* Tripti Singhal, "Maximizing GPU Utilization for Data Center Inference with NVIDIA TensorRT Inference Server", NVIDIA Webinar, 2019