

Robust Image Matching By Dynamic Feature Selection

Hao Huang¹²³

hh1811@nyu.edu

Jianchun Chen¹²

jc7009@nyu.edu

Xiang Li¹²³

xl1845@nyu.edu

Linging Wang¹²³

lingjing.wang@nyu.edu

Yi Fang*¹²³

yfang@nyu.edu

¹ NYU Multimedia and Visual
Computing Lab, USA

² New York University, USA

³ New York University Abu Dhabi, UAE

Abstract

Estimating dense correspondences between images is a long-standing image understanding task. Recent works introduce convolutional neural networks (CNNs) to extract high-level feature maps and find correspondences through feature matching. However, high-level feature maps are in low spatial resolution and therefore insufficient to provide accurate and fine-grained features to distinguish intra-class variations for correspondence matching. To address this problem, we generate robust features by dynamically selecting features at different scales. To resolve two critical issues in feature selection, *i.e.*, *how many* and *which* scales of features to be selected, we frame the feature selection process as a sequential Markov decision-making process (MDP) and introduce an optimal selection strategy using reinforcement learning (RL). We define an RL environment for image matching in which each individual action either requires new features or terminates the selection episode by referring a matching score. Deep neural networks are incorporated into our method and trained for decision making. Experimental results show that our method achieves comparable/superior performance with state-of-the-art methods on three benchmarks, demonstrating the effectiveness of our feature selection strategy.

1 Introduction

Image matching is one of the fundamental image understanding problems and serves as a building block for various applications, *i.e.*, object recognition [25], motion tracking [24], and 3D construction [1]. Image matching techniques can be further extended or adapted to other application fields including remote sensing [23], protein analysis [28] and medical imaging [2]. Existing methods tackle image matching problem by fitting a geometric transformation between the correspondences of image pairs. Previous approaches usually consist of two steps: 1) a group of hand-crafted image descriptors, *e.g.*, HOG, SIFT, or SURF,

*: indicates corresponding author.

© 2020. The copyright of this document resides with its authors.

It may be distributed unchanged freely in print or electronic forms.

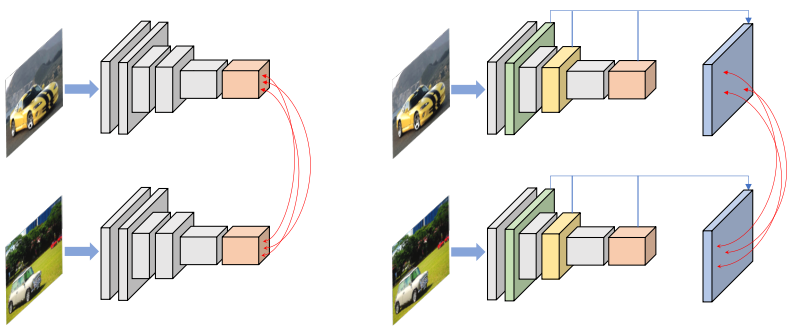


Figure 1: Comparison between image matching methods that only use high-level features (left) and ours that select and integrate features at multiple scales (right).

are pre-computed to obtain pixel-level descriptions; 2) feature matching algorithms, *e.g.*, RANSAC or Hough transformation, are then applied in an iterative manner to determine a proper geometric transformation to match pixels using their associated features. However, hand-crafted image descriptors are vulnerable to intra-class variations, *e.g.*, different texture, lighting, and material. Recently, learning-based approaches [14, 56] alleviate this problem by extracting image features through CNNs and treat image feature maps as dense image descriptors. In general, learning-based methods can be divided into two categories. Some works [16, 56] focus on predicting parameters of a global rigid or non-rigid transformation, *e.g.*, an affine or a thin-plate spline (TPS) transformation. Others [14, 57] formulate this task as a local region matching process and directly pair local regions in source images to the matched regions in target images without transformation regression.

These methods, however, generate dense correspondences only based on high-level features, *i.e.*, the outputs of the last and/or penultimate convolutional layers [56, 57]. As features produced by different levels of CNN layers contain information varying from low-level texture patterns to high-level semantic concepts, these methods fail to fully exploit multi-level image abstractions to obtain reliable matches that are robust to intra-class variations and background noise. The comparison between image matching methods that only use high-level features and multi-scale features is shown in Figure 1. However, there are two pivotal issues in utilizing different levels of features: *how many* and *which* levels of features are apt for matching. To address the issues, the previous method [29] applies beam search to select features from low levels to high levels sequentially without jumping. However, beam search is a heuristic searching strategy, and the search space is a subspace of the full solution space. If multiple optimal solutions exist, beam search may fail to find the one satisfying other requirements, *i.e.*, a minimal number of features used for matching. In this paper, we frame the feature selection problem as a sequential Markov decision-making process (MDP) and tackle it using reinforcement learning. Specifically, based on the selected features, each individual action either requires new features or terminates the selection episode by referring a matching score. The learning process is driven by reward functions. Without manually imposed prior knowledge about image pairs, the proposed method can select the optimal collection of features that are suitable for image matching. Compared with beam search, there is no strict selection order in our proposed method, *i.e.*, from low to high levels, leading to a larger search space and a higher possibility to find the optimal solution. We test the proposed method on three public datasets to demonstrate the effectiveness of our proposed feature selection strategy for robust image matching. Our paper makes three main contributions:

1. We are the first to cast the feature selection process as an MDP and adopt reinforce-

ment learning to select multiple levels of features for robust image matching.

2. We devise a simple but effective deep neural networks to fuse selected features at multiple levels and make a decision at each step, *i.e.*, either to select a new feature or to stop selection for evaluation.
3. We achieve superior/comparable results on three public benchmarks for image semantic correspondence estimation, demonstrating the effectiveness of our method.

2 Related Work

Local region matching. Methods in this category match two sets of local regions based on feature similarity. Traditional methods apply hand-crafted features with spatial regularization [17] or random sampling [9]. Due to the sparsity and vulnerability of the hand-crafted features, they are incapable of handling images containing complex scenes. Bristow *et al.* [9] use LDA-whitened SIFT descriptors, making correspondences less vulnerable to background clutters. Cho *et al.* [6] introduces an effective voting-based algorithm based on region proposals and HOG features for semantic matching and object discovery. Ham *et al.* [12] extends [9] with a local-offset matching algorithm. Recently, features extracted by CNNs have replaced hand-crafted features. Features produced by CNNs can represent high-level semantics and are robust to appearance and shape variations [8, 9, 12]. Choy *et al.* [7] proposes a similarity metric based on CNN features using a contrastive loss. Rocco *et al.* [22, 58] trains image features and correlation features for correspondence matching in an unsupervised/self-supervised setting. Min *et al.* [29] generates a pixel flow to match local regions via a Hough Voting procedure. Similarly, we also adopt CNN features for image matching.

Global image alignment. Other methods [16, 36] formulate correspondence estimation as a geometric alignment problem. Parametric models are adopted to regress parameters of a global transformation. Specifically, image correlation tensors built on image features are fed into a regression layer/network to predict transformation parameters. Some works [16, 36] utilize synthetic image pairs, while the work [37] explores to match images in a weakly supervised way. In addition, Kim *et al.* [20] introduces a recurrent spatial transformer to apply local affine transformations to images iteratively. Chen *et al.* [5] propose a global transformation independent of affine or TPS assumptions. However, these methods only apply high-level features, *i.e.*, the last and/or penultimate features, generated by CNNs to construct correlation tensors, ignoring the importance of exploiting multi-scale levels of features.

Reinforcement Learning in Vision. Our method is largely inspired by works on leveraging reinforcement learning to solve vision problems. [60] is a classical work that uses RL for the spatial attention policy in image recognition. RL is also adopted for object detection [33], video object segmentation [13], video recognition [40], object tracking in [34], scene completion [15] and point cloud parsing [27]. Our work is the first attempt to explore using RL to select a portion of hierarchical features produced by CNNs in the image matching scenario.

3 Approach

In this section, we detail our reinforcement learning method for image matching: instantiations of three basic concepts (*i.e.*, state, action, reward) in reinforcement learning (Sec. 3.2),

network architecture fusing selected features at multiple abstraction levels (Sec. 3.3) and reinforcement training process (Sec. 3.4). An overview of our method is shown in Figure 2.

3.1 Problem Formulation

Given an input image I , a backbone CNN with N convolutional layers produces a sequence of N features $S = (f_1, f_1, \dots, f_N)$ as intermediate representations, ranging from low level visual patterns to high level object semantics. We need to select a subset of rich and compact features for image matching. Once a collection of features are selected, we adopt the probabilistic Hough matching [24] to match images based on the selected features. A brief description of probabilistic Hough matching can be found in Appendix. We formulate the feature selection as a problem of maximizing the matching score function $f_m : \mathbb{S} \rightarrow R$ over the power set \mathbb{S} of the feature set S :

$$s^* = \operatorname{argmax}_{s \in \mathbb{S} \wedge s \neq \emptyset} f_m(s) . \quad (1)$$

We treat the selection process as a Markov decision process, which is applicable of modeling the discrete sequential decision making process. To reframe MDP in image matching scenario, a state set \mathcal{S} , an action set \mathcal{A} , and a reward function \mathcal{R} are defined as follows.

3.2 Overall Architecture

An overview of the proposed approach is shown in Figure 2. Specifically, it consists of two sub-networks: a backbone CNN network which produces a collection of hierarchical features, and a decision network, *i.e.*, the Q-network which decides how to select features produced by the backbone. The structure of the Q-network is detailed in Sec. 3.3.

We denote the backbone CNN with N convolutional layers as C , with each convolutional layer denoted as C_1, C_2, \dots, C_N . The features produced by the layer C_i ($i \in \{1, 2, \dots, N\}$) is denoted as $f_i \in \mathbb{R}^{c_i \times h_i \times w_i}$, where c_i is channel size and $h_i \times w_i$ is the spatial resolution of feature f_i . Note that for different levels of features, the channel size and spatial size may be different. For a pair of input images, I^s as the source image and I^t as the target image, we denote features produced by layer C_i as f_i^s and f_i^t , respectively. The goal is to find which levels of features should be selected such that the expected matching score could be maximized. We regard the problem as a sequential decision-making problem, where at each step an agent selects a new feature map from the feature set or estimates the matching score. We apply a standard reinforcement learning setting, where each episode corresponds to a match of one image pair from a dataset. Let \mathcal{S} be the state space, \mathcal{A} be the set of actions and r be reward functions, we instantiate \mathcal{S} , \mathcal{A} and \mathcal{R} as follows.

State. State $s = (I^s, I^t, \hat{\mathcal{F}}) \in \mathcal{S}$ at the K -th step consist of an image pair (I^s, I^t) and the selected set of K features $\hat{\mathcal{F}} = \{(f_i^s, f_i^t) | i = 1, 2, \dots, K\}$. The agent receives an observation $o = \hat{\mathcal{F}}$. In other words, the observation is the selected features of source and target images at the current step. This forms a partially observable Markov decision process, as non-selected features are invisible to the agent.

Action. The action set \mathcal{A} is divided into two subsets $\mathcal{A} = \mathcal{A}_f \cup \mathcal{A}_s$, as shown in Figure 2. The subset \mathcal{A}_f contains all actions to select new features, *i.e.*, each action corresponds to a feature. The subset \mathcal{A}_s has only one ‘‘terminate’’ action to end the episode and then compute a matching score. The number of actions in \mathcal{A}_f varies based on the concrete type of the backbone CNN, *e.g.*, $|\mathcal{A}_f| = 34$ for ResNet-101 as it has 34 layers. If the action $a_i \in \mathcal{A}_f$

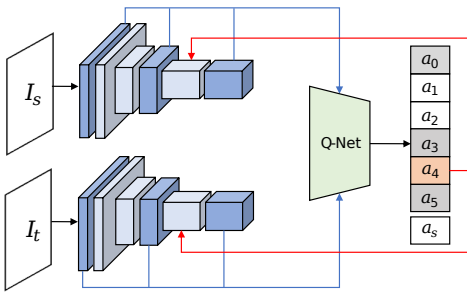


Figure 2: An overview of the proposed approach. Actions in grey color are invalid actions at the current step, as the corresponding features have been selected previously.

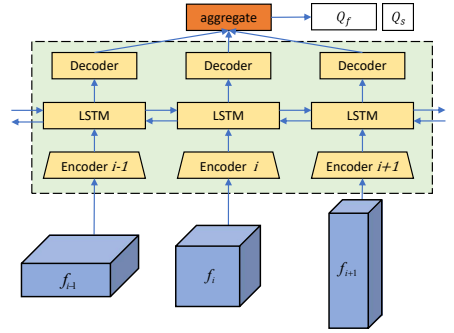


Figure 3: Q-network. It adopts the encoder-LSTM-decoder structure. The input are the current state, *i.e.*, selected features, and the output is the predicted Q-value.

is selected at the current step, the feature tuple (f_i^s, f_i^t) will be added into $\hat{\mathcal{F}}$ and the state changes for the next step.

Reward. Each action in \mathcal{A}_f indicates the corresponding feature to be selected and the agent receives a reward through this action. The reward of each action taken at the i -th step with the action a_i is defined as:

$$\mathcal{R}(a_i) = \begin{cases} -c_i, & \text{if } a_i \in \mathcal{A}_f \\ \beta V_s, & \text{if } a_i \in \mathcal{A}_s \end{cases}, \quad (2)$$

where c_i is a positive cost value and V_s is the final matching score. We define rewards for actions in \mathcal{A}_f to be negative, and the agent gets penalty according to the total number of selected features. Driven by the reward, the agent tends to select the most discriminative features as few as possible to avoid high cost. The value c_i can be either set to different values for different features f_i , or set to the same value for all features. Higher c_i forces the agents to prefer shorter episodes, *i.e.*, less selected features, to match images and vice versa. β is a hyper-parameter to re-scale the matching score V_s into a proper range. Usually, we expect the agent to use minimal features to achieve a matching score as high as possible, especially in cases where matching speed or hardware memory is a consideration.

3.3 Q-network Structure

Suppose the current state for an image pair (I^s, I^t) at step K is $\hat{\mathcal{F}} = \{(f_1^s, f_1^t), \dots, (f_K^s, f_K^t)\}$, the Q-network fuses all selected features at multiple levels to decide the next action, *i.e.*, either to select a new feature or terminate the episode to compute a matching score. The structure of the Q-network is shown in Figure 3. It consists of an encoder-LSTM-decoder structure. Notice that features produced by different convolutional layers are of different channel dimensions and spatial resolutions. Therefore, we adopt a collection of encoders to adjust all selected feature maps to the same channel dimension and spatial resolution. Specifically, suppose the backbone CNN comprises N convolutional layers, the Q-network contains N separate encoders $\{E_i\}_{i=1}^N$ with the same structure but different parameter values. Each encoder processes its corresponding feature independently if the corresponding feature is selected. As features at different levels contain distinct image abstractions, each encoder can adapt its parameters to the corresponding feature level. On the contrary, if a feature

at a certain level is not selected at the current step, the corresponding encoder would skip parameter updating. The structure of all encoders is: *Conv - ReLU - Pooling*, where the pooling layer adaptively reduces the spatial size to 1×1 regardless of the input size.

All pooled features are fed into a bi-directional LSTM to fuse different levels of information. At each time step, the LSTM produces a latent code $L(E_i(f_i))$, regarded as embedded state information for reinforcement learning. The decoder, *i.e.*, a fully connected layer, assimilates the latent code and produces Q-values for action decision. Note that unlike encoder networks, all decoders share the same parameters as the bidirectional LSTM integrates all contextual information at each time step. All Q-values of selected features are aggregated for the final decision. The aggregation operation used in our experiment is element-wise multiplication, and other permutation-invariant operations could also be applied as aggregation.

3.4 Reinforcement Training

Training a deep neural network with reinforcement learning is a challenge. To stabilize and accelerate the training process, we apply the following techniques.

Deep Q-learning. [53] proposed a combination of deep convolutional networks with a variant of Q-learning. As shown in Figure 2, the Q-network works as a “brain” to make decisions. However, training a single Q-network is usually unstable in practice. We adopt a separate *target network* with parameters θ^- introduced in [54] and denote the original Q-network as *evaluation network* with parameters θ . The target network shares the same structure of the evaluation network, but with different parameter values. The parameters θ^- are only updated every certain steps. Instead of directly copying values of θ into θ^- , we apply the method in [24] to slowly update θ^- , defined as $\theta^- := (1 - \rho)\theta^- + \rho\theta$, where ρ is a hyperparameter specifying the change ratio.

Double Q-learning. The Q-value at step t used to guide the learning of the evaluation network is defined as:

$$q_t = r_t + \gamma \max_a Q^{\theta^-}(s_{t+1}, a) , \quad (3)$$

where γ is the discount factor for future rewards. However, as indicated in [59], the max operator in Eq. 3 uses the same value for action selection and evaluation, resulting in a biased estimation of the Q-value. We apply a new formula which decouples the selection and evaluation action for Q-value assignment as:

$$q_t = r_t + \gamma Q^{\theta^-}(s_{t+1}, \operatorname{argmax}_a Q^{\theta}(s_{t+1}, a)) . \quad (4)$$

In Eq. 4, the action is decided by the evaluation network but its value is estimated by the target network.

Dueling Architecture. Following [40], our decoder consists of two parallel streams of fully-connected layers, and the latent code $L(E_i(f_i))$ produced by the LSTM is copied into two the independent streams. The first stream outputs a scalar $V(s)$ for the current state s , and the second stream outputs a vector $A(s, a)$ with $|A|$ dimensions for all actions at the state s . We refer readers to [40] for details about the decomposition of a single Q-function into two separate values. The final Q-value for an action a under the state s is defined as:

$$Q(s, a) = V(s) + A(s, a) - \frac{1}{|A|} \sum_{a'} A(s, a') . \quad (5)$$

By estimating the value of state $V(s)$ explicitly, the training process speeds up and stabilizes. **Retrace.** As we store trajectories into a memory buffer and randomly sample trajectories during each training iteration, there exists a discrepancy between the sampled trajectories and the current policy. To resolve the discrepancy, we modify the Q-value at step t based on the new estimation proposed in [18]:

$$q_t = r_t + \gamma \mathbb{E}_{a \sim \pi(s_t)} [Q(s_{t+1}, a)] + \gamma \bar{\rho}_t [q_{t+1} - Q(s_{t+1}, a_{t+1})] , \quad (6)$$

where $\bar{\rho}_t = \min(\frac{\pi(a_t|s_t)}{\mu(a_t|s_t)}, 1)$ is a truncated importance sampling between behavior policy μ that is used to generate trajectories stored in the memory buffer and the target policy π that the Q-network aims to learn. Importance sampling is a simple way to correct the discrepancy between μ and π when off-policy learning, *i.e.*, using memory buffer, is applied.

4 Experiments

In this section, we present a comprehensive implementation, evaluation, and analysis of our proposed approach on three public real image datasets.

4.1 Datasets and Metric

We compare our model to other image matching methods with both hand-crafted and CNN-based features on three datasets: PF-PASCAL [12], PF-WILLOW [12] and Caltech-101 [10]. **Datasets.** The PF-PASCAL dataset consists of 1,351 semantically related image pairs from 20 object categories. All image pairs are split into training, validation, and test sets in [12]. Manually annotated correspondences for each pair are only provided in validation and test sets. As we need matching scores as the reward of the terminal state, we further split the original validation set into two parts with a ratio of 8 : 2 as the new training and validation sets. Therefore, we use much less images for training. The PF-WILLOW dataset comprises 100 images which are grouped into 900 image pairs. All pairs are divided into four semantically related subsets. For each image, 10 keypoint annotations are provided. The Caltech-101 dataset provides images of 101 object categories with ground-truth object masks, but it does not provide ground-truth keypoint annotations. Following [17], we choose 15 image pairs for each object category and use the corresponding 1,515 image pairs for evaluation.

Metric. We adopt the percentage of correct keypoints (PCK) metric to evaluate our model on PF-PASCAL and PF-WILLOW. PCK measures the percentage of keypoints whose transformation errors are below a given threshold. The transformation error is measured as the Euclidean distance between the location of a warped keypoint and its corresponding ground-truth keypoint. The threshold is defined as $\alpha \max(h, w)$ where h and w are height and width of the object bounding box. For both datasets, we set the threshold $\alpha = 0.1$. Following [19], we adopt label transfer accuracy (LT-ACC) and intersection-over-union (IoU) to evaluate the performance on Caltech-101 dataset. Both metrics measure the number of correctly labeled pixels between ground-truth and warped masks generated by estimated correspondences.

4.2 Implementation and Training

We adopt ResNet-101 pre-trained on ImageNet [9] classification task as the backbone network. In our Q-network, the number of encoders is 34, as ResNet-101 contains 34 convolutional layers. The convolutional layers in all encoders output a 512-dimensional feature

map. The hidden size of the LSTM and the fully-connected layer size in the decoder are also set to be 512. We use features produced by the fully-connected layer of the backbone as initial states to start selection episodes, and these initial features are not applied in computing matching scores. Once the action in \mathcal{A}_s is chosen, we adopt the probabilistic Hough matching proposed in [29] to match images using selected features. During training, we fix the parameters of the backbone network and only train the Q-network. The costs for all actions in \mathcal{A}_f are fixed to 0.4 and β is set to 20. We empirically set the initial learning rate to 0.001. We adopt the Adam optimizer for optimization with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We train our model with a batch size of 16 for 3000 iterations and apply early-stopping. We keep track of the average reward on the training set and decrease the learning rate immediately once the reward fails to increase for ten times. The average reward on the validation set is used for early stopping if the reward fails to increase consecutively for twenty times.

Model		PCK ($\alpha = 0.1$)	
		PF-WILLOW	PF-PASCAL
Hand-crafted	DeepFlow [55]	0.20	0.21
	GMK [10]	0.27	0.27
	DSP [9]	0.29	0.30
	SIFTFlow [26]	0.38	0.33
	ProposalFlow [12]	0.56	0.45
	FCSS + PF-LOM [20]	0.58	0.46
CNN-based	GeoCNN (SS) [56]	0.68	0.68
	A2Net [16]	0.69	0.67
	GeoCNN (WS) [57]	0.71	0.72
	SFNet [2]	0.74	0.79
	HPFlow [29]	0.74	0.85
	Ours	0.75	0.86

Table 1: The average PCK results on PF-WILLOW and the test split of PF-PASCAL dataset with $\alpha = 0.1$. Numbers of the top-1 performance are in bold and the top-2 performance are underlined.

Methods		LT-ACC	IoU
Hand-crafted	DeepFlow [55]	0.74	0.40
	SIFTFlow [26]	0.75	0.48
	GMK [10]	0.77	0.40
	DSP [9]	0.77	0.47
	ProposalFlow [12]	0.78	0.50
	OADSC [43]	0.81	0.55
CNN-based	SCNet-AG [14]	0.79	0.51
	A2Net [16]	0.80	0.57
	FCSS + PF-LOM [20]	0.83	0.52
	GeoCNN (SS) [56]	0.83	0.61
	GeoCNN (WS) [57]	0.85	0.63
	SFNet [2]	0.88	0.67
	HPFlow [29]	0.87	0.63
	Ours	0.87	0.63

Table 2: The average quantitative results on Caltech-101 dataset. Numbers of the top-1 performance are in bold and the top-2 performance are underlined.

4.3 Results

PF-WILLOW & PF-PASCAL. In Table 1, we report the average PCK scores of our method and recent methods that are directly comparable. Note that the performance of all compared methods are taken from [12, 16, 22, 29]. As shown in Table 1, our proposed method achieves state-of-the-art performance [29]. Our model is trained on PF-PASCAL dataset and the selected features are produced by (2, 23, 25, 28) layers, which is around half of all selected layers proposed by [29]¹. For PF-WILLOW dataset, in order to verify the generalization ability of our proposed method, we directly apply the same selected layers for testing without fine-tuning. The results in Table 1 indicates: 1) fusing multiple levels of features is beneficial for image matching; 2) based on our instantiation of three core components of reinforcement learning, *i.e.*, state, action and reward, our method has the capability to learn an optimal selection policy for image matching. As we assign each selected feature a positive cost, it

¹The layers selected for PF-PASCAL dataset using ResNet-101 in [29] are (2, 17, 21, 22, 25, 26, 28).

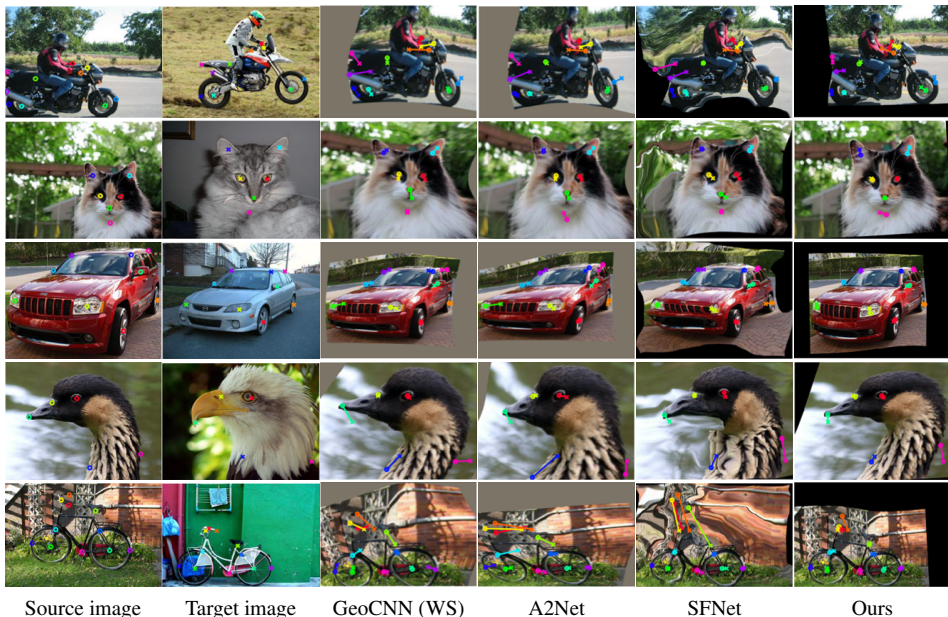


Figure 4: Examples of qualitative results from PF-PASCAL dataset. Keypoints of the source and target images are shown in circles and crosses, respectively. Compared to GeoCNN (WS) [87], A2Net [16] and SFNet [22], our method is more robust to intra-class variations.

drives our method to select features as few as possible, trading off between the final reward and the cost to receive that reward. Another possible reason leading to fewer selected features and higher scores is that there is no restriction for selection order during the policy learning process. On the contrary, beam search in [29] selects features from low levels to high levels sequentially without skipping, resulting in a smaller search space.

Caltech-101. The quantitative results for the Caltech-101 dataset are listed in Table 2. All results except for ours are taken from [12, 16, 22, 29]. Similar to PF-WILLOW dataset, we directly test our method on Caltech-101 using selected layers based on PF-PASCAL without fine-tuning. Our method achieves the same performance as [29]. Notice that the performance of SFNet [22] is better than ours and [29]. One possible explanation is that we select features from the backbone network pre-trained for image classification, while SFNet introduces two trainable convolutional layers to refine image features for image matching. Although we combine both high and low levels of features, the combined features are still not as robust as features dedicatedly trained for image matching. Therefore, incorporating feature learning into our proposed method is left for future study.

Qualitative comparison. Figure 4 visualizes the matching results between keypoints in source and target images on the test split of the PF-PASCAL dataset. We can see that our method is able to select features that are robust to (near) rigid deformations (*e.g.*, cars in the third row), local non-rigid deformations (*e.g.*, the cat’s head in the second row and the upper part of birds in the fourth row), scale changes between objects (*e.g.*, wheels of bicycles in the last row). In particular, the first example clearly demonstrates that our method establishes more discriminative correspondences, avoiding matches for non-target objects. For example, it does not match keypoints on the motorbike to keypoints on the rider between the source and target images, while all other methods mismatch some keypoints.

5 Conclusion

We propose a reinforcement learning method to dynamically select an optimal set of features at multiple levels to match images. Our method regards the collection of selected features as states and simulates an environment to generate rewards to decide actions to select the next feature. The experiments on three public benchmarks demonstrate that our method is able of selecting a small set of features for image matching. We believe applying reinforcement learning to select task-relevant visual features can be applied to various vision tasks.

Appendix

A.1 In this subsection, we briefly summarize the process of establishing keypoint correspondences using probabilistic Hough matching [29]. Given a collection of selected features, firstly, all selected features are spatially resized to the same size as the largest feature and then concatenated along the channel dimension. Given any two positions in source and target image features separately, we compute their visual similarity (cosine distance) and the spatial offset between these two positions. A Hough space is built based on offsets and discretized into different bins. Each pair of positions in source and target features cast a vote, *i.e.*, the value of visual similarity, to their corresponding offset bin. The final visual similarity between any two positions is re-weighted by summing up all votes falling into their bin. The keypoint correspondences are estimated as follows: for a keypoint in a source image, there must exist several receptive fields (RFs, corresponding to positions in image features) covering it. Each position in the source feature has a most visually similar position in the target feature corresponding to an RF in the target image. The displacements of source RFs’ centers to the source keypoint are used as weights to weighted sum up the corresponding target RFs’ centers, which is regarded as the source keypoint’s correspondence.

A.2 In this subsection, we report the experiment on the PF-PASCAL dataset by randomly selecting a collection of layers. We fix the number of layers to be selected, *i.e.*, the value of K is restricted within the range from 2 to 7, and then randomly selected K layers. For each value of K , we test twice and the result is shown in Table 3. Notice that even if the same number of K layers are used for matching, the randomly selected layers yield inferior performance compared with layers selected by the proposed method, as shown in Table 1.

# Layers	Layer Indices	PCK
2	(4, 6)	0.05
	(6, 18)	0.11
3	(17, 22, 31)	0.70
	(14, 27, 30)	0.71
4	(7, 19, 30, 31)	0.10
	(9, 15, 16, 26)	0.59
5	(6, 14, 17, 18, 22)	0.38
	(9, 16, 17, 21, 24)	0.66
6	(5, 7, 14, 27, 32, 33)	0.17
	(9, 15, 22, 23, 25, 33)	0.73
7	(6, 10, 11, 20, 28, 29, 33)	0.25
	(16, 18, 22, 26, 27, 30, 33)	0.70

Table 3: The average PCK results of randomly selected layers on PF-PASCAL dataset.

Acknowledgement

We would like to thank the reviewers for their comments and efforts towards improving our manuscript. The authors appreciate the generous support provided by Inception Institute of Artificial Intelligence in the form of NYUAD Global Ph.D. Student Fellowship.

References

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building rome in a day. *Communications of the ACM*, 54(10):105–112, 2011.
- [2] Michel A Audette, Frank P Ferrie, and Terry M Peters. An algorithmic overview of surface registration techniques for medical imaging. *Medical image analysis*, 4(3): 201–217, 2000.
- [3] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. In *ACM Transactions on Graphics (ToG)*, volume 28, page 24. ACM, 2009.
- [4] Hilton Bristow, Jack Valmadre, and Simon Lucey. Dense semantic correspondence where every pixel is a classifier. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4024–4031, 2015.
- [5] Jianchun Chen, Lingjing Wang, Xiang Li, and Yi Fang. Arbicon-net: Arbitrary continuous geometric transformation networks for image registration. In *Advances in neural information processing systems*, pages 3410–3420, 2019.
- [6] Minsu Cho, Suha Kwak, Cordelia Schmid, and Jean Ponce. Unsupervised object discovery and localization in the wild: Part-based matching with bottom-up region proposals. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1201–1210, 2015.
- [7] Christopher B Choy, JunYoung Gwak, Silvio Savarese, and Manmohan Chandraker. Universal correspondence network. In *Advances in neural information processing systems*, pages 2414–2422, 2016.
- [8] Guoxian Dai, Jin Xie, and Yi Fang. Siamese cnn-bilstm architecture for 3d shape representation learning. In *IJCAI*, pages 670–676, 2018.
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [10] Olivier Duchenne, Armand Joulin, and Jean Ponce. A graph-matching kernel for object categorization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1792–1799. IEEE, 2011.
- [11] Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.
- [12] Bumsu Ham, Minsu Cho, Cordelia Schmid, and Jean Ponce. Proposal flow: Semantic correspondences from object proposals. *IEEE transactions on pattern analysis and machine intelligence*, 40(7):1711–1725, 2017.
- [13] Junwei Han, Le Yang, Dingwen Zhang, Xiaojun Chang, and Xiaodan Liang. Reinforcement cutting-agent learning for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9080–9089, 2018.

- [14] Kai Han, Rafael S Rezende, Bumsub Ham, Kwan-Yee K Wong, Minsu Cho, Cordelia Schmid, and Jean Ponce. Scnet: Learning semantic correspondence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1831–1840, 2017.
- [15] Xiaoguang Han, Zhaoxuan Zhang, Dong Du, Mingdai Yang, Jingming Yu, Pan Pan, Xin Yang, Ligang Liu, Zixiang Xiong, and Shuguang Cui. Deep reinforcement learning of volume-guided progressive view inpainting for 3d point scene completion from a single depth image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 234–243, 2019.
- [16] Paul Hongsuck Seo, Jongmin Lee, Deunsol Jung, Bohyung Han, and Minsu Cho. Attentive semantic alignment with offset-aware correlation kernels. In *European Conference on Computer Vision*, pages 349–364, 2018.
- [17] Junhwa Hur, Hwasup Lim, Changsoo Park, and Sang Chul Ahn. Generalized deformable spatial pyramid: Geometry-preserving dense correspondence estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1392–1400, 2015.
- [18] Jaromír Janisch, Tomáš Pevný, and Viliam Lisý. Classification with costly features using deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3959–3966, 2019.
- [19] Jaechul Kim, Ce Liu, Fei Sha, and Kristen Grauman. Deformable spatial pyramid matching for fast dense correspondences. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2307–2314, 2013.
- [20] Seungryong Kim, Dongbo Min, Bumsub Ham, Sangryul Jeon, Stephen Lin, and Kwanghoon Sohn. Fcss: Fully convolutional self-similarity for dense semantic correspondence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6560–6569, 2017.
- [21] Seungryong Kim, Stephen Lin, SANG RYUL JEON, Dongbo Min, and Kwanghoon Sohn. Recurrent transformer networks for semantic correspondence. In *Advances in neural information processing systems*, pages 6126–6136, 2018.
- [22] Junghyup Lee, Dohyung Kim, Jean Ponce, and Bumsub Ham. Sfnet: Learning object-aware semantic correspondence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2278–2287, 2019.
- [23] Xiang Li, Xiaojing Yao, and Yi Fang. Building-a-nets: robust building extraction from high-resolution remote sensing images with adversarial networks. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 11(10):3680–3687, 2018.
- [24] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations*, 2016. URL <https://arxiv.org/abs/1509.02971>.

- [25] Ce Liu, Jenny Yuen, Antonio Torralba, Josef Sivic, and William T Freeman. Sift flow: Dense correspondence across different scenes. In *European conference on computer vision*, pages 28–42. Springer, 2008.
- [26] Ce Liu, Jenny Yuen, and Antonio Torralba. Sift flow: Dense correspondence across scenes and its applications. *IEEE transactions on pattern analysis and machine intelligence*, 33(5):978–994, 2010.
- [27] Fangyu Liu, Shuaipeng Li, Liqiang Zhang, Chenghu Zhou, Rongtian Ye, Yuebin Wang, and Jiwen Lu. 3dcnn-dqn-rnn: A deep reinforcement learning framework for semantic parsing of large-scale 3d point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5678–5687, 2017.
- [28] Yu-Shen Liu, Yi Fang, and Karthik Ramani. Using least median of squares for structural superposition of flexible proteins. *BMC bioinformatics*, 10(1):29, 2009.
- [29] Juhong Min, Jongmin Lee, Jean Ponce, and Minsu Cho. Hyperpixel flow: Semantic correspondence with multi-layer neural features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3395–3404, 2019.
- [30] Volodymyr Mnih, Nicolas Heess, Alex Graves, et al. Recurrent models of visual attention. In *Advances in neural information processing systems*, pages 2204–2212, 2014.
- [31] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [32] Richard A Newcombe, Steven J Lovegrove, and Andrew J Davison. Dtam: Dense tracking and mapping in real-time. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2320–2327. IEEE, 2011.
- [33] Aleksis Pirinen and Cristian Sminchisescu. Deep reinforcement learning of region proposal networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6945–6954, 2018.
- [34] Liangliang Ren, Jiwen Lu, Zifeng Wang, Qi Tian, and Jie Zhou. Collaborative deep reinforcement learning for multi-object tracking. In *European Conference on Computer Vision*, pages 586–602, 2018.
- [35] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. Deep-matching: Hierarchical deformable dense matching. *International Journal of Computer Vision*, 120(3):300–323, 2016.
- [36] Ignacio Rocco, Relja Arandjelovic, and Josef Sivic. Convolutional neural network architecture for geometric matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6148–6157, 2017.
- [37] Ignacio Rocco, Relja Arandjelović, and Josef Sivic. End-to-end weakly-supervised semantic alignment. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6917–6925, 2018.

- [38] Ignacio Rocco, Mircea Cimpoi, Relja Arandjelović, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Neighbourhood consensus networks. In *Advances in neural information processing systems*, pages 1651–1662, 2018.
- [39] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Thirtieth AAAI conference on artificial intelligence*, 2016.
- [40] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International Conference on Machine Learning*, pages 1995–2003, 2016.
- [41] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6222–6231, 2019.
- [42] Jin Xie, Meng Wang, and Yi Fang. Learned binary spectral shape descriptor for 3d shape correspondence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3309–3317, 2016.
- [43] Fan Yang, Xin Li, Hong Cheng, Jianping Li, and Leiting Chen. Object-aware dense semantic correspondence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2777–2785, 2017.