

Ray Framework

OM APURVBHAI TANK

McMaster University
Hamilton, ON, Canada
tanko1@mcmaster.ca

Abstract

This project explores the implementation of the Ray framework, an open-source distributed computing and parallel processing library designed for scalable and efficient execution of machine learning (ML) and deep learning (DL) models. Ray aimed at addressing the challenges of scaling computationally intensive Python applications, Ray offer an intuitive API that supports distributed data processing, model training and hyperparameter optimization. Ray allows absolute execution across heterogeneous multiple resources, including CPUs and GPUs.

This project leverages components like Ray Datasets, Ray Train, and Ray Tune, which showcases efficient data preprocessing, distributed training, and advanced hyperparameter optimization across real-world datasets. These implementations underscore Ray's robustness, resource efficiency, and seamless integration with popular machine learning and deep learning models. The project demonstrates the application of Ray's capabilities through practical examples involving XGBoost and Convolutional Neural Networks (CNNs). This work emphasizes Ray's potential to streamline large-scale ML models, paving the way for scalable and resilient computational solutions.

1 Introduction

In recent years, the rapid advancements in artificial intelligence (AI) and machine learning (ML) have profoundly transformed various aspects of our lives. AI-powered applications now drive innovation in industries ranging from healthcare to finance, transportation to entertainment. Machine learning, a core subset of AI, enables systems to learn from data and make intelligent predictions or decisions from personalized recommendations on streaming platforms to autonomous vehicles navigating our roads, ML applications are reshaping how we interact with technology. Due to these advancements in machine learning has led to growing demand for frameworks that can handle large-scale data, resources and computational workloads efficiently.

Ray framework is an open-source, distributed computing library specifically designed to simplify the scaling of Python-based applications in fields like machine learning, deep learning, and data science community. Unlike traditional distributed systems, which often require extensive setup and specialized expertise, Ray provides a user-friendly interface for data scientists and engineers. Its primary objective is to make distributed computing accessible and effective, ensuring that even complex ML workloads can be scaled with minimal effort.

Implementing Ray in machine learning and deep learning projects enhances performance by distributing compute-intensive tasks across multiple CPUs, GPUs, or nodes, ensuring efficient utilization of resources. Moreover, the framework's resilience against node failures and dynamic resource scheduling minimizes downtime

and operational costs. Ray's ability to integrate with popular ML libraries like TensorFlow, PyTorch, and XGBoost further extends and making it a go-to solution for scaling Python-based models.

In conclusion, as AI and ML continue to play an increasingly central role in our lives, frameworks like Ray are crucial for harnessing their full potential. By providing a scalable, efficient, and user-friendly platform for distributed computing, Ray empowers organizations to innovate faster and more effectively.

2 Problem Statement

The increasing complexity and scale of machine learning (ML) models have made efficient resource utilization and execution time a critical focus in the field of artificial intelligence (AI). While many frameworks provide robust solutions for developing ML models, they often fall short when it comes to handling distributed workloads and optimizing computational efficiency. This creates significant challenges, especially for organizations and researchers seeking to scale their models to meet growing demands without overextending computational resources.

This project thus addresses critical questions that how can distributed frameworks like Ray optimize resource utilization and execution time in ML workflows. What are the limitations of distributed systems for smaller datasets, and how do these limitations contrast with their benefits for large, computationally intensive tasks? How can advanced hyperparameter tuning techniques like PBT contribute to both resource utilization and computational efficiency?

By examining these challenges, the project demonstrates the strengths and limitations of the Ray framework in diverse ML tasks. It underscores the need for distributed solutions that can adapt to varying data scales and computational requirements, offering valuable insights into how modern ML workflows can achieve scalability and efficiency.

In this project, two distinct machine learning models were explored to evaluate the challenges of scalability, resource optimization, and execution time in distributed systems. The first involves the implementation of an XGBoost-based classification model trained on a breast cancer dataset. While the dataset is relatively small, the focus was on assessing how distributed frameworks like Ray can improve resource utilization and streamline the execution process. The second model centers around a Convolutional Neural Network (CNN) applied to the CIFAR-10 dataset, a larger dataset used for image classification tasks. The CNN model required extensive hyperparameter tuning to optimize its performance, which was managed using Ray Tune's Population-Based Training (PBT) strategy.

3 Motivation

The motivation for this project stems from the growing complexity of modern machine learning models and the challenges of optimizing resource utilization and execution efficiency. As datasets grow larger due to data-driven world, machine learning models become more computationally intensive. This creates an urgent need for a framework that not only facilitates distributed computation but also integrates seamlessly across the entire lifecycle of machine learning models, from data preprocessing to hyperparameter tuning and deployment.

One of the key reasons to focus on the Ray framework is its ability to bridge the gap between research and production. Many AI models developed in experimental environments fail to scale when deployed in real-world scenarios due to resource constraints or system incompatibilities. Ray addresses these issues by providing robust tools for distributed model training, hyperparameter tuning, data preprocessing, and model serving. Furthermore, the framework's versatility and fault-tolerant design make it ideal for modern ML workflows. Its components, such as Ray Tune for hyperparameter optimization and Ray Train for distributed training, provide specialized support for critical ML tasks. The inclusion of Ray Datasets, capable of processing massive datasets across clusters, ensures integration with the data-intensive nature of machine learning.

This project explores the potential of Ray in real-world scenarios to evaluate how its capabilities compare to existing tools in terms of scalability, performance, and ease of use. By implementing two distinct use cases—a classification problem using XGBoost and a deep learning task involving a CNN on the CIFAR-10 dataset—the project aims to assess Ray's effectiveness in handling both structured data and complex neural network models.

The practical motivation lies in understanding how frameworks like Ray can drive innovation and efficiency in AI development. For instance, the need to optimize hyperparameter tuning in resource-intensive tasks, such as training CNNs, is critical for reducing time-to-insight and improving model performance. Similarly, investigating Ray's ability to manage distributed resources effectively, even for smaller datasets like the breast cancer dataset, provides insights into the framework's adaptability across different scales of workloads. By addressing these aspects, the project seeks to highlight the transformative impact of adopting distributed frameworks like Ray in modern ML models. The findings aim to demonstrate how Ray can overcome the limitations of traditional tools, offering a scalable, efficient, and integrated solution that empowers developers to focus on innovation rather than infrastructure challenges.

4 Related Work

Distributed computing has long been a cornerstone for scaling machine learning models, enabling parallel execution and efficient utilization of computational resources. While several frameworks like Dask, PyTorch, and Ray have gained prominence for specific aspects of machine learning, other paradigms such as MPI (Message Passing Interface), also play a critical role in distributed computing. This section explores the unique contributions and limitations of these frameworks and highlights ongoing advancements in the field.

4.1 MPI

MPI is one of the oldest and most robust paradigms for distributed computing, offering low-level control over parallel execution. Unlike frameworks like Dask, PyTorch, or Ray, which abstract many complexities, MPI provides direct access to message-passing mechanisms, allowing fine-grained control over data distribution and communication across nodes. While this makes MPI highly efficient for scientific computing and custom distributed algorithms, it requires significant expertise to implement and manage. MPI does not natively support modern machine learning models and lacks higher-level abstractions for tasks like hyperparameter tuning, distributed training, or fault tolerance. Its low-level nature makes it more suitable for specialized use cases rather than general-purpose machine learning models.

4.2 PyTorch

PyTorch is a widely-adopted framework for deep learning that excels in its flexibility and dynamic computation graph capabilities. PyTorch supports rapid prototyping of deep neural networks. Its Distributed Data Parallel (DDP) module allows for efficient scaling of deep learning models across GPUs and nodes, making it a reliable choice for training large neural networks. However, PyTorch's focus is largely on model training and batch inference. It lacks comprehensive support for data preprocessing, fault-tolerant task distribution, or hyperparameter optimization within the same framework. To address these needs, developers often resort to external tools like Optuna for hyperparameter tuning or Dask for preprocessing, leading to disjointed workflows. PyTorch, therefore, excels in research and training contexts but requires significant customization for production-level distributed machine learning pipelines.

4.3 Dask

Dask is a Python-native framework designed primarily for parallel and distributed data processing step only. It integrates seamlessly with existing Python libraries, making it an excellent choice for users looking to scale operations on structured data using tools like Pandas or NumPy. However, its architecture is inherently task-centric, focusing on breaking down computations into smaller subtasks distributed across a cluster. While effective for general-purpose data analytics, Dask does not offer specialized libraries for machine learning models. For example, while Dask-ML extends the framework's capabilities for ML tasks, it lacks native support for distributed deep learning, hyperparameter optimization, or advanced model serving. This gap often necessitates integration with other tools, leading to more complex workflows and limited efficiency for end-to-end machine learning pipelines.

Although, all these frameworks and libraries provide distributed computing but didn't meet the expectations of distributed scaling, training and tuning at all stages of machine learning models. Ray represents a significant step forward in distributed computing for machine learning by handling all stages of machine learning model in a cohesive manner. Unlike MPI, which requires manual management of task distribution, or Dask, which focuses primarily on data processing, Ray provides dedicated libraries for every stage

of the ML model pipeline. Ray Train supports distributed training of machine learning models with built-in fault tolerance and dynamic resource allocation. Ray Tune streamlines hyperparameter optimization by integrating advanced search algorithms and schedulers, while Ray Serve enables scalable deployment of trained models. These capabilities make Ray a comprehensive choice for both research and production, capable of handling structured data tasks, deep learning, and hybrid workloads. However, Ray is relatively new compared to Dask or PyTorch, and its in still development phase, there is room for improvements in usability and performance for specific models.

5 Methodology

This project employed the Ray framework to demonstrate the efficiency of distributed computing in machine learning workflows. Two distinct models were implemented: an XGBoost-based classification model for the Breast Cancer dataset and a Convolutional Neural Network (CNN) for image classification using the CIFAR-10 dataset. The methodology explains the specific Ray libraries used, hardware requirements, dataset preparation, model training & tuning, and evaluation metrics for resource utilization and execution time. To achieve distributed training and hyperparameter optimization effectively, this project utilized two key libraries Ray Train and Ray Tune.

5.1 Ray Train

Ray Train is a distributed training library provided by the Ray framework, designed to simplify the dynamic scaling of machine learning (ML) training part across multiple devices such as CPUs, GPUs, or even nodes in a cluster. Its primary purpose is to enable efficient training of machine learning models by distributing workloads dynamically, ensuring optimal resource utilization.

The core advantage of Ray Train lies in its ability to abstract the complexities of distributed training. Typically, auto-scaling training models requires extensive setup, custom resource management, and manual synchronization across nodes. Ray Train eliminates these hurdles by providing a high-level API for defining and managing distributed training functions, coupled with fault-tolerant execution mechanisms. This allows developers to focus on designing models and improving performance. Figure 1 gives an idea about the components of the library.

5.1.1 Scaling Configuration: Scaling Configuration defines how resources are allocated for distributed training. Parameters like the number of workers (`num_workers`), GPU utilization (`use_gpu`), and resource constraints for each worker (`resources_per_worker`) are specified here. Resources constraints are considered for every worker present in the Ray cluster. For instance, in this project, XGBoost model utilized three CPU workers, while the CNN model leveraged 1 GPU and 2 CPU resources to accelerate training.

5.1.2 Trainer APIs (Training Functions): Trainers are wrapper classes around third-party training frameworks like XGBoost, Pytorch, and TensorFlow providing integration with core Ray actors (for distribution), Ray Tune, and Ray Datasets.

For example, XGBoostTrainer was used to manage distributed training of the XGBoost model in this project. These APIs simplify

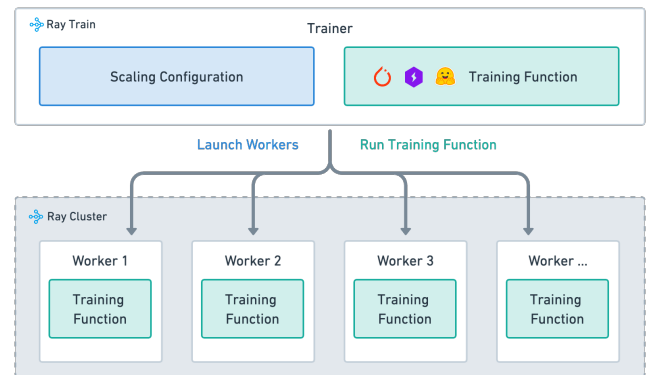


Figure 1: Components of Ray Train.

the setup process by handling data partitioning, task scheduling, and synchronization automatically.

A generic Trainer is a customisable abstraction that allows users to define a distributed training model tailored to their specific model and dataset. It requires to specify the training logic using a custom function. This flexibility is essential for scenarios where the default trainers do not meet the unique requirements of the model architecture or data processing pipeline. The custom training logic is encapsulated with a function, typically called (`train_loop_per_worker`). This function defines how each worker processes the data, trains the model, and communicates with other workers in the cluster.

5.2 Ray Tune

Ray Tune is a library within the Ray framework designed to automate hyperparameter optimization for machine learning (ML) models. Hyperparameter tuning is a critical step in ML model development, as it directly influences model performance by finding the best configurations for parameters like learning rate, batch size, dropouts and depth of trees. Ray Tune provides a highly flexible and efficient solution for conducting hyperparameter search at scale, capable of managing parallel trials, advanced search algorithms, and dynamic resource allocation.

At its core, Ray Tune streamlines the hyperparameter tuning process by enabling users to define search spaces, execute trials in parallel, and monitor performance metrics in real time. Ray Tune ensures that tuning workflows are not only faster but also more resource-efficient, making it an ideal choice for both research and production environments.

Key components of Ray Tune are as below in the image:

5.2.1 Search Spaces: Search Spaces specifies the range and sampling methods for hyperparameters. Users can define discrete or continuous ranges for parameters and specify how they are sampled—either through simple methods like grid search, random, normal distribution and choice of parameters. For example, in this project, the (`max_depth`) parameter of the XGBoost model was optimized using grid search, while the CNN used Population-Based Training (PBT) to adjust hyperparameters dynamically during training.

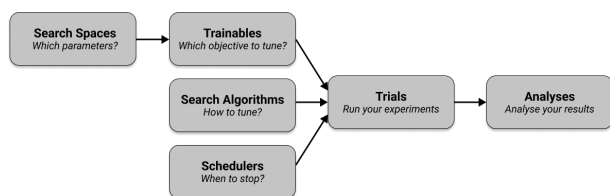


Figure 2: Components of Ray Tune.

5.2.2 Trainable: Trainable is the function or class that defines the objective function to be optimized. It encapsulates the training process, taking hyperparameters from the search space as input and returning performance metrics (e.g., accuracy, loss) for evaluation. Each trial corresponds to a specific set of hyperparameters, and the trainable function ensures that these configurations are tested under consistent conditions. For instance, in the CNN model implementation, the trainable function was responsible for setting up the model architecture, training it on a subset of the data, and evaluating its accuracy on the validation set. This modular approach allows users to integrate Ray Tune into virtually any ML model.

5.2.3 Search Algorithms: Ray Tune supports a variety of search algorithms to explore the hyperparameter space efficiently. These include traditional methods like grid search, where all possible combinations are tested, as well as more advanced algorithms such as Bayesian optimization. In this project, Population-Based Training (PBT) was employed for the CNN model. PBT combines hyperparameter search with training by periodically perturbing underperforming trials and adopting hyperparameters from better-performing ones.

5.2.4 Schedulers (Optional): Schedulers are responsible for managing the execution of trials, deciding when to pause, stop, or prioritize specific configurations based on their performance. Schedulers are particularly useful in large-scale experiments where resources are limited and many trials need to be evaluated. In this project, Population Based Training (PBT) acted as a scheduler, dynamically allocating resources to the best-performing configurations and adjusting their hyperparameters to improve results further.

5.2.5 Trials: Each trial represents an independent execution of the trainable function with a unique set of hyperparameters. Trials are managed by Ray Tune, which executes them in parallel across available resources. The number of trials and their configurations are determined by the search algorithm and scheduler. For instance, in the CNN implementation, multiple trials were run simultaneously, each testing a different combination of learning rate, dropout rate, and batch size. Each trial details can be found from the logs folder created by ray during runtime.

5.2.6 Checkpoints & Analysis: Checkpoints save the state of a trial at regular intervals, allowing interrupted trials to resume without losing progress. Additionally, checkpoints enable users to analyze the training dynamics of each trial and compare configurations effectively. The CheckpointConfig component of Ray Tune was used in the CNN implementation to save model states after every

few iterations, ensuring that training could be resumed or evaluated at any point. Ray Tune includes tools to aggregate and visualize the results of all trials, providing insights into the best hyperparameter configurations and their impact on model performance. Metrics such as loss, accuracy, and resource utilization are logged during each trial and can be visualized using built-in dashboards. These insights are invaluable for understanding the optimization process and refining the model further.

5.3 XGBoost Implementation

The XGBoost model was implemented to address a binary classification problem using the Breast Cancer dataset, which was sourced from the Kaggle repository. This structured dataset comprises diagnostic features of tumors, categorized as benign or malignant. The hardware requirements for this model are 4 CPUs used for efficient parallel computation. The memory of 8GB RAM was used to handle the dataset and distributed training. Here, GPU support was optional and not utilized due to the small dataset size. Before starting the preprocessing and training of the model, we need to install the ray module from python and import the libraries that we want to use for parallel execution and distributed training.

Data preprocessing step was performed using Ray Datasets to leverage distributed data handling. The dataset was loaded in CSV format, normalized, and split into training (80%) and testing (20%) subsets using the `train_test_split` function. Irrelevant columns were removed, and categorical labels were encoded for binary classification. The training process employed Ray Train with the XGBoost-Trainer, distributing tasks across three CPU workers to optimize resource utilization. The XGBoost model was configured with the following settings:

- **Objective:** binary:logistic for binary classification.
- **Evaluation Metrics:** Log loss and error rate to evaluate performance.
- **Boosting Rounds:** 30 iterations for training.

To enhance model performance, Ray Tune was utilized for hyperparameter tuning. A grid search was performed over the `depth` of trees (`max_depth`) parameter, identifying the optimal value to minimize log loss. Tuner, RunConfig and TuneConfig modules managed the execution of parallel trials, disabled all callbacks including Tensorboard, and add custom logs directory which dynamically allocating resources for efficient evaluation.

The trained model was then evaluated on the test subset. Metrics such as accuracy and log loss were recorded. Due to the small dataset size, distributed training did not significantly reduce execution time or improve performance of the model. However, the efficient resource utilization provided by Ray established a scalable foundation for larger datasets.

5.4 CNN Model Implementation

The Convolutional Neural Network (CNN) was designed to classify images from the CIFAR-10 dataset, a benchmark dataset for image classification comprising 60,000 32x32 RGB images divided into 10 categories. The hardware requirements for this model are 2 CPUs and 1 GPU for faster training due to the computational intensity of the CNN architecture. Before starting the training of the model, we need to install the ray module from python and import the libraries

like ray tune, train and data that we will use for distributed training across multiple nodes. The memory of 32GB RAM was used to handle the dataset and distributed training.

Preprocessing the CIFAR-10 dataset was a crucial step to prepare the data for efficient training and ensure model robustness. The CIFAR-10 dataset consists of 60,000 with 50,000 images used for training and 10,000 for testing. Pixel values were scaled to the range [0, 1] to ensure numerical stability during model training. Labels were converted into one-hot encoded vectors to align with the categorical cross-entropy loss function used during training. Ray Datasets facilitated the distributed execution of preprocessing tasks, allowing the large dataset to be efficiently partitioned and processed across multiple workers.

The CNN architecture consists of multiple convolutional layers with ReLU activation functions, followed by pooling layers for dimensionality reduction. Fully connected layers with dropout regularization were used at the end for classification. The training process was distributed across 1 GPU and 2 CPUs using Ray Train, which efficiently allocated resources to parallelize computations. The following configurations were used during training:

- **Model:** Adam optimizer which combines momentum and adaptive learning rate technique.
- **Epochs and Early Stopping:** The model was trained for up to 25 epochs. Early stopping criteria were defined based on stability in validation metrics and total execution time
- **CheckpointConfig:** Periodic checkpoints saved the state of model every 20 iterations using Ray Tune library. And (nums_to_keep) parameter limits the number of most recent checkpoints retained for trial.

Hyperparameter optimization was critical for improving the training efficiency of the CNN model. Ray Tune’s Population-Based Training (PBT) was employed to dynamically adjust hyperparameters during training, leveraging real-time feedback from trial evaluations. The initial values for hyperparameters were defined like learning rate initialized at 0.0001 and 0.01 adjusted dynamically during training, batch size tune within a range of [32, 128] and dropout rate ranged between 0.25 and 0.5 to balance regularization. Population-Based Training (PBT) iteratively identified poorly performing trials based on metrics like validation accuracy and replaced their hyperparameters with those of better-performing trials.

Multiple trials (16 trials) were executed in parallel, with Ray Tune managing resource allocation and prioritization. Metrics such as execution time, resource utilization, and validation accuracy were logged for each trial. The final model configuration was selected based on the trial that minimized training duration while maintaining acceptable performance metrics.

6 Results

The results of this project demonstrate the impact of distributed training and hyperparameter optimization on reducing execution time for XGBoost and CNN models using the Ray framework. Both the XGBoost and CNN models were evaluated based on their training time duration, resource utilization, and the scalability of the implemented solutions.

6.1 XGBoost Model Results

From Table 1, we can see results on the Breast Cancer dataset, showed an increase in execution time when using Ray, despite the implementation of distributed training and hyperparameter tuning. Without Ray, the training process took approximately 2 minutes, as the dataset was small and easily managed on a single node. The simplicity of the dataset and the lack of need for complex computation made traditional training more time-efficient. However, with Ray, the execution time increased to approximately 2.5 minutes due to the overhead of initializing the Ray cluster and distributing workloads across CPUs. The additional steps involved in Ray’s distributed processing, including partitioning the dataset and managing inter-worker communication can be considered for the results too. These steps, while beneficial for larger datasets, added unnecessary overhead for the relatively small Breast Cancer dataset, where the computation was completed quickly without the need for parallelization and distributed training.

In terms of resource utilization, both implementation used 4 CPUs, but Ray demonstrated efficient allocation and scaling of these resources. The implementation of hyperparameter tuning using Ray Tune allowed for the evaluation of five trials, optimizing the max_depth parameter. This added an extra layer of computational complexity, contributing to the slight increase in execution time. Although the model’s accuracy increased slightly from 95% to 96% with tuning, this was not a primary focus of the project. The results underscore the need to evaluate the trade-offs of distributed frameworks like Ray for small-scale datasets, where their benefits may not be fully realized.

Tune Status

Current time: 2024-12-16 19:34:23

Running for: 00:02:35.12

Memory: 6.8/7.7 GB

System Info

Using FIFO scheduling algorithm.

Logical resource usage: 4.0/4 CPUs, 0/0 GPUs

Trial Status

Trial name	status	loc	params	max_depth	iter	total time (s)	train-logsloss	train-error	valid-logsloss
XGBoostTrainer_aesid_00000	TERMINATED	127.0.0.1:20008	50	31	33.872	0.0108854	0	0.047455	
XGBoostTrainer_aesid_00001	TERMINATED	127.0.0.1:20624	5	31	22.8771	0.0103782	0	0.0403975	
XGBoostTrainer_aesid_00002	TERMINATED	127.0.0.1:15488	1	31	18.6195	0.0795789	0.0154525	0.0848822	
XGBoostTrainer_aesid_00003	TERMINATED	127.0.0.1:19428	1	31	16.2519	0.0795789	0.0154525	0.0848822	
XGBoostTrainer_aesid_00004	TERMINATED	127.0.0.1:9996	100	31	17.2518	0.0108854	0	0.047455	

Figure 3: XGBoost - Final Trial Results & Resources Utilization.

6.2 CNN Model Results

From Table 2, CNN model trained on the CIFAR-10 dataset, demonstrated significant reductions in execution time with Ray implementation. Without Ray, the model required approximately 40 minutes to complete training for 25 epochs without GPU. The computational intensity of deep learning, coupled with the large size of the CIFAR-10 dataset, created an exception in the traditional setup. In contrast, with Ray, the execution time dropped to approximately 25 minutes for total of 400 epochs, leveraging distributed training and hyperparameter optimization through Population-Based Training (PBT).

The efficiency gain was achieved by utilizing Ray’s ability to distribute training tasks across multiple workers, with one GPU and two CPUs effectively handling the workload. PBT dynamically adjusted hyperparameters such as batch size, learning rate, and

Table 1: Results Comparison: XGBoost Model (Breast Cancer dataset)

Parameters	Without Ray Implementation	With Ray Implementation
Execution Time	Approx. 2 minutes without Tuning	Approx. 2.5 minutes with Tuning
Number of Trials	1	5
Accuracy	95%	96%
Resources Utilization	4 CPUs	4 CPUs & 0 GPUs

Table 2: Results Comparison: CNN Model (CIFAR-10 dataset)

Parameters	Without Ray Implementation	With Ray Implementation
Execution Time	Approx. 40 minutes	Approx. 25 minutes with PBT tuning
Trials	1	16
Epochs	25	400 (total)
Resources Utilization	2 CPUs	2 CPUs & 1 GPUs

dropout rate during training, ensuring resource prioritization for the most promising trials. The implementation evaluated 16 hyperparameter trials, totaling 400 epochs across all trials, without a significant increase in overall execution time. This dynamic tuning enabled the model to converge faster, achieving less computation time while maintaining efficient resource utilization.

One of the key strengths of the Ray framework, highlighted in this project, is its built-in fault tolerance system, which ensures that even in the event of a node or cluster failure, the distributed setup can automatically restart from the last saved checkpoint. The significant reduction in execution time highlights Ray's scalability and suitability for large, computationally demanding tasks, showcasing its advantages over traditional single-node training.

```

Trial status: 16 TERMINATED
Current time: 2024-12-11 05:44:18. Total running time: 27min 54s
Logical resource usage: 2.0/2 CPUs, 1.0/1 GPUs (0.0/1.0 accelerator_type:T4)
Current best trial: 109d6_00013 with mean_accuracy=0.109375 and params={'epochs': 25, 'batch_size': 64, 'learning_rate': 0.0001}

```

Trial name	status	learning_rate	dropout	acc	iter	total time (s)
Cifar10Model_109d6_00000	TERMINATED	0.0001	0.25	0.09375	30	105.455
Cifar10Model_109d6_00001	TERMINATED	0.0001	0.5	0.09375	30	99.8423
Cifar10Model_109d6_00002	TERMINATED	0.01	0.25	0.078125	30	102.604
Cifar10Model_109d6_00003	TERMINATED	0.01	0.5	0.078125	30	99.8899
Cifar10Model_109d6_00004	TERMINATED	0.0001	0.25	0.078125	30	101.015
Cifar10Model_109d6_00005	TERMINATED	0.0001	0.5	0.078125	30	106.417
Cifar10Model_109d6_00006	TERMINATED	0.01	0.25	0.0859375	30	102.299
Cifar10Model_109d6_00007	TERMINATED	0.01	0.5	0.0859375	30	104.154
Cifar10Model_109d6_00008	TERMINATED	0.0001	0.25	0.0859375	30	101.566
Cifar10Model_109d6_00009	TERMINATED	0.0001	0.5	0.0859375	30	104.369
Cifar10Model_109d6_00010	TERMINATED	0.01	0.25	0.0859375	30	104.354
Cifar10Model_109d6_00011	TERMINATED	0.01	0.5	0.0859375	30	105.956
Cifar10Model_109d6_00012	TERMINATED	0.0001	0.25	0.101562	30	103.912
Cifar10Model_109d6_00013	TERMINATED	0.0001	0.5	0.109375	30	105.862
Cifar10Model_109d6_00014	TERMINATED	0.01	0.25	0.109375	30	102.69
Cifar10Model_109d6_00015	TERMINATED	0.01	0.5	0.109375	30	106.756

```

Best hyperparameters found were: {'epochs': 25, 'batch_size': 64, 'learning_rate': 0.0001, 'decay': 1e-06, 'dropout': 0.0}
Total execution time: 1074.42 seconds
(Cifar10Model_109d6_00013) Checkpoint successfully created at: Checkpoint(filesystem-local, path=/root/.ray_results/pbt_cifar10)

```

Figure 4: CNN - Final Trial Results & Resources Utilization.

7 Limitations

It is important to note that assigning more resources than what is physically available on the system can lead to potential issues. In such cases, Ray may enter an infinite loop while attempting to allocate non-existent resources, as it continuously searches for workers or hardware that are unavailable. This underscores the necessity of accurately defining the resource configuration during the setup phase, aligning the scaling parameters (e.g., `num_workers` or `use_gpu`) with the actual hardware capabilities.

For small datasets, such as the Breast Cancer dataset used in this project, this overhead can outweigh the benefits of distributed training. This limitation makes Ray less suitable for tasks involving small or computationally lightweight datasets, where single-node execution is often more time-efficient. By highlighting these challenges, the project acknowledges the contexts where Ray might not be the most suitable choice, providing a balanced and realistic assessment of its applications.

8 Conclusion

In conclusion, Ray proved to be a powerful tool for optimizing machine learning workflows, particularly for tasks involving large datasets and complex computations. Its integration of distributed training and hyperparameter tuning within a unified framework simplifies the development process and ensures efficient use of computational resources. This project demonstrated Ray's potential to streamline machine learning models at every stage, paving the way for scalable and robust solutions in modern AI applications.

9 Future Emerging Frameworks

Several emerging frameworks and research efforts aim to push beyond Ray's capabilities. For example, Horovod, developed by Uber, specializes in distributed deep learning and integrates tightly with TensorFlow, PyTorch, and Keras. Its algorithm optimizes communication between nodes, making it particularly effective for GPU-heavy workloads. Unlike Ray, which focuses on general-purpose distributed computing, Horovod is highly specialized for scaling deep learning models, offering unparalleled efficiency.

Another noteworthy advancement is TensorFlow Distributed (TFD), a native solution within TensorFlow for distributed model training. TFD provides robust support for both synchronous and asynchronous training, with built-in fault tolerance and resource scheduling. It integrates directly with TensorFlow's ecosystem, offering a cohesive solution for TensorFlow users but lacking the broader versatility of Ray for tasks outside deep learning.