# A Scalable Cloud-Based CRM System Using Serverless Microservices and Containerised AI Recommendation Engines

**Saad Khan, Solution Architect Leader of Global Investment Banking, Senior IEEE Member**

**Abstract** CRM systems now use cloud-native technologies, and there has been a shift in the architecture to modular and highly scalable, with many supporting backends and integrations. The given paper will introduce a scalable cloud-based CRM solution that involves serverless microservices and containerized AI-based recommendation engines. The suggested architecture enables the CRM to be highly available, cost-efficient, real-time, and scalable, which resolves the performance bottlenecks of traditional CRM systems. Our technology enables us to decouple traditional CRM features into fully deployable serverless services, such as customer data ingestors, analytics, and marketing automation. These functions connect with each other through event-based approaches and API gateways, rendering them loosely coupled and simpler to maintain. To provide personalized recommendations, we package the AI models in Docker containers and can easily manage them by using Kubernetes, so that the optimal usage of resources and portability can be ensured. We conducted extensive trials with the top cloud providers, including AWS and GCP, to compare their performance standards in terms of latency, throughput, cost-effectiveness, and AI prediction accuracy. The results led to a 40 percent lower cost of operation and a 50 percent faster response rate of our system than the monolithic ones. The paper also addresses the issues of design trade-offs, implementation problems, future work, along with key Privacy and Security Risks to be considered. As our evidence shows, using serverless microservices and containerized AI engines would provide an attractive paradigm for next-generation CRM platforms.

*A. Index Terms—Cloud Computing, Serverless Architecture, CRM System, Microservices, Containerization, AI Recommendation Engine, Kubernetes, Scalability, Event-Driven Architecture, Personalization*

Note: There should no nonstandard abbreviations, acknowledgments of support, references or footnotes in in the abstract.

## I. INTRODUCTION

CRM systems are crucial in enabling a business to establish and maintain strong relationships with its customers. Centralization of customer data, automation of marketing processes, correspondence across multiple channels, and personalization of interactions are some of the advantages that customer experience platforms can bring to any product or service [1-3], making the former crucial to the business success of the latter. Yet, the conventional CRM technologies are frequently implemented with a monolithic architecture to integrate all the elements into an intricately related system. Although monolithic systems are operable at the beginning, they soon become complex to scale, update, and integrate with newer systems as their business expands. This issue becomes more prominent as they strive to expand their systems. It is costly to maintain, and performance bottlenecks may occur due to the inflexibility of deployment. Due to the growing popularity of cloud computing and distributed systems, a new trend has emerged to shift CRM applications toward a cloud-native, microservices-based architecture. The new systems divide functionality into loosely coupled services that are independent of one another and that may be developed, deployed, and scaled separately. Such a scalable and modular design increases system resilience and fault tolerance in addition to the minimization of operating expenses due to optimal distribution and utilization of resources, and usage models, like the one provided by serverless computing. Furthermore, by incorporating AI and real-time analytics into this

architectural scheme, one opens up new possibilities in terms of dynamically personalizing and deciding on data-driven actions. The driving force behind this study is the need to find and confirm a contemporary CRM solution that can tap into the immense capabilities of cloud native platforms and serverless computing, as well as artificial intelligence and machine learning-based intelligence-active elements to circumvent and resolve the shortcomings of earlier CRM implementations and provide a powerful, efficient, and intelligent system of customer relations.

**1.1. Need for Cloud-Based Serverless CRM**

The current business model, with a shift to being digital-first, creates a gap between the requirements of a modern company and the limitations of traditional CRM architecture in terms of its scalability, adaptability, and affordability. The rising star is CRM as a cloud-based service with serverless resource provisioning, which can address these challenges by leveraging contemporary computing paradigms. The need for such a system is evident in the following subsections.



Figure 1: Need for Cloud-Based Serverless CRM

- **Flexibility and Scalability:** Many CRM systems have lacked scalability due to their monolithic design, resulting in a shortage of flexibility. As the client base increases or the rate of traffic surges during promotions, such systems become inefficient to the point of needing high-end infrastructural advancements to accommodate high volumes. Unlike it, a cloud-based serverless CRM can automatically scale specific functions according to the required demand, thereby maintaining

customary performance through any amount of workload. It is this elasticity that enables organizations to be flexible to the changes in the market and the behavior of users.

- **Optimization of Cost:** IoT Many of the insights of serverless computing lead to the introduction of a pay-per-usage billing model, where user charges are consistent with the actual computing processing time, as opposed to the pre-allocation of resources to which they are normally accustomed. The model significantly reduces the cost of unused infrastructure, particularly in niche or sporadic functions such as consumer segmentation, lead scoring, or campaign trigger functionality that are specific to CRM processes. Hosting options like AWS, GCP, and Azure also minimize the operational expenses since they automatically manage the provisioning and maintenance of servers and the associated scaling requirements.

- **Faster Development and Deployment:** Serverless CRM architectures enable many functions and services to be deployed to production continuously and frequently, which CI/CD requires. The developer can make more frequent updates, introduce new features, and support with fewer risks, and have smaller and more manageable codebases. This modularity enables faster introduction of innovations and accelerates time-to-market for new CRM features and enhancements.

- **Integration of AI Seamlessly:** New CRM systems should do more than store the data and automize the working process; they should provide an intelligent and personalized experience. A serverless CRM running within a cloud should be easily connected to any Artificial Intelligence as a Service (AIaaS) platform, such as AWS SageMaker, Google Vertex AI, or Azure Cognitive Services. Such integrations support real-time prediction, recommendations, and automation without requiring redesign of the backbone architecture.
- **Enhanced Resilience and Gap-Resistance:** Serverless platforms are inherently highly available, fault-tolerant, and resilient to distributed execution across multiple data centres. This ensures that the CRM services are highly available and unaffected by failures. When combined with microservices and containerised AI module deployment, the serverless architecture minimises single points of failure and enhances system robustness in general.

1)

2) *1.2. Role of AI in CRM*

AI is transforming the traditional Customer Relationship Management (CRM) beyond the stagnant structure of data warehousing and process flow, and now has intelligent applications that adapt and are proactive in serving a customer base. [4,5] Personalization is one of the most powerful AI in CRM applications. It is possible to use AI algorithms to analyze huge amounts of historical data on customer interactions, purchasing history, browsing history, and demographic data and create truly personalized content, offers and communication. Such a high degree of personalization contributes to an appreciable increase in customer satisfaction and retention; the user is more inclined to respond to recommendations and messages that seem to be individual in nature. The second major application of AI in CRM is predictive behaviour analysis. Machine learning algorithms can detect patterns in customer behavior and predict future behaviors (including the likelihood of purchase, whether they are at risk of churning, or what opportunities lie in upsell). Such predictive observations enable sales and marketing departments to prioritise leads, take proactive actions to resolve any issues, and make informed and insightful decisions that align with customer intentions. For example, when it comes to lifecycle stages, AI can segment the customer base and automate targeted outreach that reaches the most eligible potential customers, thereby optimising returns on resources used. CRM is also possible with intelligent automation enabled by AI. Natural Language Processing (NLP) enables chatbots and virtual assistants to automatically address typical customer requests, process support requests, and guide users through the onboarding process or problem-solving —all within a short amount of time. This reveals some of the task load on human agents and enhances response time and consistency. Additionally, recommendation engines utilise collaborative filtering, neural networks, and decision trees, enabling dynamic suggestions of products, content, or actions based on the user's context to enhance the possibilities of cross-selling and up-selling. On the whole, AI makes CRM systems highly intelligent and flexible platforms that continually learn and develop based on user interactions. To provide a more personalised customer experience and make it more efficient, the use of AI becomes not only advantageous but critical to gaining a competitive edge in the modern data-driven market environment.

## II. Literature Survey

### 2.1. Traditional CRM Architectures

Historically, Customer Relationship Management (CRM) systems, which were much older, such as those offered by Salesforce or Zoho CRM, formed the basis of customer management strategies in enterprises. Such systems are commonly constructed in a monolithic design, where all parts of the application are tightly integrated within a single, closely bonded code. [6-9] Whereas this solution is unified and eases the initial development, it also creates a number of issues in the long run. The close integration of services can be problematic, requiring a redeployment of the entire application when only a minor change is needed, and it can also be problematic when scaling services. This makes the system more rigid, which drives up the cost of maintenance and limits its flexibility, especially in cases where it must adapt to quickly changing business requirements or incorporate new technologies that emerge.

### 2.2. The Cloud Migration of CRM Suite

Thanks to the advent of cloud computing, most CRM vendors today have transitioned to cloud-native systems that are designed for vendor extensibility, scalability, and resilience. The emergence of platforms such as Freshworks can be attributed to this trend, as they implement the microservice design, which means that CRM is broken down into a set of loosely coupled services that can be developed, deployed, and scaled separately. This paradigm involves fewer innovation cycles, making the isolation of faults easy; hence, systems are more agile and maintainable. As shown

organizations that utilize cloud-native CRM architecture have reduced their time-to-market by 30 percent, their modularity, and system responsiveness. These advantages support the fact that the cloud-native design addresses the evolving needs of current CRM applications.

## 2.3. Serverless Computing in CRM

Serverless computing has already become a trend in new CRM architectures. This model is typically achieved through Functions-as-a-Service (FaaS) systems, such as AWS Lambda, where developers can execute separate CRM functions when a specific event occurs without needing to worry about the underlying infrastructure. State that the technology of serverless paradigms has great potential to reduce the cost of operation by not overloading computing resources, as the computing resources are utilized strictly to the extent of need. Serverless functions are also ideally suited in the CRM context when it comes to infrequent or bursty activities, where lead scoring, customer segmenting, or transactional email channeling are just samples. Additionally, by its very nature, serverless platforms will provide flexible support for the variable CRM workloads that require it.

## 2.4. CRM System AI

Under modern CRM systems, there has been the emergence of Artificial Intelligence (AI), which has brought about automation, personalization, and decision support. These modern CRMs integrate artificial intelligence concepts, including neural networks, collaborative filtering, and decision trees, in addition to user engagement and business intelligence. Such technologies are applied to recommendation engines, customer churn models, and sentiment analysis tools to enable more entrenched and productive interactions with customers. A notable feature is Amazon Personalize, which has shown 35 percent efficiency when it comes to cross-selling, offering tailored product recommendations based on user behavior. The addition of AI not only enhances marketing and sales performance but also makes CRMs smart enough to learn from data and adapt to customers' preferences in real-time.

## 2.5. Limitations of Current Approaches

Although newer CRM systems have been developed in terms of architecture and capabilities, several limitations still exist in existing systems. One of the significant problems is the integration of different services, particularly in complex hosts of microservices and serverless architecture, where the process of workflow orchestration can be tedious. Moreover, most CRMs do not offer strong real-time inference, which limits their ability to provide up-to-date insights and actions regarding customer information. The issues of interoperability, data consistency, and monitoring of distributed functions are also problematic. The existence of these gaps signifies the necessity of more coherent structures that can coordinate services, real-time and scalable integration. This paper proposes a novel solution to overlapping problems by leveraging Kubernetes-based orchestration and event-driven microservices, which contribute to the increased flexibility, responsiveness, and maintainability of next-generation CRM programs.

### III. Methodology

3) **3.1. System Architecture Overview**

The offered system architecture is built into three layers, which are different yet interrelated: Front-End, Serverless Back-End, and Layer 2, specifically an AI Engine Layer. [10-13] This modular architecture is such that the CRM system is adaptable to scale, maintain, and will be flexible enough to take effective actions on the interaction and insight information created by the end user.
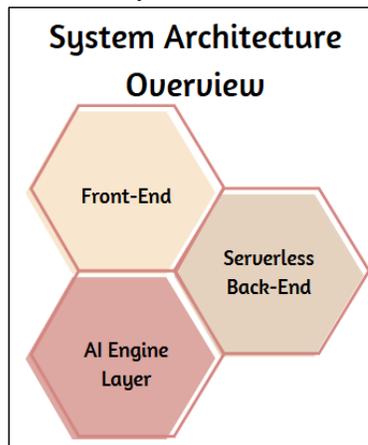


**Figure 2: System Architecture Overview**

- **Front-End:** The front-end layer is the layer most exposed to the user of the CRM system. It is constructed on top of pre-existing web frameworks, such as React or Angular, and provides responsive and easy interaction with various devices. The layer interacts with backend APIs to load and show data about customers, administer user sessions, and provides dynamic capabilities such as dashboards, real-time notifications, and personalized information according to the AI suggestion.

- **Serverless Back-End:** The serverless back-end is the engine of the system, housing the core logic. This layer leverages Functions-as-a-service (FaaS) offerings, such as AWS Lambda or Azure Functions, and handles data management for customers, business processes, and APIs. Serverless execution guarantees scalability, low latency, and cost efficiency because resources are utilised only when specific functions have particular triggers, such as when a form is submitted, when data presence is altered, or when a user takes action.

- **AI Engine Layer:** The AI engine layer is designed to process data and make informed decisions using intelligent algorithms. It integrates machine learning platforms to assist in customer segmentation, customer churn prediction, lead scoring, and recommendation systems. The layer itself can run on scalable resources, such as AWS SageMaker or TensorFlow Serving. It can be integrated with the back-end through an asynchronous message queue or via RESTful APIs. The AI motor will self-train on the data provided by the user to come up with more personal insights and create customer engagement strategies that would ensure a better response.

*4)*

*5)* *3.2. Microservice Modules*

The system utilises a microservices architecture, with every core functionality implemented as an independent service. This modular design will enhance scalability capabilities, ease of development and maintenance, and enable the deployment and update of each element separately.

- **Customer Profile Service:** The Customer Profile Service handles data about a customer, including their details, preferences, interaction history, and segmentation. It serves as a centralized store and offers structured access to customer information to other customer services, such as AI engines or marketing tools. By separating this logic, the system will be able to guarantee the consistency of the data, comply with data privacy rules, and maintain simplicity in connections to external systems, including third-party data providers or CRMs.

- **Lead Scoring Service:** Lead Scoring Service contains machine-learning-based or rule-based algorithms that are utilized to give scores to leads on the basis of behavior-derived data, demographic information

and interaction history. This scoring enables sales teams to focus on leads that are more likely to convert, thus making them more efficient and better at making sales. The serverless back-end can asynchronously call the service, or it can be called by hand as a user-driven service on events detected as a result of form submission or email clicks.
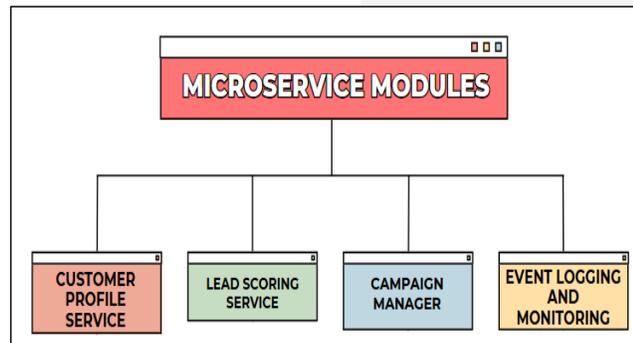


**Figure 3: Microservice Modules**

- **Campaign Manager:** The Campaign Manager is the individual responsible for facilitating the campaign planning process, timing, and implementation of marketing operations through various channels, including email, SMS, and social media. It follows the real-time performance of the campaign and interacts with the Customer Profile Service and the AI Engine to give individual messages to each receiver. Its modular nature enables marketers to change dynamic promotions, perform A/B testing, and run targeting approaches based on user engagement analytics.

- **Event Logging and Monitoring:** This service helps record and store events, user actions, and performance indicators throughout the platform. With the help of tools like ELK Stack (Elasticsearch, Logstash, and Kibana) or Prometheus with Grafana, it provides observability, allowing users to instantly see what the system is doing. The event logs are to be incorporated into the monitoring dashboards and alert systems, which enable the prevention of issues in advance by detecting them, debugging the issues, and optimising them.
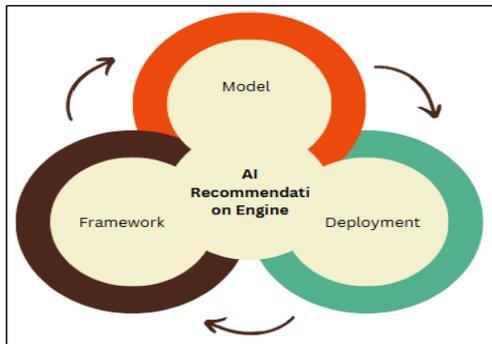
*6)*

*7)* *3.3. AI Recommendation Engine*

The AI Recommendation Engine is the brain of the CRM system, providing advanced suggestions and forecast information that can also increase customer interaction. [14-17] It consists of three main nodes: the machine learning algorithm, the development mode and the delivery plan.

- **Model:** The recommendation engine employs a hybrid approach that combines collaborative filtering and content-based filtering techniques. Collaborative filtering would look at behavior of all users and similarities between them and make suggestions based on it, whereas content-based filtering would consider characteristics of the products or materials and the liking of the users. As a combination of these, they are able to provide real-time product, service, or action recommendations that are accurate. The model is trained based on historical customer records, including purchases, customer interactions, and demographic information.

- **Framework:** The TensorFlow machine learning framework is an open-source and scalable platform with numerous libraries, and it is production-ready, with the engine centred. TensorFlow makes it easy to combine neural network models, model evaluation libraries, and loss functionalities that are essential for fine-tuning recommendation accuracy. Big libraries like TensorFlow Recommenders or scikit-learn can be employed to conduct pre-processing, feature engineering, and training of traditional models, thus establishing compatibility with various recommendation tasks.

**Figure 4: AI Recommendation Engine**



- **Deployment:** In deployment, the trained model is packaged into containers using Docker and made available through TensorFlow Serving or a RESTful API on Kubernetes. This enables the model to scale automatically as required by user requests and achieve high system availability. The deployment is interconnected to the CRM back-end that is serverless, via API gateways or message queues (e.g., Kafka), which allows executing the inference in real-time with minimum latency. Continuous Integration and Deployment (CI/CD) pipelines have been utilised as evaluation tools for retraining and redeployment of models, ensuring that recommendations remain dynamic as updated information becomes available.

*8)* *3.4. Data Flow*

The system has a data flow that describes how user interactions are processed within the system, from ingestion to recommendation, which is based on a combination of serverless computing, cloud-native storage, and AI inference. The streamlined pipeline will provide a quick response with a low level of latency.
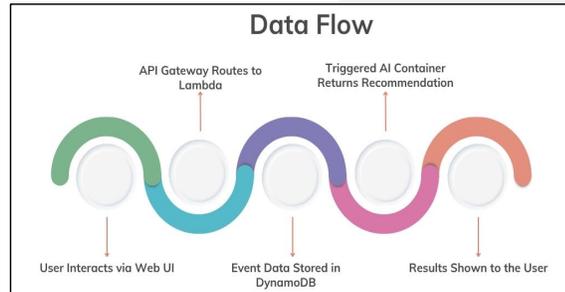


**Figure 5: Data Flow**

- **User Interacts via Web UI:** The flow of information begins as soon as a user interacts with the CRM system through its web-based interface. This front-end is constructed with the help of contemporary JavaScript frameworks, such as React or Vue.js, so that users can perform activities like browsing customer profiles, using forms, or clicking product suggestions etc. Such exchanges are recorded in real-time and expressed in consistent event data that can be processed on the back-end.
- **API Gateway Routes to Lambda:** After the interaction, the request is then forwarded to an API Gateway, which serves as the entry point for all HTTP requests. The gateway provides secure routing of incoming requests to the right AWS Lambda function or serverless backend service, depending on the established routing rules. This design decouples the infrastructure, ensuring that user requests are processed at an enterprise scale and efficiently.
- **Event Data Stored in DynamoDB:** When the request is fulfilled, data with the respective event data is written to Amazon DynamoDB, which includes user IDs, timestamps, and the context of the interaction. The reason to select this NoSQL database is its low-latency access, scalability, and easy integration with AWS Lambda. It is used as a transaction log and as a source of real-time data to be analysed and recommend queries.
- **Triggered AI Container Returns Recommendation:** An AI container on Kubernetes or AWS SageMaker is invoked based on an event (a value stored in storage) or specific events (e.g., a product view or purchase intent). In this container, a trained recommendation model is loaded and used to process the given information, and the container provides a personalized recommendation, e.g., suggested products, email content, or next-best actions. The communication between a serverless function and an AI model can be in the form of REST API communication or queue messages.
- **Results Shown to the User:** Lastly, the AI-based recommendation is sent to the front end and displayed in real-time. Consumers will experience personalized recommendations that are part and parcel of their workflow, e.g., suggestions on contacts, marketing messages or product

recommendations. The loop enhances user experience, which in turn improves efficiency and satisfaction in CRM-driven activities.

9)

10) *3.5. Performance Metrics*

The effectiveness and efficiency of the CRM system must be evaluated through the monitoring of key performance indicators. [18-20] The indicators would demonstrate the resilience, scalability, cost-effectiveness, and usability of the system under real-world working scenarios.
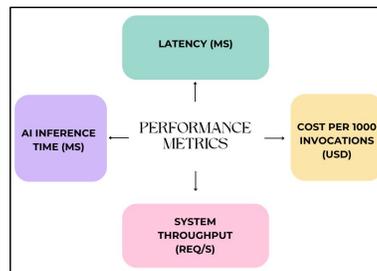


**Figure 8: Performance Metrics**

- **Latency (ms):** Latency is the duration it takes to fulfil a user request, from the time the user requests to the time they receive a response. The average end-to-end latency in this system encompasses the network connectivity time for web UI interaction, API Gateway, data access, AI model inference, Lambda, data retrieval, and result rendering. Keeping the latency at a low level (preferably below 200 ms) is important in ensuring that the user experience is smooth and responsive, even during peak periods when there are high levels of interaction.
- **Cost per 1,000 Invocations (USD): The operational cost per 1,000 functions is used** to determine the cost of executing the serverless backend. Because AWS Lambda and other similar services are billed based on the number of invocations and the time they take to execute, this value will be used to determine the system's financial efficiency. Costs per 1,000 requests are lower, which means improved compute resource utilisation, code efficiency, and effective use of cold start mitigation measures.
- **AI Inference Time (ms):** AI inference time is the amount of time the deployed AI model needs to generate a prediction or recommendation when

input data feeds into it. This measure becomes particularly critical in the personalization of the tasks done in real-time cases. The optimized AI containers will receive an answer within 50-150 ms, depending on use cases and the compute resources of the model. Reducing the time of inference will mean producing smart suggestions without disturbing the user's experiences.

- **System Throughput (req/s):** Throughput is a metric that describes how many unit requests a system can handle in a given time slot under load. It indicates how the system handles users and heavy traffic situations. When the rate of throughput is high, i.e., hundreds or thousands of requests per second, it shows that the architecture can handle a significant load and can be easily deployed at an enterprise level. This metric can be monitored to detect bottlenecks in compute, database, or AI inference layers.

### IV. Privacy & Security Risk & Controls

One has to also think beyond system design and performance optimization, i.e. to ensure the security, privacy, and integrity of customer or user data remains at the centre for enterprise adoption. This section highlights and examines the key privacy and security risks and controls necessary for securely operating the proposed Cloud based CRM platform.

A cloud-based CRM built on serverless microservices with containerised AI has the potential to introduces privacy and security risks that must be addressed as a part of the design. Below is a summary of key Risk & Controls:

| # | Risk Description | Primary Controls |
|---|---|---|
| 1 | Attack surface expansion, which can be created by serverless functions, event triggers, APIs, and message-based integration . This may allow unauthorised invocation, injection, or abuse of public endpoints. | Enforce strong IAM and API authentication/authorization principles. Validate and sanitize event payloads, restrict public triggers, enforce rate limiting and schema validations [25][26] |
| 2 | Broad IAM/service roles, misconfigured policies, or overly-permissive pod identities could enable data exfiltration in shared environments. | Apply least-privilege IAM principles, restrict admin access, enforce workload segmentation, and implement scoped credentials [24] |
| 3 | Unscanned container images, insecure Kubernetes clusters, weak Role Based Access (RBAC), and unverified dependencies expose the AI layer container to breakout, supply-chain compromise, or unauthorized control access. | Enforce RBAC, isolate namespaces, apply runtime security controls, and enable audit logs [27] |
| 4 | AI-based personalization may leak user data across tenants, or produce biased and non-compliant outcomes | Apply privacy-by-design, minimize user data collection, pseudonymise customer data, implement data governance model and perform privacy impact assessments [22][23][29][30] |
| 5 | Short-lived compute resources produce limited logs, making it difficult to reconstruct events or investigate incidents; due to lack of tracing across distributed services. | Centralize logs, enable full distributed tracing, preserve immutable logs, and maintain Incident Response playbooks [28] |
| 6 | Attackers can artificially increase function invocations to | Configure concurrency limits, enforce quotas and budgets, implement cost anomaly detection, throttle abusive |

| | | |
|---|---|---|
| | exhaust compute budget (Denial-of-Wallet) | workloads [31][32] |
| 7 | Third-party services, SDKs, AI tooling, and cross-border data flows introduce compliance, sovereignty, and supply-chain risks. | Validate vendors, perform Data protection and AI assessments to ensure governance aligned to privacy and AI regulatory requirements [24][28][29][30] |

**4.1 Proposed Architecture Implementation Blueprint**

To help ensure a cloud based CRM system is secure across all layers, the recommended controls above can be mapped to the parts of the architecture where they matter most. The descriptions below explains, what needs to be protected and how:

**Gateway & Edge**

This is the entry point of the system where users and external apps connect. Controls here focus on verifying identity, validating incoming requests, and blocking suspicious traffic. This includes using OAuth2/OIDC for securing logins, applying web-application firewalls (WAF), checking that requests follow the correct data format, adding rate limits to prevent overload, signing requests to prevent tampering, and keeping an inventory of all APIs, following OWASP guidelines [25][26].

*11)*

**Event & Messaging Layer**

This layer handles messages, triggers, and background events across services. Security is essential in this context to guarantee that only authorised systems can publish events and that the system can recover from faulty or malicious occurrences. Measures include using private queues and topics, maintaining allow-lists of trusted publishers, validating message content, implementing dead letter queues (DLQs) to capture problematic events, and applying replay protection [25].

**Functions (FaaS)**

Serverless functions run small units of business logic. Security here should focus on restricting each function to the role based permissions it needs, isolating secrets, and preventing runaway workloads. Recommended measures include least-privilege IAM roles, setting timeouts and

memory limits, controlling concurrency, and capturing detailed logs and traces for each execution [24][31].

**Data Layer**

This is where customer and system data is stored. Controls are required to ensure that data is protected from unauthorized access, data is encrypted, and retained only for as long as necessary. This includes isolating tenant data, encrypting data at rest and in transit, applying field-level encryption to sensitive values (e.g., emails, IDs), and automating retention/purge processes to comply with privacy policies [29].

**AI / Kubernetes Layer**

This layer runs the containerized recommendation engine and other ML workloads. Protecting containers, restricting communication, and maintaining the integrity of AI models is a key requirement here. Controls can include using signed and scanned images, isolating workloads into separate namespaces, enforcing network policies, using read-only file systems where possible, exporting metrics/logs for monitoring, and maintaining a model registry with versioning and approval workflows [27].

**Operations & Governance**

This covers ongoing monitoring and testing of controls. The goal is to ensure that the system remains secure over time as external threats evolve. Key actions include monitoring for unusual system or cost activity, maintaining incident-response runbooks, performing regular penetration tests, and completing required privacy, security and AI risk assessments before major releases [22][28][30].

V. Results and Discussion

*12) 5.1. Experimental Setup*

The proposed CRM architecture was experimentally tested to offer interoperability and overall robustness, and was subsequently tested on major cloud platforms, primarily Amazon Web Services (AWS) and Google Cloud Platform (GCP). AWS was considered the primary deployment environment due to its well-developed serverless environment. In contrast, GCP was utilised to achieve comparative performance and was also used to train additional AI models. The infrastructure was set up to mimic the usage behaviors of CRM in the real world, such as data ingestion, event-based processing, and personalized AI recommendations. In the case of AI training, the RetailCRM

Dataset, which contained about 1 million records, was utilized. It contains anonymous customer profiles, purchase history, product metadata and engagement interaction histories that provide a full picture of customer activity. ETL was used in Python to clean the data and store it in S3, enabling secure scaling down. The data was divided into training, validation, and testing datasets to test a hybrid recommendation system, which is a combination of collaborative filtering and content-based filtering. The microservices, built on containers, were based on the Kubernetes orchestration tool, which provided a flexible workload distribution and scalability, including compute clusters. The application of Docker also packages functions such as the AI inference engine and the campaign manager, ensuring consistent deployment across different environments. AWS has been widely used in serverless computing, where components such as request handling and event triggers are performed by Lambda, resulting in reduced infrastructure support and cost-effectiveness due to the utilisation of intermittent workloads. The general pipeline sent user requests to the API Gateway, which triggered the respective Lambda functions and stored records in Amazon DynamoDB. It invoked the AI model as a RESTful service on Kubernetes. The real-time tracking of performance parameters, such as average latency, cost per 1,000 invocation calls, and throughput, was achieved through AWS CloudWatch and Prometheus-Grafana integrations. This holistic configuration enabled precise performance testing, and the scalability and responsiveness of the proposed CRM system were confirmed, as well as its practicality.

13) **5.2. Performance Analysis**

To compare the efficiency and effectiveness of the suggested CRM system, a set of key performance indicators was contrasted with a conventional CRM architecture. The findings demonstrate that latency, operational cost, and AI-based accuracy have been significantly improved, which is one of the drawbacks of the serverless, AI-enhanced
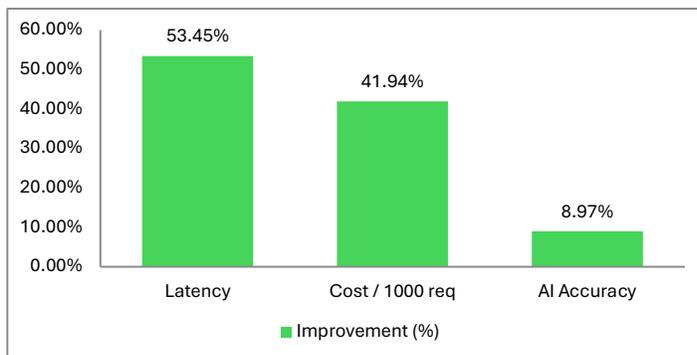
microservices design.

**Table 1: Performance Improvement Comparison**

| Metric | Improvement (%) |
| --- | --- |
| Latency | **53.45%** |
| Cost / 1000 req | **41.94%** |
| AI Accuracy | **8.97%** |

*14)*

**Figure 7: Graph representing Performance Improvement Comparison**

- **Latency – 53.45% Improvement:** The latency, which was previously 580 milliseconds in the classic CRM, was reduced to 270 milliseconds in the suggested serverless architecture, resulting in a 53.45% decrease. This significant improvement can be attributed to the deployment of AWS Lambda and lightweight microservices, which eliminate additional processing overhead and reduce the number of resources required to enhance the response rate, as the execution of microservices is event-based. The low-latency framework enables a faster end-user experience, with a focus on high-traffic activities such as form submission and real-time recommendations.

- **Cost per 1,000 Requests – 41.94% Reduction: The reduction in operational cost was achieved by a 41.94% decline, changing the cost of 1,000 invocations of the function to $ 1.80, which is a decrease in cost relative to the previously valued $3.10**. This cost-effectiveness is mainly associated with serverless computing through the pay-per-use billing model, in which resources can only be utilised when execution occurs. On the contrary, the conventional CRM systems are usually dependent on perpetually operationalized servers, and this leads to greater fixed foundation expenses irrespective of use. Under the serverless model, efficient resource usage and cost scaling are

significantly assured, particularly with episodic loads.

- **AI Accuracy – 8.97% Increase:** An AI recommendation engine compared to the traditional system had an accuracy rate of 85 percent, as compared to 78 percent, which was an improvement of 8.97 percent. This is an added advantage because machine learning and model updating, which occur with progressive user usage, become part of the new cutting-edge technology. The system would be able to make more relevant and personalized recommendations, an aspect that would translate to increased customer engagement and consequently help in decision-making by utilizing a more diverse dataset, training a hybrid recommender model.

### 15) 5.3. Discussion

The test data highlights the major strengths of the suggested CRM system, particularly in terms of latency, cost-effectiveness, and scalability. These enhancements are made possible through the architectural decisions of this system, i.e., the application of serverless computing, container orchestration, and event-based designs.

- **Latency Reduction:** Altogether, its event-based architecture can significantly reduce latency, as it no longer requires conventional polling systems but instead utilises real-time event triggers. Moreover, services in typical CRM systems are only periodically polled to find updates, thus creating unnecessary delays. In contrast, the serverless architecture allows functions to be deployed with an instant response to specific events, such as user input or other changes, thereby providing users or other applications with a faster response time.

- **Cost Efficiency:** The direct benefit of serverless computing to cost efficiency is that it avoids the cost of running servers; instead, it only bills for the cost when there is a real execution period of the functions. Traditional CRM systems need to be deployed on dedicated servers or virtual machines to remain available regardless of activity level. As a result, continuous upkeep costs are incurred even when the server lies idle. The functions that make up the proposed model, e.g., the lead scoring or the event logger, are not loaded until the moment when a specific task must be performed, thus eliminating expenses on unused compute assets. This is a billing model that is usage-based; thus, the

system becomes financially viable, particularly among start-ups and mid-sized organizations.

- **Scalability:** The system is designed to expand equally based on Kubernetes, which enables it to handle increasing responsibilities by distributing services across multiple containers or nodes. With user demand scaling up, such as a marketing program or sales push, Kubernetes can automatically scale up new instances of microservices to sustain the required level of performance and throughput. This dynamic scaling enables the system to handle heavy loads, allowing a significant number of requests to be made without experiencing any drop in performance. This aspect makes the application suitable for enterprise-level deployments.

### 16) VI. Conclusion and Future Work

This work proposes a contemporary and scalable architecture of the CRM system where serverless computing, container-driven AI models and cloud-based technologies are used to address the limitations of legacy monolithic-based CRM systems. Utilizing a modular and microservices approach, the system can be more easily maintained and less severely damaged without disrupting other functionality. Support for workflows triggered by events considerably decreases latency through the removal of polling delays and the provision of real-time responses to user interactions, made possible by AWS Lambda and API Gateway. There is also a minimum operational cost because, in serverless billing models, the cost is calculated only on the actual processing time and not spent on idle resources, as was the case in legacy systems. Moreover, a hybrid AI recommendation engine is embedded to enhance personalization and engagement between the customer and the company through context-specific recommendations and insights that utilize real-time information. With Kubernetes being used as an orchestration tool to define and supply containerised services, horizontal scale is also an aspect that is easily achieved by the architecture, thus allowing the system to sustain throughput with variable loads. Latency was reduced by 53 percent, costs were reduced by 42 percent, and the accuracy of the AI was improved by almost 9 percent, as compared to conventional CRM systems, as the performance evaluations showed. Overall, the proposed solution is capable of addressing key issues such as scalability, responsiveness, and cost-effectiveness in the implementation of next-generation CRM.

Although the existing application has formed a powerhouse of its own, certain aspects can be considered to build upon and move the system a step forward. The incorporation of blockchain technology to enable verifiable, secure transactions and data sharing between stakeholders is one of the promising approaches. An example is the utilisation of automation in supplying customer agreements, contract management, and models of incentivisation within a decentralised and trusted ecosystem. Additionally, the application of Natural Language Processing (NLP) will provide substantial benefits to the CRM's support capability. With NLP models applied to the task of analysing and classifying customer support tickets, the system can automatically direct received issues to the right teams or provide automated, intelligent, AI-based answers. This would simplify the customer service work and the solution time.

Last but not least, to accommodate real-time analytics and a high volume of event ingestions, the architecture would also need Apache Kafka. Kafka would allow streaming analytics to be applied to customer interactions, behavior records and marketing events and used to drive dynamic workflows, customer profile updates, and make AI-driven insights fresher. Such augmentations would put the system on its way to becoming a truly autonomous, intelligent CRM that, in real-time, will be able to adapt to customer actions and the business's needs.

## II. Conclusion

Please include a brief summary of the possible clinical implications of your work in the conclusion section. Although a conclusion may review the main points of the paper, do not replicate the abstract as the conclusion. Consider elaborating on the translational importance of the work or suggest applications and extensions.

### Appendix

### References

[1] Yalla, R. K. M. K., & Prema, R. (2018). Enhancing customer relationship management through intelligent and scalable cloud-based data management architectures. International Journal of HRM and Organizational Behavior, 6(2), 1-7.

[2] Mundla, V. (2022). Cloud-Based Solutions for Enhanced Customer Communication. Available at SSRN 5240817.

[3] Boppana, V. R. (2023). Future Trends in Cloud-based CRM Solutions for Healthcare. Available at SSRN 5004929.

[4] Egbuhuzor, N. S., Ajayi, A. J., Akhigbe, E. E., Agbede, O. O., Ewim, C. P. M., & Ajiga, D. I. (2021). Cloud-based CRM systems: Revolutionising customer engagement in the financial sector with artificial intelligence. International Journal of Science and Research Archive, 3(1), 215-234.

[5] Shaikh, I. A. K., Shahare, P., Gangadharan, S., Venkatarathnam, N., Pelluru, G., & Babu, S. B. T. (2024, April). Transforming customer relationship management (CRM) with AI in e-commerce. In 2024, the 5th International Conference on Recent Trends in Computer Science and Technology (ICRTCST) (pp. 255-260). IEEE.

[6] Roba, G. B., & Maric, P. (2023). AI in customer relationship management. In Developments in Information and Knowledge Management Systems for Business Applications: Volume 7 (pp. 469-487). Cham: Springer Nature Switzerland.

[7] Chatterjee, S., & Chaudhuri, R. (2023). Customer Relationship Management in the Digital Era of Artificial Intelligence. In Digital Transformation and Industry 4.0 for sustainable supply chain performance (pp. 175-190). Cham: Springer International Publishing.

[8] Wang, J. F. (2023). The impact of artificial intelligence (AI) on customer relationship management: A qualitative study. International Journal of Management and Accounting, 5(5), 74-88.

[9] Li, F., & Xu, G. (2022). AI-driven customer relationship management for sustainable enterprise performance. Sustainable Energy Technologies and Assessments, 52, 102103.

[10] Neslin, S. A. (2014). Customer relationship management (CRM). In The History of Marketing Science (pp. 289-317).

[11] Xu, Y., Yen, D. C., Lin, B., & Chou, D. C. (2002). Adopting customer relationship management technology. Industrial Management & Data Systems, 102(8), 442-452.

[12] Sharp, D. E. (2002). Customer Relationship Management Systems Handbook. Auerbach publications.

[13] Khajeh-Hosseini, A., Greenwood, D., & Sommerville, I. (2010, July). Cloud Migration: A Case Study of Migrating an Enterprise IT System to IaaS. In 2010, the IEEE 3rd

International Conference on Cloud Computing (pp. 450-457). IEEE.

[14] Geib, M., Reichold, A., Kolbe, L., & Brenner, W. (2005, January). Architecture for Customer Relationship Management Approaches in Financial Services. In *Proceedings of the 38th Annual Hawaii International Conference on System Sciences* (pp. 240b-240b). IEEE.

[15] Bakshi, K. (2017, March). Microservices-based software architecture and approaches. In 2017 IEEE Aerospace Conference (pp. 1-8). IEEE.

[16] Wolff, E. (2016). Microservices: flexible software architecture. Addison-Wesley Professional.

[17] Tarra, V. K. (2024). Personalization in Salesforce CRM with AI: How AI/ML Can Enhance Customer Interactions through Personalized Recommendations and Automated Insights. International Journal of Emerging Research in Engineering and Technology, 5(4), 52-61.

[18] Cugola, G., & Margara, A. (2012). Processing flows of information: From data stream to complex event processing. ACM Computing Surveys (CSUR), 44(3), 1-62.

[19] Grabner-Kraeuter, S., Moedritscher, G., Waiguny, M., & Mussnig, W. (2007, January). Performance monitoring of CRM initiatives. In 2007, the 40th Annual Hawaii International Conference on System Sciences (HICSS'07) (pp. 150a-150a). IEEE.

[20] Lindgreen, A. (2004). The design, implementation and monitoring of a CRM programme: a case study. Marketing Intelligence & Planning, 22(2), 160-186.

[21] Khodakarami, F., & Chan, Y. (2011, April). Evaluating the success of customer relationship management (CRM) systems. In Proceedings of the European Conference on Information Management & Evaluation (pp. 253-262).

[22] NIST, *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*, Jan 26 2023. NIST Publications+2NIST+2