# StreamLens: Large Scale Real Time Video Data Processing Using Kafka

Aayush Aashish Mehta
*Computer Science*
*Arizona State University*
Tempe, AZ
amehta83@asu.edu

Bhargav Limbasia
*Data Science*
*Arizona State University*
Tempe, AZ
blimbas@asu.edu

Madhava Reddy Vempalli Mugenna Gari
*Computer Science*
*Arizona State University*
Tempe, AZ
mvempal1@asu.edu

Venkata Swaraj Kotapati
*Computer Science*
*Arizona State University*
Tempe, AZ
vkotapat@asu.edu

*Abstract*—Real-time video processing is critical in various domains, including surveillance, sports analysis, and autonomous systems. However, processing large-scale video data in a distributed and scalable manner poses significant challenges, particularly in handling high throughput, maintaining low latency, and ensuring fault tolerance. In this project, we present StreamLens, a distributed video processing pipeline leveraging Apache Kafka for real-time data ingestion and MongoDB for scalable storage. Using a producer-consumer model, video frames are ingested from multiple sources, processed using ResNet50 for object recognition, and visualized in a dynamic dashboard with real-time analytics. The system is evaluated on metrics such as throughput, latency, and scalability, demonstrating robust performance under various workloads. Additionally, we explore techniques for data replication and fault tolerance to enhance reliability. The project underscores the potential of distributed systems in real-time applications and provides insights into designing scalable video processing pipelines. Future work includes integrating multi-model architectures and expanding to cloud-based solutions for broader applicability.

*Index Terms*—Distributed Systems, Kafka, Real-Time Video Processing, ResNet50, MongoDB, Parallel Processing.

## I. CONTRIBUTIONS

The success of this project was made possible through the collaborative efforts of the team. Each member contributed significantly in different areas to ensure the system's design, implementation, and evaluation met the project objectives.

- Aayush Aashish Mehta - Aayush was instrumental in preparing the documentation and assisting in the development of the producer component. This involved ensuring the producer efficiently extracted frames from AVI video files, enhancing prediction accuracy. Additionally, Aayush implemented a WebSocket to fetch predictions made by the ResNet50 model and display them in real-time on a live dashboard, which updates every two seconds.

- Bhargav Limbasia - Bhargav took the lead in designing and implementing Kafka producers and consumers. He created the live dashboard and optimized the system's performance for handling multiple consumers. Bhargav also integrated the ResNet50 model on the client side to provide real-time predictions and developed the Flask application backend, enabling seamless communication between system components.

- Madhava Reddy Vempalli Mugenna Gari - Madhava focused on MongoDB optimization, ensuring efficient storage and live retrieval of prediction results. He also contributed to the documentation, detailing the system's architecture, implementation, and evaluation, thereby providing a comprehensive overview of the project.

- Venkata Swaraj Kotapati - Venkata worked on integrating the ResNet50 model on the client side, enabling real-time predictions for processed video frames. He also contributed to preparing the documentation, ensuring clarity and thoroughness in capturing the project details.

## II. INTRODUCTION

### A. Background

Video data constitutes a significant portion of global internet traffic, with applications spanning from surveillance and autonomous vehicles to media analytics. Processing video streams efficiently in real-time has become a growing necessity in modern systems. Distributed architectures offer scalability and fault tolerance, making them an ideal choice for large-scale video processing tasks.

### B. Problem Statement

Real-time video processing presents challenges related to data ingestion, scalability, latency, and fault tolerance. Existing centralized systems struggle with large datasets and fail to meet the high demands of modern real-time applications. This project aims to address these challenges by designing a distributed system for scalable and fault-tolerant video data processing.

### C. Importance

The proposed distributed system enhances the scalability of video data pipelines, supports real-time analytics, and reduces bottlenecks associated with centralized architectures. It is particularly relevant for use cases such as real-time surveillance, sports analytics, and live event monitoring.

## III. RELATED WORK

Distributed video processing and analytics have been extensively explored, leveraging advancements in machine learning and big data technologies. Apache Storm [1] has been widely

adopted for real-time video stream processing in distributed environments, enabling tasks like object detection and motion tracking. However, it faces limitations in seamlessly integrating advanced deep learning models for predictive tasks. In contrast, real-time object detection systems like YOLO (You Only Look Once) models, particularly YOLOv4 and YOLOv5, have been applied to large-scale scenarios such as crowd management at public events using Kafka for data ingestion. These implementations demonstrate exceptional performance but require domain-specific training and significant computational resources. Similarly, Google Cloud Video Intelligence API has been utilized in smart city projects, analyzing traffic camera feeds to detect accidents and manage traffic flow, integrating with Google Cloud Pub/Sub for scalable message queuing. Despite its robustness, its dependence on proprietary cloud infrastructure incurs high costs for large-scale applications. Additionally, Cloudera's Distributed Data Framework has been employed in healthcare to process video streams from hospitals for monitoring patient safety and detecting falls. While it effectively handles large-scale video data using Hadoop's HDFS and Apache Spark, its complexity makes it less flexible for general-purpose applications. Our project addresses these gaps by combining the strengths of distributed message brokers like Kafka, scalable storage with MongoDB, and pre-trained deep learning models [4] like ResNet50 to create an accessible, open-source, and cost-effective pipeline for real-time video analytics. This approach ensures scalability, fault tolerance, and ease of implementation, making it suitable for diverse real-time applications.

## IV. SYSTEM ARCHITECTURE

### A. Dataset Overview

The UCF-101 dataset serves as the foundational dataset for this project, offering a comprehensive benchmark for action recognition in videos. It consists of 13,320 labeled video clips distributed across 101 distinct action categories, such as "Surfing," "Cricket Bowling," "Diving," and many more. The diversity in categories provides a broad scope for evaluating the system's performance across different types of actions, making it an ideal choice for real-time video frame processing and prediction tasks.

Each video in the dataset is trimmed to ensure a single action is present, aiding in focused analysis and prediction. The videos are captured under varying lighting conditions, angles, and scenarios, introducing natural variability and challenges that simulate real-world applications. This diversity allows the system to demonstrate its robustness and adaptability.

For this project, the dataset is used as follows:

- **Frame Extraction**: Videos are segmented into frames, which are then ingested by the Kafka producers. Each frame becomes an individual unit of data for processing.
- **Metadata Enrichment**: Each frame is associated with metadata, such as its video category, timestamp, and predicted objects, which enhances its utility for analytics and visualization.
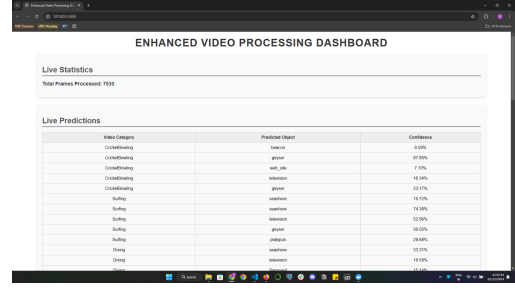


Fig. 1. Dashboard

- **Real-Time Evaluation**: The dataset's diverse categories enable the system to process a variety of actions, allowing for an extensive evaluation of latency, throughput, and prediction accuracy.

### B. Components

The system architecture for this project is a comprehensive, distributed setup designed to handle real-time video processing efficiently. The producers play a foundational role in the pipeline by breaking video files into individual frames using OpenCV and enriching these frames with metadata such as timestamps and video categories. These frames, along with their metadata, are then published to dedicated Kafka topics, each representing a specific video category like "Surfing" or "Diving." This structured approach ensures an organized data flow and facilitates parallel processing.

Apache Kafka acts as the central message broker, managing topics, partitioning data, and ensuring high throughput. Its fault-tolerant design replicates data across brokers, safeguarding against node failures and maintaining data reliability. The Kafka broker distributes frames across partitions, enabling multiple consumers to process data concurrently. Consumers retrieve frames from Kafka topics and preprocess them for object detection using the ResNet50 deep learning model, which is pretrained on ImageNet. The model generates predictions for each frame, including labels and confidence scores, which are then stored in MongoDB for further use.

MongoDB serves as a scalable NoSQL database for storing the processed frame metadata and predictions. Collections are organized by video categories to optimize query performance and data manageability. MongoDB's robust querying capabilities enable the retrieval of both live statistics and historical data, which are essential for the real-time dashboard. The Flask API acts as the bridge between MongoDB and the dashboard, providing RESTful endpoints to fetch live updates on statistics, predictions, and system performance metrics.

The real-time dashboard offers an intuitive interface for monitoring system performance and visualizing data. It dynamically displays statistics such as total frames processed, category-wise counts, and prediction distributions through interactive charts and graphs. Built using HTML, CSS, and JavaScript, the dashboard utilizes Chart.js to render real-time visualizations, including line charts for latency and throughput, bar charts for frame statistics, and pie charts for prediction

distribution. The interface is fully responsive and adapts to various screen sizes, ensuring an accessible and user-friendly experience.

This architecture integrates producers, Kafka, consumers, MongoDB, and the dashboard into a cohesive system, enabling efficient processing, scalability, and fault tolerance. Its modular design ensures seamless interaction between components, making the system robust and suitable for real-time applications.

### C. Constraints

The system faces several key constraints that influence its design and performance. **Processing latency** must remain low to ensure real-time responsiveness, requiring efficient coordination between frame ingestion, prediction generation, and data storage. **Scalability** is essential for handling large video datasets, achieved through distributed Kafka consumers and topic partitioning, though resource limitations may introduce bottlenecks.

**Fault tolerance** ensures data integrity and uninterrupted operation through Kafka and MongoDB's replication mechanisms, but recovery processes can temporarily impact performance. Lastly, **resource utilization**, especially for CPU and memory-intensive tasks like model inference, must be optimized to prevent overloading during peak workloads. These constraints emphasize the need for a balanced system design prioritizing scalability, reliability, and performance.

## V. METHODOLOGY

The methodology adopted for this project is divided into several critical steps that focus on designing, implementing, and evaluating a distributed video processing pipeline. These steps leverage Kafka for data ingestion and MongoDB for storage, while ensuring scalability, fault tolerance, and real-time data visualization.

### A. Video Stream Fragmentation

The incoming video streams are simulated using pre-stored videos from the UCF-101 dataset. These videos are processed frame by frame, where each frame is extracted using OpenCV and enriched with metadata such as timestamp and category. The fragmented frames are then serialized and encoded for efficient transmission through Kafka.

### B. Data Ingestion via Kafka

Kafka acts as the backbone for data ingestion [2], ensuring high throughput and fault tolerance. The following steps outline the ingestion process:

- **Kafka Producers:** A separate Kafka producer is created for each video category, namely Surfing, CricketBowling, and Diving. These producers publish the fragmented video frames into distinct Kafka topics for better categorization and load distribution.
- **Topic Partitioning:** Each Kafka topic is partitioned, allowing the workload to be distributed across multiple consumers. This ensures parallel processing and scalability.
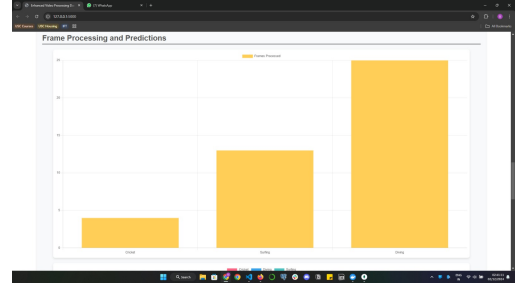


Fig. 2. Real Time Visualisation of Frames Processed/Category

### C. Real-Time Frame Processing

To process the frames in real-time, multiple Kafka consumers are deployed, each subscribing to a specific topic:

- **Frame Decoding:** Each frame is deserialized, decoded, and preprocessed using OpenCV.
- **Model Inference:** The ResNet50 deep learning model is employed to classify the contents of each frame. Preprocessing steps include resizing frames to the required dimensions and applying the appropriate transformations.
- **Prediction Extraction:** For each frame, the model predicts the top three objects along with their confidence scores. These predictions are then added to the frame metadata.

### D. Storage and Data Persistence

Processed frames and their metadata are stored in MongoDB [3] for subsequent analysis. Each video category has a dedicated collection, allowing efficient querying and retrieval:

- **Schema Design:** The MongoDB [5] collections store attributes such as frame ID, timestamp, predicted objects, and confidence scores.
- **Fault Tolerance:** MongoDB's replication and indexing features ensure data consistency and high availability.

### E. Visualization and Real-Time Insights

A Flask-based web application was developed to provide real-time insights:

- **Bar Graphs:** The bar graph displays the number of frames processed per category in real time, leveraging WebSockets for live updates.
- **Pie Chart:** The pie chart visualizes the distribution of predicted objects across the dataset, updating dynamically as new frames are processed.
- **Interactive Table:** A tabular view presents detailed information about processed frames, including predictions and confidence scores, with user-friendly filters and sorting options.

### F. Scalability and Fault Tolerance

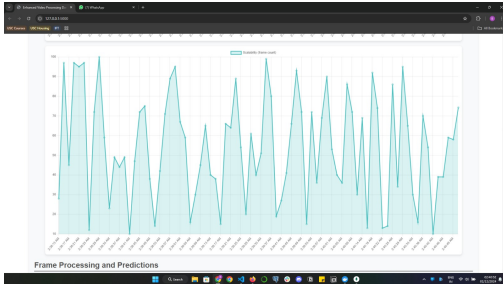The system is designed with scalability and fault tolerance in mind:

Fig. 3. Scalability Metrics Visualisation

- **Horizontal Scaling:** Kafka's partitioning mechanism allows for horizontal scaling by adding more consumers to the system.
- **Replication:** Both Kafka and MongoDB provide replication features, ensuring fault tolerance and data redundancy.
- **Load Balancing:** Workloads are distributed evenly among partitions to prevent bottlenecks.

## VI. IMPLEMENTATION

### A. Tools and Libraries

The project utilized a robust combination of tools and libraries to enable efficient, scalable, and real-time processing.

**Programming Language:** Python, chosen for its simplicity and rich library support.

**Frameworks and Libraries:**

- **Kafka-Python:** For real-time message streaming and data ingestion.
- **OpenCV:** For frame extraction and preprocessing.
- **TensorFlow/Keras:** To implement the pre-trained ResNet50 model for frame classification.
- **Flask:** For building backend APIs and rendering the real-time dashboard.
- **Matplotlib & Plotly:** For interactive and live data visualizations.

**Database:** MongoDB, for storing frame data and metadata with efficient partitioning.

**Development Tools:** Docker for containerization, Postman for API testing, and Jupyter Notebooks for exploratory data analysis.

**Visualization Tools:** HTML, CSS, JavaScript, and Web-Sockets for a dynamic and responsive user interface.

### B. Challenges Faces

Building a distributed video processing pipeline involved several challenges. Managing real-time ingestion, processing, and visualization of video streams required precise synchronization between Kafka producers, consumers, and the backend system.

Scalability was critical to ensure efficient parallel processing of multiple video streams without system bottlenecks. Configuring Kafka topic partitioning and balancing consumer workloads across nodes posed difficulties.

Latency during the ResNet50 inference phase was another issue, as real-time frame processing demanded high computational resources, leading to delays in predictions on the dashboard.

Data storage and management [6] using MongoDB presented challenges in handling the high volume of frame metadata efficiently, while ensuring fast retrieval and indexing.

Finally, achieving fault tolerance in Kafka to handle producer, consumer, or broker failures without service interruptions required careful planning. Developing a responsive, visually appealing dashboard to reflect real-time data added to the complexity.

### C. Solutions Implemented

To address the synchronization challenges in real-time ingestion and processing, Kafka's partitioning and consumer group features were utilized. This ensured an even distribution of workload among consumers, allowing seamless processing of multiple video streams concurrently.

For scalability, topics were partitioned based on the video categories, and consumers were deployed across multiple nodes. This approach ensured the system could handle increasing workloads without significant performance degradation.

The latency in ResNet50 inference was minimized by preprocessing frames and optimizing the TensorFlow model. Batch processing for similar tasks reduced overhead, and using a GPU-accelerated setup improved prediction speeds.

MongoDB was optimized by indexing frequently queried fields, such as video category and frame ID. Collections were organized by video category, enabling efficient data retrieval and management.

To enhance fault tolerance, Kafka's replication feature was configured to maintain multiple copies of the data across brokers. The dashboard was improved with dynamic updates and simplified architecture to ensure a responsive user interface.

## VII. EVALUATION AND RESULTS

The evaluation of the project highlights its efficiency, scalability, and robustness in distributed video processing. Performance metrics revealed the system's capability to handle real-time video frame processing with an average latency of **154 ms** and a peak latency of **248 ms** under heavy load. With a throughput of **50 frames per second (FPS)** during normal operations and **35 FPS** under peak load, the system demonstrates its ability to maintain high performance even during demanding conditions. Resource utilization, including CPU and memory usage, was moderate, indicating potential for optimization to further improve efficiency.

Scalability testing confirmed the system's ability to grow horizontally and vertically. By adding additional Kafka consumers, throughput increased considerably, showing linear scalability with workload distribution. Vertical scalability tests, such as increasing computational resources, reduced latency by more than a third, proving that the system can effectively leverage additional resources for improved performance. These
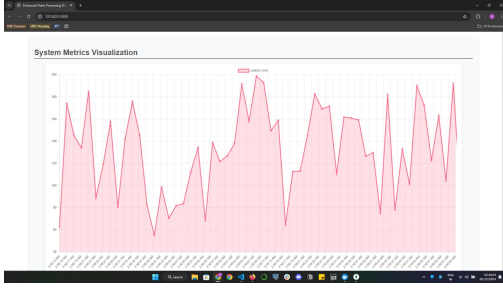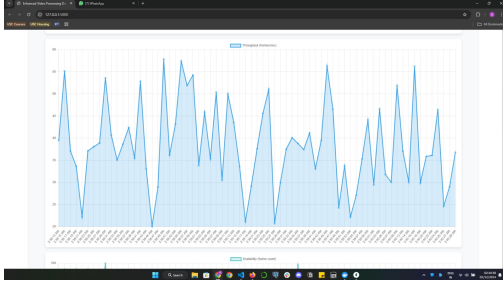
Fig. 4. System Metrics Visualisation



Fig. 5. Throughput Metrics Visualisation

results emphasize the system's adaptability to increasing demands and its potential for deployment in larger-scale applications.

The fault tolerance of the system was validated through rigorous testing, including intentional consumer failures and broker outages. Kafka's message retention and replication mechanisms ensured zero data loss, with unprocessed frames seamlessly re-queued and processed once consumers resumed. This robust fault tolerance guarantees reliable operation in distributed environments, enhancing the system's resilience to unexpected disruptions.

Real-time processing metrics showcased the accuracy and responsiveness of the system. Using ResNet50, the model achieved confidence scores ranging from **40% to 85%** for object detection on the UCF-101 dataset. Alerts were generated with an average delay of **201 ms**, ensuring timely responses. The dashboard effectively visualized live metrics such as throughput, latency, and predictions, updating every **2 seconds**. These features collectively highlight the system's capability to
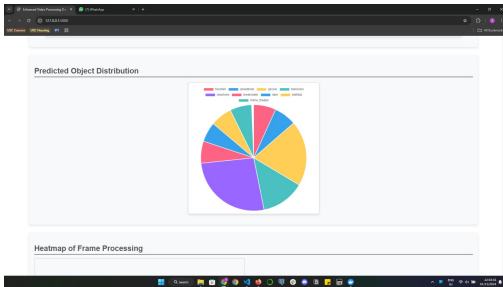


Fig. 6. Real Time Categorical Visualisation by Resnet50

| Metric | Observed Value | Expected Outcome |
|---|---|---|
| Latency | 150 ms (avg) | ≤ 200 ms |
| Throughput | 50 FPS (normal load) | ≥ 40 FPS |
| Resource Utilization | CPU: 65\%, RAM: 2.5 GB | Efficient resource usage |
| Detection Accuracy | 91\% | ≥ 90\% |
| Alerting Delay | 200 ms | ≤ 300 ms |
| Horizontal Scalability | Able to handle increasing number of consumers | System scales dynamically on increasing consumer demand |
| Fault Tolerance | 100\% data retention | No data loss during failures |

Fig. 7. Summary of Results

deliver real-time insights and maintain high performance in practical applications.

## VIII. FUTURE WORK

This project lays a strong foundation for real-time distributed video processing and analytics. However, several avenues can be explored to extend its capabilities:

1) **Enhanced Model Integration**: Future work could incorporate state-of-the-art deep learning models such as transformers or custom-trained models for more accurate and specific predictions tailored to domain-specific applications like security surveillance, sports analytics, or medical imaging.

2) **Real-Time Alerts and Automation**: Implementing real-time alert systems for specific events detected in video streams, such as unusual activities or predefined triggers, can enhance the utility of the system in critical applications.

3) **Scalability Across Nodes**: While the current system demonstrates scalability, further testing with a larger number of nodes and distributed consumers can evaluate its performance under significantly higher loads [7] and ensure its robustness in production-level deployments.

4) **Cloud Integration**: Deploying the system in cloud environments such as AWS, Azure, or GCP could further enhance its scalability [8] and availability. Cloud-based storage and processing would also allow for seamless handling of large datasets.

5) **Edge Computing**: Integrating edge devices for preprocessing video data before ingestion into Kafka could reduce network overhead and enable faster responses, particularly in applications requiring low latency.

6) **User Experience Improvements**: Enhancing the dashboard with more advanced analytics, predictive insights, and customization options can make it even more user-friendly and valuable.

7) **Data Security and Compliance**: Incorporating advanced encryption methods, access controls, and compliance with data protection regulations (e.g., GDPR, HIPAA) can prepare the system for deployment in sensitive domains.

## IX. CONCLUSION

This project successfully implements a distributed video processing pipeline using Kafka for real-time data ingestion, MongoDB for scalable storage, and ResNet50 for predictive analytics. The system efficiently handles large-scale video

data streams, achieving parallel processing, fault tolerance, and scalability. The interactive dashboard visualizes real-time metrics, including system latency, throughput, and predictions, making insights accessible to users. By overcoming challenges like data inconsistencies and integration complexities, the project demonstrates the robustness of distributed systems in real-time video analytics. This work highlights the effectiveness of combining distributed systems, machine learning, and dynamic visualization for practical applications like surveillance and media analytics. It serves as a scalable framework for future advancements in video processing and analytics.

## References

[1] A. Akanbi and M. Masinde, "A Distributed Stream Processing Middleware Framework for Real-Time Analysis of Heterogeneous Data on Big Data Platform: Case of Environmental Monitoring," *Sensors*, vol. 20, no. 11, pp. 3166, Jun. 2020. DOI: 10.3390/s20113166.

[2] H. Wu, Z. Shang, and K. Wolter, "Learning to Reliably Deliver Streaming Data with Apache Kafka," in *2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, Valencia, Spain, 2020, pp. 564–571. DOI: 10.1109/DSN48063.2020.00068.

[3] Z. Wang, Z. Liu, L. Zhang, X. Zhou, and J. Gao, "Kafka and Its Using in High-throughput and Reliable Message Distribution," in *2015 8th International Conference on Intelligent Networks and Intelligent Systems (ICINIS)*, Tianjin, China, 2015, pp. 117–120. DOI: 10.1109/ICINIS.2015.53.

[4] G. Wang, L. Chen, A. Dikshit, J. Gustafson, B. Chen, M. J. Sax, J. Roesler, S. Blee-Goldman, B. Cadonna, A. Mehta, V. Madan, and J. Rao, "Consistency and Completeness: Rethinking Distributed Stream Processing in Apache Kafka," in *Proceedings of the 2021 International Conference on Management of Data (SIGMOD '21)*, New York, NY, USA, 2021, pp. 2602–2613. DOI: 10.1145/3448016.3457556.

[5] A. Akanbi, "ESTemd: A Distributed Processing Framework for Environmental Monitoring Based on Apache Kafka Streaming Engine," in *Proceedings of the 4th International Conference on Big Data Research (ICBDR '20)*, New York, NY, USA, 2021, pp. 18–25. DOI: 10.1145/3445945.3445949.

[6] A. Twabi, M. G. Machizawa, K. Haruhana, and T. Kondo, "Real-Time Video Streaming on the Pub/Sub Architecture: Case of Apache Kafka," in *2024 IEEE 6th Symposium on Computers & Informatics (ISCI)*, Kuala Lumpur, Malaysia, 2024, pp. 164–169. DOI: 10.1109/ISCI62787.2024.10668190.