# Optimizing ML workloads with AWS Inferentia & Trainium
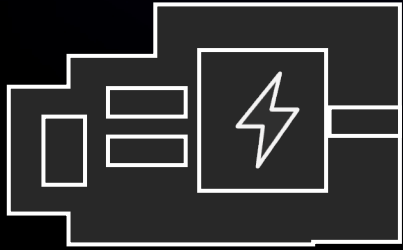
Tobias Edler von Koch

Sr.  Compiler Engineer

AWS

Ron Diamant

Sr. Principal ML Engineer

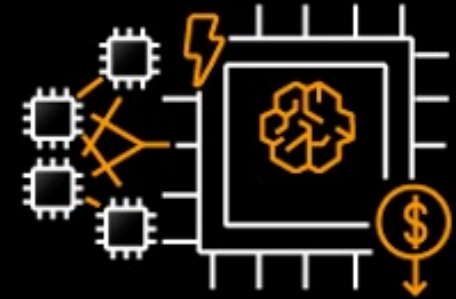AWS

aws

# Silicon innovation at AWS

## AWS Nitro System

Hypervisor, network, storage, SSD, and security

## AWS Graviton

Powerful and efficient, modern applications

## AWS Inferentia and AWS Trainium

Machine learning acceleration

# Silicon innovation at AWS

## AWS Nitro System

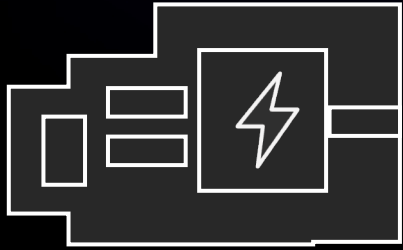Hypervisor, network, storage, SSD, and security

## AWS Graviton

Powerful and efficient, modern applications

## AWS Inferentia and AWS Trainium

Machine learning acceleration

# Silicon innovation at AWS

## AWS Nitro System

Hypervisor, network, storage, SSD, and security

## AWS Graviton

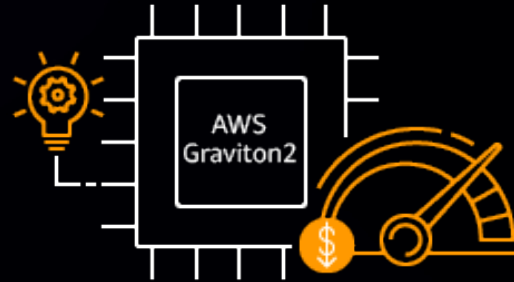Powerful and efficient, modern applications

## AWS Inferentia and AWS Trainium

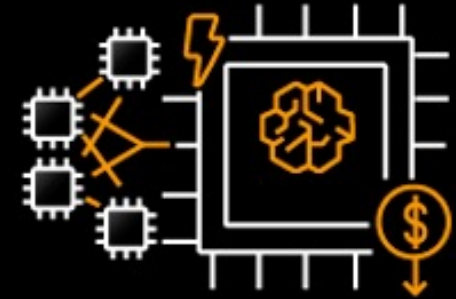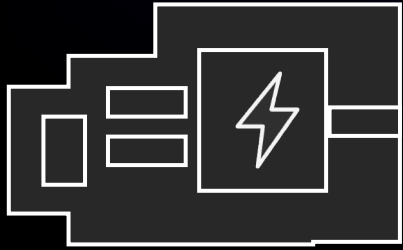Machine learning acceleration

# Silicon innovation at AWS

## AWS Nitro System

Hypervisor, network,
storage, SSD, and security

## AWS Graviton

Powerful and efficient,
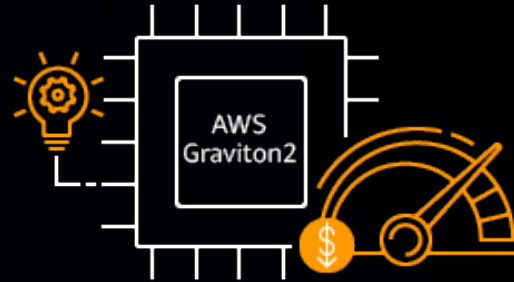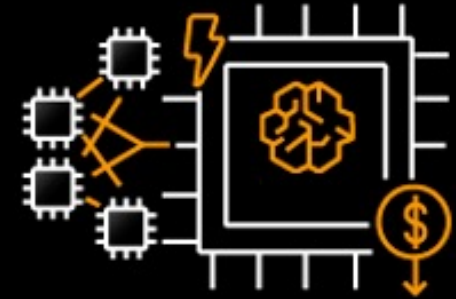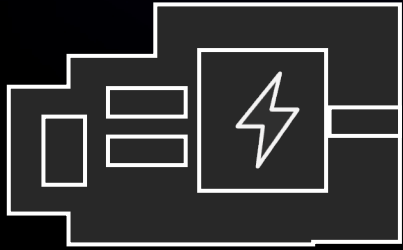modern applications

## AWS Inferentia
and AWS Trainium

Machine learning acceleration

# Machine Learning Acceleration

# Machine Learning trends

GROWTH IN MODEL COMPLEXITY
(# OF PARAMETERS)



| 100M | 1B | 10B | 100B | 1T | 10T |
|------|-----|------|-------|-----|------|
| ELMo (2018) | BERT-Large (2018) | GPT-2 (2019) | Turing NLG (2020) | GPT-3 (2020) | Switch-C (2021) | ... |

# Machine Learning trends

GROWTH IN MODEL COMPLEXITY
(# OF PARAMETERS)



| 100M | 1B | 10B | 100B | 1T | 10T |
|------|-----|------|------|-----|-----|

| ELMo (2018) | BERT-Large (2018) | GPT-2 (2019) | Turing NLG (2020) | GPT-3 (2020) | Switch-C (2021) | … |

COMPUTE-SERVER SPEED

# AWS ML Accelerators for Deep Learning

**AWS Inf1**
Powered by AWS Inferentia

**AWS Trn1**
Powered by AWS Trainium

# AWS Inferentia

- 4 Neuron Cores

- Up to 128 TOPs/chip

- Co-optimize throughput and latency

  - Large on-chip caches

  - Fast chip-to-chip interconnect

- Ease of use!

  - Supported in popular ML frameworks

  - FP16, BF16, INT8

# AWS Inferentia – Sustainable performance

**Performance-per-Watt**



Objects in an image, as detected by YOLOv4

Legend: ▮ G4dn ▮ Inf1

Measured 2-hour run, TensorFlow YOLOv4

Chart (x-axis: p5, p50, p95; y-axis: 0 to 3.5):
- p5: G4dn = 1.0, Inf1 ≈ 3.15, >3x
- p50: G4dn = 1.0, Inf1 ≈ 2.25, >2x
- p95: G4dn = 1.0, Inf1 ≈ 1.65, >1.5x

# AWS Inferentia – Up to 68% lower cost

# Co-optimizing latency and throughput



Latency vs. cost per inference

■ G4dn

# Co-optimizing latency and throughput

**Latency vs. cost per inference**



■ G4dn

# Co-optimizing latency and throughput

Latency vs. cost per inference

■ G4dn  ● Inf1

# AMAZON ALEXA

Alexa has deployed highly complex text-to-speech model that generates **human-like speech**, to support over **100 million Alexa devices globally**

With Inf1 instances, they have been able **to lower their operating costs by about 30%** over GPU instances, while achieving **25% better inference latency**

aws

# AWS Inferentia - Customer adoption

# AWS Inferentia – Customer adoption

**PyTorch BERT Throughput (Chatbot Engine)**



Inf1

G4dn

2x Higher Throughput

**Video Analysis Performance**

Inf1

G4dn

4x Higher Performance

# AWS Trainium



**AWS Trn1/Trn1n**
Powered by AWS Trainium

| Trn1 | MATH ENGINE FREQUENCY 3 GHz |
|---|---|

| BF16/FP16 3.4 PFLOPS | TF32 3.4 PFLOPS | FP32 840 TFLOPS |
|---|---|---|

| AGGREGATE ACCELERATOR MEMORY 512 GB | PEAK MEMORY BANDWIDTH 13.1 TB/sec |
|---|---|

| NEURONLINK BANDWIDTH BETWEEN CHIPS 768 GB/sec | NETWORK CONNECTIVITY 800 Gbps EFA 1600 Gpbs EFA |
|---|---|

# AWS Trainium

2 NeuronCores

Tensor, scalar, and vector engines

Dedicated collective compute engines

Embedded general purpose DSPs

   Support for custom operators

# AWS Trainium

- Rich data-type selection

# AWS Trainium

- Rich data-type selection

- Stochastic rounding



Round nearest even:

Weight update

Mid-point

$W_k$

100% of the times

Next representable FP value

Stochastic rounding:

Weight update

Mid-point

$W_k$

80% of the times

20% of the times

Next representable FP value

# AWS Trainium

- Rich data-type selection

- Stochastic rounding

- High bandwidth,
  Low latency interconnect

**INTER-CHIP INTERCONNECT**

Trn1n
1600 Gb/s
768 GB/s

P4d
400 Gb/s
600 GB/s

P3dn
100 Gb/s
300 GB/s

**NETWORK BANDWIDTH (EFA)**

# AWS Trainium

- Rich data-type selection

- Stochastic rounding

- High bandwidth,
  Low latency interconnect

- Parallelized computation and communication


Accelerator — Shared engine / Shared engine


Computation / Communication / Time

# AWS Trainium

- Rich data-type selection

- Stochastic rounding

- High bandwidth,
  Low latency interconnect

- Parallelized computation and communication

Trainium

Comm engine

Compute engine

Compute engine

Computation

Communication

Time

# AWS Neuron SDK

Supports all major frameworks

Neuron Compiler

Neuron Runtime

Developer tools

https://awsdocs-neuron.readthedocs-hosted.com

github.com/aws/aws-neuron-sdk

# End-to-end flow



User model

Distributed Training Library

ML Framework

FAL (Framework Adaptation Layer)

HLO

Neuron Compiler

NEFF

NEFF

Neuron Runtime

EFA

Neuron Dev Tools

bin

AWS Neuron

# AWS Neuron Compiler

# AWS Neuron Runtime

```
ubuntu@ip-172-31-10-131:~$ lspci
...
00:1c.0 System peripheral: Amazon.com, Inc. Device 7064 (rev 01)
00:1d.0 System peripheral: Amazon.com, Inc. Device 7064 (rev 01)
00:1e.0 System peripheral: Amazon.com, Inc. Device 7064 (rev 01)
00:1f.0 System peripheral: Amazon.com, Inc. Device 7064 (rev 01)

ubuntu@ip-172-31-10-131:~$ sudo neuron-ls
```

| PCI BDF | LOGICAL ID | NEURON CORES | MEMORY CHANNEL 0 | MEMORY CHANNEL 1 | EAST | WEST | RUNTIME ADDRESS | RUNTIME PID | RUNTIME VERSION |
|---|---|---|---|---|---|---|---|---|---|
| 0000:00:1c.0 | 0 | 4 | 4096 MB | 4096 MB | 1 | 0 | unix:/run/neuron.sock | 6311 | 1.0.7875.0 |
| 0000:00:1d.0 | 1 | 4 | 4096 MB | 4096 MB | 1 | 1 | unix:/run/neuron.sock | 6311 | 1.0.7875.0 |
| 0000:00:1e.0 | 2 | 4 | 4096 MB | 4096 MB | 1 | 1 | unix:/run/neuron.sock | 6311 | 1.0.7875.0 |
| 0000:00:1f.0 | 3 | 4 | 4096 MB | 4096 MB | 0 | 1 | unix:/run/neuron.sock | 6311 | 1.0.7875.0 |



Collective Compute

(Topology + Kernels)

# AWS Neuron Profiler

Case study – weight-sharded Transformer:

# AWS Neuron Profiler

Case study – weight-sharded Transformer:

# AWS Neuron Extensions for Training



AWS Neuron

## Framework integration

Full framework integration, JIT, Eager mode, collective compute

## Distributed training

Scale up to 10K+ devices, integration with distributed training libraries, and EFA

## Flexible and extendable

Support for custom ops, dynamic shapes, new data types, and stochastic rounding

## Fully integrated with AWS

SageMaker, EKS, ECS, ParallelCluster, Batch, AMIs

# Case study: BERT-Large pre-training

- Bring your own model

```python
1   import os
2   ...
3   import torch
4   import torch_xla
5   import torch_xla.core.xla_model as xm
6   ...
7   from transformers import BertForPreTraining
8
9   model = BertForPreTraining.from_pretrained('bert-large-uncased')
10
11  def train_loop_fn(model, optimizer, train_loader, device, epoch, global_step, training_ustep, running_loss):
12      max_grad_norm = 1.0
13      for i, data in enumerate(train_loader):
14          training_ustep += 1
15          input_ids, segment_ids, input_mask, masked_lm_labels, next_sentence_labels = data
16          outputs = model(input_ids=input_ids,
17                          attention_mask=input_mask,
18                          token_type_ids=segment_ids,
19                          labels=masked_lm_labels,
20                          next_sentence_label=next_sentence_labels)
21          loss = outputs.loss / flags.grad_accum_usteps
22          loss.backward()
23          running_loss += loss.detach()
24
25          if (training_ustep + 1) % flags.grad_accum_usteps == 0:
26              xm.mark_step()
27              running_loss_cpu = running_loss.detach().cpu().item()
28              running_loss.zero_()
29              torch.nn.utils.clip_grad_norm_(model.parameters(), max_grad_norm)
30              xm.optimizer_step(optimizer)
31              optimizer.zero_grad()
32              scheduler.step()
33              global_step += 1
34              if global_step >= flags.steps_this_run:
35                  break
36
37      return global_step, training_ustep, running_loss
```
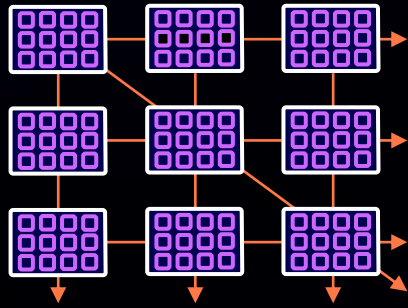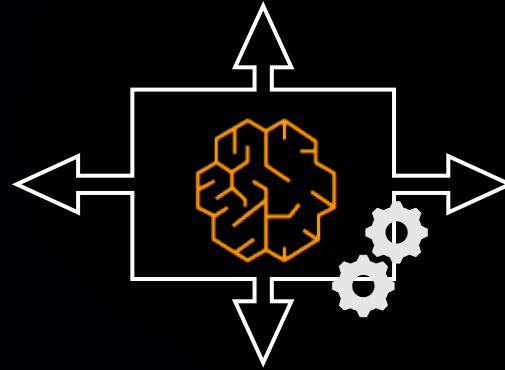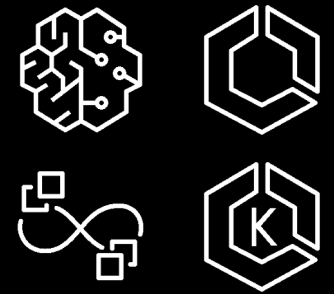
# Case study: BERT-Large pre-training
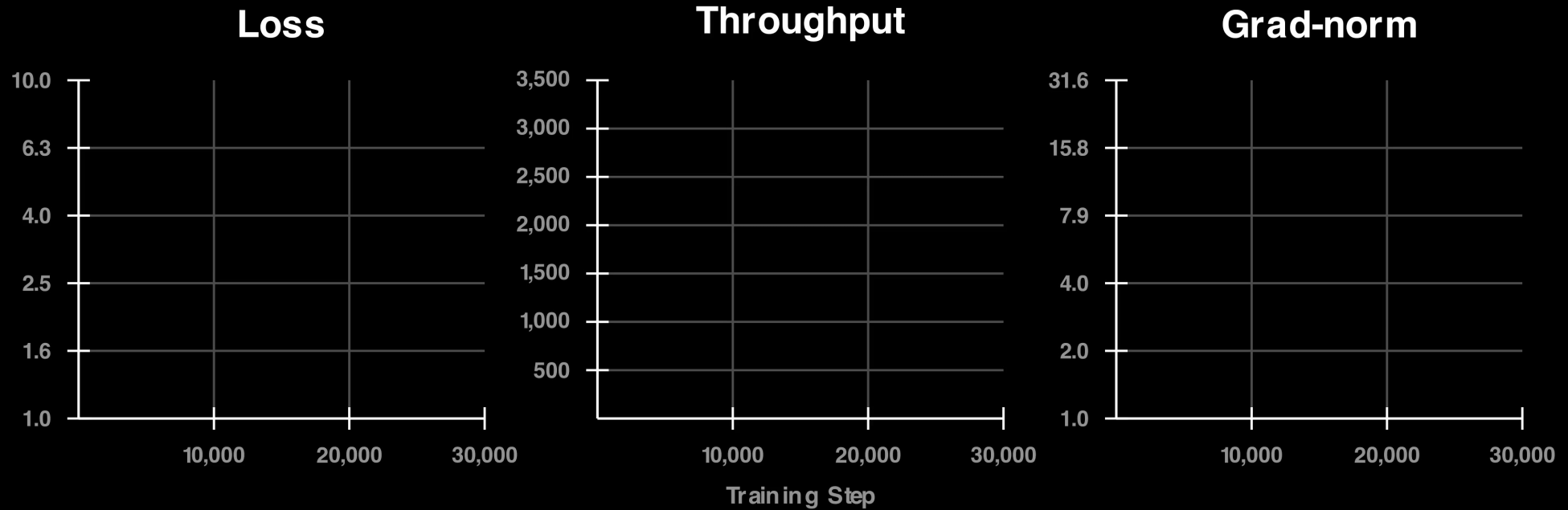
- Bring your own model

- JIT-compile to Trainium

```
1   import os
2   ...
3   import torch
4   import torch_xla
5   import torch_xla.core.xla_model as xm
6   ...
7   from transformers import BertForPreTraining
8
9   model = BertForPreTraining.from_pretrained('bert-large-uncased')
10
11  def train_loop_fn(model, optimizer, train_loader, device, epoch, global_step, training_ustep, running_loss):
12      max_grad_norm = 1.0
13      for i, data in enumerate(train_loader):
14          training_ustep += 1
15          input_ids, segment_ids, input_mask, masked_lm_labels, next_sentence_labels = data
16          outputs = model(input_ids=input_ids,
17                          attention_mask=input_mask,
18                          token_type_ids=segment_ids,
19                          labels=masked_lm_labels,
20                          next_sentence_label=next_sentence_labels)
21          loss = outputs.loss / flags.grad_accum_usteps
22          loss.backward()
23          running_loss += loss.detach()
24
25          if (training_ustep + 1) % flags.grad_accum_usteps == 0:
26              xm.mark_step()
27              running_loss_cpu = running_loss.detach().cpu().item()
28              running_loss.zero_()
29              torch.nn.utils.clip_grad_norm_(model.parameters(), max_grad_norm)
30              xm.optimizer_step(optimizer)
31              optimizer.zero_grad()
32              scheduler.step()
33              global_step += 1
34              if global_step >= flags.steps_this_run:
35                  break
36
37      return global_step, training_ustep, running_loss
```
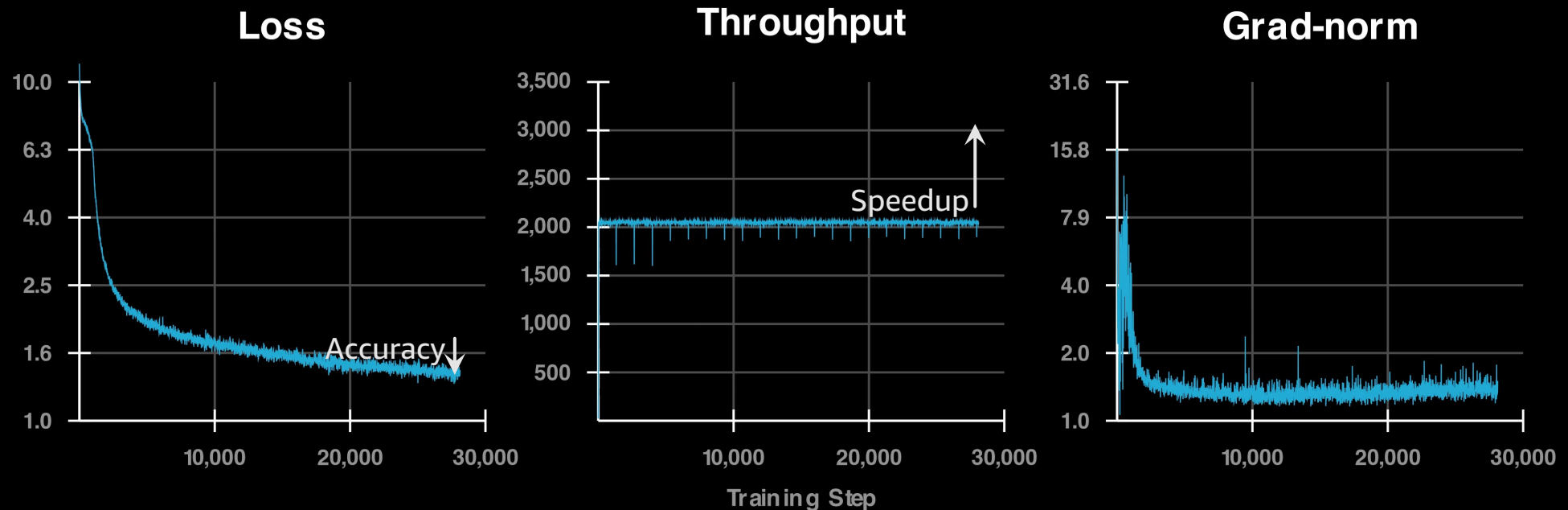
# Case study: BERT-Large pre-training

- Bring your own model

- JIT-compile to Trainium

- See it run ☺



**Loss**

**Throughput**

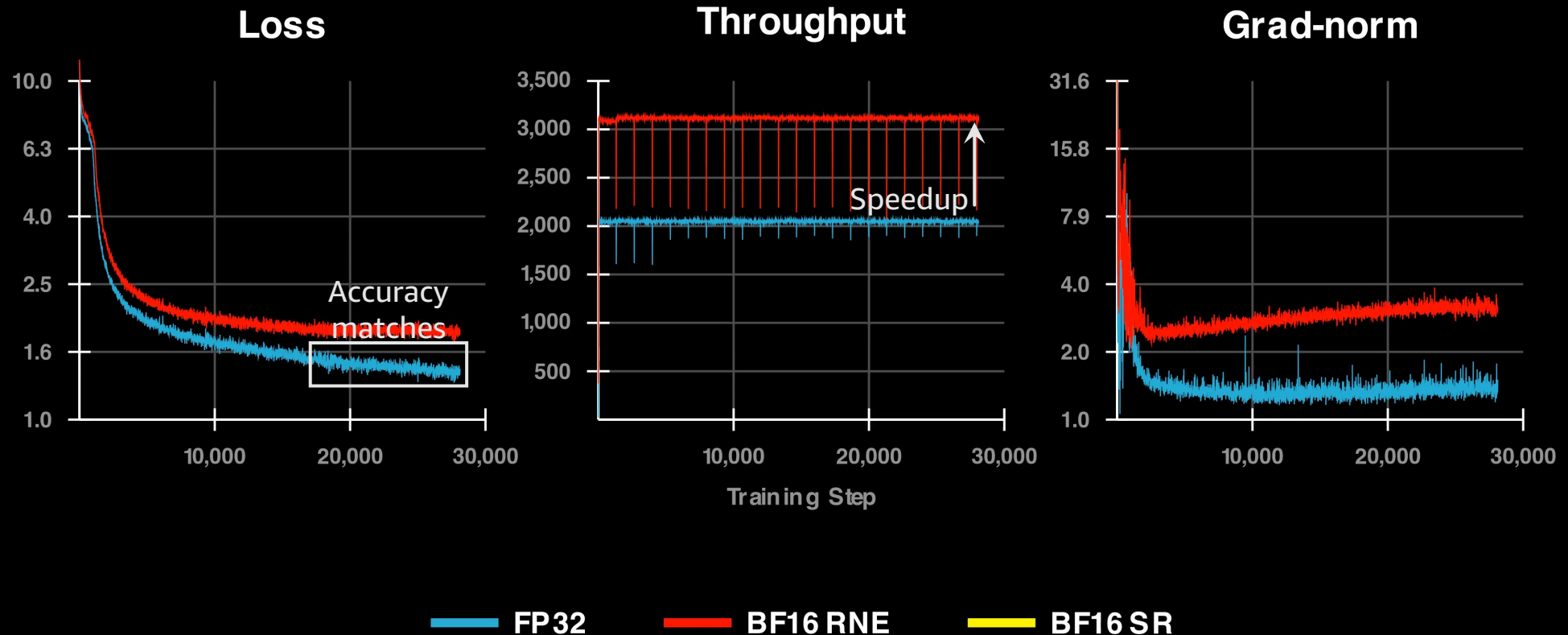**Grad-norm**

Training Step

FP32    BF16 RNE    BF16 SR

# Case study: BERT-Large pre-training

- Bring your own model

- JIT-compile to Trainium

- See it run ☺



**Loss**

**Throughput**

**Grad-norm**

Training Step

FP32 — BF16 RNE — BF16 SR

# Case study: BERT-Large pre-training

- Bring your own model

- JIT-compile to Trainium

- See it run ☺

**Loss**



Accuracy matches

**Throughput**

Speedup

Training Step

**Grad-norm**

━━ **FP32**  ━━ **BF16 RNE**  ━━ **BF16 SR**

# Amazon EC2 Trn1
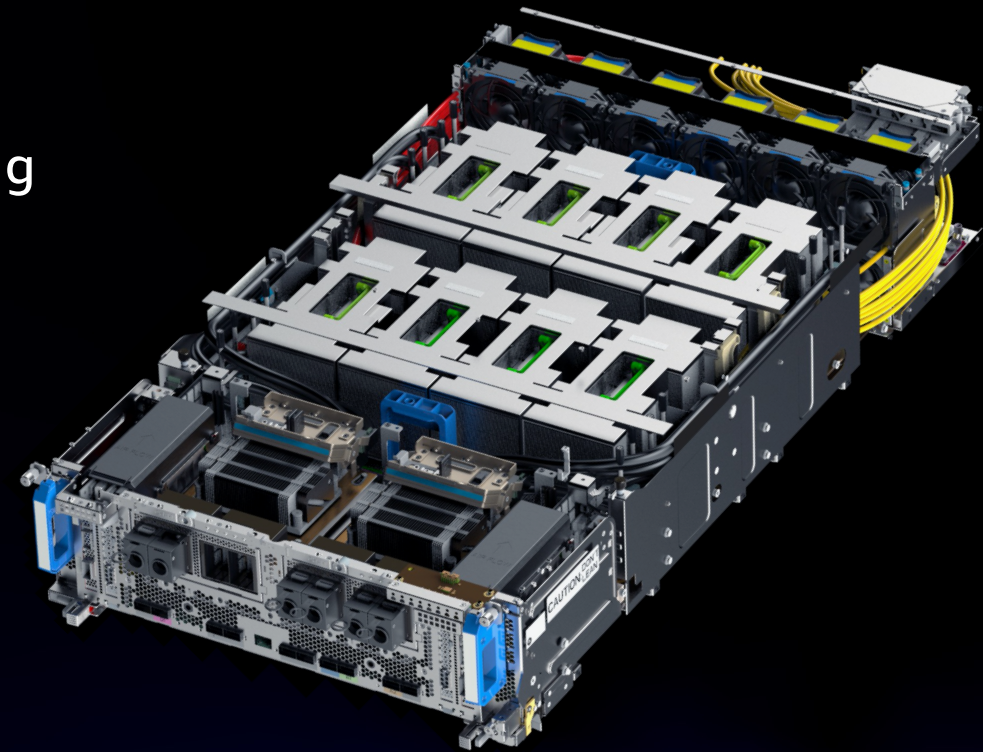
Purposely built for the **most cost-efficient DL training in the cloud** for a broad spectrum of applications

AWS is **innovating across the chips, servers, and data center** layers to provide end users with access to cutting edge hardware on-demand

**Max developer efficiency** with Neuron SDK providing full integration into PyTorch and TensorFlow

**Seamless integration with AWS services** like SageMaker, Amazon ECS, ParallelCluster and more

# Thank you!

We're hiring!    https://www.amazon.jobs/en/landing_pages/annapurna%20labs

Tobias Edler von Koch          Ron Diamant

aws

# Q&A