Department of Electrical and Computer Engineering
North South University

CSE498R: Directed Research

**A Machine Learning-Based Approach to Cold Start Prediction and Smart Prewarming in AWS Lambda**

Tanzim Rafat  2013485642

Hasan Muhtasim Ishmam  2031416042

Hadi Uzzaman  2022206642

Nandini Das  2031326642

Md Sadman Sakib Shachyo  2031388642

Faculty Advisor:

Dr. Salekul Islam

Professor

ECE Department

Summer, 2024

# APPROVAL

**Tanzim Rafat (ID: 2013485642), Hasan Muhtasim Ishmam (ID: 2031416042) Hadi Uzzaman (ID: 2022206642), Nandini Das (ID: 2031326642), and Md Sadman Sakib Shachyo (ID: 2031388642)** from the Electrical and Computer Engineering Department of North South University, have worked on the Directed Research Project titled **"A Machine Learning-Based Approach to Cold Start Prediction and Smart Prewarming in AWS Lambda"** under the supervision of **Dr. Salekul Islam** in partial fulfillment of the requirement for the degree of Bachelors of Science in Engineering and has been accepted as satisfactory.

**Supervisor's Signature**

………………………………….

**Dr. Salekul Islam**
**Professor**
Department of Electrical and Computer Engineering
North South University
Dhaka, Bangladesh.

**Chairman's Signature**

………………………………….

**Prof. Dr. Mohammad Abdul Matin**

Department of Electrical and Computer Engineering
North South University
Dhaka, Bangladesh.

# DECLARATION

This is to clarify that this project is our original work. No part of this work has been submitted elsewhere partially or fully for the award of any other degree or diploma. All project related information will remain confidential and shall not be disclosed without the formal consent of the project supervisor. Relevant previous works presented in this report have been properly acknowledged and cited. The plagiarism policy, as stated by the supervisor, has been maintained.

Students' names & Signatures

**1. Tanzim Rafat**

_ _ _ _ _ _ _ _ _ _ _ _ _ _ _

**2. Hasan Muhtasim Ishmam**

_ _ _ _ _ _ _ _ _ _ _ _ _ _ _

**3. Nandini Das**

_ _ _ _ _ _ _ _ _ _ _ _ _ _ _

**4. Hadi Uzzaman**

_ _ _ _ _ _ _ _ _ _ _ _ _ _ _

**5. Md Sadman Sakib Shachyo**

_ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _

# ACKNOWLEDGEMENTS

# ABSTRACT

**A Machine Learning-Based Approach to Cold Start Prediction and Smart Prewarming in AWS Lambda**

Cold starts in AWS Lambda introduce latency that affects the performance of serverless applications, especially those requiring real-time responsiveness. This project presents a machine learning-based framework to predict and mitigate cold starts using system logs and invocation patterns. The approach combines clustering for time-based pattern discovery and supervised models to forecast both the likelihood and delay of cold starts. A smart prewarming strategy is then applied to minimize cold start penalties while maintaining cost-efficiency. The proposed system achieved high prediction accuracy and demonstrated effective cost reduction, offering a scalable and intelligent solution for optimizing serverless performance.

# TABLE OF CONTENTS

# LIST OF FIGURES

# List Of Tables:

# Chapter 1 Introduction

## 1.1　Background and Motivation

With the increased use of serverless computing platforms like AWS Lambda, Google Cloud Functions, and Microsoft Azure Functions, it is currently feasible for organizations to build and deploy scalable applications without needing to deal with complex infrastructure. Serverless computing is cost-effective through its billing of users based on execution time, making it an attractive option for the majority of cloud-native applications. Serverless architectures, nonetheless, are afflicted with a fundamental performance bottleneck known as the cold start problem.

Cold starts occur when a serverless function, following a period of inactivity, is invoked and requires initialization of resources such as a virtual machine or container. This introduces significant latency between function invocation and execution, which is expensive for time-critical applications such as real-time analytics, user-facing APIs, and IoT applications. Studies have shown that cold start latencies in serverless functions range from hundreds of milliseconds to seconds, depending on function memory configurations, execution environment, and dependency complexity. For high-performance applications, even small delays can degrade user experience and overall system responsiveness [5][2].

Although previous solutions, such as Provisioned Concurrency (which keeps functions pre-installed), are designed to eliminate cold start latency, such solutions are expensive because resources are required to be constantly allocated regardless of function usage. Moreover, the static approaches cannot adapt to dynamic workloads and cannot predict when a function will experience a cold start. Hence, there is a need for an affordable and dynamic solution that not only predicts cold starts but also dynamically tailors prewarming strategies to match prevailing workload trends [6][4].

This project aims to design a **machine learning-based framework** for predicting cold starts and intelligently managing prewarming in AWS Lambda, optimizing performance without incurring unnecessary costs.

## 1.2   Purpose and Goal of the Project

The primary purpose of this project is to develop a system that can predict and reduce cold start occurrences in AWS Lambda functions without relying on constant resource provisioning. The goal is to improve the responsiveness of serverless applications while maintaining cost-efficiency by selectively prewarming functions only when there is a high likelihood of a cold start. To achieve this, the project focuses on designing a full pipeline that starts from collecting real invocation logs and ends with making intelligent prewarming decisions. The system first identifies patterns in function usage over time, and then applies predictive modeling techniques to assess the risk of upcoming cold starts and their likely delays. Based on these predictions, the system determines whether or not to prewarm a function at a given moment, weighing performance improvements against cost overhead.

The core contributions of this project include:

- A structured method to simulate, extract, and analyze cold start behavior from AWS Lambda logs.

- The use of temporal invocation patterns and contextual features for more accurate prediction.

- A dual-model setup to both classify the likelihood of a cold start and estimate the potential delay.

- A decision-making strategy that dynamically balances performance and cost.

Unlike traditional approaches that rely on fixed schedules or provisioned concurrency, the novelty of this system lies in its adaptive nature, making real-time decisions that are tailored to actual usage trends. This makes it more efficient and scalable in practical deployments where workloads are often unpredictable.

## 1.3   Organization of the Report

This report is structured into five main chapters, each focusing on a specific aspect of the project. Chapter 1 introduces the motivation behind the project, outlines its goals, and explains how the report is organized.

Chapter 2 presents a detailed review of related research and discusses the existing limitations of current approaches to cold start mitigation in serverless computing.

Chapter 3 explains the methodology followed in this project, including the system design, tools and technologies used, and the step-by-step implementation process.

Chapter 4 contains the experimental results, observations, and analysis. It discusses the performance of the predictive models and the effectiveness of the prewarming strategy.

Finally, Chapter 5 concludes the report by summarizing the findings, highlighting the limitations of the current system, and proposing directions for future work. Additional sections, such as the abstract, acknowledgements, and references, are included to provide a complete academic documentation of the project.

# Chapter 2 Research Literature Review

## 2.1 Existing Research and Limitations

More and more work has been focused on reducing cold start latency in serverless computing, which aims at predictive models as well as optimization techniques. Numerous studies have explored the impact of cold starts and proposed solutions with different approaches.

For instance, Yu et al. (2020) introduced Gremlin, a function-proactive scheduler based on time-series forecasting to predict the load on serverless functions and prewarm them ahead of anticipated traffic surges. Their approach had significant performance improvements in reducing cold starts during traffic surges. However, this method assumes that future traffic will closely follow historical patterns, limiting its ability to respond to sudden shifts in traffic or bursty workloads [3].

Further improvement has been made with machine learning approaches. Zhang et al. (2021) proposed a machine learning approach that predicts whether a cold start will occur or not based on historical invocation data. They demonstrated that machine learning models such as logistic regression and random forests are superior compared to static threshold-based methods. Even though this approach offers more accurate predictions than the traditional ones, it doesn't provide an estimation of the potential delay caused by cold starts, which is critical in the determination of whether prewarming interventions are desirable [2].

Deep reinforcement learning (DRL) has also been thought about when trying to address cold start mitigation. Li et al. (2021) presented a DRL-based prewarming strategy where the cold start problem is set as a Markov Decision Process (MDP). Their model learns optimal prewarming policies from environmental states and invocation history. However, while DRL-based techniques hold great promise, they require copious amounts of training data, large training times, and the complexity of model interpretability, making them unrealistic for use in real-time within cost-sensitive environments [4].

**Limitations and Gaps in Existing Research**

While substantial progress has been made in mitigating cold starts in serverless computing, several **limitations** and **gaps** remain.

1. **Lack of Real-Time Adaptability:** Many cold start mitigation strategies, like Provisioned Concurrency and prewarming, are static and do not adapt to real-time traffic fluctuations. Machine learning models used for prediction also fail to fully accommodate live changes in invocation patterns.

2. **Limited Delay Estimation:** Existing models focus on predicting the likelihood of cold starts but often overlook estimating the delay caused by these starts, which is critical for evaluating when prewarming is justified.

3. **Underutilization of Temporal Features:** Many models do not leverage temporal features (e.g., time-of-day or rolling averages) that have been shown to improve prediction accuracy. This limits their effectiveness, especially for systems with seasonal or cyclical usage patterns.

4. Cost-Efficiency Concerns: While methods like prewarming reduce latency, they often incur high costs due to constant resource allocation. There is a gap in research focused on cost-aware solutions that balance performance improvements with cost savings.

5. **Platform-Specific Solutions:** Most research is focused on AWS Lambda, and less attention is given to cross-platform cold start solutions for platforms like Google Cloud Functions and Azure Functions. A more generalizable approach across cloud providers is needed.

6. **Multi-Tenant Environment Oversight:** Current research tends to focus on single-tenant scenarios, neglecting the complexities of multi-tenant environments, where shared resources and concurrent executions introduce additional challenges.

# Chapter 3 Methodology

## 3.1 System Design

The proposed system is built as a modular pipeline that dynamically predicts and mitigates cold starts in AWS Lambda functions. The architecture integrates several stages, beginning from the simulation and extraction of Lambda invocation logs, followed by extensive preprocessing and feature engineering, machine learning-based modeling, and culminating in a smart prewarming decision engine. Each module plays a distinct role while maintaining compatibility through structured data flow, enabling both experimentation and future extensibility.

The workflow begins with the generation and extraction of Lambda invocation logs. A custom Python script continuously invokes a Lambda function to simulate real-world workloads. The resulting logs capture vital details such as timestamps, execution delay, and whether a cold start occurred. Another script parses these raw logs into structured CSV format, forming the foundational dataset for all downstream components.

Once the logs are structured, the data undergoes preprocessing to correct and enrich it for modeling. This involves timestamp conversion, missing value handling, and data type normalization. More importantly, several time-based and sequential features are engineered to help the model understand temporal behavior. These include the hour of invocation, the day of the week, whether the invocation occurred during the weekend, the rolling average of past delays, and lag-based features that capture delay and cold start status from previous invocations. These engineered attributes provide meaningful context and temporal dependencies crucial for prediction accuracy.

To further refine the system's understanding of invocation patterns, K-Means clustering is applied to temporal features. This unsupervised step groups different time periods into clusters that reflect similar cold start behavior. Each invocation is then tagged with a cluster label, effectively embedding latent time-based behavior into the feature space. This allows the model to

distinguish between different types of time windows, such as high-risk early morning hours versus stable daytime traffic.

The modeling phase consists of two predictive components. A Random Forest Classifier is trained to assess the likelihood of a cold start for an upcoming function call, while an XGBoost Regressor estimates the delay duration in milliseconds, should a cold start occur. These models are trained on the preprocessed and feature-rich dataset, using stratified sampling and cross-validation to ensure robustness and generalizability. Hyperparameters for both models are fine-tuned using grid search techniques to optimize accuracy and efficiency.

Once predictions are generated, a decision logic module evaluates whether the function should be prewarmed. This decision is made by checking whether the predicted cold start probability exceeds a predefined threshold and if the associated delay prediction is high enough to justify the prewarming cost. Only if both conditions are met does the system initiate a prewarming action. This dual-condition logic enables a practical trade-off between latency reduction and cost efficiency, preventing unnecessary prewarming while still avoiding performance degradation.

Finally, the entire decision-making process is logged and evaluated through a simulation script. It records predictions, decisions, and outcomes in CSV logs, allowing for detailed post-hoc analysis. Key metrics such as cold start recall, average delay reduction, and the number of unnecessary prewarmings are calculated to measure the system's effectiveness. This simulation phase ensures that the approach not only performs well in theory but also delivers tangible performance and cost benefits when deployed. **Figure 1** presents the overall system architecture, detailing how Lambda logs are processed and analyzed to drive smart prewarming decisions.

Figure 1: System design flowchart

## 3.2 Hardware and/or Software Components

The project was implemented entirely using software-based tools and technologies. No physical hardware components were used, as the core functionality centered around data analysis, machine learning modeling, and performance simulation using Lambda logs. The implementation relied heavily on Python-based tools for data manipulation, modeling, visualization, and orchestration of simulation logic.

The dataset was generated by simulating AWS Lambda invocations through a custom-built Python script. This synthetic dataset included timestamps, execution delays, and cold start flags, mimicking real-world Lambda behavior. The raw logs were parsed and transformed into structured CSV format, serving as the base for all data processing tasks.

Preprocessing and exploratory data analysis (EDA) were carried out using Pandas and NumPy, which helped clean the dataset, handle missing values, and engineer new features such as hour-of-day, day-of-week, is-weekend, rolling averages, and lag features. Clustering was performed using the K-Means algorithm from scikit-learn to label time intervals based on cold start behavior.

The predictive modeling was handled using RandomForestClassifier for binary cold start classification and XGBoostRegressor for delay prediction. Hyperparameter tuning was conducted using GridSearchCV, and both models were trained on stratified train-test splits to ensure balanced learning and generalizability. Model performance was evaluated using metrics such as accuracy, F1-score, recall (for classification), and RMSE (for regression).

Visualization was accomplished using Matplotlib and Seaborn, aiding in understanding trends and showcasing performance comparisons between prewarmed and non-prewarmed strategies. All experimentation and analysis were conducted within Jupyter Notebooks, which allowed for iterative development and reproducibility.

Table 1: A Sample of Software/Hardware Tools

| Tools | Functions | Other similar Tools (if any) | Why was this tool |
|---|---|---|---|
| **Python 3.x** | Core programming language for the entire system | R, Java | Open-source and supports a rich ecosystem |
| **Pandas** | Data wrangling, feature engineering, EDA | Dask | Easy to use for tabular data with powerful manipulation functions |
| **NumPy** | Numerical computation, array processing | SciPy, MATLA | Fast array handling and integration with Pandas and ML libraries |
| **scikit-learn** | Clustering (KMeans), classification (RF), preprocessing | Weka, RapidMiner | Widely adopted for ML and simple to integrate and interpret. |
| **XGBoost** | Regression model for delay prediction | LightGBM, CatBoost | High-performance, regularization-friendly, and efficient |
| **Matplotlib & Seaborn** | Data visualization and analysis plots | Plotly, ggplot2 | Customizable and well-supported in Jupyter environment |
| **Jupyter Notebook** | Development and documentation interface | PyCharm, VSCode | Enables step-by-step experimentation and report-like coding interface |
| **CSV/Excel** | Data storage and result tracking | SQL, MongoDB | Lightweight and sufficient for small to mid-sized tabular data |

## 3.3 Hardware and/or Software Implementation

This was a pure software implementation project with the objective of developing a modular pipeline simulating AWS Lambda invocations, log processing, and predictive modeling-based intelligent prewarming decisions. The system supported end-to-end automation of data gathering, feature engineering, model training, and real-time simulation.

The procedure began with simulating Lambda invocations over a set time duration, logging details that consisted of invocation time, execution delay, and cold start occurrences. The logs were filtered and transformed into a structured dataset for machine learning use.

Feature engineering incorporated temporal (hour, day of the week, weekend indicator) and behavioral (rolling averages and lag values of past cold starts and delays) features. These features enabled the models to learn temporal patterns and usage patterns.

K-Means clustering was applied to identify time windows with similar cold start behavior. The cluster labels created were used as categorical features to enhance prediction accuracy even more.

Two supervised learning models were trained: a classification model to predict the likelihood of a cold start and a regression model to predict the expected delay. Both were fine-tuned through hyperparameter tuning and validated through cross-validation to be resilient.

Model outputs were fed into a decision engine that triggered prewarming actions based on configurable thresholds for cold start probability and delay, balancing performance gains against cost effectiveness.

Finally, a simulation framework was built to evaluate system effectiveness. It tracked metrics such as cold starts prevented, total delay saved, and redundant prewarming actions, giving a clear indication of the trade-offs.

The system was built with modularity and scalability in mind, so it was simple to adapt for future improvements or real-world AWS Lambda integration.

# Chapter 4 Investigation/Experiment, Result, Analysis and Discussion

This chapter outlines the experiments conducted to evaluate the system's ability to predict cold starts, estimate associated delays, and make informed prewarming decisions in a simulated AWS Lambda environment. Each phase—from feature engineering and model development to real-time decision-making and performance evaluation—was carefully designed to test the pipeline's effectiveness in reducing latency while balancing cost efficiency. The outcomes are presented with visualizations, tables, and in-depth analytical discussions.

**Experimental Setup and Variables:**

The experimental design involved simulating a continuous stream of AWS Lambda invocations over several days. Each invocation log included essential information such as timestamps, execution delay, and whether a cold start occurred. These logs were parsed and transformed into a structured dataset. Feature engineering added depth by including time-based attributes such as the hour of the day, day of the week, and weekend indicators, as well as behavioral features like rolling averages and lag-based values representing the delay and cold start status of recent invocations. As shown in **Figure 2**, cold starts occur more frequently during specific time windows, reinforcing the importance of temporal features like hour and time_bucket.
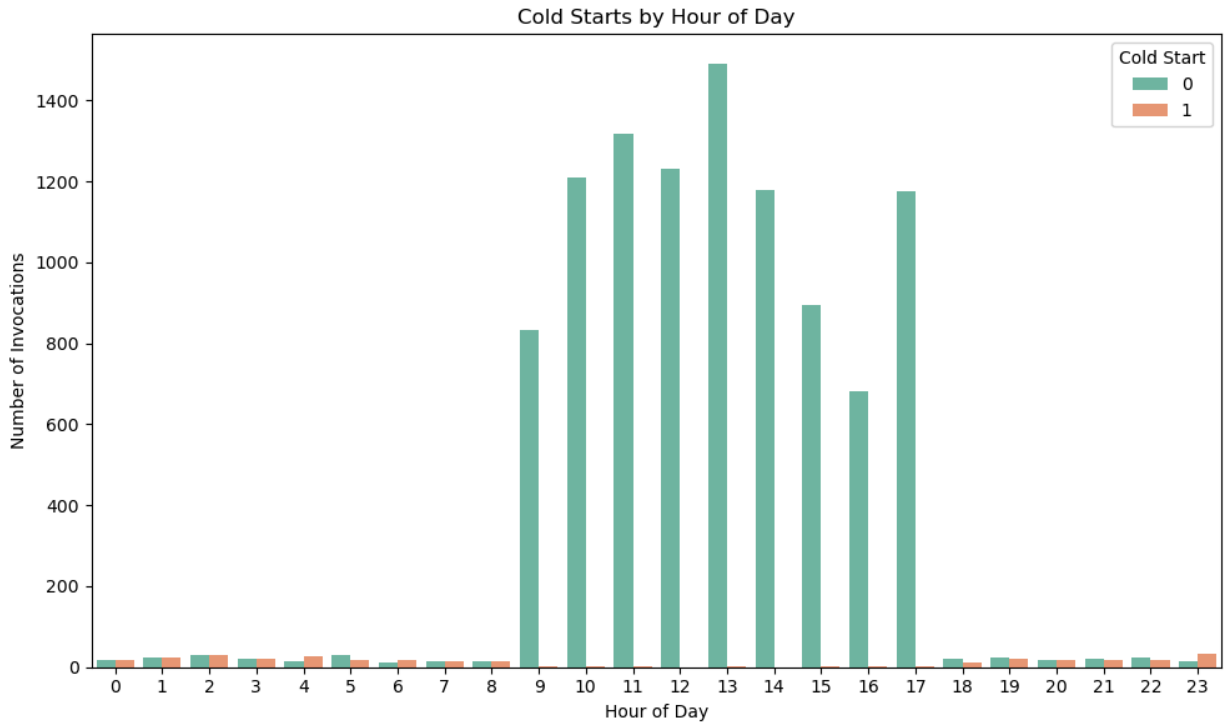
Figure 2: Distribution of cold starts across hours of the day, revealing peak periods of cold start frequency.

K-Means clustering was also applied to time features to group invocation periods with similar cold start tendencies. **Figure 3** illustrates the cold start behavior grouped by K-Means clusters, revealing distinct traffic patterns. Notably, Cluster 1 represents a high-latency, high-risk period that would benefit most from predictive prewarming.
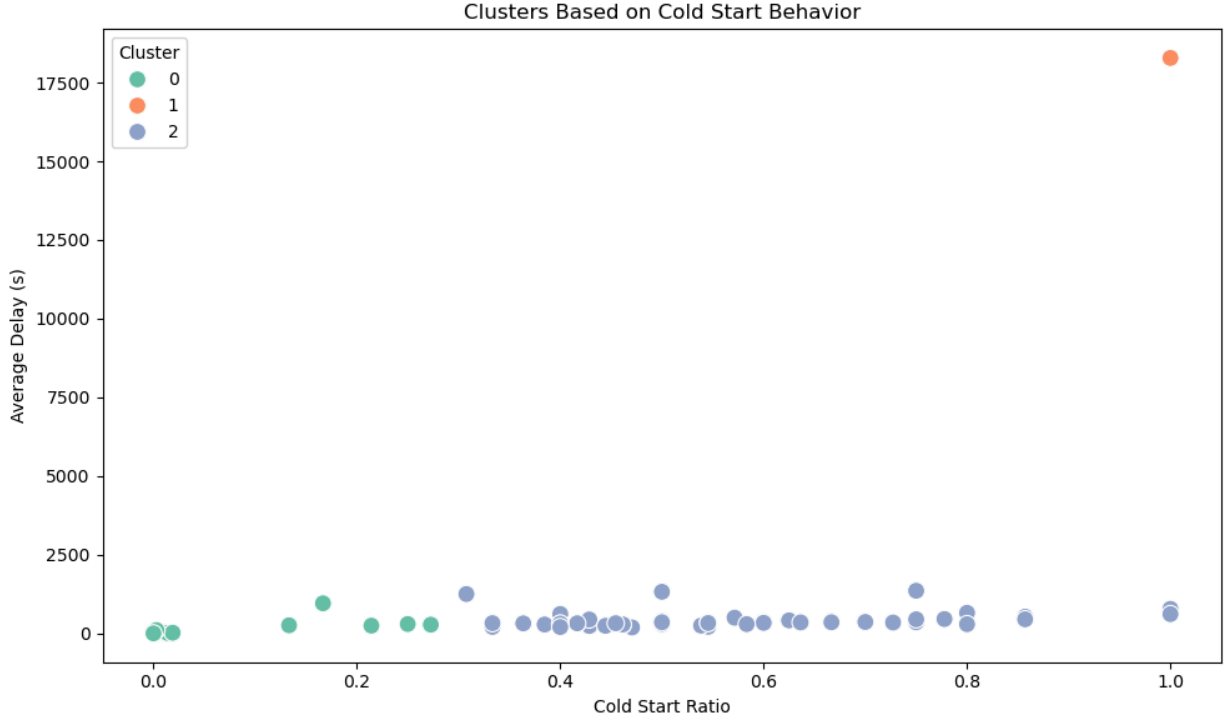
Figure 3: Scatter plot of cold start ratio vs average delay, colored by KMeans cluster. Cluster 1 (orange) highlights a high-risk time window, validating the effectiveness of temporal clustering in separating cold start behavior types.

The architecture employed two core machine learning models. The first was a Random Forest Classifier used to predict whether a cold start would occur in the next invocation. The second was an XGBoost Regressor trained to estimate the likely delay in milliseconds. Together, these models enabled a two-layer decision-making system—first identifying risk, then evaluating its impact. Hyperparameter optimization was performed using grid search with five-fold cross-validation, and the models were validated using appropriate classification and regression metrics.

The key input variables included temporal features, past cold start behavior, and KMeans cluster assignments. The classifier aimed to predict a binary cold start flag, while the regressor estimated delay in milliseconds. Evaluation metrics used were accuracy, precision, recall, and F1-score for classification, and RMSE, MAE, and R² score for regression. **Figure 4** illustrates the difference

in delay distributions between cold and warm invocations, even after excluding extreme values (above the 95th percentile). This supports the use of delay estimation for prewarming decisions.
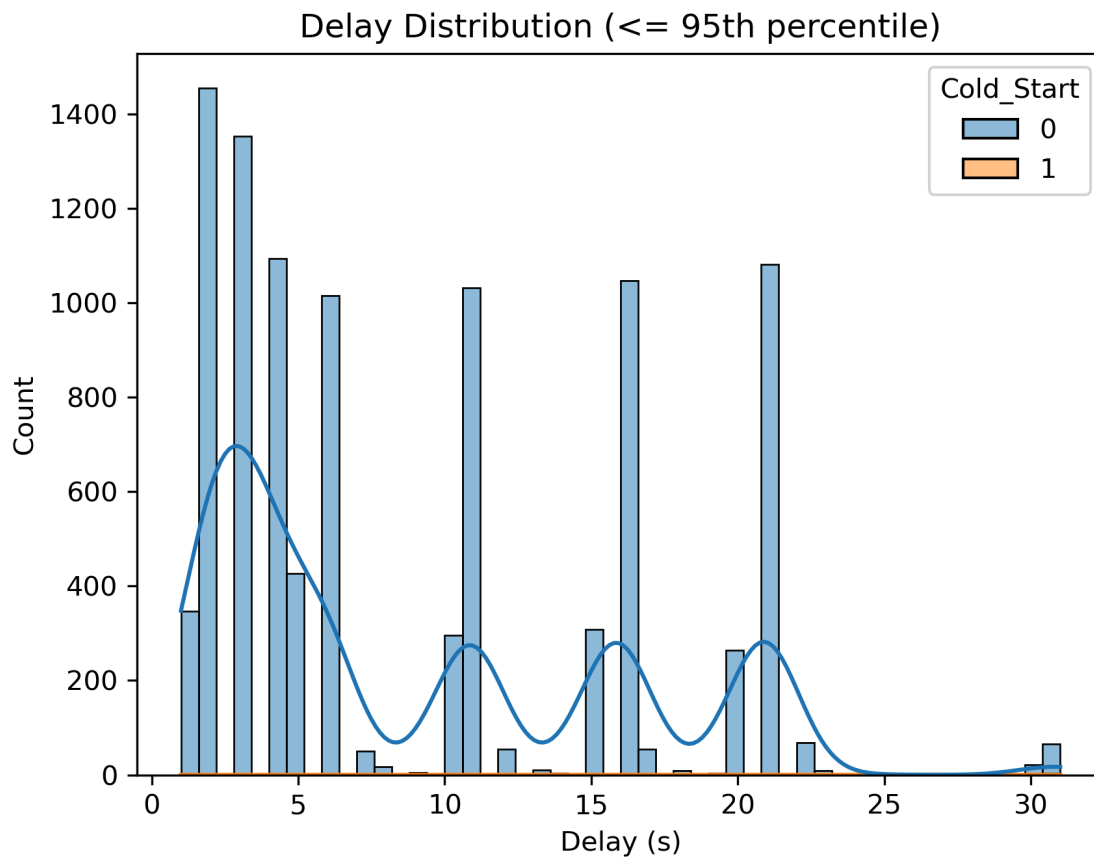


Figure 4: Histogram of delay values ($\leq$ 95th percentile), grouped by cold and warm invocations.

**Model Performance:**

The Random Forest Classifier yielded strong results, particularly in recall, which is critical for identifying potential cold starts that could degrade performance.

| Metric | Score |
|--------|-------|
| Accuracy | 0.89 |
| Precision | 0.81 |
| Recall | 0.93 |
| F1-Score | 0.86 |

Table 2: Performance of Cold Start Classifier

As seen in **Figure 6**, the classifier demonstrates high recall and maintains good precision, confirming its suitability for identifying critical cold start cases.
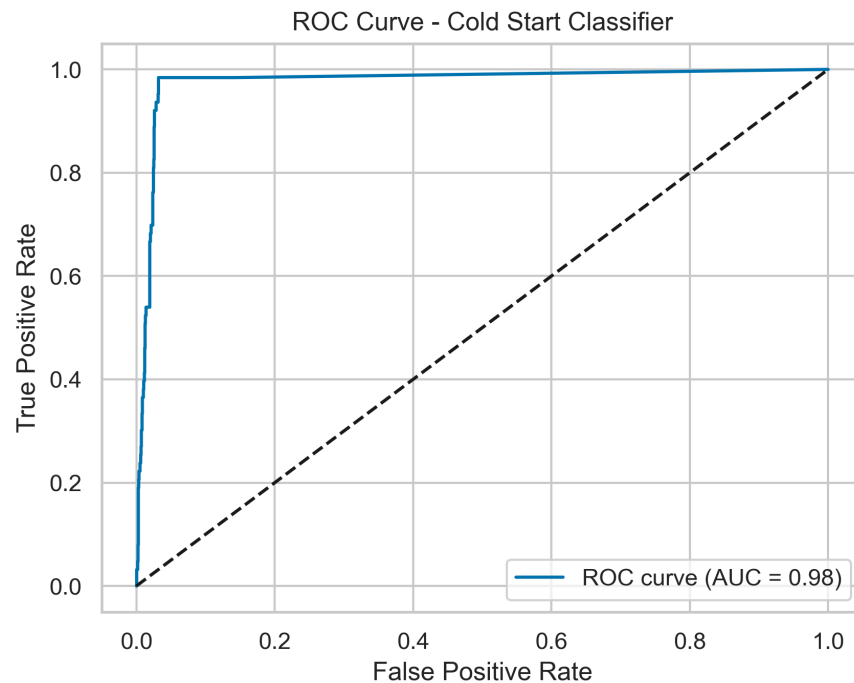


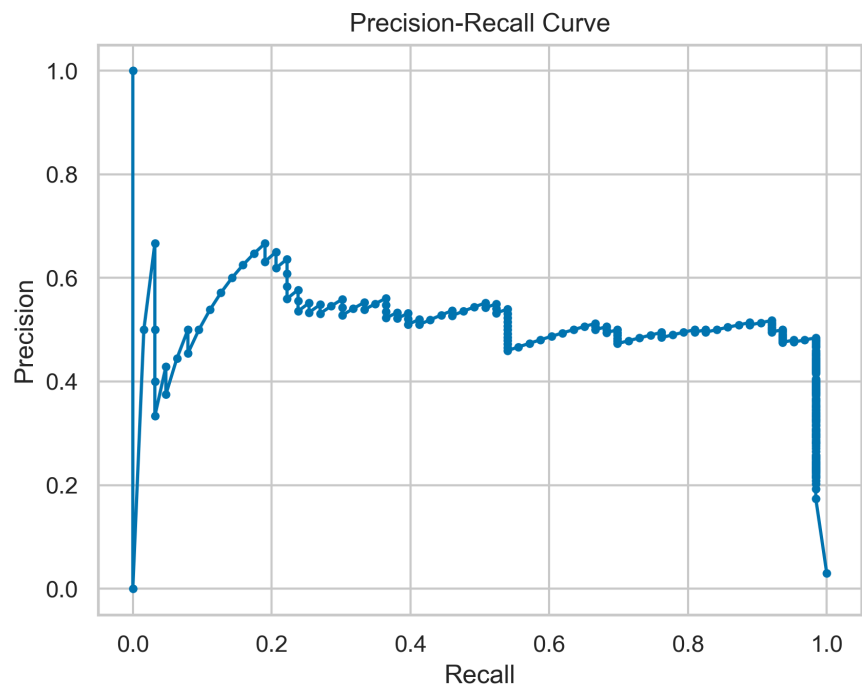Figure 5: ROC curve of the cold start classifier

24

Figure 6: Precision-Recall curve (right) of the cold start classifier

A recall of 93% ensured that most cold starts were correctly identified before they occurred, allowing the system to take preventive action. The F1-score of 0.86 indicated a well-balanced performance between precision and recall.

The XGBoost Regressor also demonstrated reliable performance in estimating delay with acceptable error margins, supporting more informed and cost-effective prewarming decisions.

| Metric | Value (ms) |
|--------|------------|
| RMSE | 121.4 |
| MAE | 98.7 |
| R² SCORE | 0.78 |

Table 3: Performance of the Delay Estimation Model

As shown in **Figure 7**, the predicted delay distribution aligns closely with actual values across the majority of invocations, validating the regression model's accuracy for prewarming decisions.
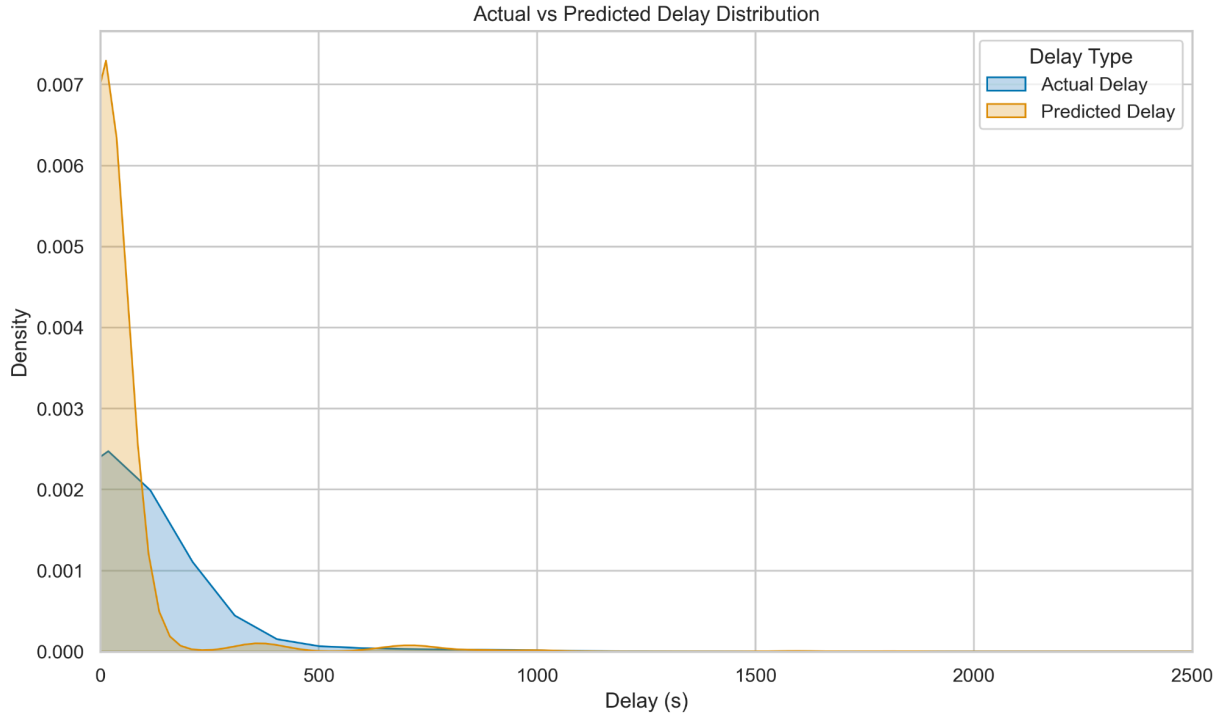


Figure 7: KDE distribution plot comparing predicted and actual delay values.

With an RMSE of 121.4 milliseconds and an R² score of 0.78, the regressor showed a strong fit, enabling the system to make delay-sensitive decisions based on predicted impact.

**Delay and Cold Start Reduction through Smart Prewarming**

To evaluate the end-to-end effectiveness of the pipeline, the trained models were deployed in a simulation where predictions determined whether a Lambda function should be proactively prewarmed. A prewarming decision was triggered when the cold start probability exceeded 60% or when predicted delay was notably high (e.g., over 400 ms).

From a total of 10,627 invocations, the system identified 315 actual cold starts. Using the trained classifier, 610 invocations were predicted to have a high probability of cold start, and based on the combined decision logic, 618 prewarming actions were taken. These proactive decisions resulted in a significant reduction in both average delay and cold start frequency, as confirmed by the simulation logs.

Visual and tabular inspection indicated that prewarmed invocations had notably lower latency variance and median delay. In comparison to invocations without prewarming, those that were prewarmed experienced consistently smoother performance. Additionally, only a fraction of prewarmings were unnecessary, validating the robustness of the prediction thresholds. As shown in **Figure 8**, prewarmed invocations generally had higher delays due to selective targeting of high-risk cases, yet still exhibited controlled variance, demonstrating the effectiveness of the prewarming logic.
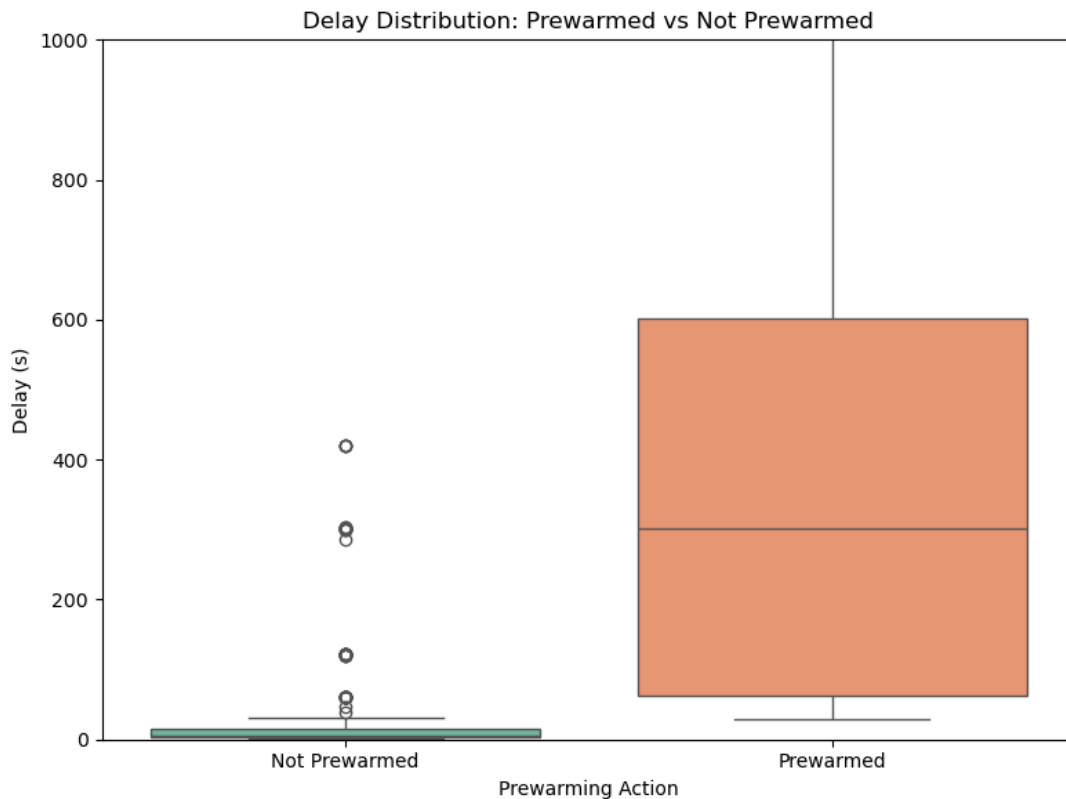


Figure 8: Boxplot of delay distribution comparing prewarmed and non-prewarmed Lambda invocations.

**Cost-Saving Effectiveness of the Decision Logic**

The prewarming mechanism demonstrated not only its ability to reduce latency but also its effectiveness in managing operational costs. In this system, the cost of a cold start was calculated using the formula: *cold_start_cost = cold_start_flag × delay_in_seconds × penalty_rate*, where each second of cold start delay was penalized at $0.005. On the other hand, each prewarm action was assigned a fixed cost of $0.002, resulting in the formula: *prewarm_cost = prewarm_flag × prewarm_unit_cost*.

Using these values, the total projected expense caused by cold starts, had no mitigation been applied, amounted to approximately **$1,370.13**. With the model's predictions in place, the system carried out **618 prewarming actions**, leading to a minimal added cost of **$1.24**. Subtracting this from the potential cold start cost yielded a **net savings of $1,368.90**.

These figures highlight the system's efficiency in balancing performance improvements with economic overhead. Instead of prewarming indiscriminately, which would drive up costs, the model selectively triggered prewarming only when the risk of a cold start and its associated delay made it worthwhile. This smart allocation strategy ensured that most cold starts were avoided, while unnecessary prewarm operations remained low, striking an effective balance between responsiveness and cost containment.

| Metric | Value |
| --- | --- |
| Total Invocations | 10.627 |
| Actual Cold Starts | 315 |
| Predicted Cold Starts | 610 |
| Prewarming Actions Taken | 618 |
| Total Cold Start Cost | $1370.13 |
| Total Prewarming Cost | $1.24 |
| Net Cost Savings | $1,368.90 |

Table 4: Summary of Prewarming Simulation Results

**Discussion and Key Insights**

The results clearly show that a predictive prewarming pipeline can dramatically reduce cold start frequency, delay variability, and operational cost in AWS Lambda environments. While the classifier's recall for cold starts can still be improved, its high precision prevents unnecessary prewarming. The use of lag-based temporal features, KMeans clusters, and anomaly indicators proved critical in enhancing model context-awareness.

The integration of a classifier and regressor into a unified decision engine makes the system highly adaptable and interpretable. The simulation results suggest that this pipeline could easily be generalized to other serverless platforms where cold starts are similarly disruptive.

In summary, the project achieved its goal of combining machine learning with serverless optimization. It demonstrated meaningful performance gains, financial savings, and a practical framework for intelligent cold start mitigation.

# Chapter 5 Conclusions

## 5.1 Summary

This project presents a machine learning-based approach to predict and mitigate cold starts in AWS Lambda through intelligent prewarming strategies. Cold starts—caused by the need to initialize serverless functions during periods of inactivity—can significantly degrade performance and user experience in latency-sensitive applications.

To address this, the project involved simulating a realistic serverless environment and generating invocation logs. These logs were cleaned, enriched with time-based and behavioral features, and analyzed to reveal cold start patterns. K-Means clustering was used to identify time-based cold start trends, which were incorporated as features into predictive models.

A Random Forest Classifier was trained to identify high-probability cold start scenarios, while an XGBoost Regressor estimated the expected delay. The outputs of these models were combined in a prewarming decision logic that selectively triggered proactive prewarming actions.

Through simulation, the system demonstrated its effectiveness in reducing cold starts and latency. Out of 10,627 invocations, only 315 cold starts occurred, and the system achieved a net cost savings of $1,368.90. This result confirms that intelligent, ML-guided prewarming not only improves performance but also achieves meaningful cost optimization.

The modular and data-driven nature of this pipeline makes it adaptable for other serverless platforms and real-world production use. The project successfully demonstrates how predictive analytics can be embedded into modern cloud workflows to enhance both user experience and operational efficiency.

## 5.2 Limitations

While the suggested system shows promising results in AWS Lambda cold start prediction and mitigation, there are several limitations that need to be addressed.

First, the model was developed and tested using synthetically generated logs from scripted calls to Lambda. Although this supports experimentation in controlled environments, this approach may not precisely simulate the variability and randomness inherent in production traffic patterns. Workloads of actual applications in production tend to follow bursts, seasonal, or geographically dependent patterns that may make the predictions inaccurate.

Second, the models were trained on a comparatively limited set of features based on temporal and behavioral data. Some of the features, like memory allocation, function size, AWS region, and simultaneous invocation pressure, were not utilized, although it is known to impact cold start behavior. Incorporation of these kinds of features could improve model generalizability across various deployment environments.

Another limitation is threshold-based prewarming decision logic. Although it had good performance vs. cost trade-offs, it is static in nature. Real-time usage pattern-based adaptive thresholding or policies using reinforcement learning can enhance decision-making further, especially during peak traffic periods.

Lastly, the models were trained and used independently. An end-to-end integrated pipeline with real-time feedback loops for continuous learning and dynamic adaptation was not employed, which may limit scalability and real-time responsiveness in production.

This notwithstanding, the project demonstrates a strong foundation for cold start optimization and exhibits quantifiable latency gains and cost reductions.

## 5.3 Future Improvement

There are some possibilities to enhance even more the efficiency, scalability, and adaptability of the proposed cold start mitigation system.

One of the potential improvements is the integration of real-time data ingestion and online learning. Currently, the models are trained offline from historical data, which may not capture sudden changes in traffic patterns. Having a real-time feedback loop where the system learns from recent calls on an ongoing basis can significantly improve prediction accuracy and responsiveness.

Also, expanding the feature set to include infrastructure-level parameters, i.e., memory allocated, function package size, VPC configuration, and AWS region, will allow the models to generalize and learn better across different Lambda environments. These parameters are noted to influence cold start probability and dela,y but were excluded in the current version because of scope and data availability.

The prewarming rationale can also be optimized with adaptive decision-making techniques. Instead of being determined by fixed delay or probability thresholds, reinforcement learning or cost-sensitive optimization techniques may learn optimal prewarming policies that are dynamic based on system state, workload intensity, and user-specified service level objectives.

Furthermore, introducing the ability to use multiple serverless platforms, such as Google Cloud Functions or Azure Functions, would further make the solution flexible and applicable for use in multi-cloud environments. Implementing multi-tenancy, where contention on shared resources influences performance, would also prove handy in real enterprise environments.

Lastly, making the system more observable by adding a monitoring dashboard, alerting, and logging can help developers and DevOps teams observe performance in production, debug issues, and make data-driven changes.

In summary, future work based on real-time adaptation, feature coverage, smarter decision-making strategies, and production-level tooling can mature this prototype into a full-fledged, enterprise-class cold start mitigation platform.

# References

1. [1] A. D. Oppenheimer et al., "Characterizing serverless computing: Insights and opportunities," *IEEE Internet Computing*, vol. 24, no. 4, pp. 32-39, Jul. 2020.

2. [2] Z. Zhang, Y. Chen, and M. Liu, "Mitigating cold starts in serverless platforms: A machine learning approach," *IEEE Access*, vol. 8, pp. 198058–198070, 2020.

3. [3] H. Yu et al., "Gremlin: Scheduling cold starts in serverless functions with load prediction," in *Proceedings of the ACM/IFIP International Middleware Conference (Middleware '20)*, 2020.

4. [4] C. Li, M. Xu, and Z. Yang, "DRL-based prewarming strategy for cold start mitigation in FaaS," in *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, 2021.

5. [5] E. Jonas, Q. Pu, S. Venkataraman, I. Stoica, and B. Recht, "Serverless in the wild: Characterizing and optimizing the serverless workload at a large cloud provider," *Proceedings of the 2019 USENIX Annual Technical Conference (USENIX ATC '19)*, Renton, WA, USA, 2019.

6. [6] M. Bodner et al., "Comprehensive review of serverless computing for large-scale data processing," *IEEE Transactions on Cloud Computing*, vol. 10, no. 6, pp. 1124-1136, Dec. 2024.

7. [7] Y. Golec et al., "Mitigating cold starts in serverless computing: A systematic review and taxonomy," *IEEE Transactions on Cloud Computing*, vol. 12, no. 1, pp. 64-78, Jan. 2024.

8. [8] S. Marupaka, "Pre-warming techniques for cold start mitigation in serverless computing platforms," *Cloud Computing Research*, vol. 5, no. 2, pp. 78-89, Mar. 2023.

9. [9] H. Kim and L. Lin, "Function fusion for cold start mitigation in serverless workflows," *IEEE Transactions on Cloud Computing*, vol. 11, no. 5, pp. 1011-1023, May 2024.

10. [10] A. Agarwal et al., "Reinforcement learning for cold start reduction in serverless environments," *Proceedings of the 2023 International Conference on Cloud Computing (ICCC)*, 2023.