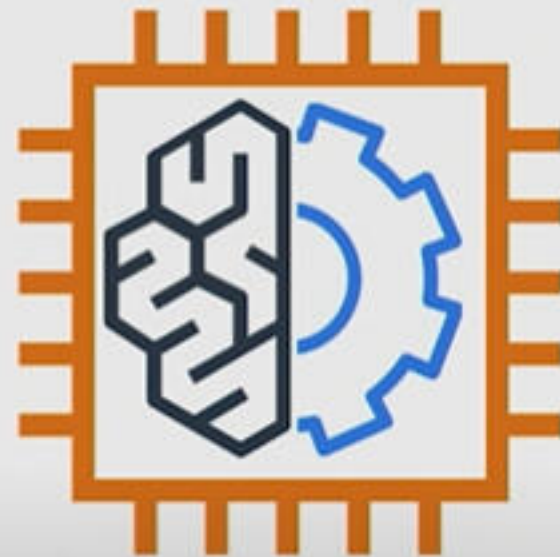


# Introduction to AWS Inferentia and Amazon EC2 Inf1 instances

Vibhav Viswanathan

# In this video

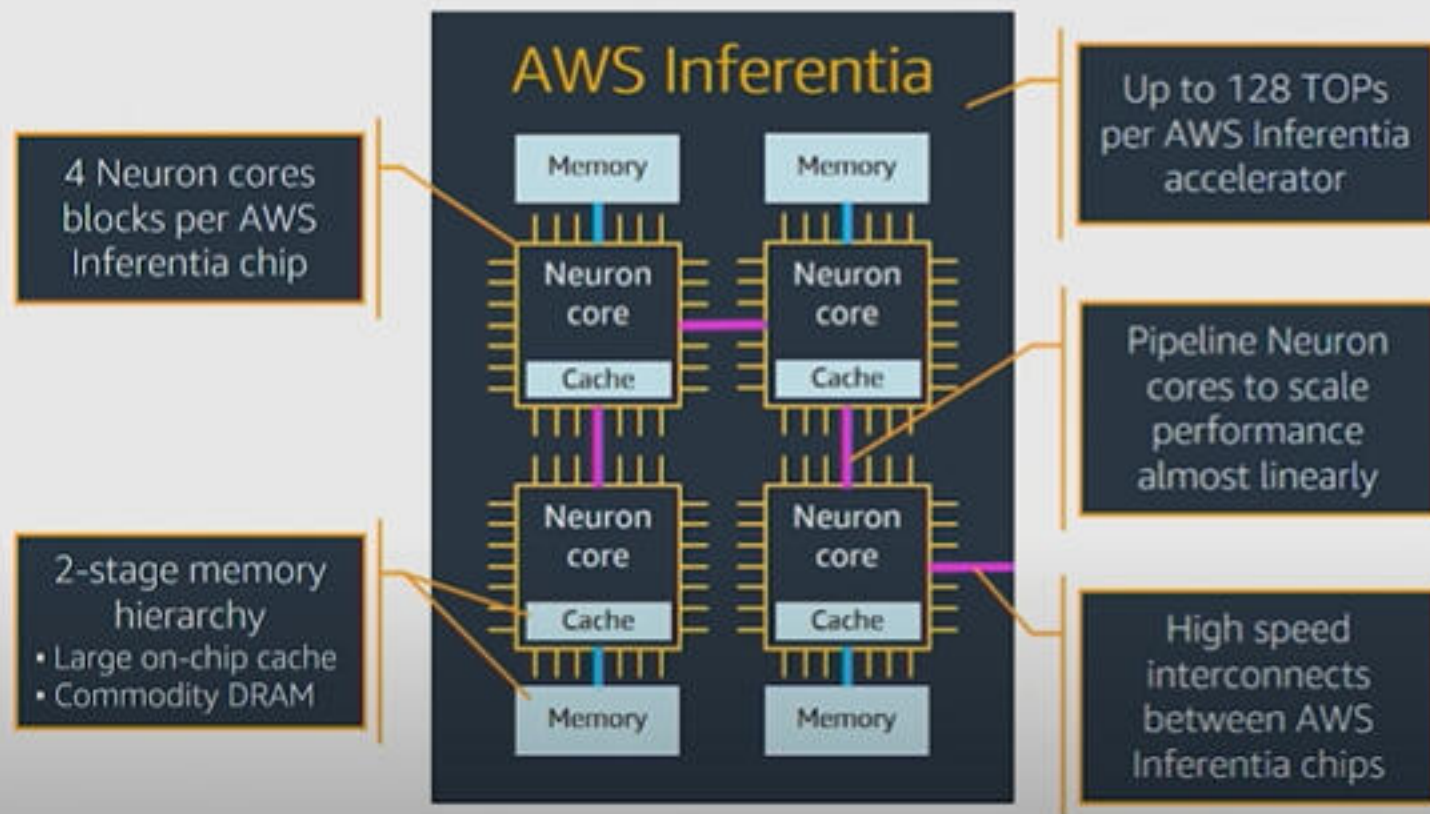
- Machine learning use cases and challenges
- AWS Inferentia overview
- Amazon EC2 Inf1 instances introduction
- Benefits
- Pricing and availability



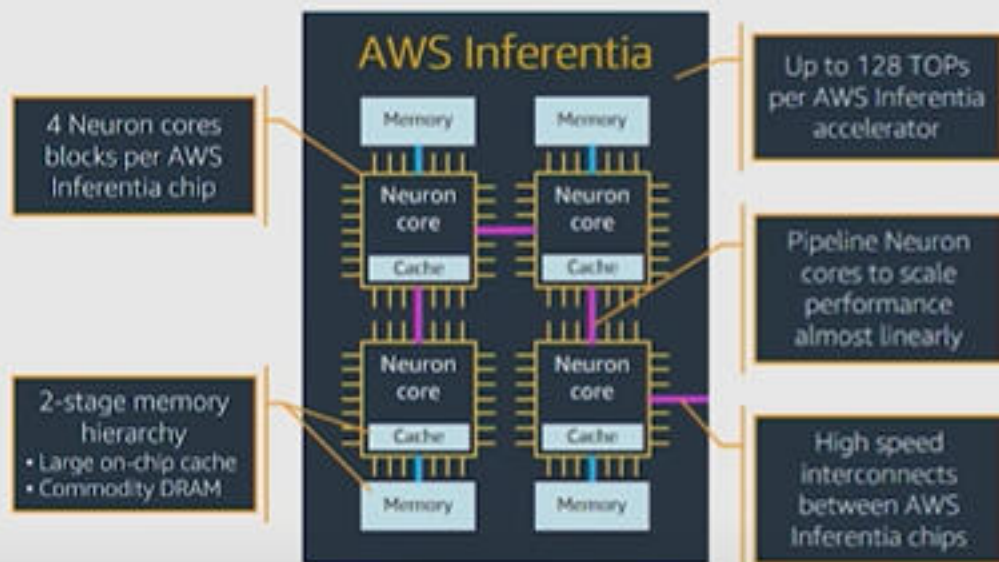
# AWS Inferentia



# AWS Inferentia accelerators



# AWS Inferentia accelerators



## Data types

- INT8
- BFloat16
- FP16 with mixed precision

## ML frameworks

- TensorFlow
  - MXNet
  - PyTorch
-





Build



Neuron Compiler  
(NCC)



Deploy

Neuron Runtime  
(NRT)



Neuron Binary  
(NBIN)



Debug/profile

Neuron Debugger/  
Profiler

```
C:\>code --version  
1.1.1
```



# AWS Neuron software



Build



Neuron Compiler  
(NCC)



Deploy

Neuron Runtime  
(NRT)



Neuron Binary  
(NBIN)



Debug/profile

Neuron Debugger/  
Profiler

```
C:\>code --version  
1.1.1
```



# AWS Neuron software



Build



Neuron Compiler  
(NCC)



Deploy

Neuron Runtime  
(NRT)



Neuron Binary  
(NBIN)



Debug/profile

Neuron Debugger/  
Profiler

```
C:\>code --version  
1.1.1
```





# AWS Neuron software



Build



Neuron Compiler  
(NCC)



Deploy

Neuron Runtime  
(NRT)



Neuron Binary  
(NBIN)



Debug/profile

Neuron Debugger/  
Profiler

```
C:\>code --version  
1.1.1
```



# AWS Neuron software



Build



Neuron Compiler  
(NCC)



Deploy

Neuron Runtime  
(NRT)



Neuron Binary  
(NBIN)



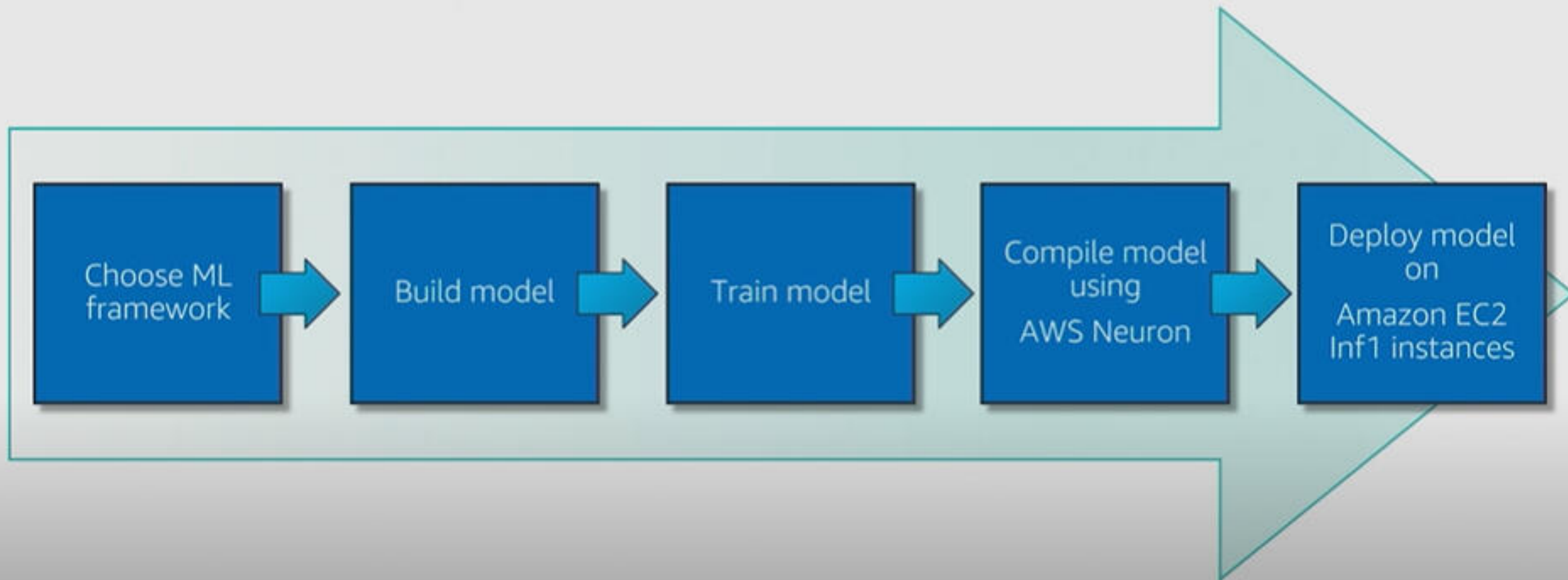
Debug/profile

Neuron Debugger/  
Profiler

```
C:\>code --version  
1.1.1
```



# AWS Neuron developer flow



# AWS Neuron highlights



Smart  
partitioning

Automatically optimize  
neural-net compute



Auto FP32  
casting

Ingest FP32 trained  
model, and cast to  
BFloat16



AWS Neuron  
Core pipeline

Super low-latency  
full bandwidth



AWS Neuron  
Core groups

Concurrently run  
multiple models



# AWS Neuron highlights



Smart  
partitioning

Automatically optimize  
neural-net compute



Auto FP32  
casting

Ingest FP32 trained  
model, and cast to  
BFloat16



AWS Neuron  
Core pipeline

Super low-latency  
full bandwidth



AWS Neuron  
Core groups

Concurrently run  
multiple models



# AWS Neuron highlights



Smart  
partitioning

Automatically optimize  
neural-net compute



Auto FP32  
casting

Ingest FP32 trained  
model, and cast to  
BFloat16



AWS Neuron  
Core pipeline

Super low-latency  
full bandwidth



AWS Neuron  
Core groups

Concurrently run  
multiple models

# Amazon EC2 Inf1 instances powered by AWS Inferentia



# Amazon EC2 Inf1 instances use:

- Custom Intel Cascade-Lake processors
- Custom AWS Inferentia ML acceleration chips
- AWS Nitro up to 100-Gbps networking

# Amazon EC2 Inf1 instance benefits



Low cost



High performance



Increased agility to meet ML needs





Up to 2,000 tera operations per second (TOPs)



# Amazon EC2 Inf1 instance benefits



Low cost



High performance



Increased agility to  
meet ML needs

# Amazon EC2 Inf1 instance benefits



Low cost



High performance



Increased agility to  
meet ML needs

# Instance sizes

Instance size	vCPUs	Memory (GiB)	AWS Inferentia chips	Chip-to-chip interconnect	Storage	Network bandwidth	Amazon EBS bandwidth
Inf1.xlarge	4	8	1	N/A	EBS only	Up to 25 Gbps	Up to 3.5 Gbps
Inf1.2xlarge	8	16	1	N/A	EBS only	Up to 25 Gbps	Up to 3.5 Gbps
Inf1.6xlarge	24	48	4	Yes	EBS only	25 Gbps	3.5 Gbps
Inf1.24xlarge	96	192	16	Yes	EBS only	100 Gbps	14 Gbps

# Instance sizes

Instance size	vCPUs	Memory (GiB)	AWS Inferentia chips	Chip-to-chip interconnect	Storage	Network bandwidth	Amazon EBS bandwidth
Inf1.xlarge	4	8	1	N/A	EBS only	Up to 25 Gbps	Up to 3.5 Gbps
Inf1.2xlarge	8	16	1	N/A	EBS only	Up to 25 Gbps	Up to 3.5 Gbps
Inf1.6xlarge	24	48	4	Yes	EBS only	25 Gbps	3.5 Gbps
Inf1.24xlarge	96	192	16	Yes	EBS only	100 Gbps	14 Gbps

# Instance sizes

Instance size	vCPUs	Memory (GiB)	AWS Inferentia chips	Chip-to-chip interconnect	Storage	Network bandwidth	Amazon EBS bandwidth
Inf1.xlarge	4	8	1	N/A	EBS only	Up to 25 Gbps	Up to 3.5 Gbps
Inf1.2xlarge	8	16	1	N/A	EBS only	Up to 25 Gbps	Up to 3.5 Gbps
Inf1.6xlarge	24	48	4	Yes	EBS only	25 Gbps	3.5 Gbps
Inf1.24xlarge	96	192	16	Yes	EBS only	100 Gbps	14 Gbps

AWS Inferentia  
chip-to-chip fast  
interconnects







## AWS Deep Learning Amazon EC2 AMIs

AWS Deep Learning AMIs

---



Amazon EKS



Amazon ECS

## AWS Deep Learning Containers

---



Amazon SageMaker

Amazon SageMaker  
managed ML service



## AWS Deep Learning Amazon EC2 AMIs

AWS Deep Learning AMIs

---



Amazon EKS



Amazon ECS

## AWS Deep Learning Containers

---



Amazon SageMaker  
managed ML service

Amazon SageMaker



## AWS Deep Learning Amazon EC2 AMIs

AWS Deep Learning AMIs

---



Amazon EKS



Amazon ECS

## AWS Deep Learning Containers

---



Amazon SageMaker  
managed ML service

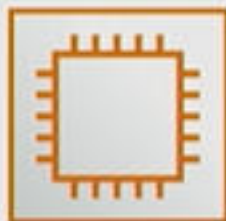
Amazon SageMaker

## Pricing

- Supports all Amazon EC2 pricing models



On-Demand instances



Reserved instances



Spot instances

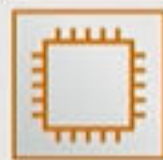
# Pricing models and availability

## Pricing

- Supports all Amazon EC2 pricing models



On-Demand instances



Reserved instances



Spot instances

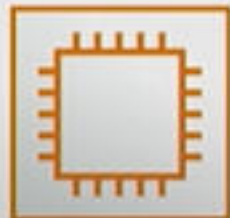


## Pricing

- Supports all Amazon EC2 pricing models



On-Demand instances



Reserved instances



Spot instances

## Availability

- Initial launch December 2019
- Deployment to all AWS Regions





© 2019 Amazon Web Services, Inc. or its affiliates. All rights reserved. This work may not be reproduced or redistributed, in whole or in part, without prior written permission from Amazon Web Services, Inc. Commercial copying, lending, or selling is prohibited. Corrections or feedback on the course, please email us at: [aws-course-feedback@amazon.com](mailto:aws-course-feedback@amazon.com). For all other questions, contact us at: <https://aws.amazon.com/contact-us/aws-training/>. All trademarks are the property of their owners.



**Certificate of Completion**  
**Hem Bahadur Gurung**

**Has successfully completed**  
**Introduction to AWS Inferentia and Amazon EC2 Inf1 Instances**

A handwritten signature in black ink, appearing to read 'Maurice Jorgensen'.

**Director, Training and Certification**

**15 minutes**

**Duration**

**10 September, 2021**

**Completion Date**



**Certificate of Completion**  
**Hem Bahadur Gurung**

**Has successfully completed**  
**Introduction to AWS Inferentia and Amazon EC2 Inf1 Instances**

A handwritten signature in black ink, which appears to read 'Maurice Jorgensen'.

**Director, Training and Certification**

**15 minutes**

**Duration**

**10 September, 2021**

**Completion Date**