



# Ray (Anyscale)

“Ray seeks to enable the development and composition of distributed applications and libraries  
*in general*”

# Simulation + Training + Serving for RL applications

Bulk-Synchronous parallel

Map-Reduce,  
Spark, Dryad

No fine-grain simulation  
or policy serving

Streaming

Naiad,  
Storm

Little support for distributed training and serving

Model-Serving

TFlow  
Serving,  
Clipper

No training and simulation

Distributed DL

TFlow,  
MXNet

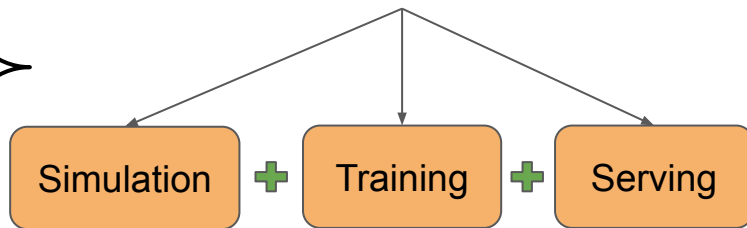
No default support for  
simulation and serving

Task-Parallel

CIEL,  
Dask

## Ray Framework

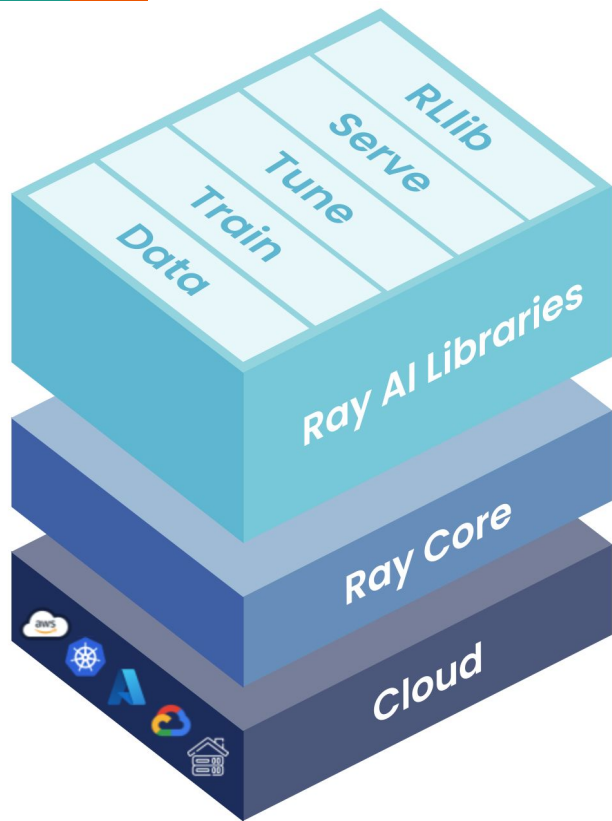
- General-Purpose
- Cluster Computing



Created at RISELab, UC Berkeley (also  
Apache Sparks, Databricks) – Pr. Ion Stoica

# Ray Layers

---

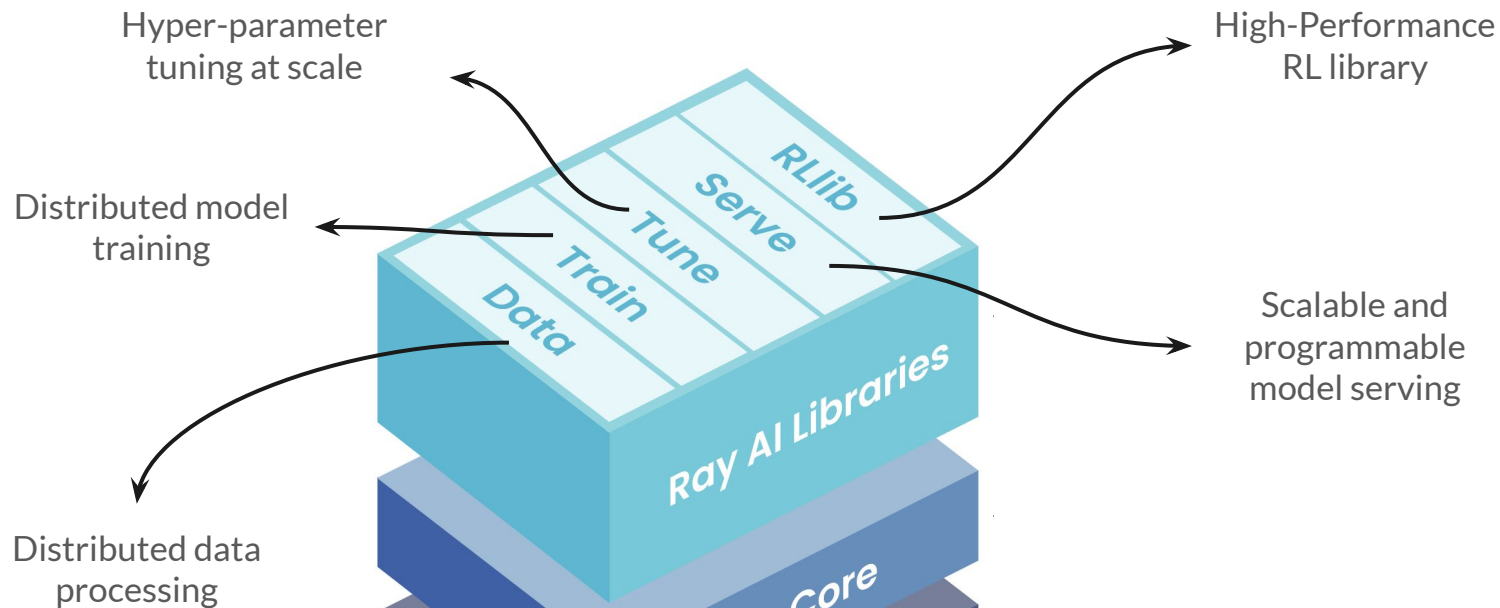


High-level libraries that enable scaling of AI workloads

Low-level distributed computing framework with concise C++ core and simple Python API

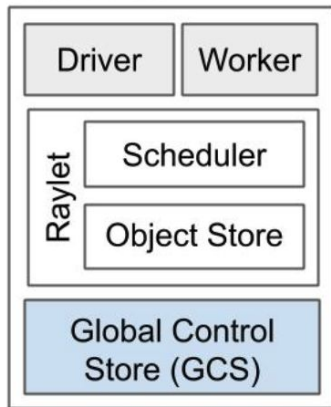
Integration with cluster environments for cloud deployment

# Ray AI Libraries

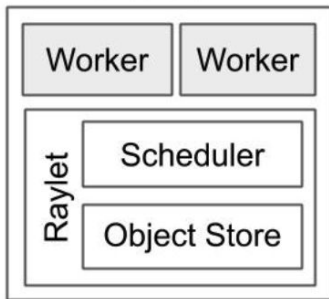


# Ray Cluster

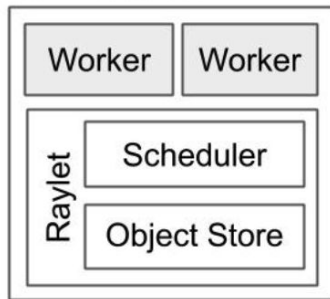
Head node



Worker node



Worker node



Separate worker nodes are instantiated in multiple-server settings.

In single-server settings Ray creates multiple **Worker** processes according to available CPUs

Node that calls `ray.init()`

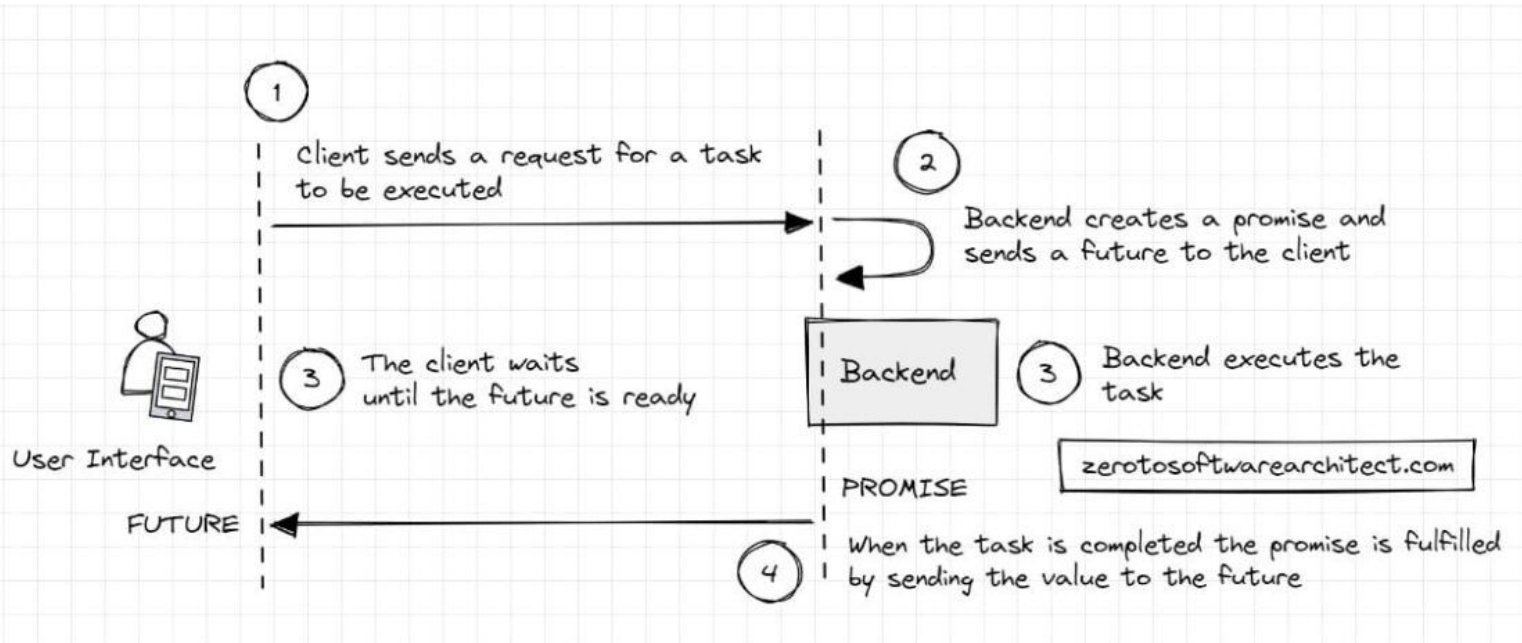
GCS:

- manages cluster-level metadata, such as the locations of actors, stored as key-value pairs that may be cached locally by workers.
- manages cluster-level operations like scheduling for placement groups, actors, cluster-node membership

Raylet:

- Scheduler responsible for resource management, task placement, and fulfilling task arguments
- Shared-memory object store (Plasma Object Store) that stores, transfers and spills large objects

# Futures & Promises

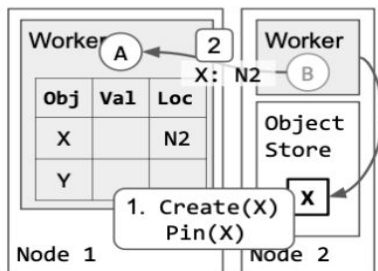
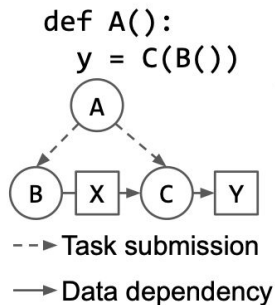


# Simple Python API with Decorators

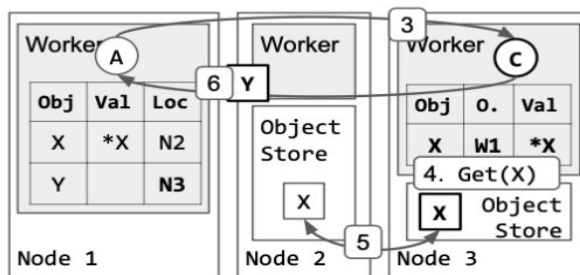
```
1  import ray
2  from random import randint
3
4  number = 12
5
6  @ray.remote
7  def create_policy():
8      return randint(1, 20)
9
10 @ray.remote
11 def update_policy(policy, rollout):
12     if rollout == 1:
13         return policy
14     else:
15         if rollout == 0:
16             return policy + 1
17         elif rollout == 2:
18             return policy - 1
```

```
20 @ray.remote
21 class Simulator(object):
22     def __init__(self):
23         # Initialize environment
24         self.env = number
25
26     def rollout(self, policy):
27         # policy-based action
28         action = policy
29         # simulate environment
30         if action < self.env:
31             return 0
32         elif action == self.env:
33             return 1
34         elif action > self.env:
35             return 2
36
37 @ray.remote
38 def train_policy():
39     # Create a policy
40     policy = create_policy.remote()
41     # Create actor
42     simulator = Simulator.remote()
43     # Do 5 steps of training
44     for _ in range(3):
45         # Rollout on actor
46         rollout = simulator.rollout.remote(policy)
47         # Update policy
48         policy = update_policy.remote(policy, rollout)
49
50     return ray.get(policy)
```

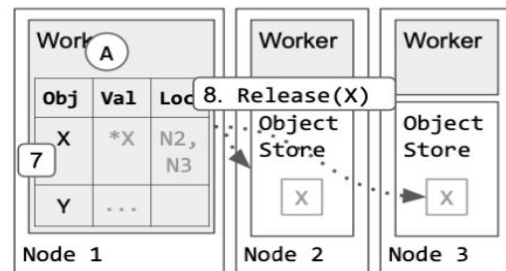
# Ownership Model



(b)



(c)

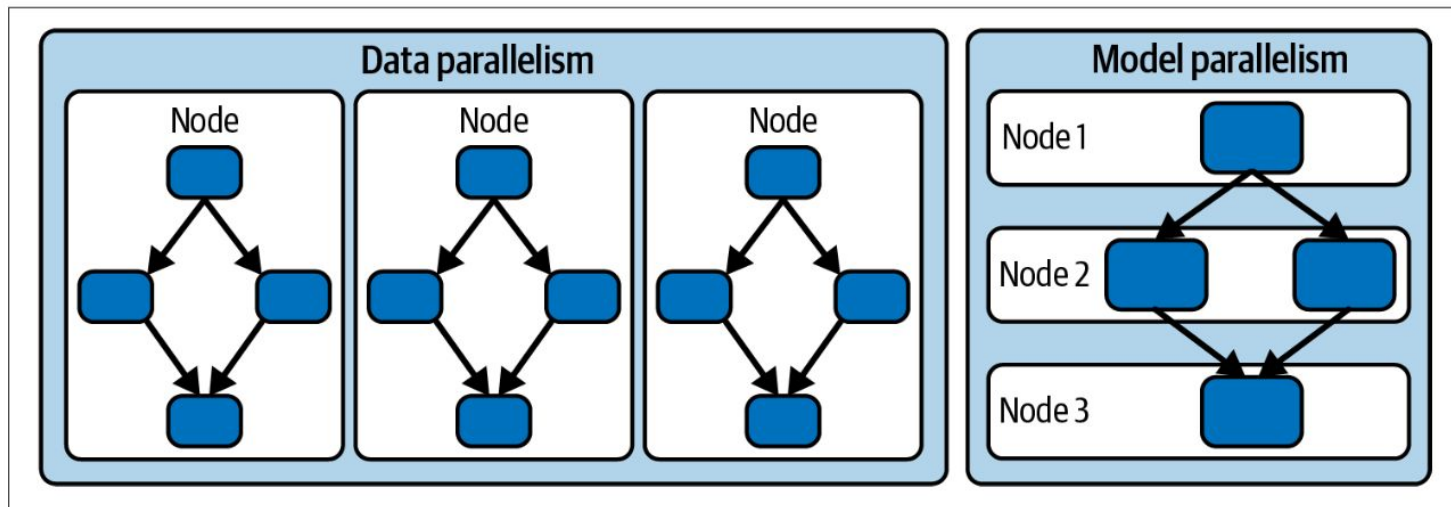


(d)

1. Large objects created and pinned at remote node (first-primary copy, deleted last)
2. Returned by reference to owner (Node 1)
3. Request of small object
- 4-5. Get(X) call to dereference large object (cached)
6. Small objects are returned by value
7. Value deletion from local table
8. Owner releases the large object

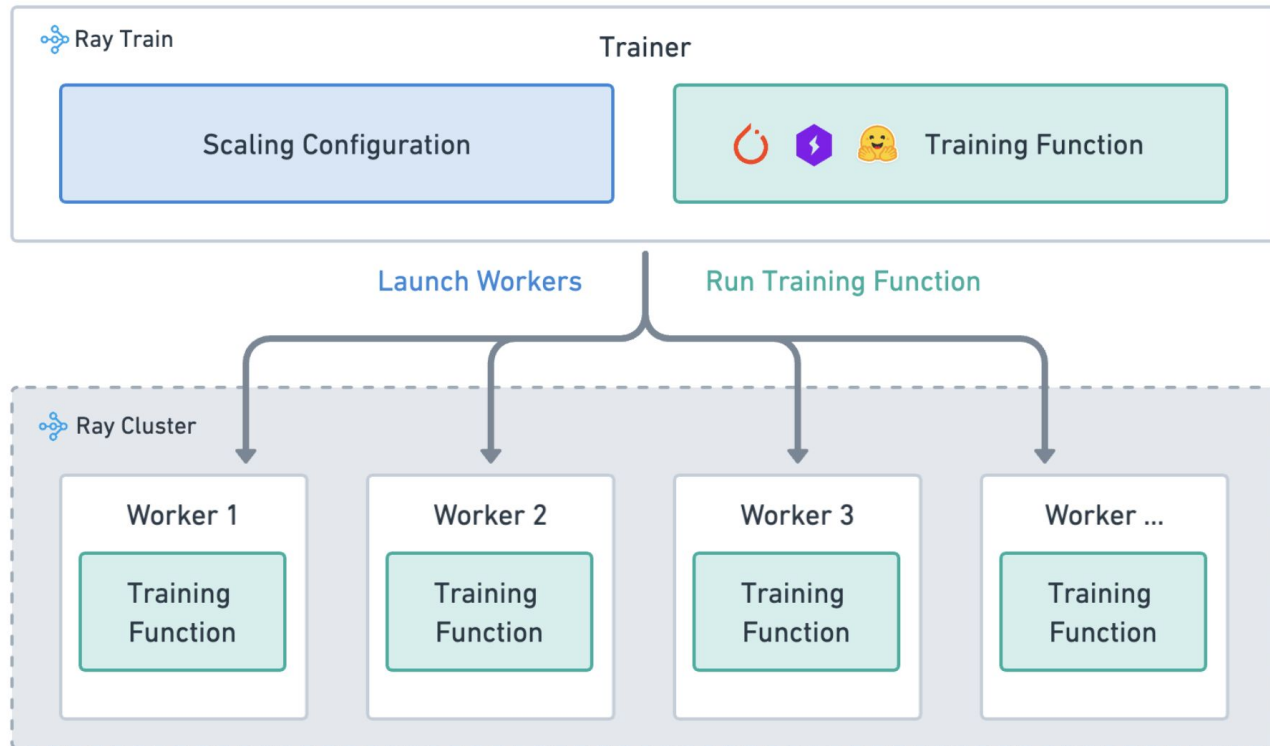


# Distributed Training



Ray Train is designed for distributed data-parallel training

# Data-Parallel Ray Workflow



## Training Function

Load model/dataset, train,  
save checkpoints and store  
metrics

# Data-Parallel Training Example – Preprocessing

```
def get_dataloaders(batch_size):
    # Transform to normalize the input images.
    transform = transforms.Compose([ToTensor(), Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])

    with FileLock(os.path.expanduser("~/data.lock")):
        # Download training data from open datasets.
        training_data = datasets.CIFAR10(
            root="~/data",
            train=True,
            download=True,
            transform=transform,
        )

        # Download test data from open datasets.
        testing_data = datasets.CIFAR10(
            root="~/data",
            train=False,
            download=True,
            transform=transform,
        )

    # Create data loaders.
    train_dataloader = DataLoader(training_data, batch_size=batch_size, shuffle=True)
    test_dataloader = DataLoader(testing_data, batch_size=batch_size)

    return train_dataloader, test_dataloader
```

# Data-Parallel Training Example – Train Function

```
def train_func_per_worker(config):
    lr = config["lr"]
    epochs = config["epochs"]
    batch_size = config["batch_size_per_worker"]

    # Get data loaders inside the worker training function.
    train_dataloader, valid_dataloader = get_dataloaders(batch_size=batch_size)

    # [1] Prepare data loader for distributed training.
    # The prepare_data_loader method assigns unique rows of data to each worker so that
    # the model sees each row once per epoch.
    # NOTE: This approach only works for map-style datasets. For a general distributed
    # preprocessing and sharding solution, see the next part using Ray Data for data
    # ingestion.
    # =====
    train_dataloader = ray.train.torch.prepare_data_loader(train_dataloader)
    valid_dataloader = ray.train.torch.prepare_data_loader(valid_dataloader)

    model = VisionTransformer(
        image_size=32, # CIFAR-10 image size is 32x32
        patch_size=4, # Patch size is 4x4
        num_layers=12, # Number of transformer layers
        num_heads=8, # Number of attention heads
        hidden_dim=384, # Hidden size (can be adjusted)
        mlp_dim=768, # MLP dimension (can be adjusted)
        num_classes=10 # CIFAR-10 has 10 classes
    )
```

```
# [2] Prepare and wrap your model with DistributedDataParallel.
# The prepare_model method moves the model to the correct GPU/CPU device.
# =====
model = ray.train.torch.prepare_model(model)

loss_fn = nn.CrossEntropyLoss()
optimizer = torch.optim.AdamW(model.parameters(), lr=lr, weight_decay=1e-2)

# Model training loop.
for epoch in range(epochs):
    if ray.train.get_context().get_world_size() > 1:
        # Required for the distributed sampler to shuffle properly across epochs
        train_dataloader.sampler.set_epoch(epoch)

    model.train()
    for X, y in tqdm(train_dataloader, desc=f"Train Epoch {epoch}"):
        pred = model(X)
        loss = loss_fn(pred, y)

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    model.eval()
    valid_loss, num_correct, num_total = 0, 0, 0
    with torch.no_grad():
        for X, y in tqdm(valid_dataloader, desc=f"Valid Epoch {epoch}"):
            pred = model(X)
            loss = loss_fn(pred, y)

            valid_loss += loss.item()
            num_total += y.shape[0]
            num_correct += (pred.argmax(1) == y).sum().item()

    valid_loss /= len(train_dataloader)
    accuracy = num_correct / num_total
```

# Data-Parallel Training Example – Train Function

---

```
# [3] (Optional) Report checkpoints and attached metrics to Ray Train.
# =====
with tempfile.TemporaryDirectory() as temp_checkpoint_dir:
    torch.save(
        model.module.state_dict(),
        os.path.join(temp_checkpoint_dir, "model.pt")
    )
    ray.train.report(
        metrics={"loss": valid_loss, "accuracy": accuracy},
        checkpoint=ray.train.Checkpoint.from_directory(temp_checkpoint_dir),
    )
    if ray.train.get_context().get_world_rank() == 0:
        print({"epoch_num": epoch, "loss": valid_loss, "accuracy": accuracy})
```

# Data-Parallel Training Example – Configuration & Results

```
def train_cifar_10(num_workers, use_gpu):
    global_batch_size = 512

    train_config = {
        "lr": 1e-3,
        "epochs": 1,
        "batch_size_per_worker": global_batch_size // num_workers,
    }

    # [1] Start distributed training.
    # Define computation resources for workers.
    # Run `train_func_per_worker` on those workers.
    # =====
    scaling_config = ScalingConfig(num_workers=num_workers, use_gpu=use_gpu)
    run_config = RunConfig(
        name=f"train_run-{uuid.uuid4().hex}",
    )
    trainer = TorchTrainer(
        train_loop_per_worker=train_func_per_worker,
        train_loop_config=train_config,
        scaling_config=scaling_config,
        run_config=run_config,
    )
    result = trainer.fit()
    print(f"Training result: {result}")

if __name__ == "__main__":
    train_cifar_10(num_workers=4, use_gpu=True)
```

```
{
  "loss": 0.35899272256967973,
  "accuracy": 0.3228,
  "timestamp": 1758600479,
  "checkpoint_dir_name": "checkpoint_000000",
  "should_checkpoint": true,
  "done": false,
  "training_iteration": 1,
  "trial_id": "d6250_00000",
  "date": "2025-09-23_00-07-59",
  "time_this_iter_s": 21.89649200439453,
  "time_total_s": 21.89649200439453,
  "pid": 2409569,
  "hostname": "yecl-gpu-cluster",
  "node_ip": "192.168.0.121",
  "config": {
    "train_loop_config": {
      "lr": 0.001,
      "epochs": 1,
      "batch_size_per_worker": 128
    }
  },
  "time_since_restore": 21.89649200439453,
  "iterations_since_restore": 1
}
```

# DeepSpeed Library



“Deep learning optimization software suite that powers unprecedented scale and speed for both training and inference”

## Training

- Speed Scale Cost
- Democratization
- MoE models
- Long sequence
- RLHF

## Inference

- Large models
- Latency
- Serving cost
- Agility

## Compression

- Model size
- Latency
- Composability
- Runnable on client devices

## Science

- Speed
- Scale
- Capability
- Diversity
- Discovery

Zero Redundancy Optimizer (ZeRO): Sharding model states, gradients and optimizer states across GPU to reduce memory usage



# ZeRO-3 Training Example – Train Function

```
23 def train_func(config):
24     """Your training function that will be launched on each worker."""
25
26     # Unpack training configs
27     set_seed(config["seed"])
28     num_epochs = config["num_epochs"]
29     train_batch_size = config["train_batch_size"]
30     eval_batch_size = config["eval_batch_size"]
31
32     # Instantiate the Model
33     model = AutoModelForSequenceClassification.from_pretrained(
34         "bert-base-cased", return_dict=True
35     )
36
37     # Prepare Ray Data Loaders
38     # =====
39     train_ds = ray.train.get_dataset_shard("train")
40     eval_ds = ray.train.get_dataset_shard("validation")
41
42     tokenizer = AutoTokenizer.from_pretrained("bert-base-cased")
43
44     def collate_fn(batch):
45         outputs = tokenizer(
46             list(batch["sentence1"]),
47             list(batch["sentence2"]),
48             truncation=True,
49             padding="longest",
50             return_tensors="pt",
51         )
52         outputs["labels"] = torch.LongTensor(batch["label"])
53         return outputs
```

```
55     train_dataloader = train_ds.iter_torch_batches(
56         batch_size=train_batch_size, collate_fn=collate_fn
57     )
58     eval_dataloader = eval_ds.iter_torch_batches(
59         batch_size=eval_batch_size, collate_fn=collate_fn
60     )
61     # =====
62
63     # Initialize DeepSpeed Engine
64     model, optimizer, _, lr_scheduler = deepspeed.initialize(
65         model=model,
66         model_parameters=model.parameters(),
67         config=deepspeed_config,
68     )
69     device = get_accelerator().device_name(model.local_rank)
70
71     # Initialize Evaluation Metrics
72     f1 = BinaryF1Score().to(device)
73     accuracy = BinaryAccuracy().to(device)
```



# ZeRO-3 Training Example – Train Function

```
75 for epoch in range(num_epochs):
76     # Training
77     model.train()
78     for batch in train_dataloader:
79         batch = {k: v.to(device) for k, v in batch.items()}
80         outputs = model(**batch)
81         loss = outputs.loss
82         model.backward(loss)
83         optimizer.step()
84         lr_scheduler.step()
85         optimizer.zero_grad()
86
87     # Evaluation
88     model.eval()
89     for batch in eval_dataloader:
90         batch = {k: v.to(device) for k, v in batch.items()}
91         with torch.no_grad():
92             outputs = model(**batch)
93             predictions = outputs.logits.argmax(dim=-1)
94
95             f1.update(predictions, batch["labels"])
96             accuracy.update(predictions, batch["labels"])
97
98     # torchmetrics will aggregate the metrics across all workers
99     eval_metric = {
100         "f1": f1.compute().item(),
101         "accuracy": accuracy.compute().item(),
102     }
103     f1.reset()
104     accuracy.reset()
```

```
106 if model.global_rank == 0:
107     print(f"epoch {epoch}:", eval_metric)
108
109 # Report checkpoint and metrics to Ray Train
110 # =====
111 with TemporaryDirectory() as tmpdir:
112     # Each worker saves its own checkpoint shard
113     model.save_checkpoint(tmpdir)
114
115     # Ensure all workers finished saving their checkpoint shard
116     torch.distributed.barrier()
117
118     # Report checkpoint shards from each worker in parallel
119     ray.train.report(
120         metrics=eval_metric, checkpoint=Checkpoint.from_directory(tmpdir)
121     )
122 # =====
```

# ZeRO-3 Training Example – Configuration

```
125 if __name__ == "__main__":
126     deepspeed_config = {
127         "optimizer": {
128             "type": "AdamW",
129             "params": {
130                 "lr": 2e-5,
131             },
132         },
133         "scheduler": {"type": "WarmupLR", "params": {"warmup_num_steps": 100}},
134         "fp16": {"enabled": True},
135         "bf16": {"enabled": False}, # Turn this on if using AMPERE GPUs.
136         "zero_optimization": {
137             "stage": 3,
138             "offload_optimizer": {
139                 "device": "none",
140             },
141             "offload_param": {
142                 "device": "none",
143             },
144         },
145         "gradient_accumulation_steps": 1,
146         "gradient_clipping": True,
147         "steps_per_print": 10,
148         "train_micro_batch_size_per_gpu": 16,
149         "wall_clock_breakdown": False,
150     }
```

```
152 training_config = {
153     "seed": 42,
154     "num_epochs": 3,
155     "train_batch_size": 16,
156     "eval_batch_size": 32,
157     "deepspeed_config": deepspeed_config,
158 }
159
160 # Prepare Ray Datasets
161 hf_datasets = load_dataset("glue", "mrpc")
162
163 # Convert HuggingFace datasets to Ray datasets using from_items
164 ray_datasets = {
165     "train": ray.data.from_items([dict(item) for item in hf_datasets["train"]]),
166     "validation": ray.data.from_items([dict(item) for item in hf_datasets["validation"]]),
167 }
168
169 trainer = TorchTrainer(
170     train_func,
171     train_loop_config=training_config,
172     scaling_config=ScalingConfig(num_workers=4, use_gpu=True),
173     datasets=ray_datasets,
174     dataset_config=DataConfig(datasets_to_split=["train", "validation"]),
175 )
176
177 result = trainer.fit()
```

# ZeRO-3 Training Example – Results

```
{
  "f1": 0.8944723606109619,
  "accuracy": 0.845588207244873,
  "timestamp": 1758598124,
  "checkpoint_dir_name": "checkpoint_000002",
  "should_checkpoint": true,
  "done": false,
  "training_iteration": 3,
  "trial_id": "3a66a_00000",
  "date": "2025-09-22_23-28-44",
  "time_this_iter_s": 12.80349349975586,
  "time_total_s": 73.79229402542114,
  "pid": 2378503,
  "hostname": "yecl-gpu-cluster",
  "node_ip": "192.168.0.121",
```

```
"config": {
  "train_loop_config": {
    "seed": 42,
    "num_epochs": 3,
    "train_batch_size": 16,
    "eval_batch_size": 32,
    "deepspeed_config": {
      "optimizer": {
        "type": "AdamW",
        "params": {
          "lr": 2e-05
        }
      },
      "scheduler": {
        "type": "WarmupLR",
        "params": {
          "warmup_num_steps": 100
        }
      },
      "fp16": {
        "enabled": true
      },
      "bf16": {
        "enabled": false
      },
```

```
    "zero_optimization": {
      "stage": 3,
      "offload_optimizer": {
        "device": "none"
      },
      "offload_param": {
        "device": "none"
      }
    },
    "gradient_accumulation_steps": 1,
    "gradient_clipping": true,
    "steps_per_print": 10,
    "train_micro_batch_size_per_gpu": 16,
    "wall_clock_breakdown": false
  },
  "time_since_restore": 73.79229402542114,
  "iterations_since_restore": 3
}
```

# Model-Parallel Training – Preparing Model

```
13 # -----
14 # 1. Define pipeline-parallel model
15 # -----
16 def create_pipeline_model(num_stages=2):
17     """Create a pipeline model with proper layer specifications"""
18     # Split the model into a sequence of layers
19     layers = [
20         LayerSpec(nn.Linear, 128, 256), # Stage 0
21         LayerSpec(nn.ReLU),             # Stage 0
22         LayerSpec(nn.Linear, 256, 512), # Stage 0
23         LayerSpec(nn.ReLU),             # Stage 0
24         LayerSpec(nn.Linear, 512, 256), # Stage 1
25         LayerSpec(nn.ReLU),             # Stage 1
26         LayerSpec(nn.Linear, 256, 128), # Stage 1
27         LayerSpec(nn.ReLU),             # Stage 1
28         LayerSpec(nn.Linear, 128, 10),  # Stage 1 (output layer)
29     ]
30
31     # Pipeline with specified number of stages
32     model = PipelineModule(
33         layers=layers,
34         num_stages=num_stages,
35         loss_fn=nn.CrossEntropyLoss(),
36         partition_method="uniform", # split layers evenly across stages
37         activation_checkpoint_interval=0, # disable for simplicity
38     )
39     return model
```

Utilize Deepspeed's *PipelineModule* to split the model into `num_stages` partitions

# Model-Parallel Training – Train Function



```
41 # -----
42 # 2. Ray Train worker function
43 # -----
44 def train_func(config):
45     """Training function for pipeline parallel model"""
46     # Get Ray Train context (Ray handles distributed initialization)
47     context = train.get_context()
48     world_rank = context.get_world_rank()
49     world_size = context.get_world_size()
50     local_rank = context.get_local_rank()
51
52     print(f"Worker {world_rank}/{world_size}, Local rank: {local_rank}")
53
54     # Set device
55     device = f"cuda:{local_rank}" if torch.cuda.is_available() else "cpu"
56     torch.cuda.set_device(local_rank)
57
58     # Initialize DeepSpeed distributed backend explicitly for PipelineModule
59     if not dist.is_initialized():
60         # Initialize with the same backend Ray Train uses
61         dist.init_process_group(
62             backend="nccl" if torch.cuda.is_available() else "gloo",
63             world_size=world_size,
64             rank=world_rank
65         )
66     print(f"Initialized distributed backend on rank {world_rank}")
67
```

```
68 # Initialize DeepSpeed's distributed backend
69 import deepspeed.comm as dist_comm
70 if not dist_comm.is_initialized():
71     dist_comm.init_distributed(
72         dist_backend="nccl" if torch.cuda.is_available() else "gloo",
73         auto_mpi_discovery=False,
74         distributed_port=29500,
75         verbose=False
76     )
77     print(f"DeepSpeed distributed backend initialized on rank {world_rank}")
78
79 # Create pipeline model
80 num_stages = config.get("num_stages", 2)
81 model = create_pipeline_model(num_stages=num_stages)
82
83 print(f"Created pipeline model with {num_stages} stages")
84
85 # Create proper DeepSpeed configuration with batch size constraints
86 deepspeed_config = config["deepspeed_config"].copy()
```

# Model-Parallel Training – Train Function

```
99     deepspeed_config.update({
100         "train_batch_size": train_batch_size,
101         "train_micro_batch_size_per_gpu": micro_batch_size,
102         "gradient_accumulation_steps": gradient_accumulation_steps,
103     })
104
105     # Initialize DeepSpeed with pipeline model
106     engine, optimizer, _, _ = deepspeed.initialize(
107         model=model,
108         config=deepspeed_config,
109         dist_init_required=False, # Ray Train handles distributed init
110     )
```

# Model-Parallel Training – Configuration

```
197 if __name__ == "__main__":
198     ray.init(num_gpus=2)
199     print("Ray initialized for pipeline-parallel training")
200
201     # DeepSpeed config for pipeline parallelism + ZeRO stage 1
202     deepspeed_config = {
203         "train_micro_batch_size_per_gpu": 4, # Smaller batch size per GPU
204         "gradient_accumulation_steps": 2,     # Accumulate gradients
205         "optimizer": {
206             "type": "Adam",
207             "params": {"lr": 1e-3, "weight_decay": 0.01},
208         },
209         "fp16": {"enabled": False}, # Keep FP32 for simplicity and stability
210         "zero_optimization": {
211             "stage": 1, # Stage 1 works well with pipeline parallelism
212             "allgather_partitions": True,
213             "reduce_scatter": True,
214             "allgather_bucket_size": 2e8,
215             "reduce_bucket_size": 2e8,
216         },
217         "steps_per_print": 5,
218         "wall_clock_breakdown": False,
219
220         # Pipeline-specific configurations
221         "pipeline": {
222             "pipe_partitioned": True,
223             "grad_partitioned": True,
224         }
225     }
```

```
226
227 # Training configuration
228 train_config = {
229     "epochs": 2, # Reduced epochs for demo
230     "num_stages": 2, # Number of pipeline stages (should match num_workers)
231     "deepspeed_config": deepspeed_config,
232 }
233
234 # Scaling configuration - one worker per pipeline stage
235 scaling_config = ScalingConfig(
236     num_workers=2, # Must match num_stages for pipeline parallelism
237     use_gpu=True,
238     resources_per_worker={"GPU": 1}, # One GPU per worker/stage
239 )
240
241 # Run configuration
242 run_config = RunConfig(
243     # verbose=2, # Enable for more detailed logging
244 )
245
246 # Create trainer
247 trainer = TorchTrainer(
248     train_func,
249     train_loop_config=train_config,
250     scaling_config=scaling_config,
251     run_config=run_config,
252 )
```



# Model-Parallel Training – Results

```
try:
    result = trainer.fit()
    print("\nPipeline-parallel training completed successfully!")
    print(f"Final metrics: {result.metrics}")

except Exception as e:
    print(f"Training failed with error: {e}")
    print("Common issues:")
    print("1. Make sure you have at least 2 GPUs available")
    print("2. Ensure DeepSpeed is properly installed")
    print("3. Check that CUDA and NCCL are working correctly")

finally:
    ray.shutdown()
    print("Ray shut down")
```

```
{
  "loss": 2.2838621139526367,
  "epoch": 1,
  "learning_rate": 0.001,
  "timestamp": 1758667846,
  "checkpoint_dir_name": "None",
  "done": true,
  "training_iteration": 2,
  "trial_id": "b853e_00000",
  "date": "2025-09-23_18-50-46",
  "time_this_iter_s": 0.2586488723754883,
  "time_total_s": 6.3532164096832275,
  "pid": 2492297,
  "hostname": "yecl-gpu-cluster",
  "node_ip": "192.168.0.121",
  "config": {
    "train_loop_config": {
      "epochs": 2,
      "num_stages": 2,
      "deepspeed_config": {
        "train_micro_batch_size_per_gpu": 4,
        "gradient_accumulation_steps": 2,
        "optimizer": {
          "type": "Adam",
          "params": {
            "lr": 0.001,
            "weight_decay": 0.01
```



# Resources



- [Ray Documentation](#)
- [Data-Parallel Pytorch Example](#)
- [ZeRO-3 DeepSpeed Example](#)
- [Ray Paper \(2017\)](#)
- [Ownership Paper \(2021\)](#)
- [Ray 2.0 Architecture Whitepaper](#)
- [DeepSpeed](#)
- [OReilly - Learning Ray: Flexible Distributed Python for Machine Learning](#)