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

## Mean ITL



Inter-Token Latency (ms)

- LLama 3.1 8B: CC 39.3, No CC 36.4
- Mistral 3.1 24B: CC 65.6, No CC 63.9
- GPT OSS 120B: CC 69.3, No CC 69.2
- LLama 3.3 70B Int4: CC 246, No CC 251

## P99 ITL



Inter-Token Latency (ms)

- LLama 3.1 8B: CC 209, No CC 209
- Mistral 3.1 24B: CC 583, No CC 577
- GPT OSS 120B: CC 351, No CC 351
- LLama 3.3 70B Int4: CC 3.6s, No CC 3.6s

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)

## Mean ITL

Inter-Token Latency (ms)

- 20.8
- 17.4
- 65.5
- 63.9
- 73.1
- 69.1
- 244
- 250

## P99 ITL

Inter-Token Latency (ms)

- 28.3
- 23.8
- 582
- 580
- 351
- 349
- 3.6s
- 3.6s

CC    No CC

LLama 3.1 8B    Mistral 3.1 24B    GPT OSS 120B    LLama 3.3 70B Int4
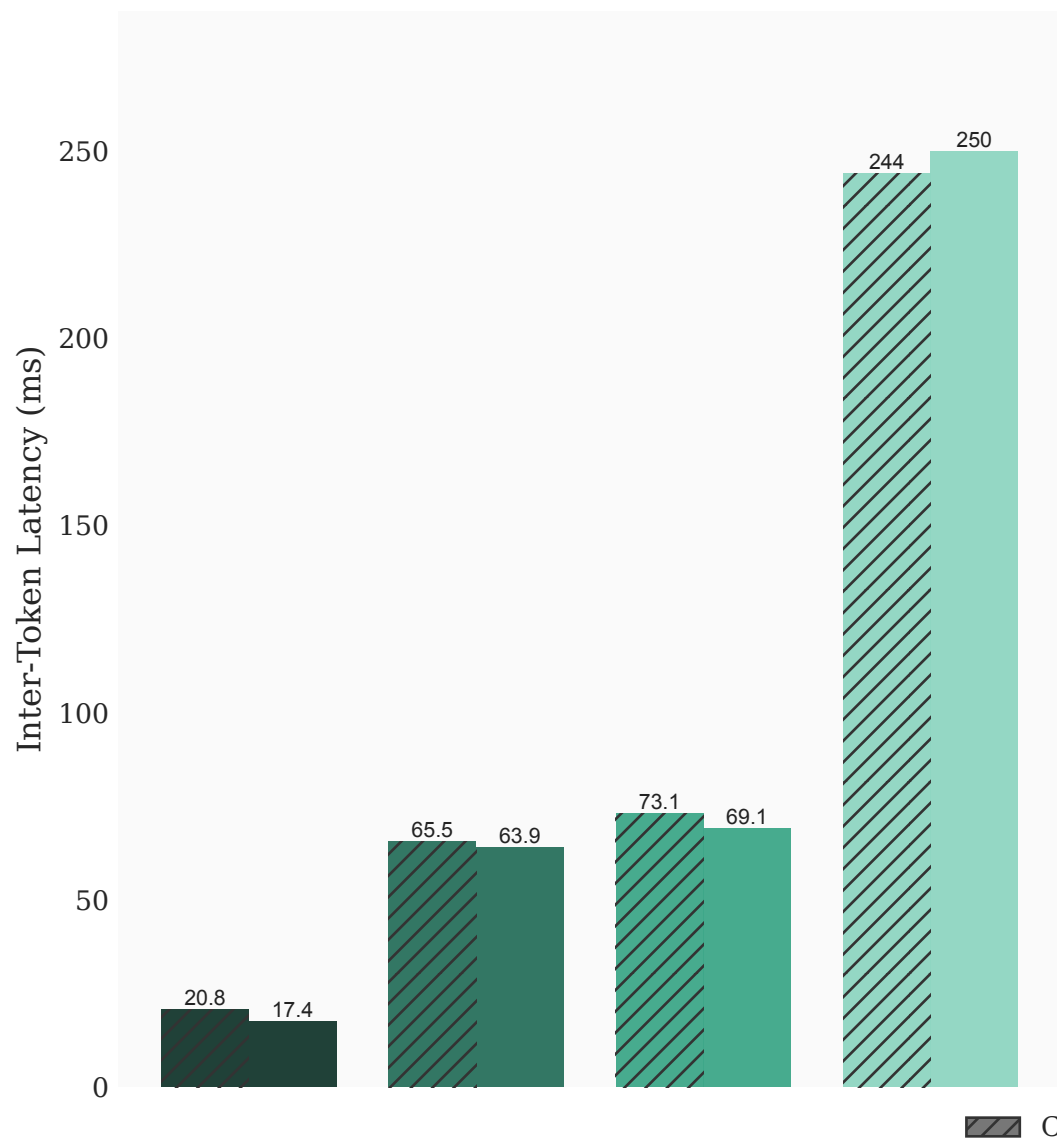
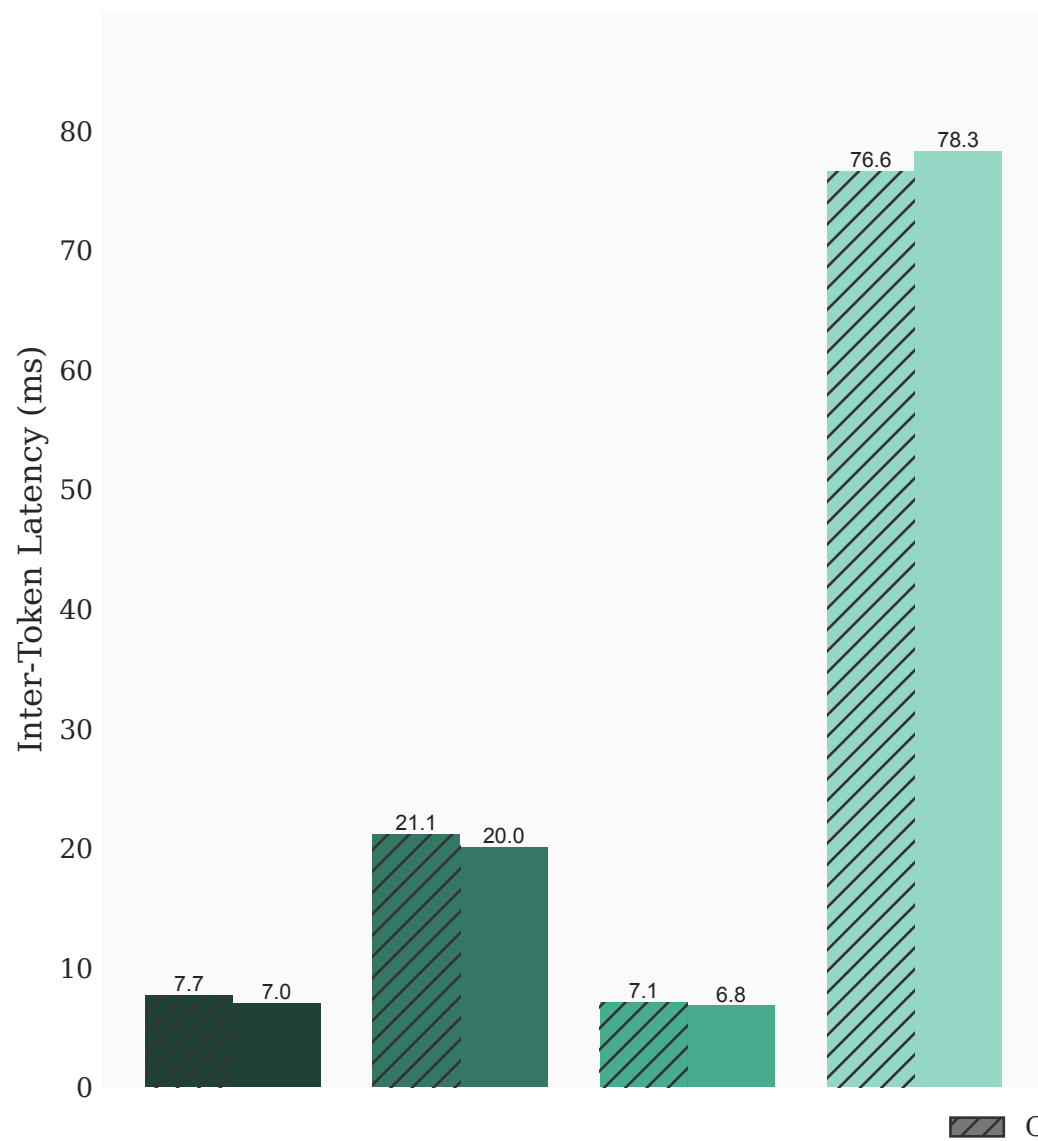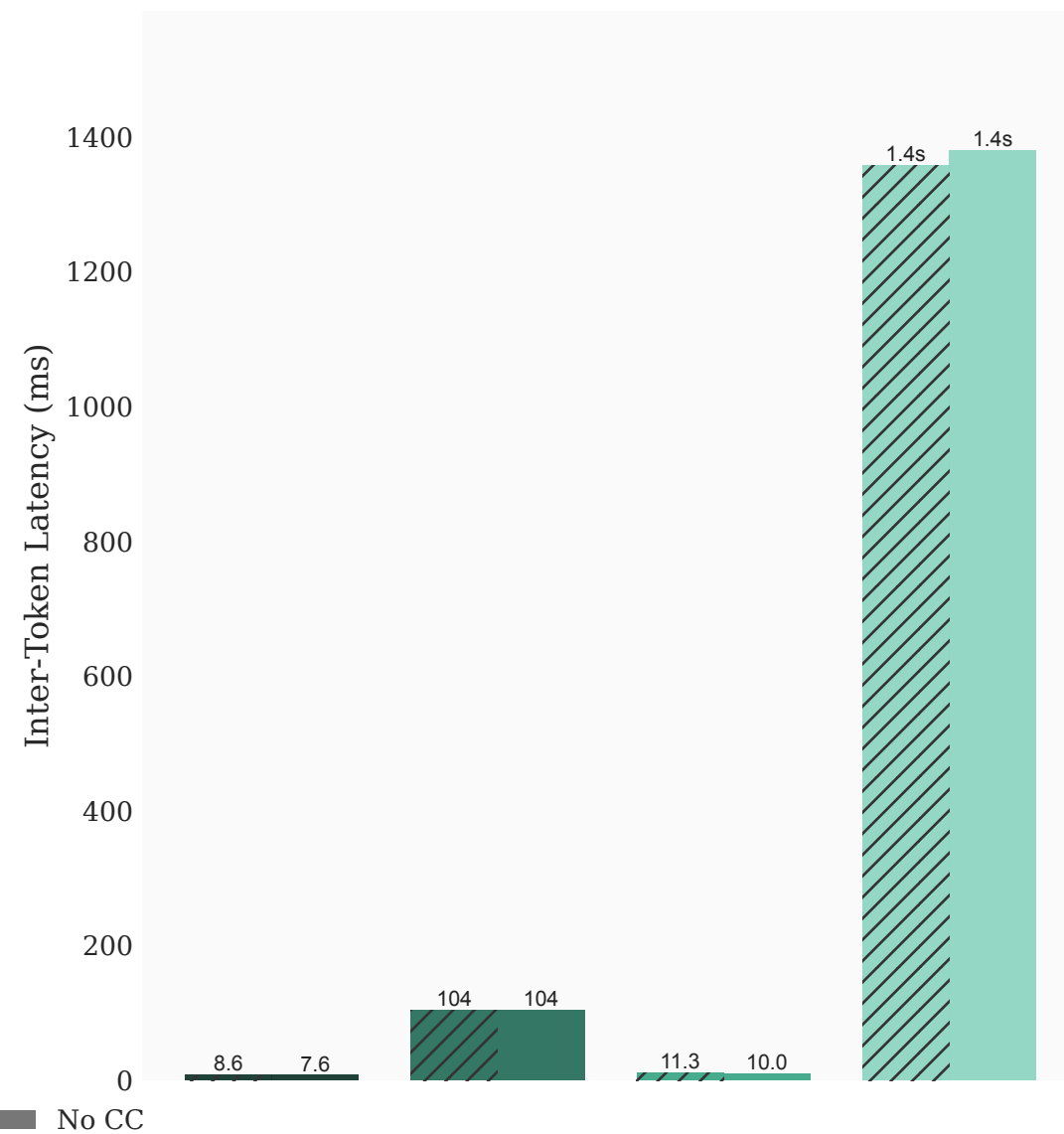**Random (1500 ⇒ 250) (1 Concurrent Requests)**

**Mean ITL**

**P99 ITL**

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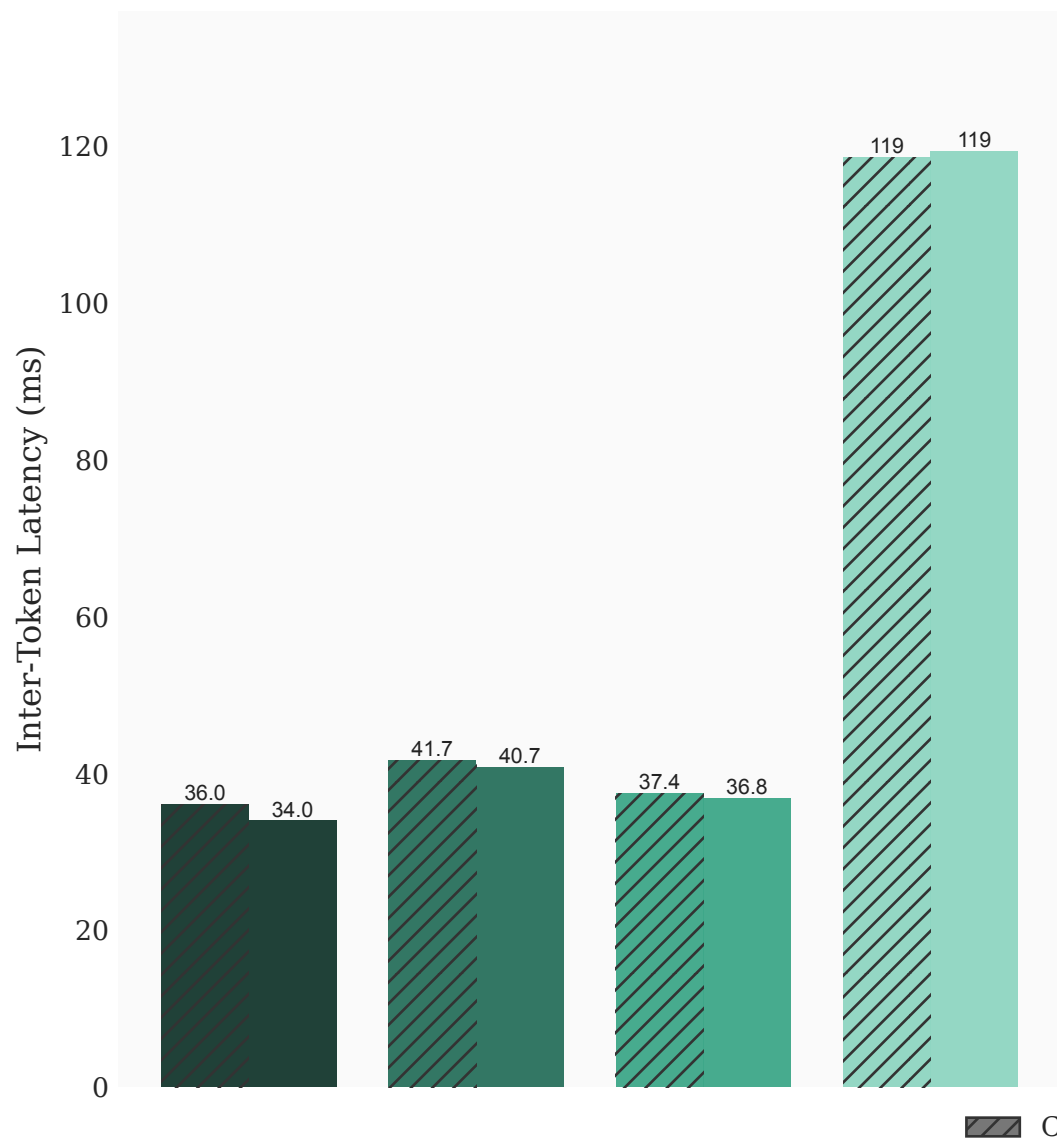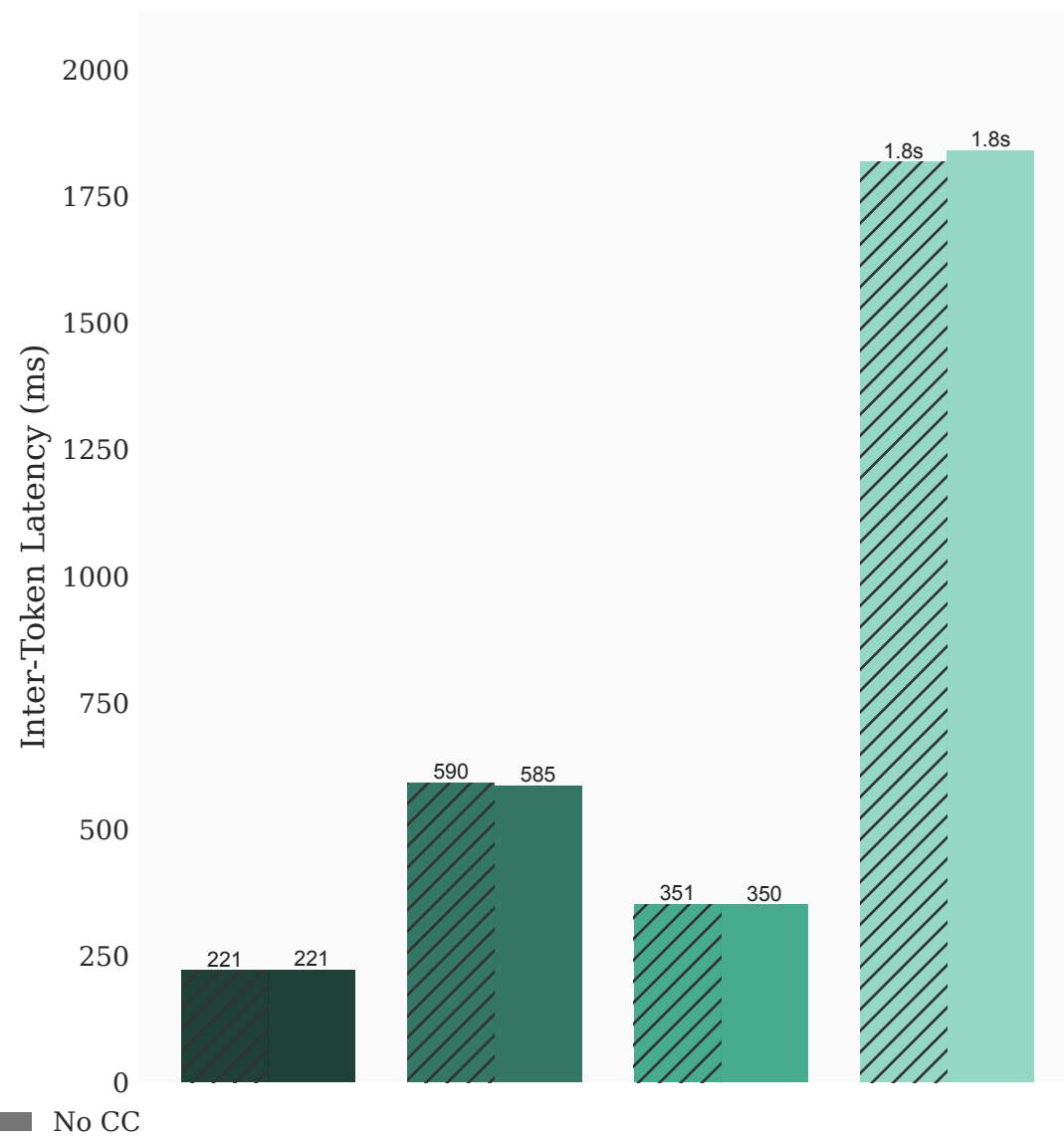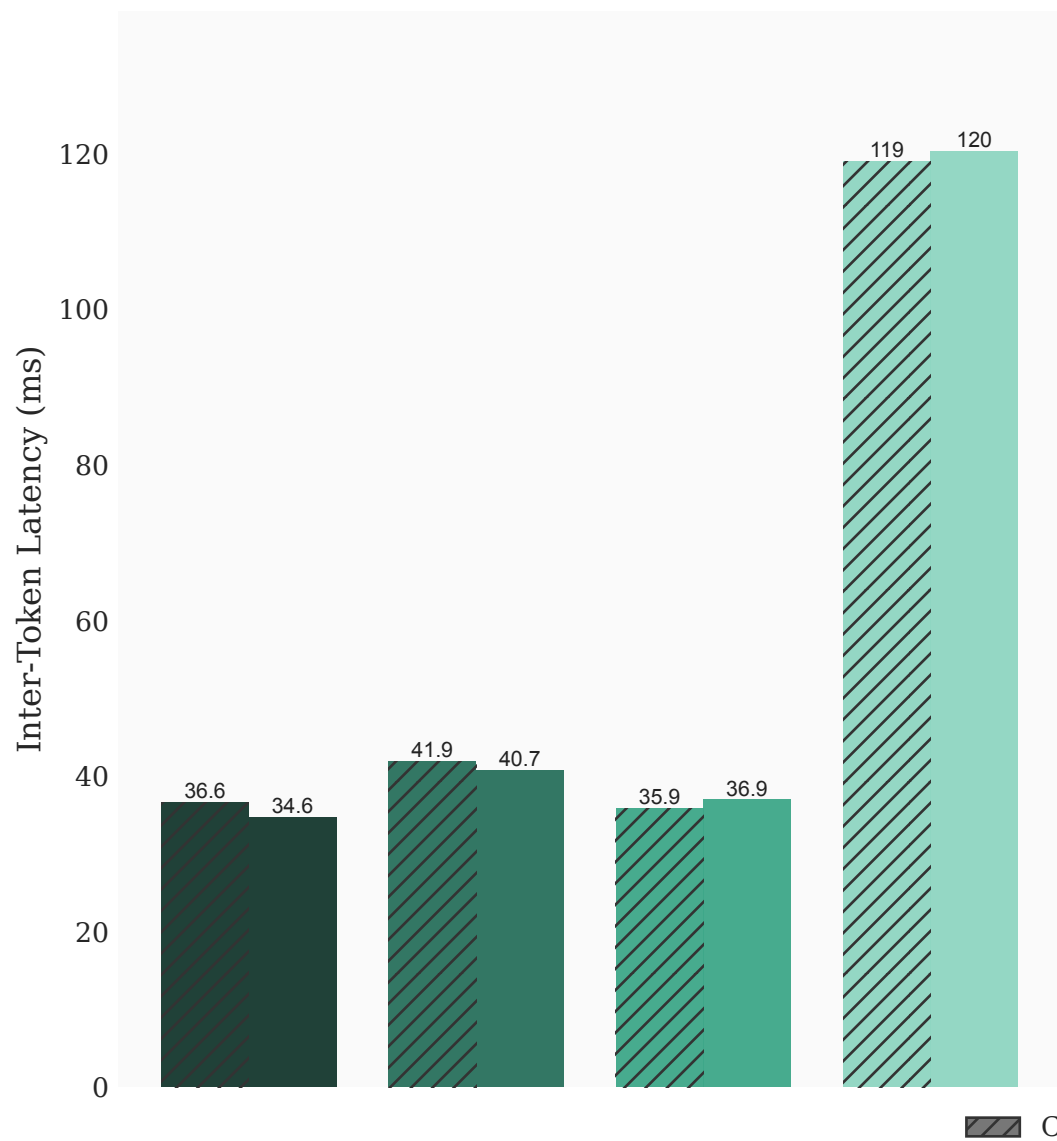
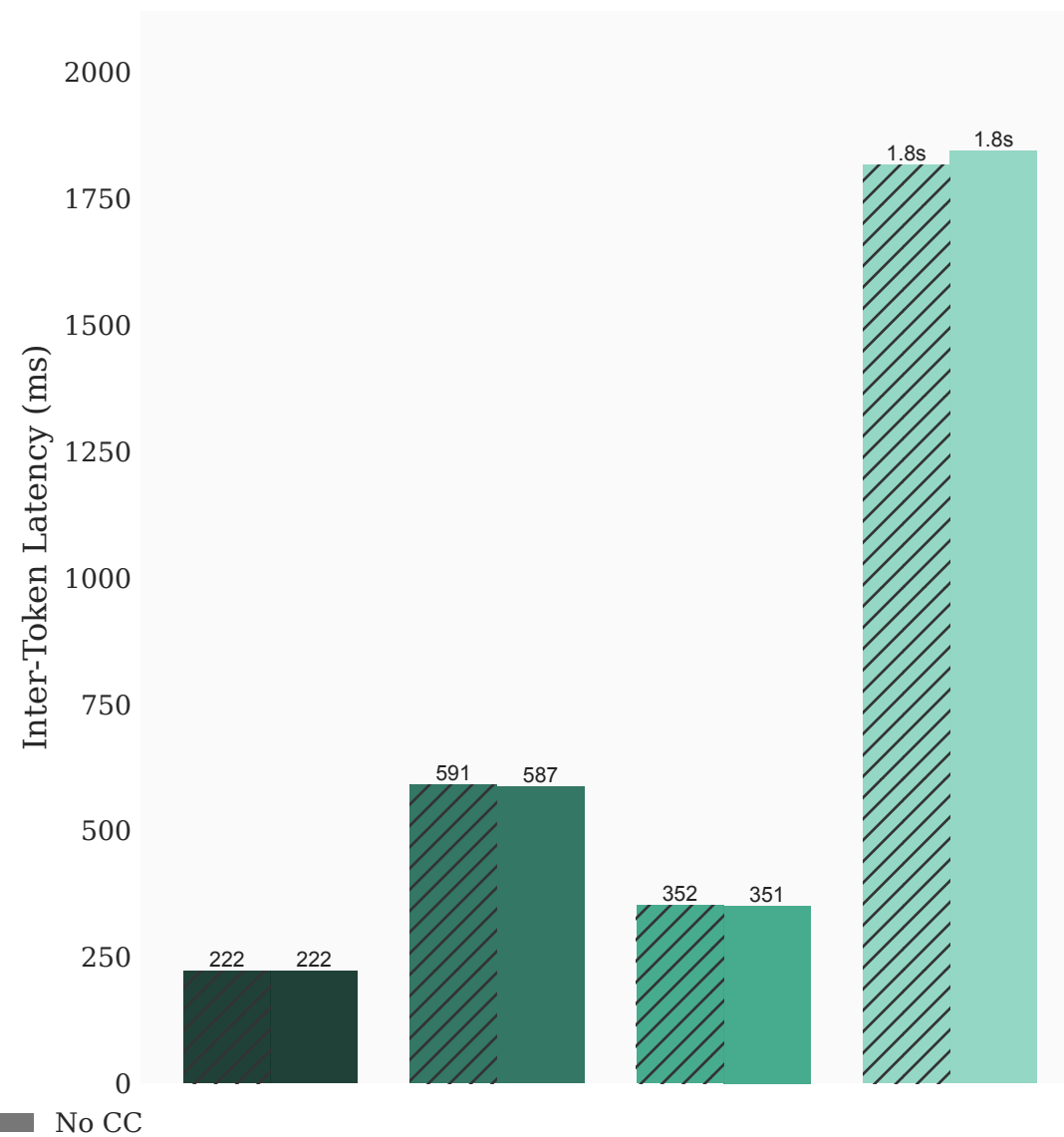LLama 3.1 8B    Mistral 3.1 24B    GPT OSS 120B    LLama 3.3 70B Int4

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

## Mean ITL



## P99 ITL



Legend: CC | No CC

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

Mean ITL values:
- 36.0, 34.0
- 41.7, 40.7
- 37.4, 36.8
- 119, 119

P99 ITL values:
- 221, 221
- 590, 585
- 351, 350
- 1.8s, 1.8s

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

## Mean ITL

**Inter-Token Latency (ms)**

- LLama 3.1 8B: CC 36.6, No CC 34.6
- Mistral 3.1 24B: CC 41.9, No CC 40.7
- GPT OSS 120B: CC 35.9, No CC 36.9
- LLama 3.3 70B Int4: CC 119, No CC 120

## P99 ITL

**Inter-Token Latency (ms)**

- LLama 3.1 8B: CC 222, No CC 222
- Mistral 3.1 24B: CC 591, No CC 587
- GPT OSS 120B: CC 352, No CC 351
- LLama 3.3 70B Int4: CC 1.8s, No CC 1.8s

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)

## Mean ITL



Inter-Token Latency (ms)

- LLama 3.1 8B: 11.1 (CC), 10.0 (No CC)
- Mistral 3.1 24B: 38.1 (CC), 35.6 (No CC)
- GPT OSS 120B: 10.6 (CC), 9.9 (No CC)
- LLama 3.3 70B Int4: 119 (CC), 121 (No CC)

## P99 ITL



Inter-Token Latency (ms)

- LLama 3.1 8B: 95.6 (CC), 12.6 (No CC)
- Mistral 3.1 24B: 291 (CC), 291 (No CC)
- GPT OSS 120B: 184 (CC), 182 (No CC)
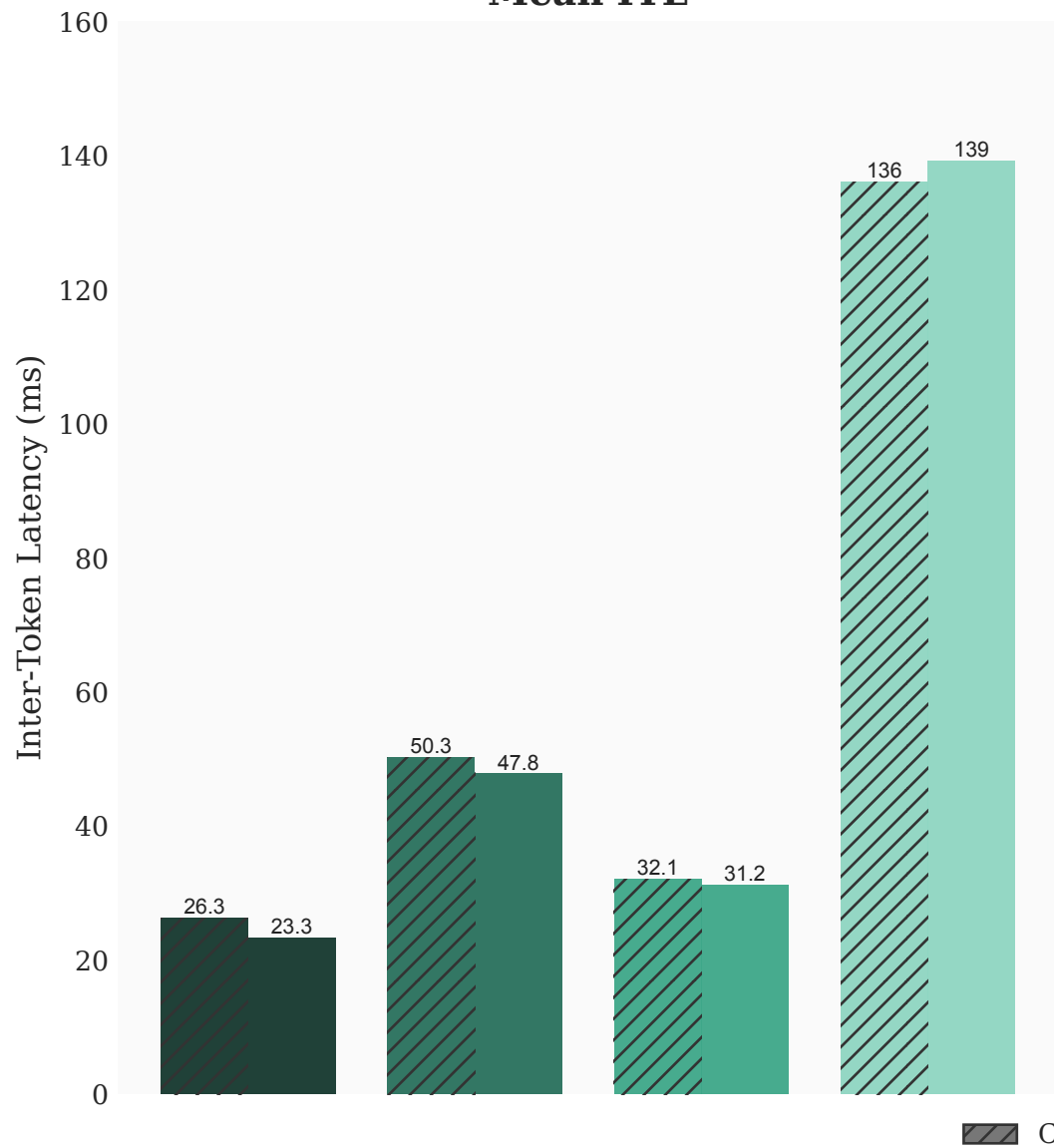- LLama 3.3 70B Int4: 1.8s (CC), 1.8s (No CC)

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)

## Mean ITL



## P99 ITL

Legend: CC | No CC

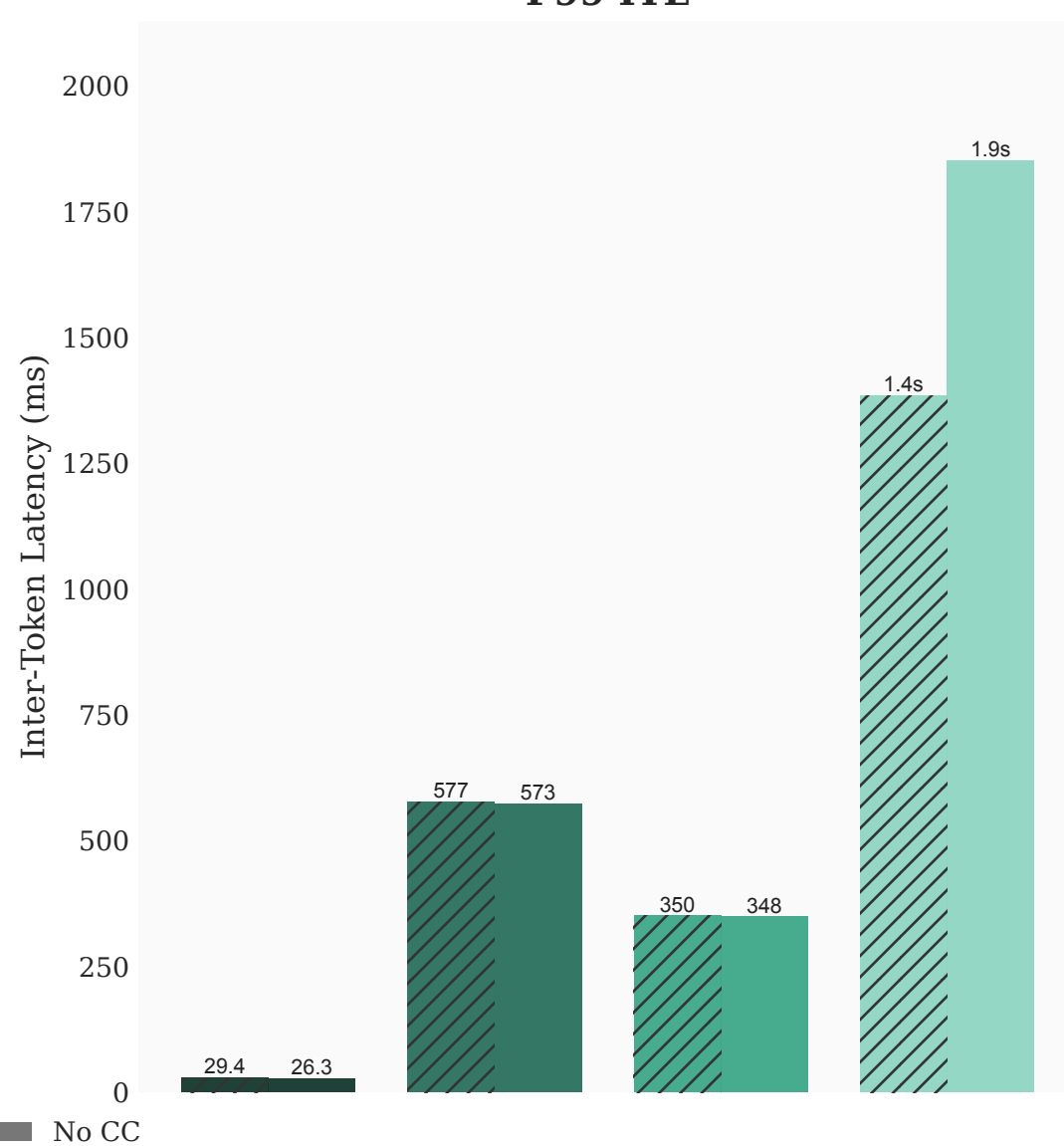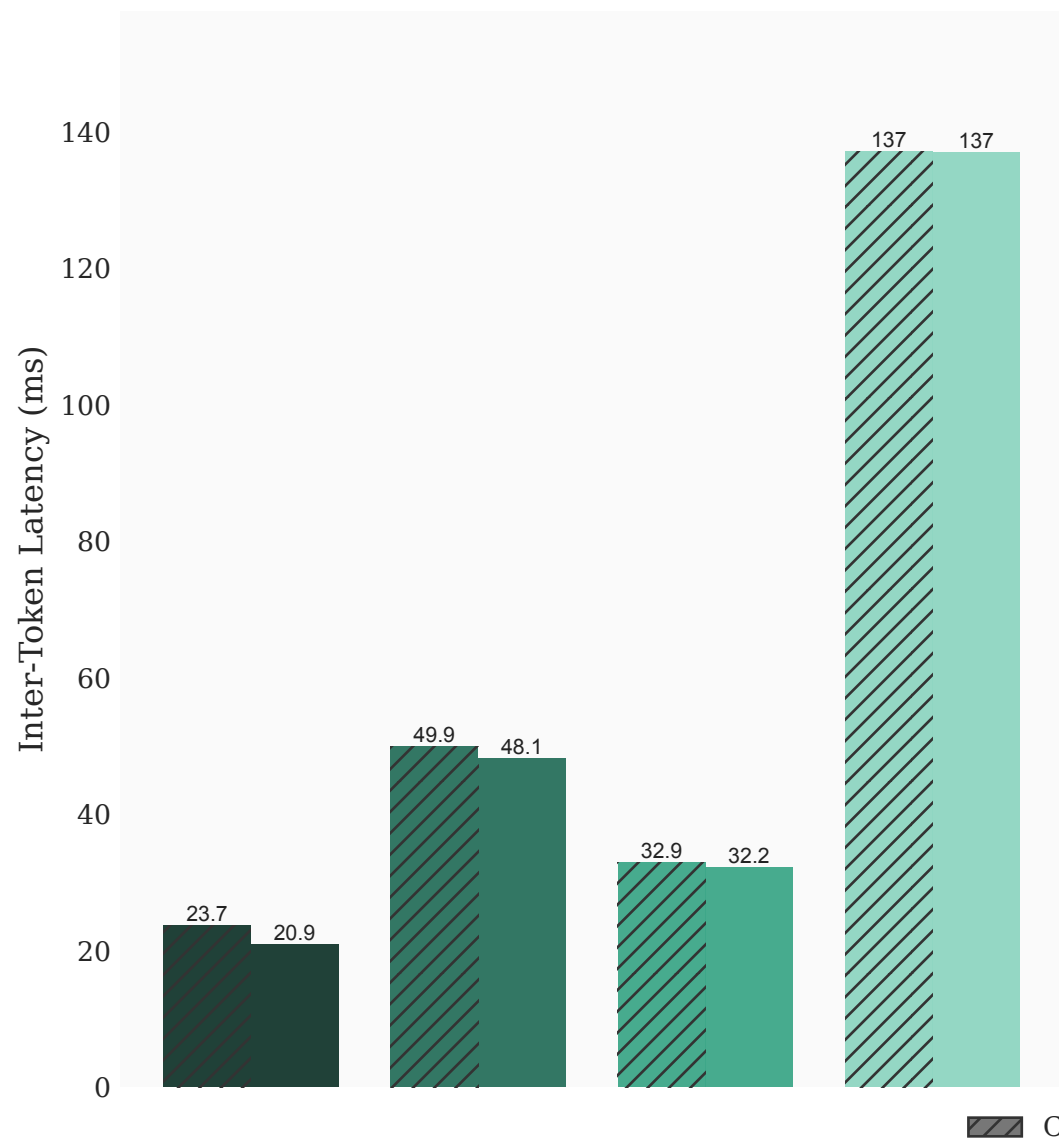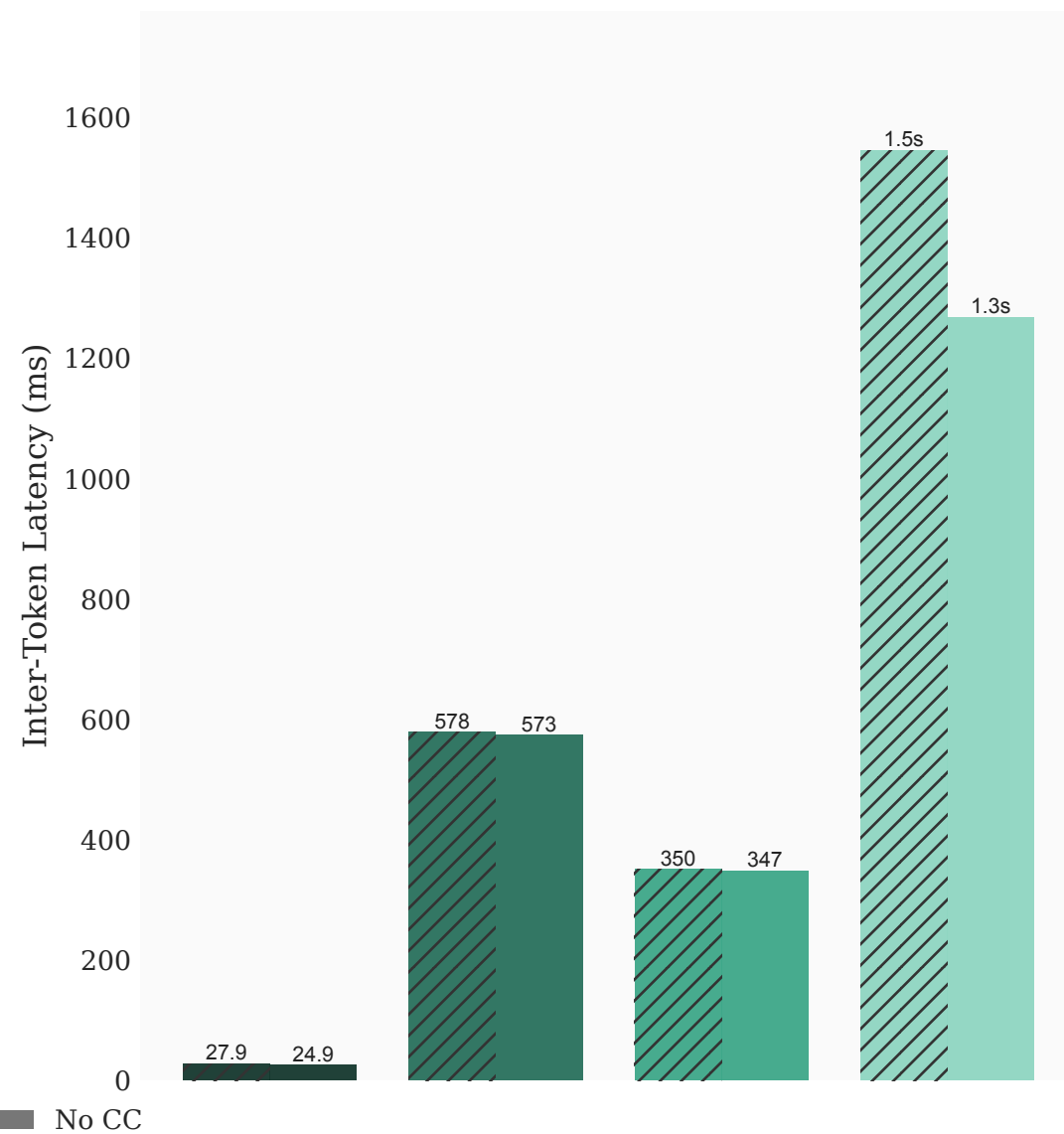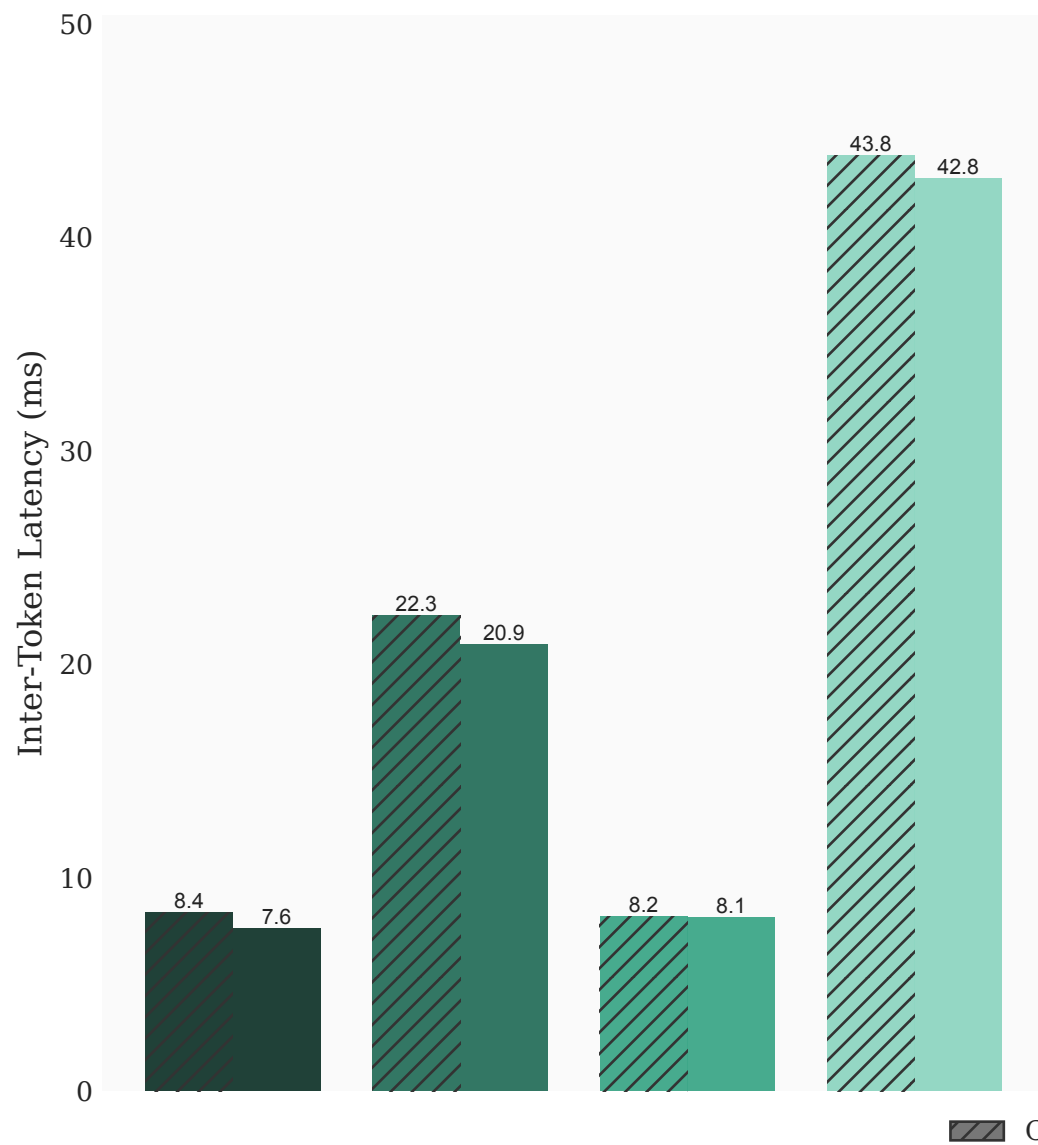LLama 3.1 8B | Mistral 3.1 24B | GPT OSS 120B | LLama 3.3 70B Int4

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

## Mean ITL



## P99 ITL



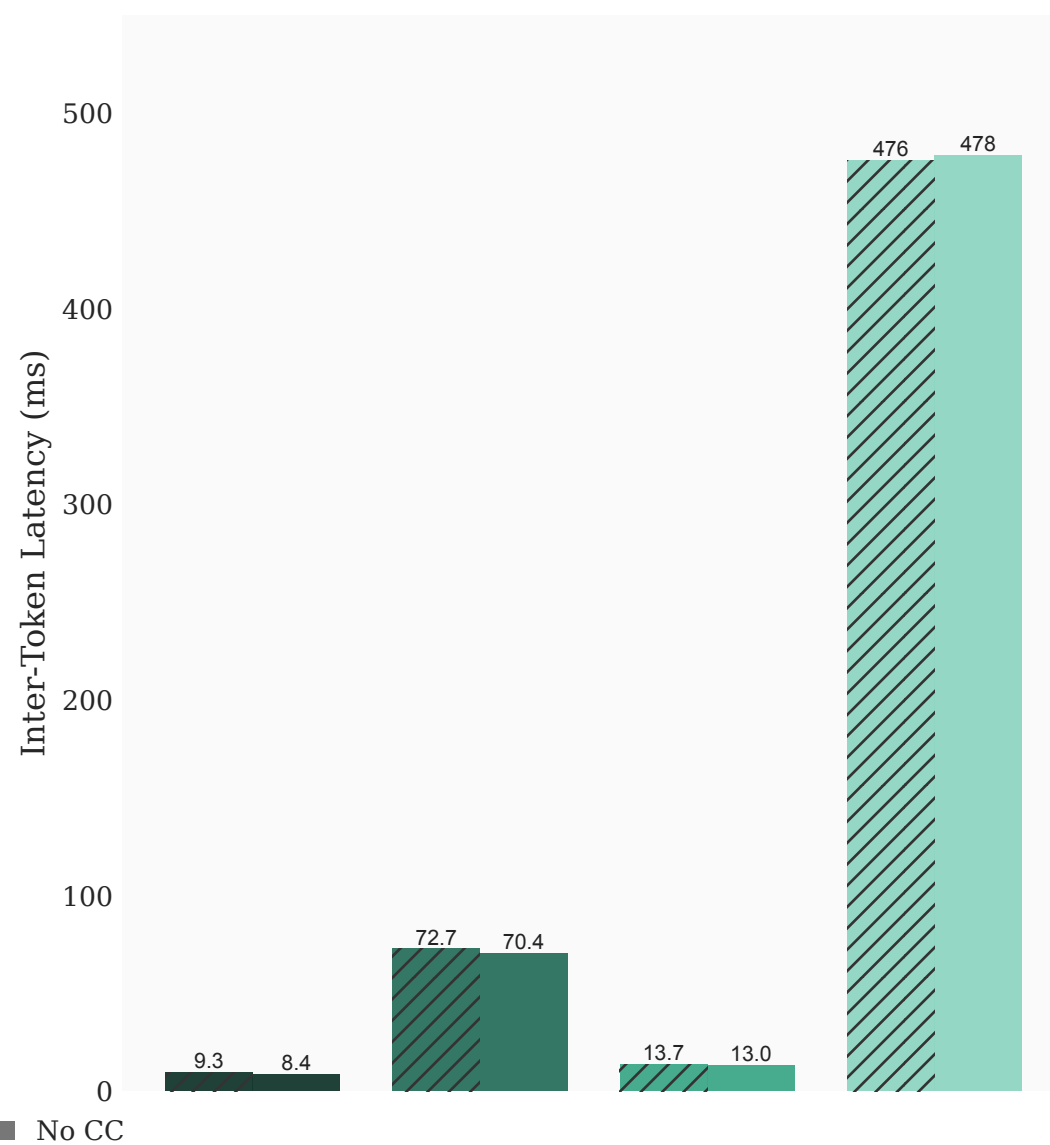Legend: CC (hatched), No CC (solid)

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

Mean ITL values:
- LLama 3.1 8B: CC 23.7, No CC 20.9
- Mistral 3.1 24B: CC 49.9, No CC 48.1
- GPT OSS 120B: CC 32.9, No CC 32.2
- LLama 3.3 70B Int4: CC 137, No CC 137

P99 ITL values:
- LLama 3.1 8B: CC 27.9, No CC 24.9
- Mistral 3.1 24B: CC 578, No CC 573
- GPT OSS 120B: CC 350, No CC 347
- LLama 3.3 70B Int4: CC 1.5s, No CC 1.3s

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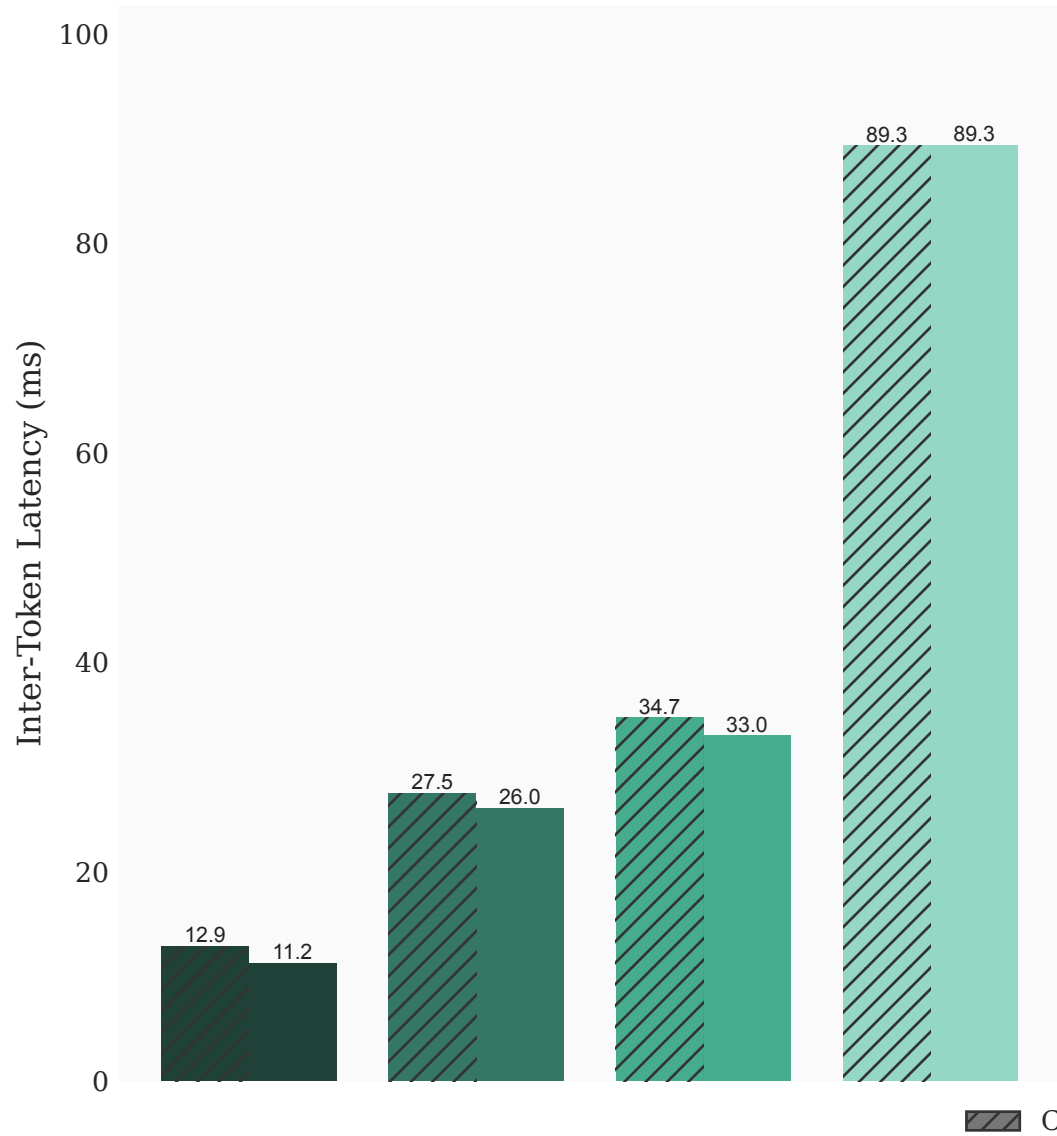
## Mean ITL



## P99 ITL

Legend: CC, No CC

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

# ShareGPT (100 Concurrent Requests)

## Mean ITL

Inter-Token Latency (ms)

- LLama 3.1 8B: CC 12.9, No CC 11.2
- Mistral 3.1 24B: CC 27.5, No CC 26.0
- GPT OSS 120B: CC 34.7, No CC 33.0
- LLama 3.3 70B Int4: CC 89.3, No CC 89.3

## P99 ITL

Inter-Token Latency (ms)

- LLama 3.1 8B: CC 39.2, No CC 33.8
- Mistral 3.1 24B: CC 26.6, No CC 231
- GPT OSS 120B: CC 212, No CC 238
- LLama 3.3 70B Int4: CC 1.7s, No CC 1.7s

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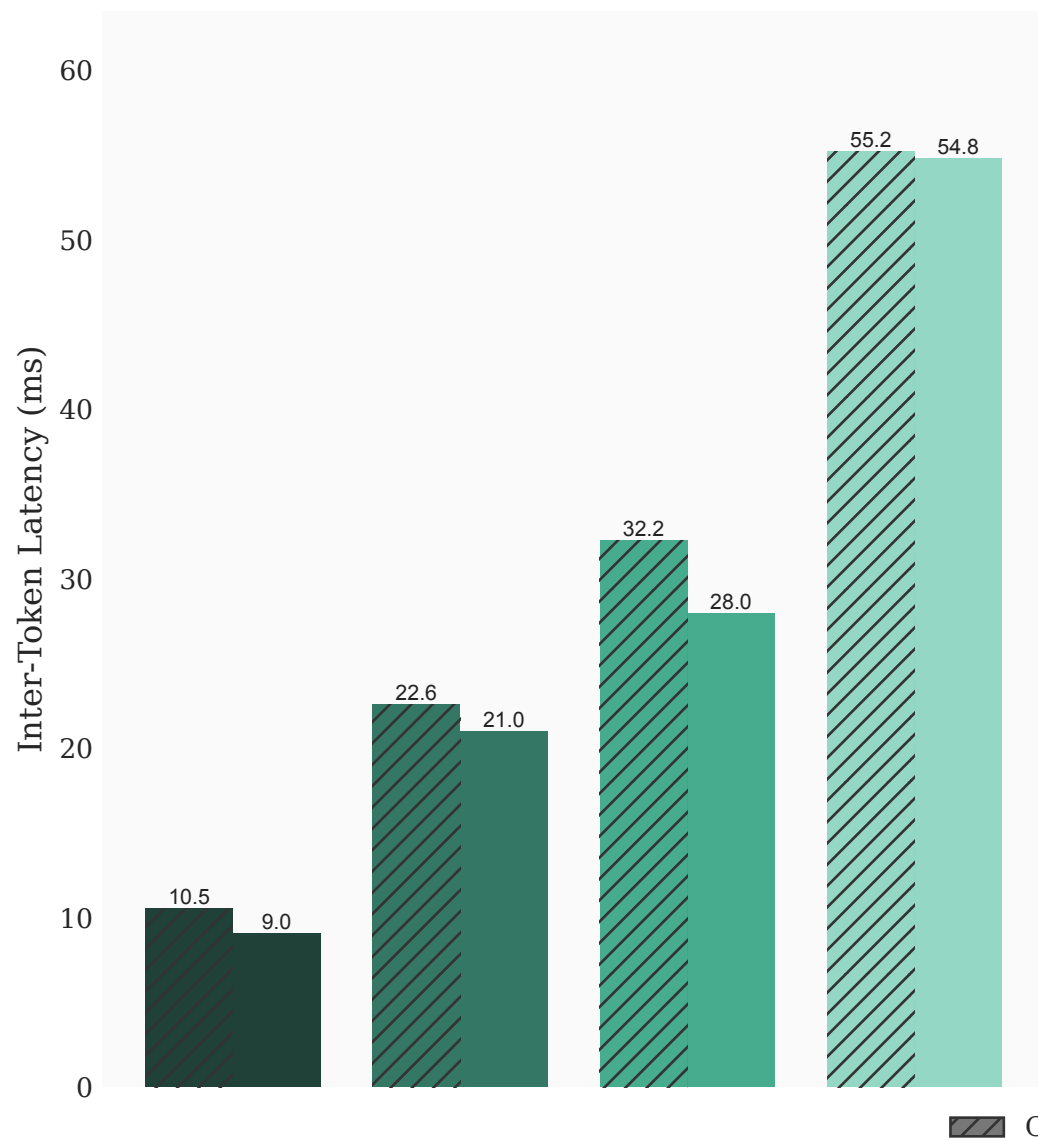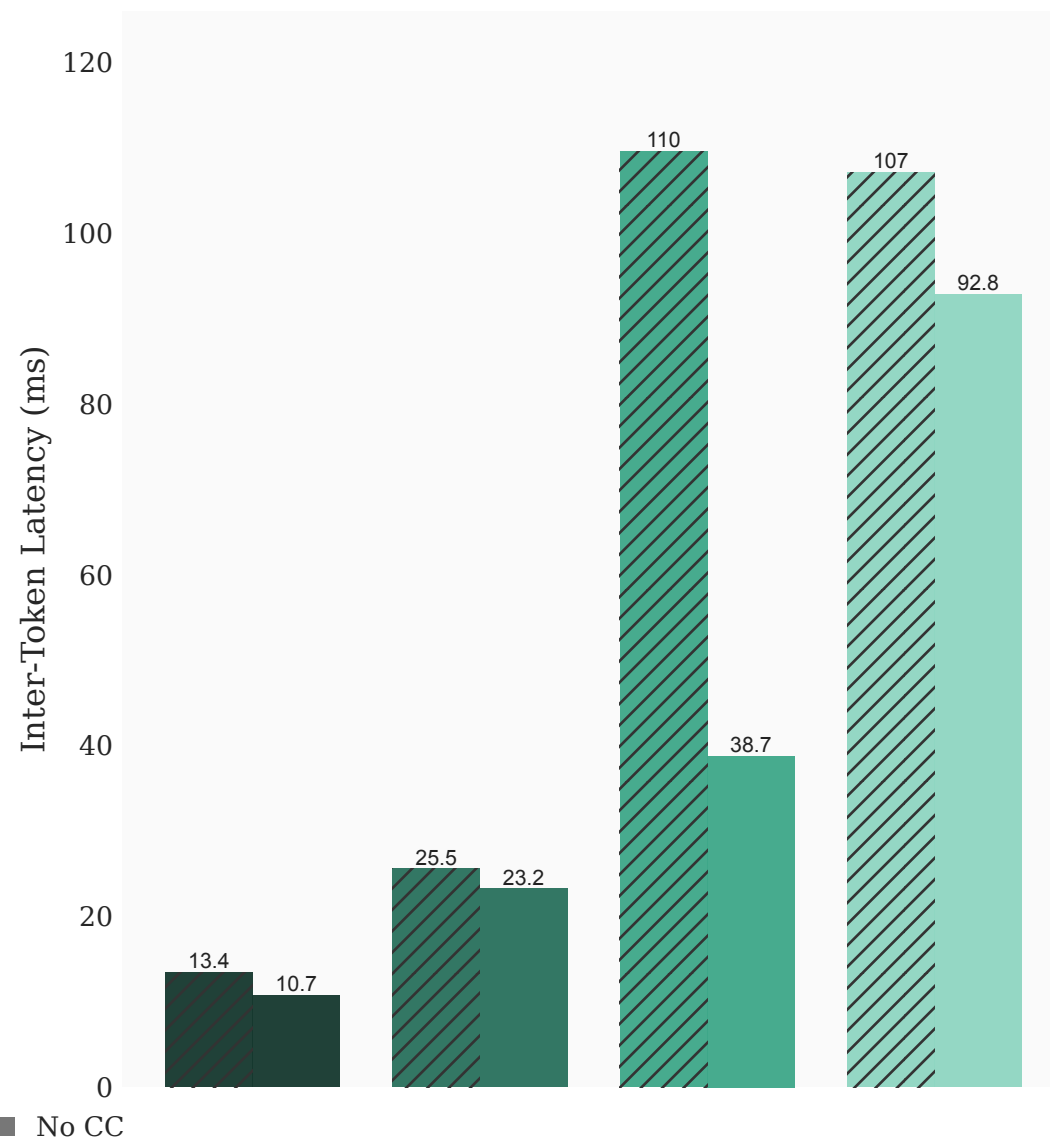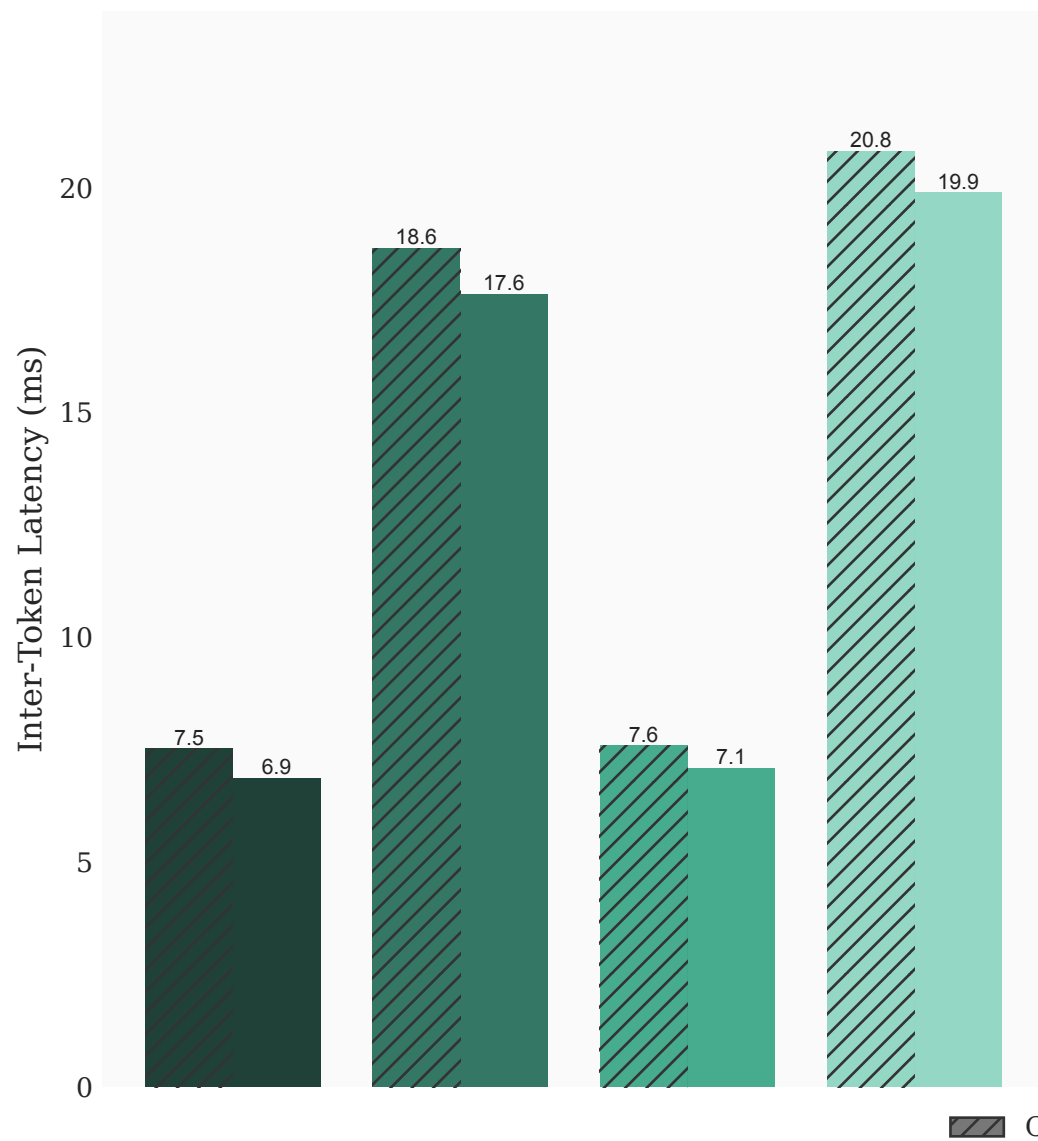
LLama 3.1 8B   Mistral 3.1 24B   GPT OSS 120B   LLama 3.3 70B Int4

# ShareGPT (50 Concurrent Requests)

## Mean ITL



## P99 ITL

Legend:
- CC
- No CC

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

Mean ITL values:
- LLama 3.1 8B: CC 10.5, No CC 9.0
- Mistral 3.1 24B: CC 22.6, No CC 21.0
- GPT OSS 120B: CC 32.2, No CC 28.0
- LLama 3.3 70B Int4: CC 55.2, No CC 54.8

P99 ITL values:
- LLama 3.1 8B: CC 13.4, No CC 10.7
- Mistral 3.1 24B: CC 25.5, No CC 23.2
- GPT OSS 120B: CC 110, No CC 38.7
- LLama 3.3 70B Int4: CC 107, No CC 92.8

# ShareGPT (1 Concurrent Requests)

## Mean ITL



## P99 ITL

Mean ITL values:
- LLama 3.1 8B: CC 7.5, No CC 6.9
- Mistral 3.1 24B: CC 18.6, No CC 17.6
- GPT OSS 120B: CC 7.6, No CC 7.1
- LLama 3.3 70B Int4: CC 20.8, No CC 19.9

P99 ITL values:
- LLama 3.1 8B: CC 8.2, No CC 7.3
- Mistral 3.1 24B: CC 20.0, No CC 18.7
- GPT OSS 120B: CC 11.5, No CC 11.1
- LLama 3.3 70B Int4: CC 23.1, No CC 21.8

Y-axis: Inter-Token Latency (ms)

Legend: CC | No CC

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

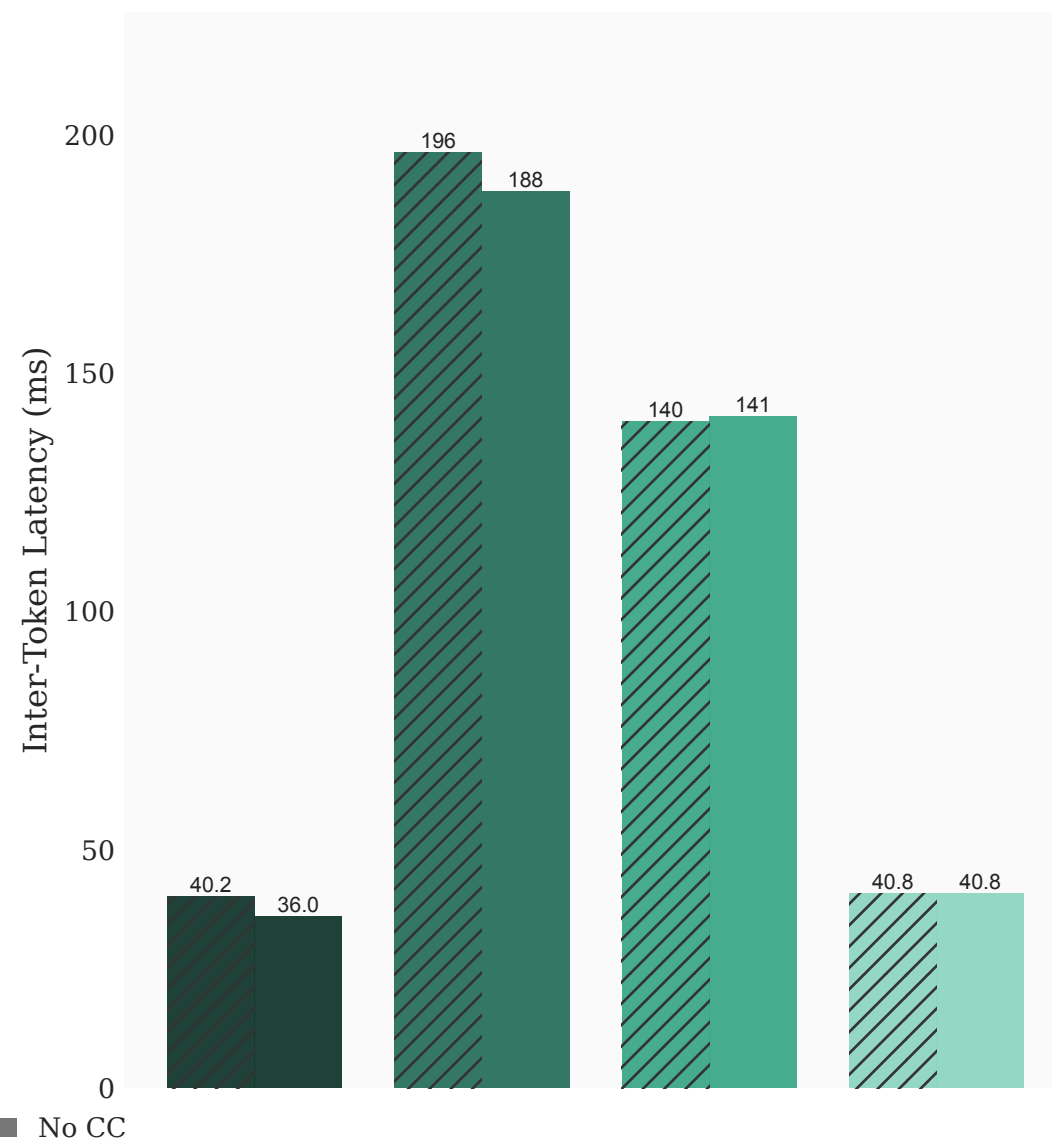# Edit 10K Characters (100 Concurrent Requests)

## Mean ITL

Inter-Token Latency (ms)

- 34.0
- 31.2
- 38.0
- 36.7
- 28.0
- 28.6
- 57.0
- 57.0

## P99 ITL

Inter-Token Latency (ms)

- 57.5
- 54.6
- 202
- 206
- 136
- 139
- 41.0
- 40.8

CC   No CC

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

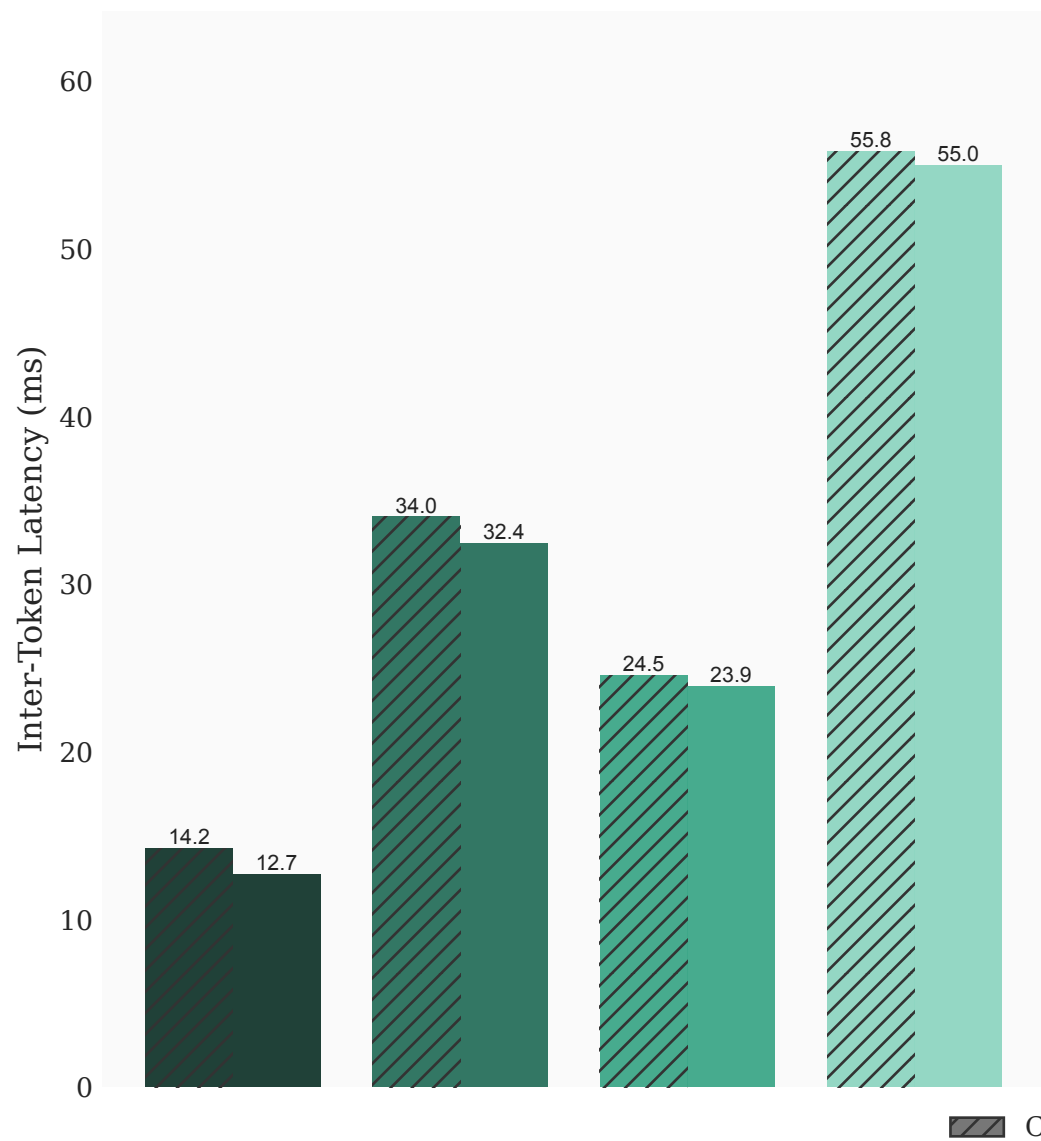# Edit 10K Characters (50 Concurrent Requests)

## Mean ITL



Inter-Token Latency (ms)

- 33.3
- 30.6
- 38.4
- 36.8
- 27.9
- 28.0
- 56.5
- 56.9

## P99 ITL



Inter-Token Latency (ms)

- 40.2
- 36.0
- 196
- 188
- 140
- 141
- 40.8
- 40.8

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)

## Mean ITL

Inter-Token Latency (ms)

- 14.2
- 12.7
- 34.0
- 32.4
- 24.5
- 23.9
- 55.8
- 55.0

## P99 ITL

Inter-Token Latency (ms)

- 18.6
- 17.0
- 210
- 207
- 139
- 139
- 438
- 201

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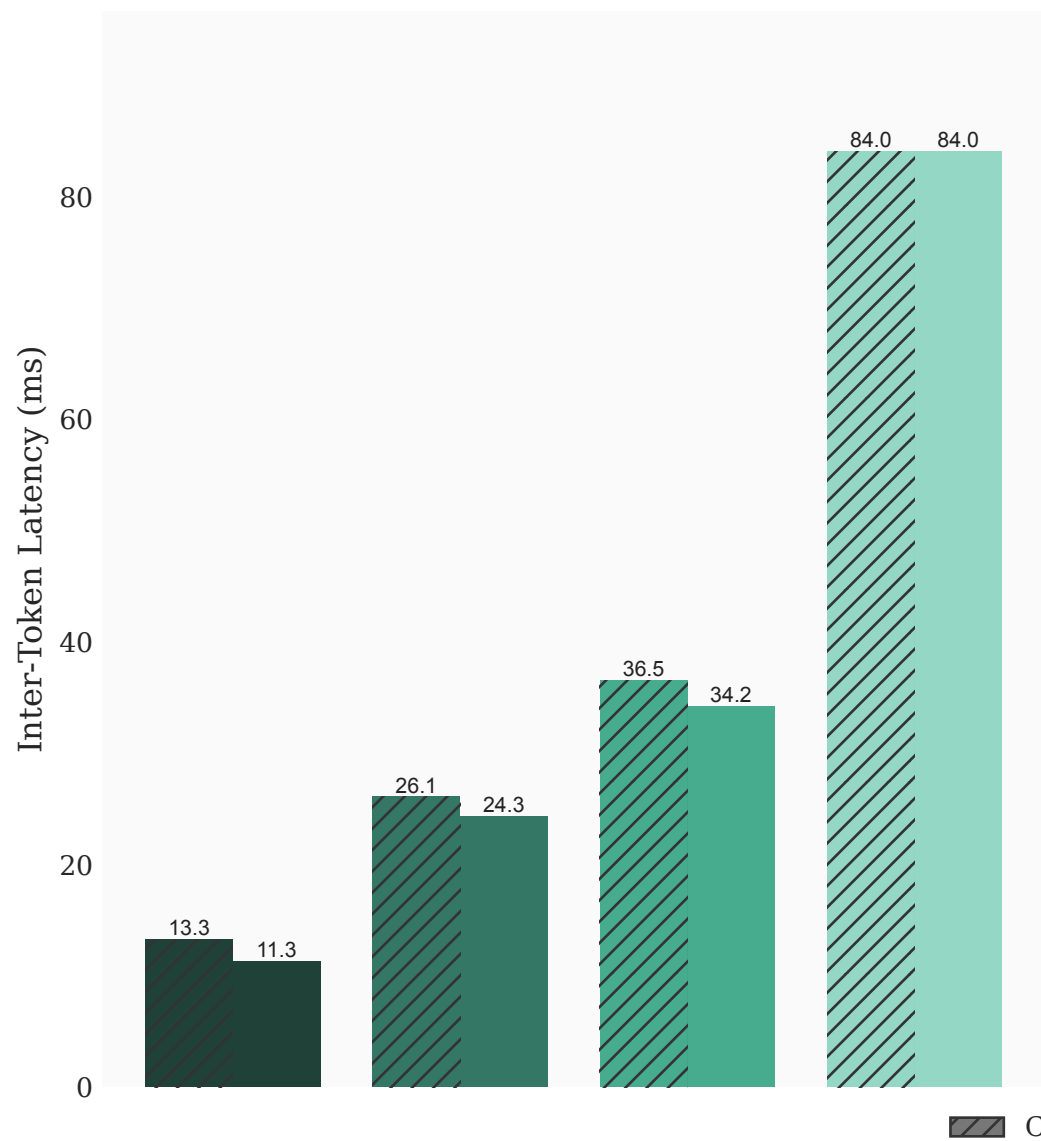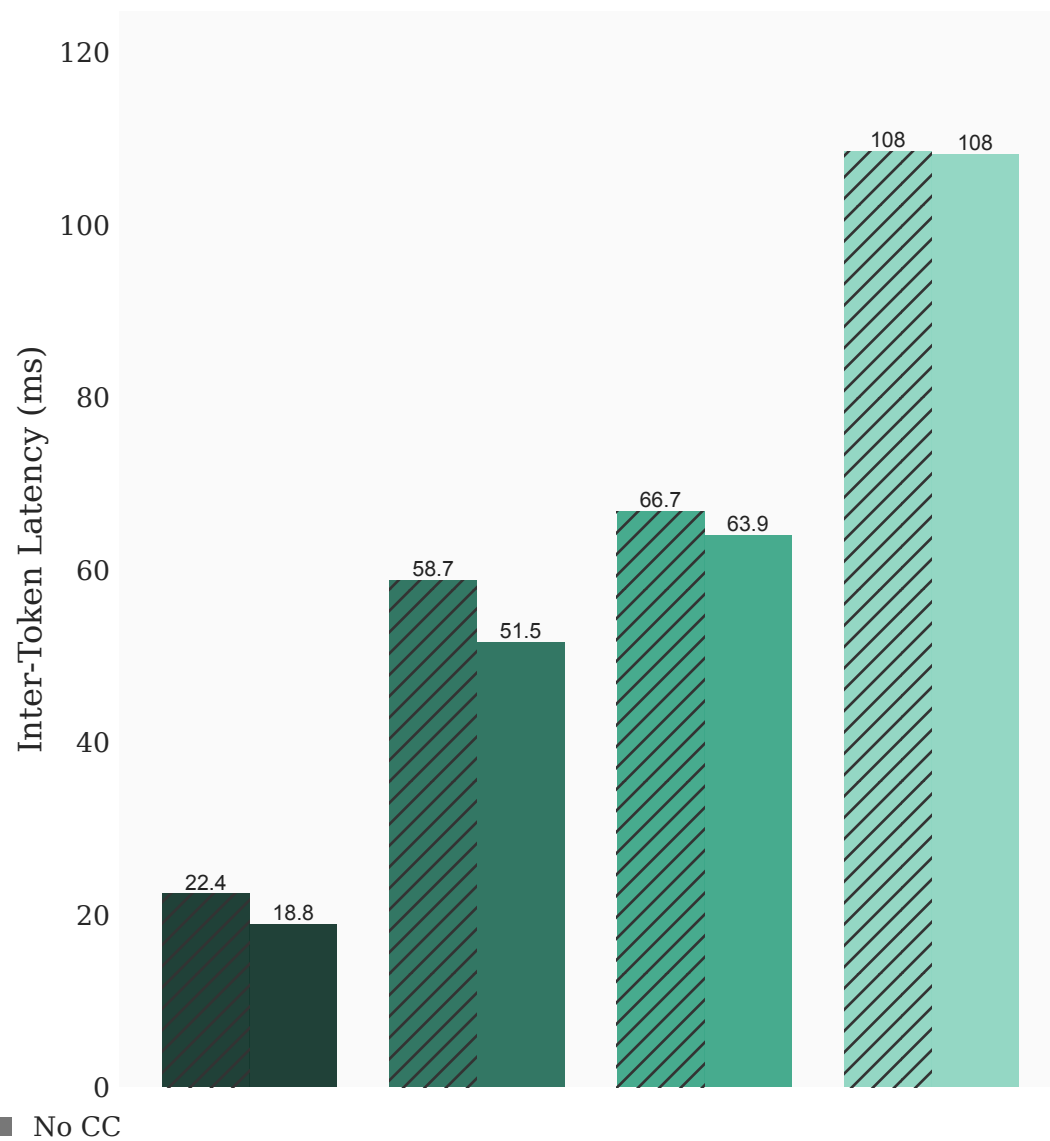
▬ LLama 3.1 8B   ▬ Mistral 3.1 24B   ▬ GPT OSS 120B   ▬ LLama 3.3 70B Int4

# Numina Math (100 Concurrent Requests)

## Mean ITL

**Inter-Token Latency (ms)**

- 13.3
- 11.3
- 26.1
- 24.3
- 36.5
- 34.2
- 84.0
- 84.0

## P99 ITL

**Inter-Token Latency (ms)**

- 22.4
- 18.8
- 58.7
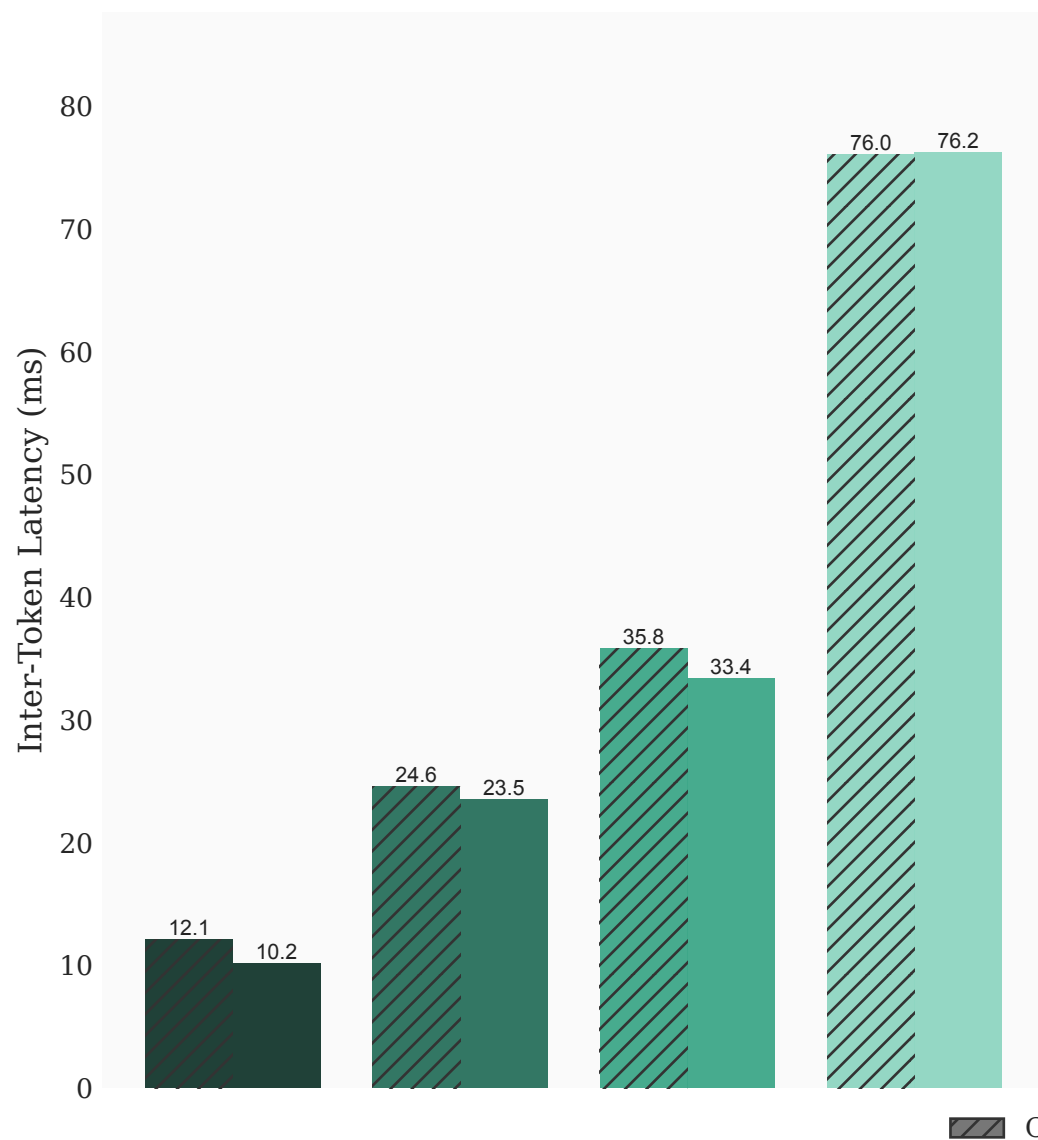- 51.5
- 66.7
- 63.9
- 108
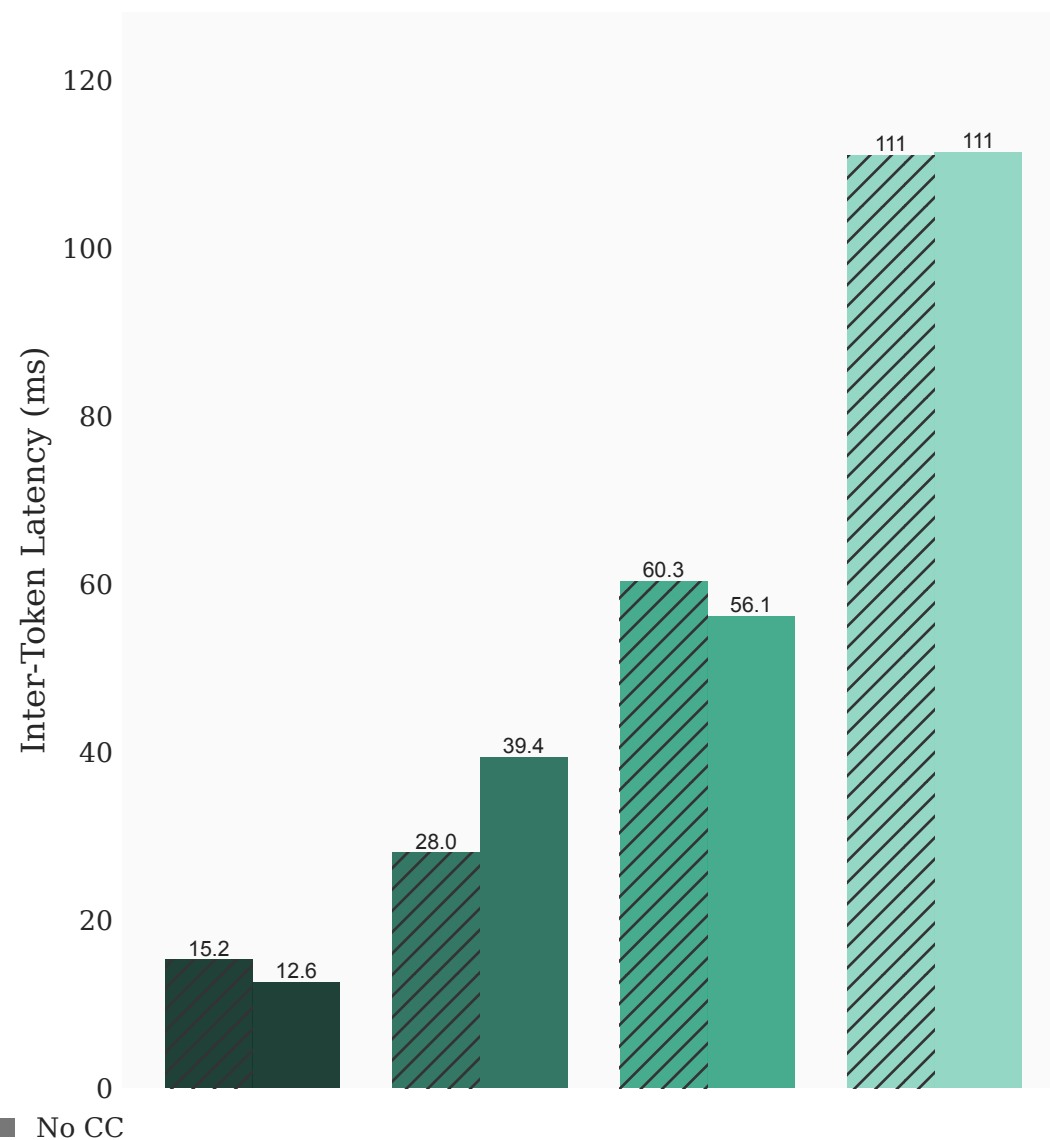- 108

CC   No CC

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

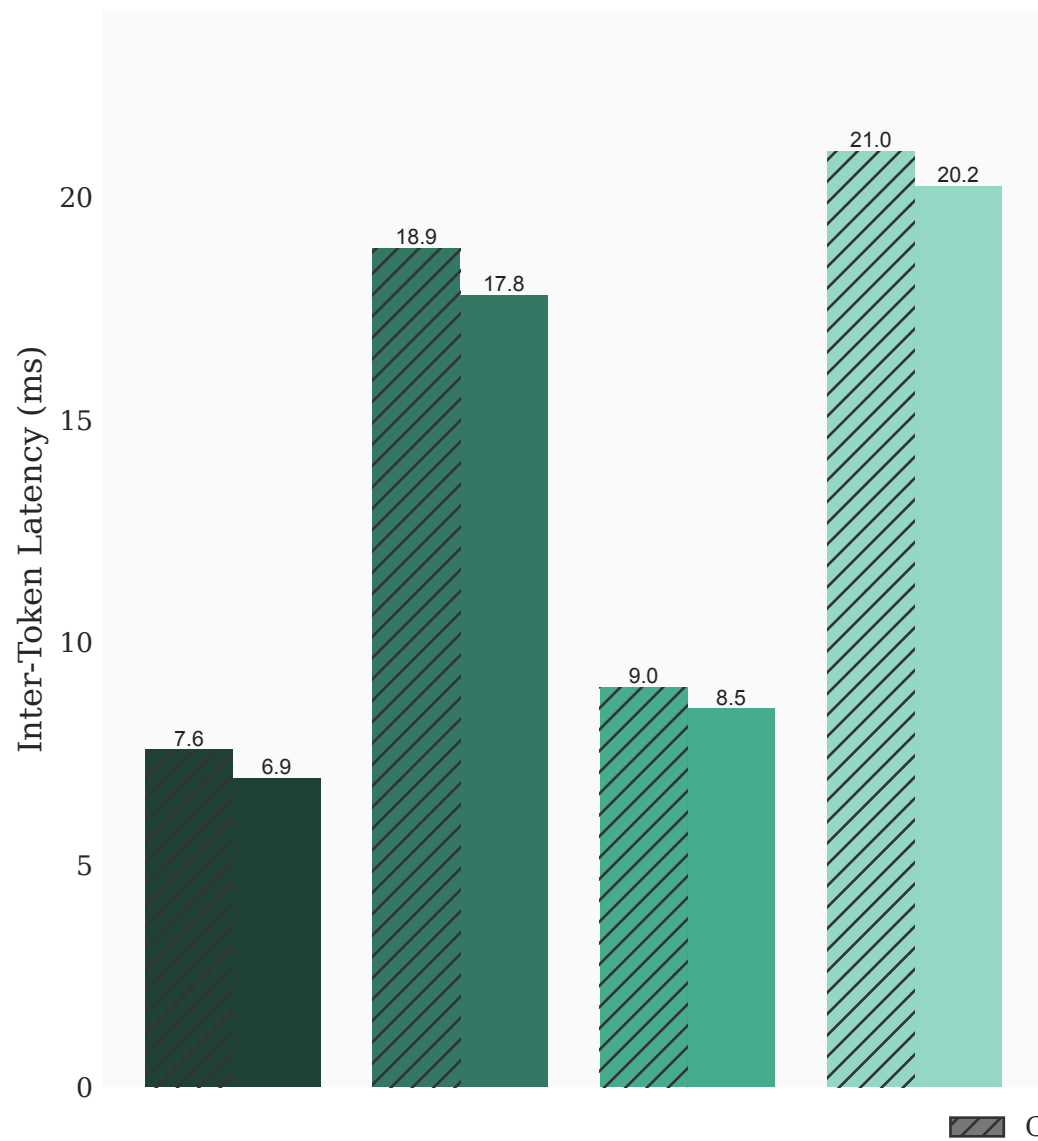# Numina Math (50 Concurrent Requests)

## Mean ITL



## P99 ITL

Legend: CC, No CC

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

# Numina Math (1 Concurrent Requests)

## Mean ITL



## P99 ITL



Legend: CC, No CC

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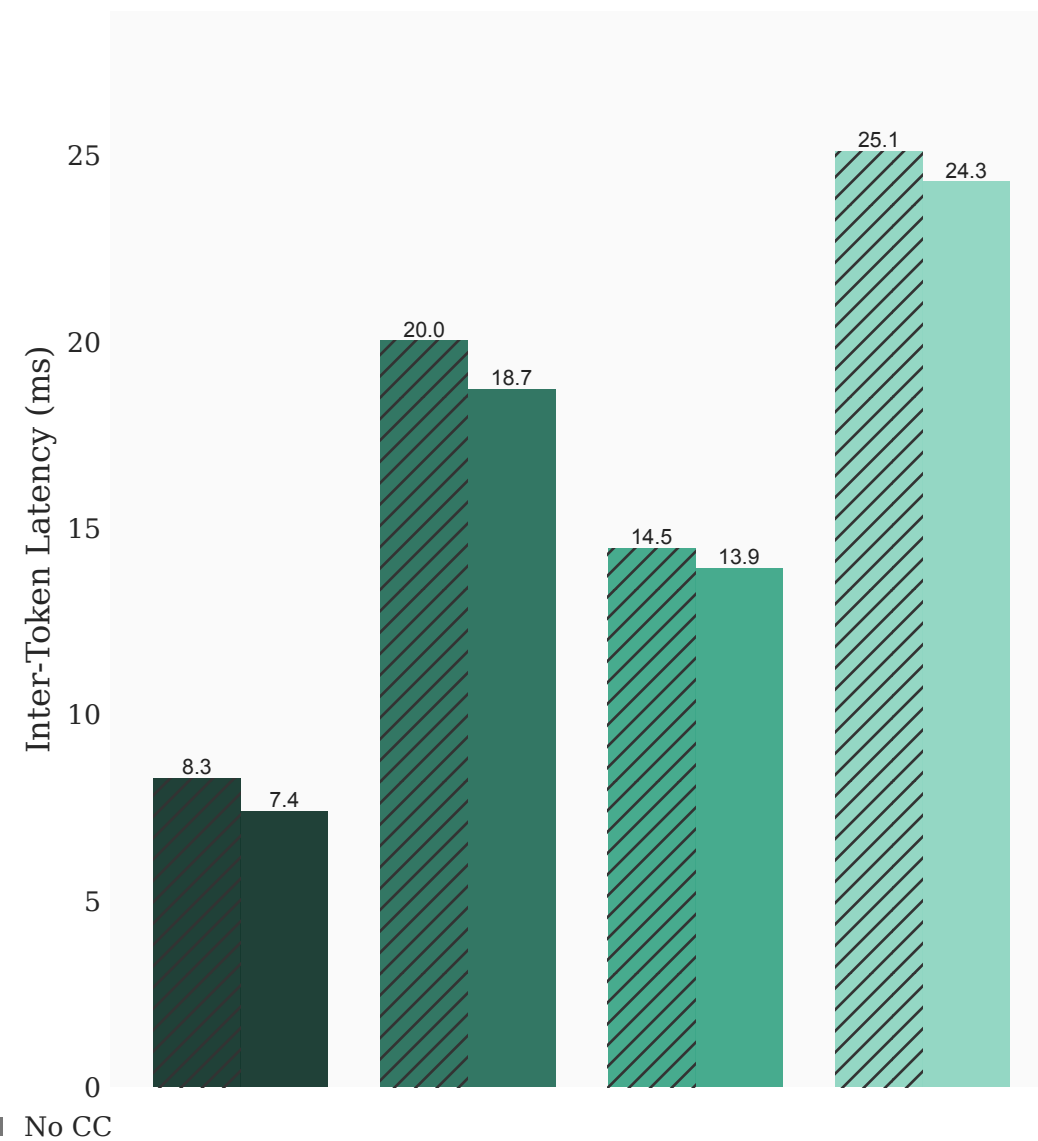