# MART: Masked Affective RepresenTation Learning via Masked Temporal Distribution Distillation

Zhicheng Zhang  Pancheng Zhao  Eunil Park  Jufeng Yang

Nankai University
成均館大學校 SUNG KYUN KWAN UNIVERSITY

## Introduction

**Motivation:** Inspired by the psychology research & empirical theory, we verify that **degree of emotion** may vary in different segment, thus introducing sentiment complementary and emotion intrinsic on temporal segments
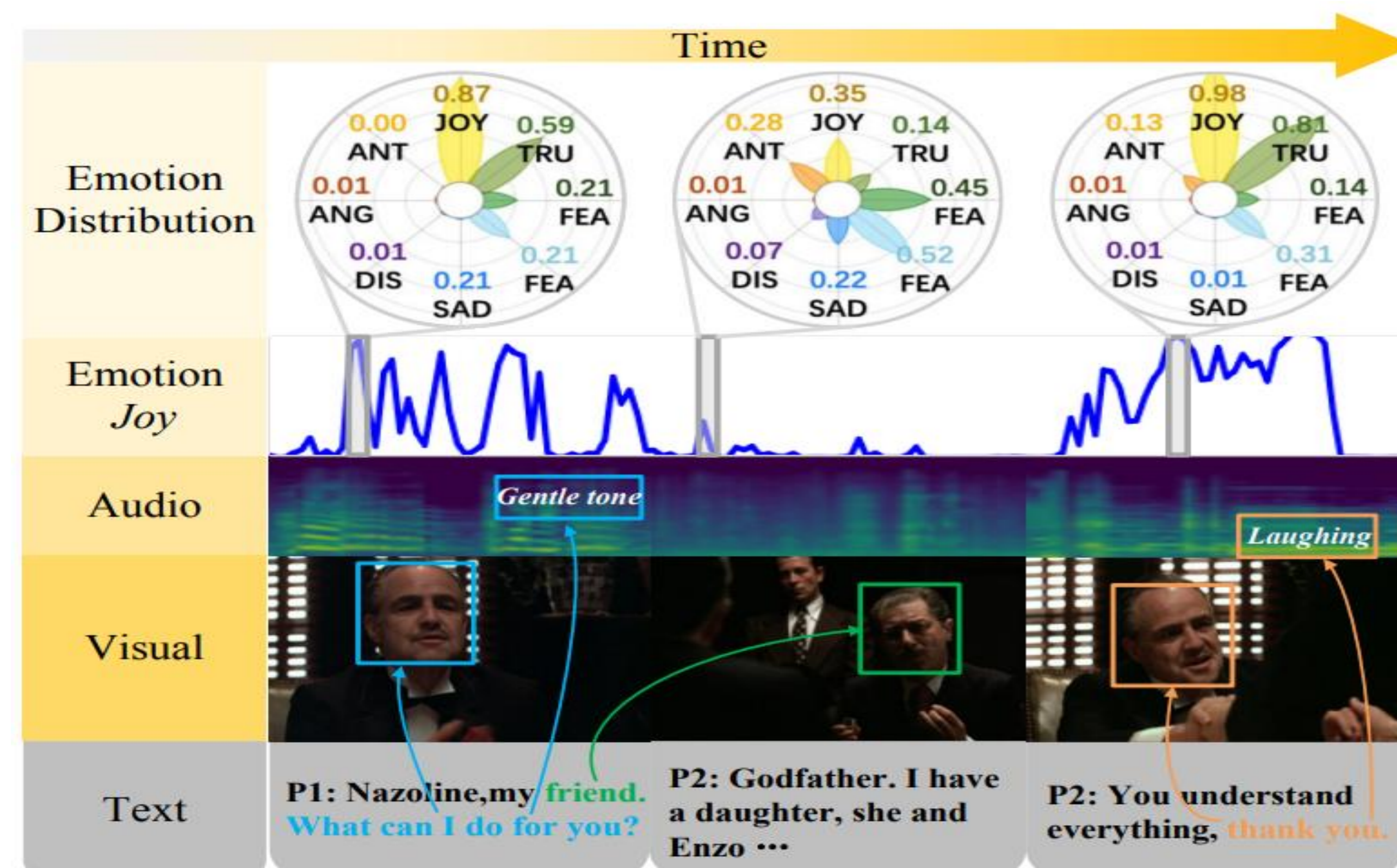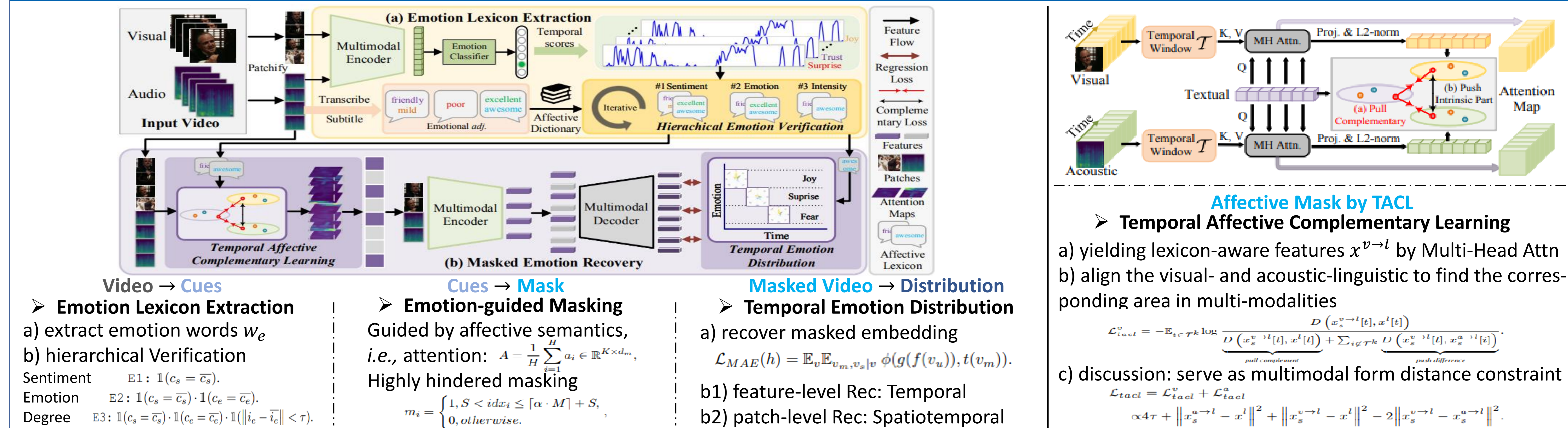


Fig.1. *Godfather* (1972) and its corresponding emotion. The colored texts indicate the emotional words given by affective lexicon and the rectangles show the related content towards emotional words.

### Contributions

- **Masked affective modeling** to exploit temporal affective cues among modalities for discriminative representation, which can be integrated into existing video emotion analysis methods as a **plug-in** module.
- Extensive experiments demonstrate effectiveness of MART **covering four downstream areas** on emotion classification and sentiment analysis.

## Methodology



**Video → Cues**

➤ **Emotion Lexicon Extraction**
a) extract emotion words $w_e$
b) hierarchical Verification

Sentiment   E1: $\mathbb{1}(c_s = \overline{c_s})$.
Emotion    E2: $\mathbb{1}(c_s = \overline{c_s}) \cdot \mathbb{1}(c_e = \overline{c_e})$.
Degree     E3: $\mathbb{1}(c_s = \overline{c_s}) \cdot \mathbb{1}(c_e = \overline{c_e}) \cdot \mathbb{1}(\|i_e - \overline{i_e}\| < \tau)$.

**Cues → Mask**

➤ **Emotion-guided Masking**
Guided by affective semantics, *i.e.,* attention: $A = \frac{1}{H}\sum_{i=1}^{H} a_i \in \mathbb{R}^{K \times d_m}$,
Highly hindered masking

$m_i = \begin{cases} 1, & S < idx_i \leq \lceil \alpha \cdot M \rceil + S, \\ 0, & otherwise. \end{cases}$

**Masked Video → Distribution**
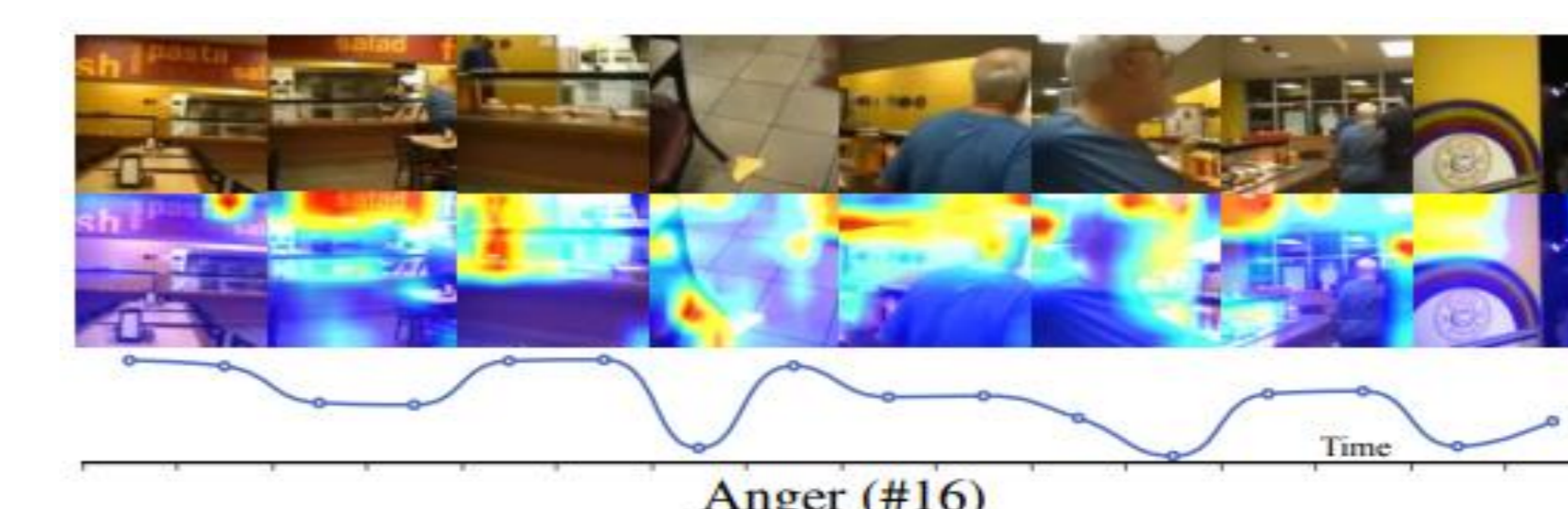
➤ **Temporal Emotion Distribution**
a) recover masked embedding
$\mathcal{L}_{MAE}(h) = \mathbb{E}_v \mathbb{E}_{v_m, v_s|v} \, \phi(g(f(v_u)), t(v_m))$.
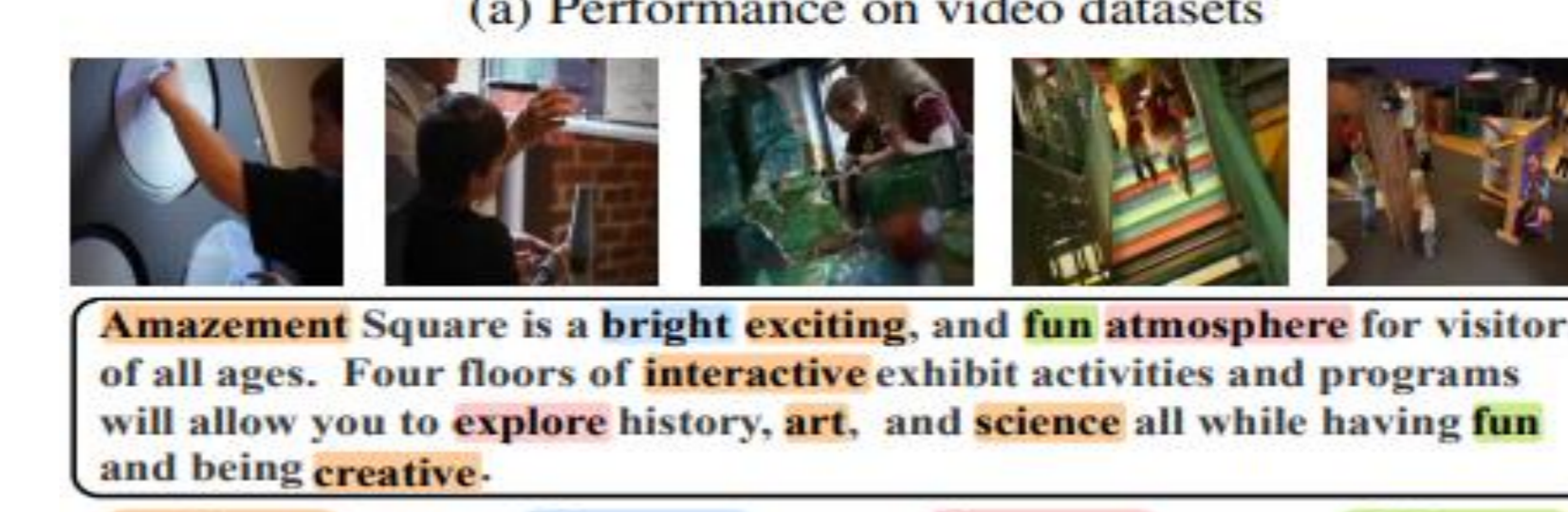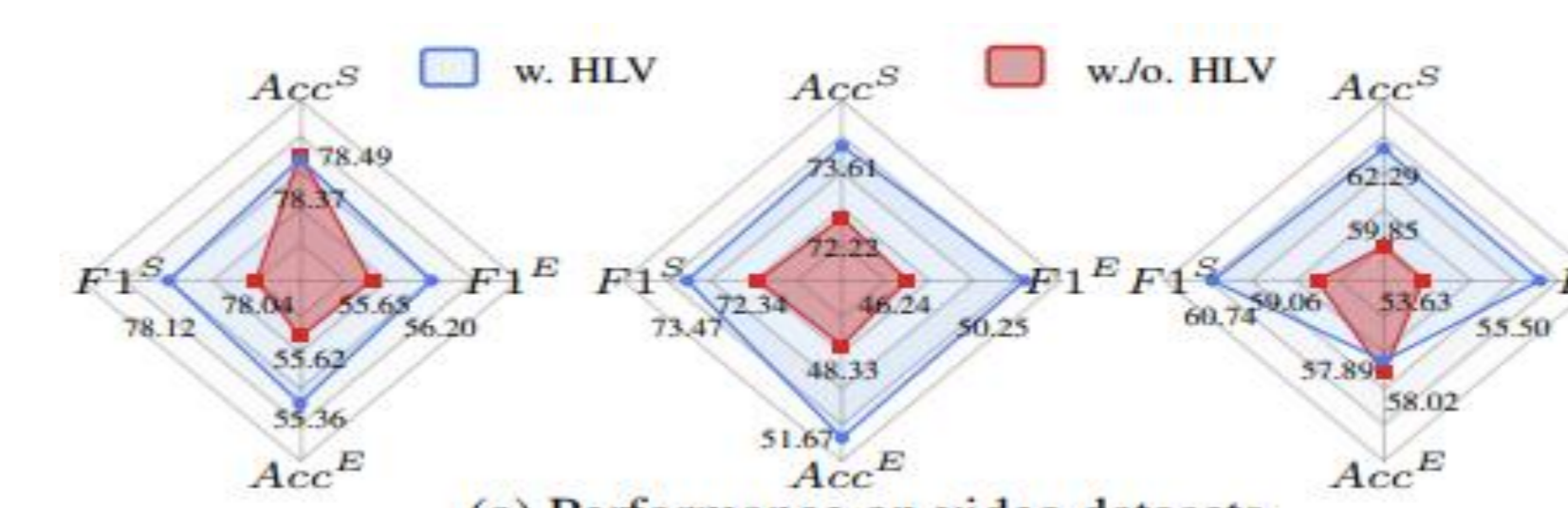b1) feature-level Rec: Temporal
b2) patch-level Rec: Spatiotemporal

**Affective Mask by TACL**

➤ **Temporal Affective Complementary Learning**
a) yielding lexicon-aware features $x^{v \to l}$ by Multi-Head Attn
b) align the visual- and acoustic-linguistic to find the corresponding area in multi-modalities

$\mathcal{L}_{tacl}^v = -\mathbb{E}_{t \in \mathcal{T}^k} \log \frac{D(x_a^{v \to l}[t], x^l[t])}{D(x_a^{v \to l}[t], x^l[t]) + \sum_{i \notin \mathcal{T}^k} D(x^{v \to l}[t], x_s^{a \to l}[i])}$

c) discussion: serve as multimodal form distance constraint
$\mathcal{L}_{tacl} = \mathcal{L}_{tacl}^v + \mathcal{L}_{tacl}^a$
$\propto 4\tau + \|x_a^{a \to l} - x^l\|^2 + \|x_s^{v \to l} - x^l\|^2 - 2\|x_s^{v \to l} - x_s^{a \to l}\|^2$.

## Performance



(a) Input video
(b) Frame masking [66]
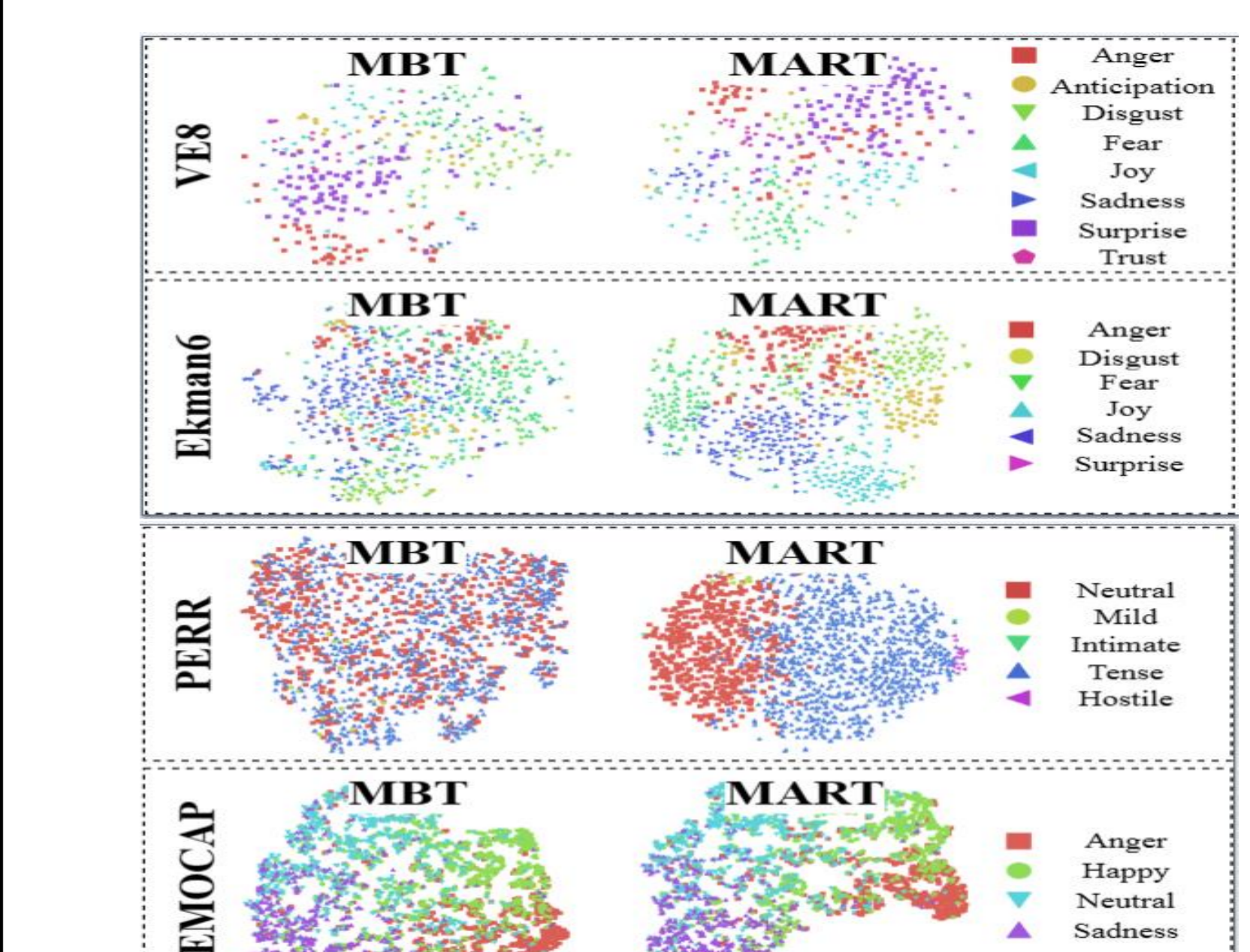(c) Random masking [22]
(d) Tube masking [62]
(e) Ours

✓ **MART** encourages video emotion analysis models to focus on affective area, unlike common mask strategy which randomly select masking area.

✓ **TACL** guides the recovery, where highly attentive areas in yellow boxes are reserved for providing affective cues to be recovered.



✓ **MART** recovers temporal distribution.

✓ **MART** shows tolerance to noisy text.



✓ **MART** extracts robust affective embedding across datasets.

**Contact**
MART
gloryzzc6@sina.com
Website  Project  Code