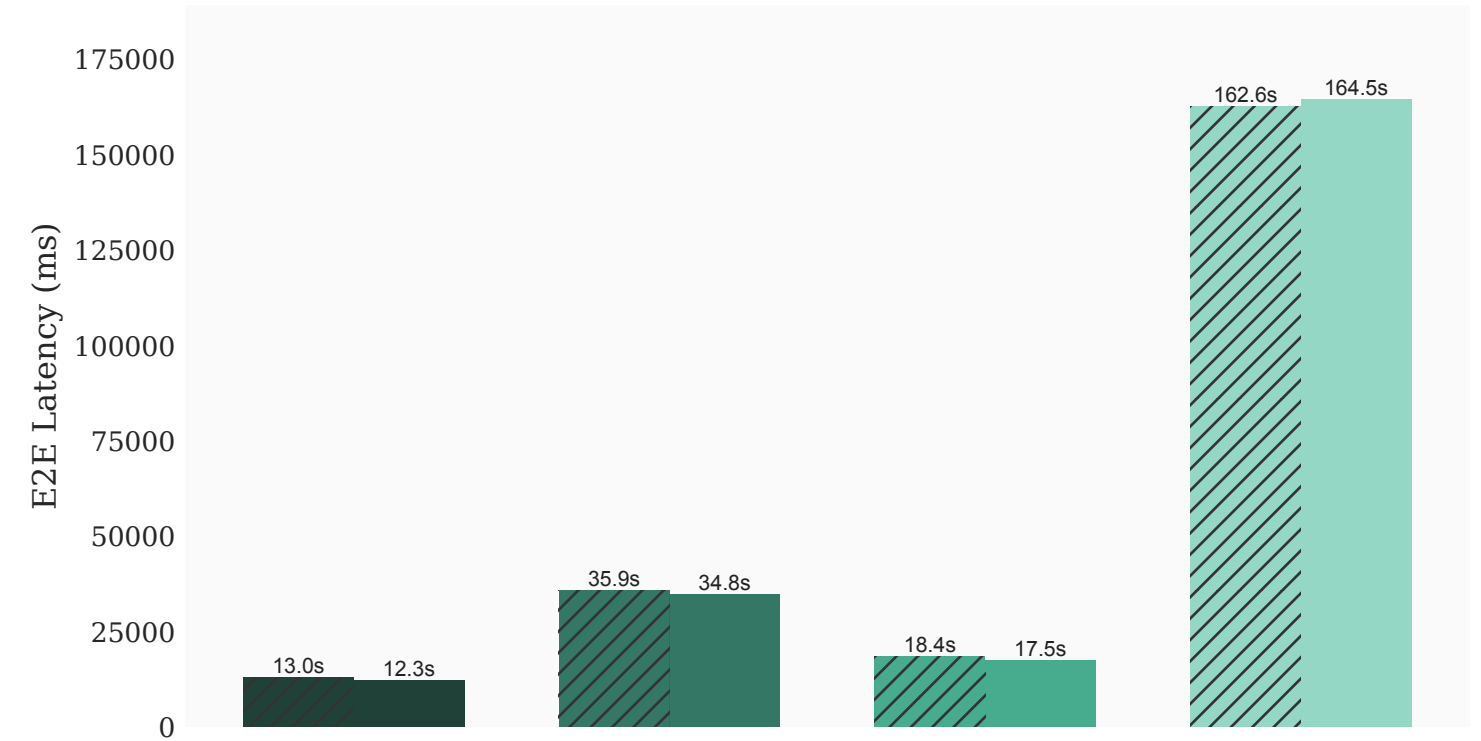- 75.3s / 75.9s

## End-to-End Latency (Mean)

E2E Latency (ms)

- 11.7s / 11.0s
- 29.6s / 28.8s
- 12.4s / 11.7s
- 118.3s / 120.1s

## Time to First Token (P99)

TTFT (ms)

- 5.5s / 5.5s
- 28.3s / 27.7s
- 11.8s / 10.9s
- 152.5s / 155.6s

## End-to-End Latency (P99)

E2E Latency (ms)

- 13.0s / 12.3s
- 35.9s / 34.8s
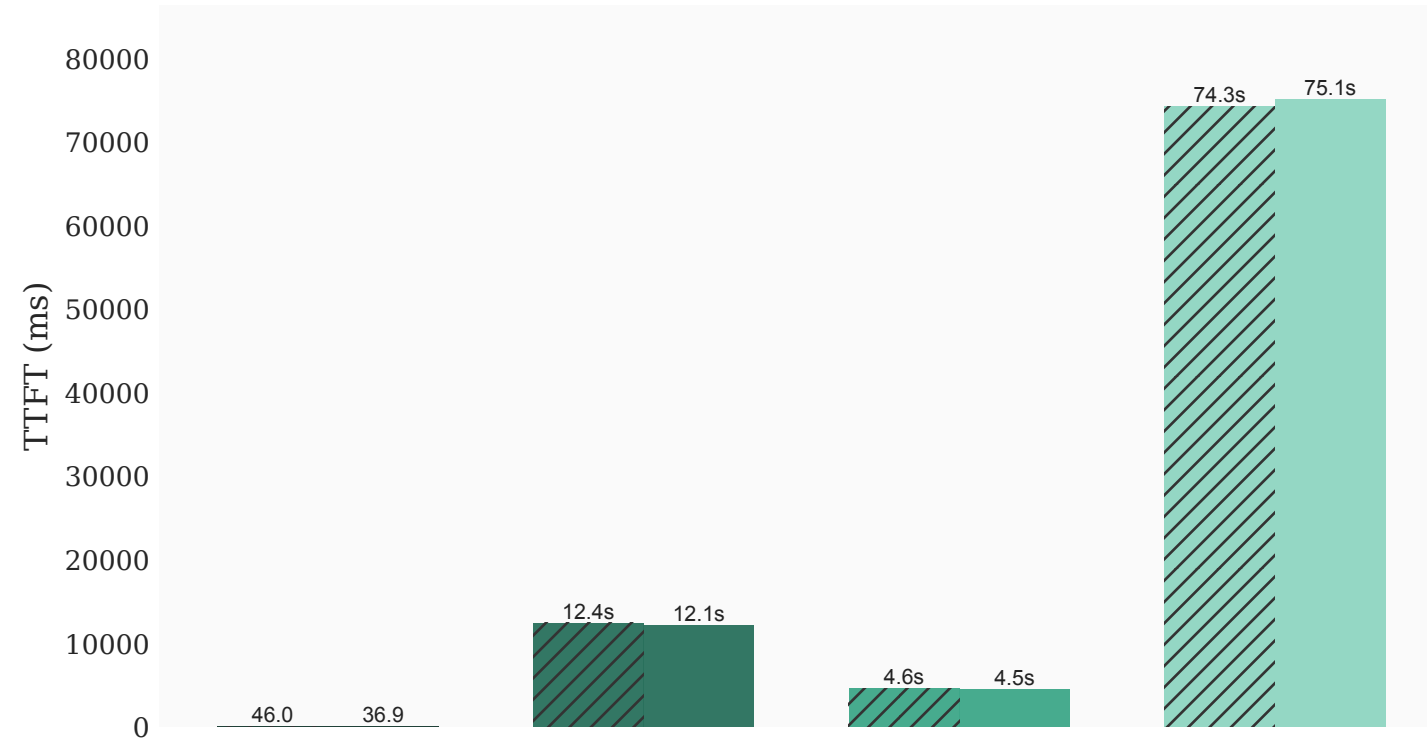- 18.4s / 17.5s
- 162.6s / 164.5s

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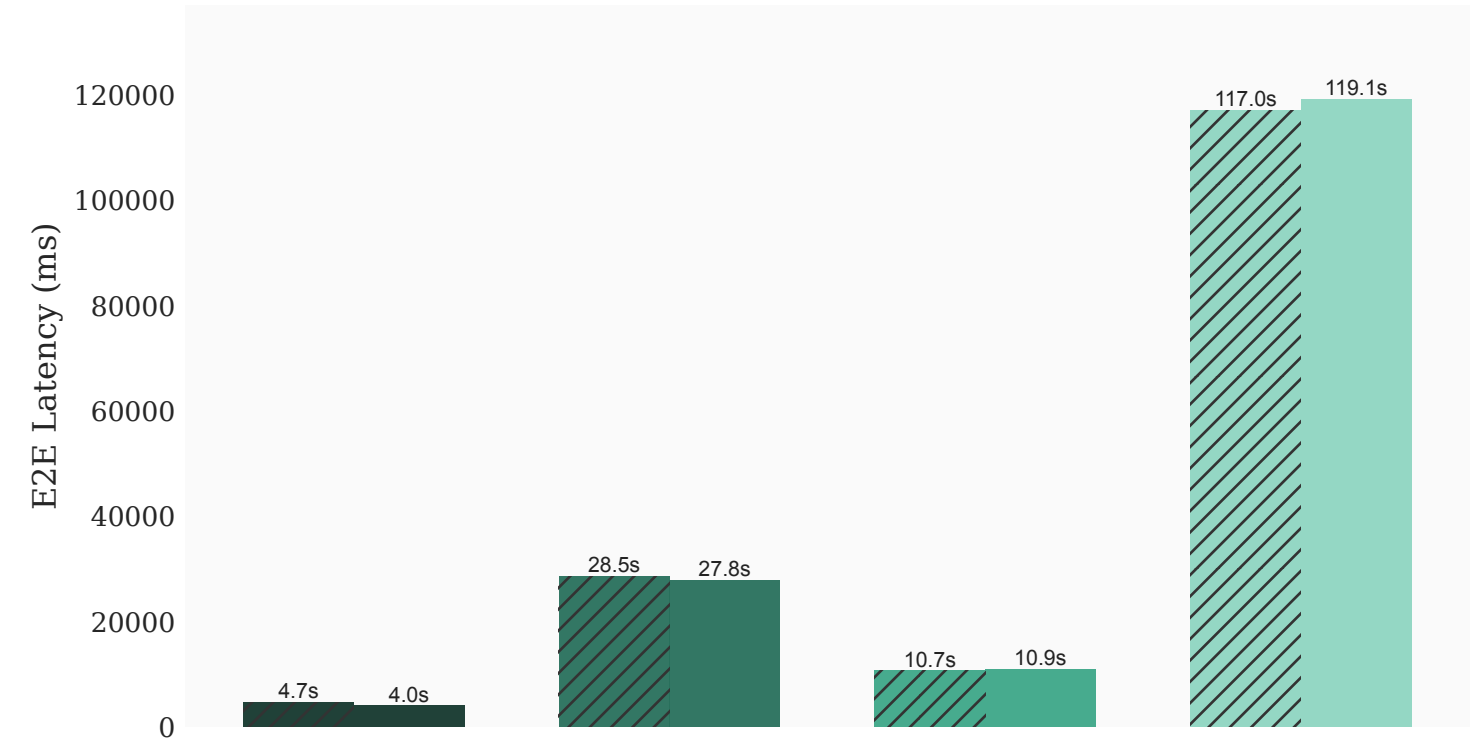
LLama 3.1 8B    Mistral 3.1 24B    GPT OSS 120B    LLama 3.3 70B Int4

# Random (1500 ⇒ 250) (Rate 50)

## Time to First Token (Mean)

| | CC | No CC |
|---|---|---|
| LLama 3.1 8B | 46.0 | 36.9 |
| Mistral 3.1 24B | 12.4s | 12.1s |
| GPT OSS 120B | 4.6s | 4.5s |
| LLama 3.3 70B Int4 | 74.3s | 75.1s |

## End-to-End Latency (Mean)

| | CC | No CC |
|---|---|---|
| LLama 3.1 8B | 4.7s | 4.0s |
| Mistral 3.1 24B | 28.5s | 27.8s |
| GPT OSS 120B | 10.7s | 10.9s |
| LLama 3.3 70B Int4 | 117.0s | 119.1s |

## Time to First Token (P99)

| | CC | No CC |
|---|---|---|
| LLama 3.1 8B | 76.8 | 60.5 |
| Mistral 3.1 24B | 26.2s | 25.7s |
| GPT OSS 120B | 9.4s | 9.0s |
| LLama 3.3 70B Int4 | 151.8s | 152.1s |

## End-to-End Latency (P99)

| | CC | No CC |
|---|---|---|
| LLama 3.1 8B | 5.6s | 4.8s |
| Mistral 3.1 24B | 34.6s | 33.6s |
| GPT OSS 120B | 17.7s | 17.6s |
| LLama 3.3 70B Int4 | 160.2s | 162.2s |

Legend: CC, No CC

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

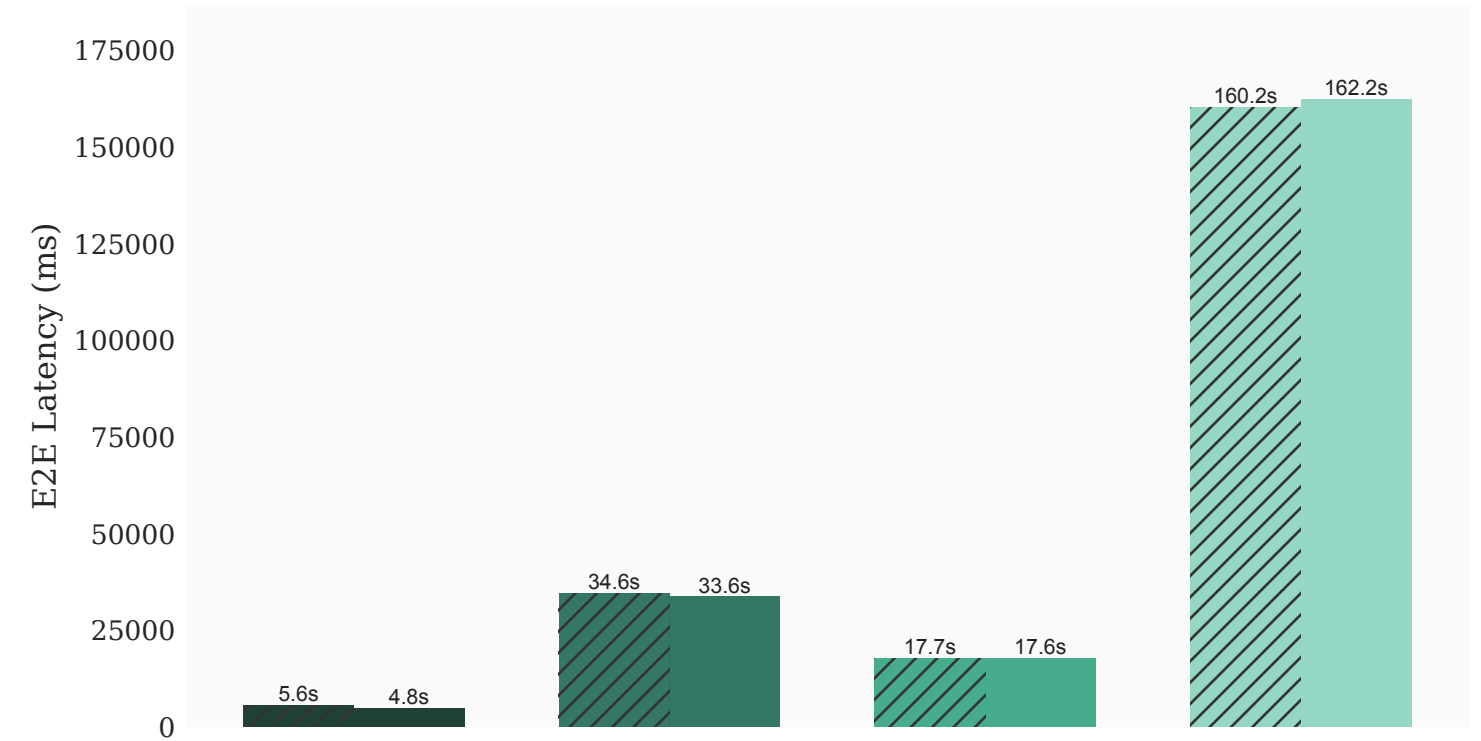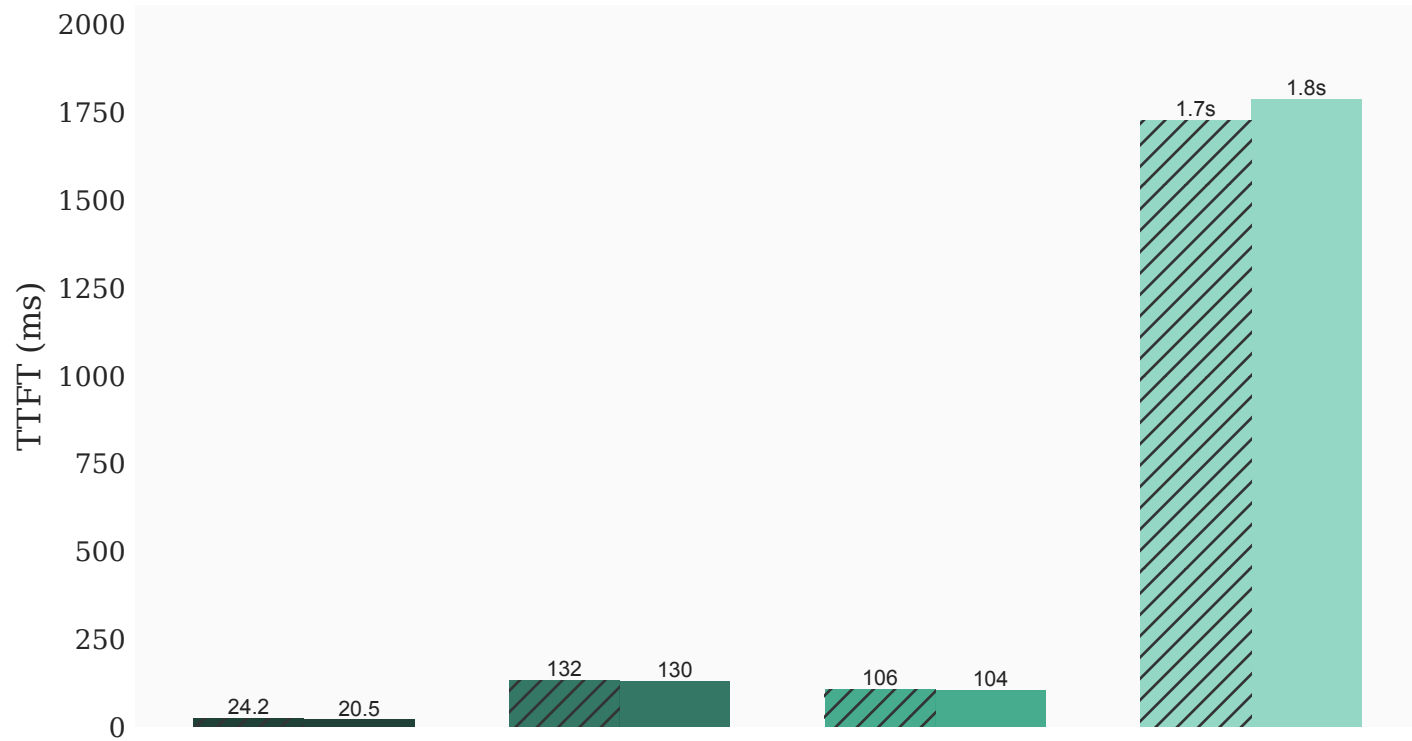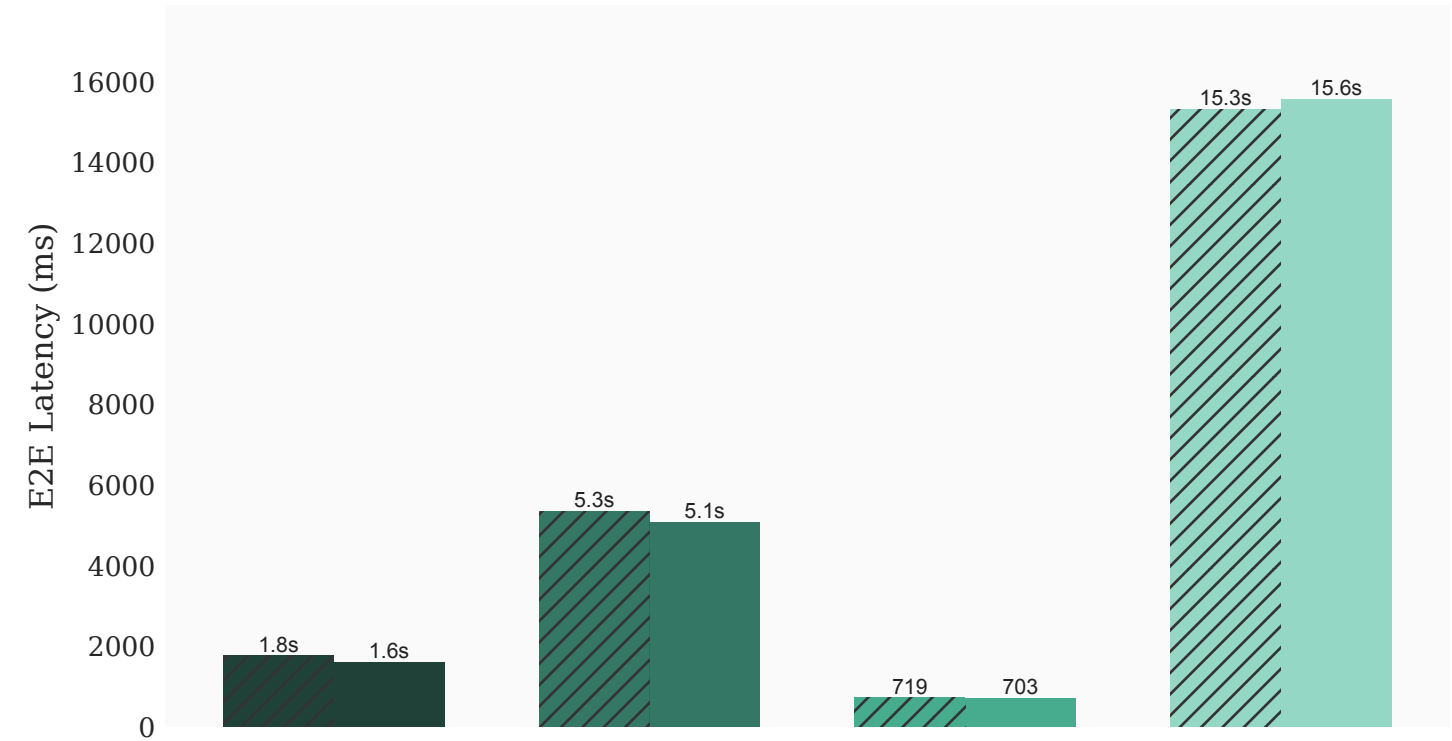# Random (1500 ⇒ 250) (Rate 1)

## Time to First Token (Mean)
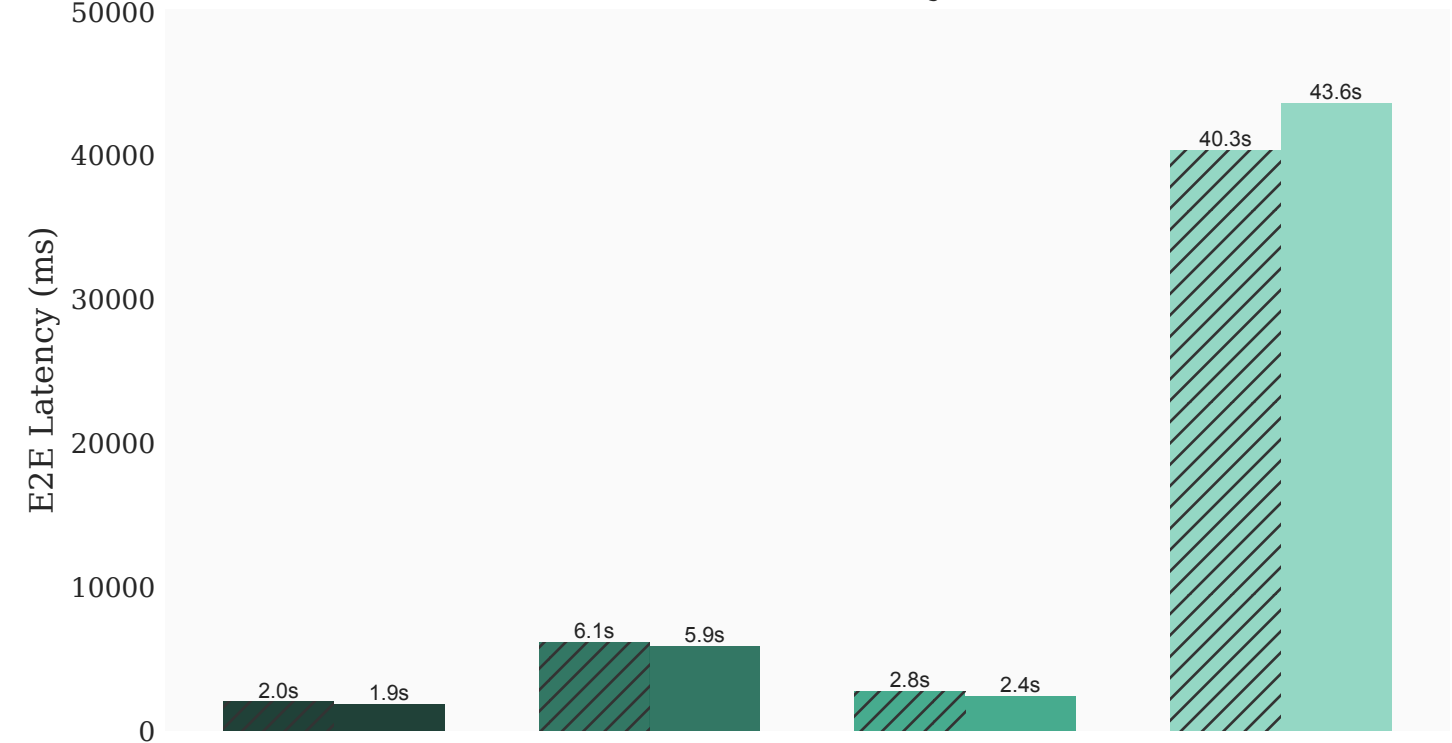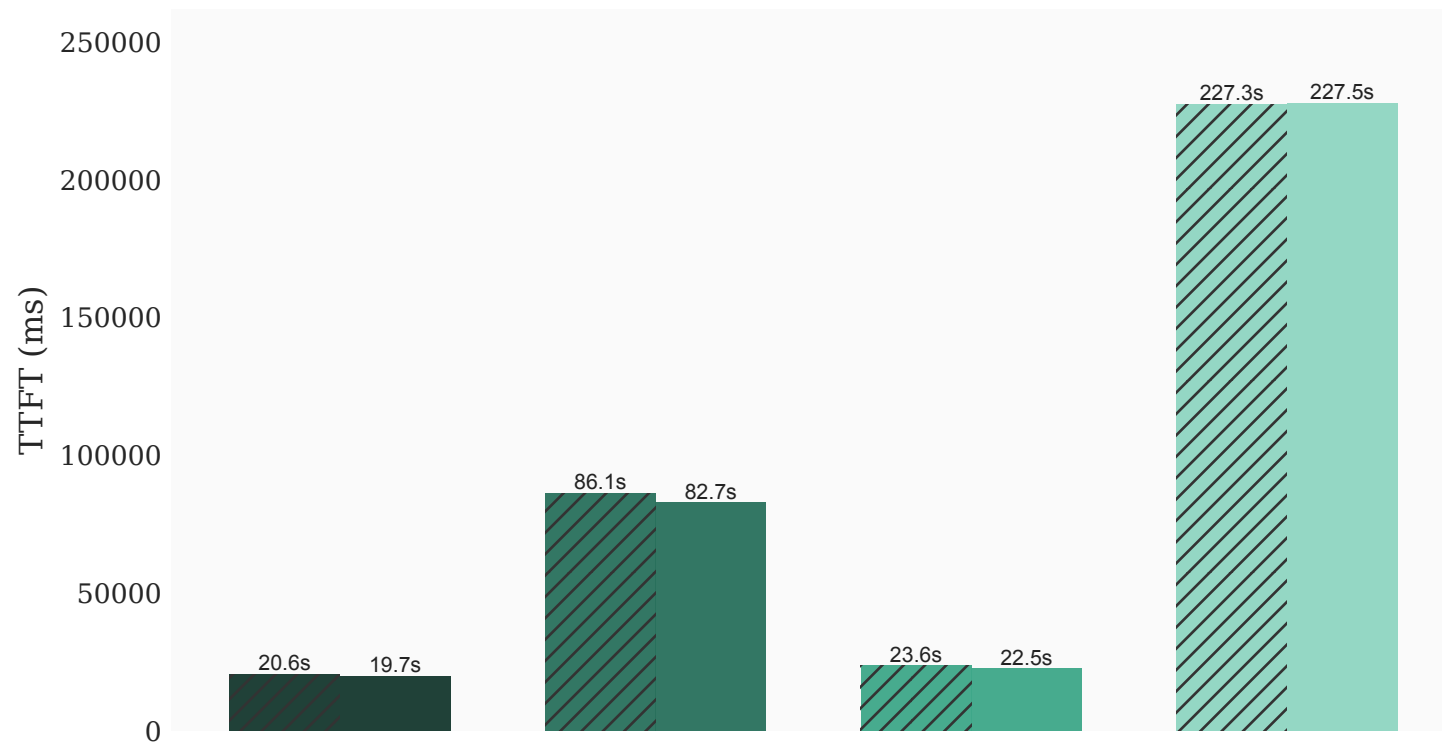
TTFT (ms)

LLama 3.1 8B: 24.2 (CC), 20.5 (No CC)
Mistral 3.1 24B: 132 (CC), 130 (No CC)
GPT OSS 120B: 106 (CC), 104 (No CC)
LLama 3.3 70B Int4: 1.7s (CC), 1.8s (No CC)

## End-to-End Latency (Mean)

E2E Latency (ms)

LLama 3.1 8B: 1.8s (CC), 1.6s (No CC)
Mistral 3.1 24B: 5.3s (CC), 5.1s (No CC)
GPT OSS 120B: 719 (CC), 703 (No CC)
LLama 3.3 70B Int4: 15.3s (CC), 15.6s (No CC)

## Time to First Token (P99)

TTFT (ms)

LLama 3.1 8B: 29.9 (CC), 25.4 (No CC)
Mistral 3.1 24B: 241 (CC), 238 (No CC)
GPT OSS 120B: 177 (CC), 179 (No CC)
LLama 3.3 70B Int4: 5.5s (CC), 5.7s (No CC)

## End-to-End Latency (P99)

E2E Latency (ms)

LLama 3.1 8B: 2.0s (CC), 1.9s (No CC)
Mistral 3.1 24B: 6.1s (CC), 5.9s (No CC)
GPT OSS 120B: 2.8s (CC), 2.4s (No CC)
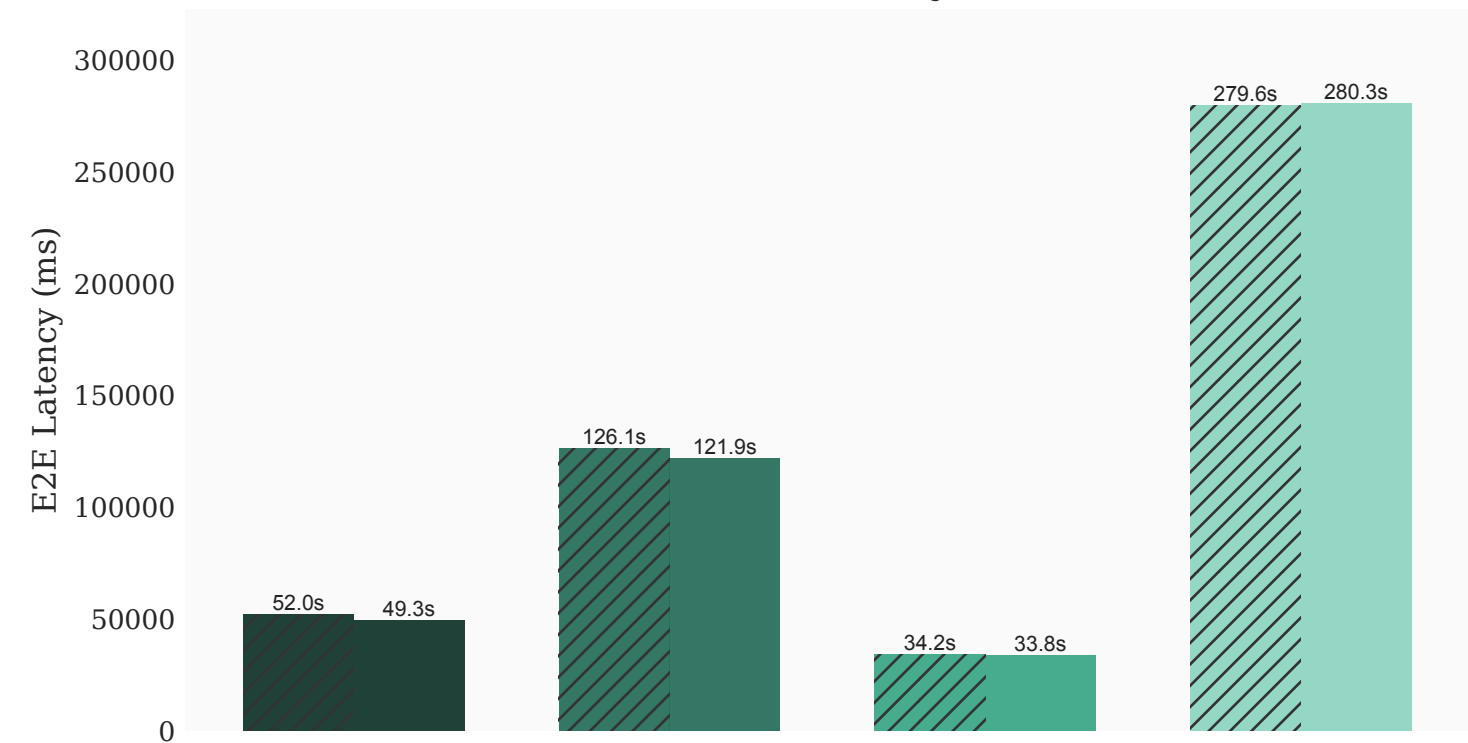LLama 3.3 70B Int4: 40.3s (CC), 43.6s (No CC)

Legend: CC, No CC

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

# Random (4000 ⇒ 1000) (Rate 100)

## Time to First Token (Mean)

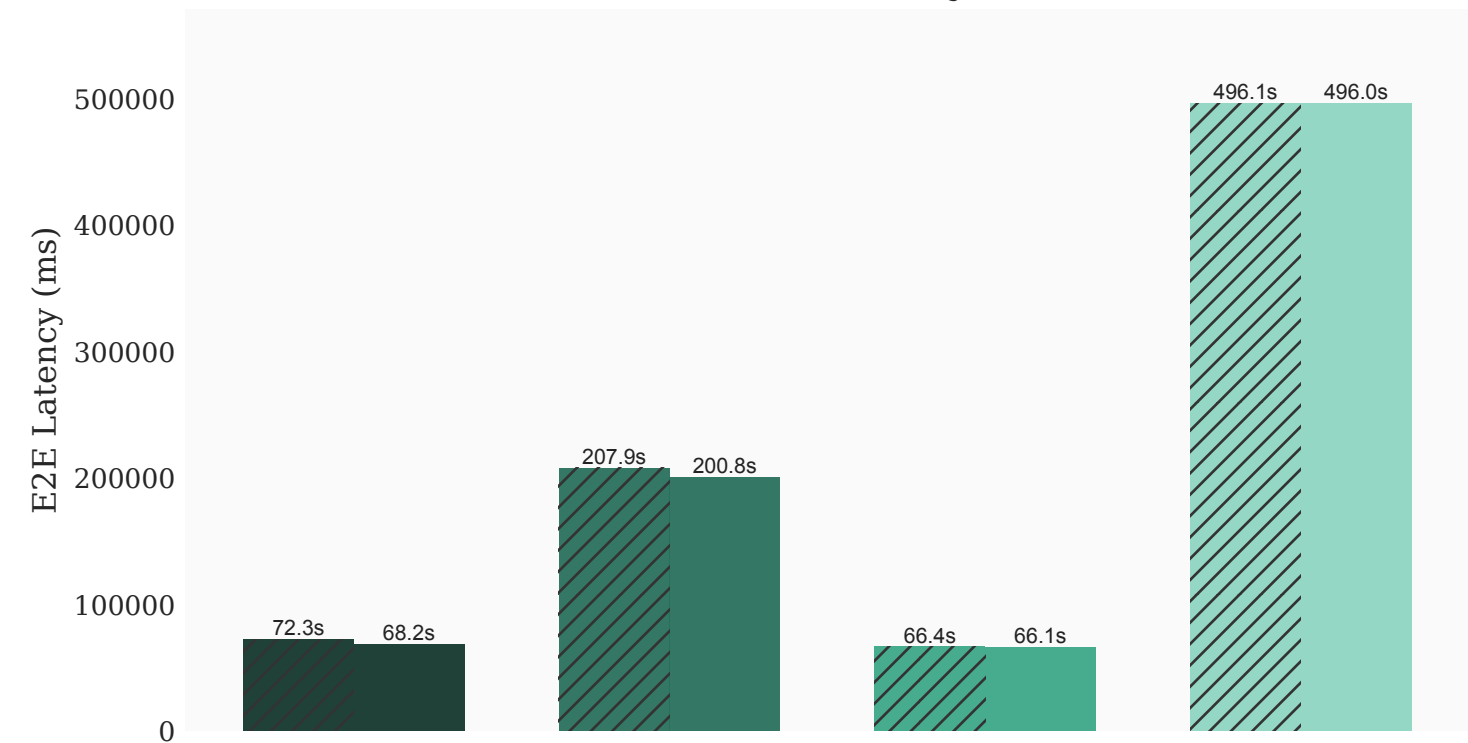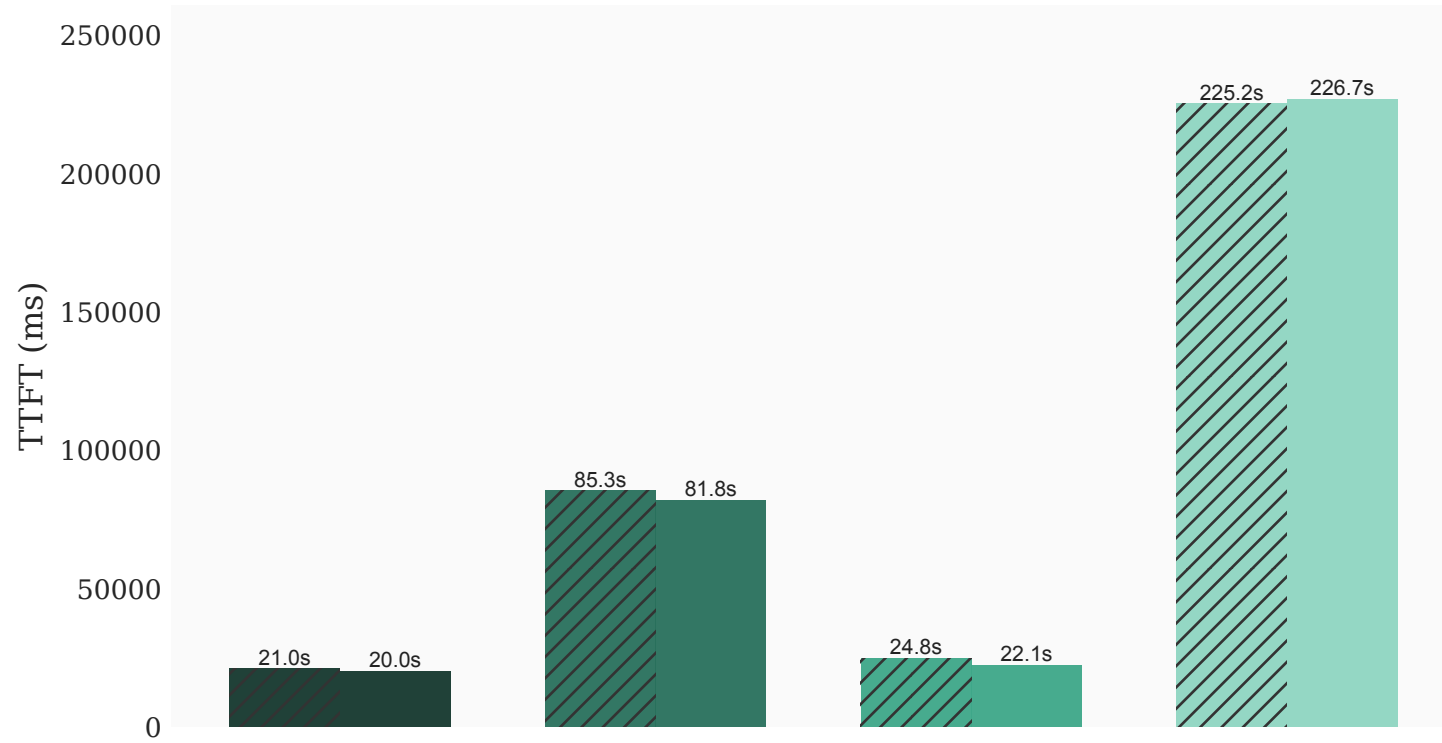TTFT (ms)

- LLama 3.1 8B: CC 20.6s, No CC 19.7s
- Mistral 3.1 24B: CC 86.1s, No CC 82.7s
- GPT OSS 120B: CC 23.6s, No CC 22.5s
- LLama 3.3 70B Int4: CC 227.3s, No CC 227.5s

## End-to-End Latency (Mean)

E2E Latency (ms)

- LLama 3.1 8B: CC 52.0s, No CC 49.3s
- Mistral 3.1 24B: CC 126.1s, No CC 121.9s
- GPT OSS 120B: CC 34.2s, No CC 33.8s
- LLama 3.3 70B Int4: CC 279.6s, No CC 280.3s

## Time to First Token (P99)

TTFT (ms)

- LLama 3.1 8B: CC 47.5s, No CC 45.5s
- Mistral 3.1 24B: CC 179.9s, No CC 175.2s
- GPT OSS 120B: CC 53.3s, No CC 52.1s
- LLama 3.3 70B Int4: CC 470.1s, No CC 469.8s

## End-to-End Latency (P99)

E2E Latency (ms)

- LLama 3.1 8B: CC 72.3s, No CC 68.2s
- Mistral 3.1 24B: CC 207.9s, No CC 200.8s
- GPT OSS 120B: CC 66.4s, No CC 66.1s
- LLama 3.3 70B Int4: CC 496.1s, No CC 496.0s

Legend: CC / No CC

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

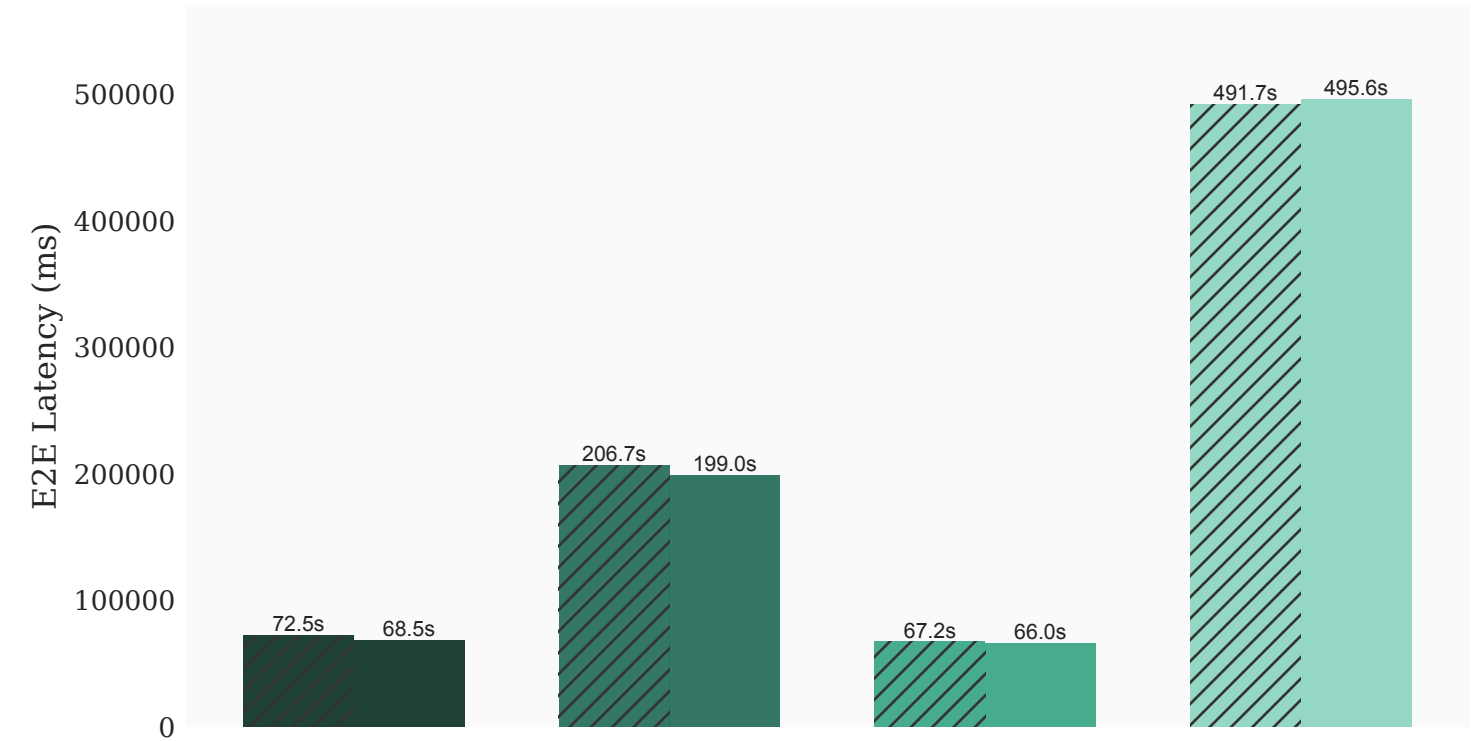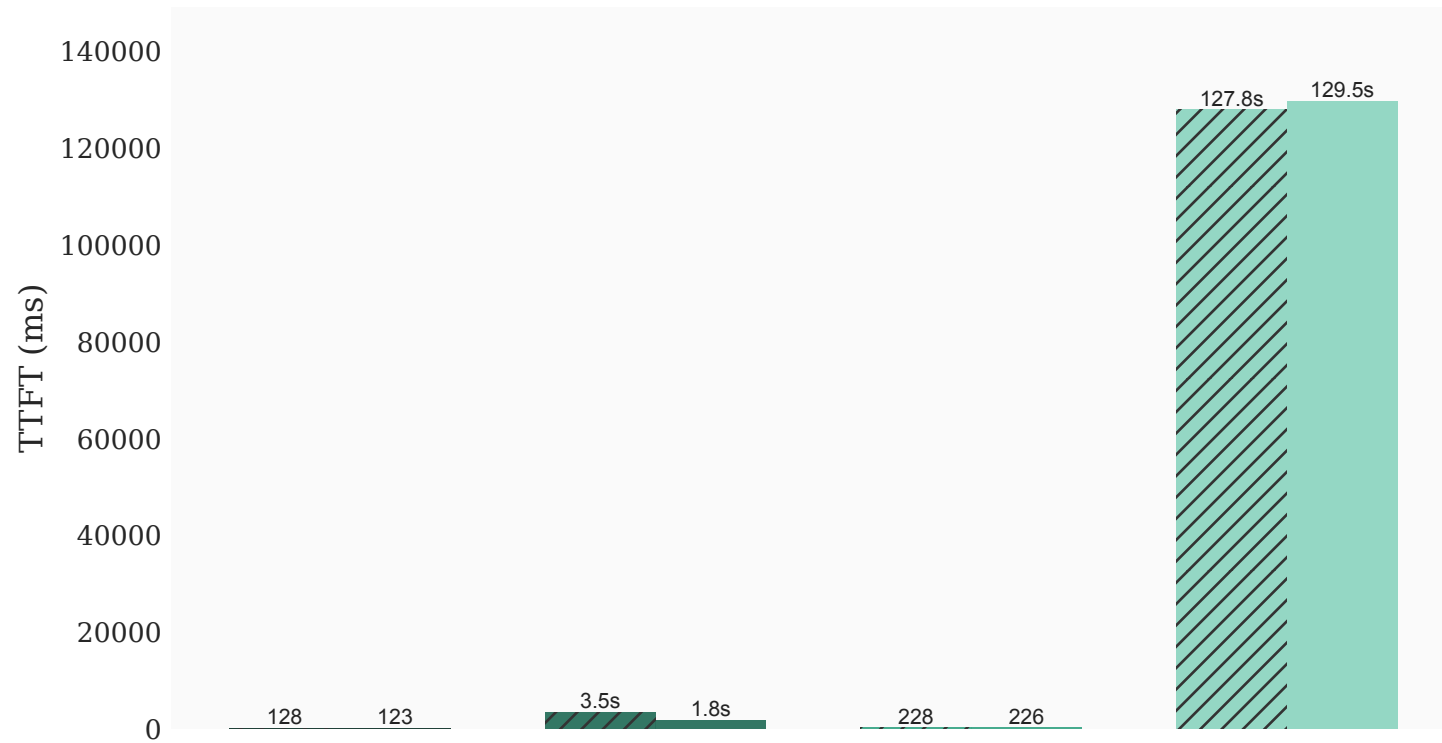# Random (4000 ⇒ 1000) (Rate 50)

## Time to First Token (Mean)
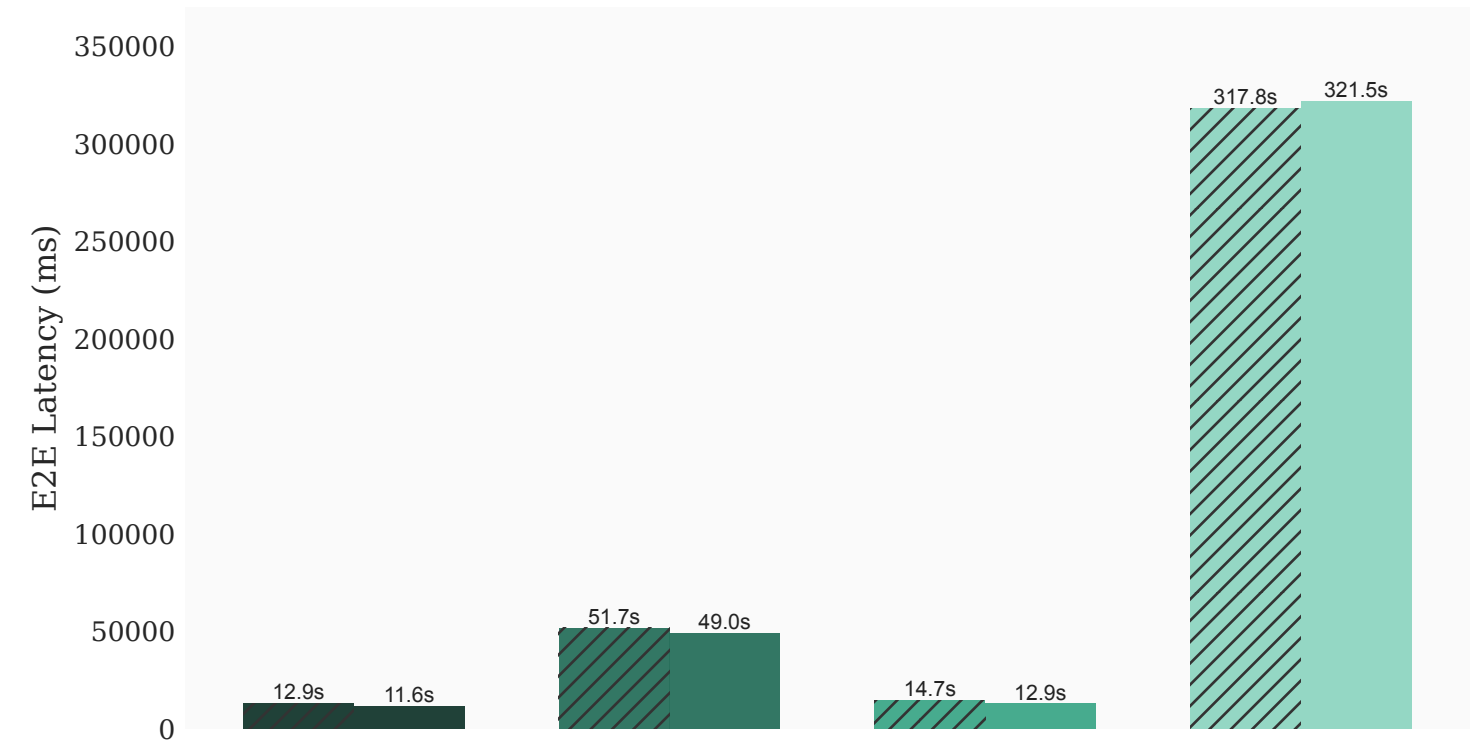
TTFT (ms)

- 21.0s, 20.0s
- 85.3s, 81.8s
- 24.8s, 22.1s
- 225.2s, 226.7s

## End-to-End Latency (Mean)

E2E Latency (ms)

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

## Time to First Token (P99)

TTFT (ms)

- 47.1s, 44.9s
- 178.6s, 173.2s
- 51.4s, 51.4s
- 464.6s, 468.7s

## End-to-End Latency (P99)

E2E Latency (ms)

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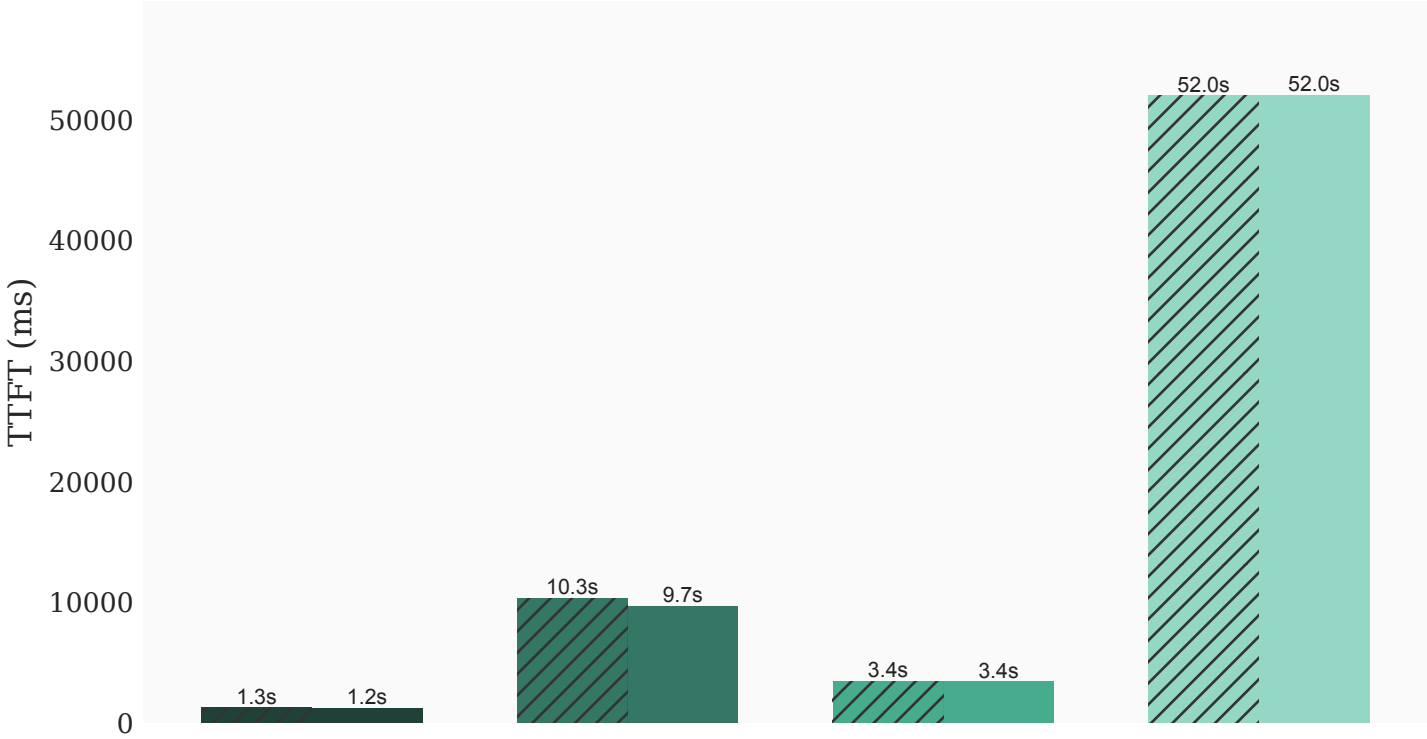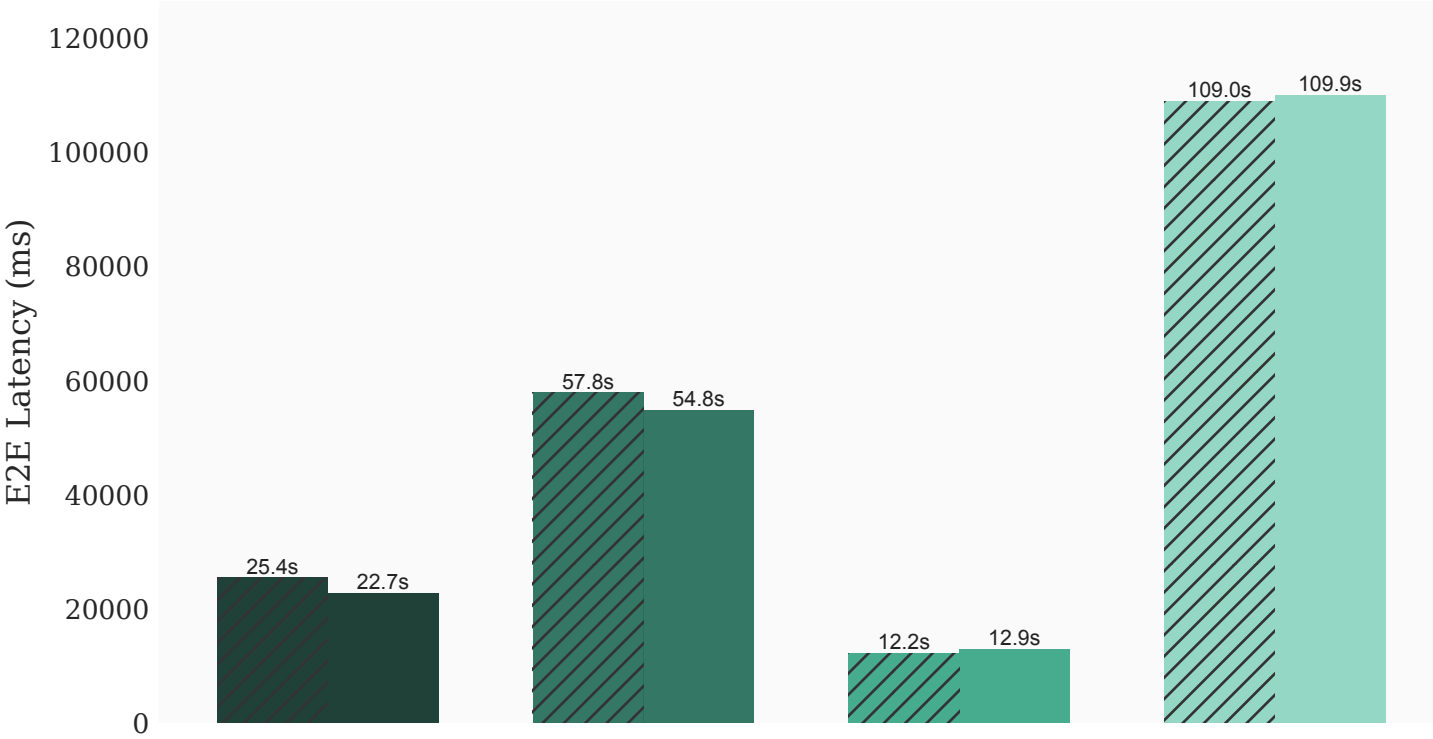- 67.2s, 66.0s
- 491.7s, 495.6s

Legend: CC / No CC

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

# Random (4000 ⇒ 1000) (Rate 1)

## Time to First Token (Mean)
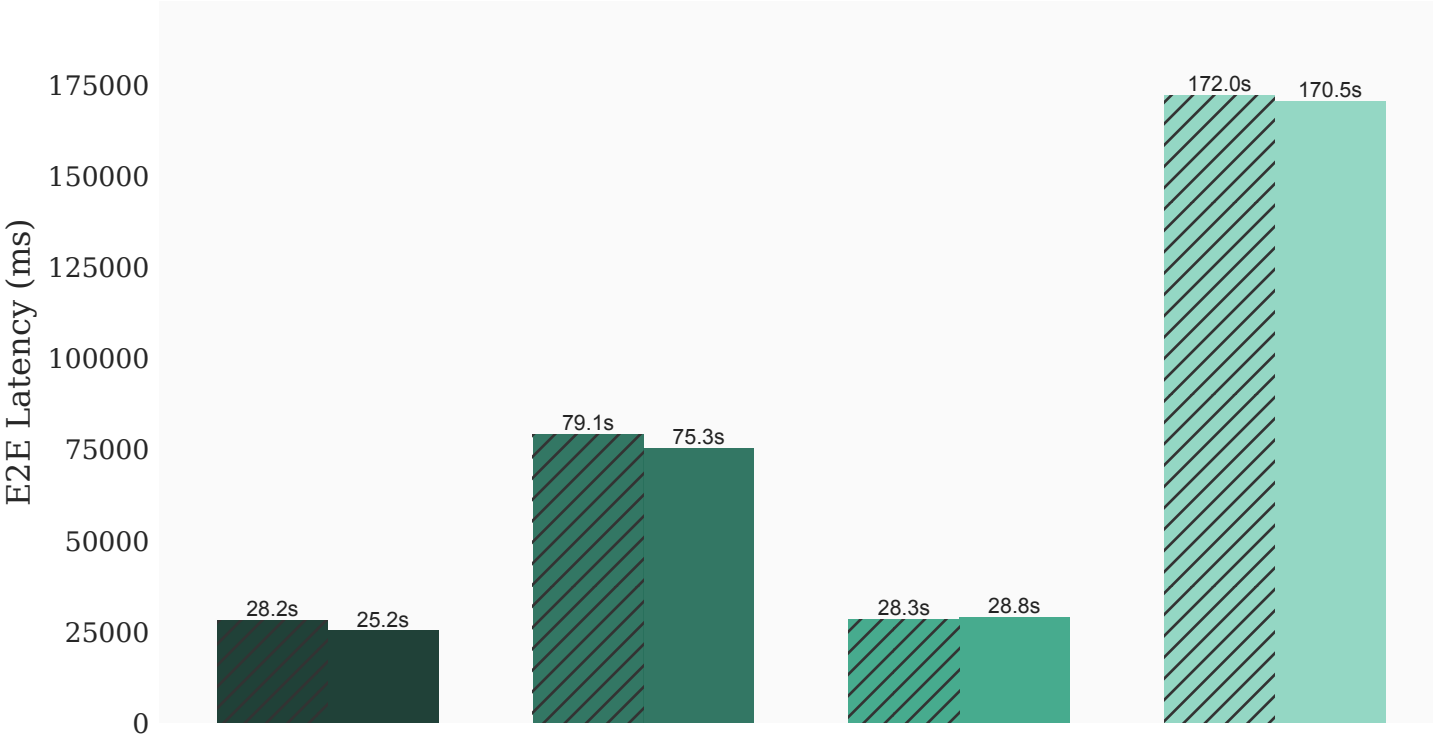
TTFT (ms)

128, 123 | 3.5s, 1.8s | 228, 226 | 127.8s, 129.5s

## End-to-End Latency (Mean)

E2E Latency (ms)

9.9s, 8.9s | 39.9s, 35.6s | 3.6s, 2.9s | 180.1s, 182.4s

## Time to First Token (P99)

TTFT (ms)

215, 217 | 10.7s, 9.0s | 443, 447 | 271.1s, 274.1s

## End-to-End Latency (P99)

E2E Latency (ms)

12.9s, 11.6s | 51.7s, 49.0s | 14.7s, 12.9s | 317.8s, 321.5s
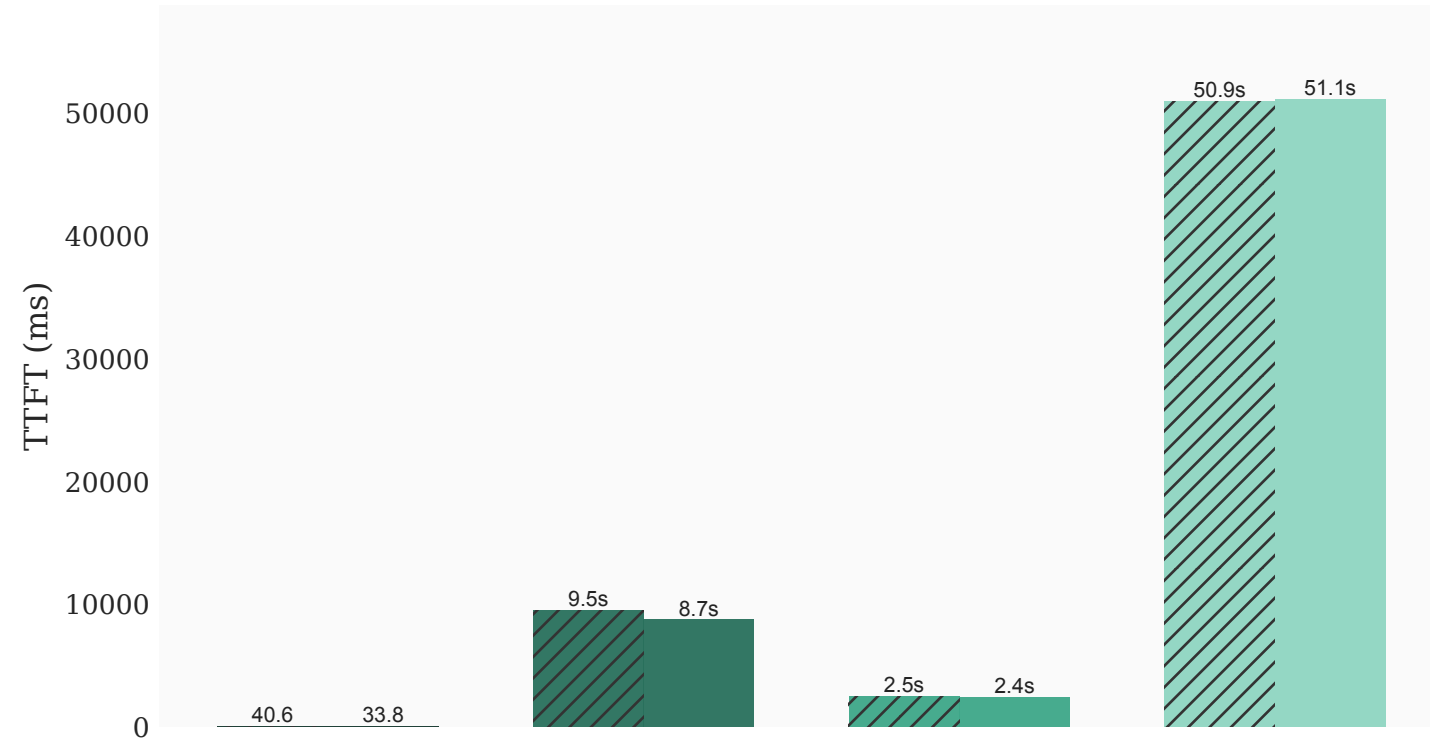
Legend: ▨ CC  ▬ No CC

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

# Random (1000 ⇒ 1000) (Rate 100)

## Time to First Token (Mean)



TTFT (ms)

- 1.3s / 1.2s — LLama 3.1 8B
- 10.3s / 9.7s — Mistral 3.1 24B
- 3.4s / 3.4s — GPT OSS 120B
- 52.0s / 52.0s — LLama 3.3 70B Int4

## End-to-End Latency (Mean)



E2E Latency (ms)

- 25.4s / 22.7s
- 57.8s / 54.8s
- 12.2s / 12.9s
- 109.0s / 109.9s

## Time to First Token (P99)



TTFT (ms)

- 2.1s / 2.0s
- 52.5s / 50.4s
- 6.6s / 6.5s
- 129.0s / 129.0s

## End-to-End Latency (P99)



E2E Latency (ms)

- 28.2s / 25.2s
- 79.1s / 75.3s
- 28.3s / 28.8s
- 172.0s / 170.5s

Legend: ▨ CC   ▤ No CC

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

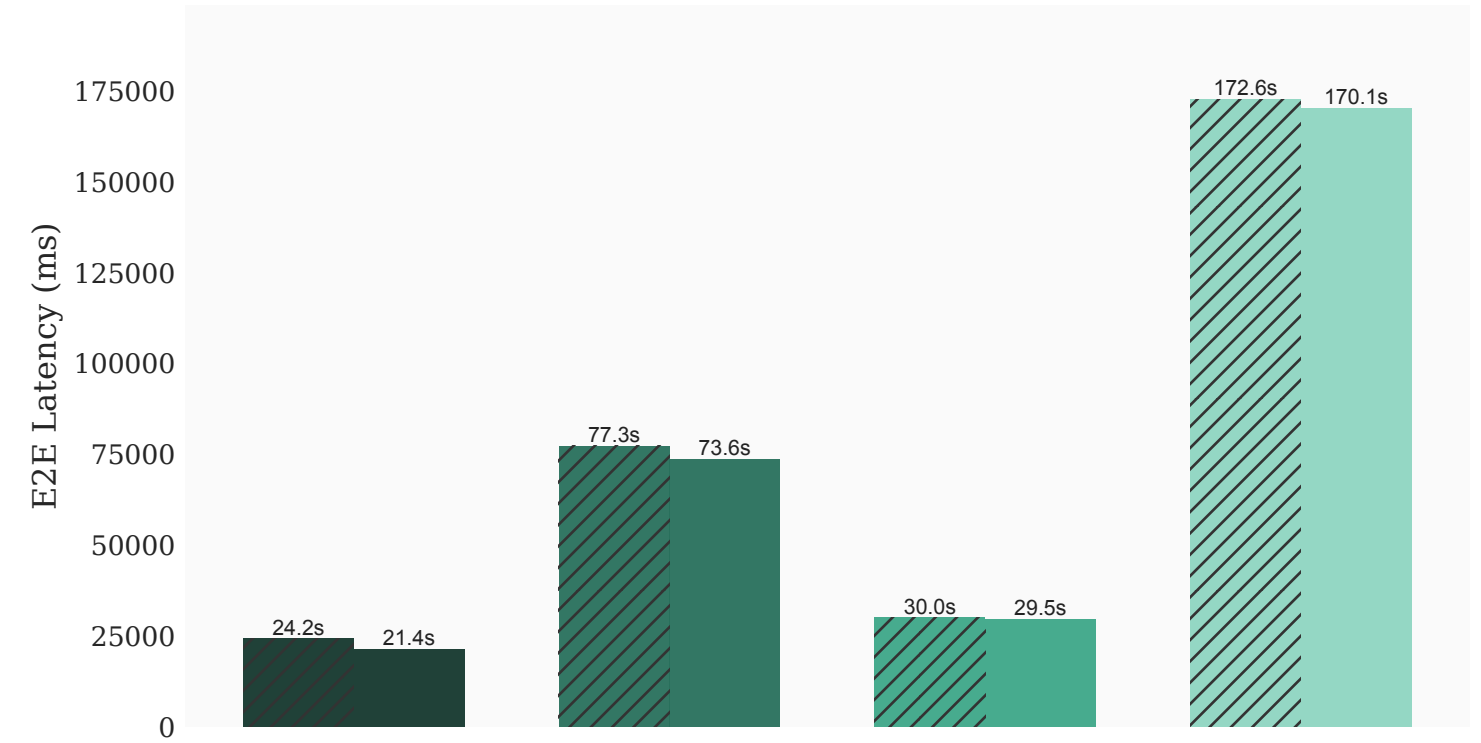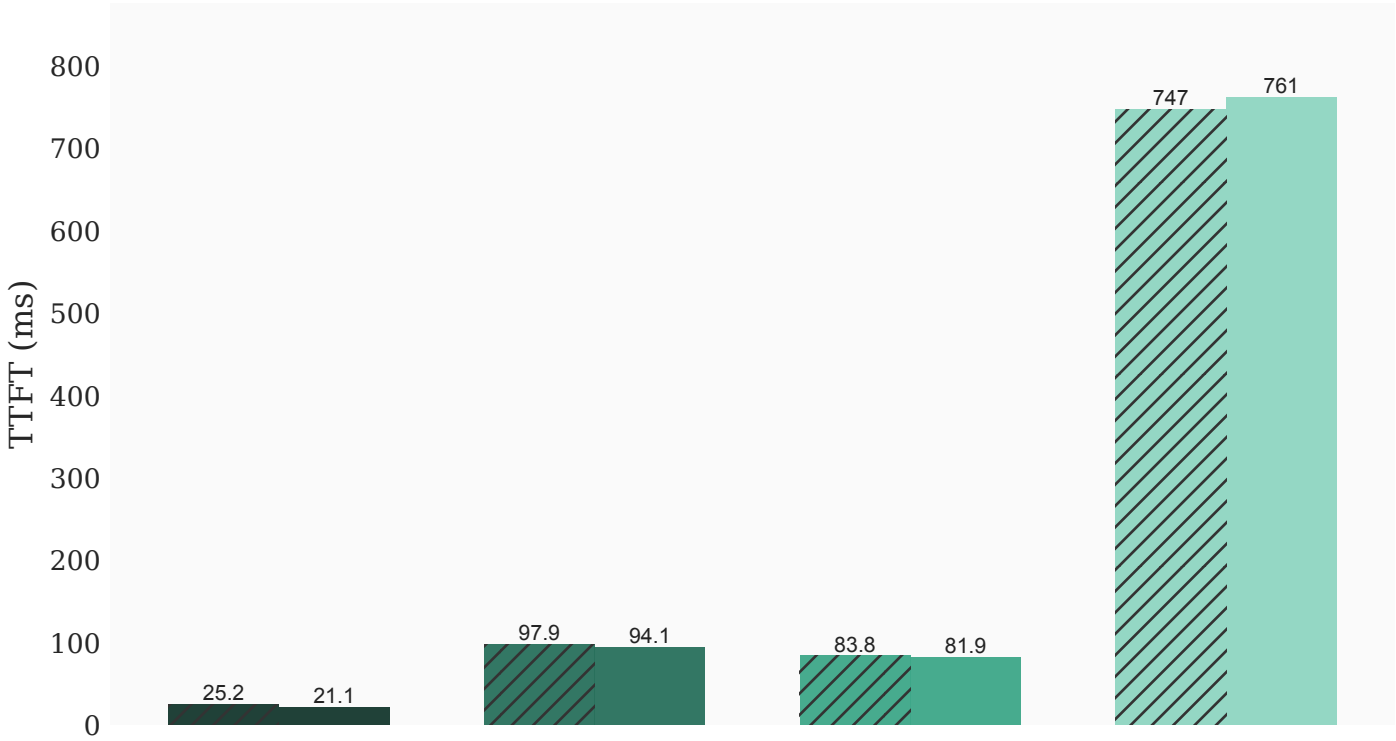# Random (1000 ⇒ 1000) (Rate 50)

## Time to First Token (Mean)

TTFT (ms)

- 40.6 / 33.8
- 9.5s / 8.7s
- 2.5s / 2.4s
- 50.9s / 51.1s

## End-to-End Latency (Mean)

E2E Latency (ms)

- 21.8s / 19.4s
- 57.0s / 54.2s
- 12.4s / 12.1s
- 108.7s / 108.9s

## Time to First Token (P99)

TTFT (ms)

- 66.9 / 56.1
- 50.4s / 48.6s
- 4.7s / 4.6s
- 126.9s / 124.0s

## End-to-End Latency (P99)

E2E Latency (ms)

- 24.2s / 21.4s
- 77.3s / 73.6s
- 30.0s / 29.5s
- 172.6s / 170.1s

Legend: ▨ CC    ▬ No CC

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

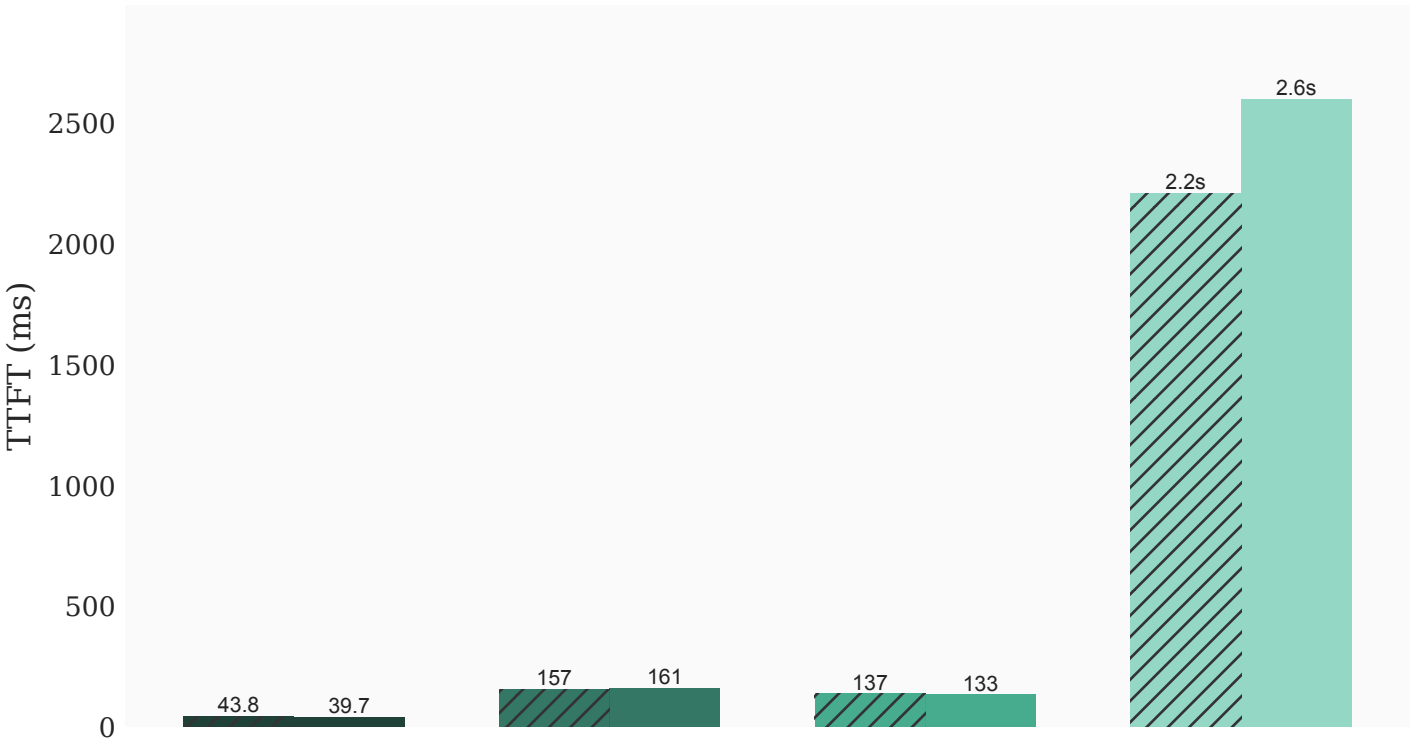# Random (1000 ⇒ 1000) (Rate 1)

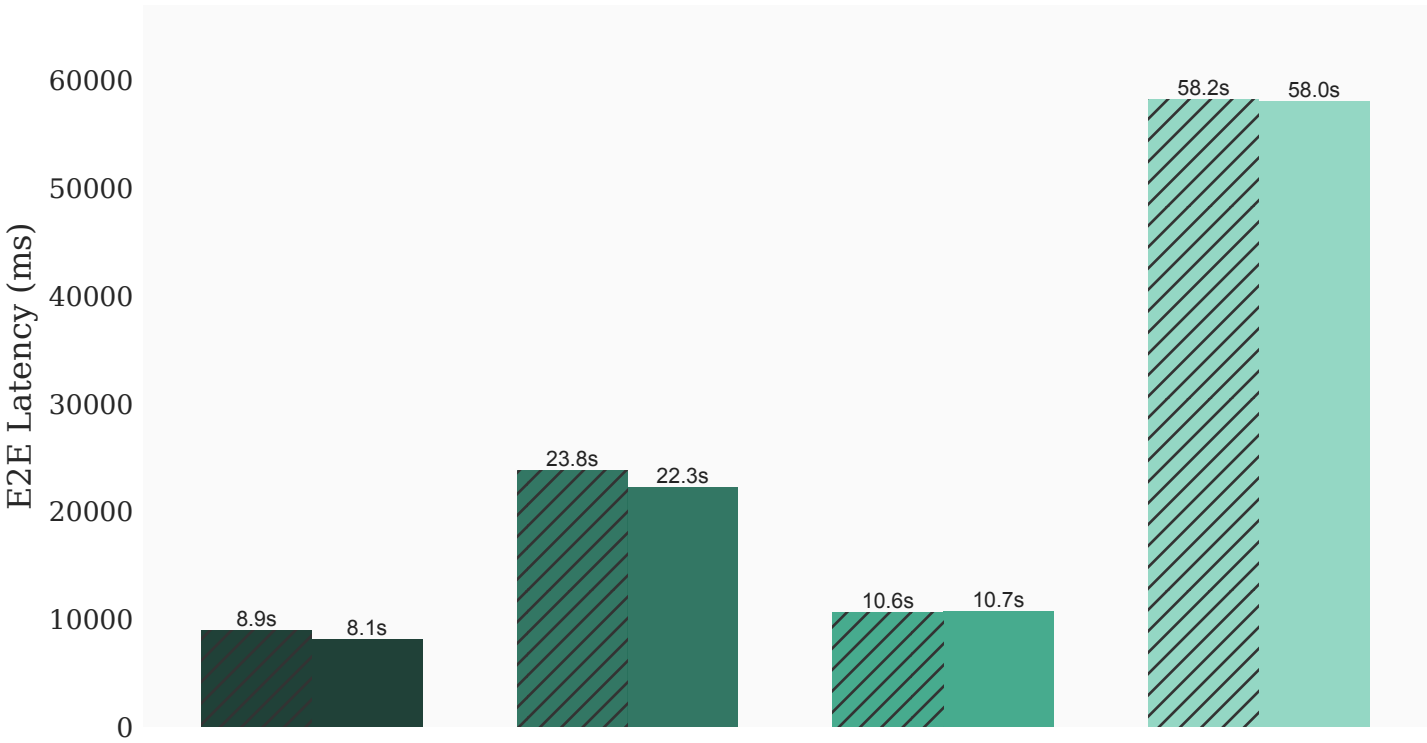## Time to First Token (Mean)



## End-to-End Latency (Mean)



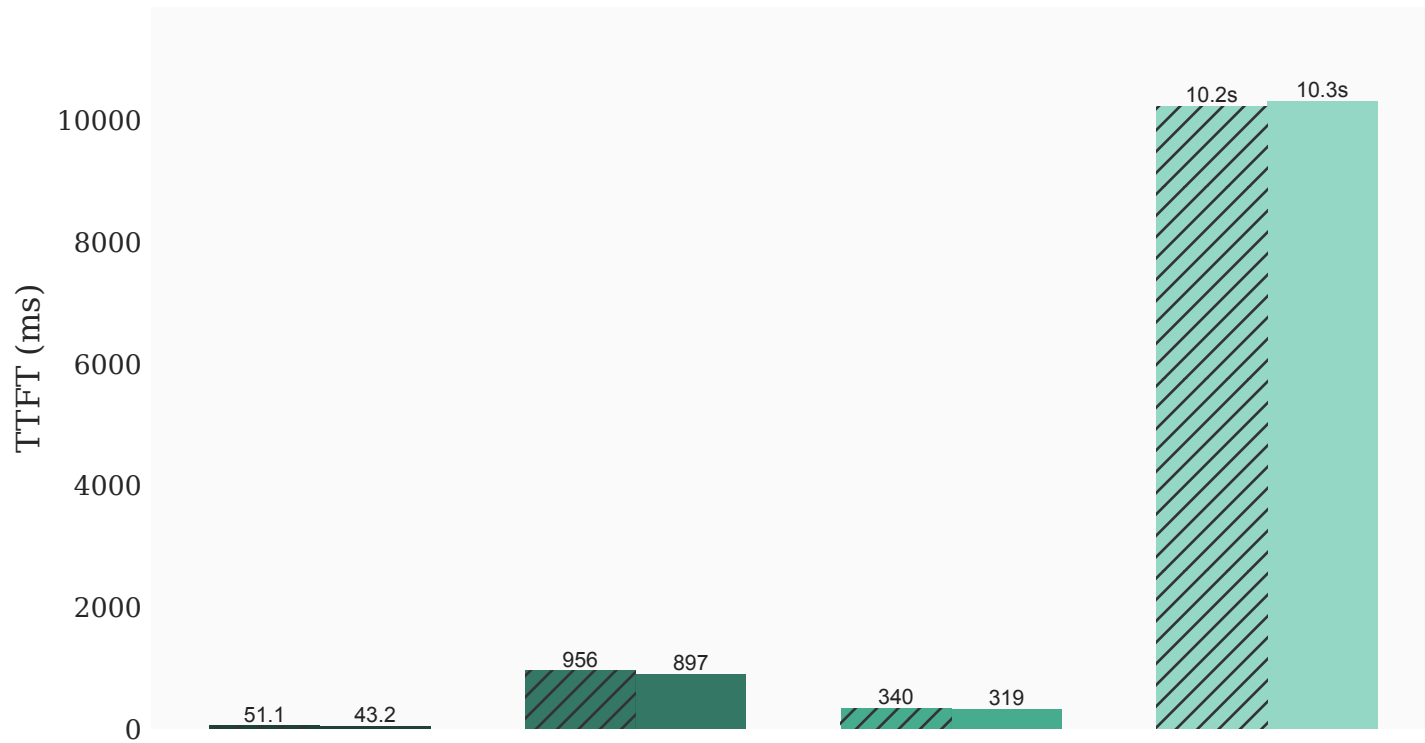## Time to First Token (P99)
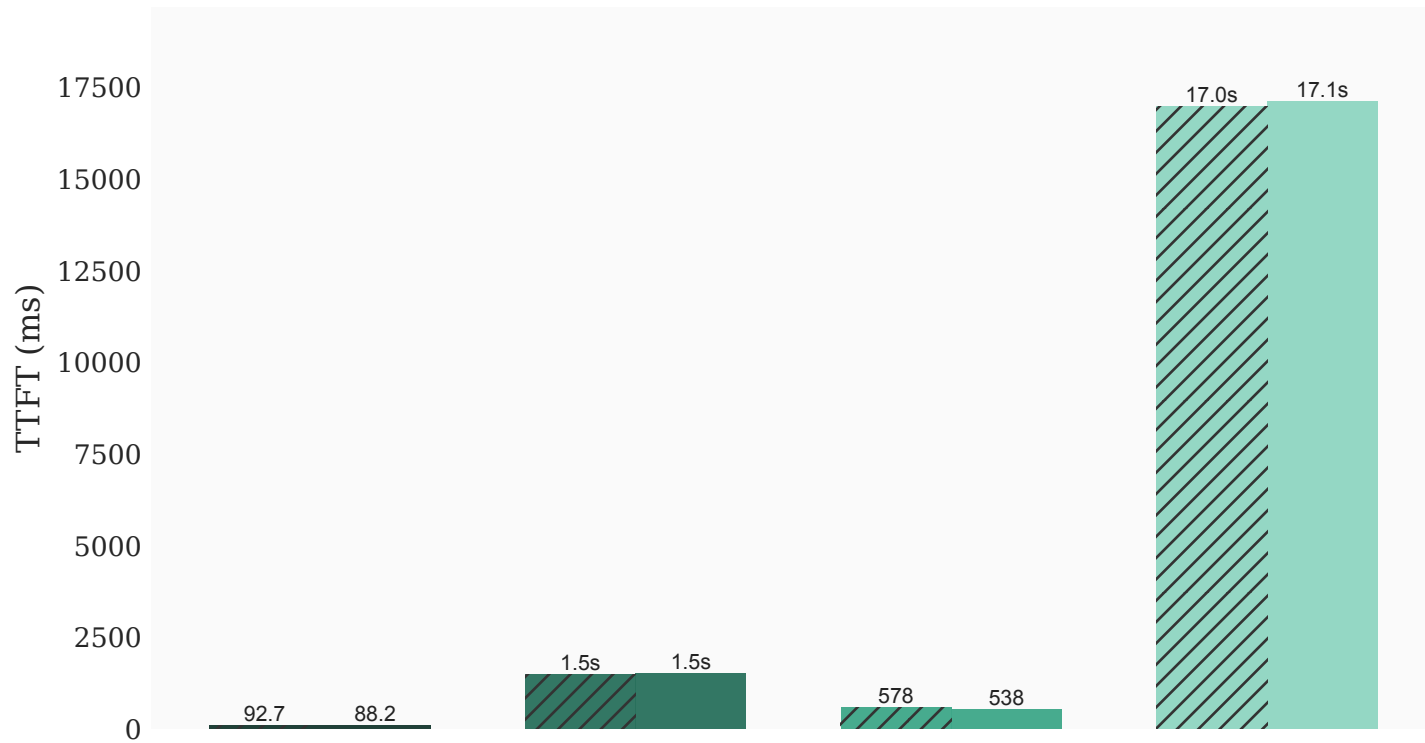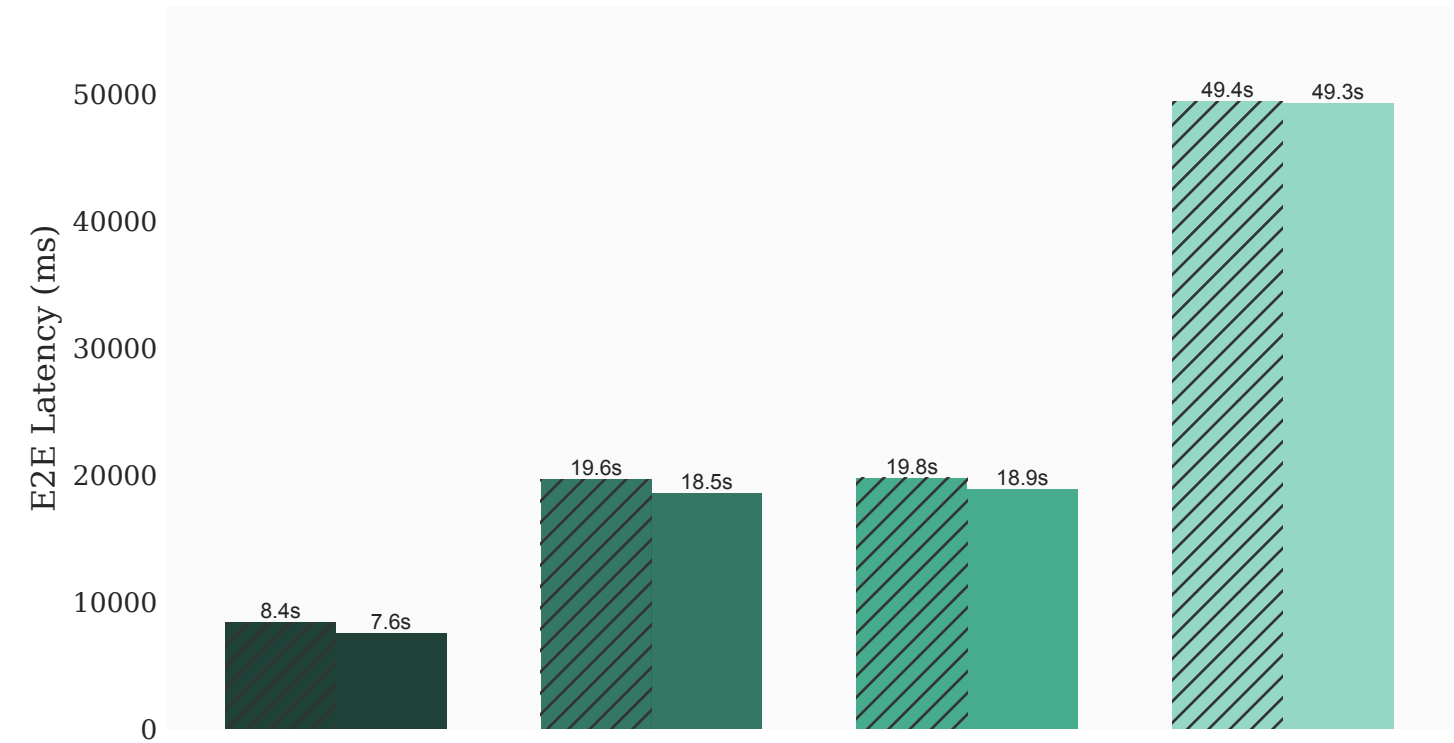


## End-to-End Latency (P99)



Legend: CC | No CC

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

## ShareGPT (Rate 100)

### Time to First Token (Mean)

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

- 51.1 / 43.2
- 956 / 897
- 340 / 319
- 10.2s / 10.3s

TTFT (ms)

### End-to-End Latency (Mean)

- 2.8s / 2.5s
- 6.5s / 6.1s
- 7.8s / 7.4s
- 29.1s / 29.2s

E2E Latency (ms)

### Time to First Token (P99)

- 92.7 / 88.2
- 1.5s / 1.5s
- 578 / 538
- 17.0s / 17.1s

TTFT (ms)

### End-to-End Latency (P99)

- 8.4s / 7.6s
- 19.6s / 18.5s
- 19.8s / 18.9s
- 49.4s / 49.3s

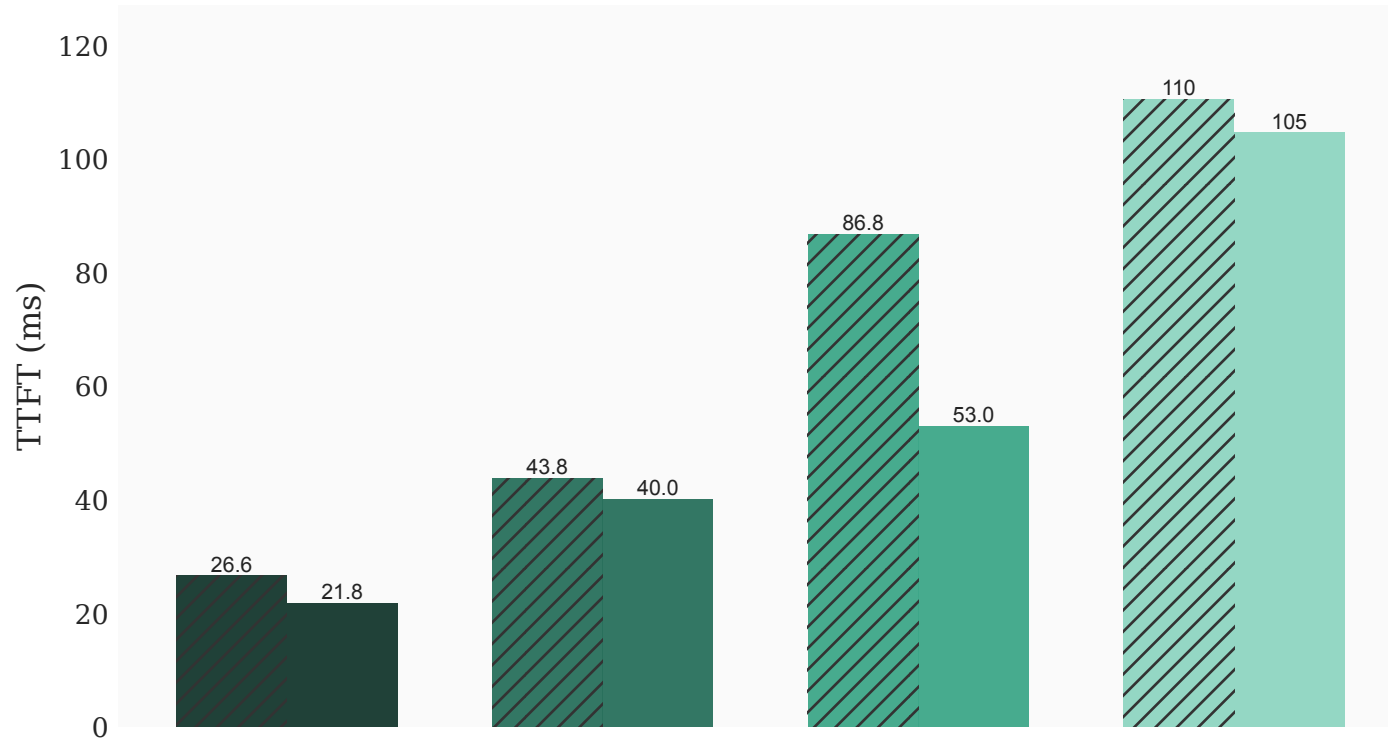E2E Latency (ms)

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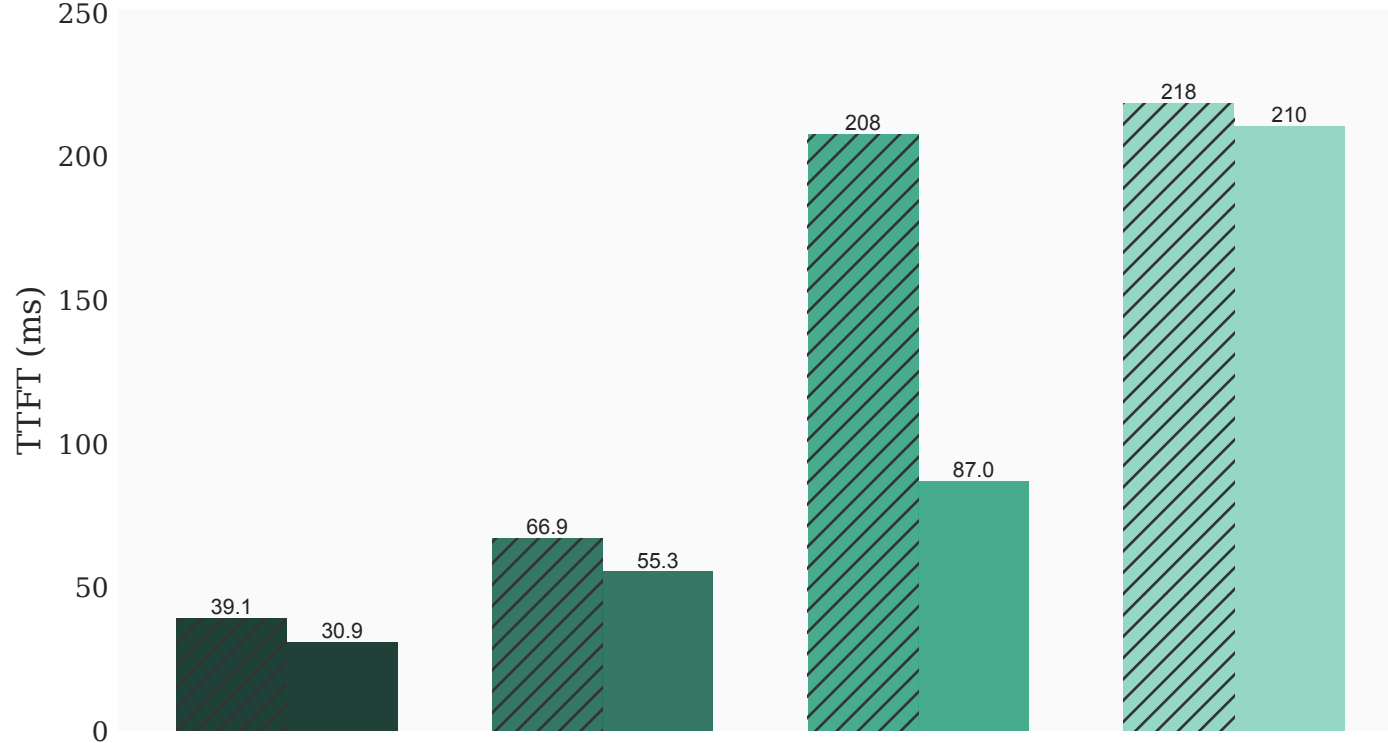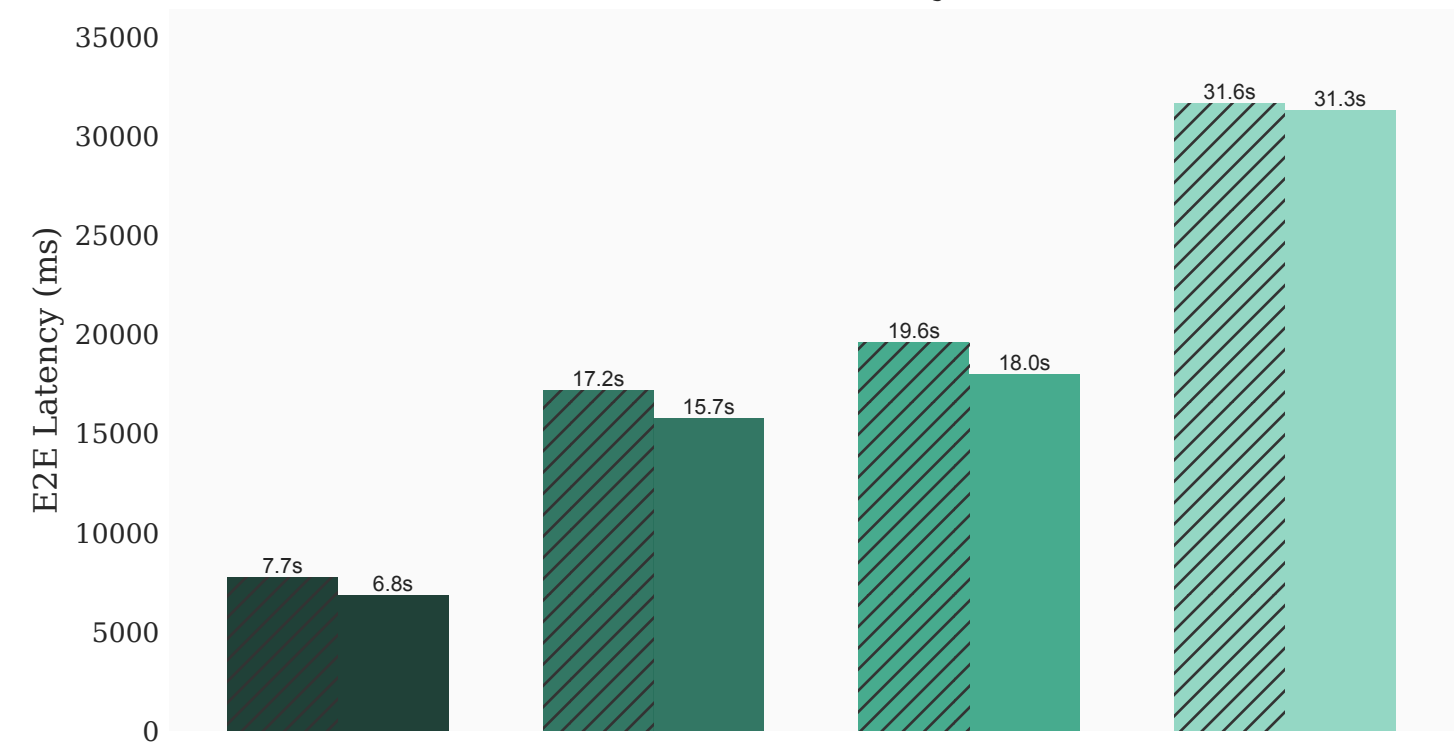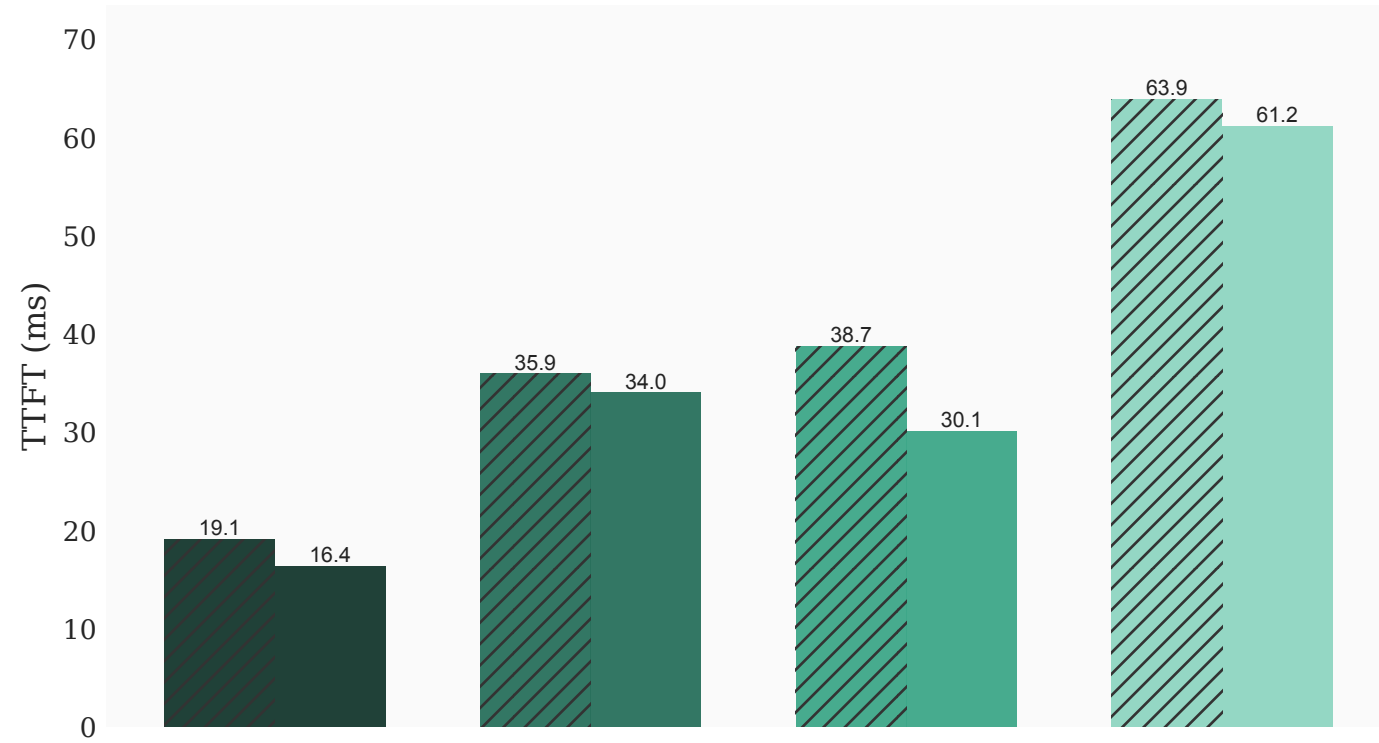
LLama 3.1 8B   Mistral 3.1 24B   GPT OSS 120B   LLama 3.3 70B Int4

# ShareGPT (Rate 50)

## Time to First Token (Mean)

Bar chart showing TTFT (ms) by model with CC (hatched) and No CC (solid) conditions:
- LLama 3.1 8B: CC 26.6, No CC 21.8
- Mistral 3.1 24B: CC 43.8, No CC 40.0
- GPT OSS 120B: CC 86.8, No CC 53.0
- LLama 3.3 70B Int4: CC 110, No CC 105

## End-to-End Latency (Mean)

Bar chart showing E2E Latency (ms):
- LLama 3.1 8B: CC 2.3s, No CC 2.0s
- Mistral 3.1 24B: CC 4.6s, No CC 4.2s
- GPT OSS 120B: CC 7.1s, No CC 6.1s
- LLama 3.3 70B Int4: CC 11.8s, No CC 11.7s

## Time to First Token (P99)

Bar chart showing TTFT (ms):
- LLama 3.1 8B: CC 39.1, No CC 30.9
- Mistral 3.1 24B: CC 66.9, No CC 55.3
- GPT OSS 120B: CC 208, No CC 87.0
- LLama 3.3 70B Int4: CC 218, No CC 210

## End-to-End Latency (P99)

Bar chart showing E2E Latency (ms):
- LLama 3.1 8B: CC 7.7s, No CC 6.8s
- Mistral 3.1 24B: CC 17.2s, No CC 15.7s
- GPT OSS 120B: CC 19.6s, No CC 18.0s
- LLama 3.3 70B Int4: CC 31.6s, No CC 31.3s

Legend: CC (hatched), No CC (solid)
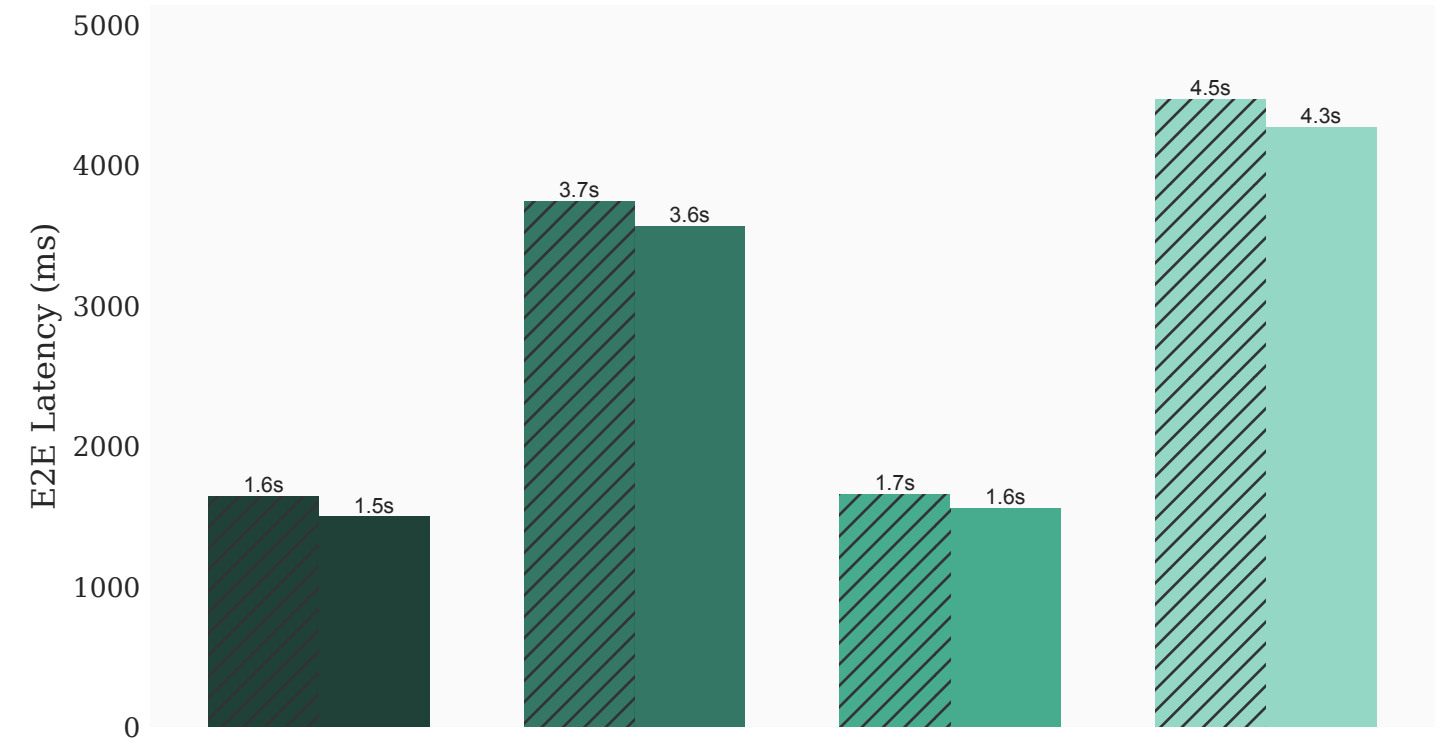
Models: LLama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, LLama 3.3 70B Int4
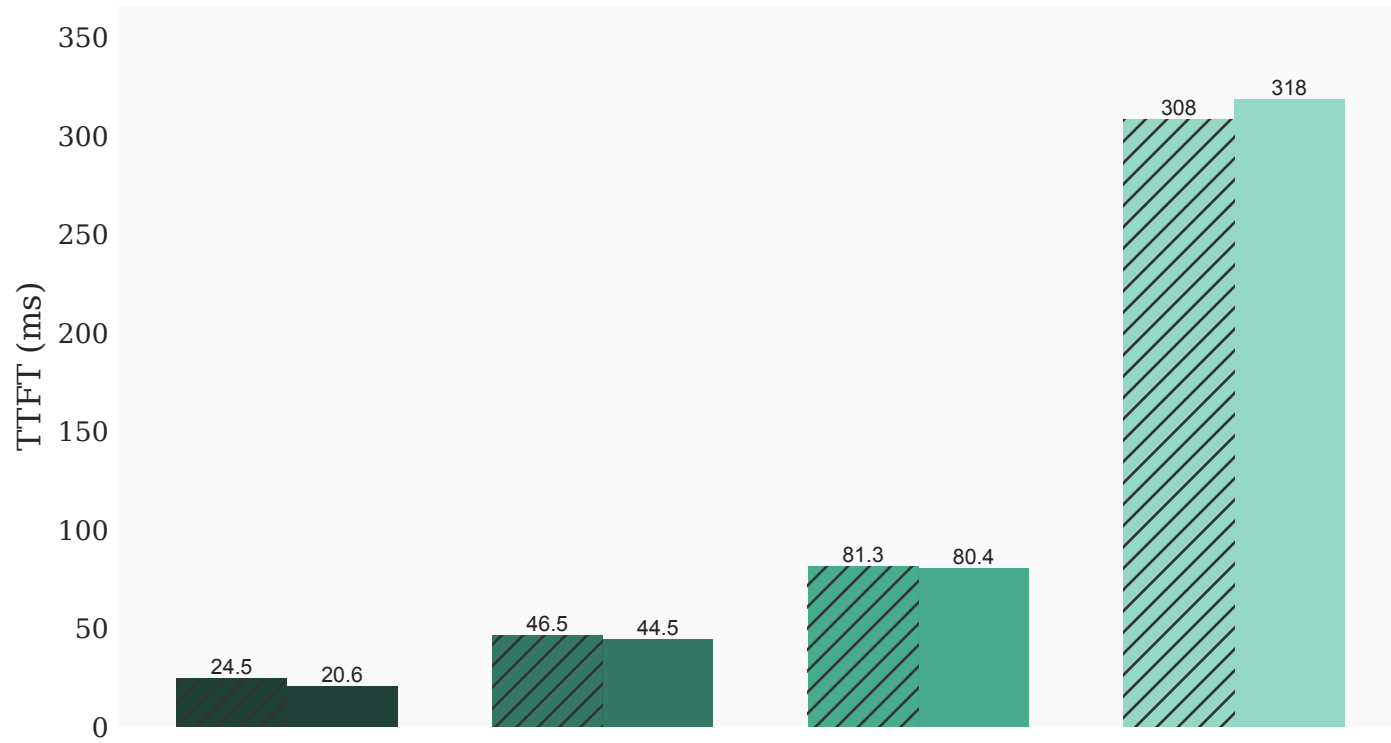
# ShareGPT (Rate 1)

## Time to First Token (Mean)



## End-to-End Latency (Mean)



## Time to First Token (P99)



## End-to-End Latency (P99)
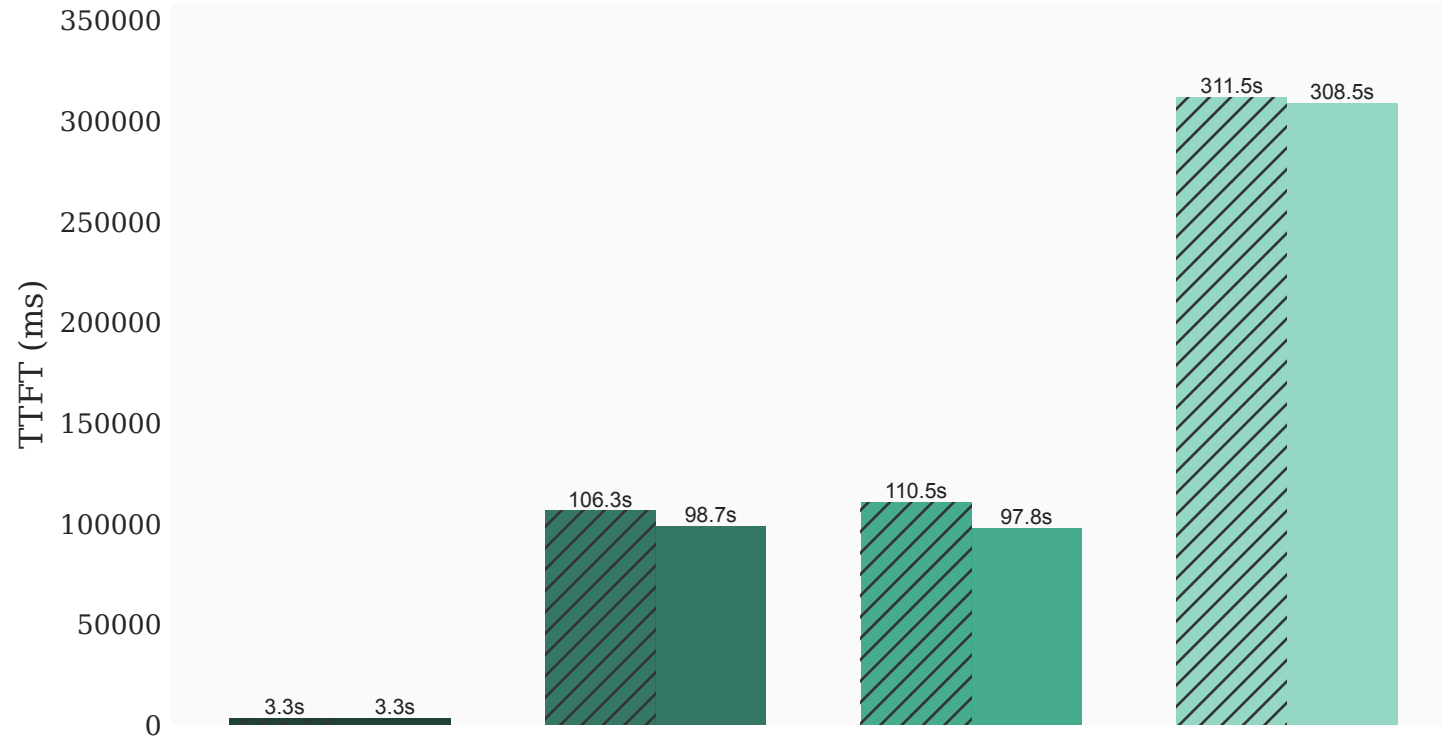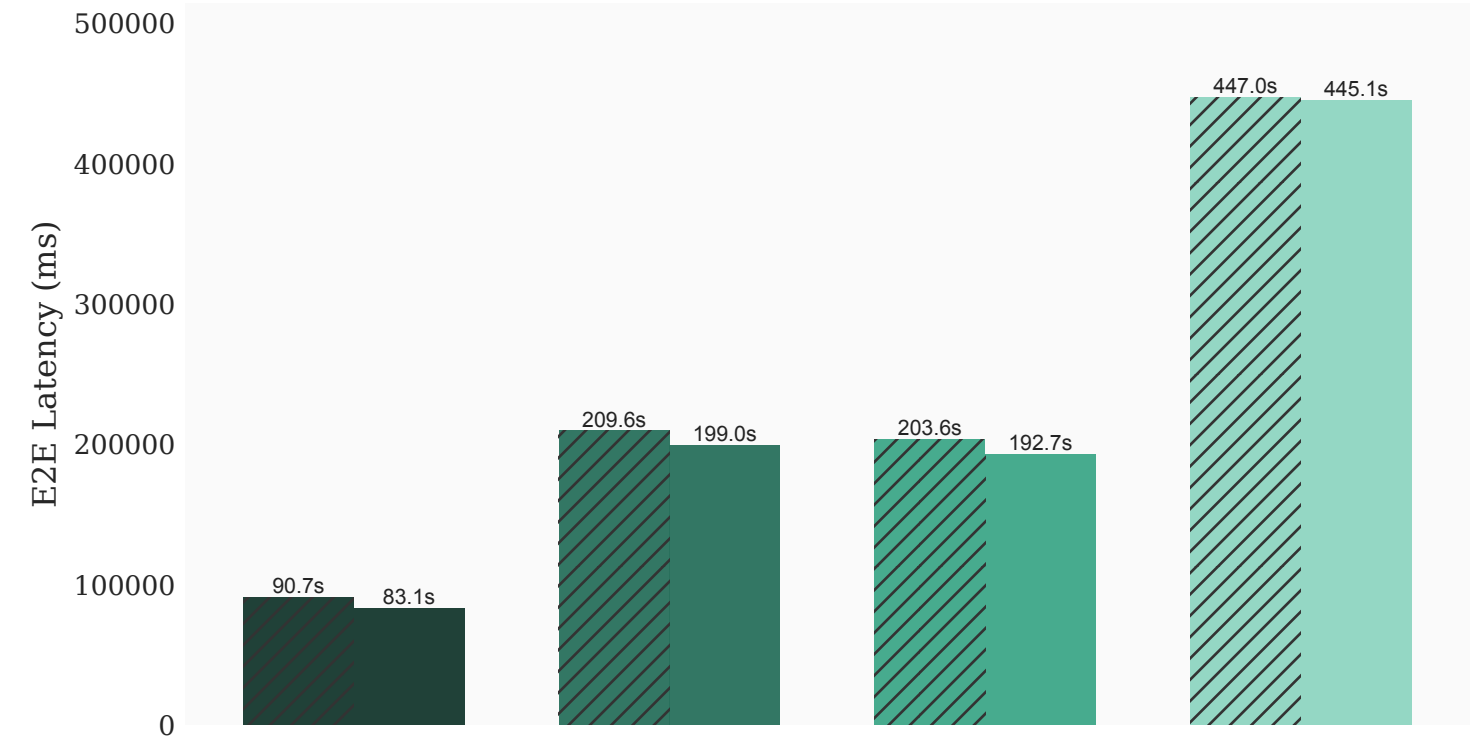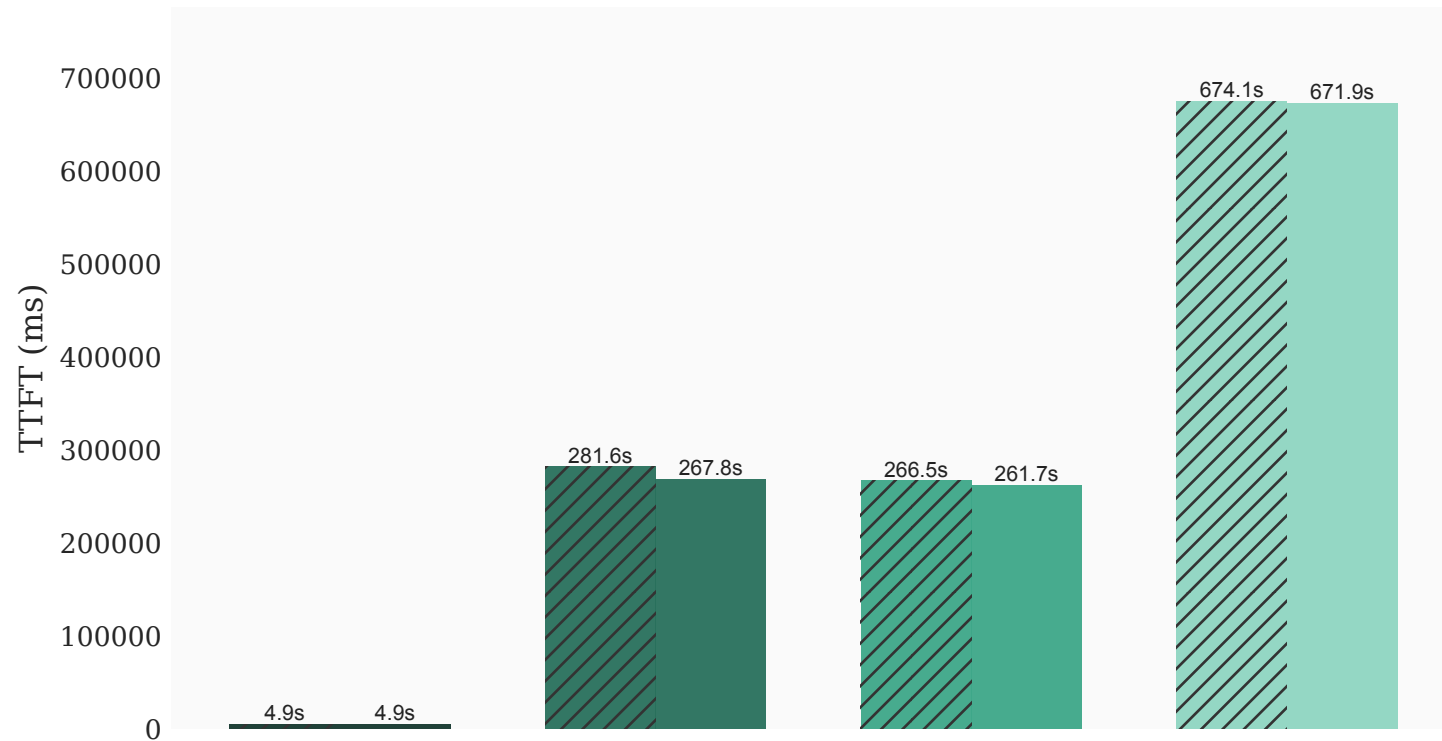


Legend: ▨ CC   ▬ No CC

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

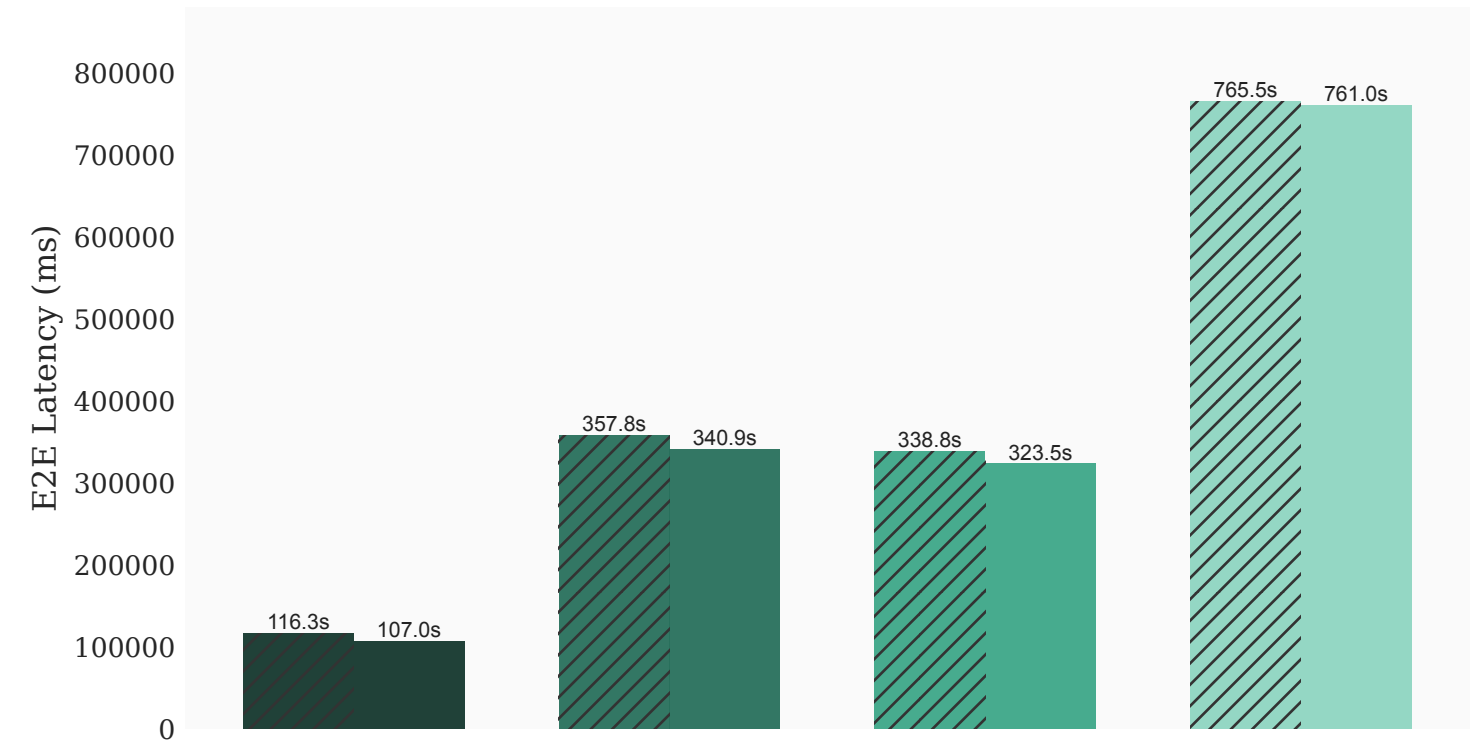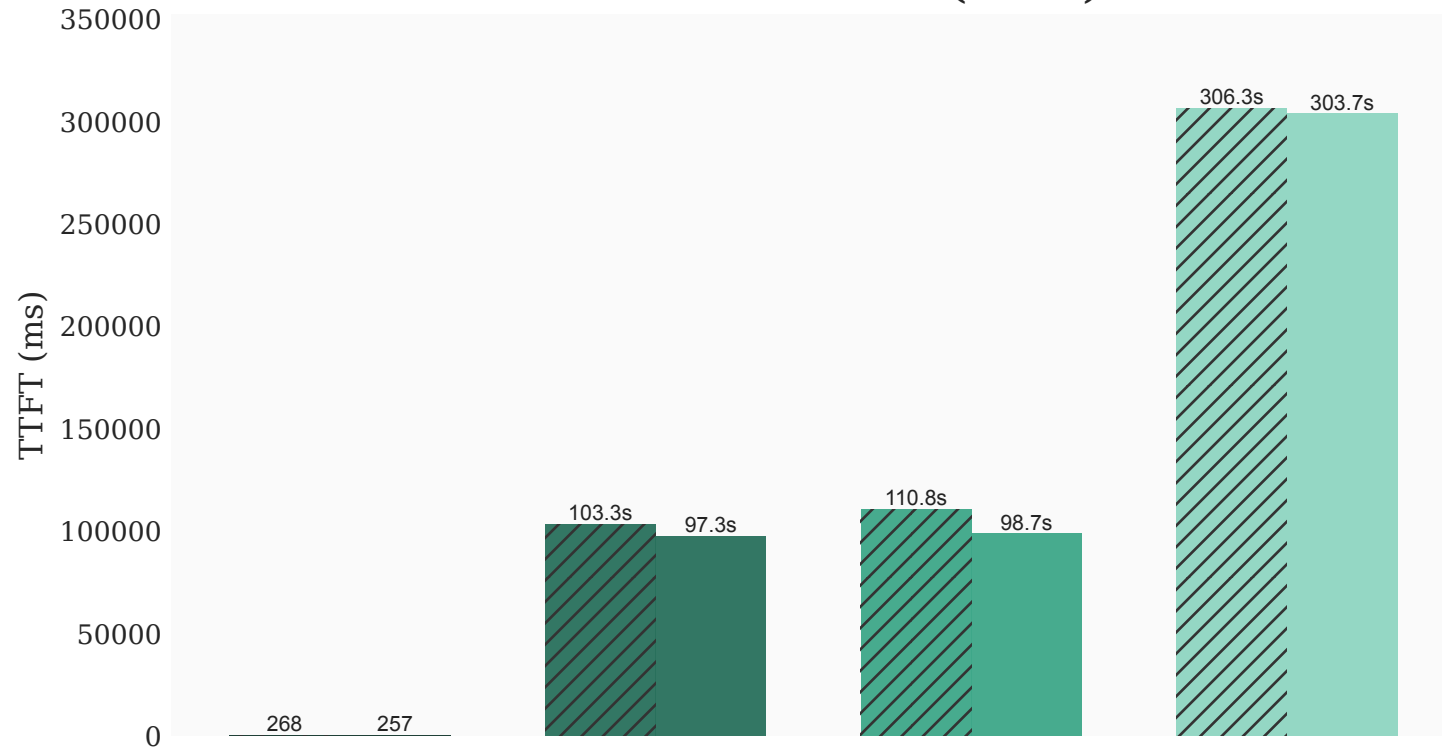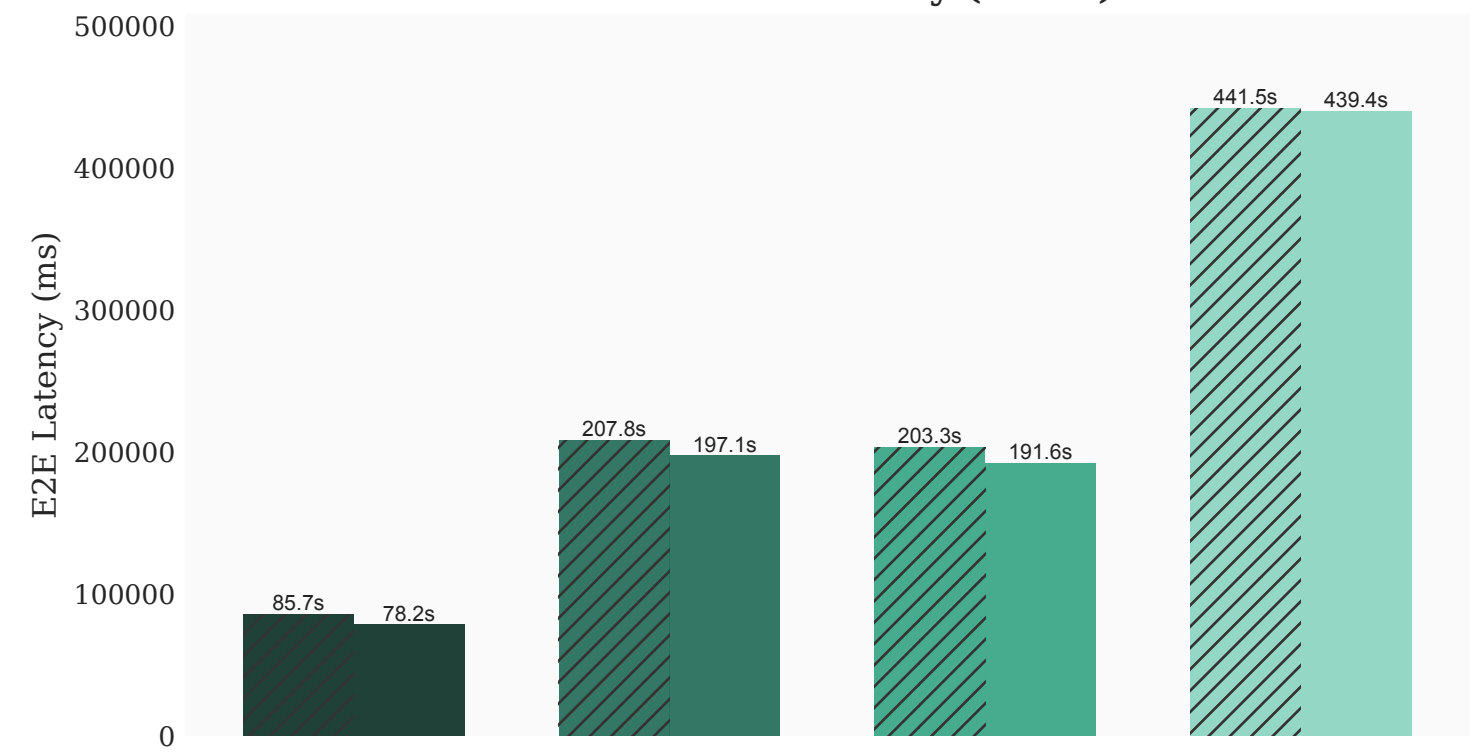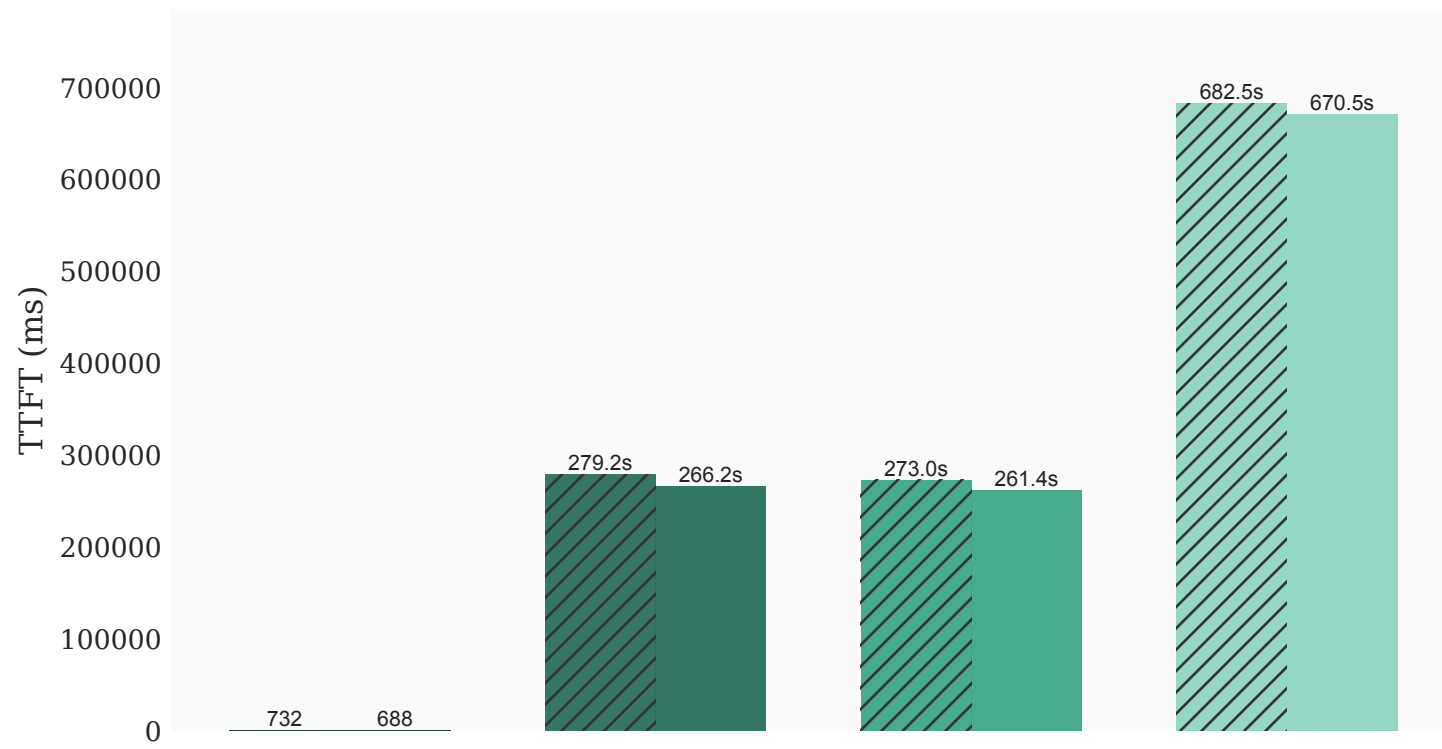# Edit 10K Characters (Rate 100)

## Time to First Token (Mean)

TTFT (ms)

- LLama 3.1 8B: 3.3s (CC), 3.3s (No CC)
- Mistral 3.1 24B: 106.3s (CC), 98.7s (No CC)
- GPT OSS 120B: 110.5s (CC), 97.8s (No CC)
- LLama 3.3 70B Int4: 311.5s (CC), 308.5s (No CC)

## End-to-End Latency (Mean)

E2E Latency (ms)

- LLama 3.1 8B: 90.7s (CC), 83.1s (No CC)
- Mistral 3.1 24B: 209.6s (CC), 199.0s (No CC)
- GPT OSS 120B: 203.6s (CC), 192.7s (No CC)
- LLama 3.3 70B Int4: 447.0s (CC), 445.1s (No CC)

## Time to First Token (P99)

TTFT (ms)

- LLama 3.1 8B: 4.9s (CC), 4.9s (No CC)
- Mistral 3.1 24B: 281.6s (CC), 267.8s (No CC)
- GPT OSS 120B: 266.5s (CC), 261.7s (No CC)
- LLama 3.3 70B Int4: 674.1s (CC), 671.9s (No CC)

## End-to-End Latency (P99)

E2E Latency (ms)

- LLama 3.1 8B: 116.3s (CC), 107.0s (No CC)
- Mistral 3.1 24B: 357.8s (CC), 340.9s (No CC)
- GPT OSS 120B: 338.8s (CC), 323.5s (No CC)
- LLama 3.3 70B Int4: 765.5s (CC), 761.0s (No CC)

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

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

# Edit 10K Characters (Rate 50)

## Time to First Token (Mean)



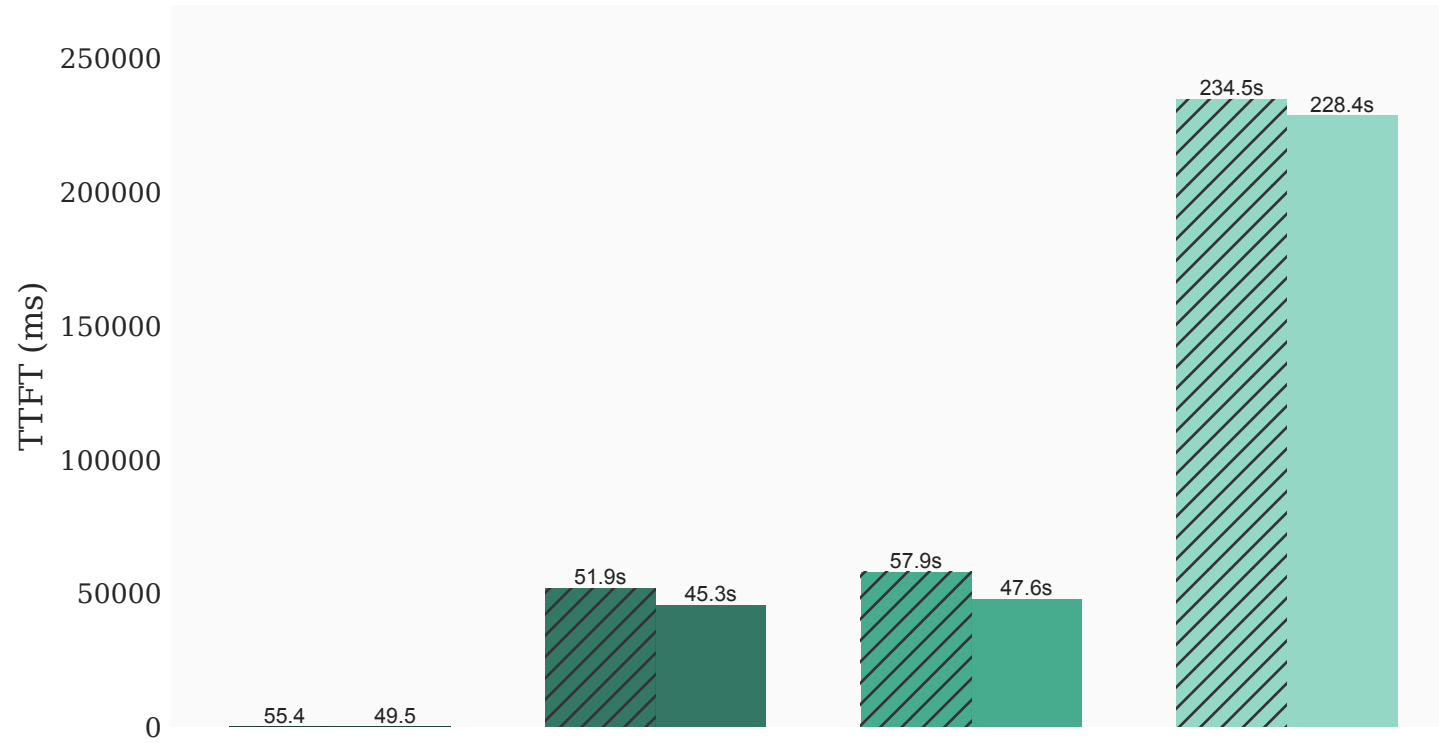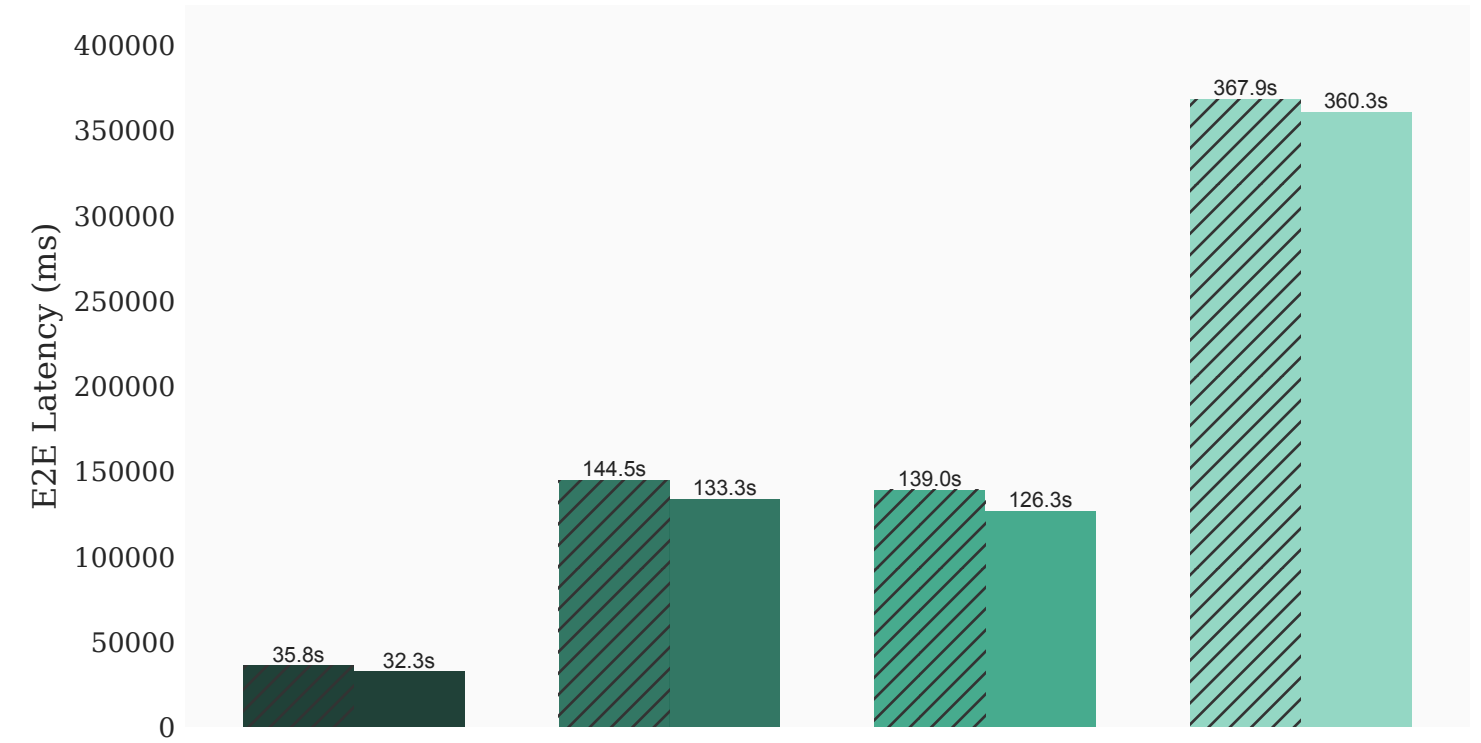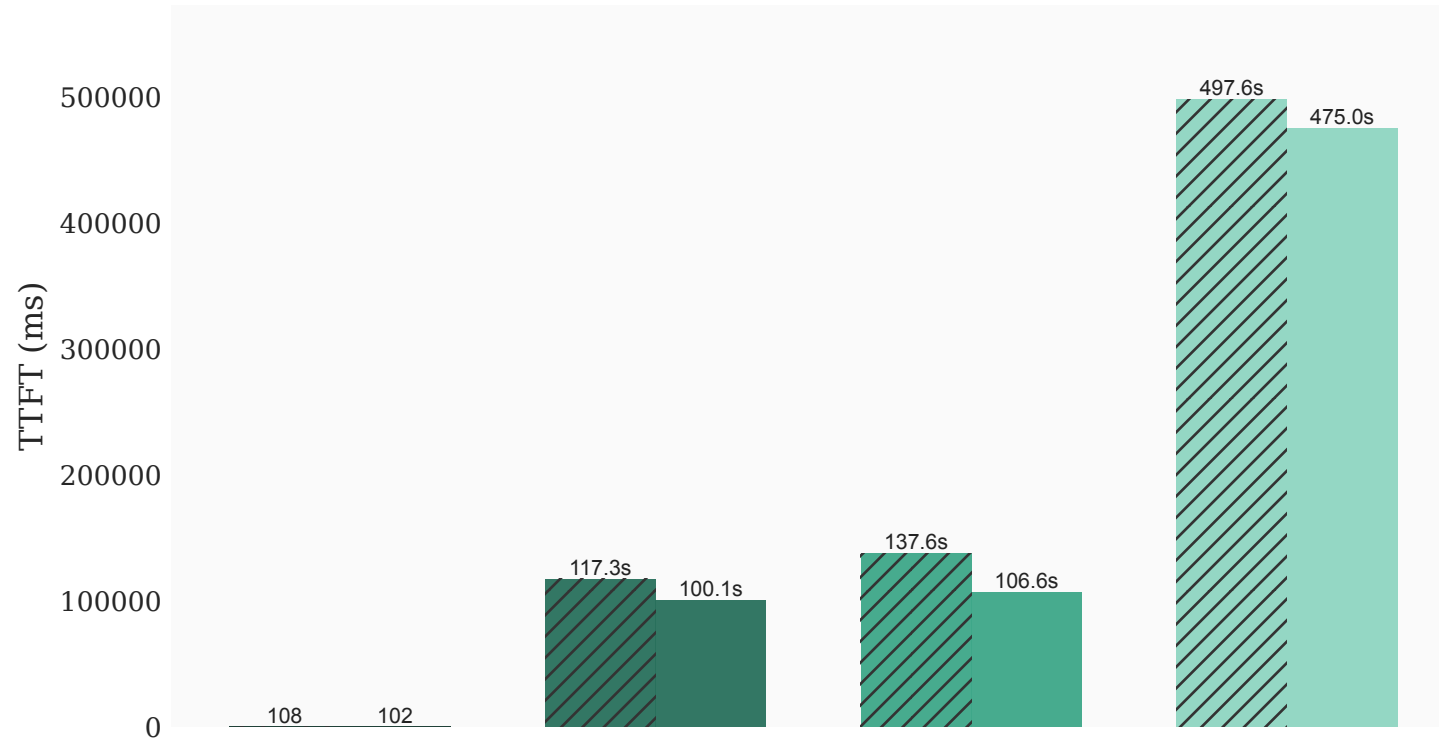TTFT (ms)

- 268, 257
- 103.3s, 97.3s
- 110.8s, 98.7s
- 306.3s, 303.7s

## End-to-End Latency (Mean)



E2E Latency (ms)

- 85.7s, 78.2s
- 207.8s, 197.1s
- 203.3s, 191.6s
- 441.5s, 439.4s

## Time to First Token (P99)



TTFT (ms)

- 732, 688
- 279.2s, 266.2s
- 273.0s, 261.4s
- 682.5s, 670.5s

## End-to-End Latency (P99)



E2E Latency (ms)

- 112.0s, 101.3s
- 355.5s, 336.1s
- 338.0s, 320.7s
- 765.8s, 753.6s

CC    No CC

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

# Edit 10K Characters (Rate 1)

## Time to First Token (Mean)
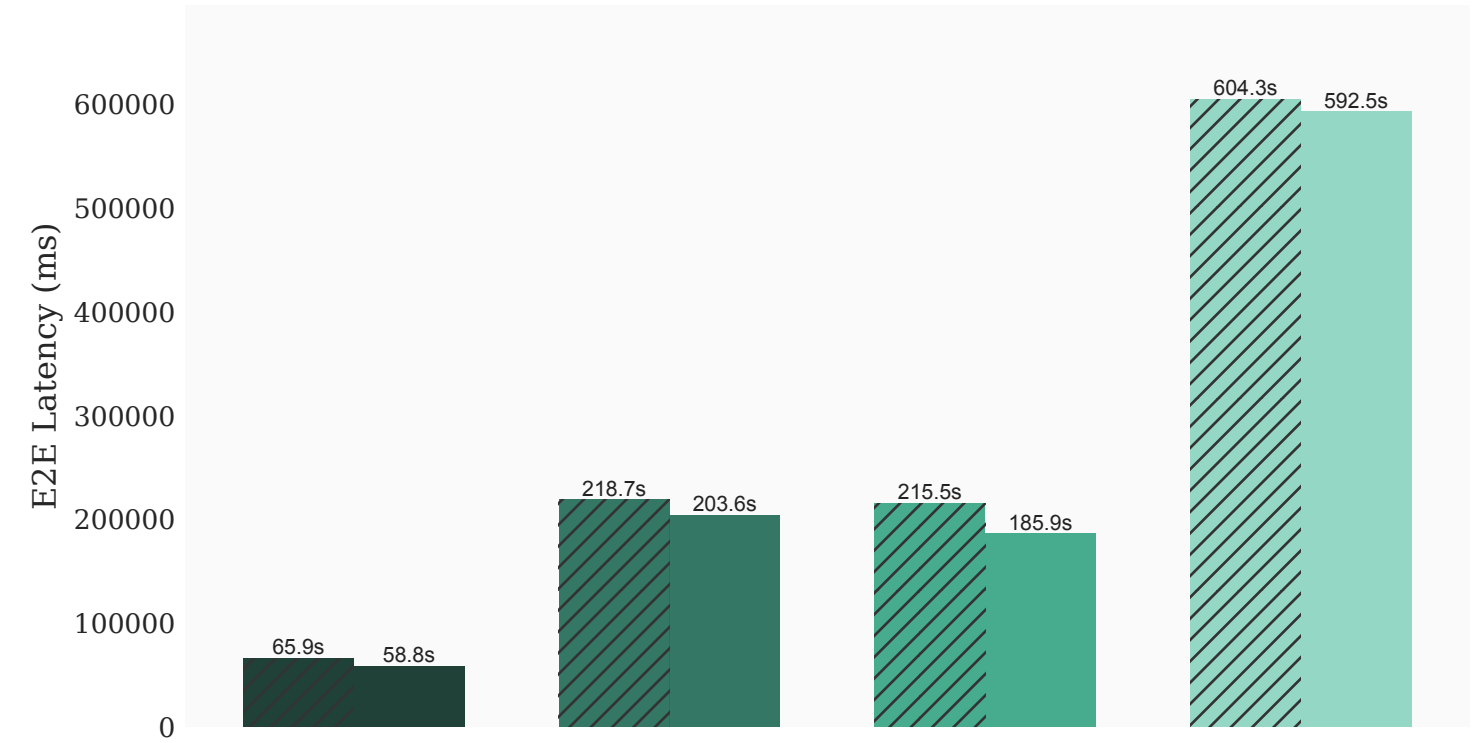
TTFT (ms)

- 55.4 / 49.5 — LLama 3.1 8B
- 51.9s (CC) / 45.3s (No CC) — Mistral 3.1 24B
- 57.9s (CC) / 47.6s (No CC) — GPT OSS 120B
- 234.5s (CC) / 228.4s (No CC) — LLama 3.3 70B Int4

## End-to-End Latency (Mean)

E2E Latency (ms)

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

## Time to First Token (P99)

TTFT (ms)

- 108 / 102 — LLama 3.1 8B
- 117.3s (CC) / 100.1s (No CC) — Mistral 3.1 24B
- 137.6s (CC) / 106.6s (No CC) — GPT OSS 120B
- 497.6s (CC) / 475.0s (No CC) — LLama 3.3 70B Int4

## End-to-End Latency (P99)

E2E Latency (ms)

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

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

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

# Numina Math (Rate 100)

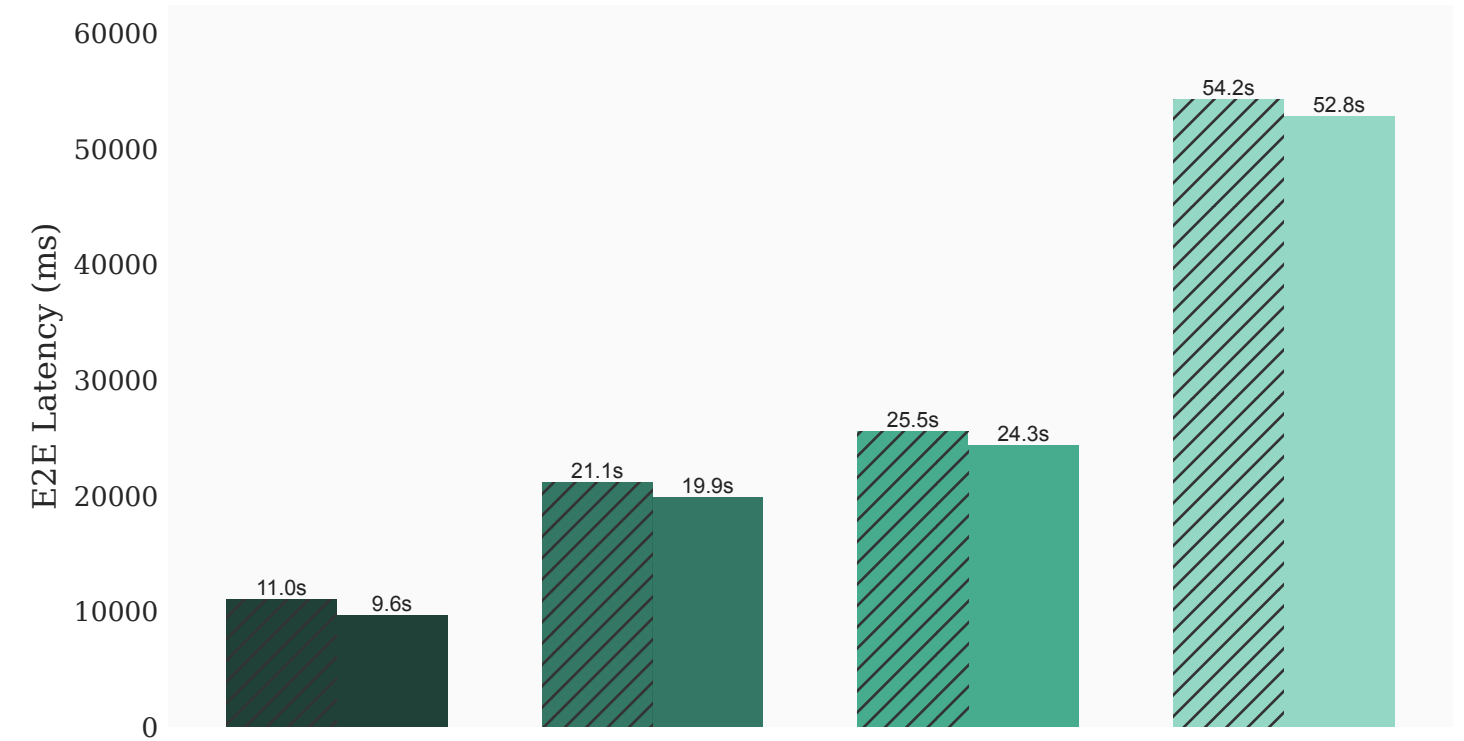## Time to First Token (Mean)



## End-to-End Latency (Mean)



## Time to First Token (P99)
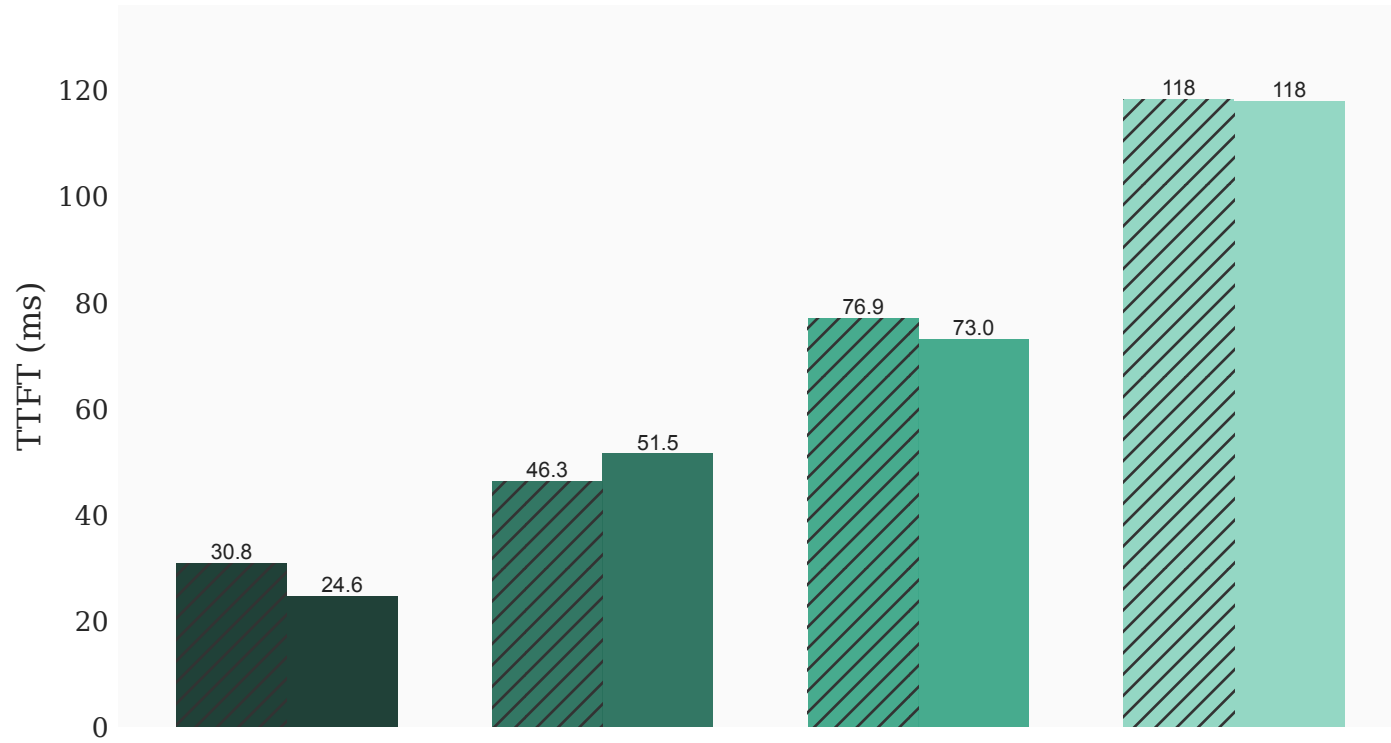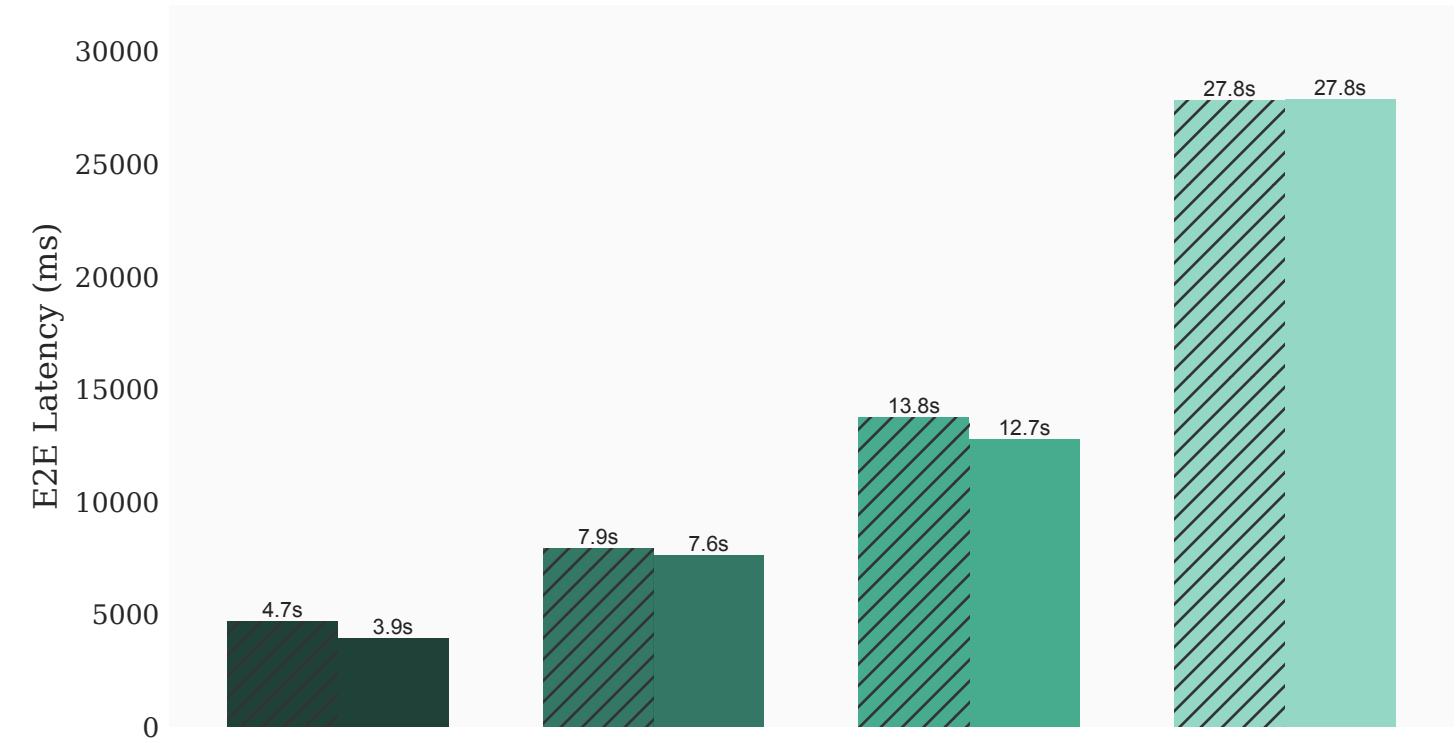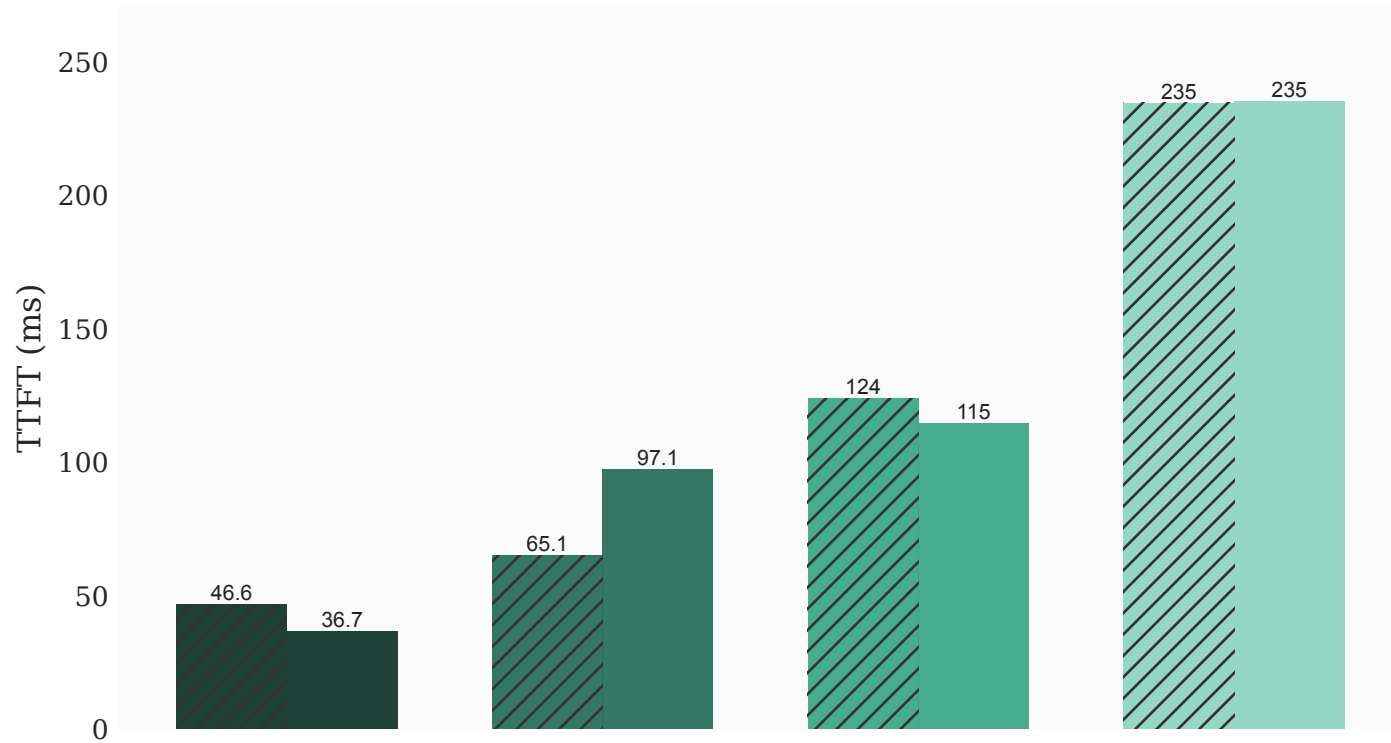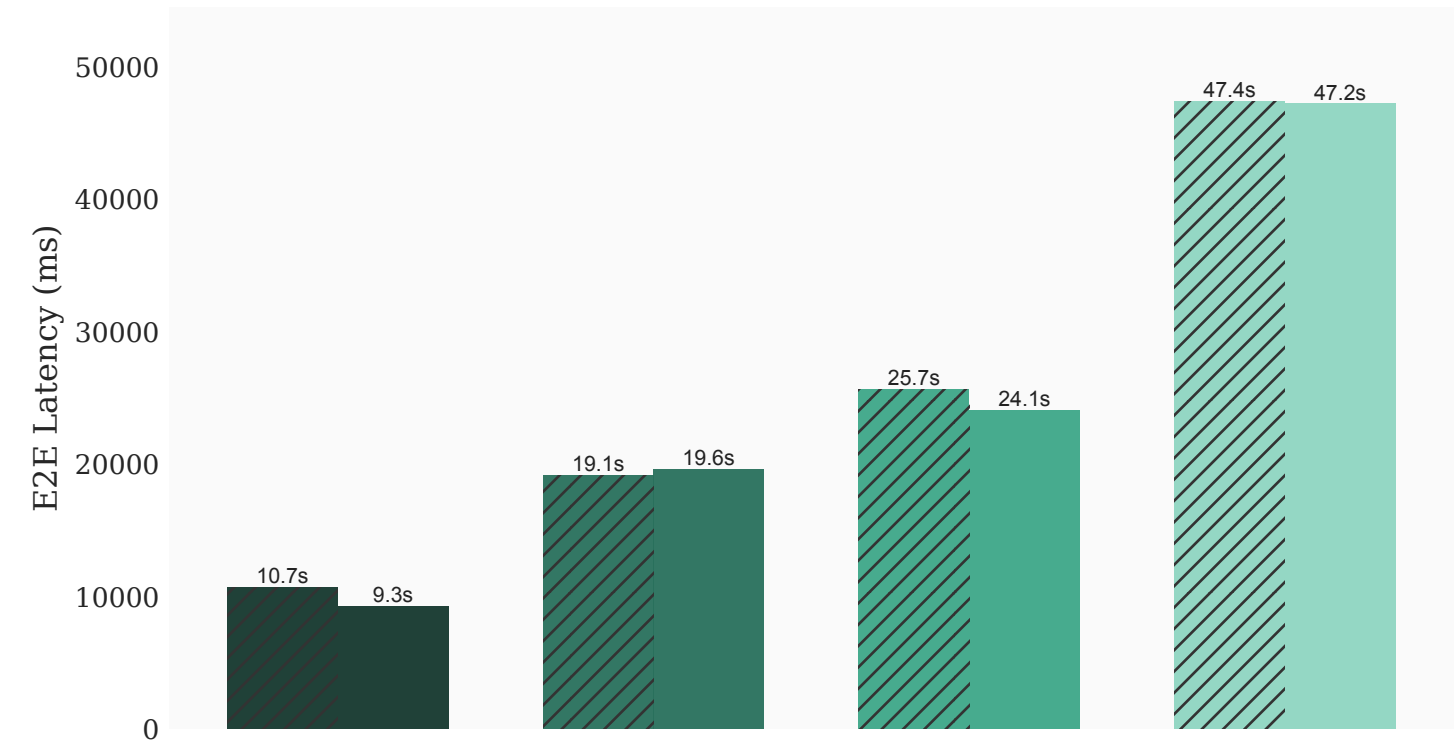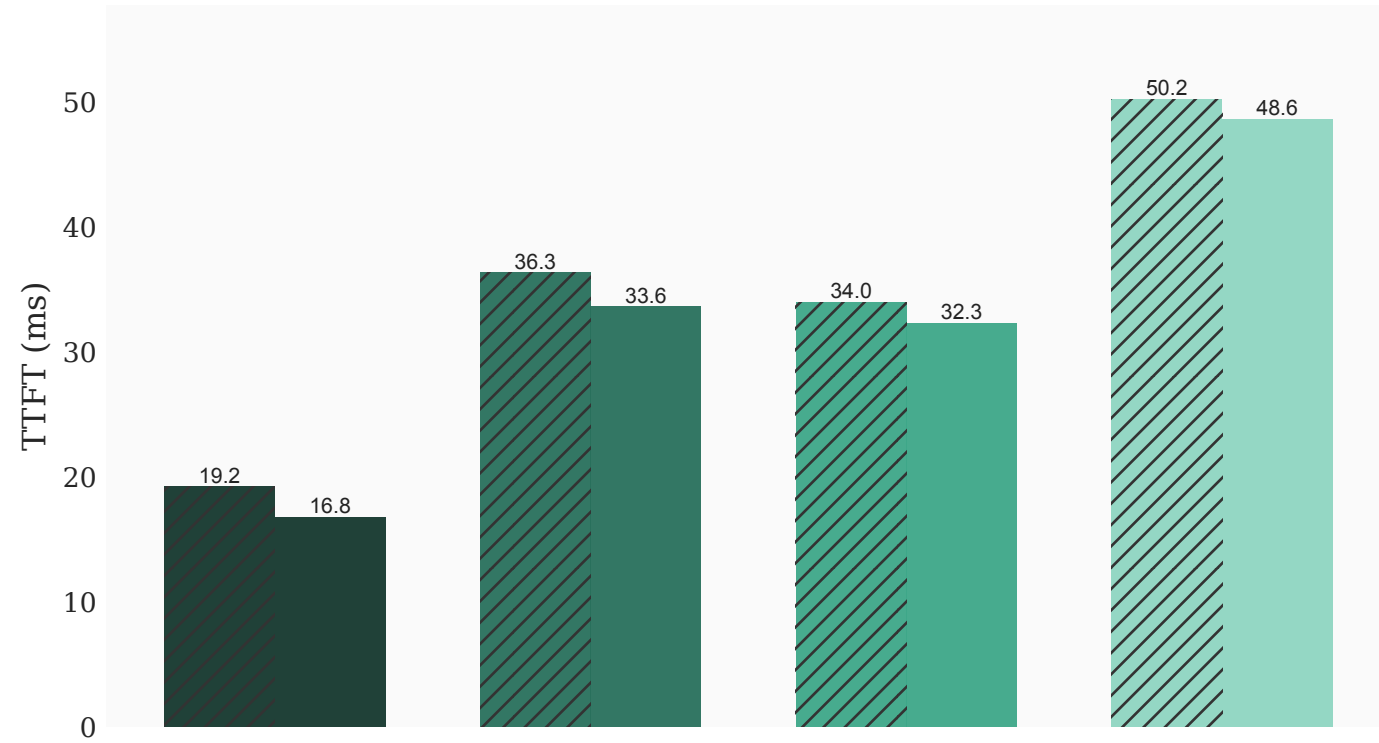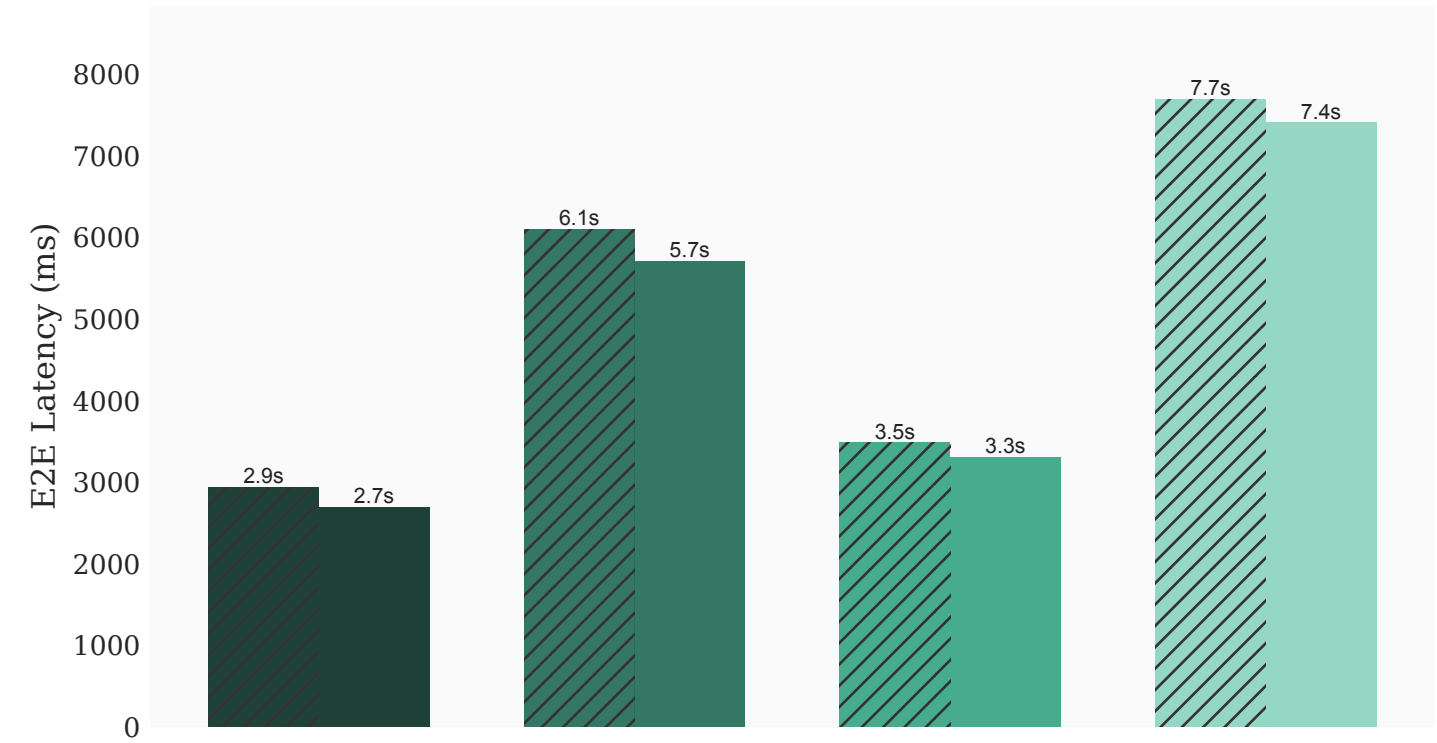


## End-to-End Latency (P99)



Legend: CC, No CC

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

# Numina Math (Rate 50)

## Time to First Token (Mean)



## End-to-End Latency (Mean)



## Time to First Token (P99)



## End-to-End Latency (P99)



CC   No CC

LLama 3.1 8B   Mistral 3.1 24B   GPT OSS 120B   LLama 3.3 70B Int4
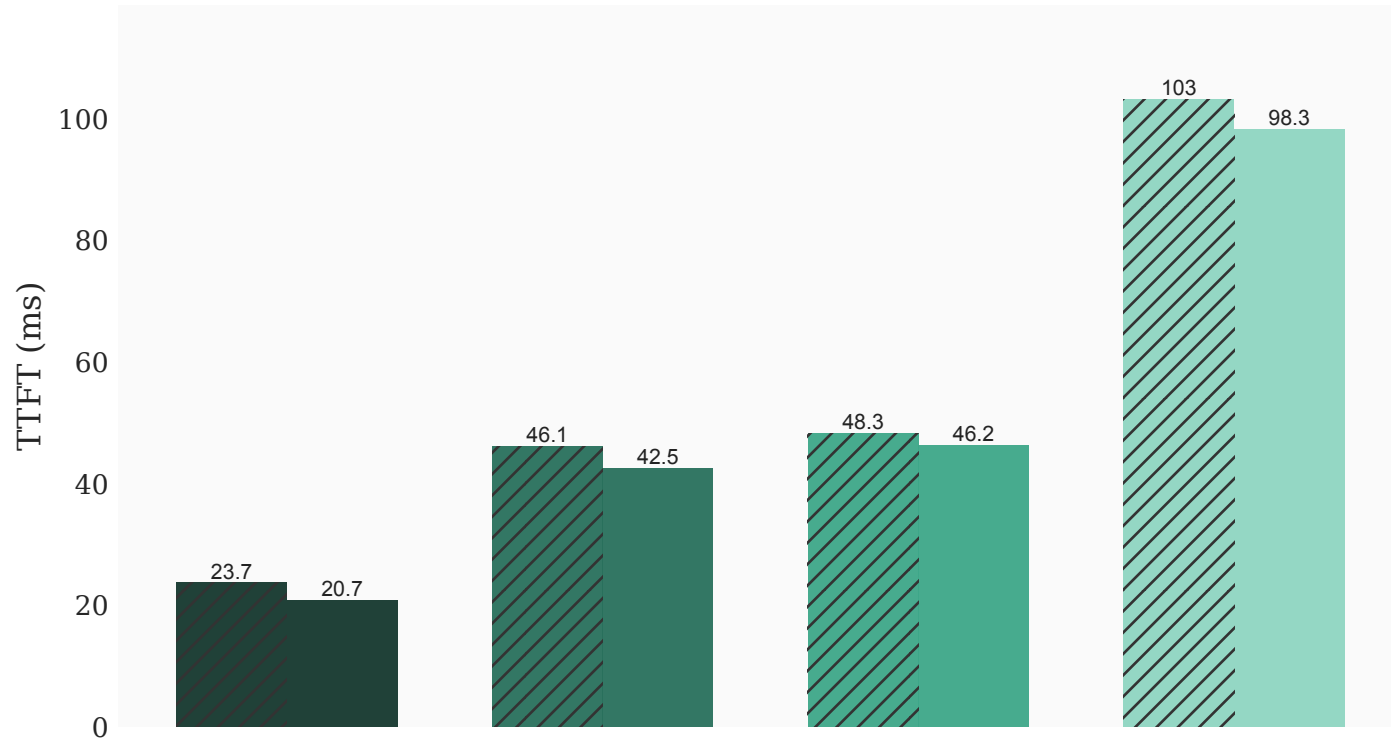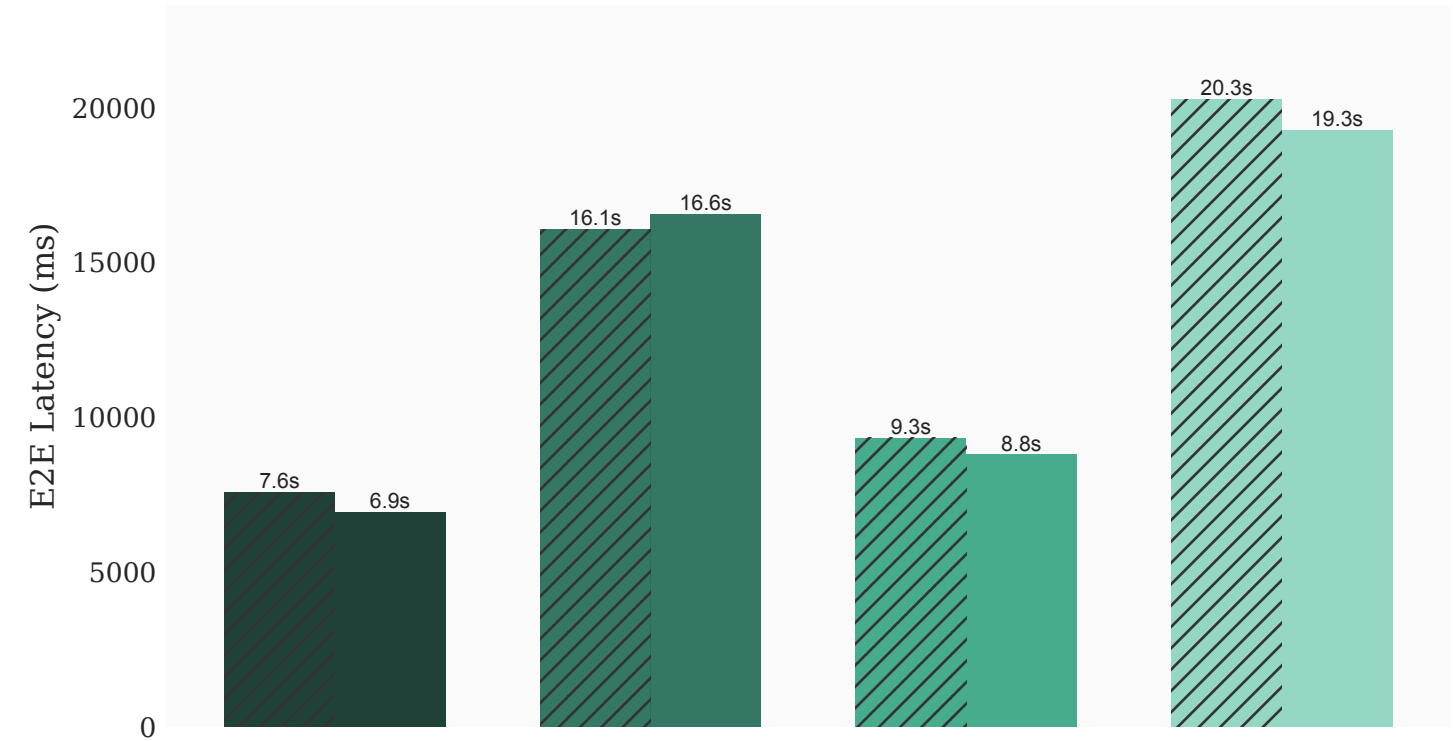
Numina Math (Rate 1)