# Deploying AI at Meta Scale

BUILDING DEVELOPER-FRIENDLY HIGH-PERFORMANCE SCALABLE SYSTEMS FOR AI + PYTORCH

Alexis Bjorlin
VP, Infrastructure Engineering

∞ Meta

# Meta

Community Statistics

## 3.65B

people use at least one of our services **monthly**, approximately

## 2.91B

**monthly** active users on Facebook

## 700M

people use Augmented Reality across our apps and devices monthly

∞ Meta

Snapshot of Data Center Footprint

>100
Data centers
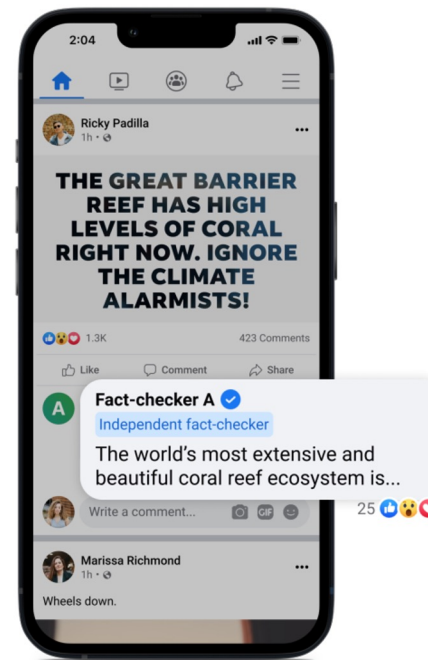
~50M
Square feet

100%
Renewable energy

New Albany, Ohio

How we **use** AI…

# Content Understanding

- Identify and eliminate inappropriate content before being viewed
- Leverages several model types, including computer vision, image classification, and natural language processing



## ~250M
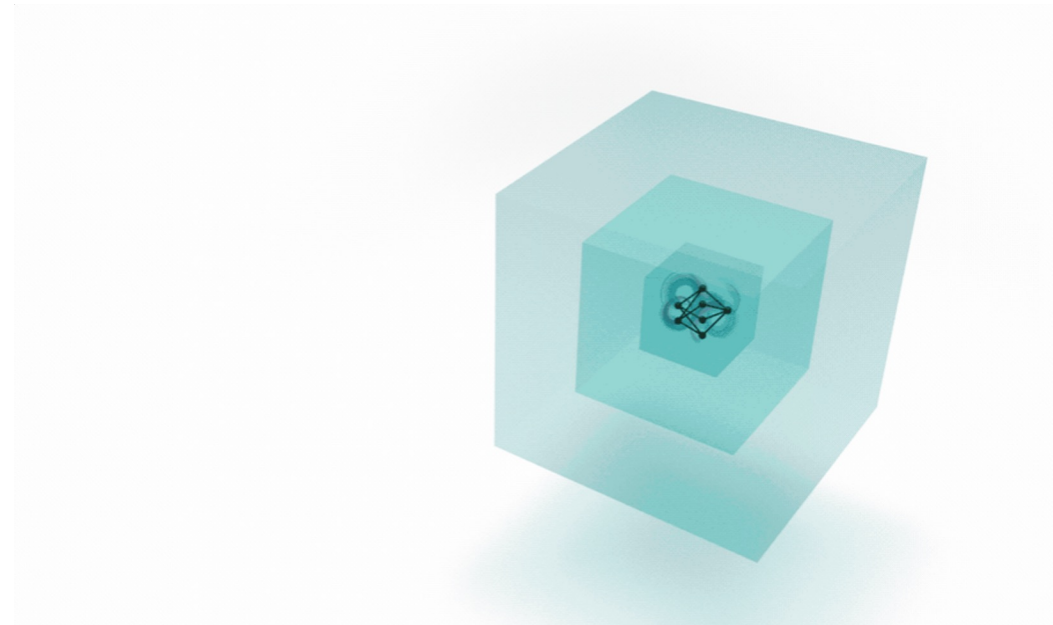Violating Content (Q2)

## 95%+
Actioned on by AI

# Natural Language Processing

- Open Pretrained Transformer (OPT-175B)
- XMLR

Sources:
https://ai.facebook.com/research/no-language-left-behind/
https://ai.facebook.com/blog/democratizing-access-to-large-scale-language-models-with-opt-175b/

# Recommendation & Personalization

- Deep Learning Recommendation Model (DLRM)
- TBSM, DHEN

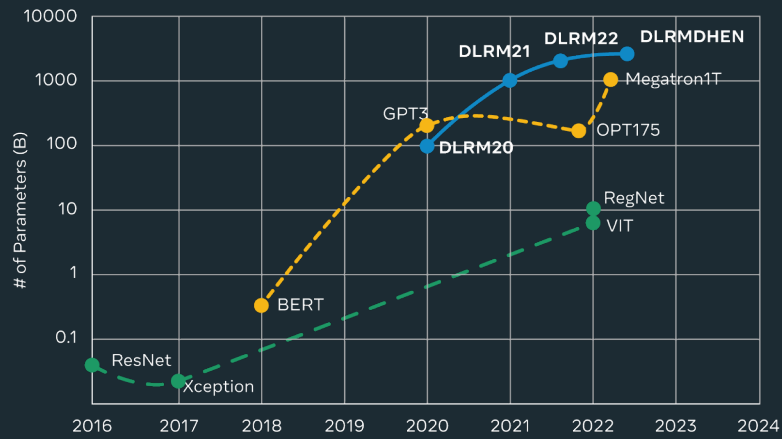Neural Net → Compute intensive

Feature Interaction → Communication intensive

Neural Net | Embedding Lookup | Embedding Lookup → Memory bandwidth intensive

Continuous Features | Categorical Features | Categorical Features → Memory capacity dominated
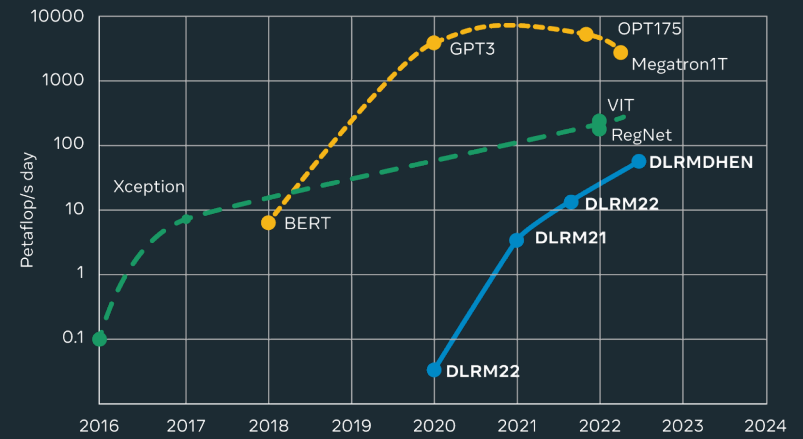
Input from network

DEEP LEARNING WORKLOADS – CHARACTERISTICS

SIZE

COMPUTE

**AI IS POWERING EVERYTHING WE DO: AI-related statistics**



# 6B
Training images



# 20B+
Translations per day



# 200T+
Predictions per day

How we **develop** AI…

# Rapid research to production

- Benefit of deploying state-of-art models fast can be huge

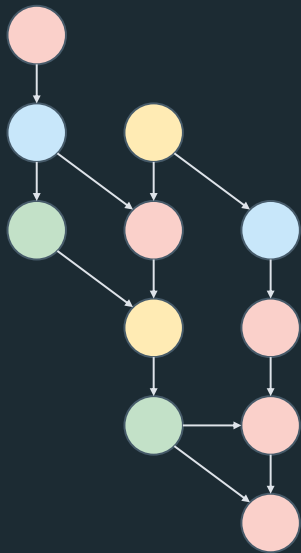Paper Implementations Grouped by Framework

# PYTORCH: A DEVELOPER-FIRST MINDSET

**Full
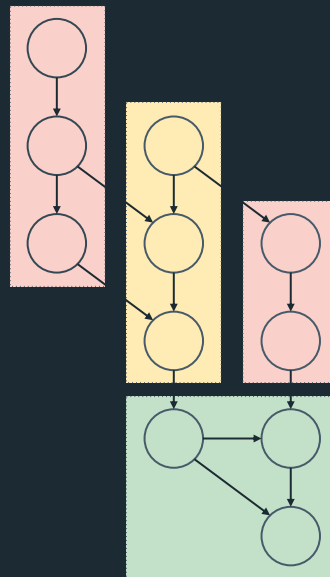Eager-Mode**

```
for epoch in range (max_epochs):
    for model_target, model_input in enumerate(training_data):
        if numpy.random.randint(100) > 90: # 10% noise
            model_target = torch.from_numpy(numpy.random.randint(2))
        model_output = dlrm(model_input)
        model_loss   = torch.nn.BCELoss(model_output, model_target)
        model_loss.backward()
        optimizer.step()
        print("BCE loss " + str(model_loss))
        matplotlib.pyplot.plot(...)
```
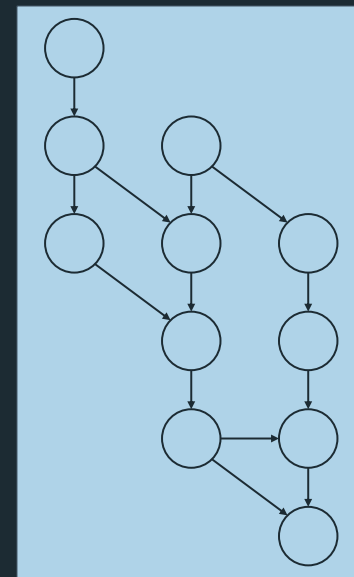
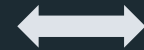DEVELOPER EFFICIENCY VS. PERFORMANCE

Full
Eager-Mode

Partial
Eager-Mode

Graph Mode

Developer
Efficiency

Model
Performance

How we **enable** AI…

**OPTIMIZING THE AI SYSTEM FOR PYTORCH**

# Programmable

Easy operator authoring → new compute primitives for model innovation.

# Dynamic

Fast operator launch.
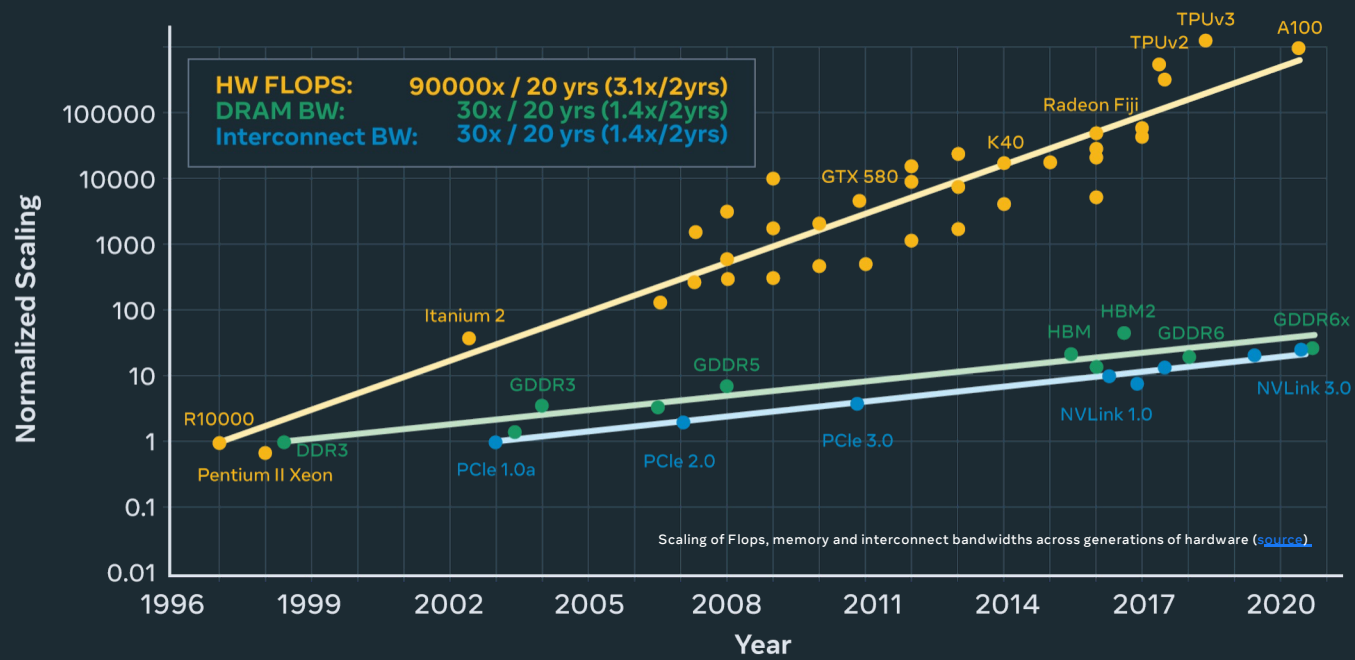Dynamic tensor shapes, memory allocation, easy to prototype

# Scalable & Tunable
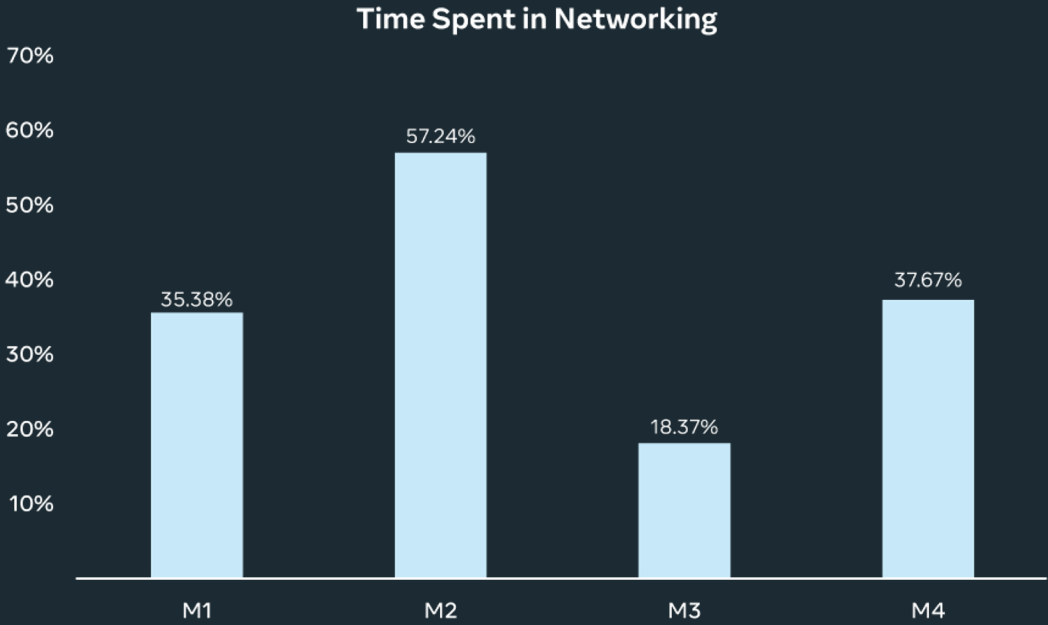
Future-proof.

Balance between compute, memory, network.

NETWORK IS LAGGING ON EVOLUTION CURVE

Scaling of Peak hardware FLOPS, and Memory/Interconnect Bandwidth

HW FLOPS: 90000x / 20 yrs (3.1x/2yrs)
DRAM BW: 30x / 20 yrs (1.4x/2yrs)
Interconnect BW: 30x / 20 yrs (1.4x/2yrs)

Scaling of Flops, memory and interconnect bandwidths across generations of hardware (source).

- Storage Cache
- Data ingestion
- Compute Nodes
- Compute fabric

## PyTorch AI Training Cluster

Storage/Cache

Preprocessing

HPC Compute (PyTorch)

Compute Network

Data Network

**The computational core of the cluster**

THE PYTORCH AI TRAINING CLUSTER

THE PYTORCH AI TRAINING CLUSTER OF THE FUTURE (2025+)
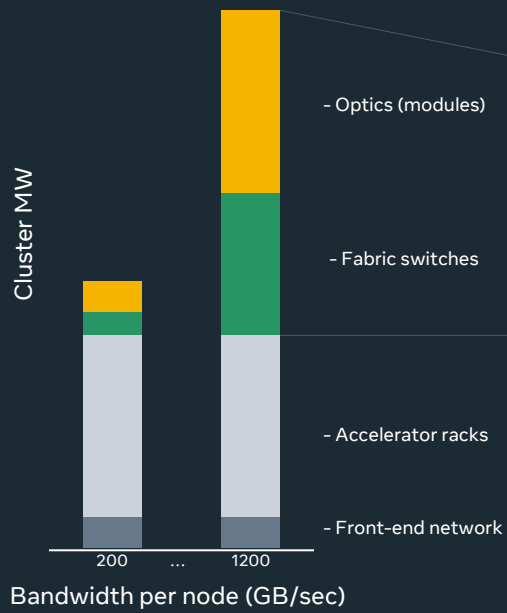
Compute Network

Data Network

~4K
Accelerators

~1TB/s
of compute
network per accelerator

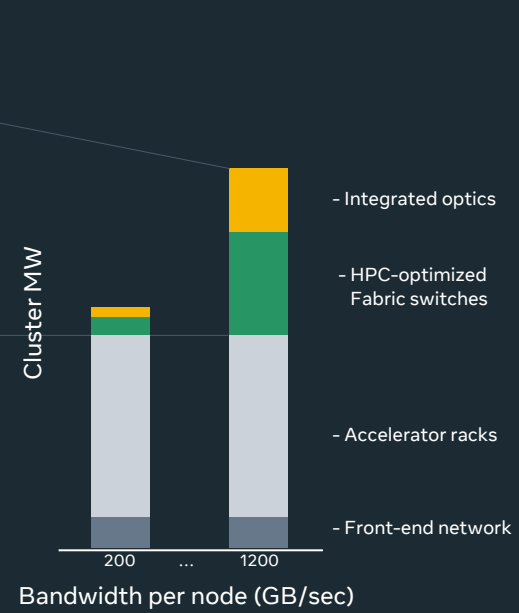The cluster fabric is **Non-blocking**

THE PROMISE OF OPTICAL IO

**Power breakdown for a cluster (optical modules)**

**Power breakdown for a cluster (CPO)**

Cluster MW

- Optics (modules)

- Fabric switches

- Accelerator racks

- Front-end network

200 ... 1200
Bandwidth per node (GB/sec)

Cluster MW

- Integrated optics

- HPC-optimized Fabric switches

- Accelerator racks

- Front-end network

200 ... 1200
Bandwidth per node (GB/sec)

# THE BENEFIT OF A PYTORCH AI TRAINING CLUSTER

## DLRM - DHEN



(Chart — Y-axis: Relative HFU, values 1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.2; X-axis: Relative Model Complexity, values 1.0, 2.4, 5.0, 7.7)

## Improving Content Relevance



(Chart — Y-axis: % Improvement in Accuracy, values 0%, 5%, 10%, 15%, 20%, 25%; X-axis: Relative Overall Complexity, values 1, 10x, 100x)

## +24%

improvement in content relevance over 2021-22

Co-design
for PyTorch

Flexibly balance compute,
memory, and network

Plan for
Innovation as
AI evolves