

Balanced Multimodal Representation Learning

Qing-Yuan Jiang

NJUST-KMG group,
The School of Computer Science and Engineering, NJUST

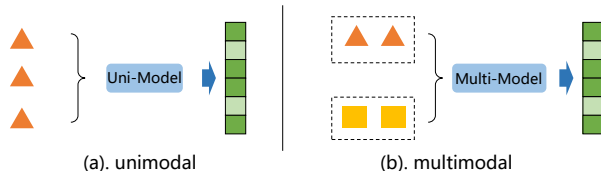
March 26, 2026



Multimodal Representation Learning

Multimodal Learning:

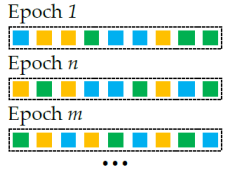
- **Goal:** fuse multimodal data to boost model performance.
- **Optimal Scenario:** maximize information extraction from multiple modalities for **better performance**.
- Taking multimodal data as inputs, and then predicting.



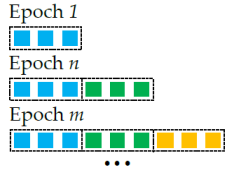
Data-Level Modality Imbalance

- The training of multimodal models is sensitive to **task difficulty** across **different learning stages**.

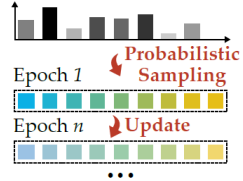
Balanced (Easy) Imbalanced (Hard)



(a1). Random sequences.



(a2). Fixed sequences.



(a3). Dynamic sequences.

Balance-Aware Sequence Sampling

Balance Score:

- Combining the correlation and information criterion:

$$s(\mathbf{x}_i) = \frac{\text{sim}(\mathbf{x}_i^a, \mathbf{x}_i^v) - \min(\mathcal{S})}{\max(\mathcal{S}) - \min(\mathcal{S})} - \frac{\ell_{\text{total}}(\mathbf{x}_i^a, \mathbf{x}_i^v, \mathbf{y}_i) - \min(\mathcal{L})}{\max(\mathcal{L}) - \min(\mathcal{L})}$$

Training scheduler:

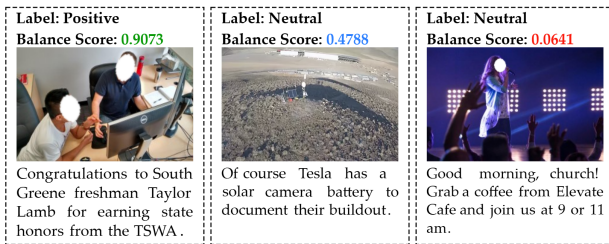
- Heuristic scheduler:** splitting training set as multiple blocks based on pace function. Then sample batch from each block.
- Learning-based scheduler:** sampling batch based on sampling probability:

$$p(\mathbf{x}_i) = \frac{e^{\hat{s}^{k+1}(\mathbf{x}_i)}}{\sum_{j=1}^n e^{\hat{s}^{k+1}(\mathbf{x}_j)}}.$$

Further Analysis for BSS

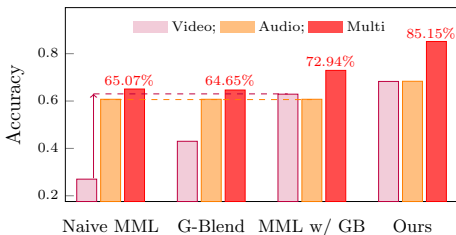
- Ablation study and visualization:

Criterion		ACC (%) / MAP (%)		
PreSim	Loss	Audio	Video	Multi
\times	\times	49.37/51.07	54.03/57.48	70.44/76.62
\times	\checkmark	52.11/54.40	54.23/57.91	72.44/79.41
\checkmark	\times	52.38/54.32	54.93/58.52	73.25/78.98
\checkmark	\checkmark	52.73/54.43	54.74/58.46	73.95/79.43



Model-Level Modality Imbalance

- Capability discrepancies across **heterogeneous models** play a crucial role in causing modality imbalance.



How to balance the *model capability* for multimodal learning?

- Enhancing the learning capabilities of weak modality:** classifier boosting.

Adaptive Unimodal Gradient Boosting (cont'd)

Since we utilize a shared encoder, we have to ensure that **all classifiers** are updated simultaneously:

$$\epsilon_o(\mathbf{x}_i^o, \mathbf{y}_i, t) = \ell \left(\mathbf{p}_{it}^o + \sum_{j=1}^{t-1} \mathbf{p}_{ij}^o, \mathbf{y}_i \right) = \ell \left(\sum_{j=1}^t \mathbf{p}_{ij}^o, \mathbf{y}_i \right).$$

Furthermore, we have to ensure that the $t - 1$ classifiers are well-trained:

$$\epsilon_p(\mathbf{x}_i^o, \mathbf{y}_i, t) = \ell \left(\sum_{j=1}^{t-1} \mathbf{p}_{ij}^o, \mathbf{y}_i \right).$$

In summary, we have:

$$L(\mathbf{x}_i^o, \mathbf{y}_i, t) = \epsilon(\mathbf{x}_i^o, \mathbf{y}_i, t) + \epsilon_o(\mathbf{x}_i^o, \mathbf{y}_i, t) + \epsilon_p(\mathbf{x}_i^o, \mathbf{y}_i, t).$$

Experiments for AUG

- Main results compared with SOTA:

Method	CREMA-D		KSounds		VGGSound		Twitter2015		Sarcasm		NVGesture	
	Acc.	MAP	Acc.	MAP	Acc.	MAP	Acc.	F1	Acc.	F1	Acc.	F1
Unimodal-1	.4583	.5879	.5412	.5669	.4655	.4701	.5863	.4333	.7181	.7073	.7822	.7833
Unimodal-2	.6317	.6861	.5562	.5837	.3494	.3478	.7367	.6849	.8136	.8056	.7863	.7865
Unimodal-3	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	.8154	.8183
Concat	.6361	.6841 [†]	.6455	.7130	.5116	.5352	.7011 [†]	.6386	.8286	.8240	.8237	.8270
Affine	.6626	.7193	.6424	.6931	.5001	.5155	.7203 [†]	.5992 [†]	.8240	.8188	.8278	.8281
ML-LSTM	.6290 [†]	.6473 [†]	.6394	.6902	.4966	.5139	.7068 [†]	.6564 [†]	.8277	.8205	.8320	.8330
Sum	.6344	.6908	.6490	.7103	.5136	.5338	.7300 [†]	.6661 [†]	.8294	.8247	.8050 [†]	.8067 [†]
Weight	.6653	.7134	.6533	.7110	.5144	.5300	.7242 [†]	.6516 [†]	.8265	.8219	.7842 [†]	.7939 [†]
MSES	.6546	.7138	.6591	.7196	.4891	.5429	.7252 [†]	.6439 [†]	.8423	.8369	.8112 [†]	.8147 [†]
G-blend	.6465	.7392	.6722	.7274	.5086	.5555	.7309 [†]	.6799 [†]	.8286	.8215	.8299	.8305
MSLR	.6868	.7412	.6756	.7282	.4987	.5415	.7232 [†]	.6382 [†]	.8439	.8378	.8237	.8284
OGM	.6612	.7372	.6582	.7159	.4829	.4978	.7058 [†]	.6435 [†]	.8360	.8293	—	—
PMR	.6659	.7058	.6675	.7274	.4647	.4866	.7357 [†]	.6636 [†]	.8310	.8256	—	—
AGM	.6733	.7807	.6791	.7388	.4711	.5198	.7261 [†]	.6502 [†]	.8306	.8293	.8279	.8284
MMParato	.7487	.8535	.7000	.7850	.5125	.5473	.7358 [†]	.6729 [†]	.8348	.8284	.8382	.8424
SMV	.7872	.8417	.6900	.7426	.5031	.5362	.7428	.6817 [†]	.8418	.8368	.8352	.8341
MLA	.7943	.8572	.7004	<u>.7945</u>	.5165	.5473	.7352 [†]	.6713 [†]	.8426	.8348	.8340	.8372
DI-MML	.8158	.8592	.7203	.7426	.5173	.5479	.7248 [†]	.6686 [†]	.8411	.8315	—	—
LFM	<u>.8362</u>	<u>.9006</u>	<u>.7253</u>	.7897	<u>.5274</u>	<u>.5694</u>	<u>.7501</u>	.7057	<u>.8497</u>	<u>.8457</u>	<u>.8436</u>	<u>.8468</u>
ReconBoost	.7557	.8140	.6855	.7662	.5097	.5387	.7442	.6832 [†]	.8437	.8317	.8386	.8434
Ours	.8515	.9103	.7263	.7901	.5301	.5826	.7512	<u>.6962</u>	.8510	.8458	.8501	.8533

Theoretical Analysis

Theorem: Convergence, Informal

Under some assumptions for the stochastic gradient $\nabla \ell(\mathbf{w}^{(k)}) \odot \mathbf{m}^{(k)}(t)$, we have:

$$\frac{1}{T} \sum_{t=1}^T \|\nabla \ell(\mathbf{w}^{(k)})\|^2 \leq \mathcal{O} \left(\frac{1 + (1 + \nu)^2}{\sqrt{T}(1 + \nu)(1 - \delta^2)} \right)$$

where $\delta \in (0, 1)$ and $\nu \geq 0$ are two constants.

Ablation Study for AMSS

- Ablation study:

Methods	Kinetics-Sound		Sarcasm-Detection	
	ACC	mAP	ACC	Mac-F1
Baseline	64.55	71.30	82.86	82.40
w/ Classifier Mask	65.37	70.87	83.40	83.05
w/ Backbone Mask	66.68	72.23	83.81	83.17
AMSS	68.96	74.89	84.14	83.69
w/ Classifier Mask+	66.80	73.80	83.44	82.93
w/ Backbone Mask+	69.15	76.13	83.77	83.10
AMSS+	72.25	79.13	84.35	83.77

Multimodal Learning with Fusion and Alignment (cont'd)

Definition of pair similarity:

$$s(\mathbf{x}_i^a, \mathbf{x}_k^v) = \frac{[\mathbf{z}_i^a]^\top \mathbf{z}_k^v}{\|\mathbf{z}_i^a\|_2 \|\mathbf{z}_k^v\|_2}.$$

Loss function for contrastive learning:

$$L_{MM}(\mathbf{X}) = -\frac{1}{2n_b} \sum_i^{n_b} \left[\log \left(\frac{e^{s(\mathbf{x}_i^a, \mathbf{x}_i^v)/\tau}}{\sum_k e^{s(\mathbf{x}_i^a, \mathbf{x}_k^v)/\tau}} \right) + \log \left(\frac{e^{s(\mathbf{x}_i^a, \mathbf{x}_i^v)/\tau}}{\sum_k e^{s(\mathbf{x}_k^a, \mathbf{x}_i^v)/\tau}} \right) \right].$$

Total loss:

$$L_{Total} = (1 - \alpha)L_{CLS}(\mathbf{X}, \mathbf{Y}) + \alpha L_{MM}(\mathbf{X}).$$

The LFA Algorithm

Algorithm 3 Multimodal Learning with Label Fitting and Alignment

Input: Training set \mathcal{X} , labels \mathcal{Y} , method.

Output: Learned parameters $\{\theta\}$ of all models.

- 1: Initialize parameters θ , parameter α , maximum iterations T , learning rate η_α .
 - 2: **for** $t = 1 \mapsto T$ **do**
 - 3: **for** $i = 1 \mapsto \text{Inner_Iters}$ **do**
 - 4: Calculate total loss L_{Total} by forward phase.
 - 5: Update parameters θ according to its gradient.
 - 6: **if** method == 'learning-based' **then**
 - 7: Calculate gradient approximation: $\nabla L_{\text{CLS}}(\theta(\alpha)) = -\nabla_{\alpha, \theta}^2 L_{\text{Total}} \cdot [\nabla_{\theta, \theta}^2 L_{\text{Total}}]^{-1} \cdot \nabla_{\theta} L_{\text{CLS}}(\mathbf{X}, \mathbf{Y})$.
 - 8: Update α according to: $\alpha = \alpha - \eta_\alpha \nabla L_{\text{CLS}}(\theta(\alpha))$.
 - 9: Clip α into $[0, 1]$: $\alpha := \max(0, \min(1, \alpha))$.
 - 10: **else if** method == 'heuristic' **then**
 - 11: Update α according to: $\alpha = 1 - e^{-1/t}$.
-

Experiments for LFA

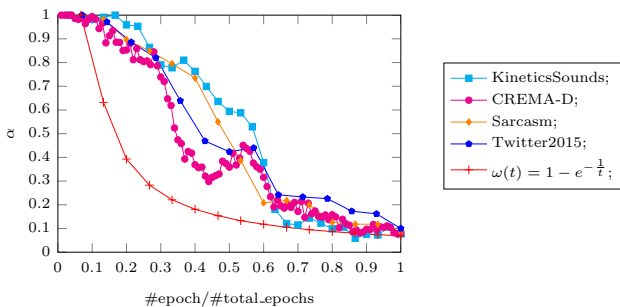
Main comparison with SOTA:

Method	KSounds		CREMA-D		Sarcasm		Twitter2015		NVGesture	
	ACC	MAP	ACC	MAP	ACC	F1	ACC	F1	ACC	F1
Unimodal-1	54.12%	56.69%	63.17%	68.61%	81.36%	80.65%	73.67%	68.49%	78.22%	78.33%
Unimodal-2	55.62%	58.37%	45.83%	58.79%	71.81%	70.73%	58.63%	43.33%	78.63%	78.65%
Unimodal-3	—	—	—	—	—	—	—	—	81.54%	81.83%
Concat	64.55%	71.31%	63.31%	68.41%	82.86%	82.43%	70.11%	63.86%	81.33%	81.47%
Affine	64.24%	69.31%	66.26%	71.93%	82.47%	81.88%	72.03%	59.92%	82.78%	82.81%
Channel	63.51%	68.66%	66.13%	71.75%	—	—	—	—	81.54%	81.57%
ML-LSTM	63.84%	69.02%	62.94%	64.73%	82.05%	70.73%	70.68%	65.64%	83.20%	83.30%
Sum	64.97%	71.03%	63.44%	69.08%	82.94%	82.47%	73.12%	66.61%	82.99%	83.05%
Weight	65.33%	71.33%	66.53%	73.26%	82.65%	82.19%	72.42%	65.16%	83.42%	83.57%
ETMC	65.67%	71.19%	65.86%	71.34%	83.69%	83.23%	73.96%	67.39%	83.61%	83.69%
MSES	64.71%	72.52%	61.56%	66.83%	84.18%	83.60%	71.84%	66.55%	81.12%	81.47%
G-Blend	67.12%	71.39%	64.65%	68.54%	83.35%	82.71%	74.35%	68.69%	82.99%	83.05%
OGM	66.06%	71.44%	66.94%	71.73%	83.23%	82.66%	74.92%	68.74%	—	—
Greedy	66.52%	72.81%	66.64%	72.64%	—	—	—	—	82.74%	82.69%
DOMFN	66.25%	72.44%	67.34%	73.72%	83.56%	82.62%	74.45%	68.57%	—	—
MSLR	65.91%	71.96%	65.46%	71.38%	84.23%	83.69%	72.52%	64.39%	82.86%	82.92%
PMR	66.56%	71.93%	66.59%	70.36%	83.61%	82.49%	74.25%	68.62%	—	—
AGM	66.02%	72.52%	67.07%	73.58%	84.28%	83.44%	74.83%	69.11%	82.78%	82.82%
MLA	70.04%	74.13%	79.43%	85.72%	84.26%	83.48%	73.52%	67.13%	83.73%	83.87%
ReconBoost	70.85%	74.24%	74.84%	81.24%	84.37%	83.17%	74.42%	68.34%	84.13%	86.32%
MMPareto	70.00%	78.50%	74.87%	75.15%	83.48%	82.84%	73.58%	67.29%	83.82%	84.24%
LFA-H	69.05%	72.97%	72.15%	80.45%	84.12%	83.98%	73.87%	69.17%	83.24%	83.87%
LFA-LB	72.53%	78.38%	83.62%	90.06%	84.97%	84.57%	75.01%	70.57%	84.36%	84.68%

Further Analysis for LFA

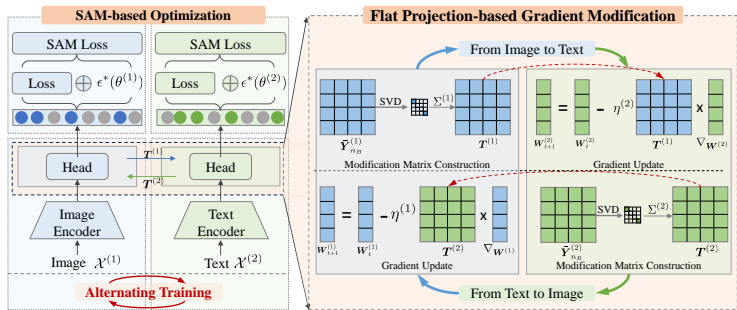
- Learning strategy and visualization:

Dataset	Modality	Constant			Stepwise				Dynamic	
		0	0.5	1	$h(0)$	$h(1)$	$h(0.05)$	$h(0.95)$	Ours-H	Ours-LB
KSounds	Multiple	64.55%	64.70%	28.67%	65.17%	66.92%	66.01%	67.41%	69.32%	72.89%
	Audio	49.17%	46.30%	34.11%	51.12%	52.34%	52.21%	53.41%	53.89%	54.32%
	Video	24.64%	44.02%	28.41%	41.21%	41.45%	42.31%	46.72%	49.18%	54.17%
CREMA-D	Multiple	63.31%	70.45%	26.49%	66.45%	70.24%	69.11%	71.45%	72.39%	84.11%
	Audio	55.65%	60.17%	33.15%	56.19%	57.38%	58.09%	60.18%	61.89%	65.13%
	Video	18.68%	42.54%	20.42%	45.14%	49.97%	46.41%	55.32%	57.14%	64.89%



Interactive Multimodal Learning via Gradient Modification

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



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

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

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:** $\nabla_{\omega} \mathcal{L}^{\text{SAM}}(\omega) = \nabla_{\omega} \mathcal{L}(\omega)|_{\omega + \epsilon^*(\omega)}$.

The IGM Algorithm

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

Algorithm 4 Multimodal Learning with Label Fitting and Alignment

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

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

- 1: Initialize gradient modification matrix, $\{\mathbf{T}^{(k)}\}_{k=1}^{(m)}$: $\forall k \in \{1, \dots, m\}, \mathbf{T}^{(k)} = \mathbf{I}$;
 - 2: **for** $i = 1 \rightarrow \text{Out_Iters}$ **do**
 - 3: **for** $j = 1 \rightarrow m$ **do**
 - 4: **for** $t = 1 \rightarrow \text{Inner_Iters}$ **do**
 - 5: Randomly construct a mini-batch $\mathcal{X}_t^{(j)}$.
 - 6: Calculate loss $L(\theta^{(j)})$ for data in $\mathcal{X}_t^{(j)}$.
 - 7: Calculate $\epsilon^*(\theta^{(j)})$ and $\nabla_{\theta^{(j)}} L^{\text{SAM}}$.
 - 8: Calculate modality index: $k = \text{mod}(j + m - 2, m) + 1$.
 - 9: Update $\theta^{(j)}$: $\theta_{t+1}^{(j)} = \theta_t^{(j)} - \eta^{(j)} \mathbf{T}^{(k)} \nabla_{\theta^{(j)}} L^{\text{SAM}}$.
 - 10: **for** $j = 1 \rightarrow n_B$ **do**
 - 11: Update cumulative variance.
 - 12: Update $\mathbf{T}^{(j)}$.
-

Experiments for IGM

- Main comparison with SOTA:

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 [†]	.6710	.7139	.7435	.6869	.8335	.8271	.8299	.8305
OGM	.6694	.7173	.6606	.7144	.7492	.6874	.8323	.8266	-	-
DOMFN	.6734	.7372	.6625	.7244	.7445	.6857	.8356	.8262	-	-
MSES	.6156 [†]	.6683 [†]	.6471	.7063	.7184 [†]	.6655 [†]	.8418	.8360	.8112 [†]	.8147 [†]
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 [†]	.6439 [†]	.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	.7413	.7352 [†]	.6713 [†]	.8426	.8348	.8373	.8387
IGM w/o SAM	<u>.8026</u>	<u>.8830</u>	<u>.7159</u>	<u>.7623</u>	<u>.7395</u>	<u>.6912</u>	<u>.8455</u>	<u>.8390</u>	<u>.8487</u>	<u>.8634</u>
IGM	.8105	.8948	.7403	.7855	<u>.7489</u>	.6917	.8468	.8392	.8693	.8703

Further Analysis of IGM

- Ablation study.

Ablation study on CREMA-D dataset.

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

Interactive enhancement analysis.

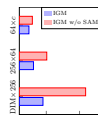
Method	Initial	Out_Iters=1		Out_Iters=2	
		Audio	Video	Audio	Video
w/o α -GM	.0325	<u>.5312</u>	.6803	.7231	.7482
w/o ν -GM	.0325	<u>.5312</u>	<u>.7023</u>	.7472	.7646
IGM	.0325	.5312	.7023	.7557	.8105

- Impact of different scope of GM, singular values, the pretrained models, and training time.

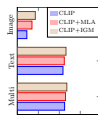
Results with different scope of GM.

Scope of GM	Accuracy	MAP
100%	75.34%	81.23%
50%	78.97%	85.58%
30%	82.97%	90.15%
1.3% (Classification head)	81.05%	89.48%
0% (w/o GM)	73.42%	81.77%

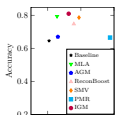
Impact of singular values, CLIP, and training time.



(a). Singular values.

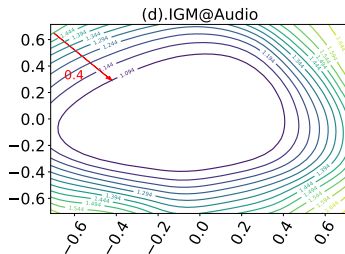
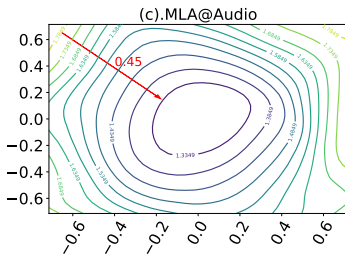
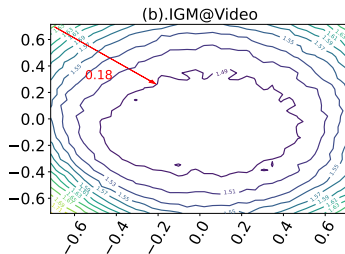
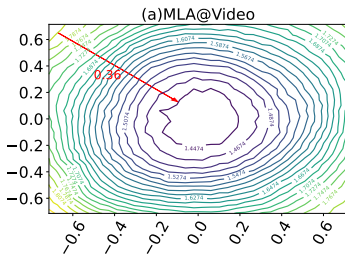


(b). Pretrained Model.



(c). Training time.

Loss Landscape Visualization



Instantaneous Probe-and-Rebalance Multimodal Learning

Existing Methods: Post-hoc Balance Learning After
“Imbalance”.

Gradient Modulation-based methods:

- Adjust updating direction during backward.
- G-blend [2], OGM [3], CGGM [4].

Interaction-based methods:

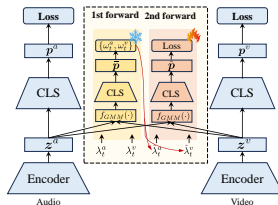
- enhance the interaction using features, logits after obtained.
- MLA [5], DI-MML [6], OLM [7]

Instantaneous Probe-and-Rebalance Multimodal Learning

IPRM: A Probe-then-Rebalance Strategy.

Two-pass forward strategy:

- Probes the strength of modality imbalance firstly.
- Facilitates multimodal learning under the balanced status secondly.



Instantaneous Probe-and-Rebalance Multimodal Learning

The 1st Forward Probing Phase:

- **Feature Extraction:** $\forall o \in \{a, v\}, \mathbf{z}_i^o = g_o(\mathbf{x}_i^o; \Theta_o)$.

- **Fusion with GMM:**

$$\bar{\mathbf{z}}_i = f_{GMM}(\bar{\mathbf{z}}_i^a, \bar{\mathbf{z}}_i^v, \lambda_t^a), \quad \bar{\mathbf{p}}_i = \text{softmax}(h(\bar{\mathbf{z}}_i)).$$

- **Evaluate Balancing Instantaneous:**

$$\omega_t^a \triangleq \frac{\mathcal{D}_{\text{KL}}(\mathcal{P}^v | \mathcal{P}; \mathcal{T}_t)}{\mathcal{D}_{\text{KL}}(\mathcal{P}^a | \mathcal{P}; \mathcal{T}_t) + \mathcal{D}_{\text{KL}}(\mathcal{P}^v | \mathcal{P}; \mathcal{T}_t)}$$

The 2nd Forward Rebalancing Phase:

- **Reassign Wight:** $\forall o \in \{a, v\}, \hat{\lambda}_t^o = \omega_t^o$.

- **Fusion with Balanced Status:**

$$\hat{\mathbf{z}}_i = f_{GMM}(\bar{\mathbf{z}}_i^a, \bar{\mathbf{z}}_i^v, \hat{\lambda}_t^a), \quad \hat{\mathbf{p}}_i = \text{softmax}(h(\hat{\mathbf{z}}_i)).$$

The IPRM Algorithm

Algorithm 5 IPRM Algorithm

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

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

- 1: Initialize parameters $\{\Theta_a, \Theta_v, \Phi, \Phi_a, \Phi_v\}$, maximum iterations M_t , learning rate η_α , modality weight $\omega_0^a = \omega_0^v = 0.5$, λ_1^a and λ_1^v ;
 - 2: **for** $t = 1 \mapsto M_t$ **do**
 - 3: Sample a mini-batch data samples \mathcal{T}_t .
 - 4: Calculate features based on $g_a(\cdot)$ and $g_v(\cdot)$.
 - 5: Calculate fused feature \bar{z} the by the first forward phase.
 - 6: Calculate the instantaneous strength score.
 - 7: Calculate the balanced weight.
 - 8: Calculate the prediction \bar{p} by the second forward phase.
 - 9: Calculate the gradients based on backward phase.
 - 10: Update the network parameters based on SGD.
 - 11: Update λ_t^a and λ_t^v .
-

Experiments for IPRM

- Main comparison with SOTA:

Dataset	Metric	Unimodal			Naive Fusion			IPRM
		A/A/R/A/I	V/V/O/V/T	D/T	Concat	Sum	Weight	
CREMA-D	Accuracy	63.17%	45.83%	N/A	63.61%	63.44%	<u>66.53%</u>	84.27%
	MAP	68.61%	58.79%	N/A	68.41% [†]	69.08%	71.34%	90.66%
KSounds	Accuracy	54.12%	55.62%	N/A	64.55%	64.90%	<u>65.33%</u>	74.37%
	MAP	56.69%	58.37%	N/A	<u>71.30%</u>	71.03%	71.10%	80.63%
NVGesture	Accuracy	78.22%	78.63%	81.54%	<u>82.37%</u>	80.50% [†]	78.42% [†]	85.89%
	Macro-F1	78.33%	78.65%	81.83%	<u>82.70%</u>	80.67% [†]	79.39% [†]	86.34%
IEMOCAP	Accuracy	58.45%	30.71%	70.55%	75.97%	<u>76.06%</u>	69.29% [†]	80.22%
	Macro-F1	58.29%	11.75%	69.93%	75.88%	<u>76.03%</u>	68.91% [†]	80.63%
Sarcasm	Accuracy	71.81%	81.36%	N/A	82.86%	<u>82.94%</u>	82.65%	85.14%
	Macro-F1	70.73%	80.56%	N/A	82.40%	<u>82.47%</u>	82.19%	84.41%

Experiments for IPRM (cont'd)

- Main comparison with SOTA:

Dataset	Metric	OGR-GB	MSLR	OGM	AGM	MMPareto	ReconBoost	MLA	LFM	IPRM
CREMA-D	Accuracy	64.65%	68.68%	66.12%	67.33%	74.87%	75.57%	79.43%	83.62%	84.27%
	MAP	73.92%	74.12%	73.72%	78.07%	85.35%	81.40%	85.72%	90.06%	90.66%
KSounds	Accuracy	67.22%	67.56%	65.82%	67.91%	70.00%	68.55%	70.04%	72.53%	74.37%
	MAP	72.74%	72.82%	71.59%	73.88%	78.50%	76.62%	79.45%	78.97%	80.63%
NVGesture	Accuracy	82.99%	82.37%	N/A	82.79%	83.82%	83.86%	83.40%	84.36%	85.89%
	Macro-F1	83.05%	82.84%	N/A	82.84%	84.24%	84.34%	83.72%	84.68%	86.34%
IEMOCAP	Accuracy	70.10%	76.69%	N/A	77.51%	77.69%	76.87%	79.31%	78.41%	80.22%
	Macro-F1	69.90%	76.77%	N/A	77.29%	77.89%	77.08%	79.73%	78.51%	80.63%
Sarcasm	Accuracy	82.86%	84.39%	83.60%	83.06%	83.48%	84.37%	84.26%	84.97%	85.14%
	Macro-F1	82.15%	83.78%	82.93%	82.93%	82.84%	83.17%	83.48%	84.57%	84.41%

Further Analysis

- Ablation Study:

Dataset	w/ L-Mixup	w/o EMA	One-Pass	IPRM
CREMA-D	75.53%	83.06%	83.47%	84.27%
KSounds	71.94%	73.91%	73.64%	74.37%
NVGesture	84.85%	85.27%	84.44%	85.89%
IEMOCAP	75.79%	78.05%	77.60%	80.22%
Sarcasm	84.52%	84.81%	84.10%	85.14%

Summary

Data-Level Modality Imbalance:

- Zhi-Hao Guan, **Qing-Yuan Jiang***, Yang Yang. Balance-aware Sequence Sampling Makes Multimodal Learning Better. *IJCAI*, 2025.

Model-Level Modality Imbalance:

- **Qing-Yuan Jiang**, Longfei Huang, Yang Yang. Rethinking Multimodal Learning from the Perspective of Mitigating Classification Ability Disproportion. *NeurIPS*, 2025 (Oral, 1.46%).
- Yang Yang, Hongpeng Pan, **Qing-Yuan Jiang**, Yi Xu, and Jinhui Tang. Learning to Rebalance Multi-Modal Optimization by Adaptively Masking Subnetworks. *TPAMI*, 2025.

Learning-Level Modality Imbalance:

- Yang Yang, Fengqiang Wan, **Qing-Yuan Jiang***, Yi Xu. Facilitating Multimodal Classification via Dynamically Learning Modality Gap. *NeurIPS*, 2024.
- **Qing-Yuan Jiang**, Zhouyang Chi, Yang Yang. Interactive Multimodal Learning via Flat Gradient Modification. *IJCAI*, 2025.
- Yang Yang, Xixian Wu, **Qing-Yuan Jiang***. Towards Equilibrium: An Instantaneous Probe-and-Rebalance Multimodal Learning Approach. *IJCAI*, 2025.

Take Home Message

Robust multimodal representation learning:

- Key Points of Efficiently Mining Multimodal Information for Large Models
- Multi-dimensional Challenges: data-, model-, and learning-level

Our contributions: A Series of Multidimensional RMML Solutions

- Data-Level: sampling sequence
- Model-Level: weak modality capability enhancement; strong modality masking
- Learning-Level: fusion and alignment simultaneously; instantaneous probe-and-rebalance fusion

Thanks

Q&A

