# Token Throughput Comparison (IB vs No IB)



**Output Tokens/sec**

Tokens/sec vs Benchmark Type — Network: No IB, IB

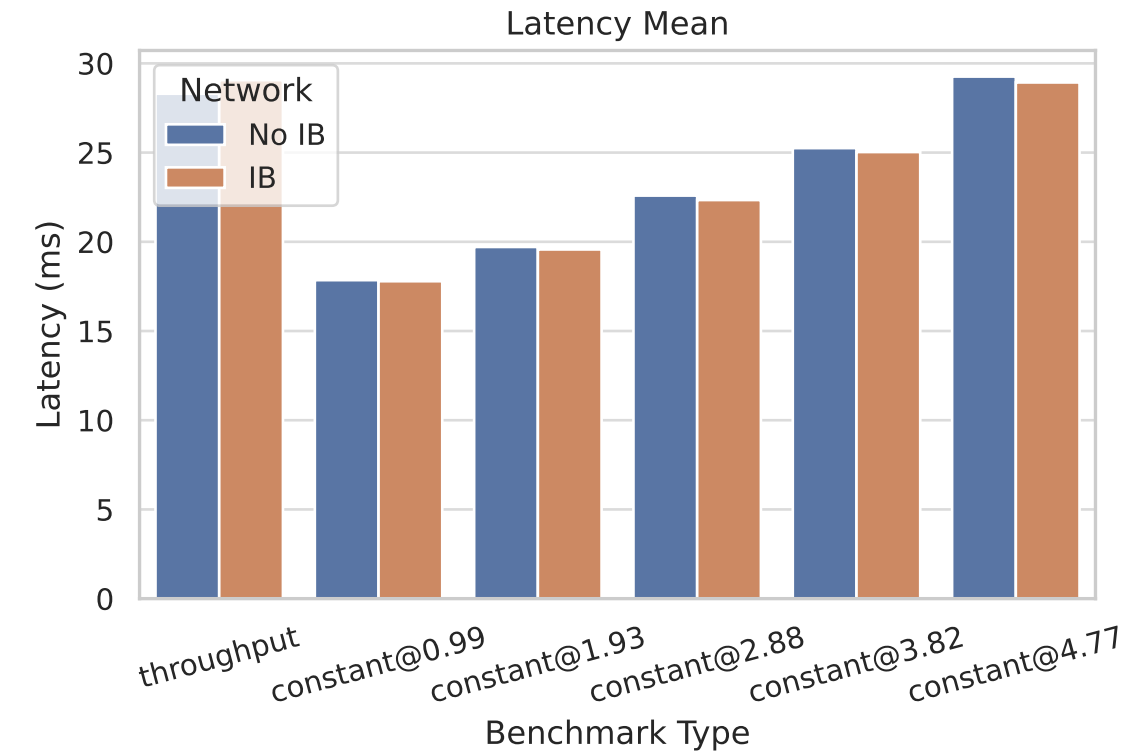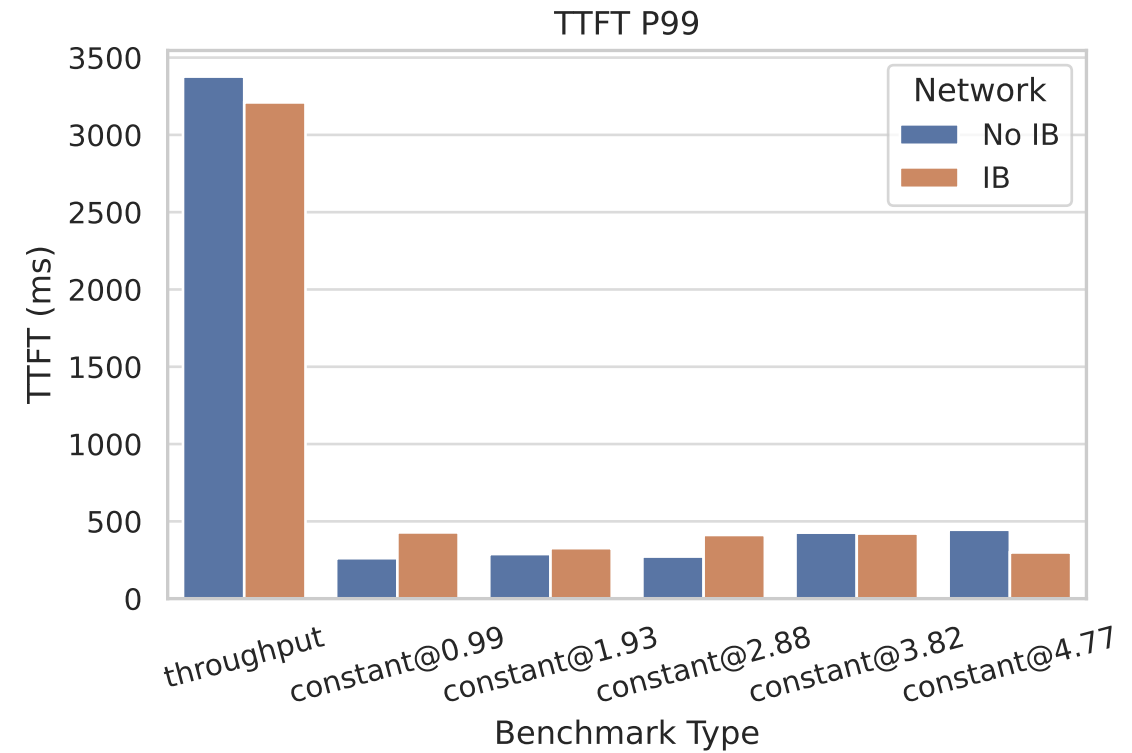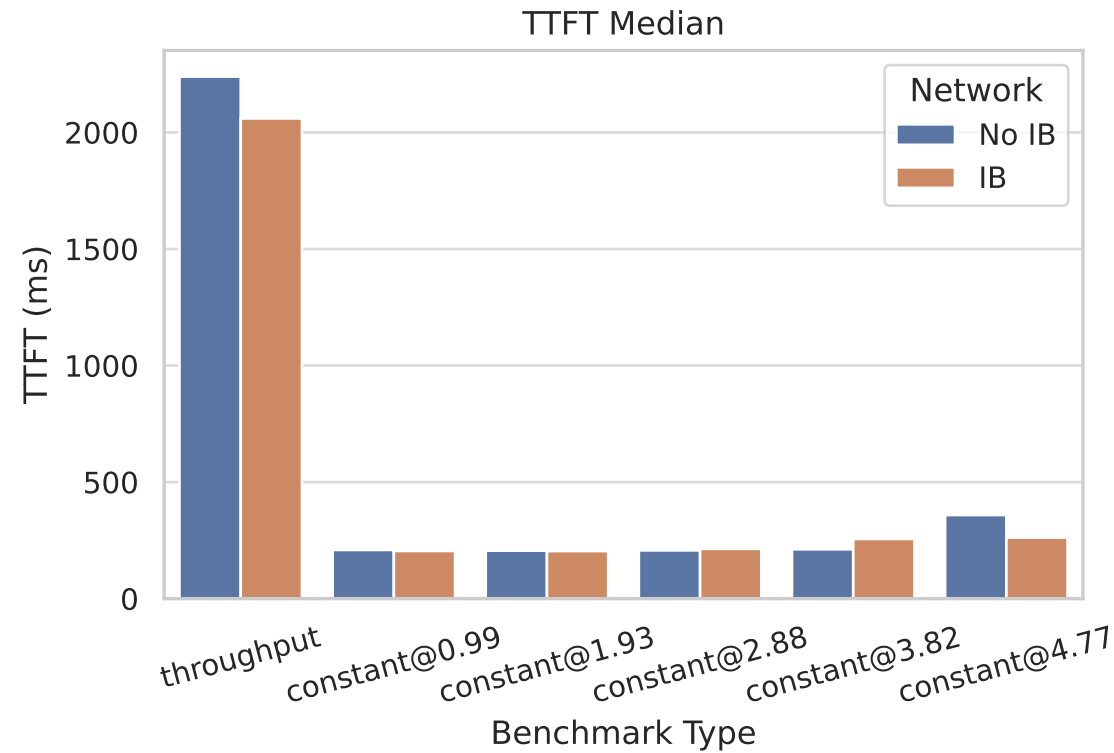Benchmark Type: throughput, constant@0.99, constant@1.93, constant@2.88, constant@3.82, constant@4.77

**Total Tokens/sec**

Tokens/sec vs Benchmark Type — Network: No IB, IB

Benchmark Type: throughput, constant@0.99, constant@1.93, constant@2.88, constant@3.82, constant@4.77

Request Latency Comparison (IB vs No IB)
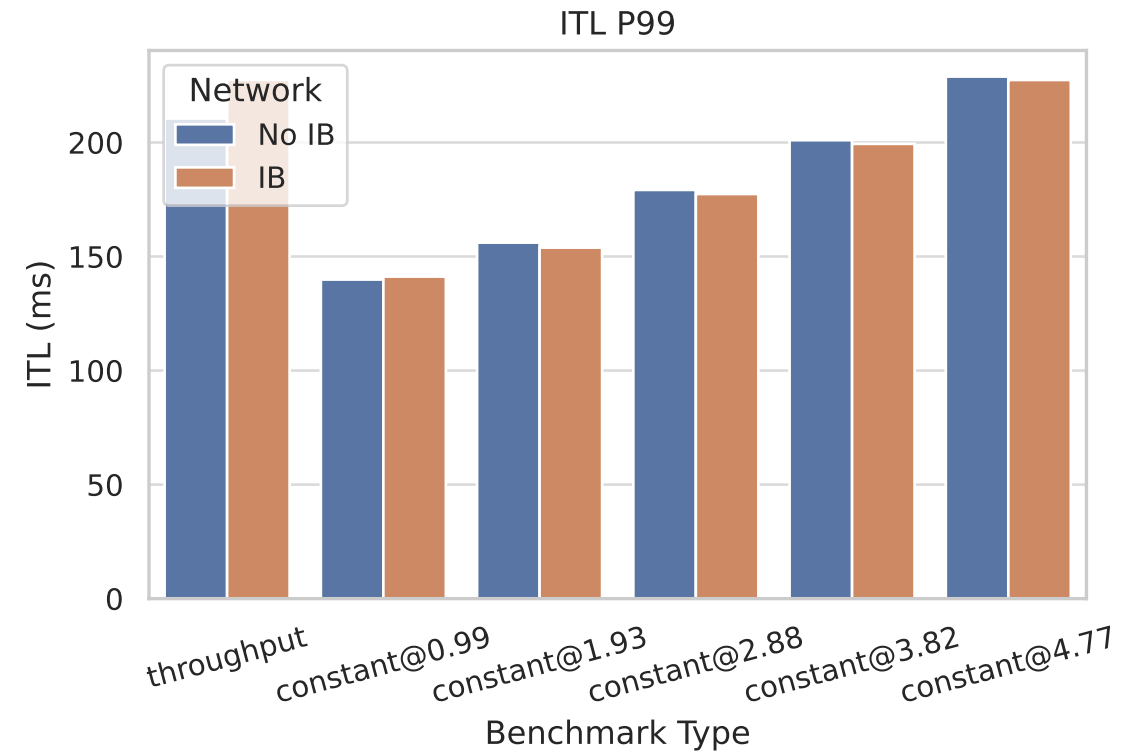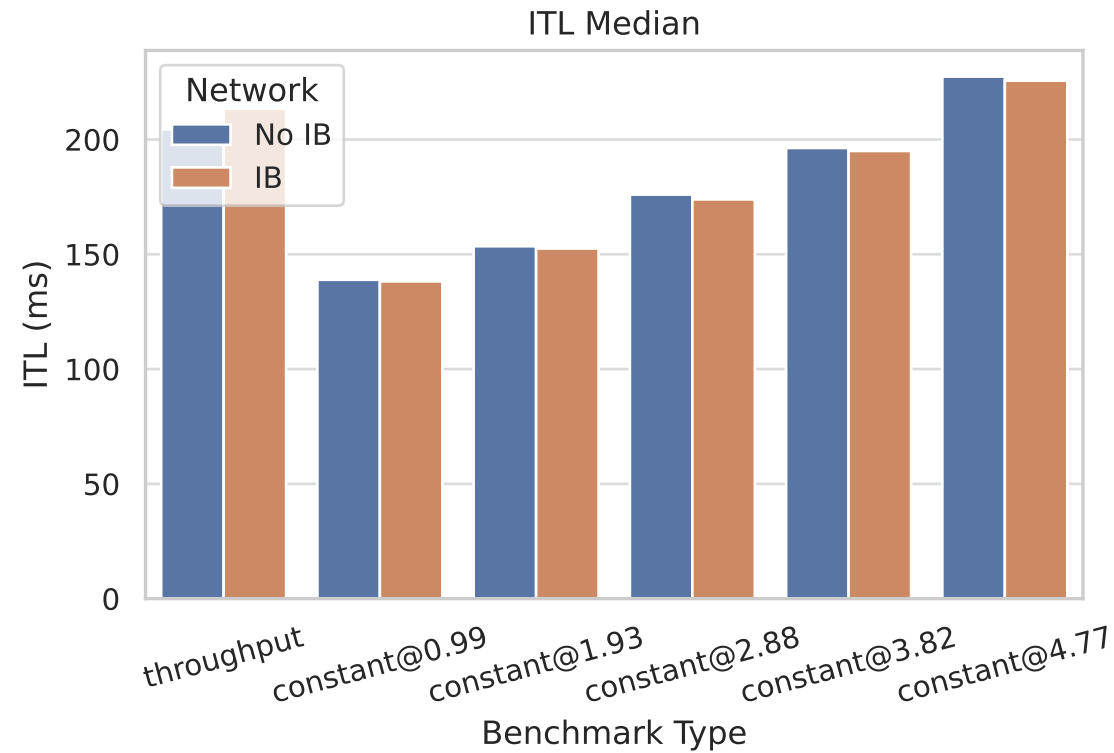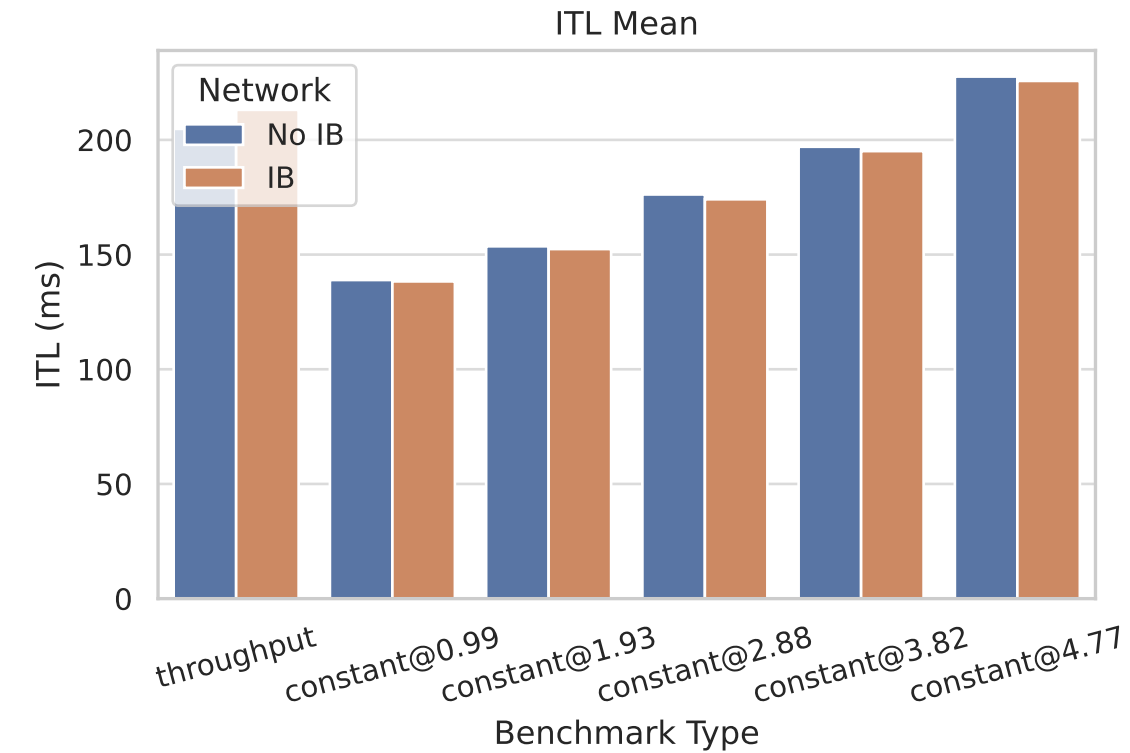
Time To First Token (TTFT) Comparison (IB vs No IB)

Inter-Token Latency (ITL) Comparison (IB vs No IB)

Time Per Output Token (TPOT) Comparison (IB vs No IB)