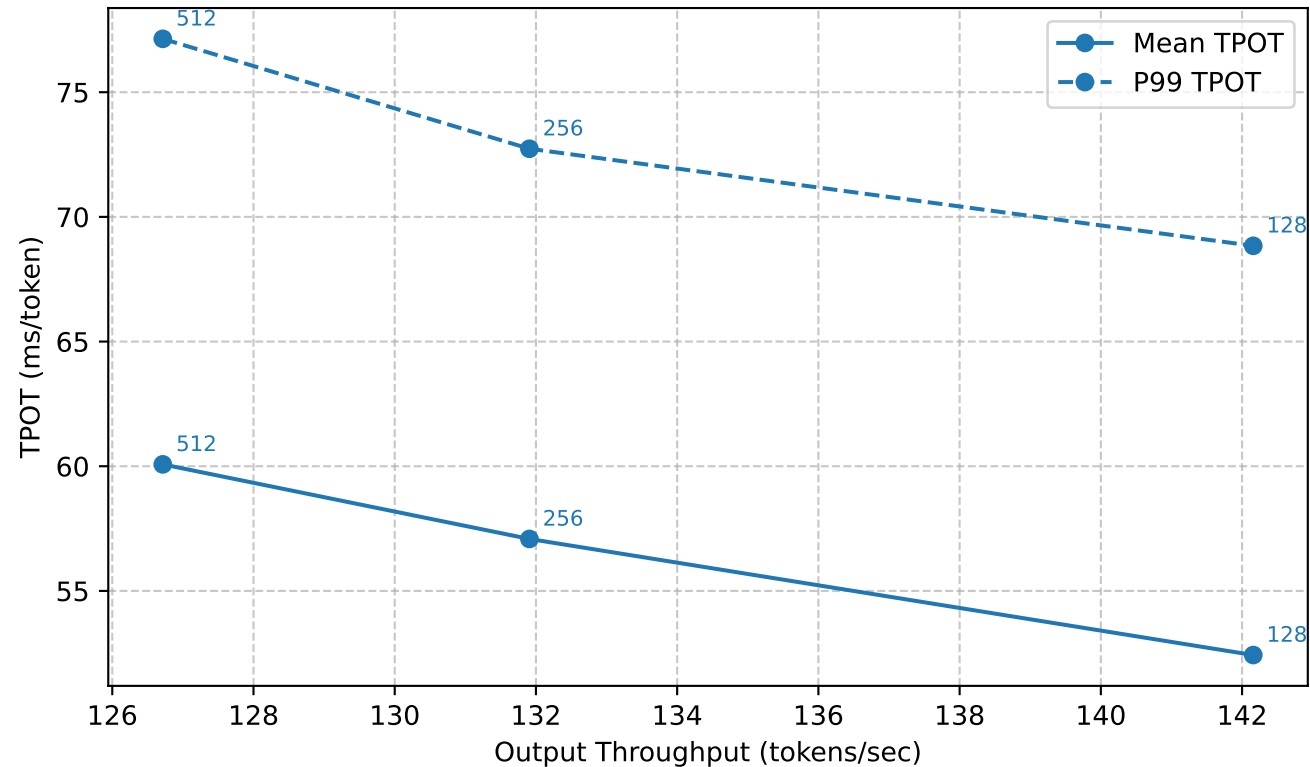


Throughput-Latency characterization  
meta-llama\_llama-3.1-70b (TP=4)  
Batch Size: 8  
Dataset: sharegpt



Throughput-Latency characterization  
meta-llama\_llama-3.1-70b (TP=4)  
Batch Size: 8  
Dataset: wildchat

