

## Lecture 08: Data-efficient Training





- Networks are trained in a supervised fashion jointly with labeled and unlabeled data.

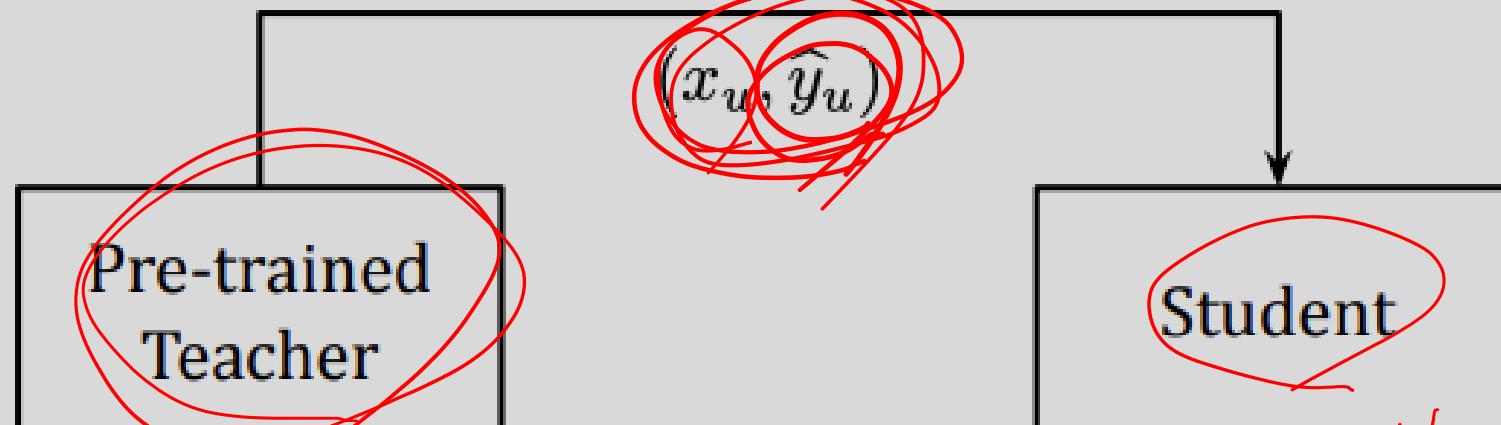


- Pseudo-Labels are target classes for unlabeled data predicted from another network as if they were true labels.

# Pseudo labels

Self-training

Pseudo-labeled data



$\theta_S^{\text{PL}} = \underset{\theta_S}{\operatorname{argmin}} \underbrace{\mathbb{E}_{x_u} [\text{CE}(T(x_u; \theta_T), S(x_u; \theta_S))]}_{:= \mathcal{L}_u(\theta_T, \theta_S)}$

$S(\cdot)$ .



Self-training with noisy student improves imagenet classification, CVPR'20.

# Self-training

**Require:** Labeled images  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  and unlabeled images  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_m\}$ .

- 1: Learn teacher model  $\theta_*^t$  which minimizes the cross entropy loss on labeled images

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, f^{noised}(x_i, \theta_*^t))$$

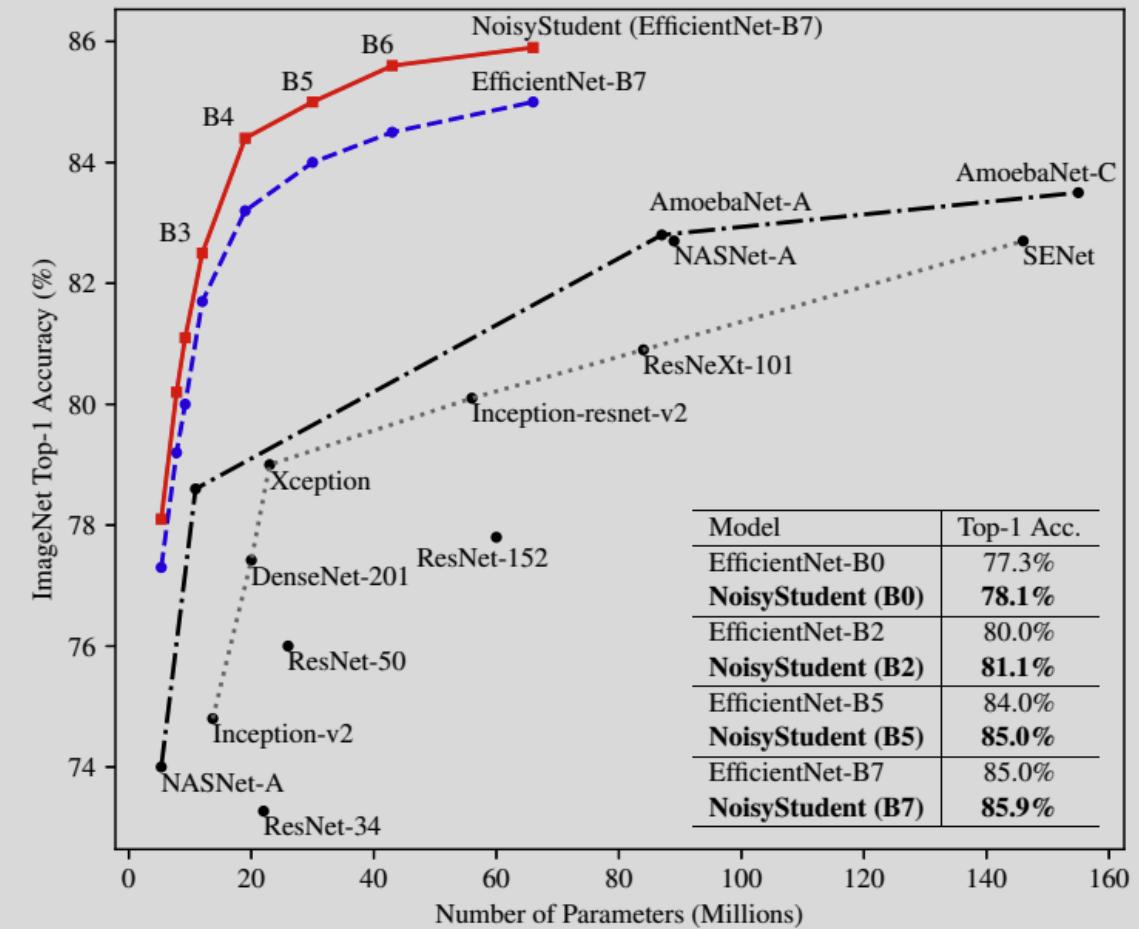
- 2: Use an unnoised teacher model to generate soft or hard pseudo labels for unlabeled images

$$\tilde{y}_i = f(\tilde{x}_i, \theta_*^t), \forall i = 1, \dots, m$$

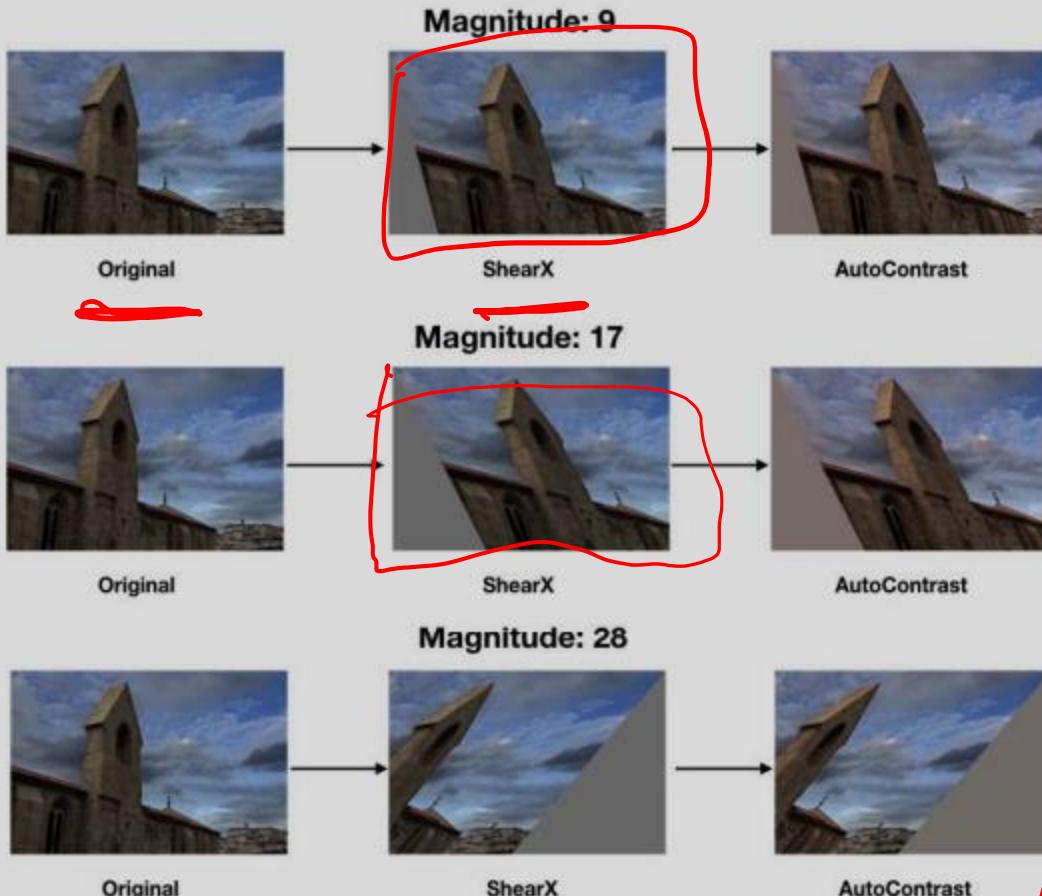
- 3: Learn an **equal-or-larger** student model  $\theta_*^s$  which minimizes the cross entropy loss on labeled images and unlabeled images with **noise** added to the student model

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, f^{noised}(x_i, \theta_*^s)) + \frac{1}{m} \sum_{i=1}^m \ell(\tilde{y}_i, f^{noised}(\tilde{x}_i, \theta_*^s))$$

- 4: Iterative training: Use the student as a teacher and go back to step 2.



# Self-training (Noise)



---

```
transforms = [  
    'Identity', 'AutoContrast', 'Equalize', 'Rotate',  
    'Solarize', 'Color',  
    'Posterize', 'Contrast', 'Brightness', 'Sharpness',  
    'ShearX', 'ShearY',  
    'TranslateX', 'TranslateY']  
  
def randaugment(N, M):  
    """Generate a set of distortions.  
  
    Args:  
        N: Number of augmentation transformations to apply  
            sequentially.  
        M: Magnitude for all the transformations.  
    """  
  
    sampled_ops = np.random.choice(transforms, N)  
    return [(op, M) for op in sampled_ops]
```

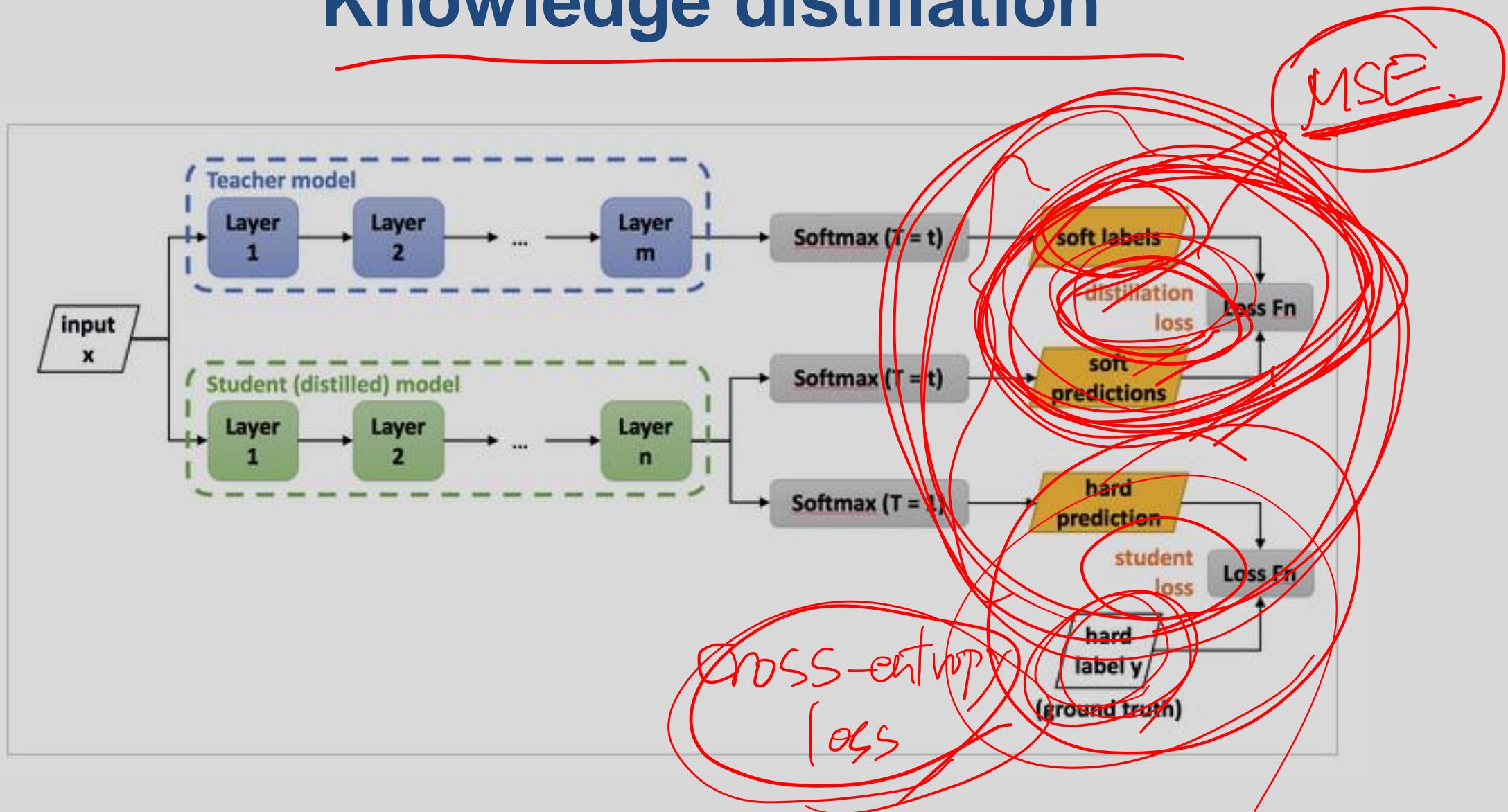
---

# Knowledge distillation

(KD),

- The teacher network provides a richer supervisory signal than the data supervision.
- KD guides the training of a student network by encouraging it to mimic some aspect of a teacher network.

# Knowledge distillation



# Knowledge distillation

ground truth

cow	dog	cat	car
0	1	0	0

hard label.

'Dark knowledge'

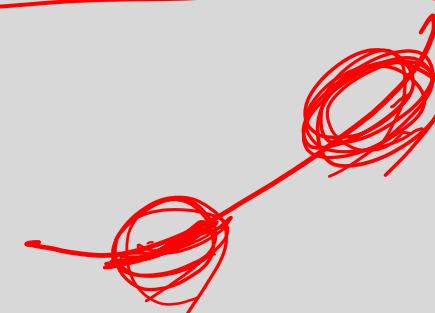
Which classes the teacher found more similar to the predicted class.

soft label. (prob.)

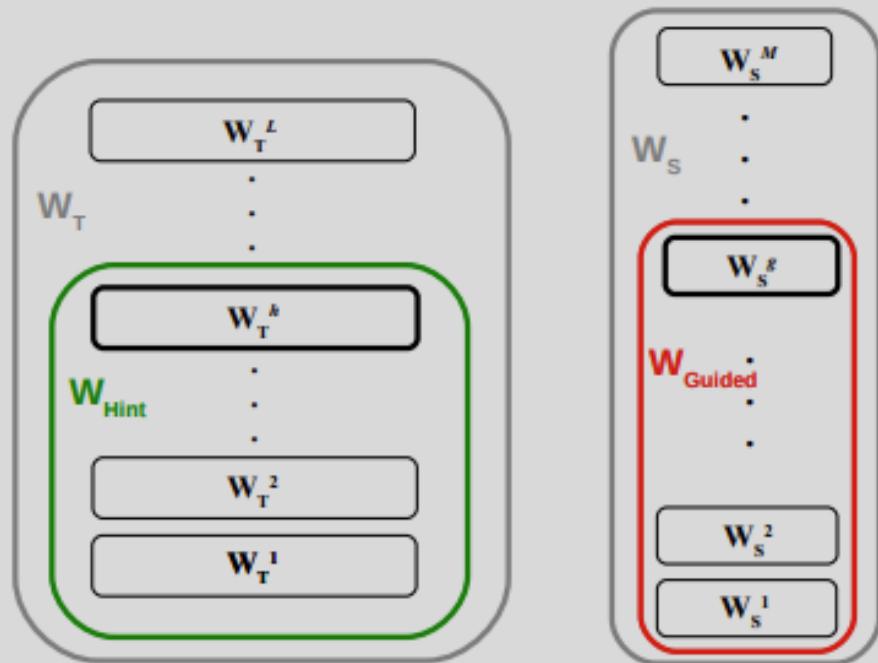
$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

q : probability  
T : temperature

Softmax.



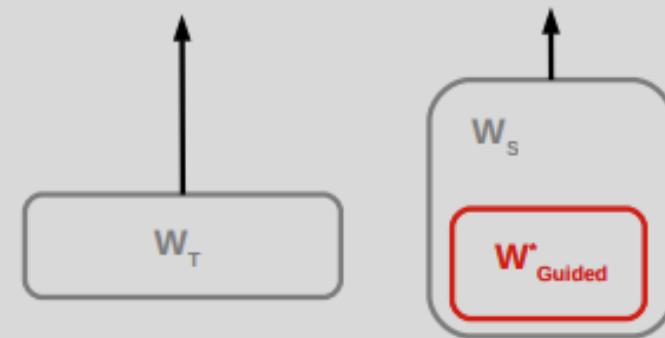
# Knowledge distillation



$$P_T^\tau = \text{softmax} \left( \frac{\mathbf{a}_T}{\tau} \right), \quad P_S^\tau = \text{softmax} \left( \frac{\mathbf{a}_S}{\tau} \right)$$

(Red annotations: circled  $\tau$  and wavy lines under the equations)

$$\mathbf{W}_S^* = \underset{\mathbf{W}_S}{\operatorname{argmin}} \mathcal{L}_{DK}(\mathbf{W}_S)$$



$$\mathcal{L}_{KD}(\mathbf{W}_S) = \mathcal{H}(y_{true}, P_S) + \lambda \mathcal{H}(P_T^\tau, P_S^\tau)$$

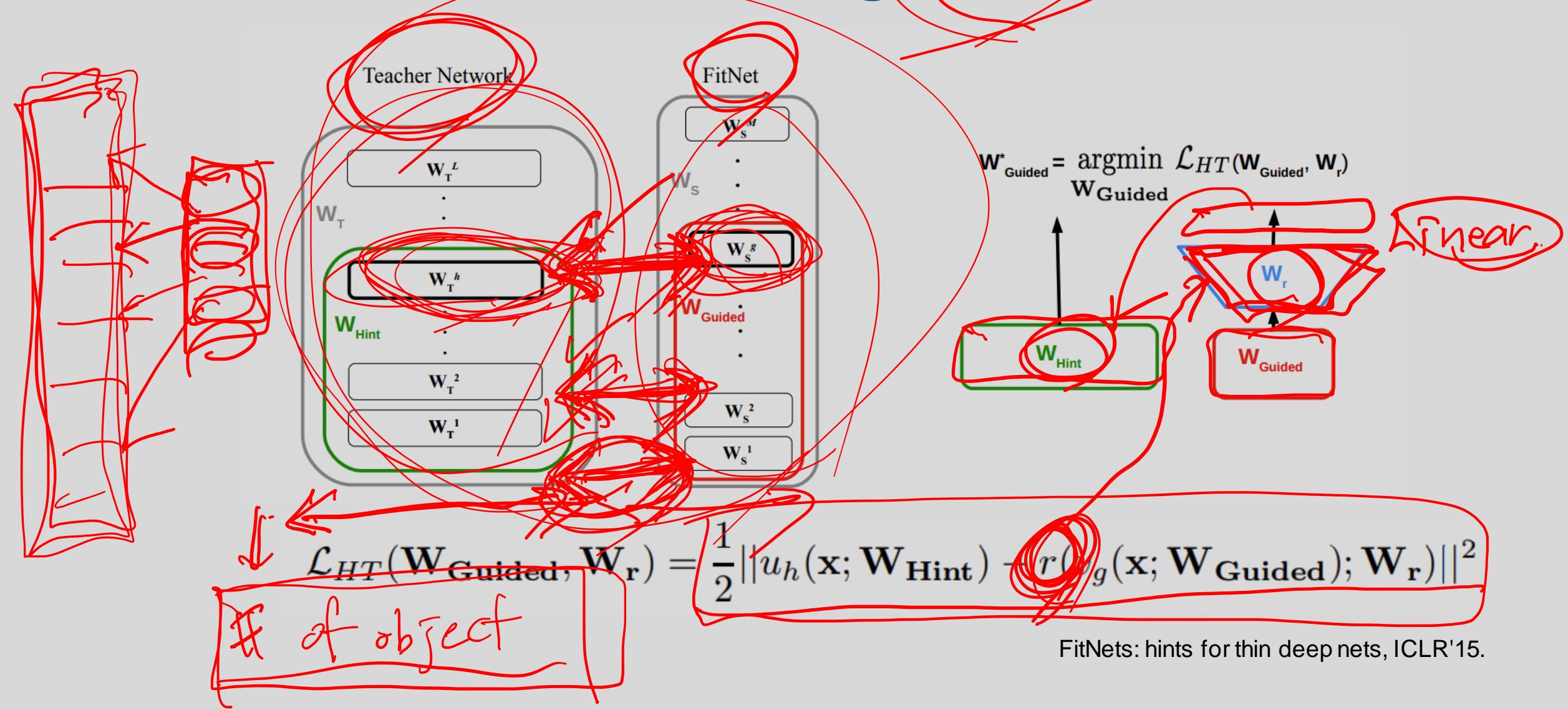
(Red annotations: circled  $\mathcal{H}$ , circled  $y_{true}$ , circled  $P_S$ , circled  $P_T^\tau$ , circled  $P_S^\tau$ , and the text 'KD Loss.' written above the equation)

Distilling the Knowledge in a Neural Network, NeurIPS'14

## Hints

- KD fails when the depth of the student network getting deeper.
- *Hint* is defined as the output of a teacher's hidden layer.

# Learning hints



FitNets: hints for thin deep nets, ICLR'15.

# Learning hints

**Input:**  $\mathbf{W}_S, \mathbf{W}_T, g, h$

**Output:**  $\mathbf{W}_S^*$

- 1:  $\mathbf{W}_{\text{Hint}} \leftarrow \{\mathbf{W}_T^1, \dots, \mathbf{W}_T^h\}$
- 2:  $\mathbf{W}_{\text{Guided}} \leftarrow \{\mathbf{W}_S^1, \dots, \mathbf{W}_S^g\}$
- 3: Initialize  $\mathbf{W}_r$  to small random values
- 4:  $\mathbf{W}_{\text{Guided}}^* \leftarrow \underset{\mathbf{W}_{\text{Guided}}}{\operatorname{argmin}} \mathcal{L}_{HT}(\mathbf{W}_{\text{Guided}}, \mathbf{W}_r)$
- 5:  $\{\mathbf{W}_S^1, \dots, \mathbf{W}_S^g\} \leftarrow \{\mathbf{W}_{\text{Guided}}^{*1}, \dots, \mathbf{W}_{\text{Guided}}^{*g}\}$
- 6:  $\mathbf{W}_S^* \leftarrow \underset{\mathbf{W}_S}{\operatorname{argmin}} \mathcal{L}_{KD}(\mathbf{W}_S)$

# Intermediate representation

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	91.61%
Teacher	~9M	90.18%
Mimic single	~54M	84.6%
Mimic single	~70M	84.9%
Mimic ensemble	~70M	85.8%
<i>State-of-the-art methods</i>		
Maxout		90.65%
Network in Network		91.2%
Deeply-Supervised Networks		91.78%
Deeply-Supervised Networks (19)		88.2%

Table 1: Accuracy on CIFAR-10

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	64.96%
Teacher	~9M	63.54%
<i>State-of-the-art methods</i>		
Maxout		61.43%
Network in Network		64.32%
Deeply-Supervised Networks		65.43%

Table 2: Accuracy on CIFAR-100

# Intermediate representation

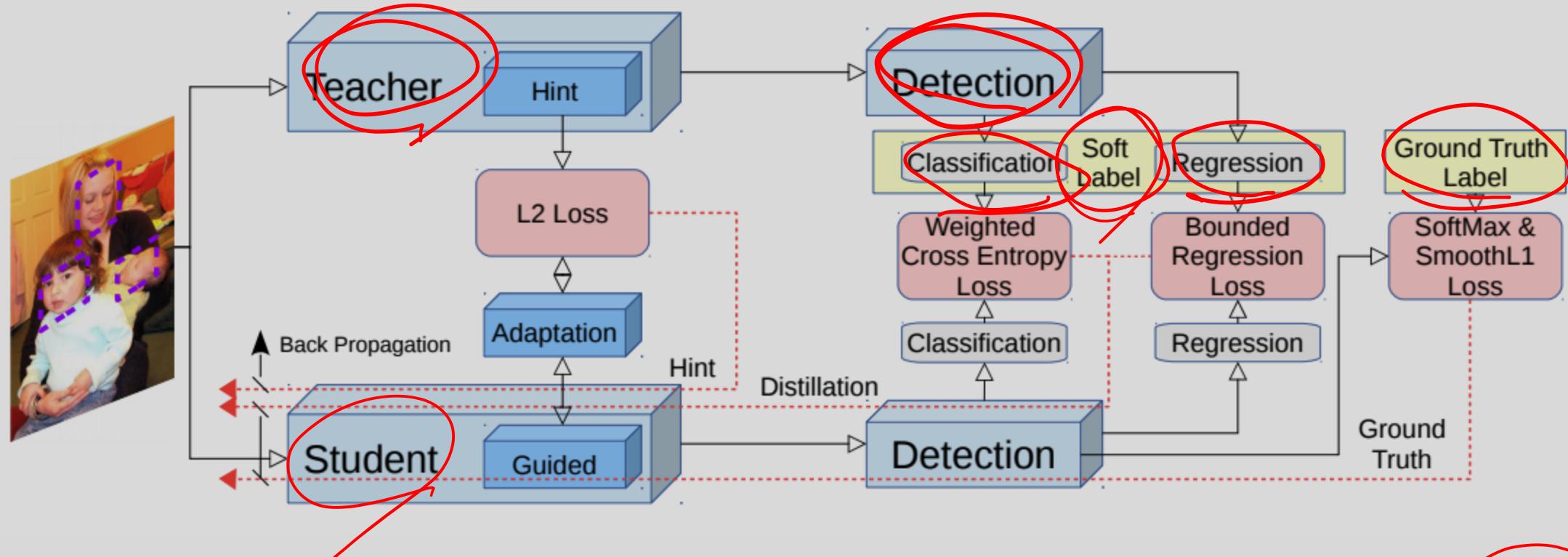
Algorithm	# params	Misclass
<i>Compression</i>		
FitNet	~1.5M	2.42%
Teacher	~4.9M	2.38%
<i>State-of-the-art methods</i>		
Maxout		2.47%
Network in Network		2.35%
Deeply-Supervised Networks		1.92%

Table 3: SVHN error

Algorithm	# params	Misclass
<i>Compression</i>		
Teacher	~361K	0.55%
Standard backprop	~30K	1.9%
KD	~30K	0.65%
FitNet	~30K	0.51%
<i>State-of-the-art methods</i>		
Maxout		0.45%
Network in Network		0.47%
Deeply-Supervised Networks		0.39%

Table 4: MNIST error

# KD for object detection



Learning Efficient Object Detection Models with Knowledge Distillation, NIPS'17.

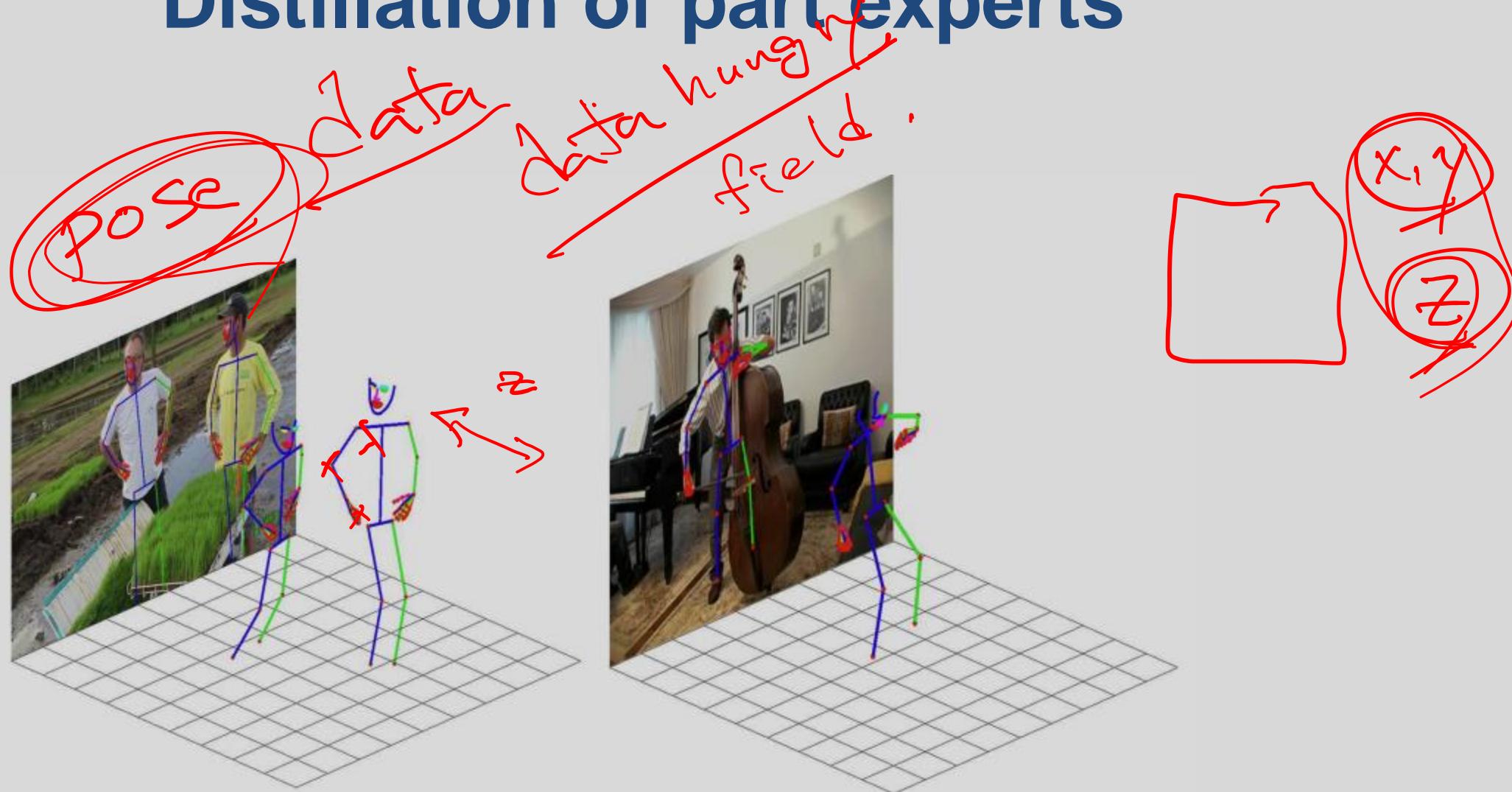
# KD for object detection

		Baseline	Distillation	Hint	Distillation + Hint
PASCAL	Trainval	79.6	78.3	80.9	83.5
	Test	54.7	58.4	58	59.4
COCO	Train	45.3	45.4	47.1	49.6
	Val	25.4	26.1	27.8	28.3

learning on different datasets with Tucker and VGG16 pair.

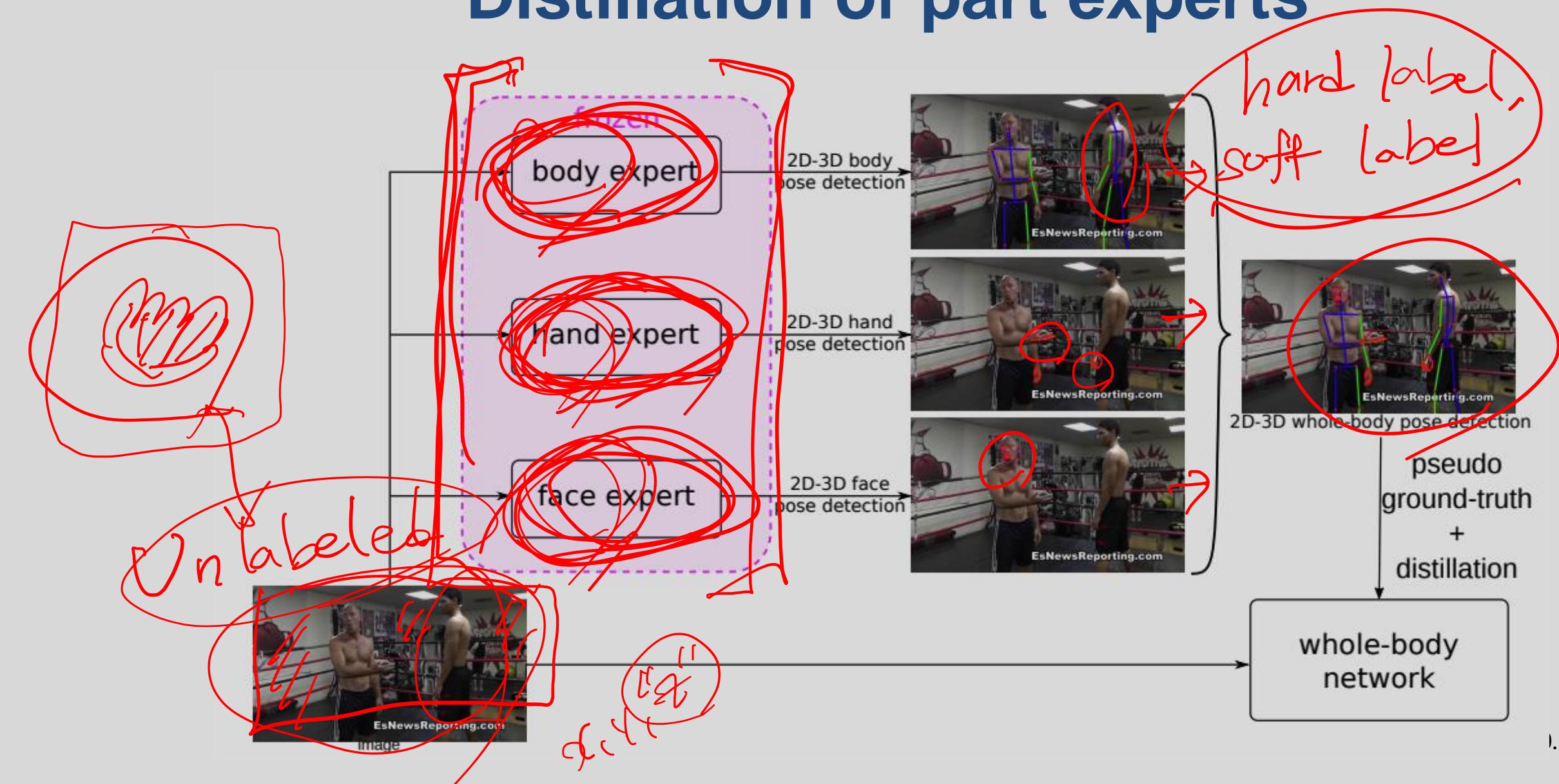
Learning Efficient Object Detection Models with Knowledge Distillation, NIPS'17.

# Distillation of part experts



DOPE: Distillation Of Part Experts for whole-body 3D pose estimation in the wild, ECCV'20.

# Distillation of part experts

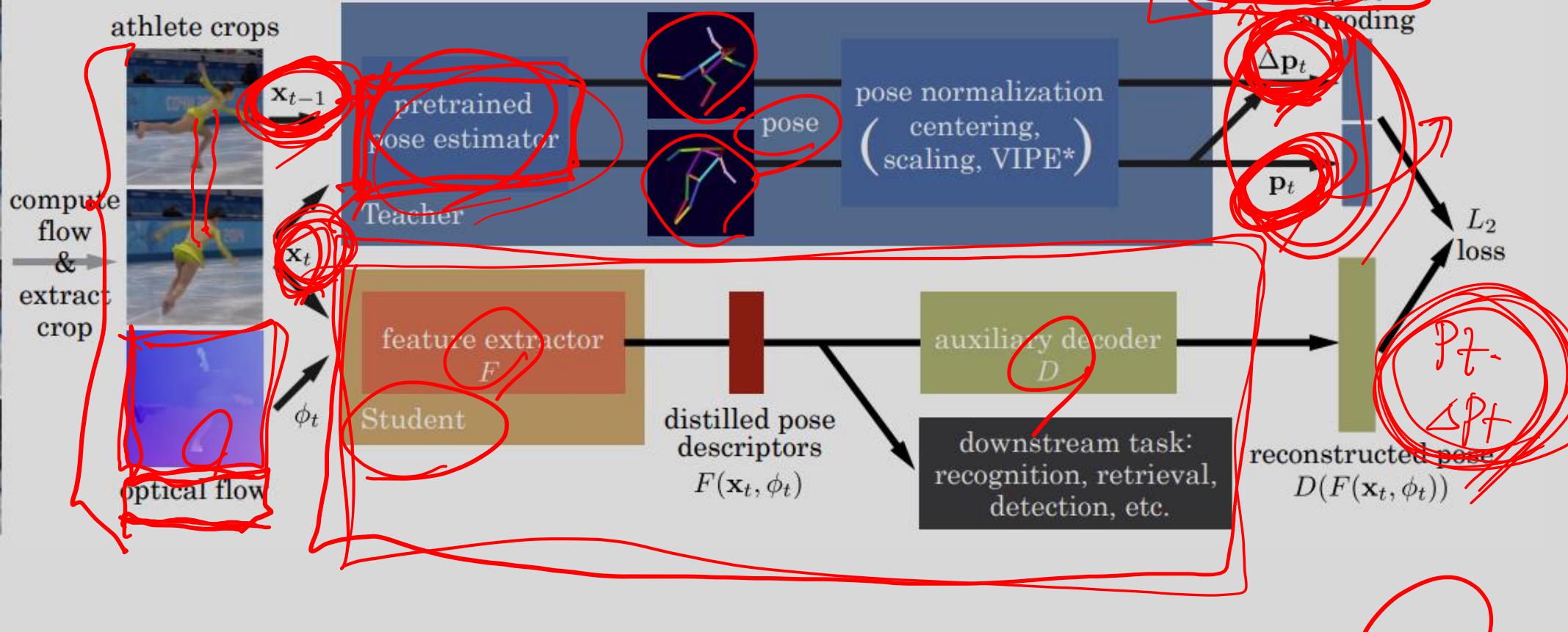


# Distillation of part experts



# Video pose distillation

each frame

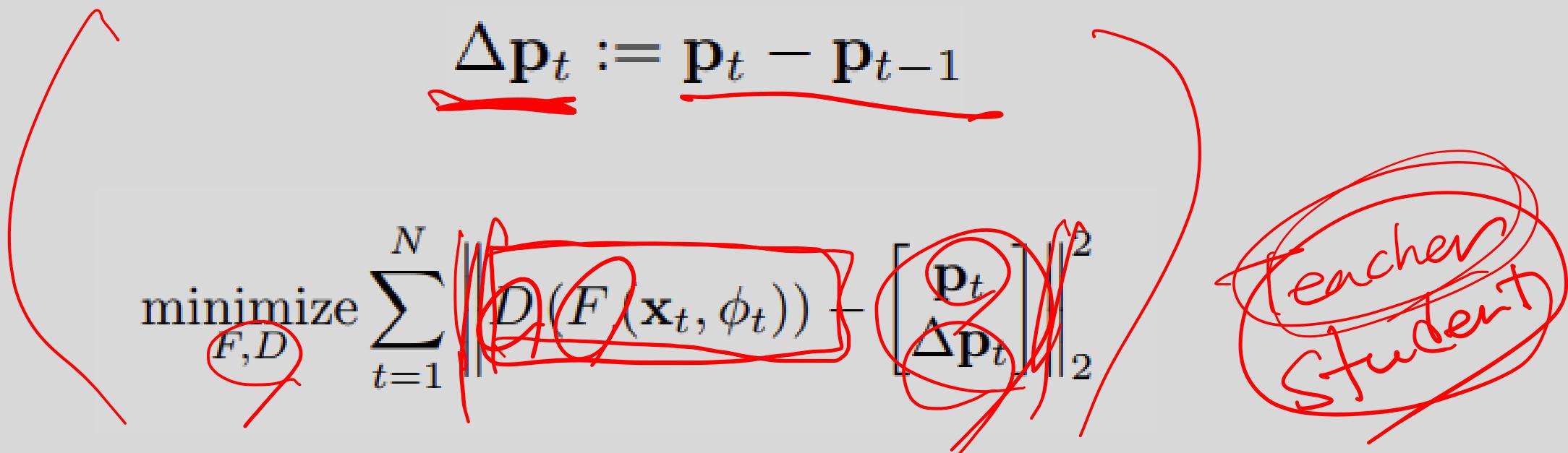


Video Pose Distillation for Few-Shot, Fine-Grained Sports Action Recognition, ICCV'21.

# Video pose distillation

$$\Delta \mathbf{p}_t := \mathbf{p}_t - \mathbf{p}_{t-1}$$
$$\min_{F, D} \sum_{t=1}^N \|D(F(\mathbf{x}_t, \phi_t)) - \begin{bmatrix} \mathbf{p}_t \\ \Delta \mathbf{p}_t \end{bmatrix}\|_2^2$$

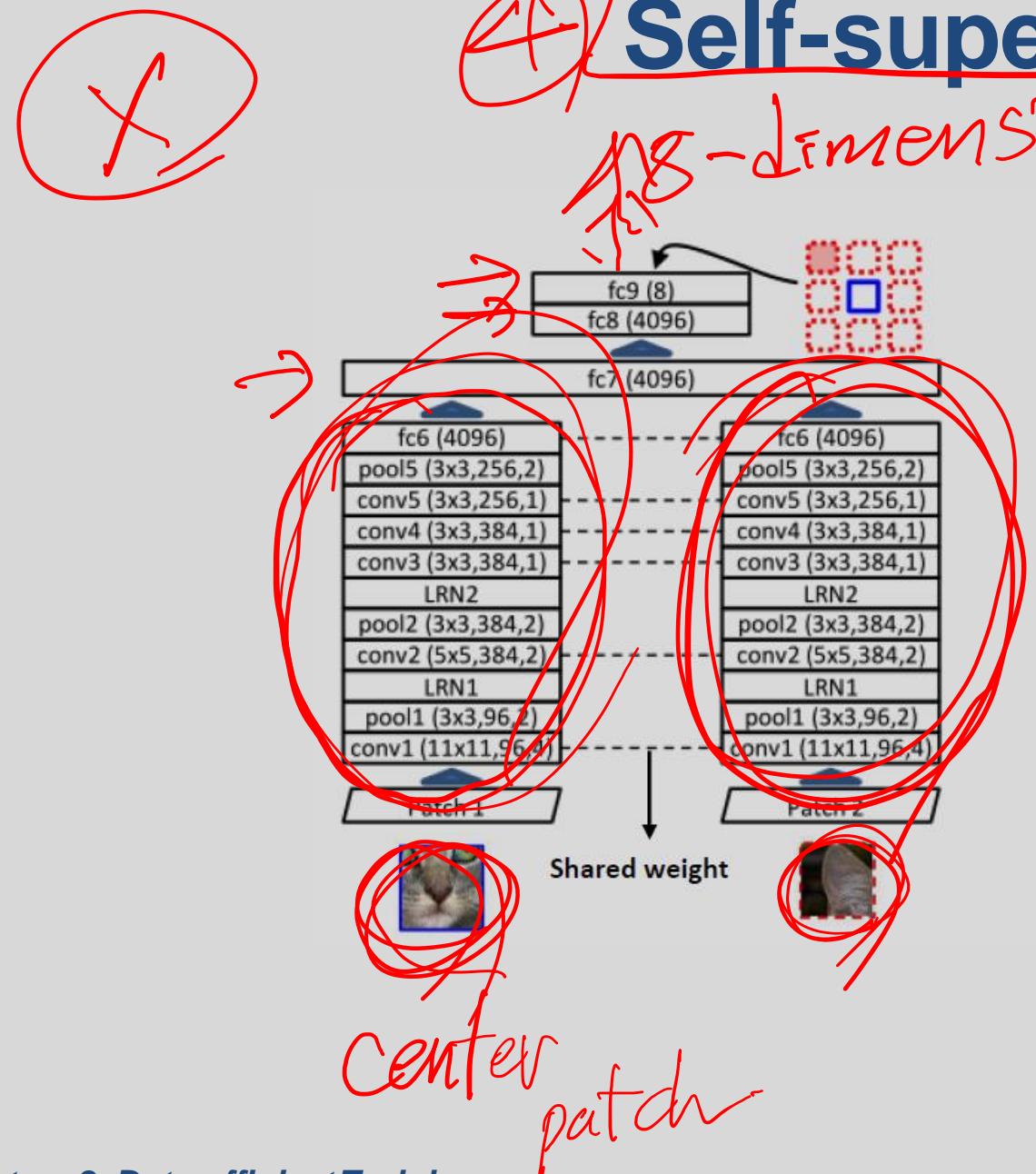
Teacher  
Student



Video Pose Distillation for Few-Shot, Fine-Grained Sports Action Recognition, ICCV'21.

# Self-supervised learning

*18-dimensional vector.*

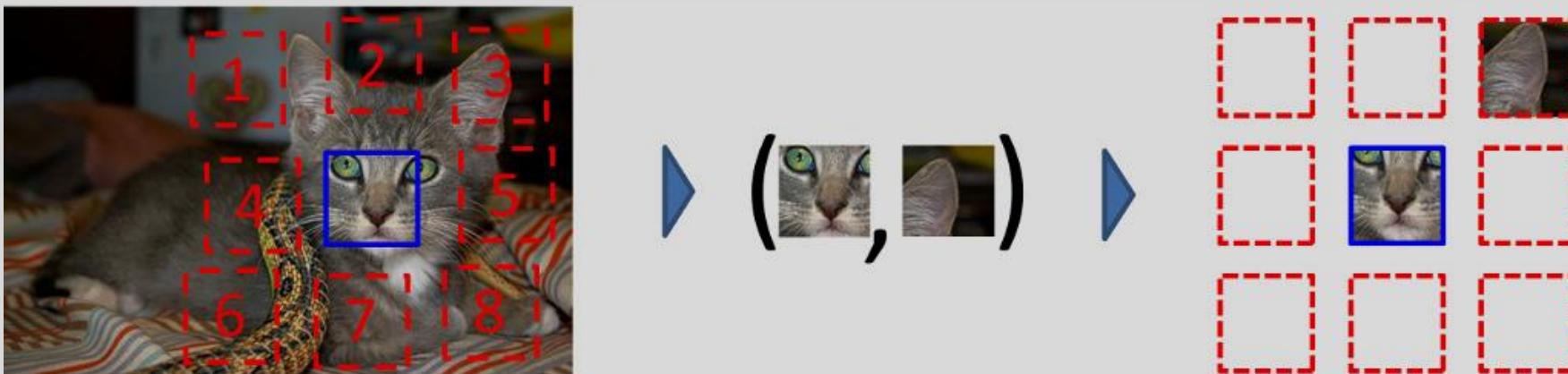


Unsupervised Visual Representation Learning by Context Prediction /ICCV 2015

# Self-supervised learning

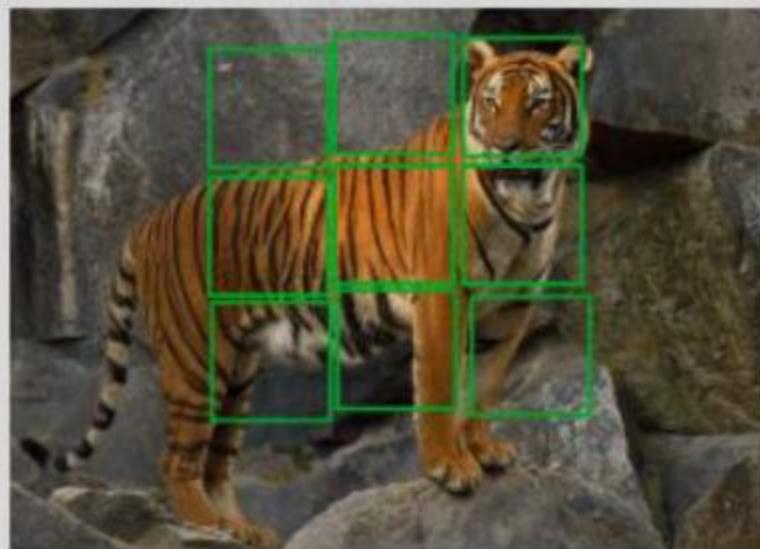
Context Prediction: Predict relative positions of patches

- You have to understand the object to solve this problem!
- Be aware of trivial solution! CNN is especially good at it

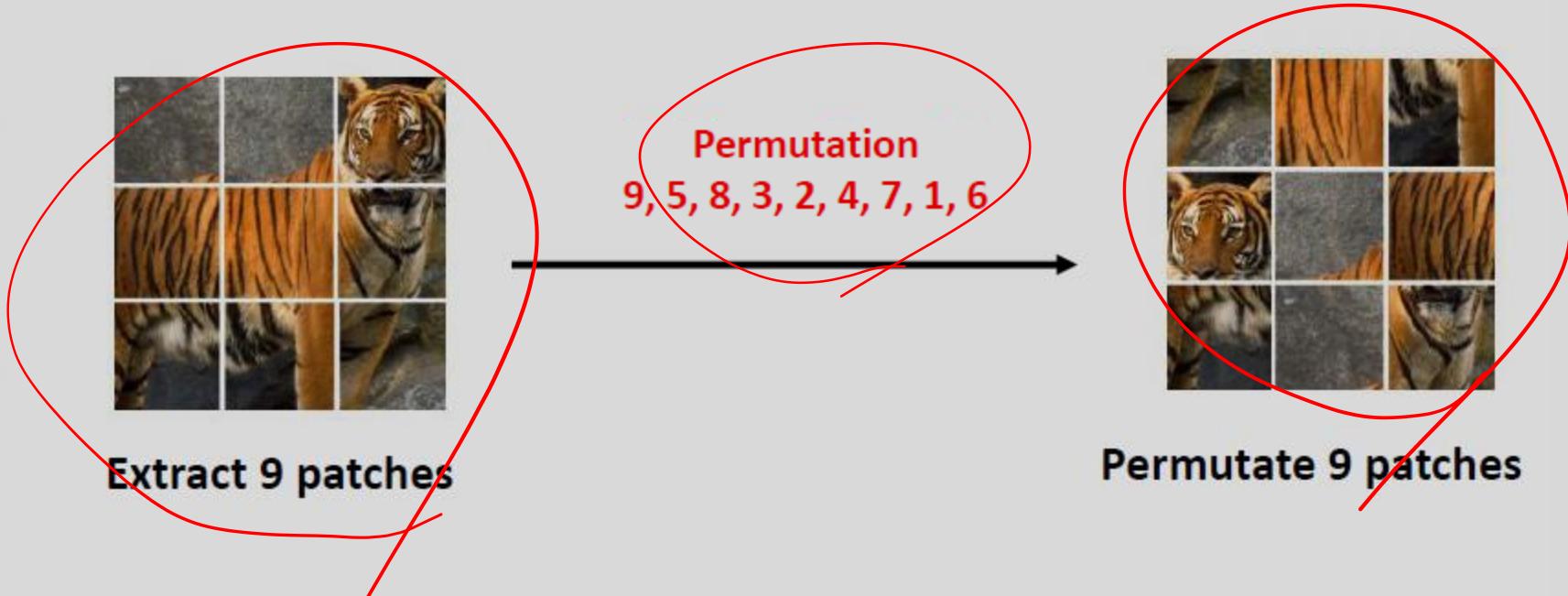


Unsupervised Visual Representation Learning by Context Prediction. ICCV 2015

# Self-supervised learning



Sample image



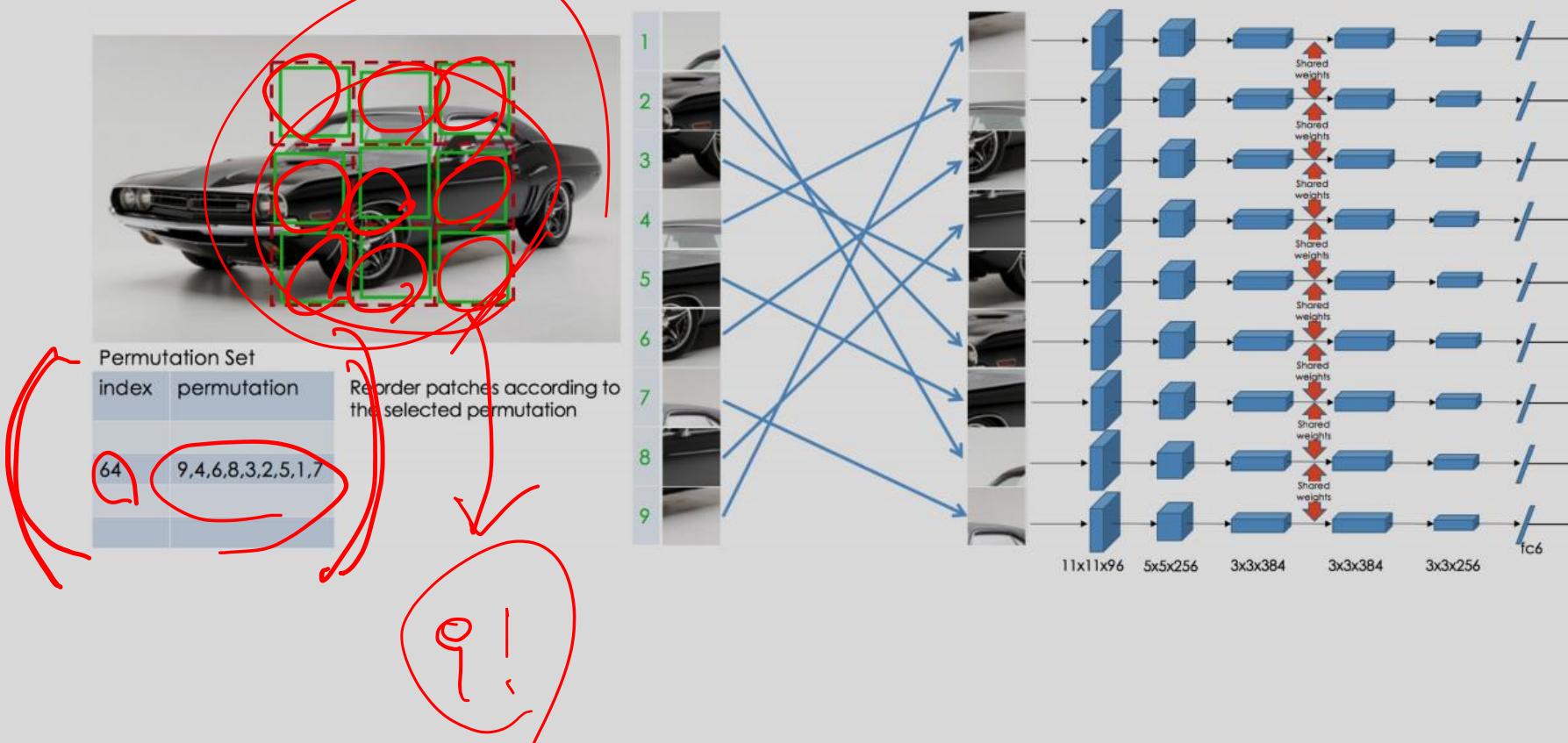
Unsupervised learning of visual representations by solving jigsaw puzzles. In ECCV2016.



# Self-supervised learning

## Solving the Jigsaw

- Use stronger supervision, solve the real jigsaw problem
- Harder problem, better ability for networks

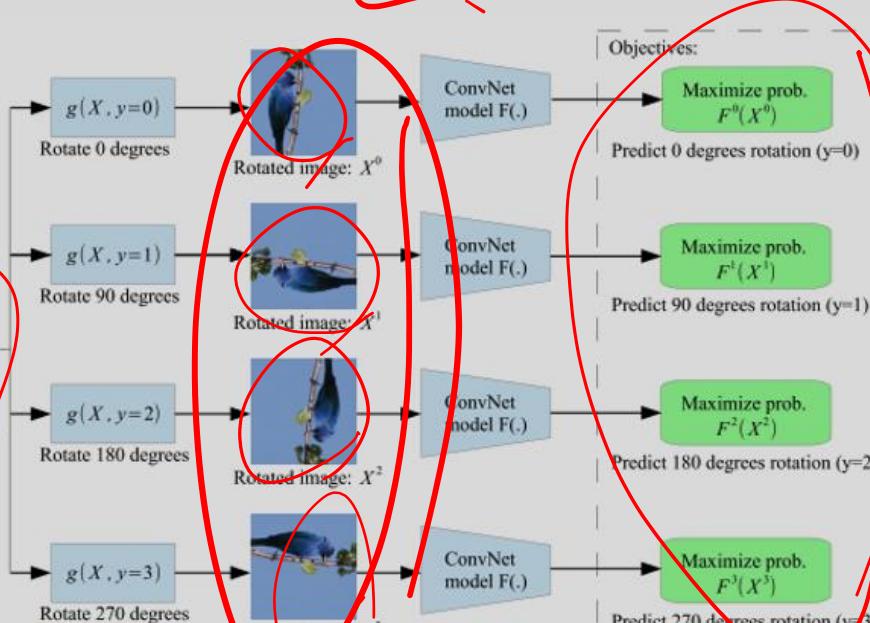


# Self-supervised learning

## Predicting the rotations

- Predict the 4 types of rotation angles.

$0^\circ, 90^\circ, 180^\circ, 270^\circ$



Method	Conv1	Conv2	Conv3	Conv4	Conv5
ImageNet labels	19.3	36.3	44.2	48.3	50.5
Random	11.6	17.1	16.9	16.3	14.1
Random rescaled [Krähenbühl et al., 2015]	17.5	23.0	24.5	23.2	20.6
Context [Doersch et al., 2015]	16.2	23.3	30.2	31.7	29.6
Context Encoders [Pathak et al., 2016b]	14.1	20.7	21.0	19.8	15.5
Colorization [Zhang et al., 2016a]	12.5	24.5	30.4	31.5	30.3
Jigsaw Puzzles [Noroozi & Favaro, 2016]	18.2	28.8	34.0	33.9	27.1
BIGAN [Donahue et al., 2016]	17.7	24.5	31.0	29.9	28.0
Split-Brain [Zhang et al., 2016b]	17.7	29.3	35.4	35.2	32.8
Counting [Noroozi et al., 2017]	18.0	30.6	34.3	32.5	25.7
(Ours) RotNet	18.8	31.7	38.7	38.2	36.5

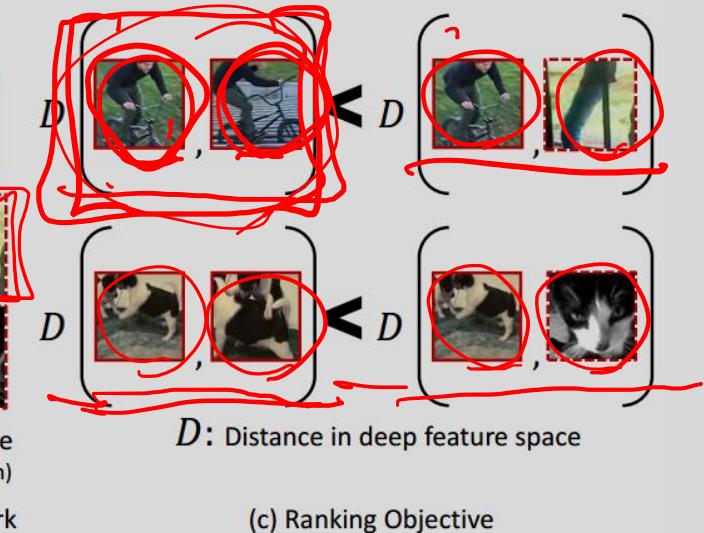
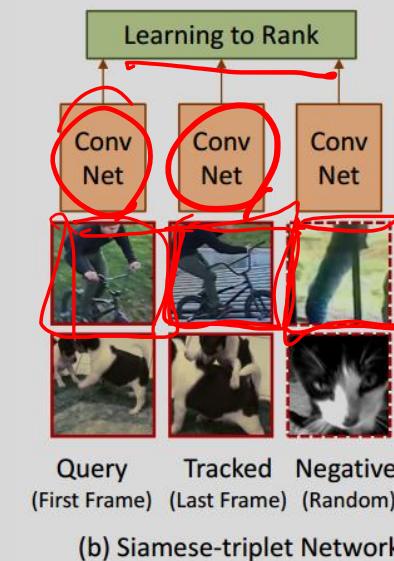
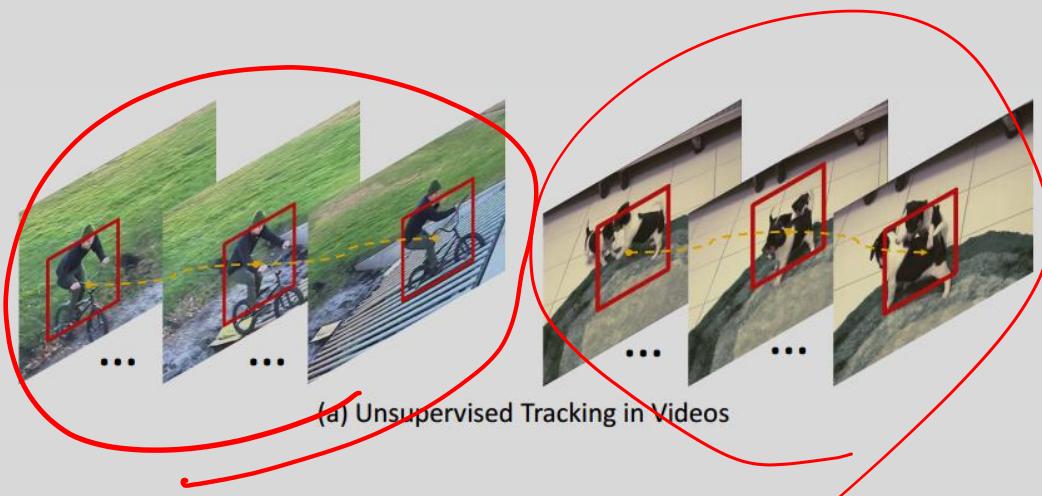
ImageNet classification top-1 accuracy

Unsupervised representation learning by predicting image rotations. In *ICLR 2018*.



# Self-supervision for video

Find corresponding pairs using visual tracking



Wang, X., & Gupta, A. (2015). Unsupervised learning of visual representations using videos. In *ICCV2015*

# Self-supervision for video

Is the temporal order of a video correct?

- Encode the cause and effect of action

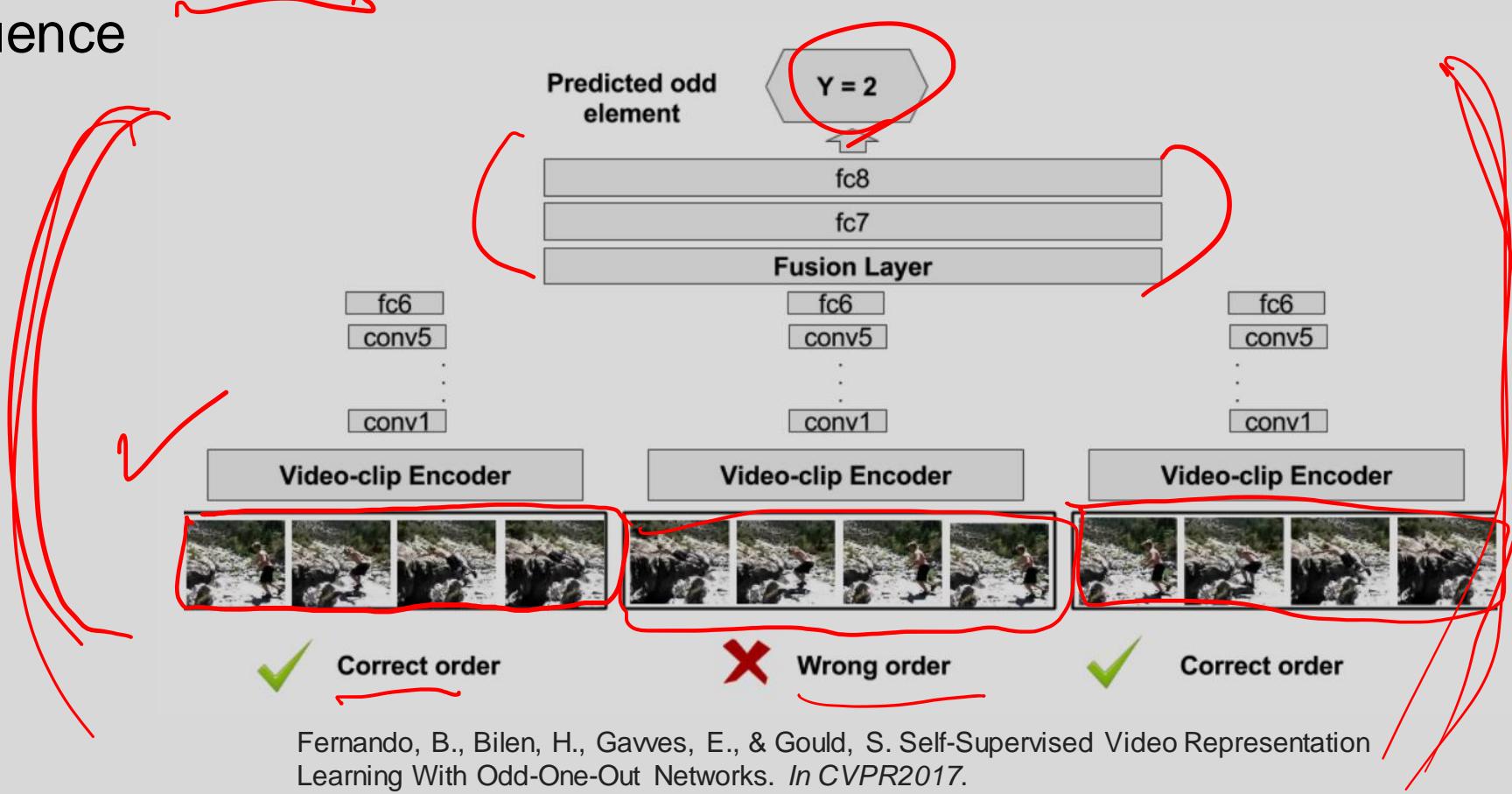


Misra, I., Zitnick, C. L., & Hebert, M. Shuffle and learn: unsupervised learning using temporal order verification. In *ECCV2016*.

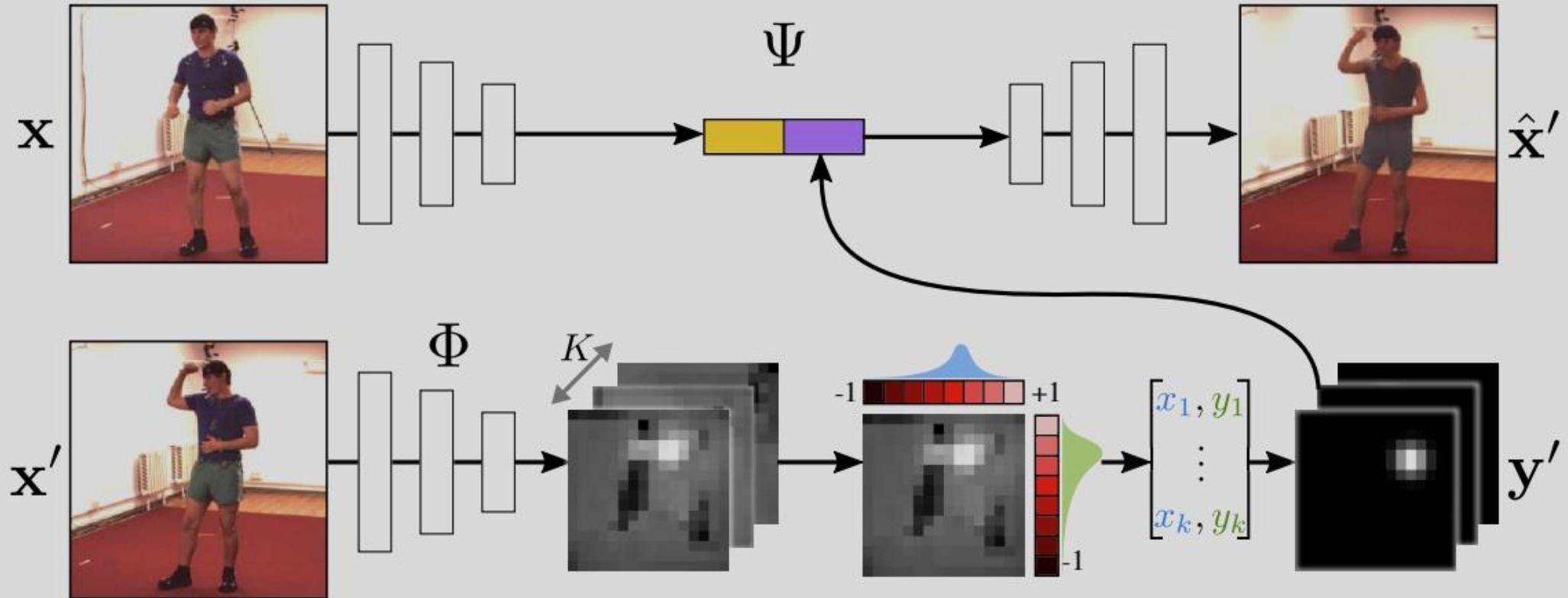
# Self-supervision for video

Is the temporal order of a video correct?

- Find the odd sequence

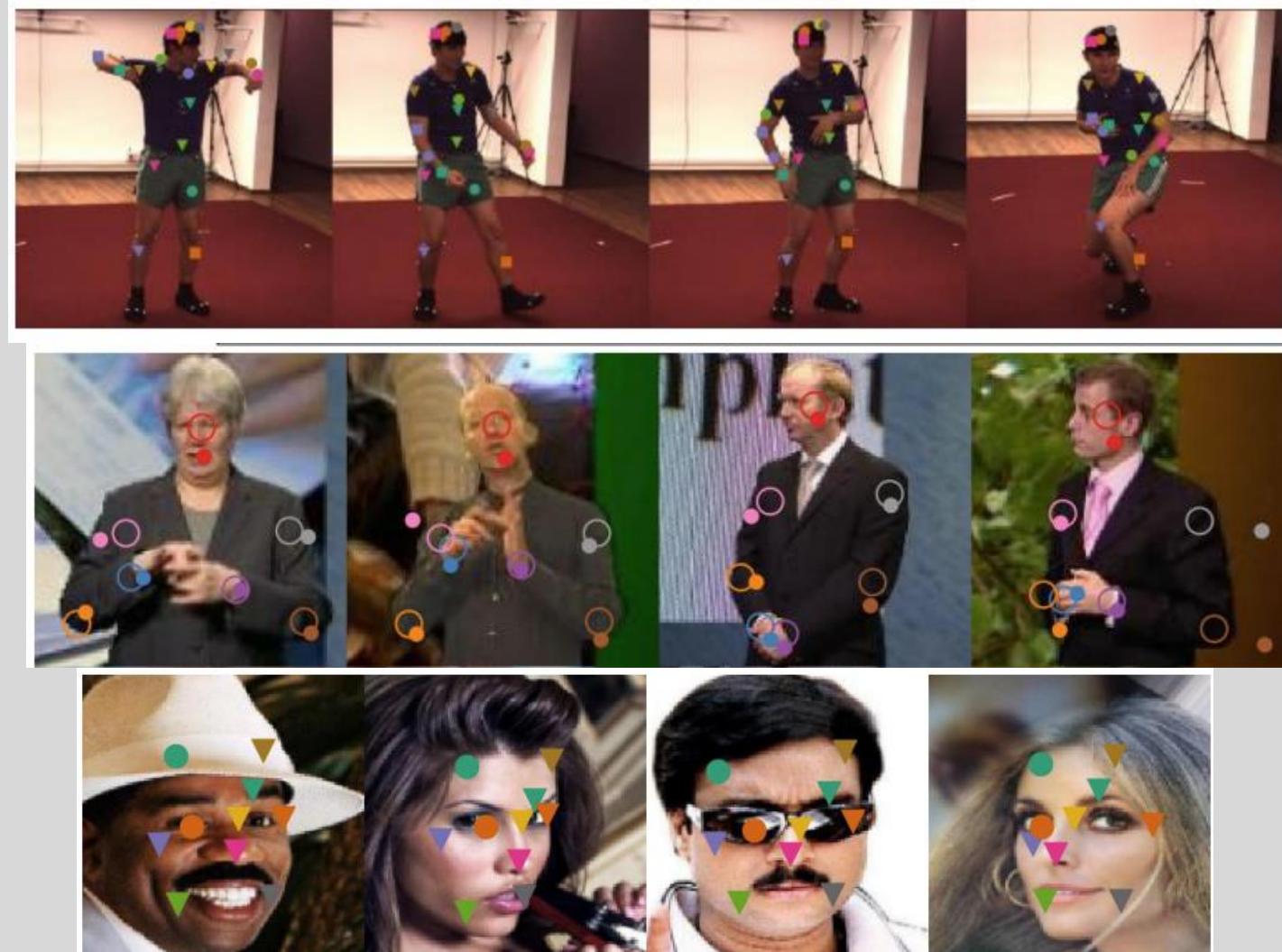


# Self-supervision for pose



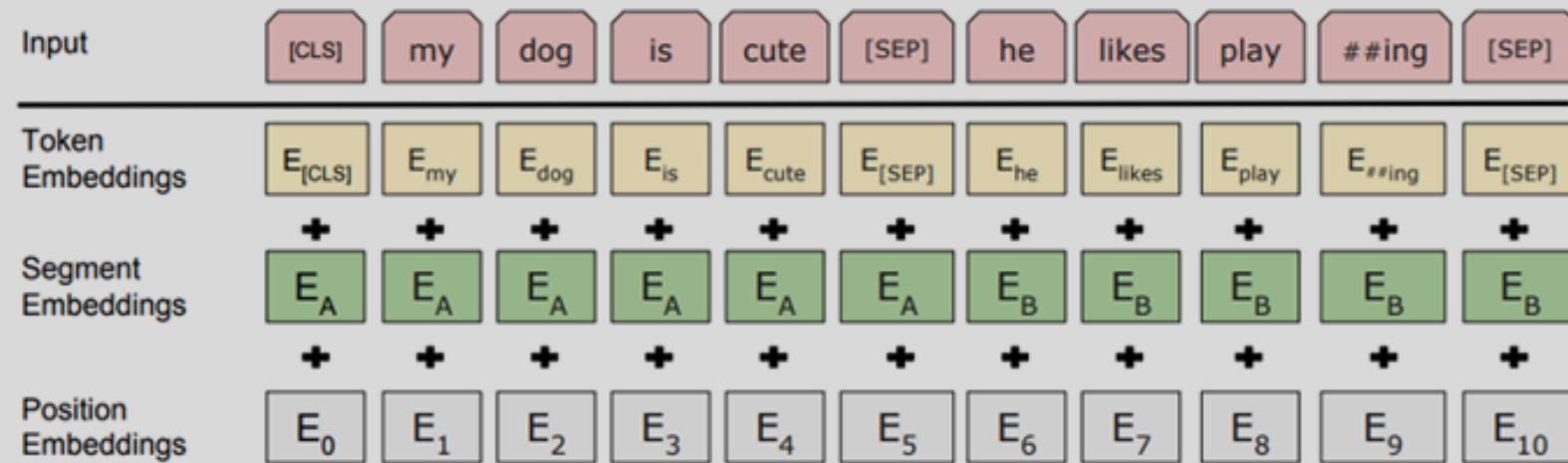
Unsupervised Learning of Object Landmarks through Conditional Image Generation, NeurIPS'18.

# Self-supervision for pose

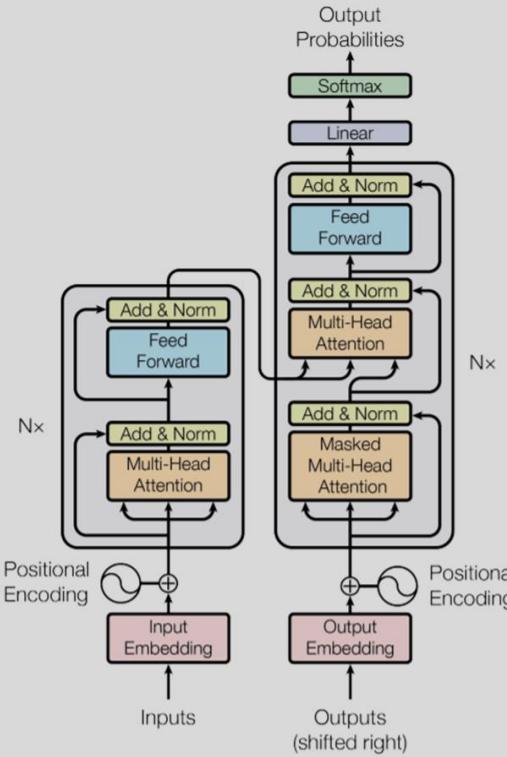


# Self-supervised learning

- Self-supervision: Learning without tagged data.
- The method could be applied to any inputs.
  - Speech, image, video, text and etc.



# Self-supervised learning



**Input** = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

**Label** = IsNext

**Input** = [CLS] the man [MASK] to the store [SEP]

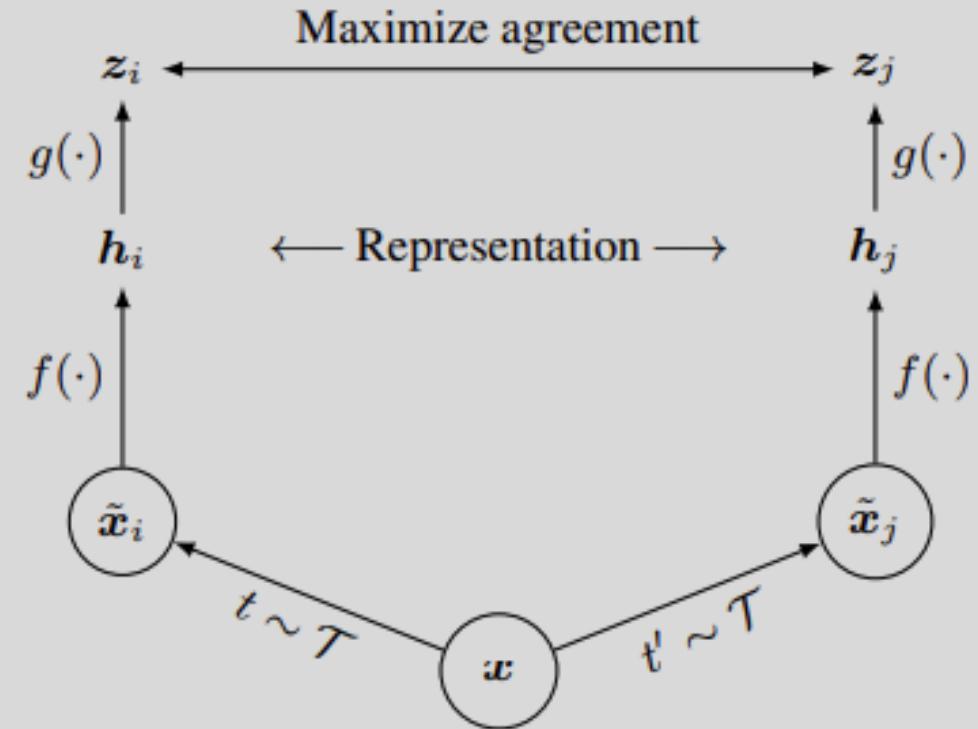
penguin [MASK] are flight ##less birds [SEP]

**Label** = NotNext

Transformer architecture is trained by 1) Masked language model, 2) Next sentence prediction

# Contrastive learning

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)}$$



# Contrastive learning



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)

(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$ 

(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

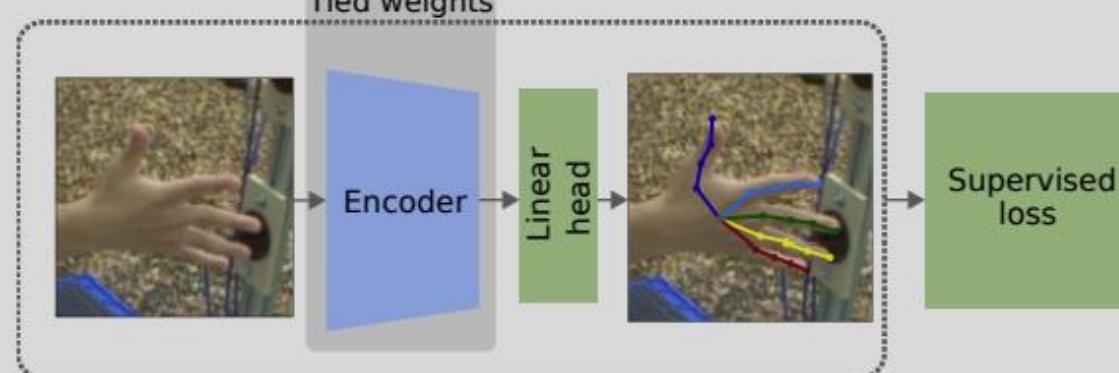
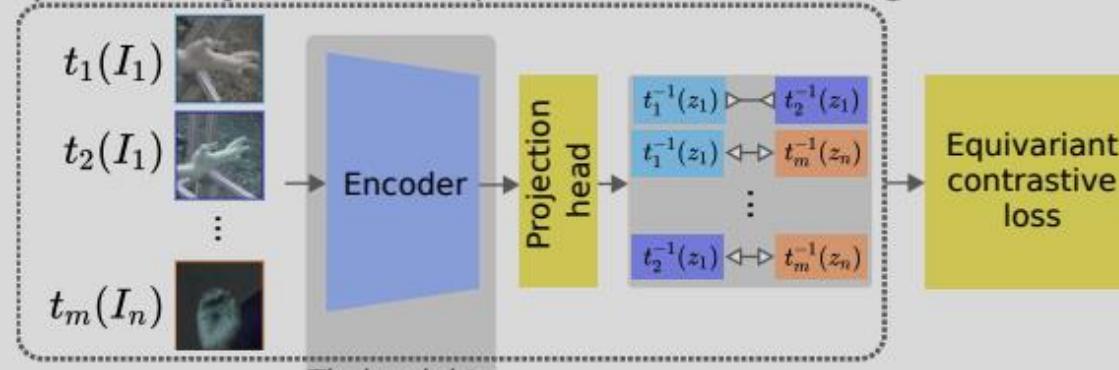
# Contrastive learning

	Food	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
SimCLR (ours)	<b>76.9</b>	<b>95.3</b>	80.2	48.4	<b>65.9</b>	60.0	61.2	<b>84.2</b>	<b>78.9</b>	89.2	<b>93.9</b>	<b>95.0</b>
Supervised	75.2	<b>95.7</b>	<b>81.2</b>	<b>56.4</b>	64.9	<b>68.8</b>	<b>63.8</b>	83.8	<b>78.7</b>	<b>92.3</b>	<b>94.1</b>	94.2
<i>Fine-tuned:</i>												
SimCLR (ours)	<b>89.4</b>	<b>98.6</b>	<b>89.0</b>	<b>78.2</b>	<b>68.1</b>	<b>92.1</b>	<b>87.0</b>	<b>86.6</b>	<b>77.8</b>	92.1	<b>94.1</b>	97.6
Supervised	88.7	98.3	<b>88.7</b>	<b>77.8</b>	67.0	91.4	<b>88.0</b>	86.5	<b>78.8</b>	<b>93.2</b>	<b>94.2</b>	<b>98.0</b>
Random init	88.3	96.0	81.9	<b>77.0</b>	53.7	91.3	84.8	69.4	64.1	82.7	72.5	92.5

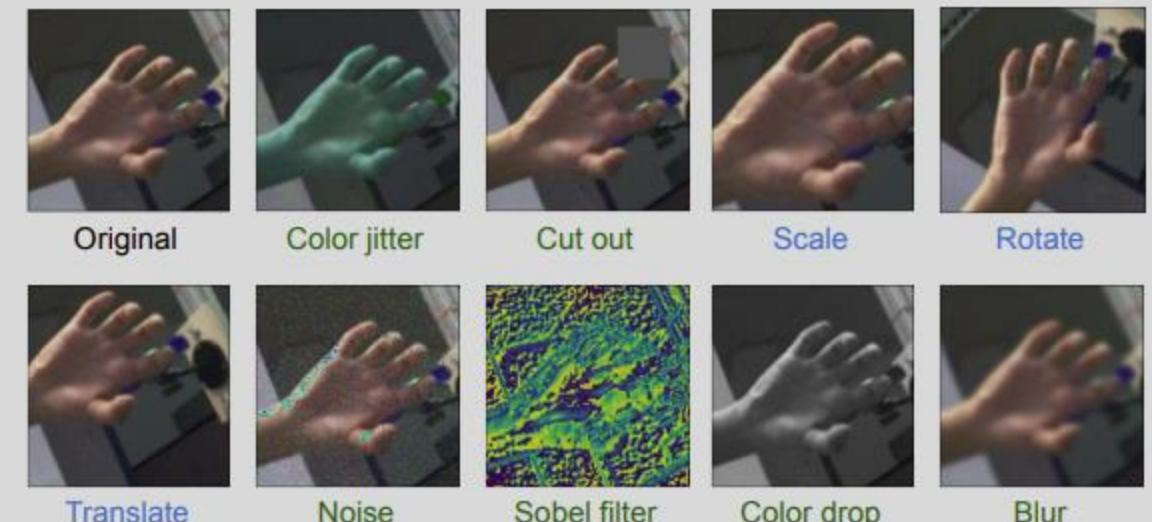
A Simple Framework for Contrastive Learning of Visual Representations, ICLR'20

# Contrastive learning for pose

## I) Self-supervised representation learning

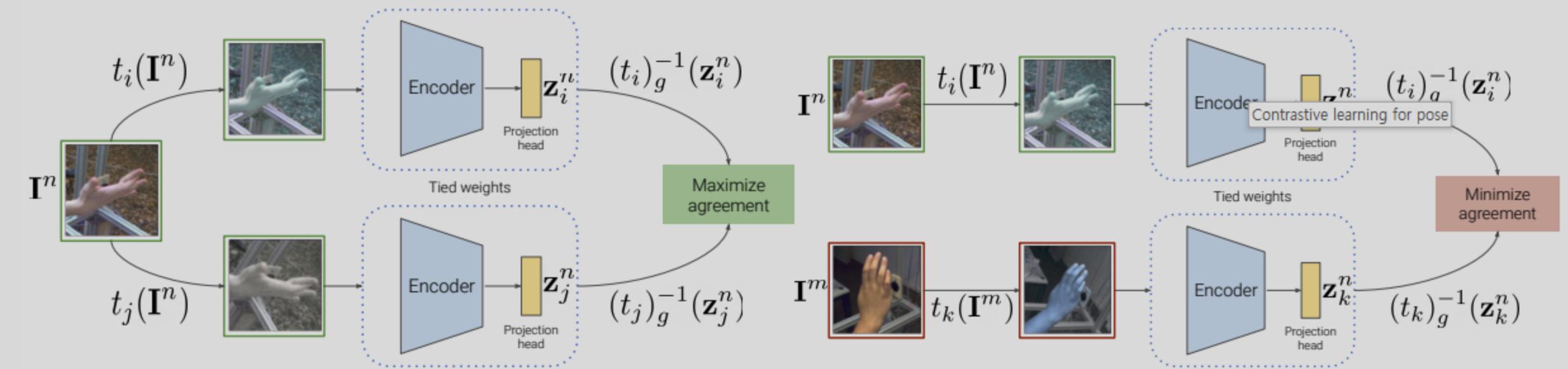


## II) Supervised hand pose estimation



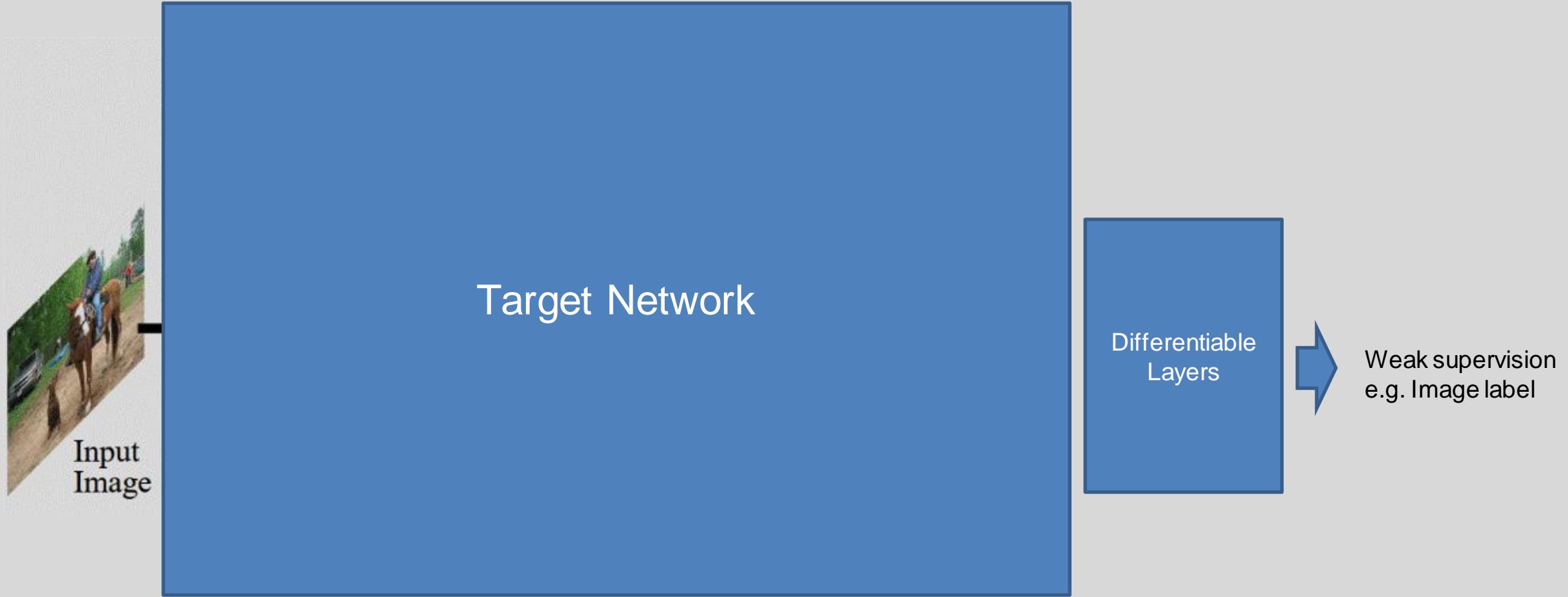
Self-Supervised 3D Hand Pose Estimation from monocular RGB via Contrastive Learning, ICCV'21

# Contrastive learning for pose

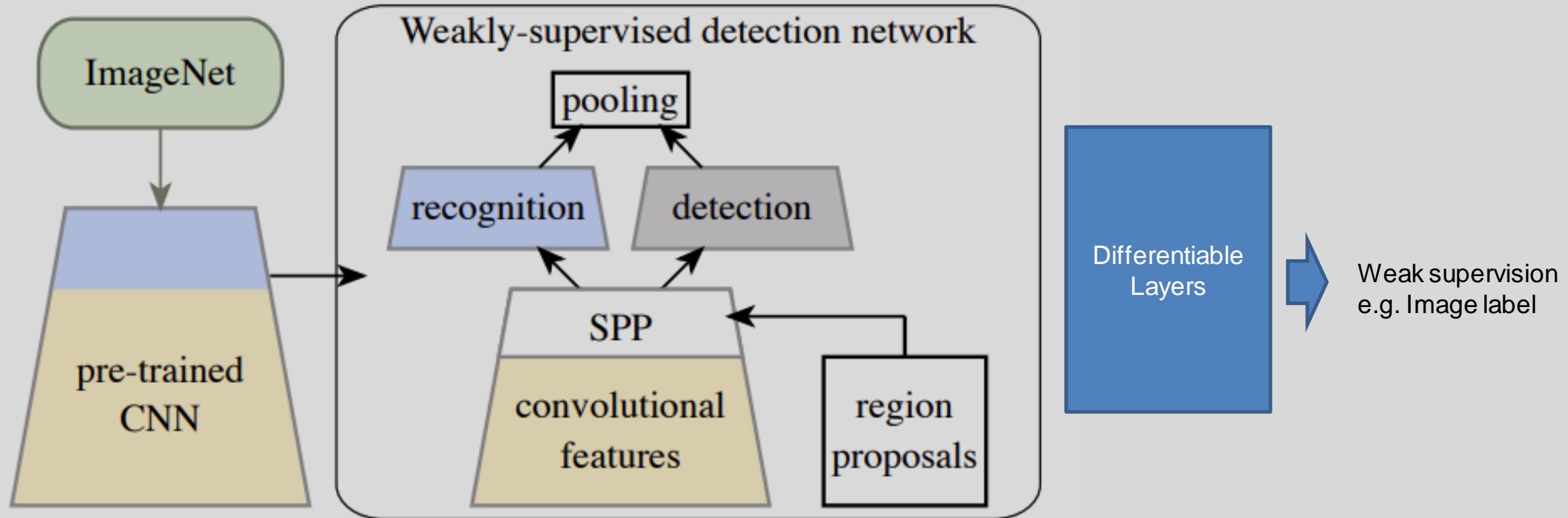


The agreement between projections from the same input image is maximized (left) and agreements amongst projections from different input images are minimized (right)

# Weakly-supervised learning



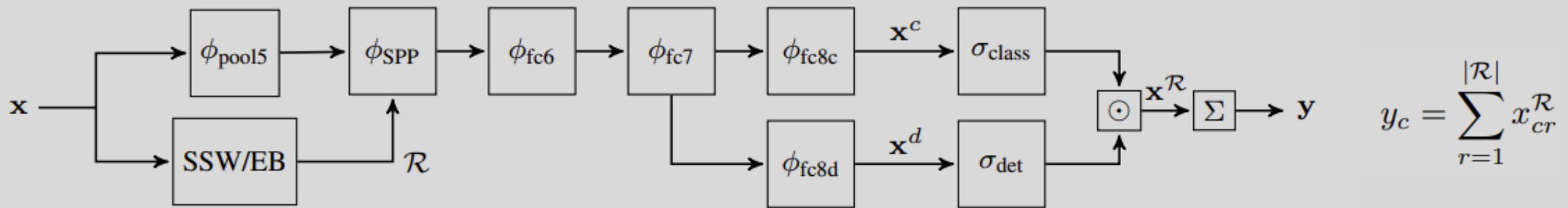
# Weakly-supervised object detection



Weakly Supervised Deep Detection Networks, CVPR'16

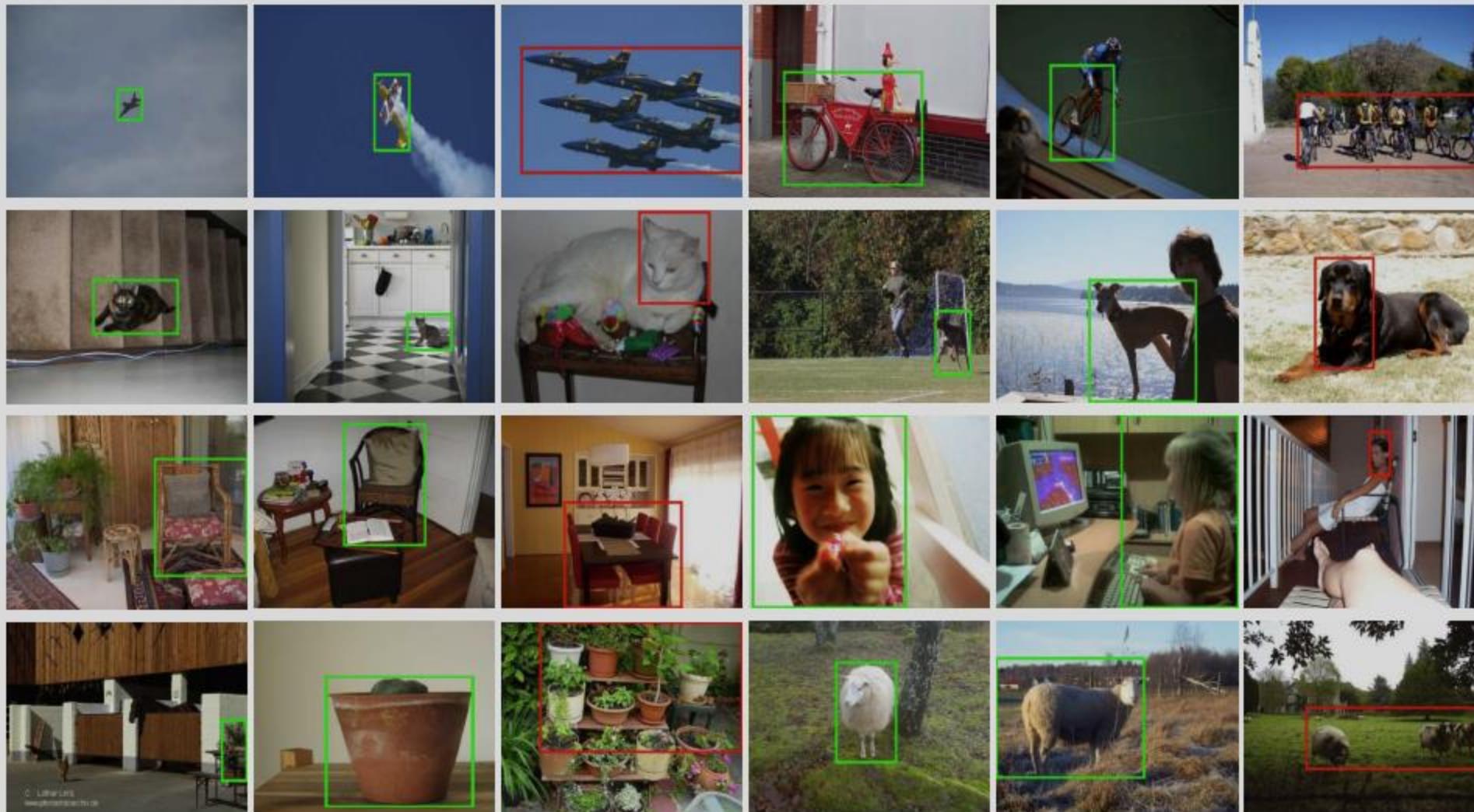
# Weakly-supervised object detection

$$[\sigma_{\text{class}}(\mathbf{x}^c)]_{ij} = \frac{e^{x_{ij}^c}}{\sum_{k=1}^C e^{x_{kj}^c}}$$

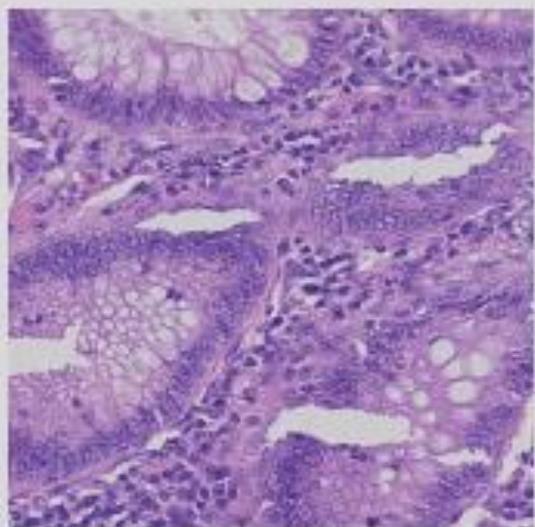


$$[\sigma_{\text{det}}(\mathbf{x}^d)]_{ij} = \frac{e^{x_{ij}^d}}{\sum_{k=1}^{|\mathcal{R}|} e^{x_{ik}^d}}$$

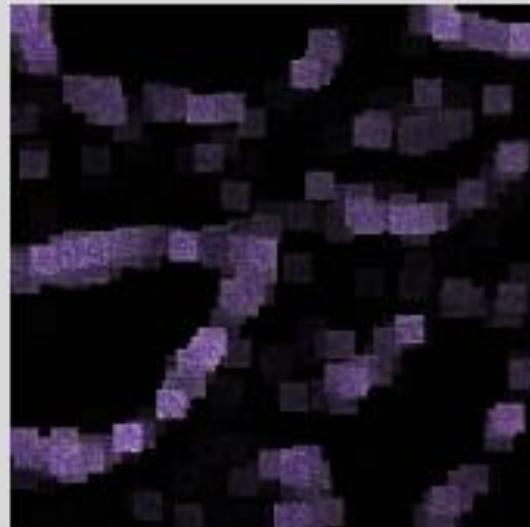
# Weakly-supervised object detection



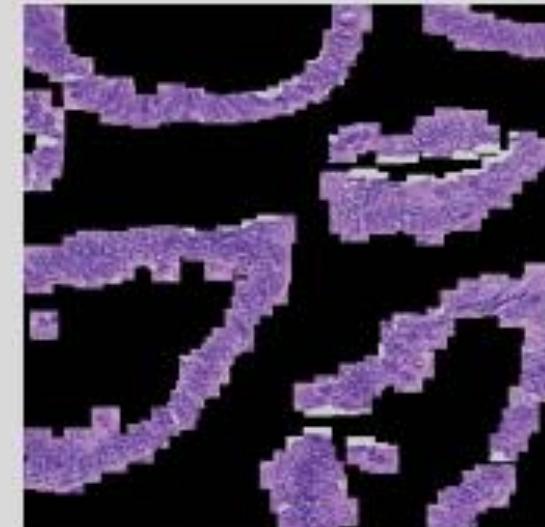
# Weakly-supervised segmentation



Original image



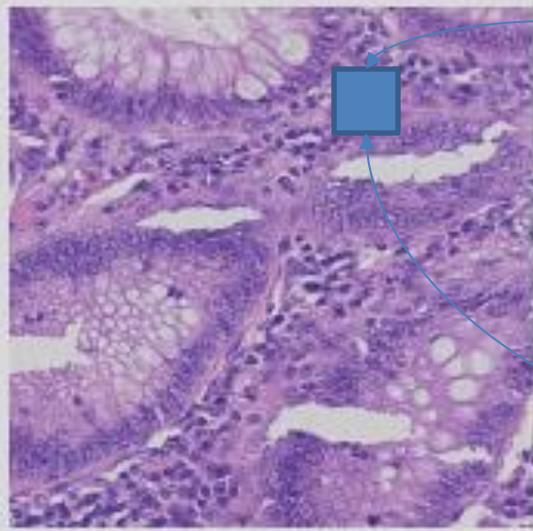
Predicted patch weights



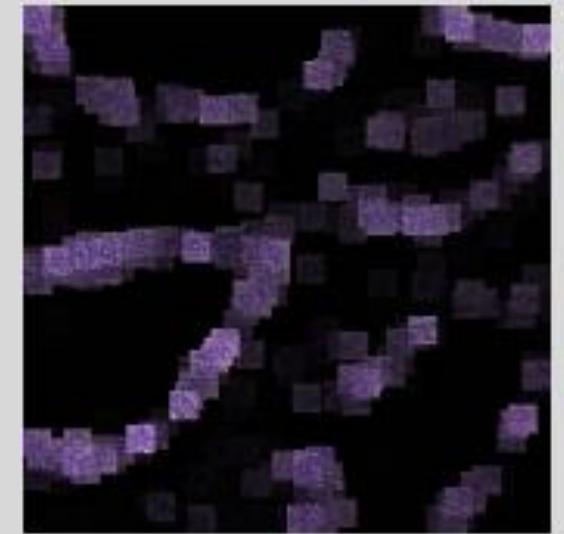
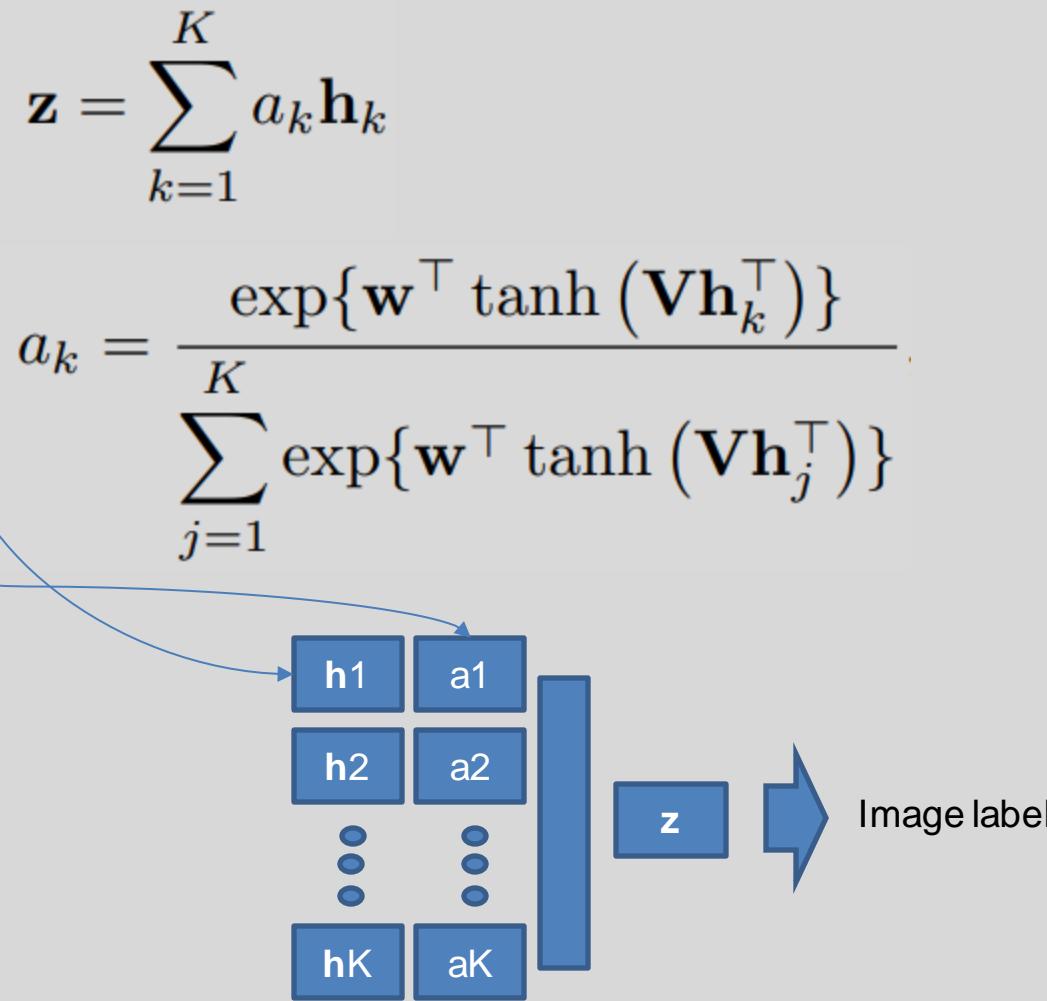
Ground-truth patches

Attention-based Deep Multiple Instance Learning, ICML'18

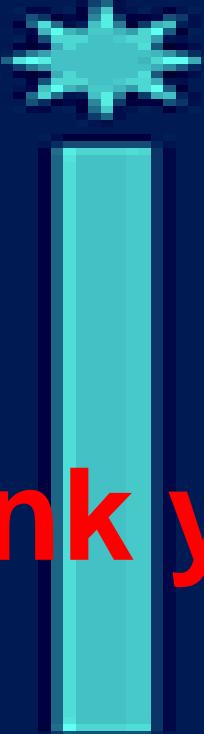
# Weakly-supervised segmentation



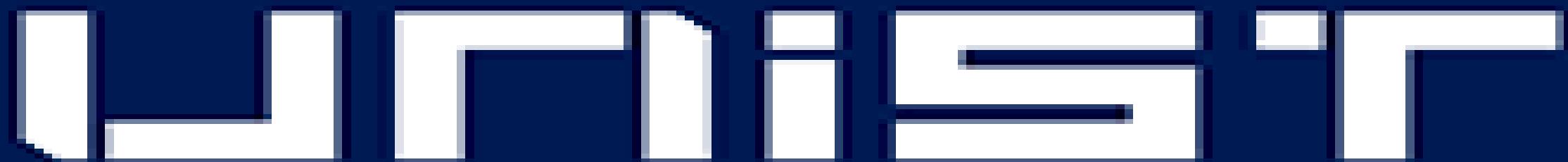
Original image



Predicted patch weights



Thank you!



ULSAN NATIONAL INSTITUTE OF  
SCIENCE AND TECHNOLOGY

2 0 0 7