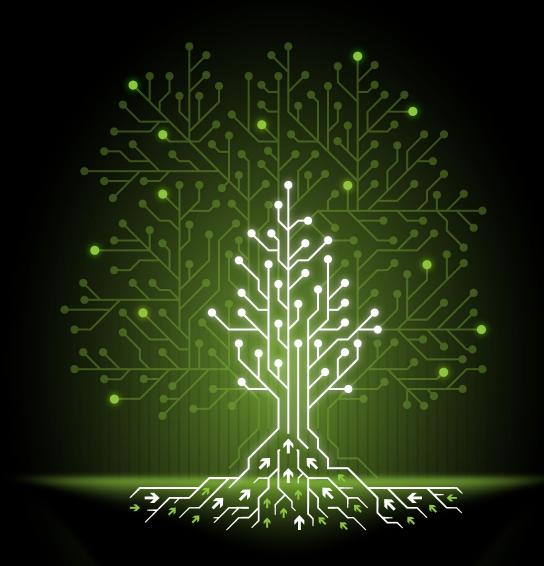
Dan Rabinovitsj

Vice President, Infrastructure Meta





Opening Al Infrastructure

Ushering In The Age Of GenAl



Al workloads are fast evolving...

Recommendation-DLRM Circa: 2022



...exponential changes

Recommendation-DLRM

Circa: 2022

Llama-65B Circa: 2023



...exponential changes

Recommendation-DLRM

Circa: 2022

Llama-65B Circa: 2023



Total Compute (PF/s)

Memory Capacity (TB)

Training Scale (GPUs)

1x

400

10

20x

and we are not done yet...

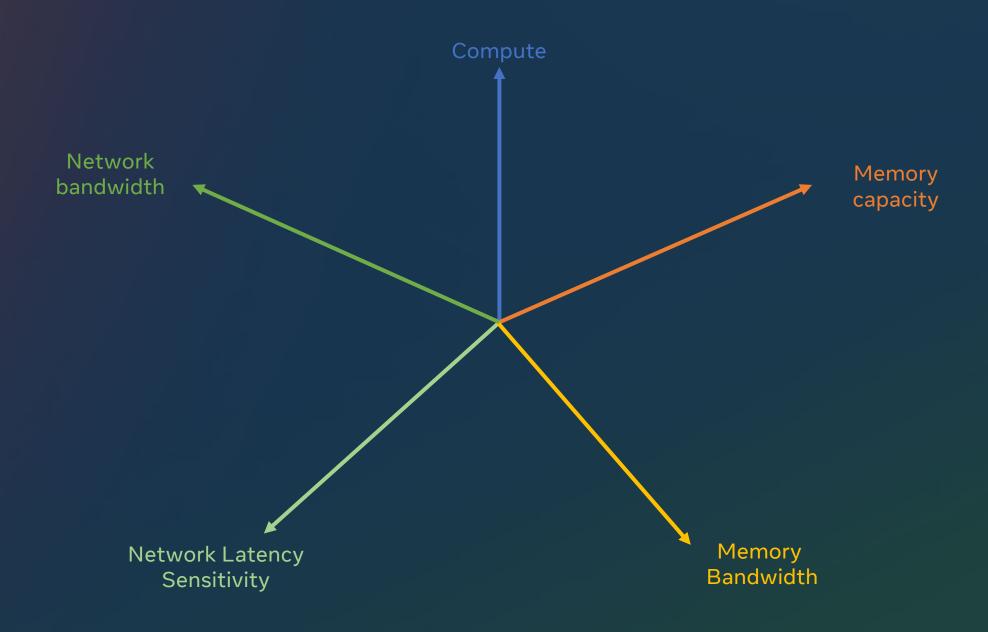
Llama-65B

Circa: 2023

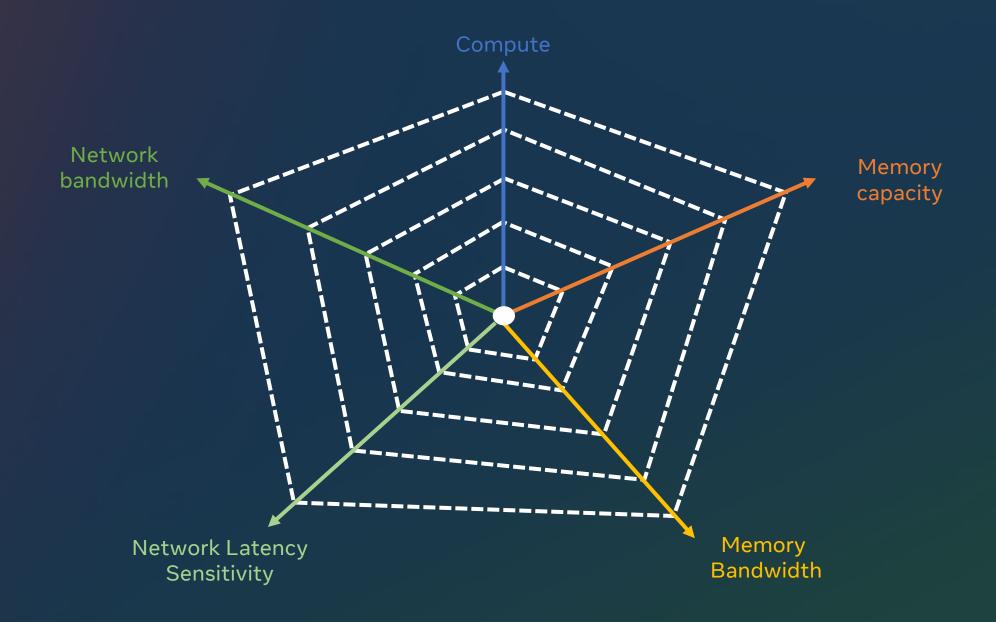
Llama-XX

Circa: 202?

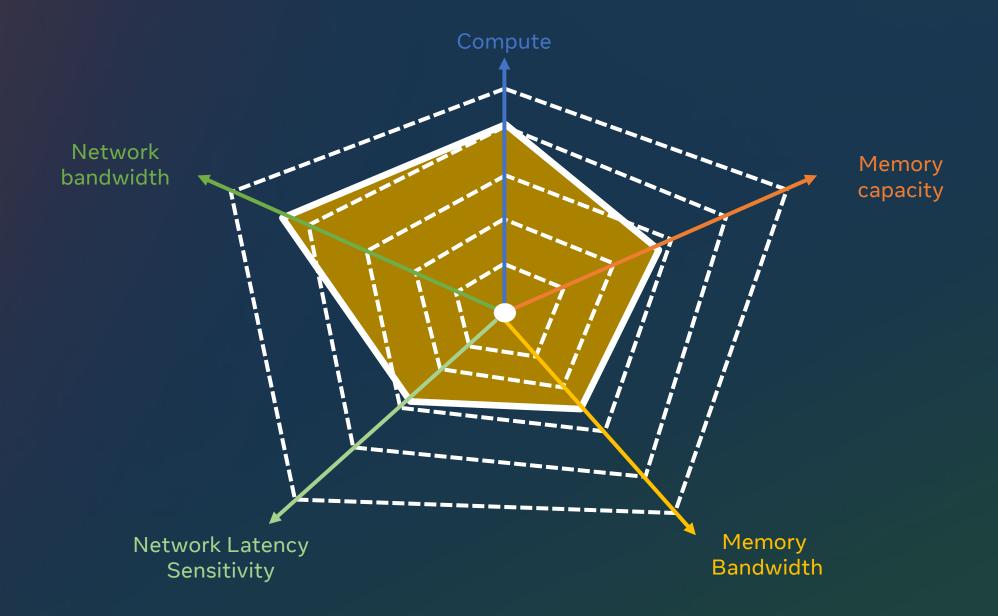
Al Systems Require Multi-Dimensional Design



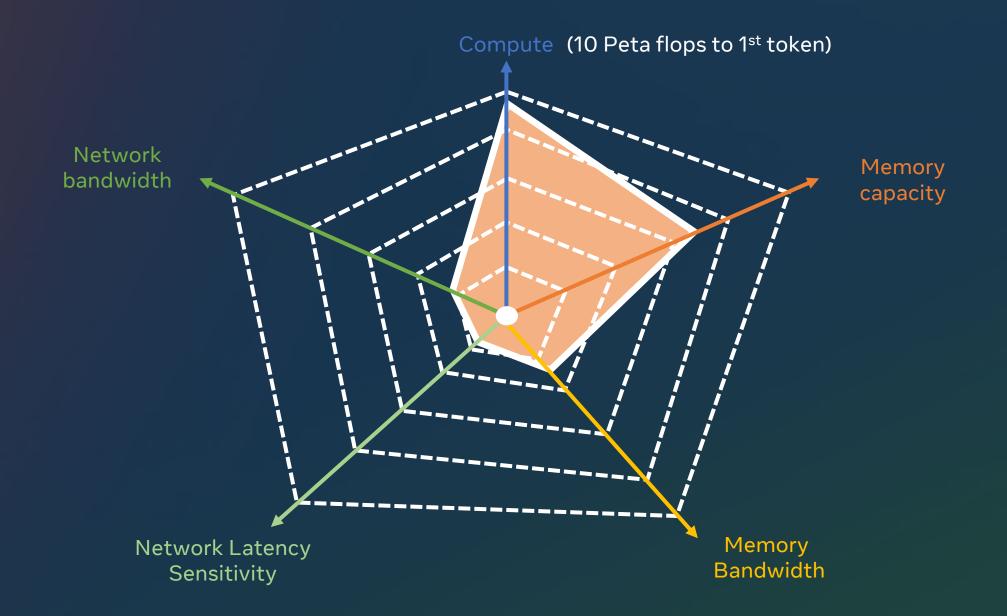
Al workloads impose diverse demands



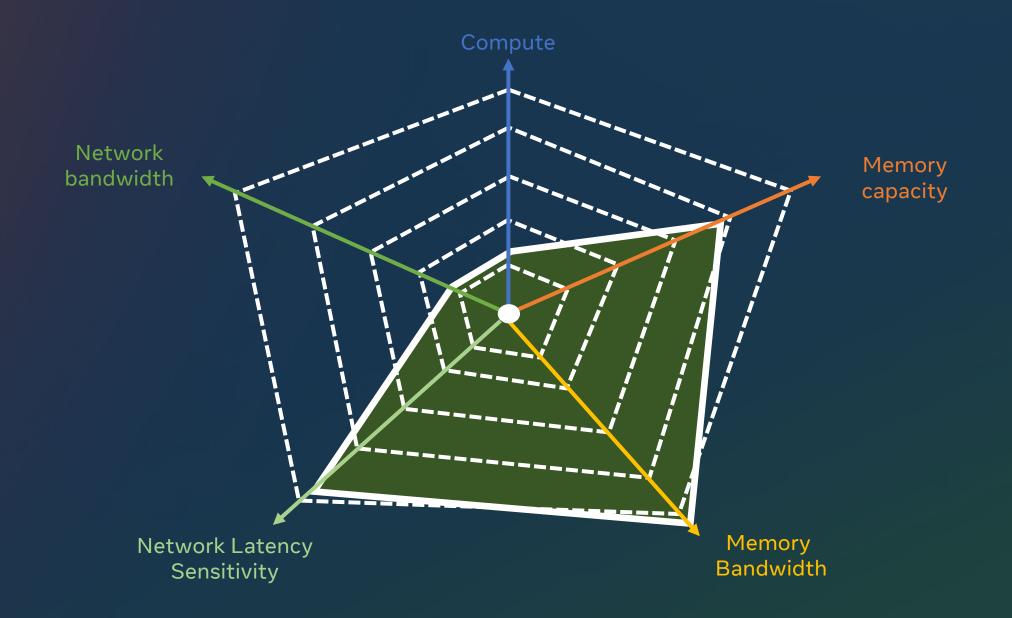
Large Language Models (LLMs) Training



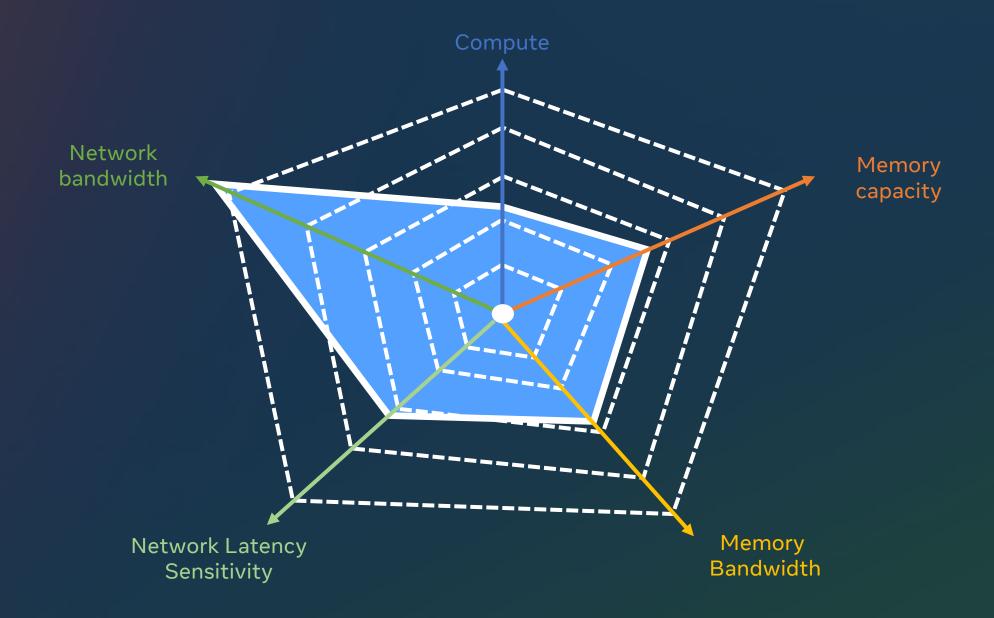
LLMs Inference Prefill



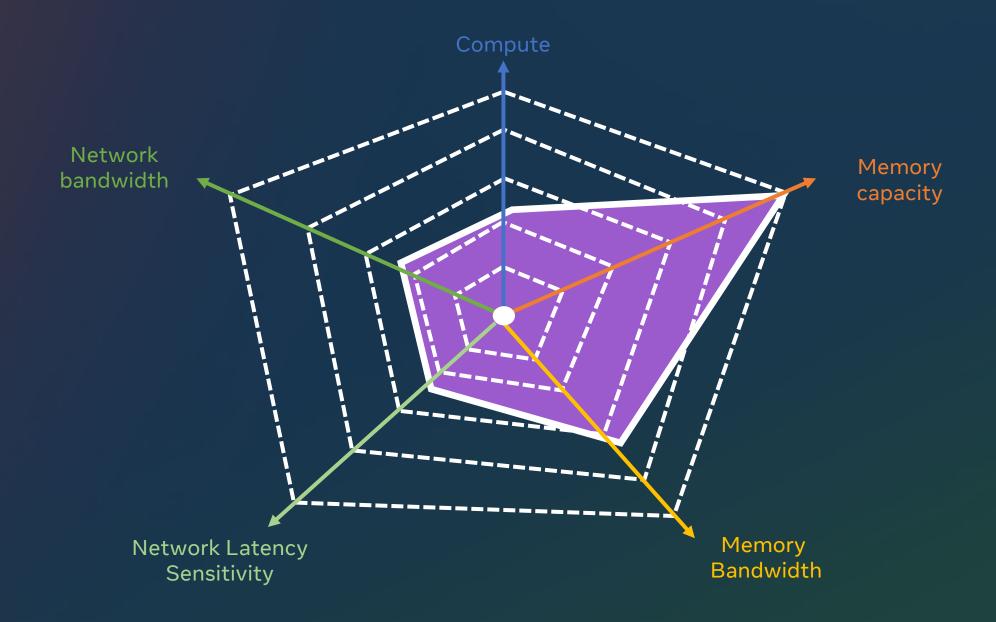
LLMs Inference Decode



Ranking and Recommendation Training

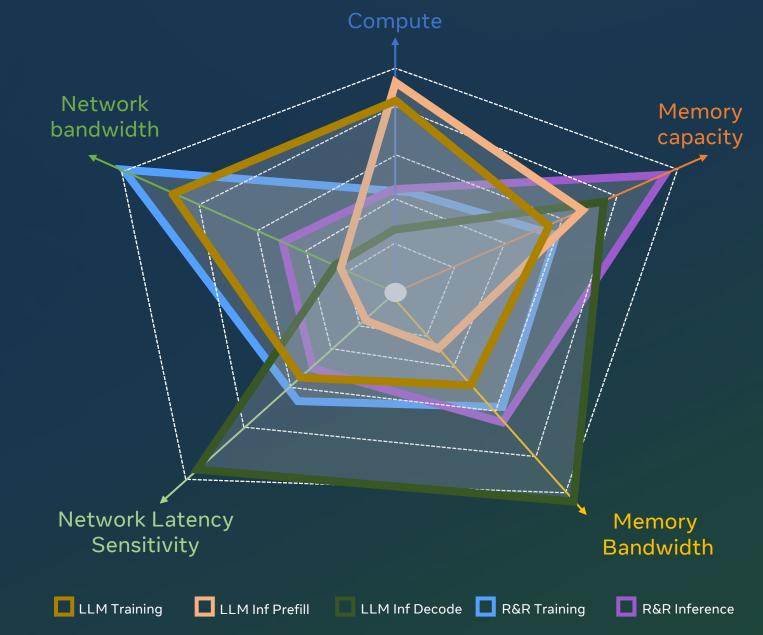


Ranking and Recommendation Inference



Workload Diversity Continues

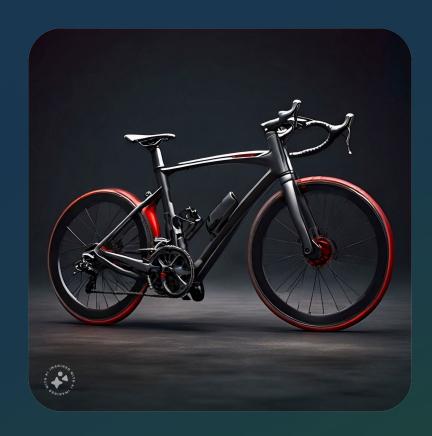
- Difficult to serve all classes of models with a single system design point
- New models & parallelism techniques put unexpected pressures on Al systems
- The next frontier of innovation is in software/hardware co-design

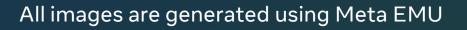


Al Systems tuned for diverse applications







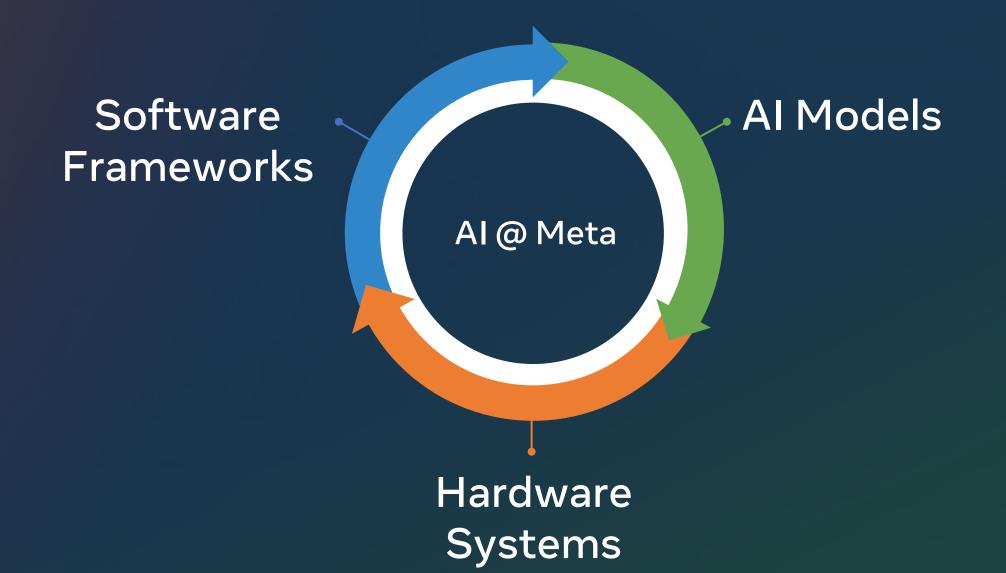




The Big Question

Is it realistic to open finely-tuned, vertically-integrated Al systems?

Al for All: The open model is already out there



Al for All: The open model is working!

```
import torch
free torch autograd import Variable
import torch nn is nn
import torch nn is nn
import torch nn functional as F

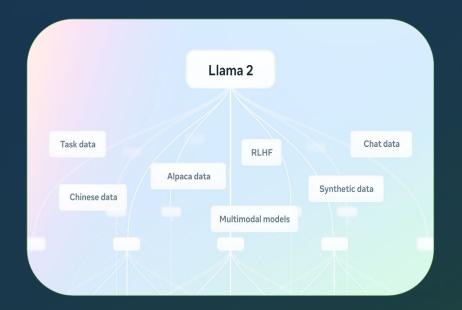
class Net(nn.Nodule):

def __init__(self):
    super(net, self).__init__()
    # input image channel, 6 output channels, 5xS square convolution
    # kennel
    self conv1 = nn.Conv2d(1, 6, 5)
    self conv2 = nn.Conv2d(6, 16, 5)
    self conv2 = nn.Conv2d(6, 16, 5)
    self fc1 = nn.linear(16, 16, 5)
    self fc2 = nn.linear(120, 84)
    self fc3 = nn.linear(120, 84)
    self fc3 = nn.linear(120, 84)
    self fc4 = nn.linear(120, 84)
    self fc3 = nn.linear(120, 84)
    self fc4 = nn.linear(120, 84)
    self fc3 = nn.linear(120, 84)
    self fc4 = nn.linear(120, 10)
    x = F.eax_pool2d F.elu self conv1(x), (2, 2))
    x = F.eax_pool2d F.elu self conv1(x), (2, 2))
    x = F.elu self fc1(x))
    x = F.elu self fc2(x))
    x = F.elu self fc2(x))
    x = relu self fc2(x))
    x = self-fc3(x)
    return x

def num_flat_features(self, x):
    size = x size()[1:] # all dimensions except the batch dimension
    num_features = 1
    for s in size:
    num_features = s
    return num_features

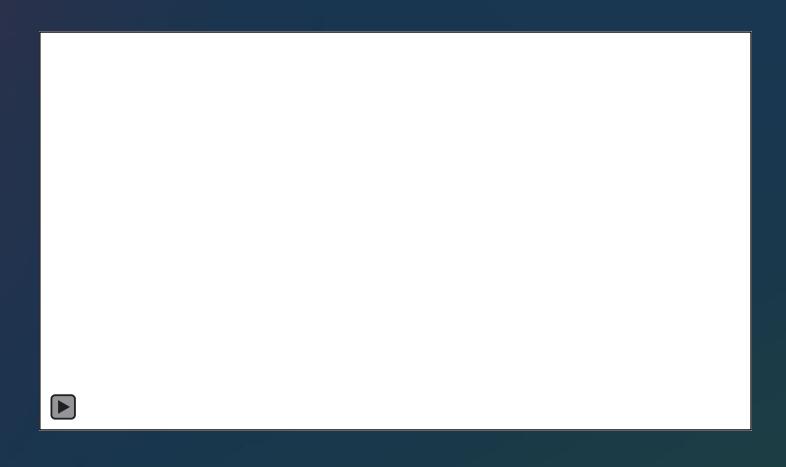
net = let()
    print(net)
```



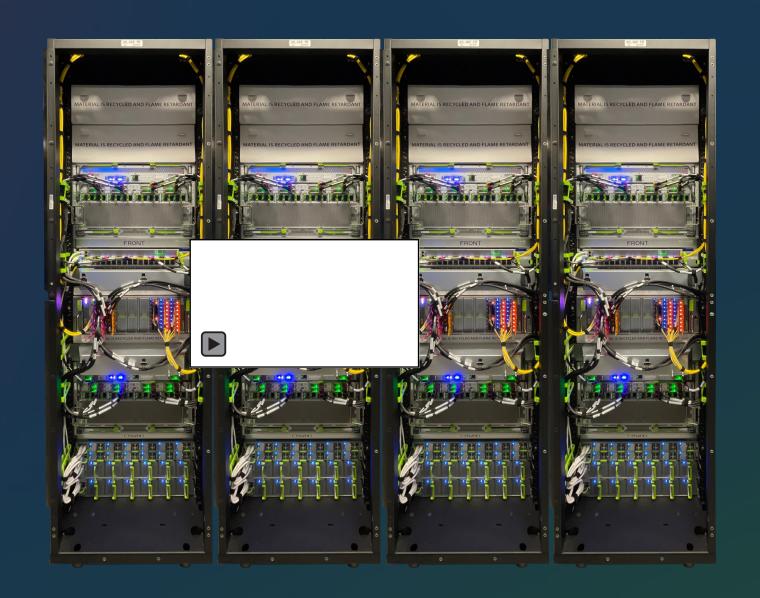




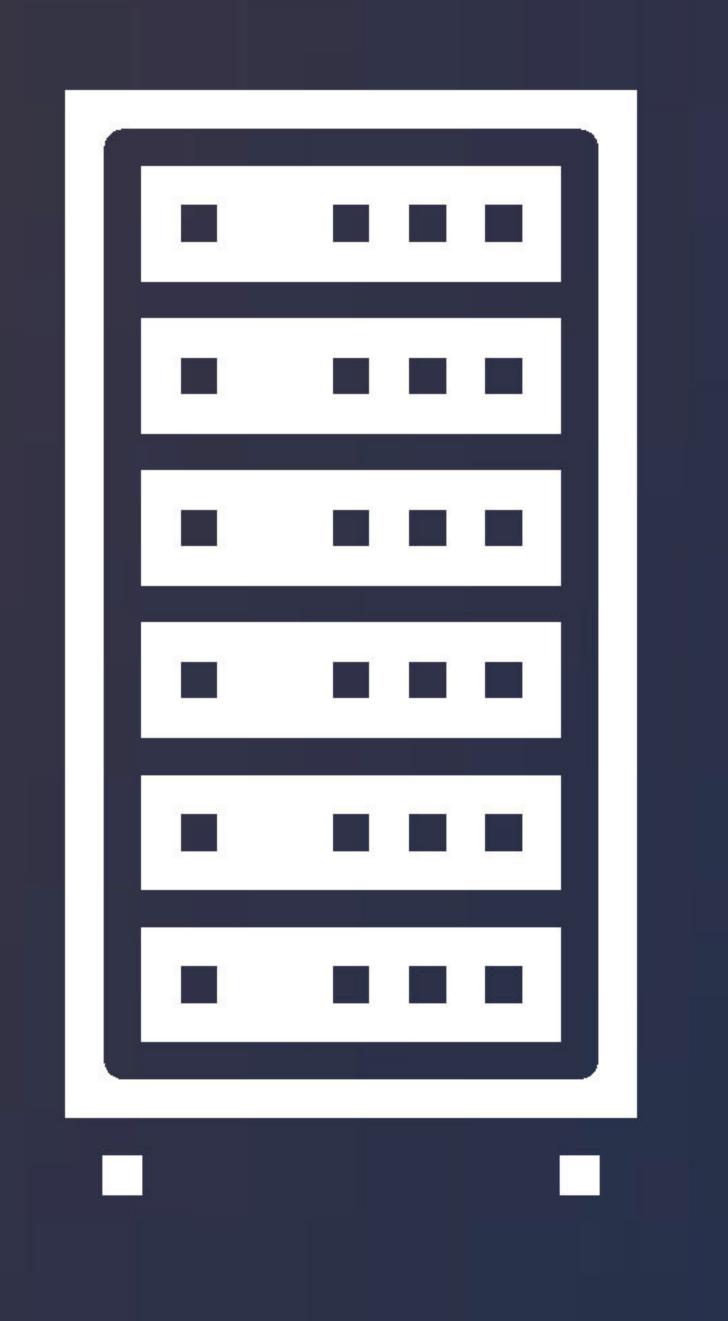
Our Al journey began with Grand Teton



Grand Teton GPU Clusters in Production



Building Open Al Infrastructure



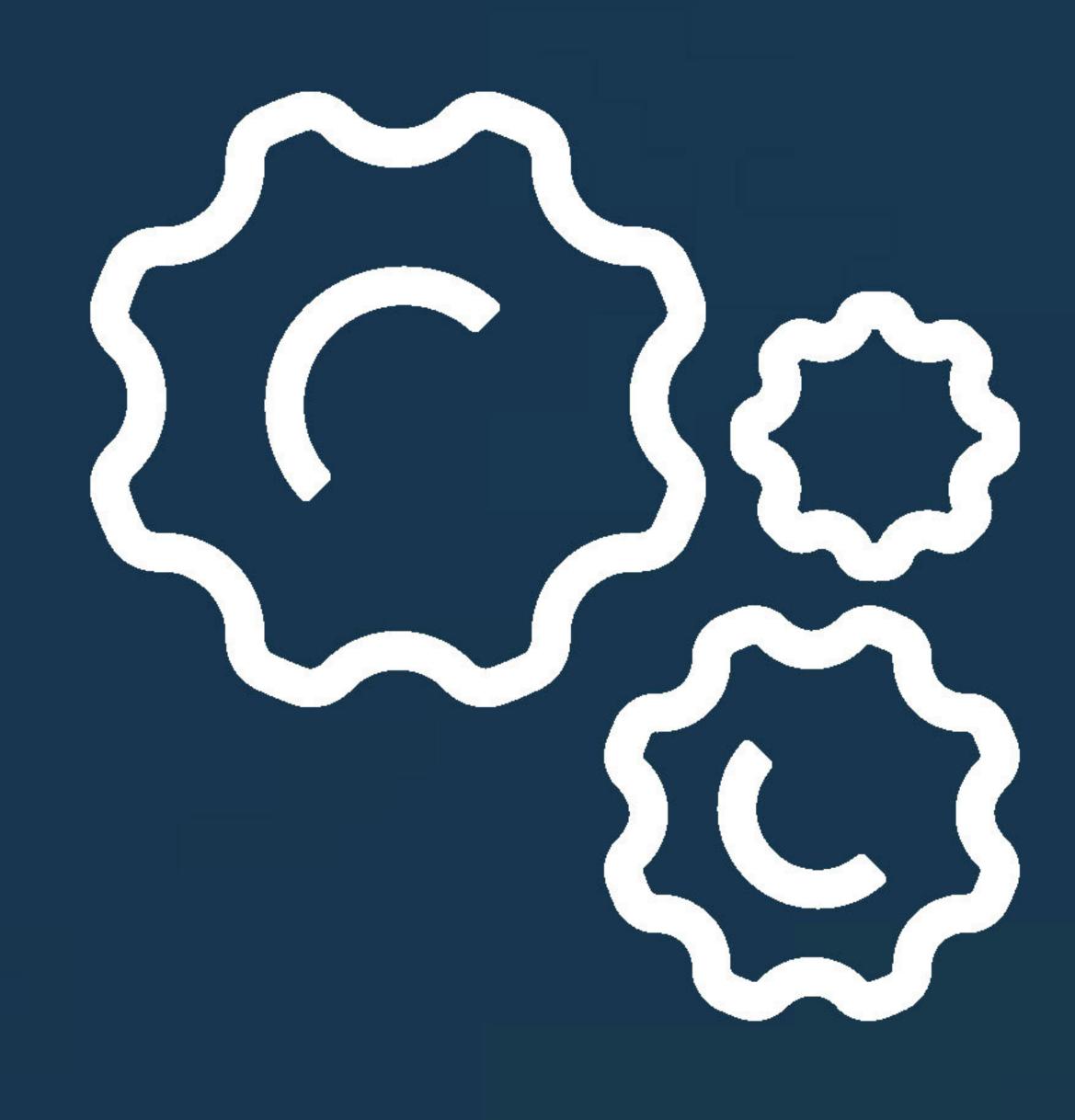
Common
Chassis & Rack
Architecture



Standard
Liquid Cooling Designs
& Blind Mate
Connectors



Open
Integration, Validation
and Test Playbooks

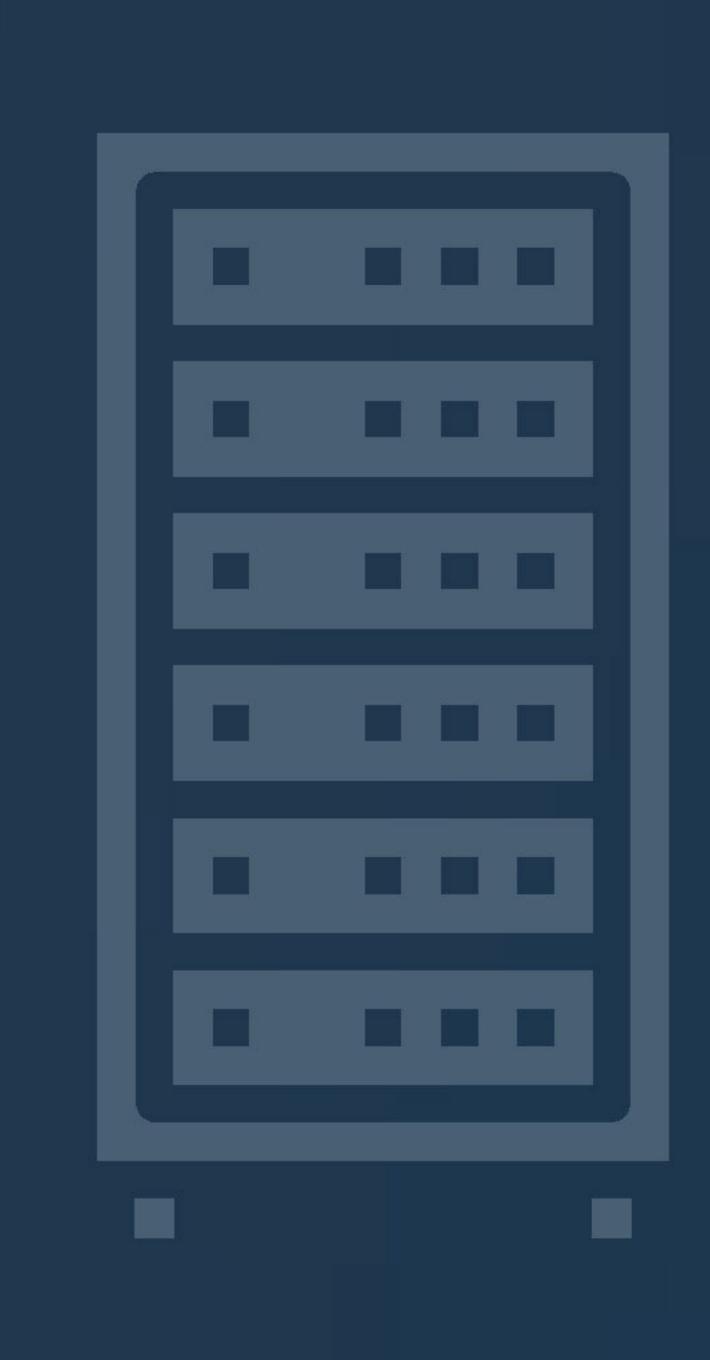


Universal
Tooling/Telemetry,
Software, and Support

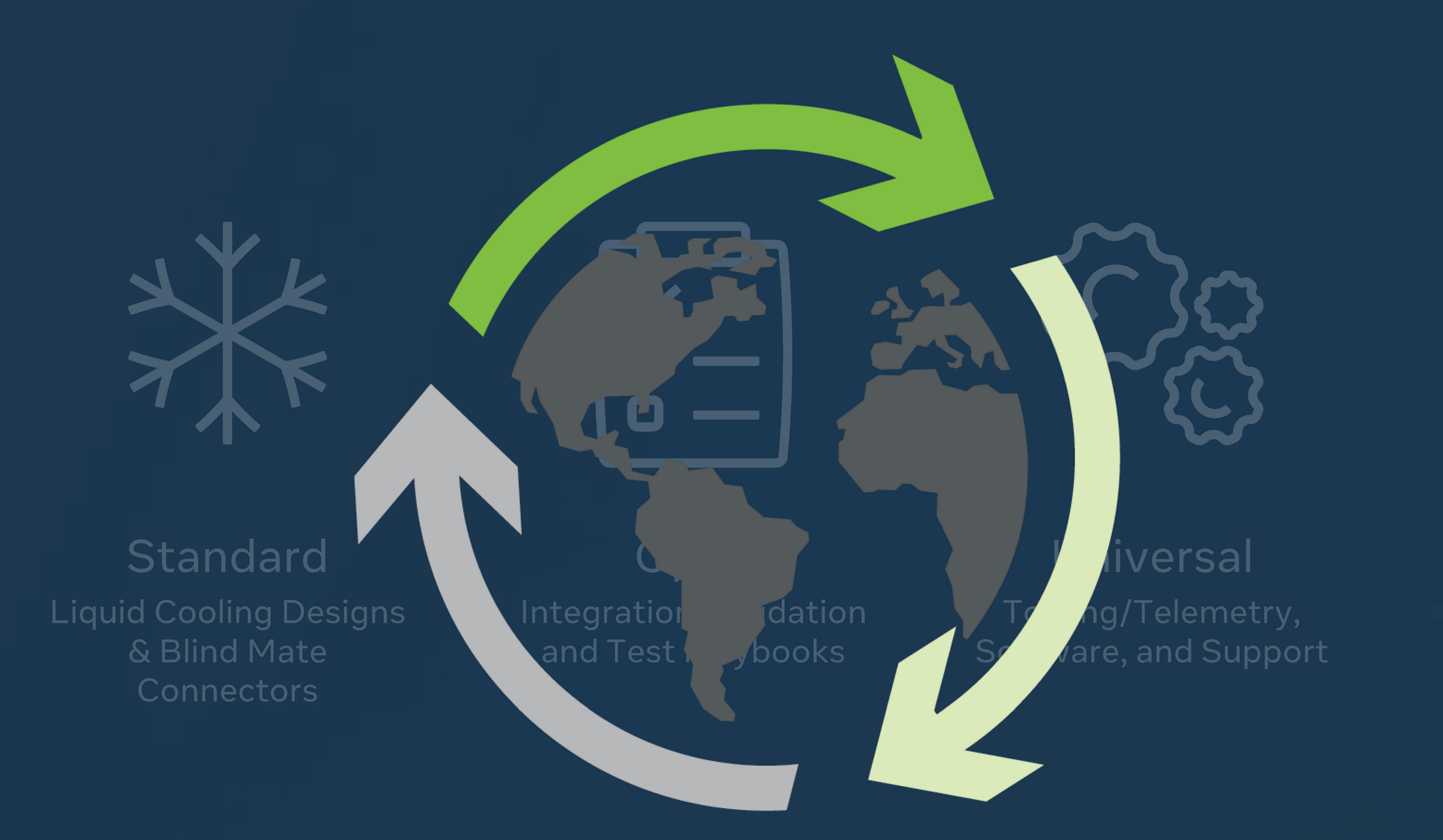


Unified
Open backend Fabric options

Building Open Al Infrastructure



Common
Chassis & Rack
Architecture





Open backend Fabric options

Takeaways







Al Infrastructure complexity is growing

OCP is needed for emerging AI Systems

Al for All requires an open community approach

<a>Meta