

An abstract graphic on the left side of the image, composed of numerous thin, wavy green lines that swirl and overlap to form a complex, organic shape. The lines are a vibrant green color against the dark blue background.

Open. Together.



OCP
SUMMIT

AI at Facebook: An Infrastructure Perspective

Whitney Zhao, Hardware Engineer, Facebook

Sam Naghshineh, Technical Program Manager, Facebook

How does FB use Machine Learning?

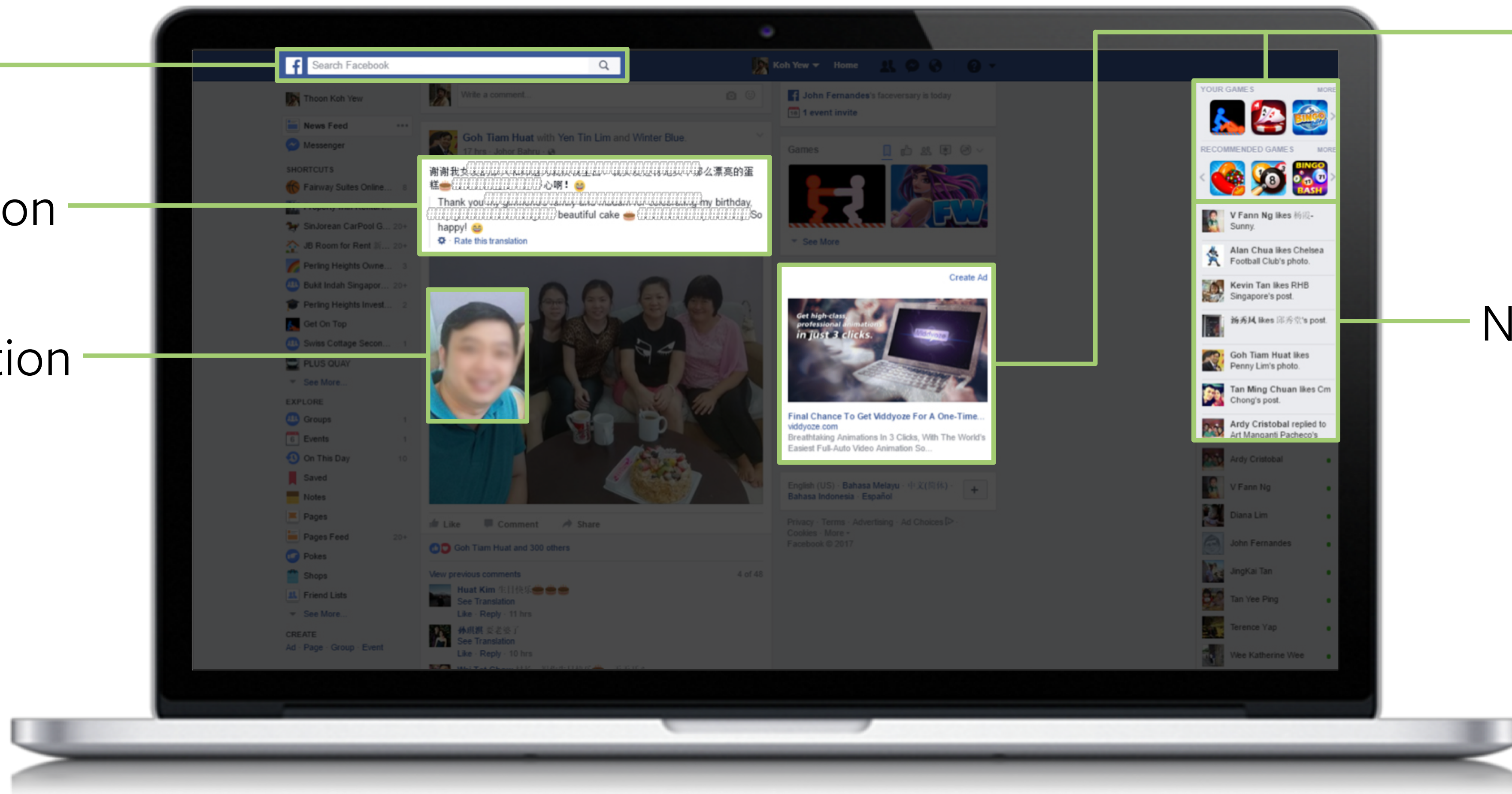
Search

Translation

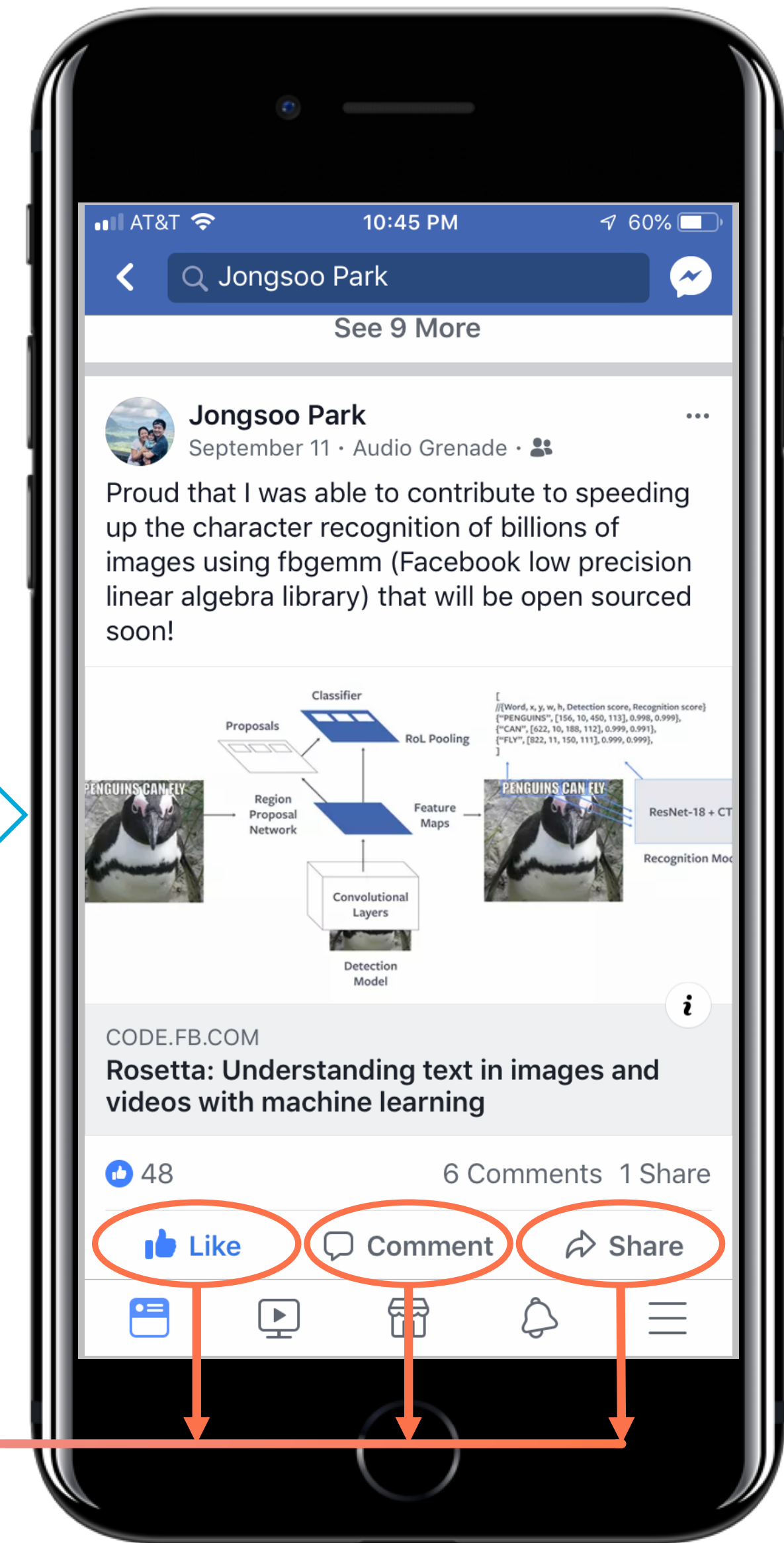
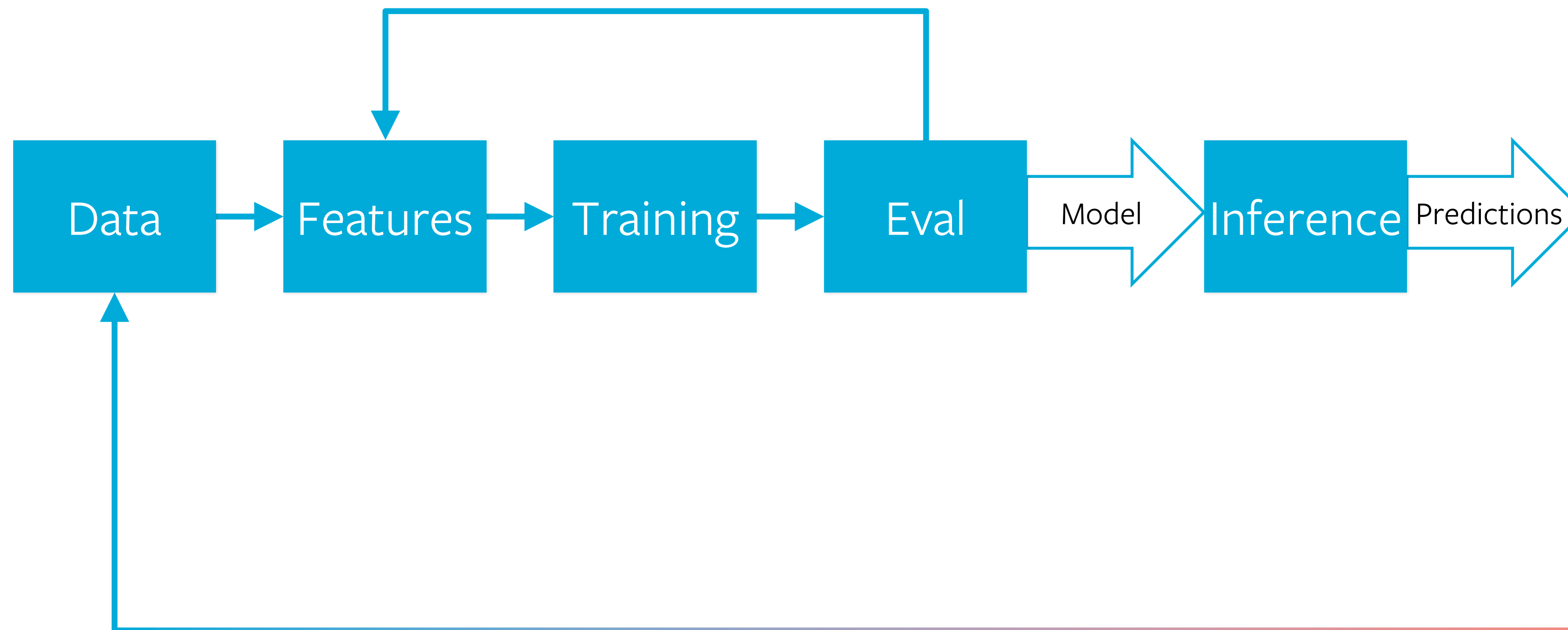
Recognition

Ads

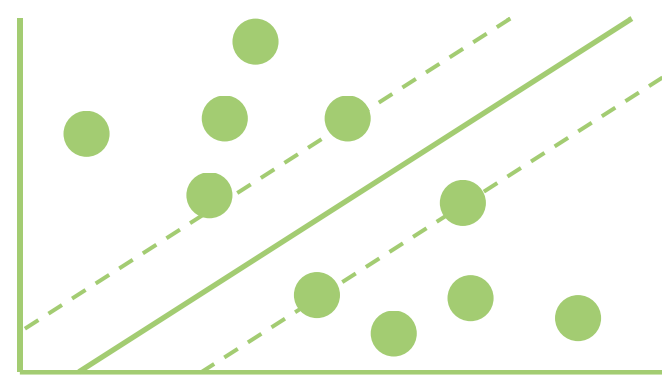
News Feed



Machine Learning Execution Flow



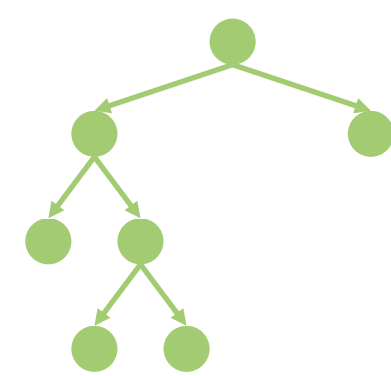
Major AI Services and Algorithms



SVM

Support Vector
Machines

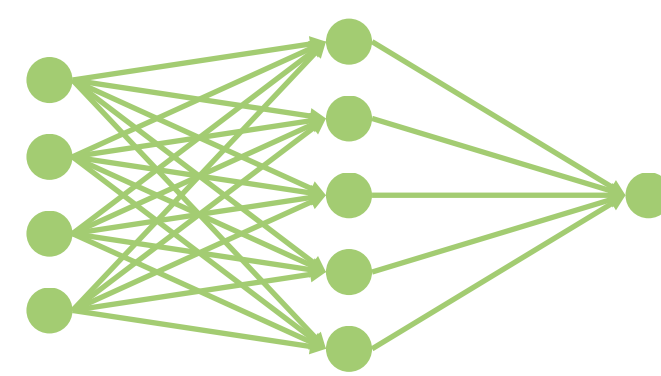
Facer



GBDT

Gradient-Boosted
Decision Trees

Sigma



MLP

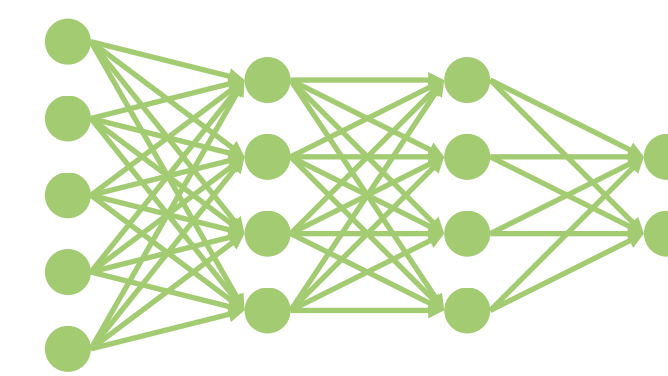
Multi-Layer
Perceptron

News Feed

Ads

Search

Sigma

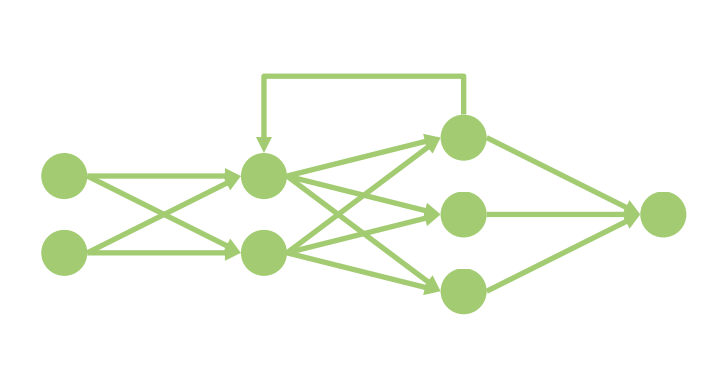


CNN

Convolutional
Neural Nets

Facer

Lumos



RNN

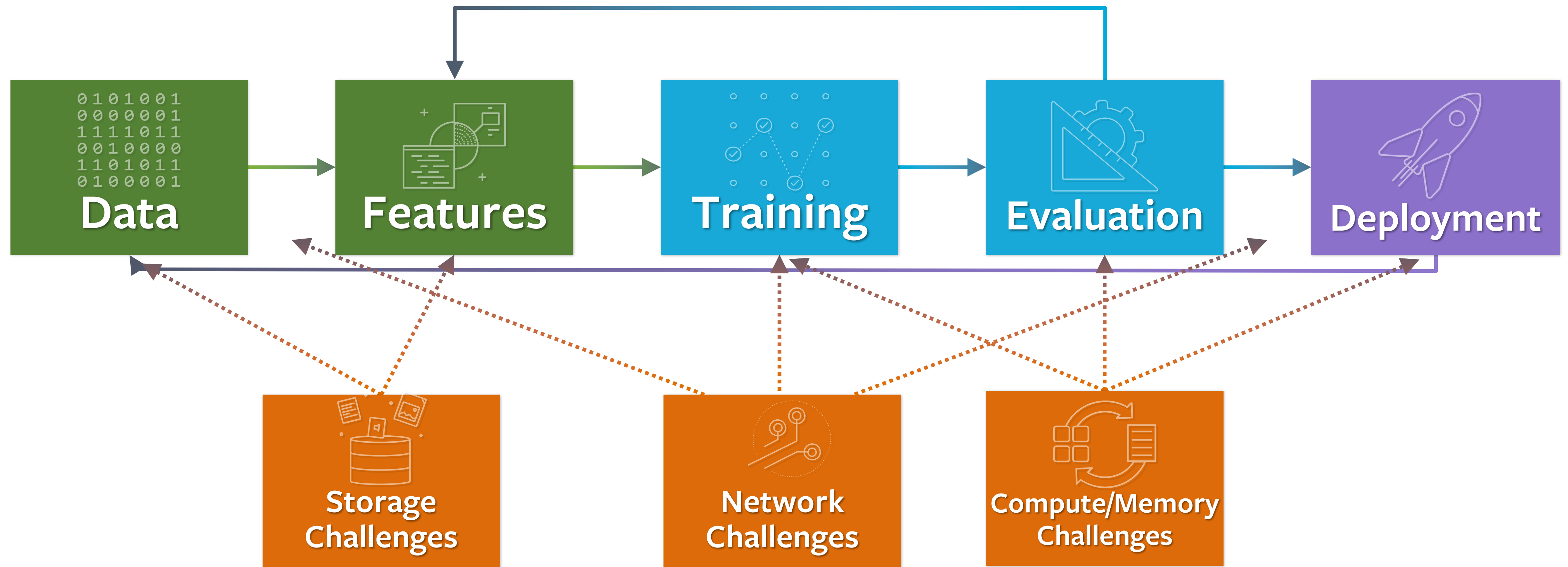
Recurrent
Neural Nets

Language
Translation

Speech
Recognition

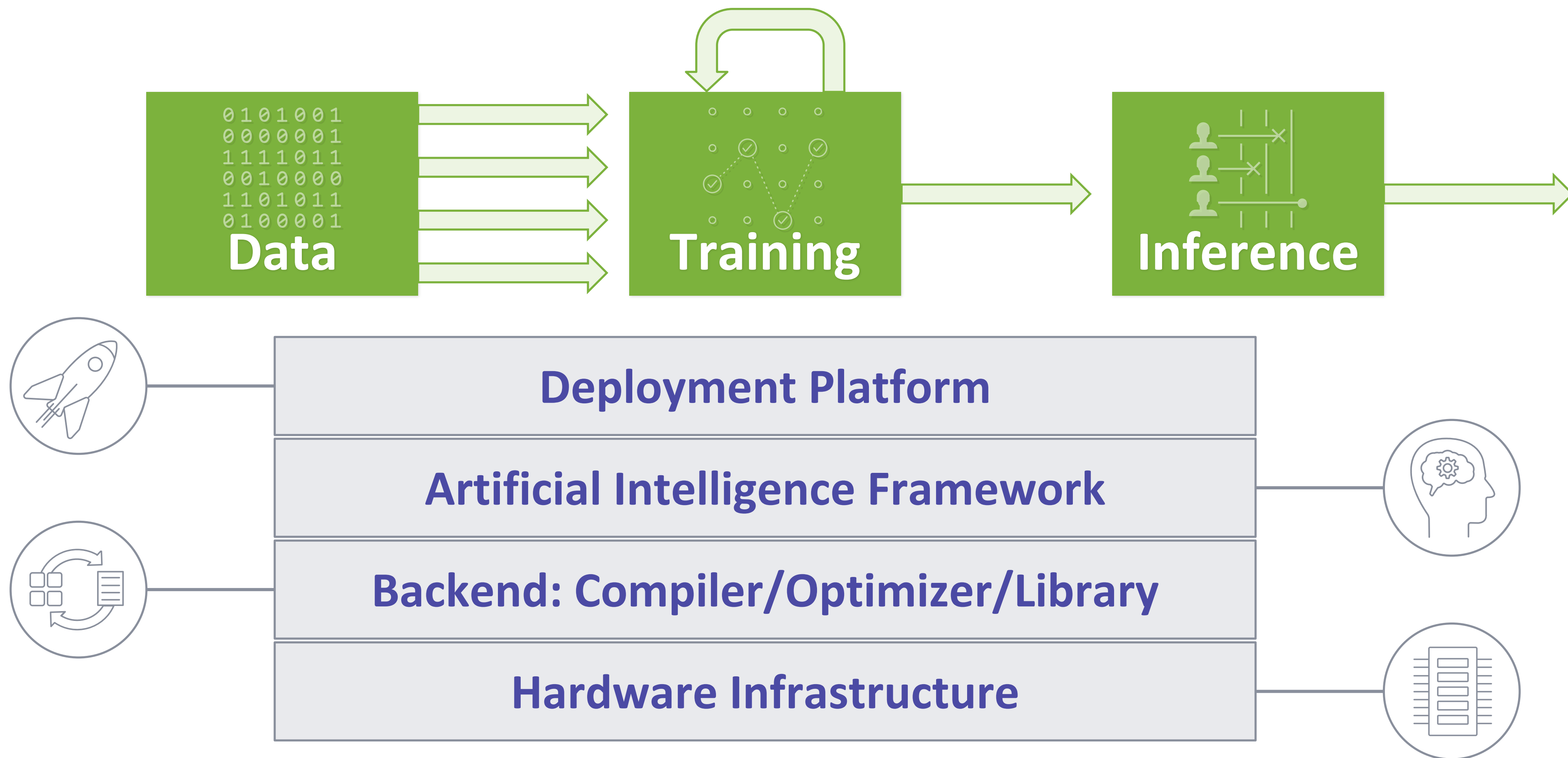
Content
Understanding

It's an Infrastructure Challenge!

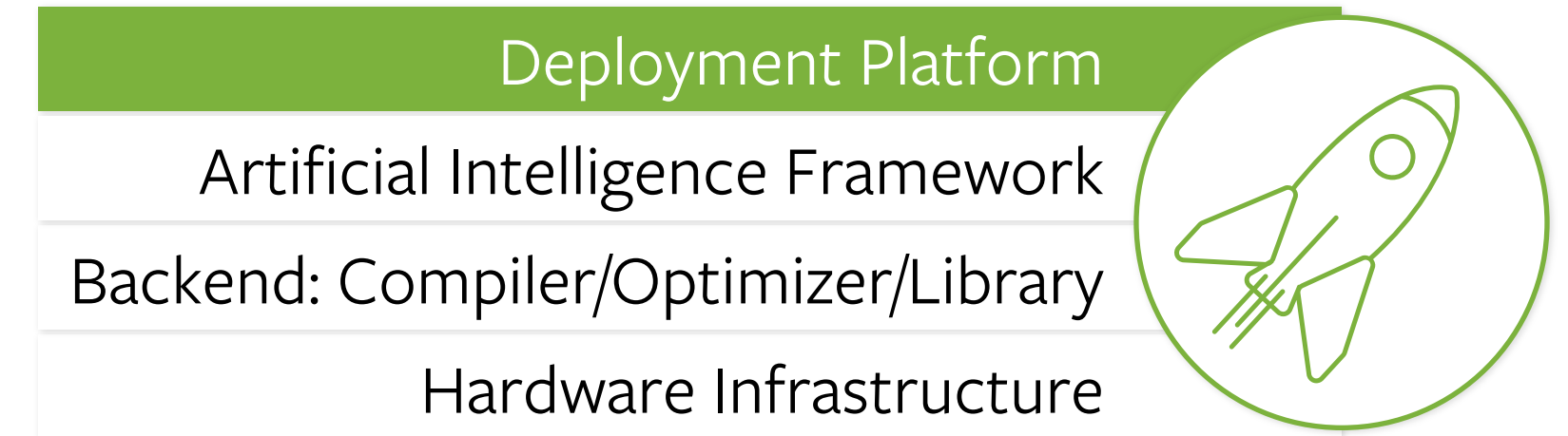


Accelerators could shift the system bottleneck from compute intensity to memory/network

FB AI Development Ecosystem

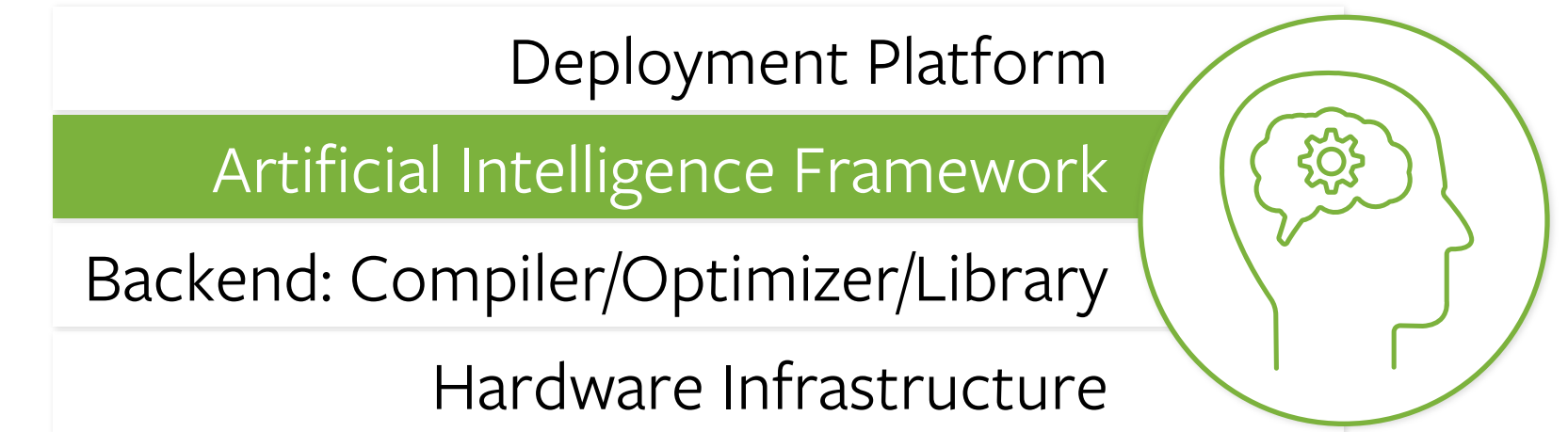


FBLearner Platform



- AI workflow for automated model management and deployment
- Programmable, reusable, distributed training pipeline
- Library of reusable ML algorithms
- Provide sharable and reusable history of past experiments

Framework



PROTOTYPING

DEPLOYMENT

From then

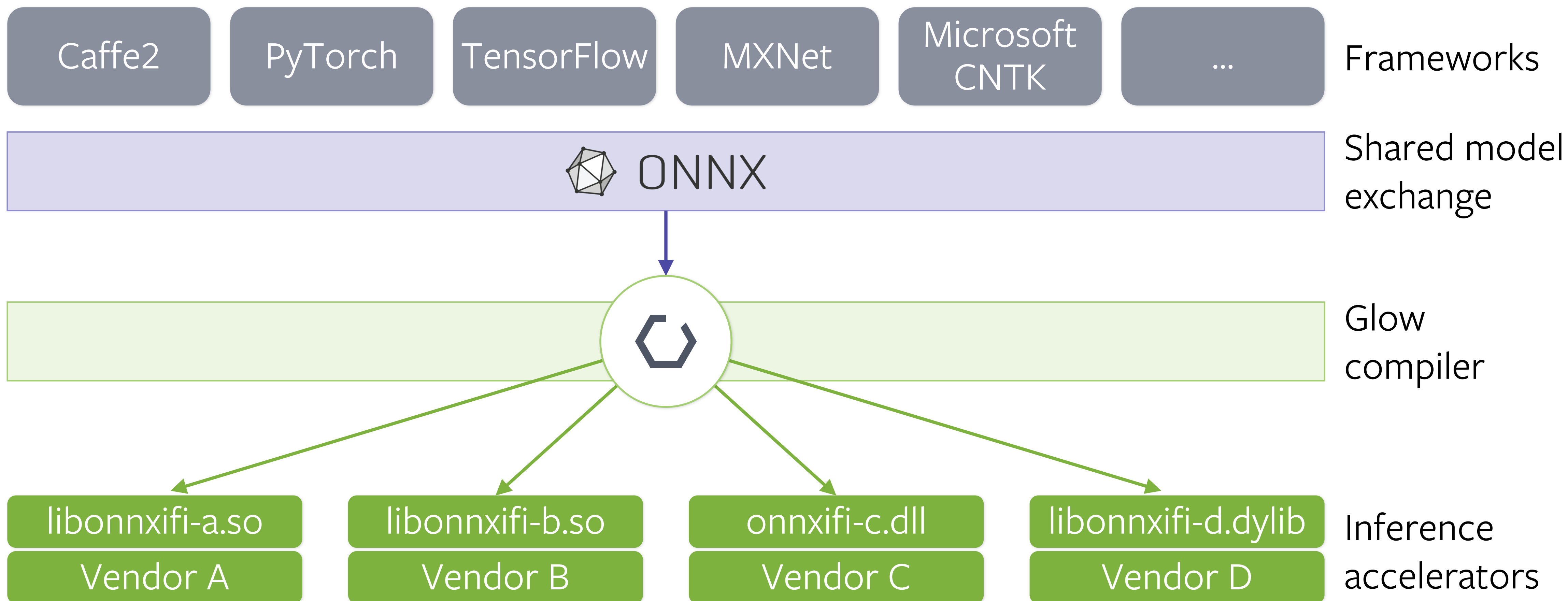
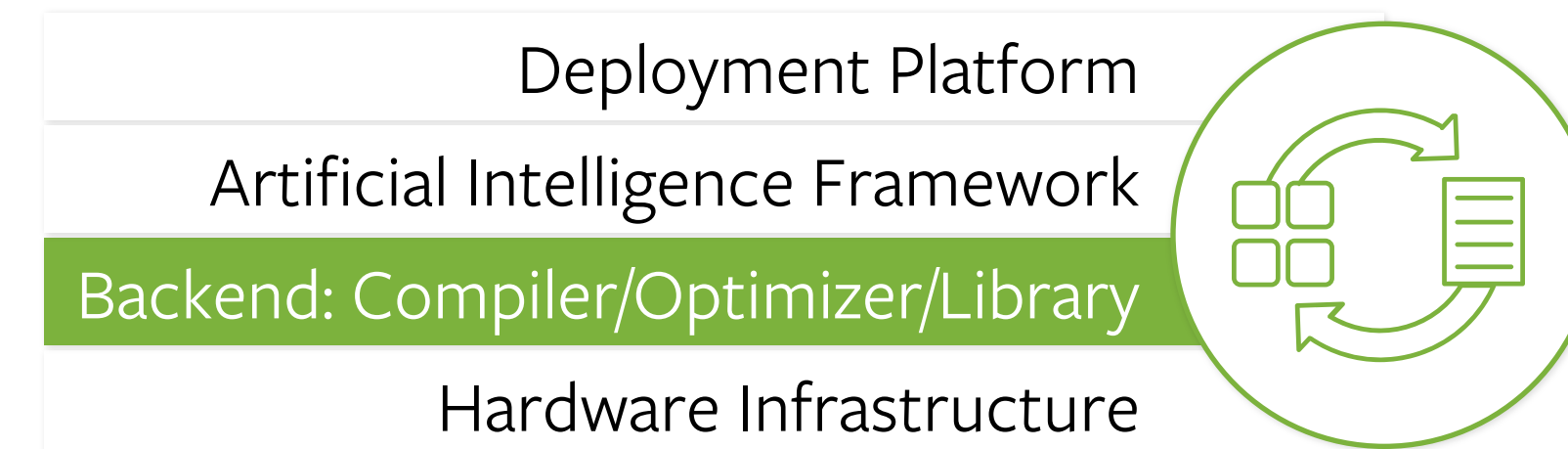
 PyTorch

 Caffe2

To now

 PyTorch

Backends

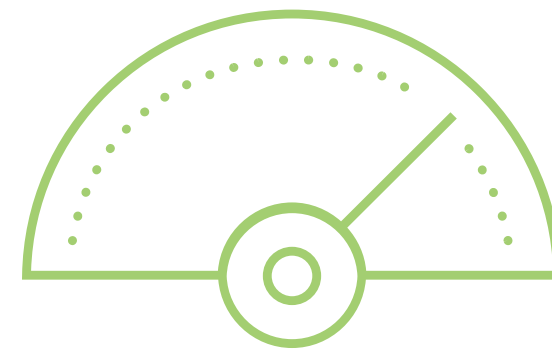


Compilers, Optimizers & Libraries



ML Compiler

Facebook Glow



Vendor Optimizers

Apple CoreML

Nvidia TensorRT

Intel Nervana nGraph

Qualcomm SNPE

...



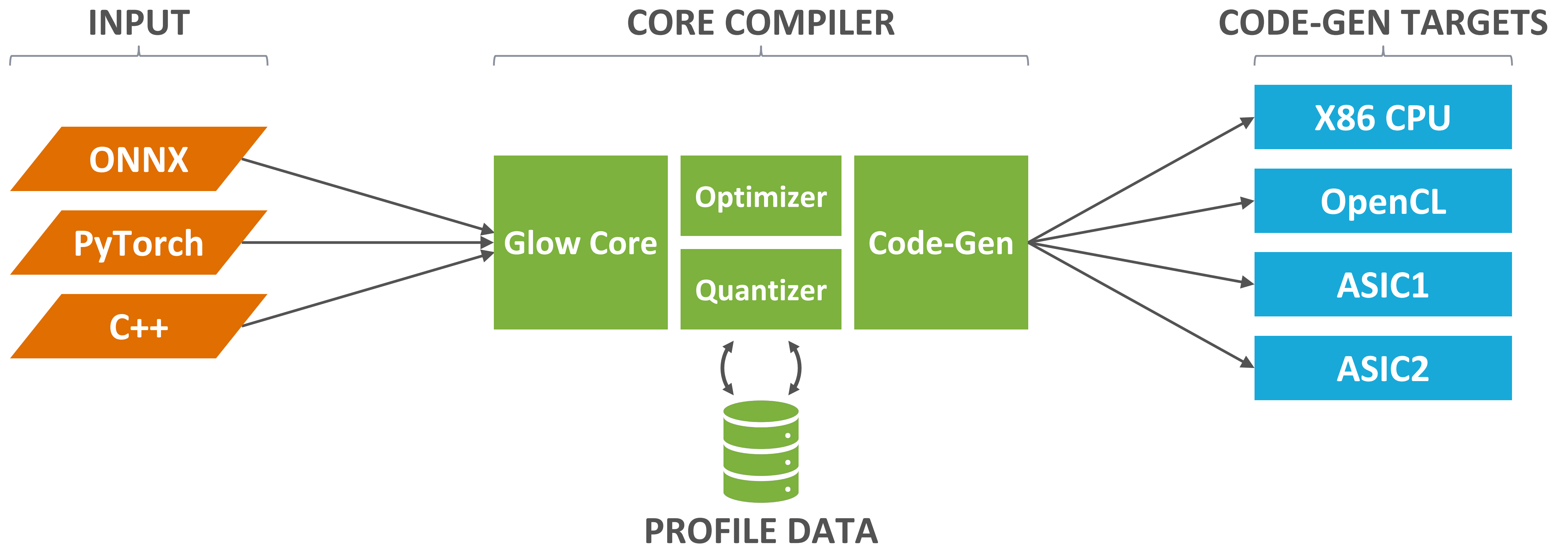
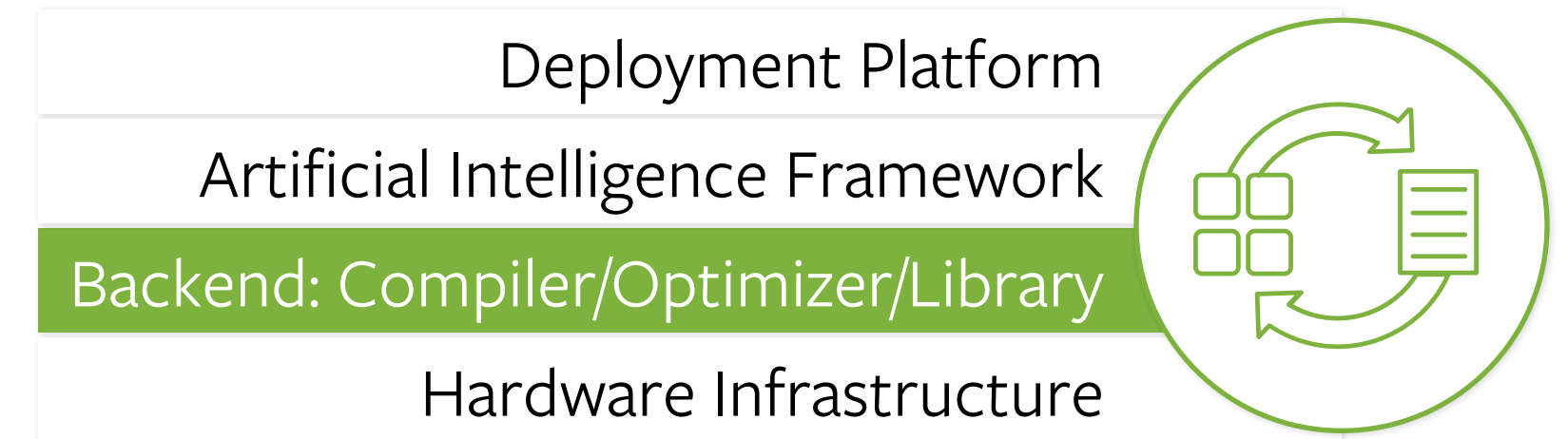
ML Libraries

QNNPACK for mobile CPU

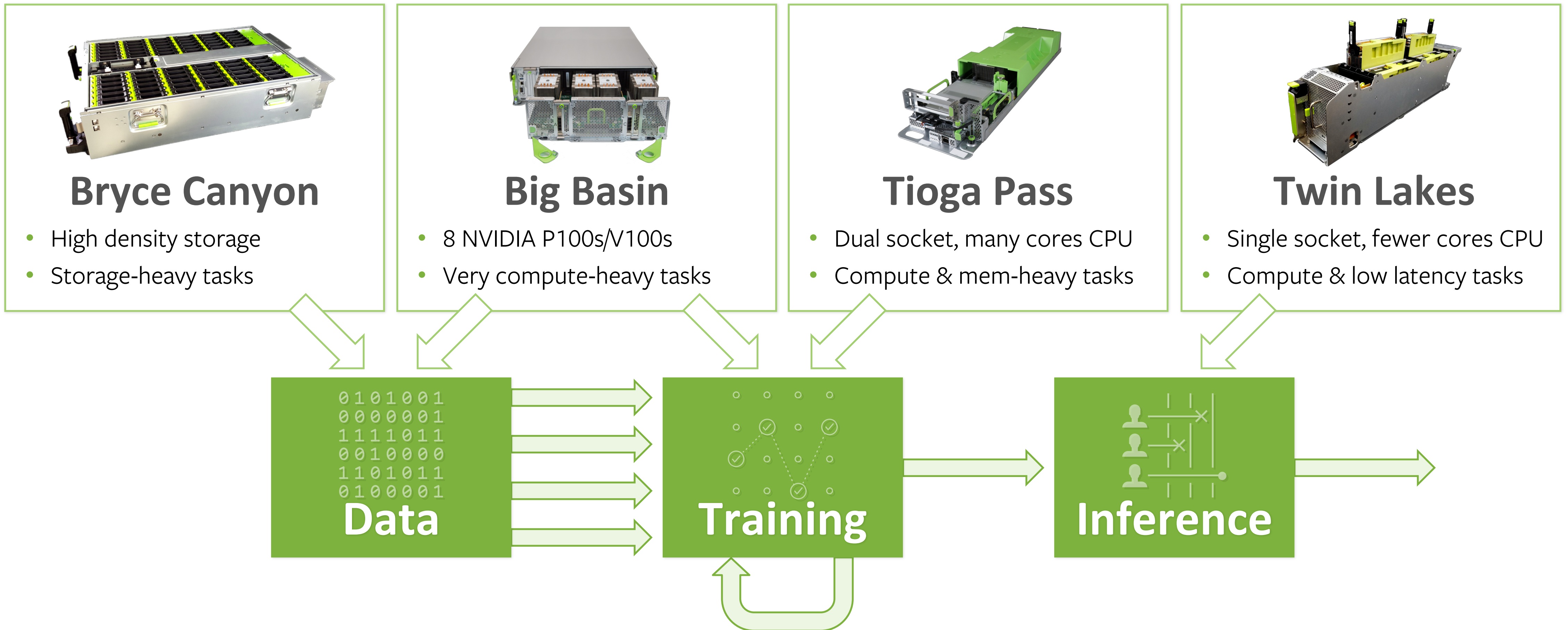
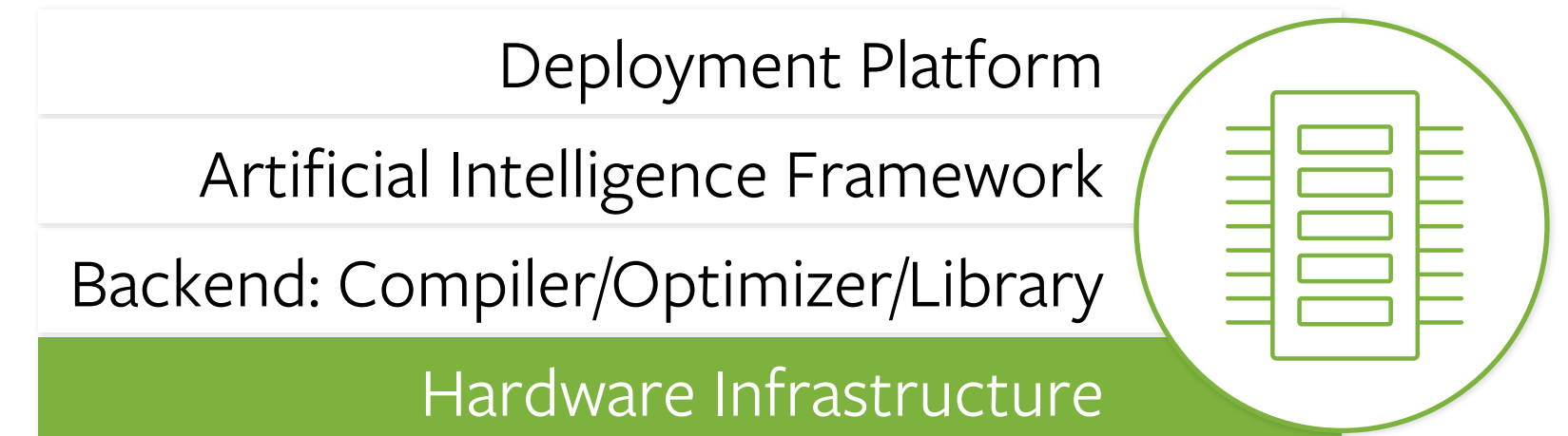
FBGEMM, Intel MKL for server
CPU

CUDNN for GPU

Compilers, Optimizers & Libraries



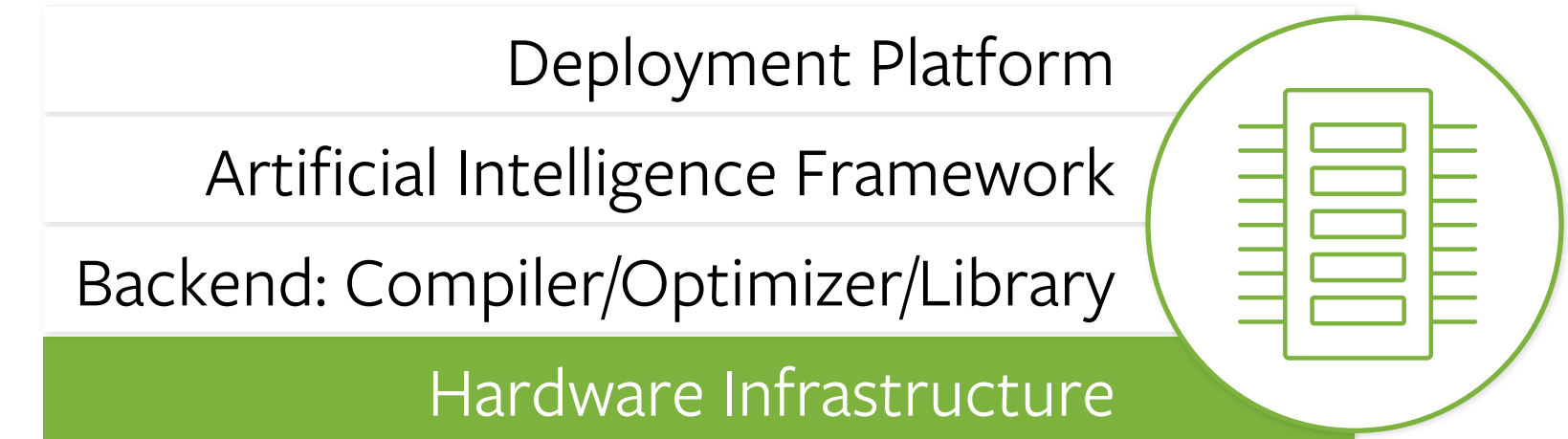
Facebook AI Hardware Today



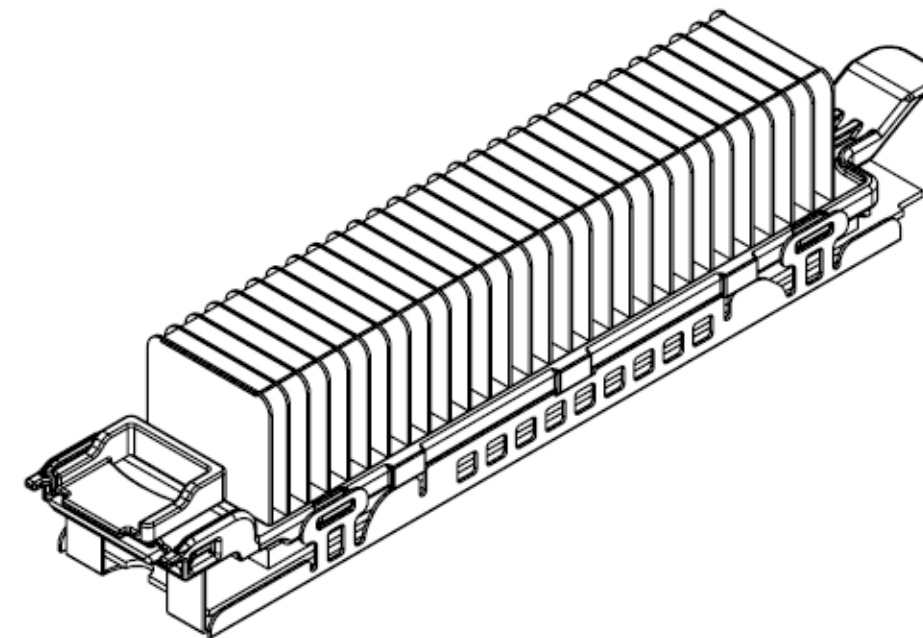
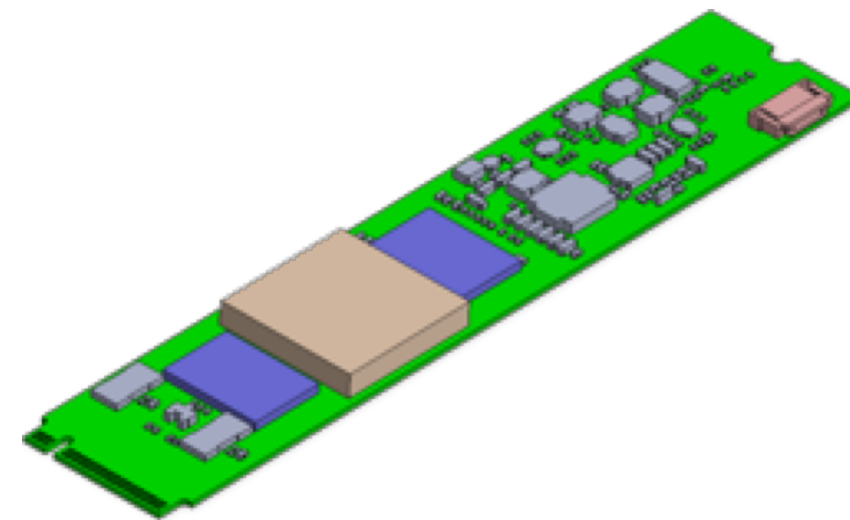
Facebook AI Hardware for Tomorrow

Inference Accelerators

King's Canyon

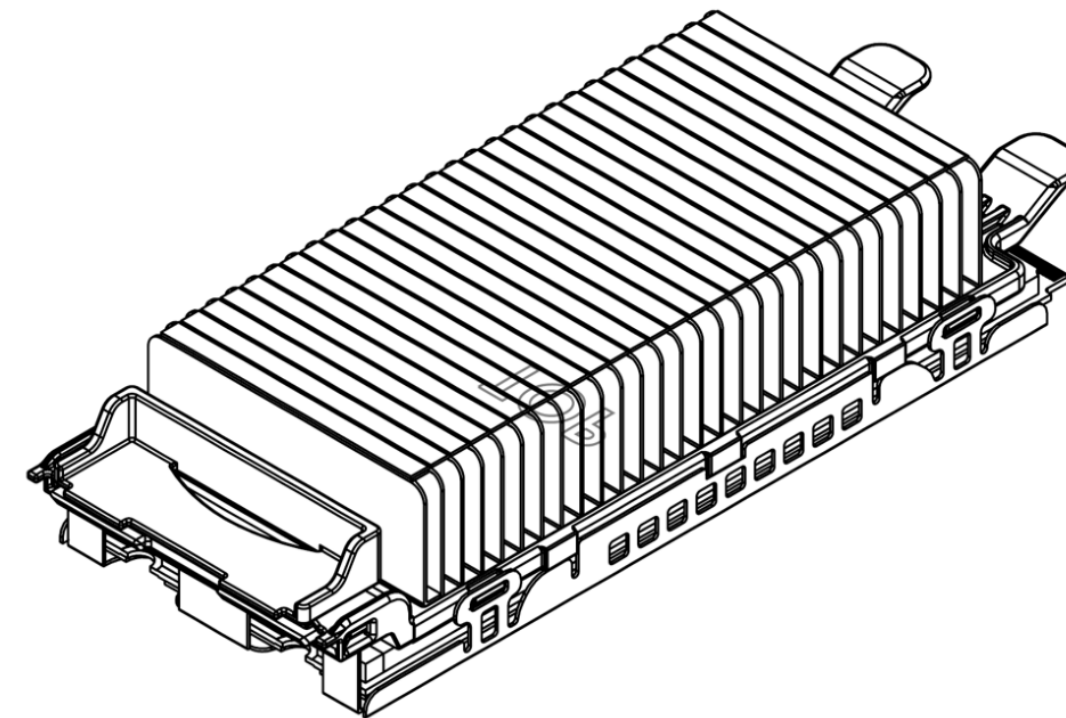
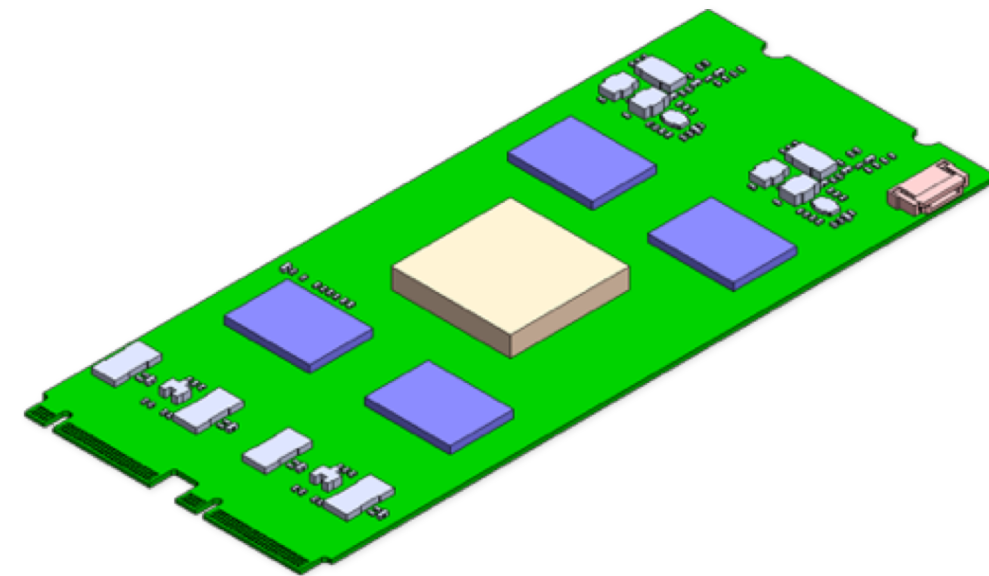


Standard M.2 Form Factor



- ASIC
- DRAM
- 12W TDP
- PCIe x4
- UART/JTAG

Dual M.2 Form Factor



- ASIC
- DRAM
- 20W TDP
- PCIe x8
- UART/JTAG

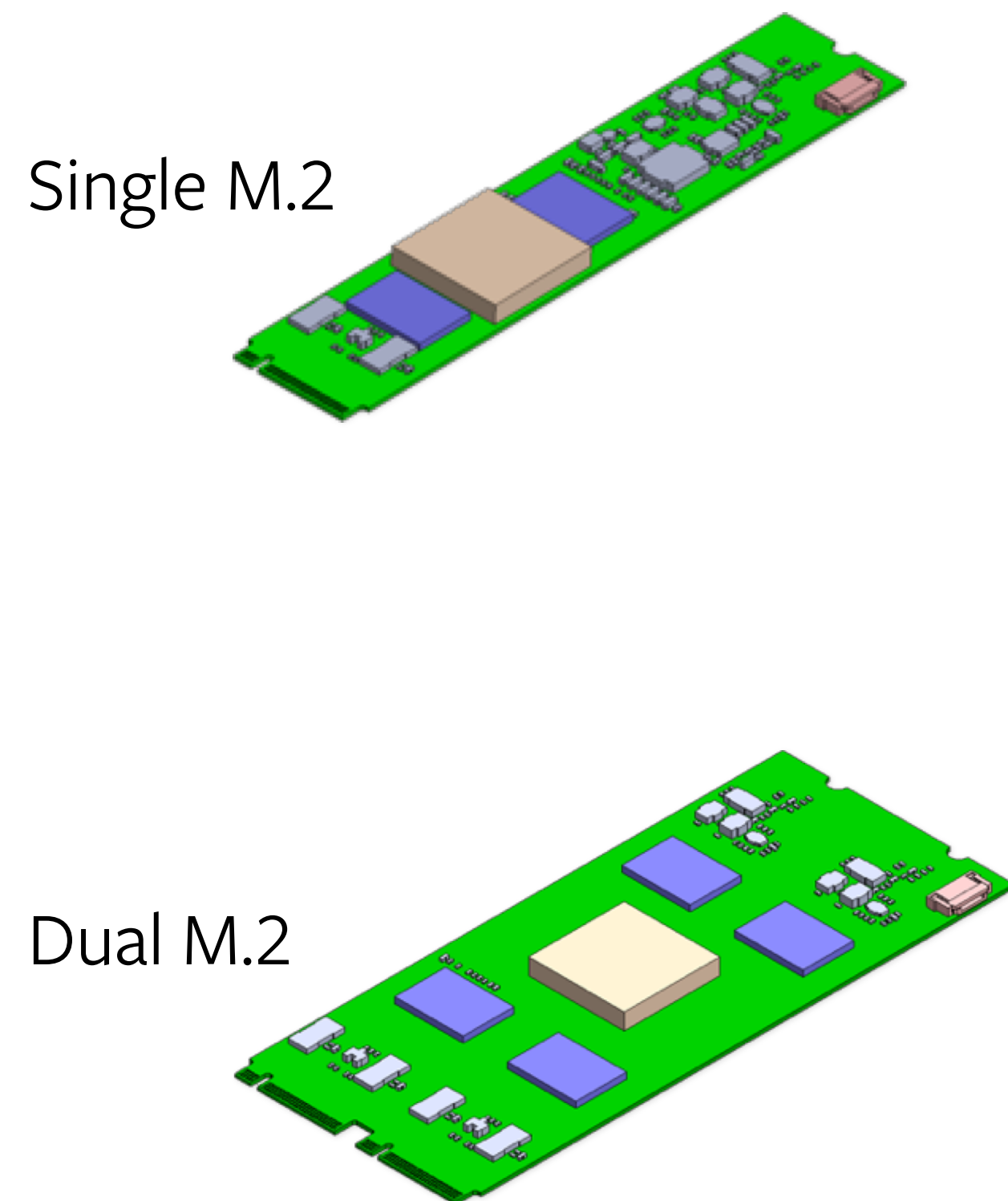
Facebook AI Hardware for Tomorrow

Inference Accelerators

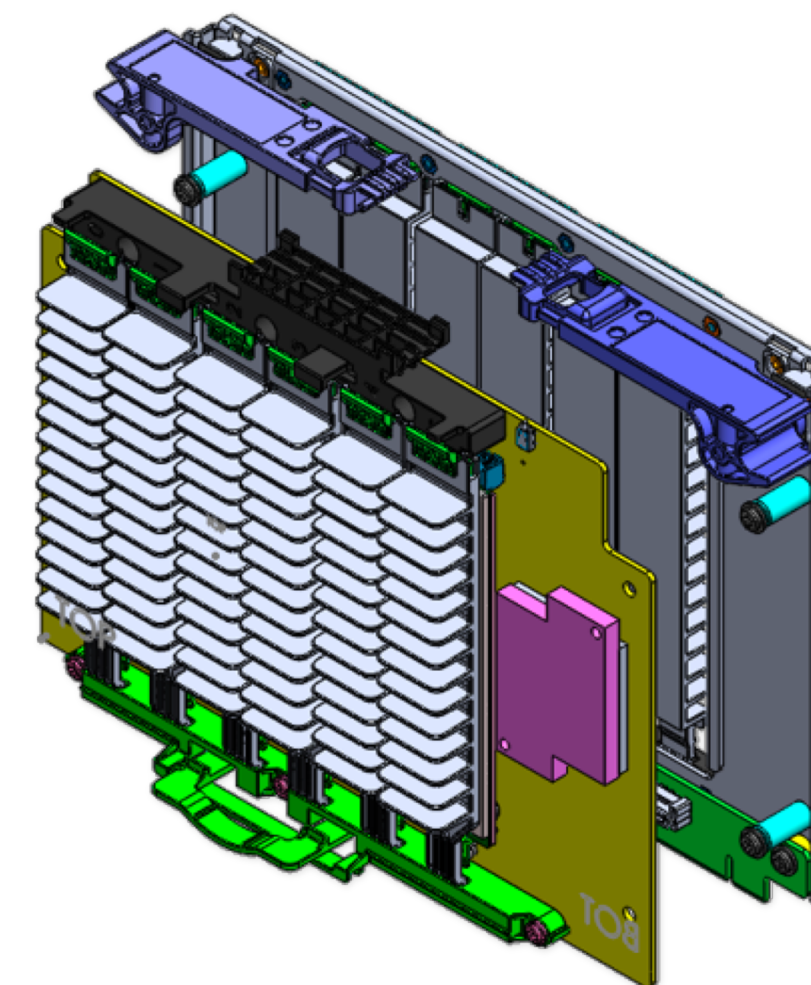
King's Canyon



Inference Module

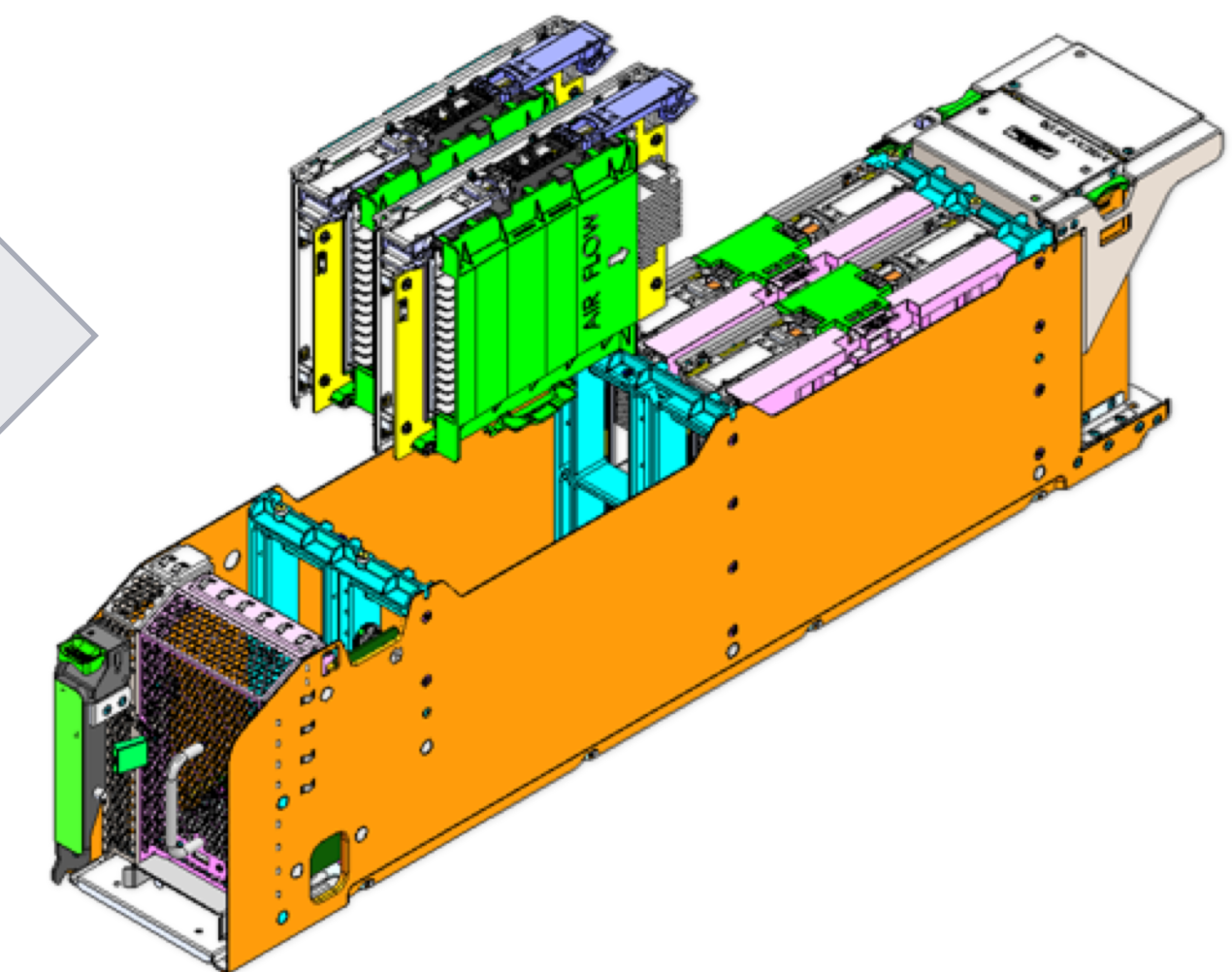


Carrier Card



Up to 12x single /
6x dual M.2

Yosemite V2

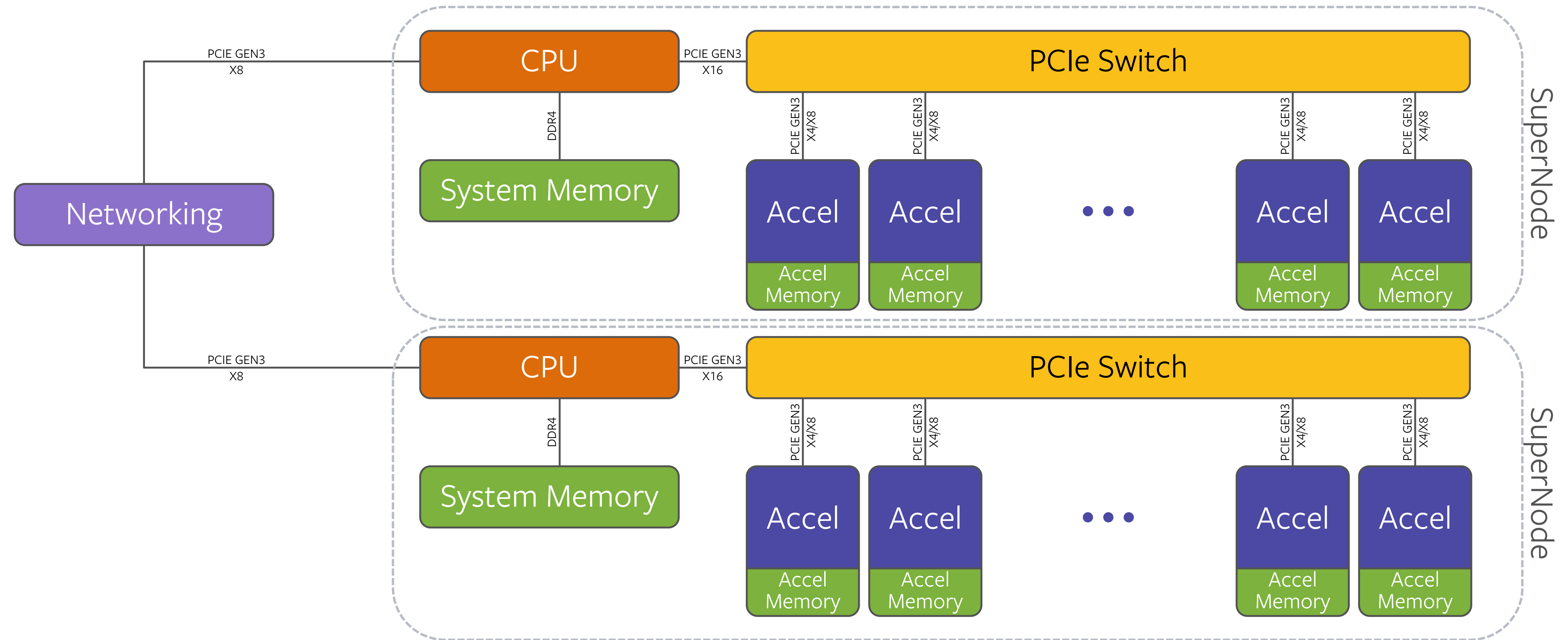
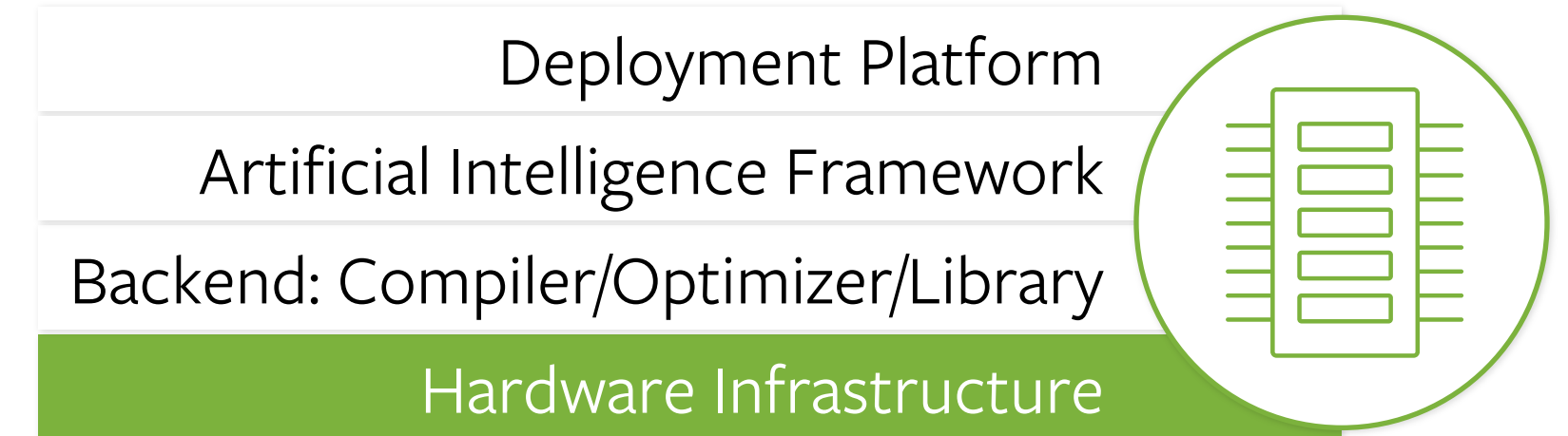


2x Twin Lakes with 2x carrier cards

Facebook AI Hardware for Tomorrow

Inference Accelerators

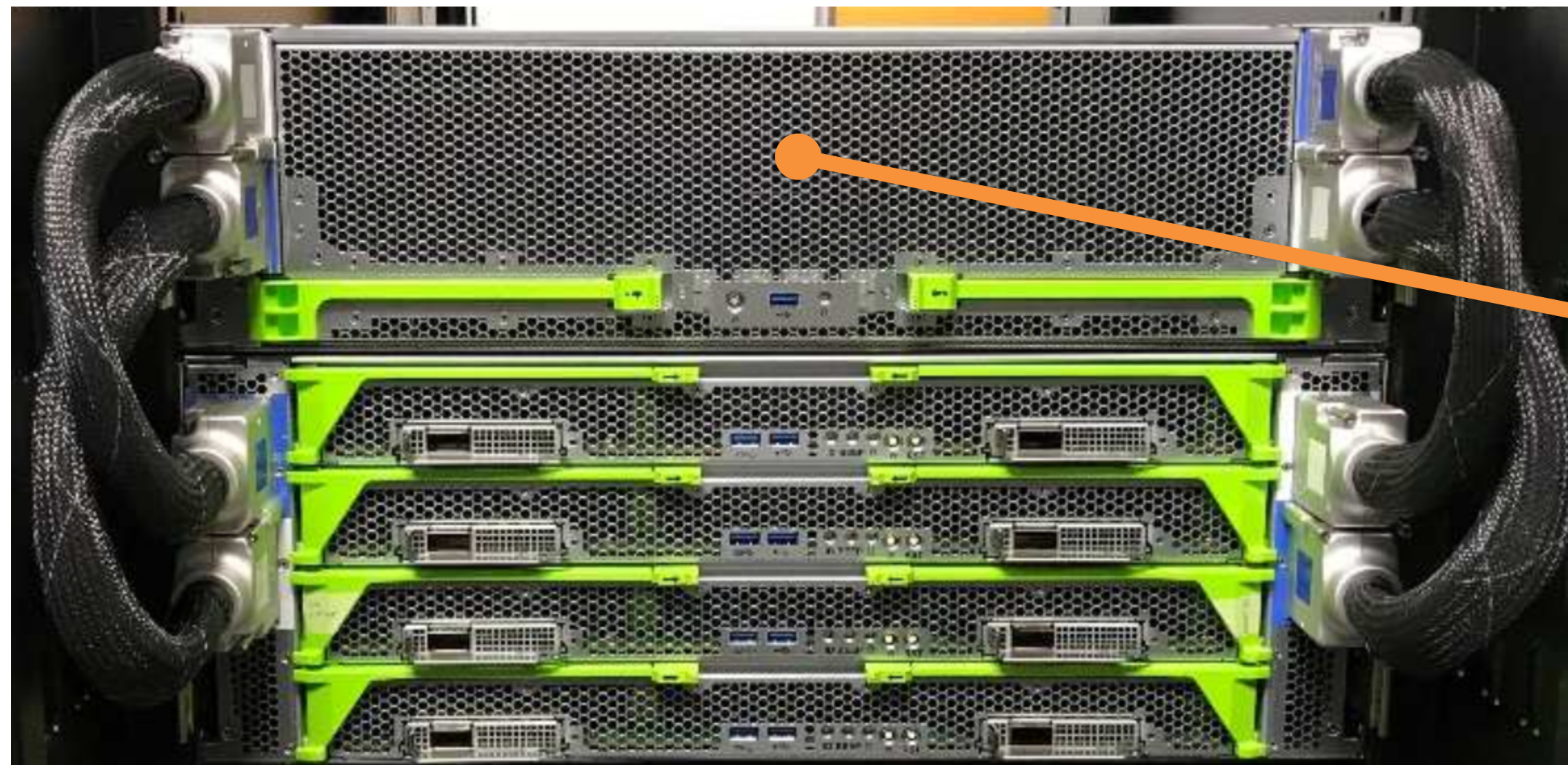
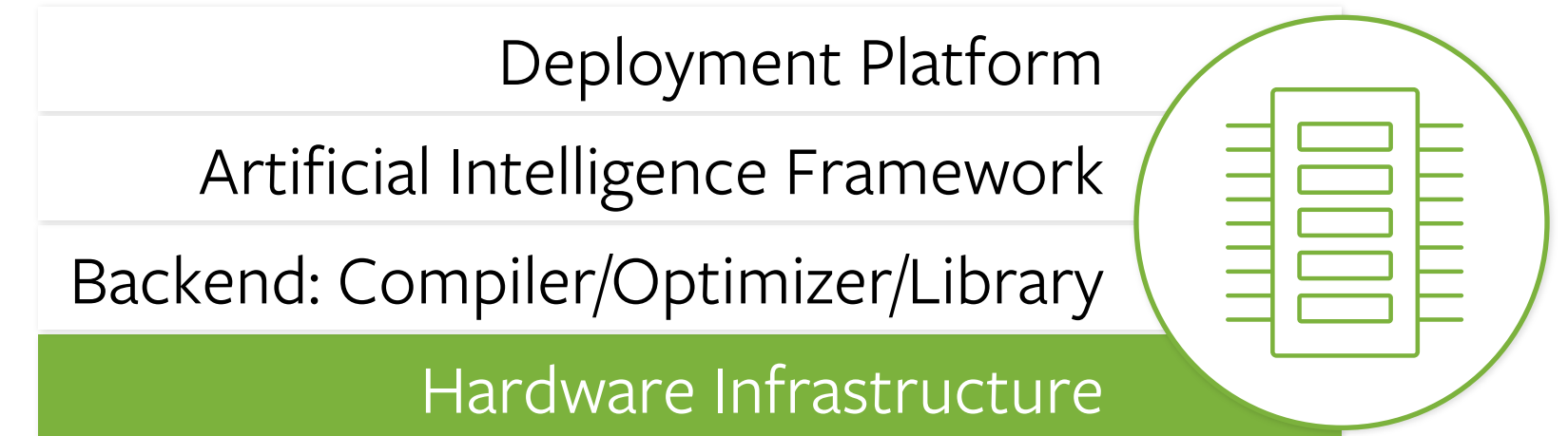
King's Canyon



Facebook AI Hardware for Tomorrow

Training

Large Memory Unified Training Platform Zion



8* Socket system with 8* Accelerators

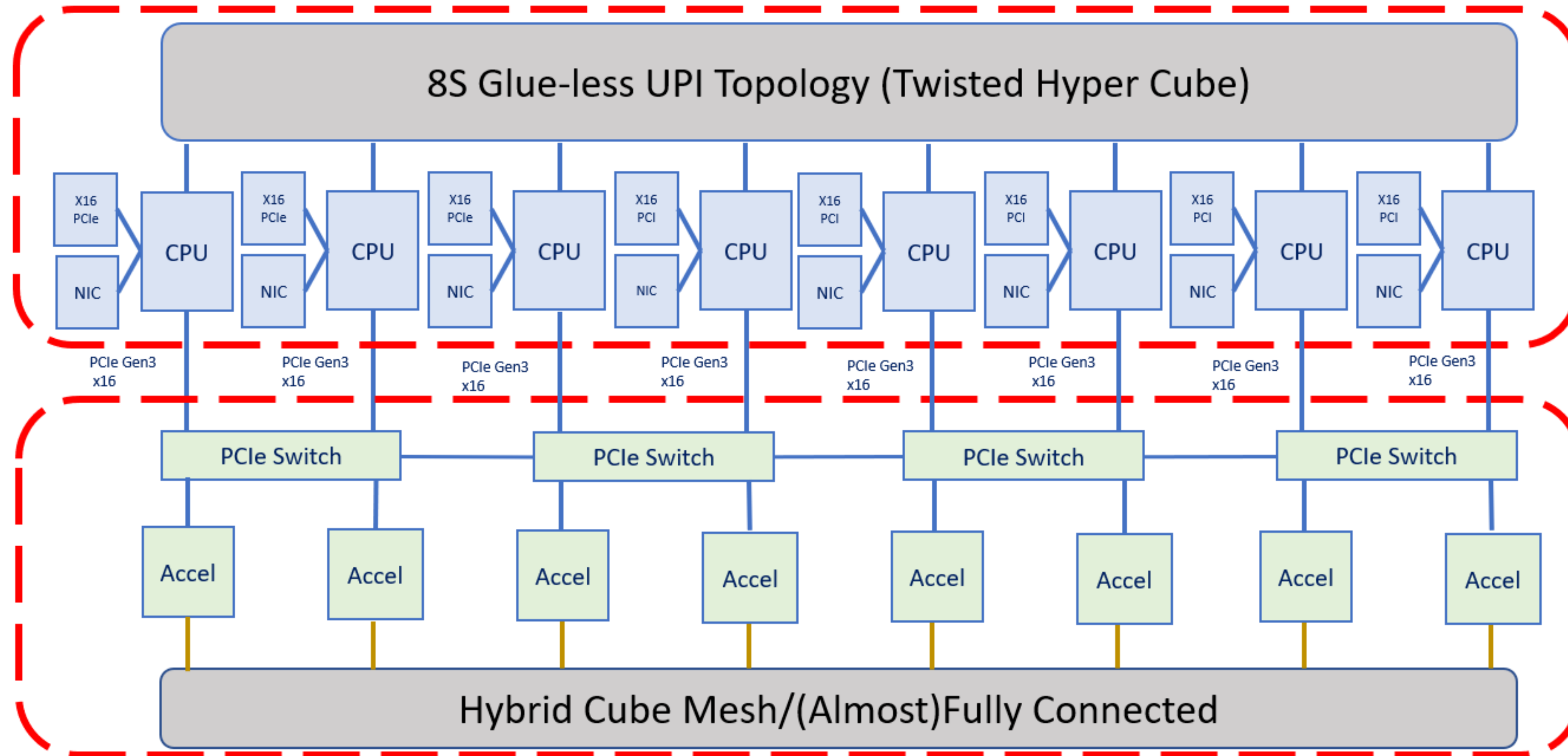
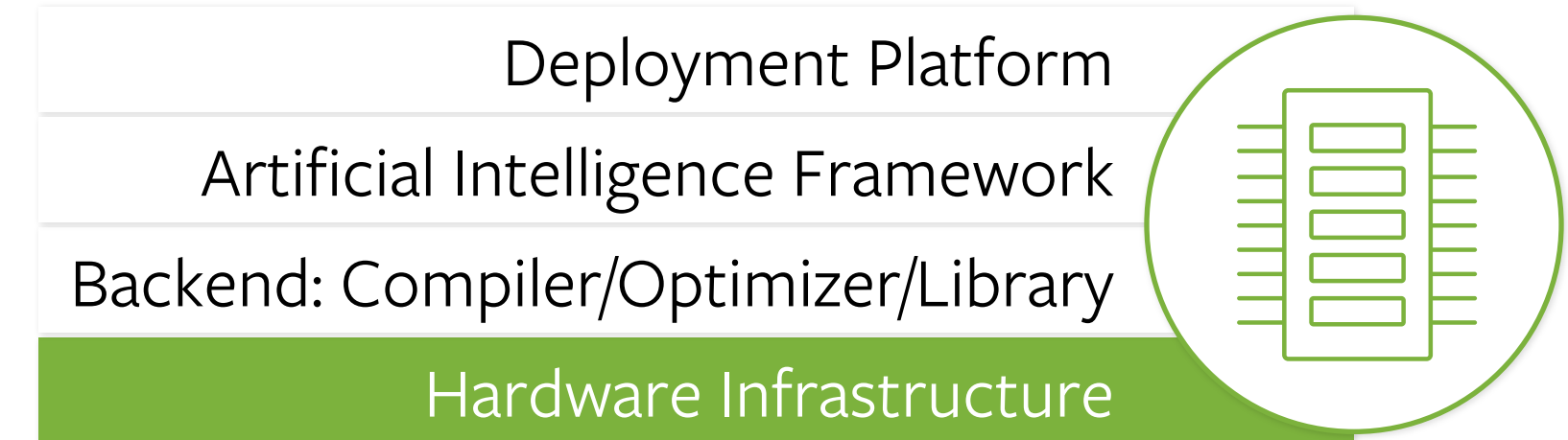


8* OCP Accelerator Modules

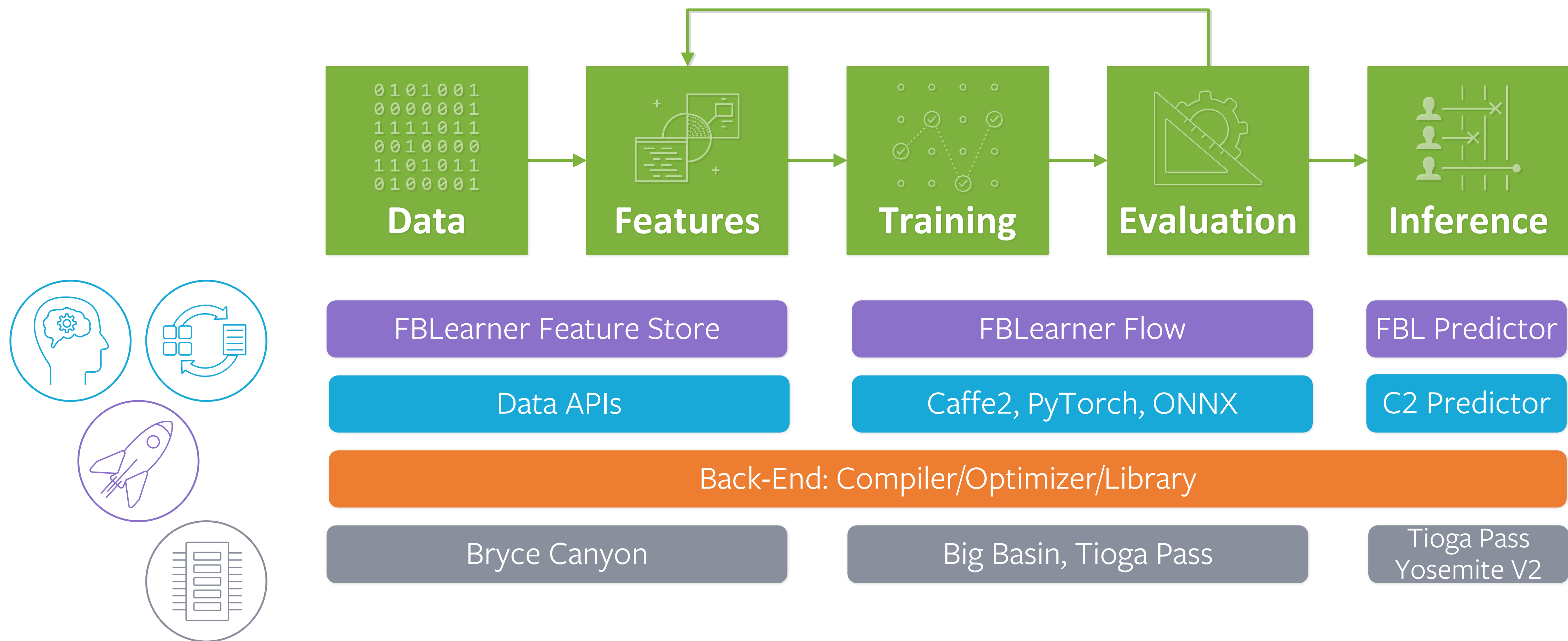
Facebook AI Hardware for Tomorrow

Training

Large Memory Unified Training Platform Zion



Putting it All Together



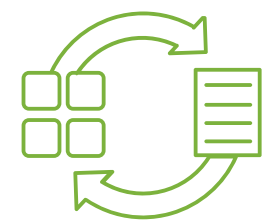


What changes when you scale to over
2 Billion People?

Scaling Challenges / Opportunities

0101001
0000001
1111011
0010000
1101011
0100001

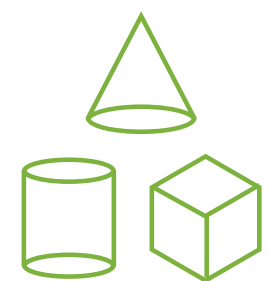
Lots of Data



Lots of Compute



Global scale



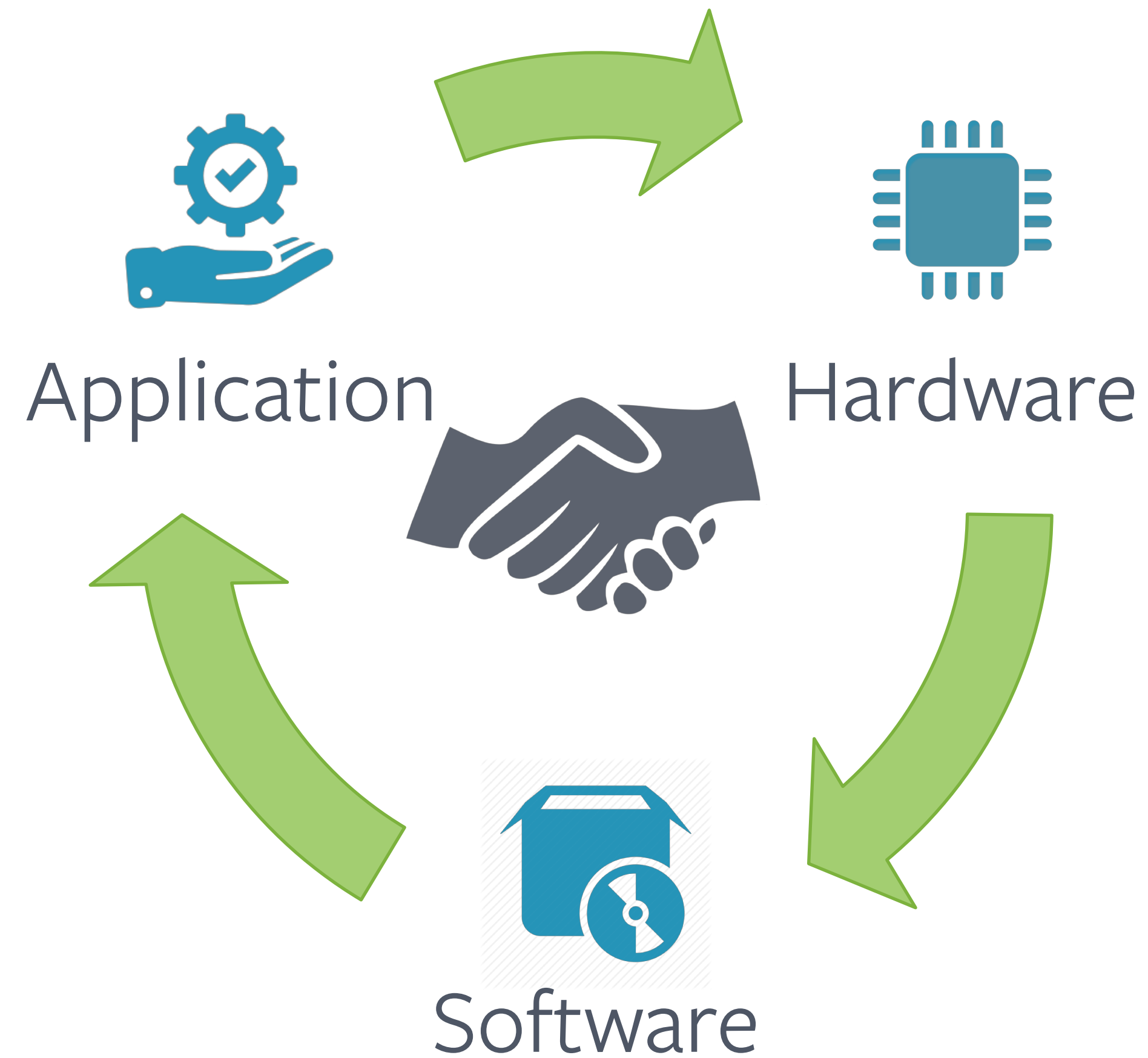
Wide variety of models



Full stack challenges



HW/SW Co-Design



**THIS JOURNEY
1% FINISHED**

Please join us for Facebook AI-related workshops on 3/15:

2:00 PM OCP Accelerator Module (OAM) System: An Open Accelerator Infrastructure Project

3:30 PM OCP Accelerator Module (OAM)



Open. Together.

OCP Global Summit | March 14–15, 2019

