Deep Semantic Matching with Foreground Detection and Cycle-Consistency

Yun-Chun Chen^{1,4} Po-Hsiang Huang¹ Li-Yu Yu¹ Jia-Bin Huang² Ming-Hsuan Yang³ Yen-Yu Lin⁴ ¹NTUEE ²Virginia Tech ³UC Merced ⁴Academia Sinica

Abstract Establishing dense semantic correspondences between object ADMETRIC. Instantisting dense semantic corresponsences between copiect instances remains a challenging problem due to background clutter, sig-nificant scale and pose differences, and large intra-class variations. In this are and pose unrecences, and arge intra-class variances. In this present an end-to-end trainable network for learning semantic denoses using only matching image pairs without manual loy-copondence annotations. To facilitate network training with this weaker form of supervision, we 1) explicitly estimate the foreground regions to supprose the effect of background clutter and 2) develop cycleconsistent losses to enforce the predicted transformations acro images to be geometrically plausible and consistent. We train the procosed model on the image pairs of the PF-PASCAL dataset and evaluate the learned model on the PF-PASCAL. PF-WILLOW, and TSS datasets experimental results show that the proposed approach achieves avorably performance compared to the state-of-the-art.

1 Introduction

Semantic matching is an important and active research topic in computer vision Previous methods such as optical flow estimation [5, 10] and stereo matching 30 rely on ner-nivel correspondence to match across images denicting the sam (a) rety on per-parate correspondence to match across images depicting the same scene or object instance. While correspondence estimation has been studied for years, there has been a growing trend to extend the idea of matching the same ching images covering different instar cotonom. This measure not only attracts a lot of attention but also facilitate category. This progress not only attracts a not of attention but also hacintate many real-world applications ranging from object recognition [16], object co-segmentation [23, 27], to 3D reconstruction [19]. However, due to the presence of background clutter, ambiguity induced by large intra-class variations, and the limited scalability of obtaining large-scale datasets with manually amountated

the minter scanonity of obtaining large-scale earliers with mannany annotated correspondences, semantic matching remains quite challenging. Conventional methods for semantic matching heavily rely on hand-crafted descriptors such as SIFT [16], HOG [4], or DAISY [29] as well as an effective goometric regularizer. However, these hand-crafted descriptors are pre-defined and cannot adapt themselves to the given visual domains, leading to the sub-optimal performance of semantic matching. Driven by the recent success of convolutional neural networks (CNNs) several learning-based annuaches e.g. [3,7,13,21,5] have been proposed for addressing the problem of semantic matching. Whil

Deep Semantic Matching with Foreground Detection and Cycle-Consistency

Transitivity Loss. The idea of forward-backward consistency between a pair of images can be further extended to the transitivity consistency across multipl images. Considering the case of three images Is, In, and Ic, we estimate three geometric transformations T_{AB} , T_{BC} , and T_{AC} . Transitivity consistency in this case states that the coordinate transformation from I_A to I_C should be path invariant Namely for any coordinate $\mathbf{n} \in \mathcal{P}_{+}$ the property $T_{n,c}(T_{+n}(\mathbf{n})) \geq$ $T_{AC}(\mathbf{p})$, holds. The transitivity loss is then expressed

 $\mathcal{L}_{T}(I_{A}, I_{B}, I_{C}; \mathcal{F}, \mathcal{G}) = \sum ||T_{BC}(T_{AB}(\mathbf{p})) - T_{AC}(\mathbf{p})||.$ (7)

3.5 Network Selection and Initialization

We try several CNN,based architectures to serve as the feature extractor F we try several CAN-based architecture to serve is the number extractor F, and finally adopt the semantic matching network proposed in [22] due to its impressive results for image alignment. The network employs the ResNet-101 [9] model. The extracted features are these concreted by laws count-29. The transformation predictor G we select is the same as that used in [21]. It contai convolutional layers followed by a fully connected layer to regress the parag ters. The transformation predictor \mathcal{G} is a cascade of two modules predicting an affine transformation and a thin plate spline (TPS) transformation. Given an image pair, the model first estimates an affine transformation with 6 degrees of freedom to obtain a rough alignment. The model then performs a second-stage geometric estimation based on the roughly aligned image pair to predict TI in termination alignment refinement. Similar to [21], we use a uniform 3 > control points for TPS, which corresponds to $3 \times 3 \times 2 = 18$ degrees reedom. We initialize the feature extractor F and the transformation pred G from the narometers pre-trained in [22] and fine-tune F and G by using the proposed objective function

4 Experimental Results

Experiments are conducted in this section. We will describe the implementation details and the experimental setting, evaluate the proposed approach, con it with the state-of-the-art, and analyze the results. More results are proin the supplementary material. The source code and the pre-trained models w be made available to the nublic

4.1 Implementation Details

We implement our model using PyTorch. The training and validation data are both obtained from the PF-PASCAL dataset [6]. All images are resized to the resolution 240×240 . We perform data augmentation by horizontal flipping, ran-dom cropping the input images, and swapping the order of images in the image pair. Our model is trained with ADAM optimizer [15] with an initial learning rate of 5 x 10⁻⁸. Our model is learned with the forward-backward consister rate of 5 × 10⁻¹. Our month is instruct with the low-word-constant consistence loss and transitivity loss first to obtain a good initialization and then is fine tuned with the full objective function. For transitivity loss, the input triplets are

2 V.C. Chen P.-H. Hunne L.-V. Yu. L.B. Hunne M.-H. Yune V.-V. Lin

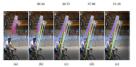


Fig. 1: (a) An image pair for semantic matching. (b)~(c) The matching results, Fig. 1: (a) An image pair or semantic matching. (b)~(c) The matching results, together with their matching errors shown above, generated by different ap-proaches, including (b) the baseline model, (c) the baseline model with fore-ground detection, (d) the baseline model with cycle-consistency checking, and (e) the proposed approach.

promising results have been shown in many of the cases, these approaches still suffer from the following limitations. The methods in [3, 7, 13, 21] require a vast amount of supervised data for training. Collecting a large-scale and diverse data. however, is expensive and labor-intensive. While weakly supervised methods such as [22] have been recently proposed to relax the issue, these approaches implicitly enforce the background features from both images to be similar. Thus, they still suffer from the unfavorable effect of background clutter.

In this work, we address the above-mentioned limitations by presenting an end-to-end trainable network for weakly supervised semantic matching, when training images covering objects of the same category are given without furthe mamual annotations (e.g., keypoint correspondences). To alleviate the negative effect of background chatter, we integrate a foreground detection module into our semantic matching network. In this way, the effect of background clutter can be mitigated by such-fing background matching. The monopoil notwork then focuses on learning the features and geometric models for better matching the detected foreground, resulting in sizable improved performance. To address the matching difficulties caused by complex image appearance and large intraclass variations, we promose to narrow down the matching space by filtering on nors, we propose to narrow down the matching space by intering ou ences with geometric inconsistency. To this end, we exploit the prop erty that correct correspondences should be cycle-consistent meaning that when natching a particular point from one image to the other and then performing reverse manning, we should arrive at the same spot. We further extend this idea splot transitivity consistency across multiple images. Fig. 1 shows an image pair and the matching results generated by using

the baseline model, the baseline model with foreground detection, the baselin model with cycle-consistency checking, and the proposed method (i.e., the base

randomly selected within a mini-batch. We sample $10 \times 10 = 100$ spatial coordi

nates for computing the forward-backward consistency loss and the transitiv

1080 graphics card.

benchmark datasets.

4.2 Evaluation Metric and Datasets

object bounding box on the image, respectively.

r comply to 0.05_0.1_and 0.15_respectively

by setting 7 to 0.05.

image-level supervision

loss. The training process takes about 2 hours on a single NVIDIA GeForce GTX

We conduct the evaluation on the PF-PASCAL [6], PF-WILLOW [6], and TSS [27]

Evaluation Metric. We evaluate the performance of the proposed method on

a semantic correspondence task. To assess the performance or an evolution of the performance of correct keypoints (PCK) metric [33] which measures the percentage

of keypoints whose reproduction errors are below the given threshold. The re-projection error is the Euclidean distance $d(\phi(\mathbf{p}), \mathbf{p}^*)$ between the locations of the warped keypoint $\phi(\mathbf{p})$ and the ground truth keypoint \mathbf{p}^* . The threshold is

defined as r.may(h w) where h and w are the height and width of the annotated

PF-PASCAL [6]. The PF-PASCAL dataset is selected from the PASCAL

2011 keypoint annotations [2] containing 1,351 semantically related image pairs from 20 object categories. For images of a category, they contain different ob-ject instances of that category with similar poses but different appearances. In

addition, the presence of background clutter makes it a challenging dataset or

semantic matching. We divide the dataset into 735 pairs for training, 308 pair

for validation, and 308 pairs for testing. Morally annotated correspondence are provided for each image pairs. However, under the weakly supervised set

ting, the keypoint annotations are only used for evaluation. We compute the PCK for each object category with τ equals to 0.1.

PF-WILLOW [6]. The PF-WILLOW dataset is composed of 100 images with

900 image pairs divided into four semantically related subsets; car, duck, mo-

torbike, and wine bottle. Each subset contains images with large intra-class variations and background clutters. For each image, there are 10 keypoint anno-

tations. We follow [7] and compute the PCK at three different thresholds with

TSS [27]. The TSS dataset comprises 400 semantically related image pairs divided into three groups, including FG3DCar, JODS, and PASCAL. FG3DCar

contains 195 image pairs of automobiles. JODS is composed of 81 image pairs of airplanes, cars, and horses. 124 image pairs of trains, cars, bases, blies, and motorbiles form the group of PASCAL. Ground truth flows for each image pair

are provided. Following [27], we densely compute the PCK over foreground object

In the following, we compare the performance of the proposed method with

the state-of-the-art approaches. Note that many of the existing methods require

manually annotated correspondences while our model can be trained using only

4.3 Experimental Results on the PF-PASCAL Dataset

Dom Semantic Matching with Forum and Detection and Carlo Consistence 3

line model with both foreground detection and cycle-consistency checking) Th numbers on the ton of Fig. 1(h), (e) are the errors measured by the average di numbers on the top of Fig. 1(b) (c) are the errors measured by the average of tance between the predicted keypoints and their corresponding ground truth It can be observed in Fig. 1(c) that the unfavorable effect of background clu ter is alleviated since most correspondences in the background are remove ter is anevanet since must correspondences in the background are removed. Fig. 1(d) shows that exploiting cycle-consistency constraints helps exclude am-biguous matching, resulting in significant reduction of matching errors. Our method simultaneously integrates foreground detection and cycle-consistency checking into semantic matching. As shown in Fig. 1(e), the quality of matching measured by the matching error is further enhanced

ing, inclusion by the matching error, is nurther emanded. The main contributions of this work are summarized as follows. First, we present an end-to-end trainable network that integrates foreground detection into semantic matching. With a module for explicit foreground detection, the pronosed network suppresses the unfavorable effect of background clutter. Second, our model implicitly tackles the ambiguity induced by vast matching sp by inferring bi-directional geometric transformations during matching. With these transformations, we evaluately enforce the inferred connetric transformations to be cycle-consistent by introducing the forward-backward consistency loss. In addition, we exploit the property of transitivity consistency and introduce the transitivity loss to further enhance the matching performance. We train our network with the image pairs of the PF-PASCAL dataset [6]. We then evaluate th proposed model on several standard benchmark datasets for semantic matching including the PF-PASCAL [6], PF-WILLOW [6], and TSS [27] datasets. Exter sive comparisons with existing semantic matching algorithms demonstrate that the proposed approach achieves the state-of-the-art performance.

2 Related Work

Semantic matching has been extensively studied in the literature. Here, we review several topics pertinent to our approach.

Semantic Correspondence. Conventional approaches to semantic matching [4,16,29] leverage hand-crafted descriptors such as SIFT or HOG along with ge motion motions models. They easily the humaint compensationers arrays imp metric matching modes. They seek the keypoint correspondences across image by optimizing a given energy function. The SIFT Flow [16] method shares a sim ilar idea with optical flow, which aligns two images in a large corpus, to establis spondeness. It further employs the SIFT descriptor to extract semantic it e.to.fine search algorithm efficient matching. Yang *et al.* [29] adopt DAISY as the descriptor to efficient nerform correspondence field estimation. Kim et al. [12] learn dense correspon dence by proposing the deformable spatial pyramid. Ham et al. [6] introduce th Proposal Flow where the HOG features and object proposals are used as the matching primitives to learn semantic correspondence. With the use of object proposal, the Proposal Flow method is robust to scaling and background clutter. Taniai et al. [27] propose a model based on hierarchical Markov random field where object co-segmentation and dense correspondences are jointly recovere

4 V.C. Chen P.-H. Hunne L.-V. Yu. L.R. Hunne M.-H. Yune V.-V. Lin

However, the aforementioned methods are established based on hand-crafted descriptors, the another limit their concrelication canability

Semantic Correspondence via Deep Learning. Convolutional neural net works have been applied to semantic matching for the superior performance of feature extraction. Choy et al. [3] propose the universal correspondence net much (UCN) along with a componendence contraction loss. Their method adout work (UCN) along with a correspondence contrastive has. There method adopts a convolutional spatial transformer for feature transformations, which makes their method robust to scaling and rotations. In [13], Kim *et al.*, propose a NN-based descriptor called fally convolutional self-similarity (FCSS) and com-Coversion descriptor visits using convolutions sub-summary (PCSS) and com-bine the descriptor with the Proposal Flow framework [6] for image matching. The SCNet developed by Han *et al.* [7] learns a geometrically plausible model for semantic correspondence by incorporating geometric consistency constraints into the loss function. While the methods in [7, 13] employ trainable descriptor or comantic correspondence feature matching is learned at the object, aronage for similarity correspondence, scatter matching is series at the opper-propose level. Consequently, these methods are not end-to-end trainable because a fo-sion step is required to produce the final results. Rocco et al. [21] present a "NN-based architecture for ecometric matching which estimates a naramet Constrained architecture for geometric matching which estimates a parameter geometric model that can be converted to dense pixel correspondence in an end-to-end trainable fashion. Although these methods [3,6,7,13,21] perform better than those based on hand-crafted features, they need supervised data (in term of manually labeled keypoint correspondence) for training. The dependency of manual more restricts the evaluation.

mamma supervision reserves are scattering. Recently, a few CNN-based methods [11, 20, 22, 28] have carried out weakly supervised semantic correspondence. Novotny *et al.* [20] propose the AncherNet supervises semance corresponsessor. Noteasy et al. (any propose to addresses) which learns a set of filters whose response is geometrically consistent across dif-ferent object instances. However, this model is combined with hand-crafted align-ment models, and therefore is not end-to-end trainable. The WarpNet [11] learns fine-grained image matching with small scale and pose variations via aligning obierts areas image through known deformation. James et al. [28] alien a set of press across images intrough known denormation, names et al. [25] angle a set o images by projecting image pixels to a common coordinate system. Althoug their method has been shown to be effective under well-controlled environments it may not generalize well to real-world scenarios. Inspired by the inlier scoring is any integrational with or the best at a standard inductory of the mast straining procedure of RANSAC, Rocco et al. [22] propose an end-to-end trainable align-ment network, which computes dense semantic correspondence while aligning two images. Our proposed method differs from [20, 22] in two ways. First, our method for they takes into account forcercound detection. Our network learns for method nurmer causes into account neeground servection. Our network neuron re-ture embedding to enhance inter-image foreground similarity while alleviating the unfavorable effects caused by complex background. Second, we introduce bi-directional transformations and further leverage cycle consistency for enforcing geometrically consistent predictions.

Foreground Detection. We review a few methods related to foreground de tection. For saliency detection, Xia et al. [31] develop a center surround infernce network that detects callent reviews of an image in an unconcretized way Zhang et al. (36) propose a deep learning framework that leverages intra-salien

Don Semantic Matching with Foreground Detection and Corla-Consistency 5

prior transfer and deep inter-saliency mining to perform co-saliency detection among groups of images. For semantic segmentation, Long et al. [17] develop a fully convolutional network (FCN) which can be trained end-to-end and picels to-pixels to segment objects of interest. The PSPNet [37] by Zhao et al. extends he idea of feature reventid for comunity commentation in a coarce.to, fine fachion He of al [9] introduce a framework that can not only detect chicate but some He et al. [8] introduce a transverse that can not only detect objects but gener-ate high-quality masks for foreground objects. The proposed method leverages foreground detection to identify the regions for matching. Thus, it can work well on images with complex background.

Cycle Consistency, Leveraging cycle consistency to regularize learning has have studied Sundarum et al [26] evaluit forward-backward consistency to bern studied. Sundaram et al. [26] exploit forward-backward consistency to tackle visual tracking. The CycleGAN [41] method proposed by Zhu et al. shares the same idea. It couples the network with an inverse mapping to deal with unpaired image-to-image translation. Meister et al. [18] propose an unsupervised learning framework, called UnFlow, which estimates bi-directional optical flow isourning framework, called Unition, which estimates to-intectional optical flow to explicitly reason about occlusion and make use of the census transform to increase robustness. While the idea of cycle consistency has been widely apniled to various vision tacks, several methods [38, 50] doare the same idea in context of semantic matching. Zhou et al. [40] tackle the problem of matc the context of semantic matching. Zhou et al. [40] tackle the problem of match-ing multiple images by jointly optimizing feature matching and enforcing cycl consistency. The FlowWeb [38] learns image alignment by establishing globally consistent dense correspondences via exploiting cycle consistency constraints However, these methods [38,40] employ hand-engineered descriptors which can not adapt to an arbitrary object category given for matching. Zhou et al. [38] establish dense correspondences by using an additional 3D CAD model to form resc, instance loop between centhetic data and real images. However, their curl as-measure noip requires four images at a time. On the contrary, we develop two isstency loss requires four images at a time. On the contrary, we develop two is functions to enforce cycle consistency and do not need additional data to guide the training. Experimental results demonstrate that by exploiting cycle sistency constraints, the recorded method immenses the performa-

3 Proposed Algorithm

In this section, we first give an overview of our approach with the developed ective function. Then, each loss adopted in the obie Finally the implementation details are recorded

3.1 Framework Overview

Let $\mathcal{D} = \{I_i\}_{i=1}^N$ denote a set of images covering instances of the same object category, where I: is the i^{ch} image and N is the number of images. Our goal is t barryot, water I_{1} as an I -mapping matrix termine the keypoint correspondences be-tween each image pair (I_{A}, I_{B}) in D without involving the object class in advance. Our formulation for semantic matching is weakly-supervised since training our model reening only weak image, level generoicien in the form of training image airs containing common objects. No ground truth correspondences are use

Deep Semantic Matching with Foreground Detection and Cycle-Consistency 13



image to a foreground pixel in the target image will not be penalized. To demonstrate the effect of masked correspondence loss \mathcal{L}_M , we compute the percentage of correctly warped pixels (i.e., pixels in the foreground/background regions the are correctly warped pixels (correctly warped into a foreground/background region). As shown in Fig. 6 our method effectively reduces the errors in matching pixels from foreground to harkemend and vice versa. The immensionent here is immertant in real-world plications but is not reflected in the metric used in the standard datasets





The ablation study shows that all of the proposed components play crucial roles in producing accurate matching results. From Fig. 4, we observe that the proposed method outperforms the best competitor [22] with a significant margin at multiple thresholds.

4.4 Experimental Results on the PF-WILLOW and TSS Datasets

To evaluate the generalization capability, we apply the learned model trained on the PF-PASCAL dataset to test directly on the PF-WILLOW and TSS datasets. Namely, our model is not fine-tuned on each of the two datasets.

Results on the PF-WILLOW Dataset. Table 3 reports the countitativ results for the PF-WILLOW dataset. We compare the performance with several recent methods [15, 71, 32, 122] as well as other methods [16, 25, 29, 34, 35] using hand-rafted features. The results are directly taken from [7] except [21, 22] For [21] and [22] we run their released code to obtain the figures. From Table , we observe that our model achieves the state-of-the-art performance with all three thresholds.

10 Y.-C. Chen, P.-H. Huang, L.-Y. Yu, J.-B. Huang, M.-H. Yang, Y.-Y. Lin

lass P	Table 1: Per
ens hile hard h	Method
10 78.4 (8.4) 10 (6.7 8.4) 14 (8.7 8.4) 14 (8.4) 14 (8.4) 15 (8.4) 15 (8.4)	ROL PELINEN UNE (4 VOL 10.40 Km A (4) VOL 10.40 Km All (4) VOL 10.40 Km All (4)
	Balla III (1995au (199

Performance Evaluation. We compare our approach with the ProposalFlow [6] the UCN [3], different versions of the SCNet [7], the CNNGeo with different fer attre extractors [21], and a weakly supervised approach proposed by Rocco at [22]. Table 1 presents the experimental results for the PF-PASCAL datase Our results show that the proposed approach compares favorably against state of-the-art methods, achieving an overall PCK of 78.0% (outperforming the pre-vious best method [22] by 3.2%). The advantage of incorporating foreground vious best method [22] by 3.2%). The advantage of mcorporating longround detection and enforcing cycle consistency constraints can be observed by com-paring our approach with ResNet-101+CNNGe0(W) [22] since both methods utilize the same feature descriptor and are trained with image, level supervisie

Deep Semantic Matching with Foreground Detection and Cycle-Consistency

F-PASCAL dataset with $\tau =$

7. Fig. 3 presents the warping results and Fig. 5 presents the qualitative results. for the PF-PASCAL dataset. To further highlight the importance of each con ponent of the proposed method, we present an ablation study of our method in



Fig. 3: Warping. We present the qualitative results of image warping that warps the source image to the target image.

Ablation Study. To analyze the importance of each loss function, we con-duct addition over-eliments on the PF-PASCAL [6] dataset. We present the PCI

2 Y.-C. Chen, P.-H. Huang, L.-Y. Yu, J.-B. Huang, M.-H. Yang, Y.-Y. Lin





Fig. 5: Semantic correspondence results on the PF-PASCAL dataset. We present the visualization of the PF-PASCAL dataset where the matche oordinates are linked with color lines.

curves depicting the performance of the ablations under different thresholds. Table 2 presents the mean PCK value of variants of our approach evaluated on the PF-PASCAL dataset with τ equals to 0.1 where "Baseline" represents ResNet-101+CNNGeo(W) [22] and "Baseline + $\mathcal{L}_M + \mathcal{L}_F + \mathcal{L}_T$ " denotes our proposed method. Our results show that both L_F and $L_T = L_F + L_T$ - dimension proposed mance when compared with the Baseline [22]. To demonstrate the effectiveness onsistency property, we visualize three examples in Fig where the red points indicate the key points and the green points represent the reprojected points. The length of the vellow line represents the distance (loss) between the corresponding points. We observe that enforcing cycle consistency property effectively encourages the network to produce consistent predictions However, the performance gain of using only first is modest. We believe that the reason is due to the evaluation protocol of datasets considers only the matchi on the foreground region. Namely, matching a background pixel in the sour



5 Conclusions





6 Y.-C. Chen, P.-H. Huang, L.-Y. Yu, J.-B. Huang, M.-H. Yang, Y.-Y. Lin



Fig. 2: Semantic matching network. Our model is compose tive CNN modules: a feature extractor F for extracting features and a tran formation predictor C for estimating the geometric transformations betwee formation predictor \mathcal{G} for estimating the geometric transformations between given image pair. The model training is driven by three loss functions, includin the masked correspondence loss \mathcal{L}_M , the forward-backward consistency loss \mathcal{L}_f and the transitivity loss \mathcal{L}_T .

To accomplish this task, we present an end-to-end trainable network which is composed of two collaborative CNN modules: the feature extractor F an he transformation predictor G. The former is a CNN model which learns an extracts the features for a given pair of images. The latter is a CNN-based regressor. For an image pair, it estimates the transformation that warps an image so that the warped image can better align the other image. Fig. 2 presents the an collaboration CNN modules in the monored network exhibitation

two collaborative CNN motions in the proposed network architecture. As shown in Fig. 2, the proposed network architecture takes an image pair for semantic matching. For each image pair (I_A, I_B) in D, they are fed into the feature extractor F to extract their feature maps f_A and f_D , respectively. then perform correlation from f_A to f_B to generate the correlation map S_{AA} The other correlation map S_{BA} is symmetrically obtained. Then the transfor nation predictor G estimates the geometric transformation T_{AB} which warps I so that the warped image I_A can align I_B . In the following section, we develo the objective function used to optimize the feature extractor F and the tran formation predictor G. After optimization, the matching between an image pa (I_A, I_B) can be performed via the predicted transformation T_{AB} or T_{BA}

3.2 Objective Function

The overall training objective consists of three loss functions. First, the masked correspondence loss \mathcal{L}_M minimizes the distance between the corresponding features based on the estimated geometric transformations. Unlike existing semanatching methods [21 22] car model predicts foreground mades to sum is maximg memory [21, 22], our mouse preserve overground masses to suppress the effect of background clutter by excluding background matching. Second the forward-backward consistency loss L_F and the transitivity loss L_T enforce the predicted transformations across multiple images to be acometrically play ble and consistent. Both losses regularize the network training. Specifically, th

14 Y.-C. Chen, P.-H. Huang, L.-Y. Yu, J.-B. Huang, M.-H. Yang, Y.-Y. Lin



Fig. 7: Cycle consistency property. We present the visualization that demonstrates the effect of forward, hardward consistency loss where the red neints ind scrates the energy of network-markward consistency loss where the rup points and cate the keypoints while the green points denote the reprojected points. Yello line represents the distance (loss) between the linked points.

dataset. The quantitative results are presented in Table 4. It can be observed that the proposed method achieves the state-of-the-art performance on two of the three groups of the TSS dataset: FG3DCar and JODS. Our results are slight were than that in [32] in the PASCAL group. However, the method in [32] as uses additional images from the PASCAL VOC 2007 dataset. We report the results for completeness. Under the same experimental settings, the propose method performs favorably against existing approaches.

We have addressed the problem of semantic matching by presenting a weaklysupervised and end-to-end trainable network. The core technical novelty lies i he explicit modeling of a foreground detection module to suppress the effect of background clutter and exploiting the cycle consistence consistence in the text of background clutter and exploiting the cycle consistency constraints so that the predicted geometric transformations are geometrically plausible and consistent across multiple images. The network training requires only training image nai with image-level supervision and thus significantly alleviates the cost of con structing and labeling large-scale training datasets. Experimental results. Among

Deep Semantic Matching with Foreground Detection and Carlo-Consistency 7

training objective is defined by

$$\mathcal{L} = \mathcal{L}_{M} + \lambda_{F} \cdot \mathcal{L}_{F} + \lambda_{T} \cdot \mathcal{L}_{T}$$

where λ_{-} and λ_{-} are hyper-parameters used to control the relative important where x_F and x_T are hyper-parameters used to control the relative importance of the respective loss functions. The details of each loss function are described in the following

3.3 Masked Correspondence Loss

To reduce the effort of background clutter and enforce only foreground regions t To reduce the energy of many context and embedded many more only neighbors to be similar, our model minimizes the masked correspondence loss. Given an im-age pair (I_A, I_B) , the feature extractor F extracts their respective feature maps $I_A \in \mathbb{R}^{k_1 \times w_M \times d}$ and $I_B \in \mathbb{R}^{k_M \times w_M \times d}$, where d is the number of channels. We conrelate f_A with f_B to generate the correlation map $S_{AB} \in \mathbb{R}^{h_A \times u_A \times h_B \times u_B}$. Each element $S_{AB}(i, j, s, t) = S_{AB}(\mathbf{p}, \mathbf{q})$ records the normalized inner product between contain $s_{AB}(r_i), s, s_i = S_{AB}(p, q)$ recover in mermion interproduct context the feature vectors stored at two spatial boottoms $[r_i, [r_i]^{(1)}$ in f_{AB} . The other correlation may $S_{AB} \in \mathbb{R}^{h \times x_{A}, s_{A}, s_{A}}$ can be yielded symmetrically. The correlation maps S_{AB} is subquoted to a tailed-dimensional tensor with dimension h_{A}, w_{A} , and $(h_{B} \times w_{B})$, i.e., $S_{AB} \in \mathbb{R}^{h \times x_{A}, s_{A}, s_{A}}$ and $(h_{B} \times w_{B})$ dimensional tensor is can be interpreted as a dame $A_{A} \times x_{A}$ and with $(h_{B} \times w_{A})$ dimensional boot features. The reshape operation is applied to S_{DA} as well. With the reshape S_{ab} we use the transformation readictor G [21] S_{AB} , we use the transformation predictory y[21] to isolutate a geometric trans-formation T_{AB} which warps I_A to \tilde{I}_A so that \tilde{I}_A aligns well to I_B . Since the correlation map $S_{AB}(\mathbf{p}, \mathbf{q})$ records the normalized inner produc

serveen two feature vectors located at p in f_A and q in f_B . Our model estimates the foreground mask $M_A \in \mathbb{R}^{h \times \times \times h}$ by

$M_A(\mathbf{p}) = \max(S_{AB}(\mathbf{p}, \mathbf{q})).$

Note that both the correlation maps S_{AB} and S_{BA} are compiled through a rectified linear unit (ReLU) to eliminate negative matching values in advance Therefore, the value of the estimated foreground masks at each pixel will be bounded between 0 and 1. Intuitively, $M_A(\mathbf{p})$ has a low value (i.e., location \mathbf{p} is likely to belong to background) if none of the feature vectors in f_B matches well with $f_A(\mathbf{p})$. Likewise, M_B can be obtained.

with $f_A(\mathbf{p})$. Linewest, $A\mathbf{p}$ can be obtained. With the estimated geometric transformation T_{AB} , we can identify and filter out geometrically inconsistent correspondences. Consider a correspondence ($\mathbf{p} \in \mathcal{P}_A, \mathbf{q} \in \mathcal{P}_B$), where \mathcal{P}_A and \mathcal{P}_B are the sets of all spatial coordinates of f_A $r_A, A_q \in r_{BJ}$, succes r_A and r_{BJ} are via fracto on a spanne constants to j_A and f_B , respectively. The distance $[T_{AB}(\mathbf{p}) - \mathbf{q}]$ represents the projection error of this correspondence with respect to transformation T_{AB} . Following [22], we introduce a correspondence mask m_A to determine if the correspondences are geometrically consistent with transformation T_{AB} . Specifically, m_A is defined by

 $m_A(\mathbf{p}, \mathbf{q}) = \begin{cases} 1, & \text{if } \|T_{AB}(\mathbf{p}) - \mathbf{q}\| \le \varphi, \\ 0, & \text{otherwise.} \end{cases}$ for $\mathbf{p} \in \mathcal{P}_A$ and $\mathbf{q} \in \mathcal{P}_B$. (3)

Deep Semantic Matching with Foreground Detection and Cycle-Consistency 1

strate that our approach performs favorably against existing semantic matching algorithms on several standard benchmarks. Moving forward, we believe that th semantic matching network can be further integrated to other computer vision sks, e.g., supporting 3D semantic object reconstruction and fine-grained visual recognitio

References

- 1. H. Bay, T. Tuytelaars, and L. Van Gool. Surf: Speeded up robust features. In
- 2. L. Bourdev and J. Malik. Poselets: Body part detectors trained using 3d human pose annotations. In *ICCV*, 2009. 5. C. B. Choy, J. Gwak, S. Savarese, and M. Chandraker. Universal correspondence
- network. In NIPS, 2016. 4. N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In
- CUDP 200 CVPR, 2005.
 5. P. Fischer, A. Dosovitskiv, E. Br. P. Häusser, C. Hazzbas, V. Golkov, P. Van der
- P. Fucher, A. Dosovitsky, E. B., P. Hauser, C. Hazrbay, V. Gollov, P. Van der Smagt, D. Cremers, and T. Brox. Flownet: Learning optical flow with convolutional networks. In *ICCV*, 2015.
 B. Ham, M. Cho, C. Schmid, and J. Ponce. Proposal flow: Semantic correspon-
- dences from object proposals. TPAMI, 2017.
 7. K. Han, R. S. Rezende, B. Ham, K.-Y. K. Wong, M. Cho, C. Schmid, and J. Ponce.
- Senset: Learning semantic correspondence. In: PCCV, 2017.
 K. He, G. Gisionari, P. Dellár, and R. Gieblick. Mask r-ran. In *ICCV*, 2017.
 K. He, X. Zhang, S. Ron, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- In CVPR, 2016.

 E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox. Florenet 2.0:
- E. Bg, N. Mayer, T. Sakin, M. Keuper, A. Dosovitskiy, and T. Beur. Florent J.R. Evolution of optical Bore estimation with doop networks. In *UVPR*, 2017.
 A. Kamazera, D. W. Jacobs, and M. Chandraker. Warpart: Weakly supervised matching for single-view reconstruction. In *CVPR*, 2018.
 J. Kim, C. Liu, F. Sha, and K. Grauman. Deformable spatial pyramid matching
- for fast dense correspondences. In CVPR, 2013. 13. S. Kim, D. Min, B. Ham, S. Jeon, S. Lin, and K. Sohn. Fosc Fully convolutional
- 14. S. Kim, D. Min, S. Lin, and K. Sohn. Dctm: Discrete-continuous transformation
- S. Kim, D. Min, S. Lin, and K. Sohn. Detin: Discrete-continuous transformation matching for semantic flow. In CVPR, 2017.
 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. arXiv,
- 16. C. Liu, J. Yuen, and A. Torralba. Sift flow: Dense correspondence across scenes
- and its applications. *TPAMI*, 2011.
 17. J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
 18. S. Meister, J. Hur, and S. Roth. Unflow: Unsupervised learning of optical flow

- S. Moister, J. Hur, and S. Roth. Unforc: Unsupervised learning of optical flow with a hidrectional cennes hose. In AAAI, 2018.
 A. Mustafa and A. Hibton. Semantically coherent co-segmentation and reconstruc-tion of dynamic scenes. In: CoVPR, 2017.
 D. Novolay, D. Larden, and A. Vedabili. Ancherenst: A weakly supervised network to learn generative-semantic inclusion for semantic matching. In: CVPR, 2017.

8 V.C. Chen, P.H. Hunne, L.V. Vu, L.B. Hunne, M.H. Vane, V.V. Lin,

where φ is the predefined threshold. Note that the correspondence mask m_A is introduced to remove the correspondences that are mometrically inconsistent introducen to remove the correspondences that are geometrically inconsistent with the transformation T_{AB} , namely those with the projection errors larger than φ . Empirically, we set φ to 1 in our experiments. Given the geometric transformation T_{AB} and correspondence mask m_A , the likelihood of each spatial location $\mathbf{p} \in \mathcal{P}_A$ being matched is computed by

> $s_A(\mathbf{p}) = \sum_{m_A(\mathbf{p}, \mathbf{q})} \cdot S_{AB}(\mathbf{p}, \mathbf{q}).$ (4)

To suppress the effect of background clutter, we incorporate the estimated foreground masks to focus on matching foreground regions. The masked corre pondence loss is defined by

$$\mathcal{L}_{M}(I_{A}, I_{B}; \mathcal{F}, \mathcal{G}) = -\left(\sum_{\mathbf{p} \in \mathcal{P}_{A}} s_{A}(\mathbf{p}) \cdot M_{A}(\mathbf{p}) + \sum_{\mathbf{q} \in \mathcal{P}_{B}} s_{B}(\mathbf{q}) \cdot M_{B}(\mathbf{q})\right).$$
 (5)

Note that the negative sign in (5) is used in the objective function, since maximizing the matching score is equivalent to minimizing the loss L

3.4 Cycle Consistency

For a pair of images I_A and I_B , the transformation predictor G estimates a geometric transformation T_{AB} , which transforms pixel coordinates from I_A to \tilde{I}_B . However, the large capacity of \mathcal{G} often leads to a circumstance where various transformations can warp I_A to \tilde{I}_A which aligns I_B very well. This implies that using the masked correspondence loss alone is not sufficient to reliably train \mathcal{G} in the weakly supervised setting since there are no supervised correspondences to the weaky supervised setting since there are no supervised correspondences to constrain the transformations. We address this issue by simultaneously estimat-ing T_{AB} and T_{BA} and enforce the predicted transformations to be geometrically plausible and consistent across multiple images. It greatly reduces the feasible space of transformations and can serve as a regularization term in training the space of transformation medictor G. To this end, we develop two loss functions whe vcle consistency constraints are done in conjunction with the proposed methor such that the model is end to end trainable. The developed loss functions are

Forward-Backward Consistency Loss. Consider the correlation maps S_{AB} and Sn + concrated from images L+ and Ln. The forward consistency states th and S_{BA} generated from images I_A and I_B . The orward consistency states that property $T_{BA}(T_{AB}(\mathbf{p})) \approx \mathbf{p}$ holds for any $\mathbf{p} \in \mathcal{P}_A$. By the same token, the backward consistency means $T_{AB}(T_{BA}(\mathbf{q})) \approx \mathbf{q}$ for any $\mathbf{q} \in \mathcal{P}_B$. The resultant forward-backward consistency loss is then defined by

 $\mathcal{L}_{F}(I_{A}, I_{B}; \mathcal{F}, \mathcal{G}) = \sum ||T_{BA}(T_{AB}(\mathbf{p})) - \mathbf{p}|| + \sum ||T_{AB}(T_{BA}(\mathbf{q})) - \mathbf{q}||, (6)$

where $||T_{BA}(T_{AB}(\mathbf{p})) - \mathbf{p}||$ is the reprojection error between coordinate \mathbf{p} and the reprojected coordinate $T_{BA}(T_{AB}(\mathbf{p}))$.

16 Y.-C. Chen, P.-H. Huang, L.-Y. Yu, J.-B. Huang, M.-H. Yang, Y.-Y. Lin

- I. Rocco, R. Arandylović, and J. Sivic. Convolutional neural network architecture for geometric matching. In *CVPR*, 2017.
 I. Rocco, R. Arandjelović, and J. Sivic. End-to-end weakly-supervised semantic

- Limm, B., Kampiloev, and J. Shor. Tarkiv-out walks-supervised smartin-slignment. In: (NYP 2018).
 M. Dafanaria, M. Jadar, J. Yangwa, E. C. N. Xampovine i joint diper discovery in the strength strength strength strength strength strength strength strength D. Shatasima, B. Shalik, A. Kuranov, M. Samo, and K. Shatasi, A. Kuranov, S. Shatasima, T. Bas, and K. Kotare, Namo yain charataria of the lange-scale image reception. arXiv:2011. Control Name productional particular bar strength strength strength strength strength strength strength strength image receptions. arXiv:2011. Strength St

- wide-baseline streee. *TPAMI*, 2010.
 30. Z.-F. Wang and Z.-G. Zheng. A region based stereo matching algorithm using cooperative optimization. In *CVPR*, 2008.
 31. C. Xia, F. Qi, and G. Shi. Bottom-up visual salency estimation with deep automooder-based sparse reconstruction. *TNNLS*, 2016.
- automoder-based sparse reconstruction. TNNLS, 2016.
 32. F. Yang, X. Li, H. Cheng, J. Li, and L. Chen. Object-aware dense semantic correspondence. In *CVPR* 2017.
- Correspondence. In Corra, 2017.
 Y. Yang and D. Ramanan. Articulated human detection with flexible mixtures of TOLING COMP.
- parts. TPAMI, 2013.
 34. K. M. Yi, E. Trulls, V. Lepetit, and P. Fua. Lift: Learned invariant feature trans-
- form. In ECCV, 2016.
 5.8. Zagorsylo and N. Komodukis. Learning to compare image patches via convo-lutional neural networks. In CVPR, 2015.
 6. D. Zhang, J. Han, J. Han, and L. Shao. Conditency detection based on intravalismcy
- prior transfer and deep intersaliency mining. TNNLS, 2016.
 37. H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jin. Pyramid scene parsing network. In
- CVPR, 2017.
 38. T. Zhou, Y. Jae Lee, S. X. Yu, and A. A. Efros. Flowweb: Joint image set alignment

- T. Zhen, Y. Jiao, Lin, S. X. Yu, and A. A. Efres. Forwards. Joint mapping halpmanet 30. T. Zhon, Y. MacLawid, M. M. Ling, Y. Hang, and A. A. Kimo. Learning dense correspondence via *Meganized cyclic consistency*. In: CVPR, 2016.
 X. Zhon, M. K. Kanalini, M. Mohimaga anticipa via fast downarding with and a structure of the struct