AWS
re:Invent

AIM404-R

# Amazon SageMaker RL: Solving business problems with RL and bandits

**Girish Dilip Patil**

Senior Architect
Amazon Web Services India

**Marc Cabocel**

Senior Architect
Amazon Web Services France

**Segolene Dessertine-panhard**

Data Scientist
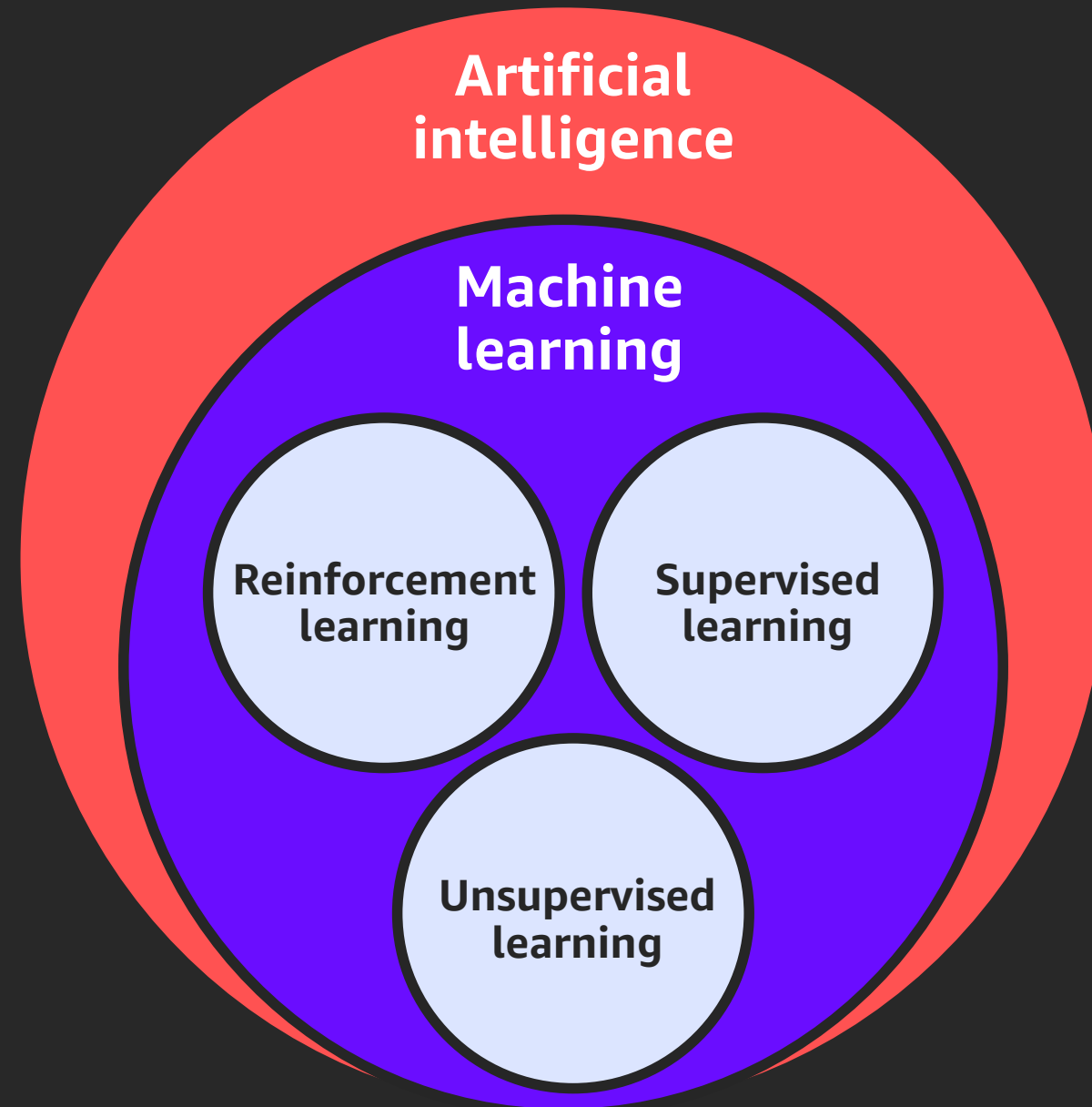Amazon Web Services France

**Anna Luo**

Applied Scientist
Amazon Web Services USA

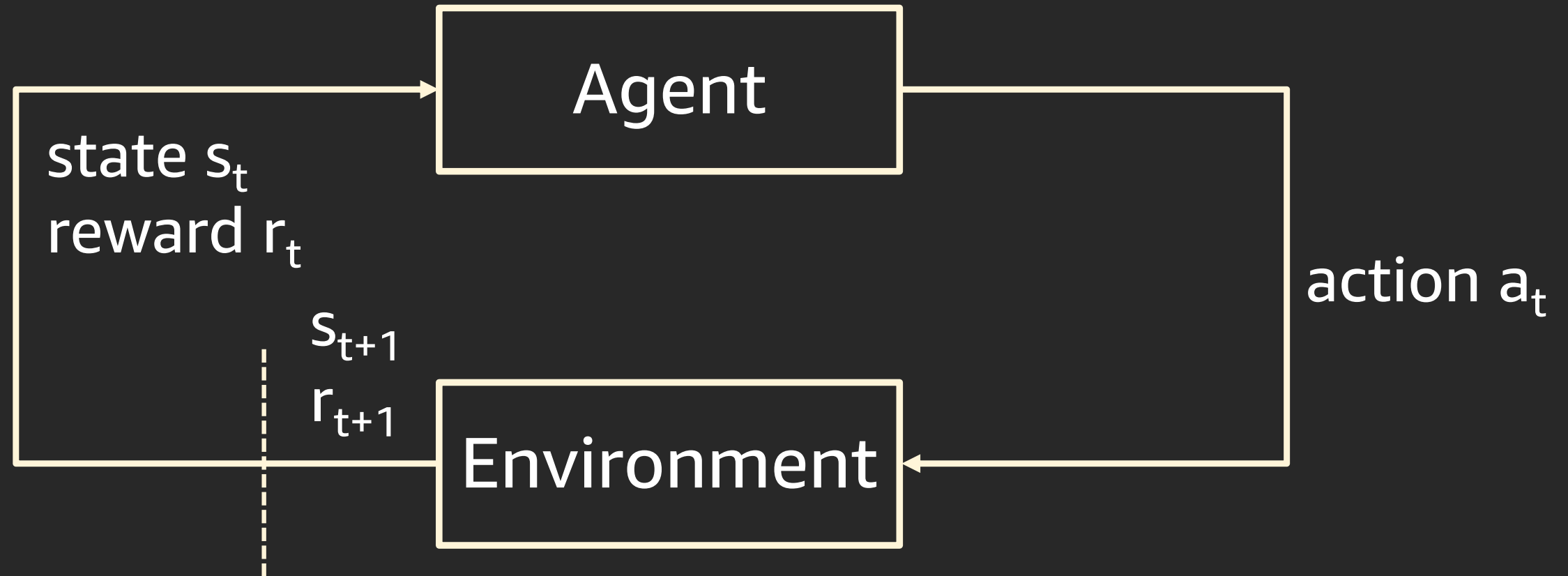AWS re:Invent

aws

# Agenda

1. A quick primer on reinforcement learning (RL)

2. An important trade-off: Explore vs. exploit

3. Amazon SageMaker RL

4. Workshop #1: Training without a simulator in a real environment

5. Workshop #2: Training with a simulator

6. Conclusion

# A quick primer on reinforcement learning

# Reinforcement learning in the broader artificial intelligence context

# Reinforcement learning



Agent

Environment

state $s_t$
reward $r_t$
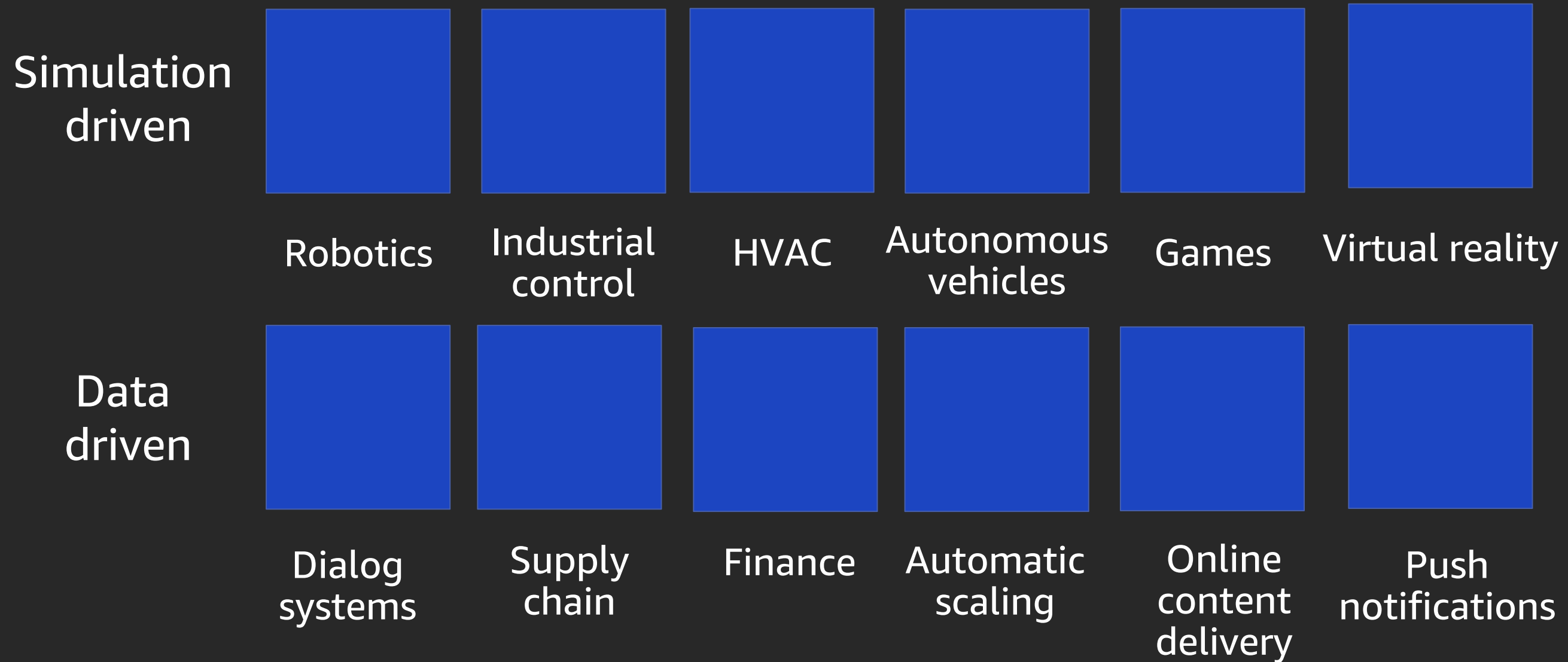
$s_{t+1}$
$r_{t+1}$

action $a_t$

Reinforcement learning is based on the reward hypothesis:

All goals can be described by the maximization of expected cumulative reward

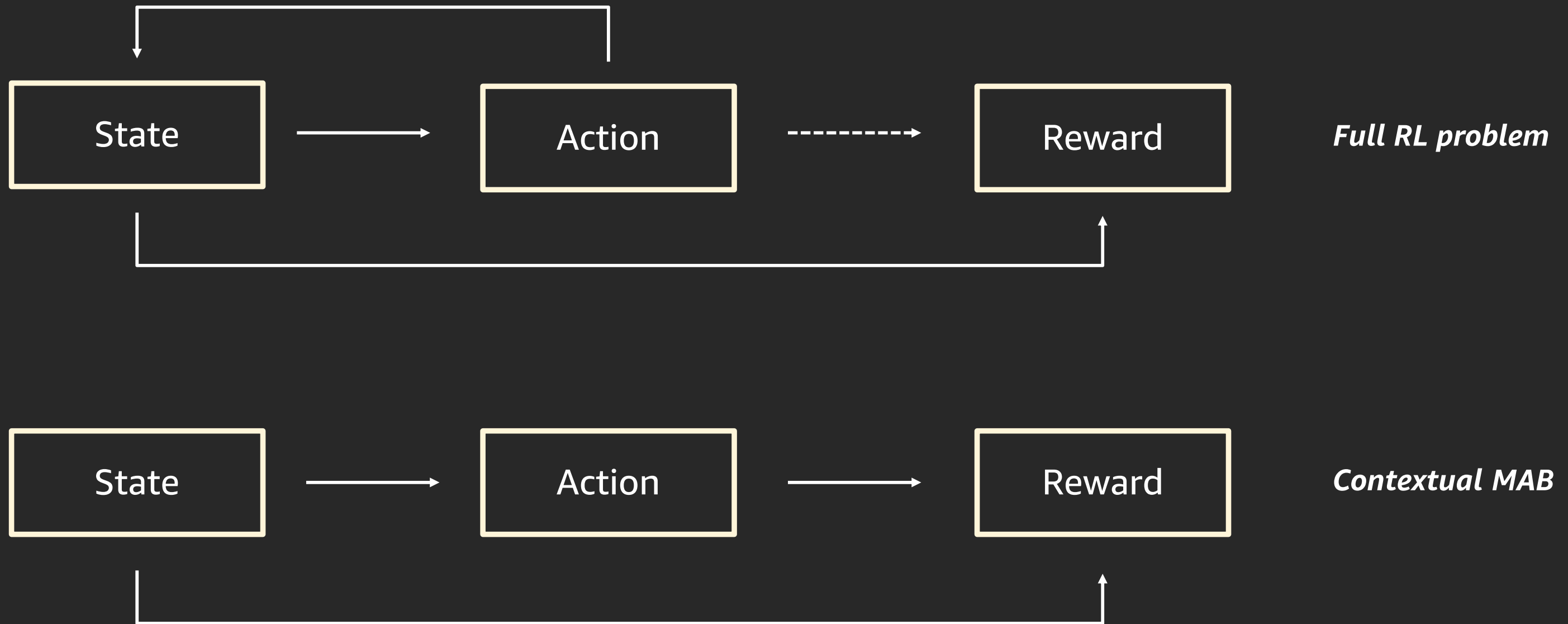# How RL differs from other variations of machine learning

o Reinforcement learning helps in learning a strategy to maximize a reward in a specific environment

o Very useful when you don't have supervised training data

o Agent learns by interacting with the environment (simulated or real)

# RL is applicable in many domains

**Simulation driven**

| Robotics | Industrial control | HVAC | Autonomous vehicles | Games | Virtual reality |

**Data driven**

| Dialog systems | Supply chain | Finance | Automatic scaling | Online content delivery | Push notifications |

# First step toward RL: Contextual multi-armed bandits

# An important trade-off: Explore vs. exploit

You need to have a balance between exploration and exploitation

# Amazon SageMaker RL

# Amazon SageMaker RL makes RL accessible

| Difficult to get started | RL agent algorithms are complex to implement | Hard to integrate environments for training | Training is computationally expensive and time-consuming | Requires trial and error & frequent tuning of hyperparameters |
|---|---|---|---|---|
| Pre-built environments for RL; numerous examples | Support for RL agent algorithms | Easy to integrate variety of simulation environments | Single/ distributed training; local/ remote environment | Local mode for debugging; automatic model tuning |

# Train RL models using state-of-the-art algorithms

## RL Toolkits that provide RL agent algorithm implementations

### RL-Coach

| DQN | PPO | HER | Rainbow | ... |

### RL-Ray RLLib

| APEX | ES | IMPALA | A3C | ... |

### Open AI Baselines

| TRPO | GAIL | ... |

...

## Amazon SageMaker Deep Learning Frameworks

| TensorFlow | MxNet | PyTorch | Chainer |

■ Amazon SageMaker supported    □ Customer BYO

\* RL Toolkits comparison

# Integrate any type of RL environment

# Amazon SageMaker RL

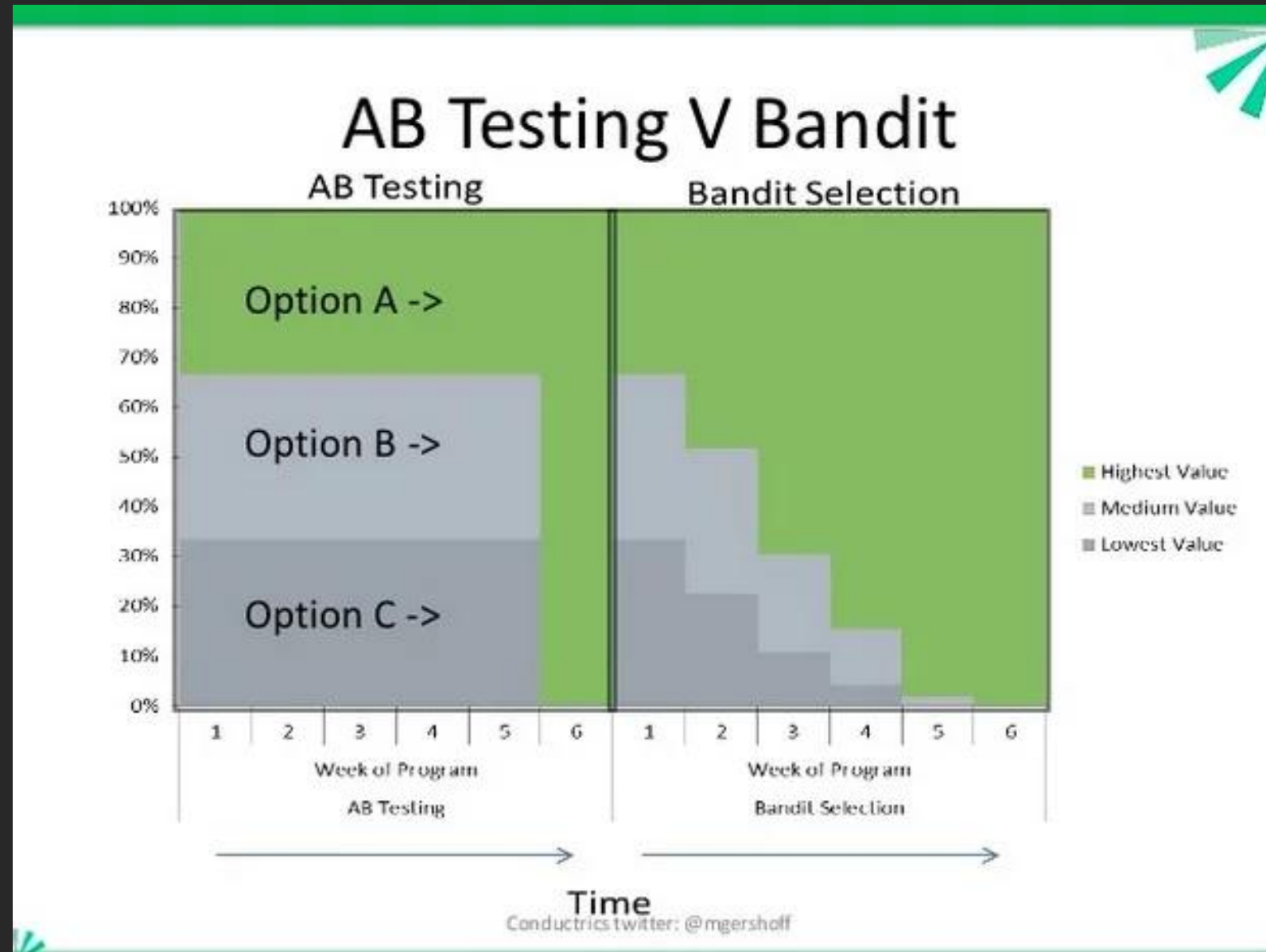# Workshop #1: Training without a simulator in a real environment

# What are the challenges

o   Feedback is delayed. It needs to be joined with inputs & actions taken to prepare next training datasets.

o   You have to learn fast. Unlike in a simulated environment, the agent doesn't have the luxury to learn from millions of episodes.

o   Training never stops
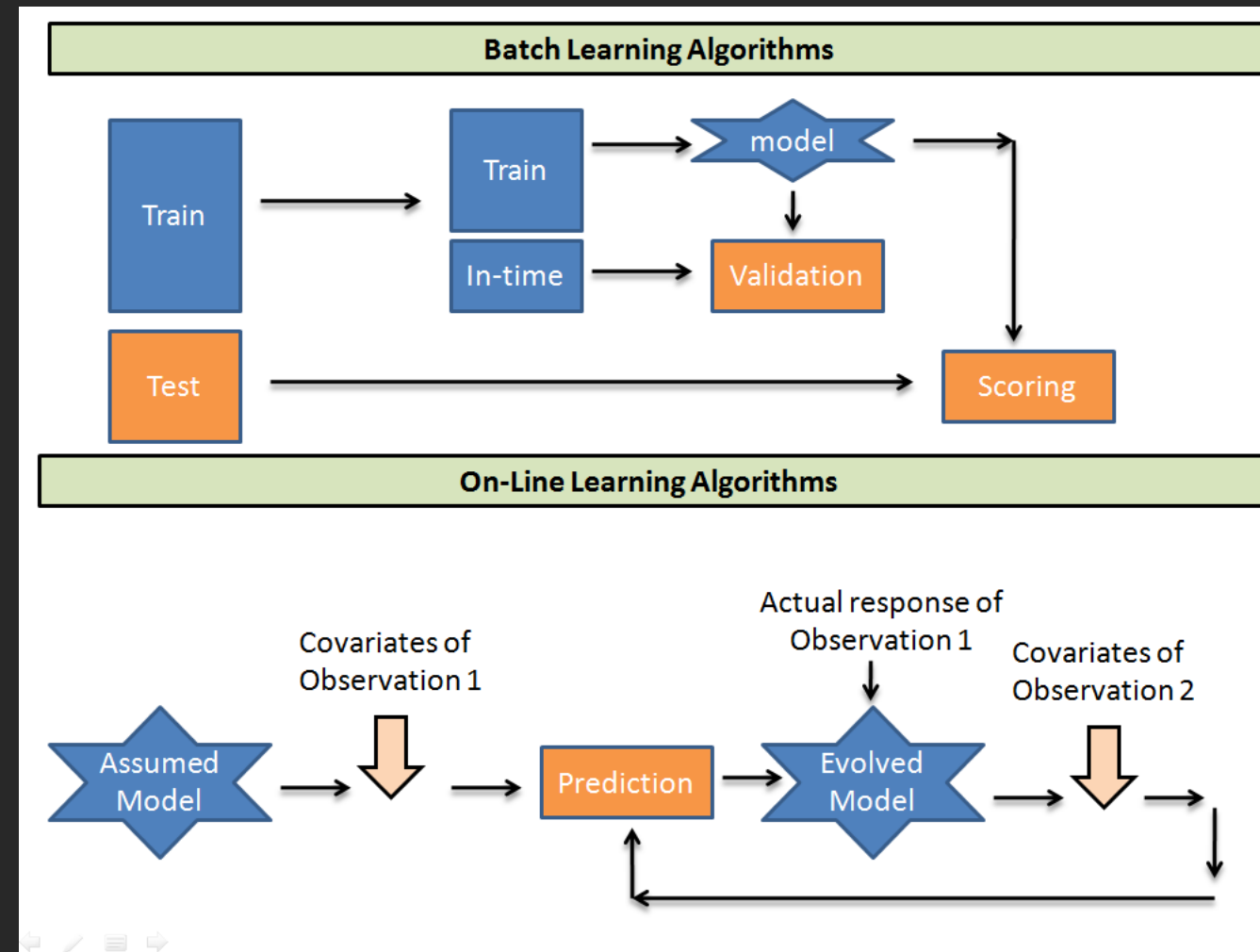
# Building a recommendation with contextual MAB



2 Recommendations (arm)

User/application
(environment)

Contextual MAB
model
(agent)

1 User features (context)

3 Implicit feedback, such as click
(reward)

# Bandits vs. A/B testing



*Image source (courtesy of Matt Gershoff)*

# Online learning



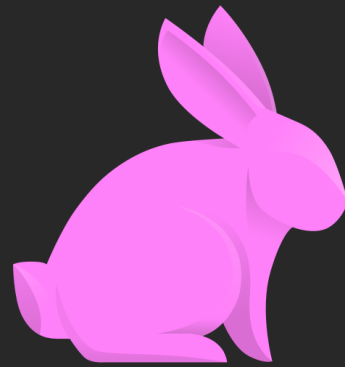Introduction to online machine learning simplified
(source https://analyticsvidhya.com)

# Adopting bandits into existing systems: Warm start

| FEATURES ▼ | CLASS ▼ |
|---|---|
| Observations set 1 | Action 3 |
| Observations set 2 | Action 2 |
| Observations set 3 | Action 1 |
| Observations set 4 | Action 2 |
| Observations set 5 | Action 2 |
| Observations set 6 | Action 1 |

| CONTEXT ▼ | Action 1/Arm 1 ▼ | Action 2/Arm 2 ▼ | Action 3/Arm 3 ▼ |
|---|---|---|---|
| Context 1 | | | Reward = 1 |
| Context 2 | | Reward = 1 | |
| Context 3 | Reward = 1 | | |
| Context 4 | | Reward = 1 | |
| Context 5 | | Reward = 1 | |
| Context 6 | Reward = 1 | | |

# Amazon SageMaker RL bandits container
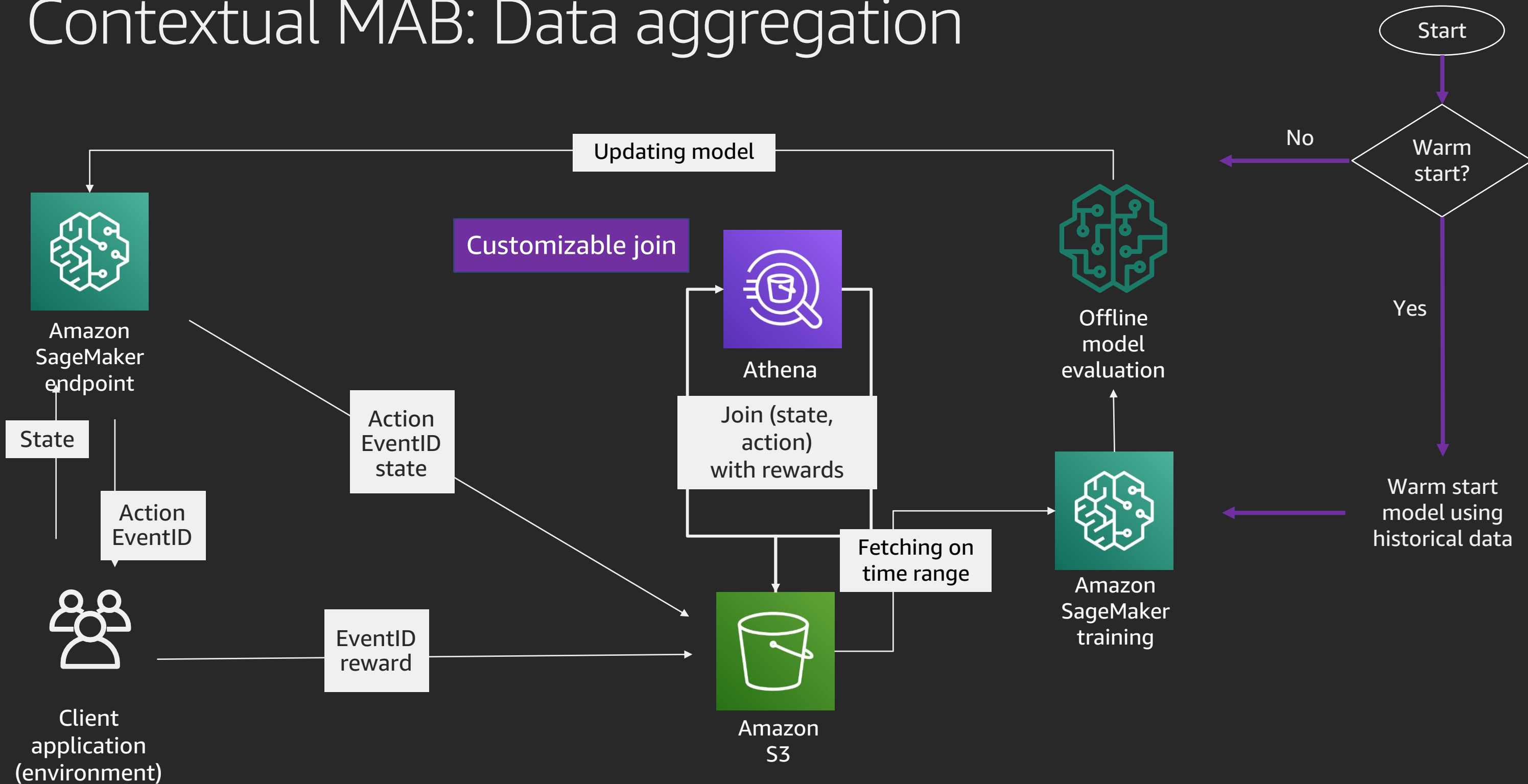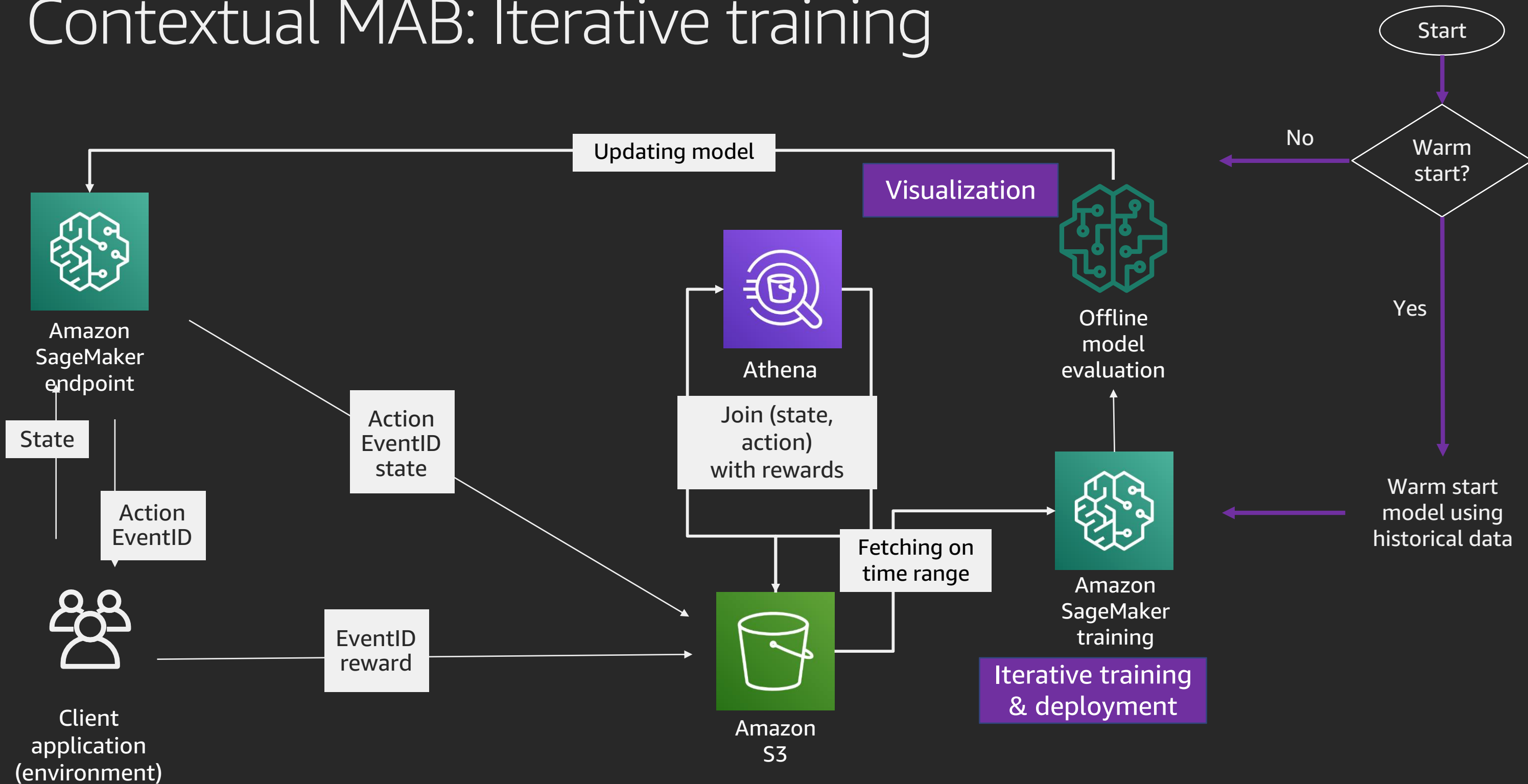


**VOWPAL WABBIT**

# Contextual MAB: Initialization



Start

Cold start

Warm start?

No

Yes

Updating model

Amazon
SageMaker
endpoint

Offline
model
evaluation

Data format:
- State
- Action
- Probability
- Reward

Fetching on
time range

Amazon
S3

Amazon
SageMaker
training

Warm start
model using
historical data

Warm start

# Contextual MAB: Data collection

Start

Warm start?

No

Updating model

Inference logging

Amazon SageMaker endpoint

State

Action EventID

Action EventID state

Yes

Offline model evaluation

Warm start model using historical data

Client application (environment)

EventID reward

Reward ingestion

Amazon S3

Fetching on time range

Amazon SageMaker training

# Contextual MAB: Data aggregation



Start

Warm start?

No

Yes

Updating model

Amazon SageMaker endpoint

State

Action EventID

Customizable join

Athena

Action EventID state

Join (state, action) with rewards

Offline model evaluation

Fetching on time range

Amazon SageMaker training

Warm start model using historical data

Client application (environment)

EventID reward

Amazon S3

# Contextual MAB: Iterative training

Start

Warm start?

No

Yes

Updating model

Visualization

Amazon SageMaker endpoint

Offline model evaluation

State

Action EventID

Action EventID state

Athena

Join (state, action) with rewards

Fetching on time range

Amazon SageMaker training

Warm start model using historical data

Client application (environment)

EventID reward

Amazon S3

Iterative training & deployment

# Contextual MAB: Evaluation

# Personalization with contextual bandits

# Configurations

```yaml
# Vowpal Wabbit container
image: "462105765813.dkr.ecr.{AWS_REGION}.amazonaws.com/sagemaker-rl-vw-container:vw-8.7.0-cpu"
# Vowpal Wabbit algorithm parameters
algor:
  algorithms_parameters:
    exploration_policy: "egreedy" # supports "egreedy", "bag", "cover"
    epsilon: 0.001 # used if egreedy is the exploration policy
    num_policies: 3 # used if bag or cover is the exploration policy
    num_arms: 7
    cfa_type: "dr" # supports "dr", "ips"
# use local mode?
local_mode: true
# if true, use the same endpoint with updated model
soft_deployment: true
```

# Reviewing the setup

## Amazon DynamoDB



## Amazon Kinesis Data Firehose

# Reviewing the setup, continued

## Amazon Athena

# Workshop #2: Training with a simulator

aws

# Training with HVAC simulator



LINK

# Conclusion

aws

Amazon SageMaker (working with other AWS services) makes it equally easy to train with and without simulation environments.

Amazon SageMaker provides containers with popular RL algorithms, and you can bring your own. This includes online learning algorithms.

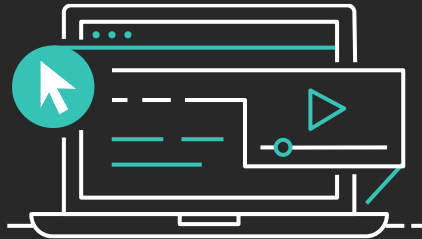Contextual bandits make experimentation very effective, and they learn rapidly.

# Conquer the newest frontier of ML: Reinforcement learning with Amazon SageMaker

# Learn ML with AWS Training and Certification

The same training that our own developers use, now available on demand

Role-based ML learning paths for developers, data scientists, data platform engineers, and business decision makers

70+ free digital ML courses from AWS experts let you learn from real-world challenges tackled at AWS

Validate expertise with the
**AWS Certified Machine Learning - Specialty** exam

Visit https://aws.training/machinelearning

aws training and certification

# Thank you!

**Girish Dilip Patil**
girpatil@amazon.com

**Marc Cabocel**
cabocel@amazon.fr

**Segolene Dessertine-panhard**

**Ann Luo**

# ! Please complete the session survey in the mobile app.