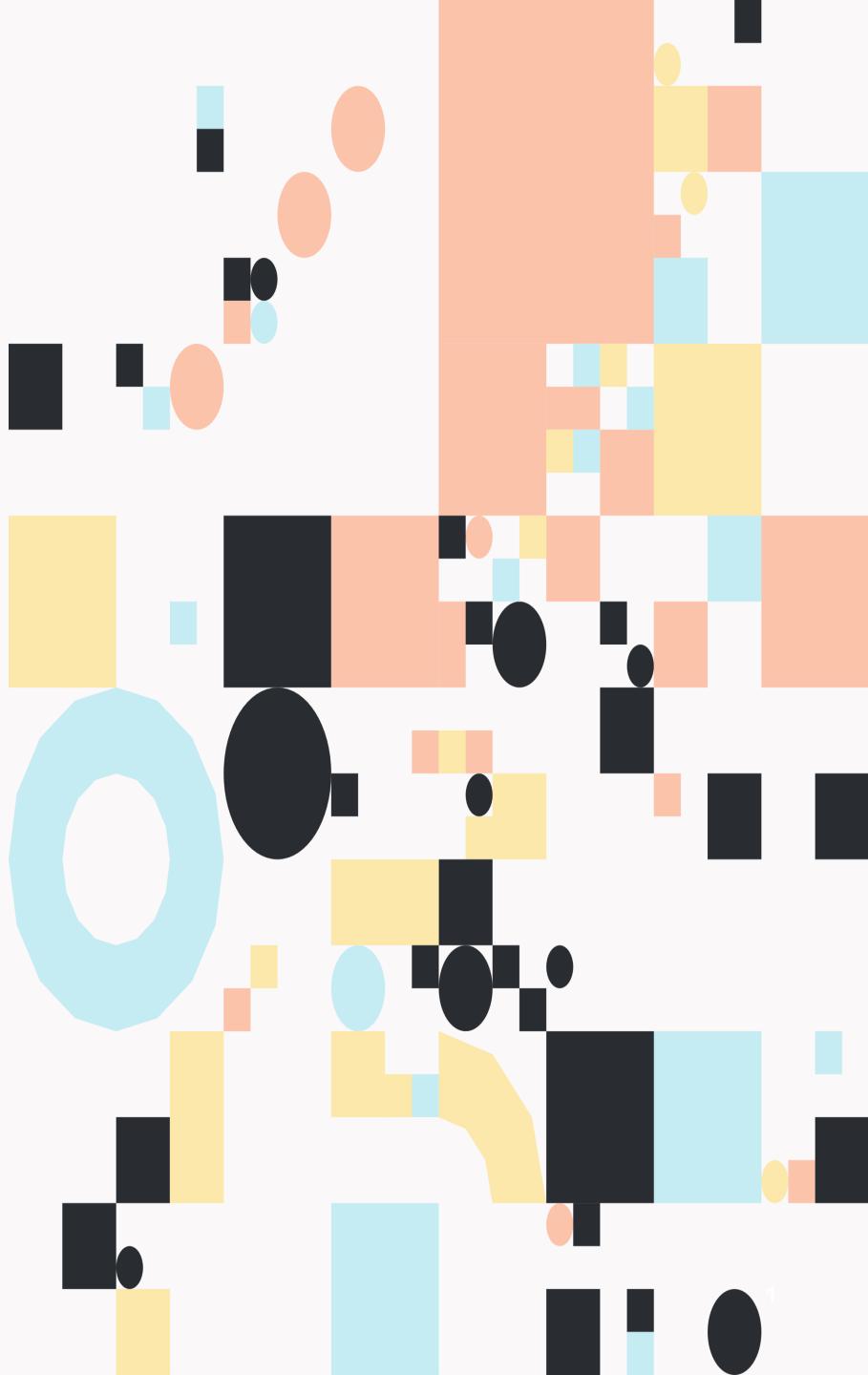


GRAPHCORE

BENCHMARKS



CONFIDENTIAL



METHODOLOGY

GPU and CPU results derived from published benchmarks and are replicated on actual hardware

IPU results derived from accurate software simulation using Poplar® tools

Training results represent forward + backward pass comparison

Best mini-batch sizes have been used for each product compared

Where results are not possible due to latency violations or other reasons, these are indicated

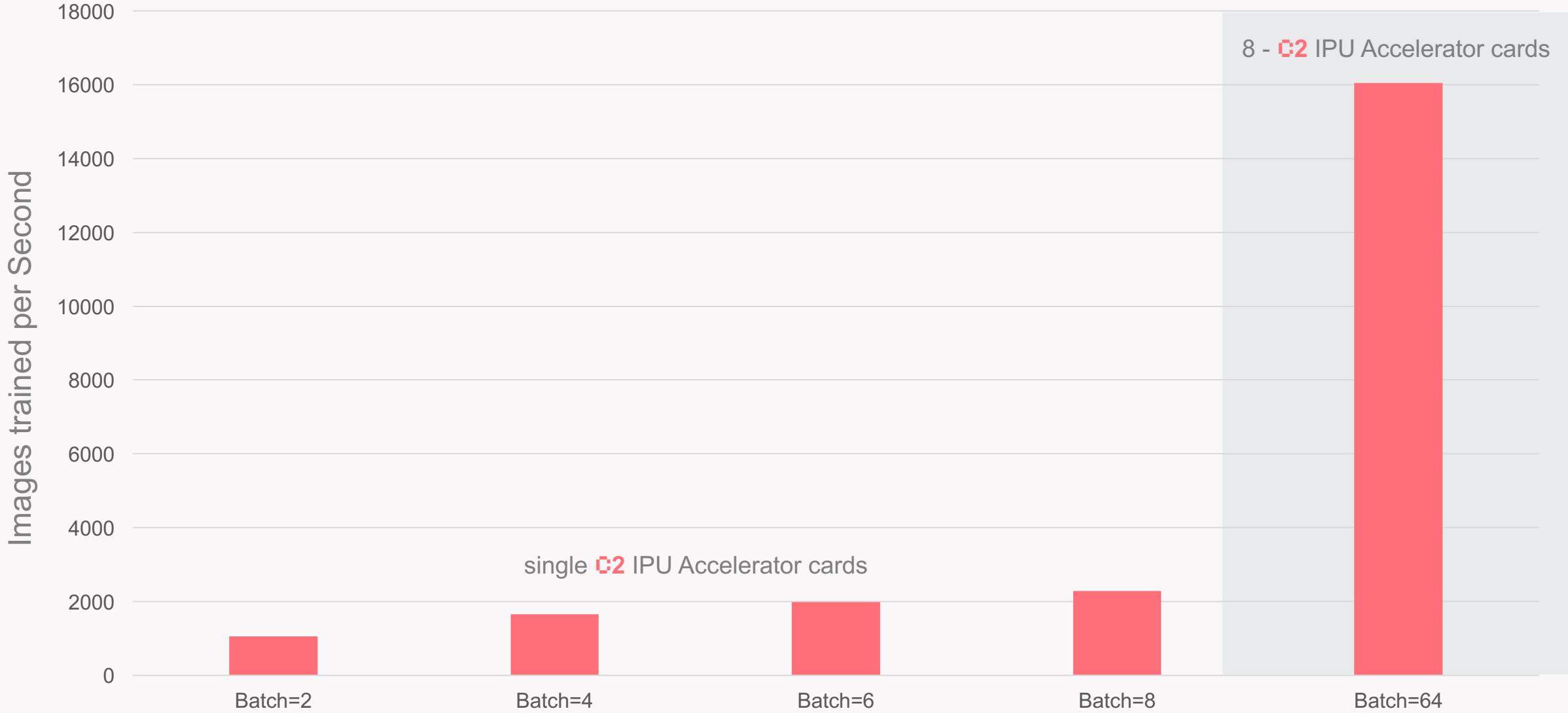
Unless otherwise stated, all comparisons are based on comparing like for like single 300W PCIe card



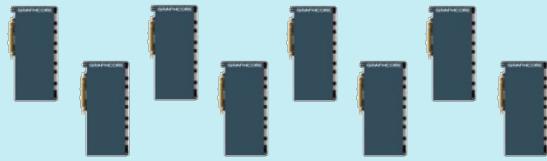
RESNET50 BENCHMARK

RESNET-50 TRAINING PERFORMANCE

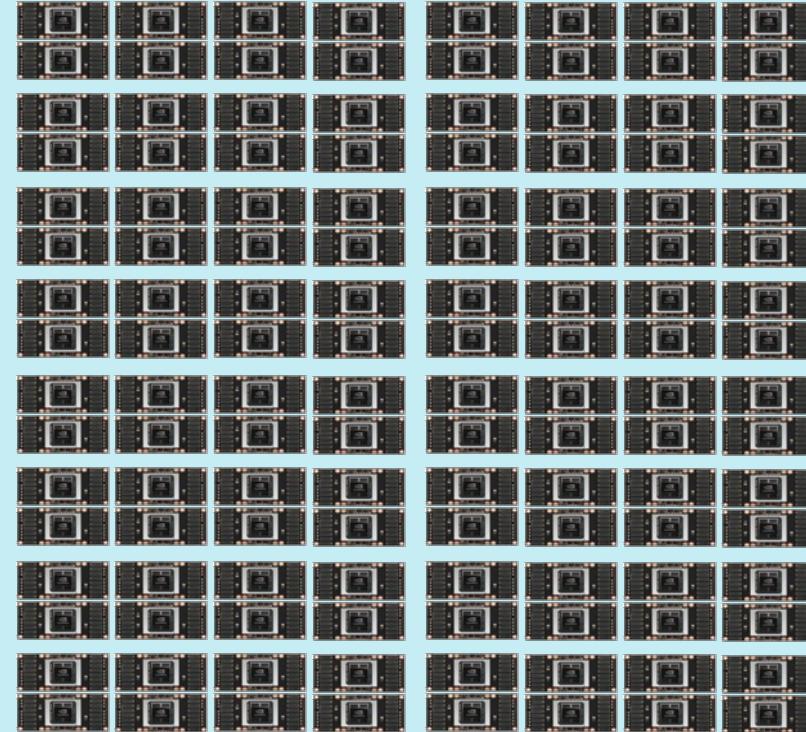
at mini-batch size of 2 to 64



RESNET-50 TRAINING COMPARISON



8 –C2 IPU Accelerator PCIe cards
(ResNet50 training @ 16,000 images/sec)

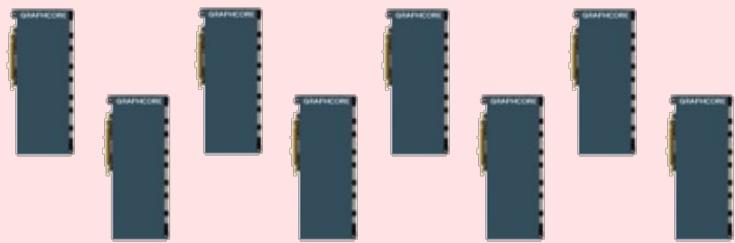


128 –NVIDIA Pascal P100
(ResNet50 training @ 16,000 images/sec)



RESNET-50 TRAINING COMPARISON

NVIDIA Volta – 2.4x vs. P100 for ResNet-50 training (*Source: Nvidia*)

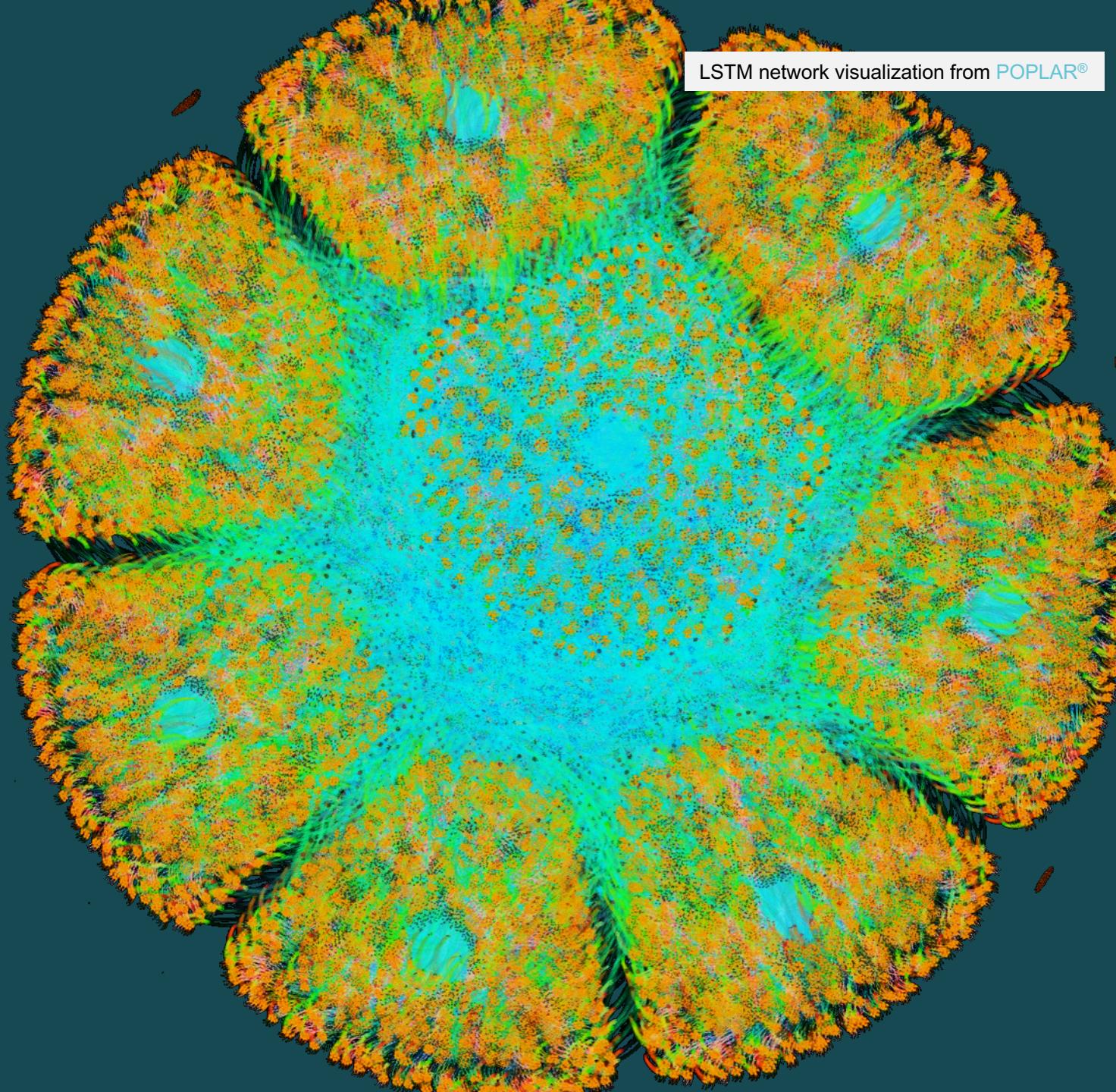


8 – C2 IPU Accelerator PCIe cards
(ResNet50 training @ 16,000 images/sec)



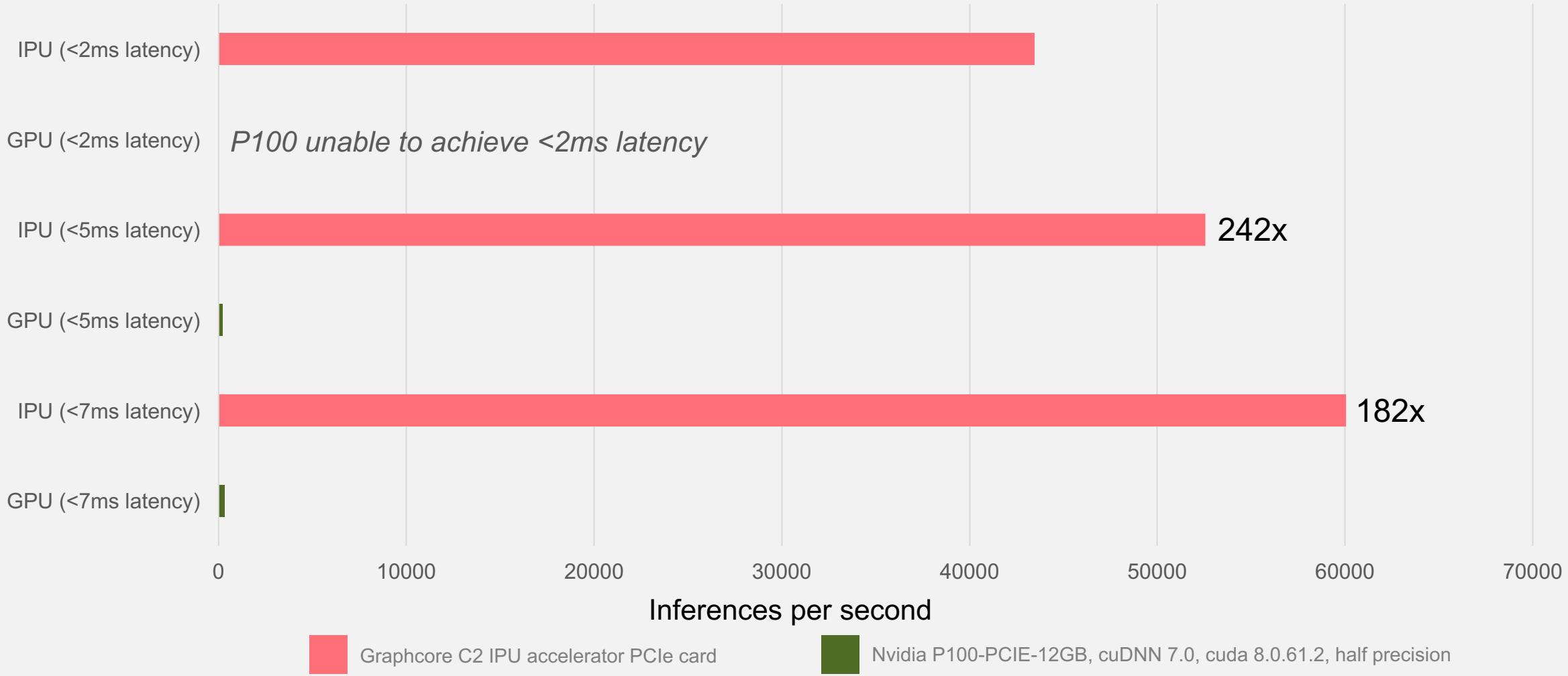
54 – NVIDIA Volta V100
(ResNet50 training @ ~16,000 images/sec)

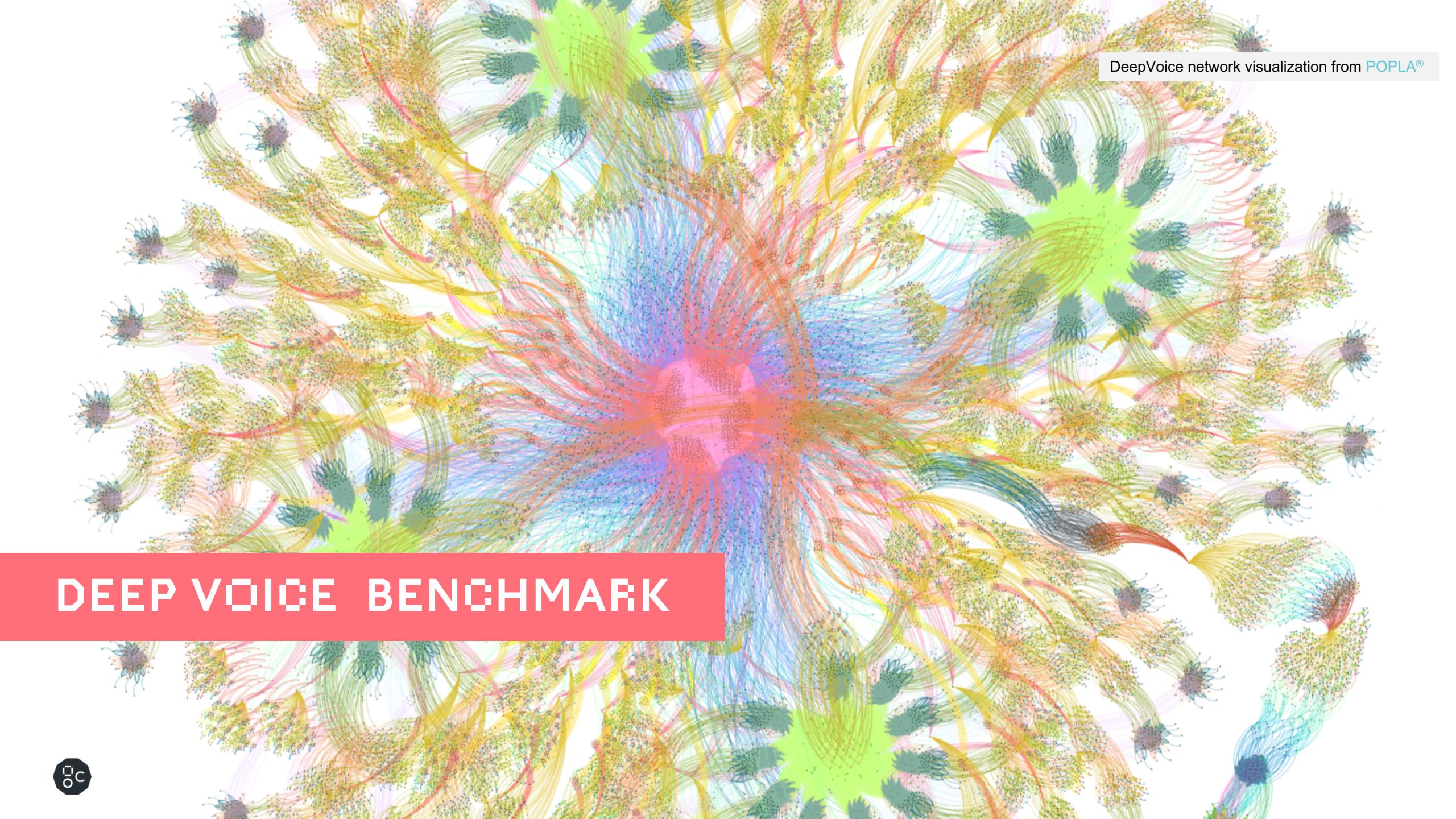
DEEPBENCH LSTM BENCHMARK



DEEPBENCH LSTM RNN: INFERENCE

Hidden Units 1536 | Time Steps 50





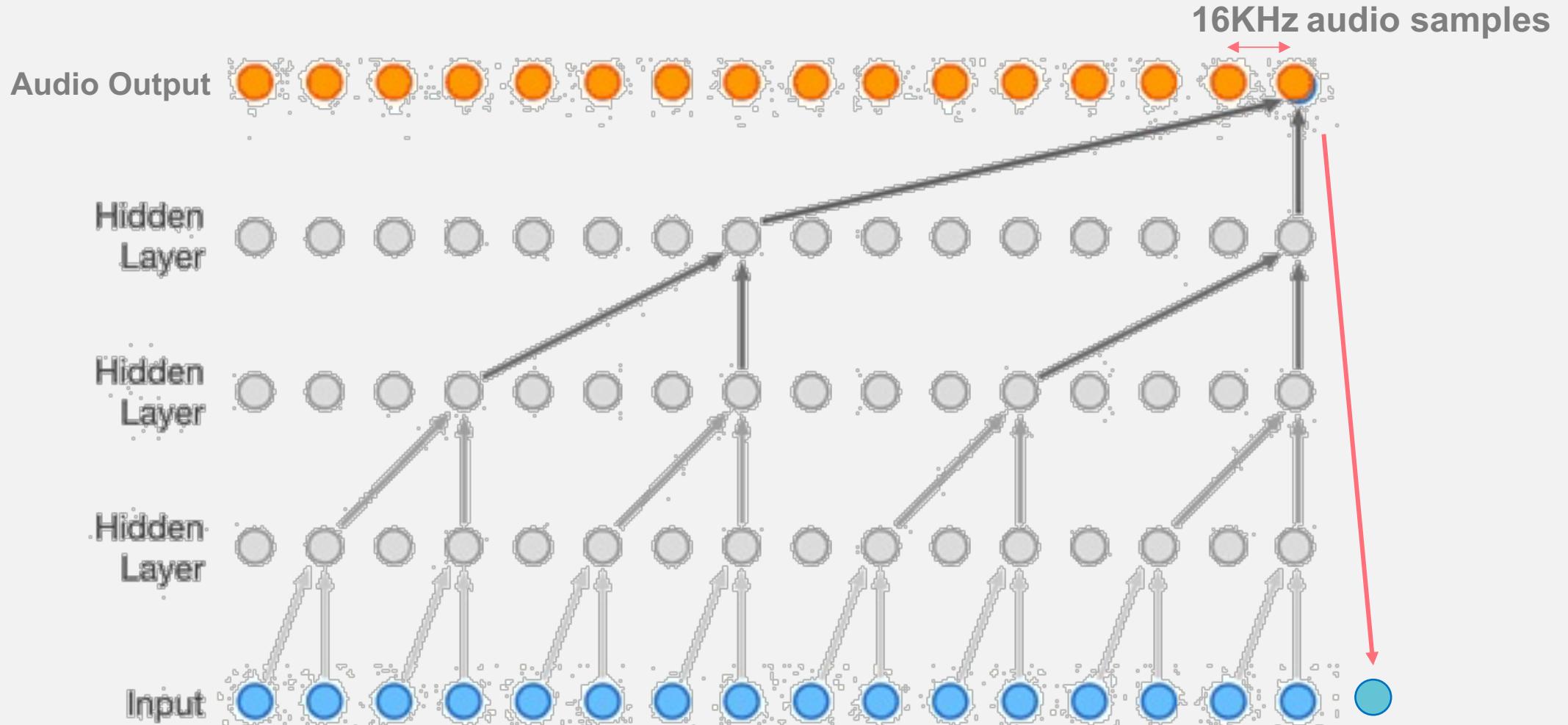
DeepVoice network visualization from POPLA®

DEEP VOICE BENCHMARK



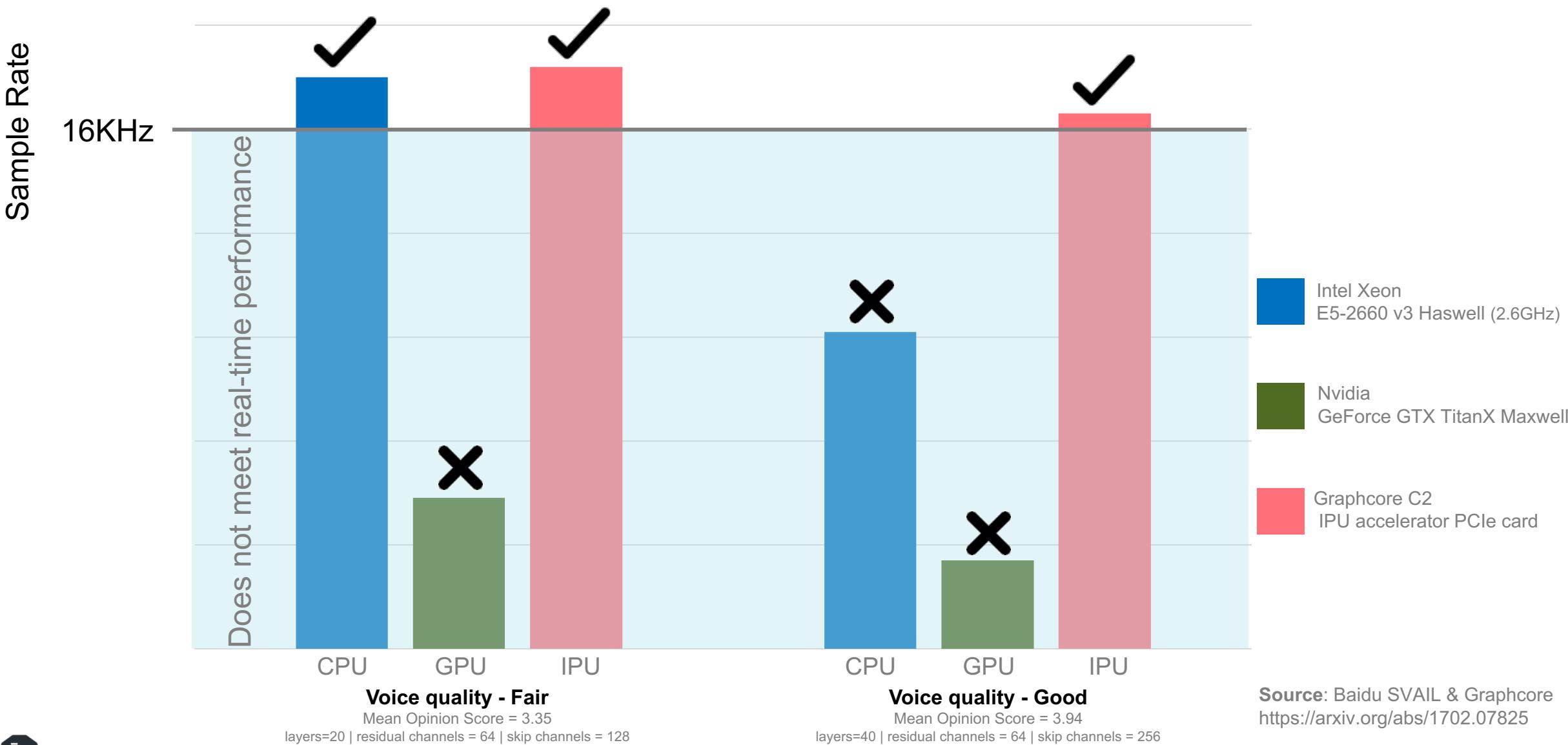
DEEP VOICE – WAVENET

Generative model for speech & audio synthesis



Source: <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>

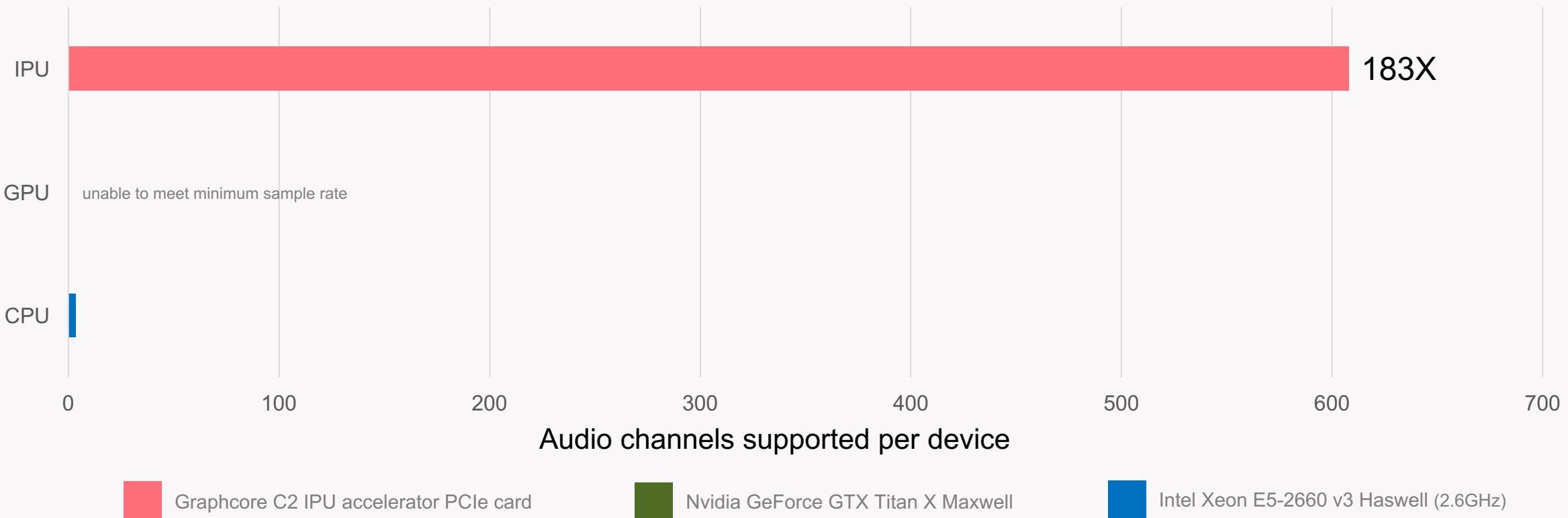
DEEPOICE: SAMPLE RATE



DEEP VOICE - #AUDIO CHANNELS

Voice quality – Fair

Mean Opinion Score = 3.35 | layers=20 | residual channels = 64 | skip channels = 128



THANK YOU

info@graphcore.ai