

# Interactive Multimodal Learning via Flat Gradient Modification

Qing-Yuan Jiang, Zhouyang Chi, and Yang Yang\*

IJCAI, 2025, Guangzhou,  
Nanjing University of Science and Technology

August 30, 2025



# 1 Introduction

## 2 Methodology

## 3 Experiments

## 4 Conclusion

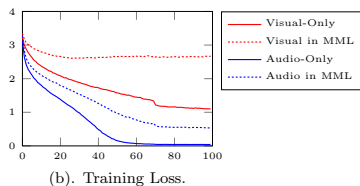
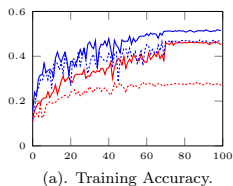
# Imbalance in Multimodal Learning

## Multimodal Learning

- **Goal:** fuse multimodal data to boost model performance.
- **Optimal Scenario:** maximize information extraction from multiple modalities for better performance.

## Modality Imbalance

- **Phenomenon:** MML **underperforms** single-modality models.
- **Strong-Weak Modality**  $\mapsto$  **Modality Imbalance**



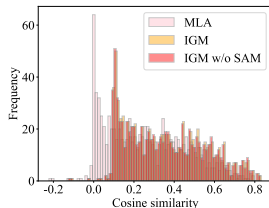
# Issue in Orthogonal Projection

**Alternating Learning:** Learning each modality one-by-one

## Poor Plasticity:

- The orthogonal projection suffers from **poor plasticity** problem, i.e., leading to **feasible gradient direction becomes narrow**.
- The poor plasticity problem results in **suboptimal solution**.

## Gradient Change:



① Introduction

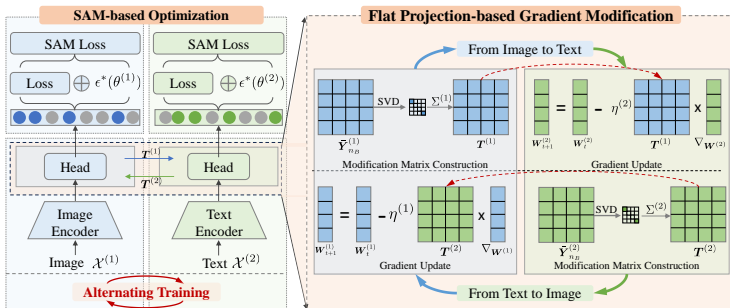
② Methodology

③ Experiments

④ Conclusion

# Overall Framework

- **Flat Projection Gradient Modification:** Project updating direction along **flat direction**.
- **SAM-based Optimization:** Smooth the objective function, enhancing the **flatness**.



# Flat Projection Gradient Modification

- **Find Flat Direction:**  $\mathbf{U}^v \mathbf{\Lambda}^v [\mathbf{V}^v]^\top = \text{svd}(\mathbf{Y}^v)$
- **Relationship of Flatness and Singular Value:** Small  $\lambda^v \mapsto$  Flat Area
- **Conduct Projection Matrix:**  
$$\Sigma^v = \exp\left(-\frac{\tau}{\lambda_{\max}^v - \lambda_{\min}^v} (\mathbf{\Lambda}^v - \lambda_{\min}^v \mathbf{I})\right)$$
- **Project during Updating:**

$$\mathbf{T}^v = \mathbf{U}^v \Sigma^v [\mathbf{V}^v]^\top$$
$$\omega_{t+1}^a = \omega_t^a - \eta \cdot \mathbf{T}^v \cdot \nabla_{\omega} \mathcal{L}$$

# SAM-based Optimization

- **Perturb the Loss:**  $\mathcal{L}^{\text{SAM}}(\omega) = \max_{\epsilon: \|\epsilon\| \leq \rho} \mathcal{L}(\omega + \epsilon).$
- **Optimal perturbation:**  $\epsilon^*(\omega) = \operatorname{argmax}_{\|\epsilon\| \leq \rho} \mathcal{L}(\omega + \epsilon).$
- **Gradient of SAM Loss:**  $\nabla_{\omega} \mathcal{L}^{\text{SAM}}(\omega) = \nabla_{\omega} \mathcal{L}(\omega)|_{\omega + \epsilon^*(\omega)}.$



# Overall Algorithm

**Updating rule:**  $\omega_{t+1}^a = \omega_t^a - \eta \cdot \mathbf{T}^\nu \cdot \nabla_{\omega^a} \mathcal{L}^{\text{SAM}}(\omega^a)$

---

**Algorithm 1** Algorithm for IGM
 

---

**Input:** Training set  $\mathcal{D}$  and labels  $\mathbf{Y}$ ;

**Output:** The learned parameters  $\{\theta^{(j)}\}_{j=1}^{(m)}$ ;

**INIT:** Initialize gradient modification matrix. Initialize

$\{\mathbf{T}^{(k)}\}_{j=1}^{(m)}$ ;  $\forall k \in \{1, \dots, m\}$ ,  $\mathbf{T}^{(k)} = \mathbf{I}$ ;

- 1: **for**  $i = 1 \rightarrow \text{Outer\_Iters}$  **do**
  - 2:   **for**  $j = 1 \rightarrow m$  **do** ▷ Main iteration.
  - 3:     **for**  $t = 1 \rightarrow \text{Inner\_Iters}$  **do**
  - 4:       Randomly construct a mini-batch  $\mathcal{X}_t^{(j)}$ .
  - 5:       Calculate loss  $L(\theta^{(j)})$  for data in  $\mathcal{X}_t^{(j)}$ .
  - 6:       Calculate  $\epsilon^*(\theta^{(j)})$  according to Eq. (7).
  - 7:       Calculate  $\nabla_{\theta^{(j)}} L^{\text{SAM}}$  according to Eq. (8).]
  - 8:       Calculate modality index:
  - 9:        $k = \text{mod}(j + m - 2, m) + 1$ .
  - 10:      Update  $\theta^{(j)}$ :  $\theta_{t+1}^{(j)} = \theta_t^{(j)} - \eta^{(j)} \mathbf{T}^{(k)} \nabla_{\theta^{(j)}} L^{\text{SAM}}$ .
  - 11:    **for**  $j = 1 \rightarrow n_B$  **do** ▷ Update  $\{\bar{\mathbf{Y}}_{n_B}^{(k)}\}$ .
  - 12:      Update cumulative variance according to Eq. (2).
  - 13:    Update  $\mathbf{T}^{(j)}$  according to Eq. (4). ▷ Update  $\mathbf{T}^{(j)}$ .
-

## ① Introduction

## ② Methodology

## ③ Experiments

## ④ Conclusion

# Settings

## Datasets:

- **CREMA-D:** 7,442 audio-video pairs with 6 categories.
- **KSounds:** 19,000 audio-video pairs with 31 action categories.
- **Twitter2015:** 5,338 image-text pairs with 3 categories.
- **Sarcasm:** 24,635 image-text pairs with 2 categories.
- **NVGesture:** 1,532 dynamic hand gestures with 3 modalities.

## Baselines:

- **Unimodal:** audio, video, image, text, RGB, OF, Depth.
- **Joint-Training:** MSES [ACPR'19], OGR-GB [CVPR'20], OGM [CVPR'22], DOMFN [MM'22], MLRSR [ACL'22], PMR [CVPR'23], AGM [ICCV'23], SMV [CVPR'24], **MMPareto** [ICML'24].
- **Alternating Learning:** ReconBoost [ICML'24] and **MLA** [CVPR'24].

# Comparison with SOTAs

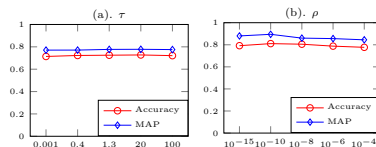
Method	CREMA-D		KSounds		Twitter2015		Sarcasm		NVGesture	
	Acc.	MAP	Acc.	MAP	Acc.	Mac-F1	Acc.	Mac-F1	Acc.	Mac-F1
Unimodal-1	.6317	.6861	.5312	.5669	.7367	.6849	.8136	.8065	.7822	.7833
Unimodal-2	.4583	.5879	.5462	.5837	.5863	.4333	.7181	.7073	.7863	.7865
Unimodal-3	-	-	-	-	-	-	-	-	.8154	.8183
OGR-GB	.6465	.6854 <sup>†</sup>	.6710	.7139	.7435	.6869	.8335	.8271	.8299	.8305
OGM	.6694	.7173	.6606	.7144	<b>.7492</b>	.6874	.8323	.8266	-	-
DOMFN	.6734	.7372	.6625	.7244	.7445	.6857	.8356	.8262	-	-
MSES	.6156 <sup>†</sup>	.6683 <sup>†</sup>	.6471	.7063	.7184 <sup>†</sup>	.6655 <sup>†</sup>	.8418	.8360	.8112 <sup>†</sup>	.8147 <sup>†</sup>
PMR	.6659	.7030	.6656	.7193	.7425	.6860	.8360	.8249	-	-
AGM	.6707	.7358	.6602	.7252	.7483	.6911	.8402	.8344	.8278	.8282
MSLR	.6546	.7138	.6591	.7196	.7252 <sup>†</sup>	.6439 <sup>†</sup>	.8423	.8369	.8286	.8292
ReconBoost	.7484	.8124	.7085	.7424	.7442	.6834	.8437	.8317	.8413	.8632
SMV	.7872	.8417	.6900	.7426	.7428	.6817	.8418	.8368	.8352	.8341
MMPareto	.7487	.8535	.7000	<u>.7850</u>	.7358	.6729	.8348	.8284	.8382	.8424
MLA	.7943	.8572	.7004	<u>.7413</u>	.7352 <sup>†</sup>	.6713 <sup>†</sup>	.8426	.8348	.8373	.8387
IGM w/o SAM	<u>.8026</u>	<u>.8830</u>	<u>.7159</u>	.7623	.7395	<u>.6912</u>	<u>.8455</u>	<u>.8390</u>	<u>.8487</u>	<u>.8634</u>
IGM	<b>.8105</b>	<b>.8948</b>	<b>.7403</b>	<b>.7855</b>	<u>.7489</u>	<b>.6917</b>	<b>.8468</b>	<b>.8392</b>	<b>.8693</b>	<b>.8703</b>

# Further Analysis

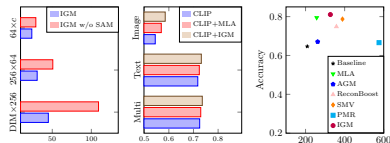
## Ablation Study:

SAM	GM	Audio	Video	Multi
✗	✗	45.83%	63.17%	64.52%
✓	✗	58.60%	64.79%	73.42%
✗	✓	60.13%	65.06%	80.26%
✓	✓	<b>61.16%</b>	<b>67.82%</b>	<b>81.05%</b>

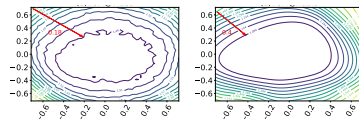
## Params. Sensitivity:



## Singular Values, Pretrained Model, and Training Time:



## Loss Landscape Visualization:



## 1 Introduction

## 2 Methodology

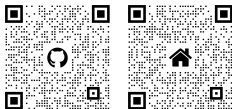
## 3 Experiments

## 4 Conclusion

# Conclusion

- A **flat-projection gradient modification** based MML method is proposed to address the **poor plasticity** issue.
- SAM optimization algorithm is integrated in the loss function to **smooth the objective function**.
- Comprehensive experiments are conducted to demonstrate the superiority and effectiveness of IGM.

## Contact Us:



# Thanks