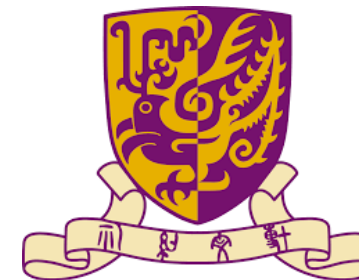


# Siamese Image Modeling for Self-Supervised Vision Representation Learning

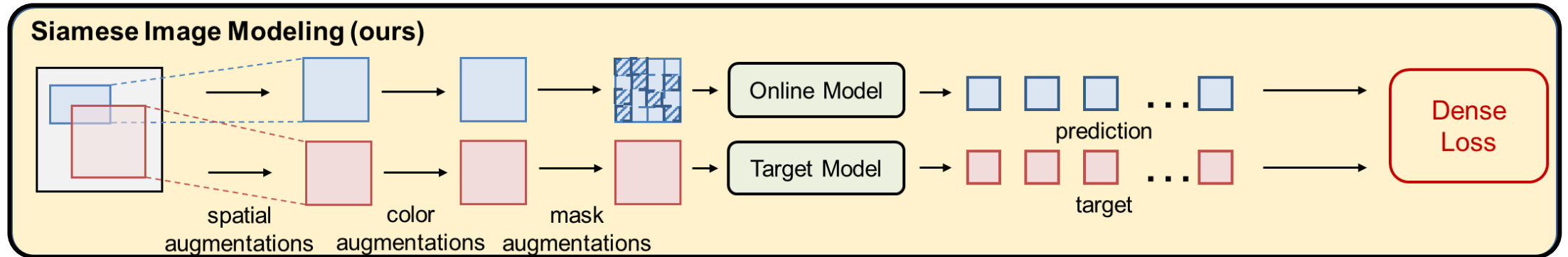
Chenxin Tao<sup>1\*</sup>, Xizhou Zhu<sup>2\*</sup>, Weijie Su<sup>3\*</sup>,  
Gao Huang<sup>1</sup>, Bin Li<sup>3</sup>, Jie Zhou<sup>1</sup>, Yu Qiao<sup>4</sup>, Xiaogang Wang<sup>5</sup>, Jifeng Dai<sup>1,4</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>SenseTime Research, <sup>3</sup>University of Science and Technology of China,  
<sup>4</sup>Shanghai Artificial Intelligence Laboratory, <sup>5</sup>The Chinese University of Hong Kong

TUE-AM-204



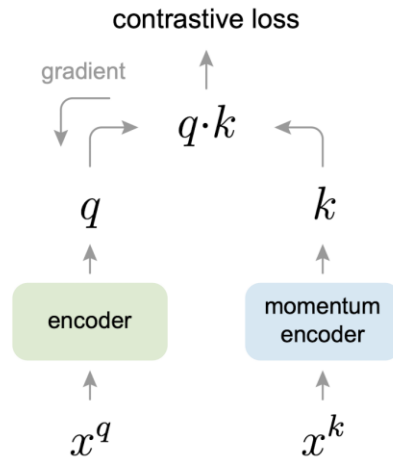
# Contribution



- A new form of self-supervised learning that can learn **semantic alignment** and **spatial sensitivity** with a single dense loss.
- Compared with MIM methods, reconstructing another view helps to obtain good semantic alignment.
- Compared with ID methods, matching the dense correspondence between two views can help to learn spatial sensitivity.
- SiameseIM is able to surpass both MIM and ID methods over a wide range of tasks. SiameseIM obtains more improvements in few-shot, long-tail and robustness-concerned scenarios.

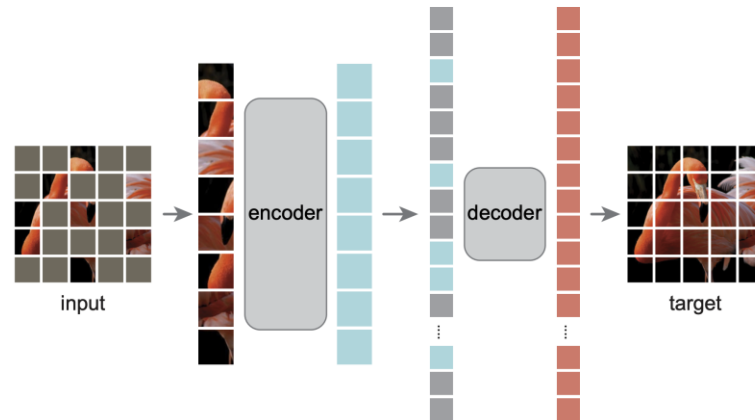
# Motivation – Two Mainstream SSL Methods

Instance  
Discrimination



- lack **spatial sensitivity**: modeling the local structure within an image

Masked  
Image Modeling

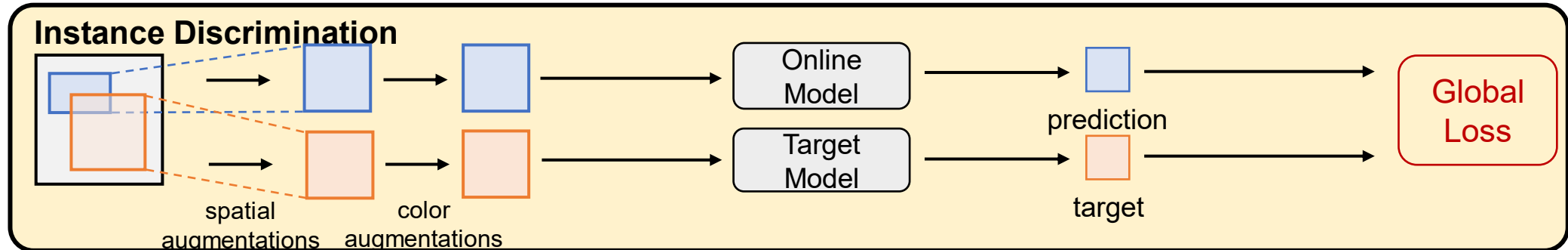


- lack **semantic alignment**: projecting semantically similar views into nearby representations

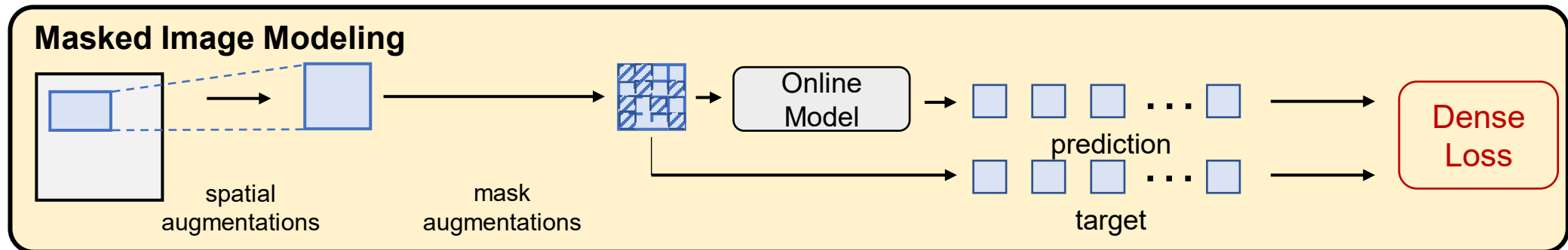
Typical Work

Characteristics

# Motivation – Two Mainstream SSL Methods



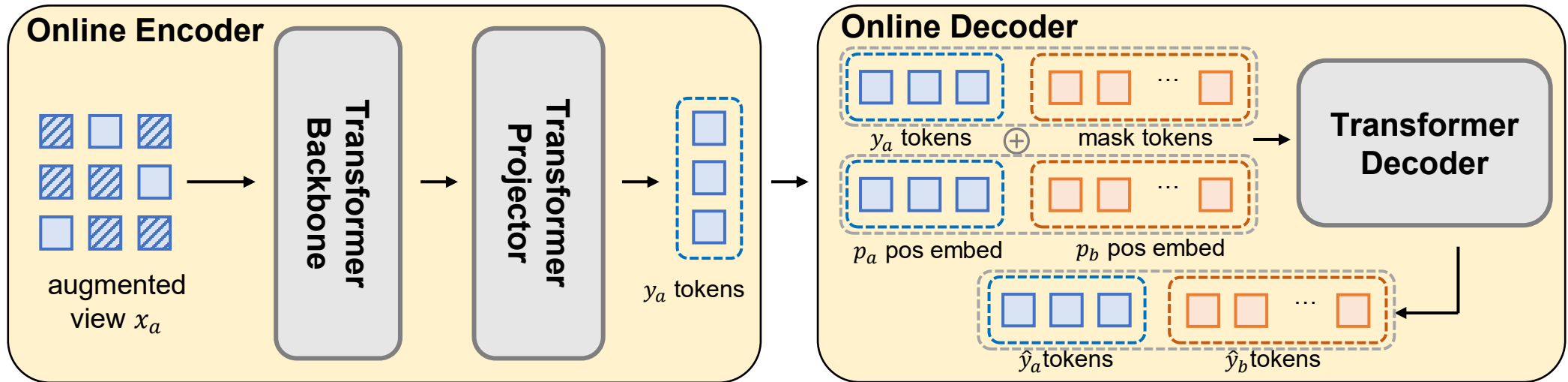
semantic alignment can be achieved by **matching different augmented view from the same image**



spatial sensitivity can be achieved by **predicting dense representations from masked images**

- We propose Siamese Image Modeling, which reconstructs the dense representations of an augmented view, based on another augmented view from the same image.

# Siamese Image Modeling



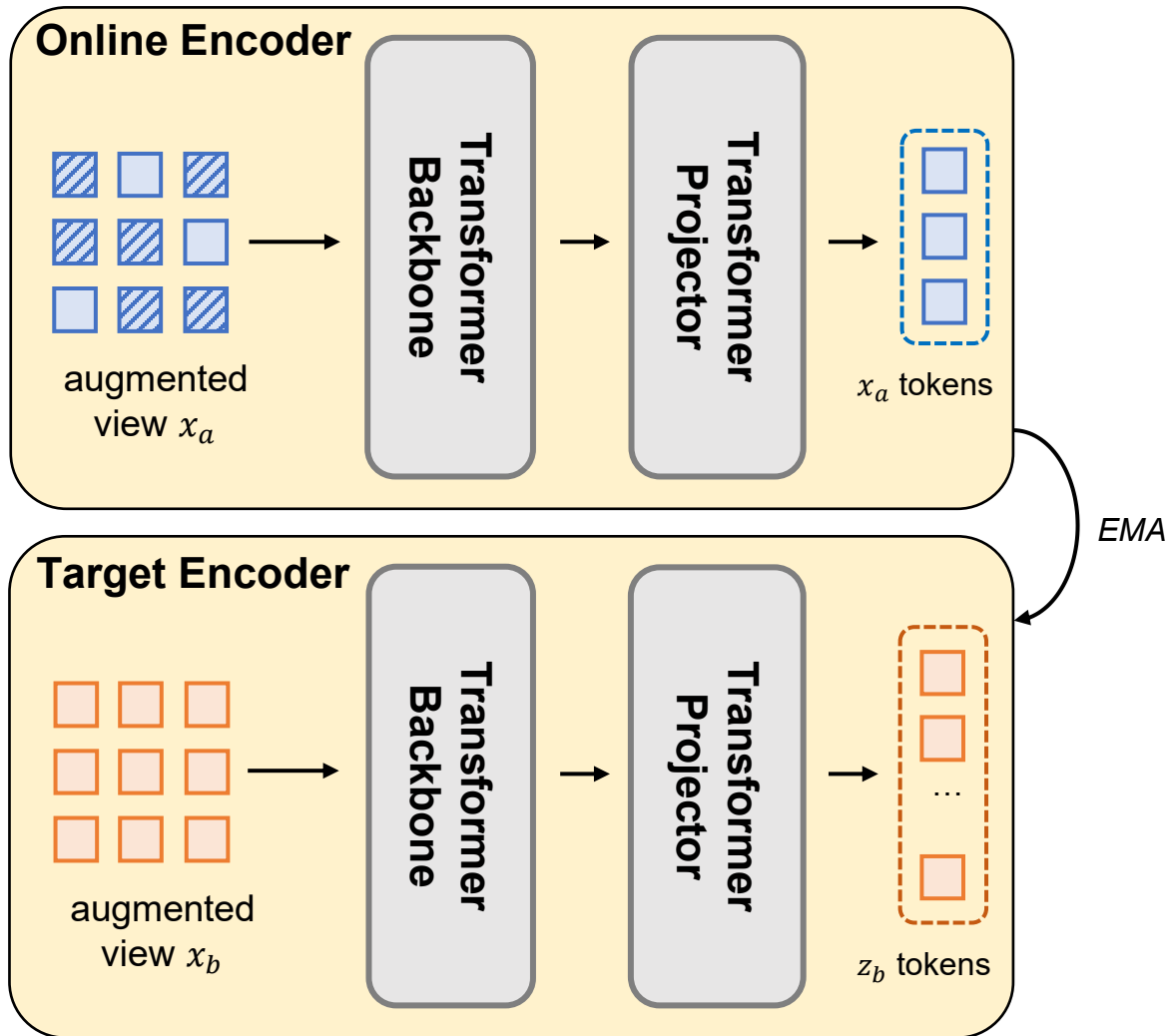
➤ Augmentation: spatial + color augmentations for both views, mask augmentation for online view

➤ Online branch makes the prediction:

$$y_b = g \left( \text{Concat} \left( y_a + p_a, \left\{ m + p_b^{(u,v)} \right\}_{u=1,v=1}^{N_h, N_w} \right) \right)$$

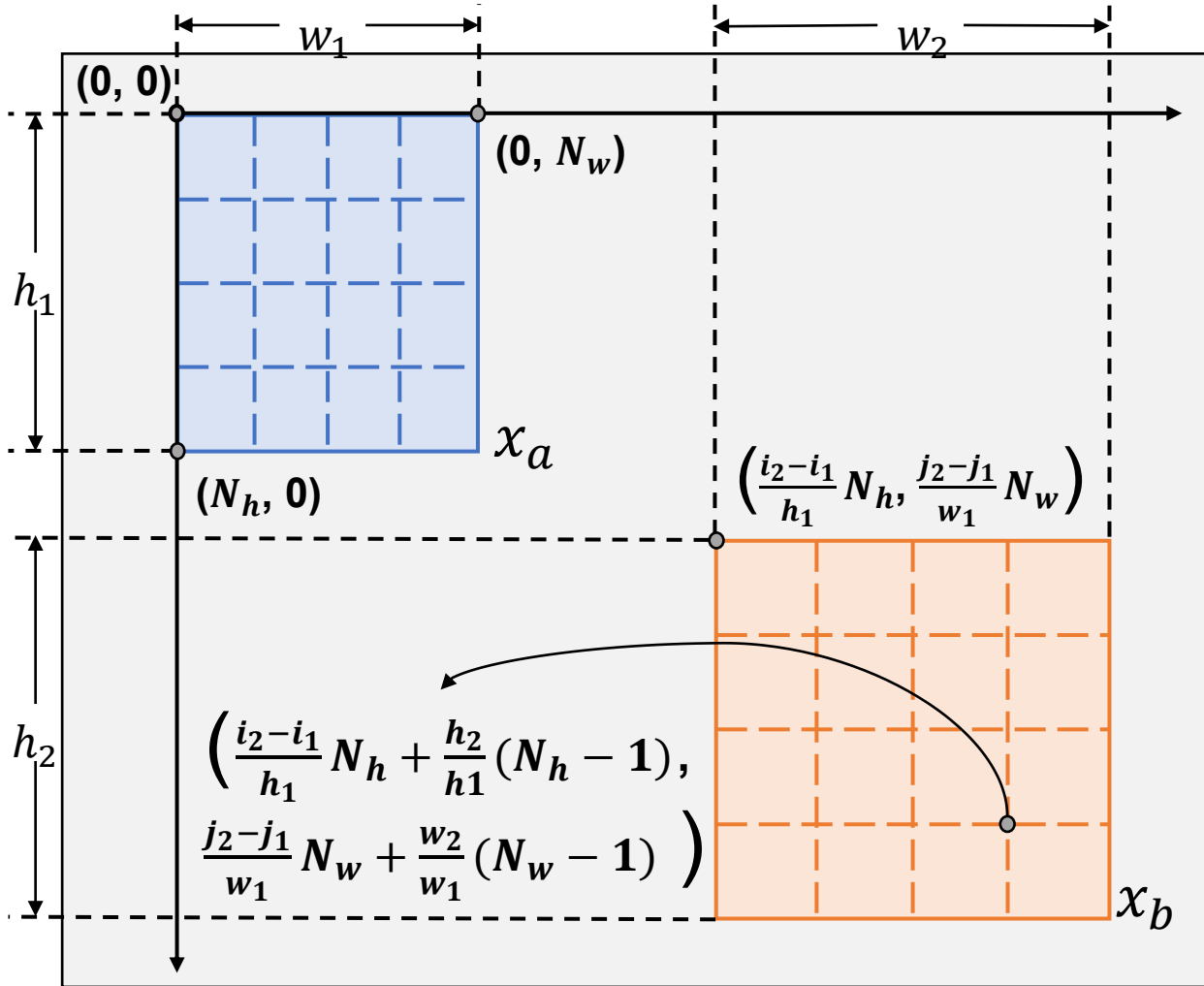
➤  $m$  indicates mask token,  $p_a, p_b$  are the positional embedding of view  $x_a$  and  $x_b$

# Siamese Image Modeling



- Target branch computes the target representation  $z_b$

# Siamese Image Modeling

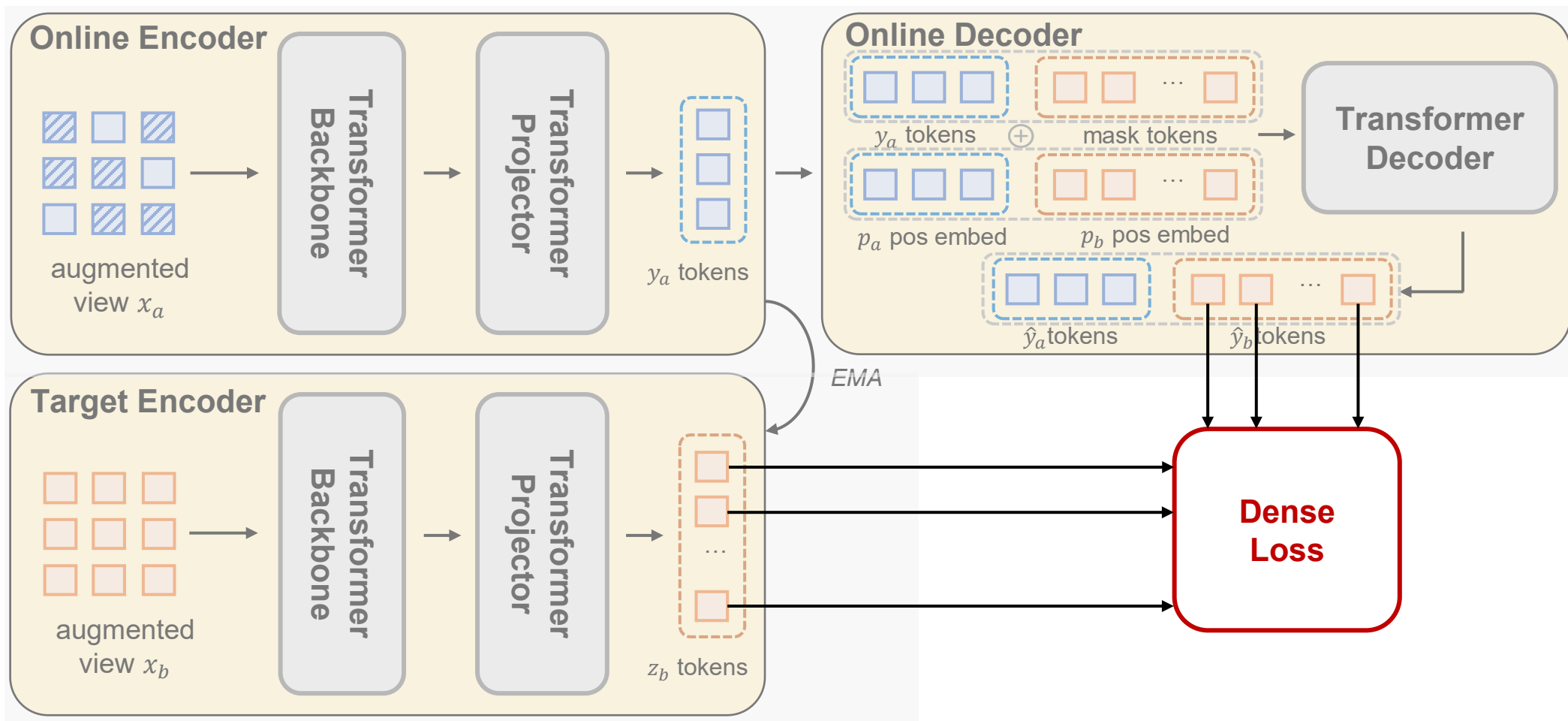


➤ Positional Embedding

$$\tilde{p}_a^{(u,v)} = (u - 1, v - 1)$$

$$\tilde{p}_b^{(u,v)} = \left( \frac{h_2}{h_1} (u - 1) + \frac{i_2 - i_1}{h_1} N_h, \frac{w_2}{w_1} (v - 1) + \frac{j_2 - j_1}{w_1} N_w \right)$$

# Siamese Image Modeling



➤ Dense Loss

$$L = \mathbb{E}_{\{y_b^i, z_b^i\}} \left[ -\|y_b^i - z_b^i\|^2 + \lambda \sum_{u \in N} (u^T y_b^i)^2 \right]$$



# Experiments – Main Results

Method	Epochs	ImageNet		
		FT	LIN	FT <sub>1%</sub>
Supervised	300	81.8	-	-
DINO*	800 <sup>†</sup>	82.8	78.2	-
iBOT*	1600 <sup>†</sup>	84.0	79.5	-
DenseCL <sup>‡</sup>	400	82.2	69.7	49.9
MoCo-v3	600 <sup>†</sup>	83.0	76.7	63.4
BEiT	800	83.2	-	-
MAE	400	83.1	62.5	-
MAE	1600	83.6	68.0	51.1
SiameseIM	400	83.7	76.8	61.8
<b>SiameseIM</b>	1600	<b>84.1 (+0.5)</b>	<b>78.0 (+1.3)</b>	<b>65.1 (+1.7)</b>

(a) Image classification.

Method	Epochs	AP <sup>b</sup>	AP <sup>b</sup> <sub>rare</sub>	AP <sup>m</sup>	AP <sup>m</sup> <sub>rare</sub>
Supervised	300	37.2	-	34.9	26.4
iBoT*	1600	36.9	29.1	34.6	28.9
DenseCL <sup>‡</sup>	400	33.8	25.1	32.1	24.6
MoCo-v3	600 <sup>†</sup>	37.3	25.5	35.3	25.8
MAE	400	38.4	25.4	36.6	25.7
MAE	1600	40.1	29.3	<b>38.1</b>	29.1
SiameseIM	400	38.5	28.9	36.1	27.7
<b>SiameseIM</b>	1600	<b>40.5 (+0.4)</b>	<b>30.9 (+1.6)</b>	<b>38.1 (+0.0)</b>	<b>30.1 (+1.0)</b>

(c) Long-tail object detection on LVIS.

Method	Epochs	COCO		ADE20k
		AP <sup>b</sup>	AP <sup>m</sup>	mIoU
Supervised	300	47.9	42.9	47.4
DINO*	800 <sup>†</sup>	50.1	43.4	46.8
iBOT*	1600 <sup>†</sup>	51.2	44.2	50.0
DenseCL <sup>‡</sup>	400	46.6	41.6	44.5
MoCo-v3	600 <sup>†</sup>	47.9	42.7	47.3
BEiT	800	49.8	44.4	47.1
MAE	400	50.6	45.1	45.0
MAE	1600	51.6	45.9	48.1
SiameseIM	400	50.7	44.9	49.6
<b>SiameseIM</b>	1600	<b>52.1 (+0.5)</b>	<b>46.2 (+0.3)</b>	<b>51.1 (+3.0)</b>

(b) Common object detection and semantic segmentation.

Method	Epochs	IN-A top-1	IN-R top-1	IN-Sketch top-1	IN-C 1-mCE
MSN*	1200 <sup>†</sup>	37.5	50.0	36.3	53.4
iBoT*	1600 <sup>†</sup>	42.4	50.9	36.9	55.5
DenseCL <sup>‡</sup>	400	30.8	43.8	29.9	48.1
MoCo-v3	600 <sup>†</sup>	32.4	49.8	35.9	55.4
MAE	1600	35.9	48.3	34.5	48.3
SiameseIM	400	38.6	51.6	37.7	55.9
<b>SiameseIM</b>	1600	<b>43.8 (+7.9)</b>	<b>52.5 (+2.7)</b>	<b>38.3 (+2.4)</b>	<b>57.1 (+1.7)</b>

(d) Robustness evaluation.

- SiameseIM is able to surpass both MIM and ID methods over a wide range of tasks.

# Experiments – Main Results

Method	Epochs	ImageNet		
		FT	LIN	FT <sub>1%</sub>
Supervised	300	81.8	-	-
DINO*	800 <sup>†</sup>	82.8	78.2	-
iBOT*	1600 <sup>†</sup>	84.0	79.5	-
DenseCL <sup>‡</sup>	400	82.2	69.7	49.9
MoCo-v3	600 <sup>†</sup>	83.0	76.7	63.4
BEiT	800	83.2	-	-
MAE	400	83.1	62.5	-
MAE	1600	83.6	68.0	51.1
SiameseIM	400	83.7	76.8	61.8
<b>SiameseIM</b>	1600	<b>84.1 (+0.5)</b>	<b>78.0 (+1.3)</b>	<b>65.1 (+1.7)</b>

(a) Image classification.

Method	Epochs	AP <sup>b</sup>	AP <sup>b</sup> <sub>rare</sub>	AP <sup>m</sup>	AP <sup>m</sup> <sub>rare</sub>
Supervised	300	37.2	-	34.9	26.4
iBoT*	1600	36.9	29.1	34.6	28.9
DenseCL <sup>‡</sup>	400	33.8	25.1	32.1	24.6
MoCo-v3	600 <sup>†</sup>	37.3	25.5	35.3	25.8
MAE	400	38.4	25.4	36.6	25.7
MAE	1600	40.1	29.3	<b>38.1</b>	29.1
SiameseIM	400	38.5	28.9	36.1	27.7
<b>SiameseIM</b>	1600	<b>40.5 (+0.4)</b>	<b>30.9 (+1.6)</b>	<b>38.1 (+0.0)</b>	<b>30.1 (+1.0)</b>

(c) Long-tail object detection on LVIS.

Method	Epochs	COCO		ADE20k
		AP <sup>b</sup>	AP <sup>m</sup>	mIoU
Supervised	300	47.9	42.9	47.4
DINO*	800 <sup>†</sup>	50.1	43.4	46.8
iBOT*	1600 <sup>†</sup>	51.2	44.2	50.0
DenseCL <sup>‡</sup>	400	46.6	41.6	44.5
MoCo-v3	600 <sup>†</sup>	47.9	42.7	47.3
BEiT	800	49.8	44.4	47.1
MAE	400	50.6	45.1	45.0
MAE	1600	51.6	45.9	48.1
SiameseIM	400	50.7	44.9	49.6
<b>SiameseIM</b>	1600	<b>52.1 (+0.5)</b>	<b>46.2 (+0.3)</b>	<b>51.1 (+3.0)</b>

(b) Common object detection and semantic segmentation.

Method	Epochs	IN-A top-1	IN-R top-1	IN-Sketch top-1	IN-C 1-mCE
MSN*	1200 <sup>†</sup>	37.5	50.0	36.3	53.4
iBoT*	1600 <sup>†</sup>	42.4	50.9	36.9	55.5
DenseCL <sup>‡</sup>	400	30.8	43.8	29.9	48.1
MoCo-v3	600 <sup>†</sup>	32.4	49.8	35.9	55.4
MAE	1600	35.9	48.3	34.5	48.3
SiameseIM	400	38.6	51.6	37.7	55.9
<b>SiameseIM</b>	1600	<b>43.8 (+7.9)</b>	<b>52.5 (+2.7)</b>	<b>38.3 (+2.4)</b>	<b>57.1 (+1.7)</b>

(d) Robustness evaluation.

- SiameseIM obtains more improvements in few-shot, long-tail and robustness-concerned scenarios.

# Experiments – Ablation Study

	target type	different views	color aug	mask type	BN/LN in proj & dec	loss norm*	loss type	loss form	FT	LIN	AP <sup>b</sup>	AP <sup>m</sup>
<i>single view with dense loss:</i>												
MAE	pixel			random	LN	MAE-like	dense	L2	83.1	62.5	46.8	42.0
(a)	pixel			random	LN	MAE-like	dense	L2	82.8	62.3	47.3	42.5
(b)	feature			random	LN	MoCo-like	dense	UniGrad	81.0	48.7	43.5	39.2
(c)	pixel		✓	random	LN	MAE-like	dense	L2	82.0	59.9	46.3	41.8
<i>multiple views with dense loss:</i>												
(d)	pixel	✓		random	LN	MAE-like	dense	L2	78.7	46.2	38.1	34.8
(e)	feature	✓		random	LN	MoCo-like	dense	UniGrad	82.9	69.6	48.5	43.4
(f)	feature	✓	✓	random	LN	MoCo-like	dense	UniGrad	83.0	73.1	47.9	43.2
(g)	feature	✓	✓	random	BN	MoCo-like	dense	UniGrad	83.2	73.6	48.7	43.7
(h)	feature	✓	✓	blockwise	BN	MoCo-like	dense	UniGrad	83.5	74.7	50.0	44.5
(i)	feature	✓	✓	blockwise	BN	MAE-like	dense	UniGrad	83.7	76.8	49.8	44.2
(j)	feature	✓	✓	blockwise	BN	MAE-like	dense	L2	83.3	76.5	49.8	44.2
<i>multiple views with global loss:</i>												
(k)	feature	✓	✓	random	BN	MoCo-like	global	UniGrad	82.7	72.0	45.9	41.4
MoCo-v3 with mask	feature	✓	✓	random	BN	MoCo-like	global	UniGrad	82.8	72.2	45.0	40.5

- Compared with MIM methods, reconstructing another view helps to obtain good semantic alignment.

# Experiments – Ablation Study

	target type	different views	color aug	mask type	BN/LN in proj & dec	loss norm*	loss type	loss form	FT	LIN	AP <sup>b</sup>	AP <sup>m</sup>
<i>single view with dense loss:</i>												
MAE	pixel			random	LN	MAE-like	dense	L2	83.1	62.5	46.8	42.0
(a)	pixel			random	LN	MAE-like	dense	L2	82.8	62.3	47.3	42.5
(b)	feature			random	LN	MoCo-like	dense	UniGrad	81.0	48.7	43.5	39.2
(c)	pixel		✓	random	LN	MAE-like	dense	L2	82.0	59.9	46.3	41.8
<i>multiple views with dense loss:</i>												
(d)	pixel	✓		random	LN	MAE-like	dense	L2	78.7	46.2	38.1	34.8
(e)	feature	✓		random	LN	MoCo-like	dense	UniGrad	82.9	69.6	48.5	43.4
(f)	feature	✓	✓	random	LN	MoCo-like	dense	UniGrad	83.0	73.1	47.9	43.2
(g)	feature	✓	✓	random	BN	MoCo-like	dense	UniGrad	83.2	73.6	48.7	43.7
(h)	feature	✓	✓	blockwise	BN	MoCo-like	dense	UniGrad	83.5	74.7	50.0	44.5
(i)	feature	✓	✓	blockwise	BN	MAE-like	dense	UniGrad	83.7	76.8	49.8	44.2
(j)	feature	✓	✓	blockwise	BN	MAE-like	dense	L2	83.3	76.5	49.8	44.2
<i>multiple views with global loss:</i>												
(k)	feature	✓	✓	random	BN	MoCo-like	global	UniGrad	82.7	72.0	45.9	41.4
MoCo-v3 with mask	feature	✓	✓	random	BN	MoCo-like	global	UniGrad	82.8	72.2	45.0	40.5

- Compared with ID methods, dense supervision can improve the spatial sensitivity.

# Siamese Image Modeling for Self-Supervised Vision Representation Learning

Thanks for watching!

Contact us:

[tcx20@mails.tsinghua.edu.cn](mailto:tcx20@mails.tsinghua.edu.cn)

