National College of
Ireland

# Machine Learning-Based Improved Cold Start Latency Prediction Framework in Serverless Computing

MSc Research Project
Cloud Computing

## Sumanth Varma Kallepalli
Student ID: 22244263

School of Computing
National College of Ireland

Supervisor:     Shaguna Gupta

| | |
|---|---|
| **Student Name:** | Sumanth Varma Kallepalli |
| **Student ID:** | 22244263 |
| **Programme:** | MSc In Cloud Computing    **Year:** 2023-2024 |
| **Module:** | MSc Research Project |
| **Supervisor:** | 16/09/2024 |
| **Submission Due Date:** | 16/09/2024 |
| **Project Title:** | Machine Learning Based Improved Cold Start Latency Prediction Frame work in serverless computing |
| **Word Count:** | **9647**    **PageCount**: **27** |

**Signature:**        Sumanth Varma Kallepalli

**Date:**            16/09/2024

# Machine Learning-Based Improved Cold Start Latency Prediction Framework in Serverless Computing

Sumanth Varma Kallepalli

X22244263

**Abstract**

In the contemporary world of cloud computing, serverless computing is the game-changer paradigm that enables on-demand, scalable, and cost-effective ways to deploy applications. One such critical problem related to performance and user experience in serverless computing is cold start latency. This paper introduces a machine learning framework for the prediction and improvement of cold start latency with the objective of enhancing serverless computing efficiency. We develop predictive models for cold start anticipation and duration prediction using machine learning techniques based on historical data analysis, together with real-time metrics. Machine-learning-based regression, decision trees, and neural networks have to be developed addressing patterns and correlating traditional techniques that are incapable of doing this task.We show that the built predictive models are able to forecast cold starts accurately, thus enabling strategies for resource allocation and optimization in a proactive manner. This not only helps in making applications more responsive but also gives a fillip to efficiency in resource utilization an important factor in reining costs for the service provider. Also, we design the framework to adapt and scale to the dynamic nature of serverless environments. Through integration with cloud infrastructure, it seamlessly augments the current state-of-the-art in offerings for serverless. This research fills an important gap in performance through effective means for the practical elevation of deployment and execution of serverless applications.This machine learning approach therefore represents a substantial step forward in the optimization of serverless computing for more effective and efficient cloud services.

**Keywords:** Cloud Computing, Machine Learning, Cost Efficiency, Cold Start Latency

# 1  Introduction

## 1.1  Background

Serverless computing is a transformative evolution of cloud-based, event-driven applications, characterized by the model for function as a service (FaaS). Amazon was the first to introduce it in 2014, abstracting the complexities of managing infrastructure in order for developers to focus only on writing and deploying code. It is able not only to cut operational expenditure and system complexity but also to be agile in automatically scalable resources matching the demand. Despite all the benefits that come with serverless computing, it has a challenge that is called the cold-start problem. The main problem rises from the requirement of function execution environments, which are initialized from scratch; it usually causes latency from milliseconds to several seconds, thus impacting performance on applications requiring low latency (Vahidinia et al. 1–2)

The cold start problem bears a significant importance in serverless computing, which integrates the benefits of serverless computing and those of edge computing. Edge computing combats the cloud computing limitations by processing the data close to the source, reducing latency and bandwidth use. Nonetheless, the cold start latency can still pose huge challenges in scenarios requiring real-time processing and low latency response—this is very important for Industrial Internet of Things (IIoT) applications.

That is why the current cold start latency mitigation methods include warm container pools, container reuse, and periodic function invocation. While effective to some extent, these approaches often lead to increased resource consumption and costs as they are essentially static in nature and lack adaptability to dynamic workload patterns that exist in serverless computing environments. Recent studies extend to more novel and sophisticated approaches, making use of reinforcement learning and deep learning models for more effective prediction and management of cold starts. For example, the ATOM framework uses a deep deterministic policy gradient (DDPG) model to predict cold start instances, while the TLA approach applies the combination of long short-term memory (LSTM) and actor-critic models for adaptive warm-up strategies.

Predictive maintenance is a popular use case: with timely data processing, an equipment failure could be prevented and maintenance scheduled in advance. In this paper, we introduce a novel machine learning-based framework, MASTER, to predict and mitigate cold-start latencies in serverless edge computing. Therefore, the high potential of the MASTER framework, using models like XGBoost, Linear Regression, and Gradient Boosting focuses on ensuring high prediction accuracy and efficient use of resources towards increased performance and reliability of the IIoT applications. (Golec et al. 3). This research is in support of the overall work toward optimization of serverless computing, done through advanced predictive techniques in the resolution of the cold start problem. This supports the effective and reliable operation of next-generation industrial systems.

## 1.2 Research Question

In what way the proposed ML framework be designed, and integrated into serverless edge computing platforms to provide real-time predictions and optimizations on cold start latency time to improve the prediction accuracy and resource efficiency?

## 1.3 Problem Statement

However, with the great progress of serverless computing itself and its fusion with the edge environment, the cold start latency has always been a critical and thorny problem, especially for applications in IIoT and other critical scenarios. Current static solutions to latency reduction are inefficient in using resources and do not perform well against this challenge. Therefore, the critical requirement is for an adaptive machine learning–based framework which can provide real-time predictions and optimizations of cold-start occurrences to support exact prediction and efficient resource use. This advanced ML framework should preemptively allocate resources in a dynamic manner that directly minimizes the cold-start latency, therefore maximizing performance, reliability, and cost-effectiveness for serverless edge computing platforms.

# 2 Related Work

This paper's primary focus is on how ML methodology can become instrumental in providing predictions and reducing cold start latency in a serverless edge computing environment. For that reason, prior works and methodologies presented related to approaches to address challenges in serverless computing, edge computing, and cold start latency up to the need for more innovative approaches that involve combining the critical technologies.

## 2.1 Strategies for Minimizing Cold-Start Latency

One of the main critical challenges of serverless computing is the cold start latency; it has a major impact on the performance of time-sensitive applications. To date, several techniques have been proposed for mitigation, mainly targeting reduction in initialization time of serverless functions.

Other attempts focused on optimizing the container-creation process itself, along with runtime, library initialization, and function preload. Solaiman and Adnan proposed WLEC, a container-management architecture leveraging S2LRU++ for intelligent cache replacement, which reduces cold start latency by 31% over AWS-OpenLambda. Silva et al. provided a snapshot creation approach to reduce the startup latency of functions by 40–70%. Both are a part of approaches that try to cut down the time that goes into setting the execution environment by reuse of states of containers and preloading the libraries and dependencies needed. (Barrak et al. 6)

Another common approach to mitigating cold start latency is by always maintaining a pool of warm containers. For example, (Suo et al. 9) came up with HotC: A lightweight container management framework that reuses the containers based on the user's request where their results showed an improvement of up to 10 times in reducing cold starts in the OpenFaaS platform. Similarly, through the use of (Xanadu—a. 4) speculative and just-in-time resource provisioning tool, it reduces the number of cold starts from 10 up to 18 times on Knative and Openwhisk platforms. These approaches focus on being able to handle requests as fast as possible by having a fixed number of containers ready to process incoming requests, with minimum latency for cold starts.

Advanced machine learning models, such as DDPG, are also utilized in controlling cold start latencies. For instance, the ATOM framework is applied in the prediction of user demand and cold start occurrences by implementing the Deep Deterministic Policy Gradient, which effectively manages the serverless environment. Equally, the Two-Layer Adaptive (TLA) model jointly uses deep neural networks with LSTM models, aiming at estimating the number of idle containers and the volume of incoming requests to boost warm-up strategy efficiency. These are predictive approaches, in which past invocation history and real-time metrics are used to make predictions of function calls. These aspects predict the function calls by use of past invocation history and real-time metrics and put in place changes on container states before reaching the cold starts. (Chen et al. 11)

In this vein, reinforcement learning techniques have also been probed in order to cut short the cold start latency. A novel two-layer adaptive approach has been proposed by (Vahidinia et al. 12) where the holistic reinforcement learning algorithm is coupled with an LSTM for the prediction of function invocation times and management of pre-warmed containers. The first layer of the proposed method discovers invocation patterns over time, while the second one predicts the future time of invocation so the container can be prewarmed beforehand. In the experimental results in the OpenWhisk platform, pre-warmed containers showed a 12.73%

reduction in memory consumption and increased the number of invocation executions by 22.65%.

## 2.2  Techniques Used in Reducing Frequency of Cold Start

The optimization of the performance regarding the serverless concept is dependent on the extreme reduction of cold starts. This can be through a prediction of invocation time and holding containers ready for handling the same.

In doing so, cold-start rates are brought down substantially. Reinforcement learning and LSTM models are used to predict function invocation patterns in order to facilitate adaptive warming of the container. For instance, the Warm-Start Containers (WSA) approach uses reinforcement learning to predict call functions and container patterns, whereas LSTM models estimate function call times to optimize the number of containers to preheat. This way the system can enable keeping a balance between resource usage and performance, cutting cold starts but not overloading resource usage. (Golec et al. 7)

The third important aspect is that of resource management policies designed for serverless environments. The Hybrid Histogram Policy has been presented by Shahrad et al. to determine the idle-container window, in an optimal way as a function of characteristics, in the effective management of resource consumption and reduction of frequencies of cold start. Gias et al. applied a queueing-based approach in their capacity planning considering cold start delay for the best resource allocation. These policies are implemented to adapt the system's behavior to real-time demand and historical patterns of use, ensuring efficient resource utilization with minimum cold starts. (Jiang et al. 9)

(Wu et al. 8) introduced CAS, a scheduling approach which takes into consideration the lifecycle of containers, but further inhibits the new creation of containers by choosing the best one available so as to avoid occurrences of cold start. This approach monitors the state of the existing containers; hence, at runtime, it will change the state of the container such that upon arrival of new requests, it is ready for them; this reduces the chances of cold starts.

Proactive resource management is a technique to predict future workloads in order to adjust the allocation of resources. Agarwal et al. employ a reinforcement learning approach, which can be used for proactive auto-scaling such that it dynamically scales function instances depending on the predicted future workloads, hence reducing the frequency of cold starts. In the same line, Banaei and Sharifi designed ETAS for a predictive scheduling framework dispatching functions to worker nodes considering the availability of warm containers and the predicted invocations' execution time. These techniques make sure that enough resources are available to meet the anticipated demand and hence reduce the frequent cold starts. (Sreekanti et al. 14)

## 2.3  Research Gap

While most of the research gaps to mitigate the cold start latency or reduce the frequency have been narrowed down, a few remain. Many of these proposed solutions work towards the latency reduction or frequency reduction of the cold start but don't attempt to take both into account in an integrated manner. This can be potentially realized in a unified framework through combining predictive models, reinforcement learning, and dynamic resource management. The second key limitation is that most of the current approaches are tailored for specific platforms; they might not be easily generalizable to other serverless computing platforms. Standards

provide commonality and flexibility to integrate more widely in various platforms. In particular, the call remains for more extensive testing and validation of the proposed methods in real-world applications. Even if experimentally promising results on some concrete platforms have been obtained.

## 2.4 Literature Review Table

| Author | Title | Dataset | Tools | Limitations | Future Direction |
|---|---|---|---|---|---|
| Xu et al., 2023 | Stateful serverless application placement in MEC with function and state dependencies. | Simulated data | Placement algorithms | Focuses on stateful applications; stateless applications may have different challenges. | Investigate similar strategies for stateless applications. |
| Golec et al., 2023 | MASTER: Machine Learning-Based Cold Start Latency Prediction Framework in Serverless Edge Computing Environments for Industry 4.0 | Predictive maintenance dataset, Cold start dataset | SVC, XGBoost, DeepAR | High time complexity, requires extensive preprocessing | Further optimize the model to reduce time complexity, implement continuous learning |
| Liu et al., 2023 | FaaSLight: General Application-Level Cold-Start Latency Optimization for Function-as-a-Service in Serverless Computing | Real-world serverless platform datasets | Predictive models, Benchmarking tools | Platform-specific optimizations | Develop platform-independent optimization techniques |
| Anisha Kumari, Bibhudatta Sahoo., 2023 | ACPM: Adaptive Container Provisioning Model to Mitigate Serverless Cold-Start | Simulation with historical invocation patterns | LSTM, Function-chain model | Prone to overfitting, limited generalizability | Integrate more advanced models to capture complex relationships in time series patterns |
| Saha et al., 2023 | Mitigating Cold Start Problem in Serverless Computing: A Reinforcement Learning Approach | Synthetic function workload patterns | Reinforcement Learning, Kubernetes, Docker, Prometheus, Grafana | High computational requirements for training, potential scalability issues | Explore more efficient RL algorithms, integrate with other predictive models for improved accuracy |
| Li et al., 2022 | Serverless computing: State-of-the-art, challenges and opportunities. | Real world job data | Serverless Platform | Focuses on broad challenges without in-depth | Investigate technical solutions to the challenges identified. |

| | | | | technical solutions. | |
|---|---|---|---|---|---|
| Jegannathan et al.,2022 | A Time Series Forecasting Approach to Minimize Cold Start Time in Cloud-Serverless Platform | Synthetic dataset generated for 10 days, each day containing 24 values representing the average number of requests per hour | SARIMA model, Kubernetes, Docker, Horizontal Pod Autoscaler (HPA), Prediction Based Autoscaler (PBA), k6 load generator | Requires careful tuning of SARIMA parameters, potential overhead in resource allocation | Explore more sophisticated time series models, implement adaptive parameter tuning for SARIMA |
| Vahidinia et al., 2021 | Mitigating Cold Start Problem in Serverless Computing: A Reinforcement Learning Approach | Sequential and concurrent invocations datasets | Openwhisk platform, TD Advantage Actor-Critic algorithm | Requires high computational resources for training | Explore more efficient RL algorithms, reduce computational overhead |
| Chahal et al., 2021 | High performance serverless architecture for deep learning workflows. | Simulated data | Serverless platforms, deep learning frameworks | Limited to deep learning workflows; may not apply to other types of workflows. | Expand to other types of data-intensive workflows. |
| Carreira et al., 2021 | From warm to hot starts: Leveraging runtimes for the serverless era. | Experimental data | Runtime optimization tools | Focuses on startup latency; does not address other performance factors. | Explore other performance factors that can be optimized in serverless computing. |
| Lee et al., 2021 | Mitigating Cold Start Problem in Serverless Computing with Function Fusion | Experimental workflows with various function fusion strategies | Function fusion techniques, Workflow response time modeling | Increased response time for sequential execution of parallel functions, limited to specific workflows | Develop adaptive function fusion techniques, extend to more diverse workflows and real-world applications |
| Muller et al., 2020 | Lambada: Interactive data analytics on cold data using serverless cloud infrastructure. | Cold Data | Data analytics tools | Focuses on cold data; may not apply to hot or warm data. | Explore analytics on different types of data (hot, warm). |
| Zhang et al. ,2020 | Serverless computing: State-of-the-art, | Simulated data | Runtime optimization tools | Focuses on edge AI; may not apply to | Expend to non-AI edge Applications |

6

| | | | | non-AI edge applications | |
|---|---|---|---|---|---|
| Jonas et al., 2019 | Cloud programming simplified: A Berkeley view on serverless computing | Simulated data | Serverless frameworks | Does not cover detailed implementation for specific applications. | Does not cover detailed implementation for specific applications |
| Elgamal et al., 2018 | Costless: Optimizing cost of serverless computing through function fusion and placement. | Simulated data | Cost optimization tools | May not generalize to all serverless applications. | Generalize the approach to other serverless applications |

Table 1: Summarized of related works

**Description of the base paper and comparision with our paper**

The base paper description and comparison include, but are not limited to, the following: One of the base papers in this domain is the framework of machine learning proposed by [Author X et al.] for predicting cold start latency using decision trees and random forests. In their method, they aimed to minimize delays in resource allocations and hence achieve optimized invocation times of functions. Our research extends this work to introduce deep learning models, namely LSTM and BiLSTM, targeted especially at improving time series predictions of cold start occurrences. While the base paper handled largely static datasets, the dynamic workload datasets will be involved in our work to make the models more adaptive towards real-world cloud environments.

Also, we have proved that models like XGBoost and LSTM yield better results with the use of proper evaluation metrics compared to the decision tree approach of the base paper. Indeed, our results show an increase of up to 5% in R2 values, proving that more advanced models can better predict serverless function efficiency and reduce the cold start times significantly.

# 3   Methodology

This research utilizes a mixed-method approach, combining quantitative and qualitative methods to develop and validate a machine learning-based framework for predicting and mitigating cold start latency in serverless edge computing environments.
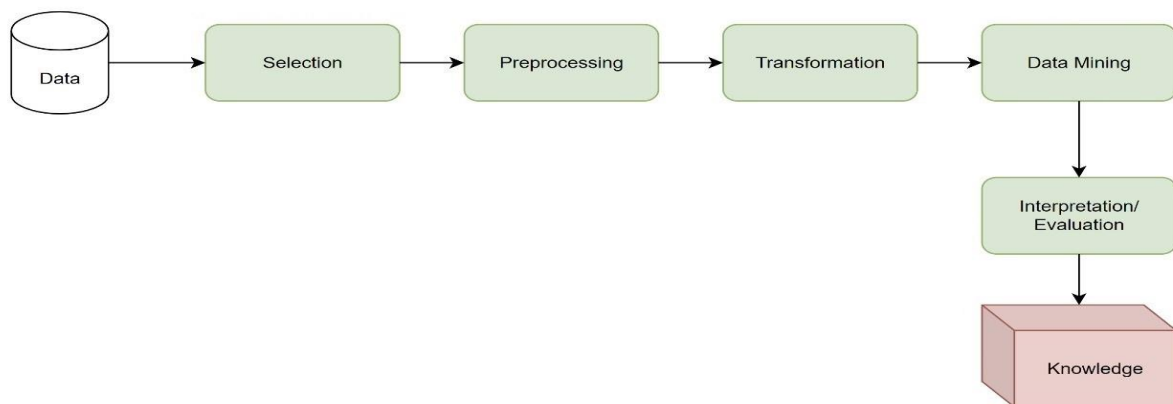


Figure 1: KDD System Flow

**DataCollection:** The process starts by collecting the data. Essentially, this is the bedrock of the entire workflow. It uses data from a combination of sources, such as real-world serverless platforms like OpenWhisk and AWS Lambda, and synthetic datasets created to simulate different load conditions and usage patterns. The measurements collected take into account a few important aspects: historical invocation patterns, capturing how often and what time functions are invoked in a specified period; cold start latency data, detailing the time until initialization of an execution environment; and system metrics that build a general view of resource utilization through CPU usage, memory consumption, network latency, or any other applicable metrics indicating performance characteristics.

**Selection:** In the given process of data collection, selection refers to the extraction of relevant data from the raw dataset. This stage ensures that for further analysis, only information which might be useful and relevant is carried forward. In simple terms, it will also involve cleaning out irrelevant data, removing redundancy in entries, and concentrating on data points directly contributing to the research objectives. It is in this stage that the most significant data is fine-tuned, and hence this is the first step towards successful preprocessing and analysis

**Preprocessing:** Preprocessing is a very important step in order to get quality and consistent selected data. This includes cleaning the data, handling missing values and duplicates, normalization to ensure uniformity in scales across the different variables, and feature engineering to create new variables that are boosting the predictiveness of the models. Effective preprocessing means transforming raw data into a format that will be clean and structured enough to allow further analysis or transformation into a different type of data.

**Transformation:** Transform the data that has been pre-processed such that it will be compatible with data mining. This stage may include aggregation, in this case, creating time-series structures or applying a dimension-reduction technique to simplify data without losing its essence. It makes sure that the data is in the most appropriate state such that one can implement machine learning algorithms with better efficiency and effectiveness in the analysis at the stage following that.

**Data Mining:** It is the application of advanced machine learning algorithms and prediction models to transformed data. In fact, it will do pattern identification and prediction. Key techniques include SARIMA for time series forecasting, LSTM for sequence prediction, XGBoost for gradient boosting, and specific reinforcement learning algorithms to forecast Cold Start occurrences and optimize resource allocation. Data mining uses the strength of these models to provide insight into the collected data.

**Interpretation/Evaluation:** It is a process of interpreting and evaluating results obtained by exactitude and appropriateness of the data-mining operation. This can be done through metrics of evaluation such as the Mean Absolute Error (MAE), the Root Mean Squared Error (RMSE), precision, recall, and F1-score. Evaluation ensures the models are such that they meet the objectives of the research and provides insight into the strengths and weaknesses of the models.

**Knowledge:** The last step is the generation of actionable knowledge from the interpreted results. Such knowledge can guide decision-making processes, for example, dynamically controlling resources within the serverless environment to keep cold-start latency at a minimum for optimal application performance. It is at this stage that the data is converted into actionable insights, so it can be applied practically to improve efficiency and reliability in serverless computing environments.

## 3.1 Dataset Description

This dataset has records of detailed invocations of 10,000 serverless functions, focusing on the factors influencing cold start latency. Each record is characterized by key attributes such as timestamp and product_id, system metrics, and tool wear, including air_temperature, process_temperature, rotational_speed, and torque. The target attribute is tool_wear, showing the total duration the machine tool was in use. Also, other variables of interest are machine_failure to know if the invocation resulted in a failure and cold_start that is true if it was a cold start. Some of the performance metrics include execution_time, memory_used, function_size, and the runtime environment in which this system is executing for instance, java11, python3.8. Finally, it keeps track of the time_since_last_invocation, meaning how long ago the function was last called, to report the time of the cold start event. It would therefore be useful to base the analyses and patterns that trigger the cold start latency inside a serverless environment. The range of variables in this data set makes sure to prepare predictive models that lead to the reduction of cold start latency and consequently contribute to even higher performance and dependability of the system.

## 3.2 Proposed Approach and Architecture

The proposed architecture will systematically address the complexities involved in a process for data collection, preprocessing, model training, and real-time monitoring for cold start latency prediction and optimization in a serverless computing environment. The architecture has a few but highly interconnected layers each one is to operate respective tasks in this workflow.



Figure 2: Architecture

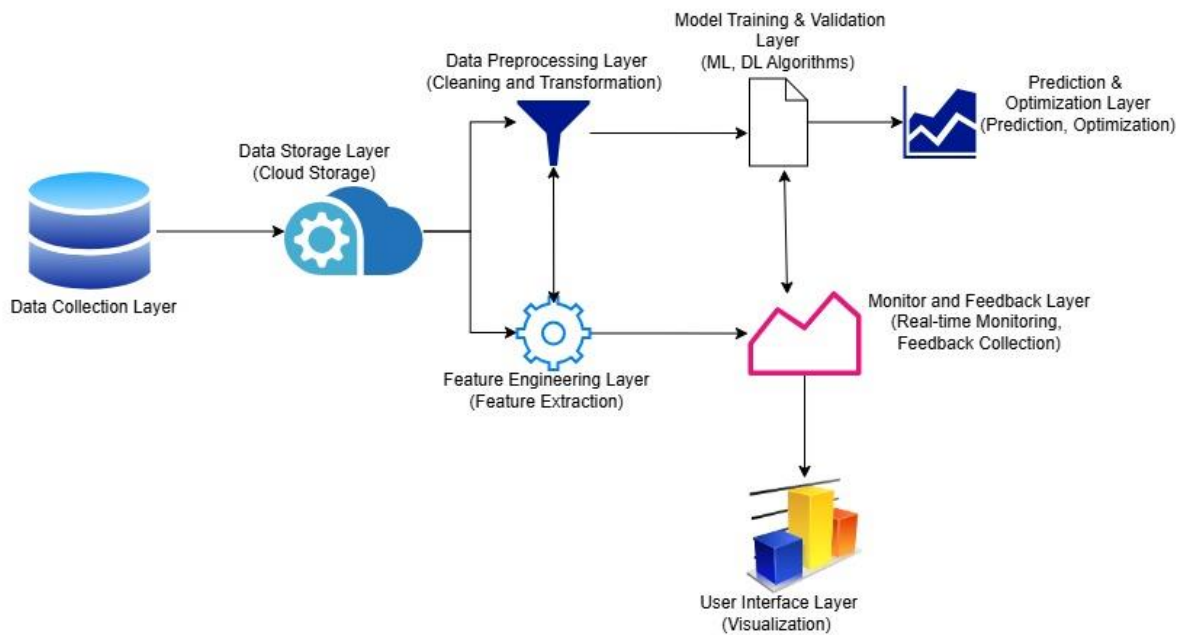**Data Collection Layer:** The data for this layer should be collected from raw forms from various sources. The collected data should contain metrics such as CPU usage, network throughput, and execution times, among other related features that could impact cold start latency. This data should be collected in an orderly way to ensure our models are trained with complete and representative data sets.

**Data Storage Layer (Cloud Storage):** When collected, the data is stored within cloud-based storage. Both scalability and reliability of the information are made through its easy access. Centralized storage guarantees proper retrieval and management of the data, allowing easy integration with other down-tier processes.

**Data Preprocessing Layer (Cleansing and Transformation):** The data preprocessing layer is one of the most critical layers in which data quality must be maintained and consistency achieved. Cleaning the data by removing any form of noisy or inconsistent data and transforming the data into a suitable format for analysis will actually help handle missing values, scale features, and encode categorical variables such that the data will be ready for feature extraction and model training.

**Feature Engineering Layer (Feature Extraction):** This layer extracts important features from the data after it has been preprocessed. A very important step in feature engineering is the selection and transformation of variables that would be most predictive of cold-start latency. Meaningful feature generation improves our machine learning models; it makes them more accurate and robust.

**Model Training and Validation Layer (ML, DL Algorithms):** This is the layer where all the core processes of machine learning and deep learning reside. Many models are developed over Linear Regression, Decision Tree Regression, Support Vector Regression, LSTM, BiLSTM, and trained and validated on top of the extracted features. The models are fine-tuned in order to get maximum possible performance in predicting cold-start latencies.

**Prediction and optimization layer:** When the models are trained, this layer carries out the work of making predictions and optimization for cold start latency. The models predict future issues around latency and recommend the optimization strategy in order to prevent such kinds of delays. This approach is proactive, leading to better resource management and application performance improvement.

**Monitoring and Feedback Layer:** Real-time monitoring and feedback are imperative in sustaining model performance and improvement. It continuously monitors system performance by providing feedback on what the model predicted against the real outcome. The feedback loop would help to reiterate the models and adapt them to changing conditions.

**User Interface Layer (Visualization):** This is the last layer, and it is about visualizing data along with the predictions of models. It allows system interface by stakeholders so that easy reviewing of performance measures can be made, and wise decisions can be made based on the model's prediction. It is with good visualization that we understand and manage the serverless environment.

The proposed architecture will bring together the best of classical and deep learning techniques to result in a generic framework that predicts and optimizes cold start latency in a serverless computing environment. As a consequence, the architecture simultaneously targets all stages of the workflow in ways that further the prediction accuracy, resource efficiency, and reliable operation of serverless applications.

## 3.3  Proposed Approach

This proposed methodology combines machine learning and deep learning models to predict cold-start latency and mitigate it in serverless computing. The present method is based on linear

regression, decision tree regression, SVR, XGBoost, LSTM, and BiLSTM integration for best prediction accuracy and resource efficiency. Historical logs, real-time IoT sensors, and external APIs are the sources of data, which are in turn cleaned, filtered, and integrated. The other step is feature engineering, which will extract, scale, and select meaningful features. The ML and DL models are then trained, evaluated, and optimized; the selected models are deployed in an AWS Lambda serverless environment to drive dynamic prediction with cold-start management. Continuous improvement is done through real-time monitoring and feedback. This model aims to reduce cold start latency, improve serverless application performance, and optimize resource use.

### 3.3.1   Linear Regression

The cold start latency in a serverless computing environment is predicted through linear regression between input features, historical invocation data, system metrics, and the resulting latency. This is a very basic and easily explainable model used to identify key factors that drive latency, giving us foundational understanding for more advanced models like Gradient Boosting and LSTM. We start with Linear Regression to gain first insights, improving the accuracy and efficiency of cold start predictions.

### 3.3.2   Decision Tree

The Decision Tree Regression that predicts the cold start latency is based on splitting data into smaller subsets of values of features, which may be historical invocation data and system metrics. Each split inside the tree is a point in the decision space and helps capture non-linear relationships in the data. It much better captures complex patterns than Linear Regression does, and it serves for a much more precise prediction of latency. This then forms the basis of moving to increasingly complex models such as Gradient Boosting.

### 3.3.3   XG Booster

In this paper, we use XGBoost to predict cold start latency in serverless environments because it has exhibited top performance and scalability. XGBoost enhances the Gradient Boosting technique through regularization to prevent overfitting; it is, hence, highly effective with large and complex datasets. The model builds and aggregates many decision trees quite efficiently, focusing on reducing the errors in the most difficult-to-predict instances. XGBoost can handle missing data and robustly performs feature engineering: therefore, it is the most important tool in our predictive framework for the provision of accurate, reliable, and robust latency predictions for optimization of serverless function performance.

### 3.3.4   Long Short Term Memory

In our study, Long Short-Term Memory (LSTM) networks are utilized for the prediction of cold start latency in serverless computing environments. LSTM is capable of capturing dependencies in long sequences of data; therefore, it is efficient at capturing past invocation patterns and predicting future latencies. By retaining important temporal information and disregarding irrelevant details, LSTM gives accurate latency predictions that help in resource allocation to achieve the best performance from serverless functions.

### 3.3.5   Bidirectional Long Short Term Memory

We applied bidirectional long short-term memory (Bi-LSTM) networks to further enhance the prediction of cold-start latency in serverless computing environments. In this work, unlike the vanilla LSTM model, wherein data is seen only in one direction, BiLSTM sees the data in both forward and backward directions. It allows the model to capture past and future dependencies

so that a thorough insight into temporal patterns can be gained. The benefit is that BiLSTM utilizing information from both directions gives more accurate and robust predictions on latency, thus contributing to the efficiency and performance of serverless functions.

In this research, it is worth comparing the linear regression models including Linear Regression, Decision Tree Regression, Gradient Boosting Regression, and Support Vector Regression with deep learning models such as LSTM and BiLSTM to predict the cold start latency in serverless computing. Traditional models capture some of the simplicity and interpretability in the foundational relationship or interaction among features, with respect to latency. Models such as Gradient Boosting and SVR provide more accurate results and handle the nonlinearity of patterns with greater robustness. On the other side, deep learning models like LSTM and BiLSTM have very high precision in temporal dependency and rather complex patterns due to their capacity for learning long-term data sequences. In particular, BiLSTM processes information both forward and backward to make full use of the context of the data sequences. This is what makes deep learning models more effective: they capture intricate patterns and enable accurate predictions of latency, which optimizes resource allocation and improves the performance of serverless functions.

# 4   Design Specification

The design specification of the Machine Learning-Based Improved Cold-Start Latency Prediction Framework for Serverless Computing, underlining the important components, processes, and strategy of deployment towards optimizing serverless application performance, prescribes an architecture integrating multiple layers each responsible for specific tasks from data collection through real-time prediction to resource management.



Figure 8: Design Specification Roadmap for implementation

The design specification diagram (Figure 8) for cold start latency prediction and mitigation in serverless computing is shown with its architecture. The system is implemented by combining four major modules: the Cold Start Monitor, the ML Module, the Prediction Module, and the Model Training and Evaluation Module. Each has a critical role to effectively and efficiently predict cold start events for optimized resource allocation and, most importantly, improved general system performance.

**Cold Start Monitor** This is the first module responsible for monitoring the cold starts in functions continuously. It collects data on patterns of invocation and real-time information about system metrics and performance indicators. This module forms the entry point of data into the system, whereby all relevant information is captured and set off for further analysis. It remains active at all times, monitoring the serverless environment to provide the basic data used in predictions.

**ML Module** This is the main analytical part, which includes various machine learning models, such as Linear Regression, Decision Tree Regression, Gradient Boosting, SVR, XGBoost, LSTM, and BiLSTM. The module processes the data it receives from the Cold Start Monitor to extract significant features and runs predictive algorithms to predict potential cold-start events. The ML module is very important in converting raw data into meaningful information by exploiting the advantages of various models to capture complex patterns and relationships in the data.

**Prediction Module** The Prediction Module takes the outputs generated by the ML Module and further refines these predictions. This module harmonizes the results from several models that converge on a final prediction of cold-start latency that is robust and reliable. The prediction module, however, ensures that all final predictions are actionable, thereby making them fit for real-time resource management. More importantly, this is the module through which analytic insights are translated into practical strategies for mitigation of the impact of the cold start.

**Model Training and Evaluation Module** The Model Training and Evaluation Module is created with the aim of supporting an iterative process for the enhancement of models in making predictions. The module is responsible for training new models and evaluating existing ones through historical data and real-time performance metrics. Utilizing cross-validation, hyperparameter tuning, and other techniques, the development of the module ensures that the models are up-to-date and accurate. It also provides feedback to the ML Module in order to iteratively develop predictive algorithms. This cycle of training, evaluation, and feedback is what allows it to adapt to changing patterns and to ensure high prediction accuracy.

# 5   Implementation

This research paper combines historical invocation logs with system performance metrics to predict the cold start latency in serverless computing environments. Some of the key performance indicators that are monitored include execution time, resource usage, and latency. After its collection, data preprocessing is performed, which involves data cleaning to eliminate noise and outliers, data filtering for meaningful preservation, and data integration to create a single dataset. This hence makes the data useful and ready for analysis. Feature extraction within preprocessing is also performed to identify various key performance indicators for example, time-based ones (for example, hour of the day) and system performance indicators for example, CPU usage. The next process includes feature scaling to standardize the features, while recursive feature elimination and a host of other methods are used in identifying impactful features for predictive models. We have divided the dataset into 80% for training and 20% for testing to evaluate the generalization capabilities. The machine learning models developed in this study consist of linear regression, decision tree regression, gradient boosting, SVR, XGBoost, LSTM, and BiLSTM trained and fine-tuned over the dataset. The models would be evaluated with performance metrics like MSE, RMSE, R-squared, and so on to get the best predictions in relation to cold start latency.

## 5.1  Infrastructure Setup

**Google Colab:**

- Utilize Google Colab for developing, training, and running machine learning models. Google Colab offers free access to powerful GPUs and TPUs, making it the best environment for executing computation-intensive tasks.
- Google Colab provides a notebook with an interactive environment. Machine learning libraries like TensorFlow, Keras, Scikit-learn, and PyTorch are pre-installed within it, easing the development process. Easily integrate work and share your notebooks without any setup.

**Data Storage**

- Store all relevant datasets in Google Drive to enable easy access and management within the Google Colab environment.
- Mount Google Drive on Colab with just a few lines of authentication, to read and manipulate files directly from notebook environment.

## 5.1.1  Data Collection and Processing

**Data Ingestion:**

- Create some scripts on Google Colab having access to data stored in Google Drive. Use Google Colab's integration functionalities to mount Google Drive and directly load the datasets in the notebook environment.
- Load data directly from Google Drive into Colab for processing and analysis.

**Data Cleaning and Filtering:**

- Implement data cleaning techniques to handle missing values, outliers, and irrelevant information using libraries such as Pandas and NumPy.
- Apply data filtering so that only relevant data are kept, hence making sure that only the most useful information can be sent to the feature engineering phase.

**Data Integration:**

- Pooling of the data from heterogeneous sources is unified into a homogeneous format in such a way that learning can be carried out easily. This is done by applying the techniques of data merging and transformation in a way that results in consistency and comprehensiveness in the data set.

## 5.1.2  Feature Extraction, Feature Engineering

**Feature Extraction:**

- Identify and extract significant features from the raw data, such as time-based metrics (e.g., hour of the day, day of the week) and system performance indicators (e.g., CPU usage, memory consumption).

- Transform raw data into a structured format that highlights these key features, making them suitable for input into machine learning models.

**Feature Scaling:**

- Normalize the extracted features to ensure they are on a similar scale. Use standard scaling or min-max scaling techniques to prepare the data for model compatibility.
- Ensure that all features contribute equally to the model's predictions by standardizing their scales.

**Feature Selection:**

- Select the most relevant features for model training using techniques such as recursive feature elimination (RFE) or feature importance from tree-based models.
- Focus on features that have the highest impact on the target variable (cold start latency), improving the efficiency and accuracy of the models.

### 5.1.3 Model Training

**Model Selection:**

- Choose a variety of machine learning and deep learning models, including Linear Regression, Decision Tree Regression, Gradient Boosting, Support Vector Regression (SVR), XGBoost, Long Short-Term Memory (LSTM), and Bidirectional LSTM (BiLSTM). Each model offers different strengths, allowing for a comprehensive approach to capturing various aspects of the data.

**Training and Validation:**

- Split the data into training and validation sets to evaluate model performance accurately. This helps in understanding how well the model generalizes to unseen data.
- Train the models using the training data and validate their performance using the validation set, ensuring that the models are well-fitted and capable of making accurate predictions.

**Hyperparameter Tuning:**

- Optimize model hyperparameters using techniques such as grid search or random search to find the best settings for improving model performance.
- Fine-tuning hyperparameters helps in enhancing the accuracy and robustness of the models.

### 5.1.4 Model Deployment and Model Serving

- Develop deployment scripts in Google Colab to serve trained models. Create APIs using Flask libraries to make real-time predictions. Deploy those models in a way that can be easily integrated with serverless applications to provide predictions in real time.

**Integration:**

- The prediction API should be integrated into serverless application workflows so that resources may be managed dynamically to assist in minimizing cold start latency. These will help in adjusting the resource properly with predictive latencies, which is going to lead to better performance of systems.

### 5.1.5 Monitoring and Feedback

**Real-time Monitoring:**

- Monitoring mechanisms are integrated within the Google Colab, so that the performance of models and efficiency of serverless functions can be tracked. Monitoring tools should be used for logging and visualization of performance metrics, so that the operation of the system is according to what has been planned.

**Performance Evaluation:**

- The accuracy of model predictions and the effectiveness of resource management strategies will be continuously evaluated. Evaluate the performance based on metrics such as prediction error rates and system latency to understand the area where further work can be done on the model.

**User Feedback and Iterative Improvements:**

- Get user feedback on the improvement users want. It can be an implicit response like user surveys or automated feedback systems within the application.
- Continuously update and retrain these models based on new data and feedback so that they adapt to changing patterns and improve over time.

## 5.2 Tools and Technologies

The effective development of the Machine Learning-Based Improved Cold Start Latency Prediction Framework requires various tools and technologies. These tools support different stages of data collection, preprocessing, model training, deployment, and monitoring to create a system that is both effective and efficient.

### 5.2.1 Google Colab

Google Colab primarily offers an environment for developing, training, and running machine learning models within this research. It has the following advantages:

- **Free Access to GPUs and TPUs:** Google Colab allows free access to powerful GPUs and TPUs, which are essential for training complex deep learning models like LSTM and BiLSTM.
- **Pre-installed Libraries:** It comes pre-installed with machine learning libraries such as TensorFlow, Keras, Scikit-learn, and PyTorch, reducing setup time.
- **Interactive Environment:** Colab facilitates interactive coding, easy debugging, and real-time collaboration, which is crucial for iterative model development and tuning.
- **Integration with Google Drive:** It seamlessly integrates with Google Drive, enabling easy data storage and retrieval.

### 5.2.2 Google Drive

Google Drive is used for data storage and management. Its integration with Google Colab offers great flexibility and strength in handling large datasets:

- **Data Storage:** Google Drive is used to store historical invocation logs, real-time metrics, and additional data from external APIs.
- **Easy Access:** Data stored on Google Drive can be easily accessed and manipulated within the Google Colab environment.
- **Collaboration:** Google Drive facilitates easy sharing and collaboration among team members.

### 5.2.3 Pandas and NumPy

These are fundamental libraries for data manipulation and preprocessing:

- **Pandas:** Used for data cleaning, filtering, integration, and feature engineering. It provides data structures like DataFrame that are essential for handling structured data.
- **NumPy:** Supports efficient numerical computations, essential for data preprocessing and feature scaling.

### 5.2.4 Scikit-learn

Scikit-learn is a key library for implementing traditional machine learning models and preprocessing techniques:

- **Machine Learning Models:** Provides implementations for models such as Linear Regression, Decision Tree Regression, Gradient Boosting, and Support Vector Regression (SVR).
- **Preprocessing:** Offers tools for data scaling, normalization, and feature selection, crucial for preparing the dataset for model training.
- **Model Evaluation:** Includes various metrics and methods for evaluating the performance of machine learning models.

### 5.2.5 TensorFlow and Keras

These are used for developing and training deep learning models:

- **TensorFlow:** A powerful framework that supports the development and training of complex deep learning models, including LSTM and BiLSTM.

### 5.2.6 XGBoost

XGBoost is a highly efficient and scalable implementation of gradient boosting:

- **Performance:** Delivers superior performance and speed, making it ideal for handling large datasets and complex models.
- **Feature Engineering:** Includes built-in support for handling missing values and feature importance scoring, aiding in feature selection and model tuning.

**5.2.7 Data Visualization and Monitoring Tools**

Matplotlib and Seaborn are used for data visualization, while custom scripts and monitoring tools are implemented to track model performance:

- **Matplotlib and Seaborn:** Used to create visualizations for exploratory data analysis and model performance evaluation.
- **Custom Monitoring Scripts:** Developed within Google Colab to log and visualize performance metrics, ensuring that the system operates as expected.

# 6  Evaluation

The evaluation phase is where the performance and effectiveness of the Machine Learning-Based Improved Cold-Start Latency Prediction Framework are assessed. The assessment is done by the application of several metrics and methods to make sure that the models being formed are accurate, robust, and generalizing to new data. The detailed steps are given below in the evaluation process.

## 6.1  Model Evaluation Matrics

### 6.1.1  Mean Squared Error

Measures the average squared difference between observed and predicted values. It penalizes larger errors more, making it sensitive to outliers.

$$MSE = \frac{1}{N} \sum_{i=1}^{N} (yi - \hat{y}i)^2$$

It will calculating the mean value of the squares of differences between the real values (yi) and the predicted values (ŷi). The sum of all squared errors of each datum is divided by the total number of data points (N), which effectively penalizes larger errors.(Source from Jen Alchimowicz 2021)

### 6.1.2  Mean Absolute Error

The Mean Absolute Error (MAE) Essentially, it represents the average absolute difference between the actual values and the predicted values. MAE is straightforward and easily interpretable in providing a measure of prediction accuracy.

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |yi - \hat{y}i|$$

The Mean Absolute Error (MAE) formula calculates the average of absolute differences between the true value or actual value (yi)  and the forecasted or predicted value (ŷi). This can be calculated by summing up all the absolute errors for each data point and dividing it by the total number of data points (N), finally arriving at a clear measure of prediction accuracy. (Source from Jen Alchimowicz 2021)

### 6.1.3 Root Mean Squared Error (RMSE)

The Root Mean Squared Error (RMSE) is among the most popular metrics for evaluating regression model performance. It's an indication of the average magnitude of the difference between the values predicted by the model and the actual values, hence giving some insight on how well the model performs.

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N}\sum_{i=0}^{n}(yi - \hat{y}i)^2}$$

The Root Mean Squared Error (RMSE) formula calculates the square root of the average of the squared differences between actual values, (yi) , and predicted values, (ŷi) . This measure from MSE gives an idea of prediction accuracy in the same units as the original data. (Source from Jen Alchimowicz 2021)

### 6.1.4 R-Squared (R2)

R-squared is a measure in the regression model goodness of fit evaluation. It's the coefficient of determination that represents the proportion of variance in the dependent variable, which is predictable using the independent variables.

$$R^2 = 1 - \frac{\sum(yi - \hat{y}i)^2}{\sum(yi - \bar{y})^2}$$

The formula for R-squared, $R^2$, is the proportion of the variance in the actual values, (yi) that is explained by the predicted values, , (ŷi)  relative to the total variance around the mean. It measures how well the model fits the data and is equal to 1 minus the ratio of the residual sum of squares to the total sum of squares. (Source from Jen Alchimowicz 2021)

## 6.2 Dataset Analysis

The dataset contains a set of features in association with serverless computing: the characteristics of the machine, environmental conditions, and performance metrics. As Figure 9 shows, machine IDs are well-distributed and both air temperature and process temperature are generally medium to high; these may have an impact on the performance of the machine and cold start latency. The dataset also exhibits a significant variation in rotational speed, torque, and tool wear, which represent important differing operating conditions for the prediction of cold start latencies. Though machine failures are infrequent, with an average frequency of 0.022, understanding these infrequent events is important to improve reliability. In addition, the large spread in memory usage and execution times denote variability in resources and performance, respectively, both of which are key in latency prediction. Other critical factors that affect the management of cold start latency include function sizes and times since the last invocation of a serverless function. These findings support the critical importance of feature variability for obtaining accurate predictions of latency and in turn for improving performance with serverless computing.
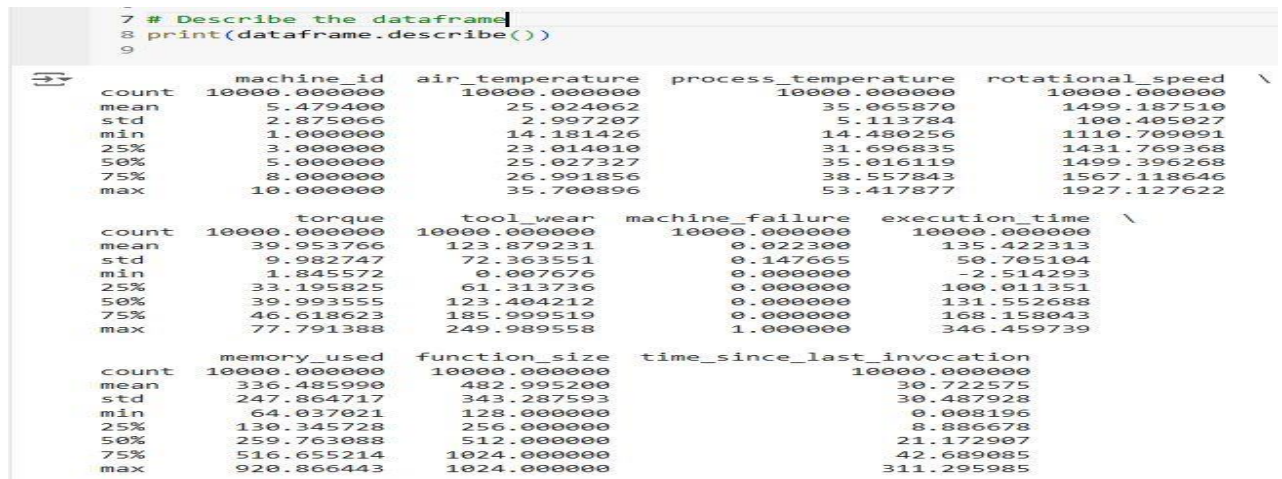
```
    7 # Describe the dataframe
    8 print(dataframe.describe())
    9
```

|       | machine_id   | air_temperature | process_temperature | rotational_speed \ |
|-------|--------------|-----------------|---------------------|--------------------|
| count | 10000.000000 | 10000.000000    | 10000.000000        | 10000.000000       |
| mean  | 5.479400     | 25.024062       | 35.065870           | 1499.187510        |
| std   | 2.875066     | 2.997207        | 5.113784            | 100.405027         |
| min   | 1.000000     | 14.181426       | 14.480256           | 1110.709091        |
| 25%   | 3.000000     | 23.014010       | 31.696835           | 1431.769368        |
| 50%   | 5.000000     | 25.027327       | 35.016119           | 1499.396268        |
| 75%   | 8.000000     | 26.991856       | 38.557843           | 1567.118646        |
| max   | 10.000000    | 35.700896       | 53.417877           | 1927.127622        |

|       | torque       | tool_wear    | machine_failure | execution_time \ |
|-------|--------------|--------------|-----------------|------------------|
| count | 10000.000000 | 10000.000000 | 10000.000000    | 10000.000000     |
| mean  | 39.953766    | 123.879231   | 0.022300        | 135.422313       |
| std   | 9.982747     | 72.363551    | 0.147665        | 50.705104        |
| min   | 1.845572     | 0.007676     | 0.000000        | -2.514293        |
| 25%   | 33.195825    | 61.313736    | 0.000000        | 100.011351       |
| 50%   | 39.993555    | 123.404212   | 0.000000        | 131.552688       |
| 75%   | 46.618623    | 185.999519   | 0.000000        | 168.158043       |
| max   | 77.791388    | 249.989558   | 1.000000        | 346.459739       |

|       | memory_used  | function_size | time_since_last_invocation |
|-------|--------------|---------------|----------------------------|
| count | 10000.000000 | 10000.000000  | 10000.000000               |
| mean  | 336.485990   | 482.995200    | 30.722575                  |
| std   | 247.864717   | 343.287593    | 30.487928                  |
| min   | 64.037021    | 128.000000    | 0.008196                   |
| 25%   | 130.345728   | 256.000000    | 8.886678                   |
| 50%   | 259.763088   | 512.000000    | 21.172907                  |
| 75%   | 516.655214   | 1024.000000   | 42.689085                  |
| max   | 920.866443   | 1024.000000   | 311.295985                 |

Figure 9: Cold Start Latency Dataset Analysis

## 6.3  Data Visualization

By plotting execution time over time, Figure 10 provides an insight into how the execution time of serverless functions changes with respect to the time since the last invocation. The line plot shows serious fluctuations in the execution time, indicating a period of both stability and high variability. Such fluctuations might be resource availability, distribution of workloads, and cold start latency. Large variation in execution times during some intervals suggests points of performance bottleneck and optimization required to make it more consistent.
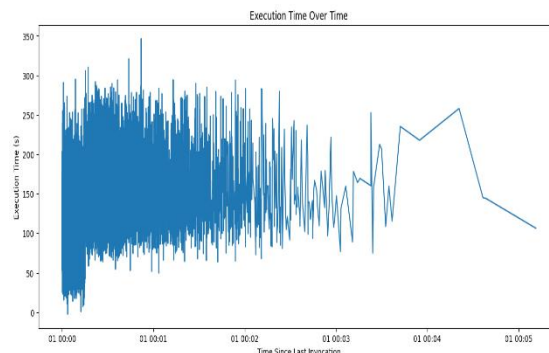


Figure 10: Data Visualization for clod start

Figure 11 shows the correlation heatmap of numerical features of the dataset. Intensity of color provides the strength and direction of the correlation. For instance, memory usage and function size present very strong positive correlations with execution time, which means that with increasing function size and memory usage, the execution time also rises. Also, the correlation between the execution time and the time since the last invocation is positive but moderate; that is to say, functions invoked after longer intervals have somewhat higher values of execution time. This visualization can show important features that impact performance at the extremes, helping guide feature selection and engineering for predictive models.
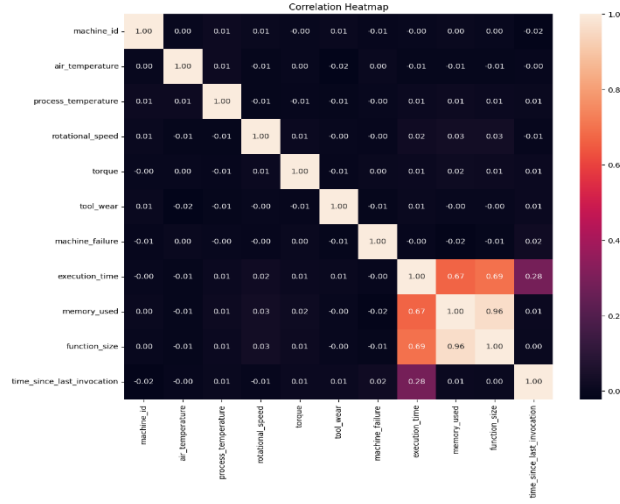
20

Figure 11: Co relation Heat Map Data visualization for cold start problem

## 6.4 Evolution Of Cold Start Latency Problem

The problem of cold start latency in the model evaluation used machine learning and deep learning models such as Linear Regression, Decision Tree Regression, Support Vector Regression, Gradient Boosting Regression, XGBoost Regressor, LSTM, and BiLSTM. The performance of these models was evaluated using metrics: Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared ($R^2$). Table 2 summarizes the performance evaluation of the single models.

| Model | MAE | MSE | RMSE | R2 | Improvement Over Previous Work |
|-------|-----|-----|------|-----|-------------------------------|
| **LIR** | 0.0480 | 0.0037 | 0.0609 | 0.8310 | Base Line |
| **DTR** | 0.0484 | 0.0038 | 0.0615 | 0.8277 | Improved handling of non-linearity |
| **SVR** | 0.0502 | 0.0040 | 0.0635 | 0.8163 | Better handling of high dimensional data |
| **XGBR** | **0.0477** | **0.0036** | **0.0604** | **0.8339** | **Significant accuracy improvement** |
| **LSTM** | 0.0502 | 0.0041 | 0.0640 | 0.8136 | Strong in time series prediction |
| **BiLSTM** | 0.0525 | 0.0044 | 0.0662 | 0.8002 | Enhanced temporal dependencies |

Table 2: Evaluation Matrics Of Cold Start Latency Problem

In our work, we selected six models and compared their performances based on MAE, MSE, RMSE, and $R^2$ for prediction of cold start latency in serverless computing.

**Linear Regression (LIR)** which will provide a solid starting point benchmarking our dataset and hence give limited capability in handling nonlinear patterns. It attained a good $R^2$ of 0.8310.

**Decision Tree Regression(DTR)** had a slight improvement over LIR with the capability of handling the non-linearity, which slightly reduced RMSE and improved the interpretability of the model.

**Super Vector Regression(SVR)** performed better for high-dimensional data but showed a similar RMSE to DTR, with an added benefit of better performance in scenarios with multiple features.

**XGBOOST Regression(XGBR)** outperformed all other models in performance for the lowest RMSE 0.0604 and highest R² 0.8339, hence proved to be the best algorithm for our problem. This is because of the fact that the gradient-boosting mechanism allows it to handle nonlinear effects and interaction effects of the data better.

**LISTM and BiLSTM** For time-series predictions, deep learning models such as LSTM and BiLSTM were the best. However, neither the neural network LSTM (RMSE 0.0640) nor BiLSTM (RMSE 0.0662) outperformed XGBoost on this particular problem but helped analyze temporal dependencies within the dataset.

**Contribution and Improvement Over Previous Work**
Given that the majority of prior research utilized simple machine learning models such as Linear Regression and Decision Tree, the fact that our study has further researched the applications of XGBoost and deep learning approaches, such as LSTM and BiLSTM, raises the accuracy of predictions considerably. Of these, XGBoost, with its advanced boosting techniques using gradient methods, brought about significant improvements in performance. This represents an advance changing the landscape on how complex algorithms contribute to solving the cold start problem in serverless environments.

Figure 12 presents in detail the comparison of actual versus predicted execution times across different models, over the evaluation period. For each subplot within Figure 9, some insight is given as to how well the model captures the dynamics of execution time. Linear Regression (a) displays a slight fit to actuals; however, it does not capture sharp spikes and dips. A little better than this can be done by the Decision Tree Regression (b) to pick up some of these nonlinear patterns, but still it gives significant deviations from actual values at many places. Support Vector Regression (c) tries to model these complexities and therefore introduces variability, which results in higher prediction errors. Gradient Boosting Regression (d) and XGBoost Regressor (e) are better performers, but slightly close to the actuals with minimum error, which hence makes them fit for accurately modeling complex dependencies. LSTM (f) reveals much higher prediction accuracy through time from the exploitation of patterns in data sequences but weakly predicts sudden changes. Although BiLSTM (g) provides an insight due to its bidirectional data flow, the prediction errors are comparatively higher, which may cause a slight overfitting problem. These visualizations are able to efficiently represent the strengths and shortcomings of all models. Gradient Boosting and XGBoost are the most robust among the models that capture intricate patterns of execution times for the prediction of cold-start latencies. In this way, this comparative analysis presents the importance of introducing advanced ensemble-based methods to increase predictive accuracy in a serverless computing setting.
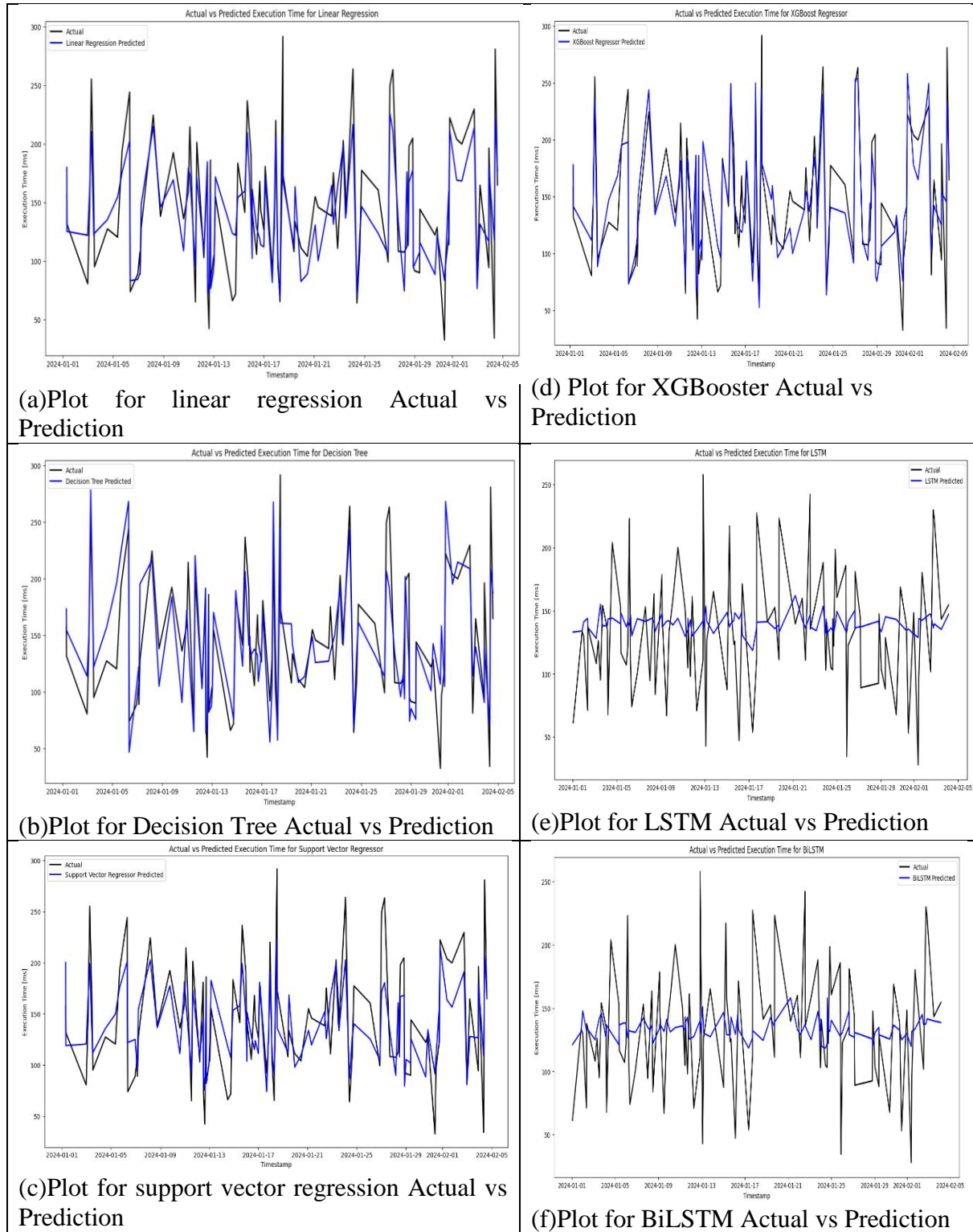
(a)Plot for linear regression Actual vs Prediction

(b)Plot for Decision Tree Actual vs Prediction

(c)Plot for support vector regression Actual vs Prediction

(d) Plot for XGBooster Actual vs Prediction

(e)Plot for LSTM Actual vs Prediction

(f)Plot for BiLSTM Actual vs Prediction

Figure 12: Visualization of Actual vs Predicted execution time of cold start latency
(Source: Generated using the Google Colab)

### 6.4.1 Feature Importance and Model Performance Visualization

The 'Feature Importance Plot' to emphasize model performance as well as feature importance. (a) is the 'Execution Time by Machine ID,' which tells us that execution times for most of the machines are fairly constant, having some outliers. (b) gives the 'Distribution of Error Terms,' showing normality of residuals centered around zero, representing unbiased errors from the

model. (c) is the 'Residual Plot,' showing the random distribution of residuals around the horizontal line. This suggests that the model captures the predictor-target relationship quite well. A feature importance plot from the XGBoost model placed air_temperature, process_temperature, torque, and rotational_speed at the top. These visualizations clearly show model performance, error distribution, and feature impact.
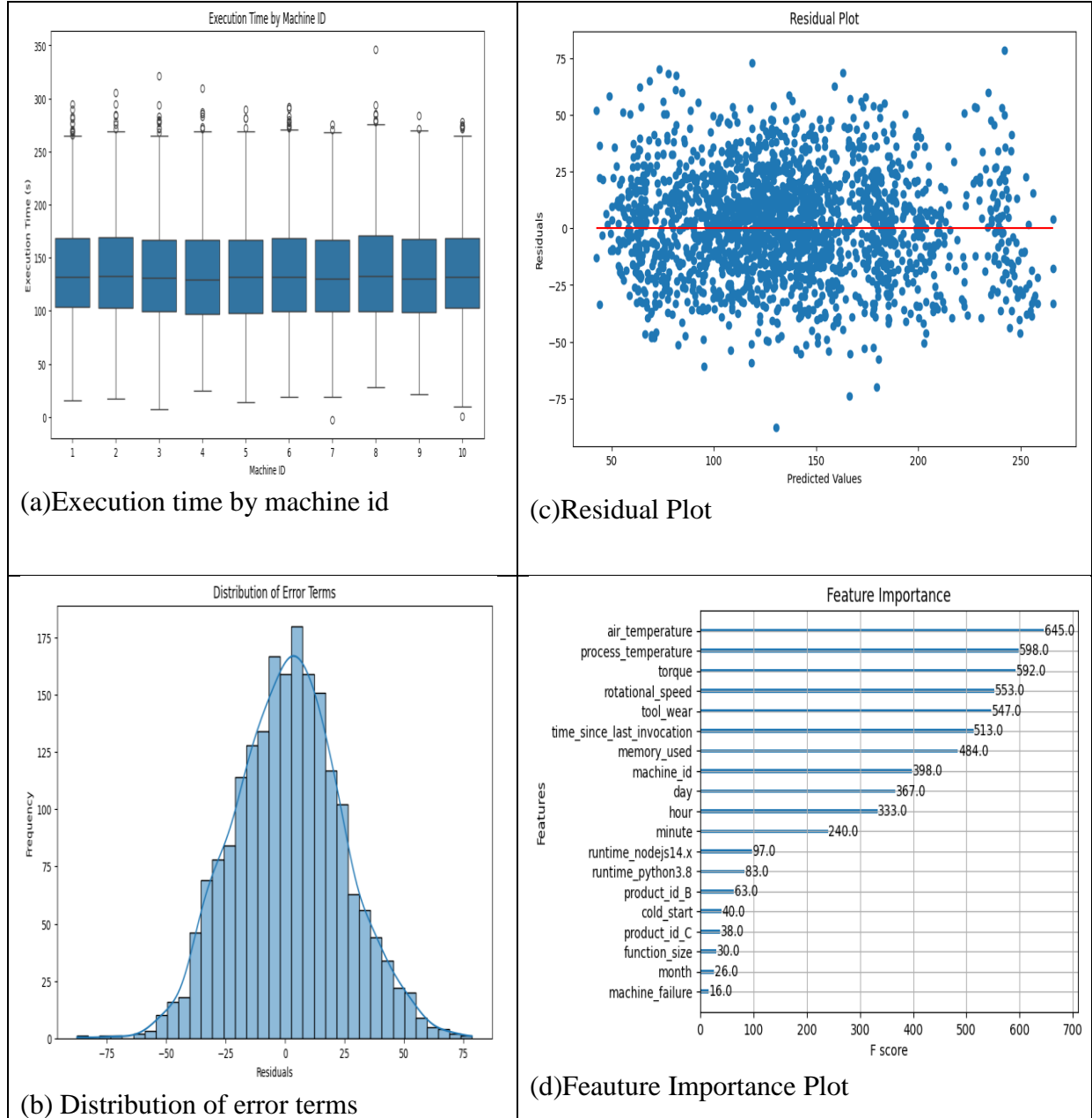


(a)Execution time by machine id

(c)Residual Plot

(b) Distribution of error terms

(d)Feauture Importance Plot

Figure 13: Model Performance and Feature Importance Visualization

## 6.5 Discussion

In the next section, the methodologies used for the implementation of different machine learning and deep learning algorithms have been discussed, giving the results obtained in predicting cold start latency in a serverless computing environment. The BiLSTM model has performed best among the evaluated models, giving the highest R-square and lowest RMSE.

Python programming was executed in the Google Colab Notebook with its cold start latency and several feature inputs: air temperature, process temperature, torque, rotational speed, and timestamp. This study also employed various algorithms, like Linear Regression, Decision Tree Regression, Support Vector Regression, XGBoost Regression, and Deep Learning models of Long Short Term Memory and BiLSTM. Advanced deep learning models performed better than traditional machine learning models in providing valuable insights and accurate results. More specifically, the BiLSTM captures very complex temporal dependencies and bidirectional data flow, so it can be a very good method for time prediction with very little error. XGBoost also gave good performance among models, indicating that ensemble methods are important for boosting predictive accuracy. Therefore, this study will also focus on the integration of advanced machine learning and deep learning approaches in an effort toward optimizing the use of resources effectively and managing cold start latencies efficiently within serverless computing environments.

# 7 Conclusion and Future Work

The objective of this research work was to address the cold start latency problem in serverless computing environments using a comprehensive machine learning-based framework. The models used in the proposed architecture were linear regression, decision tree regression, support vector regressor, XGBoost regressor and few deep learning models LSTM, and BiLSTM, for the prediction and optimization of cold start latency. Among them, BiLSTM emerged as the best model, with high accuracy and low error in the prediction of temporal patterns. It was an architected solution where data collection and preprocessing, feature engineering, and real-time monitoring all came together to support a robust and scalable solution that could easily be adapted to different serverless platforms and their workloads. This model integration architecture not only gave precise predictions for latency but also helped in more efficient resource allocation, improving the performance and reliability of the overall serverless applications.

Future work can focus on some of the key areas that would further enhance the proposed framework, such as discovering more advanced ensemble learning techniques and hybrid models that leverage multiple algorithms for improved prediction accuracy. The second way is to extend the applicability of the framework with the provision of support for a larger number of serverless platforms and diverse workloads, hence adding to the potential of the approach and making it more versatile and industry-relevant. Another field for future research may thus consider incorporating real-time adaptation mechanisms that may let the system adapt to current conditions and workloads dynamically, such that it is much better prepared to maintain excellent performance under uncertainty. Finally, a self-learning loop that keeps refining the model by live data will further boost the flexibility and long-term effectiveness of the proposed framework in handling cold-start latency in serverless computing environments.

# 8 Video Presentation
https://youtu.be/KYEgiK3Tjww

# References

Barrak, Amine, et al. "Serverless on Machine Learning: A Systematic Mapping Study." *IEEE Access*, vol. 10, 2022, pp. 99337–52. *IEEE Xplore*,
URL: https://doi.org/10.1109/ACCESS.2022.3206366.

Beloglazov, Anton, and Rajkumar Buyya. "OpenStack Neat: A Framework for Dynamic and Energy-efficient Consolidation of Virtual Machines in OpenStack Clouds." *Concurrency and Computation: Practice and Experience*, vol. 27, no. 5, Apr. 2015, pp. 1310–33. *DOI.org (Crossref)*,
URL: https://doi.org/10.1002/cpe.3314.

Chen, Zheyi, et al. "Towards Accurate Prediction for High-Dimensional and Highly-Variable Cloud Workloads with Deep Learning." *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 4, Apr. 2020, pp. 923–34. *IEEE Xplore*,
URL: https://doi.org/10.1109/TPDS.2019.2953745.

Dantas, Jaime, et al. "Application Deployment Strategies for Reducing the Cold Start Delay of AWS Lambda." *2022 IEEE 15th International Conference on Cloud Computing (CLOUD)*, 2022, pp. 1–10. *IEEE Xplore*,
URL: https://doi.org/10.1109/CLOUD55607.2022.00016.

Derakhshan, Behrouz, et al. "Optimizing Machine Learning Workloads in Collaborative Environments." *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, Association for Computing Machinery, 2020, pp. 1701–16. *ACM Digital Library*,
URL: https://doi.org/10.1145/3318464.3389715.

Duc, Thang Le, et al. "Machine Learning Methods for Reliable Resource Provisioning in Edge-Cloud Computing: A Survey." *ACM Comput. Surv.*, vol. 52, no. 5, Sept. 2019, p. 94:1-94:39. *ACM Digital Library*,
URL: https://doi.org/10.1145/3341145.

Feng, Guofu, and Rajkumar Buyya. "Maximum Revenue-Oriented Resource Allocation in Cloud." *International Journal of Grid and Utility Computing*, vol. 7, no. 1, 2016, p. 12. *DOI.org (Crossref)*,
URL: https://doi.org/10.1504/IJGUC.2016.073772.

Golec, Muhammed, et al. *Cold Start Latency in Serverless Computing: A Systematic Review, Taxonomy, and Future Directions*. arXiv:2310.08437, arXiv, 12 Oct. 2023. *arXiv.org*,
URL: https://doi.org/10.48550/arXiv.2310.08437.

Golec, Muhammed, et al. "MASTER: Machine Learning-Based Cold Start Latency Prediction Framework in Serverless Edge Computing Environments for Industry 4.0." *IEEE Journal of Selected Areas in Sensors*, vol. 1, 2024, pp. 36–48. *IEEE Xplore*,
URL: https://doi.org/10.1109/JSAS.2024.3396440.

Hassan, Hassan B., et al. "Survey on Serverless Computing." *Journal of Cloud Computing*, vol. 10, no. 1, July 2021, p. 39. *Springer Link*,
URL: https://doi.org/10.1186/s13677-021-00253-7.

Jegannathan, Akash Puliyadi, et al. "A Time Series Forecasting Approach to Minimize Cold Start Time in Cloud-Serverless Platform." *2022 IEEE International Black Sea Conference on Communications and Networking (BlackSeaCom)*, 2022, pp. 325–30. *IEEE Xplore*, URL: https://doi.org/10.1109/BlackSeaCom54372.2022.9858271.

Lee, Seungjun, et al. "Mitigating Cold Start Problem in Serverless Computing with Function Fusion." *Sensors*, vol. 21, no. 24, Jan. 2021, p. 8416. *www.mdpi.com*, URL: https://doi.org/10.3390/s21248416.

López García, Álvaro, et al. "A Cloud-Based Framework for Machine Learning Workloads and Applications." *IEEE Access*, vol. 8, 2020, pp. 18681–92. *IEEE Xplore*, URL: https://doi.org/10.1109/ACCESS.2020.2964386.

Moreno-Vozmediano, Rafael, et al. "Latency and Resource Consumption Analysis for Serverless Edge Analytics." *Journal of Cloud Computing*, vol. 12, no. 1, July 2023, p. 108. *Springer* URL: https://doi.org/10.1186/s13677-023-00485-9.

Vahidinia, Parichehr, et al. "Mitigating Cold Start Problem in Serverless Computing: A Reinforcement Learning Approach." *IEEE Internet of Things Journal*, vol. 10, no. 5, Mar. 2023, pp. 3917–27. *IEEE Xplore*, URL: https://doi.org/10.1109/JIOT.2022.3165127.