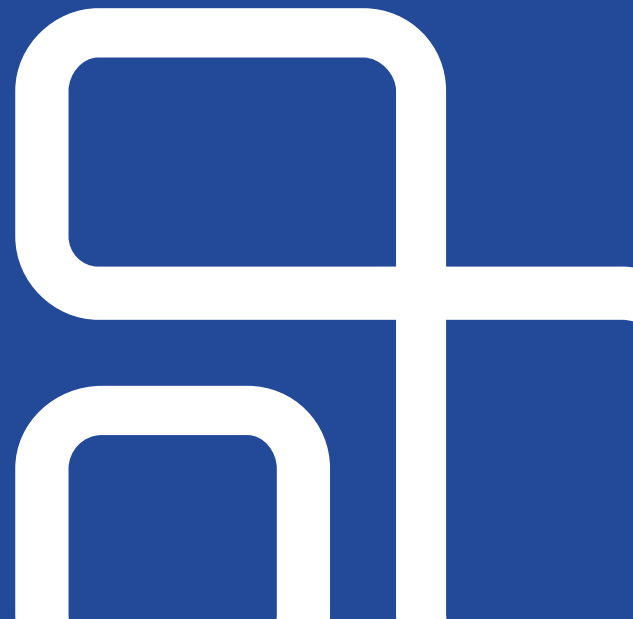


# Scaling AI/ML Workloads with Ray Ecosystem

---

**Jules S. Damji, @2twitme**

Lead Developer Advocate, Ray Team @ Anyscale  
Data Council, Austin, TX March 23, 2022



## Overview

- **Why & What Ray & Ray Ecosystem**
  - **Ray Architecture & Components**
  - **Ray Core APIs**
  - **Ray Native ML Libraries**
    - Ray Tune, XGBoost-Ray
  - **Demo**
    - Scaling ML workloads
  - **Q & A**
-

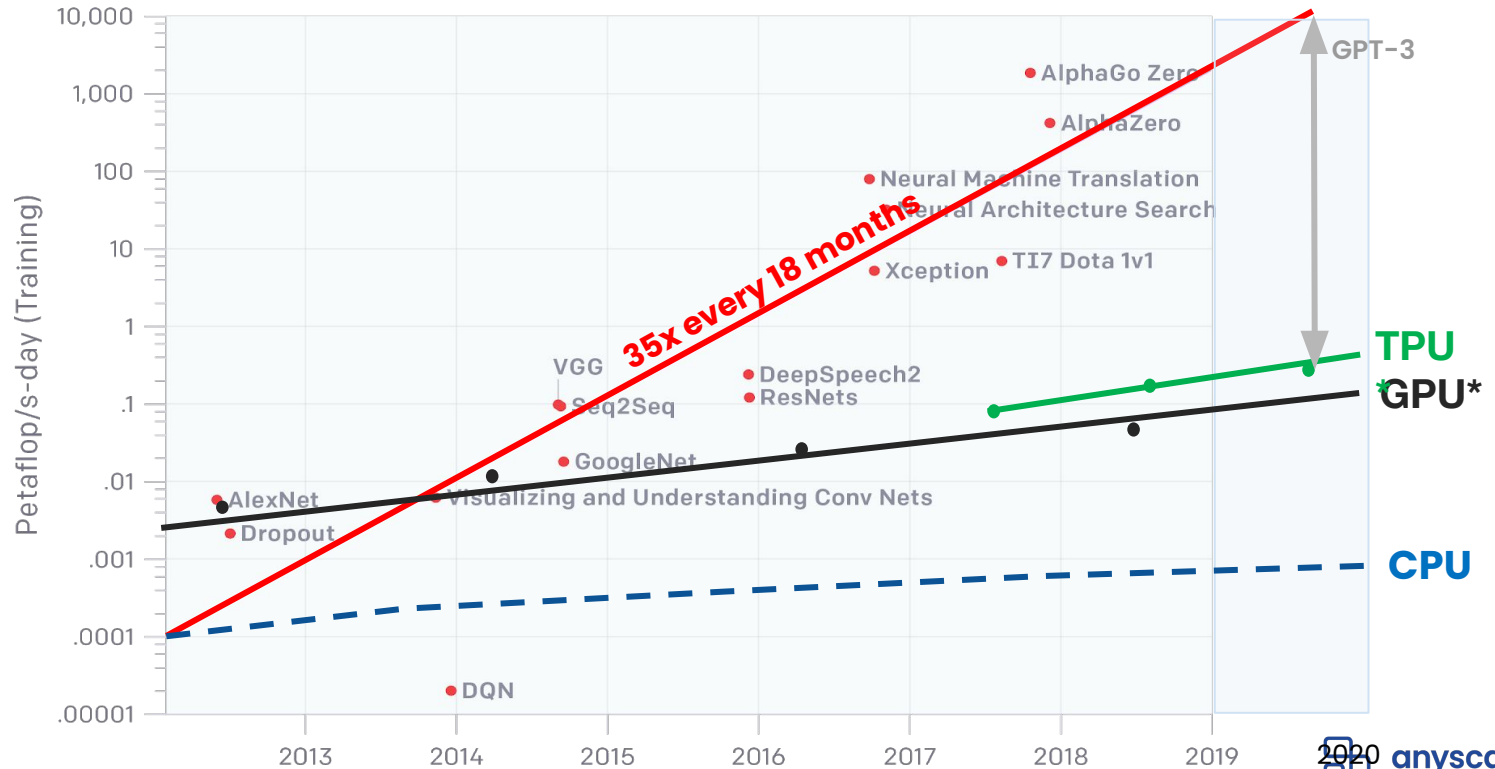
# Why Ray?

- **Machine learning is pervasive in every domain**
- Distributed machine learning is becoming a necessity
- Distributed systems is notoriously hard

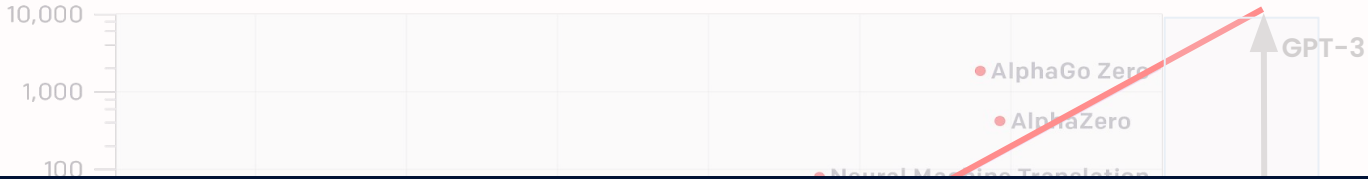
# Why Ray?

- Machine learning is pervasive in every domain
- **Distributed machine learning is becoming a necessity**
- Distributed systems is notoriously hard

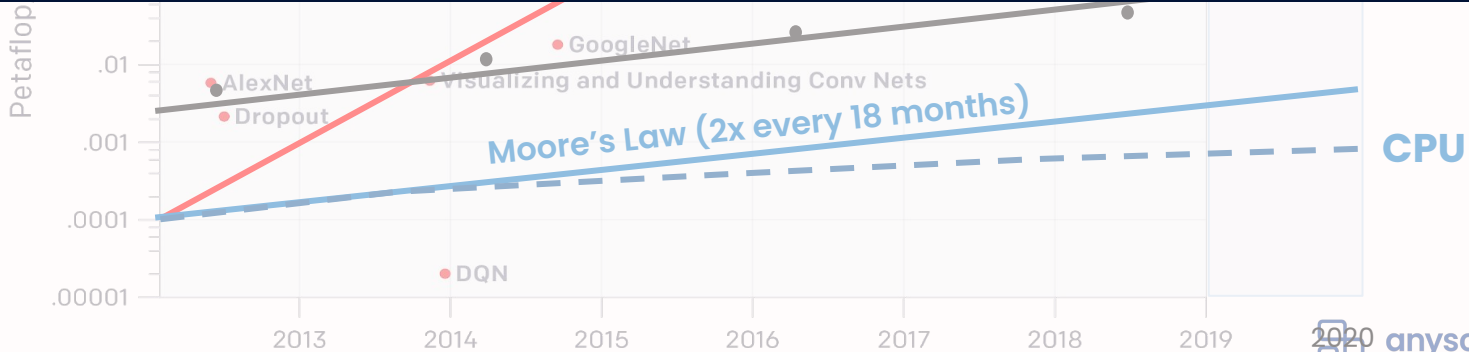
# Specialized hardware is also not enough



# Specialized hardware is also not enough



No way out but to distribute!




# Why Ray?

- Machine learning is pervasive in every domain
- Distributed machine learning is becoming a necessity
- **Distributed systems and programming are notoriously hard**

# Existing solutions have many tradeoffs

Ease of development

**Serverless**



**Limitations**

- 01. Cloud specific
- 02. Stateless only
- 03. No GPUs/TPUs
- 04. Runtime limit


**Stitch together existing frameworks**



**Hard to**

- 01. Develop
- 02. Deploy
- 03. Manage

**Clean slate**



**Expensive to develop**

- 01. Time
- 02. People

Generality



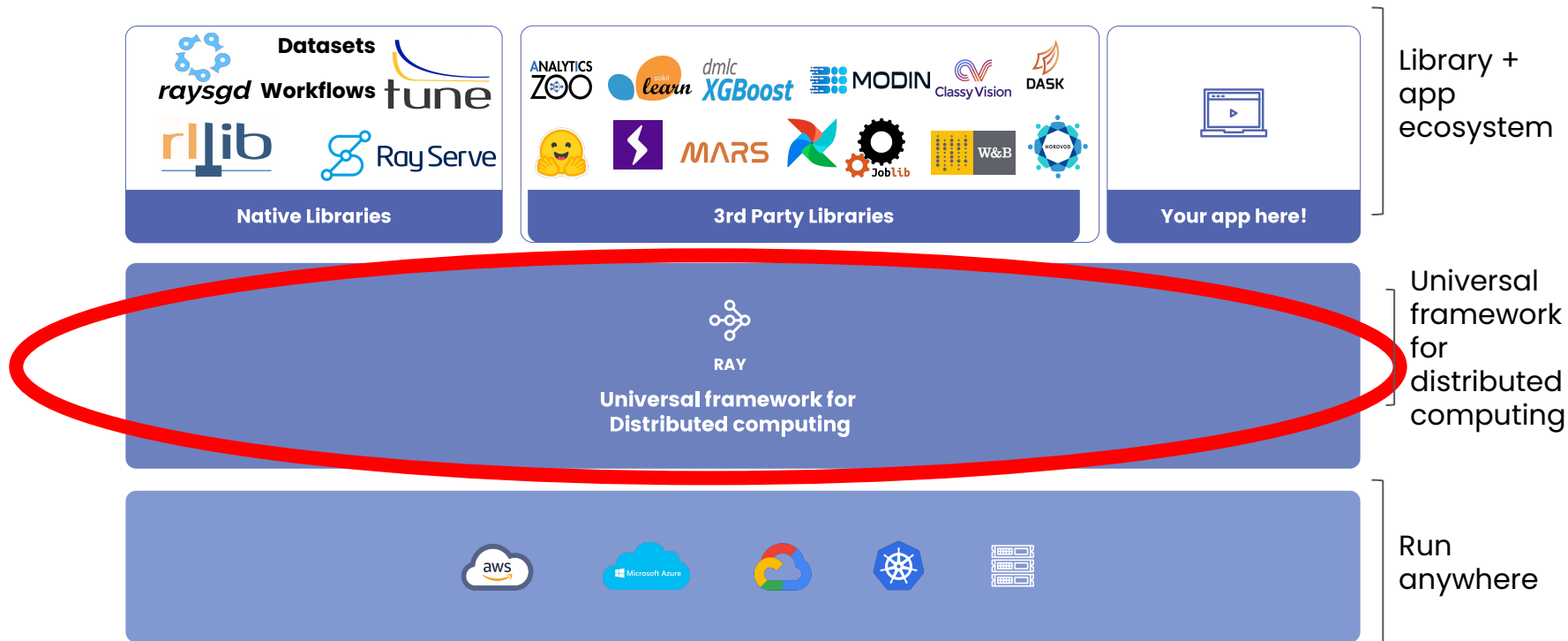
# Why Ray?

- Machine learning is pervasive in every domain
- Distributed machine learning is becoming a necessity
- Distributed systems are notoriously hard

## **Ray's vision:**

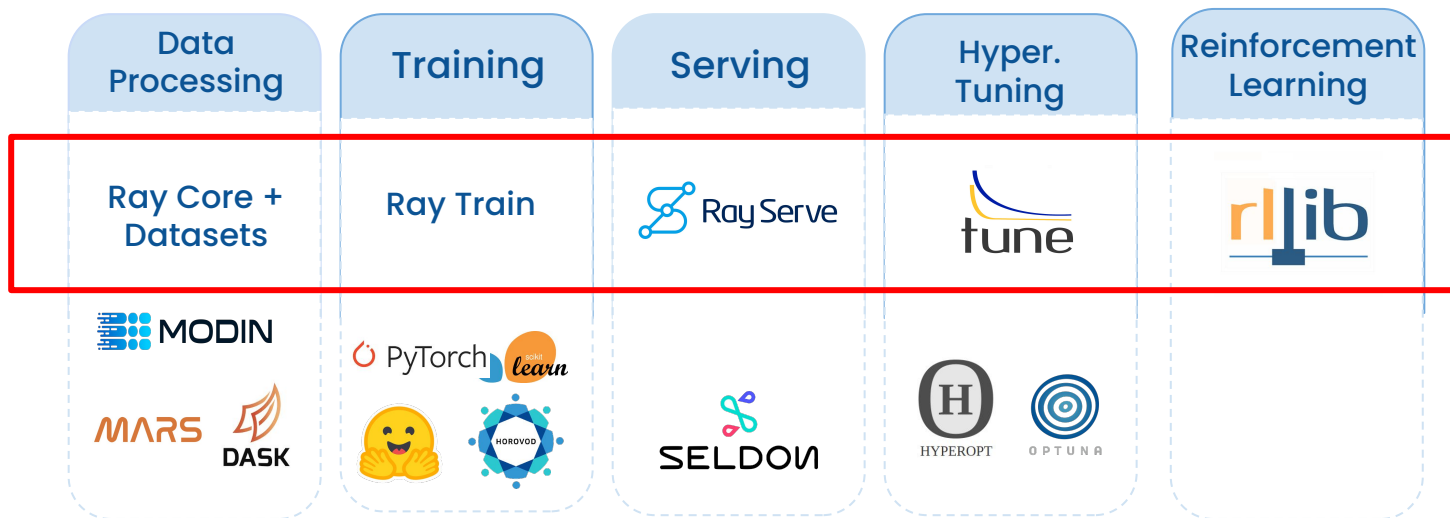
Make distributed computing accessible to every developer

# The Ray Layered Cake and Ecosystem



# Rich ecosystem for scaling ML workloads

Built-in  
“batteries  
included”  
libraries



Only use the libraries you need!

# Companies scaling ML with Ray

amazon



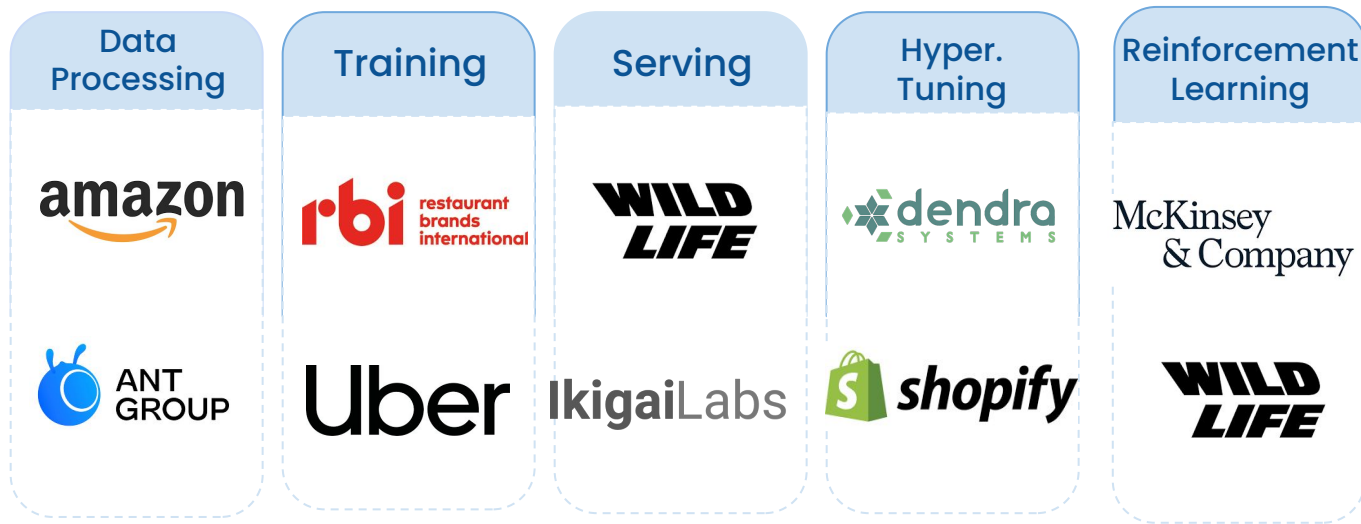
Uber



McKinsey  
& Company



# Companies scaling ML with Ray



- <https://eng.uber.com/horovod-ray/>
- <https://www.anyscale.com/blog/wildlife-studios-serves-in-game-offers-3x-faster-at-1-10th-the-cost-with-ray>
- <https://www.ikigailabs.com/blog/how-ikigai-labs-serves-interactive-ai-workflows-at-scale-using-ray-serve>

# Ray's approach for scaling ML

(Traditional, non-parallelized vanilla Python)



Runs on



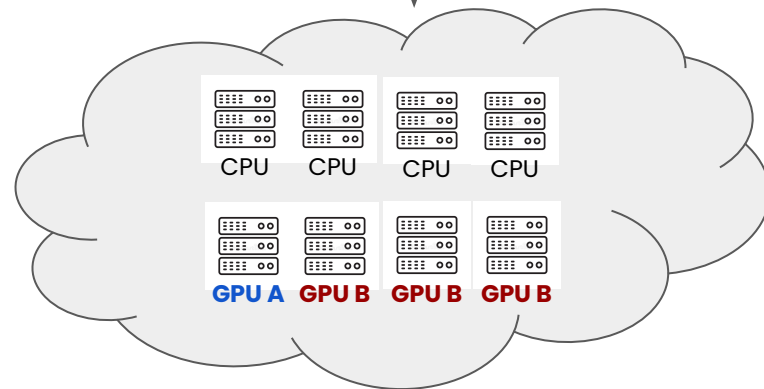
(with as few code changes as possible)

**Same Python code runs on laptop as infinite cloud!**



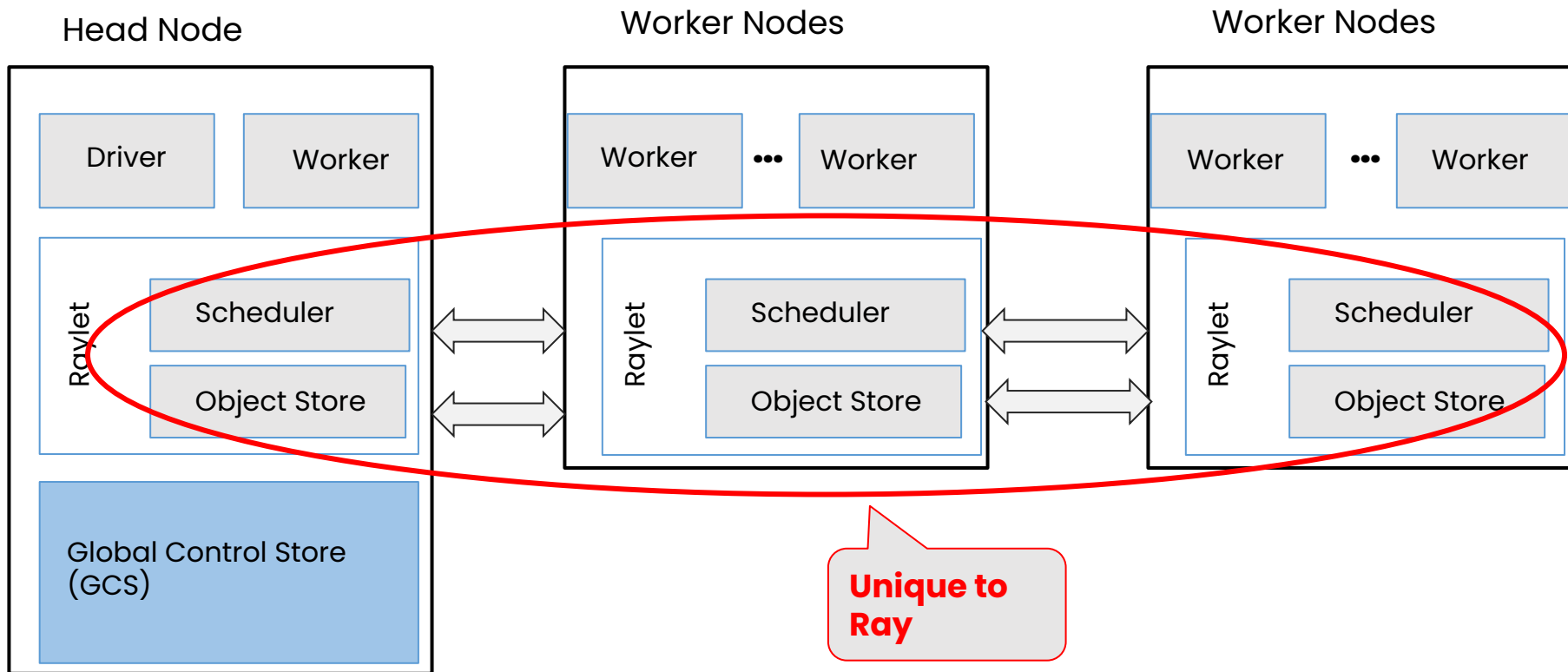
Ray-ified Python

Runs on

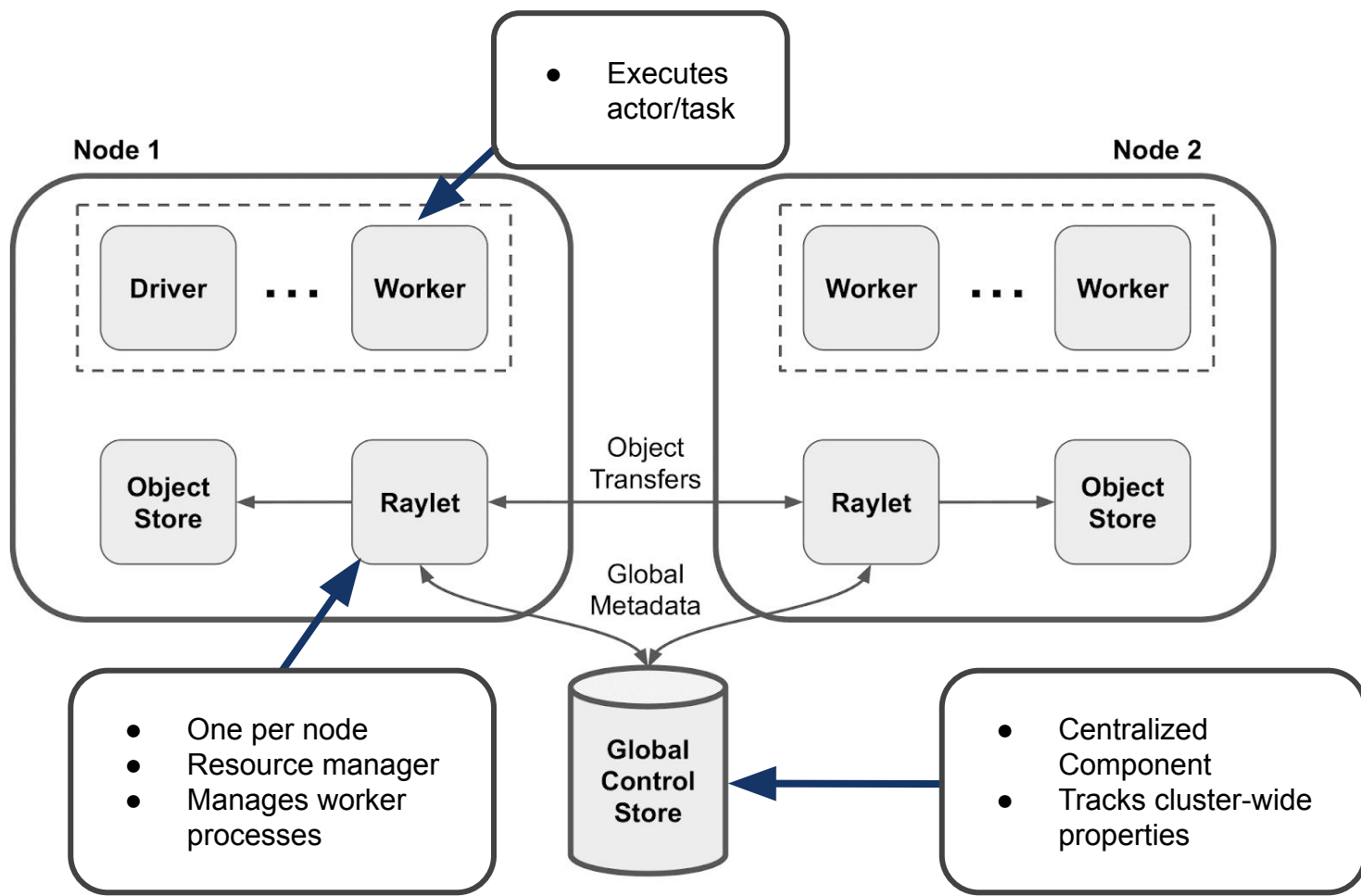


# Ray Architecture & Components

# What does Ray Cluster Looks Like ...







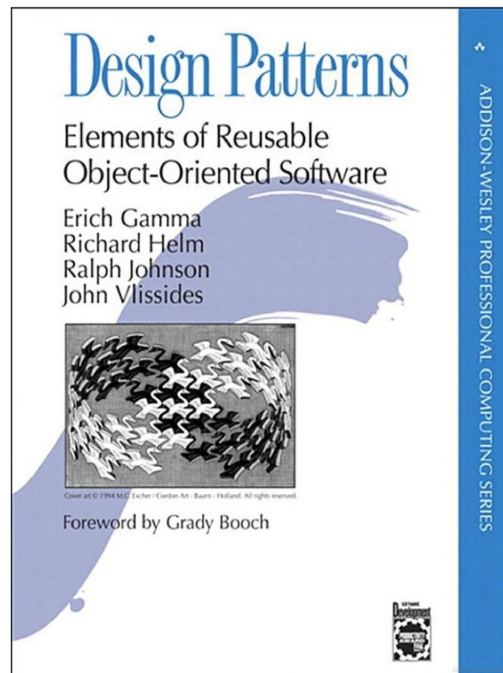
# Ray Distributed Design Patterns & APIs



# Ray Basic Design Patterns

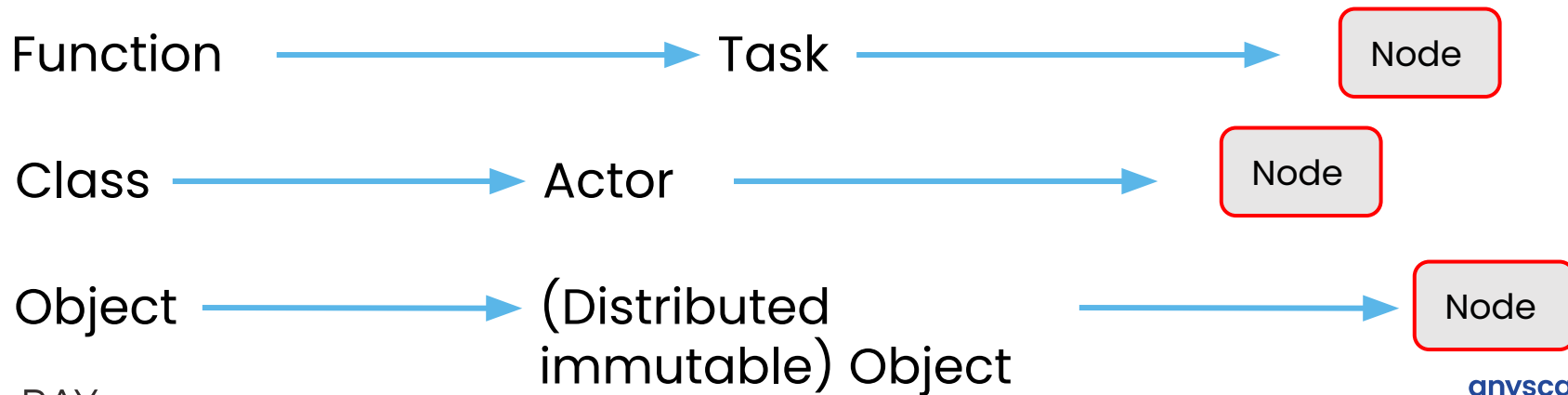
- **Ray Parallel Tasks**
  - Functions as stateless units of execution
  - Functions distributed across a clusters as tasks
- **Ray Objects or Futures**
  - Distributed (immutable) Object stored in cluster
  - Retrievable when available
  - Enable asynchronous execution of
- **Ray Actors**
  - Stateful service on a cluster
  - Message passing and maintains state

1. [Patterns for Parallel Programming](#)
2. [Ray Design Patterns](#)
3. [Ray Distributed Library Integration Patterns](#)





# Python → Ray Basic Patterns



## Function → Task

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a
```

```
@ray.remote
def add(a, b):
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)
id2 = read_array.remote(file2)
id = add.remote(id1, id2)
sum = ray.get(id)
```

## Class → Actor

```
@ray.remote(num_gpus=1)
class Counter(object):
    def __init__(self):
        self.value = 0
    def inc(self):
        self.value += 1
        return self.value
```

```
c = Counter.remote()
id4 = c.inc.remote()
id5 = c.inc.remote()
```

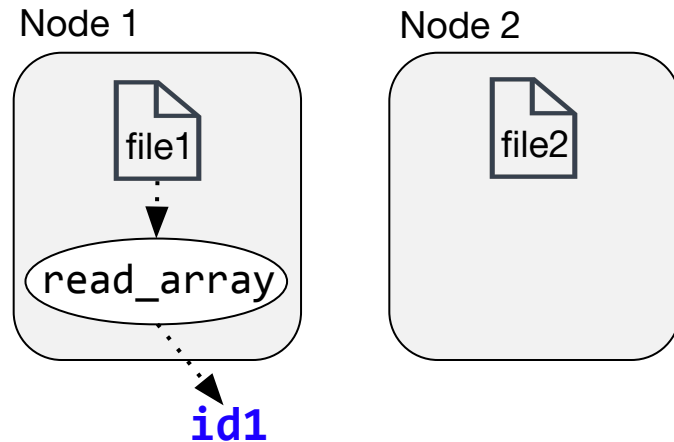
# Task API

Blue variables are `ObjectRef` IDs  
(similar to futures)

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a
```

```
@ray.remote
def add(a, b):
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)
id2 = read_array.remote(file2)
id = add.remote(id1, id2)
sum = ray.get(id)
```



Return `id1` (future) immediately,  
before `read_array()` finishes

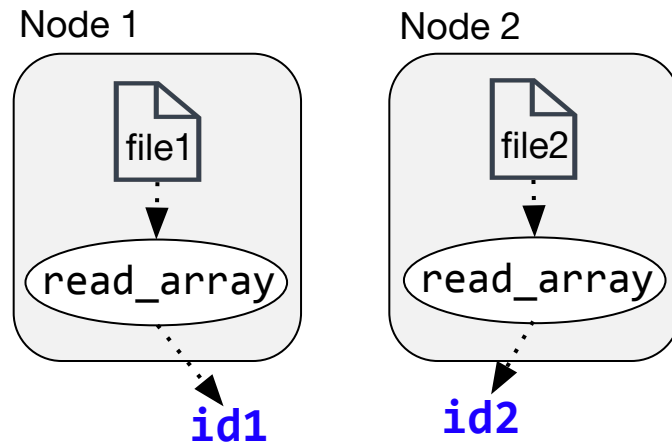
# Task API

```
@ray.remote
def read_array(file):
    # read ndarray "a"
    # from "file"
    return a
```

```
@ray.remote
def add(a, b):
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)
id2 = read_array.remote(file2)
id = add.remote(id1, id2)
sum = ray.get(id)
```

Blue variables are Object IDs  
(similar to futures)



Dynamic task graph:  
build at runtime

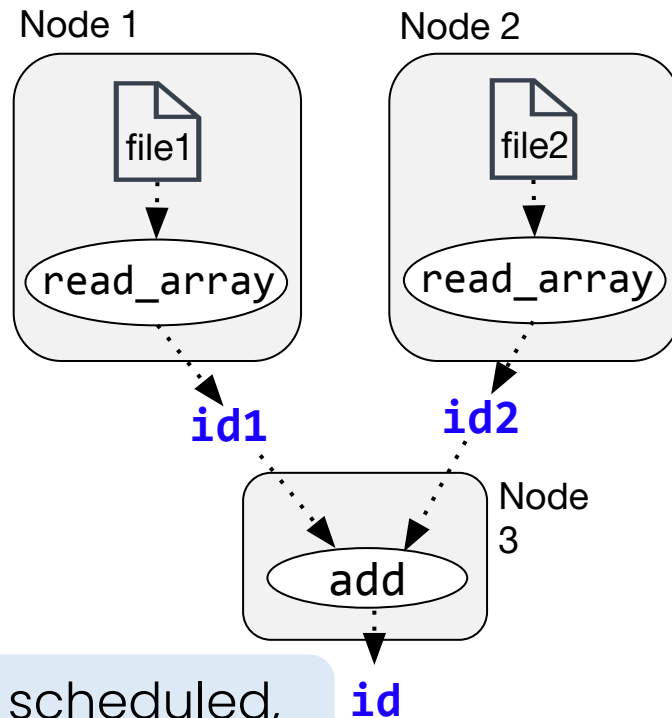
# Task API

```
@ray.remote  
def read_array(file):  
    # read ndarray "a"  
    # from "file"  
    return a
```

```
@ray.remote  
def add(a, b):  
    return np.add(a, b)
```

```
id1 = read_array.remote(file1)  
id2 = read_array.remote(file2)  
id = add.remote(id1, id2)  
sum = ray.get(id)
```

Blue variables are Object IDs  
(similar to futures)



Every task scheduled,  
but not finished yet



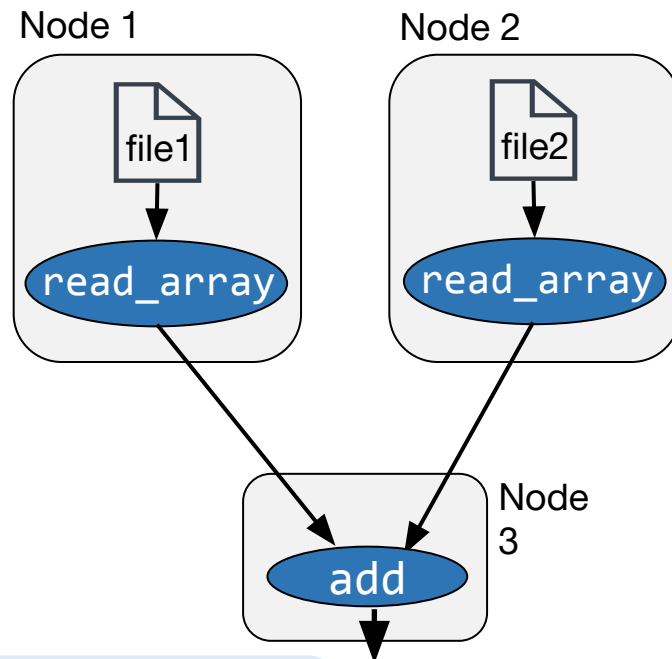
# Task API

```
@ray.remote  
def read_array(file):  
    # read ndarray "a"  
    # from "file"  
    return a
```

```
@ray.remote  
def add(a, b):  
    return np.add(a, b)
```

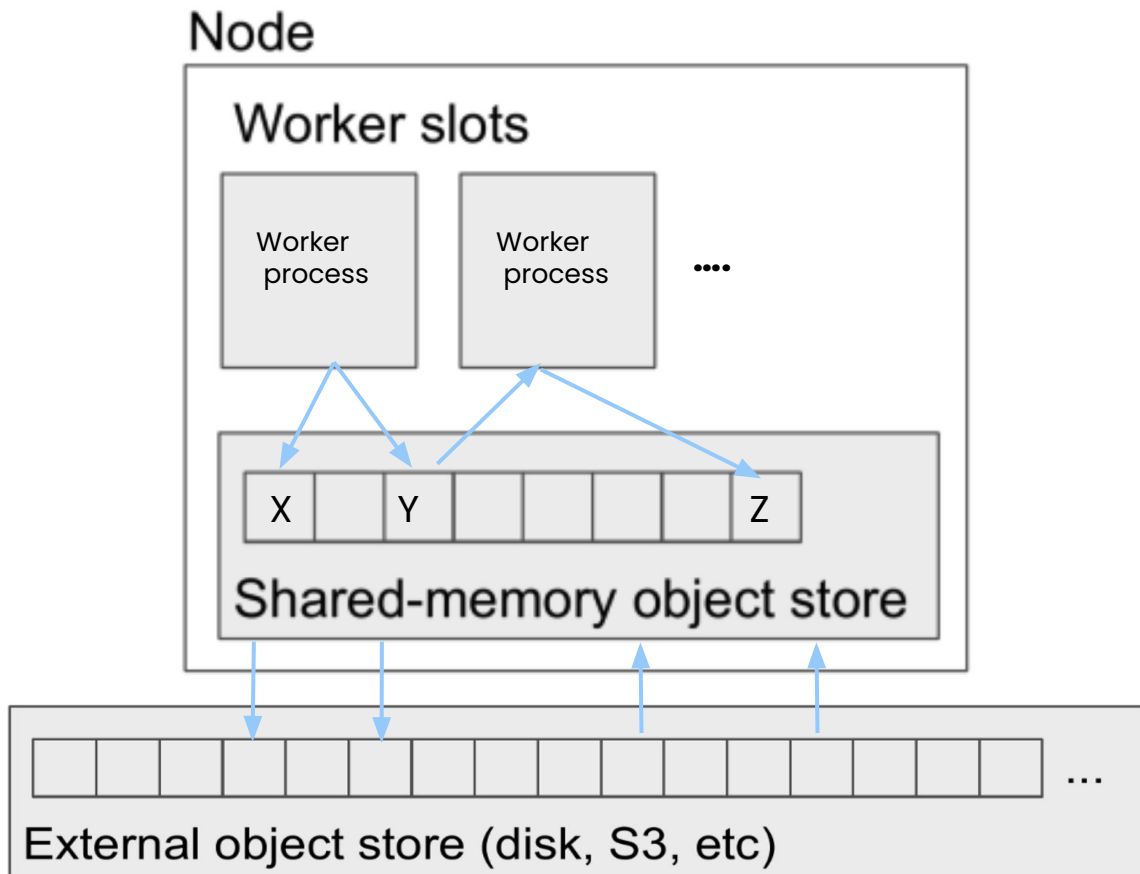
```
id1 = read_array.remote(file1)  
id2 = read_array.remote(file2)  
id = add.remote(id1, id2)  
sum = ray.get(id)
```

Blue variables are Object IDs  
(similar to futures)



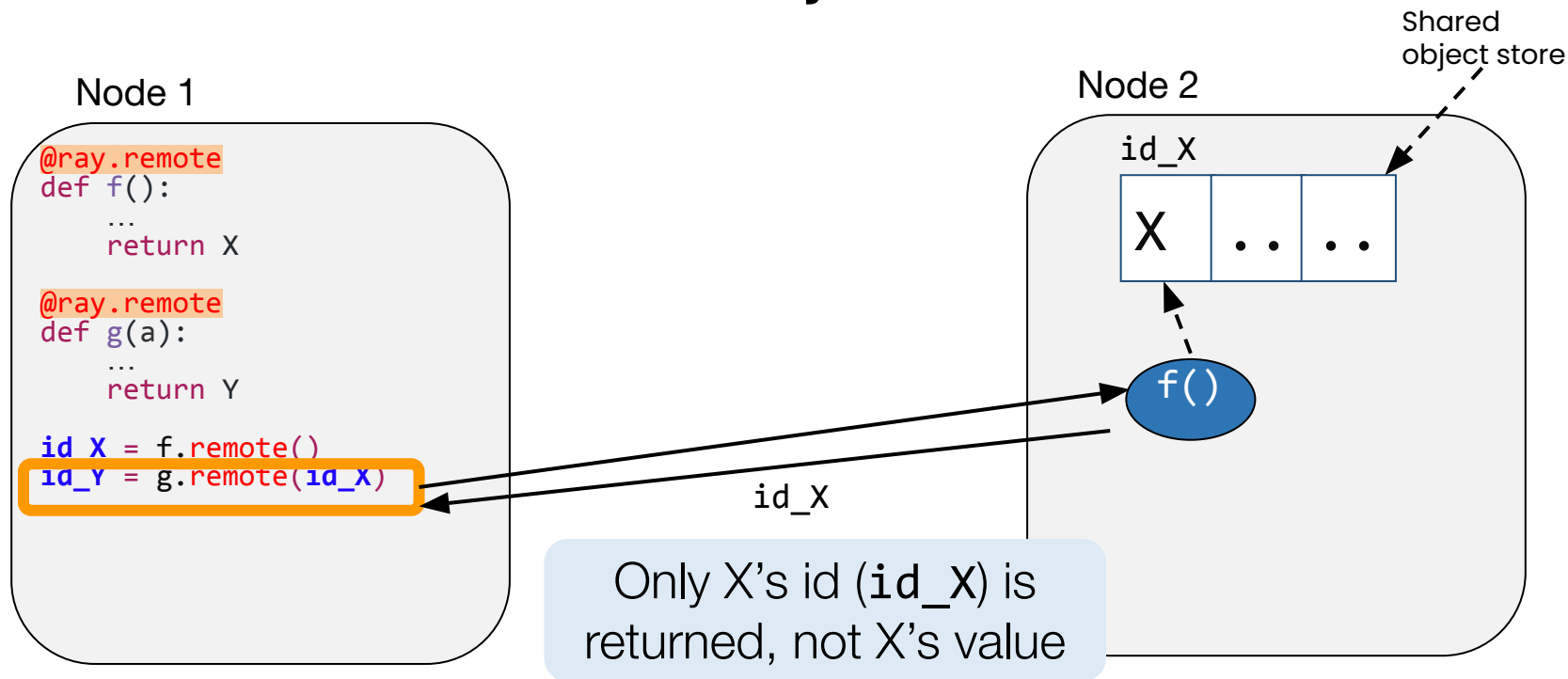
Task graph executed to compute sum

# Distributed Immutable object store

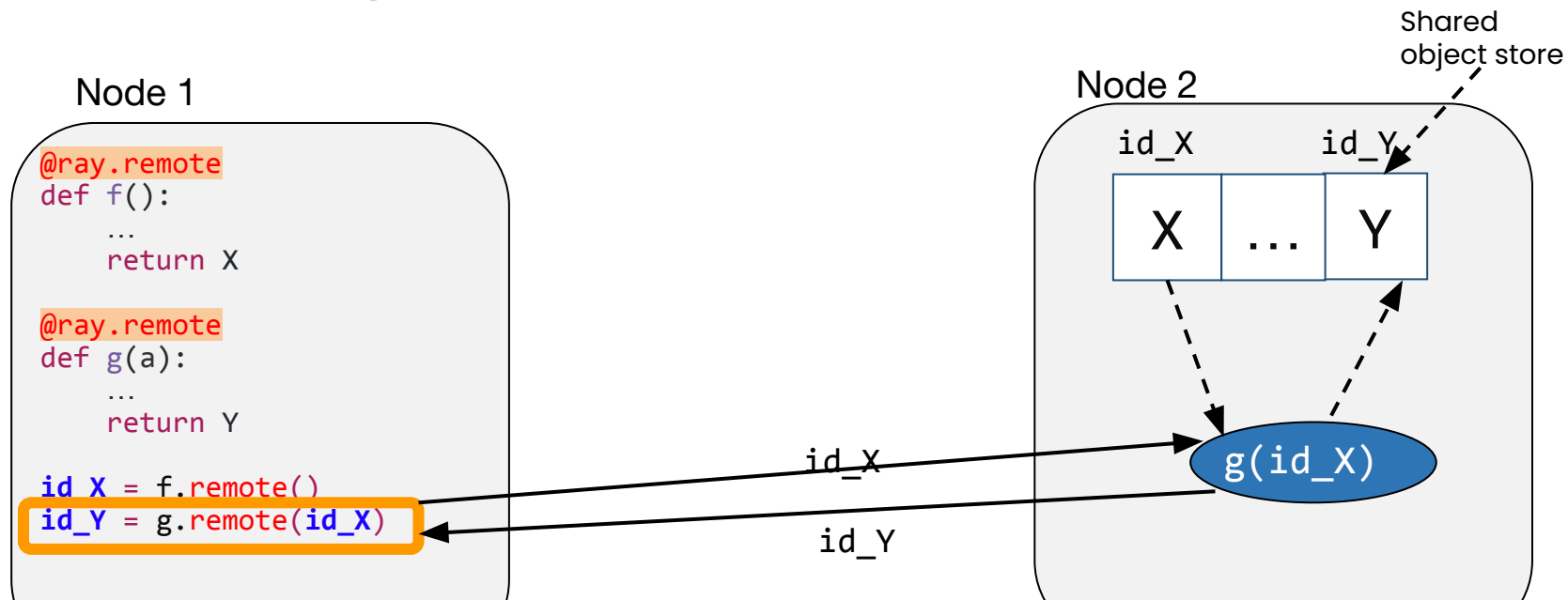


Spill over to  
external storage

# Distributed Immutable object store



# Distributed object store



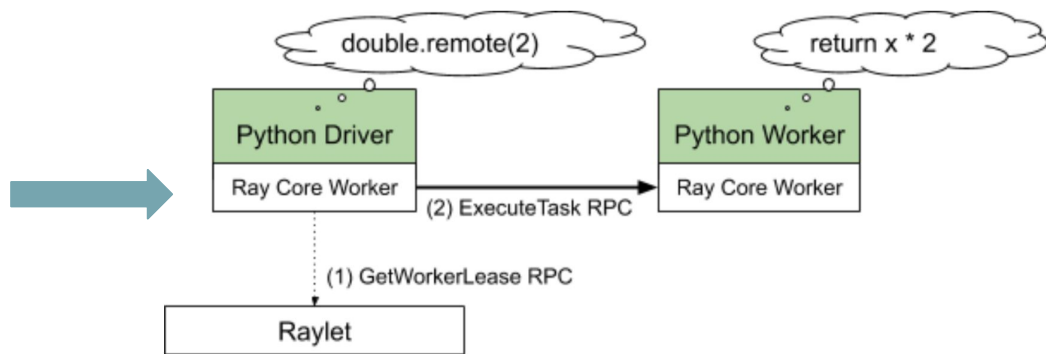
`g(id_X)` is scheduled on same node, so X is never transferred

# How Raylet Schedules Tasks

## Basic Ray Task Call

```
@ray.remote
def double(x):
    return x * 2

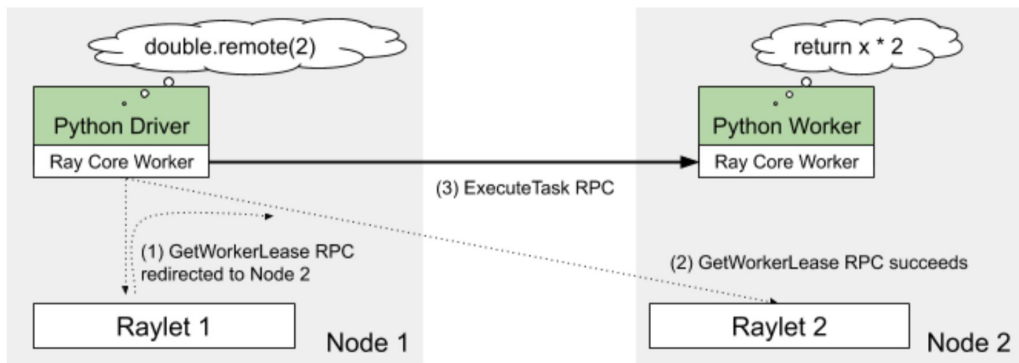
fut1 = double.remote(2)
assert ray.get(fut1) == 4
```



Components in green boxes represent Python code. Components in white boxes are part of the Ray common runtime written in C++. Joined boxes represent a process. Any Python driver or worker can call into the Ray C++ core worker library to execute further tasks. In this figure, all processes are running on the same machine. Ray uses gRPC as a unified communication layer for both local and remote procedure calls.

# Scaling to Multiple Nodes

1. The driver asks Raylet 1 for a worker to execute `double`. It has no free workers, but Raylet 1 knows Raylet 2 has free resources, and redirects the request to Raylet 2.
2. The driver sends `ExecuteTask` to the remote Python worker leased from Raylet 2 over gRPC.

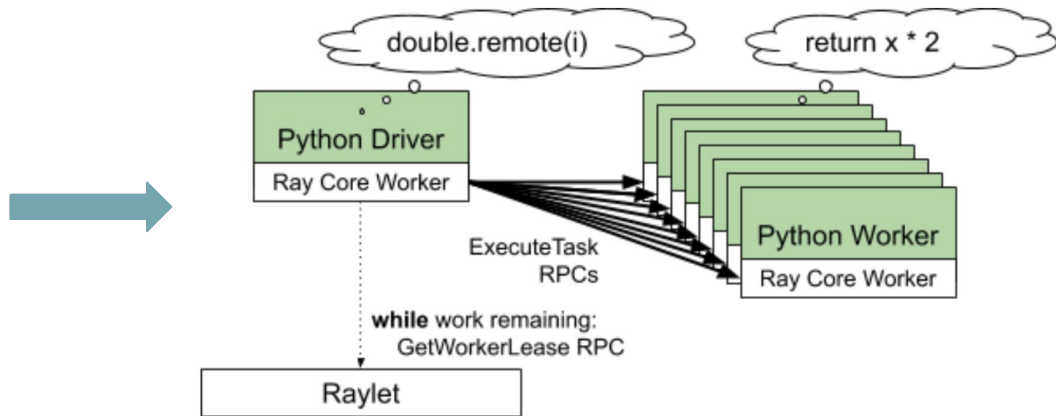


Tasks are sent to remote workers if there are no local resources available, transparently scaling Ray applications out to multiple nodes.

# Caching Scheduling Decisions

```
futures = [double.remote(i)
for i in range(10000)]
ray.get(futures)
```

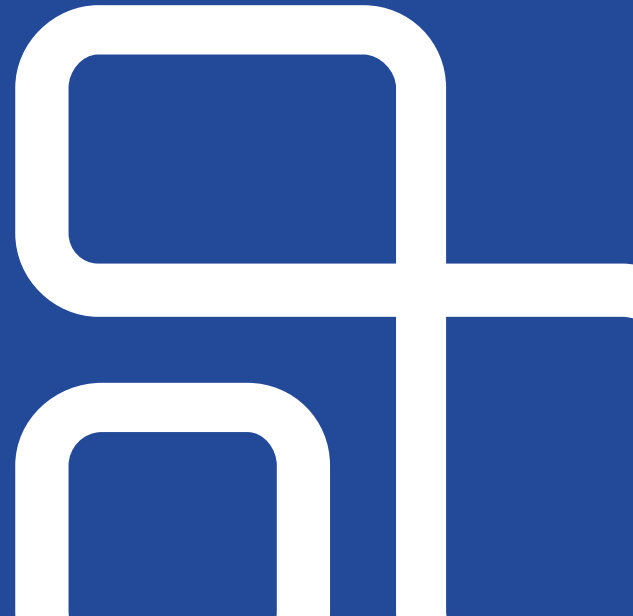
```
# [0, 2, 4, 6 ...]
```



Once a scheduling decision is made by the Raylet, the worker returned can be reused for other tasks with the same resource requirements and input dependencies. This amortizes scheduling RPC overhead when executing many similar tasks. To avoid unfair monopolization of workers when there are multiple processes trying to submit tasks, callers are only allowed to reuse workers within a few hundred milliseconds of initial grant.

# Ray Ecosystem

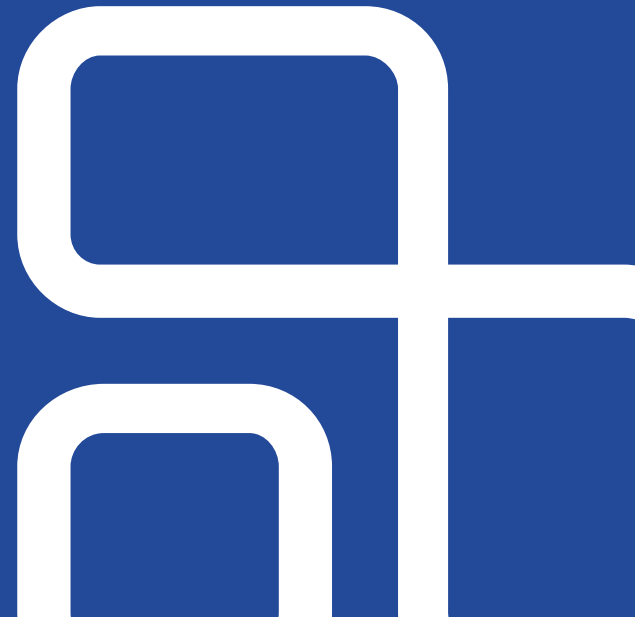
- Ray Tune
- XGBoost-Ray





# Ray Tune

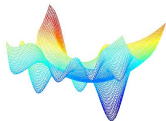
---



# Ray Tune – For distributed HPO

- Efficient algorithms that enable running trials in parallel
- Effective orchestration of distributed trials
- Easy to use APIs
- Interoperable with Ray Train and Ray Datasets
- Saves cost (early stopping bad trials)

Cutting edge  
optimization algorithms



Compatible with ML  
ecosystem



```
tune.run(train_model)
```

Minimal code changes to  
work in distributed  
settings

Single Process

Multi-process/  
Multi-GPU

Multi-Node



# Ray tune

## Search Algorithms (tune.suggest)

Tune's Search Algorithms are wrappers around open-source optimization libraries for efficient hyperparameter selection. Each library has a specific way of defining the search space - please refer to their documentation for more details.

You can utilize these search algorithms as follows:

```
from ray.tune.suggest.hyperopt import HyperOptSearch
tune.run(my_function, search_alg=HyperOptSearch(...))
```

### Summary

SearchAlgorithm	Summary	Website	Code Example
Random search/grid search	Random search/grid search		<a href="#">tune_basic_example</a>
AxSearch	Bayesian/Bandit Optimization	<a href="#">[Ax]</a>	<a href="#">ax_example</a>
BlendSearch	Blended Search	<a href="#">[Bs]</a>	<a href="#">blendsearch_example</a>
CFO	Cost-Frugal hyperparameter Optimization	<a href="#">[Cfo]</a>	<a href="#">cfo_example</a>
DragonflySearch	Scalable Bayesian Optimization	<a href="#">[Dragonfly]</a>	<a href="#">dragonfly_example</a>
SkoptSearch	Bayesian Optimization	<a href="#">[Scikit-Optimize]</a>	<a href="#">skopt_example</a>
HyperOptSearch	Tree-Parzen Estimators	<a href="#">[HyperOpt]</a>	<a href="#">hyperopt_example</a>
BayesOptSearch	Bayesian Optimization	<a href="#">[BayesianOptimization]</a>	<a href="#">bayesopt_example</a>
TuneBOHB	Bayesian Opt/HyperBand	<a href="#">[BOHB]</a>	<a href="#">bohb_example</a>
NevergradSearch	Gradient-free Optimization	<a href="#">[Nevergrad]</a>	<a href="#">nevergrad_example</a>
OptunaSearch	Optuna search algorithms	<a href="#">[Optuna]</a>	<a href="#">optuna_example</a>
ZOOptSearch	Zeroth-order Optimization	<a href="#">[ZOOpt]</a>	<a href="#">zoopt_example</a>
SigOptSearch	Closed source	<a href="#">[SigOpt]</a>	<a href="#">sigopt_example</a>
HEBOSearch	Heteroscedastic Evolutionary Bayesian Optimization	<a href="#">[HEBO]</a>	<a href="#">hebo_example</a>

[https://docs.ray.io/en/latest/tune/api\\_docs/suggestion.html#tune-search-alg](https://docs.ray.io/en/latest/tune/api_docs/suggestion.html#tune-search-alg)

## Trial Schedulers (tune.schedulers)

In Tune, some hyperparameter optimization algorithms are written as "scheduling algorithms". These Trial Schedulers can early terminate bad trials, pause trials, clone trials, and alter hyperparameters of a running trial.

All Trial Schedulers take in a `metric`, which is a value returned in the result dict of your Trainable and is maximized or minimized according to `mode`.

```
tune.run(..., scheduler=Scheduler(metric="accuracy", mode="max"))
```

### Summary

Tune includes distributed implementations of early stopping algorithms such as [Median Stopping Rule](#), [HyperBand](#), and [ASHA](#). Tune also includes a distributed implementation of [Population Based Training \(PBT\)](#) and [Population Based Bandits \(PB2\)](#).

#### Tip

The easiest scheduler to start with is the [ASHAScheduler](#) which will aggressively terminate low-performing trials.

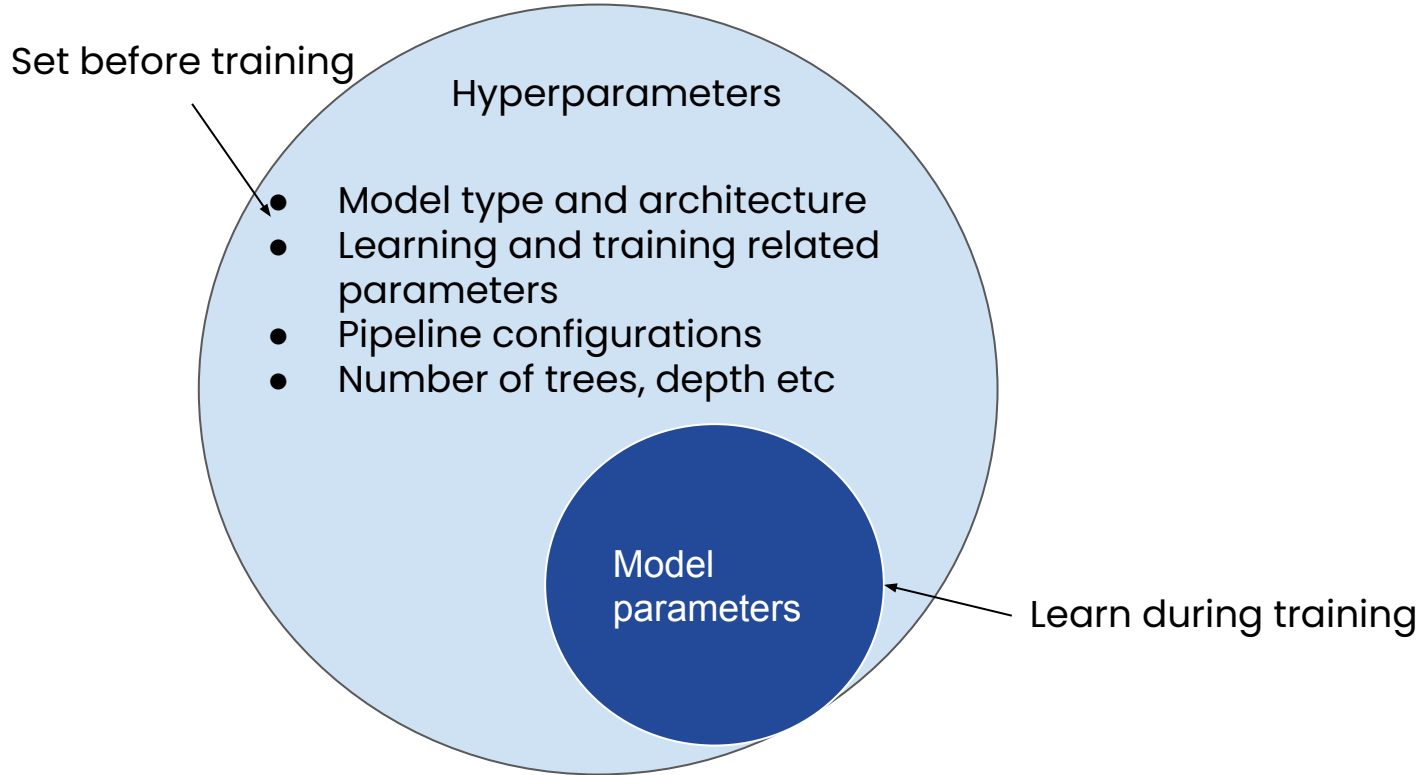
When using schedulers, you may face compatibility issues, as shown in the below compatibility matrix. Certain schedulers cannot be used with Search Algorithms, and certain schedulers are required to be implemented.

Schedulers can dynamically change trial resource requirements during tuning. This is currently implemented in [ResourceChangingScheduler](#), which can wrap around any other scheduler.

Scheduler	Need Checkpointing?	SearchAlg Compatible?	Example
ASHA	No	Yes	<a href="#">Link</a>
Median Stopping Rule	No	Yes	<a href="#">Link</a>
HyperBand	Yes	Yes	<a href="#">Link</a>
BOHB	Yes	Only TuneBOHB	<a href="#">Link</a>
Population Based Training	Yes	Not Compatible	<a href="#">Link</a>
Population Based Bandits	Yes	Not Compatible	<a href="#">Basic Example, PPO example</a>

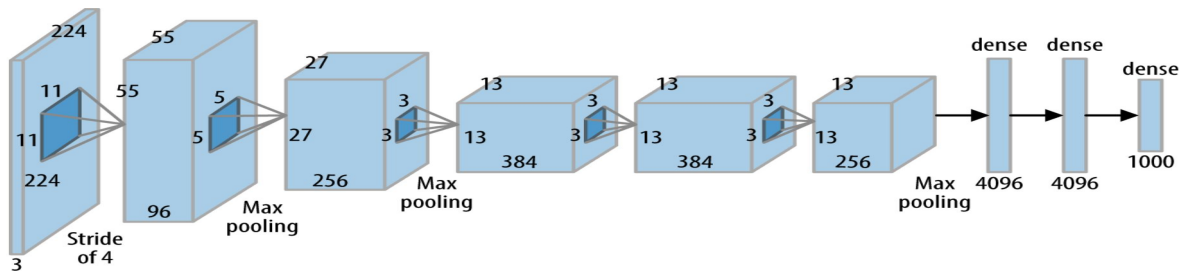
[https://docs.ray.io/en/latest/tune/api\\_docs/schedulers.html#tune-schedulers](https://docs.ray.io/en/latest/tune/api_docs/schedulers.html#tune-schedulers)

# Hyperparameters



# Hyperparameter tuning

“choosing a set of optimal hyperparameters for a learning algorithm”



**Example:** what network structure is best for your binary classification problem?

How many layers? What kinds of layers? Learning rate schedule?

Every number here is a hyperparameter!

# HPO Challenges at scale

- Time consuming and costly
- Use Resources (GPUs/CPU) at lower costs
- Fault-tolerance and elasticity



+



# Ray Tune - HPO algorithms

- Over 15+ algorithms natively provided or integrated
- Easy to swap out different algorithms with no code change

---

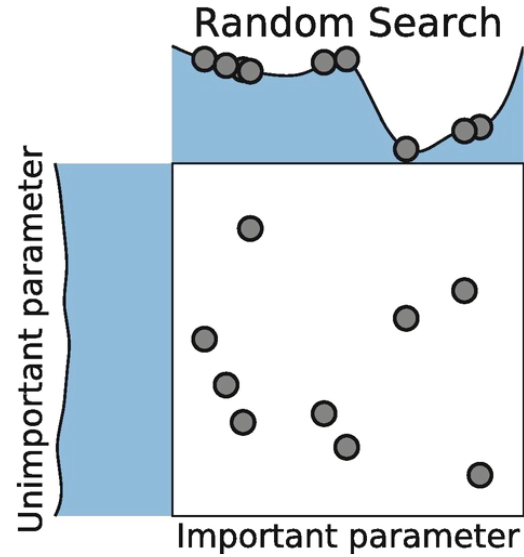
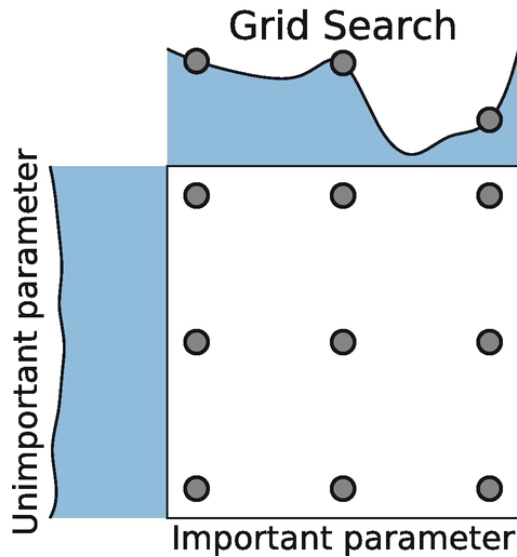
01 Exhaustive  
Search

02 Bayesian  
Optimization

03 Advanced  
Scheduling

# Exhaustive Search

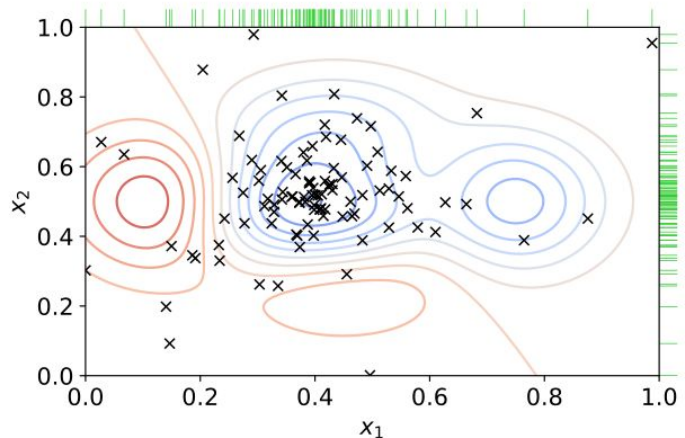
- Easily parallelizable, easy to implement
- Inefficient, compute intensive





# Bayesian optimization

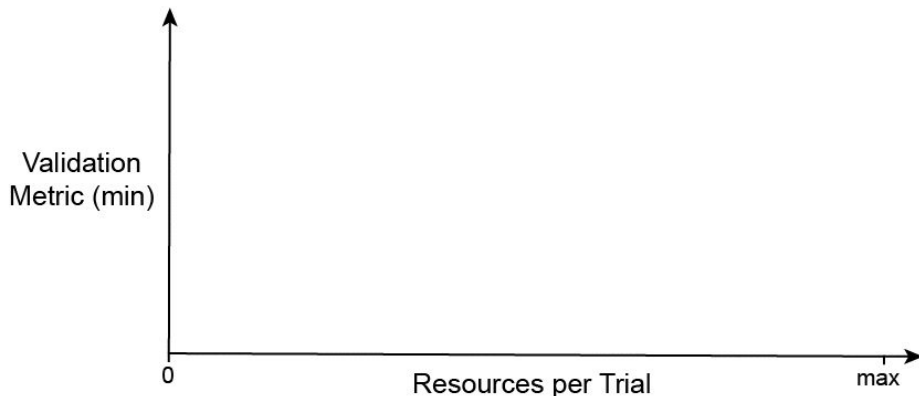
- Uses results from previous combinations (trials) to decide which trial to try next
- Inherently sequential
- Popular libraries:
  - hyperopt
  - Optuna
  - Scikit-optimize
  - Nevergrad



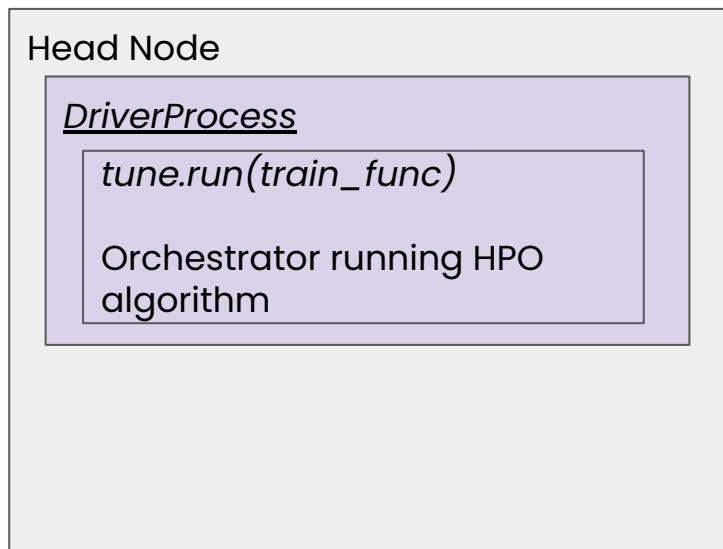
[https://www.wikiwand.com/en/Hyperparameter\\_optimization](https://www.wikiwand.com/en/Hyperparameter_optimization)

# Advanced Scheduling – Early stopping

- Fan out parallel trials during the initial exploration phase
- Use intermediate results (epochs, trees, samples) to prune underperforming trials, saving time and computing resources
- Median stopping, ASHA/Hyperband
- Can be combined with Bayesian Optimization (BOHB)



# Ray Tune - *distributed* HPO



```
from ray import tune
```

Easily define your training function

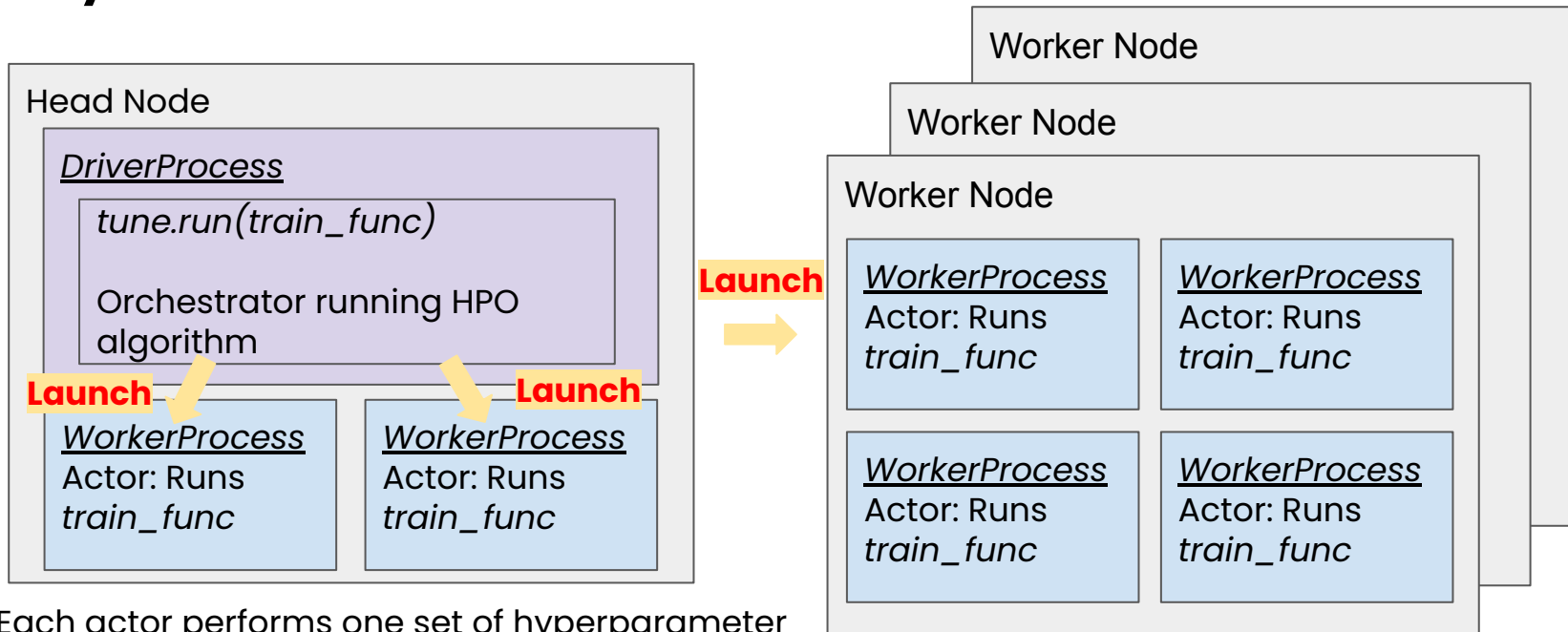
```
def train_func(config):  
    model = ConvNet(config)  
    for i in range(epochs):  
        current_loss = model.train()  
        tune.report(loss=current_loss)
```

```
tune.run(  
    train_func,  
    config={"alpha": tune.uniform(0.001,  
0.1)},  
    num_samples=100,  
    scheduler="asha",  
    search_alg="optuna")
```

Just use `tune.run(..)`

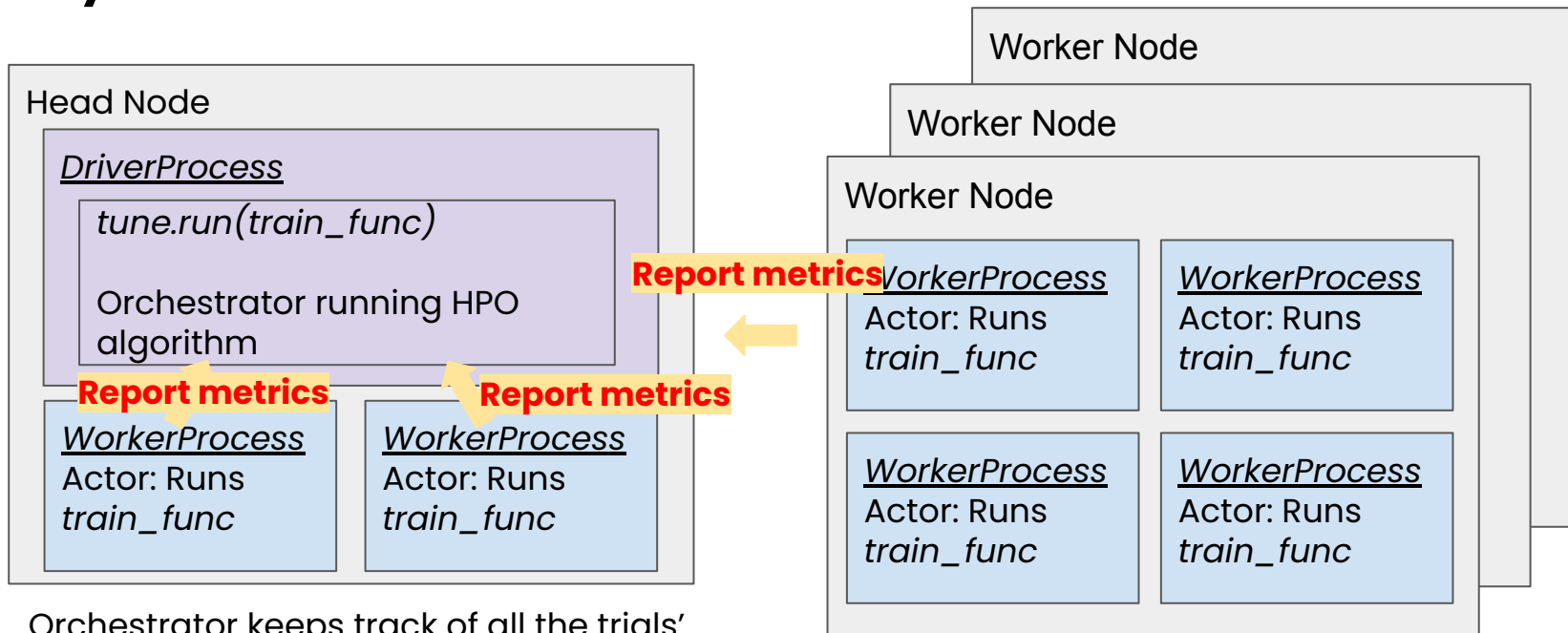
Easily specify hyperparameter ranges to search over

# Ray Tune - *distributed* HPO



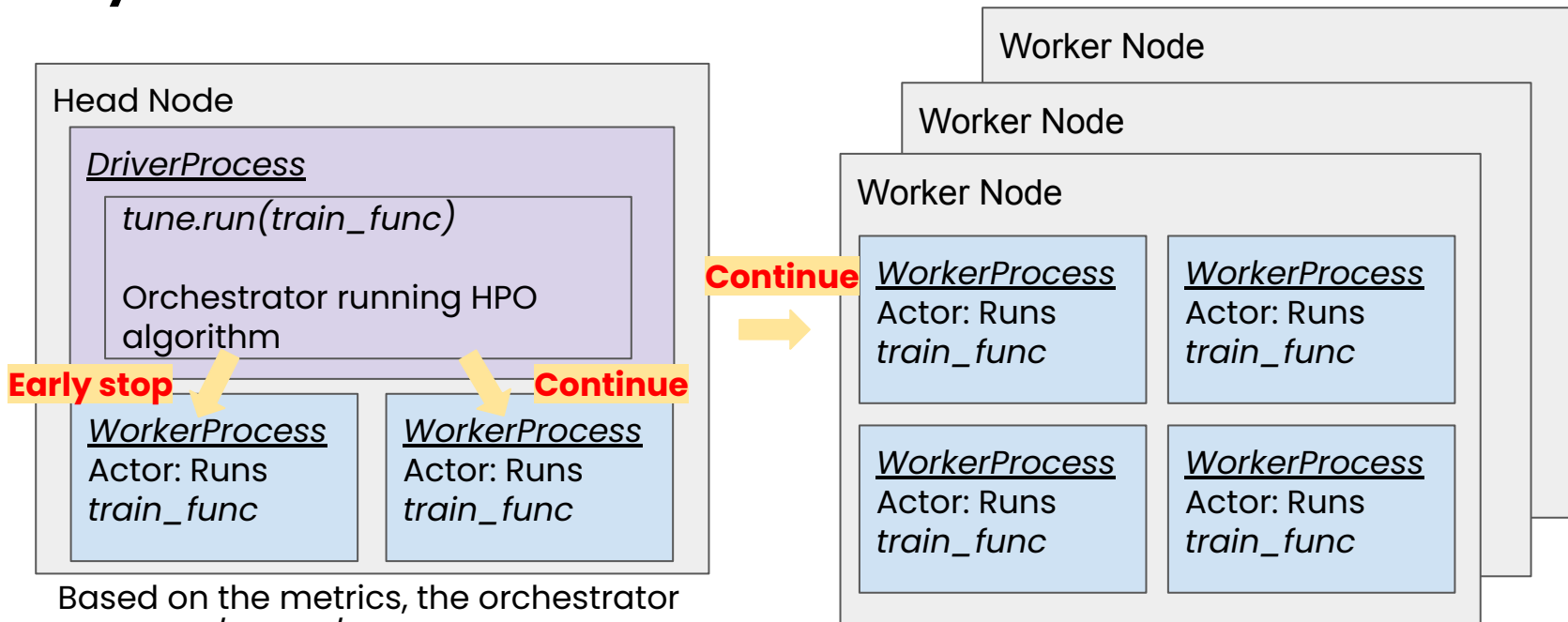
Each actor performs one set of hyperparameter combination evaluation (a trial)

# Ray Tune - *distributed* HPO



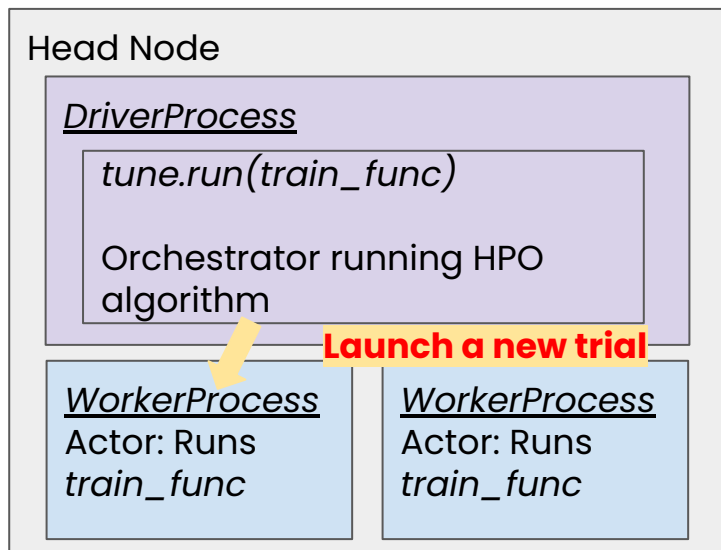
Orchestrator keeps track of all the trials' progress and metrics.

# Ray Tune - *distributed* HPO

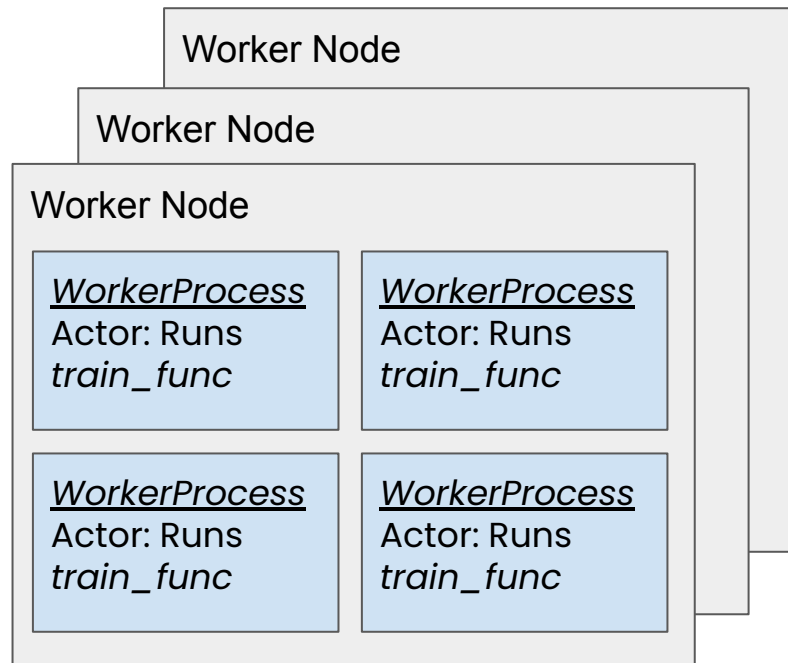


Based on the metrics, the orchestrator may stop/pause/mutate trials or launch new trials when resources are available.

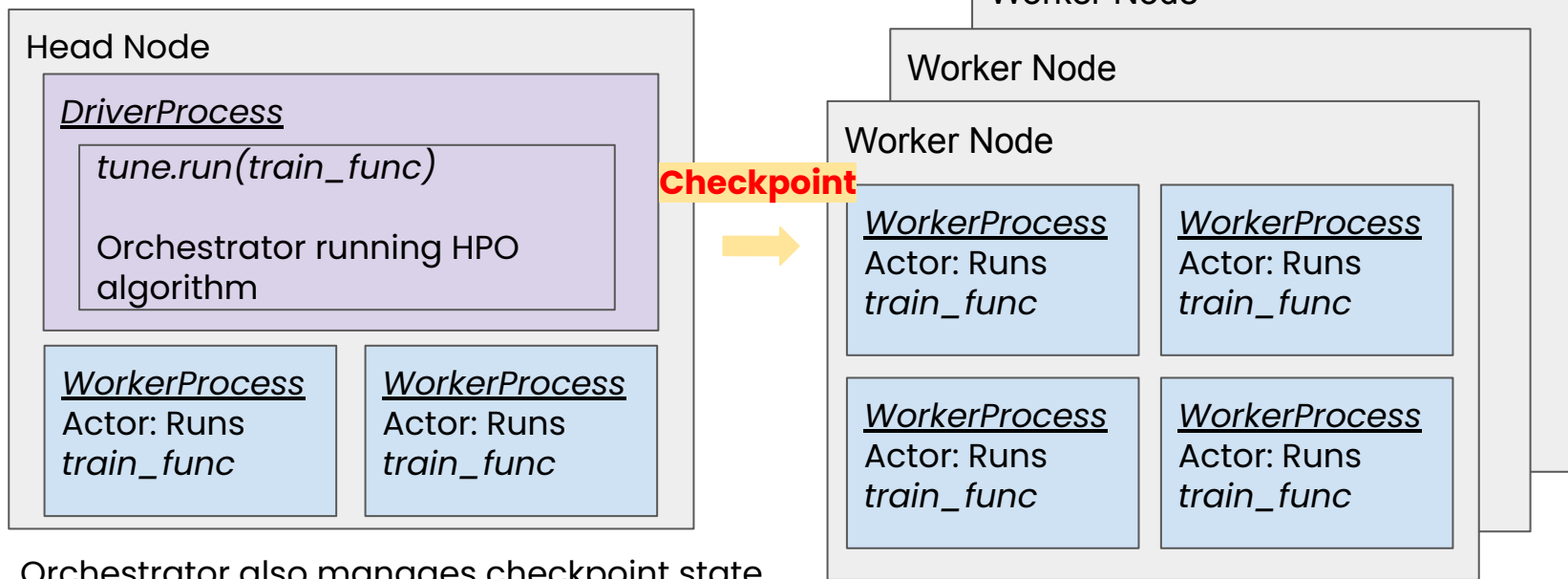
# Ray Tune - *distributed* HPO



Resources are repurposed to explore new trials.



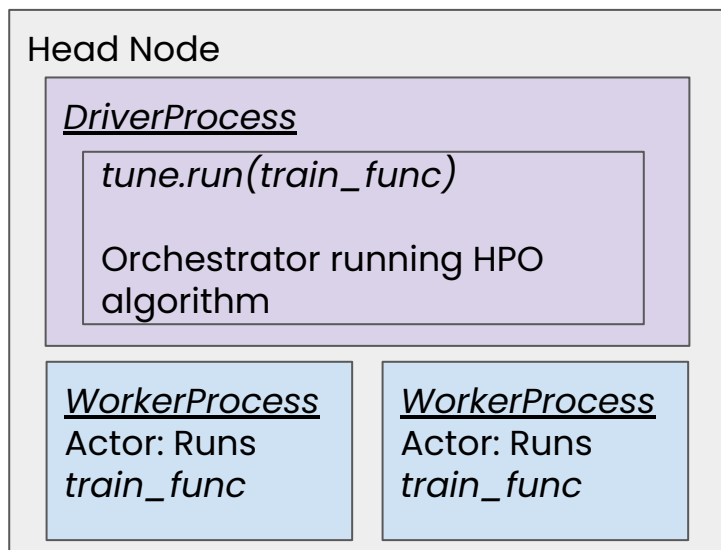
# Ray Tune - *distributed* HPO



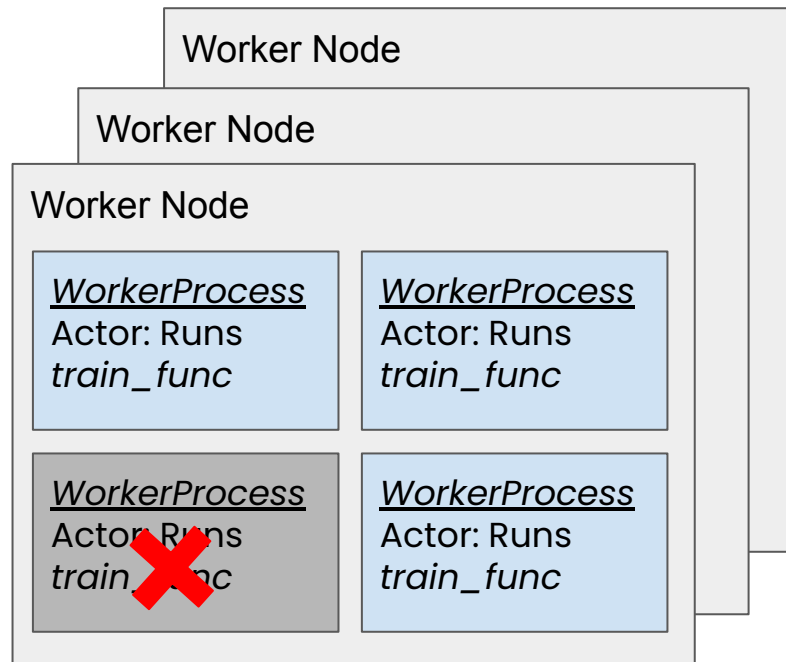
Orchestrator also manages checkpoint state.



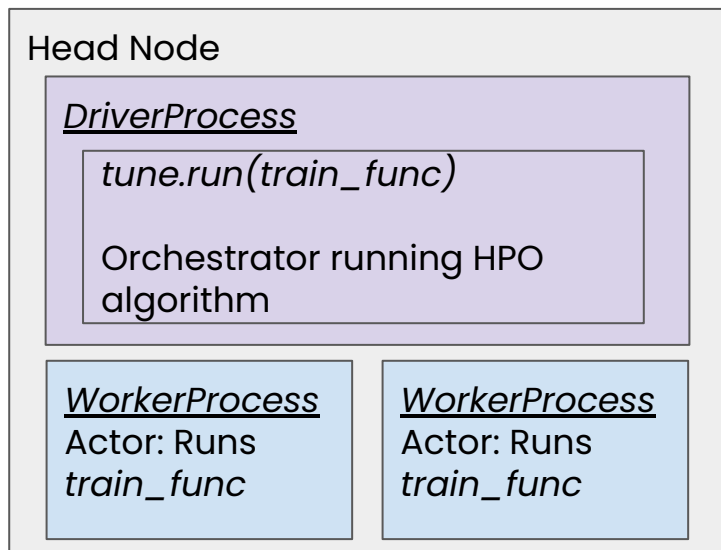
# Ray Tune - *distributed* HPO



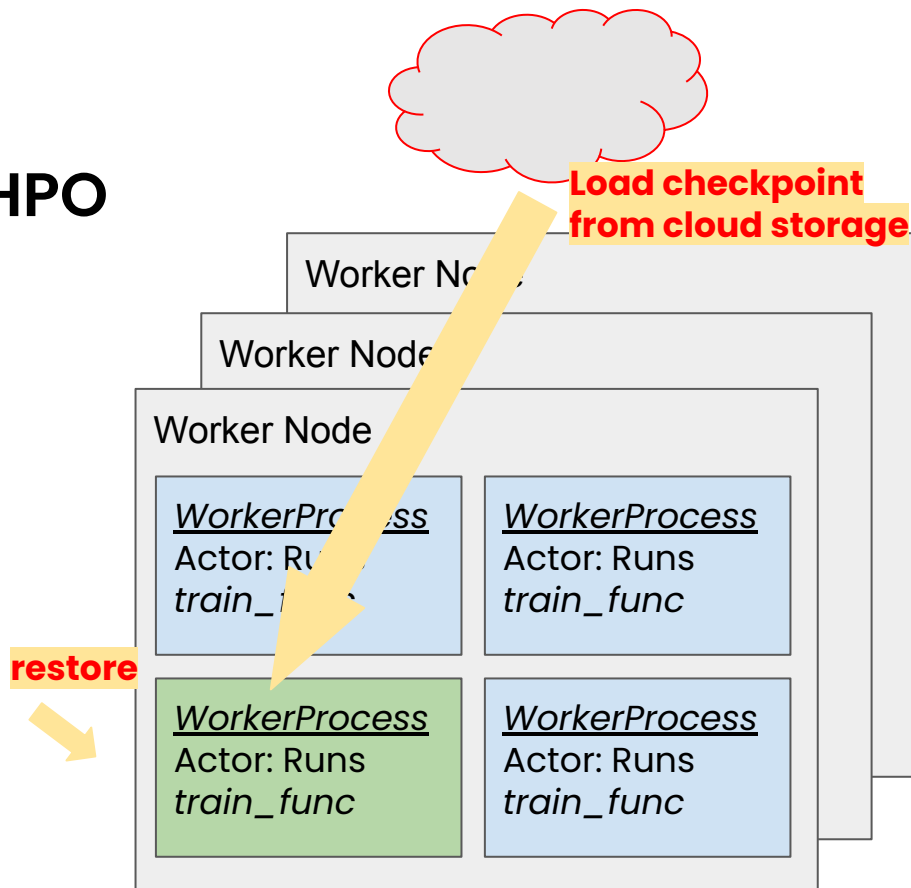
Some worker process crashes.



# Ray Tune - *distributed* HPO

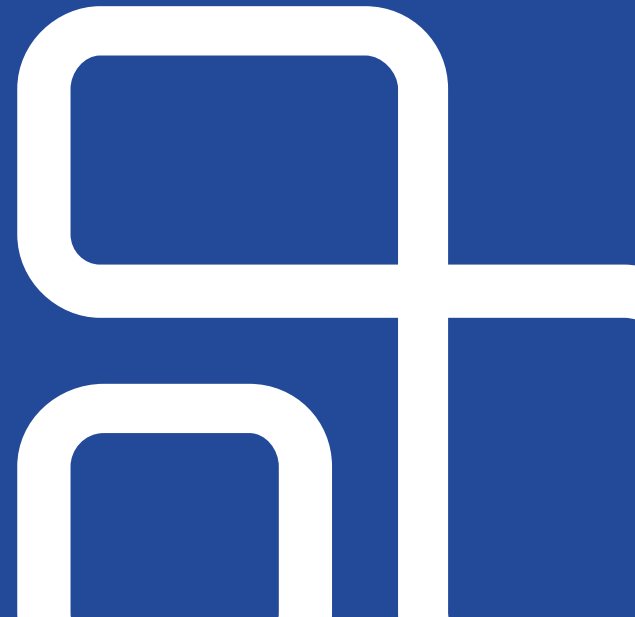


New actor comes up fresh and the crashed trial is restored from remote checkpoint.



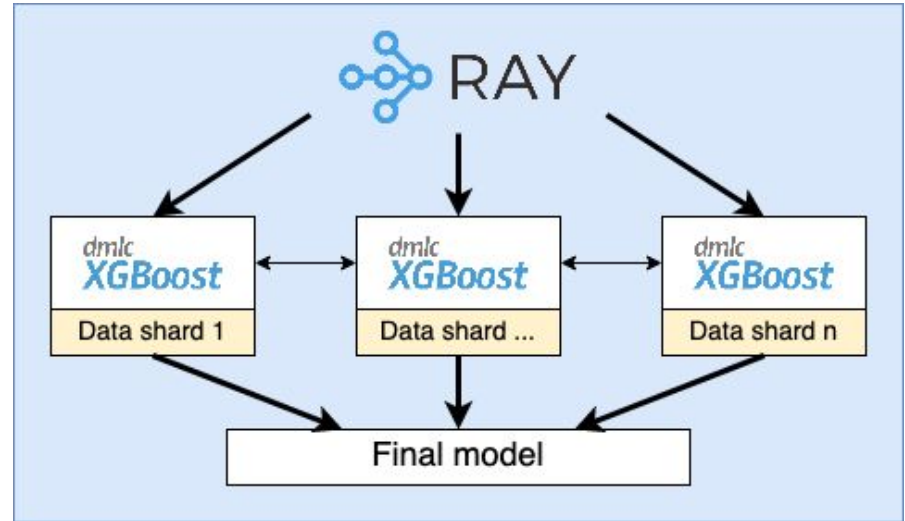
# XGBoost-Ray

- Design & Features
- 



# XGBoost-Ray

- Distributed XGBoost-Ray – Drop-in replacement for XGBoost
- Fault tolerance & Elastic training
- Integration with Ray Datasets and Ray Tune



- [https://github.com/ray-project/xgboost\\_ray](https://github.com/ray-project/xgboost_ray)
- <https://docs.ray.io/en/latest/xgboost-ray.html>

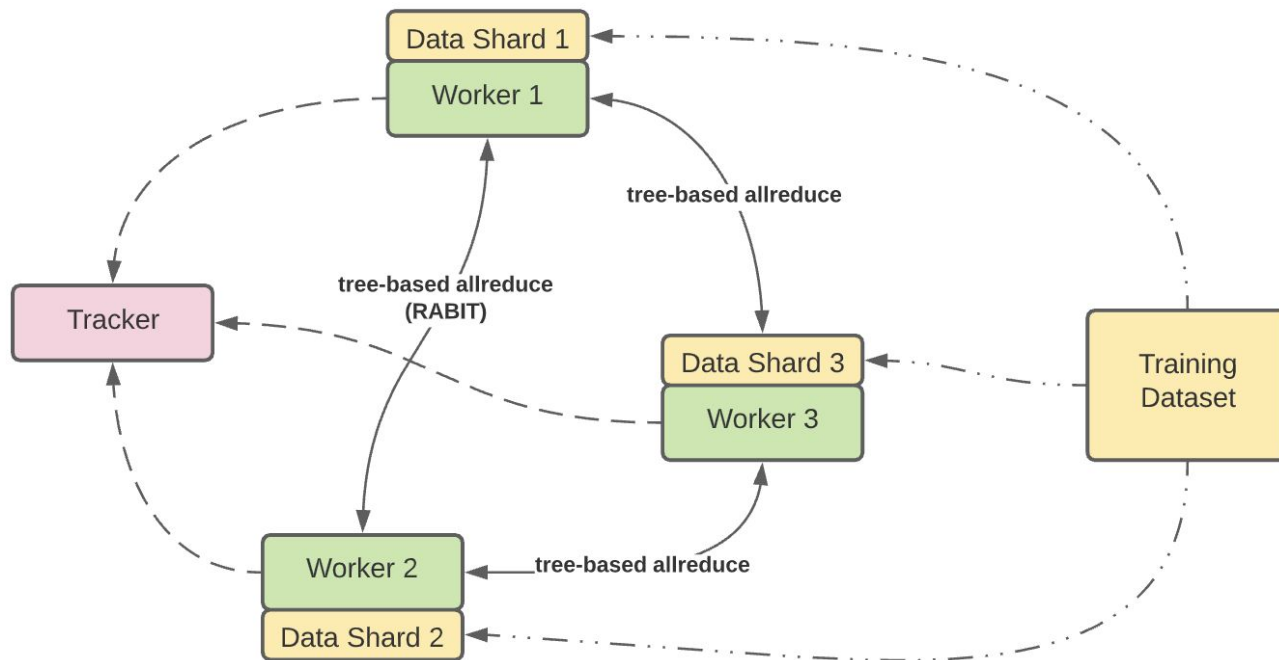
# Motivation

- There are existing solutions for distributed XGBoost
  - E.g. Apache Spark, Dask, Kubernetes etc
- But most existing solutions have shortcomings:
  - Dynamic computation graphs
  - Fault tolerance handling
  - GPU support
  - Integration with hyperparameter tuning libraries

# XGBoost-Ray

- Ray actors for stateful training workers
- Advanced fault tolerance mechanisms
- Full (multi) GPU support
- Locality-aware distributed data loading
- Integration with Ray Tune

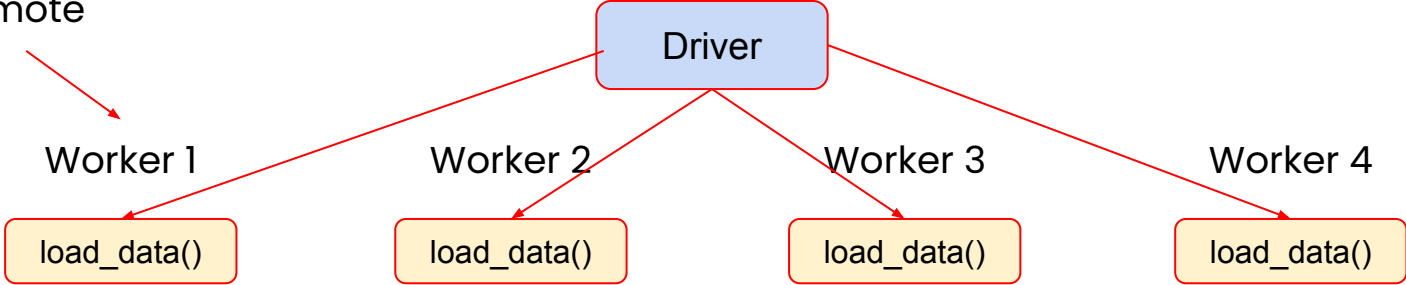
# Distributed XGBoost Architecture



# Architecture

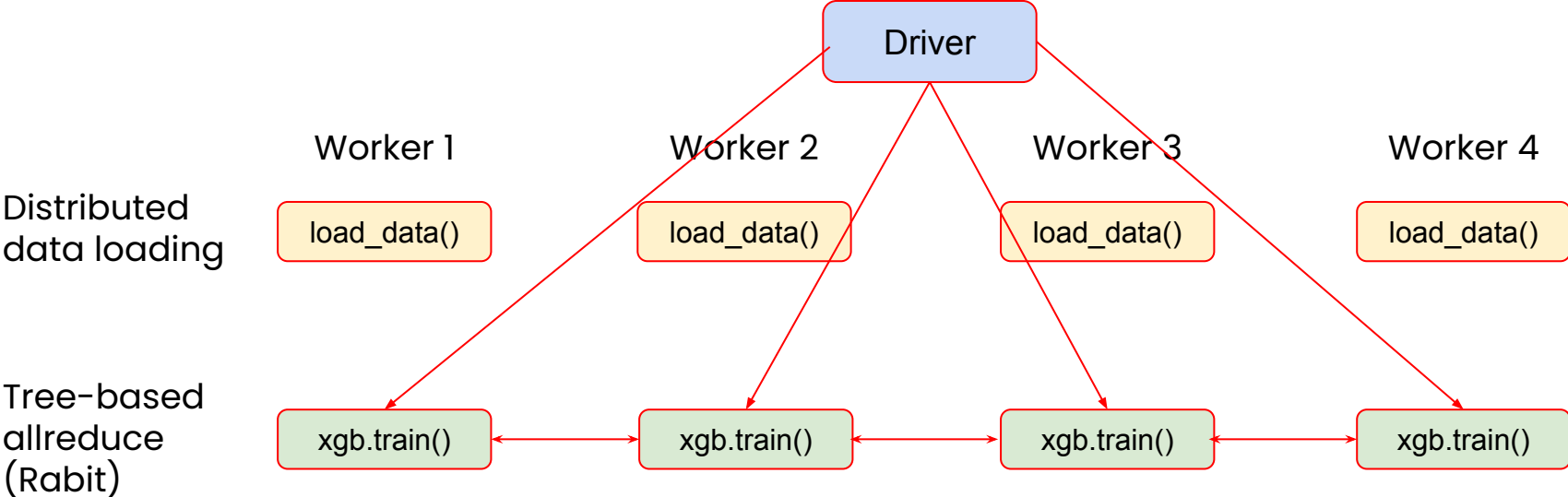
@ray.remote  
Actors

Distributed  
data loading

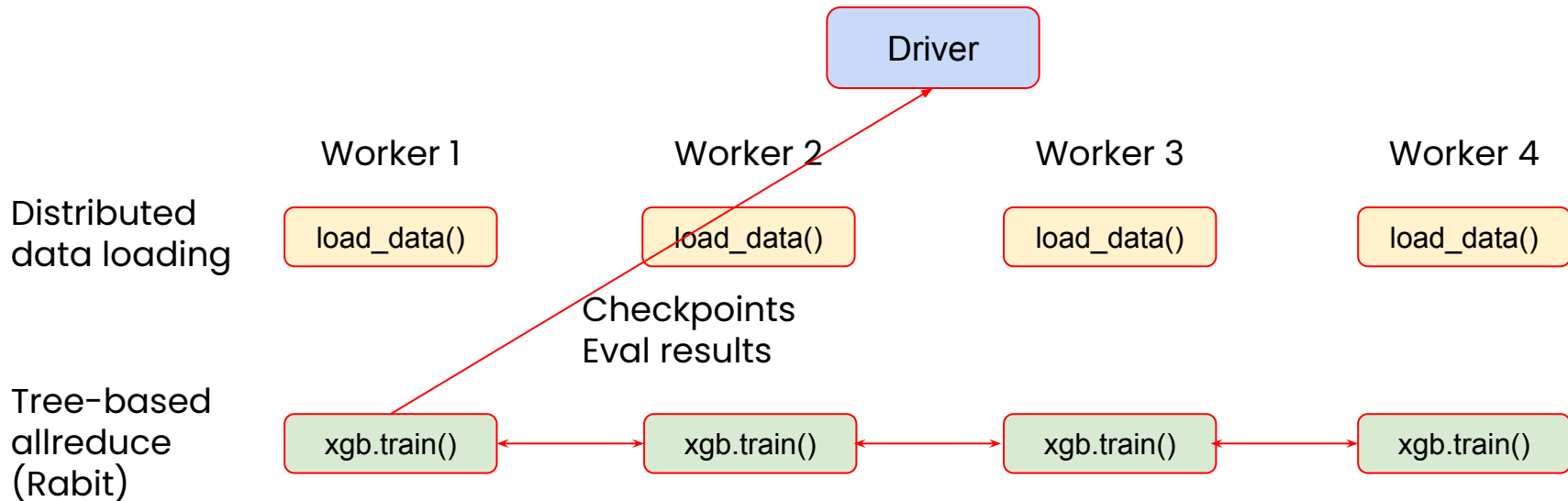




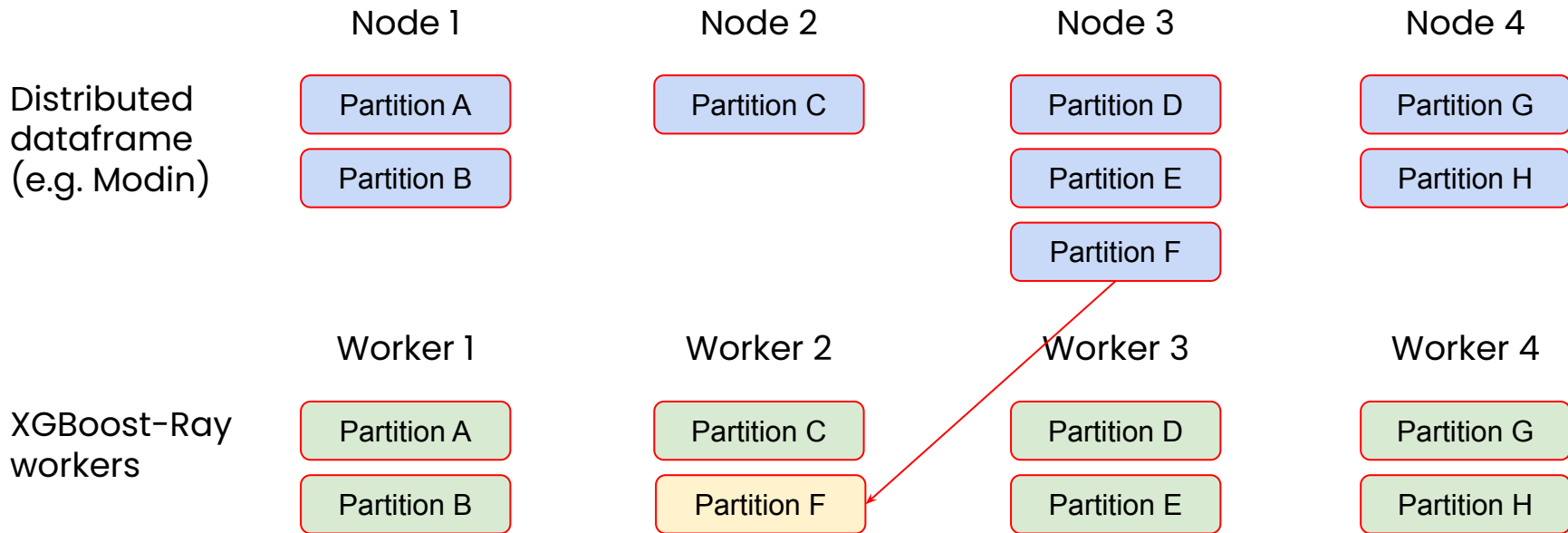
# Architecture



# Architecture



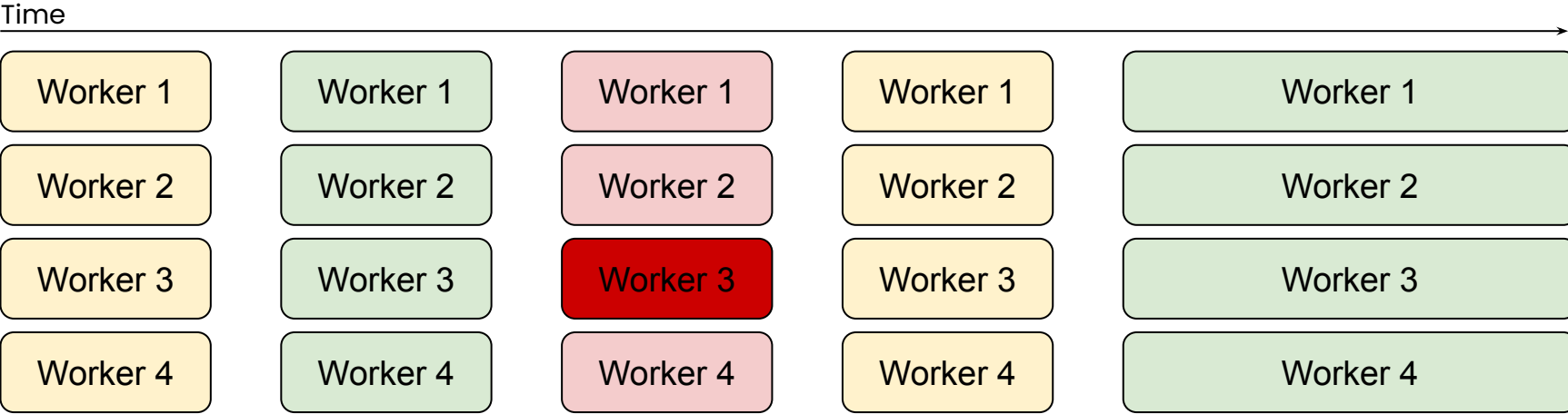
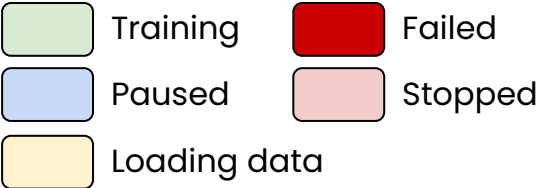
# Distributed data loading



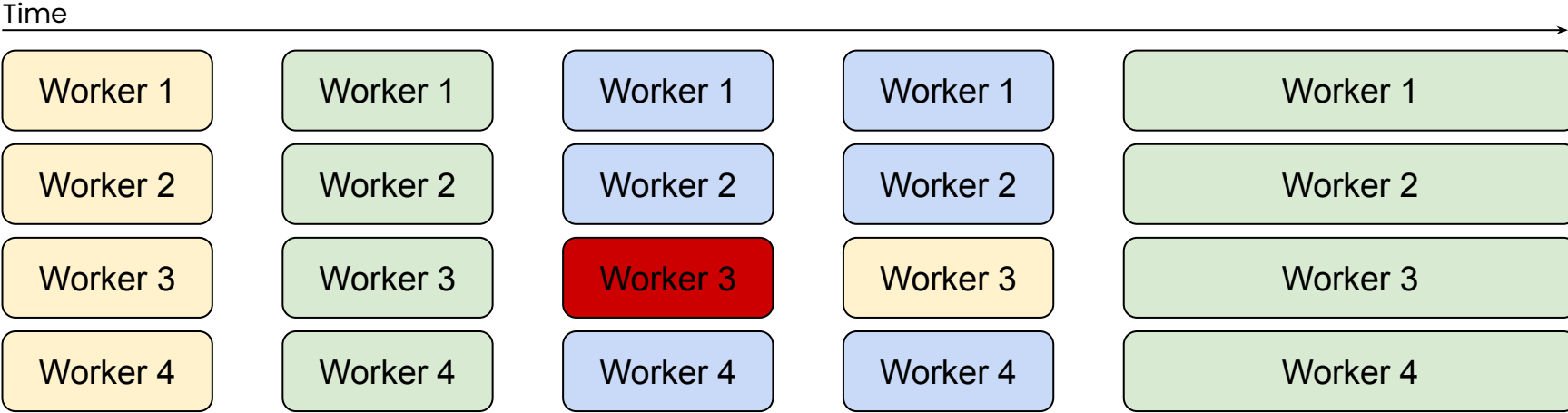
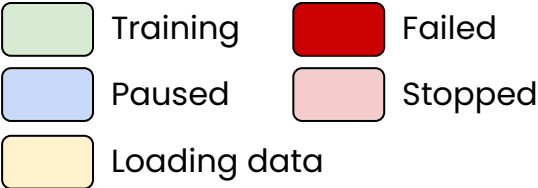
# Fault tolerance strategies

- In distributed training, some worker nodes are bound to fail eventually
- **Default:** Simple (cold) restart from last checkpoint
- **Non-elastic** training (warm restart):  
Only failing worker restarts
- **Elastic training:** Continue training with fewer workers until failed actor is back

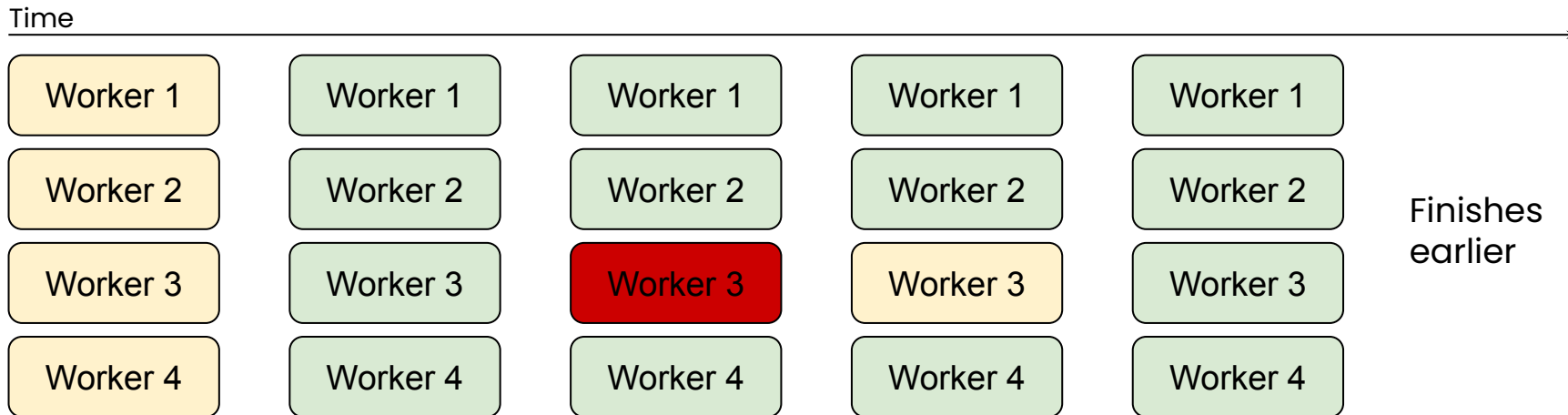
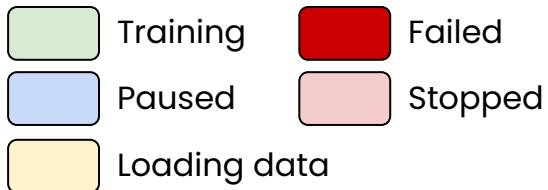
# Fault tolerance: Simple (cold) restart



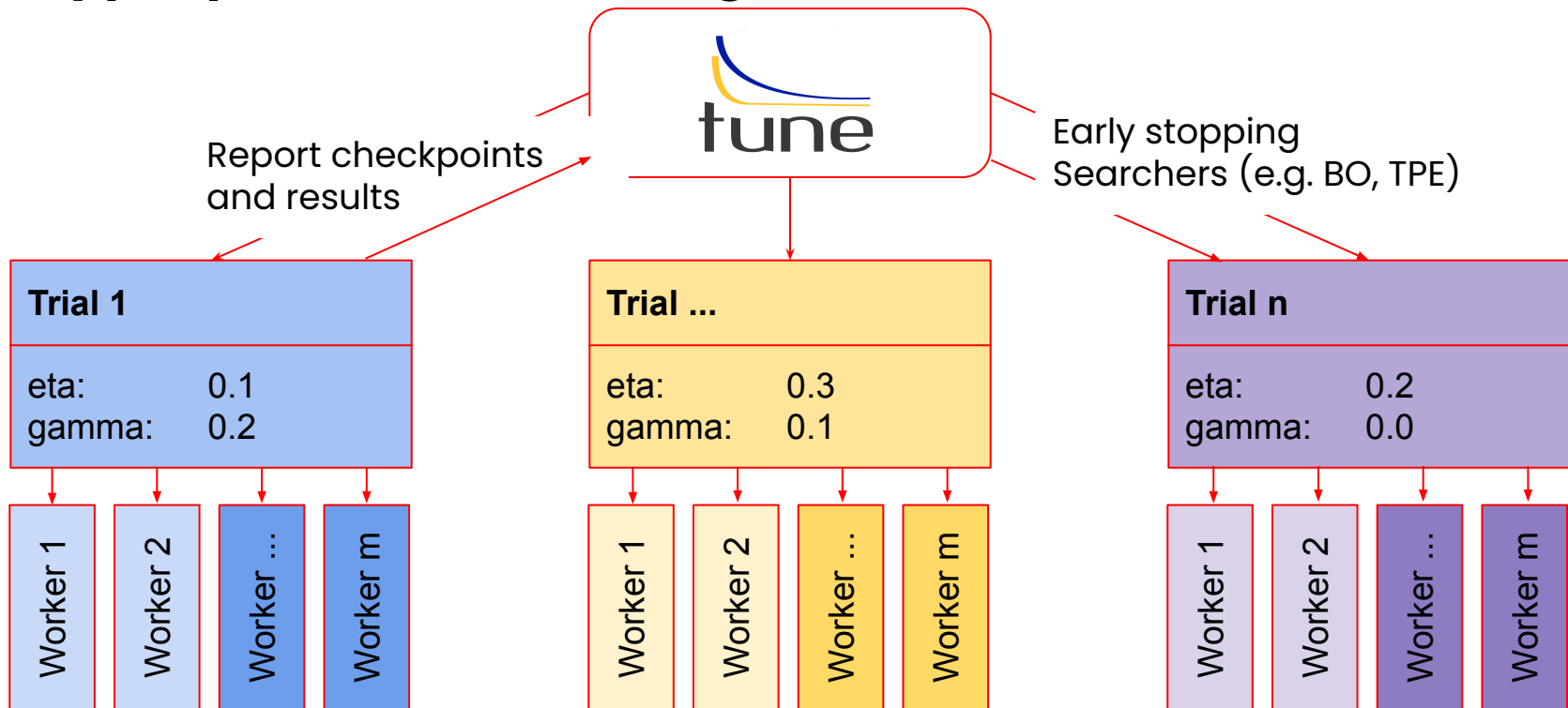
# Fault tolerance: Non-elastic training (warm restart)



# Fault tolerance: Elastic training



# Hyperparameter tuning





# Simple API example

```
from sklearn.datasets import load_breast_cancer
from xgboost_ray import RayDMatrix, RayParams, train

train_x, train_y = load_breast_cancer(return_X_y=True)
train_set = RayDMatrix(train_x, train_y)

bst = train(
    {"objective": "binary:logistic"},
    train_set,
    ray_params=RayParams(num_actors=2)
)
bst.save_model("trained.xgb")
```

# Takeaways

- **Distributed computing is a necessity & norm**
- **Ray's vision: make distributed programming simple**
  - **Don't have to be distributed systems expert. Just use `@ray.remote` :)**
- **Scale your ML workloads with Ray Libraries**

MARCH 29 - VIRTUAL - FREE

# Production RL Summit

*A reinforcement learning  
event for practitioners*

Register: <https://tinyurl.com/mr9rd32h>



ORGANIZED BY  **anyscale**



**Sergey Levine**

  
Berkeley  
UNIVERSITY OF CALIFORNIA



**Ben Kasper**

  
RIOT  
GAMES



**Sumitra Ganesh**

J.P.Morgan



**Adam Kelloway**

  
DOW<sup>®</sup>



**Marc Weber**

**SIEMENS**



**Volkmar Sterzing**

**SIEMENS**

MARCH 29 - VIRTUAL

# Production RL Summit

*A reinforcement learning  
event for practitioners*

Register: <https://tinyurl.com/mr9rd32h>



~~\$75~~ \$30

Use code DCRL2022

ORGANIZED BY  **anyscale**

HANDS-ON TUTORIAL

## Contextual Bandits & RL with RLlib

Learn how to apply cutting edge RL in production with RLlib.

### Tutorial covers:

- Brief overview of RL concepts.
- Train and tune contextual bandits and SlateQ algorithm
- Offline RL using cutting-edge algos
- Deploy RL models into a live service

~~\$75~~ \$30 (use code DCRL2022)



Instructor:

Sven Mika, Lead maintainer, RLlib



RAY SUMMIT 2021 • JUNE 22 – 24

**RAY SUMMIT 2022**

AUGUST 23 & 24TH | SAN FRANCISCO

# Call for Papers is Now Open!

Submit your talk at  
[anyscale.com/ray-summit-2022](https://anyscale.com/ray-summit-2022)

**DON'T WAIT!**  
CFP closes  
April 11th

 **anyscale**

# Start learning Ray and contributing ...

**Getting Started:** `pip install ray`

## **Documentation (docs.ray.io)**

*Quick start example, reference guides, etc*

## **Join Ray Meetup**

*Revived in Jan 2022. Next meetup March 2nd.*

*Meetup each month and publish recording to the members*

<https://www.meetup.com/Bay-Area-Ray-Meetup/>

## **Forums (discuss.ray.io)**

*Learn / share with broader Ray community, including core team*

## **Ray Slack**

*Connect with the Ray team and community*

## **Social Media (@raydistributed, @anyscalecompute)**

*Follow us on Twitter and linkedIn*

## **GitHub**

*Check out sources, file an issue, become a contributor, give us a **Star** :)*

<https://github.com/ray-project/ray>

# Thank you!

---

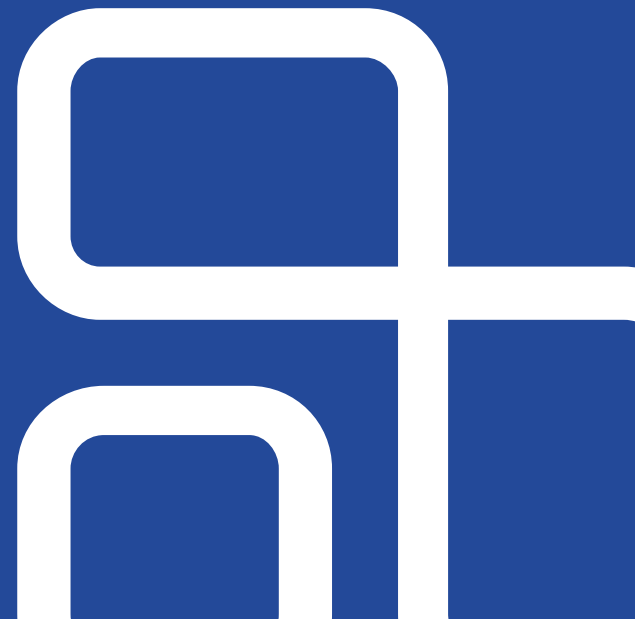
Let's stay in touch:

jules@anyscale.com

<https://www.linkedin.com/in/dmatrix/>



@2twitme



# VIDEO

---



