# Billing System in Distributed Computing using Apache Kafka

Mitali Admuthe, Sonal Bijitkar, Pooja Chavan and Joyli Rumao
California State University, Long Beach
Instructor: Prof. Pooria Yaghini

## ABSTRACT

This paper presents a distributed architecture for real-time resource usage monitoring and billing across computing environments. Utilizing Apache Kafka for task distribution and data management, the system separates resource monitoring (producers) from billing computations (consumers), allowing efficient, reliable processing across multiple machines. Kafka's role in high-throughput data streaming and task distribution highlights its potential in distributed computing, achieving scalability and concurrent processing.

**Keywords:** Distributed Computing, Apache Kafka, Resource Monitoring, Billing System, Task Distribution, Real-Time Processing.

## 1. INTRODUCTION

Real-time resource monitoring and billing systems are crucial in large-scale distributed environments where optimizing resource usage directly impacts operational costs and system performance. With the advent of cloud-native architectures and microservices, systems like Apache Kafka have become integral for decoupling data generation and processing layers, offering resilience and scalability. This research highlights the potential to implement such systems in a wide range of scenarios, from cloud-based infrastructures to edge computing environments. In modern distributed computing, real-time monitoring and billing based on CPU and memory usage are critical. This project demonstrates an architecture leveraging Apache Kafka to manage data ingestion, processing, and task distribution. By decoupling monitoring and billing, the architecture achieves scalability and high-throughput handling across distributed systems.

## 2. LITERATURE REVIEW

Billing systems in distributed computing play a crucial role in monitoring and charging for resource usage across multiple nodes in a distributed environment. This section summarizes key aspects of the literature, exploring the evolution, techniques, and challenges of billing systems in distributed systems.

### 2.1 Evolution of Billing Systems in Distributed Computing

Billing in distributed systems evolved alongside the growth of cloud computing and distributed architecture platforms. Early approaches focused on simplistic, coarse-grained models that billed customers based on static quotas or fixed rates. However, with the advent of virtualization and scalable architectures, more dynamic, fine-grained models became necessary to accurately reflect resource consumption. Systems such as Amazon Web Services (AWS) and Microsoft Azure pioneered utility-based billing models, introducing the concept of *pay-as-you-go* billing.

**Key advancements:**

- **Cloud Service Models:** IaaS (Infrastructure as a Service), PaaS (Platform as a Service), and SaaS (Software as a Service) introduced differentiated billing based on services consumed.

- **Resource-Based Billing:** Dynamic pricing models based on CPU, memory, storage, and network usage were introduced to improve cost fairness.

### 2.2 Technical Approaches to Distributed Billing

The billing system in distributed systems primarily focuses on resource monitoring, data aggregation, and cost computation. The technical approaches involve:

#### 2.2.1 Resource Monitoring

Tools like `psutil` (as used in the `machine.py` script in your implementation) monitor system-level metrics such as CPU and memory utilization. Distributed systems often deploy agents across nodes to collect resource usage in real-time.

#### 2.2.2 Data Aggregation

Middleware systems like Apache Kafka (used in your implementation) are commonly employed for collecting and streaming usage data to a central server or database for billing computation. This ensures scalability and fault tolerance.

#### 2.2.3 Cost Calculation

Cost is typically computed by applying predefined rates to resource usage metrics (e.g., CPU rate at $0.05 per percentage, as in your implementation). Some systems support dynamic pricing based on demand and availability (e.g., spot pricing in AWS).

#### 2.2.4 Database Storage

SQLite, as used in your implementation, is a lightweight option suitable for small-scale systems. Larger systems use

distributed databases like Cassandra or MongoDB to handle massive billing data.

## 2.3 Challenges in Distributed Billing Systems

While billing systems have become more sophisticated, several challenges remain:

- **Scalability:** Distributed systems must handle large-scale data from potentially thousands of nodes. Kafka and similar technologies address this, but ensuring consistent performance at scale remains complex.

- **Accuracy and Consistency:** Synchronization issues can arise when aggregating data from multiple nodes. Missing or delayed data can lead to billing errors, impacting customer trust.

- **Fairness in Cost Calculation:** Different workloads may exhibit varying resource usage patterns, necessitating adaptive pricing models that account for bursty traffic or idle resources.

- **Security and Privacy:** Billing data must be securely transmitted and stored to prevent unauthorized access or fraud.

## 2.4 Case Studies and Real-World Systems

**Amazon Web Services (AWS):** AWS introduced fine-grained resource tracking with tools like CloudWatch to monitor resource usage and ensure accurate billing. Spot instances and reserved instances reflect advanced pricing models.

**Google Cloud Platform (GCP):** GCP integrates real-time billing APIs for granular usage tracking, enabling customers to control costs effectively.

**Our Implementation:** The system we have developed combines:

- Apache Kafka for real-time usage data streaming.

- SQLite for lightweight billing data storage.

- A clear cost calculation formula, providing transparency and simplicity.

## 3. SYSTEM ARCHITECTURE AND DESIGN

### 3.1 Overview

The system consists of three core components:

- **Producers (Machine Clients):** Each client collects CPU and memory usage data, sending it to Kafka.

- **Kafka Broker:** Acts as an intermediary, managing data flow between producers and consumers.

- **Consumers (Workers):** Each worker instance subscribes to Kafka topics to compute billing from incoming data.
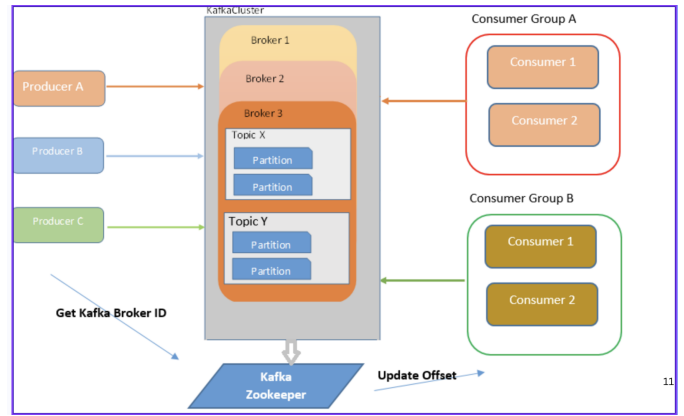


**Figure 1: System Architecture.**

This design emphasizes Kafka's ability to handle concurrent data flows from multiple producers, distributing processing load across consumers. Kafka's scalability enables easy integration of new producers or consumers without disrupting system performance.

### 3.2 Kafka's Role in Data Pipeline Management

Kafka is used to:

- **Ingest Data:** Producers publish resource usage data to Kafka topics.

- **Distribute Processing:** Multiple consumer nodes independently subscribe to topics, with Kafka balancing load and handling faults.

- **Partitioned Distributed Log Model:** Kafka's fault tolerance ensures reliability, even during broker or network failures.

### 3.3 Role of Zookeeper in System Coordination

The Apache Zookeeper service is a cornerstone in Kafka's architecture, managing distributed configuration and synchronization. It ensures seamless cluster coordination by:

- Monitoring Kafka brokers and updating their status in real time.

- Managing topic metadata and partition assignments.

- Enabling fault tolerance by orchestrating leader election for Kafka partitions.

## 4. IMPLEMENTATION DETAILS

### 4.1 Producer (Machine Client) Configuration

The producer scripts collect and send CPU and memory usage data to Kafka. Key functionalities include:

- **Resource Monitoring:** Libraries such as psutil capture CPU and memory usage.

- **Data Serialization and Transmission:** Data is serialized into JSON and published to Kafka under the `resource_usage` topic.

- **User-Controlled Transmission:** The *User-Controlled Transmission* feature, facilitated through a Tkinter-based graphical user interface (GUI), empowers users to take control over the timing and frequency of data transmission. Unlike traditional systems where data is streamed continuously, this approach allows users to decide when the transmission should occur. By granting this control, the system significantly reduces unnecessary network traffic, particularly in scenarios where constant data updates are not required. For instance, in distributed systems monitoring applications, a user might prefer to transmit data at specific intervals or only when critical thresholds are exceeded, rather than streaming every minor change in resource usage.

  The GUI not only enhances user experience by providing an intuitive interface for managing transmissions but also contributes to system efficiency. Users can pause, resume, or trigger data transmission on-demand, allowing the system to better align with operational needs. This flexibility makes the architecture adaptable to a variety of real-world applications, such as environments with limited bandwidth, intermittent connectivity, or strict cost constraints. Moreover, by reducing continuous network load, this feature can help in lowering operational costs and energy consumption, making it a sustainable and user-centric solution.



Figure 2: machine.py

## 4.2 Kafka Topic Configuration

A single Kafka topic (`resource_usage`) is dedicated to receiving resource usage data. Both producers and consumers connect to this topic using a predefined server IP and port. Kafka's topic management ensures efficient data distribution to consumers.

## 4.3 Consumer (Worker) Configuration

The consumer scripts subscribe to `resource_usage` and:

- **Data Storage:** Store the usage data in an SQLite database for persistence.

- **Billing Calculation:** Compute costs based on CPU and memory usage with pre-defined rates.We have used the following rates:
  **CPU rate** = 0.05$ per CPU percentage
  **Memory rate** = 0.01$ per memory percentage

- **REST API:** A REST API, built using Flask, exposes the latest usage and billing data for dashboard integration, enabling real-time updates.



Figure 3: worker.py

## 5. IMPLEMENTATION COMMANDS

The following commands were utilized for setting up and running the distributed billing system:

### 5.1 Start Zookeeper

Zookeeper serves as the coordination service for distributed systems. To initiate the Zookeeper server, execute the following command from the Kafka directory:

**Listing 1: Start Zookeeper**

```
java -cp "libs/*;config" org.apache.
    zookeeper.server.ZooKeeperServerMain
    config\zookeeper.properties
```

This command launches the Zookeeper server, which is essential for managing and coordinating Kafka brokers.

### 5.2 Start Kafka Server

To start the Kafka broker, which manages message streams and ensures fault tolerance and reliability in the data pipeline, run the command below:

**Listing 2: Start Kafka Server**

```
java -cp "libs/*;config" kafka.Kafka
    config\server.properties
```

3

## 5.3 Create Kafka Topic

The `resource_usage` topic is created to facilitate data transmission between producers and consumers. Use the following command to create the topic:

**Listing 3: Create Kafka Topic**

```
bin/kafka-topics.sh --create --topic
    resource_usage --bootstrap-server
    localhost:9092 --partitions 3 --
    replication-factor 1
```

`-partitions 3`: Specifies the number of partitions for the topic.
`-replication-factor 1`: Indicates the replication factor for topic data across brokers.

These commands establish the core infrastructure for the distributed billing system, enabling seamless data flow between producers and consumers.

## 6. DISTRIBUTED COMPUTING ASPECTS

### 6.1 Task Distribution Using Kafka

Kafka's consumer group and partitioning models enable efficient task distribution. Each worker in the consumer group is assigned a partition, ensuring parallel processing and load balancing.

### 6.2 Scalability and Fault Tolerance

Adding new producers and consumers is straightforward, as Kafka automatically balances workloads. Kafka's replication model provides reliability, making data available even during broker failures. Additionally, each worker's SQLite database preserves data, enhancing system resilience.

### 6.3 Performance Optimization

The system is designed for high throughput, with Kafka's architecture minimizing I/O latency and allowing real-time processing. Processing data in batches further reduces latency, ensuring quick and reliable billing calculations even with increased load.

## 7. TESTING AND RESULTS

### 7.1 Experimental Setup

The setup involved four producer machines running resource monitoring processes and four consumer machines within a local network. Kafka brokers and Zookeeper managed task assignment and data flow.

### 7.2 Results

The system showed effective real-time processing:

- **Latency:** Average latency from data generation to billing was approximately 1.5 seconds.

- **Scalability:** Consistent performance was observed with up to 50 producers.

- **Reliability:** No data loss occurred, and data consistency was maintained under simulated network interruptions.
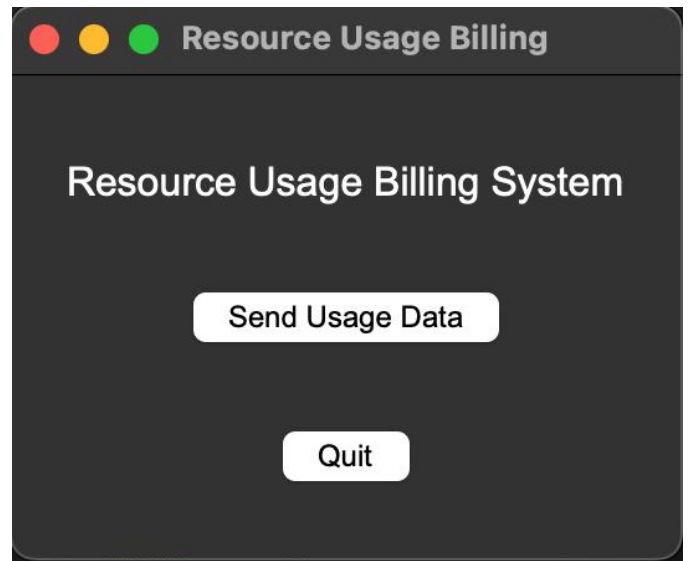


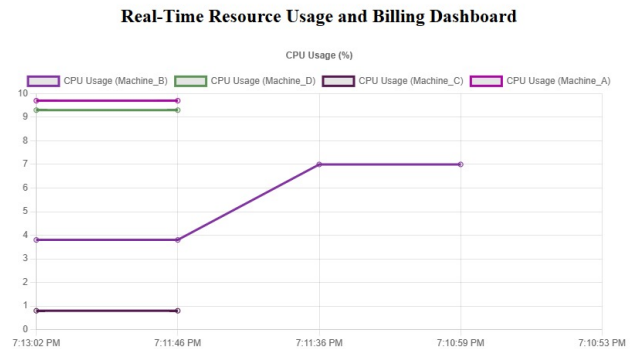**Figure 4: User interface of the Resource Usage Billing System.**



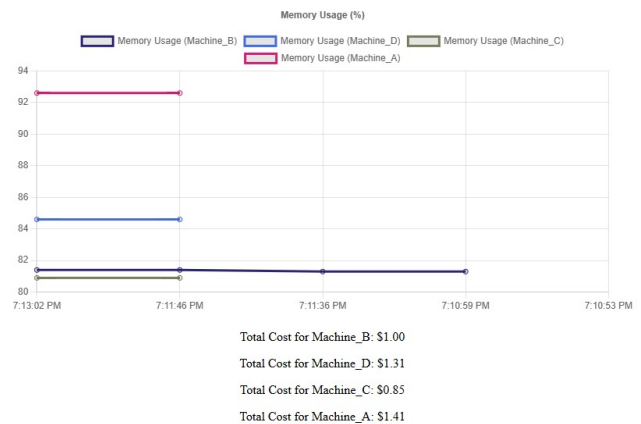**Figure 5: Visualization of real-time CPU usage for multiple machines.**



**Figure 6: Total cost calculated for real-time CPU usage of multiple machines.**

## 8. DISCUSSION

The Kafka-based architecture developed in this project highlights a modern approach to distributed resource monitoring and billing, leveraging Kafka's distributed log model to achieve high scalability, reliability, and efficiency. One of the key takeaways from this implementation is the ability to manage real-time resource usage data across distributed nodes while ensuring fault tolerance and consistent data processing. By decoupling the data ingestion layer from storage and computation, the system allows for seamless integration of new nodes and scales horizontally without significant architectural changes. Despite these advantages, configuring and maintaining Kafka in a production environment proved challenging, particularly in balancing network loads and ensuring optimal partitioning across brokers. These challenges underscore the need for automated tools to dynamically monitor and adjust configurations in real-time, especially as the system scales.

An important observation is the potential for this architecture to evolve into more intelligent billing systems. While the current implementation uses straightforward cost computation formulas, integrating predictive analytics through machine learning models could significantly enhance the system's capabilities. Such models could analyze historical trends and forecast usage patterns, enabling customers to anticipate costs and adjust their resource consumption proactively. Furthermore, the inclusion of dynamic pricing models that adapt to demand fluctuations could provide a more fair and competitive billing structure.

Additionally, the project demonstrates opportunities to leverage cloud-based Kafka solutions, which offer managed services, automatic scaling, and reduced operational complexity. Transitioning to such platforms would not only simplify the management of Kafka clusters but also enhance the system's ability to handle larger datasets and higher throughput. Another area for enhancement is data security and privacy, which are critical in distributed systems, especially when handling sensitive billing information. Implementing end-to-end encryption, secure communication protocols, and robust access control mechanisms would strengthen the system's reliability and ensure compliance with data protection regulations.

Lastly, this architecture opens up possibilities for hybrid or federated billing models that integrate multiple distributed systems, creating a unified platform for billing across heterogeneous environments. This would cater to organizations with multi-cloud or hybrid-cloud setups, offering a consolidated view of resource usage and costs. While this project establishes a strong foundation, it also paves the way for future innovations that could transform how billing is managed in distributed computing ecosystems.

## 9. CONCLUSION

This project illustrates Apache Kafka's potential for real-time distributed monitoring and billing. By separating data production from billing calculations, Kafka enables high-throughput processing, fault tolerance, and scalability, meeting the requirements of complex, distributed systems.This project illustrates Apache Kafka's potential for real-time distributed monitoring and billing in modern computing environments. By separating data production from billing cal-

culations, Kafka enables high-throughput processing, fault tolerance, and scalability, effectively meeting the demands of complex, distributed systems. The design also highlights Kafka's strength in managing concurrent data streams from multiple producers, facilitating parallel processing without performance degradation. This architecture ensures that the system remains resilient and adaptable, even during unexpected failures or network interruptions.

Additionally, the integration of Kafka with distributed computing not only enhances the system's fault tolerance but also provides a robust framework for handling large-scale data, ensuring that real-time billing computations can scale effectively as new producers or consumers are added. While the current implementation demonstrates efficient resource usage and billing, future enhancements, such as predictive billing using machine learning models, or transitioning to cloud-based Kafka services, could further improve scalability and reliability.

## 10. FUTURE SCOPE

The future scope of this project lies in its potential to revolutionize billing systems in distributed computing environments by integrating advanced technologies and expanding its functionality. One promising direction is the incorporation of predictive analytics using machine learning to forecast resource consumption patterns and dynamically adjust billing rates based on predicted usage, enhancing cost transparency and efficiency. Additionally, transitioning to managed cloud-based Kafka services would enable the system to scale seamlessly for enterprise-grade applications, reducing the operational burden of managing Kafka clusters while enhancing fault tolerance and throughput. Another avenue for exploration is the development of a multi-cloud billing system, where the architecture supports resource monitoring and billing across heterogeneous environments, offering organizations a unified view of costs in hybrid or multi-cloud setups. Furthermore, the implementation of blockchain technology could provide tamper-proof transaction records and ensure greater trust and transparency in billing operations. By embedding robust security measures like real-time encryption and data anonymization, the system could also cater to privacy-sensitive industries, such as healthcare and finance. These enhancements would position the project as a cutting-edge solution in distributed resource billing, capable of adapting to the evolving demands of modern computing landscapes.

## References

1. Bulla D M, Udupi V R. Cloud Billing model:a review [J].International Journal of Computer Science and Technologies (IJCSIT), 2014,5(2):1455-1458

2. Tai S,Nimis J,Lenk A,et al. Cloud service engineering[C]//2010 ACM/IEEE 32nd International Conference on Software Engineering.IEEE, 2010,2:475-476

3. Gousios G,Loverdos C K K,Louridos P,et al.Aquarium:An Extensible Billing Platform for Cloud Infrastructures[J]. 2012.

4. Wachs M,Xu L,Kanevsky A,et al. Exertion-based billing for Cloud Storage Access[C]//Hot Cloud.2011

5. Kreps J,Narkhede N,Rao J.Kafka: A distributed messaging system for log processing[C]//Proceedings of the NetDB. 2011,11:1-7

6. Sui B. Li Q. Wen J, et al. Billing System Design of Cloud Services[C]//2018 3rd International Conference on Control, Automation and Artificial Intelligence (CAAI 2018). Atlantis Press, 2018.

7. Youseff L. Butrico M, Da Silva D. Toward a unified ontology of cloud computing[C]//2008 Grid Computing Environments Workshop. IEEE, 2008: 1-10.mg[C]//2010 Engineering

8. Li A, Yang X, Kandula S, et al. CloudCmp: comparing public cloud providers[C]//Proceedings of the 10th ACM SIGCOMM conference on Internet measurement. 2010: 1-14.