# Random (1500 ⇒ 250) (100 Concurrent Requests)
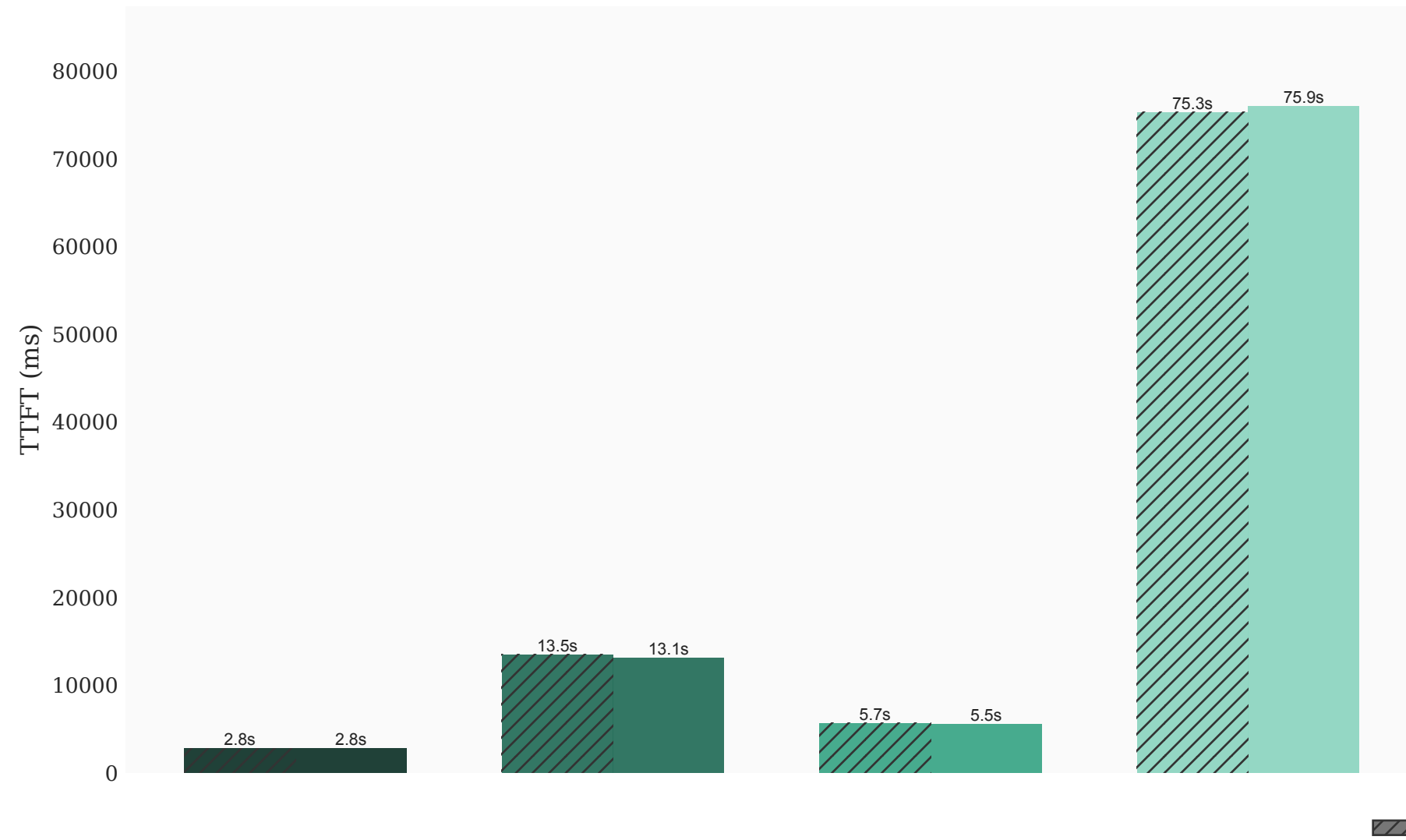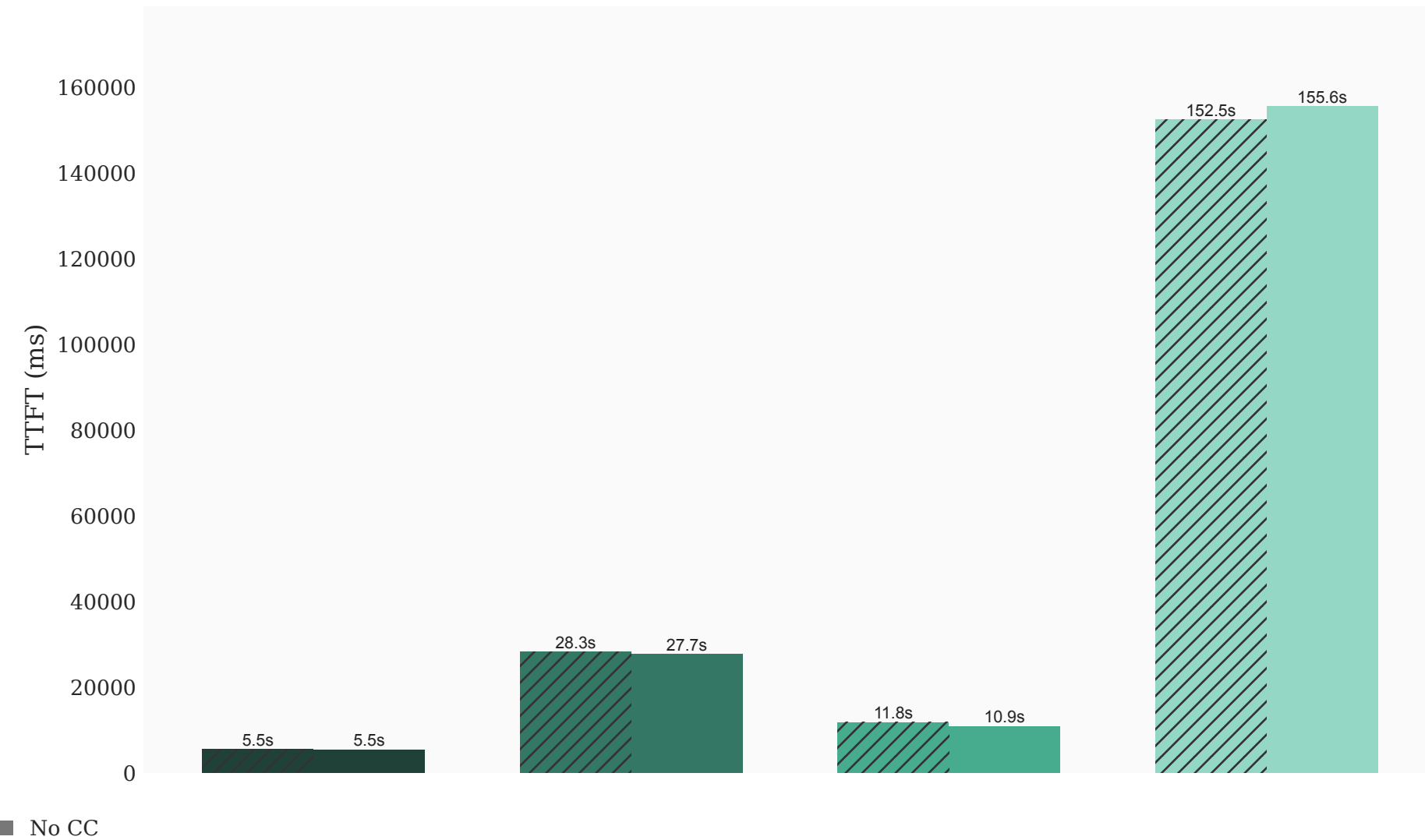
## Time to First Token (Mean)



## Time to First Token (P99)

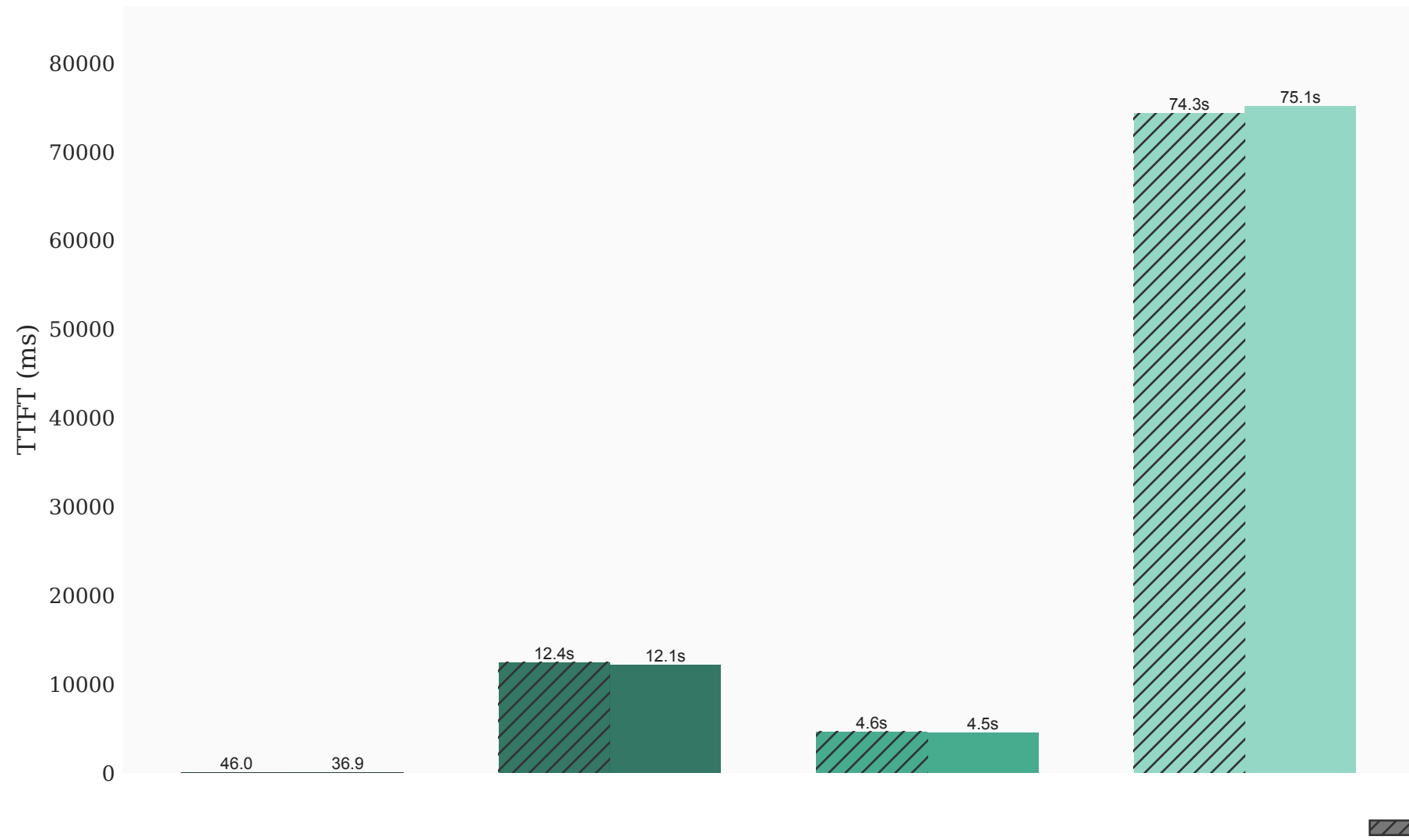Legend: ▨ CC   ▬ No CC

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

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

## Time to First Token (Mean)

TTFT (ms)

46.0 36.9

12.4s 12.1s

4.6s 4.5s

74.3s 75.1s

## Time to First Token (P99)

TTFT (ms)

76.8 60.5

26.2s 25.7s

9.4s 9.0s

151.8s 152.1s

CC  No CC
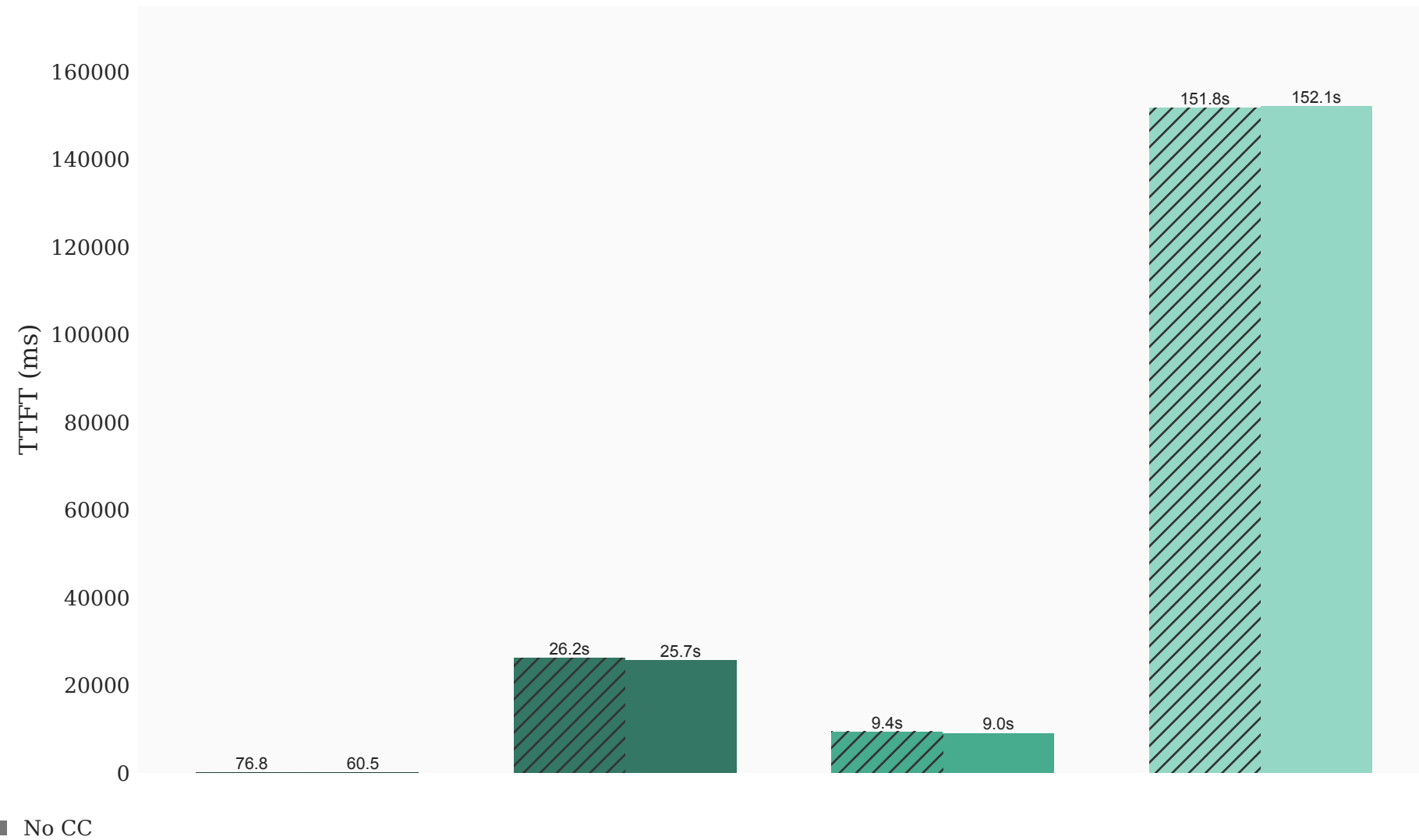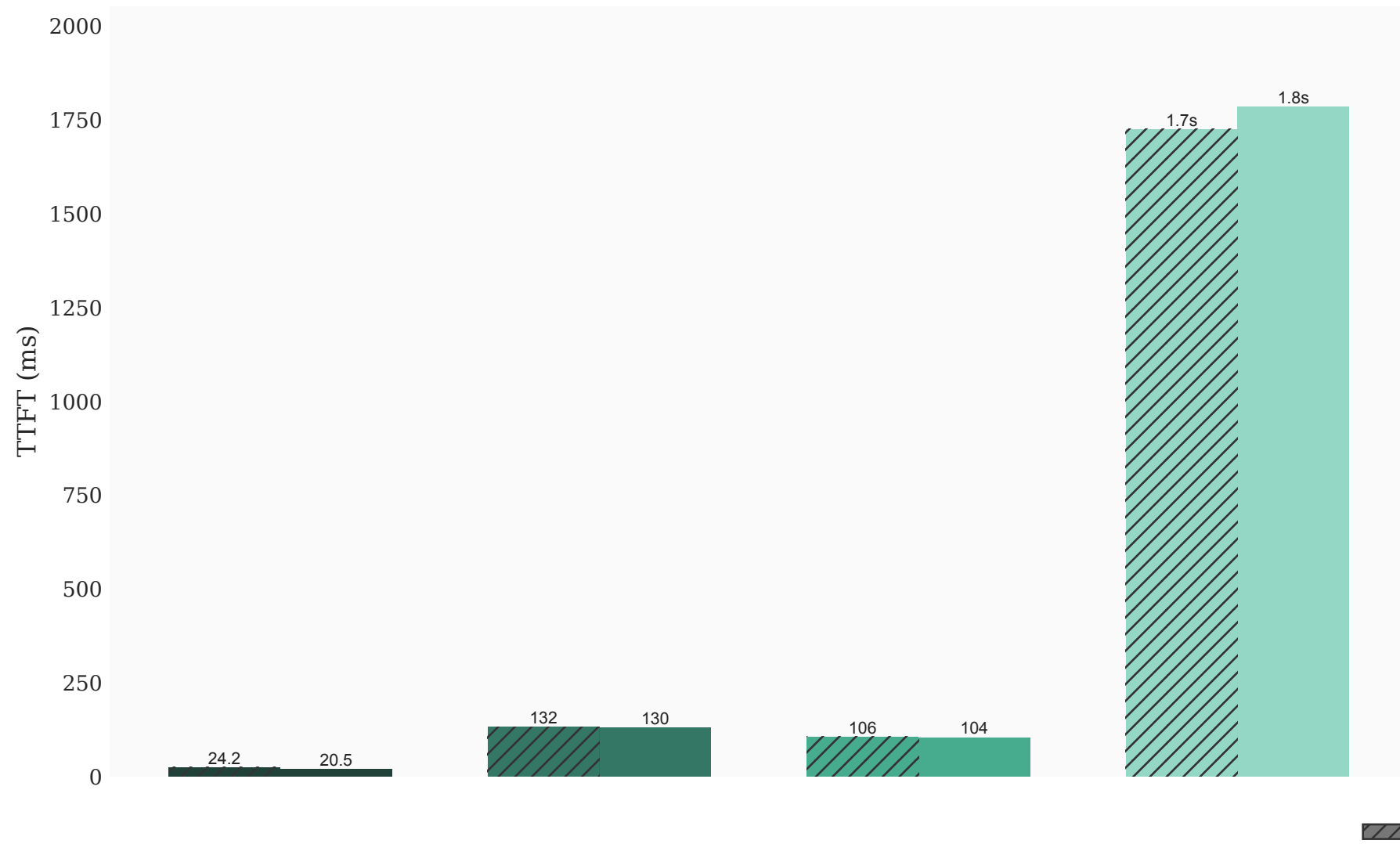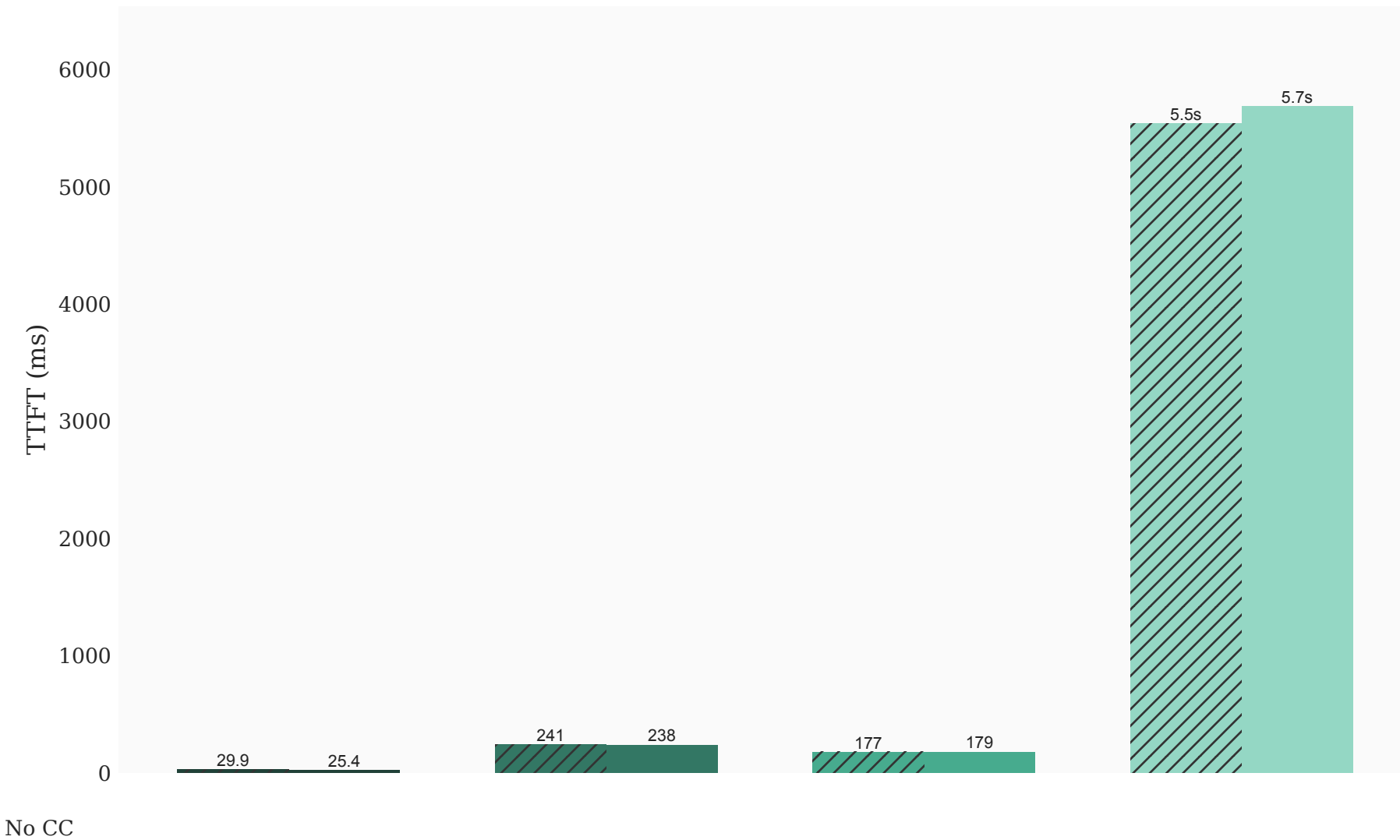
■ LLama 3.1 8B  ■ Mistral 3.1 24B  ■ GPT OSS 120B  ■ LLama 3.3 70B Int4

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

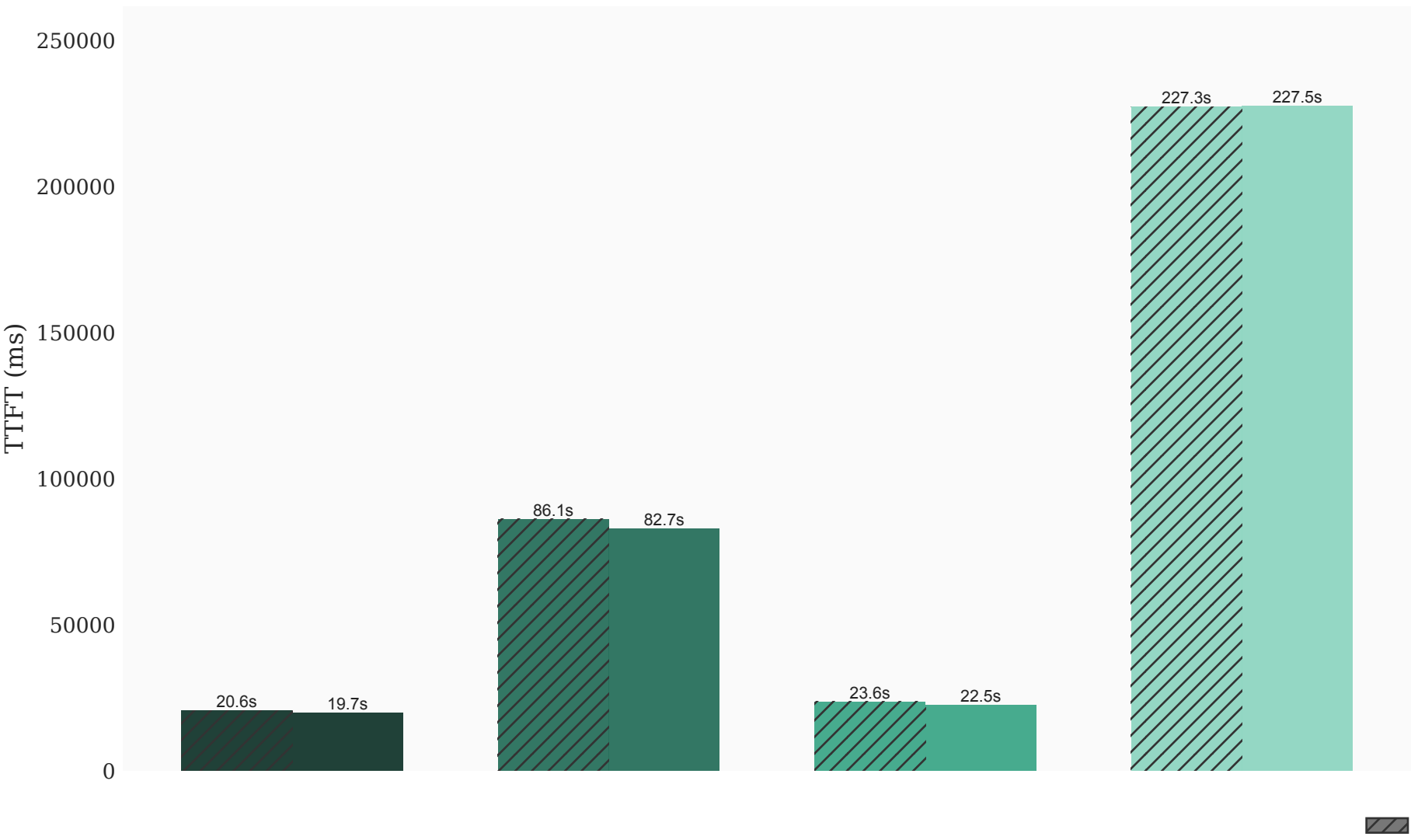## Time to First Token (Mean)

## 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

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

## Time to First Token (Mean)



## Time to First Token (P99)

Mean chart values: 20.6s, 19.7s, 86.1s, 82.7s, 23.6s, 22.5s, 227.3s, 227.5s

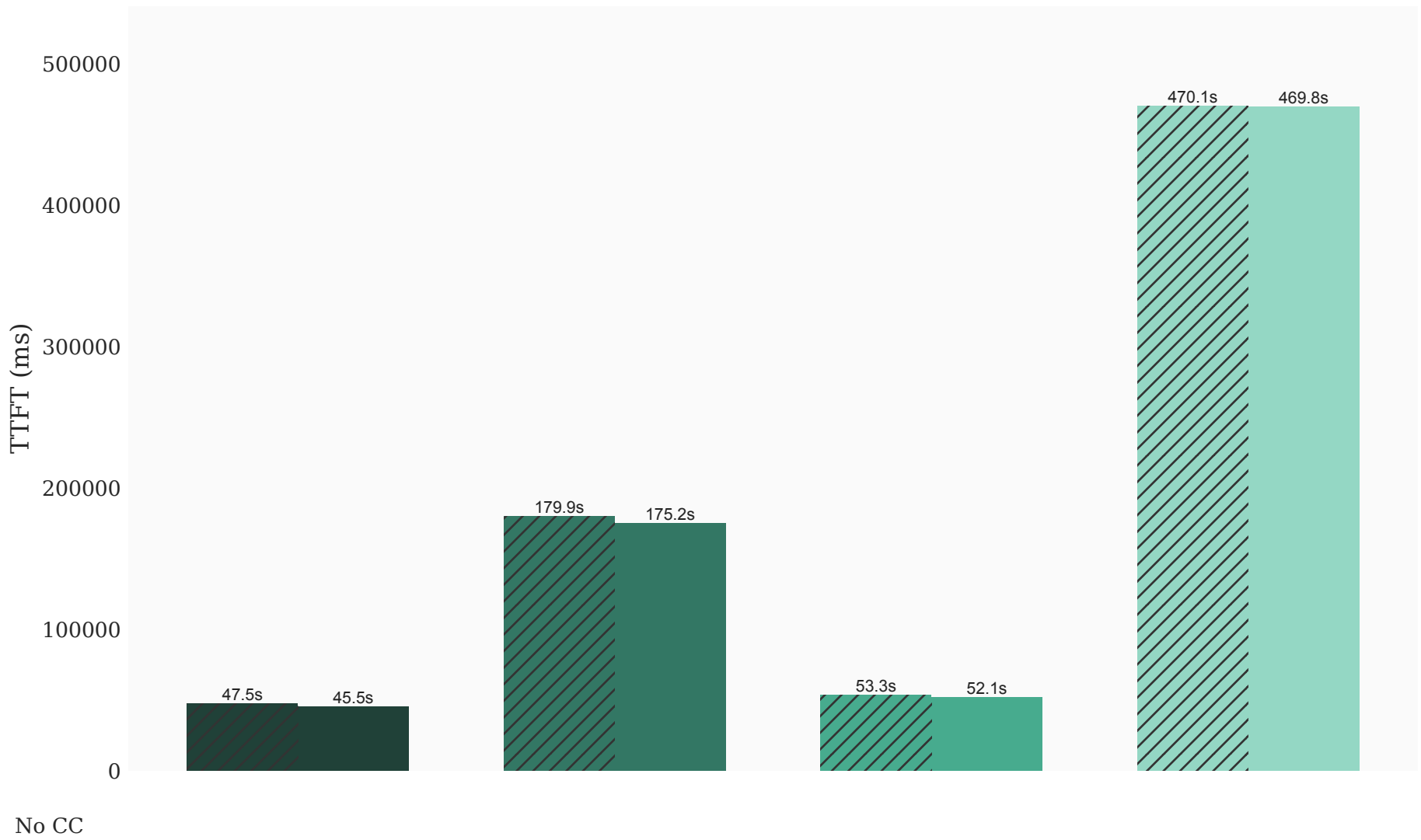P99 chart values: 47.5s, 45.5s, 179.9s, 175.2s, 53.3s, 52.1s, 470.1s, 469.8s

Legend: CC, No CC

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

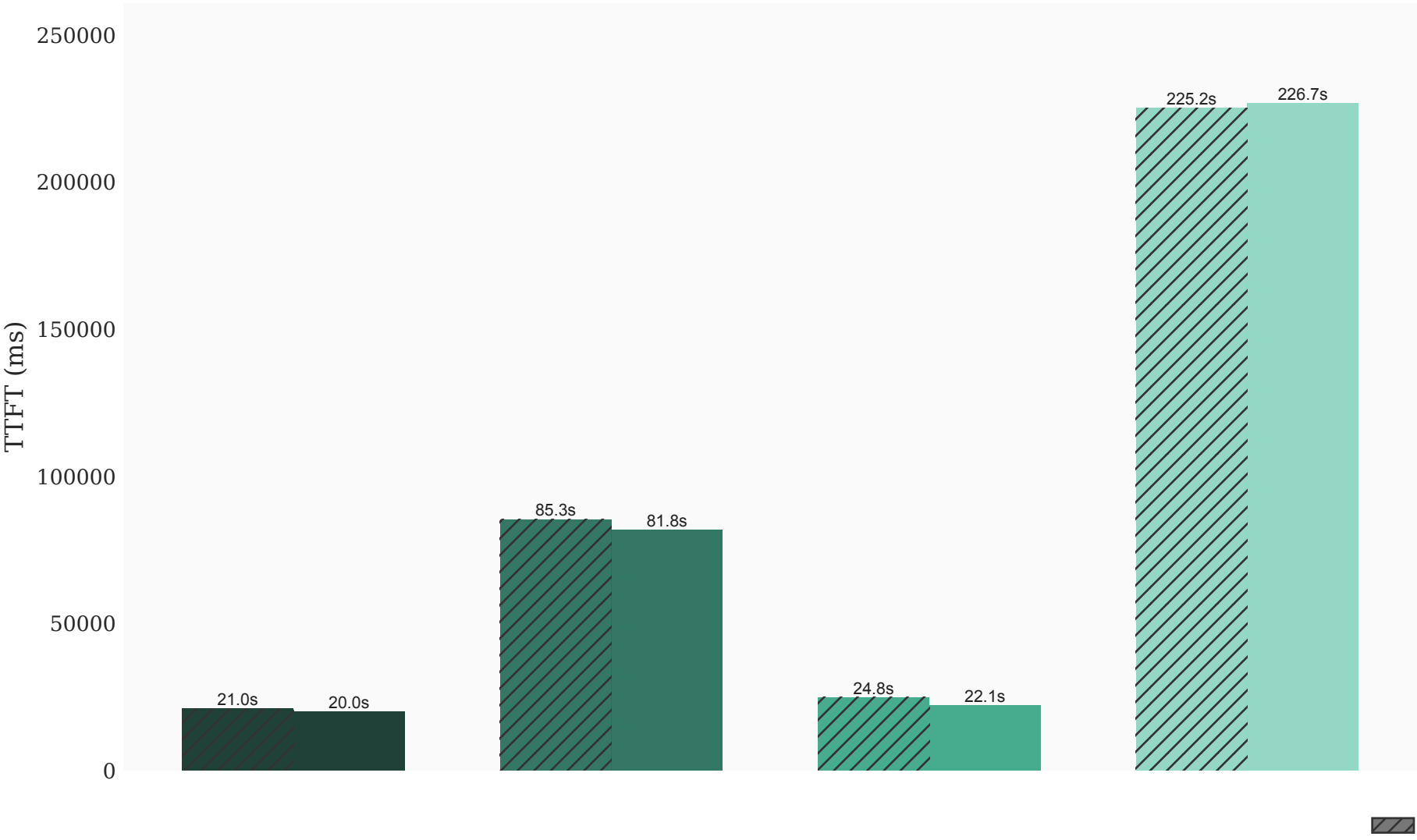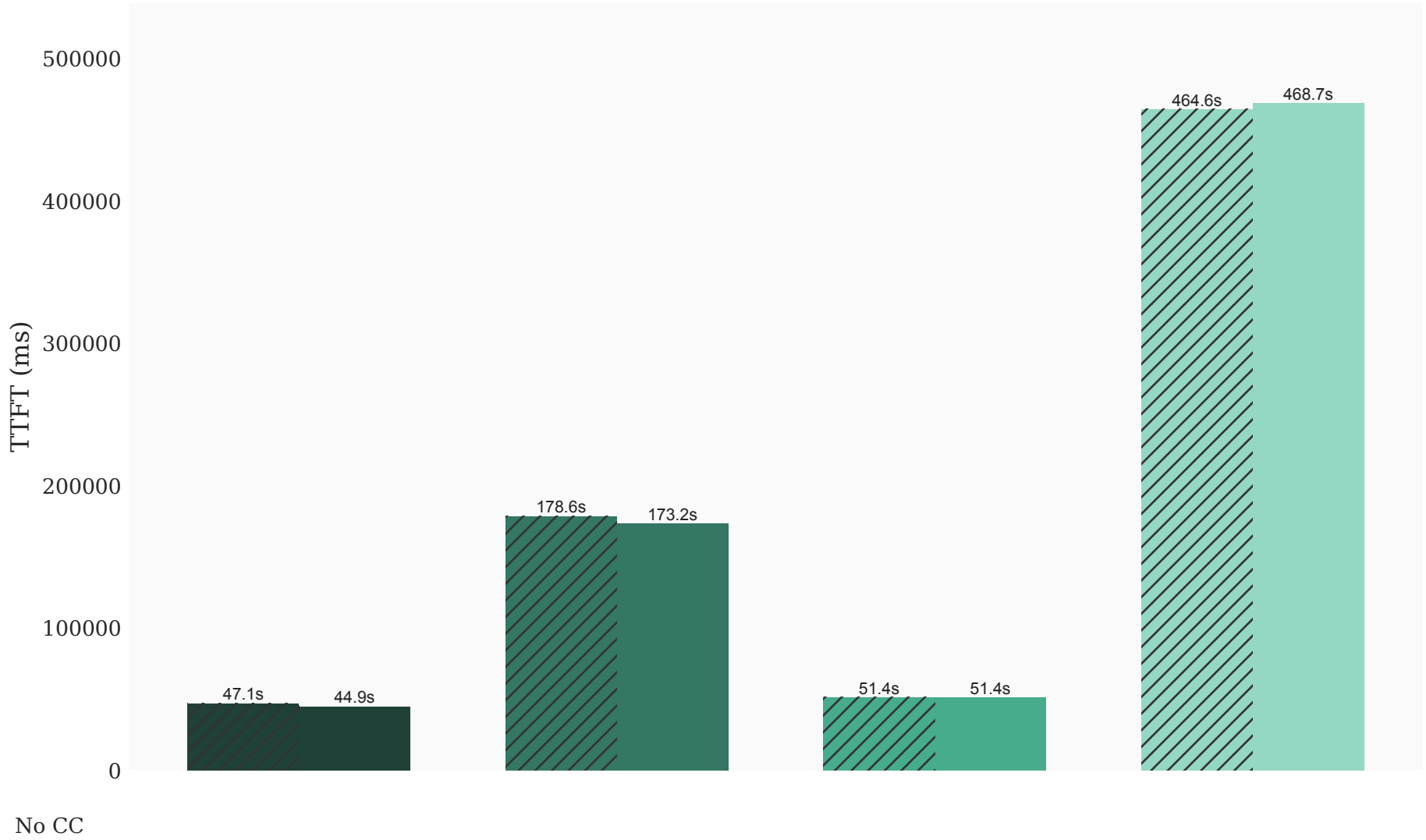# Random (4000 ⇒ 1000) (50 Concurrent Requests)

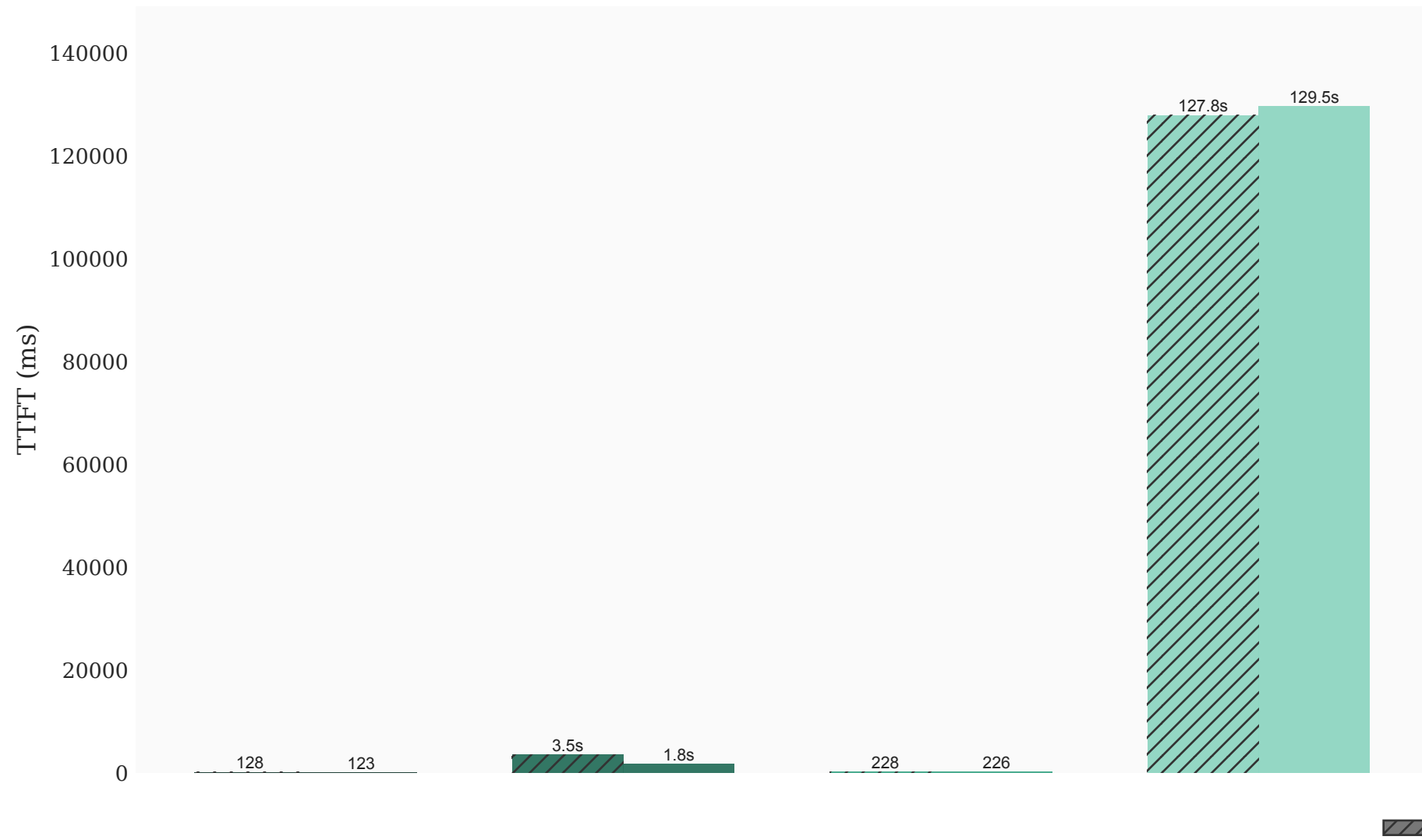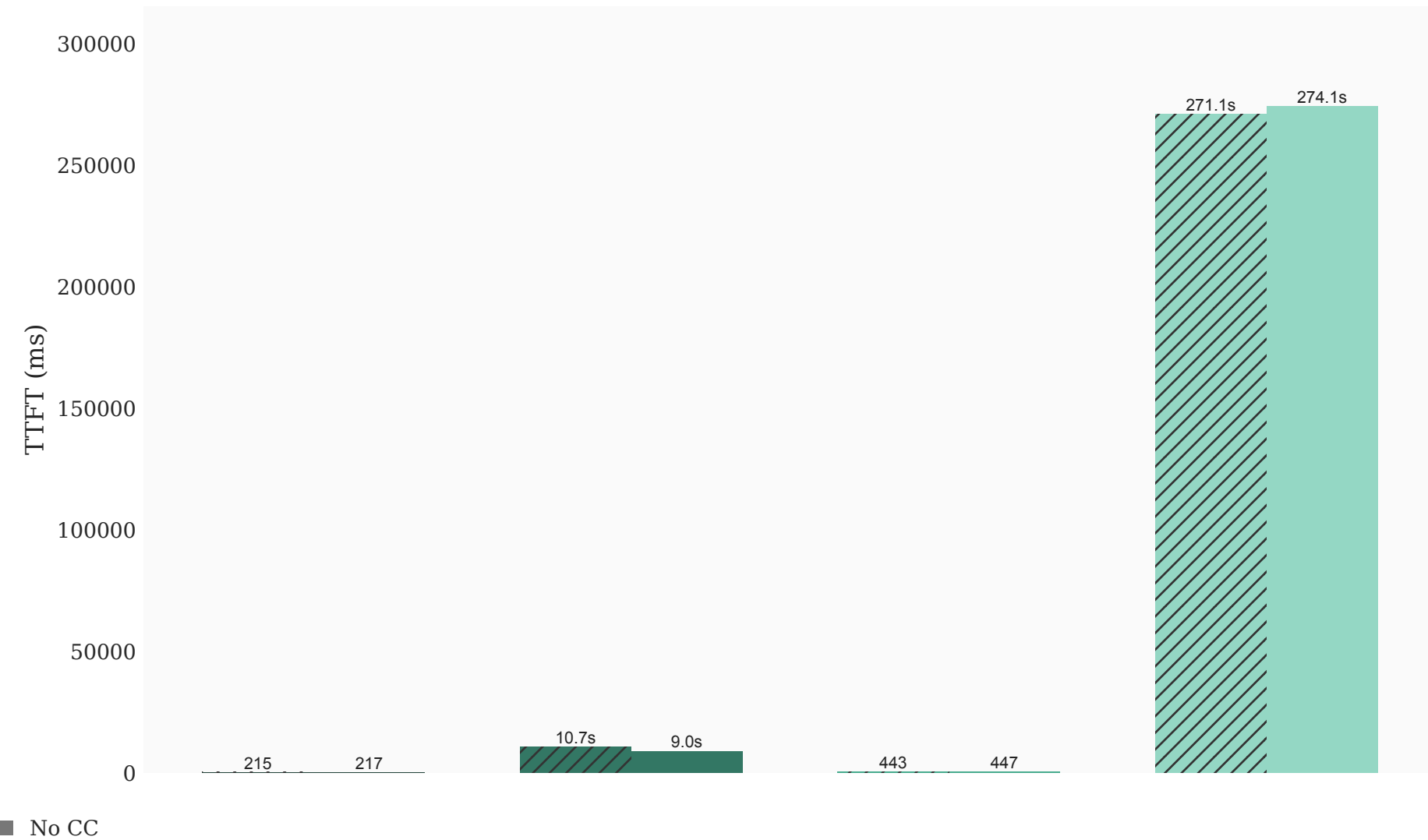## Time to First Token (Mean)



## Time to First Token (P99)

Mean values: 21.0s, 20.0s, 85.3s, 81.8s, 24.8s, 22.1s, 225.2s, 226.7s

P99 values: 47.1s, 44.9s, 178.6s, 173.2s, 51.4s, 51.4s, 464.6s, 468.7s

Legend: CC, No CC

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)

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

## Time to First Token (P99)

TTFT (ms)

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

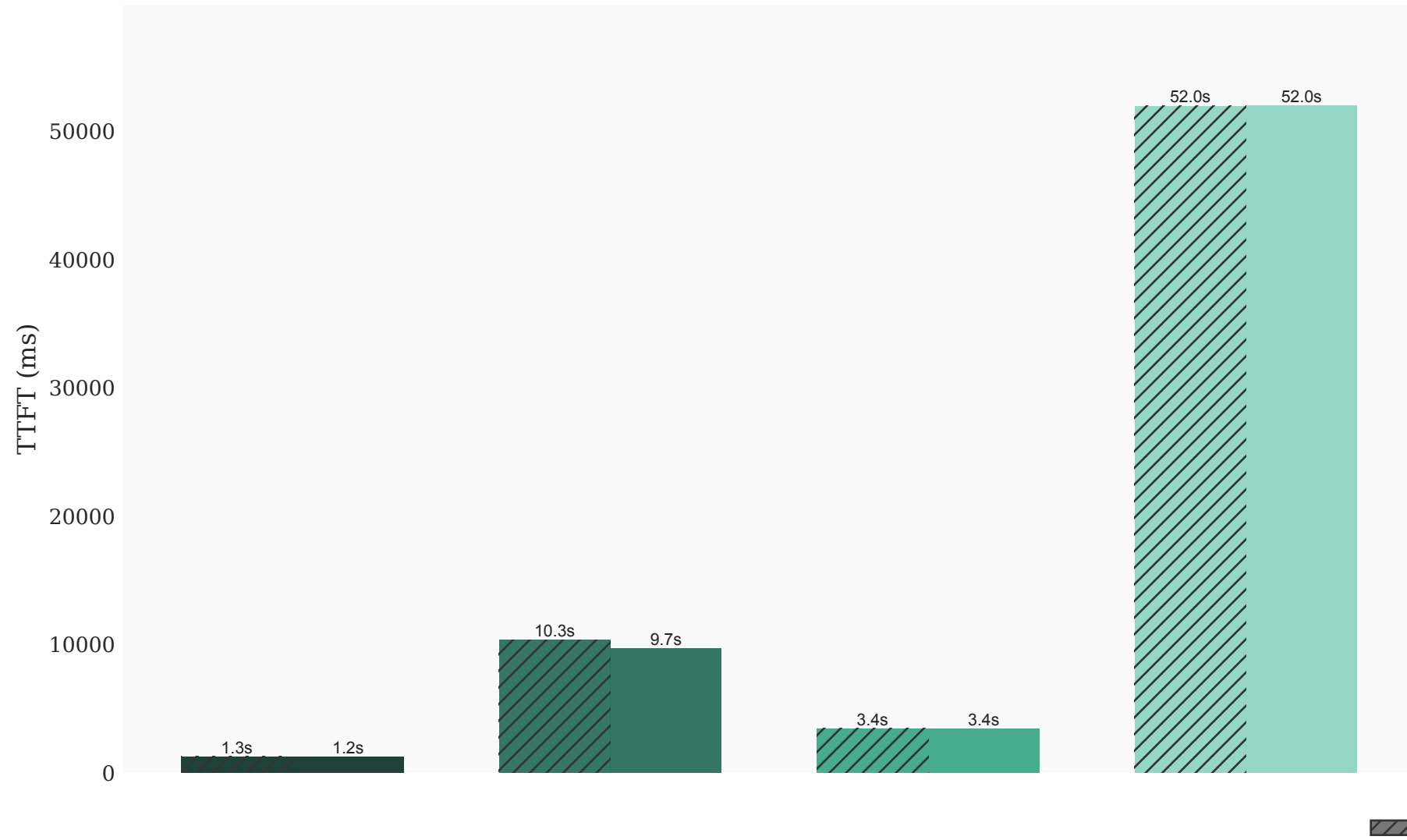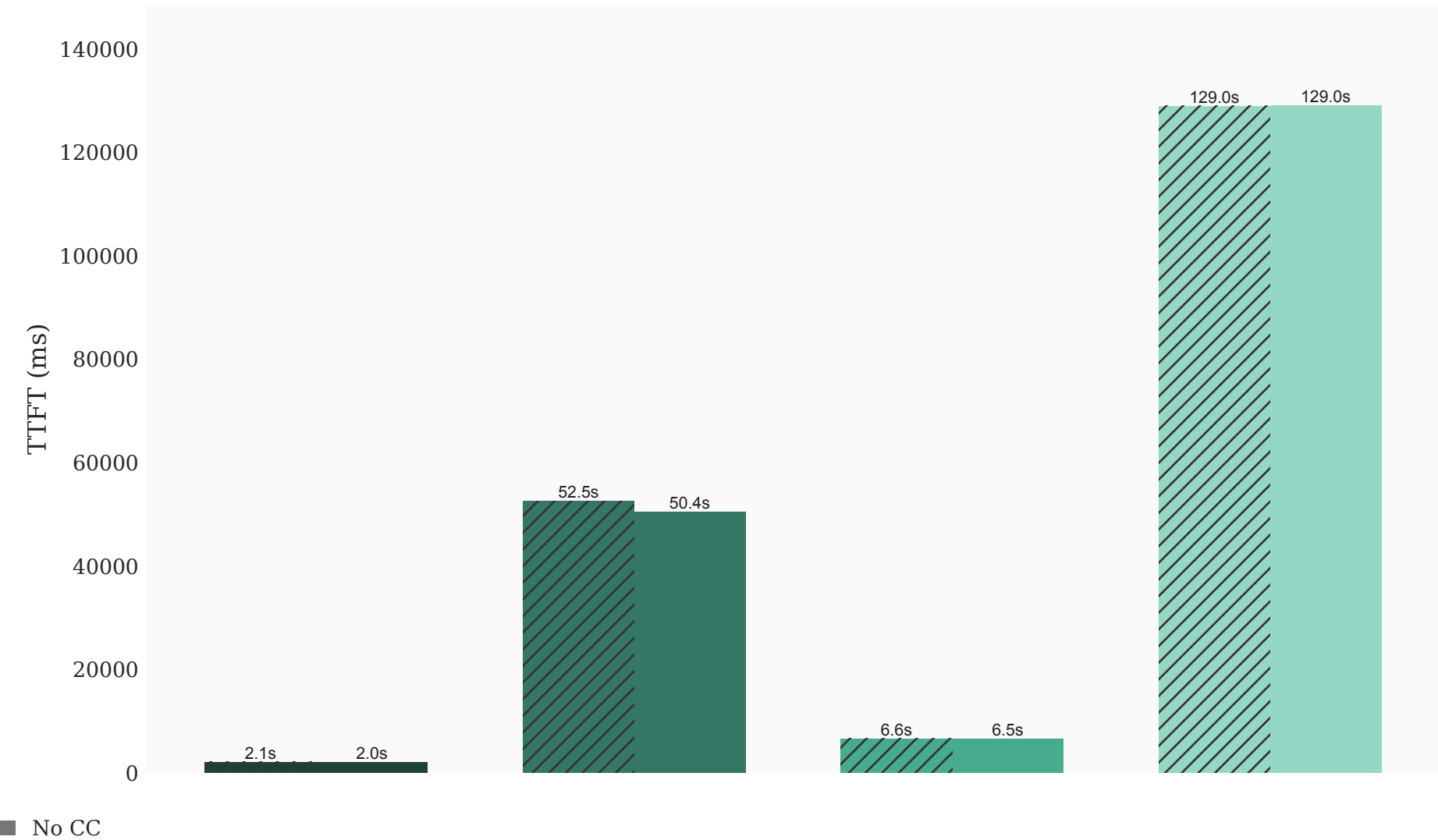# Random (1000 ⇒ 1000) (100 Concurrent Requests)

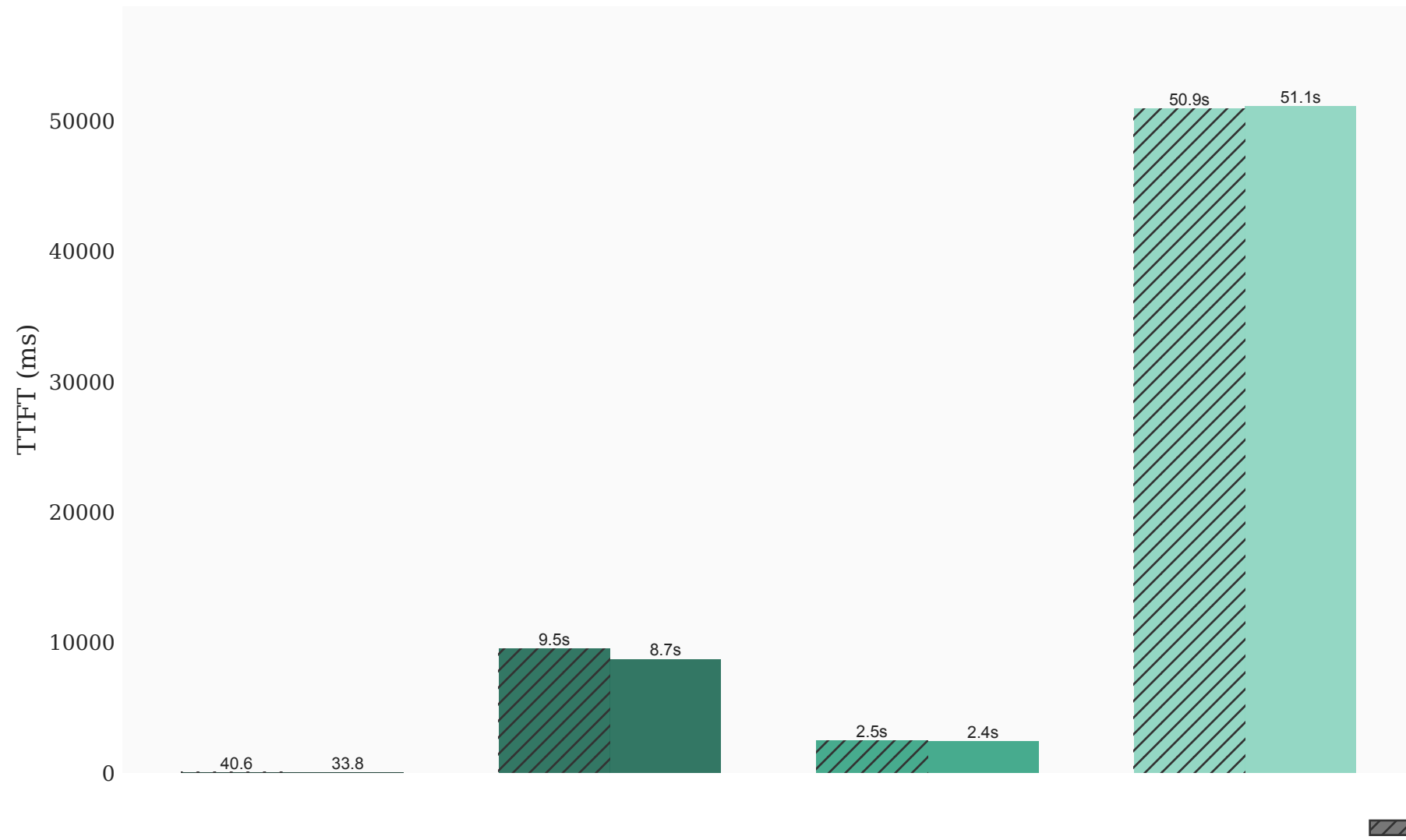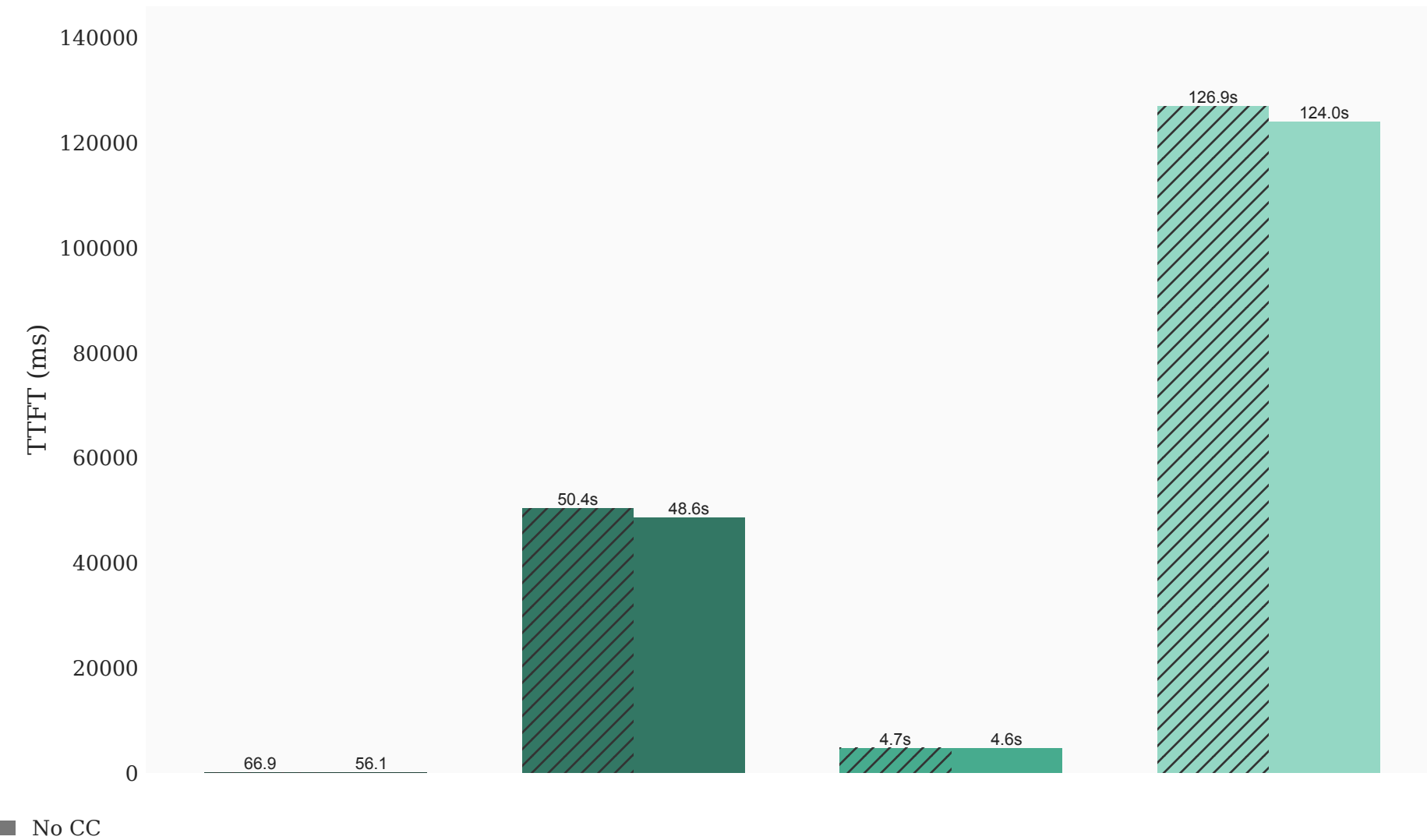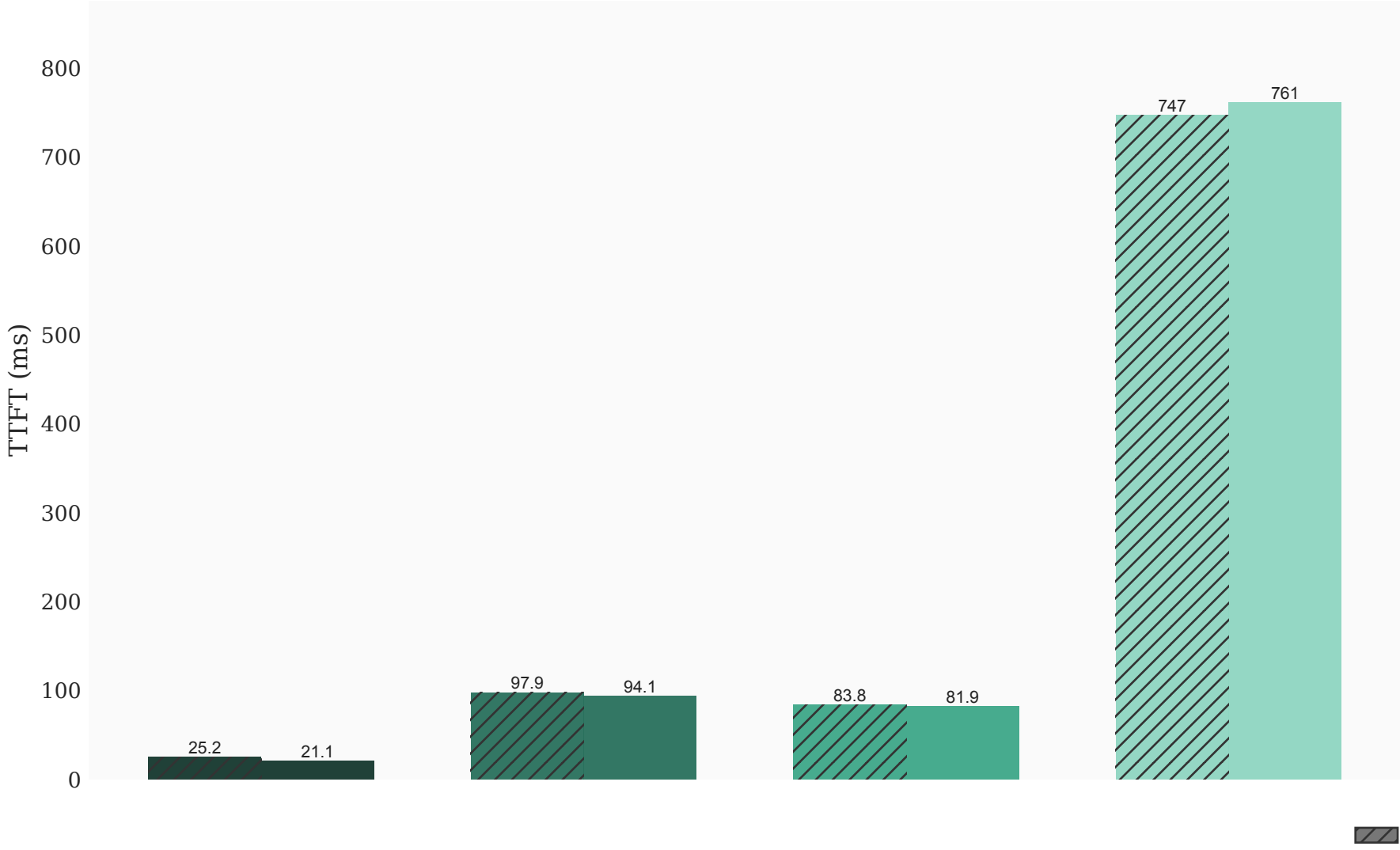## Time to First Token (Mean)



## Time to First Token (P99)

Legend: CC, No CC

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

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

**Time to First Token (Mean)**

**Time to First Token (P99)**

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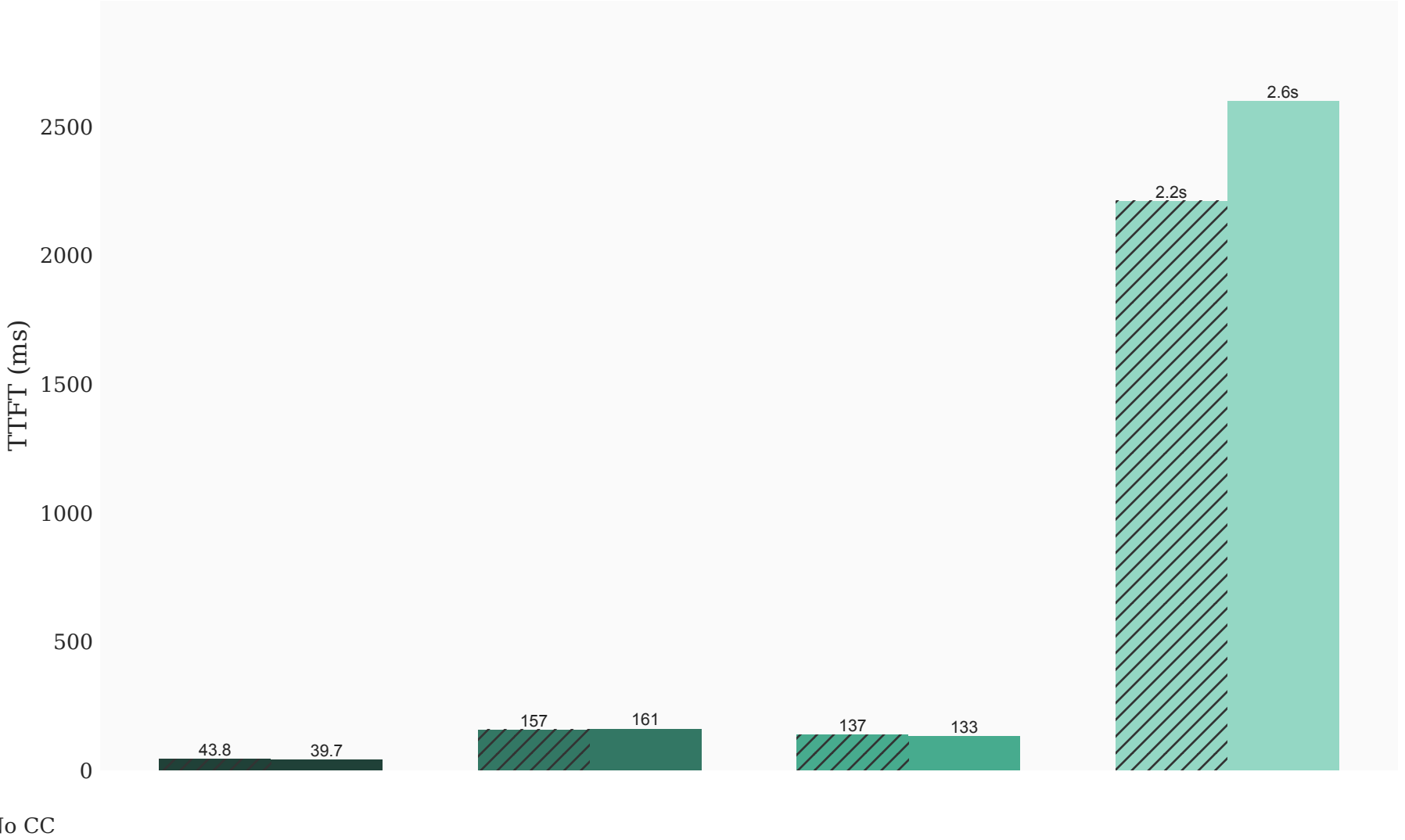
LLama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, LLama 3.3 70B Int4

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

## Time to First Token (Mean)



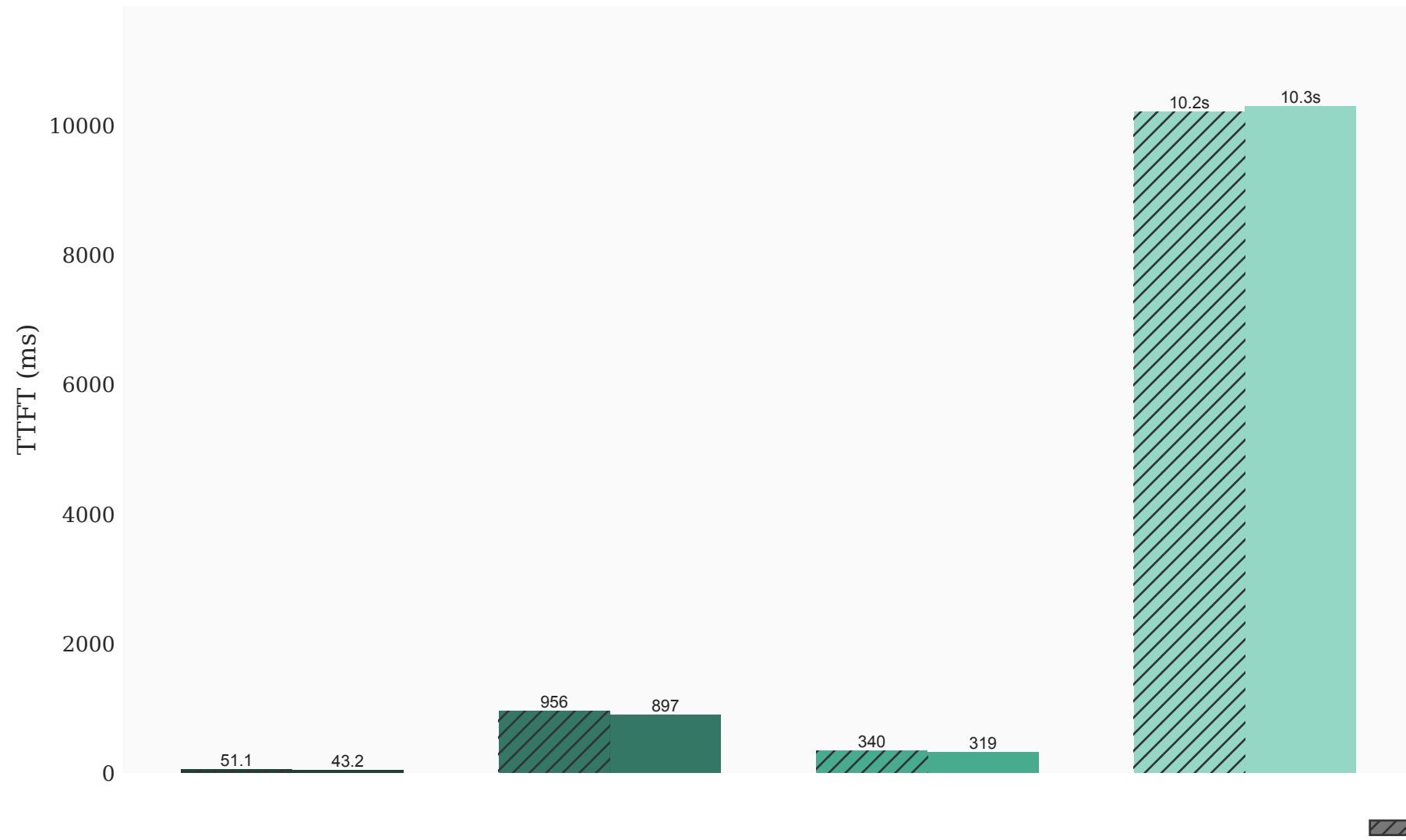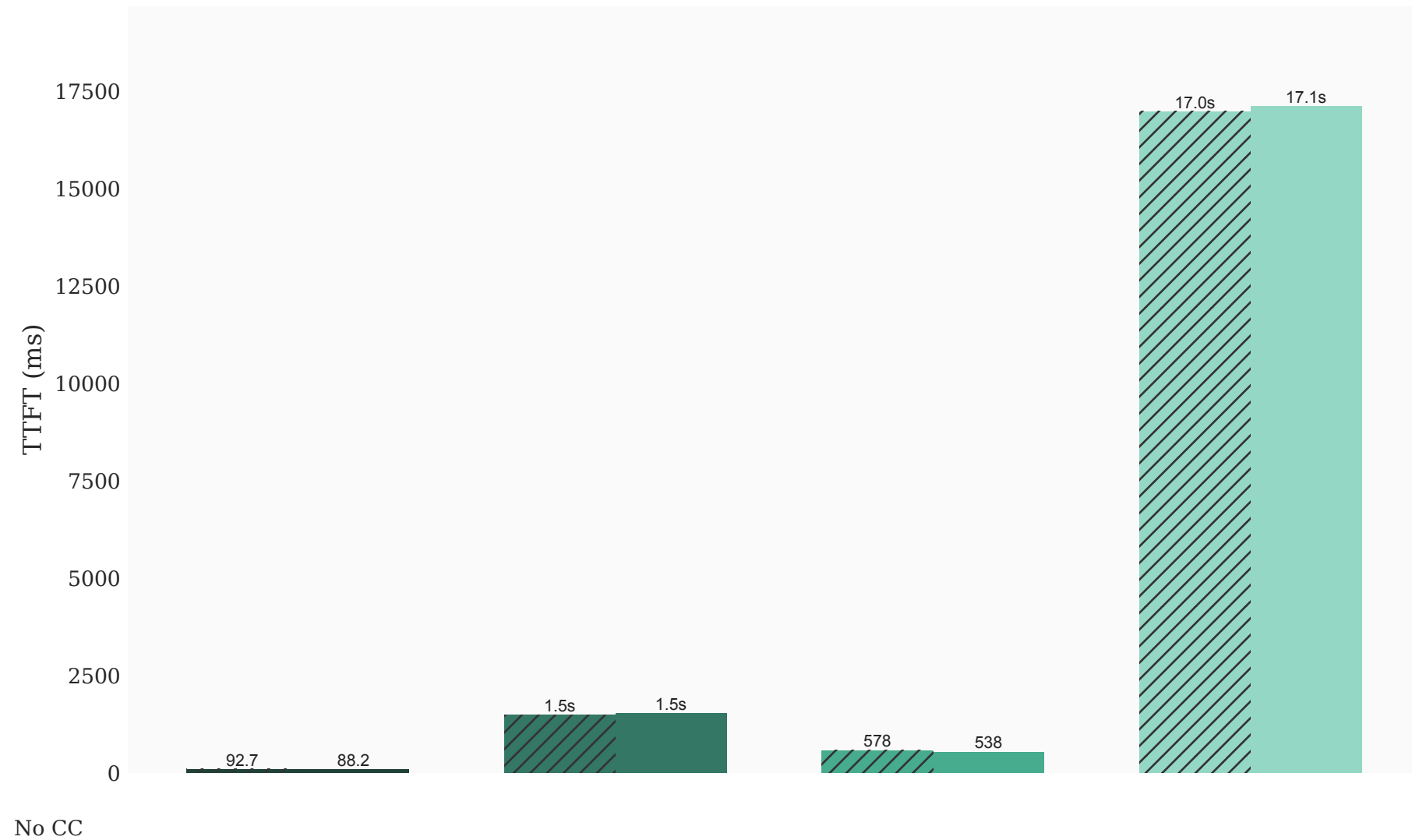## Time to First Token (P99)

Legend: CC | No CC

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

# ShareGPT (100 Concurrent Requests)

## Time to First Token (Mean)



## Time to First Token (P99)

| | 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)
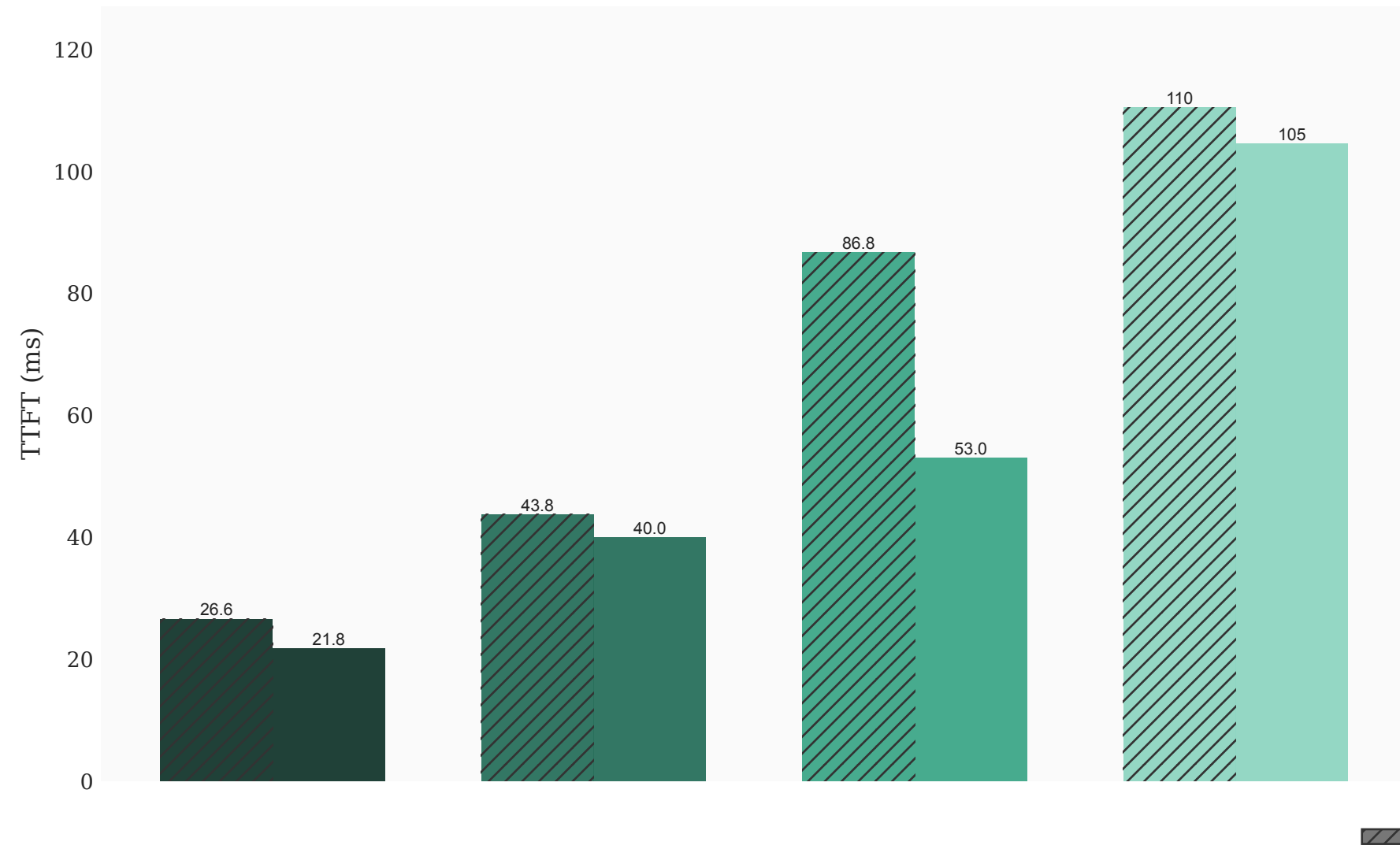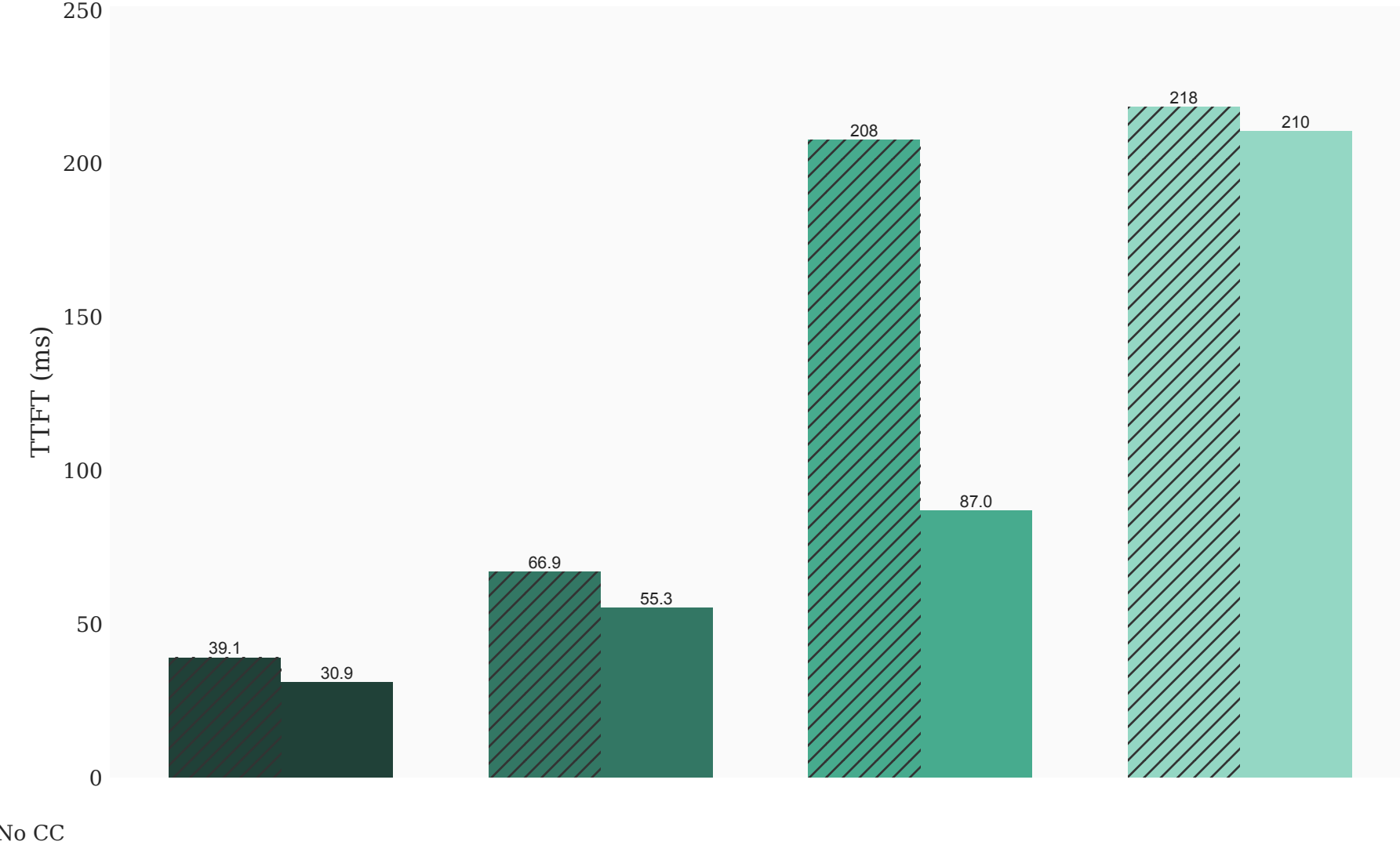


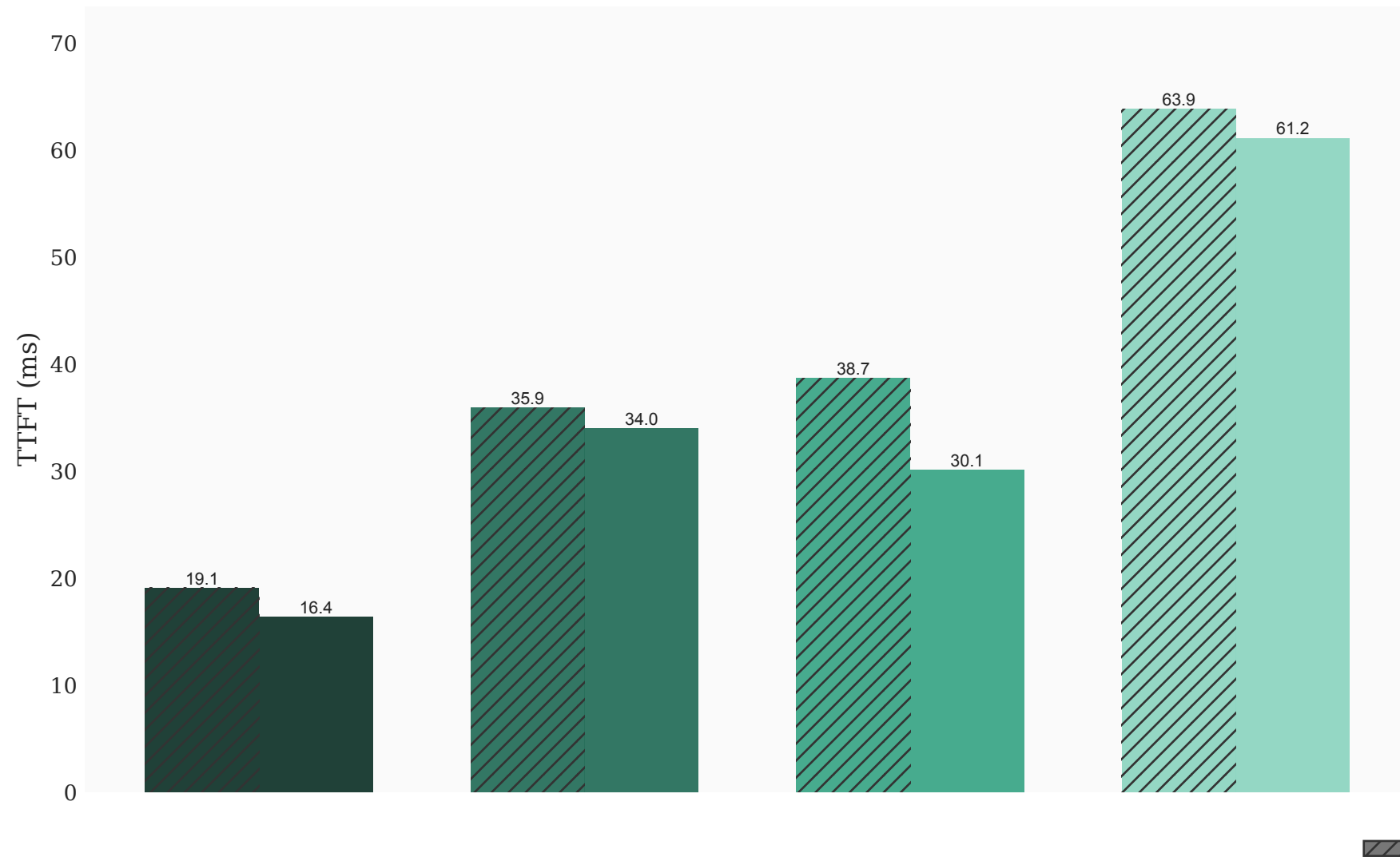## Time to First Token (P99)

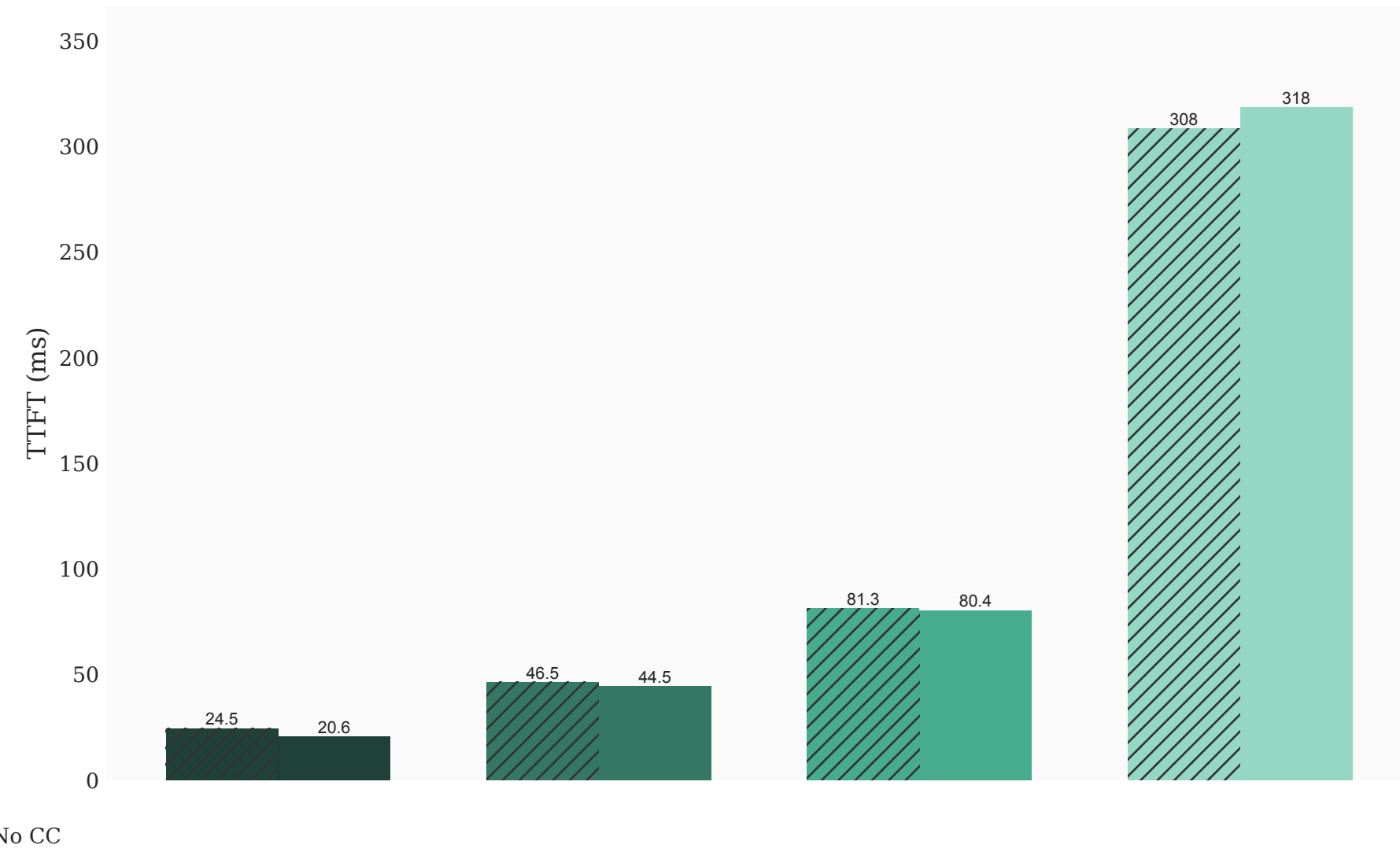Legend: CC (hatched), No CC (solid)

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)



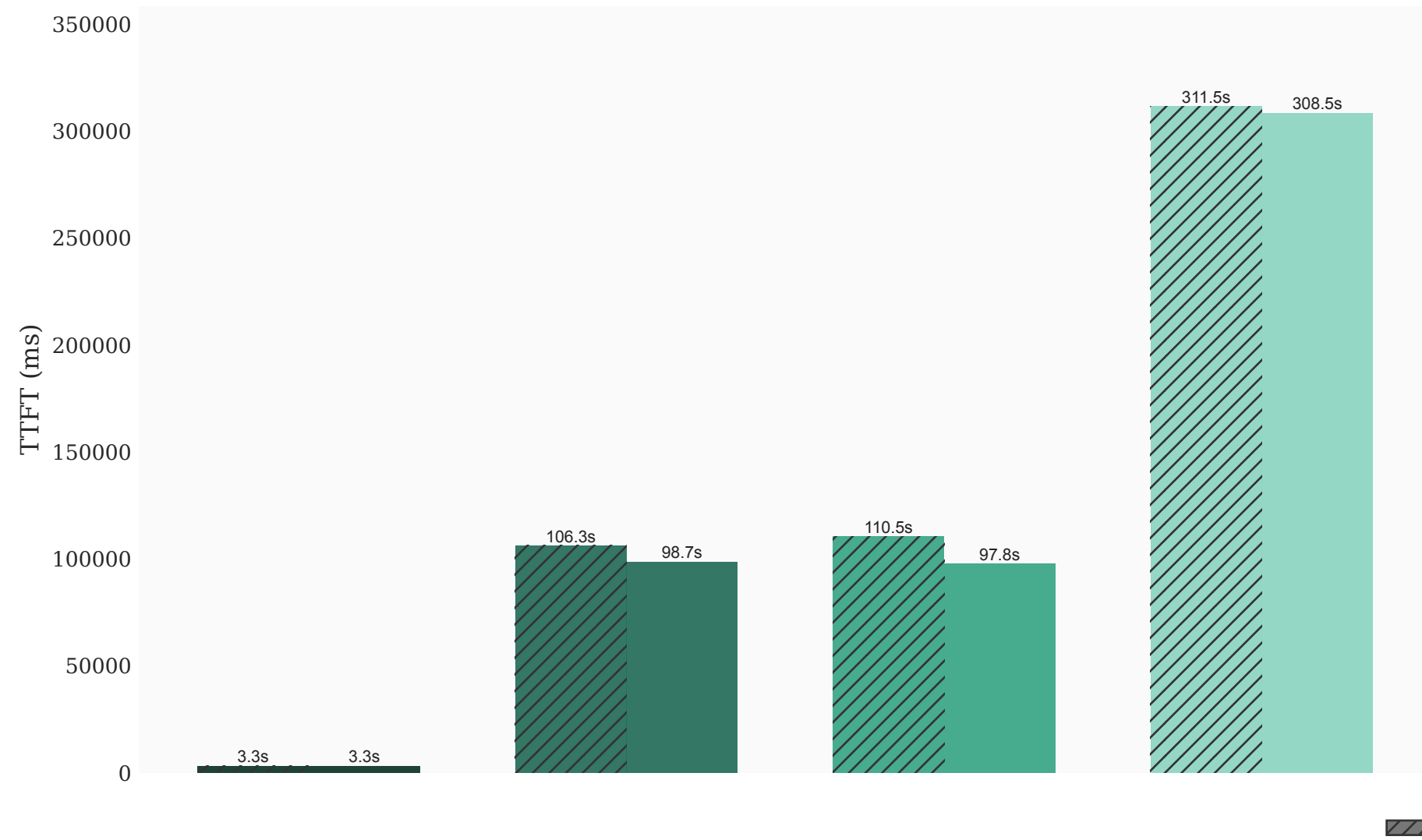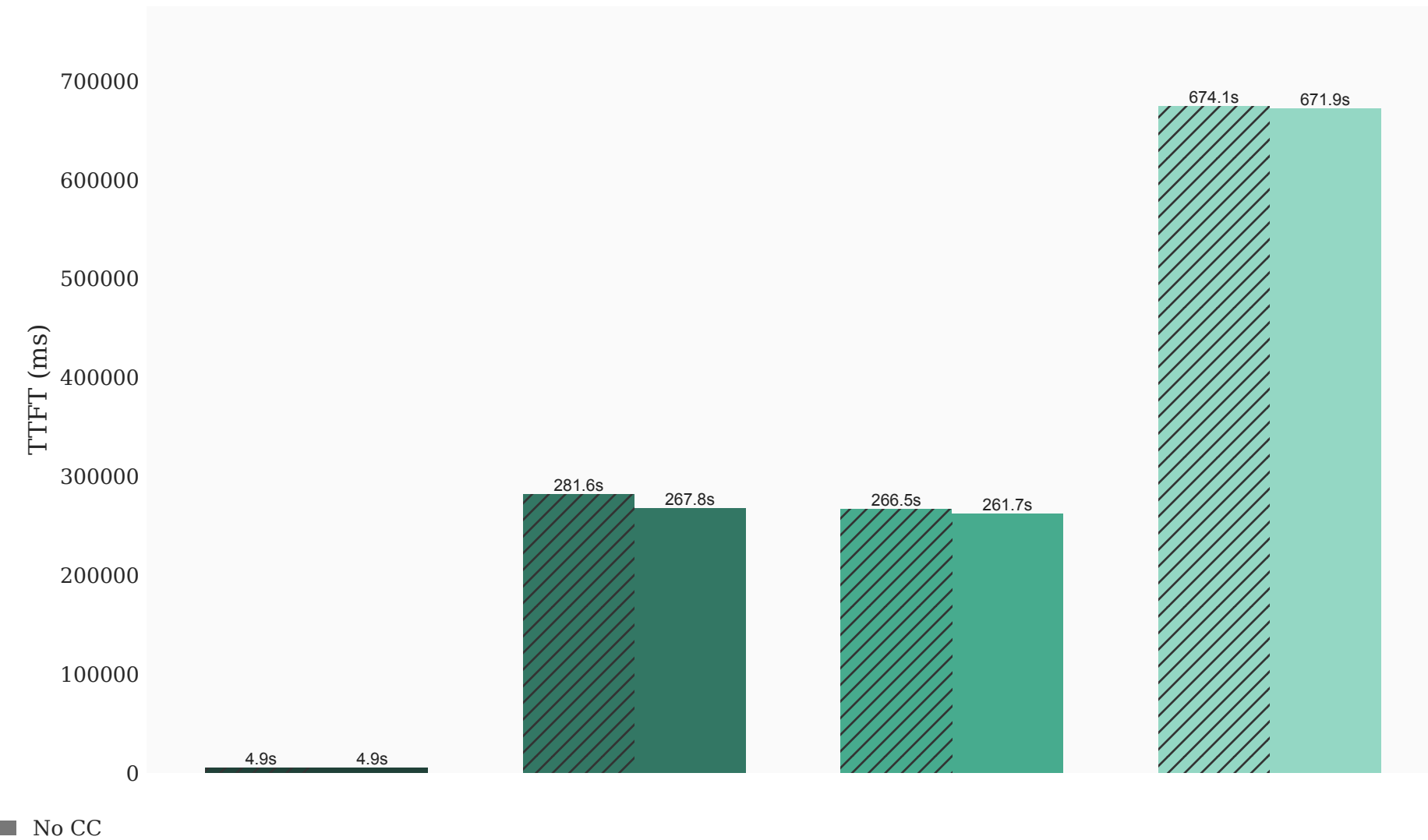Legend: CC / No CC

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)



## 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

# Edit 10K Characters (50 Concurrent Requests)

## Time to First Token (Mean)



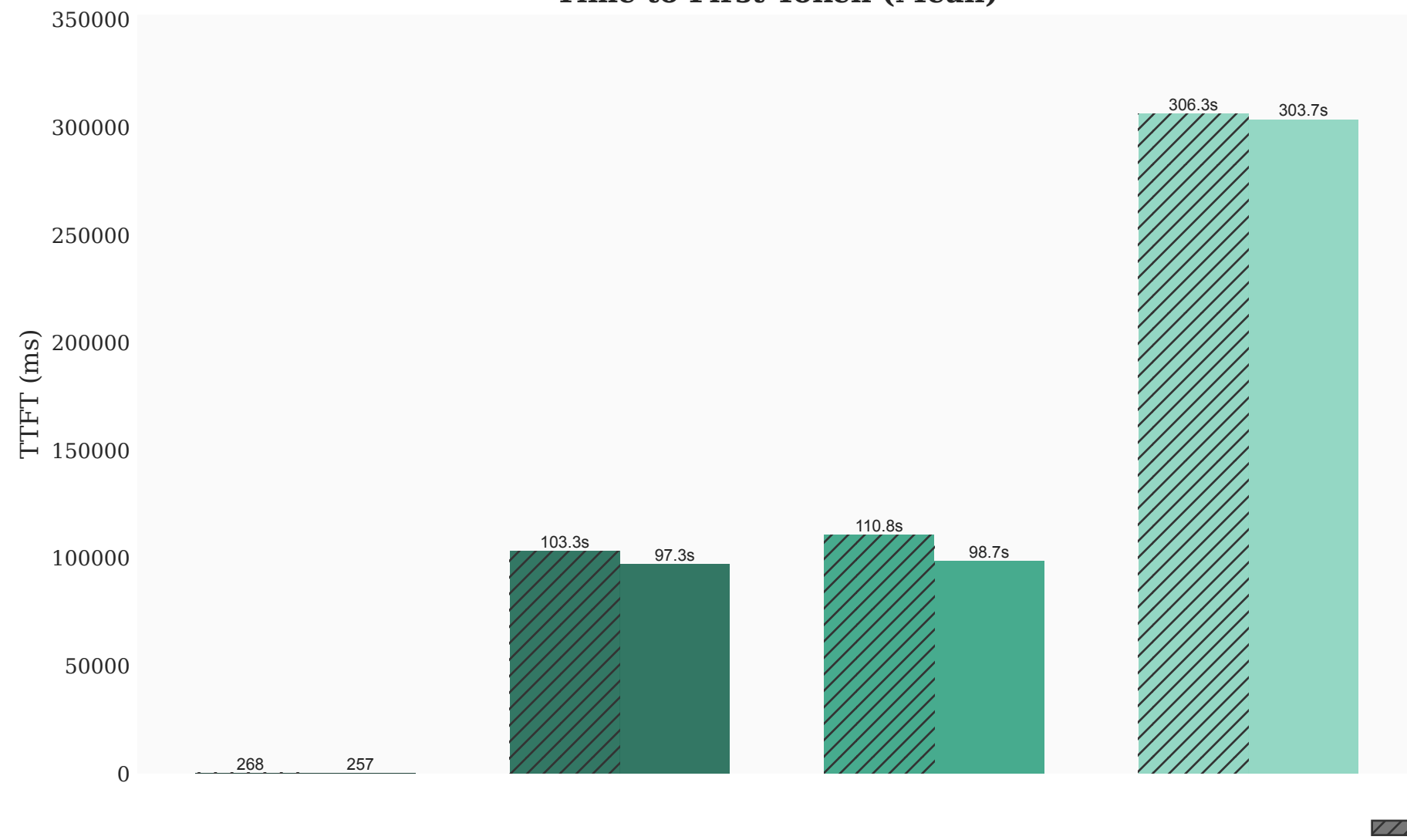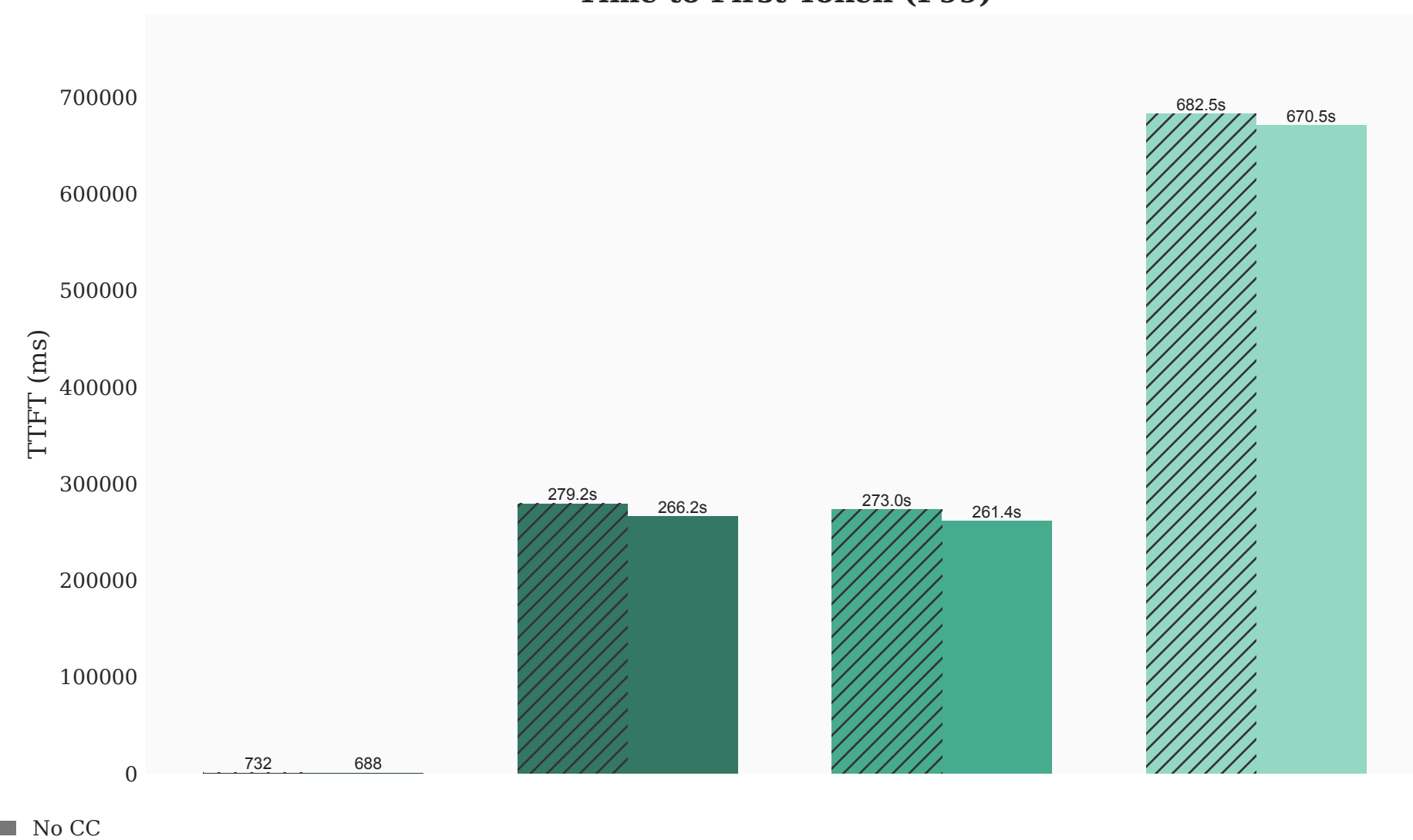## Time to First Token (P99)

Legend: ▨ CC  ▬ No CC

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

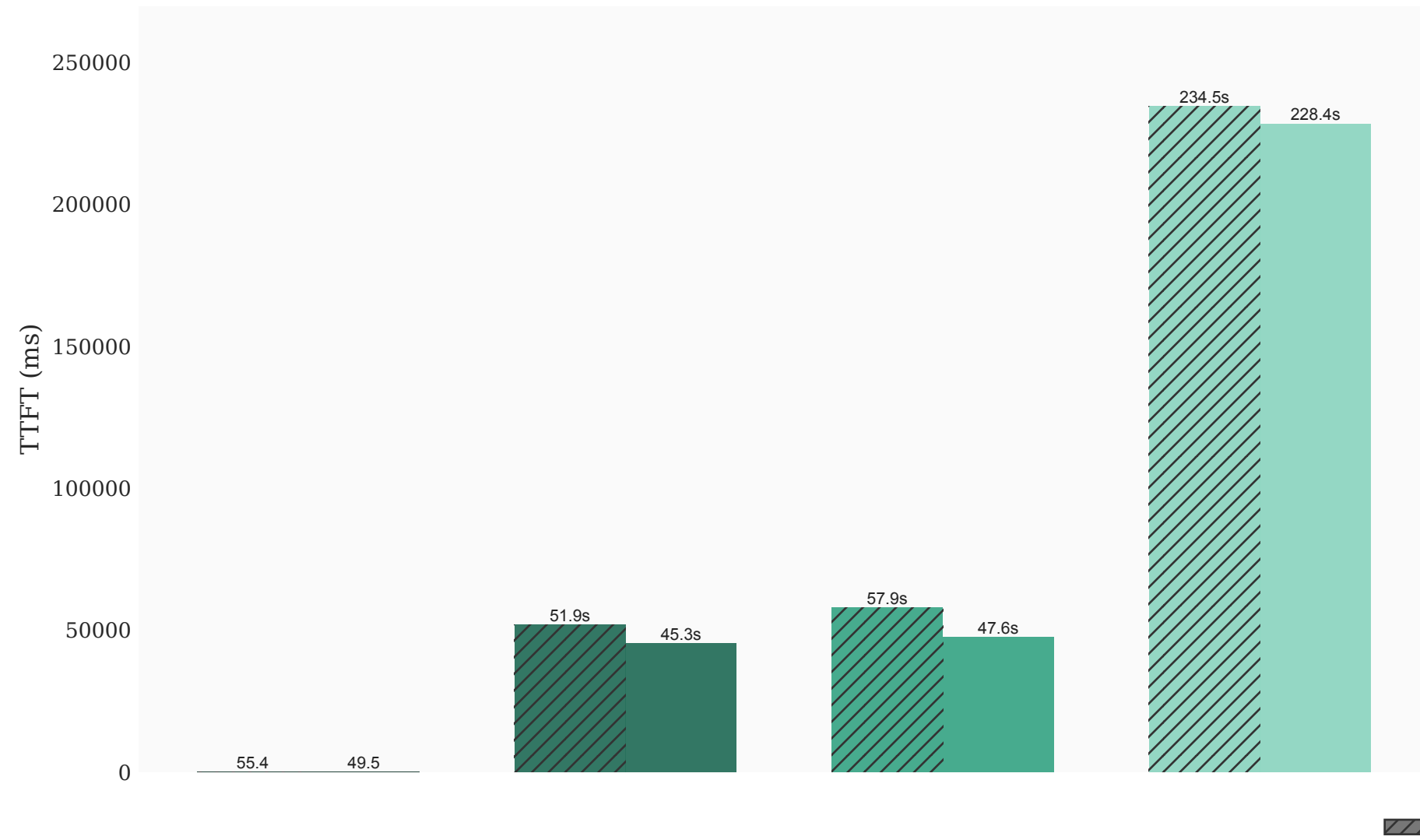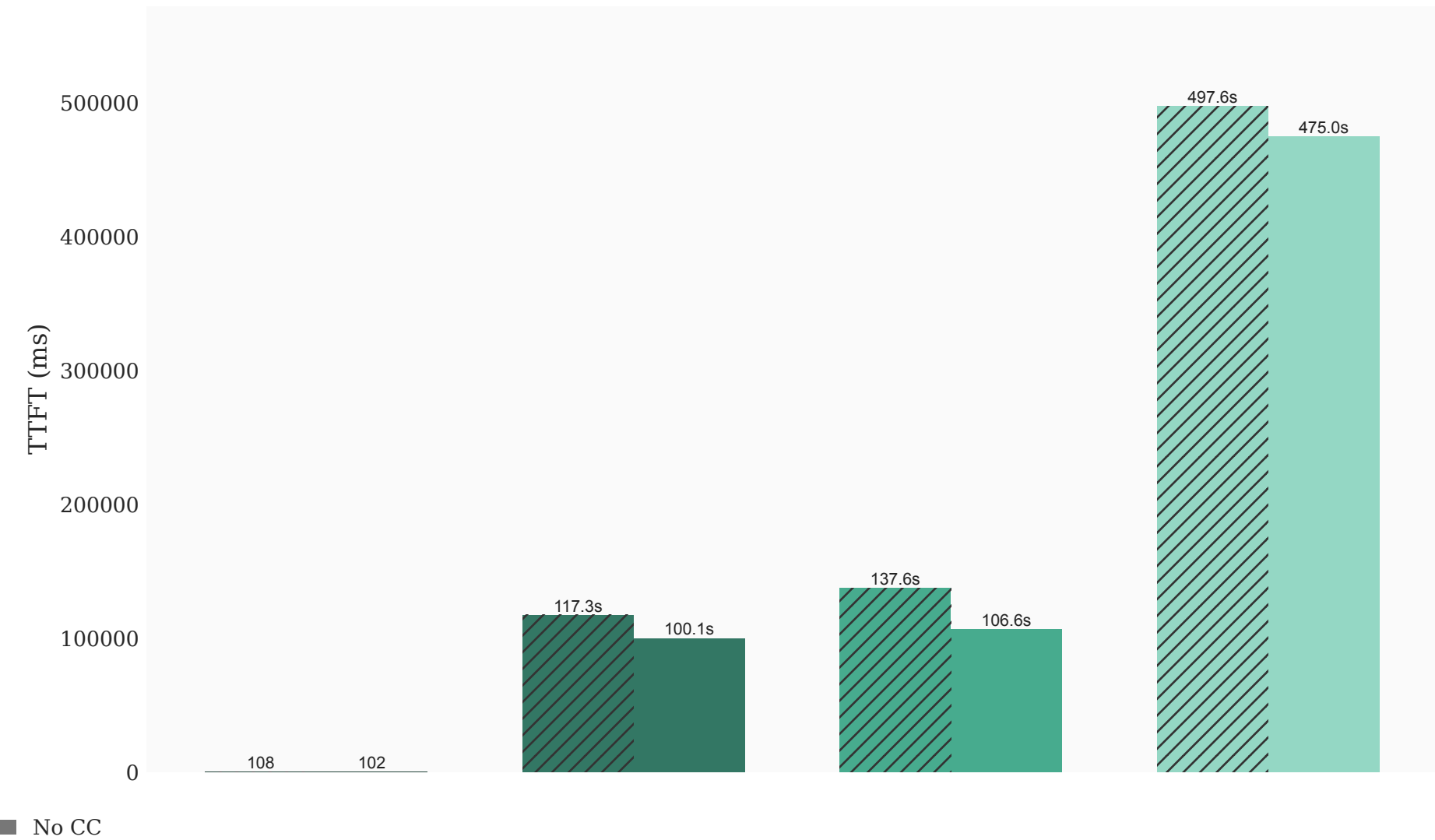# Edit 10K Characters (1 Concurrent Requests)

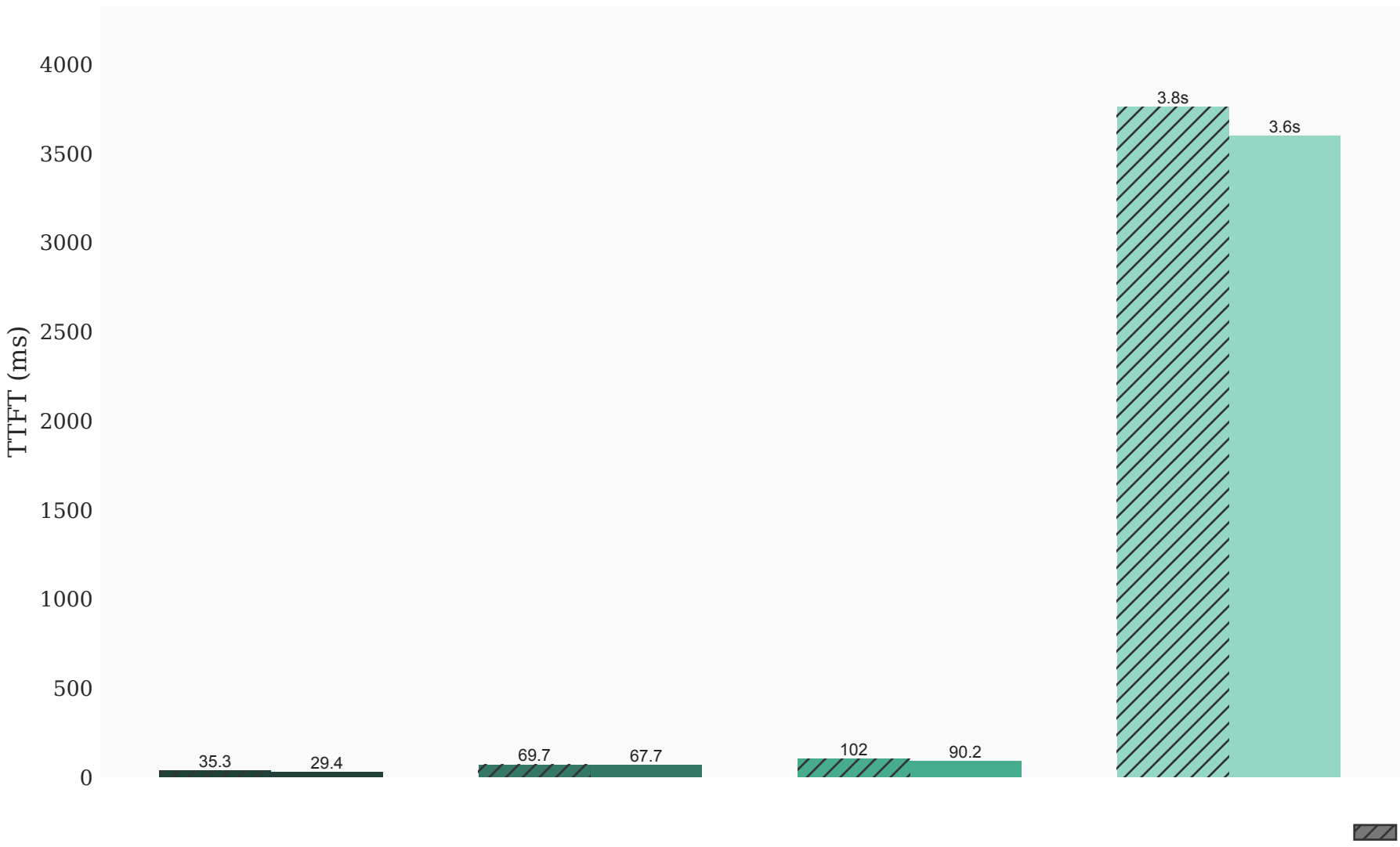## Time to First Token (Mean)



## 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:
- 55.4, 49.5
- 51.9s, 45.3s
- 57.9s, 47.6s
- 234.5s, 228.4s

Time to First Token (P99) values:
- 108, 102
- 117.3s, 100.1s
- 137.6s, 106.6s
- 497.6s, 475.0s

# Numina Math (100 Concurrent Requests)

## Time to First Token (Mean)



## Time to First Token (P99)

TTFT (ms)

Mean values: 35.3, 29.4, 69.7, 67.7, 102, 90.2, 3.8s, 3.6s

P99 values: 64.7, 48.7, 127, 124, 156, 147, 6.0s, 5.7s

CC    No CC

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

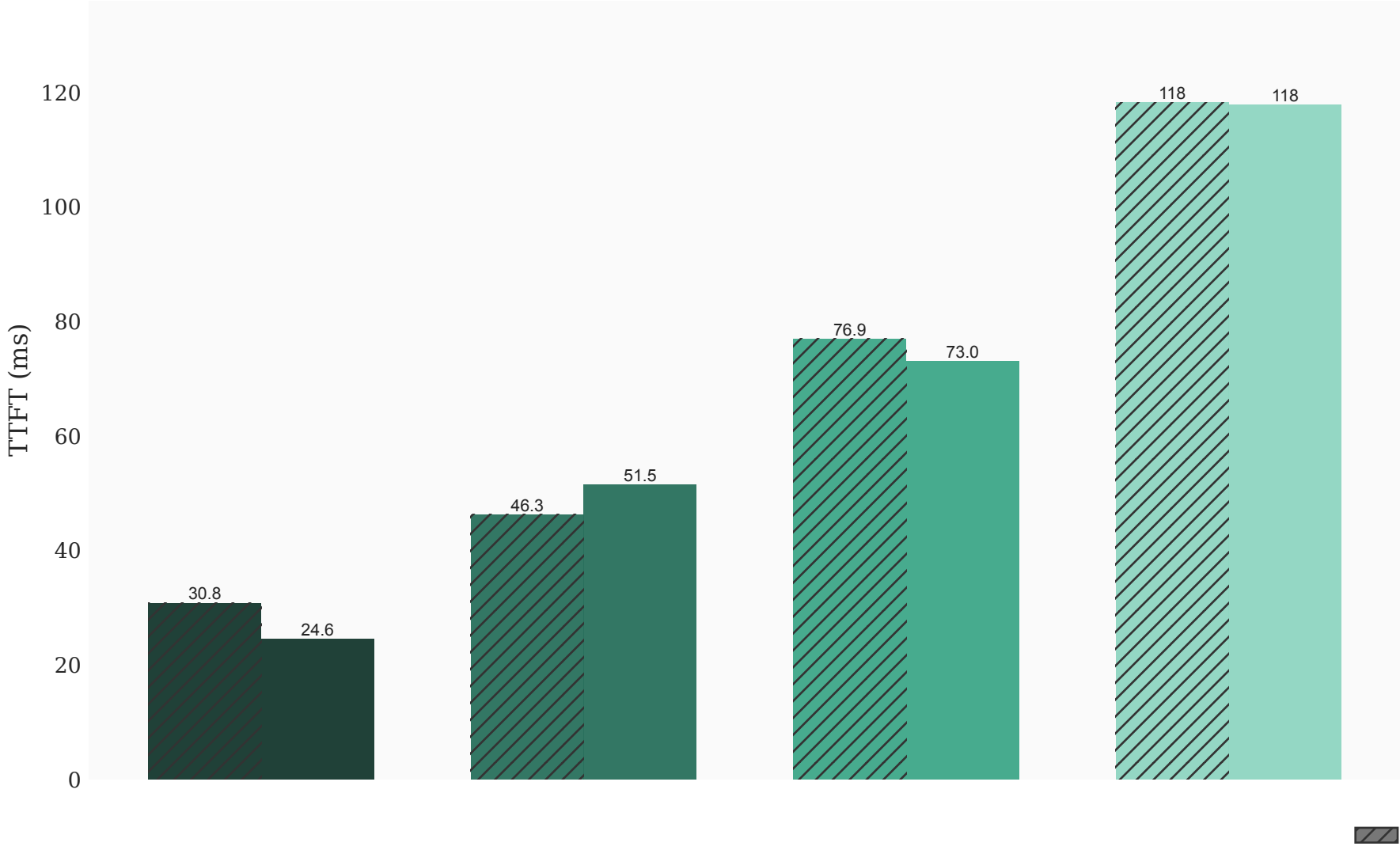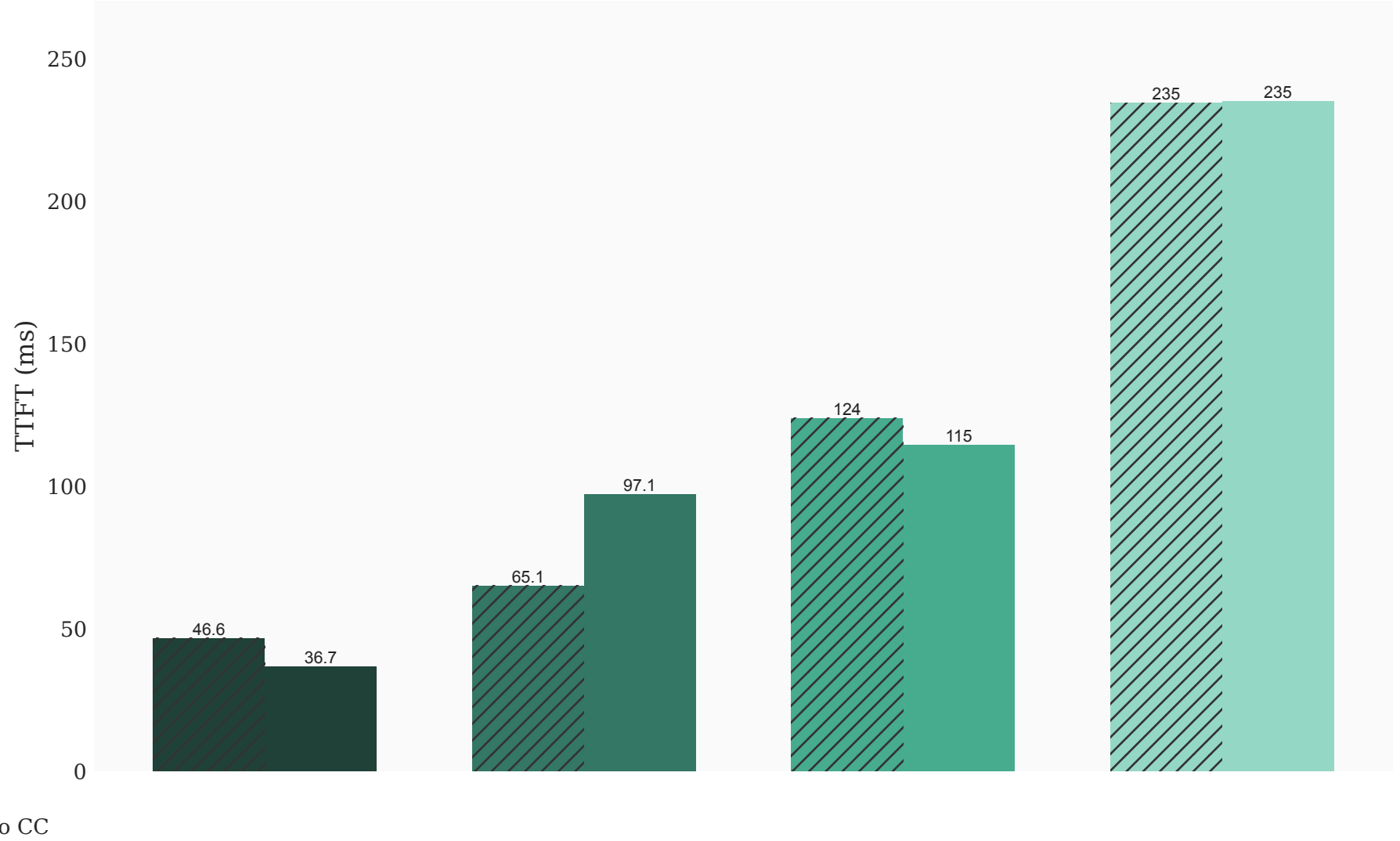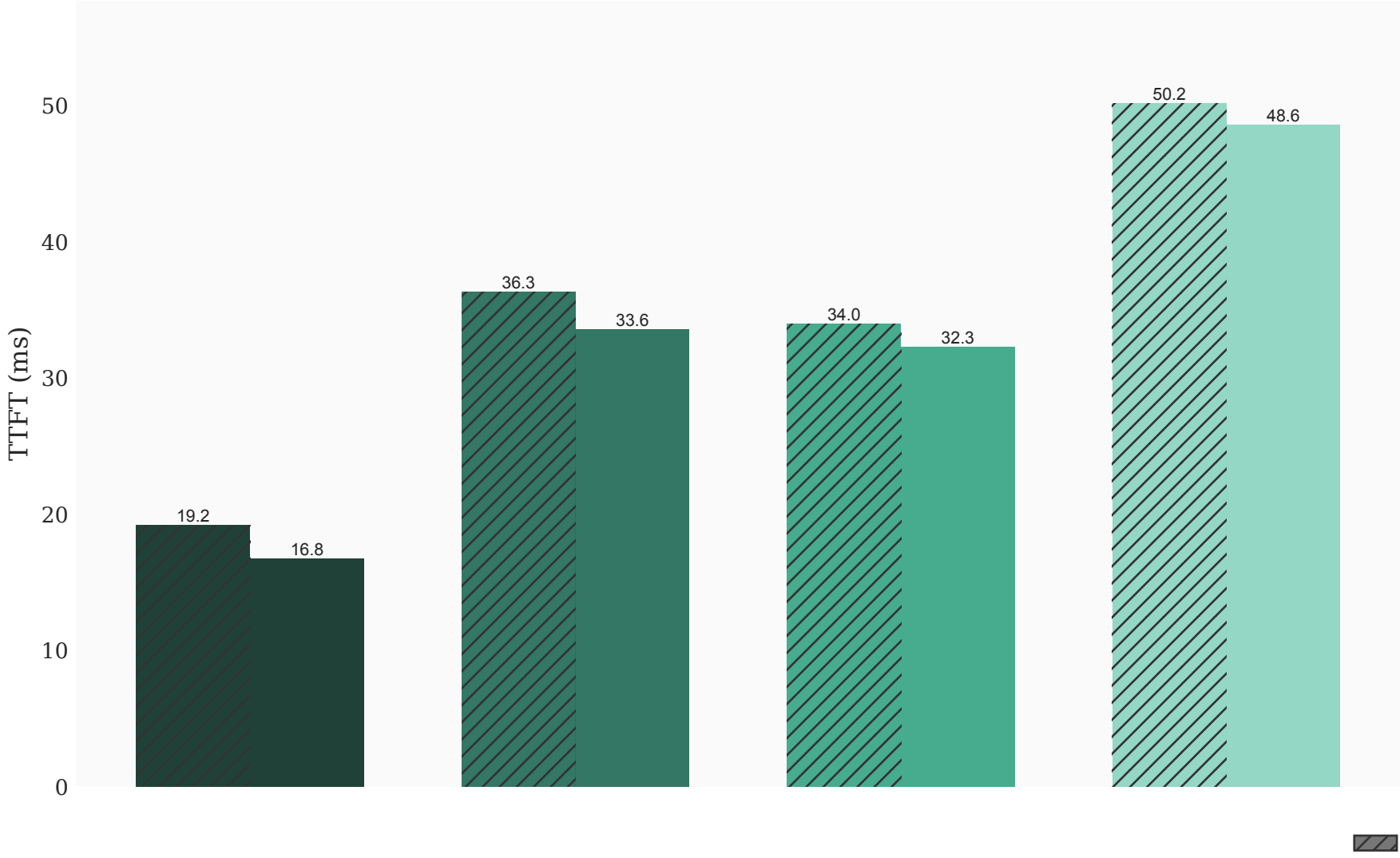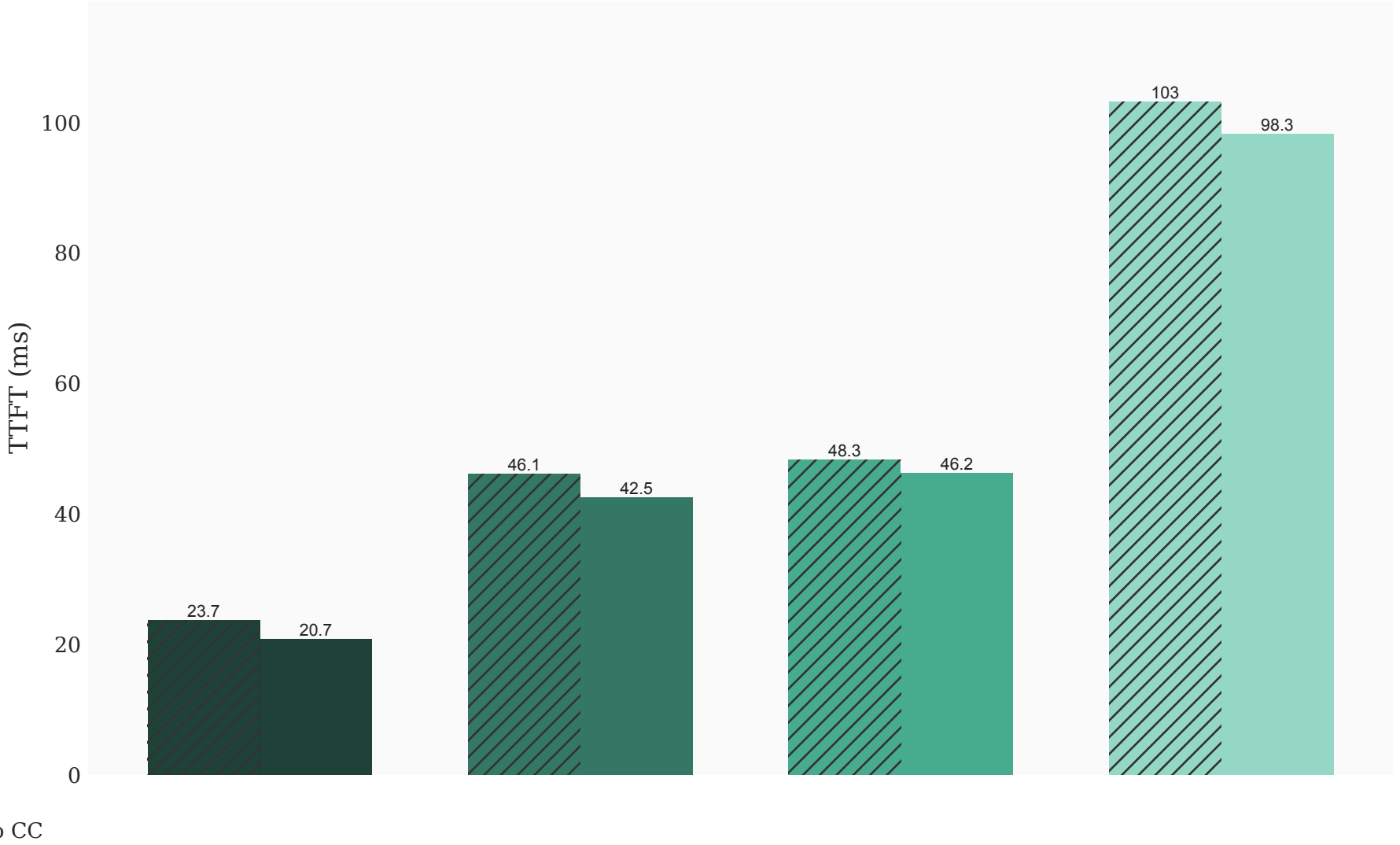LLama 3.1 8B · Mistral 3.1 24B · GPT OSS 120B · LLama 3.3 70B Int4

Mean values: 30.8, 24.6, 46.3, 51.5, 76.9, 73.0, 118, 118

P99 values: 46.6, 36.7, 65.1, 97.1, 124, 115, 235, 235

Numina Math (1 Concurrent Requests)