



H-Unique

Lend a Hand To Fight Child Abuse



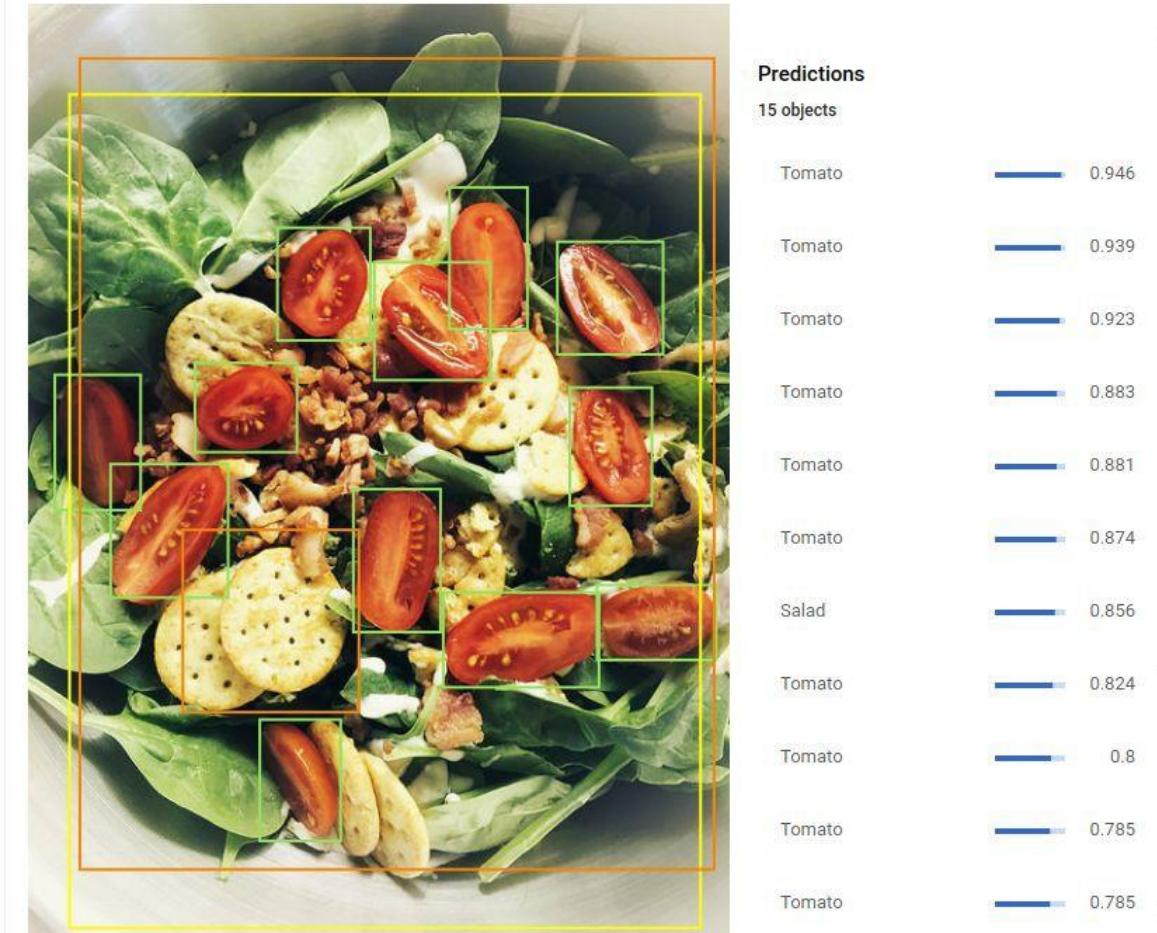
Lancaster
University



Weakly Supervised Learning in Computer Vision

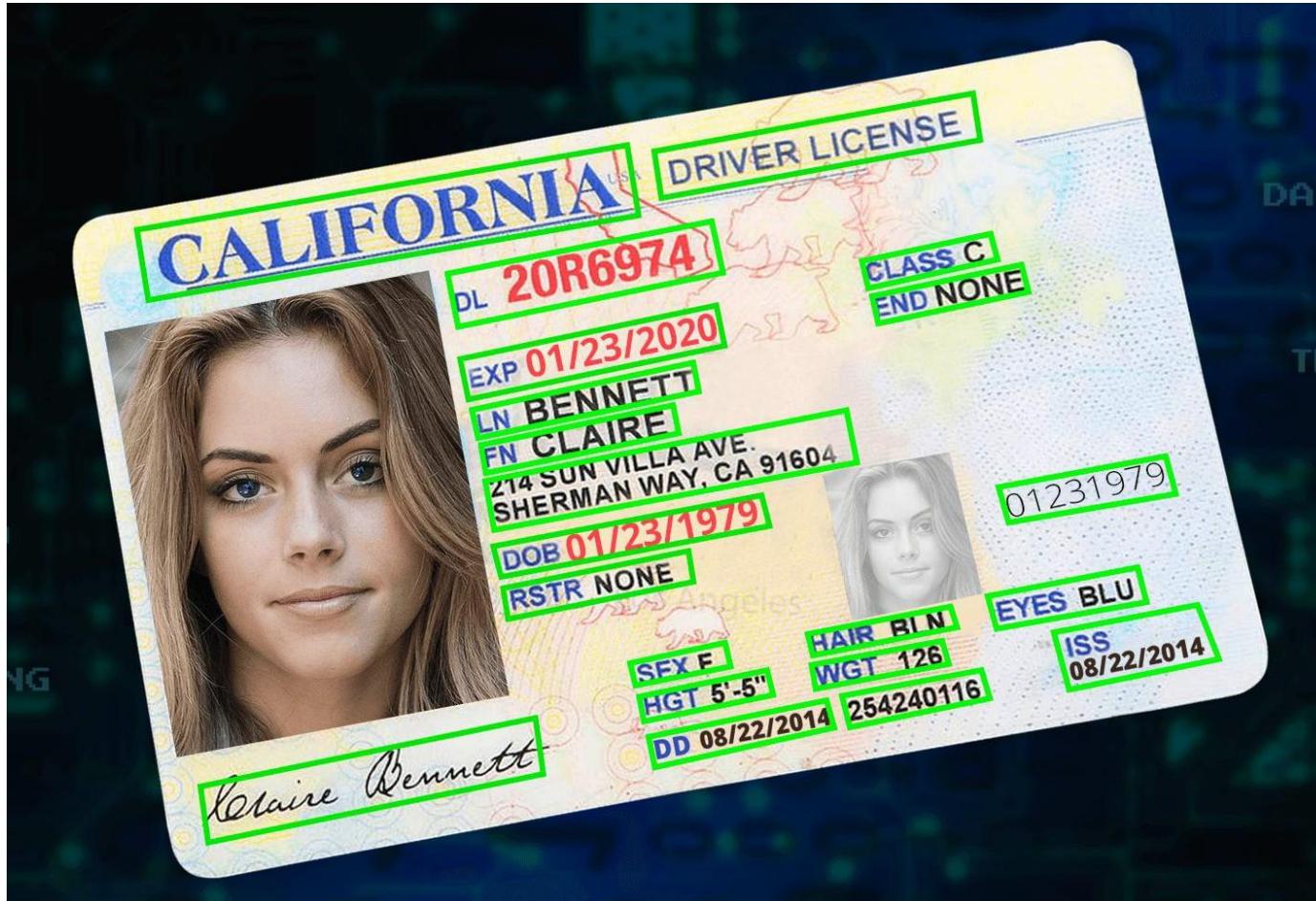
Speaker: Xinyu Yang
(Senior Research Associate)

Computer vision is working in real life!



<https://cloud.google.com/vision/automl/object-detection/docs>

Computer vision is working in real life!



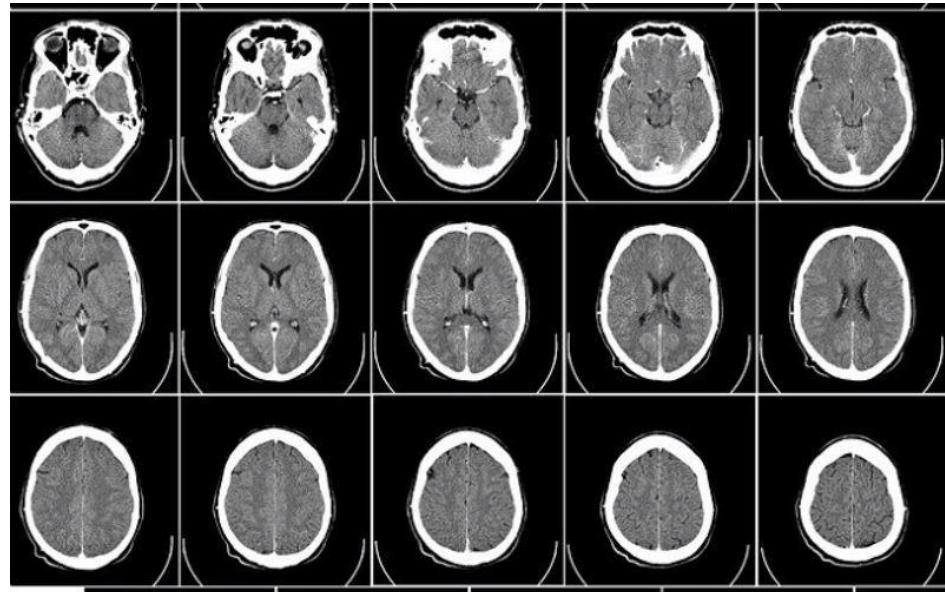
<https://mobidev.biz/blog/ocr-machine-learning-implementation>

Computer vision is working in real life!



<https://www.biometricupdate.com/202001/airport-biometrics-predictions-deployments-upgrades-and-plans-for-future-services>

Computer vision is working in real life!



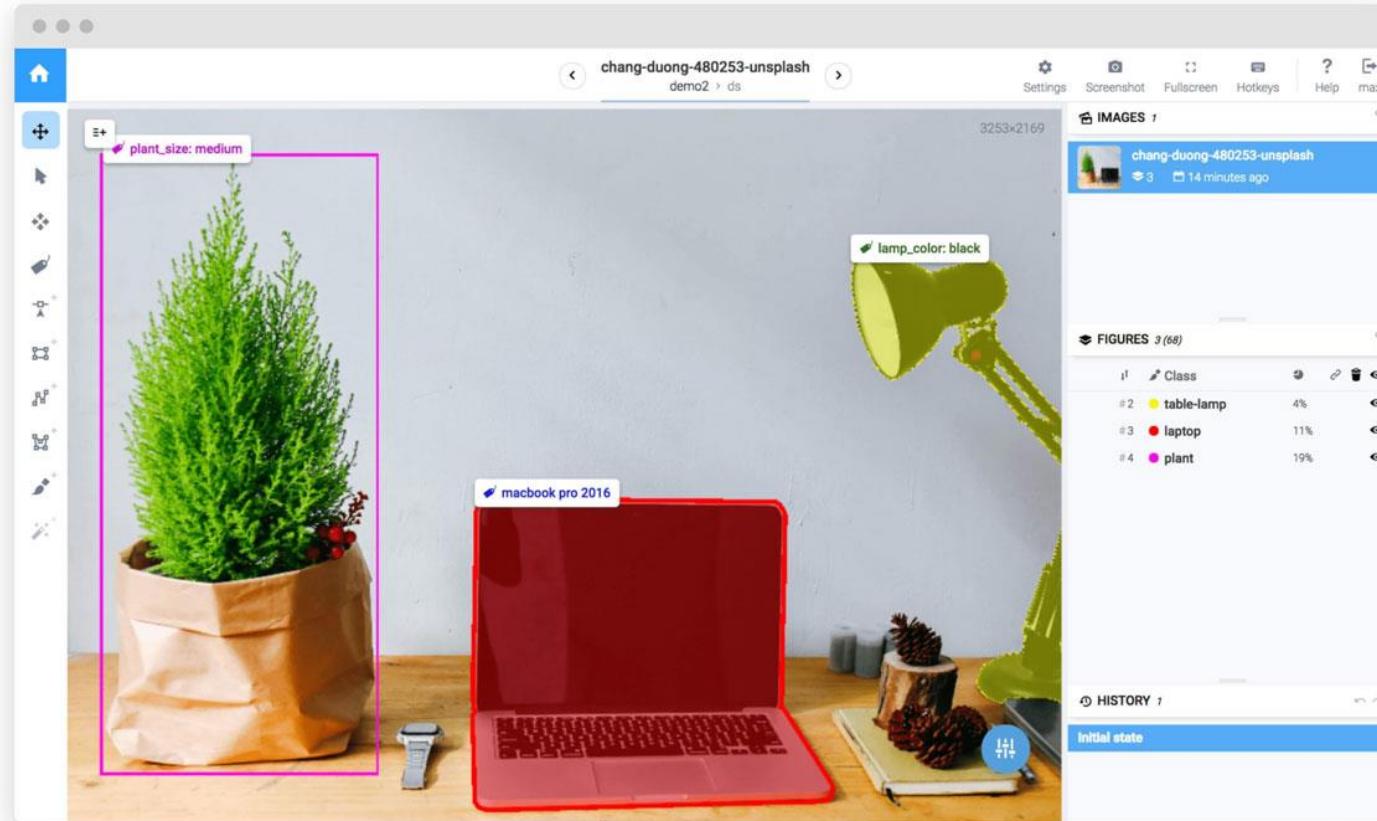
Medical AI



Self-driving cars

Introduction

Behind the success... Huge annotation costs.



ImageNet1K

Multi-label
annotation takes
26 seconds / image
14M images in total!

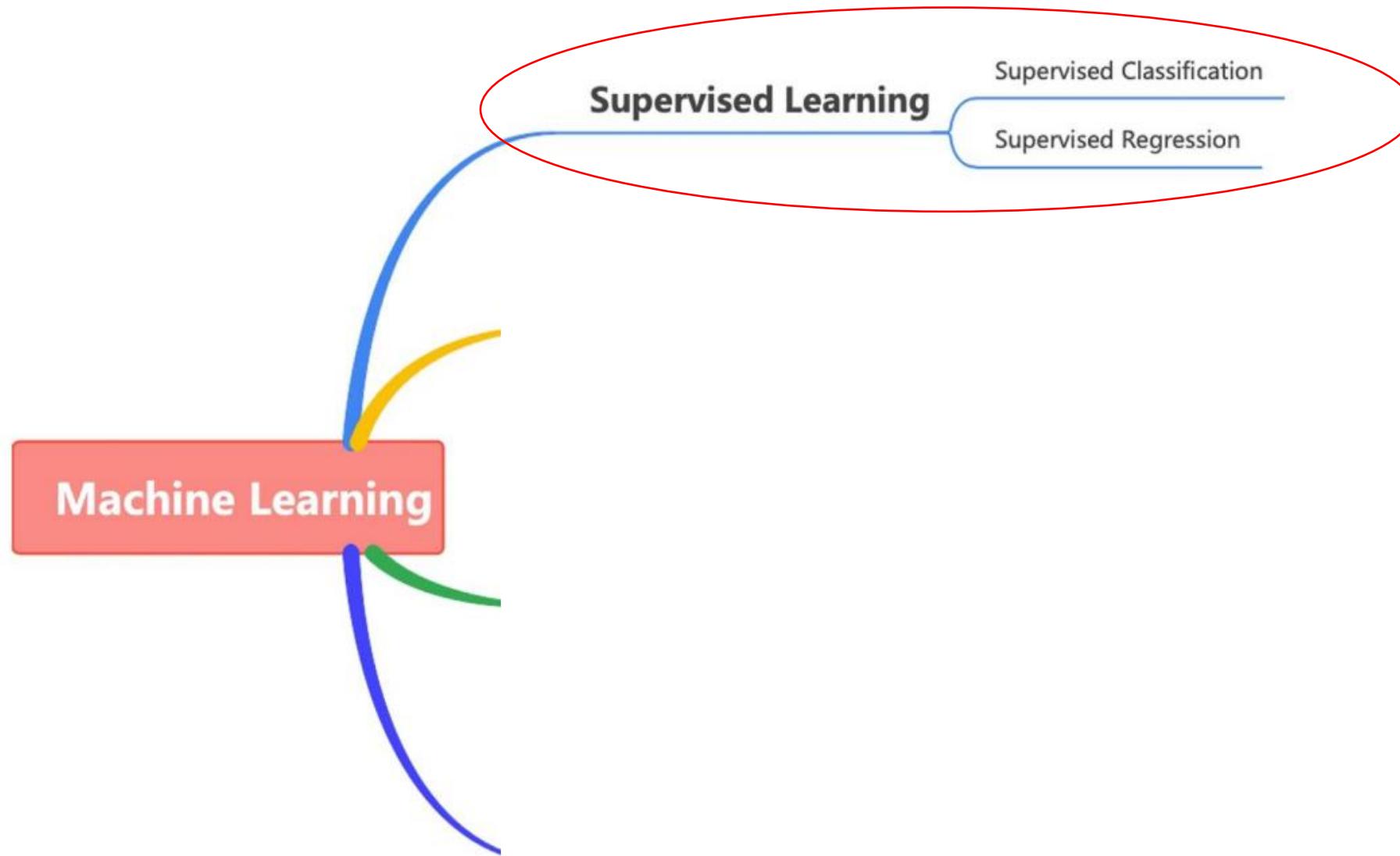
Behind the success… Huge annotation costs.



Cityscapes

1.5 hour / image.

Supervision Types



Supervised Classification

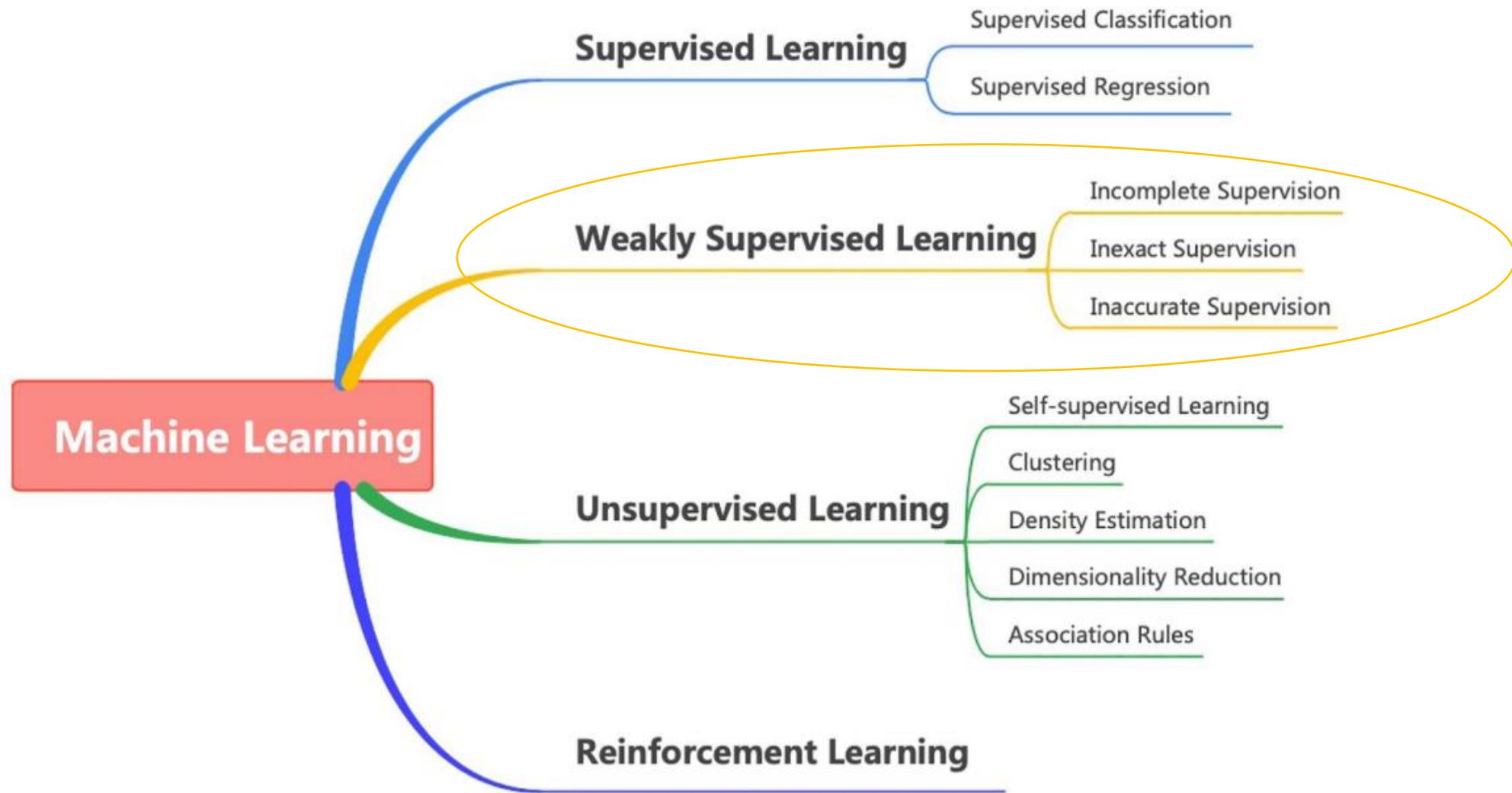


Label: Yes

Is there a cat?

Supervision Types

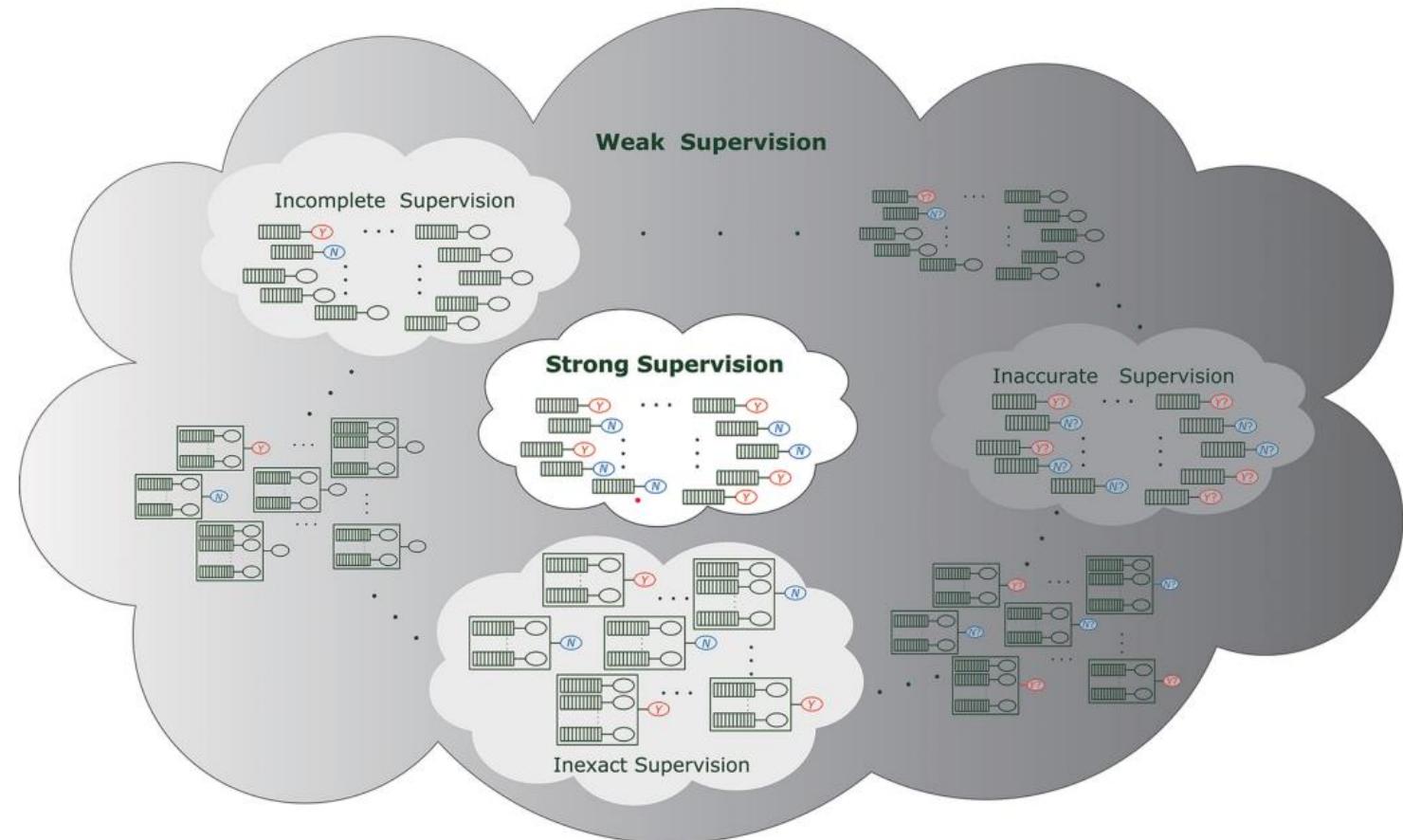
10



Weak Supervision Types

11

- Incomplete supervision:
labeled + unlabeled
- Inaccurate supervision:
noisy labels
- Inexact supervision:
coarse labels

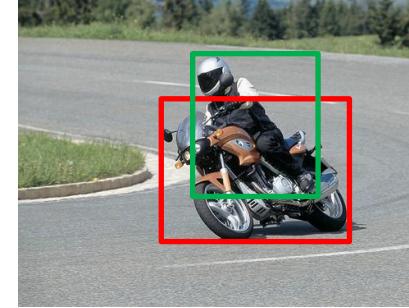
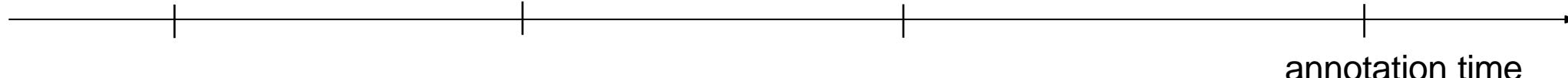


Why Weak Supervision

Weak supervision: How to save costs on labelling.



{motorbike, person}

{motorbike (point),
person (point)}{motorbike (b-box),
person (b-box)}{motorbike (pixel labels),
person (pixel labels)}1 sec
per class2.4 sec
per instance10 sec
per instance78 sec
per instance

Weakly Supervised Learning in Computer Vision



Classification



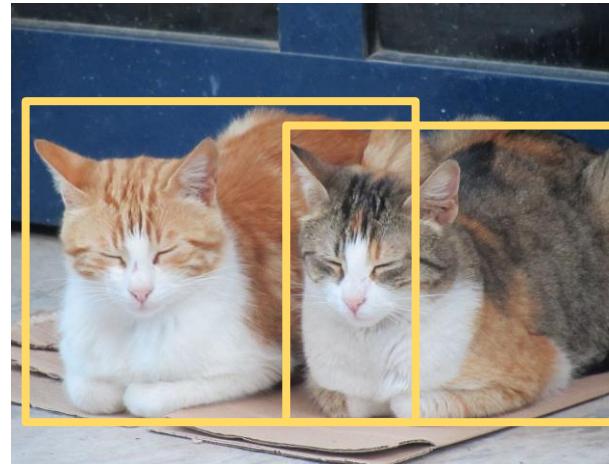
Is there a cat?

Classification



Is there a cat?

Detection



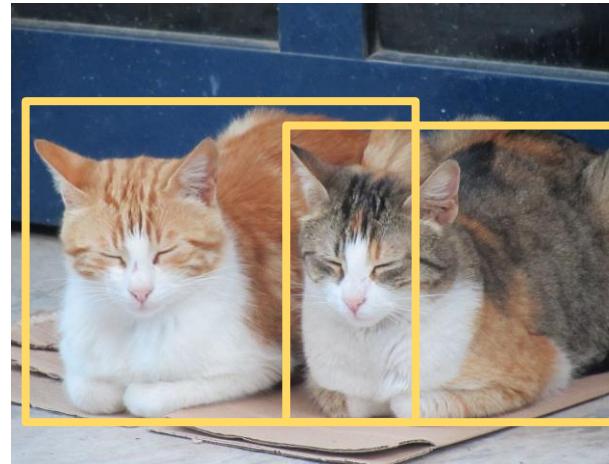
Where is the cat?

Classification



Is there a cat?

Detection



Where is the cat?

Semantic
segmentation



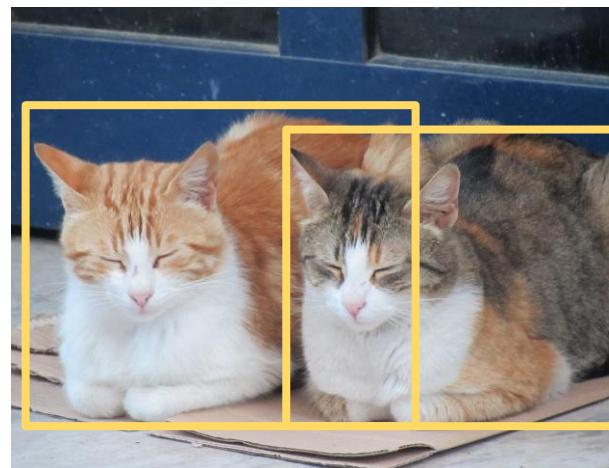
Which pixel is cat ?

Classification



Is there a cat?

Detection



Where is the cat?

Semantic
segmentation



Which pixel is cat ?

Instance
segmentation



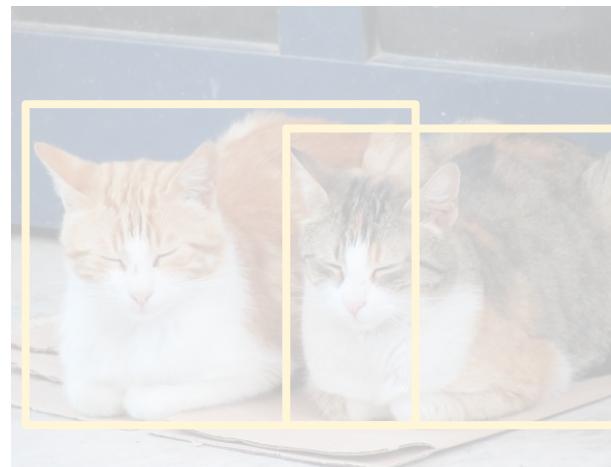
Which are the
cat's pixels?

Classification



Is there a cat?

Detection



Where is the cat?

Semantic
segmentation



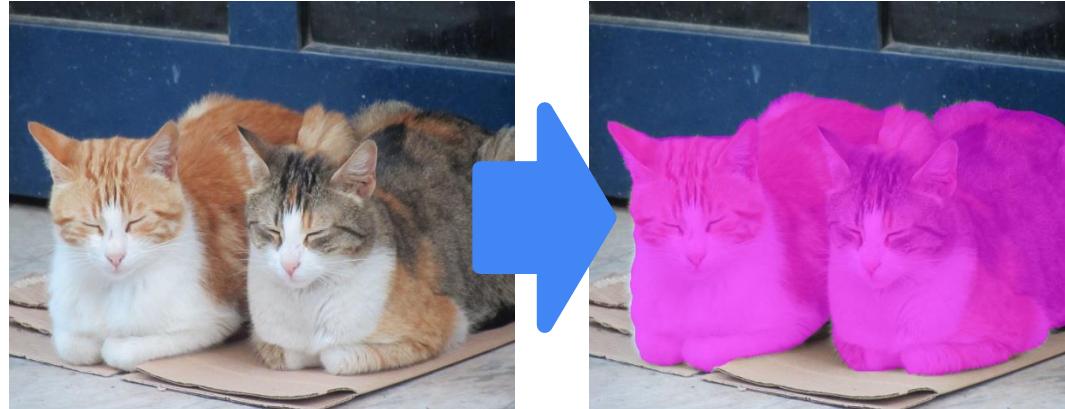
Which pixel is cat ?

Instance
segmentation



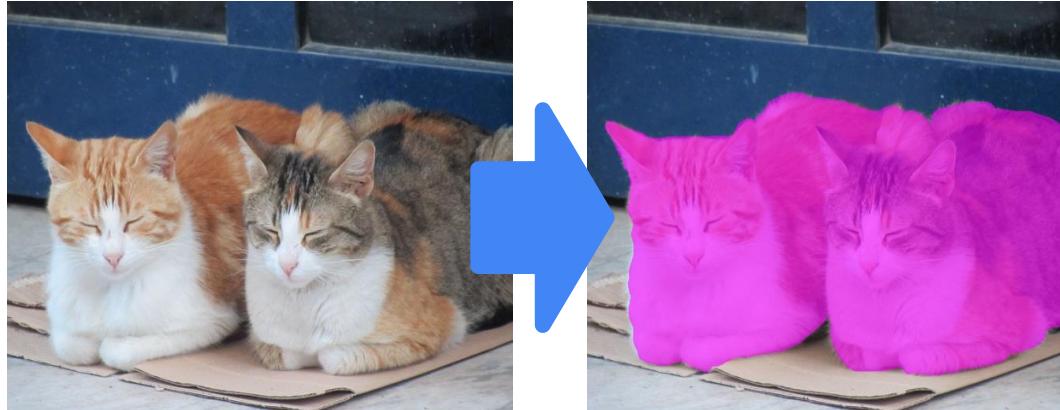
Which are the
cat's pixels?

Fully supervised → **Directly** supervised



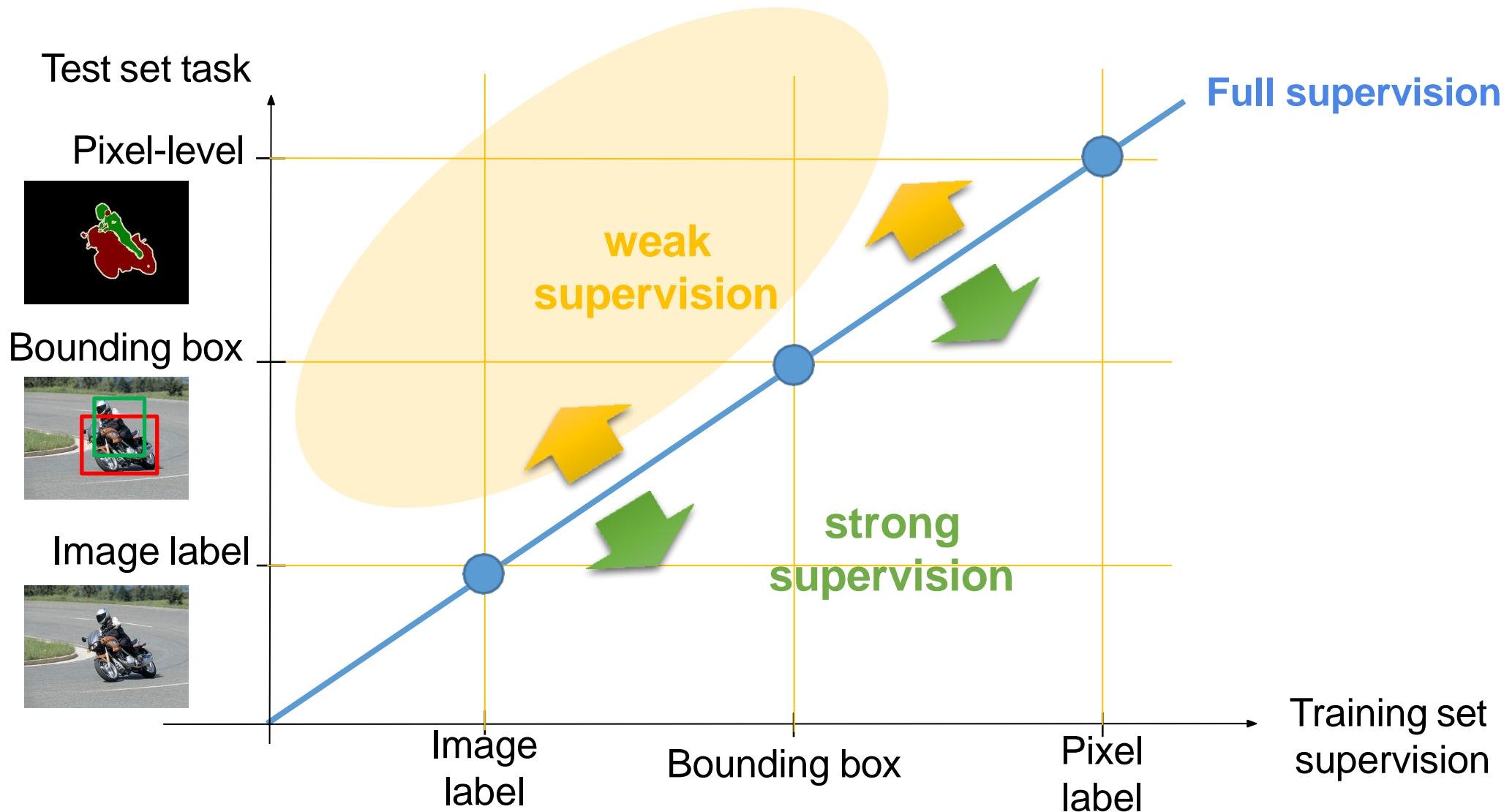
For each input, the desired output is provided

Fully supervised → **Directly** supervised



Weakly supervised → **Indirectly** supervised





Priors + Hints =

Weakly Supervised Algorithms

1. Size of the Object



Which cat is in the right size?

2. Shape of the Object



Which cat is in the right shape?

3. Object location



Which cat is in the right location?

4. Colour Contrast



Which pixels separate cats?

1. Image Labels



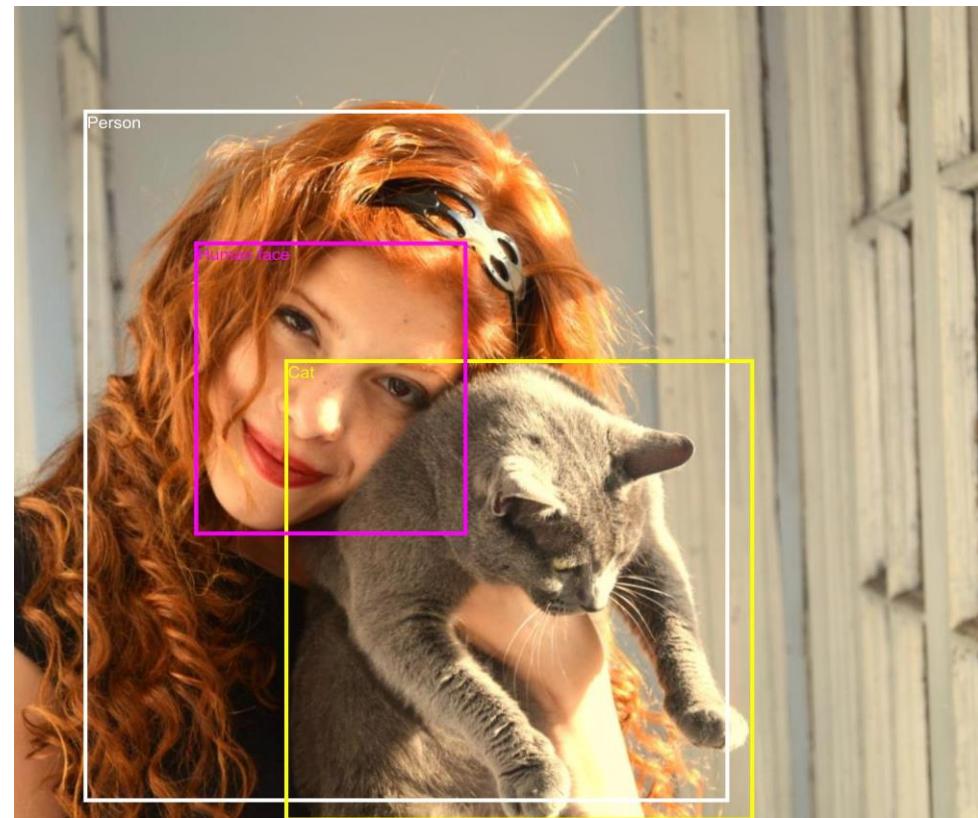
Person, cat, door

2. Image Captions

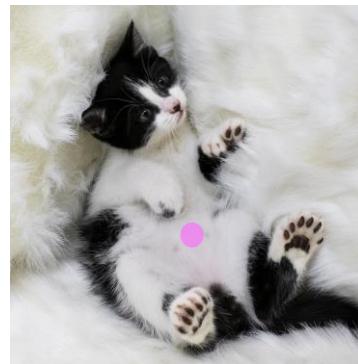
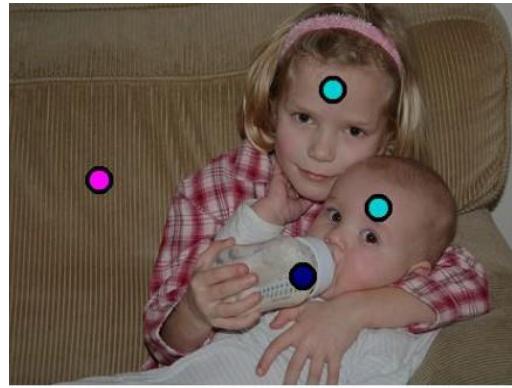


Kattie holds the cat next to the door

3. Object Bounding Boxes



4. Object Points



Click inside the objects

5. Scribbles



Image-level weakly-supervised semantic segmentation

32

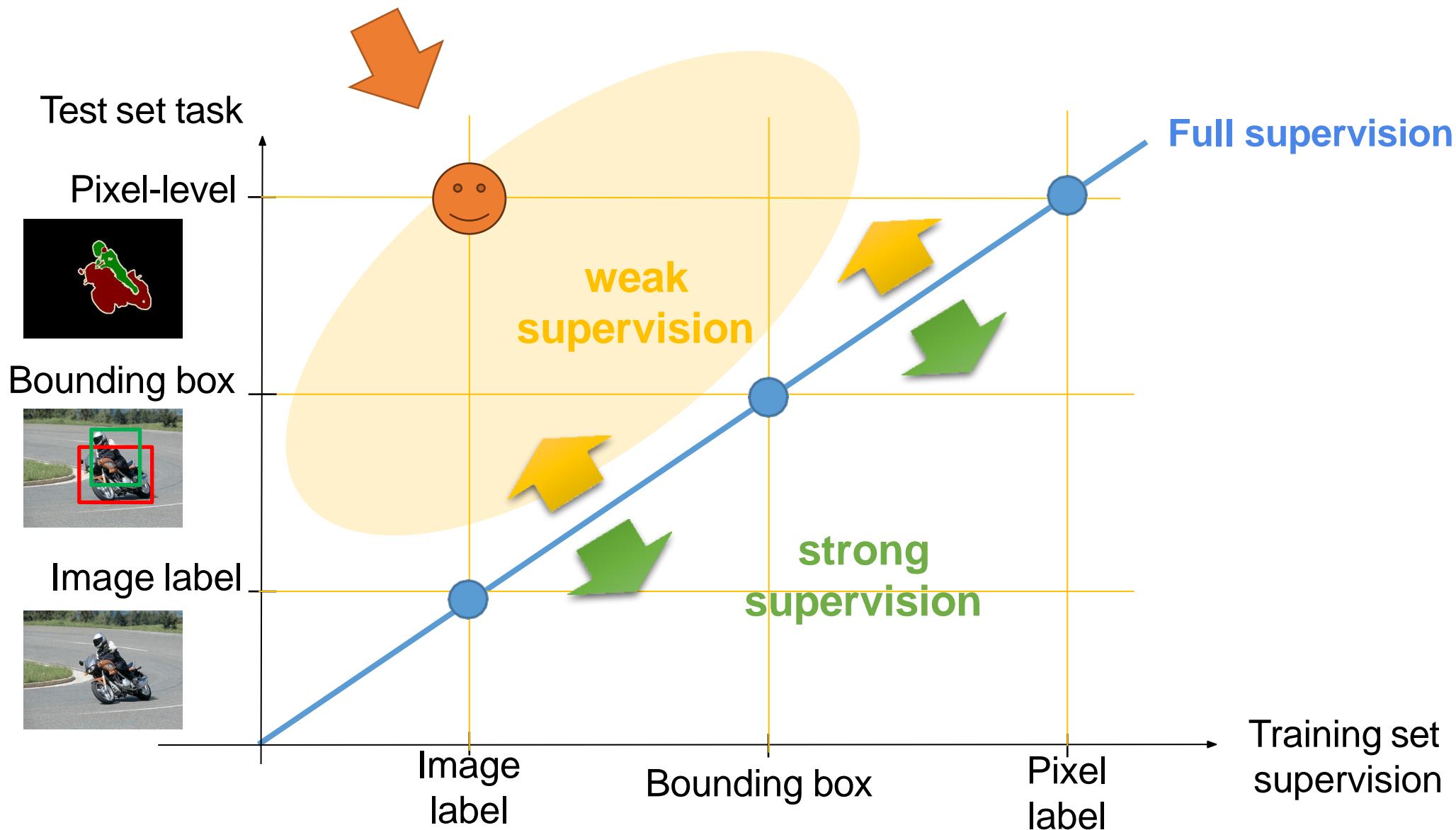
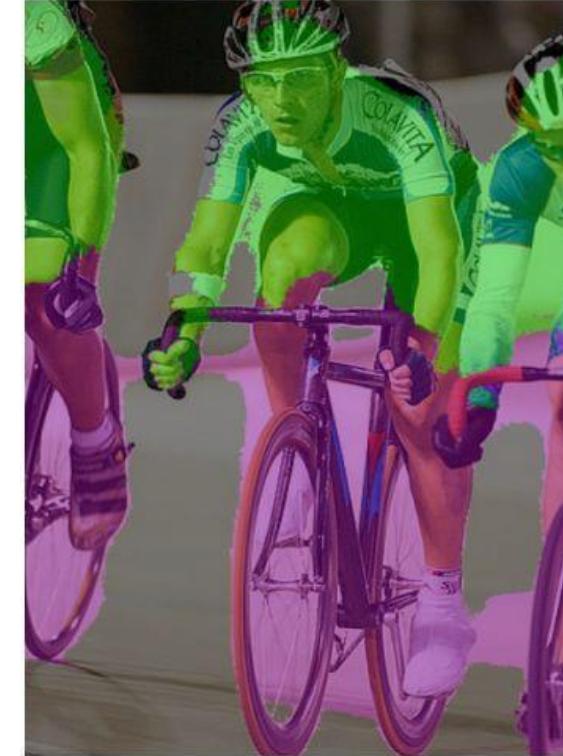
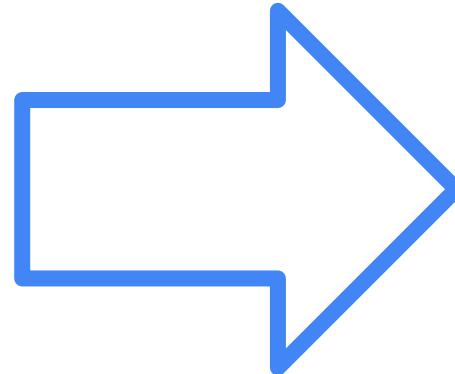




Image-level labels
at training time



Pixel-level labels
at test time

Early Works

34

Smilkov et al. 2017

Step 1: Train an image classifier

Step 2: **Add small random noise** to image

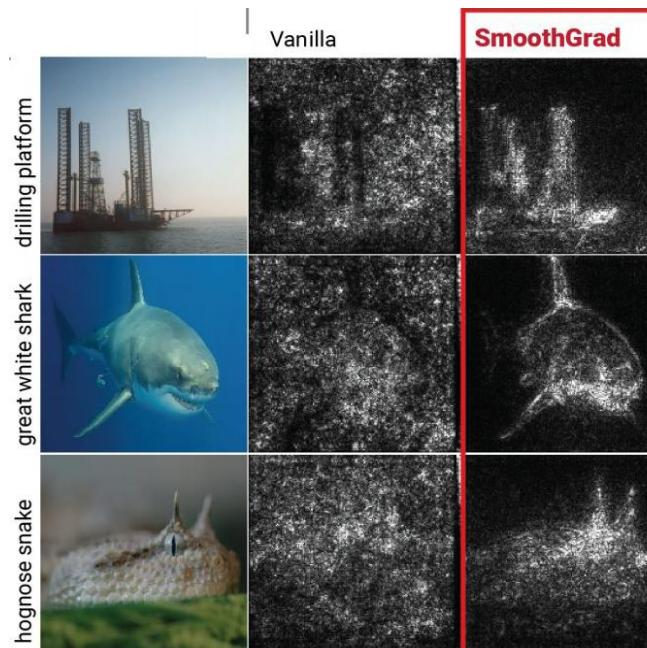
Step 3: given known label,

compute

$$\frac{\partial S_c}{\partial I} \Big|_{I_0}$$

("how much does changes in the input image affect the score?")

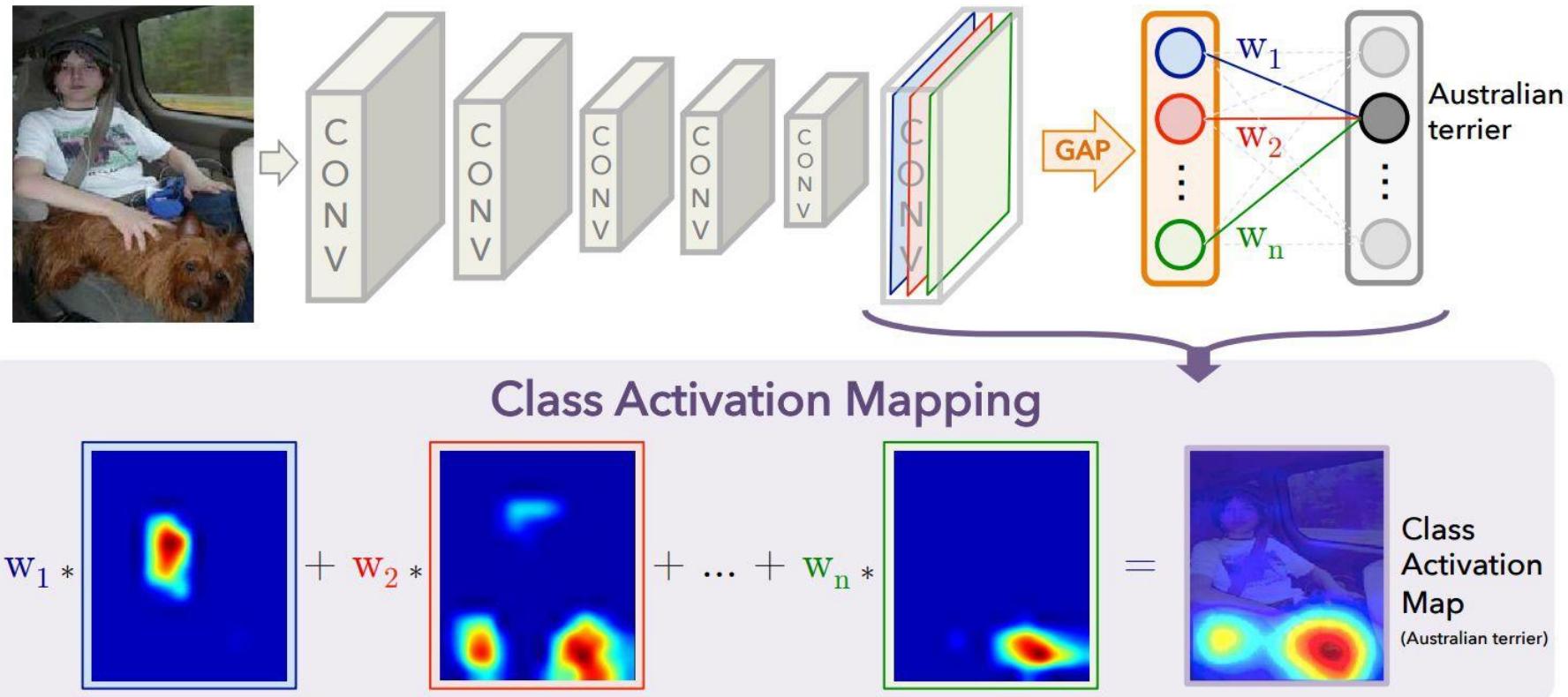
Step 4: **Repeat 50 times step 2&3, average result**



Early Works

35

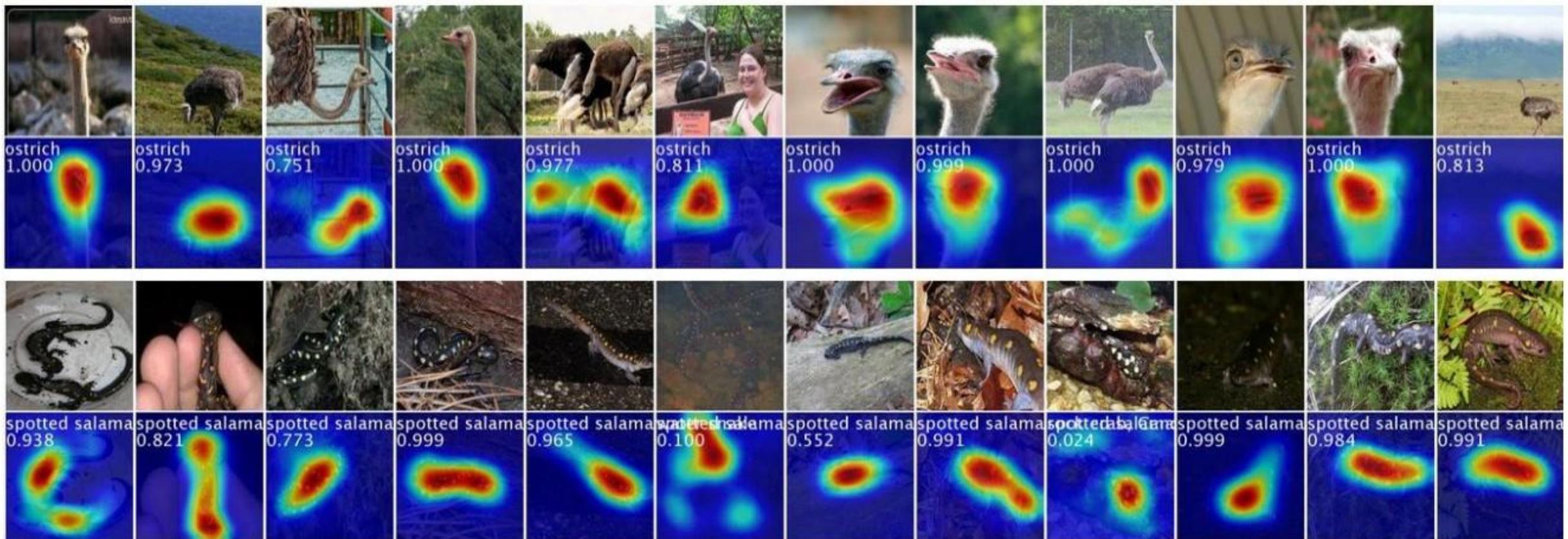
Zhou et al. 2016



$$\text{CAM: } M_c(x, y) = \sum_k w_k^c f_k(x, y).$$

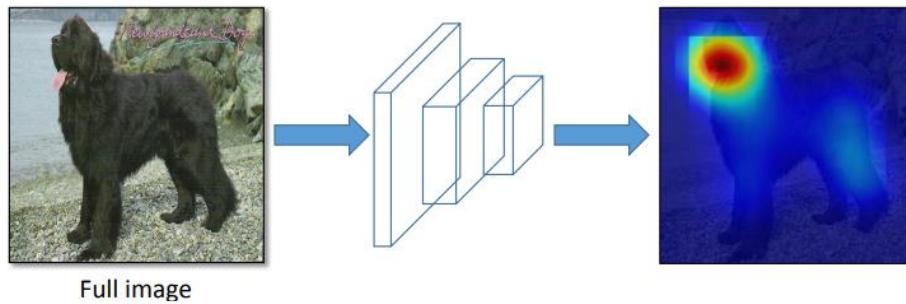
Good, but not good enough:

1. Resulting masks are not sharp
2. Focused on discriminative area only



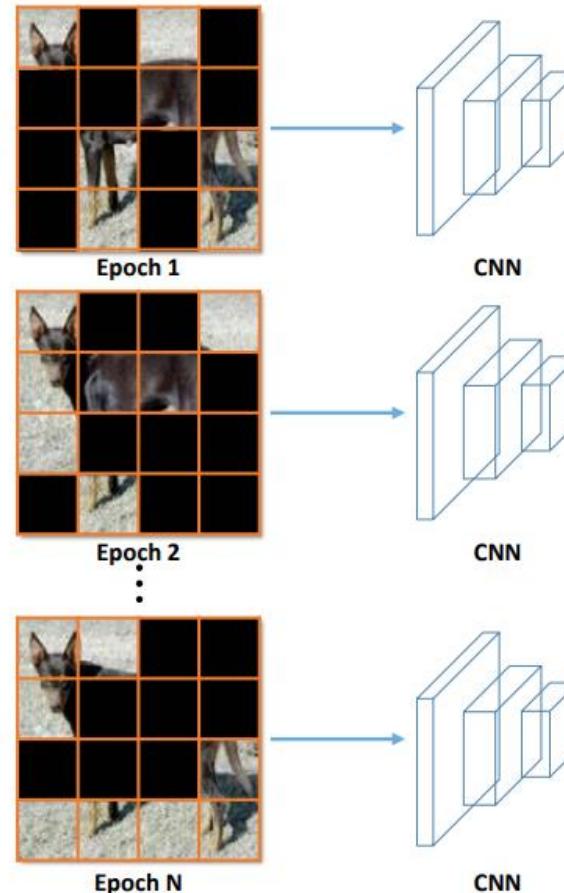
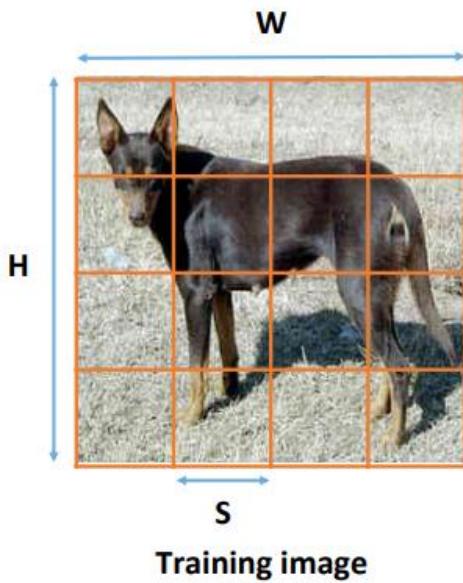
Improvement

37



Hide-and-Seek Algorithm

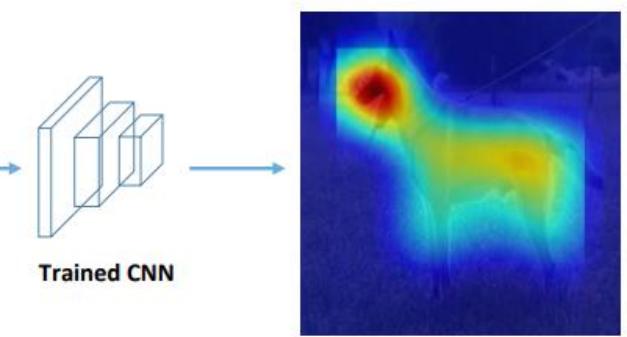
Training phase



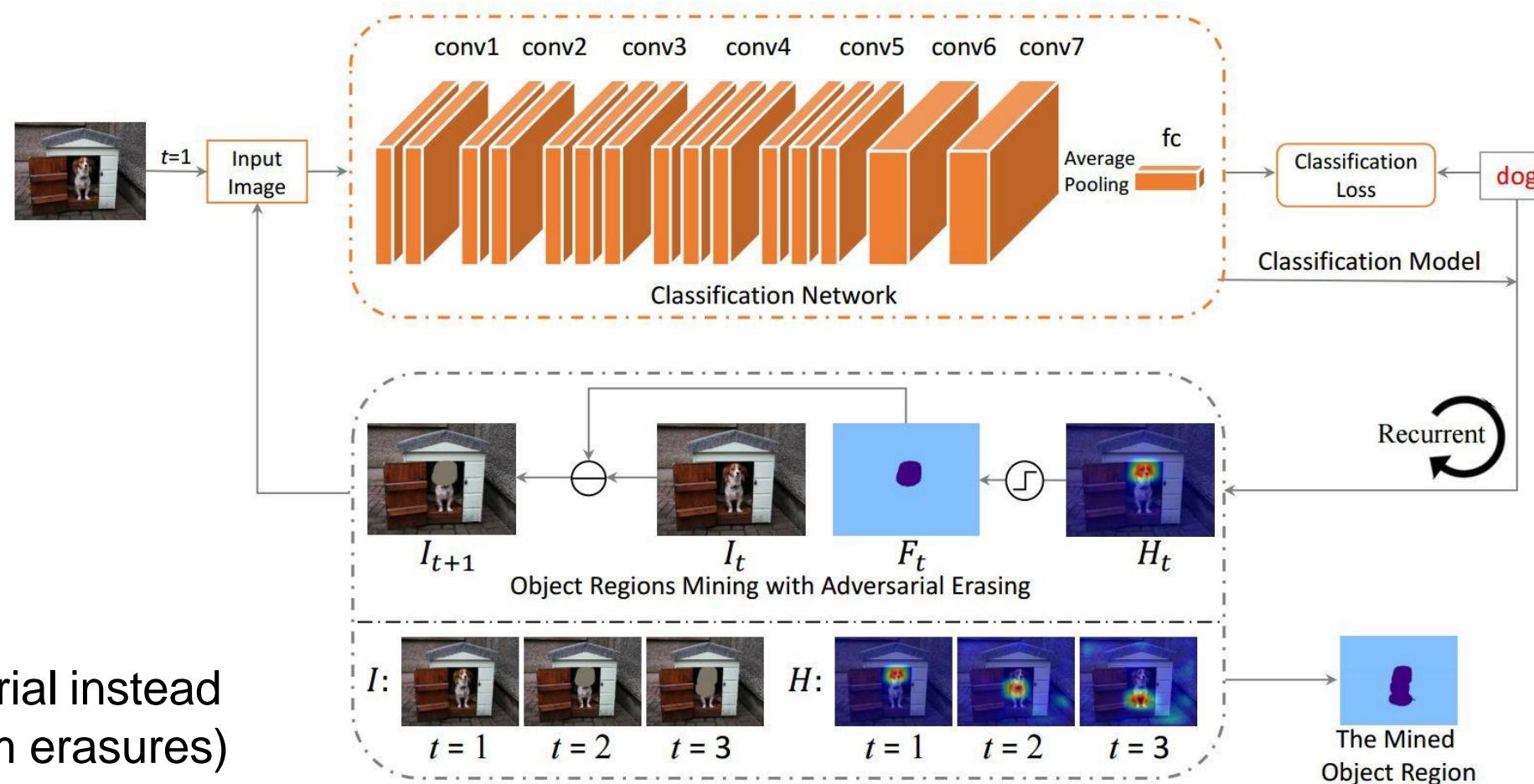
Testing phase



Test image
(no hidden patches)

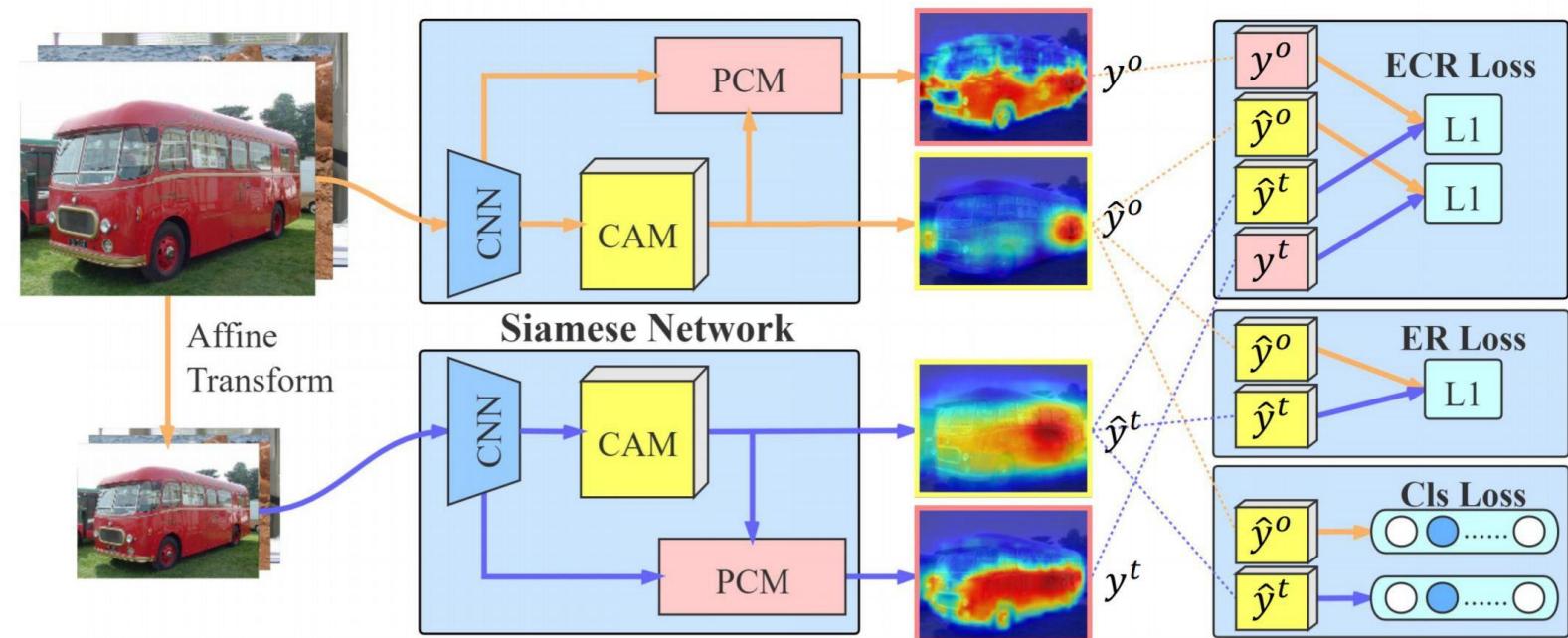
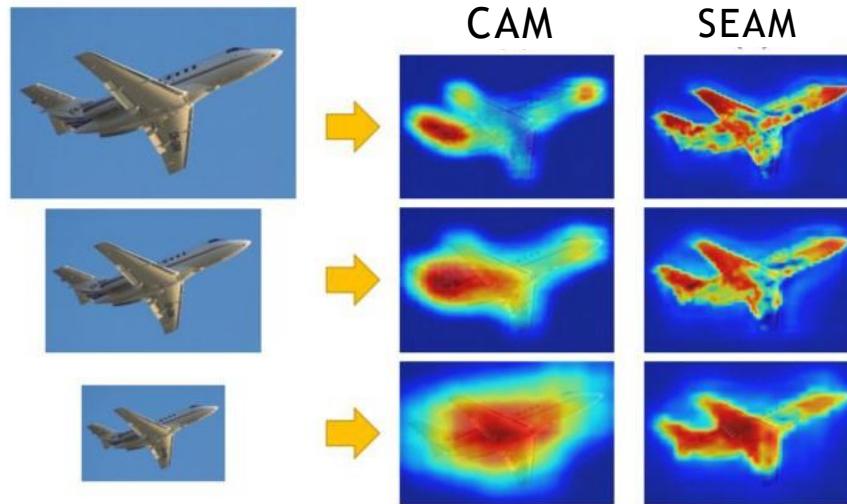


Adversarial Erasing



Invariant to Geo-transformation

"Activation maps should be equivariant to the input"

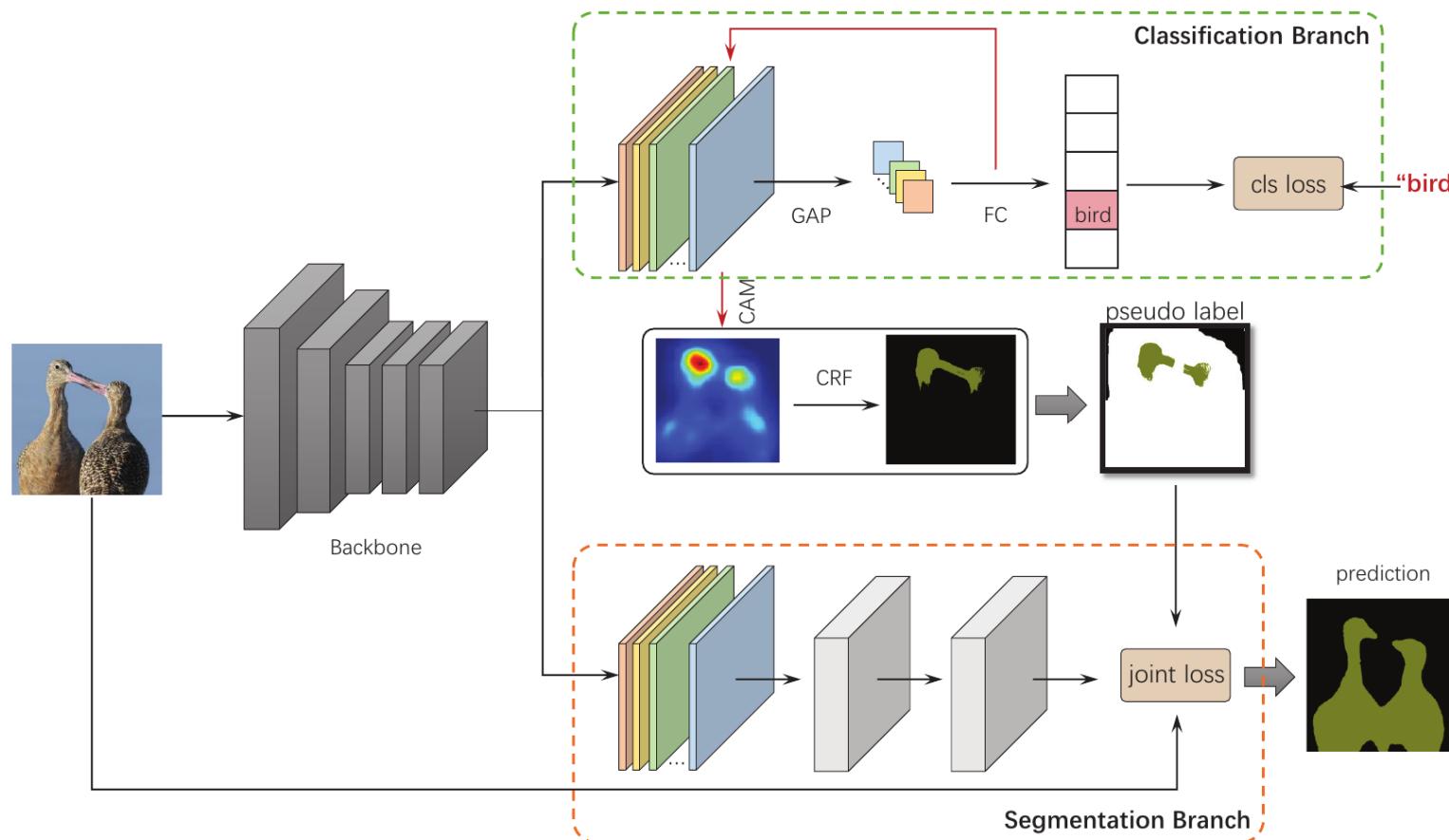


$$\mathcal{R}_{ER} = \|F(A(I)) - A(F(I))\|_1.$$

How to do weakly-supervised segmentation

- Step1: train an image classifier.
- Step2: generate activation maps.
- Step3: refine the activation maps and produce pseudo-labels
- Step4: train segmentation model using pseudo-labels

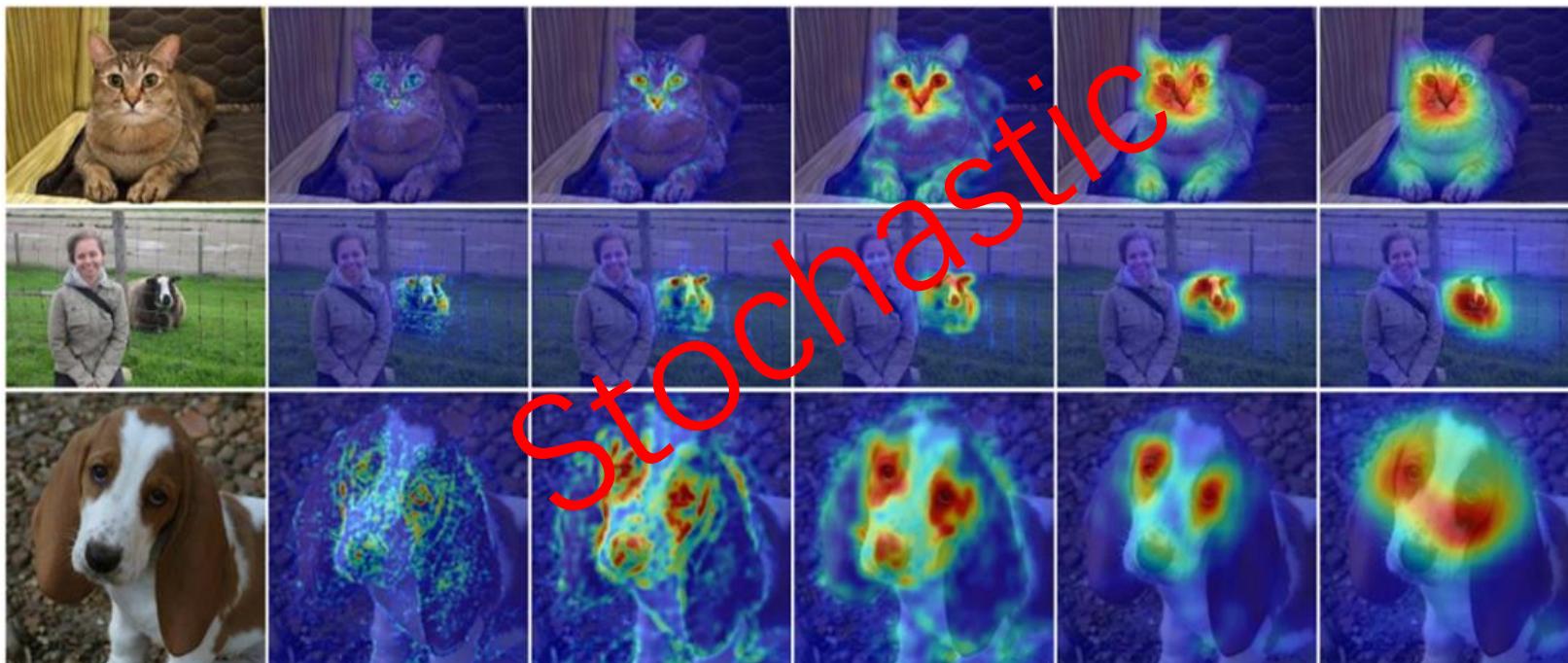
Unified Model for weakly supervised segmentation



Further improvement

CAM is

- Incomplete activation
- Not robust to input variation

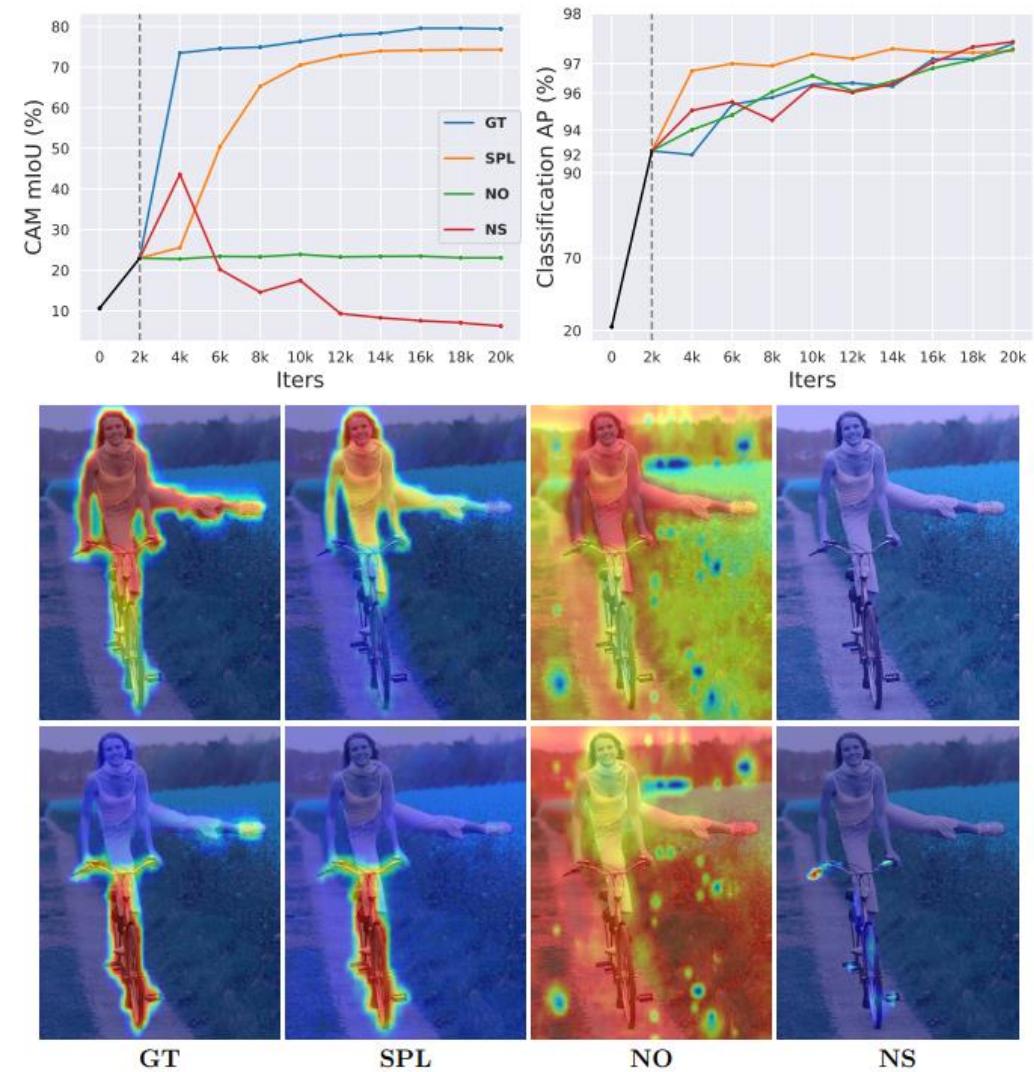


Proof of Concept

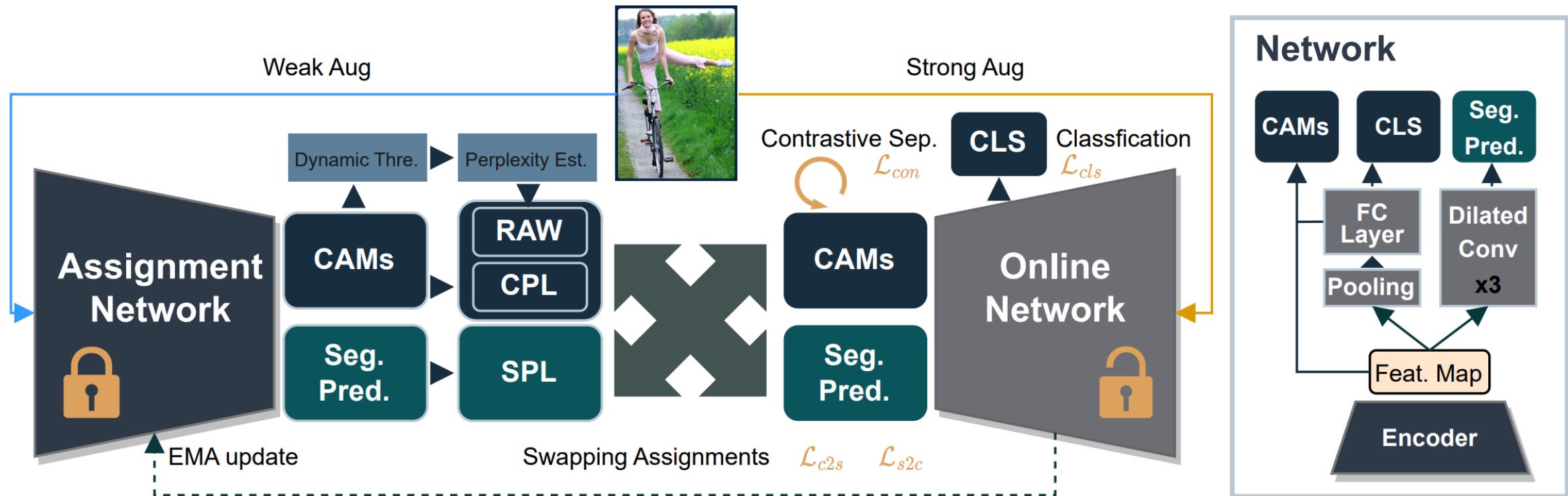
- CAM is differentiable
- But not fully optimized
- Thus it shows stochastic behaviours

Experiments:

- GT: guided by ground truth (pixel-label)
 - NO: no guidance
 - NS: guided by random noise.
-
- But cannot use pixel-label
 - Segmentation pseudo-label (SPL) as the substitute



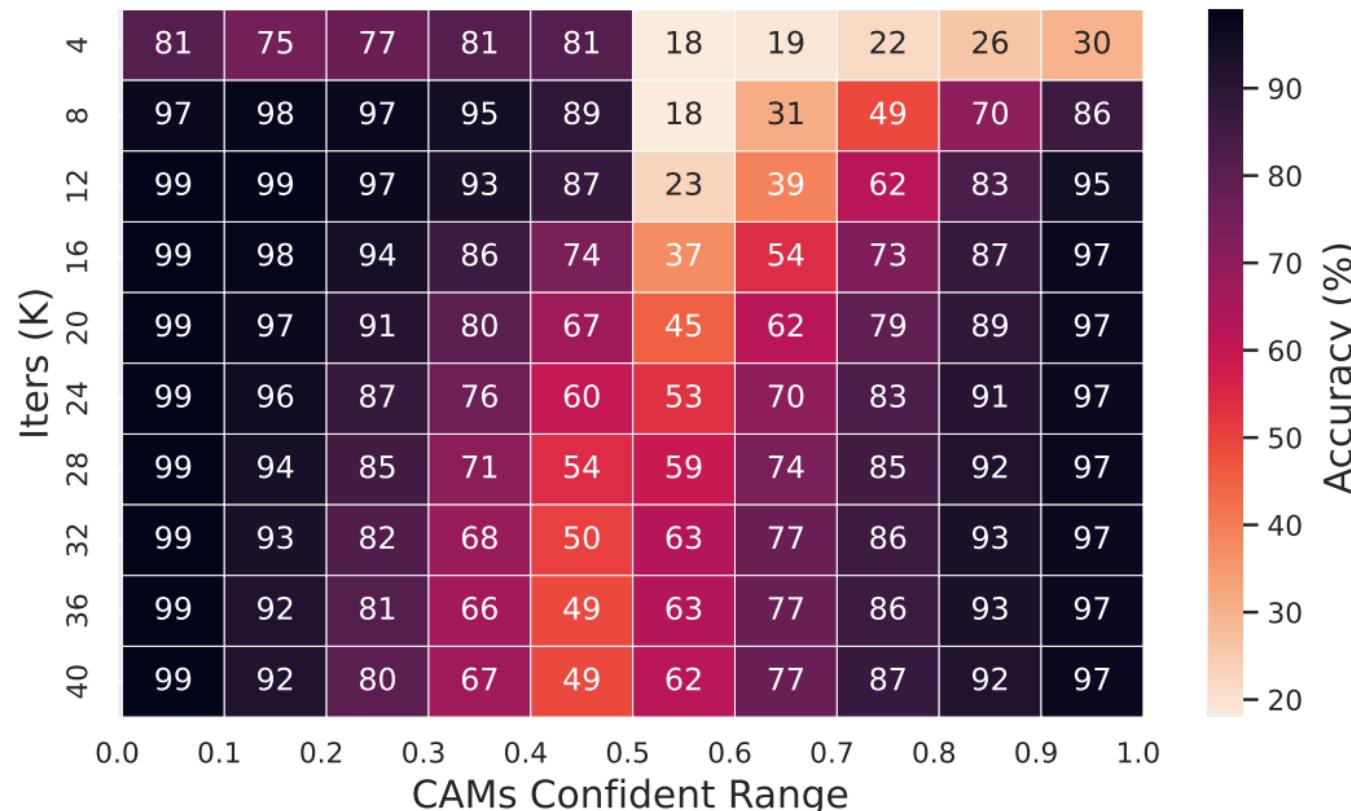
Our Recent Work



Proposed Solution

Reliability based Adaptive Weighting

- Pseudo-label is noisy label
- Need to be accessed before usage



Reliability based Adaptive Weighting (RAW)

- We can define the perplexity as

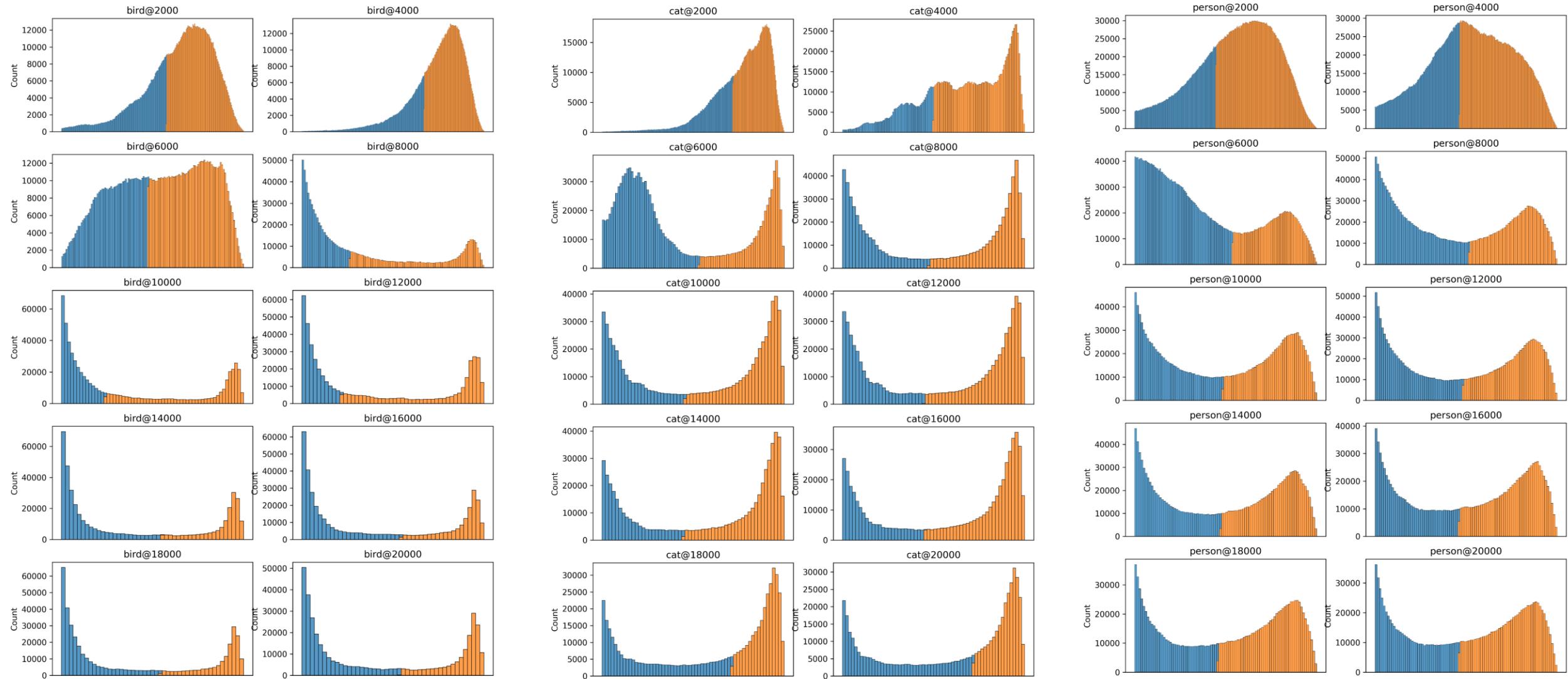
$$\mathcal{P}_{x,y} = \begin{cases} \left[-\log \left(\lambda_\alpha \frac{\max(\mathcal{M}'_{x,y}) - \xi}{1 - \xi} \right) \right]^{\lambda_\beta} & \text{if } \max(\mathcal{M}'_{x,y}) \geq \xi, \\ \left[-\log \left(\lambda_\alpha \frac{\xi - \max(\mathcal{M}'_{x,y})}{\xi} \right) \right]^{\lambda_\beta} & \text{if } \max(\mathcal{M}'_{x,y}) < \xi, \end{cases}$$

- RAW is the normalized reciprocal of perplexity

$$\mathcal{W}_{x,y}^{\text{raw}} = |\mathcal{R}| / \sum_{i,j \in \mathcal{R}} \frac{1}{\mathcal{P}_{i,j} \mathcal{P}_{x,y}},$$

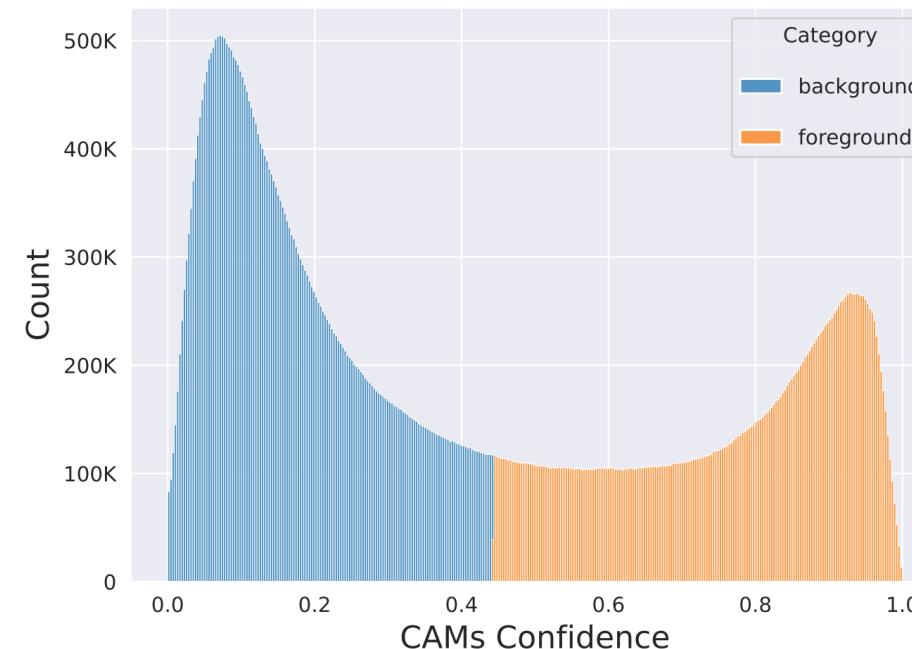
Dynamic Threshold

- To generate CAMs pseudo-labels (CPL) a threshold is need.
- Previous works use a fixed threshold.



Dynamic Threshold

- To generate CAMs pseudo-labels (CPL) a threshold is need.
- Previous works use a fixed threshold.
- The bi-modal distribution



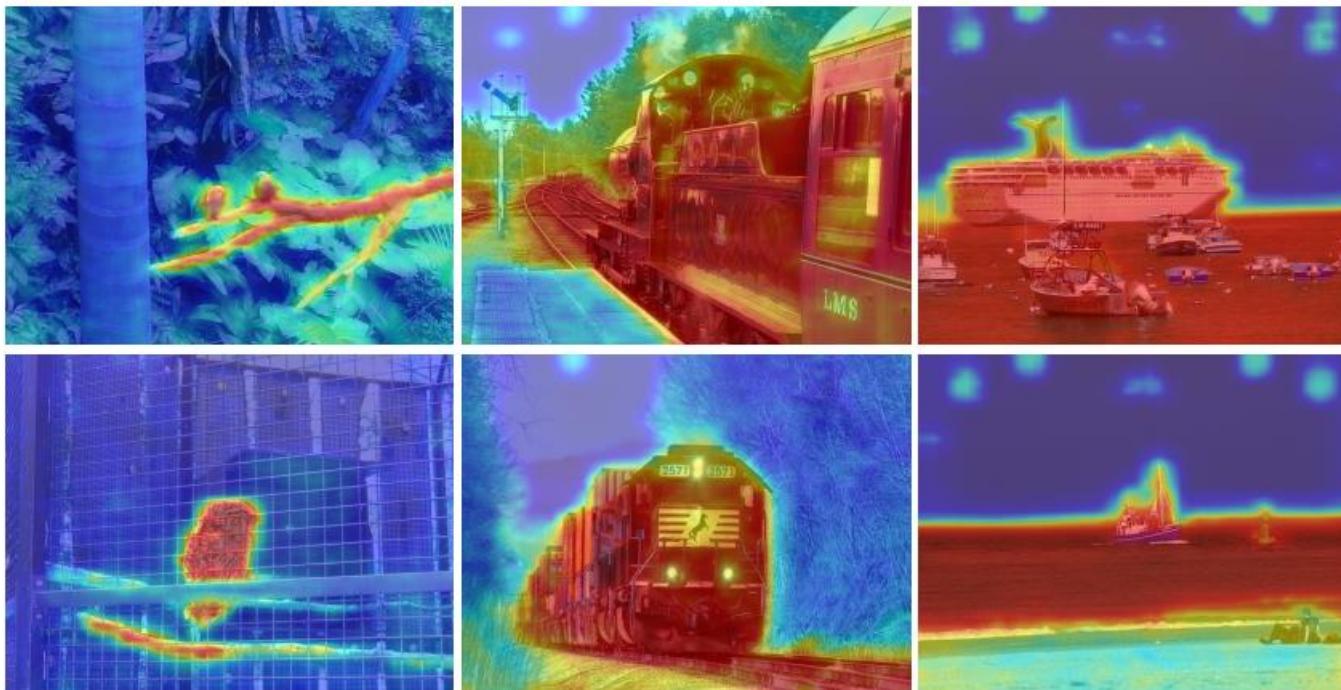
Dynamic Threshold

- To generate CAMs pseudo-labels (CPL) a threshold is need.
- Previous works use a fixed threshold.
- The bi-modal distribution
- Propose to use Gaussian Mixture Model (GMM) to find optimal threshold

$$\begin{aligned}\xi^* = \operatorname{argmax}_{\xi} & \prod_{x \in \{\mathcal{M}' \geq \xi\}} \tilde{\pi}_{fg} \mathcal{N}(x | \tilde{\mu}_{fg}, \tilde{\Sigma}_{fg}) \\ & + \prod_{x \in \{\mathcal{M}' < \xi\}} \tilde{\pi}_{bg} \mathcal{N}(x | \tilde{\mu}_{bg}, \tilde{\Sigma}_{bg}) ,\end{aligned}$$

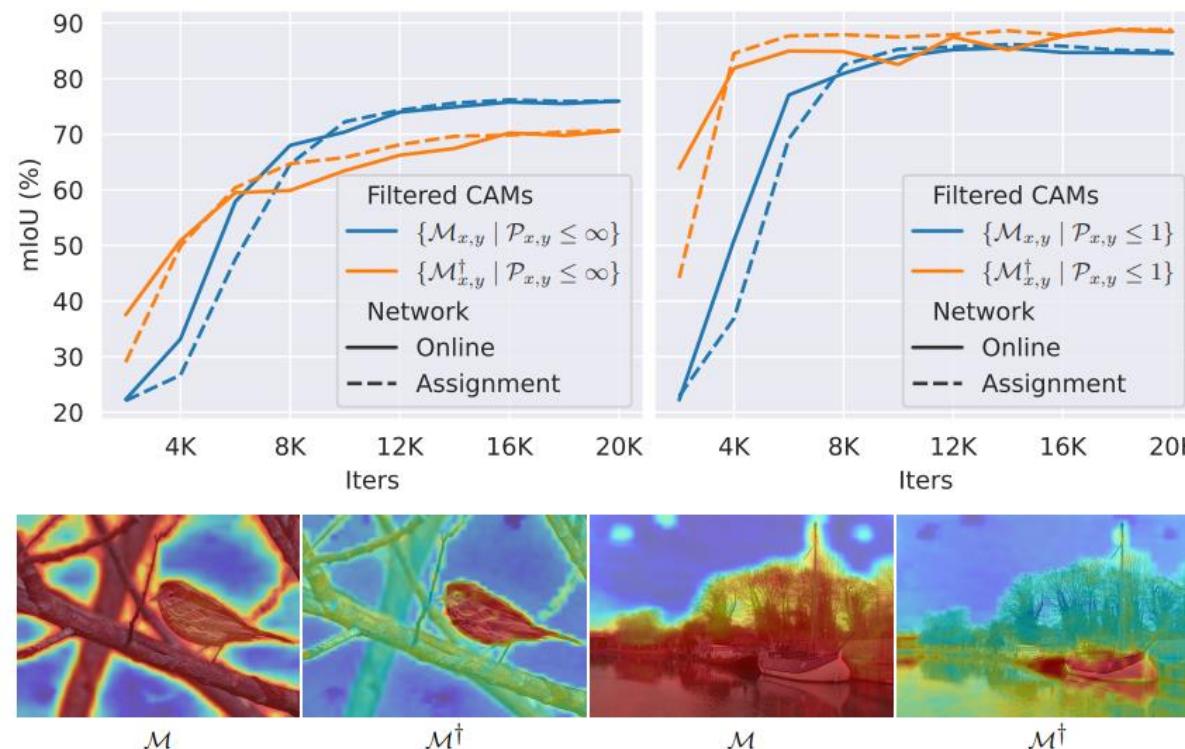
Contrastive Separation in CAMs

- Coexistence Problem in CAMs



Contrastive Separation in CAMs

- Auxiliary CAMs (from early stage of the model)



Contrastive Separation in CAMs

- Use Auxiliary CAMs filter by low perplexity to define some positive regions and negative regions

$$\begin{aligned}\mathcal{R}_{i,j}^+ &= \left\{ (x, y) \mid \mathcal{P}_{x,y} \leq \epsilon, \hat{y}_{x,y}^{\text{CPL}} = \hat{y}_{i,j}^{\text{CPL}}, x \neq i, y \neq j \right\}, \\ \mathcal{R}_{i,j}^- &= \left\{ (x, y) \mid \mathcal{P}_{x,y} \leq \epsilon, \hat{y}_{x,y}^{\text{CPL}} \neq \hat{y}_{i,j}^{\text{CPL}} \right\},\end{aligned}$$

- Contrastively separate these regions in CAMs by minimising the distance among the positives and maximising the distance among negatives.

$$\begin{aligned}\mathcal{L}_{con} = - \frac{1}{|\{(x, y) \mid \mathcal{P}_{x,y} \leq \epsilon\}|} \sum_{i,j \in \{(x, y) \mid \mathcal{P}_{x,y} \leq \epsilon\}} \frac{1}{|\mathcal{R}_{i,j}^+|} \sum_{x,y \in \mathcal{R}_{i,j}^+} \\ \left[\log \frac{\exp(l_d(\mathcal{M}_{i,j}, \mathcal{M}_{x,y}) / \tau)}{\exp(l_d(\mathcal{M}_{i,j}, \mathcal{M}_{x,y}) / \tau) + \sum_{n,m \in \mathcal{R}_{i,j}^-} \exp(l_d(\mathcal{M}_{i,j}, \mathcal{M}_{n,m}) / \tau)} \right]\end{aligned}$$

Results

Results on VOC & COCO

- State-of-the-art performance compared with the recent methods.
- Our single-stage model outperforms multi-stage models and models that use additional modalities.

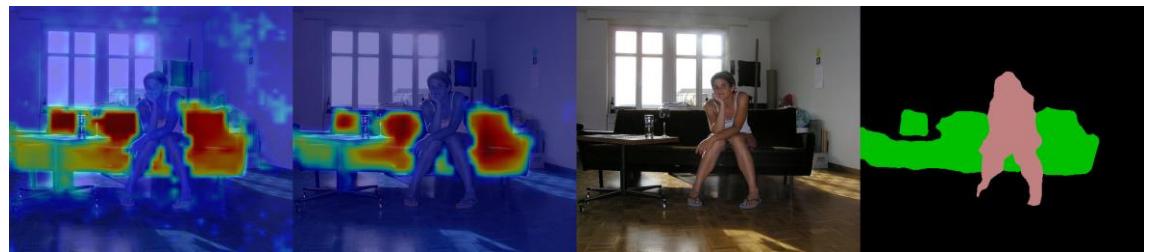
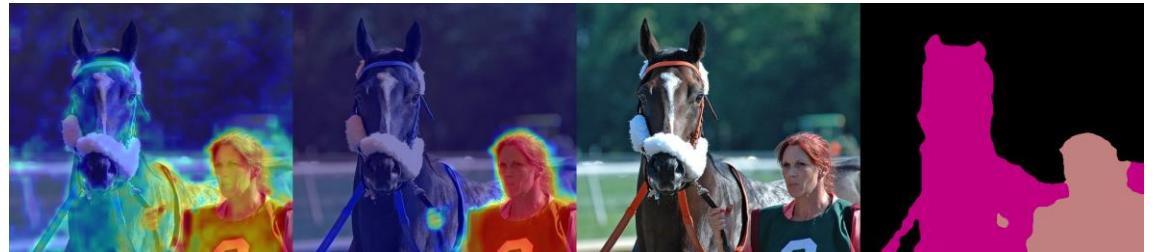
Methods	Sup.	Net.	VOC		COCO
			val	test	val
<i>Supervised Upperbounds.</i>					
Deeplab [10] <small>TPAMI'2017</small>	\mathcal{F}	R101	77.6	79.7	–
WideRes38 [59] <small>PR'2019</small>	\mathcal{F}	WR38	80.8	82.5	–
ViT-Base [17] <small>ICLR'2021</small>	\mathcal{F}	ViT-B	80.5	81.0	–
UperNet-Swin [42] <small>ICCV'2021</small>	\mathcal{F}	SWIN	83.4	83.7	–
<i>Multi-stage Methods.</i>					
L2G [26] <small>CVPR'2022</small>	$\mathcal{I} + \mathcal{S}$	R101	72.1	71.7	44.2
Du <i>et al.</i> [18] <small>CVPR'2022</small>	$\mathcal{I} + \mathcal{S}$	R101	72.6	73.6	–
CLIP-ES [40] <small>CVPR'2023</small>	$\mathcal{I} + \mathcal{L}$	R101	73.8	73.9	45.4
ESOL [37] <small>NeurIPS'2022</small>	\mathcal{I}	R101	69.9	69.3	42.6
BECO [48] <small>CVPR'2023</small>	\mathcal{I}	R101	72.1	71.8	45.1
Mat-Label [57] <small>ICCV'2023</small>	\mathcal{I}	R101	73.0	72.7	45.6
CoSA-MS	\mathcal{I}	R101	76.5	75.3^[1]	50.9
Xu <i>et al.</i> [63] <small>CVPR'2023</small>	$\mathcal{I} + \mathcal{L}$	WR38	72.2	72.2	45.9
W-OoD [33] <small>CVPR'2022</small>	\mathcal{I}	WR38	70.7	70.1	–
MCT [62] <small>CVPR'2022</small>	\mathcal{I}	WR38	71.9	71.6	42.0
ex-ViT [66] <small>PR'2023</small>	\mathcal{I}	WR38	71.2	71.1	42.9
ACR-ViT [31] <small>CVPR'2023</small>	\mathcal{I}	WR38	72.4	72.4	–
MCT+OCR [15] <small>CVPR'2023</small>	\mathcal{I}	WR38	72.7	72.0	42.0
CoSA-MS	\mathcal{I}	WR38	76.6	74.9^[2]	50.1
ReCAM [14] <small>CVPR'2022</small>	\mathcal{I}	SWIN	70.4	71.7	47.9
LPCM [12] <small>CVPR'2023</small>	\mathcal{I}	SWIN	73.1	73.4	48.3
CoSA-MS	\mathcal{I}	SWIN	81.4	78.4^[3]	53.7
<i>Single-stage (End-to-end) Methods.</i>					
1Stage [3] <small>CVPR'2020</small>	\mathcal{I}	WR38	62.7	64.3	–
RRM [67] <small>AAAI'2020</small>	\mathcal{I}	WR38	62.6	62.9	–
AFA [50] <small>CVPR'2022</small>	\mathcal{I}	ViT-B1	66.0	66.3	38.9
RRM [67] [†] <small>AAAI'2020</small>	\mathcal{I}	ViT-B	63.1	62.4	–
ViT-PCM [49] <small>ECCV'2022</small>	\mathcal{I}	ViT-B	69.3	–	45.0
ToCo [51] <small>CVPR'2023</small>	\mathcal{I}	ViT-B	71.1	72.2	42.3
CoSA	\mathcal{I}	ViT-B	76.2	75.1 ^[4]	51.0
CoSA*	\mathcal{I}	ViT-B	76.4	75.2^[5]	51.1

Results on VOC

Aux. CAM CAM Input Image Segmentation



Aux. CAM CAM Input Image Segmentation

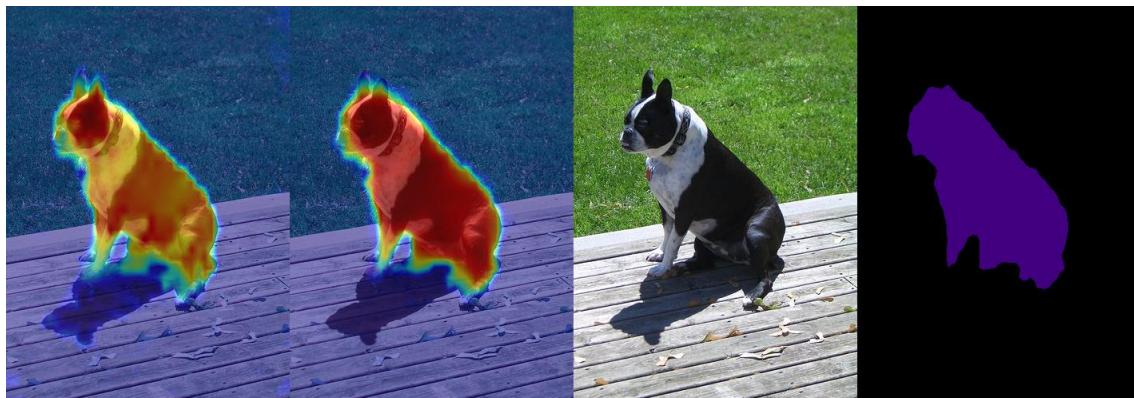
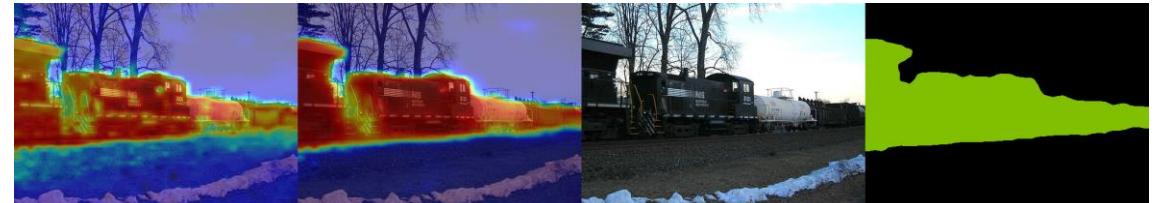


Results on VOC

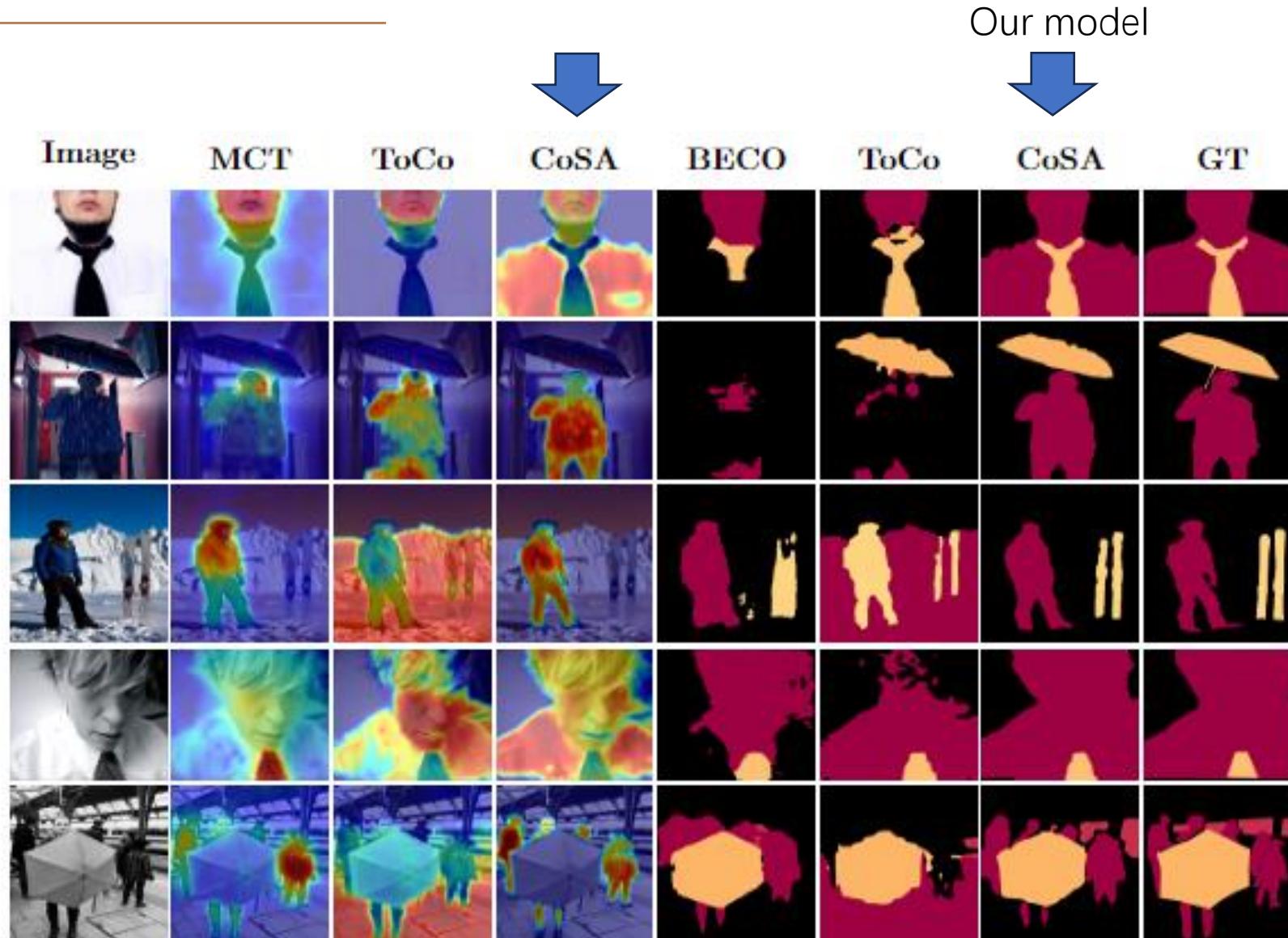
Aux. CAM CAM Input Image Segmentation



Aux. CAM CAM Input Image Segmentation



Results on COCO



H-Unique Datasets



The task is to segment the interesting part in a hand:

Tattoo



Ring



Lunule

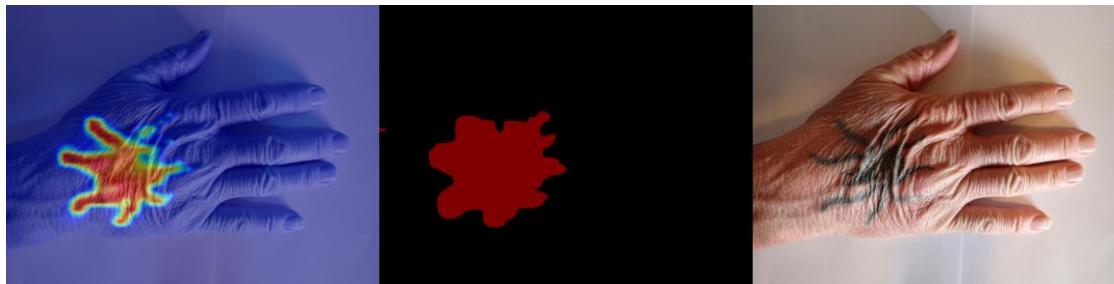


Polished nail



Credit to Ricki for collecting
and annotating data

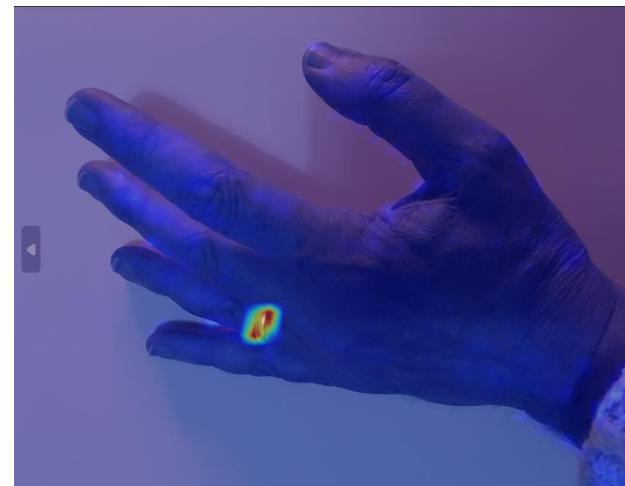
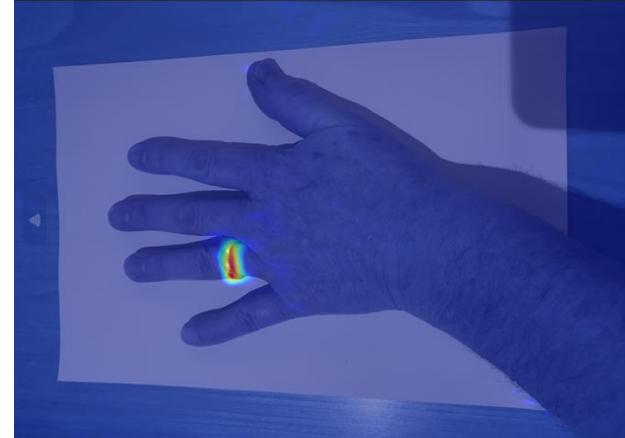
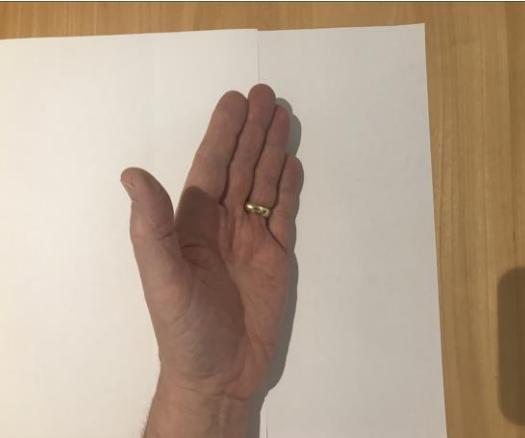
Results on Tattoo segmentation



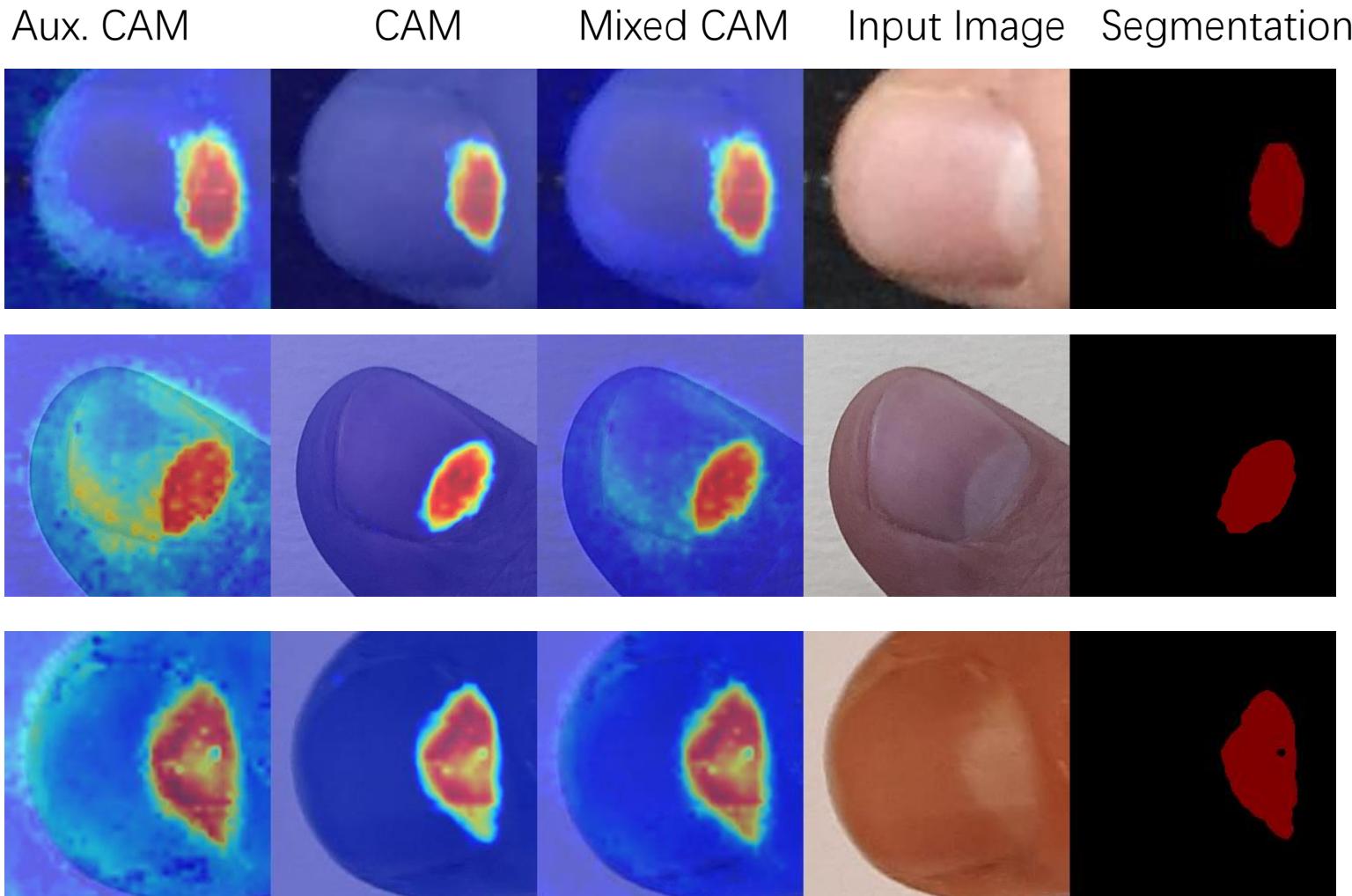
Results on Tattoo segmentation



Results on Ring segmentation

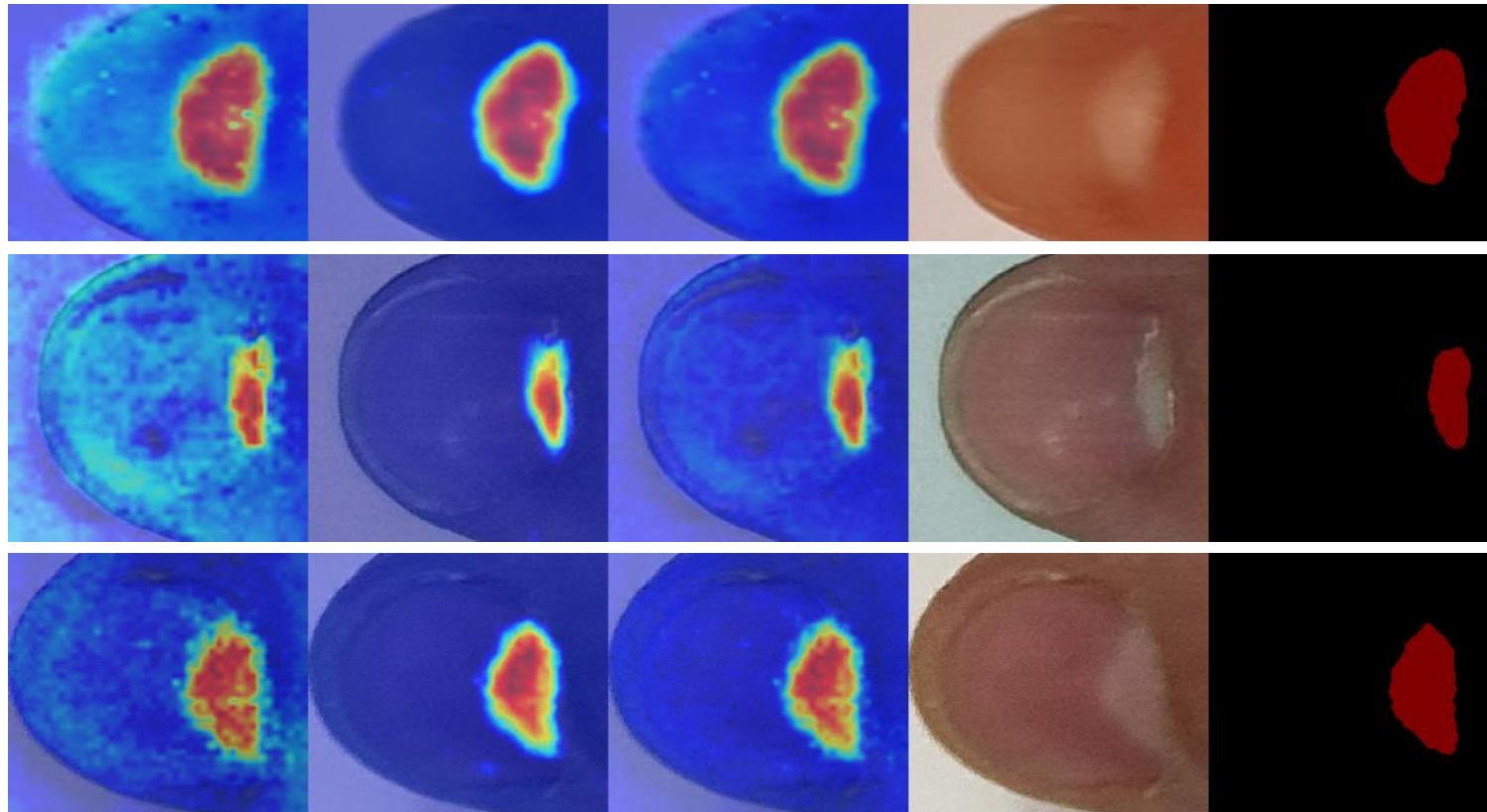


Results on Lunule segmentation

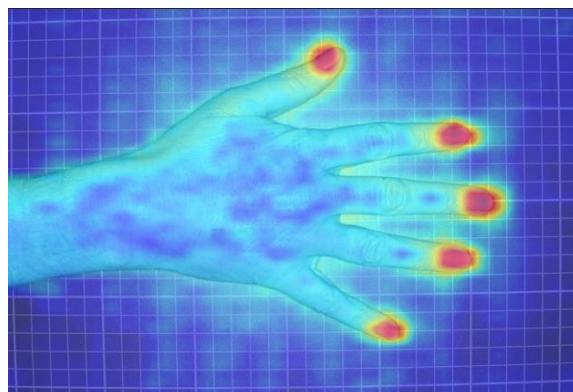
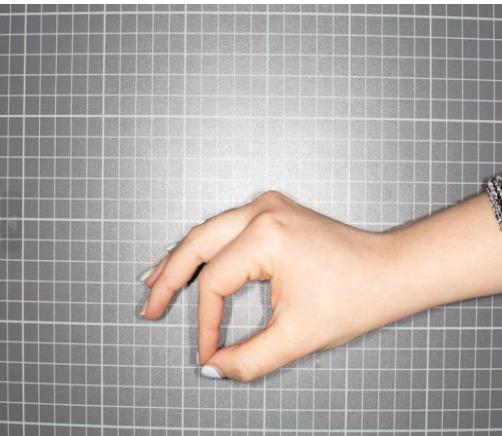
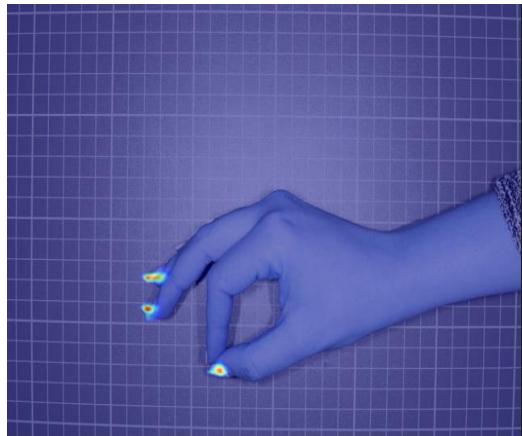
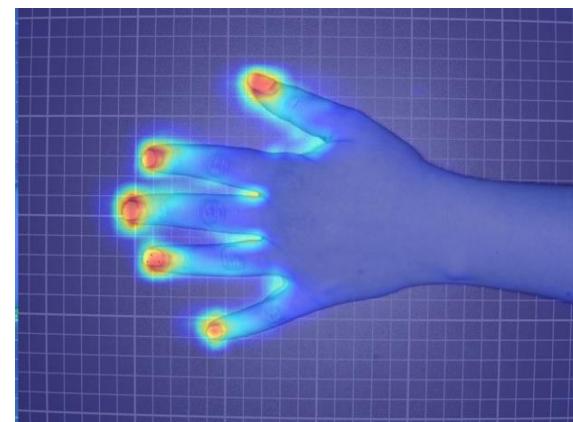
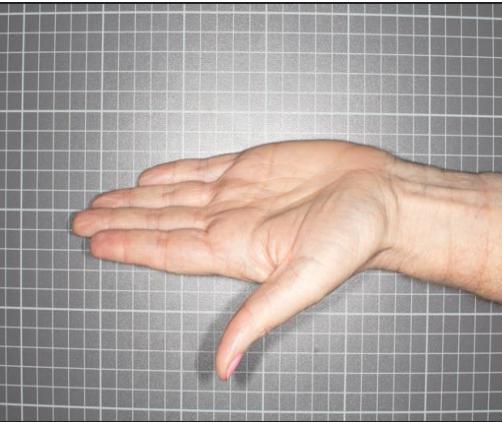
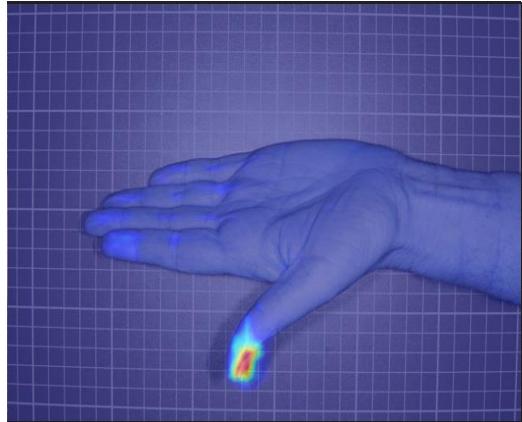


Results on Lunule segmentation

Aux. CAM CAM Mixed CAM Input Image Segmentation



Results on Polished Nail segmentation



- Computer vision requires lots of data; data is expensive.
- Weakly-supervised learning is an active research area and required in practice.
- You use priors and hints to develop weakly-supervised learning algorithm.
- It is a good time to join the field and make contributions in this field!

Thanks