# HPC for DL

Serafina Di Gioia
PostDoctoral Researcher @ ICTP

# From LISP to the DL revolution…



ARTIFICIAL INTELLIGENCE
Early artificial intelligence stirs excitement.

MACHINE LEARNING
Machine learning begins to flourish.

DEEP LEARNING
Deep learning breakthroughs drive AI boom.

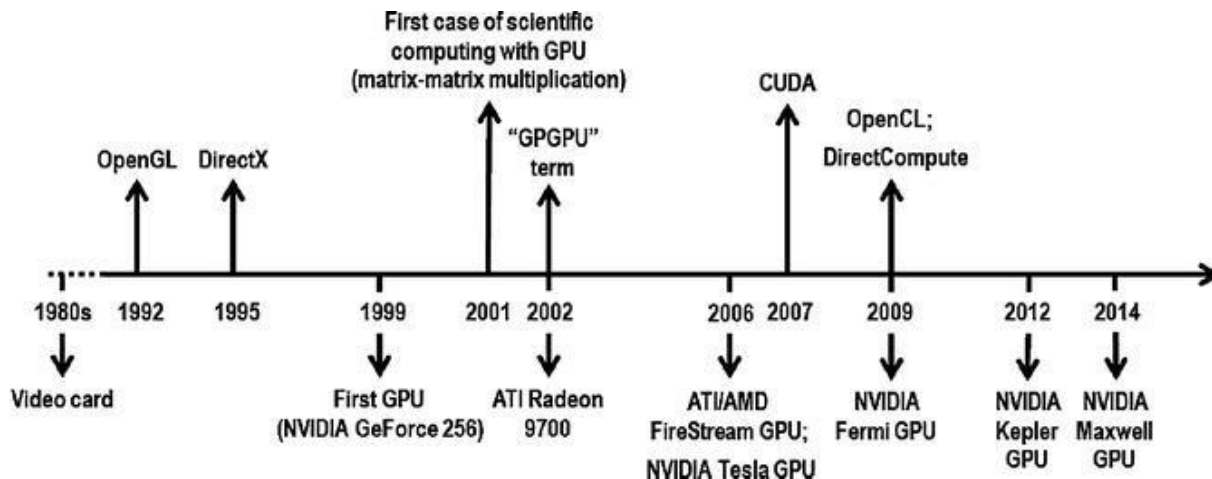1950's    1960's    1970's    1980's    1990's    2000's    2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

# Main ingredients for DL breakthrough

- large datasets available (e.g IMAGENET)
- GPUs development (in particular, CUDA introduction)
- increased involvement of developers from CV and scientific communities

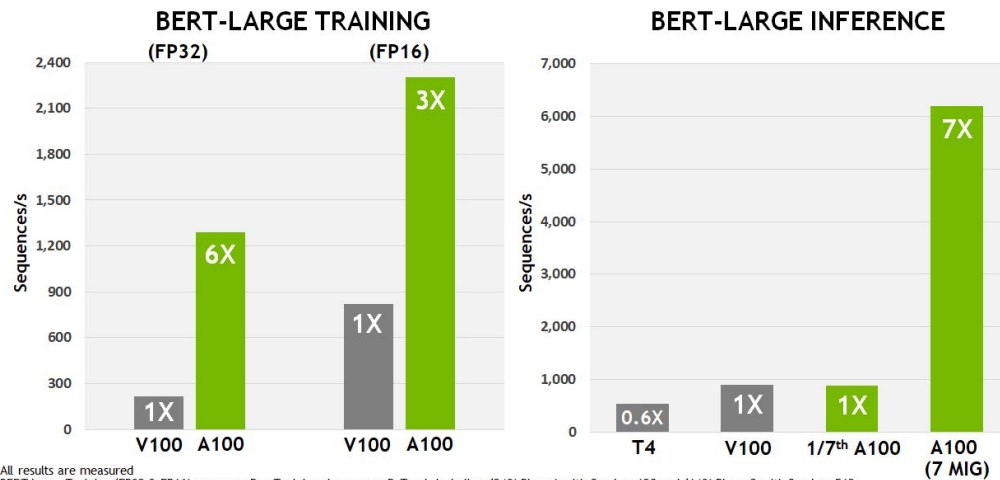The DL era starts few years after that CUDA came to light

# Incredible performance increase of NVIDIA

| NVIDIA Accelerator Specification Comparison | | | |
|---|---|---|---|
| | A100 | V100 | P100 |
| FP32 CUDA Cores | 6912 | 5120 | 3584 |
| Boost Clock | ~1.41GHz | 1530MHz | 1480MHz |
| Memory Clock | 2.4Gbps HBM2 | 1.75Gbps HBM2 | 1.4Gbps HBM2 |
| Memory Bus Width | 5120-bit | 4096-bit | 4096-bit |
| Memory Bandwidth | 1.6TB/sec | 900GB/sec | 720GB/sec |
| VRAM | 40GB | 16GB/32GB | 16GB |
| Single Precision | 19.5 TFLOPs | 15.7 TFLOPs | 10.6 TFLOPs |
| Double Precision | 9.7 TFLOPs (1/2 FP32 rate) | 7.8 TFLOPs (1/2 FP32 rate) | 5.3 TFLOPs (1/2 FP32 rate) |
| INT8 Tensor | 624 TOPs | N/A | N/A |
| FP16 Tensor | 312 TFLOPs | 125 TFLOPs | N/A |
| TF32 Tensor | 156 TFLOPs | N/A | N/A |
| Interconnect | NVLink 3 12 Links (600GB/sec) | NVLink 2 6 Links (300GB/sec) | NVLink 1 4 Links (160GB/sec) |
| GPU | GA100 (826mm2) | GV100 (815mm2) | GP100 (610mm2) |
| Transistor Count | 54.2B | 21.1B | 15.3B |
| TDP | 400W | 300W/350W | 300W |
| Manufacturing Process | TSMC 7N | TSMC 12nm FFN | TSMC 16nm FinFET |
| Interface | SXM4 | SXM2/SXM3 | SXM |
| Architecture | Ampere | Volta | Pascal |

## UNIFIED AI ACCELERATION

### BERT-LARGE TRAINING
(FP32)    (FP16)

Sequences/s

V100: 1X    A100: 6X    V100: 1X    A100: 3X

### BERT-LARGE INFERENCE

Sequences/s

T4: 0.6X    V100: 1X    1/7th A100: 1X    A100 (7 MIG): 7X

All results are measured
BERT Large Training (FP32 & FP16) measures Pre-Training phase, uses PyTorch including (2/3) Phase1 with Seq Len 128 and (1/3) Phase 2 with Seq Len 512,
V100 is DGX1 Server with 8xV100, A100 is DGX A100 Server with 8xA100, A100 uses TF32 Tensor Core for FP32 training
BERT Large Inference uses TRT 7.1 for T4/V100, with INT8/FP16 at batch size 256. Pre-production TRT for A100, uses batch size 94 and INT8 with sparsity

# Why GPUs for DL?

1) **Neural networks are embarrassingly parallel algorithm**
2) **most of the operations performed in DL models can be rewritten as matrix multiplications**
3) **big datasets require to perform big matrix computation (**extremely slow on CPU with respect to GPU**)**
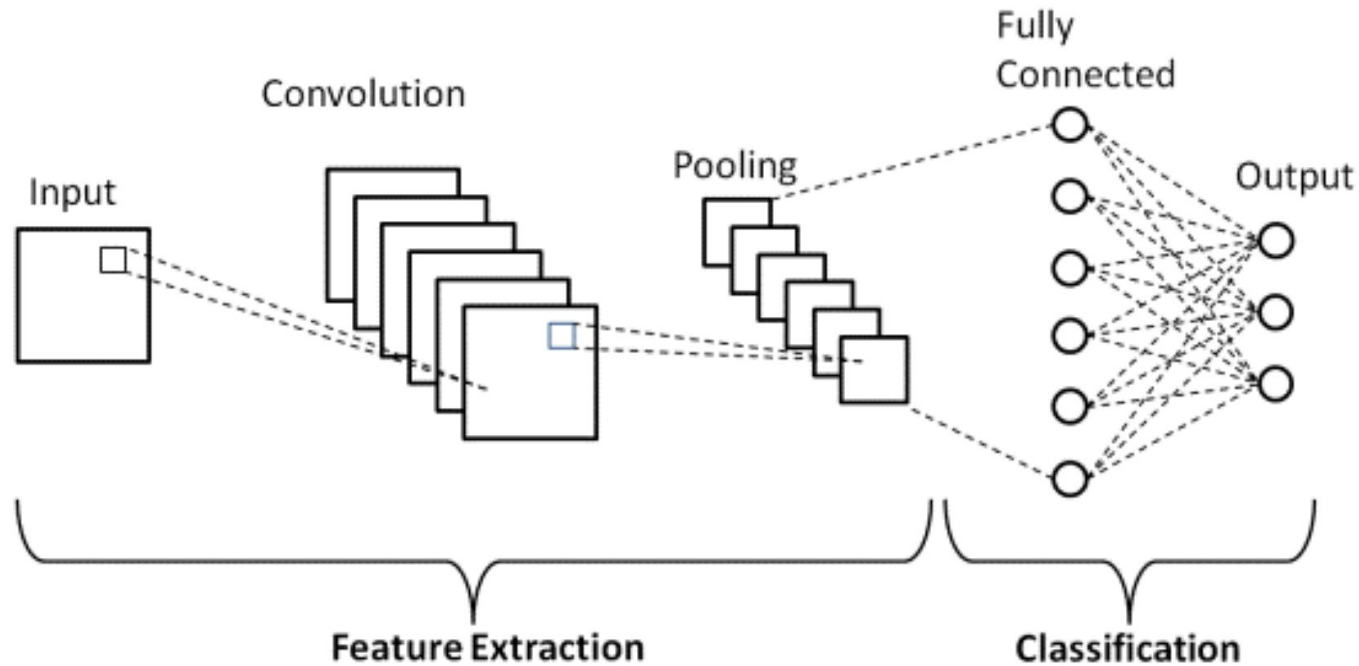4) **well established libraries, with specific classes for ML objects (e.g. cuDNN, more recently tensorRT)**

**and we know that GPUs are very good in solving specific parallel tasks (e.g matrix multiplication) , thanks to**

- +1000 cores (>100K threads)
- **SIMD / SIMT**
- **high memory bandwidth**
- **newer GPUs have also tensor cores (particularly suited to tensor ops typical of NNs), and mixed precision**

**However, also GPUs have limitations:**

- ○ GPUs might not be as efficient for extreme sparse networks, due to the overhead of managing sparse data structures.
- ○ Some specialized sparse operations might not be as optimized as dense operations on GPUs.

# DL for image classification: Convolutional Neural Networks (CNN)

# Everything in DL turns to be a matrix multiplication at the end…

Let's take a look at basic element of CNN: convolution layer

Consider the case where we are applying (2,2) kernel



to a (3,3) matrix:

The convolution can be rewritten as

# Analyzing the computational workload of DL models



training — 80-95%

data encoding — 2-5%

testing/inference — 2-5%

Forward pass
*We predict a label*

Loss Function

Optimizer (SGD, ...)

@Thom_Wolf

# Different strategies for Multi-GPUs training/evaluation

We can identify 5 different categories of parallelism

model parallelism

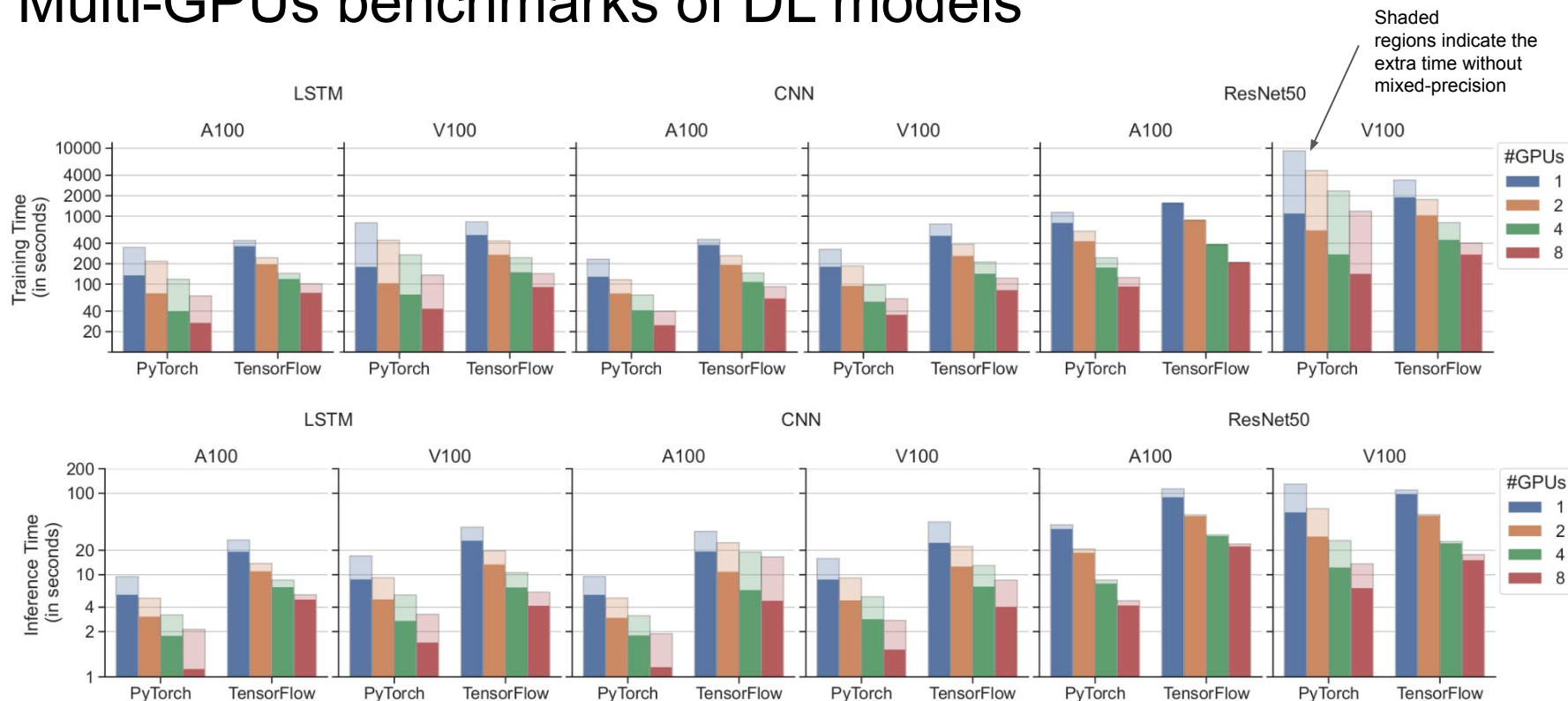tensor parallelism

data parallelism

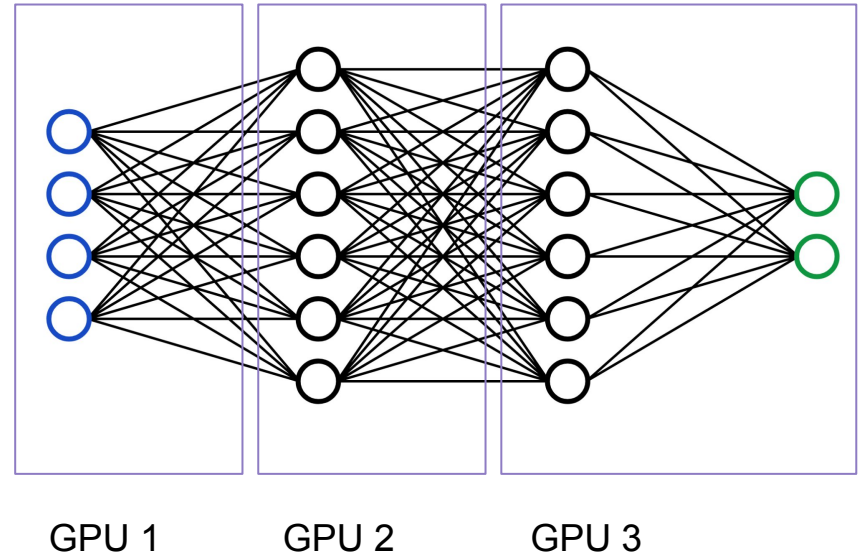sequence parallelism

pipeline parallelism

# Multi-GPUs benchmarks of DL models



Shaded regions indicate the extra time without mixed-precision

# Model parallelism

In this parallelism framework we choose to put different layers of the NN on different GPUs

to work around GPU memory limits



GPU 1          GPU 2          GPU 3
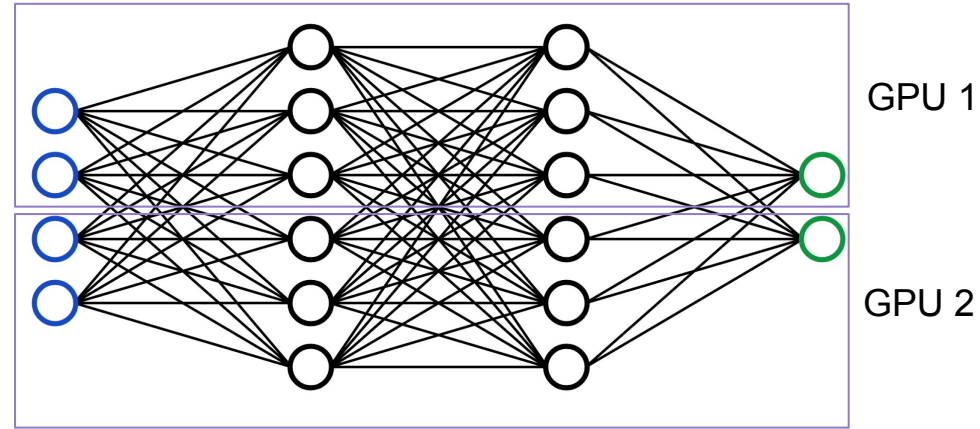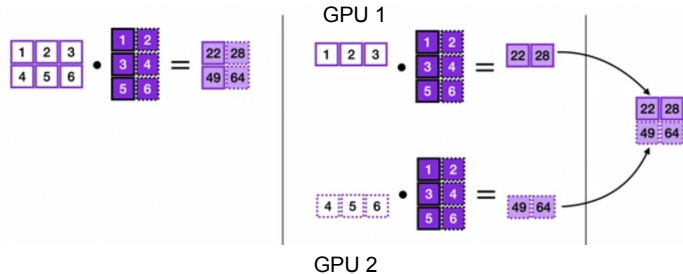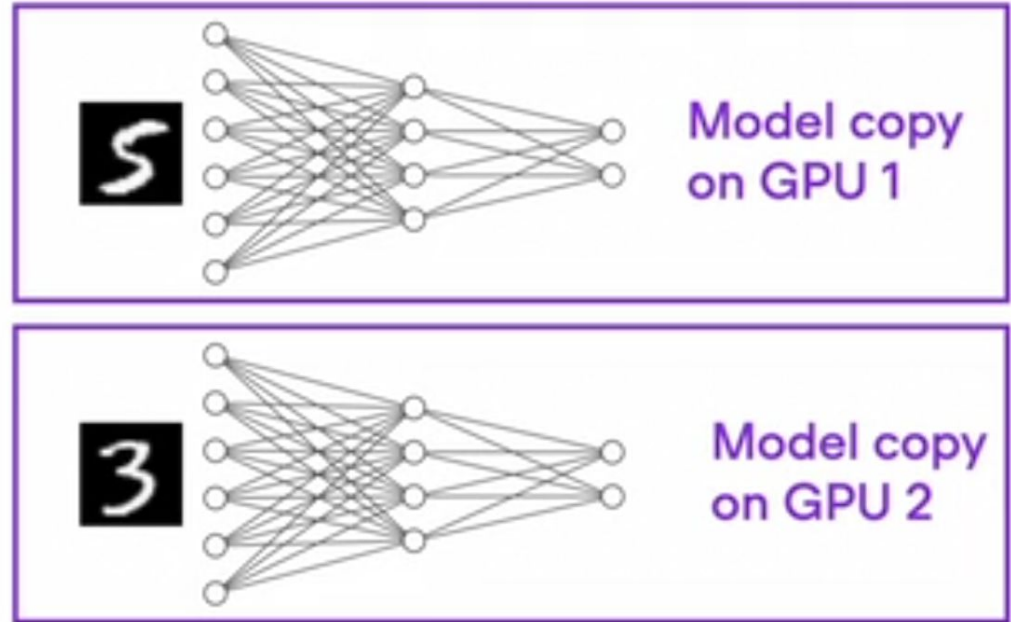
# Tensor parallelism

In this framework we split the tensor operation done at each layer among different GPUs
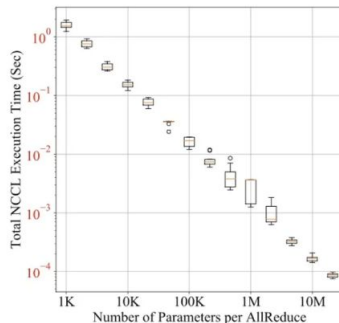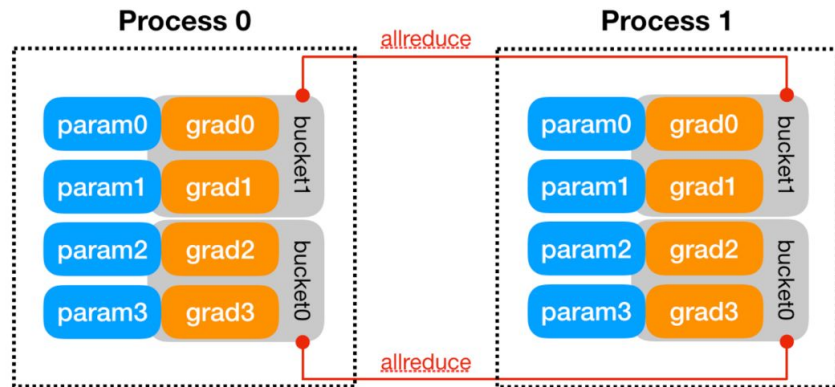
similarly to what we would have done for matmul
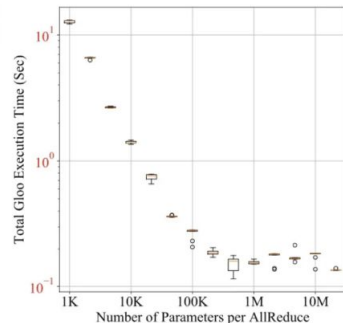
# Data Parallelism

In this framework we split batches to train DL model into different GPUs



Model copy on GPU 1

Model copy on GPU 2

# Distributed Data Parallel (DDP)



(a) NCCL  (b) GLOO



**Algorithm 1:** DistributedDataParallel

**Input:** Process rank $r$, bucket size cap $c$, local model $net$

1 **Function** constructor($net$):
2      **if** $r=0$ **then**
3          broadcast $net$ states to other processes
4      init buckets, allocate parameters to buckets in the reverse order of net.parameters()
5      **for** $p$ **in** net.parameters() **do**
6          acc $\leftarrow$ $p$.grad_accumulator
7          acc $\rightarrow$ add_post_hook(autograd_hook)

8 **Function** forward($inp$):
9      out $= net$($inp$)
10      traverse autograd graph from out and mark unused parameters as ready
11      **return** out

12 **Function** autograd_hook($param\_index$):
13      get bucket $b_i$ and bucket $offset$ using $param\_index$
14      get parameter $var$ using $param\_index$
15      view $\leftarrow b_i.narrow(offset, var.size())$
16      view.copy_(var.grad)
17      **if** $all\ grads\ in\ b_i\ are\ ready$ **then**
18          mark $b_i$ as ready
19      launch AllReduce on ready buckets in order
20      **if** $all\ buckets\ are\ ready$ **then**
21          block waiting for all AllReduce ops

# Common guidelines by Pytorch documentation

1. Use DistributedDataParallel (DDP), if your model fits in a single GPU but you want to easily scale up training using multiple GPUs.

   - Use torchrun, to launch multiple pytorch processes if you are using more than one node.
   - See also: Getting Started with Distributed Data Parallel

2. Use FullyShardedDataParallel (FSDP2) when your model cannot fit on one GPU.

   - See also: Getting Started with FSDP2

3. Use Tensor Parallel (TP) and/or Pipeline Parallel (PP) if you reach scaling limitations with FSDP2.

   - Try our Tensor Parallelism Tutorial
   - See also: TorchTitan end to end example of 3D parallelism
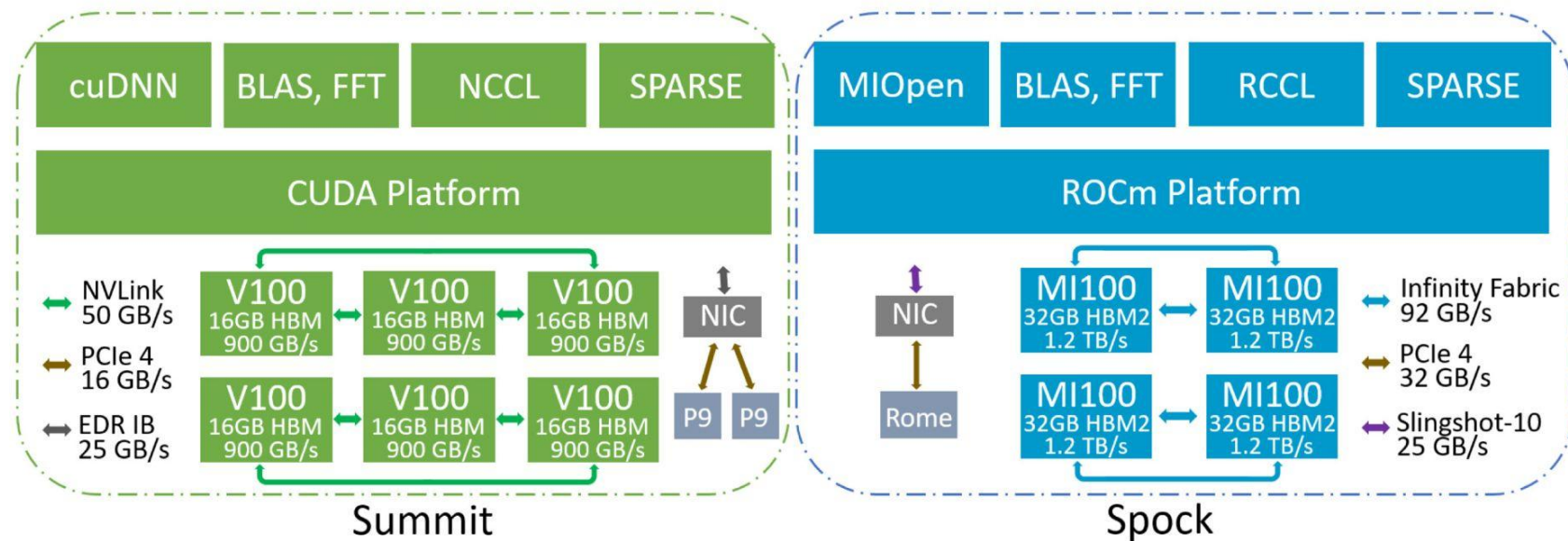
# More complex strategies for DL training

Sequence parallelism and pipeline parallelism frameworks

are obtained combining the previous approaches,

and are typically applied to DL models dealing with

spatio-temporal data.
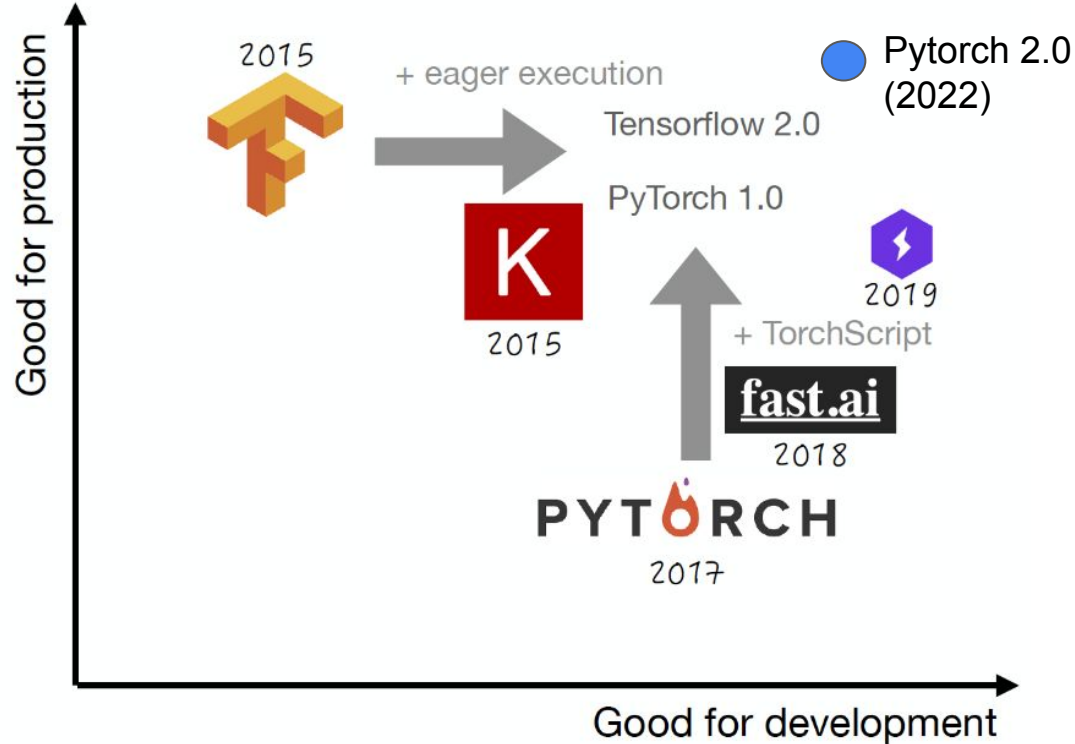
# What's behind Pytorch/Tensorflow?

# Backend of torch.distributed

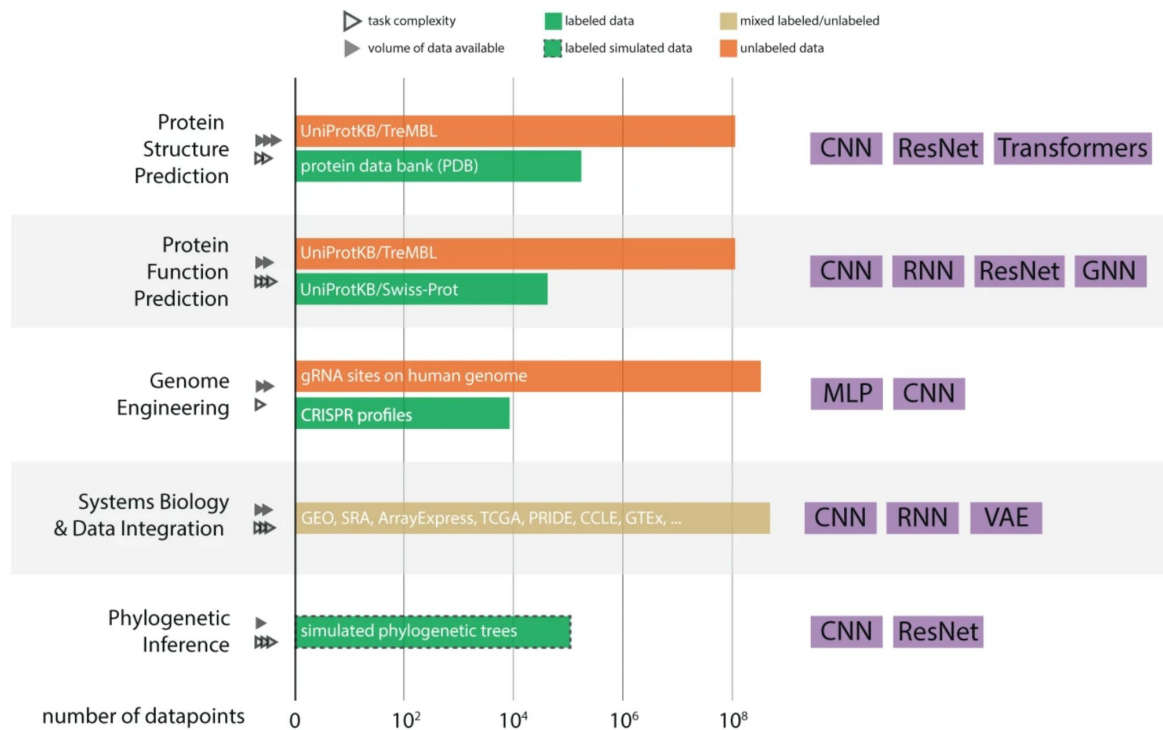| Backend | `gloo` | | `mpi` | | `nccl` | |
|---|---|---|---|---|---|---|
| Device | CPU | GPU | CPU | GPU | CPU | GPU |
| send | ✓ | ✗ | ✓ | ? | ✗ | ✓ |
| recv | ✓ | ✗ | ✓ | ? | ✗ | ✓ |
| broadcast | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| all_reduce | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| reduce | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| all_gather | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| gather | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| scatter | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| reduce_scatter | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| all_to_all | ✓ | ✓ | ✓ | ? | ✗ | ✓ |
| barrier | ✓ | ✗ | ✓ | ? | ✗ | ✓ |

# Python most used libraries/frameworks for DL

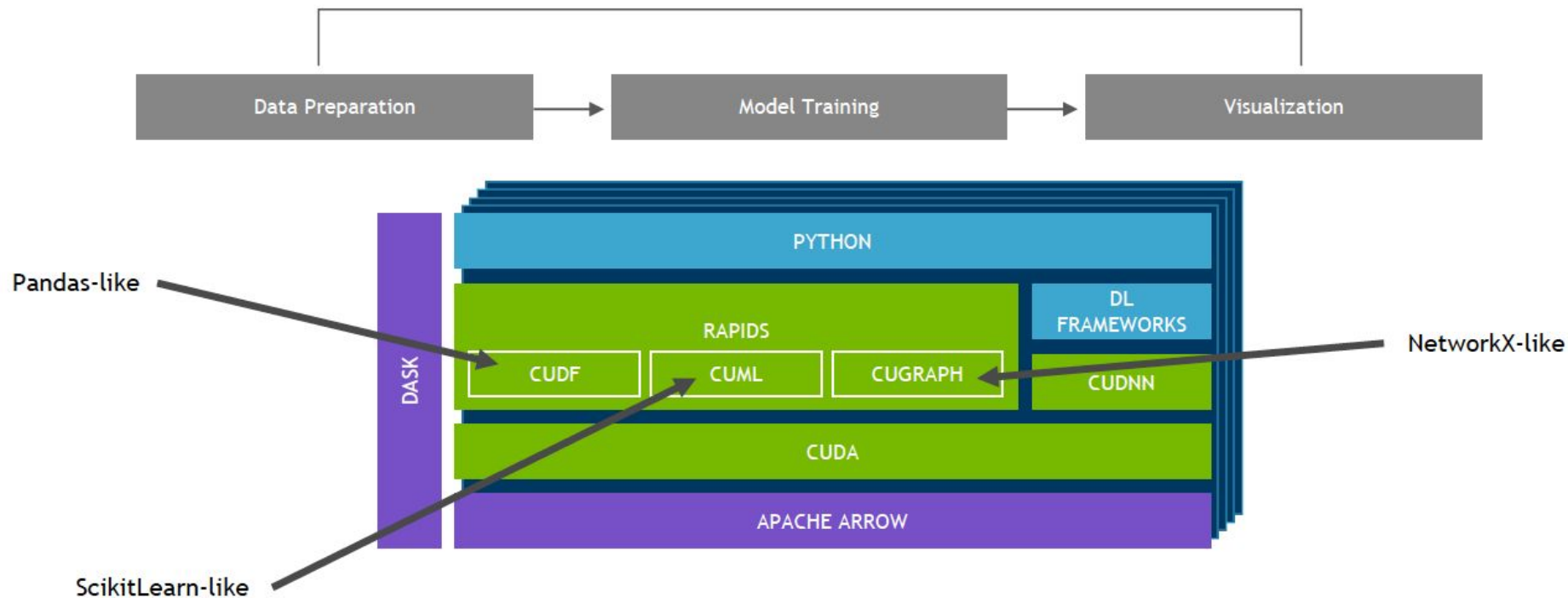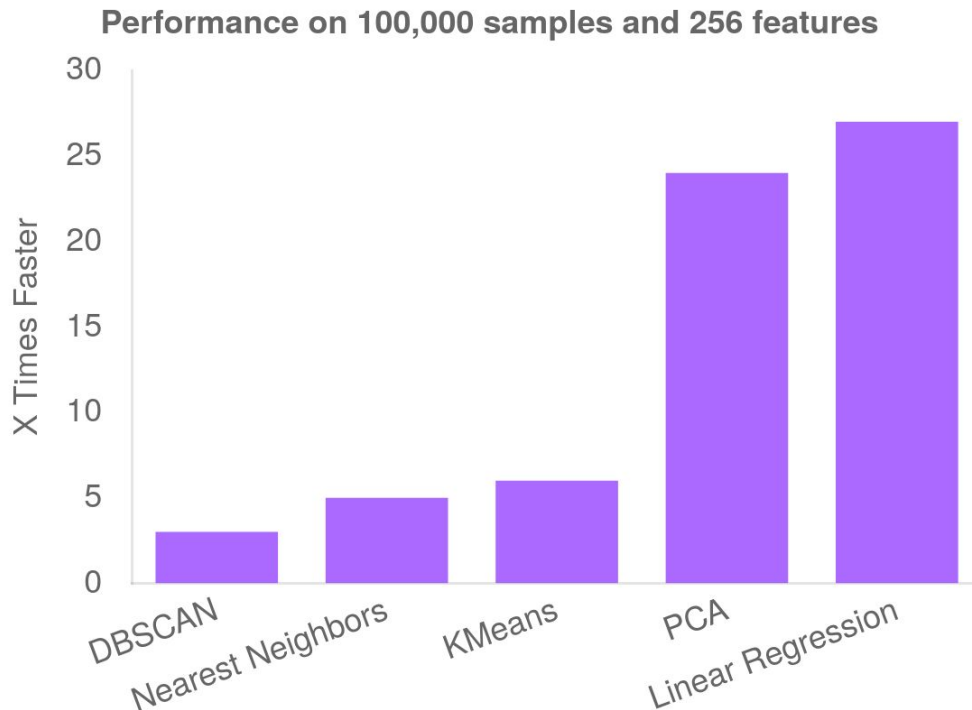# Do we need GPUs also for other ML tasks?

# Typical sizes of DL data sets



The increasing complexity of the new datasets, typical of big data epoch motivates the need for GPU-based libraries for feature analysis and data preprocessing

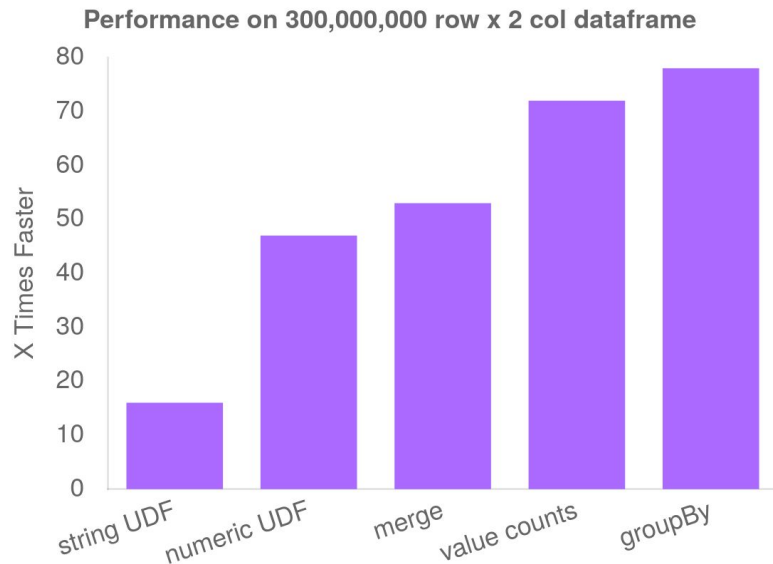# GPU-based libraries outside of Pytorch

# GPU for classical ML

**Performance on 100,000 samples and 256 features**



* Benchmark on AMD EPYC 7642 (using 1x 2.3GHz CPU core) w/ 512GB and NVIDIA A100 80GB (1x GPU) w/ scikit-learn v1.2 and cuML v23.02

SMR 4054 - HPC, AI and Regional Climate Modeling

# GPUs for data preprocessing

**Performance on 300,000,000 row x 2 col dataframe**



* Benchmark on AMD EPYC 7642 (using 1x 2.3GHz CPU core) w/ 512GB and NVIDIA A100 80GB (1x GPU) w/ pandas v1.5 and cuDF v23.02

# References

Milan Jain, Sayan Ghosh, Sai Pushpak Nandanoori, 2022, *Workload Characterization of a Time-Series Prediction System for Spatio-Temporal Data*

Junqi Yin et. al, 2021, *Comparative evaluation of deep learning workloads for leadership-class systems*

NVIDIA Booklet on GPU development, 2021

*Deep Learning: A Comprehensive Overview on Techniques, Taxonomy, Applications and Research Directions*