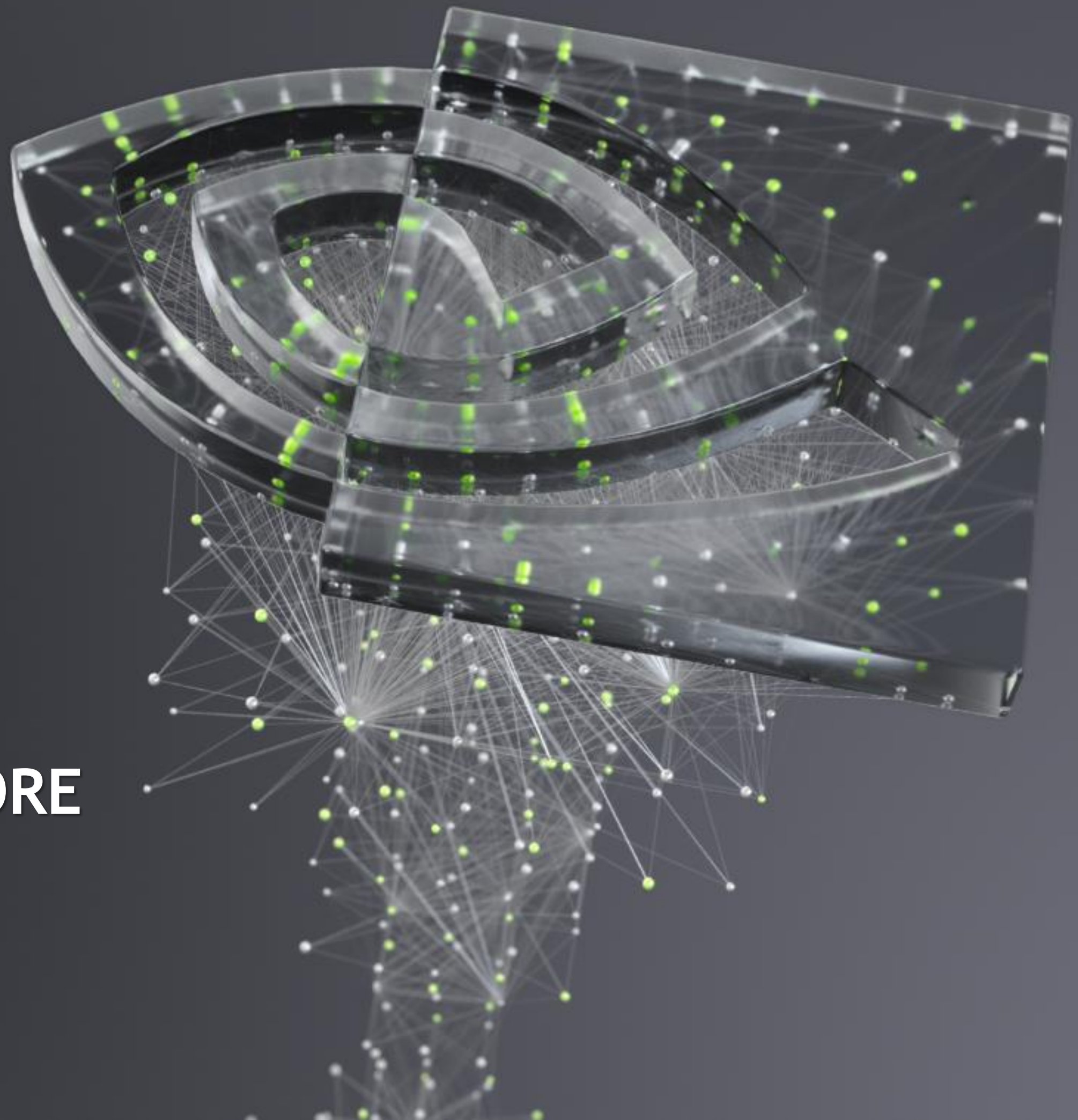# THE TRADER
# OF TOMORROW:
## SMARTER, FASTER, AND MORE PYTHONIC

John Ashley, October 2020

*Timing is everything.*

NVIDIA GTC

DON'T MISS JENSEN HUANG'S GTC KEYNOTE

October 5, 06:00 PDT (UTC-7)

SAVE THE DATE ›

*Timing is everything.*

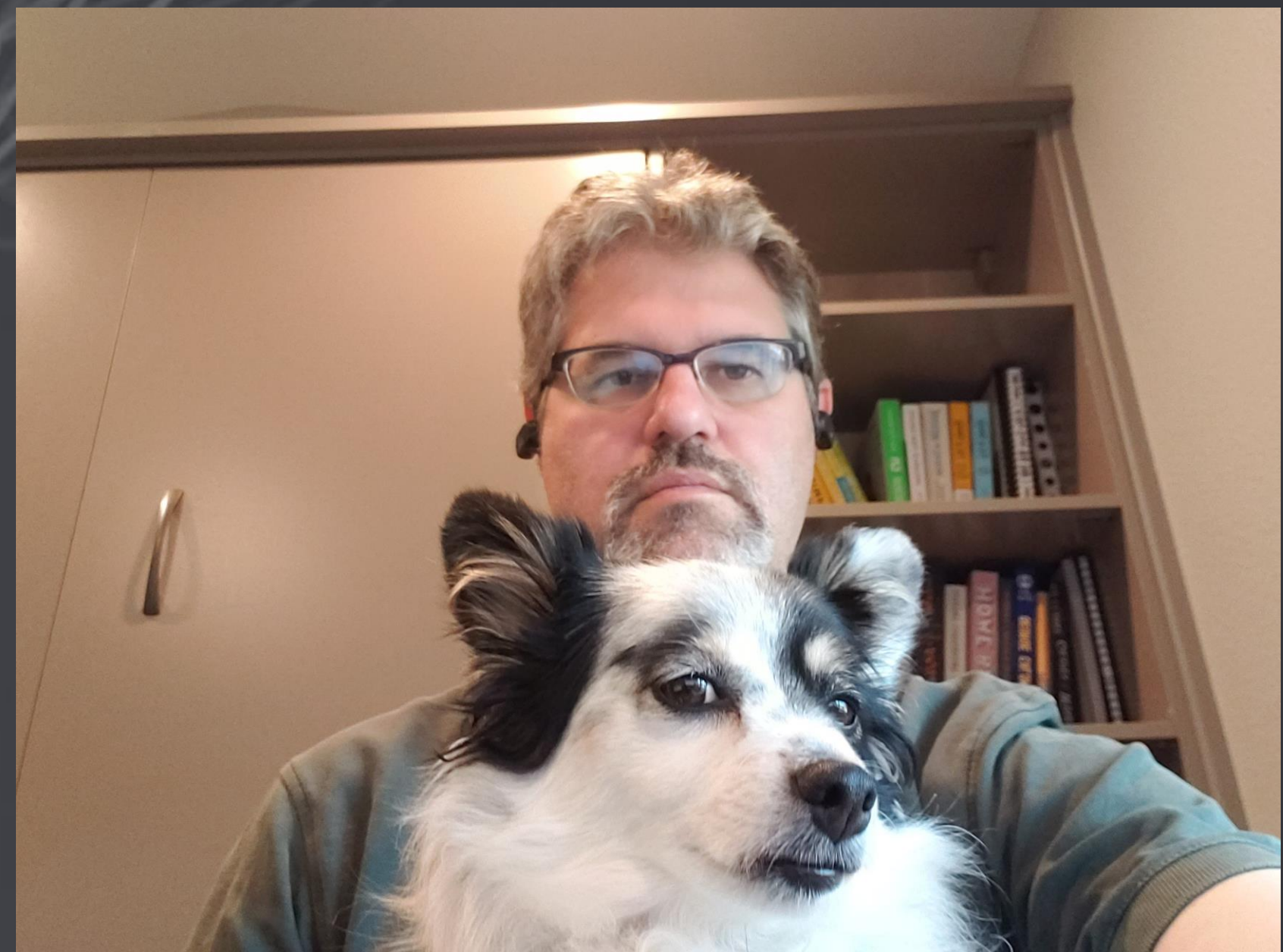DON'T MISS JENSEN HUANG'S
GTC KEYNOTE

October 5, 06:00 PDT (UTC-7)

SAVE THE DATE >

*Timing is everything.*

*Portions (all) of this talk were pre-recorded in front of a live studio audience.*

*Timing is everything.*

*Portions of this talk were pre-recorded in front of a live studio audience.*

# The Trader of the Future

## Smarter

NLP, BERT, and domain specific language models; or

"The more you learn, the more you earn." – Warren Buffet

## Faster

Optimizing AI models for inference; or

"Simplify, then add lightness." – Colin Chapman

## More Pythonic

Pointers to various useful bits of accelerated Python; or

"Every sufficiently advanced LISP application will eventually reimplement Python." – Hodgson's Law
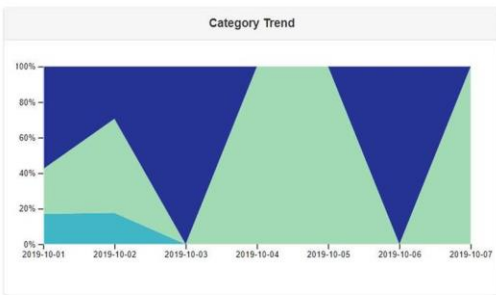
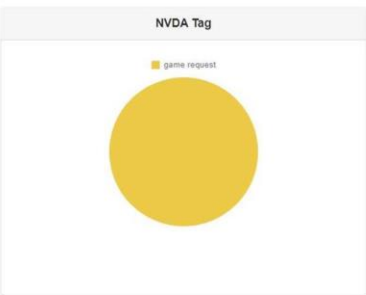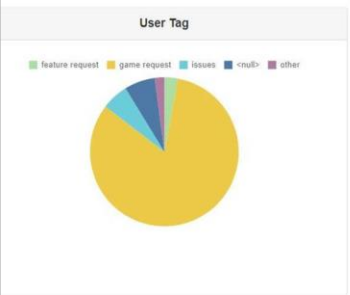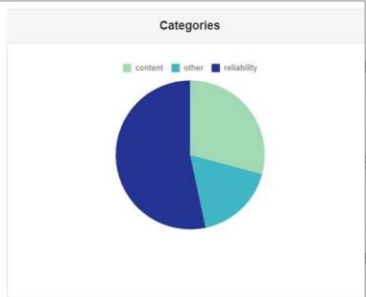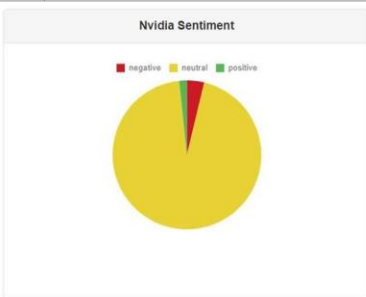"The more you learn, the more you earn."
– Warren Buffet

SMARTER

# UNDERSTAND GEFORCE NOW USERS

Thousands of gamer comments every day

# LANGUAGE UNDERSTANDING IMPROVEMENT
## Reaching human level

NOT A STAC BENCHMARK

**GLUE Aggregate Score**

Detect grammatical errors

Predict if movie review is positive or negative

Decide if an abstract correctly summarizes an article

Sentence-level Semantic equivalence

Basic reading comprehension

Pronoun disambiguation



Chart: GLUE Aggregate Score vs. time (Jan-2016 to Dec-2019). Y-axis 50–95.

Human Level (~86)

Data points: SkipThought (~61), InferSent (~64), DisSent (~62), ELMo (~66), GenSen (~66), BERT (~80.5), MT-DNN (~87.5), XLNet (~88.5), RoBERTa (~88.5), ALBERT (~89.5)

https://gluebenchmark.com/

9

NVIDIA.

# NATURAL LANGUAGE UNDERSTANDING

## BERT universal language model

**Input: Two sentences with 15% of words masked out**

1 = "Initially he supported himself and his ███ by farming on a plot ██ family land."

2 = "██ in turn attracted the attention of ██ *St.* ████ *Post-Dispatch*, which sent a reporter to Murray to █████ review Stubblefield's wireless ██████."

**Output 1: Reconstruct missing words**

family, of
this, the, Louis, personally, telephone

**Output 2: Is two the next sentence after one?**

NOT_NEXT_SENTENCE

https://arxiv.org/abs/1810.04805

# THE POWER OF TRANSFER LEARNING

## Domain Specific ASR - KENSHO & NVIDIA

### Evaluation WER vs Iteration

fine-tune    from-scratch

- Jasper trained on domain specific financial data outperformed all leading ASR models

- Fine tuning was faster and had more accuracy than training from scratch

- Enables quick start and has many benefits for consulting engagements

# DOMAIN SPECIFIC -- BIOBERT

## Context & Specialized Knowledge Matter

**Table 8.**

Biomedical question answering test results

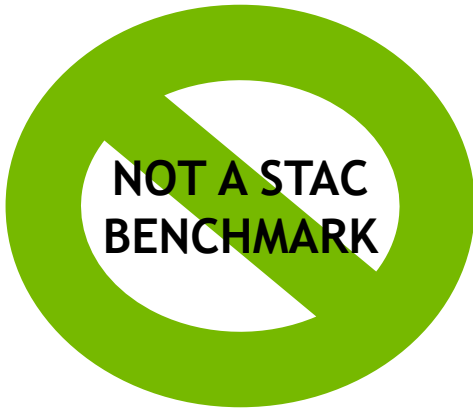| Datasets | Metrics | SOTA | BERT (Wiki + Books) | BioBERT v1.0 (+ PubMed) | (+ PMC) | (+ PubMed + PMC) | BioBERT v1.1 (+ PubMed) |
|---|---|---|---|---|---|---|---|
| BioASQ 4b | S | 20.01 | 27.33 | 25.47 | 26.09 | **28.57** | 27.95 |
|  | L | 28.81 | 44.72 | 44.72 | 42.24 | **47.82** | 44.10 |
|  | M | 23.52 | 33.77 | 33.28 | 32.42 | **35.17** | 34.72 |
| BioASQ 5b | S | 41.33 | 39.33 | 41.33 | 42.00 | 44.00 | **46.00** |
|  | L | 56.67 | 52.67 | 55.33 | 54.67 | 56.67 | **60.00** |
|  | M | 47.24 | 44.27 | 46.73 | 46.93 | 49.38 | **51.64** |
| BioASQ 6b | S | 24.22 | 33.54 | **43.48** | 41.61 | 40.37 | 42.86 |
|  | L | 37.89 | 51.55 | 55.90 | 55.28 | **57.77** | **57.77** |
|  | M | 27.84 | 40.88 | 48.11 | 47.02 | 47.48 | **48.43** |

*Notes:* Strict Accuracy (S), Lenient Accuracy (L) and Mean Reciprocal Rank (M) scores on each dataset are reported.

ngc.nvidia.com/catalog/resources/nvidia:biobert_for_tensorflow

NGC | CATALOG

Resources: nvidia:biobert_for_tensorflow

## BioBERT for TensorFlow

| Publisher | Application | Version | Created | Modified |
|---|---|---|---|---|
| NVIDIA | NLP | - | November 5, 2019 | September 24, 2020 |

| Framework | Model Format | Precision |
|---|---|---|
| TensorFlow | TensorFlow CKPT | FP16, FP32 |

**Description**
BERT for biomedical text-mining

**Labels**
Conversational AI   Deep Learning Examples   NLP   NLU   Natural Language Processing   Natural Language Understanding

Wget Resource   CLI Command

$ Latest version not available

Overview   Setup   Quick Start Guide   Advanced   Performance   Version History   File Browser   Release Notes   Related Collections

*This resource is a subproject of bert_for_tensorflow. Visit the parent project to download the code and get more information about the setup.*

In the original BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding paper, pre-training is done on Wikipedia and Books Corpus, with state-of-the-art results demonstrated on SQuAD (Stanford Question Answering Dataset) benchmark.

Meanwhile, many works, including BioBERT, SciBERT, NCBI-BERT, ClinicalBERT (MIT), ClinicalBERT (NYU, Princeton), and others at BioNLP'19 workshop, show that additional pre-training of BERT on large biomedical text corpus such as PubMed results in better performance in biomedical text-mining tasks.

This repository provides scripts and recipe to adopt the NVIDIA BERT code-base to achieve state-of-the-art results in the following biomedical text-mining benchmark tasks:

https://ngc.nvidia.com/catalog/resources/nvidia:biobert_for_tensorflow

# HOW TO BUILD YOUR OWN DOMAIN SPECIFIC ASR MODELS

https://ngc.nvidia.com/catalog/containers/nvidia:**nemo_asr_app_img**

**Pre-trained Quartznet model** LibriSpeech (old fiction books) → **Finetune Acoustic Model** with Wall Street Journal data (modern business news) → **Train Language Model** with Wall Street Journal → **Compare Models'** performance → **Export model** for Deployment



**Word Error Rate**

78% improvement!

NOT A STAC BENCHMARK

13

"Simplify, then add lightness." – Colin Chapman

FASTER

# FASTER TREES

Forest Inference Library

## DOES IT MOVE?

No → Should it? → No → No Problem

No → Should it? → Yes → [WD-40]

Yes → Should it? → Yes → No Problem

Yes → Should it? → No → [duct tape]



1000 trees, depth 10, Higgs

Predictor
- xgb_cpu
- xgb_gpu
- treelite
- FIL

$\mu$s per row - log

n_rows

Inference time scaling as batch size increases

https://medium.com/rapids-ai/rapids-forest-inference-library-prediction-at-100-million-rows-per-second-19558890bc35

# NVIDIA Fraud Detection Example Using PaySim Dataset

Accelerating Inferencing Using Forest Inferencing Library (FIL) on NVIDIA GPUs

35X Faster Using FIL



1M rows, 1000 trees, depth=10

CPU and GPU performance across datasets

Image credit: https://medium.com/rapids-ai/rapids-forest-inference-library-prediction-at-100-million-rows-per-second-19558890bc35
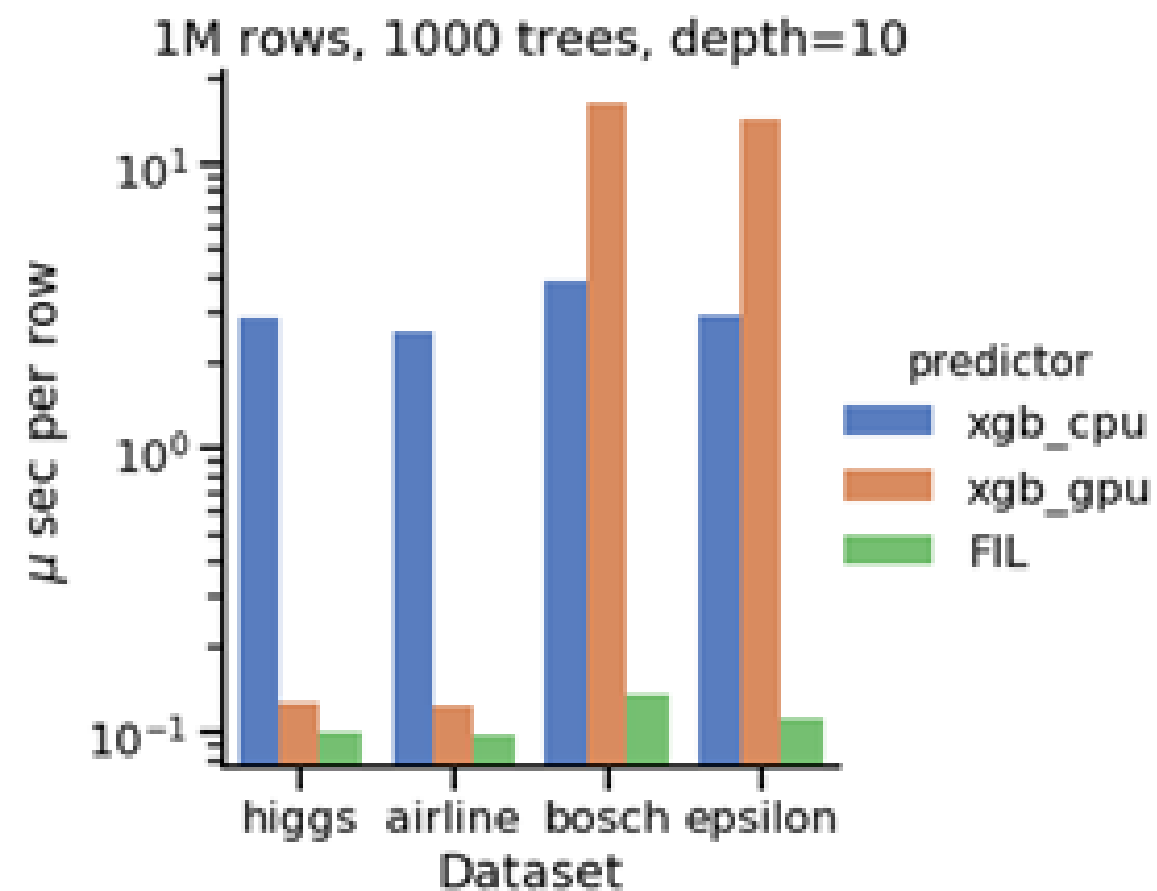
```python
from cuml import ForestInference
# Load the classifier previously saved with xgboost
# model_save()
import sklearn.datasets
model_path = 'xgb.model'

# Generate random sample data
fm = ForestInference.load(model_path,
output_class=True)

# Generate predictions (as a gpu array)
X_test, y_test =
sklearn.datasets.make_classification()

fil_preds_gpu = fm.predict(X_test.astype('float32'))
```
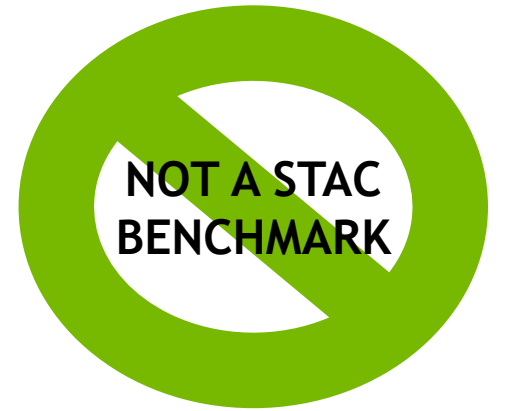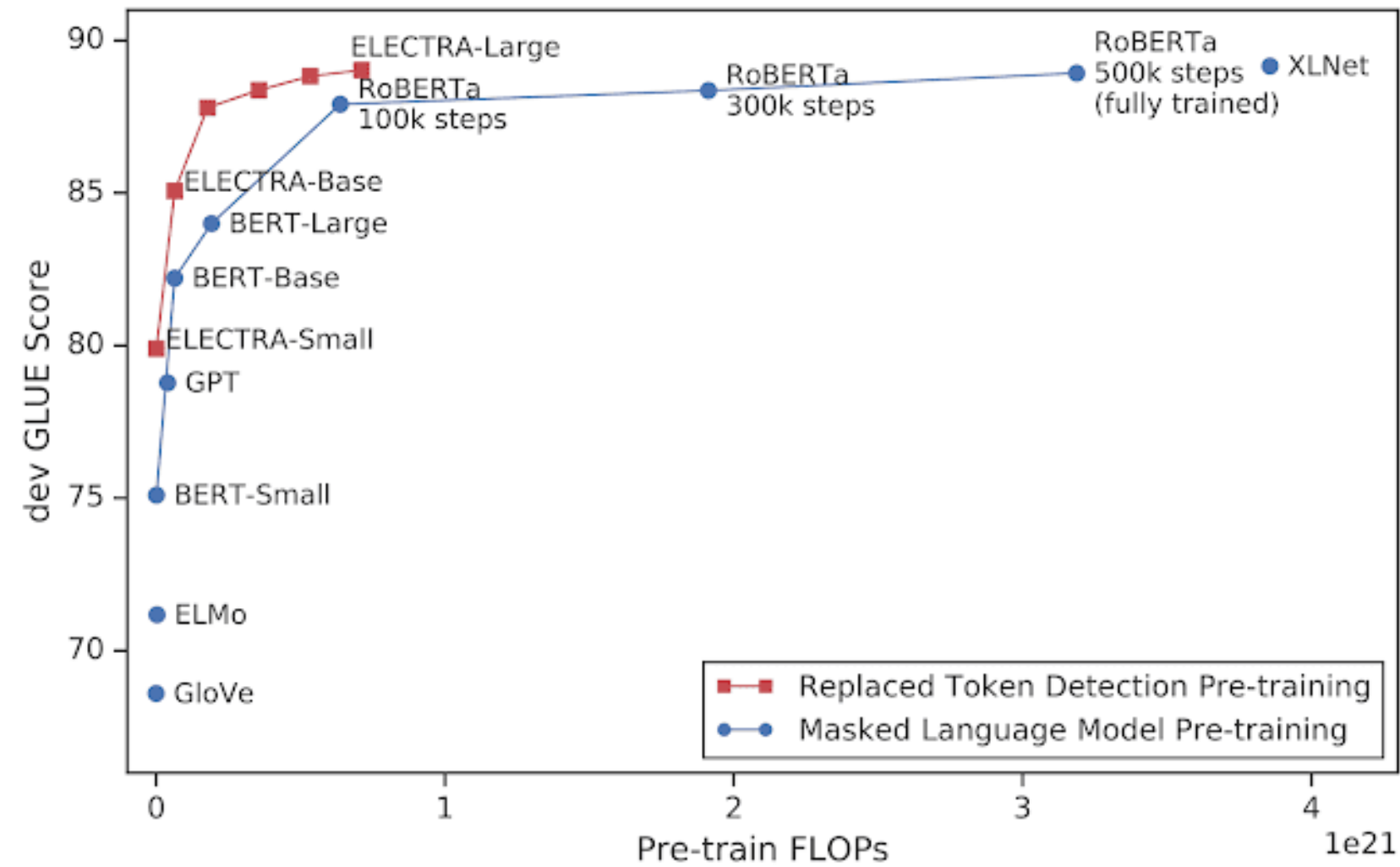
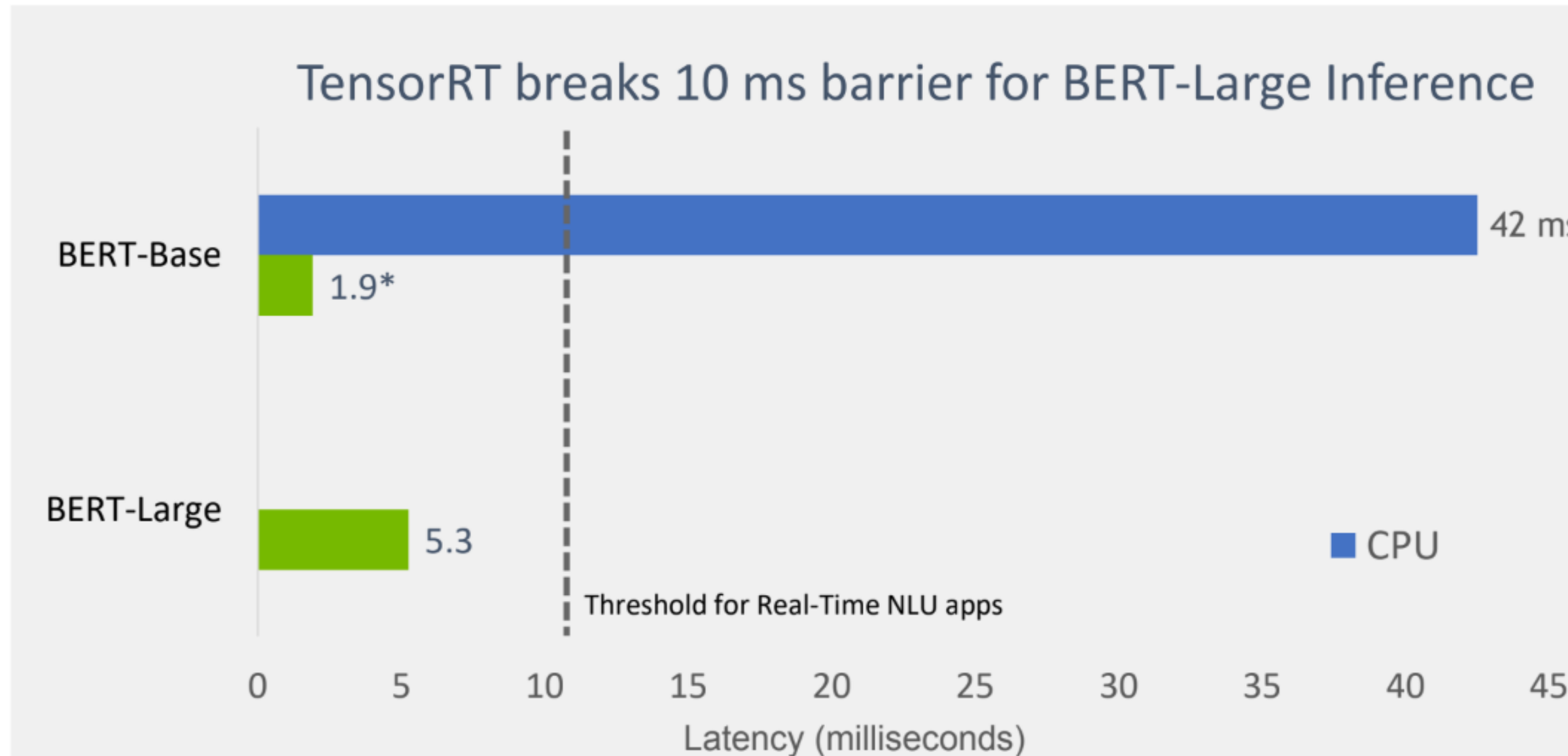# NLP MODELS ARE LARGE

## The Inference cost is high



NOT A STAC BENCHMARK

# BERT-LARGE INFERENCE IN 5.3 ms

## Makes Real-Time Natural Language Understanding Possible

### TensorRT breaks 10 ms barrier for BERT-Large Inference

**NOT A STAC BENCHMARK**

**BERT-Base**
- 42 ms (CPU)
- 1.9*

**BERT-Large**
- 5.3

Threshold for Real-Time NLU apps

■ CPU
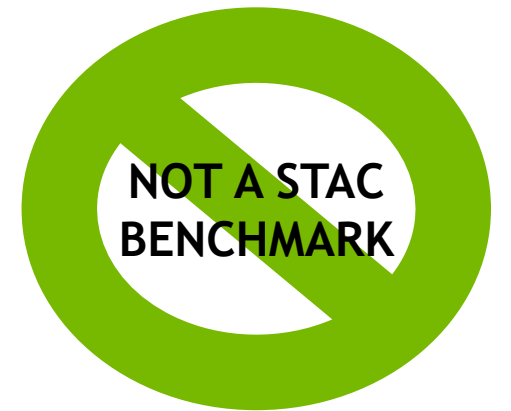
Latency (milliseconds) — 0, 5, 10, 15, 20, 25, 30, 35, 40, 45

\* In our tests, OpenVINO Release 2019 R2 did not execute BERT-Large and exited with an error

BERT Sample Code in TensorRT Repo
Jupyter Python Notebook
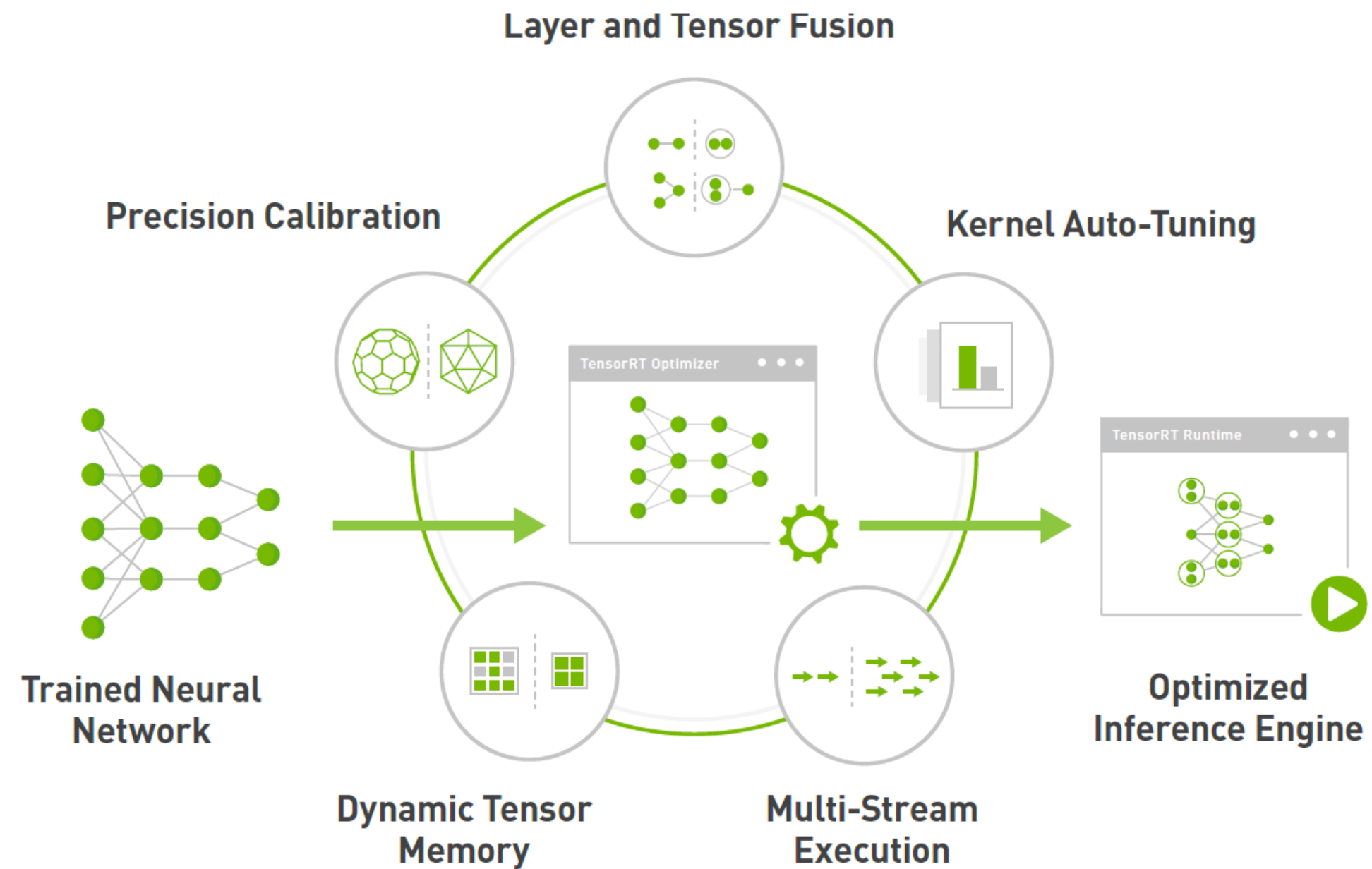Blog: Real-Time Natural Language Understanding with BERT Using TensorRT

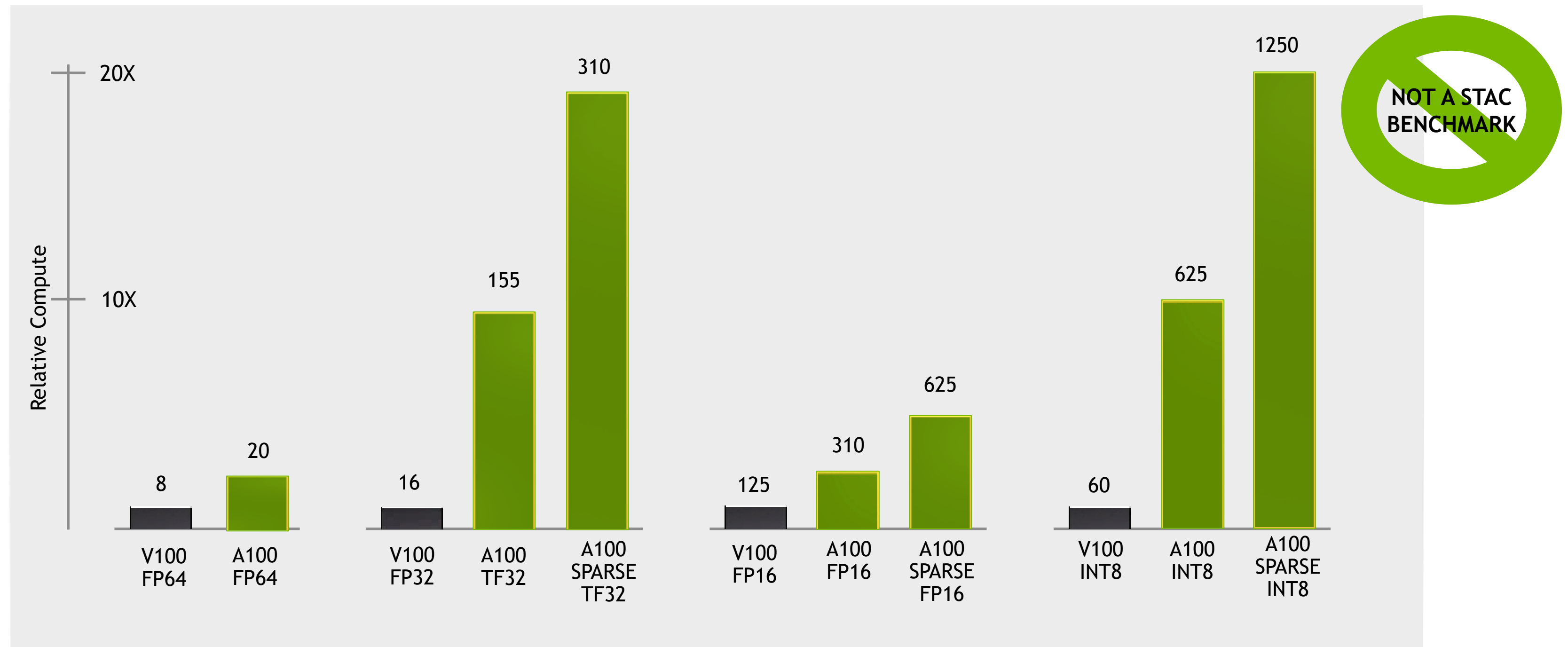developer.nvidia.com/tensorrt

18

NVIDIA.

# TENSORRT

## Simplification and the addition of lightness.

# INCREASING IMPORTANCE OF PRUNING AND QUANTIZATION

## Hardware acceleration for reduced precision arithmetic and sparsity

# QUANTIZATION

## The idea



FP32
(pre-quantized)

quantize

INT8
(quantized)

dequantize

FP32
(dequantized)

# PRUNING
## The idea

The opportunity:

- Reduced memory bandwidth

- Reduced memory footprint

- Acceleration (especially in presence of hardware acceleration)

Tambe, T., Yang, E. Y., Wan, Z., Deng, Y., Reddi, V. J., Rush, A., ... & Wei, G. Y. (2019). AdaptivFloat: A Floating-point based Data Type for Resilient Deep Learning Inference. *arXiv preprint arXiv:1909.13271.*

# CUSTOM PLUGINS
## Self-attention layer



Self-Attention Layer
(Before optimizations)

Self-Attention Layer
(With optimizations through TensorRT)

# TensorRT 7 FAMILY

## ASR, NLU & TTS | 1000+ Kernels | FP32, FP16, INT8



**Compiler Supports RNNs, Transformers and CNNs**

**20+ ONNX Ops & Dynamic Shapes Enhancements Accelerating Speech**

ASR With Jasper Example

NLU With BERT Example

TTS With Tacotron 2+Waveglow Blog & Example

**Get Started with ASR, NLU, TTS Today**

# TensorRT ONNX PARSER

## High-Performance Inference for ONNX Models

Optimize and deploy models from ONNX-supported frameworks to production

Apply TensorRT optimizations to any ONNX framework (Caffe 2, Microsoft Cognitive Toolkit, MxNet & PyTorch)

**Import TensorFlow and Keras through converters (tf2onnx, keras2onnx)**

Use with C++ and Python apps

20+ New Ops in TensorRT 7

Support for Opset 11 (See List of Supported Ops)

# TF-TRT = TF +TRT

Optimize TF inference while still using the TF ecosystem

# HOW TO USE?
## TF-TRT 2.x Workflow

TensorFlow + TensorRT

Existing workflow

Additional steps

**TF-TRT**

| SavedModel | → | LoadModel | → | Predict |

**TF-TRT**

SavedModel → LoadModel → Convert to TF-T$_{RT}$ → Predict

**TF-TRT**

SavedModel → LoadModel → Convert to TF-TRT → **Pre-build TRT engine** (Optional) → SavedModel

# INT8 TF-TRT API IN TENSORFLOW 2.0
## TF-TRT API

```python
from tensorflow.python.compiler.tensorrt import trt_convert as trt
conversion_params = trt.TrtConversionParams(
    precision_mode=trt.TrtPrecisionMode.INT8)


converter = trt.TrtGraphConverterV2(
    input_saved_model_dir=input_saved_model_dir,
    conversion_params=conversion_params)


converter.convert(calibration_input_fn=my_input_fn)


#optionally build TRT engines before deployment
converter.build(input_fn=my_input_fn)


converter.save(output_saved_model_dir)
```

Jupyter notebook example: https://github.com/tensorflow/tensorrt/blob/master/tftrt/examples/image-classification/TF-TRT-inference-from-saved-model.ipynb

## Developers' Software Optimizations Deliver Better Performance on the Same Hardware

**Monthly DL Framework Updates & Stack Optimizations Drive Performance**

**cuDNN** - Highly tuned standard training routines

**cuBLAS** - Highly tuned matrix multiplication

**DALI** – Moves compute intensive pre-processing to GPUs

**NCCL** – Faster training across multi-GPU architecture

**Framework** – Latest versions w/ newest features and superior perf

MONTHLY UPDATES DELIVER FASTER TRAINING PERFORMANCE

NOT A STAC BENCHMARK



- v20.05 (V100)   ■ v20.07 (V100)   ■ v20.07 (A100)

BERT-Large and ResNet-50 v1.5 Training performance with TensorFlow on a single node 8x V100 (32GB) & A100 (40GB). Mixed Precision.
Batch size for BERT: 10 (V100), 24 (A100), ResNet: 512 (V100, v20.05), 256 (v20.07)
DLRM Training performance with PyTorch on 1x V100 & 1x A100. Mixed Precision. Batch size 32768. DRLM trained with v20.03 and v20.07

"Every sufficiently advanced LISP application will eventually reimplement Python." – Hodgson's Law

MORE PYTHONIC

# CAN WE HAVE FAST DEVELOPMENT AND FAST EXECUTION?

Yes, if we leverage the whole Python ecosystem

**https://rapids.ai/**

## OPEN SOURCE

CONTRIBUTORS

ADOPTERS

# FRACTIONAL DIFFERENCING
## Easy & Fast

GPU-accelerated

Fractional Differencing for Time Series Stationarity

Ritchie Ng, Jie Fu, Tat-Seng Chua

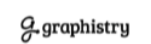| | 100k | 1m | 10m | 100m |
|---|---|---|---|---|
| | | | | |
| | | | | |
| Speed-up 1x T4 vs 8x vCPUs | 6.38 x | 67.38 x | 237.66 x | 328.08 x |
| Speed-up 1x V100 vs 8x vCPUs | 9.87 x | 83.76 x | 281.15 x | 411.72 x |

```python
from numba import cuda

def moving_dot_product_kernel(in_data, out, window_size, weights):
    ...
    # [Single loop] Compute fractional differencing values
    for i in range(cuda.threadIdx.x + window_size - 1, in_data.size, cuda.blockDim.x):
        # Compute dot product of preceding window_size rows
        rolling_dot_product = 0.0
        k = 0
        for j in range(i - window_size + 1, i + 1):
            rolling_dot_product += in_data[j] * weights[k][0]
            k += 1

        out[i] = rolling_dot_product

def frac_diff_gpu(df, d, floor=1e-3):
    ...
    gdf_raw = cudf.from_pandas(df).reset_index(drop=True)
    gdf_raw.columns = ['in_data']
    ...
    # Bring weights to GPU
    gdf_weights = cudf.DataFrame()
    ...
    threads_per_block = 518
    ...
    # Get fractionally differenced time series through GPU function
    gdf_raw_fd = gdf_raw.apply_chunks(moving_dot_product_kernel,
                incols=['in_data'],
                outcols=dict(out=np.float64),
                kwargs=dict(window_size=weights_window_size, weights=weights),
                chunks=list(range(0, data_length, trunk_size)) + [data_length],
                tpb=threads_per_block)

# Bring to CPU for normal manipulation
df_raw_fd = gdf_raw_fd.to_pandas().dropna().iloc[:-1, 1]

return df_raw_fd, weights
```

NOT A STAC BENCHMARK

https://www.researchgate.net/publication/335159299_GFD_GPU_Fractional_Differencing_for_Rapid_Large-scale_Stationarizing_of_Time_Series_Data_while_Minimizing_Memory_Loss

https://github.com/ritchieng/fractional_differencing_gpu/blob/master/notebooks/gpu_fractional_differencing.ipynb
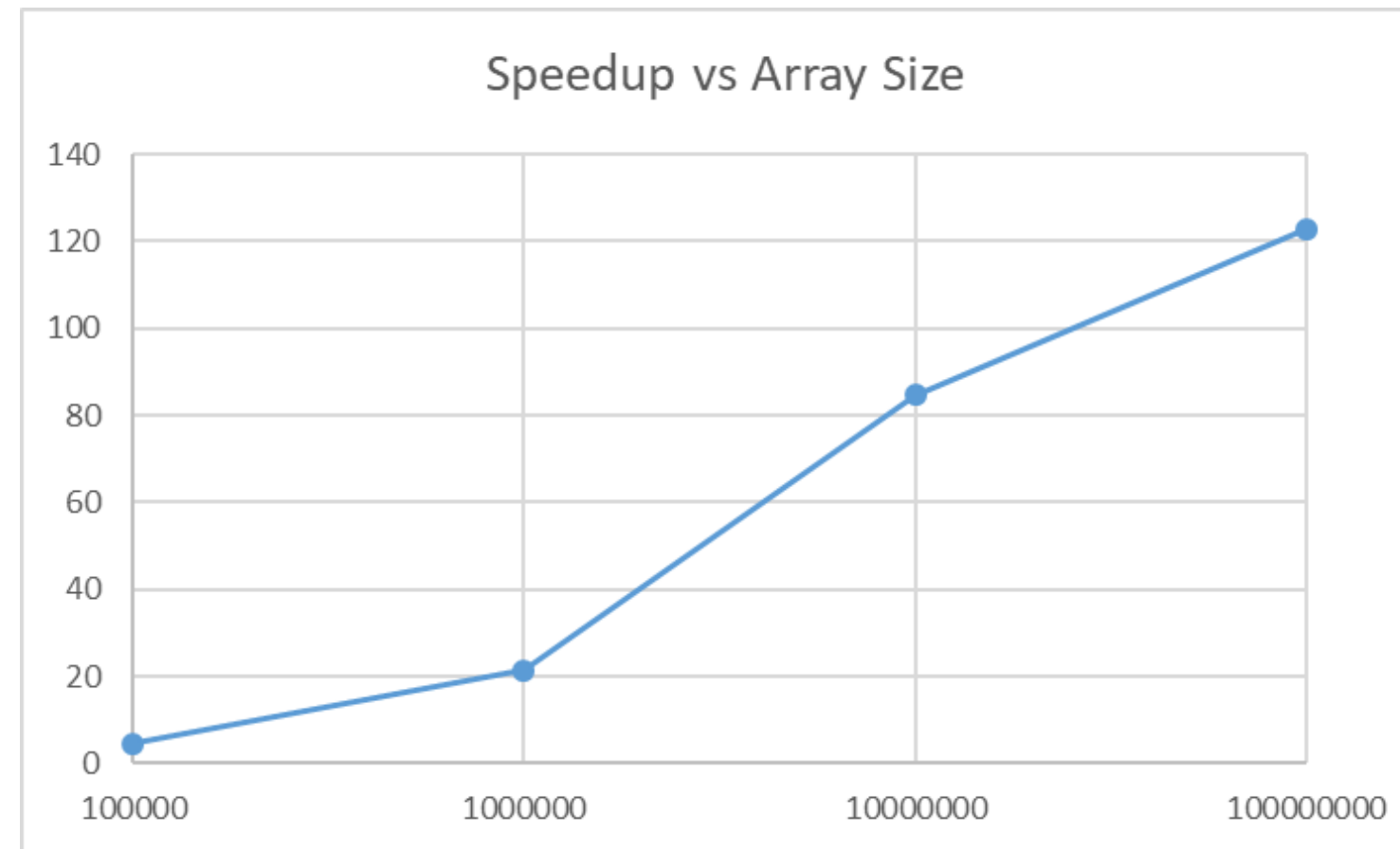
# FRACTIONAL DIFFERENCING

## With Numba JIT, even faster!

NOT A STAC
BENCHMARK

```
# allocate the output array
    gpu_out = numba.cuda.device_array_like(gpu_in)

    …

    # call the conv kernel
    kernel[(number_of_blocks,),
        (number_of_threads,),
        0,
        shared_buffer_size * 8](gpu_in,
                weights,
                gpu_out,
                window,
                array_len,
                thread_tile,
                min_periods)
    return gpu_out, weights_out
```

**Speedup vs Array Size**



```
@cuda.jit
def kernel(in_arr, weight_arr, out_arr, window,
        arr_len, thread_tile, min_size):
…
shared = cuda.shared.array(shape=0,
                dtype=numba.float64)
  …
  # copy the weights into the shared
  for j in range(0, window, block_size):
      element_id = tx + j
      if (((tx + j) < window) and (element_id < window)):
          shared[thread_tile * block_size + window - 1 + tx +
              j] = weight_arr[tx + j]
      cuda.syncthreads()
  # slice the shared memory for each threads
  start_shared = tx * thread_tile
  his_len = min(window - 1,
          starting_id + tx * thread_tile)
  # slice the global memory for each threads
  start = starting_id + tx * thread_tile
  end = min(starting_id + (tx + 1) * thread_tile, arr_len)
  sub_outarr = out_arr[start:end]
  sub_len = end - start
  conv_window(shared, his_len, sub_outarr,
          window, sub_len,
          window - 1 + start_shared,
          thread_tile * block_size + window - 1,
          min_size)
```

https://medium.com/rapids-ai/fast-fractional-differencing-on-gpus-using-numba-and-rapids-part-1-b271a6b68b41

NVIDIA.

THANK YOU!