aws

# Generative AI on AWS

*Self-managed and Managed ML workloads*

Kanwaljit Khurmi

Principal AI/ML Solutions Architect

# GenAI is transforming AI

Image generation, transformation, upscaling



Seamless transformation

**Text to Image:** Generated by Stable Diffusion 2.0 This interior does not exist

**4x**

Upscaling

# Challenges with training large-scale models

**Hardware**

**Health checks**

**Orchestration**

**Data**

**Scaling up**

**Cost**

# AWS AI/ML stack for Self-Managed and Managed Envs

Traditional machine learning (ML)

Training and inference

Deep learning (DL)

Inference

Training

| M6a | M7g | **C5** | C7g | **R5** | R6a | R7g | **Inf1** | **Inf2** | **G5g** | **Trn1** | **P4d P4de** | **DL1** | **and more ..** |

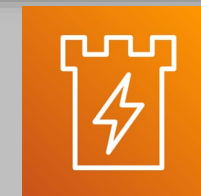intel

AMD
EPYC CPU

annapurnalabs
an amazon company
Graviton CPU
Inferentia and Trainium chip

NVIDIA®
A100, V100, T4 GPUs

AWS Nitro Enclaves

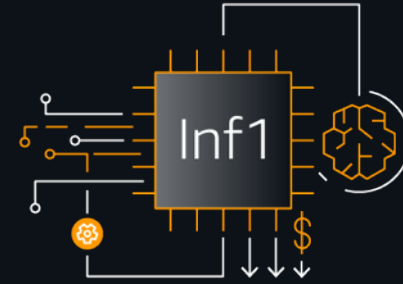Elastic Fabric Adapter

FSx
Amazon FSx for Lustre

Broadest and deepest compute for AI/ML

**AWS Inferentia1**

Inf1

Up to 70% lower cost per inference

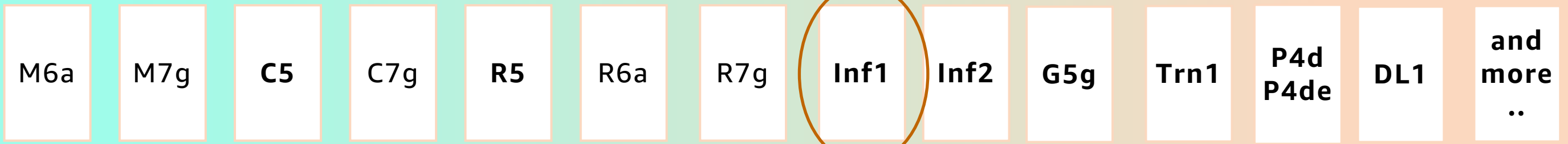Lowest cost inference in the cloud for running deep learning models

Traditional machine learning (ML)

Training and inference

Deep learning (DL)

Inference

Training

| M6a | M7g | **C5** | C7g | **R5** | R6a | R7g | **Inf1** | **Inf2** | **G5g** | **Trn1** | **P4d P4de** | **DL1** | **and more ..** |

intel

**AMD**
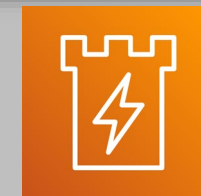EPYC CPU

**annapurnalabs** an **amazon** company
Graviton CPU
Inferentia and Trainium chip

**NVIDIA.**
A100, V100, T4 GPUs
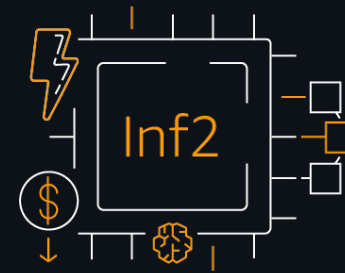
AWS Nitro Enclaves

Elastic Fabric Adapter

FSx
Amazon FSx for Lustre

Broadest and deepest compute for AI/ML

**AWS Inferentia2**

Inf2

Up to 40% better price performance for Generative AI

High performance at the lowest cost per inference for LLMs and diffusion models

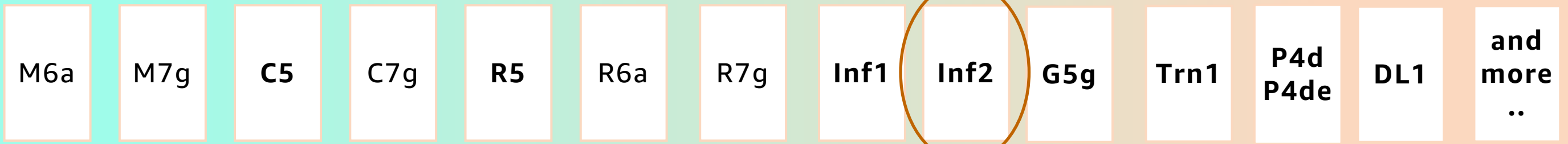Traditional machine learning (ML)

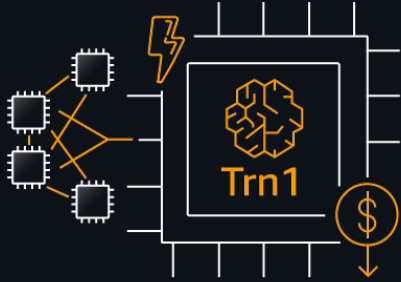Training and inference

Deep learning (DL)

Inference

Training

| M6a | M7g | **C5** | C7g | **R5** | R6a | R7g | **Inf1** | **Inf2** | **G5g** | **Trn1** | **P4d P4de** | **DL1** | **and more ..** |

intel®

**AMD**
EPYC CPU

annapurnalabs
an amazon company
Graviton CPU
Inferentia and Trainium chip

**NVIDIA**®
A100, V100, T4 GPUs

AWS Nitro Enclaves

Elastic Fabric Adapter

FSx
Amazon FSx for Lustre

Broadest and deepest compute for AI/ML

**AWS Trainium**

Trn1

Up to 50% cost-to-train savings

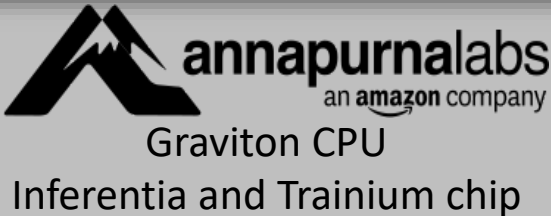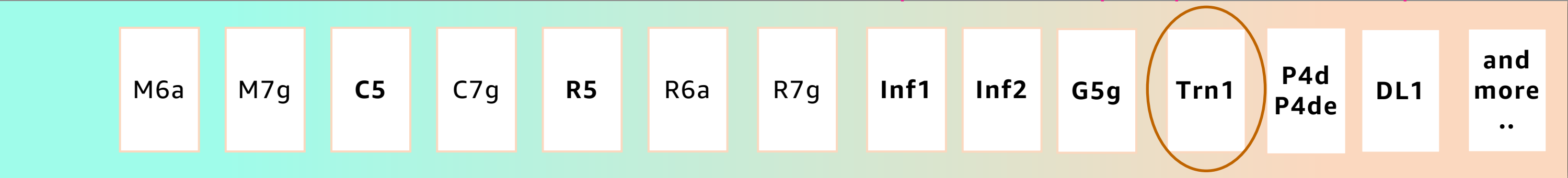The most cost efficient for high-performance training of LLMs and diffusion models

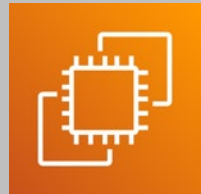Traditional machine learning (ML)

Training and inference

Deep learning (DL)

Inference

Training

| M6a | M7g | **C5** | C7g | **R5** | R6a | R7g | **Inf1** | **Inf2** | **G5g** | **Trn1** | **P4d P4de** | **DL1** | **and more ..** |

(intel)

**AMD**
EPYC CPU

annapurnalabs
an amazon company
Graviton CPU
Inferentia and Trainium chip

NVIDIA.
A100, V100, T4 GPUs

AWS Nitro Enclaves

Elastic Fabric Adapter

FSx
Amazon FSx for Lustre

Broadest and deepest compute for AI/ML

| Amazon EC2 | AWS ParallelCluster | Amazon EKS | Amazon ECS | AWS Fargate | more …. | Amazon SageMaker |

Self-Managed → Managed

**Broadest and deepest compute services for AI/ML**

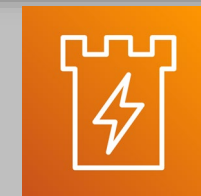| M6a | M7g | **C5** | C7g | **R5** | R6a | R7g | **Inf1** | **Inf2** | **G5g** | **Trn1** | **P4d P4de** | **DL1** | **and more ..** |

intel

**AMD**
EPYC CPU

annapurnalabs
an amazon company
Graviton CPU
Inferentia and Trainium chip

**NVIDIA.**
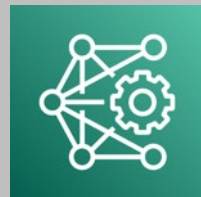A100, V100, T4 GPUs

AWS Nitro Enclaves
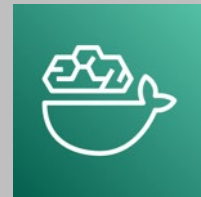
Elastic Fabric Adapter
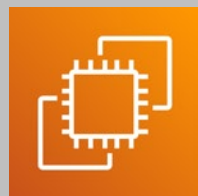
FSx
Amazon FSx for Lustre

Broadest and deepest compute for AI/ML

AWS Deep Learning AMIs
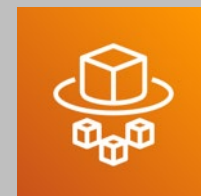
AWS Deep Learning Containers

Amazon EC2

AWS ParallelCluster

Amazon EKS

Amazon ECS

AWS Fargate

more ....

Amazon SageMaker

Self-Managed → **Broadest and deepest compute services for AI/ML** → Managed

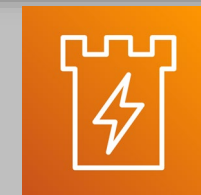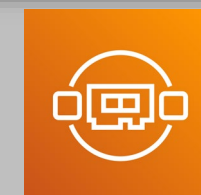| M6a | M7g | **C5** | C7g | **R5** | R6a | R7g | **Inf1** | **Inf2** | **G5g** | **Trn1** | **P4d P4de** | **DL1** | **and more ..** |

intel

AMD
EPYC CPU

annapurnalabs
an amazon company
Graviton CPU
Inferentia and Trainium chip
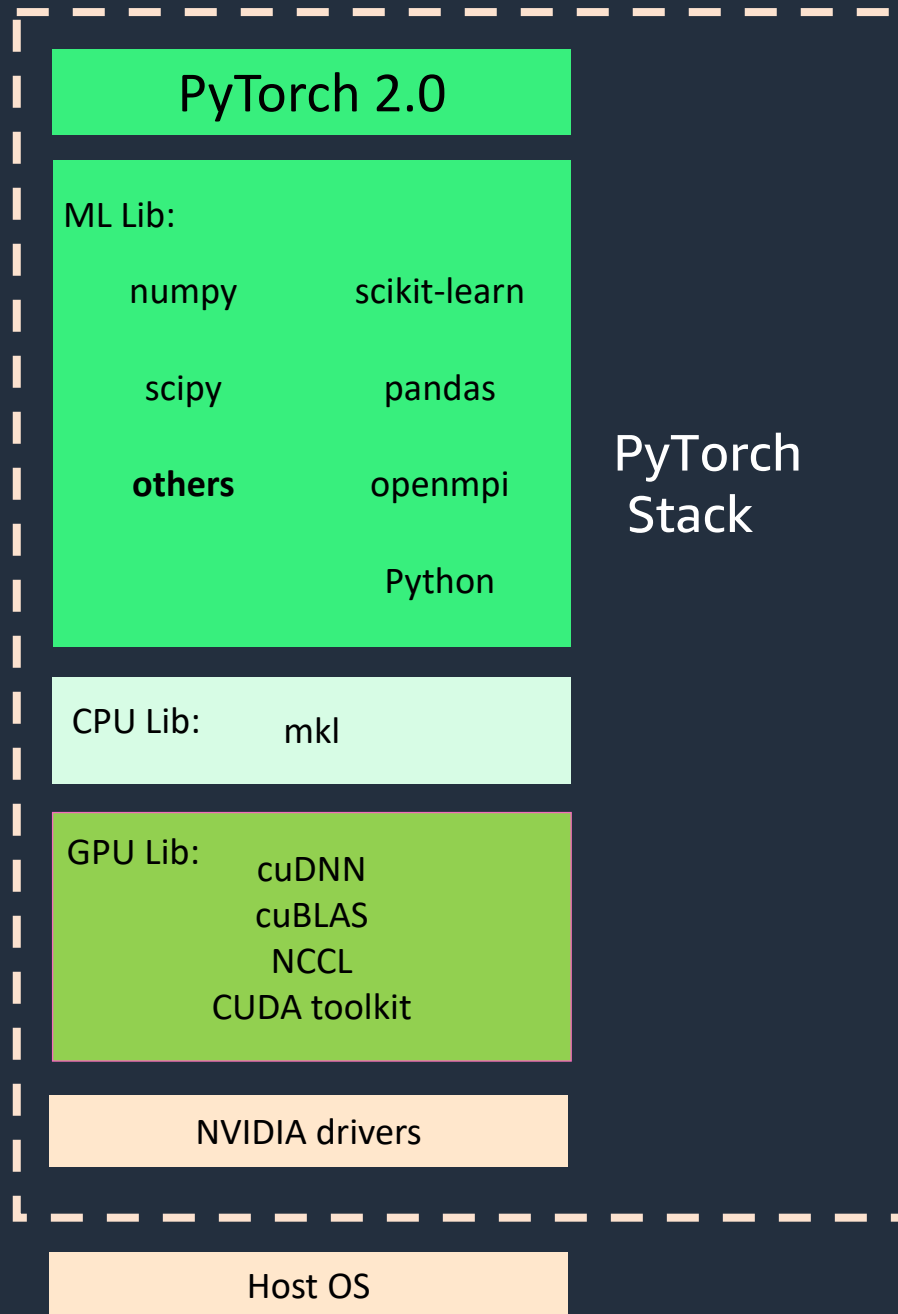
NVIDIA®
A100, V100, T4 GPUs

AWS Nitro Enclaves

Elastic Fabric Adapter

FSx
Amazon FSx for Lustre

Broadest and deepest compute for AI/ML

# Why use DLAMI and DLC ?



PyTorch 2.0

ML Lib:

numpy          scikit-learn

scipy          pandas

**others**         openmpi

Python

PyTorch Stack

CPU Lib:    mkl

GPU Lib:    cuDNN
            cuBLAS
            NCCL
            CUDA toolkit

NVIDIA drivers

Host OS

Build on your own

OR

AWS Managed

**Use AWS DLAMI**

**Use AWS DLC**

Host OS

**Pre-packaged**

**Easy to use**

**Higher performance**

**Zero cost**

# Train and deploy Stable Diffusion model using Hybrid Archtecture

```
Train  →  Inference on Inf2 or G5
```

**ML FRAMEWORKS**

AWS Deep Learning AMIs

AWS Deep Learning Containers

AWS Deep Learning Containers

**ML SERVICES**

Amazon EC2

Amazon SageMaker

Amazon SageMaker

**ML INFRASTRUCTURE**

P4d

G5

Inf2

aws

# Demo

# References

- AWS Neuron Documentation: https://awsdocs-neuron.readthedocs-hosted.com/en/latest/index.html

- Hugging Face Optimum Neuron: https://huggingface.co/docs/optimum-neuron/index

- AWS DLC: https://docs.aws.amazon.com/deep-learning-containers/

- AWS DLAMI: https://docs.aws.amazon.com/dlami/

aws

# Thank you!

Kanwaljit Khurmi

kkhurmi@amazon.com