



# SageMaker Training



**Scott Pletcher**

INSTRUCTOR

**Machine Learning = Mathematics**

**Computers are good with math.**

**Graphical Processing Units are optimized for math.**

**GPUs are not the most efficient math engine.**



# Why GPUs?

Most Performant

Least Performant



# Why GPUs?

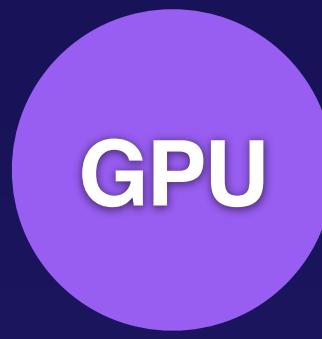
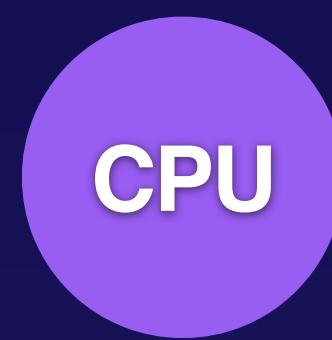
Least Flexible



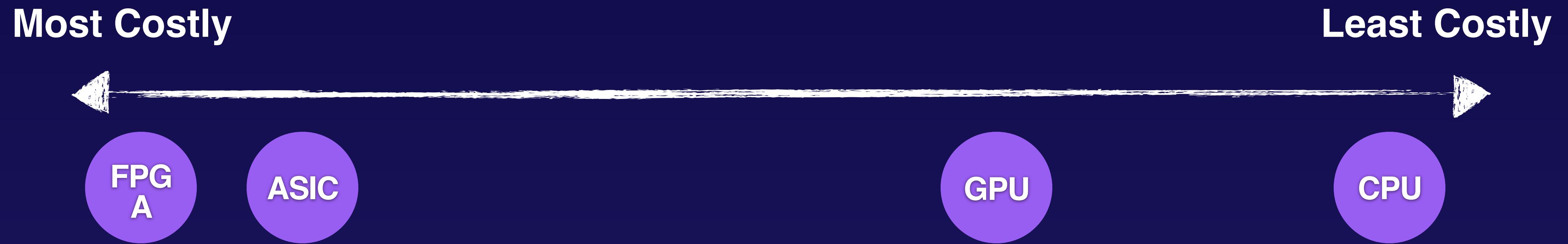
FPG  
A

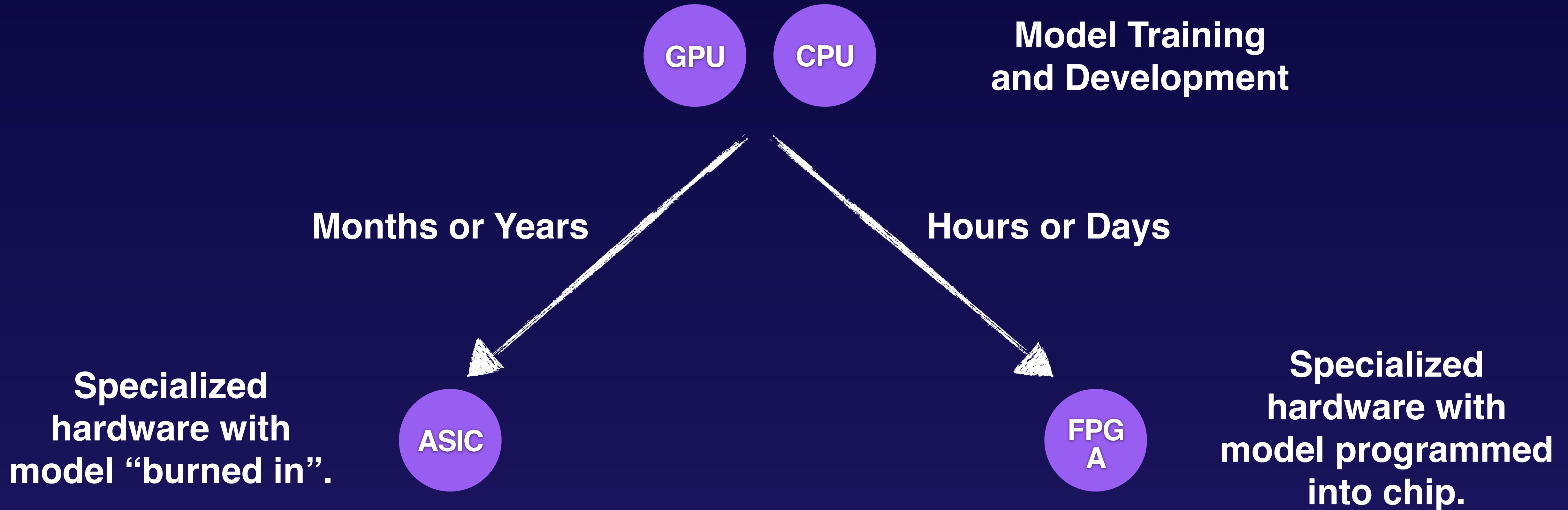


Most Flexible



# Why GPUs?





# 51% Attack

**Group of miners  
controlling more than  
51% of the hash rate.**

- Prevent new transactions
- Reverse completed transactions



# Outsmart the Machines





High-Level Python  
library provided by  
Amazon SageMaker



Use the SDK for Python

CreateTrainingJob API

## High-Level Python Library



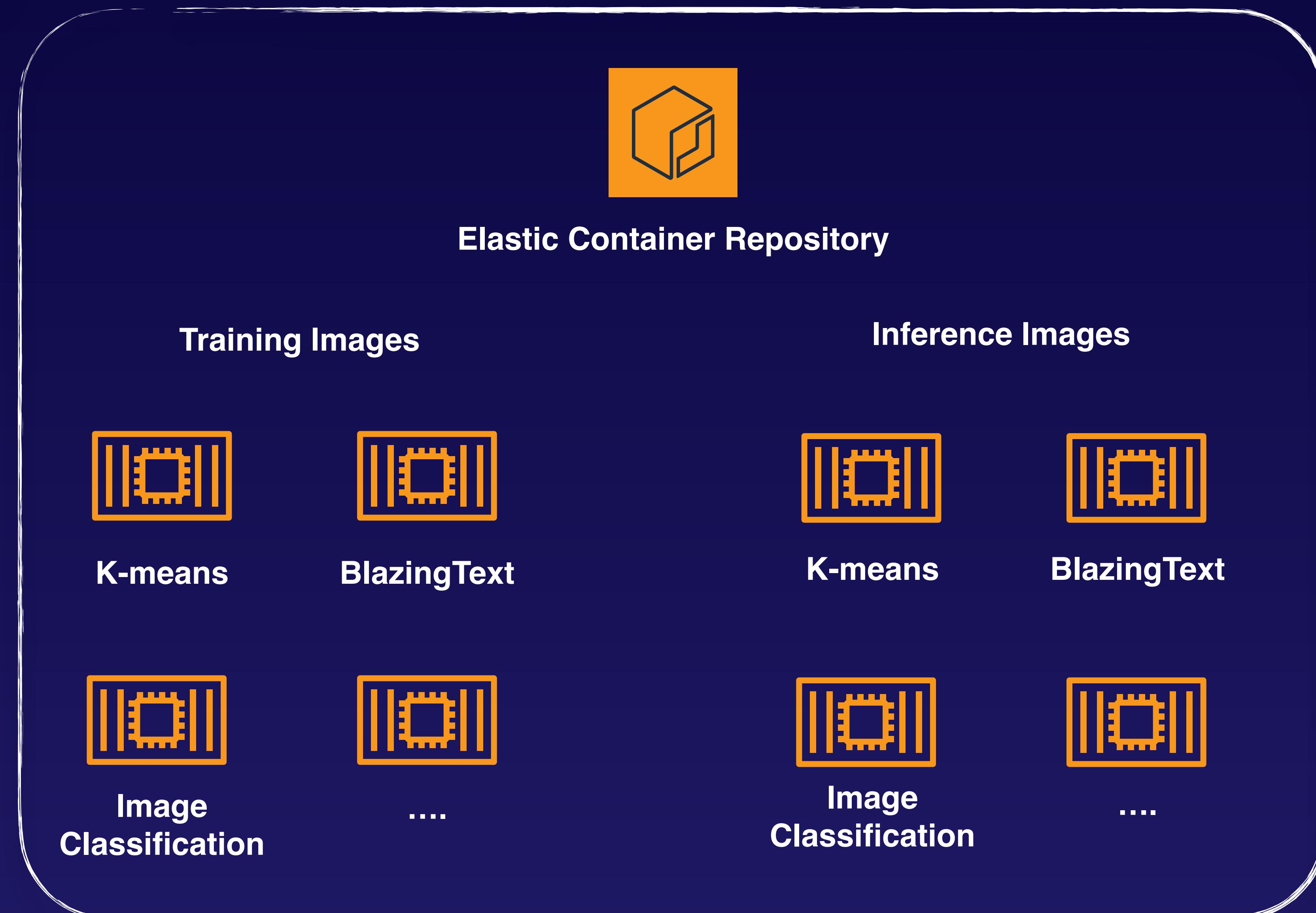
```
from sagemaker import KMeans

data_location = 's3://{} /kmeans_highlevel_example/data'.format(bucket)
output_location = 's3://{} /kmeans_highlevel_example/output'.format(bucket)

kmeans = KMeans(role = role,
                 train_instance_count = 2,
                 train_instance_type = 'ml.c4.8xlarge',
                 output_path = output_location,
                 k = 10,
                 data_location = data_location)

kmeans.fit(kmeans.record_set(train_set[0]))
```

From a Jupyter Notebook, you can issue commands to launch a training job.





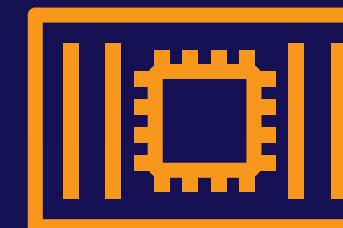
### Elastic Container Repository

**CreateTrainingJob API Call**

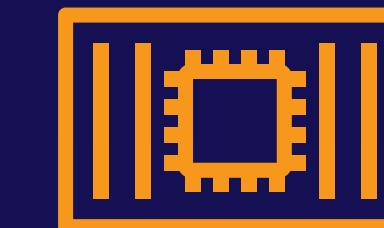
Training Images

**CreateModel API Call**

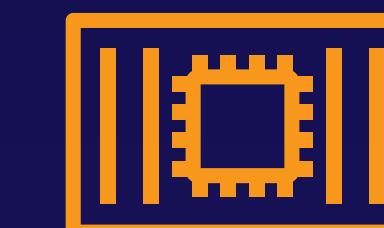
Inference Images



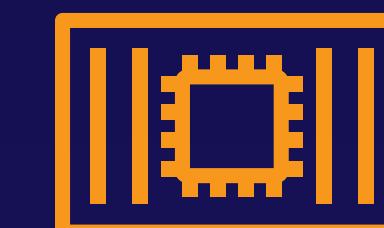
K-means



BlazingText



K-means



BlazingText

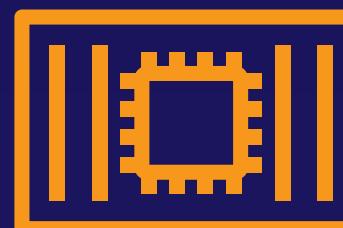
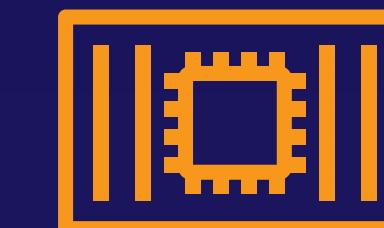


Image  
Classification



....

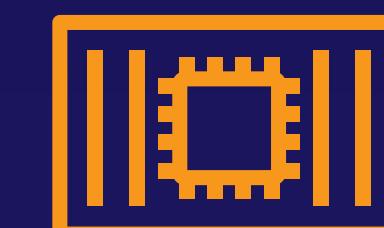
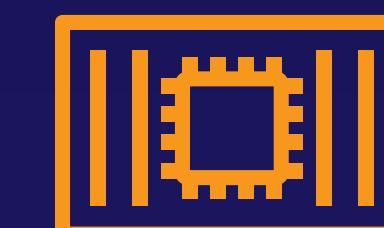


Image  
Classification



....

# Training



Region	ECR Path
us-west-1	632365934929.dkr.ecr.us-west-1.amazonaws.com
us-west-2	174872318107.dkr.ecr.us-west-2.amazonaws.com
ap-south-1	991648021394.dkr.ecr.ap-south-1.amazonaws.com
...	...

## Algorithms Provided by Amazon SageMaker: Common Parameters

The following table lists parameters for each of the algorithms provided by Amazon SageMaker.

Algorithm Name	Channel Name	Training Image and Inference Image Registry Path	Training Input Mode	File Type	Instance Class
BlazingText	train	<ecr_path>/blazingtext:<tag>	File or Pipe	Text file (one sentence per line with space-separated tokens)	GPU (single instance only) or CPU
DeepAR Forecasting	train and (optionally) test	<ecr_path>/forecasting-deepar:<tag>	File	JSON Lines or Parquet	GPU or CPU
Factorization Machines	train and (optionally) test	<ecr_path>/factorization-machines:<tag>	File or Pipe	recordIO-protobuf	CPU (GPU for dense data)
Image Classification	train and validation, (optionally) train_lst, validation_lst, and model	<ecr_path>/image-classification:<tag>	File or Pipe	recordIO or image files (.jpg or .png)	GPU
IP Insights	train and (optionally) validation	<ecr_path>/ipinsights:<tag>	File	CSV	GPU or GPU
k-means	train and (optionally) test	<ecr_path>/kmeans:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPUCommon (single GPU device on one or more instances)
k-nearest-neighbor (k-NN)	train and (optionally) test	<ecr_path>/knn:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPU (single GPU device on one or more instances)

**Named input source for algorithm consumption.**

# Training

[AWS Documentation](#) » [Amazon SageMaker](#) » [Developer Guide](#) » [Using Built-in Algorithms with Amazon SageMaker](#) » [Algorithms Provided by Amazon SageMaker: Common Information](#)  
 » [Algorithms Provided by Amazon SageMaker: Common Parameters](#)

## Algorithms Provided by Amazon SageMaker: Common Parameters

The following table lists parameters for each of the algorithms provided by Amazon SageMaker.

Algorithm Name	Channel Name	Training Image and Inference Image Registry Path	Training Input Mode	File Type	Instance Class
BlazingText	train	<ecr_path>/blazingtext:<tag>	File or Pipe	Text file (one sentence per line with space-separated tokens)	GPU (single instance only) or CPU
DeepAR Forecasting	train and (optionally) test	<ecr_path>/forecasting-deepar:<tag>	File	JSON Lines or Parquet	GPU or CPU
Factorization Machines	train and (optionally) test	<ecr_path>/factorization-machines:<tag>	File or Pipe	recordIO-protobuf	CPU (GPU for dense data)
Image Classification	train and validation, (optionally) train_lst, validation_lst, and model	<ecr_path>/image-classification:<tag>	File or Pipe	recordIO or image files (.jpg or .png)	GPU
IP Insights	train and (optionally) validation	<ecr_path>/ipinsights:<tag>	File	CSV	GPU or GPU
k-means	train and (optionally) test	<ecr_path>/kmeans:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPUCommon (single GPU device on one or more instances)
k-nearest-neighbor (k-NN)	train and (optionally) test	<ecr_path>/knn:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPU (single GPU device on one or more instances)

**Path to the ECR image for a given algorithm.**

**If using the SageMaker Python Library, it will automatically know most paths.**

# Training

[AWS Documentation](#) » [Amazon SageMaker](#) » [Developer Guide](#) » [Using Built-in Algorithms with Amazon SageMaker](#) » [Algorithms Provided by Amazon SageMaker: Common Information](#)  
 » [Algorithms Provided by Amazon SageMaker: Common Parameters](#)

## Algorithms Provided by Amazon SageMaker: Common Parameters

The following table lists parameters for each of the algorithms provided by Amazon SageMaker.

Algorithm Name	Channel Name	Training Image and Inference Image Registry Path	Training Input Mode	File Type	Instance Class
BlazingText	train	<ecr_path>/blazingtext:<tag>	File or Pipe	Text file (one sentence per line with space-separated tokens)	GPU (single instance only) or CPU
DeepAR Forecasting	train and (optionally) test	<ecr_path>/forecasting-deepar:<tag>	File	JSON Lines or Parquet	GPU or CPU
Factorization Machines	train and (optionally) test	<ecr_path>/factorization-machines:<tag>	File or Pipe	recordIO-protobuf	CPU (GPU for dense data)
Image Classification	train and validation, (optionally) train_lst, validation_lst, and model	<ecr_path>/image-classification:<tag>	File or Pipe	recordIO or image files (.jpg or .png)	GPU
IP Insights	train and (optionally) validation	<ecr_path>/ipinsights:<tag>	File	CSV	GPU or GPU
k-means	train and (optionally) test	<ecr_path>/kmeans:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPUCommon (single GPU device on one or more instances)
k-nearest-neighbor (k-NN)	train and (optionally) test	<ecr_path>/knn:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPU (single GPU device on one or more instances)

**Note the Tag attribute. It is a form of versioning for the repository.**

**Use :1 version tag to ensure the stable version.**

**Use :latest version tag for up-to-date but potentially not backward compatible.**

**Use :1 for production purposes.**

# Training

[AWS Documentation](#) » [Amazon SageMaker](#) » [Developer Guide](#) » [Using Built-in Algorithms with Amazon SageMaker](#) » [Algorithms Provided by Amazon SageMaker: Common Information](#)  
 » [Algorithms Provided by Amazon SageMaker: Common Parameters](#)

## Algorithms Provided by Amazon SageMaker: Common Parameters

The following table lists parameters for each of the algorithms provided by Amazon SageMaker.

Algorithm Name	Channel Name	Training Image and Inference Image Registry Path	Training Input Mode	File Type	Instance Class
BlazingText	train	<ecr_path>/blazingtext:<tag>	File or Pipe	Text file (one sentence per line with space-separated tokens)	GPU (single instance only) or CPU
DeepAR Forecasting	train and (optionally) test	<ecr_path>/forecasting-deepar:<tag>	File	JSON Lines or Parquet	GPU or CPU
Factorization Machines	train and (optionally) test	<ecr_path>/factorization-machines:<tag>	File or Pipe	recordIO-protobuf	CPU (GPU for dense data)
Image Classification	train and validation, (optionally) train_lst, validation_lst, and model	<ecr_path>/image-classification:<tag>	File or Pipe	recordIO or image files (.jpg or .png)	GPU
IP Insights	train and (optionally) validation	<ecr_path>/ipinsights:<tag>	File	CSV	GPU or GPU
k-means	train and (optionally) test	<ecr_path>/kmeans:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPUCommon (single GPU device on one or more instances)
k-nearest-neighbor (k-NN)	train and (optionally) test	<ecr_path>/knn:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPU (single GPU device on one or more instances)

Type of input data accepted.

Recall that Pipe has better performance.

# Training

[AWS Documentation](#) » [Amazon SageMaker](#) » [Developer Guide](#) » [Using Built-in Algorithms with Amazon SageMaker](#) » [Algorithms Provided by Amazon SageMaker: Common Information](#)  
 » [Algorithms Provided by Amazon SageMaker: Common Parameters](#)

## Algorithms Provided by Amazon SageMaker: Common Parameters

The following table lists parameters for each of the algorithms provided by Amazon SageMaker.

Algorithm Name	Channel Name	Training Image and Inference Image Registry Path	Training Input Mode	File Type	Instance Class
BlazingText	train	<ecr_path>/blazingtext:<tag>	File or Pipe	Text file (one sentence per line with space-separated tokens)	GPU (single instance only) or CPU
DeepAR Forecasting	train and (optionally) test	<ecr_path>/forecasting-deepar:<tag>	File	JSON Lines or Parquet	GPU or CPU
Factorization Machines	train and (optionally) test	<ecr_path>/factorization-machines:<tag>	File or Pipe	recordIO-protobuf	CPU (GPU for dense data)
Image Classification	train and validation, (optionally) train_lst, validation_lst, and model	<ecr_path>/image-classification:<tag>	File or Pipe	recordIO or image files (.jpg or .png)	GPU
IP Insights	train and (optionally) validation	<ecr_path>/ipinsights:<tag>	File	CSV	GPU or GPU
k-means	train and (optionally) test	<ecr_path>/kmeans:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPUCommon (single GPU device on one or more instances)
k-nearest-neighbor (k-NN)	train and (optionally) test	<ecr_path>/knn:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPU (single GPU device on one or more instances)

**Format of data accepted by algorithm.**

**Recall that recordIO protobuf has best performance.**

# Training

[AWS Documentation](#) » [Amazon SageMaker](#) » [Developer Guide](#) » [Using Built-in Algorithms with Amazon SageMaker](#) » [Algorithms Provided by Amazon SageMaker: Common Information](#)  
 » [Algorithms Provided by Amazon SageMaker: Common Parameters](#)

## Algorithms Provided by Amazon SageMaker: Common Parameters

The following table lists parameters for each of the algorithms provided by Amazon SageMaker.

**Type of instance required for algorithm.**

**Note that some are GPU only or CPU only.**

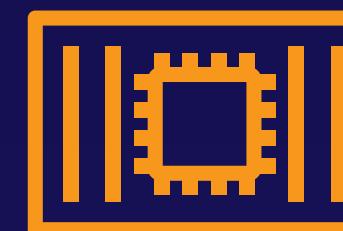
**Example:**  
**XGBoost implements an open source algorithm optimized for CPUs**

Algorithm Name	Channel Name	Training Image and Inference Image Registry Path	Training Input Mode	File Type	Instance Class
BlazingText	train	<ecr_path>/blazingtext:<tag>	File or Pipe	Text file (one sentence per line with space-separated tokens)	GPU (single instance only) or CPU
DeepAR Forecasting	train and (optionally) test	<ecr_path>/forecasting-deepar:<tag>	File	JSON Lines or Parquet	GPU or CPU
Factorization Machines	train and (optionally) test	<ecr_path>/factorization-machines:<tag>	File or Pipe	recordIO-protobuf	CPU (GPU for dense data)
Image Classification	train and validation, (optionally) train_lst, validation_lst, and model	<ecr_path>/image-classification:<tag>	File or Pipe	recordIO or image files (.jpg or .png)	GPU
IP Insights	train and (optionally) validation	<ecr_path>/ipinsights:<tag>	File	CSV	GPU or GPU
k-means	train and (optionally) test	<ecr_path>/kmeans:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPUCommon (single GPU device on one or more instances)
k-nearest-neighbor (k-NN)	train and (optionally) test	<ecr_path>/knn:<tag>	File or Pipe	recordIO-protobuf or CSV	CPU or GPU (single GPU device on one or more instances)



## Elastic Container Repository

### Training Images



K-means

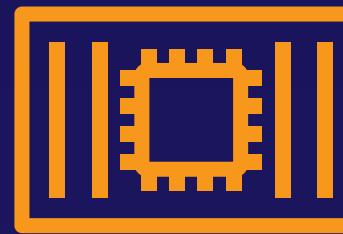
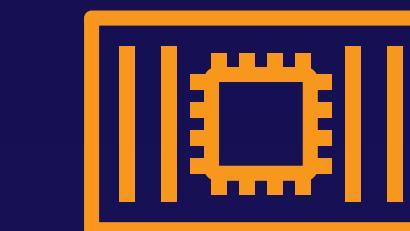
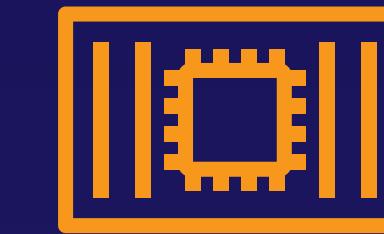


Image  
Classification

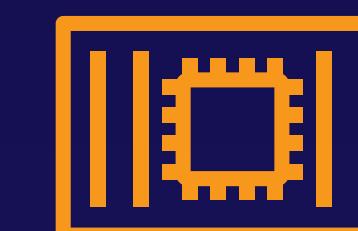


BlazingText

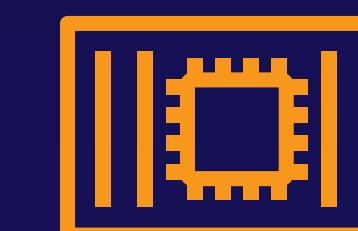


....

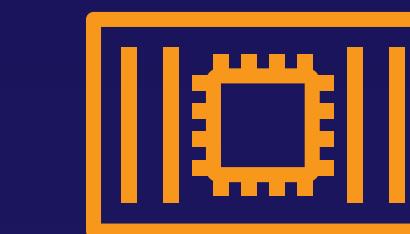
### Inference Images



K-means

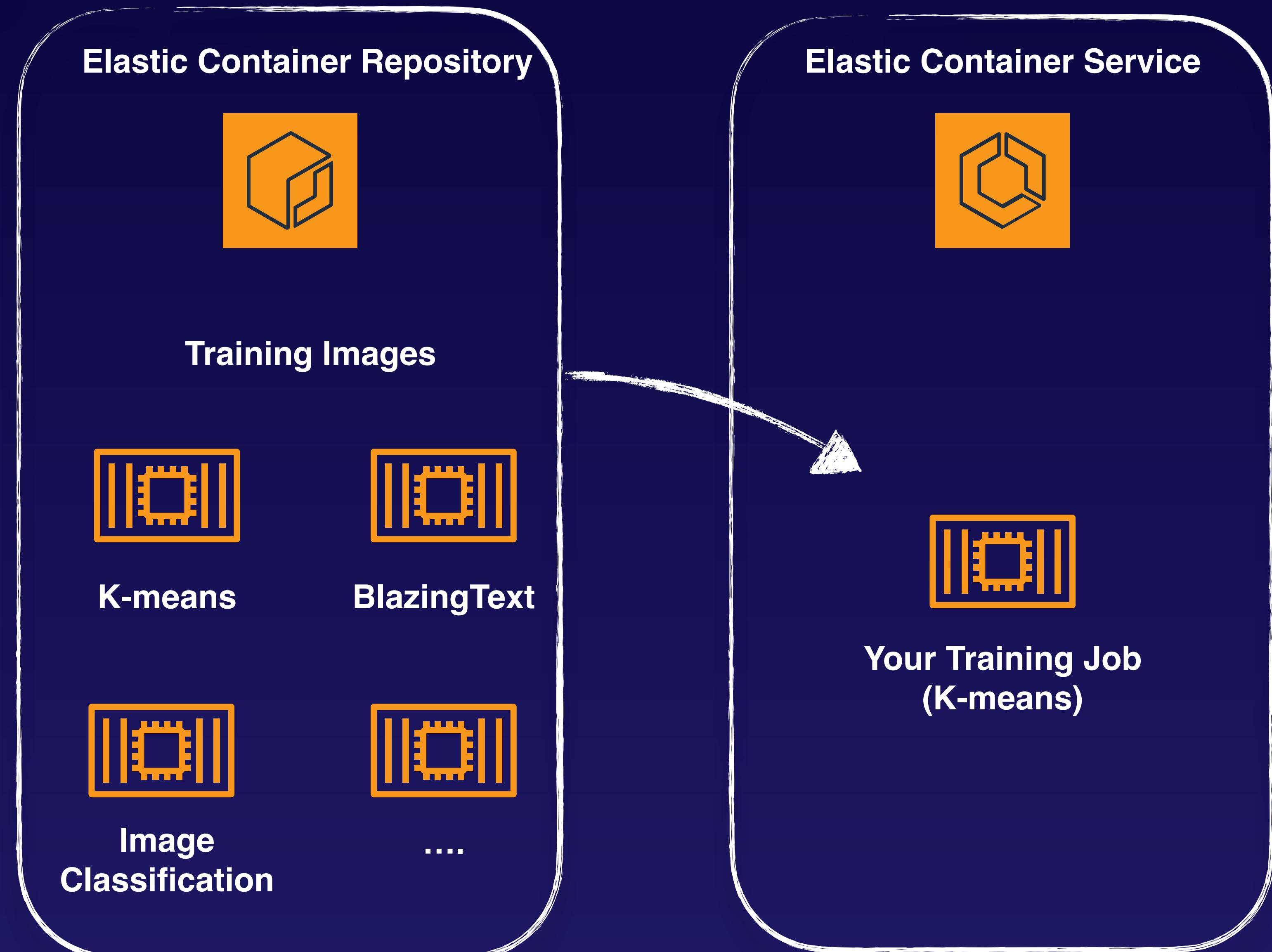


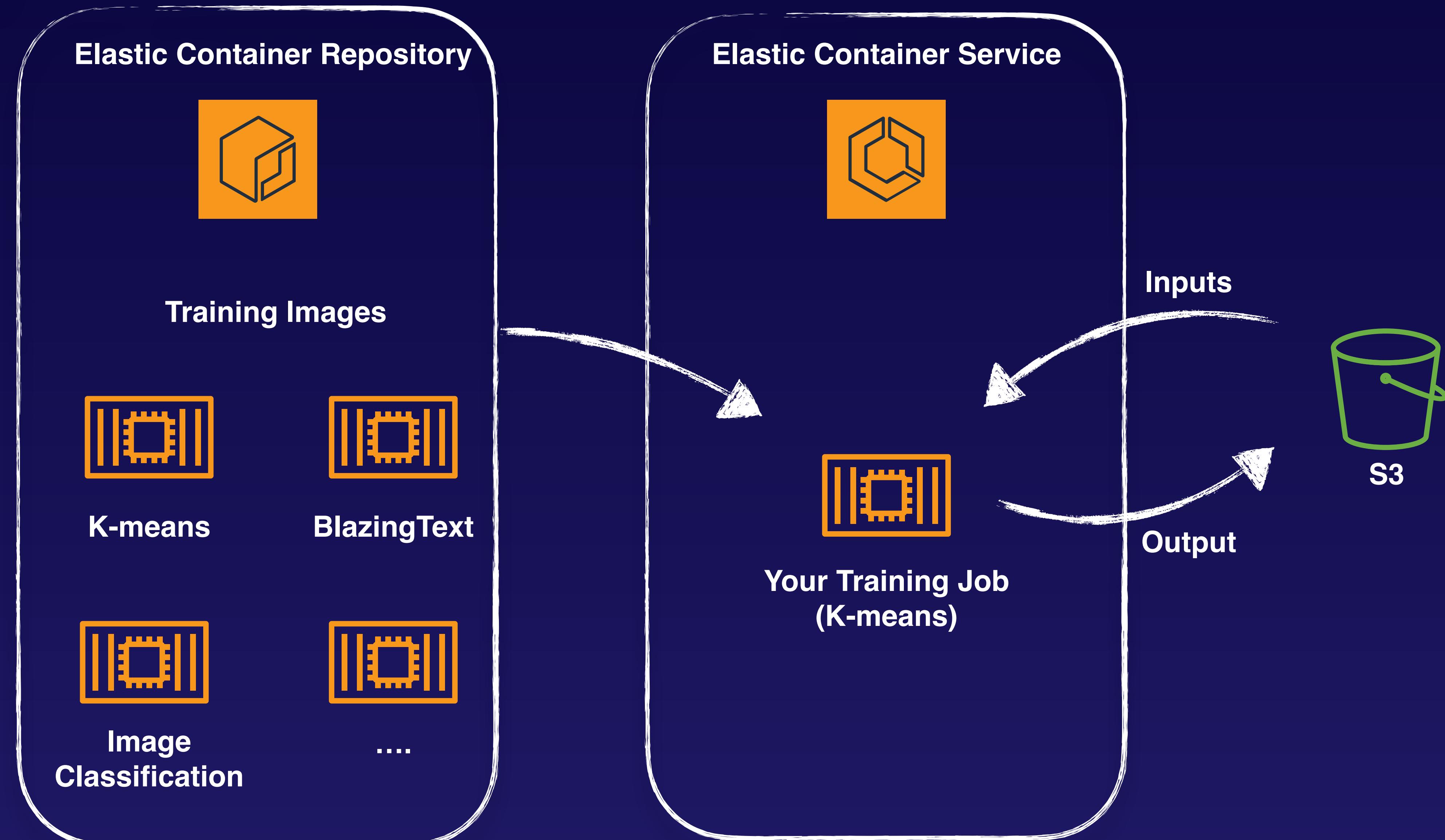
BlazingText



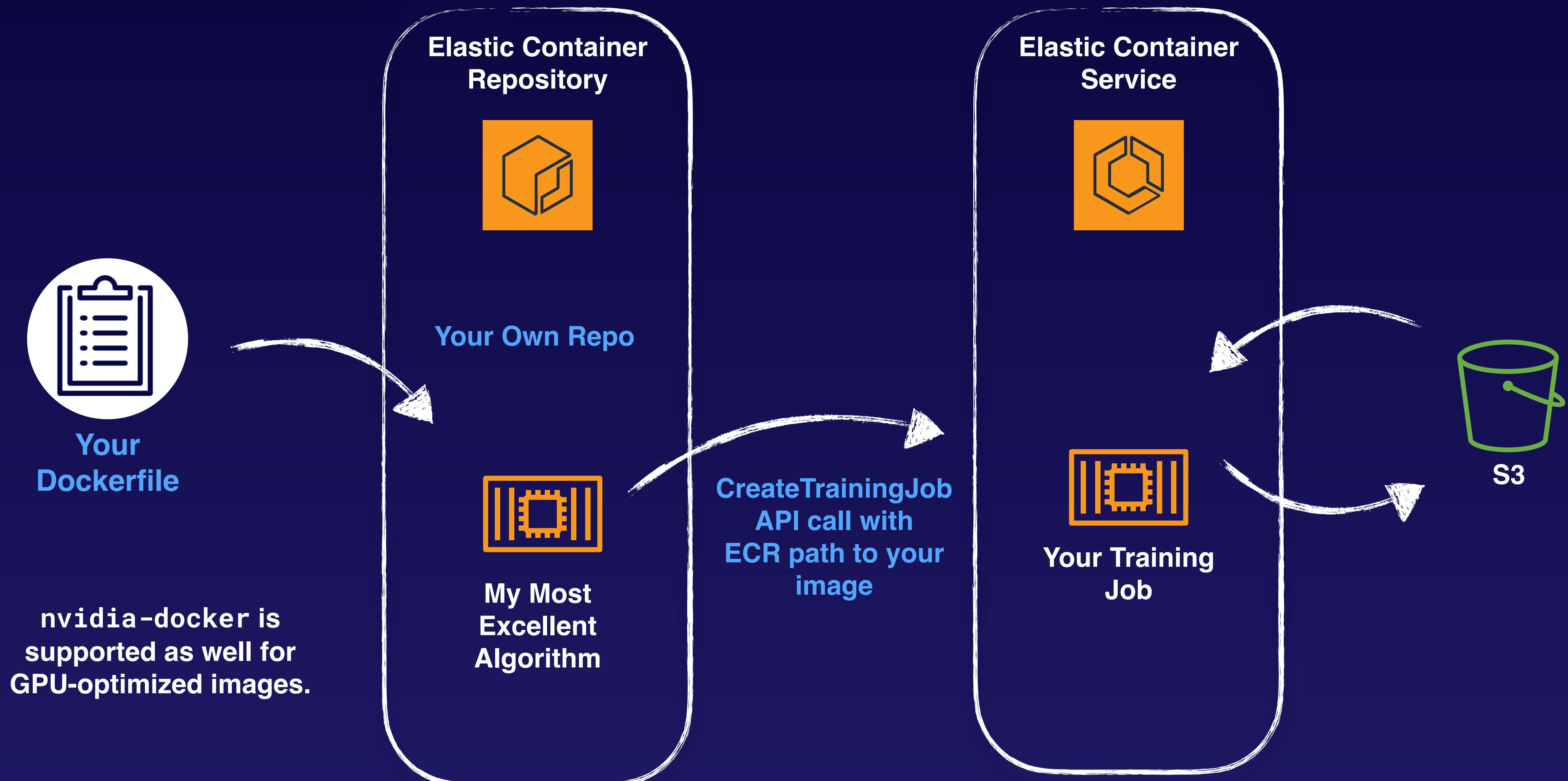
....

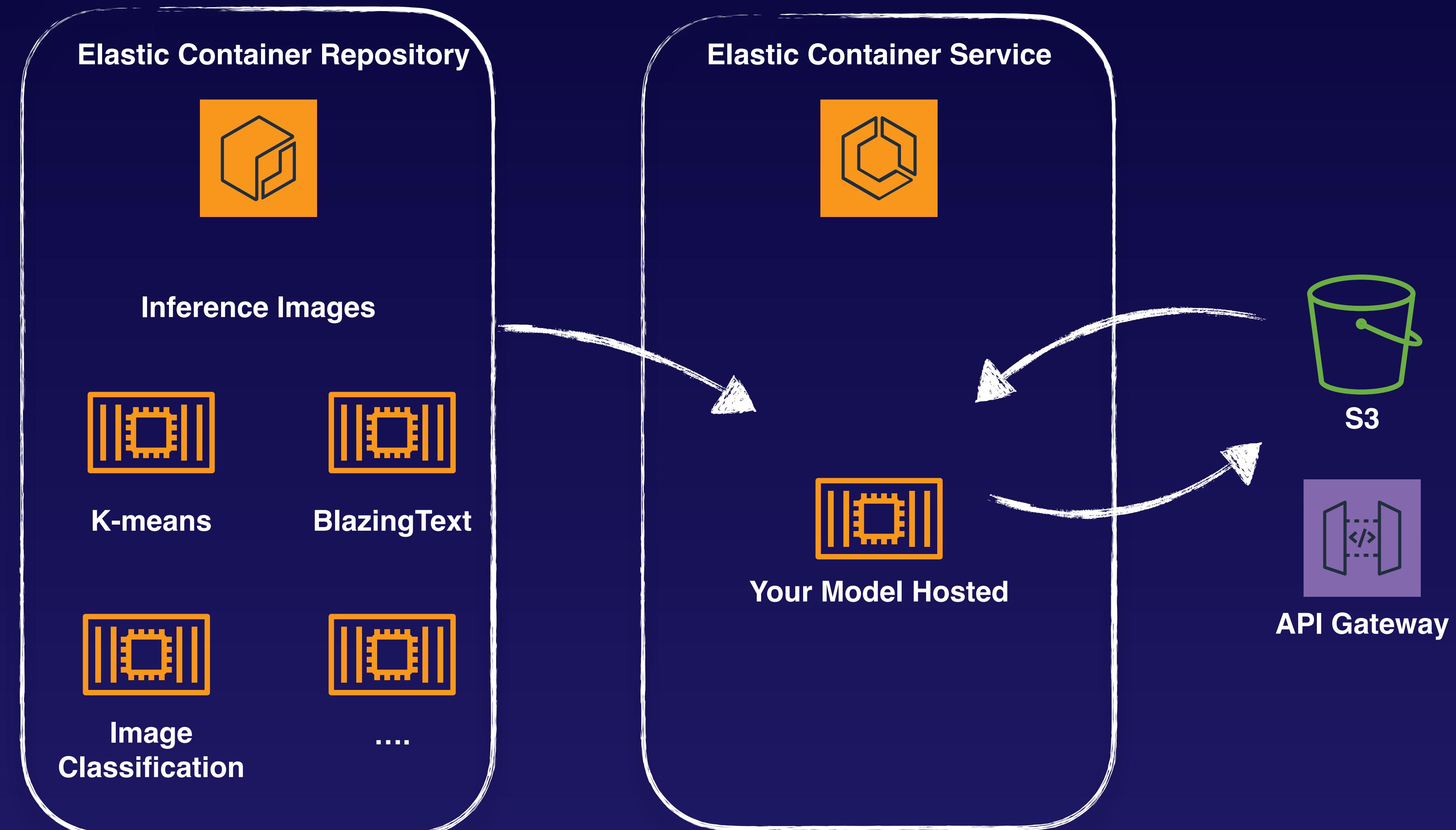
Image  
Classification





# Custom Algorithms





## Example of a Binary Classification Process

0	Luke	Rebel
1	JarJar	Empire
0	Han	Rebel
1	Vadar	Empire
0	Yoda	Rebel

## Example of a Binary Classification Process

0	Luke	Rebel	Does “Luke” = Evil?	N
1	JarJar	Empire	Does “JarJar” = Evil?	Y
0	Han	Rebel	Does “Han” = Evil?	N
1	Vadar	Empire	Does “Vadar” = Evil?	Y
0	Yoda	Rebel	Does “Yoda” = Evil?	N

## Example of a Binary Classification Process

0	Luke	Rebel	Does “Luke” = Evil?	N
1	JarJar	Empire	Does “JarJar” = Evil?	Y
0	Han	Rebel	Does “Han” = Evil?	N
1	Vadar	Empire	Does “Vadar” = Evil?	Y
0	Yoda	Rebel	Does “Yoda” = Evil?	N

Weak Correlation

## Example of a Binary Classification Process

0	Luke	Rebel	Does “Rebel” = Evil?	N
1	JarJar	Empire	Does “Empire” = Evil?	Y
0	Han	Rebel	Does “Rebel” = Evil?	N
1	Vadar	Empire	Does “Empire” = Evil?	Y
0	Yoda	Rebel	Does “Rebel” = Evil?	N

## Example of a Binary Classification Process

0	Luke	Rebel	Does “Rebel” = Evil?	N
1	JarJar	Empire	Does “Empire” = Evil?	Y
0	Han	Rebel	Does “Rebel” = Evil?	N
1	Vadar	Empire	Does “Empire” = Evil?	Y
0	Yoda	Rebel	Does “Rebel” = Evil?	N

Very Strong Correlation

## Example of a Binary Classification Process

0	Luke	Rebel	Does “Luke” & “Rebel” = Evil?	N
1	JarJar	Empire	Does “JarJar” & “Empire” = Evil?	Y
0	Han	Rebel	Does “Han” & “Rebel” = Evil?	N
1	Vadar	Empire	Does “Vadar” & “Empire” = Evil?	Y
0	Yoda	Rebel	Does “Yoda” & “Rebel” = Evil?	N

Just as good as Rebel/Empire...

If parameter = “Empire”, we predict “Evil (1)”.

### Testing Set

1	Jabba	Empire
1	Maul	Empire
0	Leia	Rebel

Given “Empire”, I predict “Evil (1)”

Given “Empire”, I predict “Evil (1)”

Given “Rebel”, I predict “Not Evil (0)”



100% Precision! Ship it!



## CloudWatch

### Information Logged

- Arguments provided
- Errors during training
- Algorithm accuracy statistics
- Timing of the algorithm

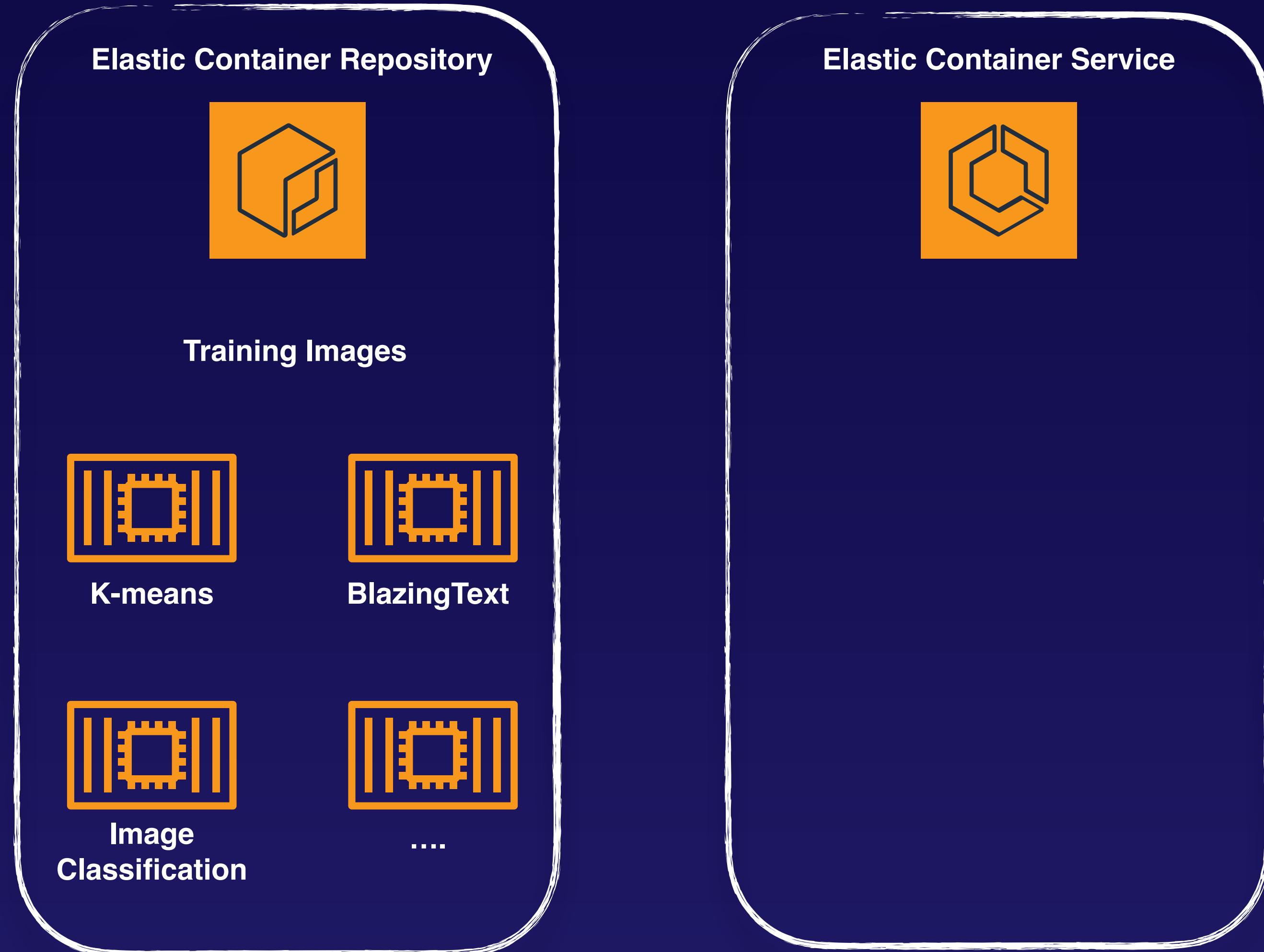


## CloudWatch

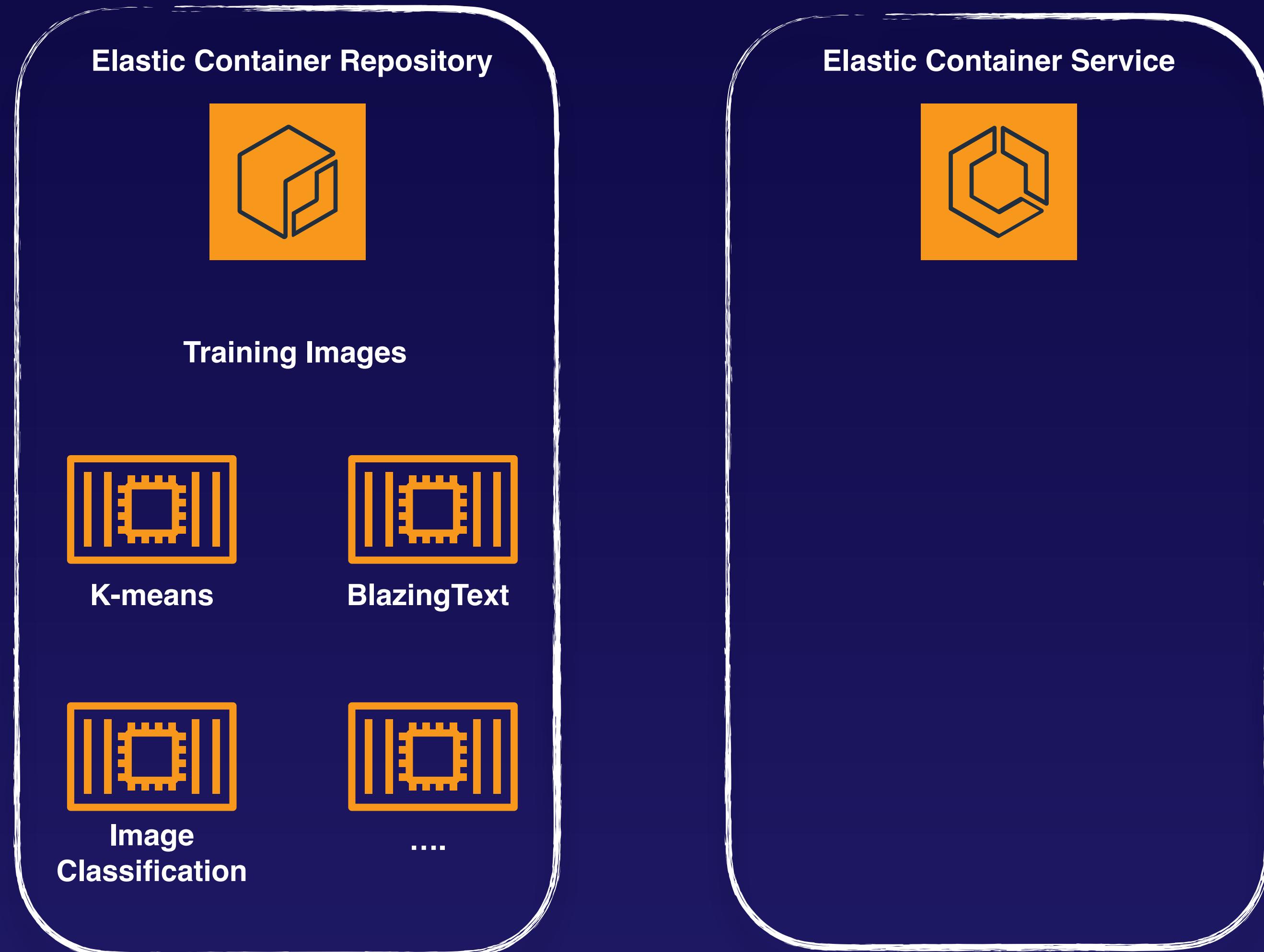
### Common Errors

- **Error in specifying a hyperparameter such as an extra one.**
- **Invalid value for a hyperparameter**
- **Incorrect protobuf file format**

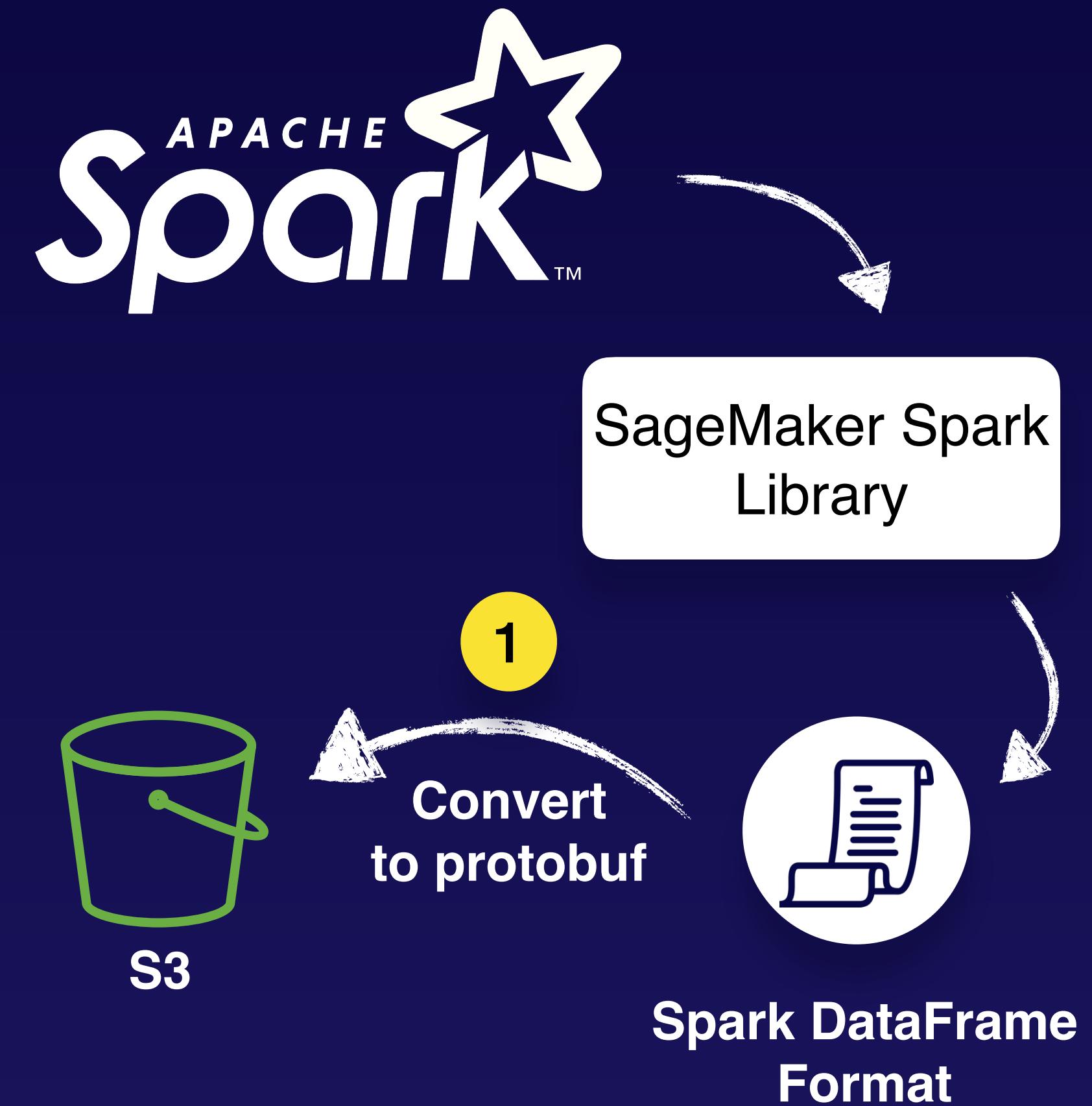
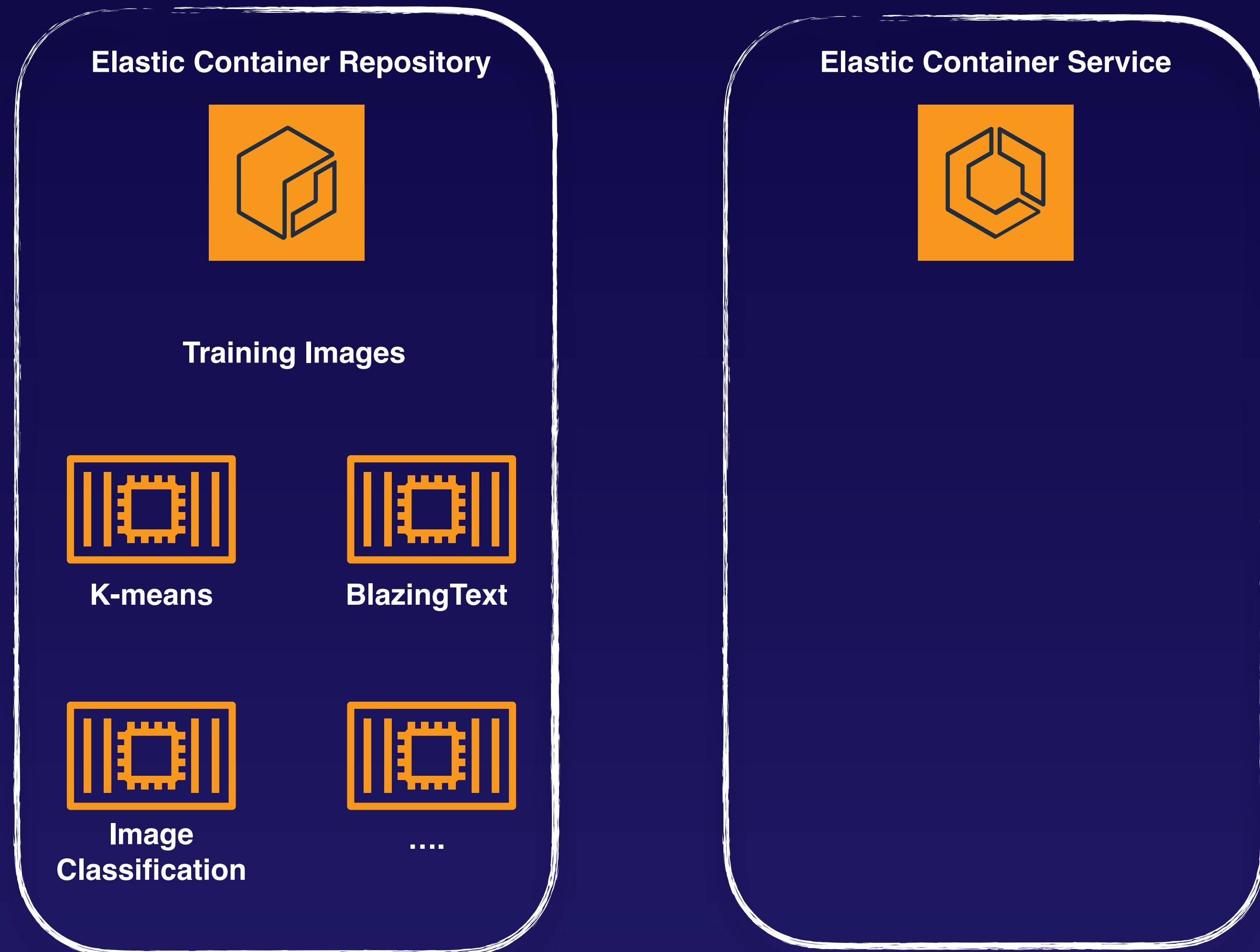
# Training with Apache Spark



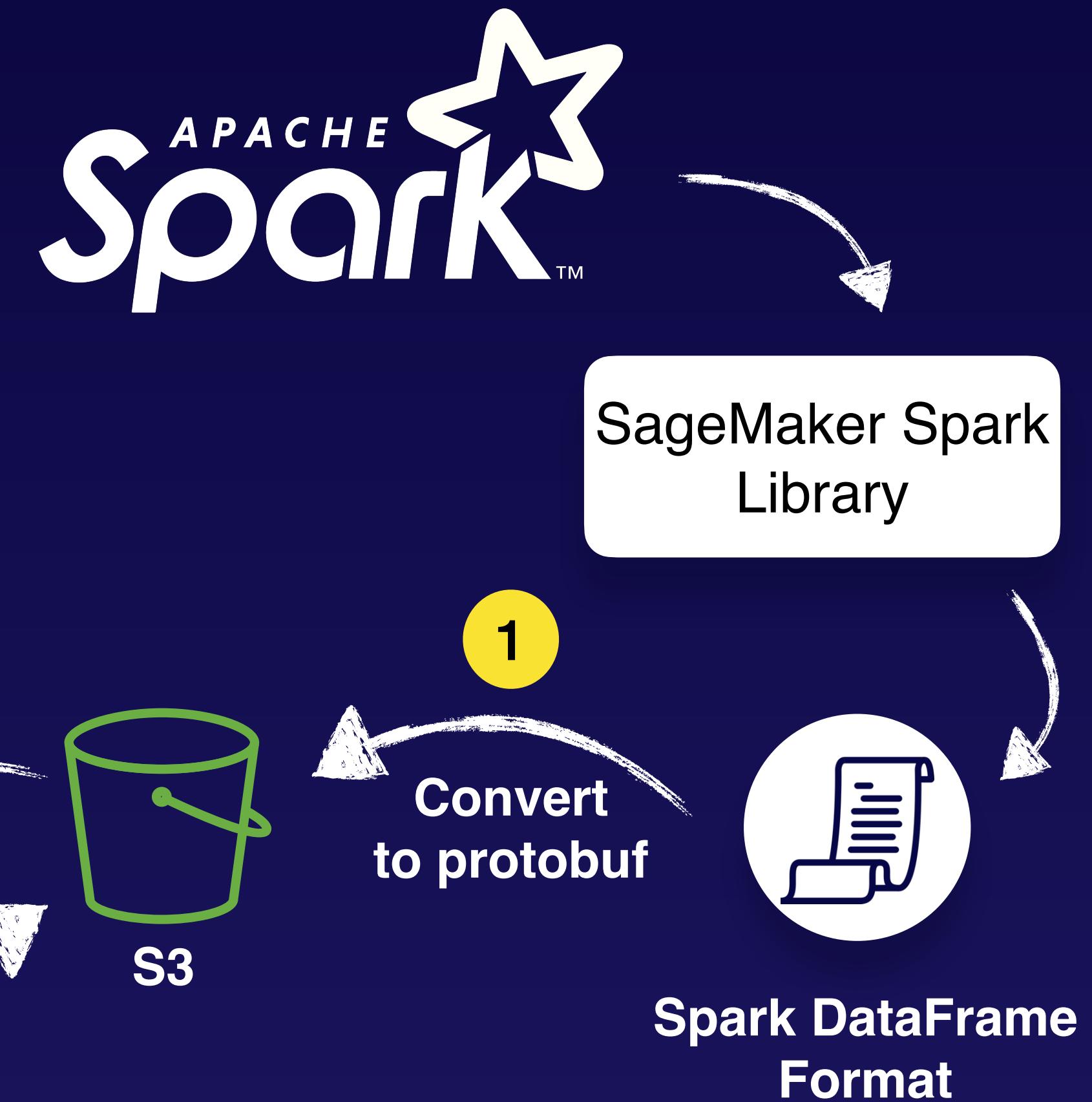
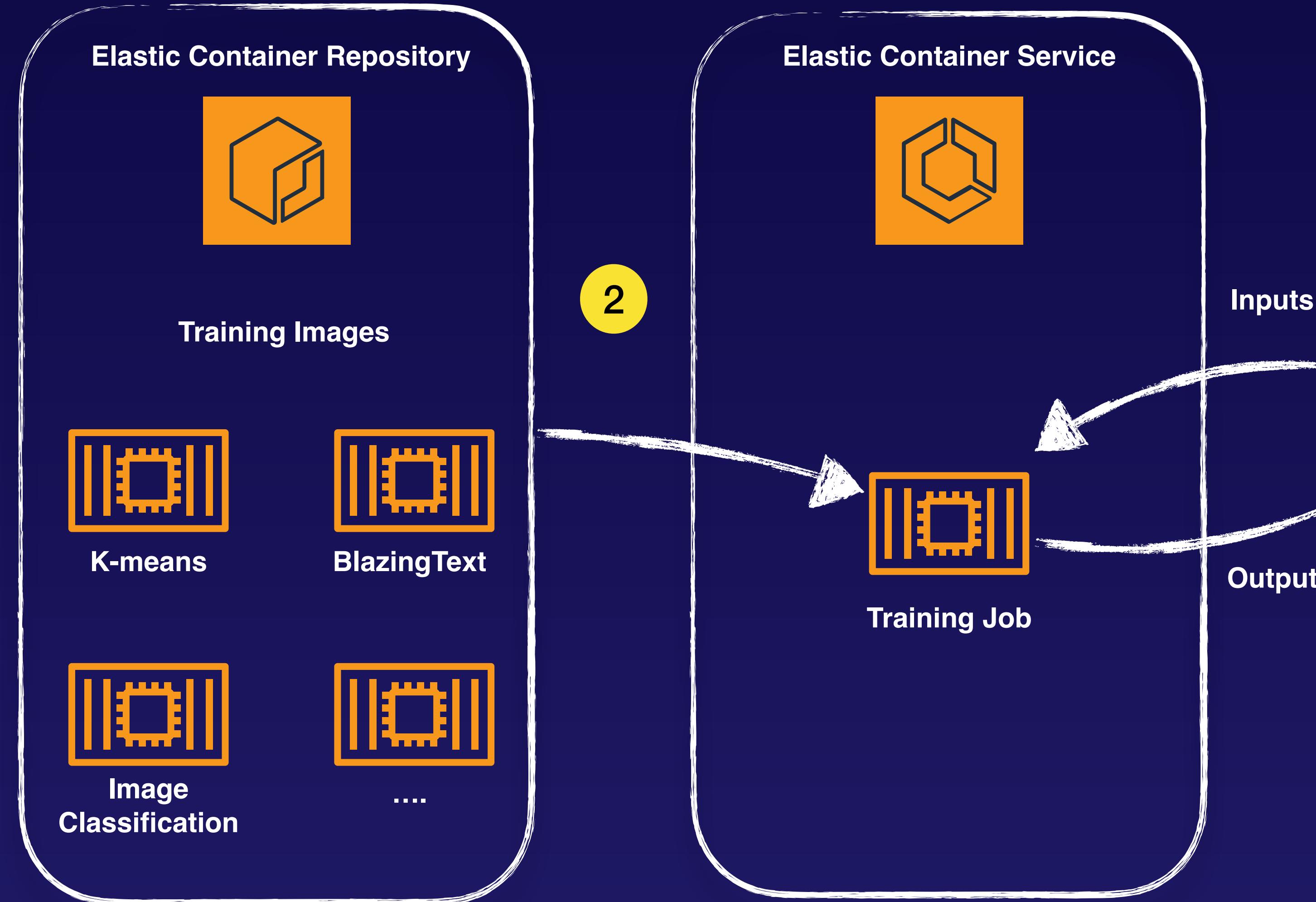
# Training with Apache Spark



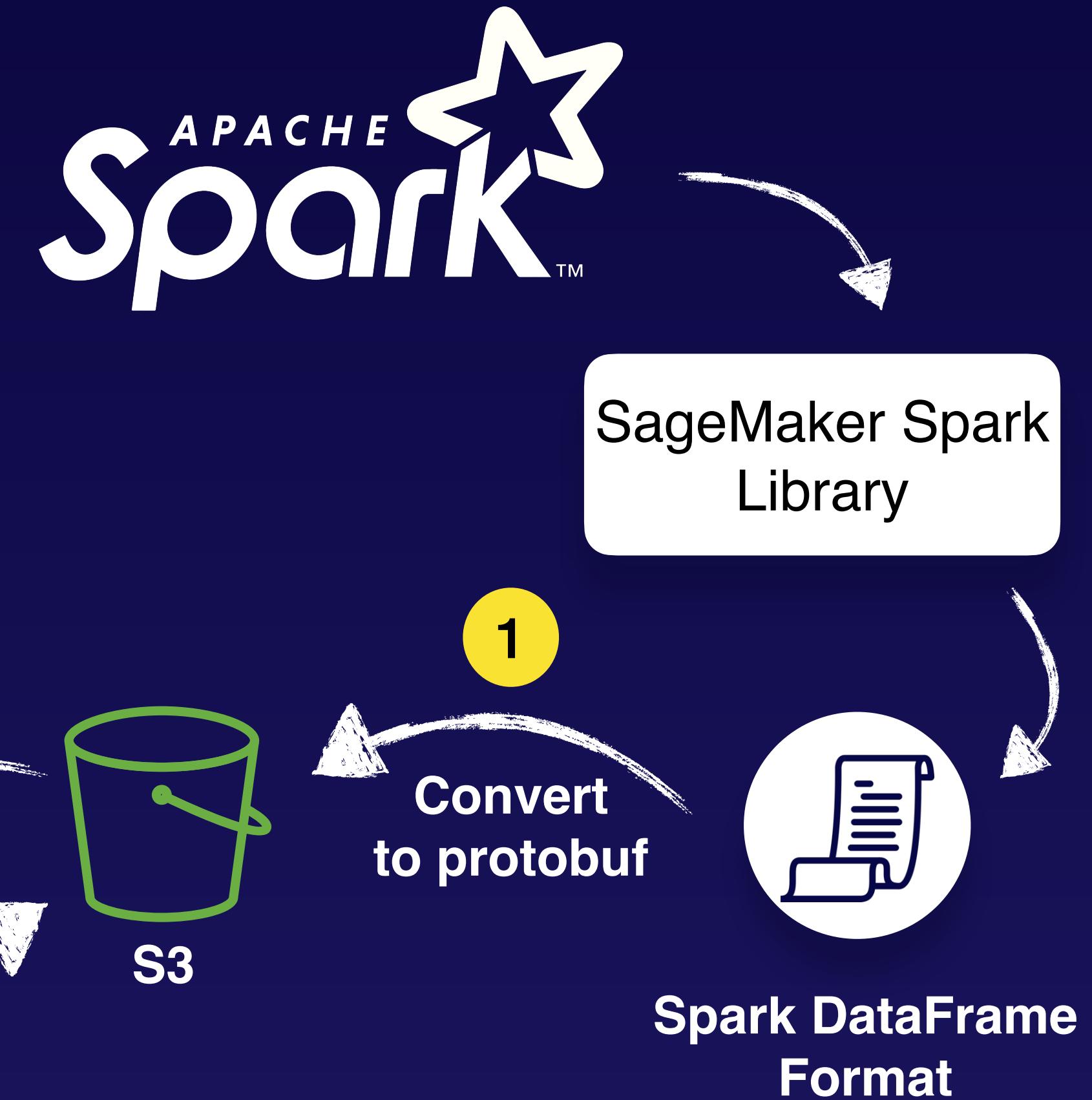
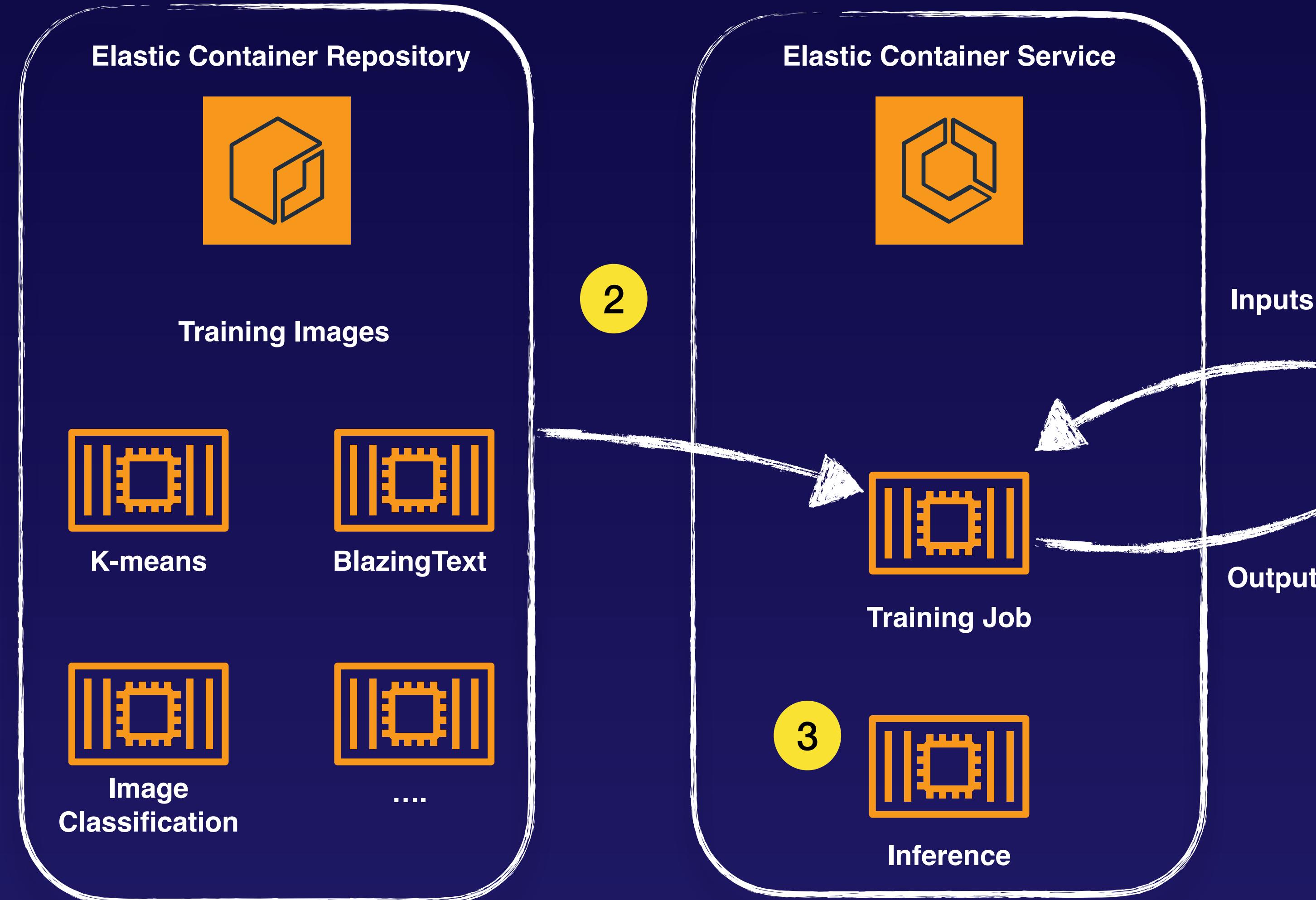
# Training with Apache Spark



# Training with Apache Spark



# Training with Apache Spark



# Training with Apache Spark

