# Introduction To Ray

A Distributed Computing Framework

Bhagirathi Hegde

Sarath Srinivas

# Simple Batch Inference

```python
15    def predict(image_batch):
16        label_list = []
17        for image in image_batch:
18            try:
19                img = Image.open(image)
20                processor = ViTImageProcessor.from_pretrained("google/vit-base-patch16-224")
21                model = ViTForImageClassification.from_pretrained("google/vit-base-patch16-224")
22                inputs = processor(images=img, return_tensors="pt")
23                outputs = model(**inputs)
24                logits = outputs.logits
25                predicted_class_idx = logits.argmax(-1).item()
26                label_list.append(model.config.id2label[predicted_class_idx])
27  >         except ValueError as e: ···
30        return label_list
31
32
33    if __name__ == "__main__":
34        images = glob.glob(f"{DATA_FOLDER}*.JPEG")[:300]
35        # Split into chunks of 15 images each
36        image_batches = [images[i:i + 15] for i in range(0, len(images), 15)]
37        st = time.perf_counter()
38        results = [predict(image_batch) for image_batch in image_batches]
39        en = time.perf_counter()
40        print("Inference Throughput (images/sec): ", len(images) / (en - st))
```

# Program Execution

```
ray@805d0f20ec54:~/ray-scaling-experiments$ python3 scripts/simple_batch_inference.py
Inference Throughput (images/sec):  0.7533543410439415
```

Single processor

```
ray@805d0f20ec54:~/ray-scaling-experiments$ python3 scripts/simple_batch_inference.py
Inference Throughput (images/sec):  2.1059700257584897
```

Multiprocessing on 4 cores

```
ray@805d0f20ec54:~/ray-scaling-experiments$ python3 scripts/simple_batch_inference.py
2024-05-19 22:35:18,377 INFO packaging.py:530 -- Creating a file package for local directory '/home/ray/ray-scaling-exper
2024-05-19 22:35:18,463 INFO packaging.py:358 -- Pushing file package 'gcs://_ray_pkg_ab7f7809f7bd9c09.zip' (8.36MiB) to
2024-05-19 22:35:20,838 INFO packaging.py:371 -- Successfully pushed file package 'gcs://_ray_pkg_ab7f7809f7bd9c09.zip'.
Inference Throughput (images/sec):  10.026132188471529
```

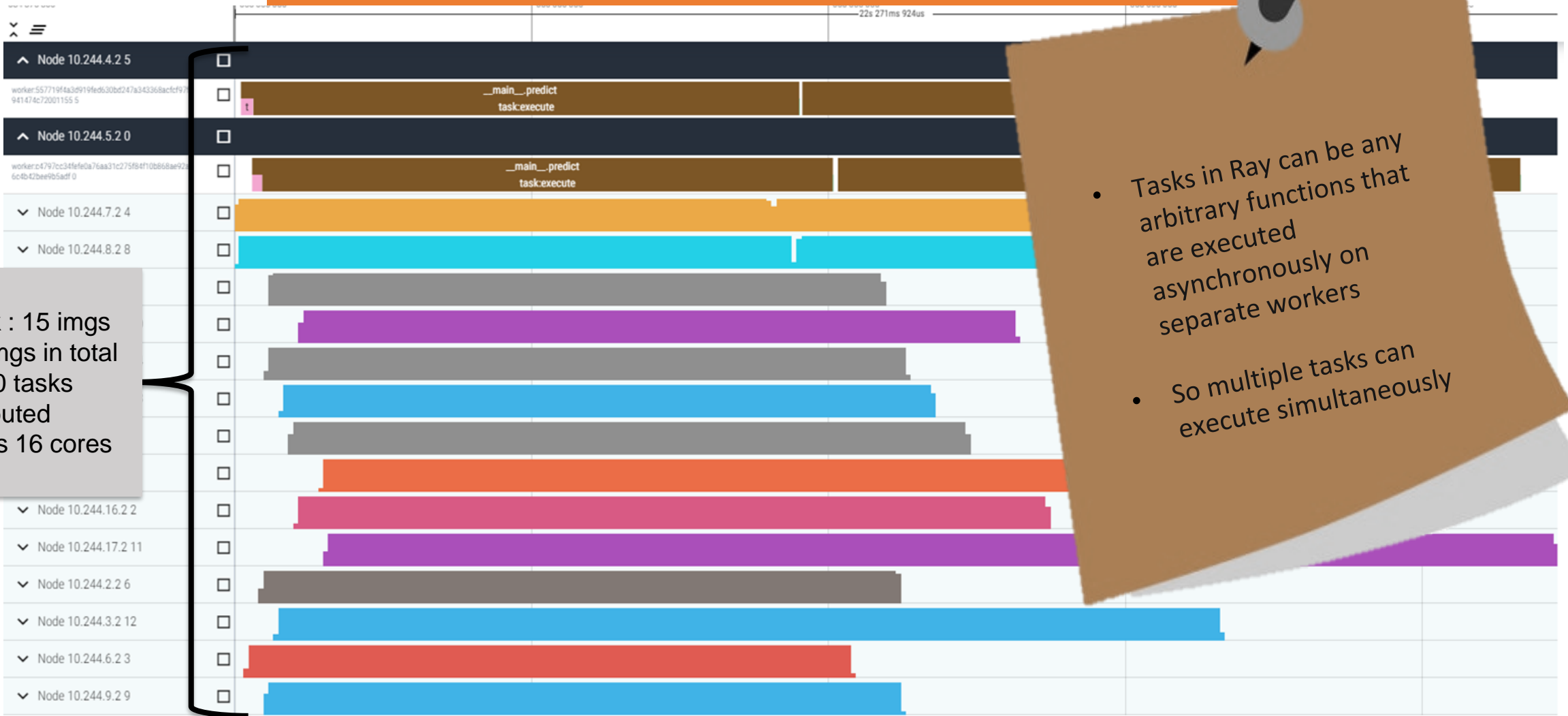Ray cluster with 16  cores (multiple machines)

How to scale?

# Distributed Execution

Ray distributes tasks across a cluster of machines, accelerating workflows
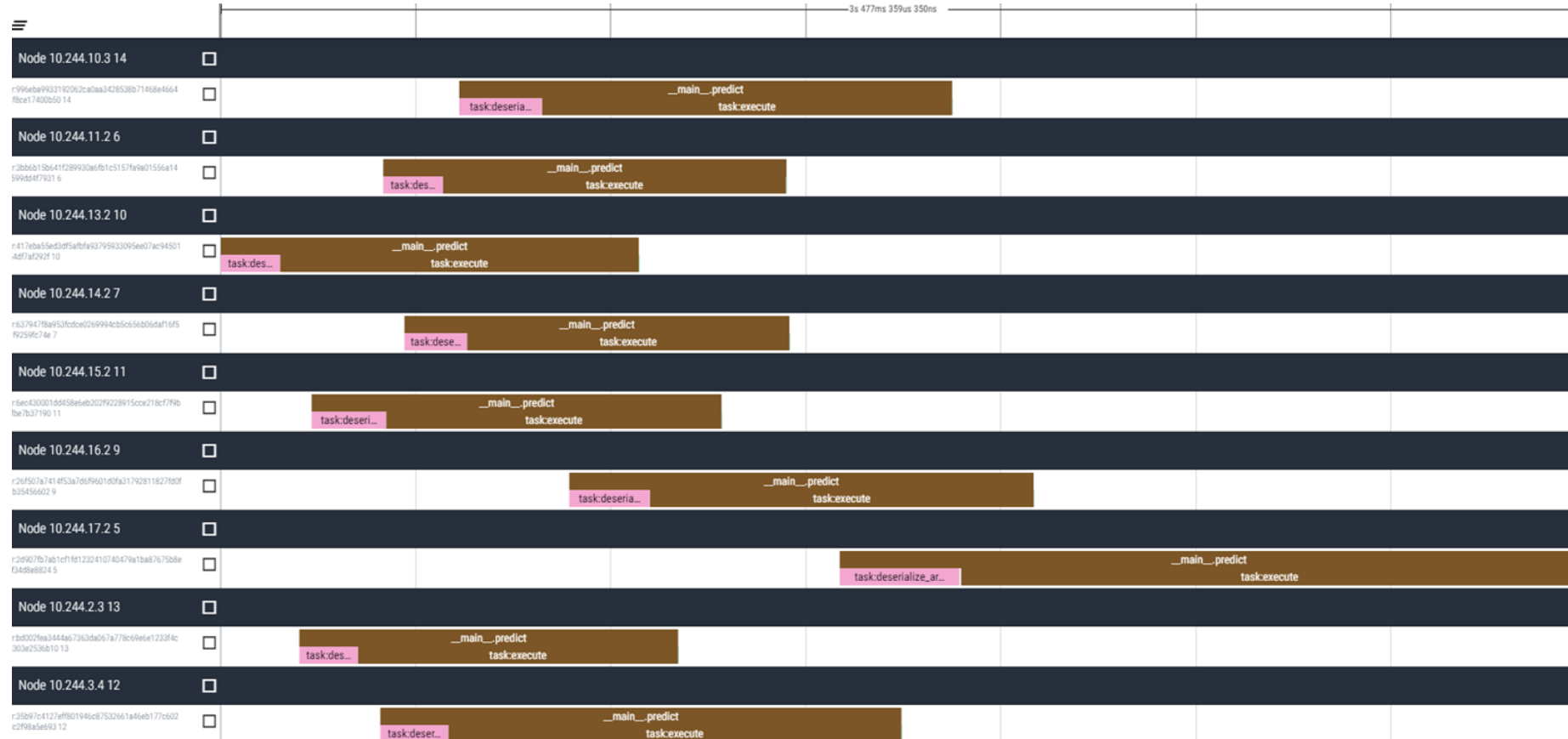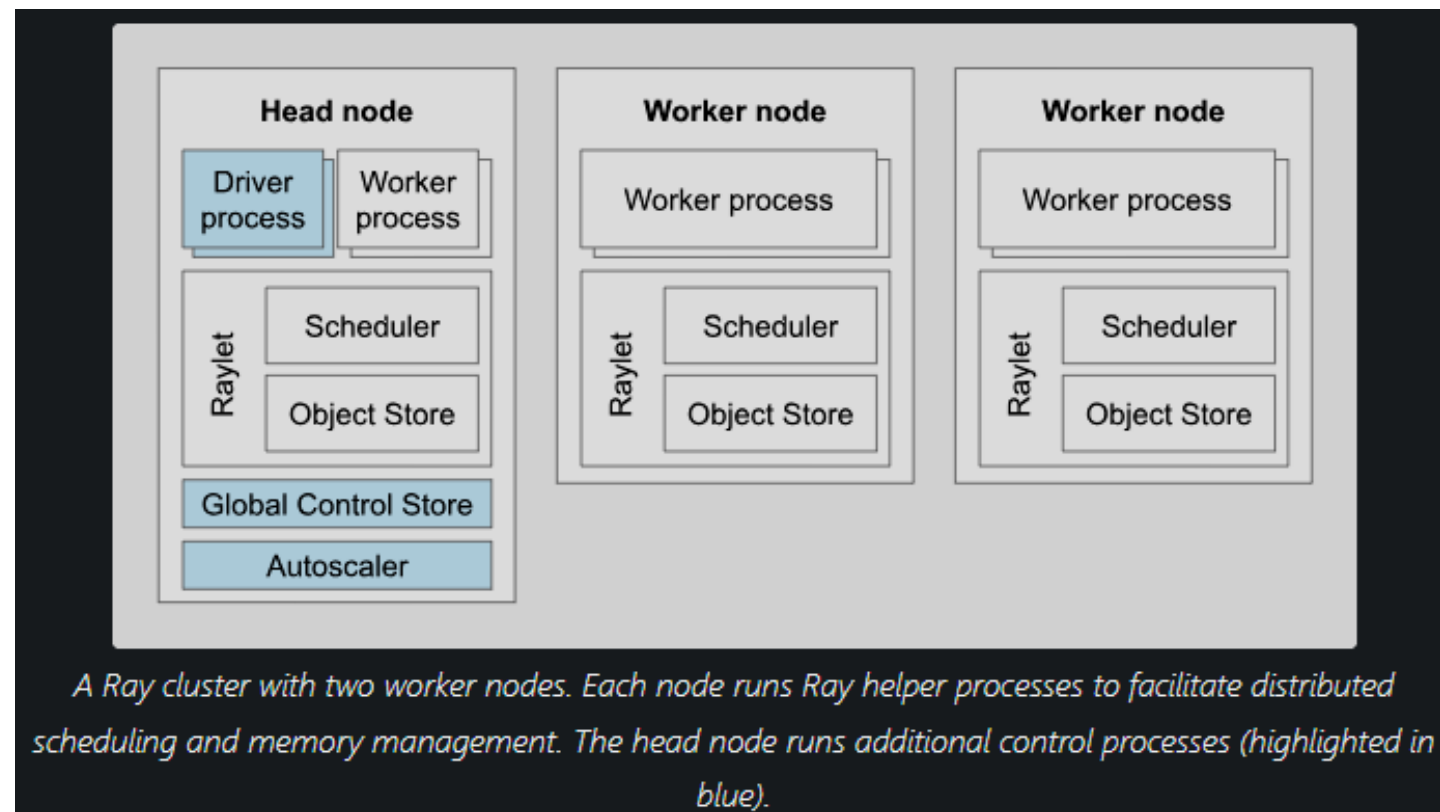
1 task : 15 imgs
300 imgs in total
So, 20 tasks
distributed
across 16 cores

- Tasks in Ray can be any arbitrary functions that are executed asynchronously on separate workers

- So multiple tasks can execute simultaneously

# Distribution overhead

Amortized for big enough tasks



```
ray@805d0f20ec54:~/ray-scaling-experiments$ python3 scripts/simple_batch_inference.py
Inference Throughput (images/sec):  1.5786928062833752
```

# The only code changes

# Ray cluster



A Ray cluster with two worker nodes. Each node runs Ray helper processes to facilitate distributed scheduling and memory management. The head node runs additional control processes (highlighted in blue).

# Starting a ray cluster is very simple

```yaml
cluster_name: default

provider:
    type: local
    head_ip: 20.106.179.167
    worker_ips: [20.127.236.169, 20.127.238.96]

# How Ray will authenticate with newly launched nodes.
auth:
    ssh_user: ray-admin
    ssh_private_key: ~/.ssh/id_rsa

# Command to start ray on the head node. You don't need to change this.
head_start_ray_commands:
    - ray stop
    - ulimit -c unlimited && ray start --head --port=6379 --num-cpus=0

# Command to start ray on worker nodes. You don't need to change this.
worker_start_ray_commands:
    - ray stop
    - ray start --address=$RAY_HEAD_IP:6379  --num-cpus=NUM_CPUS_VAR
```
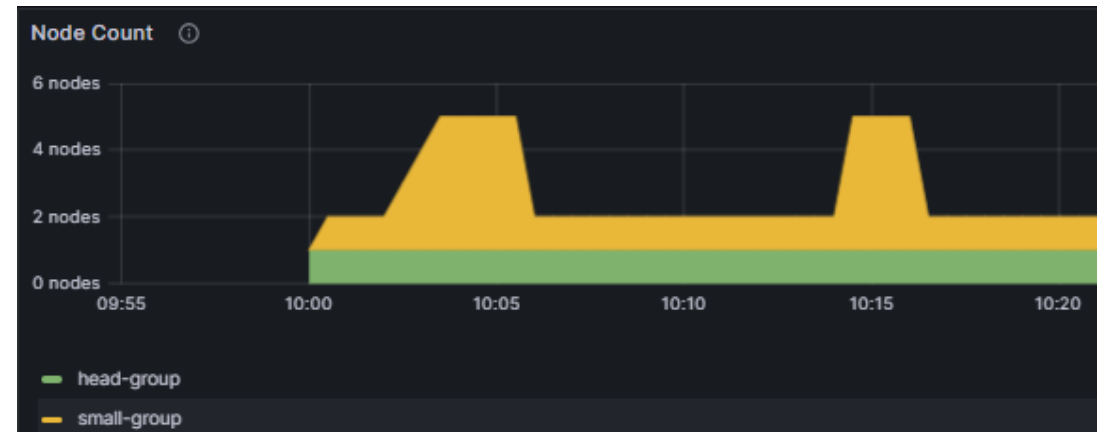
`ray-admin@ray-head-vm:~/ray-scaling-experiments/deploy$ ray up cluster.yaml`

# Autoscaler with kuberay

Ray scales up and down resources dynamically based on requirements/requests

# Run Hybrid Workloads

Task Timelines: Blocking vs. Non-Blocking sleep

# Benefits of Ray over other distributed computing options

**RAY**
- Python-centric
- Flexible and general-purpose
- Scales efficiently
- Optimized for real-time and low-latency applications

vs

**Spark**
- MapReduce centric, dataflow-based – limited task flexibility
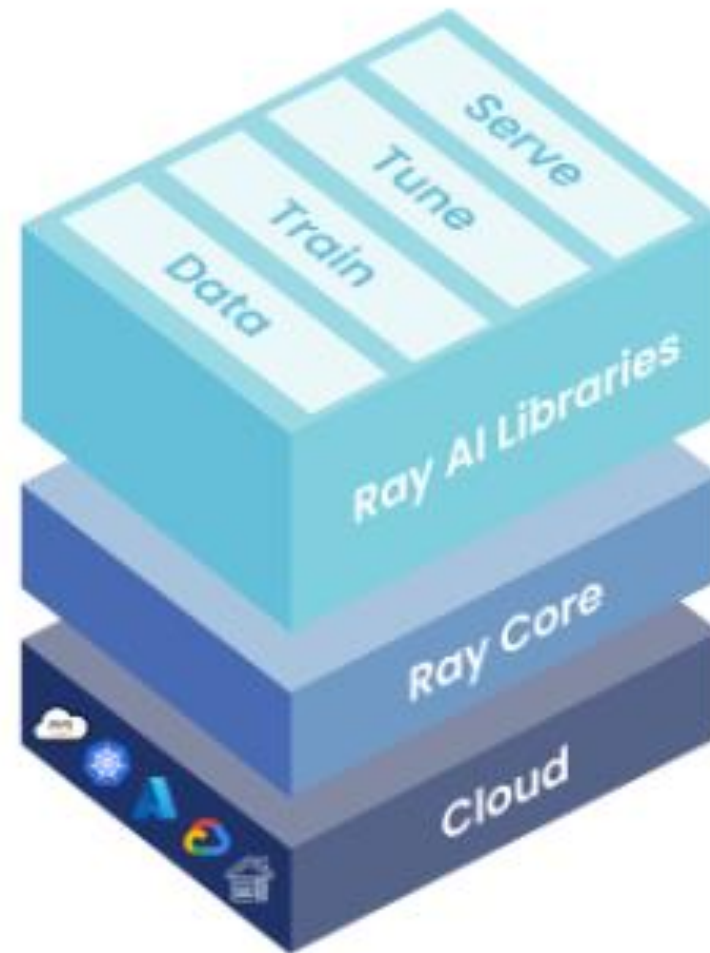- Natively scala

Complex, Limited task management

**dask**
- Mostly Data Parallel/ Mapreduce oriented

# Ray Frameworks

high-level libraries that enable simple scaling of AI workloads

a low-level distributed computing framework with a concise core and Python-first API

Serve
Tune
Train
Data
Ray AI Libraries
Ray Core
Cloud

# Summary

- Ray is a distributed computing framework for scaling up Python applications.
- Benefits:
  - Enables parallel execution of tasks across multiple machines or cores.
  - Simplifies building scalable and fault-tolerant applications.
  - Offers a user-friendly API for distributed computing.

From program on laptop to a high-performance distributed application with relatively few additional lines of python code

# Resources

- [Ray documentation](#)
- [Our GitHub repo](#)