

Distributed Machine Learning with Ray

Scaling AI Applications in Python with the Ray
Framework

Learning Outcomes

Learning Outcomes

- Basics of distributed machine learning

Learning Outcomes

- Basics of distributed machine learning
- Use cases that Ray is especially well-suited to solve

Learning Outcomes

- Basics of distributed machine learning
- Use cases that **Ray** is especially well-suited to solve
- Differences between **Ray** vs. other frameworks

Learning Outcomes

- Basics of distributed machine learning
- Use cases that **Ray** is especially well-suited to solve
- Differences between **Ray** vs. other frameworks
- Internals of **Ray** and its execution model

Learning Outcomes

- Basics of distributed machine learning
- Use cases that **Ray** is especially well-suited to solve
- Differences between **Ray** vs. other frameworks
- Internals of **Ray** and its execution model
- Context of the larger **Ray** ecosystem.

Getting the Materials

<https://sr.ht/~hyphaebeast/ray-live-training/>

Do you have any experience with machine learning?

**Have you used a distributed computing framework?
(like Spark, Flink, Hadoop, etc.)**

Why Scale?

A Brief History

In terms of size [of transistors] you can see that we're approaching the size of atoms which is a fundamental barrier, but it'll be two or three generations before we get that far—but that's as far out as we've ever been able to see.

— Gordon Moore (2006)

A Brief History

We have another 10 to 20 years before we reach a fundamental limit. By then they'll be able to make bigger chips and have transistor budgets in the billions.

— Gordon Moore (2006)

Big Data

Big Data

- **Long Data:** Large # of instances/records/rows/etc.

Big Data

- **Long Data:** Large # of instances/records/rows/etc.
- **Wide Data:** High dimensionality

Big Models

Big Models

- **Large Parameter Spaces:** Models with large number of parameters to learn

Big Models

- **Large Parameter Spaces:** Models with large number of parameters to learn
- **Prediction Cascades:** Complex pipeline of multiple prediction tasks

Small Latencies

Small Latencies

- Real-time predictions (i.e. website ad serving)

Small Latencies

- Real-time predictions (i.e. website ad serving)
- Prediction Cascades... again (i.e. self-driving cars)

What Scale?

What Scale?

→ Memory

What Scale?

- Memory
- Computation

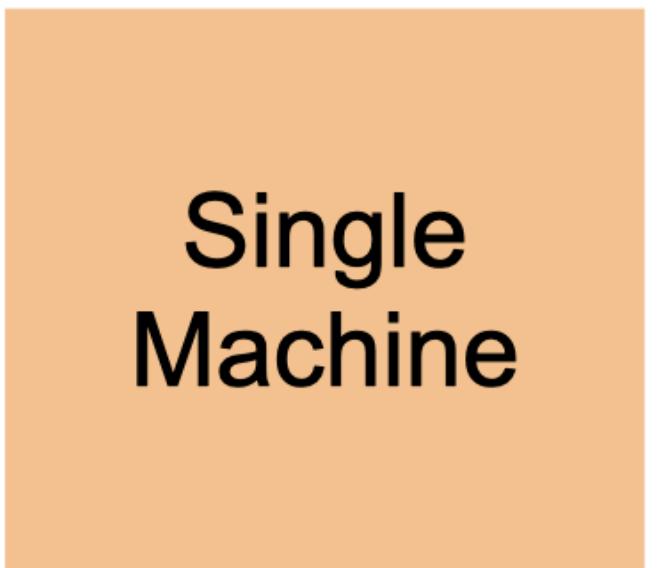
Introduction to distributed systems

Distributed programming is the art of solving the same problem that you can solve on a single computer using multiple computers.

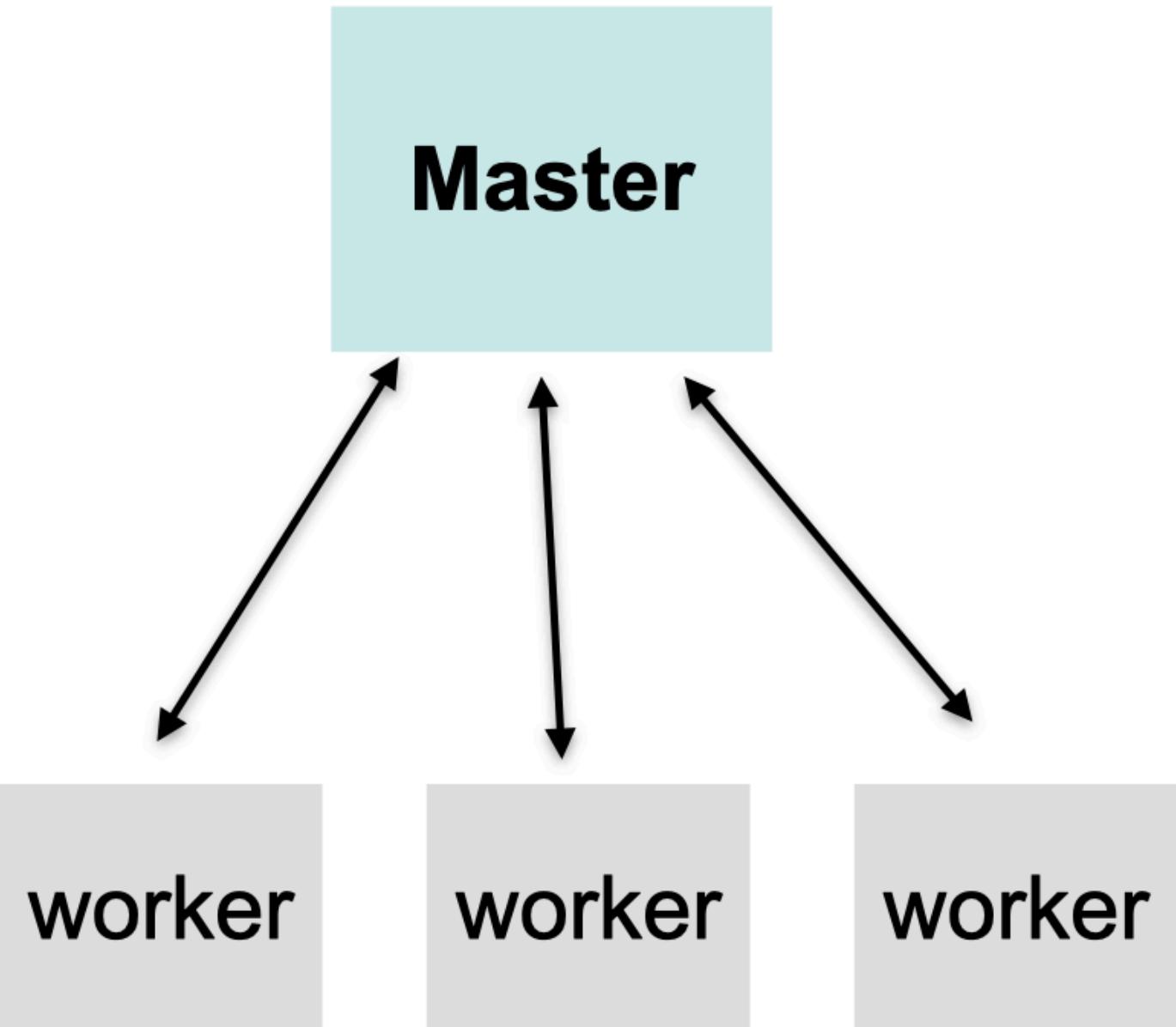
— Mikito Takada¹

¹Distributed systems: for fun and profit

Local



Distributed



How to Distribute

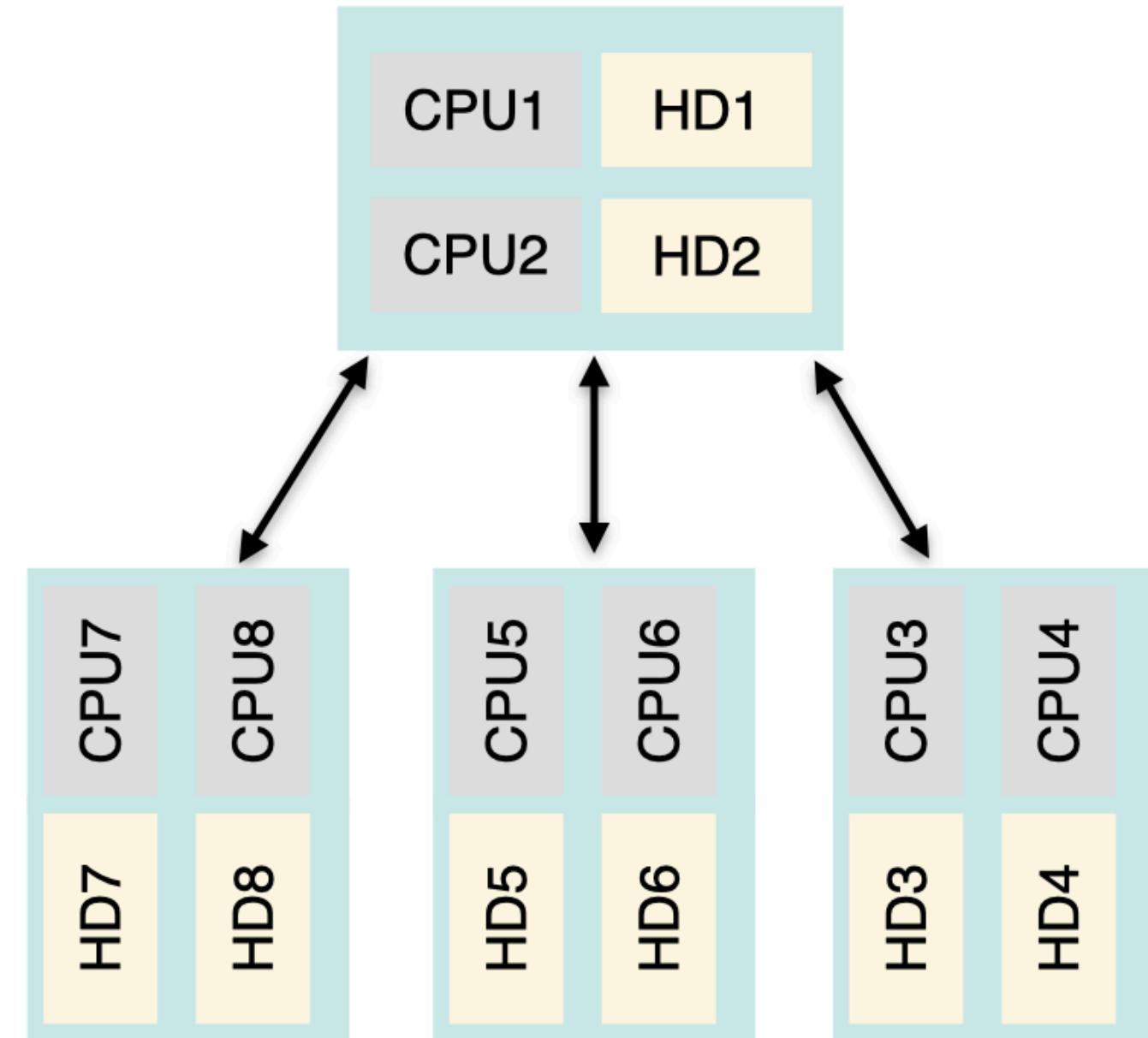
Local



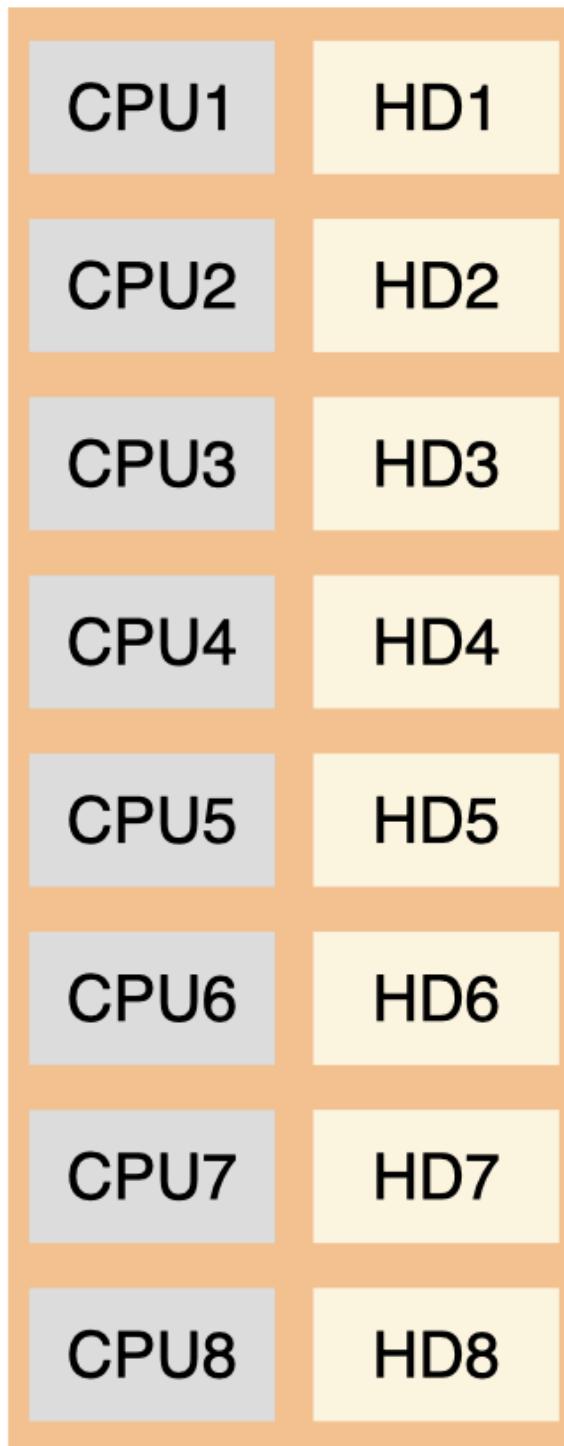
Storage: Hard drives

Computation: CPUs

Distributed

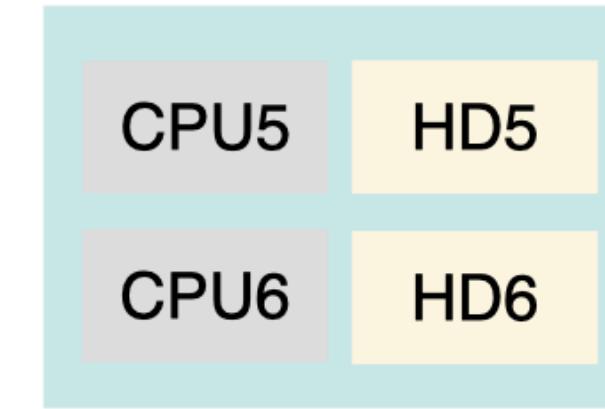
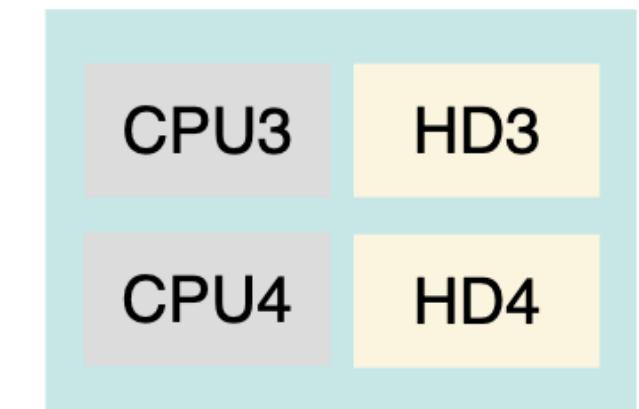
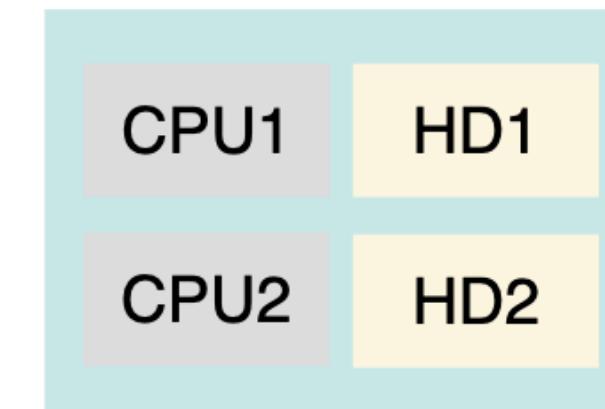


Scale Up versus Out



Up

Out



Streaming Algorithms

we won't talk about these now....

Parallelism

Paralellism

- **Data Parallelism:** Split instances between machines

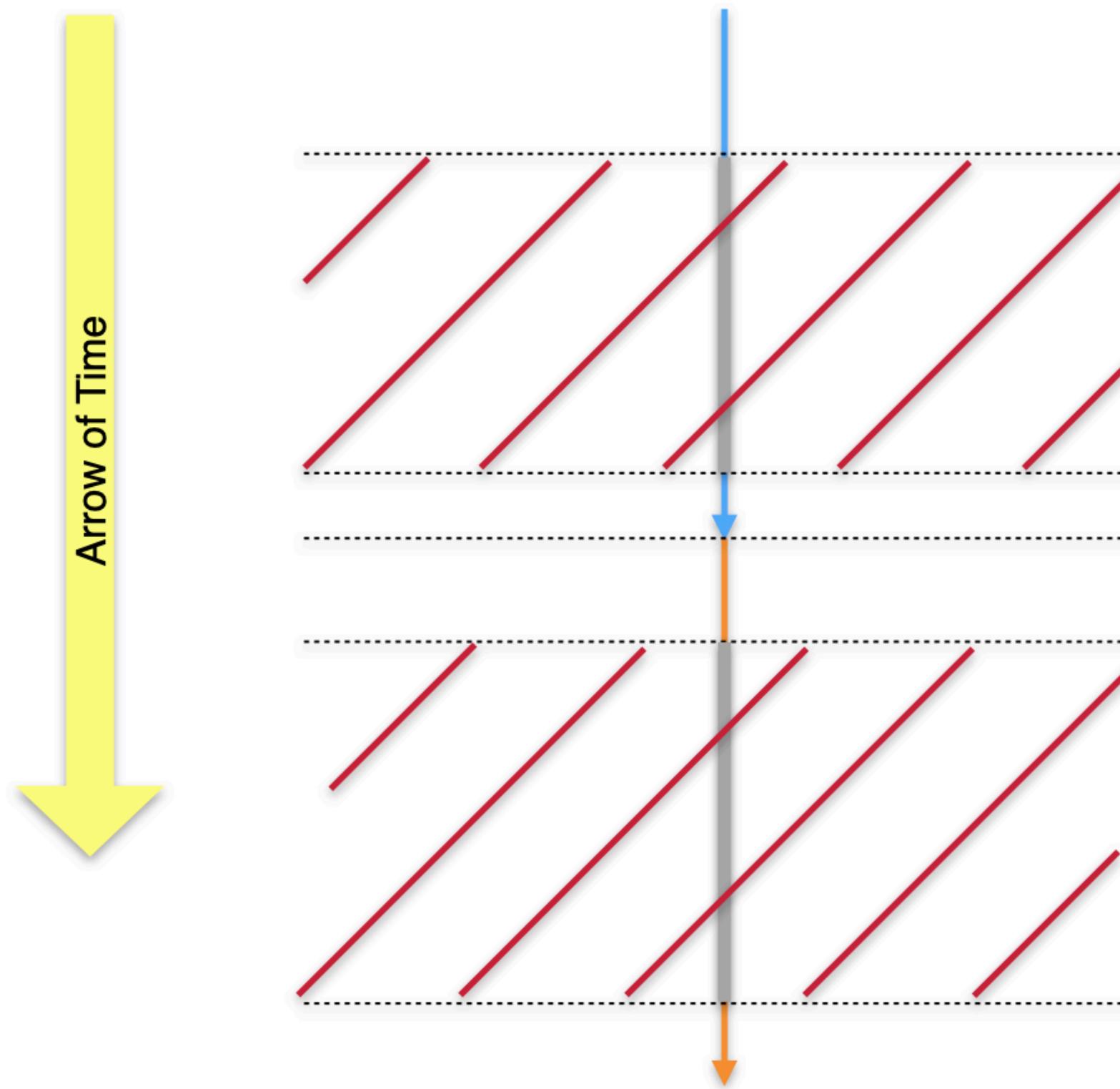
Paralellism

- **Data Paralellism:** Split instances between machines
- **Model Paralellism:** Split parameters between machines

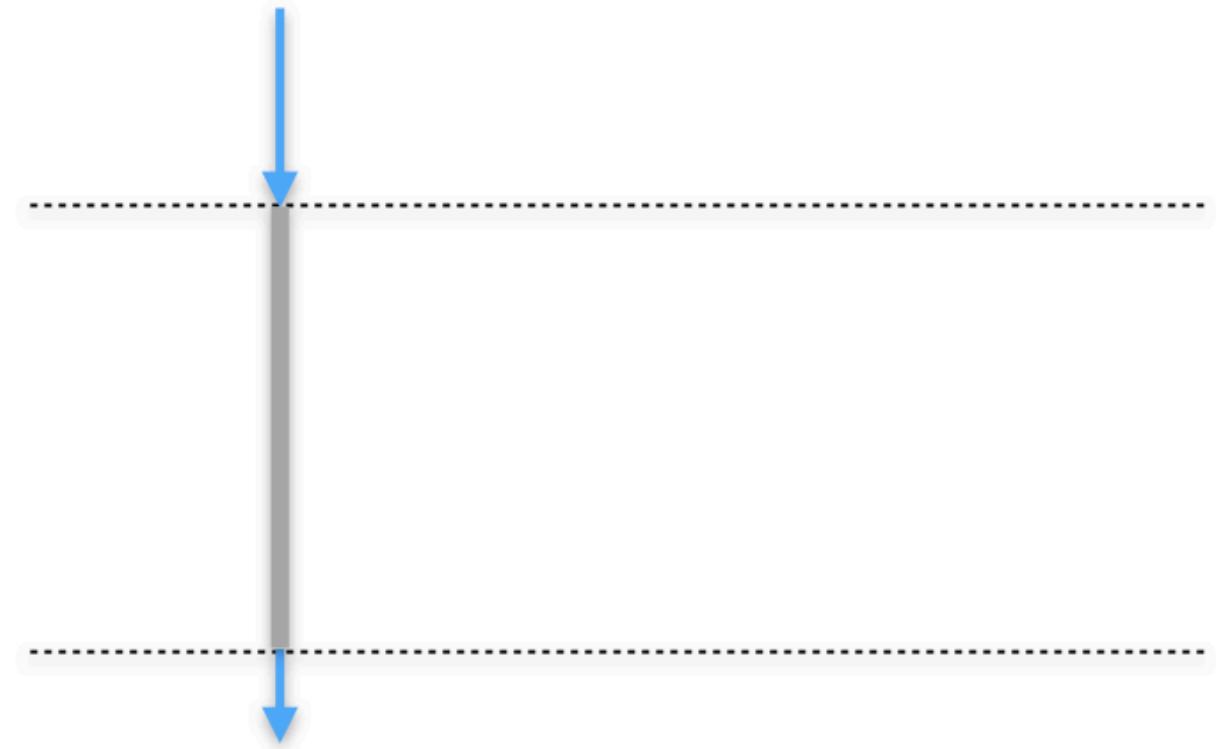
Paralellism

- **Data Paralellism:** Split instances between machines
- **Model Paralellism:** Split parameters between machines
- **Task Paralellism:** Split algorithms between machines

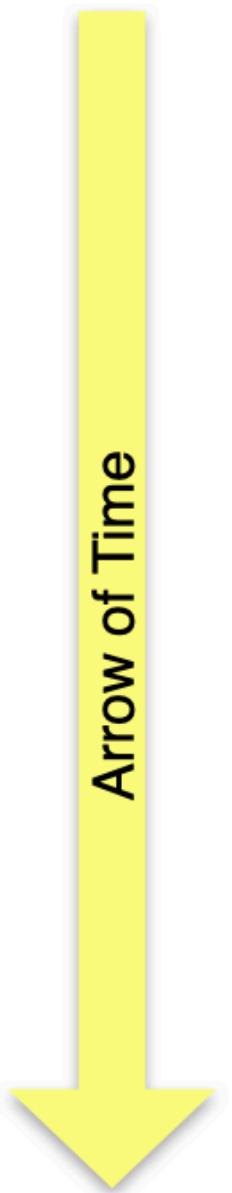
Sequential



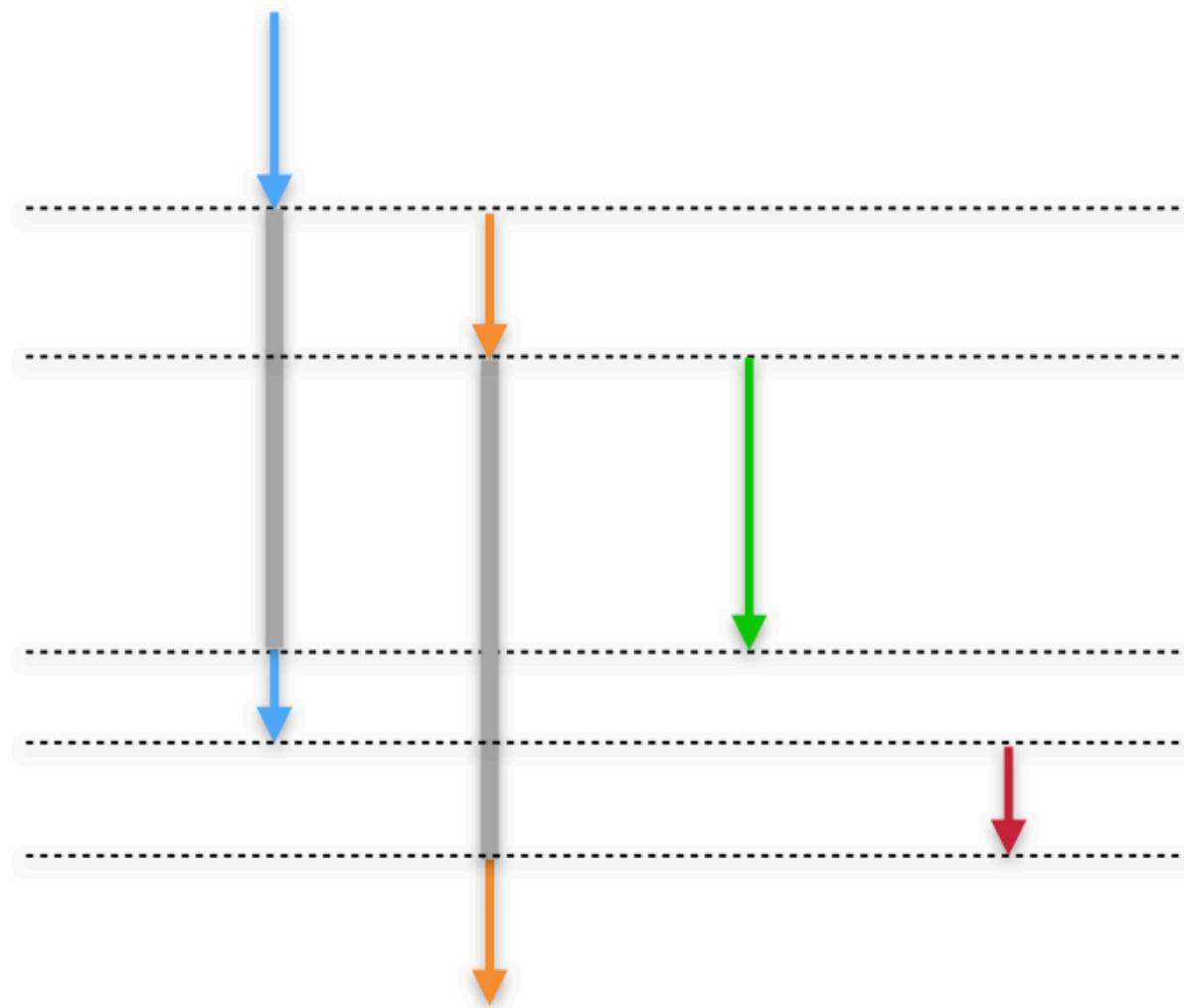
Concurrent



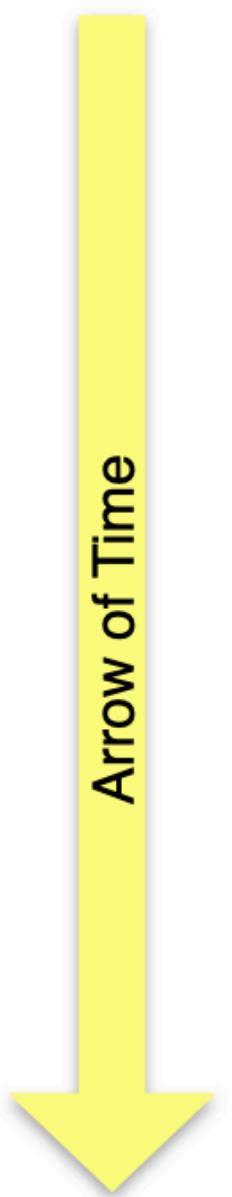
Parallel



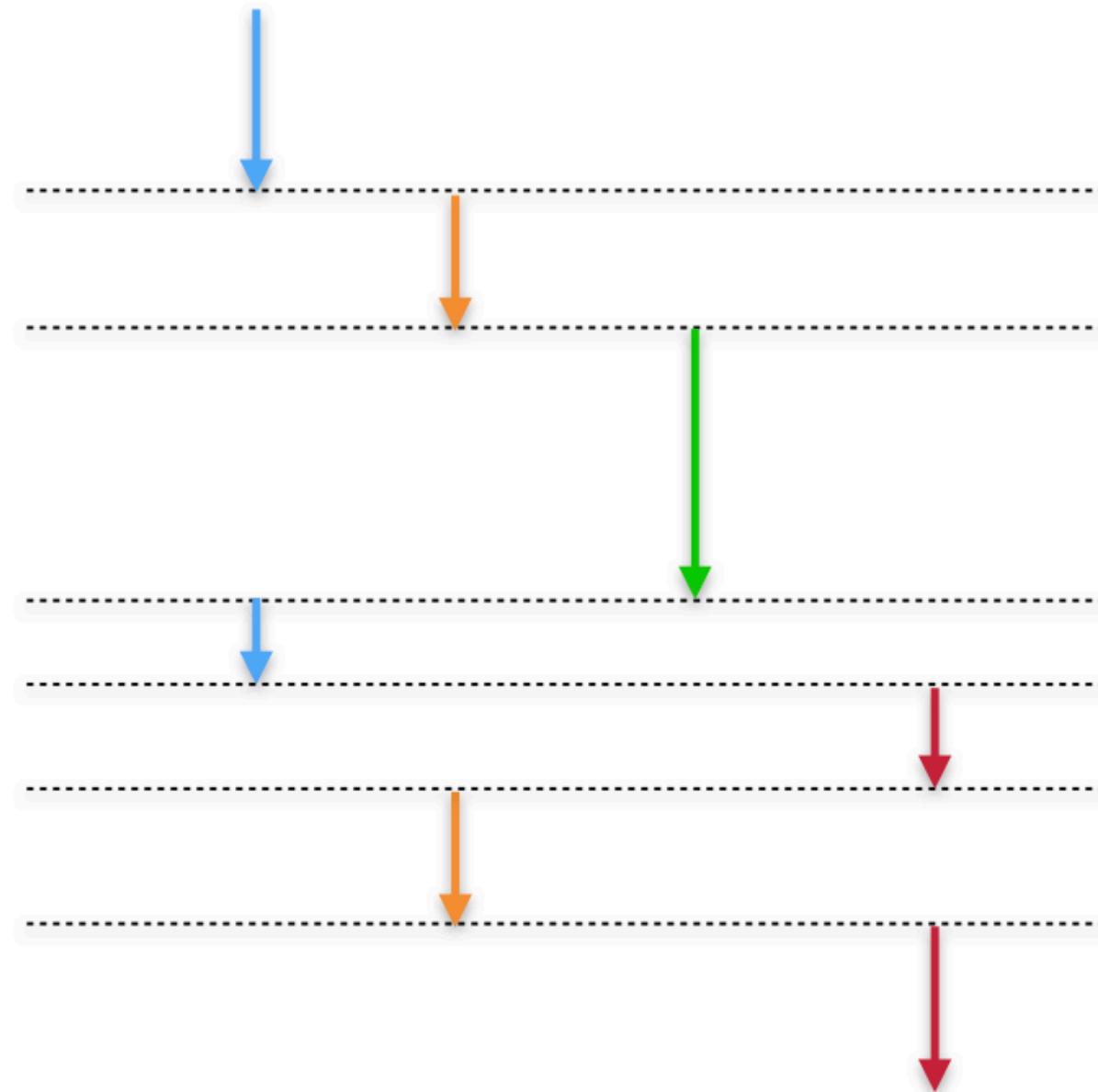
Concurrent



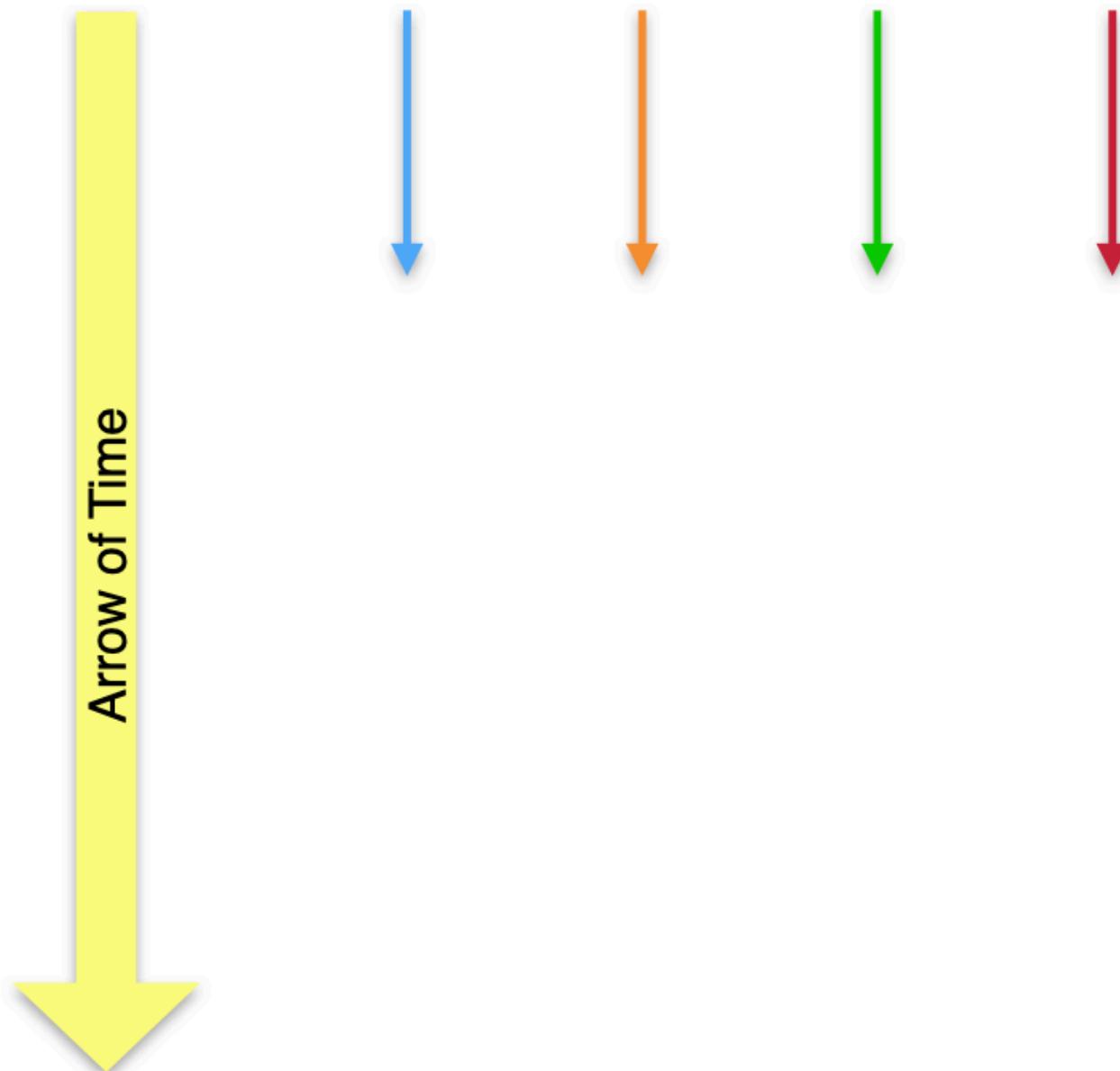
Parallel



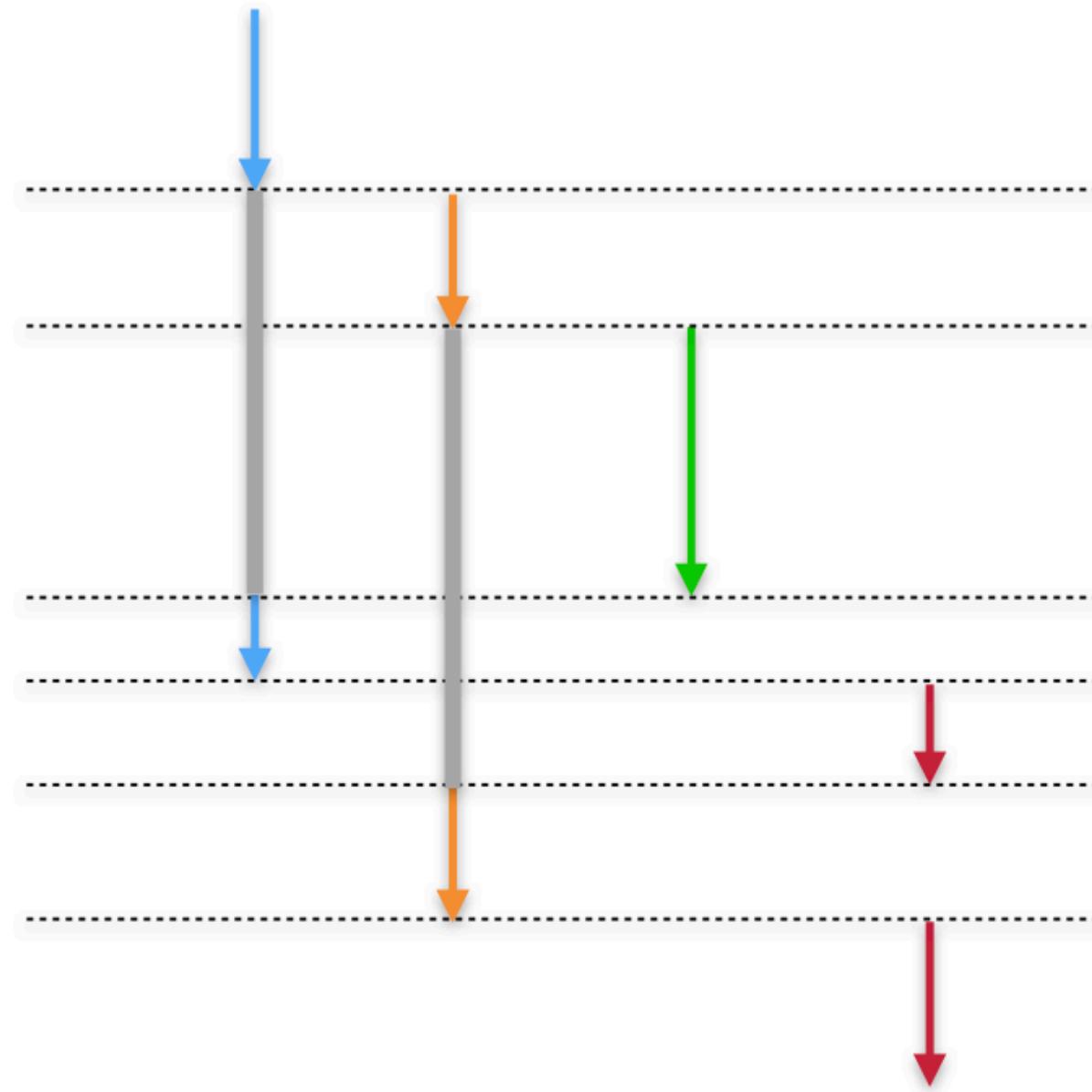
Concurrent



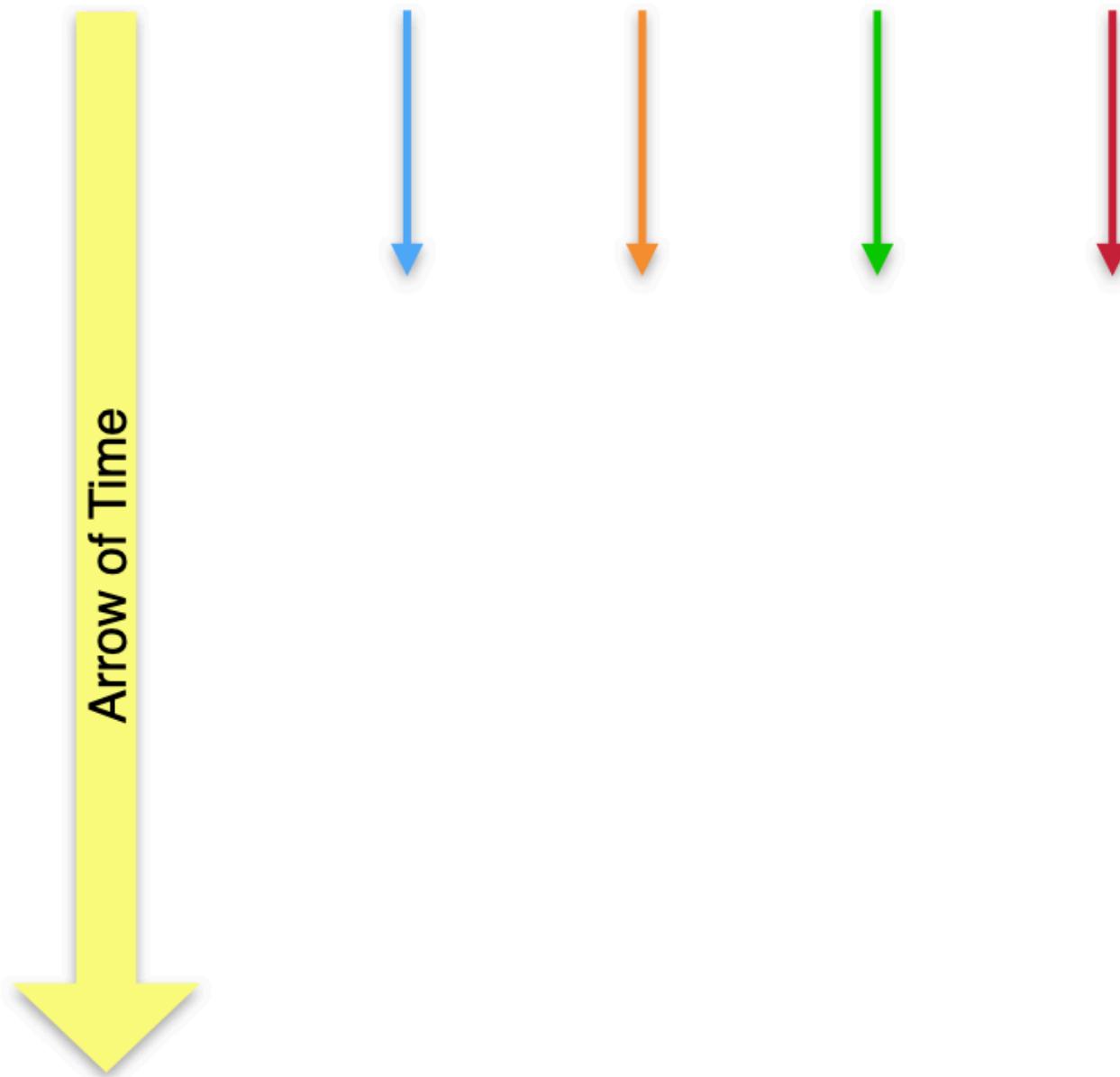
Parallel



Concurrent



Parallel



Problem Solved?

Problem Solved?

- *Coordination*: what happens when and to whom? And how do you synchronize these events

Problem Solved?

- **Coordination:** what happens when and to whom? And how do you synchronize these events
- **Communication:** how do the different nodes in your system talk to one another? And how do you talk to your nodes?

Problem Solved?

- **Coordination:** what happens when and to whom? And how do you synchronize these events
- **Communication:** how do the different nodes in your system talk to one another? And how do you talk to your nodes?
- **Fault Tolerance:** is your system robust to network and machine failures (because they will happen)?

This is what Ray helps us with



- **Coordination:** what happens when and to whom? And how do you synchronize these events
- **Communication:** how do the different nodes in your system talk to one another? And how do you talk to your nodes?
- **Fault Tolerance:** is your system robust to network and machine failures (because they will happen)?

CAP Theorem²

Consistency: all nodes see the same data at the same time

Availability: every request receives a response about whether it succeeded or failed

Partition tolerance: the system continues to operate despite arbitrary partitioning due to network failures

²Choose two: except you can't really sacrifice **P**!

Break

Why Ray?

Why Ray?

- *Distributed computing framework* meant for a new generation of AI applications.

Why Ray?

- *Distributed computing framework* meant for a new generation of AI applications.
- Emerged from *RL applications*

Why Ray?

- *Distributed computing framework* meant for a new generation of AI applications.
- Emerged from *RL applications*
- *More generic* programming model than Spark/Flink/Hadoop/etc.

Why Ray?

- *Distributed computing framework* meant for a new generation of AI applications.
- Emerged from *RL applications*
- More generic programming model than Spark/Flink/Hadoop/etc.
- More fault tolerant than MPI

multiprocess meets RPC....

The Ray programming model

The Ray programming model

- Tasks (remote functions)

The Ray programming model

- Tasks (remote functions)
- Actors (remote classes)

Tasks

```
@ray.remote  
def index(url):  
    return requests.get(url)  
  
futures = [index.remote(site) for site in internet]  
print(ray.get(futures))
```

love code

Actors

```
@ray.remote
class Child(object):
    def __init__(self):
        self.name = random.choice(baby_names)
        self.age = 1

    def grow(self):
        self.age += 1

    def greet(self):
        return (
            f'My name is {self.name} '
            f'and I am {self.age} years old'
        )

children = [Child.remote() for i in range(4)]
[c.grow.remote() for c in children]
futures = [c.greet.remote() for c in children]
print(ray.get(futures))
```

love code

Break

Distributed Training Architectures

Distributed Training Architectures

- Parallelized Stochastic Gradient Descent (2010)

Distributed Training Architectures

- Parallelized Stochastic Gradient Descent (2010)
- Hogwild! (2011)

Distributed Training Architectures

- Parallelized Stochastic Gradient Descent (2010)
- Hogwild! (2011)
- Downpour SGD (DistBelief) (2012)

Distributed Training Architectures

- Parallelized Stochastic Gradient Descent (2010)
- Hogwild! (2011)
- Downpour SGD (DistBelief) (2012)
- AllReduce (2014)

Distributed Training Architectures

- Parallelized Stochastic Gradient Descent (2010)
- Hogwild! (2011)
- Downpour SGD (DistBelief) (2012)
- AllReduce (2014)
- Parameter Server (2014)

Will parallelization help me?

Amdahl's Law

$$\text{Speedup} = \frac{1}{(1 - F) + \frac{F}{S}}$$

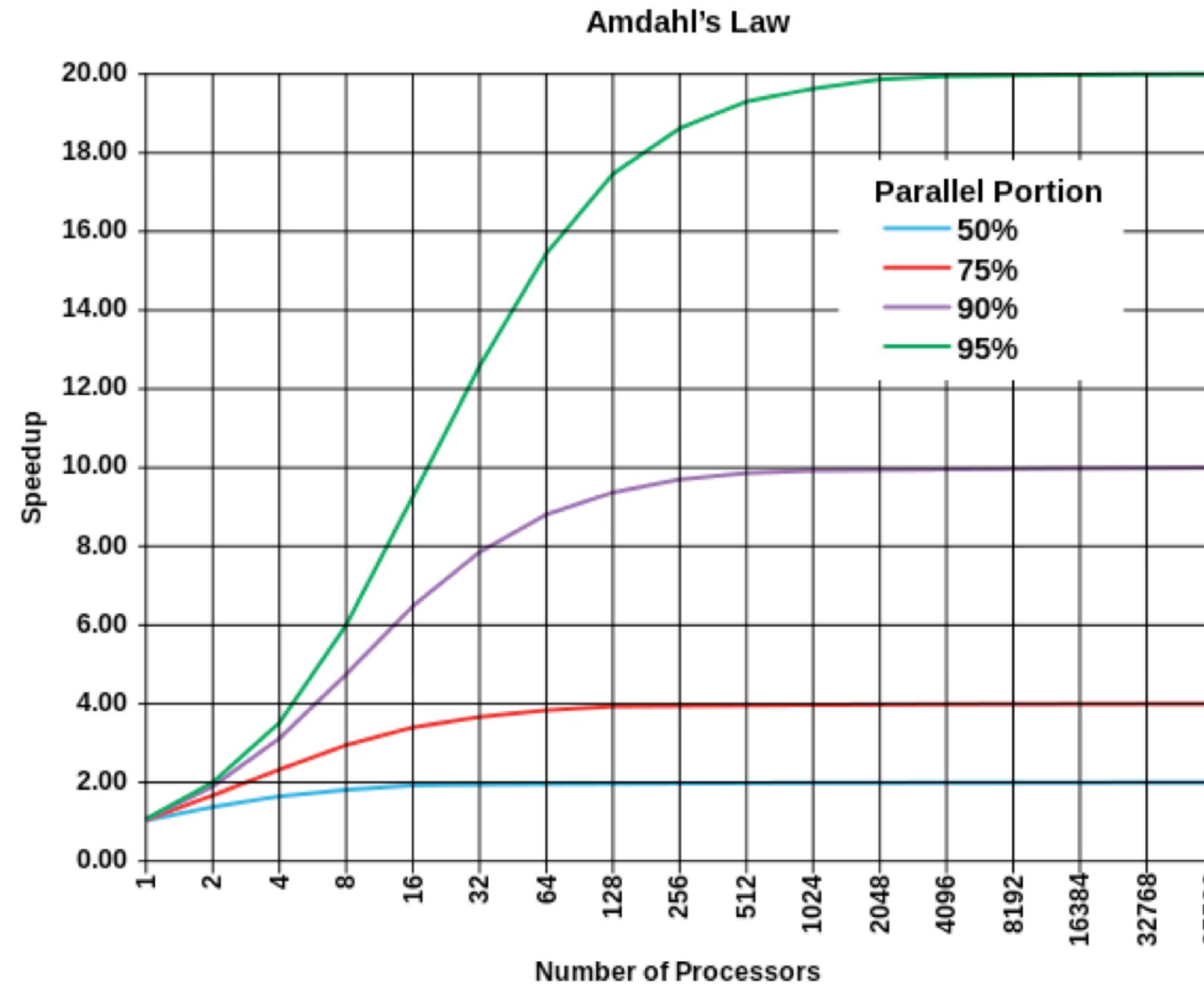
Fraction that can be parallelized

Fraction that cannot be parallelized

Number of Processors

The diagram shows the Amdahl's Law formula: Speedup = 1 / ((1 - F) + (F / S)). Red annotations provide clarity: 'Fraction that can be parallelized' points to the term F/S; 'Fraction that cannot be parallelized' points to the term 1 - F; and 'Number of Processors' points to the variable S.

Amdahl's Law



Amdahl's Law

20x improvement
25% of the time

$$\text{Speedup} = \frac{1}{(1 - 0.25) + \frac{0.25}{20}}$$

Fraction that cannot be parallelized

Fraction that can be parallelized

Number of Processors

The diagram shows the Amdahl's Law formula: Speedup = 1 / ((1 - fraction_parallelizable) + (fraction_parallelizable / number_processors)). Red annotations explain each term: 'Fraction that cannot be parallelized' points to 1 - 0.25; 'Fraction that can be parallelized' points to 0.25; and 'Number of Processors' points to 20.

Amdahl's Law

20x improvement
25% of the time

Only 1.3x Speedup!

$$\text{Speedup} = \frac{1}{0.75 + \frac{0.25}{20}} = 1.31$$

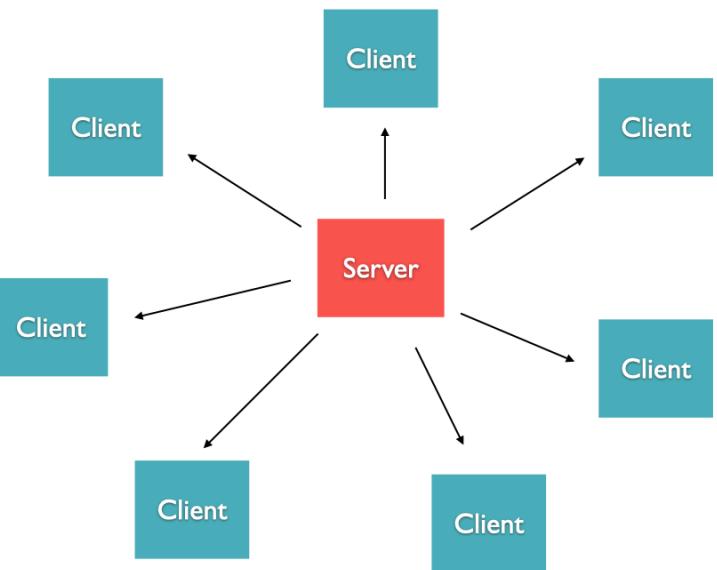
Even with 20 processors

Moral?

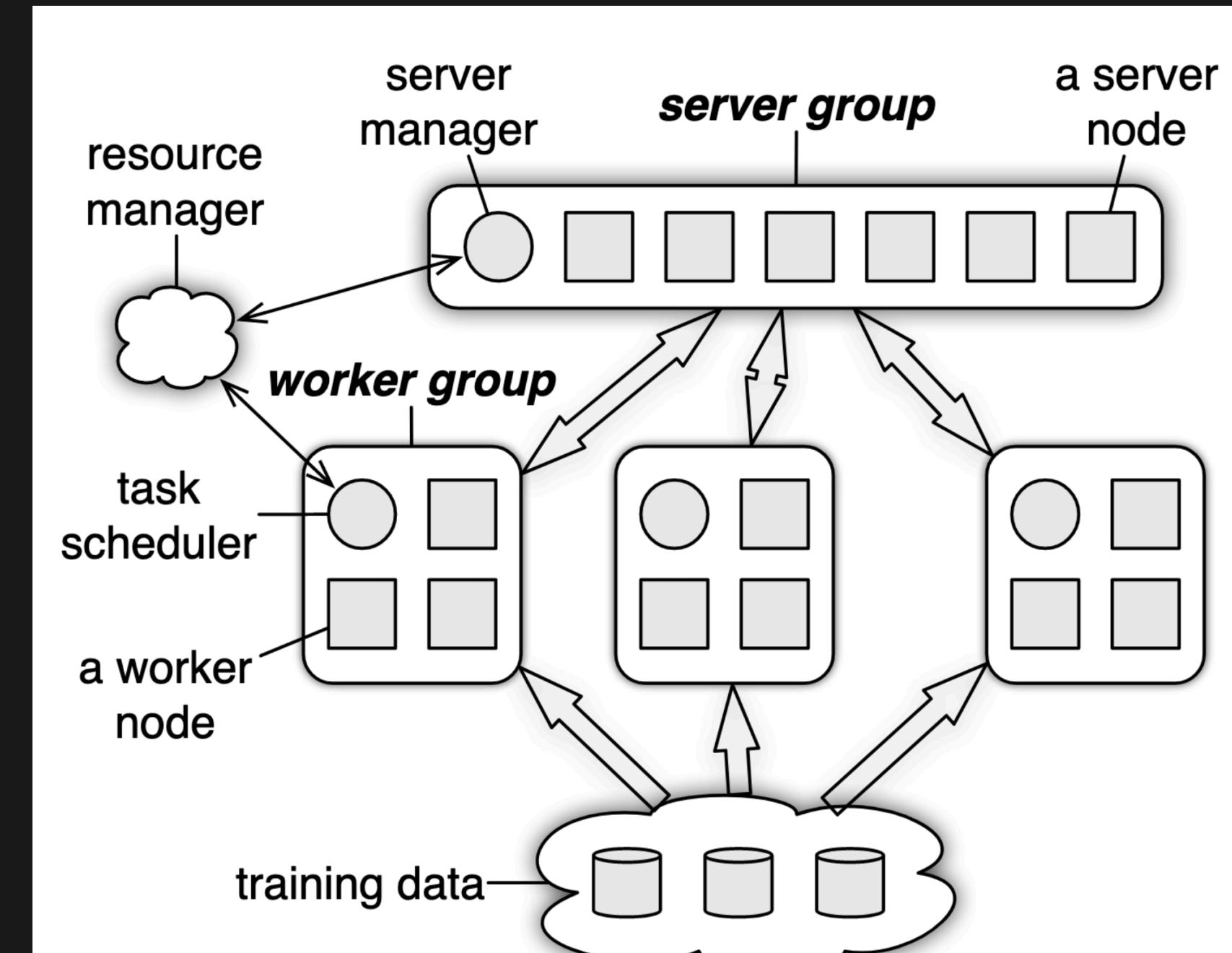
It doesn't how many machines you
throw at a problem... at some point it
is just not going to run any faster

Distributed training with the parameter server

Parameter Server

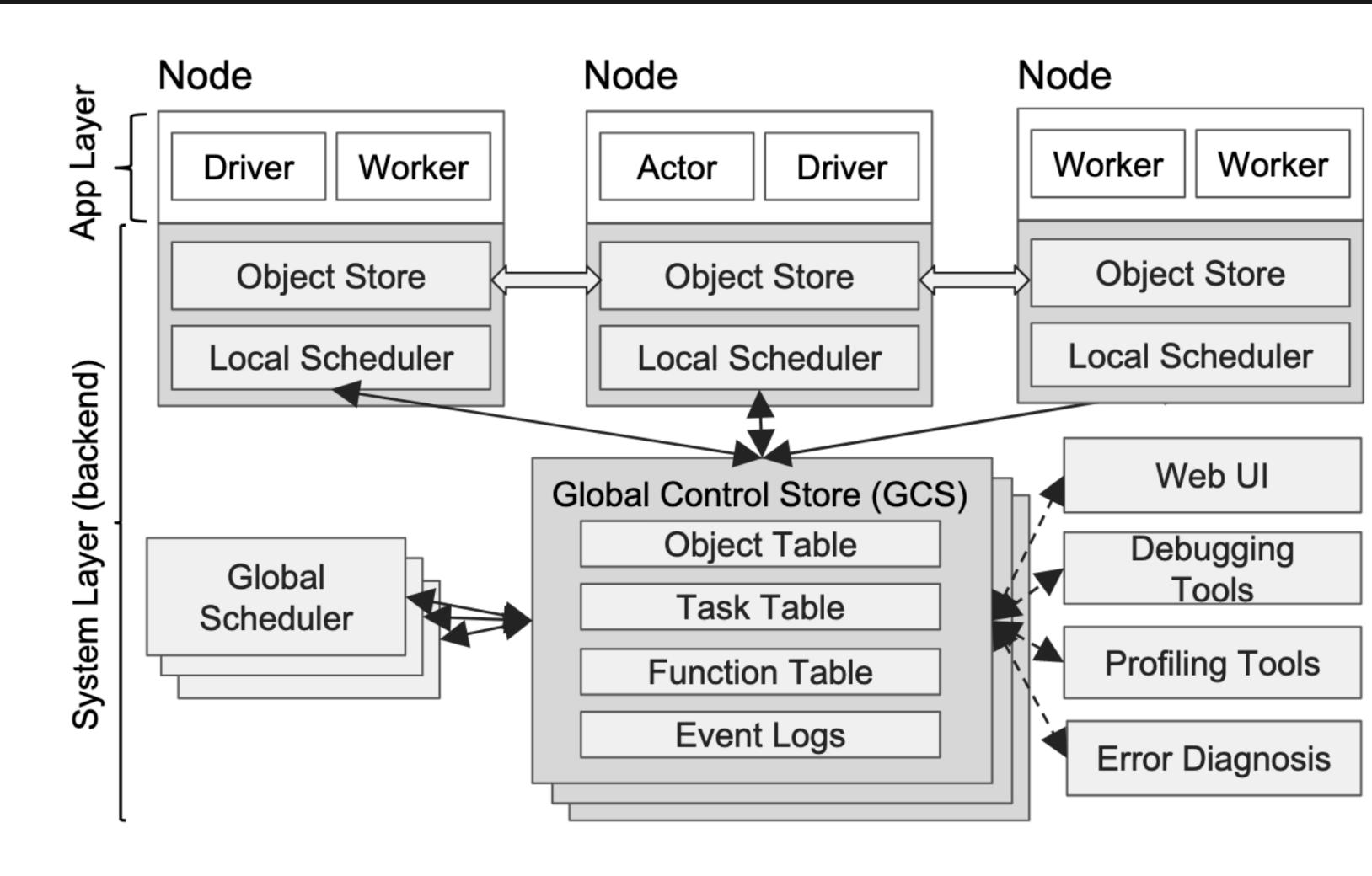


System rather than
algorithm³



³ [Scaling Distributed Machine Learning with the Parameter Server](#)

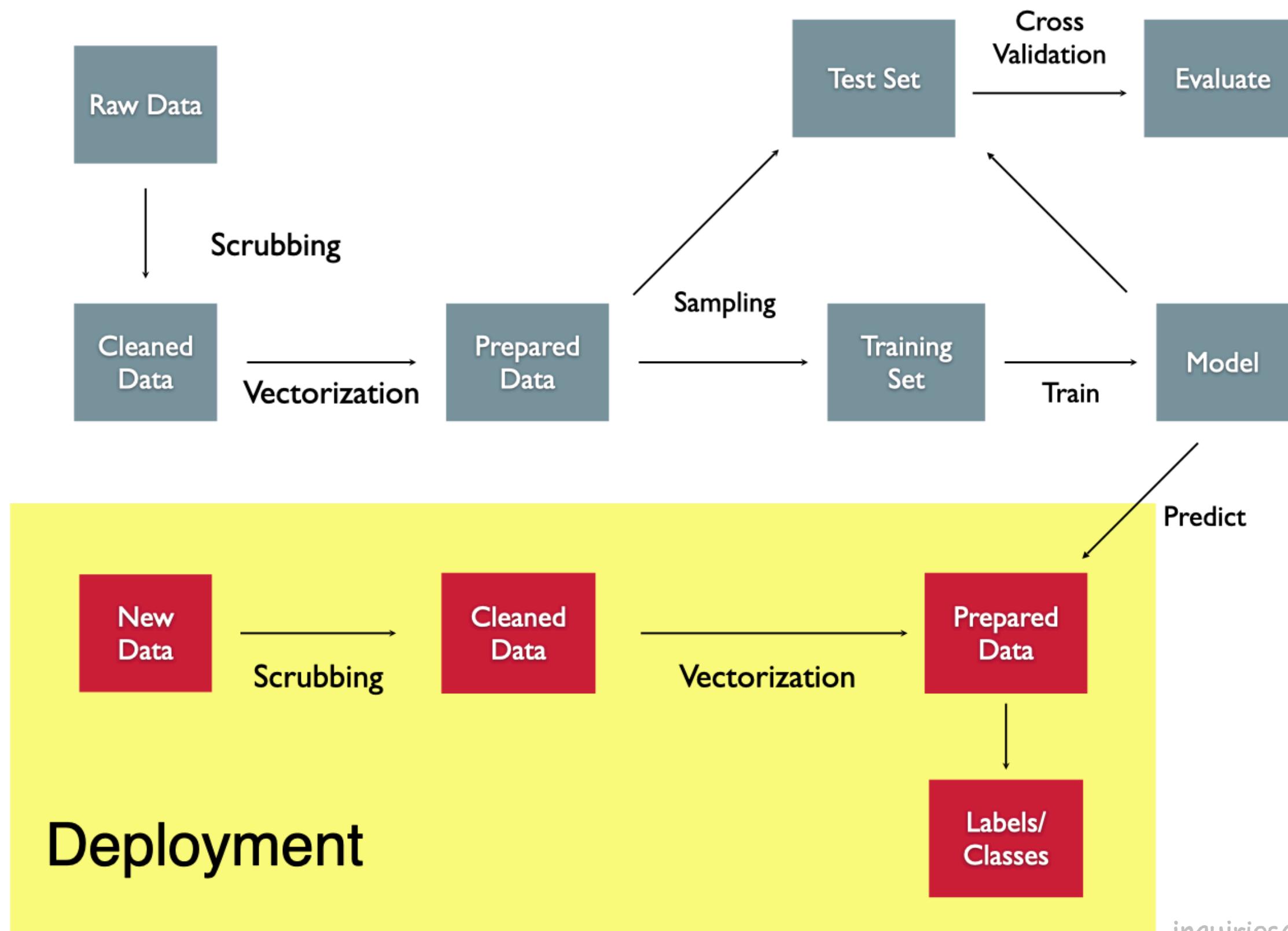
Ray internals⁴



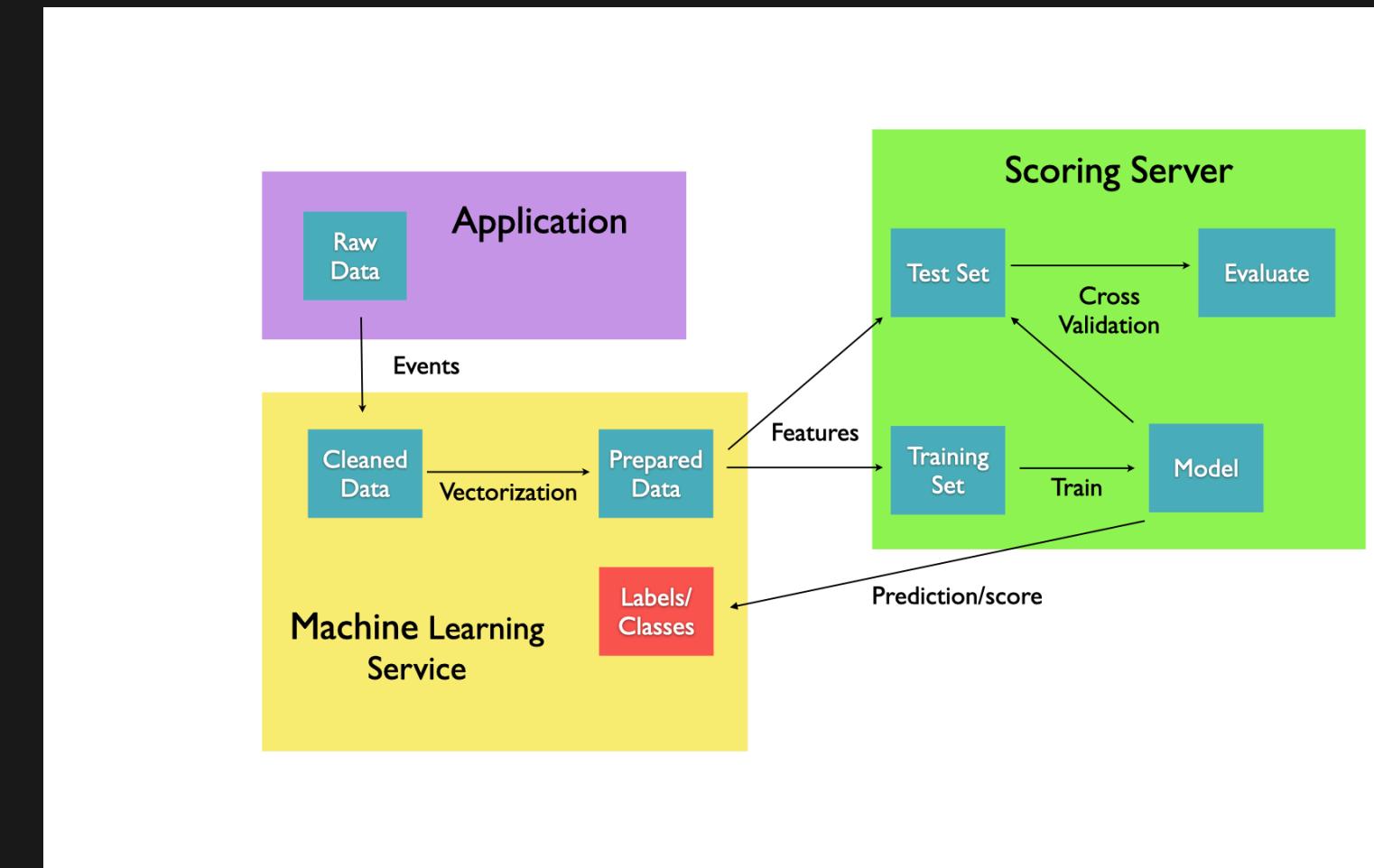
⁴ [Ray: A Distributed Framework for Emerging AI Applications](#)

love code

Deploying ML

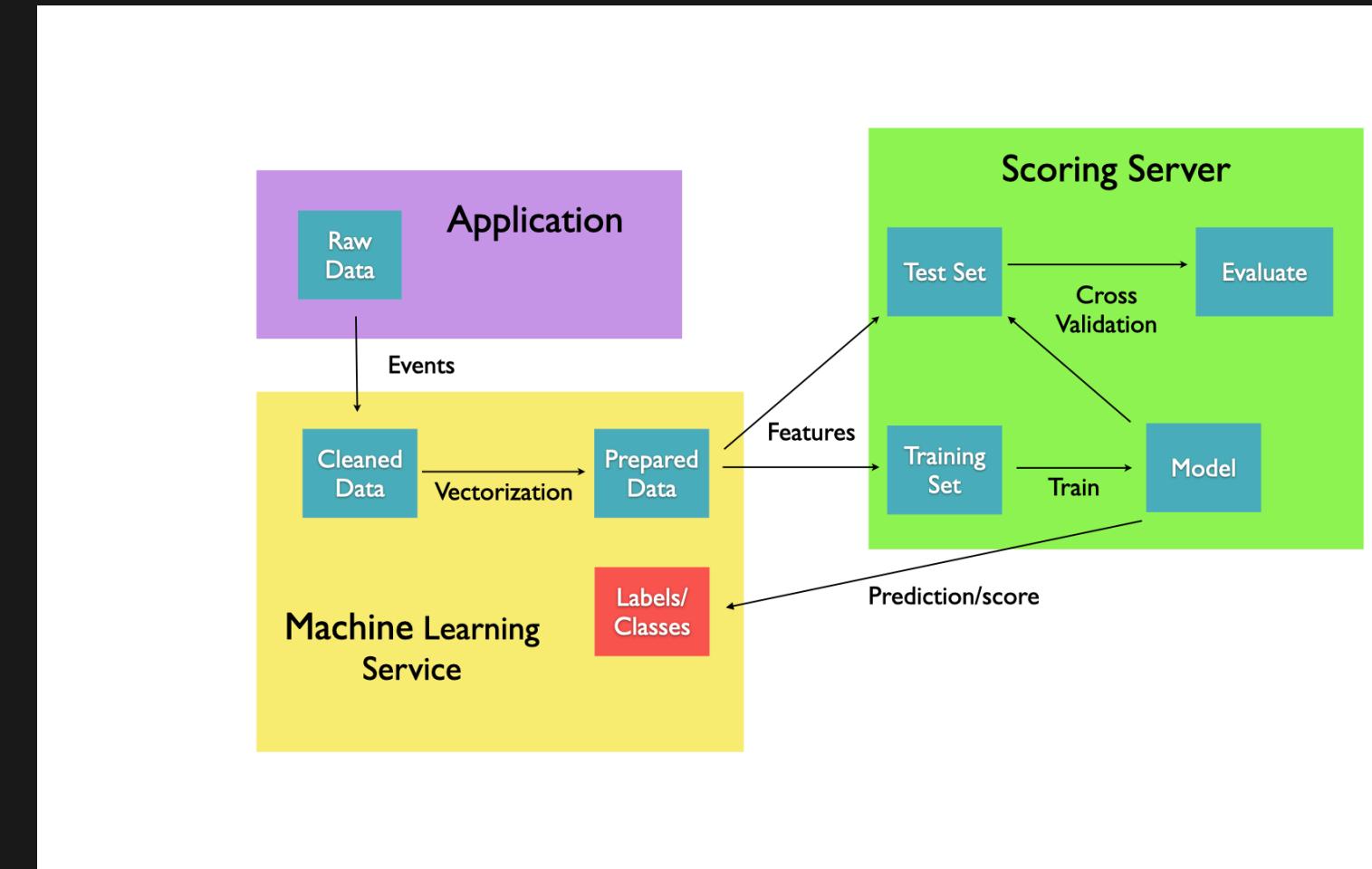


ML infrastructure is varied



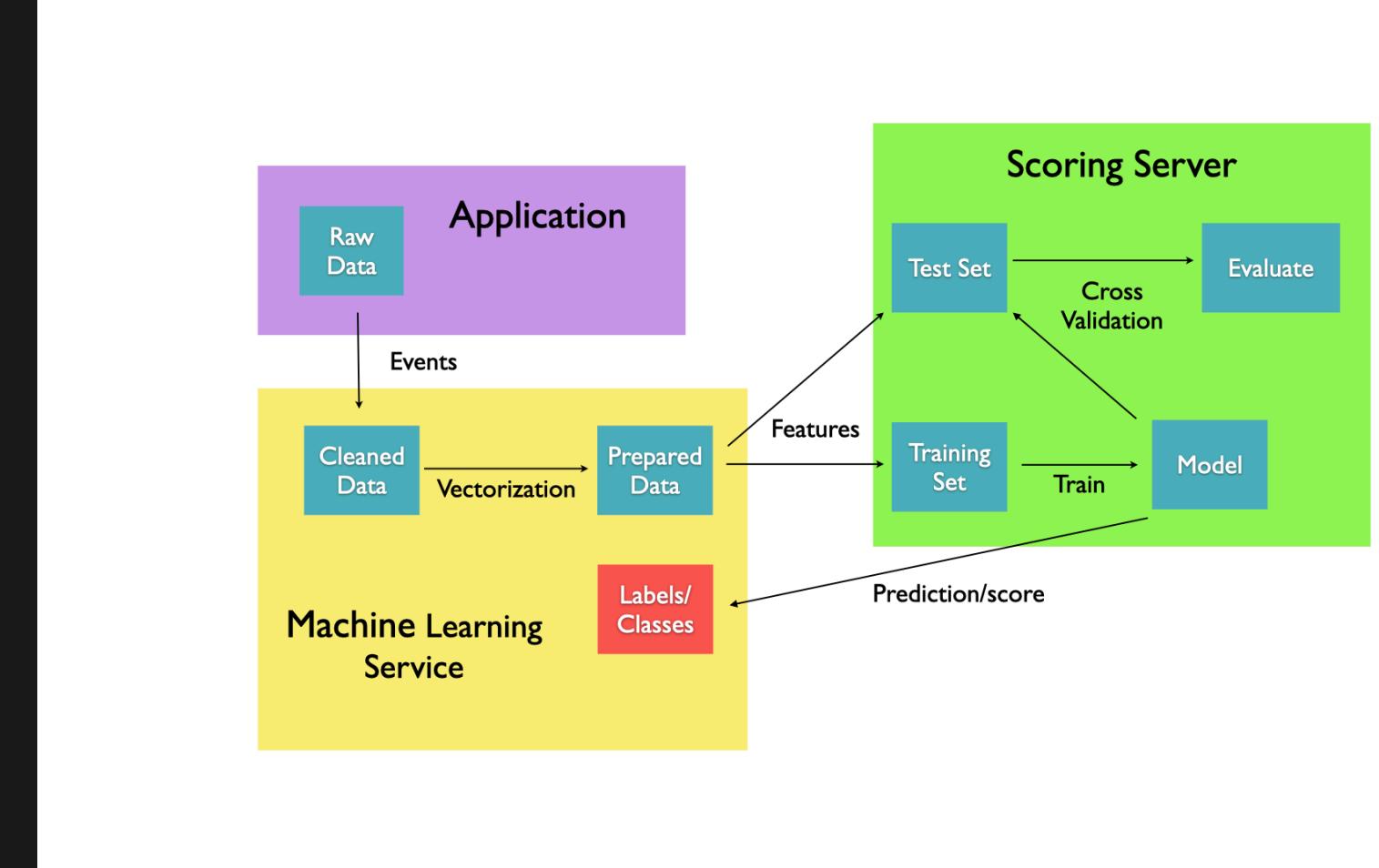
ML infrastructure is varied

→ MLaaS (SageMaker, Google AutoML, etc.)



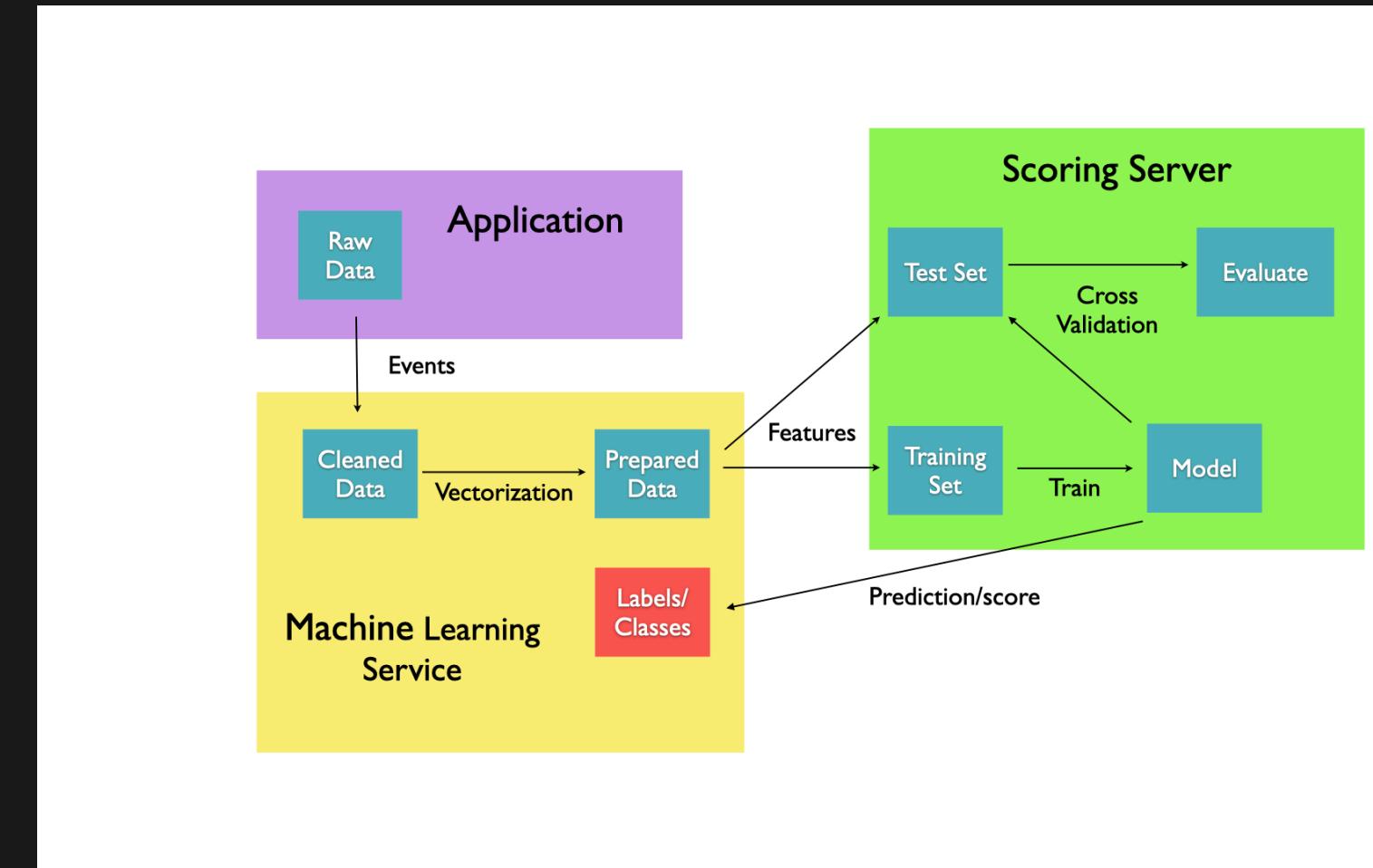
ML infrastructure is varied

- MLaaS (SageMaker, Google AutoML, etc.)
- Containers (Fargate, GKE, etc.)



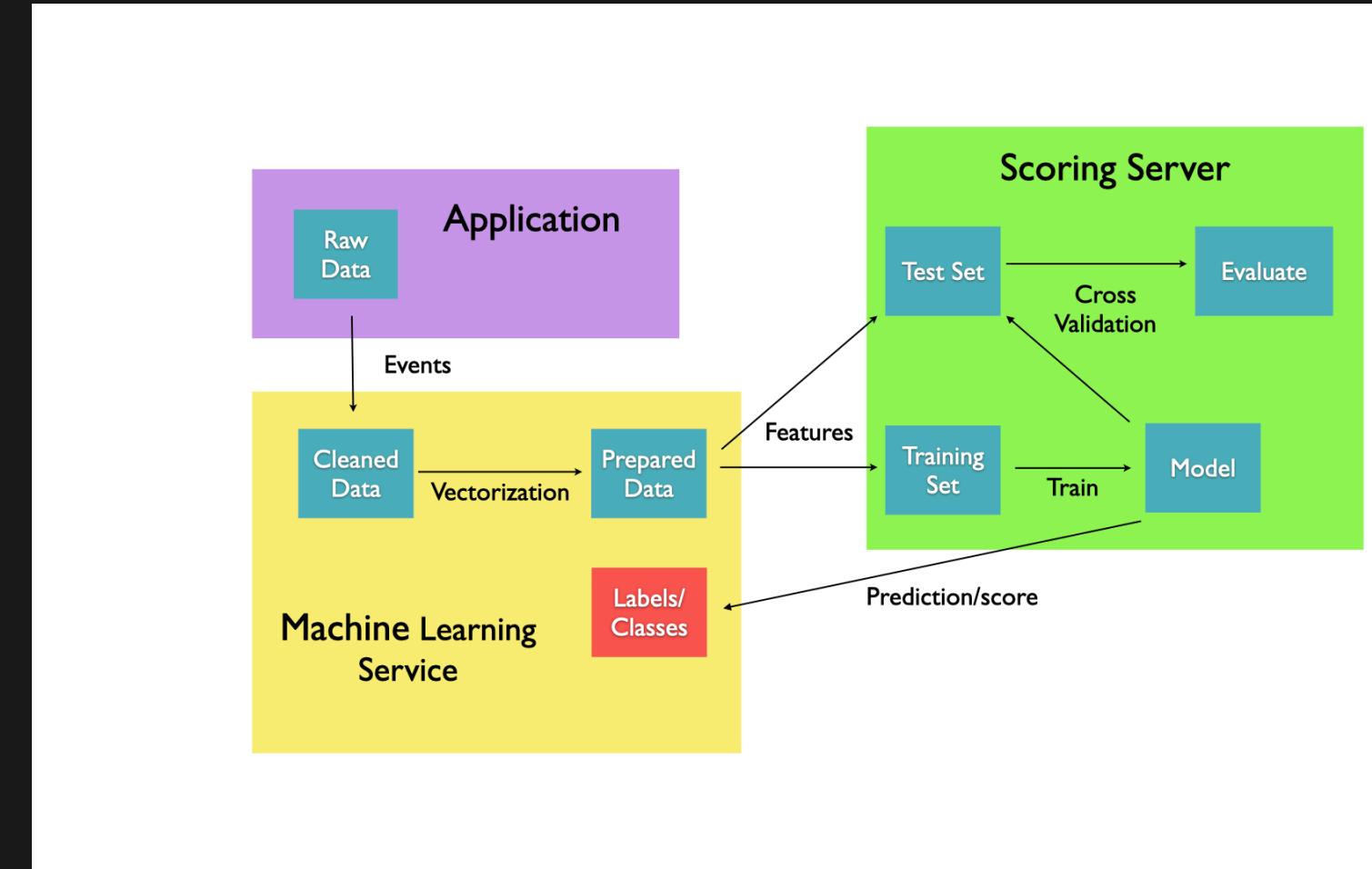
ML infrastructure is varied

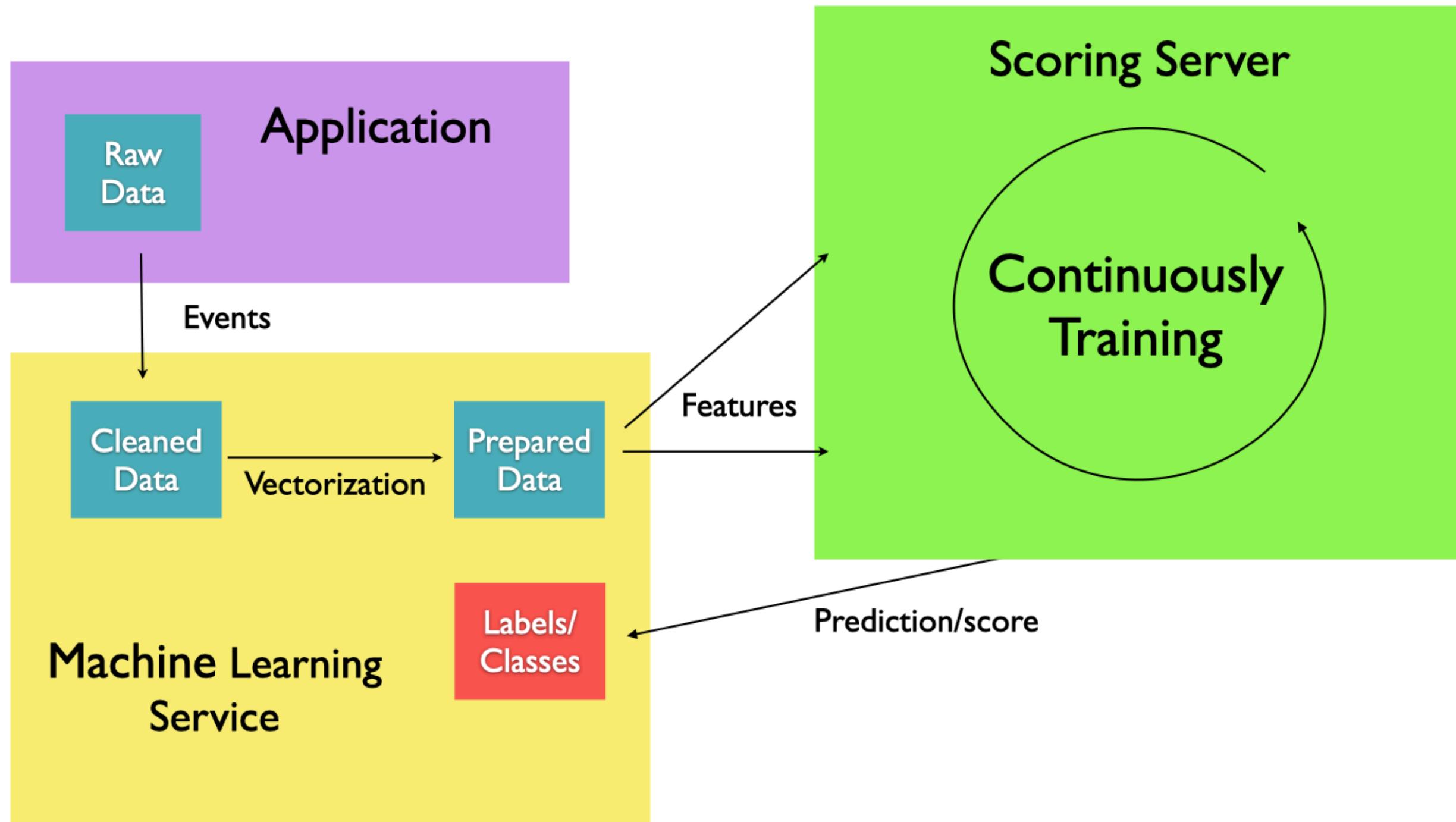
- MLaaS (SageMaker, Google AutoML, etc.)
 - Containers (Fargate, GKE, etc.)
 - IaaS (Ec2, GCP, etc.)

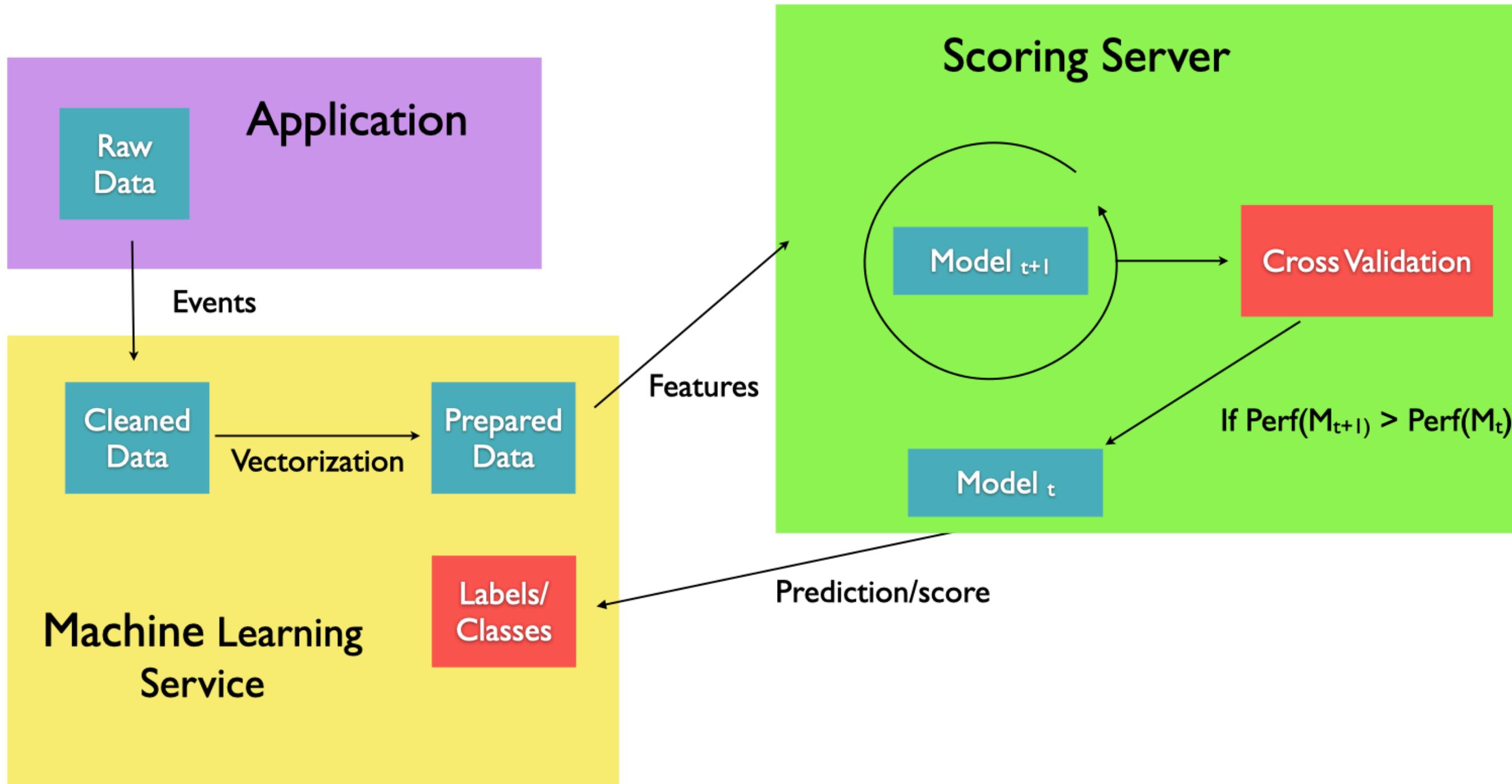


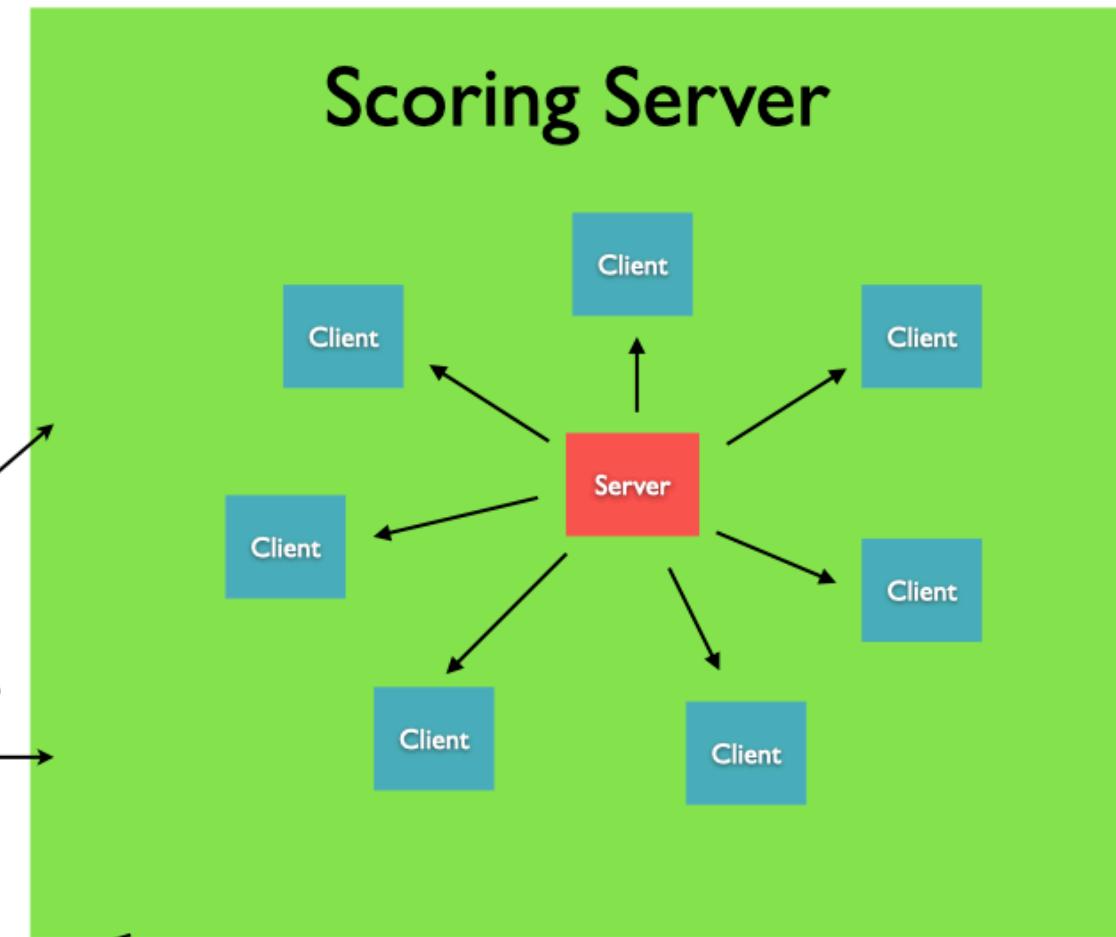
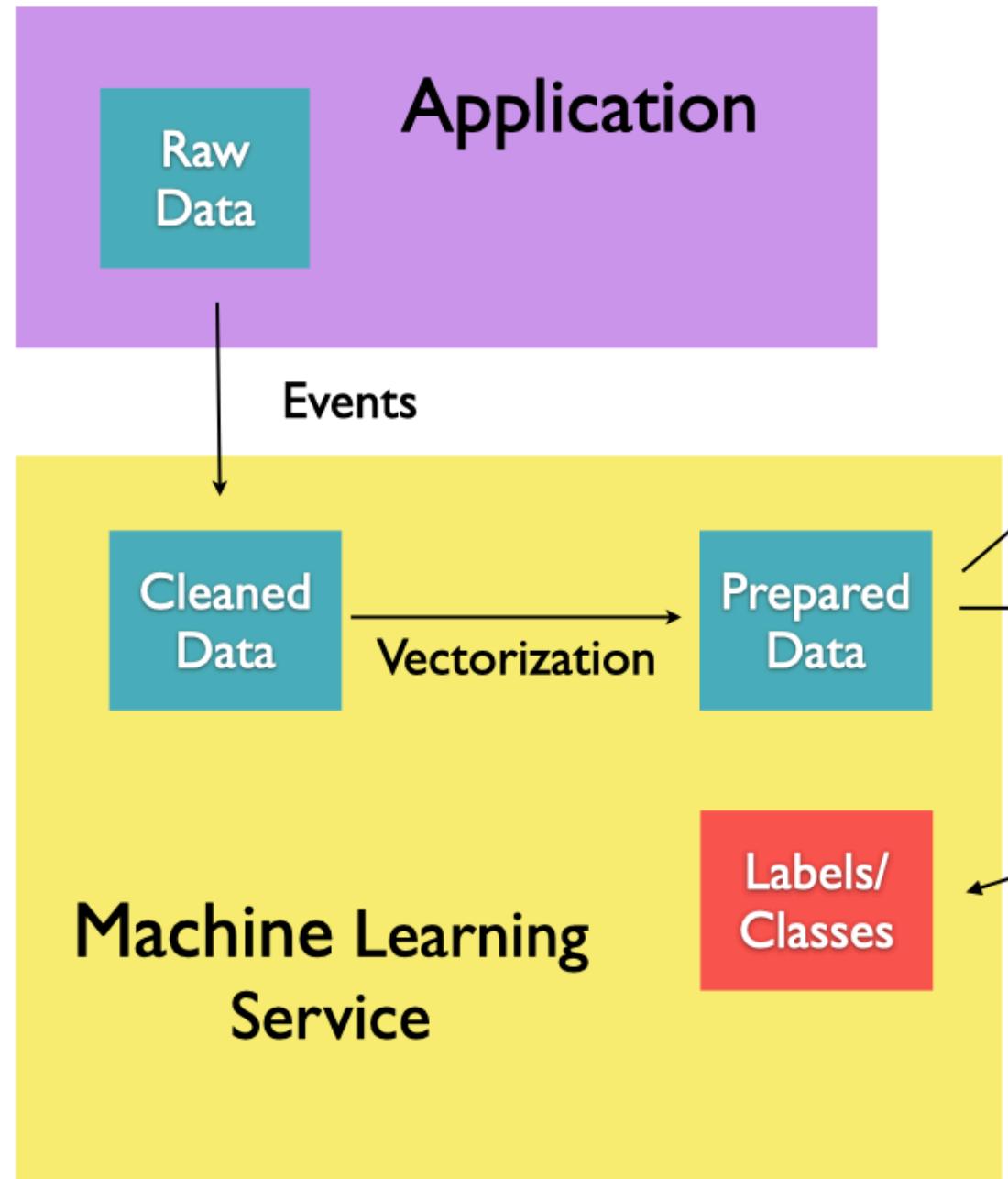
ML infrastructure is varied

- MLaaS (SageMaker, Google AutoML, etc.)
- Containers (Fargate, GKE, etc.)
- IaaS (Ec2, GCP, etc.)
- On premise/data center









Next Steps

Learning More DistML

Scaling up Machine Learning

SCALING UP MACHINE LEARNING

Parallel and Distributed Approaches

EDITED BY

RON BEKKERMAN
MIKHAIL BILENKO
JOHN LANGFORD

Future Live Training

Future Live Training

- Bayesian Hyperparameter Optimization (Sept. 10)

Q & A

Materials, mailing list, questions, feedback 

<https://sr.ht/~hyphaebeast/ray-live-training/>