



# PROFILING DEEP LEARNING NETWORK

Jack Han | Solutions Architect | [jahan@nvidia.com](mailto:jahan@nvidia.com)



# AGENDA

Profiling Deep Learning Network

Introduction to Nsight Systems

PyTorch NVTX Annotation

Examples and Benefits

Mixed Precision with BERT

TensorFlow NVTX Plugins

Plugin Usage Example

# HOW TO SPEED-UP NETWORK

- ▶ Use the latest GPU and more GPU?
- ▶ Wait NVIDIA to release the new GPU or newer library?
- ▶ **Optimize Neural Network Computing or Something**
  - ▶ Need to understand your network operation

We are talking about this

# PROFILING NEURAL NETWORK

## Profiling with cProfiler + Snakeviz

 **Soumith Chintala**   
@soumithchintala

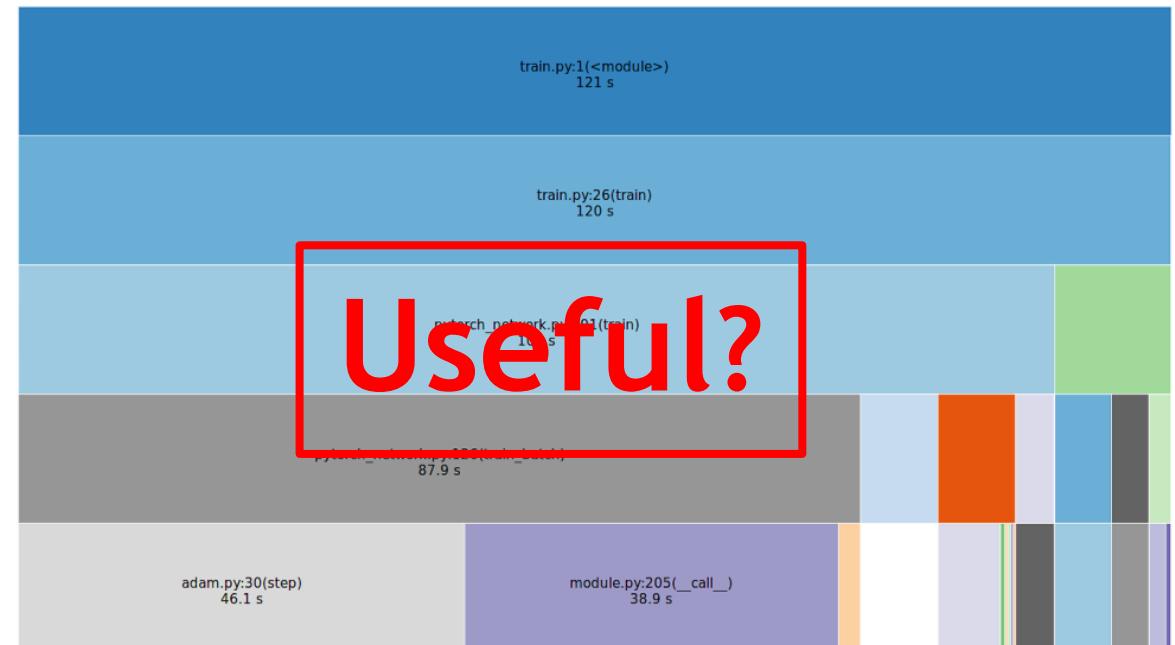
Following ▾

Replying to @sleepinyourhat @PyTorch

I've found this helpful when profiling pytorch code: [jiffyclub.github.io/snakeviz/](https://jiffyclub.github.io/snakeviz/)  
It takes cProfile outputs and gives much nicer viz

7:01 AM - 27 May 2017

5 Retweets 46 Likes 

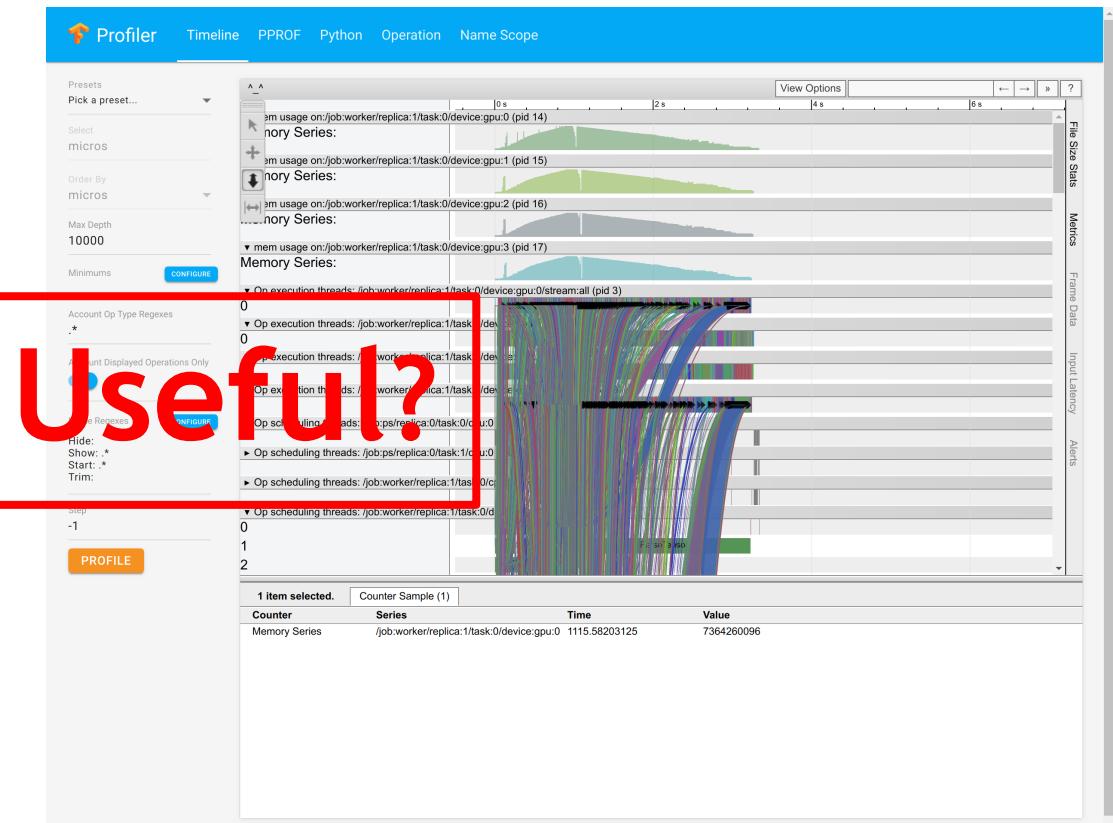


```
python -m cProfile -o 100_percent_gpu_utilization.prof train.py
snakeviz 100_percent_gpu_utilization.prof
```

# TENSORFLOW PROFILER

<https://github.com/tensorflow/profiler-ui>

- ▶ Show timeline and can trace network
- ▶ Still difficult to understand



Useful?

# NVIDIA PROFILING TOOLS

Nsight Systems



Nsight Compute



Nsight Visual Studio Edition

Nsight Eclipse Edition



NVIDIA Profiler (nvprof)

NVIDIA Visual Profiler (nvvprof)



CUPTI (CUDA Profiling Tools Interface)

# NSIGHT PRODUCT FAMILY

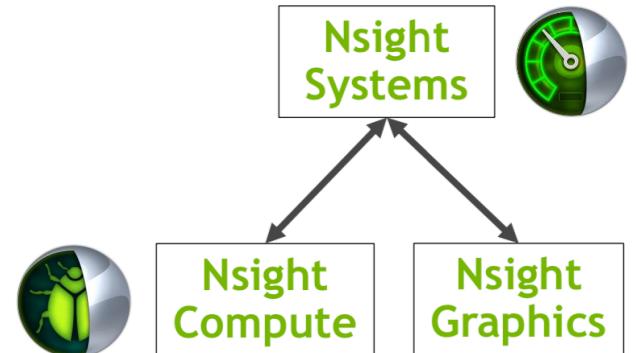
- ▶ Standalone Performance Tools

- ▶ **Nsight Systems** system-wide application algorithm tuning
- ▶ **Nsight Compute** Debug/optimize specific CUDA kernel
- ▶ **Nsight Graphics** Debug/optimize specific graphics

- ▶ IDE plugins

- ▶ **Nsight Eclipse Edition/Visual Studio** editor, debugger, some perf analysis

## Workflow

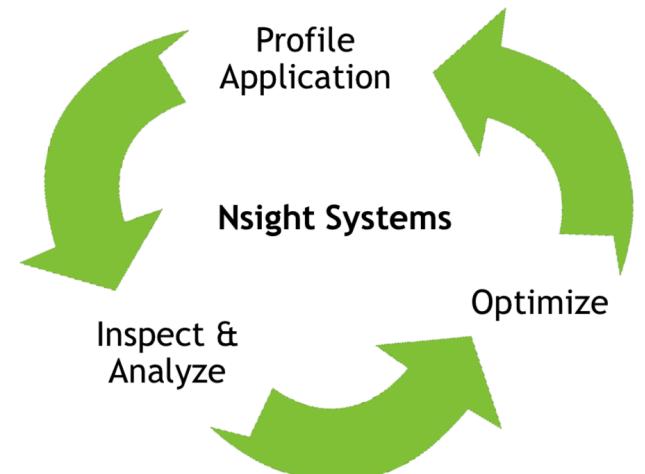


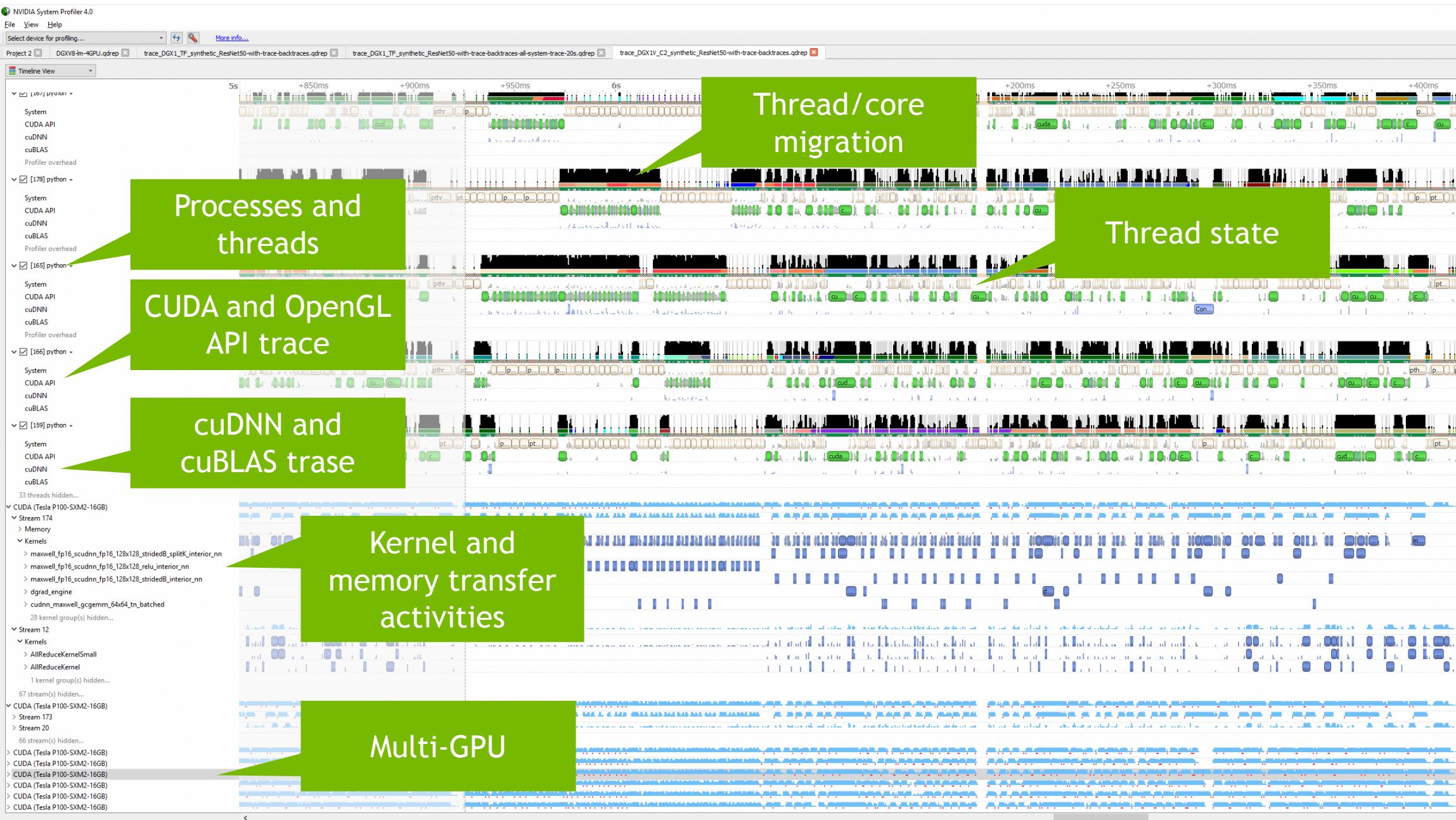


# NSIGHT SYSTEMS

## Overview

- ▶ Profile **System-wide** application
  - ▶ Multi-process tree, GPU workload trace, etc
- ▶ Investigate your workload across multiple CPUs and GPUs
  - ▶ CPU algorithms, utilization, and thread states
  - ▶ GPU streams kernels, memory transfers, etc
  - ▶ NVTX, CUDA & Library API, etc
- ▶ Ready for Big Data
  - ▶ docker, user privilege (linux), cli, etc





CPU (80)

Threads (78)

[1221583] Caffe2F

NVTX

CUDA API

NVTX Tracing

[1221583] Caffe2F

NVTX

CUDA API

Weight... ConvGradient [1...

cudaEventSynchronize

[1221589] Caffe2F

NVTX

CUDA API

Spati...

cud... cudaEv...

[1221587] Caffe2F

[1221582] Caffe2F

73 threads hidden...

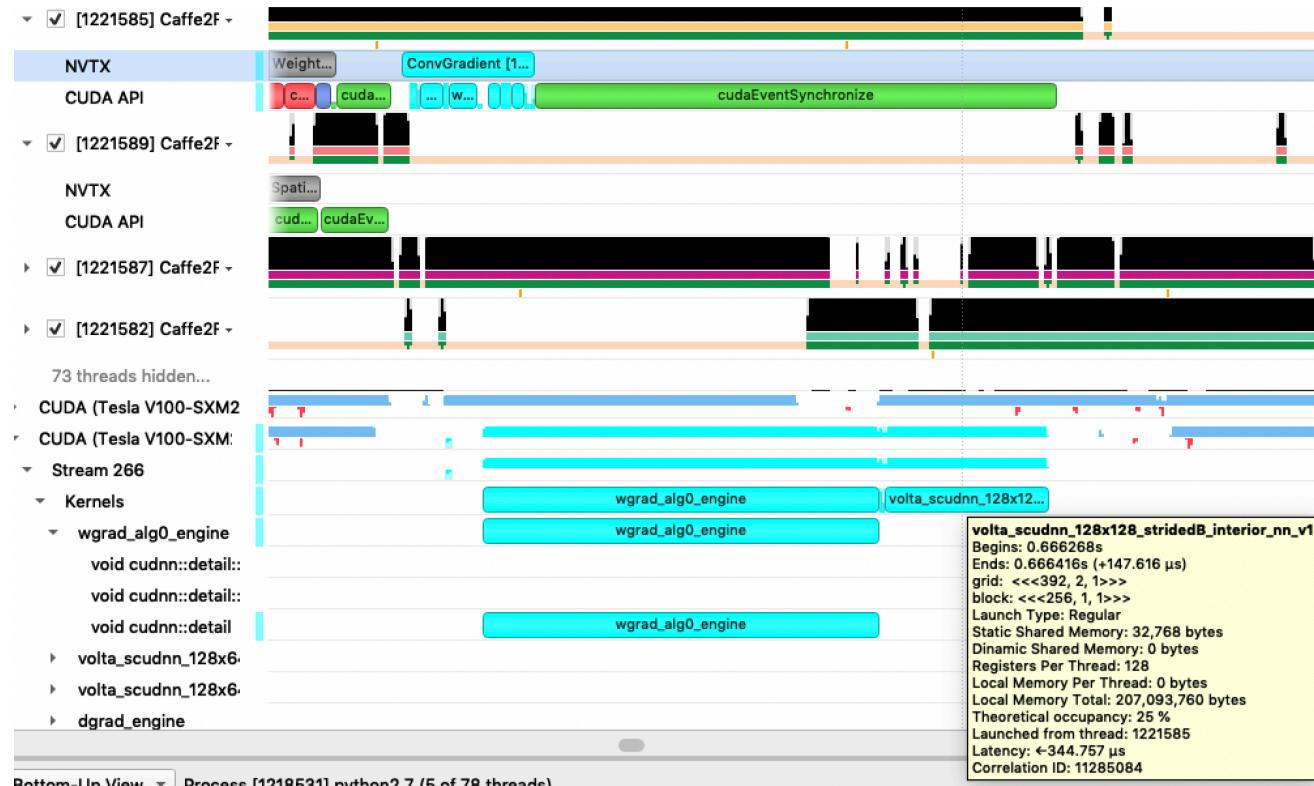
CUDA (Tesla V100-SXM2

CUDA (Tesla V100-S)

CUDA Kernel running  
Time: 0.666129s

# TRANSITIONING TO PROFILE A KERNEL

## Dive into kernel analysis



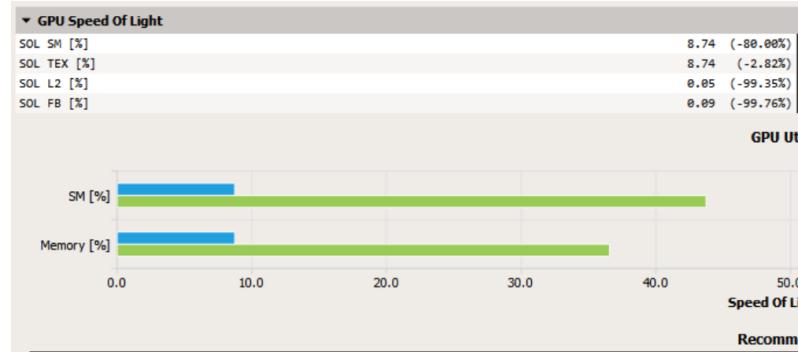


# NVIDIA NSIGHT COMPUTE

## Next Generation Kernel Profiler

- ▶ Interactive CUDA API debugging and kernel profiling
- ▶ Fast Data Collection
- ▶ Graphical and multiple kernel comparison reports
- ▶ Improved Workflow and Fully Customizable (Baselining, Programmable UI/Rules)
- ▶ Command Line, Standalone, IDE Integration
- ▶ Platform Support
  - ▶ OS: Linux(x86,ARM), Windows, OSX (host only)
  - ▶ GPUs: Pascal, Volta, Turing

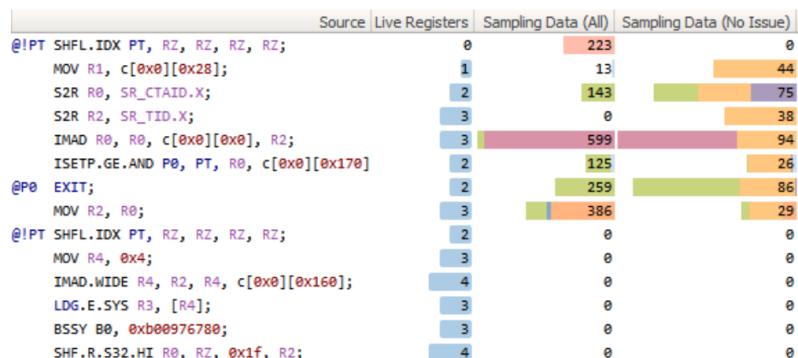
### Kernel Profile Comparisons with Baseline



### Metric Data

inst_executed [inst]	16,528.00	16,528.00	-	13,476.00	13,476.00	-
l1tex_sol_pct [%]			14.33			n/a
launch_block_size			128.00			128.00
launch_function_pcs			47,611,587,968.00			12,273,728.00
launch_grid_size			4,132.00			3,369.00
launch_occupancy_limit_blocks [block]			32.00			32.00
launch_occupancy_limit_registers [register]			21.00			21.00
launch_occupancy_limit_shared_mem [bytes]			384.00			384.00
launch_occupancy_limit_warp [warps]			16.00			16.00
launch_occupancy_per_block_size			3,638.00			3,638.00
launch_occupancy_per_register_count			5,792.00			5,792.00
launch_occupancy_per_shared_mem_size			2,260.00			2,260.00
launch_registers_per_thread [register/thread]			17.00			17.00
launch_shared_mem_config_size [bytes]			49,152.00			49,152.00
launch_shared_mem_per_block_dynamic [bytes/block]			0.00			0.00
launch_shared_mem_per_block_static [bytes/block]			20.00			20.00
launch_thread_count [thread]			528,896.00			431,232.00
ltc_sol_pct [%]			3.23			42.11
memory_access_size_type [bytes]			6.93			7.18
	2.00	32.00	32.00	32.00	32.00	32.00

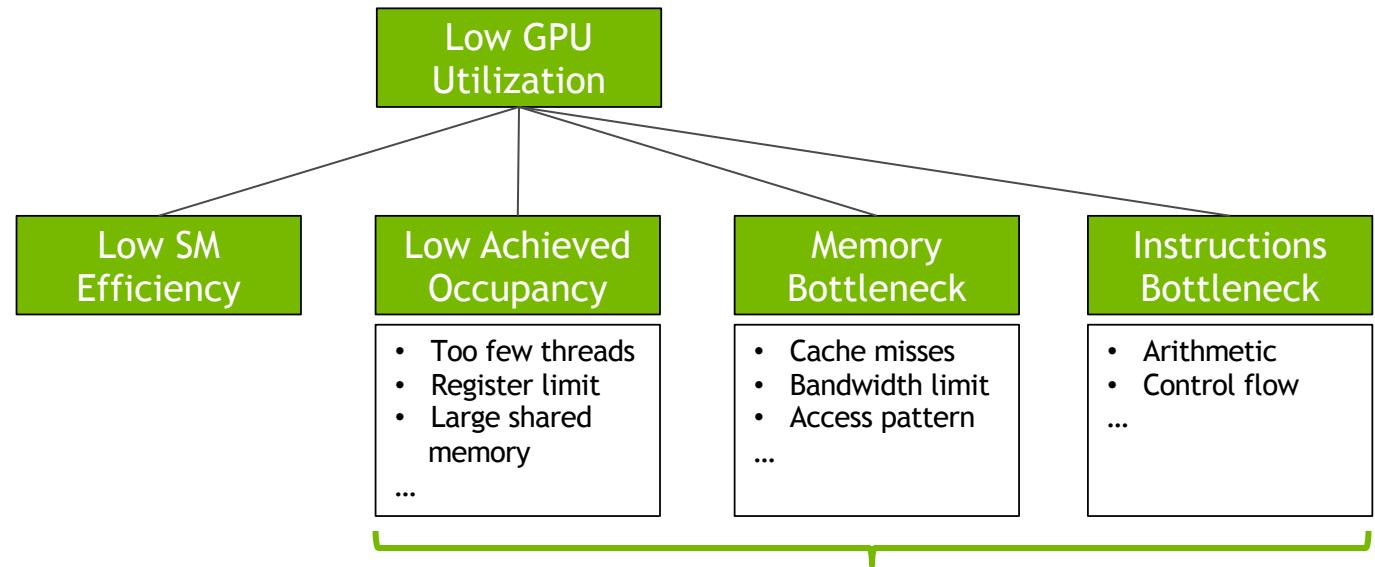
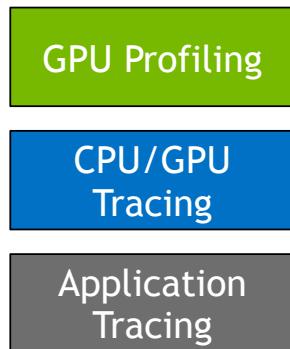
### Source Correlation



# PROFILING GPU APPLICATION

## Focusing GPU Computing

How to measure



NVIDIA (Visual) Profiler / Nsight Compute

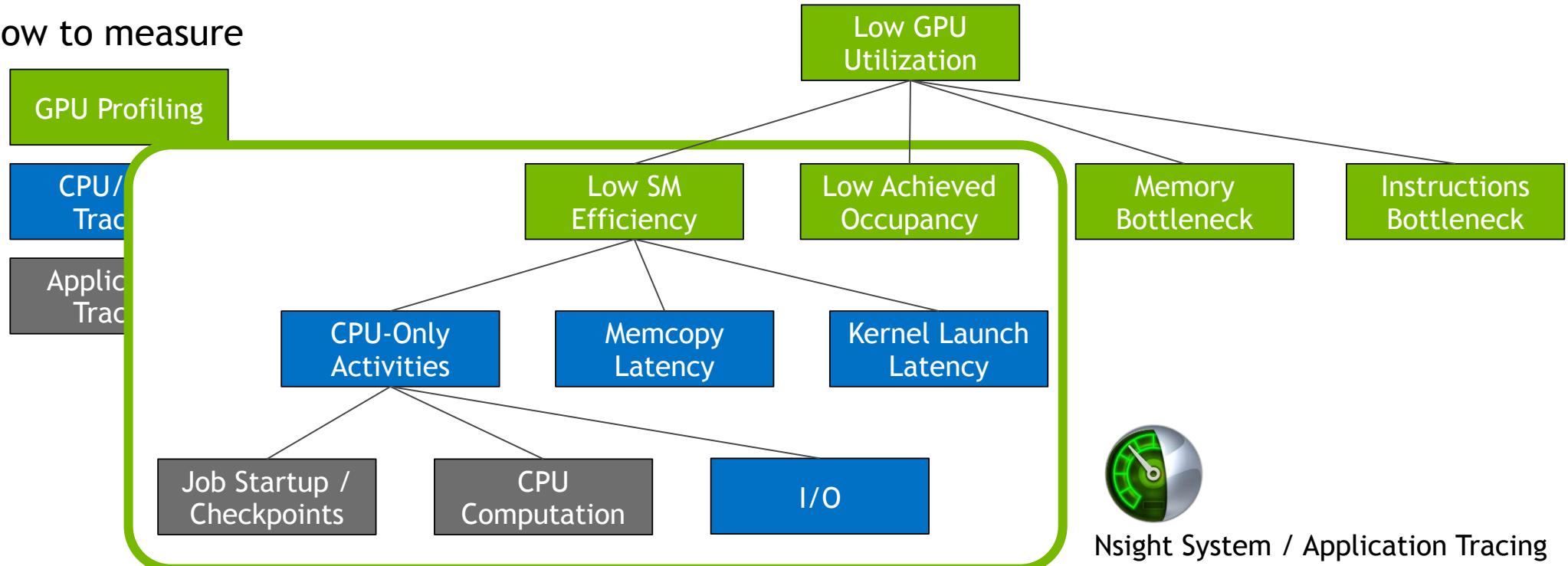


NVIDIA Supports them with cuDNN, cuBLAS, and so on

# PROFILING GPU APPLICATION

## Focusing System Operation

How to measure



# NSIGHT SYSTEMS PROFILE

## Profile with CLI

```
$ nsys profile -t cuda,osrt,nvtx,cudnn,cublas \
    -y 60 \
    -d 20 \
    -o baseline \
    -f true \
    -w true \
    python main.py
```

← APIs to be traced  
← Delayed profile (sec)  
← Profiling duration (sec)  
← Output filename  
← Overwrite when it's true  
← Display  
← Execution command

cuda - GPU kernel

osrt - OS runtime

nvtx - NVIDIA Tools Extension

cudnn - CUDA Deep NN library

cublas - CUDA BLAS library

[https://docs.nvidia.com/nsight-systems/#nsight\\_systems/2019.3.6-x86/06-cli-profiling.htm](https://docs.nvidia.com/nsight-systems/#nsight_systems/2019.3.6-x86/06-cli-profiling.htm)

# DISABLING ASYNCHRONOUS OPERATION

## Elimination of CPU interop effects

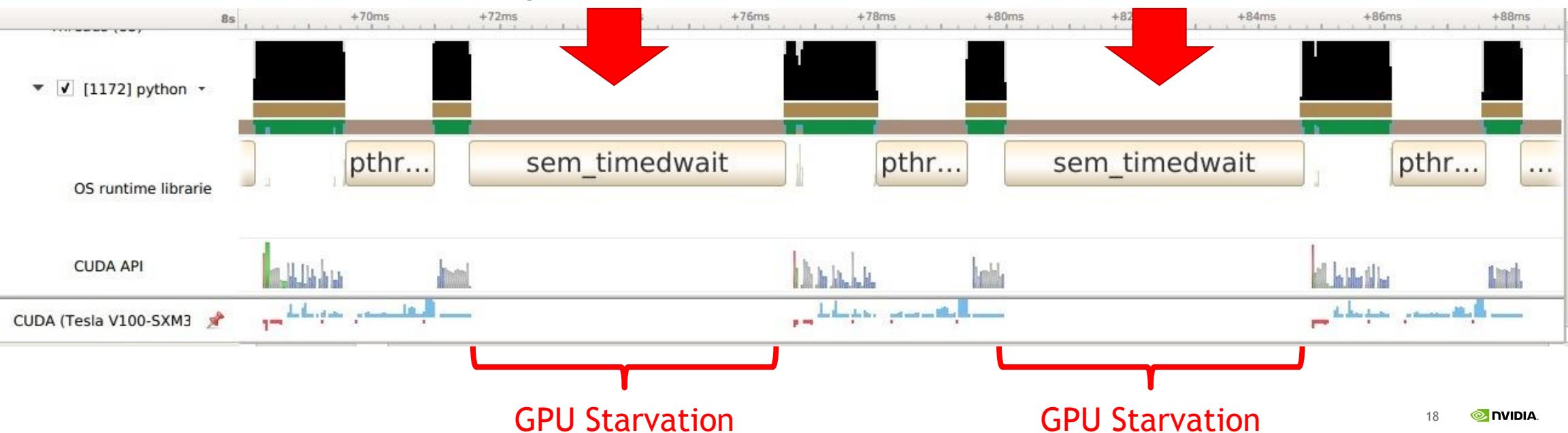
```
export CUDA_LAUNCH_BLOCKING=1
```



# PROFILING PYTORCH MODEL

# BASELINE PROFILE

- ▶ MNIST Training: 89 sec, <5% utilization
- ▶ CPU waits on a semaphore and starves the GPU!



# NVTX ANNOTATIONS

```
def train(args, model, device, train_loader, optimizer, epoch):
    model.train()
    for batch_idx, (data, target) in enumerate(train_loader):

        data, target = data.to(device), target.to(device)

        optimizer.zero_grad()
        output = model(data)
        loss = F.nll_loss(output, target)
```

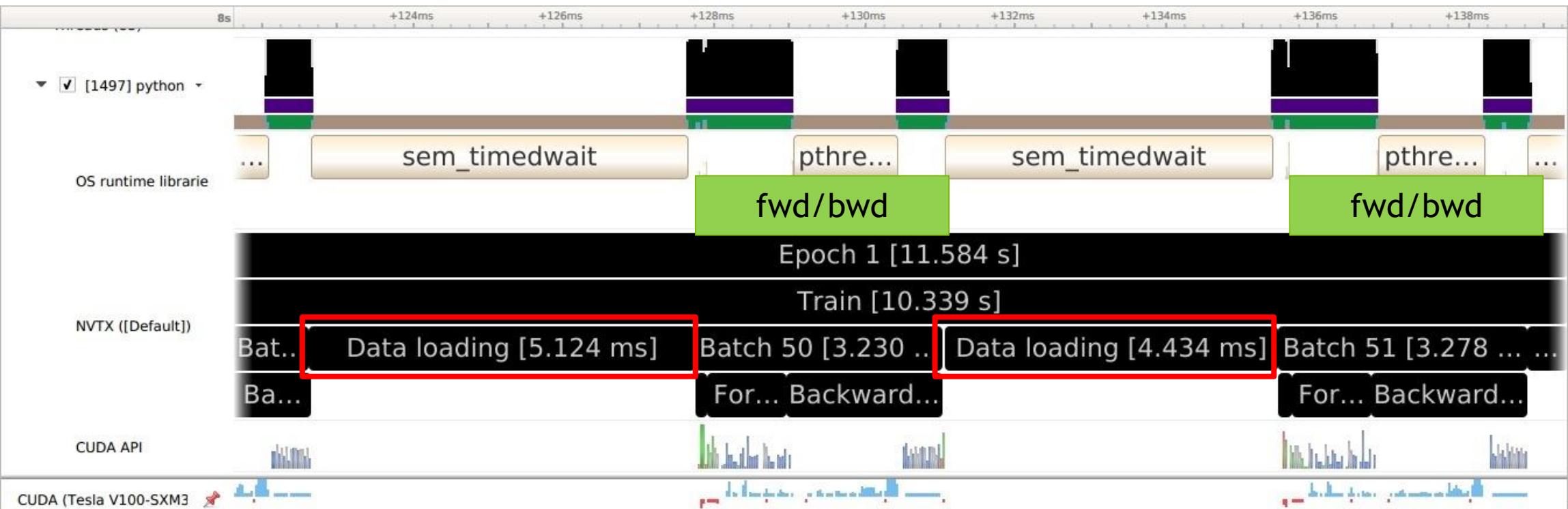
# NVTX ANNOTATIONS

```
import torch.cuda.nvtx as nvtx
def train(args, model, device, train_loader, optimizer, epoch):
    model.train()
    for batch_idx, (data, target) in enumerate(train_loader):
        nvtx.range_push("Batch " + str(batch_idx))

        nvtx.range_push("Copy to device")
        data, target = data.to(device), target.to(device)
        nvtx.range_pop()

        nvtx.range_push("Forward pass")
        optimizer.zero_grad()
        output = model(data)
        loss = F.nll_loss(output, target)
        nvtx.range_pop()
        nvtx.range_pop()
```

# BASELINE PROFILE (WITH NVTX)



- ▶ GPU is idle during **data loading**
- ▶ Data is loaded using a single thread. This starves the GPU!

5.1ms

# OPTIMIZE SOURCE CODE

- ▶ Data loader was configured to use 1 worker thread

```
kwargs = {'num_workers': 1, 'pin_memory': True if use_cuda else {}}
```

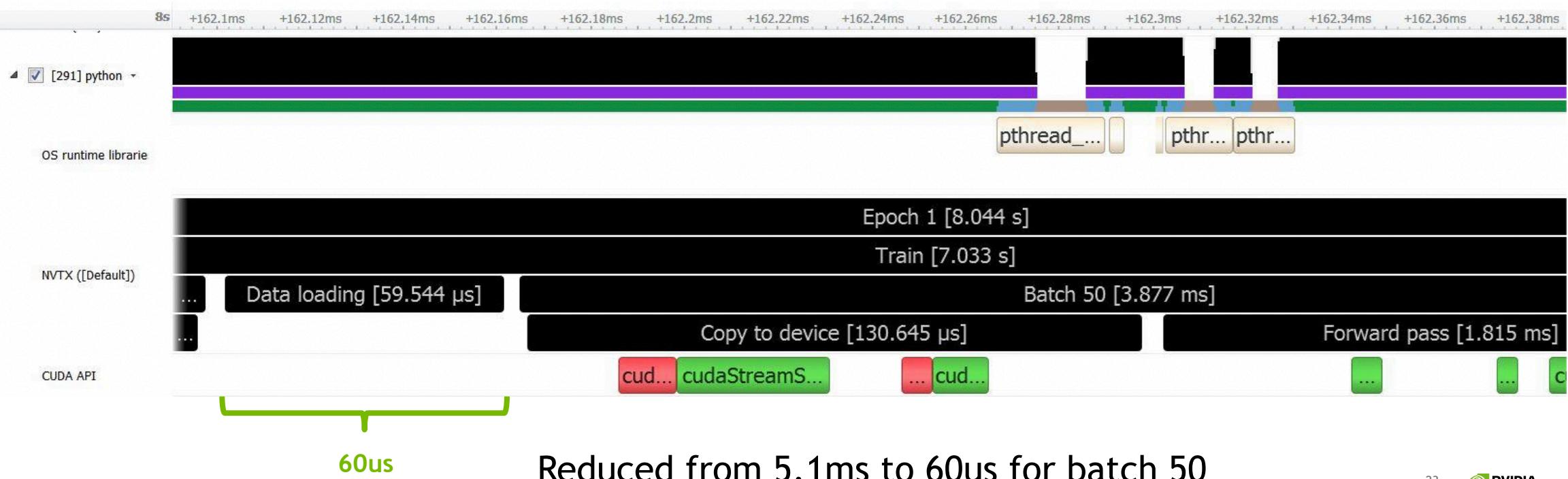


- ▶ Let's switch to using 8 worker threads:

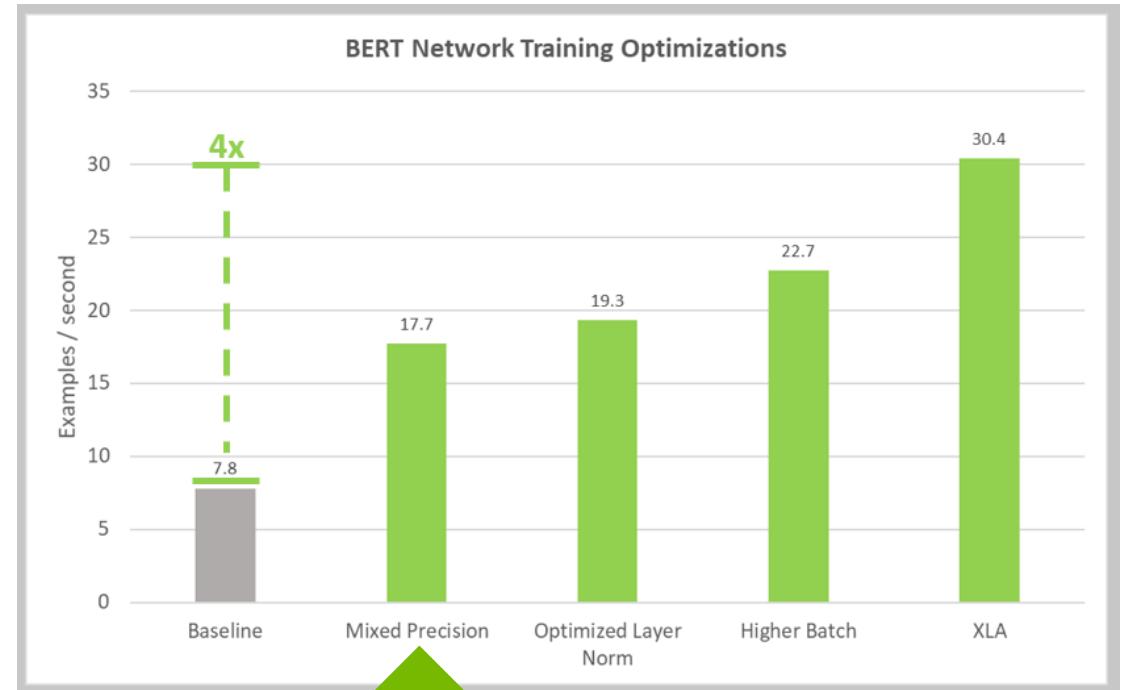
```
kwargs = {'num_workers': 8, 'pin_memory': True if use_cuda else {}}
```

# AFTER OPTIMIZATION

- ▶ Time for data loading reduced for each bath



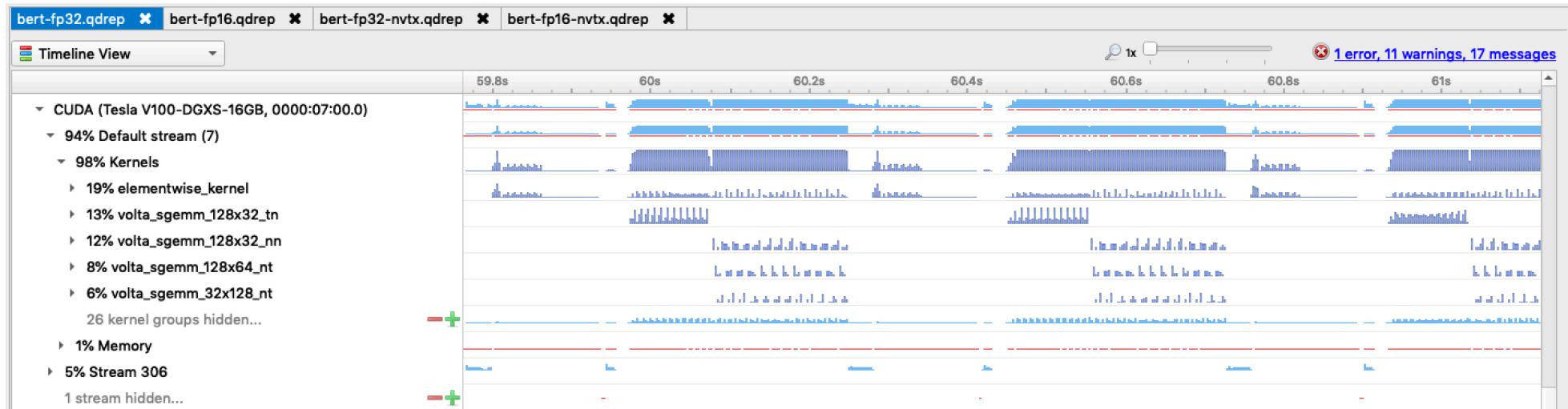
# REAL EXAMPLE



How?

# INITIAL PROFILE

## No NVTX



- ▶ Difficult to understand → no useful

# NVTX TAGGING

```
loss = model(input_ids, segment_ids, input_mask, start_positions, end_positions)          Foward  
...  
  
if args.fp16:  
    optimizer.backward(loss)  
else:  
    loss.backward()  
if (step + 1) % args.gradient_accumulation_steps == 0:  
if args.fp16:  
    # modify learning rate with special warm up BERT uses  
    # if args.fp16 is False, BertAdam is used and handles this automatically  
    lr_this_step = args.learning_rate * warmup_linear.get_lr(global_step, args.warmup_proportion)  
    for param_group in optimizer.param_groups:  
        param_group['lr'] = lr_this_step  
optimizer.step()  
optimizer.zero_grad()  
global_step += 1
```

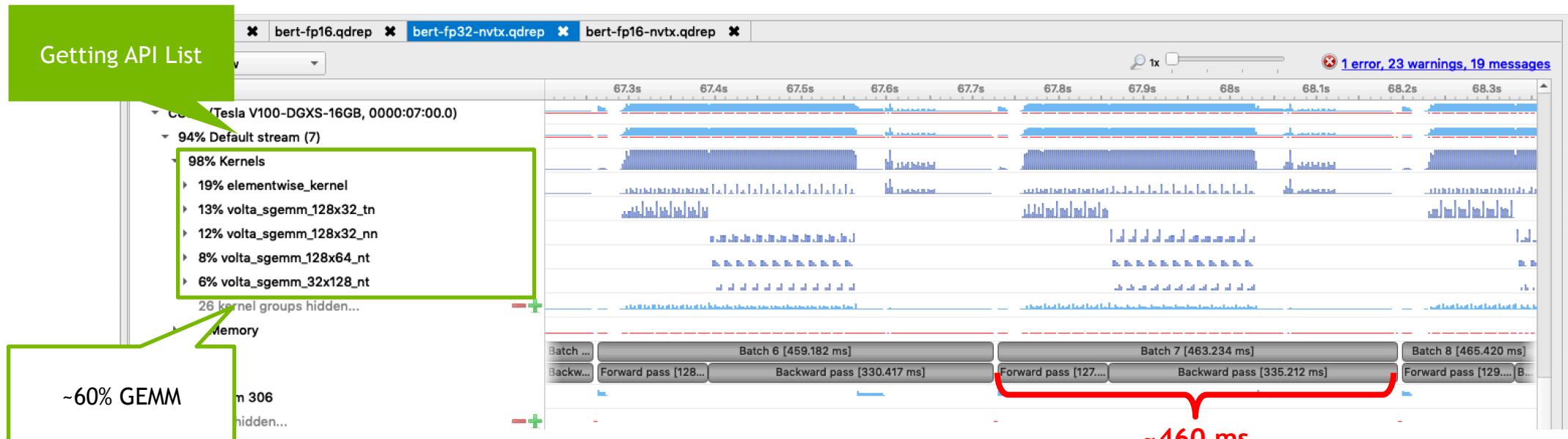
# NVTX TAGGING

```
nvtx.range_push("Batch " + str(step))
nvtx.range_push("Forward pass")
loss = model(input_ids, segment_ids, input_mask, start_positions, end_positions)      Foward
nvtx.range_pop()

...
nvtx.range_push("Backward pass")                                                       Backward
if args.fp16:
    optimizer.backward(loss)
else:
    loss.backward()
if (step + 1) % args.gradient_accumulation_steps == 0:
if args.fp16:
    # modify learning rate with special warm up BERT uses
    # if args.fp16 is False, BertAdam is used and handles this automatically
    lr_this_step = args.learning_rate * warmup_linear.get_lr(global_step, args.warmup_proportion)
    for param_group in optimizer.param_groups:
        param_group['lr'] = lr_this_step
optimizer.step()
optimizer.zero_grad()
global_step += 1
nvtx.range_pop()
nvtx.range_pop()
```

# BERT FP32 BENCHMARK

## HuggingFace's pretrained BERT



4 V100 GPUs w/ NVLINK, Batch size: 32, max\_seq\_length: 512

# BERT FP16 BENCHMARK

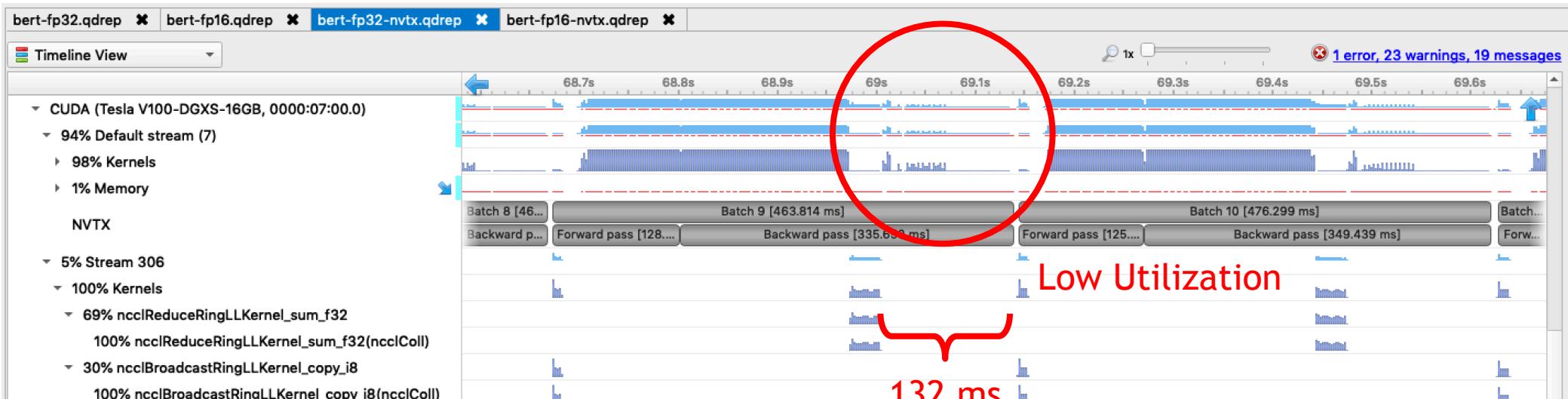
## HuggingFace's pretrained BERT



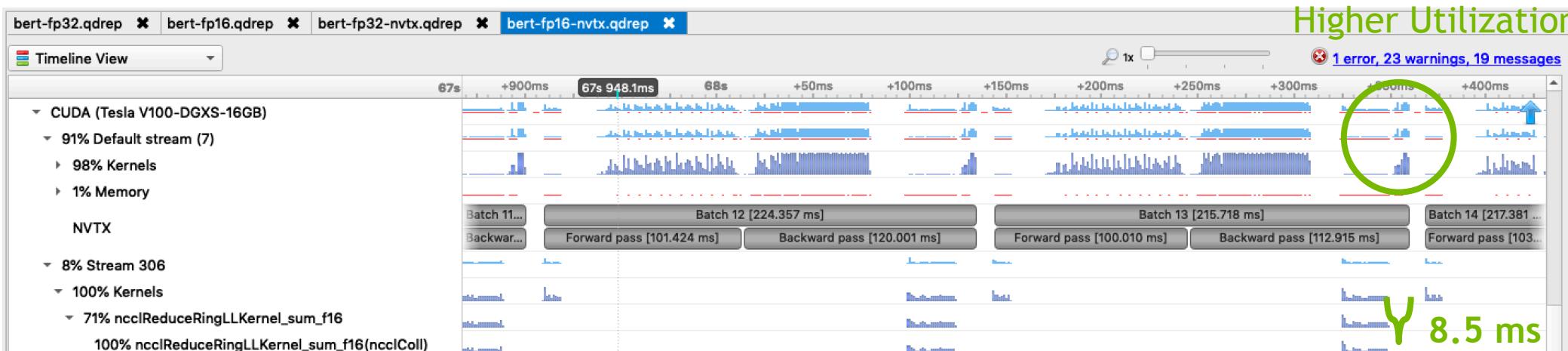
4 V100 GPUs w/ NVLINK, Batch size: 32, max\_seq\_length: 512

# APEX OPTIMIZER OPTIMIZATION

FP32



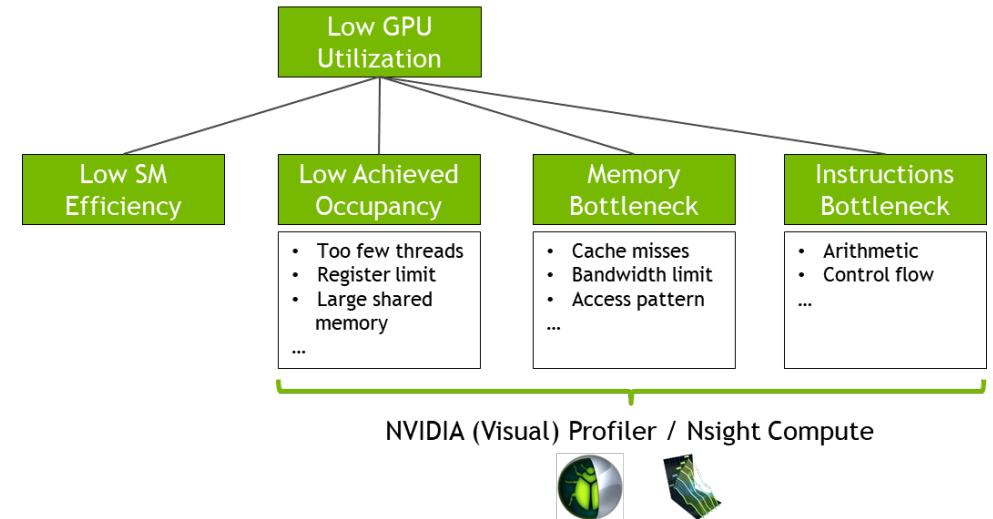
FP16

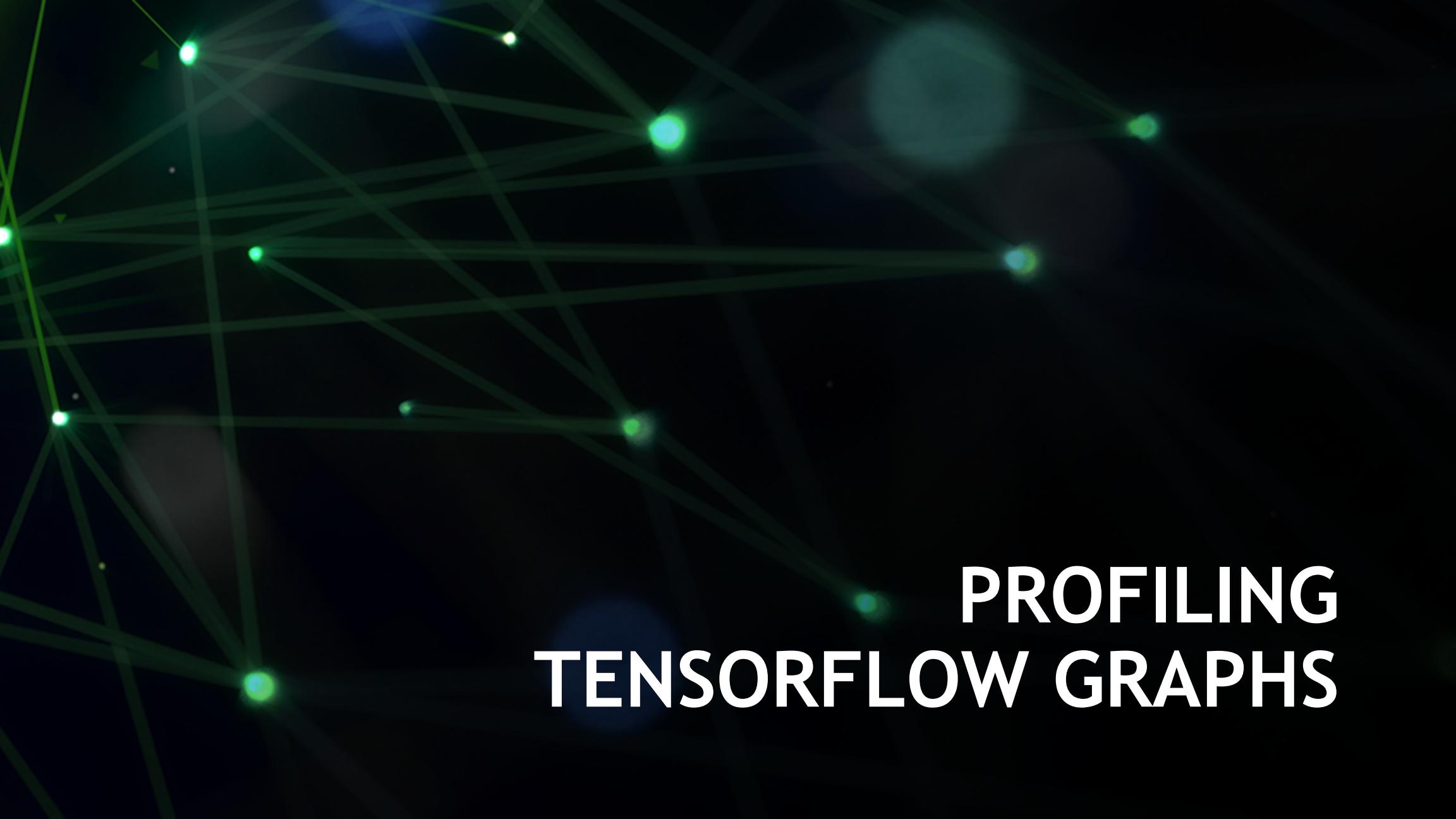


# PROFILING GPU APPLICATION

## Reminding

- ▶ You can investigate the performance bottleneck in each layer
- ▶ NVIDIA provides such inefficient layers / graphs optimization support
  - ▶ For example,
    - ▶ Like this, adam\_cuda\_kernel in APEX
    - ▶ cudnnMultiHeadAttnForward() in cuDNN
    - ▶ BERT speed-up





# PROFILING TENSORFLOW GRAPHS

# NVTX WITH TENSORFLOW

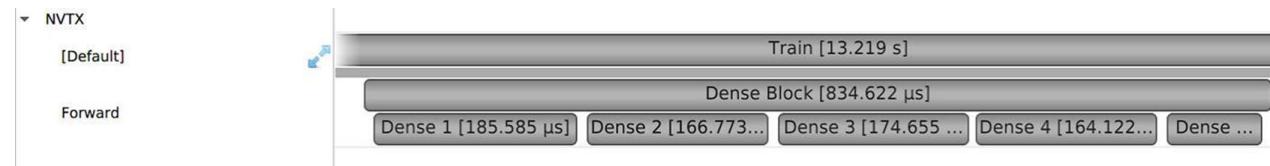
- ▶ The TF\_DISABLE\_NVTX\_RANGES variable enables and disables NVTX ranges in TensorFlow
  - ▶ NVTX collect is enabled by default
  - ▶ Available in NGC Container
- ▶ Having NVTX plugin for TF
  - ▶ Import nvtx plugin
  - ▶ Annotation in python code
  - ▶ Get data from through Profiler

```
1  import tensorflow as tf
2
3  import nvtx.plugins.tf as nvtx_tf
4
5  # Option 1: use decorators
6  @nvtx_tf.layers.trace(message="Dense Block", domain_name="Forward",
7  grad_domain_name="Gradient")
8 ▼ def dense_block(net):
9      net = tf.layers.dense(net, 1, activation=tf.nn.relu, name='dense_1')
10     net = tf.layers.dense(net, 64, activation=tf.nn.relu, name='dense_2')
11     net = tf.layers.dense(net, 128, activation=tf.nn.relu, name='dense_3')
12     net = tf.layers.dense(net, 256, activation=tf.nn.relu, name='dense_4')
13     net = tf.layers.dense(net, 512, activation=tf.nn.relu, name='dense_5')
14     return net
15
16
17 ▼ def build_model(inputs)
18
19     x = dense_block(inputs)
20
21     # Option 2: wrap parts of your graph with nvtx layers
22     x, nvtx_context = nvtx_tf.layers.start(x, message="logits-softmax")
23     x = tf.layers.dense(x, 512, activation=None, name='logits')
24     x = tf.nn.softmax(x, name='probs')
25     x = nvtx_tf.layers.end(x, nvtx_context)
26
27     return x
28
```

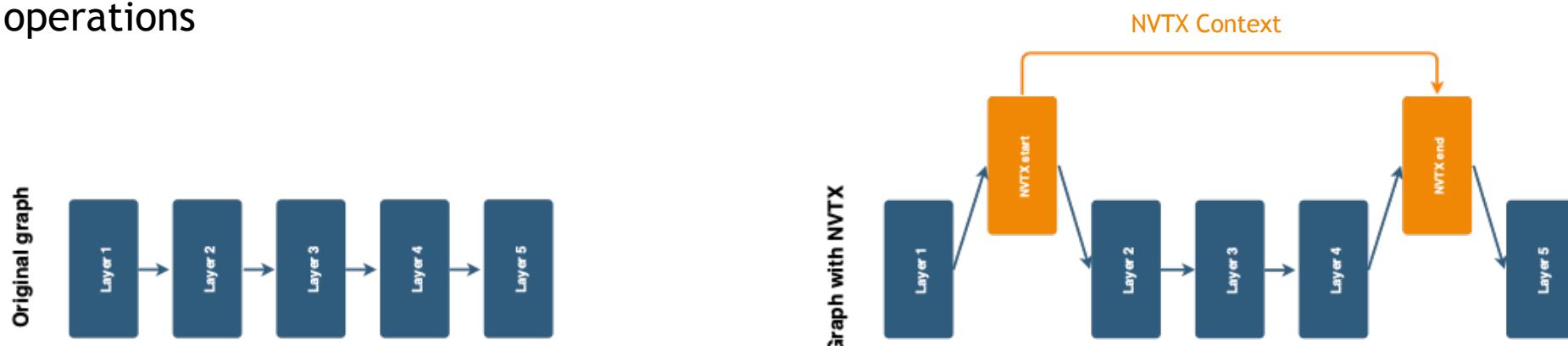
# NVTX PLUGINS FOR DEEP LEARNING

## Profiling TensorFlow for the graphs

- Allows users to add their own NVIDIA Tools Extension (NVTX) events and time ranges to a TensorFlow graph



- Ranges are added by wrapping regions of the computation graph with start and end operations



# SETUP

## Install NVTX plugins

```
pip install nvtx-plugins-tf
```

No need to install this plugin since NGC 19.08

# HOW TO USE

## ▶ Adding Markers

```
import nvtx.plugins.tf as nvtx_tf

x, nvtx_context = nvtx_tf.ops.start(x, message='Dense 1-2', \
                                      domain_name='Forward', grad_domain_name='Gradient', enabled=ENABLE_NVTX, trainable=True)
x = tf.layers.dense(x, 1000, activation=tf.nn.relu, name='dense_1')
x = tf.layers.dense(x, 1000, activation=tf.nn.relu, name='dense_2')
x = nvtx_tf.ops.end(x, nvtx_context)
x = tf.layers.dense(x, 1000, activation=tf.nn.relu, name='dense_3')
```

## ▶ Function Detector

```
import nvtx.plugins.tf as nvtx_tf

@nvtx_tf.ops.trace(message='Dense Block', domain_name='Forward',
                    grad_domain_name='Gradient', enabled=ENABLE_NVTX, trainable=True)
def dense_block(x):
    x = tf.layers.dense(x, 1000, activation=tf.nn.relu, name='dense_1')
    x = tf.layers.dense(x, 1000, activation=tf.nn.relu, name='dense_2')
    return x
```

# WORKING WITH SESSION/TF.ESTIMATOR

```
from nvtx.plugins.tf.estimator import NVTXHook

nvtx_callback = NVTXHook(skip_n_steps=1, name='Train')
training_hooks=[ ]
training_hooks.append(nvtx_callback)

with tf.train.MonitoredSession(hooks=[nvtx_callback]) as sess:

tf.estimator.Estimator(hooks=training_hooks, ...)
```

# PROFIING MPI PROCESS

- ▶ To profile everything, putting the data into one file

```
nsys [nsys options] mpirun [mpi options]
```

- ▶ To profile everything, putting the data into each node into a separated file

```
mpirun [mpi options] nsys profile [nsys options]
```

# COMMAND EXAMPLES

## NGC TensorFlow ResNet50-1.5

- ▶ 1 GPU profiling

```
nsys profile -y 40 -d 20 -f true -w true -o tf_amp \  
-t osrt,nvtx,cuda,cublas,cudnn \  
python ./main.py --mode=training_benchmark --warmup_steps 200 \  
--num_iter 500 --iter_unit batch --results_dir=results --batch_size 64
```

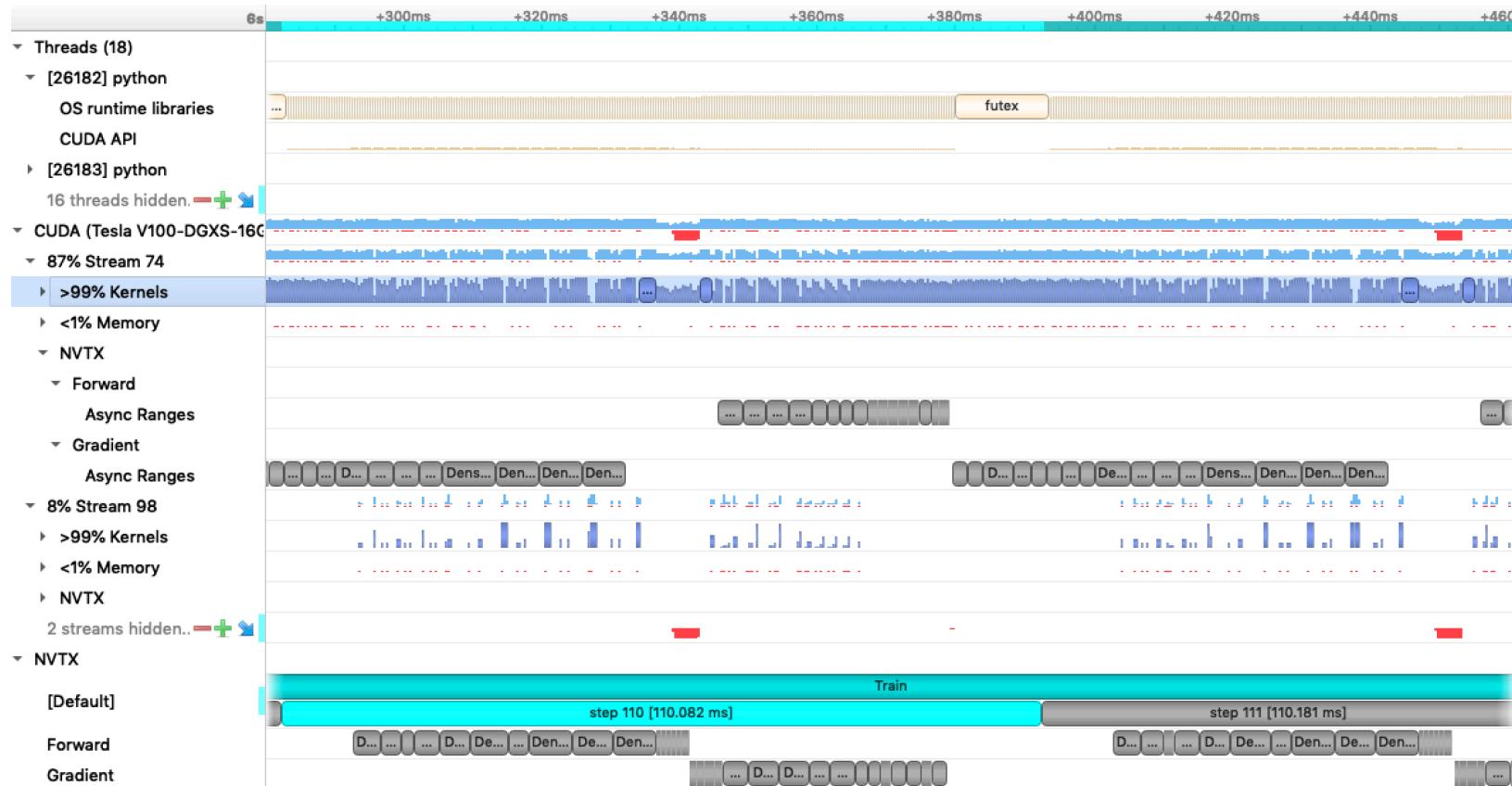
- ▶ 4 GPU profiling

```
nsys profile -y 40 -d 20 -f true -w true -o tf_amp_4gpu \  
-t osrt,nvtx,cuda,cublas,cudnn \  
mpiexec --allow-run-as-root --bind-to socket -np 4 \  
python ./main.py --mode=training_benchmark --warmup_steps 200 \  
--num_iter 500 --iter_unit batch --results_dir=results --batch_size 64
```

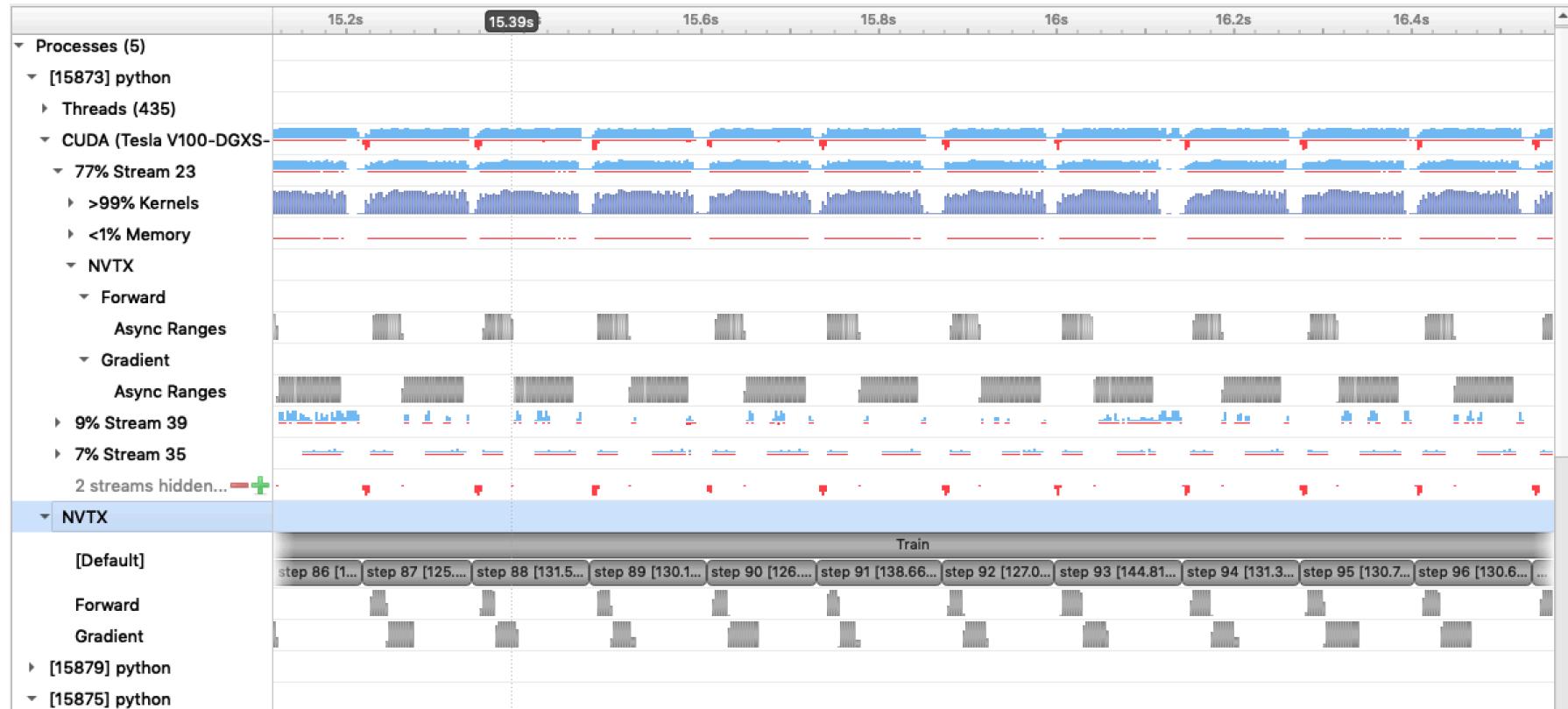
# TF WITH NVTX



# AMP ENABLED



# MULTI-GPU





NVIDIA®

