

# Beyond Appearance: a Semantic Controllable Self-supervised Learning Framework for Human-Centric Visual Tasks

*Weihua Chen, Xianzhe Xu, Jian Jia, Hao Luo, Yaohua Wang, Fan Wang, Rong Jin, Xiuyu Sun*

Paper ID: 8746

Session Tag: WED-PM-257

# SOLIDER

## Beyond Appearance: a **S**emantic **C**ontrollable **S**elf-supervised **L**earning Framework for Human-Centric Visual Tasks

*Weihua Chen, Xianzhe Xu, Jian Jia, Hao Luo, Yaohua Wang, Fan Wang, Rong Jin, Xiuyu Sun*

Paper ID: 8746

Session Tag: WED-PM-257

# 1. MOTIVATION

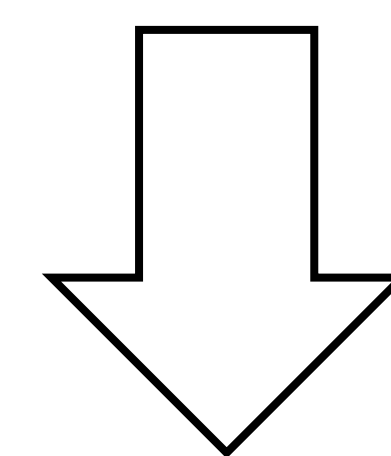
## Background

1. Human-centric tasks are **crucial in application**;
2. **Massive unlabeled** human images available;
3. **Semantic information** is important.
4. Different downstream tasks require **different needs** of semantic.

# 1. MOTIVATION

## Background

1. Human-centric tasks are **curial in application**;
2. **Massive unlabeled** human images available;
3. **Semantic information** is important.
4. Different downstream tasks require **different needs** of semantic.



## Goal

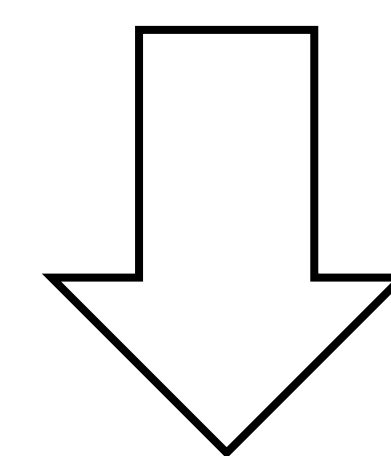
1. Learn **a general human representation** from **massive unlabeled** data.
2. The pre-trained model can benefit to downstream tasks and meet their **different semantic requirements**. (**One pretrained model adapt to All.**)



# 1. MOTIVATION

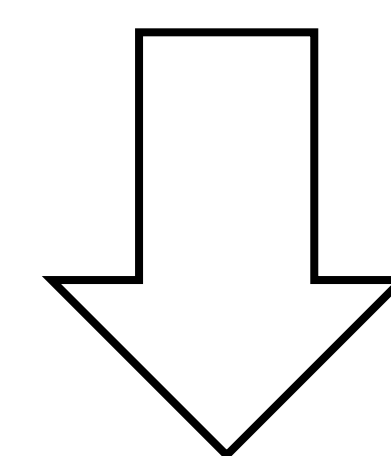
## Background

1. Human-centric tasks are **curial in application**;
2. **Massive unlabeled** human images available;
3. **Semantic information** is important.
4. Different downstream tasks require **different needs** of semantic.



## Goal

1. Learn **a general human representation** from **massive unlabeled** data.
2. The pre-trained model can benefit to downstream tasks and meet their **different semantic requirements**. (**One pretrained model adapt to All.**)



## Main Contributions

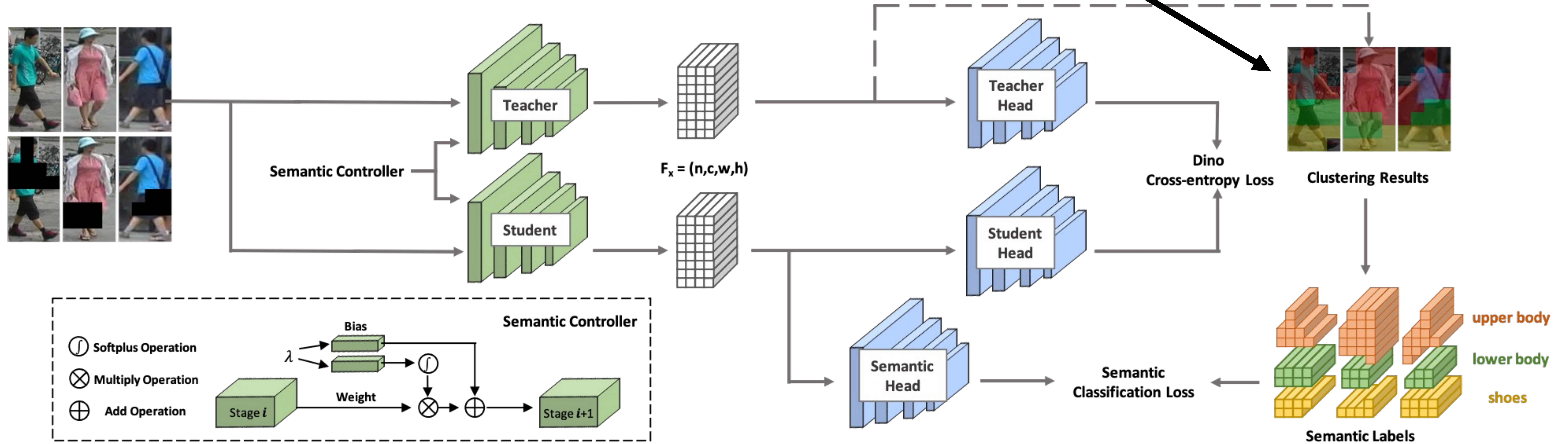
1. Build a **semantic supervision from human prior** to train SOLIDER.
2. Design a **semantic controller** to adjust the semantic rate in the pre-trained model.



# 2. METHOD

## Contribution 1. Semantic Supervision from Human Prior

The framework of SOLIDER

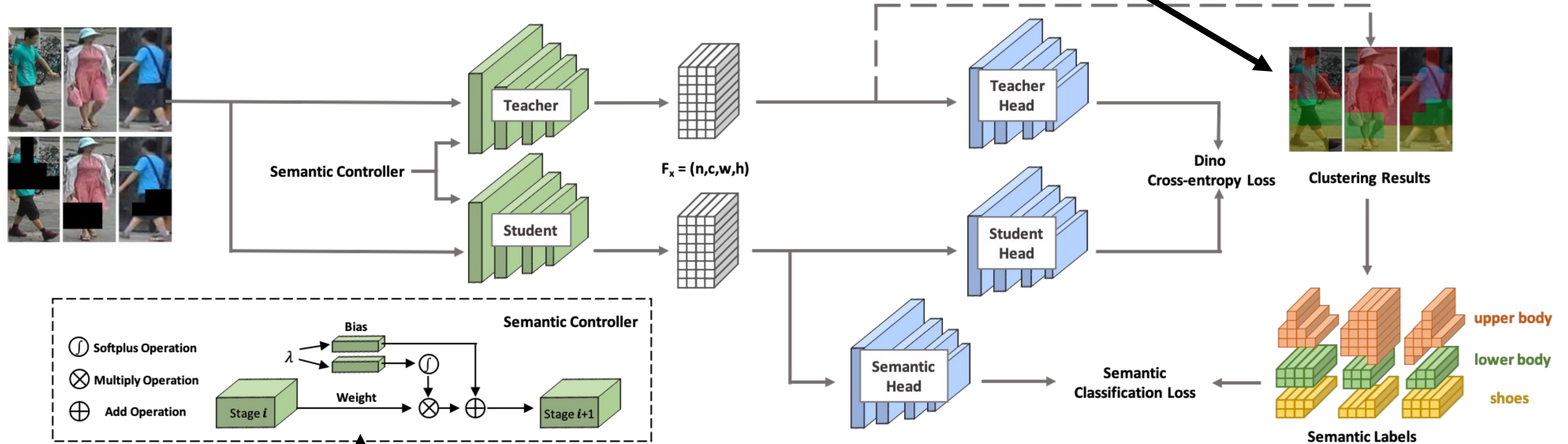




# 2. METHOD

## Contribution 1. Semantic Supervision from Human Prior

The framework of SOLIDER



## Contribution 2. Semantic Controller



# 3. EXPERIMENTS

## 3.1. Can SOLIDER learn semantic ?

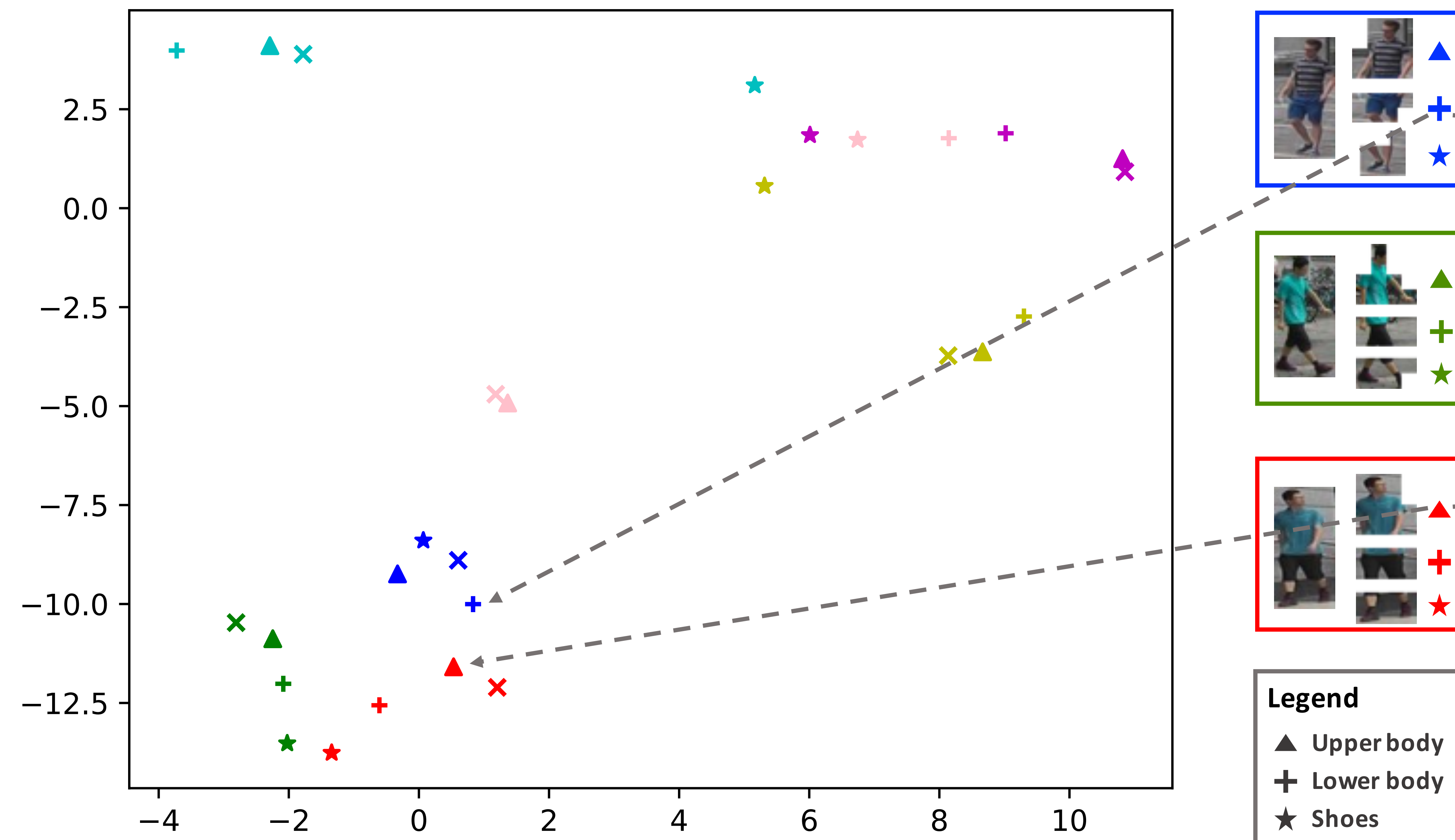
**DINO space:**

Different parts from same person are closer.

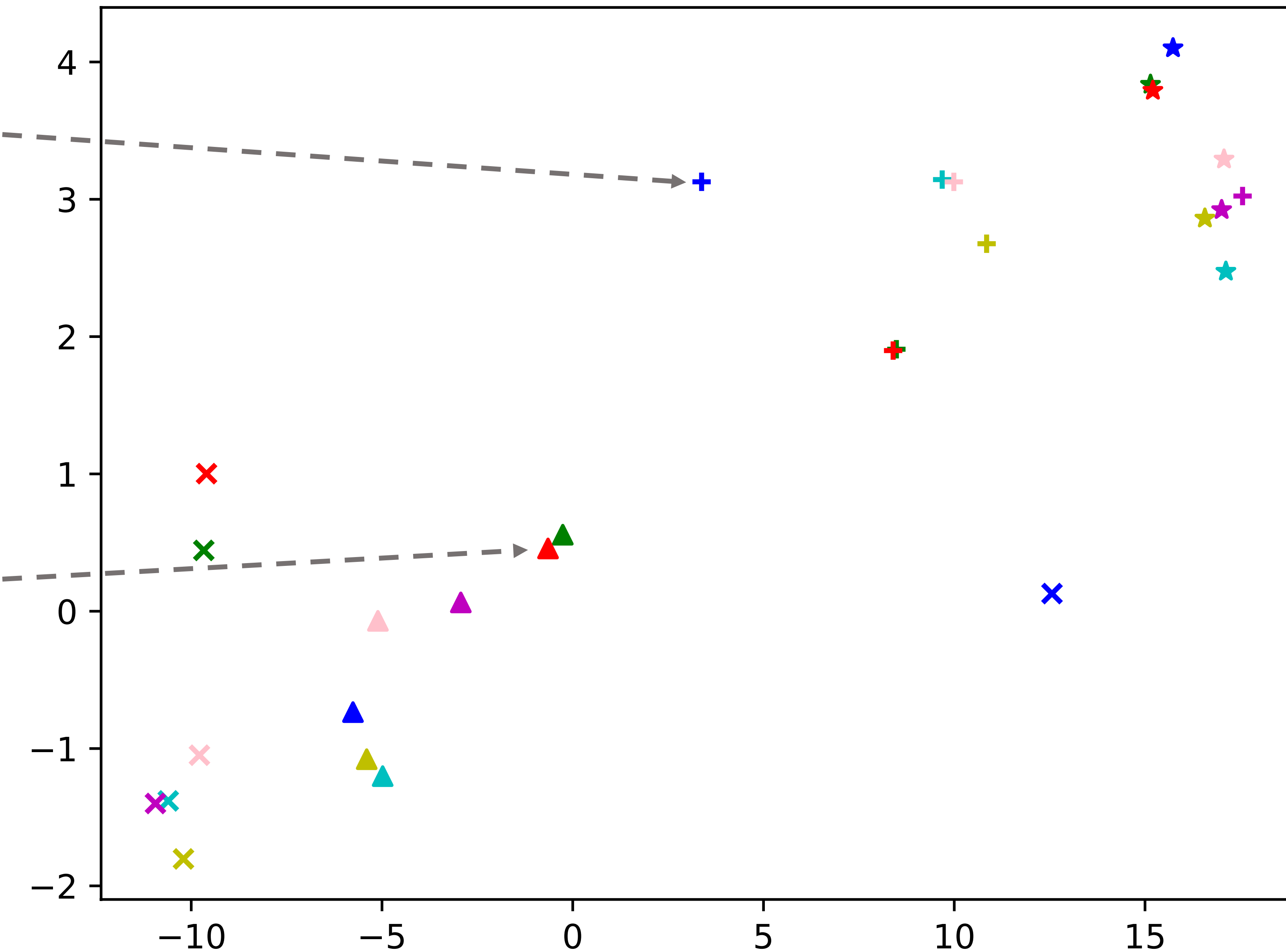
**SOLIDER space:**

Same semantic parts from different people are closer.

After DINO



After SOLIDER





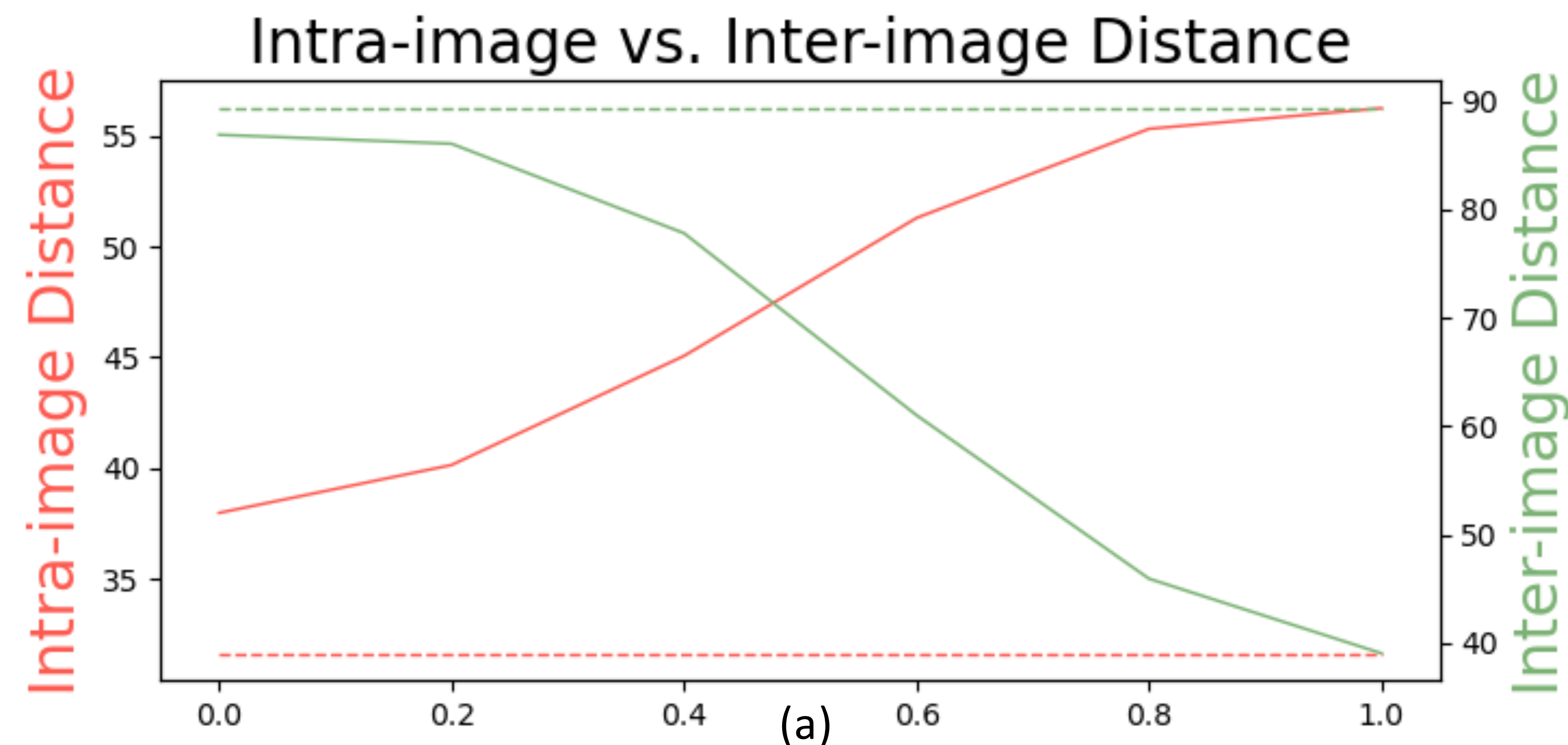
# 3. EXPERIMENTS

## 3.2. Can $\lambda$ control semantic rate ?

**[Exp 1]** With  $\lambda$  increased:

**Observation 1.** Different parts from same person (intra-distance) gets further.

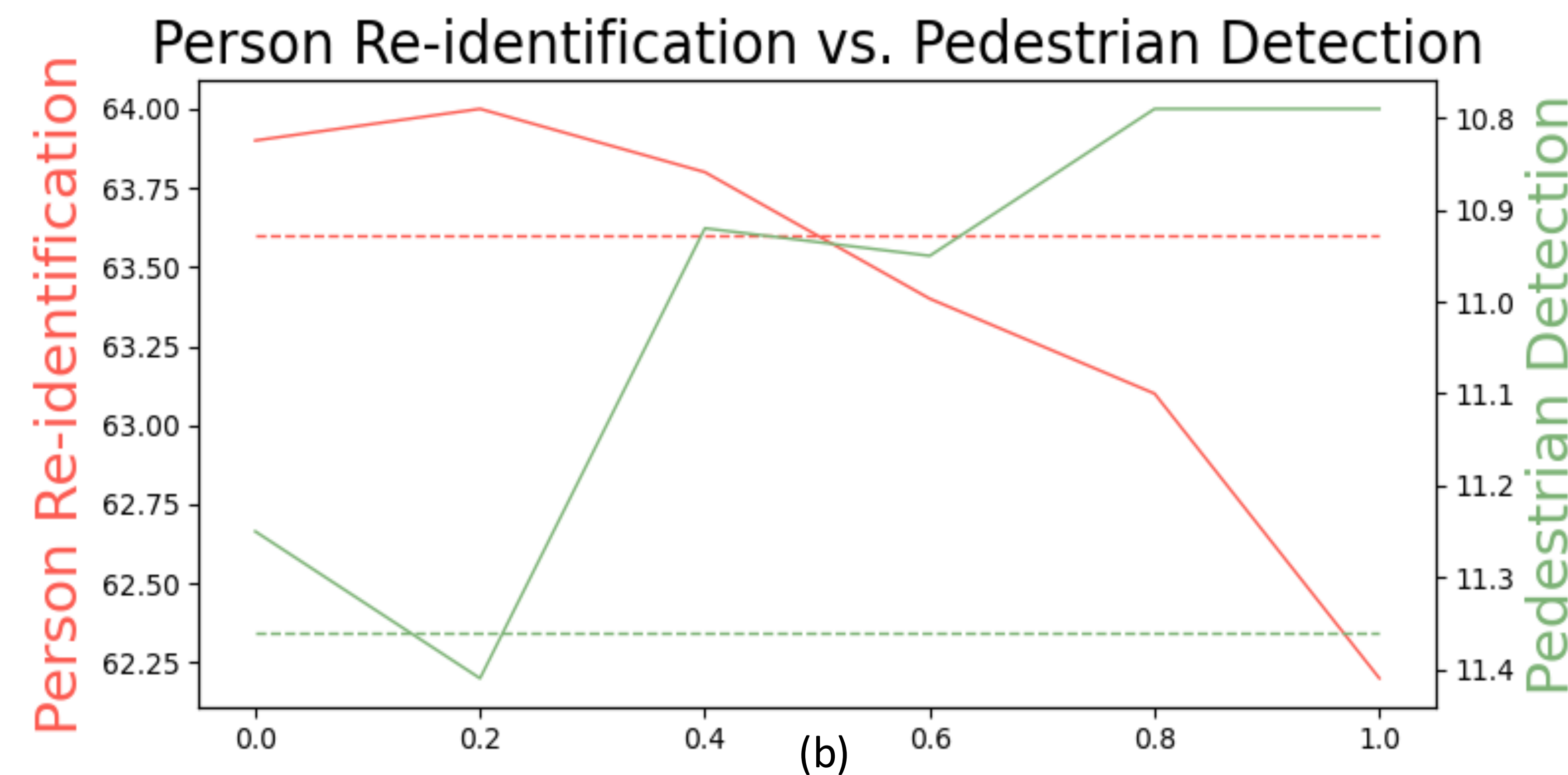
**Observation 2.** Same semantic parts from different people (inter-distance) gets closer.



**[Exp 2]** With  $\lambda$  increased:

**Observation 1.** The person re-identification performance is lower.

**Observation 2.** The pedestrian detection performance is higher.





# 3. EXPERIMENTS

## 3.3. How is Ablation Study and SOTA Comparison?

### Ablation Study

Pretrain Methods		Sup	DINO [6]	+ Clustering	+ Clustering&Controller
Pretrain Data		ImageNet	LUP1M	LUP1M	LUP1M
<b>Person Re-identification</b> mAP/Rank1 ↑ TransReID [36]	Market1501	78.1/90.2	89.6/95.9	89.5/95.5	<b>89.9/96.1</b>
	MSMT17	49.7/73.6	63.3/83.2	61.6/82.2	<b>63.9/83.8</b>
<b>Attribute Recognition</b> mA ↑ RethinkPAR [38]	PETA <sub>zs</sub>	72.86	73.64	73.90	<b>74.20</b>
	RAP <sub>zs</sub>	72.10	73.04	73.16	<b>73.21</b>
	PA100k	80.67	82.98	82.98	<b>84.15</b>
<b>Person Search</b> mAP/Rank1 ↑ SeqNet [45]	CUHK-SYSU	93.0/94.1	93.6/94.3	93.6/94.1	<b>94.0/94.7</b>
	PRW	50.0/84.4	52.9/84.7	53.0/84.0	<b>54.1/85.0</b>
<b>Pedestrian Detection</b> MR <sup>-2</sup> (R/HO) ↓ CSP [73]	CityPerson	11.6/43.8	11.4/43.1	11.1/41.7	<b>10.8/40.7</b>
<b>Human Parsing</b> mIOU ↑ SCHP [43]	LIP	51.10	54.45	55.25	<b>55.45</b>
<b>Pose Estimation</b> AP/AR ↑ HRFormer [83]	COCO	72.4/78.2	73.1/78.5	73.4/78.7	<b>74.4/79.7</b>

### SOTA Comparison

<b>Person Re-identification</b> mAP/Rank1 ↑		SCSN [16]	ABDNet [9]	TransReID [36]	UP-ReID [81]	PASS [97]	SOLIDER		
	Market1501	88.5/95.7	88.3/95.6	89.5/95.2	91.1/97.1	93.3/ <b>96.9</b>	Swin-T	Swin-S	Swin-B
	MSMT17	58.5/83.8	60.8/82.3	69.4/86.2	63.3/84.3	74.3/89.7	91.6/96.1	93.3/96.6	<b>93.9/96.9</b>
<b>Attribute Recognition</b> mA ↑		MsVAA [61]	VAC [30]	ALM [66]	JLAC [65]	RethinkPAR [38]	SOLIDER		
	PETA <sub>zs</sub>	71.53	71.91	73.01	73.60	71.62	Swin-T	Swin-S	Swin-B
	RAP <sub>zs</sub>	72.04	73.70	74.28	76.38	72.32	74.37	76.21	<b>76.43</b>
	PA100k	80.41	79.16	80.68	82.31	81.61	74.23	76.84	<b>77.06</b>
<b>Person Search</b> mAP/Rank1 ↑		NAE+ [7]	AlignPS+ [79]	TCTS [68]	SeqNet [45]	GLCNet [93]	SOLIDER		
	CUHK-SYSU	92.1/92.9	94.0/94.5	93.9/95.1	94.8/95.7	<b>95.8/96.2</b>	Swin-T	Swin-S	Swin-B
	PRW	44.0/81.1	46.1/82.1	46.8/87.5	47.6/87.6	47.8/ <b>87.8</b>	94.9/95.7	95.5/95.8	94.9/95.5
<b>Pedestrian Detection</b> MR <sup>-2</sup> (R/HO) ↓		RepLoss [71]	CSP [73]	NMS-Loss [56]	ACSP [70]	PedesFormer [33]	SOLIDER		
	CityPerson	13.2/56.9	11.0/49.3	10.8/-	9.3/46.3	<b>9.2/36.9</b>	Swin-T	Swin-S	Swin-B
<b>Human Parsing</b> mIOU ↑		JPPNet [46]	BraidNet [50]	CE2P [60]	PCNet [87]	SCHP [43]	SOLIDER		
	LIP	51.37	54.40	53.10	57.03	59.36	Swin-T	Swin-S	Swin-B
<b>Pose Estimation</b> AP/AR ↑		CPN [17]	SimpleBase [75]	TokenPose [44]	HRNet [63]	HRFormer [83]	SOLIDER		
	COCO	68.6/-	74.3/79.7	75.8/80.9	76.3/81.2	<b>77.2/82.0</b>	Swin-T	Swin-S	Swin-B
							74.4/79.6	76.3/81.3	76.6/81.5



## 4. OPEN CODE

Project Website: <https://github.com/tinyvision/SOLIDER>



a **S**emantic **c**Ontrollable **se**lf-supervised **D**l**Ea**Rning framework  
from **Alibaba Group**

- State of the Art Pedestrian Attribute Recognition on PA-100K
- State of the Art Person Re-Identification on MSMT17 (using additional training data)
- Ranked #3 Person Re-Identification on Market-1501 (using additional training data) Ranked #4 Person Search on CUHK-SYSU
- State of the Art Person Search on PRW Ranked #6 Pedestrian Detection on CityPersons (using additional training data)
- State of the Art Semantic Segmentation on LIP val Ranked #2 Pose Estimation on COCO

### Updates

- [2023/04/24: Codes of attribute recognition task is released!] **NEW**
  - Training details of our pretrained model on downstream person attribute recognition task is released.
- [2023/03/28: Codes of 3 downstream tasks are released!]
  - Training details of our pretrained model on 3 downstream human visual tasks, including person re-identification, person search and pedestrian detection, are released.
- [2023/03/13: SOLIDER is accepted by CVPR2023!]
  - The paper of SOLIDER is accepted by CVPR2023, and its official pytorch implementation is released in this repo.