

Deep Learning Models for Weakly-Supervised Object Localization and Segmentation

Soufiane Belharbi, Eric Granger and Ismail Ben Ayed

ICPR 2022 Tutorial
August 21, 2022
Montreal, Canada

Overview

1) Introduction

2) Review of WSOL Methods

- a) bottom-up and top-down methods
- b) case studies

Coffee Break (pause at 10h for 30 mins)

3) Review of WSSS Methods

4) Applications of WSOL / WSSS

5) Key Challenges and Future Directions



LABORATOIRE
D'IMAGERIE, DE VISION
ET D'INTELLIGENCE
ARTIFICIELLE



Belharbi, Soufiane: Post-doc. LIVIA, Dept. of Systems Engineering, ETS Montreal, Canada.

Areas: Machine Learning, Deep Learning, Computer Vision.



Granger, Eric: Professor, LIVIA, Dept. of Systems Engineering, ETS Montreal, Canada.

Areas: Machine Learning, Pattern Recognition, Computer Vision, Information Fusion, Affective Computing, Biometrics, Video Surveillance



Ben Ayed, Ismail: Professor, LIVIA, Dept. of Systems Engineering, ETS Montreal, Canada.

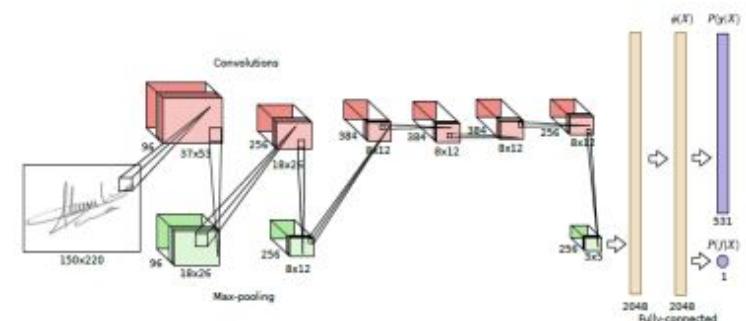
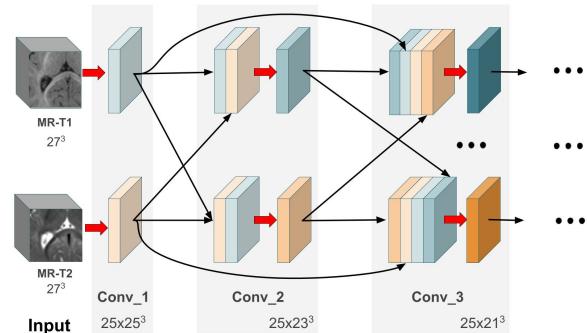
Areas: Computer Vision, Machine Learning, Pattern Recognition, Optimization, Medical Image Analysis, Information Theory

<http://liviamtl.ca/>



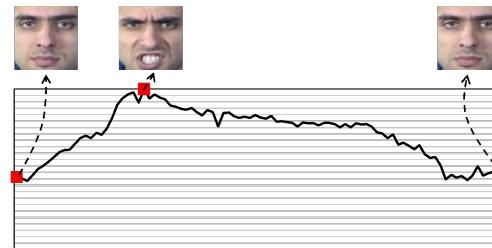
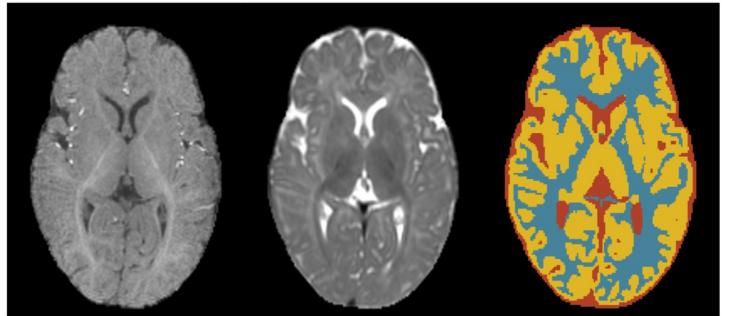
Research Axes:

- machine learning
- computer vision: perception in 2D and 3D scenes
- pattern recognition in static and dynamically-changing environments
- information fusion
- optimization of complex systems



Application Areas:

- analysis of medical, aerial images
- video analytics and surveillance
- biometrics (face, voice and signature)
- document analysis
- affective computing



Part 1

Introduction

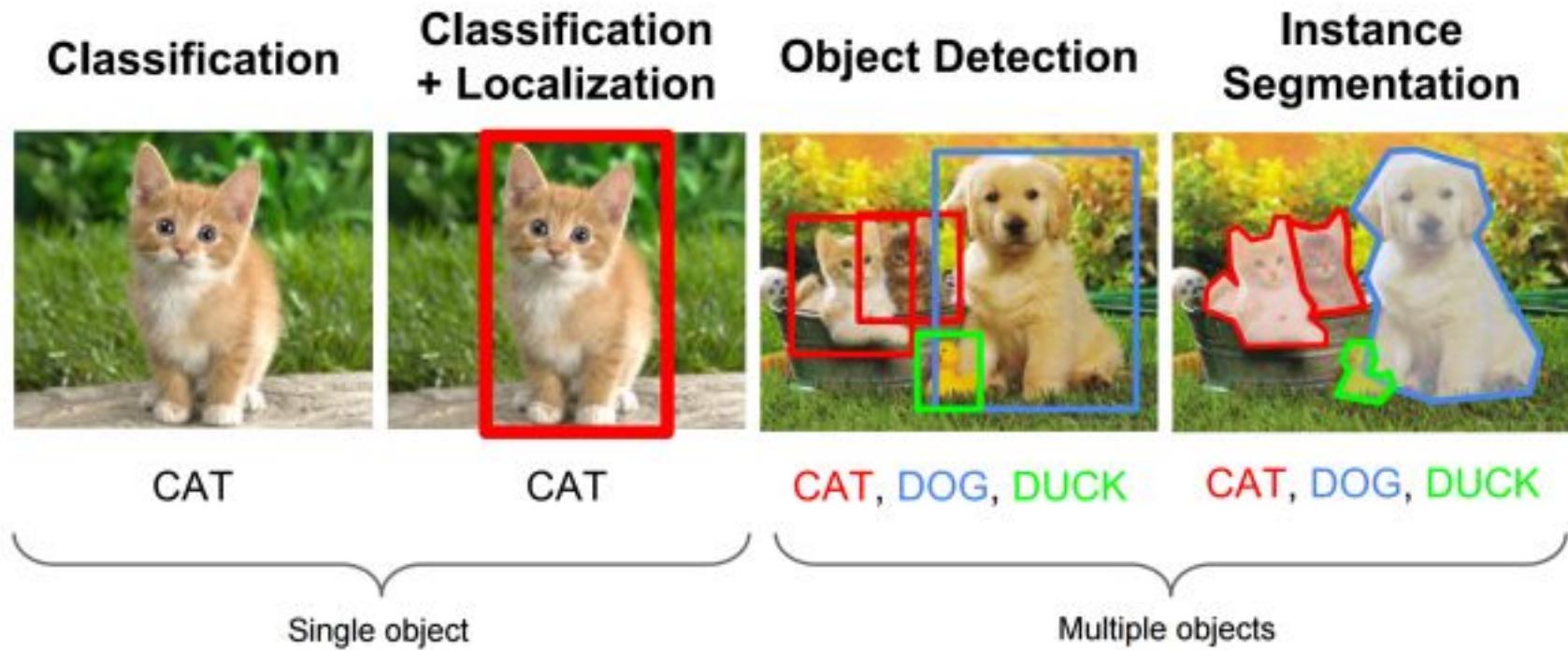
- Weakly-Supervised Learning
- Focus of this Tutorial

Visual Recognition Tasks

Differences between localization, detection and segmentation

- OL: aims to locate the main (or most visible) object in an image
- OD/IS: aims to find all the instances of objects and their boundaries

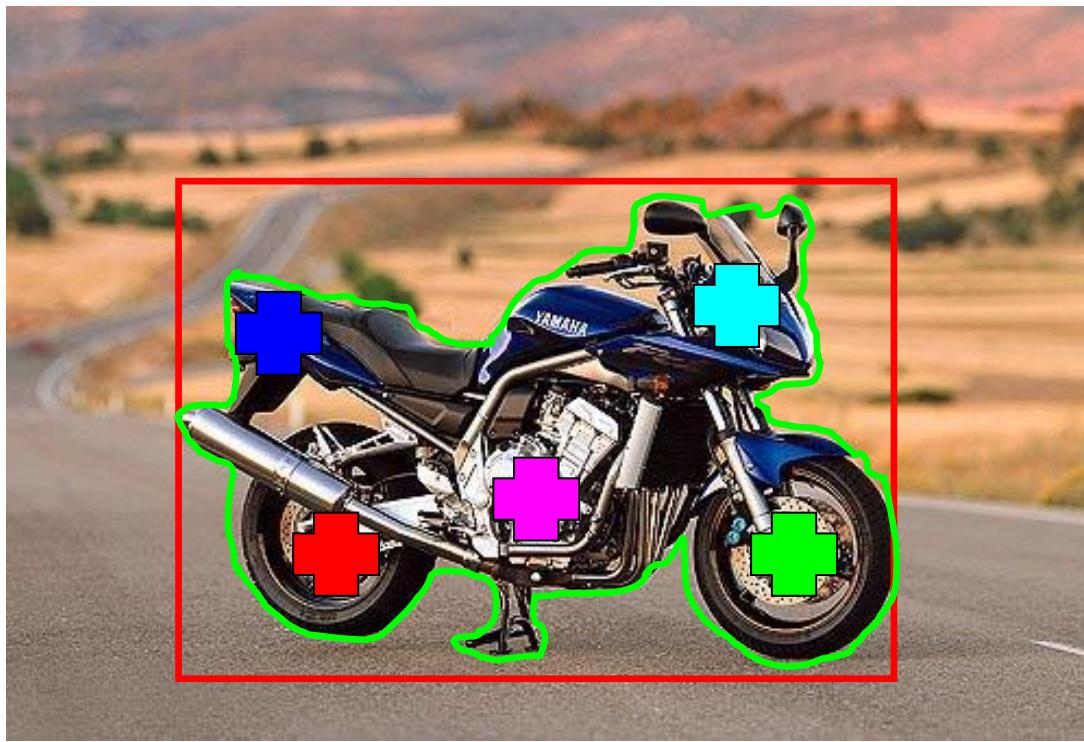
Semantic segmentation associates every pixel of an image with a class label



Visual Recognition Tasks

- **Supervised learning:** images in the training set are annotated with the “correct answer” that the model is expected to produce

Contains a motorbike



Levels of annotations:

- Image labels: for a classification task
- Patches or bounding boxes: for a detection or localization tasks
- Points: for a point-wise localization task
- Pixel regions: for a segmentation task

Visual Recognition Tasks

- Why learn from data with *weak annotations*?

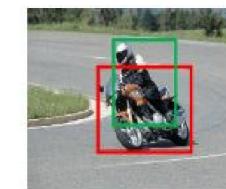
To optimize large DL models end-to-end, and this requires much data to determine all the model parameters

Collection and annotation of labeled training data generally costly

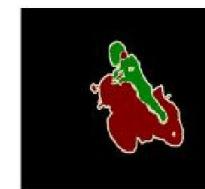
- full supervision may involve outlining objects, patches or marking points, pixels with category labels, etc.
- labels are ambiguous in some applications



{motorbike, person}
motorbike (point),
person (point)}



{motorbike (b-box),
person (b-box)}



{motorbike (pixel labels),
person (pixel labels)}

1 sec
per class

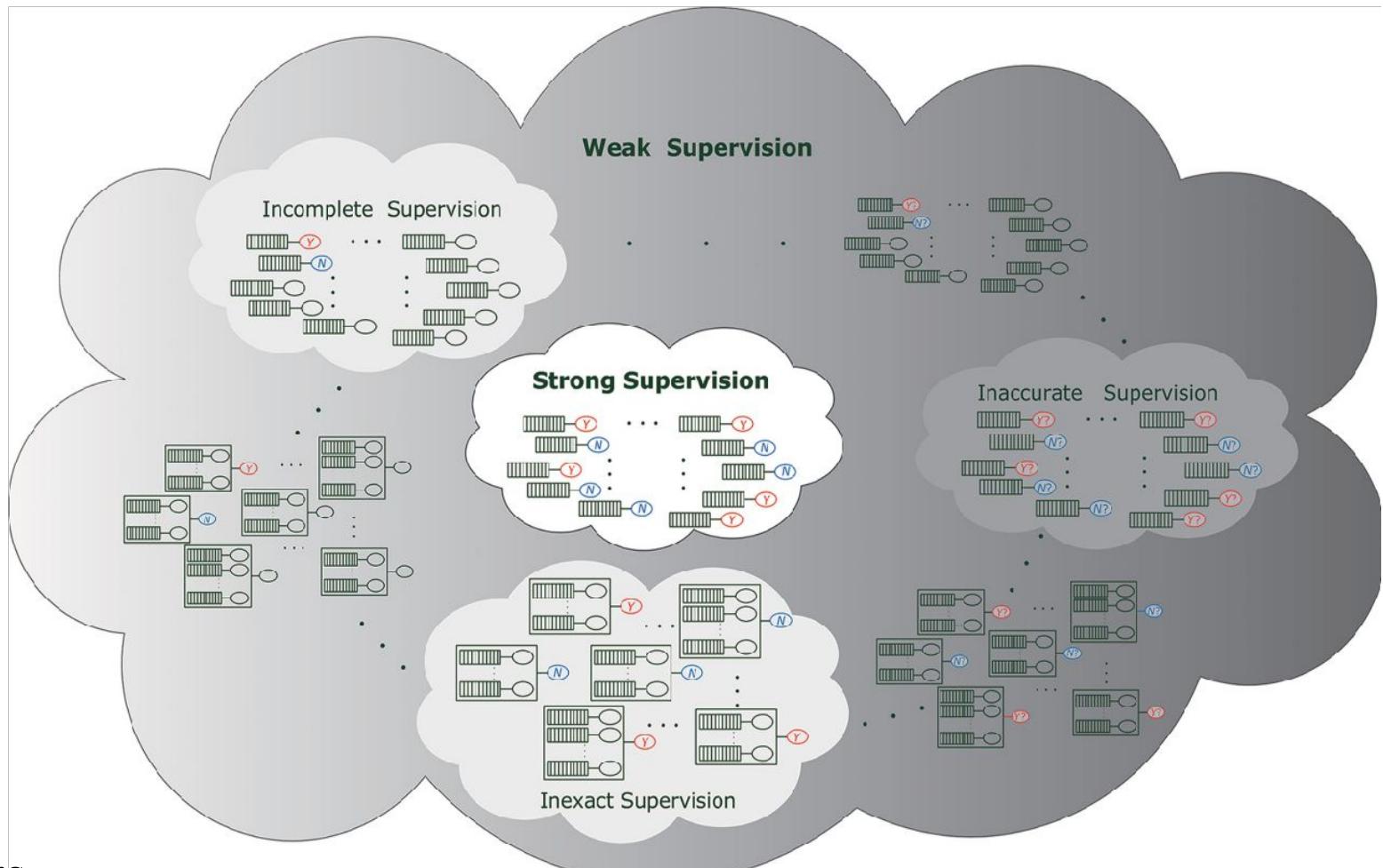
2.4 sec
per instance

10 sec
per instance

78 sec
per instance



Weak Supervised Learning Scenarios



- bars = vectors
- red/blue ovals = labels
- "?" = inaccurate labels

Source: Z. Zhou. 'A brief introduction to weakly supervised learning.' *National Science Review*, 5(1):44–53, 2018.

Weak Supervised Learning Scenarios

1) **Incomplete supervision**: when only a small subset of training data has labels, although unlabelled data is abundant

WSL techniques:

- AL (active learning)
- SSL (semi-supervised learning)

2) **Inexact supervision**: when training on labelled data with coarse labels

WSL technique:

- MIL (multiple instance learning)

3) **Inaccurate supervision**: when labels may suffer from errors or noise

WSL techniques:

- data-editing methods
- crowdsourcing with majority vote

Weak Supervised Learning Techniques

1) Incomplete supervision:

- **AL (active learning):** query an expert to label most relevant samples
- **SSL (semi-supervised learning):** involves training a model using both fully labeled and unlabeled examples

2) Inexact supervision:

- **MIL (multiple-instance learning):** uses training examples grouped into sets (*bags*). Supervision is provided only for an entire set of instances.

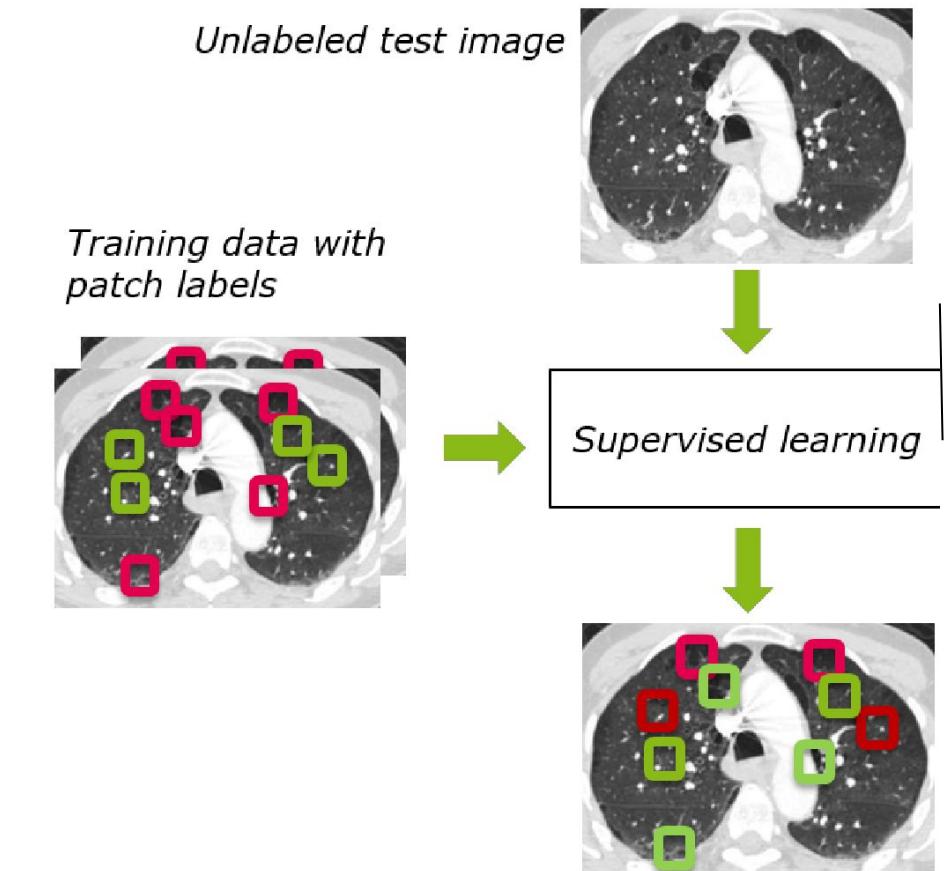
3) Inaccurate supervision:

- **Data-editing methods:** determine outlier annotations
- **Crowdsourcing with majority vote:** synthesis of responses from a large population of annotators

Weak Supervised Learning Techniques

Example: classifying healthy (green) vs emphysema (red) patches of tissue in chest CT images

- 1) *Supervised learning:* labeled healthy and abnormal patches available

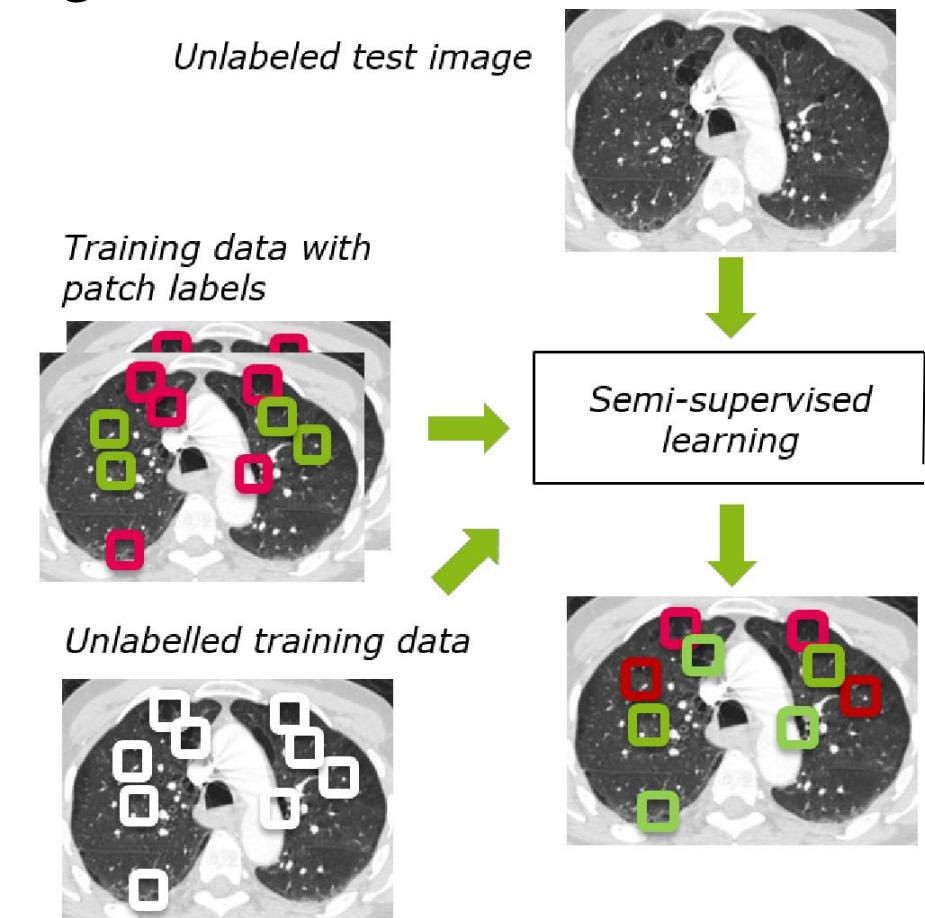


Source: V. Cheplygina et al., ‘Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis,’ *Medical Image Analysis*, 2019.

Weak Supervised Learning Techniques

Example: classifying healthy (green) vs emphysema (red) patches of tissue in chest CT images

2) *Semi-supervised learning*: in addition to healthy and abnormal patches, unlabeled patches are available.

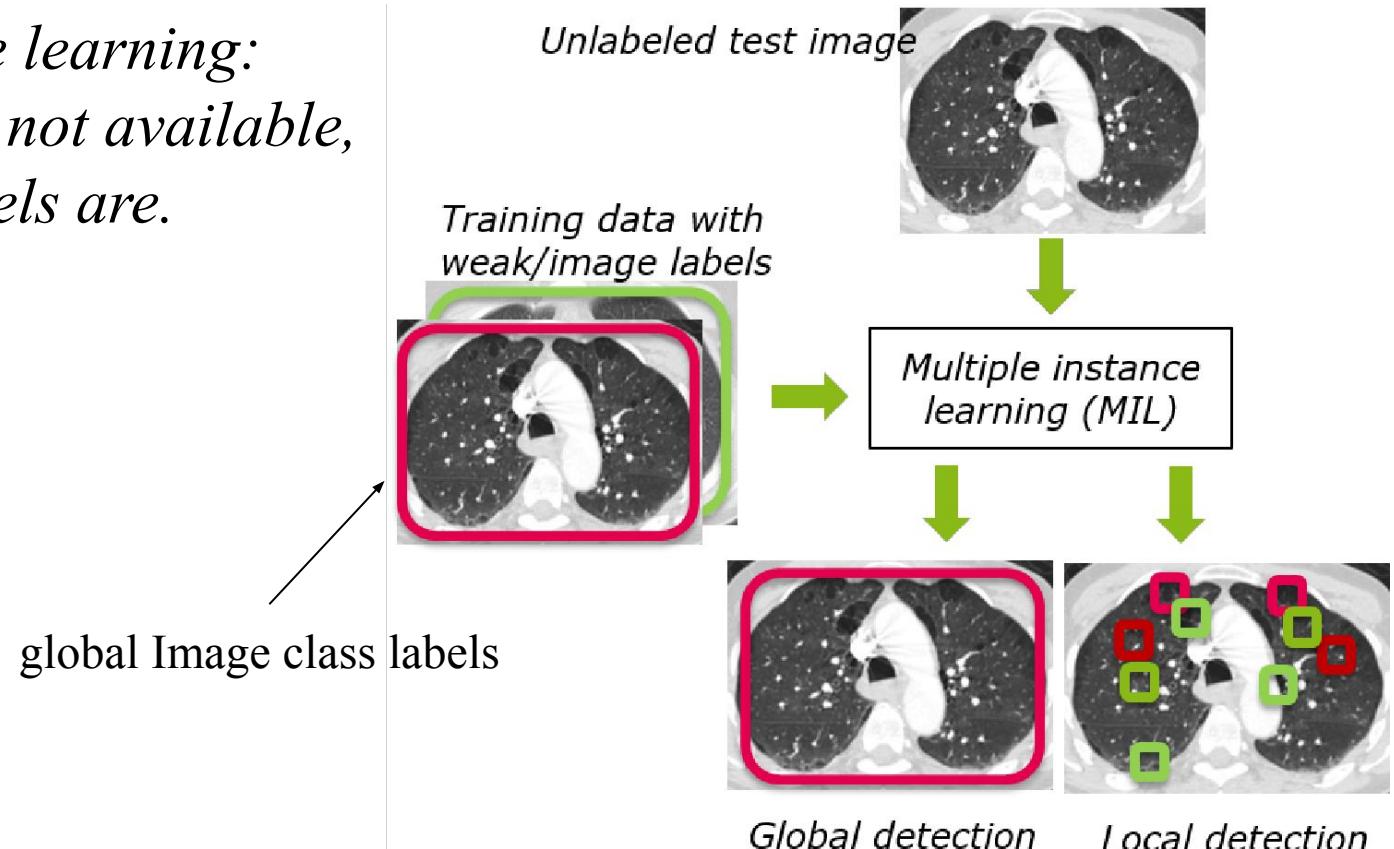


Source: V. Cheplygina et al., ‘Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis,’ *Medical Image Analysis*, 2019.

Weak Supervised Learning Techniques

Example: classifying healthy (green) vs emphysema (red) patches of tissue in chest CT images

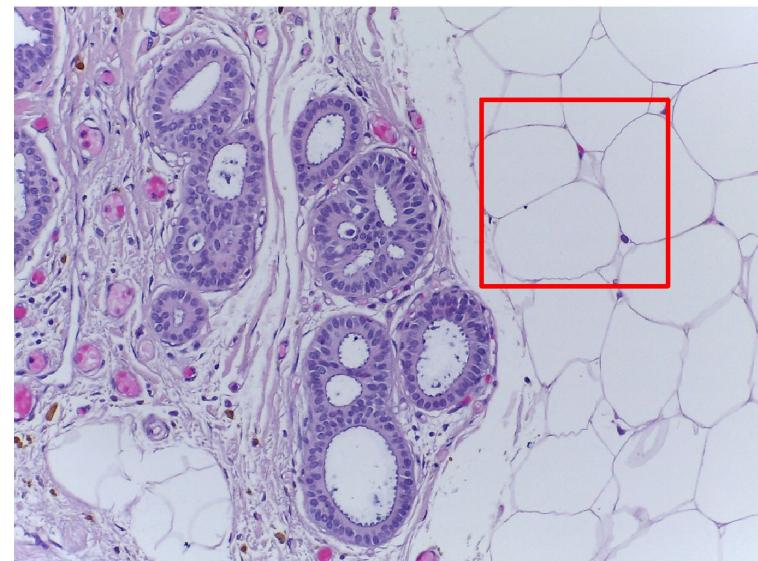
3) *Multiple instance learning:*
labeled patches are not available,
but Image-level labels are.



Source: V. Cheplygina et al., ‘Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis,’ *Medical Image Analysis*, 2019.

Computer Aided Diagnosis

- **Example:** MIL in computer aided diagnosis
given a large histological image, predict if a subject is diseased or healthy and locate regions of interest
- **Image (bag)** = set of segments or patches (instances) with global annotation
- **Database** - example of an image from the ICIAR BACH Challenge 2018

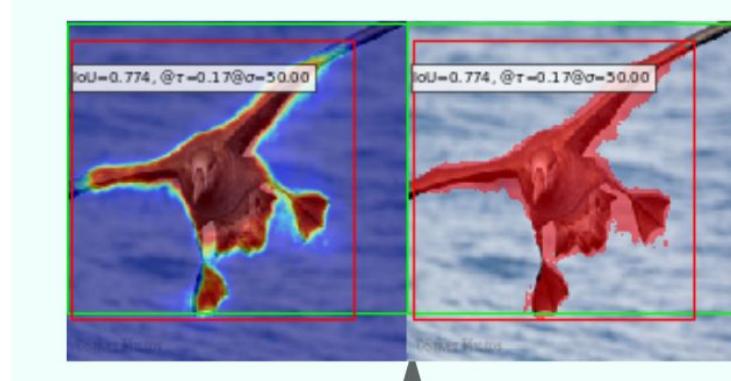


Focus of this Tutorial

Section 2:

Weakly-Supervised Object Localization (WSOL):

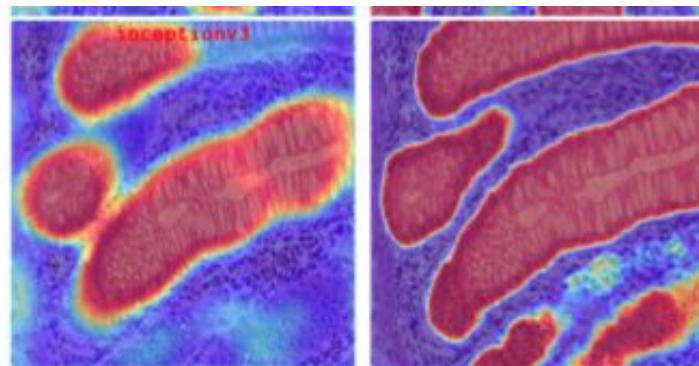
Localizing objects based on training data with global image-class labels



Section 3:

Weakly-Supervised Semantic Segmentation (WSSS)

segmenting objects based on training data with global image-class labels



Part 2

Review of WSOL Methods

- **WSOL literature: bottom-up and top-down methods**
- **Case studies:**
 - a) F-CAM for improved interpolation
 - b) Transformer-based models

Part 2. Review of WSOL methods: Setup

Supervised object localization



- Regression task
- One object
- Supervision: bounding box

Part 2. Review of WSOL methods: Setup

Supervised object localization



- Regression task
- One object
- Supervision: bounding box [high cost, prevents scaling up]

Is there other type of CHEAP SUPERVISION ?

Part 2. Review of WSOL methods: Setup

Supervised object localization



- Regression task
- One object
- Supervision: bounding box [high cost, prevents scaling up]

Is there other type of CHEAP SUPERVISION ?

Yes, global image class !!! – Weak supervision

Part 2. Review of WSOL methods: Setup

Weakly Supervised object localization: WSOL



Input

WSOL model

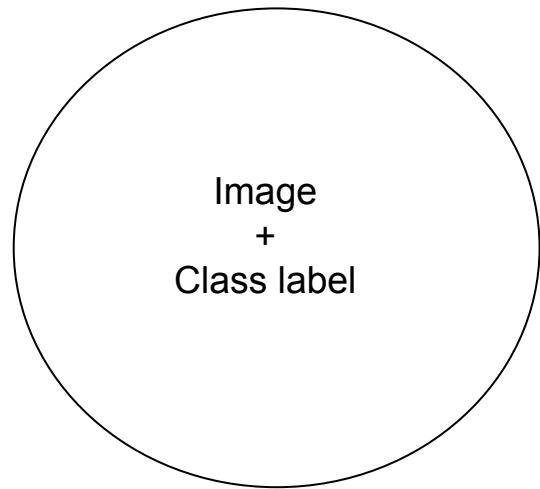


Bounding box + image class

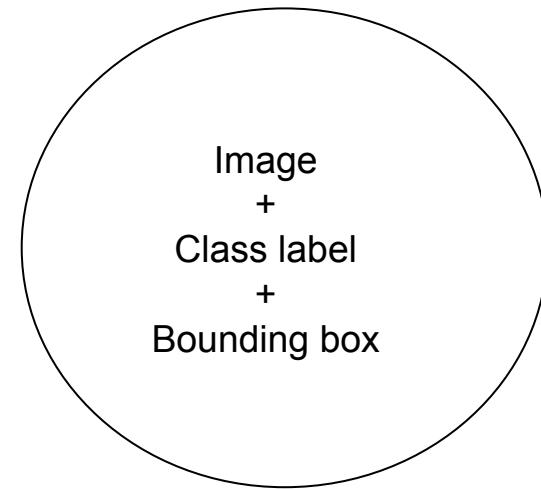
- One object
- Supervision: Image class
- Output: Bounding box + image class

Part 2. Review of WSOL methods: Setup

Weakly Supervised object localization: WSOL



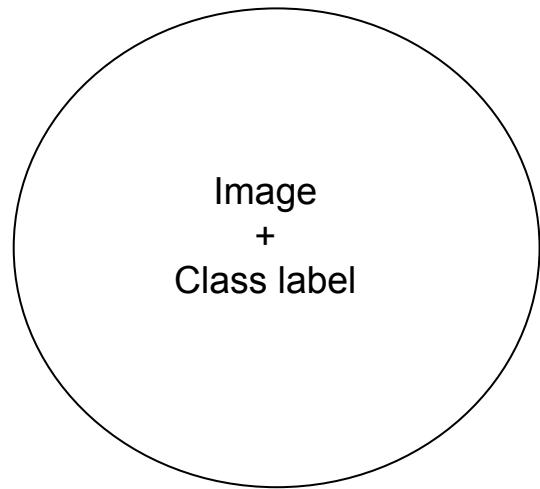
Train set



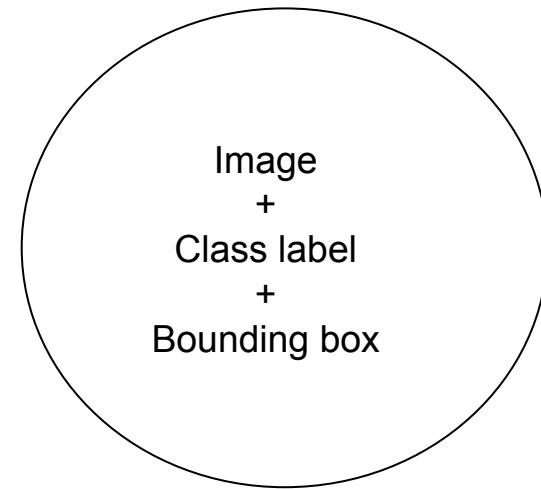
Test set

Part 2. Review of WSOL methods: Setup

Weakly Supervised object localization: WSOL



Train set



Test set

Evaluate: classification, localization,
classification + localization

Datasets:

- CUB: birds species (200 classes)
- Imagenet-1k: common objects (1k classes)

Part 2. Review of WSOL methods: Setup

Weakly Supervised object localization: WSOL

**How a bounding box is produced in
WSOL?**

Part 2. Review of WSOL methods: Setup

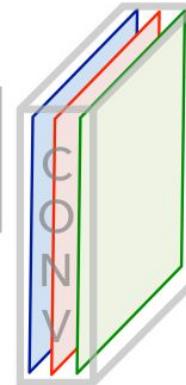
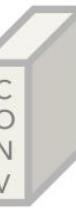
How a bounding box is produced in WSOL?

- Deep learning methods
- Class Activation Maps (CAMs)

Part 2. Review of WSOL methods: Setup

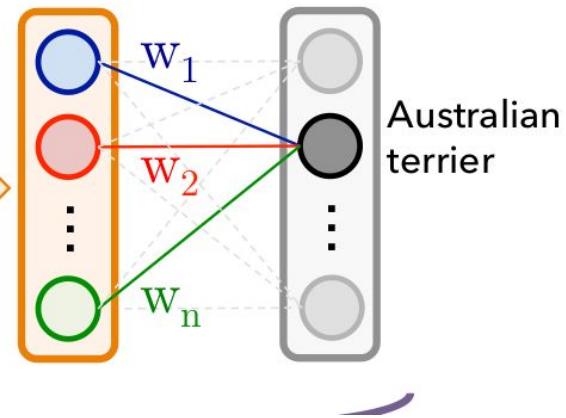
How a bounding box is produced in WSOL?

CAMs



CAMs

Per-class spatial map

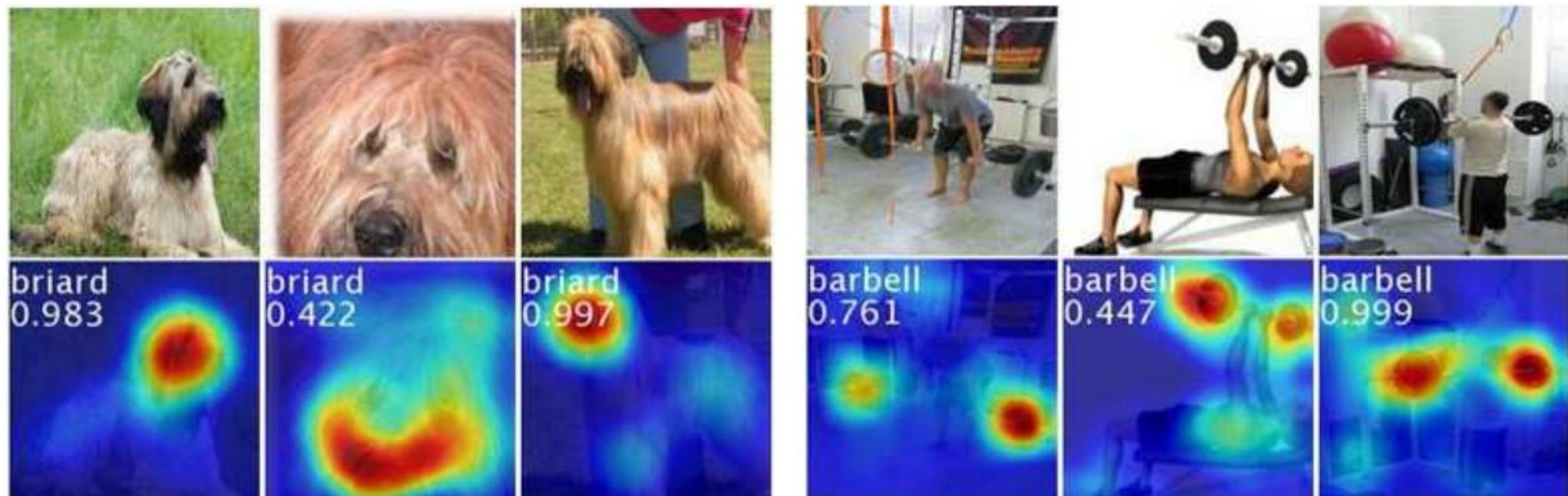


Class
Activation
Map
(Australian terrier)

Part 2. Review of WSOL methods: Setup

How a bounding box is produced in WSOL?

CAMs



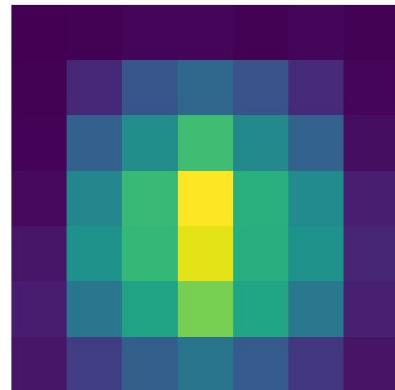
Part 2. Review of WSOL methods: Setup

How a bounding box is produced in WSOL?

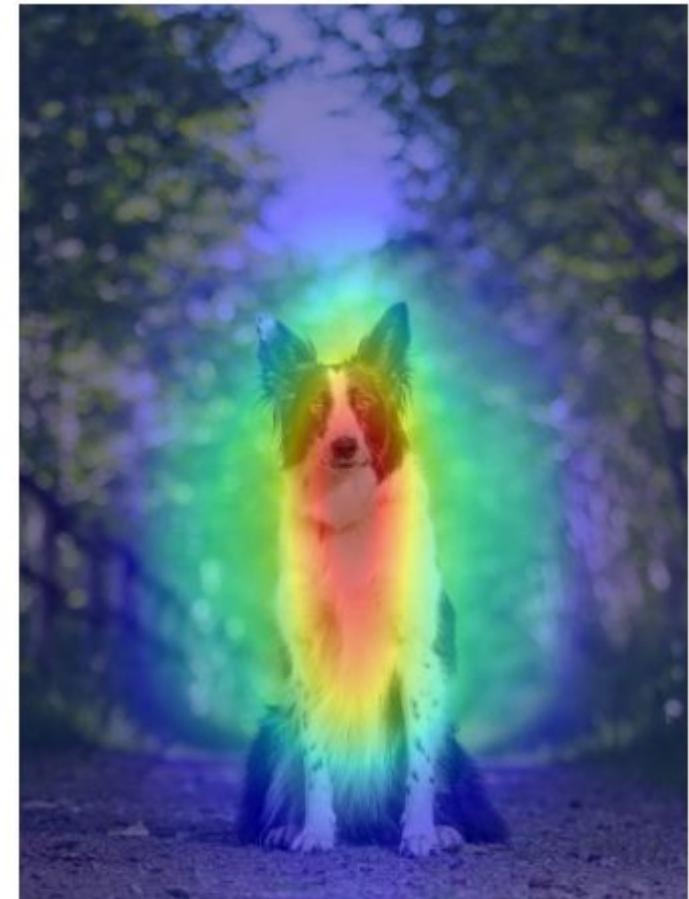
CAMs



Input



Original CAM



Interpolated CAM

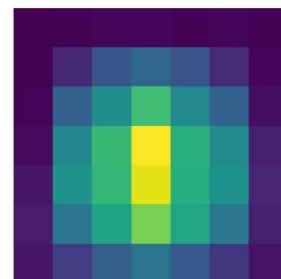
Part 2. Review of WSOL methods: Setup

How a bounding box is produced in WSOL?

CAMs



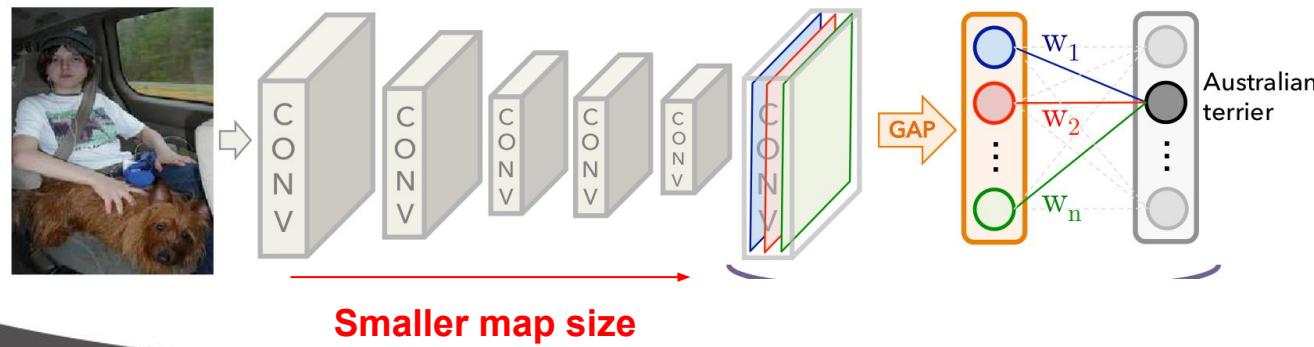
Input



Original
CAM



Interpolated CAM

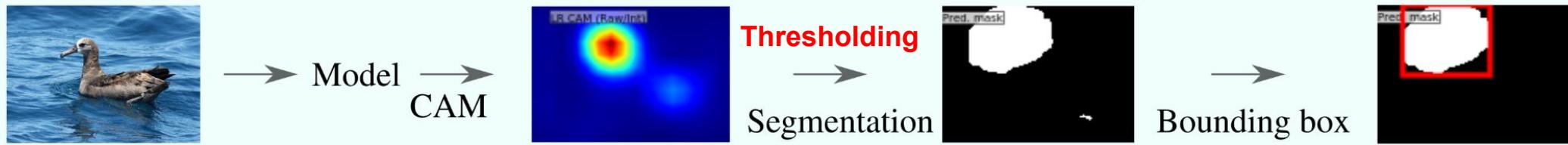


Part 2. Review of WSOL methods: Setup

How a bounding box is produced in WSOL?

CAMs

Image
processing,
contours,
tightest bbx

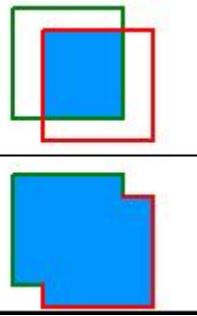


Part 2. Review of WSOL methods: Setup

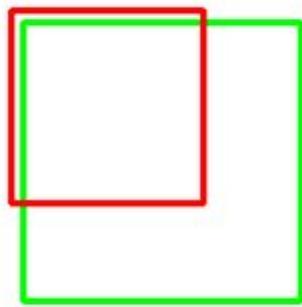
WSOL: Evaluation/metrics

Intersection Over Union

$$IOU = \frac{\text{area of overlap}}{\text{area of union}}$$

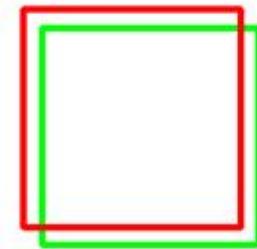


IoU: 0.4034



Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

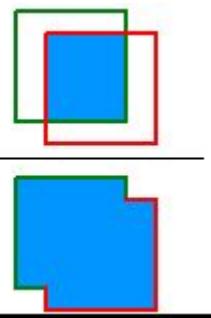
Part 2. Review of WSOL methods: Setup

WSOL: Evaluation/metrics

Standard **localization** metric:

- **CorLoc**: Correct localization
=
- **GT-known** localization accuracy

1 if $\text{IOU} \geq \delta = 0.5$ else 0

$$IOU = \frac{\text{area of overlap}}{\text{area of union}} = \frac{\text{blue area}}{\text{red + green + blue areas}}$$


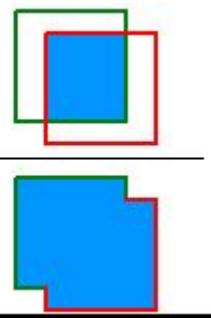
Part 2. Review of WSOL methods: Setup

WSOL: Evaluation/metrics

Standard **localization** metric:

- **CorLoc**: Correct localization
=
- **GT-known** localization accuracy

1 if $\text{IOU} \geq \delta = 0.5$ else 0

$$\text{IOU} = \frac{\text{area of overlap}}{\text{area of union}} = \frac{\text{blue area}}{\text{red + green + blue areas}}$$


In CVPR 2020, 2 new **localization** metrics:

- MaxBoxAcc
- MaxBoxAccV2

Take in consideration CAM thresholding

(localization via CAMs is threshold-dependent)

Evaluating Weakly Supervised Object Localization Methods Right. Choe et al. CVPR 2020.

Evaluating Weakly Supervised Object Localization Methods Right. Choe et al. CVPR 2020.

Part 2. Review of WSOL methods: Setup

WSOL: Evaluation/metrics

WSOL localization metric: **MaxBoxAcc**

$$\text{BoxAcc}(\tau, \delta) = \frac{1}{N} \sum_n 1_{\text{IoU}\left(\underbrace{\text{box}(s(\mathbf{x}^{(n)}), \tau), B^{(n)})}_{\text{Tightest bbox around largest connected component}}\right) \geq \delta}$$

Normalized CAM
G-truth bbox
IoU threshold
CAM threshold

Part 2. Review of WSOL methods: Setup

WSOL: Evaluation/metrics

WSOL localization metric: **MaxBoxAcc**

$$\text{BoxAcc}(\tau, \delta) = \frac{1}{N} \sum_n 1_{\text{IoU}\left(\text{box}(s(\mathbf{X}^{(n)}), \tau), B^{(n)}\right) \geq \delta}$$

$$\text{MaxBoxAcc}(\delta) := \max_{\tau} \text{BoxAcc}(\tau, \delta)$$

**τ in [0, 1], step: 0.001.
δ = 0.5**



Part 2. Review of WSOL methods: Setup

WSOL: Evaluation/metrics

WSOL localization metric: **MaxBoxAccV2**

Account for variable object sizes

$$1/3 \sum_{\delta} MaxBoxAcc(\delta), \quad \delta \in \{0.3, 0.5, 0.7\}$$

MaxBoxAccV2 more difficult than MaxBoxAcc

Part 2. Review of WSOL methods: Setup

WSOL: Evaluation/metrics

WSOL localization and classification metric: **top1-localization, top5-localization**

Top1-localization

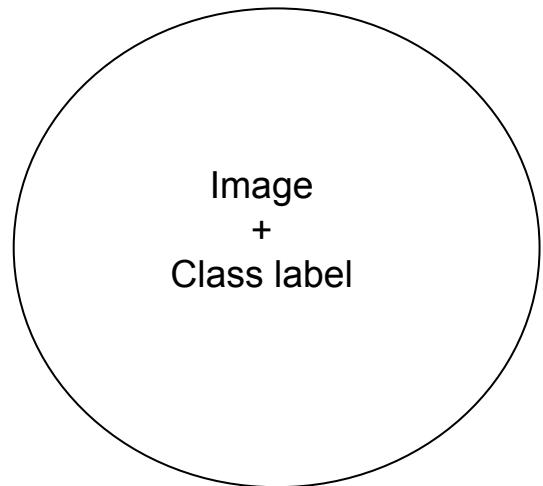
MaxBoxAcc(0.5) = 1 and true class = **top1** prediction

Top5-localization

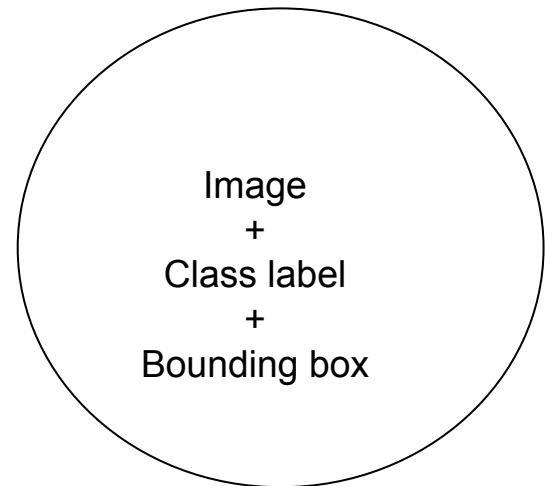
MaxBoxAcc(0.5) = 1 and true class in **top5** predictions

Part 2. Review of WSOL methods: Setup

WSOL: Model selection (validation / early stopping)



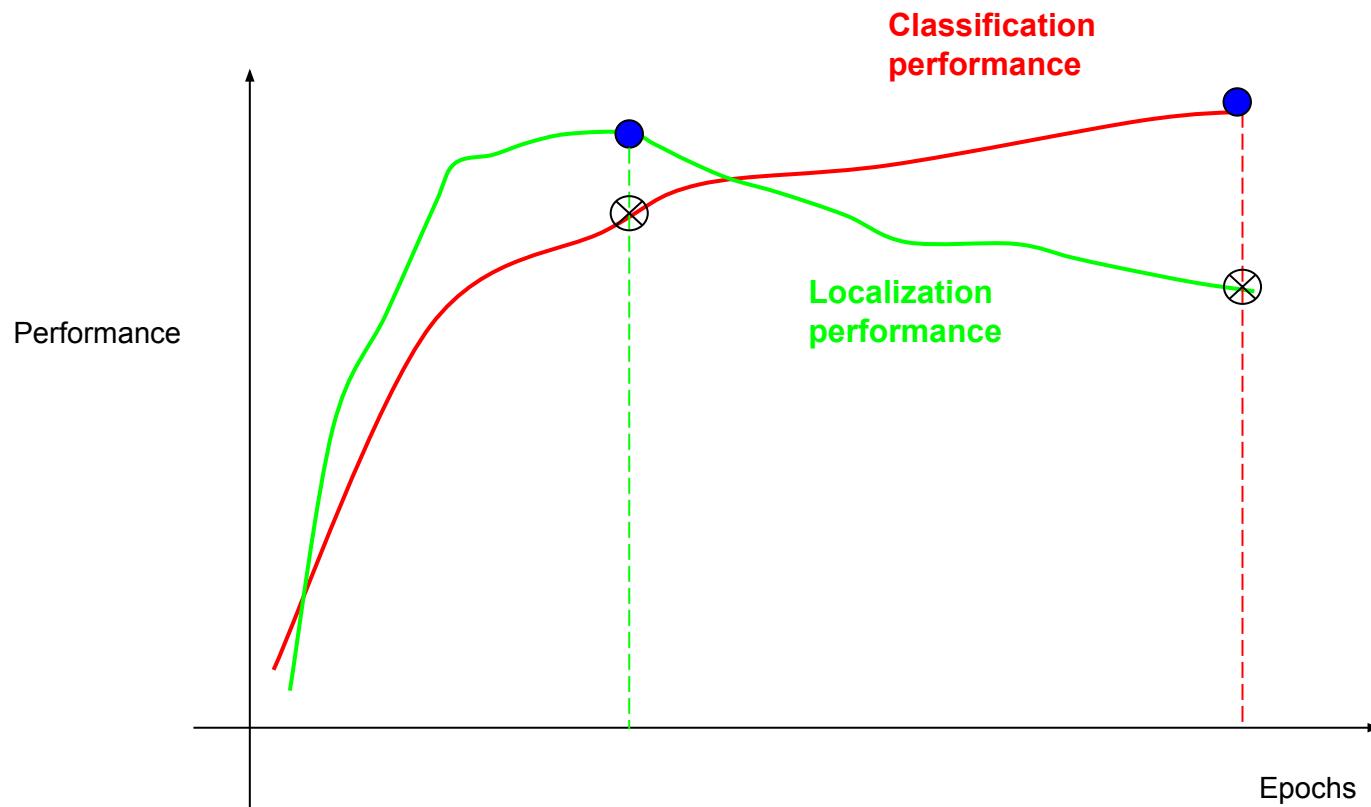
Train set



Test set

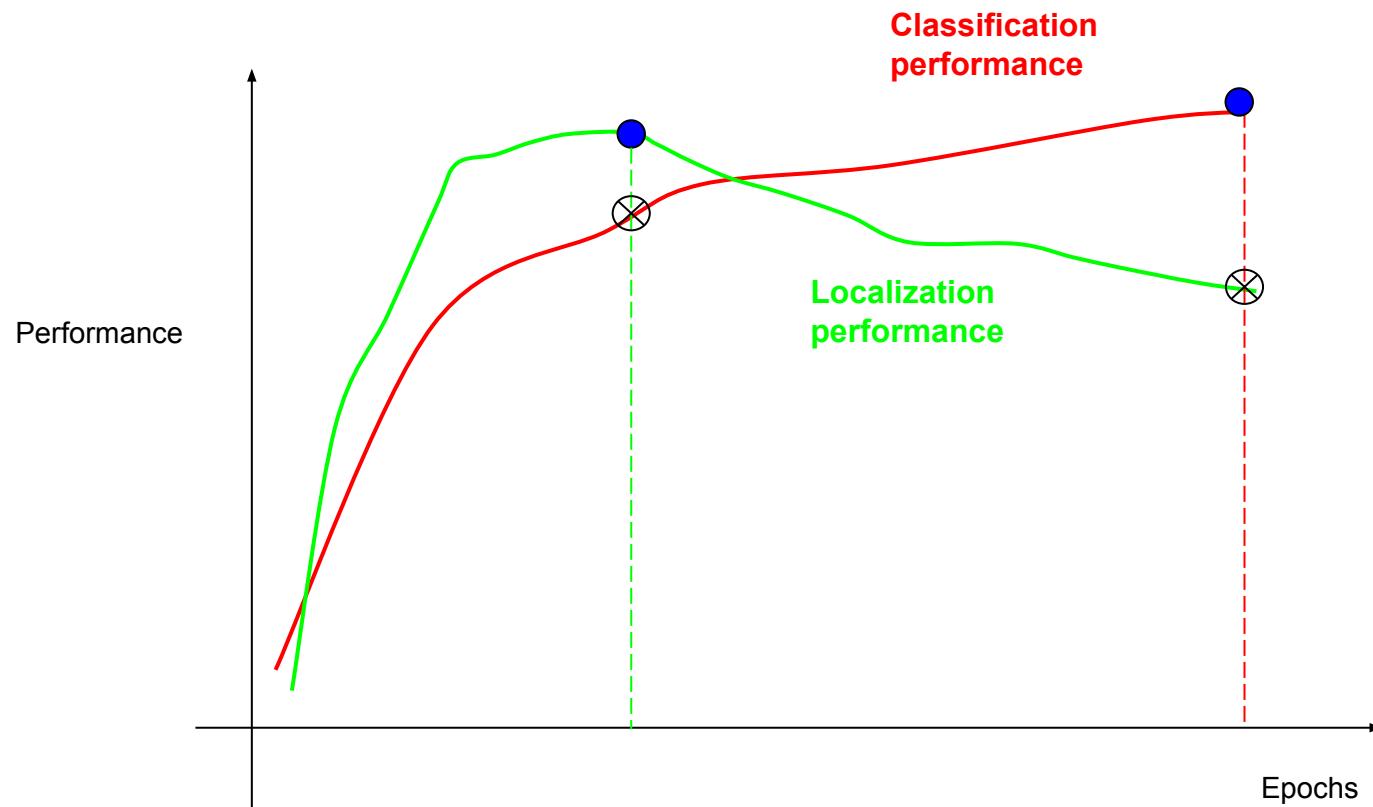
Part 2. Review of WSOL methods: Setup

WSOL: Model selection (validation / early stopping)



Part 2. Review of WSOL methods: Setup

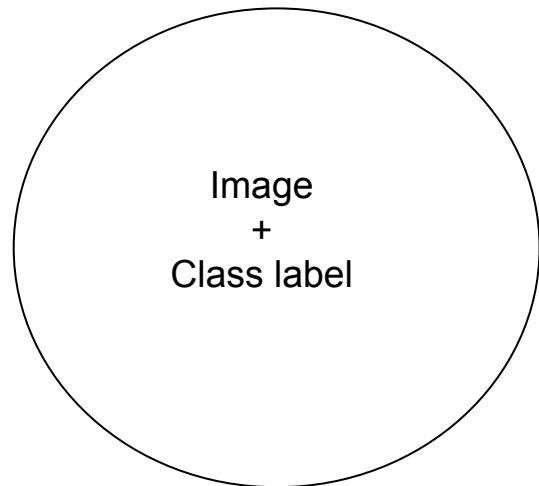
WSOL: Model selection (validation / early stopping)



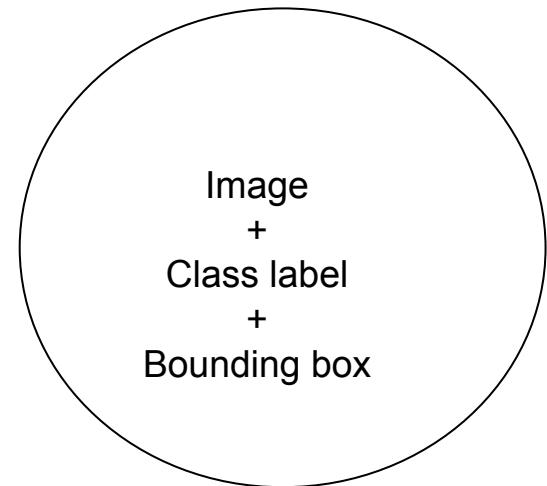
Classification, localization: antagonist tasks (WSOL)

Part 2. Review of WSOL methods: Setup

WSOL: Model selection (validation / early stopping)



Train set

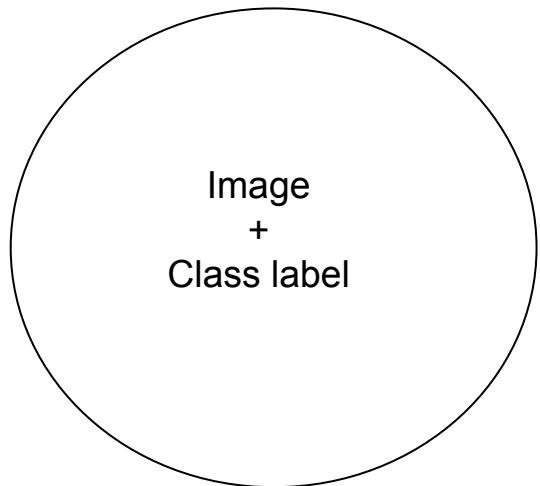


Test set

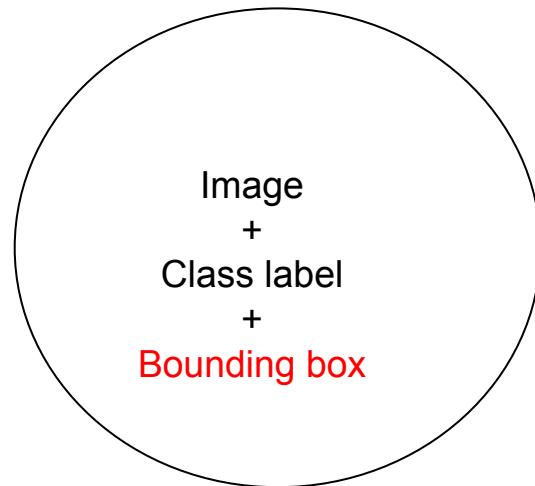
User bias over test set !!!

Part 2. Review of WSOL methods: Setup

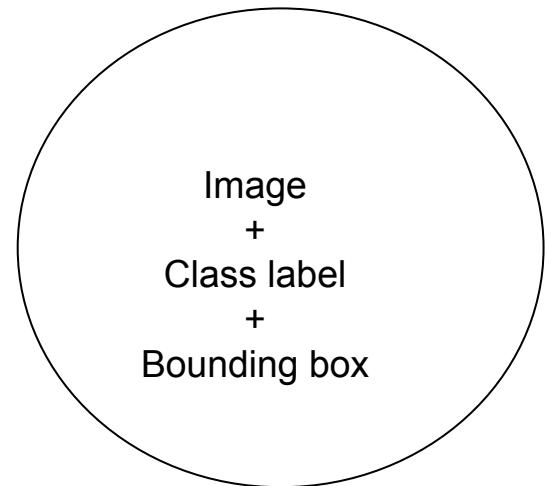
WSOL: Model selection (validation / early stopping)



Train set



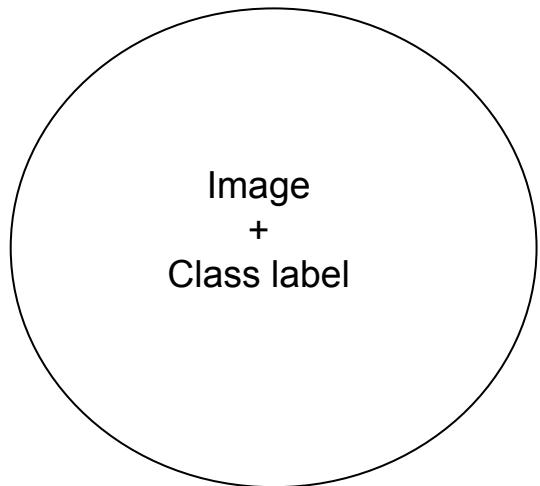
Valid set



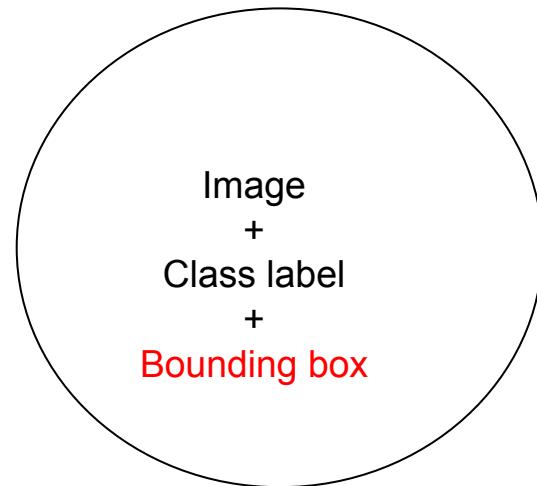
Test set

Part 2. Review of WSOL methods: Setup

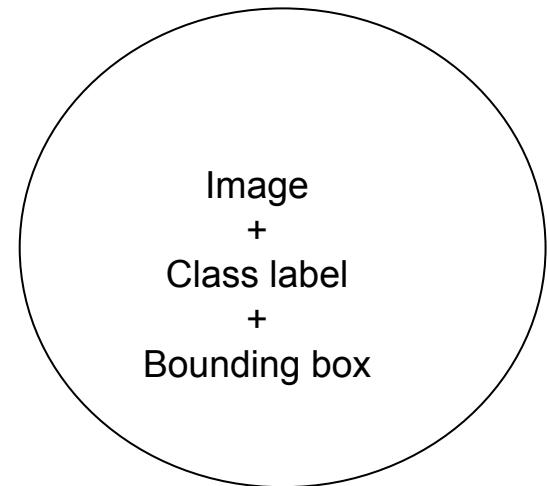
WSOL: Model selection (validation / early stopping)



Train set



Valid set

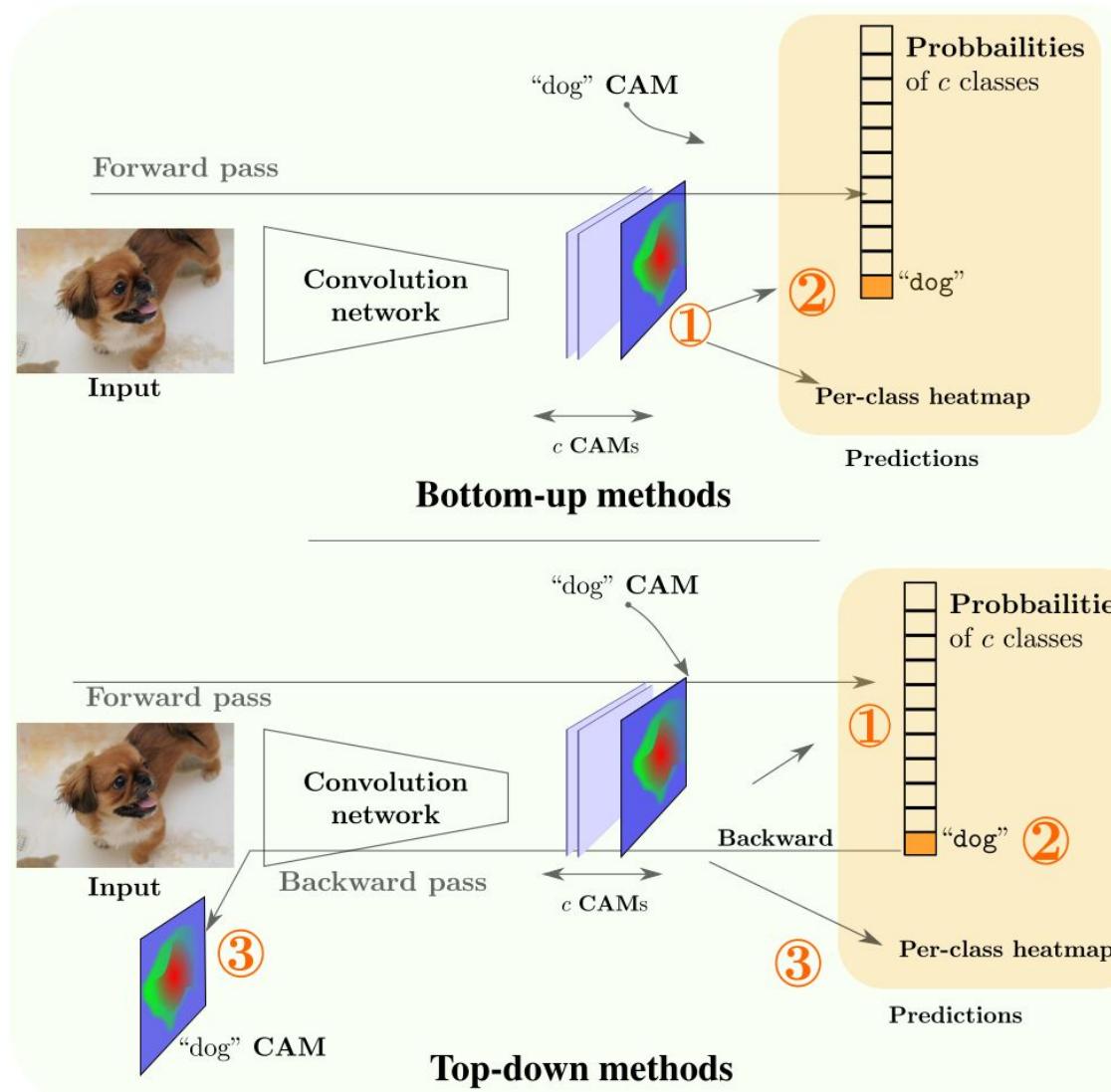


Test set

- Allow more fair comparison between methods
- Realistic?

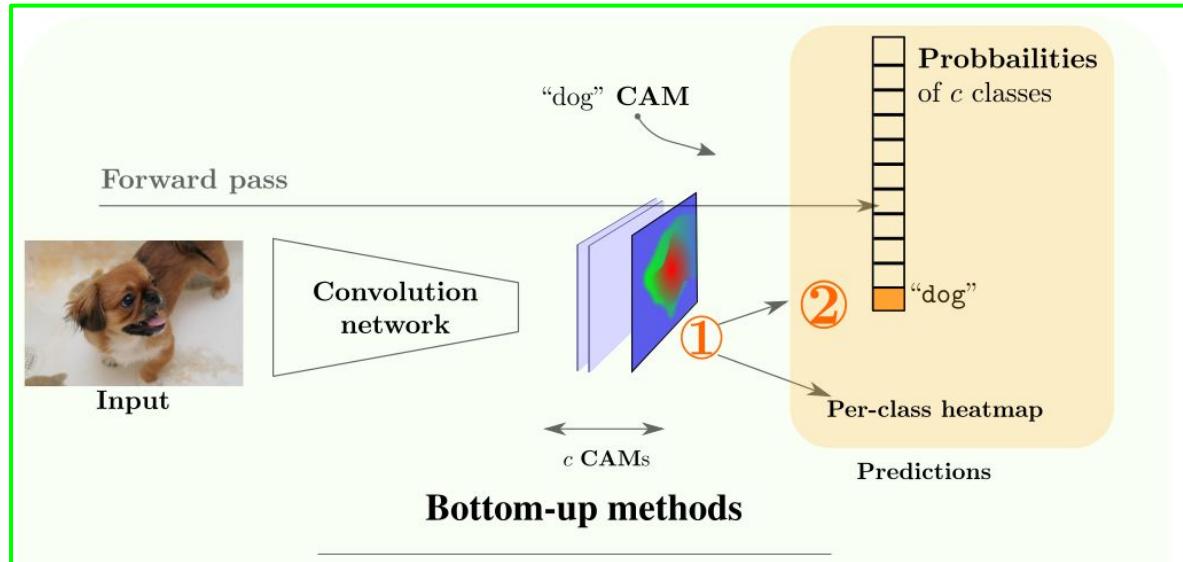
Part 2. Review of WSOL methods: Literature

Taxonomy

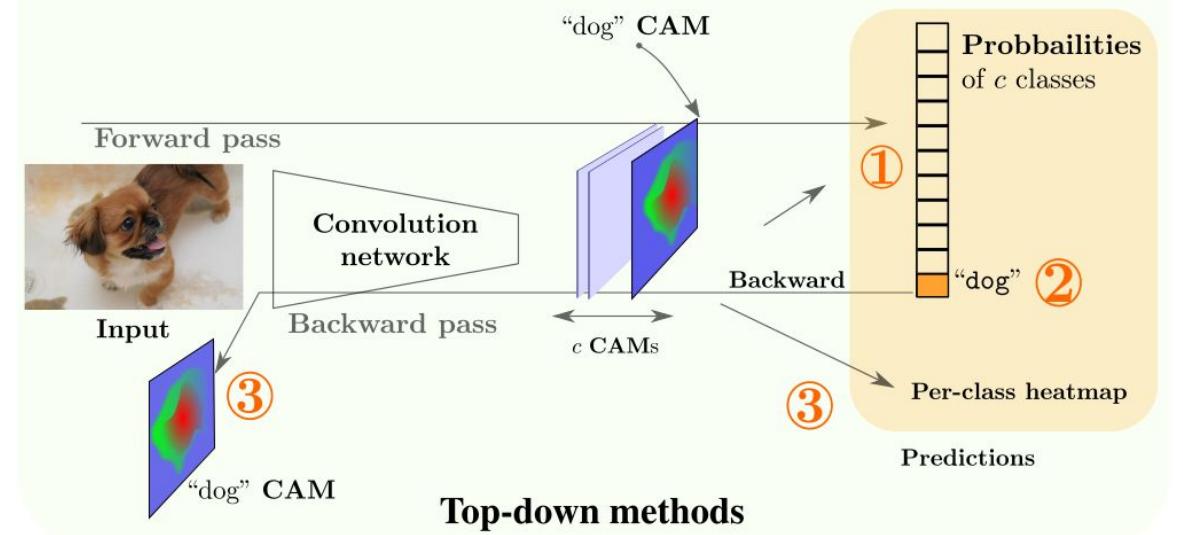


Part 2. Review of WSOL methods: Literature

Taxonomy

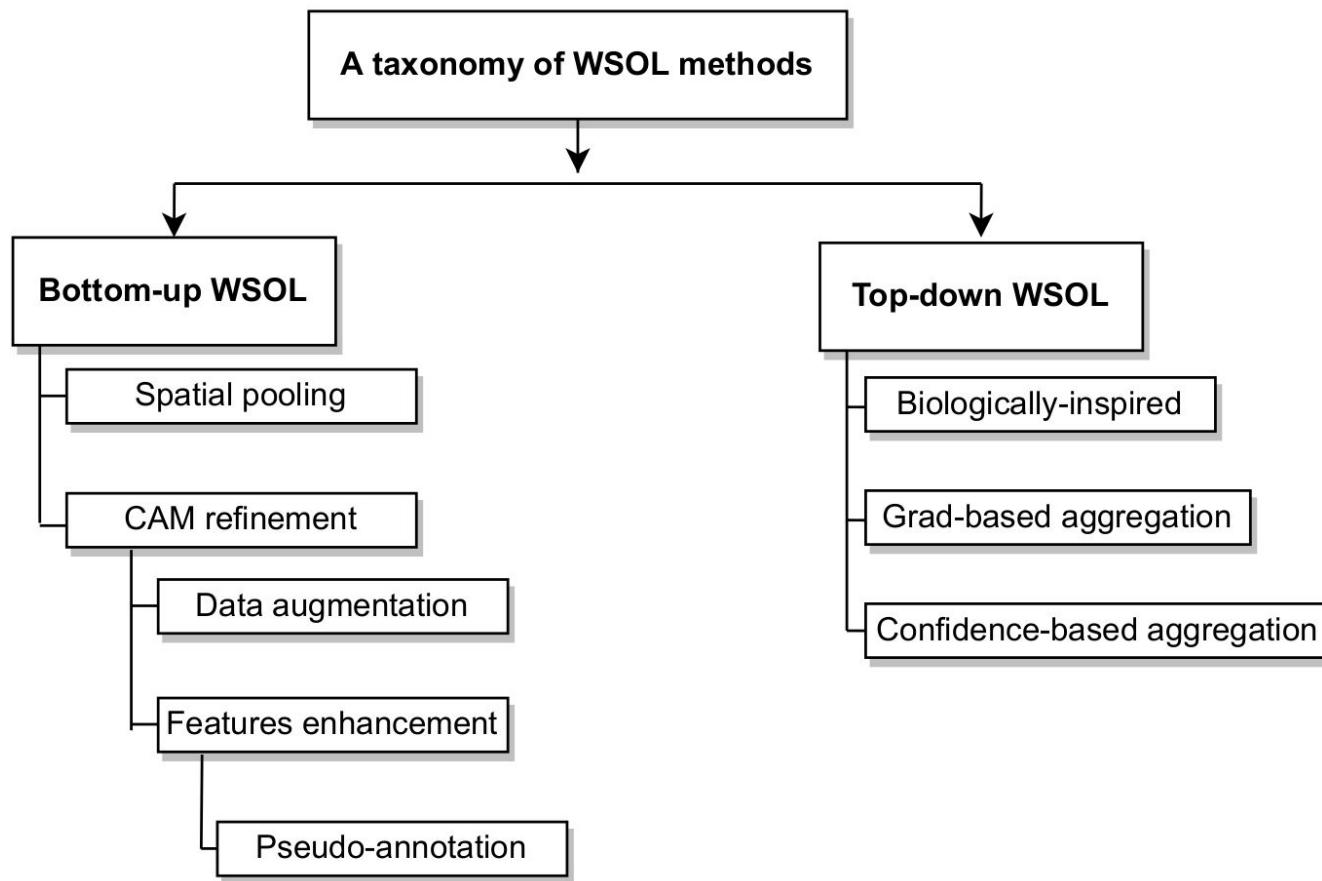


More common
in WSOL



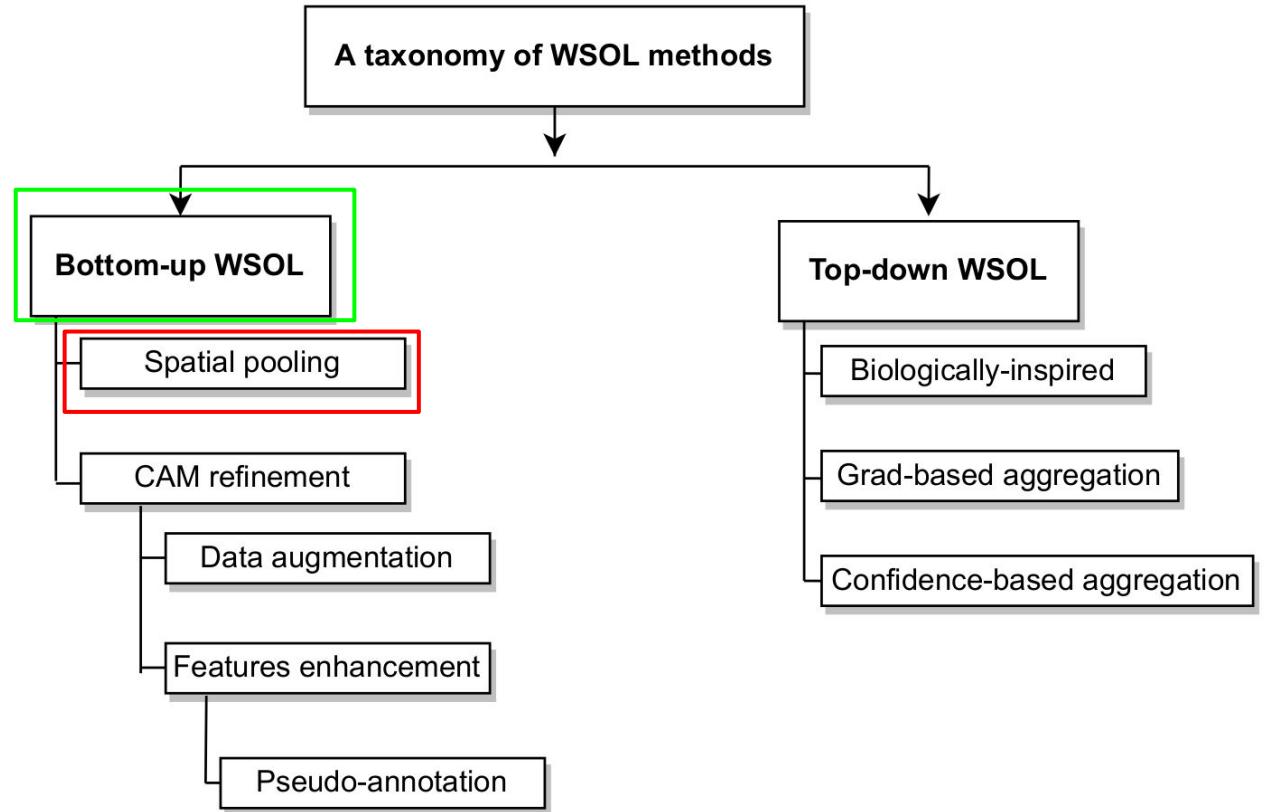
Part 2. Review of WSOL methods: Literature

Taxonomy



Part 2. Review of WSOL methods: Literature

Taxonomy



Part 2. Review of WSOL methods: Literature

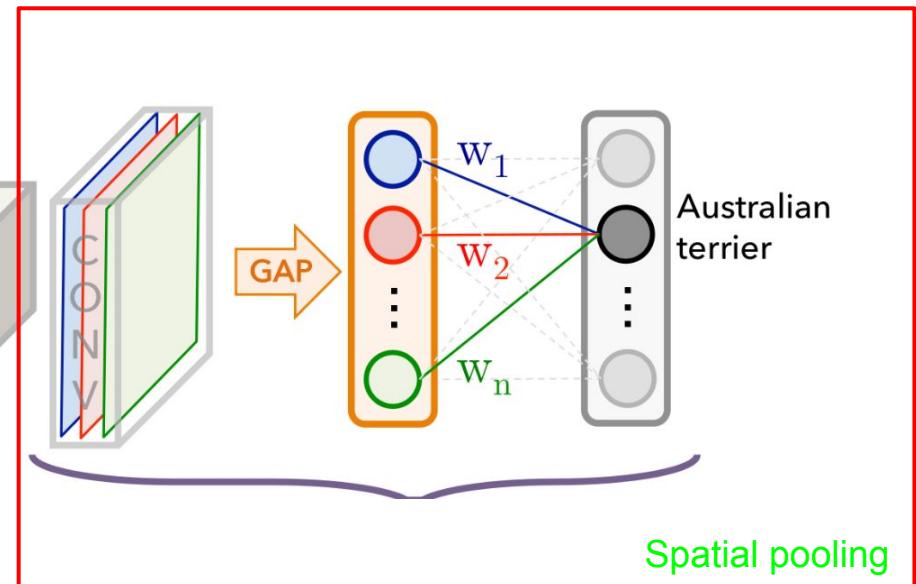
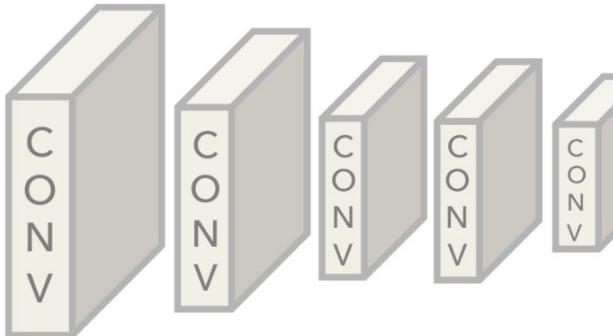
Bottom-up WSOL

Spatial pooling

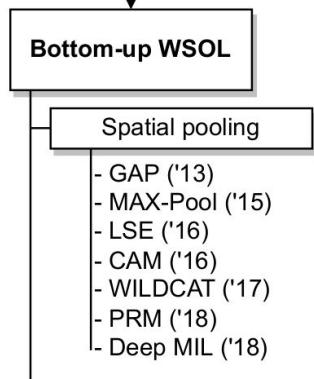
- GAP ('13)
- MAX-Pool ('15)
- LSE ('16)
- CAM ('16)
- WILDCAT ('17)
- PRM ('18)
- Deep MIL ('18)

Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps



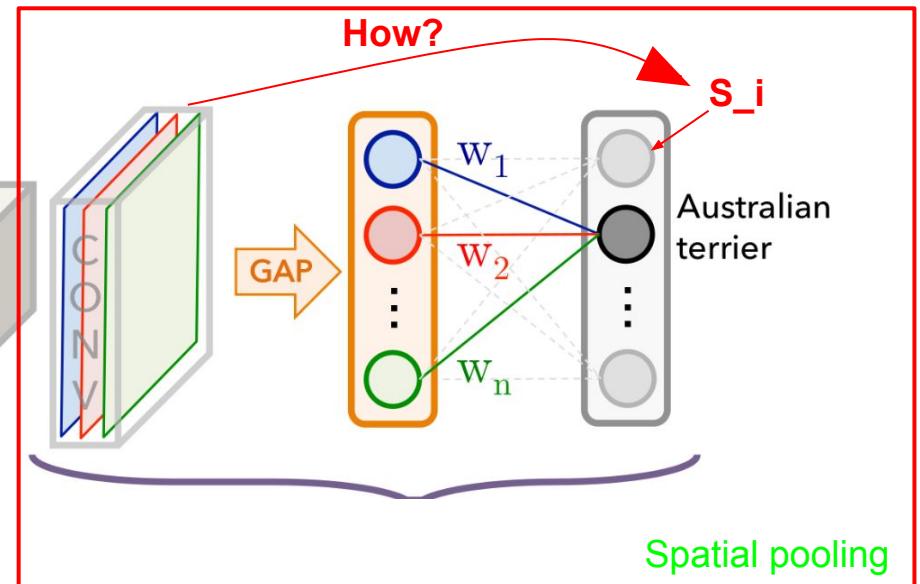
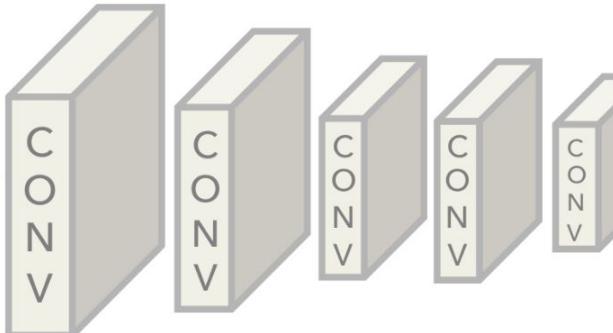
Part 2. Review of WSOL methods: Literature



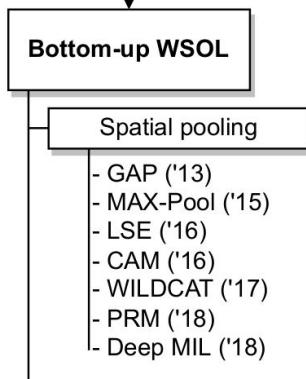
Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps

$$\Pr(y = i | \mathbf{x}) = \text{softmax}(\mathbf{s})_i = \frac{\exp(s_i)}{\sum_{j=1}^C \exp(s_j)} .$$



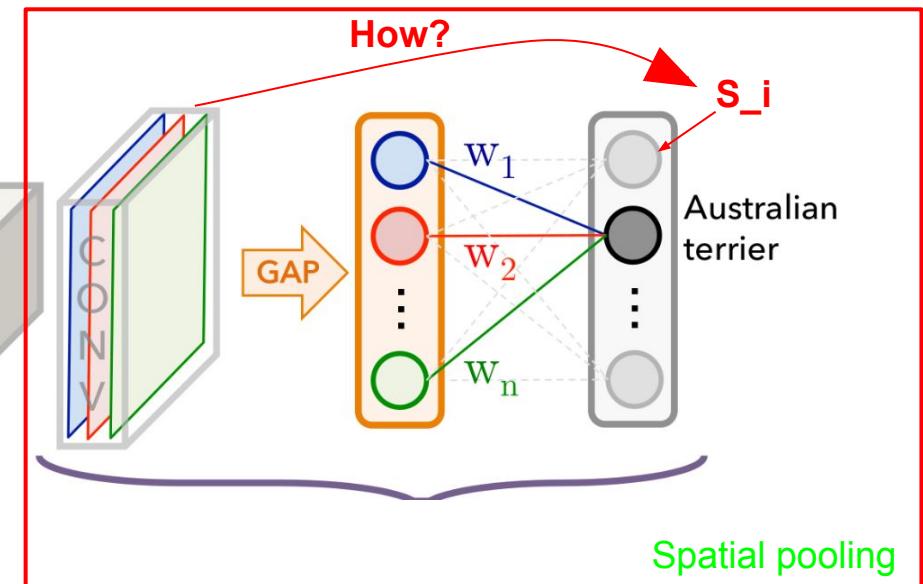
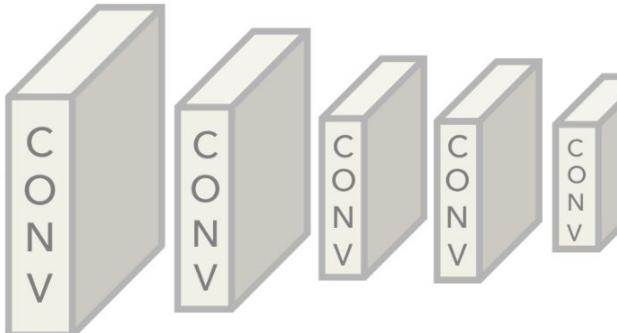
Part 2. Review of WSOL methods: Literature



Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps

$$\Pr(y = i | \mathbf{x}) = \text{softmax}(\mathbf{s})_i = \frac{\exp(s_i)}{\sum_{j=1}^C \exp(s_j)} .$$



Spatial pooling: Helps to excite ROI to emerge in CAMs

Part 2. Review of WSOL methods: Literature

Bottom-up WSOL

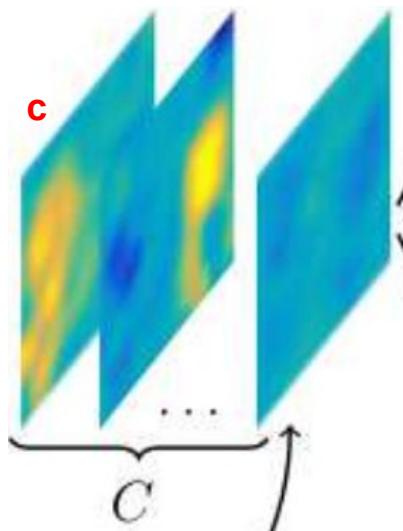
Spatial pooling

- GAP ('13)
- MAX-Pool ('15)
- LSE ('16)
- CAM ('16)
- WILDCAT ('17)
- PRM ('18)
- Deep MIL ('18)

Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps

M^c

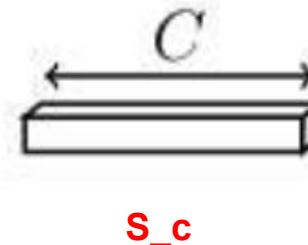


GAP: Global Average Pooling

All locations take part in scoring!!!!

$$s_c = \frac{1}{h \times w} \sum_{i,j} M_{i,j}^c$$

Classification



s_c

- ✓ cat
- ✓ dog
- ✗ bus

Part 2. Review of WSOL methods: Literature

Bottom-up WSOL

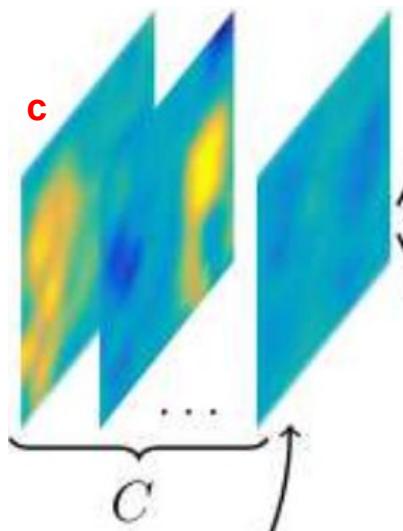
Spatial pooling

- GAP ('13)
- MAX-Pool ('15)
- LSE ('16)
- CAM ('16)
- WILDCAT ('17)
- PRM ('18)
- Deep MIL ('18)

Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps

M^c

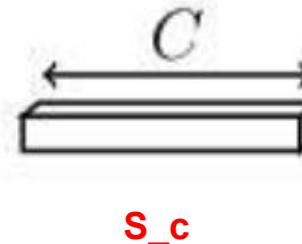


GAP: Global Average Pooling

All locations take part in scoring!!!!

$$s_c = \frac{1}{h \times w} \sum_{i,j} M_{i,j}^c$$

Classification



- ✓ cat
- ✓ dog
- ✗ bus

Only small part of the object activates!!!

Part 2. Review of WSOL methods: Literature

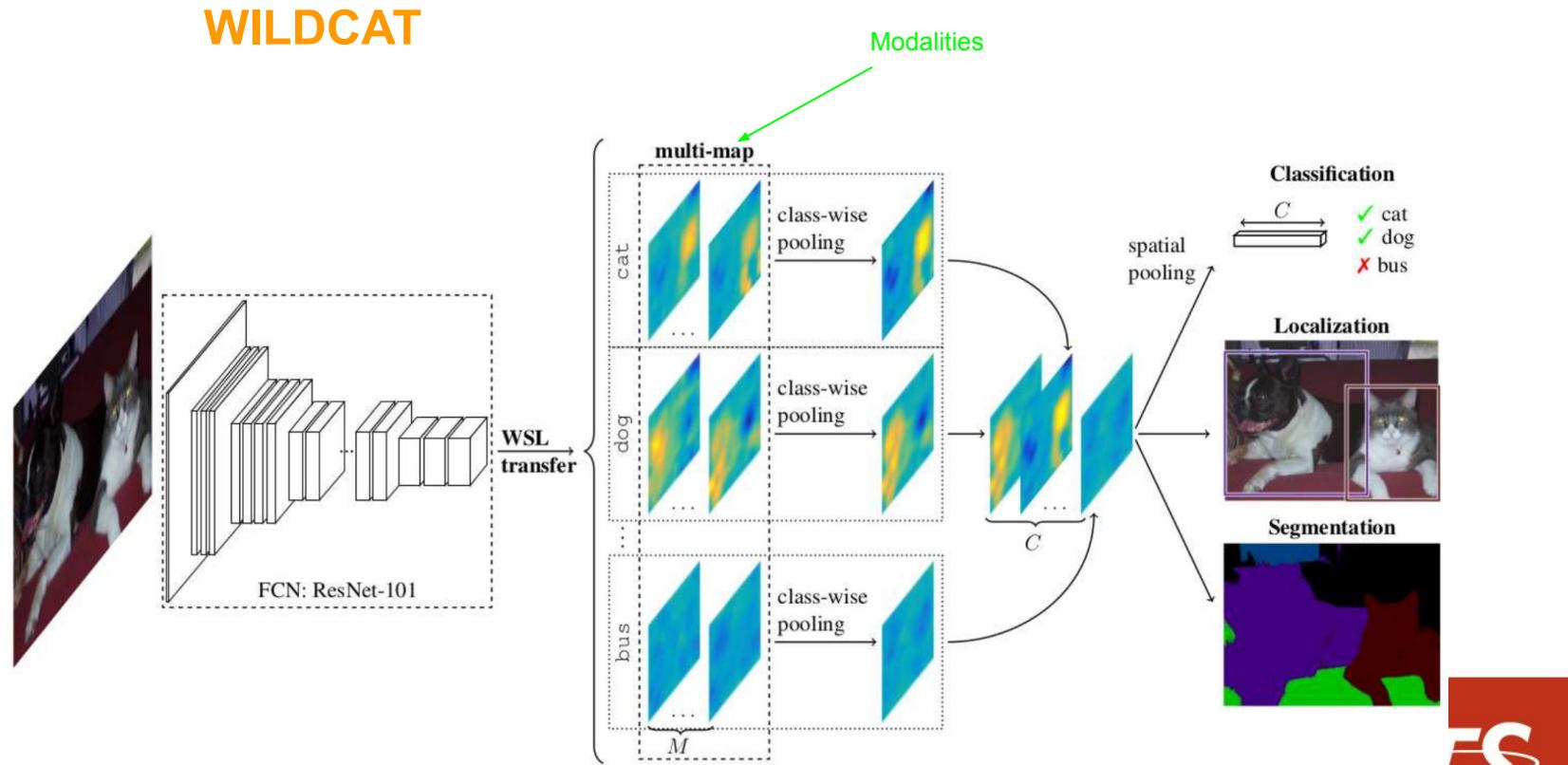
Bottom-up WSOL

Spatial pooling

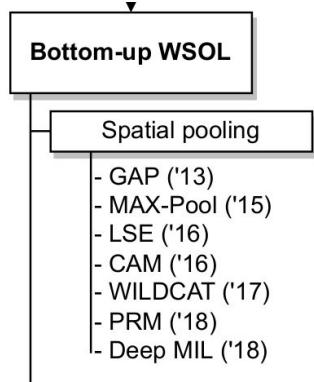
- GAP ('13)
- MAX-Pool ('15)
- LSE ('16)
- CAM ('16)
- WILDCAT ('17)
- PRM ('18)
- Deep MIL ('18)

Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps

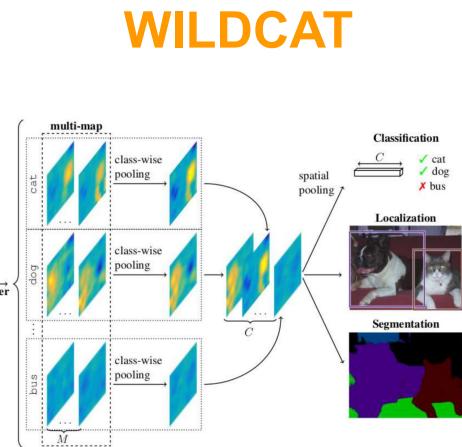


Part 2. Review of WSOL methods: Literature

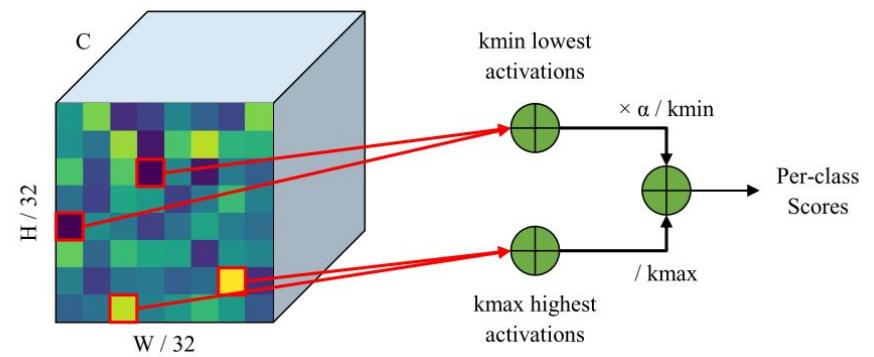


Taxonomy: Bottom-up Spatial pooling

Spatial pooling function = function to compute posterior per-class probability from spatial maps



$$s^c = \max_{\mathbf{h} \in \mathcal{H}_{k+}} \frac{1}{k^+} \sum_{i,j} h_{i,j} \bar{z}_{i,j}^c + \alpha \left(\min_{\mathbf{h} \in \mathcal{H}_{k-}} \frac{1}{k^-} \sum_{i,j} h_{i,j} \bar{z}_{i,j}^c \right)$$



Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up
Spatial pooling

Bottom-up WSOL

Spatial pooling

- GAP ('13)
- MAX-Pool ('15)
- LSE ('16)
- CAM ('16)
- WILDCAT ('17)
- PRM ('18)
- Deep MIL ('18)

Main issue:
**Minimal coverage of objects
(most discriminative parts only)**



Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up
Spatial pooling

Bottom-up WSOL

Spatial pooling

- GAP ('13)
- MAX-Pool ('15)
- LSE ('16)
- CAM ('16)
- WILDCAT ('17)
- PRM ('18)
- Deep MIL ('18)

Main issue:
**Minimal coverage of objects
(most discriminative parts only)**



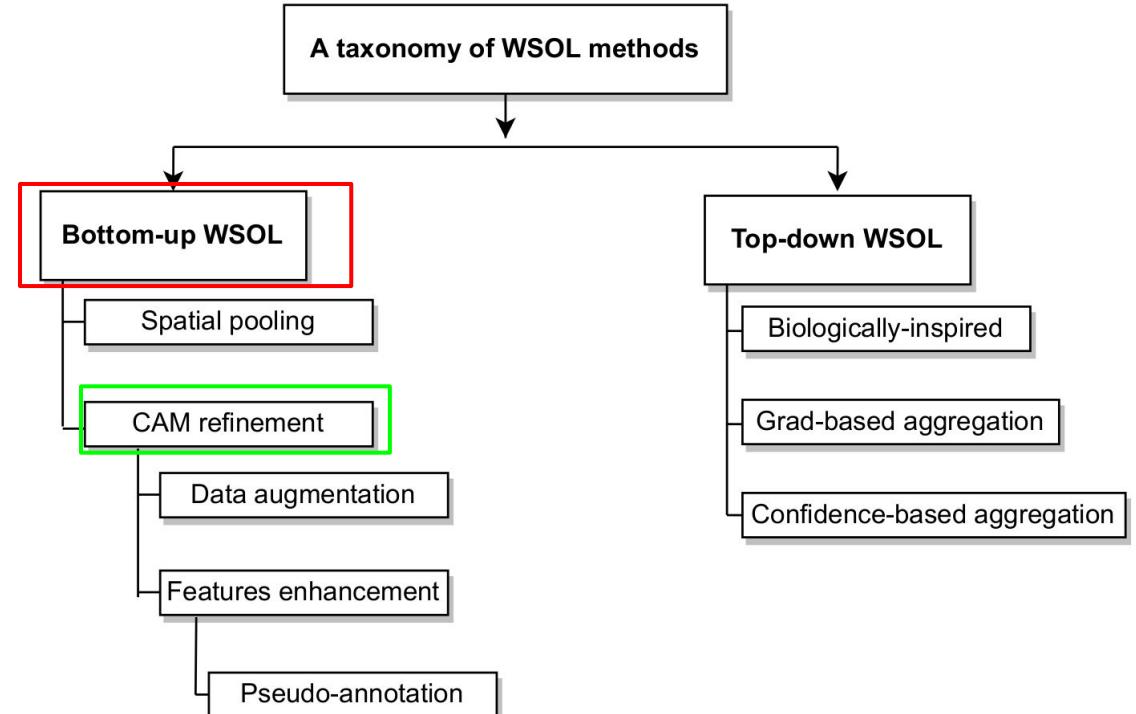
Solution: Refine the CAM!!!

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up CAM refinement

CAM refinement:

- Data augmentation
- Features enhancement
- Pseudo-annotation



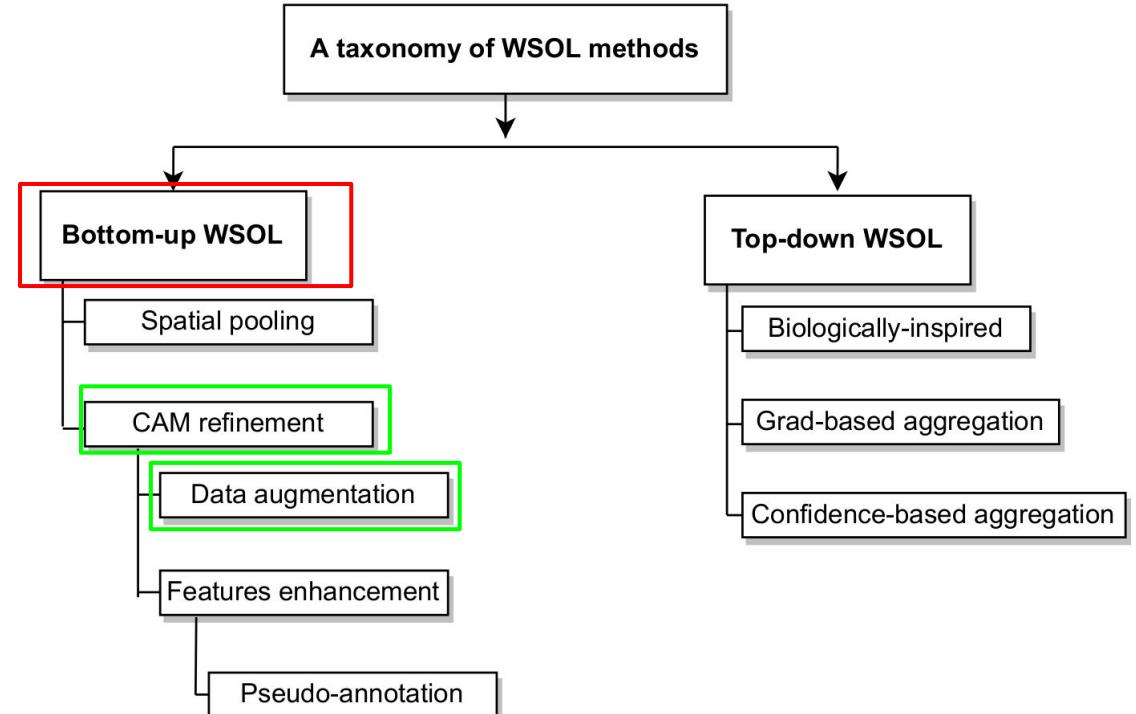
Goal: how to recover full object?

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up CAM refinement

CAM refinement:

- **Data augmentation**
- Features enhancement
- Pseudo-annotation



Goal: how to recover full object?

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Data augmentation

- HaS ('17)
- SPN ('17)
- AE ('17)
- Two-Phase ('17)
- ACoL ('18)
- GAIN ('18)
- CutMix ('19)
- ADL ('19)
- RecMin ('19)
- PuzzleMix ('20)
- MEIL ('20)
- GC-Net ('20)
- SaliencyMix ('21)
- ScoreMix ('22)
- MAXMIN ('22)

- Uses simple pooling functions
- Data augmentation: **prevent model from overfitting over single small part of the object**
- **Mine** comple objects by **perturbing image / features: Information suppression (erasing)**

Part 2. Review of WSOL methods: Literature

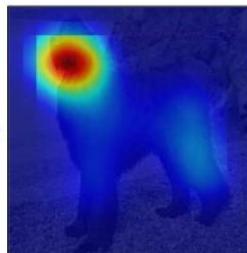
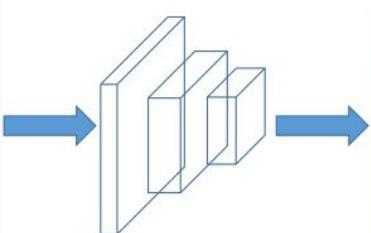
Taxonomy: Bottom-up

CAM refinement: Data augmentation

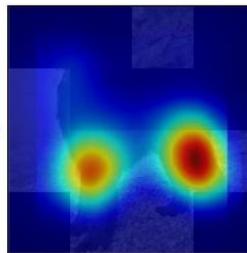
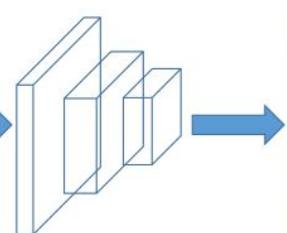
HaS: Hide and Seek



Full image



Randomly hidden patches



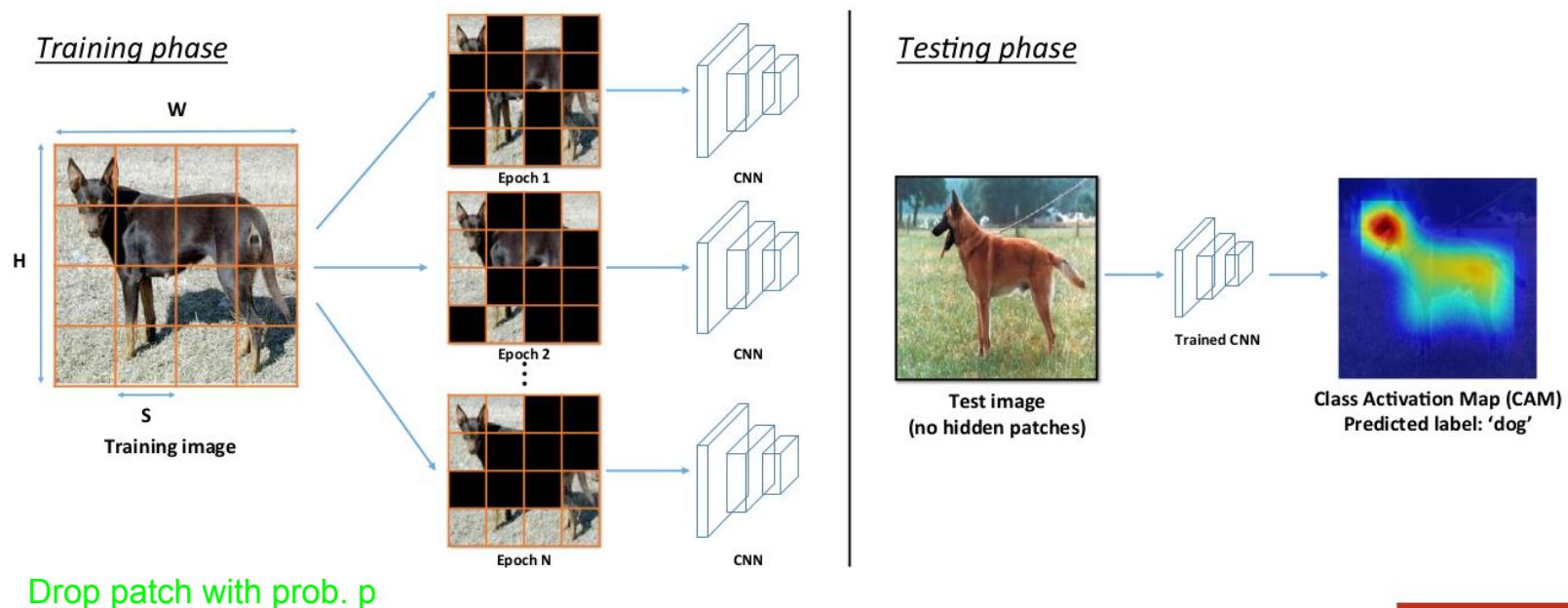
- HaS ('17)
- SPN ('17)
- AE ('17)
- Two-Phase ('17)
- ACoL ('18)
- GAIN ('18)
- CutMix ('19)
- ADL ('19)
- RecMin ('19)
- PuzzleMix ('20)
- MEIL ('20)
- GC-Net ('20)
- SaliencyMix ('21)
- ScoreMix ('22)
- MAXMIN ('22)

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Data augmentation

HaS: Hide and Seek



- HaS ('17)
- SPN ('17)
- AE ('17)
- Two-Phase ('17)
- ACoL ('18)
- GAIN ('18)
- CutMix ('19)
- ADL ('19)
- RecMin ('19)
- PuzzleMix ('20)
- MEIL ('20)
- GC-Net ('20)
- SaliencyMix ('21)
- ScoreMix ('22)
- MAXMIN ('22)

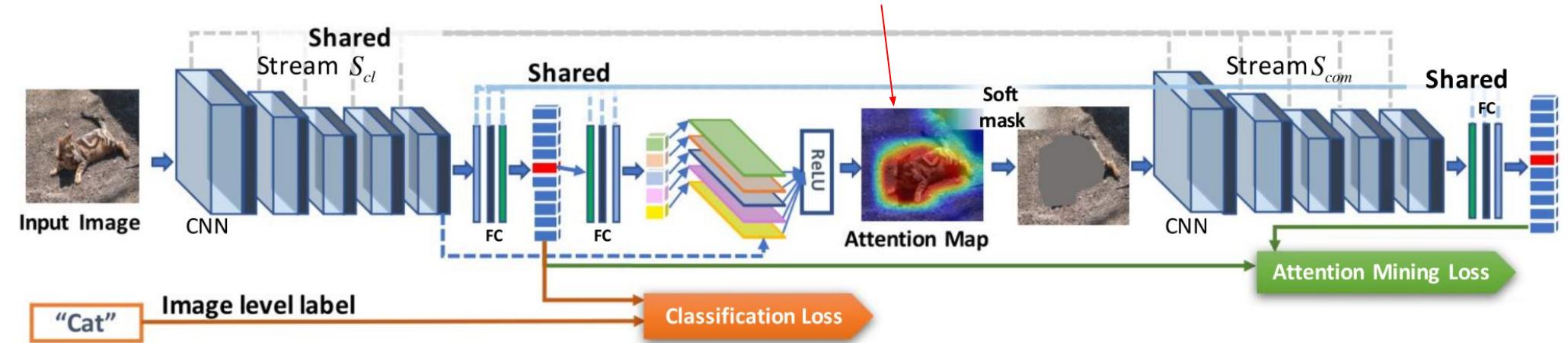
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Data augmentation

GAIN: Erasing (image)

Heated CAM

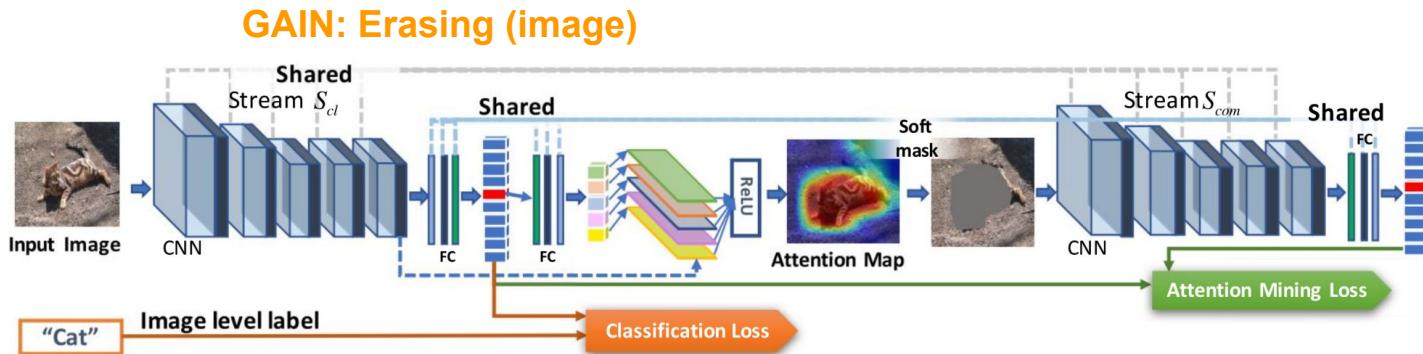


- HaS ('17)
- SPN ('17)
- AE ('17)
- Two-Phase ('17)
- ACoL ('18)
- GAIN ('18)
- CutMix ('19)
- ADL ('19)
- RecMin ('19)
- PuzzleMix ('20)
- MEIL ('20)
- GC-Net ('20)
- SaliencyMix ('21)
- ScoreMix ('22)
- MAXMIN ('22)

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Data augmentation



Train loss

$$L_{self} = L_{cl} + \alpha L_{am},$$

cross-entropy

$$L_{am} = \frac{1}{n} \sum_c s^c(I^{*c}),$$

- HaS ('17)
- SPN ('17)
- AE ('17)
- Two-Phase ('17)
- ACoL ('18)
- GAIN ('18)
- CutMix ('19)
- ADL ('19)
- RecMin ('19)
- PuzzleMix ('20)
- MEIL ('20)
- GC-Net ('20)
- SaliencyMix ('21)
- ScoreMix ('22)
- MAXMIN ('22)

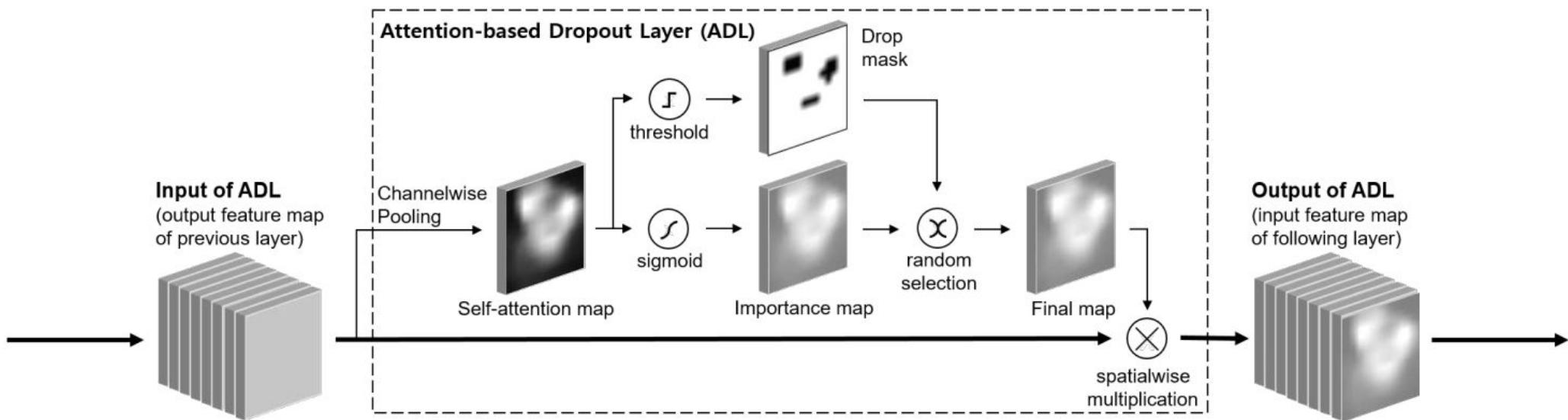
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Data augmentation

ADL: Erasing spatial features

Mining ROI via dropout!



Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Features enhancement

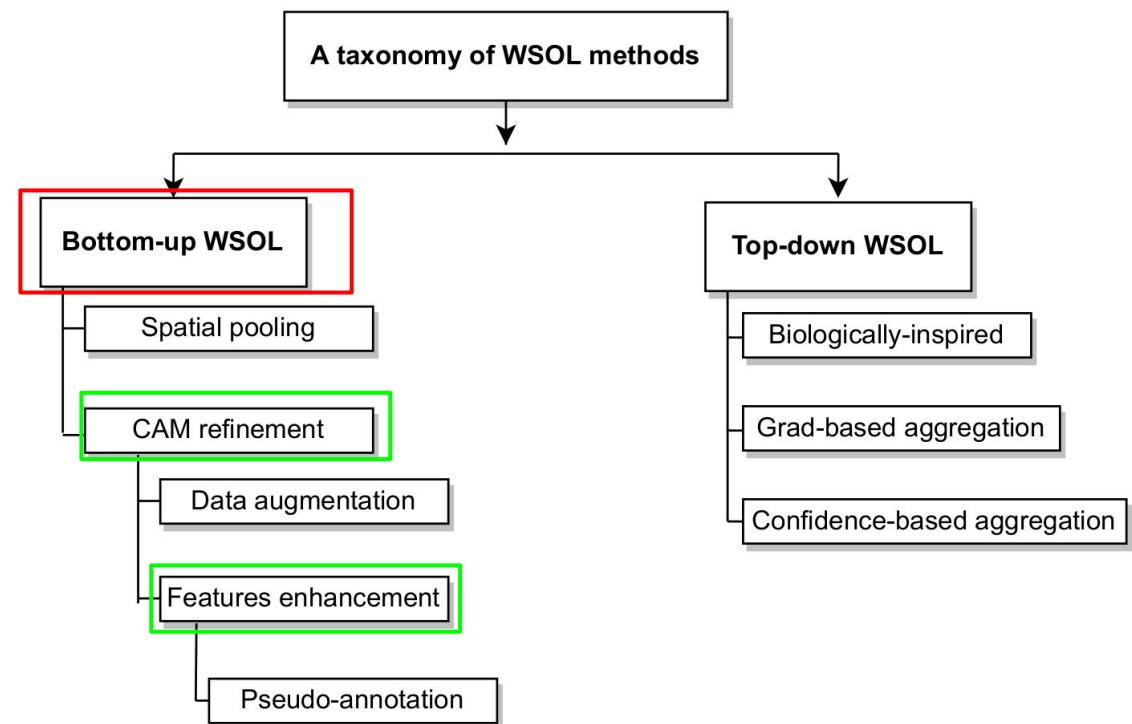
CAM refinement:

- Data augmentation
- **Features enhancement**
- Pseudo-annotation

Change architecture

+

Use low features



Part 2. Review of WSOL methods: Literature

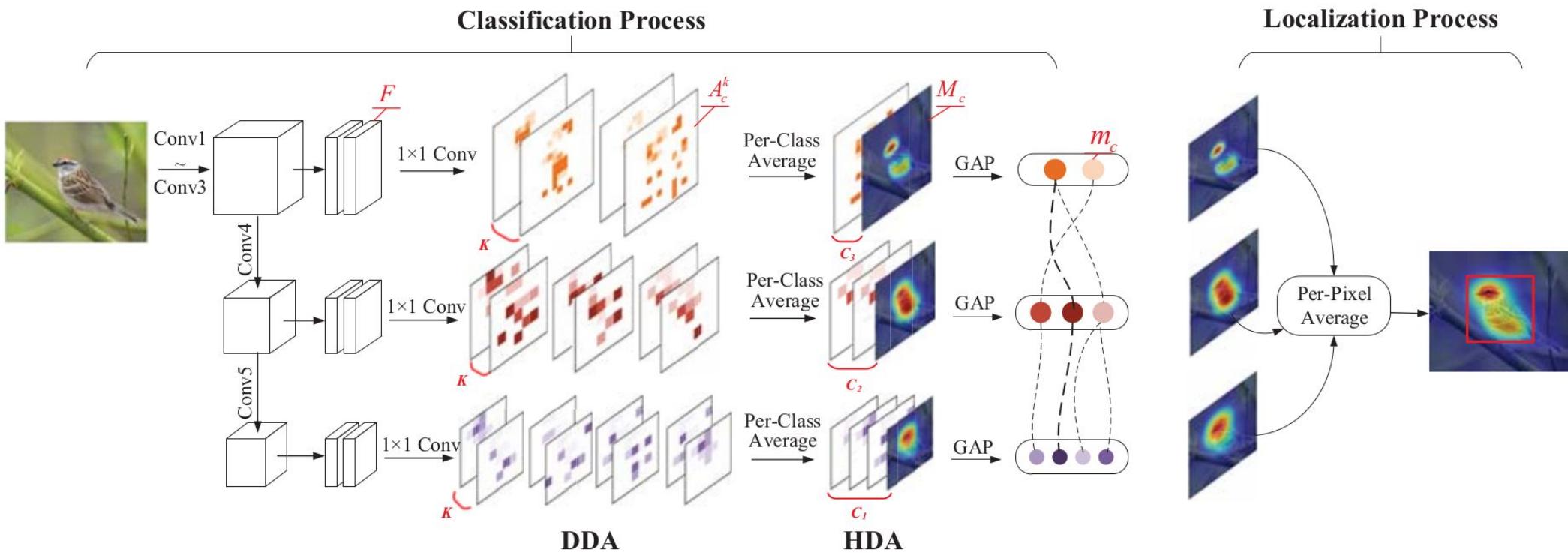
Taxonomy: Bottom-up

CAM refinement: Features enhancement

Features enhancement

- MDC ('18)
- FickleNet ('19)
- DANet ('19)
- NL-CCAM ('20)
- I2C ('20)
- ICL ('20)
- CSTN ('20)
- TS-CAM ('21)

DANet



Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

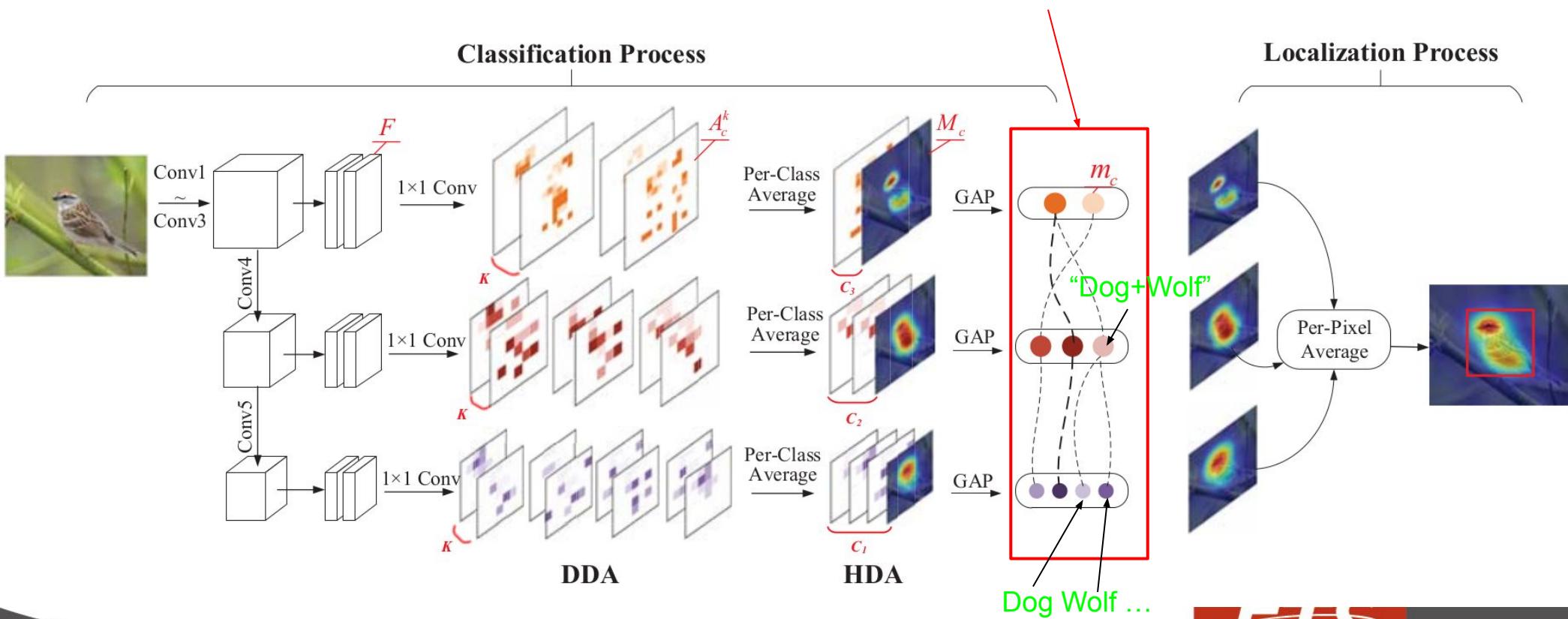
CAM refinement: Features enhancement

Features enhancement

- MDC ('18)
- FickleNet ('19)
- DANet ('19)
- NL-CCAM ('20)
- I2C ('20)
- ICL ('20)
- CSTN ('20)
- TS-CAM ('21)

Define:
a hierarchy of classes,
and
new parent-classes

DANet



Part 2. Review of WSOL methods: Literature

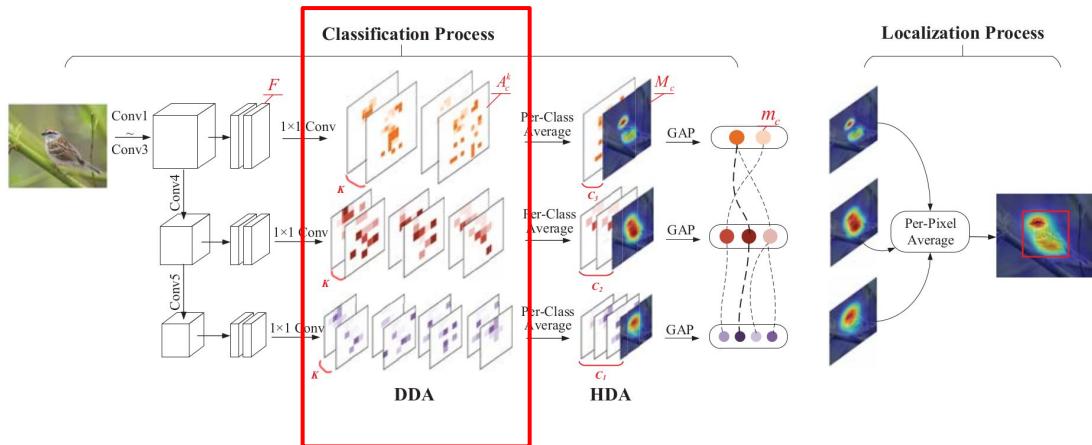
Features enhancement

- MDC ('18)
- FickleNet ('19)
- DANet ('19)
- NL-CCAM ('20)
- I2C ('20)
- ICL ('20)
- CSTN ('20)
- TS-CAM ('21)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

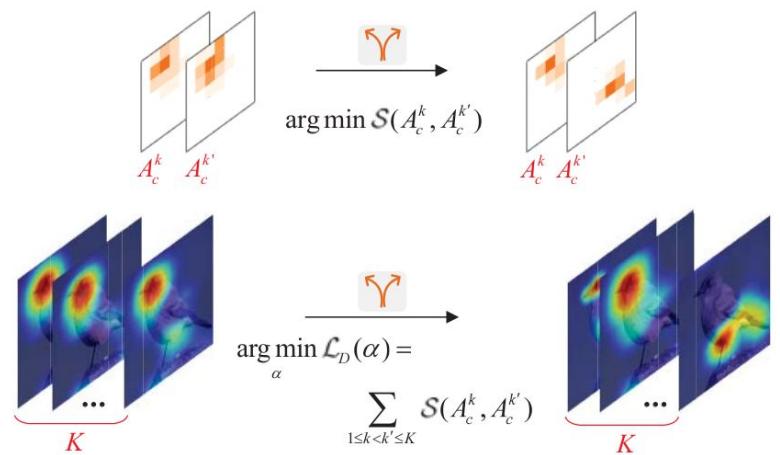
DANet



$$\arg \min_{\alpha} \mathcal{L}_D(\alpha) = \sum_{1 \leq k < k' \leq K} \mathcal{S}(A_c^k, A_c^{k'}),$$

$$\mathcal{S}(A_c^k, A_c^{k'}) = \frac{A_c^k \cdot A_c^{k'}}{\|A_c^k\| \cdot \|A_c^{k'}\|},$$

DDA loss: Discrepant Divergent Activation



Part 2. Review of WSOL methods: Literature

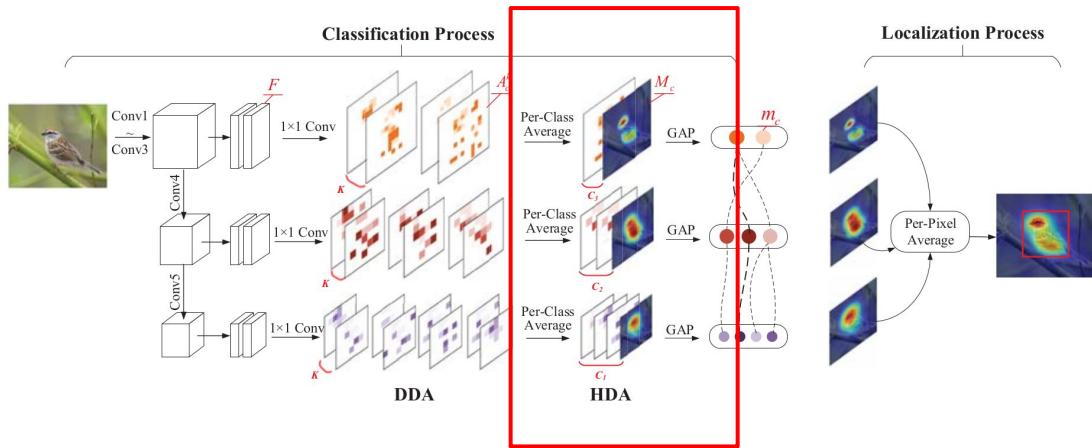
Features enhancement

- MDC ('18)
- FickleNet ('19)
- DANet ('19)
- NL-CCAM ('20)
- I2C ('20)
- ICL ('20)
- CSTN ('20)
- TS-CAM ('21)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

DANet

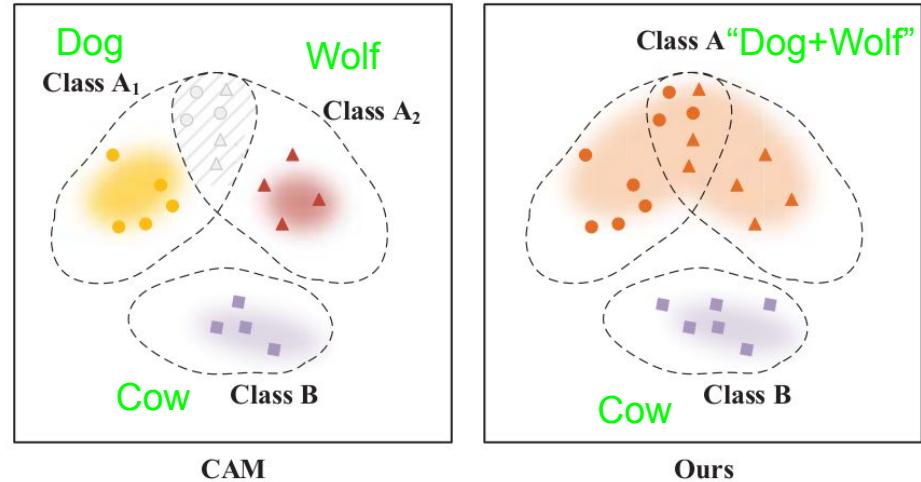


$$\arg \min_{\alpha} \mathcal{L}_H(\alpha) = \sum_h \mathcal{L}_h(\alpha) = - \sum_h \frac{1}{C^h} \sum_c y_c^h \log(p_c^h),$$

Standard classification loss

HDA loss: Hierarchical Divergent Activation

Hierarchical Divergent Activation (HDA)



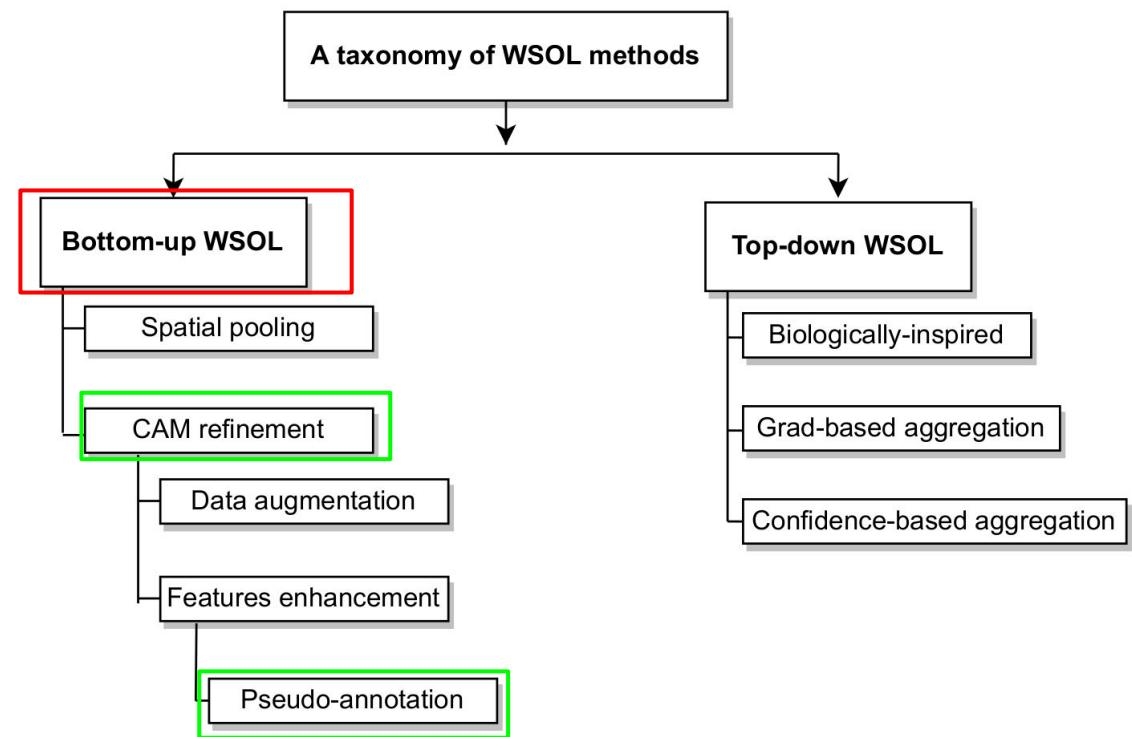
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up

CAM refinement: Features enhancement

CAM refinement:

- Data augmentation
- Features enhancement
- **Pseudo-annotation**



Part 2. Review of WSOL methods: Literature

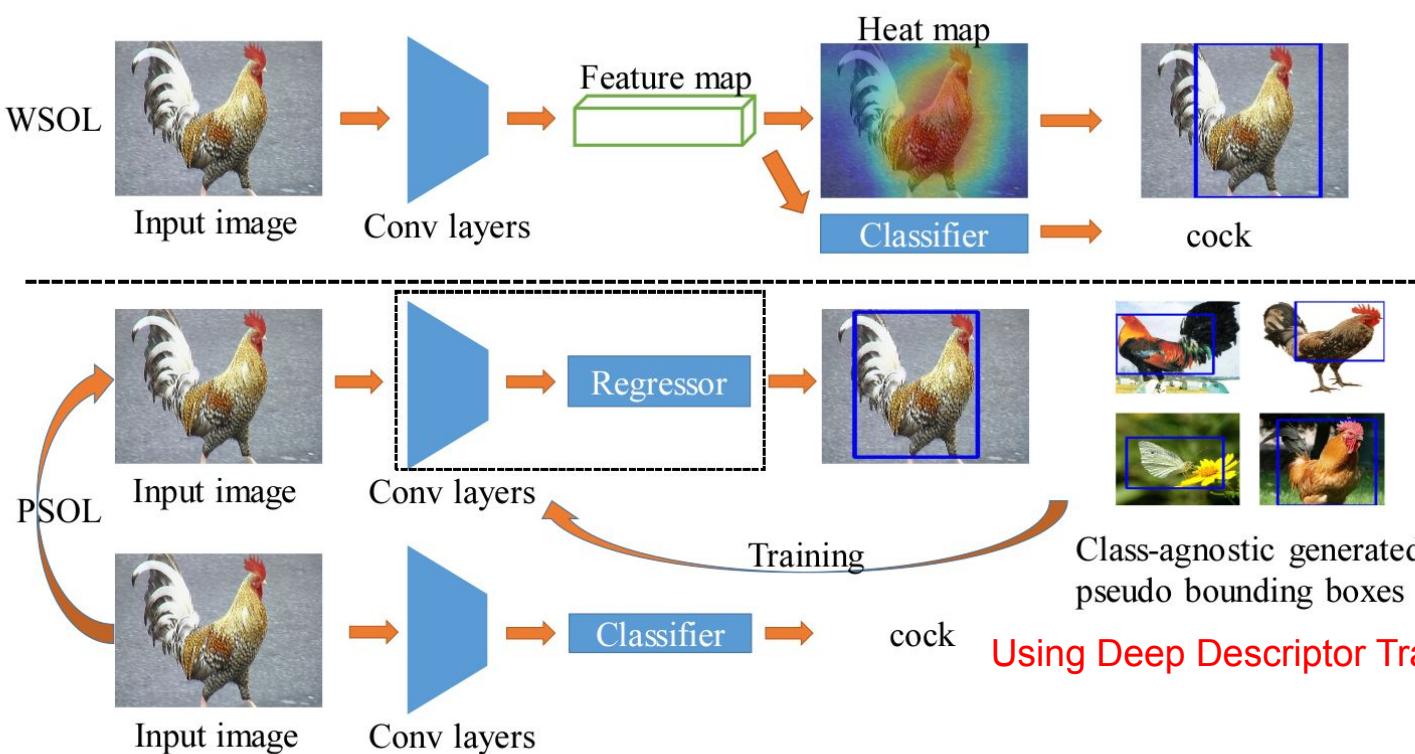
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

PSOL: Pseudo-Supervised Object Localization



Le génie pour l'industrie

[1] Zhang, C., Cao, Y., and Wu, J. (2020a). Rethinking the route towards weakly supervised object localization. In CVPR.

[2] X.-S. Wei, C.-L. Zhang, J. Wu, C. Shen, and Z.-H. Zhou. Unsupervised object discovery and co-localization by deep descriptor transformation. PR, 88:113–126, 2019.

Part 2. Review of WSOL methods: Literature

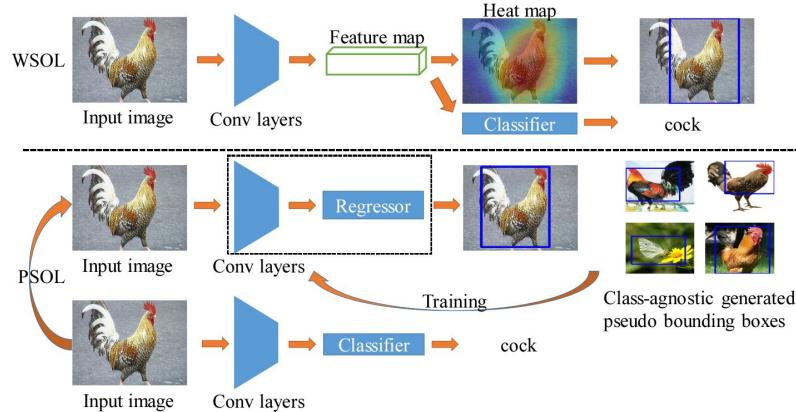
Taxonomy: Bottom-up

CAM refinement: Features enhancement

Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

PSOL: Pseudo-Supervised Object Localization



In WSOL:
Classification task
And
Localization task

Are antagonist
→ need to be separated

Algorithm 1 Pseudo Supervised Object Localization

Input: Training images I_{tr} with class label L_{tr}

Output: Predicted bounding boxes b_{te} and class labels L_{te} on testing images I_{te}

- 1: Generate pseudo bounding boxes \tilde{b}_{tr} on I_{tr}
- 2: Train a localization CNN F_{loc} on I_{tr} with \tilde{b}_{tr}
- 3: Train a classification CNN F_{cls} on I_{tr} with L_{tr}
- 4: Use F_{loc} to predict b_{te} on I_{te}
- 5: Use F_{cls} to predict L_{te} on I_{te}
- 6: **Return:** b_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

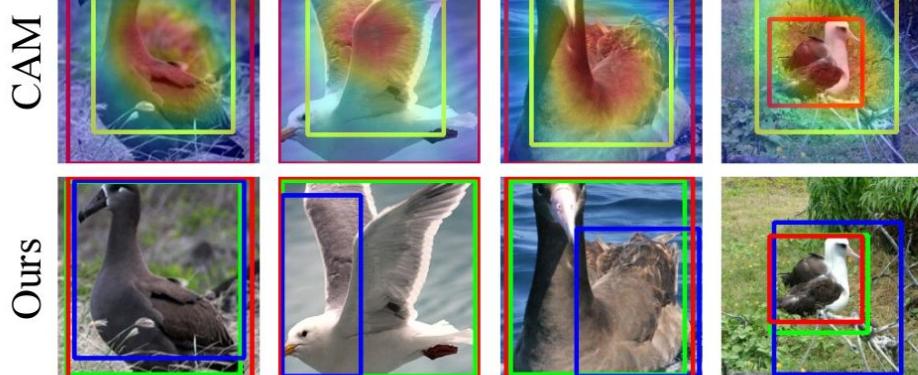
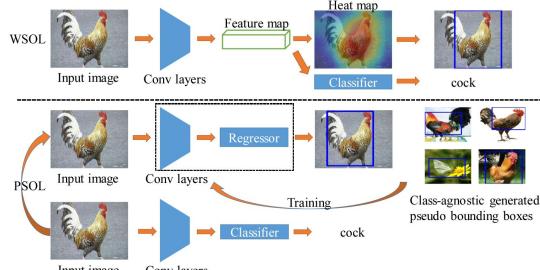
Taxonomy: Bottom-up

CAM refinement: Features enhancement

Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

PSOL: Pseudo-Supervised Object Localization



(a) CUB-200-2011



(b) ImageNet-1k

Part 2. Review of WSOL methods: Literature

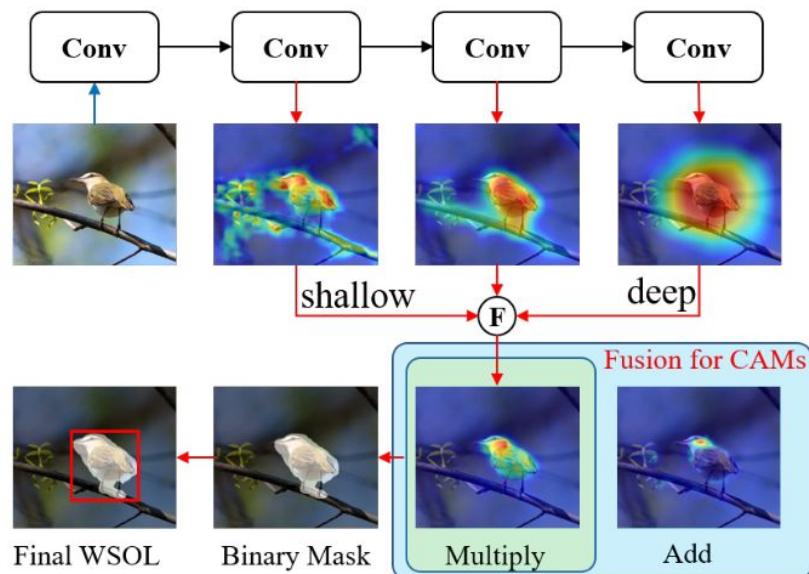
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

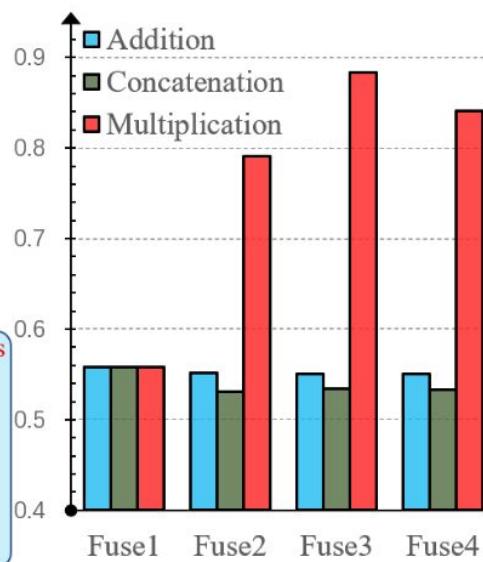
Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



(a) Improved CAMs from multiplicative fusion using shallow features for WSOL



(b) Different fusion strategies from different layers

Shallow features
are useful for
localization!

Part 2. Review of WSOL methods: Literature

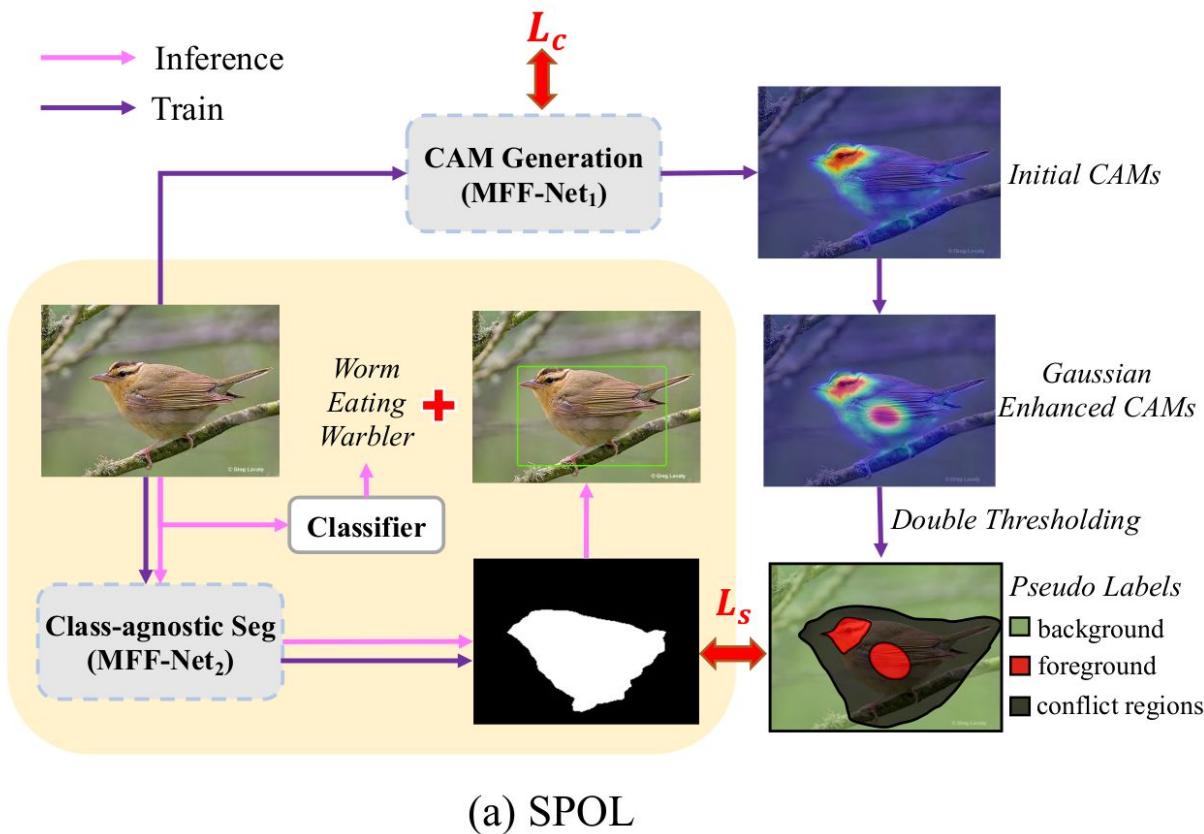
Taxonomy: Bottom-up

CAM refinement: Features enhancement

Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

SPOL: Shallow Pseudo-supervised Object Localization



Part 2. Review of WSOL methods: Literature

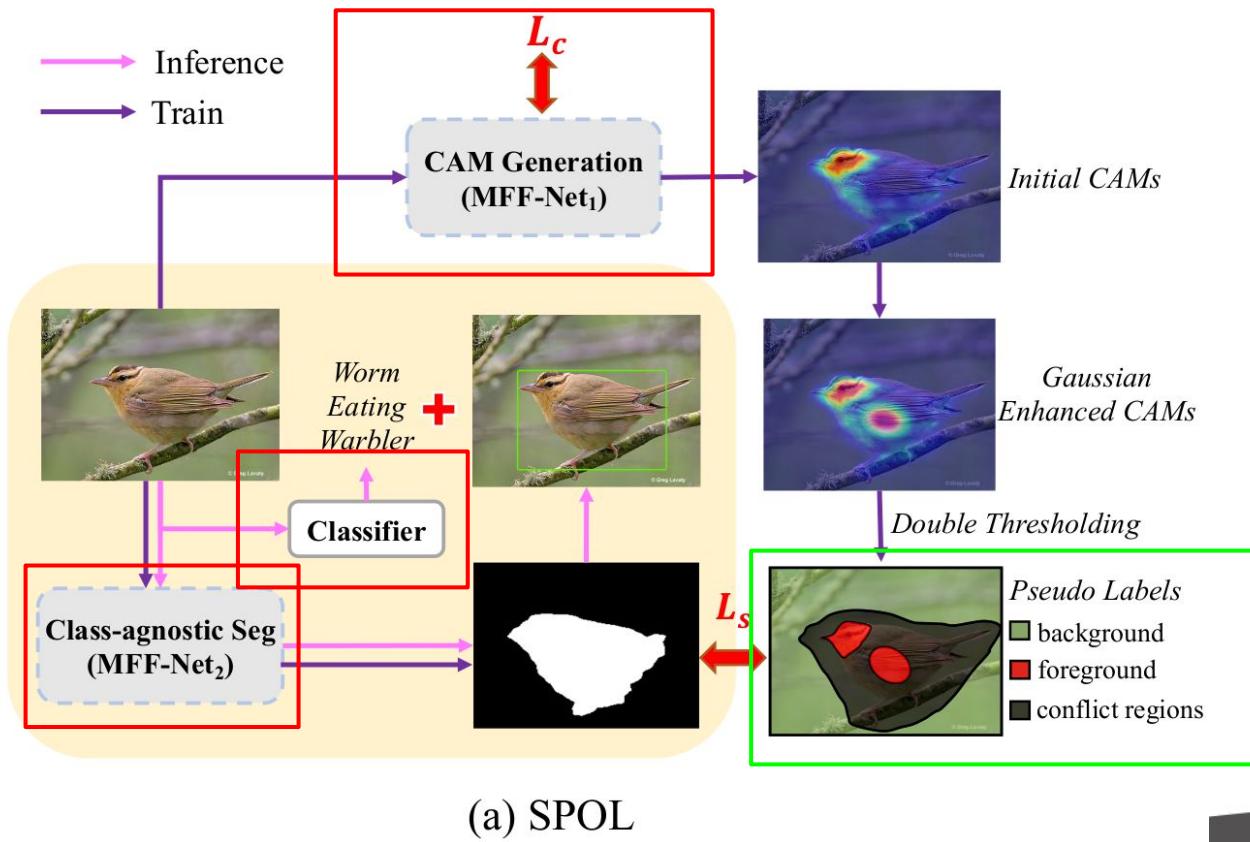
Taxonomy: Bottom-up

CAM refinement: Features enhancement

Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

SPOL: Shallow Pseudo-supervised Object Localization



Train loss:

L_c : standard cross-entropy

L_s : partial cross-entropy over only foreground and background

Part 2. Review of WSOL methods: Literature

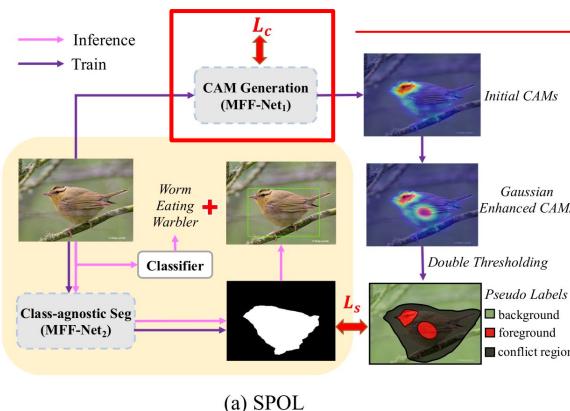
Taxonomy: Bottom-up

CAM refinement: Features enhancement

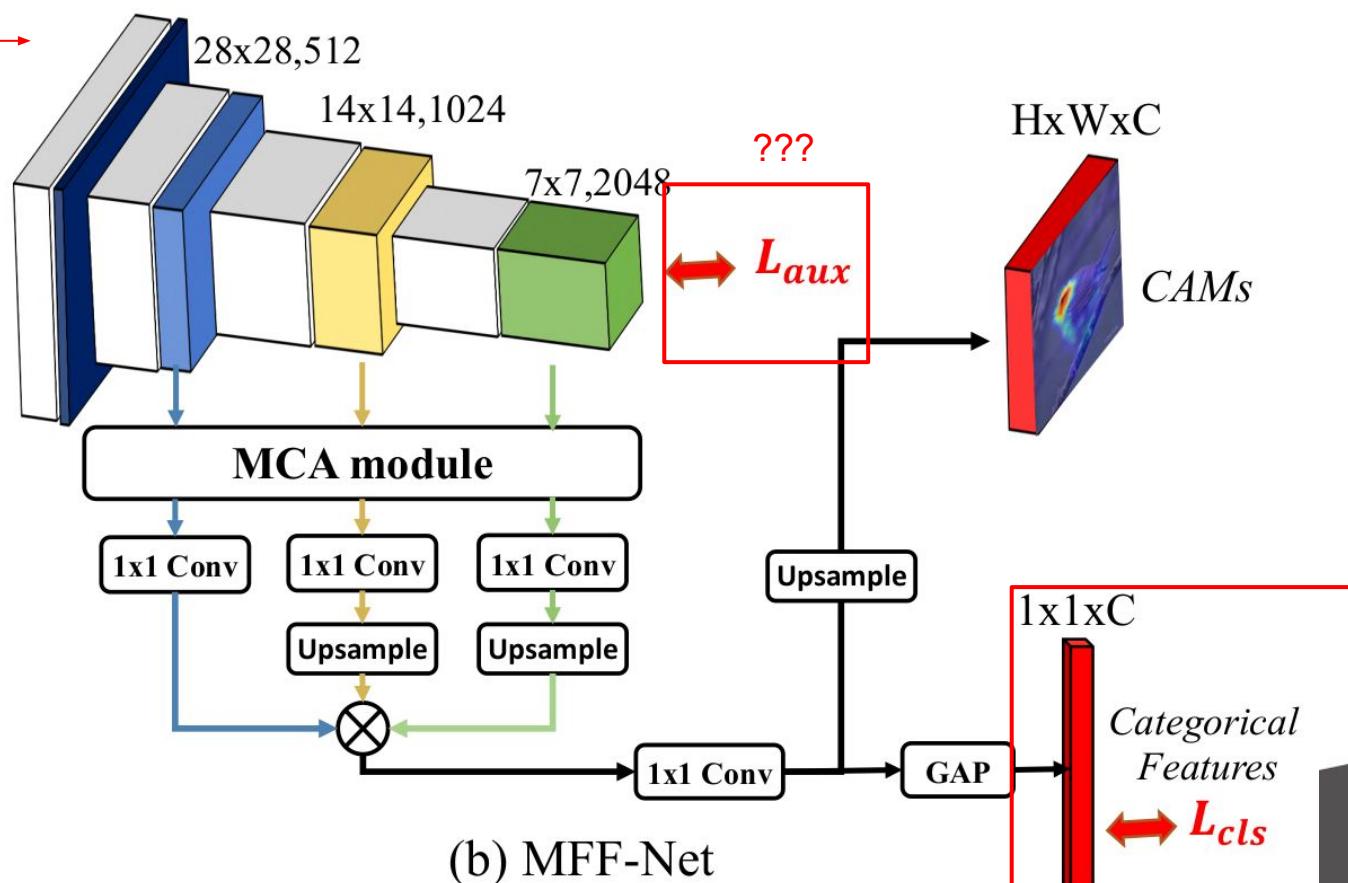
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

SPOL: Shallow Pseudo-supervised Object Localization



Heavy usage of low features!



Part 2. Review of WSOL methods: Literature

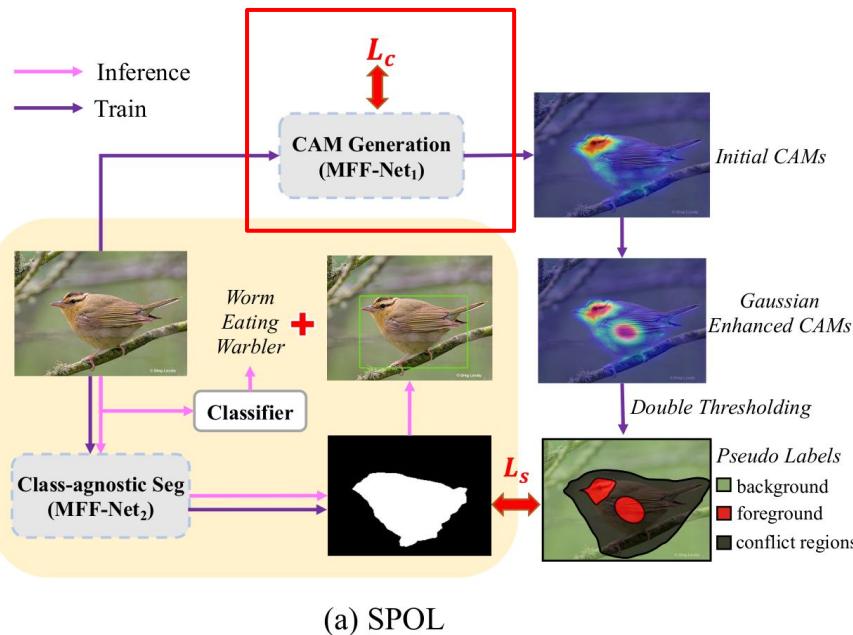
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

- 1 // Training Phase
- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}
- 6 // Inference Phase
- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}
- 10 **Return:** B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

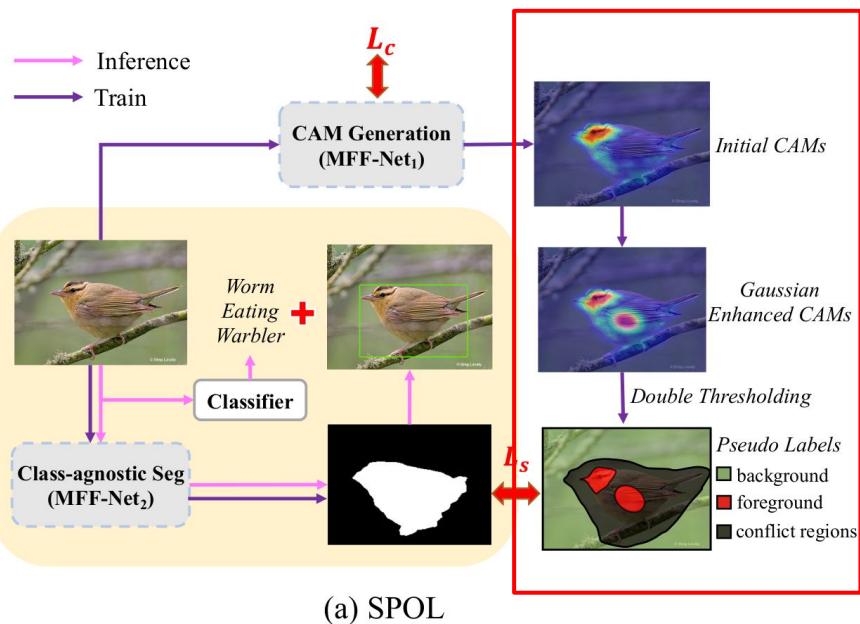
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

- 1 // Training Phase
- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}
- 6 // Inference Phase
- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}
- 10 **Return:** B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

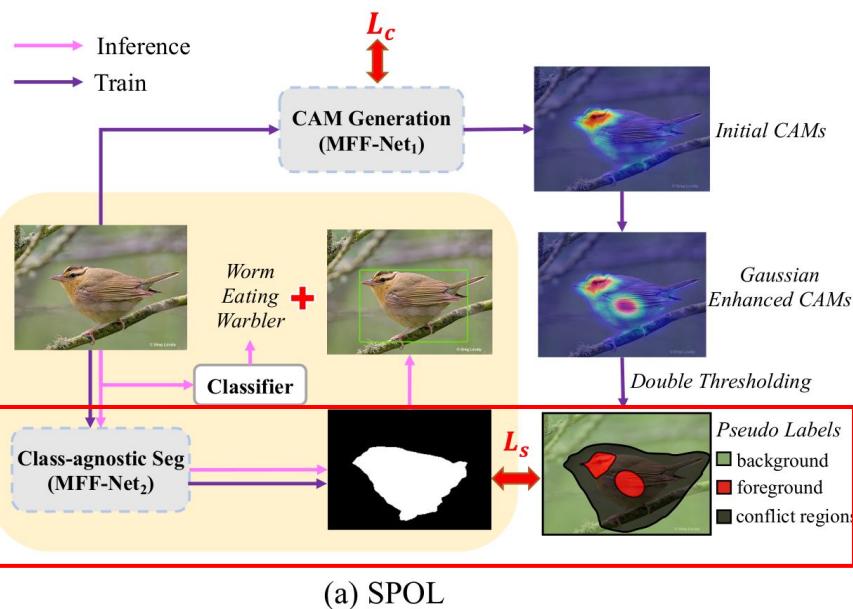
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

1 // Training Phase

- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}

6 // Inference Phase

- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}

10 Return: B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

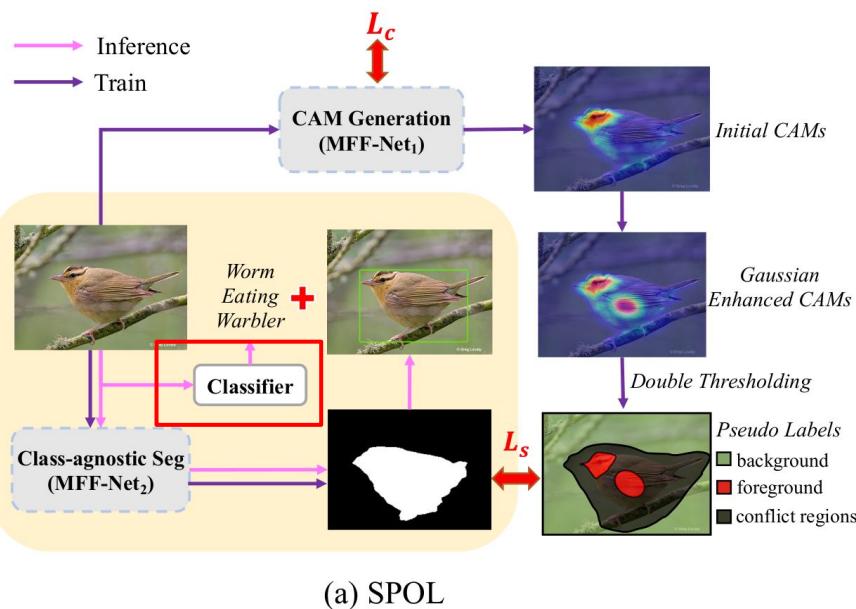
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

1 // Training Phase

- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}

6 // Inference Phase

- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}

10 Return: B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

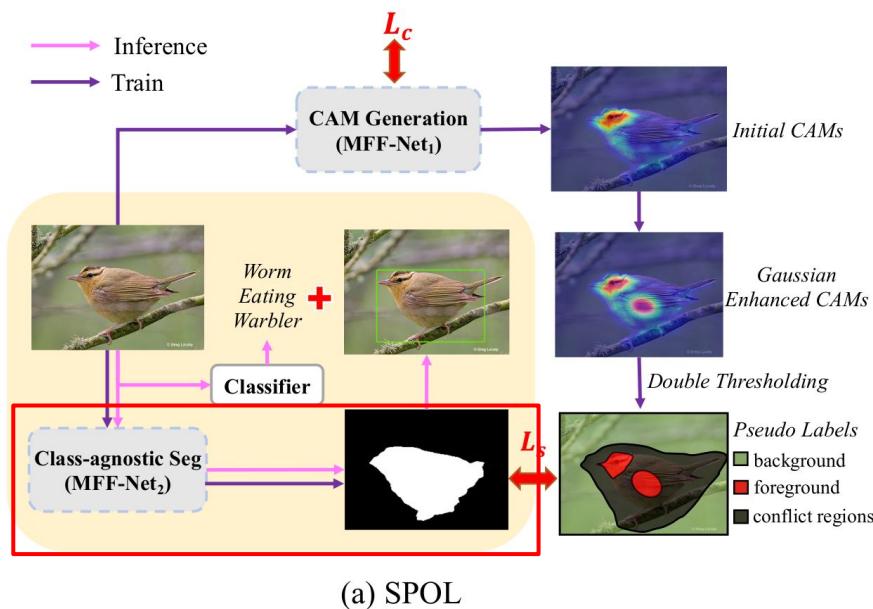
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

- 1 // Training Phase
- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}

6 // Inference Phase

- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}
- 10 **Return:** B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

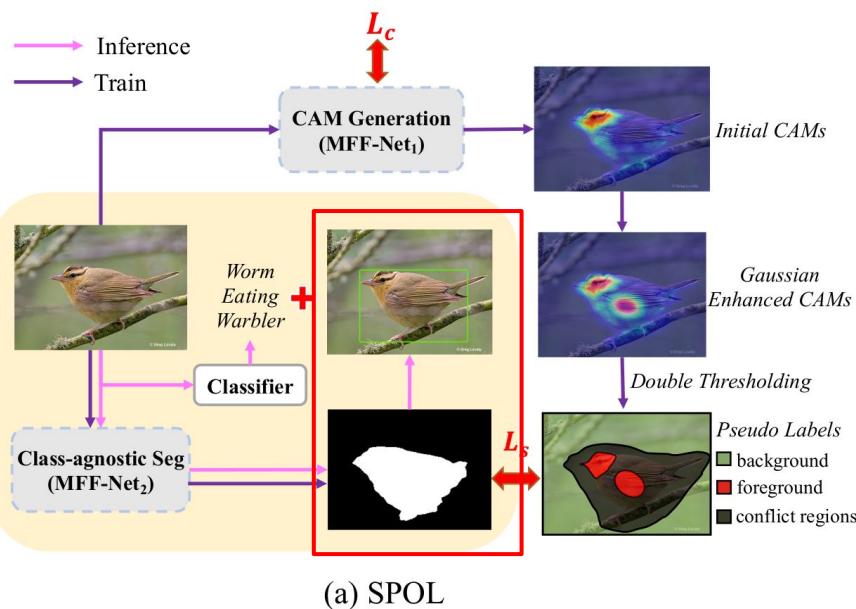
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

- 1 // Training Phase
- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}

6 // Inference Phase

- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}
- 10 **Return:** B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

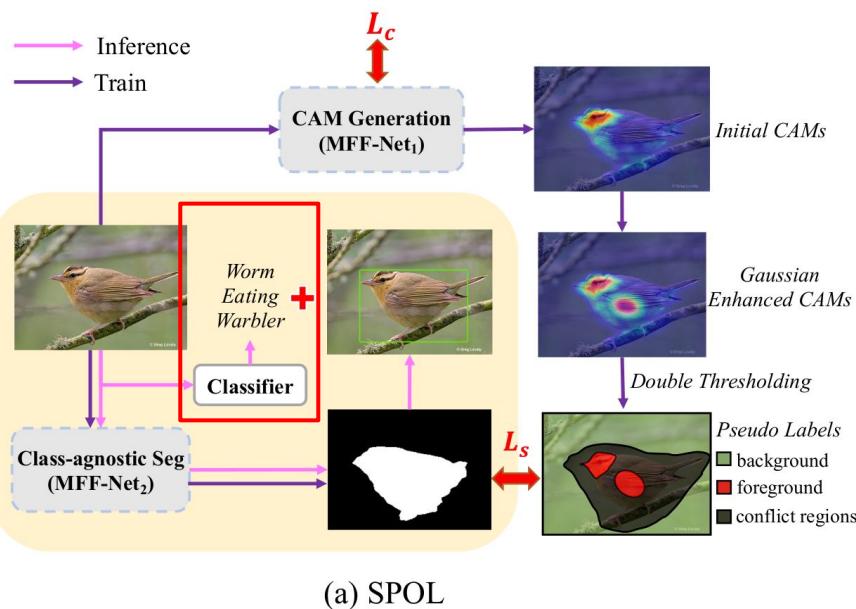
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization



Algorithm 1: SPOL

Input: Training images I_{tr} with class label L_{tr}
Output: Predicted bounding boxes B_{te} and class labels L_{te} on testing images I_{te}

- 1 // Training Phase
- 2 Train MFF-Net₁ F_w on I_{tr} with L_{tr}
- 3 Use F_w to generate pseudo label M_{tr} on I_{tr}
- 4 Train MFF-Net₂ F_s on I_{tr} for Seg. with M_{tr}
- 5 Train a classifier F_c on I_{tr} with L_{tr}

6 // Inference Phase

- 7 Use F_s to predict M_{te} on I_{te}
- 8 Extract object bounding box B_{te} from M_{te}
- 9 Use F_c to predict L_{te} on I_{te}

10 **Return:** B_{te}, L_{te}

Part 2. Review of WSOL methods: Literature

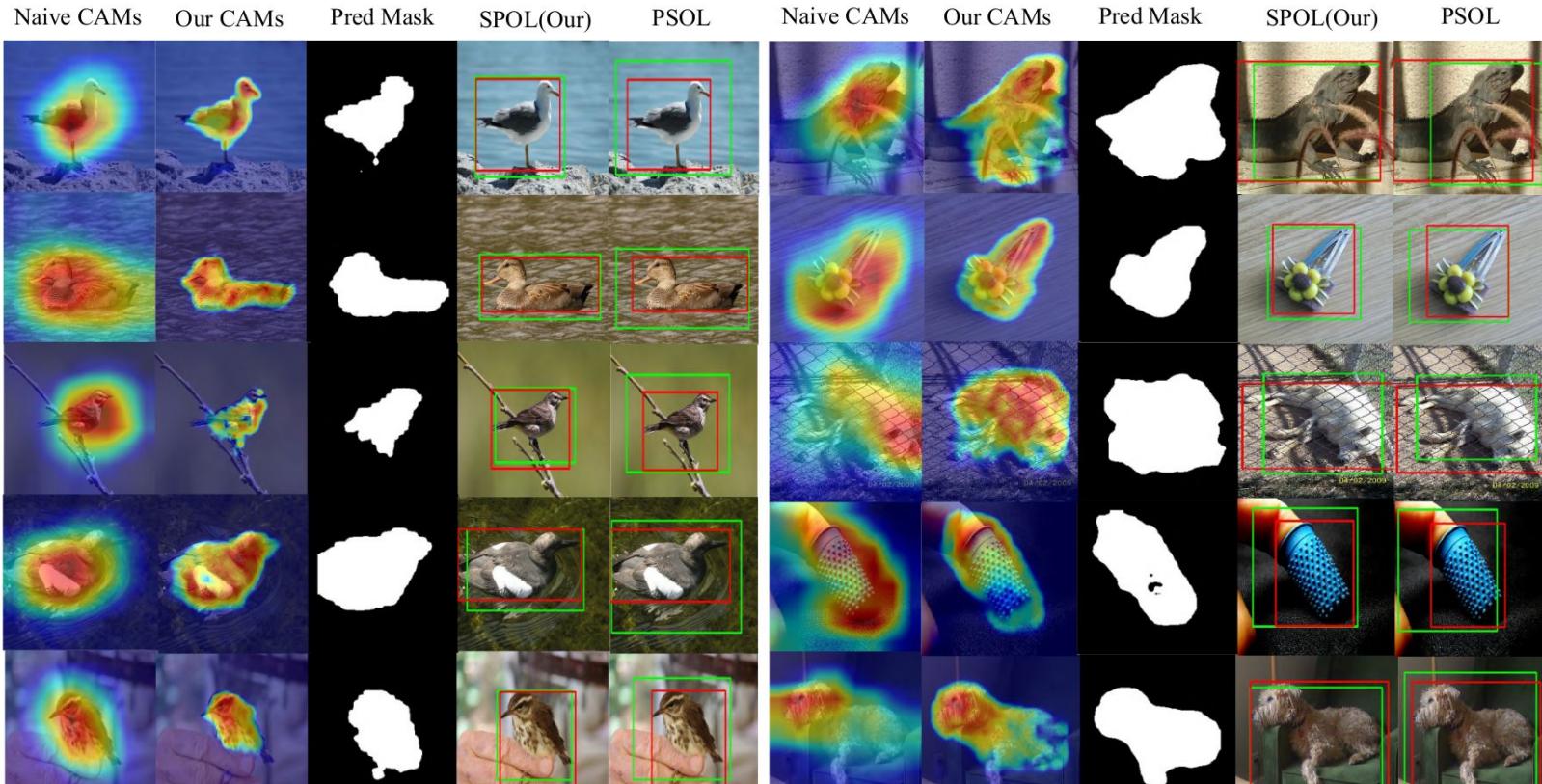
Pseudo-annotation

- SPG ('18)
- Pair-Sim ('20)
- PSOL ('20)
- SPOL ('21)
- F-CAM ('22)
- NEGEV ('22)

Taxonomy: Bottom-up

CAM refinement: Features enhancement

SPOL: Shallow Pseudo-supervised Object Localization

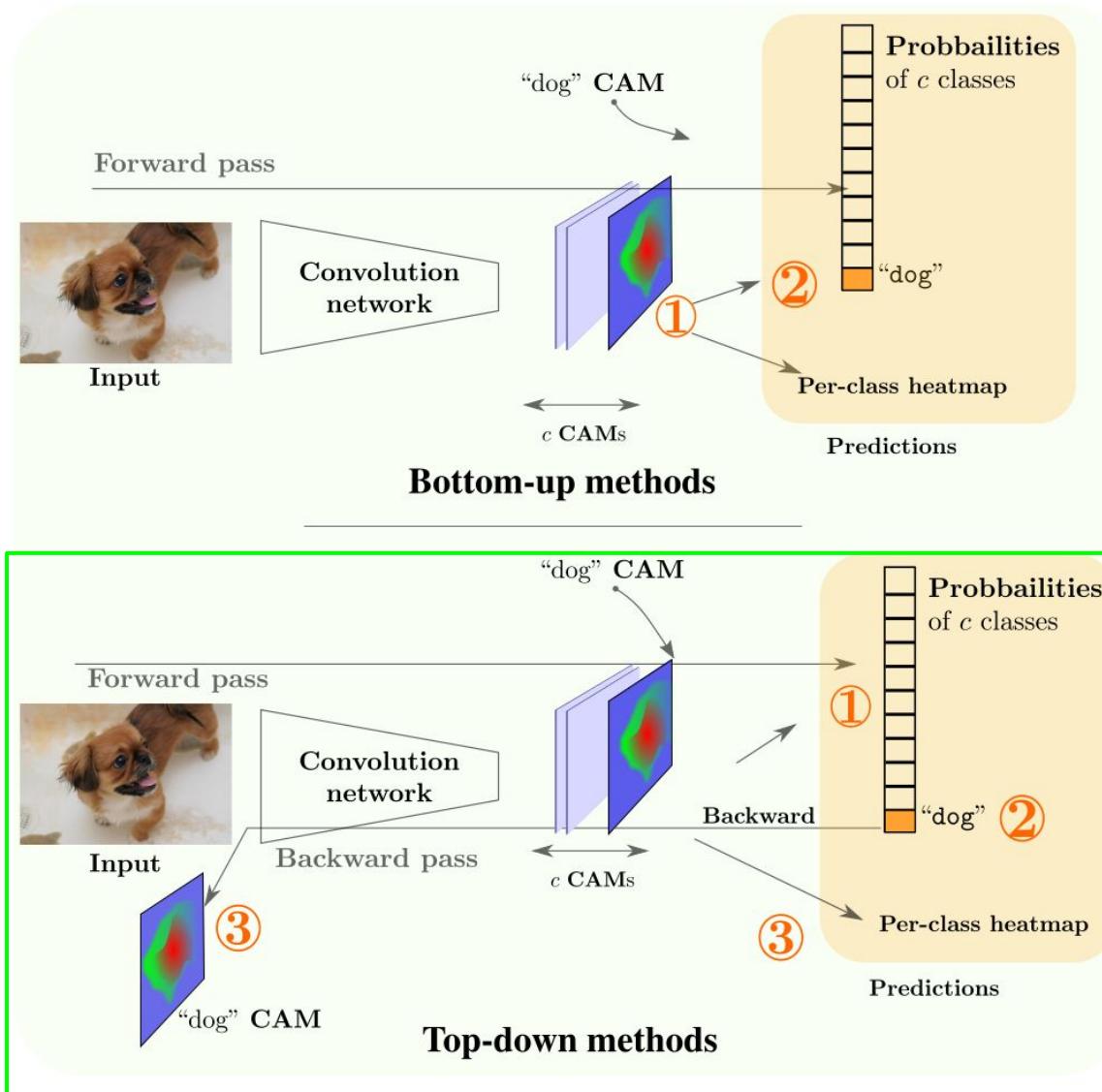


(a) CUB200 dataset

(b) ImageNet-1K dataset

Part 2. Review of WSOL methods: Literature

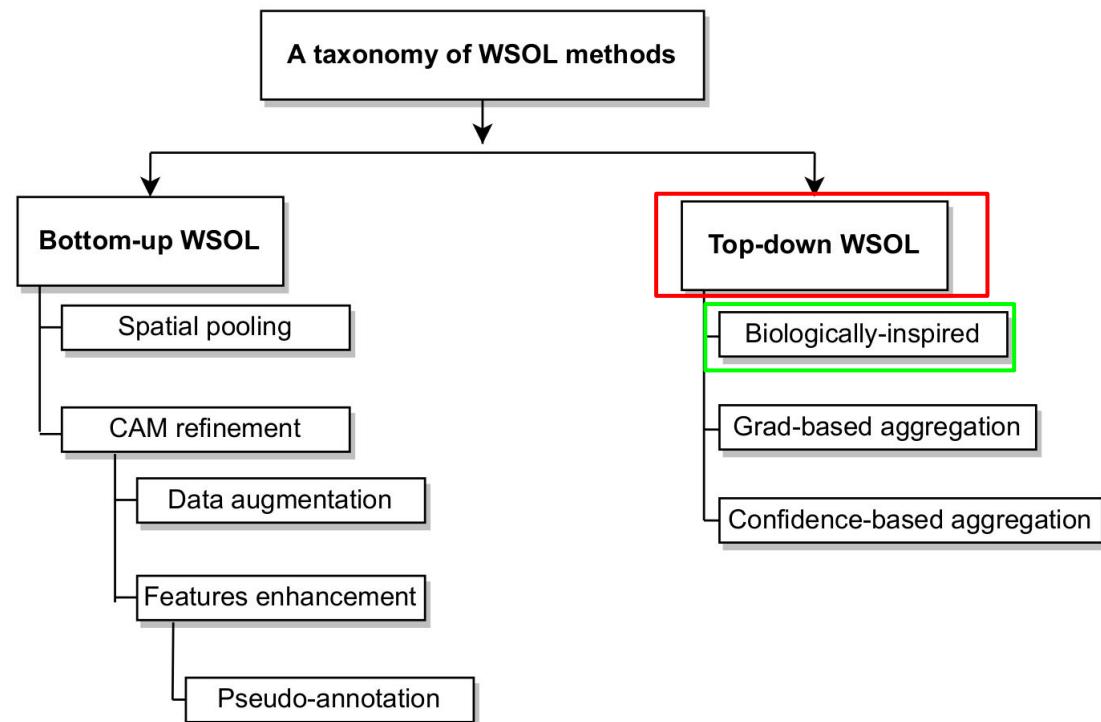
Taxonomy



Part 2. Review of WSOL methods: Literature

Taxonomy

- Cognitive science
- Human visual attention (top-down mechanism)



Part 2. Review of WSOL methods: Literature

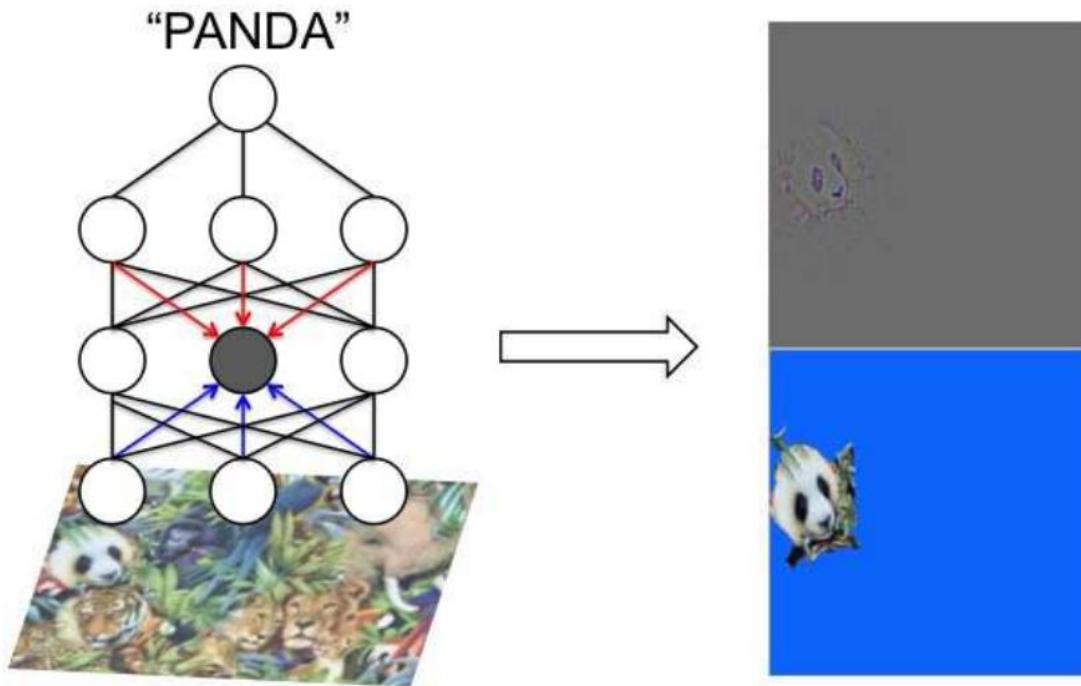
Biologically-inspired

- Feedback layer ('15)
- Excitation-backprop ('18)

Taxonomy: Top-Down Biologically-inspired

Feedback layer

Use bottom-up
input image
top-down semantic
label to infer
hidden neuron
activation



Le génie pour l'industrie

Part 2. Review of WSOL methods: Literature

Biologically-inspired

- Feedback layer ('15)
- Excitation-backprop ('18)

Taxonomy: Top-Down Biologically-inspired

Feedback layer

Localization over
input neurons

Bounding box
localization



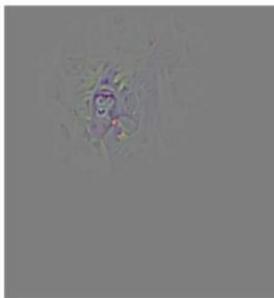
(a) Input Image



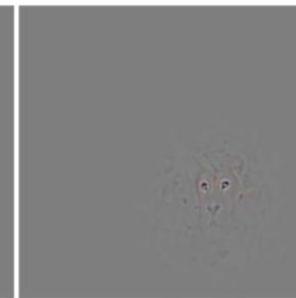
(b) Panda



(c) Tiger



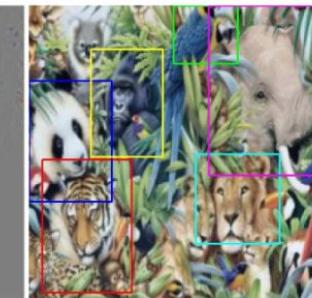
(d) Gorilla



(e) Lion



(f) Elephant



(g) Localization

Part 2. Review of WSOL methods: Literature

Biologically-inspired

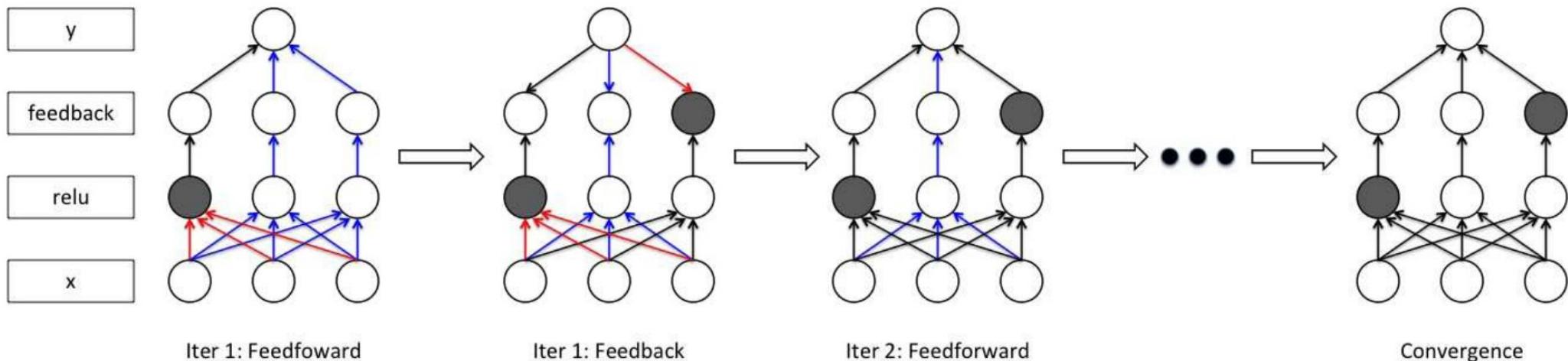
- Feedback layer ('15)
- Excitation-backprop ('18)

Taxonomy: Top-Down Biologically-inspired

Feedback layer

Feedback layer: network of binary control gates z

Iterative optimization process to find z



Part 2. Review of WSOL methods: Literature

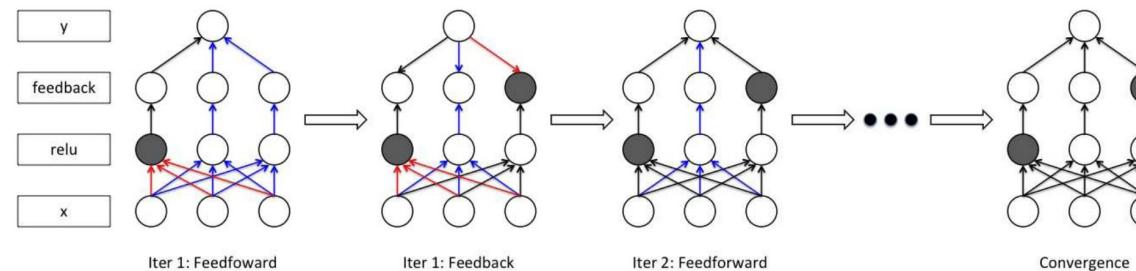
Biologically-inspired

- Feedback layer ('15)
- Excitation-backprop ('18)

Taxonomy: Top-Down Biologically-inspired

Feedback layer

Feedback layer: network of binary control gates \mathbf{z}



Promote selectivity

Iterative optimization
process to find \mathbf{z}
(relaxed)

$$\begin{aligned} \max_{\mathbf{z}} \quad & s_k(I, \mathbf{z}) - \lambda \|\mathbf{z}\|_1 \\ \text{s.t.} \quad & 0 \leq z_{i,j,c}^l \leq 1, \forall l, i, j, c \end{aligned}$$

Gradient ascent:

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \alpha \cdot \left(\frac{\partial s_k}{\partial \mathbf{z}} |_{\mathbf{z}_t} - \lambda \right)$$

Part 2. Review of WSOL methods: Literature

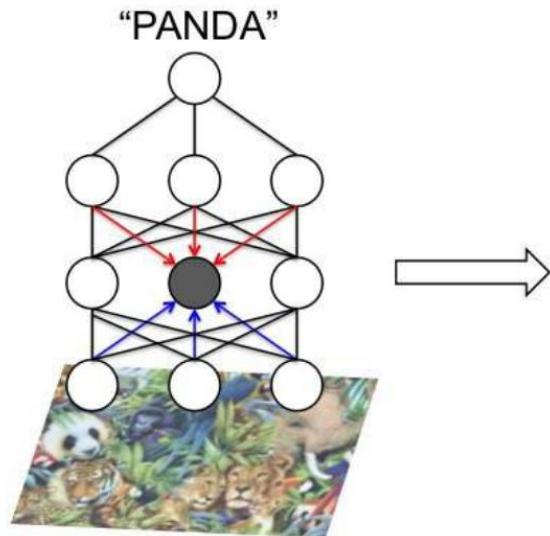
Biologically-inspired

- Feedback layer ('15)
- Excitation-backprop ('18)

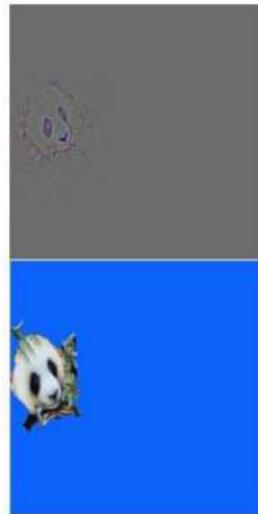
Taxonomy: Top-Down Biologically-inspired

Feedback layer

Feedback layer: network of binary control gates \mathbf{z}



Promote selectivity



Iterative optimization process to find \mathbf{z} (relaxed)

$$\begin{aligned} \max_{\mathbf{z}} \quad & s_k(I, \mathbf{z}) - \lambda \|\mathbf{z}\|_1 \\ \text{s.t.} \quad & 0 \leq z_{i,j,c}^l \leq 1, \forall l, i, j, c \end{aligned}$$

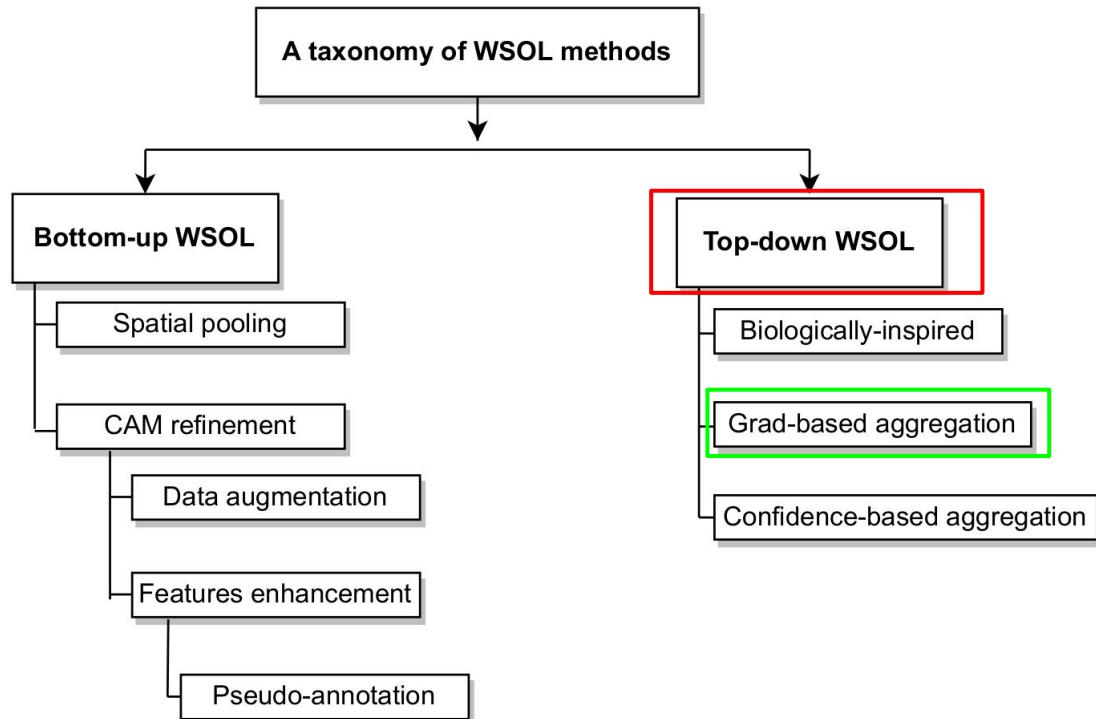
Gradient ascent:

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \alpha \cdot \left(\frac{\partial s_k}{\partial \mathbf{z}} |_{\mathbf{z}_t} - \lambda \right)$$

Part 2. Review of WSOL methods: Literature

Taxonomy

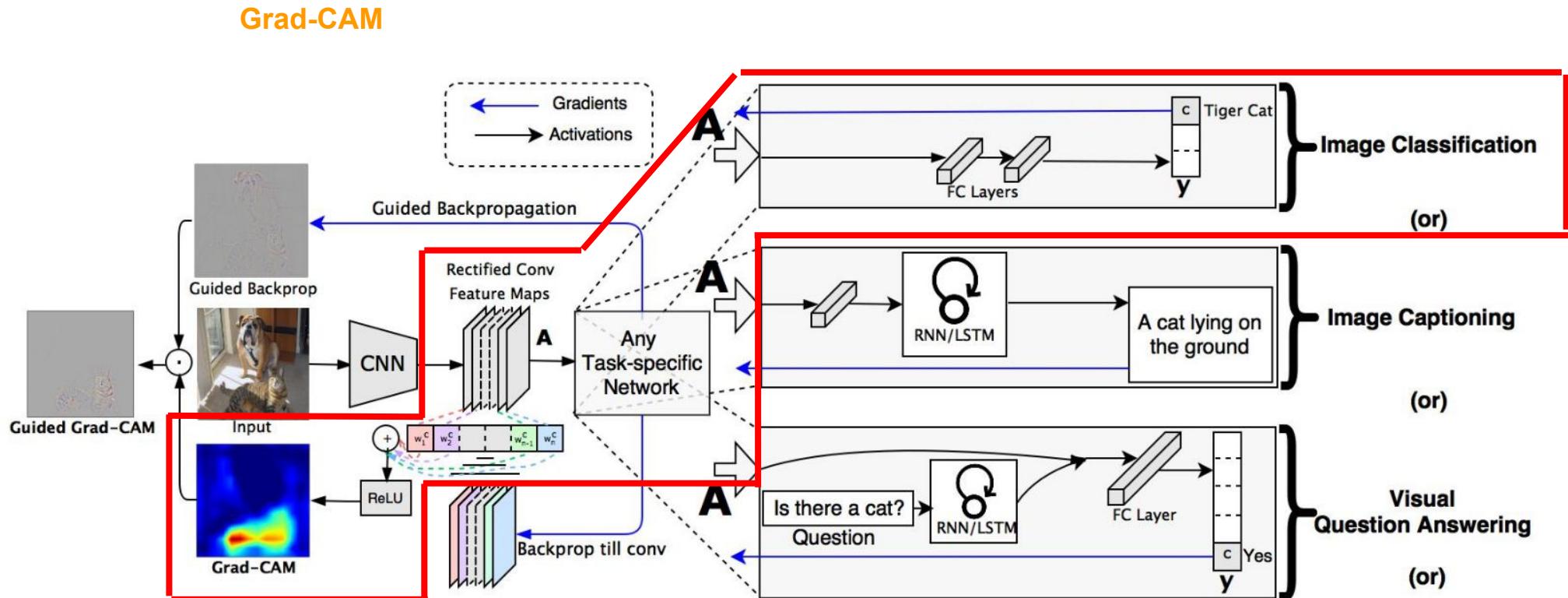
- Aggregation of feature maps using gradient



Part 2. Review of WSOL methods: Literature

Taxonomy: Top-Down Grad-based aggregation

- Grad-based aggregation
 - Grad-CAM ('17)
 - Grad-CAM++ ('18)
 - Smooth-Grad-CAM++ ('19)
 - XGrad-CAM ('20)
 - LayerCAM ('21)

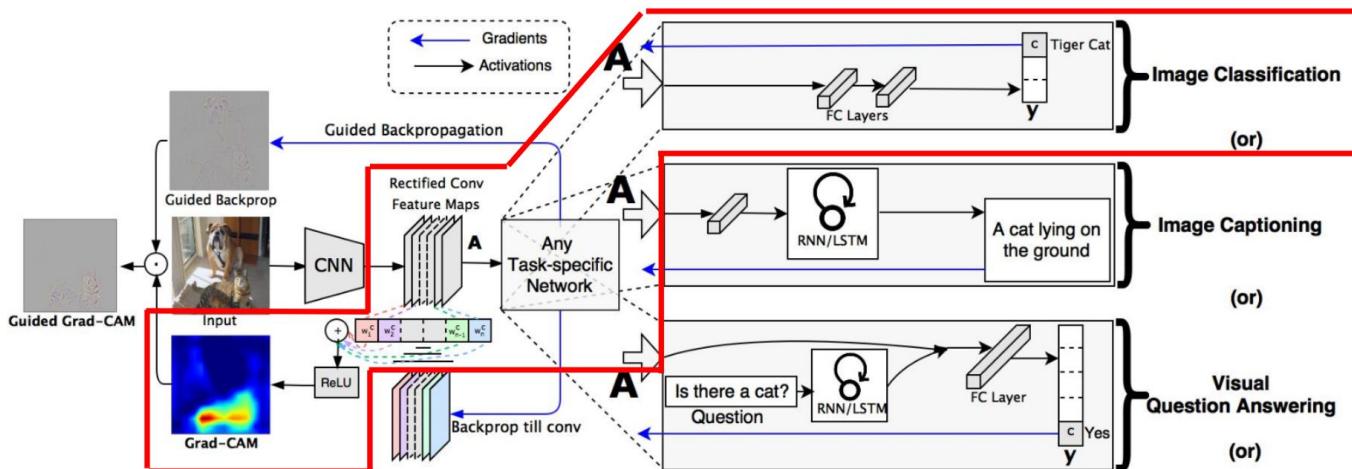


Part 2. Review of WSOL methods: Literature

Taxonomy: Top-Down Grad-based aggregation

- Grad-based aggregation
 - Grad-CAM ('17)
 - Grad-CAM++ ('18)
 - Smooth-Grad-CAM++ ('19)
 - XGrad-CAM ('20)
 - LayerCAM ('21)

Grad-CAM



$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\underbrace{\sum_k \alpha_k^c A^k}_{\text{linear combination}} \right)$$

$$\alpha_k^c = \overbrace{\frac{1}{Z} \sum_i \sum_j}^{\text{global average pooling}} \underbrace{\frac{\partial y^c}{\partial A_{ij}^k}}_{\text{gradients via backprop}}$$

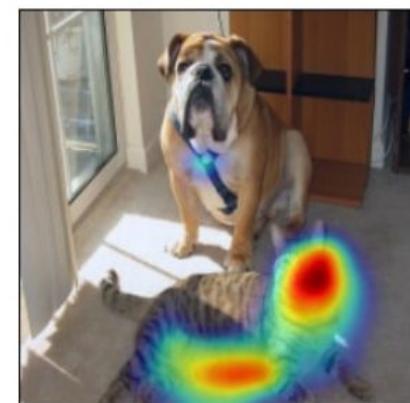
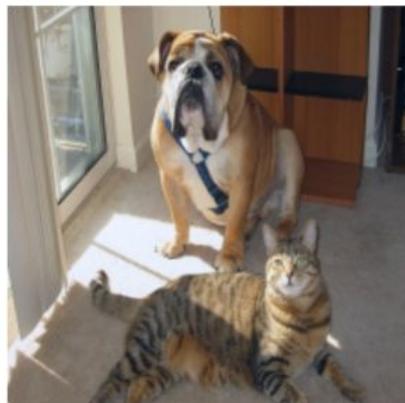
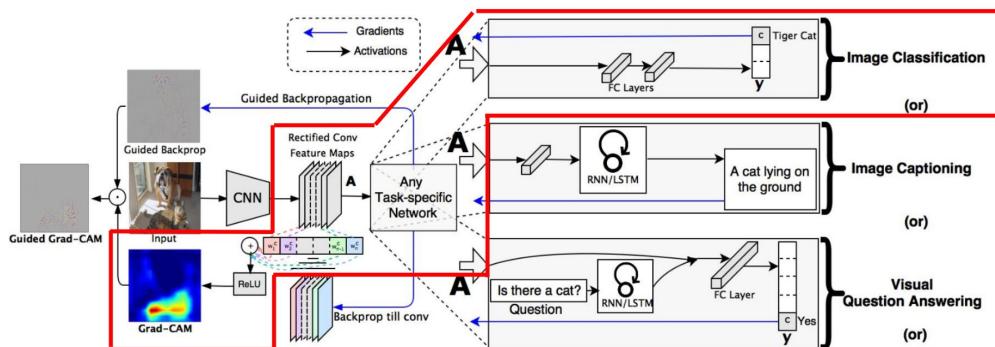
Part 2. Review of WSOL methods: Literature

Taxonomy: Top-Down Grad-based aggregation

Grad-based aggregation

- Grad-CAM ('17)
- Grad-CAM++ ('18)
- Smooth-Grad-CAM++ ('19)
- XGrad-CAM ('20)
- LayerCAM ('21)

Grad-CAM



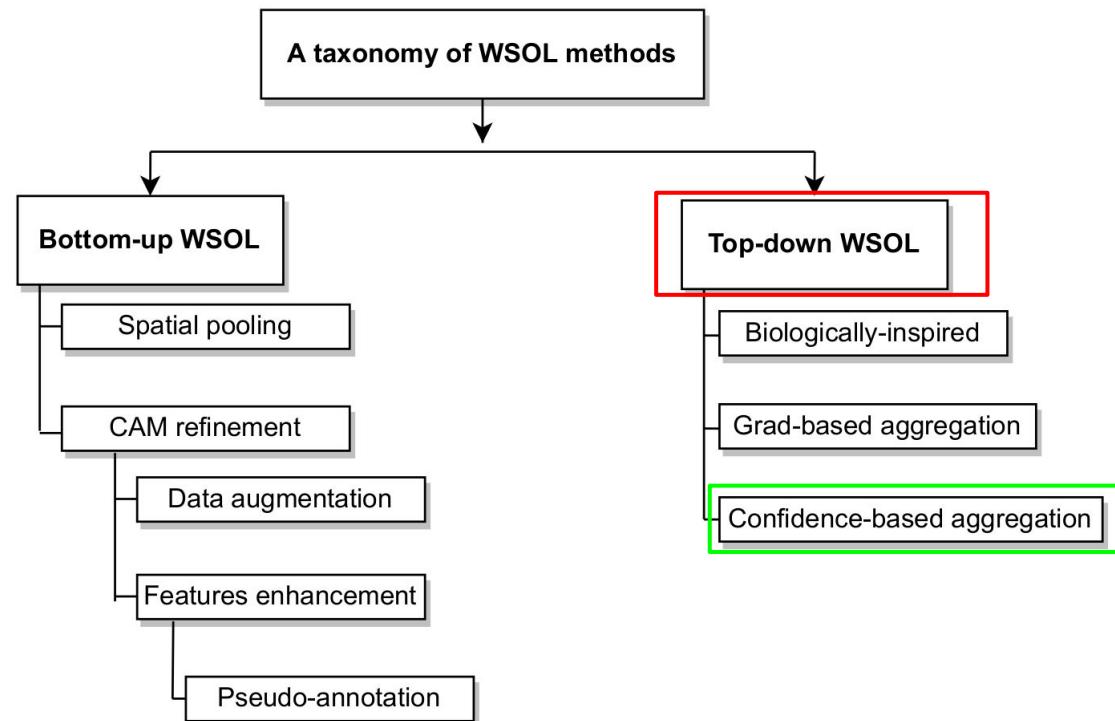
Dog

Cat

Part 2. Review of WSOL methods: Literature

Taxonomy

- Aggregation of feature maps using classification confidence



Part 2. Review of WSOL methods: Literature

Confidence-based aggregation

- Score-CAM ('20)
- SS-CAM ('20)
- IS-CAM ('20)
- Ablation-CAM ('20)

Taxonomy: Top-Down

Confidence-based aggregation

Score-CAM

$$\mathbf{M}_{Score-CAM}^c = \text{ReLU} \left(\sum_k \alpha_k^c \mathbf{F}_k^l \right),$$



In function of classifier response

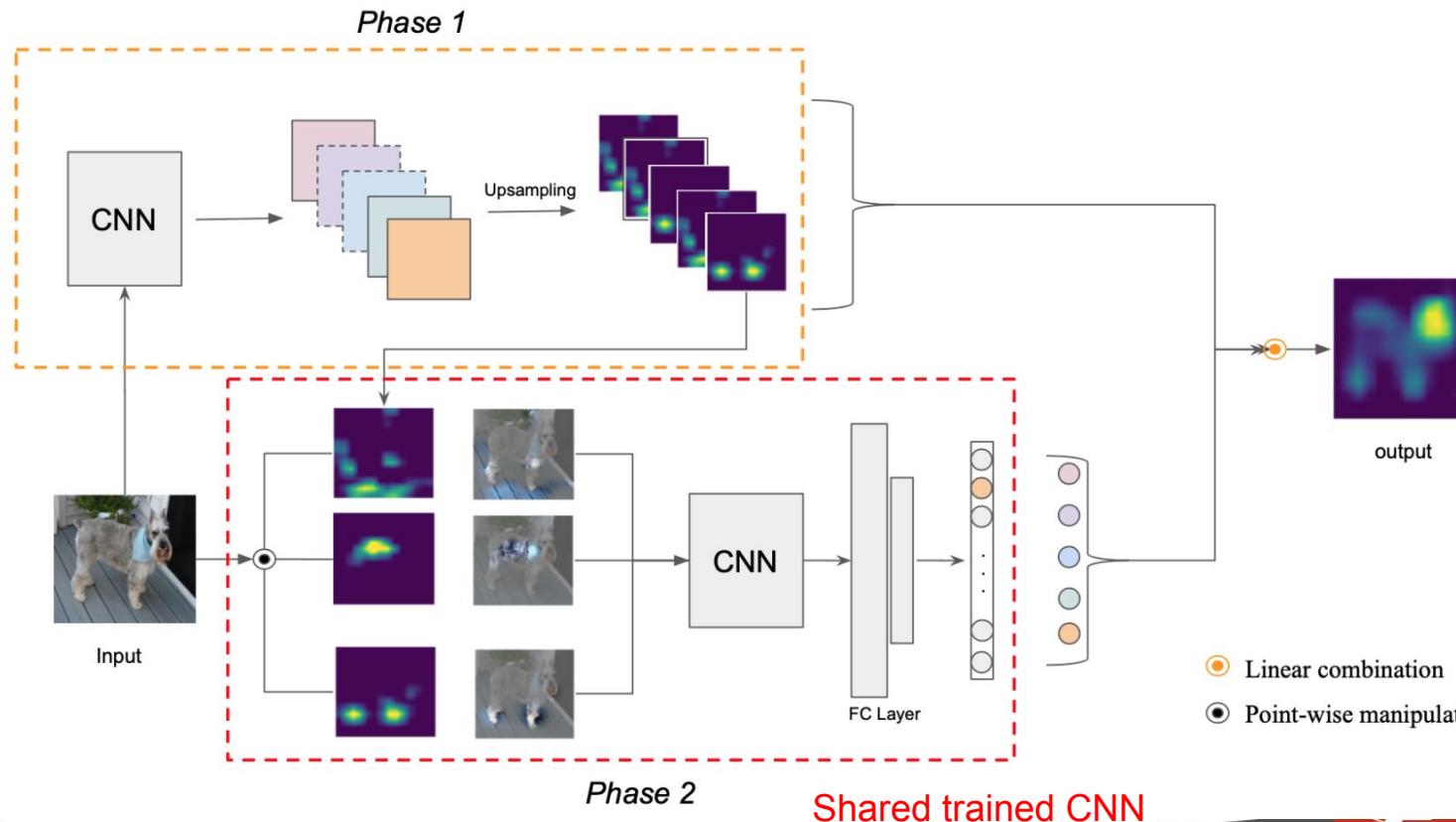
Part 2. Review of WSOL methods: Literature

Confidence-based aggregation

- Score-CAM ('20)
- SS-CAM ('20)
- IS-CAM ('20)
- Ablation-CAM ('20)

Taxonomy: Top-Down Confidence-based aggregation

Score-CAM

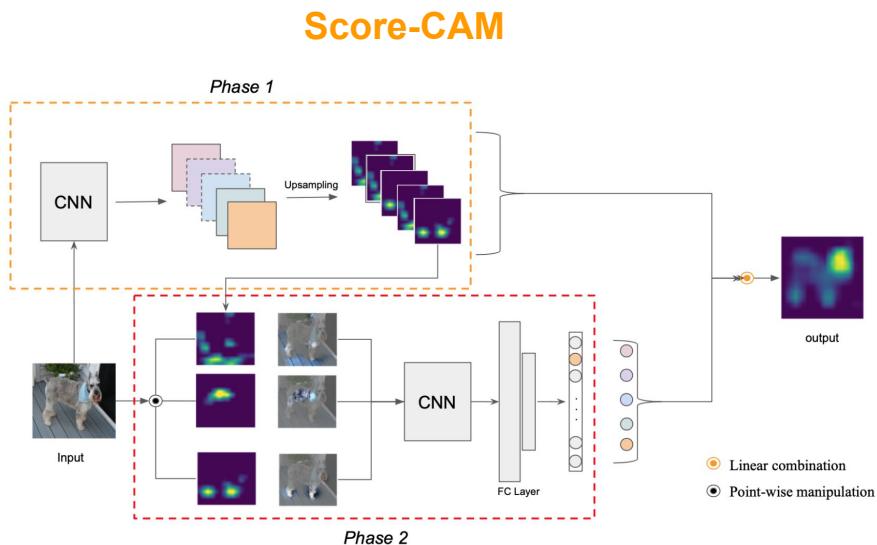


Part 2. Review of WSOL methods: Literature

Confidence-based aggregation

- Score-CAM ('20)
- SS-CAM ('20)
- IS-CAM ('20)
- Ablation-CAM ('20)

Taxonomy: Top-Down Confidence-based aggregation



$$\mathbf{M}_{Score-CAM}^c = \text{ReLU} \left(\sum_k \alpha_k^c \mathbf{F}_k^l \right),$$

$$\alpha_k^c = C(A_l^k)$$

Channel-wise Increase of Confidence (CIC)

$$C(A_l^k) = f(\underline{X \circ H_l^k}) - f(X_b)$$

Perturbed input
image

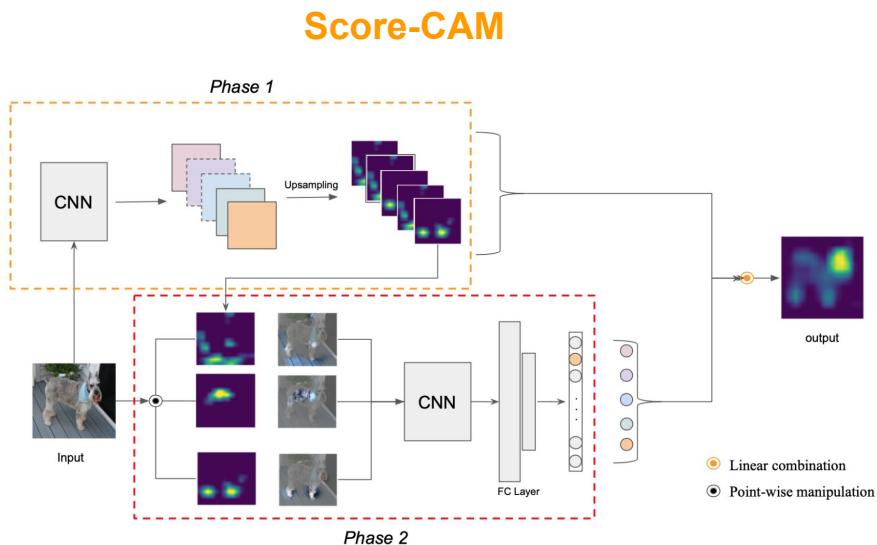
Reference
image

Part 2. Review of WSOL methods: Literature

Confidence-based aggregation

- Score-CAM ('20)
- SS-CAM ('20)
- IS-CAM ('20)
- Ablation-CAM ('20)

Taxonomy: Top-Down Confidence-based aggregation



Normalization:

$$\alpha_k^c = \frac{\exp(C(A_l^k))}{\sum_k \exp(C(A_l^k))}$$

$$\mathbf{M}_{Score-CAM}^c = ReLU \left(\sum_k \alpha_k^c \mathbf{F}_k^l \right),$$

$$\alpha_k^c = C(A_l^k)$$

Channel-wise Increase of Confidence

$$C(A_l^k) = f(\underline{X \circ H_l^k}) - f(X_b)$$

Perturbed input
image

Reference
image

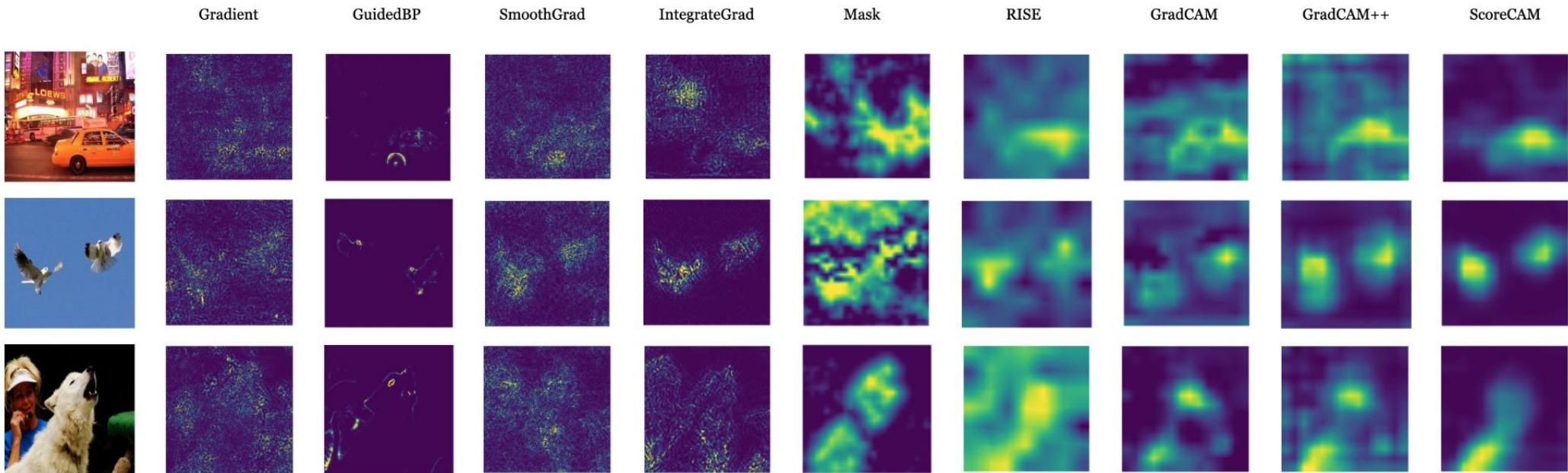
Part 2. Review of WSOL methods: Literature

Confidence-based aggregation

- Score-CAM ('20)
- SS-CAM ('20)
- IS-CAM ('20)
- Ablation-CAM ('20)

Taxonomy: Top-Down Confidence-based aggregation

Score-CAM



Le génie pour l'industrie

Part 2. Review of WSOL methods: Literature

Confidence-based aggregation

- Score-CAM ('20)
- SS-CAM ('20)
- IS-CAM ('20)
- Ablation-CAM ('20)

Taxonomy: Top-Down Confidence-based aggregation

Are these methods practical?

Backbones (encoders) Methods	VGG16				Inception				ResNet50			
	#PCL	#NFM	SFM	#PDEC	#PCL	#NFM	SFM	#PDEC	#PCL	#NFM	SFM	#PDEC
Details	≈19.6	1024	28x28	≈23.1	≈25.6	1024	28x28	≈5.7	≈23.9	2048	28x28	≈9
CAM* [58]		.2ms				.2ms				.3ms		
GradCAM [32]		7.7ms				21.1ms				27.8ms		
GradCAM++ [7]		23.5ms				23.7ms				28.0ms		
Smooth-GradCAM [25]		62.0ms				150.7ms				136.2ms		
XGradCAM [12]		2.9ms				19.2ms				14.2ms		
LayerCAM [15]		3.2ms				18.2ms				17.9ms		
Mean		16.6ms				38.8ms				37.4ms		
ours + STDCL		6.2ms				25.5ms				18.5ms		
ACoL [55]		12.0ms				19.2ms				24.9ms		
SPG [56]		11.0ms				18ms				23.9ms		
ADI [9]		6.4ms				16.0				14.4ms		
ScoreCAM [44]		1.9sec				3.4sec				9.3sec		
SSCAM [24]		1min45sec				2min16sec				5min49sec		
IS-CAM [23]		30.1sec				39.0sec				1min39sec		

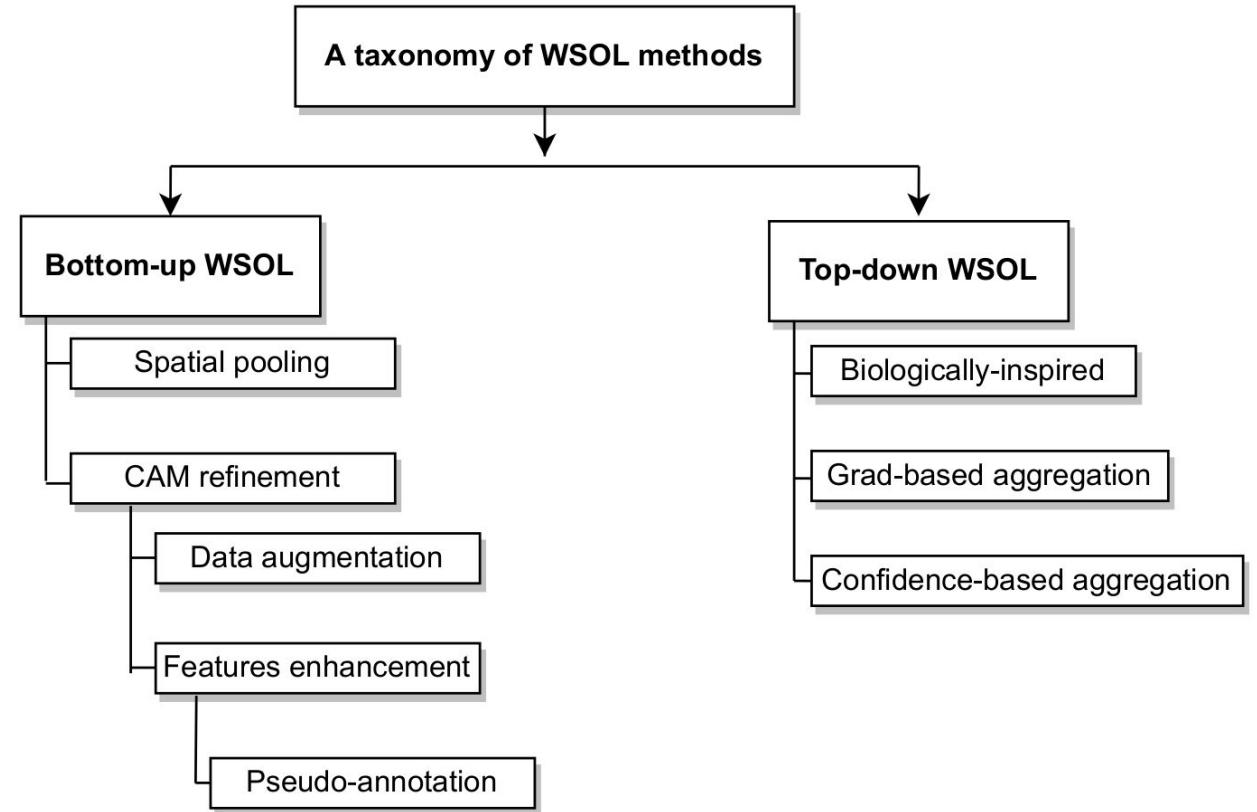
Time necessary to build a full size CAM over an idle **Tesla P100 GPU** for one random RGB image of size **224 × 224** with 200 classes. Methods SSCAM [24] ($N = 35$, $\sigma = 2$), IS-CAM [23] ($N = 10$), IS-CAM [23] ($N = 10$) are evaluated with batch size 32 with their original hyper-parameters (N , and σ).



Le génie pour l'industrie

Part 2. Review of WSOL methods: Literature

Taxonomy



Completed.

Part 2. Review of WSOL methods: Literature

**Taxonomy: Bottom-up
NON-CAM methods**

CSTN: Convolutional STN

Are there NON-CAM WSOL methods?

Part 2. Review of WSOL methods: Literature

**Taxonomy: Bottom-up
NON-CAM methods**

CSTN: Convolutional STN

Are there NON-CAM WSOL methods?

**Yes, but very limited.
why?**

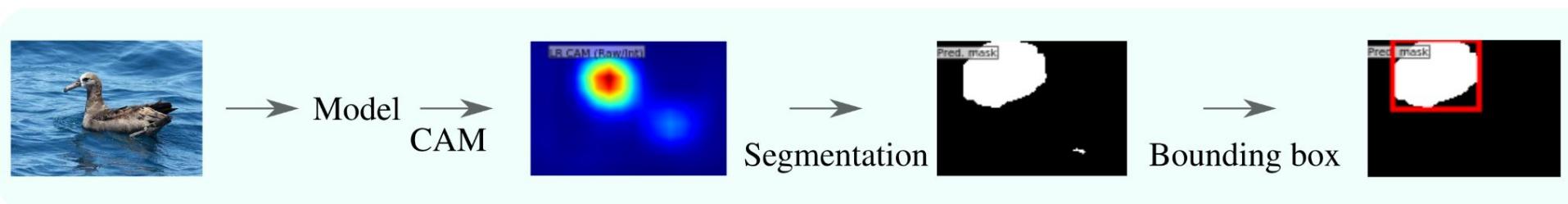
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

CSTN: Convolutional STN

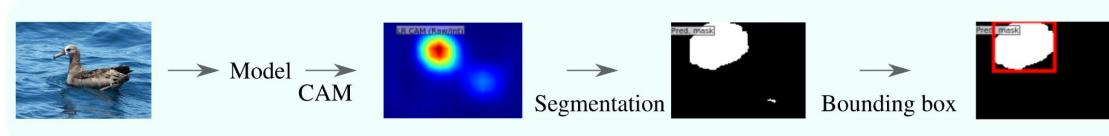
Are there NON-CAM WSOL methods?

Yes, but very limited.
why?



Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods



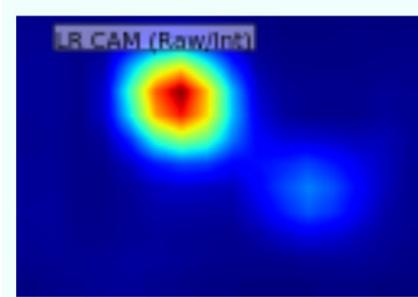
CSTN: Convolutional STN

Are there NON-CAM WSOL methods?

Yes, but very limited.
why?



1: Model



CAM
Soft-segmentation
(Dense localization)

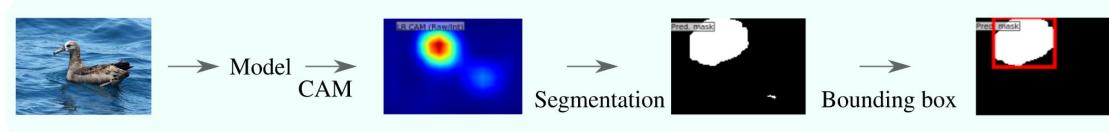


Bounding box
localization
(Coarse localization)

Input, image class

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods



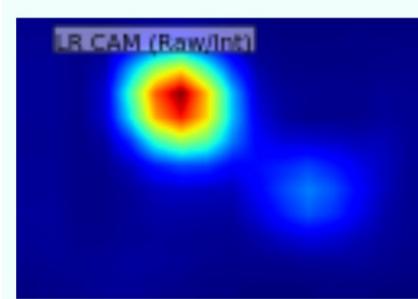
CSTN: Convolutional STN

Are there NON-CAM WSOL methods?

Yes, but very limited.
why?



1: Model



CAM
Soft-segmentation
(Dense localization)

Why not?

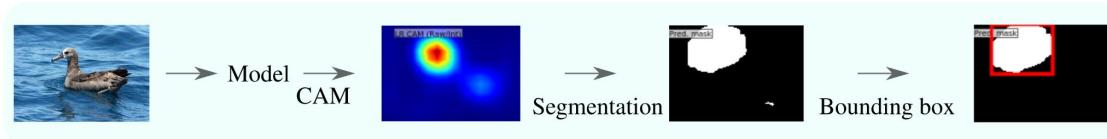


Bounding box
localization
(Coarse localization)

Input, image class

Part 2. Review of WSOL methods: Literature

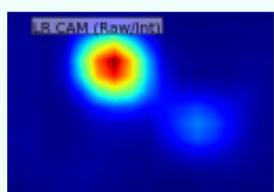
Taxonomy: Bottom-up NON-CAM methods



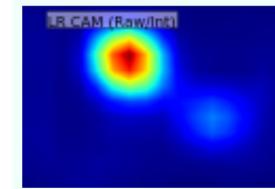
CSTN: Convolutional STN

Are there NON-CAM WSOL methods?

Yes, but very limited.
why?



1: Model



2: Image post-processing
(non-differentiable)

CAM
Soft-segmentation
(Dense localization)



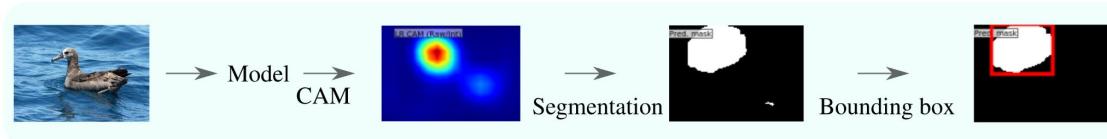
Why not?

Bounding box localization
(Coarse localization)

- **CAM:** emerging property from convolution over visible pixels in image [--> easy]

Part 2. Review of WSOL methods: Literature

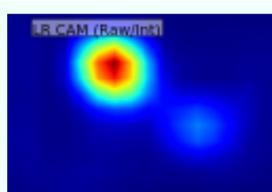
Taxonomy: Bottom-up NON-CAM methods



CSTN: Convolutional STN

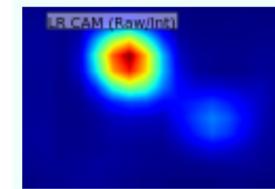
Are there NON-CAM WSOL methods?

Yes, but very limited.
why?



Input, image class

1: Model



CAM
Soft-segmentation
(Dense localization)



2: Image post-processing
(non-differentiable)

Why not?

Bounding box localization
(Coarse localization)

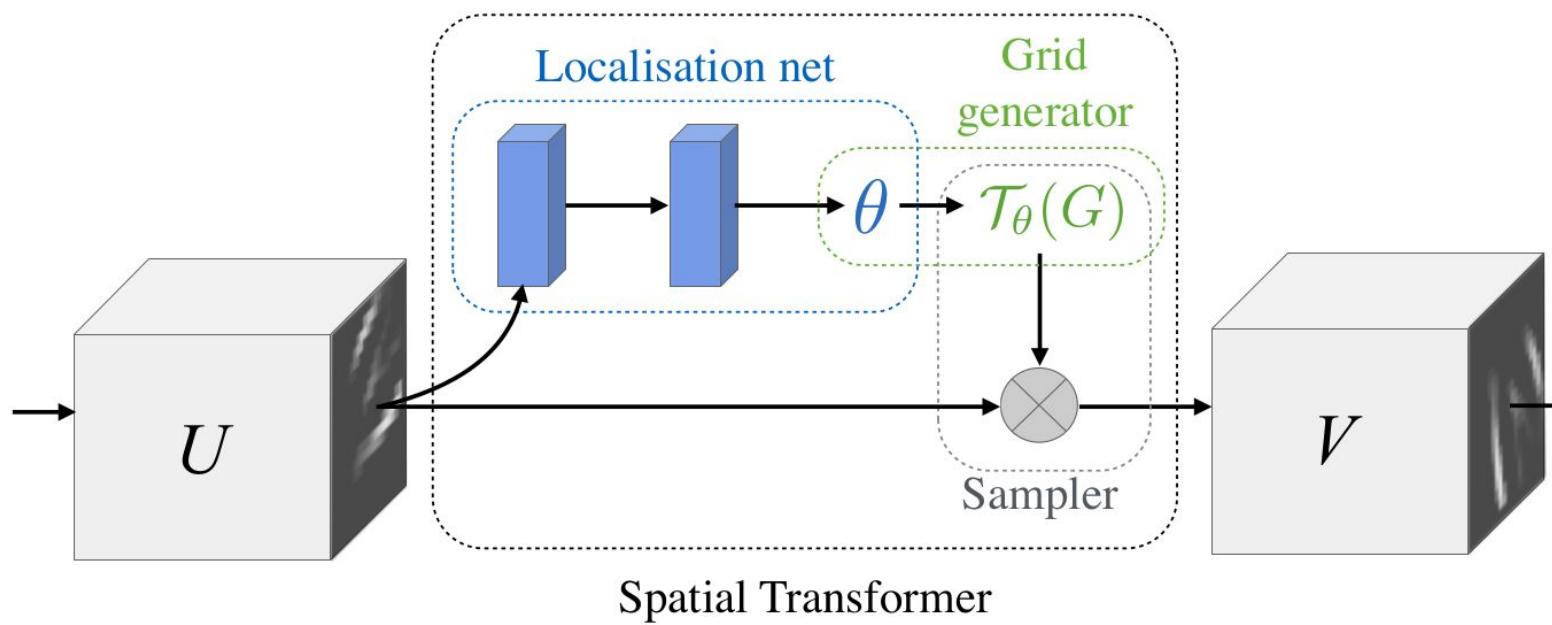
- **CAM:** emerging property from convolution over visible pixels in image [--> easy]
- **Bounding box:** abstract concept (invisible in image) [--> more difficult]



Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

CSTN: Convolutional STN



Spatial Transformer Networks (STN)

M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial transformer networks," in eurIPS, 2015.

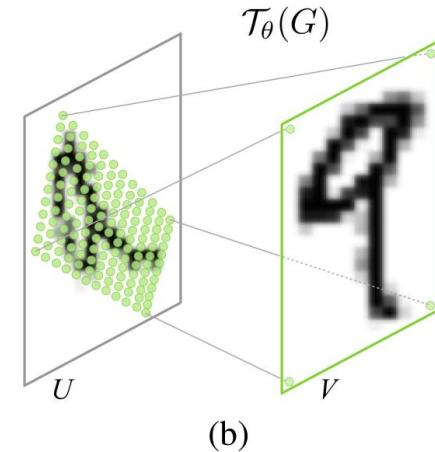
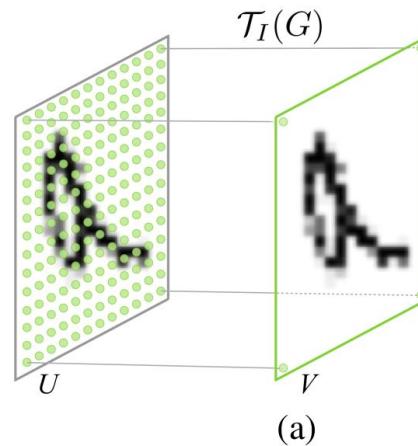
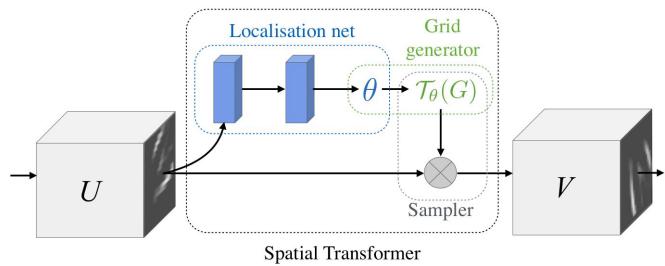
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

Spatial Transformer Networks (STN)

M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial transformer networks," in eurIPS, 2015.

CSTN: Convolutional STN



Model/Invariant to affine transformations:
translation, scale, rotation, ...

Differentiable
'cropping'

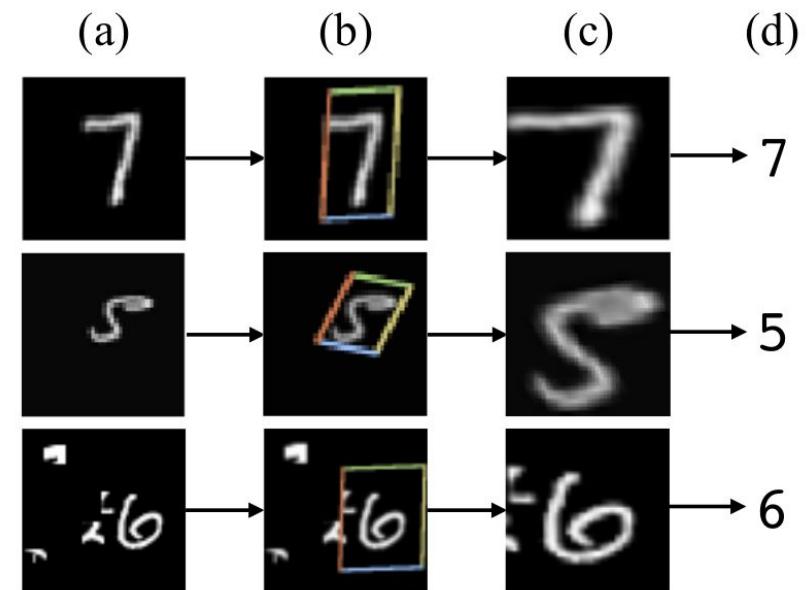
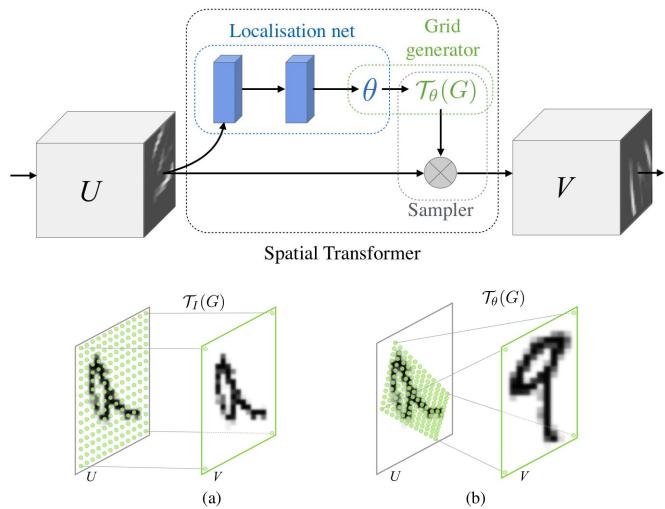
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

Spatial Transformer Networks (STN)

M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial transformer networks," in eurIPS, 2015.

CSTN: Convolutional STN

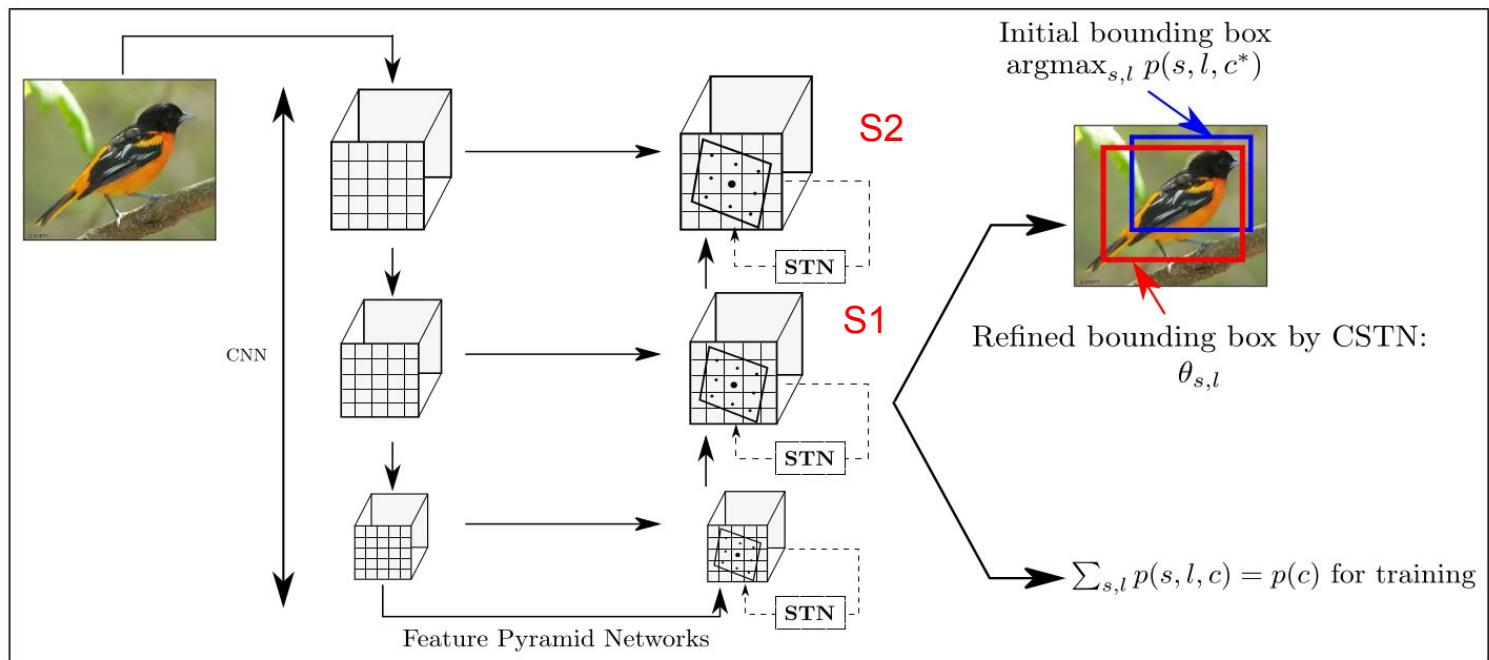


Localization

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

CSTN: Convolutional STN



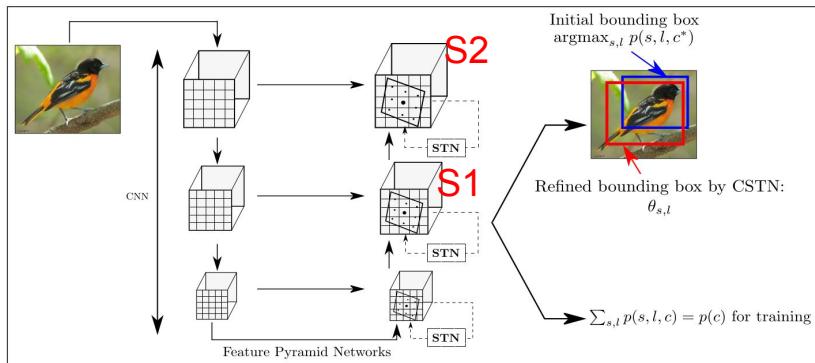
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

Train loss

CSTN: Convolutional STN

$$L(x, y) = L_{cls}(x, y) + \lambda L_\theta + \alpha L_{scale}(x)$$



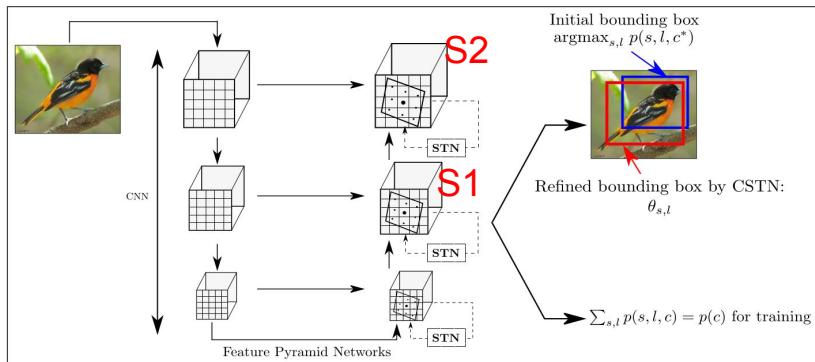
Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

Train loss

CSTN: Convolutional STN

$$L(x, y) = L_{cls}(x, y) + \lambda L_\theta + \alpha L_{scale}(x)$$

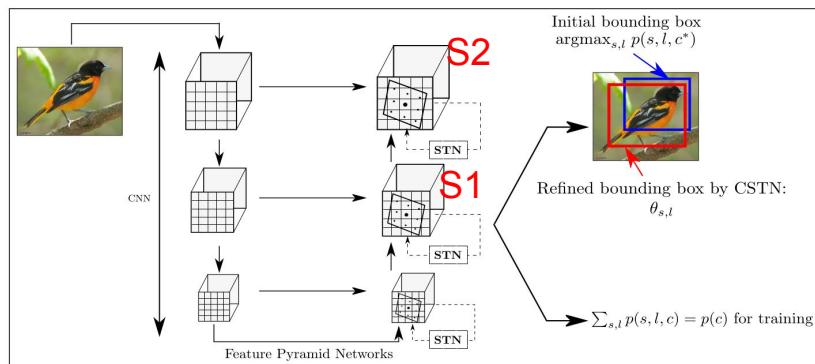


Standard
classification
loss

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

CSTN: Convolutional STN



$$L(x, y) = L_{cls}(x, y) + \lambda L_\theta + \alpha L_{scale}(x)$$

Train loss

Deal with
degenerate
transformations

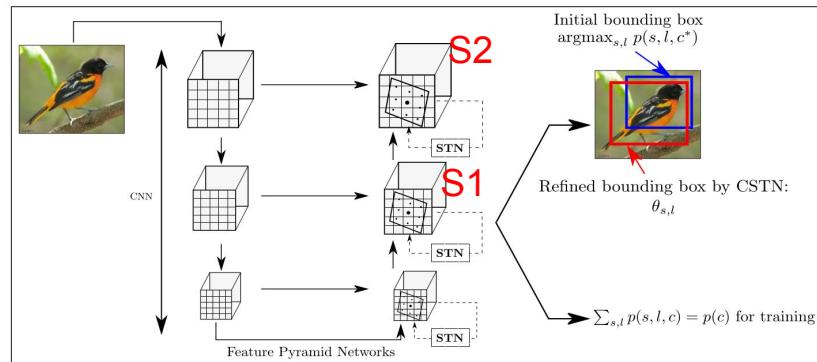
$$L_\theta = \sum_{s \in S} \sum_{i=1}^{h_s \times w_s} \|\theta_{ref} - \theta_i\|^2 .$$

Reference
transformation
(identity)

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

CSTN: Convolutional STN



$$L(x, y) = L_{cls}(x, y) + \lambda L_\theta + \alpha L_{scale}(x)$$

Train loss

Deal with scale issue:
Large objects are selected at low level layers (small part only) → allow top layer to select full object

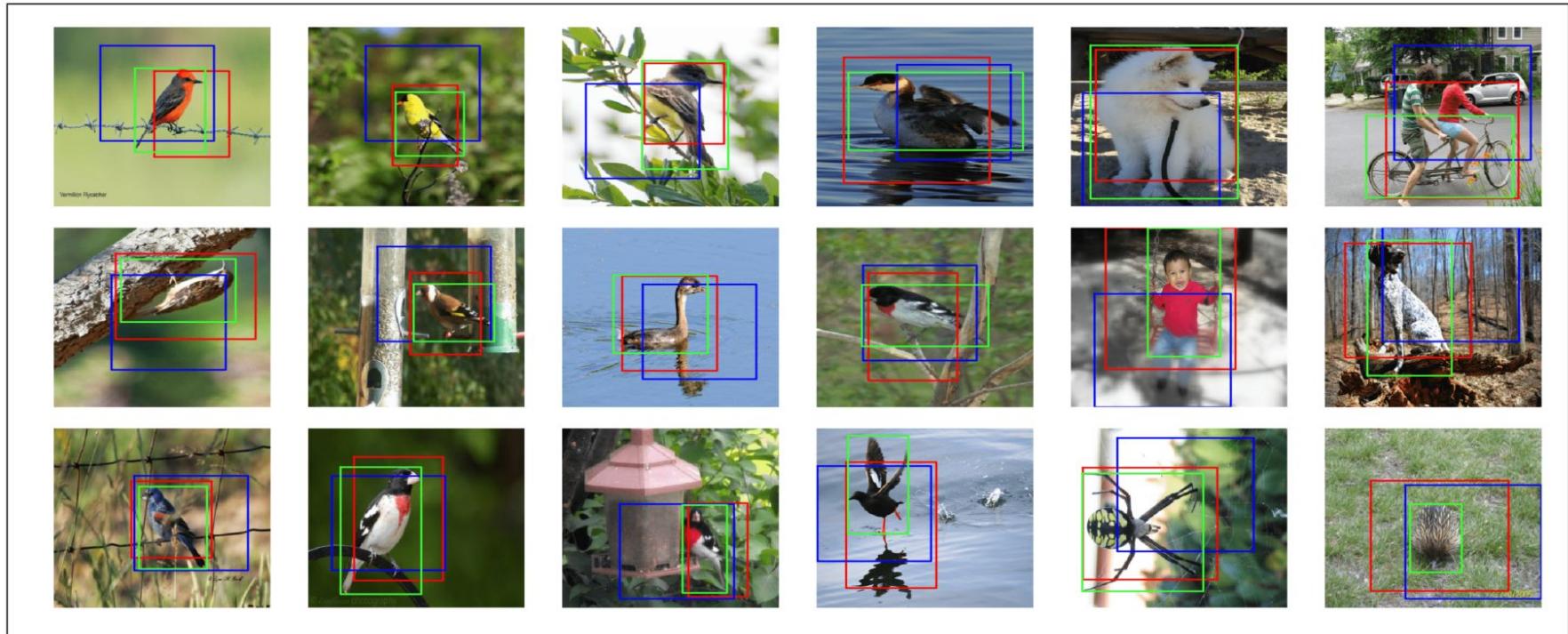
$$L_{scale}(x) = \max \left(0, \max_l p(s = s_1, l, c = c^* | x) - \max_l (p(s = s_2, l, c = c^* | x)) \right)$$

Part 2. Review of WSOL methods: Literature

Taxonomy: Bottom-up NON-CAM methods

Green: ground truth
Blue: without STN
Red: with STN

CSTN: Convolutional STN



Part 2

Review of WSOL Methods

- WSOL literature: bottom-up and top-down methods
- Case studies:
 - a) F-CAM for improved interpolation
 - b) Transformer-based models

Case Study (a): F-CAM for Improved Interpolation

F-CAM: Full Resolution Class Activation Maps via Guided Parametric Upscaling

Soufiane Belharbi, Aydin Sarraf, Marco Pedersoli, Ismail Ben Ayed,
Luke McCaffrey, Eric Granger

WACV 2022: Winter Conf. on Applications of Computer Vision



LIVIA LABORATOIRE
D'IMAGERIE, DE VISION
ET D'INTELLIGENCE
ARTIFICIELLE



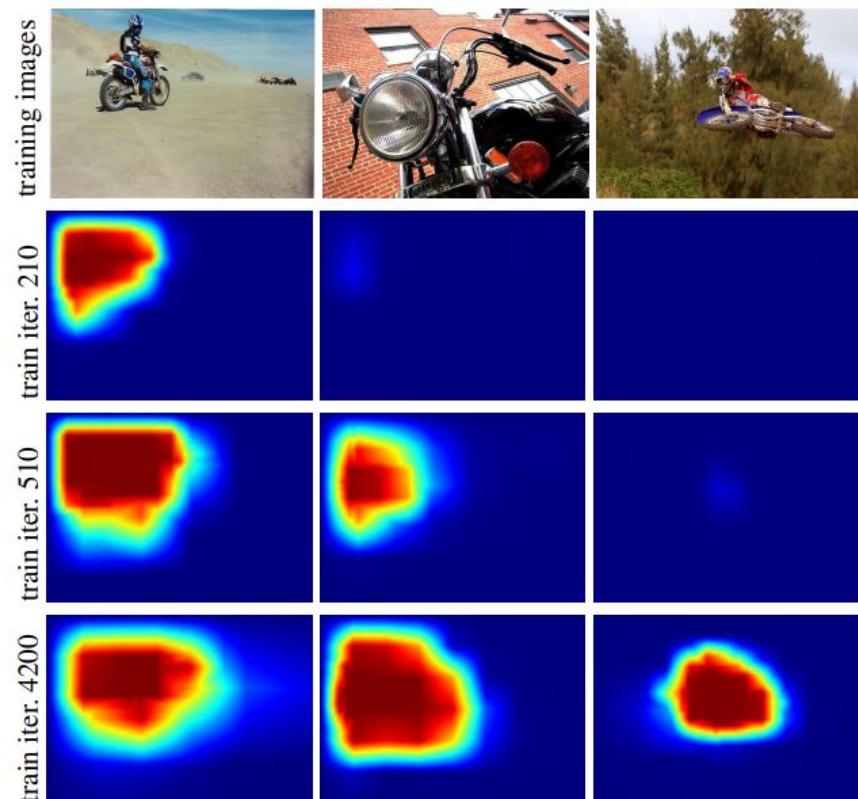
Case Study (a): F-CAM for Improved Interpolation

- **A Challenge with CAMs:** low resolution (due to convolution and pooling) has negative impact on localization performance

Standard interpolated from CAM of 8x8 resolution (downscale factor of 32)

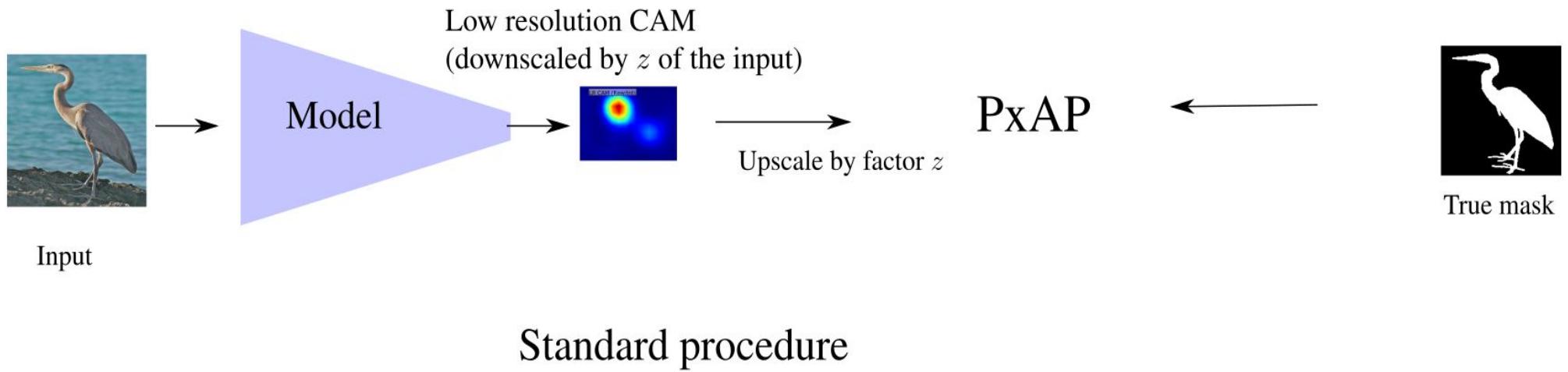


Standard interpolated CAMs



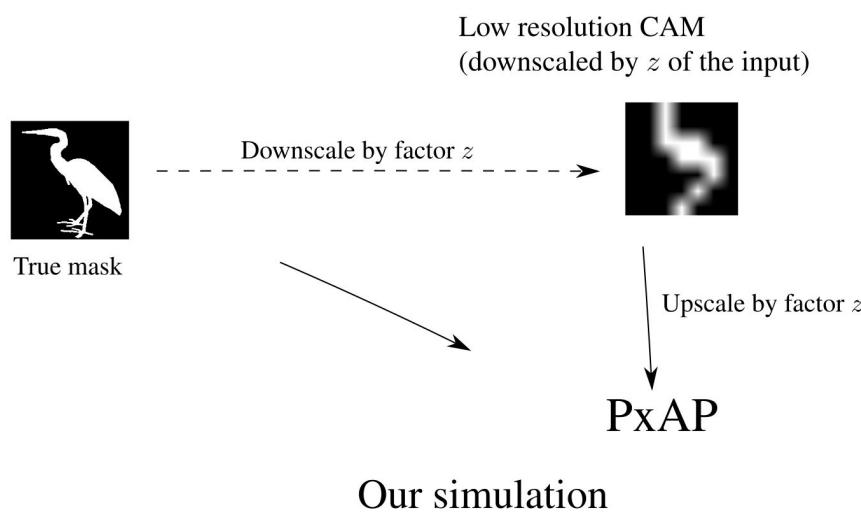
Case Study (a): F-CAM for Improved Interpolation

- **Challenges:** Evaluating the impact on localization performance of CAM size (CUB dataset)
 - the CNN produces low-resolution CAMs that are interpolated (by a factor z) to return to the input image size
 - PxAP measures localization accuracy: it evaluates interpolated images against the true segmentation masks

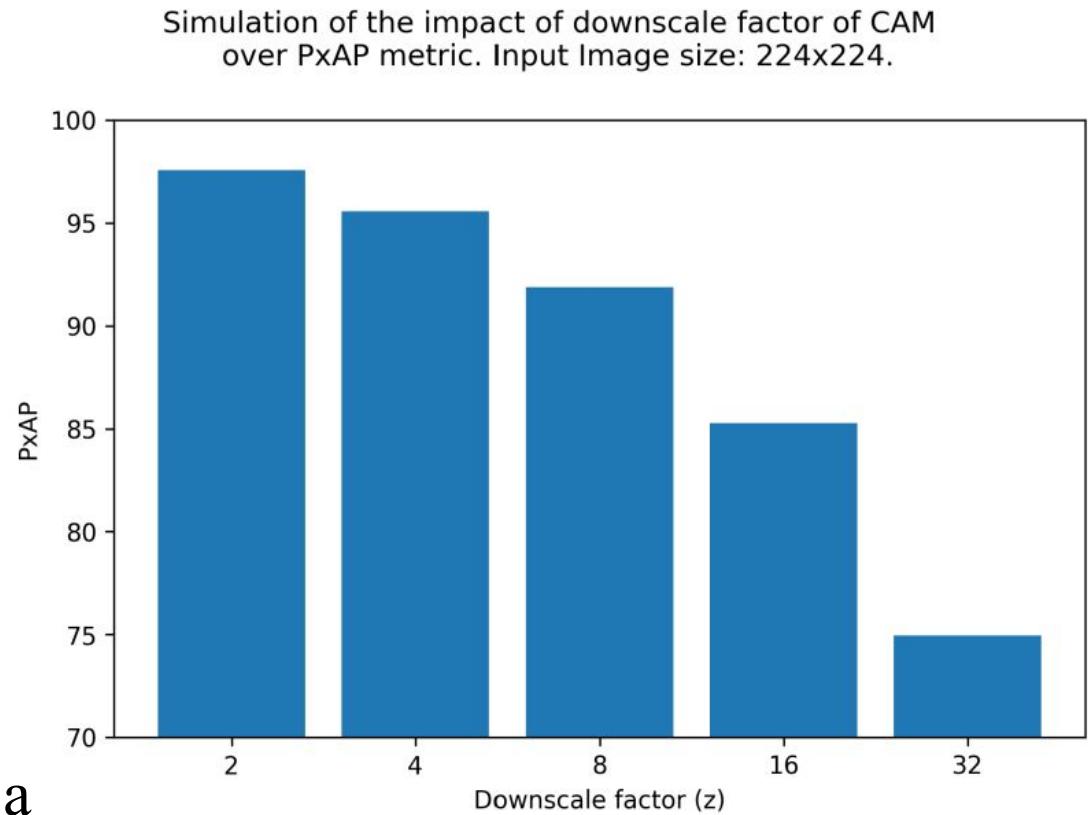


Case Study (a): F-CAM for Improved Interpolation

- **Challenges:** Impact of CAMs size on localization performance



Results: increasing the downscaling factor (z) leads to a considerable decline in localization accuracy



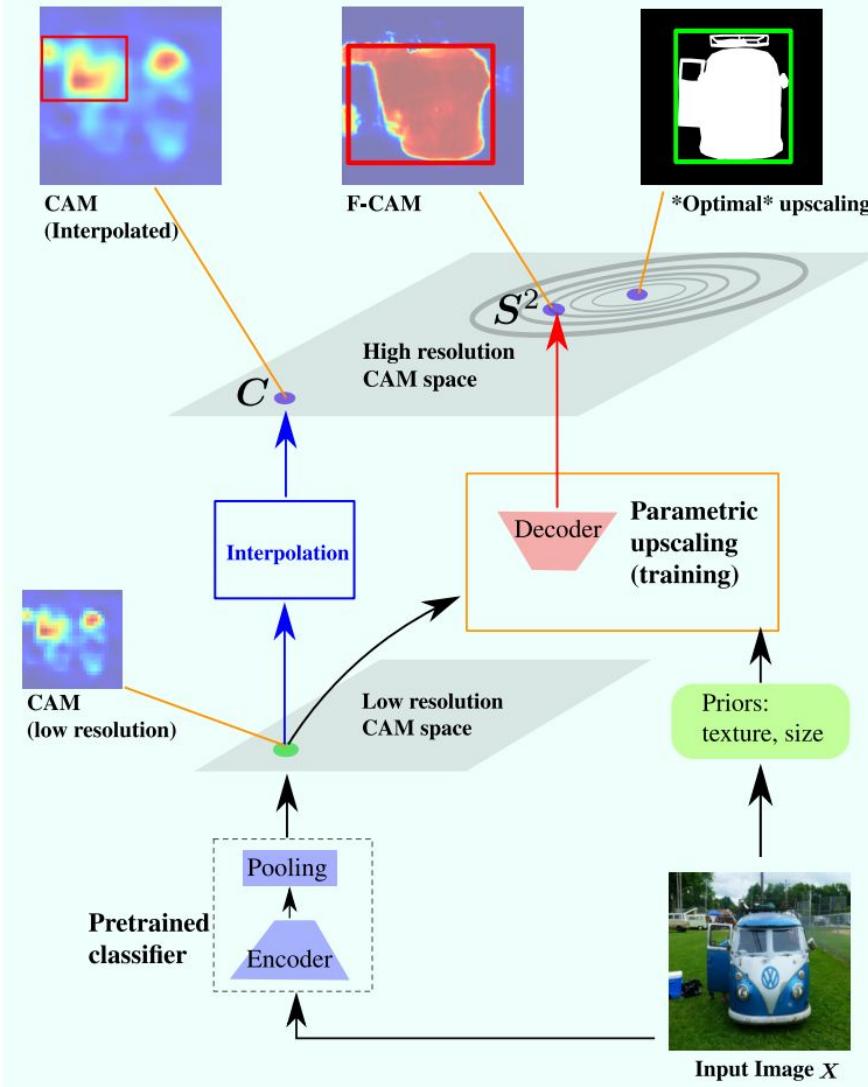
Case Study (a): F-CAM for Improved Interpolation

- **Interpolation:** common, but it does not consider statistical properties of an object, such as color and texture or its shape
- **Litterature:** alternatives to avoid bi-cubic interpolation for producing higher resolution CAMs - *learnable upscaling*
 - residual dilation networks [Yu, CVPR 2017]
 - an end-to-end weakly supervised semantic segmentation approach that upscaled the feature maps [Zhang, AAAI 2020]
 - U-Net architecture to reconstruct the image [Tagaris. ICIP 219]

These methods either produce smaller CAMs, are difficult to scale to large number of classes, or require costly post-processing

Case Study (a): F-CAM for Improved Interpolation

- Proposed F-CAM with Guided Parametric Upscaling



Encoder: any pre-trained CNN classifier,

L_c = classification loss (supervised)

Decoder: trained to perform parametric upscaling

L_D = pixel alignment loss (unsupervised)

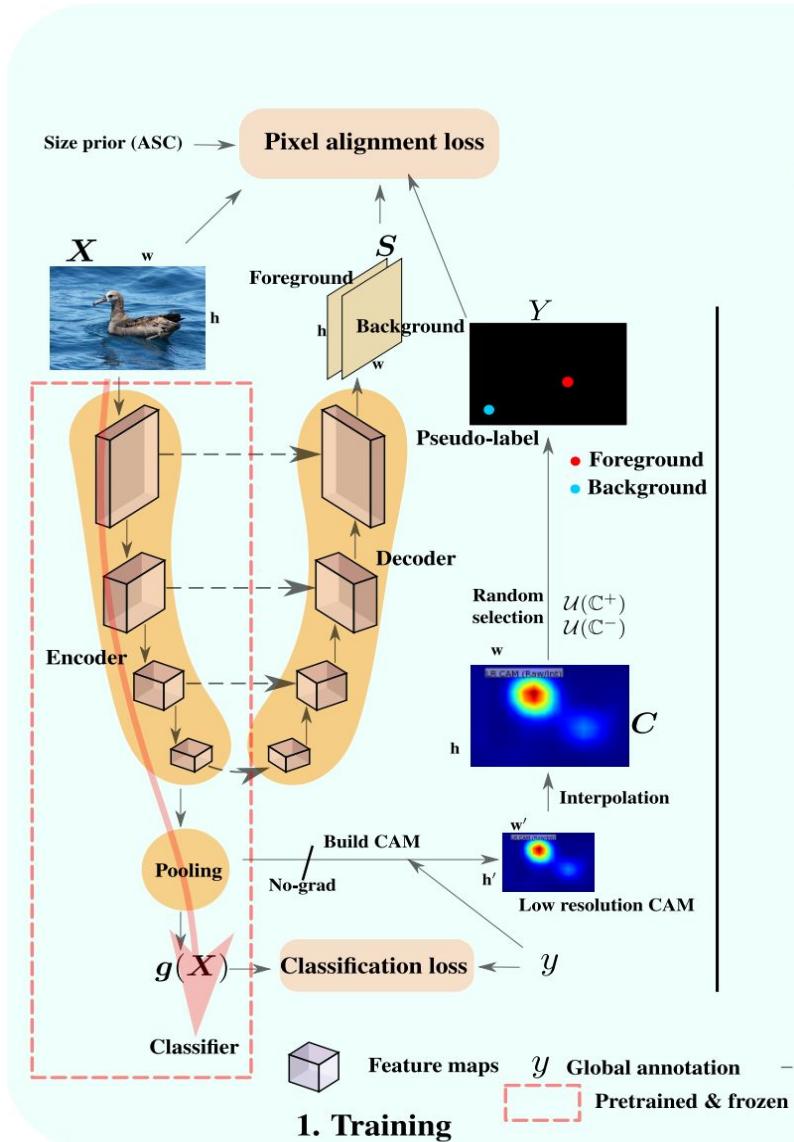
$$= \text{SR (CAM)} + \text{CRF (image)} + \text{ASC (size)}$$

where

- SR: seeds (positive/negative evidence at pixel level)
- CRF: image properties
- ASC: unsupervised size constraint

Case Study (a): F-CAM for Improved Interpolation

- Proposed F-CAM: training models the foreground and background



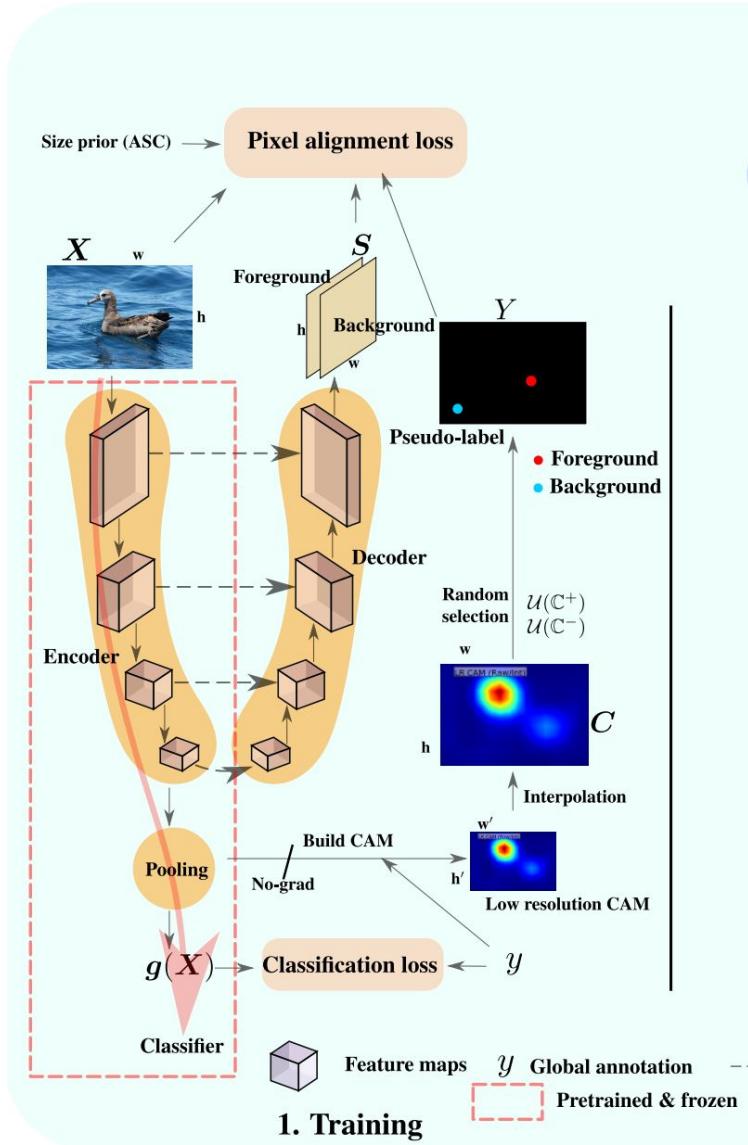
Overall loss for end-to-end training

$$\begin{aligned}
 L_c &= L_{CE} \\
 L_D &= L_{SR} + L_{CRF} \\
 \min_{\theta} \quad & -\log(g(X)[y]) + \alpha \sum_{p \in \Omega'} H(Y_p, S_p) + \lambda \mathcal{R}(S, X), \\
 \text{s.t.} \quad & \sum S^r \geq 0, \quad r \in \{1, 2\},
 \end{aligned}$$

ASC: area size constraint

Case Study (a): F-CAM for Improved Interpolation

- Proposed F-CAM: training models the foreground and background



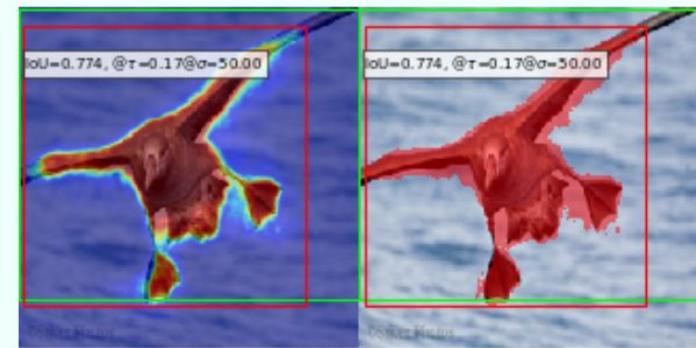
L_{SR} : Random sampling from *noisy evidence*

- use stochastic sampling (extreme dropout) to avoid overfitting wrong labels,
- give the model enough time to allow consistent ROI to emerge
- at each SGD step, select 2 new random pixels (foreground and background)

Case Study (a): F-CAM for Improved Interpolation

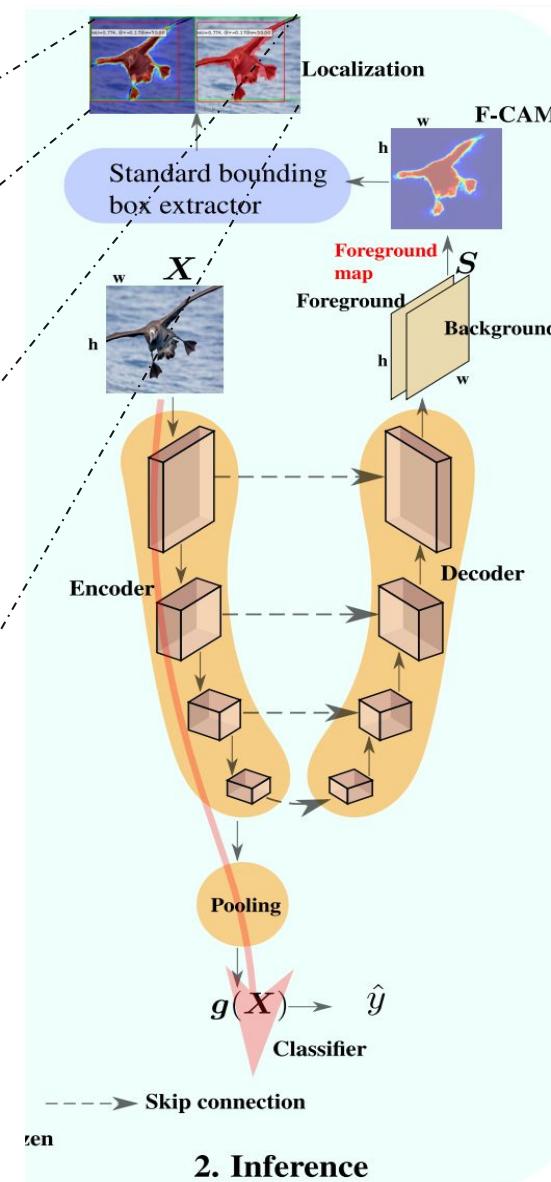
- Proposed F-CAM:

At inference time



F-CAM

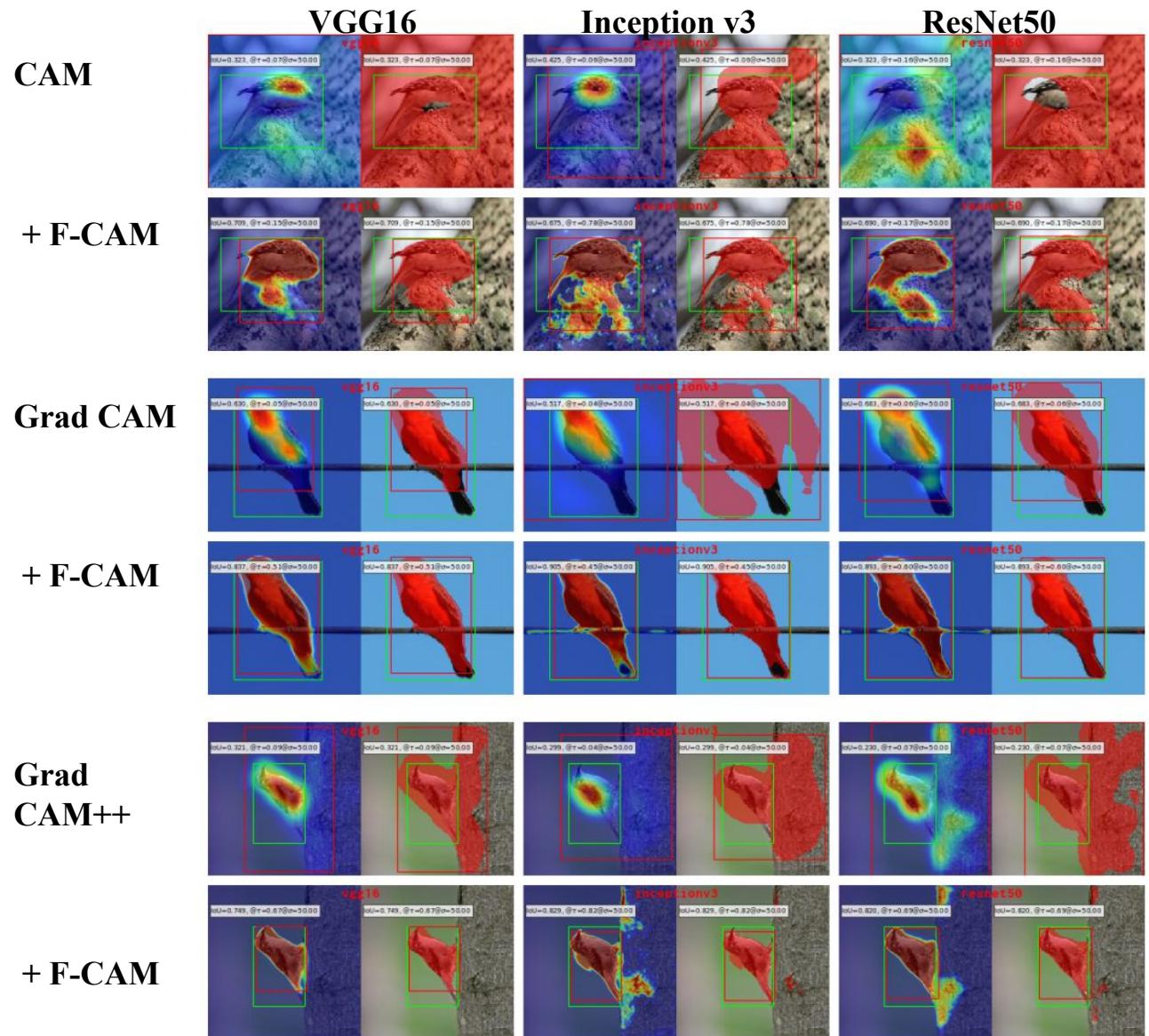
Mask after
thresholding



Case Study (a): F-CAM for Improved Interpolation

- Experiments:

Visual results on images from the CUB dataset



Case Study (a): F-CAM for Improved Interpolation

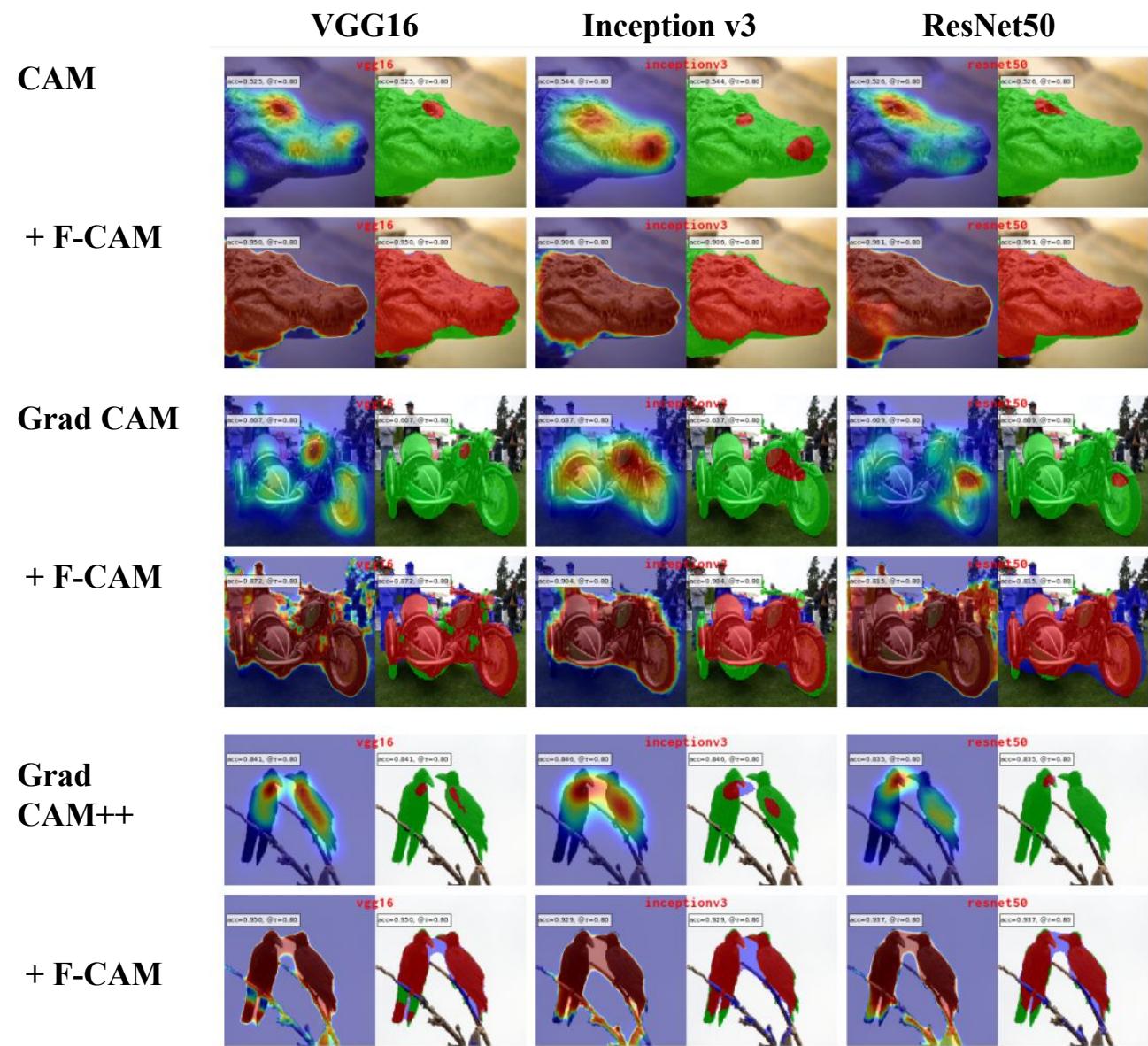
- Experiments:

Visual results on images from the OpenImages dataset

Compare CAMs to segmentation masks

Pixel classification

- Red: true positive pixels
- Blue: false positive pixels
- Green: false negative pixels



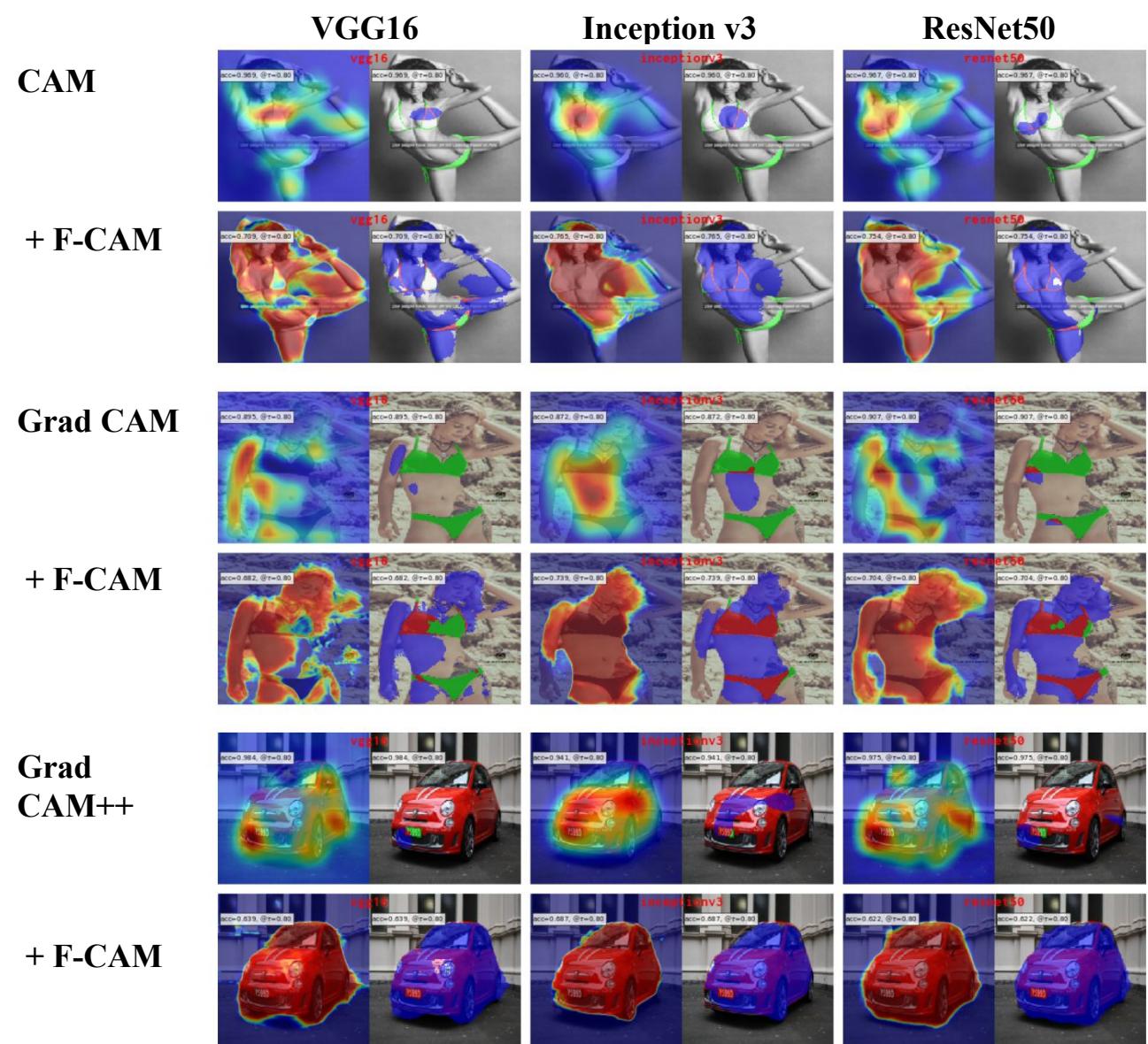
Case Study (a): F-CAM for Improved Interpolation

- **Experiments:**

Visual results on images from the OpenImages dataset

Failure cases

- due to sampling errors
- performance is tied to quality of CAM from pre-trained CNN classifier



Case Study (a): F-CAM for Improved Interpolation

- Experiments:

Localization accuracy
on the CUB and
OpenImages datasets

3 CNNs

Performance measures:

- MaxBoxAcc (CUB)
- PxAP (OpenImages)

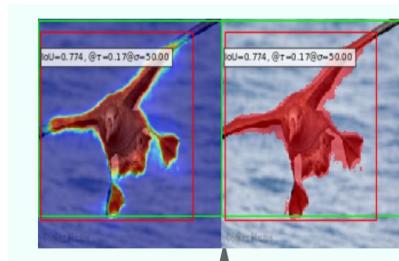
Methods	CUB (MaxBoxAcc)				OpenImages (PxAP)			
	VGG	Inception	ResNet	Mean	VGG	Inception	ResNet	Mean
CAM [57] (cvpr,2016)	71.1	62.1	73.2	68.8	58.1	61.4	58.0	59.1
HaS [34] (iccv,2017)	76.3	57.7	78.1	70.7	56.9	59.5	58.2	57.8
ACoL [53] (cvpr,2018)	72.3	59.6	72.7	68.2	54.7	63.0	57.8	58.4
SPG [54] (eccv,2018)	63.7	62.8	71.4	66.0	55.9	62.4	57.7	58.6
ADL [9] (cvpr,2019)	75.7	63.4	73.5	70.8	58.3	62.1	54.3	58.2
CutMix [51] (eccv,2019)	71.9	65.5	67.8	68.4	58.2	61.7	58.7	59.5
Best WSOL	76.3	65.5	78.1	70.8	58.3	63.0	58.7	59.5
FSL baseline	86.3	94.0	95.8	92.0	61.5	70.3	74.4	68.7
Center baseline	59.7	59.7	59.7	59.7	45.8	45.8	45.8	45.8
CSTN [22] (icpr,2020)	Resnet101 [14]: 76.0				–	–	–	–
TS-CAM [13] (corr,2021)	Deit-S [39]: 83.8				–	–	–	–
MEIL [21] (cvpr,2020)	73.8	–	–	–	–	–	–	–
DANet [47] (iccv,2019)	67.7	67.03	–	–	–	–	–	–
SPOL [44] (cvpr,2021)	–	–	96.4	–	–	–	–	–
CAM* [57] (cvpr,2016)	61.6	58.8	71.5	63.9	53.0	62.7	56.8	57.5
GradCAM [32] (iccv,2017)	69.3	62.3	73.1	68.2	59.6	63.9	60.1	61.2
GradCAM++ [7] (wacv,2018)	84.1	63.3	81.9	76.4	60.5	64.0	60.2	61.5
Smooth-GradCAM++ [25] (corr,2019)	69.7	66.9	76.3	70.9	52.2	61.7	54.3	56.0
XGradCAM [12] (bmvc,2020)	69.3	60.9	72.7	67.6	59.0	63.9	60.2	61.0
LayerCAM [15] (ieee,2021)	84.3	66.5	85.2	78.6	59.5	63.5	61.1	61.3
CAM* [57] + ours	87.3	82.0	90.3	86.5	67.8	71.9	72.1	70.6
GradCAM [32] + ours	87.5	84.4	90.5	87.4	68.6	70.0	70.9	69.8
GradCAM++ [57] + ours	91.5	84.6	91.0	89.0	64.8	67.1	66.3	66.0
Smooth-GradCAM++ [57] + ours	89.1	86.8	90.7	88.8	60.3	65.4	64.4	63.3
XGradCAM [57] + ours	86.8	84.4	90.4	88.8	68.7	71.3	70.4	70.1
LayerCAM [57] + ours	91.0	85.3	92.4	89.7	64.3	64.9	65.3	64.8
Best WSOL + ours	91.5	86.8	92.4	89.7	68.7	71.9	72.1	70.6

Table 1: Performance on MaxBoxAcc and PxAP metrics.

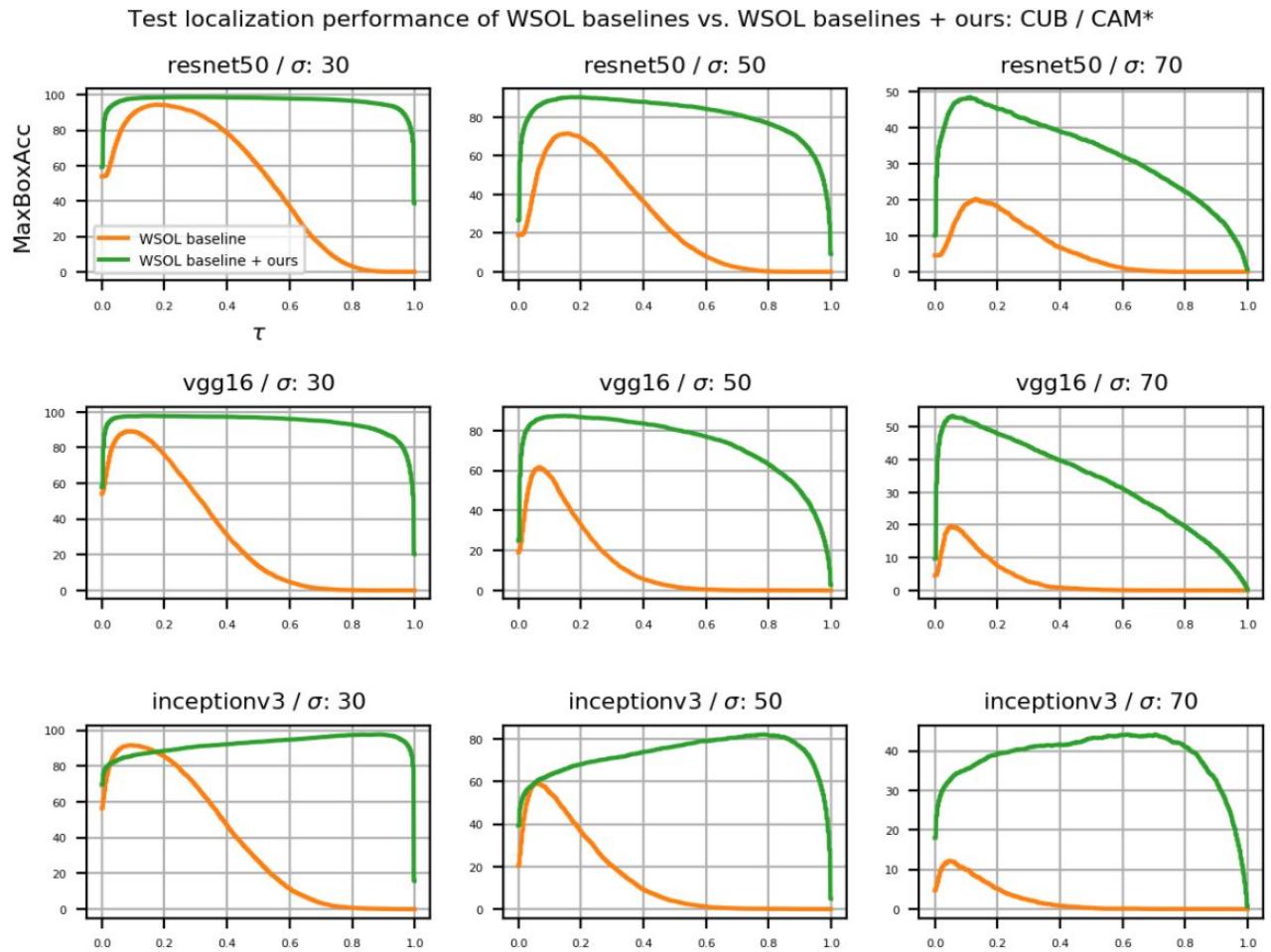
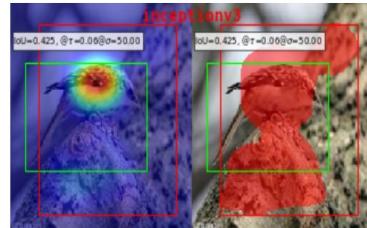
Case Study (a): F-CAM for Improved Interpolation

- Experiments:
sensitivity to
threshold values
on the CUB
dataset

CAM + ours



CAM

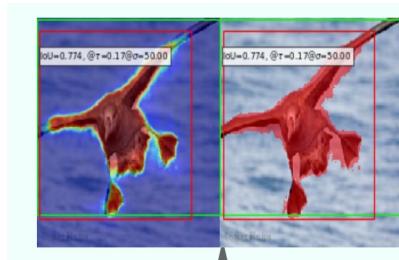


(a) MaxBoxAcc metric.

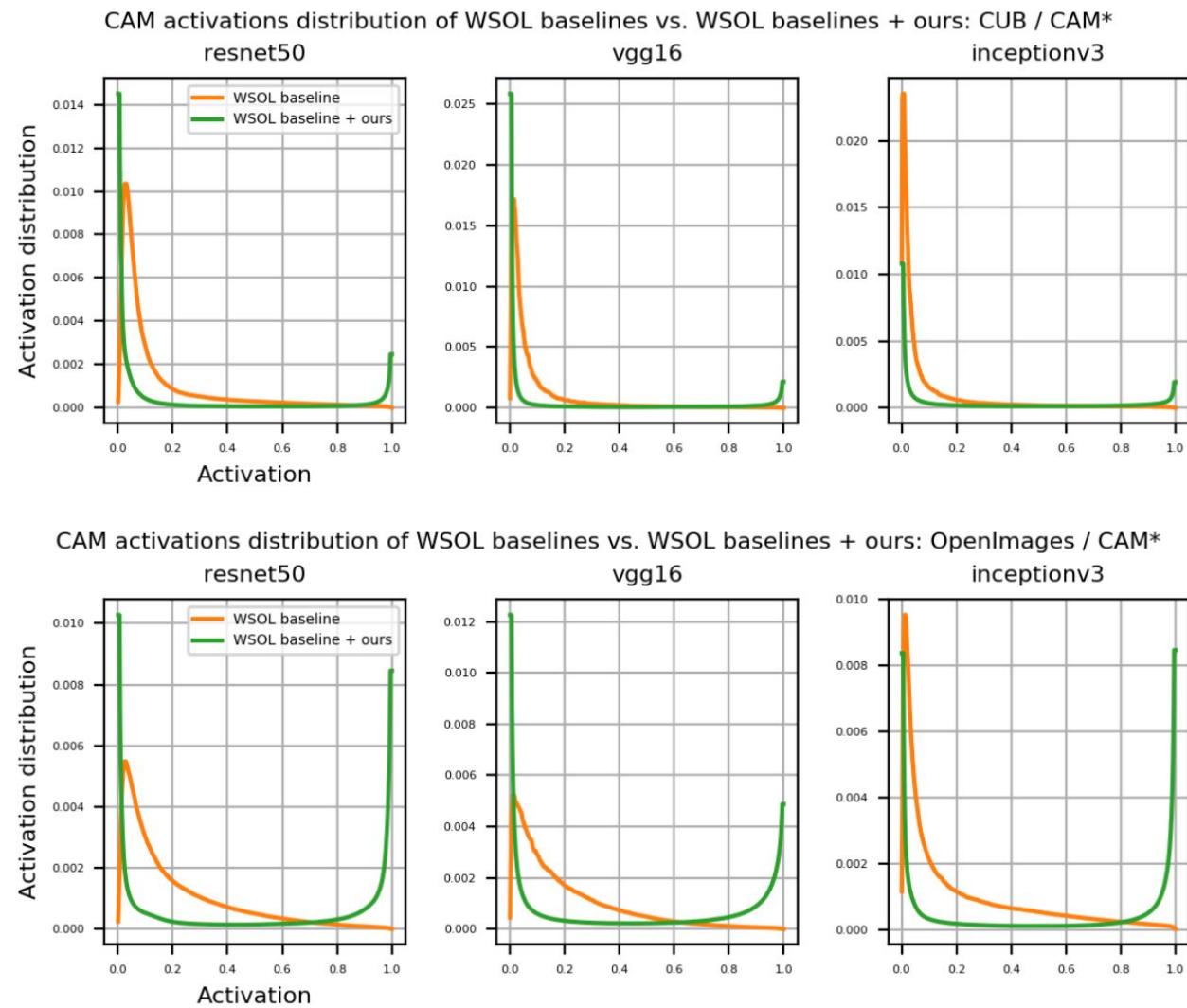
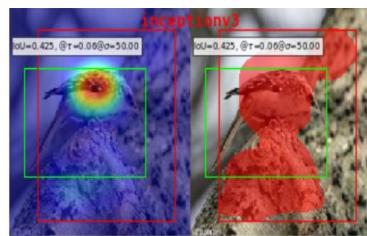
Case Study (a): F-CAM for Improved Interpolation

- **Experiments:**
sensitivity to
threshold values
on the CUB
dataset

CAM + ours



CAM



Case Study (a): F-CAM for Improved Interpolation

- **Experiments:** Ablation study on the impact on loss components

Methods	CUB (MaxBoxAcc)				OpenImages (PxAP)			
	VGG	Inception	ResNet	Mean	VGG	Inception	ResNet	Mean
CAM* [58]	61.6	58.8	71.5	63.9	53.0	62.7	56.8	57.5
CAM* [58] + SR	84.2	73.0	82.2	79.8	64.5	64.1	63.8	64.1
CAM* [58] + SR + ASC	82.9	74.1	83.2	80.0	63.9	63.4	62.0	63.1
CAM* [58] + SR + CRF	84.6	78.9	86.1	83.2	66.3	68.3	67.5	67.3
CAM* [58] + SR + CRF + ASC	87.3	82.0	90.3	86.5	67.8	71.9	72.1	70.6
Improvement	+25.7	+23.2	+18.8	+22.5	+14.8	+9.2	+15.3	+12.8

Case Study (a): F-CAM for Improved Interpolation

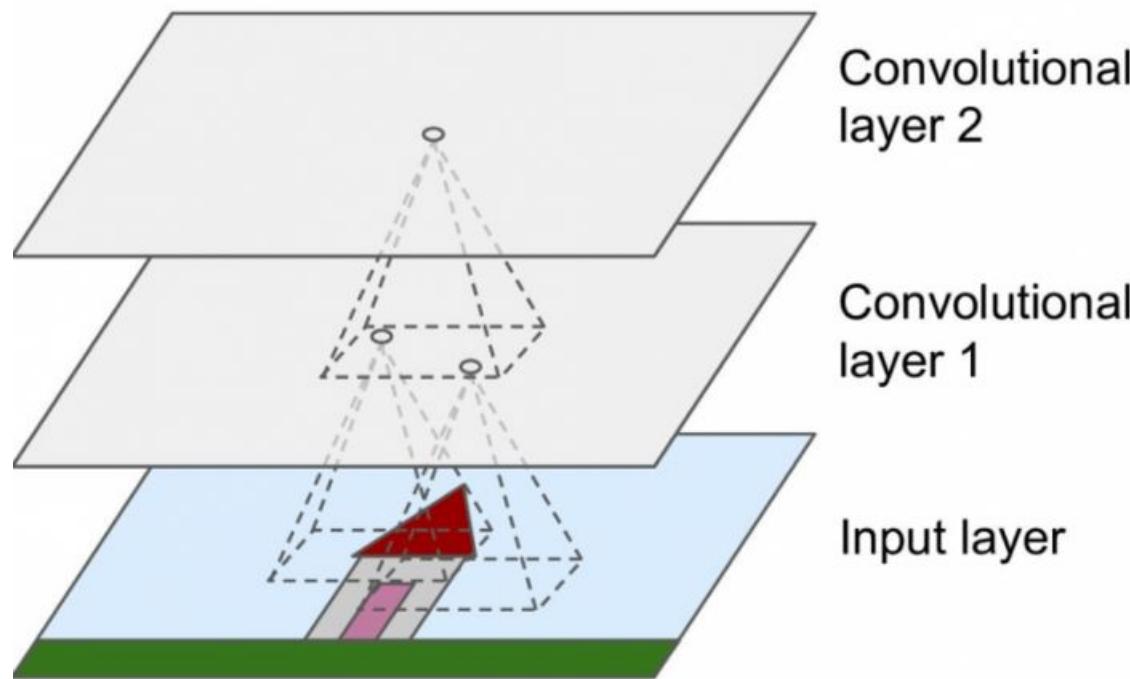
- **Experiments:** Complexity
adding the decoder for upscaling was a competitive runtime during inference

Backbones (encoders)	VGG16				Inception				ResNet50			
Methods	#PCL	#NFM	SFM	#PDEC	#PCL	#NFM	SFM	#PDEC	#PCL	#NFM	SFM	#PDEC
Details	≈19.6	1024	28x28	≈23.1	≈25.6	1024	28x28	≈5.7	≈23.9	2048	28x28	≈9
CAM* [58]		.2ms			.2ms			.3ms				
GradCAM [32]		7.7ms			21.1ms			27.8ms				
GradCAM++ [7]		23.5ms			23.7ms			28.0ms				
Smooth-GradCAM [25]		62.0ms			150.7ms			136.2ms				
XGradCAM [12]		2.9ms			19.2ms			14.2ms				
LayerCAM [15]		3.2ms			18.2ms			17.9ms				
Mean		16.6ms			38.8ms			37.4ms				
ours + STDCL		6.2ms			25.5ms			18.5ms				
ACoL [55]		12.0ms			19.2ms			24.9ms				
SPG [56]		11.0ms			18ms			23.9ms				
ADL [9]		6.4ms			16.0			14.4ms				
ScoreCAM [44]		1.9sec			3.4sec			9.3sec				
SSCAM [24]		1min45sec			2min16sec			5min49sec				
IS-CAM [23]		30.1sec			39.0sec			1min39sec				

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

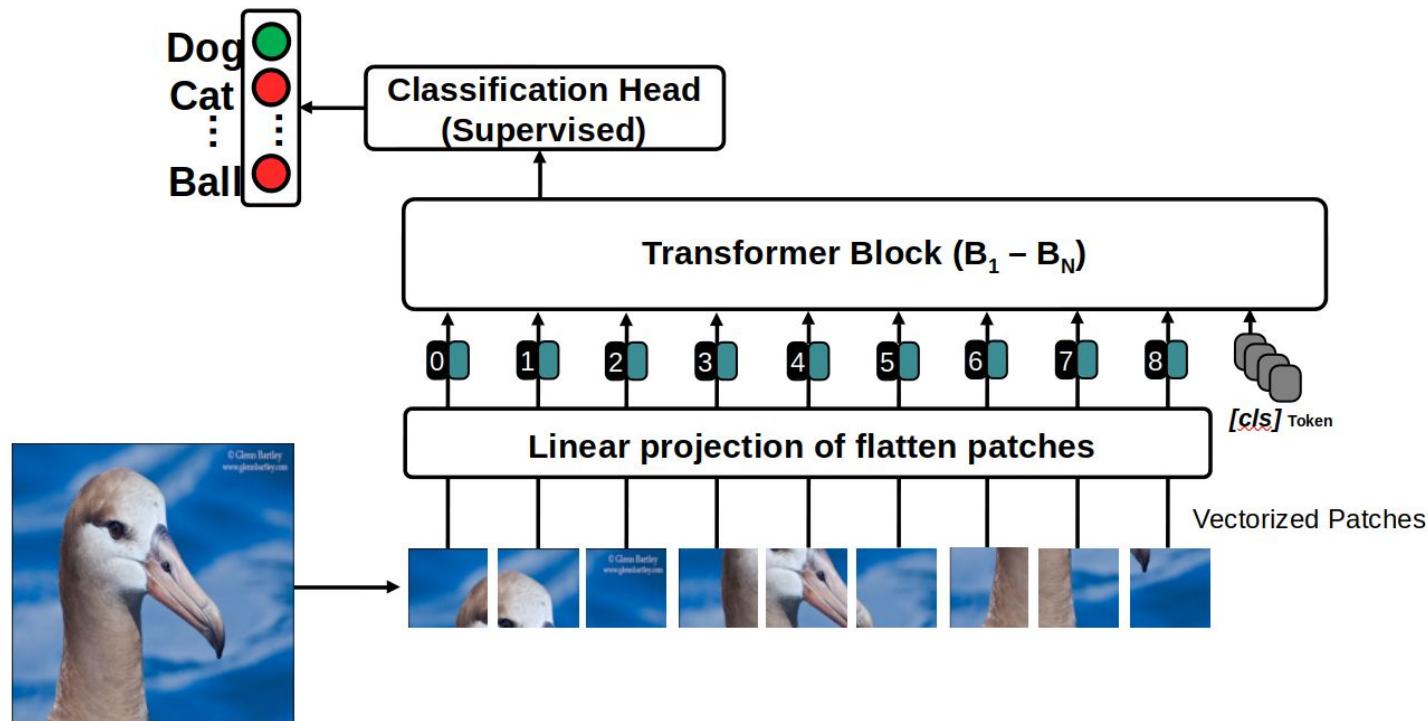
- Small local receptive field of CNN



Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Long receptive field of transformers: Long range

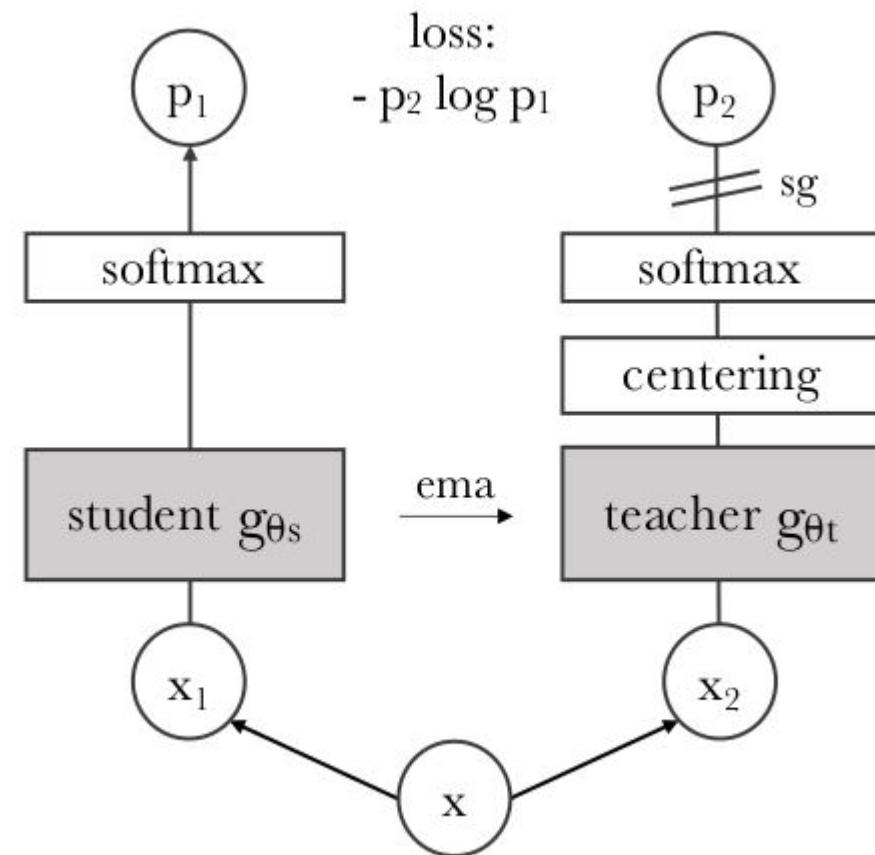


Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Dino: self-distillation

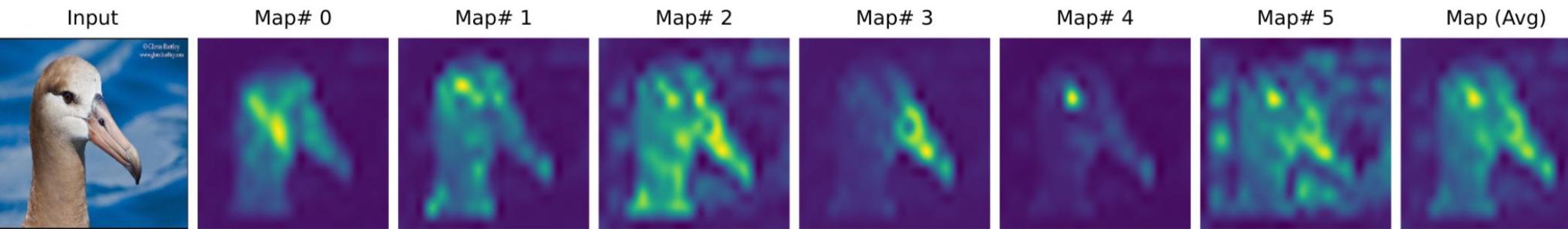
ViT: Vision Transformers



Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Dino: self-distillation

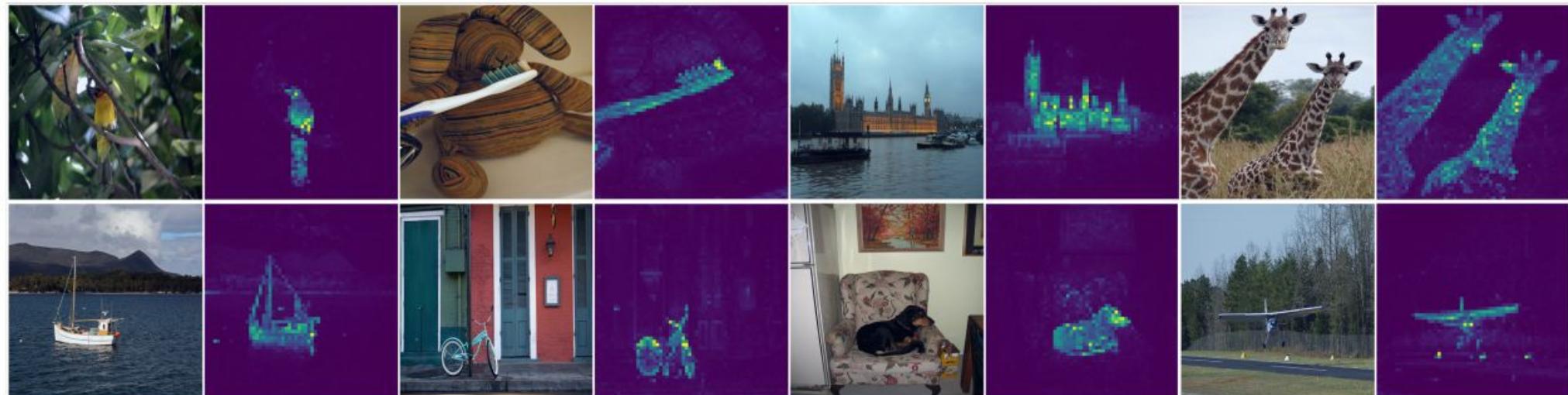


Localization at **output 6 heads.**
Which head to select?

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Dino: self-distillation



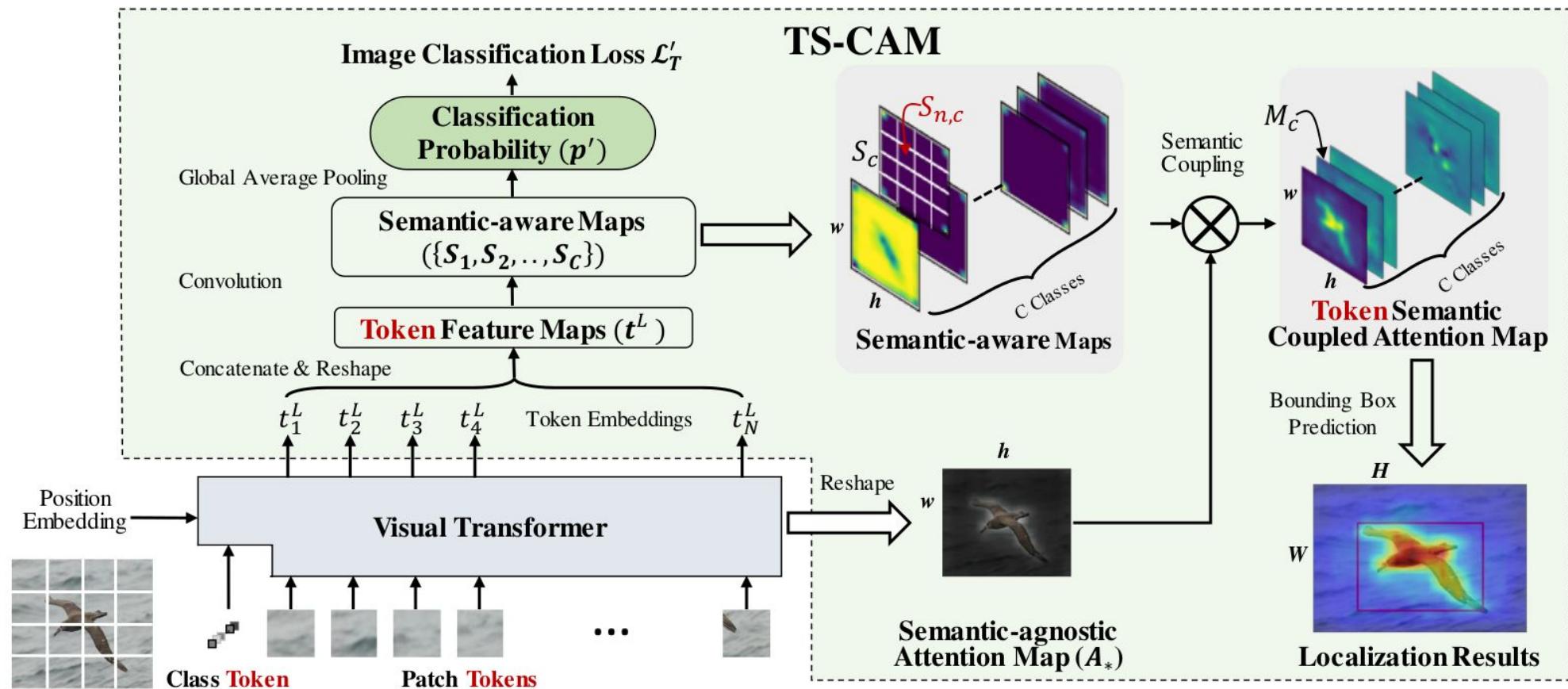
Localization at Dino's **output heads**:
Best head selected using ground truth

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- TS-CAM: Token Semantic CAM

How to build CAM?



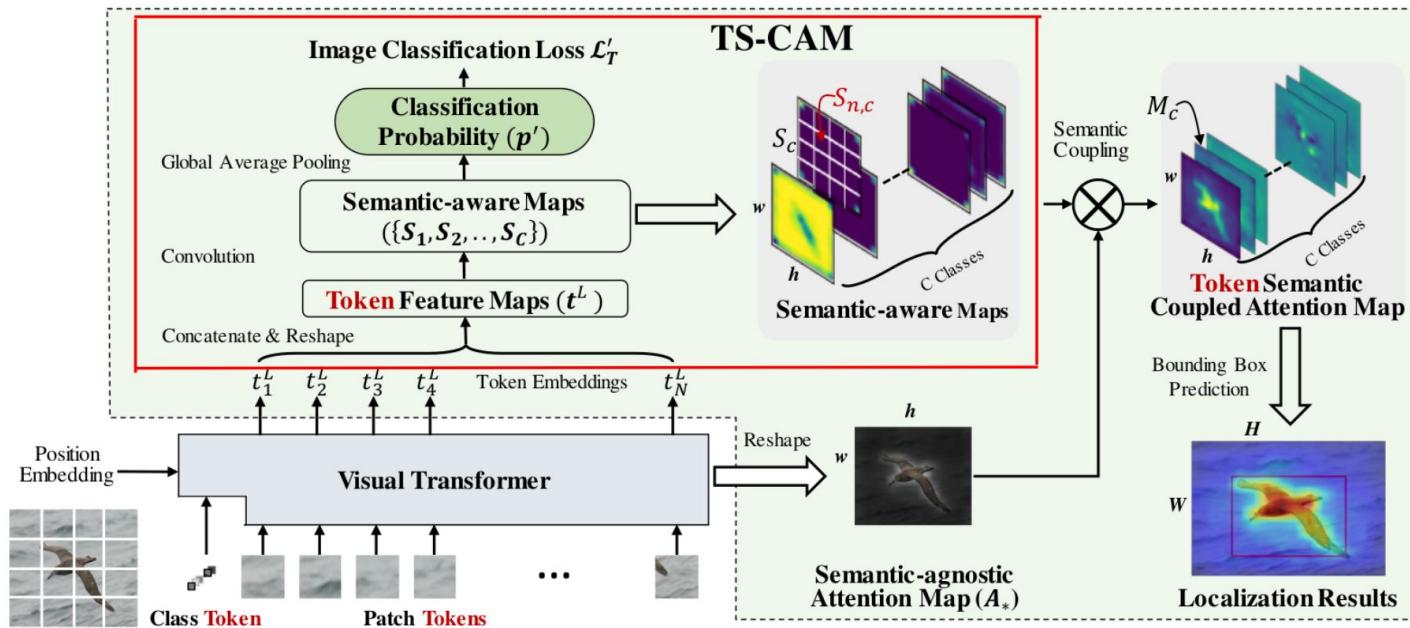
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- TS-CAM: Token Semantic CAM

How to build CAM?

CAMs



$$S_c = \sum_d t_d^L * k_{c,d},$$

3x3 conv

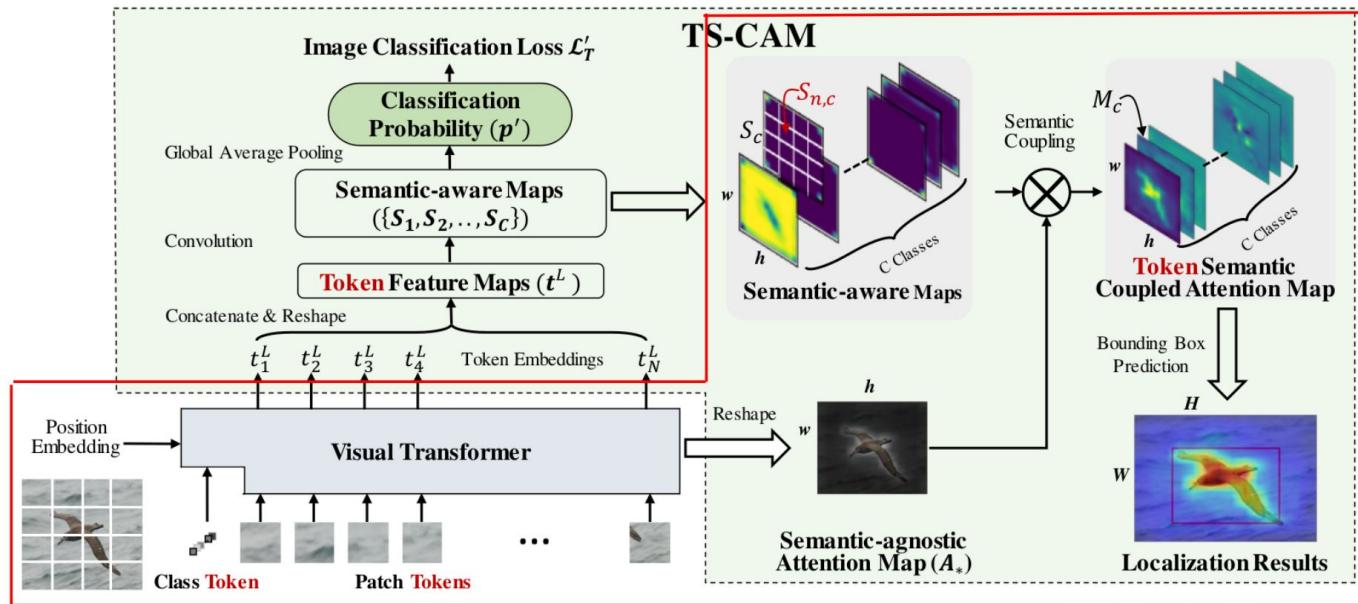
$$\begin{aligned} \mathcal{L}'_T &= -\log p'_y \\ &= -\log \frac{\exp(\sum_n S_{n,y}/N)}{\sum_c \exp(\sum_n S_{n,c}/N)}, \end{aligned}$$

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

● TS-CAM: Token Semantic CAM

How to build CAM?



Cross-layers
attention:

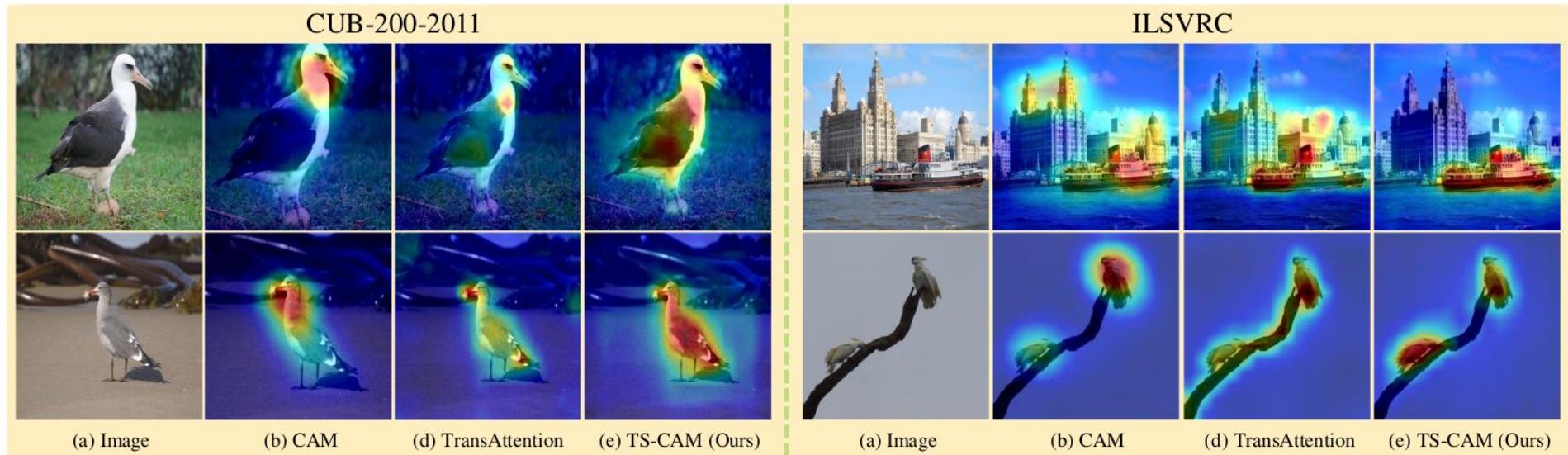
$$A_* = \frac{1}{L} \sum_l A_*^l,$$

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- TS-CAM: Token Semantic CAM

How to build CAM?

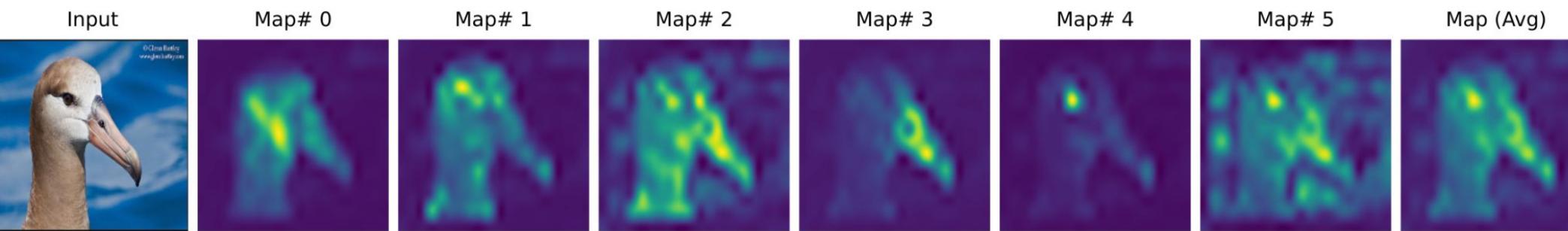


Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work

How to build CAM?



Localization at **output 6 heads.**
Which head to select?

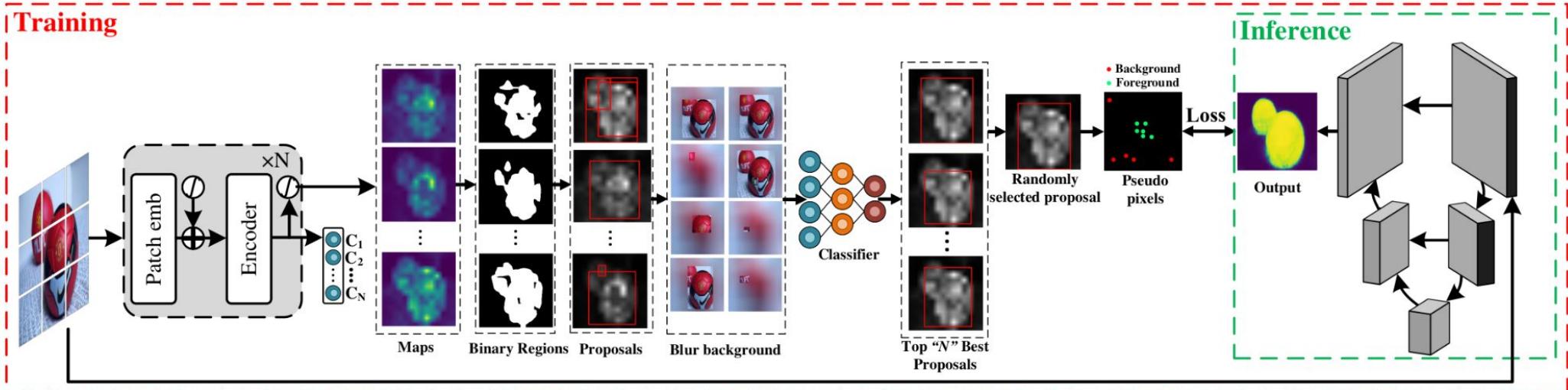
→ Fuse heads!

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?



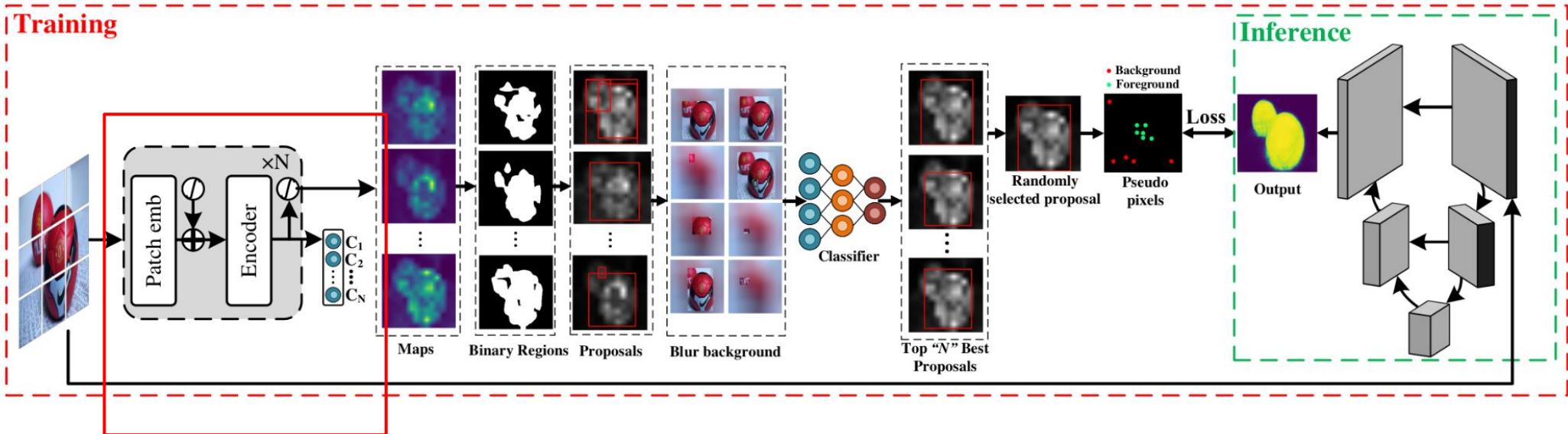
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



Dino

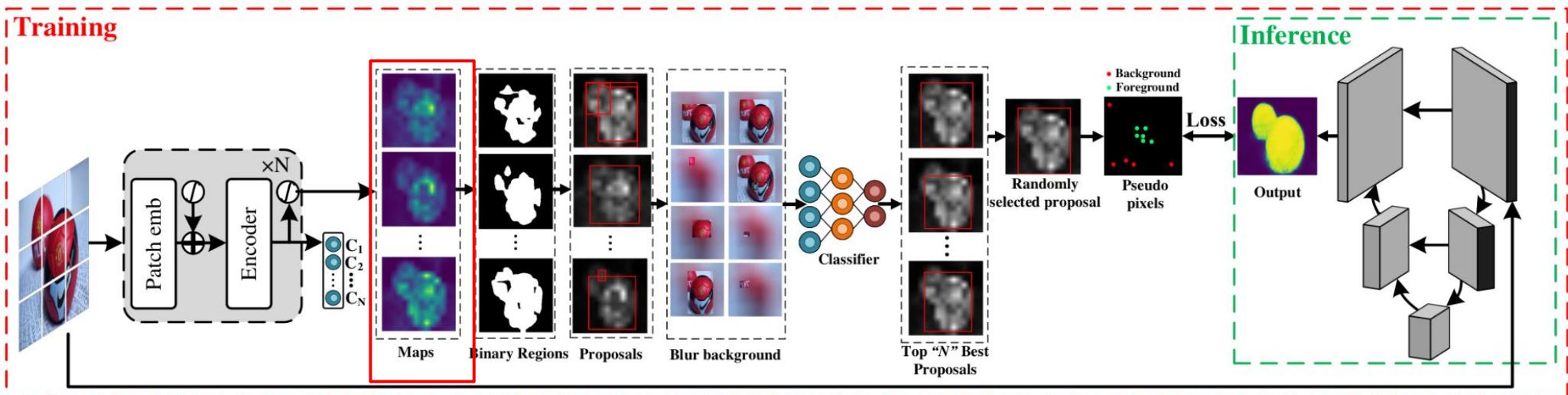
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



K heads

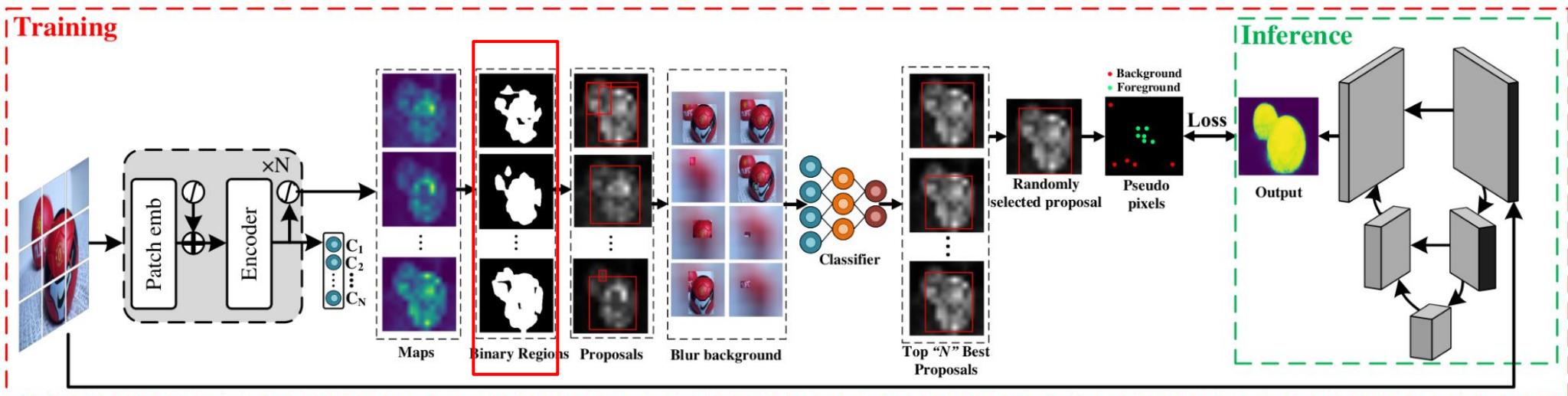
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



Proposals
generation

Case Study (b): Transformer-based models for WSOL

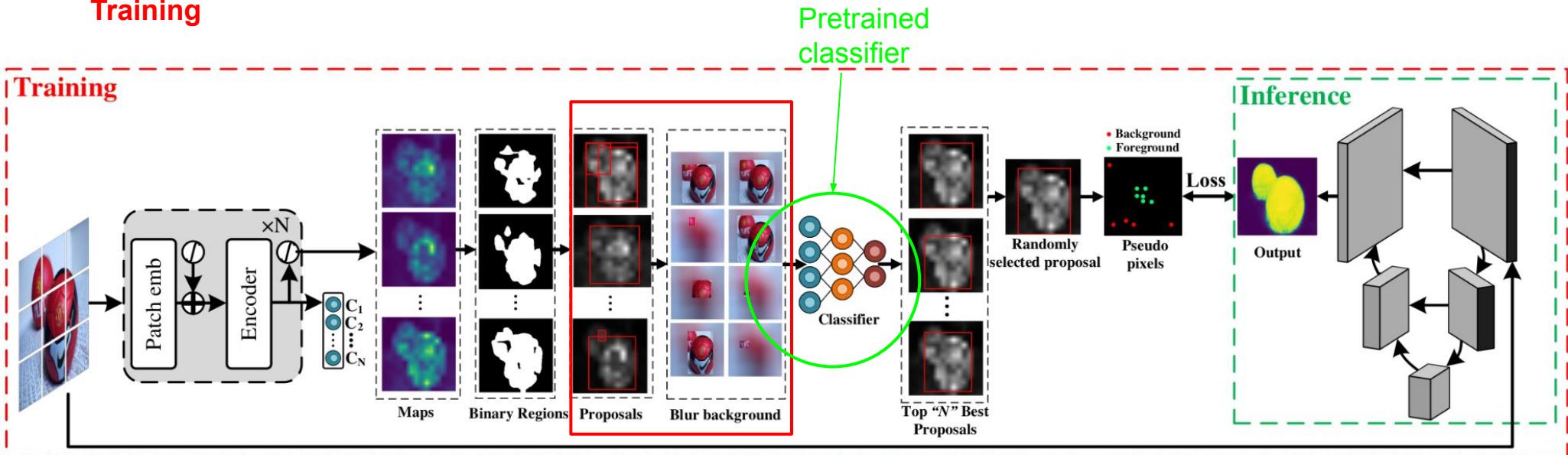
(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training

Training



Discriminative
Proposals selection

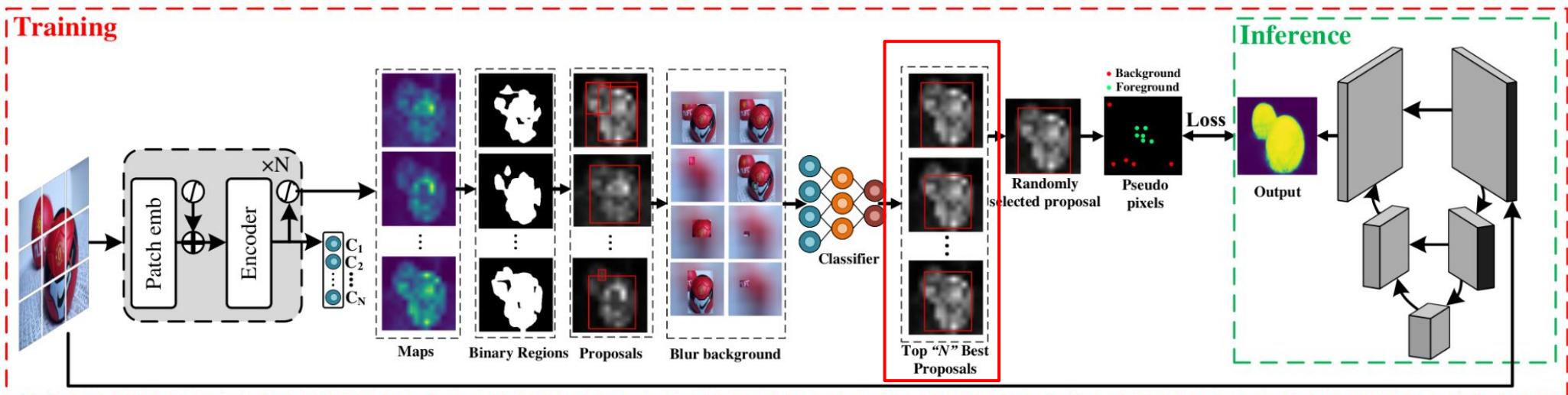
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



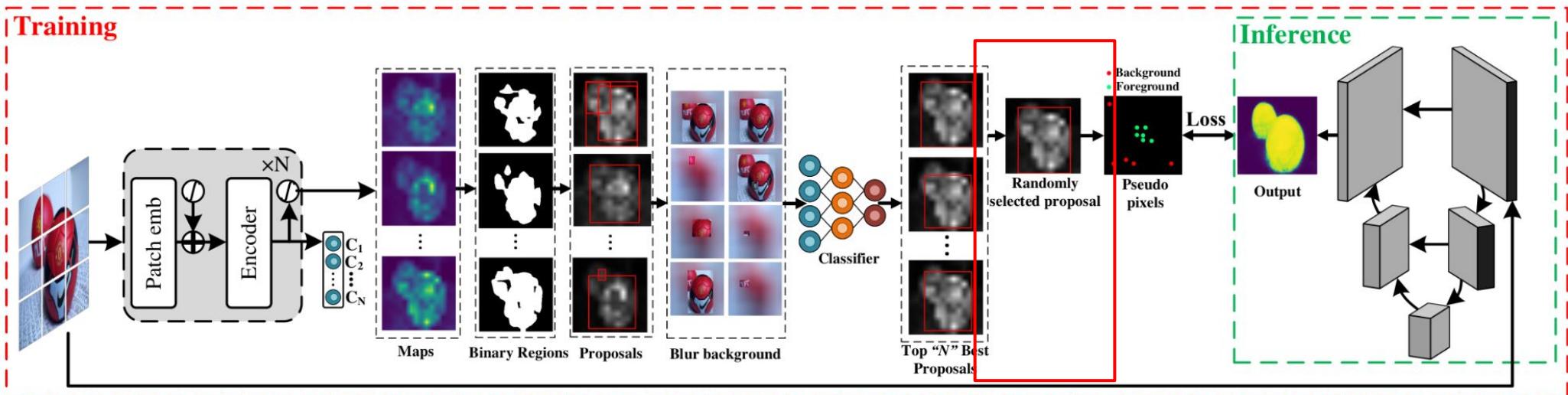
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



Random selection of a
proposal from the top-n
proposals

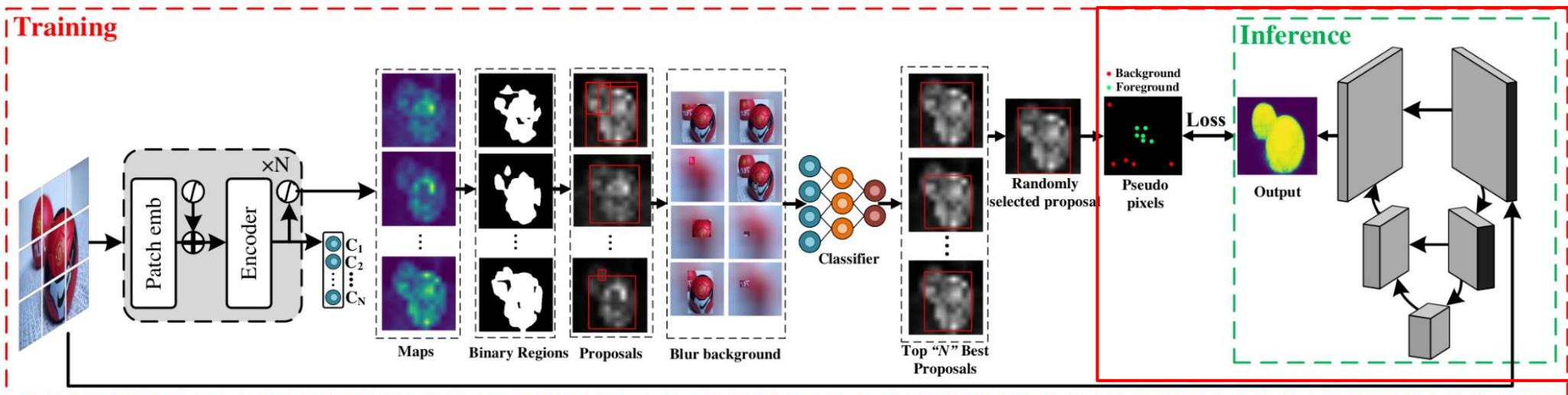
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



Training U-Net like net
using sampled seeds
(FG/BG) from proposal

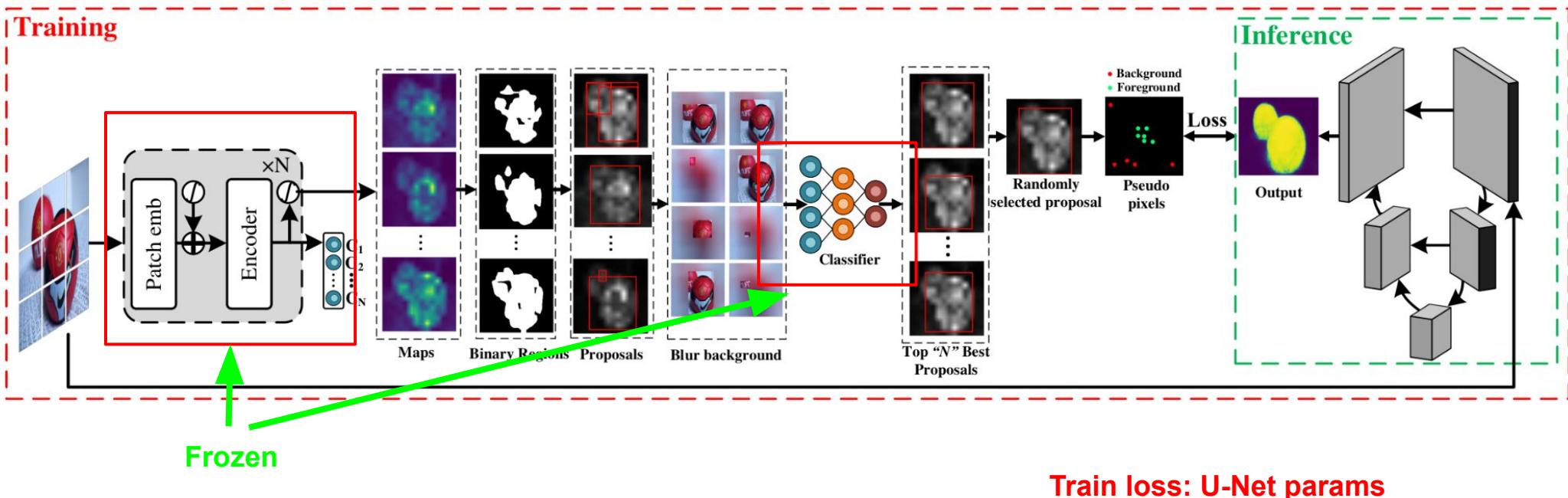
Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Training



$$\mathcal{L}_{Total} = \min_{\theta} \lambda_{CLS} + \lambda_{CPA} \mathcal{L}_{CPA} + \lambda_{CRF} \mathcal{L}_{CRF}$$

Case Study (b): Transformer-based models for WSOL

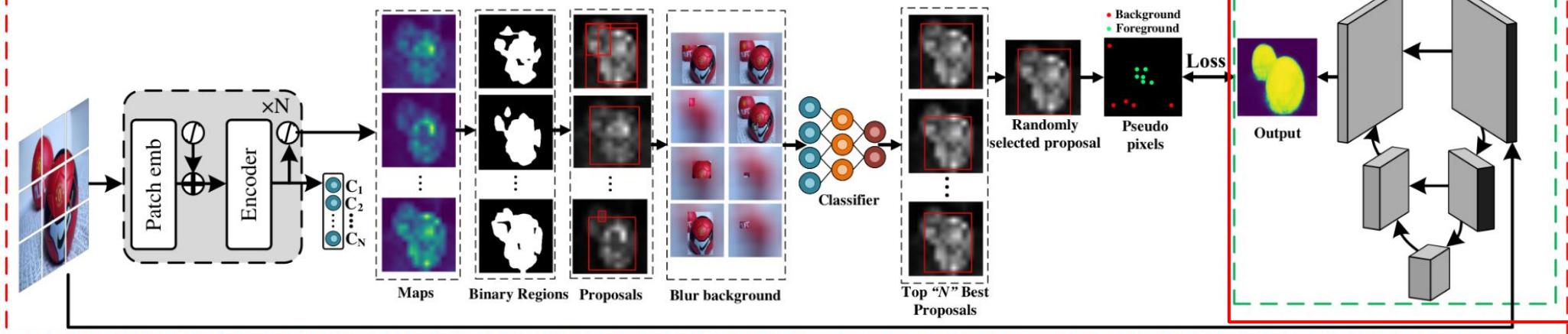
(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

Inference

Training



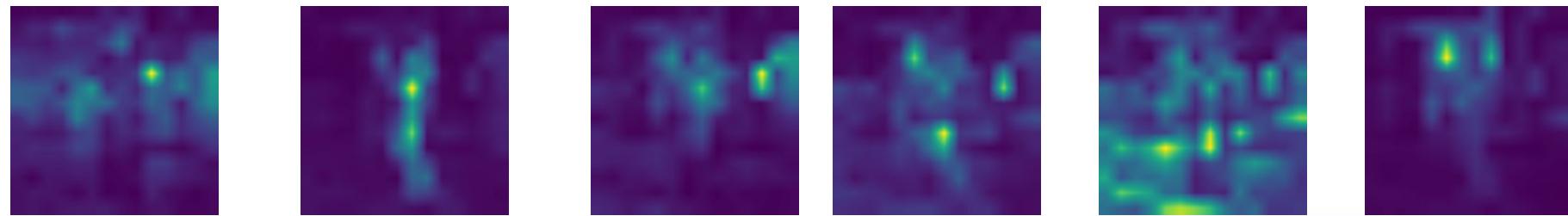
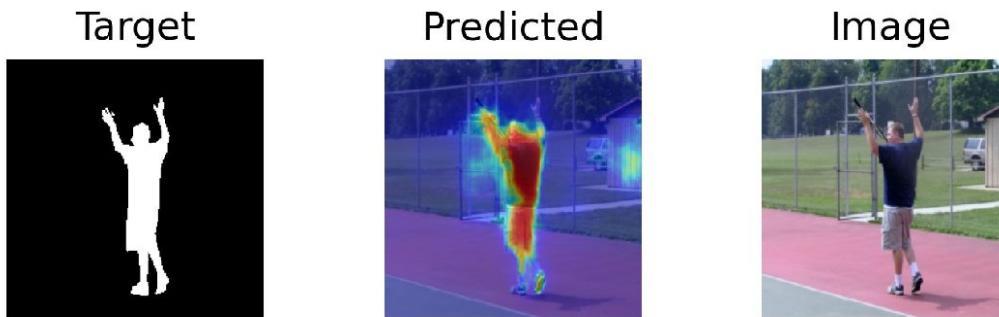
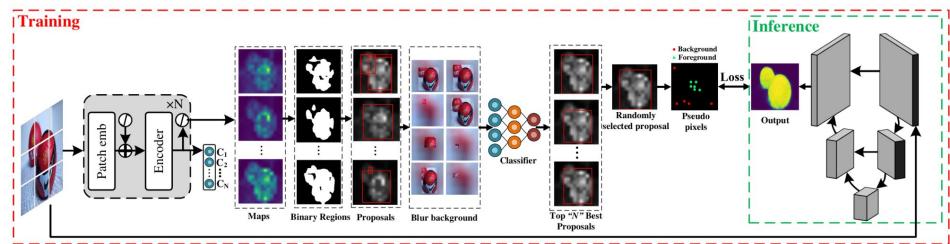
Localization: Using
only input image and
U-Net.

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?

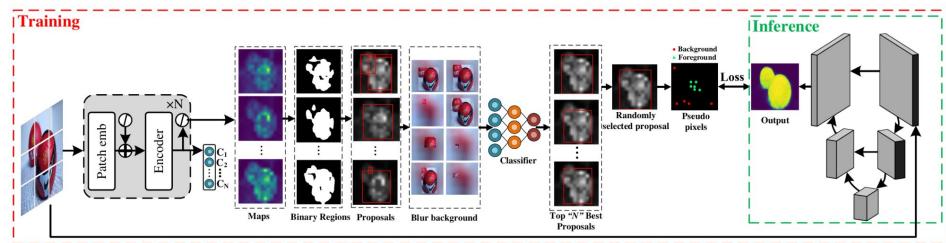


Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?



Model prediction

Target



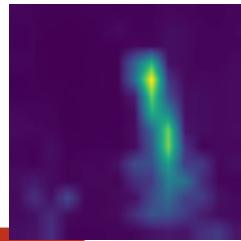
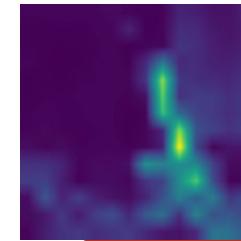
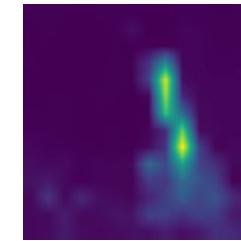
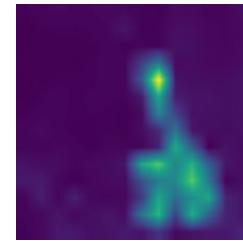
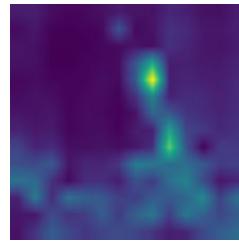
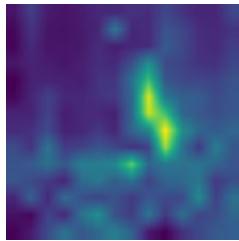
Predicted



Image



Dino heads

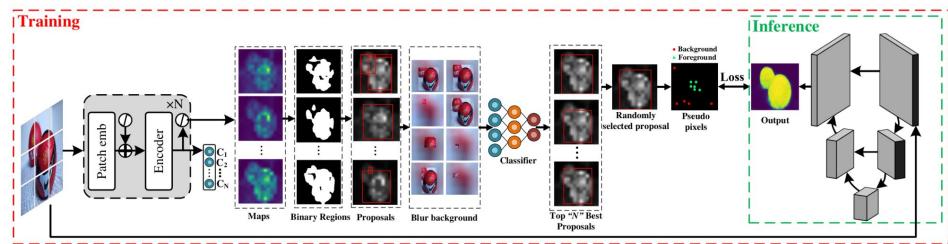


Case Study (b): Transformer-based models for WSOL

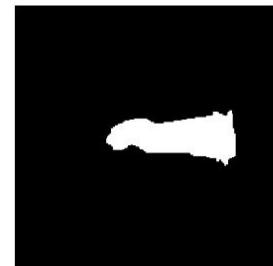
(Ongoing work)

- Ongoing work: fusion of heads

How to build CAM?



Target



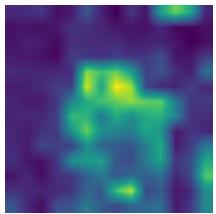
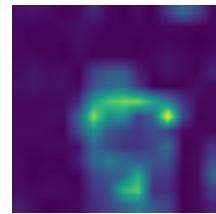
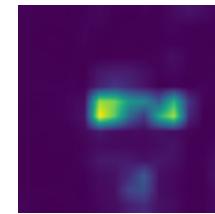
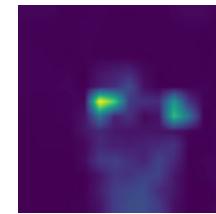
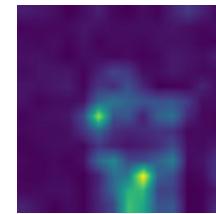
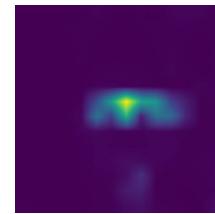
Predicted



Image



Model prediction

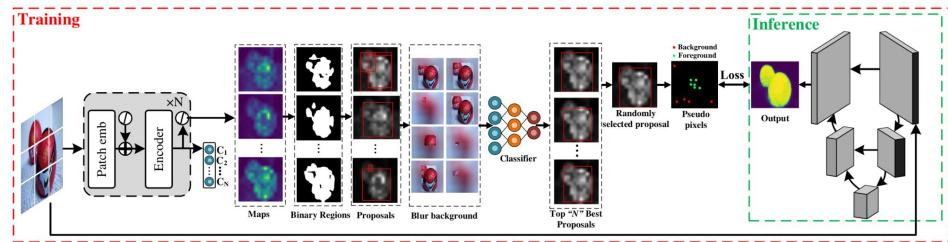


Dino heads

Case Study (b): Transformer-based models for WSOL

(Ongoing work)

- Ongoing work: fusion of heads



Paper + code: will be available early september.

Online Resources

F-CAM (WACV 2022):

- paper: <https://arxiv.org/abs/2109.07069>
- code: <https://github.com/sbelharbi/fcam-wsol>

F-CAM extension to histology (MIDL 2022):

- paper: <https://arxiv.org/abs/2201.02445>
- code: <https://github.com/sbelharbi/negev>

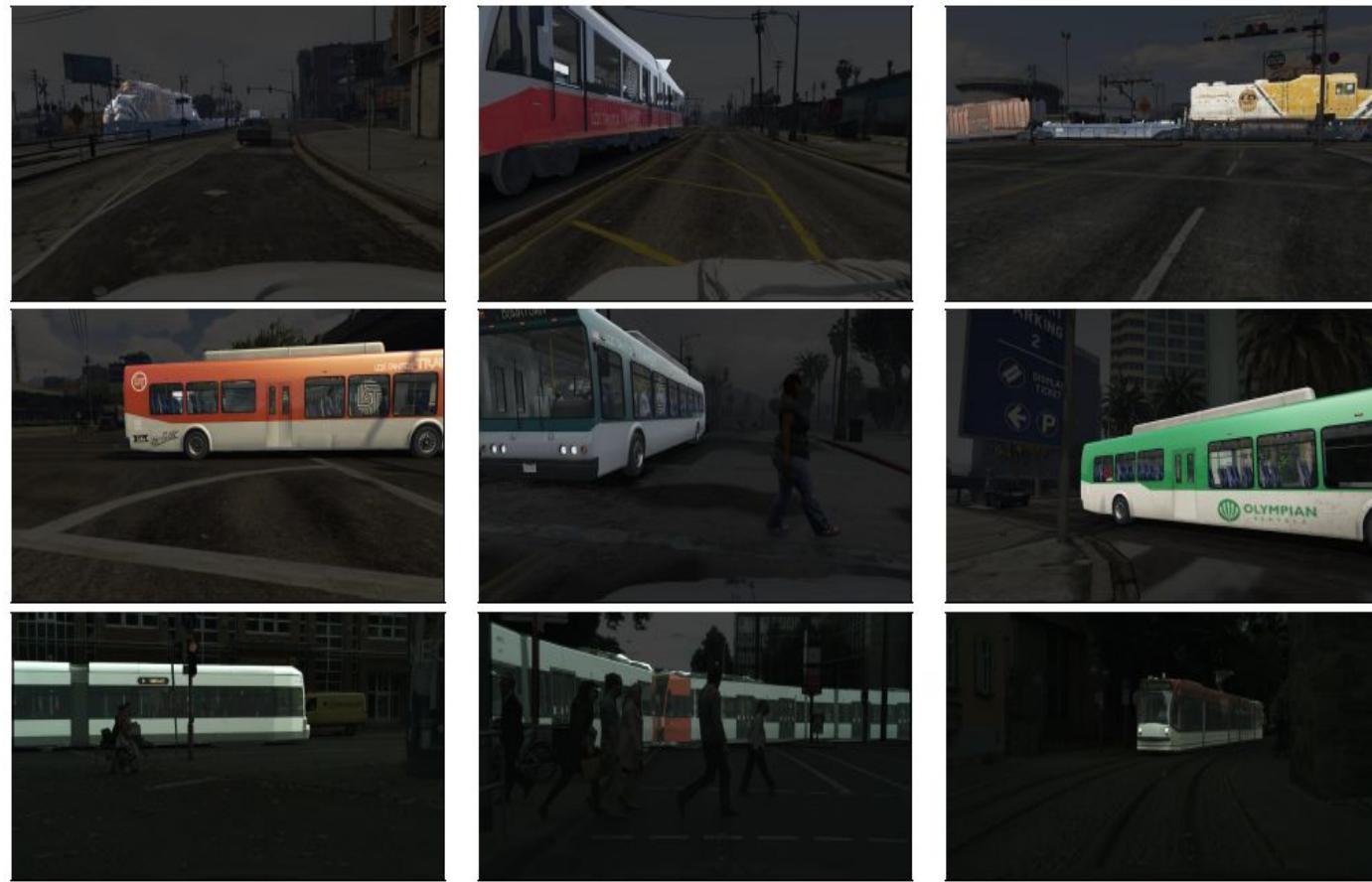
Part 3:

Review of Loss Functions and Settings for WSSS methods

Dense pixel-wise labels are very expensive



Domain shifts make things worse



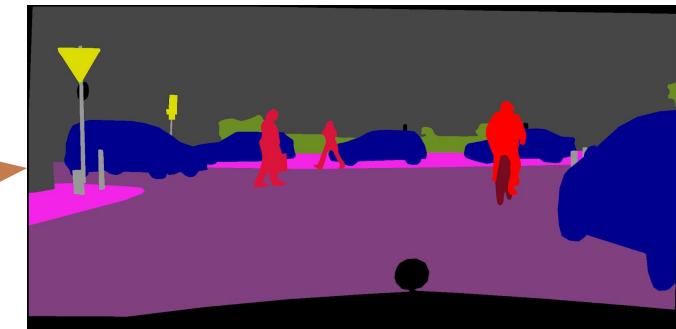
“train”
GTA

“bus”
GTA

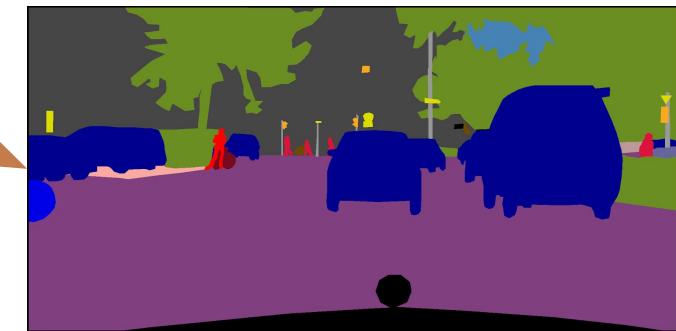
“train”
Cityscapes

Figures from [Zhang et al., A Curriculum Domain Adaptation Approach to the Semantic Segmentation of Urban Scenes TPAMI 2019]

Domain shifts make things worse



Frankfurt



Zurich

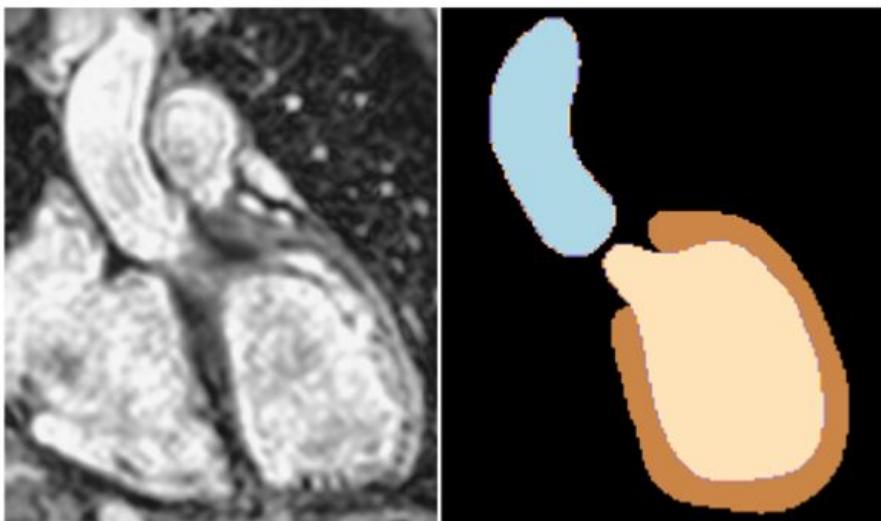
road	sidewalk	building	wall	fence	pole	traffic light	traffic sign	vegetation	terrain
sky	person	rider	car	truck	bus	train	motorcycle	bicycle	unlabeled

Cityscapes (5000 images): labeling of 1 image takes 90 min at average

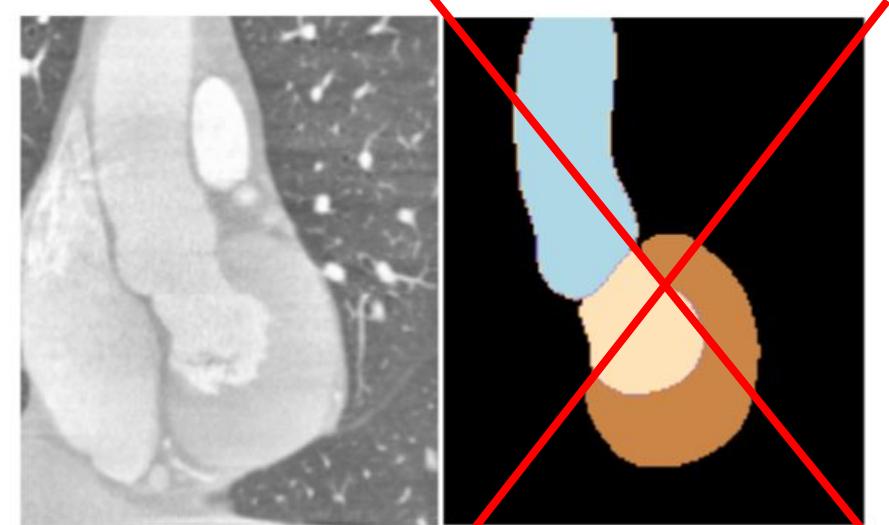
[Cordt et al., CVPR 2016]

Domain shifts make things worse

Source domain (MRI)



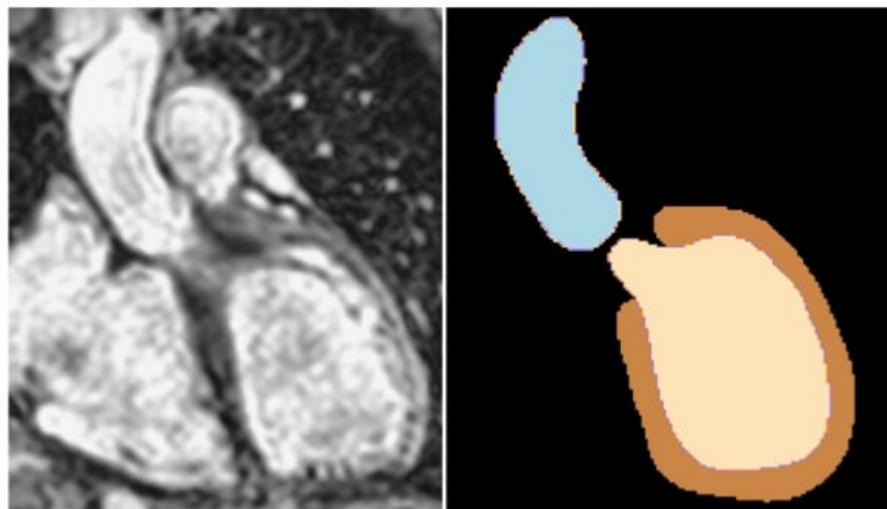
Target domain (CT)



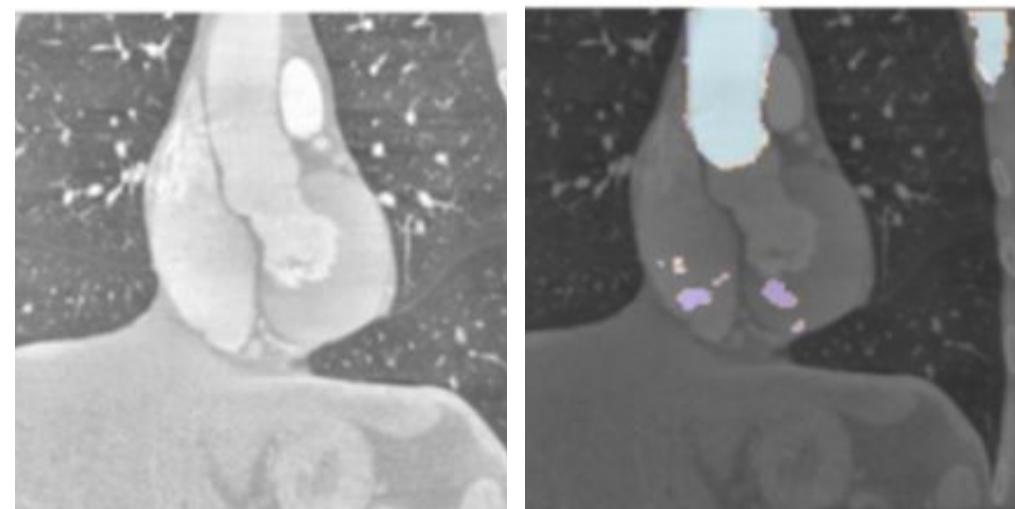
[Images from Bateson et al., Constrained Domain Adaptation for Image Segmentation, TMI'21]

Domain shifts make things worse

Source domain (MRI)

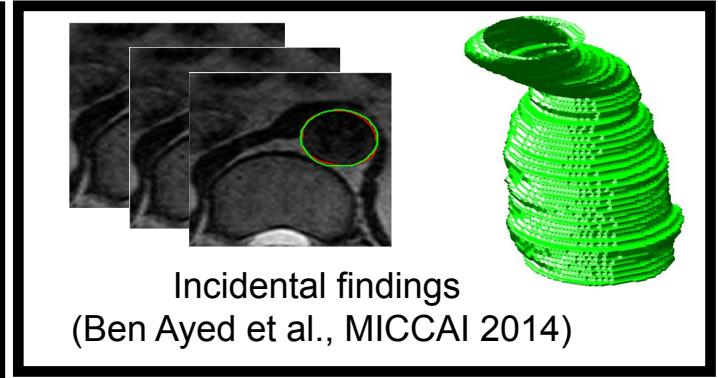
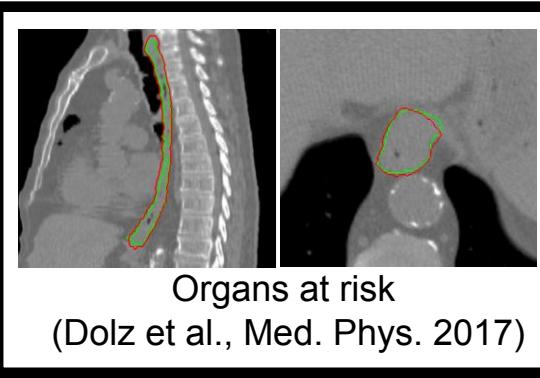
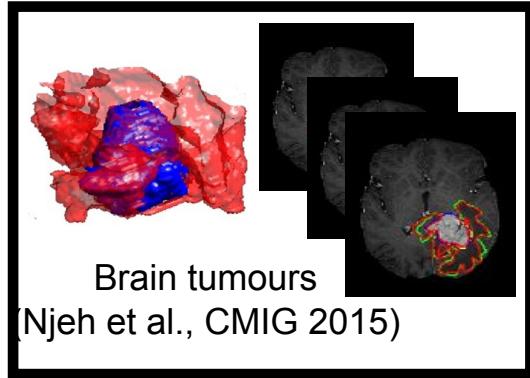
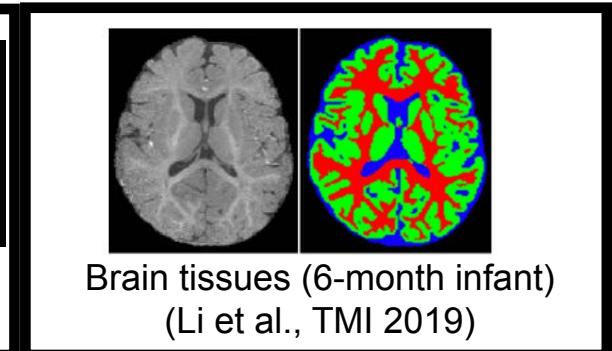
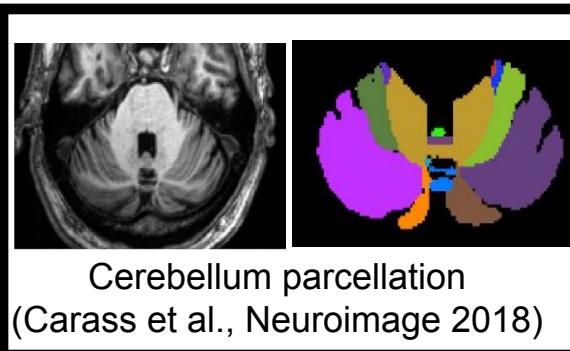
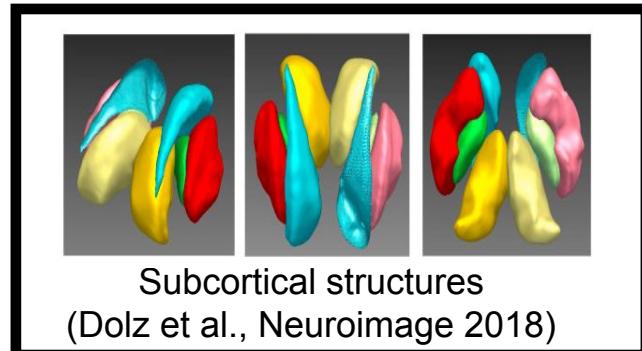


No adaptation
(bad generalization to the target)



[Images from Bateson et al., Constrained Domain Adaptation for Image Segmentation, TMI'21]

...and more complex is some applications (e.g medical image analysis)



Dense labels are more complicated in medical imaging:

Not anywhere close to the 10k images of Pascal VOC and the 5k of Cityscapes

Crowdsourcing?

Select all images with
esophagus
Click verify once there are none left.

VERIFY

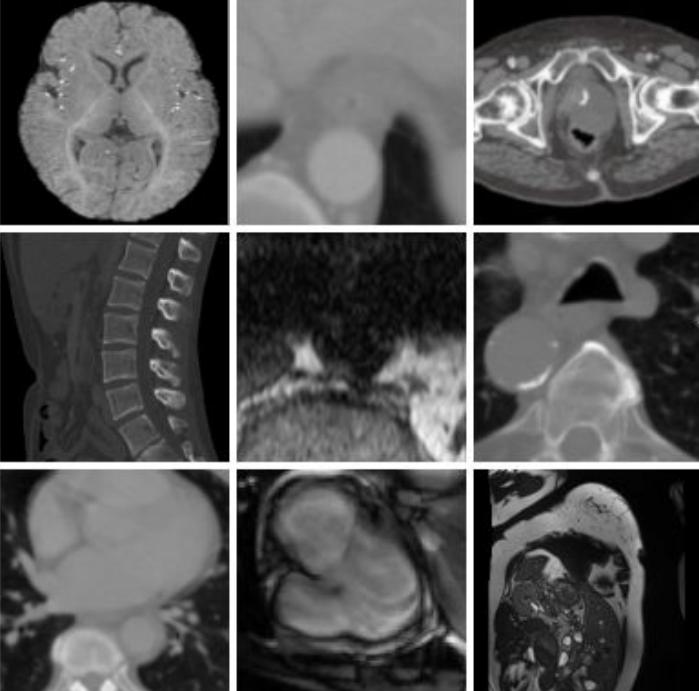


Dense labels are more complicated in medical imaging:

Not anywhere close to the 10k images of Pascal VOC and the 5k of Cityscapes

Crowdsourcing?

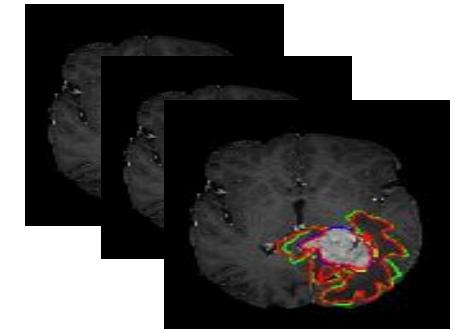
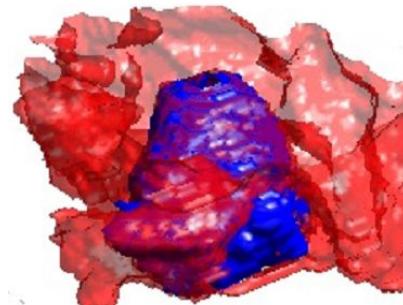
Select all images with
esophagus
Click verify once there are none left.



VERIFY

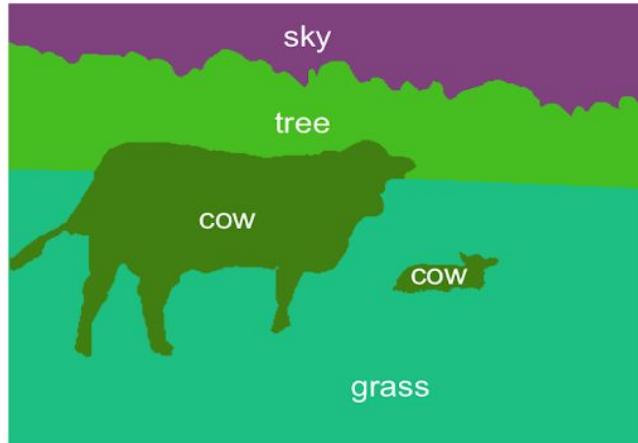
□

Dense 3D annotations: several hours
(of radiologist time)

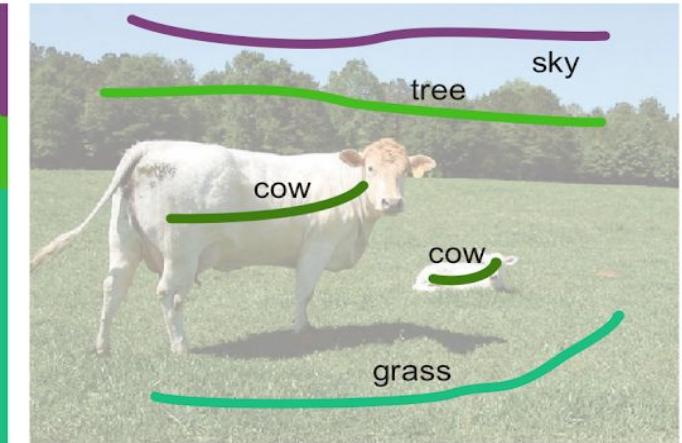


Semi-supervised learning (SSL)

A lot of unlabeled data, and only a fraction of points are labeled



Full annotations



Semi-supervised

Figures from Lin et al. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation, CVPR 2016

Forms of weak supervision in segmentation

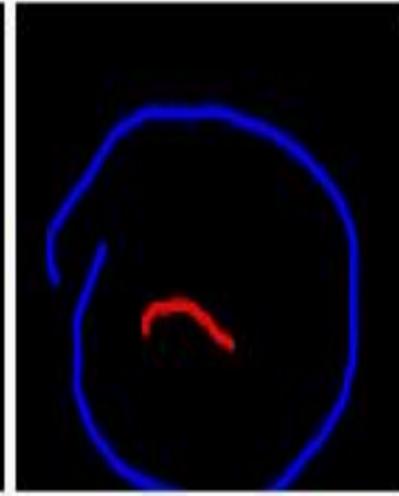


Car
Parking
Sky
No person

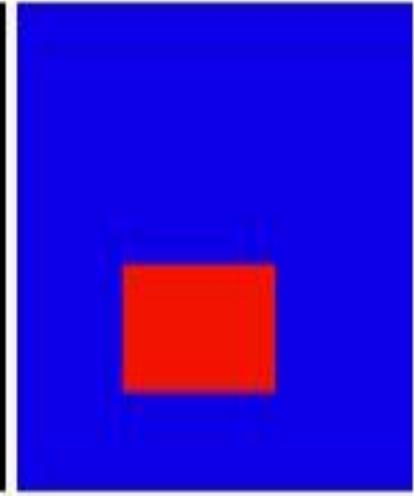
Image tags
(MIL)



Points
(SSL)



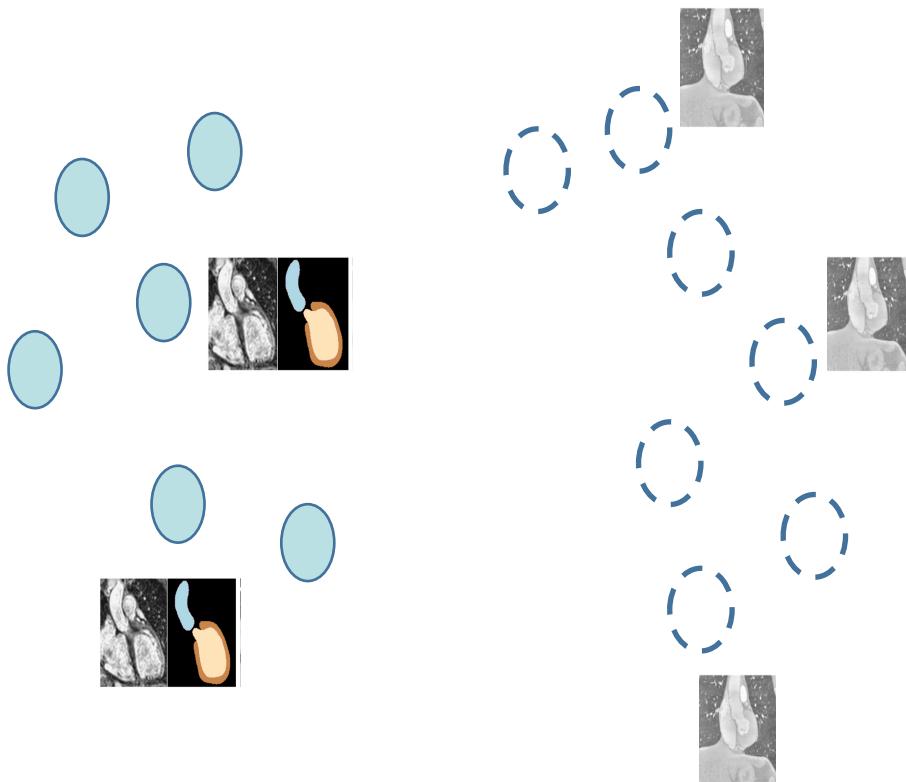
Scribbles
(SSL)



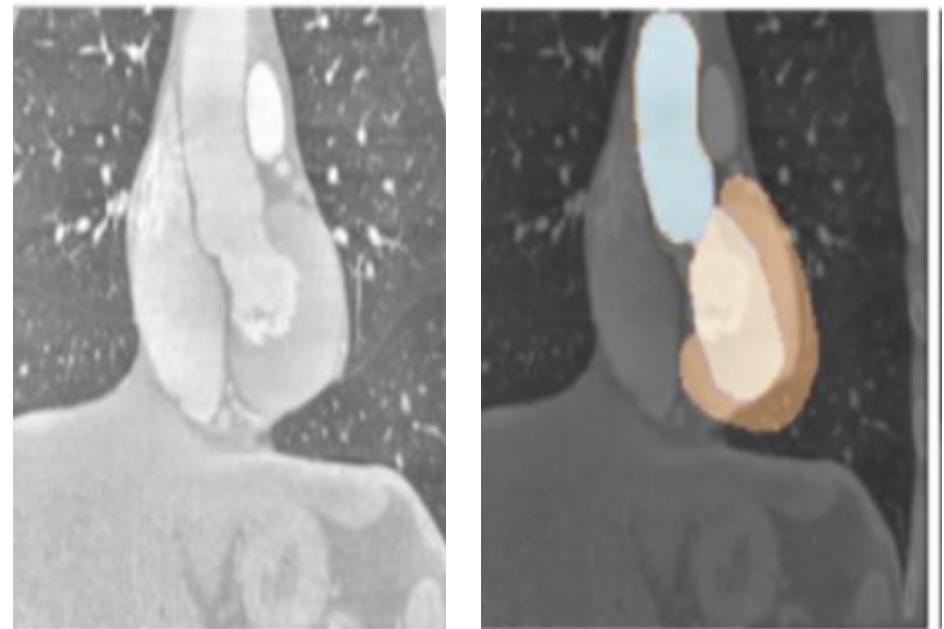
Boxes
(MIL)

Closely related problem: Unsupervised Domain Adaptation (UDA)

Training on both
labeled and unlabeled data

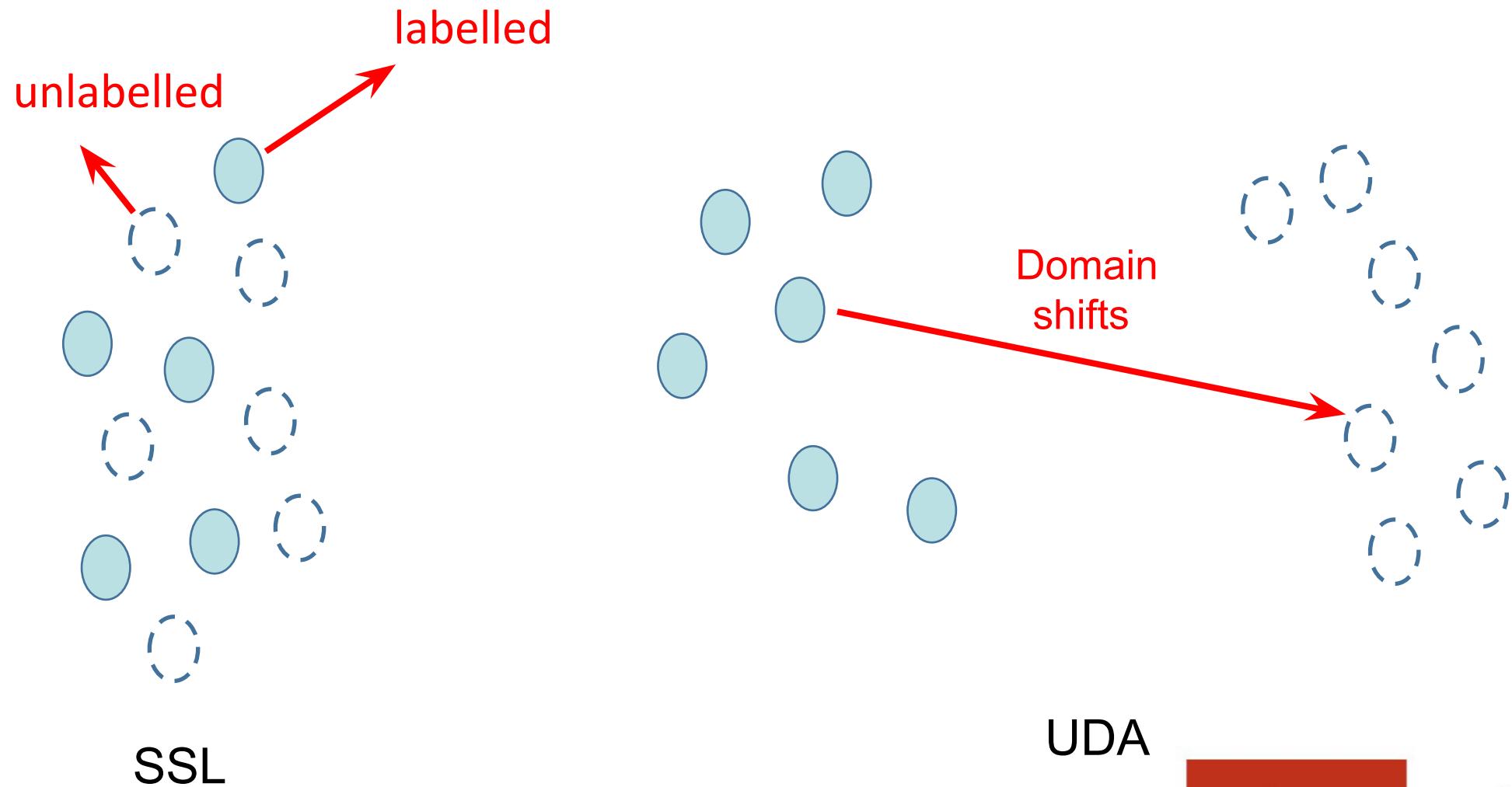


Corrected result with UDA



[Images from Bateson et al., Constrained Domain Adaptation for Image Segmentation, TMI'21]

UDA = SSL +Domain Shifts



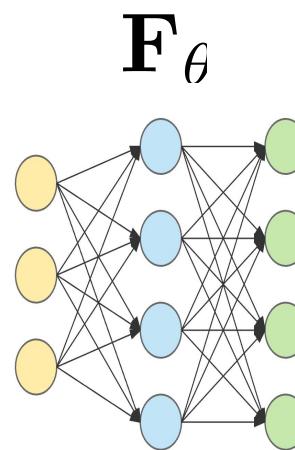
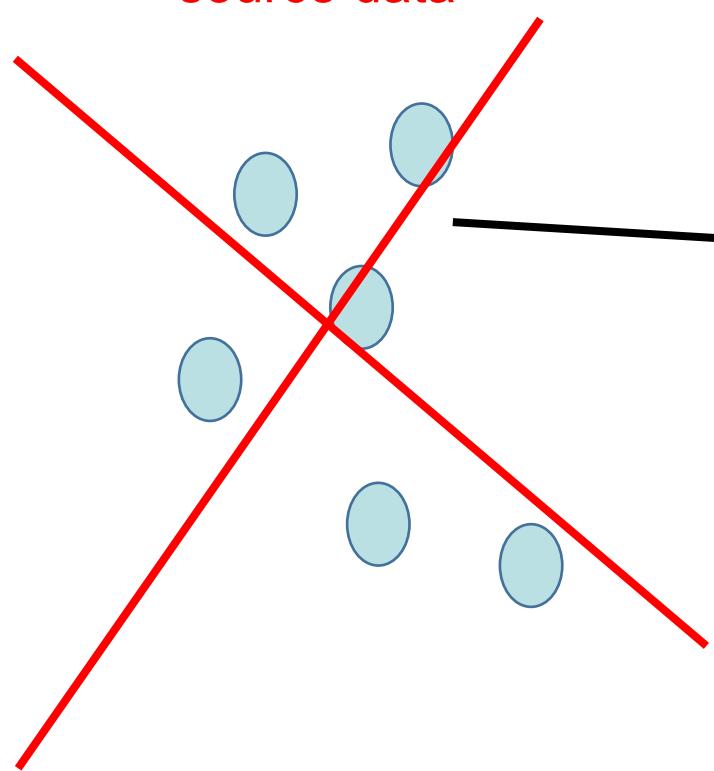
Test-time adaptation (TTA)

UDA without access to the source data

No access to the labeled
source data

just access to the
model

Unlabeled target data



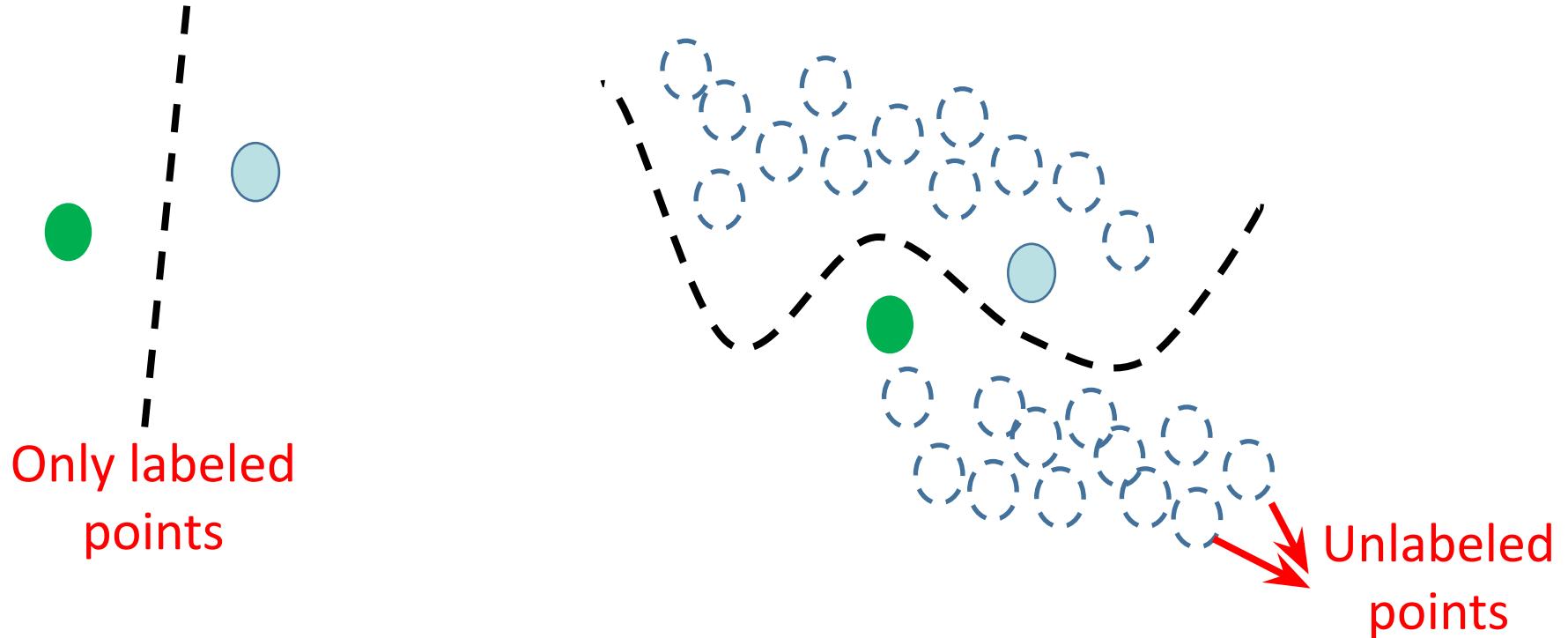
SSL/MIL/UDA/TTA unsupervised loss functions in a nutshell

Leveraging unlabeled data with priors

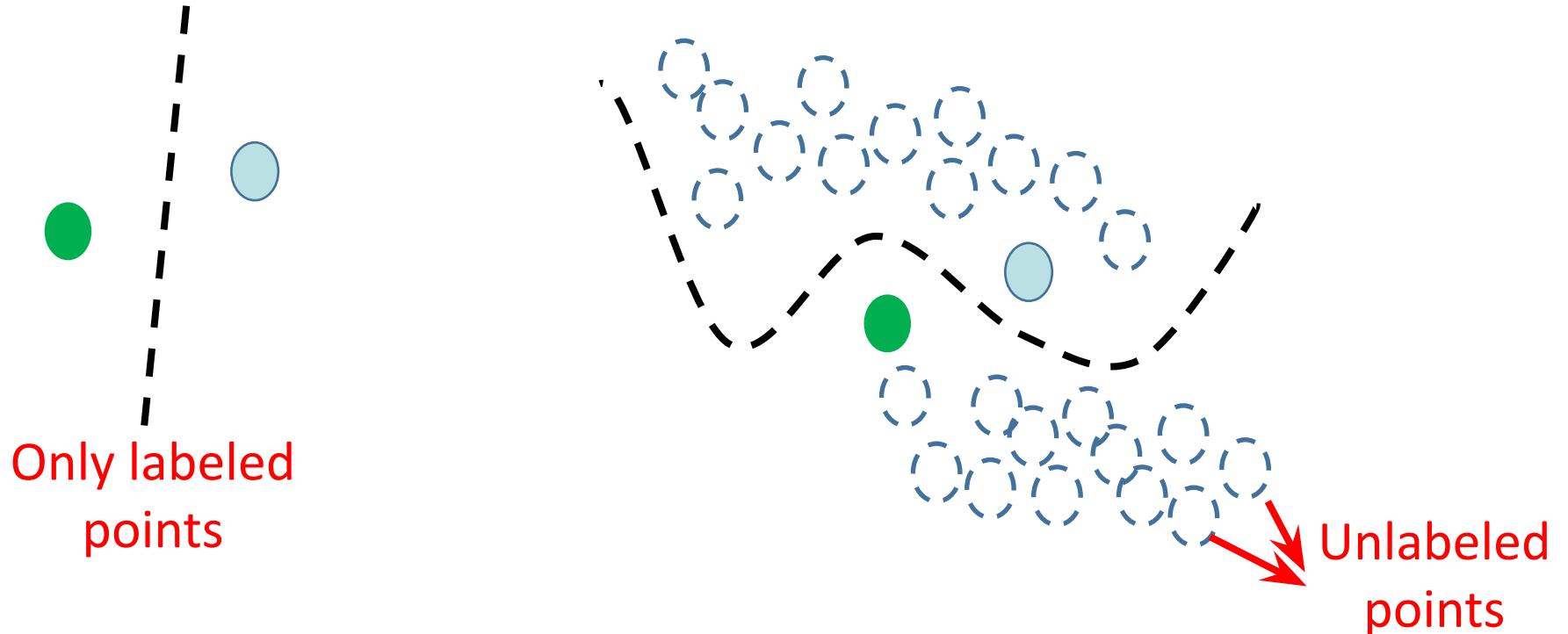
- Structure-driven priors: Regularization
- Knowledge-driven priors (e.g., anatomical constraints)
- Invariance priors (e.g., contrastive learning)
- Multi-modal priors (e.g., text info associated with the images)

Unsupervised Manifold Regularization

In SSL: Learning from both unlabeled and labeled data

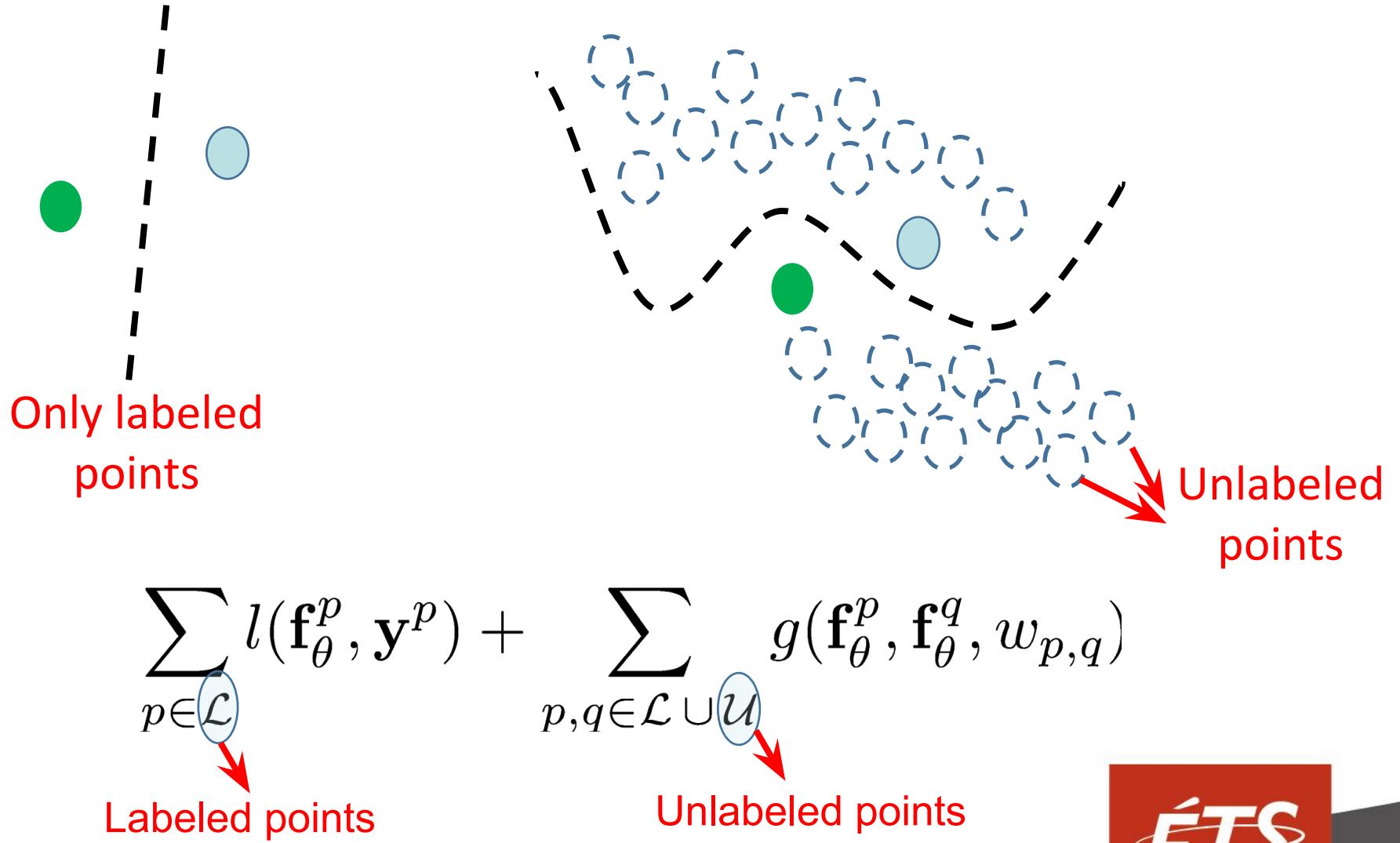


In SSL: Learning from both unlabeled and labeled data



$$\sum_{p \in \mathcal{L}} l(\mathbf{f}_\theta^p, \mathbf{y}^p) + \sum_{p,q \in \mathcal{L} \cup \mathcal{U}} g(\mathbf{f}_\theta^p, \mathbf{f}_\theta^q, w_{p,q})$$

In SSL: Learning from both unlabeled and labeled data



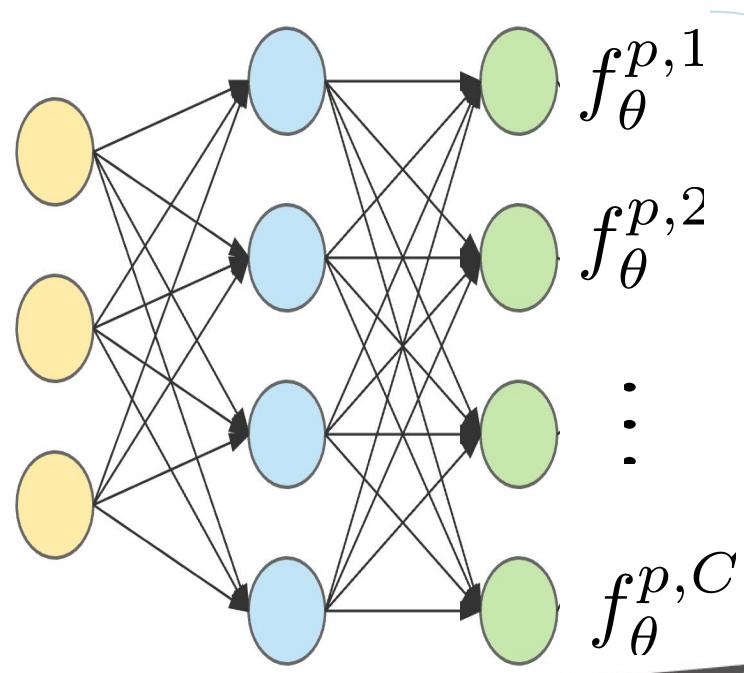
In SSL: Learning from both unlabeled and labeled data

$$\sum_{p \in \mathcal{L}} l(\mathbf{f}_\theta^p, \mathbf{y}^p) + \sum_{p,q \in \mathcal{L} \cup \mathcal{U}} g(\mathbf{f}_\theta^p, \mathbf{f}_\theta^q, w_{p,q})$$

e.g.: cross-entropy

e.g.: softmax outputs

Labels



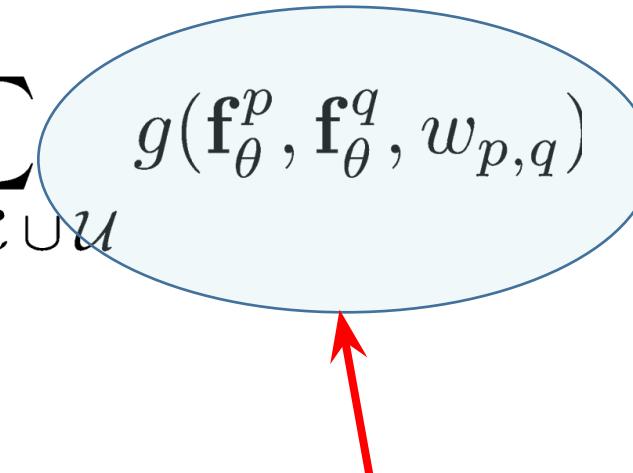
Laplacian regularization: Standard in classical SSL

$$\sum_{p \in \mathcal{L}} l(\mathbf{f}_\theta^p, \mathbf{y}^p) + \sum_{p,q \in \mathcal{L} \cup \mathcal{U}} g(\mathbf{f}_\theta^p, \mathbf{f}_\theta^q, w_{p,q})$$

Diagram illustrating the loss function for semi-supervised learning:

- The first term, $\sum_{p \in \mathcal{L}} l(\mathbf{f}_\theta^p, \mathbf{y}^p)$, represents the supervised loss for labeled data. It is composed of:
 - \mathbf{f}_θ^p : Model output for labeled sample p .
 - \mathbf{y}^p : True label for labeled sample p .
 - $l(\cdot, \cdot)$: Loss function, with examples like cross-entropy.
- The second term, $\sum_{p,q \in \mathcal{L} \cup \mathcal{U}} g(\mathbf{f}_\theta^p, \mathbf{f}_\theta^q, w_{p,q})$, represents the unlabeled loss or regularization term. It is composed of:
 - \mathbf{f}_θ^p and \mathbf{f}_θ^q : Model outputs for unlabeled samples p and q .
 - $w_{p,q}$: Regularization weight.
 - $g(\cdot, \cdot, \cdot)$: Regularization function, with examples like Laplacian regularization.

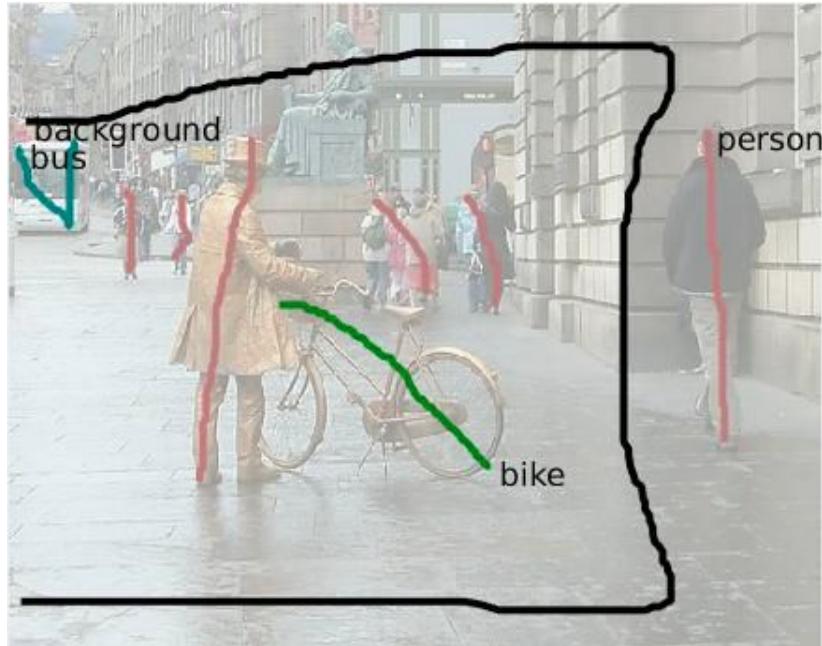
Laplacian regularization: Standard in classical SSL

$$\sum_{p \in \mathcal{L}} l(\mathbf{f}_\theta^p, \mathbf{y}^p) + \sum_{p,q \in \mathcal{L} \cup \mathcal{U}} g(\mathbf{f}_\theta^p, \mathbf{f}_\theta^q, w_{p,q})$$

$$w_{p,q} \|\mathbf{f}_\theta^p - \mathbf{f}_\theta^q\|^2$$

- [Weston et al., Deep Learning via semi-supervised embedding, ICML 2008]
- [Belkin et al., Manifold regularization: a geometric framework for learning from Labeled and Unlabeled Examples, JMLR 2006]
- [Zhu et al., Semi-Supervised Learning Using Gaussian Fields and Harmonic Functions, ICML 2003]

Semi-supervision loss in segmentation

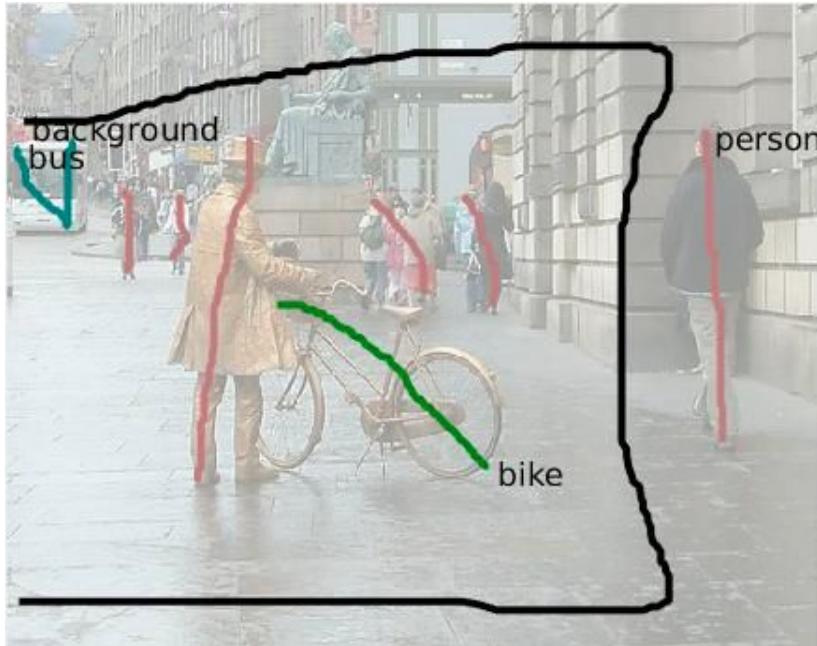
$$\min_{\theta} \sum_{p \in \mathcal{L}} l(\mathbf{y}^p, \mathbf{s}_{\theta}^p) + \sum_{p,q \in \mathcal{L} \cup \mathcal{U}} w_{p,q} \|\mathbf{s}_{\theta}^p - \mathbf{s}_{\theta}^q\|$$



[Tang et al., On regularized losses for weakly supervised segmentation, ECCV 2018]

Semi-supervision loss in segmentation

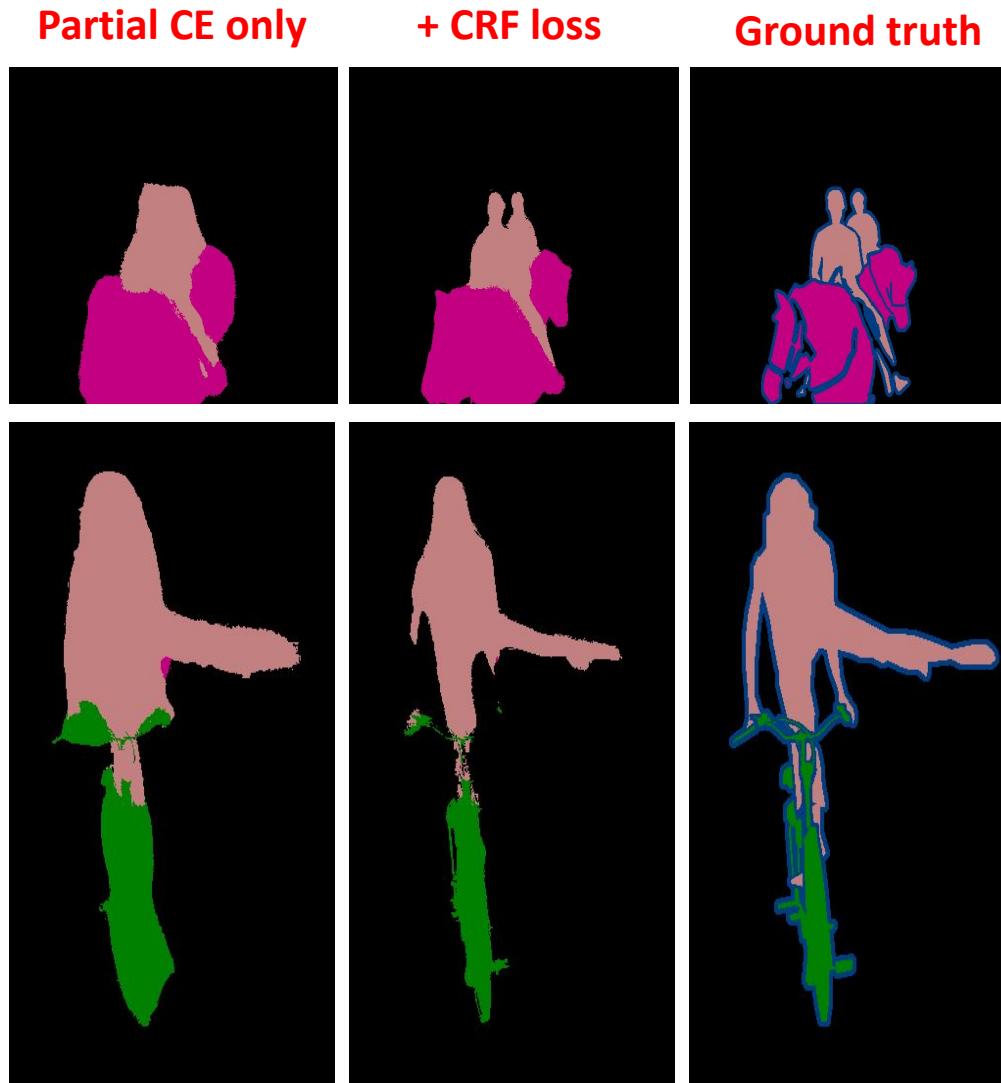
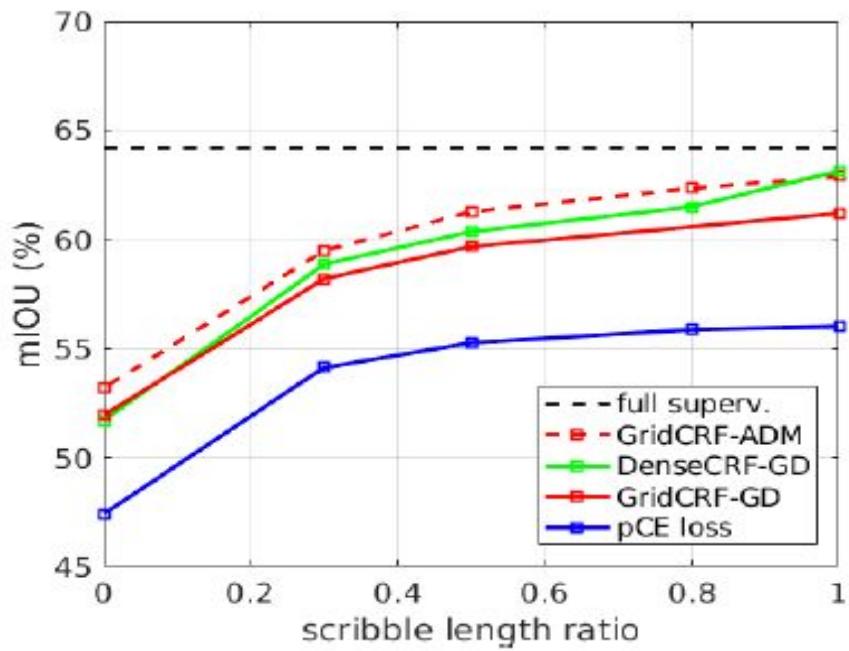
$$\min_{\theta} \sum_{p \in \mathcal{L}} l(\mathbf{y}^p, \mathbf{s}_{\theta}^p) + \sum_{p,q \in \mathcal{L} \cup \mathcal{U}} w_{p,q} \|\mathbf{s}_{\theta}^p - \mathbf{s}_{\theta}^q\|$$



On the vertices of the simplex (binary variables), this is exactly the popular **Potts model in CRFs**

[Tang et al., On regularized losses for weakly supervised segmentation, ECCV 2018]

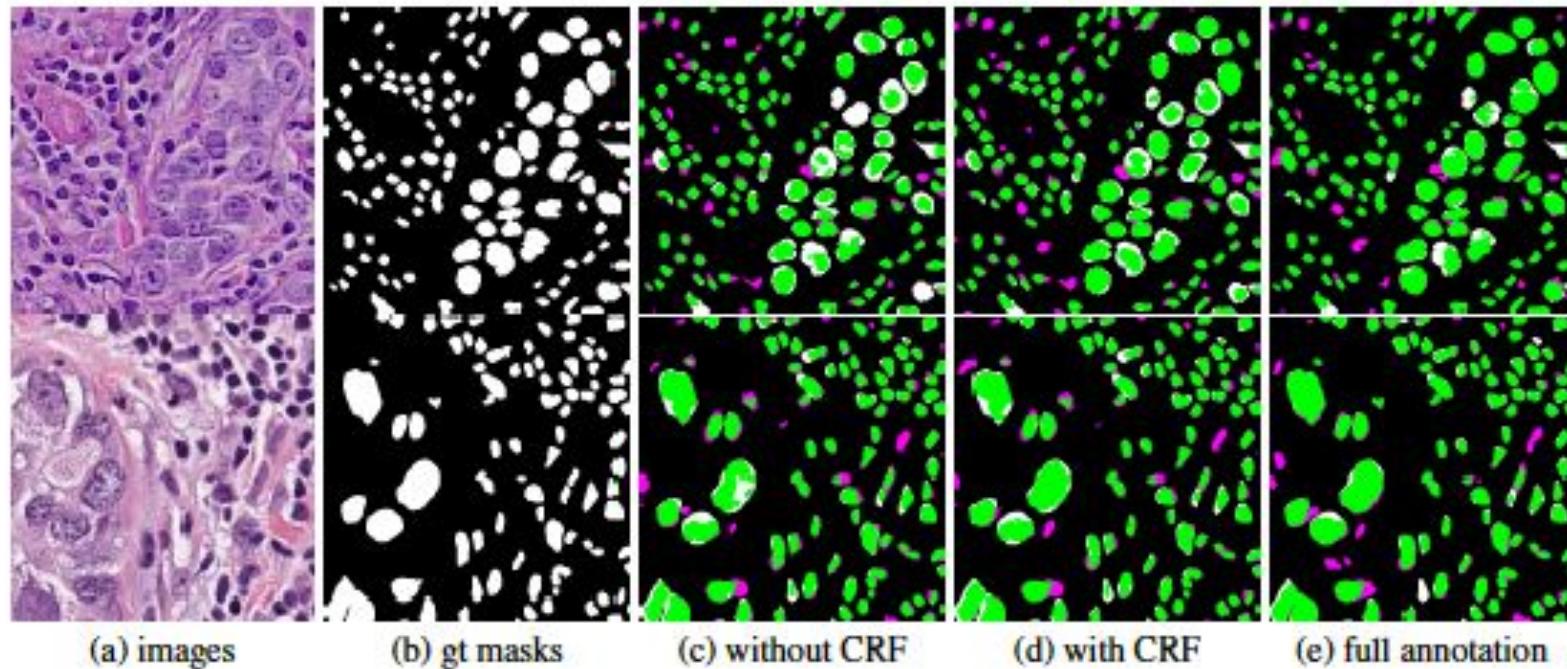
Semi-supervision loss in segmentation



[Tang et al., On regularized losses for weakly supervised segmentation,
ECCV 2018]

Quite used in medical imaging

White (FN); Magenta (FP); Green (TP)

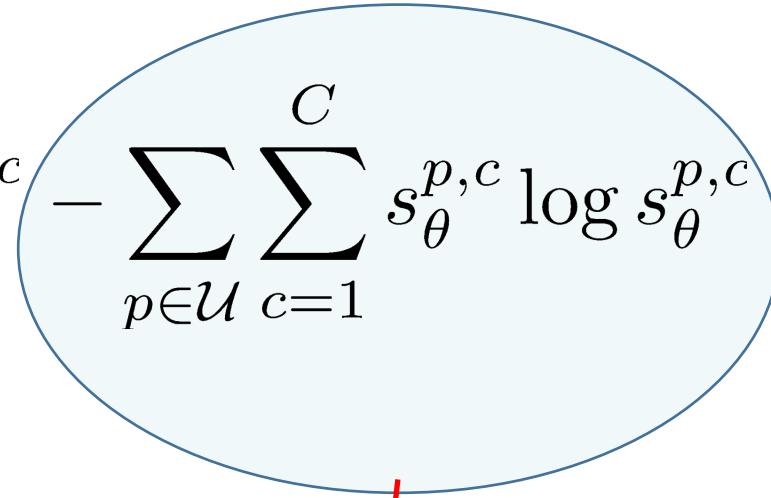


- Figures from Qu et al., Weakly Supervised Deep Nuclei Segmentation using Points Annotation in Histopathology Images, MIDL 2019 [\[Histology, point annotation\]](#)
- Ji et al., Scribble-Based Hierarchical Weakly Supervised Learning for Brain Tumor Segmentation, MICCAI 2019
[Brain tumor images, scribble annotations]

Unsupervised Entropy Regularization

Entropy minimization in SSL

$$\min_{\theta} - \sum_{p \in \mathcal{L}} \sum_{c=1}^C y^{p,c} \log s_{\theta}^{p,c} - \sum_{p \in \mathcal{U}} \sum_{c=1}^C s_{\theta}^{p,c} \log s_{\theta}^{p,c}$$



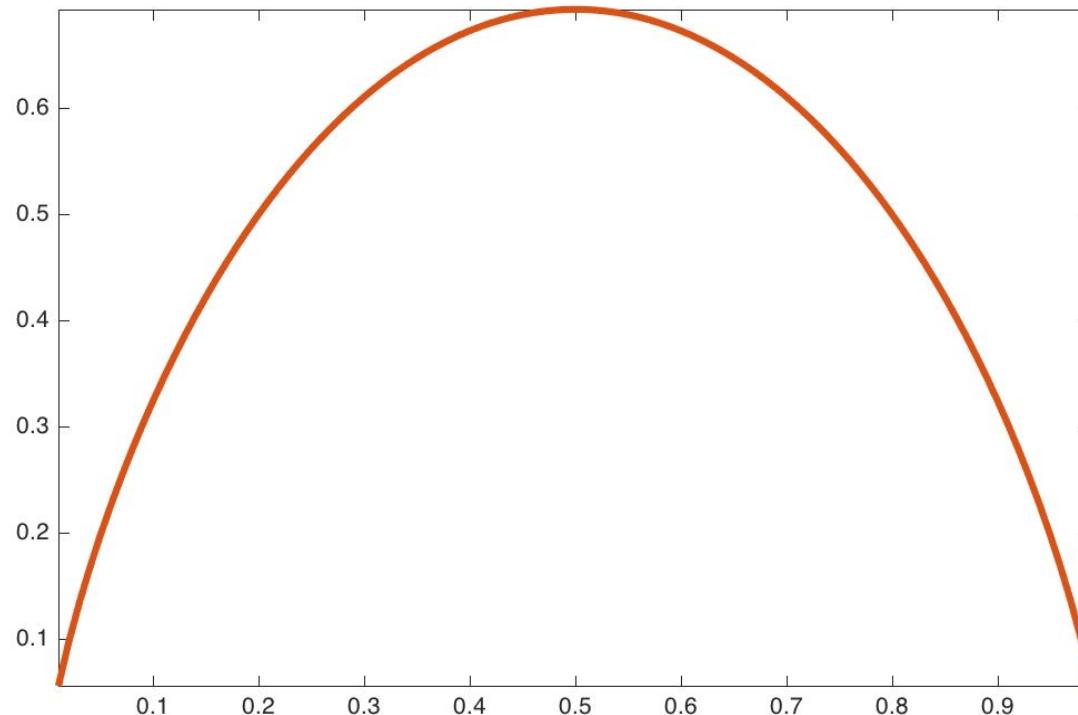
Shannon Entropies: “unsupervised cross-entropies (with unknown labels)”

- Grandvalet & Bengio, Semi-supervised learning by entropy minimization, NIPS 2005
- Gomes et al., Discriminative clustering by regularized information maximization, NIPS 2010

Effect of the Unsupervised Entropy (Why Is It Good for SSL)

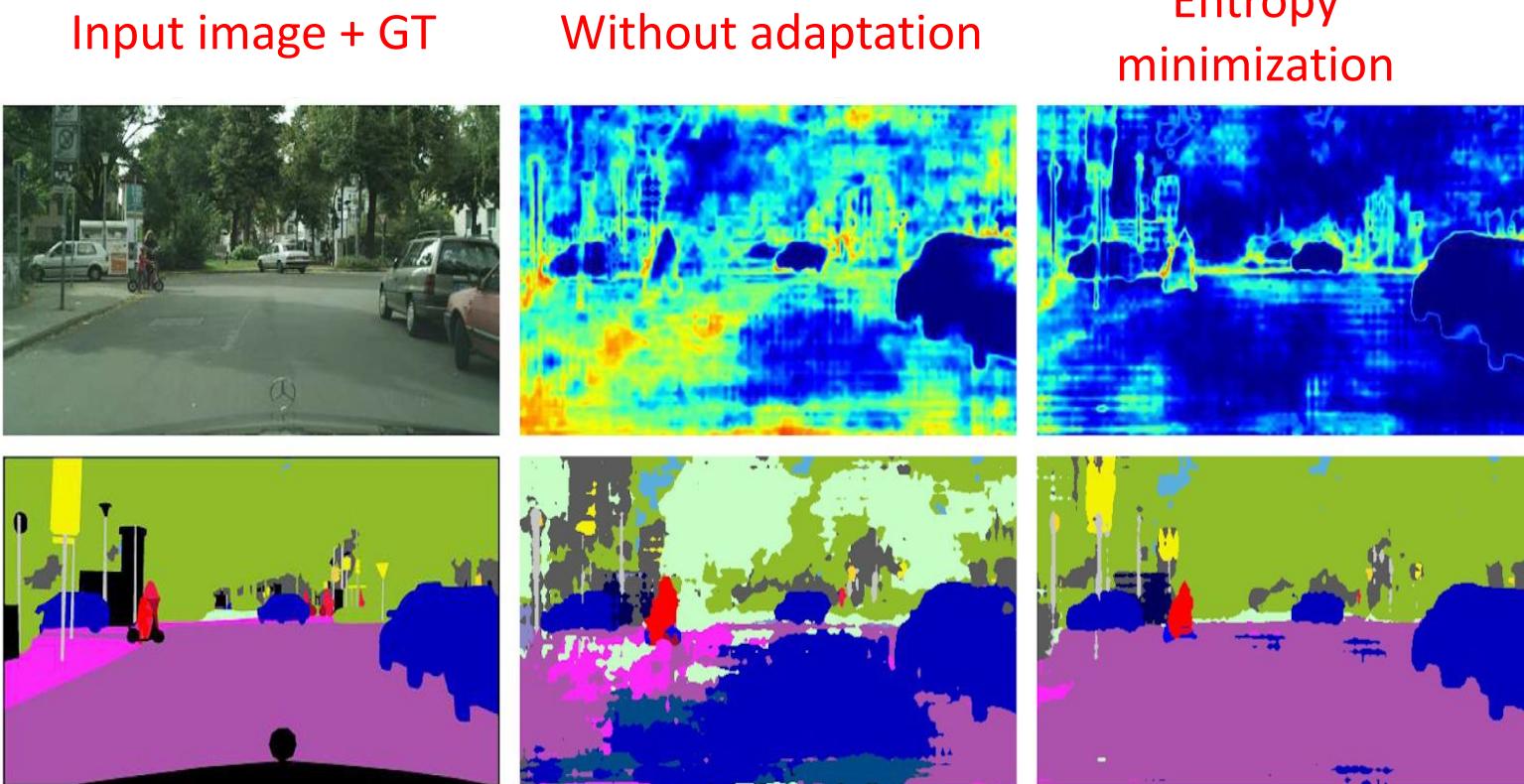
It makes predictions confident (just like the cross-entropy)

$$-s_\theta^p \log s_\theta^p - (1 - s_\theta^p) \log(1 - s_\theta^p)$$



Entropy Minimization in UDA

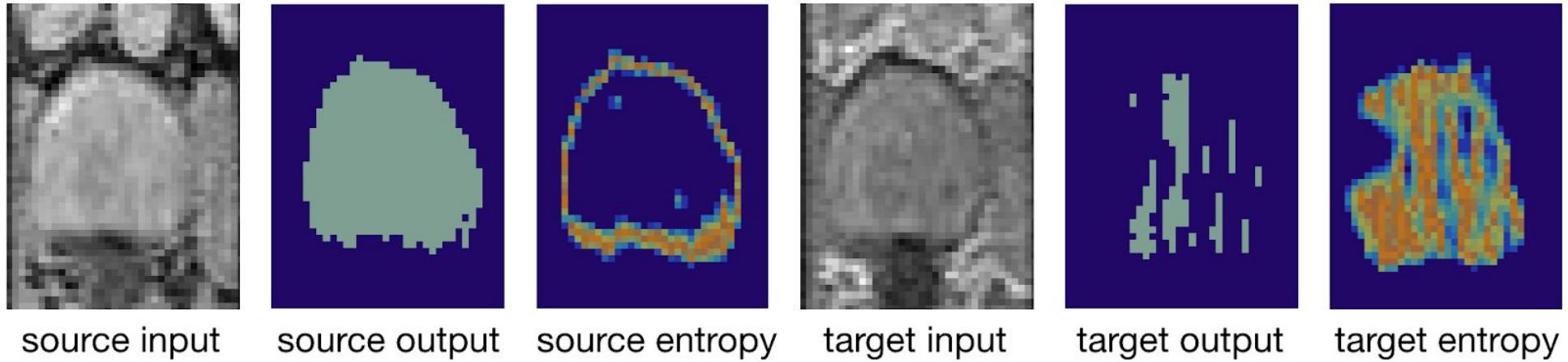
It makes predictions confident (just like the cross-entropy)



Images from Vu et al., ADVENT: Adversarial Entropy Minimization for Domain Adaptation in Semantic Segmentation, CVPR 2019

Entropy Minimization in UDA

It makes predictions confident (just like the cross-entropy)



source input

source output

source entropy

target input

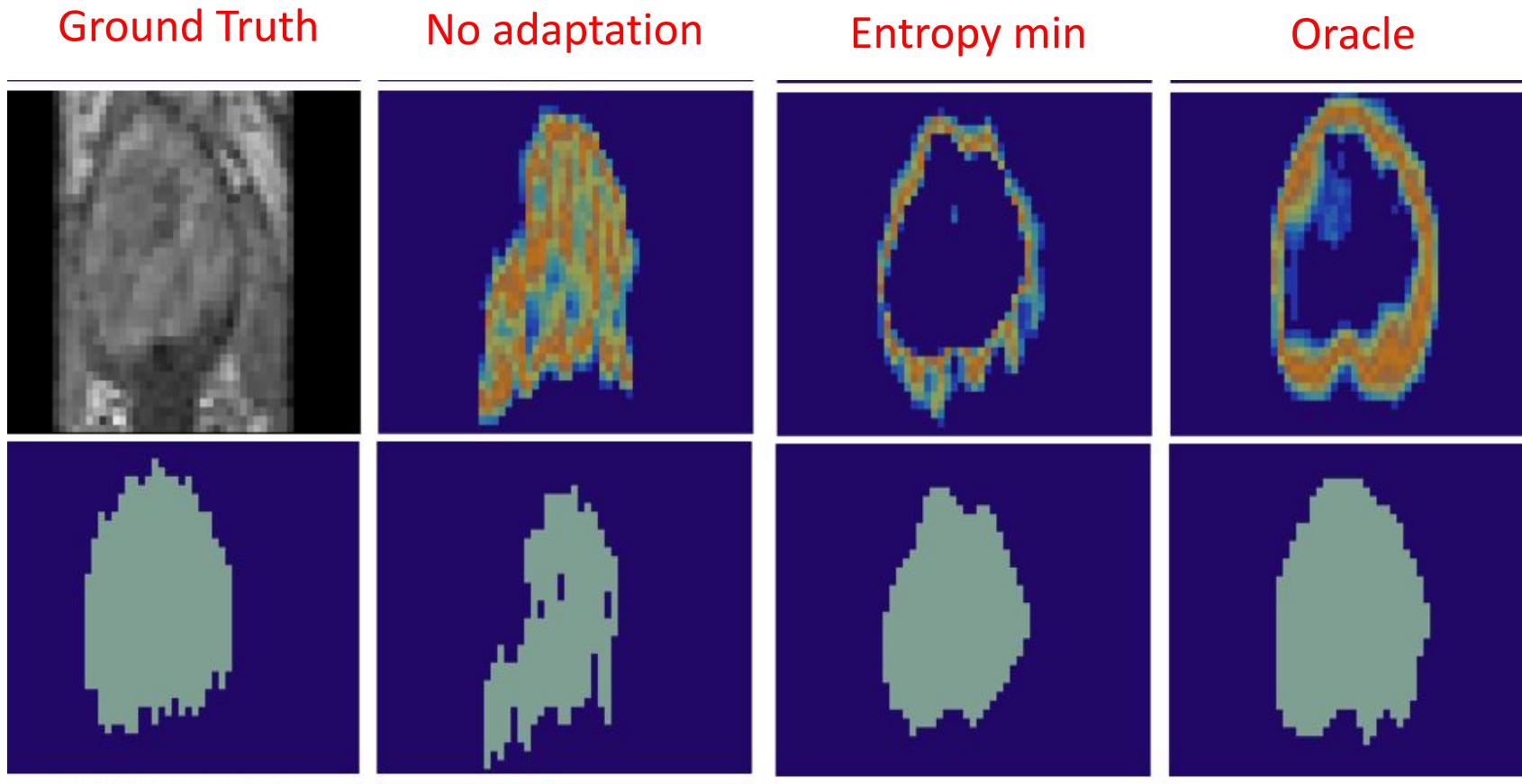
target output

target entropy

Images from Bateson et al., Source-relaxed domain adaptation for segmentation, MICCAI 2020

Entropy Minimization in UDA

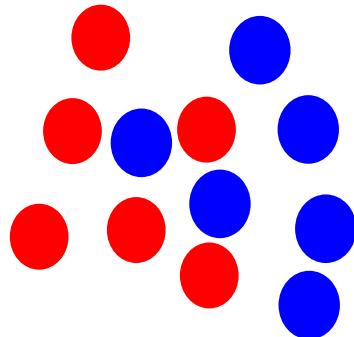
It makes predictions confident (just like the cross-entropy)



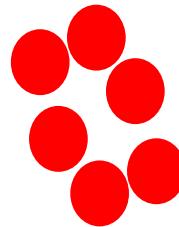
Images from Bateson et al., Source-relaxed domain adaptation for segmentation, MICCAI 2020

Why Entropy Minimization is good

It increases the margin between the classes



*High entropy
(low confidence)*



*Low entropy
(high confidence)*

PAUSE (30 mins)

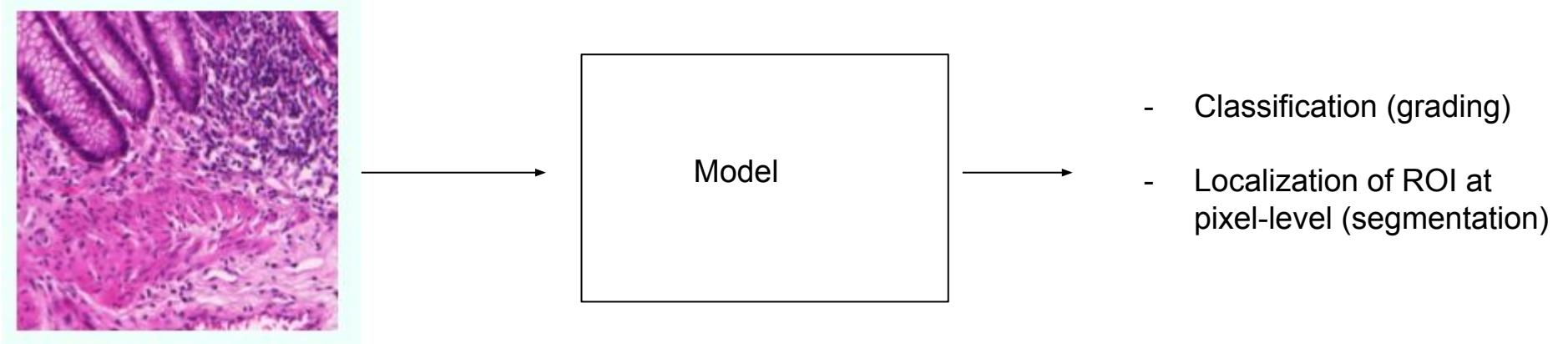
Part 4:

Applications of WSOL / WSSS

- (a) Medical Cancer Grading and ROI Localization in Histology
- (b) Weakly-Supervised Video Object Localization
- (c) Person ReID: Embedding Networks
- (d) Medical Semantic Segmentation

(a) Cancer Grading and Localization in Histology:

- Task

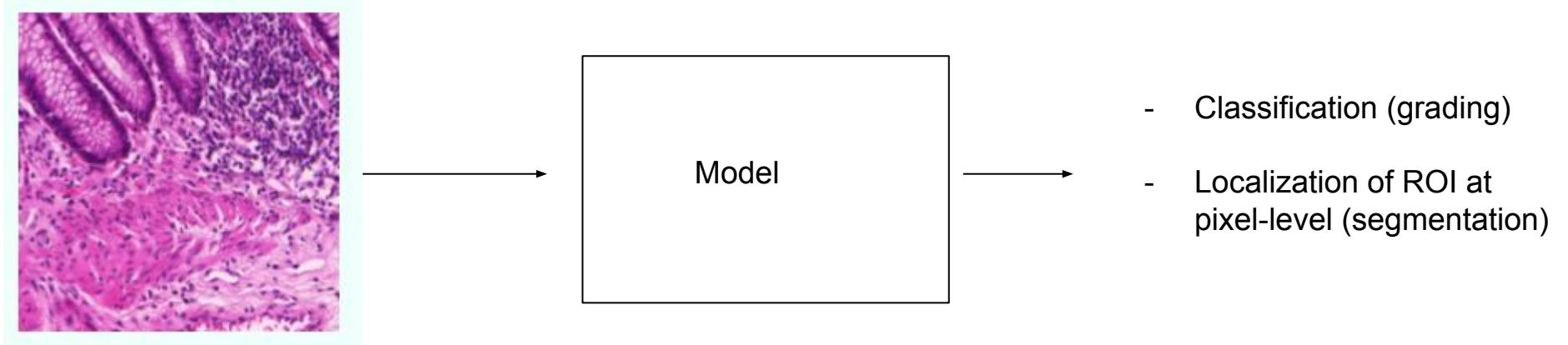


Input: histology image

Microscopic data for
cancer diagnostic

(a) Cancer Grading and Localization in Histology:

- Task



Input: histology image

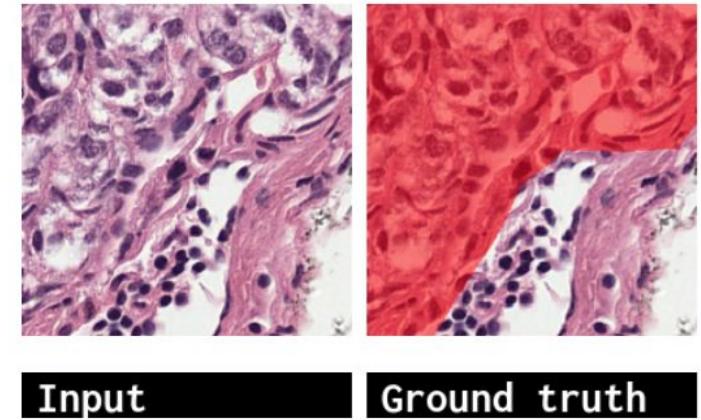
Microscopic data for
cancer diagnostic

Available supervision: global image grad (class)

(a) Cancer Grading and Localization in Histology:

- **Histology data challenges**

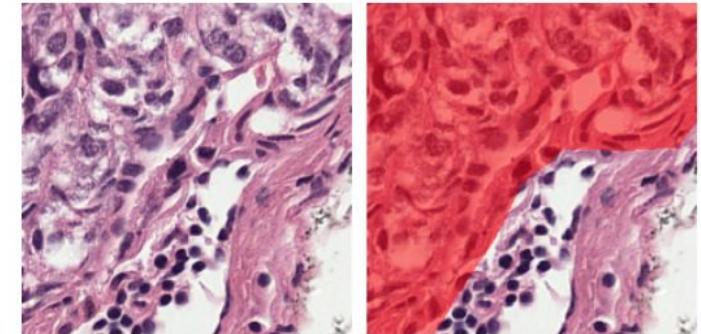
- Object non-salient (foreground is similar to background)



(a) Cancer Grading and Localization in Histology:

- **Histology data challenges**

- Object non-salient (foreground is similar to background)
- ROI with arbitrary shapes – no common global structure – opposite to natural scene images.



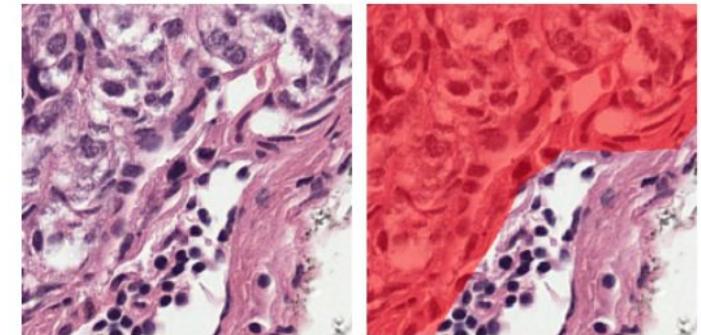
Input

Ground truth

(a) Cancer Grading and Localization in Histology:

- **Histology data challenges**

- Object non-salient (foreground is similar to background)
- ROI with arbitrary shapes – no common global structure – opposite to natural scene images.
- Stain variation – Hematoxylin and Eosin (H&E)

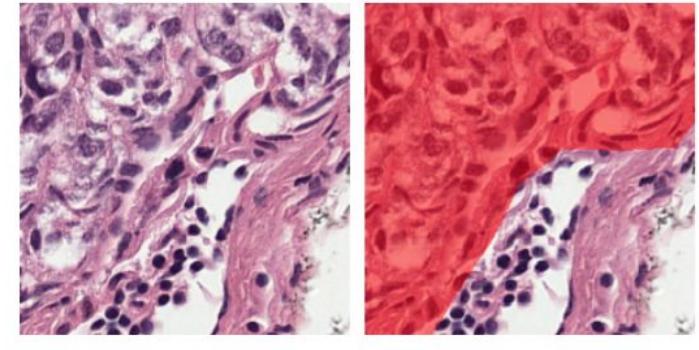


Input **Ground truth**

(a) Cancer Grading and Localization in Histology:

- Presented work

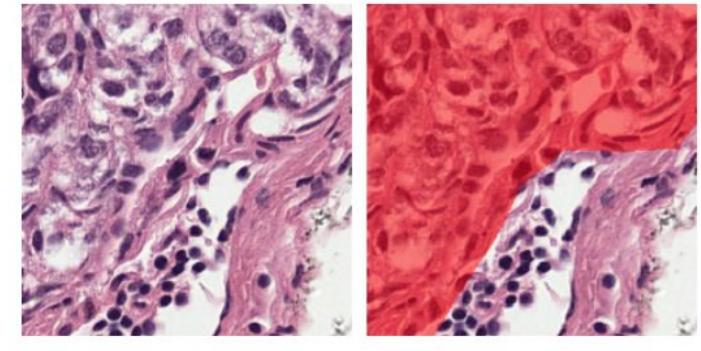
- Belharbi, S., Rony, J., Dolz, J., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty.** IEEE Transactions on Medical Imaging, 41:702–714.
- Belharbi, S., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Negative evidence matters in interpretable histology image classification.** In Medical Imaging with Deep Learning (MIDL).



(a) Cancer Grading and Localization in Histology:

- Presented work

- Belharbi, S., Rony, J., Dolz, J., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty.** IEEE Transactions on Medical Imaging, 41:702–714.
- Belharbi, S., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Negative evidence matters in interpretable histology image classification.** In Medical Imaging with Deep Learning (MIDL).



(a) Cancer Grading and Localization in Histology:

Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty

Transactions on Medical Imaging, 2022

Code:

<https://github.com/sbelharbi/deep-wsl-histo-min-max-uncertainty>

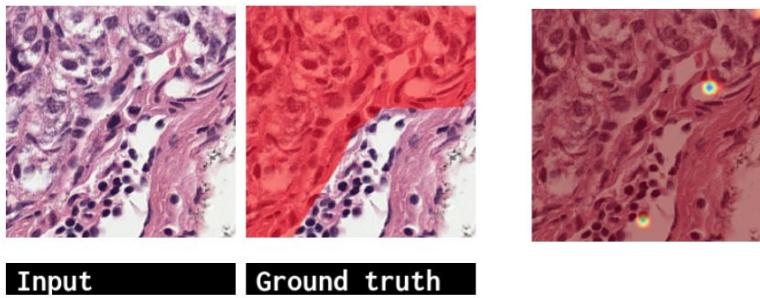
(a) Cancer Grading and Localization in Histology:

- Issue

- Non-salient object → visual similarity between foreground/background



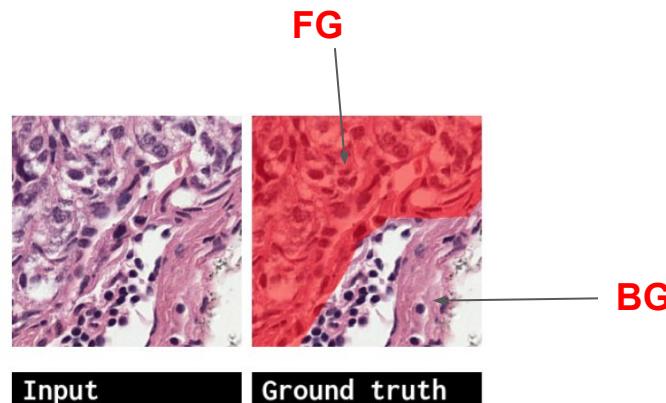
High false positives/negatives



(a) Cancer Grading and Localization in Histology:

- Constrain CAMs

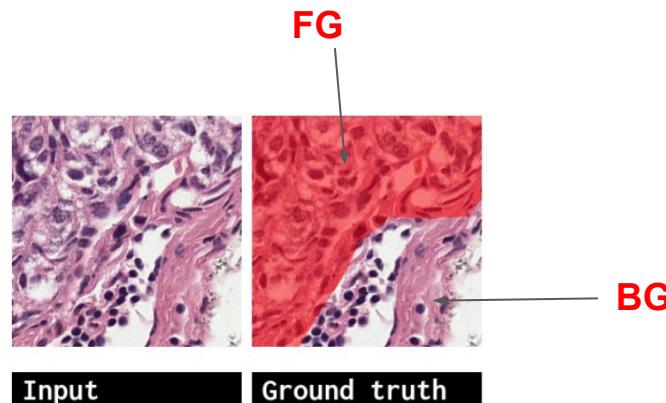
- Explicit modeling foreground/background map



(a) Cancer Grading and Localization in Histology:

- Constrain CAMs

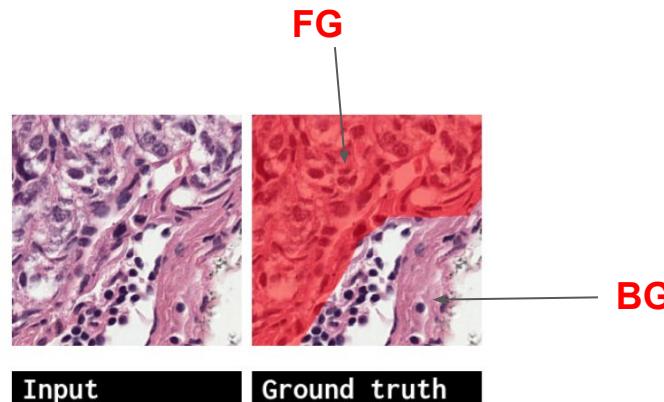
- Explicit modeling foreground/background map
- Constrain the presence of both FG / BG using size constraints



(a) Cancer Grading and Localization in Histology:

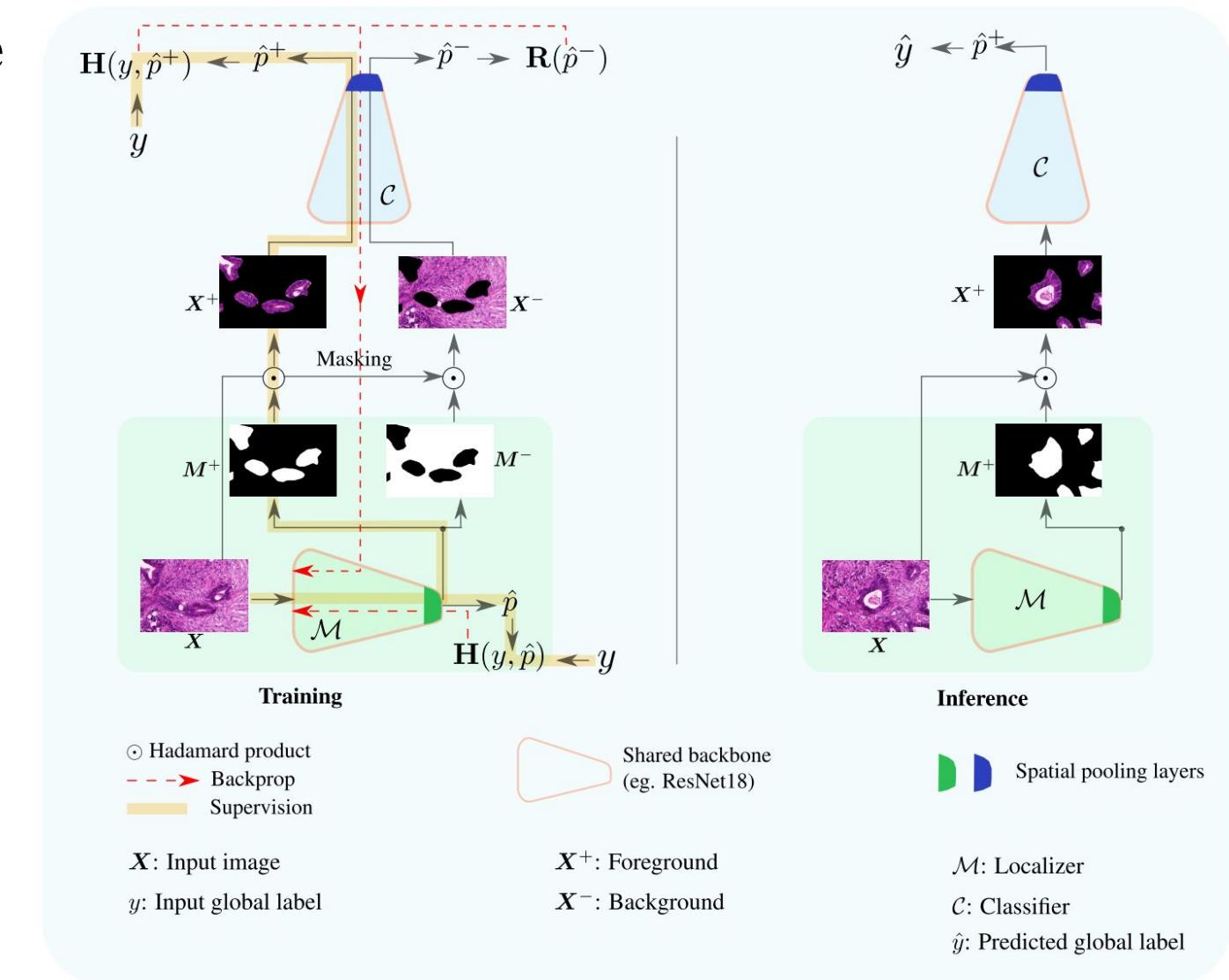
- Constrain CAMs

- Explicit modeling foreground/background map
- Constrain the presence of both FG / BG using size constraints
- Ensure that each map is consistent using classifier response.



(a) Cancer Grading and Localization in Histology:

- Our architecture



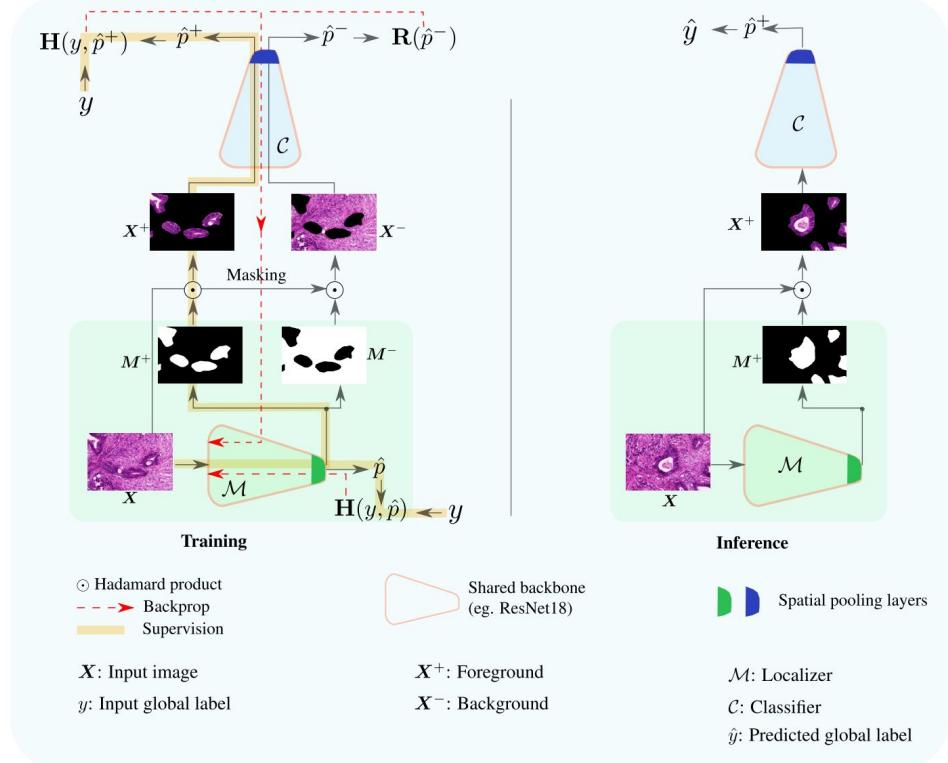
(a) Cancer Grading and Localization in Histology:

- Training loss

$$\min_{\theta_C} \quad \mathbf{H}(p, \hat{p}^+) + \lambda \mathbf{R}(\hat{p}^-) - \frac{1}{t} [\log s^+ + \log s^-],$$



Maximize classifier
Response over FG



(a) Cancer Grading and Localization in Histology:

- Training loss

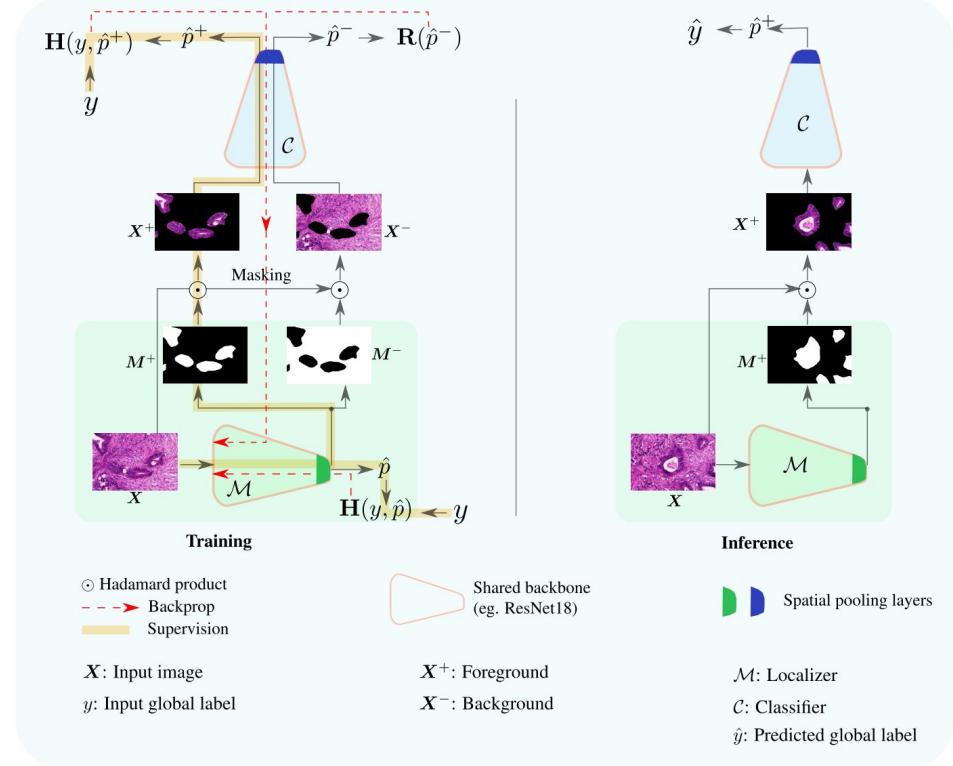
$$\min_{\theta_C} \mathbf{H}(p, \hat{p}^+) + \lambda \mathbf{R}(\hat{p}^-) - \frac{1}{t} [\log s^+ + \log s^-],$$

The BG has no
discriminative regions
left.
Max uncertainty

$$\mathbf{R}(\hat{p}^-) = -\mathbf{H}(\hat{p}^-); \quad \text{or} \quad \mathbf{R}(\hat{p}^-) = \mathbf{H}(q, \hat{p}^-),$$

**Explicit Entropy
Maximization (EEM)**

**Surrogate for explicit Entropy
Maximization (SEM).** q: uniform dist.



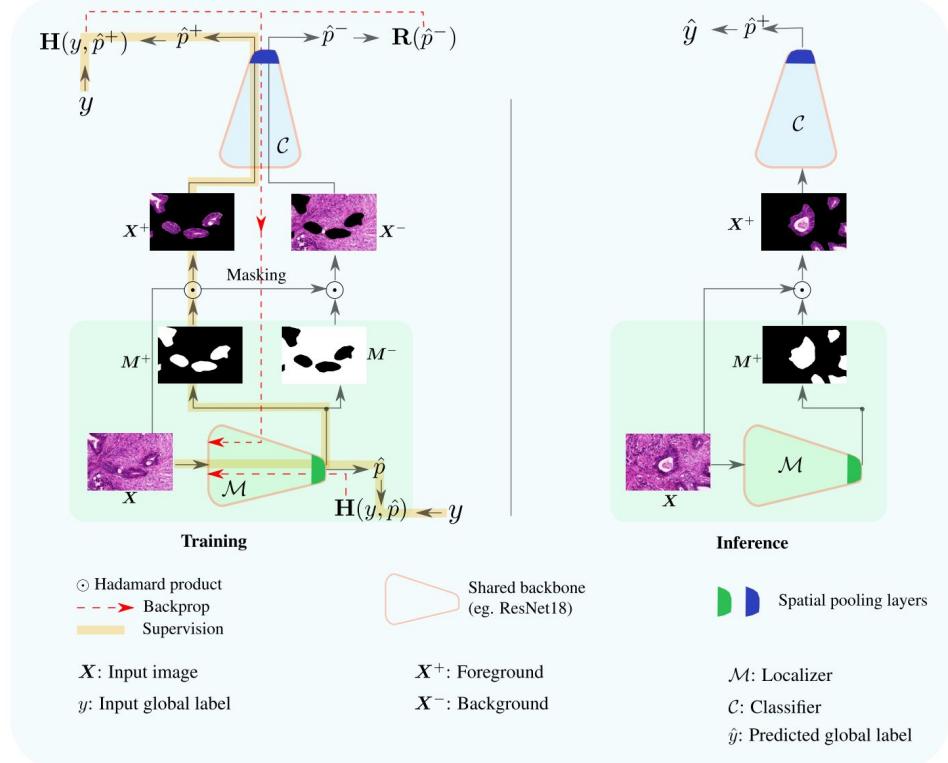
(a) Cancer Grading and Localization in Histology:

- Training loss

ASC: Absolute Size Constraint

$$\min_{\theta_C} \mathbf{H}(p, \hat{p}^+) + \lambda \mathbf{R}(\hat{p}^-) - \frac{1}{t} [\log s^+ + \log s^-],$$

Ensure both FG/BG are present: max size.



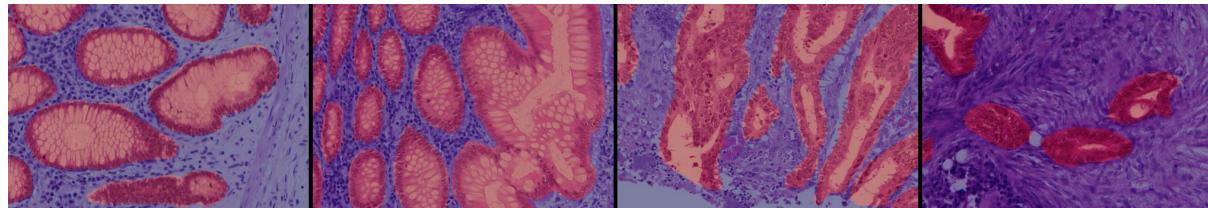
$$s^+ = \sum_{z \in \Omega} M^+(z), \quad s^- = \sum_{z \in \Omega} M^-(z)$$

Sizes

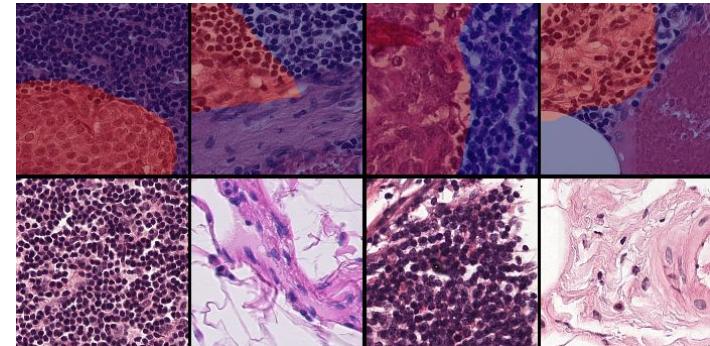
(a) Cancer Grading and Localization in Histology:

- **Experiments**

- Task: classify and localize ROI
- 2 public datasets: GlaS, Camelyon16 patches.



GlaS: colon cancer diagnosis



Camelyon16 patches: breast cancer

(a) Cancer Grading and Localization in Histology:

• Results

Method	Image level		Pixel level	
	Cl. error (%)	F1 ⁺ (%)	F1 ⁻ (%)	
All-ones (Lower-bound)	--	66.01	00.00	
PN [39]	--	65.52	24.08	
ERASE [80]	7.50	65.60	25.01	
CAM-Max [52]	1.25	66.00	26.32	
CAM-LSE [57, 74]	1.25	66.05	27.93	
Grad-CAM [64]	0.00	66.30	21.30	
CAM-Avg [88]	0.00	66.90	17.88	
Wildcat [20]	1.25	67.21	22.96	
Deep MIL [33]	2.50	68.52	41.34	
Ours (EEM)	0.00	72.11	69.07	
Ours (SEM)	0.00	71.94	69.23	
U-Net [60] (Upper-bound)	--	90.19	88.52	

GlaS

Method	Image level		Pixel level	
	Cl. error (%)	F1 ⁺ (%)	F1 ⁻ (%)	
All-ones (Lower-bound)	--	59.44	00.00	
PN [39]	--	31.15	37.36	
ERASE [80]	8.61	31.30	42.48	
CAM-Max [52]	10.06	48.28	81.92	
CAM-LSE [57, 74]	1.51	64.31	63.78	
Grad-CAM [64]	2.40	62.78	79.05	
CAM-Avg [88]	2.40	62.75	79.05	
Wildcat [20]	1.48	62.73	72.59	
Deep MIL [33]	1.93	59.01	36.94	
Ours (EEM)	6.26	67.98	88.80	
Ours (SEM)	6.95	68.26	88.55	
U-Net [60] (Upper-bound)	--	71.11	89.68	

Camelyon16

(a) Cancer Grading and Localization in Histology:

- Visual results

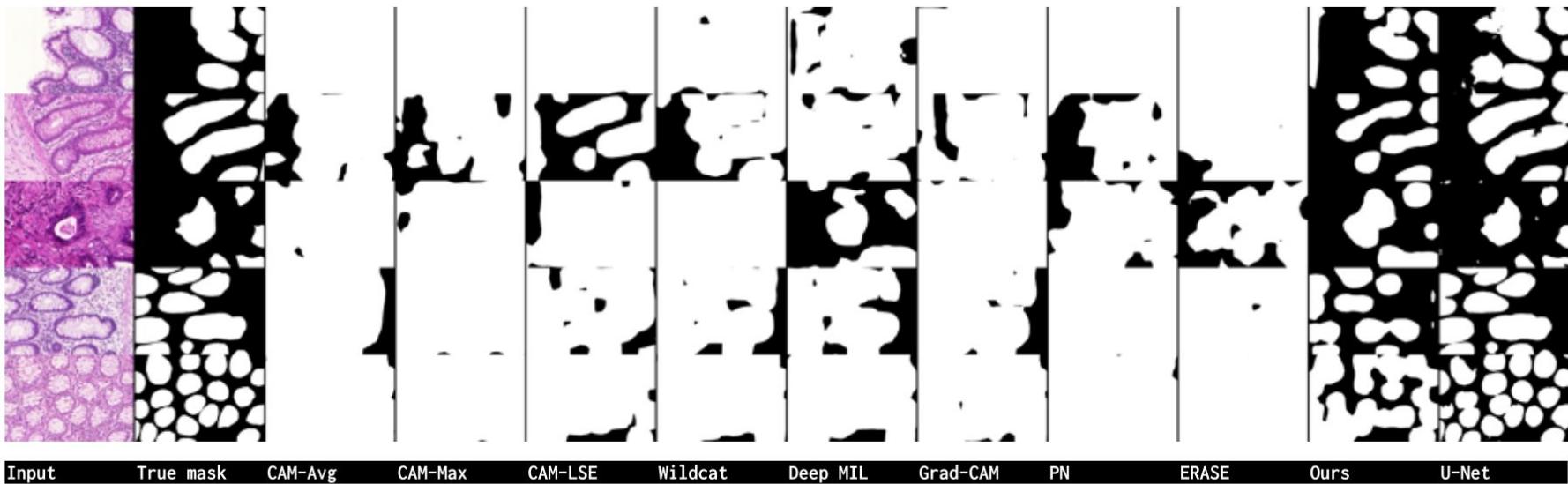


Figure 2: **GlaS dataset:** Qualitative results of the predicted binary mask for each method on several GlaS test images. Our method, referred to as *Ours*, is the SEM version with the ASC regularization term. (Best visualized in color.)

(a) Cancer Grading and Localization in Histology:

- Visual results

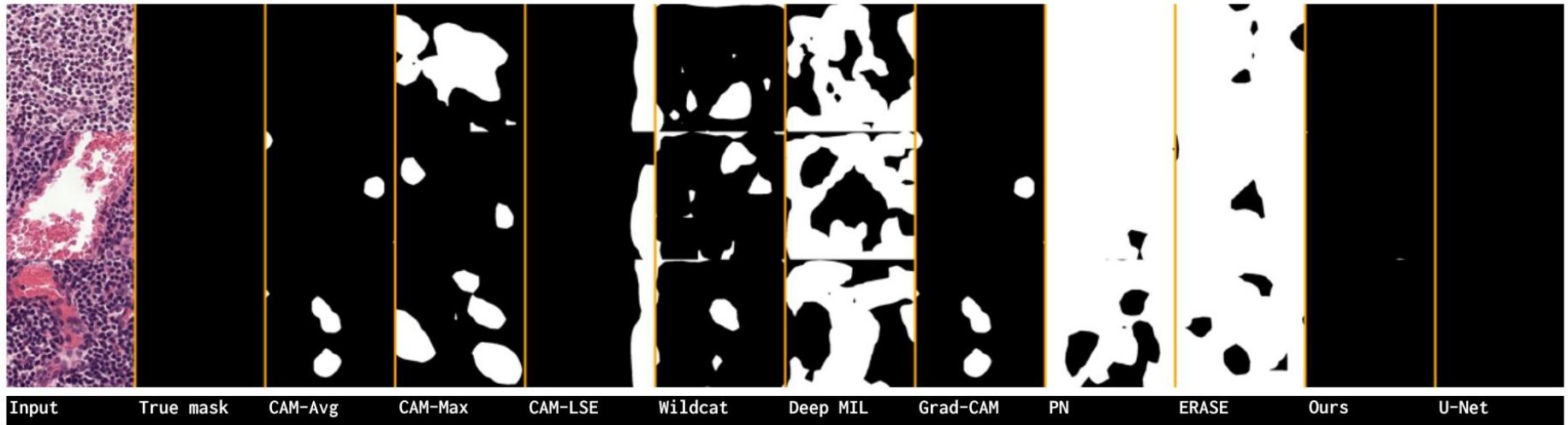
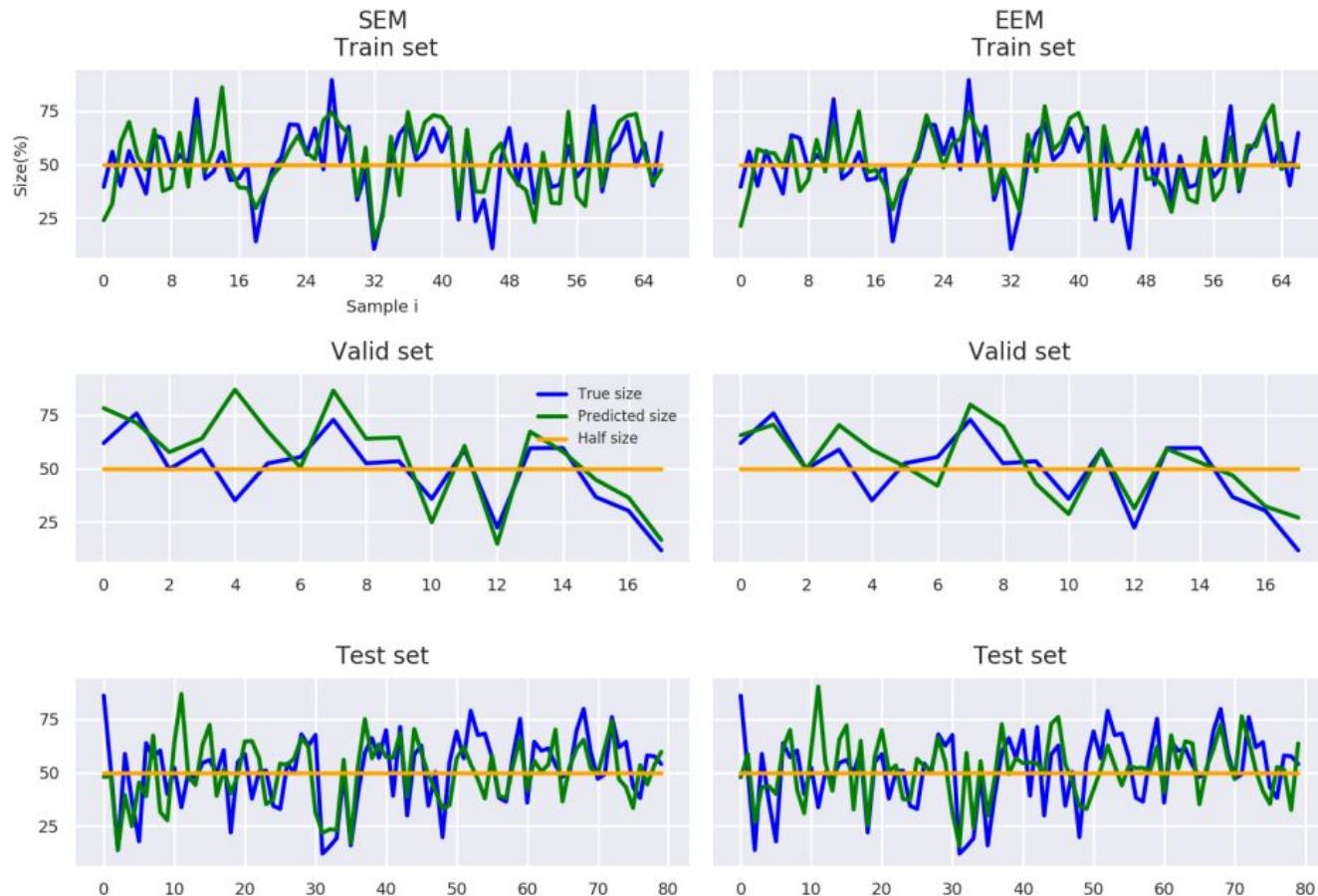


Figure 3: **Camelyon16-P512 benchmark:** Examples of mask predictions over **normal** samples from the testing set. White pixels indicate metastatic regions, while black pixels indicate normal tissue. This illustrates false positives. Note that normal samples do not contain any metastatic regions. **Ours** is SEM version with the ASC regularization. (Best visualized in color.)

(a) Cancer Grading and Localization in Histology:

- Size constraint

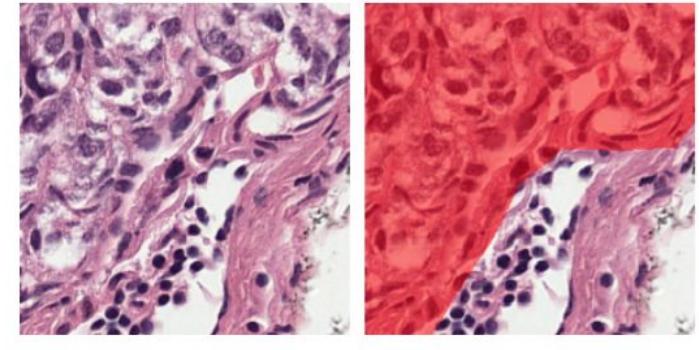


$$\min_{\theta_c} \mathbf{H}(p, \hat{p}^+) + \lambda \mathbf{R}(\hat{p}^-) - \frac{1}{t} [\log s^+ + \log s^-],$$

(a) Cancer Grading and Localization in Histology:

- Presented work

- Belharbi, S., Rony, J., Dolz, J., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty.** IEEE Transactions on Medical Imaging, 41:702–714.
- Belharbi, S., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Negative evidence matters in interpretable histology image classification.** In Medical Imaging with Deep Learning (MIDL).



(a) Cancer Grading and Localization in Histology:

Negative evidence matters in interpretable histology image classification

Medical Imaging with Deep Learning (MIDL), 2022

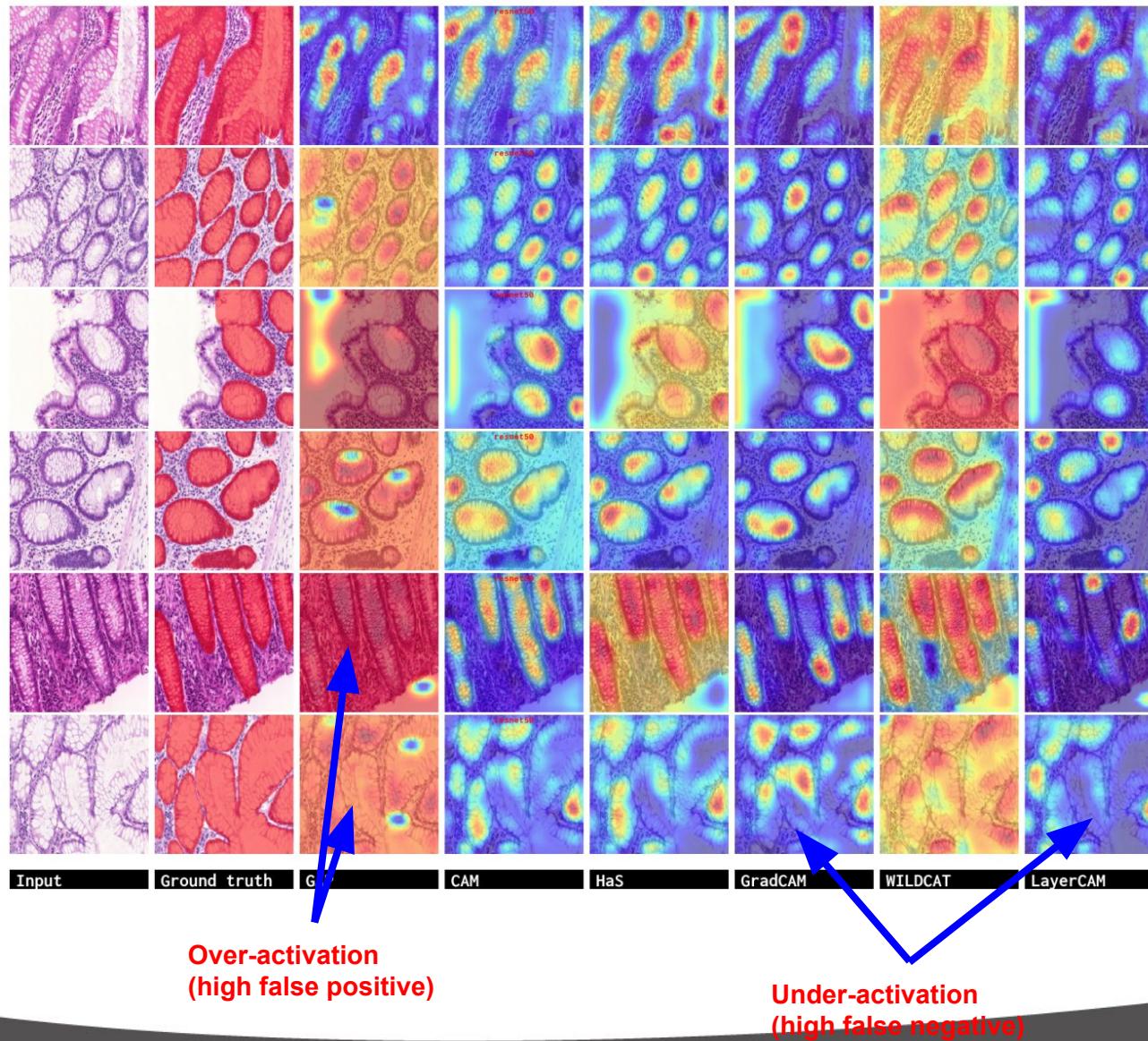
Code:

<https://github.com/sbelharbi/negev>



(a) Cancer Grading and Localization in Histology:

- CAMs' challenges in histology images



(a) Cancer Grading and Localization in Histology:

- **Using negative knowledge**

- To reduce mis-predictions, **guide the CAM learning with available Negative knowledge.**

Negative knowledge = all what is not ROI.

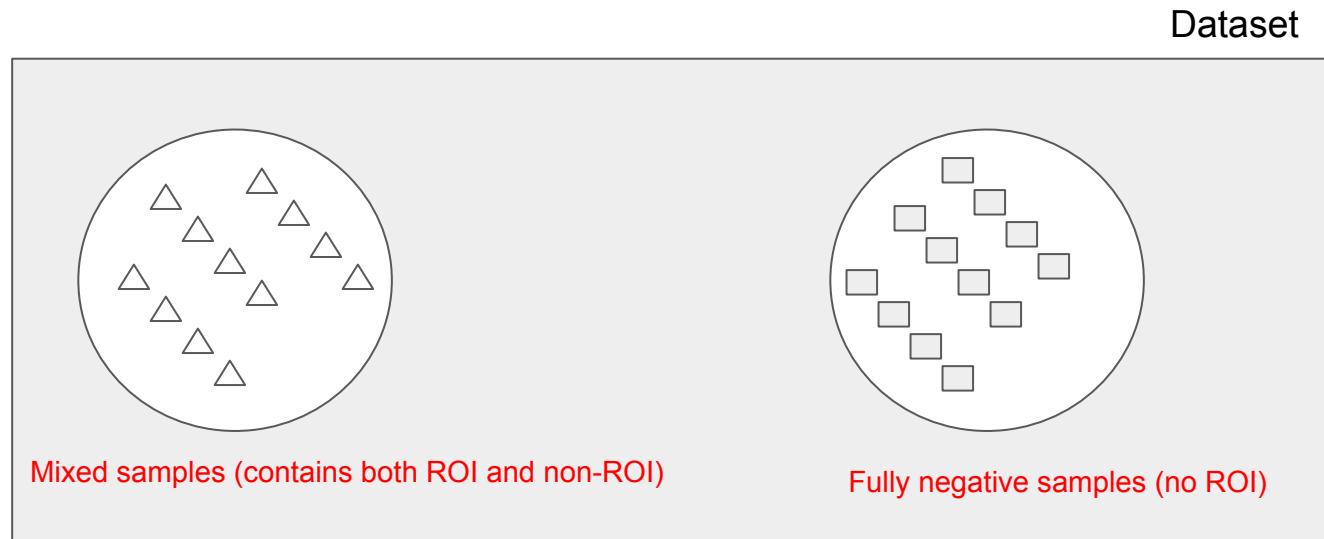
(a) Cancer Grading and Localization in Histology:

- **Using negative knowledge**
 - 2 sources of negative knowledge

(a) Cancer Grading and Localization in Histology:

- Using negative knowledge

- 2 sources of negative knowledge
 - 1 - Naturally occurring in dataset

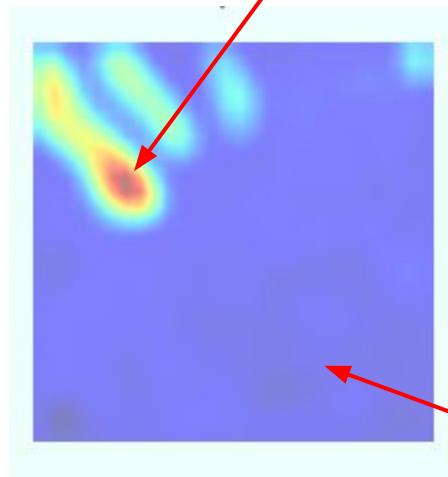
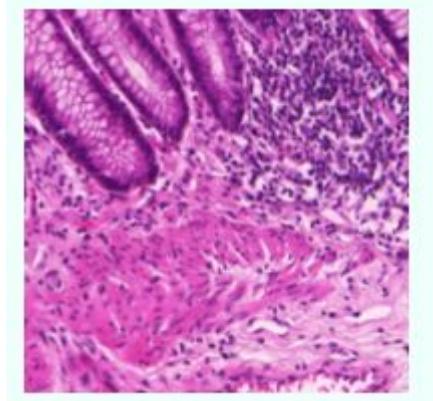


(a) Cancer Grading and Localization in Histology:

- Using negative knowledge

- 2 sources of negative knowledge

2 - Low activation in CAMs



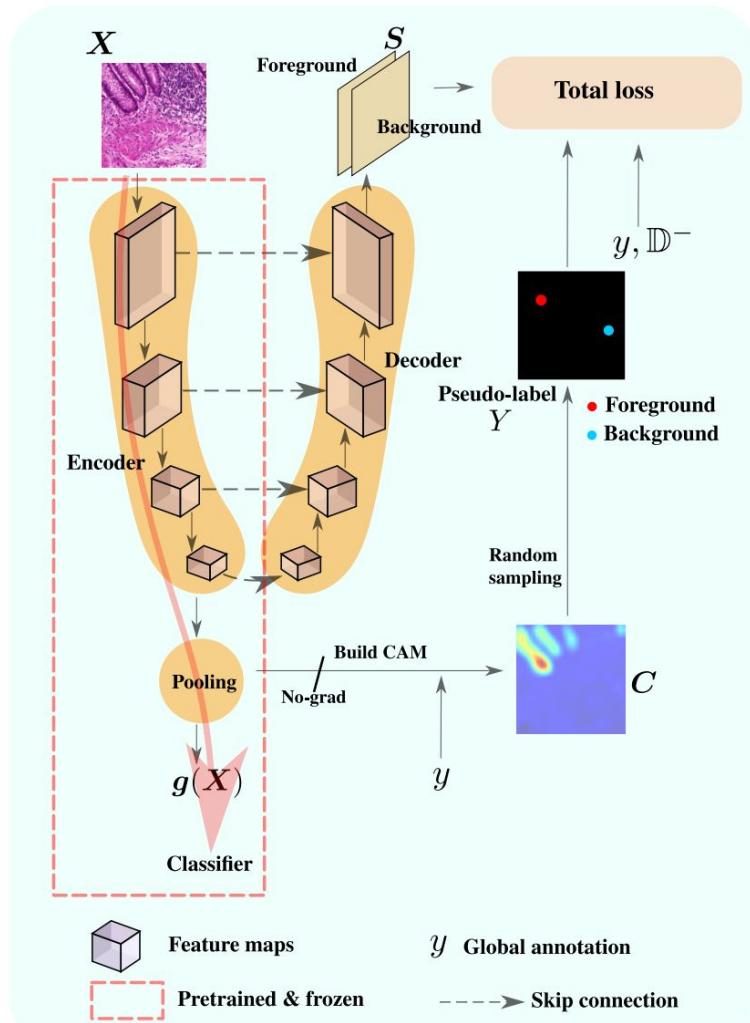
More likely foreground

More likely background

(a) Cancer Grading and Localization in Histology:

- **Architecture**

- Requires only image class for training

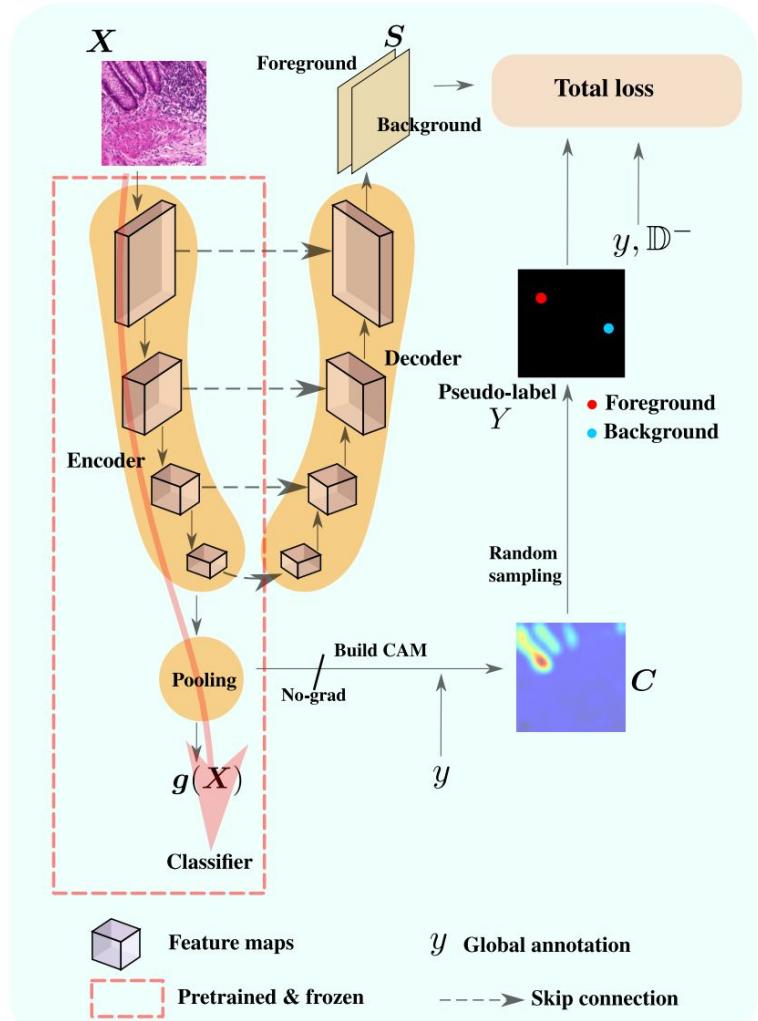


(a) Cancer Grading and Localization in Histology:

- Training

1- Exploit CAM positive/negative information

$$\min_{\theta} \sum_{p \in \{\mathbb{C}^+ \cup \mathbb{C}^-\}} H(Y_p, S_p) .$$

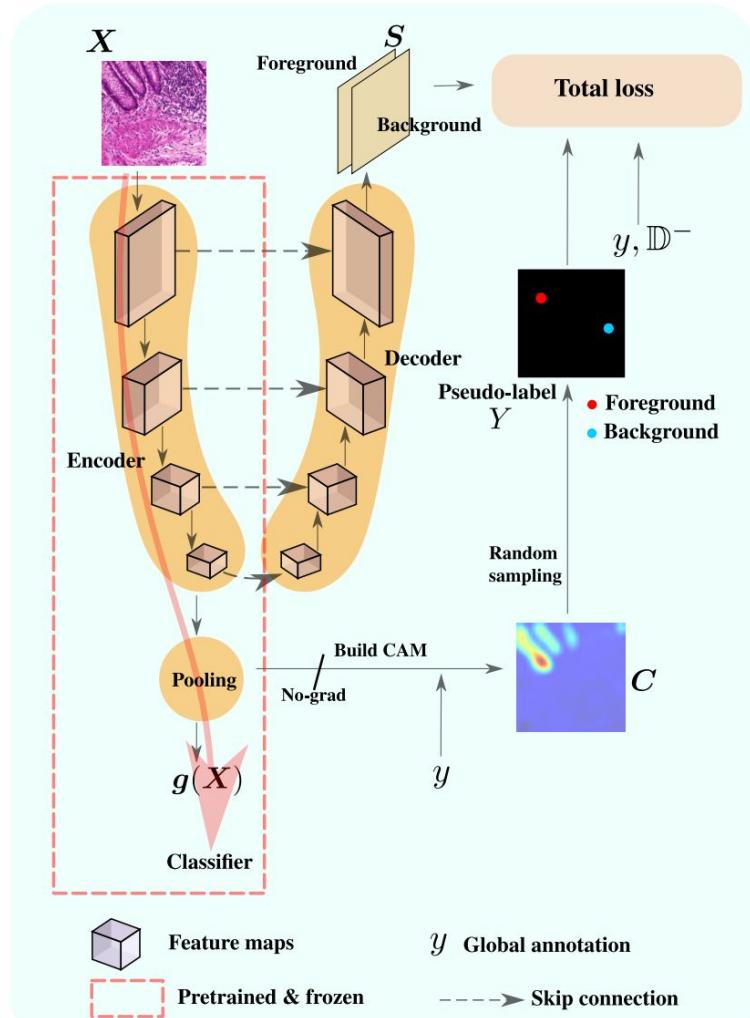


(a) Cancer Grading and Localization in Histology:

- Training

2- Fully negative samples

$$\min_{\theta} \sum_{p \in \Omega} -\log(1 - S_p^0), \forall X \in \mathbb{D}^- .$$

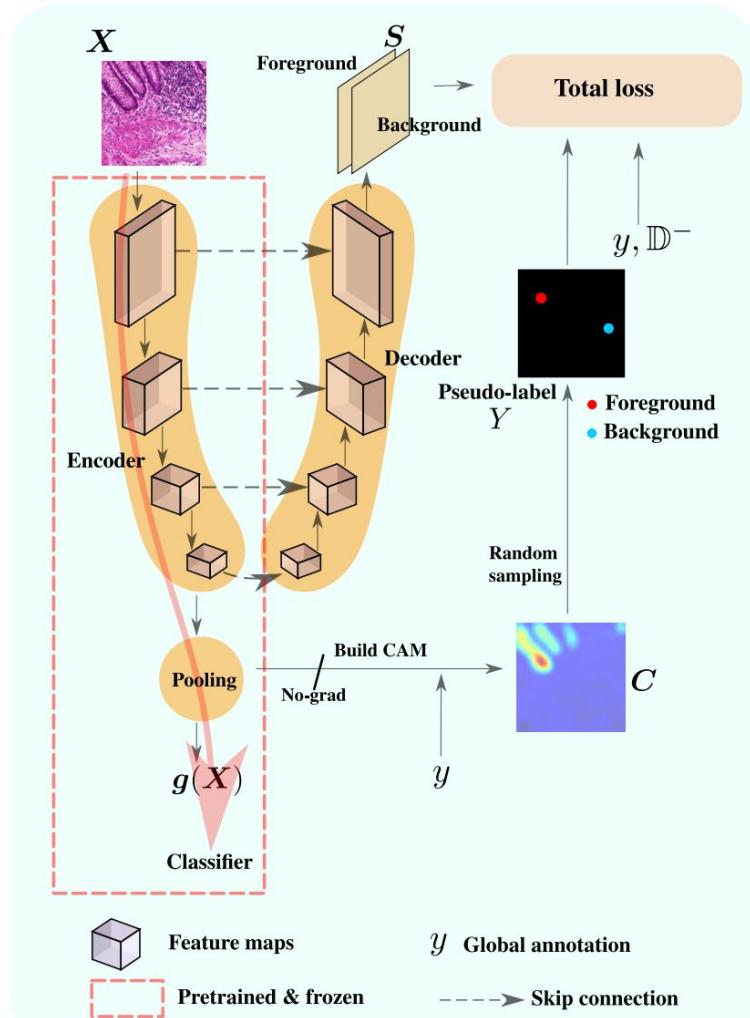


(a) Cancer Grading and Localization in Histology:

- Training

Total adaptive loss

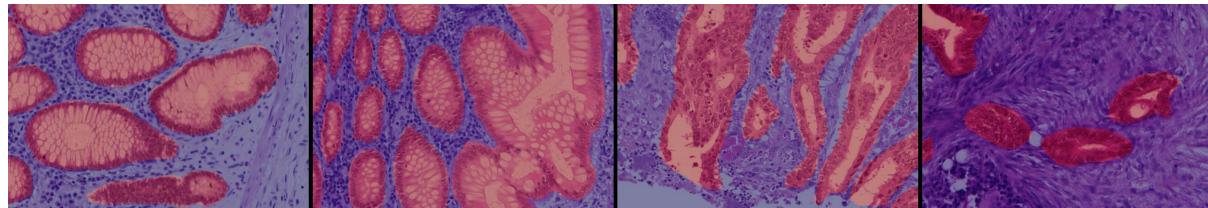
$$\min_{\theta} \quad \mathbb{1}_{X \in \mathbb{D}^-} \left(\sum_{p \in \Omega} -\log(1 - S_p^0) \right) + (1 - \mathbb{1}_{X \in \mathbb{D}^-}) \left(\lambda \sum_{p \in \{\mathbb{C}^+ \cup \mathbb{C}^-\}} H(Y_p, S_p) \right),$$



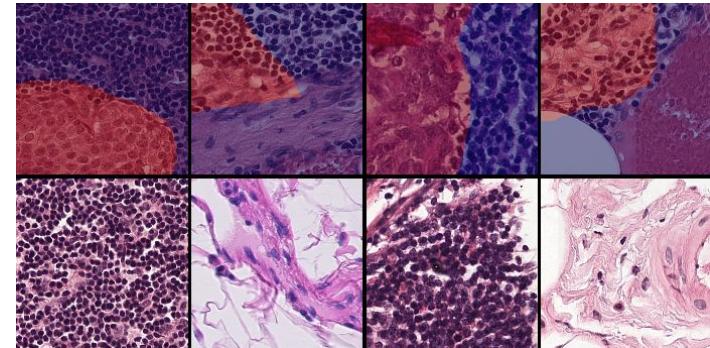
(a) Cancer Grading and Localization in Histology:

- **Experiments**

- Task: classify and localize ROI
- 2 public datasets: GlaS, Camelyon16 patches.



GlaS: colon cancer diagnosis



Camelyon16 patches: breast cancer

(a) Cancer Grading and Localization in Histology:

- Results

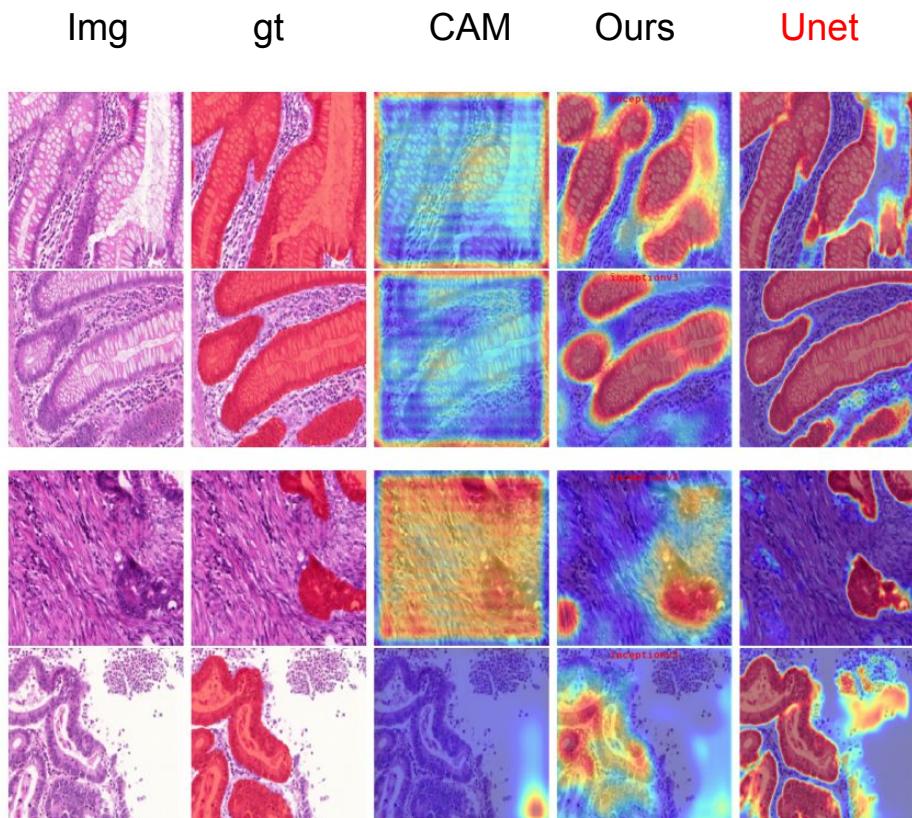
Metric	GlaS				CAMELYON16			
	VGG	Inception	ResNet	Mean	VGG	Inception	ResNet	Mean
WSL								
GAP (Lin et al., 2013) (corr,2013)								
GAP (Lin et al., 2013) (corr,2013)	58.5	57.5	56.2	57.4	37.5	24.6	43.7	35.2
MAX-POOL (Oquab et al., 2015) (cvpr,2015)	58.5	57.1	46.2	53.9	42.1	40.9	20.2	34.4
LSE (Sun et al., 2016) (cvpr,2016)	63.9	62.8	59.1	61.9	63.1	29.0	42.1	44.7
CAM (Zhou et al., 2016) (cvpr,2016)	68.5	50.5	64.4	61.1	25.4	48.7	27.5	33.8
HaS (Singh and Lee, 2017) (iccv,2017)	65.5	65.4	63.5	64.8	25.4	47.1	29.7	34.0
GradCAM (Selvaraju et al., 2017) (iccv,2017)	75.7	56.9	70.0	67.5	40.2	34.4	29.1	34.5
WILDCAT (Durand et al., 2017) (cvpr,2017)	56.1	54.9	60.1	57.0	44.4	31.4	31.0	35.6
ACoL (Zhang et al., 2018a) (cvpr,2018)	63.7	58.2	54.2	58.7	31.3	39.3	31.3	33.9
SPG (Zhang et al., 2018b) (eccv,2018)	63.6	58.3	51.4	57.7	45.4	24.5	22.6	30.8
GradCAM++ (Chattopadhyay et al., 2018) (wacv,2018)	76.1	65.7	70.7	70.8	41.3	43.9	25.8	37.0
Deep MIL (Ilse et al., 2018) (icml,2018)	66.6	61.8	64.7	64.3	53.8	51.1	57.9	54.2
PRM (Zhou et al., 2018) (cvpr,2018)	59.8	53.1	62.3	58.4	46.0	41.7	23.2	36.9
ADL (Choe and Shim, 2019) (cvpr,2019)	65.0	60.6	54.1	59.9	19.0	46.0	46.0	37.0
CutMix (Yun et al., 2019) (eccv,2019)	59.9	50.4	56.7	55.6	56.4	44.9	20.7	40.6
Smooth-GradCAM (Omeiza et al., 2019) (corr,2019)	71.3	67.6	75.5	71.4	35.1	31.6	25.1	30.6
XGradCAM (Fu et al., 2020) (bmvc,2020)	73.7	66.4	62.6	67.5	40.2	33.0	24.4	32.5
LayerCAM (Jiang et al., 2021) (ieee,2021)	67.8	66.1	70.9	68.2	34.1	25.0	29.1	29.4
NEGEV (ours)	81.3	70.1	82.0	77.8	70.3	53.8	52.6	58.9
Fully supervised								
U-Net (Ronneberger et al., 2015)(miccai,2015)	96.8	95.4	96.4	96.2	83.0	82.2	83.6	82.9

Table 1: PxAP performance over GlaS and CAMELYON16 test sets.

(a) Cancer Grading and Localization in Histology:

• Results

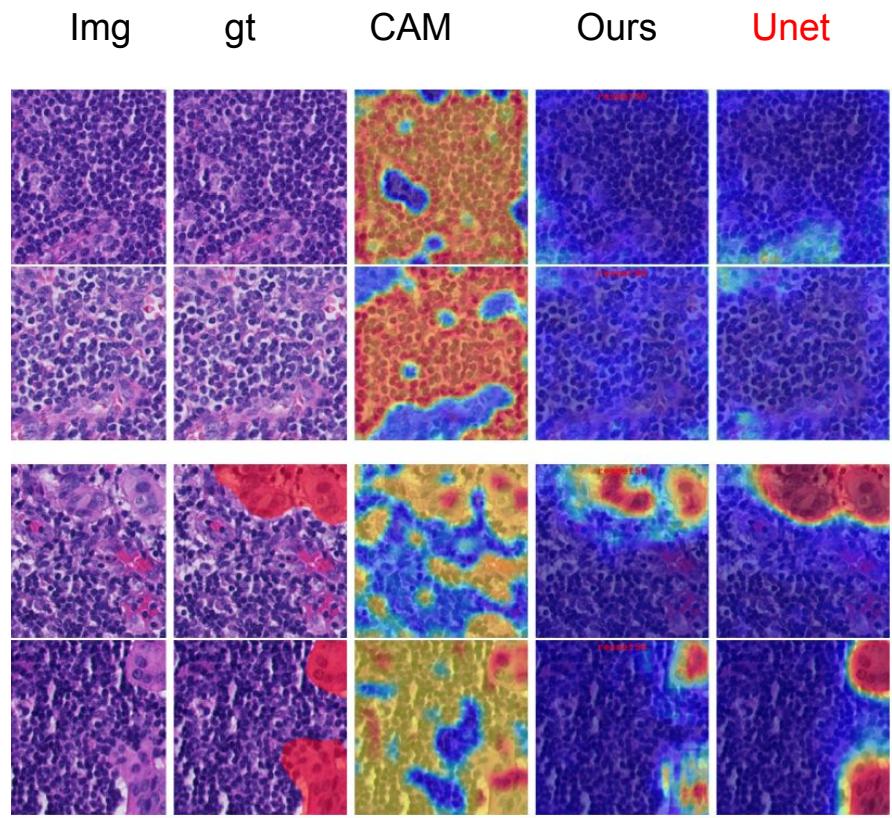
Full-sup



GlaS

Camelyon16

Full-sup



(a) Cancer Grading and Localization in Histology:

- **Ablations**

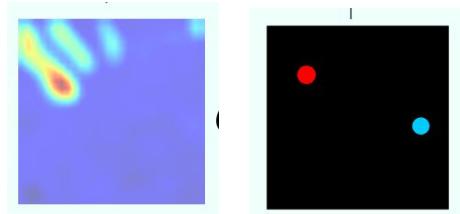
Methods	GlaS				CAMELYON16			
	VGG	Inception	ResNet	Mean	VGG	Inception	ResNet	Mean
CAM (Zhou et al., 2016)	68.5	50.5	64.4	61.1	25.4	48.7	27.5	20.3
Ours + C ⁺	81.3	53.3	81.3	71.9	38.1	36.5	30.8	35.1
Ours + C ⁺ + C ⁻	81.3	70.1	82.0	77.8	38.1	35.3	30.2	34.5
Ours + C ⁺ + C ⁻ + D ⁻	—	—	—	—	70.3	53.8	52.6	58.9
Improvement	+12.8	+19.6	+17.6	+16.6	+44.9	+5.0	+25.1	+25.0

Impact of different terms

$$\begin{aligned} \min_{\theta} \quad & \mathbb{1}_{\mathbf{X} \in \mathbb{D}^-} \left(\sum_{p \in \Omega} -\log(1 - S_p^0) \right) \\ & + (1 - \mathbb{1}_{\mathbf{X} \in \mathbb{D}^-}) \left(\lambda \sum_{p \in \{\mathbb{C}^+ \cup \mathbb{C}^-\}} \mathbf{H}(Y_p, S_p) \right), \end{aligned}$$

(a) Cancer Grading and Localization in Histology:

- **Ablations**



Methods	GlaS				CAMELYON16			
	VGG	Inception	ResNet	Mean	VGG	Inception	ResNet	Mean
CAM (Zhou et al., 2016)	68.5	50.5	64.4	61.1	25.4	48.7	27.5	33.8
Ours ($n = 1$, random selection)	81.3	70.1	82.0	77.8	70.3	53.8	52.6	58.9
Ours ($n = 1$, static selection)	77.7	60.3	76.5	71.5	57.5	47.4	42.8	49.2
Performance drop	-3.6	-9.8	-5.5	-6.3	-12.8	-6.4	-9.8	-9.6

Fixed vs random seeds selection

(a) Cancer Grading and Localization in Histology:

- **Ablations**

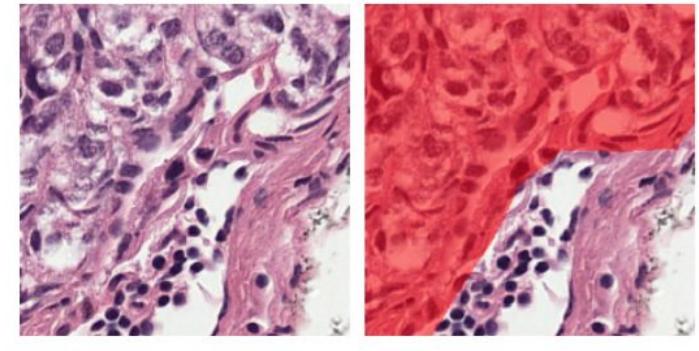
n	GlaS				CAMELYON16			
	VGG	Inception	ResNet	Mean	VGG	Inception	ResNet	Mean
1	81.3	70.1	82.0	77.8	70.3	53.8	52.6	58.9
2	81.3	52.9	81.3	71.8	69.7	51.1	47.2	56.0
3	81.3	52.9	81.3	71.8	69.7	51.9	47.2	56.2
4	81.3	55.0	81.3	72.5	69.7	50.0	47.2	56.6
5	81.3	52.9	81.3	71.8	69.7	53.4	47.2	56.7
10	81.3	53.7	81.3	72.1	69.7	52.6	47.2	56.5
20	81.3	52.9	82.2	72.1	69.7	51.3	47.2	56.0
50	81.3	52.0	81.3	71.5	69.7	53.8	50.3	57.9
100	81.3	53.4	81.3	72.0	69.7	50.5	47.2	55.8
500	81.3	52.9	81.3	71.8	69.7	51.5	47.6	56.2
1k	81.3	53.7	81.3	72.1	69.7	51.2	48.5	56.4
2k	81.3	53.0	81.3	71.8	69.7	51.5	47.2	56.1
3k	81.3	54.2	81.3	72.2	69.7	50.4	48.5	56.2
4k	81.3	52.9	81.3	71.8	69.7	52.9	47.2	56.6
5k	81.3	53.2	82.7	72.4	69.7	51.4	47.7	56.2
10k	81.3	52.9	81.3	71.8	69.7	52.1	47.2	56.3
<hr/>								
CAM (Zhou et al., 2016)	68.5	50.5	64.4	61.1	25.4	48.7	27.5	33.8

How many pixels to sample?

(a) Cancer Grading and Localization in Histology:

- Presented work

- Belharbi, S., Rony, J., Dolz, J., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty.** IEEE Transactions on Medical Imaging, 41:702–714.
- Belharbi, S., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Negative evidence matters in interpretable histology image classification.** In Medical Imaging with Deep Learning (MIDL).



Completed.

Part 4:

Applications of WSOL / WSSS

- (a) Medical Cancer Grading and ROI Localization in Histology
- (b) Weakly-Supervised Video Object Localization**
- (c) Person ReID: Embedding Networks
- (d) Medical Semantic Segmentation

(b) Weakly-Supervised Video Object Localization:
(Ongoing work)

**TCAM: Temporal Class Activation Maps for Object
Localization in Weakly-Labeled Unconstrained Videos**

Will be available early Sept. (paper + code)



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Task

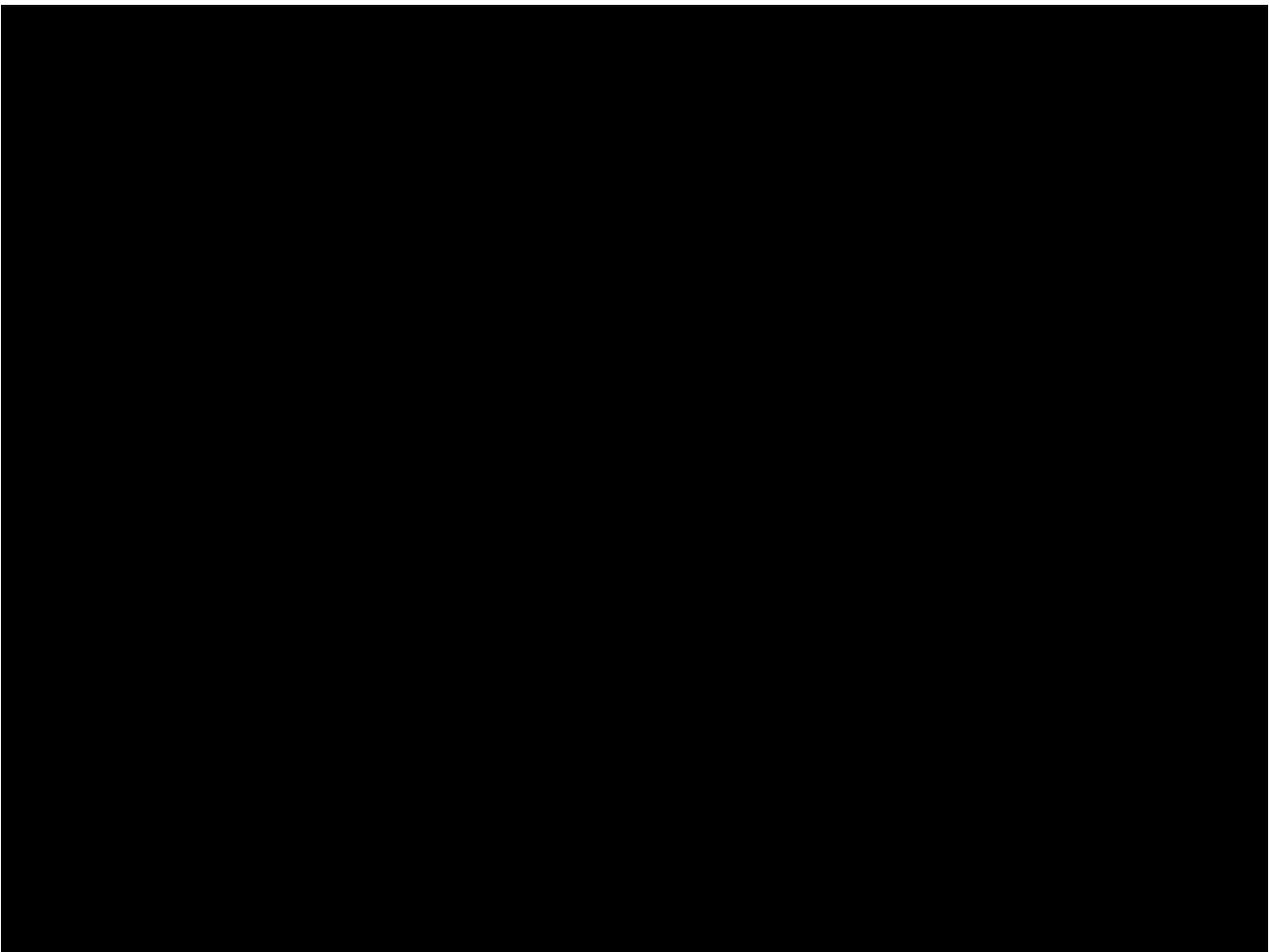
Moving objects

Camera motions

View-point changes

Decoding artifacts

Editing effects



Video:

https://docs.google.com/file/d/1sF9cynslvdcS_pAtDispmyxUPGa2Y7dg/preview

Unconstrained videos
(Dataset: YouTube-Objects v1.0)



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

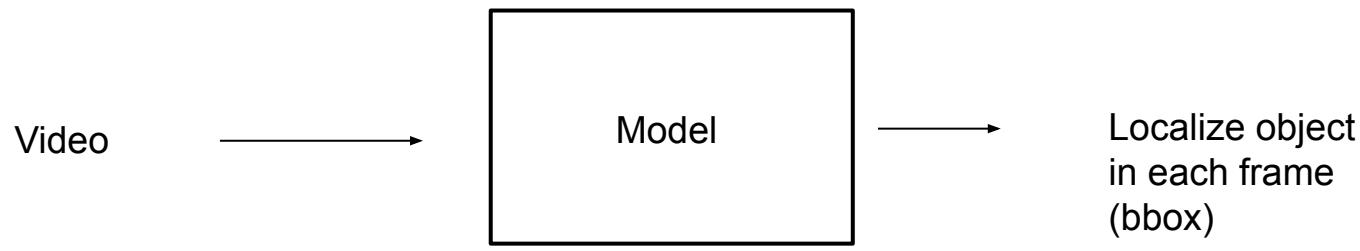
- **Task: Supervision**

- Labeling all frames via bbox is expensive
- Use **weak supervision: global video tag (cheap!)**

Global video tag: main object class in the video (**not necessarily present in all frames**)

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Task



Model is trained using weak labels (global video tag)

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Task**

Video object localization helps:

- Localizing object of interest in video
- Video content understanding
- Improve subsequent tasks: video summarization, event detection, video object detection, video tracking, ...

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Current state-of-the-art**

- Stagnated (<= 2020)

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Current state-of-the-art**

- Stagnated (<= 2020)
- Multiple sequential, independent stages:
 1. Generate spatio-temporal segments/proposals (visual and motion cues)
 2. Identify prominent object
 3. Refine

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Current state-of-the-art**

- Stagnated (<= 2020)
- Multiple sequential, independent stages:
 1. Generate spatio-temporal segments/proposals (visual and motion cues)
 2. Identify prominent object
 3. Refine
- Video tags are used only to cluster video

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Current state-of-the-art**

- Stagnated (<= 2020)
- Multiple sequential, independent stages:
 1. Generate spatio-temporal segments/proposals (visual and motion cues)
 2. Identify prominent object
 3. Refine
- Video tags are used only to cluster video
- ROI are not necessarily discriminative

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Current state-of-the-art**

- Stagnated (<= 2020)
- Multiple sequential, independent stages:
 1. Generate spatio-temporal segments/proposals (visual and motion cues)
 2. Identify prominent object
 3. Refine
- Video tags are used only to cluster video
- ROI are not necessarily discriminative
- Motion cues (optical flow) is not necessarily discriminative, is noisy and requires further post-processing

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Current state-of-the-art**

- Stagnated (<= 2020)
- Multiple sequential, independent stages:
 1. Generate spatio-temporal segments/proposals (visual and motion cues)
 2. Identify prominent object
 3. Refine
- Video tags are used only to cluster video
- ROI are not necessarily discriminative
- Motion cues (optical flow) is not necessarily discriminative, is noisy and requires further post-processing
- Localization is done by solving an optimization problem over a cluster of videos or single video (slow inference time, build model per-class/video)

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Proposal**

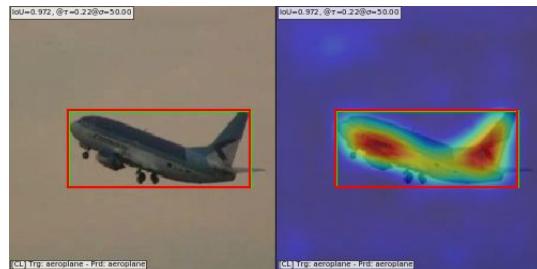
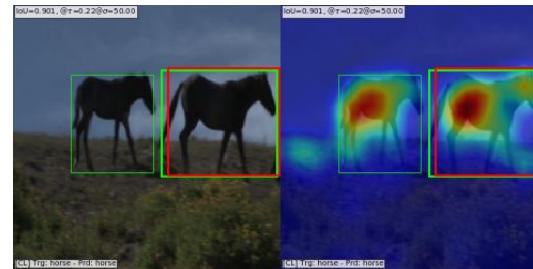
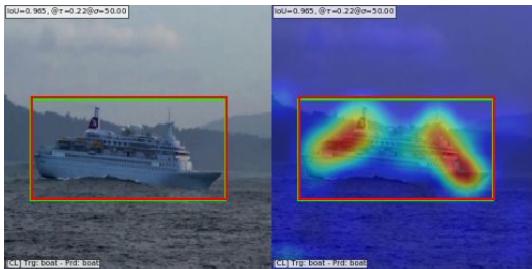
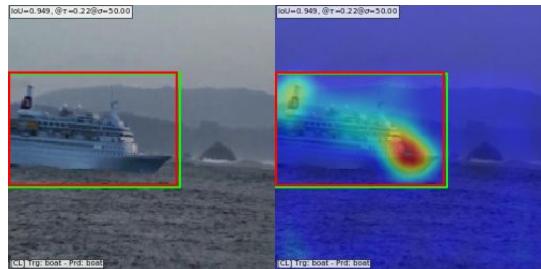
Leverage CAMs for weakly supervised video object localization



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- ## Proposal

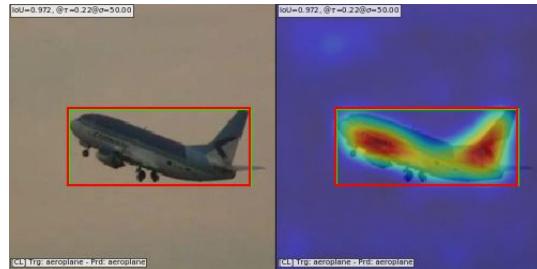
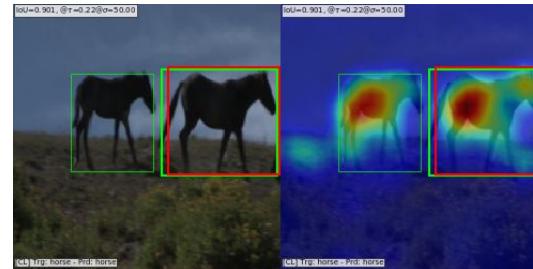
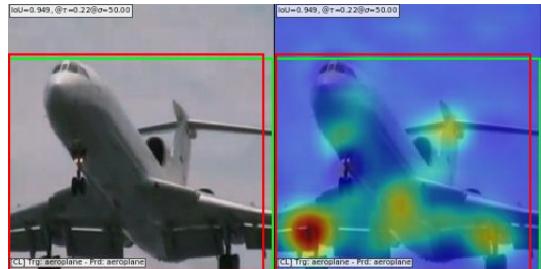
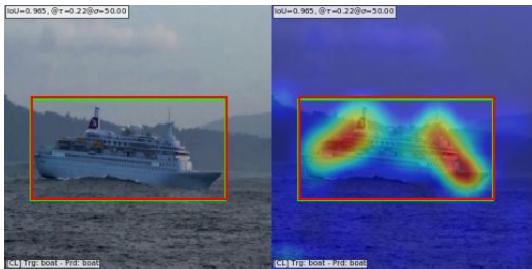
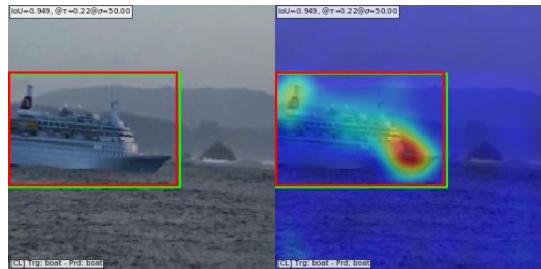
WSOL CAM methods trained on **still images** yield descent performance



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- ## Proposal

WSOL CAM methods trained on **still images** yield descent performance



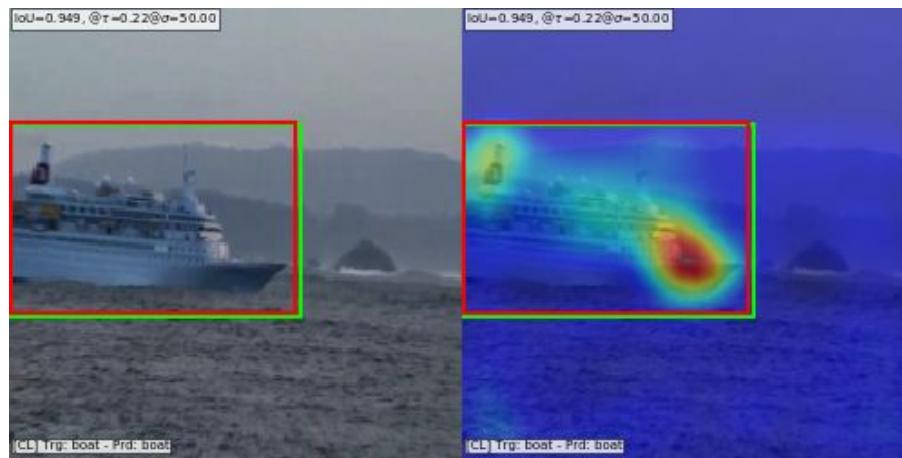
Do not account for
spatio-temporal
dependency!

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Proposal**

Observation:

Slight variation in consecutive frames leads to variation in CAMs



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Proposal**

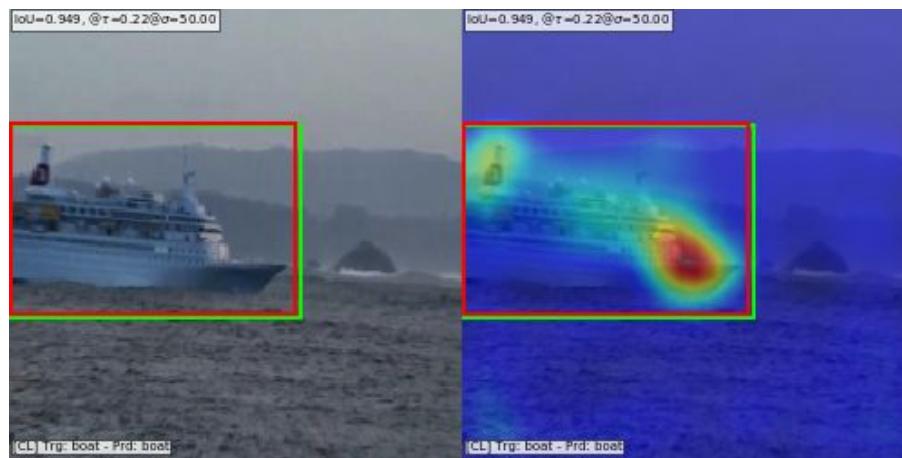
Observation:

Slight variation in consecutive frames leads to variation in CAMs

—>

Aggregate consecutive CAMs to build complete CAM

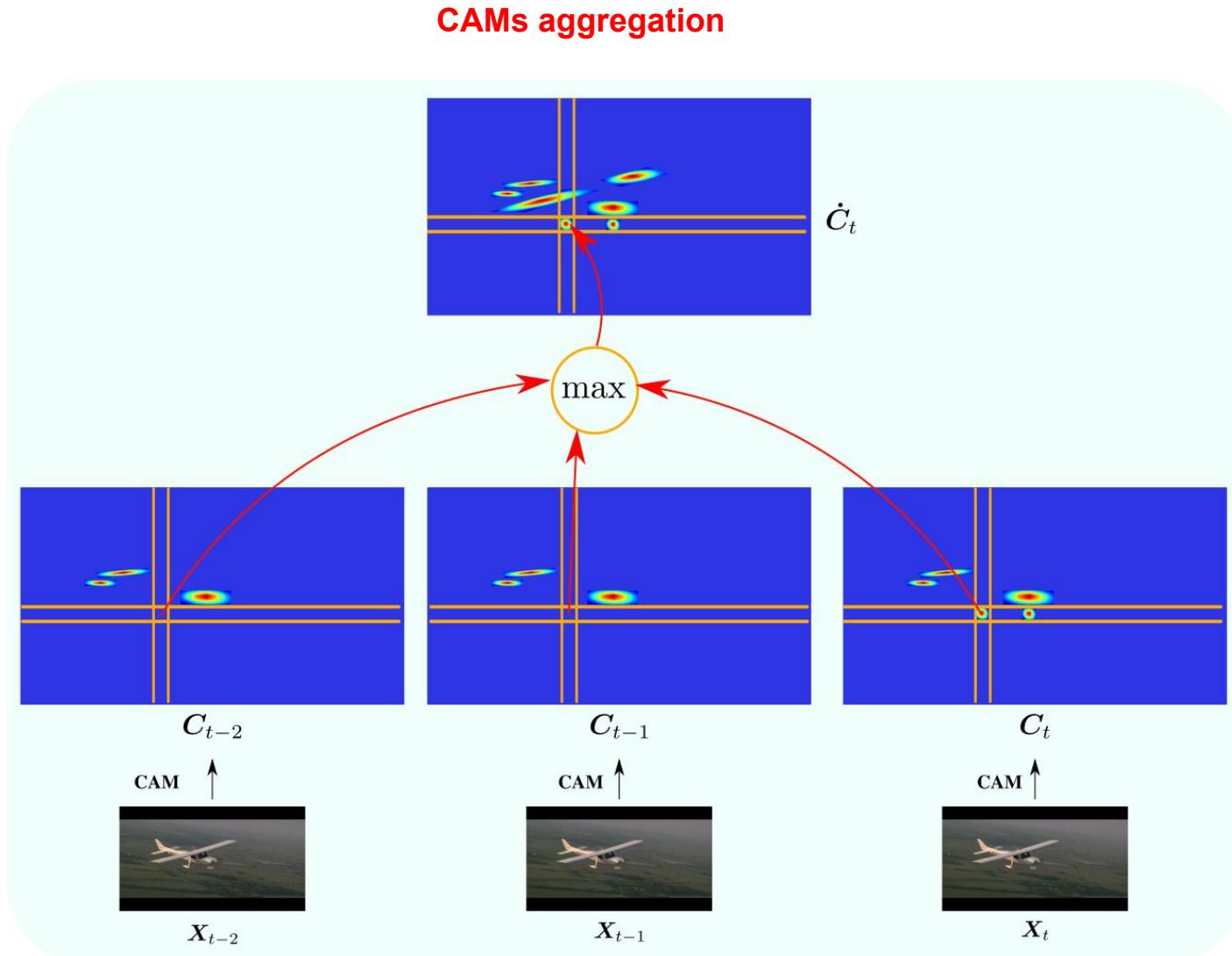
And use it to sample pseudo-labels



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

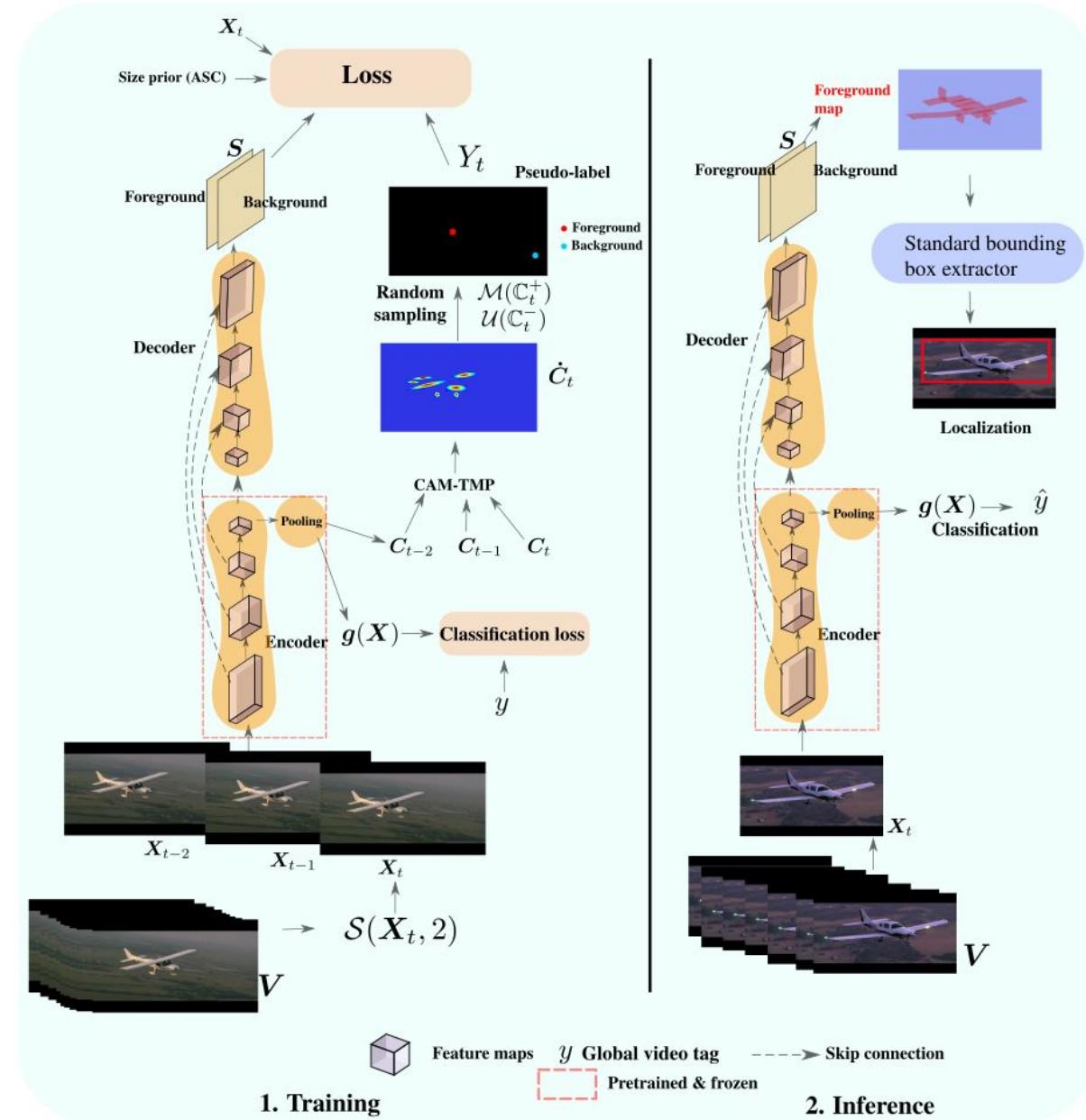
- **Proposal**

**CAM-Temporal Max Pooling
(CAM-TMP)**



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Proposal**

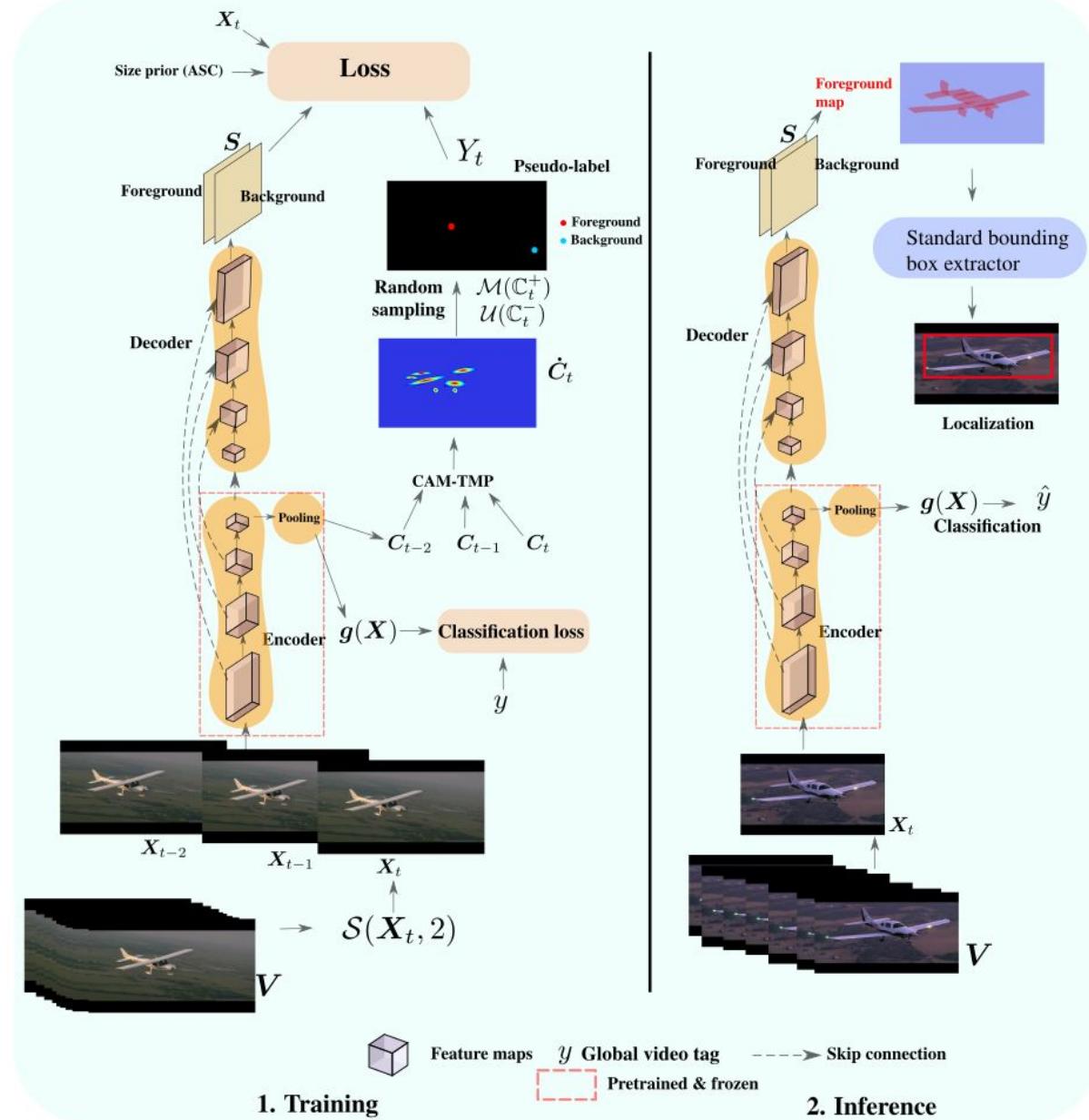


(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Proposal**

Training:
accounts for spatio-temporal dependency

$$\begin{aligned} \min_{\theta} \quad & \sum_{p \in \Omega'_t} H_p(Y_t, S_t) + \lambda \mathcal{R}(S_t, X_t), \\ \text{s.t.} \quad & \sum S^r(t) \geq 0, \quad r \in \{0, 1\}, \end{aligned}$$



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

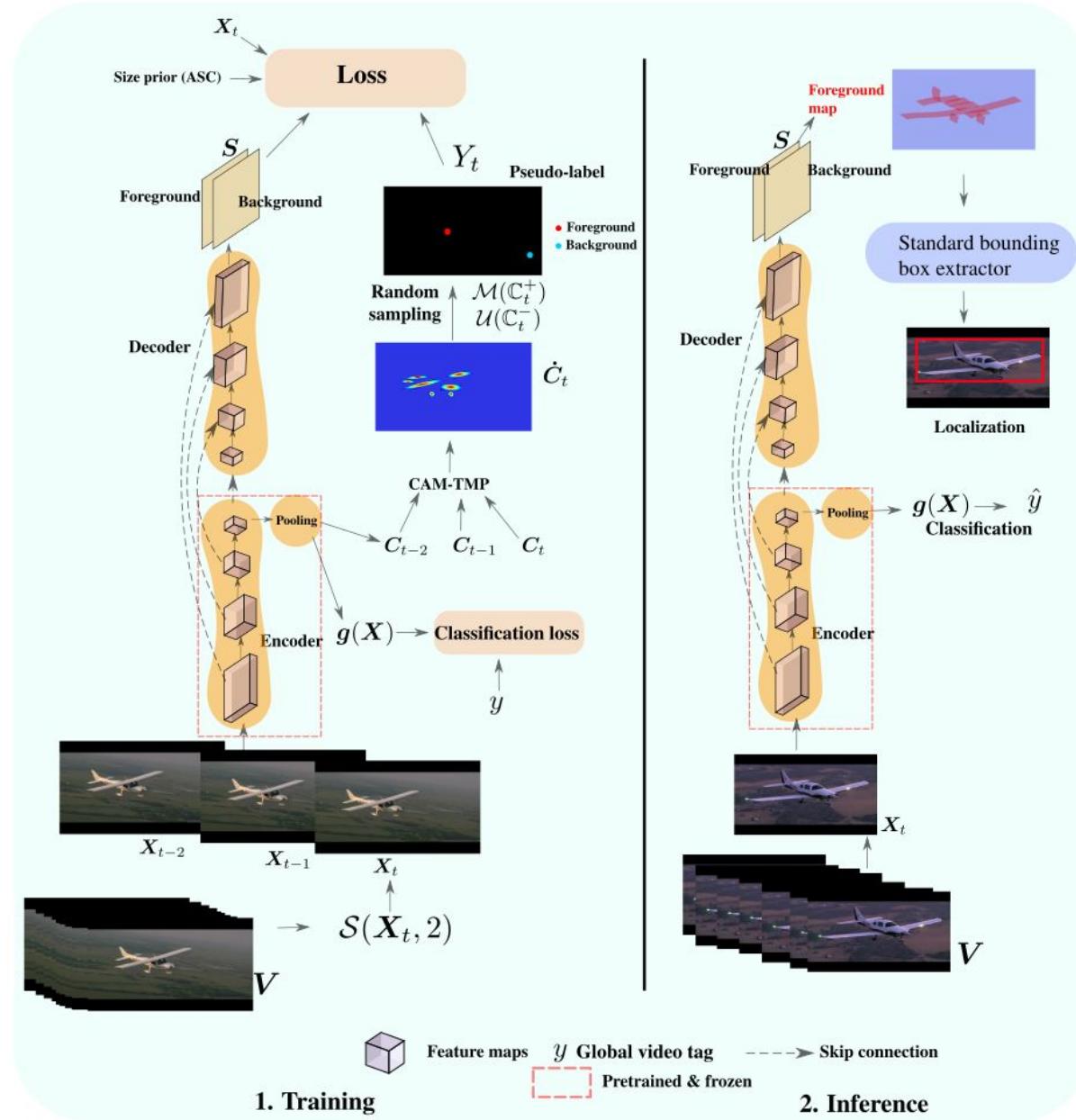
- **Proposal**

Training:
accounts for spatio-temporal dependency

$$\begin{aligned} \min_{\theta} \quad & \sum_{p \in \Omega'_t} H_p(Y_t, S_t) + \lambda \mathcal{R}(S_t, X_t), \\ \text{s.t.} \quad & \sum S^r(t) \geq 0, \quad r \in \{0, 1\}, \end{aligned}$$

Methods	CorLoc
Layer-CAM [22] (ieee,2021)	63.0
Ours + \mathbb{C}^+ + \mathbb{C}^-	68.5
Ours + \mathbb{C}^+ + \mathbb{C}^- + CRF	69.6
Ours + \mathbb{C}^+ + \mathbb{C}^- + ASC	66.2
Ours + \mathbb{C}^+ + \mathbb{C}^- + CRF + ASC	70.5
Ours + \mathbb{C}^+ + \mathbb{C}^- + CRF + ASC + CAM-TMP	72.8
Improvement	+9.8

Ablation

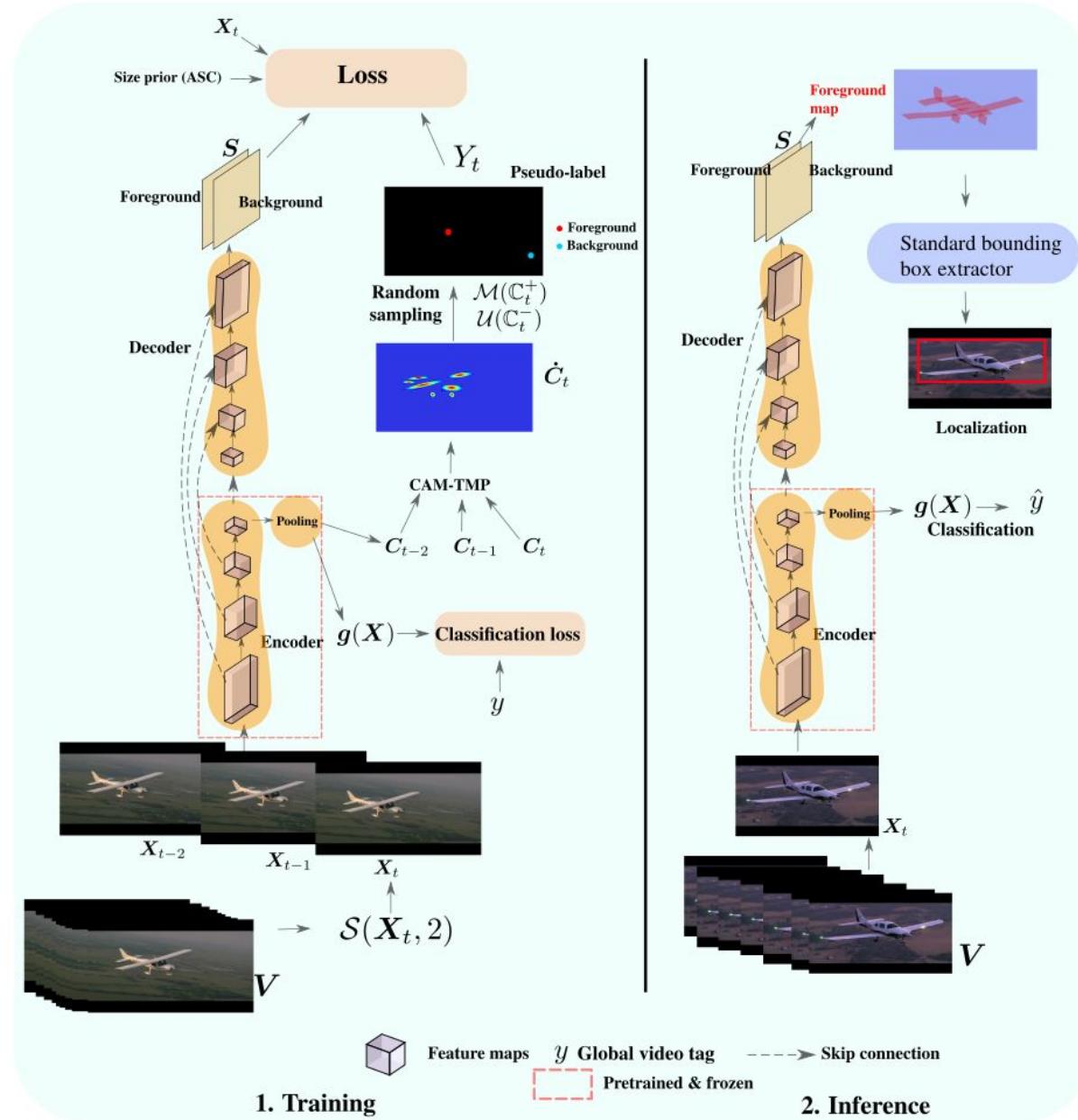


(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Proposal**

Inference:

Independent frames → fast



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- **Results**



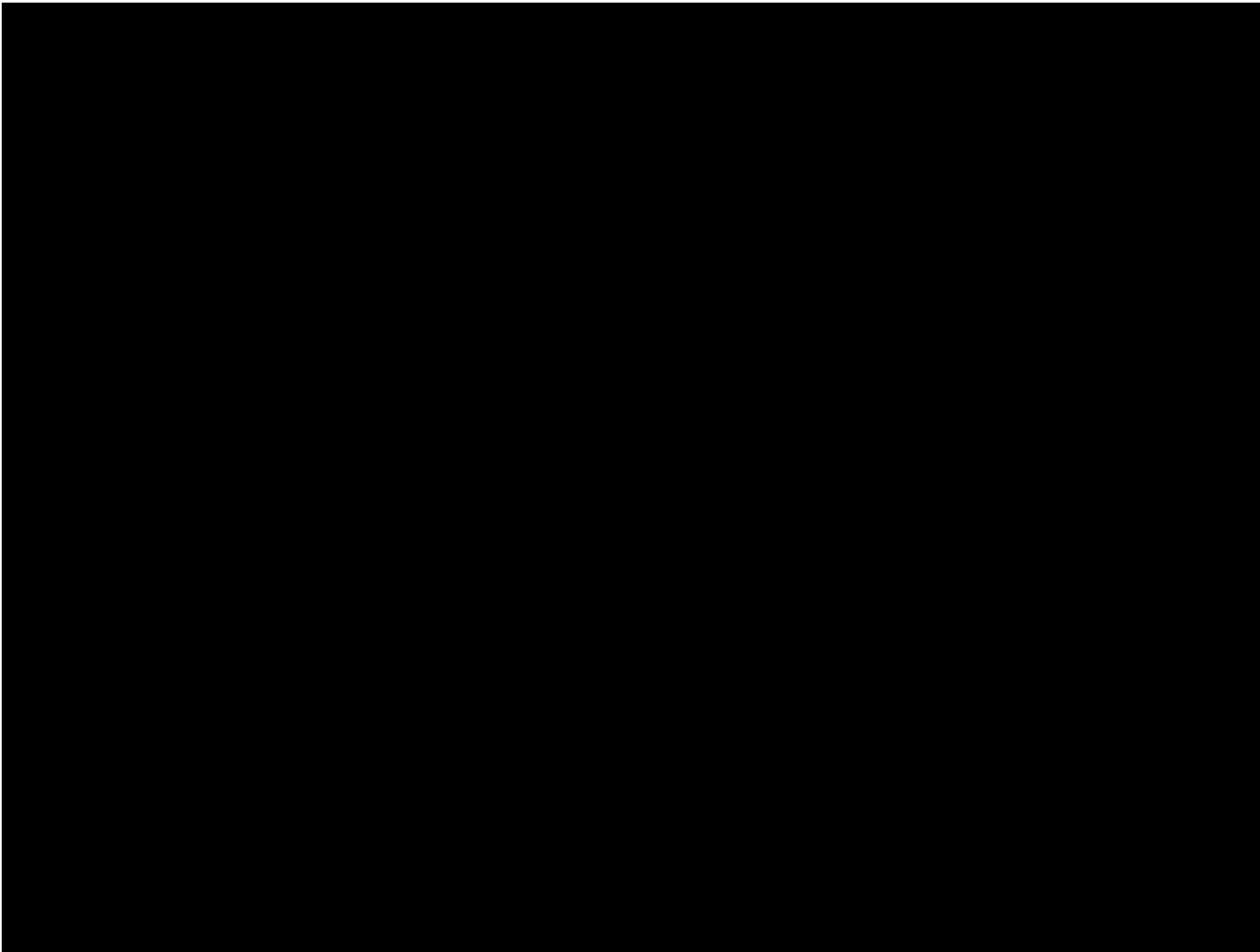
Video:
[https://docs.google.com/file/d/1KbrQu35oX2NpoH8
NiKDt0cY5Qy1kqbwo/preview](https://docs.google.com/file/d/1KbrQu35oX2NpoH8NiKDt0cY5Qy1kqbwo/preview)

“Plane”



(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Results



Video:

https://docs.google.com/file/d/1_UP0Dwdp7Tl4qs84BPvK2_rN7ENPysmK/preview

“Horse”

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Results



Video:

https://docs.google.com/file/d/1wCxm1votCm_1M-cENBq414ZA0tnZXCC5/preview

“Car”

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Results



Video:

<https://docs.google.com/file/d/1dckaUlkaaqPSyeQyIcbdNFyTTguVG-vQ/preview>

“Car”

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Results



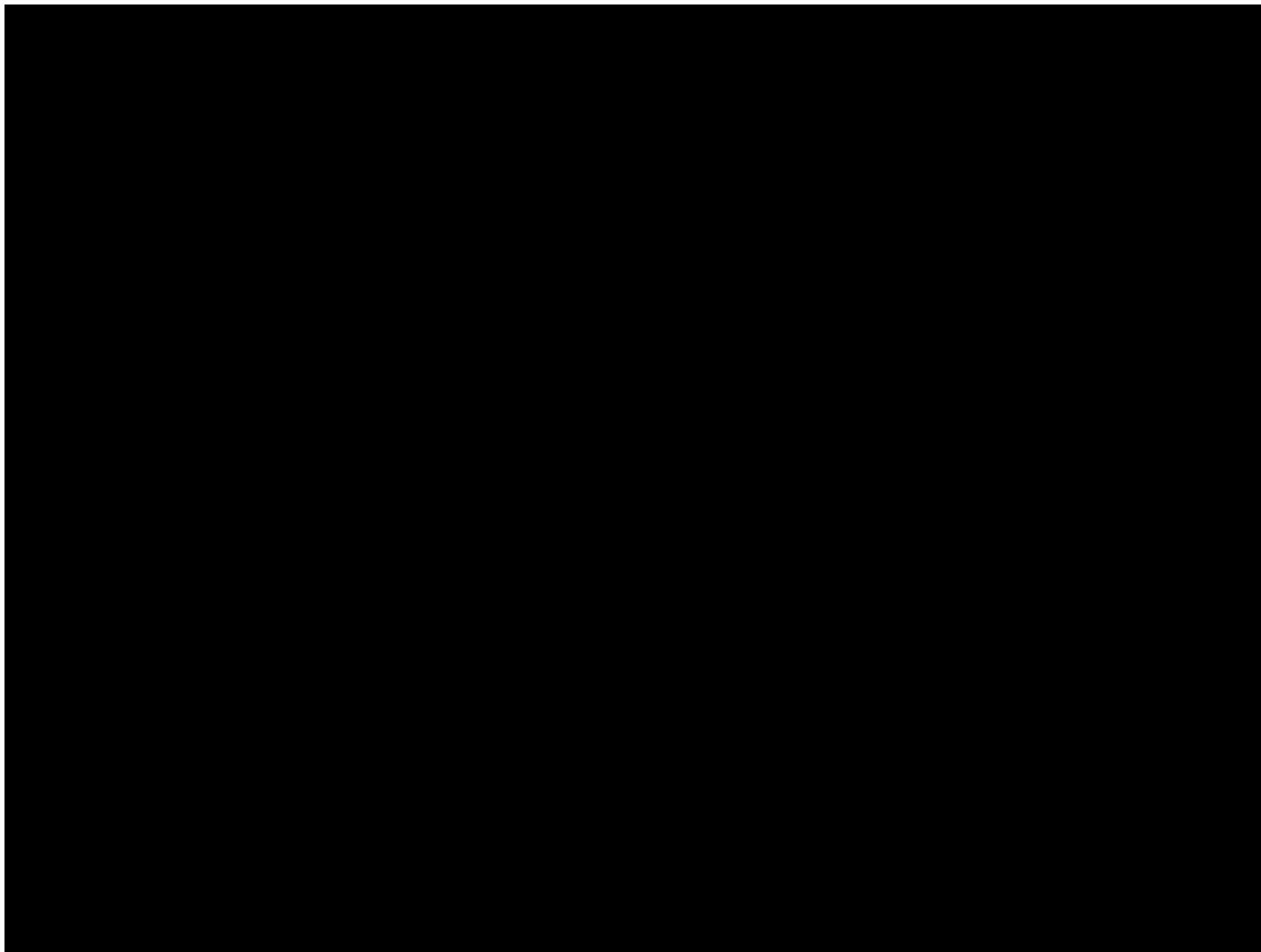
Video:

<https://docs.google.com/file/d/1XkZBEjtB-I-bu5nK2gSSsKH-FfASqOQ-/preview>

“Car”

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

- Results



Video:

https://docs.google.com/file/d/1Gf8y2SUNEKhqWk_kxARlo0xY65kJPorXq/preview

“Car”

(b) Weakly-Supervised Video Object Localization: (Ongoing work)

TCAM: Temporal Class Activation Maps for Object Localization in Weakly-Labeled Unconstrained Videos

Will be available early Sept. (paper + code)

Completed.



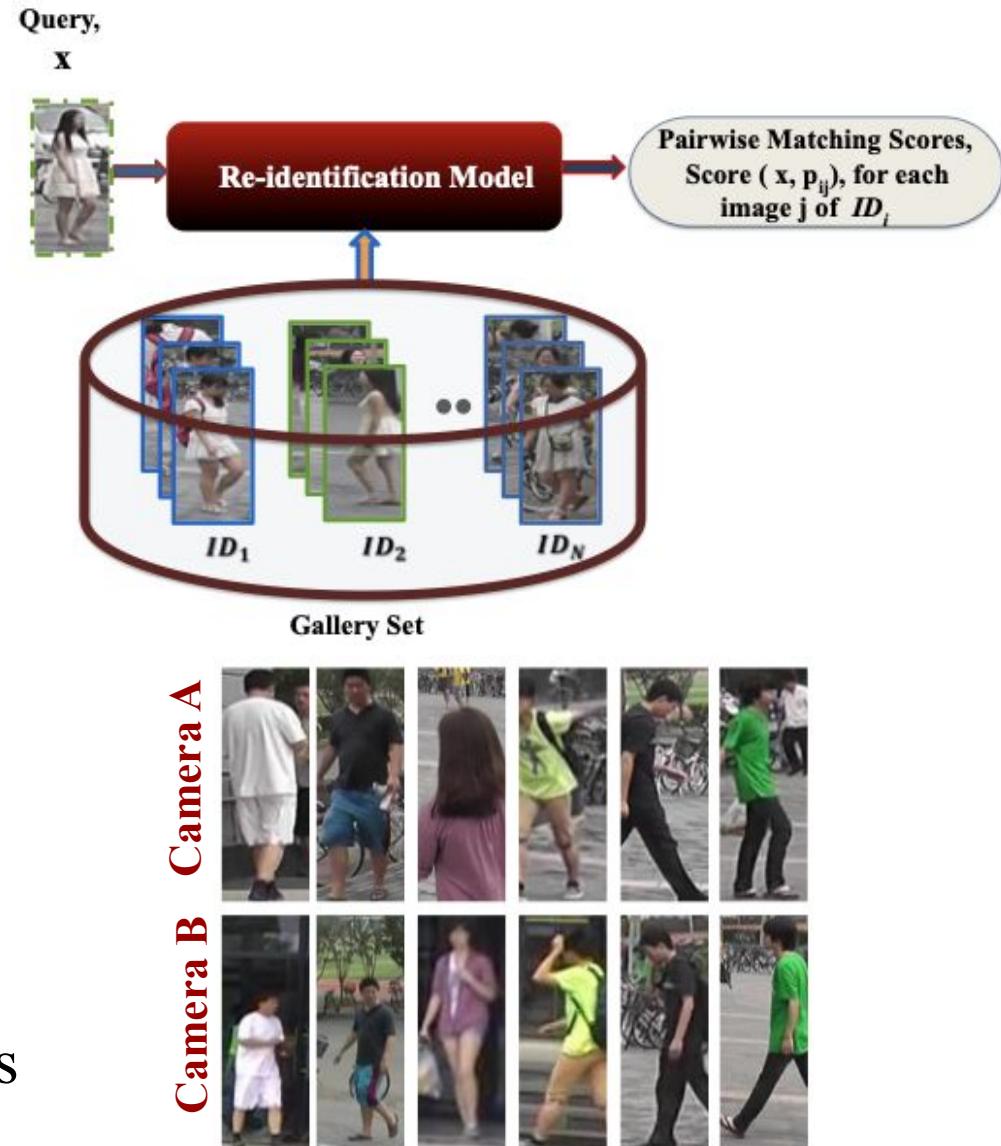
(c) Person ReID: Embedding Networks

Person Re-Identification:

recognize individuals captured over a distributed set of non-overlapping video camera views

Challenges in real-world scenarios:

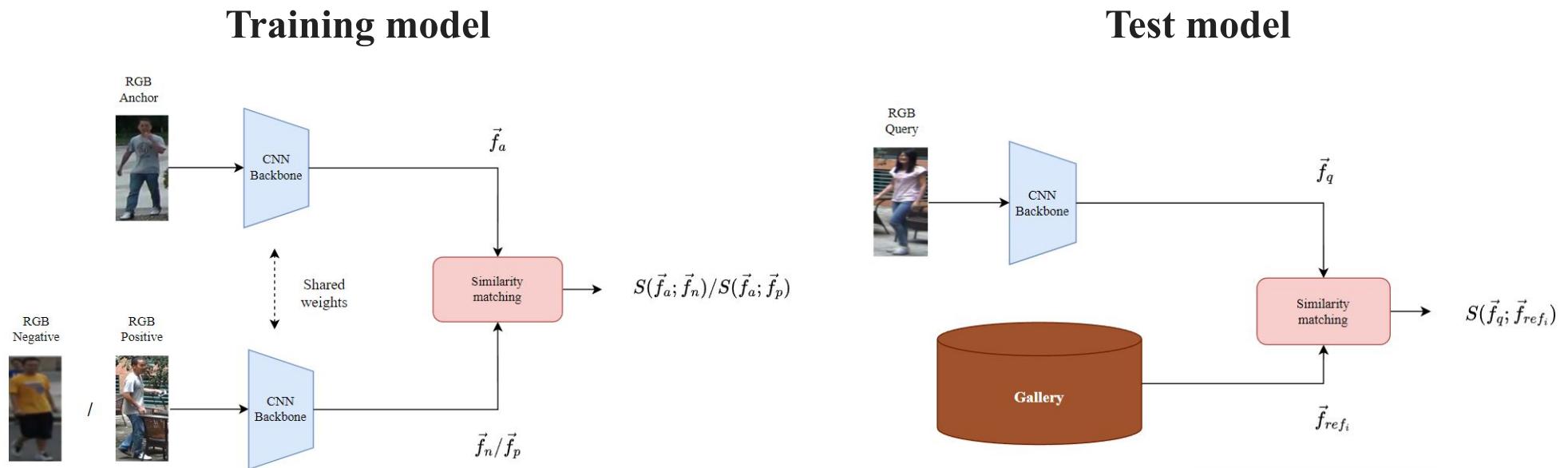
- low resolution, scales,
- motion blur
- occlusions, poor viewpoint
- variations in pose and viewpoint
- variation in illumination
- open set scenarios
- pedestrians with similar clothes
- misalignment over different views



(c) Person ReID: Embedding Networks

Deep Siamese CNNs for pairwise similarity matching:

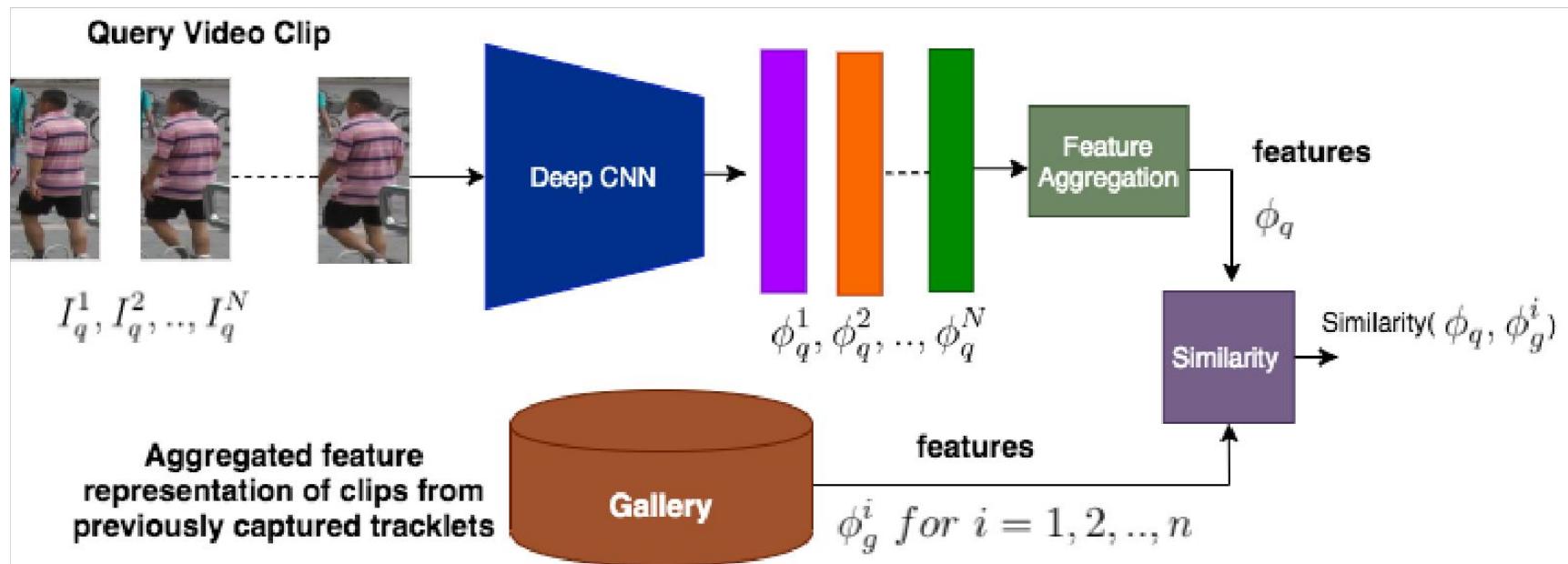
- **matching**: given query and reference images, they assess similarity, using Euclidean, cosine distance, etc., between feature representations
- **metric learning**: specific losses, like triplet, contrastive, and magnet losses are used to learn coherent embeddings



(c) Person ReID: Embedding Networks

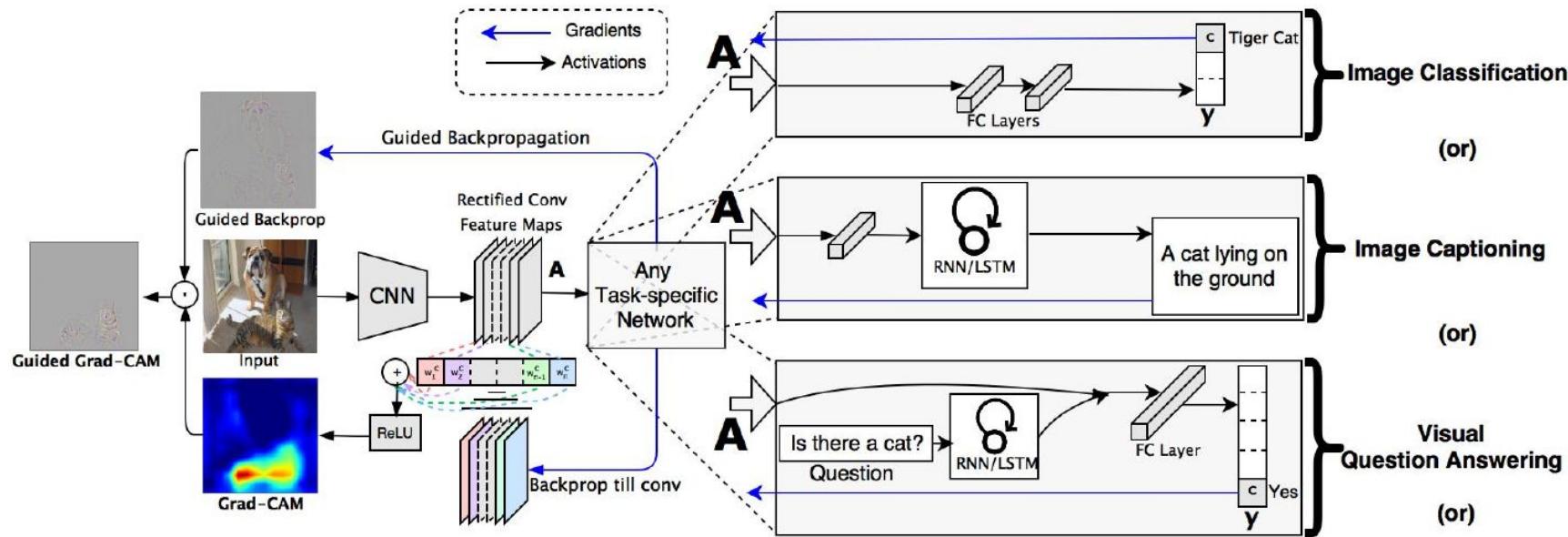
Deep Siamese CNNs for pairwise similarity matching:

- well adapted for person recognition in video surveillance
 - metric learning without training on reference data of the persons being sought
- can also be used for video person ReID, based on video clips



(c) Person ReID: Embedding Networks

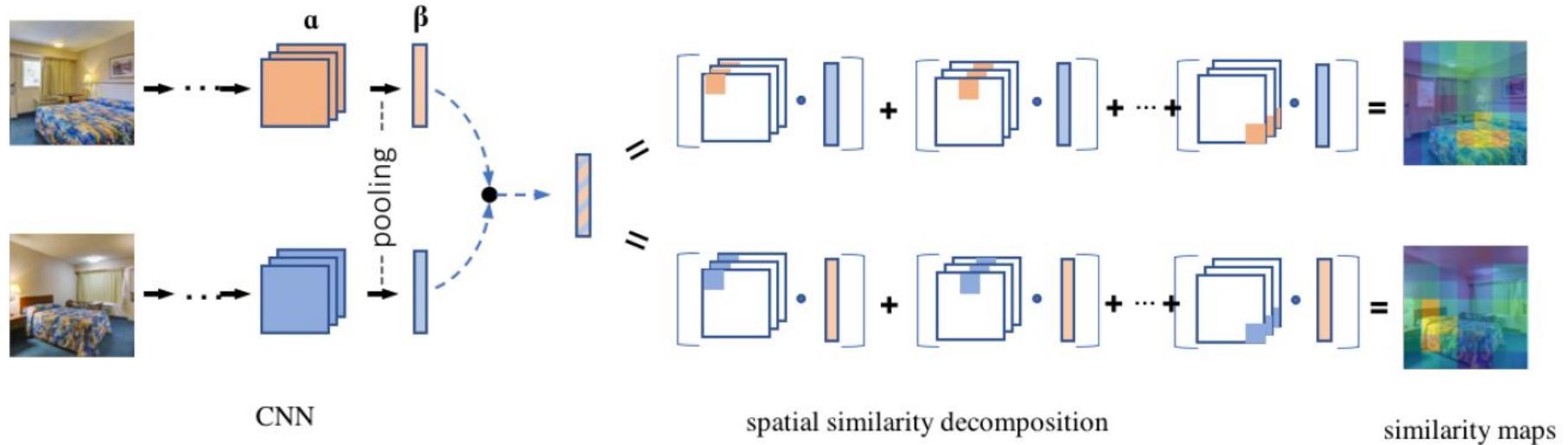
Grad-CAM: First extension of CAMs to diverse CNN models



- **Principle:** the gradient for a given class is used after the last conv layer to produce weights for the localization map
- **Not adapted** to embedding networks since it does not produce class-related gradients at inference time - classes not encountered by the model

(c) Person ReID: Embedding Networks

Visualizing Similarity Networks:

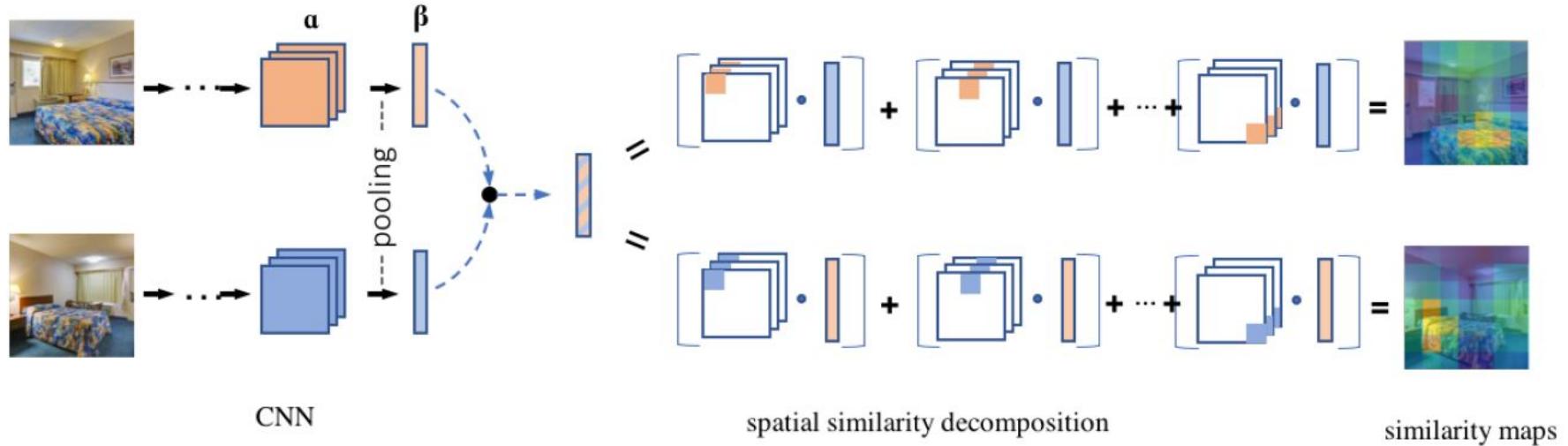


- Produces spatial similarity maps for a pair of images $I^{(i)}$ and $I^{(j)}$
 - In these maps, the cosine similarity between two image feature vectors $\beta^{(i)}$ and $\beta^{(j)}$ is spatially decomposed to highlight the contribution of image regions to the overall pairwise similarity:

$$s(\boldsymbol{\beta}^{(i)}, \boldsymbol{\beta}^{(j)}) = \frac{\boldsymbol{\beta}^{(i)} \cdot \boldsymbol{\beta}^{(j)}}{\|\boldsymbol{\beta}^{(i)}\| \|\boldsymbol{\beta}^{(j)}\|}$$

(c) Person ReID: Embedding Networks

Visualizing Similarity Networks:



- **Approach** – the similarity measure $s(\beta^{(i)}, \beta^{(j)})$ is spatially decomposed to observe the influence of each part of the image in this measure
- Computing the similarity maps depends on the pool operation.
- **Example:** using global average pooling $\beta = \frac{1}{K^2} \sum_{x,y} \alpha_{(x,y)}$

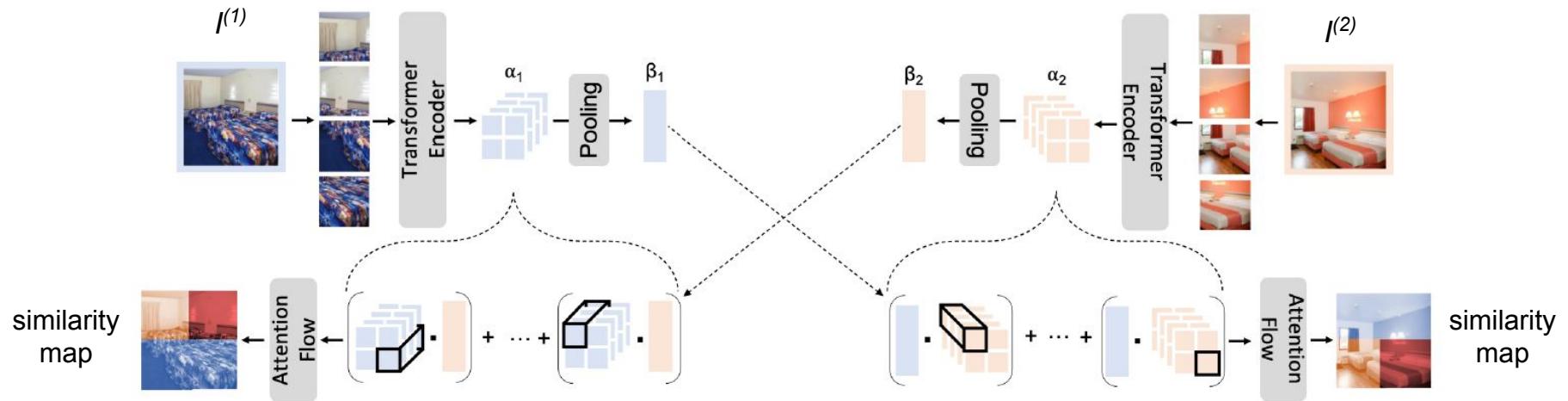
$$s(\beta^{(i)}, \beta^{(j)}) = \frac{\beta^{(i)} \cdot \beta^{(j)}}{\|\beta^{(i)}\| \|\beta^{(j)}\|} = \frac{\alpha_{(1,1)}^{(i)} \cdot \beta^{(j)} + \dots + \alpha_{(K,K)}^{(i)} \cdot \beta^{(j)}}{Z}$$

Z: normalizing factor $K^2 \|\beta^{(i)}\| \|\beta^{(j)}\|$

arrange terms spatially
and visualize as a
similarity map

(c) Person ReID: Embedding Networks

Visualizing Similarity Networks: extend the same idea for interpretation of transformer embedding networks

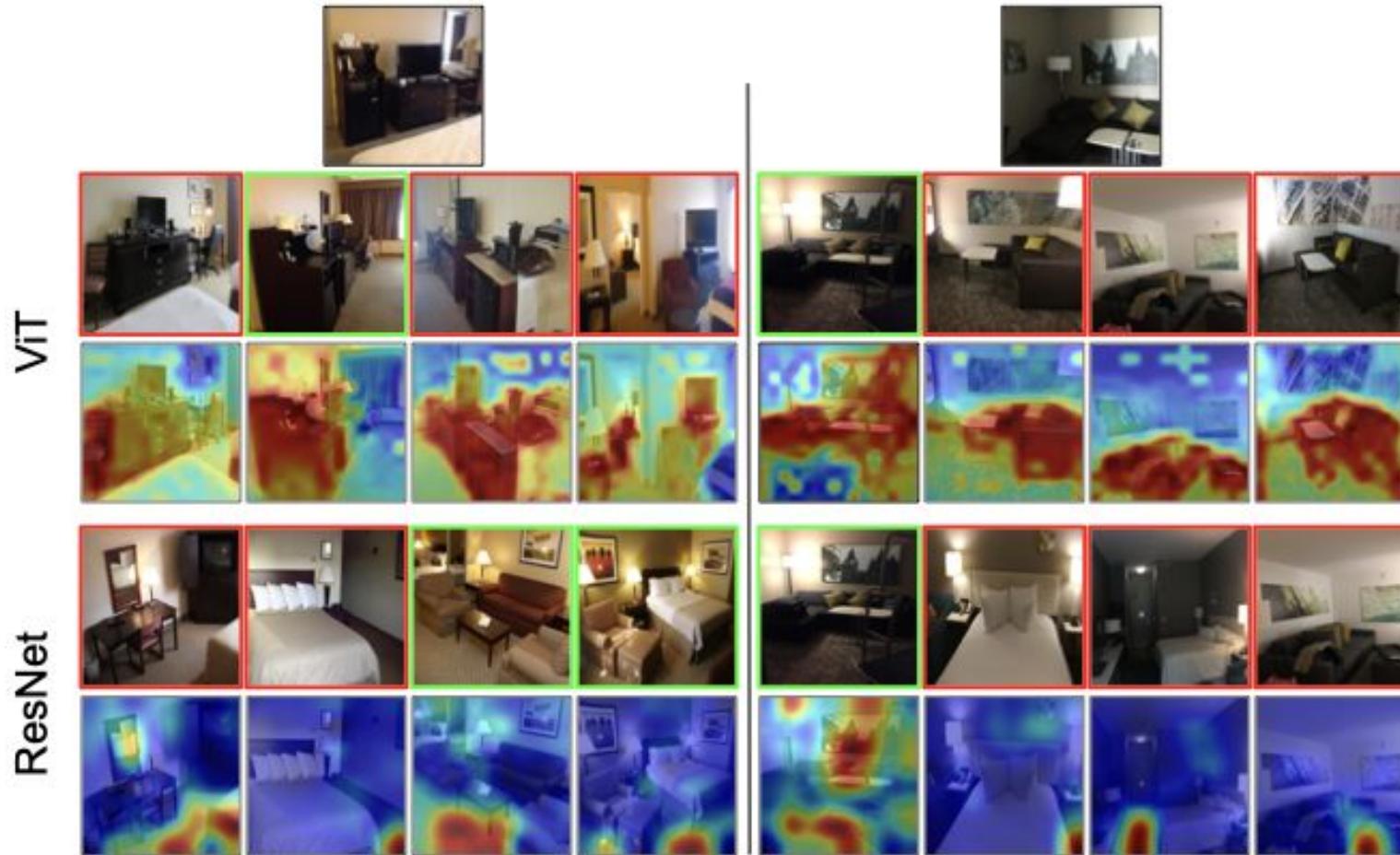


- Instead of using the feature map of the last conv. layer to decompose the similarity $s(\beta^{(i)}, \beta^{(j)})$, it relies the last output tokens
- **Idea** – Use rollout algorithm to approximately match token activation in the similarity measure with the correct spatial image location

$$s(\beta^{(i)}, \beta^{(j)}) = \frac{\beta^{(i)} \cdot \beta^{(j)}}{\|\beta^{(i)}\| \|\beta^{(j)}\|}$$

(c) Person ReID: Embedding Networks

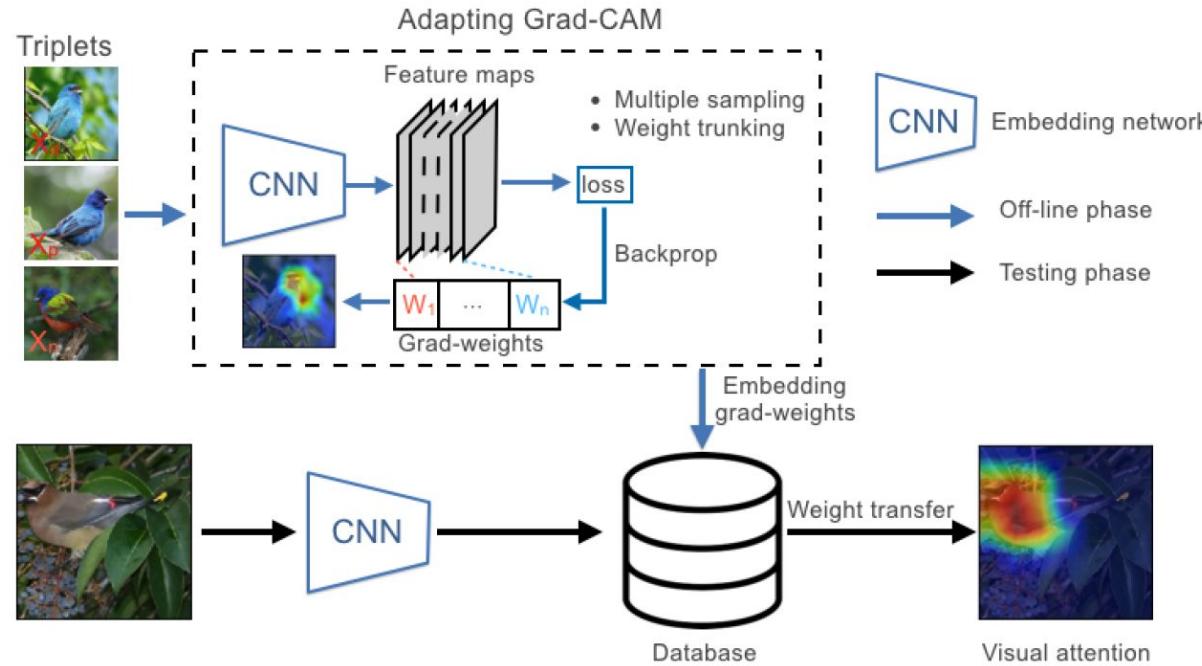
Visualizing Similarity Networks: extend the same idea for interpretation of transformer embedding networks



(c) Person ReID: Embedding Networks

Adapted Grad-CAM for Embedding Nets:

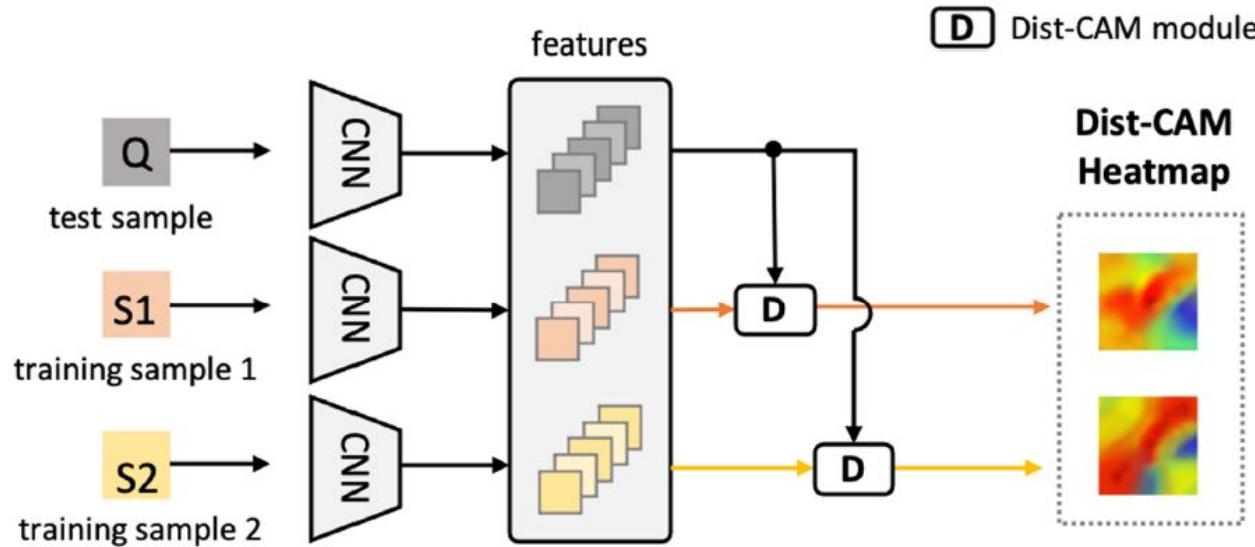
Replace the per-class gradient by the triplet loss gradient $g(A^k) = \frac{\partial \mathcal{L}_{tri}}{\partial A^k}$.



- query the nearest neighbor embedding in the training set for a given test set embedding
- apply the grad-weights of the nearest neighbor to produce the CAM visualization

(c) Person ReID: Embedding Networks

Distance-based CAM: Produce a map from the closest training sample weights

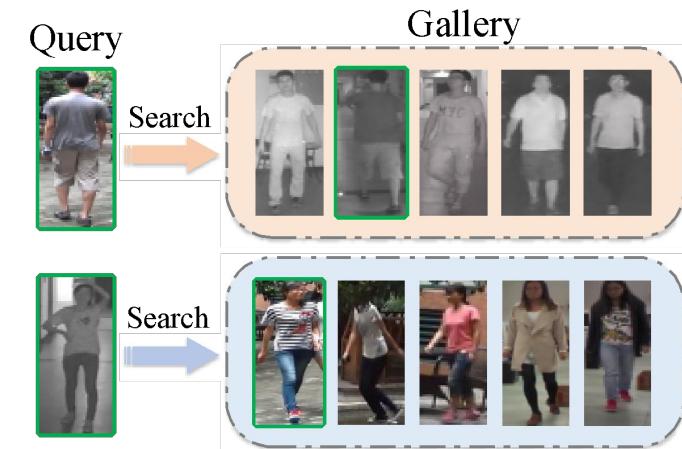


- query the closest embedding in the training set for a given test set embedding
- adapt the FC layer weights of the closest training sample regarding each channel distance from the test to the training feature map.

(c) Person ReID: Embedding Networks

Visible-Infrared (cross-modality) Person ReID:

- RGB can provide high quality color information (dependant on illumination)
- IR works well at night, or under poor lighting conditions



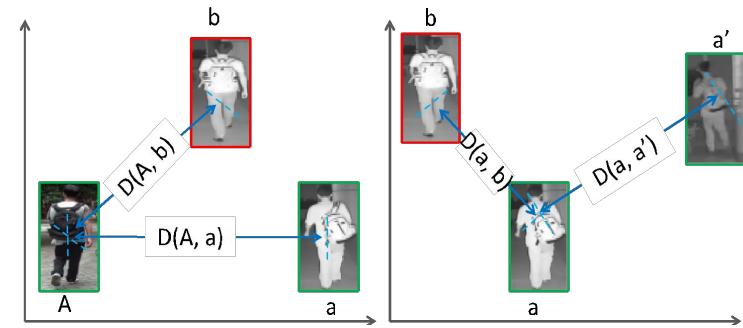
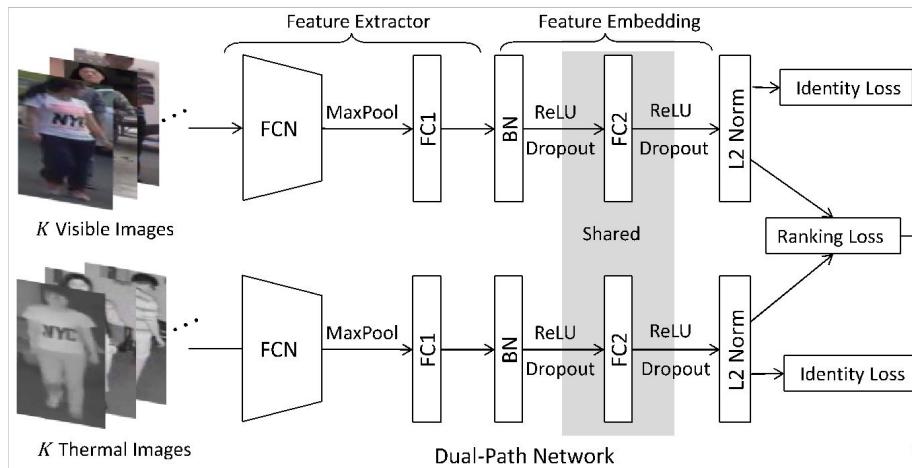
Objectives: match a same person across a network of different spectrum cameras (RGB and IR) - *large domain shift*



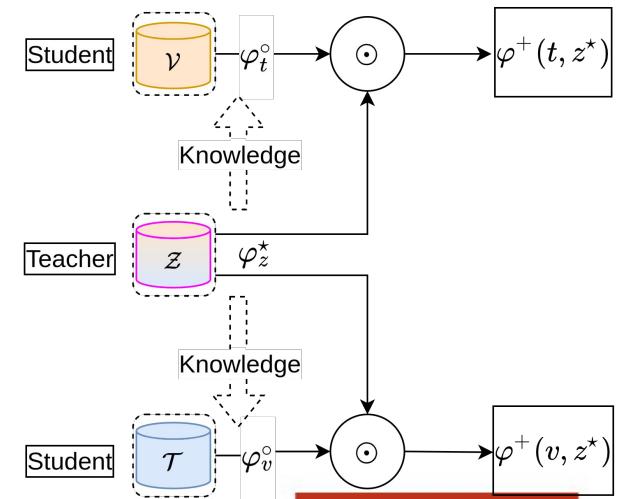
(c) Person ReID: Embedding Networks

Visible-Infrared (cross-modality) Person ReID:

- **BDTR framework:** dual-path network for feature extraction and bi-directional dual-constrained top-ranking loss



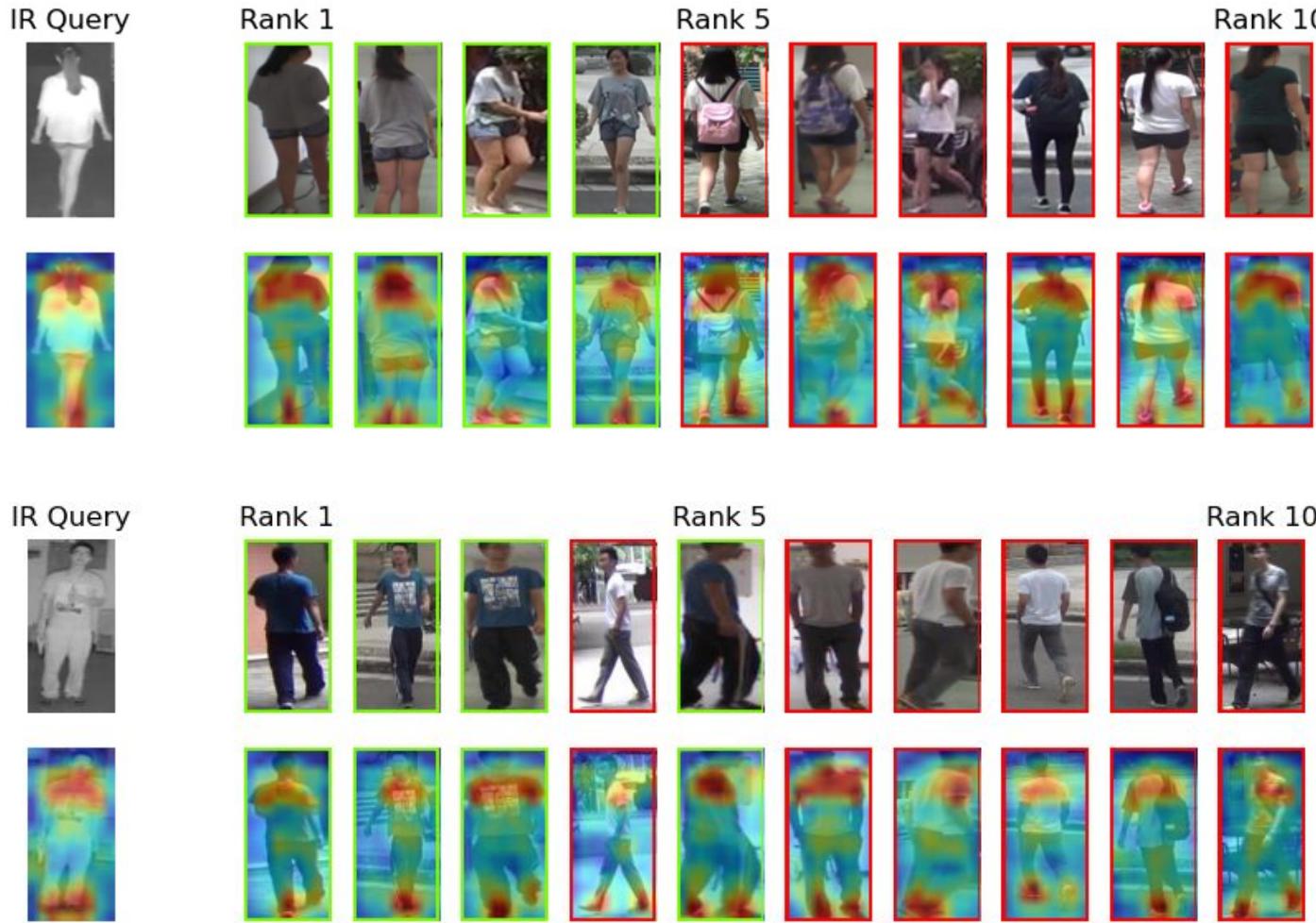
- Network trained using privileged information from intermediate virtual domains (teacher, Z), between infrared and visual (students, V and T)



(c) Person ReID: Embedding Networks

Visible-infrared Person ReID:

Visual results on SYSU dataset using “similarity CAM” (Stylianou, 2019)

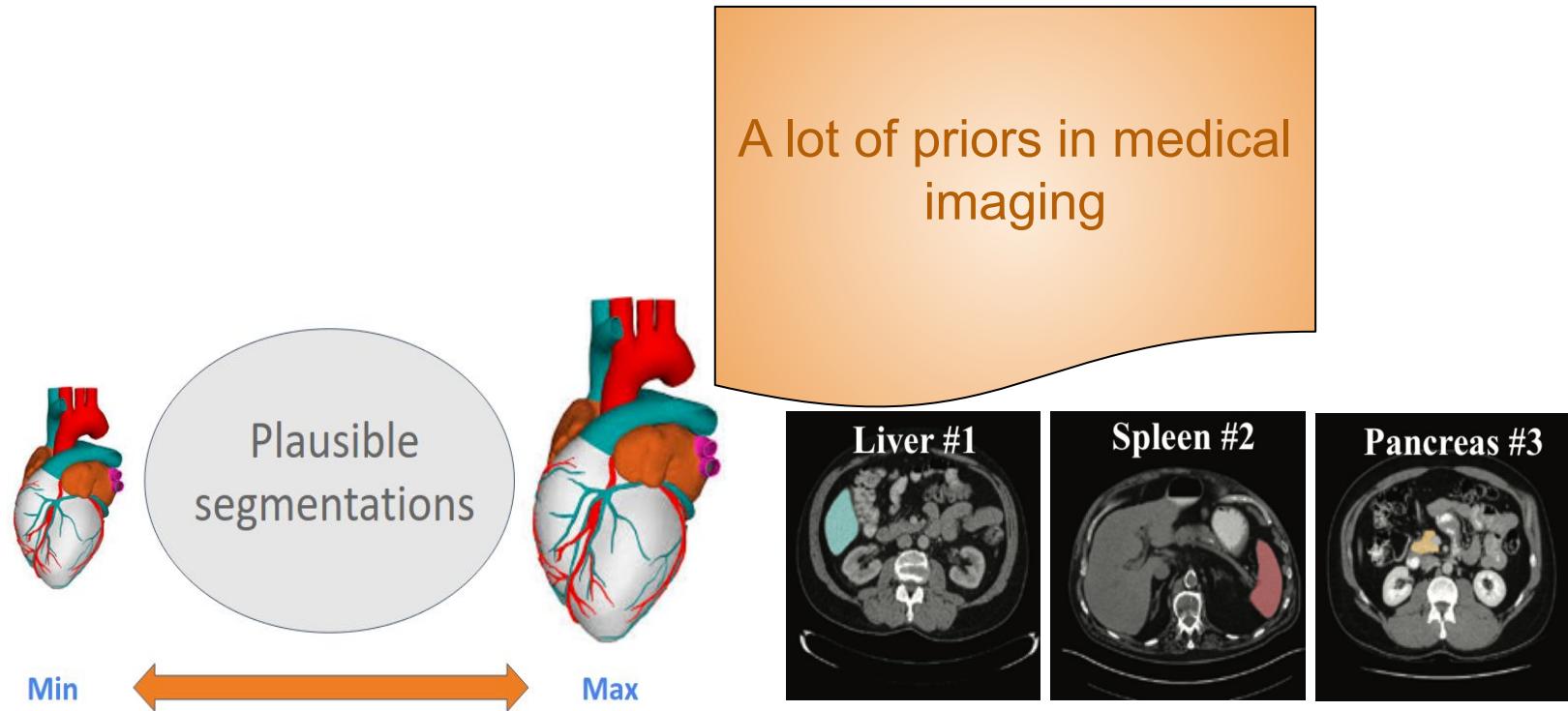


Medical Semantic Segmentation

Source: .

Anatomical Constraints

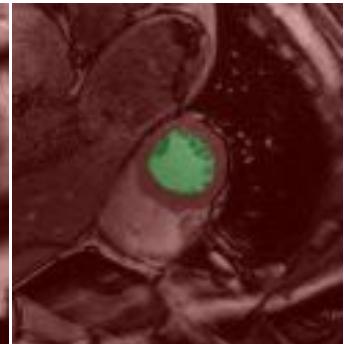
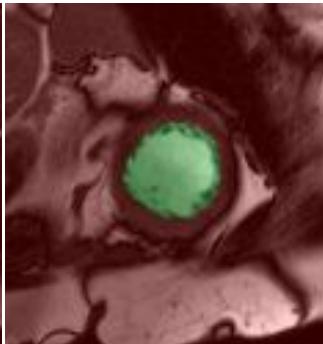
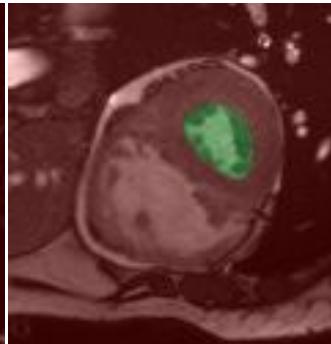
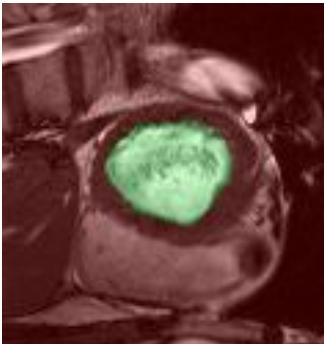
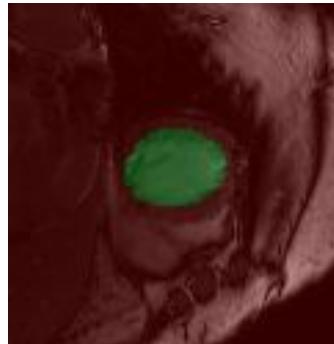
Data meets domain knowledge



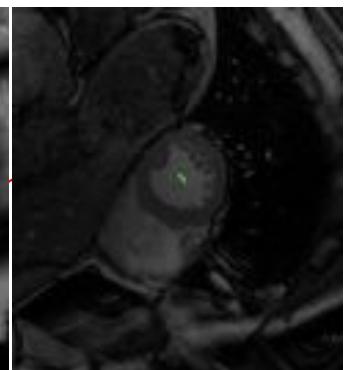
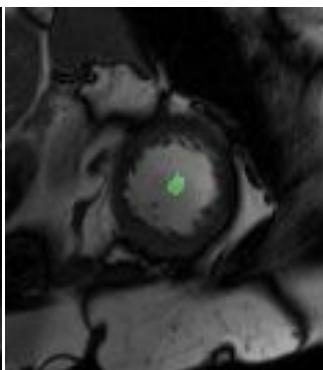
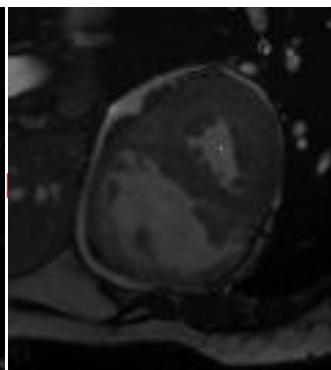
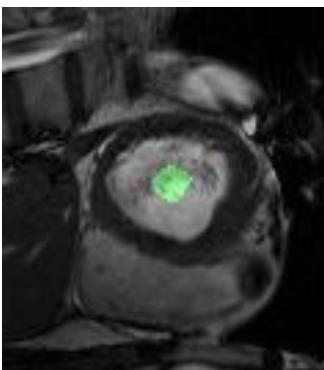
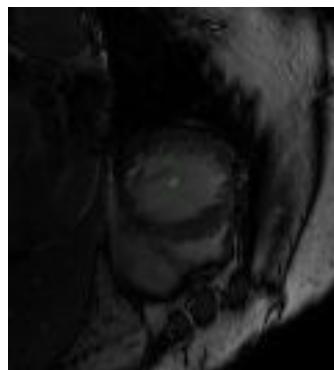
Anatomical priors (e.g., shapes)

Partially labeled data
(e.g., exploiting organ relationships)

Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints



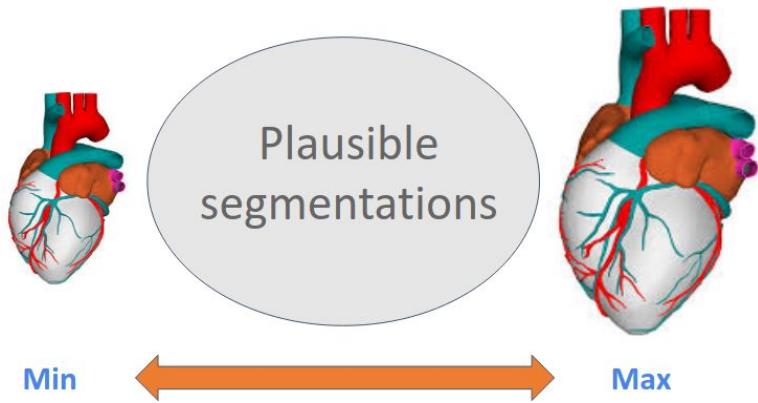
Full annotations



Partial annotations for cross-entropy

Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints

$$\min \leq \sum_{p \in \Omega} s_{\theta}^{p,c} \leq \max$$

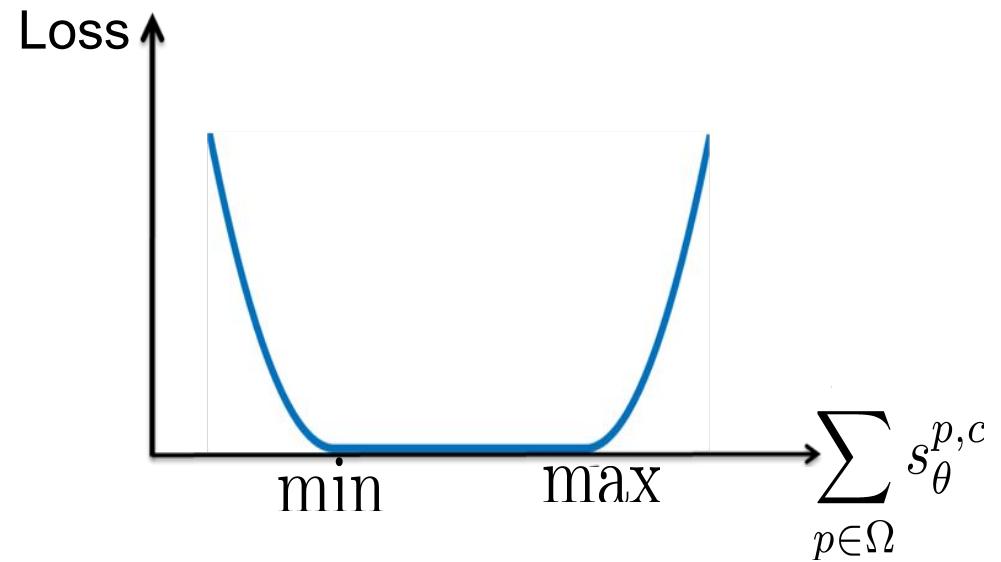
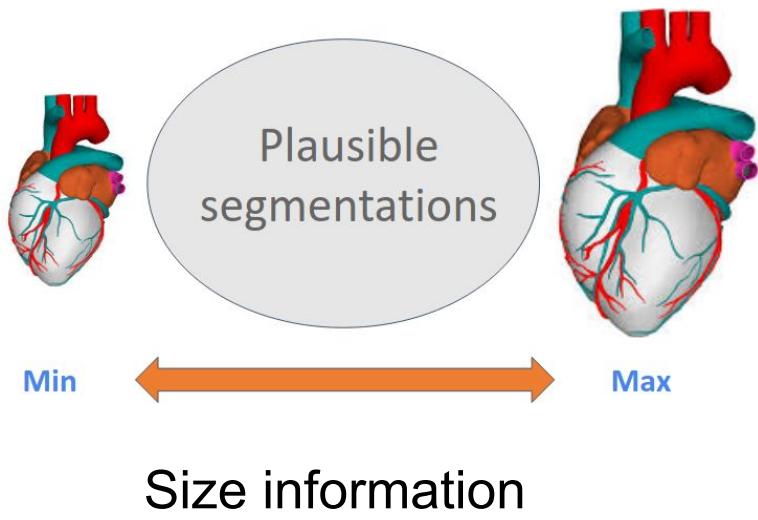


Size information

Kervadec et al., Constrained-CNN Losses for Weakly Supervised Segmentation,
Media'19

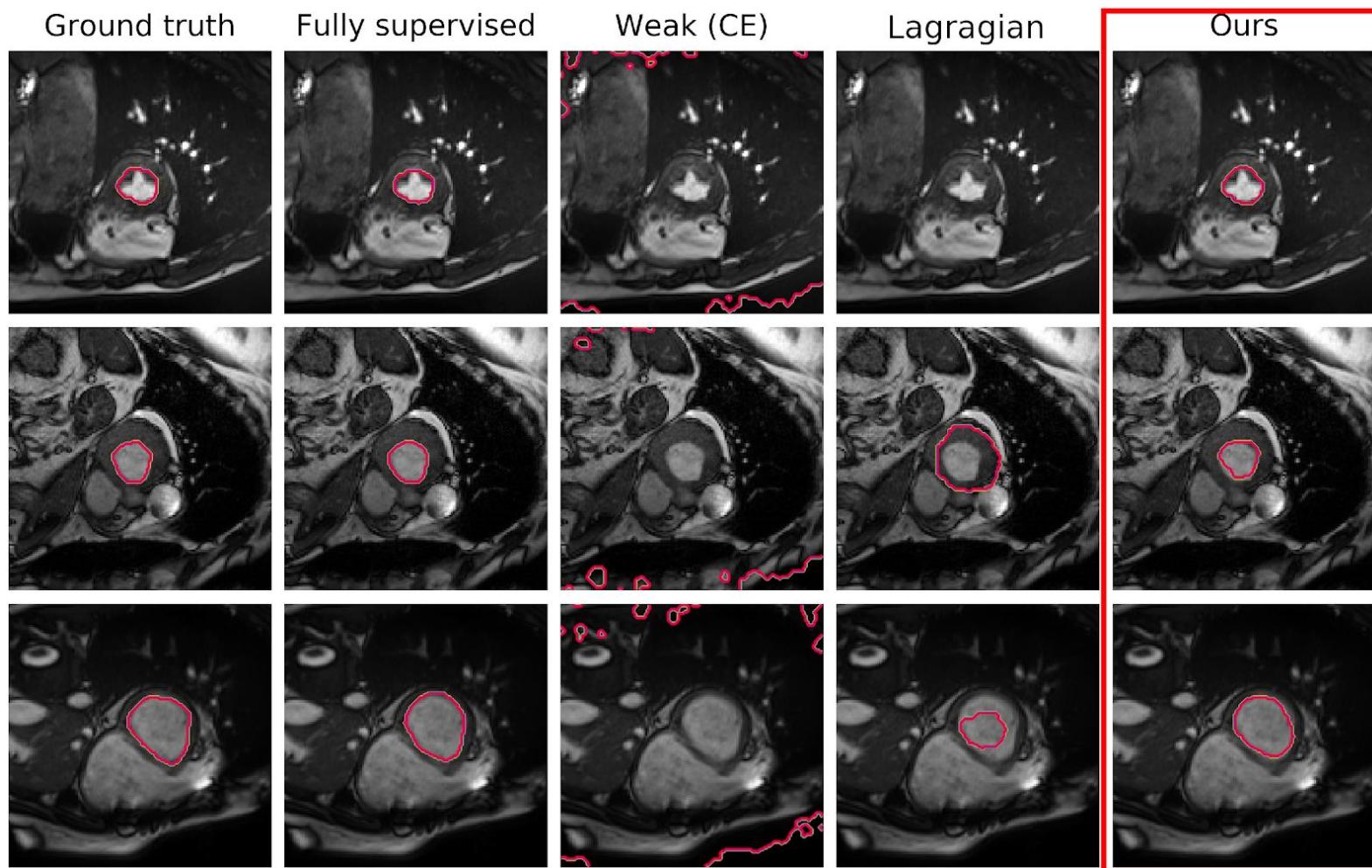
Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints

$$\min \leq \sum_{p \in \Omega} s_{\theta}^{p,c} \leq \max$$



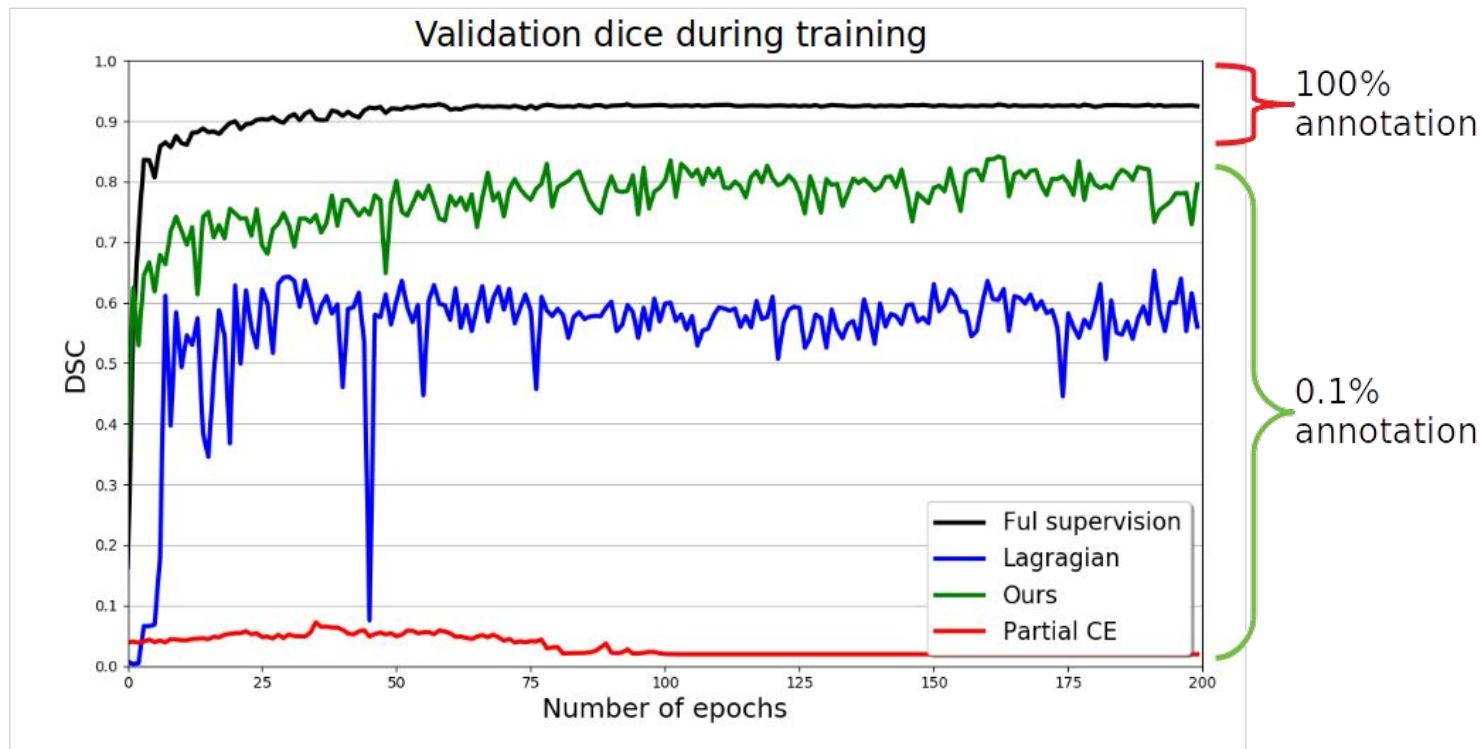
Kervadec et al., Constrained-CNN Losses for Weakly Supervised Segmentation,
Media'19

Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints



Kervadec et al., Constrained-CNN Losses for Weakly Supervised Segmentation,
Media'19

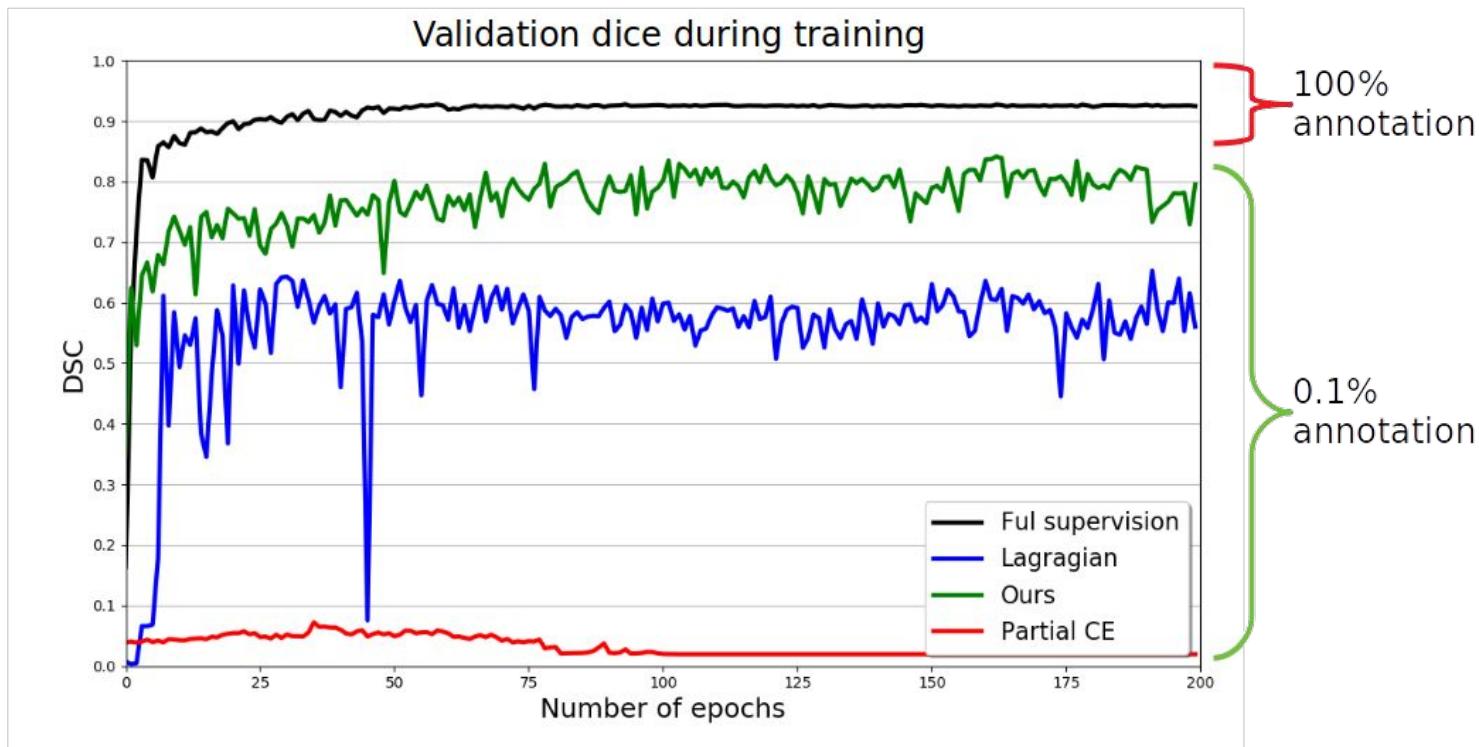
Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints



Kervadec et al., Constrained-CNN Losses for Weakly Supervised Segmentation,
Media'19

Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints

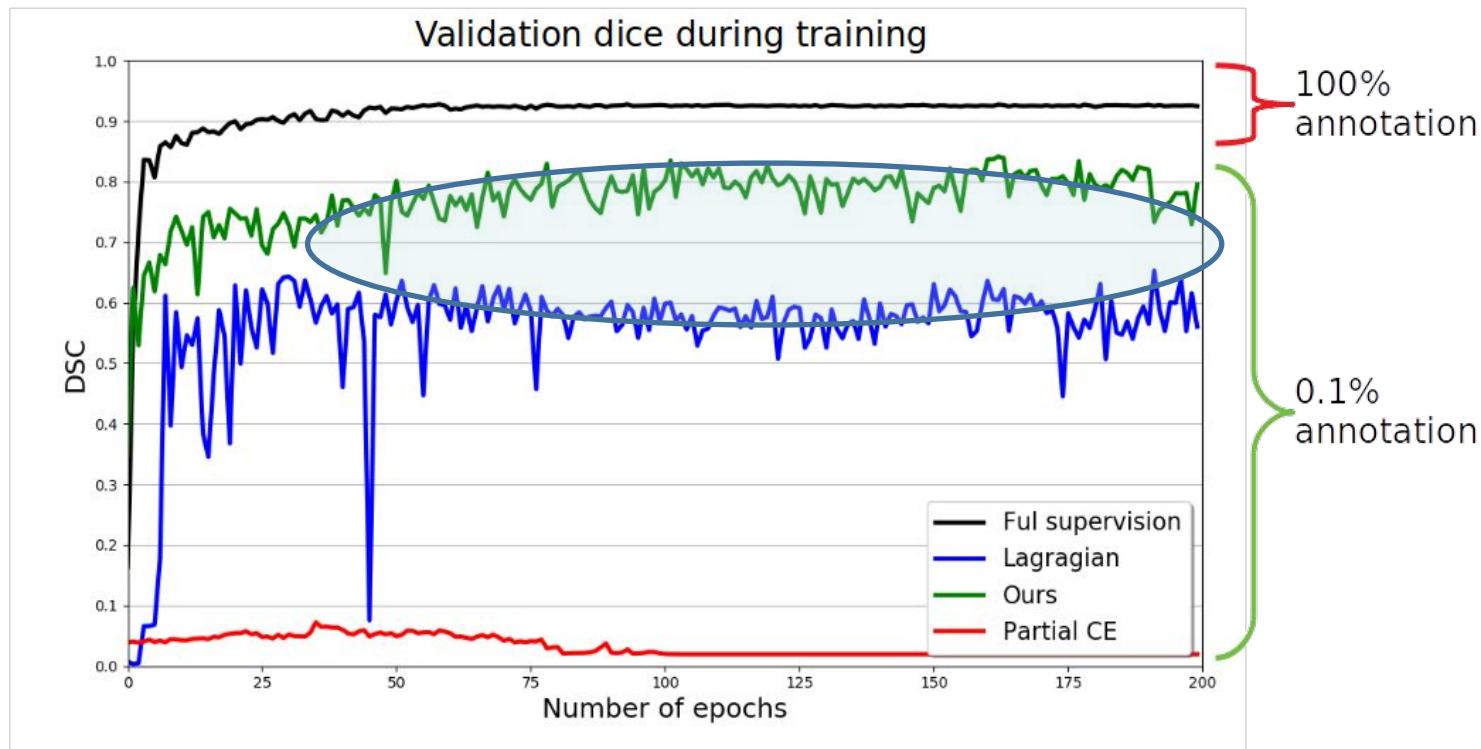
The exciting part: 90% of full supervision Dice with 0.1% of labels



Kervadec et al., Constrained-CNN Losses for Weakly Supervised Segmentation,
Media'19

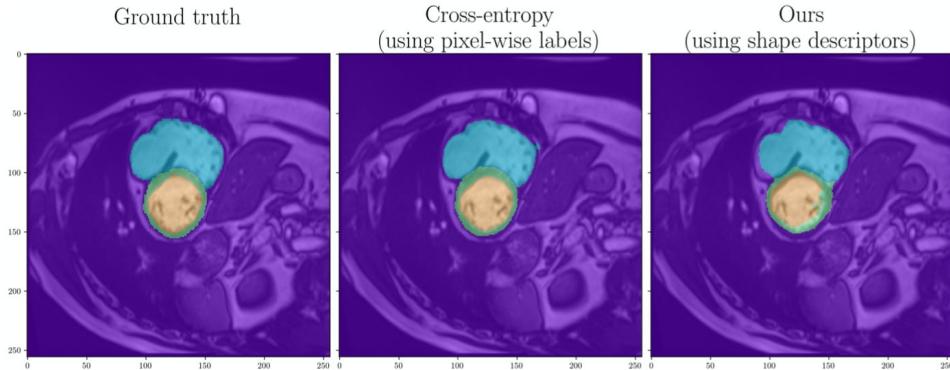
Example: Left Ventricle Segmentation in Cardiac MRI with Volumetric Constraints

The surprising part: Lagrangian optimization is much worse than a simple penalty



Kervadec et al., Constrained-CNN Losses for Weakly Supervised Segmentation,
Media'19

Beyond size: Exploring shape priors



(a) A visual comparison of the different supervision methods on the ACDC dataset.

Pixel	Label
0	RV
1	BACKGROUND
2	LV
\vdots	
65536	BACKGROUND

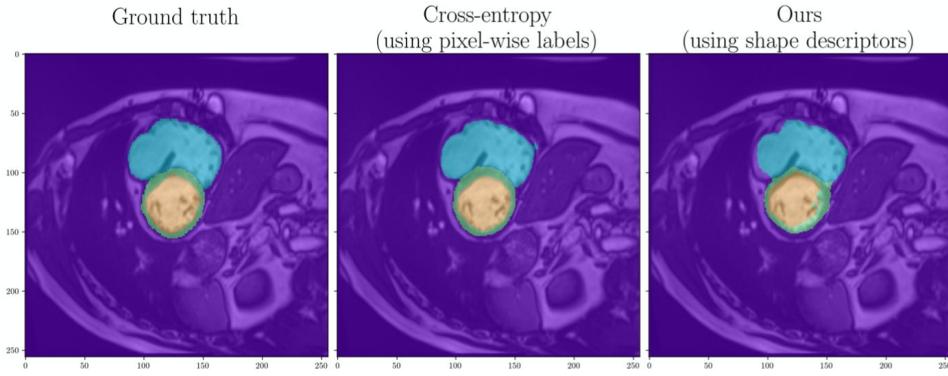
(b) Pixel-wise labels
(65k discrete values)

Shape descriptor	Class		
	RV	MYO	LV
Object volume \mathfrak{V}	3100	800	1600
Centroid location \mathfrak{C}	(125, 80)	(125, 125)	
Avg. dist. to centroid \mathfrak{D}	(20, 15)	(15, 20)	(10, 10)
Object length \mathfrak{L}	750	1000	500

(c) Shape descriptors
(16 continuous values)

Kervadec et al., Beyond pixelwise supervision: A few shape descriptors might be surprisingly good! MIDL'21

Beyond size: Exploring shape priors



(a) A visual comparison of the different supervision methods on the ACDC dataset.

Pixel	Label
0	RV
1	BACKGROUND
2	LV
⋮	
65536	BACKGROUND

(b) Pixel-wise labels
(65k discrete values)

Shape descriptor (in pixels)	Class		
	RV	MYO	LV
Object volume \mathfrak{V}	3100	800	1600
Centroid location \mathfrak{C}	(125, 80)	(125, 125)	
Avg. dist. to centroid \mathfrak{D}	(20, 15)	(15, 20)	(10, 10)
Object length \mathfrak{L}	750	1000	500

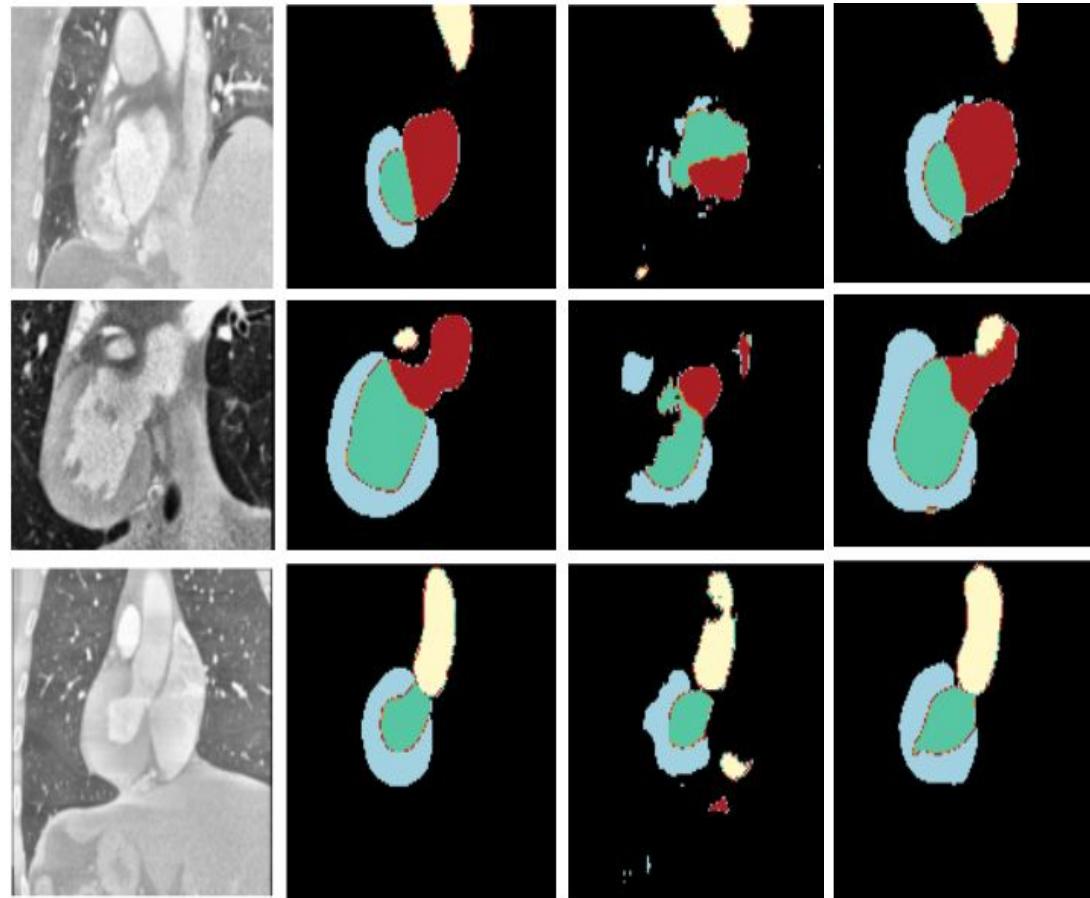
(c) Shape descriptors
(16 continuous values)

$$\mu_{i,j}^c = \sum_p s_\theta^{p,c} x_p^i y_p^j$$

↓
Spatial coordinates

Kervadec et al., Beyond pixelwise supervision: A few shape descriptors might be surprisingly good! MIDL'21

Shape moments: Also powerful in test-time adaptation



*Loss optim. w.r.t scale and
bias param. of batch norm
layers*

Bateson et al., Test-Time Adaptation with Shape Moments for Image Segmentation,
MICCAI 2022

Part 5:

Key challenges and future directions

Part 5. Key challenges and future directions

Take away message



- WSOL is mostly done via CAMs

Part 5. Key challenges and future directions

Take away message



- WSOL is mostly done via CAMs
- Bottom-up methods are dominant

Part 2. Review of WSOL methods: Literature

Take away message



- WSOL is mostly done via CAMs
- Bottom-up methods are dominant
- What currently works better:
 - Leveraging low level features
 - Pseudo-labels

Part 5. Key challenges and future directions

Take away message

Ongoing issues of CAMs



Part 5. Key challenges and future directions

Take away message

Ongoing issues of CAMs

- Cover full discriminative objects



Part 5. Key challenges and future directions

Take away message

Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)



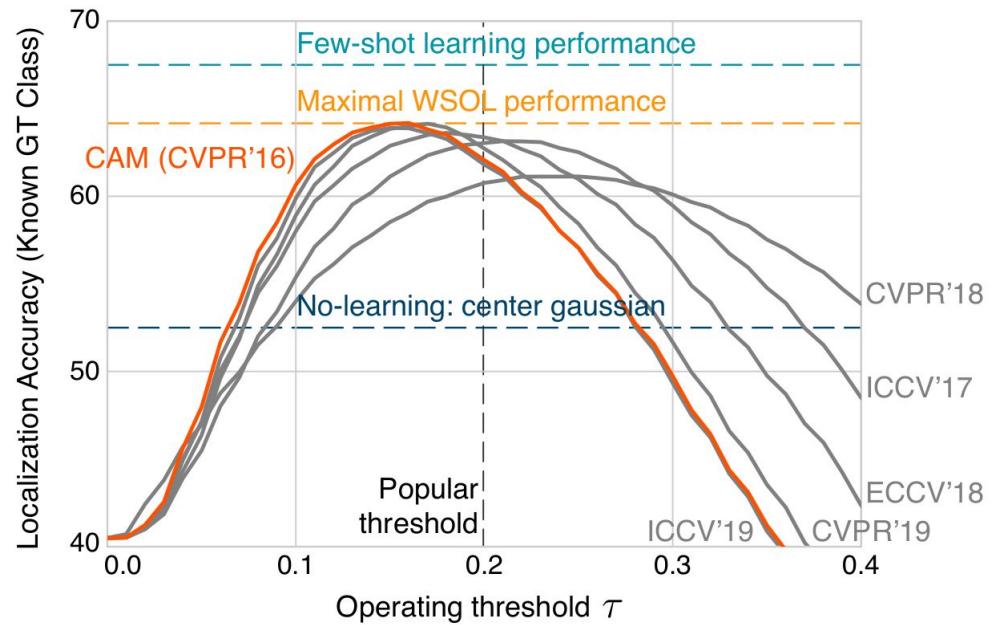
Part 5. Key challenges and future directions

Take away message



Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence



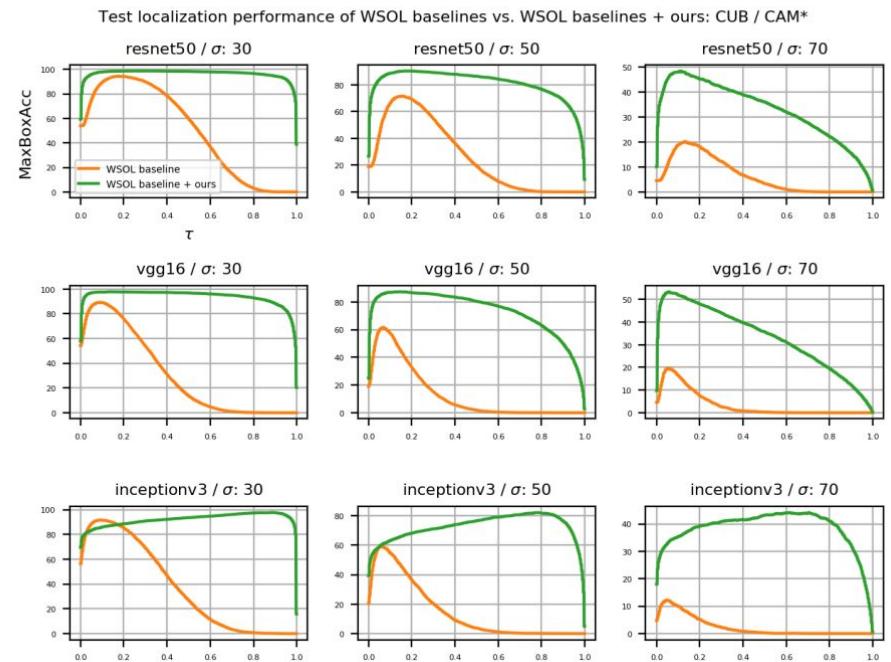
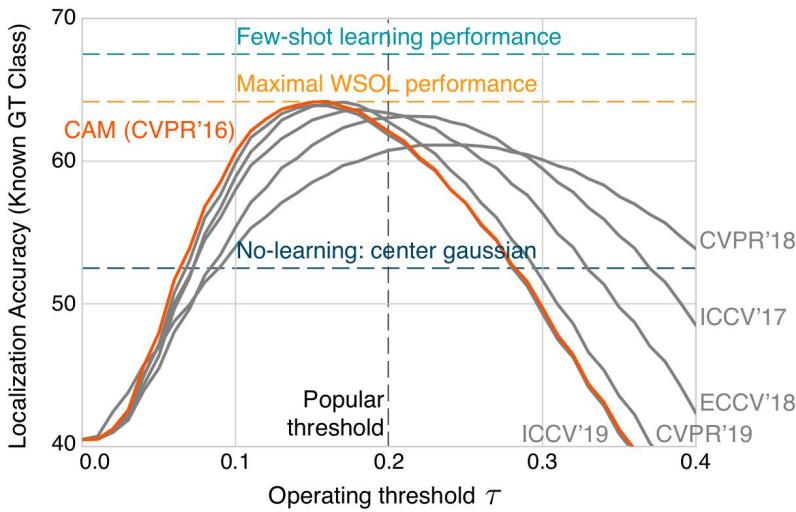
Part 5. Key challenges and future directions

Take away message



Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence

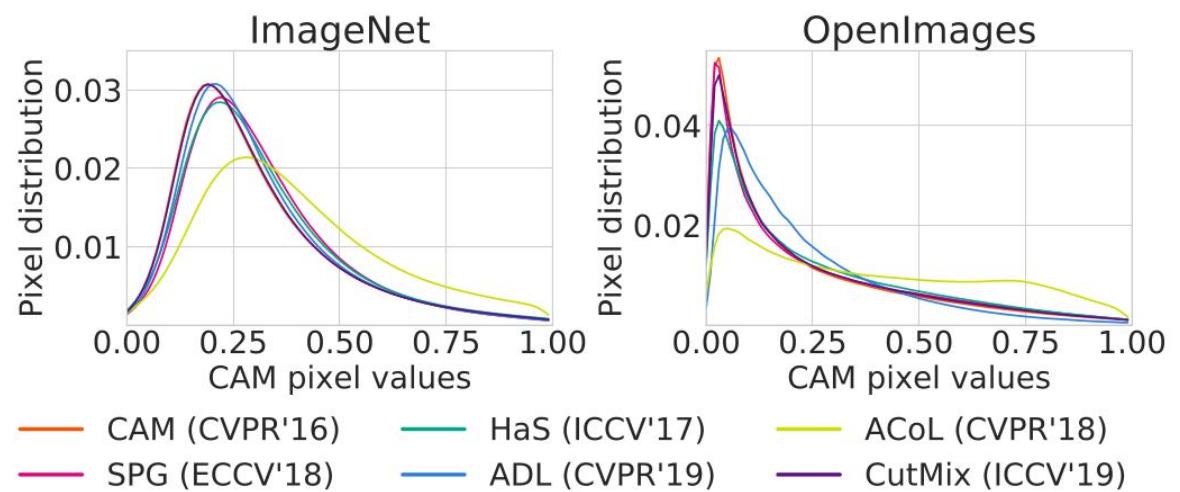


Part 5. Key challenges and future directions

Take away message

Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence
- Uncertainty



Part 5. Key challenges and future directions

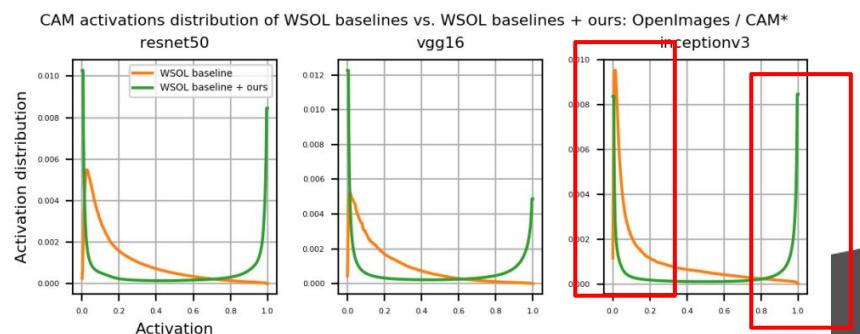
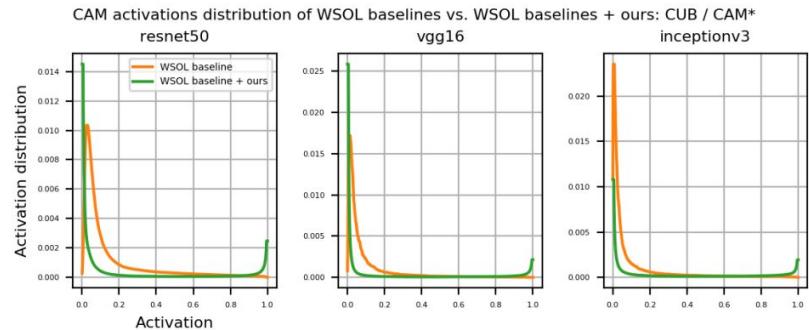
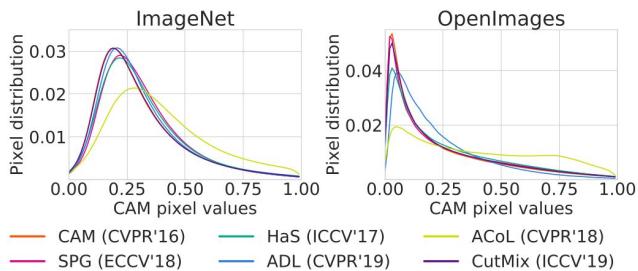
Take away message



Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence
- Uncertainty

2 modes



Part 5. Key challenges and future directions

Take away message

Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence
- Uncertainty
- Tiny objects



Part 2. Review of WSOL methods: Literature

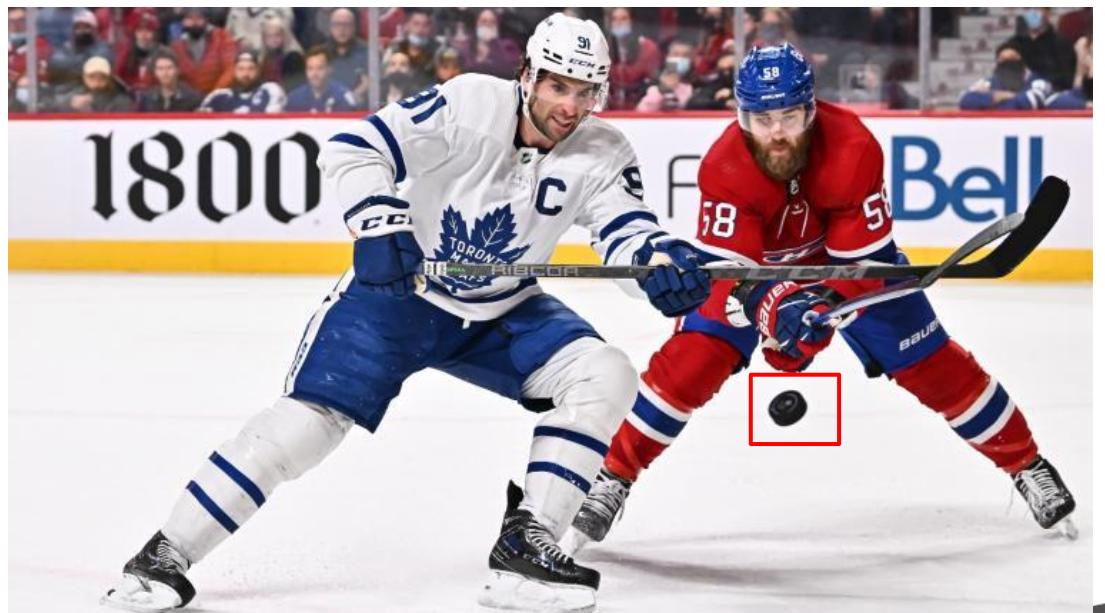
Take away message

Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence
- Uncertainty
- Tiny objects



Class: 'Puck'



Part 5. Key challenges and future directions

Take away message



Ongoing issues of CAMs

- Cover full discriminative objects
- Deal with background (complex scene, non-salient objects)
- Threshold dependence
- Uncertainty
- Tiny objects
- Objects co-occurrence

Class: 'Puck'



Part 5. Key challenges and future directions

Take away message

WSOL datasets: Saturation

- CUB dataset: ~97% MaxBoxAcc
- Imagenet-1k dataset: ~70% MaxBoxAcc

Code

- **Deep Weakly-Supervised Learning Methods for Classification and Localization in Histology Images: A Survey.** Rony et al. 2022. **Code:** https://github.com/jeromerony/survey_wsl_histology
- **F-cam: Full resolution class activation maps via guided parametric upscaling.** Belharbi, S., Sarraf, A., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). In WACV. **Code:** <https://github.com/sbelharbi/fcam-wsol>
- **Negative evidence matters in interpretable histology image classification.** Belharbi, S., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). In Medical Imaging with Deep Learning (MIDL). **Code:** <https://github.com/sbelharbi/negev>
- **Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty.** Belharbi, S., Rony, J., Dolz, J., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). IEEE Transactions on Medical Imaging, 41:702–714. **Code:** <https://github.com/sbelharbi/deep-wsl-histo-min-max-uncertainty>
- **Convolutional stn for weakly supervised object localization and beyond.** Meethal, A., Pedersoli, M., Belharbi, S., and Granger, E. (2020). In ICPR. **Code:** <https://github.com/akhilpm/ConvSTN>
- **TCAM: Temporal Class Activation Maps for Object Localization in Weakly-Labeled Unconstrained Videos.** Belharbi, S., Ben Ayed, I., McCaffrey, L., and Granger, E.. 2022. **Code (soon!):** <https://github.com/sbelharbi/tcam-wsol-video>
- **Holistic Guidance for Occluded Person Re-Identification.** M. Kiran, R G. Praveen, L. T. Nguyen-Meidine, S. Belharbi, L.-A. Blais-Morin, E. Granger. BMVC 2021. **Code:** <https://github.com/madhukiranets/HolisticGuidanceOccReID2>
- **Deep Active Learning for Joint Classification & Segmentation with Weak Annotator.** S. Belharbi, I. Ben Ayed, L. McCaffrey, E. Granger. WACV 2021. **Code:** <https://github.com/sbelharbi/deep-active-learning-for-joint-classification-and-segmentation-with-weak-annotator>
- **Constrained-CNN Losses for Weakly Supervised Segmentation.** Kervadec et et al, Media'19. **Code:** https://github.com/LIVIAETS/SizeLoss_WSS
- **Beyond pixelwise supervision: A few shape descriptors might be surprisingly good!.** Kervadec et et al., MIDL'21. **Code:** https://github.com/hkervadec/shape_descriptors
- **Evaluating Weakly Supervised Object Localization Methods Right.** J. Choe, S. Joon Oh, S. Lee, S. Chun, Z. Akata, H. Shim. CVPR 2020. **Code:** <https://github.com/clovaai/wsolevaluation>

References

- Bearman, Amy, et al. "**What's the point: Semantic segmentation with point supervision.**" ECCV 2016.
- Z. Zhou. '**A brief introduction to weakly supervised learning.**' National Science Review, 5(1):44–53, 2018.
- V. Cheplygina et al., '**Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis,**' Medical Image Analysis, 2019.
- **Learning Deep Features for Discriminative Localization.** Zhou et al. CVPR, 2016
- **Evaluating Weakly Supervised Object Localization Methods Right.** Choe et al. CVPR 2020.
- **Deep Weakly-Supervised Learning Methods for Classification and Localization in Histology Images: A Survey.** Rony et al. 2022.
- Lin, M., Chen, Q., and Yan, S. (2013). **Network in network.** coRR, abs/1312.4400.
- Durand, T., Mordan, T., Thome, N., et al. (2017). **Wildcat: Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation.** In CVPR.
- Singh, K. and Lee, Y. (2017). **Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization.** In ICCV.
- Li, K., Wu, Z., Peng, K., et al. (2018). **Tell me where to look: Guided attention inference network.** In CVPR.
- Choe, J. and Shim, H. (2019). **Attention-based dropout layer for weakly supervised object localization.** In CVPR.
- Xue, H., Liu, C., Wan, F., Jiao, J., Ji, X., and Ye, Q. (2019). **Danet: Divergent activation for weakly supervised object localization.** In ICCV.
- Zhang, C., Cao, Y., and Wu, J. (2020a). **Rethinking the route towards weakly supervised object localization.** In CVPR.
- X.-S. Wei, C.-L. Zhang, J. Wu, C. Shen, and Z.-H. Zhou. **Unsupervised object discovery and co-localization by deep descriptor transformation.** PR, 88:113–126, 2019.
- Wei, J., Wang, Q., Li, Z., Wang, S., Zhou, S. K., and Cui, S. (2021). **Shallow feature matters for weakly supervised object localization.** In CVPR.
- Cao, C., Liu, X., Yang, Y., et al. (2015). **Look and think twice: Capturing top-down visual attention with feedback convolutional neural networks.** In ICCV.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). **Grad-cam: Visual explanations from deep networks via gradient-based localization.** In ICCV.
- Wang, H., Wang, Z., Du, M., Yang, F., Zhang, Z., Ding, S., Mardziel, P., and Hu, X. (2020). **Score-cam: Score-weighted visual explanations for convolutional neural networks.** In CVPR workshop.
- **Spatial Transformer Networks (STN)**, M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, "Spatial transformer networks," in eurIPS, 2015.

References

- Meethal, A., Pedersoli, M., Belharbi, S., and Granger, E. (2020). **Convolutional stn for weakly supervised object localization and beyond.** In ICPR.
- Belharbi, S., Sarraf, A., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **F-cam: Full resolution class activation maps via guided parametric upscaling.** In WACV.
- F. Yu, V. Koltun, and T. Funkhouser, **Dilated residual networks**, CVPR 2017
- Oquab, M., et al., **Is object localization for free?-weakly-supervised learning with CNNs.** In CVPR 2015
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., & Joulin, A. (2021). **Emerging properties in self-supervised vision transformers.** In ICCV.
- Gao, W., Wan, F., Pan, X., Peng, Z., Tian, Q., Han, Z., ... & Ye, Q. (2021). **Ts-cam: Token semantic coupled attention map for weakly supervised object localization.** In ICCV..
- Figures from [Zhang et al., **A Curriculum Domain Adaptation Approach to the Semantic Segmentation of Urban Scenes**, TPAMI 2019
- Bateson et al., **Constrained Domain Adaptation for Image Segmentation**, TMI'21
- Lin et al. **Scribblesup: Scribble-supervised convolutional networks for semantic segmentation**, CVPR 2016
- Tang et al., **On regularized losses for weakly supervised segmentation**, ECCV 2018
- Qu et al., **Weakly Supervised Deep Nuclei Segmentation using Points Annotation in Histopathology Images**, MIDL 2019
- Ji et al., **Scribble-Based Hierarchical Weakly Supervised Learning for Brain Tumor Segmentation**, MICCAI 2019
- Grandvalet & Bengio, **Semi-supervised learning by entropy minimization**, NIPS 2005
- Gomes et al., **Discriminative clustering by regularized information maximization**, NIPS 2010
- **ADVENT: Adversarial Entropy Minimization for Domain Adaptation in Semantic Segmentation**, CVPR 2019
- **Source-relaxed domain adaptation for segmentation**, MICCAI 2020
- Belharbi, S., Rony, J., Dolz, J., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Deep interpretable classification and weakly-supervised segmentation of histology images via max-min uncertainty.** IEEE Transactions on Medical Imaging, 41:702–714.
- Belharbi, S., Pedersoli, M., Ben Ayed, I., McCaffrey, L., and Granger, E. (2022). **Negative evidence matters in interpretable histology image classification.** In Medical Imaging with Deep Learning (MIDL).
- A Bhuiyan, et al., **Pose Guided Gated Fusion for Person Re-identification**, WACV 2020.
- R. Selvaraju, et al., **Grad-CAM: Visual explanations from deep networks via gradient-based localization**. CVPR 2017.
- A Stylianou, R Souvenir, and R Pless, **Visualizing deep similarity networks**, WACV 2019.

References

- S Black, et al., "Visualizing Paired Image Similarity in Transformer Networks." WACV 2022
- L Chen, et al., Adapting grad-cam for embedding networks. WCCV 2020.
- Shen, Yeqing, et al. Distance-Based Class Activation Map for Metric Learning. PRCV 2021.
- M Ye, et al., Bi-directional center-constrained top-ranking for visible thermal person re-identification. IEEE TIFS 2019
- M Alehdaghi, et al., Visible-Infrared Person Re-Identification Using Privileged Intermediate Information. ECCVw 2022
- Kervadec et et al., Constrained-CNN Losses for Weakly Supervised Segmentation, MedIA'19
- Kervadec et et al., Beyond pixelwise supervision: A few shape descriptors might be surprisingly good! MIDL'21
- Bateson et et al., Test-Time Adaptation with Shape Moments for Image Segmentation, MICCAI 2022
- **TCAM: Temporal Class Activation Maps for Object Localization in Weakly-Labeled Unconstrained Videos.** 2022 [soon]

Discussion

Slides

These slides will be available at:

<https://sbelharbi.github.io/publications/icpr-tutorial-wsl-2022/slides.pdf>