



# Algorithms in Amazon SageMaker

**Zohar Karnin, Amazon AI Labs**

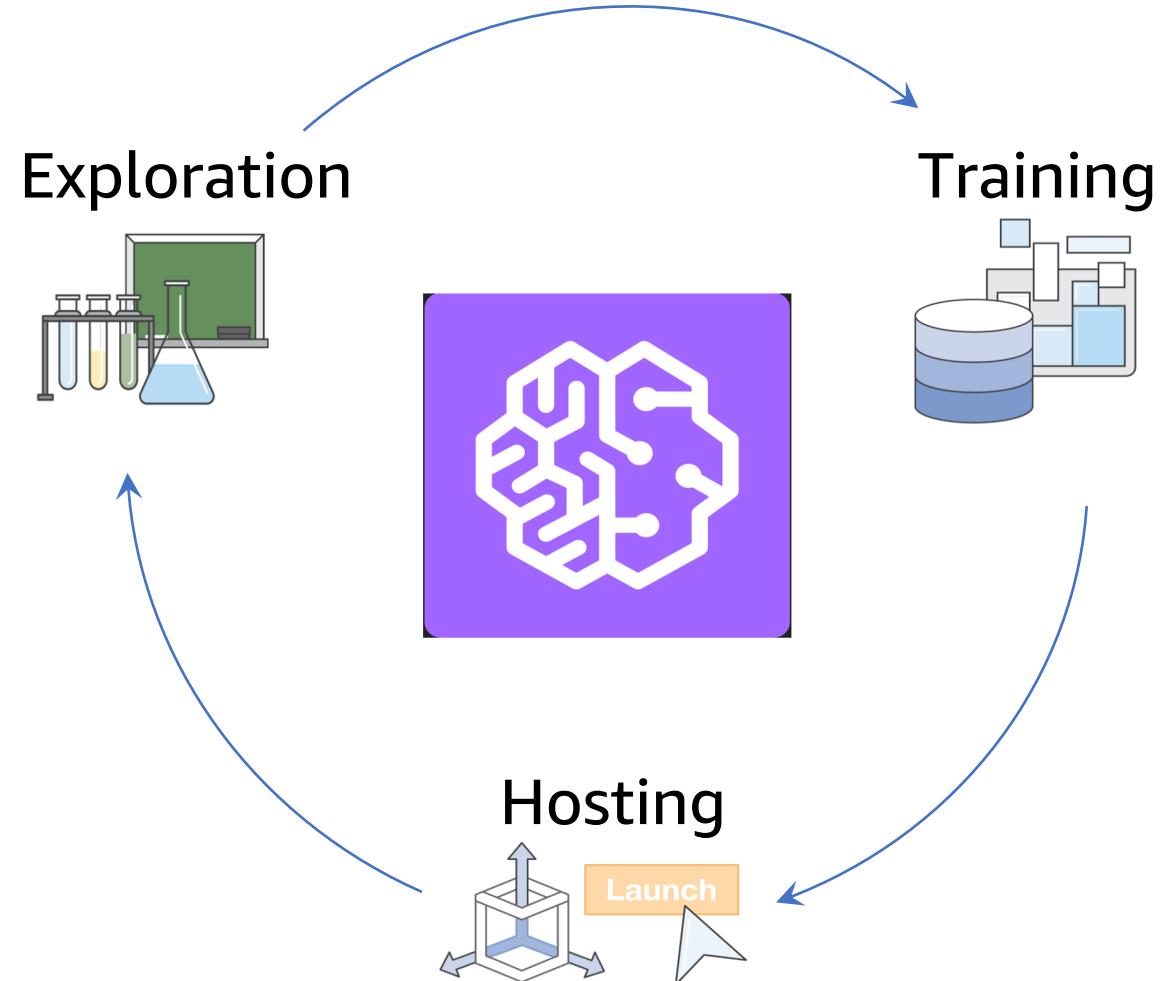
Edo Liberty, Bing Xiang, Baris Cuskon, Ramesh Nallapati, Phillip Gautier, Madhav Jha, Ran Ding, Tim Januschowski, David Selinas, Bernie Wang, Jan Gasthaus, Laurence Rouesnel, Amir Sadoughi, Piali Das, Julio Delgado Mangas, Yury Astashonok, Can Balioglu, Saswata Chakravarty, Alex Smola

- 1) ML Algorithms in The Cloud - New Challenges
- 2) SageMaker Algorithms - Architecture and Data Flow
- 3) Science of Streaming Algorithms – Advantages and Challenges
- 4) SageMaker Algorithms – Accurate, Fast, Scalable, and Easy to Use.
- 5) Deep dive – SageMaker K-Means

# ML Algorithms in The Cloud – New Challenges



# Lifecycle of a Machine Learning Project



# Small Data - Machine Learning



# Our Customers use ML at massive scale!



"Our data warehouse is 100TB and we are **processing 2TB daily**. We're running mostly gradient boosting (trees), LDA and K-Means clustering and collaborative filtering."

Shahar Cizer Kobrinsky, VP Architecture



"We process **3 million ad requests a second**, 100,000 features per request. That's 250 trillion per day. Not your run of the mill Data science problem!"

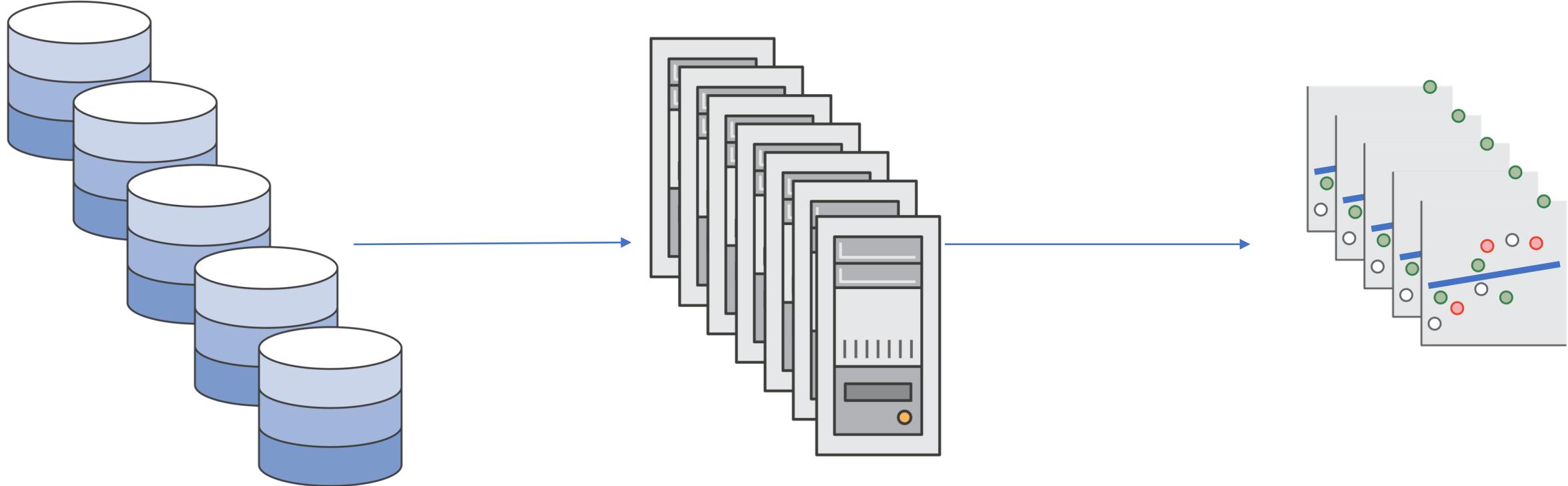
Bill Simmons, CTO



"We collect **160M events daily** in the ML pipeline and run training over the last 15 days and need it to complete in one hour. Effectively there's **100M features in the model**"

Valentino Volonghi, CTO

# Large Scale Machine Learning



# Large Scale Machine Learning

*dmlc*  
**XGBoost**



Spark + H<sub>2</sub>O  
SPARKLING  
**WATER**

Microsoft  
**CNTK**

The scikit-learn logo features two overlapping circles: a blue one on the left and an orange one on the right, both containing the word "learn". Above the orange circle is the text "scikit".

Spark  
MLlib

The mxnet logo consists of the word "mxnet" in a bold, white, sans-serif font, centered within a dark grey horizontal bar.

The GraphLab logo features a magenta silhouette of a dog running to the left, with the word "GraphLab" in a black serif font below it.



The Caffe2 logo features a stylized coffee cup icon with a plus sign above it, followed by the word "Caffe2" in a bold, dark font.

The Apache Hadoop logo features a yellow elephant icon followed by the word "hadoop" in a blue, lowercase, sans-serif font.

The Apache Mahout logo features a yellow elephant icon with a small figure riding on its back, followed by the word "mahout" in a blue, lowercase, sans-serif font.

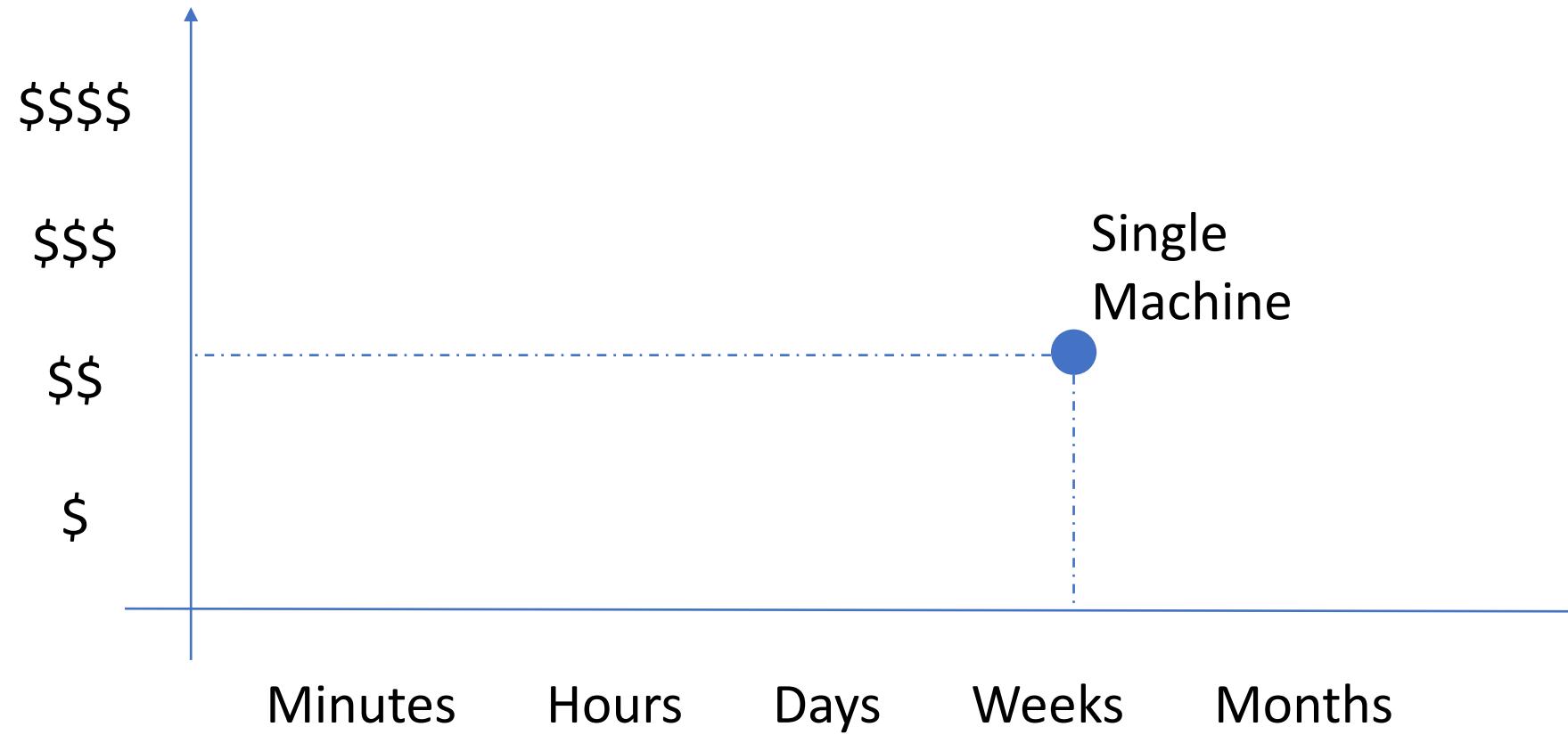


The Amazon logo, featuring the word "amazon" in a lowercase, italicized font with a red arrow underneath.

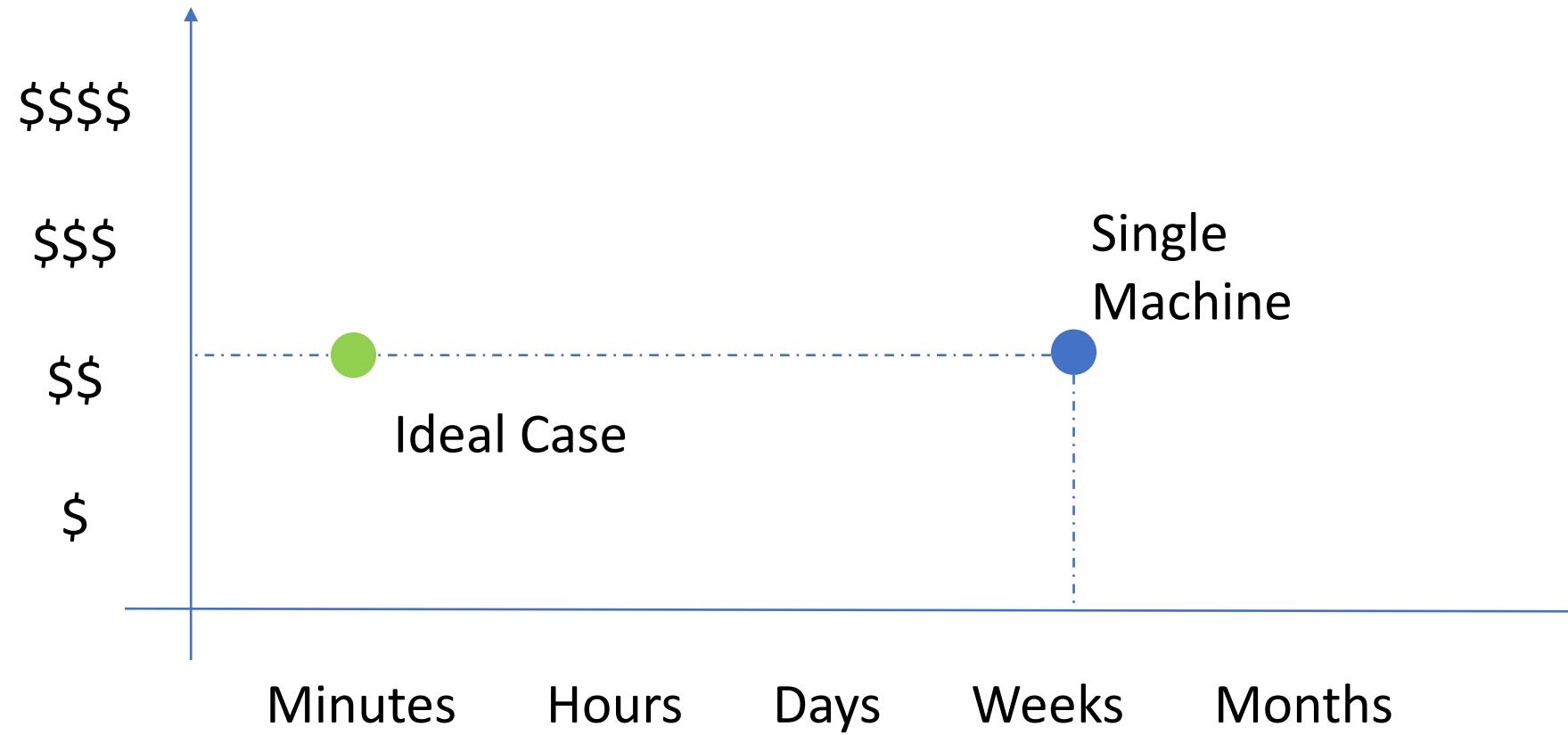
© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

The AWS logo, featuring the word "aws" in a lowercase, bold, dark font with an orange arrow underneath.

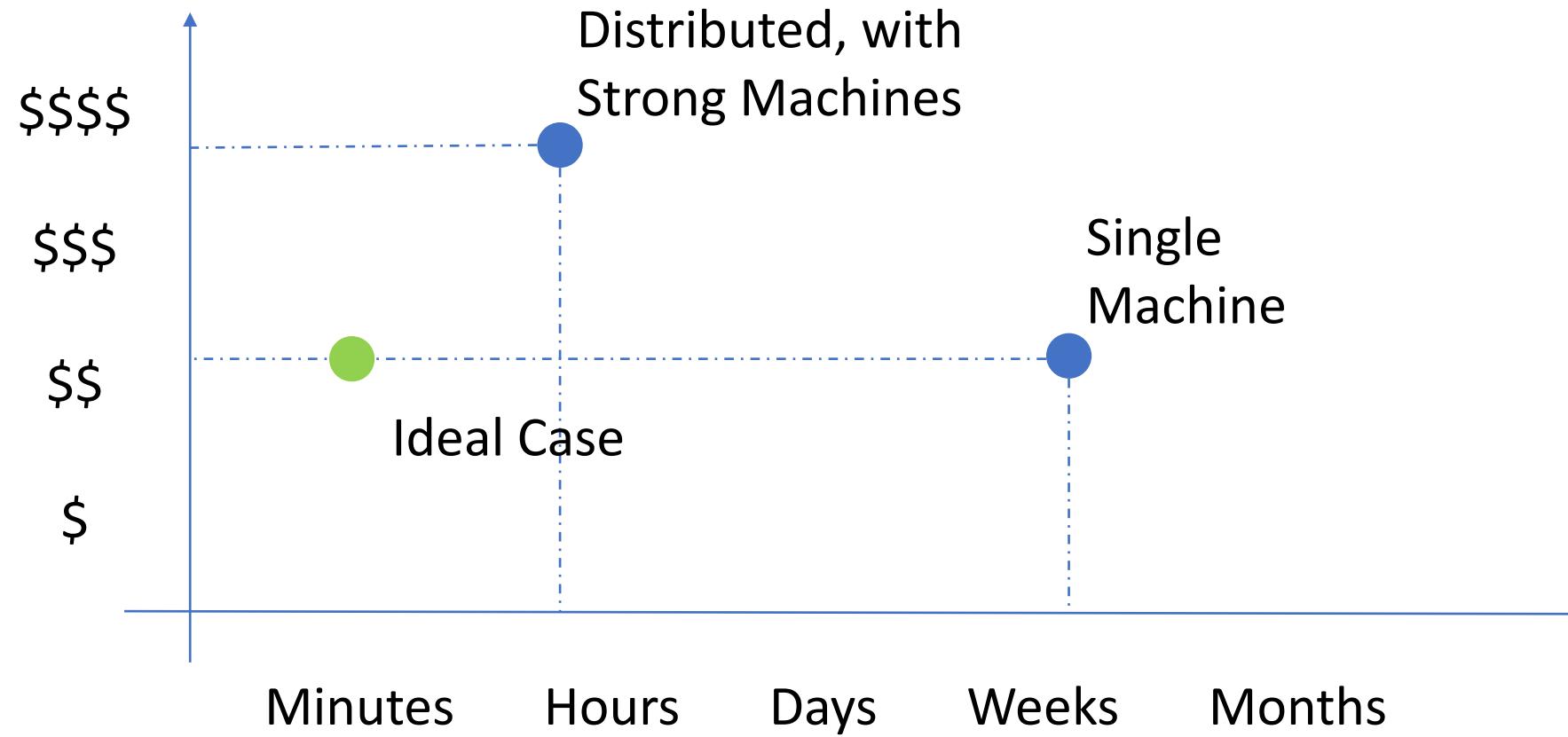
# Cost vs. Time



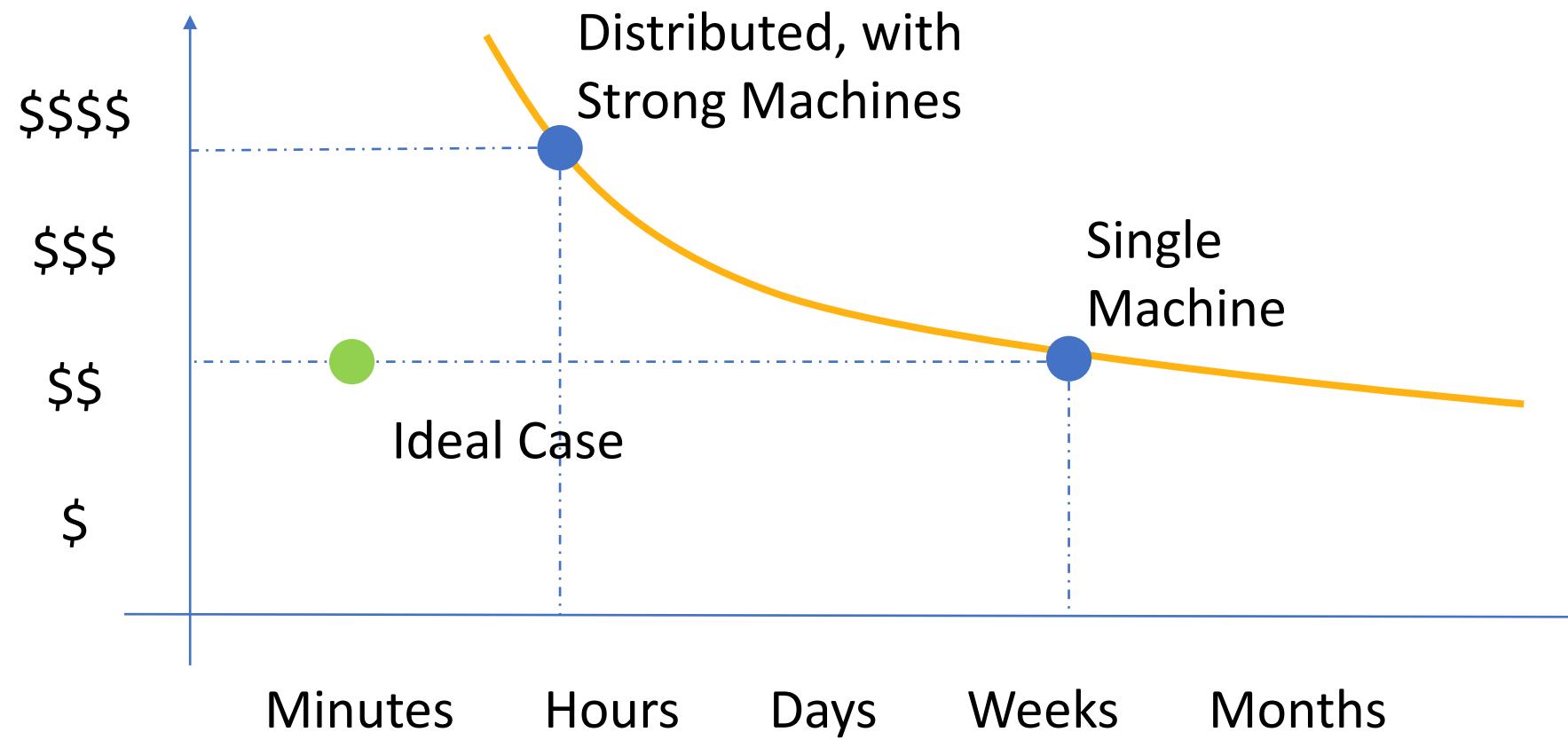
# Cost vs. Time



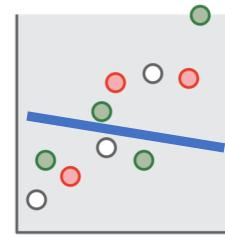
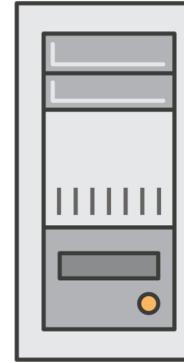
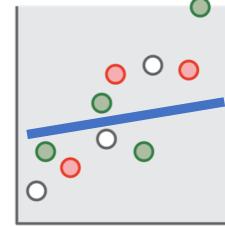
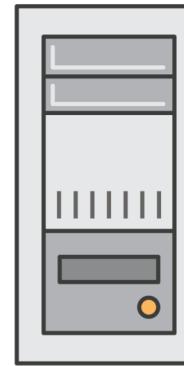
# Cost vs. Time



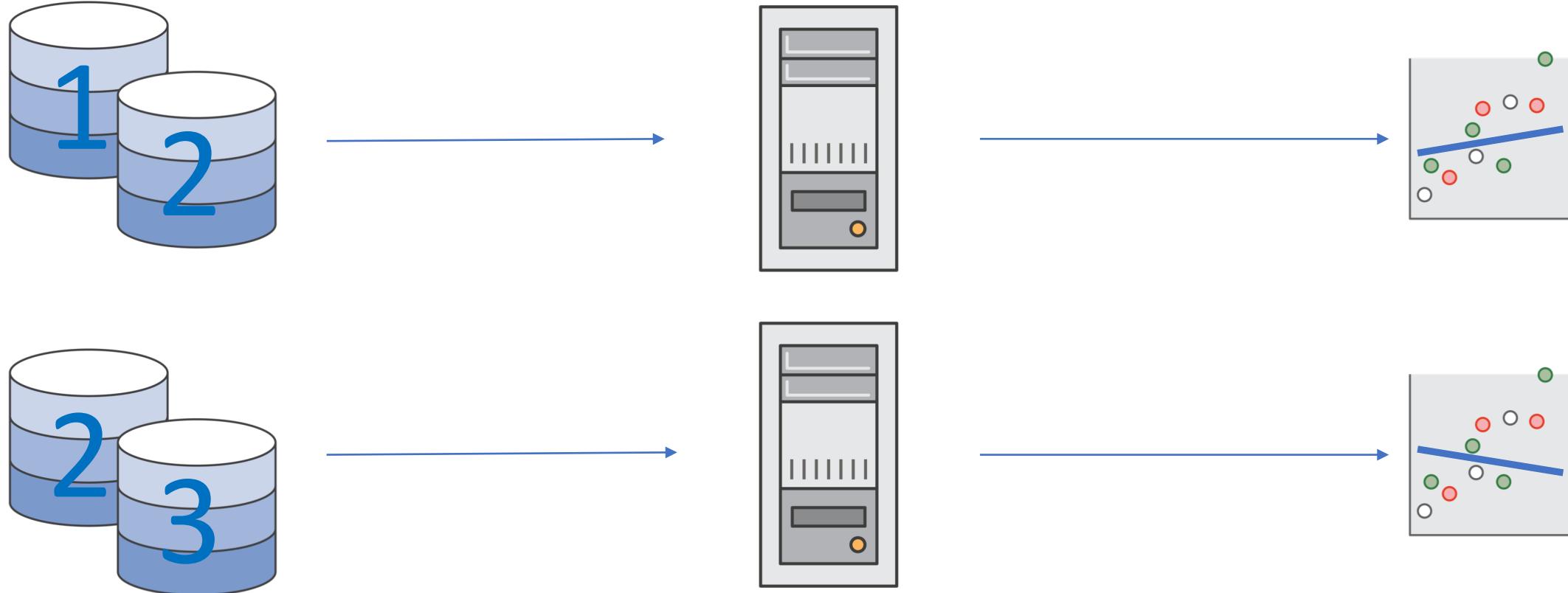
# Cost vs. Time



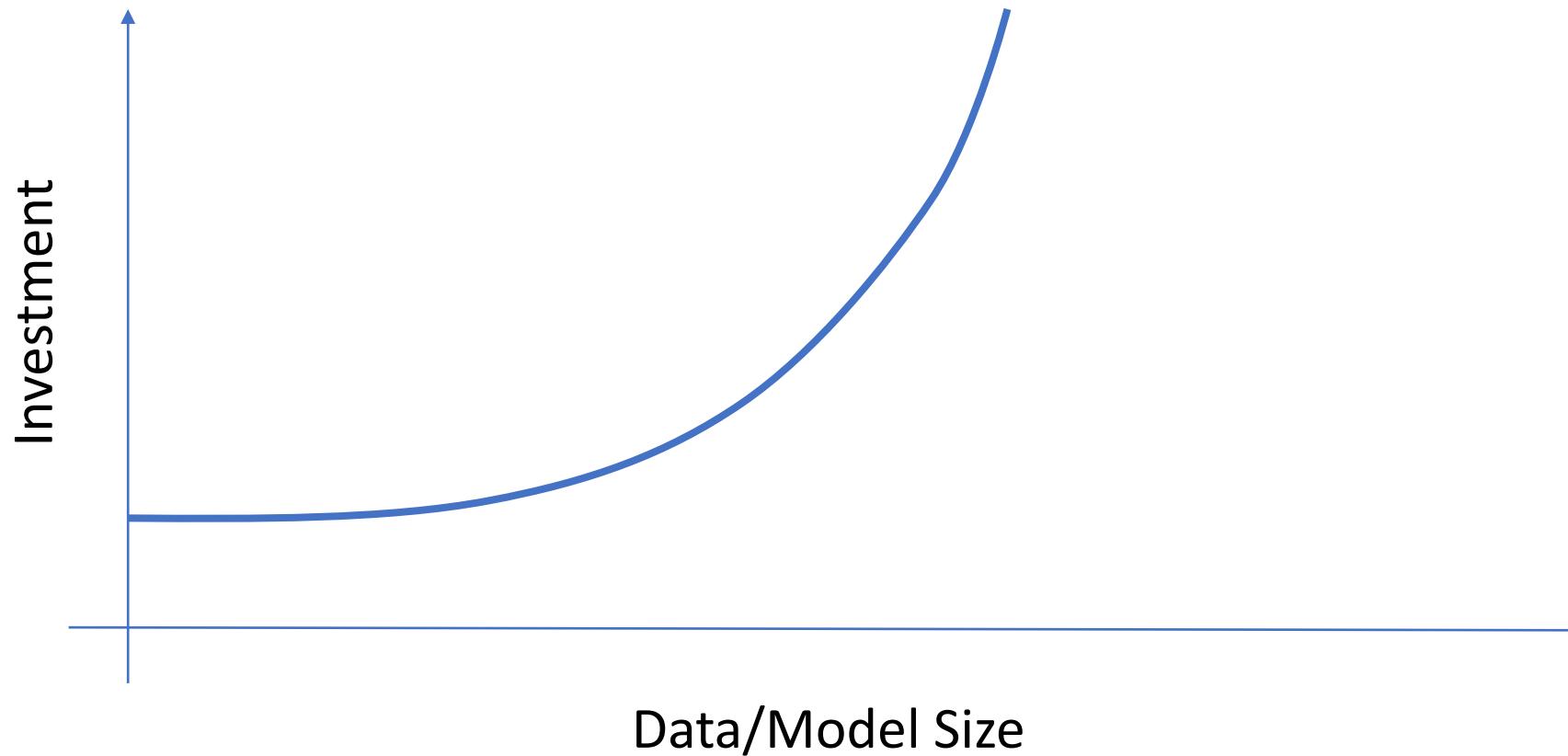
# Model Selection



# Incremental Training



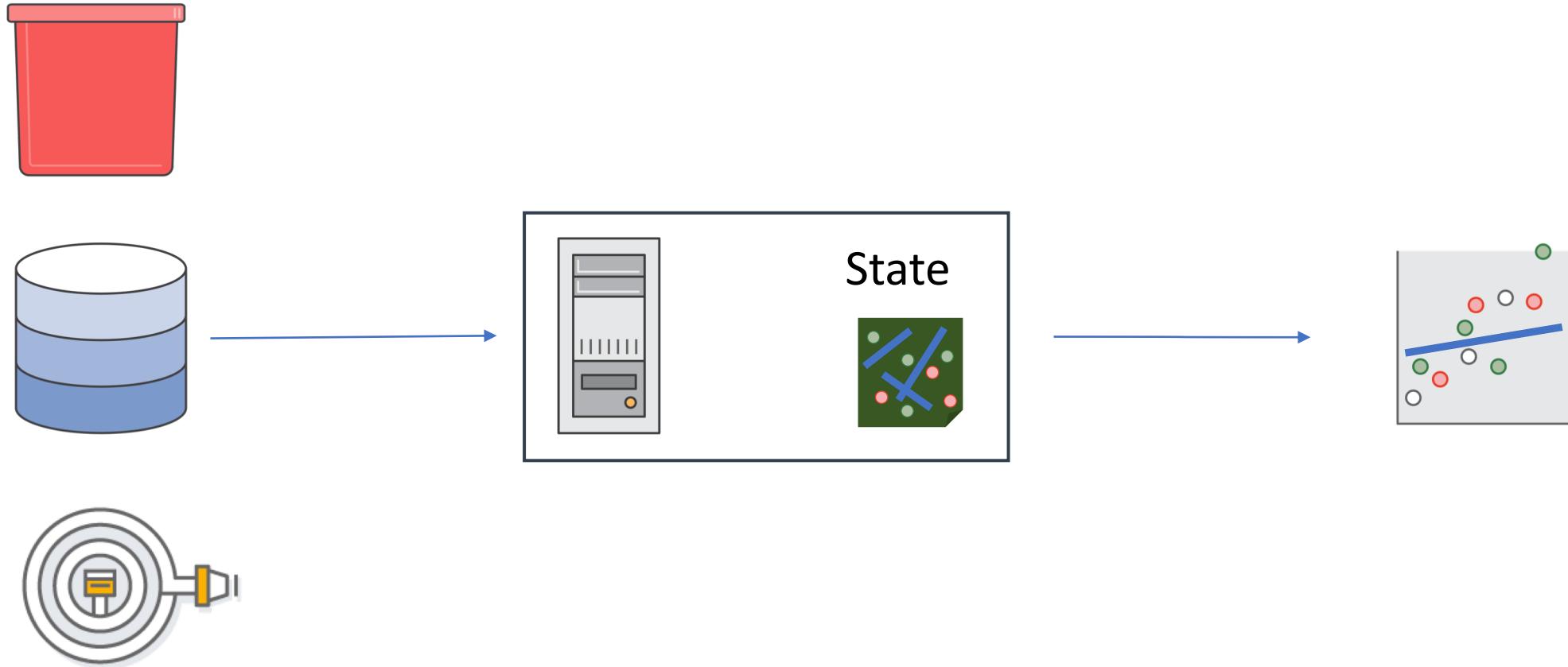
# Production Readiness



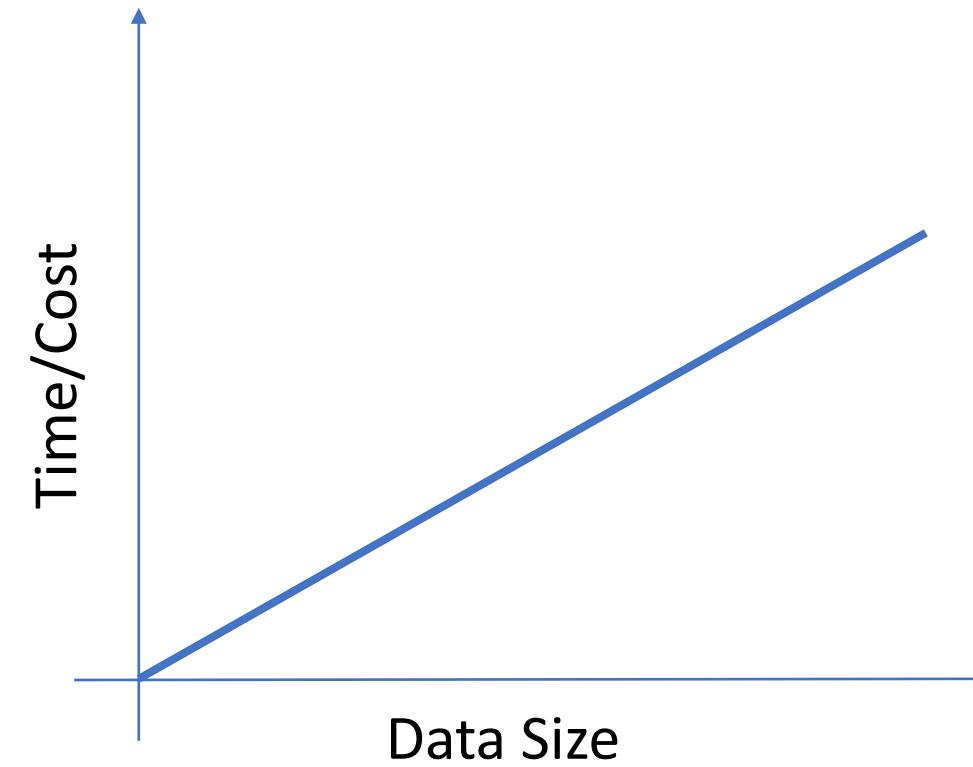
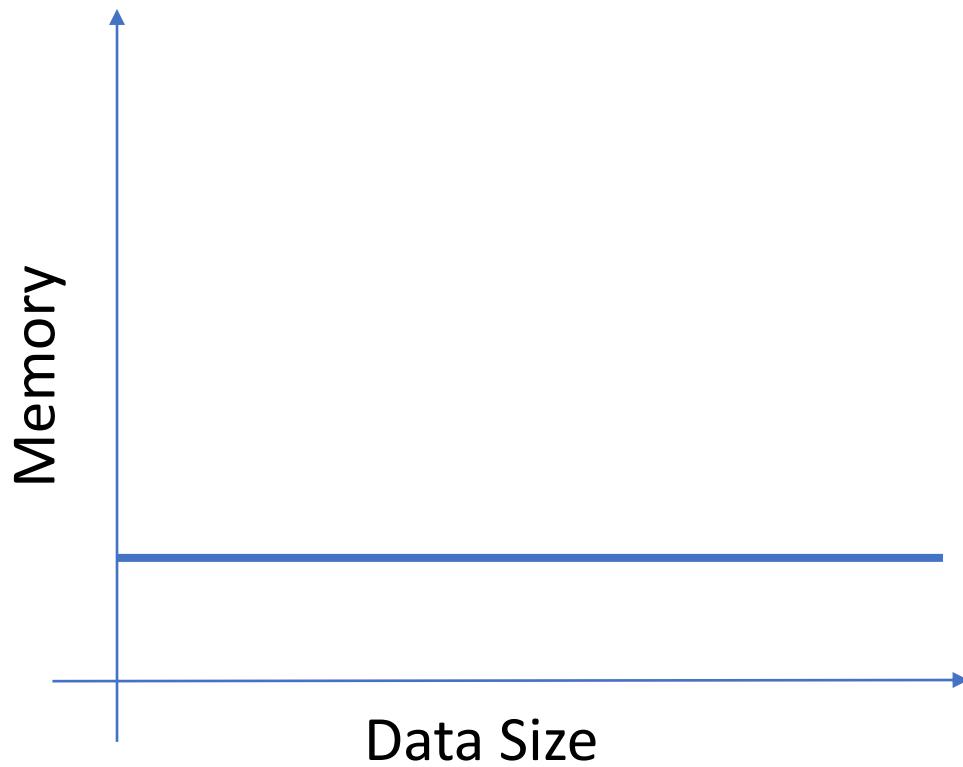
# SageMaker Algorithms - Architecture and Data Flow



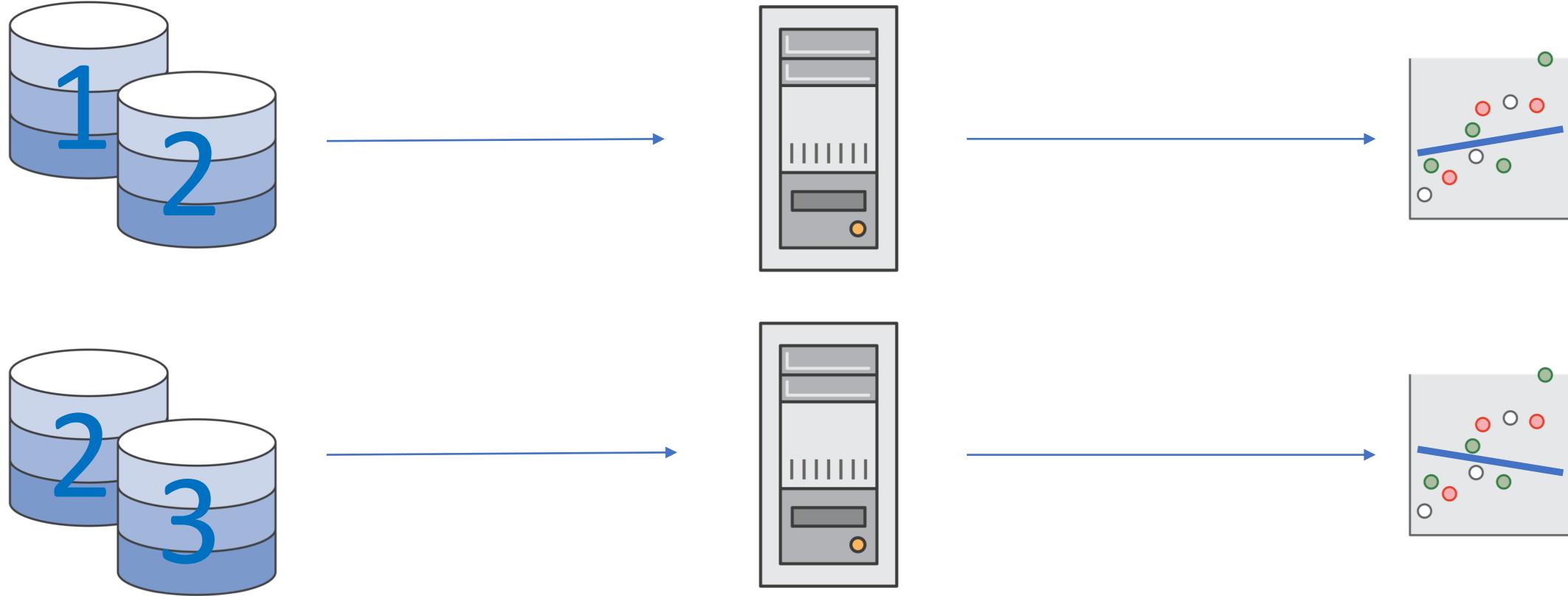
# Streaming



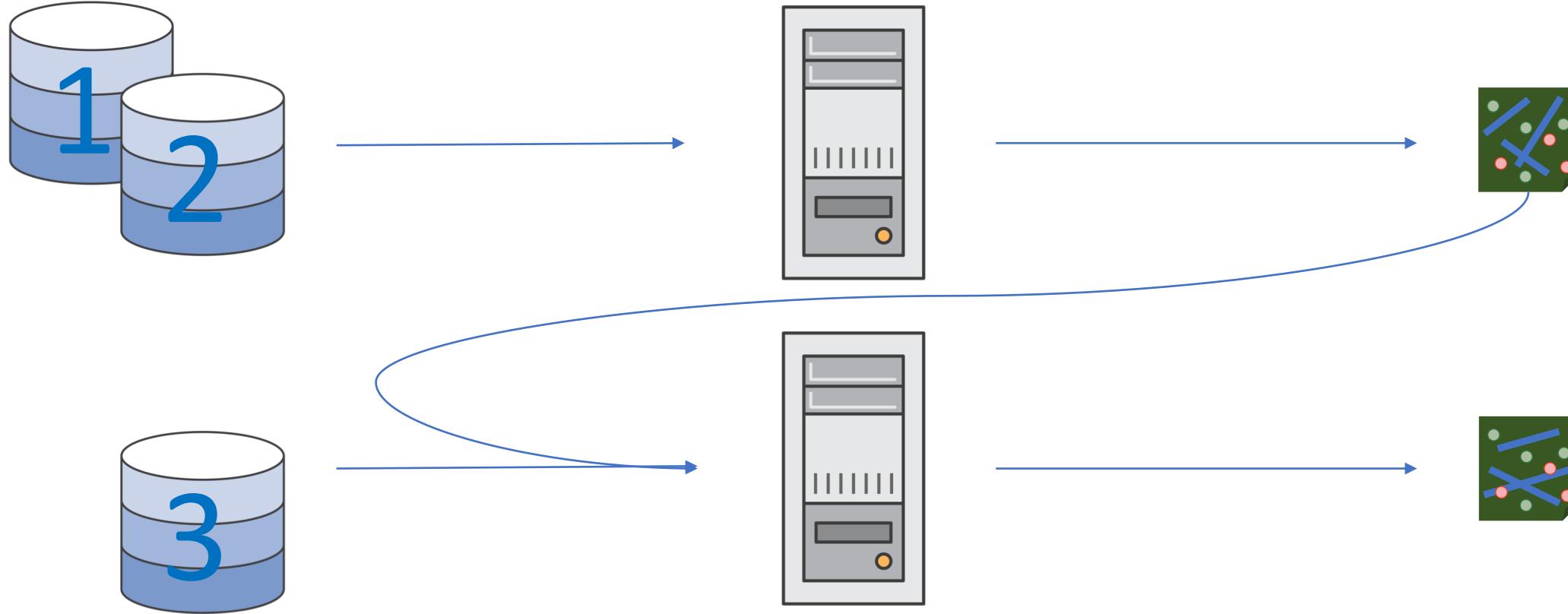
# Streaming



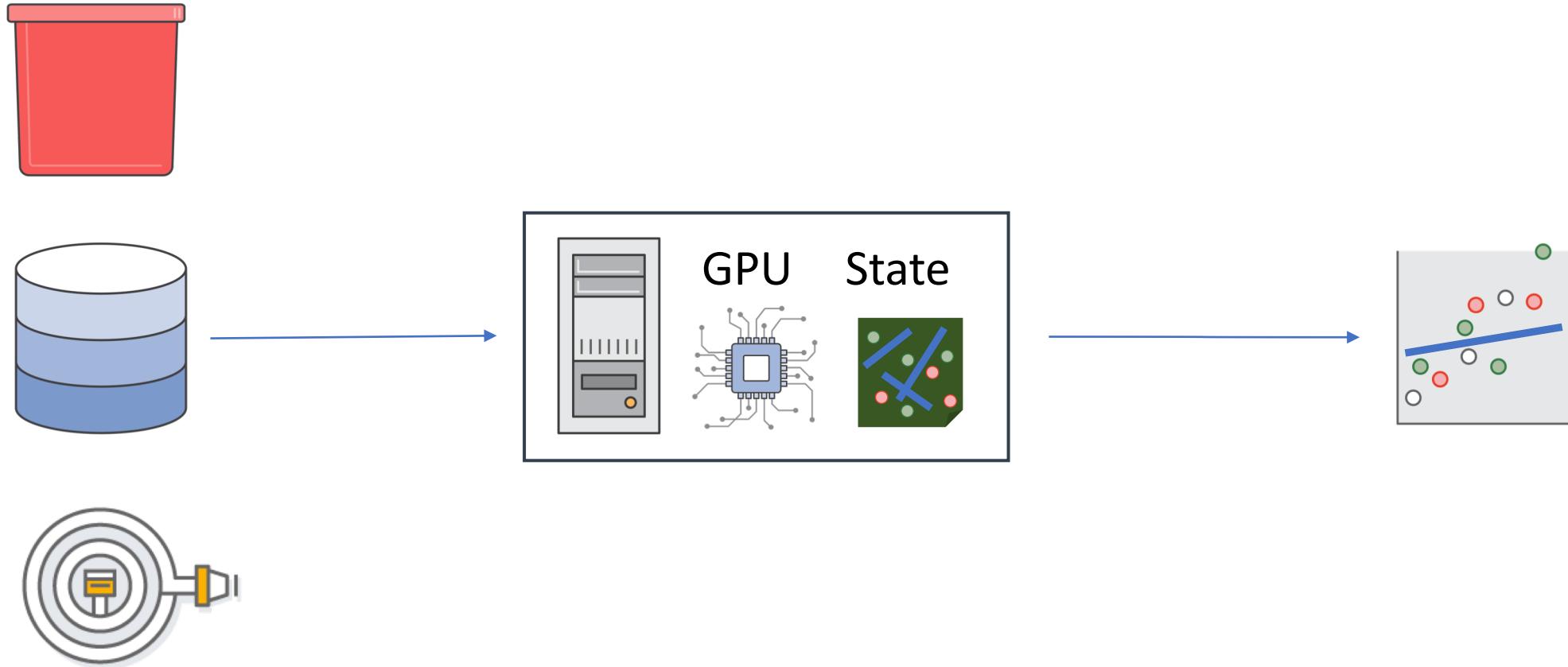
# Incremental Training



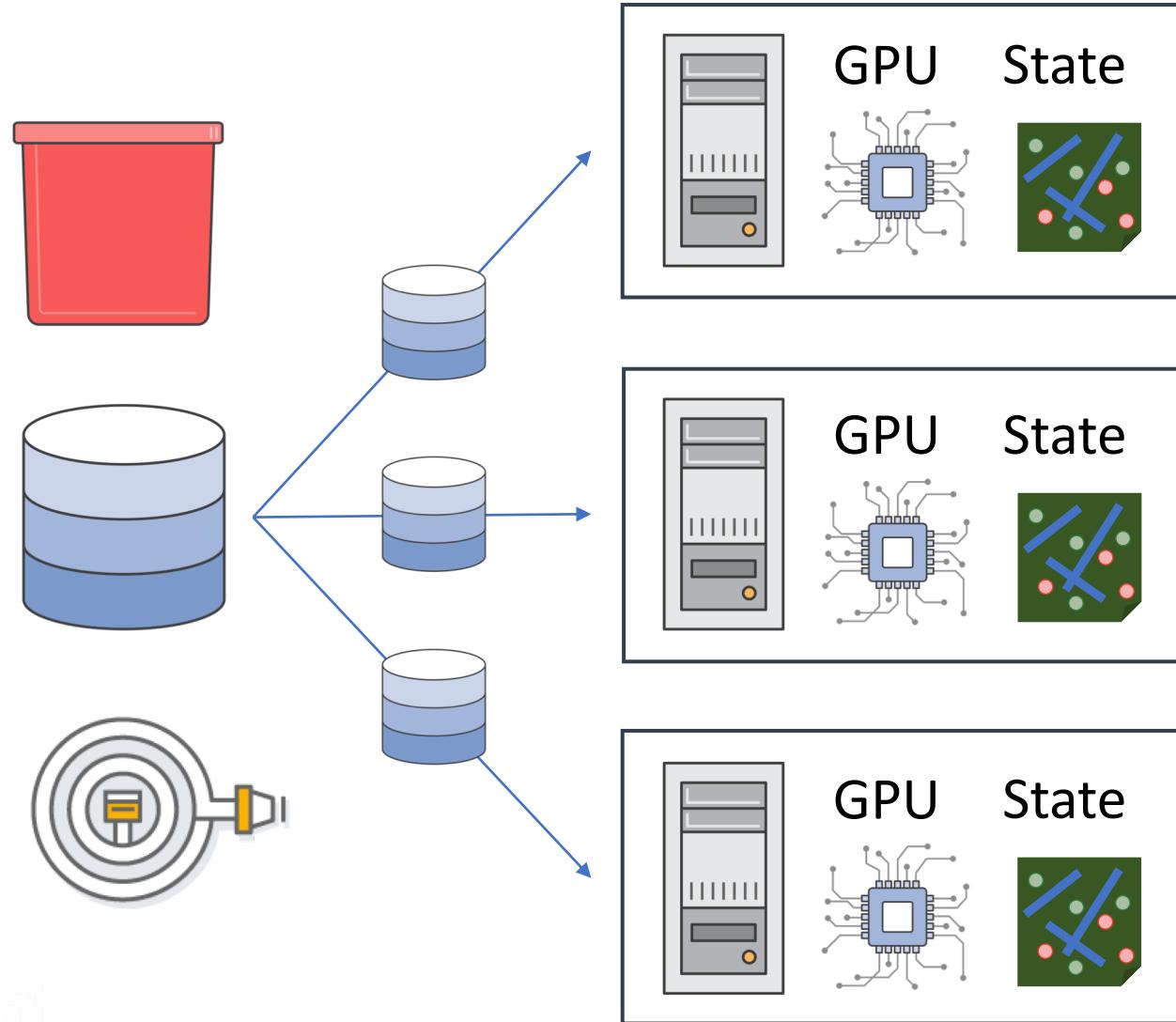
# Incremental Training



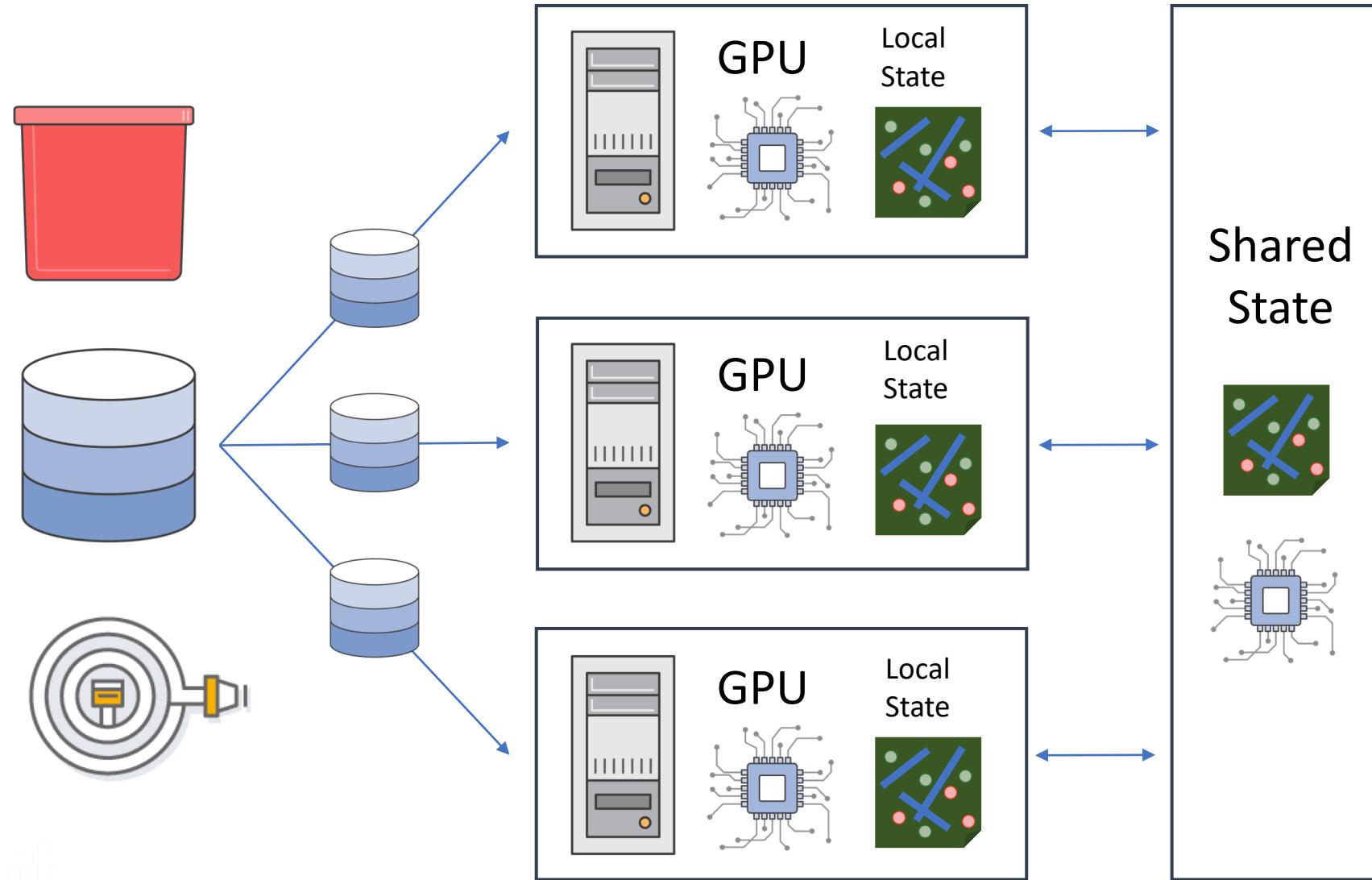
# GPU/CPU



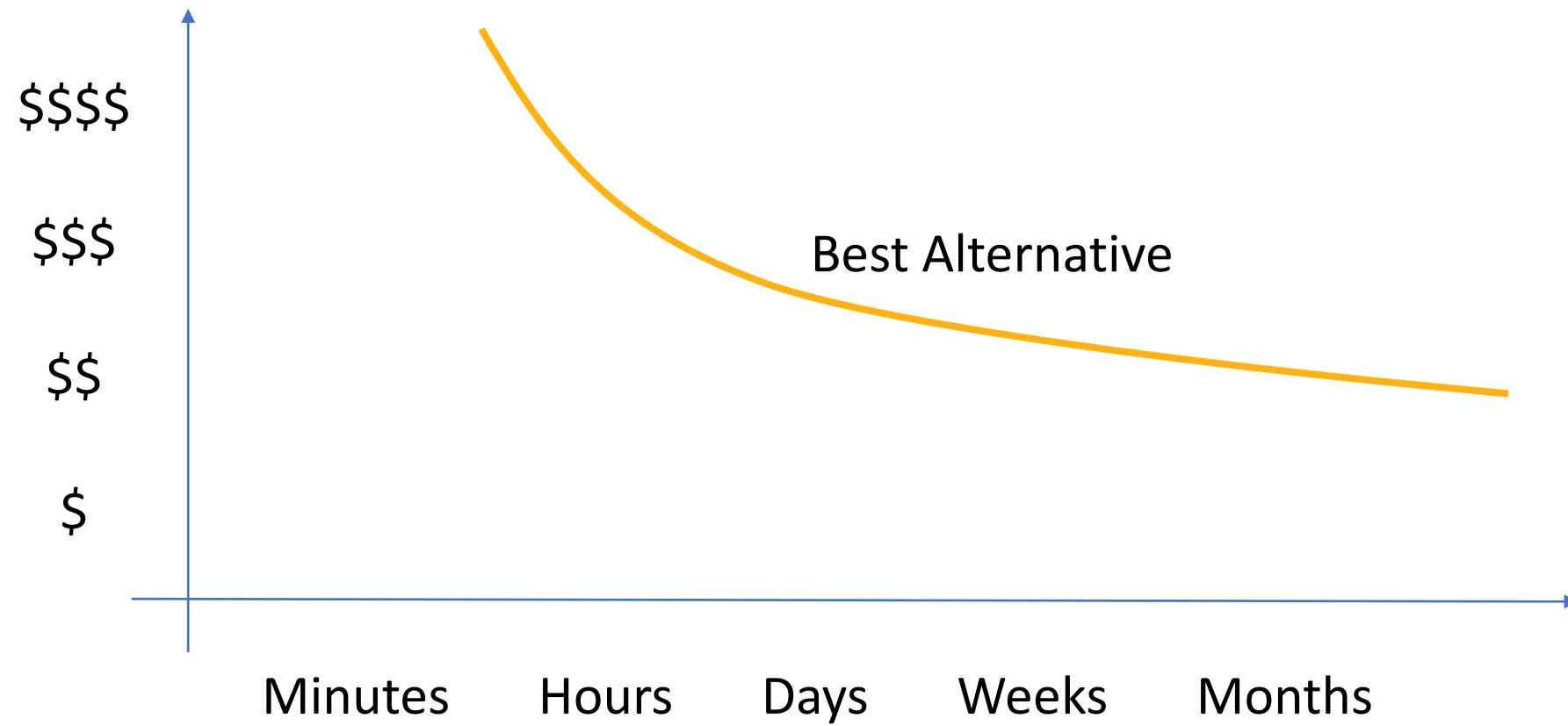
# Distributed



# Parameter Server – distributed (k,v) store.



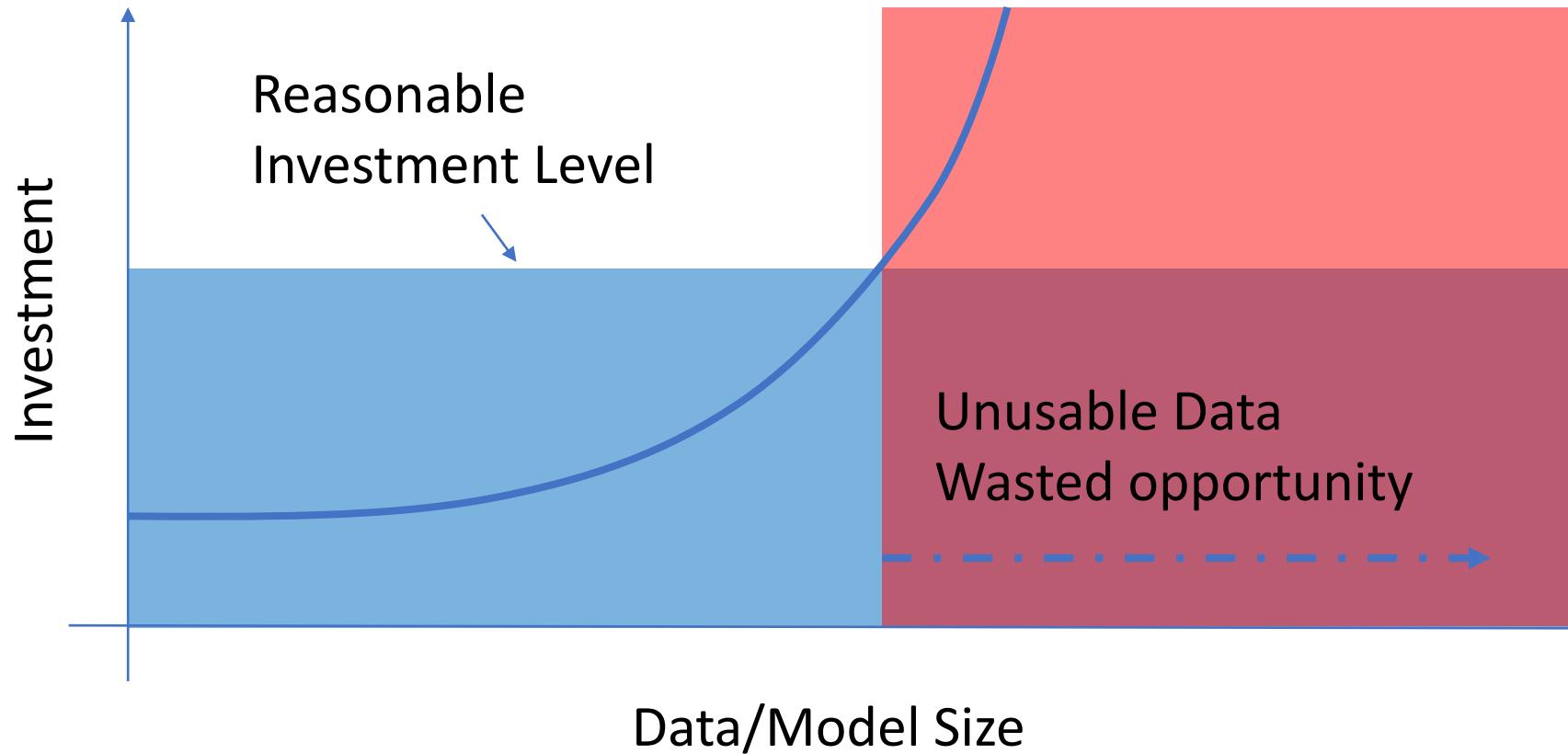
# Cost vs. Time



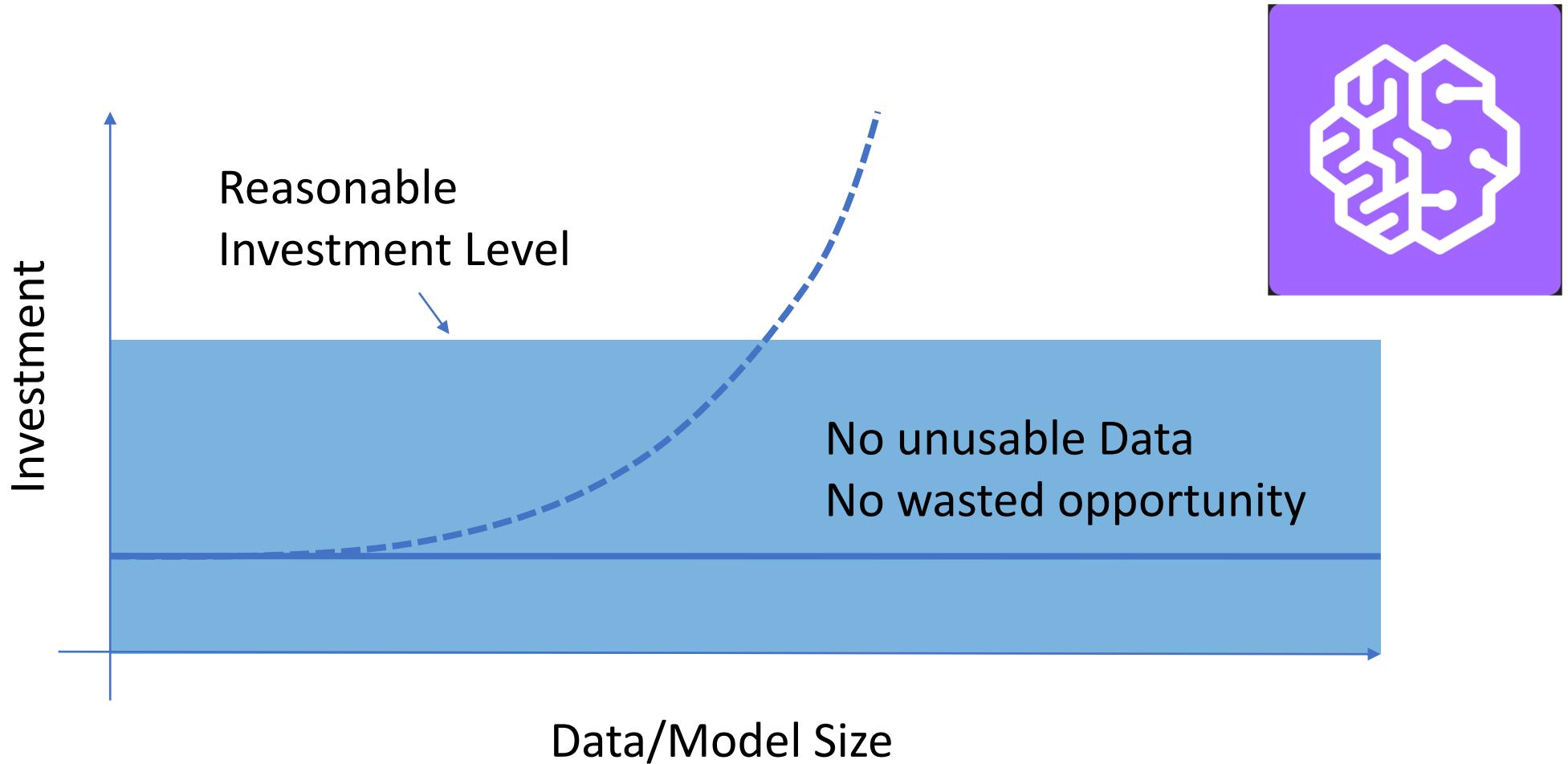
# Cost vs. Time



# Production Readiness



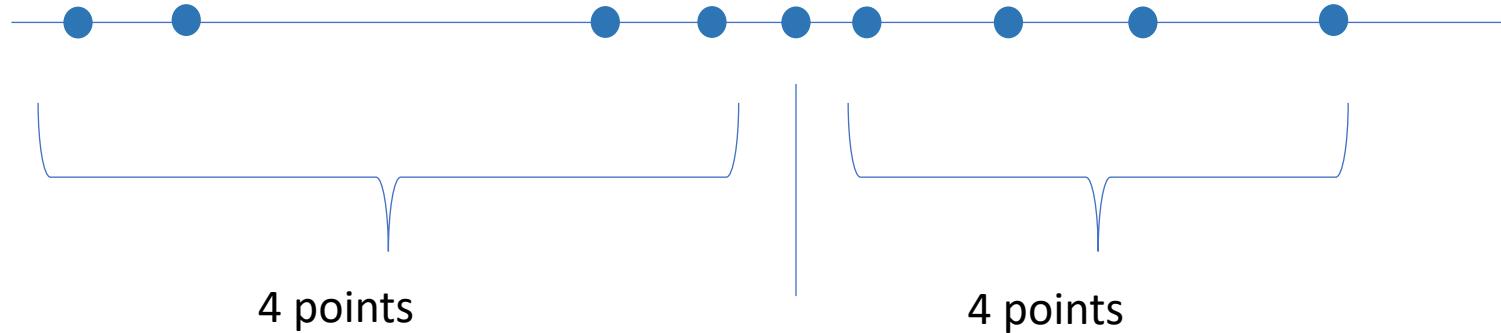
# Production Readiness



# Science of Streaming Algorithms – Advantages and Challenges



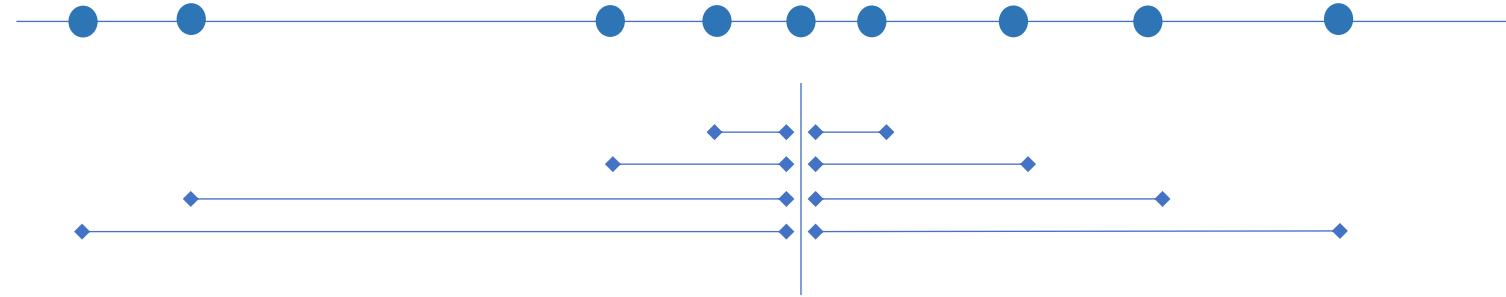
# Simple Problems Are Unsolvable



Finding the exact median in a stream is impossible!

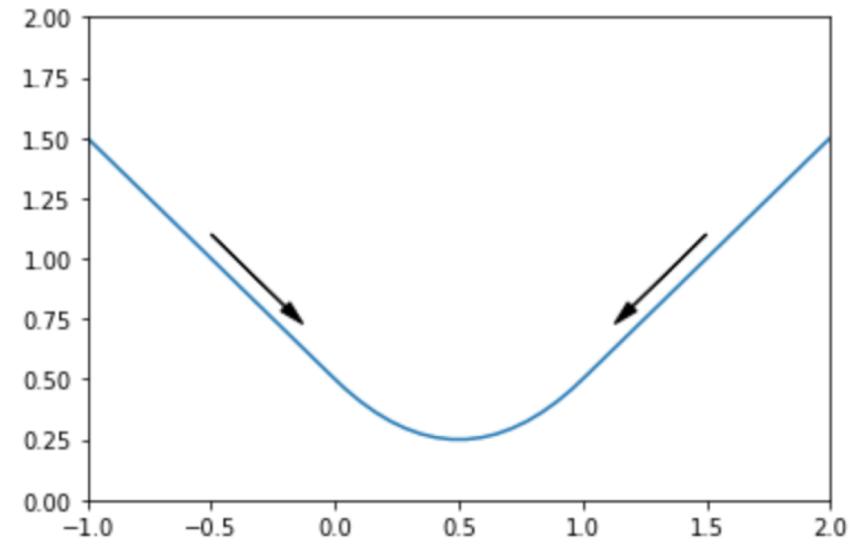
- After seeing half the items, each one of them might still be the median.
- The algorithm must remember all of them.
- It cannot have a fixed memory footprint.

# Gradient Descent



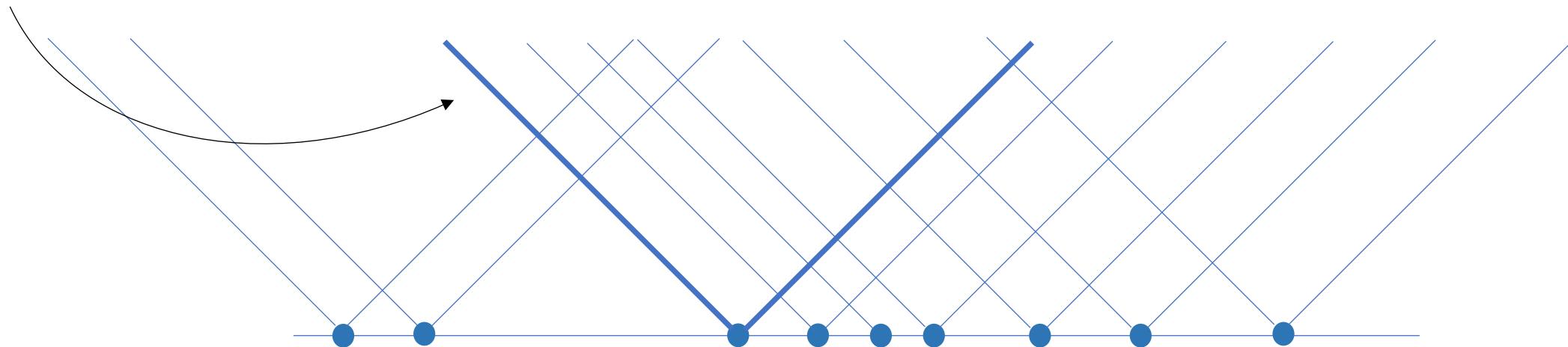
$$f(x) = \frac{1}{n} \sum_{i=1}^n |x - x_i|$$

$$m = \arg \min_x f(x)$$



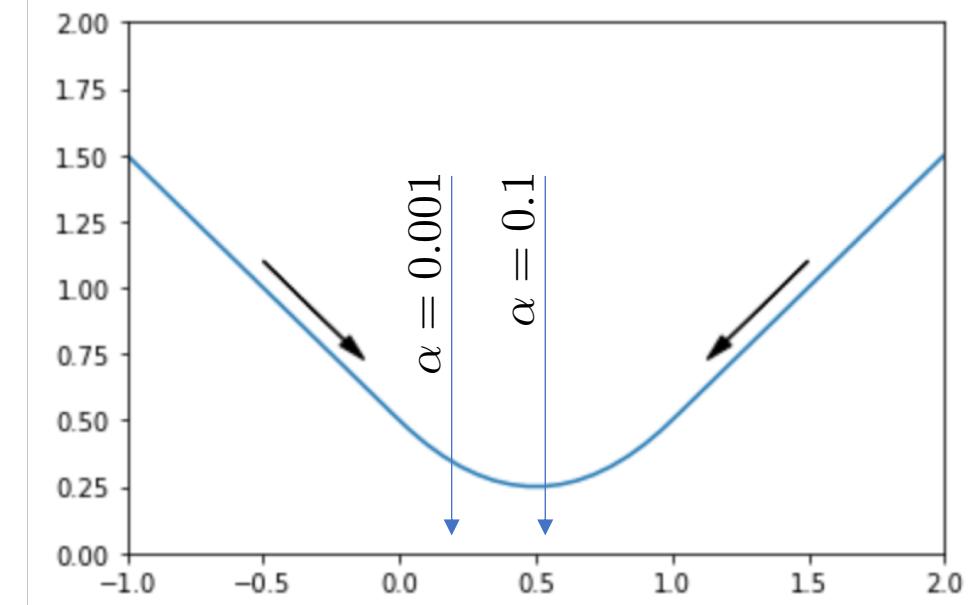
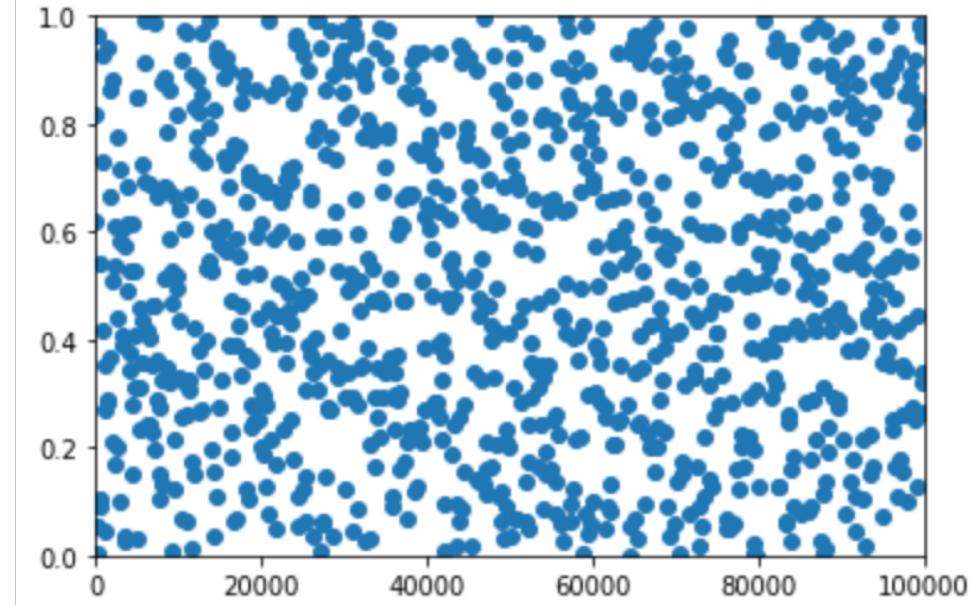
# Stochastic Gradient Descent

$$f_i = |x_i - x| \implies \mathbb{E}_i[f_i] = f \implies \mathbb{E}_i[f'_i] = f'$$



$$m_t = \begin{cases} m_{t-1} + \alpha/\sqrt{t} & \text{if } m_{t-1} < x_i \\ m_{t-1} - \alpha/\sqrt{t} & \text{if } m_{t-1} > x_i \end{cases}$$

# SGD – Parameter Tuning



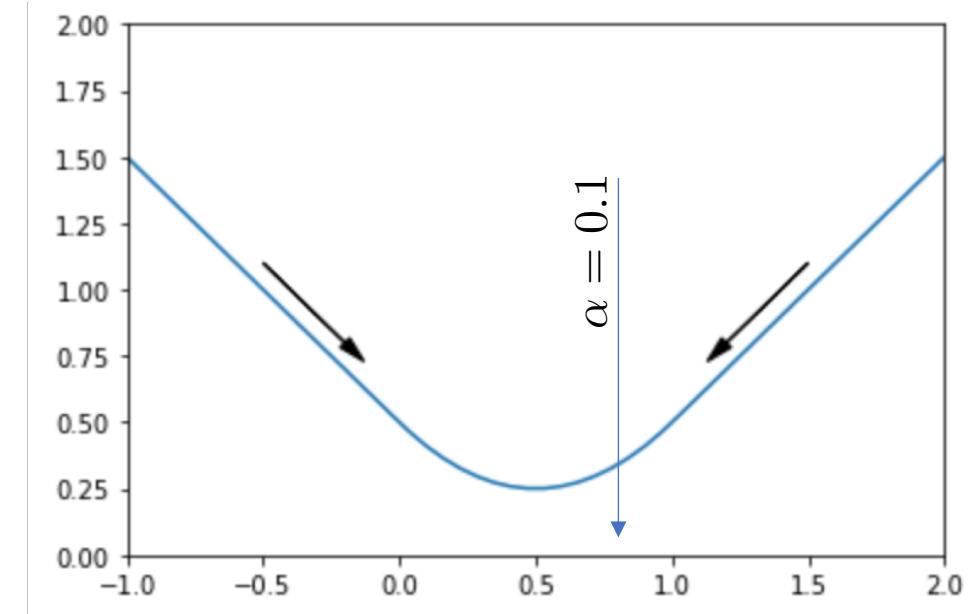
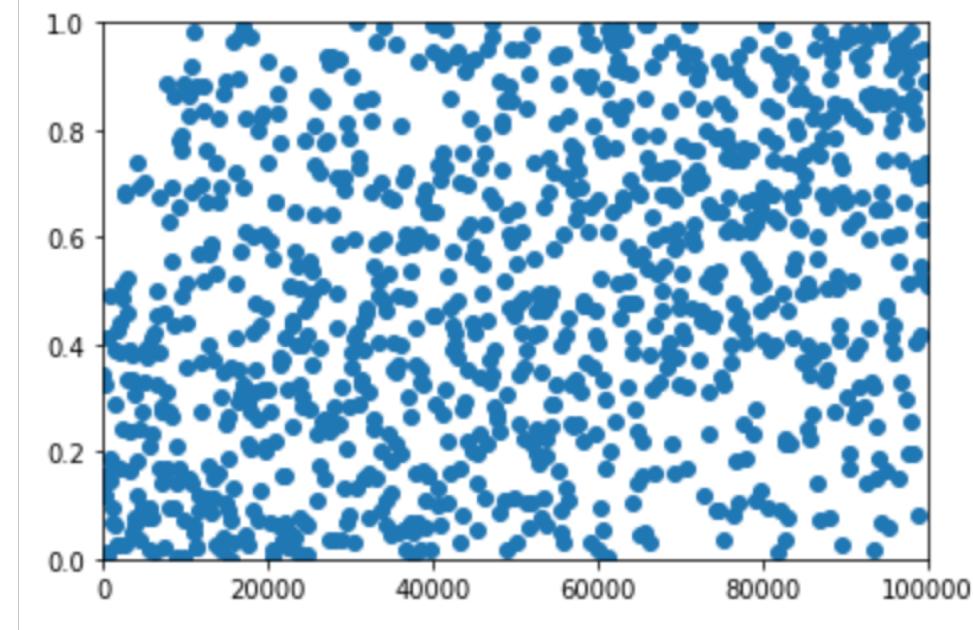
$$m_t = \begin{cases} m_{t-1} + \alpha/\sqrt{t} & \text{if } m_{t-1} < x_i \\ m_{t-1} - \alpha/\sqrt{t} & \text{if } m_{t-1} > x_i \end{cases}$$

Frugal Streaming for Estimating Quantiles: One (or two) memory suffices: Qiang Ma, S. Muthukrishnan, Mark Sandler

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



# SGD – Distribution Drift



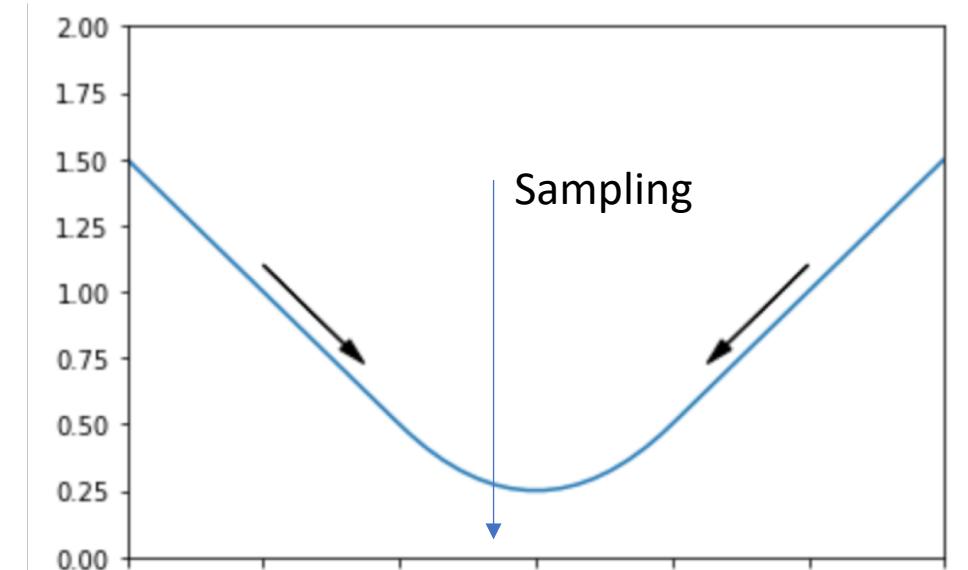
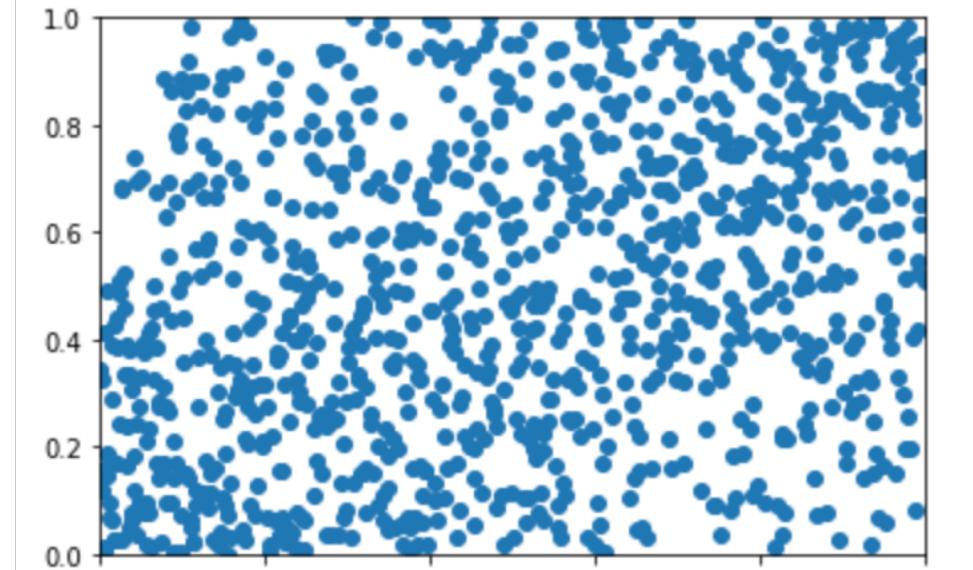
$$m_t = \begin{cases} m_{t-1} + \alpha/\sqrt{t} & \text{if } m_{t-1} < x_i \\ m_{t-1} - \alpha/\sqrt{t} & \text{if } m_{t-1} > x_i \end{cases}$$

Frugal Streaming for Estimating Quantiles: One (or two) memory suffices: Qiang Ma, S. Muthukrishnan, Mark Sandler

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



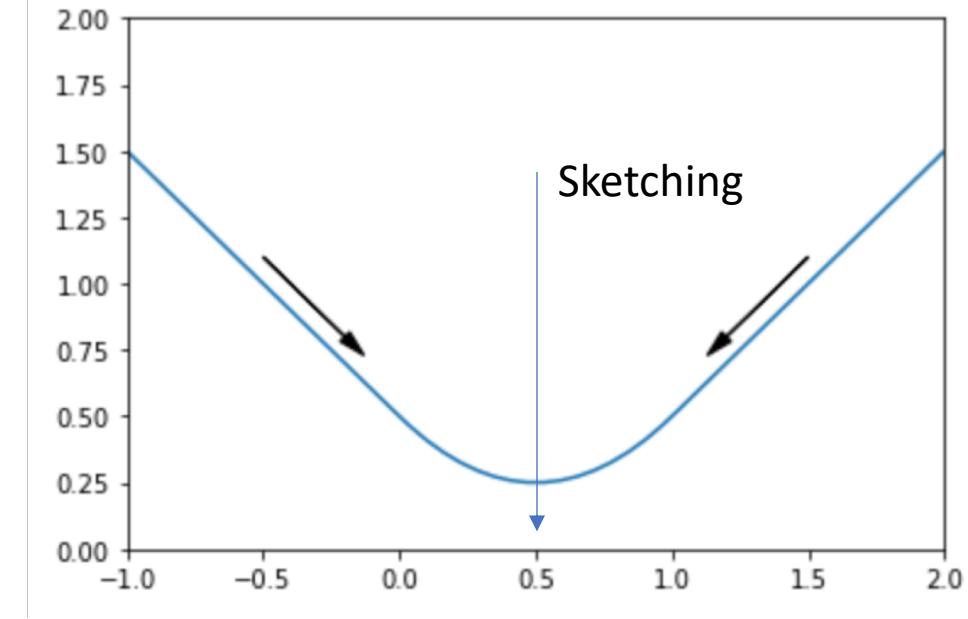
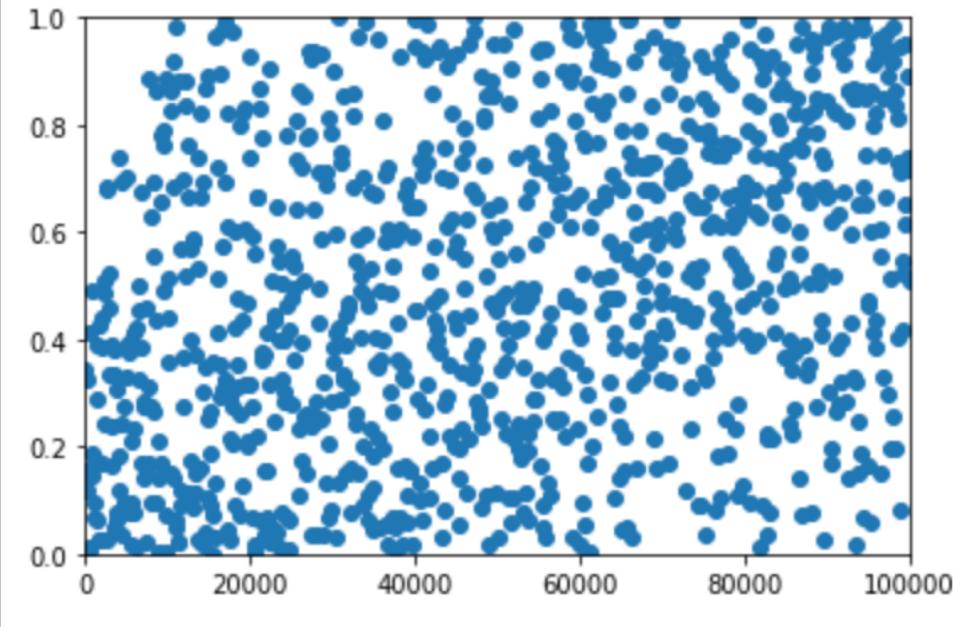
# Median - Sampling Algorithm



## Sampling Algorithm:

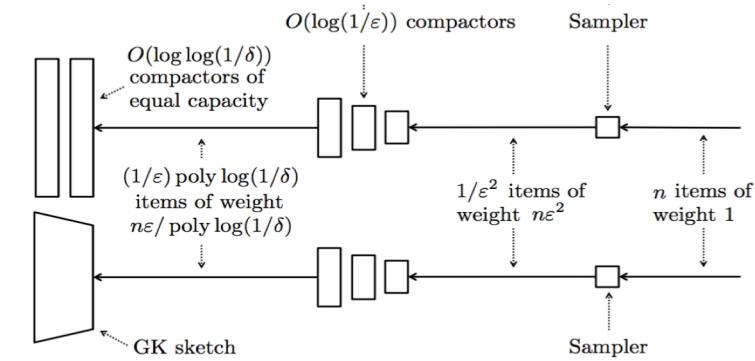
- 1) Reservoir Sample  $k$  points from the data
- 2) Return the median of the sample

# Median – Sketching Algorithm

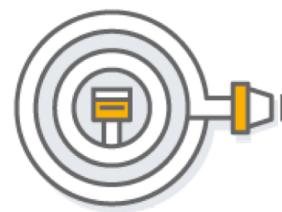


## Sketching Algorithm

- 1) Too complex to explain here...
- 2) Optimal Quantile Approximation in Streams;  
Zohar Karnin, Kevin Lang, Edo Liberty



# Framework for Streaming Algorithms



$x_1, x_2, \dots, x_n$



$$\frac{1}{n} \sum_{i=1}^n x_i$$

# Average: Single Machine

```
running_sum = 0.0  
items_seen = 0      } Initialize()
```

```
for x in stream:  
    running_sum += x  
    items_seen += 1      } update(x)
```

```
average = running_sum / items_seen } finalize()
```

# Average: Single Machine

```
state.initialize()  
for x in stream:  
    state.update(x)  
return state.finalize()
```

# Average: Distributed

```
In parallel for w in workers:  
    state[w].initialize()  
    for x in stream[w]:  
        state[w].update(x)
```

# Average: Distributed

In parallel for w in workers:

    state[w]. initialize()

    for x in stream[w]:

        state[w]. update(x)

running\_sum = 0      }  
items\_seen = 0      } initialize()

for w in workers:

    running\_sum += state[w]. running\_sum

    items\_seen += state[w]. items\_seen

} merge(state[w])

average = running\_sum / items\_seen      } finalize()



# Average: Distributed

```
In parallel for w in workers:
```

```
    state[w].initialize()
```

```
    for x in stream[w]:
```

```
        state[w].update(x)
```

```
master_state.initialize()
```

```
for w in workers:
```

```
    master_state.merge(state[w])
```

```
return master_state.finalize()
```

# Distributed Learning

Fully mergeable

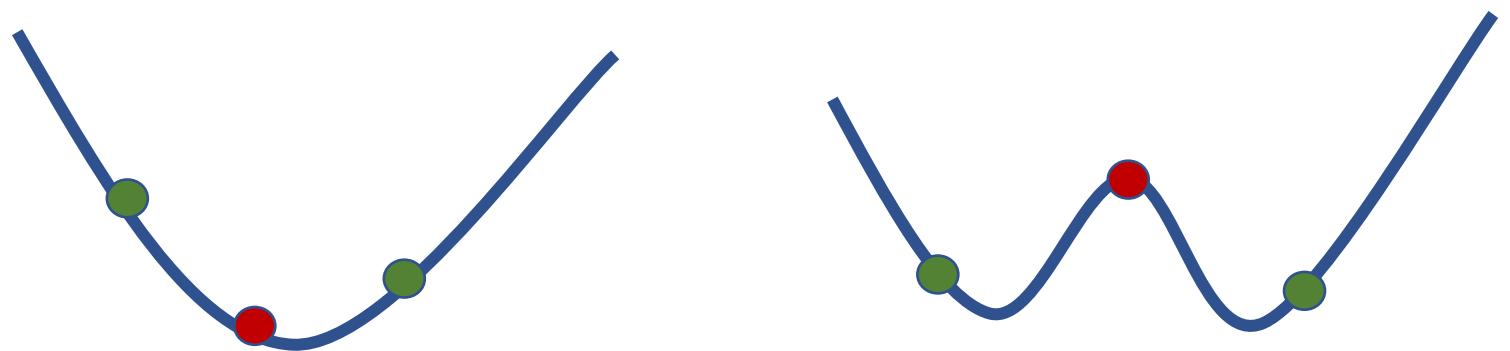


PCA:  
requires  $E[x_t x_t^T]$

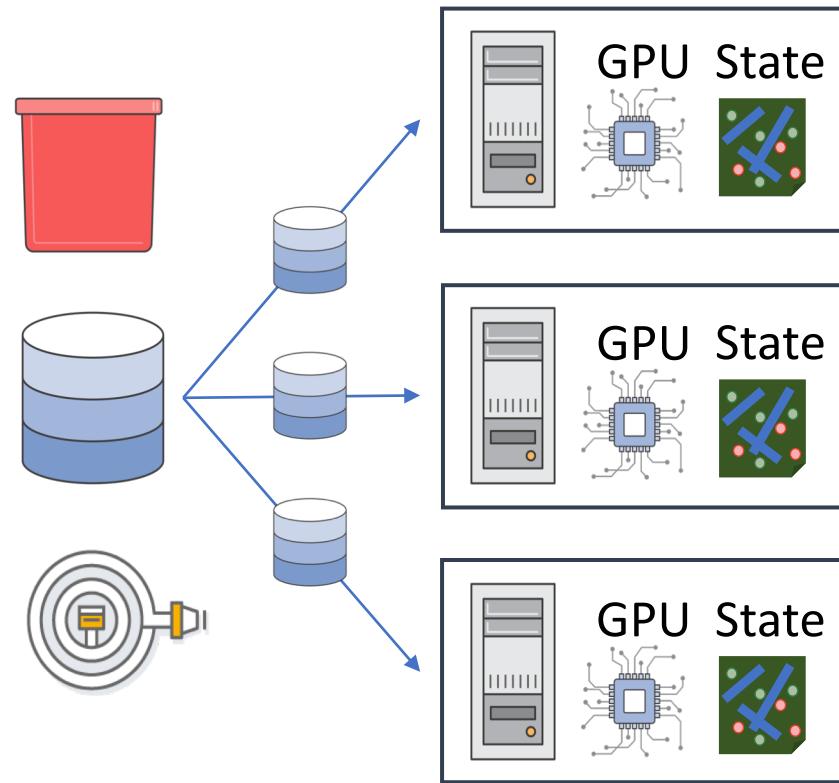
Convex SGD

Not mergeable

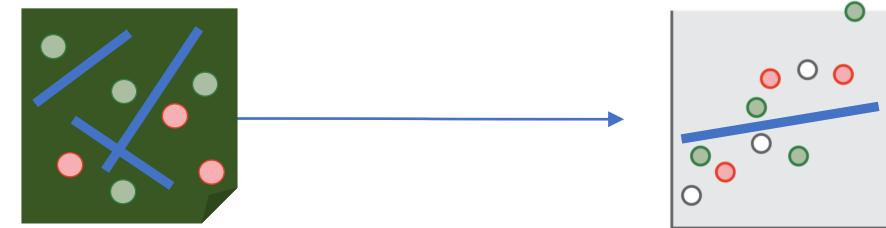
Non-convex SGD



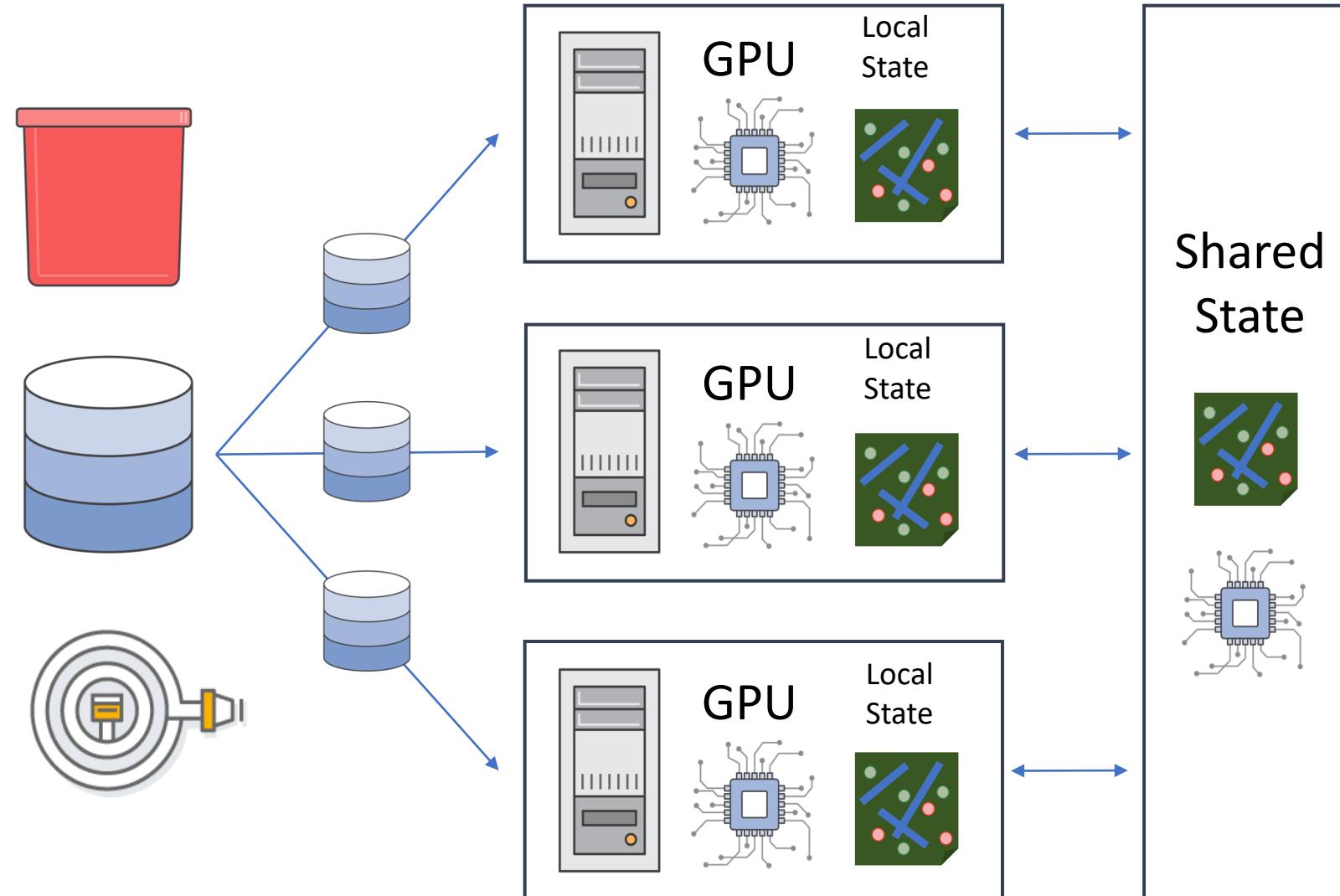
# Mergeable



Merged  
State



# Not mergeable



# Distributed SGD

```
local_state.initialize()
```

Obtain recent model weights from parameter server

```
for batch in stream:
```

```
    local_state.pull()
```

```
    delta = local_state.update(batch)
```

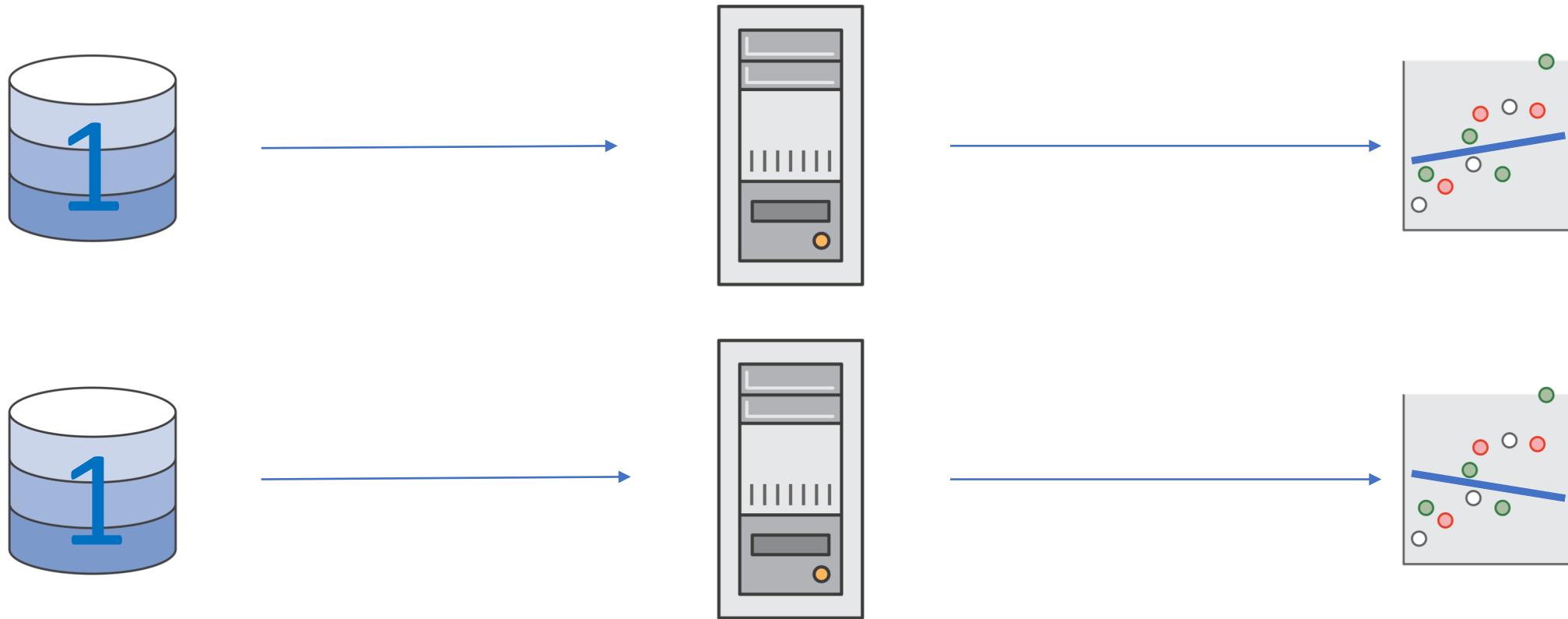
```
    local_state.push(delta)
```

Compute gradient on current batch

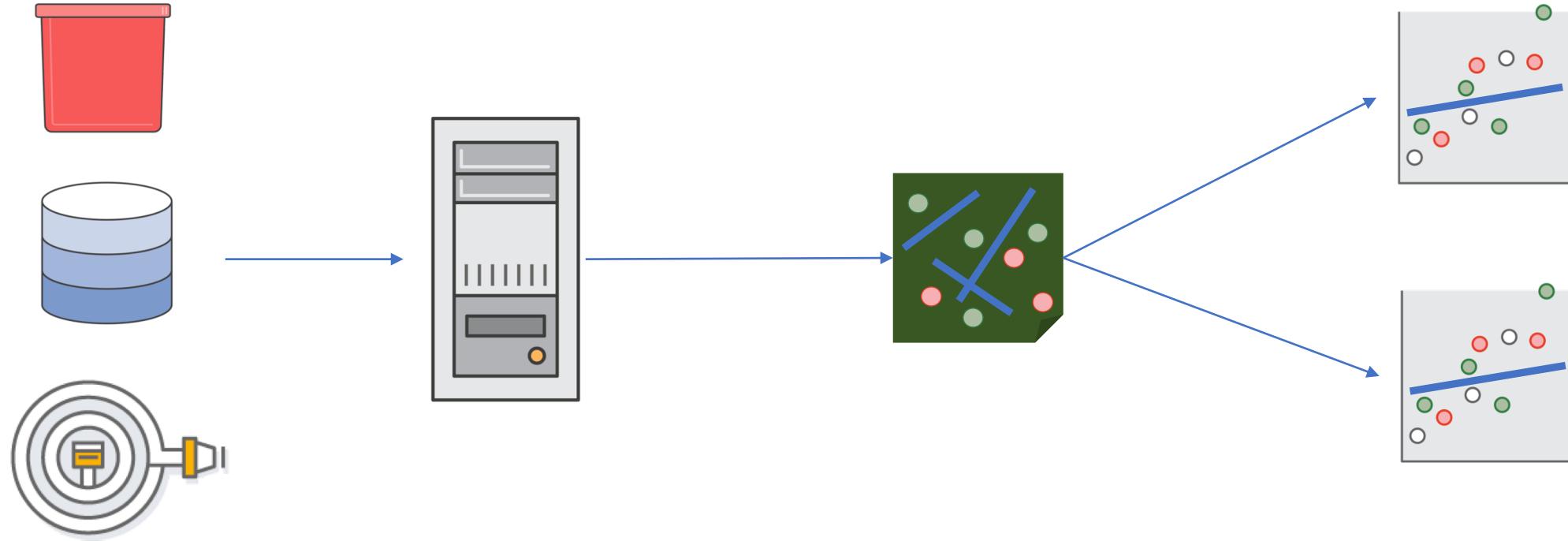
```
return local_state.finalize()
```

Send (compressed) gradient to parameter server

# HPO, State > Model



# HPO, State > Model



Logistic Regression: Tune threshold for Acc, F1, p@r, etc.  
K-Means: state.finalize(k)

SageMaker Algorithms – Accurate,  
Fast, Scalable, and Easy to Use.

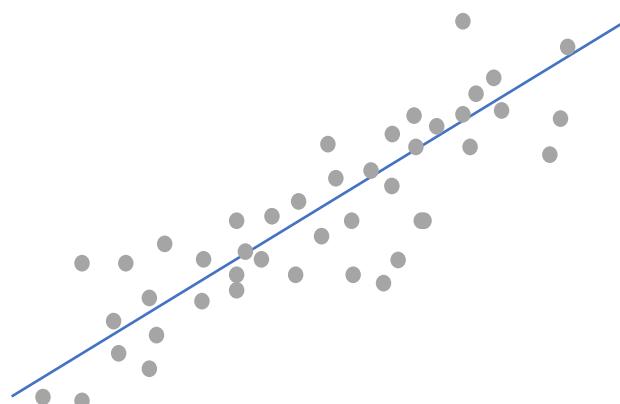


# Algorithms- Example Usage

Algorithm	Function	Example Usage
Linear Learner		
Boosted Decision Trees (XGBoost)	Classification and regression, these are the most popular ML algorithms used today.	<ul style="list-style-type: none"><li>• Estimating click probability for online advisements (for a customer)</li><li>• Directing a customer's inbound phone call to relevant agents</li><li>• Deciding whether a login event is legitimate.</li></ul>
Factorization Machines		
K-Nearest Neighbor		
K-means	Clustering	<ul style="list-style-type: none"><li>• Grouping similar events/document/images together</li></ul>
PCA	Principal Component Analysis	<ul style="list-style-type: none"><li>• Reduce Dimensionality of data</li><li>• Explore main factors/trends in data</li><li>• Visualization</li></ul>
Neural Topic Modelling		
Spectral LDA	Topic Modeling	<ul style="list-style-type: none"><li>• Maps documents into distribution over topics</li><li>• Discover dominant topics in your text corpus</li></ul>
Blazing Text	Word Embedding	<ul style="list-style-type: none"><li>• Feature Engineering for text</li></ul>
DeepAR	Time-series Forecasting	<ul style="list-style-type: none"><li>• Predict the number of page views you'll get in an hour (and the number of servers you'll need to host them!)</li></ul>
Image Classification	Classification of Images	<ul style="list-style-type: none"><li>• Detect quality assurance issues in manufactured goods using images.</li></ul>
Sequence to Sequence	Learn mapping between pairs of sequences	<ul style="list-style-type: none"><li>• Translating text between different languages.</li></ul>

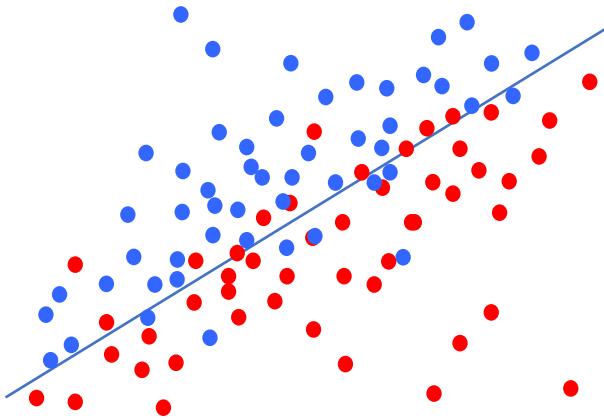
# Linear Learner

Regression:  
Estimate a real  
valued function



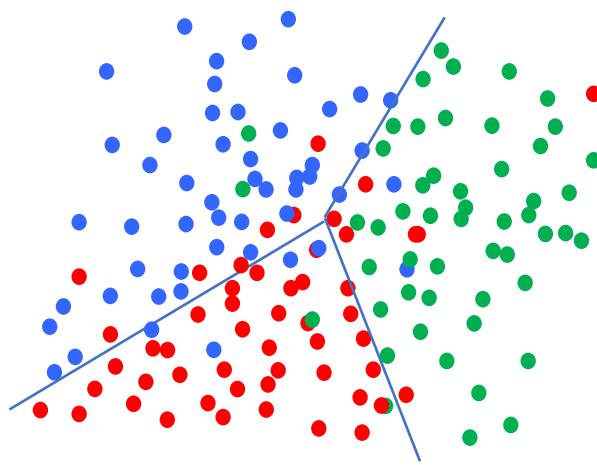
$$\hat{y} = \langle x, w \rangle + b$$

Binary Classification:  
Predict a 0/1 class



$$\hat{y} = \langle x, w \rangle > t ? 0 : 1$$

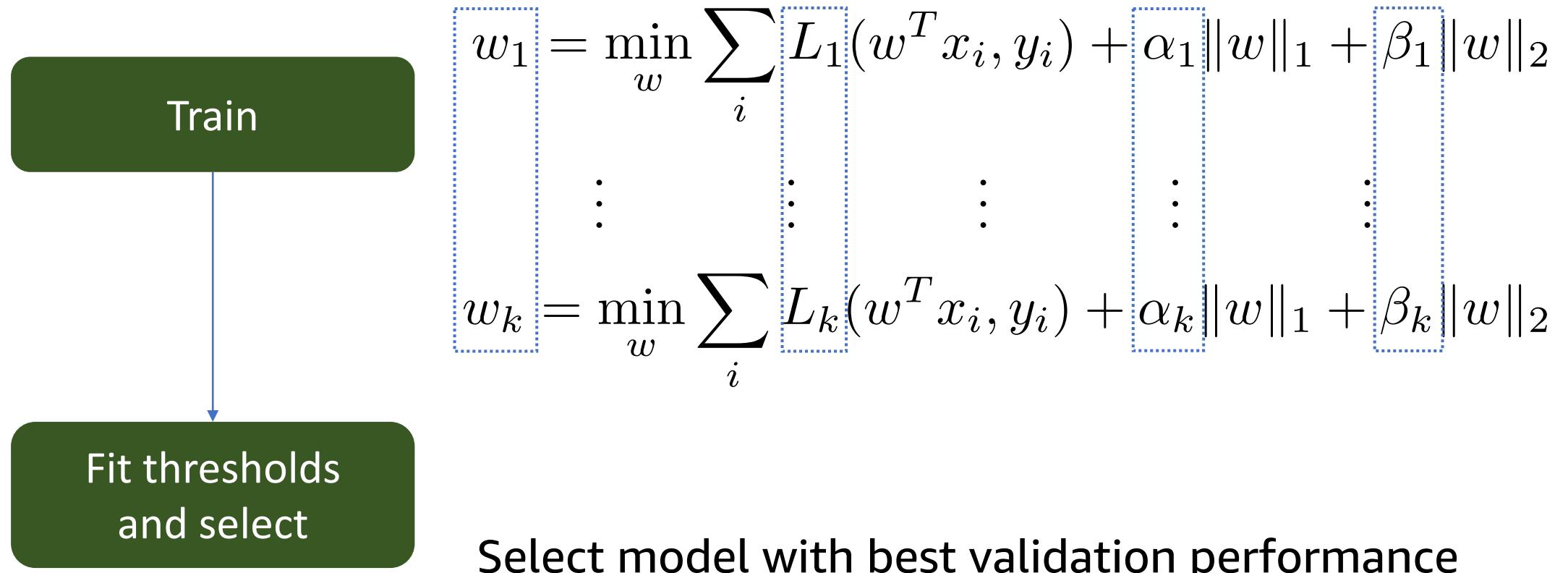
Multiclass Classification:  
Predict one of k classes



$$\hat{y} = \operatorname{argmax}_i \langle x, w_i \rangle$$

# Linear Learner

>8x speedup over naïve parallel training!

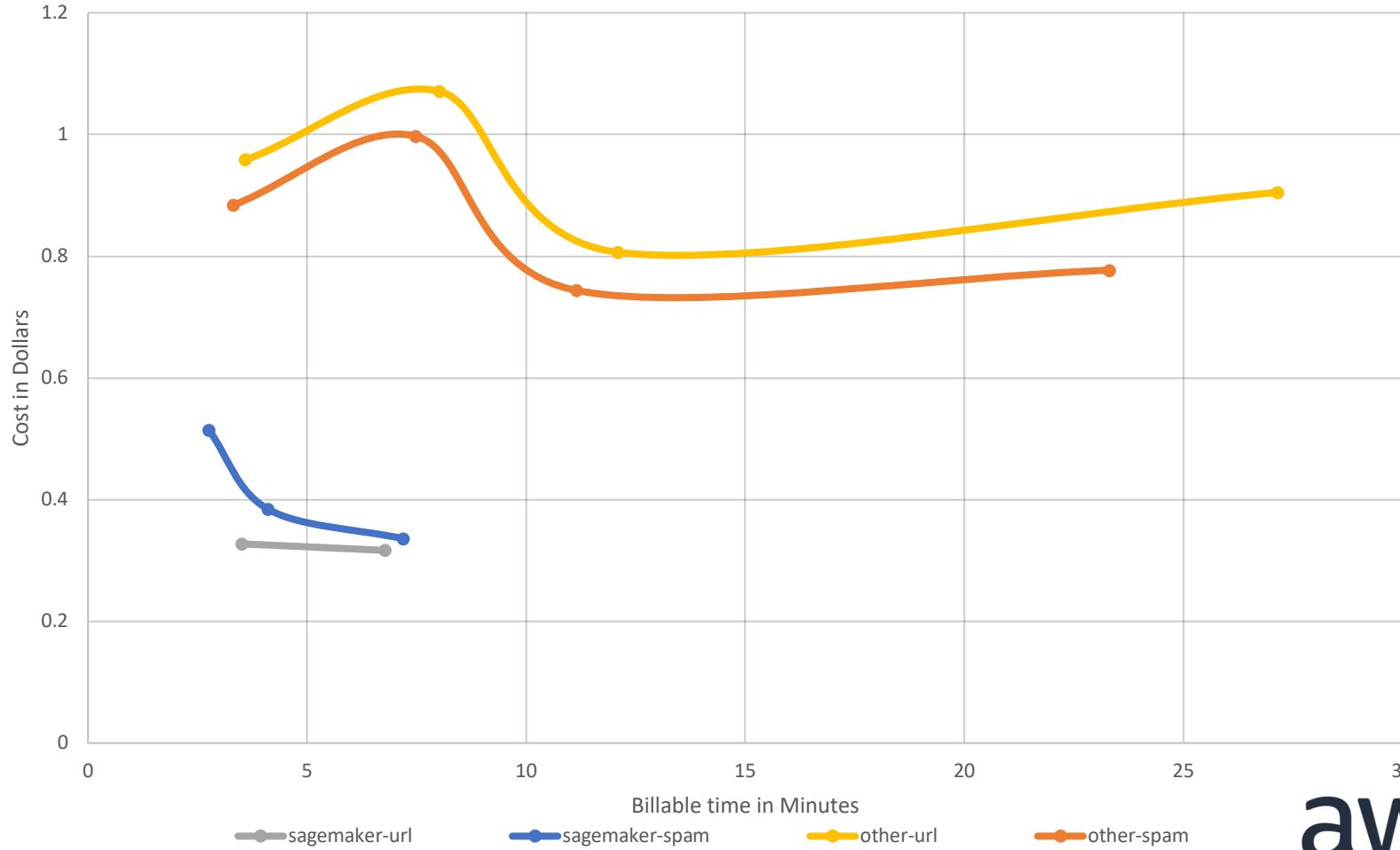


# Linear Learner

Regression (mean squared error)	
SageMaker	Other
1.02	1.06
1.09	1.02
0.332	0.183
0.086	0.129
83.3	84.5

Classification (F1 Score)	
SageMaker	Other
0.980	0.981
0.870	0.930
0.997	0.997
0.978	0.964
0.914	0.859
0.470	0.472
0.903	0.908
0.508	0.508

30 GB datasets for web-spam and web-url classification

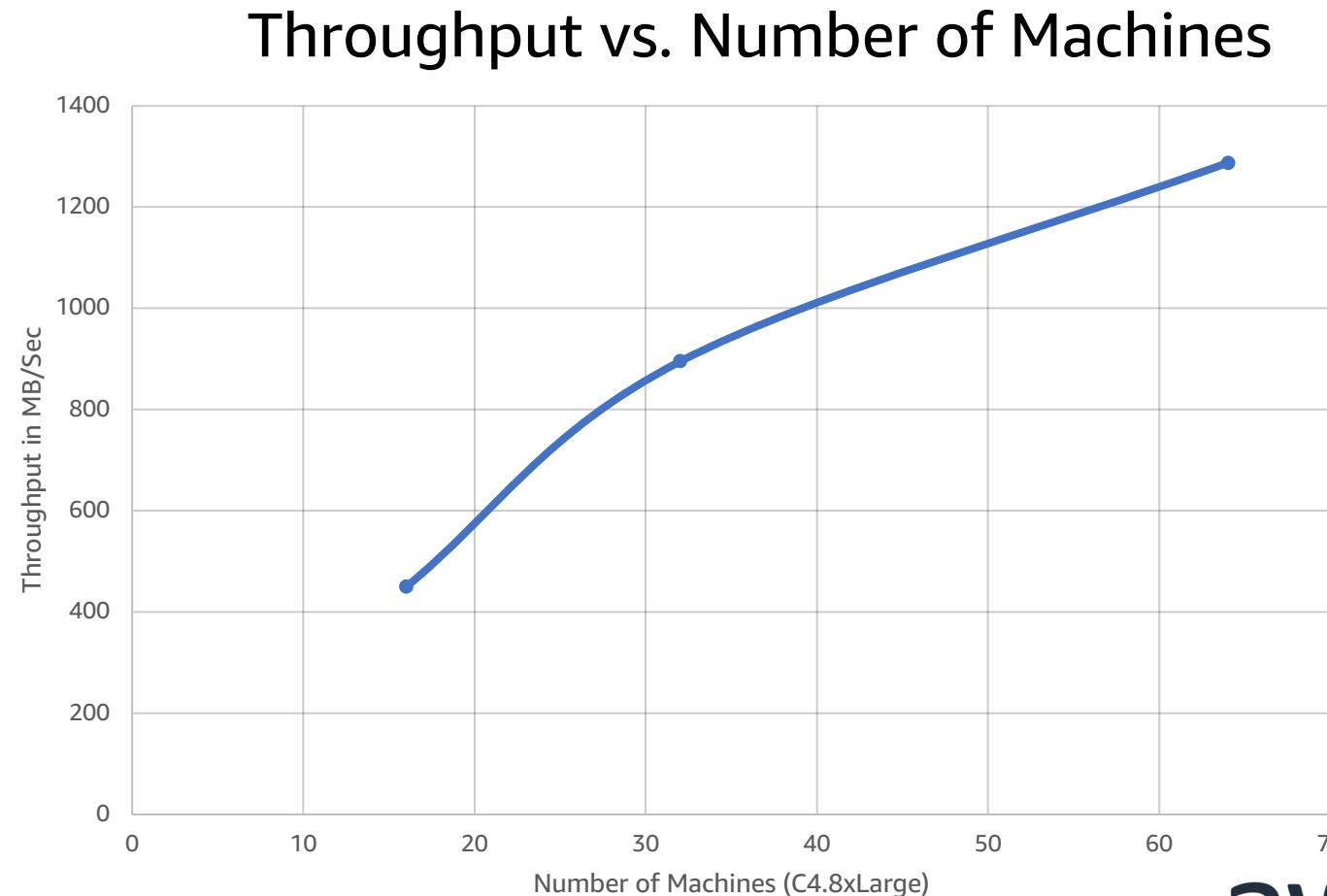


# Boosted Decision Trees

XGBoost is one of the most commonly used implementations of boosted decision trees in the world.

It is now available in Amazon SageMaker!

*dmlc*  
**XGBoost**

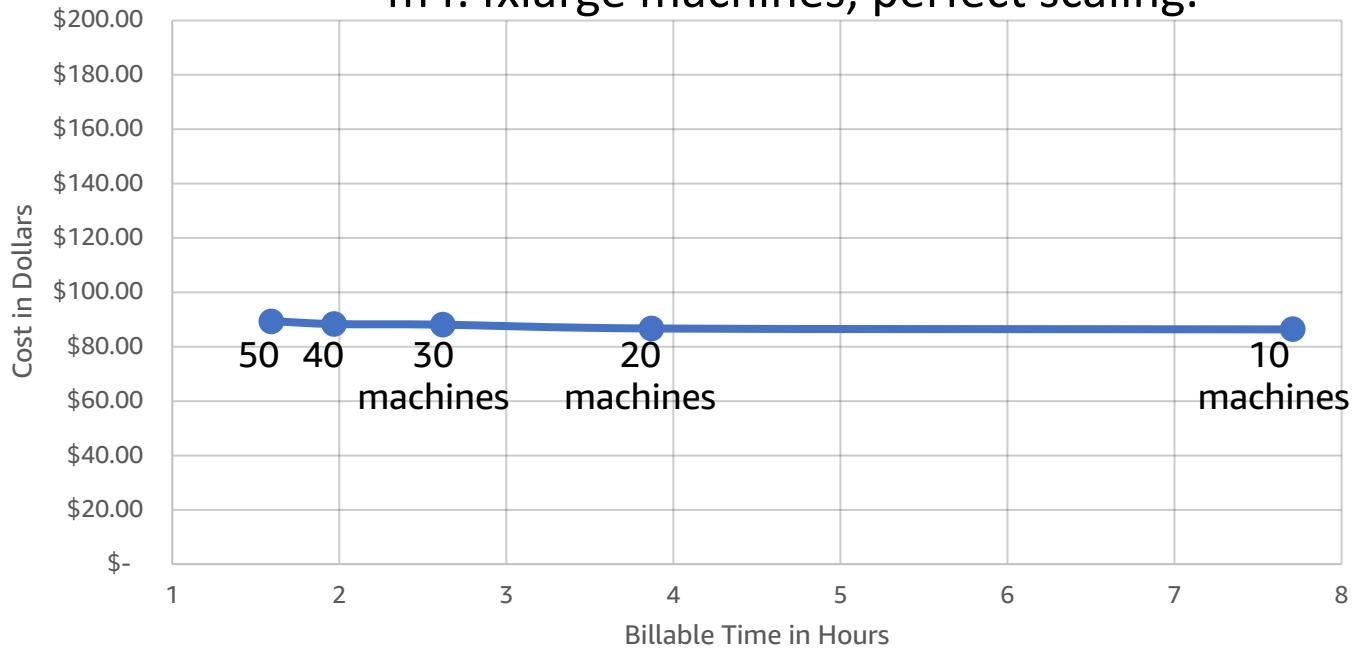


# Factorization Machines

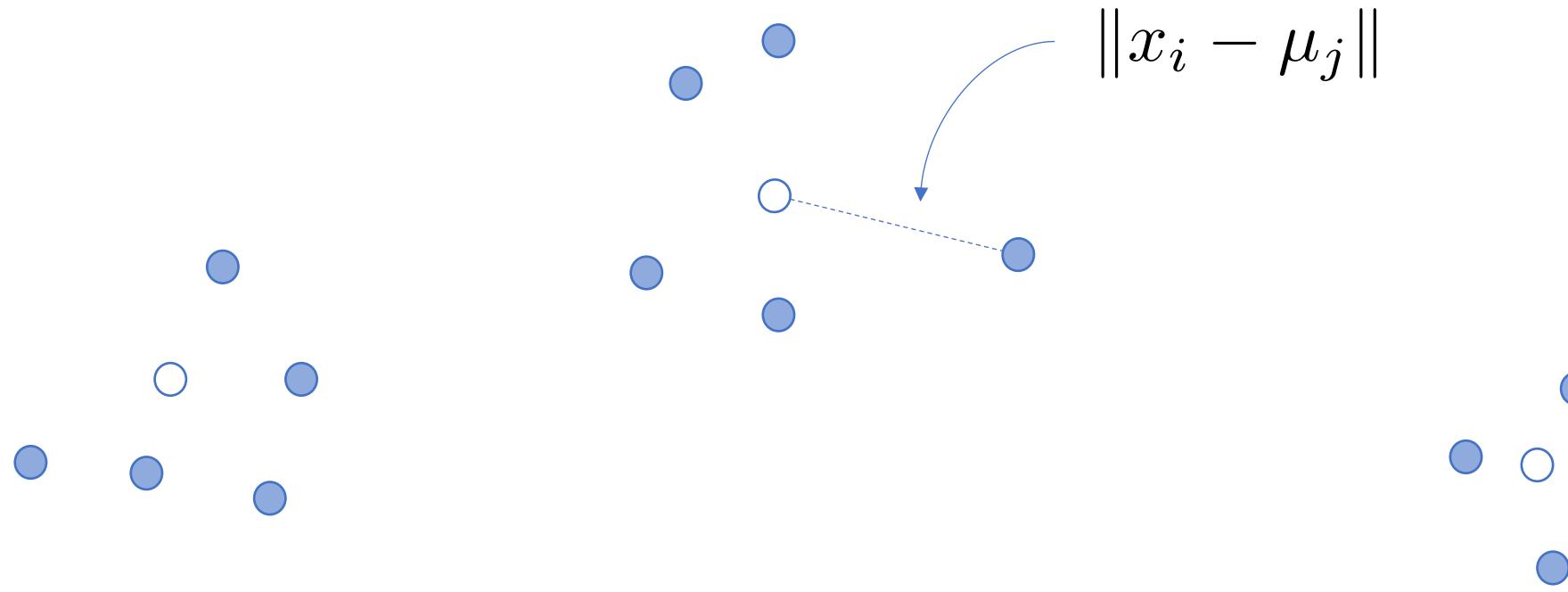
$$\hat{y} = w_0 + \langle w_1, x \rangle + \sum_{i,j>i} x_i x_j \langle v_i, v_j \rangle$$

	Log_loss	F1 Score	Seconds
SageMaker (single pass)	0.494	0.277	820
Other (10 Iter)	0.516	0.190	650
Other (20 Iter)	0.507	0.254	1300
Other (50 Iter)	0.481	0.313	3250

Click Prediction 1 TB advertising dataset,  
m4.4xlarge machines, perfect scaling.



# K-Means Clustering



# K-Means Clustering

Method	Accurate?	Passes	Efficient Tuning	Comments
Lloyds [1]	Yes*	5-10	No	
K-Means ++ [2]	Yes	k+5 to k+10	No	scikit-learn
K-Means   [3]	Yes	7-12	No	spark.ml
Online [4]	No	1	No	
Streaming [5,6]	No	1	No	Impractical
Webscale [7]	No	1	No	spark streaming
Coresets [8]	No	1	Yes	Impractical
<b>SageMaker</b>	<b>Yes</b>	<b>1</b>	<b>Yes</b>	

[1] Lloyd, IEEE TIT, 1982

[2] Arthur et. al. ACM-SIAM, 2007

[3] Bahmani et. al., VLDB, 2012

[4] Liberty et. al., 2015

[5] Shindler et. al, NIPS, 2011

[6] Guha et. al, IEEE Trans. Knowl. Data Eng. 2003

[7] Sculley, WWW, 2010

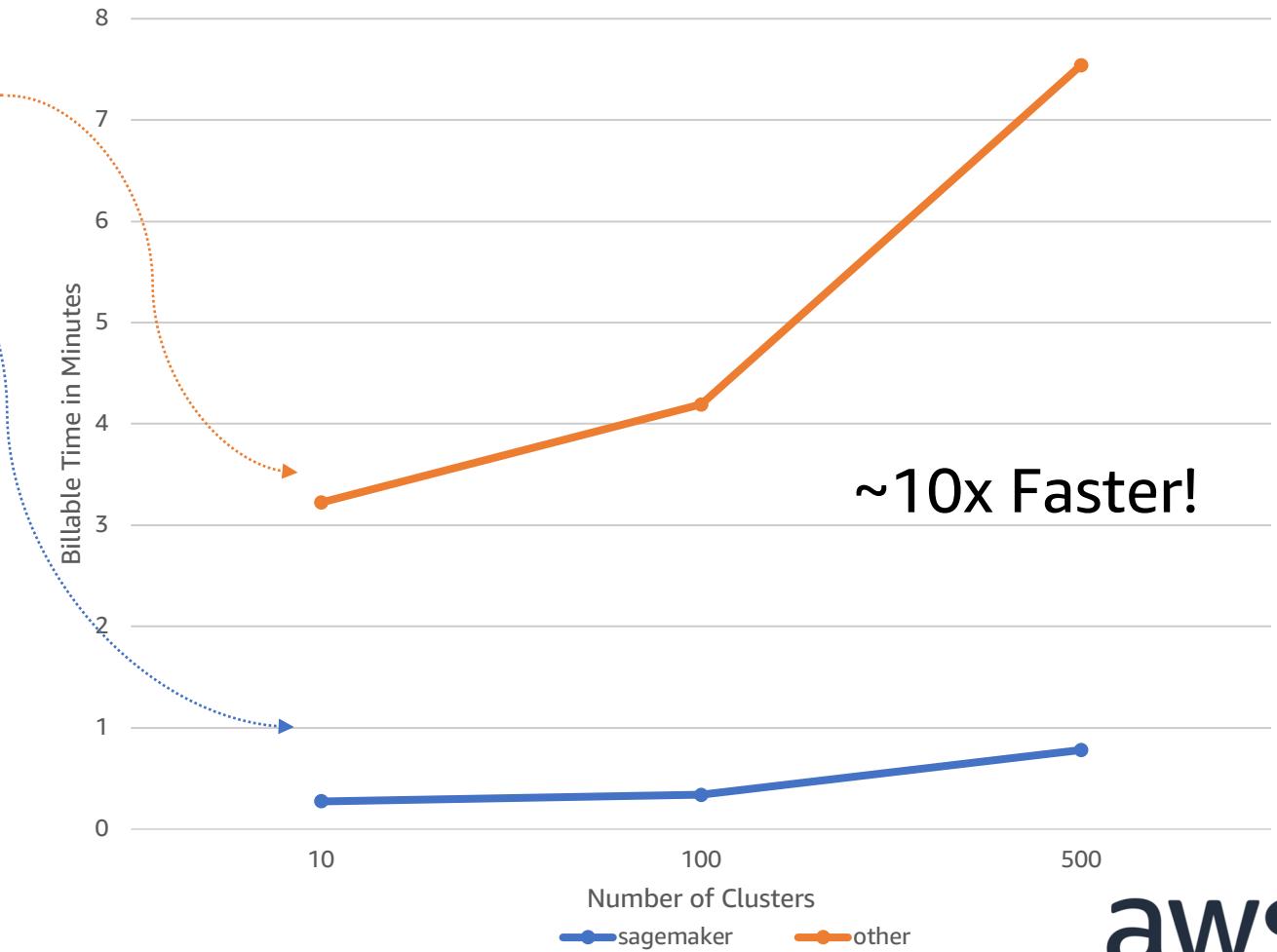
[8] Feldman et. al.



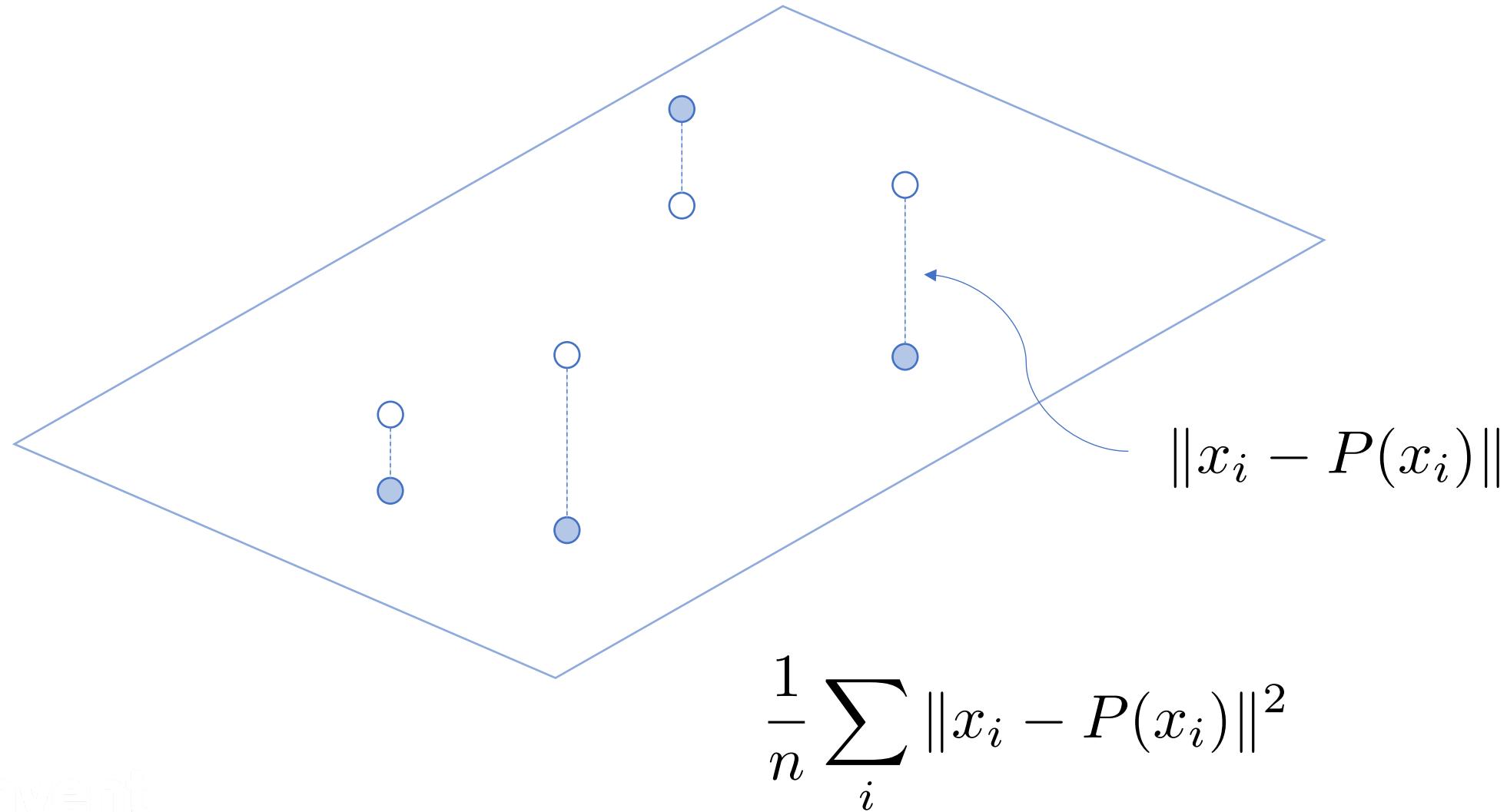
# K-Means Clustering

	k	SageMaker	Other
Text 1.2GB	10	1.18E3	1.18E3
	100	1.00E3	9.77E2
	500	9.18.E2	9.03E2
Images 9GB	10	3.29E2	3.28E2
	100	2.72E2	2.71E2
	500	2.17E2	Failed
Videos 27GB	10	2.19E2	2.18E2
	100	2.03E2	2.02E2
	500	1.86E2	1.85E2
Advertising 127GB	10	1.72E7	Failed
	100	1.30E7	Failed
	500	1.03E7	Failed
Synthetic 1100GB	10	3.81E7	Failed
	100	3.51E7	Failed
	500	2.81E7	Failed

## Running Time vs. Number of Clusters

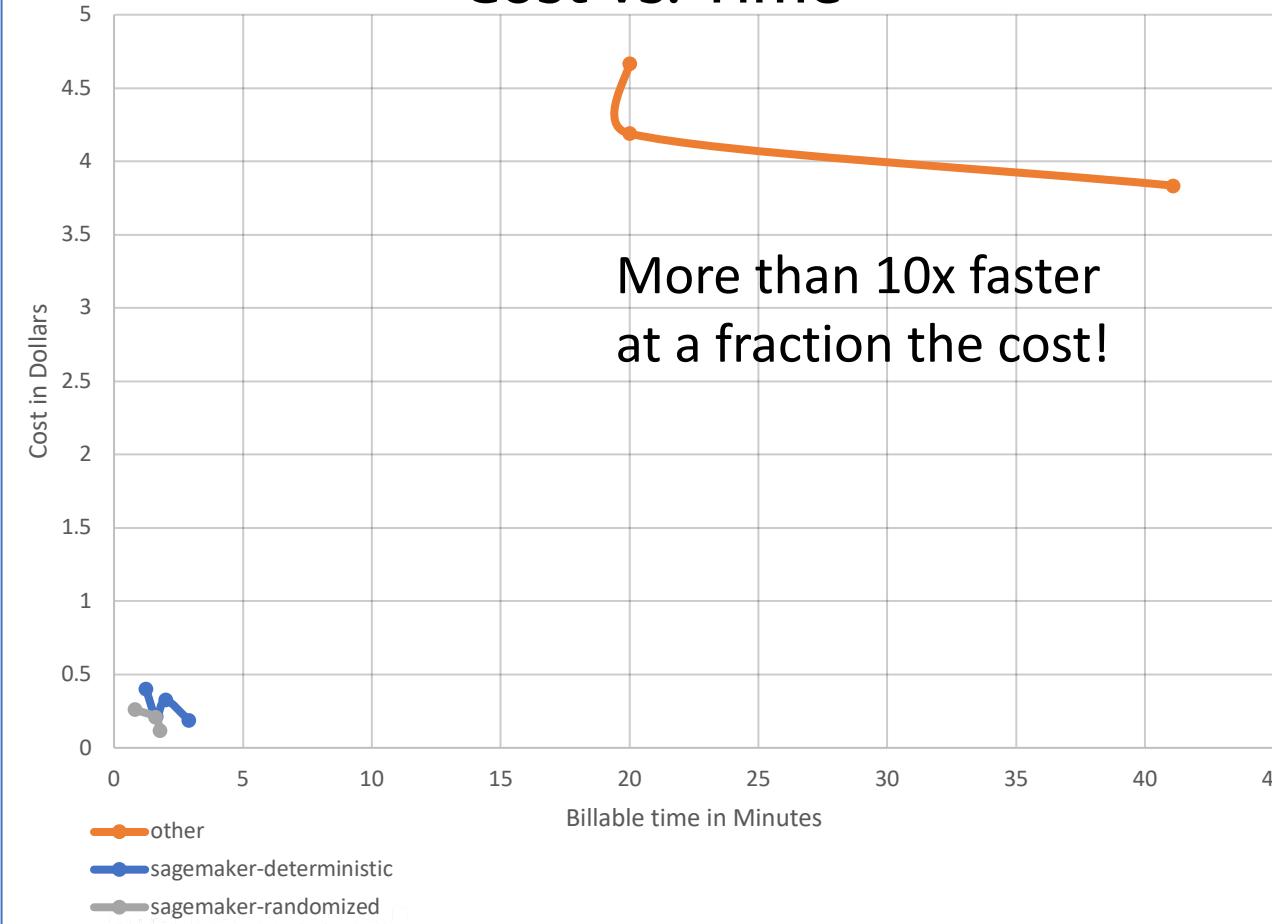


# Principal Component Analysis (PCA)

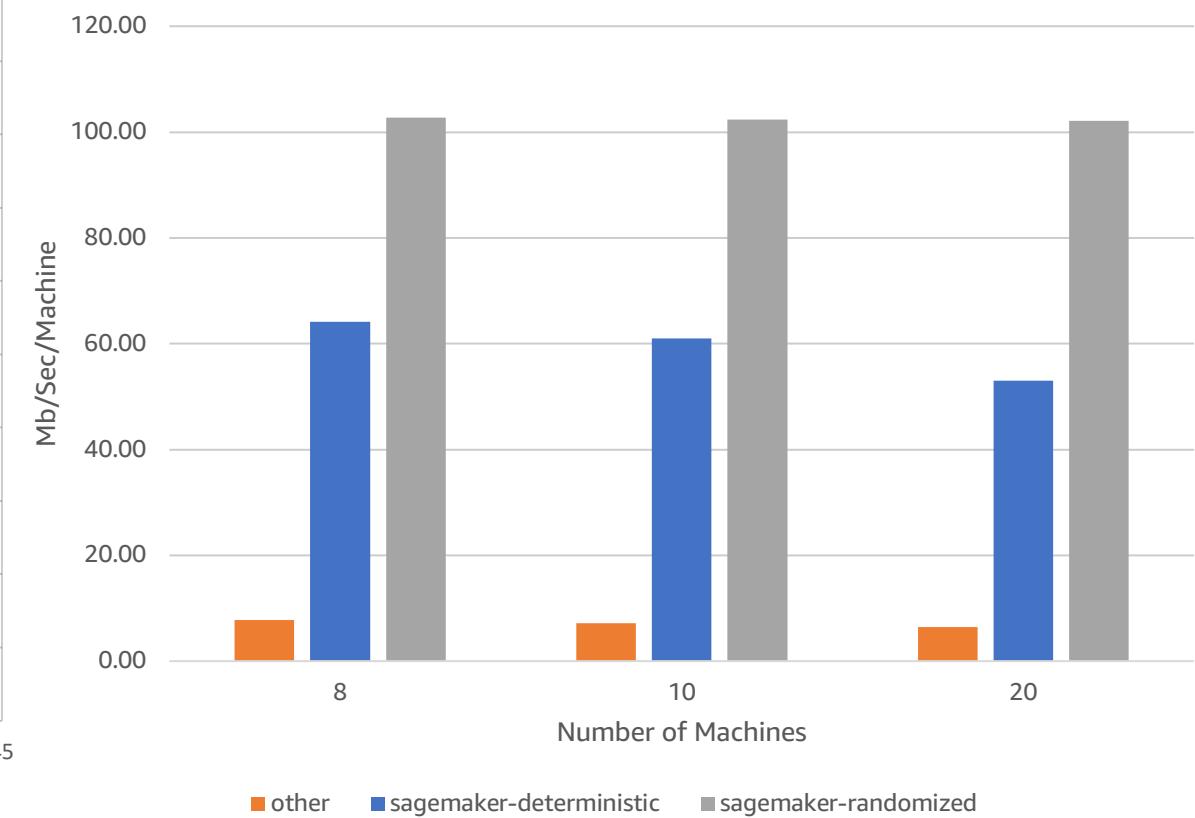


# Principal Component Analysis (PCA)

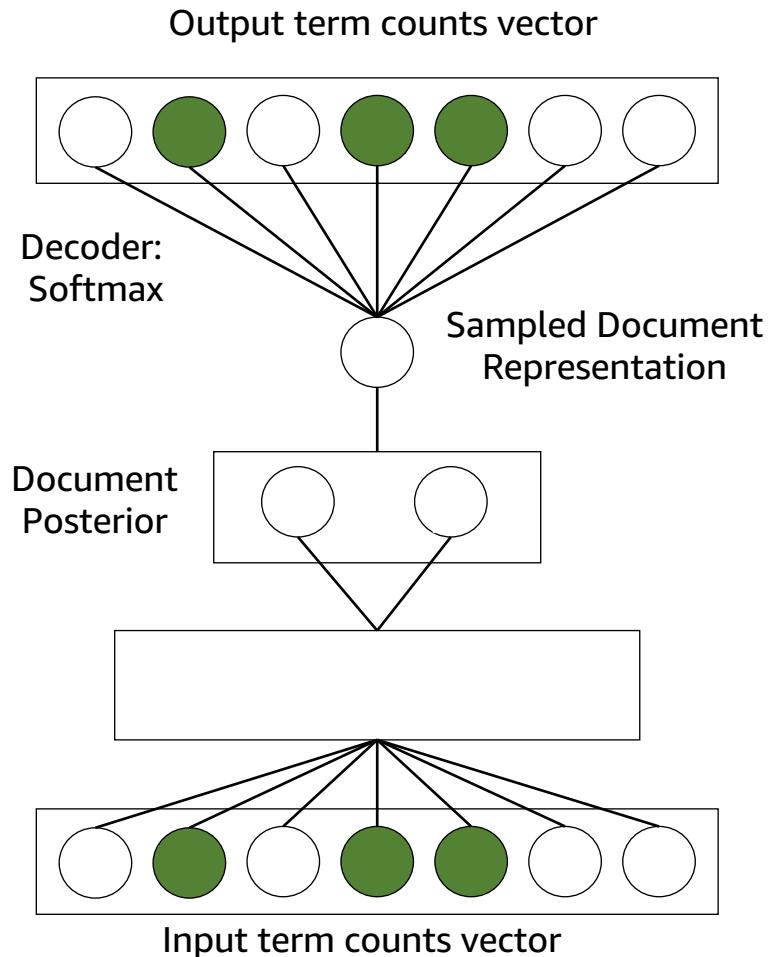
## Cost vs. Time



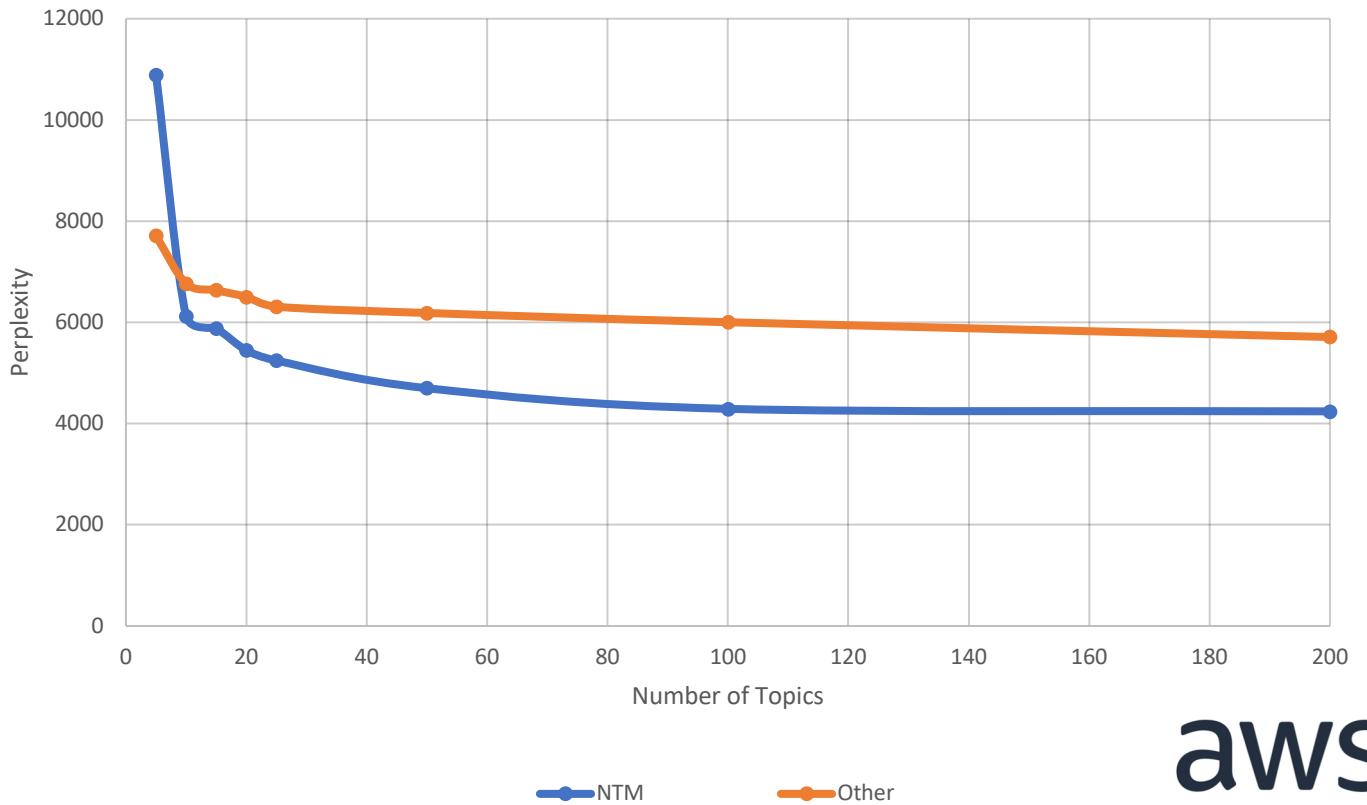
## Throughput and Scalability



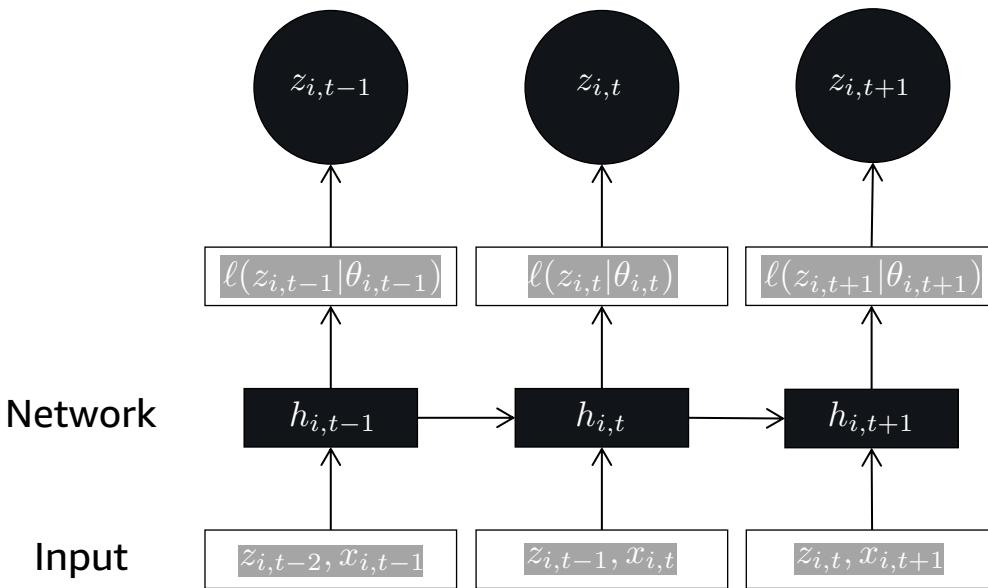
# Neural Topic Modeling



- Perplexity vs. Number of Topic
- (~200K documents, ~100K vocabulary)



# Time Series Forecasting



		Mean absolute percentage error		P90 Loss	
		DeepAR	R	DeepAR	R
traffic	Hourly occupancy rate of 963 bay area freeways	<b>0.14</b>	0.27	<b>0.13</b>	0.24
electricity	Electricity use of 370 homes over time	<b>0.07</b>	0.11	<b>0.08</b>	0.09
pageviews	10k Page view hits of websites	<b>0.32</b>	<b>0.32</b>	0.44	<b>0.31</b>
	180k	<b>0.32</b>	0.34	<b>0.29</b>	NA

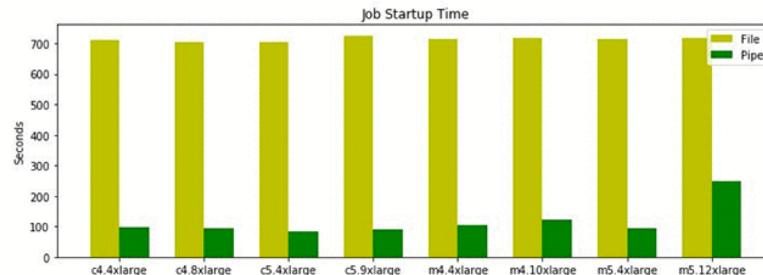
One hour on p2.xlarge, \$1

# Pipe Mode (launched May 23<sup>rd</sup>)

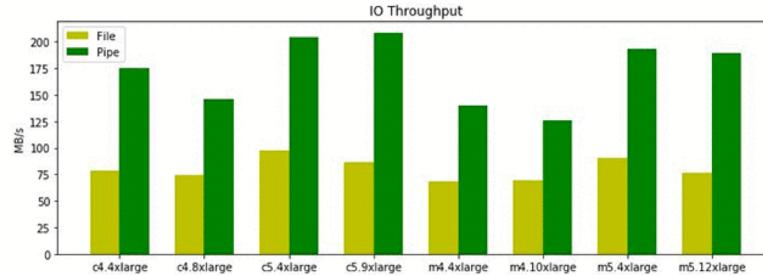
Job Execution Time



Job Startup Time



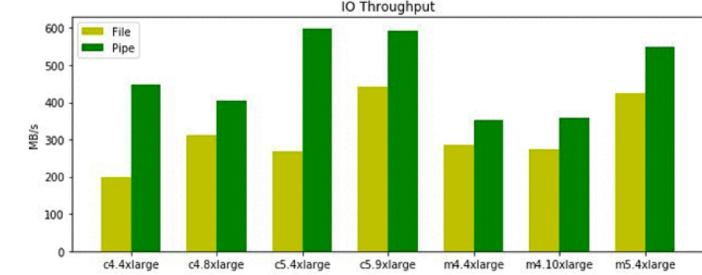
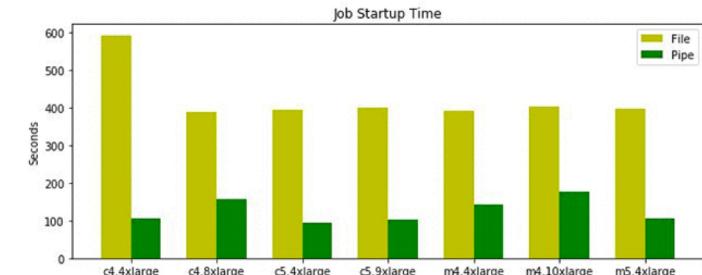
Throughput



PCA



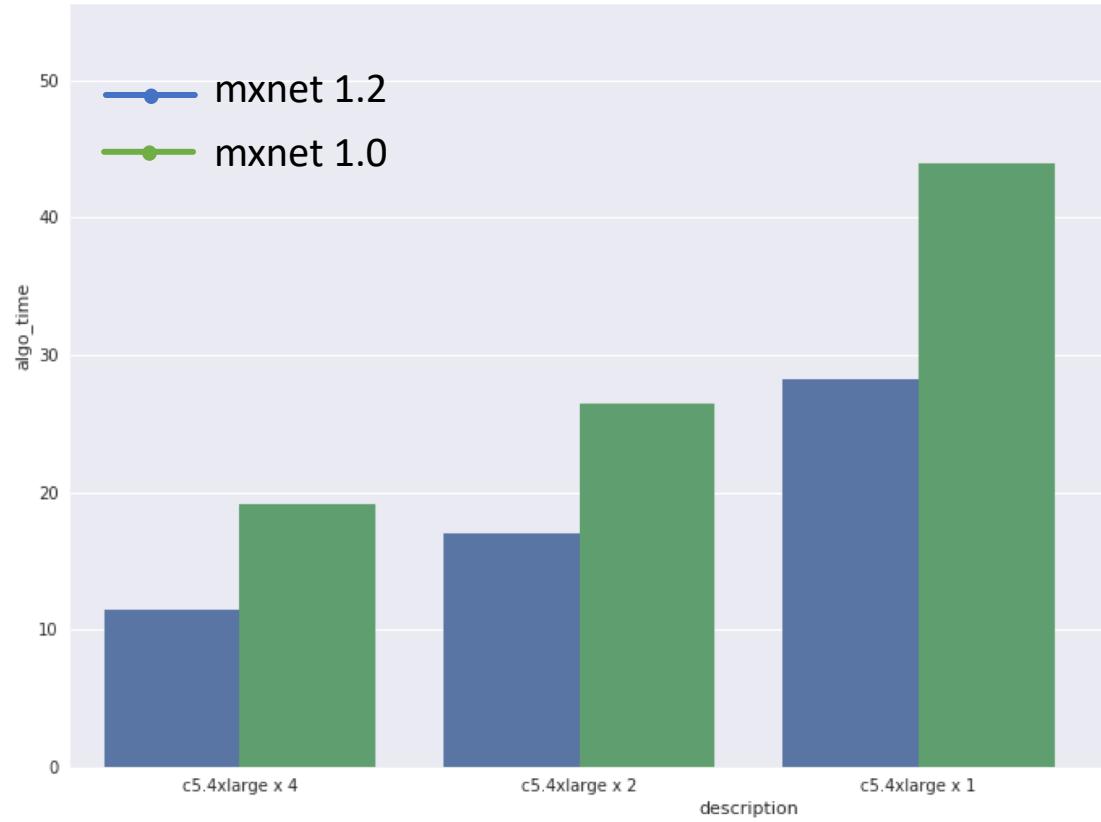
© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved.



K-Means



# Mxnet as an engine



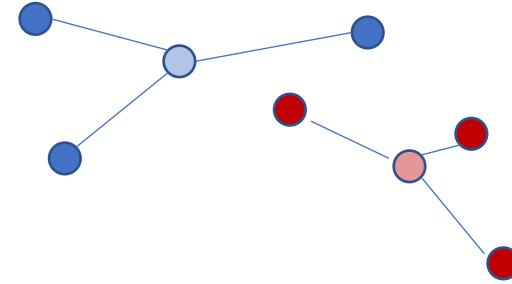
# Deep Dive: SageMaker K-Means



# Lloyds Algorithms

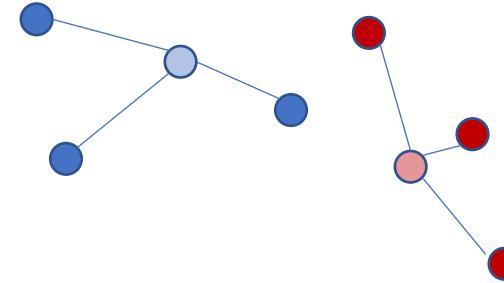
$$\min_{c,A} \sum_{i=1}^n |x_i - c_{A(i)}|^2$$

Fix A, optimize c



$$c_j = E_{x \in A^{-1}(j)}[x]$$

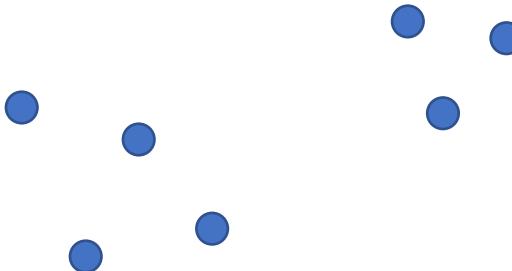
Fix c, optimize A



$$A(i) = \min_j |x_i - c_j|^2$$

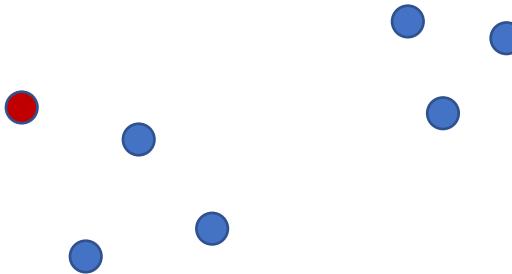
# K-Means++

$$\Pr[c_{next} = x] \propto \min_j |x - c_j|^2$$



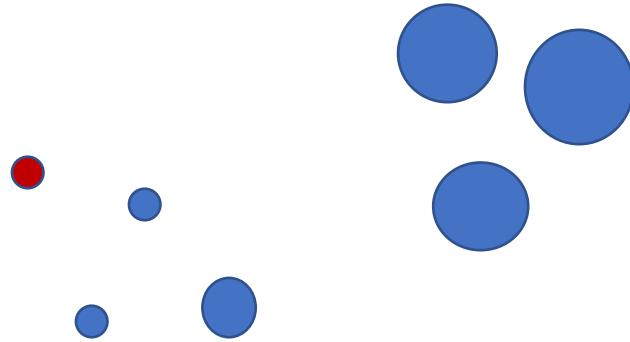
# K-Means++

$$\Pr[c_{next} = x] \propto \min_j |x - c_j|^2$$



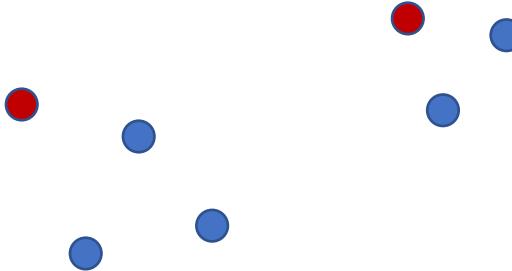
# K-Means++

$$\Pr[c_{next} = x] \propto \min_j |x - c_j|^2$$



# K-Means++

$$\Pr[c_{next} = x] \propto \min_j |x - c_j|^2$$

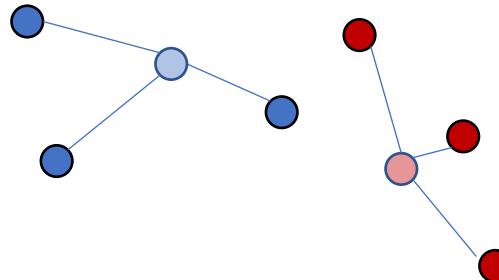


# Online K-Means

For batch in data:

- Assign points in batch to clusters

- Move centers to new average

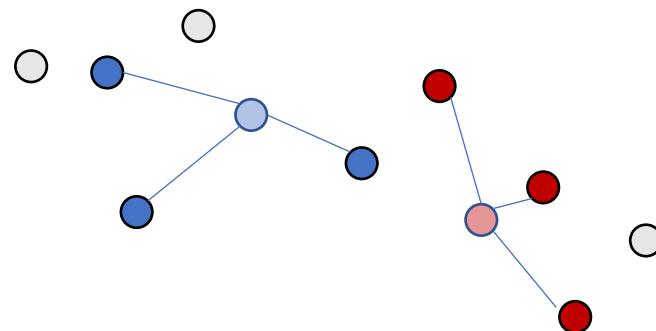


# Online K-Means

For batch in data:

- Assign points in batch to clusters

- Move centers to new average

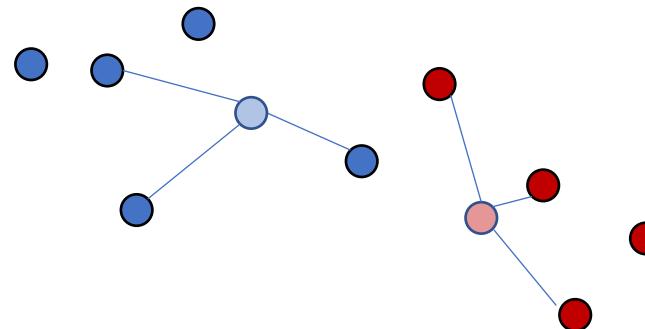


# Online K-Means

For batch in data:

- Assign points in batch to clusters

- Move centers to new average

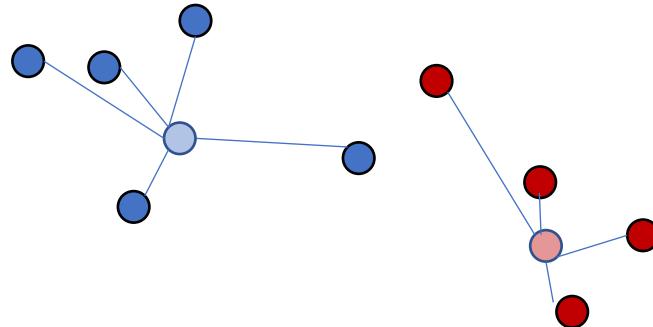


# Online K-Means

For batch in data:

- Assign points in batch to clusters

- Move centers to new average



# Online K-Means

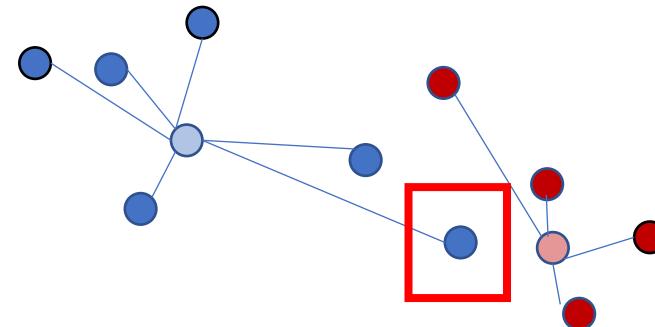
For batch in data:

    Assign points in batch to clusters

    Move centers to new average

## Cons

- Sub-optimal Convergence



# Online K-Means

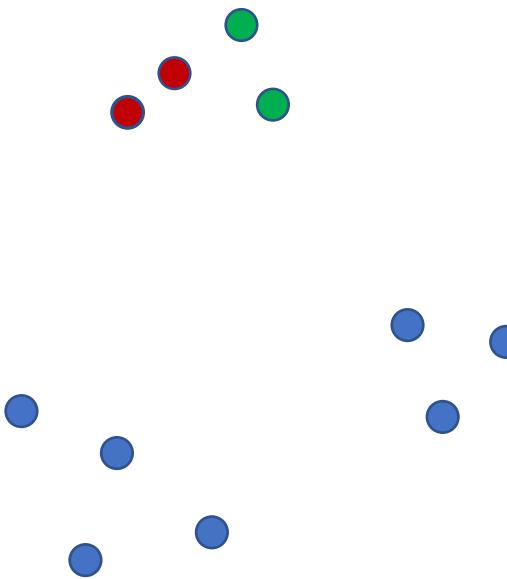
For batch in data:

    Assign points in batch to clusters

    Move centers to new average

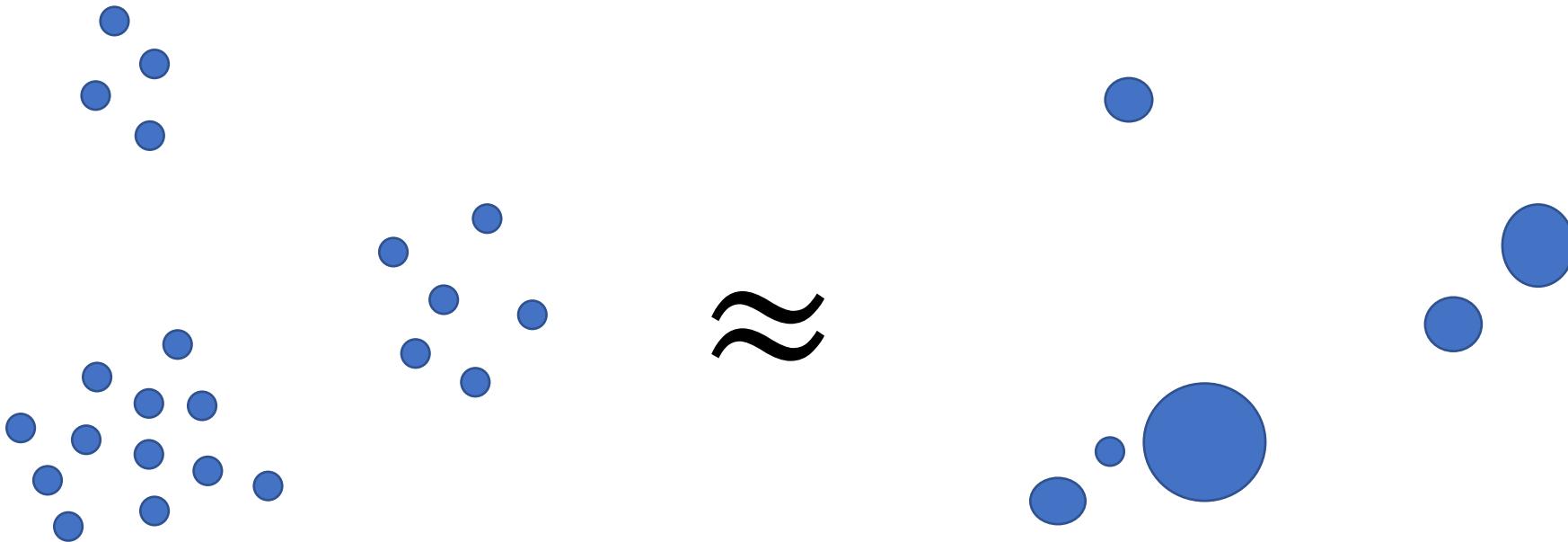
Cons

- Sub-optimal Convergence
- Sensitive to initialization



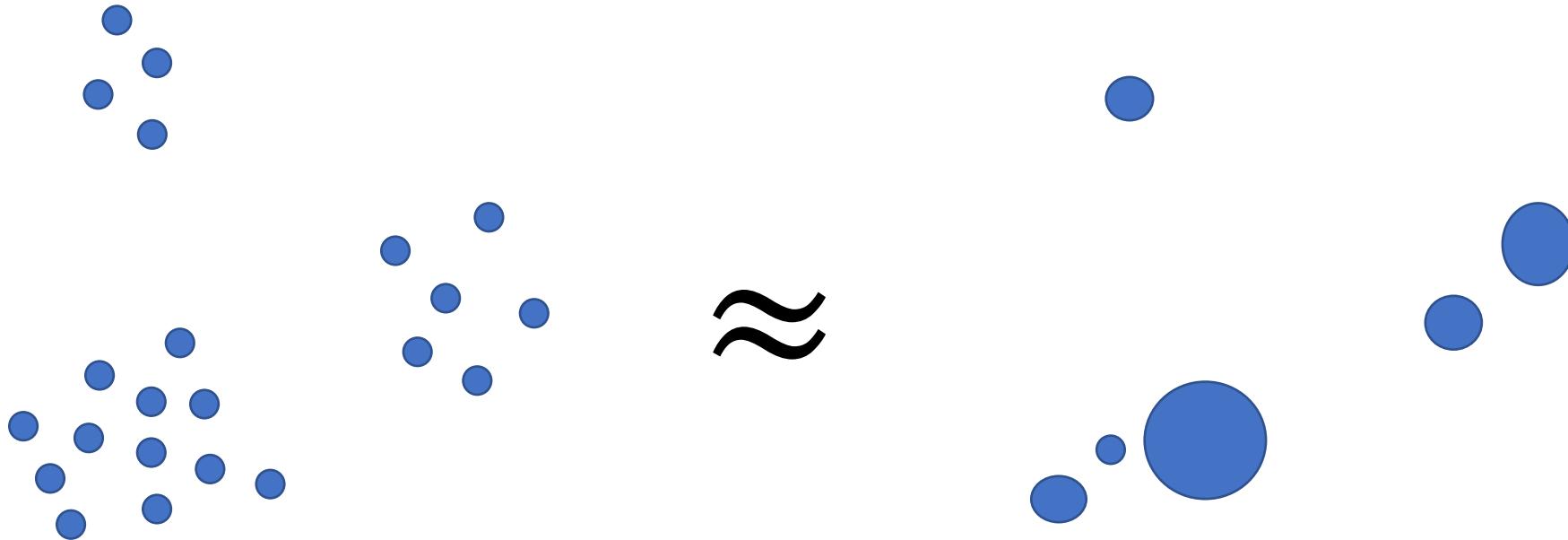
# Core Set

Represent  $n$  points as  $r \ll n$  weighted points



# Core Set

Represent  $n$  points as  $r \ll n$  weighted points

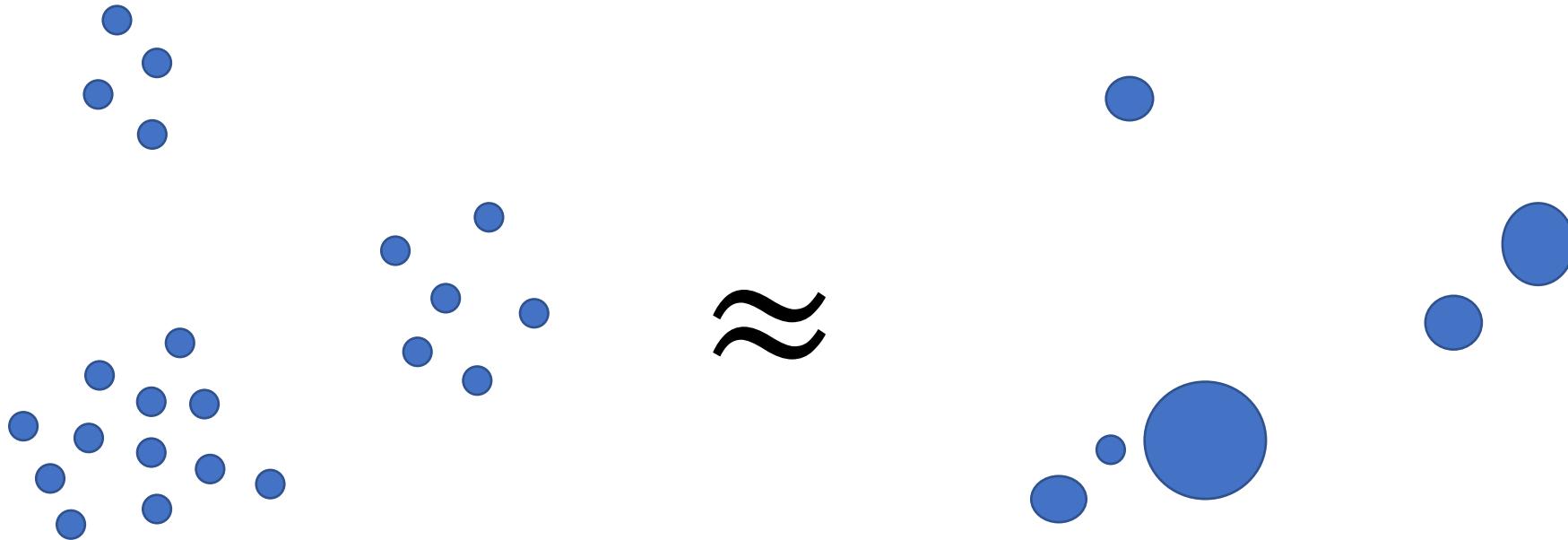


Practical Heuristic

For  $r > k$ , centers of (crude)  $r$ -means  $\approx$  core set for  $k$ -means

# Core Set

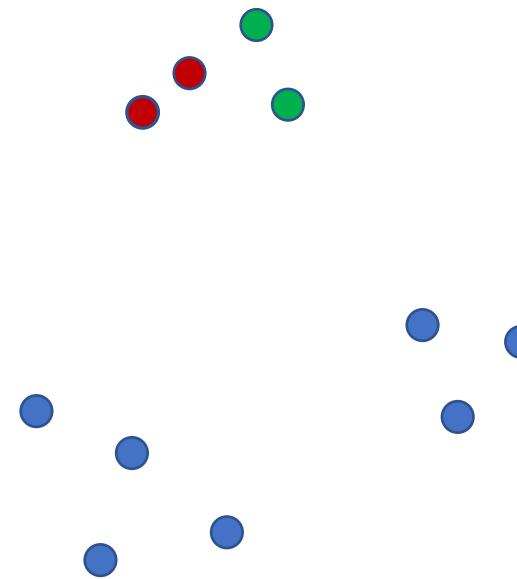
Represent  $n$  points as  $r \ll n$  weighted points



Solves issue #1 in online K-Means  
(suboptimal convergence)

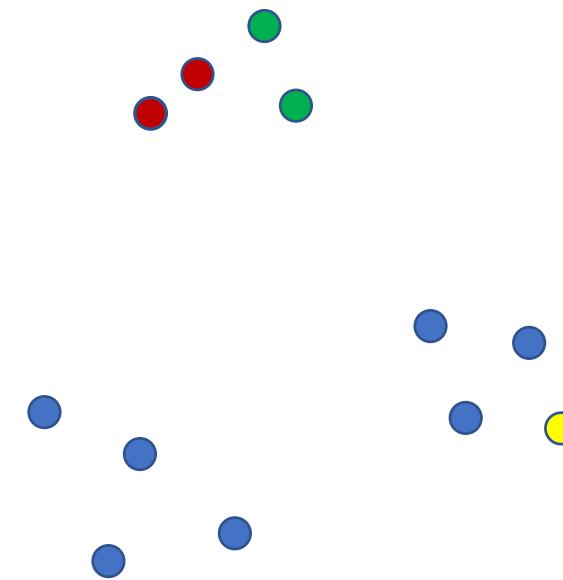
# Adding Centers Over Time

Add more centers over time, based on  $\text{cost} = \min_j |x - c_j|^2$



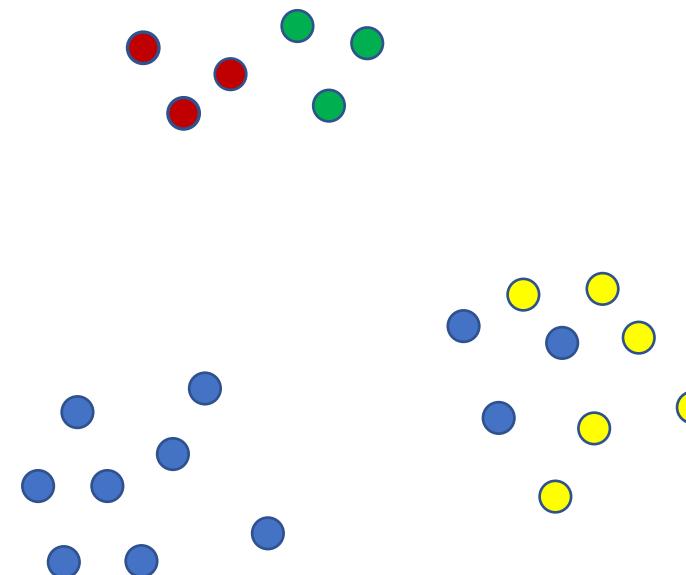
# Adding Centers Over Time

Add more centers over time, based on  $\text{cost} = \min_j |x - c_j|^2$



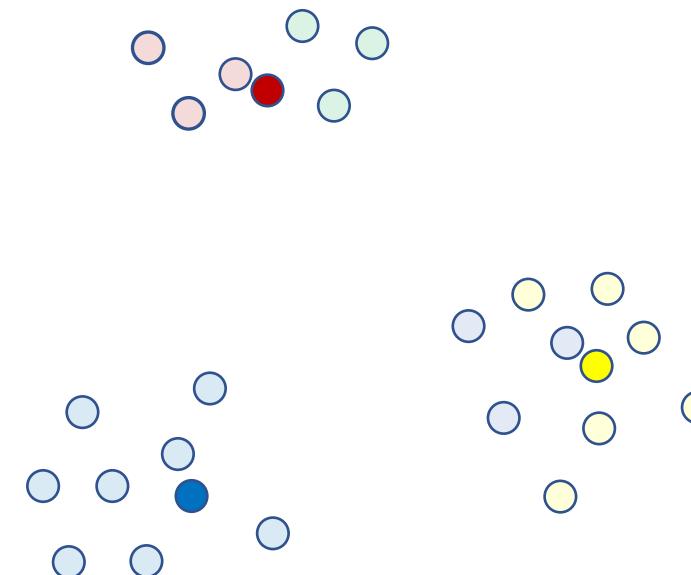
# Adding Centers Over Time

Add more centers over time, based on  $\text{cost} = \min_j |x - c_j|^2$



# Adding Centers Over Time

Add more centers over time, based on  $\text{cost} = \min_j |x - c_j|^2$



# Adding Centers Over Time

Add more centers over time, based on cost =  $\min_j |x - c_j|^2$

Overcomes poor initializations

# Adding Centers Over Time

Add more centers over time, based on cost =  $\min_j |x - c_j|^2$

Overcomes poor initializations

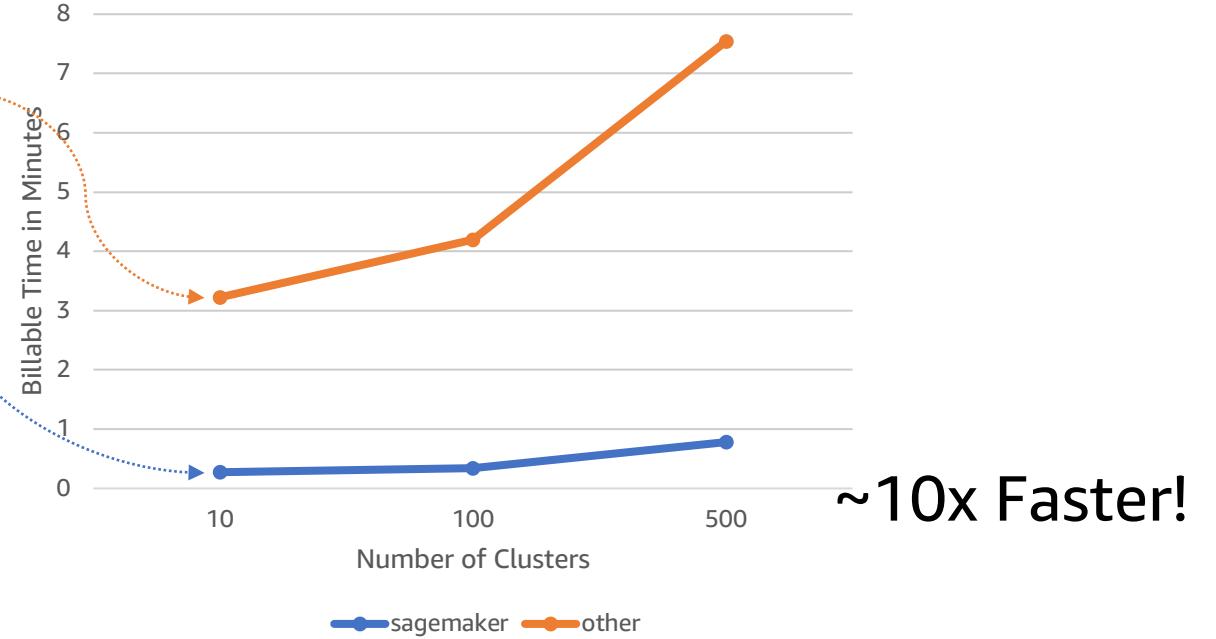
We already have r slots...

Start with r centers, regularly throw away least useful centers

# Experiments

	k	SageMaker	Other
Text 1.2GB	10	1.18E3	1.18E3
	100	1.00E3	9.77E2
	500	9.18.E2	9.03E2
Images 9GB	10	3.29E2	3.28E2
	100	2.72E2	2.71E2
	500	2.17E2	Failed
Videos 27GB	10	2.19E2	2.18E2
	100	2.03E2	2.02E2
	500	1.86E2	1.85E2
Advertising 127GB	10	1.72E7	Failed
	100	1.30E7	Failed
	500	1.03E7	Failed
Synthetic 1100GB	10	3.81E7	Failed
	100	3.51E7	Failed
	500	2.81E7	Failed

## Running Time vs. Number of Clusters



~10x Faster!

Other = Lloyds + kmeans ||  
Requires 7-12 passes over the data

Online k-means achieves significantly worse loss

# Using Amazon SageMaker Algorithms on AWS



# From Amazon SageMaker Notebooks

Hardware

```
import boto3
import sagemaker

sess = sagemaker.Session()

pca = sagemaker.estimator.Estimator(containers[boto3.Session().region_name],
                                      role,
                                      train_instance_count=1,
                                      train_instance_type='ml.c4.xlarge',
                                      output_path=output_location,
                                      sagemaker_session=sess)

pca.set_hyperparameters(feature_dim=50000,
                        num_components=10,
                        subtract_mean=True,
                        algorithm_mode='randomized',
                        mini_batch_size=200)

pca.fit({'train': s3_train_data})
```

Start Training

Host model

```
pca_predictor = pca.deploy(initial_instance_count=1,
                            instance_type='ml.c4.xlarge')
```

# From Command Line

Algorithm



Input Data



Hardware



```
profile=<your_profile>
arn_role=<your_arn_role>
training_image=382416733822.dkr.ecr.us-east-1.amazonaws.com/kmeans:1
training_job_name=cluterizing_text_documents_`date '+%Y_%m_%d_%H_%M_%S'` \
aws --profile $profile \
--region us-east-1 \
sagemaker create-training-job \
--training-job-name $training_job_name \
--algorithm-specification TrainingImage=$training_image,TrainingInputMode=File \
--hyper-parameters k=10,feature_dim=1024,mini_batch_size=1000 \
--role-arn $arn_role \
--input-data-config '{"ChannelName": "train", "DataSource": {"S3DataSource": {"S3DataType": "S3Prefix", "S3Uri": "s3://kmeans_demo/train", "S3DataDistributionType": "ShardedByS3Key"}, "CompressionType": "None", "RecordWrapperType": "None"}' \
--output-data-config S3OutputPath=s3://training_output/$training_job_name \
--resource-config InstanceCount=2,InstanceType=m1.c4.8xlarge,VolumeSizeInGB=50 \
--stopping-condition MaxRuntimeInSeconds=3600
```

# SageMaker + Spark =

```
# Python/PySpark Example
from sagemaker_pyspark import SageMakerEstimator

features = spark.read.parquet('s3://<bucket>/<dataset>')

algorithm = SageMakerEstimator(
    trainingImage=ntm_container,
    modelImage=ntm_container,
    trainingInstanceType='ml.p3.8xlarge',
    trainingInstanceCount=16,
    endpointInstanceType='ml.c5.2xlarge',
    endpointInitialInstanceCount=4,
    hyperParameters={
        "num_topics": "100",
        "feature_dim": 250000,
        "mini_batch_size": "10000",
    },
    sagemakerRole=IAMRole(role_arn)
)

model = algorithm.fit(features)
```

# SageMaker + Spark =

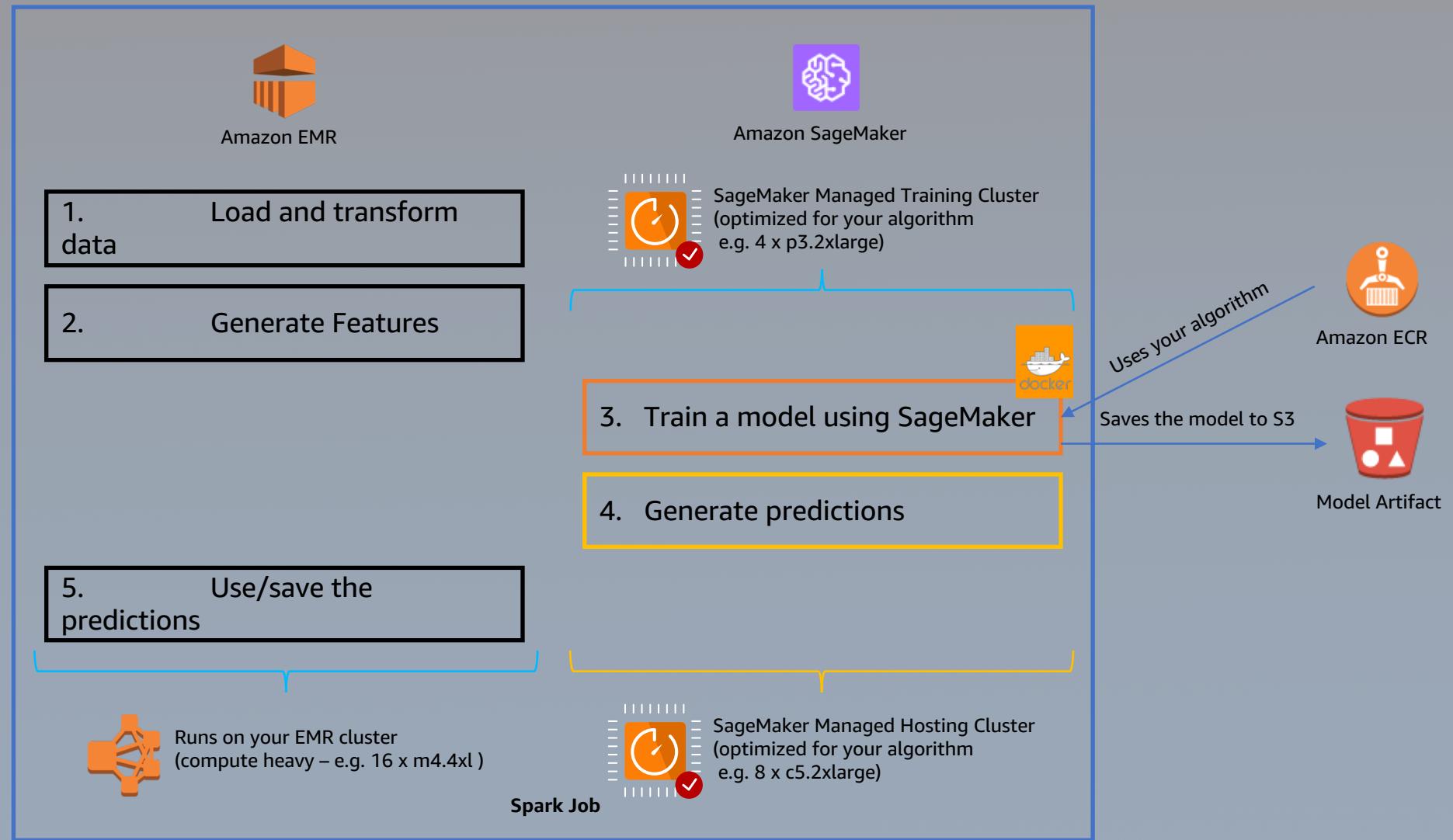
```
// Scala Example
import com.amazonaws.services.sagemaker.sparksdk.{IAMRole, SageMakerEstimator}

val features = spark.read.parquet("s3://<bucket>/<dataset>")

val algorithm = new SageMakerEstimator(
    trainingImage = ntm_container,
    modelImage = ntm_container,
    trainingInstanceType = "ml.p3.8xlarge",
    trainingInstanceCount = 16,
    endpointInstanceType = "ml.c5.2xlarge",
    endpointInitialInstanceCount = 4,
    hyperParameters = Map(
        "num_topics" -> "100",
        "feature_dim" -> "250000",
        "mini_batch_size" -> "10000"
    ),
    sagemakerRole = IAMRole(roleArn)
)

val model = estimator.fit(features)
```

# SageMaker + Spark =



# Amazon SageMaker - Try It Out

AWS Services Resource Groups Admin/libertye-Isengard @ 90... Oregon Support

Amazon SageMaker X

Amazon SageMaker > Dashboard

Overview Hide

The screenshot shows the Amazon SageMaker dashboard. On the left, there's a sidebar with links for Dashboard, Notebook instances, Jobs, Resources, Models, Endpoint configuration, and Endpoints. The main area has a title 'Overview' with four sections: 'Notebook instance', 'Jobs', 'Models', and 'Endpoint'. Each section contains a brief description and a call-to-action button. The 'Notebook instance' section features a 'Create notebook instance' button. The 'Jobs' section features a 'View jobs' button. The 'Models' section features a 'View models' button. The 'Endpoint' section features a 'View endpoints' button.

Dashboard

Notebook instances

Jobs

Resources

Models

Endpoint configuration

Endpoints

Amazon SageMaker > Dashboard

Overview Hide

Notebook instance

Explore AWS data in your notebooks, and use algorithms to create models via training jobs.

Create notebook instance

Jobs

Track training jobs at your desk or remotely. Leverage high-performance AWS algorithms.

View jobs

Models

Create models for hosting from job outputs, or import externally trained models into Amazon SageMaker.

View models

Endpoint

Deploy endpoints for developers to use in production. A/B Test model variants via an endpoint.

View endpoints