**Indergraduate Programme** 

Student Name:

Major:

Computer Science

LI Shibiao

Project Mentor:

Prof. Fang Yulin

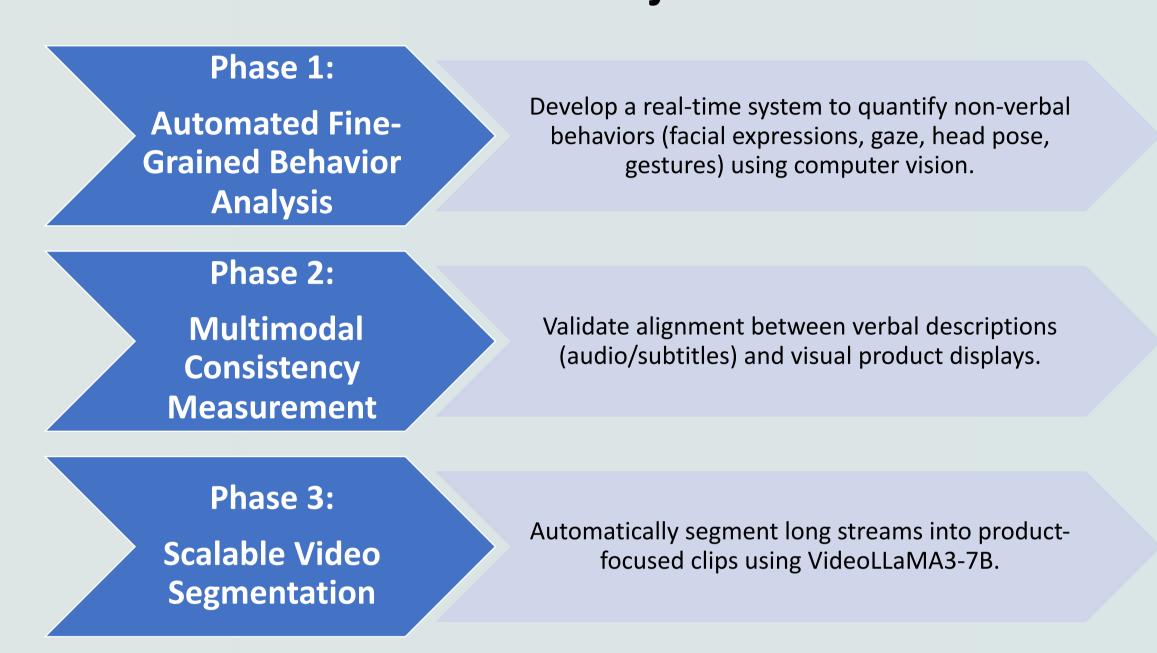
Department: Dept. of Computer Science

# Al-driven video Analysis in Live Commerce: From Automated Video Processing to Multimodal Recognition

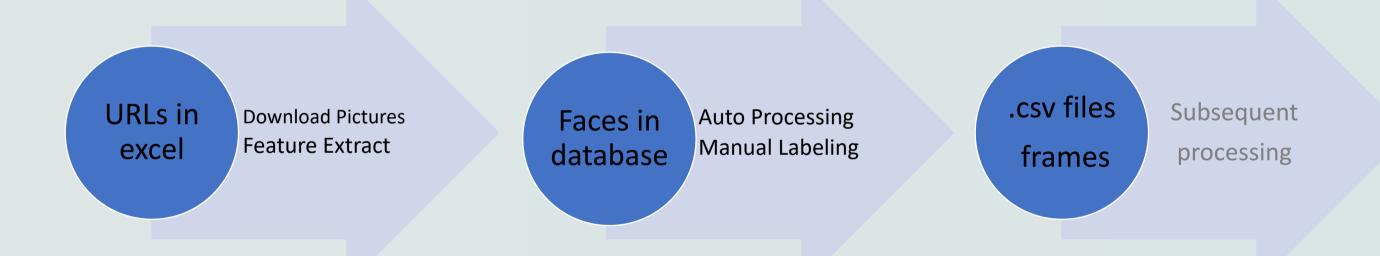
#### Introduction

The rise of live commerce has transformed product marketing, yet the impact of streamers' non-verbal behaviors (e.g., facial expressions, gestures) on consumer decisions remains understudied. Current methods rely on manual analysis, which is time-consuming and inconsistent. This project addresses this gap by developing an Al-driven pipeline to automatically extract and analyze streamer behaviors from hours of live-stream footage. By integrating computer vision, multimodal models, and human-Al collaboration, the system streamlines research workflows, enabling scalable behavioral studies in live commerce.

#### **Research objectives**

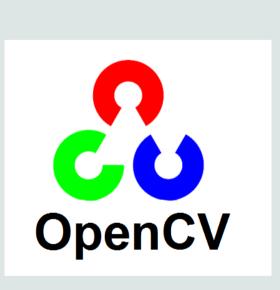


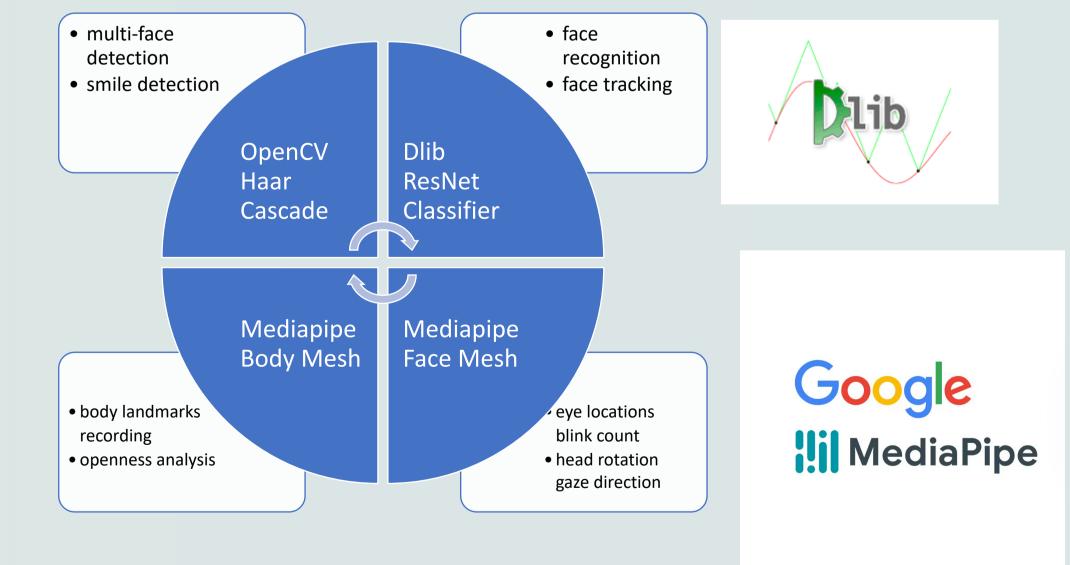
## Methods & Features Phase 1: Behavior Analysis Pipeline



Phase 1 begins by downloading streamer headshots from URLs listed in an Excel file, which are stored in a structured director. These images train a facial recognition database using Dlib and MediaPipe. During video processing, **1 frame every 20 frames** is analyzed: known faces are automatically tracked, and their behaviors (smile, gaze, pose) are logged to a CSV file. Unknown faces or low-confidence results trigger manual validation via an interface, where users can refine labels. The final CSV files and extracted frames feed into downstream analysis, balancing automation efficiency with human oversight for accuracy.

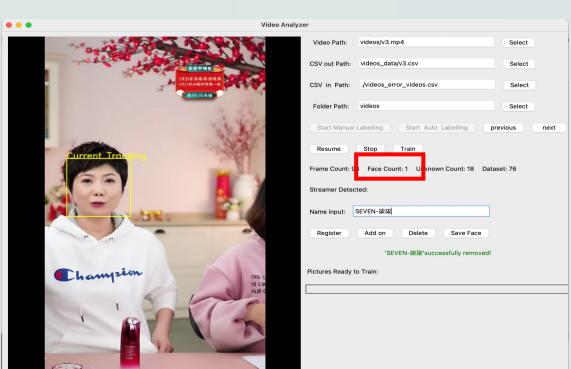
#### **Key Components of the Tech Stack of phase 1**



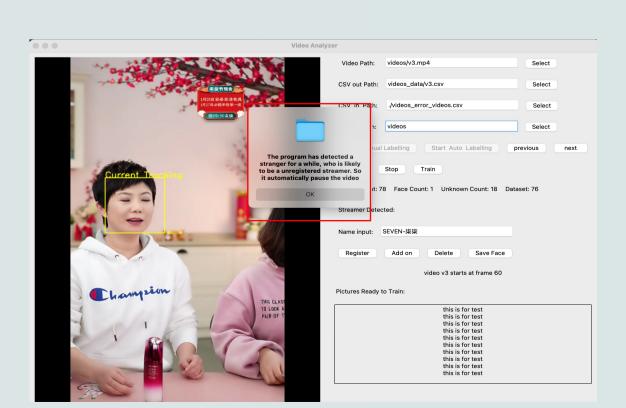


#### Adaptive Face Recognition with Human-in-the-Loop Validation

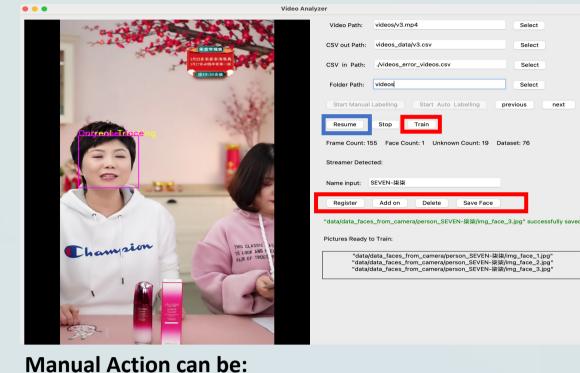
When the system detects an unregistered face (e.g., a new streamer or guest), processing is automatically paused, prompting manual intervention.



**Detection & Counting**: When an unregistered face (e.g., guest or new streamer) is detected, an *unknown* counter increments.



**Threshold Trigger**: If the counter reaches a threshold (e.g., 20 consecutive frames), processing pauses and alerts the user for manual intervention.

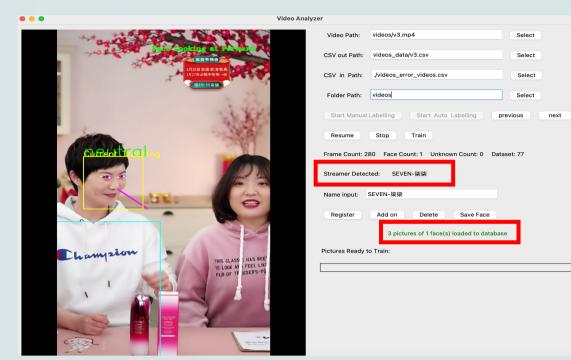


Either **Resume**: Clicking "Resume" clears the counter

(e.g., "steamer\_1").

and continues automated processing.

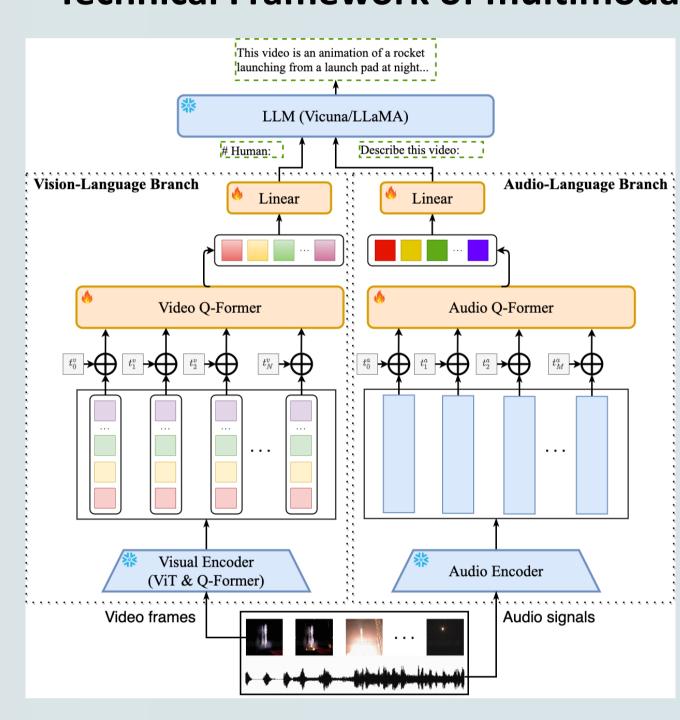
Or **Register**: Users can screenshot the frame, crop the face, and register it into the database with a unique ID



**Seamless Integration**: Once registered, the system immediately recognizes the new face in subsequent frames, enabling continuous tracking without reprocessing.

#### Phase 2 & 3: Multimodal Consistency and Customizable AI Analysis

#### **Technical Framework of multimodal AI**



The system employs a multimodal AI architecture to achieve precise video understanding and temporal localization. A Vision Transformer (ViT) processes raw video frames to capture spatial features (e.g., product placement, hand gestures), while a Video Q-Former refines these features by modeling temporal dependencies across frames, enabling detection of dynamic actions like unboxing or product comparisons. Cross-modal alignment is achieved via AudioCLIP/VideoCLIP embeddings, which project visual and audio features into a shared latent space.

For temporal precision, a **transformer-based decoder** jointly analyzes fused multimodal features and positional encodings to predict segment boundaries. This decoder leverages causal attention masks to enforce chronological constraints, achieving a timestamp error margin of **<0.5 seconds**(validated on 200 annotated streams).

Output: Structured CSV files include:

- Segment ID, Start/End Time (s),
- Product Category (e.g., cosmetics, electronics),

#### **Adaptive Processing and Prompt Engineering**

This dual-configuration approach—adjusting **FPS/token limits** and **customizing prompts**—optimizes the trade-off between computational efficiency, analysis granularity, and task-specific accuracy. It allows researchers to prioritize speed for large datasets or precision for critical segments, while domain-tailored prompts ensure outputs align with study objectives.

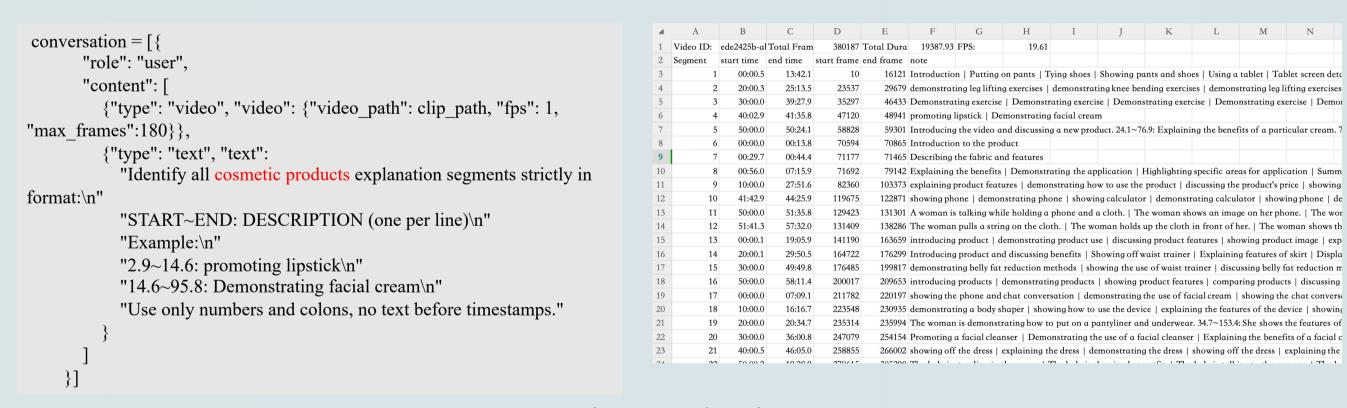
#### 1. FPS and Token Adjustment

- Long Streams (8+ hours): Set fps=0.5 and max\_frames=300 to sample 1 frame every 2 seconds, reducing GPU memory
- usage by 60% while retaining key events (e.g., product reveals).
   Detailed Analysis: Increase to fps=5 and max\_frames=900 for high-stakes segments (e.g., limited-time promotions) to capture subtle gestures (e.g., hand movements, facial expressions).

#### 2. Customizable Prompt Design

Additionally, the program supports customizable prompts. Users can specify focal points for the AI, such as prioritizing the detection of cosmetic products or electronic devices. It can also retrieve items from a specified list, and even generate annotations in Chinese for optimized outputs.

See two examples below:



### Prompt with prioritized product



Prompt with specified list and Chinese note

#### **Results & Conclusions**

The framework has successfully processed **6,000+ video clips** (1.2TB) in Phase 1 and **215GB of raw footage** in Phase 2/3, isolating **12,400 high-value segments** focused on product explanations (averaging 5.2 segments per hour). By automating content filtering (e.g., removing ads and irrelevant chatter), manual review time was reduced by approximately **85%**, slashing analysis time from about 1.5 to 5 hours per video. Leveraging GPU-accelerated processing and adaptive sampling, the system combines automated efficiency with human validation for complex scenarios (e.g., multi-host streams), while customizable prompts ensure adaptability across industries. Future developments aim to enable real-time analysis and expand compatibility with platforms like TikTok and Amazon Live, transforming unstructured video data into strategic tools for optimizing live-commerce strategies.