



PUZZLE-CAM: IMPROVED LOCALIZATION VIA MATCHING PARTIAL AND FULL FEATURES (+RecurSeed and EdgePredictMix)

2023.01.12

Presenter : Taewoong Kang
twk@deepnoid.com

Overview

- **Weakly Supervised Semantic Segmentation (WSSS)**
- 이전은 CAM을 이용해 segmentation mask를 얻으려 refine 한 논문들이 많았다.
- 이 논문은 **reconstructing regularization** with a puzzle module을 포함한 CAM으로 Semantic Segmentation하는 것이 목표이다.

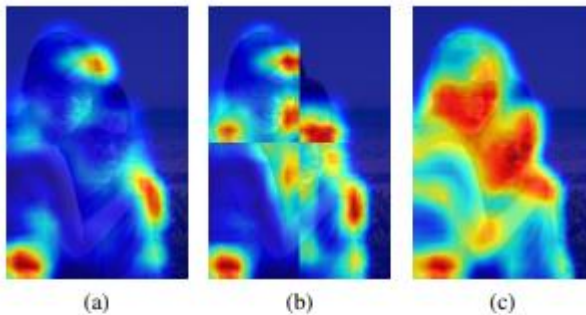


Fig. 1: CAMs generated from tiled and original images: (a) conventional CAMs from the original image, (b) generated CAMs from the tiled images, and (c) predicted CAMs by the proposed Puzzle-CAM.

- Puzzle-CAM architecture

- Comparison of CAM and Puzzle-CAM

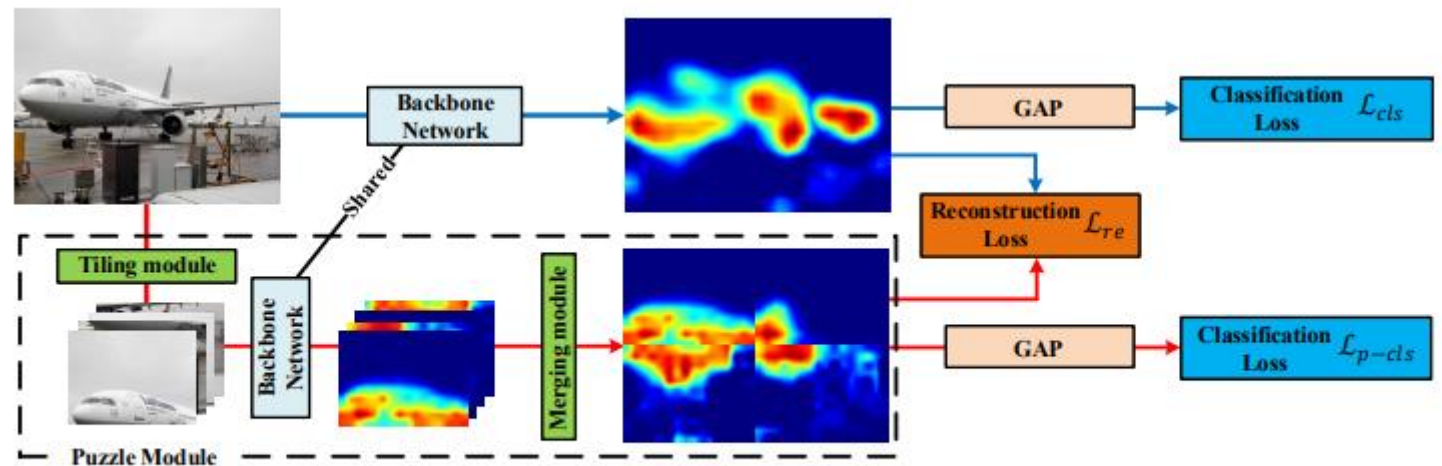


Fig. 2: The overall architecture of the proposed Puzzle-CAM showing the integration of reconstructing regularization and the puzzle module.

Overview

- Weakly supervised **object localization**
- Identify the importance of the image regions by **projecting back the weights** of the output layer on the convolutional feature maps
- 기존 Network와 유사한 Classification 정확도 (1~2% 하락)와 더불어 Localization까지 가능하다.
- (어디에 집중해서 Classification을 하는지와 연관)

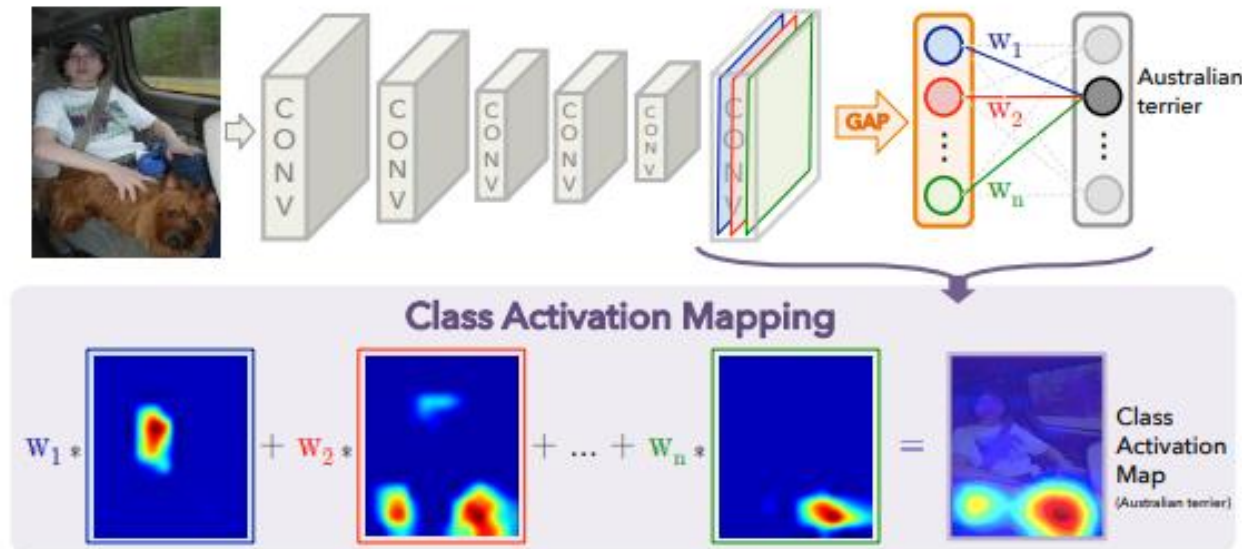
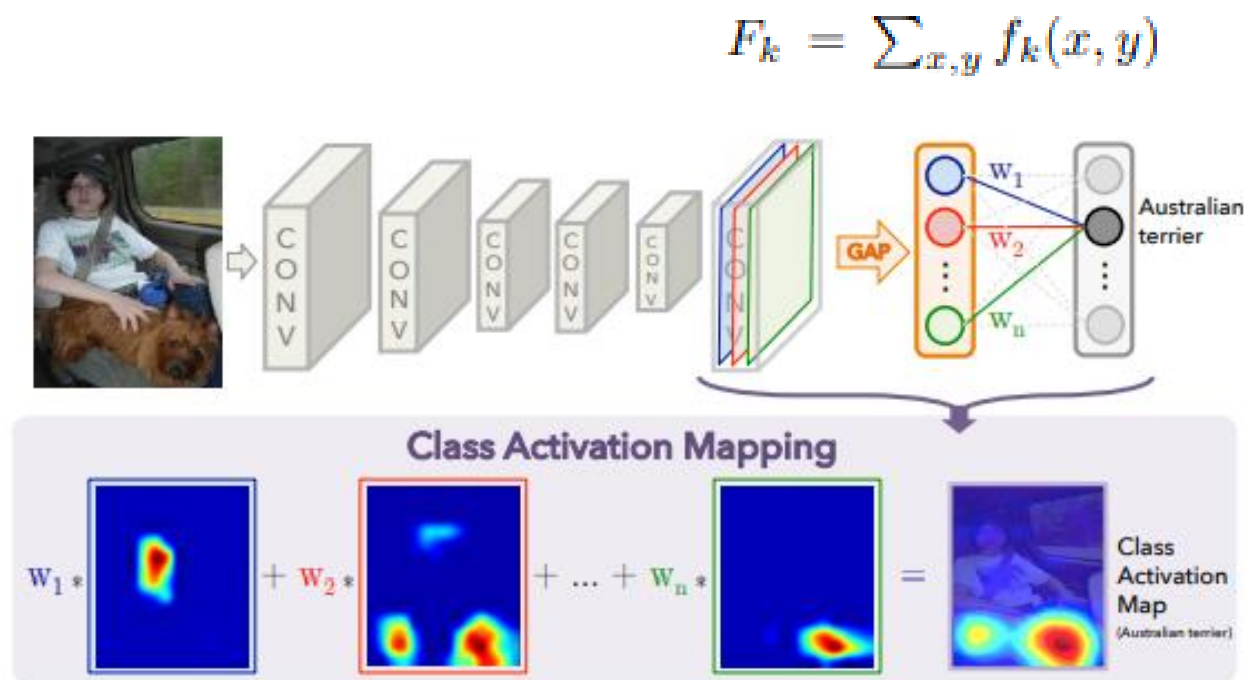


Figure 2. Class Activation Mapping: the predicted class score is mapped back to the previous convolutional layer to generate the class activation maps (CAMs). The CAM highlights the class-specific discriminative regions.

- Basic CAM Flow

Architecture

- Fully Connected Layer 이전에 GAP Layer 사용 (just before the final out-put layer)
= FC Layer → GAP + FC softmax Layer



$$S_c = \sum_k w_k^c \sum_{x,y} f_k(x, y) = \sum_{x,y} \sum_k w_k^c f_k(x, y)$$

$$M_c(x, y) = \sum_k w_k^c f_k(x, y)$$

- Basic CAM Flow

Figure 2. Class Activation Mapping: the predicted class score is mapped back to the previous convolutional layer to generate the class activation maps (CAMs). The CAM highlights the class-specific discriminative regions.

Localization

- Use a simple thresholding technique to segment the heatmap
- CAM의 max value의 20%이상만 사용.
- Segmentation map에서 largest connected component를 cover하는 Bbox 생성
 - But Top-5 test error 37.1
 - Annotated Bounding box로 train 되지 않았다는 것이 주목할 만한 점이다.

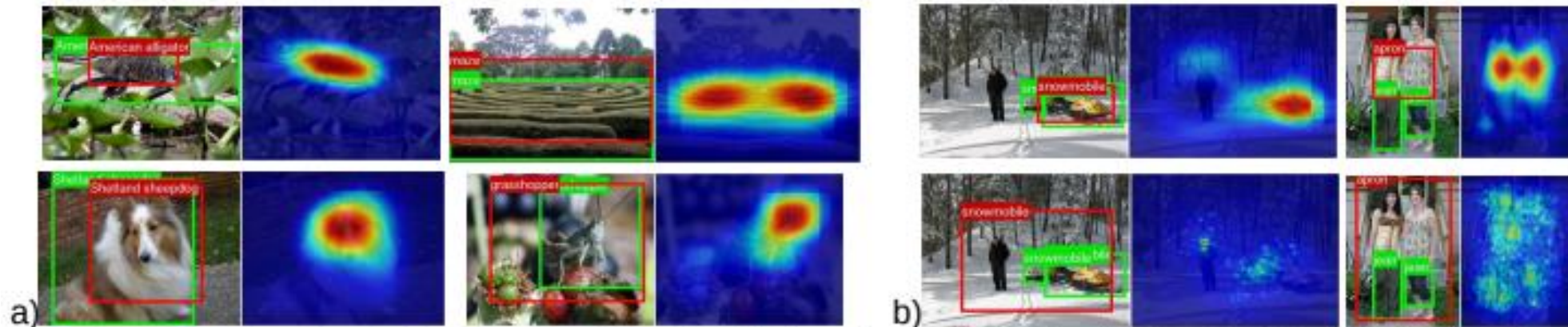


Figure 6. a) Examples of localization from GooleNet-GAP. b) Comparison of the localization from GooleNet-GAP (upper two) and the backpropagation using AlexNet (lower two). The ground-truth boxes are in green and the predicted bounding boxes from the class activation map are in red.

Limit

- Focused on small parts of the semantic objects to efficiently classify them
-> prevent the segmentation models from learning pixel-level semantic knowledge.

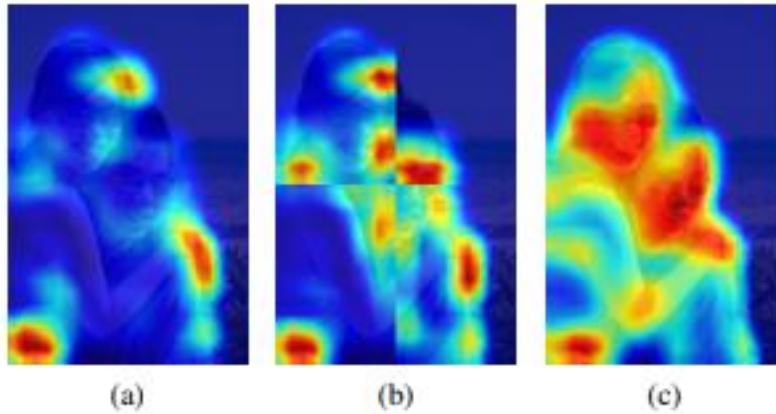


Fig. 1: CAMs generated from tiled and original images: (a) conventional CAMs from the original image, (b) generated CAMs from the tiled images, and (c) predicted CAMs by the proposed Puzzle-CAM.

Overview

- Image-level supervision (**WSSS**)
- Reconstructing regularization** with a puzzle module

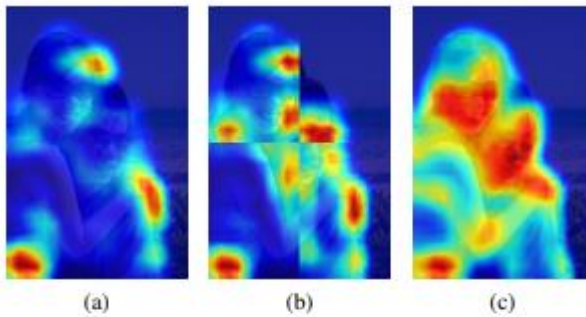


Fig. 1: CAMs generated from tiled and original images: (a) conventional CAMs from the original image, (b) generated CAMs from the tiled images, and (c) predicted CAMs by the proposed Puzzle-CAM.

- Puzzle-CAM architecture

- Comparison of CAM and Puzzle-CAM

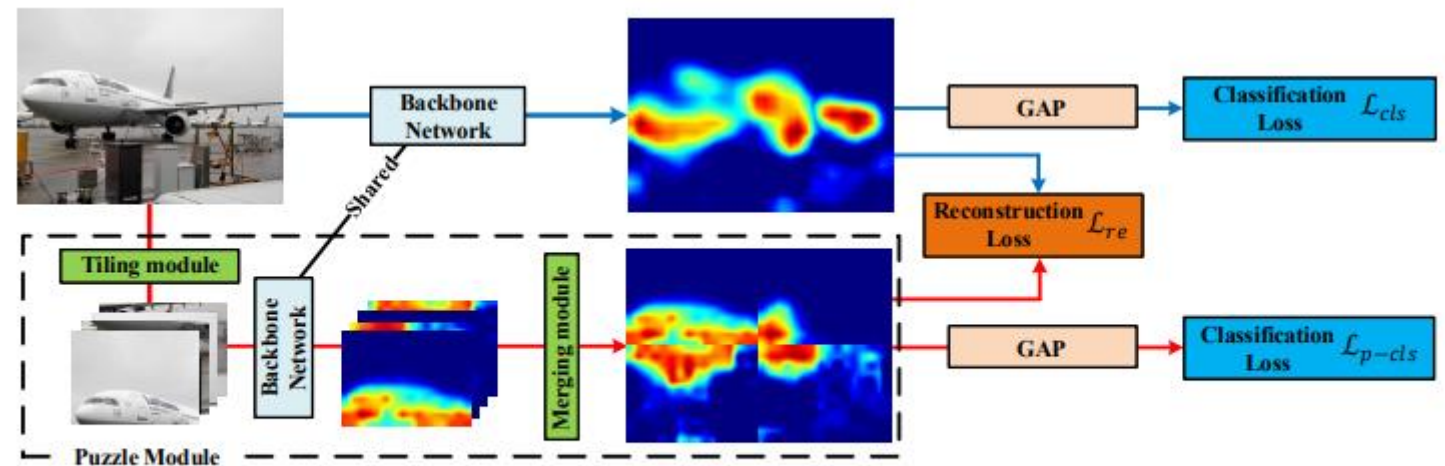


Fig. 2: The overall architecture of the proposed Puzzle-CAM showing the integration of reconstructing regularization and the puzzle module.

Overview

- 원래 이미지에 대한 CAM과, Puzzle Module을 이용한 CAM을 각각 구함
- Puzzle Module은 본 이미지를 4개로 나눈 후 각 puzzle에 대한 CAM을 구한 후 한 개로 합침. (두 개의 CAM은 같은 크기)
- 두 개의 CAM에 대한 **Classification Loss**, 두 CAM의 차이로 구한 **Reconstruction Loss**를 이용해 학습

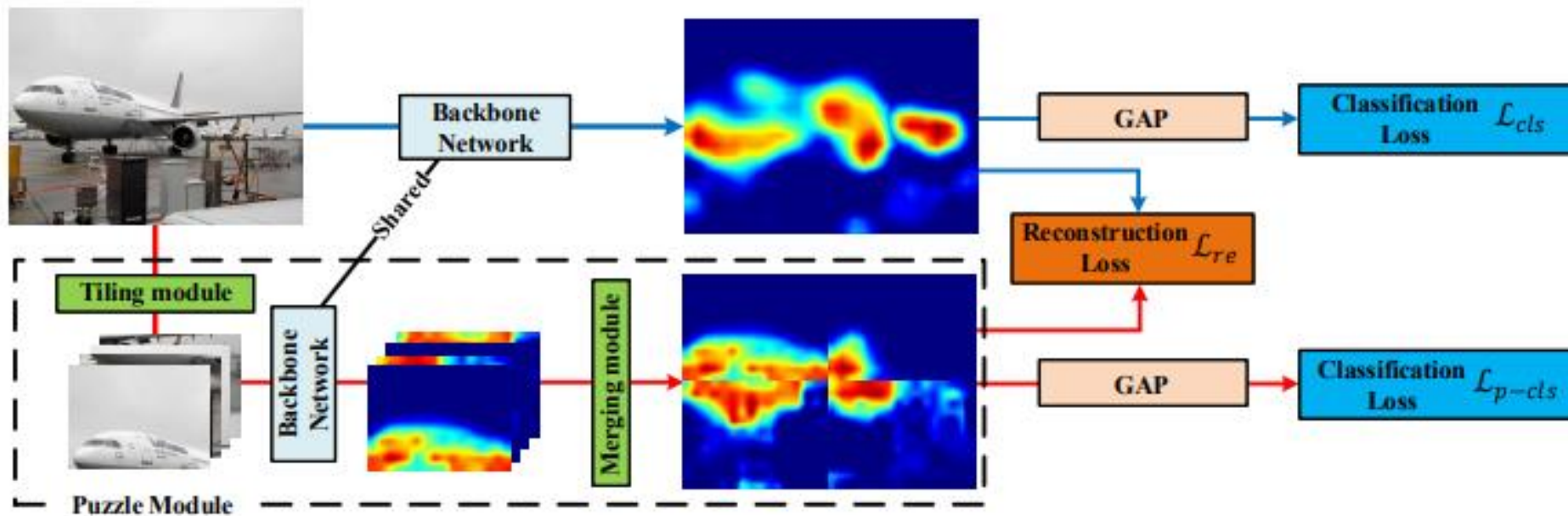
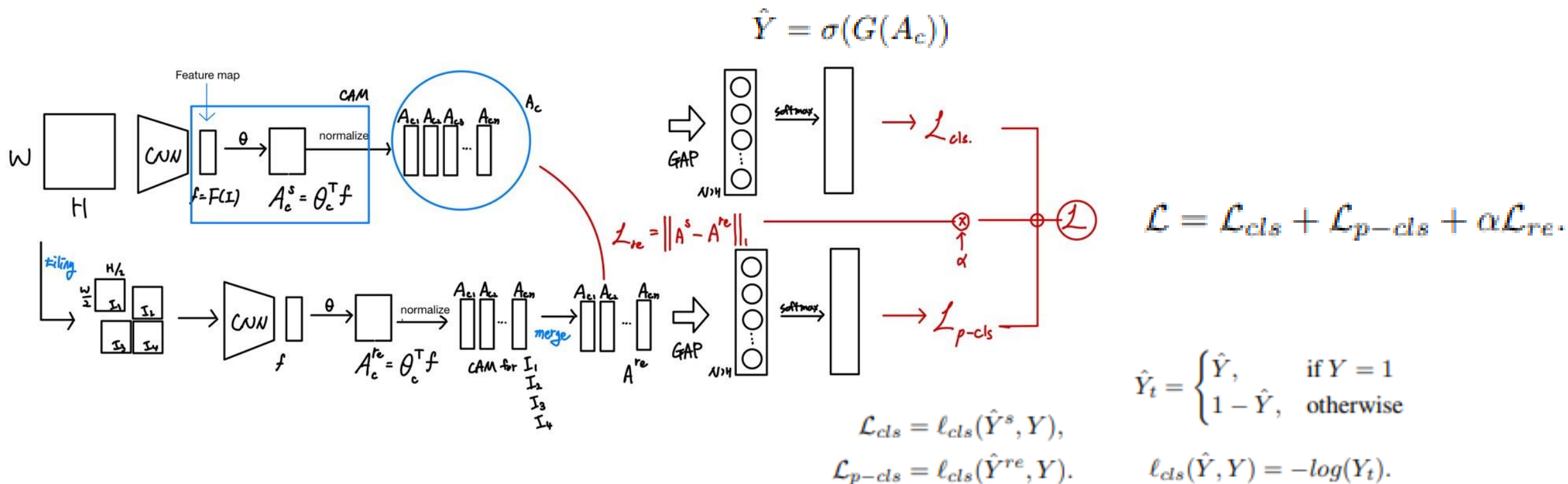


Fig. 2: The overall architecture of the proposed Puzzle-CAM showing the integration of reconstructing regularization and the puzzle module.

Loss function

- 원래 이미지에 대한 CAM과, Puzzle Module을 이용한 CAM을 각각 구함
- Puzzle Module은 본 이미지를 4개로 나눈 후 각 puzzle에 대한 CAM을 구한 후 한 개로 합침. (두 개의 CAM은 같은 크기)
- 두 개의 CAM에 대한 Classification Loss, 두 CAM의 차이로 구한 Reconstruction Loss를 이용해 학습



Result

Method	Backbone	Supervision	val	test
AffinityNet [4]	Wide-ResNet-38	\mathcal{I}	61.7	63.7
DSRG [12]	ResNet-101	$\mathcal{I} + \mathcal{S}$	61.4	63.2
SeeNet [13]	ResNet-101	$\mathcal{I} + \mathcal{S}$	63.1	62.8
IRNet [4]	ResNet-50	\mathcal{I}	63.5	64.8
FickleNet [6]	ResNet-101	$\mathcal{I} + \mathcal{S}$	64.9	65.3
ICD [17]	ResNet-101	\mathcal{I}	64.1	64.3
SEAM [5]	Wide-ResNet-38	\mathcal{I}	64.5	65.7
Ours (Puzzle-CAM)	ResNeSt-101	\mathcal{I}	66.9	67.7
Ours (Puzzle-CAM)	ResNeSt-269	\mathcal{I}	71.9	72.2

여러 WSSS 기법 들과의 비교

L_{cls}	L_{p-cls}	L_{re}	mIoU (%)
✓			47.82
✓	✓		47.70
✓		✓	49.21
✓	✓	✓	51.53

Loss 값 유무에 따른 성능의 변화



Fig. 4: Qualitative segmentation results on the PASCAL VOC 2012 *val* set. Top: original images. Middle: ground truth. Bottom: prediction of the segmentation model trained using the pseudo-labels from Puzzle-CAM.

Limit

- Represent partial regions for large-scale objects
- For small-scale objects, over activation causes them to deviate from the object edges

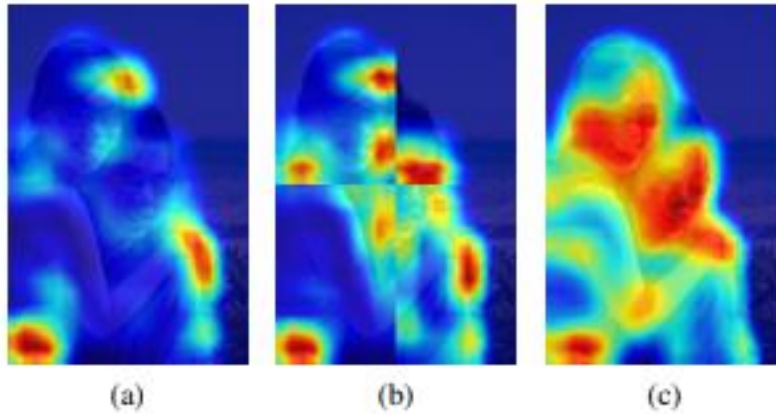
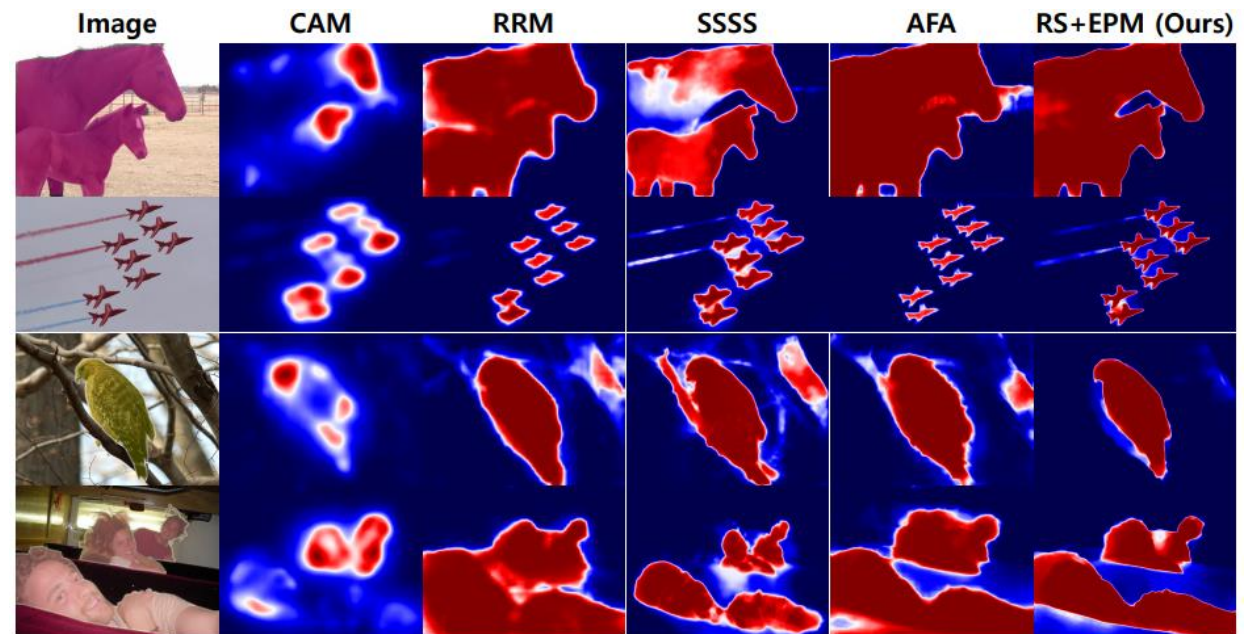


Fig. 1: CAMs generated from tiled and original images: (a) conventional CAMs from the original image, (b) generated CAMs from the tiled images, and (c) predicted CAMs by the proposed Puzzle-CAM.

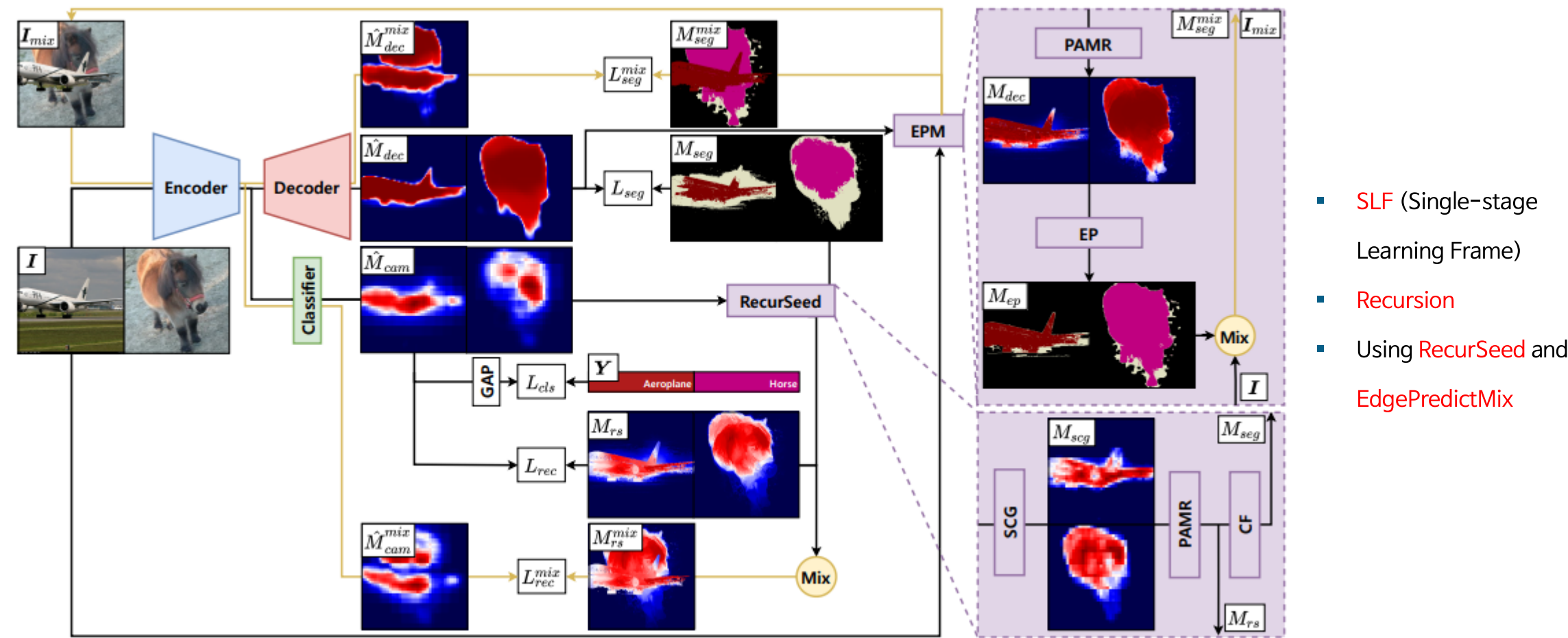
Overview

- To solve two problem (in WSSS)
- **FPN**-WSSS-IL \rightarrow mismatching of FNs and FPs
(SCG considering high-order correlation reduces FNs, but increases the number of FPs)
 \rightarrow Using Recurseed (both SCG and PAMR)

- **IBDA**-WSSS-IL \rightarrow the simple synthesis (Cutout, CutMix, SaliencyGrafting, CDA and ClassMix) without any refinements using predicted masks inevitably accelerates the ambiguity of mixed result due to insufficient spread
 \rightarrow using EPM



Overview



Algorithm Flow

- Self-correlation 을 통해 CAM을 보정

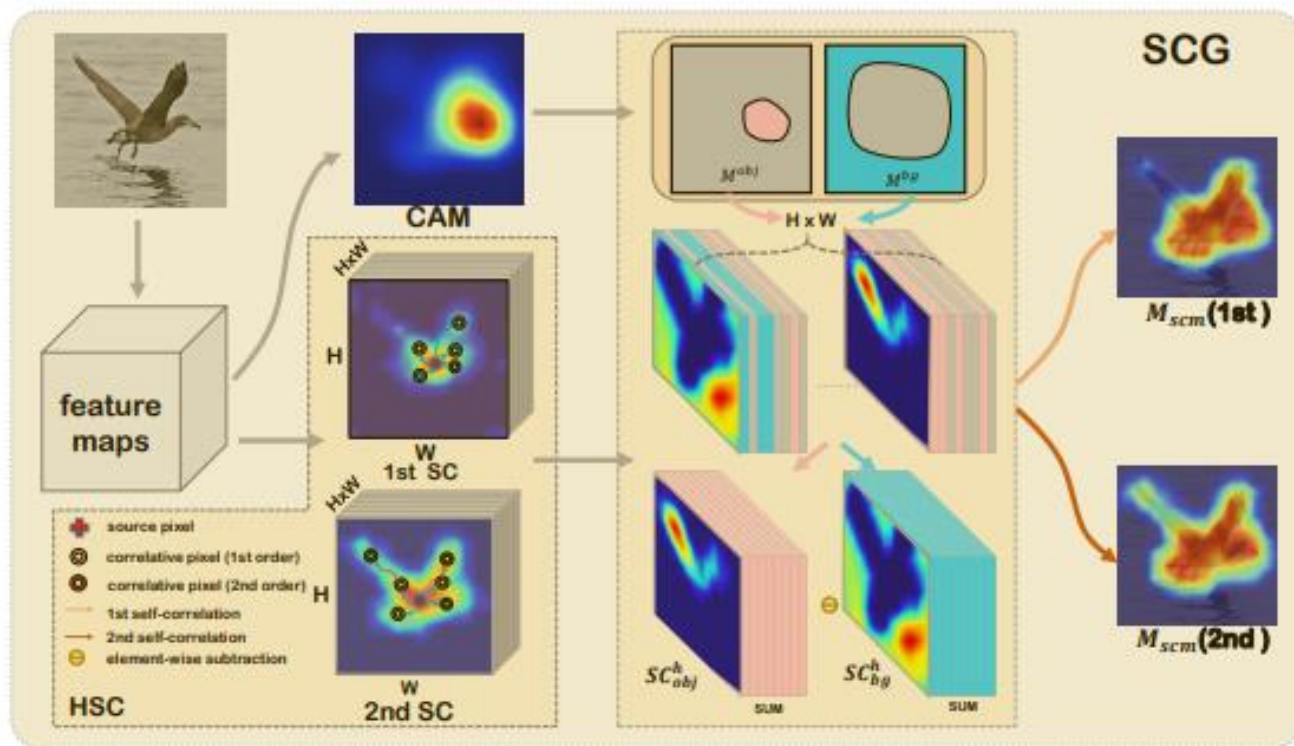
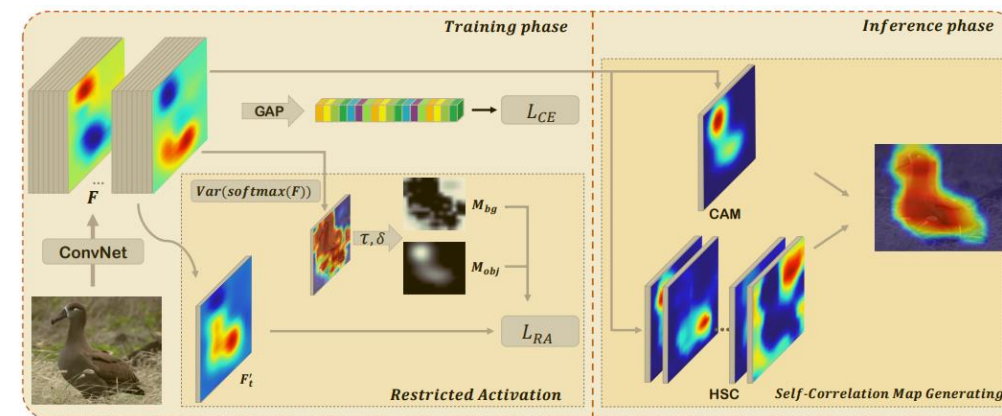
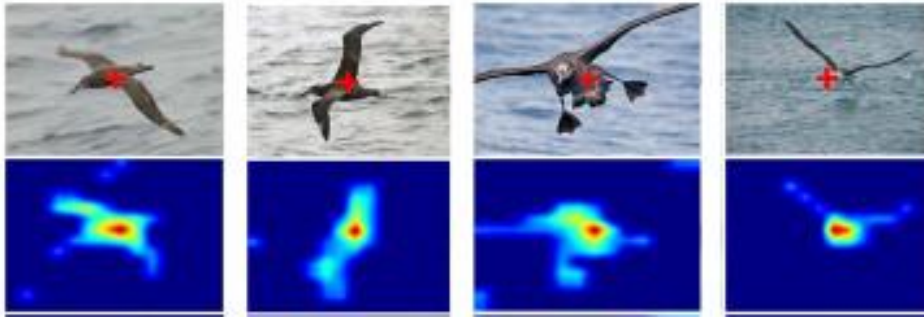


Figure 3. Pipeline of the proposed SCG module. Here we show examples of using first- and second-order SC to obtain final localization maps, respectively.



Algorithm Flow

First-order Self-correlation

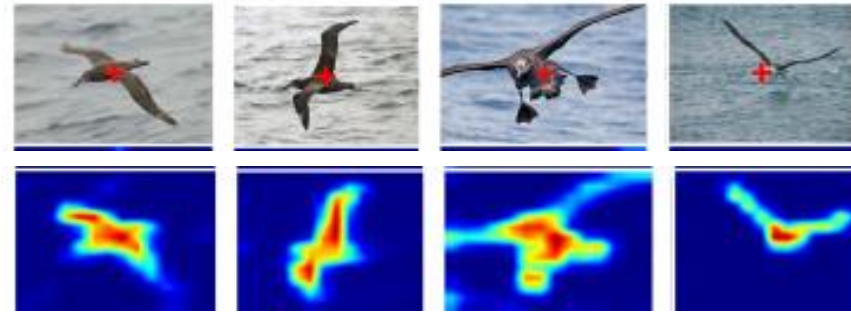


$$S(f_i, f_j) = \frac{f_i^T f_j}{\|f_i\| \cdot \|f_j\|},$$

$$SC^1(f) = [SC^1(f)_{i,j}],$$

$$\text{where } SC^1(f)_{i,j} = \text{ReLU}(S(f_i, f_j)).$$

Second-order Self-correlation



$$S^2(f_i, f_j) = \frac{1}{(HW)} \sum_{k \in \Omega} S(f_i, f_k) \cdot S(f_k, f_j),$$

$$\hat{S}^2(f_i, f_j) = \frac{S^2(f_i, f_j) - \min_{k \in \Omega} S^2(f_i, f_k)}{\max_{k \in \Omega} S^2(f_i, f_k) - \min_{k \in \Omega} S^2(f_i, f_k)},$$

$$SC^2(f) = [\hat{S}^2(f_i, f_j)|_{i,j}].$$

Algorithm Flow

f_i, f_j : CAM features of specific index
 $f_i, f_j \in R^{C \times 1}$

First-order Self-correlation

$$SC_l^1[i, j, :, :] = \text{ReLU}\left(\frac{f_i^{l\top} f_j^l}{\|f_i^l\| \|f_j^l\|}\right),$$

Second-order Self-correlation

$$SC_l^2[i, :, :, :] = \text{mnmx}_j \left(\text{avg}_k \left[\frac{f_i^{l\top} f_k^l}{\|f_i^l\| \|f_k^l\|} \odot \frac{f_k^{l\top} f_j^l}{\|f_k^l\| \|f_j^l\|} \right] \right),$$

merge

$$HSC = \frac{1}{L} \sum_{l=1}^L \max(SC_l^1, SC_l^2)$$

Adapt to CAM

$$\begin{aligned} M_{scg} &= SCG(M_{cam}) \\ &= \text{ReLU}(K_{scg}(\{M_{cam}\}_{>\delta_h}) - K_{scg}(\{M_{cam}\}_{<\delta_l})). \end{aligned}$$

$$K_{scg}(\{M\}_{\leq \delta}) = \frac{1}{|\{M\}_{\leq \delta}|} \sum_{(i,j) \in \{M\}_{\leq \delta}} HSC[i, j, :, :].$$

Algorithm 1 Localization algorithm of SCG.

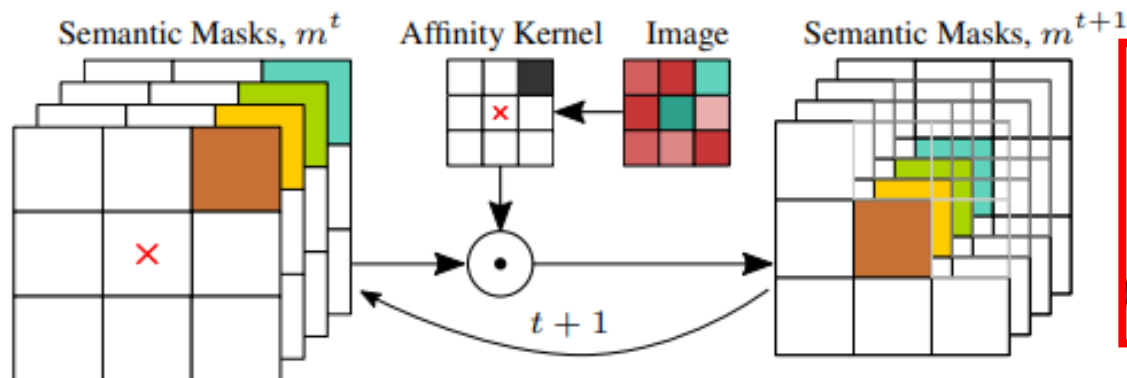
Input: Coarse localization map $M_{cam} \in \mathbb{R}^{H \times W}$; feature map $f \in \mathbb{R}^{H \times W \times C}$; threshold δ_h and δ_l ;

Output: Final localization map M_{scg} ;

- 1: Obtain high-order self-correlation $HSC \in \mathbb{R}^{HW \times HW}$
- 2: Reshape $HSC \in \mathbb{R}^{H \times W \times H \times W} \leftarrow \text{reshape}(HSC)$
- 3: Discover the coarse object region $M_{cam}^{obj} \leftarrow M_{cam} > \delta_h$
- 4: Extract object HSC $HSC_{obj} \leftarrow G(HSC, M_{cam}^{obj})$
- 5: Obtain the object map $M_{scg}^{obj} \leftarrow \text{sum}(HSC_{obj})$
- 6: Discover background region $M_{cam}^{bg} \leftarrow M_{cam} < \delta_l$
- 7: Extract background HSC $HSC_{bg} \leftarrow G(HSC, M_{cam}^{bg})$
- 8: Obtain the background map $M_{scg}^{bg} \leftarrow \text{sum}(HSC_{bg})$
- 9: Obtain localization map $M_{scg} \leftarrow (M_{scg}^{obj} - M_{scg}^{bg})_{(>0)}$
return M_{scg} ;

Algorithm Flow

- Local Consistency**: nearby regions sharing the same appearance should be assigned to the same class
- Reduced the computational complexity by narrowing the affinity kernel computation to regions of contiguous pixels
- Pixel-wise mask prediction $m_{:,i,:} \in (0, 1)^{(C+1) \times h \times w}$ (+1 for the background class)



$$M_{pamr} = PAMR(M_{cam}; \mathcal{W}) = G_T(M_{cam})$$

$$G_0(M_{cam}) = M_{cam}$$

$$G_t(M_{cam})[i, j] = \sum_{(k, n) \in \mathcal{N}(i, j)} \alpha_{i, j, k, n}$$

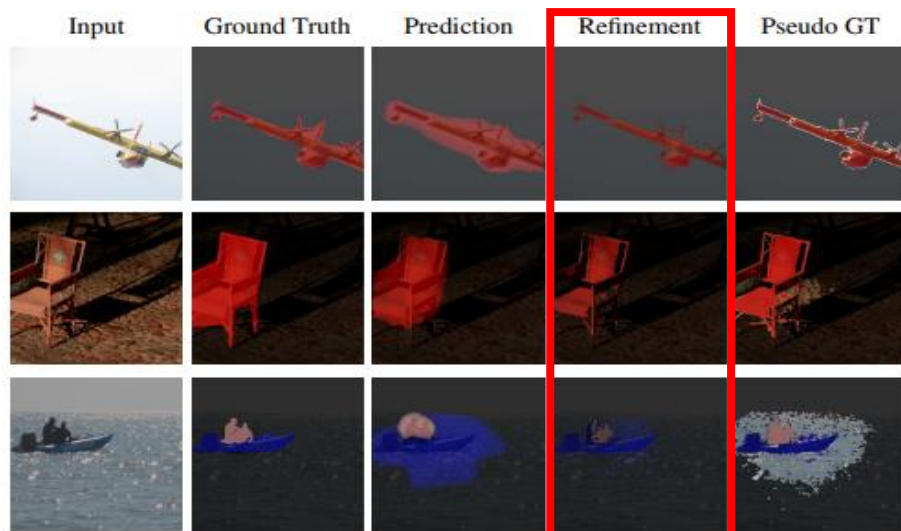
$$m_{:,i,j}^t = \sum_{(l,n) \in \mathcal{N}(i,j)} \alpha_{i,j,l,n} \cdot m_{:,l,n}^{t-1}$$

$$\alpha_{i,j,k,n} = \frac{\exp[\bar{k}(i,j,k,n)]}{\sum_{(q,r) \in \mathcal{N}(i,j)} \exp[k(i,j,q,r)]}$$

$$k(i, j, k, n) = -|I_{i,j} - I_{k,n}| / \sigma_{i,j}^2$$

$\bar{k}(\cdot)$ is the average affinity value $k(\cdot)$ across the RGB channels

$\sigma_{i,j}$ denotes the standard deviation of the image intensity computed locally for the affinity kernel.

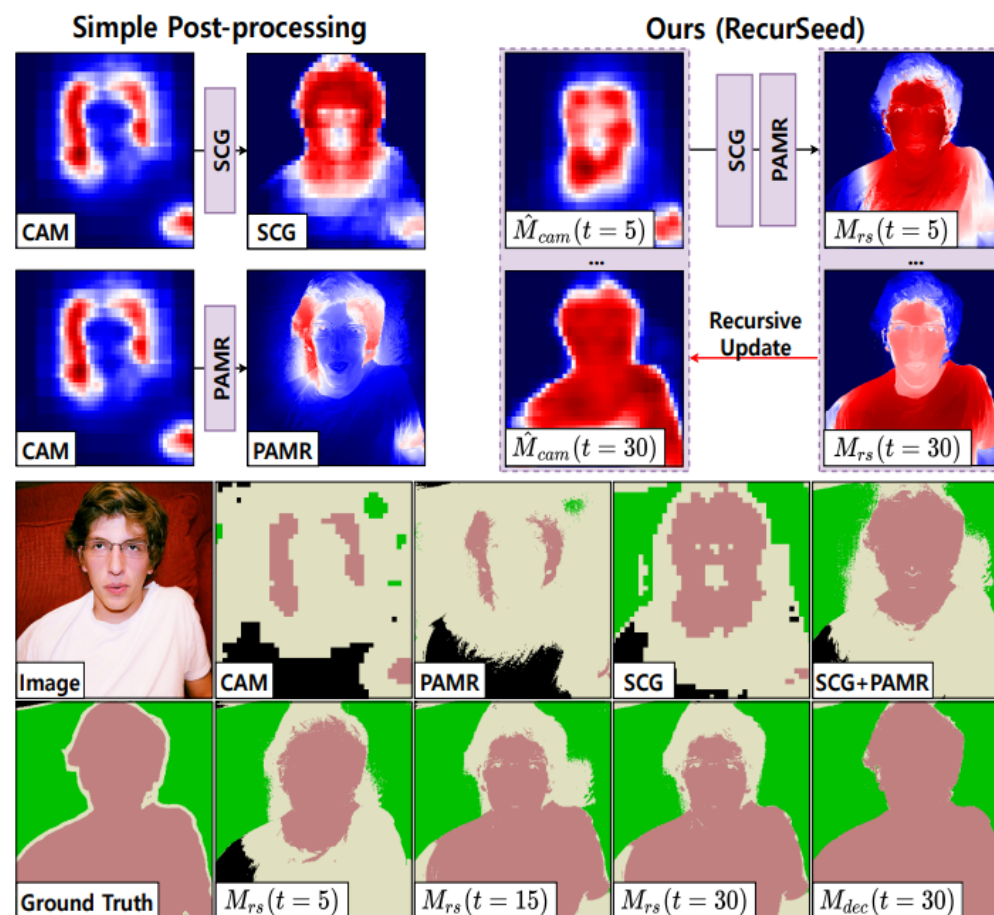


Algorithm Flow

- The performance of SCG depends highly on feature maps owing to its ability to update the CAM from SCG and PAMR recursively
- The proposed RS **gradually updates the initial CAM** to remedy the shortcoming of the SCG and PAMR

$$M_{rs}(t) = PAMR(SCG(M_{cam}(t)); \mathcal{W})$$

$$M_{rs}(t) \approx M_{cam}(t+1) \text{ for every epoch } t \in \{0 : T\}$$



Loss function

To achieve the objective of $M_{rs}(t) \approx M_{cam}(t+1)$

$$\begin{aligned} \theta_t = \theta_{t-1} - \eta \frac{\partial}{\partial \theta} \mathbb{E}_{\mathbf{I}} [\mathcal{L}_{cls}(\hat{Y}_{cls}(\tau), \mathbf{Y}; \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1} \\ - \eta \frac{\partial}{\partial \theta} \mathbb{E}_{\mathbf{I}} [\mathcal{L}_{seg}(\hat{M}_{dec}(\tau), M_{seg}(\tau); \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1} \\ - \eta \frac{\partial}{\partial \theta} \mathbb{E}_{\mathbf{I}} [\mathcal{L}_{rec}(\hat{M}_{cam}(\tau), M_{rs}(\tau); \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1} \end{aligned}$$

where $\mathbf{I} = \{I^1, I^2, \dots, I^B\}$

$$\hat{Y}_{cls}(t) = \sigma(GAP(\hat{M}_{cam}(t)))$$

Classification
(multi-label soft margin loss)

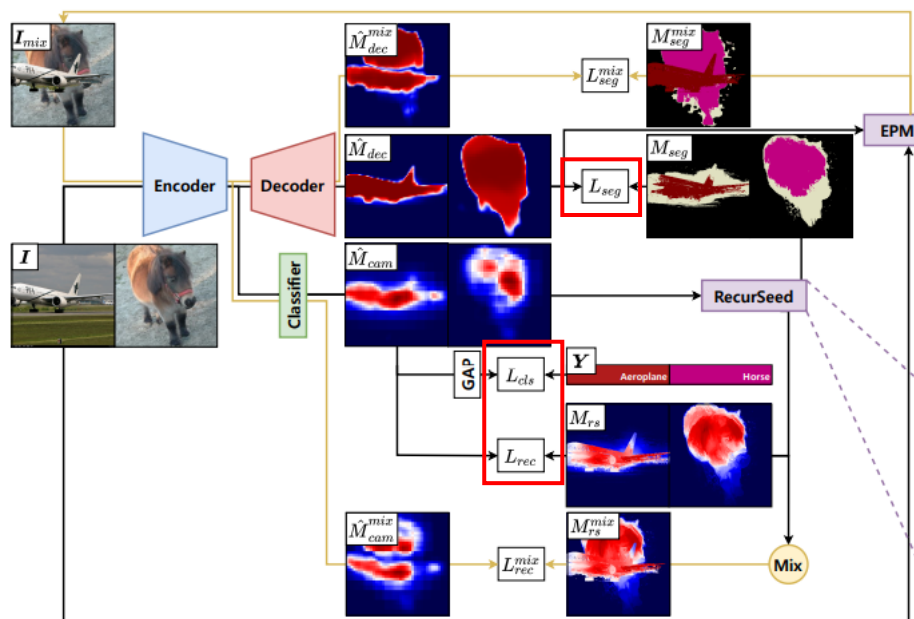
$$\hat{M}_{dec}(t) = D_{\theta_t^{dec}}(f_{enc}(t))$$

Decoder
(cross-entropy loss)

$$f_{enc}(t) = E_{\theta_t^{enc}}(\mathbf{I}),$$

$$\hat{M}_{cam}(t) = A_{\theta_t^{cls}}(f_{enc}(t))$$

Encoder
(L1 loss)

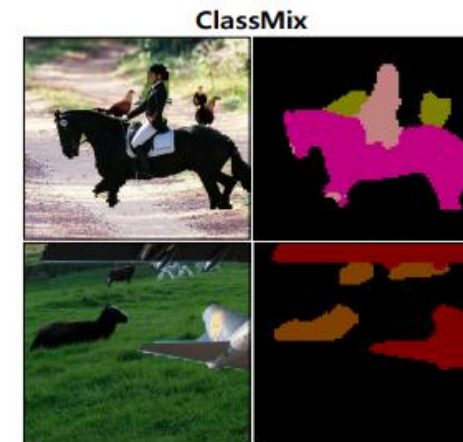
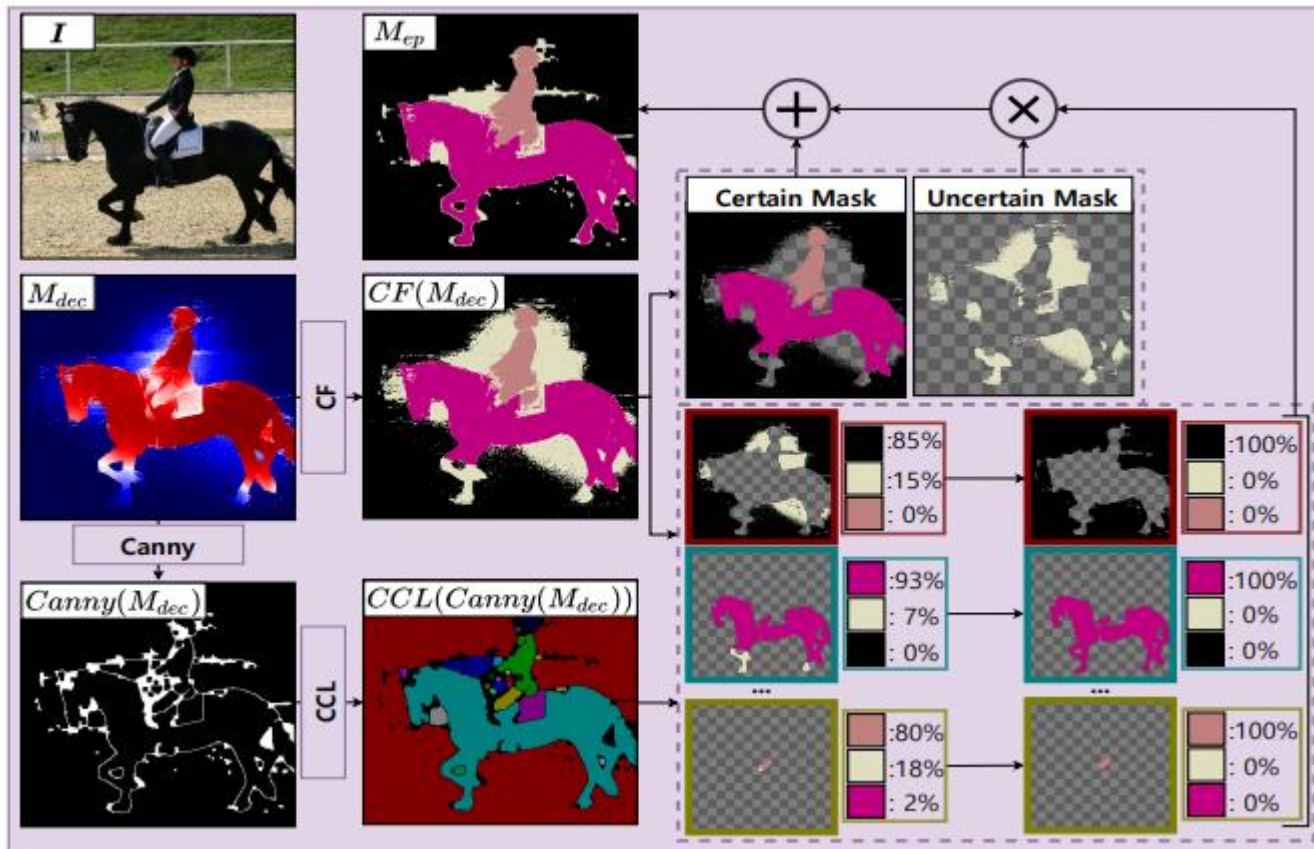


$$M_{seg}(t) = CF(M_{rs}(t))$$

$$= \begin{cases} \underset{c \in \mathcal{C}}{\operatorname{argmax}} (M_{rs}^c(t)[:, i, j]) & \text{if } \max_{c \in \mathcal{C}} (M_{rs}^c(t)[:, i, j]) > \delta_{fg}, \\ 0 & \text{if } \max_{c \in \mathcal{C}} (M_{rs}^c(t)[:, i, j]) < \delta_{bg}, \\ 255 & \text{otherwise.} \end{cases}$$

Overview

- **Mixes two images and pseudo masks refined by EP**, which disentangles foreground and background regions by using edge information in the per-pixel class probability domain.
- Leading to sample diversification and significantly improving performance for WSSS



Overview

- 1. perform the mask refinement by using absolute and relative per-pixel probability values
- 인접한 per-pixel probability values의 상대적 차이를 이용한다.
- Deriving edge from the mask, obtaining super pixels from the edge
- Singling out the most dominant class within each super pixel → boundary aware mask
- Extract edge → **canny**, Extract super pixel → Connected-component labeling (**CCL**)

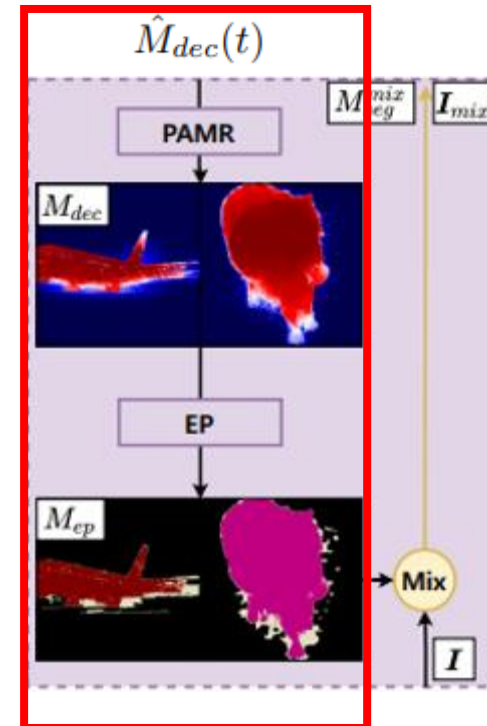
$$\hat{M}_{dec}(t)$$

$$M_{dec}^i(t) = PAMR(\boxed{D_{\theta_t^{dec}}(E_{\theta_t^{enc}}(I_i)); \mathcal{W}})$$

$$M_{ep}^i(t) = EP(M_{dec}^i(t))$$

$$\begin{aligned} \mathcal{R}_c^i(t) &= \{(k, n) \mid M_{ep}^i(t)[c, k, n] > \delta_{fg}\} \\ \mathcal{M}_{fg}^i &= \mathbb{1}[\cup_{c \in \mathcal{C}} \mathcal{R}_c^i(t)], \end{aligned}$$

Extract the union of all EP-refined foregrounds for image



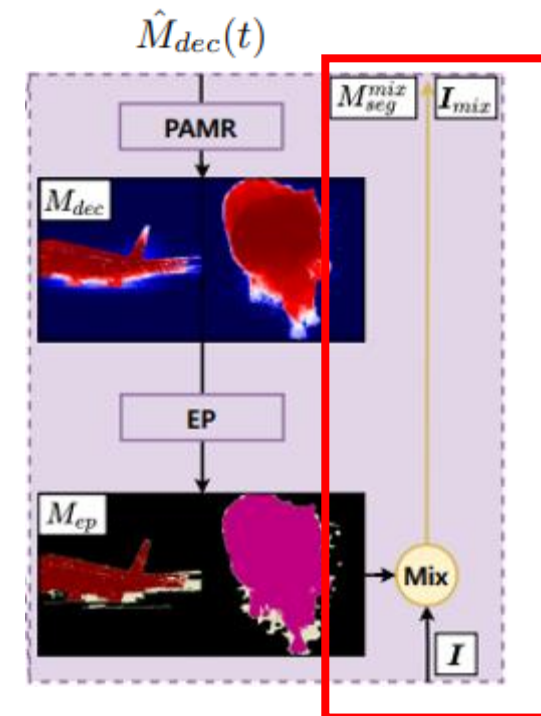
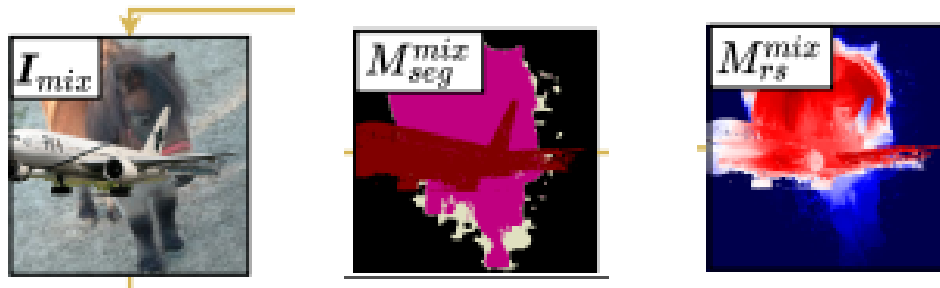
Overview

- 2. **blend** two EP-refined masks (and their corresponding original images)
-for two arbitrary indices i, j

$$I_{i \rightarrow j} = I_i \odot \mathcal{M}_{fg}^i + I_j \odot (1 - \mathcal{M}_{fg}^i)$$

$$M_{i \rightarrow j}^{seg}(t) = M_{ep}^i(t) \odot \mathcal{M}_{fg}^i + M_{ep}^j(t) \odot (1 - \mathcal{M}_{fg}^i)$$

$$M_{i \rightarrow j}^{rs}(t) = M_{rs}^i(t) \odot \mathcal{M}_{fg}^i + M_{rs}^j(t) \odot (1 - \mathcal{M}_{fg}^i)$$



Loss function

- make the network train the mixed images and labels

$$\begin{aligned} \theta_t = & \theta_{t-1} - (*) \\ & - \eta \frac{\partial}{\partial \theta} \mathbb{E}_I [\mathcal{L}_{seg}^{\text{mix}}(\hat{M}_{dec}^{\text{mix}}(\tau), M_{seg}^{\text{mix}}(\tau); \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1}, \\ & - \eta \frac{\partial}{\partial \theta} \mathbb{E}_I [\mathcal{L}_{rec}^{\text{mix}}(\hat{M}_{cam}^{\text{mix}}(\tau), M_{rs}^{\text{mix}}(\tau); \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1}, \end{aligned}$$

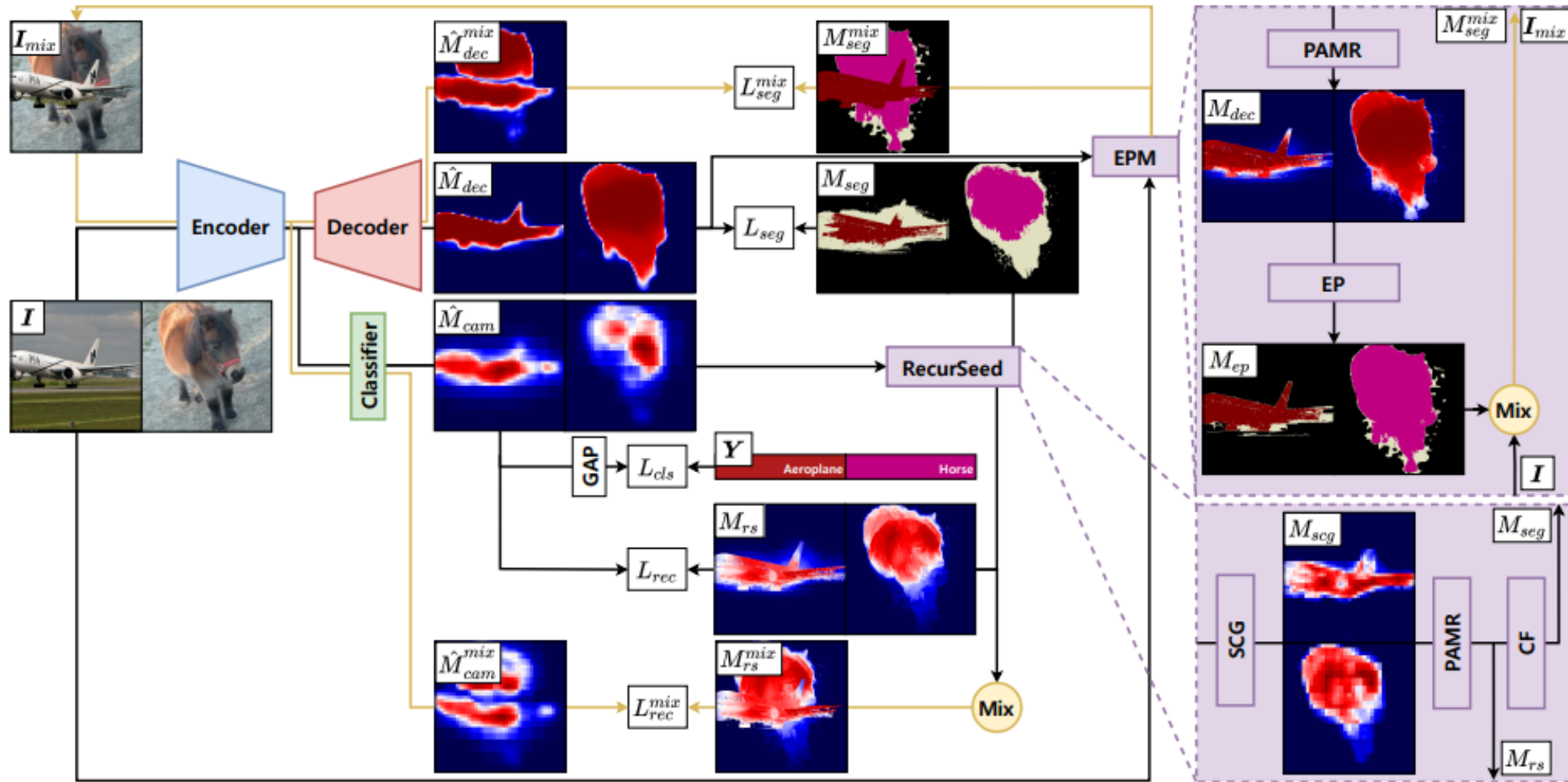
= Cross Entropy Loss

$$\begin{aligned} & \eta \frac{\partial}{\partial \theta} \mathbb{E}_I [\mathcal{L}_{cls}(\hat{Y}_{cls}(\tau), \mathbf{Y}; \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1} \\ & \eta \frac{\partial}{\partial \theta} \mathbb{E}_I [\mathcal{L}_{seg}(\hat{M}_{dec}(\tau), M_{seg}(\tau); \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1} \\ & \eta \frac{\partial}{\partial \theta} \mathbb{E}_I [\mathcal{L}_{rec}(\hat{M}_{cam}(\tau), M_{rs}(\tau); \theta)] \Big|_{\theta=\theta_{\tau}, \tau=t-1}, \end{aligned}$$

= RecurSeed Loss function

$$\begin{aligned} \mathbf{I}_{mix} &= \{I_{i \rightarrow j} \mid j \sim \text{unif}(\{1 : B\}/i) \text{ for } i \in \{1 : B\}\} \\ M_{seg}^{\text{mix}}(t) &= \{M_{i \rightarrow j}^{\text{seg}}(t) \mid I_{i \rightarrow j} \in \mathbf{I}_{mix}\} \\ \hat{M}_{dec}^{\text{mix}}(t) &= D_{\theta_t^{dec}}(E_{\theta_t^{cam}}(\mathbf{I}_{mix})) \\ M_{rs}^{\text{mix}}(t) &= \{M_{i \rightarrow j}^{rs}(t) \mid I_{i \rightarrow j} \in \mathbf{I}_{mix}\} \\ \hat{M}_{cam}^{\text{mix}}(t) &= A_{\theta_t^{cls}}(E_{\theta_t^{cam}}(\mathbf{I}_{mix})) \end{aligned}$$

Overview



Result

Method	Backbone	Sup.	VOC		COCO
			val	test	val
Single stage:					
EM ICCV'15 [17]	VGG16	\mathcal{I}	38.2	39.6	-
RRM AAAI'20 [19]	WR38	\mathcal{I}	62.6	62.9	-
SSSS CVPR'20 [20]	WR38	\mathcal{I}	62.7	64.3	-
AFA CVPR'22 [43]	MiT-B1	\mathcal{I}	66.0	66.3	38.9
Ours (single-stage, RS)	R50	\mathcal{I}	66.5	67.9	40.0
Ours (single-stage, RS+EPM)	R50	\mathcal{I}	69.5	70.6	42.2
Multiple stages:					
DSRG CVPR'18 [22]	R101	$\mathcal{I}+S$	61.4	63.2	26.0*
FickleNet CVPR'19 [6]	R101	$\mathcal{I}+S$	64.9	65.3	-
NSRM CVPR'21 [51]	R101	$\mathcal{I}+S$	68.3	68.5	-
AuxSegNet ICCV'21 [44]	WR38	$\mathcal{I}+S$	69.0	68.6	33.9
EDAM CVPR'21 [33]	R101	$\mathcal{I}+S$	70.9	70.6	-
DRS AAAI'21 [52]	R101	$\mathcal{I}+S$	71.2	71.4	-
CLIMS CVPR'22 [34]	R50	$\mathcal{I}+\mathcal{D}$	69.3	68.7	-
W-OoD CVPR'22 [45]	R101	$\mathcal{I}+\mathcal{D}$	70.7	70.1	-
EPS CVPR'21 [53]	R101	$\mathcal{I}+S$	70.9	70.8	35.7*
L2G CVPR'22 [14]	R101	$\mathcal{I}+S$	72.1	71.7	44.2
RCA CVPR'22 [12]	R101	$\mathcal{I}+S$	72.2	72.8	36.8*
PPC CVPR'22 [35]	R101	$\mathcal{I}+S$	72.6	73.6	-
PSA CVPR'18 [5]	WR38	\mathcal{I}	61.7	63.7	-
IRNet CVPR'19 [31]	R50	\mathcal{I}	63.5	64.8	-
SEAM CVPR'20 [29]	WR38	\mathcal{I}	64.5	65.7	31.9
AdvCAM CVPR'21 [7]	R101	\mathcal{I}	68.1	68.0	-
CSE ICCV'21 [54]	WR38	\mathcal{I}	68.4	68.2	36.4
CPN ICCV'21 [30]	WR38	\mathcal{I}	67.8	68.5	-
RIB NIPS'21 [32]	R101	\mathcal{I}	68.3	68.6	43.8
ReCAM CVPR'22 [55]	R101	\mathcal{I}	68.5	68.4	-
ADEHE CVPR'22 [41]	R101	\mathcal{I}	68.6	68.9	-
AMR AAAI'22 [56]	R101	\mathcal{I}	68.8	69.1	-
URN AAAI'22 [57]	R101	\mathcal{I}	69.5	69.7	40.7
SIPE CVPR'22 [13]	R101	\mathcal{I}	68.8	69.7	40.6
AMN CVPR'22 [36]	R101	\mathcal{I}	69.5	69.6	44.7
MCTformer CVPR'22 [39]	WR38	\mathcal{I}	71.9	71.6	42.0
SANCE CVPR'22 [42]	R101	\mathcal{I}	70.9	72.2	44.7†
Ours (multi-stage, RS)	R101	\mathcal{I}	72.8	72.8	45.8
Ours (multi-stage, RS+EPM)	R101	\mathcal{I}	74.4	73.6	46.4

RS	SCG	PAMR	mIoU	FP	FN
✓			58.0	0.268	0.165
✓	✓		59.3	0.225 (↓ 0.043)	0.194 (↑ 0.029)
✓	✓		65.2	0.216	0.143
✓	✓	✓	65.9	0.210 (↓ 0.006)	0.141 (↓ 0.002)
✓	✓		*67.4	*0.196	*0.141
✓	✓	✓	*70.7	*0.171 (↓ 0.025)	*0.134 (↓ 0.007)

* denotes the decoder map result.

Method	Backbone	F_1	mIoU
RecurSeed	R50	94.7	70.7
RecurSeed + *CutMix [25]	R50	95.6	68.5
RecurSeed + *SaliencyGrafting [26]	R50	96.8	68.6
RecurSeed + *CDA [27]	R50	96.0	69.0
RecurSeed + *ClassMix [28]	R50	94.6	71.2
RecurSeed + EdgePredictMix	R50	95.2	75.2

Result

Method	Backbone	CAM (%)	CAM +RW (%)	CAM+RW +dCRF (%)
AffinityNet [4]	ResNet-50	47.82	58.10	59.70
Puzzle-CAM	ResNet-50	51.53	64.16	64.70
Puzzle-CAM	ResNeSt-50	57.59	69.48	69.91
Puzzle-CAM	ResNeSt-101	61.85	71.92	72.46
Puzzle-CAM	ResNeSt-269	62.45	74.14	74.67

Puzzle-CAM with MLF

Method	Backbone	Seed	RW
SEAM [29]	WR38	55.4	63.6
IRNet [31]	R50	48.8	66.3
CSE [54]	WR38	56.0	66.9
CDA [27]	R50	50.8	67.7
CPN [30]	WR38	57.4	67.8
CONTA [58]	R50	48.8	67.9
AMR [56]	R50	56.8	69.7
AdvCAM [7]	R50	55.6	69.9
PPC [35]	WR38	61.5	70.1
RIB [32]	R50	56.5	70.6
Ours (RS)	R50	70.7	74.8
Ours (RS+EPM)	R50	75.2	76.7

RS & EPM with MLF

- WSSS에서 PASCAL VOC 2012, MS COCO 2014 benchmark SOTA 달성
- RS와 EPM은 encoder와 decoder을 포함한 SLF에 범용적으로 적용할 수 있으며, 더 높은 성능을 위해 Backbone을 Upgrade할 수 있다.

Reference

- **Learning Deep Features for Discriminative Localization**

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba Computer Science and Artificial Intelligence Laboratory, MIT

- **PUZZLE-CAM: IMPROVED LOCALIZATION VIA MATCHING PARTIAL AND FULL FEATURES**

Sanghyun Jo GYNetworks In-Jae Yu* KAIST School of Computing

- **RecurSeed and EdgePredictMix: Single-stage Learning is Sufficient for Weakly-Supervised Semantic Segmentation**

Sanghyun Jo^{1*}, In-Jae Yu^{2*}, Kyungsu Kim^{3,4} †

- **Unveiling the Potential of Structure Preserving for Weakly Supervised Object Localization**

Xingjia Pan^{1,3*} Yingguo Gao^{1*} Zhiwen Lin^{1*} Fan Tang² † Weiming Dong^{3,4,5} Haolei Yuan¹ Feiyue Huang¹ Changsheng Xu^{3,4,5}

- **Single-Stage Semantic Segmentation from Image Labels**

Nikita Araslanov Stefan Roth

Q & A

감사합니다