# FULLY SHARDED DATA PARALLEL

GPUs goes brrrr

Mert Bozkir
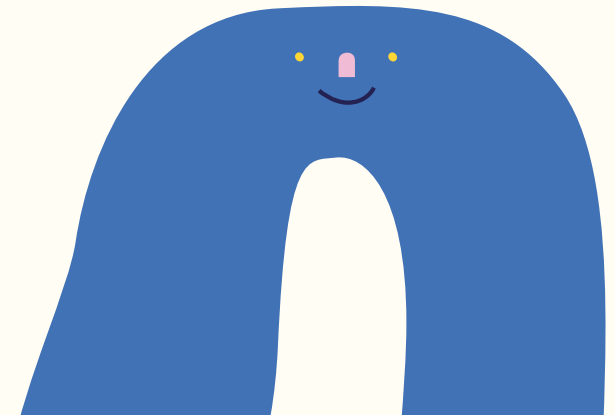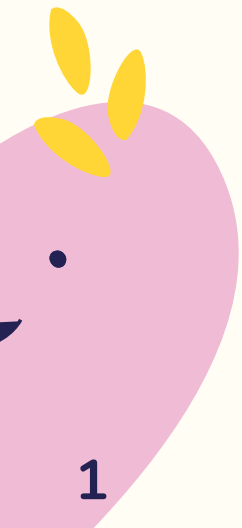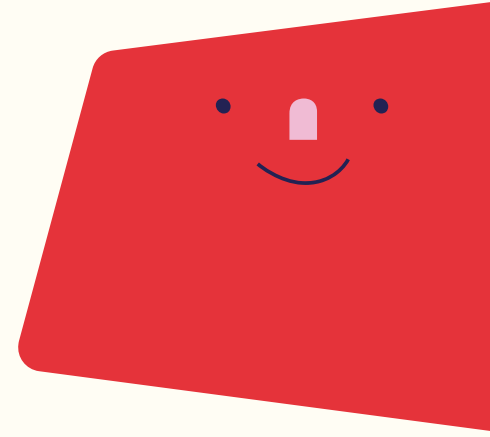
# TABLE OF CONTENT

1

# INTRODUCTION

## Self-taught ML Engineer

- a not-drop-out CS student

- crazy enthusiast about LLMs! 🔥

- a slow but concentrated learner

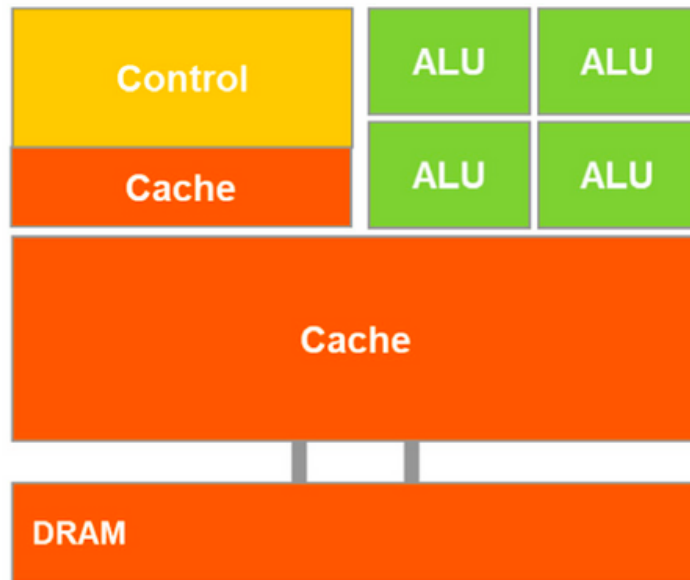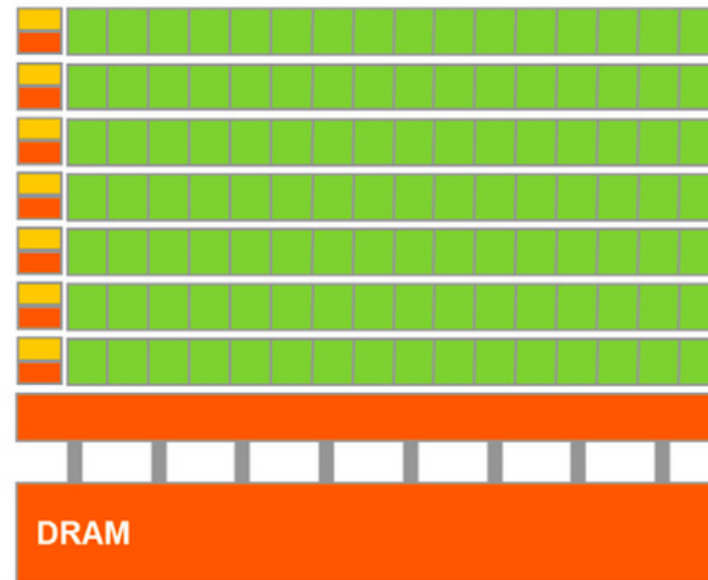# HOW & WHY PARALLELIZATIONS WORKS

Decrease the time,

Increase the compute
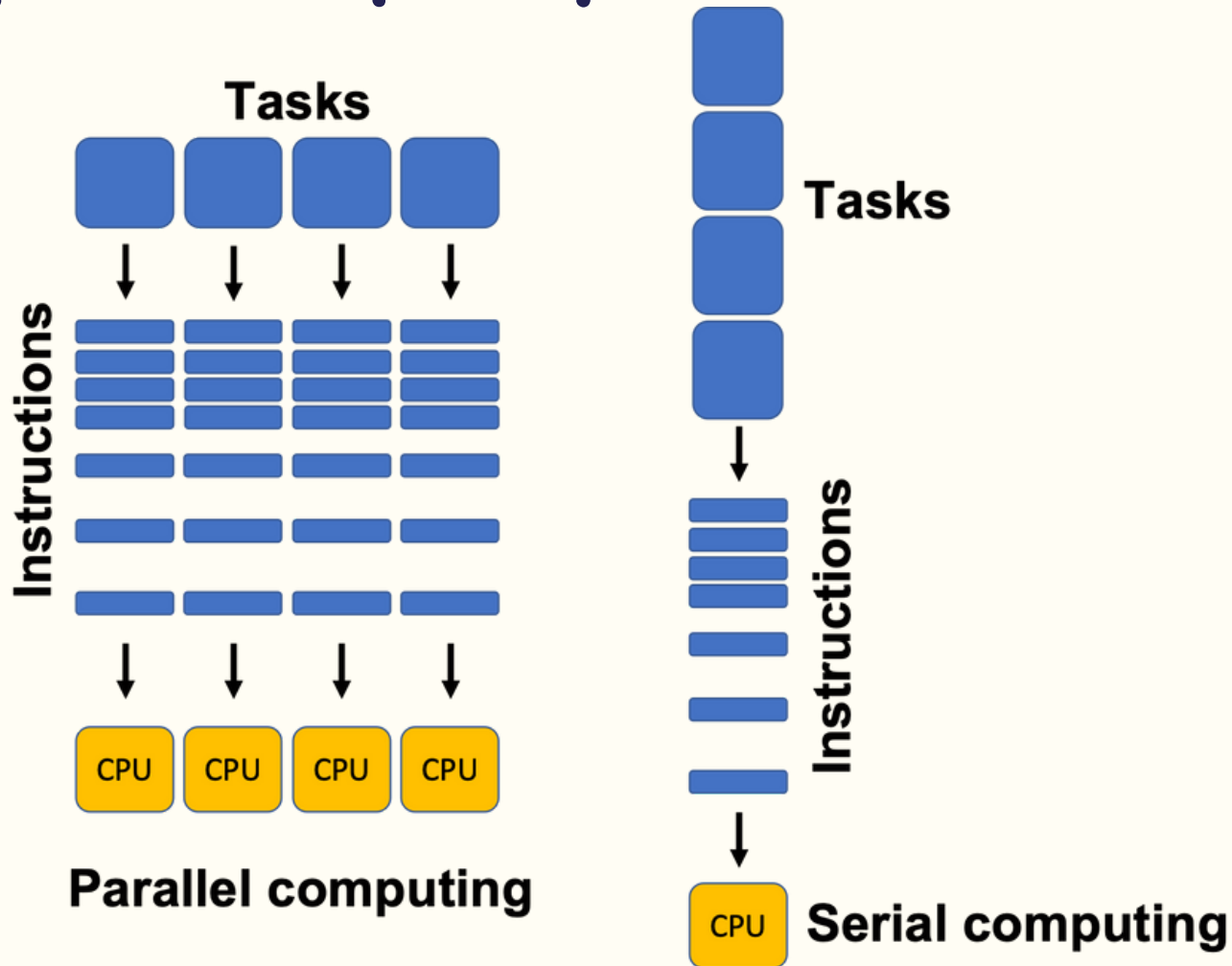
# HOW & WHY PARALLELIZATIONS WORKS



CPU vs GPU Architectures!

# HOW & WHY PARALLELIZATIONS WORKS

Tasks

Instructions

CPU CPU CPU CPU

**Parallel computing**

Tasks

Instructions

CPU **Serial computing**

# HOW & WHY PARALLELIZATIONS WORKS



6

# HOW & WHY PARALLELIZATIONS WORKS

# DATA OR MODEL PARALLELIZATION

## Inter-Layer (Pipeline) Parallelism

- Inference:
    - Maximizes GPU utilization and Throughput
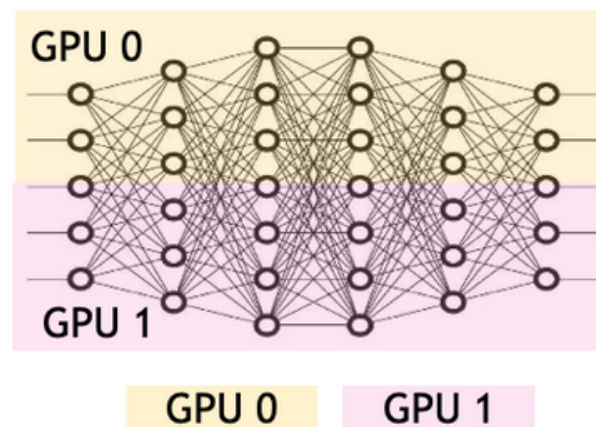    - Can be used easily with TRITON
- Split contiguous sets of layers across multiple GPUs



DGX #1    DGX #2

## Intra-Layer (Tensor) Parallelism

- Split individual layers across multiple devices
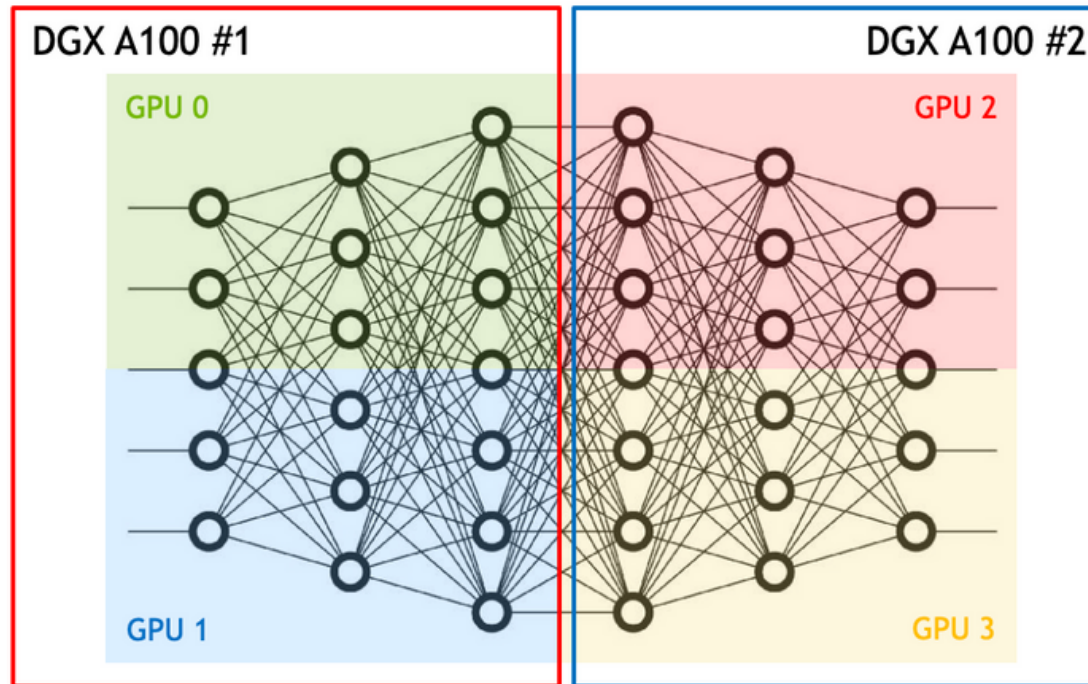- Inference:
    - Minimizes latency



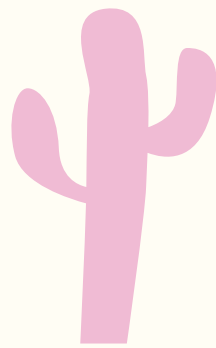GPU 0    GPU 1

# DATA OR MODEL PARALLELIZATION



## MODEL PARALLELISM
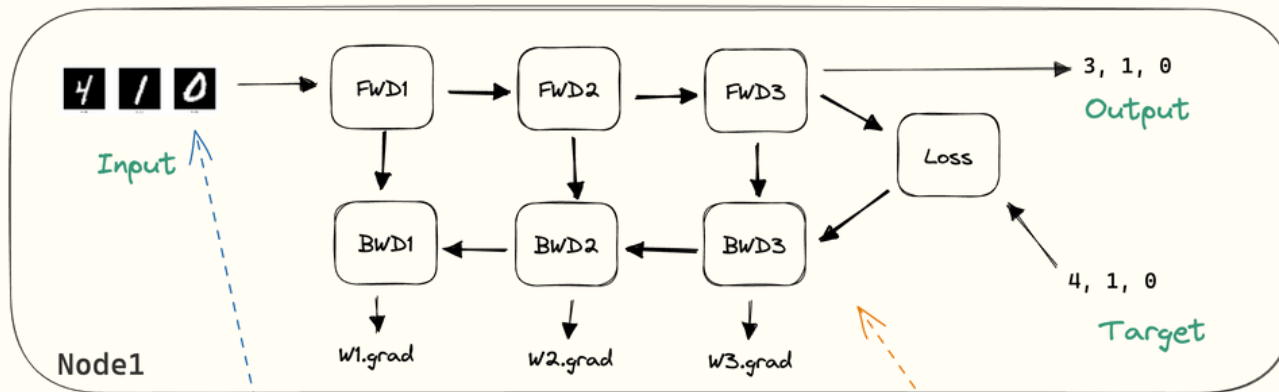Combined Model Parallelism. Multiple GPUs in Multiple DGXs.

DGX A100 #1

GPU 0

GPU 1

DGX A100 #2

GPU 2

GPU 3

Inter + Intra Parallelism

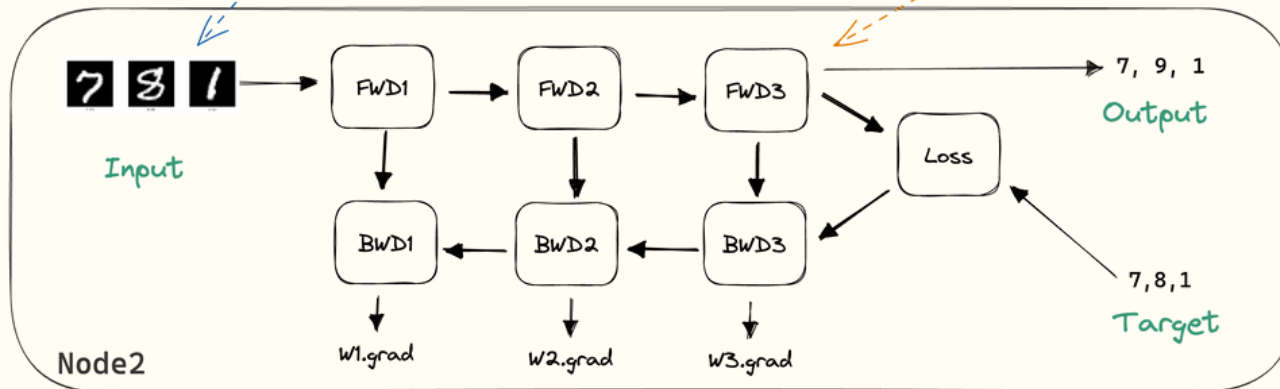# DISTRIBUTED DATA PARALLEL

Data parallel training with 2 compute nodes
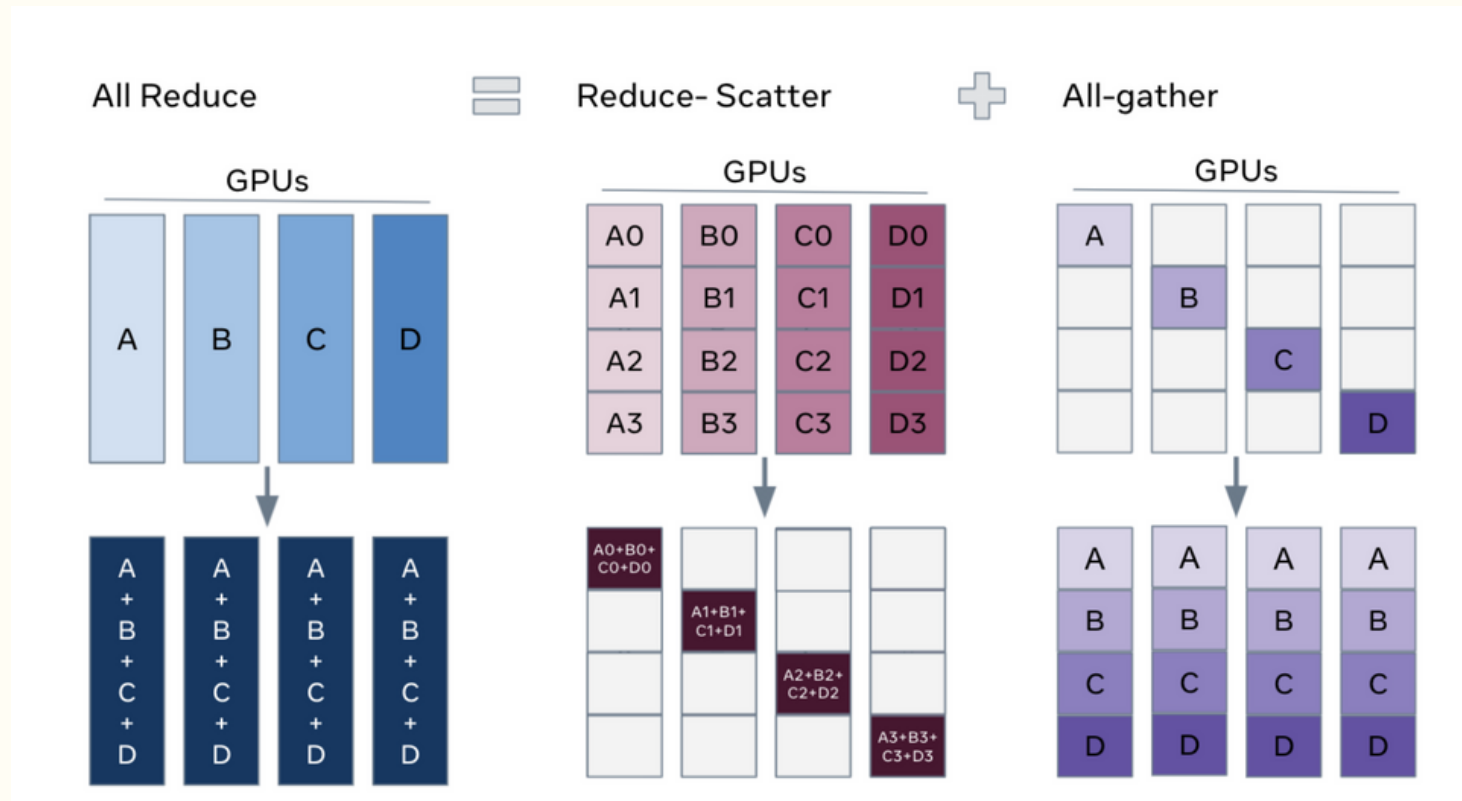


Distributed Sampler

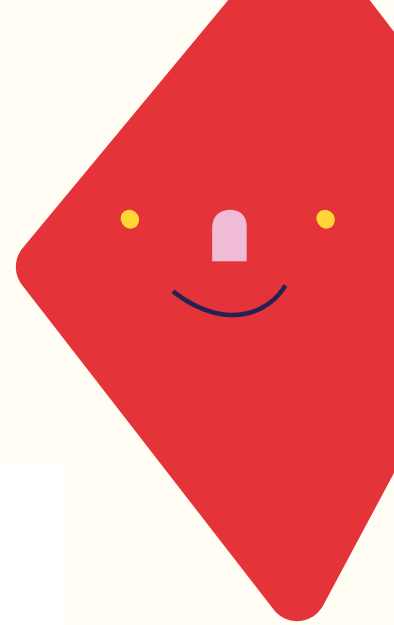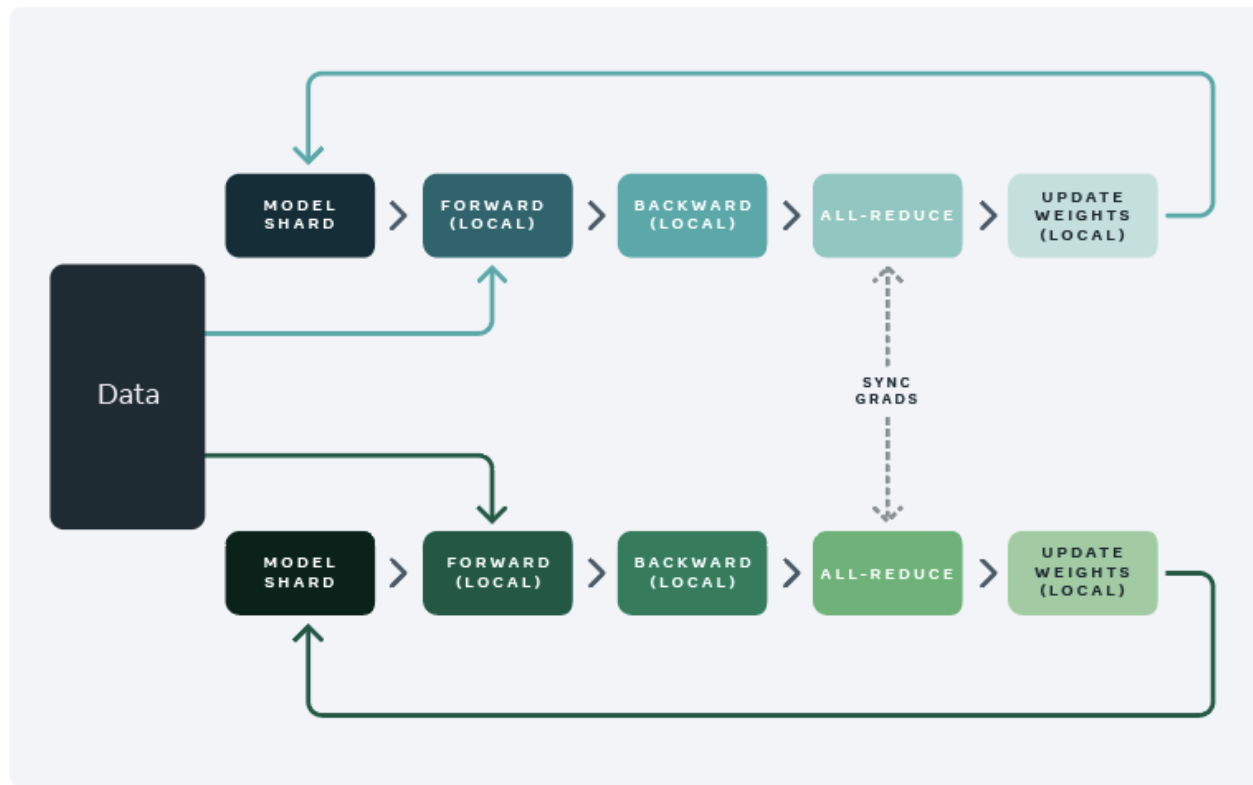# DISTRIBUTED DATA PARALLEL



it takes more GPU memory than it needs because the model weights and optimizer states are replicated across all DDP workers.
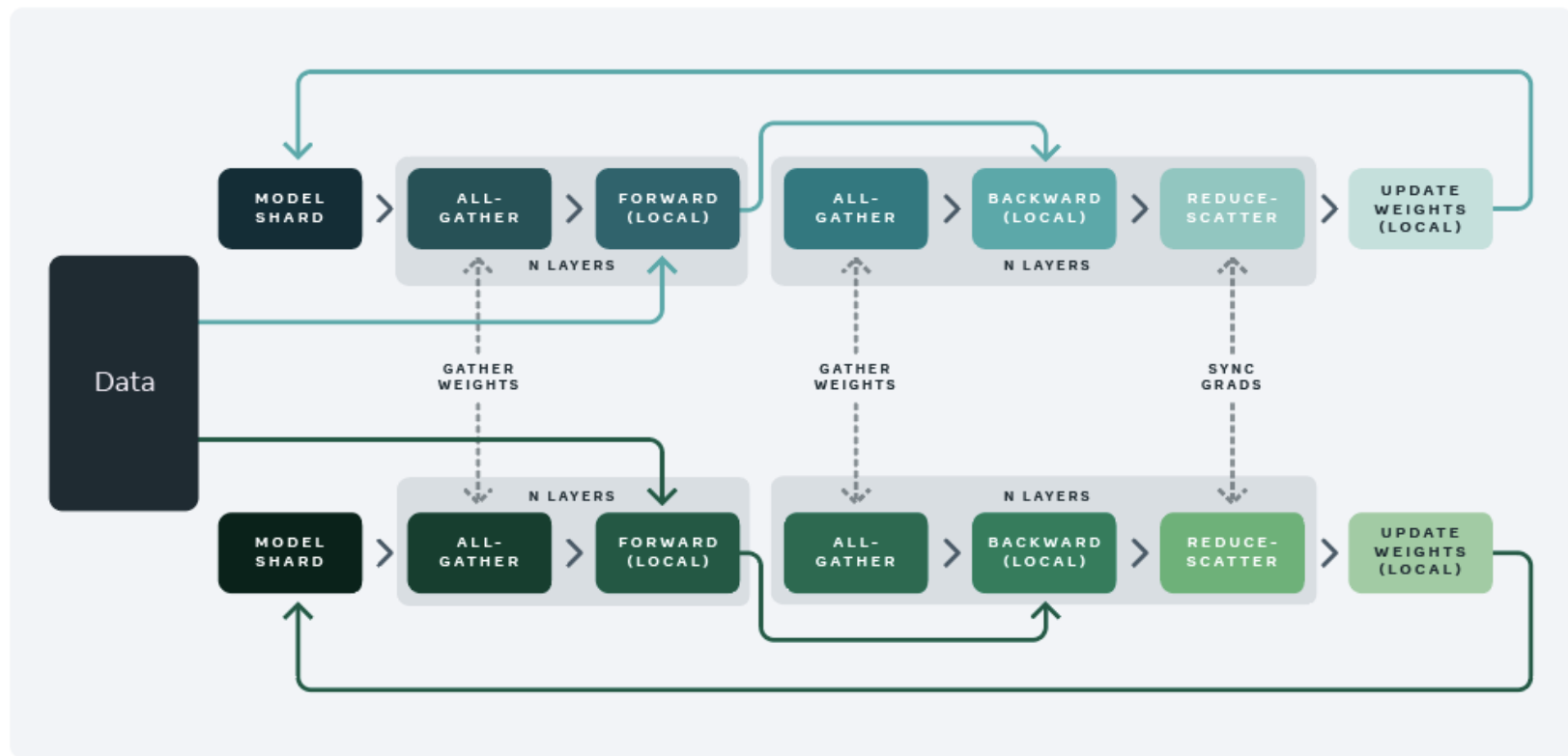
# FULLY SHARDED DATA PARALLEL
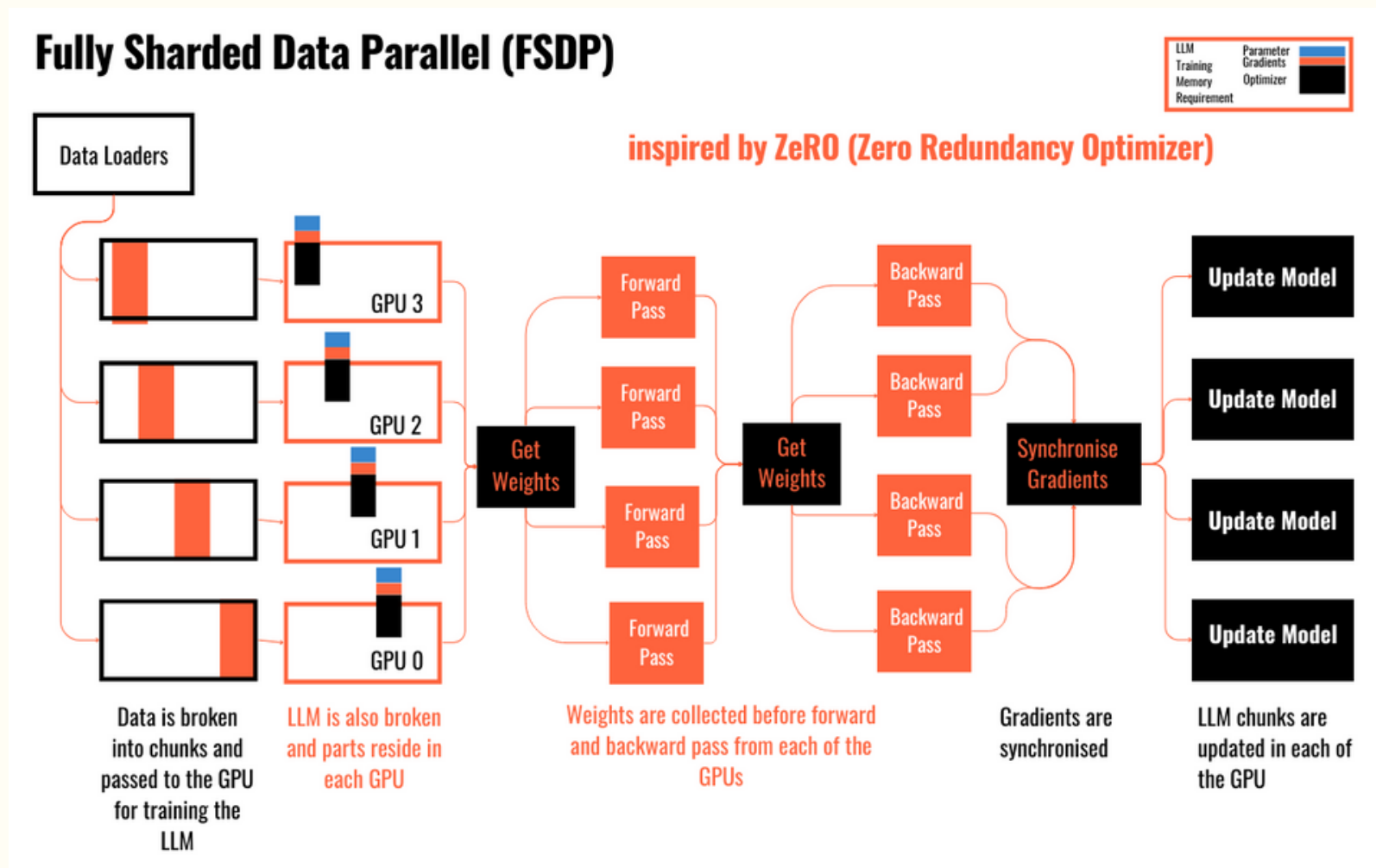


Standard data parallel training
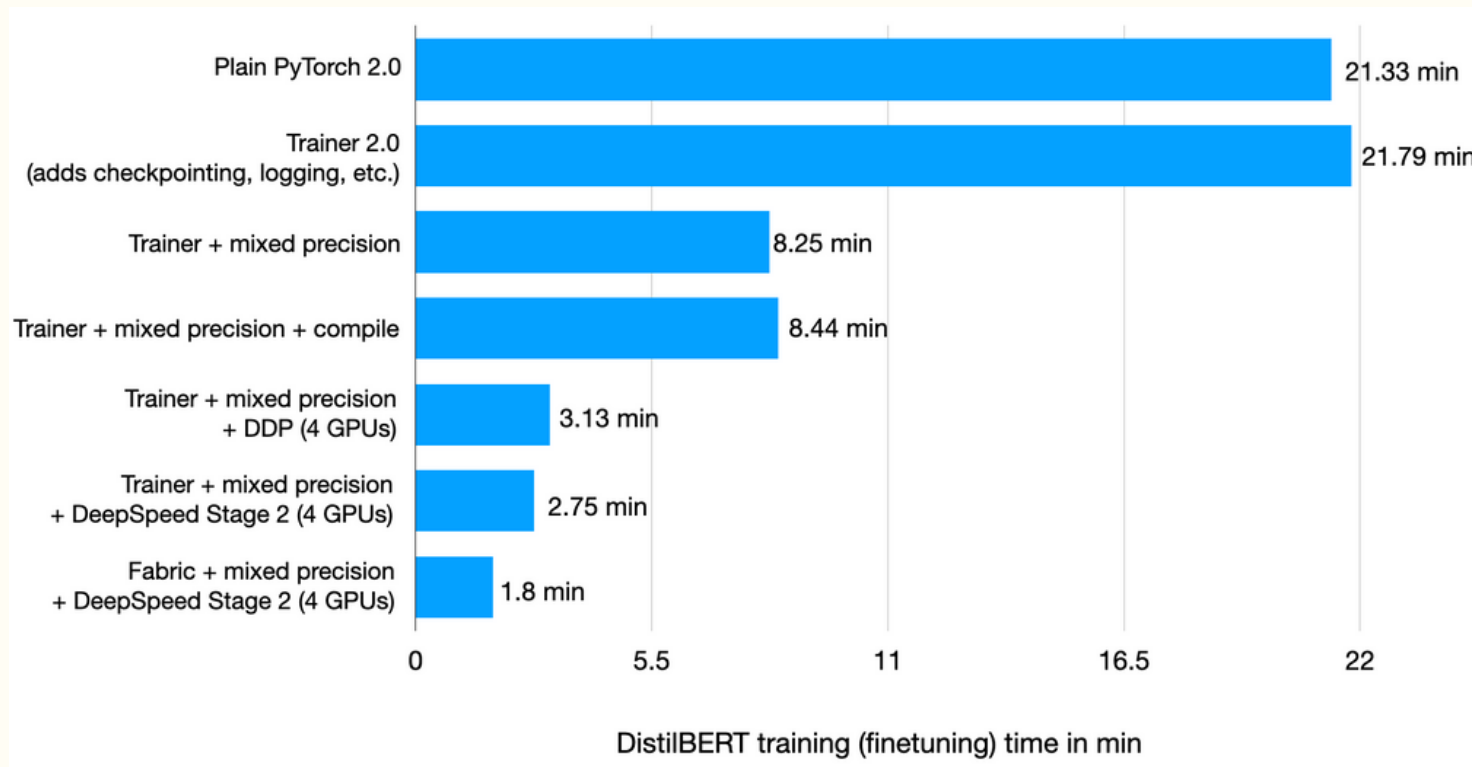
# FULLY SHARDED DATA PARALLEL



Fully sharded data parallel training

# FULLY SHARDED DATA PARALLEL

# EXAMPLE WITH DEEPSPEED + MULTI-GPUS



University of Amsterdam: https://uvadlc-notebooks.readthedocs.io

# RESOURCES

- [Fully Sharded Data Parallel: faster AI training with fewer GPUs](#)

- [Getting Started with Fully Sharded Data Parallel (FSDP)](#)

- [Introducing PyTorch Fully Sharded Data Parallel (FSDP) API](#)

- [Multi GPU Fine tuning with DDP and FSDP](#)

# LET'S GET JAMMIN'

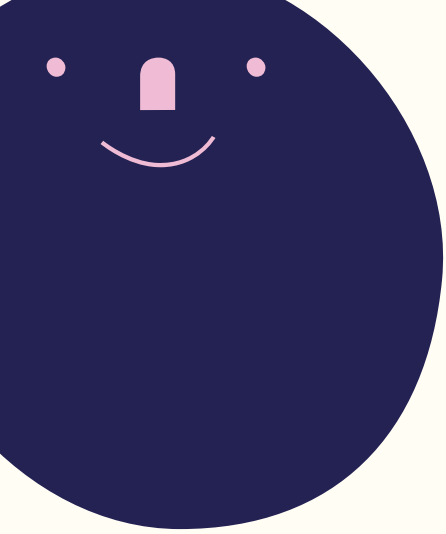## For questions, comments, and feedbacks

**Email**

mert.bozkirr@gmail.com

**Website**

www.mertbozkir.com

**Linkedin**

https://linkedin.com/in/mertbozkir

# THANK YOU

for listening