



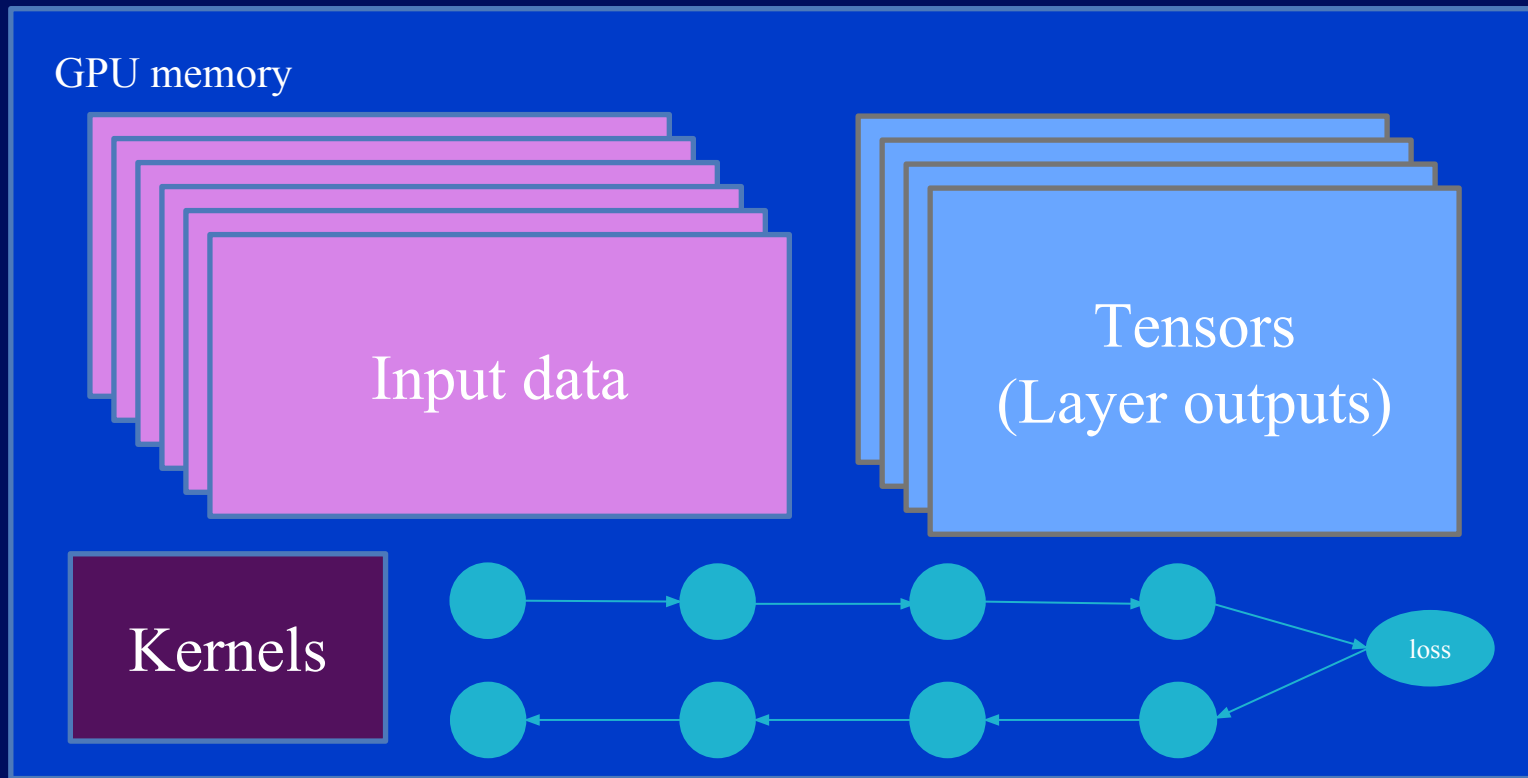
Using Tensor Swapping and NVLink to Overcome GPU Memory Limits with TensorFlow

Sam Matzek

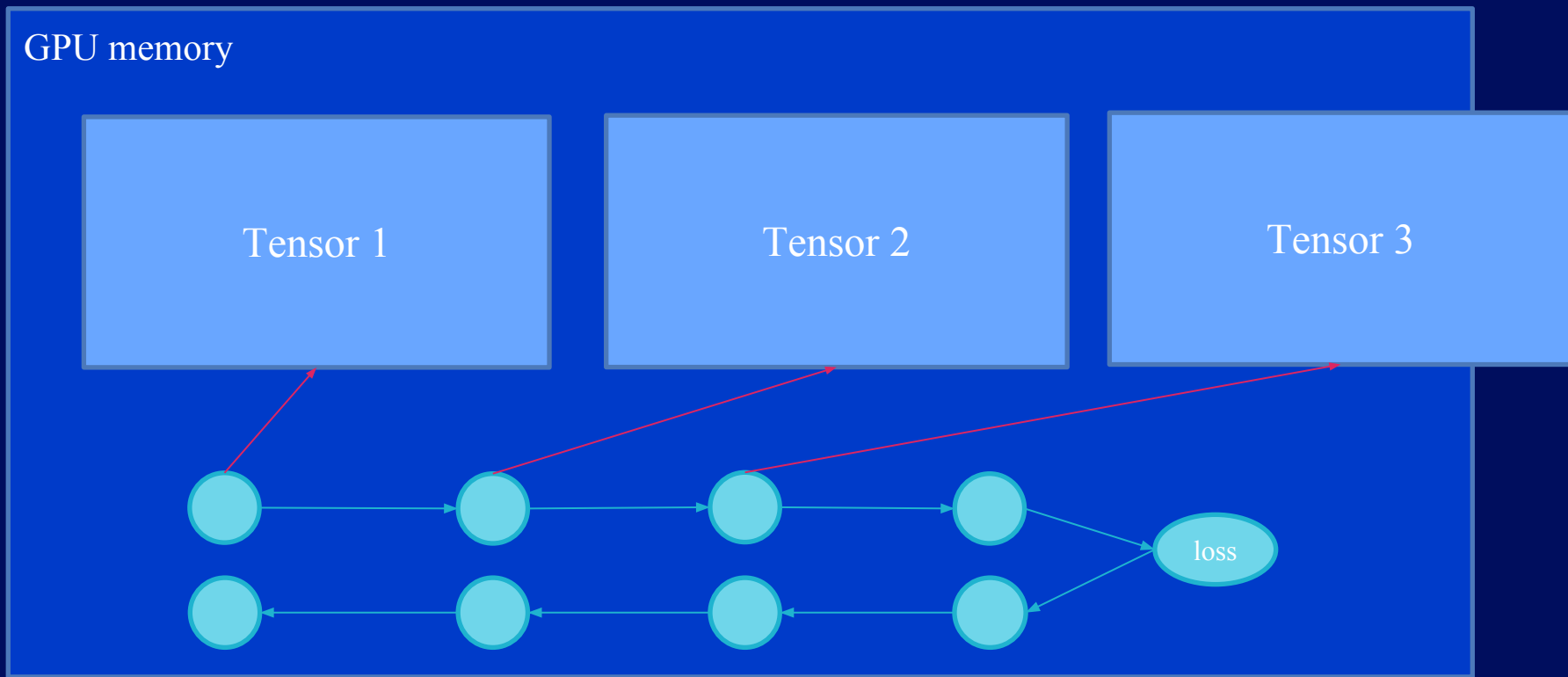
Deep learning is memory constrained

- GPUs have limited memory
- Neural networks are growing deeper and wider
- Amount and size of data to process is always growing

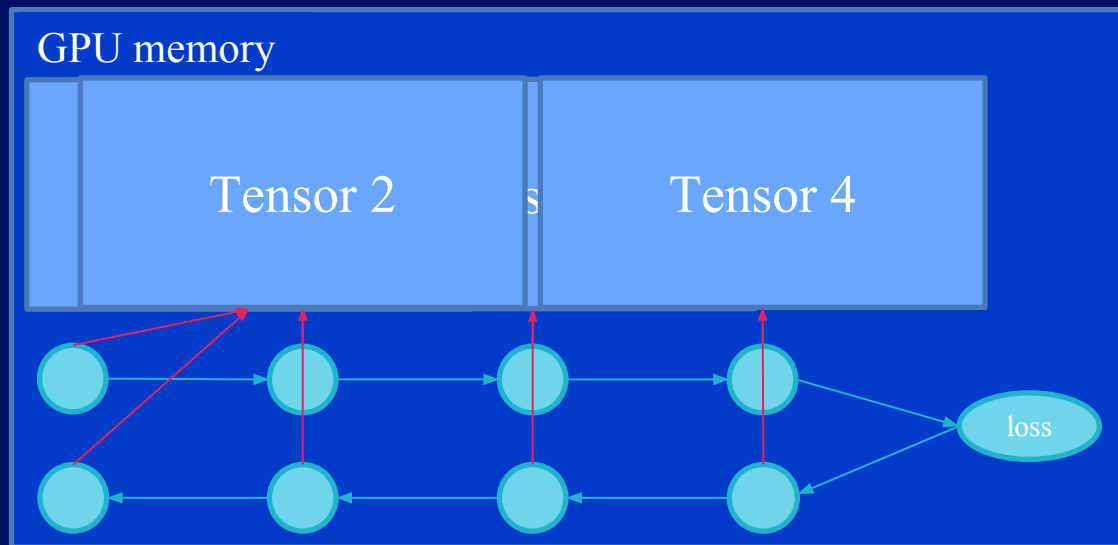
GPU Memory Usage



Model Training in GPU Memory

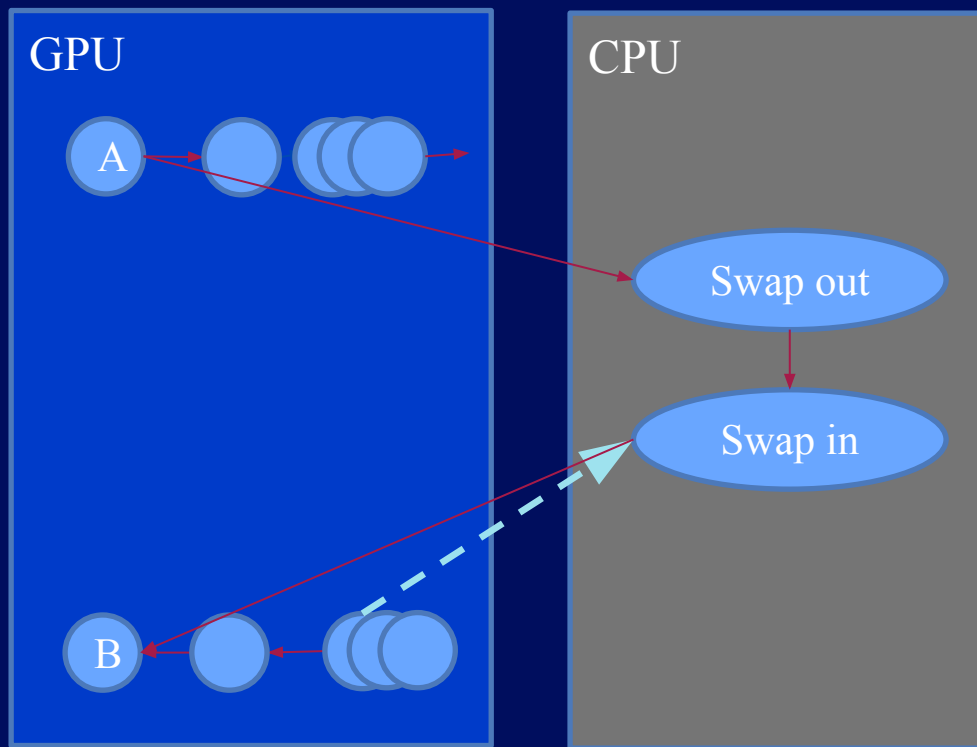


Model Training with Tensor Swapping



System memory

TensorFlow Large Model Support Graph Modifications



Enabling TensorFlow Large Model Support

Keras API

```
from tensorflow_large_model_support import LMS
lms = LMS()
lms.batch_size = 1
# ...
model.fit_generator(generator=training_gen,
                    callbacks=[lms])
```

Estimator API

```
from tensorflow_large_model_support import LMS
lms = LMS()
# ...
mnist_classifier.train(input_fn=train_input_fn, steps=20000,
                      hooks=[logging_hook, lms])
```

What's possible with Large Model Support?

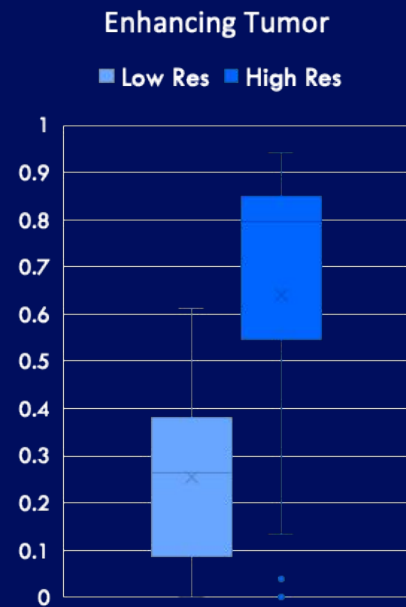
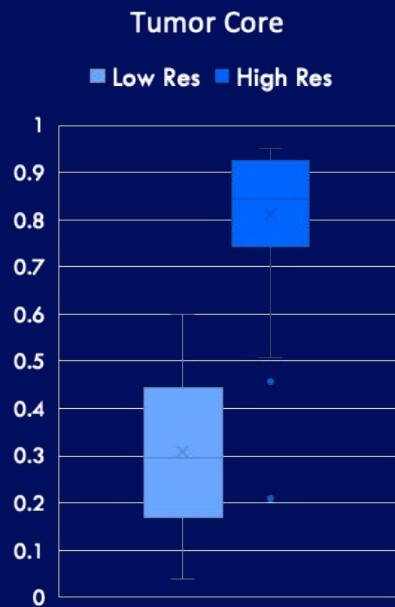
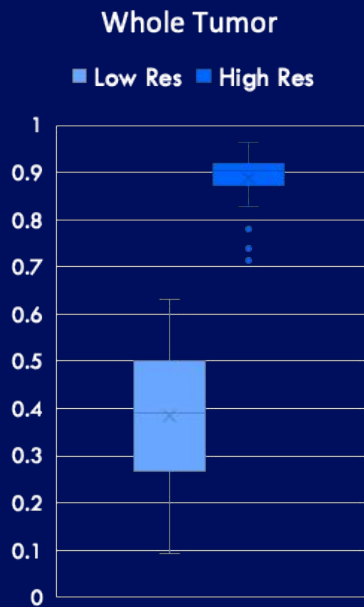
- 10x image resolution - Keras ResNet50
- 10x image resolution - DeepLabV3 2D image segmentation
- 5x MRI resolution - 3D U-Net 3D image segmentation

3D U-Net image segmentation

- 3D U-Net generally has high memory usage requirements
- International Multimodal Brain Tumor Segmentation Challenge (BraTS)
- Existing Keras model with TensorFlow backend

Effect of 2x resolution on Dice Coefficients

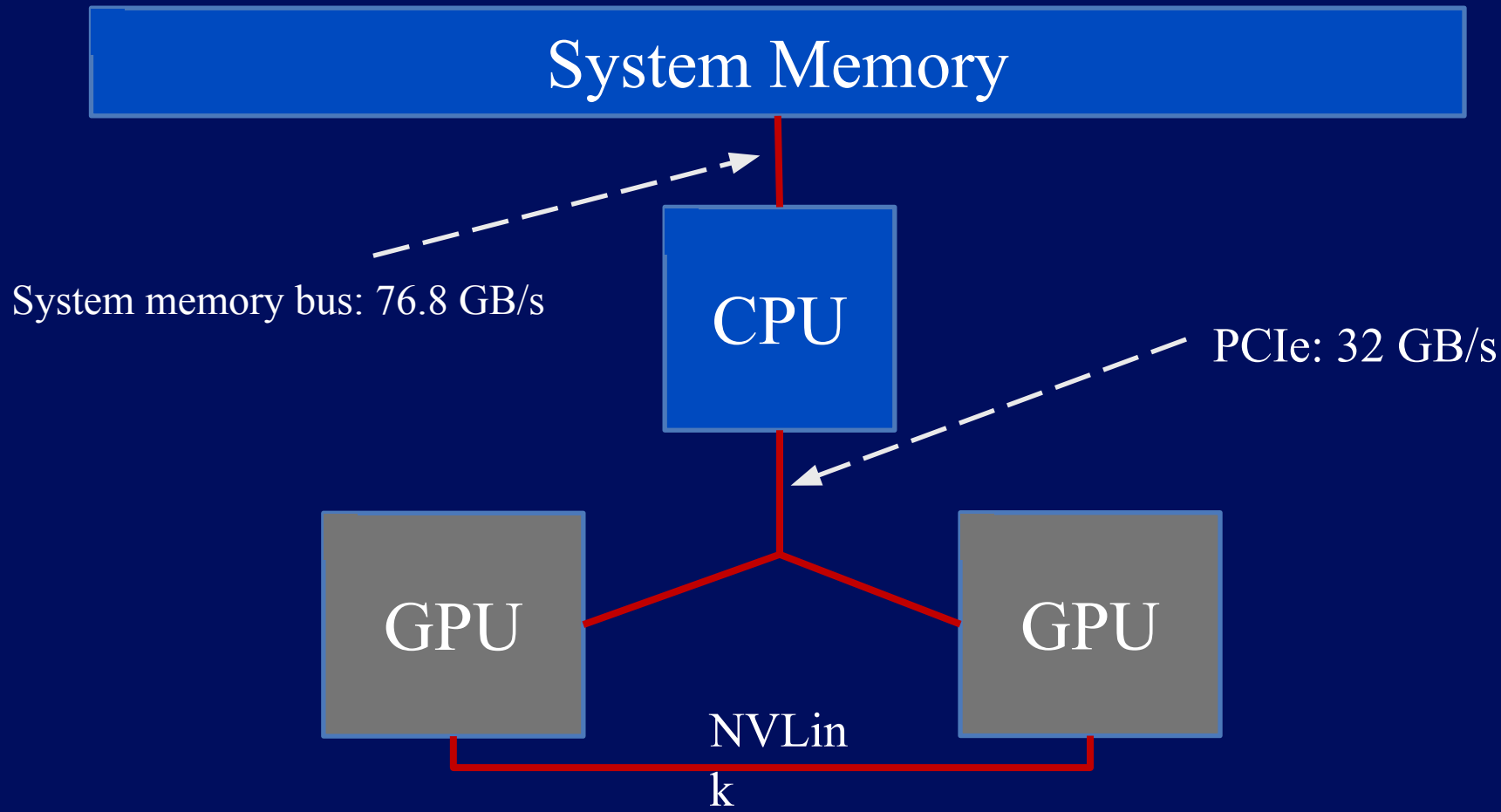
(higher is better)



“Swapping makes everything slow”



Typical GPU connectivity



POWER9 CPU to GPU connectivity

System Memory

CPU

GPU

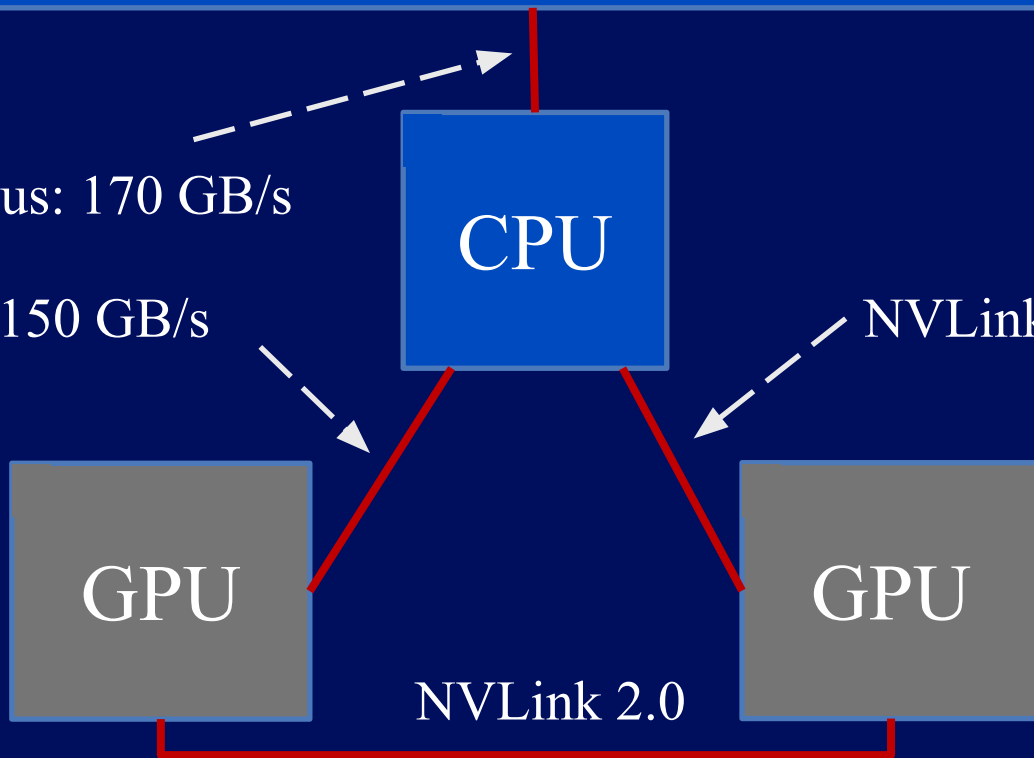
GPU

NVLink 2.0

System memory bus: 170 GB/s

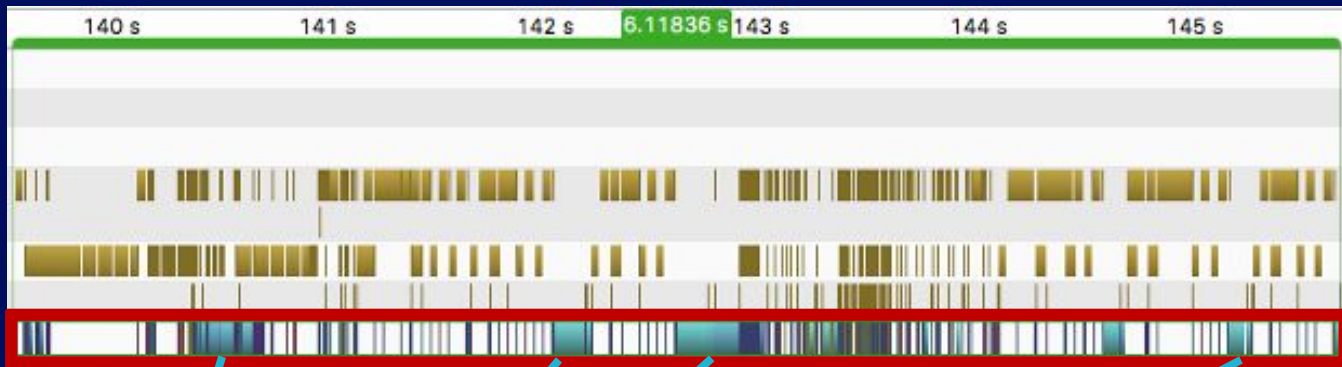
NVLink 2.0: 150 GB/s

NVLink 2.0: 150 GB/s

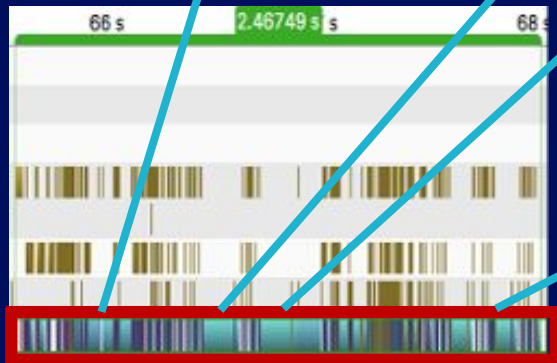


Effects of NVLink 2.0 on Large Model Support

PCIe connected GPU training one high res 3D MRI with large model support

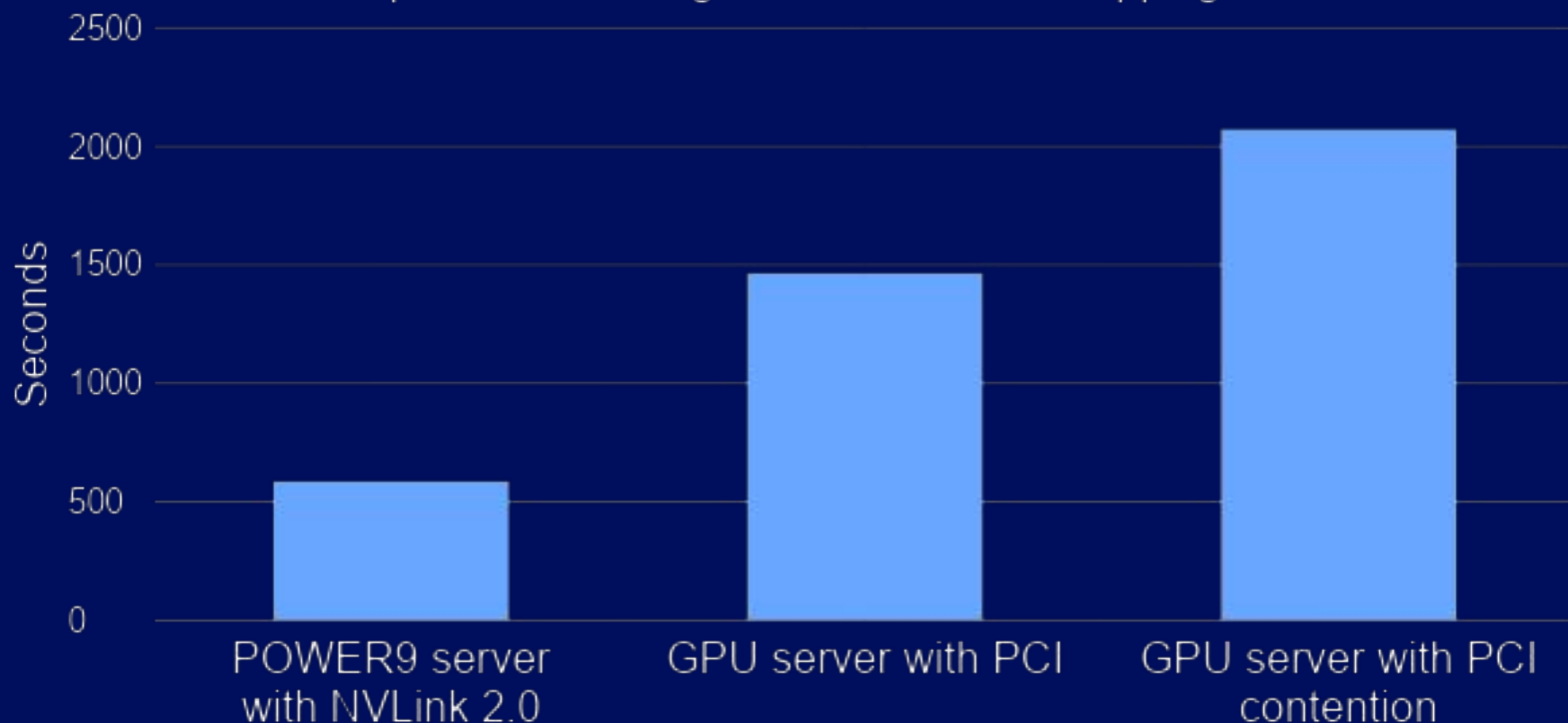


NVLink 2.0 connected GPU training one high res 3D MRI with large model support

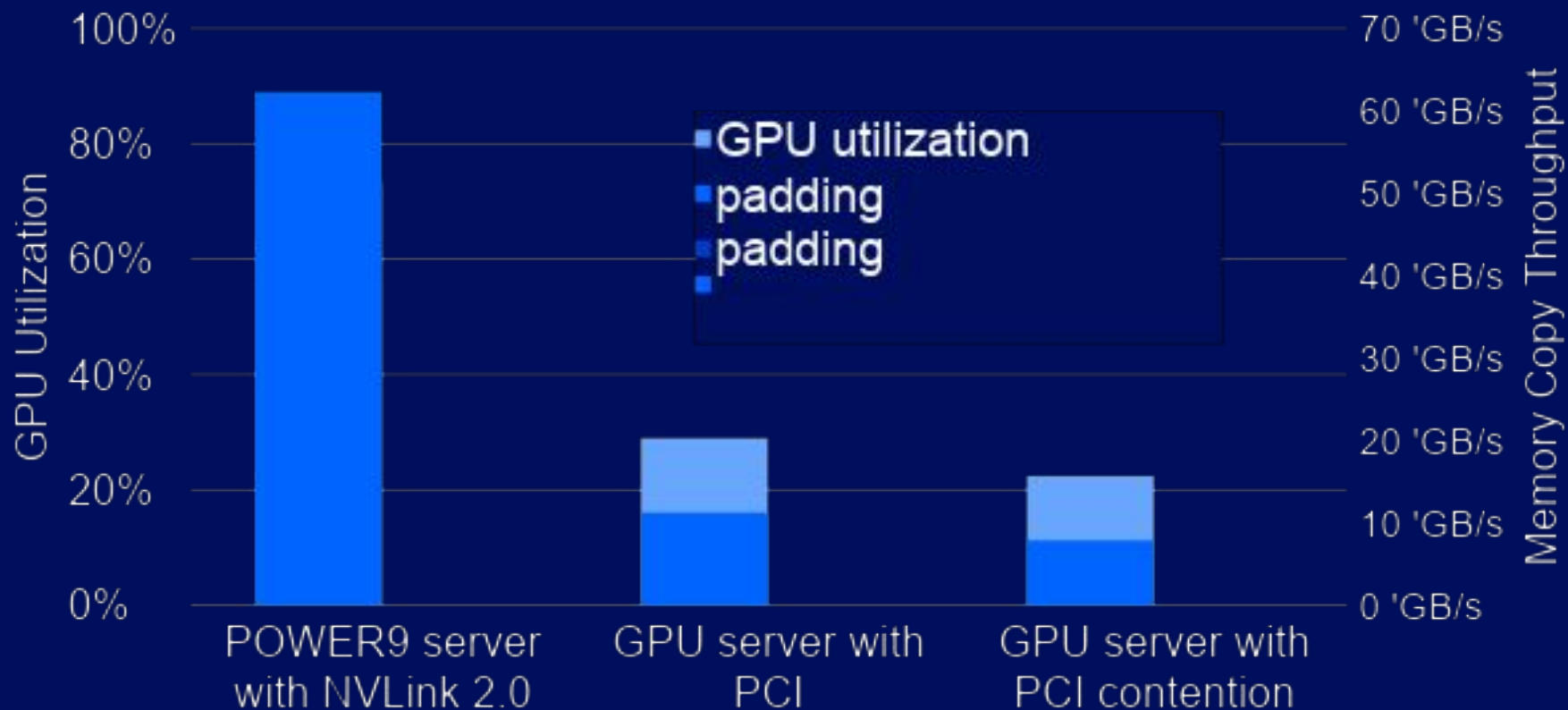


Effects of NVLink 2.0 on epoch times

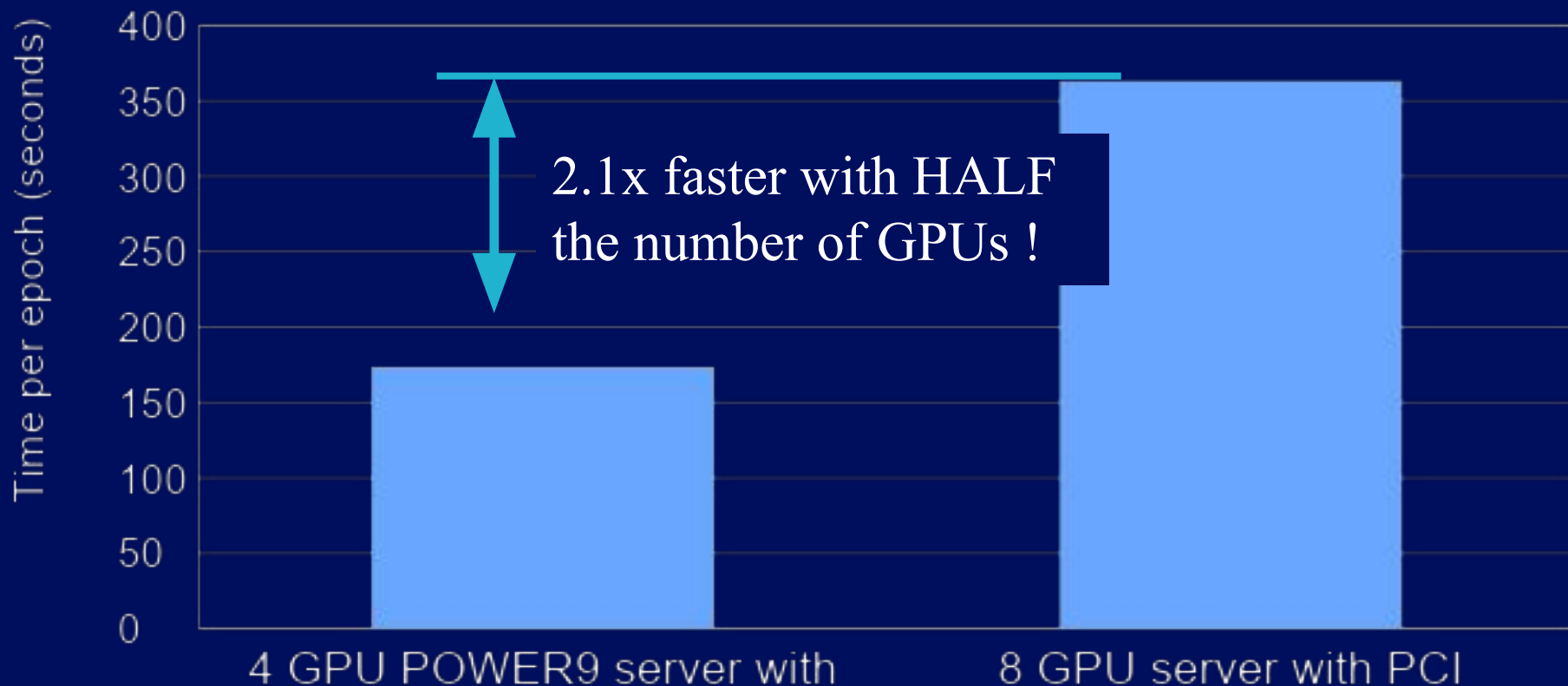
Epoch times at high resolution with swapping



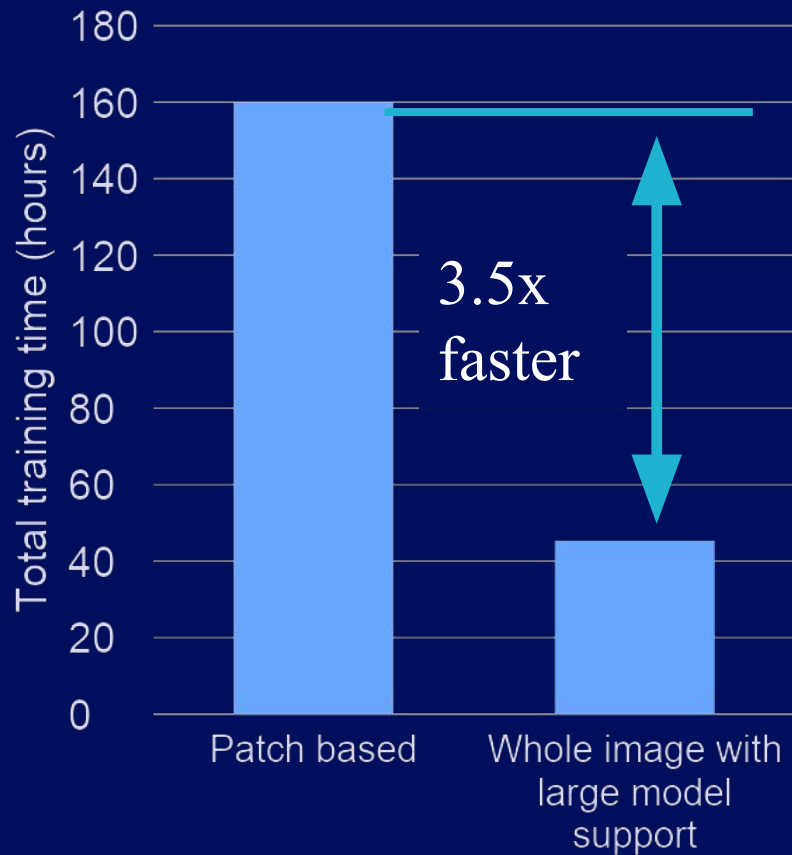
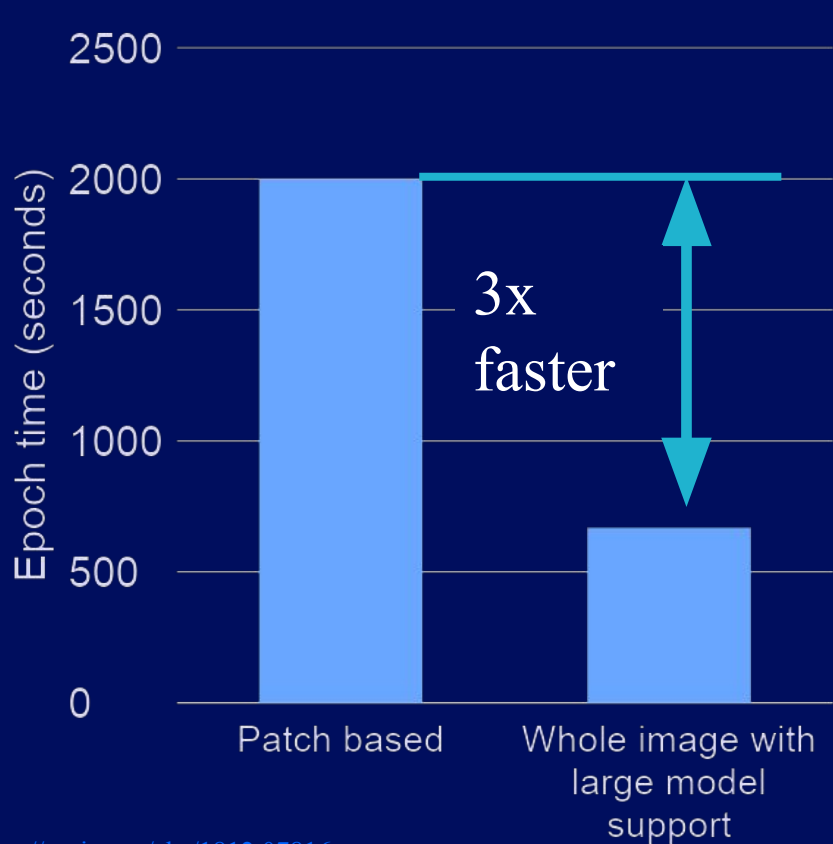
Effects of NVLink 2.0 on GPU Utilization



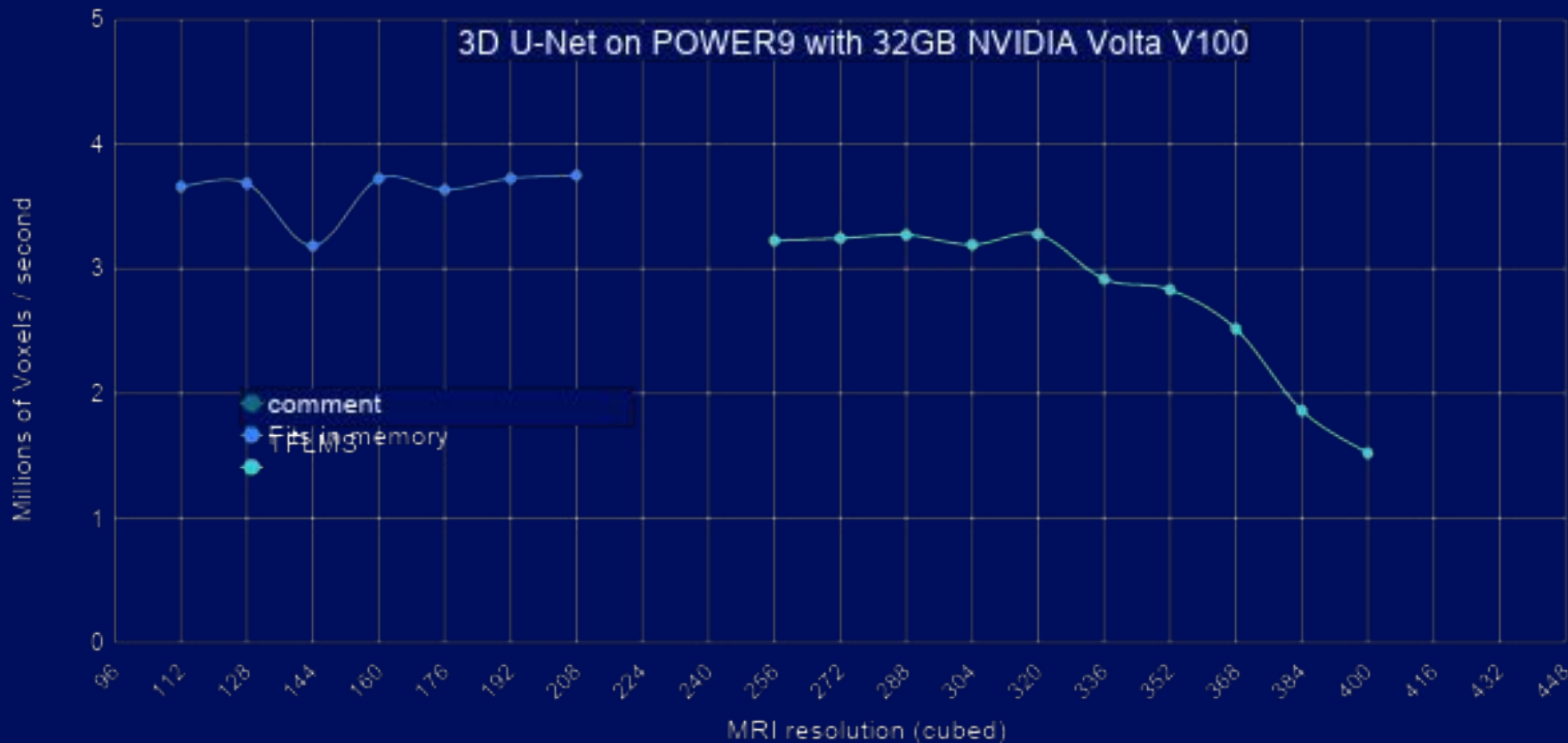
Multi-GPU model training with NVLink 2.0



Patches versus whole image



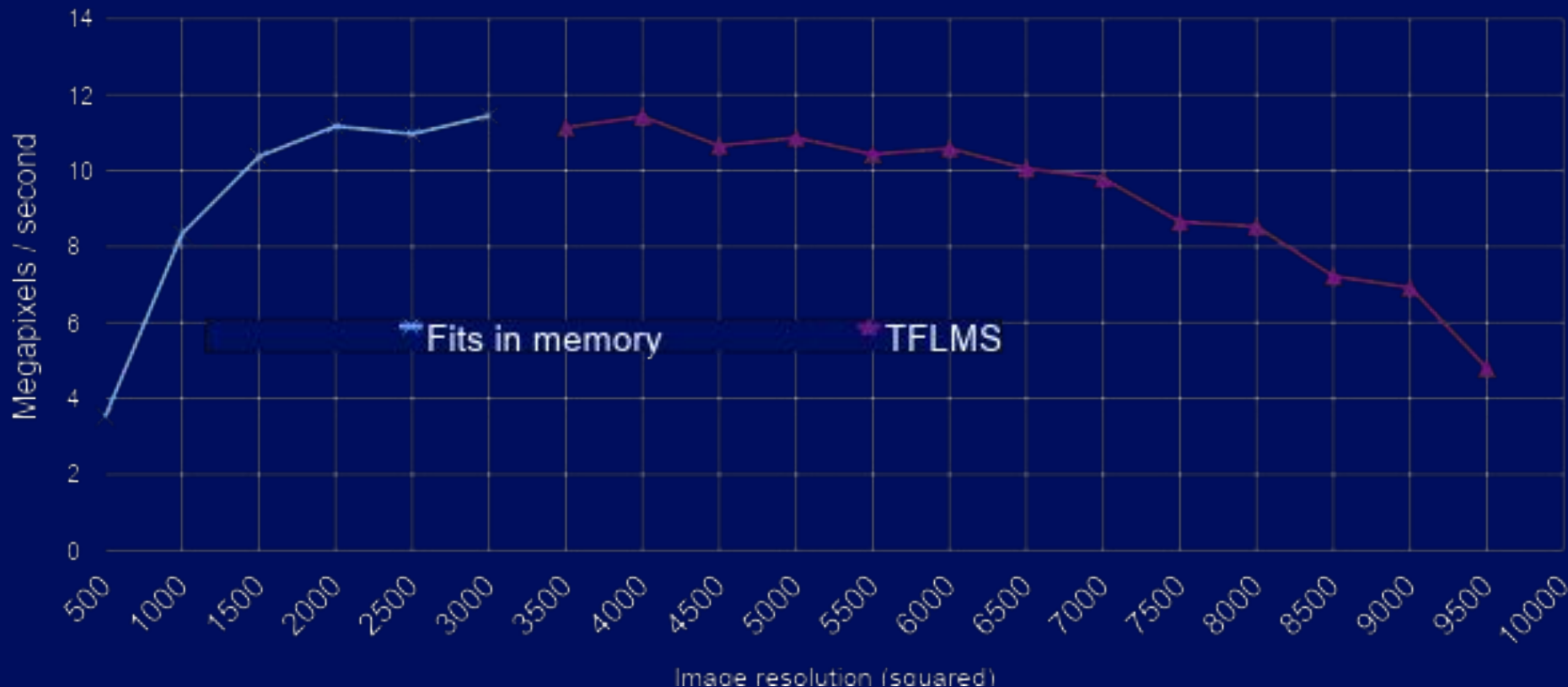
Overhead of Large Model Support with NVLink 2.0



Measured with TFLMS v2.0.0 on TensorFlow 1.13, CUDA 10.1, cuDNN 7.5

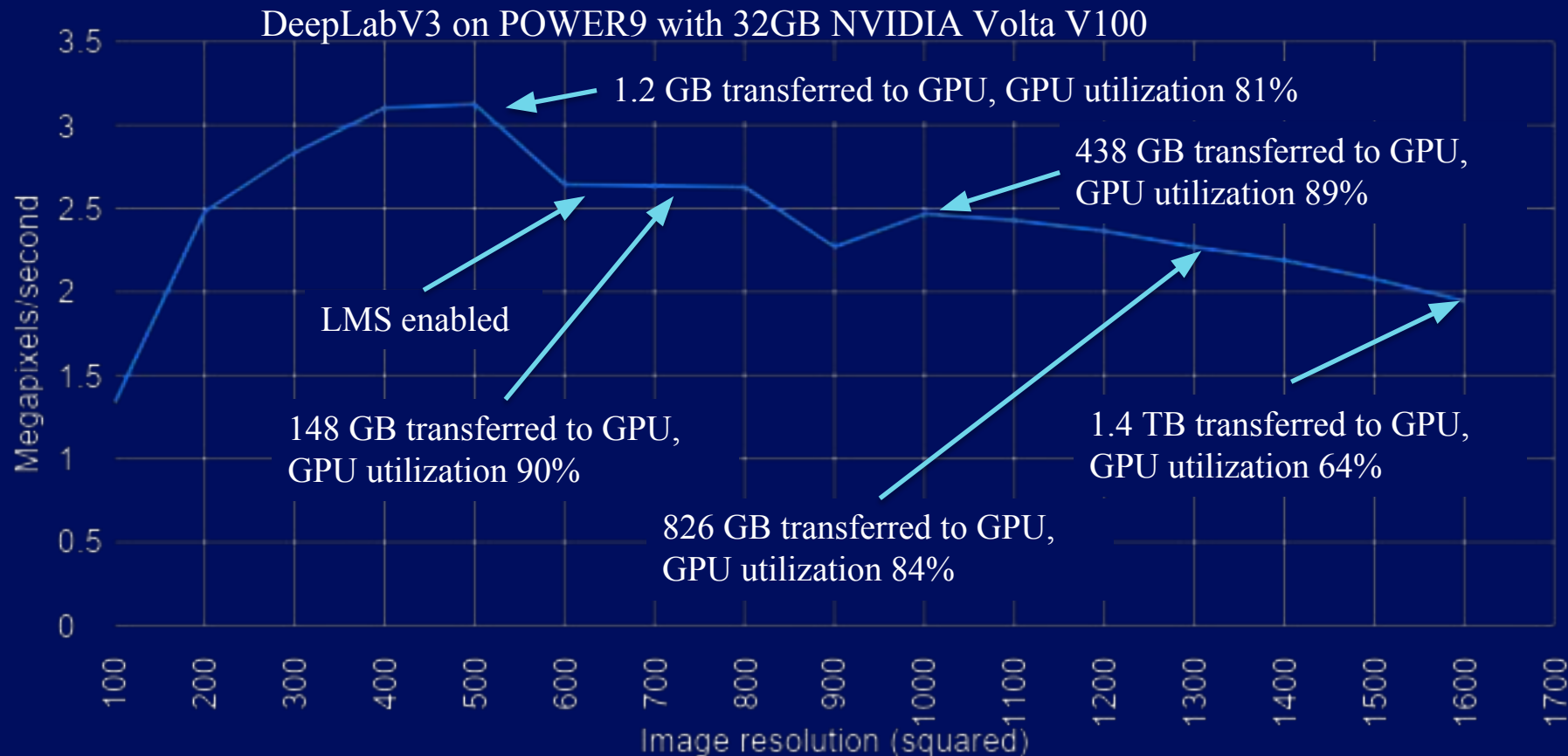
Overhead of Large Model Support with NVLink 2.0

ResNet50 on POWER9 with 32GB NVIDIA Volta V100



Measured with TFLMS v2.0.0 on TensorFlow 1.13, CUDA 10.1, cuDNN 7.5

Overhead of Large Model Support with NVLink 2.0



Large Model Support with NVLink 2.0

- Tensor swapping can be used to overcome GPU memory limits
- Allows training of:
 - deeper models
 - higher resolution data
 - larger batch sizes
- NVLink 2.0 between CPU and GPU allow tensor swapping with minimal overhead

More information

TensorFlow Large Model Support

<https://github.com/IBM/tensorflow-large-model-support>

TFLMS: Large Model Support in TensorFlow by Graph Rewriting

<https://arxiv.org/pdf/1807.02037.pdf>

TensorFlow Large Model Support Case Study

<https://developer.ibm.com/linuxonpower/2018/07/27/tensorflow-large-model-support-case-study-3d-image-segmentation/>

Performance of 3DUnet Multi GPU Model for Medical Image Segmentation using TensorFlow Large Model Support

<http://ibm.biz/3dunet-tflms-multigpu>

Fast and Accurate 3D Medical Image Segmentation with Data-swapping Method

<https://arxiv.org/abs/1812.07816>

Data-parallel distributed training of very large models beyond GPU capacity

<https://arxiv.org/abs/1811.12174>

POWER9 server with NVLink 2.0 connections between CPU and GPU (IBM AC922):

<https://www.ibm.com/us-en/marketplace/power-systems-ac922>