# HEAR ME OUT (& THINK): MAESTRO, A MULTIMODAL AGENTIC MODEL WITH EFFICIENT, SYNERGISTIC TEXT-REASONING OPTIMISATION FRAMEWORK

**Members:**
Felicia Tan Ee Shan, Low Li Ying Amy
(Raffles Institution)

**Mentor:**
Kuek Yong Jie Adriel
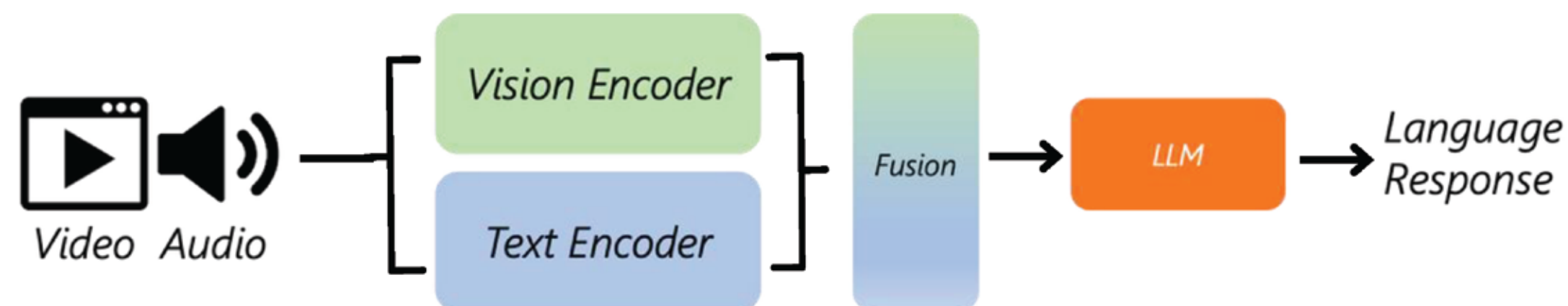(DSO National Laboratories)

## INTRODUCTION

**Specific Use-Case: Hateful Video Classification**

**500 hours** of videos/min
**1 Billion** users

1. Demands Inference Reasoning
2. Multimodal Reasoning
3. Temporal Understanding
4. Adaptability

## CURRENT VLMs



Video Audio → Vision Encoder / Text Encoder → Fusion → LLM → Language Response
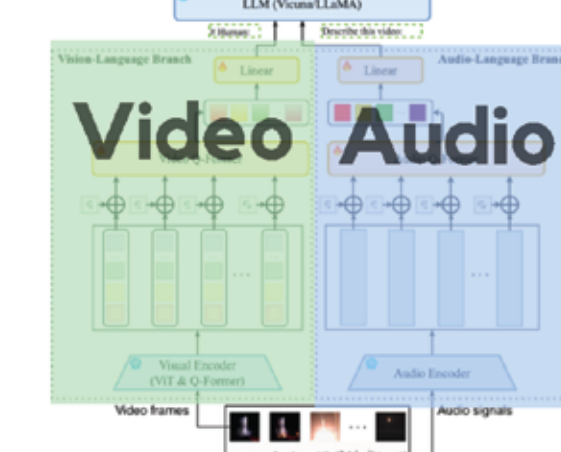
### High Computational Overhead

**Exhaustive frame-by-frame analysis** — **Use of Q-formers / transformers**

Achieving fine-grained understanding in VLMs often necessitates processing a vast number of video frames, leading to high computational costs that scale significantly with input length

### Missing integration of audio modality

**Separate Streams** — **Neglects Time Alignment**
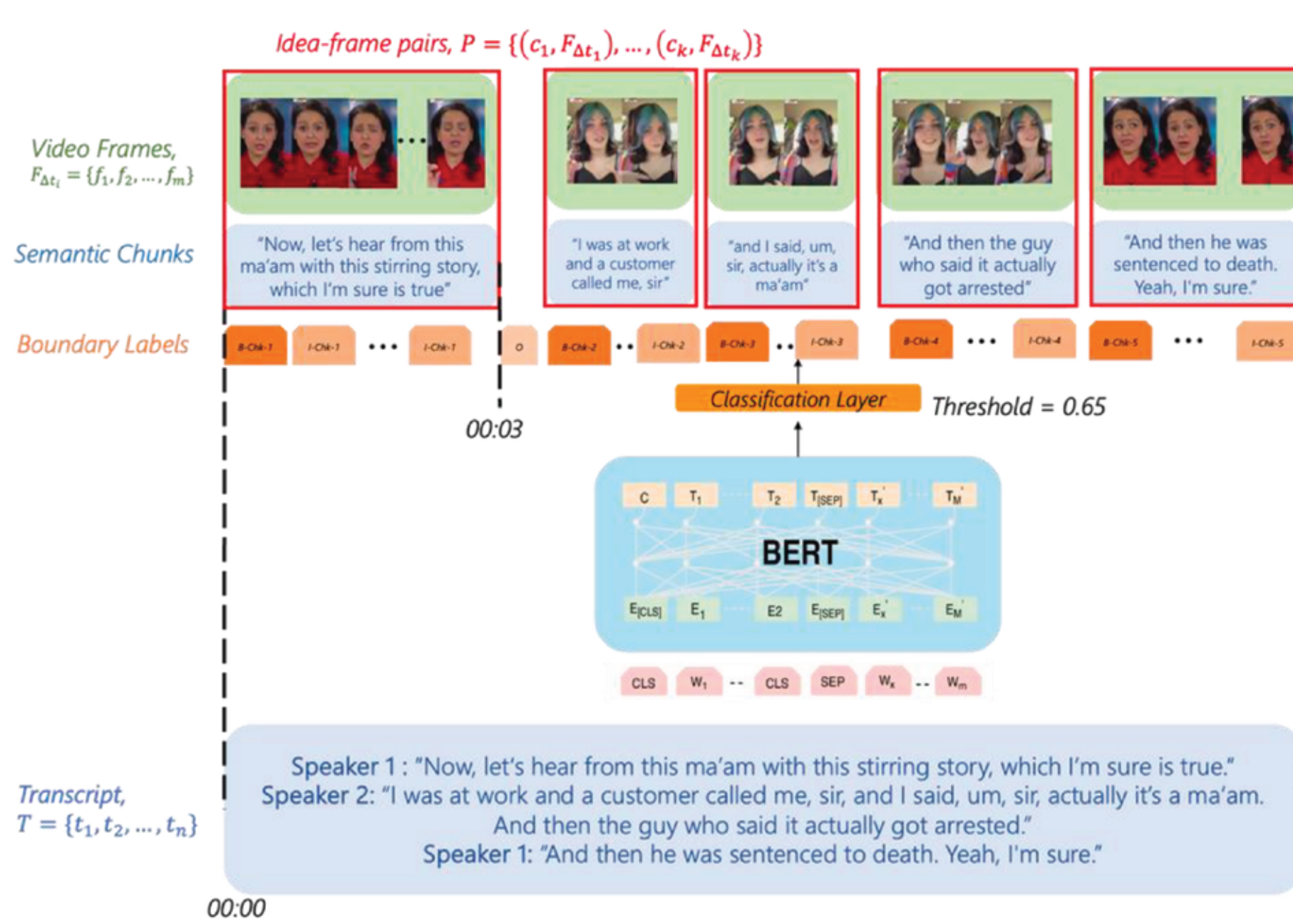
Video / Audio

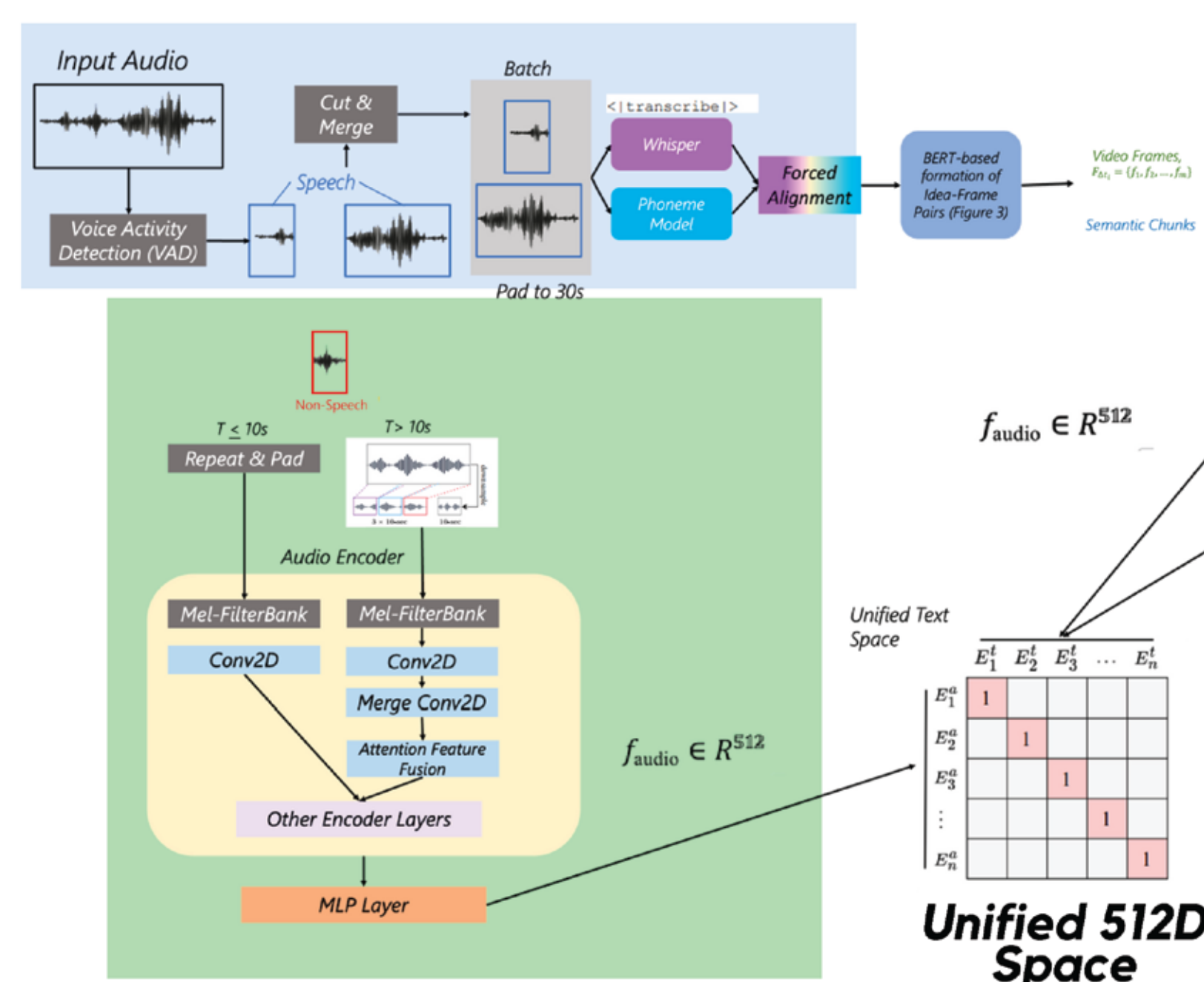### Weak/Static Reasoning

**Static Reasoning Pipeline**

USER: Is this video hateful? Explain why or why not
Video-LLaMA: Based on the video descriptions, it is not clear whether the video is hateful or not. However, the video shows a large crowd of people at an airport, and some of them appear to be excited or happy. *The video also shows a person with a phone, which could be a sign of annoyance or frustration*
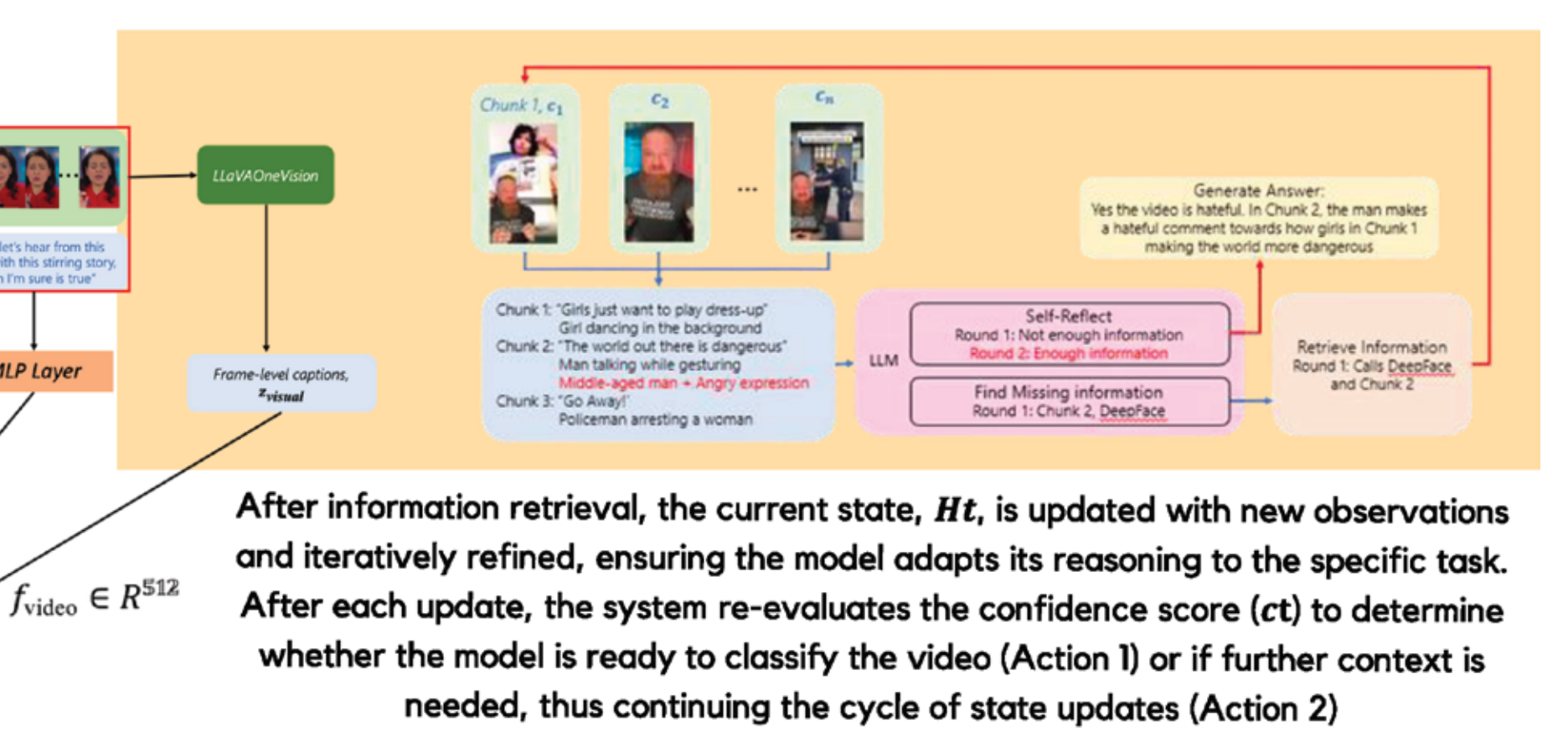
## METHODOLOGY

### Transcript Chunking → Unified Modality Alignment → Global-Local Reasoning Loop



Idea-frame pairs, $P = \{(c_1, F_{\Delta t_1}), ..., (c_k, F_{\Delta t_k})\}$

Video Frames, $F_{\Delta t} = \{f_1, f_2, ..., f_m\}$

Semantic Chunks

Boundary Labels

Classification Layer — Threshold = 0.65

BERT

Speaker 1 : "Now, let's hear from this ma'am with this stirring story, which I'm sure is true."
Speaker 2: "I was at work and a customer called me, sir, and I said, um, sir, actually it's a ma'am. And then the guy who said it actually got arrested."
Speaker 1: "And then he was sentenced to death. Yeah, I'm sure."

Transcript, $T = \{t_1, t_2, ..., t_n\}$

Input Audio — Cut & Merge — Whisper — Phoneme Model — Forced Alignment — BERT-based formation of Idea-Frame Pairs (Figure 3)

Voice Activity Detection (VAD) — Speech — Pad to 30s

$f_{audio} \in \mathbb{R}^{512}$

$f_{video} \in \mathbb{R}^{512}$

Unified 512D Space

After information retrieval, the current state, $H_t$, is updated with new observations and iteratively refined, ensuring the model adapts its reasoning to the specific task. After each update, the system re-evaluates the confidence score ($ct$) to determine whether the model is ready to classify the video (Action 1) or if further context is needed, thus continuing the cycle of state updates (Action 2)

**1 Agent, 3 Tools** — ChatGPT-3.5

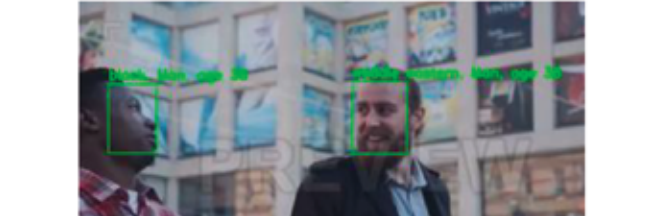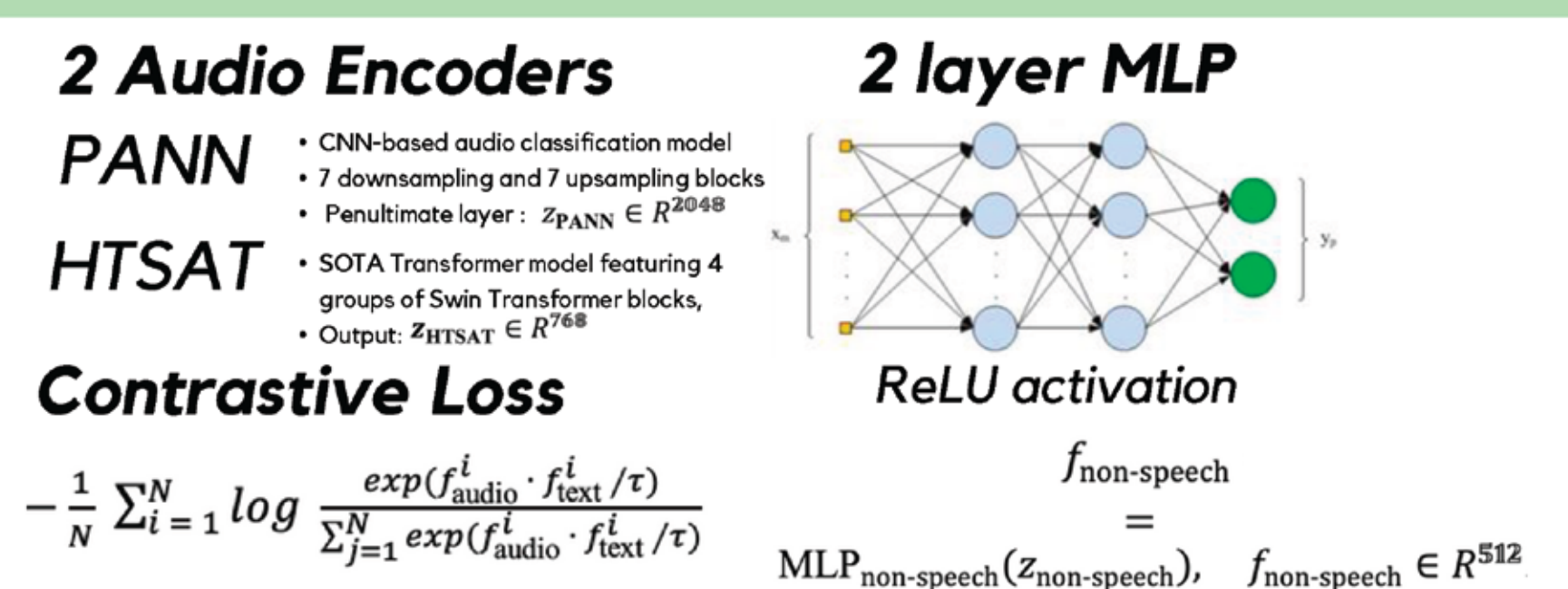| YOLO | LLaVAOneVision | videoDeepFace |
|---|---|---|
| Object & Symbol Detection | Action Recognition • Accurate Descriptions • Weak Reasoning | Race, Age, Gender, Emotion |

### Video

Image adapted from LLaVA-OneVision (Li et al.)

Qwen-2 $f_\phi$ — Language Response $X_a$

2-layer MLP — SigLIP $g_\psi$ — $Z_v$, $H_v$, $X_v$, Video — $H_q$, $X_q$, Language Instruction

**SigLIP Vision Encoder** — **Qwen-2** offers various model size and exhibits strong language capabilities to date among publicly available checkpoints

### Speech Audio

**Voice Activity Detection (VAD)**
Acoustic Features — Binary Labels
$f_\theta : A = \{a_1, ..., a_T\} \rightarrow y = \{y_1, ..., y_T\}$
$y_t \in \{0,1\}$ — Speech present

**Dynamic Time Warping** — Aligns phonemes temporally

Distance Matrix

Word boundaries derived from $t_{start}$ and $t_{end}$

**Forced Phoneme Alignment**
$L_i \in \mathbb{R}^{|C_{T_i}| \times T}$

**Min-Cut Strategy** — Overly long segments — Split at the point of minimum voice activation score

**Neighbour Merging** — Overly short segments — Below $\tau = |A_{train}|$

**Maintains Temporal Consistency**

### MAESTRO-CLAP (Non-Speech Audio)

**2 Audio Encoders**
PANN — • CNN-based audio classification model • 7 downsampling and 7 upsampling blocks • Penultimate layer : $z_{PANN} \in \mathbb{R}^{2048}$
HTSAT — • SOTA Transformer model featuring 4 groups of Swin Transformer blocks, • Output: $z_{HTSAT} \in \mathbb{R}^{768}$

**2 layer MLP** — ReLU activation

$f_{non-speech} = MLP_{non-speech}(z_{non-speech})$, $f_{non-speech} \in \mathbb{R}^{512}$

**Contrastive Loss**
$$-\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(f_{audio}^i \cdot f_{text}^i / \tau)}{\sum_{j=1}^{N} \exp(f_{audio}^i \cdot f_{text}^j / \tau)}$$

**Projected into unified space**

**Output,** $f_{output} \in \mathbb{R}^{512}$

## RESULTS

### MultiHateClip (Hateful Videos)

| Model | Acc | Binary M-F1 | F1(O) | R(O) | P(O) |
|---|---|---|---|---|---|
| mBERT | 0.57 | 0.57 | 0.52 | 0.42 | 0.68 |
| GPT-4 | 0.81 | 0.79 | 0.73 | 0.69 | 0.78 |
| Qwen | 0.72 | 0.71 | 0.65 | 0.57 | 0.75 |
| MFCC | 0.54 | 0.50 | 0.36 | 0.33 | 0.40 |
| Wav2Vec | 0.53 | 0.48 | 0.64 | 0.50 | 0.90 |
| ViViT | 0.73 | 0.73 | 0.68 | 0.57 | 0.86 |
| Vit | 0.63 | 0.58 | 0.44 | 0.46 | 0.45 |
| VLM | 0.70 | 0.64 | 0.48 | 0.59 | 0.41 |
| GPT-4V | 0.81 | 0.79 | 0.73 | 0.72 | 0.73 |
| Qwen-VL | 0.62 | 0.61 | 0.56 | 0.46 | 0.72 |
| TI ⊙ AI ⊙ V1 | 0.75 | 0.74 | 0.67 | 0.61 | 0.77 |
| MAESTRO (Ours) | 0.96 | 0.95 | 0.93 | 0.87 | 1.0 |

Achieves **SOTA** by ≥ 15%

Correct Reasoning — Localisation Ability
Accurate Description — Correct Label

Hallucinations — Incorrect Reasoning

### Industry Benchmarks (General VQA)

| Model | Modality | MSRVTT-QA | MSVD-QA | ActivityNet-QA |
|---|---|---|---|---|
| QueST | V | 34.6 | 34.6 | - |
| ClipBERT | V | 37.4 | - | - |
| JustAsk | V | 41.5 | 46.3 | 38.9 |
| GIT | V | 42.7 | 55.1 | - |
| MERLOT | V | 43.1 | - | 41.4 |
| Singularity | V | 43.5 | - | 43.1 |
| Clover | V | 43.9 | 51.9 | - |
| VideoChat | V | 45.0 | 56.3 | 26.5 |
| Video-ChatGPT | V | 49.3 | 64.9 | 35.2 |
| VALOR | V,A | 46.7 | 56.4 | 44.8 |
| FrozenBiLM | V | 47.0 | 54.8 | 43.2 |
| Valley | V | 45.7 | 65.4 | 42.9 |
| Video-LLaMA | V,A | 29.6 | 51.6 | 12.4 |
| PandaGPT | V,A | 25.5 | 42.1 | 14.5 |
| LLaVA-OneVision-7B | V | 49.8 | 51.7 | 56.6 |
| MacawLLM | V,A | 25.5 | 42.1 | 14.5 |
| MAESTRO (Ours) | V,A | 82.0 | 86.9 | 87.2 |

Achieves **SOTA** by 32.2%, 21.5%, 30.6%
Advances general multimodal understanding of VLMs

### MAESTRO-CLAP

| Model | Multiclass (Zero-shot) | | | | |
|---|---|---|---|---|---|
| | Acc | M-F1 | F1(O) | R(O) | P(O) |
| ZAC | 0.21 | 0.23 | 0.20 | 0.26 | 0.33 |
| LAION-CLAP | 0.32 | 0.28 | 0.30 | 0.34 | 0.31 |
| MAESTRO-CLAP (Ours) | 0.82 | 0.89 | 0.84 | 0.92 | 0.90 |

ZAC — LAION-CLAP — MAESTRO-CLAP (Ours)

Strong Zero-Shot Abilities
Suited for performance in the wild