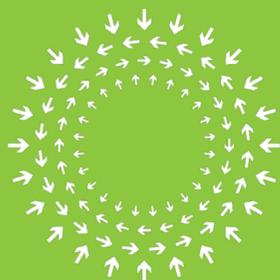# OCP SUMMIT

March 20-21
2018
San Jose, CA

**OPEN**
Compute Project

# AI HARDWARE INFRASTRUCTURE AT FACEBOOK

Xiaodong Wang
Kevin Lee

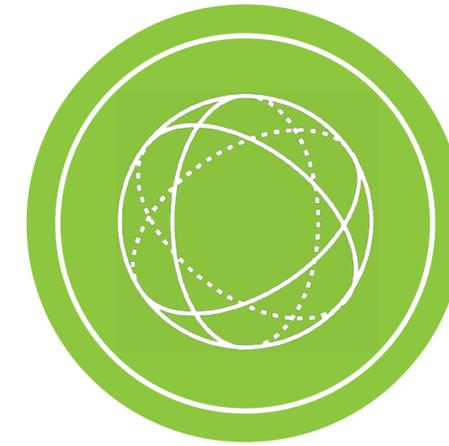Today's AI Hardware Infrastructure

New Hardware Announcement

What's Next
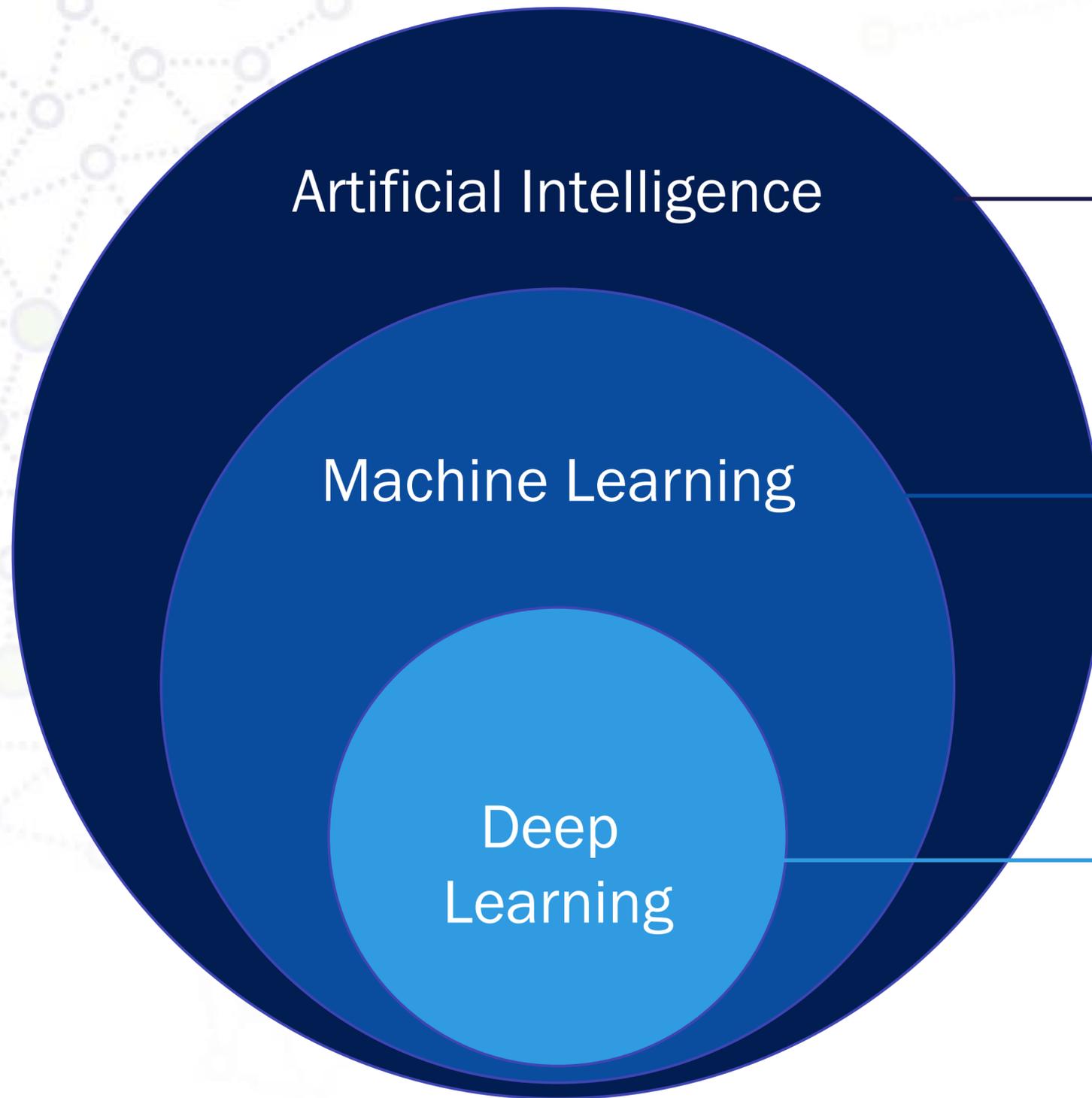
Applied Machine Learning

AI Research

AI Infrastructure

OPEN
Compute Project

Facebook
Open Source

Artificial Intelligence — Program that can sense, reason, act, and adapt like humans

Machine Learning — Program that can perform actions without explicitly being programmed

Deep Learning — Subset of machine learning in which multilayered neural networks learn
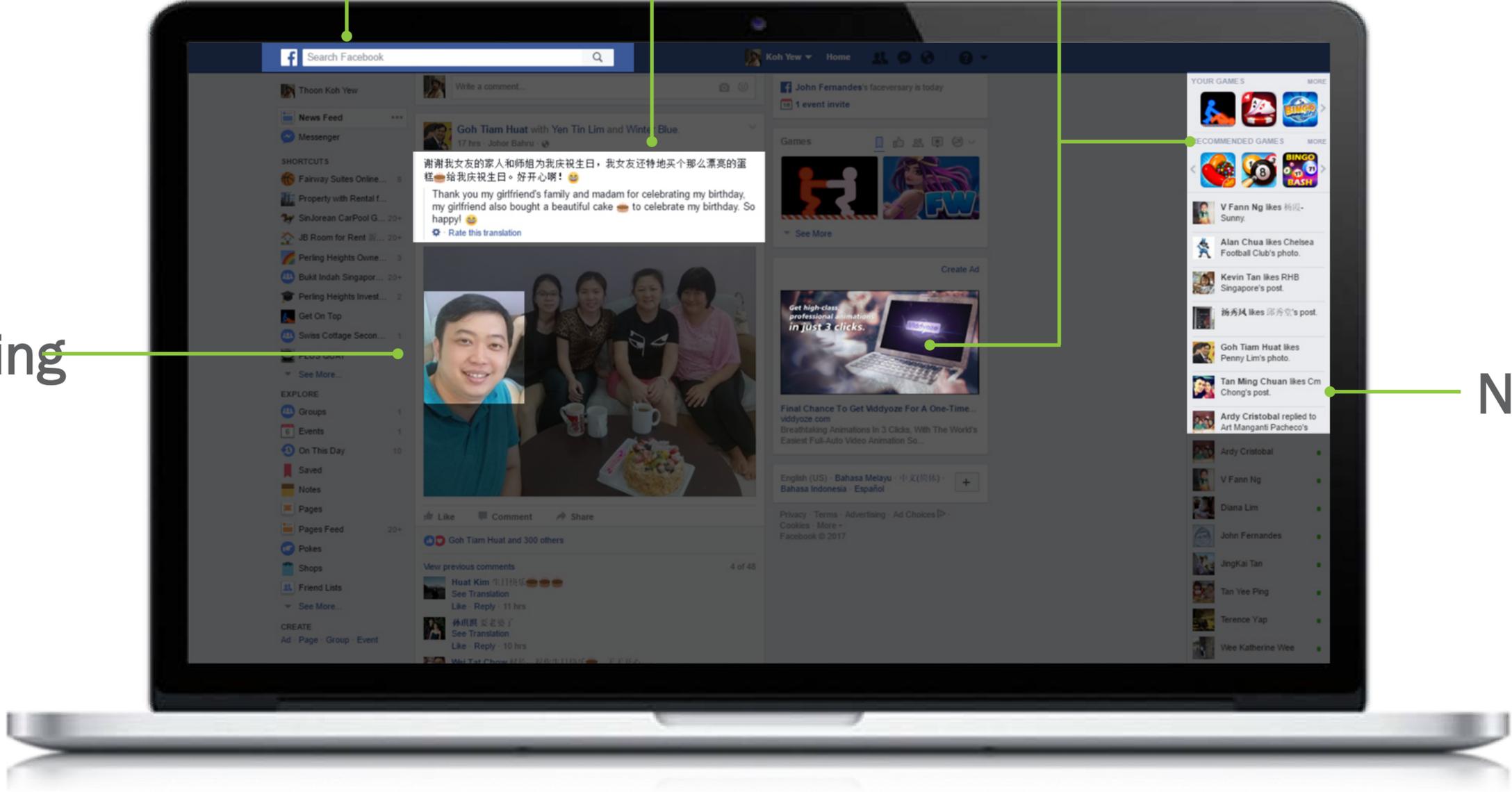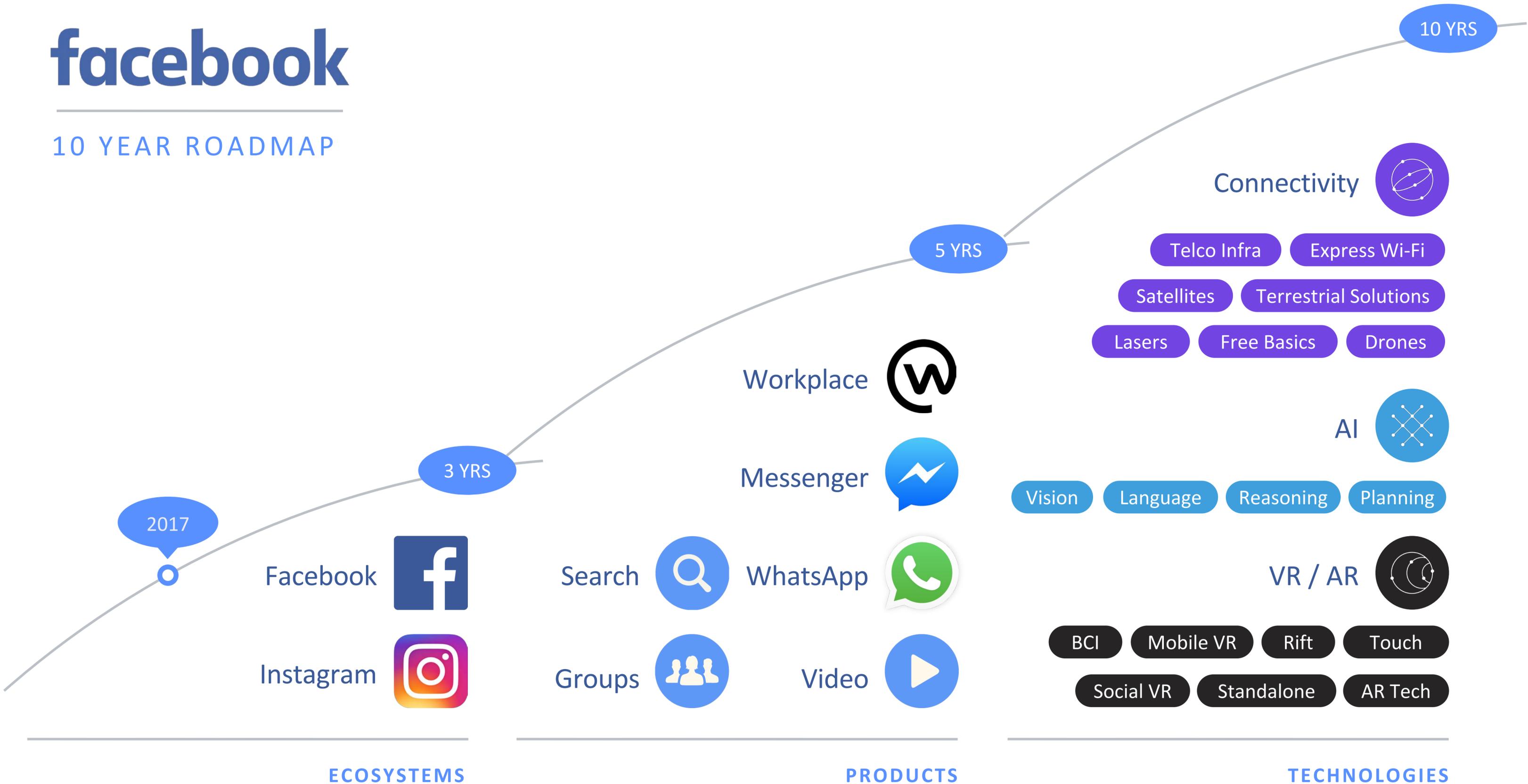
# Services and Use Cases

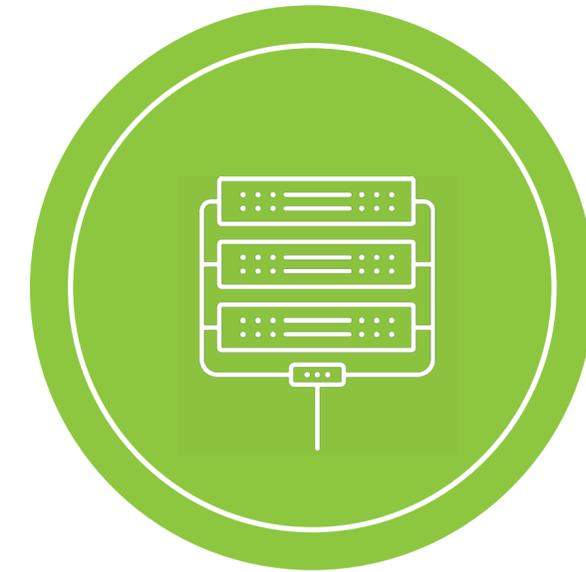Search    Translation    Ads

Face tagging

News Feed

# Why Now?

Research

Data

Compute

| Internal Platforms & Toolkits | Frameworks & Infrastructure |
|---|---|
| FBLearner Feature Store | Caffe2 |
| FBLearner Flow | PYTORCH |
| FBLearner Predictor | ONNX |

# How Are We Leveraging Machine Learning at FB?

Bryce Canyon

Big Basin

Tioga Pass

Twin Lakes



| Data | Features | Training | Evaluation | Inference |
|------|----------|----------|------------|-----------|

FBLearner
Feature Store

FBLearner
Flow

FBLearner
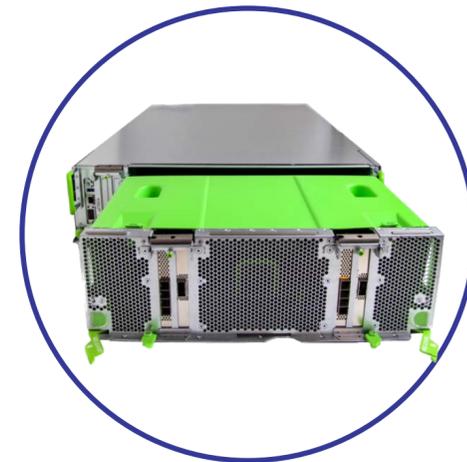Predictor

2012
Experimentation

2013
HP SL270s G8

2015
Big Sur

2016
Big Basin Pascal

2018
Announcing Today

Today's AI Hardware
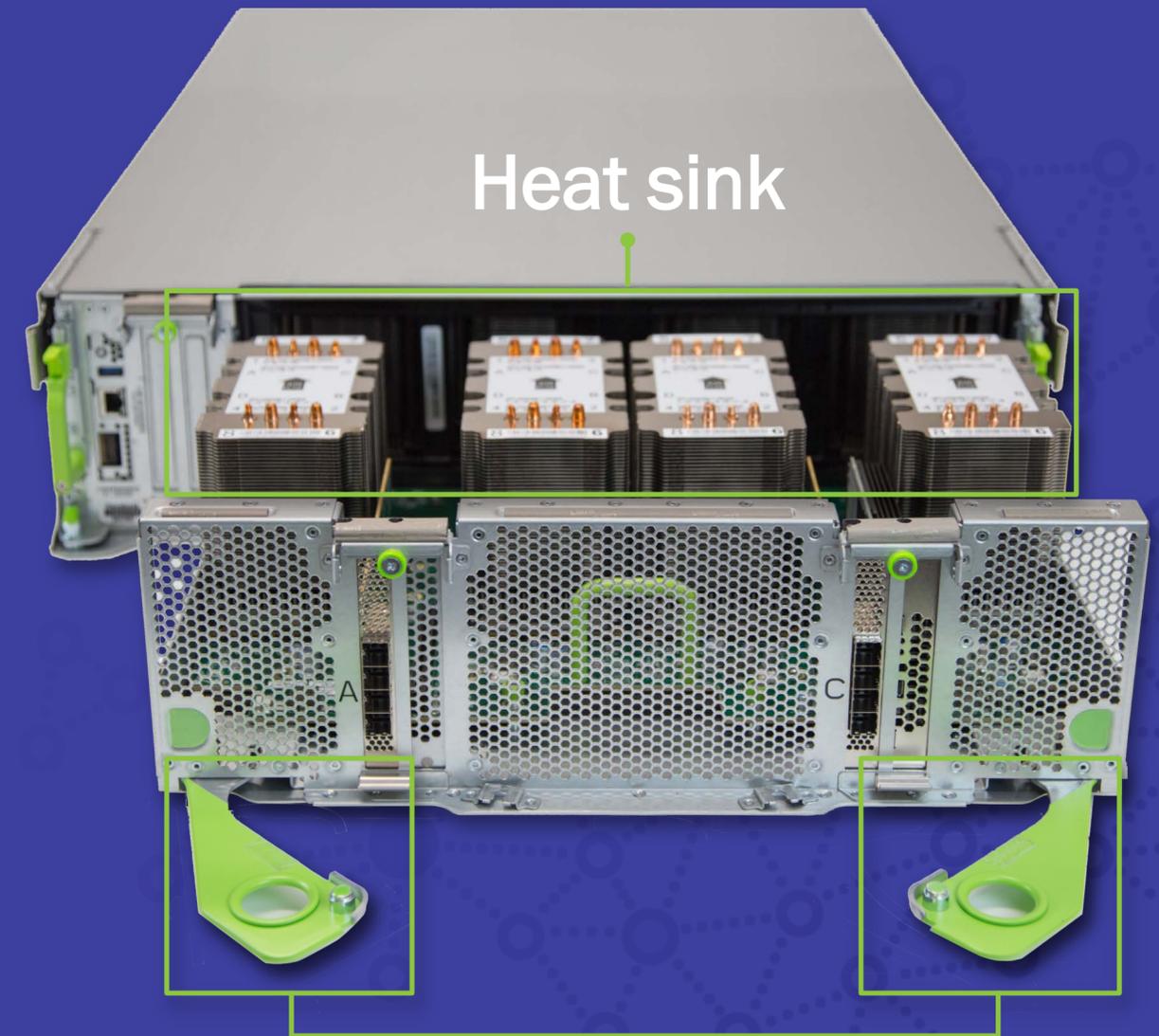Infrastructure

Announcing
Big Basin Volta

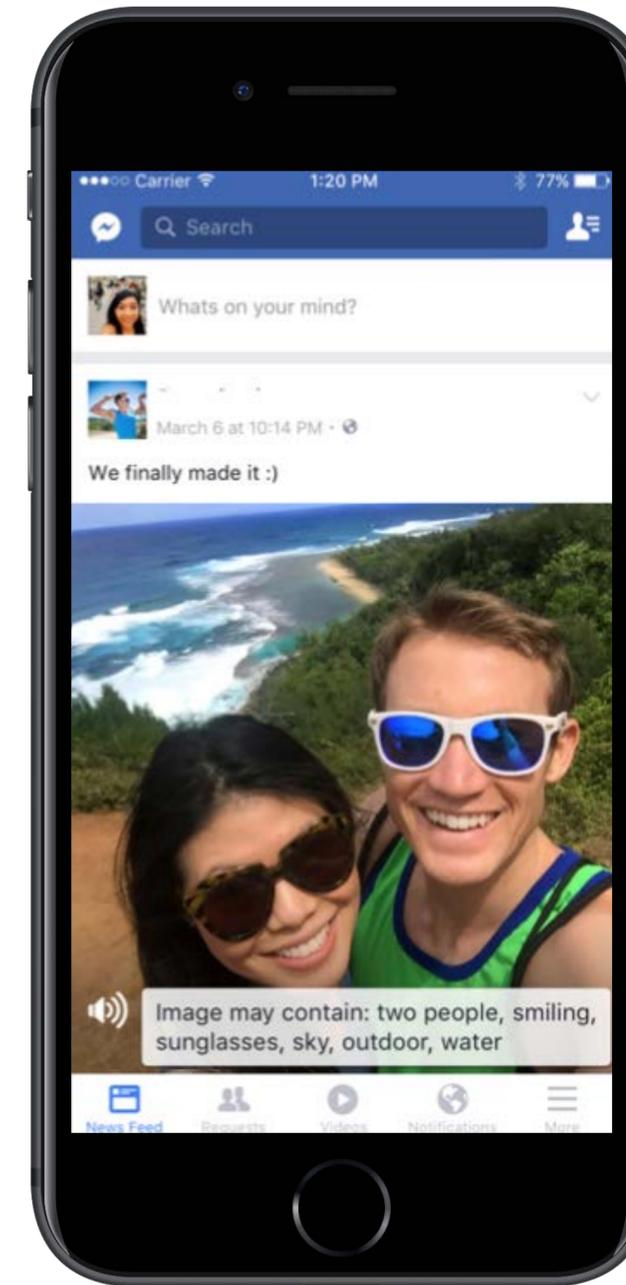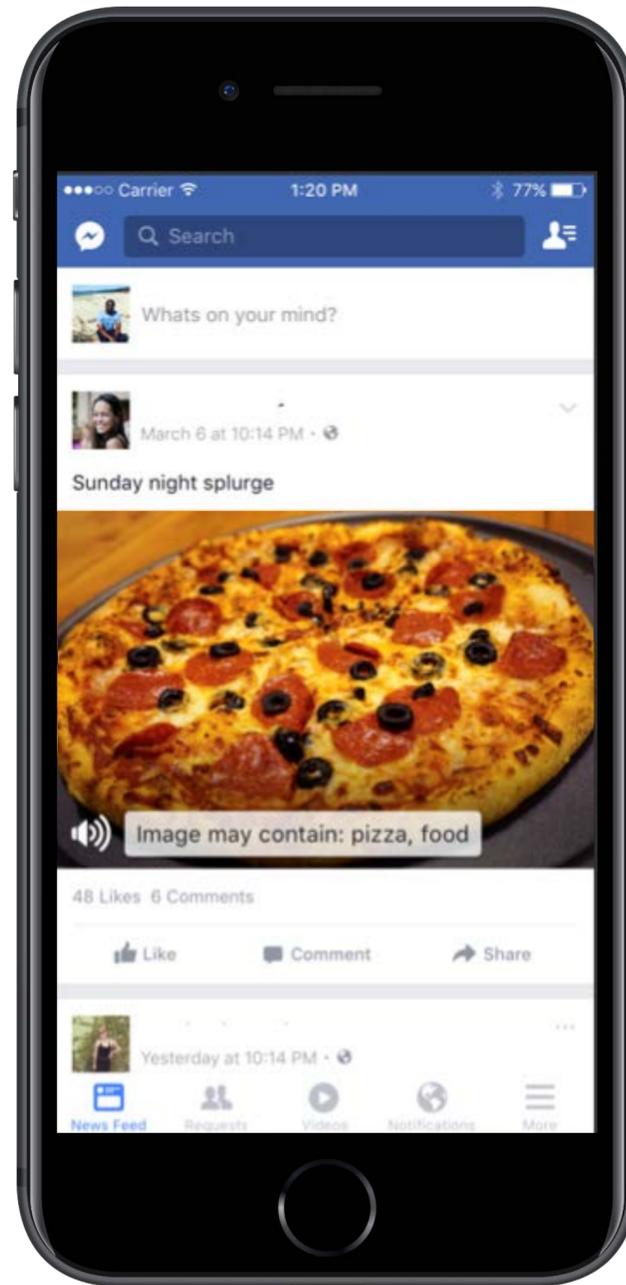What's Next

Big Basin Volta at OCP 2018

# Big Basin Volta at OCP 2018

- Paired with OCP Tioga Pass platform
- Increased GPU performance
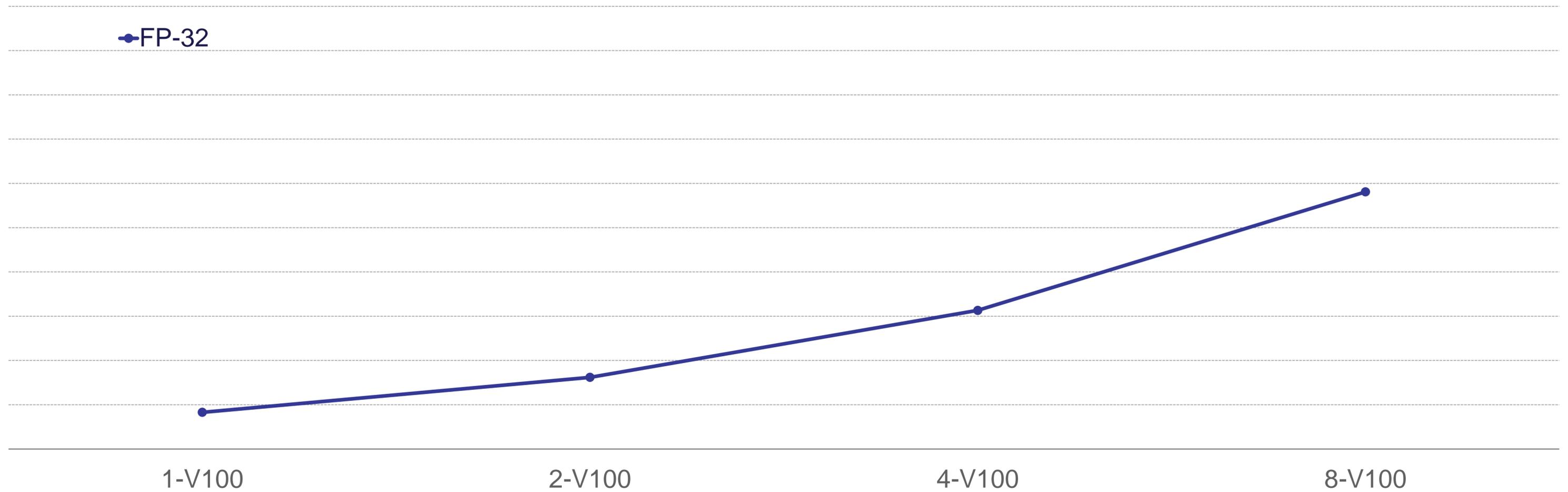- Improved airflow and serviceability
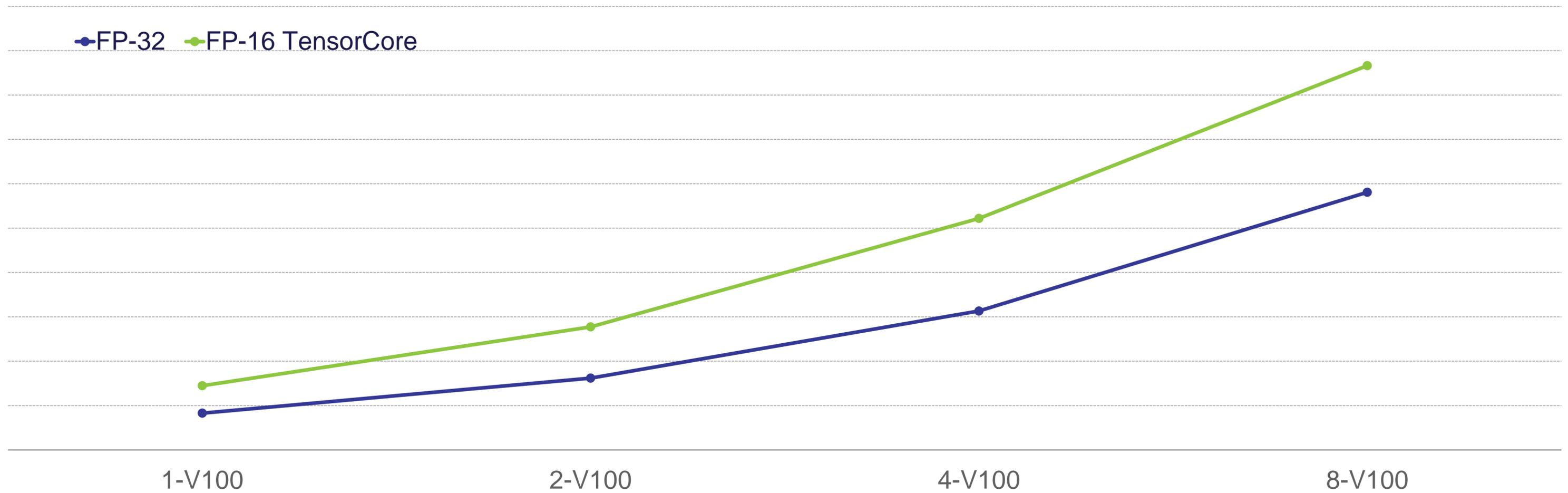


Heat sink

Enlarged handles

# Computer Vision at FB

# Computer Vision Performance

## Multi-GPU speedup



FP-32

1-V100          2-V100          4-V100          8-V100

# Computer Vision Performance

## High-bandwidth FP-16 TensorCore



Legend: FP-32, FP-16 TensorCore

X-axis: 1-V100, 2-V100, 4-V100, 8-V100

# Machine Translation

## Better Translation Quality



Phrase-based statistical approach

Neural network approach

Today's AI Hardware Infrastructure

Announcing Big Basin Volta

What's Next

Consolidated Hardware Design

Max Operating Efficiency

Disaggregated Design

Vendor Agnostic Platform

# THIS JOURNEY 1% FINISHED

Please join us Wednesday at 11:30am for Big Basin Volta workshop!