

# Bidirectional Multi-Step Domain Generalization for Visible-Infrared Person Re-Identification

# Mahdi Alehdagi<sup>1</sup>, Pourya Shamsolmoali<sup>2</sup>, Rafael M. O. Cruz<sup>1</sup> and Eric Granger<sup>1</sup>

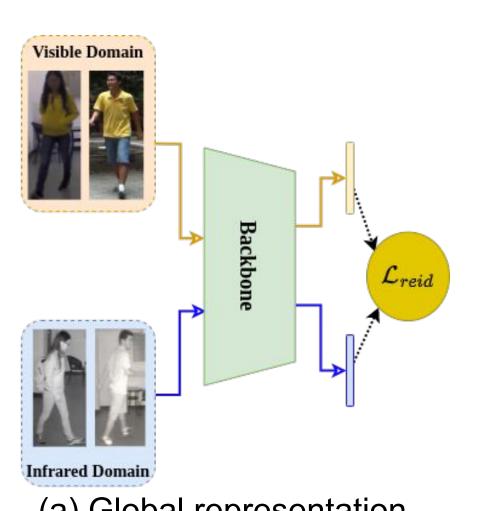


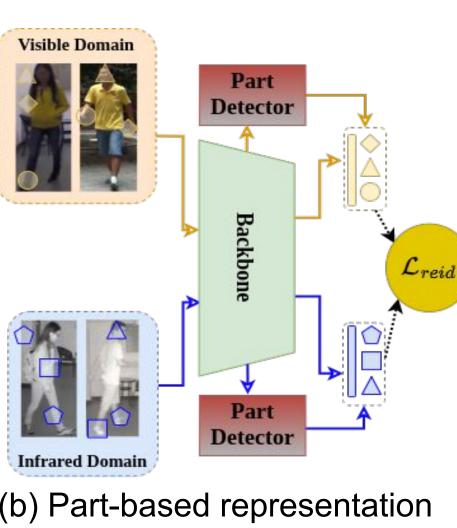
<sup>1</sup>LIVIA, ILLS, Dept. of Systems Engineering, ETS Montreal, Canada, <sup>2</sup>Dept. of Computer Science, University of York, UK mahdi.alehdaghi.1@ens.etsmtl.ca, pshams55@gmail.com, {rafael.menelau-cruz, eric.granger}@etsmtl.ca

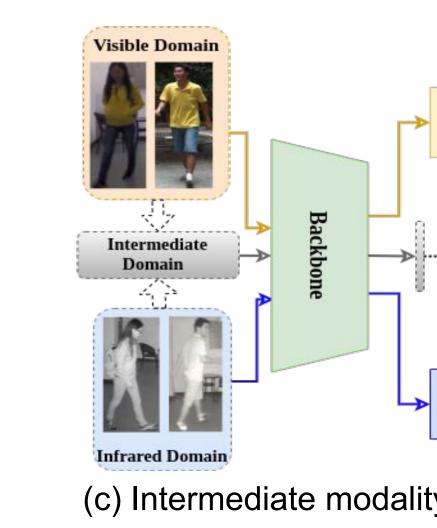


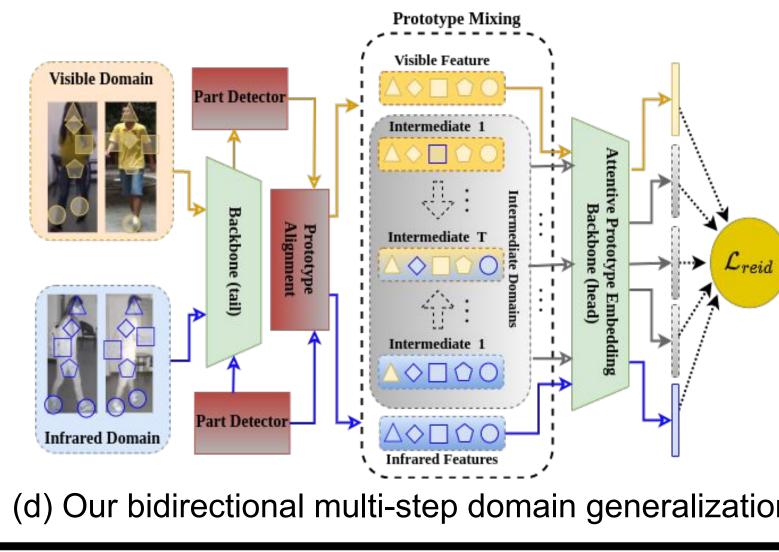
#### -Introduction

- V-I ReID is a challenging task due to significant discrepancies across V and I modalities
- State-of-the-art methods use global or part-based representations, but struggle with modality alignment
- Bridging modality gaps effectively requires leveraging intermediate feature spaces









## Bidirectional Multi-step Learning -

• Mixing Prototypes: creates intermediate feature representations for each modality by gradually increasing mixing rate

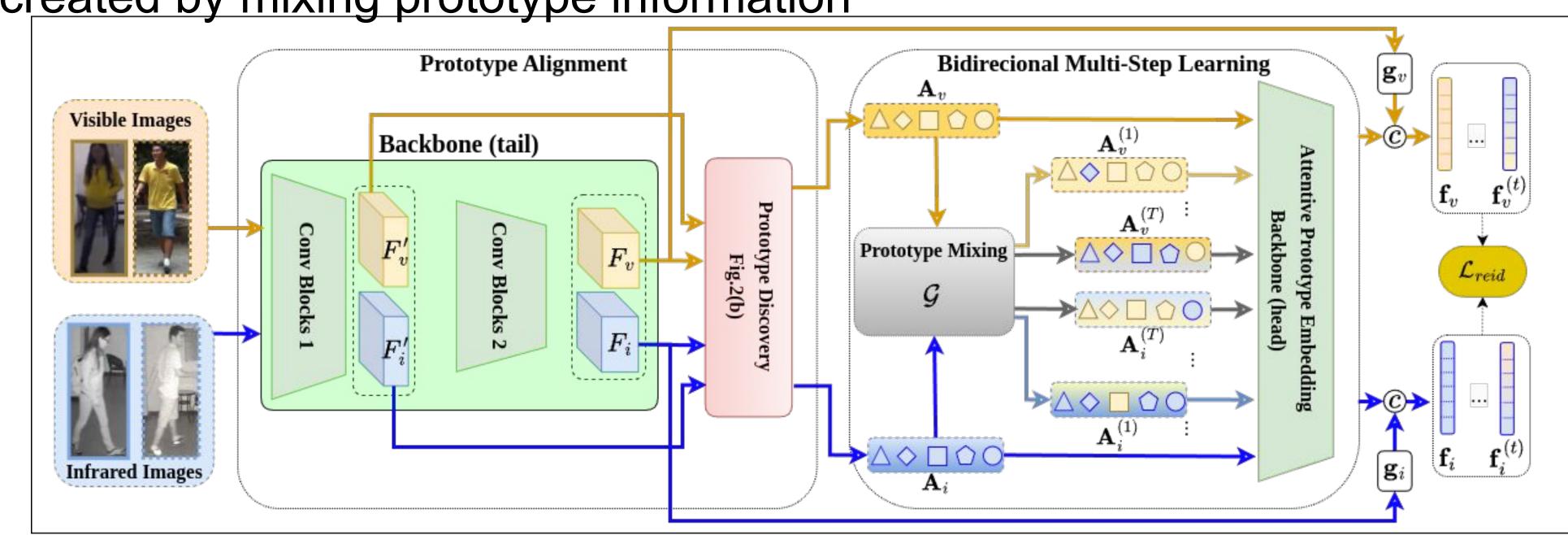
$$\mathbf{A}_m^{(t)} = \mathcal{G}(\mathbf{A}_m, \mathbf{A}_{ ilde{m}}, t) ext{ where } \mathbf{A}_m = [\mathbf{p}_m^1; \mathbf{p}_m^2; \dots; \mathbf{p}_m^K], \ \mathcal{G}(\mathbf{A}_m, \mathbf{A}_{ ilde{m}}, t) = [\mathbf{p}_{r(m, ilde{m}, t)}^1, \dots, \mathbf{p}_{r(m, ilde{m}, t)}^K] ext{ where } r(m, ilde{m}, t) = egin{cases} m & t/T \leq \mathcal{U}(0, 1) \ ilde{m} & ext{else} \end{cases}$$

- Joint feature representation:  $\mathbf{f}_m = [\mathcal{F}(\mathbf{A}_m); \mathbf{g}_m] \text{ and } \mathbf{f}_m^{(t)} = [\mathcal{F}(\mathbf{A}_m^{(t)}); \mathbf{g}_m]$ 
  - $\mathcal{L}_{ ext{bce}} = \mathcal{L}_{ ext{ce}}(\mathbf{f}_v) + \mathcal{L}_{ ext{ce}}(\mathbf{f}_v^{(t)}) + \mathcal{L}_{ ext{ce}}(\mathbf{f}_i) + \mathcal{L}_{ ext{ce}}(\mathbf{f}_i^{(t)})$  $\mathcal{L}_{ ext{re}} = \mathcal{L}_{ ext{bce}} + \mathcal{L}_{ ext{bcc}}$

# Proposed Bidirectional Multi-Step Domain Generalization

### BMDG Training Architecture:

- Prototype Alignment: extracts body part prototype representations from V and I images
- Bidirectional Multi-Step Learning: extracts discriminant features using multiple intermediate domains created by mixing prototype information

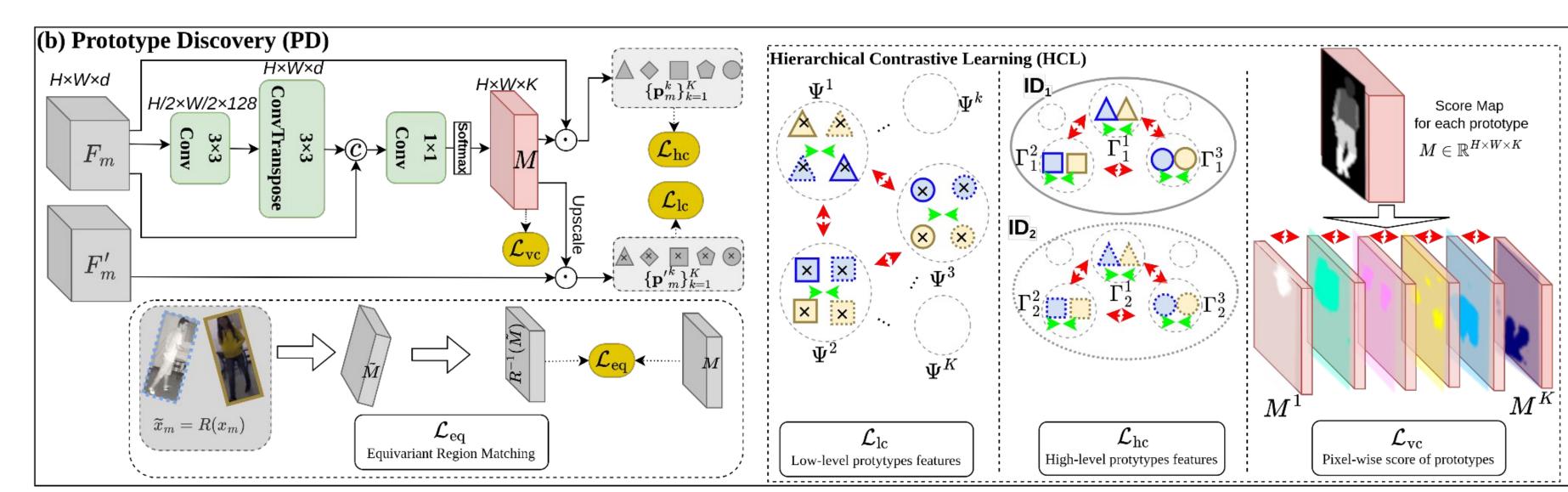


• Goal: maximize the similarity of extracted features of same person

$$\max\{S(\mathbf{f}_v^n,\mathbf{f}_v^{l,(t)})+S(\mathbf{f}_i^n,\mathbf{f}_i^{l,(t)})\cdot o\}$$
 where  $o$  = 1 if  $y_v^n=y_i^l$  else  $o$  = -1 and  $t\in\{1...T\}$  is step.

## -Prototype Alignment

- Prototype discovery (PD): mines prototypes from spatial features
- Hierarchical contrastive learning (HCL): encourages the prototypes to focus on similar semantics for all individuals without losing ID-discriminative information.



 $\bullet \; \mathsf{Losses} : \mathcal{L}_{\mathrm{lc}} = -\sum_{k=1}^K \sum_{\hat{\mathbf{p}}' \in \Psi^k} \log \frac{e^{\mathbf{p}'^k \cdot \hat{\mathbf{p}}'/\tau}}{e^{\mathbf{p}'^k \cdot \hat{\mathbf{p}}'/\tau} + \sum_{q \neq k} e^{\mathbf{p}'^k \cdot \hat{\mathbf{p}}'^q/\tau}}, \; \mathcal{L}_{\mathrm{hc}} = -\sum_{y=1}^{C_y} \sum_{k=1}^K \sum_{\hat{\mathbf{p}} \in \Gamma_y^k} \log \frac{e^{\mathbf{p}^k \cdot \hat{\mathbf{p}}/\tau}}{e^{\mathbf{p}^k \cdot \hat{\mathbf{p}}/\tau} + \sum_{q \neq k} e^{\mathbf{p}^k \cdot \hat{\mathbf{p}}'^q/\tau}}.$ 

$$\mathcal{L}_{ ext{eq}} = \sum_{k=1}^K \| M_m^k - R^{-1}( ilde{M}_m^k) \| \qquad \mathcal{L}_{ ext{vc}} = \sum_{k=1}^K \sum_{q=k+1}^K \sum_{u \in U} [\| M_m^k - M_m^q \|]_u,$$

# Comparison with State-of-Art Methods

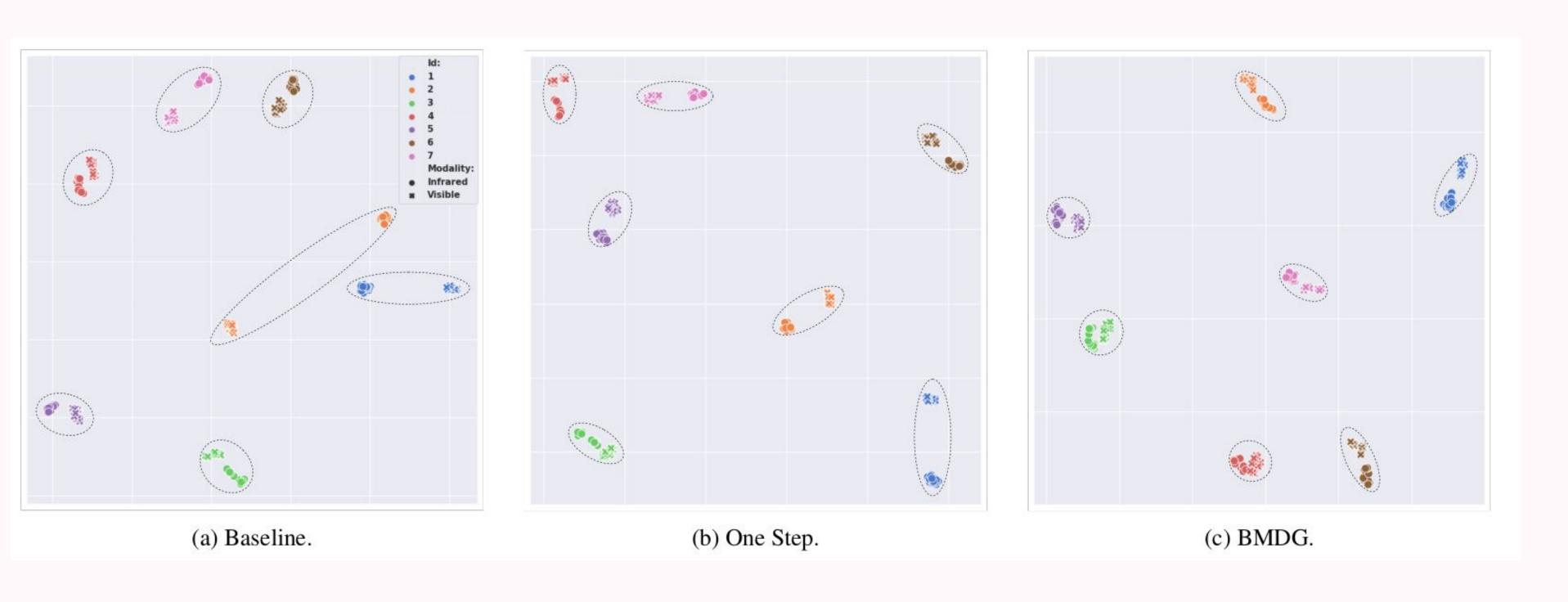
Table 1. Accuracy of the proposed BMDG and state-of-the-art methods on the SYSU-MM01 (single-shot setting) and RegDB datasets.

	Family		SYSU-MM01				RegDB				
			All S	earch	Indoor	Search	Visible	$\rightarrow$ Infrared	Infrared	$\rightarrow$ Visible	
1	Method Venue		R1	mAP	R1	mAP	R1	mAP	R1	mAP	
	AGW	TPAMI'20	47.50	47.65	54.17	62.97	70.05	50.19	70.49	65.90	
bal	CAJ	ICCV'21	69.88	66.89	76.26	80.37	85.03	79.14	84.75	77.82	
Glok	RAPV-T	ESA'23	63.97	62.33	69.00	75.41	86.81	81.02	86.60	80.52	
	G2DA	PR'23	63.94	60.73	71.06	76.01	_	-	-	-	
pa	DDAG	ECCV'20	54.74	53.02	61.02	67.98	69.34	63.46	68.06	90.31	
Jase	DDAG	ICCV'23	75.90	77.03	82.20	80.01	91.07	91.45	92.09	92.01	
Ŧ		ICCV'23	74.66	71.73	79.69	83.68	94.51	88.67	93.64	87.61	
Pa	PartMix	CVPR'23	<u>77.78</u>	74.62	81.52	84.83	84.93	82.52	85.66	82.27	
ate	SMCL	ICCV'21	67.39	61.78	68.84	75.56	83.93	79.83	83.05	78.57	
edi	RPIG	ECCVw'22	71.08	67.56	82.35	82.73	87.95	82.73	86.80	81.26	
E	G2DA	PR'23	63.94	60.73	71.06	76.01	-	_	-		
nte	SEFL	CVPR'23	75.18	70.12	78.40	81.20	91.07	85.23	92.18	86.59	
	BMDG	<b>-</b>	78.08	78.22	83.59	86.35	94.76	92.21	94.56	93.07	
	SEFL	No. of the Control of	75.18	70.12	78.40	81.20	CONTRACTOR CO.	V-850-350 htt	0.50		

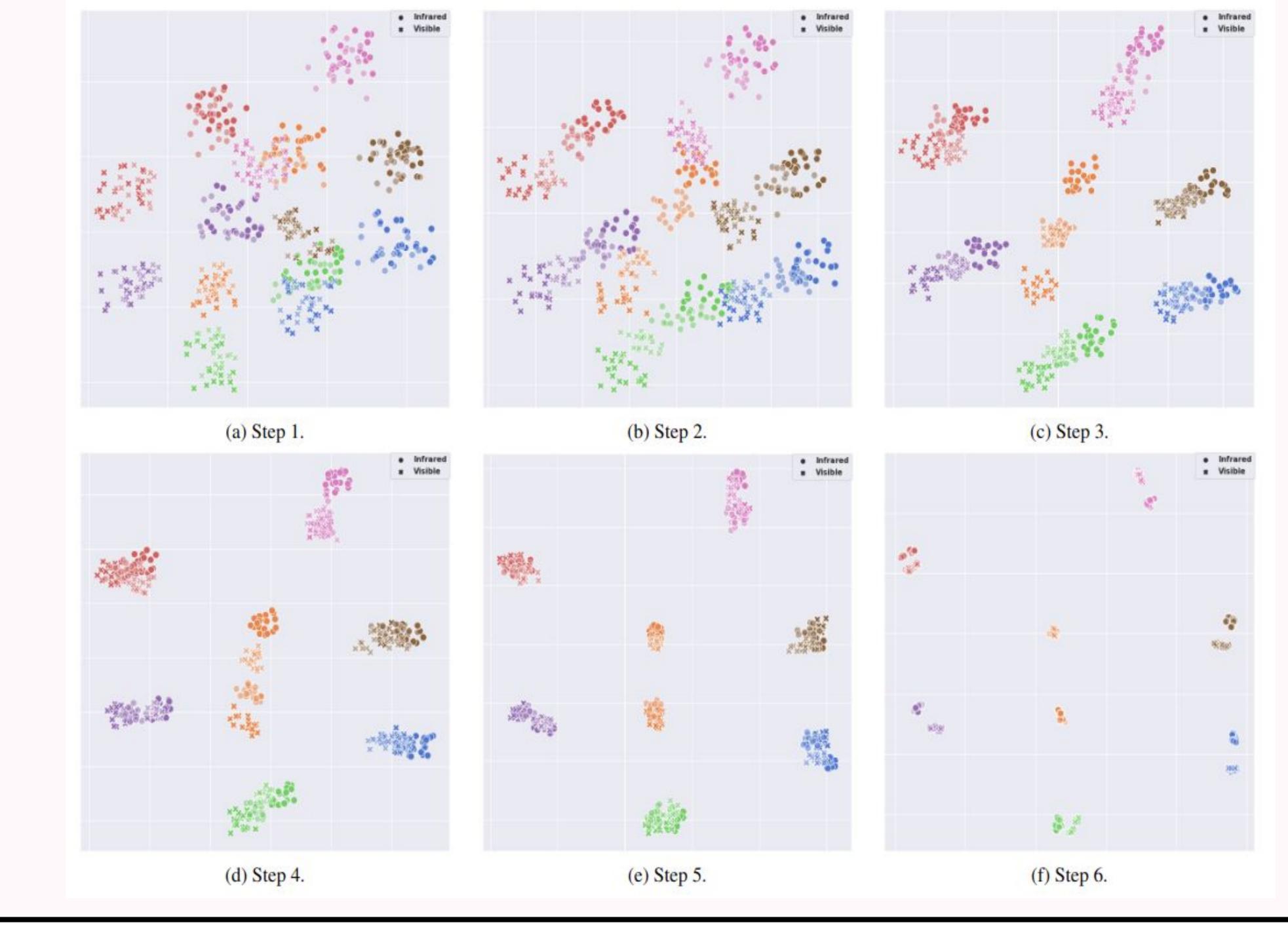
- BMDG outperforms state-of-the-art methods in various cross-modal settings
- It improves robustness and accuracy by gradually reducing the modality gap

# -Visual Results (t-SNE Plots)

Comparing BMDG against One-Step and Baseline approaches:



Impact of using our prototype mixing for 6 intermediate steps (T):



#### -Ablation Studies -

Table 2. Impact on rank-1 and mAP accuracy using different settings for multi-step domain generalization.

	Number of part prototypes $(K)$								
Settings	4		5		6		7		
	R1	mAP	R1	mAP	R1	mAP	R1	mAP	
Single step									
One $(V \rightarrow I)$	71.0	66.5	71.5	67.1	72.3	67.5	71.6	66.8	
One $(I \rightarrow V)$	70.8	66.3	71.2	67.0	72.6	67.6	71.0	66.7	
Bidirectional	74.0	71.8	74.1	71.9	75.4	72.8	73.5	70.2	

Our bi-directional approach avoids forgetting source modalities since it is not biased to a specific modality

Table 3. Accuracy of part-based ReID methods with BMDG on the SYSU-MM01.

Method	R1 (%)	mAP (%)
DDAG	53.62	52.71
DDAG with BMDG	55.36 (+1.741)	54.05 (+1.341)
MPANet	66.24	62.89
MPANet with BMDG	68.74 (+2.501)	64.25 (+1.361)
SAAI	71.87	68.16
SAAI with BMDG	73.69 (+1.821)	70.08 (+1.921)

BMDG can be integrated into any part-based V-I ReID method to improve performance

Table 4. mAP accuracy of BMDG for different values of K and T.

T	Number of part prototypes $(K)$								
1	3	4	5	6	7	10			
0	65.98	67.03	67.23	68.11	67.55	65.66			
1	67.42	68.72	69.28	69.46	68.32	69.28			
2	69.51	70.08	70.72	71.14	69.97	71.25			
3	71.44	71.69	71.82	72.02	71.06	71.67			
4	7-1	71.98	72.15	72.86	71.19	69.00			
6	-	-	-	72.40	71.17	69.54			
10	-	_	-	-	-	69.46			

Accuracy grows with T for a number of prototypes K

#### -Conclusion

- We propose a BMDG training architecture, allowing to learn the discriminative and complementary set of prototypes from both modalities
- These prototypes create intermediate domains by mixing modalities, reducing the domain gap
- BMDG achieves state-of-the-art performance on cross-modal evaluation settings on several challenging V-I ReID datasets