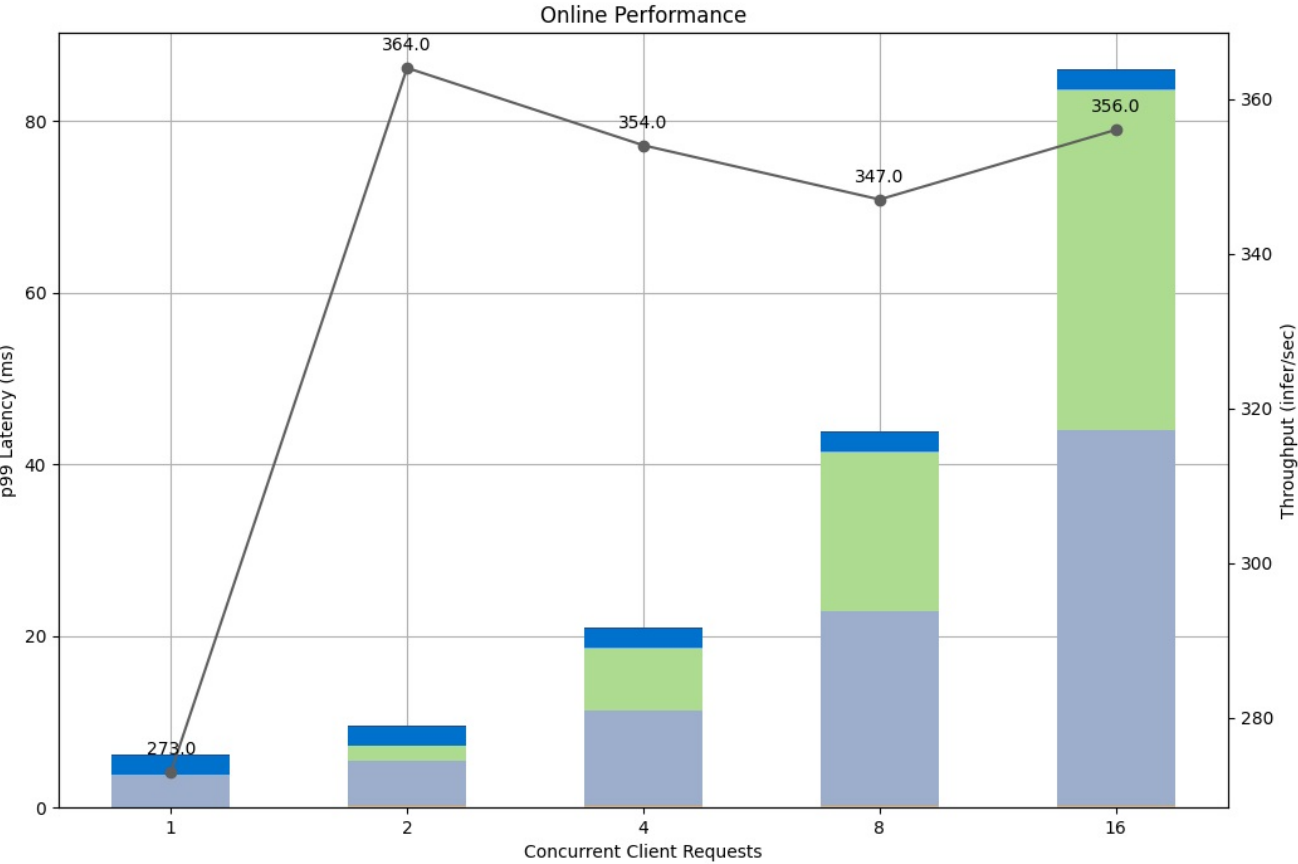
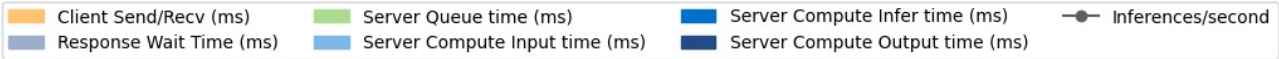
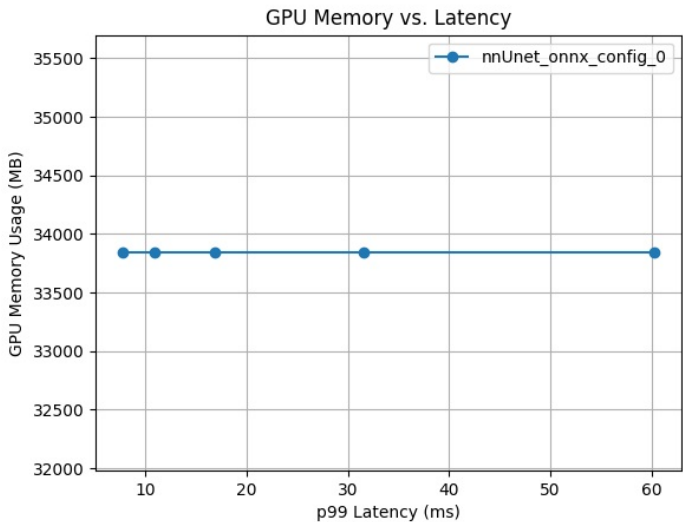


Detailed Report

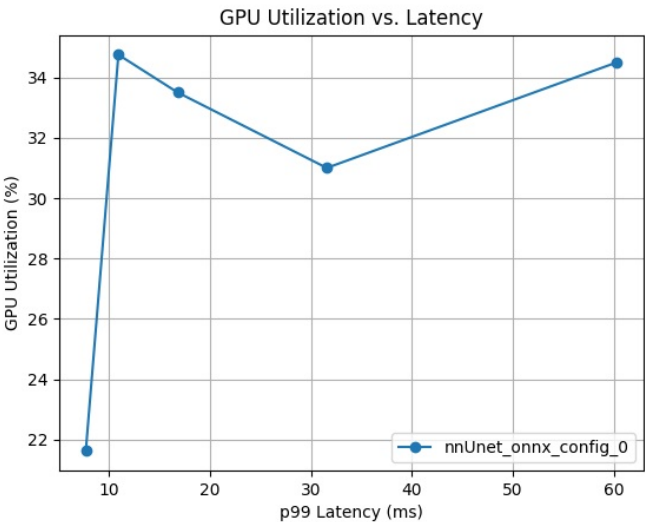
Model Config: nnUnet_onnx_config_0



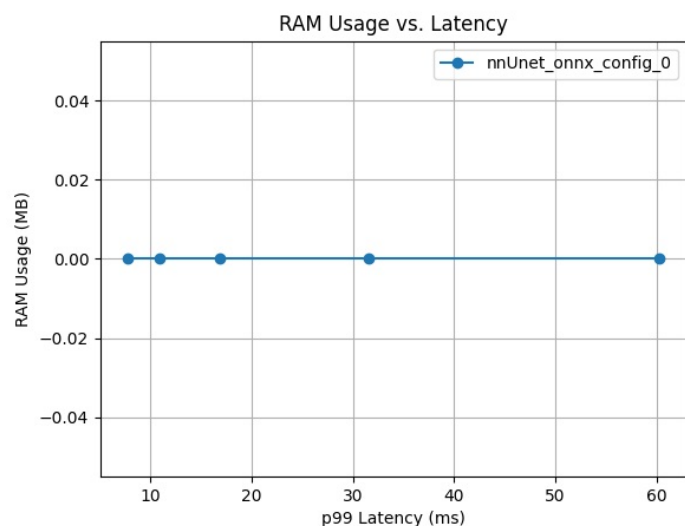
Latency Breakdown for Online Performance of nnUnet_onnx_config_0



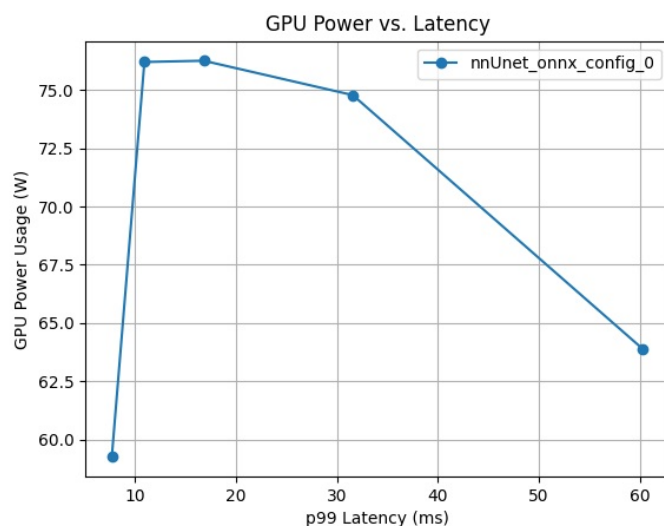
GPU Memory vs. Latency curves for config nnUnet_onnx_config_0



GPU Utilization vs. Latency curves for config nnUnet_onnx_config_0



RAM Usage vs. Latency curves for config nnUnet_onnx_config_0



GPU Power vs. Latency curves for config nnUnet_onnx_config_0

Request Concurrency	p99 Latency (ms)	Client Response Wait (ms)	Server Queue (ms)	Server Compute Input (ms)	Server Compute Infer (ms)	Throughput (infer/sec)	Max CPU Memory Usage (MB)	Max GPU Memory Usage (MB)	Average GPU Utilization (%)
16	60.301	43.806	39.659	0.121	2.149	356.0	0	33839.0	34.5
8	31.571	22.697	18.501	0.122	2.24	347.0	0	33839.0	31.0
4	16.81	11.101	7.28	0.134	2.161	354.0	0	33839.0	33.5
2	10.943	5.379	1.674	0.108	2.145	364.0	0	33839.0	34.8
1	7.717	3.575	0.052	0.112	2.232	273.0	0	33839.0	21.6

The model config "nnUnet_onnx_config_0" uses 1 GPU instance(s) with a max batch size of 1 and has dynamic batching enabled. 5 measurement(s) were obtained for the model config on GPU(s) NVIDIA A100-PCIE-40GB with memory limit(s) 39.4 GB. This model uses the platform onnxruntime_onnx.

The first plot above shows the breakdown of the latencies in the latency throughput curve for this model config. Following that are the requested configurable plots showing the relationship between various metrics measured by the Model Analyzer. The above table contains detailed data for each of the measurements taken for this model config in decreasing order of throughput.