

CATs++: Boosting Cost Aggregation with Convolutions and Transformers

Seokju Cho*, *Student Member, IEEE*, Sunghwan Hong*, *Student Member, IEEE*,
and Seungryong Kim†, *Member, IEEE*

Abstract—Cost aggregation is a process in image matching tasks that aims to disambiguate the noisy matching scores. Existing methods generally tackle this by hand-crafted or CNN-based methods, which either lack robustness to severe deformations or inherit the limitation of CNNs that fail to discriminate incorrect matches due to limited receptive fields and inadaptability. In this paper, we introduce Cost Aggregation with Transformers (CATs) to tackle this by exploring global consensus among initial correlation map with the help of some architectural designs that allow us to benefit from global receptive fields of self-attention mechanism. To this end, we include appearance affinity modeling, which helps to disambiguate the noisy initial correlation maps. Furthermore, we introduce some techniques, including multi-level aggregation to exploit rich semantics prevalent at different feature levels and swapping self-attention to obtain reciprocal matching scores to act as a regularization. Although CATs can attain competitive performance, it may face some limitations, *i.e.*, high computational costs, which may restrict its applicability only at limited resolution and hurt performance. To overcome this, we propose CATs++, an extension of CATs. Concretely, we introduce early convolutions prior to cost aggregation with a transformer to control the number of tokens and inject some convolutional inductive bias, then propose a novel transformer architecture for both efficient and effective cost aggregation, which results in apparent performance boost and cost reduction. With the reduced costs, we are able to compose our network with a hierarchical structure to process higher-resolution inputs. We show that the proposed method with these integrated outperforms the previous state-of-the-art methods by large margins. Codes and pretrained weights are available at: <https://ku-cvlab.github.io/CATs-PlusPlus-Project-Page/>

Index Terms—Semantic visual correspondence, cost aggregation, efficient transformer

1 INTRODUCTION

ESTABLISHING dense correspondences across semantically similar images can facilitate many Computer Vision applications, including semantic segmentation [1], [2], [3], object detection [4], and image editing [5], [6], [7], [8], [9]. Unlike classical dense correspondence problems such as stereo matching, geometric matching, or optical flow that consider visually similar images taken under the geometrically constrained settings [10], [11], [12], [13], semantic correspondence poses additional challenges from large intra-class appearance and geometric variations [14], [15], [16] caused by the unconstrained settings of the given image pair.

Recent approaches [17], [18], [19], [20], [21], [22], [23], [24], [25], [26], [27], [28] addressed these challenges by carefully designing deep Convolutional Neural Networks (CNNs)-based models analogously to the classical matching pipeline [29], [30], namely feature extraction, cost aggregation, and flow estimation. Several works [21], [22], [25], [27], [31], [32] focused on the feature extraction stage, as it has been demonstrated that the more powerful feature representation the model learns, the more robust matching is obtained [27], [31], [32]. However, solely relying on the matching similarity between features without any prior, *e.g.*, spatial smoothness often suffers from the challenges due to ambiguities generated by repetitive patterns or background clutters [17], [31], [33]. On the other hand, some methods [17], [18], [24], [33], [34], [35] focused on the flow estimation stage either by designing additional CNN as an ad-hoc regressor that predicts the parameters of a single global transformation [17], [18], finding confident matches from correlation maps [33], [36],

or directly feeding the correlation maps into the decoder to infer dense correspondences [24]. However, these methods highly rely on the quality of the initial correlation maps [26].

To address this, the latest methods [19], [21], [23], [28], [37], [38], [39] have focused on the second stage, highlighting the importance of cost aggregation. Since the quality of correlation maps is of prime importance, they proposed to refine the matching scores by formulating the task as Optimal Transport problem [23], [25], re-weighting matching scores by Hough space voting for geometric consistency [21], [22], or utilizing high-dimensional 4D or 6D convolutions to find locally consistent matching points [19], [28], [37], [39], [40]. Although formulated variously, these methods either use hand-crafted techniques that are neither learnable nor robust to severe deformations, nor inherit the limitation of CNNs, which their limited receptive fields make it hard to discriminate incorrect matching points, and they lack an ability to adapt to the input contents due to the inherent characteristic of convolution, *i.e.*, a fixed and shared kernel across all the input pixels.

In this work, we also focus on the cost aggregation stage and propose a novel cost aggregation network to tackle the aforementioned issues. Our network, called Cost Aggregation with Transformers (CATs), is based on transformer [41], [42], which is renowned for its global receptive field and ability to flexibly adapt to consider pairwise interactions among all the input tokens. By considering all the matching scores computed between features of input images globally, our aggregation network explores global consensus and thus refines the ambiguous or noisy matching scores effectively. Specifically, based on the observation that desired correspondence should be aligned at discontinuities with the appearance of images [43], [44], we concatenate an appearance

* Equal contribution

† Corresponding author

• S. Cho, S. Hong, and S. Kim are with Korea University, Seoul, Korea.

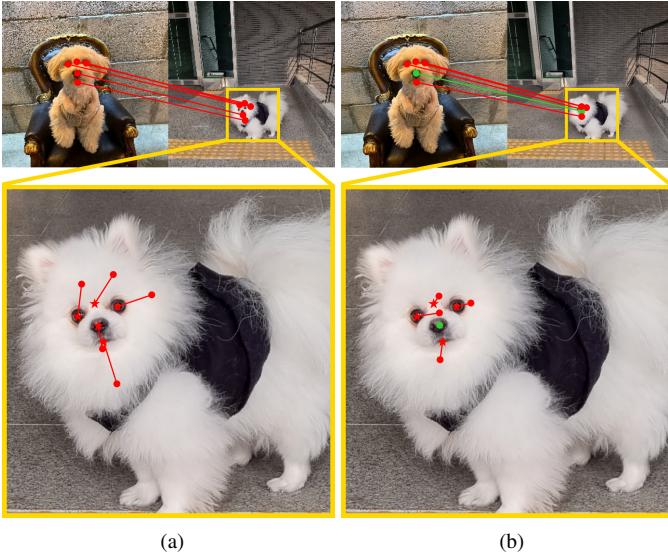


Fig. 1: **Qualitative comparison of (a) CATs and (b) CATs++.** We show that our proposed methods perform surprisingly well even for images in the wild. Note that green and red circles denote the correctly predicted points and incorrectly predicted points, respectively, and the red star denotes the ground-truth keypoint. CATs++ finds highly accurate and fine-grained matching points, even surpassing CATs, thanks to its hierarchical architecture as well as its enhanced cost aggregation components.

embedding to the correlation map, which helps to disambiguate the correlation map within the transformer aggregator. To benefit from hierarchical feature representations, following [22], [24], [33], we use a stack of correlation maps constructed from multi-level features, and propose to effectively aggregate the scores across the multi-level correlation maps. Furthermore, we consider the bidirectional nature of a correlation map to leverage the correlation map from both source and target directions, obtaining reciprocal scores by swapping the pair of dimensions of the correlation map to allow global consensus. In addition to these, we provide residual connections around aggregation networks to ease the learning process.

Although competitive performance can be attained by CATs, it may face some limitations, *i.e.*, high computational costs induced by the use of full attention adopted by standard transformers [41], which restrict its applicability only at limited resolution and result in rather limited performance. To address these, we propose CATs++ to not only alleviate the computational burden but also to enhance the cost aggregation for the improved performance, as exemplified in Fig. 1. Specifically, we introduce early convolutions to reduce the costs by controlling the number of tokens and their dimensions and inject some convolutional inductive bias prior to cost aggregation by transformers, which also leads to apparent performance boost as this allows the model capable of aggregating both local and global interactions. Furthermore, we introduce a novel design for transformer, tailored for cost aggregation, which includes reformulation of standard Query-Key-Value (QKV) projection and Feed-Forward Network (FFN). This is achieved by including appearance affinity at the intermediate process of QKV projection to reduce computational burden and utilizing 4D convolutions rather than the linear projection in FFN to not only strengthen the power to consider locality but also to flexibly control the input dimensions. With the reduced costs

from the reformulation, we compose a hierarchical architecture, allowing the coarser levels to guide the cost aggregation at finer levels. Lastly, unlike CATs, we refrain from arbitrarily selecting different combinations of feature maps [21], [39], [45] for each dataset, which adds extra burden, but rather use all the feature maps as done in [3] to exploit richer semantic representations.

With these combined, we demonstrate our methods on standard benchmarks [14], [15], [16] for semantic correspondence. Experimental results on various benchmarks prove the effectiveness of the proposed models over the latest methods for semantic correspondence, clearly outperforming and attaining state-of-the-art for all the benchmarks. We also provide an extensive ablation study to validate and analyze components. Finally, we show that the proposed method is robust to both domain and class shift, and also works surprisingly well for a few-shot segmentation task as one of the possible applications.

This work is extended from the preliminary version of [45] by (I) leveraging both convolutions and transformers to take the best of two by strengthening the ability to consider locality interactions prior to global consideration of pairwise interactions by transformers and reduce the computational loads, (II) introducing a novel and efficient transformer architecture that can significantly boost the performance and reduce the computational costs, (III) designing a hierarchical architecture that aggregates the cost volume constructed at higher resolution, (IV) setting a new state-of-the-art performance for standard benchmarks [14], [15], [16] in semantic correspondence, making an 8.5%p increase from the best-published results, and (V) providing additional experimental results, ablation studies and visualizations.

2 RELATED WORK

2.1 Semantic Correspondence

Methods for semantic correspondence generally follow the classical matching pipeline [29], [30], including feature extraction, cost aggregation, and flow estimation. Most early efforts [6], [14], [46] leveraged the hand-crafted features which are inherently limited in capturing high-level semantics. Though using deep CNN-based features [17], [18], [31], [33], [34], [35], [47] has become increasingly popular thanks to their invariance to deformations, without a means to refine the matching scores independently computed between the features, the performance would be rather limited. Recently, MMNet [48] attempts to use transformers [41] to refine the features, focusing on the feature extraction stage, with slight changes made to self-attention computation to compute local self-attention over feature maps to aggregate the local feature values, while we take a different approach, aggregating the cost volume and leveraging both convolutions and transformers to not only inject some convolutional inductive bias and reduce computation cost but also to exploit the merits of both, *i.e.*, locality bias of convolution, adaptability and global receptive fields of transformers.

To alleviate the limited performance by the use of matching scores without aggregation, several methods focused on the flow estimation stage. Rocco et al. [17], [18] proposed an end-to-end network to predict global transformation parameters from the matching scores, and their success inspired many variants [8], [34], [35]. Among the subsequent works, RTNs [35] attempted to obtain semantic correspondences through an iterative process of estimating spatial transformations. DGC-Net [20], Semantic-GLU-Net [24], and DMP [49] utilize a CNN-based decoder

to directly find correspondence fields. Recently, PDC-Net [50] proposed a flexible probabilistic model that jointly learns the flow estimation and its uncertainty. Arguably, directly regressing correspondences from the initial matching scores highly relies on the quality of them, which motivates aggregation of matching scores for more reliable correspondences.

Recent numerous methods [19], [21], [22], [23], [25], [27], [28] thus have focused on the cost aggregation stage to refine the initial matching scores. Among hand-crafted methods, SCOT [23] formulates semantic correspondence as an Optimal Transport problem and attempts to solve two issues, namely many to one matching and background matching. HPF [21] first computes appearance matching confidence using hyperpixel features and then uses Regularized Hough Matching (RHM) algorithm for cost aggregation to enforce geometric consistency. DHPF [22], which replaces the feature selection algorithm of HPF [21] with trainable networks, also uses RHM. However, these hand-crafted techniques for refining the matching scores are neither learnable nor robust to severe deformations.

As learning-based approaches, NC-Net [19] utilizes 4D convolution to achieve local neighborhood consensus by finding locally consistent matches, and its variants [37], [39] proposed more cost-efficient methods. PMNC [51] utilizes a classical algorithm, PatchMatch [52], and 4D convolutions to tackle the semantic correspondence task, and proposes a PatchMatch-based optimization strategy to iteratively refine the correspondence field. GO-Cor [26] proposed an aggregation module that directly improves the correlation maps. GSF [38] formulated a pruning module to suppress false positives of correspondences in order to refine the initial correlation maps. CHM [28] goes one step further, proposing a learnable geometric matching algorithm that utilizes 6D convolution. CHMNet [40] extends CHM [28] by introducing an efficient kernel decomposition with center-pivot neighbors. Although outstanding performance with the help of additional data augmentation [45] and multi-level features [21], [22], [45] is achieved, CHMNet [40] yet suffers from limitations of CNN-based architectures for cost aggregation.

2.2 Visual Correspondence Applications

Establishing visual correspondences has been a cornerstone of many Computer Vision applications. To name a few, object tracking [53], [54], video object segmentation (VOS) [55], [56], [57], few-shot segmentation (FSS) [3], [58], visual localization [37], [59], structure-from-motion (SfM) [60], and image retrieval [61], [62] are the tasks that establishing accurate correspondence fields remain as one of the milestones for their goals. More specifically, object tracking [63] and video object segmentation [64] need space-time correspondences between objects of interests for detection or segmentation across different time steps, and FSS [65] requires segmenting objects of unseen classes by finding correspondences given a few annotated support images and a query image. SfM [60] and visual localization [62] are the 3D computer vision tasks to estimate three-dimensional structures by finding correspondences across multiple two-dimensional image sequences. In this paper, we also show how the proposed method is applicable to FSS and visual localization tasks, and evaluate its effectiveness.

2.3 Transformers in Vision

Transformers [41], the *de facto* standard for Natural Language Processing (NLP) tasks, has recently imposed significant impact

on various tasks in Computer Vision fields such as image classification [42], [66], object detection [67], [68], tracking and matching [27], [69]. ViT [42], the first work to propose an end-to-end transformer-based architecture for the image classification task, successfully extended the receptive field, owing to its self-attention nature that can capture the global relationships between features. For those works addressing visual correspondence, LoFTR [27] uses a cross and self-attention module to refine the feature maps conditioned on both input images, and formulate the hand-crafted aggregation layer with dual-softmax [19], [70], and Optimal Transport [25] to infer correspondences. In another work, COTR [71] takes coordinates as input and addresses dense correspondence tasks without the use of a correlation map. Unlike these, for the first time, we propose a transformer-based cost aggregation module.

3 METHODOLOGY

3.1 Motivation and Overview

Let us denote a pair of images as source I_s and target I_t , which represent semantically similar images, and features extracted from I_s and I_t as D_s and D_t , respectively. Here, our goal is to establish a dense correspondence field $F(i)$ between two images that are defined for each pixel i , which warps I_t towards I_s .

Estimating the correspondence with sole reliance on matching similarities between D_s and D_t is often challenged by ambiguous matches due to the repetitive patterns or background clutters [17], [24], [31], [33], [49]. To address this, numerous methods proposed cost aggregation techniques that try to refine the initial matching similarities either by formulating the task as Optimal Transport problem [23], [25], using regularized Hough Matching to re-weight the costs [21], [22] or adopting 4D or 6D convolutions [19], [28], [37], [39], [40]. However, these methods either use hand-crafted techniques that are weak to severe deformations or only consider long-range pairwise interactions implicitly by stacking the series of convolution blocks.

To overcome these, we present transformer-based cost aggregation networks that effectively integrate information present in all matching costs as a whole by considering pairwise interactions explicitly with the help from global receptive fields of attention mechanism, as illustrated in Fig. 2. In the following, we explain feature extraction and cost computation and then explain the proposed transformer aggregator. Finally, we introduce additional techniques to further boost the performance to complete the design of cost aggregation networks with transformers.

3.2 Feature Extraction and Cost Computation

To extract dense feature maps from images, as shown in [20], [21], [22], [23], [24], we leverage multi-level features to capture hierarchical semantic feature representations. To this end, we use multi-level features from different levels of convolutional layers to construct a stack of correlation maps. Specifically, we use CNNs¹ that produce a sequence of L feature maps, and D^l represents a feature map at the l -th level. As done in [21], [23], we use a different combination of multi-level features depending on the dataset trained on, e.g., PF-PASCAL [15] or SPair-71k [16]. Given a sequence of feature maps, we resize all the selected feature maps

¹ Note that we could also use transformer-based feature backbone networks, e.g., ViT [42], which we explore the influence of different feature backbone networks in Section 5.4.

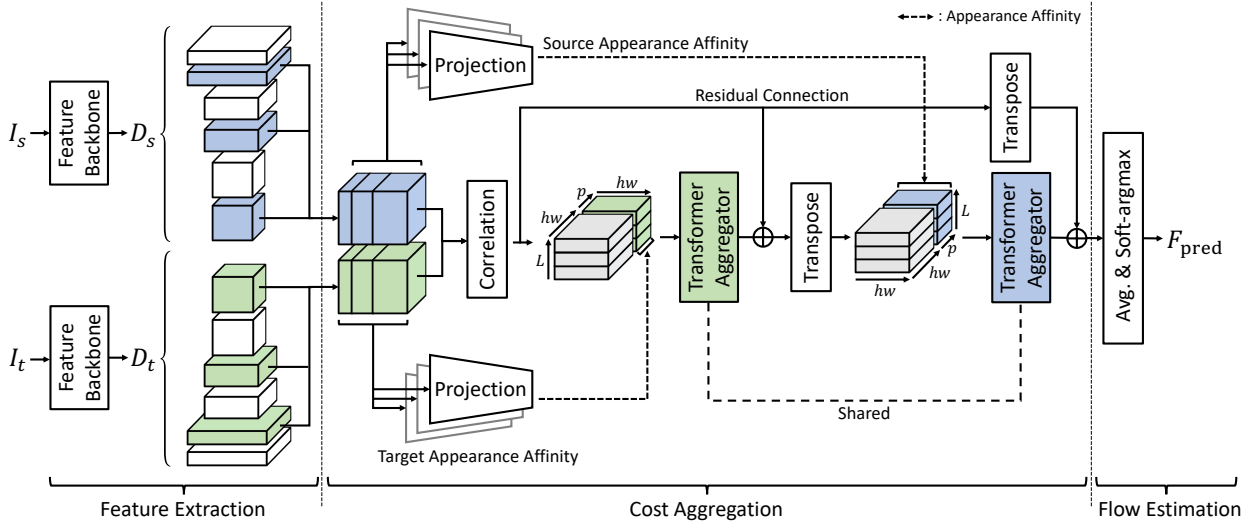


Fig. 2: **Overall network architecture of CATs.** Analogously to the classical matching pipeline, our networks consist of feature extraction, cost aggregation, and flow estimation modules. We first extract multi-level dense features and construct a stack of correlation maps. We then concatenate with embedded features and feed into the transformer-based cost aggregator to obtain a refined correlation map. The flow is then inferred from the refined correlation map.

to $\mathbb{R}^{h \times w \times c}$, with height h , width w , and c channels. Note that we down-sample the features to fixed h and w for efficiency. The resized features first undergo $l-2$ normalization and a correlation map is computed using the inner product between them followed by an activation function such that:

$$C^l(i, j) = \text{ReLU} \left(\frac{D_s^l(i) \cdot D_t^l(j)}{\|D_s^l(i)\| \|D_t^l(j)\|} \right), \quad (1)$$

where i and j denote 2D spatial positions of feature maps, respectively, and $\text{ReLU}(\cdot)$ denotes ReLU activation function [72]. In this way, all pairwise feature matches are computed and stored. We subsequently concatenate the computed correlation maps to obtain a stacked correlation map $\mathcal{C} \in \mathbb{R}^{hw \times hw \times L}$. However, raw matching scores contain numerous ambiguous matching points as exemplified in Fig. 3, which results in inaccurate correspondences. To remedy this, we propose cost aggregation networks in the following that aim to refine the ambiguous or noisy matching scores.

3.3 Preliminary: Transformers

Before moving on to the proposed transformer aggregator, we start by explaining Vision Transformer (ViT) encoder [41]. Note that the proposed method does not utilize a decoder, which we omit its explanation in this section. Renowned for its global receptive fields, one of the key elements of transformer [41] is the self-attention mechanism, which enables finding the correlated input tokens by first feeding into scaled dot product attention function (self-attention layer [41]), normalizing with Layer Normalization (LN) [73], and passing the normalized values to an MLP, which is also known as Feed-Forward Network (FFN). Formulating the overall process of vision transformer [42], $y = \mathcal{T}(I)$, where y serves as the image representation obtained by prepending a learnable embedding prior to feeding into the encoder, it first undergoes Query-Key-Value projection as:

$$\begin{aligned} Q &= \mathcal{P}_Q(\text{LN}(I' + E_{\text{pos}})), \\ K &= \mathcal{P}_K(\text{LN}(I' + E_{\text{pos}})), \\ V &= \mathcal{P}_V(\text{LN}(I' + E_{\text{pos}})), \end{aligned} \quad (2)$$

where \mathcal{P}_Q , \mathcal{P}_K and \mathcal{P}_V denote query, key and value linear projection, respectively; I' denotes token embeddings; E_{pos} denotes positional embedding. Then the obtained Q , K , and V undergo a self-attention layer:

$$\hat{z} = \text{SA}(Q, K, V) = \text{softmax}(QK^T)V, \quad (3)$$

where $\text{SA}(\cdot)$ denotes Self-Attention. Subsequently, the output undergoes LN and residual connection followed by FFN that consists of two linear transformations with a GELU [74] activation in between, as shown in Fig. 7, which is formulated as:

$$\begin{aligned} z &= \hat{z} + z_{\text{pos}}, \\ x &= \text{LN}(z), \\ \hat{x} &= \mathcal{P}_{\text{FFN}}^2(\text{GELU}(\mathcal{P}_{\text{FFN}}^1(x))), \\ y &= \hat{x} + z, \end{aligned} \quad (4)$$

where $z_{\text{pos}} = I' + E_{\text{pos}}$, and $\mathcal{P}_{\text{FFN}}^1$ and $\mathcal{P}_{\text{FFN}}^2$ are the first and the second position-wise linear transformations within the feed-forward network, respectively.

3.4 Transformer Aggregator

Several works [27], [42], [67], [68] have shown that given images or features as input, transformers [41] integrate the global context information by learning to find the attention scores for all pairs of tokens. In this paper, we take a different approach, leveraging the transformers to integrate the pairwise relationships across matching scores to not only discover global consensus but also to carefully consider the pairwise interactions among tokens for an effective cost aggregation in a global manner. Unlike ViT [42], we do not use a class token and obtain a refined cost $\hat{\mathcal{C}}$ as output by feeding the stacked raw cost $\mathcal{C} \in \mathbb{R}^{hw \times hw \times L}$ to the transformer \mathcal{T} , consisting of QKV projections, SA, LN, and FFN:

$$\hat{\mathcal{C}} = \mathcal{T}(\mathcal{C}). \quad (5)$$

The standard transformers receive input as a 1D sequence of token embeddings. To this end, we also treat each spatial location at either source or target as 1D token. We visualize the refined correlation map with self-attention in Fig. 5, where the ambiguities are

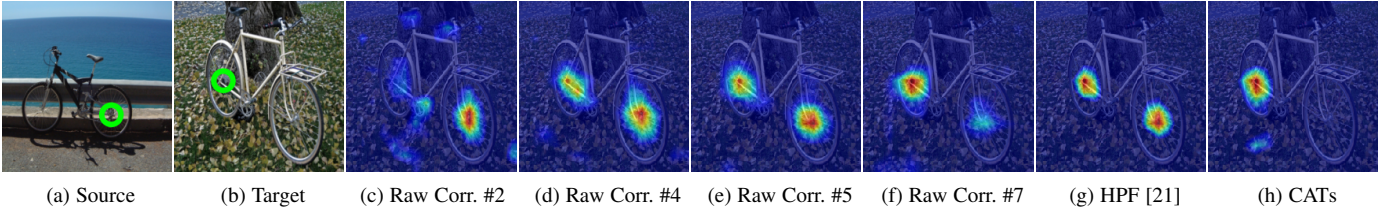


Fig. 3: **Visualization of multi-level aggregation:** (a) source, (b) target images, (c)-(f) raw correlation maps, respectively, and final correlation maps by (e) HPF [21] and (f) CATs. Note that HPF and CATs utilize the same feature maps. Compared to HPF, CATs successfully embrace richer semantics in different levels of feature maps.

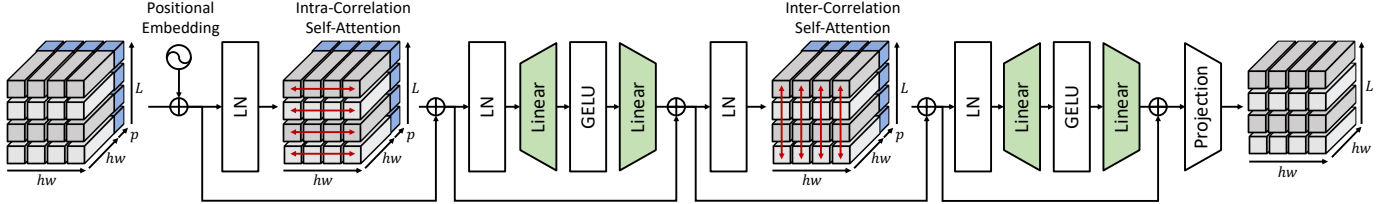


Fig. 4: **Illustration of Transformer aggregator.** Given correlation maps \mathcal{C} with projected features, transformer aggregator consisting of intra- and inter-correlation self-attention with LN and MLP refines the inputs not only across spatial domains but across levels.

significantly resolved. In the following, we introduce appearance affinity modeling to disambiguate the noisy correlation maps and multi-level aggregation to effectively aggregate across spatial- and level-dimension.

3.4.1 Appearance Affinity Modelling

When only matching costs are considered for aggregation, self-attention layer processes the correlation map disregarding the noise involved in the correlation map, which may lead to inaccurate correspondences. Rather than solely relying on a raw correlation map, we additionally provide an appearance embedding from input features to disambiguate the correlation map. The intuition behind is that visually similar points in an image, e.g., color or feature, have similar correspondences, as proven in stereo matching literature, e.g., Cost Volume Filtering (CVF) [10], [12].

To provide appearance affinity, we propose to concatenate embedded features projected from input features with the correlation map. We first feed the features D^l into the linear projection networks, and then concatenate the output along the corresponding dimension, so that the correlation map is augmented such that $[\mathcal{C}^l, \mathcal{P}^l(D^l)] \in \mathbb{R}^{hw \times (hw+p)}$, where $[\cdot]$ denotes concatenation, \mathcal{P} denotes linear projection networks, and p is channel dimension of the embedded feature. The transformer can be then formulated as

$$\hat{\mathcal{C}} = \mathcal{T}([\mathcal{C}^l, \mathcal{P}^l(D^l)]_{l=1}^L). \quad (6)$$

Specifically, within the transformer \mathcal{T} , this augmented correlation map then undergoes QKV projection:

$$\begin{aligned} Q &= \mathcal{P}_Q(\text{LN}([\mathcal{C}^l, \mathcal{P}^l(D^l)]_{l=1}^L + E_{\text{pos}})), \\ K &= \mathcal{P}_K(\text{LN}([\mathcal{C}^l, \mathcal{P}^l(D^l)]_{l=1}^L + E_{\text{pos}})), \\ V &= \mathcal{P}_V(\text{LN}([\mathcal{C}^l, \mathcal{P}^l(D^l)]_{l=1}^L + E_{\text{pos}})), \end{aligned} \quad (7)$$

where for E_{pos} , we let the networks be aware of the positional information by providing a learnable embedding [42] rather than a fixed [41] as shown in Fig. 4. Note that $\mathcal{P}_Q, \mathcal{P}_K$ and \mathcal{P}_V process the stacked correlation map \mathcal{C} , while \mathcal{P} processes at each l -th level of \mathcal{C} . After the QKV projections, following the process of a standard transformer, we feed the output into the self-attention layer, which is formulated in Eq. 3, to aggregate the matching costs

concatenated with appearance embedding to reason about their relationships, and then they undergo FFN, which is formulated in Eq. 4. Subsequently, the aggregated augmented correlation map is passed to the linear projection networks, to retain the size of original correlation \mathcal{C} , which we obtain $\hat{\mathcal{C}}$ as illustrated in Fig. 4.

3.4.2 Multi-level Aggregation

For each l -th level of stacked correlation maps $\mathcal{C} = [\mathcal{C}^l]_{l=1}^L$, level-wise matching scores are available, and we aim to employ all of them present across different levels. In order to employ both hierarchical and spatial semantic representations, we introduce a multi-level aggregation technique, which aggregates the stack of augmented correlation maps in a spatial- and level-wise manner. In this way, the aggregation networks can consider interactions among matching scores at a particular level, as well as across levels.

Specifically, as shown in Fig. 4, a stack of L augmented correlation maps, $[\mathcal{C}^l, \mathcal{P}^l(D^l)]_{l=1}^L \in \mathbb{R}^{hw \times (hw+p) \times L}$, undergo the transformer aggregator. For each l -th augmented correlation map, we aggregate with self-attention layer across all the points in the augmented correlation map, and we refer to this as *intra*-correlation self-attention. This is shown as a red arrow, which is a level-wise aggregation. In addition, subsequent to this, the correlation map undergoes *inter*-correlation self-attention across multi-level dimensions, which is a spatial-wise aggregation. Contrary to HPF [21] that concatenates all the multi-level features to compute a stack of correlation maps and aggregate using RHM, which only limitedly consider the multi-level similarities, within the inter-correlation self-attention layer of the proposed model, the similar matching scores are explicitly explored across multi-level dimensions. In this way, we can embrace richer semantics in different levels of feature maps, as shown in Fig. 3.

3.5 Cost Aggregation with Transformers

By leveraging the transformer aggregator, we present a cost aggregation framework with the following additional techniques to improve the performance, and by combining them, we complete the architecture of CATs.



Fig. 5: **Visualization of self-attention:** Each attention map attends to different aspects, and CATs aggregates the costs, leveraging rich semantic representations.

3.5.1 Swapping Self-Attention

In order to impose mutually consistent matching scores, we argue that reciprocal scores should be used to infer confident correspondences. As the correlation map contains bidirectional matching scores, from both target and source perspectives, we can leverage matching similarities from both directions in order to obtain more reciprocal scores as done similarly in other works [19], [33].

As shown in Fig. 2, we first feed the augmented correlation map to the aforementioned transformer aggregator. Then we transpose the output, swapping the pair of dimensions in order to concatenate with the embedded feature from the other image and feed it into the subsequent aggregator. Note that we share the parameters of the transformer aggregators to obtain reciprocal scores. Formally, we define the whole process as follows:

$$\begin{aligned} \mathcal{S} &= \mathcal{T}([\mathcal{C}^l, \mathcal{P}^l(D_t^l)]_{l=1}^L), \\ \hat{\mathcal{C}} &= \mathcal{T}([\mathcal{S}^l]^T, \mathcal{P}^l(D_s^l)]_{l=1}^L), \end{aligned} \quad (8)$$

where $\mathcal{C}^T(i, j) = \mathcal{C}(j, i)$ denotes swapping the pair of dimensions corresponding to the source and target images; \mathcal{S} denotes the intermediate correlation map before swapping the axis. Note that NC-Net [19] proposed a similar procedure, but instead of processing serially, they separately process the correlation map and its transposed version and add the outputs, which this procedure is designed to produce a correlation map invariant to the particular order of the input images. Unlike this, we process the correlation map serially, first aggregating the one pair of dimensions and then further aggregating the other pair. In this way, the subsequent attention layer is given cost maps with more consistent matching scores as an input, allowing further reduction of inconsistent matching scores. We justify our design choice in the ablation study.

3.5.2 Residual Connection

At the initial phase when the correlation map is fed into the transformers, noisy score maps are inferred due to randomly initialized parameters, which could complicate the learning process. To stabilize the learning process and provide a better initialization for the matching, we employ the residual connection. Specifically, we enforce the cost aggregation networks to estimate the residual correlation by adding residual connection around aggregation networks. Now Eq. 8 changes to:

$$\begin{aligned} \mathcal{S} &= \mathcal{T}([\mathcal{C}^l, \mathcal{P}^l(D_t^l)]_{l=1}^L) + \mathcal{C}, \\ \hat{\mathcal{C}} &= \mathcal{T}([\mathcal{S}^l]^T, \mathcal{P}^l(D_s^l)]_{l=1}^L) + \mathcal{C}^T. \end{aligned} \quad (9)$$

3.6 Training Objective

As in [21], [22], [28], we assume that the ground-truth keypoints are given for each pair of images. We first average the stack of refined correlation maps $\hat{\mathcal{C}} \in \mathbb{R}^{hw \times hw \times L}$ along level dimension and then transform it into a dense flow field F_{pred} using soft-argmax operator [33]. Subsequently, we compare the predicted dense flow field with the ground-truth flow field F_{GT} obtained by following the protocol of [21] using input keypoints. For the training objective, we utilize Average End-Point Error (AEPE) [20], computed by averaging the Euclidean distance between the ground-truth and estimated flow. We thus formulate the objective function as $\mathcal{L} = \|F_{\text{GT}} - F_{\text{pred}}\|_2$.

4 BOOSTING COST AGGREGATION WITH CONVOLUTIONS AND TRANSFORMERS

4.1 Motivation and Overview

Despite its deliberate design and effectiveness, CATs may struggle with high computational costs induced by the use of standard transformers [41], and this could restrict the cost aggregation to be performed only at limited resolutions. More specifically, as stated in some works [41], [42], [75] that computational complexity for self-attention layer is quadratic to token length, and those of QKV projection and FFN are quadratic to feature dimension, processing high-dimensional correlation maps with standard transformers [41] inevitably face large computational burden.

To alleviate this, we propose to introduce early convolutions prior to cost aggregation with transformers to control the number of tokens to reduce the costs and additionally inject some convolutional inductive bias to enhance the subsequent cost aggregation, which results in an apparent performance boost. Moreover, to reduce the computational loads from QKV projection and FFN, we introduce a novel transformer architecture for efficient cost aggregation that reformulates the two components by including appearance embedding inside the QKV projection and introducing convolutions to replace linear projections, which not only results in a costs reduction as intended but also an apparent performance boost. Thanks to the reduced costs, we are permitted to process the inputs of higher resolutions and manage to build a hierarchical architecture. Moreover, we deviate from arbitrarily selecting different combinations of feature maps [21] as was done in CATs, but rather use all the feature maps to exploit richer semantic representations. With these combined, we introduce an extension of CATs, namely CATs++.

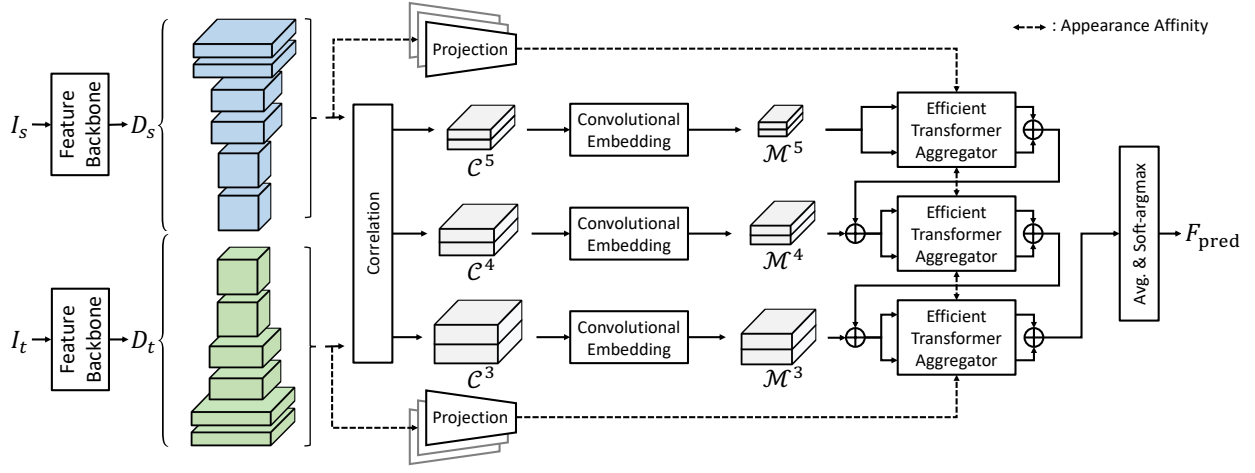


Fig. 6: **Overall network architecture of CATs++.** The networks of CATs++ differ from CATs in that all the feature maps are used to compute cost volumes rather than selecting different combinations for different datasets, and early convolution along with the efficient transformer aggregator helps to reduce computational costs and improve performance. Note that we colored the extracted feature maps to indicate that we use all the intermediate feature maps and highlight the difference to Fig. 2.

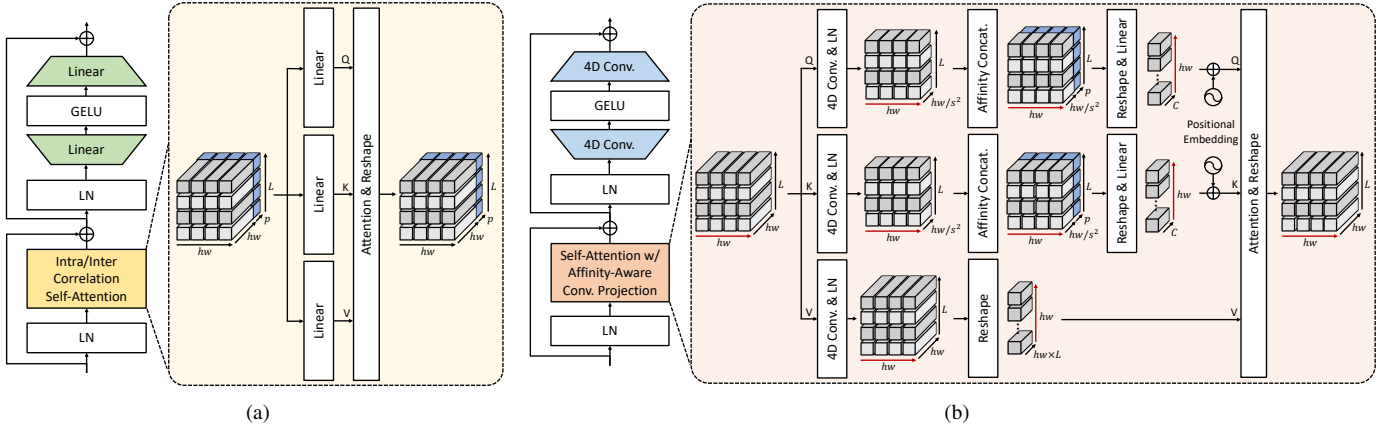


Fig. 7: **Intuition of the proposed components:** (a) standard transformer-based aggregator (in CATs), (b) efficient transformer-based aggregator (in CATs++). CATs++ replaces the query, key and value projections and feed-forward networks of standard transformer aggregator. With this novel design, we achieve significant performance boost as well as meaningful costs reductions.

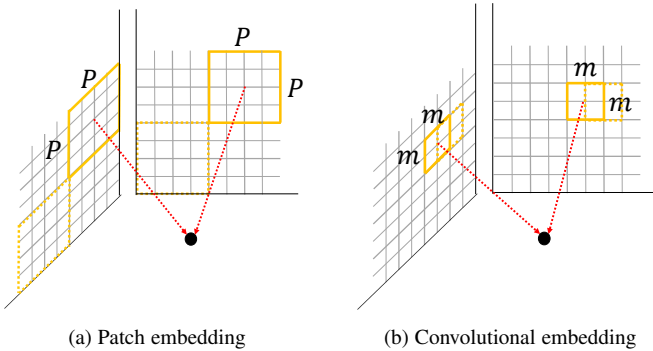


Fig. 8: **Comparison of different embedding strategy.** Unlike patch embedding [42] that utilizes a large kernel for non-overlapping convolutions, our convolutional embedding uses a series of small kernels for overlapping convolutions.

4.2 Feature Extraction and Cost Computation

To exploit rich semantics present at different feature levels, we also leverage multi-level features as done in CATs and [21], [22], [33], [40]. However, taking a slightly different approach, we use all the

intermediate feature maps as done in [3]. Specifically, similar to CATs, we produce a sequence of L feature maps, but unlike how CATs selected a specific number of feature maps for each dataset, which is an extra burden, CATs++ use all the feature maps. In this way, not only we do not need to select different combinations of feature maps for each dataset we use, but also we can exploit richer semantics by using all the intermediate feature maps.

For the cost computation, following [3], given a pair of feature maps, D_s^l and D_t^l , we compute a correlation map using the inner product between $l-2$ normalized features. We then collect correlation maps of the same spatial sizes and denote the subset as $\{C^l\}_{l \in \mathcal{L}^q}$, where \mathcal{L}^q is a subset of CNN layer indices $\{1, \dots, L\}$ at some pyramid layer q . Lastly, we concatenate all the collected $\{C^l\}_{l \in \mathcal{L}^q}$ to obtain a hypercorrelation $C^q \in \mathbb{R}^{h_s w_s \times h_t w_t \times |\mathcal{L}^q|}$, where h_s, w_s and h_t, w_t are height and width of feature maps of source and target images, respectively.

4.3 Convolutional Embedding Module

Without a means to control the high computational costs induced by processing the hypercorrelation C^q , the amount of time, resources, and memory required would not be negligible. Considering the quadratic complexity of a standard transformer [41],

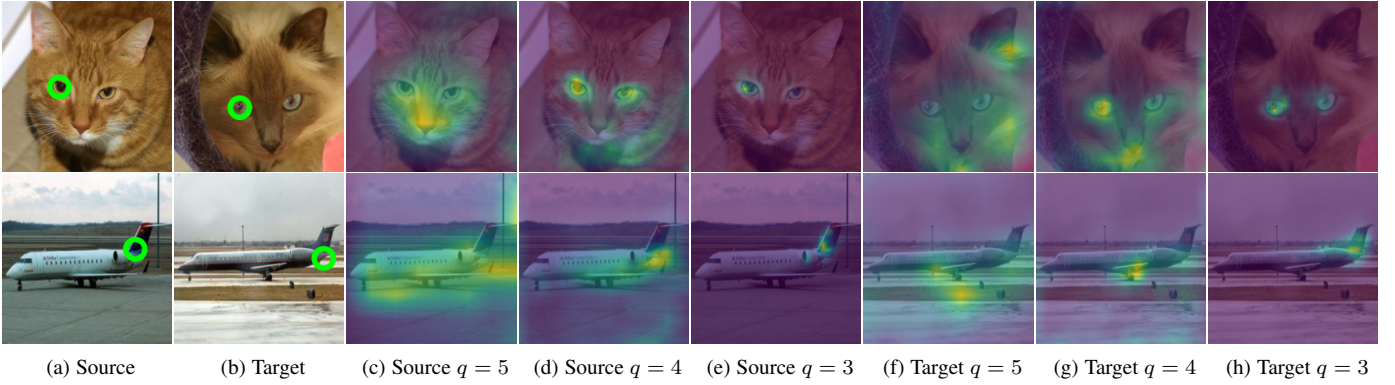


Fig. 9: **Visualization of self-attention.** Note that CATs++ utilizes hierarchical semantic representations, which the finer level self-attention is shown in the smaller region, and higher q denotes coarser level.

simply treating all the spatial dimensions as tokens seems unwise. A better approach and the most straightforward approach, perhaps, would be to reduce the number of tokens by adopting 4D spatial pooling across source and target dimensions. However, this strategy risks losing some information, which we ought to avoid and will be verified in Section 5.4.1. An alternative option could be to patchify the hypercorrelation by splitting it into non-overlapping tensors and embed with a large learnable kernel as done in ViT [42]. Naïve implementation to this could be extending 2D patch embedding to 4D as shown in Fig. 8 (a). However, as claimed in [58], non-overlapping operations only provides limited inductive bias, which only the relatively lower translation equivariance is achieved without replacing such operations. Moreover, disregarding window boundaries due to non-overlapping kernels may hurt overall performance due to discontinuity.

To compensate for the issues, analogously to VAT [58], we use a small early convolution layer to process the correlation map prior to transformer aggregation. Given a hypercorrelation \mathcal{C}^q , we consider receptive fields of 4D window, *i.e.*, $m \times m \times m \times m$, and build a tensor $\mathcal{M}^q \in \mathbb{R}^{\hat{h}_s \times \hat{w}_s \times \hat{h}_t \times \hat{w}_t \times |\mathcal{L}^q|}$, where \hat{h} and \hat{w} are the processed sizes as illustrated in Fig. 8 (b). Note that we reduce both spatial dimensions within hypercorrelation for computational efficiency, while VAT [58] preserves its spatial dimension of query image to obtain a fine-grained segmentation map. During this operation, we locally mix the multi-level similarity vector at each 4D position, which acts as a local inter-correlation aggregation. In this way, we enhance the subsequent cost aggregation with transformers by injecting some convolutional bias to mitigate the above issues and allowing both local and global aggregation.

4.4 Efficient Transformer Aggregator

Although the convolutional embedding module adjusts the number of tokens and feature dimensions to some extent, directly applying standard transformers [41] is still challenging. This is due to what some works [41], [42], [75] stated that the standard self-attention mechanism has quadratic complexity with respect to the number of tokens, and QKV projections and FFN have quadratic complexity with respect to the feature dimensionality. To address this, recently, numerous works [76], [77], [78], [79], [80] have proposed to reduce the computational burden by introducing an efficient transformer or self-attention, *i.e.*, token pruning [81], linear self-attention [77] and locality-sensitive hashing [76]. However, although these methods effectively alleviate the computational burden by aiming to either control the number of tokens or reduce

the complexity of self-attention, in our context, QKV projections and FFN relatively impose more computational burden as we flatten one of the spatial dimension, *i.e.*, the spatial dimension of source, and level dimension to treat them as feature dimension for a token. This has been relatively underexplored [82], and we propose a novel design specialized in matching cost aggregation that addresses the aforementioned issues.

4.4.1 Affinity-Aware Convolutional Projection

As illustrated in Fig. 7, the proposed QKV projection differs from that of CATs that follows the standard approach in that it uses convolutions to first process high resolution input to control the subsequent computation by reducing the spatial dimension from hw to hw/s^2 , where s denotes a stride size. By introducing convolutions prior to linear projection, not only we provide efficiency by allowing feature dimension controllable, but also strengthen the power to consider locality. While [83] tries a similar approach by applying depth-wise separable convolutions to reduce the number of tokens, here we aim to reduce the feature dimension. Subsequently, we concatenate appearance embedding to the intermediate outputs of query and queue projections, introducing appearance affinity inside the transformer architecture. It should be noted that linear projection is included only within the Query and Key projections as we aim to not only utilize and preserve the spatial resolution of the raw correlation map but also to further reduce the computation as we find that it is sufficient to only include appearance affinity information within the attention score computation QK^T rather than within the $(QK^T) \cdot V$ computation in order to aggregate both the appearance affinity along with the matching scores, which the QK projection is responsible for. Then we further reduce the feature dimension for Q and K projection by feeding them into linear projection, reducing self-attention computational loads. In this way, we benefit from the reduced computation, eliminating one-third of computational loads of value projection, had the linear projection been included.

Concretely, the process can be expressed as follows:

$$\begin{aligned} Q^q &= \mathcal{P}_Q^q([\text{LN}(\mathcal{Q}_Q^q(\text{LN}(\mathcal{M}^q))), \mathcal{P}^q(D^q)]) + E_{\text{pos}}^q, \\ K^q &= \mathcal{P}_K^q([\text{LN}(\mathcal{Q}_K^q(\text{LN}(\mathcal{M}^q))), \mathcal{P}^q(D^q)]) + E_{\text{pos}}^q, \\ V^q &= \text{LN}(\mathcal{Q}_V^q(\text{LN}(\mathcal{M}^q))), \end{aligned} \quad (10)$$

where \mathcal{Q} denotes convolutional projection; E_{pos}^q denotes positional embedding at q -th layer. Note that projections are performed at each q -th layer.

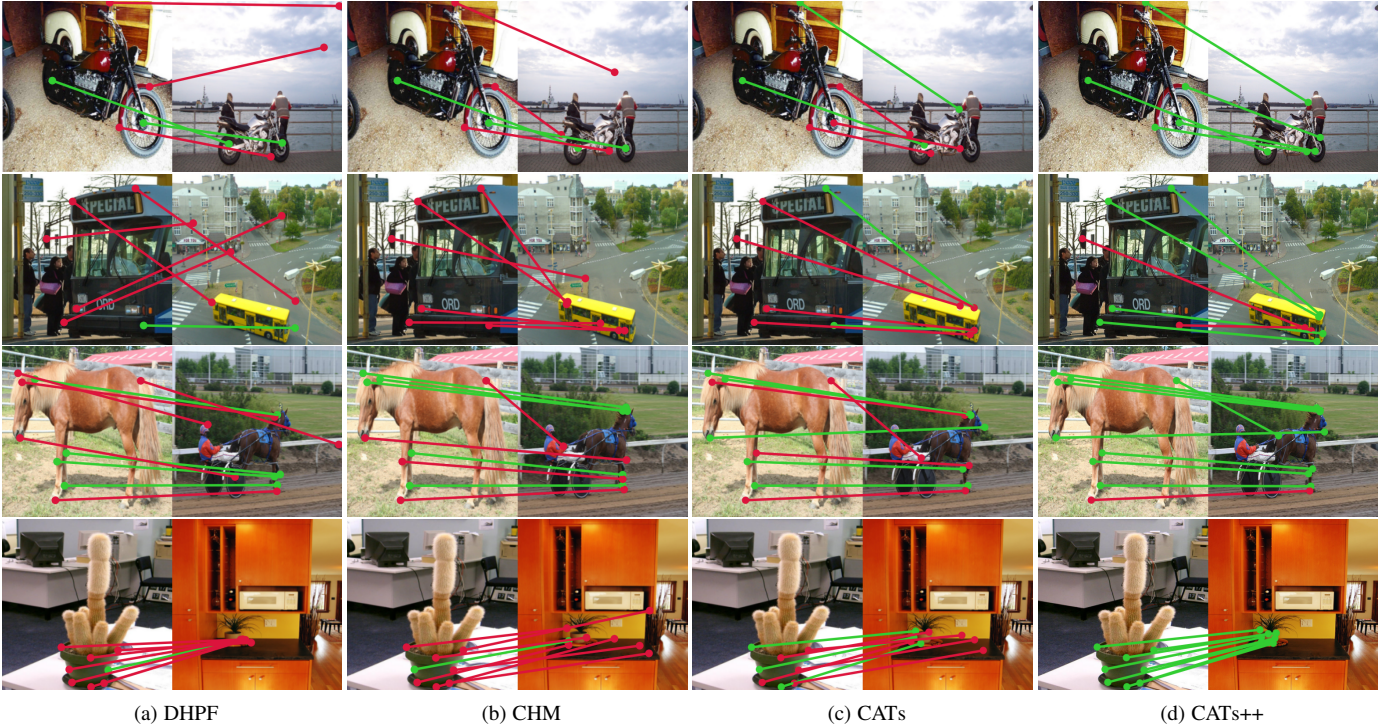


Fig. 10: **Qualitative results on SPair-71k [16]**: keypoints transfer results by DHPF [22], CHM [28], CATs, and CATs++. Note that the green and red line denote correct and wrong prediction, respectively, with respect to the ground-truth.

4.4.2 Volumetric Convolutional Feed-Forward

Given outputs of QKV projections, we follow the standard self-attention computation to obtain \hat{z} as in Eq. 3. Note that we find it relatively less influential to change it to other efficient self-attention computation methods [76], [77], [84], and we thus use standard self-attention [42] for simplicity. Subsequent to self-attention computation, in order to model additional locality bias, reduce the number of parameters for better generalization and provide more memory efficiency, we replace linear transformations and ReLU activation with 4D convolutions and GELU, respectively, as shown in Fig. 7. Moreover, it should be noted that flattening the correlation volume inevitably requires a linear projection with a large feature dimension, which results in a significant computation load. This provides efficiency benefits to using convolution over a large linear projection layer. The process is defined as:

$$\begin{aligned}
 z &= \hat{z} + \mathcal{M}^q, \\
 x &= \text{LN}(z), \\
 \hat{x} &= \mathcal{Q}_{\text{FFN}}^2(\text{GELU}(\mathcal{Q}_{\text{FFN}}^1(x))), \\
 y &= \hat{x} + z,
 \end{aligned} \tag{11}$$

where $\mathcal{Q}_{\text{FFN}}^1$ and $\mathcal{Q}_{\text{FFN}}^2$ are the first and the second convolutional projections within the feed-forward network, respectively. We find that this simple change brings a surprisingly apparent performance boost when combined with our proposed architecture, which will be discussed in Section 5.4.5.

4.5 Cost Aggregation with Pyramidal Processing

Combining all the proposed components of CATs++, we manage to reduce the costs and process higher resolution inputs with pyramidal architecture. Specifically, analogous to [3], [12], [20], [24], [36], [49] and inspired by several works [3], [19], [24], [49],

we also utilize the coarse-to-fine approach through a pyramidal processing as illustrated in Fig. 6 to allow matching scores processed at previous levels to guide the aggregation process at finer levels by employing residual connections. In this way, we enforce the networks to focus on learning complementary matching scores [48] as well as let the coarse level aggregation guide the finer level aggregation.

Overall, we define the whole process of our pyramidal cost aggregation as follows:

$$\begin{aligned}
 \hat{\mathcal{C}}^q &= \mathcal{T}_+^q(\mathcal{M}^q, \mathcal{P}^q(D^q)) + \mathcal{T}_+^q((\mathcal{M}^q)^T, \mathcal{P}^q(D^q)), \\
 \mathcal{M}^{q-1} &\leftarrow \text{up}(\hat{\mathcal{C}}^q) + \mathcal{M}^{q-1},
 \end{aligned} \tag{12}$$

where \mathcal{T}_+ denotes efficient transformer aggregator and $\text{up}(\cdot)$ denotes bilinear upsampling.

4.6 Comparison to Concurrent Work

Recently, a concurrent work, Volumetric Aggregation with Transformers (VAT) [58], also proposed a Transformer-based architecture to tackle few-shot segmentation task. Although VAT shares a similar high-level idea, which is to combine convolutions and Transformer, and evaluates its effectiveness on both few-shot segmentation and semantic correspondence tasks, there exist apparent differences, which we summarize in the following.

First of all, the motivations behind the designs of each architecture differ as their main target tasks are different, *i.e.*, few-shot segmentation and semantic correspondence. Although they may share similar challenges, *e.g.*, background clutters and intra-class variations, at high-level, VAT [58] designs its architecture to enhance its generalization power to handle new classes beyond training samples tailored for few-shot segmentation, while CATs++ redesigns the components of the standard Transformer for efficient computations tailored for semantic correspondence.

In this regard, VAT [58] and our CATs++ formulated different transformer aggregators. It should be noted that the components prior to transformer aggregator, *i.e.*, hypercorrelation and early convolutions, are the designs that both architectures share, which means that the subsequent modules lead to notable variations in effectiveness and efficiency. VAT [58] adopts Swin Transformer [85] as its transformer aggregator and extends to handle 4D cost volumes with 4D local windows, while CATs++ redesigns the standard transformer to a more efficient version. To address the complexity, VAT [58] adopts shifting window multi-head self-attention [85] that is designed to reduce the computational complexity of a global multi-head self-attention module. Unlike this, a key contribution that differentiates the proposed method to other Transformer-based architecture is that CATs++ reformulates the standard QKV projections and FFN, and proposes affinity-aware convolutional projections and volumetric convolutional feed-forward network, aiming to reduce the memory and computational complexity that are typically high when the input is high dimensional correlation maps.

More importantly, we argue that the relative positioning bias that Swin Transformer [85] include is what makes it excel at predicting segmentation map for unseen classes. We observe that relative positioning bias allows greater generalization power, which will be evidenced in Table 5 in experiments, which was also observed in Swin Transformer [85], and this led to performance improvement for few-shot segmentation benchmark. However, for the semantic matching task, where generalization power has a relatively smaller influence to performance than to few-shot segmentation, we observe CATs++ outperforms VAT. This can be explained by the special design of CATs++ that integrates convolutions and appearance embedding with the proposed transformer aggregator, which not only allows an improved efficiency, but also a clear difference in architecture that leads to an apparent performance difference to VAT [58] for semantic correspondence task.

5 EXPERIMENTS

5.1 Implementation Details

For the backbone feature extractor, we use ResNet-101 [86] pre-trained on ImageNet [87] for both CATs and CATs++, while other backbone features can also be used. For the hyperparameters of the transformer encoder, we set the number of encoders as 1 for all the transformer aggregators in both CATs and CATs++ based on the ablation study, which we chose to set it as 1 for simplicity. For CATs, we resize the spatial size of the input image pairs to 256×256 and a sequence of selected features is resized to $h = 16$ and $w = 16$, while we resize image pairs to 512×512 for CATs++. For pyramidal processing in CATs++, we employ pyramidal layers of $q = 3, 4, 5$: from conv_3_x to conv_5_x , similar to [3]. AdamW [88] optimizer with an initial learning rate of $3e - 5$ for the aggregators and $3e - 6$ for the backbone features are used, which we gradually decrease during training. To apply data augmentation [89], [90] with predetermined probabilities to input images at random. Specifically, 50% of the time, we randomly crop the input image, and independently for each augmentation function used in [89], we set the probability for applying the augmentation as 20%. We implemented our network using PyTorch [91].

5.2 Experimental Settings

In this section, we conduct comprehensive experiments for semantic correspondence, by evaluating our approach through comparisons to state-of-the-art methods including CNNGeo [17], A2Net [34], WeakAlign [18], NC-Net [19], RTNs [35], SFNet [33], HPP [21], DCC-Net [92], ANC-Net [39], DHPF [22], SCOT [23], GSF [38], and CHM [28], CATs [45], PMNC [51], MMNet [48], and CHMNet [40].

5.2.1 Datasets

We consider SPair-71k [16] which provides a total of 70,958 image pairs with extreme and diverse viewpoints, scale variations, and rich annotations for each image pair, *e.g.*, keypoints, scale difference, truncation and occlusion difference, and clear data split. Previously, for semantic matching, most of the datasets are limited to a small quantity with similar viewpoints and scales [14], [15]. As our network relies on transformer which requires a large number of data for training, SPair-71k [16] makes the use of transformer in our model feasible. We also consider PF-PASCAL [15] containing 1,351 image pairs from 20 categories and PF-WILLOW [14] containing 900 image pairs from 4 categories, each dataset providing corresponding ground-truth annotations.

5.2.2 Evaluation Metric

For evaluation on SPair-71k [16], PF-PASCAL [15], and PF-WILLOW [14], we employ a percentage of correct keypoints (PCK), computed as the ratio of estimated keypoints within the threshold from ground-truths to the total number of keypoints. Given predicted keypoint k_{pred} and ground-truth keypoint k_{GT} , we count the number of predicted keypoints that satisfy the following condition: $d(k_{\text{pred}}, k_{\text{GT}}) \leq \alpha \cdot \max(H, W)$, where $d(\cdot)$ denotes Euclidean distance; α denotes a threshold. We evaluate on PF-PASCAL with α_{img} , SPair-71k, and PF-WILLOW with α_{bbox} . H and W denote height and width of the object bounding box or the entire image, respectively. As stated in [40], we additionally report results of $\alpha_{\text{bbox-kp}}$ for PF-WILLOW [14] for a fair comparison.

5.3 Semantic Correspondence Results

For a fair comparison, we follow the evaluation protocol of [21] for SPair-71k [16], in which our network is trained on the training split and evaluated on the test split. Similarly, for PF-PASCAL [15] and PF-WILLOW [14], following the common evaluation protocol of [21], [22], [35], [92], [93], we train our network on the training split of PF-PASCAL [15] and then evaluate on the test split of PF-PASCAL [15] and PF-WILLOW [14]. All the results of other methods are reported under the identical setting.

Fig. 10 visualizes qualitative results for extremely challenging image pairs. We observe that our methods are capable of suppressing noisy scores and finding accurate correspondences in cases where large scale and geometric variations are present. Table 1 summarizes quantitative results on SPair-71k [16], PF-PASCAL [15] and PF-WILLOW [14]. In order to ensure a fair comparison, we note whether each method leverages multi-level features and fine-tunes the backbone networks. We additionally denote the types of cost aggregation.

We first compare the proposed methods with those that do not fine-tune backbone networks. Comparing CATs† and CATs++† to other methods, we find that not only CATs† outperforms

TABLE 1: **Quantitative evaluation on standard benchmarks [14], [15], [16].** Higher PCK is better. The best results are in bold, and the second best results are underlined. CATs \dagger (or CATs $++\dagger$) denotes CATs (or CATs $++$) without fine-tuning feature backbone. Note that ResNet-101 is used as a backbone feature extractor as a default setting unless stated otherwise, *i.e.*, MMNet [48]. *Feat.-level:* Feature-level, *FT. feat.:* Fine-tune feature.

Methods	Data Aug.	Feat.-level	FT. feat.	Cost Aggregation	SPair-71k [16]	PF-PASCAL [15]			PF-WILLOW [14]			
					PCK @ α_{bbox} 0.1	PCK @ α_{img} 0.05	0.1	0.15	PCK @ α_{bbox} 0.05	0.1	0.05	PCK @ $\alpha_{\text{bbox-kp}}$ 0.1
WTA	\times	Single	\times	-	25.7	35.2	53.3	62.8	30.1	52.9	24.7	46.9
CNNGeo [17]	\times	Single	\times	-	20.6	41.0	69.5	80.4	-	-	36.9	69.2
A2Net [34]	\times	Single	\times	-	22.3	42.8	70.8	83.3	-	-	36.3	68.8
WeakAlign [18]	\times	Single	\times	-	20.9	49.0	74.8	84.0	-	-	37.0	70.2
RTNs [35]	\times	Single	\times	-	25.7	55.2	75.9	85.2	-	-	41.3	71.9
SFNet [33]	\times	Multi	\times	-	-	53.6	81.9	90.6	-	-	46.3	74.0
MMNet-FCN [48]	\times	Multi	\checkmark	-	50.4	81.1	91.6	95.9	-	-	-	-
NC-Net [19]	\times	Single	\checkmark	4D Conv.	20.1	54.3	78.9	86.0	-	-	33.8	67.0
DCC-Net [92]	\times	Single	\times	4D Conv.	-	55.6	82.3	90.5	-	-	43.6	73.8
HPF [21]	\times	Multi	-	RHM	28.2	60.1	84.8	92.7	-	-	45.9	74.4
GSF [38]	\times	Multi	\times	2D Conv.	36.1	65.6	87.8	95.9	-	-	49.1	78.7
ANC-Net [39]	\times	Single	\times	4D Conv.	-	-	86.1	-	-	-	-	-
DHPF [22]	\times	Multi	\times	RHM	37.3	75.7	90.7	95.0	49.5	77.6	-	71.0
SCOT [23]	\times	Multi	-	OT-RHM	35.6	63.1	85.4	92.7	-	-	<u>47.8</u>	<u>76.0</u>
CHM [28]	\times	Single	\checkmark	6D Conv.	46.3	80.1	91.6	94.9	52.7	79.4	-	69.6
PMNC [51]	\times	Multi	\checkmark	4D Conv.	50.4	<u>82.4</u>	90.6	-	-	-	-	-
CHMNet [40]	\checkmark	Multi	\checkmark	6D Conv.	51.3	81.3	<u>92.9</u>	-	53.8	79.3	-	69.3
VAT [58]	\checkmark	Multi	\checkmark	4D Conv. + Transformer	<u>55.5</u>	78.2	92.3	96.2	52.8	81.6	-	-
CATs \dagger	\checkmark	Multi	\times	Transformer	42.4	67.5	89.1	94.9	46.6	75.6	37.3	65.7
CATs	\checkmark	Multi	\checkmark	Transformer	49.9	75.4	92.6	<u>96.4</u>	50.3	79.2	40.7	69.0
CATs $++\dagger$	\checkmark	Multi	\times	4D Conv. + Transformer	50.0	73.0	89.6	95.0	<u>54.1</u>	78.5	44.5	69.9
CATs $++$	\checkmark	Multi	\checkmark	4D Conv. + Transformer	59.8	84.9	93.8	96.8	56.7	<u>81.2</u>	47.0	72.6

TABLE 2: **Ablation study of CATs.**

Components		SPair-71k $\alpha_{\text{bbox}} = 0.1$	
		FT. feat.	\times
(I)	Baseline	26.8	46.7
(II)	+ Appearance Modelling	33.5	46.3
(III)	+ Multi-level Aggregation	35.9	47.0
(IV)	+ Swapping Self-Attention	38.8	47.6
(V)	+ Residual Connection	42.4	49.9

TABLE 3: **Ablation study of CATs $++$.**

Components		SPair-71k $\alpha_{\text{bbox}} = 0.1$	
		FT. feat.	\times
(I)	Baseline	22.7	49.8
(II)	+ Hypercorrelation	20.2	50.2
(III)	+ Convolutional Embedding	27.9	54.1
(IV)	+ Efficient Transformer Aggregator	42.1	55.3
(V)	+ Appearance Modelling	50.0	59.8
(VI)	(IV) - Convolutional Embedding	47.4	57.8
(VII)	(V) - Convolutional Embedding	48.5	59.1
(VIII)	(V) - Hypercorrelation	46.7	58.1
(IX)	(VIII) - Convolutional Embedding	45.6	55.7

DHPF [22], by 5.1%*p* on SPair-71k, but also CATs $++\dagger$ further boosts the performance to achieve 50.0%, which makes it beyond comparison. Interestingly, the results obtained by CATs $++\dagger$ are almost on par with the previous state-of-the-art method, CHMNet [40], even without fine-tuning the backbone. With the backbone networks fine-tuned, CATs shows highly competitive performance, achieving similar performance to previous state-of-the-art methods, thanks to its ability to explore global consensus and powerful data augmentation technique to fulfill the need of data-hungry transformer. It should be noted that CATs $++$ sets a new state-of-the-art on almost all the benchmarks, and it surpasses CHMNet by 8.5%*p* on SPair-71k, which clearly demonstrates

its effectiveness and highlights that cost aggregation is of prime importance, and leveraging both convolutions and transformers in a right way clearly makes stronger cost aggregation. Note that VAT [58] also achieves highly competitive results, even surpassing the results of proposed method on PF-WILLOW [14] at $\alpha = 0.1$. Nevertheless, it is shown that the proposed approach outperforms for all other benchmarks, highlighting the superiority of the proposed efficient transformer to Swin Transformer [85] in VAT [58] for semantic correspondence task.

Now comparing CATs and CATs $++$, we observe a significant performance boost for both backbone fine-tuned and frozen. Specifically, when the backbone is frozen, 7.6%*p* increase is observed while there is a 9.9%*p* increase when backbone is fine-tuned. This clearly demonstrates the effectiveness extension of CATs. Also, we observe from the results of PF-WILLOW [14] that more fine details are predicted correctly as well as the generalization power is enhanced. This is confirmed by how CATs $++\dagger$ is almost on par with CATs on PF-WILLOW even though their performance gap on PF-PASCAL shows a large gap.

Additionally, as stated in CHMNet [40], for a fair comparison, we report results of PCK @ $\alpha_{\text{bbox-kp}}$ for PF-WILLOW [14]. It is notable that our methods generally report lower PCK on PF-WILLOW [14] compared to other state-of-the-art methods. We conjecture that this is due to the provision of sparsely annotated data and the use of transformers, since sparse annotations can hurt the generalization power [40]. Moreover, it is well known that the transformers lack generalization power with a small dataset size [42]. When we evaluate on PF-WILLOW, we infer with the model trained on the training split of PF-PASCAL [15] with sparse keypoint annotations, which only contains 1,351 image pairs, and as an only relatively small quantity of image pairs is available within the PF-PASCAL training split, the proposed method shows low generalization power. Note that to compensate for this issue, we provided data augmentation techniques, which clearly helped to increase the generalization power. From this,

TABLE 4: **Component analysis of efficient transformer aggregator.** A.A.C: Affinity-aware Convolutional. V.C: Volumetric Convolutional.

Component	SPair-71k $\alpha_{\text{bbox}} = 0.1$	Memory [GB]	Num. of param.
Baseline	36.7	2.2	24.9M
(I) A.A.C QKV projection	46.6	2.0	14.4M
(II) V.C feed-forward	40.0	2.0	15.9M
(I+II) Efficient transformer	50.0	1.9	5.5M

we suspect that a provision of more data could help to address the generalization issue, and perhaps the weak supervisions other than sparse keypoints could further enhance its generalization power.

5.4 Ablation study

In this section, we show an ablation analysis to validate critical components we made to design our architecture, and provide analysis on the use of different backbone features, the depth of transformer encoder, comparison among different cost aggregators, memory/run-time, and data augmentation. We train all the variants on the training split of SPair-71k [16] when evaluating on SPair-71k, and train on PF-PASCAL [15] for evaluating on PF-PASCAL. We measure the PCK, and each ablative experiment is conducted under the same experimental setting for a fair comparison.

5.4.1 Effectiveness of each component

Table 2 shows the analysis of key components of CATs. We analyze four key components for this ablation study, which include appearance modeling, multi-level aggregation, swapping self-attention, and residual connection. As shown in Table 3, we also conduct an ablation study on key components of CATs++, which we consider hypercorrelation, convolutional embedding, efficient transformer aggregator, and appearance modeling. For both analyses, we evaluate on SPair-71k benchmark by progressively adding each key component. In this experiments, we report the results for the backbone models with and without fine-tuned. Moreover, to explore how the proposed modules are co-affected, we evaluate the models equipped with different combinations of the modules, and the results are reported in Table 3 from VI to IX.

In Table 2, we first start with CATs by defining the model without any of the key components as a baseline model, which simply feeds the correlation map into the self-attention layer. From I to V, we observe a consistent improvement in performance when each component is added. The results of I show relatively lower PCK, indicating that simply utilizing transformers does not yield high performance, which highlights the importance of each of our proposed components. Interestingly, II shows a large improvement in performance, which we find that the appearance modeling enables the model to refine the ambiguous or noisy matching scores and the use of both appearance and cost volume allows the transformer to relate the pairwise relationships and appearance, helping to find more accurate correspondences. Although a relatively small increase of PCK for III, it shows that the proposed model successfully aggregates the multi-level correlation maps with the help of intra- and inter- correlation self-attention. Furthermore, IV and V show an apparent increase by helping the training process, clearly confirming the significance of both components.

In Table 3, for CATs++, we define the baseline model that does not utilize appearance affinity, only uses feature maps of

the last index at each p -th layer, replaces convolutional embedding with simple 4D max-pooling with linear projection, and utilizes raw correlation map without cost aggregation. In the results, we observe that each component except II contributes to a surprisingly large performance boost. Especially from III to IV, the performance gap is $14.2\%p$, which clearly demonstrates the effectiveness of the proposed efficient transformer aggregator. It should be noted that V is actually included within our proposed efficient transformer aggregator, which means that the proposed efficient transformer by itself brings a $22.1\%p$ increase. Interestingly, we observe a slight drop in performance when we attempt to use hypercorrelation as shown in II. We conjecture that although the use of multi-level features generally helps to find better correspondences [21], [22], [40], [45], without a means to aggregate the costs, the use of multi-level features rather make the matching process more complex.

We also observe consistent improvements in performance even with backbone network fine-tuned when each module is cumulatively stacked. However, the results for VI and VIII deviate from this trend, where the PCK of VIII with backbone network fine-tuned is higher than that of VI despite lower PCK for the backbone model without fine-tuning. Also, we find that although we generally observe consistent performance improvements as we stack up the components, II in Table 2 shows a slight performance drop. This may seem trivial for only a 0.4% drop, but these could indicate that fine-tuning the backbone network may have had larger influence than introducing a new component. Another notable result is VII, which shows slightly lower PCK to the best PCK reported, indicating that although convolutional embedding can help to achieve higher performance, highly competitive results can be attained even without convolutional embedding.

5.4.2 Effectiveness of each component in efficient transformer aggregator

In this ablation study, we provide an analysis of the proposed efficient transformer aggregator of CATs++ with respect to performance, memory, and number of learnable parameters. For this study, we first define the baseline model; it simply replaces the efficient transformer aggregator with standard transformer [41]. We investigate the effectiveness of the proposed components by replacing each component in a standard transformer with our proposed component, as illustrated in Table 4. Note that for calculating both memory consumption and the number of parameters, we exclude memory and parameters of backbone networks and convolutional embedding in order to emphasize the changes made by adopting each of our proposed efficient transformer aggregators.

Replacing QKV projection with affinity-aware convolutional projection, we observe that the PCK improves by $9.9\%p$. We also observe a similar increase when we replace feed-forward network with the proposed volumetric convolutional feed-forward. This apparent performance boost demonstrates their effectiveness, and this is further confirmed when two components are combined, which our efficient transformer aggregator obtains PCK of $50.0\%p$, a $13.3\%p$ higher results than that of the baseline model. We also report memory consumption and the number of learnable parameters to validate their efficiency. Overall, each component reduces the memory consumption by approximately 0.2 GB, which makes roughly 0.3 GB or 11% memory reduction in total. What surprised us is that the number of learnable parameters

TABLE 5: **Ablation study of cost aggregation methods.** Note that we abandon some of our key contributions in order to implement 3D convolutions. *RPB*: *Relative Positioning Bias*.

Aggregator	SPair-71k $\alpha_{\text{bbox}} = 0.1$	FSS-1000 mIoU (1 shot)	Memory [GB]
MLP	34.4	-	2.0
MLP-Mixer [94]	39.1	-	2.2
3D convolutions [95]	30.6	-	2.3
Center-pivot 4D convolutions [40]	36.9	-	1.7
Standard Transformer [41]	42.4	80.8	1.9
Standard Transformer [41], [85] (w/ RPB)	42.6	81.8	1.9

is reduced by almost 80%, which may have resulted in better generalization power.

5.4.3 Why transformers for cost aggregation?

As our architecture is based on transformers, perhaps it is necessary to compare the effectiveness of transformers to other methods that are capable of relating the tokens of given inputs, *i.e.*, MLP or convolutions. Especially, MLP and MLP mixer [94] also enjoy from global receptive fields, which is similar to transformers, so we believe that this comparison is necessary to justify our choice to select transformers over others. To this end, we compare the transformer aggregator with all the other cost aggregation methods, and the results are summarized in Table 5. For this experiment, we simply replace the self-attention module with other aggregator methods without touching any of our contributions. Note that, unlike other methods, 3D convolution is not compatible with our contributions as the different structure of the spatial axis makes them infeasible to directly replace the module, which we have to abandon some of the proposed components, *e.g.*, multi-level aggregation. However, we included 3D convolutions to evaluate the performance as a cost aggregator. Lastly, we adjusted the number of parameters for a fair comparison and report the memory consumption for different cost aggregators to indicate the efficiency.

From the results, using simple MLP yields surprisingly competitive results, thanks to its global cost aggregation, but its relatively poor performance highlights the importance of pairwise relationships, as MLP is only responsible for channel-mixing operation. This is similar to convolution, as MLP also lacks the adaptability to input pixels prevents from accurate matching. Using MLP-Mixer [94], which includes both channel and token mixing, yields highly competitive results as well, but compared to CATs, the performance is relatively poor, which the inadaptability can be one of the reasons. Given this, we argue that the transformer better considers the pairwise relationships than MLP-based aggregators and shows its superiority to better take global context in a correlation map into account. Now comparing the results to convolution-based aggregators [40], [95], we observe that the performance gap is quite large for 3D convolutions. This is because we abandoned some of our contributions for this experiment. The use of center-pivot 4D convolution surprisingly performs well, but compared to transformers, the gap seems large. With such differences, CATs obtaining better performance over other methods indicates its superiority and advantage to explicitly learn pairwise relationships with global attention range, process the input as a whole, and remove locality constraint which highlights the importance of our contribution and the effectiveness of transformer aggregator.

To compare the efficiency of different aggregators, we also note the memory consumption for each. Summarizing the results, we observe that memory consumption gap is quite large between

TABLE 6: **Ablation study of feature backbone of CATs.**

Feature Backbone	SPair-71k $\alpha_{\text{bbox}} = 0.1$	PF-PASCAL $\alpha_{\text{img}} = 0.1$
DeiT-B _{single} [66]	32.1	76.5
DeiT-B _{all} [66]	38.2	87.5
DINO w/ ViT-B/16 _{single} [97]	39.5	88.9
DINO w/ ViT-B/16 _{all} [97]	42.0	88.9
ResNet-101 _{single} [86]	37.4	87.3
ResNet-101 _{multi} [86]	42.4	89.1
ViT w/ CLIP _{single} [98]	36.8	85.3
ViT w/ CLIP _{all} [98]	47.9	89.0
ResNet-101 w/ CLIP _{single} [98]	30.8	77.2
ResNet-101 w/ CLIP _{all} [98]	45.3	89.4

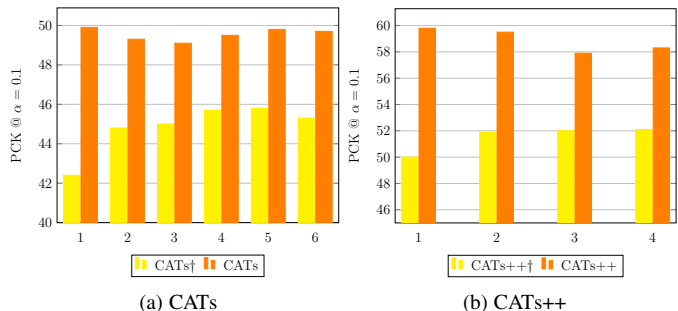


Fig. 11: **Effects of varying the number of encoders.**

3D convolution and center-pivot 4D convolution, but others show relatively similar memory consumption, and transformer requires the second lowest memory. Although center-pivot 4D convolutions shows the best efficiency thanks to its deliberate design, considering the large performance gap and relatively small memory consumption gap, we claim that the proposed transformer-based cost aggregator shows its superiority.

Additionally, we include the results of the model that employs Relative Positioning Bias (RPB) [85]. In this experiment, we aim to investigate the reasoning behind the performance difference between this work and that of VAT [58]. To this end, we include relative positioning bias within the self-attention computation and left other components untouched. In practice, we did not find notable performance changes for semantic matching benchmarks, but rather, we observed a notable performance boosts for FSS-1000 [96] dataset. Note that we also observed a similar tendency for CATs++, except that we observed a slight performance drop when evaluated on SPair-71k [16], confirming a relatively weaker influence of relative positioning bias on a semantic matching task. We thus find that the relative positioning bias can enhance the generalization power, which would explain some portion of performance gap between this work and VAT [58].

5.4.4 Does different feature backbone matter?

As shown in Table 6, we explore the impact of different feature backbones on the performance on SPair-71k [16] and PF-PASCAL [15]. We report the results with backbone networks frozen. For subscript *single*, we use a single-level feature, while for subscript *all*, every feature map from 12 layers is used for cost construction. For ResNet-101 subscript *multi*, we use the best layer subset provided by [21], while we use hypercorrelation approach for subscript *all*.

Summarizing the results, we consistently observe that leveraging multi-level features shows apparent improvements in performance, proving the effectiveness of multi-level aggregation. It is

TABLE 7: Effects of augmentation.

	Augment.	SPair-71k $\alpha_{\text{bbox}} = 0.1$	
DHPF [22]	✗	37.3	-
DHPF [22]	✓	39.4	-
CHMNet [40]	✗	-	47.0
CHMNet [40]	✓	-	51.3
CATs	✗	35.1	43.5
CATs	✓	42.4	49.9
CATs++	✗	47.3	54.0
CATs++	✓	50.0	59.8

worth noting that DINO, which is more excel at dense tasks than DeiT-B, outperforms DeiT-B when applied to semantic matching. This indicates that the learned representation has a large influence on the final performance. The interesting results we observed are the ones we obtained with CLIP-based backbone. Given ViT w/ CLIP_{all} or ResNet-101 w/ CLIP_{all} as backbone networks, we find that the performance is enhanced compared to the results with conventional ResNet-101. For ViT w/ CLIP_{all} we were surprised to observe such a significant performance boost. Even for PF-PASCAL, the performance boost occurs, confirming that the use of CLIP backbone shows an apparent performance boost.

5.4.5 Effects of varying the number of the encoders

As done in numerous works [27], [39], [42], [67] that utilize transformers, we can also stack more encoders to increase their capacity and validate the effectiveness by varying the number of encoders. We show the effects of varying the number of encoders in Figure 11. It is natural that choosing a higher number of encoders inevitably increases the memory consumption and run-time, but it obtains a larger model capacity in return. This is confirmed by the results when the number of encoder is set to 1, where both CATs[†] and CATs^{++†} report the lowest PCK. This indicates that the increase of the capacity of transformer aggregator clearly helps to boost the performance, which even surpasses the results reported in Table 1. Interestingly, showing completely contrary results, when the backbone networks are fine-tuned, the best PCKs are obtained when the number of encoders is set to 1. We suspect that trying to optimize both the backbone networks and the cost aggregator already risks overfitting, and by increasing the capacity of our cost aggregator, the model starts to overfit.

5.5 Analysis

5.5.1 Data augmentation

Transformer is well known for lacking some of the inductive bias and its data-hungry nature thus necessitates a large quantity of training data to be fed [41], [42]. Recent methods [66], [85], [99] that employ transformers to address Computer Vision tasks have empirically shown that data augmentation techniques have a positive impact on performance. Also, to compensate for low generalization power caused by the provision of sparsely annotated data, *i.e.*, keypoints, we provide a means to address it. Moreover, in the correspondence task, the question of to what extent can data augmentation affect performance has not yet been properly addressed. To this end, from the experiments, we empirically find that data augmentation has consistent positive impacts on performance in semantic correspondence.

TABLE 8: Comparison of serial/parallel processing.

	Serial	Parallel	SPair-71k $\alpha_{\text{bbox}} = 0.1$
CATs [†]	✓	✗	42.4
CATs [†]	✗	✓	38.3
CATs [†]	✓	✓	43.3
CATs	✓	✗	49.9
CATs	✗	✓	48.3
CATs	✓	✓	49.4
CATs ^{++†}	✓	✗	45.5
CATs ^{++†}	✗	✓	50.0
CATs ^{++†}	✓	✓	46.6
CATs ⁺⁺	✓	✗	54.8
CATs ⁺⁺	✗	✓	59.8
CATs ⁺⁺	✓	✓	55.9

TABLE 9: GPU memory and run-time comparison. Inference time for aggregator is denoted by (·) and subscript *coarse* represents the coarsest layer.

	Aggregation	Memory [GB]	Run-time [ms]
NC-Net [19]	4D Conv.	1.2	193.3 (166.1)
SCOT [23]	OT-RHM	4.6	146.5 (81.6)
DHPF [22]	RHM	<u>1.6</u>	57.7 (29.5)
CHM [28]	6D Conv	<u>1.6</u>	<u>47.2</u> (38.3)
CATs	Transformer	1.9	34.5 (7.4)
CATs ⁺⁺	4D Conv. + Transformer	3.1	110.2 (60.6)
CATs ^{++<i>coarse</i>}	4D Conv. + Transformer	1.2	57.4 (3.1)

In Table 7, we compared the PCK performance between our variants, DHPF [22] and CHMNet [40]. We note if the model is trained with augmentation. For a fair comparison, we evaluate all the methods trained on SPair-71k [16] using strong supervision, which assumes that the ground-truth keypoints are given. For CHMNet, we use the results reported in [40]. The results show that compared to DHPF and CHMNet, two typical examples of CNN-based cost aggregation methods, data augmentation technique has a larger influence on CATs in terms of performance. This demonstrates that not only we ease the data-hunger problem inherent in transformers, but also found that applying augmentations for matching has positive effects. Specifically, this technique allowed the performance boosts for DHPF and CHMNet by 2.1%p and 4.3%p, respectively, but it is further boosted for CATs by 6.4%p. This is further confirmed by the results of CATs⁺⁺, showing that the augmentation technique brings apparent performance improvements by 5.8%p.

5.5.2 Serial/Parallel processing

It is apparent that Equation 9 is not designed for an order-invariant output. Different from NC-Net [19], we design CATs in a way that we let the correlation map undergo the self-attention module in a serial manner. We conducted a simple experiment to compare the difference between each approach. Moreover, we also provide a quantitative results comparison of CATs⁺⁺ to justify our choice. For the experimental setting of serial processing, we sequentially aggregate the correlation map with a shared aggregator, while for parallel processing, we transpose the given input and let the original and transposed correlation maps undergo the same aggregator and add them subsequently. The results are summarized in Table 8. We argue that although CATs may not support order invariance, adopting serial processing can obtain higher PCK as it has a better capability to reduce inconsistent matching scores by additionally processing the already processed cost map, which we finalize the architecture to include serial processing. However,

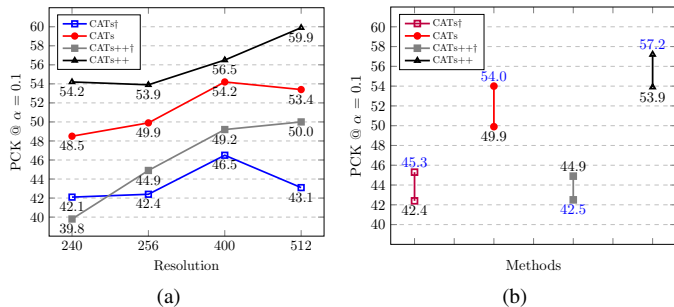


Fig. 12: **Ablation study of image resolution:** (a) as varying the image resolutions, (b) comparison between proposed methods at resolution set to 256. Note that we add results for CLIP-ResNet-101 [98], which is written as blue color.

interestingly, CATs++ clearly shows different results, which the parallel processing surpasses serial processing by a large margin. Also, we observe that employing both serial and parallel processing does not yield higher performance than parallel processing. We conjecture that whether the network makes use of convolutions yields different results. More specifically, while token-mixing convolutional operations mix both spatial dimensions of correlation maps, self-attention only considers either source or target spatial dimension, and this makes self-attention suitable to serial processing and convolutions suitable to parallel processing. Note that through parallel processing, CATs++ aggregates local contexts of both spatial dimensions and the following self-attention impart them to all pixels, which allows the best results. We thus argue that employing serial processing for CATs++ may complicate the learning process by performing redundant aggregations, while it can benefit CATs.

5.5.3 GPU memory consumption and run-time comparison

In Table 9, we show the memory and run-time comparison to NC-Net [19], SCOT [23], DHPF [22], and CHM [28]. For a fair comparison, the results are obtained using a single NVIDIA GeForce RTX 2080 Ti GPU and Intel Core i7-10700 CPU. We measure the GPU memory consumption given a single batch. We also measure the inference time for both the process without counting feature extraction and the whole process. Specifically, we note the type of aggregation method in the table, and the reported results represent the memory consumption and running-time of the aggregator, for instance, convolutional embedding module and efficient transformer aggregator for CATs++. Thanks to transformers’ fast computation nature, compared to other methods, CATs show much faster inference time, allowing real-time inference. This can benefit several applications which finding correspondences act as one of the milestones, *e.g.*, object tracking and video object segmentation [64]. Taking an example, an autonomous vehicles require real-time detection, tracking or segmentation of objects for a decision making, and a fast inference time can reduce delays during a decision making process, greatly reducing probability of potential incidents. Note that although CATs++ may suffer from relatively slower inference speed to other works [22], [28], less than 80 ms slower run time is a minor sacrifice for the better performance.

For the memory consumption, we find that compared to other cost aggregation methods including 4D, 6D convolutions, OT-RHM, and RHM, the transformer aggregator shows comparable efficiency in terms of computational cost. Note that NC-Net and

CHM show relatively lower memory consumption to others as they utilize a single feature map while all other methods utilize multi-level feature maps. Also, it is worth noting that CATs++ requires the largest memory consumption. This is because of the use of multi-level features as well as the processing of inputs at higher resolutions, which we show the reduced memory consumption and run-time at a coarse level of CATs++ to confirm this. Although CATs++ managed to reduce computational costs as summarized in Table 9, processing the correlation maps at higher resolutions inherently requires much larger costs in exchange for large performance gain. Further enhancing the efficiency by balancing with the performance would be a promising direction, which we leave as future work.

5.5.4 Resolution of input images

We find that different baseline methods [19], [21], [22], [23], [28], [33], [40], [45], [48], [50], [51] for semantic correspondence does not use the fixed input resolutions at resolution according to the implementations authors provide. This means that comparing the proposed method to other baseline methods that use different resolutions may make the comparison unfair. We argue that to make a fair comparison, the training resolutions can be different among works, but the evaluation resolutions should be the same. To this end, we provide the results of our proposed methods at different resolutions to make fair comparisons to previous methods first, then we suggest a benchmark for semantic correspondence task for future works.

As shown in Fig. 12 (a), we observe that the higher resolution generally yields higher performance. Specifically, at resolution 400, CATs yields the best results, while at resolution 512, CATs++ yields the best results. This shows that the higher resolution generally helps to gain higher performance. However, it can be seen at resolution 512 that the performance of CATs drops severely. We suspect that this is the limitation of cost aggregation with a standard transformer, which its performance can not infinitely scale as the resolution increases. As seen in Fig. 12 (a), we argue that direct comparison is possible only when the results for different input sizes are given. Now taking this into account, we suggest a fixed benchmark for future works in semantic correspondence task.

We suggest the evaluation sizes be consistent across all the works for a fair comparison, which we decide as 256, and the results of the proposed methods at 256 evaluation resolution are shown in Fig. 12 (b). As shown, the performance of CATs++ dropped by a quite large margin. More specifically, we observe 6.0% drop in performance for CATs++ when the evaluation size is changed from 512 to 256. This should make a fair comparison to other works evaluated at 256. Now we introduce a new finding to narrow this gap, which we adopt a new feature backbone, ResNet-101 initialized by pre-trained weights released by CLIP [98]. We were surprised to observe that by simply replacing the conventional ResNet-101 pre-trained on ImageNet [87] to that of CLIP, apparent performance boosts can be made. This can be a valuable finding that could bring an apparent performance boost even for other tasks, thanks to the improved training scheme and datasets by CLIP.

5.5.5 Capturing fine details

As an extension to CATs, we proposed CATs++, a novel cost aggregation approach that enjoys both reduced costs and boosted performance, and composed to have a hierarchical architecture,

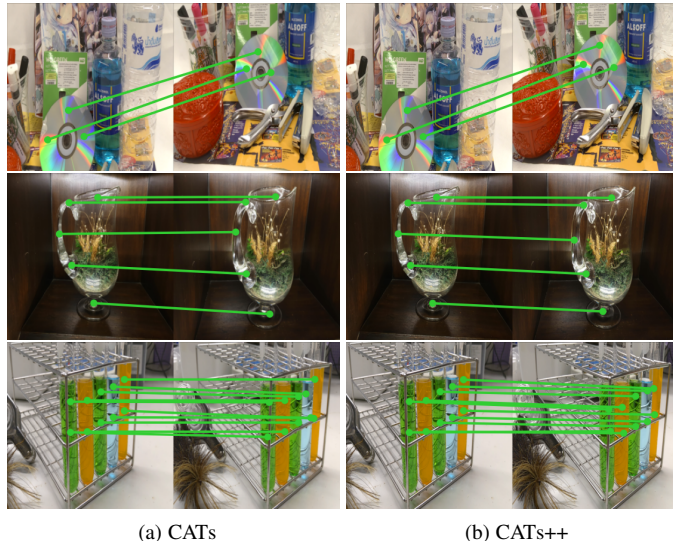


Fig. 13: **Qualitative results on non-lambertian objects.** The proposed methods can obtain accurate correspondence fields between non-lambertian objects.

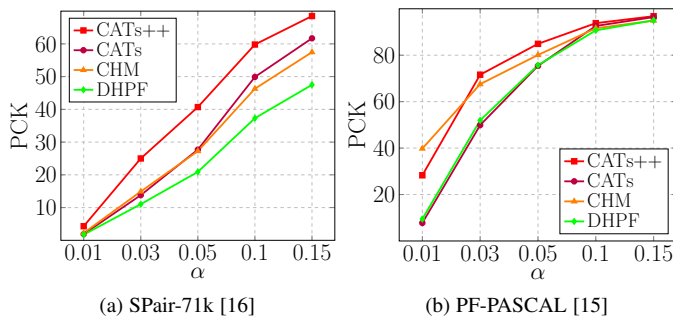


Fig. 14: **Quantitative results as varying the alpha threshold on SPair-71k [16] and PF-PASCAL [15].** CATs++ captures fine-details surprisingly well compared to other methods (including CATs), thanks to its hierarchical architecture.

which as a result, it allowed the model to excel at capturing fine-details as shown in Fig. 14. As shown in Fig. 14, where we evaluate DHPF [22], CHM [28], CATs and CATs++ on SPair-71k [16] and PF-PASCAL [15] by varying the alpha threshold used for the PCK evaluation, thanks to its hierarchical architecture, CATs++ successfully captures fine details, accurately finding the matching points, which CATs struggles to find. Note that CATs++ consistently outperforms all other methods except one case where CHM outperforms at $\alpha = 0.01$ for PF-PASCAL. Interestingly, when the α threshold gets lower, CATs perform worse than other methods, indicating that it struggles to capture the fine details, which CATs++ successfully overcomes this limitation.

5.5.6 Extension to few-shot segmentation

We also show that the proposed methods perform well at a rudimentary level even for few-shot segmentation task. Few-shot segmentation aims to reduce reliance on labeled data and only a few support images with their associated masks are given for obtaining the segmentation maps for a query image. For this experiment, we simply replace soft-argmax with a 2D decoder as done in [3] to produce a segmentation map and we show that CATs and CATs++ both also achieve competitive results on FSS-1000 [96]. This is because few-shot segmentation shares a similar

TABLE 10: **Quantitative comparison on FSS-1000 [96].**

Backbone feature	Methods	mIoU		Memory [GB]	# of learnable params.
		1-shot	5-shot		
ResNet50 [86]	FSOT [101]	82.5	83.8	-	-
	HSNet [3]	<u>85.5</u>	<u>87.8</u>	-	2.6M
	VAT [58]	90.1	90.7	-	3.2M
	CATs++	85.2	85.4	-	5.5M
ResNet101 [86]	DAN [100]	85.2	88.1	-	-
	HSNet [3]	86.5	88.5	2.2	2.6M
	VAT [58]	90.3	90.8	3.8	3.3M
	CATs	80.8	80.9	1.9	4.6M
	CATs++	85.1	85.3	2.6	5.5M

condition to semantic correspondence task, as it’s also challenged by intra-class variations and background clutter.

In Table 10, we summarize the quantitative results on FSS-1000 [96] dataset. FSS-1000 [96], a dataset specifically designed for few-shot segmentation, consists of 1000 object classes. Following [96], 1000 classes are categorized into 3 splits for training, validation and testing, which consist of 520, 240 and 240 classes, respectively. For the evaluation metric, we employ mean intersection over union (mIoU), which averages over all IoU values for all object classes.

We observed that CATs++ performs almost on par with other methods [3], [100], [101] that are specifically designed for few-shot segmentation task. However, we find that the overall performance is generally lower than those of other works, especially when compared to VAT [58]. The reason for the performance gap can be explained by the difference in generalization power. As CATs only adopts Transformer as its aggregator, it may suffer from lower generalization power caused by the lack of inductive bias. While other works, including HSNet [3], VAT [58] and CATs++, avoid this by utilizing convolutions for translation equivariance or adopting Swin Transformer [85] that benefits from relative positioning bias. Although HSNet [3] and CATs++ benefit from convolutional inductive bias just like VAT [58] does, the best performance by VAT [58] can be explained by the joint use of both convolutions and Swin Transformer [85], which allowed stronger generalization power compared to HSNet [3] and CATs++ that only utilize convolutions within their architectures. Nevertheless, our approach showed that it is competent at both tasks despite its relatively lacking generalization power, and even outperforms other works in semantic matching.

Finally, comparing the memory consumption and the learnable parameters to HSNet [3] and VAT [58], we observe that CATs requires much less memory consumption than all other approaches. It should also be noted that CATs++ is much more lightweight than VAT [58], which clearly demonstrates its efficiency. However, we observe that the number of learnable parameters for CATs and CATs++ are larger than other works, which is known to have an inverse relation to generalization power [102], which may explain the performance gap when the task is to predict segmentation maps of unseen classes.

5.5.7 Extension to Visual Localization

Finally, we show that our approach also benefits the task of visual localization and evaluate on Aachen dataset [106]. Visual localization aims to determine location from images by estimating the absolute 6 DoF pose of a query image with respect to the corresponding 3D scene model. To evaluate, we submit the results to the online evaluation benchmark [106].

Note that the proposed approach outputs a dense flow field, which makes it require non-trivial implementations and strategies

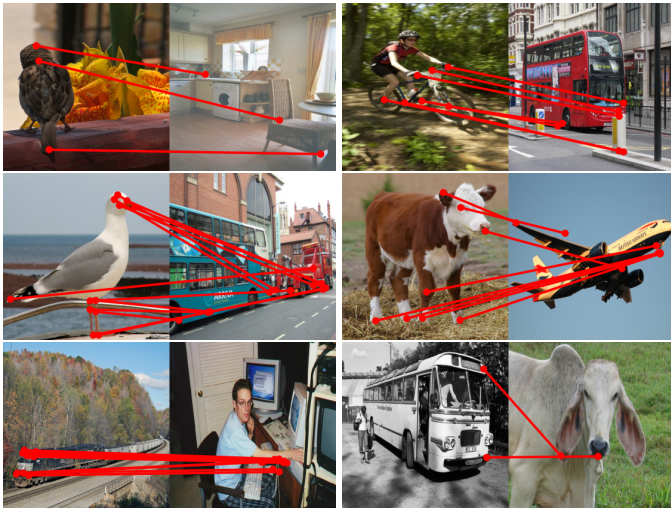


Fig. 15: **Failure cases.** Given images with irrelevant objects, the proposed method may struggle.

to apply to SfM 3D reconstruction without hurting the performance. To alleviate this issue, we employ SuperPoint [32] to be responsible for the keypoint detection stage, and we then follow Hloc [107] outdoor localization pipeline to obtain the final localization results. To train our networks, we employ large-scale outdoor dataset, MegaDepth [108]. Given output flow maps from our networks, we up-sample the flow map to original image size and utilize keypoint coordinates already obtained from SuperPoint [32] to find the corresponding keypoint at the other image.

The results are summarized in Table 11. We observe that the proposed approach struggles to find accurate correspondences given high-resolution images, failing to localize with fine-details. Unlike other works, which include PDCNet [50], DualRC-Net [59] and SuperGlue [25], our approach outputs the flow maps at relatively low resolution, *i.e.*, 32×32 . This prevents from capturing fine-details, an important aspect for the visual localization task. Nevertheless, although a new state-of-the-art is not attained, a promising future direction is revealed, which is to capture fine details. As a future direction, a module that captures the missing details of the predicted flow can be designed. For example, local cropping around the coarse correspondence may be exploited to find fine details.

5.5.8 Limitations

An apparent limitation of the proposed method is that it directly employs correlation maps, which are computationally expensive to process. This limits the resolutions of the input images the model

Method type	Method	0.5m, 2°	0.5m, 5°	5m, 10°
Sparse	D2-Net [103]	74.5	86.7	100.0
	R2D2 [104]	69.4	86.7	94.9
	R2D2 [104] (K=20k)	76.5	90.8	100.0
	SuperPoint [32]	73.5	79.6	88.8
	SuperPoint [32] + SuperGlue [25]	79.6	90.8	100.0
	Patch2Pix	78.6	88.8	99.0
Dense-to-sparse	Sparse-NCNet [37]	76.5	84.7	98.0
	DualRC-Net [59]	79.6	88.8	100.0
Dense	RANSAC-flow [105] + Superpoint [50]	74.5	87.8	100.0
	PDC-Net [50]	76.5	85.7	100.0
	PDC-Net [50] + SuperPoint [32]	80.6	87.8	100.0
	CATs++ + SuperPoint [32]	62.2	73.5	93.9

TABLE 11: **Visual localization on the Aachen day-night dataset [106].**

can handle, enforcing the model to down-sample, then find the correspondences. Moreover, the proposed approaches assume that a pair of images with common objects, *i.e.*, semantically similar, are given as inputs, making its applicability limited to only when the pair of images with semantically similar objects are given, as shown in Fig. 15. Although such pairs can be obtained via image retrieval process, this may be one of the apparent limitations.

6 CONCLUSION

In this paper, we have proposed, for the first time, transformer-based cost aggregation networks for semantic correspondence which enables aggregating the matching scores computed between input features, dubbed CATs. We have made several architectural designs in the network architecture, including appearance affinity modelling, multi-level aggregation, swapping self-attention, and residual connection. We have shown that our method surpasses the current state-of-the-art in several benchmarks. Moreover, we extended CATs by introducing early convolutions and efficient transformer aggregator to reduce the computation costs and boost the performance, dubbed CATs++. We demonstrated the effectiveness of the proposed method on several benchmarks, which our method outperforms current state-of-the-art with large margin. We also conducted extensive ablation studies to validate our choices and explore its capacity and introduced expansion to other task, few-shot segmentation.

Acknowledgements. This research was supported by the MSIT, Korea (IITP-2022-2020-0-01819, ICT Creative Consilience program), and National Research Foundation of Korea (NRF-2021R1C1C1006897).

REFERENCES

- [1] M. Rubinstein, A. Joulin, J. Kopf, and C. Liu, “Unsupervised joint object discovery and segmentation in internet images,” in *CVPR*, 2013.
- [2] T. Tanai, S. N. Sinha, and Y. Sato, “Joint recovery of dense correspondence and cosegmentation in two images,” in *CVPR*, 2016.
- [3] J. Min, D. Kang, and M. Cho, “Hypercorrelation squeeze for few-shot segmentation,” *arXiv preprint arXiv:2104.01538*, 2021.
- [4] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, “Feature pyramid networks for object detection,” in *CVPR*, 2017.
- [5] R. Szeliski, “Image alignment and stitching: A tutorial,” *Foundations and Trends® in Computer Graphics and Vision*, 2006.
- [6] C. Liu, J. Yuen, and A. Torralba, “Sift flow: Dense correspondence across scenes and its applications,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 978–994, 2010.
- [7] J. Liao, Y. Yao, L. Yuan, G. Hua, and S. B. Kang, “Visual attribute transfer through deep image analogy,” *arXiv:1705.01088*, 2017.
- [8] S. Kim, D. Min, S. Jeong, S. Kim, S. Jeon, and K. Sohn, “Semantic attribute matching networks,” in *CVPR*, 2019.
- [9] W. Peebles, J.-Y. Zhu, R. Zhang, A. Torralba, A. Efros, and E. Shechtman, “Gan-supervised dense visual alignment,” *arXiv preprint arXiv:2112.05143*, 2021.
- [10] A. Hosni, C. Rhemann, M. Bleyer, C. Rother, and M. Gelautz, “Fast cost-volume filtering for visual correspondence and beyond,” *PAMI*, 2012.
- [11] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, “FlowNet 2.0: Evolution of optical flow estimation with deep networks,” in *CVPR*, 2017.
- [12] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, “Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume,” in *CVPR*, 2018.
- [13] T.-W. Hui, X. Tang, and C. Change Loy, “Liteflownet: A lightweight convolutional neural network for optical flow estimation,” in *CVPR*, 2018.
- [14] B. Ham, M. Cho, C. Schmid, and J. Ponce, “Proposal flow,” in *CVPR*, 2016.
- [15] —, “Proposal flow: Semantic correspondences from object proposals,” *IEEE transactions on pattern analysis and machine intelligence*, 2017.

- [16] J. Min, J. Lee, J. Ponce, and M. Cho, "Spair-71k: A large-scale benchmark for semantic correspondence," *arXiv preprint arXiv:1908.10543*, 2019.
- [17] I. Rocco, R. Arandjelovic, and J. Sivic, "Convolutional neural network architecture for geometric matching," in *CVPR*, 2017.
- [18] I. Rocco, R. Arandjelović, and J. Sivic, "End-to-end weakly-supervised semantic alignment," in *CVPR*, 2018.
- [19] I. Rocco, M. Cimpoi, R. Arandjelović, A. Torii, T. Pajdla, and J. Sivic, "Neighbourhood consensus networks," in *NeurIPS*, 2018.
- [20] I. Melekhov, A. Tiulpin, T. Sattler, M. Pollefeys, E. Rahtu, and J. Kannala, "Dgc-net: Dense geometric correspondence network," in *WACV*, 2019.
- [21] J. Min, J. Lee, J. Ponce, and M. Cho, "Hyperpixel flow: Semantic correspondence with multi-layer neural features," in *ICCV*, 2019.
- [22] —, "Learning to compose hypercolumns for visual correspondence," in *ECCV*, 2020.
- [23] Y. Liu, L. Zhu, M. Yamada, and Y. Yang, "Semantic correspondence as an optimal transport problem," in *CVPR*, 2020.
- [24] P. Truong, M. Danelljan, and R. Timofte, "Glu-net: Global-local universal network for dense flow and correspondences," in *CVPR*, 2020.
- [25] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superglue: Learning feature matching with graph neural networks," in *CVPR*, 2020.
- [26] P. Truong, M. Danelljan, L. V. Gool, and R. Timofte, "Gocor: Bringing globally optimized correspondence volumes into your neural network," in *NeurIPS*, 2020.
- [27] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, "Loft: Detector-free local feature matching with transformers," *arXiv preprint arXiv:2104.00680*, 2021.
- [28] J. Min and M. Cho, "Convolutional hough matching networks," *arXiv preprint arXiv:2103.16831*, 2021.
- [29] D. Scharstein and R. Szeliski, "A taxonomy and evaluation of dense two-frame stereo correspondence algorithms," *International journal of computer vision*, 2002.
- [30] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *CVPR*. IEEE, 2007.
- [31] S. Kim, D. Min, B. Ham, S. Jeon, S. Lin, and K. Sohn, "Fcsc: Fully convolutional self-similarity for dense semantic correspondence," in *CVPR*, 2017.
- [32] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *CVPR*, 2018.
- [33] J. Lee, D. Kim, J. Ponce, and B. Ham, "Sfnet: Learning object-aware semantic correspondence," in *CVPR*, 2019.
- [34] P. H. Seo, J. Lee, D. Jung, B. Han, and M. Cho, "Attentive semantic alignment with offset-aware correlation kernels," in *ECCV*, 2018.
- [35] S. Kim, S. Lin, S. R. Jeon, D. Min, and K. Sohn, "Recurrent transformer networks for semantic correspondence," in *NeurIPS*, 2018.
- [36] S. Jeon, S. Kim, D. Min, and K. Sohn, "Parn: Pyramidal affine regression networks for dense semantic correspondence," in *ECCV*, 2018.
- [37] I. Rocco, R. Arandjelović, and J. Sivic, "Efficient neighbourhood consensus networks via submanifold sparse convolutions," in *ECCV*, 2020.
- [38] S. Jeon, D. Min, S. Kim, J. Choe, and K. Sohn, "Guided semantic flow," in *ECCV*. Springer, 2020.
- [39] S. Li, K. Han, T. W. Costain, H. Howard-Jenkins, and V. Prisacariu, "Correspondence networks with adaptive neighbourhood consensus," in *CVPR*, 2020.
- [40] J. Min, S. Kim, and M. Cho, "Convolutional hough matching networks for robust and efficient visual correspondence," *arXiv preprint arXiv:2109.05221*, 2021.
- [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [42] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [43] K. Zhang, Y. Fang, D. Min, L. Sun, S. Yang, S. Yan, and Q. Tian, "Cross-scale cost aggregation for stereo matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1590–1597.
- [44] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach, and A. Bry, "End-to-end learning of geometry and context for deep stereo regression," in *ICCV*, 2017.
- [45] S. Cho, S. Hong, S. Jeon, Y. Lee, K. Sohn, and S. Kim, "Cats: Cost aggregation transformers for visual correspondence," in *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- [46] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *CVPR Workshops*, 2005.
- [47] C. B. Choy, J. Gwak, S. Savarese, and M. Chandraker, "Universal correspondence network," *NeurIPS*, vol. 29, pp. 2414–2422, 2016.
- [48] D. Zhao, Z. Song, Z. Ji, G. Zhao, W. Ge, and Y. Yu, "Multi-scale matching networks for semantic correspondence," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3354–3364.
- [49] S. Hong and S. Kim, "Deep matching prior: Test-time optimization for dense correspondence," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 9907–9917.
- [50] P. Truong, M. Danelljan, L. Van Gool, and R. Timofte, "Learning accurate dense correspondences and when to trust them," *arXiv preprint arXiv:2101.01710*, 2021.
- [51] J. Y. Lee, J. DeGol, V. Frago, and S. N. Sinha, "Patchmatch-based neighborhood consensus for semantic correspondence," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13 153–13 163.
- [52] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph.*, 2009.
- [53] S. Gao, C. Zhou, C. Ma, X. Wang, and J. Yuan, "Aiatrack: Attention in attention for transformer visual tracking," 2022.
- [54] B. Yan, Y. Jiang, P. Sun, D. Wang, Z. Yuan, P. Luo, and H. Lu, "Towards grand unification of object tracking," *arXiv preprint arXiv:2207.07078*, 2022.
- [55] X. Wang, A. Jabri, and A. A. Efros, "Learning correspondence from the cycle-consistency of time," in *CVPR*, 2019.
- [56] Z. Lai and W. Xie, "Self-supervised learning for video correspondence flow," *arXiv preprint arXiv:1905.00875*, 2019.
- [57] A. Jabri, A. Owens, and A. Efros, "Space-time correspondence as a contrastive random walk," *NeurIPS*, 2020.
- [58] S. Hong, S. Cho, J. Nam, S. Lin, and S. Kim, "Cost aggregation with 4d convolutional swin transformer for few-shot segmentation," *arXiv preprint arXiv:2207.10866*, 2022.
- [59] X. Li, K. Han, S. Li, and V. A. Prisacariu, "Dual-resolution correspondence networks," *arXiv:2006.08844*, 2020.
- [60] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.
- [61] F. Tan, J. Yuan, and V. Ordonez, "Instance-level image retrieval using reranking transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12 105–12 115.
- [62] S. Lee, H. Seong, S. Lee, and E. Kim, "Correlation verification for image retrieval," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 5374–5384.
- [63] A. Yilmaz, O. Javed, and M. Shah, "Object tracking: A survey," *Acm computing surveys (CSUR)*, vol. 38, no. 4, pp. 13–es, 2006.
- [64] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung, "A benchmark dataset and evaluation methodology for video object segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 724–732.
- [65] N. Dong and E. P. Xing, "Few-shot semantic segmentation with prototype learning," in *BMVC*, vol. 3, no. 4, 2018.
- [66] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers and distillation through attention," *arXiv preprint arXiv:2012.12877*, 2020.
- [67] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *ECCV*. Springer, 2020.
- [68] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.
- [69] P. Sun, Y. Jiang, R. Zhang, E. Xie, J. Cao, X. Hu, T. Kong, Z. Yuan, C. Wang, and P. Luo, "Transtrack: Multiple-object tracking with transformer," *arXiv preprint arXiv:2012.15460*, 2020.
- [70] M. J. Tyszkiewicz, P. Fua, and E. Trulls, "Disk: Learning local features with policy gradient," *arXiv preprint arXiv:2006.13566*, 2020.
- [71] W. Jiang, E. Trulls, J. Hosang, A. Tagliasacchi, and K. M. Yi, "COTR: Correspondence Transformer for Matching Across Images," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021.

- [72] A. F. Agarap, “Deep learning using rectified linear units (relu),” *arXiv preprint arXiv:1803.08375*, 2018.
- [73] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [74] D. Hendrycks and K. Gimpel, “Gaussian error linear units (gelus),” *arXiv preprint arXiv:1606.08415*, 2016.
- [75] D. Marin, J.-H. R. Chang, A. Ranjan, A. Prabhu, M. Rastegari, and O. Tuzel, “Token pooling in vision transformers,” *arXiv preprint arXiv:2110.03860*, 2021.
- [76] N. Kitaev, Ł. Kaiser, and A. Levskaya, “Reformer: The efficient transformer,” *arXiv preprint arXiv:2001.04451*, 2020.
- [77] S. Wang, B. Z. Li, M. Khabsa, H. Fang, and H. Ma, “Linformer: Self-attention with linear complexity,” *arXiv preprint arXiv:2006.04768*, 2020.
- [78] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, “Transformers are rnns: Fast autoregressive transformers with linear attention,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 5156–5165.
- [79] C. Wu, F. Wu, T. Qi, Y. Huang, and X. Xie, “Fastformer: Additive attention can be all you need,” *arXiv preprint arXiv:2108.09084*, 2021.
- [80] J. Yang, C. Li, P. Zhang, X. Dai, B. Xiao, L. Yuan, and J. Gao, “Focal self-attention for local-global interactions in vision transformers,” *arXiv preprint arXiv:2107.00641*, 2021.
- [81] Y. Rao, W. Zhao, B. Liu, J. Lu, J. Zhou, and C.-J. Hsieh, “Dynamicvit: Efficient vision transformers with dynamic token sparsification,” *arXiv preprint arXiv:2106.02034*, 2021.
- [82] S. Jaszczur, A. Chowdhery, A. Mohiuddin, Ł. Kaiser, W. Gajewski, H. Michalewski, and J. Kanerva, “Sparse is enough in scaling transformers,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [83] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, “Cvt: Introducing convolutions to vision transformers,” *arXiv preprint arXiv:2103.15808*, 2021.
- [84] Y. Xiong, Z. Zeng, R. Chakraborty, M. Tan, G. Fung, Y. Li, and V. Singh, “Nyströmformer: A nyström-based algorithm for approximating self-attention,” in *Proceedings of the... AAAI Conference on Artificial Intelligence*. AAAI Conference on Artificial Intelligence, vol. 35, no. 16. NIH Public Access, 2021, p. 14138.
- [85] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” *arXiv preprint arXiv:2103.14030*, 2021.
- [86] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016.
- [87] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [88] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *arXiv:1711.05101*, 2017.
- [89] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, “Randaugment: Practical automated data augmentation with a reduced search space,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [90] A. Buslaev, V. I. Iglovikov, E. Khvedchenya, A. Parinov, M. Druzhinin, and A. A. Kalinin, “Albumentations: fast and flexible image augmentations,” *Information*, 2020.
- [91] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch,” 2017.
- [92] S. Huang, Q. Wang, S. Zhang, S. Yan, and X. He, “Dynamic context correspondence network for semantic alignment,” in *ICCV*, 2019.
- [93] K. Han, R. S. Rezende, B. Ham, K.-Y. K. Wong, M. Cho, C. Schmid, and J. Ponce, “Scnet: Learning semantic correspondence,” in *ICCV*, 2017.
- [94] I. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, D. Keysers, J. Uszkoreit, M. Lucic, and A. Dosovitskiy, “Mlp-mixer: An all-mlp architecture for vision,” 2021.
- [95] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, and M. Paluri, “A closer look at spatiotemporal convolutions for action recognition,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 6450–6459.
- [96] X. Li, T. Wei, Y. P. Chen, Y.-W. Tai, and C.-K. Tang, “Fss-1000: A 1000-class dataset for few-shot segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [97] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging properties in self-supervised vision transformers,” *arXiv preprint arXiv:2104.14294*, 2021.
- [98] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” *arXiv preprint arXiv:2103.00020*, 2021.
- [99] H. Touvron, M. Cord, A. Sablayrolles, G. Synnaeve, and H. Jégou, “Going deeper with image transformers,” *arXiv preprint arXiv:2103.17239*, 2021.
- [100] H. Wang, X. Zhang, Y. Hu, Y. Yang, X. Cao, and X. Zhen, “Few-shot semantic segmentation with democratic attention networks,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*. Springer, 2020, pp. 730–746.
- [101] W. Liu, C. Zhang, H. Ding, T.-Y. Hung, and G. Lin, “Few-shot segmentation with optimal transport matching and message flow,” *arXiv preprint arXiv:2108.08518*, 2021.
- [102] I. V. Tetko, D. J. Livingstone, and A. I. Luik, “Neural network studies, 1. comparison of overfitting and overtraining,” *J. Chem. Inf. Comput. Sci.*, vol. 35, pp. 826–833, 1995.
- [103] M. Dusmanu, I. Rocco, T. Pajdla, M. Pollefeys, J. Sivic, A. Torii, and T. Sattler, “D2-net: A trainable cnn for joint description and detection of local features,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 8092–8101.
- [104] J. Revaud, P. Weinzaepfel, C. De Souza, N. Pion, G. Csurka, Y. Cabon, and M. Humenberger, “R2d2: repeatable and reliable detector and descriptor,” *arXiv preprint arXiv:1906.06195*, 2019.
- [105] X. Shen, F. Darmon, A. A. Efros, and M. Aubry, “Ransac-flow: generic two-stage image alignment,” *arXiv:2004.01526*, 2020.
- [106] T. Sattler, W. Maddern, C. Toft, A. Torii, L. Hammarstrand, E. Stenborg, D. Safari, M. Okutomi, M. Pollefeys, J. Sivic *et al.*, “Benchmarking 6dof outdoor visual localization in changing conditions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8601–8610.
- [107] P.-E. Sarlin, C. Cadena, R. Siegwart, and M. Dymczyk, “From coarse to fine: Robust hierarchical localization at large scale,” in *CVPR*, 2019.
- [108] Z. Li and N. Snavely, “Megadepth: Learning single-view depth prediction from internet photos,” in *Computer Vision and Pattern Recognition (CVPR)*, 2018.



Seokju Cho received his B.S. degree in Computer Science from Yonsei University in 2022, currently pursuing a Ph.D. degree at Korea University. His primary research interest is learning visual correspondences and its applications such as few-shot semantic segmentation and 3D reconstruction.



Sunghwan Hong received his BS degree in Computer Sciences and Engineering from Korea University in 2021, where he is currently pursuing his PhD degree. His current research focuses on visual correspondences and its applications such as few-shot learning and visual localization. He is also interested in 3D vision, multi-modal learning, NeRF and talking head generations.



Seungryong Kim received the B.S. and Ph.D. degrees from the School of Electrical and Electronic Engineering from Yonsei University, Seoul, Korea, in 2012 and 2018, respectively. From 2018 to 2019, he was Post-Doctoral Researcher in Yonsei University, Seoul, Korea. From 2019 to 2020, he has been Post-Doctoral Researcher in School of Computer and Communication Sciences at École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland. Since 2020, he has been an assistant professor with the Department of Computer Science and Engineering, Korea University, Seoul. His current research interests include 2D/3D computer vision, computational photography, and machine learning.