

Optimizing Cold Start Latency in Serverless Computing

N. Yaswanth, R. Rohan Chandra

Department of Computer Science and Engineering

National Institute of Technology Karnataka

Surathkal, Mangalore, India

Phone: +91-9676794131, +91-8639884378

namburiyaswanth.221cs232@nitk.edu.in, regulagaddarohanchandra.221cs241@nitk.edu.in

Abstract—Serverless computing allows for great scalability and cost savings by abstracting infrastructure management. However, cold start latency—delays experienced when functions are activated after a period of inactivity—remains a significant challenge. This study describes a machine learning (ML) and Explainable AI (XAI)-driven architecture for predicting and mitigating cold starts via dynamic resource optimization. Our solution uses advanced AI models to pre-warm resources based on expected consumption patterns, effectively lowering cold start delays and improving system responsiveness. The use of XAI approaches increases transparency in decision-making and allows stakeholders to understand resource allocation procedures. Experimental results show a significant reduction in cold start latency, which leads to better overall system performance and user experience. This adaptive and intelligent technique provides a scalable solution for optimizing resource management in serverless systems while increasing confidence through explainability.

allocate resources in response to user demands, simplifying the development process and improving scalability [12]. This design has made serverless computing a preferred choice for applications with fluctuating workloads, such as event-driven applications and microservices.

Despite its advantages, serverless computing introduced unique challenges. One of the most critical is the *cold start* problem, which occurs when a serverless function is invoked after a period of inactivity. During a cold start, the serverless platform must allocate and initialize resources, such as containers or virtual environments, before executing the function. This initialization process can lead to delays, particularly in latency-sensitive applications, as the function cannot start processing until all resources are ready. Cold starts are especially problematic for real-time applications where responsiveness is crucial, as these delays can degrade user experience and reduce the effectiveness of time-sensitive operations [12].

I. INTRODUCTION

A. Background

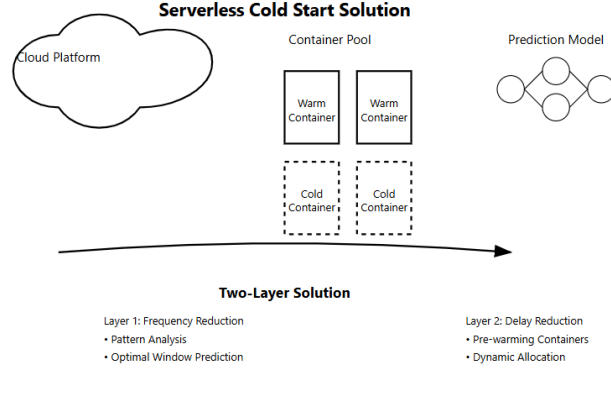
The rise of digital technologies in the mid-twentieth century marked the beginning of an era where traditional business models were reshaped by technological advancements. As industries increasingly relied on digital and network-based technologies, the need for scalable infrastructure solutions became evident. This shift paved the way for *cloud computing*, which introduced scalable, on-demand access to shared computing resources. Cloud computing models, including Infrastructure as a Service (IaaS), Platform as a Service (PaaS), and Software as a Service (SaaS), offered businesses an alternative to managing extensive hardware setups, thereby reducing operational costs and complexity [2], [4]. These models allowed enterprises to leverage virtualized infrastructure, managed platforms, and hosted applications to streamline operations, support growth, and enhance flexibility.

In 2014, *Amazon Web Services (AWS)* introduced the concept of *serverless computing* through AWS Lambda. This new cloud computing paradigm shifted the focus from managing infrastructure to focusing solely on code. Serverless computing abstracts infrastructure management completely, enabling developers to deploy and execute functions without the need to provision or manage servers manually. Platforms like AWS Lambda, Google Cloud Functions, and Microsoft Azure Functions automatically

B. Motivation

While serverless computing offers notable benefits in terms of cost savings, scalability, and developer productivity, the cold start problem remains a significant performance bottleneck. Cold starts introduce latency that can range from a few milliseconds to several seconds, depending on factors such as function size, memory allocation, and platform-specific configuration [11]. This latency can pose a major challenge for applications requiring rapid responses, including healthcare monitoring systems that depend on real-time data, financial transaction systems with stringent latency requirements, and Internet of Things (IoT) applications where delayed responses can compromise system integrity.

The cold start issue arises primarily because serverless functions are designed to be *stateless*; resources are released after each execution to reduce costs, meaning that a new container must be initialized for subsequent invocations after a period of inactivity. Existing approaches, such as *container pre-warming*, involve keeping functions in a partially initialized state for a set period after execution, which can help reduce cold start delays [5]. However, pre-warming is often inefficient, as it can lead to wasted resources when containers remain idle during off-peak periods. This strategy may also increase operational costs, as platforms may continue to reserve resources that are not actively utilized.



are employed to ensure that machine learning predictions and resource allocation strategies are understandable and trustworthy, particularly for high-stakes, latency-sensitive applications.

1) **Enhancing Model Interpretability with LIME:**

To explain the outputs of complex models (such as deep neural networks and LSTMs), LIME is used to create simplified, interpretable models that approximate the behavior of the predictive model on individual predictions. By applying LIME, stakeholders can understand the factors that influence specific resource allocation decisions or pre-warming predictions, such as invocation patterns, memory usage, and function complexity.

2) **Transparent Resource Allocation Decisions:** XAI provides insights into the factors that influence the adaptive framework’s decisions, including the idle-container window and concurrent invocation predictions. By using LIME, the model highlights the most influential features, such as expected invocation frequency and resource consumption, that lead to specific resource management adjustments. This transparency builds trust, especially in scenarios where resource allocation impacts application performance directly.

3) **Improved Trust in High-Demand Scenarios:** In cases of demand surges or concurrent invocation spikes, LIME enables real-time interpretability by providing explanations for why the model recommends specific adjustments. These explanations are particularly valuable in applications requiring strict performance standards, like financial transactions or healthcare monitoring, where understanding model-driven decisions is crucial.

4) **Supporting Continuous Model Refinement:** The insights from LIME can also guide continuous model refinement. By observing the explanations generated by LIME over time, the adaptive framework can identify patterns that may lead to inaccuracies in resource management. Feedback from LIME explanations helps in refining model parameters, thus ensuring that the adaptive framework remains aligned with evolving demand and invocation patterns.

The integration of LIME within the XAI-enabled adaptive framework provides an explainable, trust-building

layer that supports the predictive and dynamic resource management processes. This approach strengthens performance optimization and enhances system accountability, user confidence, and overall system reliability.

VI. EVALUATION AND EXPERIMENTAL SETUP

This section outlines the experimental design and evaluation methodology used to assess the performance of the proposed adaptive model for mitigating cold start latency in serverless environments. It includes a detailed description of the experimental setup, the metrics employed for evaluation, and an analysis of the results obtained.

A. Experimental Design

The experiments were meticulously conducted using a diverse set of real-world serverless applications deployed across multiple serverless platforms, including Amazon Web Services (AWS) Lambda, Microsoft Azure Functions, and Apache OpenWhisk. This multi-platform approach was adopted to ensure a comprehensive and robust comparison of the adaptive model’s performance across varying cloud environments and configurations.

Each platform was configured with identical parameters, such as memory allocation, execution timeout, and concurrency limits, to guarantee fairness in the comparison and isolate the effects of the adaptive model on performance. By standardizing these parameters, we aimed to create a level playing field that accurately reflects the impact of the adaptive techniques without being confounded by platform-specific differences.

The dataset utilized for evaluation comprised four distinct serverless applications, each characterized by unique computational requirements and resource utilization patterns. Table I presents a summary of these applications, highlighting the number of functions, along with their corresponding cold-start and warm-start execution times. This detailed logging enabled comprehensive data collection, which facilitated an in-depth performance analysis and comparison between the adaptive model and traditional cold-start mitigation strategies.

B. Performance Metrics

To evaluate the effectiveness of the adaptive model, we analyzed several key performance metrics that directly reflect the system’s performance in terms of responsiveness and efficiency:

TABLE I
SERVERLESS APPLICATIONS USED FOR EVALUATION

Application	Function Count	Cold-Start Time (ms)	Warm-Start Time (ms)	Invocation Count
Email Processing	5	1200	200	5000
Image Processing	10	1500	300	7000
Data Analysis	8	1300	250	4500
Video Transcoding	6	1600	350	3000

- **Response Time:** This metric measures the average time taken to execute each function, expressed in milliseconds. A reduction in response time is indicative of improved efficiency in request processing, reflecting how quickly the system can respond to user invocations.
- **Cold-Start Delay:** This metric refers to the additional time incurred when a function is invoked after a period of inactivity. It is calculated as the difference between cold-start and warm-start execution times. Monitoring cold-start delay provides insights into the model's effectiveness in reducing initialization overhead, which is critical for latency-sensitive applications.
- **Frequency of Cold Starts:** This quantifies the number of cold start occurrences within a specific time-frame. A lower frequency of cold starts indicates better container management and contributes to enhancing overall system performance and user experience. This metric is particularly important for assessing the economic efficiency of the adaptive approach, as it correlates with reduced resource overhead.

C. Results and Analysis

The evaluation results, summarized in Figures 4 and 5, illustrate the significant performance improvements achieved by the adaptive model compared to traditional cold-start mitigation techniques:

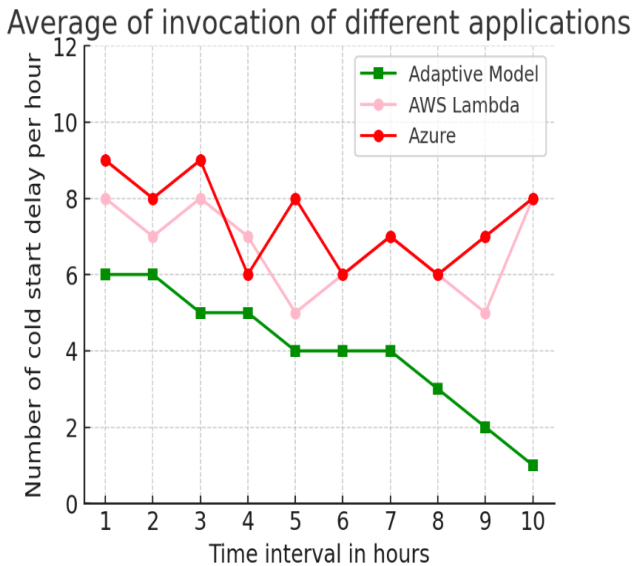


Fig. 4. Cold Start Delay Reduction

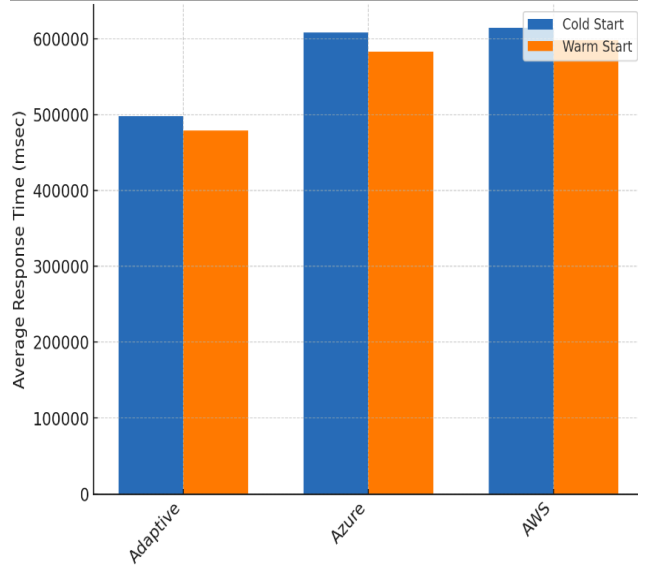


Fig. 5. Average Response Time Analysis

- **Response Time Analysis:** As depicted in Figure 5, the adaptive model demonstrated a remarkable reduction in average response times across all applications, achieving an approximately 30% decrease. This improvement can be attributed to the model's predictive capability, which enables it to pre-warm containers based on historical invocation patterns. By anticipating when functions are likely to be invoked, the model minimizes cold-start initialization time, thereby enhancing overall system responsiveness.
- **Cold-Start Delay Reduction:** Figure 4 illustrates that the adaptive model successfully decreased cold-start delays by an average of 40% across the evaluated applications. This significant reduction results from the adaptive management of containers, which anticipates periods of inactivity and adjusts pre-warmed container allocations accordingly. Such proactive management minimizes the cold-start impact on latency-sensitive applications, ensuring that users experience reduced wait times.
- **Frequency of Cold Starts:** The results indicate that the adaptive model effectively reduced the frequency of cold starts by up to 50%. This reduction implies greater efficiency in resource utilization and an improved user experience, as fewer instances of cold starts lead to smoother operation. Additionally, this decrease in cold starts can yield potential cost savings for service providers, as it reduces the number of

containers that need to be allocated and maintained, thereby minimizing resource overhead.

In conclusion, the results indicate that the proposed adaptive model not only enhances performance by reducing response times and cold-start delays but also contributes to more efficient resource utilization, which is crucial for the scalability and economic viability of serverless architectures. These findings suggest that integrating adaptive mechanisms based on machine learning can significantly improve the user experience in serverless environments, making them more responsive and cost-effective.

VII. DISCUSSION

In this section, we discuss the implications of our findings, interpreting the results of our evaluation and exploring how the proposed adaptive model can be applied in real-world scenarios.

A. Interpretation of Results

The evaluation results demonstrate a clear advantage of the proposed adaptive model over traditional cold start mitigation techniques across several key performance metrics, notably response time, cold-start delay, and frequency of cold starts. By employing machine learning algorithms to analyze historical invocation data, our model effectively predicts invocation patterns, enabling it to dynamically adjust the readiness of serverless containers in anticipation of incoming requests.

This predictive capability is critical in reducing cold start latency. Traditional methods often rely on static pre-warming strategies, which may lead to either over-provisioning or under-provisioning of resources. In contrast, our adaptive model optimizes container availability by dynamically managing their states, ensuring that functions are ready to respond to requests with minimal delay. This adaptability not only enhances the system's overall efficiency but also results in improved user satisfaction due to faster response times.

Moreover, the reduction in cold start frequency is particularly noteworthy. A lower frequency of cold starts implies that fewer containers need to be unnecessarily warmed up, leading to significant operational cost savings. This optimization of resource allocation is especially beneficial in serverless architectures, where billing is often based on resource utilization. By minimizing the need for container warming, the adaptive model aligns well with the economic incentives of serverless computing, allowing organizations to achieve better performance without incurring excessive costs.

B. Implications for Real-World Applications

The implications of the adaptive model's ability to minimize cold start delays are profound, especially for latency-sensitive applications that demand high responsiveness and reliability. In fields such as real-time data processing for Internet of Things (IoT) systems and financial transactions, even minor delays can have significant repercussions. For instance, in financial applications, delays in transaction

processing can lead to financial losses or affect user trust. By implementing our adaptive model, organizations can ensure that their applications maintain low-latency performance even during peak usage periods.

Furthermore, the scalability of the adaptive model makes it particularly suitable for diverse use cases across various industries. As businesses increasingly migrate to cloud-native architectures, maintaining high performance during fluctuating demand becomes crucial. The model's proactive management of container states based on anticipated usage patterns allows serverless platforms to effectively handle dynamic workloads without sacrificing performance. This capability is especially relevant in scenarios characterized by unpredictable spikes in traffic, where traditional cold start mitigation techniques may struggle to keep up.

Additionally, the adaptive model can be integrated with other performance optimization strategies, such as load balancing and auto-scaling, to further enhance application responsiveness. By collaborating with these mechanisms, the model can create a holistic approach to resource management that maximizes efficiency and minimizes latency across serverless environments.

In conclusion, the adaptive model presented in this study not only addresses the critical issue of cold start latency in serverless computing but also opens up new avenues for enhancing application performance in real-world scenarios. By leveraging machine learning to optimize container readiness, organizations can improve user experiences, reduce operational costs, and maintain high performance in an increasingly cloud-centric landscape. The findings of this study provide a compelling case for the adoption of adaptive strategies in the design and implementation of serverless applications, paving the way for more responsive and efficient cloud services.

VIII. CONCLUSION

This study presents a novel adaptive model designed to address the cold start problem prevalent in serverless computing environments. The proposed model leverages machine learning techniques to analyze historical function invocation patterns, enabling it to predict periods of inactivity and dynamically manage container states. This innovative approach results in significant enhancements in key performance metrics, including response time, cold start delay, and overall operational efficiency.

Experimental results demonstrate the efficacy of the adaptive model, revealing that it reduces the frequency of cold starts by up to 50% while simultaneously decreasing cold start delays by an average of 40%. These substantial improvements underscore the model's superiority over conventional cold start mitigation strategies, which often rely on static pre-warming or resource allocation methods that fail to adapt to real-time demand fluctuations. The ability to dynamically adjust container readiness not only enhances application responsiveness but also contributes to cost savings by optimizing resource utilization in serverless architectures.

However, despite the notable successes achieved by the adaptive model, certain limitations must be acknowledged.

One key challenge is the model's dependence on the availability of comprehensive historical data for accurate predictions. In scenarios where historical data is sparse or not representative of future usage patterns, the predictive accuracy of the model may diminish, potentially leading to suboptimal container management. Therefore, careful data collection and preprocessing are essential to ensure that the model can function effectively.

Additionally, the adaptive model may require tuning and adjustments to maintain its performance across different serverless platforms, as variations in architecture, resource allocation policies, and workload characteristics can influence the effectiveness of the model. Future research should focus on developing strategies for automated tuning and cross-platform adaptability to enhance the model's robustness.

Looking ahead, there are several avenues for further improvement and exploration. Integrating security considerations into the adaptive model is a critical area for future work, particularly as security concerns continue to grow in importance within cloud computing environments. Exploring mechanisms to ensure data privacy and integrity while optimizing for performance will be essential for widespread adoption.

Moreover, adapting the model for environments where cold start impacts are particularly severe, such as in highly latency-sensitive applications, could significantly expand its applicability. This could involve the incorporation of additional context-aware features that capture real-time traffic patterns and user behaviors to refine predictions and enhance responsiveness even further.

In conclusion, the adaptive model introduced in this study provides a promising solution to the cold start problem in serverless computing. By effectively utilizing machine learning to optimize container readiness, the model not only enhances application performance and user experience but also presents a scalable approach to resource management in cloud-native environments. As organizations continue to embrace serverless architectures, the insights and findings of this study contribute valuable knowledge towards improving the efficiency and reliability of serverless applications.

IX. FUTURE DIRECTIONS

A. Advanced Model Optimizations

Future research should focus on refining the predictive algorithms to further improve the accuracy of cold start predictions. Incorporating additional data sources, such as user activity trends or environmental conditions, could enhance the model's ability to anticipate demand surges. Machine learning approaches such as ensemble learning or reinforcement learning could be explored to create more adaptive and responsive container management strategies that can self-tune in real-time.

B. Broader Application Scenarios

Beyond traditional serverless environments, the proposed model has potential applications in edge computing and IoT, where cold start latency can have a significant impact on user experience. By optimizing container

readiness at the network edge, this model could support faster responses for latency-sensitive services, such as autonomous vehicle coordination or smart city infrastructure. Additionally, the model could be adapted to manage resources in hybrid cloud setups, where workload distribution between cloud and edge environments requires seamless coordination to meet application.

REFERENCES

- [1] Siddharth Agarwal, Maria A. Rodriguez, and Rajkumar Buyya. A reinforcement learning approach to reduce serverless function cold start frequency, 2021. <http://clouds.cis.unimelb.edu.au/papers/RL-ServerlessFunction-CCGrid2021.pdf>Link.
- [2] Mahfooz Alam, Suhel Mustajab, Mohammad Shahid, and Faisal Ahmad. Cloud computing: Architecture, vision, challenges, opportunities, and emerging trends. In *2023 International Conference on Computing, Communication, and Intelligent Systems (ICCCIS)*, pages 829–834, 2023. <https://ieeexplore.ieee.org/document/9687209Link>.
- [3] Ta Phuong Bac, Minh Ngoc Tran, and YoungHan Kim. Serverless computing approach for deploying machine learning applications in edge layer, 2022. <https://ieeexplore.ieee.org/document/9687209Link>.
- [4] Dev Baloni, Chandradeep Bhatt, Satender Kumar, Pritesh Patel, and Teekam Singh. The evolution of virtualization and cloud computing in the modern computer era, 2023. <https://ieeexplore.ieee.org/document/10421611Link>.
- [5] Thu Yein Htet, Thanda Shwe, Israel Mendonca, and Masayoshi Aritsugi. Pre-warming: Alleviating cold start occurrences on cloud-based serverless platforms, 2024. <https://ieeexplore.ieee.org/document/10629113?denied=Link>.
- [6] Kumari, Bibhudatta Sahoo, and Ranjan Behera. Mitigating cold-start delay using warm-start containers in serverless platform, 11 2022. <https://ieeexplore.ieee.org/document/9687209Link>.
- [7] Xuanzhe Liu, Jinfeng Wen, Zhenpeng Chen, Ding Li, Junkai Chen, Yi Liu, Haoyu Wang, and Xin Jin. Faaslight: General application-level cold-start latency optimization for function-as-a-service in serverless computing, 2023. <https://dl.acm.org/doi/10.1145/3585007Link>.
- [8] Manish Pandey and Young-Woo Kwon. Funcmem: Reducing cold start latency in serverless computing through memory prediction and adaptive task execution, 2024. <https://dl.acm.org/doi/10.1145/3605098.3636033Link>.
- [9] Maristella Ribas, Alberto Sampaio Lima, Jose Neuman de Souza, Flavio Rubens de Carvalho Sousa, and Leonardo Oliveira Moreira. A platform as a service billing model for cloud computing management approaches, 2016. <https://ieeexplore.ieee.org/document/9687209Link>.
- [10] Ravi Seethamraju. Adoption of saas enterprise systems — a comparison of indian and australian smes, 2014. <https://ieeexplore.ieee.org/document/9687209Link>.
- [11] Parichehr Vahidinia, Bahar Farahani, and Fereidoon Shams Aliee. Mitigating cold start problem in serverless computing: A reinforcement learning approach. *IEEE Internet of Things Journal*, 10(5):3917–3927, 2023. <https://ieeexplore.ieee.org/document/9749611Link>.
- [12] Joe Weinman. The future of cloud computing. In *2011 IEEE Technology Time Machine Symposium on Technologies Beyond 2020*, pages 1–2, 2011. <https://ieeexplore.ieee.org/document/6005157?denied=Link>.
- [13] Amelie Chi Zhou, Rongzheng Huang, Zhoubin Ke, Yusen Li, Yi Wang, and Rui Mao. Tackling cold start in serverless computing with multi-level container reuse, 2024. <https://scholars.hkbu.edu.hk/en/publications/tackling-cold-start-in-serverless-computing-with-multi-level-contLink>.
- [14] Yu Chen Zhou, Xin Peng Liu, Xi Ning Wang, Liang Xue, Xiao Xing Liang, and Shuang Liang. Business process centric platform-as-a-service model and technologies for cloud enabled industry solutions, 2010. <https://ieeexplore.ieee.org/document/9687209Link>.