

CycleSegNet: Object Co-Segmentation with Cycle Refinement and Region Correspondence

Chi Zhang*, Guankai Li*, Guosheng Lin, Qingyao Wu, Rui Yao

Abstract—Image co-segmentation is an active computer vision task that aims to segment the common objects from a set of images. Recently, researchers design various learning-based algorithms to undertake the co-segmentation task. The main difficulty in this task is how to effectively transfer information between images to make conditional predictions. In this paper, we present CycleSegNet, a novel framework for the co-segmentation task. Our network design has two key components: a region correspondence module which is the basic operation for exchanging information between local image regions, and a cycle refinement module, which utilizes ConvLSTMs to progressively update image representations and exchange information in a cycle and iterative manner. Extensive experiments demonstrate that our proposed method significantly outperforms the state-of-the-art methods on four popular benchmark datasets — PASCAL VOC dataset, MSRC dataset, Internet dataset, and iCoseg dataset, by 2.6%, 7.7%, 2.2%, and 2.9%, respectively.

Index Terms—deep learning, co-segmentation, cycle refinement, attention

1 INTRODUCTION

Image co-segmentation is an active computer vision topic with a long research history, which aims to segment the common objects jointly from a set of images. Image co-segmentation algorithms have shown their usages in various computer vision tasks, such as image retrieval [47], 3D reconstruction [37], photo collections [41], image matching [6], [58], [59], and video object tracking [32], [33], [41]. Recently, data-driven deep neural networks based methods attract wide interest in the literature. The powerful deep neural networks along with the challenging evaluation benchmarks built upon large-scale public datasets have brought this task to a new era and make it more challenging. Although deep neural networks have shown remarkable success in many other segmentation tasks, such as semantic segmentation [7], [13], [24], [32], [34], [60], [61], interactive segmentation [44], and instance segmentation [17], [55], [65], these models are not directly applicable to the co-segmentation tasks, as the outputs are conditioned on the pairwise or group-wise relations between input images.

Many existing CNN based methods [1], [5], [29], [33] solve the image co-segmentation problems by employing a pair of parameter-shared Siamese networks to generate feature representations of two images and using various methods to transfer information between them to make conditional predictions. The intuitions behind these methods

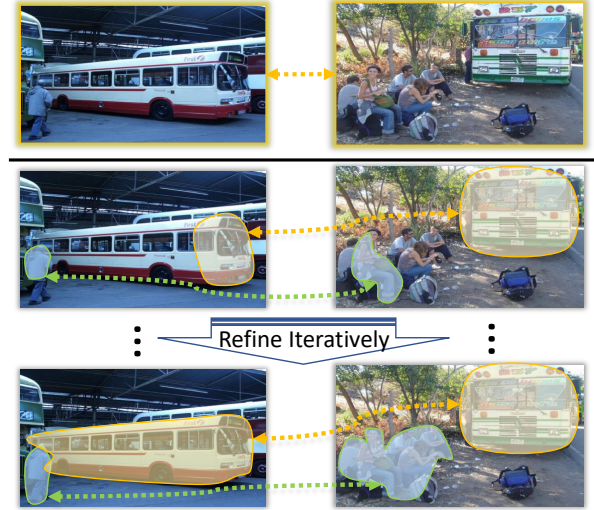


Fig. 1: An illustration of our motivation. Previous methods exchange image-level semantics as cues for image co-segmentation (up), while our method establishes correspondence between local image regions and refines the predictions iteratively (down).

are that in order to segment the common objects between images, the cues about what objects are existing in individual images must be exchanged for predictions. The challenge here is how to effectively transfer useful information based on image representations. Directly copying the image representations is unrealistic, as images have structured representations and no correspondence information is provided. Hence, simply multiplying or concatenating structured image representations would not work when no alignment information is provided. To tackle this problem, many recent works [1], [5] bypass the correspondence problem by squeezing the image structures and exchange information in the form of a global image-level representation which is usu-

* indicates equal contribution.

Corresponding author: Guosheng Lin.

- Chi Zhang, Guankai Li, Guosheng Lin are with School of Computer Science and Engineering, Nanyang Technological University (NTU), Singapore 639798 (email: guankai001@e.ntu.edu.sg, chi007@e.ntu.edu.sg, gslin@e.ntu.edu.sg).
- Qingyao Wu is with School of Software Engineering, South China University of Technology, and Pazhou Lab, Guangzhou, China (email: qyw@scut.edu.cn)
- R. Yao is with the School of Computer Science and Technology, China University of Mining and Technology, Xuzhou 221116, China (e-mail: ruiyao@cumt.edu.cn)

ally achieved by the global average pooling. As a result, the prediction of each pixel location can refer to the same global representations that are exchanged. However, squeezing the image structures into a global representation inevitably loses local discriminative information which can be useful to locate common objects from the scene. Moreover, as the image content often has a complex composition, including the cluttered background and objects from different classes, the image-level representation may also introduce noise into the cues for conditional predictions. In this paper, we propose to solve the aforementioned problems from two perspectives, which are illustrated in Figure 1.

First, we argue that a desired deep image co-segmentation algorithm should maintain the image structures and the exchanged information should contain local discriminative information as the cues for conditional predictions. To achieve this goal, we design a region correspondence module (RCM) that utilizes the attention mechanism to establish the correspondence between local regional representations in different images. In the RCM, each pixel location can attend to its most relevant regions in other images, and thus the prediction of each pixel refers more to the relevant regions based on the attention values. In comparison with the cues generated by global average poolings, the cues in our method include more local discriminative features and contain less noise for prediction.

Second, we argue that the task of image co-segmentation can be modeled as an iterative optimization process where the predictions as well as the exchanged cues can be refined progressively. We take inspiration from how we humans locate common objects from two images — given two images that contain the same objects, humans may observe two images alternatively and every time an image is watched, humans remember what exists in one image and try to find them in other images. Gradually, those exclusive image contents are out of consideration and the cues about the common objects are getting clearer. Based on such intuition, we design a cycle refinement module (CRM) to refine feature representations for predictions in a cycle and iterative manner. We employ a pair of parameter-shared convolutional LSTMs (ConvLSTMs) [45] to achieve this goal, which is operated at the bottleneck of an encoder-decoder network. At each time step, the LSTM cell takes the exchanged information from the opposite branch as the input and updates its own embedding, which is then sent to the opposite branch as the exchanged cues in the next step. In this way, the cycle refinement module can continually exchange information between two images and refine the image representations. As a result, both the image representations for predictions and cues about the common objects can be refined iteratively, and they can benefit from each other in this process. Specifically, a clear cue received from other images about the common objects can definitely improve the image representation that is used for conditional prediction, and at the same time, a good image representation that contains more information about the common objects and less about cluttered background produces better cues.

An advantage of our design is that our network for co-segmentation can handle the input with paired images as well as a group of images with the same model and

network parameters. This is because the key operation in our network design is based on the attention mechanism and does not require a fixed number of pixel locations on each side. Therefore, we can transform the network-like relations between grouped inputs to a one-versus-rest relation to undertake the group segmentation task with the same model.

Finally, we develop a model variant that exploits multi-level features to better infer the common objects in two images. As is often observed in the CNN visualization literature [32], [57], high-level features contain more class-related information while middle-level features correspond to object parts that may be shared across different categories. It is therefore useful to establish region connections at multiple semantic levels to improve predictions. We apply our proposed modules to the features at different stages independently and then fuse their results during the decoding process. We show that our network can better locate the pixels belonging to the common objects by utilizing multi-level features.

To validate the effectiveness of our design, we implement various experiments on multiple datasets to investigate each component in our network. Extensive experiments demonstrate that our design significantly outperforms existing models and the baselines. The main contributions of this work are summarized as follows:

- We propose a region correspondence module as a basic operation to exchange information between images for the co-segmentation task. The proposed module maintains local discriminative representations and utilizes the attention mechanism to directly establish correspondence between image regions.
- We propose a cycle refinement module that employs ConvLSTMs to progressively refine network predictions as well as the exchanged information.
- We demonstrate that our model can handle inputs of paired images and grouped images with the same network and parameters.
- We develop a model variant that exploits multi-level features to undertake the co-segmentation task. The multi-level version of our network can better locate common objects and improve the performance.
- Experiments on four popular benchmarks — Pascal VOC dataset, MSRC dataset, Internet dataset, and iCoseg dataset, demonstrate that our proposed network on both co-segmentation and group segmentation tasks significantly outperforms the baseline methods and achieves new state-of-the-art performance with remarkable advantages.

2 RELATED WORK

Object co-segmentation. Early works often address object co-segmentation problems by comparing the visual features of paired images, such as foreground color histogram [41], SIFT [36], and saliency [4]. Joulin *et al.* [25] show that finding discriminative features can be well adapted to solve the object co-segmentation tasks. Rubio *et al.* [43] train SVM classifiers based on a Gaussian Mixture Model to capture the correspondence between different regions of input images.

Rubinstein *et al.* [42] utilize dense correspondences and visual saliency to find out the sparsity and visual variability of the common objects in the whole image database.

In the past few years, with the assistance of deep learning, some CNN-based networks have been proposed to solve the co-segmentation task. For instance, Quan *et al.* [38] propose a manifold ranking algorithm by combining the features obtained from the VGG network and hand-crafted features for superpixels to implement object co-segmentation. Li *et al.* [29] utilize a pure fully convolutional neural network to solve the object co-segmentation task. Their method has a Siamese encoder-decoder architecture and utilizes a mutual correlation layer to segment the common objects in a pair of images. Chen *et al.* [5] propose a semantic attention layer to spotlight feature channels, which is the first work that leverages the channel attentions for object co-segmentation tasks. Zhang *et al.* [66] design a spatial modulator for group-wise mask learning and a semantic modulator for co-category classification. Li *et al.* [27] propose a co-attention recurrent unit to generate the group representation containing the synergetic information, which is broadcasted to each individual image to facilitate the inferring of common objects. Hsu *et al.* [19] explore an unsupervised co-segmentation setting where no mask supervision is provided, which is different from our task setting where the ground-truth masks are provided. To enable the learning of the deep neural networks without mask labels, they design two loss functions that aim to minimize inter-image object discrepancy and maximize intra-image figure-ground separation. Image degradation is a common problem when an image segmentation model is applied to real world situations, while models trained with clean images may generate poor predictions. To solve this problem, Guo *et al.* [14] propose a novel Dense-Gram Network based on the Gram matrix to effectively reduce the gap.

Iterative refinement in segmentation. Many previous works use iterative designs to optimize segmentation results. For example, McIntosh *et al.* [35] adopt convolutional LSTMs to control computation budgets by adjusting the number of iteration steps. Wang *et al.* [50] propose to use saliency maps to iteratively refine segmentation results under the weakly supervised setting. Lin *et al.* [30] propose to use graphical models to optimize the mask learned by scribble supervision for semantic segmentation. In [20] and [61], iterative structures are also used to improve few-shot segmentation results. Compared with previous iterative structures, the main difference in our design is that our two branches can benefit from each other and meanwhile refine their own predictions in a cycle manner.

Co-saliency detection. Co-saliency detection [12], [15], [51], [62]–[64] is a closely related topic which aims to identify the common and salient objects from a group of images. Zhang *et al.* [62] explore intra-saliency prior transfer and deep inter-saliency mining for co-saliency detection. Wei *et al.* [51] introduce the group-wise feature representation learning and the collaborative learning to address the co-salient object discovery problem. To enable fast common information learning, Fan *et al.* [12] propose a network to simultaneously embed the appearance and semantic features through a co-attention projection strategy. Zhang *et al.* [67] utilize multi-level CNN features to improve the RGB-T

salient object detection, which shares similar spirits with our model variants. Li *et al.* [28] employ attention mechanisms to integrate cross-modal and cross-level complementarity from multi-modal data for RGB-D salient object detection. We also utilize attention mechanisms in our design, but the attention is used to establish regional connections between images, which has a different purpose. Zhang *et al.* [68] design a feature fusion network to undertake the RGB-T saliency detection task, which includes multi-scale, multi-modality, and multi-level feature fusion modules. Yu *et al.* [54] explore a novel problem setting that aims to identify common saliency within an image and propose a bottom-up framework to address the problem. Compared with co-saliency detection, object co-segmentation must identify objects from multiple categories from a complex scene, which are not necessarily the salient areas in images. For example, the buses and the people in Figure 1 are both common objects in the images.

3 METHOD

In this section, we present our framework for the image co-segmentation task. We begin with the description of our model in the case of paired input images. Each image is encoded and decoded with a parameter-shared Siamese network. The network design has an encoder-decoder structure, and the two branches exchange information at the network bottlenecks. We first describe our model version that only exploits the features at the last layer of the encoder for information exchange, which is shown in Figure 2. Then, we describe how to extend our network to the group-segmentation task without learning new parameters. Finally, we describe the model variant that uses multi-level features.

3.1 Cycle Refinement Module

The cycle refinement module (CRM) aims to update the image embeddings by incorporating information from the other image. The overall structure of the cycle refinement module can be found in Figure 2. In each step, the CRM implicitly compares the co-occurrent semantics of two images and lets the representations focus on common objects progressively. To achieve this goal, we employ a ConvLSTM at the network bottleneck to undertake the tasks of information exchange and representation updating. LSTM has proven its success in numerous computer vision and NLP tasks. Its feedback connections and memory cells make it well-suited for tasks with sequential inputs. The ConvLSTM further replaces the linear transformations in LSTM with convolutional kernels, such that the LSTM owns large field-of-views and is suited for processing image data. A typical ConvLSTM cell has the following structures:

$$i_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + b_i), \quad (1)$$

$$f_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + b_f), \quad (2)$$

$$o_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + b_o), \quad (3)$$

$$\tilde{C}_t = i_t \odot \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + b_c) \quad (4)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t, \quad (5)$$

$$H_t = o_t \odot \tanh(C_t), \quad (6)$$

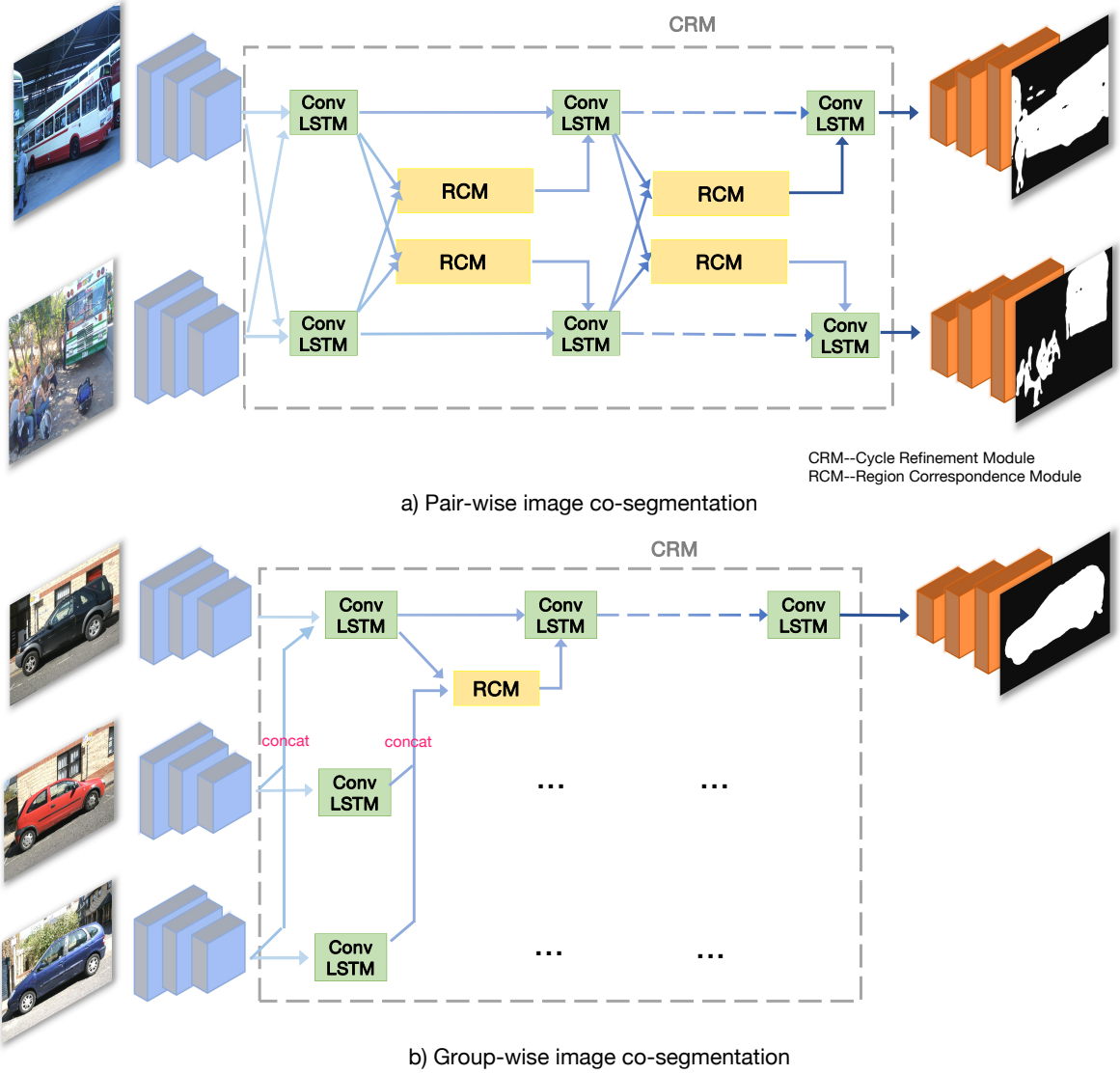


Fig. 2: Our network for image co-segmentation with paired (a) or grouped (b) input images. The main components in our network include a Cycle Refinement Module (CRM) for feature updating and a Region Correspondence Module (RCM) for information exchange.

where i_t is the input gate, which controls the activation of the new input information; f_t is the forget gate that clears the past cell status; o_t is the output gate, controlling whether the latest cell output C_t will be propagated to the final state H_t ; W is convolutional kernels; $*$ denotes the convolution operator and \odot denotes the Hadamard product. At each time step t , the ConvLSTM cell takes in an input X_t and updates the hidden states H_t and cell state C_t , which are both initialized by the image representations generated by the encoders. The updated cell state C_t is then sent to the region correspondence module (described in Section 3.2), which compares the cells from both sides and generates the inputs of the next step:

$$X_t^A = \text{RCM}(C_{t-1}^A, C_{t-1}^B), \quad (7)$$

$$X_t^B = \text{RCM}(C_{t-1}^B, C_{t-1}^A). \quad (8)$$

At the last step, the hidden state H is decoded by the decoder and generates the predicted mask of the common objects. Based on this recurrent process, both the quality of the transferred representations and the accuracy of predictions can be improved gradually. The number of steps for information exchange is flexible. In our experiment, we investigate the influence of the step numbers on the performance and the computation cost. Our proposed CRM for image co-segmentation has a symmetric structure design and all the operations on both sides are parameter-shared.

3.2 Region Correspondence Module

As we have seen in the cycle refinement module, the key component that exchanges information between two images is the region correspondence module. Previous works often achieve this goal through a global representation as the exchanged information — the exchanged global representation is then fused with the dense image representations

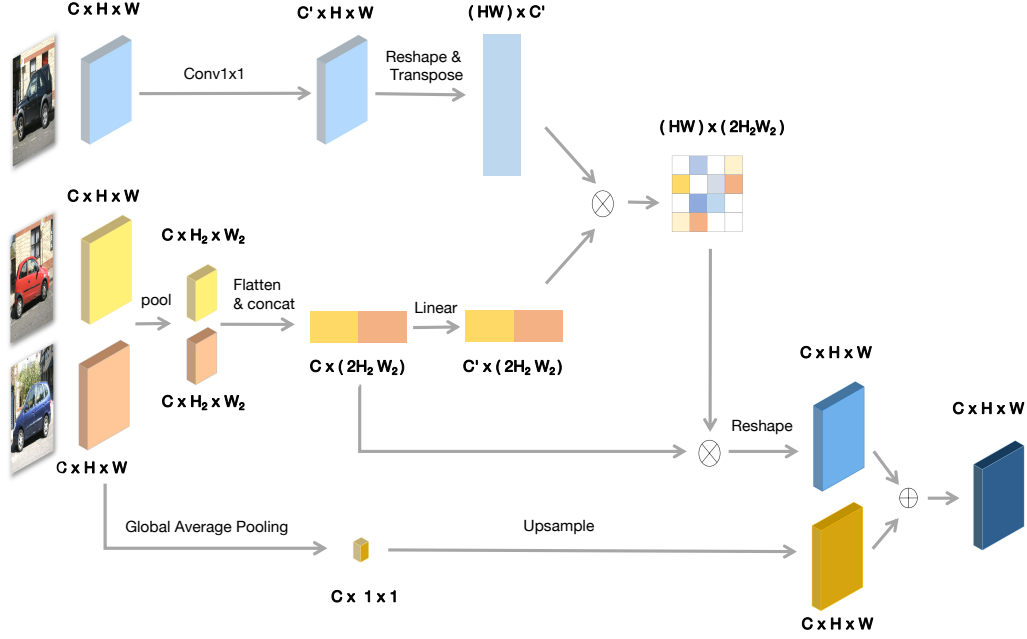


Fig. 3: Architecture of the Region Correspondence Module. \oplus indicates element-wise addition and \otimes indicates matrix multiplication.

by multiplication or concatenation. However, as image segmentation is a dense prediction task, local discriminative representations are important information to inference the common object regions. Therefore, we maintain the local representations of images and use attention mechanisms to establish the regional connections between two images. The architecture of RCM is illustrated in Figure 3. Specifically, the input to the region correspondence module is the current representations of two images $C_t^A \in \mathbb{R}^{H_A \times W_A \times C}$ and $C_t^B \in \mathbb{R}^{H_B \times W_B \times C}$, which are the cell states in the ConvLSTMs. When we want to transfer information from image I_B to image I_A , we first get the regional representations of image I_B by applying the ROI pooling to the representation C_t^B which downsamples the image representations to a fixed spatial size $\hat{H} \times \hat{W}$:

$$C_t'^B = \text{ROI}(C_t^B). \quad (9)$$

The regional representations can provide context information of local regions, which are proven useful in many previous segmentation works [69]. We apply both the ROI average pooling and the ROI max pooling to $C_t'^B$, and fuse their results by convolutions to get $C_t'^B$. The operations here share some similarities with the Channel Attention Block proposed in [52], where the global average pooling and global max pooling are used together to generate channel attentions. Then we compute the similarity between each feature point in C_t^A and $C_t'^B$ by dot product, which is implemented in parallel by matrix multiplication:

$$S_t^{AB} = W^A(C_t^A)(W^B(C_t'^B))^T, \quad (10)$$

where $S_t^{AB} \in \mathbb{R}^{H^A \times W^A \times \hat{H} \times \hat{W}}$, and W is the linear transformation function followed by ReLU non-linearity, implemented

as 1×1 convolutions. For notational convenience, we omit the reshape operations in Equation 10. Intuitively, the function W projects the feature representations into a space for computing feature similarity and the affinity matrix S_t^{AB} can reflect the correlations between local representations. Then, we normalize the affinity matrix by softmax and use the normalized affinity matrix to query features from the regional representations of image I^B :

$$X_t^{AB} = \text{softmax}(S_t^{AB})C_t'^B, \quad (11)$$

where $X_t^{AB} \in \mathbb{R}^{H_A \times W_A \times C}$. We also incorporate the global statistics of image I^B into the exchanged information by global average pooling, and upsample it to the same spatial size of image I^A . The final output of region correspondence module is

$$\text{RCM}(C_t^A, C_t^B) = (X_t^{AB} + \text{Up}(\text{AvgPool}(C_t^B)))/2. \quad (12)$$

When transferring information from image I^A to image I^B , we can simply swap the notation of A and B in the above equations, and apply the same operations. As we can see above, the exchanged information for each pixel location is different. Every pixel can attend all the local regions in the other image and selectively extract information from different regions.

3.3 Group Segmentation

Group segmentation is a special case of the co-segmentation task, where the goal is to segment the common objects from a group of images. An advantage of our network is that the image co-segmentation model for paired input images can also be applied to the group segmentation task with only minor modifications and without learning new parameters.

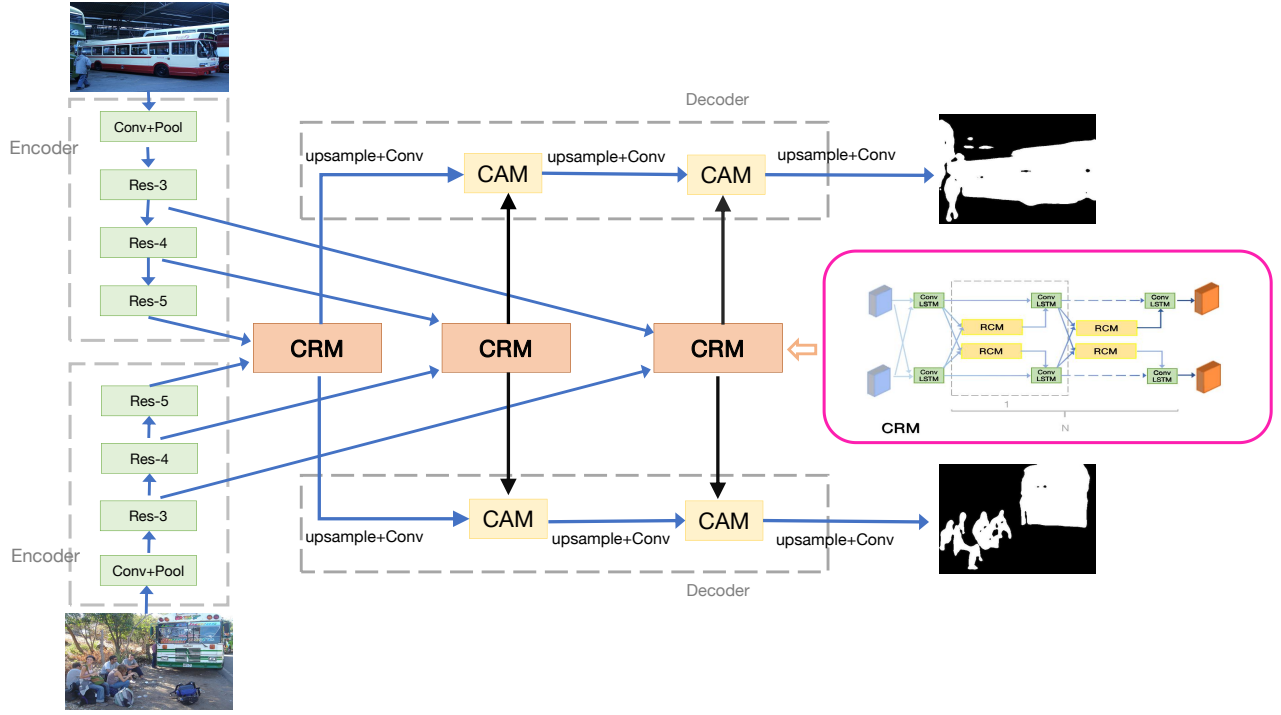


Fig. 4: Our model variant utilizing multi-level features in the encoder for image co-segmentation. We apply the **CRM** module (described in Section 3.2 and illustrated in Figure 2) to features from different layers independently and fuse the resulting features using the Channel Attention module (CAM) proposed in [53] during the upsampling stage.

Since all the operations in the region correspondence module do not require a fixed spatial size of the representations, we can use the region correspondence module to directly transfer information from all other images to the target image by treating all other images as a big virtual image. Figure 2 (b) illustrates how our co-segmentation network can be used to undertake the group segmentation task. For example, when a group segmentation task has three input images, the prediction of an image should refer to information from the other two images. We achieve this goal by constructing an affinity matrix that contains the similarity between the regional representations in the target image and all other images. We can therefore selectively extract information from all other images based on the attention distribution. Specifically, the Equation 9 and Equation 10 become:

$$C_t'^{BC} = \text{Flatten}(\text{ROI}(C_t^B)) \parallel \text{Flatten}(\text{ROI}(C_t^C)), \quad (13)$$

$$S_t^{A_BC} = W^A(C_t^A)(W^B(C_t'^{BC}))^T, \quad (14)$$

where Flatten is the operation that reshapes the 2D tensor to 1D tensor and \parallel denotes the concatenation operator. To make it more clear, the shapes of the tensors above are listed here: $C_t'^{BC} \in \mathbb{R}^{(\hat{H}\hat{W} + \hat{H}\hat{W}) \times C}$ and $S_t^{A_BC} \in \mathbb{R}^{H^A W^A \times (\hat{H}\hat{W} + \hat{H}\hat{W})}$. As a result, we can reuse all the operations in the region correspondence module to undertake the group segmentation task. The above operations transform the network-like relations between individual images to a one-versus-rest relation, such that we can use the co-segmentation model to solve the group segmentation problem. We can repeat such operations to make predictions for each of the images in the group using the same parameters.

3.4 Multi-Level Features for Co-Segmentation

As the goal of co-segmentation is to segment common objects in two images, the similar regions in two images are important cues for prediction. We observe that such similarity may exist at multiple semantic levels. For example, if both the training and testing images contain the class *car*, high-level features that correspond to object categories are useful for locating such a category. However, the testing images also contain novel classes that are never seen at training time. In this case, middle-level features that correspond to object parts would be useful, as they are more likely to be shared across classes. Based on such intuition, we design a model variant that exploits multi-level features in the encoder to undertake the co-segmentation task. We apply our cycle refinement module to features on different layers individually, and fuse their representations in the upsampling stage, as shown in Figure 4. We adopt the feature fusion module proposed in [53] to fuse the representations from different levels. The multi-level version can also apply to the group-segmentation task with multiple input images. Multi-level features have proven useful in many vision tasks, such as semantic segmentation [40] and object detection [39]. Our experiments in Section 4.3 show that using multi-level features to reason the common object region is helpful in the co-segmentation task.

4 EXPERIMENT

4.1 Implementation Details.

Network. The best performance of our network is obtained when we employ the ResNet34 [18] as the encoder network and use multi-level features from the last three layers; the

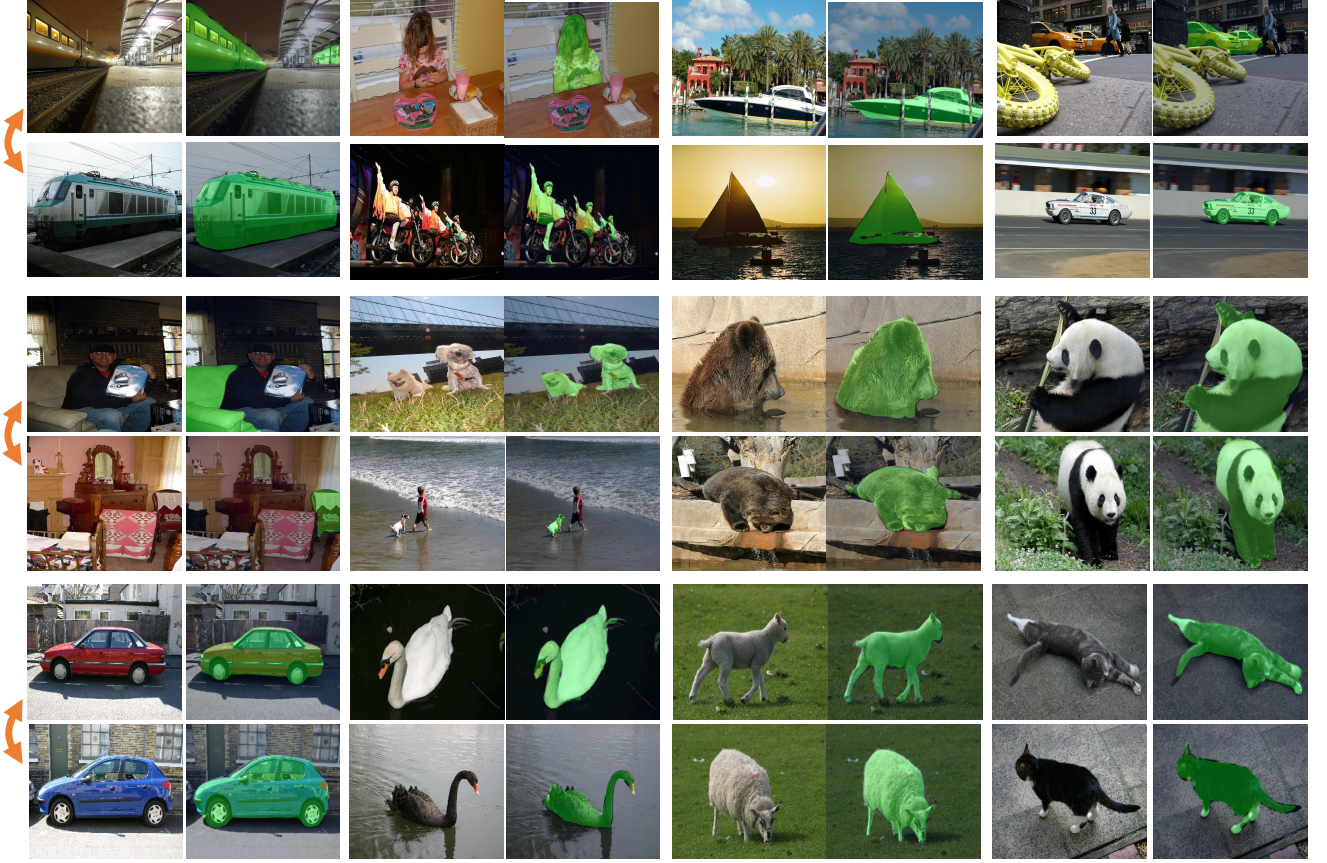


Fig. 5: Qualitative results of our network on PASCAL VOC dataset, MSRC dataset, Internet dataset, and iCoseg dataset.

number of iterative steps is set as 7; the ROI region size is set as 2×2 . If not further specified, we use them as the default setting in the experiments. The encoder is pretrained on ImageNet dataset [9]. For the model version utilizing single-level features, we use the same decoder structure with [5] which contains five blocks of bilinear upsampling layers and convolutional layers.

Training. We use Adam algorithm [26] with the learning rate of $1e-5$ to optimize the whole network in an end-to-end manner and set the weight decay as 0.0005. For both training and evaluation, we resize every image from the datasets to the spatial size of 512×512 . We set the batch size as 4 and train our network on the NVIDIA 1080TI GPU for 40K iterations. We use Lovász-Softmax proposed in [3] as the loss function, as we find it can yield slightly better performance than the pixel-wise cross-entropy loss. The overall loss function is given as:

$$\mathcal{L} = \frac{1}{C} \sum_c \Delta_{J_c}^-(m(c)), \quad (15)$$

$$m_i(c) = \begin{cases} 1 - p_i(c) & \text{if } c = y_i^*(c), \\ p_i(c) & \text{otherwise,} \end{cases}$$

where C denotes the number of classes, $\Delta_{J_c}^-(\cdot)$ is the Lovász extension of the Jaccard Index, $m(c)$ is a vector of pixel errors for class c , $y_i^* \in \{-1, 1\}$ is the ground truth label of pixel i for class c , and $p_i(c) \in [0, 1]$ is the predicted probability of pixel i for class c . Please refer to [3] for more details about Lovász-Softmax loss function.

4.2 Datasets and Evaluation Metric

As previous models use different training datasets in this task, we report the performance of our model with different training datasets for a fair comparison. Specifically, in [5], [29], the training set is constructed based on the training set of PASCAL VOC2012 dataset while in [27], [66], the training set is based on COCO dataset [31]. We evaluate our proposed method and compare it with existing methods on four widely-used benchmark datasets for object co-segmentation, including PASCAL VOC dataset [10], MSRC dataset [46], Internet dataset [42], and iCoseg dataset [2].

PASCAL VOC. PASCAL VOC2012 dataset contains 20 foreground object classes and one background class and it has 1,464 training images and 1,449 validation images in total. Following [5], [29], [49], we split the original validation set into a validation set (724 images) and a test set (725 images) for the co-segmentation task. For PASCAL VOC 2010, we follow the setting in [27], [66], where a total of 1,037 images of 20 objects classes from PASCAL VOC 2010 dataset are used for evaluation.

MSRC. Following [29], [66], we use the same subset of MSRC dataset. This subset includes 7 classes: bird, car, cat, cow, dog, plane, and sheep. Each class contains 10 images.

Internet. Internet dataset contains images of three categories including airplane, car, and horse. Following previous works [5], [19], we use the same subset of the Internet dataset where each class has 100 images.

iCoseg. iCoseg dataset consists of 643 images from 38 categories. Large variances of viewpoints and deformations

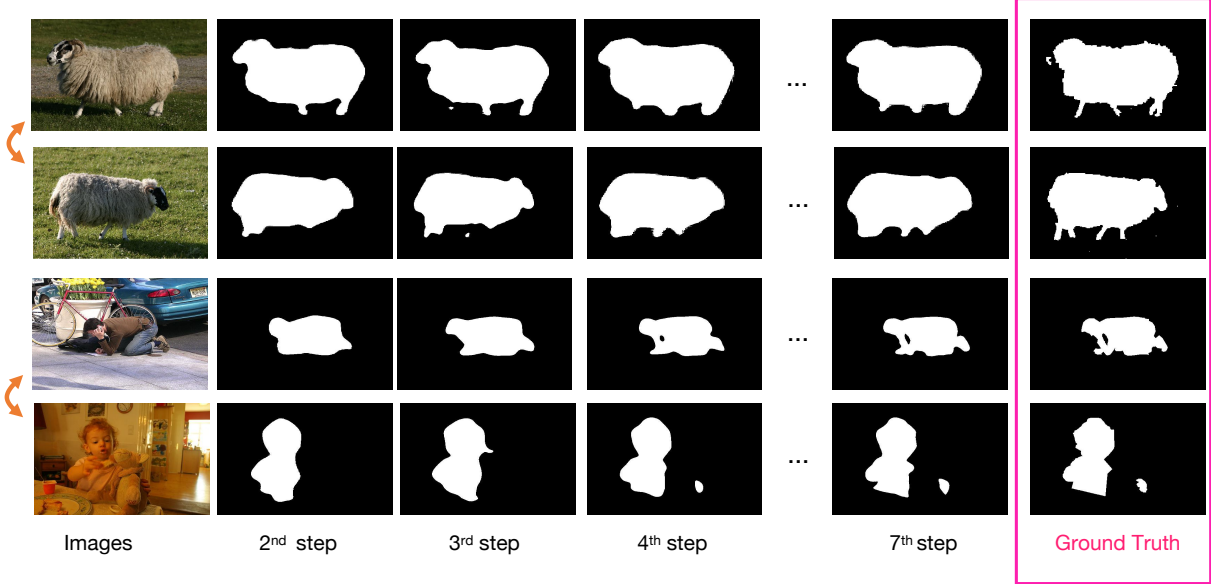


Fig. 6: Predictions from different refinement steps. We decode the representations in each time step and the result shows that our proposed cycle refinement module can consistently improve the prediction.

M_{cat}	M_{mul}	M_{cross}	M_{cycle}	\mathcal{P} (%)	\mathcal{J} (%)
				77.2	58.6
✓				81.3	62.9
	✓			83.6	64.2
		✓		87.1	68.2
		✓	✓	91.8	70.4

TABLE 1: Ablative analysis of the proposed model. “ M_{cat} ”, “ M_{mul} ” and “ M_{cross} ” denote three kinds of ways to exchange information between images. “ M_{cycle} ” denotes our proposed cycle refinement module. Our designs (“ M_{cross} ” and “ M_{cycle} ”) outperform the baselines significantly. Please refer to Section 4.3 for analysis.

Num. of steps (N)	\mathcal{P} (%)	\mathcal{J} (%)	Time(s)
$N=2$	92.2	71.9	0.085
$N=3$	93.8	73.1	0.088
$N=4$	94.2	74.0	0.096
$N=5$	94.8	74.5	0.105
$N=6$	95.3	75.2	0.109
$N=7$	95.8	75.4	0.114

TABLE 2: Our network with different refinement steps N in the Cycle Refinement Module. As the refinement step increases, the performance grows consistently.

are present in this dataset. A subset that contains 8 classes is used to evaluate the generalizability of our proposed method.

We evaluate the performance of co-segmentation models using two common evaluation metrics: *Precision* (\mathcal{P}) and *Jaccard Index* (\mathcal{J}). We report the performance under both metrics on these four co-segmentation benchmark datasets in our experiments. All the experiments for analysis are trained and evaluated on PASCAL VOC 2012 dataset.

4.3 Analysis

In this part, we investigate the effectiveness of each component in our proposed CycleSegNet by ablative analysis. We first create a baseline model which is a foreground prediction network where two branches do not exchange any information. Then we incorporate another two baseline methods M_{cat} and M_{mul} that exchange information in the form of global representations. Concretely, M_{cat} is to upsample the global vector and fuse it with the feature maps by concatenation, while M_{mul} is to fuse by multiplication. Then we add each of our proposed modules to the baseline

model one by one. All the models are tested with the single-level feature encoder. The results are shown in Table 1. As we can see, both our proposed modules are very effective in improving the performance over the baseline networks. Specifically, our region correspondence module (M_{cross}) significantly outperforms both baseline methods M_{cat} and M_{mul} by 5.8% and 3.5% in terms of Precision and 5.3% and 4.0% in terms of Jaccard, which shows that our design that exchanges information between local regions is more effective than previous methods using global representations. The cycle refinement module further boosts the Precision and Jaccard score by 4.7% and 2.2%, respectively.

Number of the refinement step. We report the performance of our network under different refinement steps N in Table 2. We also compare the average time to inference one pair of images. As we can see from the table, when the number of the refinement steps increases, the performance in terms of both Precision and Jaccard increases consistently while the increment in inference time is marginal. Particularly, after 7 steps of refinement, the network can improve the initial predictions by 3.6% and 3.5% in terms of Precision and Jaccard, respectively.

Region Size	$\mathcal{P}(\%)$	$\mathcal{J}(\%)$
Raw	92.9	73.1
2×2	95.8	75.4
3×3	94.0	74.1
4×4	95.1	74.8
5×5	93.6	74.0
6×6	94.8	74.4
7×7	94.6	74.3

TABLE 3: Comparison of different ROI pooling sizes in the region correspondence module. The optimal result is obtained when the region size is 2×2 .

res_5	res_4	res_3	$\mathcal{P}(\%)$	$\mathcal{J}(\%)$
✓			91.8	70.4
✓	✓		93.5	72.5
✓	✓	✓	95.8	75.4
vgg_5	vgg_4	vgg_3	$\mathcal{P}(\%)$	$\mathcal{J}(\%)$
✓			91.2	69.8
✓	✓		93.1	71.8
✓	✓	✓	94.9	74.2

TABLE 4: Performance of our method utilizing features from different stages in the encoder network. We conduct experiments with ResNet-34 and VGG16 as the network encoders. Using multi-level features can effectively boost network performance.

Input Size	Strategy a)			Strategy b)			Strategy c)			Strategy d)		
	Car	Airplane	Horse	Car	Airplane	Horse	Car	Airplane	Horse	Car	Airplane	Horse
$k = 2$	85.8	77.2	74.1	84.4	78.3	73.9	85.9	77.3	74.0	85.7	77.3	74.1
$k = 3$	86.6	78.1	75.0	84.8	78.8	74.1	86.6	78.1	75.0	86.8	78.5	75.1
$k = 4$	-	-	-	85.6	78.9	75.1	86.9	78.6	75.4	87.0	78.8	75.5
$k = 5$	-	-	-	86.5	79.3	75.4	86.9	78.6	75.5	87.0	78.8	75.5
$k = 6$	-	-	-	86.4	79.0	75.4	86.9	78.7	75.4	86.9	78.7	75.5
$k = 7$	-	-	-	86.3	78.6	75.2	86.7	78.7	75.4	86.8	78.7	75.4
$k = 8$	-	-	-	86.1	78.9	75.3	86.7	78.7	75.5	86.8	78.7	75.5

TABLE 5: The results (Jaccard) of group segmentation with different sampling strategies and different numbers of input images on Internet dataset. Our network can handle more than two images at a time, which generates better results than our network with paired input images. Please refer to Section 4.3 for the description and analysis of different strategies.

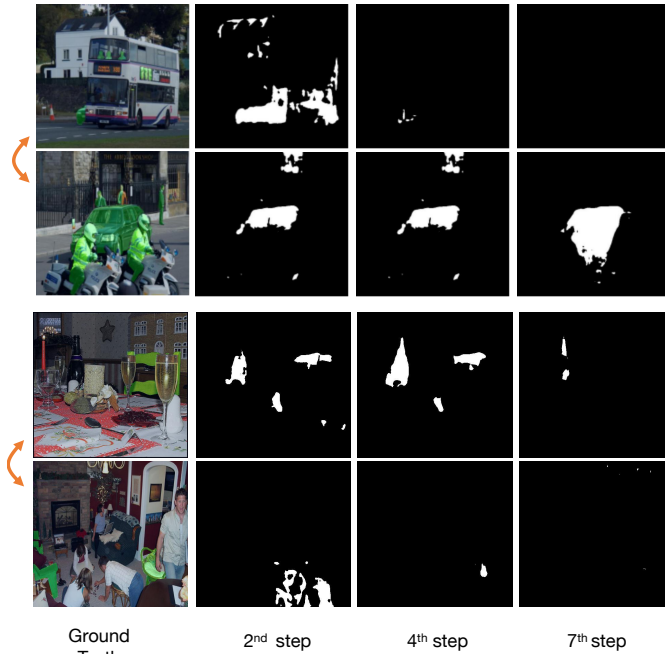


Fig. 7: Failure cases. The ground truth is denoted with green masks in the RGB images. When the common objects demonstrate a large variance in appearance, the network may fail to establish meaningful connections between local regions, and the cycle refinement module can hardly further refine the predictions.

The size of the regional representations. In the region correspondence module, we use the ROI pooling to get the regional representations of images. We next investigate how the region size influences the performance in Table 3. As we can see, the best performance is observed when ROI pooling outputs a 2×2 feature map. Directly using the original representations without ROI pooling does not yield better performance than regional representations. A possible reason is that the feature points in the original feature representations have a relatively small effective field-of-view and thus lack the expressive power of abstract and semantic concepts. Using the ROI pooling can enlarge the field-of-view and provide context information in local regions, which is helpful in computing the correlation scores between regions for information exchange.

Multi-level features. We compare the performance of our network using features from different layers in the encoder network in Table 4. As can be observed, multi-level features can boost the performance of our network with both the VGG16 backbone and the ResNet34 backbone. Specifically, using the features from the last three layers is better than using the features from the last two layers. This indicates that middle-level features are helpful and can provide important cues in the co-segmentation task.

Group Segmentation. We present a detailed analysis about group segmentation in this part. The group segmentation aims to predict the common objects in a group of images. We apply different CNN based solutions in previous works to our network and then compare their performance. Since the networks in most previous works for co-segmentation can only handle paired inputs, they often adopt various sampling strategies and fusion methods to undertake the group segmentation tasks, which can be applied to our network as well. An advantage of our

VOC2012(pair)	Training Set	P(%)	J(%)
CA [5]	VOC2012	-	59.2
FCA [5]	VOC2012	-	59.4
CSA [5]	VOC2012	-	59.8
DOCS [29]	VOC2012	94.2	64.5
CycleSegNet(Ours)	VOC2012	95.8	75.4

(a) Pairwise co-segmentation results on the PASCAL VOC2012 dataset.

VOC2010 (group)	Training Set	P(%)	J(%)
Quan <i>et al.</i> [38]	-	89.0	52.0
Jerripothula <i>et al.</i> [23]	-	85.2	45.0
Jerripothula <i>et al.</i> [21]	-	80.1	40.0
Li <i>et al.</i> [27]	COCO	94.1	63.0
Zhang <i>et al.</i> [66]	COCO	94.9	71.0
CycleSegNet(Ours)	COCO	96.8	73.6

(b) Group segmentation results on PASCAL VOC2010 dataset.

MSRC	Training Set	group		pair	
		P(%)	J(%)	P(%)	J(%)
Vicente <i>et al.</i> [47]	-	90.2	70.6	-	-
Wang <i>et al.</i> [48]	-	92.2	-	-	-
Rubinstein <i>et al.</i> [42]	-	92.2	74.7	-	-
Faktor <i>et al.</i> [11]	-	92.0	77.0	-	-
Chen <i>et al.</i> [5]	VOC2012	-	73.9	96.6	76.5
Li <i>et al.</i> [29]	VOC2012	95.4	82.9	-	-
CycleSegNet(Ours)	VOC2012	97.9	87.2	97.3	86.8
Zhang <i>et al.</i> [66]	COCO	95.2	81.9	-	-
CycleSegNet(Ours)	COCO	97.6	89.6	97.2	88.2

(c) Results on the MSRC dataset.

Internet	Training Set	group				pair			
		Airplane	Car	Horse	Avg.J(%)	Airplane	Car	Horse	Avg.J(%)
Rubinstein <i>et al.</i> [42]	-	55.8	64.4	51.6	57.3	-	-	-	-
Quan <i>et al.</i> [38]	-	56.3	66.8	58.1	60.4	-	-	-	-
Jerripothula <i>et al.</i> [23]	-	61.0	71.0	60.0	64.0	-	-	-	-
Chen <i>et al.</i> [8]	-	65.0	82.0	63.0	70.0	-	-	-	-
Yuan <i>et al.</i> [56]	VOC2012	66.0	72.0	65.0	67.7	-	-	-	-
Li <i>et al.</i> [29]	VOC2012	65.4	82.8	69.4	72.6	-	-	-	-
Chen <i>et al.</i> [5]	VOC2012	-	-	-	-	71.4	79.9	68.0	73.1
CycleSegNet(Ours)	VOC2012	78.8	87.0	75.5	80.4	77.2	85.8	74.1	79.0
Zhang <i>et al.</i> [66]	COCO	69.6	82.5	70.2	74.1	-	-	-	-
Li <i>et al.</i> [27]	COCO	83.0	93.0	76.0	84.0	-	-	-	-
CycleSegNet(Ours)	COCO	84.1	91.6	82.9	86.2	83.2	88.5	80.3	84.0

(d) Results on Internet dataset.

iCoseg (pair)	Training Set	bear2	brownbear	cheetah	elephant	helicopter	hotballoon	panda1	panda2	Avg.J(%)
Li <i>et al.</i> [29]	VOC2012	88.3	92.0	68.8	84.6	79.0	91.7	82.6	86.7	84.2
Chen <i>et al.</i> [5]	VOC2012	88.3	91.5	71.3	84.4	76.5	94.0	91.8	90.3	86.0
CycleSegNet(Ours)	VOC2012	92.1	94.8	84.5	90.2	80.2	95.8	94.5	94.1	90.8

(e) Pairwise co-segmentation results on iCoseg dataset.

iCoseg (group)	Training Set	bear2	brownbear	cheetah	elephant	helicopter	hotballoon	panda1	panda2	Avg.J(%)
Rubinstein <i>et al.</i> [42]	-	65.3	73.6	69.7	68.8	80.3	65.7	75.9	62.5	70.2
Jerripothula <i>et al.</i> [22]	-	70.1	66.2	75.4	73.5	76.6	76.3	80.6	71.8	73.8
Faktor <i>et al.</i> [11]	-	72.0	66.2	75.4	73.5	76.6	76.3	80.6	71.8	73.8
Jerripothula <i>et al.</i> [23]	-	67.5	72.5	78.0	79.9	80.0	80.2	72.2	61.4	70.4
Zhang <i>et al.</i> [66]	COCO	91.1	89.6	88.6	90.0	76.4	94.2	90.4	87.5	89.2
CycleSegNet(Ours)	COCO	92.8	94.1	89.1	91.6	85.2	95.1	94.8	94.1	92.1

(f) Group segmentation results on iCoseg dataset.

TABLE 6: Comparison with the state-of-the-art performance of pairwise co-segmentation and group segmentation tasks on four benchmark datasets: PASCAL VOC dataset, MSRC dataset, Internet dataset and iCoseg dataset. Our model achieves new state-of-the-art performance on all datasets.

network is that our network can handle more than two images at a time, as is described earlier, so the number of input images is flexible. However, directly sending all images in the group to our network is unrealistic due to the limited GPU memory. We therefore sample a small group of k images at a time and send them to our network for predictions. Then, we fuse all the predictions of an image as the final predicted mask. Suppose that there are N images to be segmented in a group segmentation task and our network handles a tuple of k images at a time. We compare the following strategies to undertake group segmentation:

a) Following [1], we sample all possible k -element tuples

to make predictions and average all the corresponding confidence maps of an image as the final prediction. This method can make use of all images in the group to make predictions for each image.

b) Following [66], we randomly divide all images into N/k small groups and each group has k images. In this case, each image is only tested once and the whole evaluation process is quick.

c) Following [29], to predict the mask of an image, we randomly sample 5 groups of images from the other $N - 1$ images, where each group has $k - 1$ images. Then, the target image joins each of these groups to make a prediction with

our network, and we average 5 predictions to generate the final mask. The process is repeated for all N images.

d) Following [49], to predict the mask of an image, we uniformly split the other $N - 1$ images into T groups, where $T = (N - 1)/(k - 1)$ and each group has $k - 1$ images. Then, the target image joins each of these groups to make a prediction with our network, and we average T predictions to generate the final mask. The process is repeated for all N images.

Besides these solutions, there are also some previous works using graphs [16], [38] or recurrent models [49] to undertake group segmentation, which, however, can not be applied to our network directly. Therefore, we compare the four solutions above with our network. We conduct experiments on the subset of the Internet dataset, where each class has 100 images. For the solution **a**, there would be too many possible tuples for testing when k is large. Therefore, we only test the cases where $k = 2$ and $k = 3$. The result is shown in Table 5. As we can see, predicting with more than two input images is consistently better than predicting with paired input images with all four strategies, which shows the superiority of our network design for group segmentation. The best tuple size is 4 or 5, while further increasing the tuple size can not boost the performance.

Visualization. We present the visualization of our network predictions to qualitatively evaluate our design in Figure 5. As we can see, our model can well segment the common objects from the complex scenes in two images. To further observe how our cycle refinement module improves the predictions, we visualize the predictions in each refinement step, as shown in Figure 6. We can see from the result that the predicted masks become more accurate with more refinement steps. We also present some failure cases in Figure 7. As is shown, when the objects in different images demonstrate a large appearance variance, such as huge differences in object sizes and poses, the network fails to establish meaningful correspondence between local regions for predictions, and the CRM can not gradually improve the results.

4.4 Comparison with the State-of-the-Arts

Finally, we compare our proposed network with the state-of-the-art (SOTA) methods on four popular benchmark datasets: PASCAL VOC dataset, MSRC dataset, Internet dataset, and iCoseg dataset. We report the model performance for group segmentation as well as the pairwise co-segmentation. The results are shown in Table 6. As can be seen, our method outperforms all previous methods on different benchmarks with significant margins. In particular, on the challenging PASCAL VOC2012 dataset, the performance of our method outperforms the SOTA pairwise co-segmentation result by **10.9%** in terms of Jaccard (Table 6 (a)); on the MSRC dataset, when the models are trained with PASCAL VOC 2012 dataset, our performance of group segmentation and pairwise co-segmentation outperform the SOTA result by **4.3%** and **10.3%**, respectively, in terms of Jaccard (Table 6 (c)); on Internet dataset, the performance of our model trained with PASCAL VOC 2012 dataset on the group segmentation task and pairwise co-segmentation task outperforms the SOTA result by **7.8%** and **5.9%**, respectively, in terms of Jaccard (Table 6(d)); on iCoseg dataset,

our pairwise co-segmentation performance outperforms the SOTA result by **4.8%**, in terms of Jaccard (Table 6(e));

5 CONCLUSIONS

In this paper, we propose a novel and effective approach for the image co-segmentation task. The proposed region correspondence module directly exchanges information between local regions from different images, which demonstrates advantages over the baseline methods that transfer global image representations. The cycle refinement module that employs ConvLSTMs to progressively exchange information between images and update image representations can consistently improve the network predictions. Our algorithm can handle the input with paired images as well as a group of images with the same network and parameters. The multi-level feature encoder can further boost the network performance effectively. Experiment results on four object co-segmentation datasets — PASCAL VOC dataset, MSRC dataset, Internet dataset, and iCoseg dataset demonstrate that our proposed method significantly outperforms the existing methods and achieves new state-of-the-art performance.

ACKNOWLEDGMENTS

This research is supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG-RP-2018-003), the MOE Tier-1 research grants: RG28/18 (S), RG22/19 (S) and RG95/20, and the National Natural Science Foundation of China (No. 61772530). This work is also supported by National Natural Science Foundation of China (NSFC) 61876208, Key-Area Research and Development Program of Guangdong Province 2018B010108002, and Central Universities of China under Grant D2192860.

REFERENCES

- [1] Banerjee, S., Hati, A., Chaudhuri, S., Velmurugan, R.: Cosegnet: Image co-segmentation using a conditional siamese convolutional network. In: Proc. 28th Int. Joint Conf. Artif. Intell. pp. 673–679 (2019)
- [2] Batra, D., Kowdle, A., Parikh, D., Luo, J., Chen, T.: icoseg: Interactive co-segmentation with intelligent scribble guidance. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 3169–3176. IEEE (2010)
- [3] Berman, M., Rannen Triki, A., Blaschko, M.B.: The lovász-softmax loss: a tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4413–4421 (2018)
- [4] Chang, K.Y., Liu, T.L., Lai, S.H.: From co-saliency to co-segmentation: An efficient and fully unsupervised energy minimization model. In: CVPR 2011. pp. 2129–2136. IEEE (2011)
- [5] Chen, H., Huang, Y., Nakayama, H.: Semantic aware attention based deep object co-segmentation. In: Asian Conference on Computer Vision. pp. 435–450. Springer (2018)
- [6] Chen, H.Y., Lin, Y.Y., Chen, B.Y.: Co-segmentation guided hough transform for robust feature matching. IEEE transactions on pattern analysis and machine intelligence 37(12), 2388–2401 (2015)
- [7] Chen, X., Zhang, C., Lin, G., Han, J.: Compositional prototype network with multi-view comparison for few-shot point cloud semantic segmentation. arXiv preprint arXiv:2012.14255 (2020)
- [8] Chen, Y.C., Lin, Y.Y., Yang, M.H., Huang, J.B.: Show, match and segment: Joint learning of semantic matching and object co-segmentation. arXiv preprint arXiv:1906.05857 (2019)

- [9] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
- [10] Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. *International journal of computer vision* **111**(1), 98–136 (2015)
- [11] Faktor, A., Irani, M.: Co-segmentation by composition. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1297–1304 (2013)
- [12] Fan, D.P., Li, T., Lin, Z., Ji, G.P., Zhang, D., Cheng, M.M., Fu, H., Shen, J.: Re-thinking co-salient object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021)
- [13] Fu, J., Liu, J., Wang, Y., Zhou, J., Wang, C., Lu, H.: Stacked deconvolutional network for semantic segmentation. *IEEE Transactions on Image Processing* (2019)
- [14] Guo, D., Pei, Y., Zheng, K., Yu, H., Lu, Y., Wang, S.: Degraded image semantic segmentation with dense-gram networks. *IEEE Transactions on Image Processing* **29**, 782–795 (2019)
- [15] Han, J., Cheng, G., Li, Z., Zhang, D.: A unified metric learning-based framework for co-saliency detection. *IEEE Transactions on Circuits and Systems for Video Technology* **28**(10), 2473–2483 (2017)
- [16] Han, J., Quan, R., Zhang, D., Nie, F.: Robust object co-segmentation using background prior. *IEEE Transactions on Image Processing* **27**(4), 1639–1651 (2017)
- [17] He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
- [18] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
- [19] Hsu, K.J., Lin, Y.Y., Chuang, Y.Y.: Co-attention cnns for unsupervised object co-segmentation. In: *IJCAI*. pp. 748–756 (2018)
- [20] Hu, T., Yang, P., Zhang, C., Yu, G., Mu, Y., Snoek, C.G.: Attention-based multi-context guiding for few-shot semantic segmentation. In: Proceedings of the AAAI conference on artificial intelligence. vol. 33, pp. 8441–8448 (2019)
- [21] Jerripothula, K.R., Cai, J., Lu, J., Yuan, J.: Object co-skeletonization with co-segmentation. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3881–3889. IEEE (2017)
- [22] Jerripothula, K.R., Cai, J., Meng, F., Yuan, J.: Automatic image co-segmentation using geometric mean saliency. In: 2014 IEEE International Conference on Image Processing (ICIP). pp. 3277–3281. IEEE (2014)
- [23] Jerripothula, K.R., Cai, J., Yuan, J.: Image co-segmentation via saliency co-fusion. *IEEE Transactions on Multimedia* **18**(9), 1896–1909 (2016)
- [24] Jing, L., Chen, Y., Tian, Y.: Coarse-to-fine semantic segmentation from image-level labels. *IEEE Transactions on Image Processing* **29**, 225–236 (2019)
- [25] Joulin, A., Bach, F., Ponce, J.: Discriminative clustering for image co-segmentation. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 1943–1950. IEEE (2010)
- [26] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
- [27] Li, B., Sun, Z., Li, Q., Wu, Y., Hu, A.: Group-wise deep object co-segmentation with co-attention recurrent neural network. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 8519–8528 (2019)
- [28] Li, C., Cong, R., Kwong, S., Hou, J., Fu, H., Zhu, G., Zhang, D., Huang, Q.: Asif-net: Attention steered interweave fusion network for rgb-d salient object detection. *IEEE transactions on cybernetics* **51**(1), 88–100 (2020)
- [29] Li, W., Jafari, O.H., Rother, C.: Deep object co-segmentation. In: Asian Conference on Computer Vision. pp. 638–653. Springer (2018)
- [30] Lin, D., Dai, J., Jia, J., He, K., Sun, J.: Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3159–3167 (2016)
- [31] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
- [32] Liu, W., Zhang, C., Lin, G., HUNG, T.Y., Miao, C.: Weakly supervised segmentation with maximum bipartite graph matching. In: Proceedings of the 28th ACM International Conference on Multimedia. pp. 2085–2094 (2020)
- [33] Liu, W., Zhang, C., Lin, G., Liu, F.: Crnet: Cross-reference networks for few-shot segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4165–4173 (2020)
- [34] Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)
- [35] McIntosh, L., Maheswaranathan, N., Sussillo, D., Shlens, J.: Recurrent segmentation for variable computational budgets. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. pp. 1648–1657 (2018)
- [36] Mukherjee, L., Singh, V., Peng, J.: Scale invariant cosegmentation for image groups. In: CVPR 2011. pp. 1881–1888. IEEE (2011)
- [37] Mustafa, A., Hilton, A.: Semantically coherent co-segmentation and reconstruction of dynamic scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 422–431 (2017)
- [38] Quan, R., Han, J., Zhang, D., Nie, F.: Object co-segmentation via graph optimized-flexible manifold ranking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 687–695 (2016)
- [39] Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems. pp. 91–99 (2015)
- [40] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)
- [41] Rother, C., Minka, T., Blake, A., Kolmogorov, V.: Cosegmentation of image pairs by histogram matching-incorporating a global constraint into mrfs. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06). vol. 1, pp. 993–1000. IEEE (2006)
- [42] Rubinstein, M., Joulin, A., Kopf, J., Liu, C.: Unsupervised joint object discovery and segmentation in internet images. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1939–1946 (2013)
- [43] Rubio, J.C., Serrat, J., López, A., Paragios, N.: Unsupervised co-segmentation through region matching. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 749–756. IEEE (2012)
- [44] Sakinis, T., Milletari, F., Roth, H., Korfiatis, P., Kostandy, P., Philbrick, K., Akkus, Z., Xu, Z., Xu, D., Erickson, B.J.: Interactive segmentation of medical images through fully convolutional neural networks. *arXiv preprint arXiv:1903.08205* (2019)
- [45] Shi, X., Chen, Z., Wang, H., Yeung, D.Y., Wong, W.K., Woo, W.C.: Convolutional lstm network: A machine learning approach for precipitation nowcasting. *arXiv preprint arXiv:1506.04214* (2015)
- [46] Shotton, J., Winn, J., Rother, C., Criminisi, A.: Textonboost: Joint appearance, shape and context modeling for multi-class object recognition and segmentation. In: European conference on computer vision. pp. 1–15. Springer (2006)
- [47] Vicente, S., Rother, C., Kolmogorov, V.: Object cosegmentation. In: CVPR 2011. pp. 2217–2224. IEEE (2011)
- [48] Wang, F., Huang, Q., Guibas, L.J.: Image co-segmentation via consistent functional maps. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 849–856 (2013)
- [49] Wang, W., Lu, X., Shen, J., Crandall, D.J., Shao, L.: Zero-shot video object segmentation via attentive graph neural networks. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9236–9245 (2019)
- [50] Wang, X., You, S., Li, X., Ma, H.: Weakly-supervised semantic segmentation by iteratively mining common object features. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1354–1362 (2018)
- [51] Wei, L., Zhao, S., Bourahla, O.E.F., Li, X., Wu, F.: Group-wise deep co-saliency detection. *arXiv preprint arXiv:1707.07381* (2017)
- [52] Woo, S., Park, J., Lee, J.Y., So Kweon, I.: Cbam: Convolutional block attention module. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 3–19 (2018)
- [53] Yu, C., Wang, J., Peng, C., Gao, C., Yu, G., Sang, N.: Learning a discriminative feature network for semantic segmentation. In: Pro-

- ceedings of the IEEE conference on computer vision and pattern recognition. pp. 1857–1866 (2018)
- [54] Yu, H., Zheng, K., Fang, J., Guo, H., Feng, W., Wang, S.: Co-saliency detection within a single image. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 32 (2018)
 - [55] Yu, J.G., Li, Y., Gao, C., Gao, H., Xia, G.S., Yu, Z.L., Li, Y.: Exemplar-based recursive instance segmentation with application to plant image analysis. *IEEE Transactions on Image Processing* **29**, 389–404 (2019)
 - [56] Yuan, Z.H., Lu, T., Wu, Y.: Deep-dense conditional random fields for object co-segmentation. In: *IJCAI*. pp. 3371–3377 (2017)
 - [57] Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: *European conference on computer vision*. pp. 818–833. Springer (2014)
 - [58] Zhang, C., Cai, Y., Lin, G., Shen, C.: Deepemd: Differentiable earth mover’s distance for few-shot learning (2020)
 - [59] Zhang, C., Cai, Y., Lin, G., Shen, C.: Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 12203–12213 (2020)
 - [60] Zhang, C., Lin, G., Liu, F., Guo, J., Wu, Q., Yao, R.: Pyramid graph networks with connection attentions for region-based one-shot semantic segmentation. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 9587–9595 (2019)
 - [61] Zhang, C., Lin, G., Liu, F., Yao, R., Shen, C.: Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5217–5226 (2019)
 - [62] Zhang, D., Han, J., Han, J., Shao, L.: Cosaliency detection based on intrasaliency prior transfer and deep intersaliency mining. *IEEE transactions on neural networks and learning systems* **27**(6), 1163–1176 (2015)
 - [63] Zhang, D., Han, J., Li, C., Wang, J., Li, X.: Detection of co-salient objects by looking deep and wide. *International Journal of Computer Vision* **120**(2), 215–232 (2016)
 - [64] Zhang, D., Meng, D., Han, J.: Co-saliency detection via a self-paced multiple-instance learning framework. *IEEE transactions on pattern analysis and machine intelligence* **39**(5), 865–878 (2016)
 - [65] Zhang, H., Tian, Y., Wang, K., Zhang, W., Wang, F.Y.: Mask ssd: An effective single-stage approach to object instance segmentation. *IEEE Transactions on Image Processing* **29**(1), 2078–2093 (2019)
 - [66] Zhang, K., Chen, J., Liu, B., Liu, Q.: Deep object co-segmentation via spatial-semantic network modulation. *arXiv preprint arXiv:1911.12950* (2019)
 - [67] Zhang, Q., Huang, N., Yao, L., Zhang, D., Shan, C., Han, J.: Rgb-t salient object detection via fusing multi-level cnn features. *IEEE Transactions on Image Processing* **29**, 3321–3335 (2019)
 - [68] Zhang, Q., Xiao, T., Huang, N., Zhang, D., Han, J.: Revisiting feature fusion for rgb-t salient object detection. *IEEE Transactions on Circuits and Systems for Video Technology* (2020)
 - [69] Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2881–2890 (2017)