

TransFlow: Transformer as Flow Learner

Yawen Lu^{1,2}, Qifan Wang³, Siqi Ma⁴, Tong Geng⁵, Yingjie Victor Chen¹, Huaijin Chen⁶, and Dongfang Liu^{2*}

¹Purdue University

²Rochester Institute of Technology

³Meta AI

⁴Westlake University

⁵University of Rochester

⁶Vayu Robotics

Abstract

Optical flow is an indispensable building block for various important computer vision tasks, including motion estimation, object tracking, and disparity measurement. In this work, we propose TransFlow, a pure transformer architecture for optical flow estimation. Compared to dominant CNN-based methods, TransFlow demonstrates three advantages. First, it provides more accurate correlation and trustworthy matching in flow estimation by utilizing spatial self-attention and cross-attention mechanisms between adjacent frames to effectively capture global dependencies; Second, it recovers more compromised information (e.g., occlusion and motion blur) in flow estimation through long-range temporal association in dynamic scenes; Third, it enables a concise self-learning paradigm and effectively eliminate the complex and laborious multi-stage pre-training procedures. We achieve the state-of-the-art results on the Sintel, KITTI-15, as well as several downstream tasks, including video object detection, interpolation and stabilization. For its efficacy, we hope TransFlow could serve as a flexible baseline for optical flow estimation.

1. Introduction

With the renaissance of connectionism, rapid progress has been made in optical flow. Till now, most of the state-of-the-art flow learners were built upon Convolutional Neural Networks (CNNs) [5, 7, 16, 32, 40, 44]. Despite their diversified model designs and tantalizing results, existing CNN-based methods commonly rely on

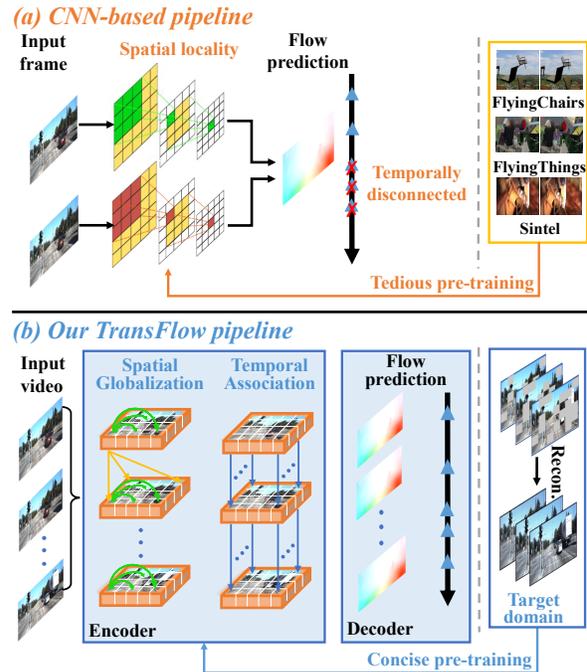


Figure 1. **Conceptual comparison of flow estimation methods.** Existing CNN-based methods regress flow via local spatial convolutions, while TransFlow relies on *Transformer* to perform *global matching* (both *spatial* and *temporal*). The demo video can be found in <https://youtu.be/xbnyj9wspqA>.

spatial locality and compute displacements from all-pair correlation volume for flow prediction (Fig. 1(a)). Very recently, the vast success of Transformer [42, 47, 48]

stimulates the emergence of attention-based paradigms for various tasks in language, vision, speech, and more. It appears to form a unanimous endeavor in the deep learning community to develop unified methodologies for solving problems in different areas. Towards unifying methodologies, less inductive biases [42] are introduced for a specific problem, which urges the models to learn useful knowledge purely from the input data.

Jumping on the bandwagon of unifying architecture, we study applying Vision Transformer [13] to the task of optical flow. The following question naturally arises: *What are the major limitations of existing CNN-based approaches?* Tackling this question can provide insights into the model design of optical flow, and motivate us to rethink the task from an attention-driven view. First, the concurrent CNN-based methods demonstrate inefficiency in modeling *global spatial dependencies* due to the intrinsic locality of the convolution operation. It usually requires a large number of CNN layers to capture the correlations between two pixels that are spatially far away. Second, CNN-based flow learners generally model the flow between only two consecutive frames, and fail to explore *temporal associations* in the neighboring contexts, resulting in weak prediction in the presence of significant photometric and geometric changes. Third, the existing training strategy usually requires a *tedious pipeline*. Performance guarantees heavily rely on excessive pre-training on extra datasets (e.g., FlyingChairs [14], FlyingThings [33], etc). Without adequate pre-training procedures, the model typically converges with large errors.

In order to craft a Transformer architecture for optical flow that pursues performance guarantees, the question becomes more fundamental: *How to address these limitations using Transformer?* As responses to this question, we articulate the technical contributions to address each of the above limitations:

- We introduce *spatial attention* mechanism that effectively captures global dependencies and achieves precise correlation and reliable matching for flow estimation. Essentially, the spatial attention in Transformer enables effective contextual cue propagation from coherent regions to the surroundings with heavy-tailed noise, motion blurs, and large displacements, which significantly prevents performance degradation in flow estimation.
- We explicitly model *temporal association* in dynamic scenes using multi-frame features extracted in the designed Transformer encoder. The correspondences among different frames are learned to generate the final estimated flow. One advantage

is that when a region in a frame is occluded or blurred, neighboring frames can effectively recover the missing information based on the learned temporal association.

- We design a concise *self-supervised pre-training* module that effectively eliminates the complex and laborious multi-stage pre-training procedures. In particular, extended from MAE [17], we develop a masking strategy during the training to adaptively mask out visual tokens and learn strong pixel representations by reconstructing clean signals from corrupted inputs. We demonstrate that the simple architecture results in a powerful entrance model, achieving stronger performance compared with SOTA baselines [19, 23, 32, 41, 60].

In summary, we re-formalize typical optical flow estimation within a *pure transformer architecture* — *TransFlow* (Fig. 1(b)), which factorises pixel-wise flow learners with spatial dependencies and temporal associations to increase performance guarantees. Remaining of the work is organized as follows: §2 provides literature review on the concurrent flow estimation methods. §3 describes the model architecture of TransFlow. In §4, we detail the configuration setup and experimental settings. Concretely, §4.1, shows that our method achieves impressive results in popular flow estimation datasets (e.g. Sintel [4] and KITTI 2015 [34]) and outperforms recent leading approaches; In §4.2, with a set of diagnostic experiment, our extensive experimental settings verify the effectiveness of our method. In §4.3, we demonstrate the transferability and generalizability of our method for modeling object motion for downstream tasks (i.e., video object detection, video interpolation, and video stabilization), which can benefit from our method without bells and whistles. In the end, we make conclusions in §5 to highlight that this work is expected to pave the way for future research in this area.

2. Important Knowns and Gap

Optical Flow Learners. Traditionally, optical flow was formulated as an energy optimization problem for maximizing similarity between image pairs [2, 3, 39]. More recently, visual similarity is computed via the computationally expensive correlation of high-dimensional features encoded by convolutional neural networks [14, 18, 19, 21, 40, 41]. FlowNet [14] was the first end-to-end CNN-based network, which uses a coarse and refined branch for optical flow estimation. Its successive work, FlowNet2.0 [21], adopted a stacked architecture with warping operation, leading to improved performance. Following the coarse-to-fine strategy, PWC-Net [40] de-

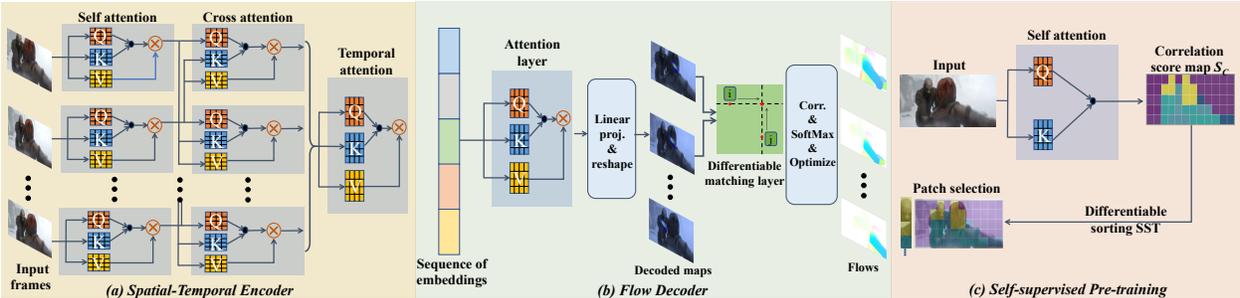


Figure 2. **Overall model architecture of TransFlow.** It consists of three major components, a spatial-temporal encoder, a flow decoder, and a self-supervised pre-training module. The spatial-temporal encoder jointly performs spatial globalization and temporal association among patch tokens. The flow decoder decodes the feature maps for multiple frames and generates the final optical flow. The pre-training module is designed to learn effective image representation in a self-supervised manner.

veloped a framework composed of stacked image pyramids, image warping and cost volumes; Hofinger et al., [18] replaced the image warping with a sampling-based strategy to improve the cost volume construction; Teed and Deng et al., [41] proposed to build a 4D cost volume for matching between all pairs of pixels and added a recurrent decoder for propagation. However, feature maps generated in these methods are usually suffering from a limited receptive field and high susceptibility to outliers, making them unsuitable as effective structures for learning global motion clues. Essentially, our method is conceptually different from these pioneer arts. We design a transformer structure that takes advantage of both self- and cross- attentions and temporal association for effective global matching. Moreover, we demonstrate it is possible to achieve competitive results without costly training pipelines, using self-supervised learning.

Attention in Optical Flow. While becoming the standard for natural language processing tasks [10, 12], a flurry of research has successfully introduced Transformers to computer vision. Inspired by successes in image classification [13, 57], multiple recent architectures have been trying to combine CNN-based architectures with self-attention, including detection [6, 9], image restoration [49], video inpainting [29] and flow estimation [38]. Recently, there have been several attempts to apply Transformer structures to boost the performance of optical flow. Generally for these works, attention is applied in tandem with CNNs to compensate for the absence of image-specific inductive bias [29, 50, 52, 55]. A stack of Transformer blocks are added between CNN encoder and decoder for preventing blurry edges [29] and a combination of light-weight self-attention and convolutions are unitized to improve the inconsistent segmentation output [55]. The most relevant work to ours is FlowFormer [19]. It embeds the 4D cost volume into

a latent feature with a transformer-based structure and decodes the latent feature with a convolutional recurrent network. There are some major differences between this work and ours. First, they only model two consecutive frames but ignore long-range temporal correlations. Second, we enable efficient and effective pre-training in optical flow, which fully explores the potential of the transformer model to rely on the target datasets only. More comparisons of these two are in the experiments.

3. TransFlow

The task of optical flow is to estimate a series of dense displacement fields from a sequence of consecutive frames. The overview of the TransFlow model architecture is illustrated in Fig. 2. Our TransFlow is a transformer model that consists of three major components, a spatial-temporal encoder, a flow decoder and a self-supervised pre-training module. The spatial-temporal encoder jointly performs spatial globalization and temporal association to effectively capture the correlations among frames and propagate global flow features. The flow decoder decodes the feature maps for multiple frames which are then used to generate the estimated optical flow. The pre-training module is designed to learn the effective image pixel representation in a self-supervised manner, which eliminates the complex and laborious multi-stage pre-training procedures widely used in previous approaches.

3.1. Problem Definition

The input is a sequence of frames $X \in \mathbb{R}^{T \times H \times W \times C}$, which consists of T frames and C channels with (H, W) as the resolution. Following the work in ViT [13], we split each frame into N fixed-size non-overlapping patches \mathbf{x}_p , where $p \in \{1, 2, \dots, N\}$, $h \times w$ is patch size, and $N = \frac{H}{h} \times \frac{W}{w}$ ensuring the N patches span the entire

frame. The purpose of TransFlow is to output a sequence of feature map m for each frame, which is then used to generate the per-pixel displacement field f between any source and target frames.

3.2. Spatial-Temporal Encoder

Spatial Globalization The existing CNN-based flow learners demonstrate inefficiency in modeling global spatial dependencies due to the intrinsic locality of the convolution operation. However, the global spatial correlation is important information which enables effective contextual cue propagation from coherent regions to the surroundings with heavy-tailed noise, motion blurs, and large displacements, preventing performance degradation in estimating the optical flow.

In this work, we apply a spatial attention mechanism between two consecutive frames to capture the global spatial dependencies among the pixels. In particular, similar to ViT [13], each patch \mathbf{x}_p (cf. Table 2b for patch size) is first converted into a d -dimensional embedding vector $\mathbf{e}_p \in \mathbb{R}^d$ with a projection matrix \mathbf{W}_e . The final input sequence of patch embeddings is denoted as:

$$\begin{aligned} \mathbf{z}^0 &= [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_N]; \\ \mathbf{e}_p &= \mathbf{W}_e \cdot \mathbf{x}_p + \mathbf{p}_p, \end{aligned} \quad (1)$$

where \mathbf{p}_p is a set of learnable position embeddings (cf. Table 2c) to retain the positional information, which is significant to motion clues. The patch tokens are passed through a series of Multi-head Self-Attention (MSA), Multi-head Cross-Attention (MCA) and MLP layers:

$$\begin{aligned} \mathbf{y}^\ell &= \text{MLP}(\text{MSA}(\mathbf{z}^{\ell-1})) + \mathbf{z}^{\ell-1}; \\ \mathbf{z}^\ell &= \text{MLP}(\text{MCA}(\mathbf{y}^\ell)) + \mathbf{y}^\ell. \end{aligned} \quad (2)$$

Essentially, the self-attention is used to capture the global pixel dependencies within the same frame, while the cross-attention is designed to communicate the information between two adjacent frames. The self-attention and cross-attention are defined as:

$$\begin{aligned} \text{MSA}(\mathbf{z}) &= \text{softmax}(\mathbf{Q} \cdot \mathbf{K}^T / \sqrt{d}) \cdot \mathbf{V}; \\ \text{MCA}(\mathbf{z}) &= \text{softmax}(\mathbf{Q} \cdot \mathbf{K}'^T / \sqrt{d}) \cdot \mathbf{V}', \end{aligned} \quad (3)$$

where $\mathbf{Q} = \mathbf{z} \cdot \mathbf{W}^Q$, $\mathbf{K} = \mathbf{z} \cdot \mathbf{W}^K$ and $\mathbf{V} = \mathbf{z} \cdot \mathbf{W}^V$ are the query, key and value embedding matrices in MSA. $\mathbf{K}' = \mathbf{z}' \cdot \mathbf{W}'^K$ and $\mathbf{V}' = \mathbf{z}' \cdot \mathbf{W}'^V$ are the key and value embedding matrices in MCA.

The output from the attention layers is the refined correlation features, and we interleave the self-attention and cross-attention layers by L times. Through the joint aggregation, we benefit from the feature aggregation via

local frame in self-attention, which is further facilitated via adjacent perspectives in cross-attention, as depicted in Fig. 2 (a).

Temporal Association Previous approaches for estimating optical flow from each pair of adjacent frames are less effective as they ignore the inherent nature of long-range temporal associations. The motion estimation in discontinuous and occluded regions cannot be well modelled under modern architectures. To better capture high-level temporal information in the flow tokens, we learn the token embeddings by jointly modeling the temporal association with the spatial attention described above. As a result, each transformer layer can measure long-range interaction between input embeddings. Specifically, given a sequence of attentioned features from video clips consisting of T frames (see related experiments in Table 2d), we iteratively choose one as query and the rest as key features to compute the temporal attention using Softmax, which is similar as Eq. 3. Resulting in a d -dimensional embedded feature volume $\mathbf{z} \in \mathbb{R}^{T \times h \times w \times d}$, the feature volume is then passed to the following transformer decoding block. By learning temporal features in this way, temporal information can be accumulated into each frame to capture temporal associations across frames, which is shown in Fig. 2 (a).

3.3. Flow Decoder

Different from the traditional Transformer decoder, our decoder is designed to decode the feature maps of all frames (Fig. 2 (b)). These decoded features are then utilized in obtaining the final flows. Therefore, our decoder aims to generate multiple feature maps at the same time with a fixed sequence length, instead of autoregressive decoding. There are two major advantages in such a design. First, simultaneous decoding allows us to remove the encoder-decoder cross-attention in the traditional Transformer decoders. Second, beam search is no longer needed, which makes our decoding process much more efficient. Therefore, in this work, we adopt a structurally symmetric design with the Transformer encoder. In other words, our decoder has the same self-attention architecture as our encoder except that the input to the decoder is the latent cost embedding from the encoder.

Given the decoded feature maps between two consecutive frames, we compare the feature similarity by computing the correlation following [46]. To enable the end-to-end training, we apply the differentiable matching layer [54] to identify the correspondence from the adjacent frames. The final flow f can then be generated from the correspondences. During training, we further conduct an additional occlusion detection [15] by performing a forward consistency checking and considering

pixels to be occluded if the mismatching in both frames is too large. Consequently, the occlusion areas M_{occ} is computed as $M_{occ} = f_D(I_s - I_t(x + f))$, where f_D can be any function that measures the photometric distance. f is the estimated forward optical flow. I_s and I_t are the source and target images/frames. The overall objective can be formulated as:

$$L = \sum_{i=1}^R (1 - M_{occ}) \gamma^{(i-R)} \|f_{gt} - f\|_1, \quad (4)$$

where R is the total number of the training iterations. γ is a hyperparameter that controls the weight of the loss among different iterations. f_{gt} stands for the provided ground-truth flow map.

3.4. Self-supervised Pre-training

The performance of the existing flow learners heavily relies on excessive pre-training on extra synthetic datasets, followed by fine-tuning on the target domain. Without adequate pre-training procedures on large-scale data, the model typically converges with large errors. Therefore, it is an important task to design an efficient and effective pre-training strategy that improves the downstream optical flow task.

Inspired by the recent MAE [17], we introduce a masking strategy in self-supervised pre-training that adaptively masks out patch tokens and learns pixel representations by reconstructing clean signals from corrupted inputs. Specifically, we learn a score map for patch selection to choose the most informative patches as masked tokens under a determined ratio, as opposed to randomly masking in [17] or uniformly masking in [28]. In our diagnostic experiments (*cf.* Table 2e and 2f), we will demonstrate that the capability of our self-learning paradigm in recovering crucial regions can be enhanced. More specifically, we adopt multiple layers of self-attention blocks taking all the patch token embeddings as input. The attention map is then calculated as the correlation between the query embedding from the image token Q and all key embeddings across all patches K . The correlations are then followed by a Softmax activation to generate the correlation score map S_c , as depicted in Eq. 5. The correlation score map output from the final layer of the attention blocks will be utilized to guide our strategic masking learning:

$$S_c = \text{softmax}(Q \cdot K^T) \quad (5)$$

The obtained correlation score map S_c is then modeled as a ranking problem to be sorted in an ascending order to select the most informative tokens for masking, as in Fig. 2 (c). In order to prevent the discrete property

of the *argsort* operation, we instead utilize the soft sort operation in [36] denoting as $SST(\cdot)$:

$$SST(\cdot) = \text{softmax}\left(\frac{|\text{sort}(S_c)\mathbf{1}^T - \mathbf{1}(S_c)^T|}{\tau}\right) \quad (6)$$

where $|\cdot|$ calculates element-wise absolute and τ is the temperature constant that is set to 0.1 to control the degree of approximation. With the differentiable sorting, we are able to identify and retain the most significant token candidates and adaptively learn the score map as network weights in conjunction with our primary task.

4. Experiments

Datasets. Existing flow estimation approaches require a tedious training pipeline which first pre-train the models on FlyingChairs (“C”) [14] and FlyingThings (“T”) [33], and then fine-tune the trained models on Sintel (“S”) [4] and KITTI 2015 (“K”) [34]. Without the progressive steps, the flow estimation performance will get a significant degradation. Simplifying the cumbersome procedures, we rely on training optical flow task on the target domain without excessive pre-training stages. MPI-Sintel [4] dataset is rendered based on animated movies and is split into *Clean* and *Final* pass. KITTI-15 [34] contains 200 training and 200 testing road scenes with sparse ground truth flow, where images are captured via stereo cameras. For datasets provide only pairwise flow (e.g., KITTI-15), we access raw data in the self-supervised pre-training.

Implementation Details. We stack 12 transformer blocks in the encoder to adaptively learn the feature encoding. To keep the resolution to be the same as the input, we adopt the convex upsampling technique in [41] to upsample the prediction. The model is first pre-trained in a self-learning paradigm with a learning rate of $1e-4$ and then the entire network is continuously trained on the target domain with a batch size of 6 and learning rate of $12.5e-5$ for 140K steps. For the hyperparameters, γ is set to 0.8 and the masking ratio is 50%. The detailed diagnostic experiments of these hyperparameters are provided in §4.2.

Evaluation Metrics. The main evaluation metrics, used by the Sintel datasets, is the average end-point error (*AEPE*), which denotes the average pixel-wise flow error. The KITTI dataset adopts *F1-epe (%)* and *F1-all (%)*, which refers to the percentage of flow outliers over all pixels on foreground regions and entire image pixels.

4.1. Comparison to the State-of-the-Art Methods

Quantitative Evaluations. We compare our approach with existing supervised flow estimation methods on

Training Data	Method	Sintel (train)		KITTI-15 (train)		Sintel (test)		KITTI-15 (test)
		clean	final	F1-epe	F1-all	clean	final	F1-all
C+T	PWC-Net [CVPR18] [40]	2.55	3.93	10.35	33.7	-	-	-
	HD3 [CVPR19] [56]	3.84	8.77	13.17	24.0	-	-	-
	LiteFlowNet [TPAMI20] [20]	2.24	3.78	8.97	25.9	-	-	-
	RAFT [ECCV20] [41]	1.43	2.71	5.04	17.4	-	-	-
	FM-RAFT [ECCV21] [25]	1.29	2.95	6.80	19.3	-	-	-
	GMA [ICCV21] [24]	1.30	2.74	4.69	17.1	-	-	-
	Separable Flow [ICCV21] [58]	1.30	2.59	4.60	15.9	-	-	-
	FlowID [ICCV21] [53]	1.98	3.27	5.59	22.95	-	-	-
	AGFlow [AAAI22] [31]	1.31	2.69	4.82	17.0	-	-	-
	KPA-Flow [CVPR22] [30]	1.28	2.68	4.46	15.9	-	-	-
	Flowformer [ECCV22] [19]	<u>1.01</u>	<u>2.40</u>	<u>4.09</u>	<u>14.72</u>	-	-	-
Ours	0.93	2.33	3.98	14.4	-	-	-	
C+T+S+K (+H)	PWC-Net [CVPR18] [40]	-	-	-	-	4.39	5.04	9.60
	HD3 [CVPR19] [56]	1.87	1.17	1.31	4.1	4.79	4.67	6.55
	LiteFlowNet [TPAMI20] [20]	1.35	1.78	1.62	5.58	4.54	5.38	9.38
	RAFT [ECCV20] [41]	0.77	1.20	0.64	1.5	2.08	3.41	5.27
	FM-RAFT [ECCV21] [25]	0.86	1.75	0.75	2.1	1.77	3.88	6.17
	Separable Flow [ICCV21] [58]	0.71	1.14	0.68	1.57	1.99	3.27	4.89
	FlowID [ICCV21] [53]	(0.84)	(1.25)	-	(1.6)	(2.24)	(3.81)	(6.27)
	KPA-Flow [CVPR22] [30]	(0.60)	(1.02)	(0.52)	(1.10)	(1.35)	(2.36)	(4.60)
	Flowformer [ECCV22] [19]	(0.48)	(0.74)	(0.53)	(1.11)	(1.16)	(2.09)	(4.68)
	Ours	(0.42)	(0.69)	(0.49)	(1.05)	(1.06)	(2.08)	(4.32)

Table 1. **Quantitative comparisons with state-of-the-arts.** We follow existing works to compare the results on two standard benchmarks Sintel and KITTI-15. "C+T" denotes training only on FlyingChairs and FlyingThings datasets and testing on others for the generalization ability. "C+T+S+K(+H)" denotes training on mixed datasets and testing on Sintel and KITTI-15 for evaluation. Recent works [19, 30, 53] including HD1K [26] dataset for training are marked with brackets in results. Our self-learning paradigm helps to get superior results by avoiding tedious pre-training stages on "C/T" and simplifying the training pipeline. The best and second best results are highlighted in bold and underlined. See §4.1 for details.

the most popular optical flow benchmarks (*i.e.* Sintel and KITTI). Without tedious multi-stage flow estimation pre-training on synthetic benchmarks FlyingChairs and FlyingThings, our designated framework beats existing state-of-the-art methods, as demonstrated by the quantitative results in Table 1. As shown, for generalization ability, we train our TransFlow on the FlyingChairs and FlyingThings (C+T) and directly evaluate it on the Sintel and KITTI-15 without further fine-tuning. Our TransFlow depicts the best result with the smallest errors among all compared methods on both datasets. Specifically on the Sintel dataset, we achieve **0.93** and **2.33** *AEPE* on the clean and final pass, which is **0.50** and **0.38** lower than the widely used method RAFT [41]. On the KITTI-15 dataset, we reduce the *F1-all* error by **17.2%** of RAFT [41].

When following the introduced self-learning paradigm on the target datasets and evaluate on the Sintel test set, our method achieves a **1.06** and **2.08** *AEPE* on the Sintel clean and final pass, which is **49%** and **39%** lower than RAFT [41], respectively. Similarly on the KITTI-15 benchmark, our approach performs a **4.32** *F1-all* score in errors, which is **0.95** and **0.36** lower than recent RAFT [41] and FlowFormer [19]. However, RAFT [41] and FlowFormer [19] both require a multi-stage flow estimation pre-trainings before training on C+T+S+K or C+T+S+K+H datasets. Compared

to those existing approaches, our proposed pipeline delivers superior performance with significantly simpler and more effective steps.

Qualitative Evaluations. We sampled test samples from Sintel `val` set and provide the corresponding optical flow estimation of the state-of-the-art FlowFormer [19] and our TransFlow in Fig. 3. As shown, our TransFlow shows superior capability in distinguishing occluded regions and performing clearer boundaries on small and thin objects, benefiting from our consideration in the spatial globalization and temporal associations. These enhancements demonstrate the efficacy of our designed structures in spatial attentions, temporal associations, and novel strategic masking strategy in improving flow reasoning.

4.2. Diagnostic Experiments

Core Components. First, we study the efficacy of the core components of our algorithm in Table 2a to analyze their contributions to the final results. As shown, It can be seen that solely self-attention can yield limited performance as a baseline of our framework, while a combination of both self- and cross- attentions boost the performance thanks to the effective local feature aggregation between two views. The design of temporal association among multiple frames successfully

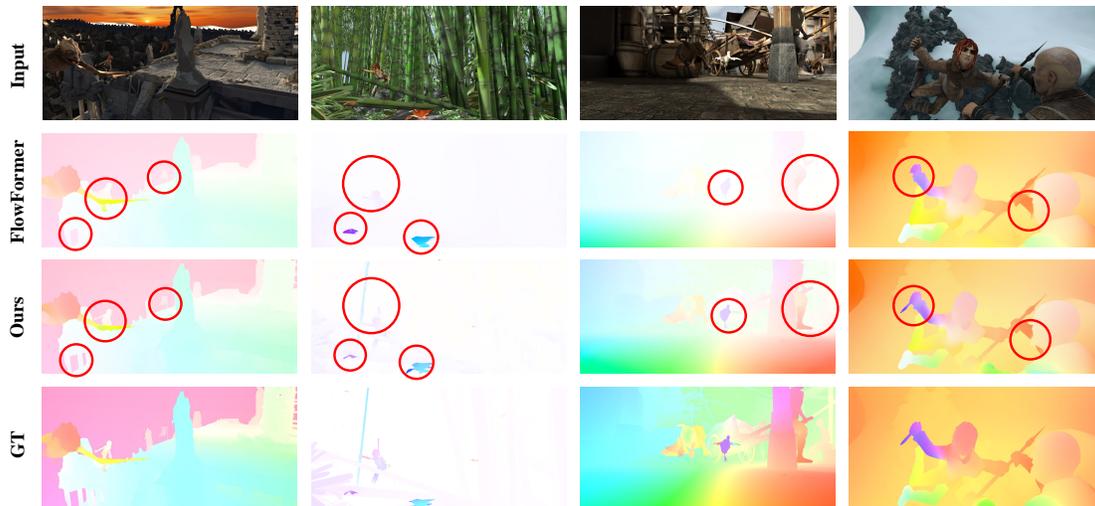


Figure 3. **Qualitative results on Sintel v1 set.** Given the target frame, we show the results of the state-of-the-art FlowFormer [19], our TransFlow results, and the provided ground truth flow. \circ highlights comparing details. See §4.1 for details.

alleviates the potential ambiguities in un-smooth and occluded regions, therefore bringing further improvement. The strategic masking reduces the reliance on multi-stage pre-training and benefits our self-learning paradigm. The occlusion consistency on the flow loss has a similar effect with the temporal association and adds additional performance gains to the results. From Table 2a, each component contributes to the improvement of performance, clearly depicting the effectiveness of our proposed components.

Patch Settings. We empirically evaluate the end-point error by adjusting the patch size (*cf.* §3.2) in our TransFlow. As shown in Table 2b, when patch size are increasingly set from 4×4 to 8×8 , the performance are slightly improved (**0.44** and **0.71** \rightarrow **0.42** and **0.69**) in *EPE* and *F1-all* on Sintel and KITTI datasets, respectively. Nevertheless, when further enlarging the patch size to 16×16 , the performance drops while the cost of computing continues to drop. The reason is that larger patch sizes lead to bigger kernel regions, resulting in the loss of global and long-range context information which is crucial for flow propagation.

Positional Embedding. The efficacy of different positional embeddings (*cf.* §3.2) is rarely discussed in previous works. Consequently, we compare the performance of flow estimation under different positional embeddings (*e.g.* Abs/Rel, Learnable/Fixed). As depicted in Table 2c, we observe that the learnable Abs Pos. achieves a slightly better result than the fixed Abs Pos. and a recent Positional Encoding Generator (PEG) [8], while showing a larger improvement than relative Pos.

We suppose that the fixed sin-cos Abs Pos. can encode the flow features almost as well as the learnable Abs Pos. and PEG Pos., and positional embedding is indispensable in our setting. In addition, we believe the degradation from relative Pos. is due to the fact that object motion requires more absolute position encoding in order to locate and learn the motion, and global information is also more important than local relative information in this task, which is consistent with our claim.

Temporal Length. Table 2d shows that as we increase the number of frames (*cf.* §3.2) fed into our temporal module, the error gets decrease since the network is able to incorporate longer temporal context and to avoid temporal artifacts and discontinuous estimation in the flow. However, Table 2d also demonstrates that the accuracy will become saturated once the number of temporal length is sufficient enough to cover visible motion. Considering that when increasing the temporal length from 5 to 7, there is no discernible difference in performance while the computational cost will increase correspondingly. Therefore, we choose 5-frame length as input.

Sampling Strategy. Table 2e shows the effect of various sampling strategies for masking (*cf.* §3.4). We compare our strategic masking with block-wise masking [1], random masking [17] and uniform masking [28]. Under the same masking ratio, it can be seen the compared samplings have different levels of degradation compared to ours. The naive block-wise masking and random masking may destroy the tokens of vital regions of the original image that are required for object motion, whereas uniform masking may disregard the sig-

(a) Contribution of each core component (cf. §3).						(b) Patch size settings (cf. §3.2).								
Self		Cross		Temporal		Masking		Occ		Sintel (val)		KITTI (val)		
✓	✓	✓	✓	✓	✓	Clean	Final	F1-all	Clean	Final	F1-all	Clean	Final	F1-all
✓						0.58	0.81	1.22						
✓	✓					0.55	0.78	1.18						
✓	✓					0.49	0.73	1.18						
✓	✓		✓			0.44	0.70	1.07						
✓	✓		✓	✓	✓	0.42	0.69	1.05						

Patch size	Sintel (val)		KITTI (val)	
	Clean	Final	F1-all	F1-all
4 × 4	0.44	0.71	1.11	
8 × 8	0.42	0.69	1.05	
14 × 14	0.46	0.77	1.14	
16 × 16	0.59	0.85	1.21	

(c) Different Pos. embeddings (cf. §3.2).				(d) Temporal Length (cf. §3.2).			(e) Sampling methods (cf. §3.4).			(f) Masking ratio (cf. §3.4).				
Pos. Embed	Sintel (val)		KITTI (val)	Temporal length	Sintel (val)		Sintel (val)	KITTI (val)	Masking ratio	Sintel (val)			KITTI (val)	
	Clean	Final	F1-all		Clean	Final				Clean	Final	Clean	Final	F1-all
Fixed Abs Pos.	0.43	0.71	1.06	2 frame	0.47	0.75	Block	0.48	0.74	1.10	30%	0.45	0.71	1.08
Learnable Abs Pos.	0.42	0.69	1.05	3 frame	0.44	0.73	Random	0.50	0.76	1.13	50%	0.42	0.69	1.05
PEG Pos.	0.42	0.70	1.07	5 frame	0.42	0.69	Uniform	0.45	0.72	1.09	70%	0.46	0.71	1.09
Learnable Rel Pos.	0.47	0.76	1.10	7 frame	0.43	0.68	Strategic	0.42	0.69	1.05	90%	0.49	0.77	1.16

Table 2. A set of diagnostic experiments. The adopted algorithm designs and settings are marked in red. See §4.2 for details.

nificance and relationship between tokens. On the contrary, our sampling has the ability to learn pixel representations effectively, which validates our claim.

Masking Ratio. Table 2f illustrates the effect of varying masking ratios (cf. §3.4). It depicts that a suitable masking ratio (50% for ours) outperforms other settings with notable advantages. Such an empirical advantage can be explained by that the higher masking ratio may discard too much necessary information for self-learning paradigm via reconstruction to learn an effective image representation, whereas a low masking ratio may not be sufficient to increase the reconstruction difficulty and, consequently, the quality of predicting a flow map.

4.3. Downstream Tasks

High-quality flow estimation plays a crucial role in many video-based downstream tasks. We show here quantitatively that TransFlow generalizes well and can help further improve the state-of-the-art of various video-based tasks, including video object detection, interpolation, and stabilization.

	Method	Backbone	<i>mAP</i> (%)
	Video Object Detection	RDN [11]	ResNet-50
RDN+ours		ResNet-50	80.4 (3.7↑)
SELSA [51]		ResNet-101	80.3
SELSA+ours		ResNet-101	82.9 (2.6↑)
PTSEFormer [43]		ResNet-101	87.4
PTSEFormer+ours		ResNet-101	89.1 (1.7↑)
Video Interpolation	Method	PSNR	SSIM
	SuperSloMo [22]	28.52	0.891
	SuperSloMo+ours	28.81 (0.29↑)	0.905 (0.014↑)
	IFR-Net [27]	29.84	0.920
	IFR-Net+ours	30.02 (0.18↑)	0.932 (0.012↑)
Video Stabilization	Method	Distortion	Stability
	StabNet [45]	0.83	0.75
	StabNet+ours	0.85 (0.02↑)	0.79 (0.04↑)
	PWStableNet [59]	0.79	0.80
	PWStableNet+ours	0.82 (0.03↑)	0.82 (0.02↑)

Table 3. Quantitative comparison of downstream video task performance with our TransFlow. See §4.3 for details.

Video Object Detection. We conduct our experiments on the ImageNet VID dataset [37] containing over 1M frames for training and more than 100k frames for validation. As shown in Table 3, adding TransFlow encoder feature in RDN [11], SELSA [51] and PTSEFormer [43] results in **3.7%**, **2.6%** and **1.7%** improvement in the mean average precision (*mAP*).

Video Frame Interpolation. To evaluate our model for 8× interpolation, we train SuperSloMo [22] and IFR-Net [27] on GoPro [35] training set with our TransFlow encoder features embedded, and test the trained model on GoPro testing set. As shown in Table 3, the updated model outperform original methods with 2 input frames in both *PSNR* and *SSIM* (*i.e.* **0.29** dB higher results than SuperSloMo and **0.18** dB higher than IFR-Net).

Video Stabilization. We follow the training configurations of StabNet [45] and PWStableNet [59] and aggregate the learned features from the TransFlow encoder and the original encoder together for the later regressor. On the DeepStab [45] dataset which contains 61 pairs of stable and unstable videos, TransFlow feature-added method achieves a higher *Distortion Value* (*D*) and *Stability Score* (*S*) than the ones without it.

5. Conclusion

We propose TransFlow, a pure transformer architecture for optical flow estimation. Extensive empirical analysis demonstrates that TransFlow establishes new records for public benchmarks. We trust that this work has the potential to provide valuable insights into the applicability of Transformer in more broad vision tasks.

Acknowledgement

This work is supported by the National Science Foundation under Award No. 2242243.

References

- [1] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR*, 2022. 7
- [2] Michael J Black and Padmanabhan Anandan. A framework for the robust estimation of optical flow. In *ICCV*, 1993. 2
- [3] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *IJCV*, 61(3):211–231, 2005. 2
- [4] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, 2012. 2, 5
- [5] Zhiwen Cao, Dongfang Liu, Qifan Wang, and Yingjie Chen. Towards unbiased label distribution learning for facial pose estimation using anisotropic spherical gaussian. In *ECCV*, 2022. 1
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 3
- [7] Zhiyuan Cheng, James Liang, Hongjun Choi, Guan hong Tao, Zhiwen Cao, Dongfang Liu, and Xiangyu Zhang. Physical attack on monocular depth estimation with optimal adversarial patches. In *ECCV*, 2022. 1
- [8] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Conditional positional encodings for vision transformers. *Arxiv preprint 2102.10882*, 2021. 7
- [9] Yiming Cui, Liqi Yan, Zhiwen Cao, and Dongfang Liu. Tf-blender: Temporal feature blender for video object detection. In *ICCV*, 2021. 3
- [10] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *ACL*, 2019. 3
- [11] Jiajun Deng, Yingwei Pan, Ting Yao, Wengang Zhou, Houqiang Li, and Tao Mei. Relation distillation networks for video object detection. In *ICCV*, 2019. 8
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 3
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 2, 3, 4
- [14] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015. 2, 5
- [15] Ziyue Feng, Liang Yang, Longlong Jing, Haiyan Wang, YingLi Tian, and Bing Li. Disentangling object motion and occlusion for unsupervised multi-frame monocular depth. In *ECCV*, 2022. 4
- [16] Mathias Gehrig, Mario Millhäusler, Daniel Gehrig, and Davide Scaramuzza. E-raft: Dense optical flow from event cameras. In *3DV*, 2021. 1
- [17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 2, 5, 7
- [18] Markus Hofinger, Samuel Rota Bulò, Lorenzo Porzi, Arno Knapitsch, Thomas Pock, and Peter Kontschieder. Improving optical flow on a pyramid level. In *ECCV*, 2020. 2, 3
- [19] Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Flowformer: A transformer architecture for optical flow. In *ECCV*, 2022. 2, 3, 6, 7
- [20] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-flownet: A lightweight convolutional neural network for optical flow estimation. In *CVPR*, 2018. 6
- [21] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, 2017. 2
- [22] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slo-mo: High quality estimation of multiple intermediate frames for video interpolation. In *CVPR*, 2018. 8
- [23] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *ICCV*, 2021. 2
- [24] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *ICCV*, 2021. 6
- [25] Shihao Jiang, Yao Lu, Hongdong Li, and Richard Hartley. Learning optical flow from a few matches. In *CVPR*, 2021. 6
- [26] Daniel Kondermann, Rahul Nair, Katrin Honauer, Karsten Krispin, Jonas Andrusis, Alexander Brock, Burkhard Gusefeld, Mohsen Rahimimoghaddam, Sabine Hofmann, Claus Brenner, et al. The hci benchmark suite: Stereo and flow ground truth with uncertainties for urban autonomous driving. In *CVPRW*, 2016. 6
- [27] Lingtong Kong, Boyuan Jiang, Donghao Luo, Wenqing Chu, Xiaoming Huang, Ying Tai, Chengjie Wang, and Jie Yang. Ifrnet: Intermediate feature refine network for efficient frame interpolation. In *CVPR*, 2022. 8
- [28] Xiang Li, Wenhai Wang, Lingfeng Yang, and Jian Yang. Uniform masking: Enabling mae pre-training for pyramid-based vision transformers with locality. *arXiv preprint arXiv:2205.10063*, 2022. 5, 7

- [29] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fuseformer: Fusing fine-grained information in transformers for video inpainting. In *ICCV*, 2021. 3
- [30] Ao Luo, Fan Yang, Xin Li, and Shuaicheng Liu. Learning optical flow with kernel patch attention. In *CVPR*, 2022. 6
- [31] Ao Luo, Fan Yang, Kunming Luo, Xin Li, Haoqiang Fan, and Shuaicheng Liu. Learning optical flow with adaptive graph reasoning. In *AAAI*, 2022. 6
- [32] Kunming Luo, Chuan Wang, Shuaicheng Liu, Haoqiang Fan, Jue Wang, and Jian Sun. Upflow: Upsampling pyramid for unsupervised optical flow learning. In *CVPR*, 2021. 1, 2
- [33] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, 2016. 2, 5
- [34] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene flow. *ISPRS Annals*, 2:427, 2015. 2, 5
- [35] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, 2017. 8
- [36] Sebastian Prillo and Julian Eisenschlos. Softsort: A continuous relaxation for the argsort operator. In *ICML*, 2020. 5
- [37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 8
- [38] Xiuchao Sui, Shaohua Li, Xue Geng, Yan Wu, Xinxing Xu, Yong Liu, Rick Goh, and Hongyuan Zhu. Craft: Cross-attentional flow transformer for robust optical flow. In *CVPR*, 2022. 3
- [39] Deqing Sun, Stefan Roth, and Michael J Black. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *IJCV*, 106(2):115–137, 2014. 2
- [40] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018. 1, 2, 6
- [41] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 2, 3, 5, 6
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017. 1, 2
- [43] Han Wang, Jun Tang, Xiaodong Liu, Shanyan Guan, Rong Xie, and Li Song. Ptseformer: Progressive temporal-spatial enhanced transformer towards video object detection. In *ECCV*, 2022. 8
- [44] Jianyuan Wang, Yiran Zhong, Yuchao Dai, Kaihao Zhang, Pan Ji, and Hongdong Li. Displacement-invariant matching cost learning for accurate optical flow estimation. In *NIPS*, 2020. 1
- [45] Miao Wang, Guo-Ye Yang, Jin-Kun Lin, Song-Hai Zhang, Ariel Shamir, Shao-Ping Lu, and Shi-Min Hu. Deep online video stabilization with multi-grid warping transformation learning. *TIP*, 28(5):2283–2292, 2018. 8
- [46] Qianqian Wang, Xiaowei Zhou, Bharath Hariharan, and Noah Snavely. Learning feature descriptors using camera pose supervision. In *ECCV*, 2020. 4
- [47] Wenguan Wang, Cheng Han, Tianfei Zhou, and Dongfang Liu. Visual recognition with deep nearest centroids. *NeurIPS*, 2022. 1
- [48] Wenguan Wang, James Liang, and Dongfang Liu. Learning equivariant segmentation with instance-unique querying. *NeurIPS*, 2022. 1
- [49] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *CVPR*, 2022. 3
- [50] Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. In *ICLR*, 2019. 3
- [51] Haiping Wu, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. Sequence level semantics aggregation for video object detection. In *ICCV*, 2019. 8
- [52] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *ICCV*, 2021. 3
- [53] Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-resolution optical flow from 1d attention and correlation. In *ICCV*, 2021. 6
- [54] Haofei Xu and Juyong Zhang. Aanet: Adaptive aggregation network for efficient stereo matching. In *CVPR*, 2020. 4
- [55] Chenglin Yang, Yilin Wang, Jianming Zhang, He Zhang, Zijun Wei, Zhe Lin, and Alan Yuille. Lite vision transformer with enhanced self-attention. In *CVPR*, 2022. 3
- [56] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *CVPR*, 2019. 6
- [57] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*, 2021. 3
- [58] Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable flow: Learning motion cost volumes for optical flow estimation. In *ICCV*, 2021. 6
- [59] Minda Zhao and Qiang Ling. Pwstabilenet: Learning pixel-wise warping maps for video stabilization. *TIP*, 29:3582–3595, 2020. 8

- [60] Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. Maskflownet: Asymmetric feature matching with learnable occlusion mask. In *CVPR*, 2020. [2](#)