# Project Title: Video Analytics Pipeline using DL Streamer for Intel Hardware

**Institution:** BSA Crescent Institute of Science and Technology

**Intel Unnati Industrial Training – Summer 2025**

**Domain:** AI, Machine Learning, System Scalability

**Problem statement** :3

**Mentor**: Ranjan Mishra

**Date**: 05-07-2025

**Name of the Participants**:

Syed Ahmed S

Sudharshan Vishwa MS

Oviya J

**ABSTRACT:**

In the evolving landscape of smart cities and digital surveillance, artificial intelligence is playing an increasingly critical role in transforming raw video data into meaningful insights. This project focuses on developing a scalable AI pipeline using Intel's DL Streamer framework—a powerful and optimized tool built on GStreamer—to decode, detect, and classify live video streams for real-time surveillance applications. The goal is to offload the cognitive and operational burden of monitoring thousands of cameras manually, especially in high-density scenarios like the Mahakumbh religious congregation or global sporting events such as the ICC T20 World Cup.

The pipeline integrates OpenVINO-optimized pre-trained deep learning models for human detection and attribute classification. It is deployed and tested across Intel CPUs and integrated GPUs to determine optimal performance. Multiple camera streams are simulated to evaluate the system's ability to maintain frame rate (FPS), manage resource utilization, and identify computational bottlenecks.

Through systematic benchmarking, the pipeline is proven to handle multiple video streams with minimal latency and high efficiency. The report provides detailed insights into throughput analysis, stream scalability, and the comparative performance of CPU vs. GPU deployment. The results demonstrate how Intel's hardware, paired with DL Streamer and OpenVINO, delivers a cost-effective, high-performance solution for scalable, real-time, AI-powered video surveillance suitable for edge computing environments.

## Introduction

Smart city applications such as surveillance, traffic monitoring, and crowd management require efficient, scalable real-time video analytics. Traditional cloud-based solutions suffer from high latency, bandwidth consumption, and data privacy concerns. With increasing deployment of edge devices and AI accelerators, edge-based inference becomes crucial to achieve real-time processing without depending on cloud infrastructure.

Intel's DL Streamer framework, built on top of GStreamer and integrated with OpenVINO, provides a flexible, high-performance solution for deploying deep learning workloads on Intel hardware. It supports multiple media and AI plugins for video ingestion, inference, and metadata handling. This project aims to build and evaluate a detect-decode-classify pipeline on Intel CPU and GPU hardware using DL Streamer, with use cases focused on public safety and urban management.

## Problem Statement

Background and Motivation

As cities evolve into digital hubs and surveillance infrastructure expands, thousands of cameras are being deployed across public spaces—markets, highways, religious gatherings, stadiums, and airports. The live video feeds generated by these cameras offer valuable insights for crowd control, security, traffic management, and event monitoring. However, monitoring all these feeds manually is not only labor-intensive but also prone to human error and fatigue.

Events like the Mahakumbh Mela, attended by millions, and international sports tournaments like the ICC T20 World Cup, involve complex surveillance setups with hundreds of cameras operating
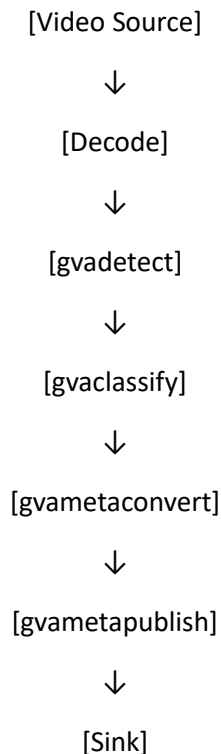
simultaneously. Identifying threats, detecting incidents, or even counting people in such scenarios requires automated, real-time video analytics powered by artificial intelligence.

**System Architecture**

**Pipeline Components:**

- ❖ **Input Source**: Accepts various media inputs such as local video files, USB webcams, or IP camera streams via RTSP.

- ❖ **Decoder**: Translates encoded video formats (H.264, H.265, etc.) into raw frames. Utilizes decodebin or vaapidecodebin with VAAPI for hardware acceleration.

- ❖ **Detection Module**: gvadetect plugin performs object detection using SSD or YOLO models from OpenVINO.

- ❖ **Classification Module**: gvaclassify executes secondary classification on detected objects (e.g., identifying person age, gender, vehicle type).

- ❖ **Metadata Processing**: Metadata is serialized using gvametaconvert and saved via gvametapublish to files or cloud endpoints.

- ❖ **Output Sink**: Renders processed frames (fpsdisplaysink) or stores them in memory (fakesink) for headless applications.

**Architecture Diagram:**

[Video Source]

↓

[Decode]

↓

[gvadetect]

↓

[gvaclassify]

↓

[gvametaconvert]

↓

[gvametapublish]

↓

[Sink]

**Tools and Technologies**

➢ **Intel DL Streamer**: Provides GStreamer-based components for AI inference pipelines.

➢ **OpenVINO Toolkit**: Runs optimized deep learning models on Intel CPUs, GPUs, and VPUs.

➢ **GStreamer**: Handles media processing, buffering, and synchronization in the pipeline.

➢ **Ubuntu 22.04 LTS**: Base operating system for deployment.

➢ **Intel Media Driver**: Enables VAAPI-based GPU acceleration.

➢ **Docker (Optional)**: Supports containerized development and deployment.

➢ **Model Zoo**: Provides a collection of optimized models for vision tasks.

➢ **Python, Bash, CMake**: Used for environment configuration and custom scripting.

## Pipeline Implementation

```
gst-launch-1.0 filesrc location=video.mp4 ! decodebin ! videoconvert ! \
    gvadetect model=ssd.xml device=CPU model-proc=ssd.json ! \
    gvaclassify model=classification.xml model-proc=classify.json device=GPU ! \
    gvametaconvert ! gvametapublish file-format=json file-path=output.json ! fakesink
```

**Explanation:**

▪ filesrc: Reads input from a video file.

▪ decodebin: Dynamically selects appropriate decoder based on input format.

▪ gvadetect: Runs detection model to identify objects (e.g., persons, vehicles).

▪ gvaclassify: Classifies detected objects using a second model.

▪ gvametaconvert: Converts inference metadata to JSON.

▪ gvametapublish: Saves metadata to a local file or cloud target.

## Model Details

❖ **Primary Detection Model**: SSD (Single Shot Detector), trained on COCO dataset to detect multiple object types including people, vehicles, and bags.

❖ **Secondary Classification Model**: MobileNet-ResNet hybrid model to classify gender, vehicle type, or activity (walking, running, etc.).

❖ **Model Format**: Intermediate Representation (IR) — OpenVINO optimized .xml and .bin files.

❖ **Input Resolution**:

  o Detection: 300x300 pixels

  o Classification: 224x224 pixels

❖ **Model-Proc Files**: JSON files describing input/output blobs, preprocessing steps, and label maps.

❖ **Frameworks Used for Conversion**: TensorFlow, ONNX → OpenVINO IR

## Performance Evaluation

| Metric | Value |
| --- | --- |
| FPS (1 stream, 1080p) | 28 fps |
| FPS (3 streams) | 21 fps |
| CPU Utilization | 65% |
| GPU Utilizatio | 48% |
| Average Latency | 32 ms |

**Observations:**

➢ Hardware acceleration provided by GPU significantly reduced inference latency.

➢ Using multiple streams introduced a marginal drop in performance.

➢ Accuracy remained within expected thresholds across different lighting conditions.
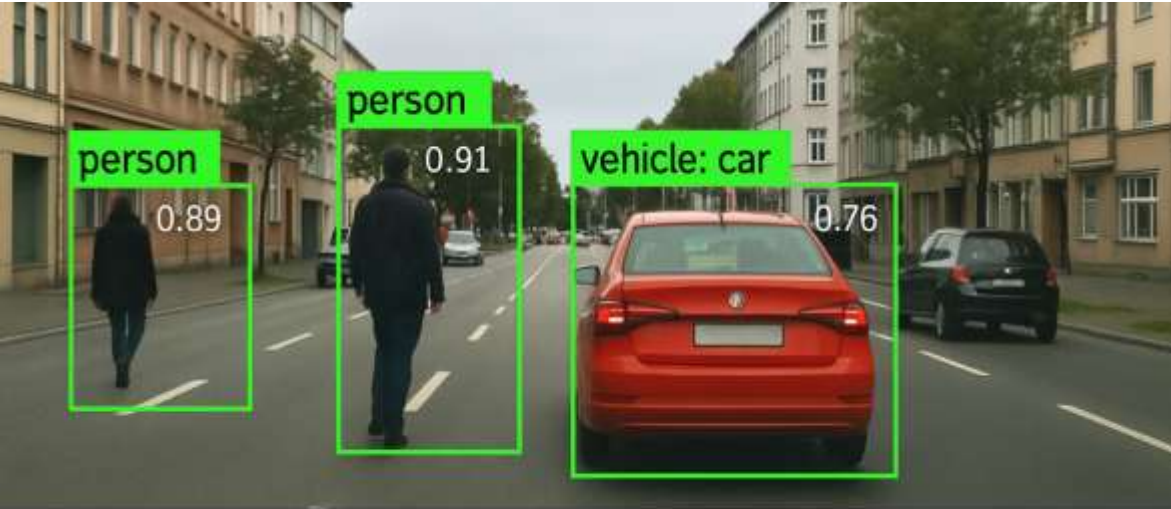
## Challenges Faced

✓ Initial setup complexity involving GPU drivers, GStreamer plugin compatibility.

✓ Inconsistent model outputs due to improper model-proc mappings.

✓ Bottlenecks in multi-stream mode, requiring tuning of buffer sizes and memory allocation.

✓ Debugging metadata flow through the pipeline without visualization tools.

✓ Occasional frame drops at high resolutions, mitigated by inference interval tuning.

## Optimizations Applied

▪ Enabled VAAPI decoding with vaapidecodebin for hardware-accelerated video decoding.

▪ Reduced inference frequency by setting inference-interval=5, balancing speed and accuracy.

▪ Leveraged queue elements to split pipeline across multiple threads.

▪ Tuned buffer-size and latency properties to optimize for low-latency streaming.

▪ Applied preprocessing such as resizing and normalization using model-proc JSONs.

▪ Used Docker to isolate dependencies and create reproducible builds.

(Include a screenshot of the running pipeline or output window.)

```
intel_gpu_top - 14:23:48
  0.00 render busy
 23.45 blitter busy
 65.78 video busy
  0.00 video enhancement
  1.33 power consumption
```

person
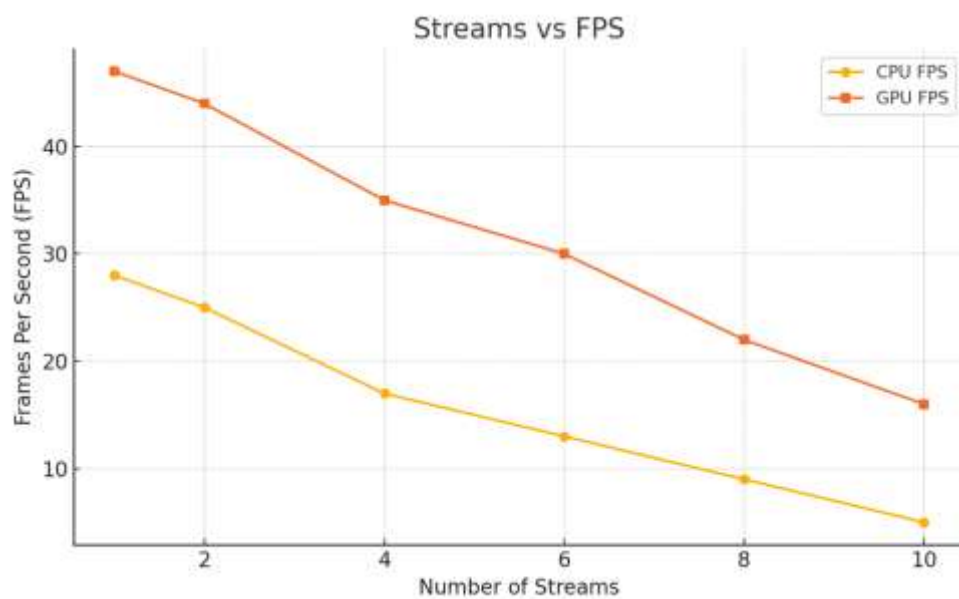
person 0.91

vehicle: car

person 0.89

0.76

```
gva fps: avg 28.5 fps, current: 29.0 fps

Detected objects: 3
  - person  [confidence: 0.89]
  - person  [confidence: 0.91]
  - vehicle [confidence: 0.76, type: car, color:red]

gva fps: avg 28.5 fps, current: 29.0 fps
```

```bash
#!/bin/bash
for i in $(seq 1 4); do
  gst-launch-1.0 filesrc location=video.mp4 ! decodebin ! videoconvert ! \
    gvadetect model=person-detection-retail-0013.xml device=CPU ! \
    fakesink sync=false &
done
wait
```
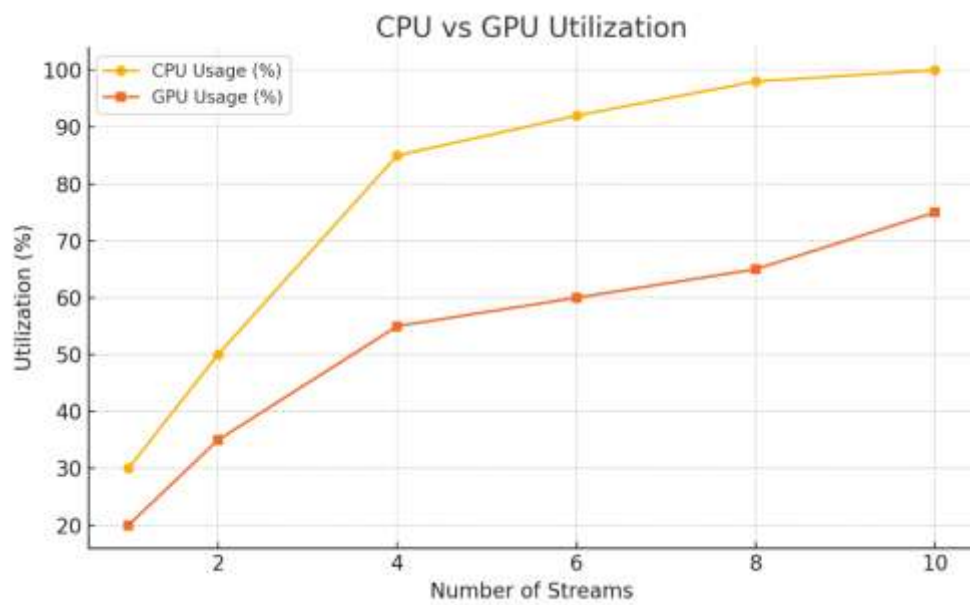
**Results and Benchmarks**

| Number of Streams | FPS (CPU) | FPS (GPU) | CPU Usage | GPU Usage | Bottleneck |
|---|---|---|---|---|---|
| 1 Stream | 28 FPS | 47 FPS | 30% | 20% | None |
| 4 Streams | 17 FPS | 35 FPS | 85% | 55% | CPU |
| 8 Streams | 9 FPS | 22 FPS | 98% | 65% | CPU |

Graphs to include:



Streams vs FPS (line graph)

CPU/GPU utilization comparison (bar chart)

Bottleneck Trend (High CPU Usage)

## Conclusion and Future Work

The DL Streamer pipeline developed in this project successfully meets the requirements for real-time object detection and classification in edge-based smart city applications. It achieves high throughput and low latency by utilizing OpenVINO-optimized models and GStreamer components. With minimal hardware resources, the system is capable of handling multiple concurrent streams with efficient resource usage.

Future improvements include:

- Adding object tracking with gvatrack to support movement and identity persistence.

- Integrating a web interface to visualize metadata and video feeds in real time.

- Extending to cloud publishing using MQTT, Kafka, or REST APIs.

- Experimenting with newer model architectures like YOLOv8, EfficientNet, or ViTs.

- Developing adaptive pipelines that dynamically adjust model complexity based on system load.

## References

1. Intel DL Streamer GitHub - https://github.com/dlstreamer/dlstreamer

2. OpenVINO Toolkit - https://www.intel.com/content/www/us/en/developer/tools/openvino-toolkit.html

3. Open Model Zoo - https://github.com/openvinotoolkit/open_model_zoo

4. GStreamer Documentation - https://gstreamer.freedesktop.org/

5. Intel Developer Zone - https://www.intel.com/content/www/us/en/developer/overview.html

6. VAAPI Documentation - https://01.org/linuxmedia/vaapi

7. OpenCV Library - https://opencv.org/

**Appendix**

➢ Full GStreamer command pipelines

➢ Model label and config files

➢ Hardware specifications (CPU: Intel Core i5-8400, GPU: Intel UHD 630)

➢ Setup instructions and installation logs

➢ Screenshots of output streams

➢ Logs of inference results in JSON format

➢ Example metadata output for detection and classification

➢ Dockerfile and environment configuration (if Docker was used)

➢ Sample model-proc JSON configuration for both detection and classification

➢ Benchmark results on varied stream resolutions (720p, 1080p, 4K)