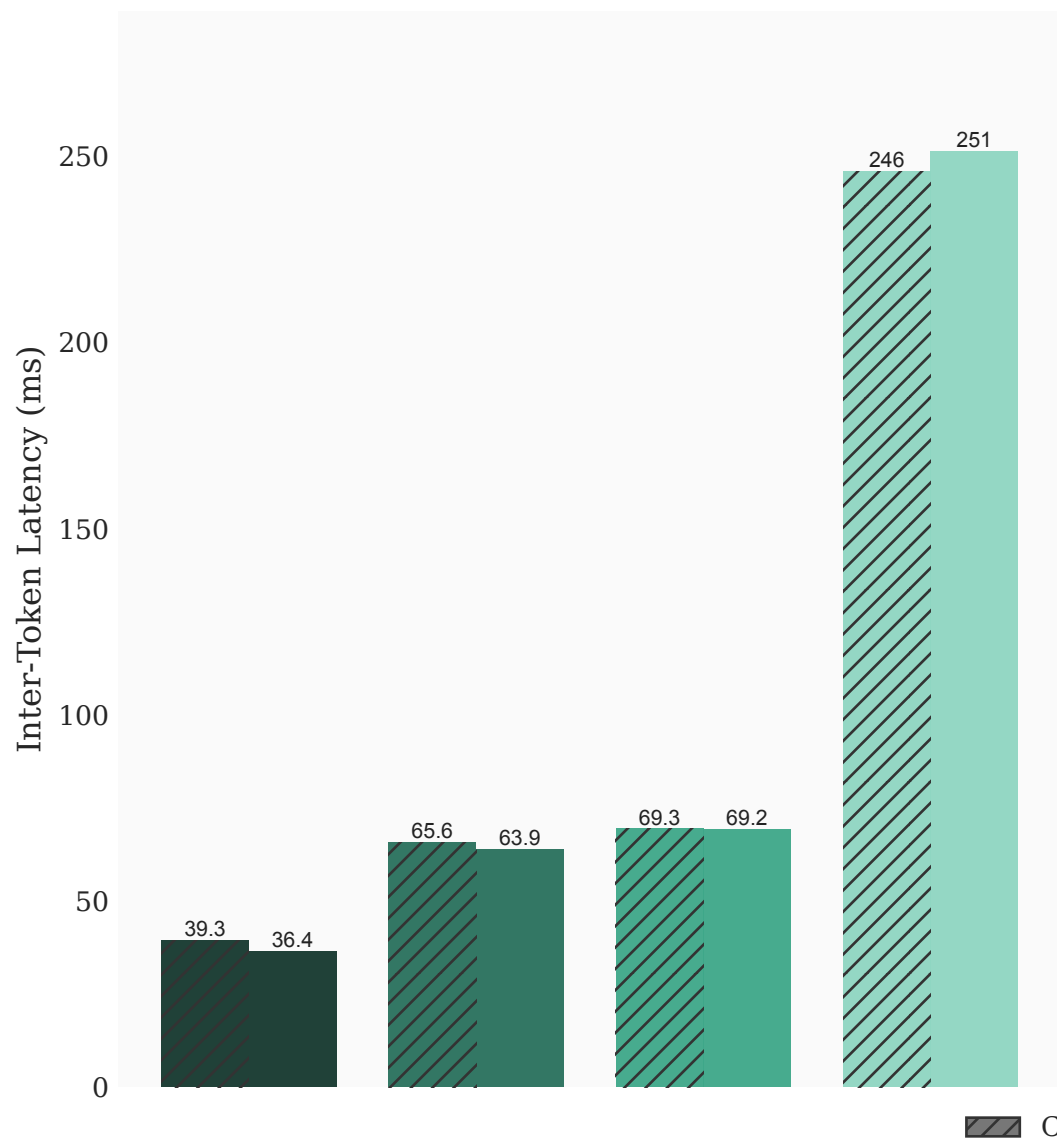
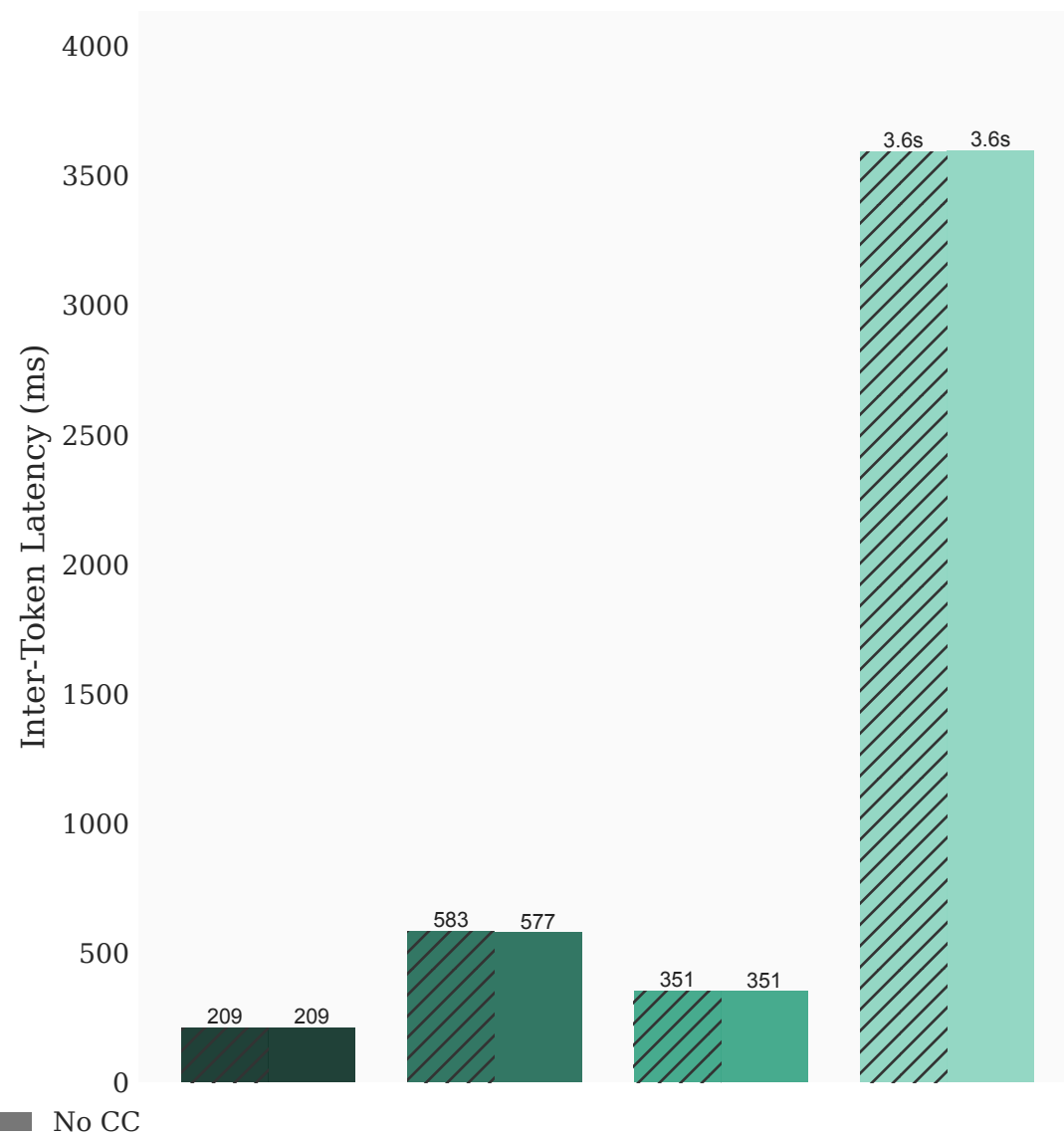


Random (1500 \Rightarrow 250) (100 Request Rate)

Mean ITL



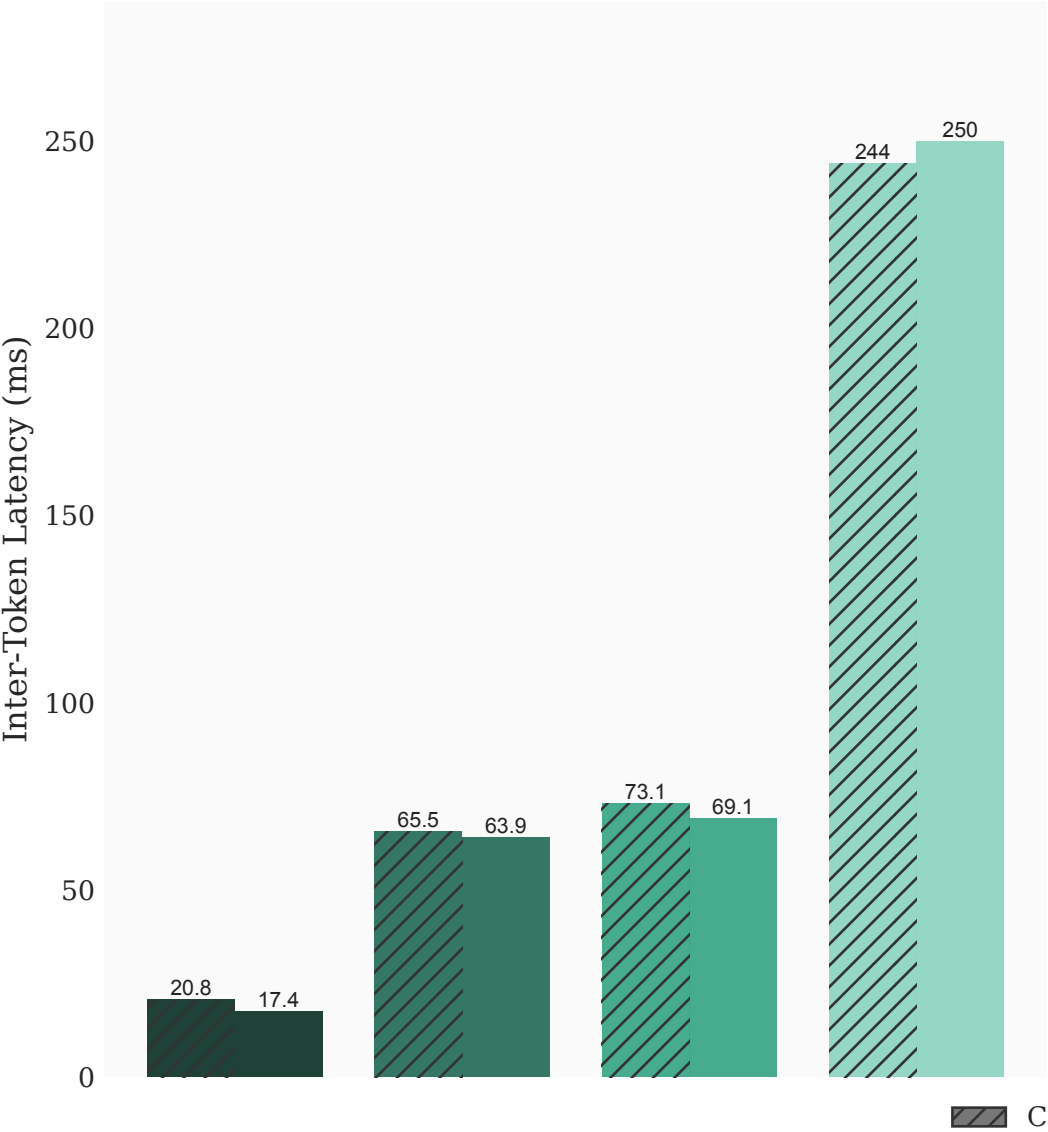
P99 ITL



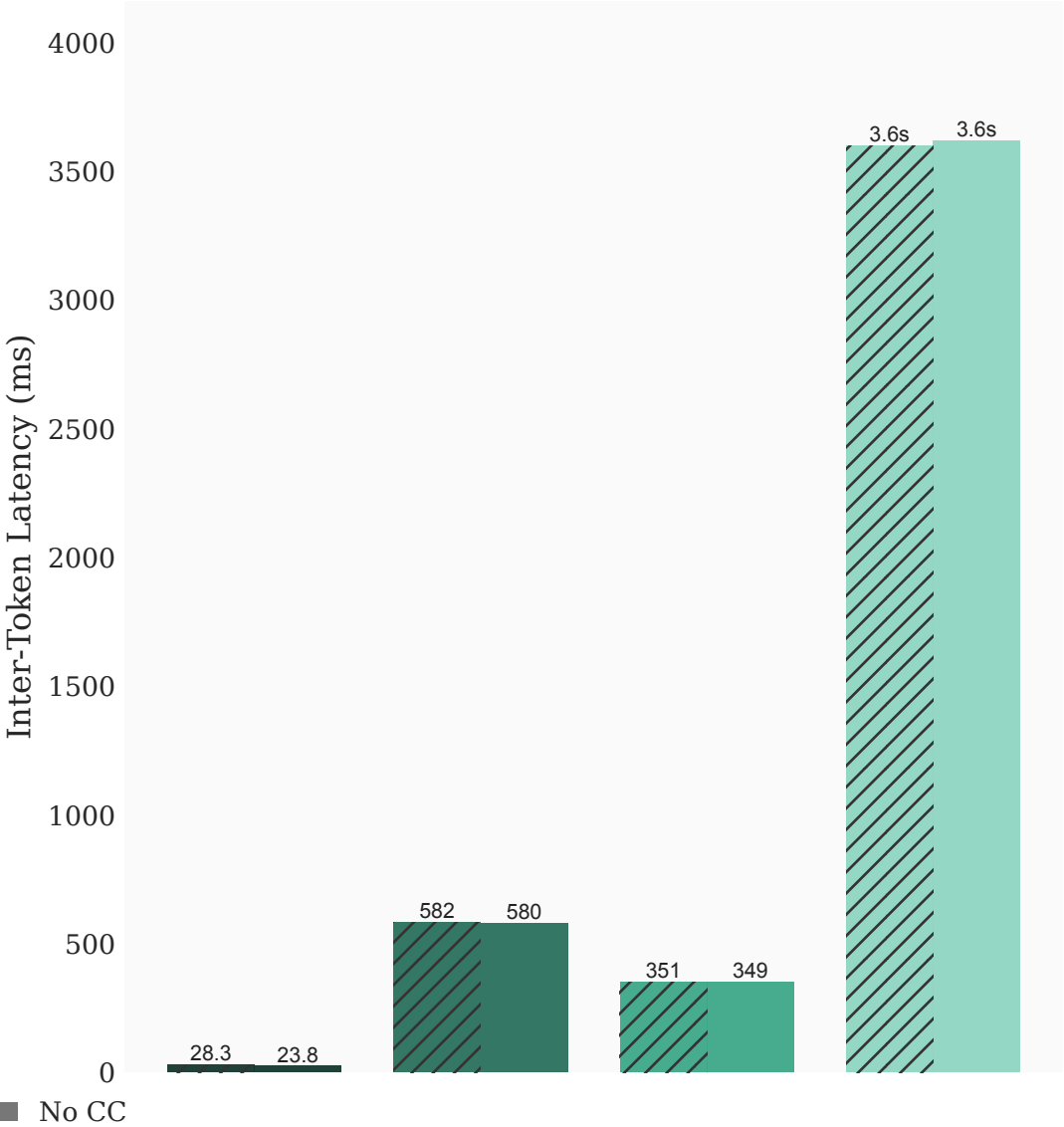
Legend: CC (hatched), No CC (solid). Models: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4.

Random (1500 \Rightarrow 250) (50 Request Rate)

Mean ITL



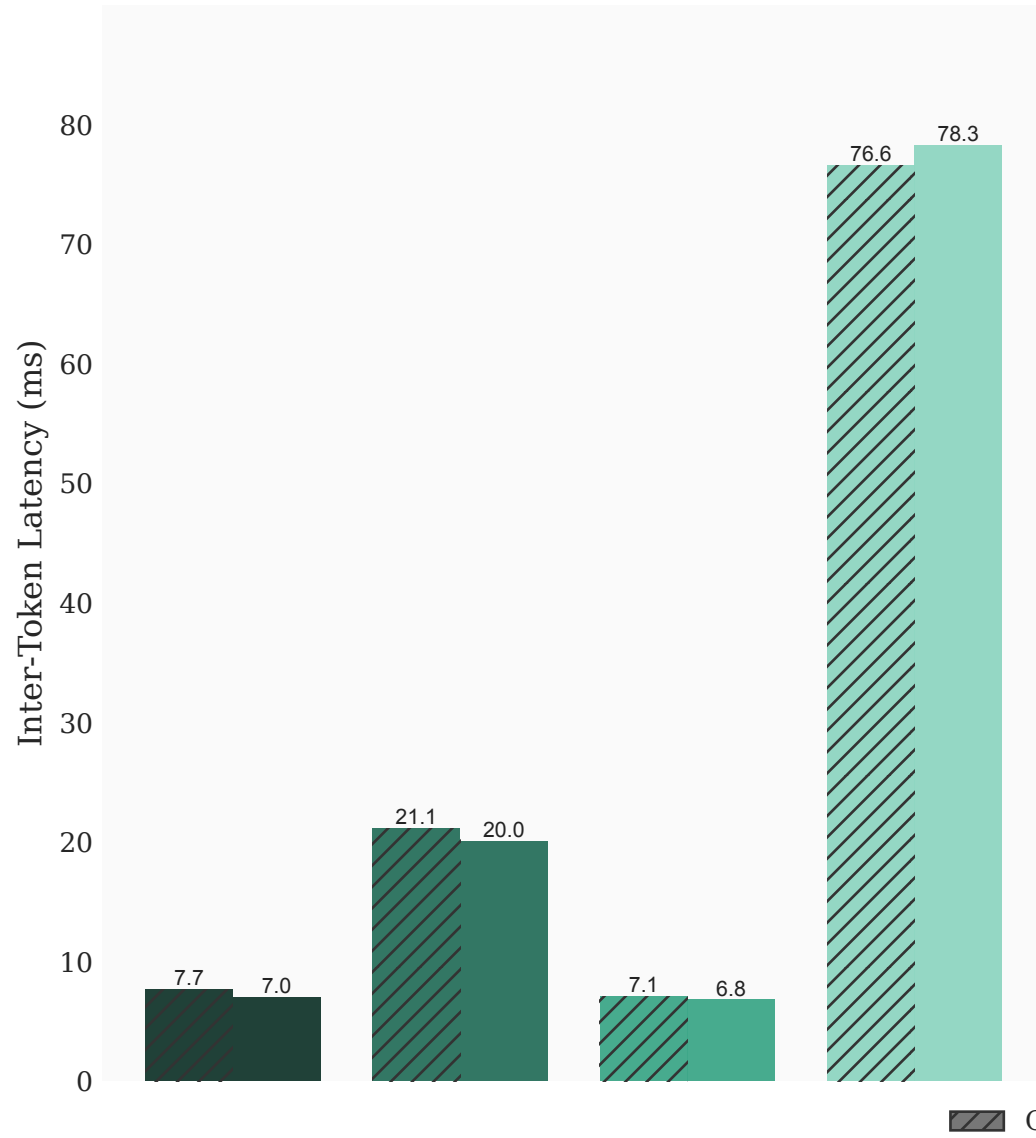
P99 ITL



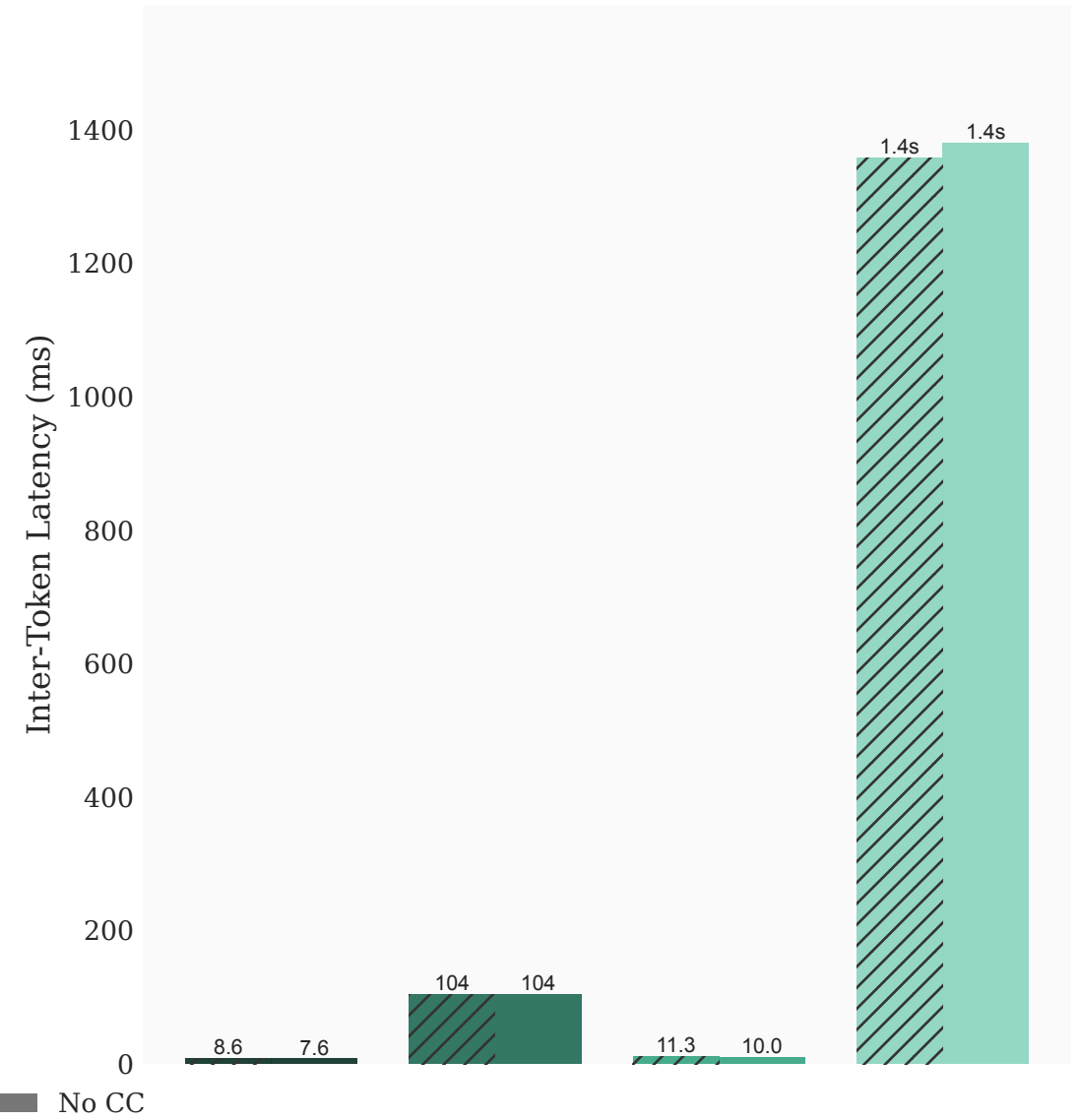
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Random (1500 \Rightarrow 250) (Single Request)

Mean ITL



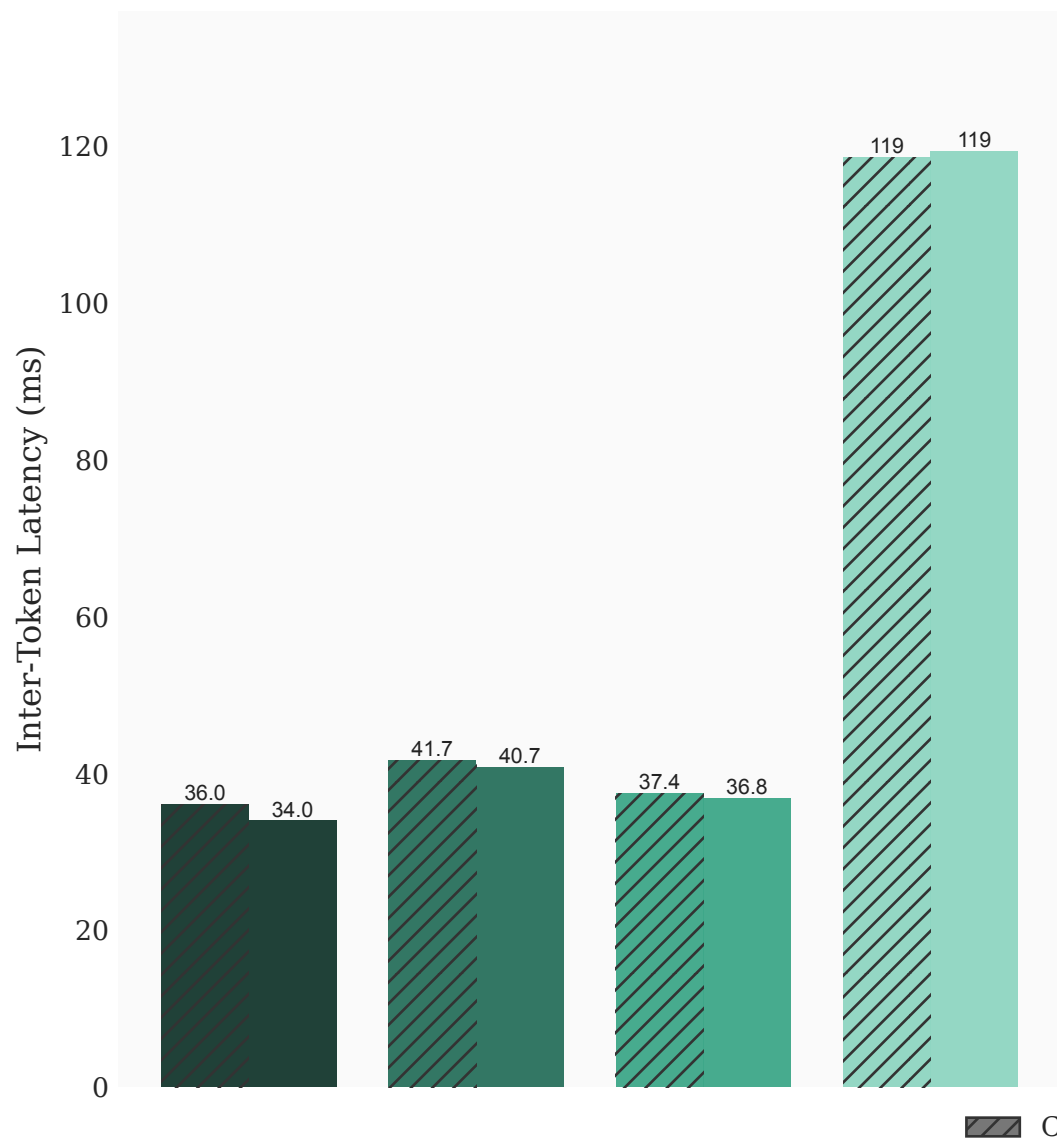
P99 ITL



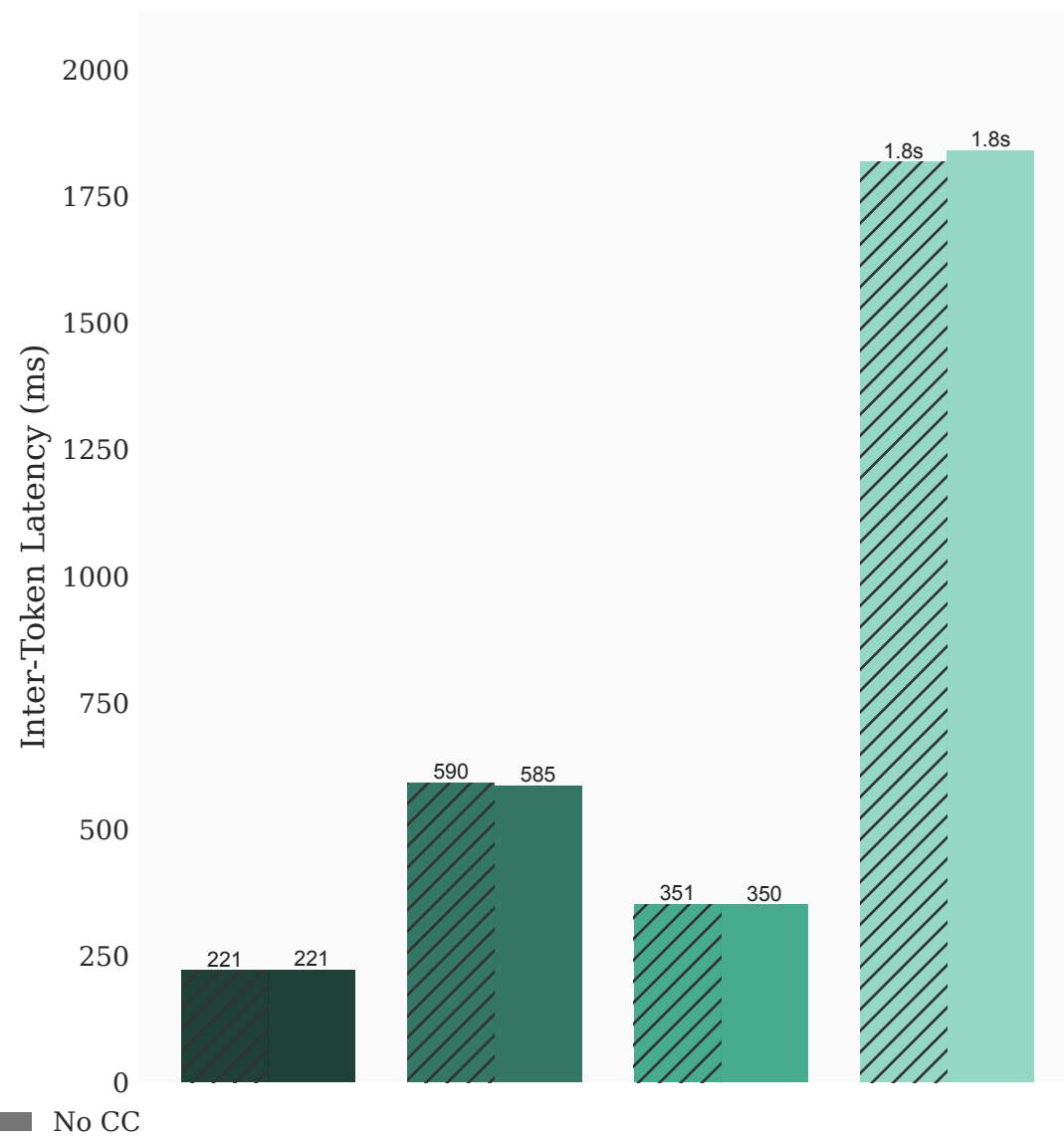
Legend: CC (hatched), No CC (solid). Models: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4.

Random (4000 \Rightarrow 1000) (100 Request Rate)

Mean ITL



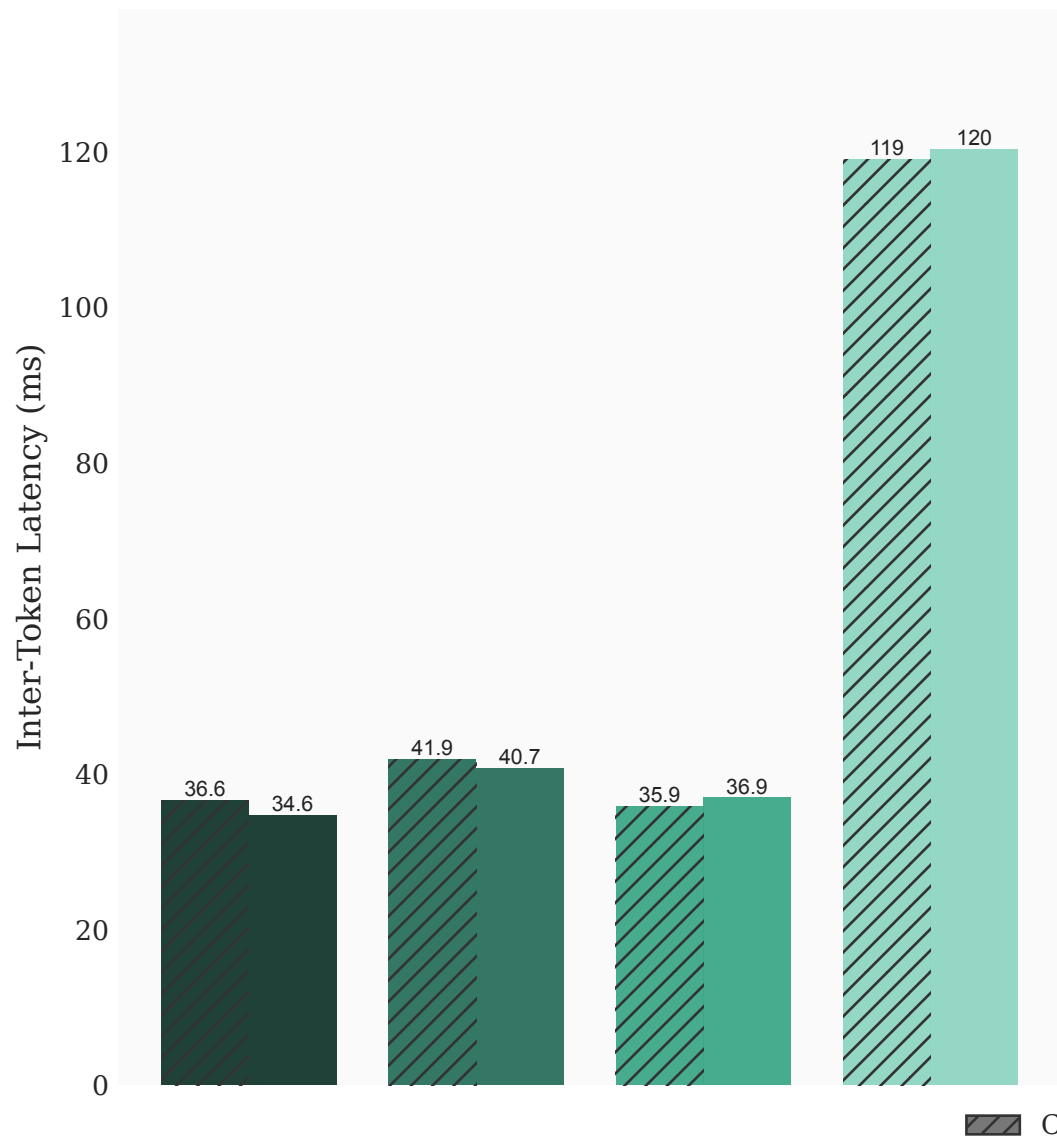
P99 ITL



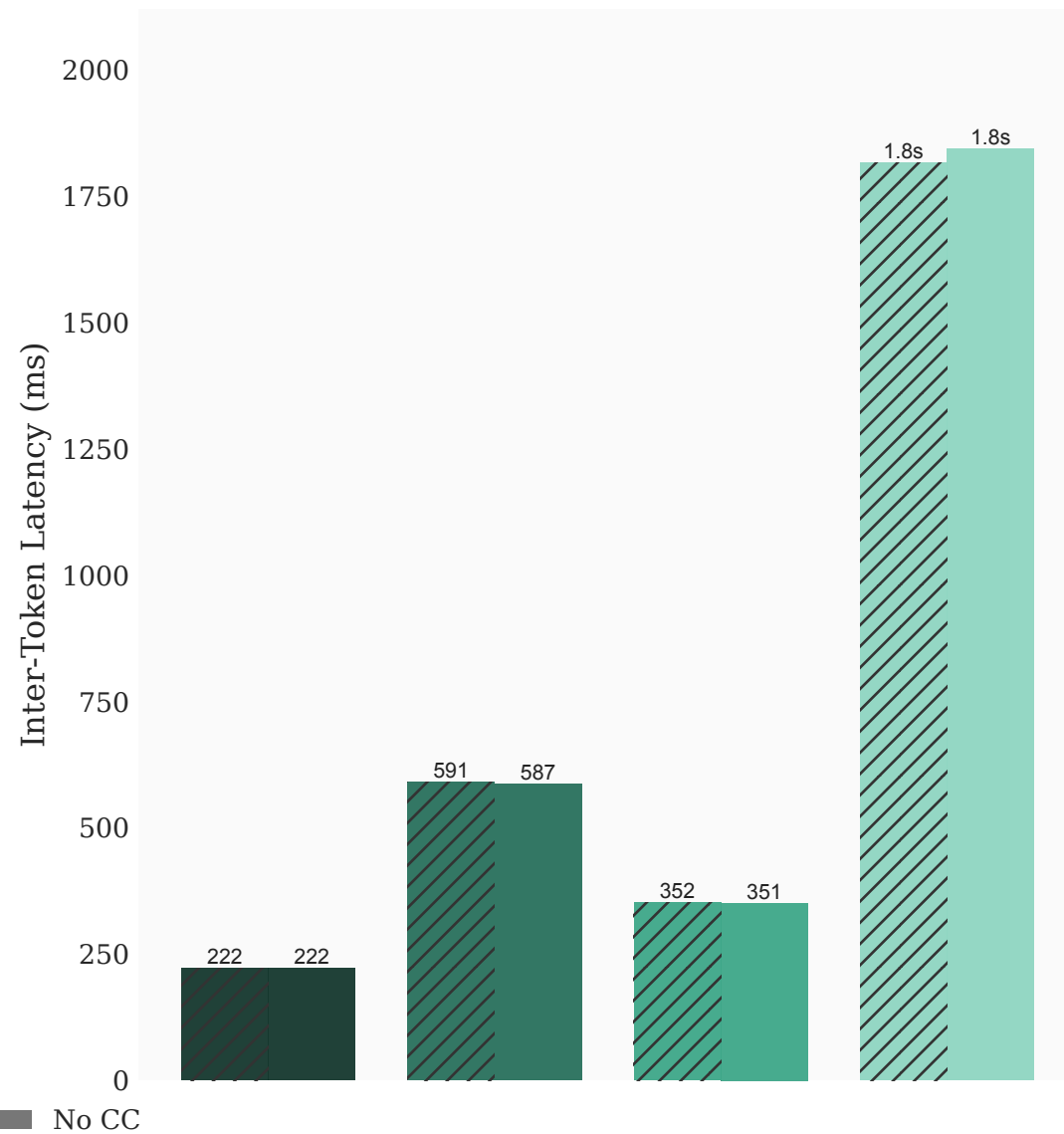
LLama 3.1 8B Mistral 3.1 24B GPT OSS 120B LLama 3.3 70B Int4

Random (4000 \Rightarrow 1000) (50 Request Rate)

Mean ITL



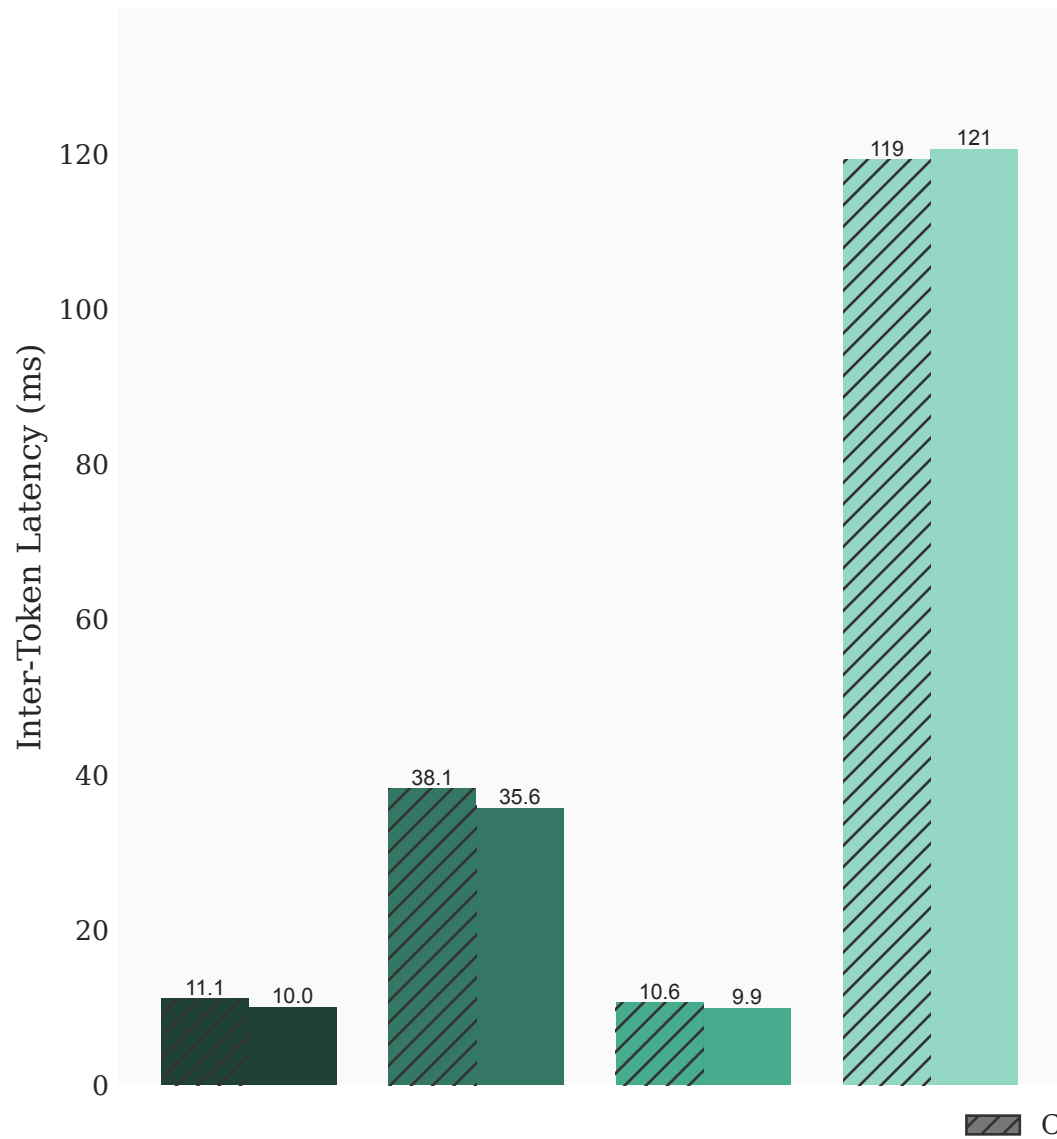
P99 ITL



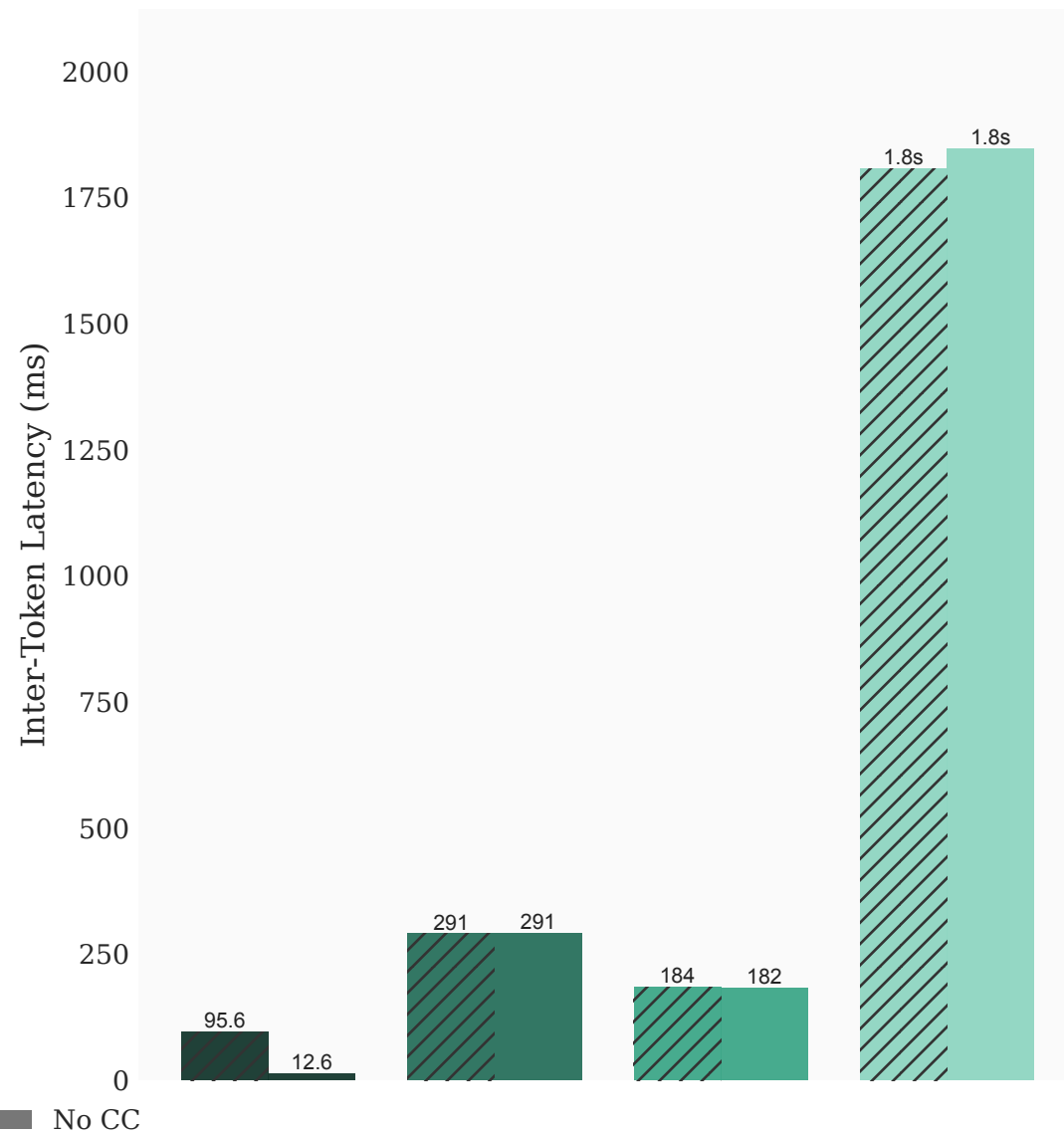
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Random (4000 \Rightarrow 1000) (Single Request)

Mean ITL



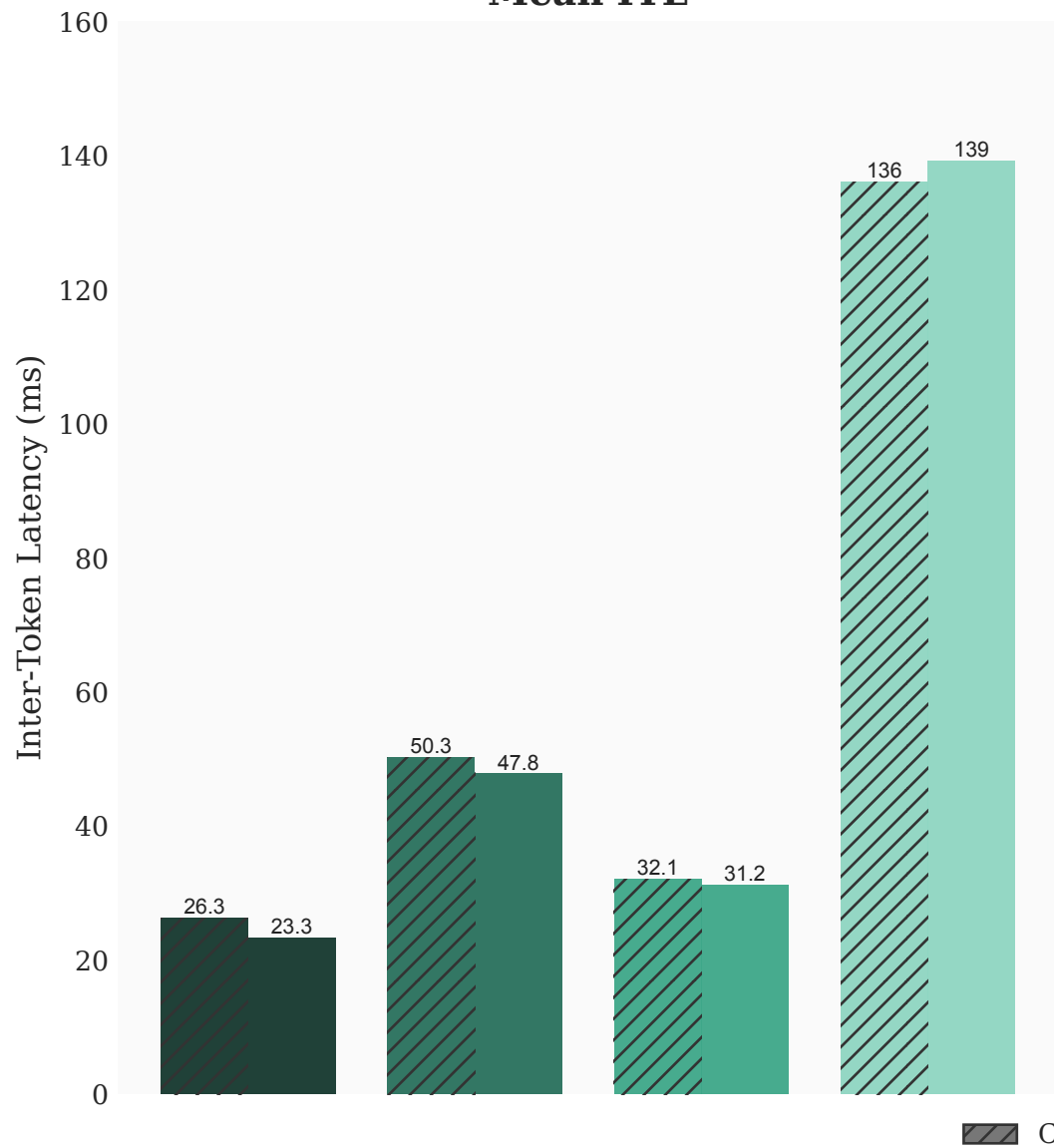
P99 ITL



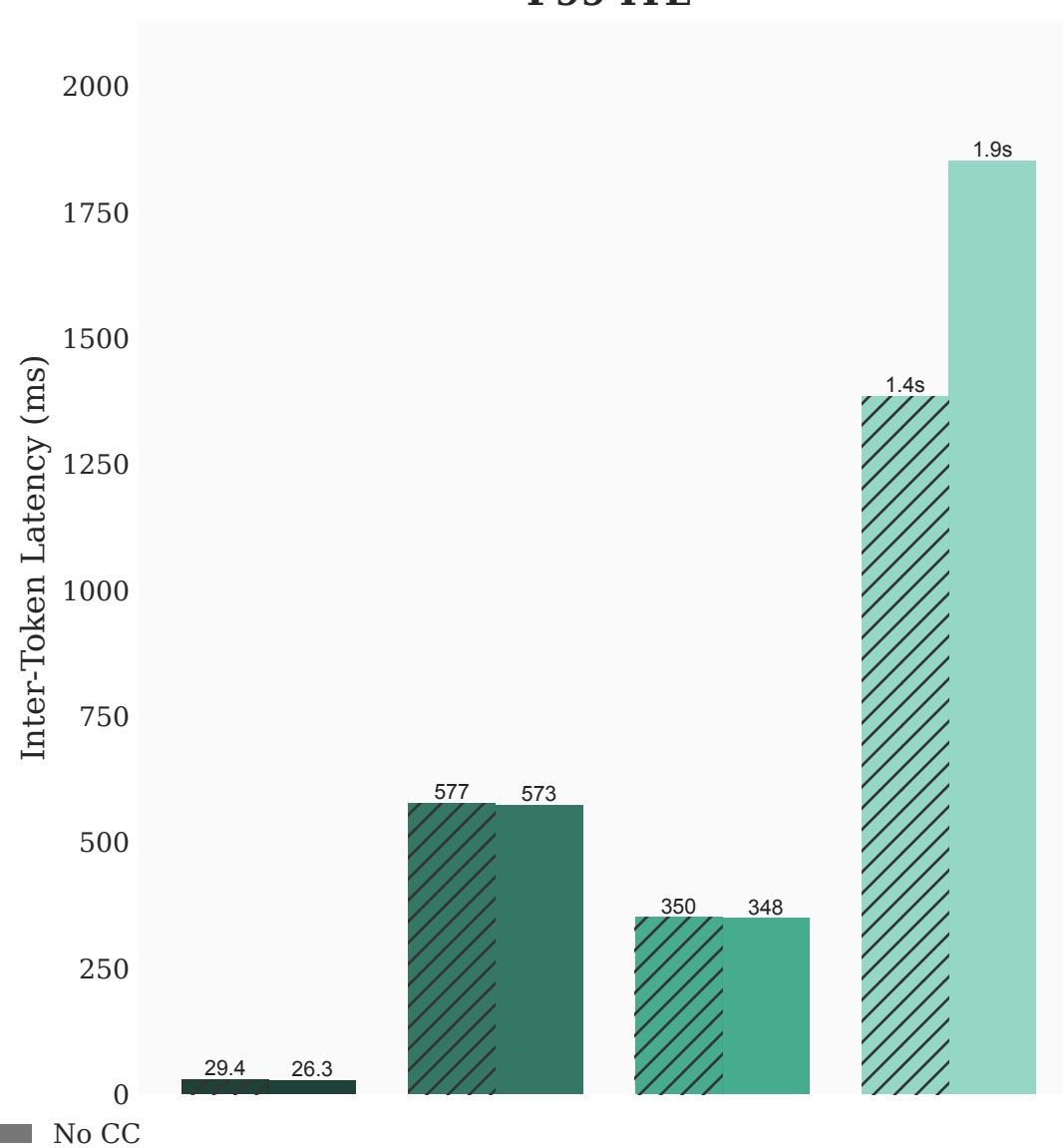
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Random (1000 \Rightarrow 1000) (100 Request Rate)

Mean ITL



P99 ITL

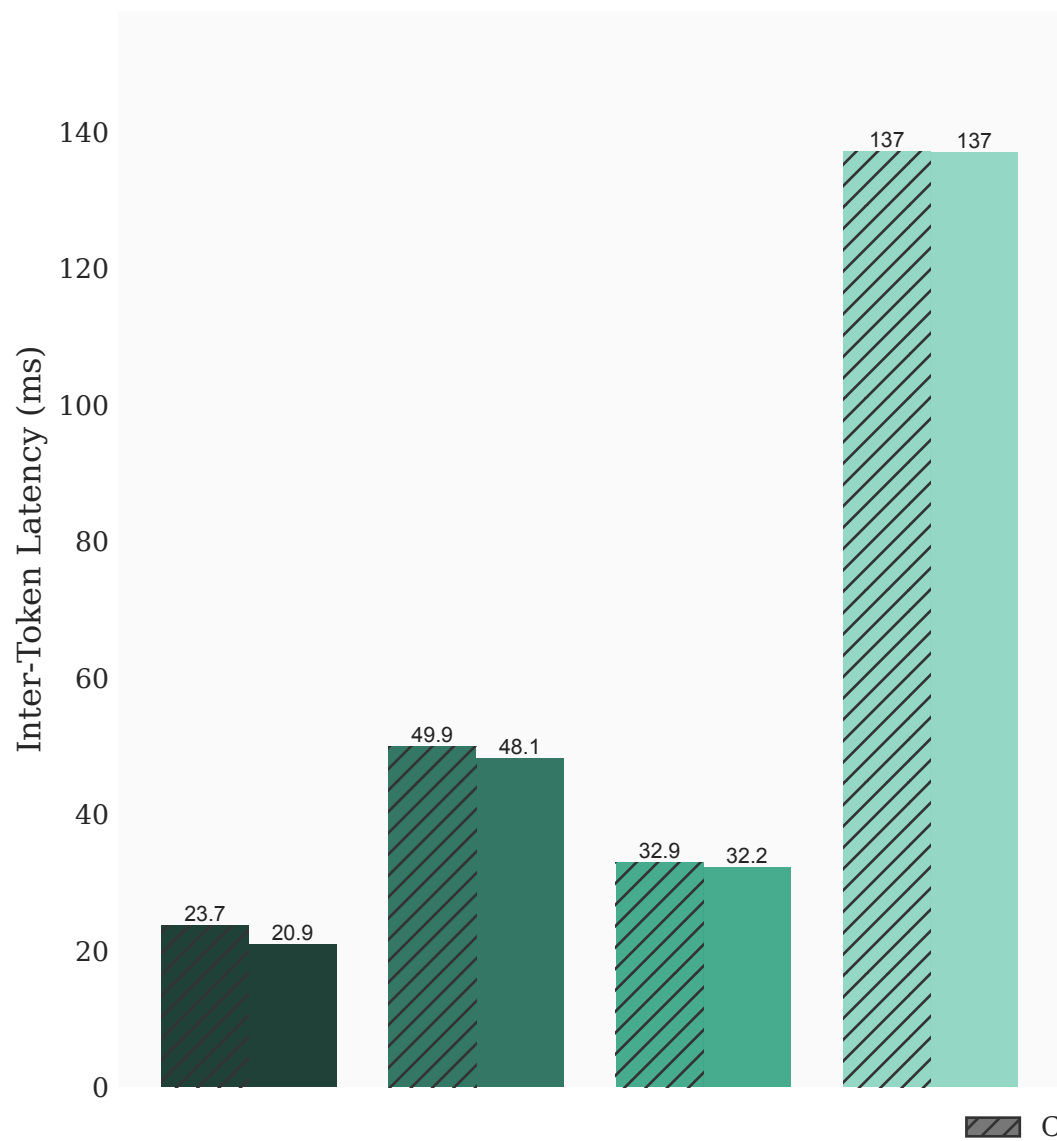


■ CC ■ No CC

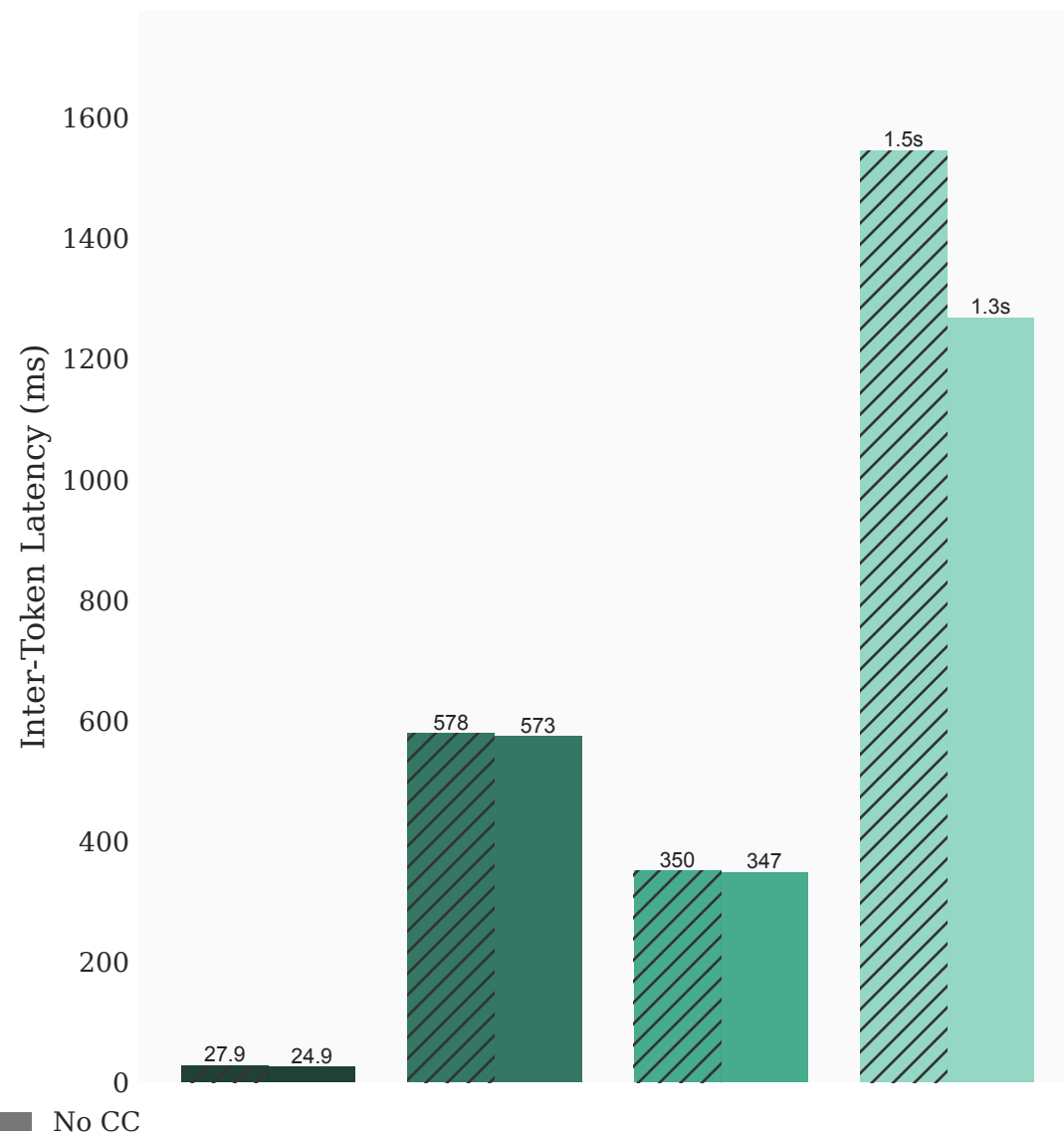
■ LLama 3.1 8B ■ Mistral 3.1 24B ■ GPT OSS 120B ■ LLama 3.3 70B Int4

Random (1000 \Rightarrow 1000) (50 Request Rate)

Mean ITL



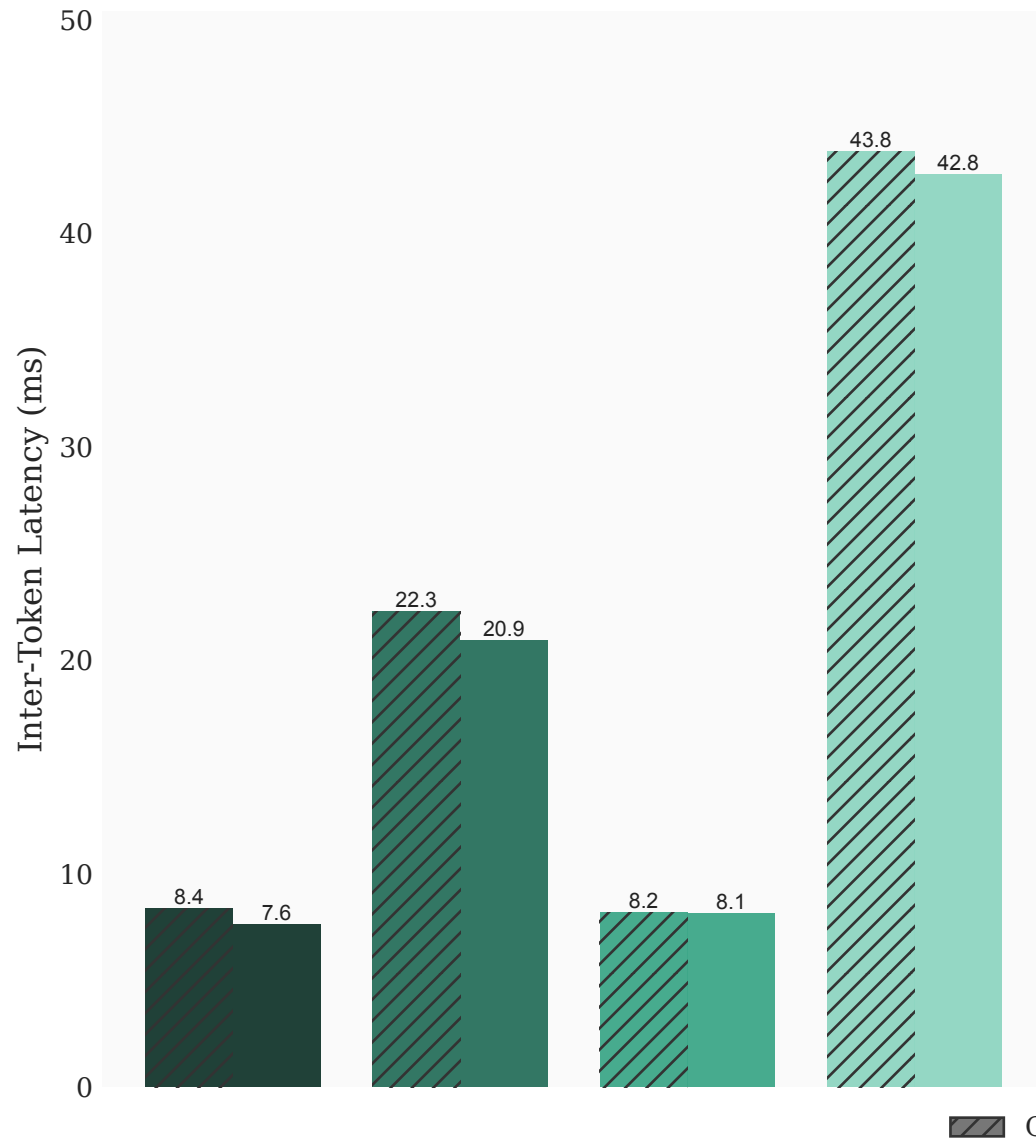
P99 ITL



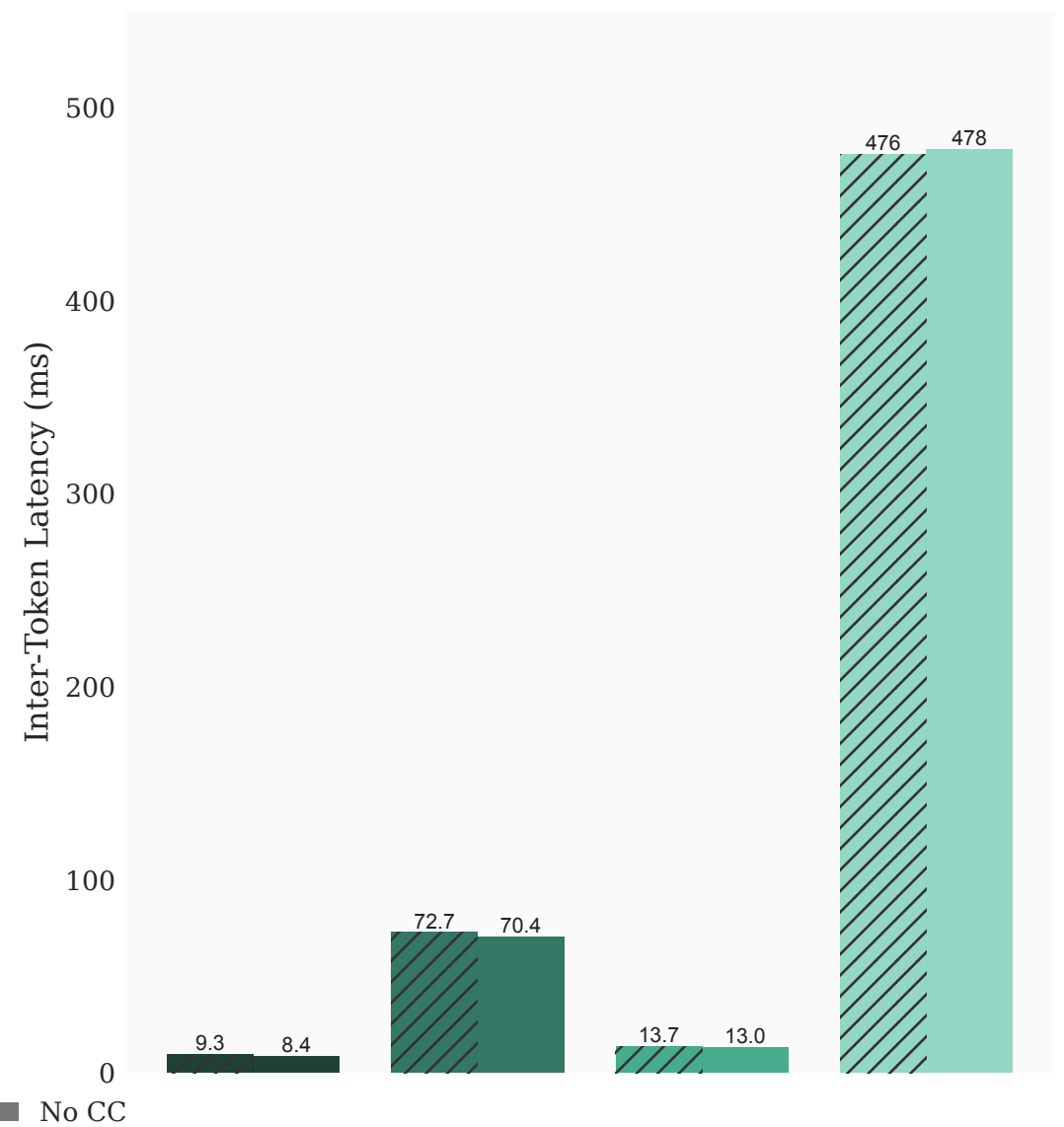
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Random (1000 \Rightarrow 1000) (Single Request)

Mean ITL



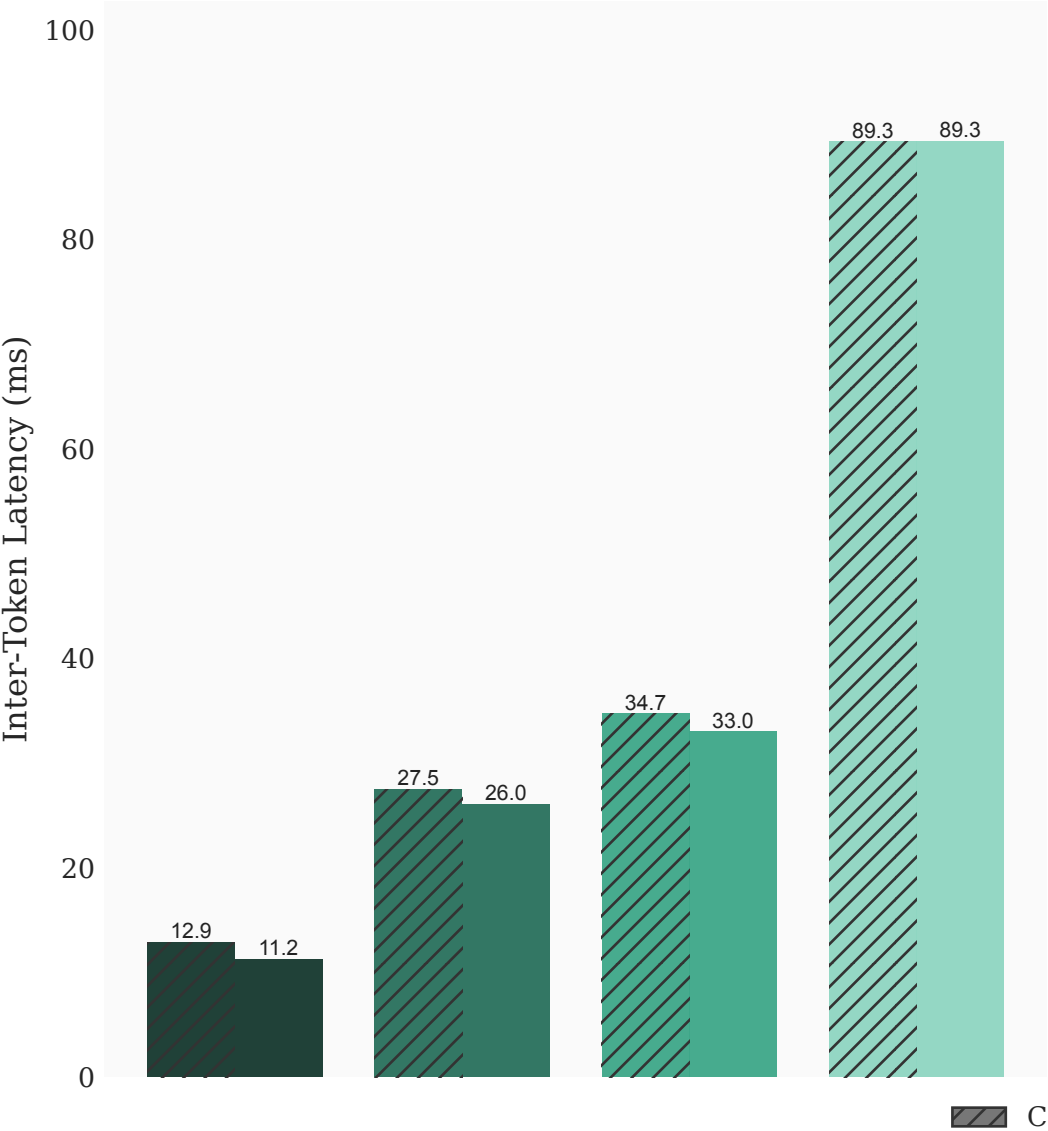
P99 ITL



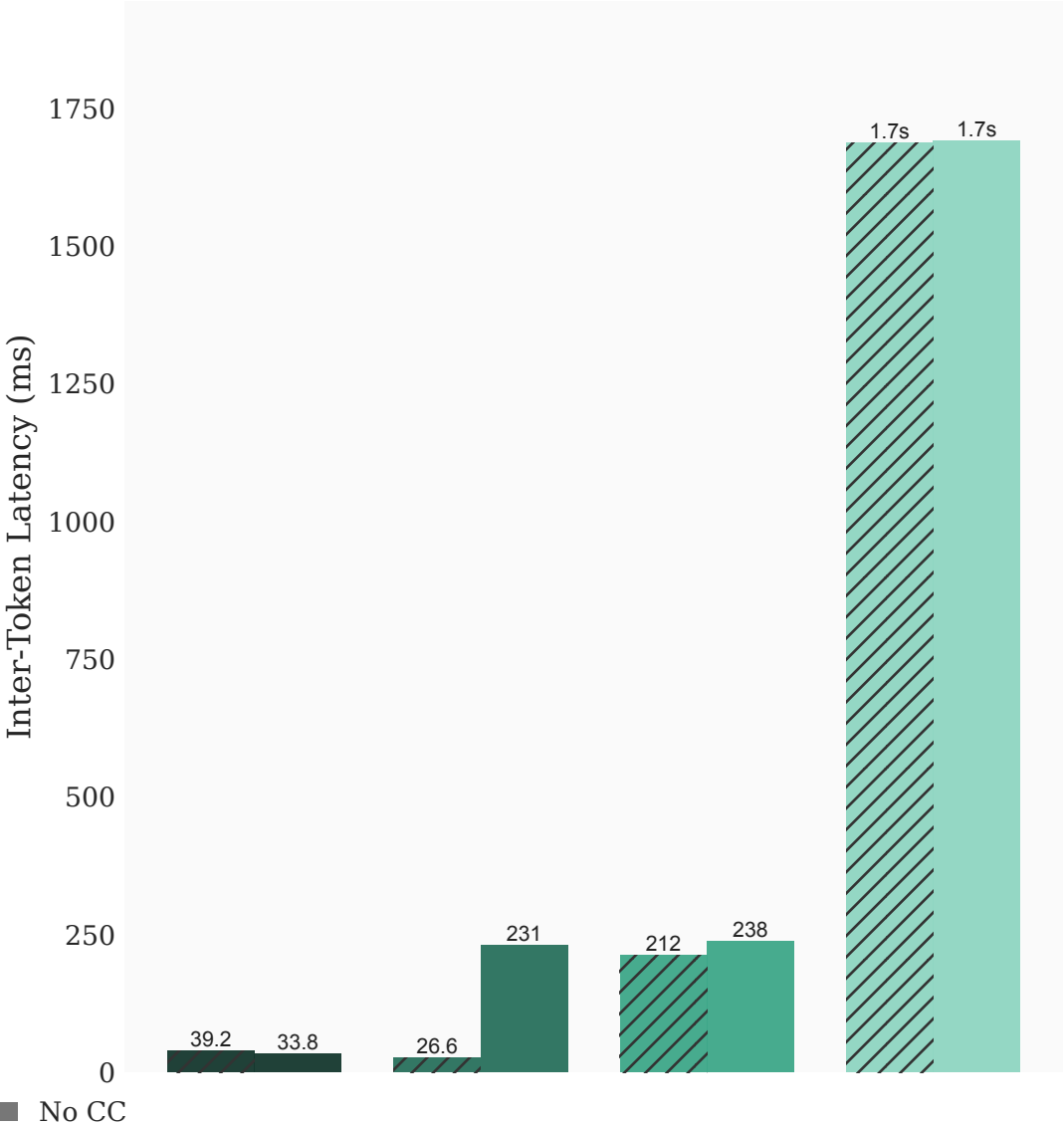
Legend: CC (hatched), No CC (solid), Llama 3.1 8B (dark teal), Mistral 3.1 24B (medium teal), GPT OSS 120B (light teal), Llama 3.3 70B Int4 (very light teal)

ShareGPT (100 Request Rate)

Mean ITL



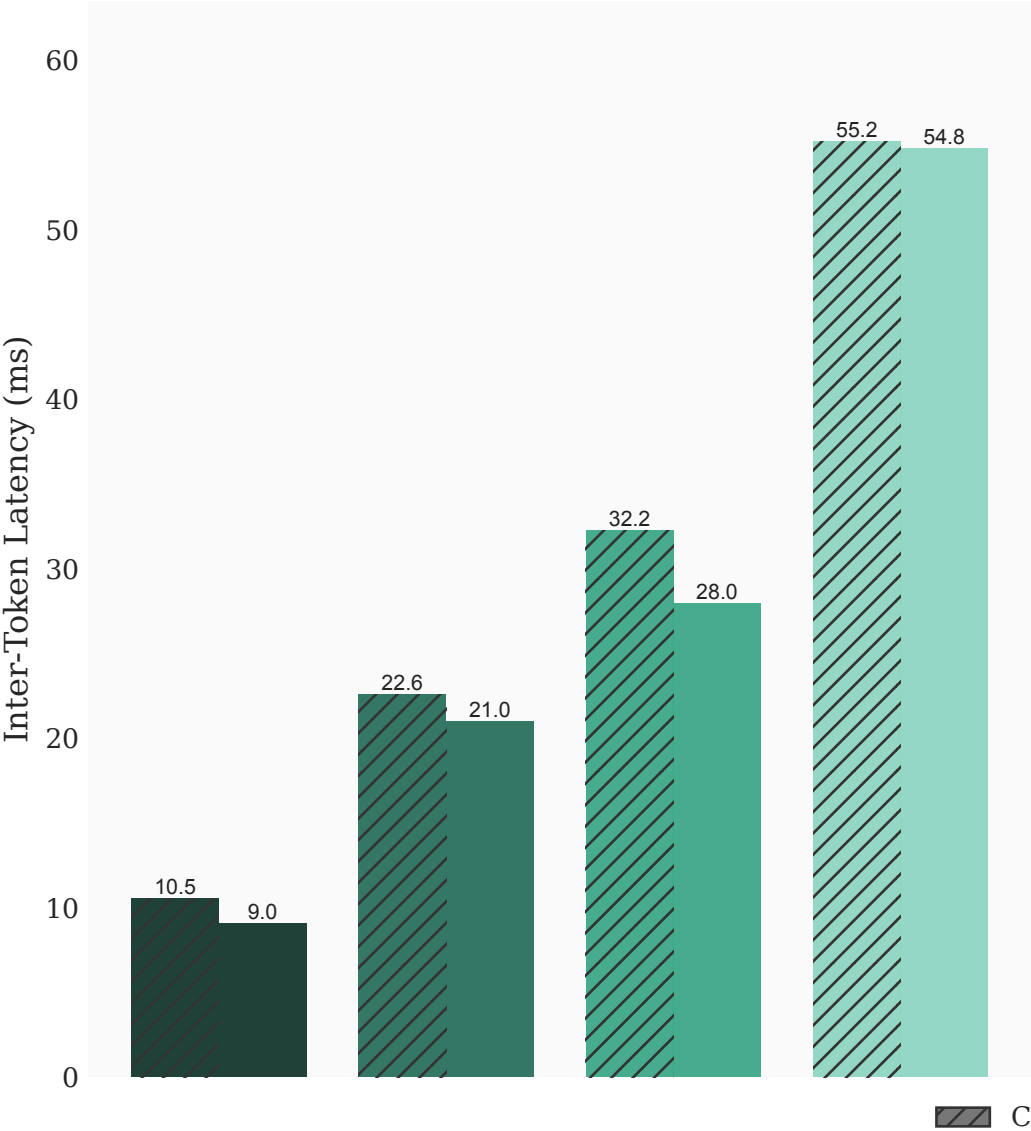
P99 ITL



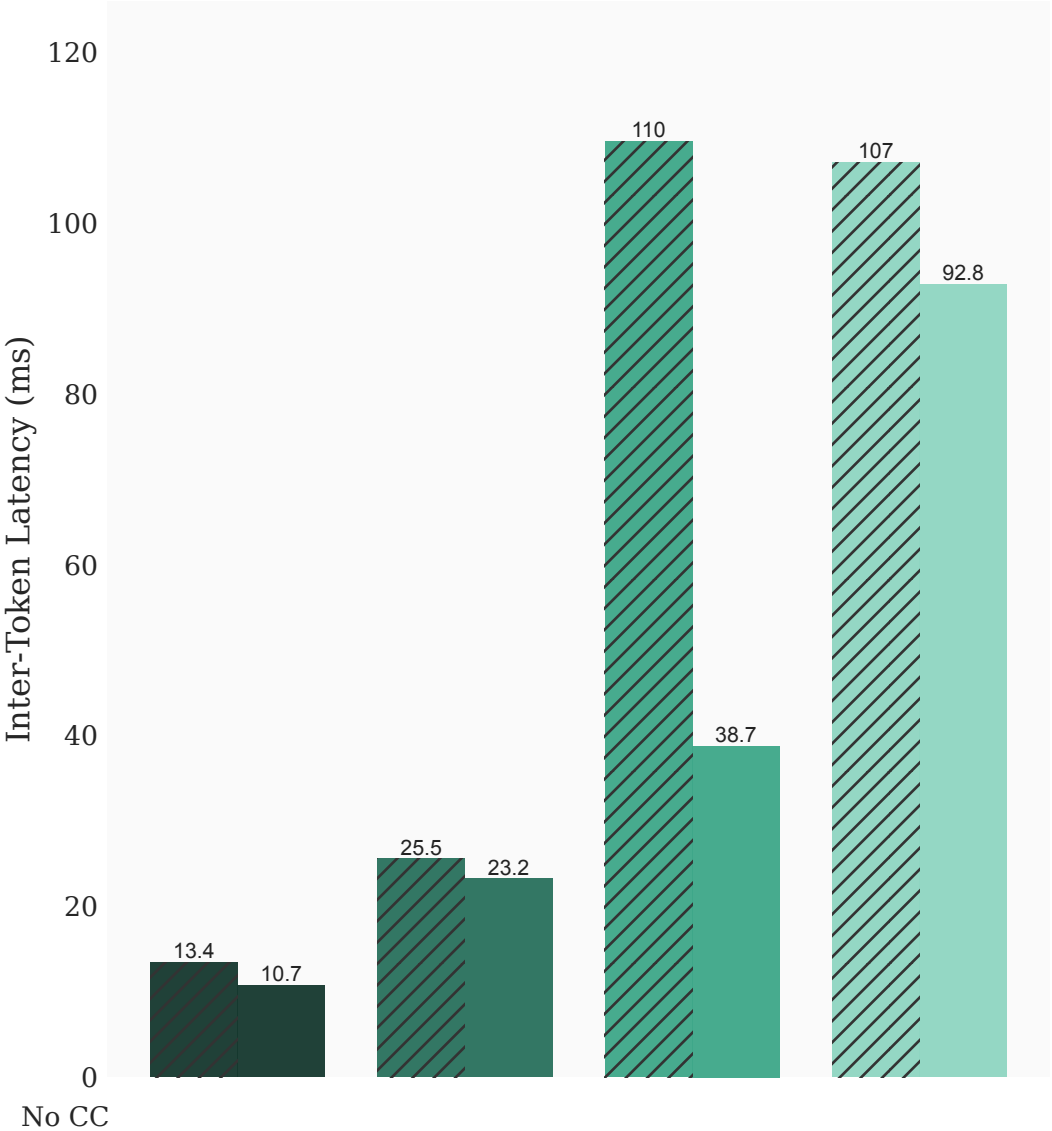
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

ShareGPT (50 Request Rate)

Mean ITL



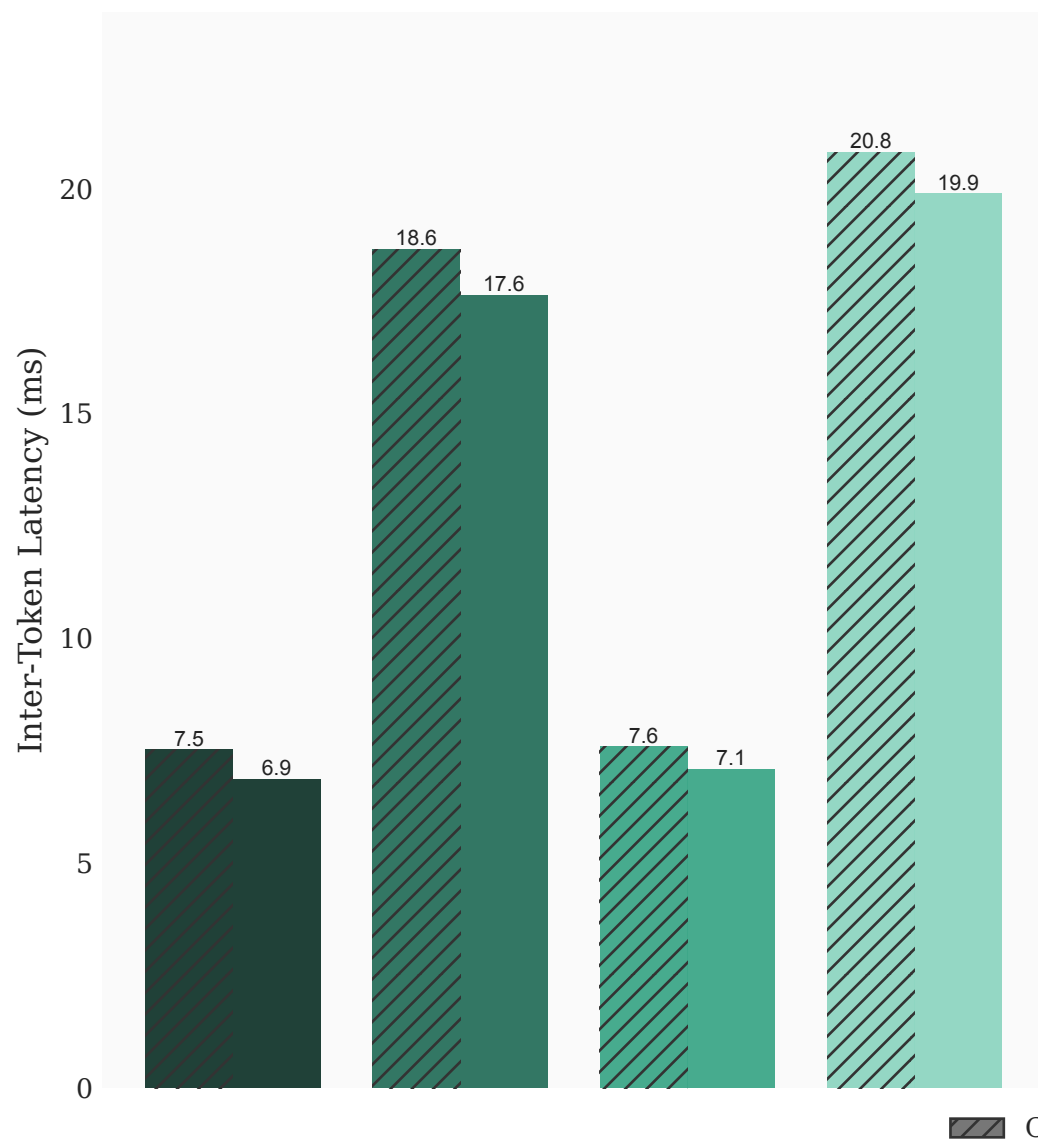
P99 ITL



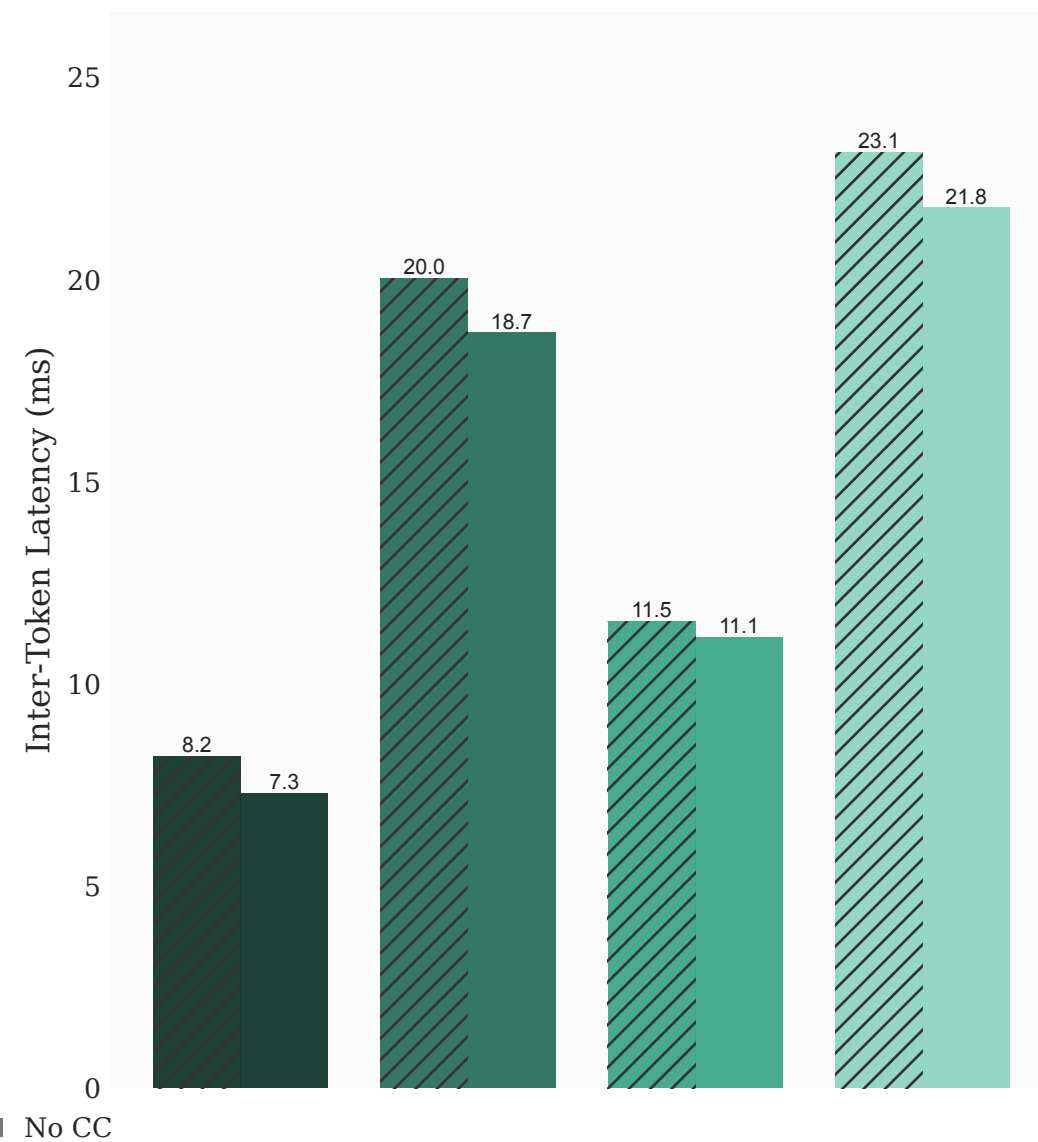
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

ShareGPT (Single Request)

Mean ITL



P99 ITL

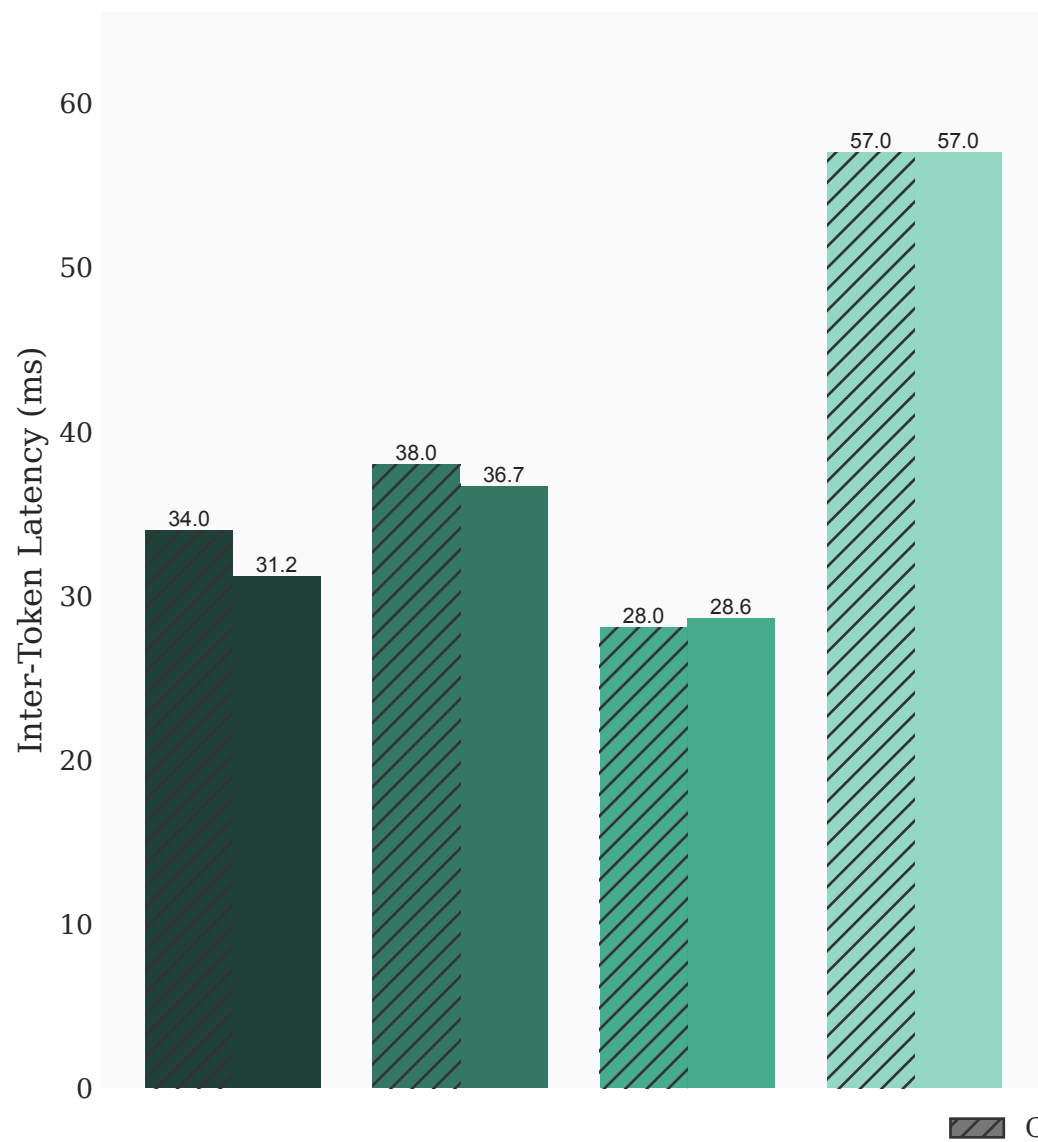


Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

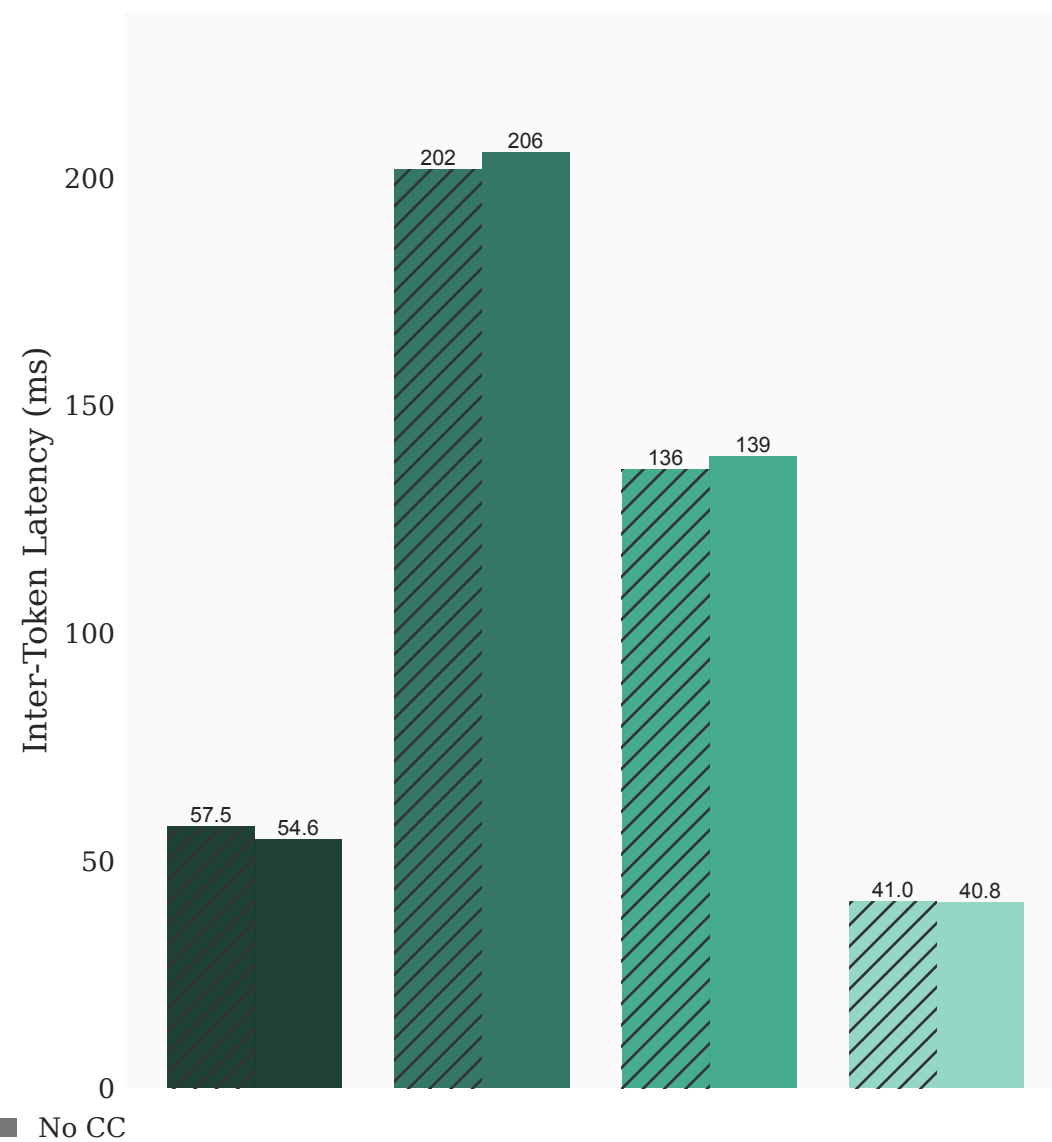
Legend: CC, No CC

Edit 10K Characters (100 Request Rate)

Mean ITL



P99 ITL

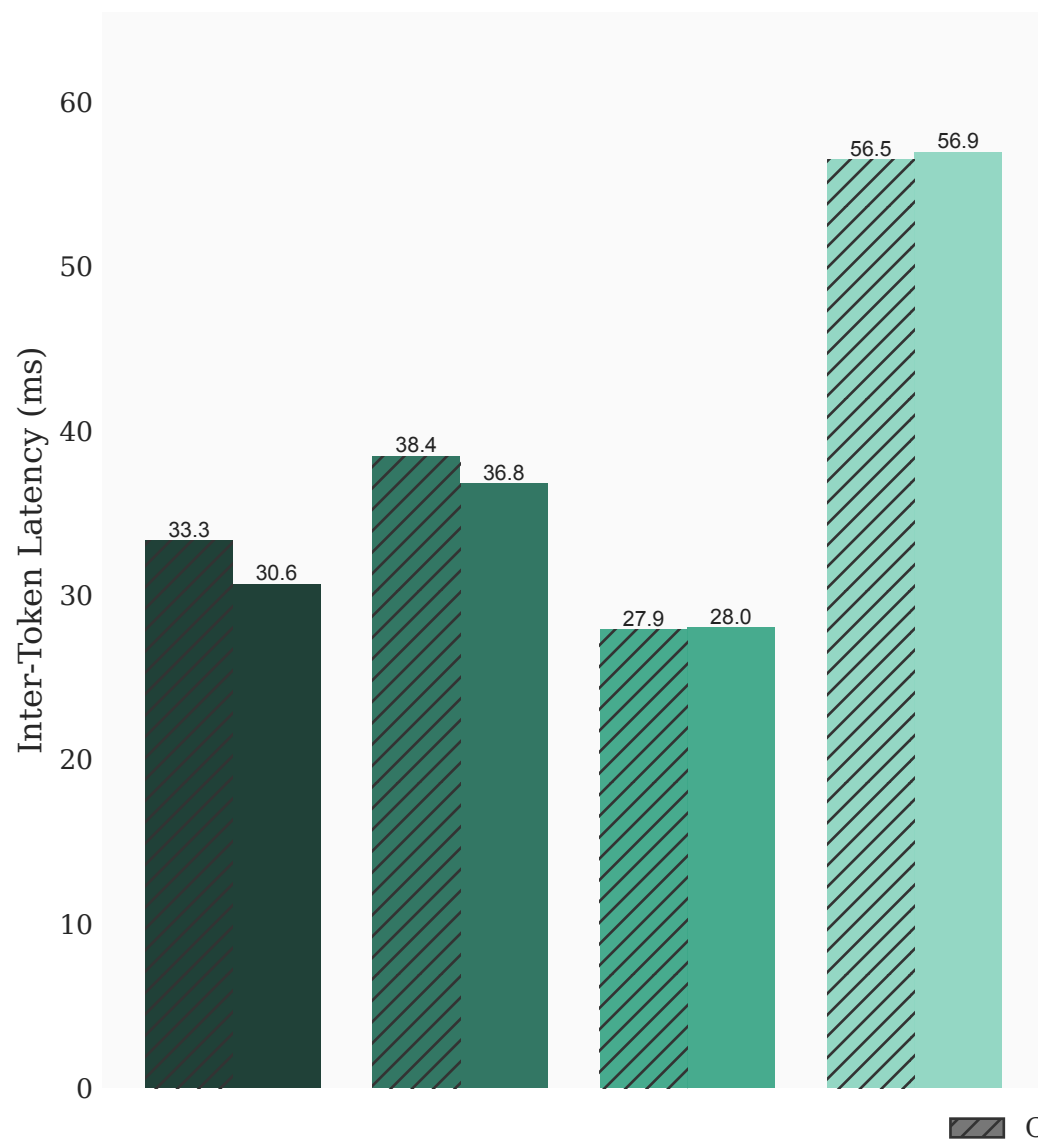


Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

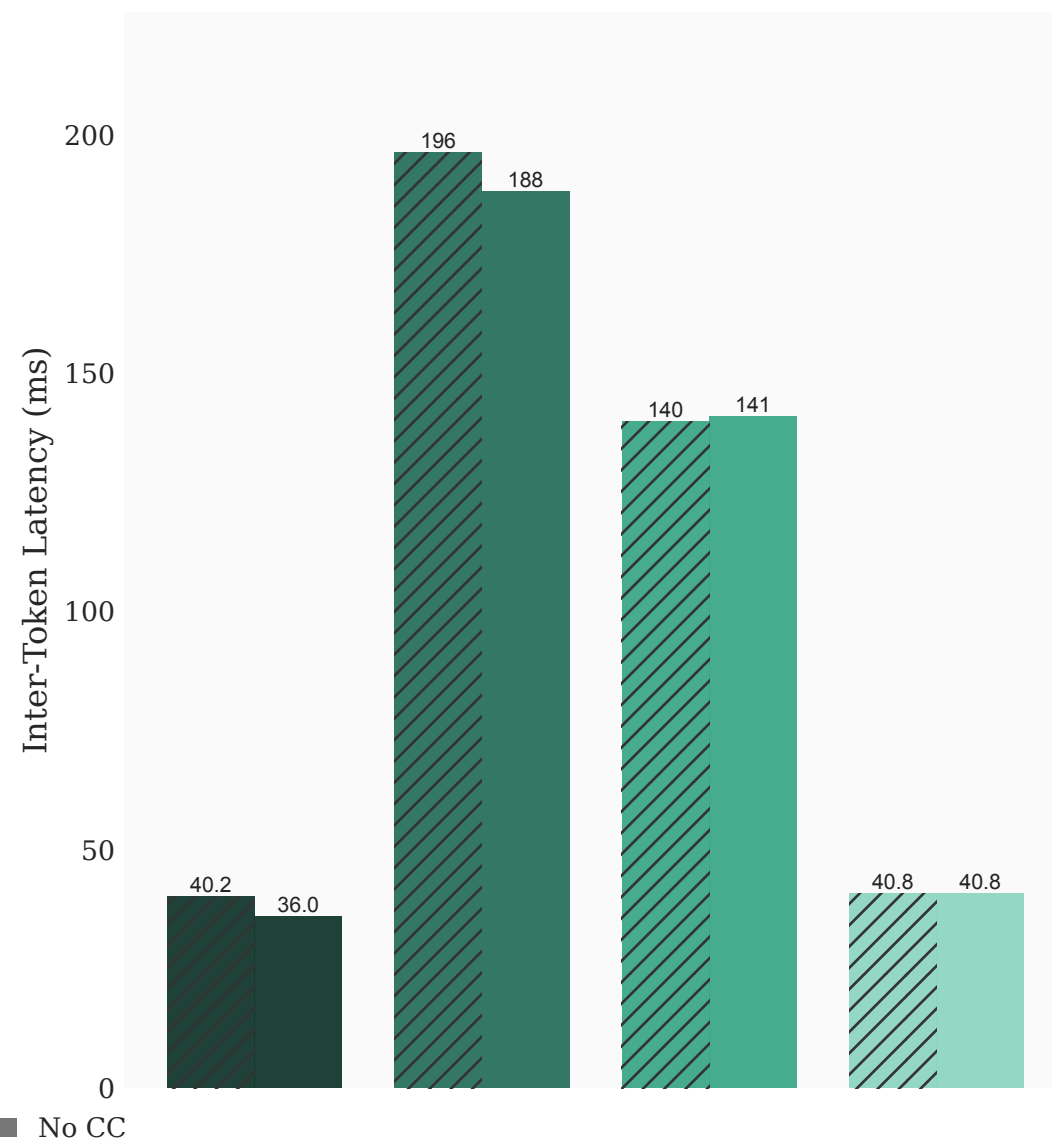
Legend: CC, No CC

Edit 10K Characters (50 Request Rate)

Mean ITL



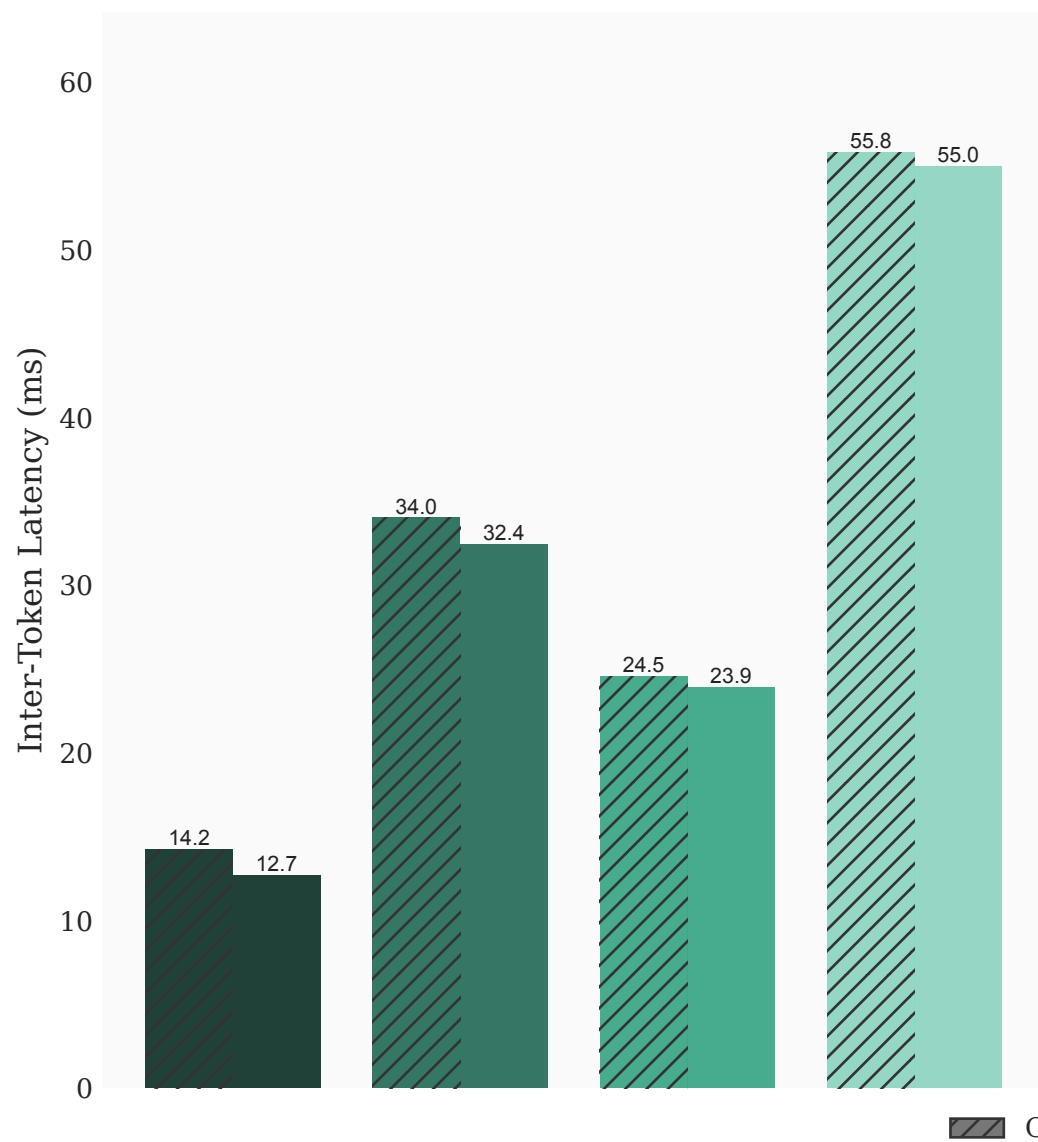
P99 ITL



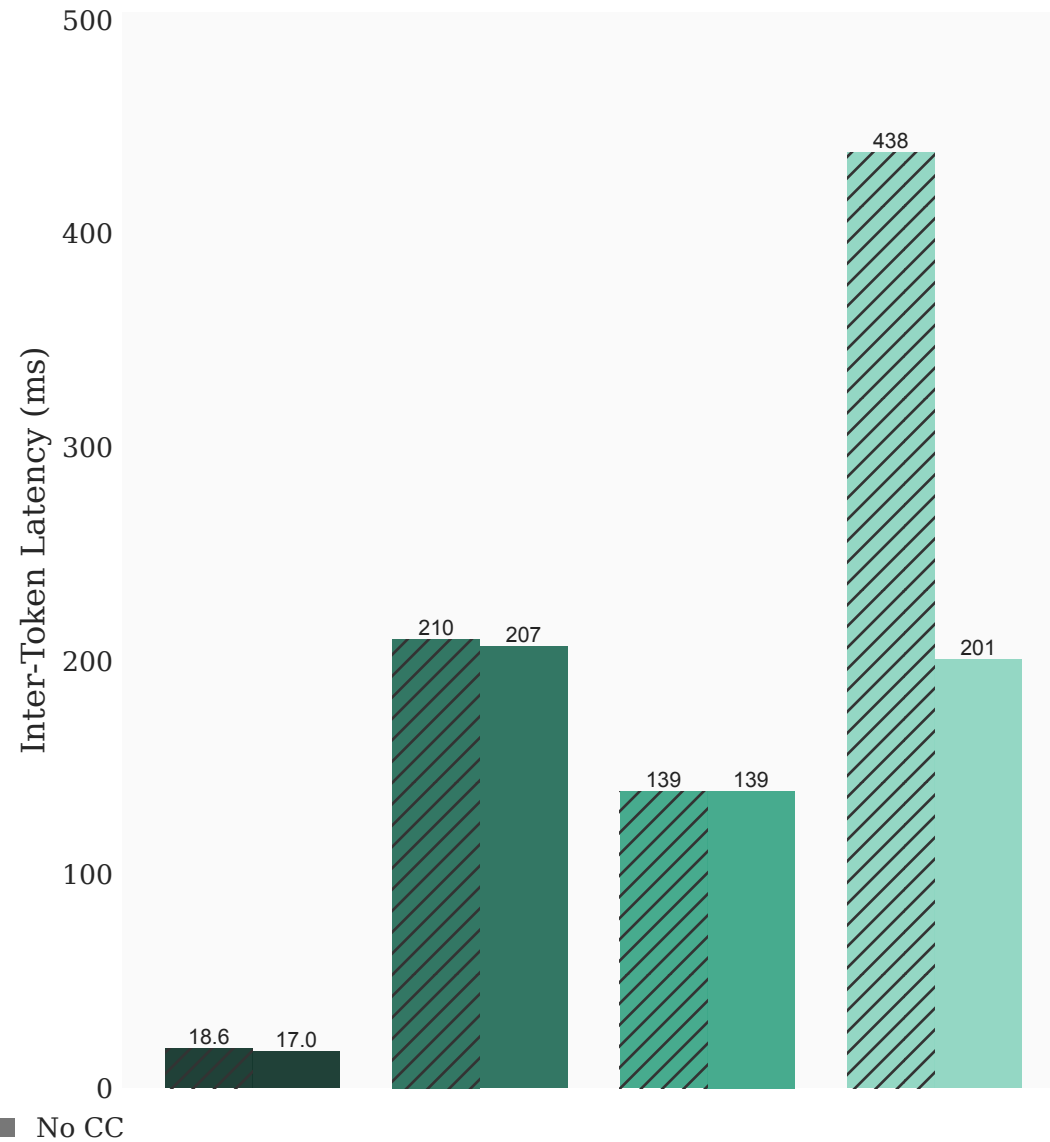
LLama 3.1 8B Mistral 3.1 24B GPT OSS 120B LLama 3.3 70B Int4

Edit 10K Characters (Single Request)

Mean ITL



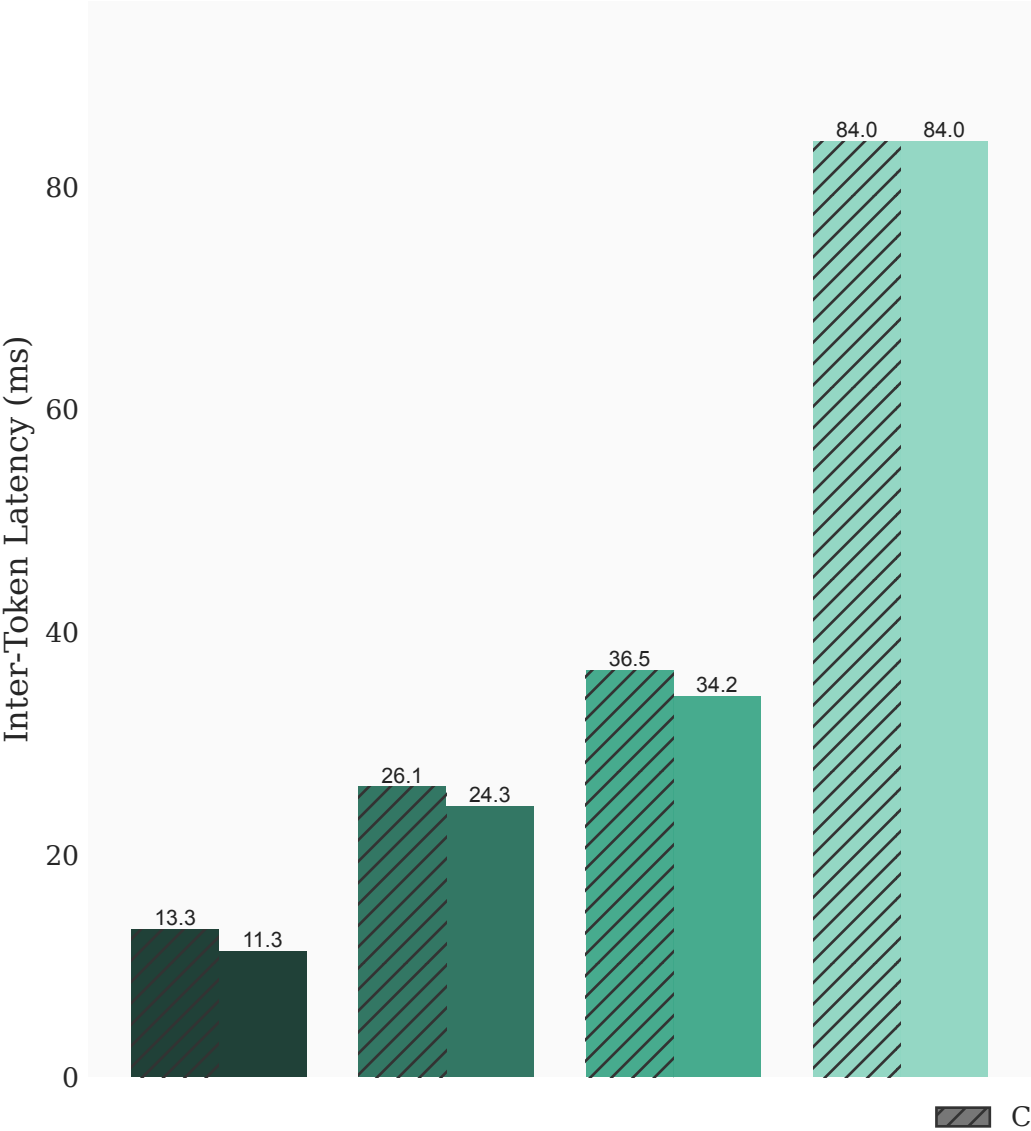
P99 ITL



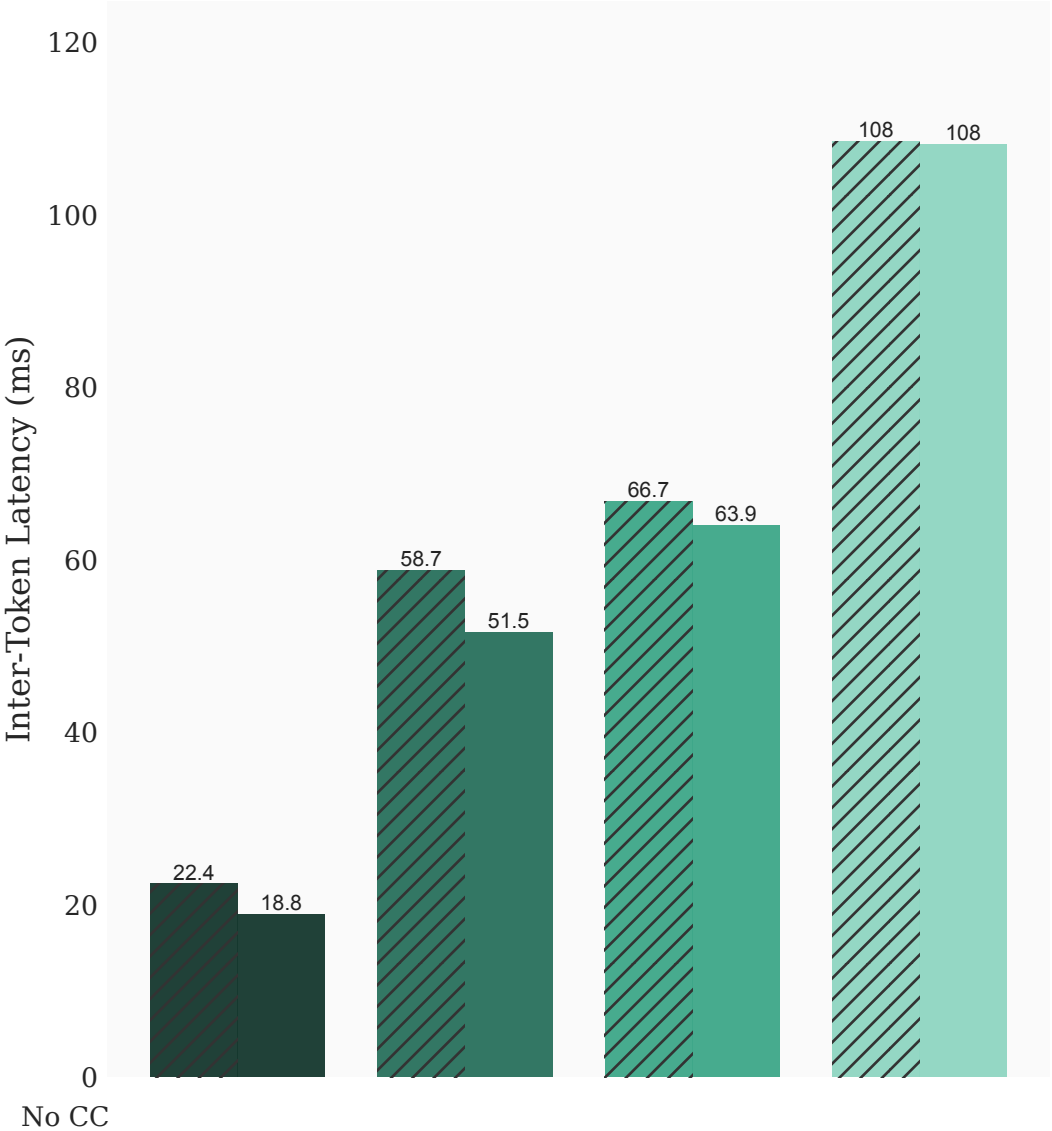
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Numina Math (100 Request Rate)

Mean ITL



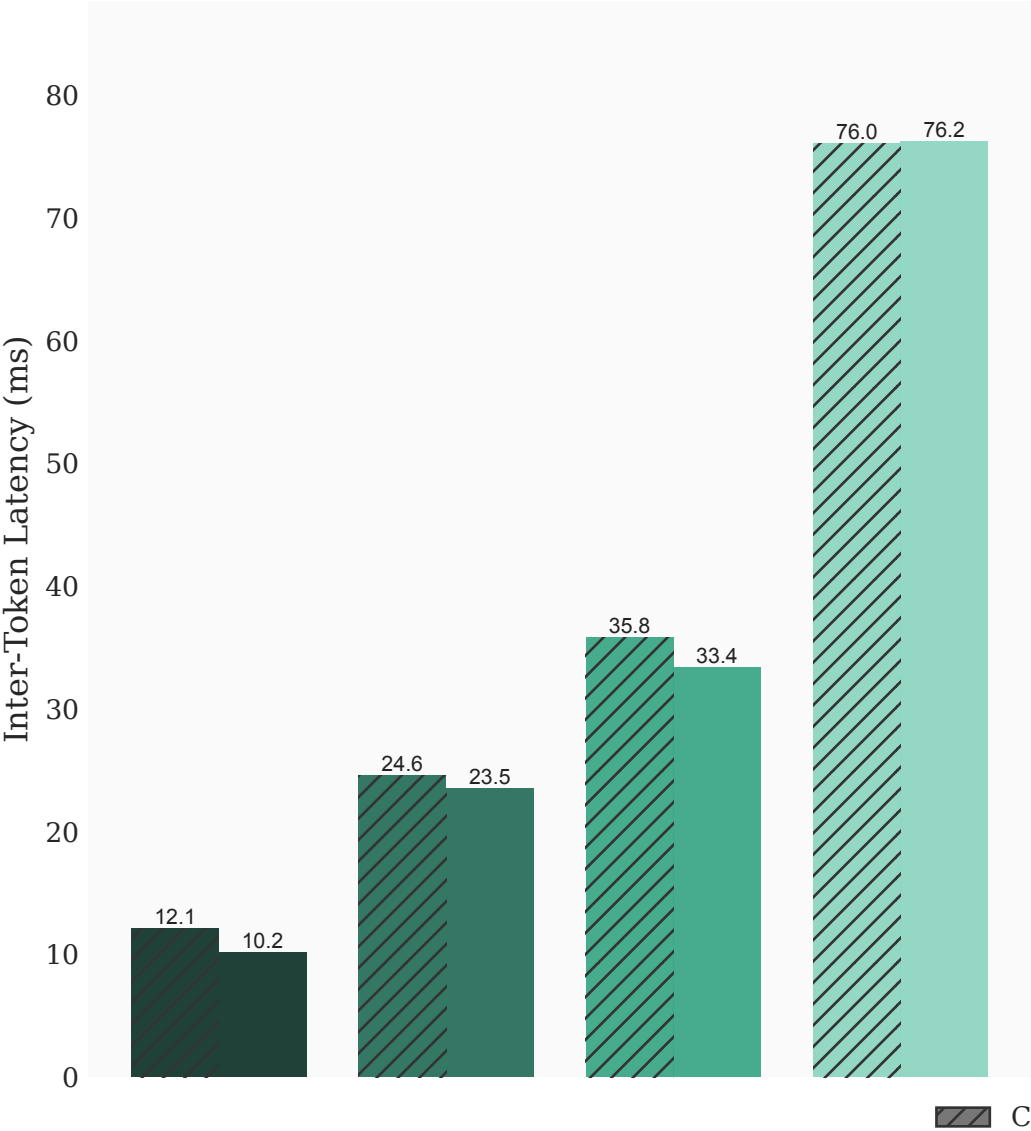
P99 ITL



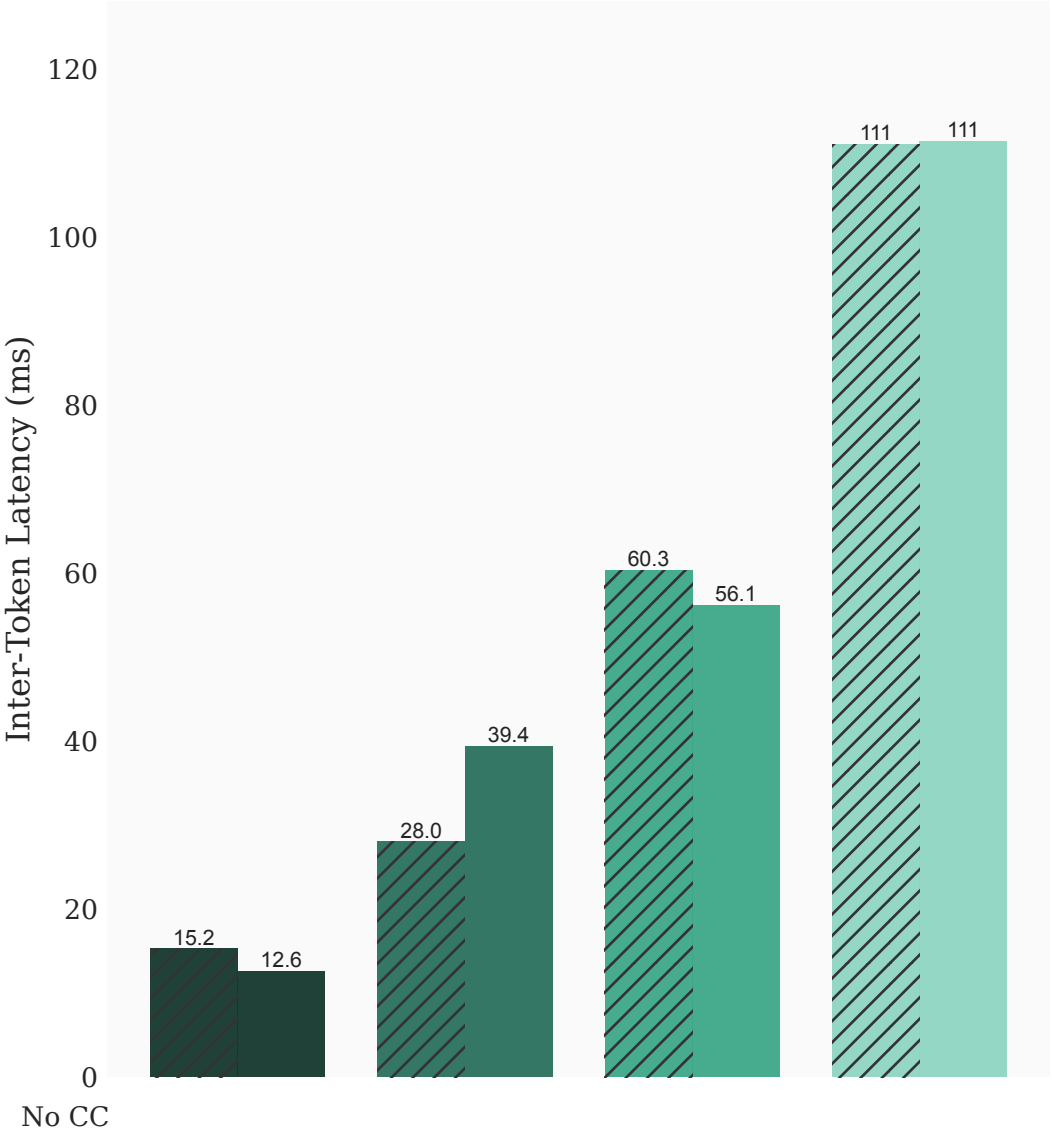
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Numina Math (50 Request Rate)

Mean ITL



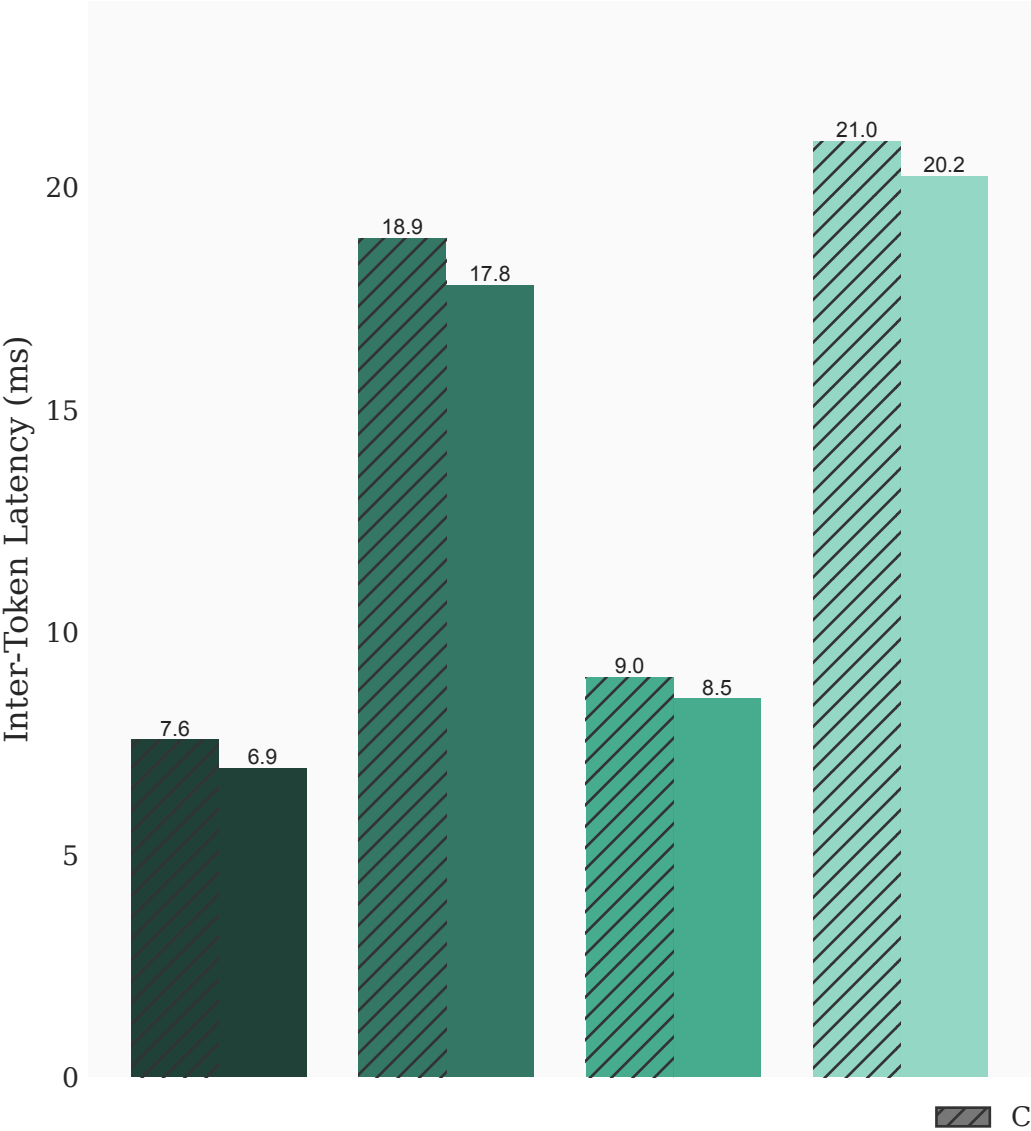
P99 ITL



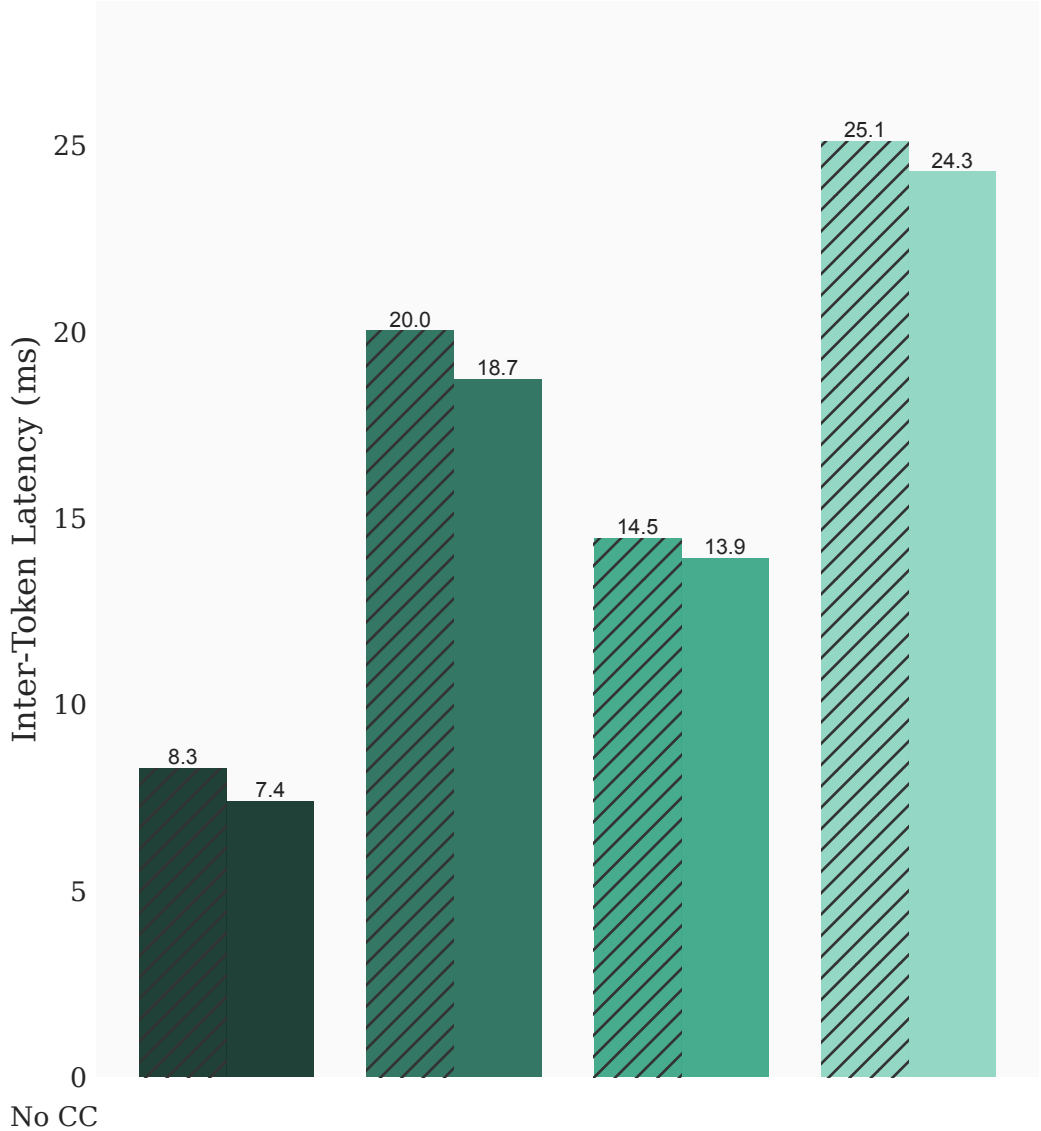
Legend: Llama 3.1 8B, Mistral 3.1 24B, GPT OSS 120B, Llama 3.3 70B Int4

Numina Math (Single Request)

Mean ITL



P99 ITL



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