

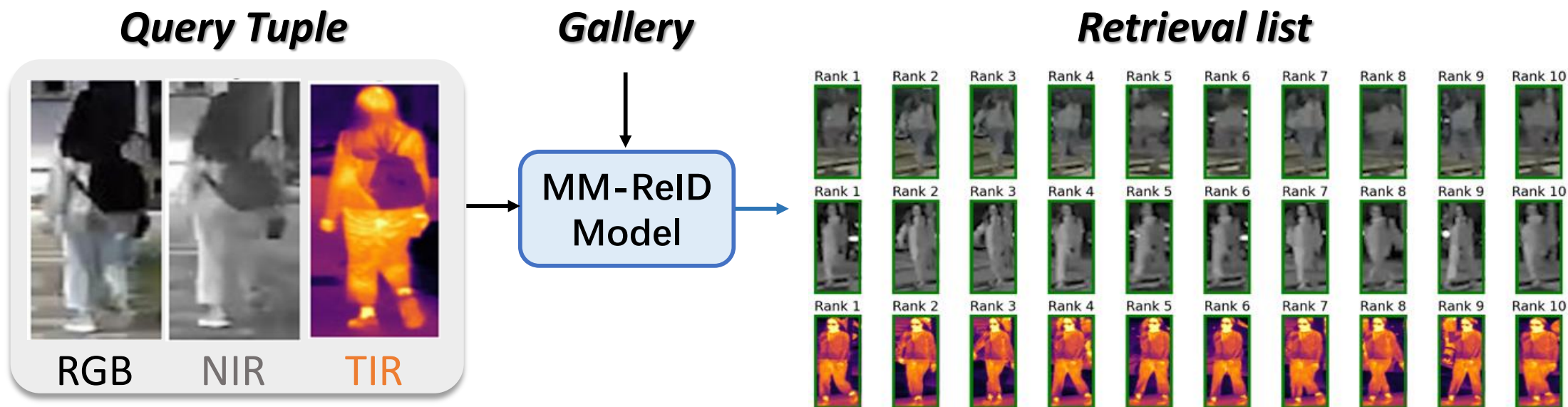
Miss-ReID: Delivering Robust Multi-Modality Object Re-Identification Despite Missing Modalities

Ruida Xi

State Key Laboratory of Electromechanical Integrated
Manufacturing of High-Performance Electronic Equipment, Xidian University

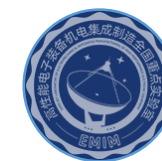
Paper: <https://nips.cc/virtual/2025/poster/119663>

Multi-Modality Object Re-Identification

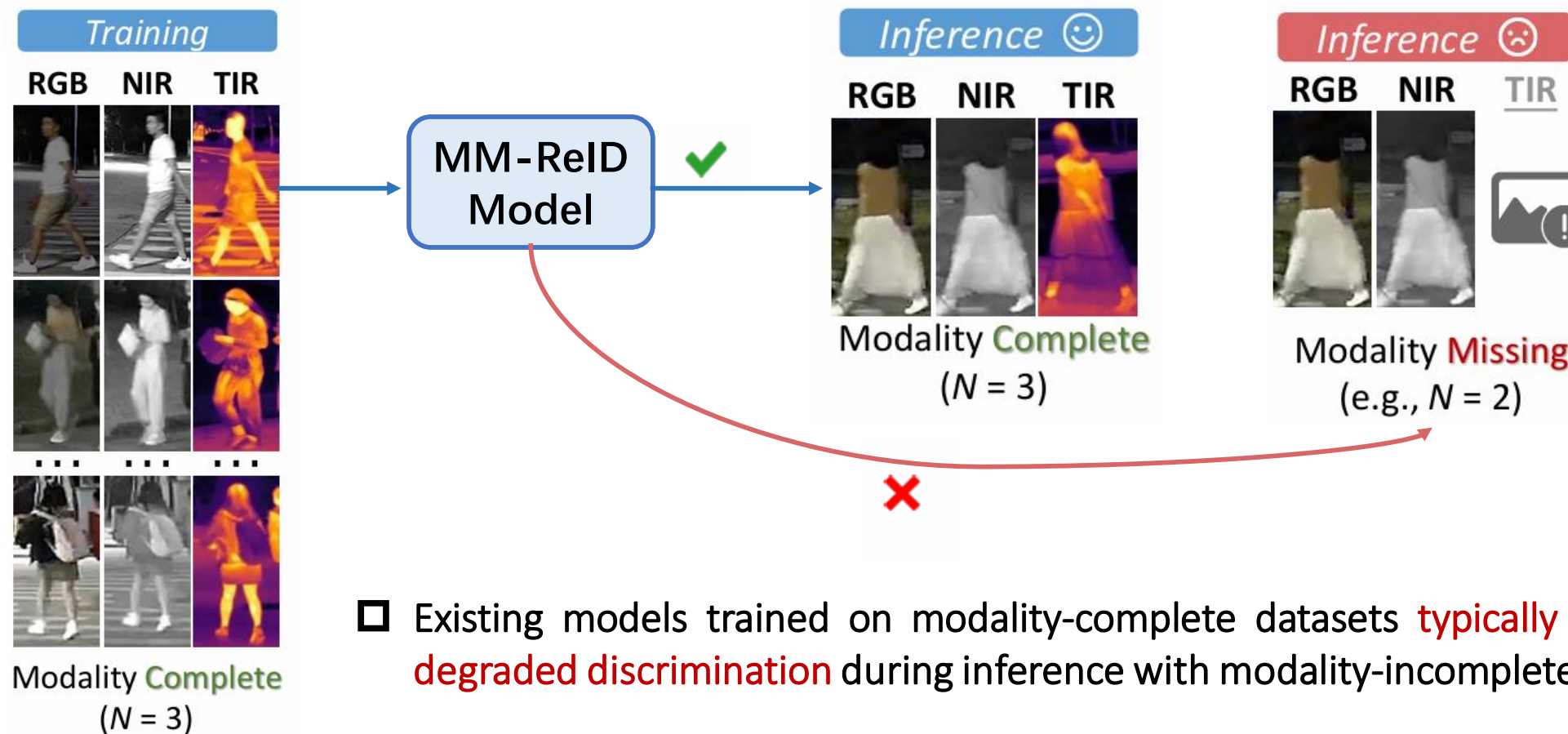


Multi-modality object Re-IDentification (ReID) targets to retrieve special objects by integrating complementary information from diverse visual sources.

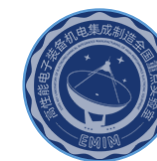
Problem



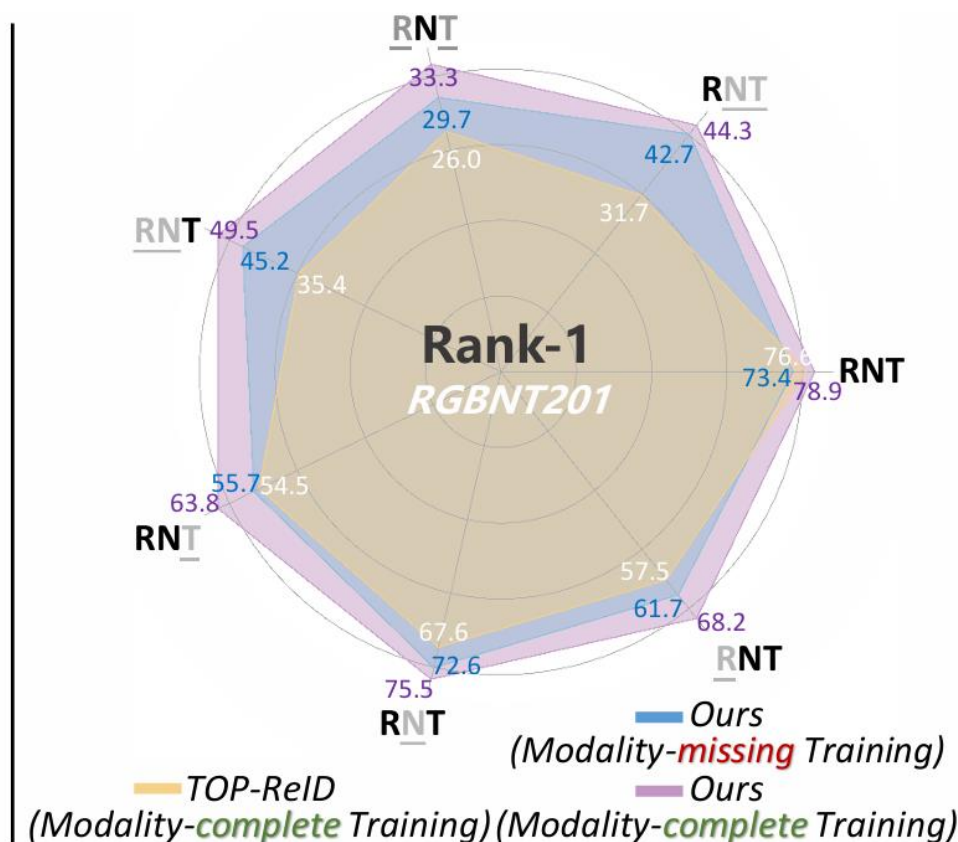
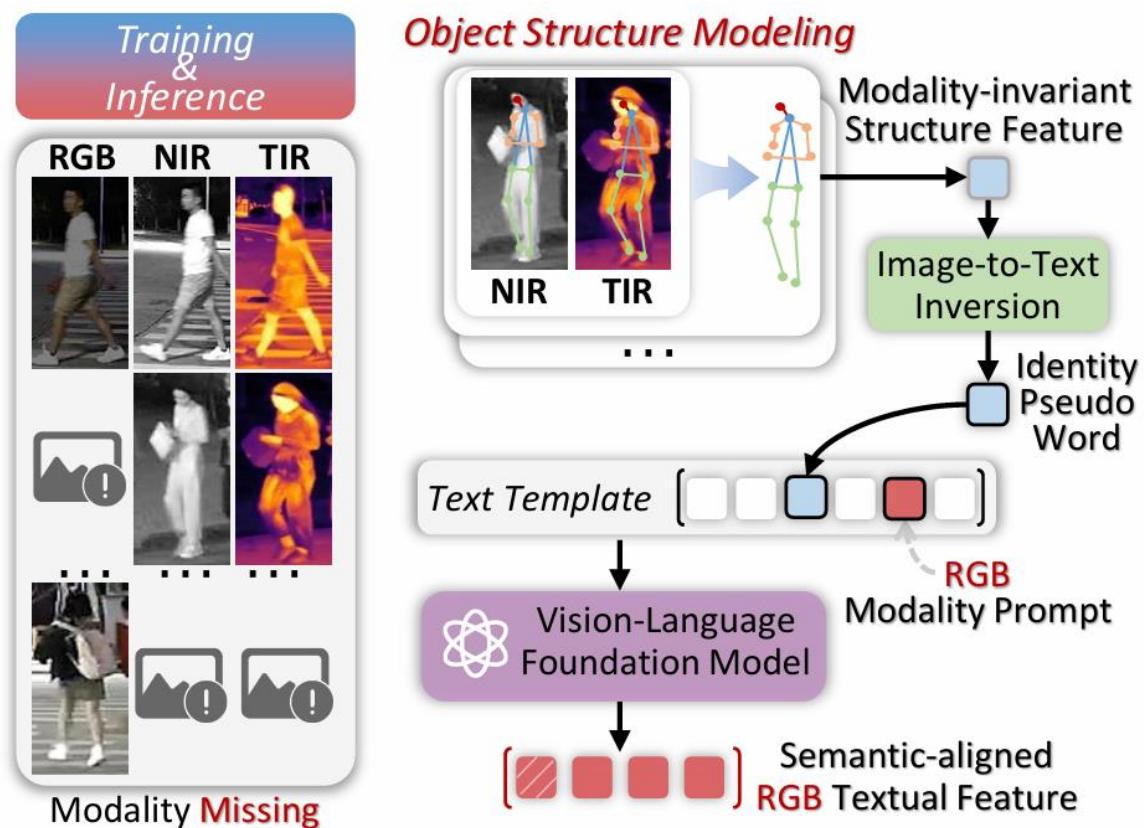
- Though exhibiting promising performance, existing multi-modality ReID methods **typically rely on an assumption regarding the modality completeness**, which may not hold in practice owing to privacy protections, sensor failures or security requirements.



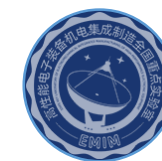
Motivation



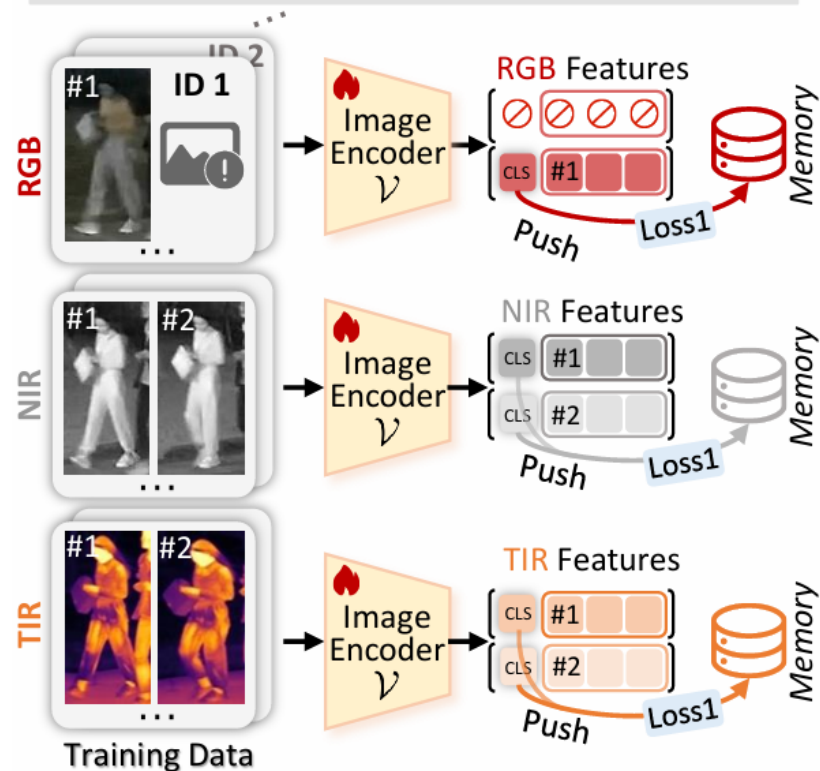
- By harnessing VLMs' open-world vision-text alignment, text-derived semantic features may effectively compensate for incomplete visual information, enabling robust solutions for modality-missing training and inference.



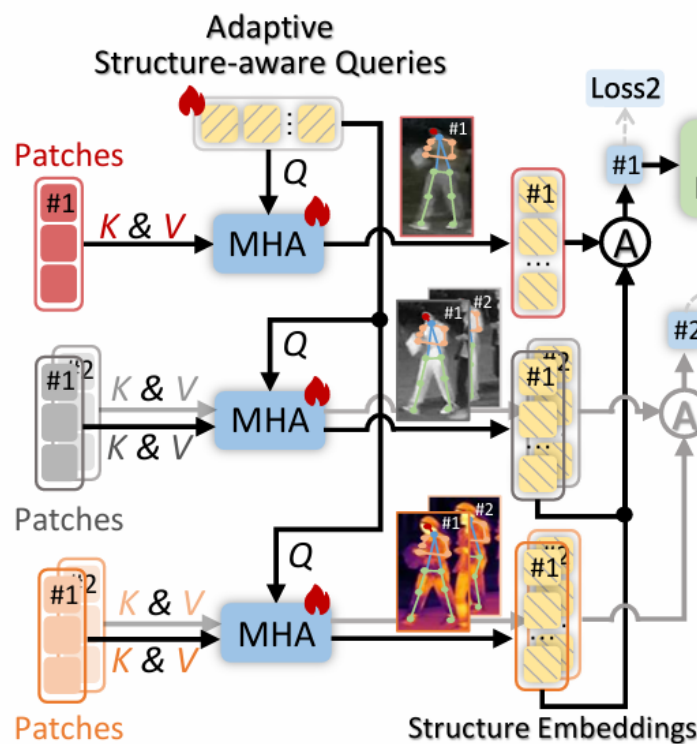
Framework



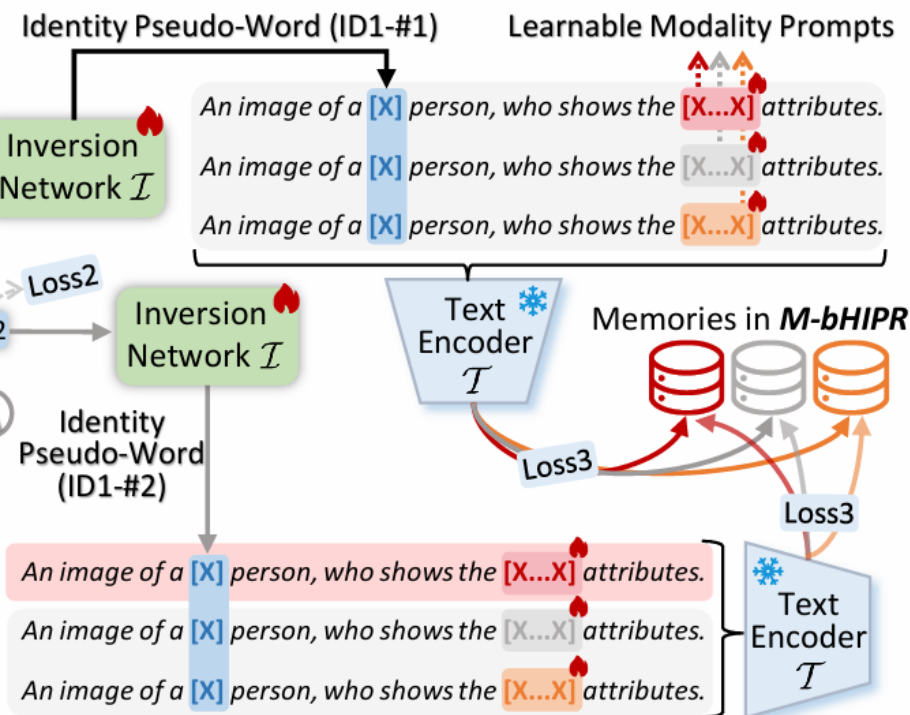
1. Memory-based Heterogeneous Identity Prototype Representation (*M-bHIPR*)



2. Modality-invariant Object Structure Modeling (*M-iOSM*)



3. Language-driven Missing Modality Completion (*L-dMMC*)



#1: The 1-st sample with Identity(ID) 1

MHA Multi-Head Attention

A Average Pooling

Learnable

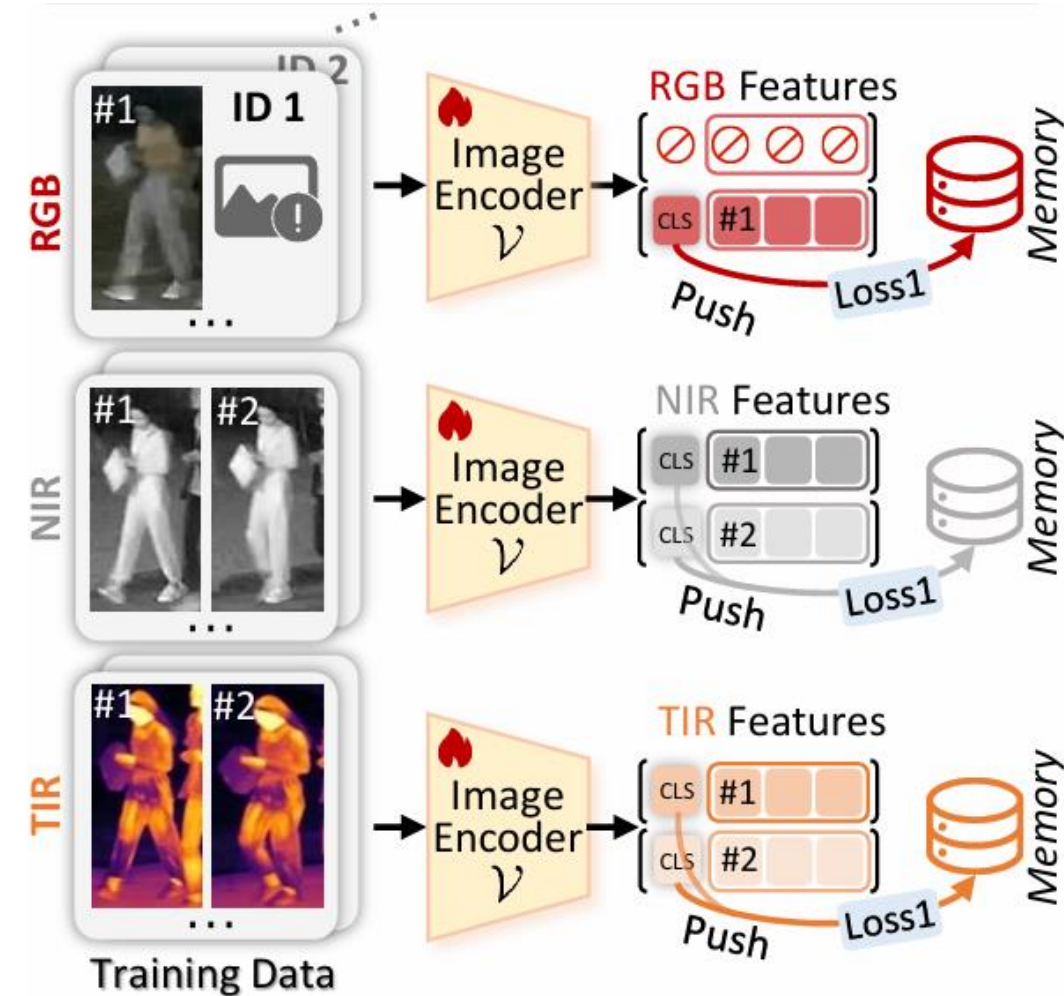
Frozen

Loss1 $\mathcal{L}_{Tri}^{m,V} + \mathcal{L}_{CE}^{m,V} + \mathcal{L}_{Nce}^m$

Loss2 $\mathcal{L}_{Tri}^S + \mathcal{L}_{CE}^S$

Loss3 $\mathcal{L}_{Tri}^{m,T} + \mathcal{L}_{CE}^{m,T} + \mathcal{L}_{T2V}^m$

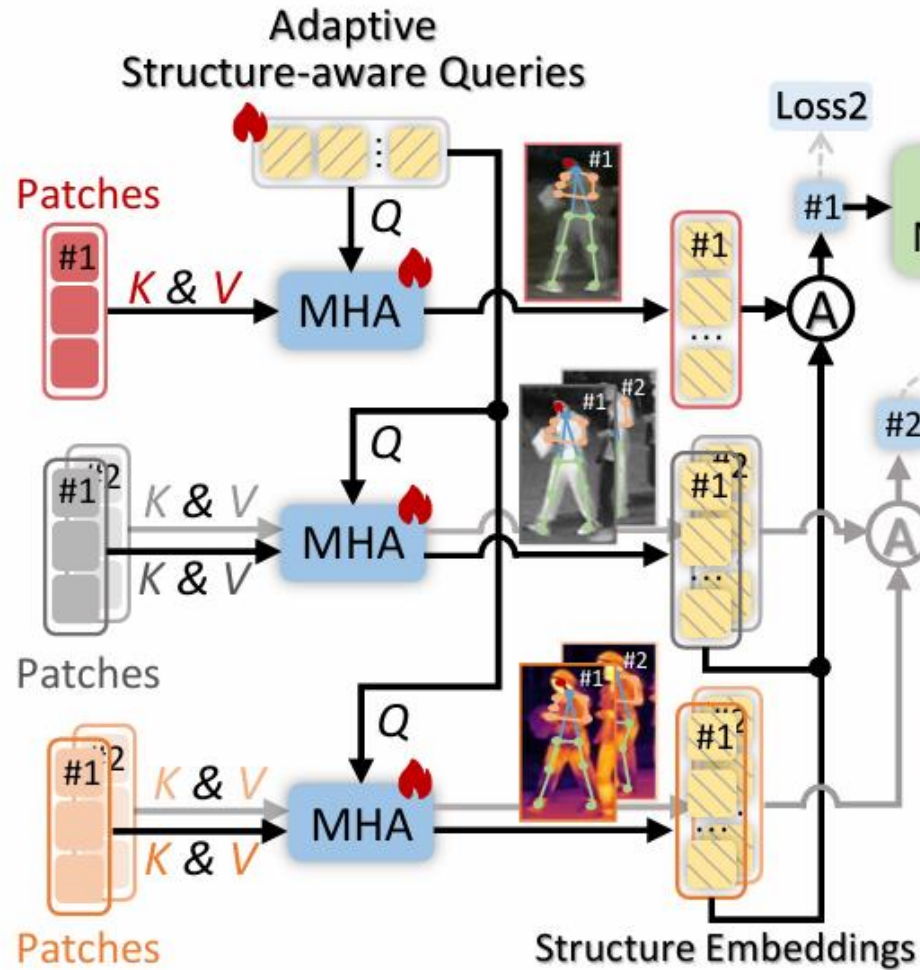
- Miss-ReID is the first work to handle multi-modality ReID under more general modality-missing scenarios encountered during both training and inference.



Memory-based Heterogeneous Identity Prototype Representation



Firstly, M-bHIPR extracts diverse visual features from **accessible modalities**, and then builds **modality-specific memory banks** to store heterogeneous prototypes for each individual identity, ensuring the preservation of multi-modality characteristics.



Modality-invariant
Object Structure Modeling



Afterwards, M-iOSM employs **structure-aware query interactions** to dynamically distill modality-invariant object structures from existing localized visual patches.

Language-driven Missing Modality Completion



- ✓ By leveraging the **textual inversion** technique, the extracted visual structural features are further reversed into pseudo-word tokens that encapsulate the identity-relevant structural semantics with L-dMMC module.
- ✓ Ultimately, the inverted tokens, integrated with diverse learnable **modality prompts**, are embedded into crafted textual templates to form the personalized linguistic descriptions for diverse modalities.
- ✓ Benefiting from VLMs' inherent vision-text alignment capability, L-dMMC produces the textual embeddings to substitute the absent visual cues.

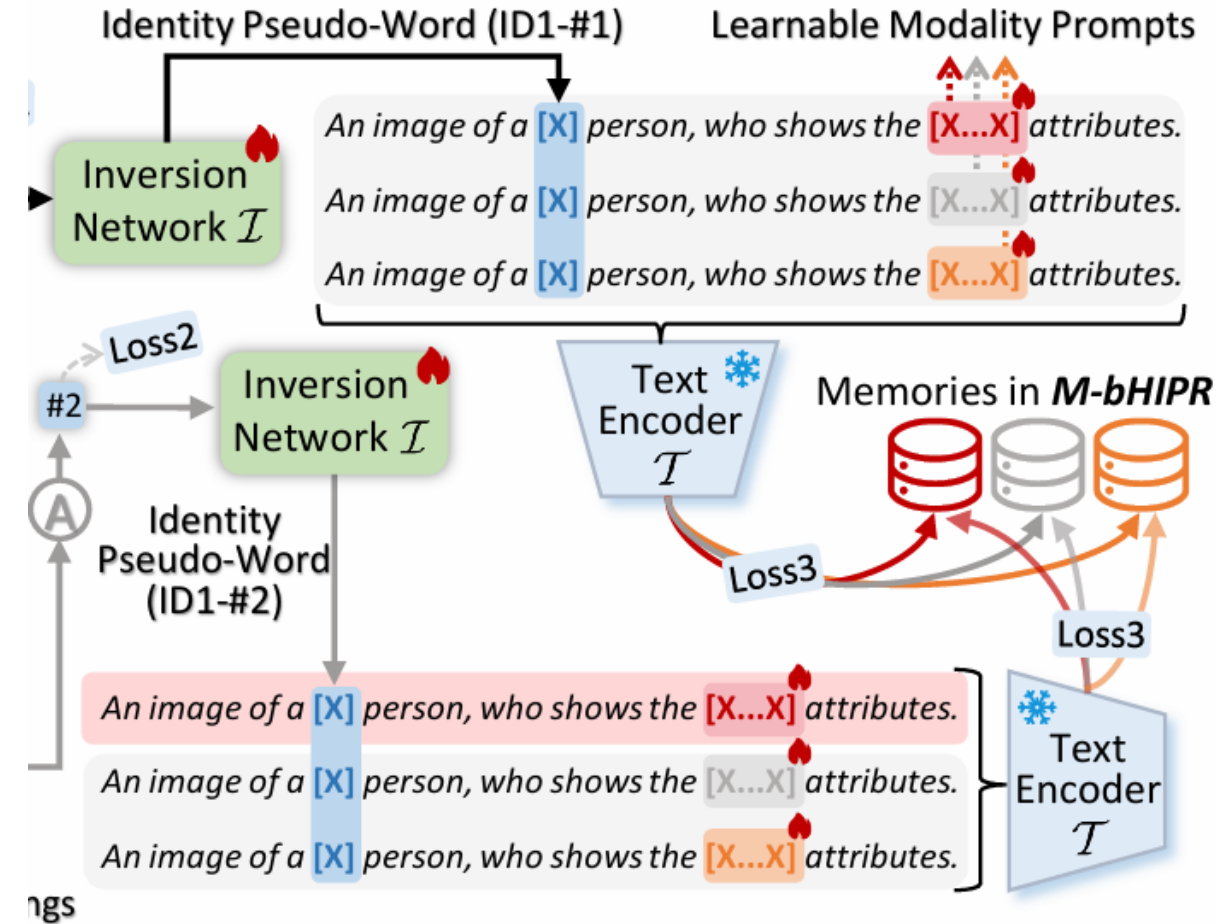


Table 1: The impacts of various components. We report the comparison results between different combinations (Model **B** – **F**) and the baseline (Model **A**) under both **modality-complete** and **-missing** training settings on RGBNT201. Here, ‘**Modality Complete**’ represents learning the modality-complete data during **training**, and ‘ $\eta = (0.1, 0.1, 0.1)$ ’ denotes randomly abandoning 10% RGB images, 10% NIR images, and 10% TIR images during **training**. The evaluations are both conducted across six modality-missing scenarios, and mean mAP and R-1 are reported below.

Index	Modules			Complexity		Modality Complete		$\eta = (0.1, 0.1, 0.1)$	
	M-bHIPR	M-iOSM	L-dMMC	Params	FLOPs	Mean mAP	Mean R-1	Mean mAP	Mean R-1
A	✗	✗	✗	86.4M	34.3G	48.9	50.4	46.4	47.0
B	✓	✗	✗	86.4M	34.3G	51.1(+2.2)	51.4(+1.0)	47.4(+1.0)	48.0(+1.0)
C	✗	✓	✗	86.4M	34.3G	50.2(+1.3)	52.1(+1.7)	46.9(+0.5)	48.7(+1.7)
D	✓	✓	✗	86.4M	34.3G	53.3(+4.4)	54.1(+3.7)	49.4(+3.0)	49.8(+2.8)
E	✗	✓	✓	89.6M	43.6G	52.1(+3.2)	53.4(+3.0)	47.4(+1.0)	49.7(+2.7)
F	✓	✓	✓	89.6M	43.6G	54.6(+5.7)	55.7(+5.3)	50.1(+3.7)	51.3(+4.3)

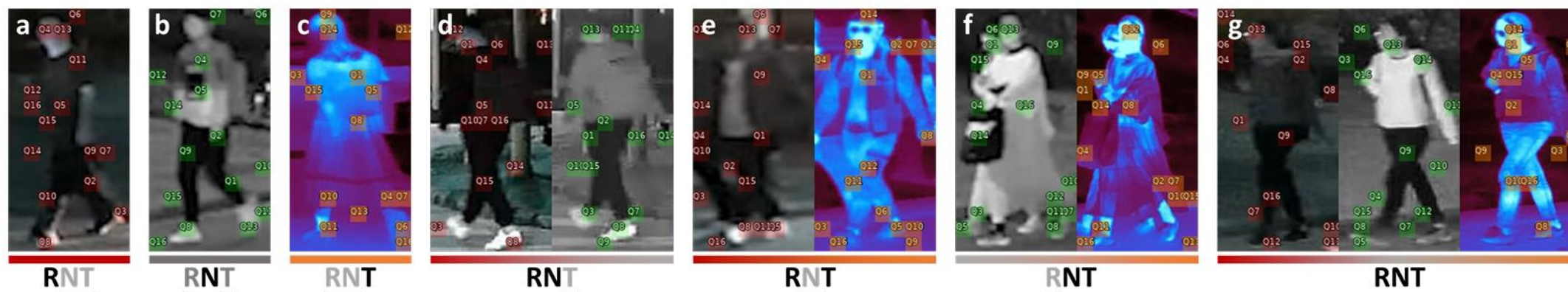
Table 2: Performance comparisons under **modality-missing** situations that only occur at the **inference** phase of multi-modality person ReID on RGBNT201. † denotes the model that is trained using both images and their corresponding text annotations. The best results are labeled with **boldface**. ↓x.x% and ↓x.x% highlight the lowest mAP and R-1 drop rates, respectively. ‘–’ indicates that the metric is unpublished.

Methods	RNT		<u>RNT</u>		<u>RNT</u>		<u>RNT</u>		<u>RNT</u>		<u>RNT</u>		<u>RNT</u>		Mean	
	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1
PCB [ECCV 2018]	32.8	28.1	23.6	24.2	24.4	25.1	19.9	14.7	20.6	23.6	11.0	6.8	18.6	14.4	19.7	18.1
			↓28.0%	↓13.9%	↓25.6%	↓10.7%	↓39.3%	↓47.7%	↓37.2%	↓16.0%	↓66.5%	↓75.8%	↓43.3%	↓48.8%	↓39.9%	↓35.6%
TOP-ReID [AAAI 2024]	72.3	76.6	54.4	57.5	64.3	67.6	51.9	54.5	35.3	35.4	26.2	26.0	34.1	31.7	44.4	45.4
			↓24.8%	↓24.9%	↓11.1%	↓11.7%	↓28.2%	↓28.9%	↓51.2%	↓53.8%	↓63.8%	↓66.1%	↓52.8%	↓58.6%	↓38.6%	↓40.7%
DeMo [AAAI 2025]	79.0	82.3	63.3	65.3	72.6	75.7	56.2	54.1	45.6	46.5	26.3	24.9	40.3	38.5	50.7	50.8
			↓19.9%	↓20.7%	↓8.1%	↓8.0%	↓28.9%	↓34.3%	↓42.3%	↓43.5%	↓66.7%	↓69.7%	↓49.0%	↓53.2%	↓35.8%	↓38.3%
IDEA [†] [CVPR 2025]	80.2	82.1	62.9	–	71.5	–	58.4	–	43.3	–	27.1	–	39.9	–	50.5	–
			↓21.6%	↓–%	↓10.8%	↓–%	↓27.2%	↓–%	↓46.0%	↓–%	↓66.2%	↓–%	↓50.2%	↓–%	↓37.0%	↓–%
Miss-ReID [Ours]	76.9	78.9	66.6	68.2	72.4	75.5	63.2	63.8	47.2	49.5	34.5	33.3	43.9	44.3	54.6	55.7
			↓13.4%	↓13.6%	↓5.9%	↓4.3%	↓17.8%	↓19.1%	↓38.6%	↓37.3%	↓55.1%	↓57.8%	↓42.9%	↓43.9%	↓29.0%	↓29.4%

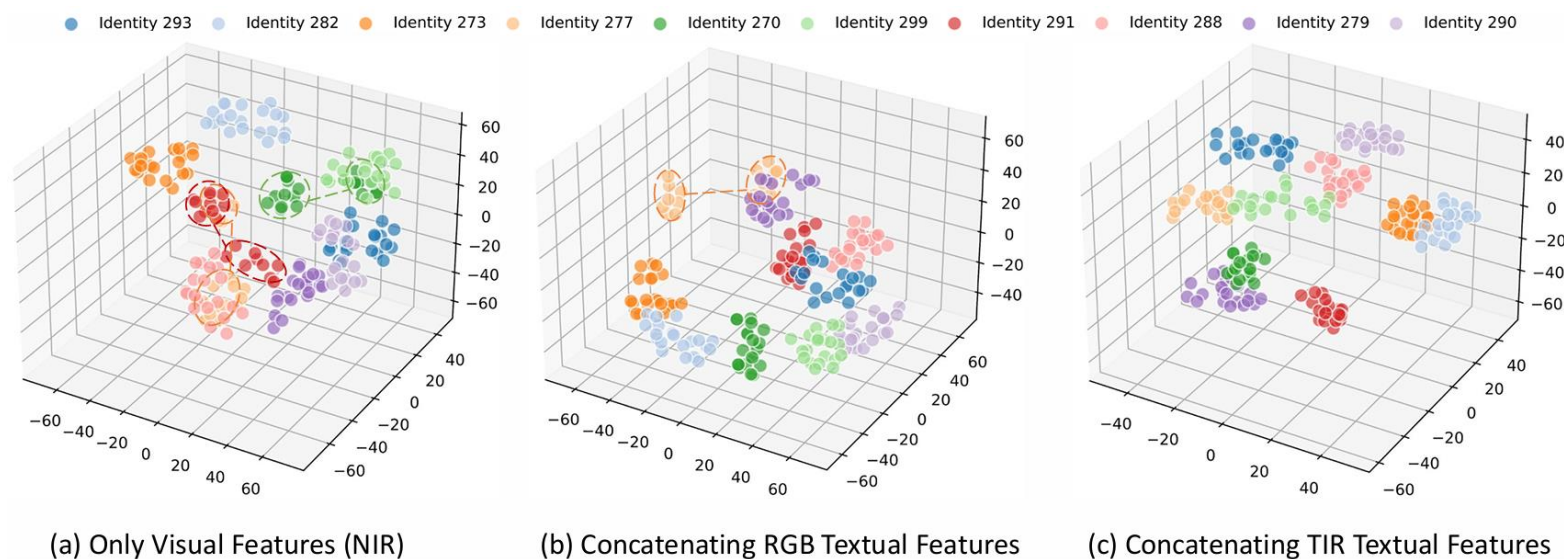
Performance Analysis of Miss-ReID under Varying Tri-modality Missing Rates

Table 3: Performance comparisons of setting different **tri-modality missing rates** on RGBNT201. Each tuple $(\eta_{rgb}, \eta_{nir}, \eta_{tir})$ represents the proportion of randomly abandoned RGB, Near-Infrared, and Thermal-Infrared images during **training**.

Tri-Modality Missing Rate η	(0.0, 0.0, 0.0)		(0.1, 0.1, 0.1)		(0.3, 0.3, 0.3)		(0.5, 0.5, 0.5)		(0.1, 0.3, 0.5)		(0.5, 0.3, 0.1)	
	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1	mAP	R-1
RNT	76.9	78.9	72.3	73.4	68.4	71.2	68.2	72.8	69.6	72.2	67.6	67.6
<u>RNT</u>	66.6	68.2	61.3	61.7	57.6	58.3	56.7	58.4	56.1	58.4	57.6	58.4
<u>RNT</u>	72.4	75.5	68.8	72.6	65.8	69.5	63.6	65.1	66.9	69.7	65.7	67.0
<u>RNT</u>	63.2	63.8	55.3	55.7	52.3	56.5	52.3	54.4	53.2	57.3	50.2	50.7
<u>RNT</u>	47.2	49.5	42.8	45.2	44.9	47.5	41.6	40.8	41.1	40.3	47.1	47.6
<u><u>RNT</u></u>	34.5	33.3	30.9	29.7	26.8	26.3	26.5	22.2	27.1	28.1	26.5	24.6
<u>RNT</u>	43.9	44.3	41.5	42.7	43.0	45.5	42.7	46.4	43.9	47.0	39.2	38.8
Mean	54.6	55.7	50.1	51.3	48.4	50.6	47.3	47.9	48.1	50.1	47.7	47.8



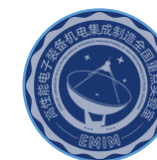
Visualizations of the attentive regions towards 16 well-learned structure-aware queries.



The feature distributions of 10 random identities.

Experiments

Retrieval Result

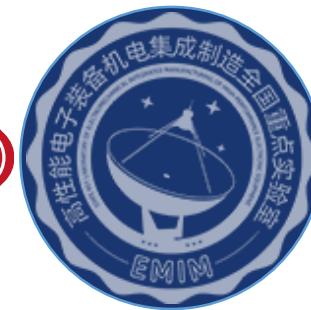


(a) Baseline

(b) Miss-ReID



NEURAL INFORMATION
PROCESSING SYSTEMS



Thanks for watching!

Ruida Xi