

BEIJING INSTITUTE OF TECHNOLOGY

Few Shot Segmentation

方致远 张泽康 2021/5/14

Motivation



Defects of Deep-learning-based Segmentation

- Eager for large amount of samples of target
- Cannot segment unseen objects





The aims of Few-shot Segmentation

Segment unseen object with only a few references after training



Implementation



■ Divide classes *C* of dataset into seen C_{seen} and unseen C_{unseen} , $C_{seen} \cap C_{unseen} = \emptyset$

$$D_{train} = \{(S_i, Q_i)\}_{i=1}^{N_{train}}, \quad D_{test} = \{(S_i, Q_i)\}_{i=1}^{N_{test}}$$

- $S_i = \{ (I_{c,k}, M_{c,k}) \mid c \in c_{seen} \text{ for training } else \ c \in c_{unseen} \}, \ Q_i = \{ (I_{c,n}, M_{c,n}) \}$
- (S_i, Q_i) called episode is the input to the model during training/testing stage



Implementation









Dataset:

Pascal VOC 2012

MS COCO

Metric:

Mean Intersection over Union (mIoU)

 $\blacksquare mIoU = \frac{TP}{FP + FN + TP}$







- SG-One: Similarity Guidance Network for One-Shot Semantic Segmentation (2018)
- PANet: Few-Shot Image Semantic Segmentation with Prototype Alignment (2019)
- Prior Guided Feature Enrichment Network for Few-Shot Segmentation (2020)
- Adaptive Prototype Learning and Allocation for Few-Shot Segmentation (2021)
- Edge-Labeling Graph Neural Network for Few-shot Learning (2019)



SG-One: Similarity Guidance Network for One-Shot Semantic Segmentation

德以明理 学以特工





Key points:

- Masked Average Pooling strategy
- Build relationship between support and query with cosine similarity
- Unified network structure



SG-One





德以明理 学以特乙

SG-One



K-shot Testing:

- Ensemble the K segmentation mask with equation:
 - $\hat{Y}_{x,y} = max(\hat{Y}^1_{x,y}, \hat{Y}^2_{x,y}, \cdots, \hat{Y}^K_{x,y})$
- Average the K representative vectors

MEAN IOU RESULTS OF ONE-SHOT SEGMENTATION ON THE PASCAL-51 DATASET. THE BEST RESULTS ARE IN BOLD

Methods	PASCAL-5 ⁰	PASCAL-5 ¹	PASCAL-5 ²	PASCAL-5 ³	Mean
1-NN	25.3	44.9	41.7	18.4	32.6
LogReg	26.9	42.9	37.1	18.4	31.4
Siamese	28.1	39.9	31.8	25.8	31.4
OSVOS [37]	24.9	38.8	36.5	30.1	32.6
OSLSM [15]	33.6	55.3	40.9	33.5	40.8
co-FCN [16]	36.7	50.6	44.9	32.4	41.1
SG-One(Ours)	40.2	58.4	48.4	38.4	46.3



PANet: Few-Shot Image Semantic Segmentation with Prototype Alignment







Key points:

- Propose a non-parametric metric learning to extract knowledge and perform segmentation
- Design a Prototype Alignment structure to fully exploit knowledge from the support



/ PANet









Method	1-shot					5-shot					Δ	#Params
	split-1	split-2	split-3	split-4	Mean	split-1	split-2	split-3	split-4	Mean	Mean	
OSLSM [21]	33.6	55.3	40.9	33.5	40.8	35.9	58.1	42.7	39.1	43.9	3.1	272.6M
co-FCN [16]†	36.7	50.6	44.9	32.4	41.1	37.5	50.0	44.1	33.9	41.4	0.3	34.2M
SG-One [28]	40.2	58.4	48.4	38.4	46.3	41.9	58.6	48.6	39.4	47.1	0.8	19.0M
PANet-init	30.8	40.7	38.3	31.4	35.3	41.6	52.7	51.6	40.8	46.7	11.4	14.7M
PANet	42.3	58.0	51.1	41.2	48.1	51.8	64.6	59.8	46.5	55.7	7.6	14.7M



Prior Guided Feature Enrichment Network for Few-Shot Segmentation

德以明理 学以特工





PFENet mainly wants to solve following two problems:

- High-level features using in a few-shot model may cause performance drop
 - Using fixed high-level features to yield the prior mask
- Spatial Inconsistency
 - Proposing a FPN-like Feature Enrichment Module with residual blocks





德以明理 学以特工



德以明理 学以特工





Methods	1-Shot						5-Shot					Parame
Tvic mous	Fold-0	Fold-1	Fold-2	Fold-3	Mean	Fold-0	Fo	ld-1	Fold-2	Fold-3	Mean	1 aranns
VGG-16 Backbone												
OSLSM ₂₀₁₇ [33]	33.6	55.3	40.9	33.5	40.8	35.9	5	8.1	42.7	39.1	44.0	276.7M
co-FCN ₂₀₁₈ [29]	36.7	50.6	44.9	32.4	41.1	37.5	5	0.0	44.1	33.9	41.4	34.2M
SG-One ₂₀₁₈ [58]	40.2	58.4	48.4	38.4	46.3	41.9	5	8.6	48.6	39.4	47.1	19.0M
AMP ₂₀₁₉ [35]	41.9	50.2	46.7	34.7	43.4	41.8	5	5.5	50.3	39.9	46.9	34.7M
PANet ₂₀₁₉ [45]	42.3	58.0	51.1	41.2	48.1	51.8	6	4.6	59.8	46.5	55.7	14.7M
FWBF ₂₀₁₉ [28]	47.0	59.6	52.6	48.3	51.9	50.9	6	2.9	56.5	50.1	55.1	-
Ours	56.9	68.2	54.4	52.4	58.0	59.0	6	9.1	54.8	52.9	59.0	10.4M
				ResNe	et-50 Bacl	kbone						
CANet ₂₀₁₉ [54]	52.5	65.9	51.3	51.9	55.4	55.5	6	7.8	51.9	53.2	57.1	19.0M
PGNet ₂₀₁₉ [53]	56.0	66.9	50.6	50.4	56.0	54.9	6	7.4	51.8	53.0	56.8	17.2M
Ours	61.7	69.5	55.4	56.3	60.8	63.1	7	0.7	55.8	57.9	61.9	10.8M
				ResNe	t-101 Bac	kbone						
FWBF ₂₀₁₉ [28]	51.3	64.5	56.7	52.2	56.2	54.8	6	7.4	62.2	55.3	59.9	-
Ours	60.5	69.4	54.4	55.9	60.1	62.8	7	0.4	54.9	57.6	61.4	10.8M
Methods	Backhone			1-Sho	ot					5-Sho	ot	
	Duckeone	Fold-0	Fold-	1 Fold-	2 Fold	l-3 Me	ean	Fold-0	Fold-1	Fold-	2 Fold	-3 Mean
			Cl	ass mIoU	Evaluati	on						
FWBF2019 [28]	VGG-16	18.4	16.7	19.6	25.	4 20).0	20.9	19.2	21.9	28.4	4 22.6
Ours	VGG-16	33.4	36.0	34.1	32.	8 34	l.1	35.9	40.7	38.1	36.3	1 37.7
PANet ₂₀₁₉ † [45]	VGG-16	-	-	-	-	20).9	-	-	-	-	29.7
Ours†	VGG-16	35.4	38.1	36.8	34.	7 36	5.3	38.2	42.5	41.8	38.9	9 40.4
FWBF2019 [28]	ResNet-101	19.9	18.0	21.0	28.	9 21	1.2	19.1	21.5	23.9	30.1	1 23.7
Ours	ResNet-101	34.3	33.0	32.3	30.	1 32	2.4	38.5	38.6	38.2	34.3	3 37.4
Ours†	ResNet-101	36.8	41.8	38.7	36.	7 38	3.5	40.4	46.8	43.2	40.5	5 42.7

德以明理 学以特之















(a) SGC is adaptive to object scale variation



(b) GPA is adaptive to object shape variation





Figure 2. Overall architecture of the proposed Adaptive Superpixel-guided Network.



Super-pixel

Serve super-pixel centroids as prototypes

- concatenate the coordinates of each pixel with the support feature map
- Get foreground feature map with mask
- Iteration

$$D = \sqrt{(d_f)^2 + (d_s/r)^2},$$
 (1)

where d_f , d_s are the Euclidean distance for features and coordinate values, and r is a weighting factor. We filter





Prototype Allocation



Figure 4. Illustration of proposed Guided Prototype Allocation.



Adaptability

Compute the number of prototypes with the number of pixels in the foreground

$$N_{sp} = \min\left(\left\lfloor \frac{N_m}{S_{sp}} \right\rfloor, N_{max}\right),$$

K-shot setting

Directly get all prototypes from K-shot

$$N_{sp} = max \sum_{i=1}^{k} N_{sp}^{i}$$



Backbone	Methods	1-shot										
		s-0	s-1	s-2	s-3	mean	s-0	s-1	s-2	s-3	mean	4
	OSLSM [23]	33.60	55.30	40.90	33.50	40.80	35.90	58.10	42.70	39.10	43.95	3.15
	co-FCN [21]	36.70	50.60	44.90	32.40	41.10	37.50	50.00	44.10	33.90	41.40	0.30
VCC 16	AMP [24]	41.90	50.20	46.70	34.40	43.40	40.30	55.30	49.90	40.10	46.40	3.00
VUU-10	SG-One [42]	40.20	58.40	48.40	38.40	46.30	41.90	58.60	48.60	39.40	47.10	0.80
	PANet [32]	42.30	58.00	51.10	41.20	48.10	51.80	64.60	59.80	46.50	55.70	7.60
	FWB [20]	47.04	59.64	52.61	48.27	51.90	50.87	62.86	56.48	50.09	55.08	3.18
	CANet [†] [40]	52.50	65.90	51.30	51.90	55.40	55.50	67.80	51.90	53.20	57.10	1.70
	PGNet [†] [39]	56.00	66.90	50.60	50.40	56.00	57.70	68.70	52.90	54.60	58.50	2.50
	RPMMs [34]	55.15	66.91	52.61	50.68	56.34	56.28	67.34	54.52	51.00	57.30	0.96
ResNet50	SimPropNet [8]	54.82	67.33	54.52	52.02	57.19	57.20	68.50	58.40	56.05	60.04	2.85
	PPNet [18]	47.83	58.75	53.80	45.63	51.50	58.39	67.83	64.88	56.73	61.96	10.46
	PFENet [29]	61.70	69.50	55.40	56.30	60.80	63.10	70.70	55.80	57.90	61.90	1.10
	ASGNet (ours)	58.84	67.86	56.79	53.66	59.29	63.66	70.55	64.17	57.38	63.94	4.65
	FWB [20]	51.30	64.49	56.71	52.24	56.19	54.84	67.38	62.16	55.30	59.92	3.73
ResNet101	DAN† [31]	54.70	68.60	57.80	51.60	58.20	57.90	69.00	60.10	54.90	60.50	2.30
	PFENet [29]	60.50	69.40	54.40	55.90	60.10	62.80	70.40	54.90	57.60	61.40	1.30
	ASGNet (ours)	59.84	67.43	55.59	54.39	59.31	64.55	71.32	64.24	57.33	64.36	5.05



德以明理 学以特工



Graph Neural Network (GNN)

 $\mathbf{H} = F(\mathbf{H}, \mathbf{X})$ $\mathbf{O} = G(\mathbf{H}, \mathbf{X}_N)$

GCN

$$\mathbf{X}_{s}^{k+1} = \sigma \left(\left(\mathbf{I}_{|\mathcal{V}^{s}|} + \mathbf{D}_{s}^{-\frac{1}{2}} \mathbf{A}_{s} \mathbf{D}_{s}^{-\frac{1}{2}} \right) \mathbf{X}_{s}^{k} \mathbf{W}_{s}^{k} \right)$$

[1]Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. ICLR, 2017.







Not suitable for few-shot learning





(a) Graph Construction

(c) Query Node Label Prediction

(b-2) EGNN: Edge Feature Update

德以明理 学以特之



Initial edge label

$$\mathbf{e}_{ij}^{0} = \begin{cases} [1||0], & \text{if } y_{ij} = 1 \text{ and } i, j \leq N \times K, \\ [0||1], & \text{if } y_{ij} = 0 \text{ and } i, j \leq N \times K, \\ [0.5||0.5], & \text{otherwise}, \end{cases}$$

Node update

$$\mathbf{v}_{i}^{\ell} = f_{v}^{\ell}(\left[\sum_{j} \tilde{e}_{ij1}^{\ell-1} \mathbf{v}_{j}^{\ell-1} || \sum_{j} \tilde{e}_{ij2}^{\ell-1} \mathbf{v}_{j}^{\ell-1}\right]; \theta_{v}^{\ell}),$$

$$\bar{e}_{ij1}^{\ell} = \frac{f_e^{\ell}(\mathbf{v}_i^{\ell}, \mathbf{v}_j^{\ell}; \theta_e^{\ell}) e_{ij1}^{\ell-1}}{\sum_k f_e^{\ell}(\mathbf{v}_i^{\ell}, \mathbf{v}_k^{\ell}; \theta_e^{\ell}) e_{ik1}^{\ell-1} / (\sum_k e_{ik1}^{\ell-1})}, \quad (4)$$

$$\bar{e}_{ij2}^{\ell} = \frac{(1 - f_e^{\ell}(\mathbf{v}_i^{\ell}, \mathbf{v}_j^{\ell}; \theta_e^{\ell})) e_{ij2}^{\ell-1}}{\sum_k (1 - f_e^{\ell}(\mathbf{v}_i^{\ell}, \mathbf{v}_k^{\ell}; \theta_e^{\ell})) e_{ik2}^{\ell-1} / (\sum_k e_{ik2}^{\ell-1})}, \quad (5)$$

$$\mathbf{e}_{ij}^{\ell} = \bar{\mathbf{e}}_{ij}^{\ell} / \|\bar{\mathbf{e}}_{ij}^{\ell}\|_{1}, \quad (6)$$

Edge update

德以明理 学以特之



Algorithm 1: The process of EGNN for inference 1 Input: $\mathcal{G} = (\mathcal{V}, \mathcal{E}; \mathcal{T})$, where $\mathcal{T} = \mathcal{S} \bigcup \mathcal{Q}$, $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{N \times K}, \mathcal{Q} = \{\mathbf{x}_i\}_{i=N \times K+1}^{N \times K+T}$ 2 Parameters: $\theta_{emb} \cup \{\theta_v^{\ell}, \theta_e^{\ell}\}_{\ell=1}^L$ 3 Output: $\{\hat{y}_i\}_{i=N \times K+1}^{N \times K+T}$ 4 Initialize: $\mathbf{v}_i^0 = f_{emb}(\mathbf{x}_i; \theta_{emb}), \mathbf{e}_{ij}^0, \forall i, j$ 5 for $\ell = 1, \dots, L$ do

12 end

```
/* Query node label prediction */
13 \{\hat{y}_i\}_{i=N \times K+1}^{N \times K+T} \leftarrow \texttt{Edge2NodePred}(\{y_i\}_{i=1}^{N \times K}, \{\mathbf{e}_{ij}^L\})
```

德以明理 学以特之



Figure 3: Detailed network architectures used in EGNN. (a) Embedding network f_{emb} . (b) Feature (node) transformation network f_v^{ℓ} . (c) Metric network f_e^{ℓ} .



(a) *mini*ImageNet

Model	Trans.	5-Way 5-Shot	(b) <i>tiered</i> ImageNet		
Matching Networks [2]	No	55.30			
Reptile [46]	No	62.74	Model	Trans.	5-Way 5-Shot
Prototypical Net [3]	No	65.77	Reptile [46]	No	66.47
GNN [6]	No	66.41	Prototypical Net [3]	No	69.57
EGNN	No	66.85	EGNN	No	70.98
MAML [4]	BN	63.11	MAML [4]	BN	70.30
Reptile $+$ BN [46]	BN	65.99	Reptile + BN [46]	BN	71.03
Relation Net [5]	BN	67.07	Relation Net [5]	BN	71.31
MAML+Transduction [4]	Yes	66.19	MAML+Transduction [4]	Yes	70.83
TPN [12]	Yes	69.43	TPN [12]	Yes	72.58
TPN (Higher K) [12]	Yes	69.86	EGNN+Transduction	Yes	80.15
EGNN+Transduction	Yes	76.37			



Transductive inference

Classify the entire test set at once to alleviate the low-data problem



[2]Liu Y, Lee J, Park M, et al. Learning to propagate labels: Transductive propagation network for few-shot learning[C]//7th International Conference on Learning Representations, ICLR 2019. 2019.

Edge-Labeling Graph Neural Network for Few-



- **11 2**- -----