# Random (1500 ⇒ 250) (100 Concurrent Requests)

## Time to First Token (Mean)



## Time to First Token (P99)

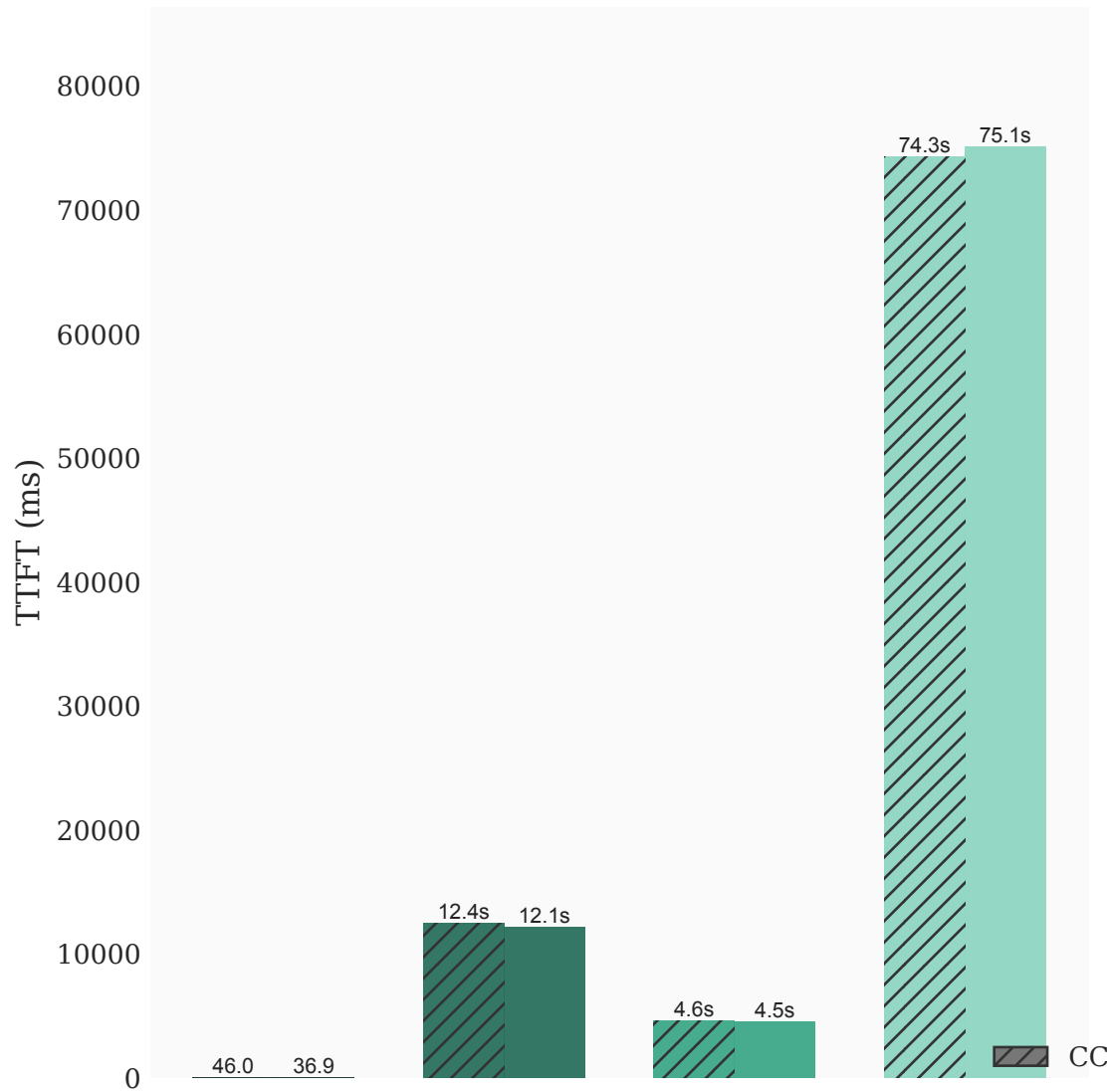Legend: CC (hatched) / No CC

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

Time to First Token (Mean):
- LLama 3.1 8B: 2.8s / 2.8s
- Mistral 3.1 24B: 13.5s / 13.1s
- GPT OSS 120B: 5.7s / 5.5s
- LLama 3.3 70B Int4: 75.3s / 75.9s

Time to First Token (P99):
- LLama 3.1 8B: 5.5s / 5.5s
- Mistral 3.1 24B: 28.3s / 27.7s
- GPT OSS 120B: 11.8s / 10.9s
- LLama 3.3 70B Int4: 152.5s / 155.6s

# Random (1500 ⇒ 250) (50 Concurrent Requests)

## Time to First Token (Mean)

TTFT (ms)

74.3s  75.1s

12.4s  12.1s

4.6s  4.5s

46.0  36.9

▨ CC  ▬ No CC

## Time to First Token (P99)

TTFT (ms)

151.8s  152.1s

26.2s  25.7s

9.4s  9.0s

76.8  60.5

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

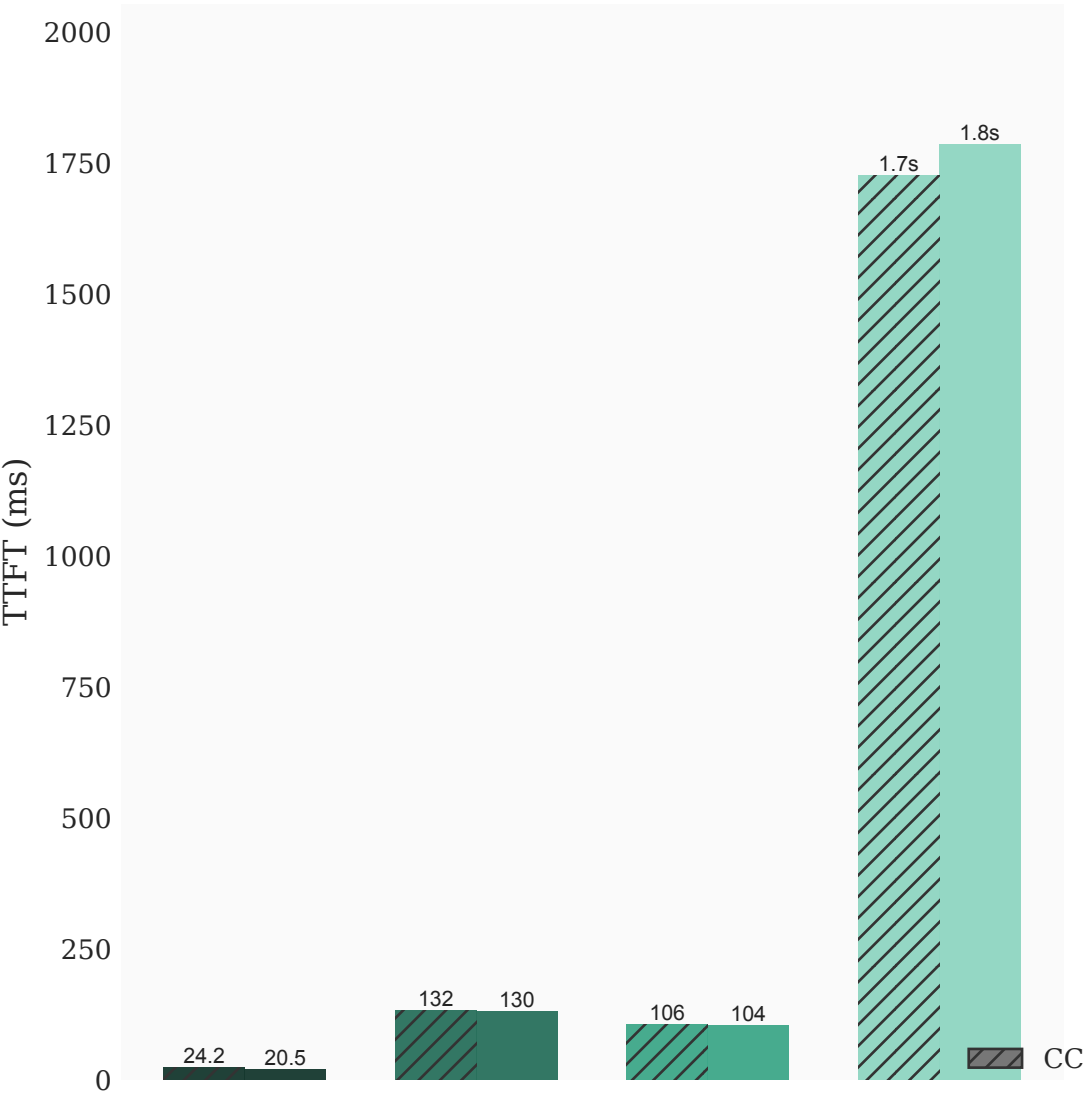# Random (1500 ⇒ 250) (1 Concurrent Requests)

## Time to First Token (Mean)

TTFT (ms)

- 24.2 / 20.5
- 132 / 130
- 106 / 104
- 1.7s / 1.8s

CC    No CC

## Time to First Token (P99)

TTFT (ms)

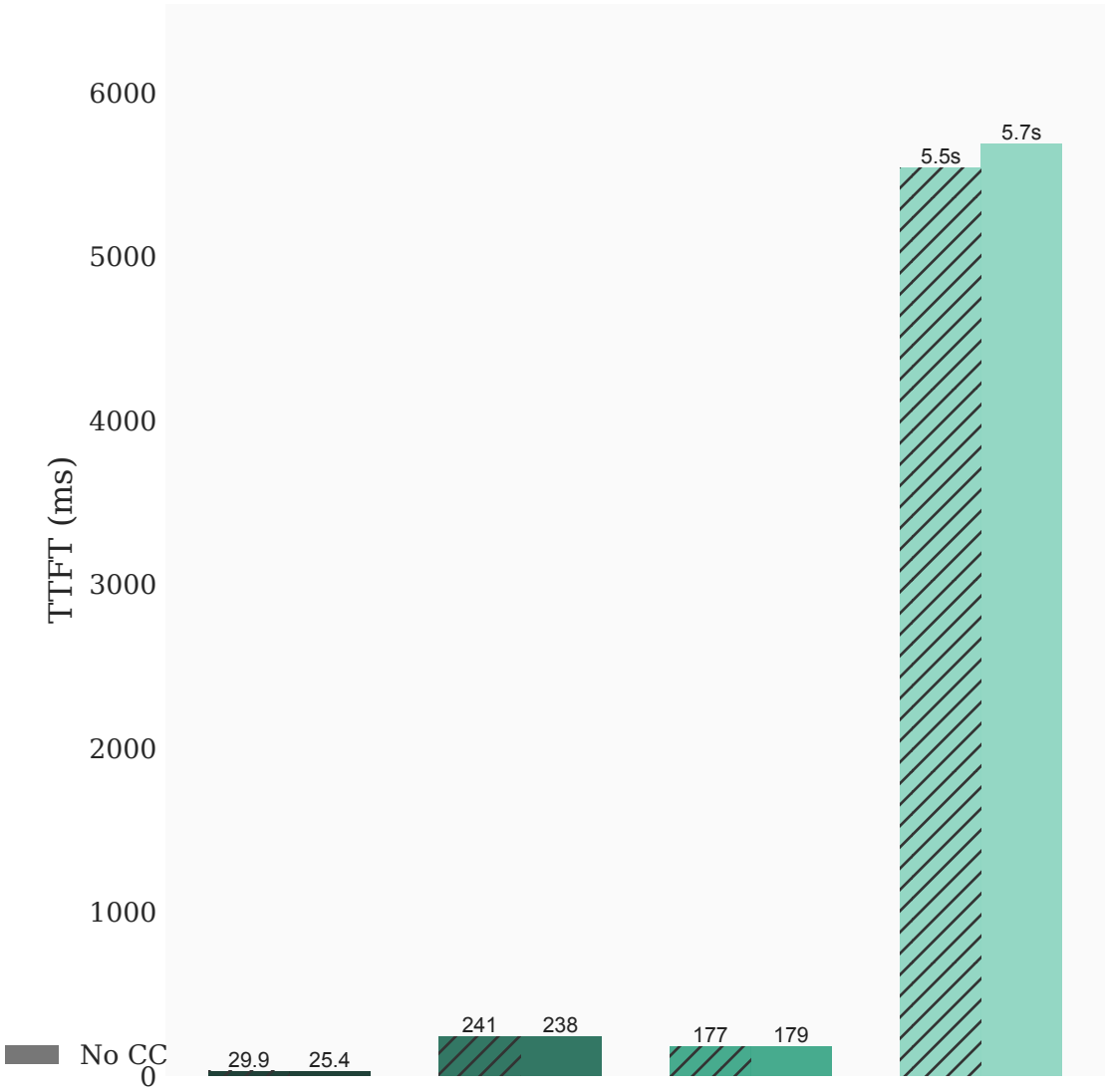- 29.9 / 25.4
- 241 / 238
- 177 / 179
- 5.5s / 5.7s

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

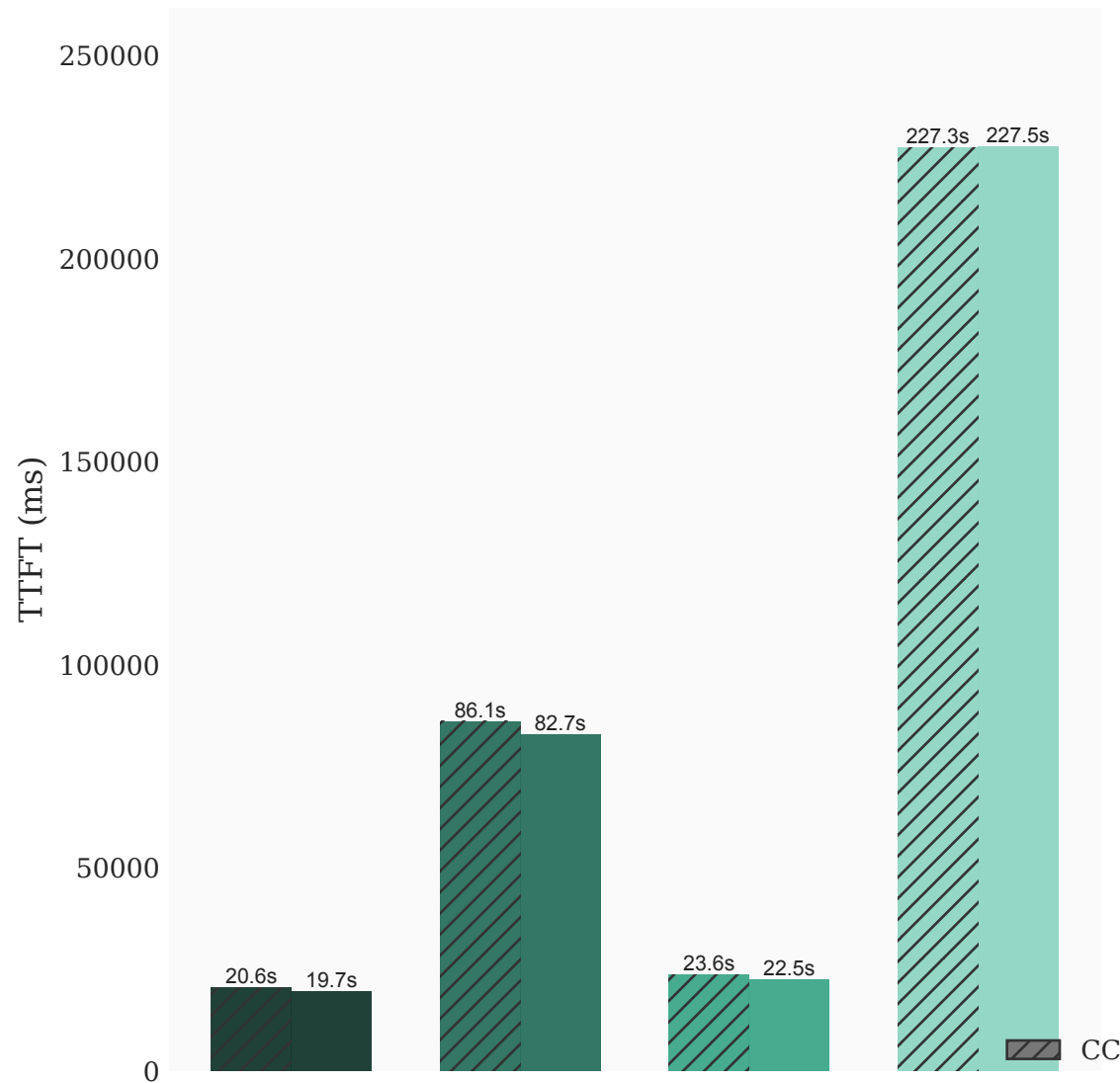# Random (4000 ⇒ 1000) (100 Concurrent Requests)

## Time to First Token (Mean)



TTFT (ms)

- 20.6s / 19.7s
- 86.1s / 82.7s
- 23.6s / 22.5s
- 227.3s / 227.5s

## Time to First Token (P99)



TTFT (ms)

- 47.5s / 45.5s
- 179.9s / 175.2s
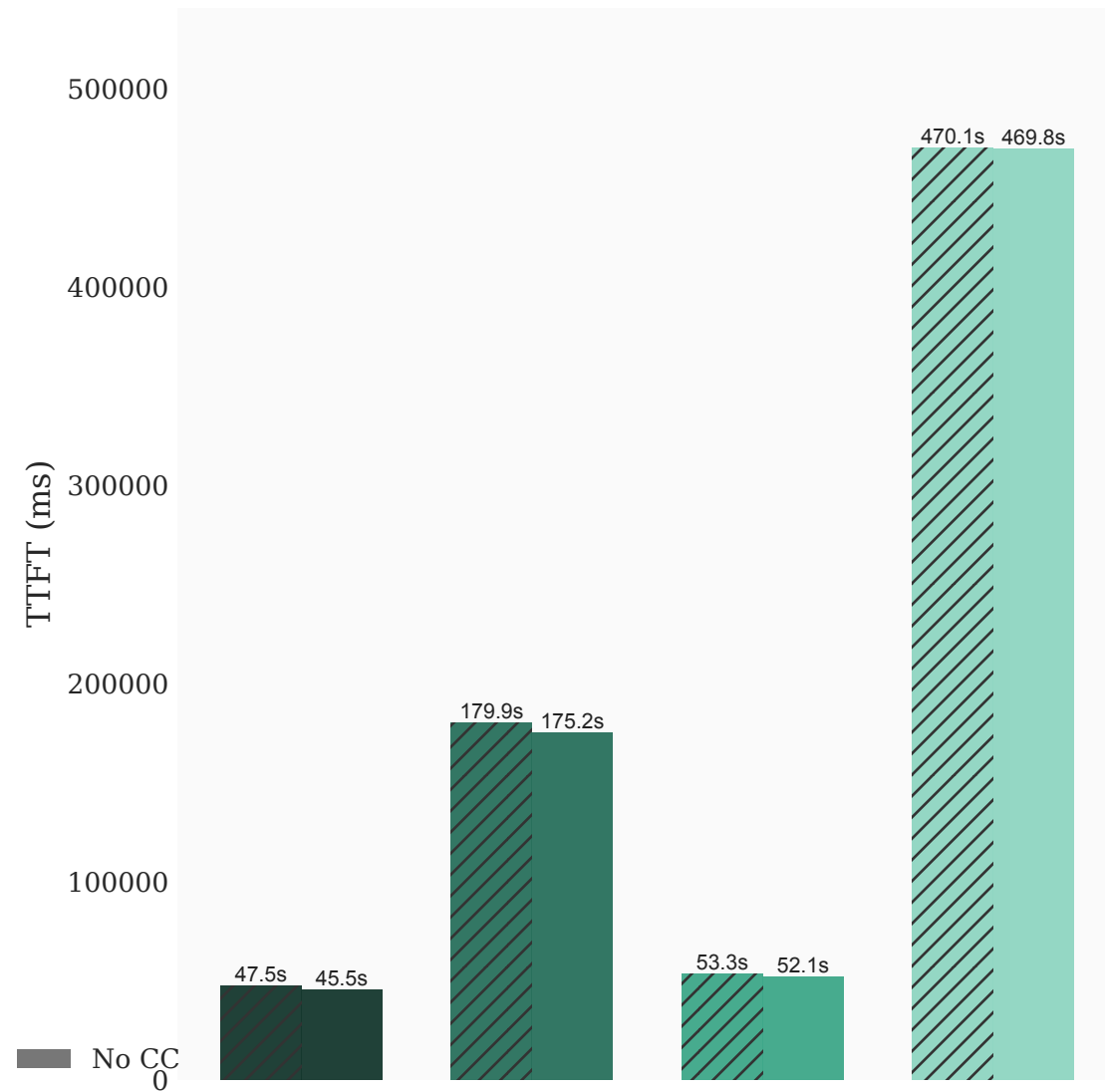- 53.3s / 52.1s
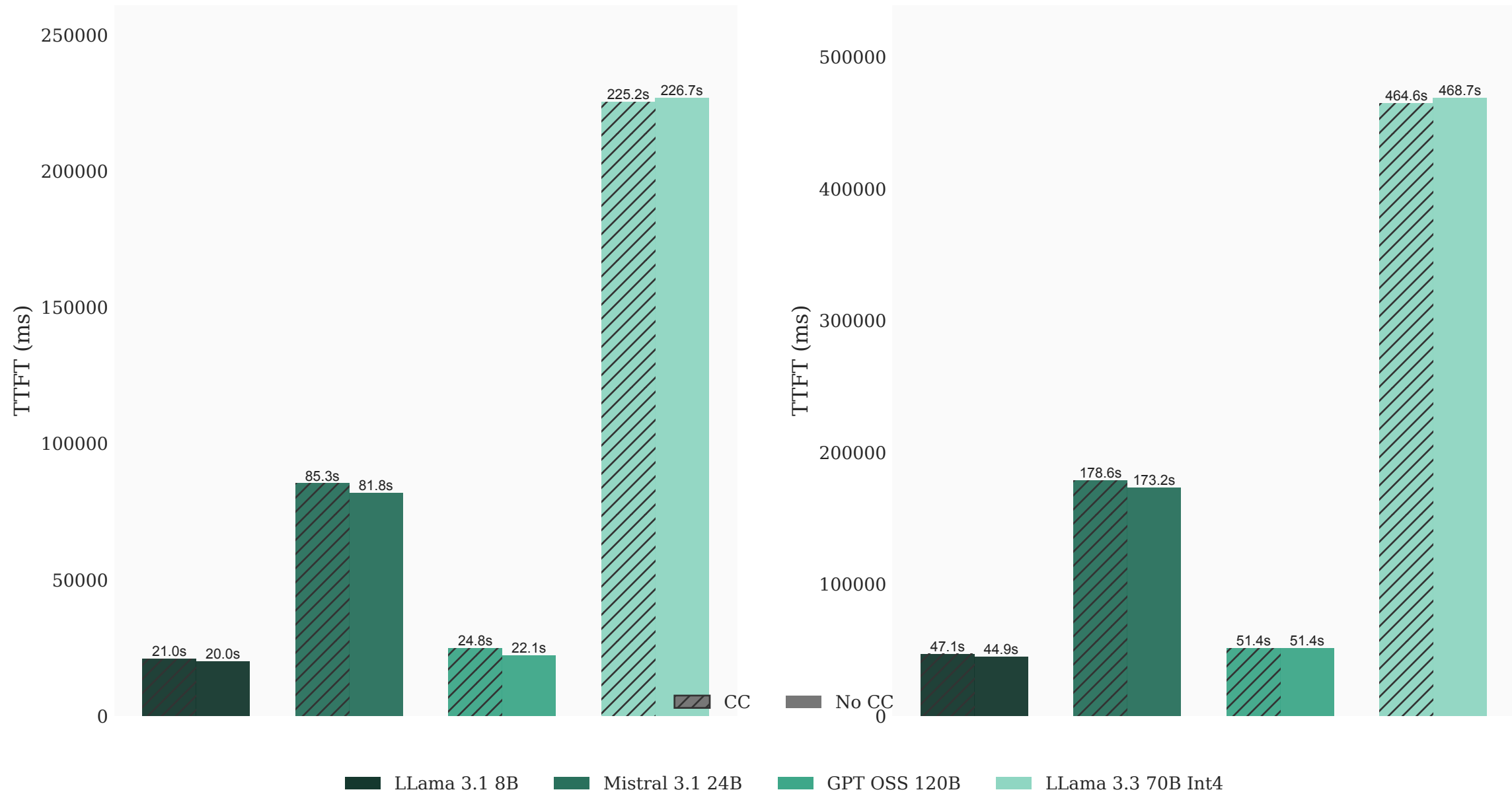- 470.1s / 469.8s

Legend: ▨ CC   ▬ No CC

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

# Random (4000 ⇒ 1000) (50 Concurrent Requests)

## Time to First Token (Mean)



TTFT (ms)

- 21.0s / 20.0s (LLama 3.1 8B)
- 85.3s / 81.8s (Mistral 3.1 24B)
- 24.8s / 22.1s (GPT OSS 120B)
- 225.2s / 226.7s (LLama 3.3 70B Int4)

CC / No CC

## Time to First Token (P99)



TTFT (ms)

- 47.1s / 44.9s (LLama 3.1 8B)
- 178.6s / 173.2s (Mistral 3.1 24B)
- 51.4s / 51.4s (GPT OSS 120B)
- 464.6s / 468.7s (LLama 3.3 70B Int4)

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

# Random (4000 ⇒ 1000) (1 Concurrent Requests)

## Time to First Token (Mean)
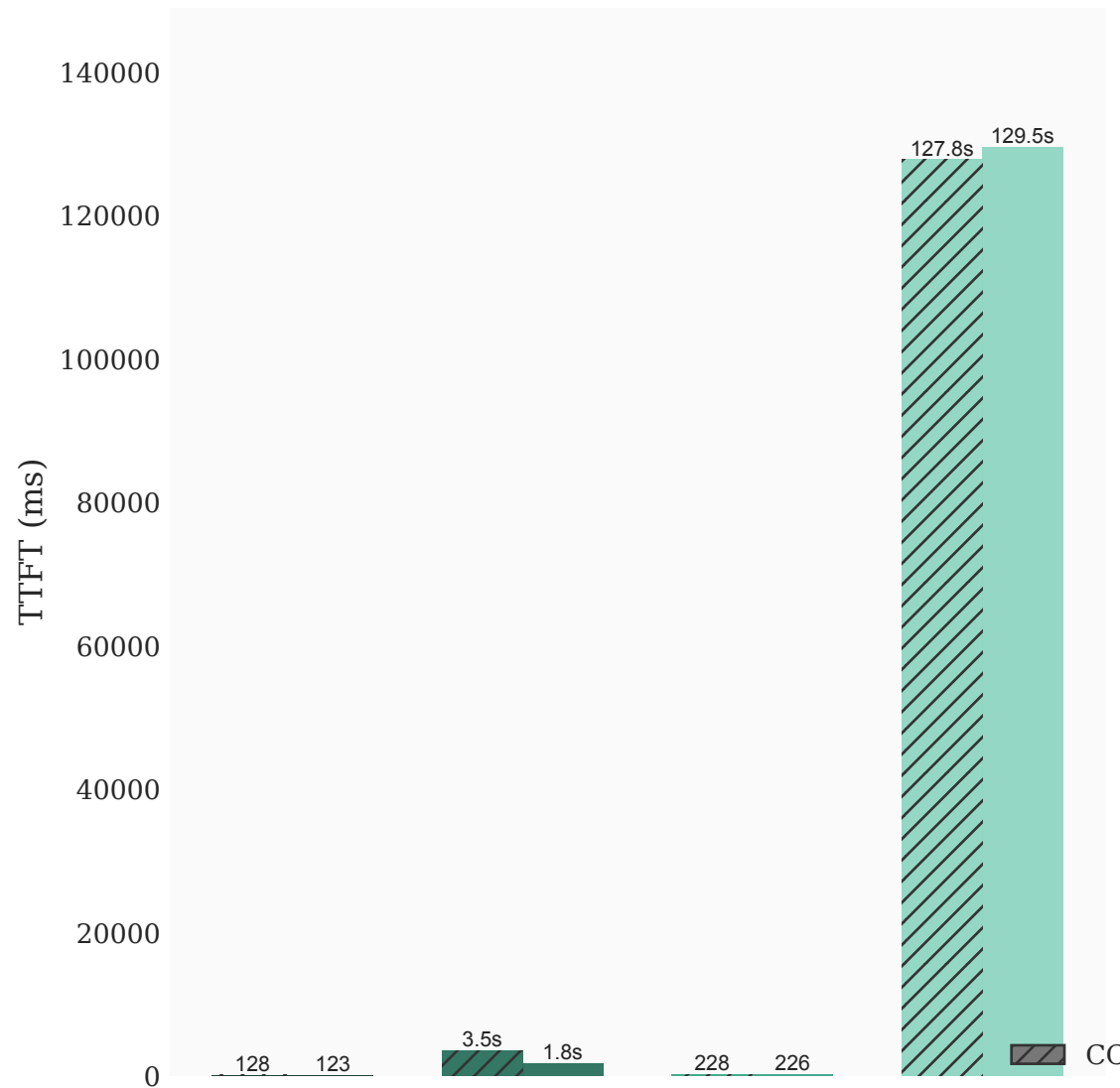


TTFT (ms)

- 140000
- 120000
- 100000
- 80000
- 60000
- 40000
- 20000
- 0

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

## Time to First Token (P99)



TTFT (ms)

- 300000
- 250000
- 200000
- 150000
- 100000
- 50000
- 0

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

CC   No CC

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

# Random (1000 ⇒ 1000) (100 Concurrent Requests)

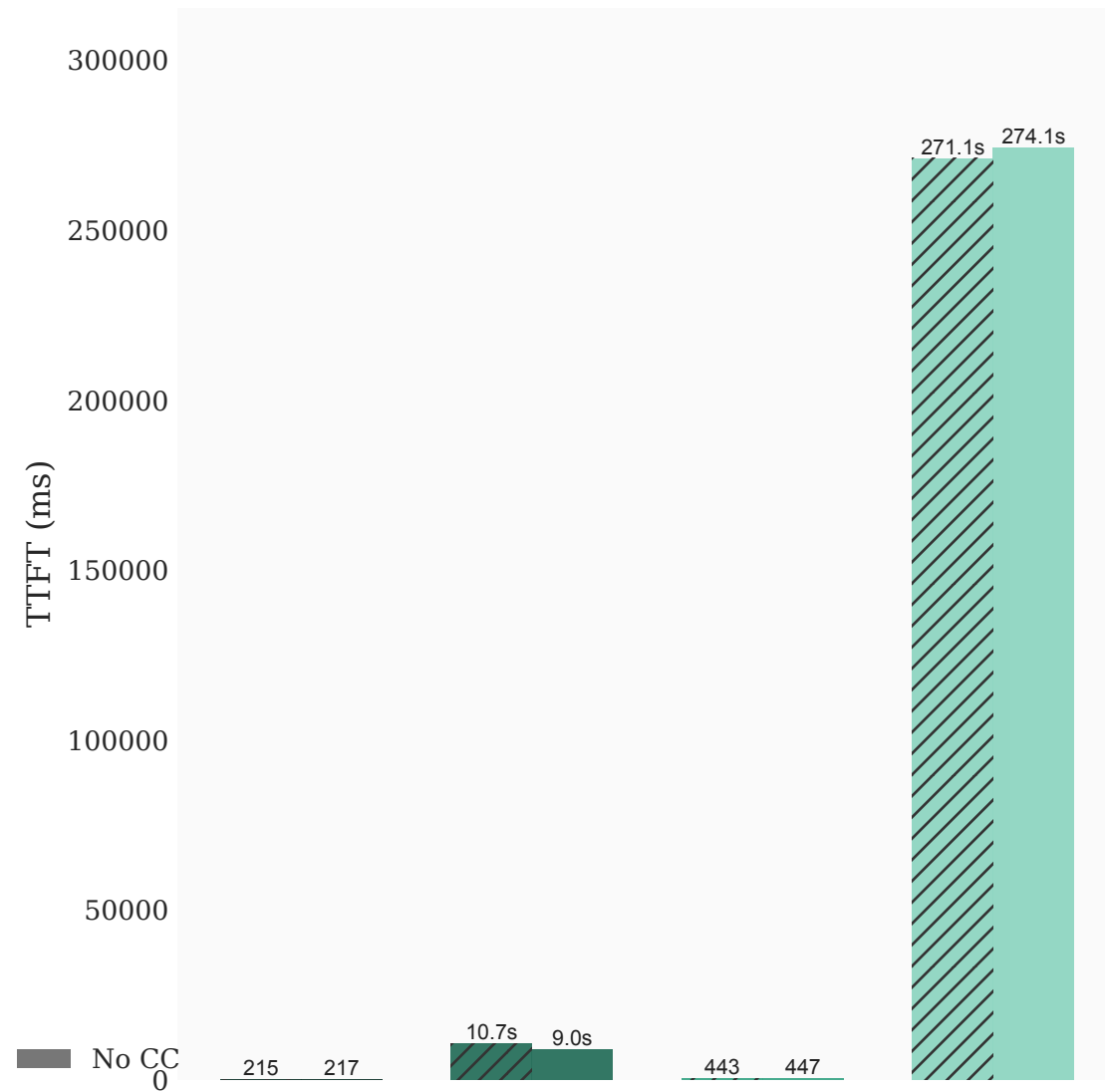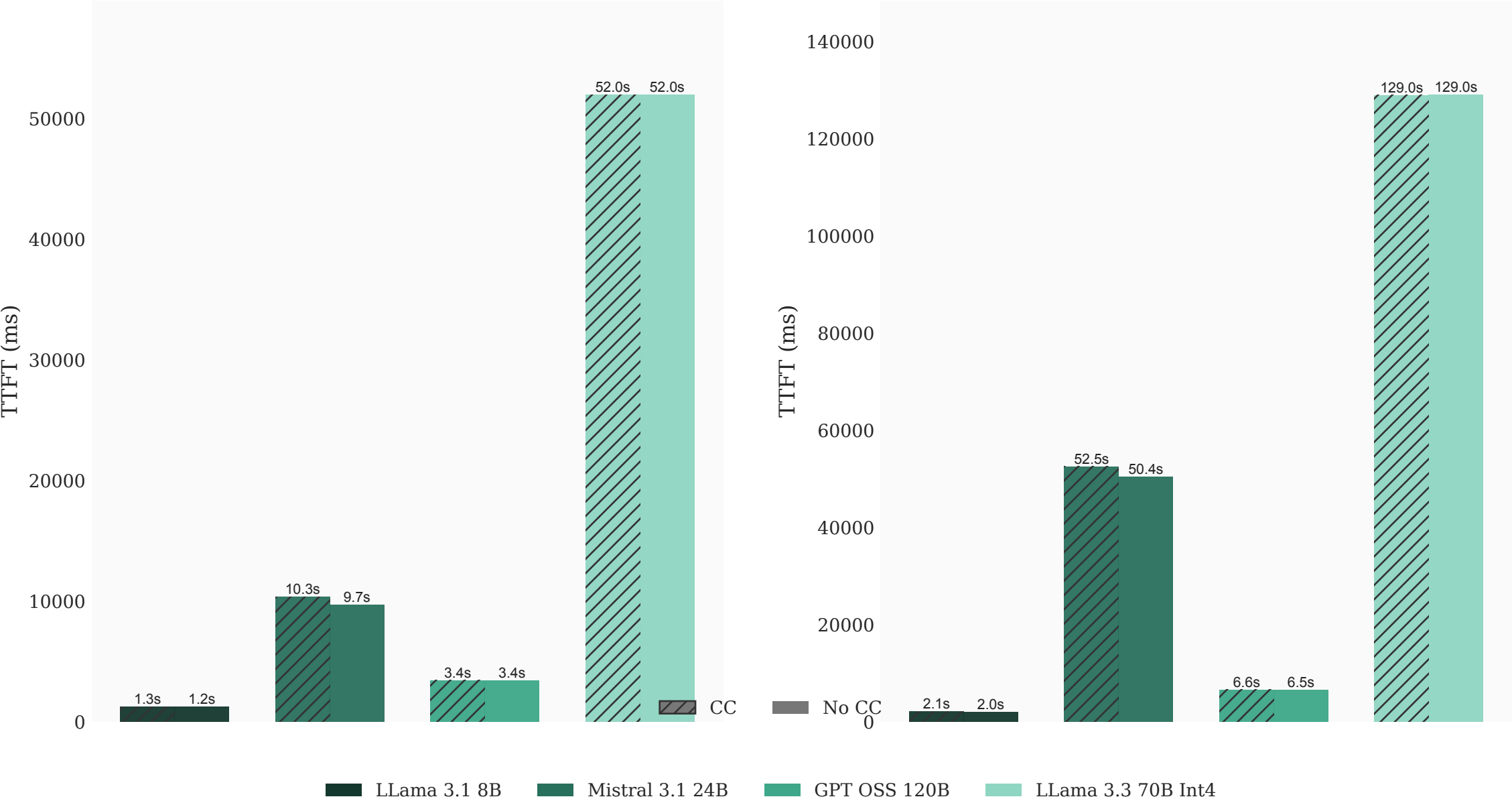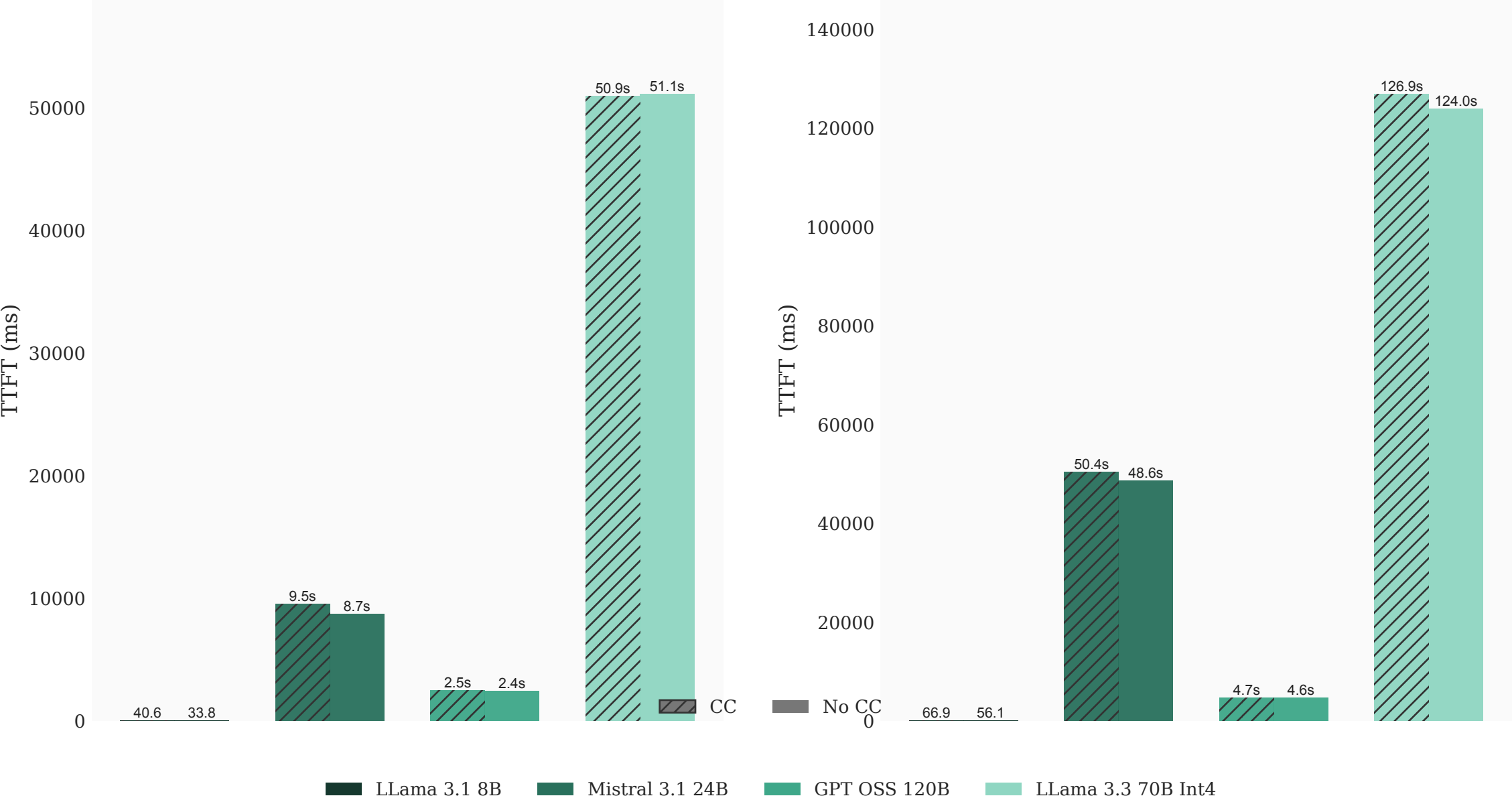## Time to First Token (Mean)



## Time to First Token (P99)

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

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

Time to First Token (Mean) values:
- LLama 3.1 8B: 1.3s / 1.2s
- Mistral 3.1 24B: 10.3s / 9.7s
- GPT OSS 120B: 3.4s / 3.4s
- LLama 3.3 70B Int4: 52.0s / 52.0s

Time to First Token (P99) values:
- LLama 3.1 8B: 2.1s / 2.0s
- Mistral 3.1 24B: 52.5s / 50.4s
- GPT OSS 120B: 6.6s / 6.5s
- LLama 3.3 70B Int4: 129.0s / 129.0s

# Random (1000 ⇒ 1000) (50 Concurrent Requests)

## Time to First Token (Mean)

## Time to First Token (P99)



TTFT (ms)

Mean values: 40.6, 33.8, 9.5s, 8.7s, 2.5s, 2.4s, 50.9s, 51.1s

P99 values: 66.9, 56.1, 50.4s, 48.6s, 4.7s, 4.6s, 126.9s, 124.0s

Legend: ▨ CC    ▬ No CC

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

# Random (1000 ⇒ 1000) (1 Concurrent Requests)

## Time to First Token (Mean)



## Time to First Token (P99)

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

CC, No CC

Mean values: 25.2, 21.1, 97.9, 94.1, 83.8, 81.9, 747, 761

P99 values: 43.8, 39.7, 157, 161, 137, 133, 2.2s, 2.6s

# ShareGPT (100 Concurrent Requests)

## Time to First Token (Mean)

TTFT (ms)

- 51.1 / 43.2 — LLama 3.1 8B
- 956 / 897 — Mistral 3.1 24B
- 340 / 319 — GPT OSS 120B
- 10.2s / 10.3s — LLama 3.3 70B Int4

## Time to First Token (P99)

TTFT (ms)

- 92.7 / 88.2 — LLama 3.1 8B
- 1.5s / 1.5s — Mistral 3.1 24B
- 578 / 538 — GPT OSS 120B
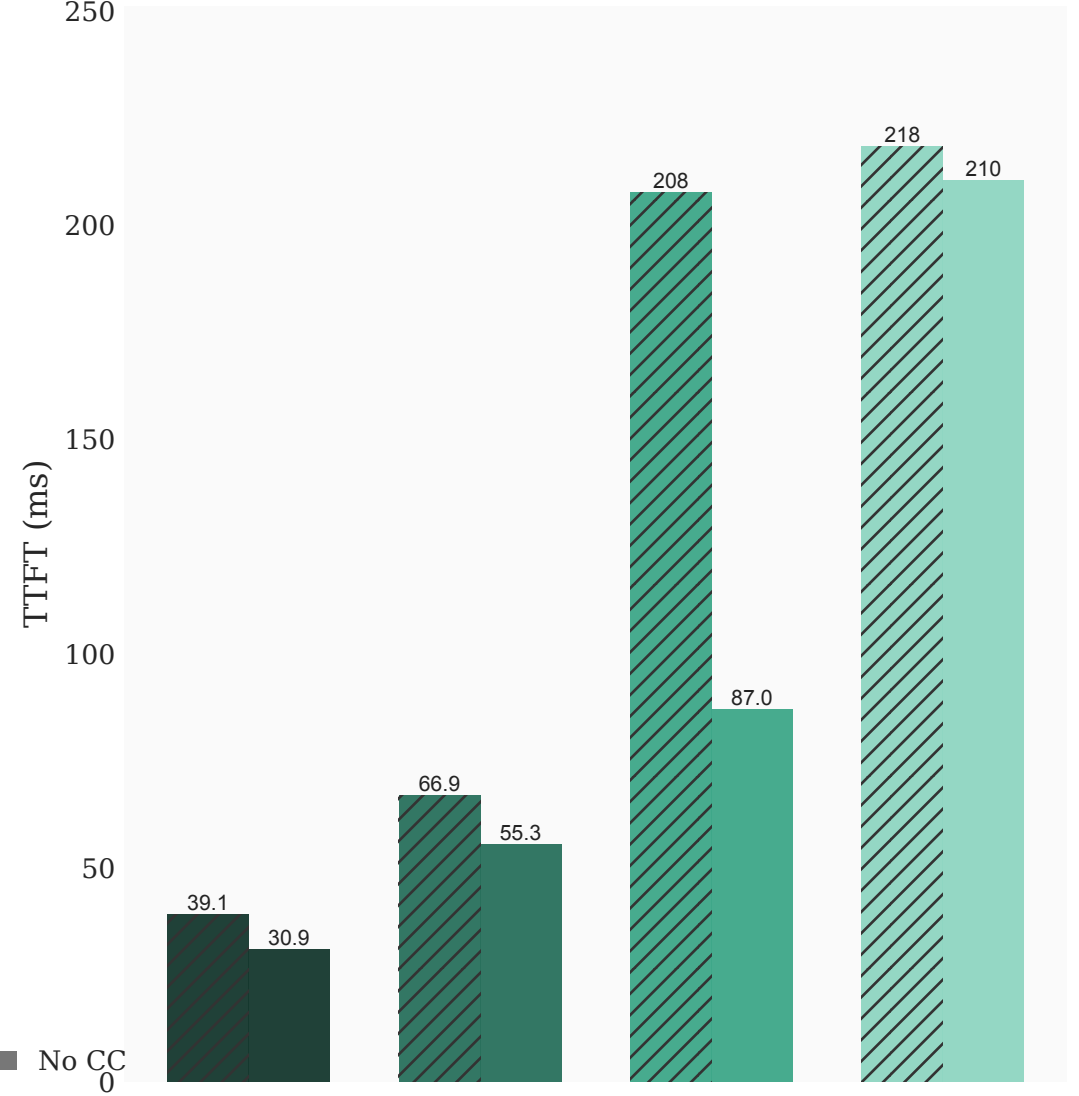- 17.0s / 17.1s — LLama 3.3 70B Int4

Legend: CC | No CC

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

# ShareGPT (50 Concurrent Requests)

## Time to First Token (Mean)



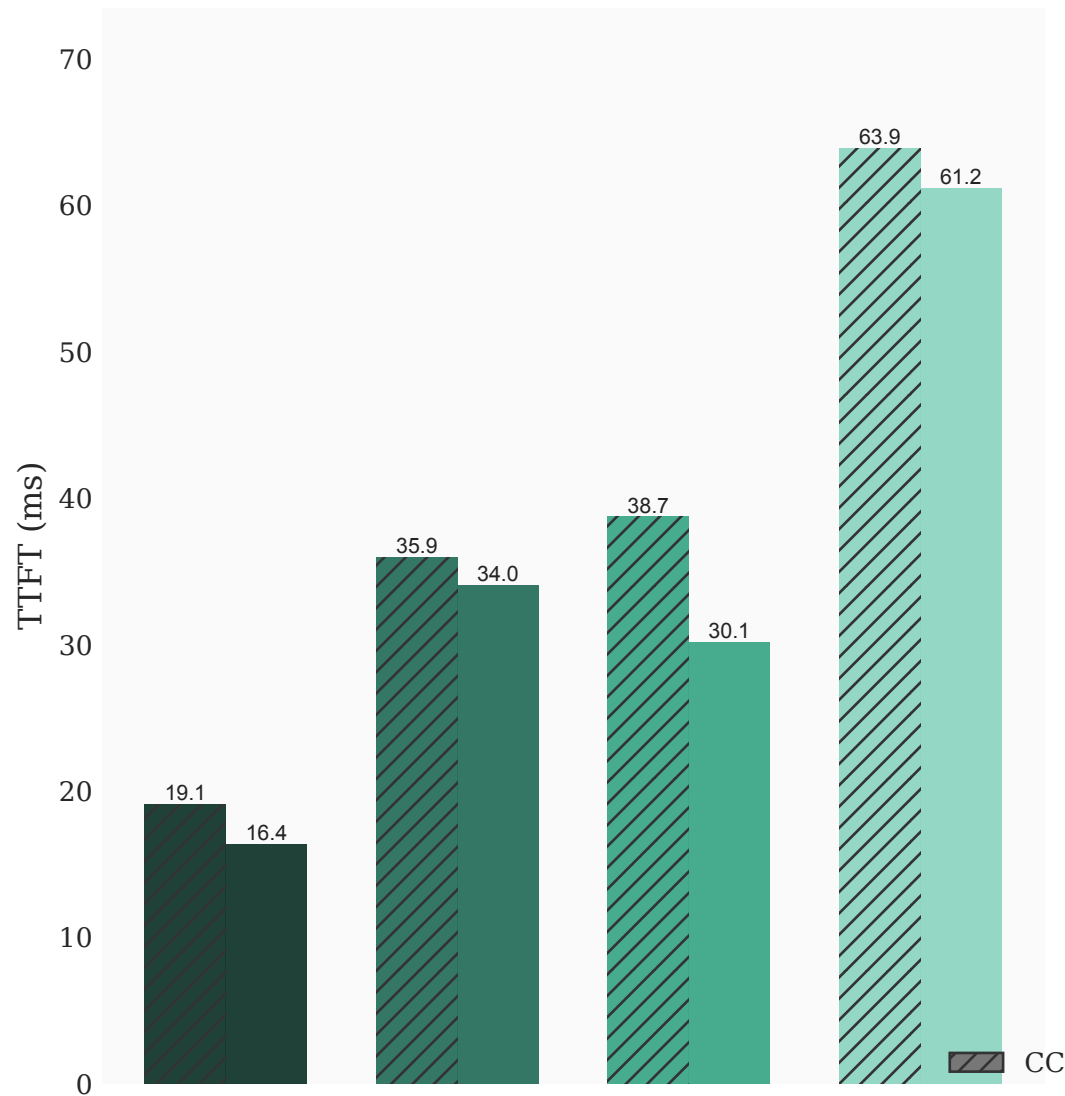## Time to First Token (P99)

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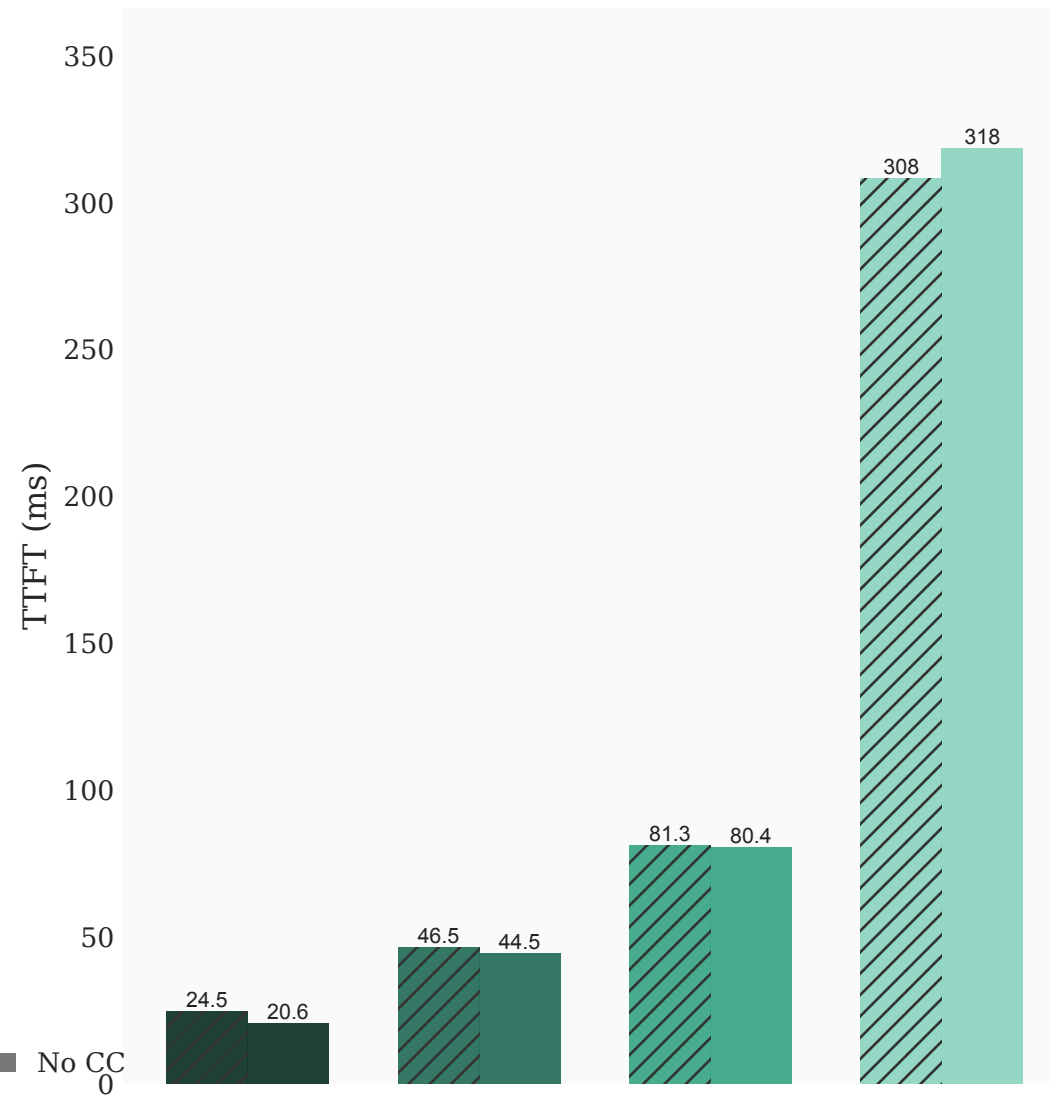
LLama 3.1 8B  ▬  Mistral 3.1 24B  ▬  GPT OSS 120B  ▬  LLama 3.3 70B Int4

# ShareGPT (1 Concurrent Requests)
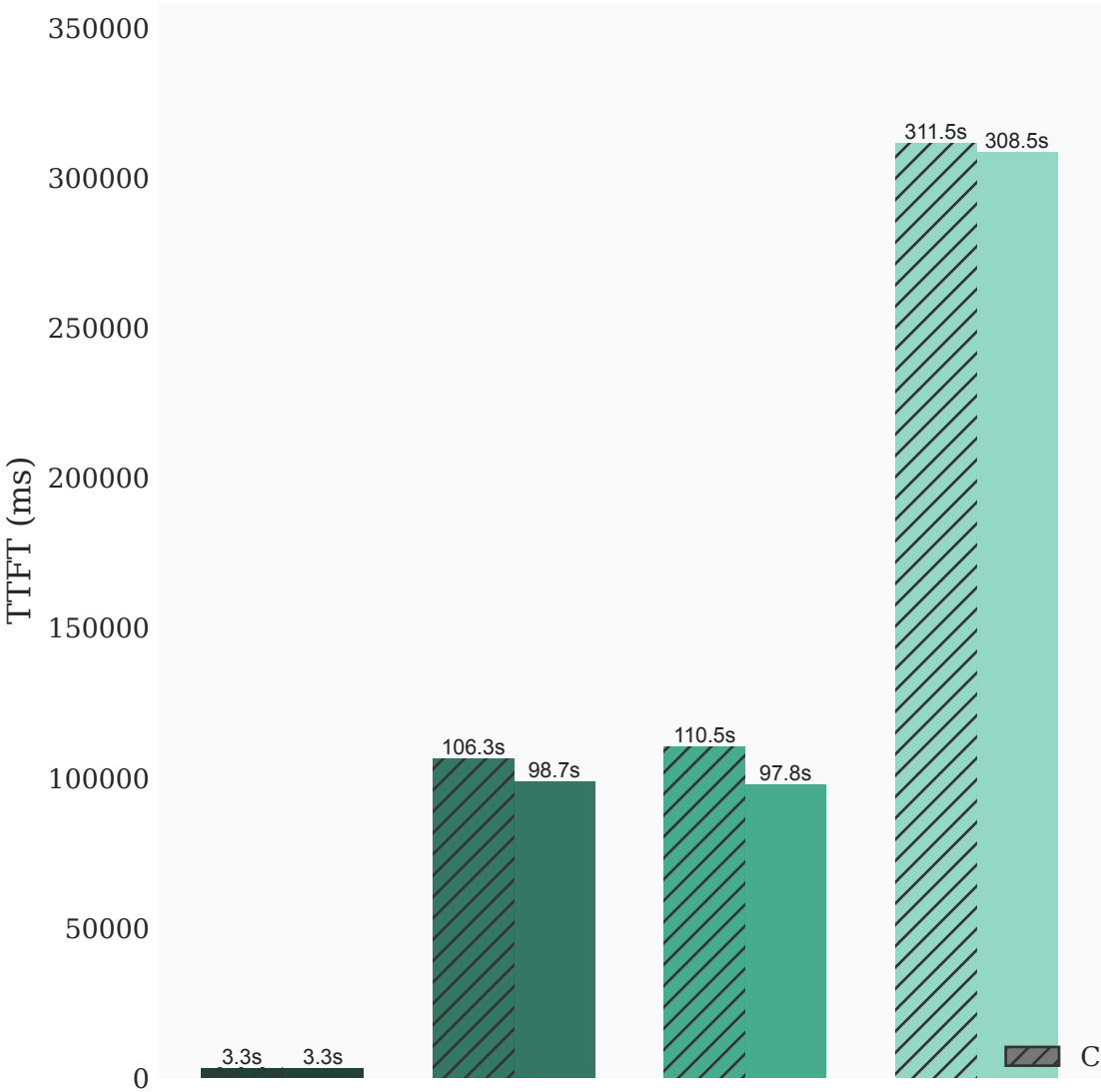
## Time to First Token (Mean)



## Time to First Token (P99)

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

# Edit 10K Characters (100 Concurrent Requests)

## Time to First Token (Mean)

TTFT (ms)

- 3.3s 3.3s
- 106.3s 98.7s
- 110.5s 97.8s
- 311.5s 308.5s

CC  No CC

## Time to First Token (P99)

TTFT (ms)

- 4.9s 4.9s
- 281.6s 267.8s
- 266.5s 261.7s
- 674.1s 671.9s
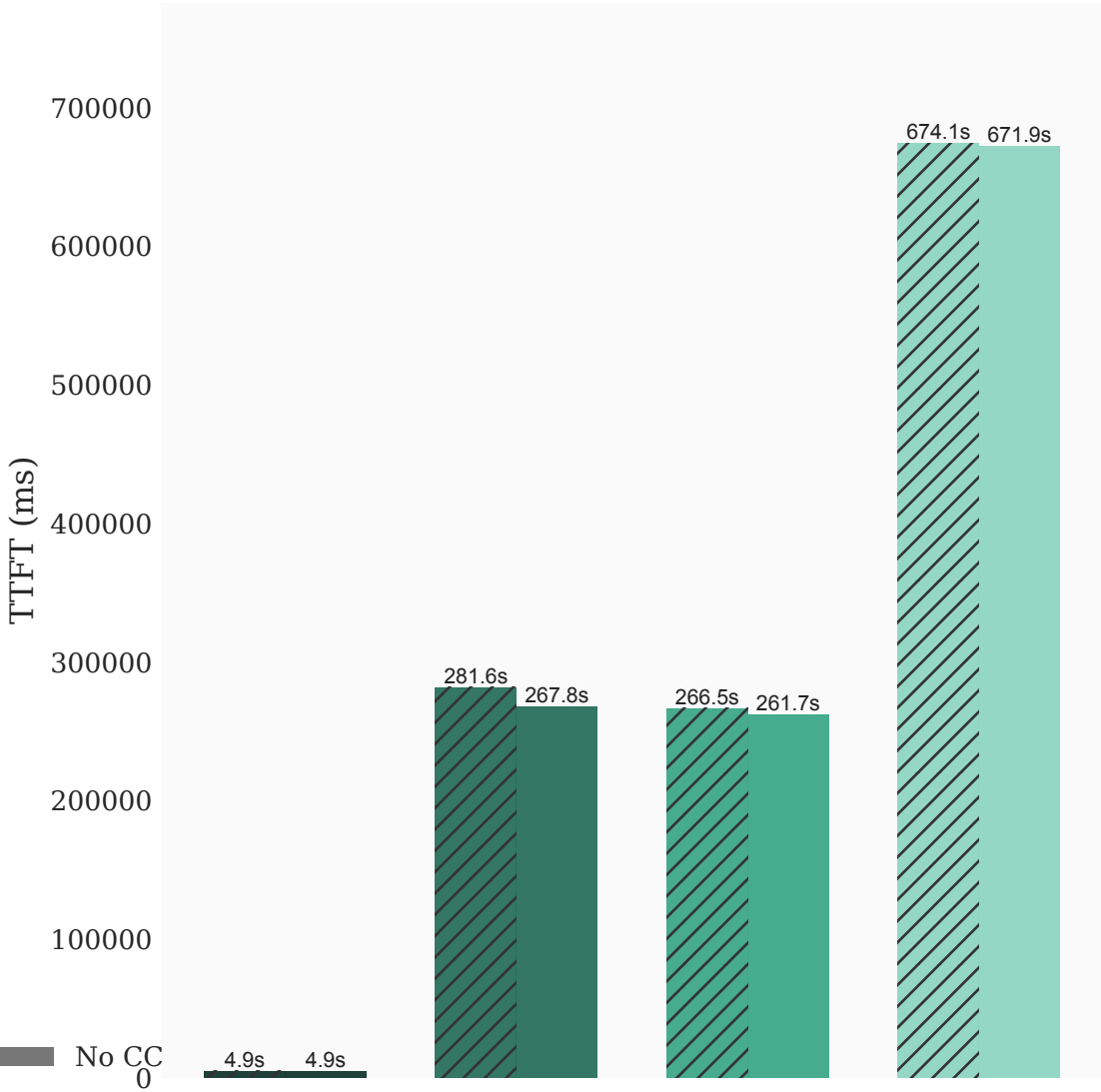
LLama 3.1 8B   Mistral 3.1 24B   GPT OSS 120B   LLama 3.3 70B Int4

# Edit 10K Characters (50 Concurrent Requests)

## Time to First Token (Mean)

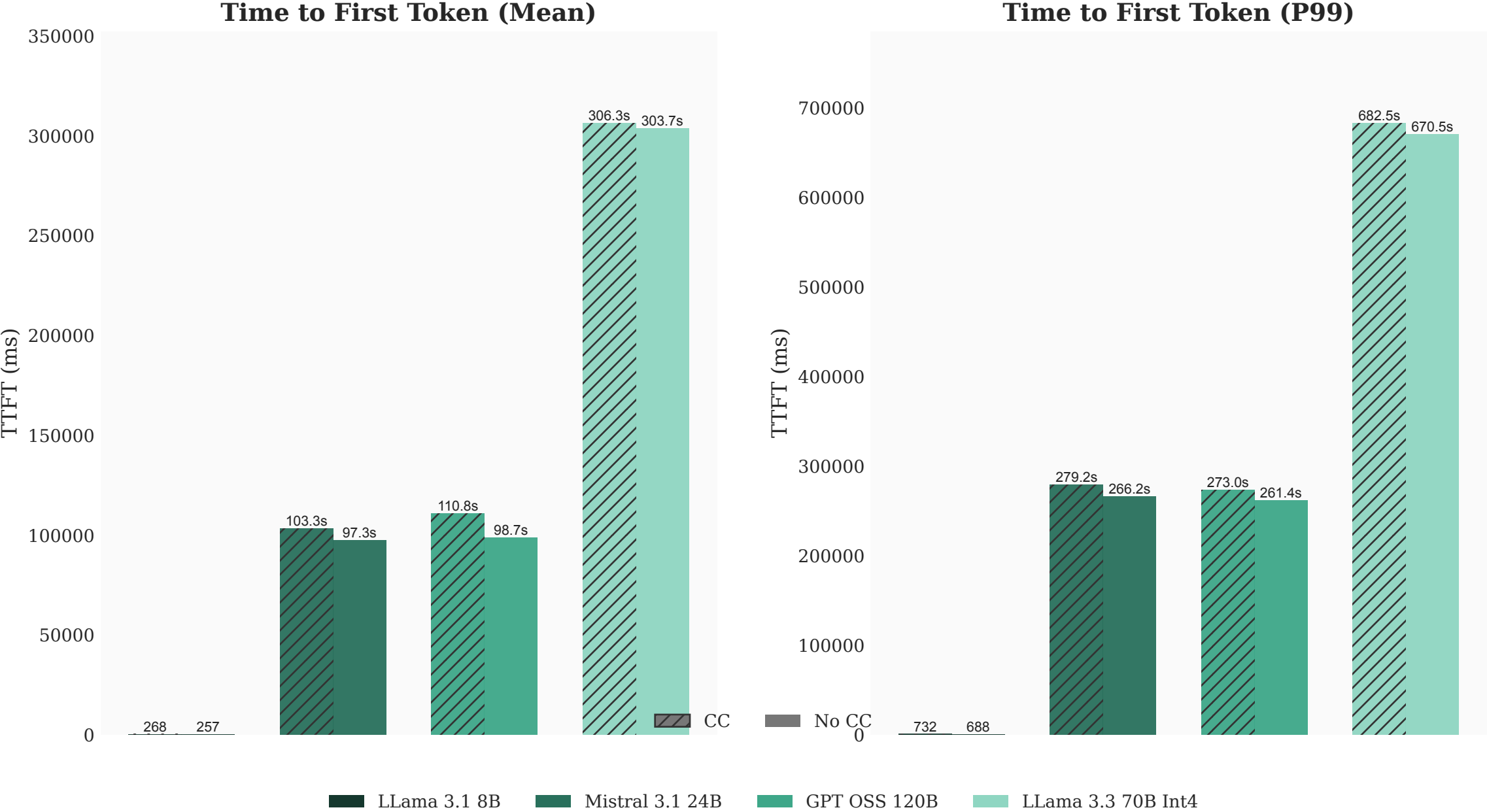TTFT (ms)

- 350000
- 300000
- 250000
- 200000
- 150000
- 100000
- 50000
- 0

268 | 257
103.3s | 97.3s
110.8s | 98.7s
306.3s | 303.7s

CC | No CC

## Time to First Token (P99)

TTFT (ms)

- 700000
- 600000
- 500000
- 400000
- 300000
- 200000
- 100000
- 0

732 | 688
279.2s | 266.2s
273.0s | 261.4s
682.5s | 670.5s

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

# Edit 10K Characters (1 Concurrent Requests)

## Time to First Token (Mean)

TTFT (ms)

- 250000
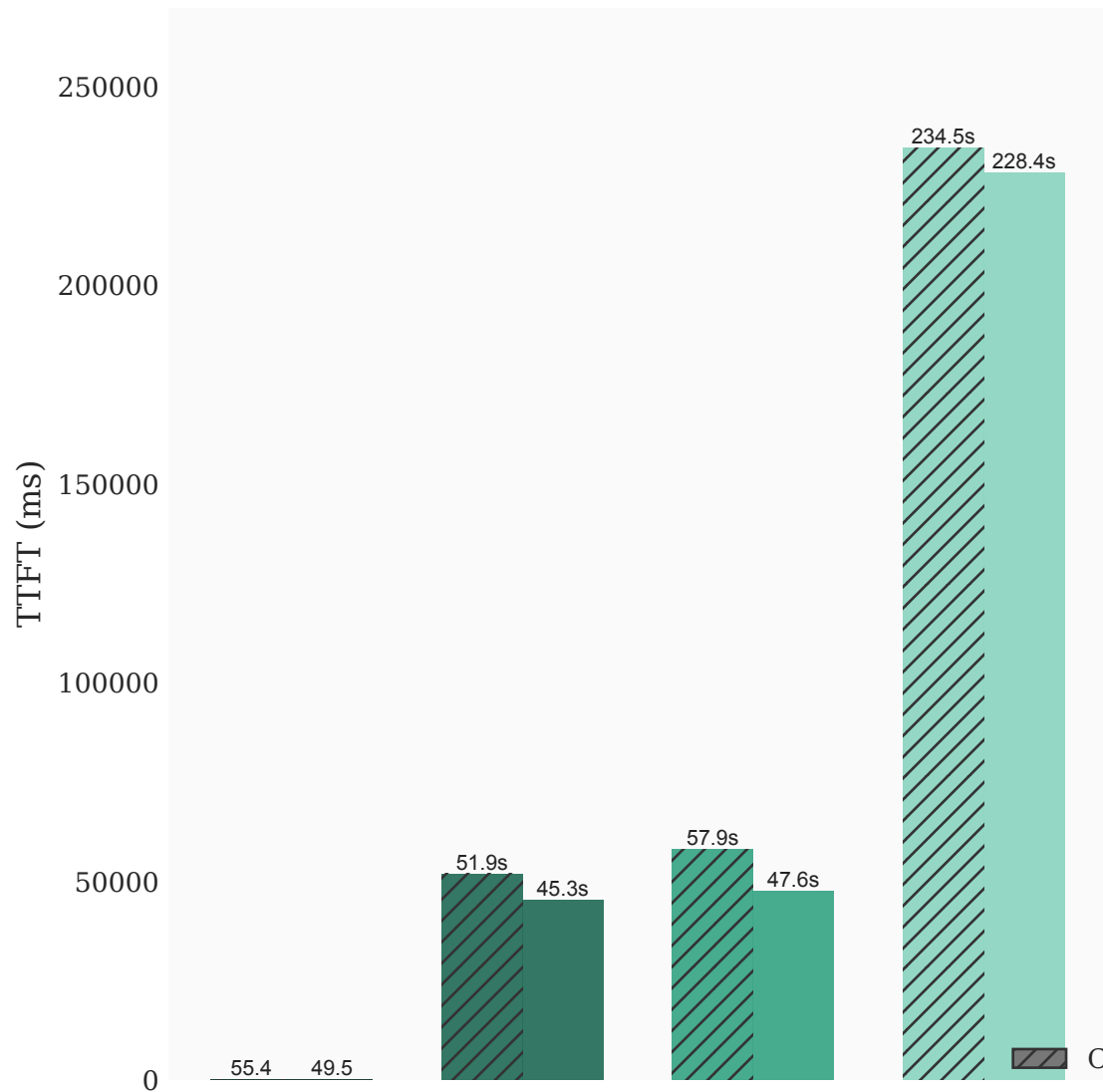- 200000
- 150000
- 100000
- 50000
- 0

55.4 49.5
51.9s 45.3s
57.9s 47.6s
234.5s 228.4s

CC    No CC

## Time to First Token (P99)

TTFT (ms)

- 500000
- 400000
- 300000
- 200000
- 100000
- 0

108 102
117.3s 100.1s
137.6s 106.6s
497.6s 475.0s
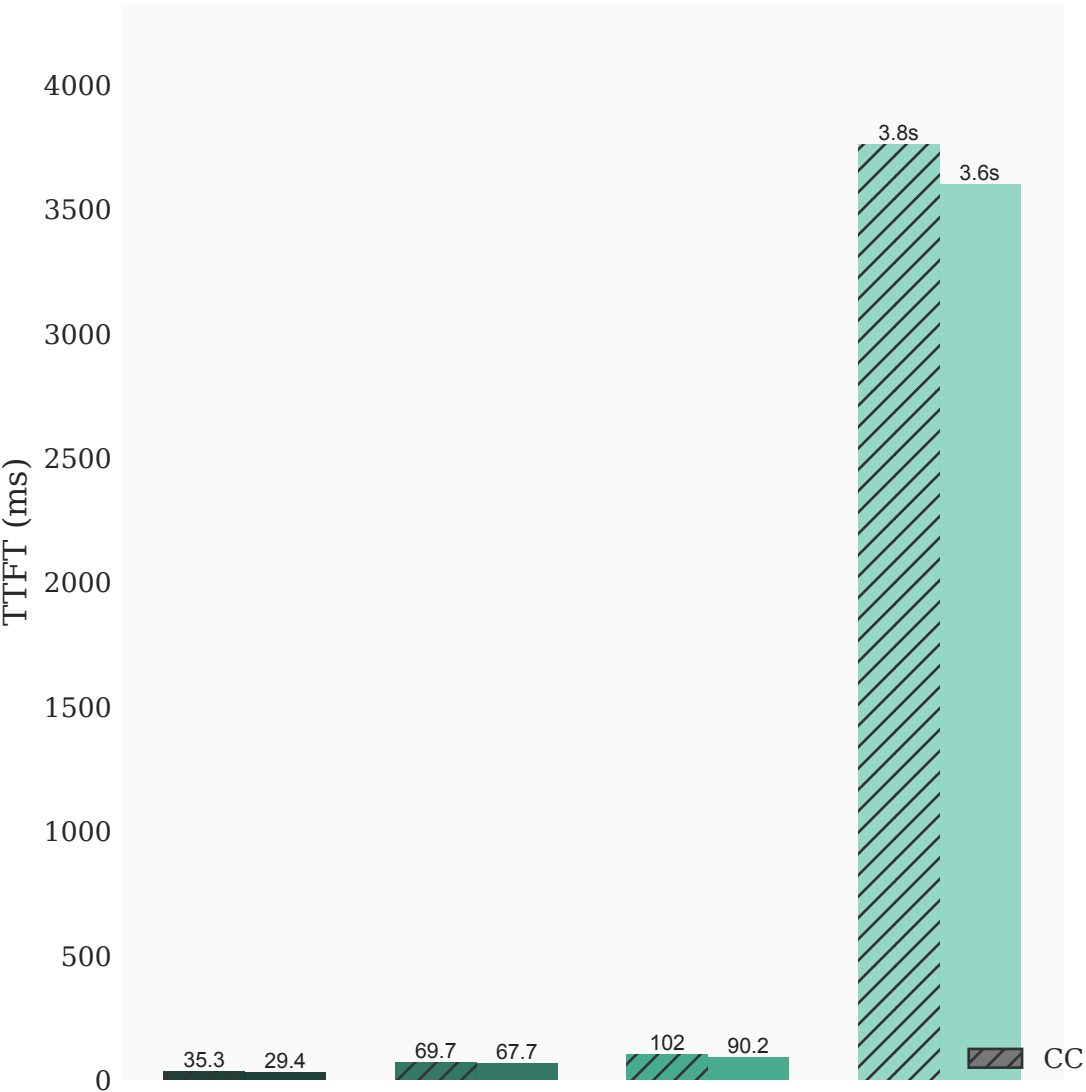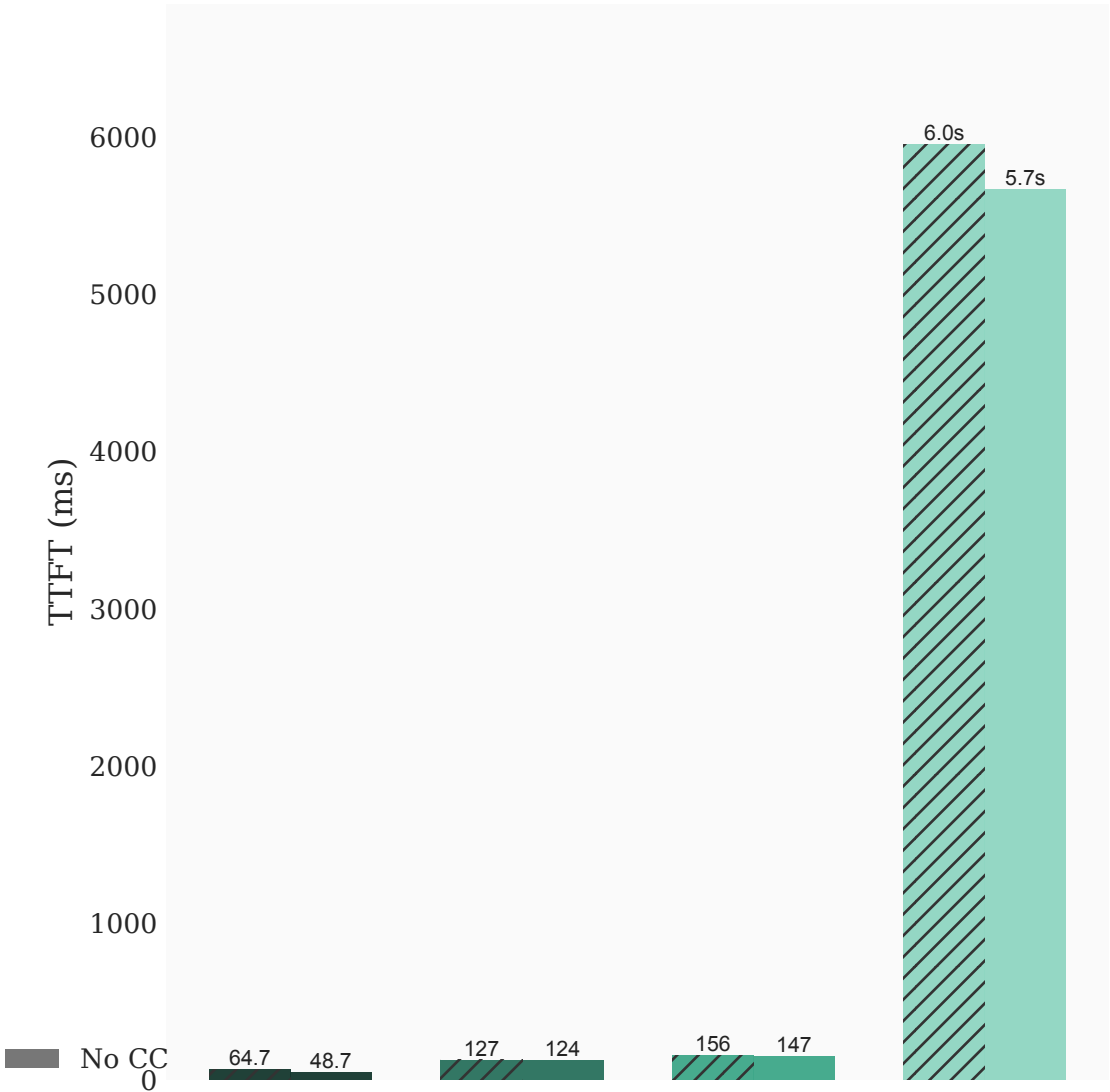
■ LLama 3.1 8B    ■ Mistral 3.1 24B    ■ GPT OSS 120B    ■ LLama 3.3 70B Int4

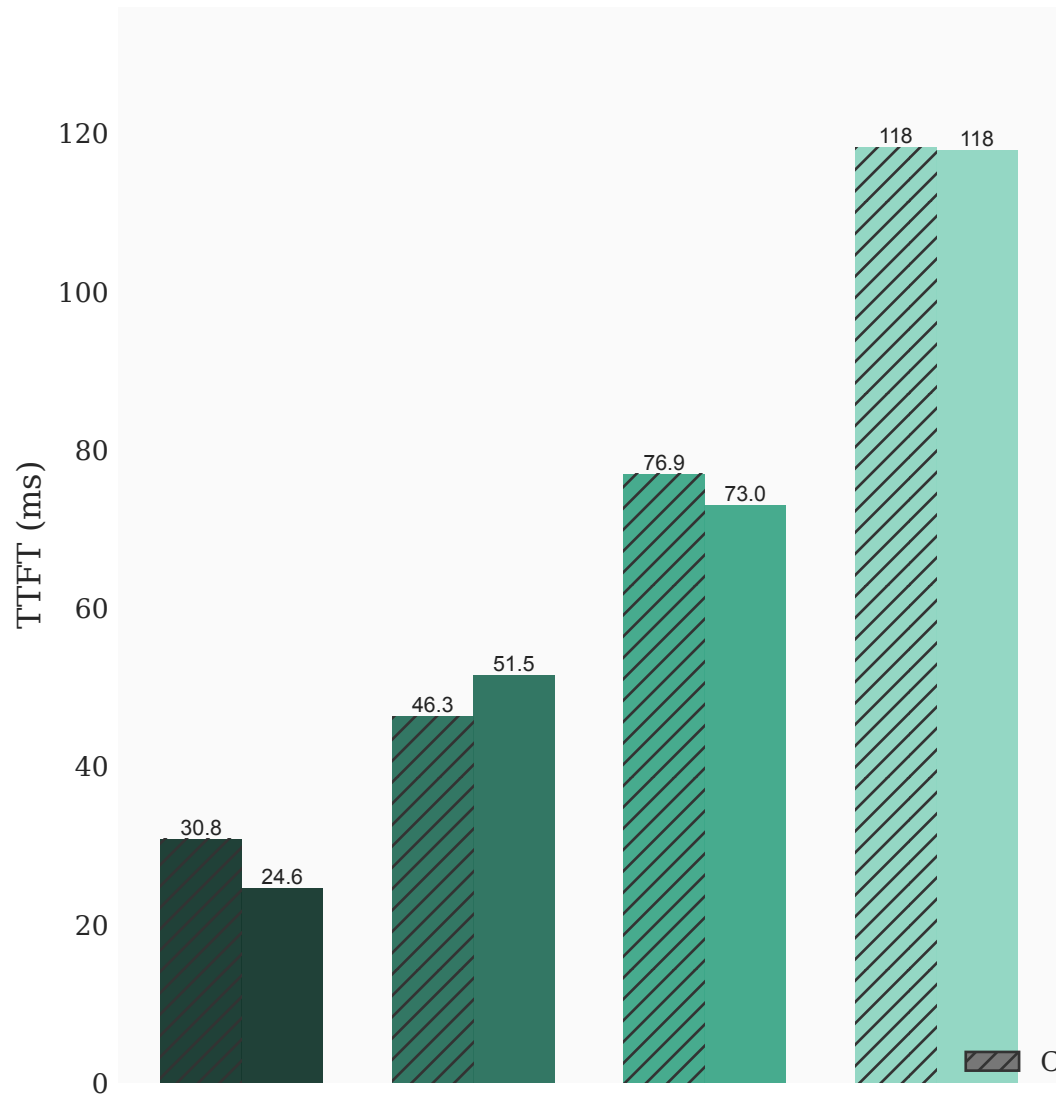# Numina Math (100 Concurrent Requests)

## Time to First Token (Mean)



Bar chart, TTFT (ms) y-axis:
- LLama 3.1 8B: CC 35.3, No CC 29.4
- Mistral 3.1 24B: CC 69.7, No CC 67.7
- GPT OSS 120B: CC 102, No CC 90.2
- LLama 3.3 70B Int4: CC 3.8s, No CC 3.6s

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

## Time to First Token (P99)

Bar chart, TTFT (ms) y-axis:
- LLama 3.1 8B: CC 64.7, No CC 48.7
- Mistral 3.1 24B: CC 127, No CC 124
- GPT OSS 120B: CC 156, No CC 147
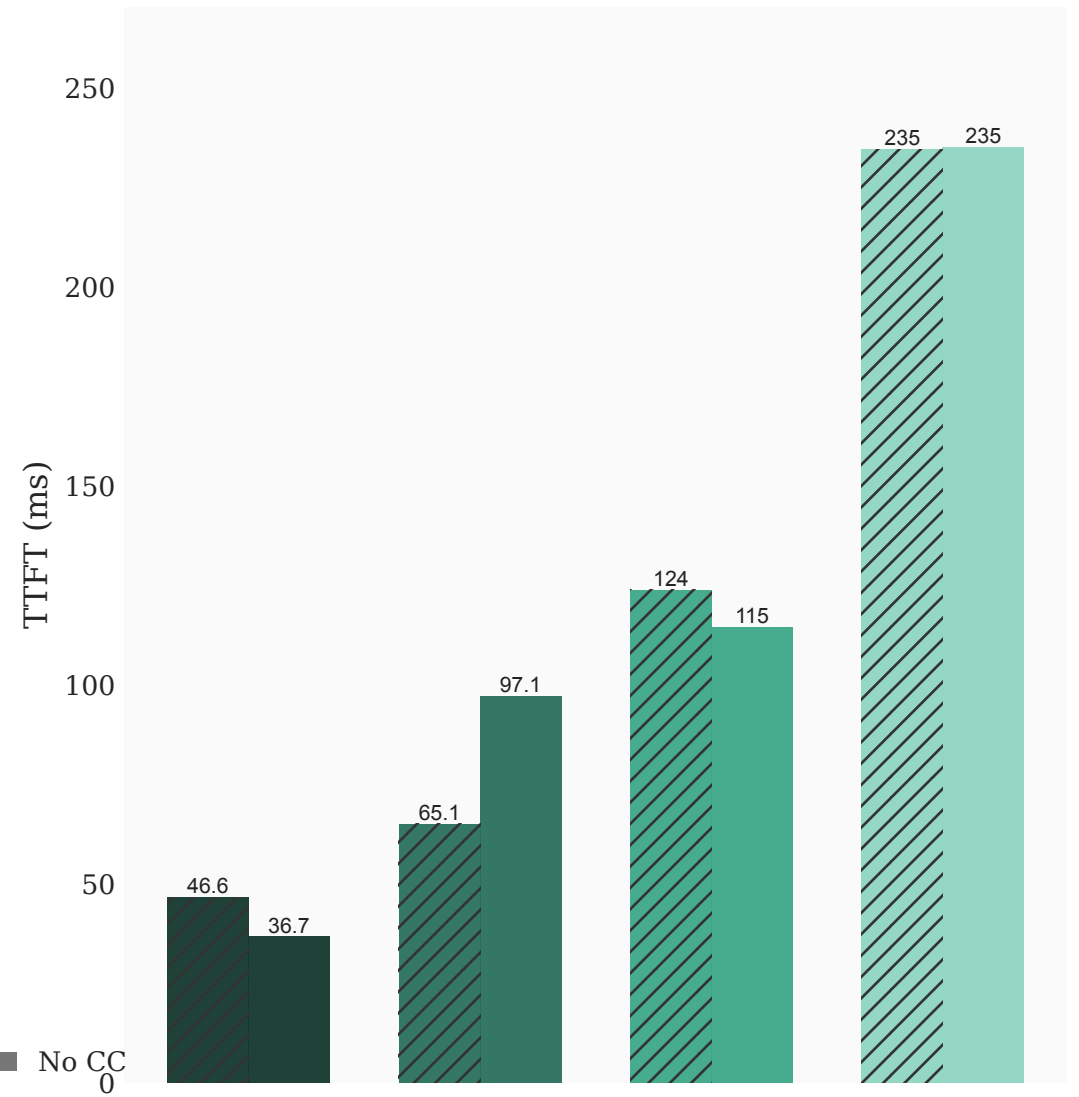- LLama 3.3 70B Int4: CC 6.0s, No CC 5.7s

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

# Numina Math (50 Concurrent Requests)

## Time to First Token (Mean)



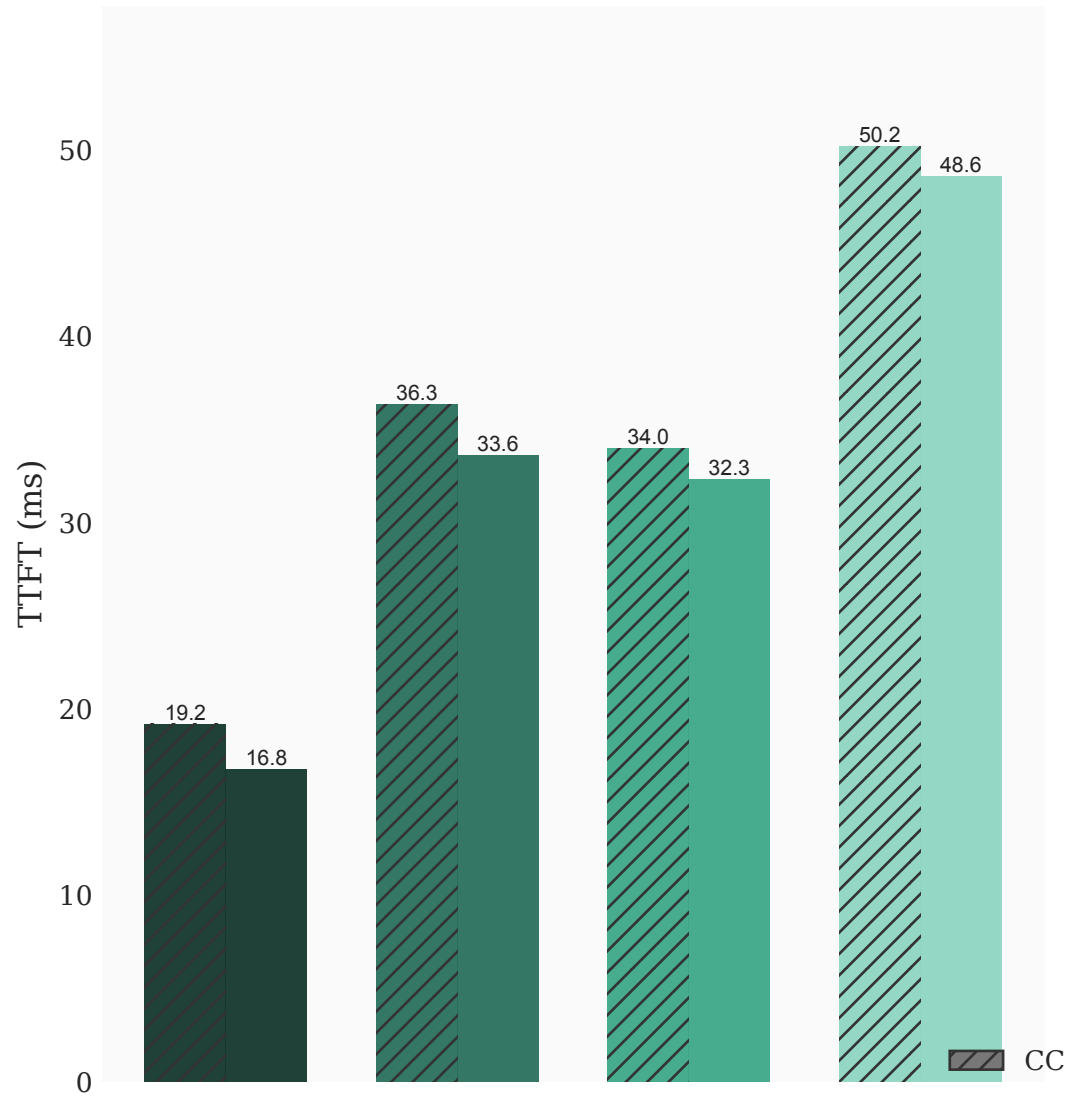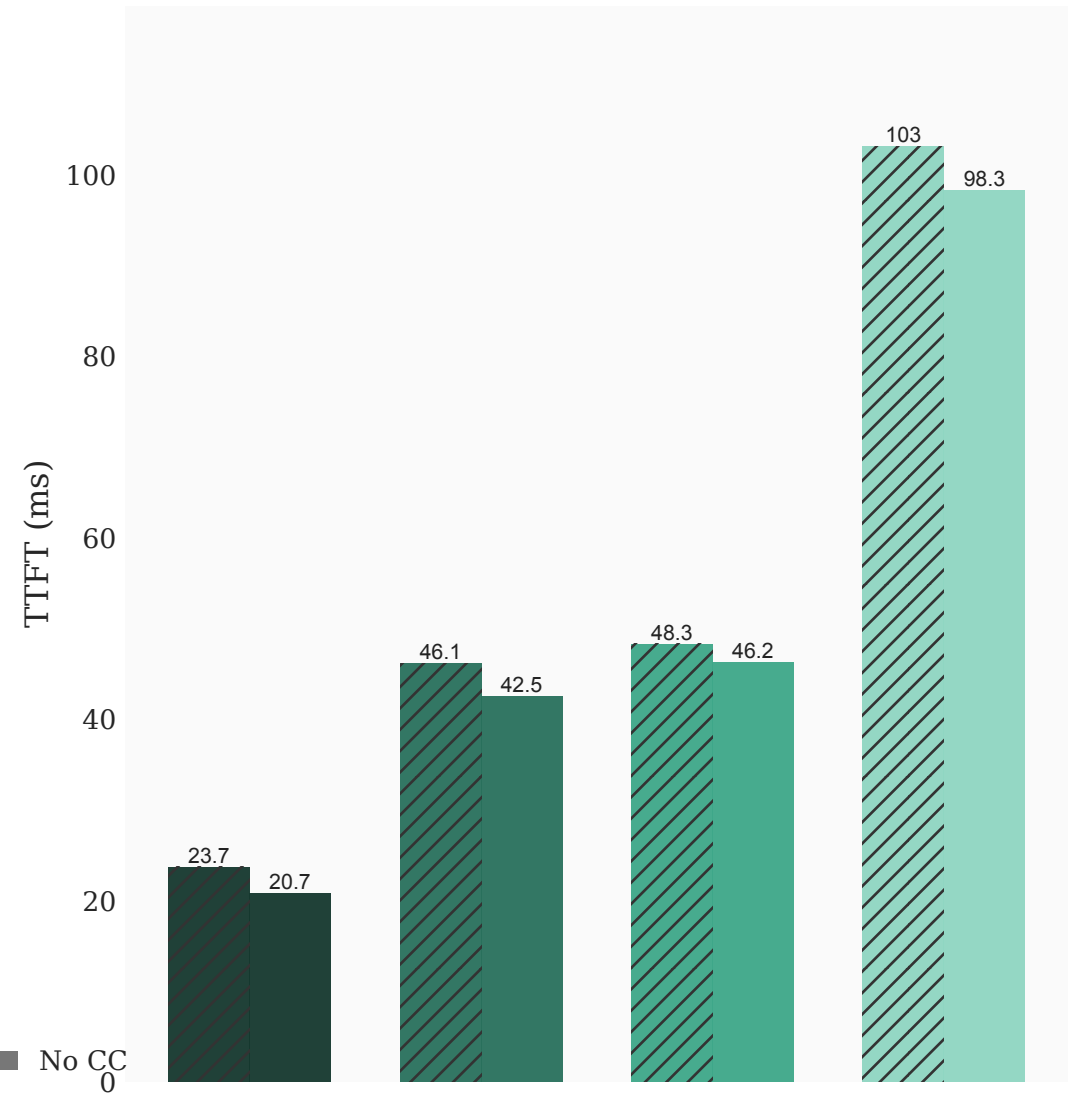## Time to First Token (P99)

**Legend:** ▨ CC   ▬ No CC

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

Time to First Token (Mean) values: 30.8, 24.6, 46.3, 51.5, 76.9, 73.0, 118, 118

Time to First Token (P99) values: 46.6, 36.7, 65.1, 97.1, 124, 115, 235, 235

# Numina Math (1 Concurrent Requests)

## Time to First Token (Mean)



## Time to First Token (P99)



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