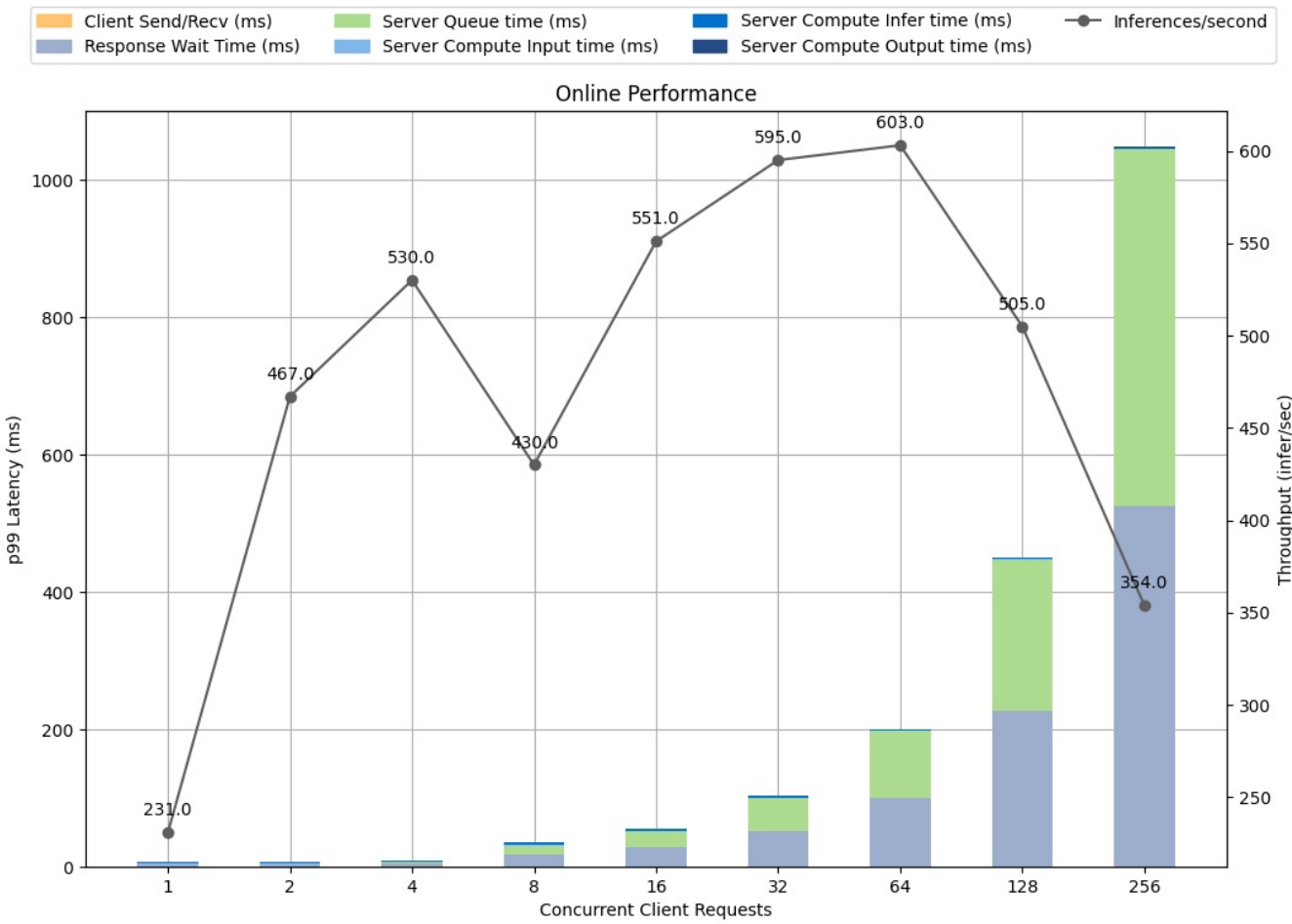
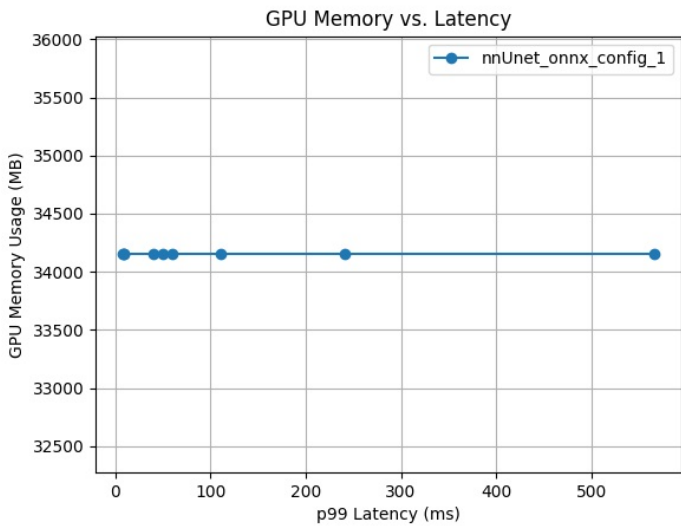


# Detailed Report

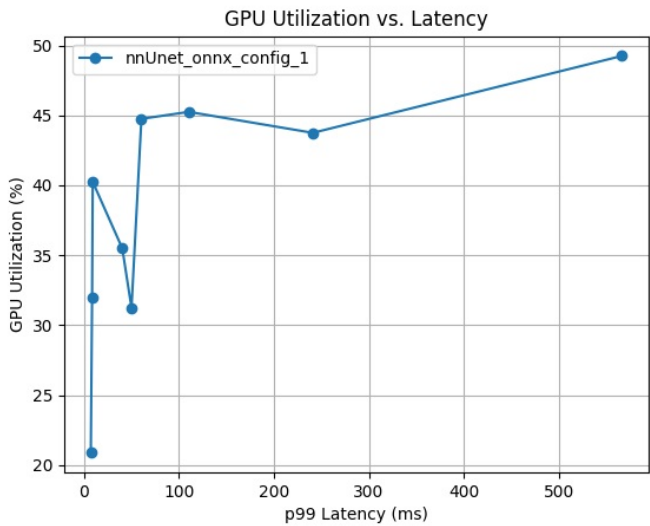
## Model Config: nnUnet\_onnx\_config\_1



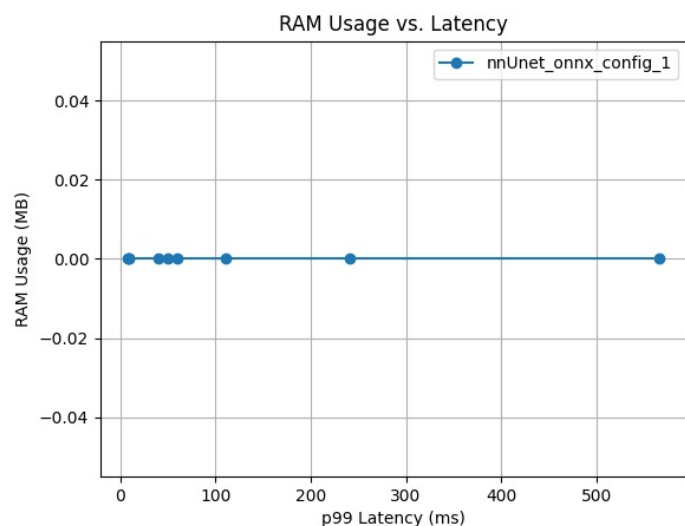
Latency Breakdown for Online Performance of nnUnet\_onnx\_config\_1



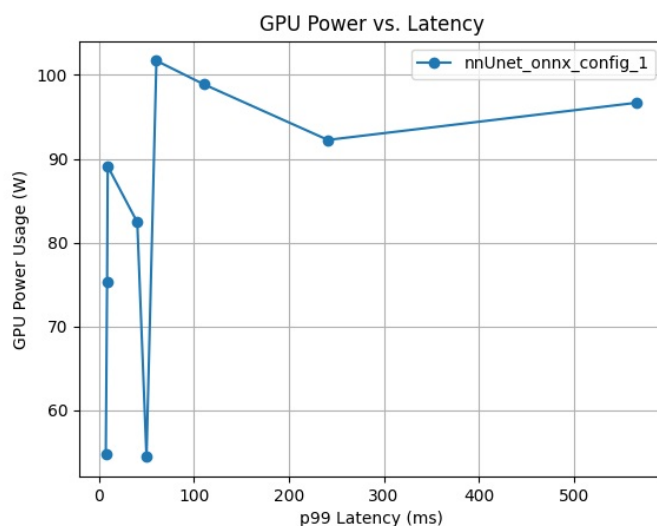
GPU Memory vs. Latency curves for config nnUnet\_onnx\_config\_1



GPU Utilization vs. Latency curves for config nnUnet\_onnx\_config\_1



RAM Usage vs. Latency curves for config nnUnet\_onnx\_config\_1



GPU Power vs. Latency curves for config nnUnet\_onnx\_config\_1

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 (%)
256	566.621	524.643	519.288	0.21	3.088	354.0	0	34150.0	49.2
128	241.032	225.888	221.303	0.21	2.807	505.0	0	34150.0	43.8
64	110.395	100.741	96.619	0.131	2.595	603.0	0	34150.0	45.2
32	60.514	52.192	47.94	0.127	2.693	595.0	0	34150.0	44.8
8	50.087	18.314	12.776	0.153	3.732	430.0	0	34150.0	31.2
16	40.356	28.444	23.583	0.152	2.945	551.0	0	34150.0	35.5
2	9.25	4.202	0.086	0.129	2.693	467.0	0	34150.0	40.2
4	8.585	5.53	0.837	0.134	3.047	530.0	0	34150.0	32.0
1	7.226	4.201	0.046	0.148	2.72	231.0	0	34150.0	20.9

The model config "nnUnet\_onnx\_config\_1" uses 2 GPU instance(s) with a max batch size of 1 and has dynamic batching enabled. 9 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.