

StreamFlow: Streamlined Multi-Frame Optical Flow Estimation for Video Sequences

Shangkun Sun^{1,2} Jiaming Liu³ Thomas H. Li¹ Huaxia Li³ Guoqing Liu⁴ Wei Gao¹

¹School of Electronic and Computer Engineering, Peking University

²Peng Cheng Laboratory, ³Xiaohongshu Inc., ⁴Minieye Inc.

Abstract

Occlusions between consecutive frames have long posed a significant challenge in optical flow estimation. The inherent ambiguity introduced by occlusions directly violates the brightness constancy constraint and considerably hinders pixel-to-pixel matching. To address this issue, multi-frame optical flow methods leverage adjacent frames to mitigate the local ambiguity. Nevertheless, prior multi-frame methods predominantly adopt recursive flow estimation, resulting in a considerable computational overlap. In contrast, we propose a streamlined in-batch framework that eliminates the need for extensive redundant recursive computations while concurrently developing effective spatio-temporal modeling approaches under in-batch estimation constraints. Specifically, we present a Streamlined In-batch Multi-frame (SIM) pipeline tailored to video input, attaining a similar level of time efficiency to two-frame networks. Furthermore, we introduce an efficient Integrative Spatio-temporal Coherence (ISC) modeling method for effective spatio-temporal modeling during the encoding phase, which introduces no additional parameter overhead. Additionally, we devise a Global Temporal Regressor (GTR) that effectively explores temporal relations during decoding. Benefiting from the efficient SIM pipeline and effective modules, StreamFlow not only excels in terms of performance on the challenging KITTI and Sintel datasets, with particular improvement in occluded areas but also attains a remarkable 63.82% enhancement in speed compared with previous multi-frame methods. *Code* will be available soon.

1. Introduction

Optical flow estimation, which aims to model the per-pixel correspondence between two consecutive frames, is a fundamental task in computer vision. It has various downstream applications, such as video compression [18, 20], object tracking [6, 16], and autonomous driving [4, 31]. Despite significant advancements in optical flow estimation in

recent years, occlusion remains an issue that has not been fully resolved. In particular, we consider occlusion as the disappearance of pixels in the current frame in the next frame [14], which violates the brightness consistency constraint and leads to great local ambiguity, significantly disrupting per-pixel matching.

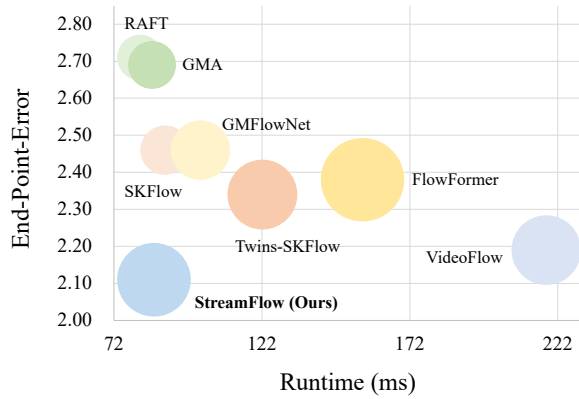


Figure 1. Comparison between performance, runtime, and parameters. A larger bubble represents more parameters. Models are trained via (C+)T schedule and validated on the Sintel final pass.

To alleviate this issue, prior research [9, 14, 38, 39, 41, 46] has proposed various approaches based on a two-frame setup. More recently, there has been a growing interest in exploring temporal cues across multiple frames [5, 15, 21, 32]. Multi-frame optical flow methods utilize information from preceding and subsequent frames to better describe the temporal continuity of pixel motion, leading to a more accurate estimation of occluded motion. Nonetheless, when dealing with video inputs, previous multi-frame flow frameworks suffer from a considerable degree of redundant computation overlap, resulting in suboptimal efficiency, as exemplified in Fig. 2. For instance, TransFlow [21] devises a pure transformer architecture based on cross-frame attention and leverages self-supervised pre-training to better optimize the spatio-temporal modules. However, the compu-

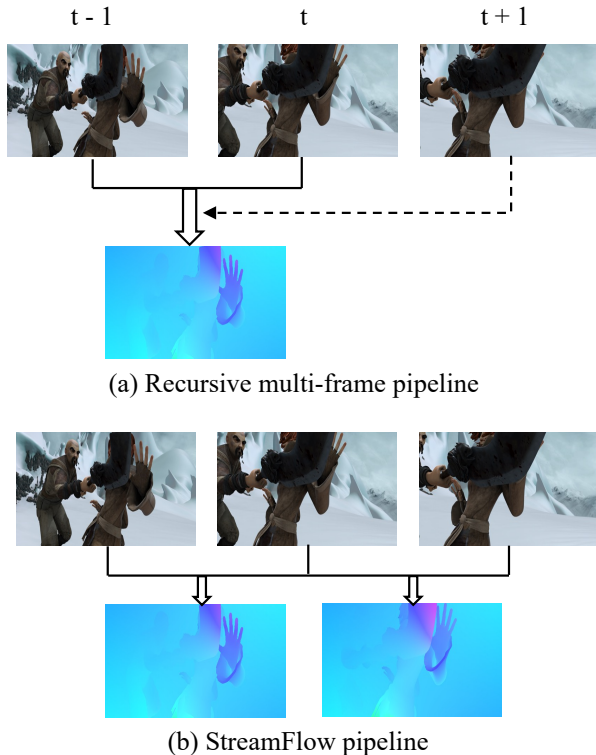


Figure 2. Comparison between different pipelines. Recursive methods leverage multi-frame for estimating two-frame flow, entailing substantial redundancy, while StreamFlow estimates multi-frame flows in-batch and eliminates overlapping computation.

tation of cross-frame attention still remains pairwise overlapping, and the pure transformer scheme is not advantageous in real-time applications. On the other hand, VideoFlow [32] additionally predicts bidirectional flows and wins a remarkable performance gain. It successfully avoids redundant pairwise computations for bidirectional flows but still necessitates recursive estimation when predicting multiple unidirectional flows.

This gives rise to a core question: *Is it possible to design a multi-frame pipeline that mitigates overlapping computations for video sequences while still effectively exploiting temporal cues and maintaining high efficiency in training and inference?*

In this work, we propose StreamFlow, a streamlined multi-frame optical flow estimation method tailored for video inputs. StreamFlow is made efficient through the Streamlined In-batch Multi-frame (SIM) pipeline, which avoids repetitive, overlapping computations when predicting unidirectional flows for video sequences. Furthermore, StreamFlow also explores the challenge of effectively modeling spatio-temporal cues under the constraint of non-overlapping in-batch estimation. StreamFlow proposes a

parameter-efficient Integrative Spatio-temporal Coherence (ISC) modeling module during encoding, and a Global Temporal Regressor (GTR) to decode all flows. Notably, these modules are quite lightweight, and StreamFlow attains comparable efficiency compared to two-frame methods with remarkable accuracy, as illustrated in Fig. 1. Without self-supervised pre-training and the aim of bidirectional flows, StreamFlow achieves superior performance on Sintel and KITTI datasets, especially on the occluded regions.

In summary, our contributions are as follows:

- We propose a Streamlined In-batch Multi-frame (SIM) pipeline for optical flow estimation, which eliminates the repetitive overlapping computation when computing unidirectional flows for video inputs.
- Under the constraint of a non-overlapping pipeline, we specifically designed the Integrative Spatio-temporal Coherence (ISC) module, which introduces no additional parameters and effectively exploits spatio-temporal cues.
- For the SIM pipeline, we devise a Global Temporal Regressor (GTR) during decoding to further exploit temporal cues with modest additional computation cost.
- The proposed StreamFlow achieves superior performances on multiple benchmarks, particularly in occluded regions with comparable efficiency compared with two-frame methods, resulting in substantial improvements in optical flow estimation.

2. Related work

Two-frame optical flow. Optical flow estimation in the form of a supervised learning task has been performed by FlowNet [8] using Convolutional Neural Networks (CNN). The encoder-decoder architecture of FlowNet predicts flow from coarse-to-fine using the hierarchy of the flow pyramid. Thereafter, a number of refined coarse-to-fine approaches [10–13, 36, 37, 42, 45] emerged. The flow pyramid is constructed for the coarse-to-fine approach, which predicts the flow based on the flow guidance at a higher pyramid level. However, the flow guidance is often too coarse to capture small motions delicately and creates errors in later estimation. RAFT [39] recently introduced an iterative all-pairs flow transform technique, which enables the prediction of high-resolution flow and recurrent refinement of the residual flow estimation. RAFT positively addresses the challenges of small motions and has consequently received high interest and performance in the field, inspiring numerous follow-up works [14, 22, 23, 38, 41, 46].

Occlusions handling. Occlusion poses a great challenge to optical flow networks. It directly violates the brightness consistency constraint, which supposes pixels between adjacent frames remain the same brightness during the motion. The ambiguity brought by occlusions seriously interferes with the per-pixel matching as two-frame networks heavily

rely on local evidence. Previous two-frame works mainly resolve the occluded pixels via multi-scale searching [36] or non-local modeling [9, 14, 33, 38, 41, 46]. These methods resolve the absent information to a certain extent. Nevertheless, in situations with severe occlusions, it becomes difficult to make up for the lack of local evidence without temporal cues, and the performance of two-frame networks remains limited in such scenarios.

Multi-frame optical flow. Exploiting temporal cues in optical flow estimation is an effective way to recover the occluded motion. Previous works [1, 5, 15, 21, 26, 30, 32, 40] propose various approaches to fuse temporal cues, such as leveraging previously predicted motion feature, optical flow, or contextual information. For instance, ContinualFlow [26] uses previous flow priors to estimate the current occlusion map. STaRFlow [1] pass extracted features from different in multiple scales, jointly with occlusion maps. [39] proposes a warm-start strategy to initialize the original flow with the past flow before prediction. MFCFlow [5] and MFR [15] propose to leverage previously estimated motion features during decoding via feed-forward CNNs and self-similarity modeling, respectively. Nevertheless, these methods obtain a recursive strategy when handling video sequences, which divide the input sequence into lots of overlapping groups and take huge repeated computations. TransFlow [21] decodes all flows simultaneously and achieves impressive results. However, it needs self-supervised pre-training on the flow datasets to help the temporal modeling modules converge. Besides, its pure transformer architecture and the overlapping computation when calculating cross-frame attention do not have advantages in terms of time. VideoFlow [32] additionally predicts the bi-direction flow to help the uni-direction flow estimation and win remarkable performance gain. Nevertheless, it still follows the recursive method to predict multiple unidirectional flows with the cost of predicting bidirectional flows. In contrast, StreamFlow is proposed to avoid redundant, overlapping computation for consecutive unidirectional flow predictions while exploring efficient and effective temporal modules design under such a pipeline.

3. Methodology

In this Section, we introduce StreamFlow, an efficient and effective in-batch framework for multi-frame optical flow estimation. The key components of StreamFlow consist of three parts: (1) The Streamlined In-batch Multi-frame (SIM) pipeline for efficient multi-frame estimation. (2) Integrative Spatio-temporal Coherence (ISC) modeling that is specifically designed for spatio-temporal modeling in the encoder of the SIM pipeline. (3) Global Temporal Regressor (GTR) that learns temporal relations for the SIM pipeline during decoding. We will first give an overview of

our methods in Sec. 3.1, and then introduce each module in Sec. 3.2, Sec. 3.3, and 3.4, respectively. In the end, we discuss the loss function design in Sec. 3.5.

3.1. Overview

The overall framework of StreamFlow is illustrated in Figure 3. For the basic encoder and decoder, similar to VideoFlow [32], StreamFlow adopts the Twins transformer [7] as the encoder and utilizes the motion encoder and updater in SKFlow [38] during decoding. The overall iterative-refinement design that adopts an iterative decoder is the paradigm proposed in RAFT [39] and followed by a lot of subsequent works [9, 14, 33, 35, 38, 38]. Input frames are first passed to two feature encoders that share the same architecture to extract the correlation feature and contextual feature, respectively. Then, the multi-scale all-pairs correlation vector is calculated based on the correlation feature. Namely, given feature embeddings \mathbf{e}_1 and \mathbf{e}_2 from the target frame and the reference frame, respectively:

$$\mathbf{c}^l(i, j, m, n) = \frac{1}{2^{2l}} \sum_u^{2^l} \sum_v^{2^l} \langle \mathbf{e}_1(i, j), \mathbf{e}_2(2^l m + u, 2^l n + v) \rangle, \quad (1)$$

where the derived $\mathbf{c}^l(i, j, m, n)$ is the average over the correlation in the local $2^l \times 2^l$ window. l denotes the l th correlation level. u and v are the horizontal and vertical pixel motions, respectively. $\langle \cdot, \cdot \rangle$ refers to the dot product function. In summary, $\mathbf{c}^l(i, j, m, n)$ means the cost volume vector of \mathbf{e}_1 and \mathbf{e}_2 pooled with the $2^l \times 2^l$ kernel.

Then, the iterative decoder refines the flows via several updates. As depicted in Fig. 3, flows are initialized to zeros. The derived multi-scale correlation vector, extracted context feature, and the initialized flows are passed to the decoder, and then the refinement is conducted.

3.2. Streamlined in-batch multi-frame pipeline

As shown in Fig. 2, previous multi-frame networks mainly compute recursively for the video inputs, resulting in a great deal of overlapping computation. Specifically, frames are divided into groups, and the flow between each frame in sequence is predicted recursively before processing the next group. The issue here lies in the overlap between frames within a group, where the same optical flow between overlapped frames would be calculated repeatedly. In contrast, StreamFlow is equipped with a Streamlined In-batch Multi-frame (SIM) Pipeline that tries to avoid redundancy. In the SIM pipeline, frames are divided into non-overlapping groups except for the initial frame. And in the same group, the repetitive computation is greatly reduced. First, each frame and its embeddings are stored in the memory bank so that the feature extraction and correlation construction are conducted only once. Besides, the spatio-temporal modeling methods are also designed specifically for non-

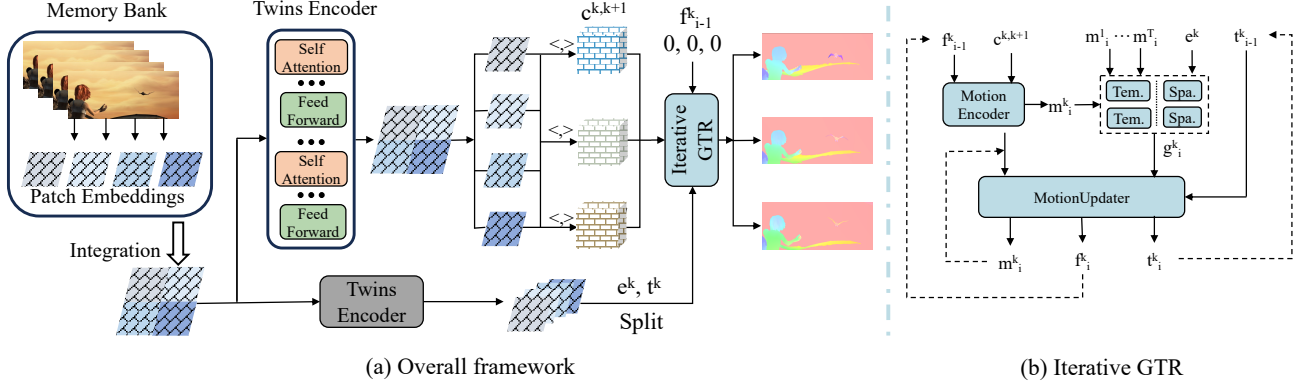


Figure 3. Overview of StreamFlow. (a) illustrates the overall framework and \langle, \rangle denotes the dot-product operation. (b) depicts the detailed module design of the GTR decoder.

overlapping computation, which will be given detailed discussion in Sec. 3.3 and Sec. 3.4. The pipeline is comparable to two-frame methods in latency with more accuracy and modest additional computation, as illustrated in Fig. 1.

3.3. Integrative spatio-temporal coherence

During the encoding process, we propose an Integrative Spatio-temporal Coherence (ISC) modeling method, especially for the SIM pipeline. Our design principles for temporal modeling modules in the decoder encompass two facets: firstly, adherence to the design criteria of the SIM pipeline, with a focus on minimizing pair-wise overlap operations, such as the computation of cross-frame attention between every pair of consecutive frames. Secondly, the modules should be efficient enough and not impede the overall speed of the network.

Therefore, we design the ISC method, which introduces no additional parameters and overlapping computation while learning spatio-temporal relations efficiently and effectively. The ISC method inherently takes the original modules in Twins. Specifically, after deriving patch embeddings from consecutive, ISC integrates temporally contiguous multiple input embeddings into a large feature embedding along the spatial dimension. Subsequently, it models the derived spatio-temporal graph using self-attention mechanisms and feed-forward layers in Twins, which could be formulated as,

$$\mathbf{x}_c^i = \text{Integration}_{t=1}^T(\mathbf{x}_{t,c}^j), \quad (2)$$

$$f(a_i, b_j) = \frac{\exp(a_i^T b_j / \sqrt{d})}{\sum_{j=1}^N \exp(a_i^T b_j / \sqrt{d})} \quad (3)$$

$$\mathbf{y}_c^i = f(\mathbf{q}(\mathbf{x}_c^i), \mathbf{k}(\mathbf{x}_c^i))\mathbf{v}(\mathbf{x}_c^i), \quad (4)$$

$$\mathbf{x}_c^i = \mathbf{x}_c^i + \mathbf{W}_{\text{proj}}\mathbf{y}_c^i, \quad (5)$$

$$(6)$$

where $f(\cdot)$ is the attention function which conducts dot-product and softmax operation, $\mathbf{x}_{(t,c)}^j$ is the j th vector along spatial dimension at channel c of the t th frame. \mathbf{q} , \mathbf{k} and \mathbf{v} is the derived query, key, and value vector \mathbf{W}_{proj} is the projection matrix. By leveraging the derived spatio-temporal graph, the spatial and temporal relations are learned effectively, and no additional parameters are involved.

3.4. Global temporal regressor

As for the decoder, we propose a Global Temporal Regressor (GTR) to predict and refine the predicted flows. Compared with the previous widely used decoder [14, 22, 23, 38, 39, 44], GTR introduces the temporal modeling module to exploit temporal cues from consecutive frames. Different from VideoFlow [32] that concatenates motion features along a temporal dimension and implicitly learns temporal relations or TransFlow [21] that applies a transformer symmetric to the encoder, the core of GTR is super convolution kernels [38] and a lightweight temporal transformer block. The input correlation vectors, initialized flows, and contextual features are first passed into a motion encoder to derive motion features and then extracted for temporal and spatial features, which could be formulated as:

$$\mathbf{m}_i^k = \text{MotionEncoder}(\mathbf{f}_{i-1}^k, \mathbf{c}^{k,k+1}), \quad (7)$$

$$\mathbf{r}_i = \text{TemLayer}_{j=1}^T(\mathbf{m}_i^j), \quad (8)$$

$$\mathbf{s}_i^k = \text{SpaCrossAttn}(\mathbf{m}_i^j, \mathbf{e}^j), \quad (9)$$

$$\mathbf{g}_i^k = \text{Concat}(\mathbf{r}_i, \mathbf{s}_i^k), \quad (10)$$

$$\mathbf{t}_i^k, \mathbf{m}_i^k, \Delta \mathbf{f}_i^k = \text{MotionUpdater}(\mathbf{m}_i^k, \mathbf{g}_i^k, \mathbf{t}_{i-1}^k), \quad (11)$$

$$\mathbf{f}_i^k = \mathbf{f}_{i-1}^k + \Delta \mathbf{f}_i^k \quad (12)$$

where \mathbf{m}_i^k is the derived motion feature of frame k at the i th update and \mathbf{f}_{i-1}^k denote the flow of frame k after $i-1$ th refinement. $\mathbf{c}_{k,k+1}$ denotes the correlation vector between frame k and $k+1$. *MotionEncoder* is the same motion

encoder in the decoder of SKFlow [38]. \mathbf{r}_i denotes the temporal feature embedding extracted from the motion features of all frames. Notably, the caching mechanism of the MemoryBank is employed, thus necessitating the calculation of \mathbf{r}_i only once for different frames. *TempLayer* is a lightweight temporal-learning layer that consists of temporal attention and feed-forward layers. \mathbf{e}^k refers to the feature embedding of frame k . Note that e and c are not updated during the refinement. Inspired by the success of cross-attention mechanism in GMA [14], *SpaCrossAttn* utilizes \mathbf{m}_i^k and \mathbf{e}^k to perform cross-attention. \mathbf{t} denotes the extracted contextual information, which would be updated during each refinement. In practice, the decoder estimates the residual of flow $\Delta\mathbf{f}_i^k$. And the final flow \mathbf{f}^k is updated via $\Delta\mathbf{f}_i^k$ during each refinement.

3.5. Supervision

StreamFlow adopts the overall loss in the same group as the total loss function. For each flow, StreamFlow adopts the same loss function as successful two-frame networks. Namely, the weighted sum for the predicted flows at different refinements. During both the training and the fine-tuning process, the supervision could be formulated as follows:

$$\mathcal{L} = \sum_{k=1}^T \sum_{i=1}^N \theta^{N-i} \|\mathbf{f}_i^k - \mathbf{f}_{gt}^k\|_1, \quad (13)$$

where \mathbf{f}_i^k refers to the flow of frame k at the i th refinement. T and N are the number of frames and refinements, respectively. θ denotes the weights on corresponding estimated flows. \mathbf{f}_{gt} is the ground truth flow and $\|\cdot\|_1$ means the l_1 distance between ground truth and our predicted flow. In practice, N is set to 12, θ is set to 0.8, the same as previous works [14, 32, 38, 39] for a fair comparison.

4. Experiments

Experimental setup. In this study, we evaluate our StreamFlow model on the Sintel [3] and KITTI [25] datasets, following previous works [9, 38, 39]. In previous works, models are initially pre-trained on the FlyingChairs [8] and FlyingThings [24] datasets using the ‘‘C+T’’ schedule and then are subsequently fine-tuned using the ‘‘C+T+S+K+H’’ schedule on Sintel and KITTI datasets. In specific, for Sintel, models are fine-tuned on a combination of FlyingThings, Sintel, KITTI, and HD1K [17]. After fine-tuning on Sintel, models are further fine-tuned using the KITTI dataset for the evaluation of KITTI.

Implementation details. Our StreamFlow method is built with PyTorch [27] library, and our experiments are conducted on the NVIDIA A100 GPUs. During training, we

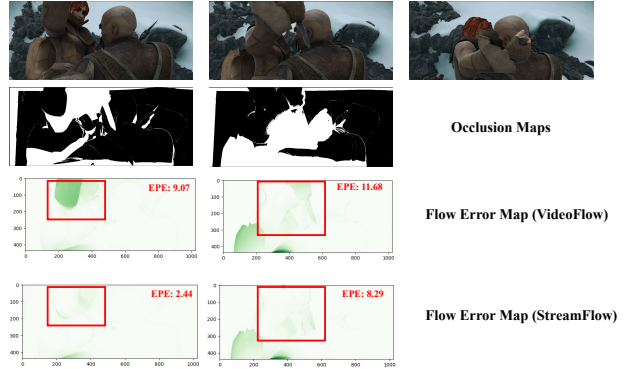


Figure 4. Visualizations of the performance on the occluded regions. StreamFlow achieves comparable performance even with advanced methods. All models are trained on the FlyingThings dataset. A darker color in the flow error map denotes a higher estimation error compared with ground truth.

adopt the AdamW [19] optimizer and the one-cycle learning rate policy [34], following previous works [14, 38, 39]. During training, the number of refinements in the decoder is set to 12, following previous works. Given the absence of multi-frame data information in the Chairs dataset, we follow VideoFlow [32] to directly train on the FlyingThings dataset in the first stage. The remaining training configurations remain consistent with prior works [14, 32, 38, 39]. The temporal and non-temporal modeling modules are concurrently trained.

4.1. Quantitative Results

From Table 1, we can learn that StreamFlow achieves superior performance on Sintel and KITTI. After being pre-trained on the FlyingThings dataset, StreamFlow demonstrates strong generalization ability across datasets. Given the leading performance of previous methods, StreamFlow could further reduce the end-point error by 0.16 and 0.08 on the challenging Sintel clean and final pass, respectively. On KITTI, StreamFlow outperforms the previous state-of-the-art method with 0.11 and 17.65% lower EPE and Fl-all metric. Notably, without self-supervised pre-training or bi-directional flows, StreamFlow attains remarkable accuracy and efficiency on the challenging Sintel and KITTI benchmarks after the (C)+T and the +S+K+H schedule.

4.2. Occlusion Analysis

In this section, we validate if StreamFlow could help improve the performance on the occlusions. We compare StreamFlow with its base two-frame model Twins-SKFlow, which strengthens SKFlow [38] with the Twins [7] encoder. Evaluations are conducted on the matched and unmatched areas of the challenging Sintel test dataset. The matched

Training Data	Method	Sintel (train)		KITTI-15 (train)		Sintel (test)		KITTI-15 (test)
		Clean	Final	Fl-epe	Fl-all	Clean	Final	Fl-all
(C+)T	HD3 [43]	3.84	8.77	13.17	24.0	-	-	-
	VCN [42]	2.21	3.68	8.36	25.1	-	-	-
	FlowNet2 [13]	2.02	3.54	10.08	30.0	3.96	6.02	-
	RAFT [39]	1.43	2.71	5.04	17.4	-	-	-
	CRAFT [35]	1.27	2.79	4.88	17.5	-	-	-
	GMA [14]	1.30	2.74	4.69	17.1	-	-	-
	SKFlow [38]	1.22	2.46	4.27	15.5	-	-	-
	FlowFormer [9]	1.00	2.45	4.09	14.7	-	-	-
	GAFlow [23]	1.02	2.45	3.98	15.0	-	-	-
	TransFlow [21]	<u>0.93</u>	2.33	3.98	<u>14.4</u>	-	-	-
	VideoFlow-BOF [32]	1.03	<u>2.19</u>	3.96	15.3	-	-	-
Ours	0.87	2.11	3.85	12.6	-	-	-	
(C+)T+S+K+H	LiteFlowNet2 [11]	(1.30)	(1.62)	(1.47)	(4.8)	3.48	4.69	7.74
	IRR-PWC [12]	(1.92)	(2.51)	(1.63)	(5.3)	3.84	4.58	7.65
	MaskFlowNet [45]	-	-	-	-	2.52	4.17	6.10
	Separable Flow[44]	(0.69)	(1.10)	(0.69)	(1.6)	1.50	2.67	4.64
	PWC-Fusion [37]	-	-	-	-	3.43	4.57	7.17
	StarFlow [1]	-	-	-	-	2.72	3.71	7.65
	RAFT* [39]	(0.76)	(1.22)	(0.63)	(1.5)	1.61	2.86	5.10
	GMA* [14]	(0.62)	(1.06)	(0.57)	(1.2)	1.39	2.47	5.15
	GMFlow [41]	-	-	-	-	1.74	2.90	9.32
	GMFlowNet [46]	(0.59)	(0.91)	(0.64)	(1.5)	1.39	2.65	4.79
	AGFlow* [22]	(0.65)	(1.07)	(0.58)	(1.2)	1.43	2.47	4.89
	SKFlow* [38]	(0.52)	(0.78)	(0.51)	(0.9)	1.28	2.27	4.84
	FlowFormer [9]	(0.48)	(0.74)	(0.53)	(1.1)	1.16	2.09	4.68
	MFRFlow [15]	(0.64)	(1.04)	(0.54)	(1.1)	1.55	2.80	5.03
	MFCFlow [5]	(0.56)	(0.89)	(0.55)	(1.1)	1.49	2.58	5.00
	TransFlow [21]	(0.42)	(0.69)	<u>(0.49)</u>	(1.05)	1.06	2.08	<u>4.32</u>
VideoFlow-BOF [32]	<u>(0.37)</u>	<u>(0.54)</u>	<u>(0.52)</u>	<u>(0.85)</u>	1.00	1.71	4.44	
Ours	(0.28)	0.38	0.47	0.77	<u>1.04</u>	<u>1.87</u>	4.24	

Table 1. Quantitative results on Sintel and KITTI. The average End-Point Error (EPE) is reported as the evaluation metric if not specified. * refers to the warm-start strategy [39] that use the previous flow for initialization. Bold and underlined metrics denote the method that ranks 1st and 2nd, respectively. Our method achieves superior performance on different benchmarks.

areas denote regions visible in adjacent frames and the unmatched areas refer to regions visible only in one of two adjacent frames. Our models are trained using the T+S+H+K schedule. We could learn that StreamFlow attains remarkable improvements on occluded areas, as shown in Tab. 3. We also visualize the performance on occluded regions in Fig. 4. On the challenging Sintel final test set, StreamFlow attains the improvement of 10.77% and 11.83% on unmatched and matched regions, respectively. On the clean pass, StreamFlow improves the performance by 15.53%, 15.56%, and 15.45% on unmatched, matched, and overall regions. We could learn that StreamFlow improves not only the flow estimation in unmatched regions but also the estimation in matched regions.

4.3. Ablations

In this section, we verify the effectiveness of StreamFlow designs, as shown in Tab. 2. For a fair comparison, all models in the same experiment are trained under the same settings on the FlyingThings dataset. Then we evaluate each method on Sintel and KITTI. Below we will introduce each experiment in more detail.

SIM pipeline. We test the efficiency of the vanilla recursive pipeline and our SIM pipeline. Recursive methods utilize multi-frames to predict the flow of the current two frames and bring substantial redundant computation, while the SIM pipeline estimates multiple flows concurrently and

Experiment	Method	Sintel				KITTI		Param (M)	Latency (ms)
		Clean	Final	Occ (Albedo)	Noc (Albedo)	Fl-epe	Fl-all		
SIM pipeline	w/o	1.03	2.34	7.69	0.35	4.64	14.70	12.49	122.18
	<u>w/</u>	1.03	2.34	7.69	0.35	4.64	14.70	12.49	84.59
Temporal modules	w/o	1.03	2.34	7.69	0.35	4.64	14.70	12.49	84.59
	Temporal attn	0.96	2.31	7.38	0.35	4.38	14.96	14.14	91.17
	Pseudo 3D conv	1.05	2.36	7.60	0.38	4.46	15.20	13.48	87.41
	3D conv	0.98	2.34	7.63	0.33	4.57	15.59	16.03	93.05
	<u>ISC</u>	0.97	2.29	7.11	0.32	4.14	14.16	12.49	88.35
Additional params	w/o	0.97	2.29	7.11	0.32	4.14	14.16	12.49	84.59
	w/	0.98	2.24	7.33	0.31	4.15	13.94	13.77	89.29
	<u>Ours</u>	0.93	2.15	7.06	0.31	3.92	12.36	13.77	89.76
GTR module	w/o	0.97	2.29	7.11	0.32	4.14	14.16	12.49	88.35
	<u>w/</u>	0.93	2.15	7.06	0.31	3.92	12.36	13.77	89.76
ISC module	w/o	1.01	2.19	7.23	0.33	4.06	13.95	13.77	86.02
	<u>w/</u>	0.93	2.15	7.06	0.31	3.92	12.36	13.77	89.76
Number of frames	3	0.93	2.15	7.06	0.31	3.92	12.36	13.77	89.76
	<u>4</u>	0.87	2.11	6.24	0.31	3.85	12.62	14.25	85.53

Table 2. Ablations on our proposed design. All models are trained using the "C+T" schedule and validated on Sintel. The number of refinements is 12 for all methods. The settings used in our final model are underlined.

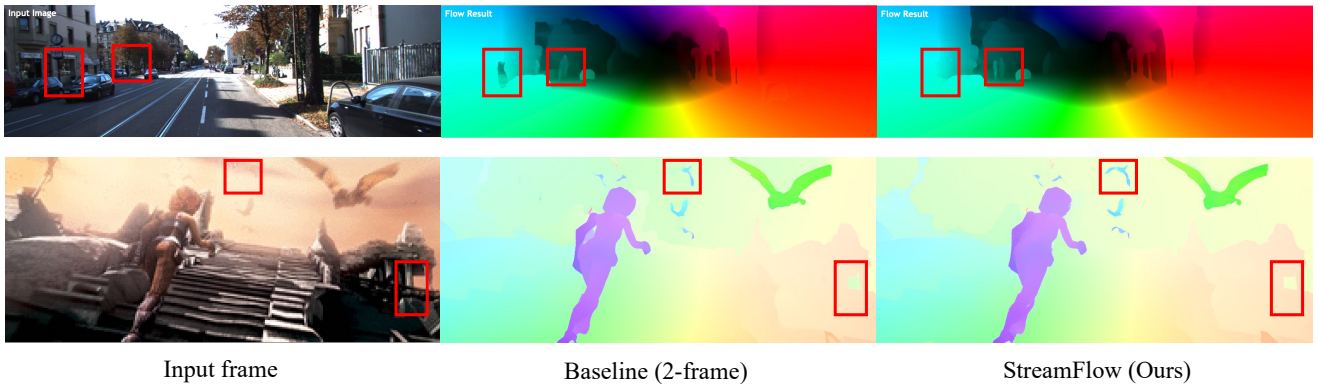


Figure 5. Visualizations of results on Sintel and KITTI test sets. Differences are highlighted with red bounding boxes. StreamFlow achieves fewer artifacts on both synthetic and real-world scenes.

minimizes the overlapping calculation. As shown in Tab. 2, the SIM pipeline brings great gain in efficiency.

Temporal modules. In this part, we explore the performance and efficiency of different temporal modeling methods in the flow encoder. Temporal attn refers to applying a temporal attention layer after each spatial self-attention modeling in Twins. Pseudo conv [29] denotes stacking 1D convolution layers in the temporal dimension to imitate 3D convolutions at minimal cost. We also apply 3D convolutions at the end of the flow encoder to learn temporal relations. As shown in Tab. 2, our ISC module achieves a

good trade-off between efficiency and effectiveness. The improvements achieved by other methods are not as pronounced. We hypothesize that the limited volume of optical flow data impedes the efficient training of the spatio-temporal module from scratch to accomplish good optimization. For comparison, VideoFlow does not apply temporal modeling modules in the encoder, and TransFlow [21] applies self-supervised pre-training for better optimization.

Additional params. In this part, we aim to determine whether the performance gain is due to the additional parameters or the effective temporal modeling method. To this

Method	Clean			Final		
	Unm.	Mat.	All	Unm.	Mat.	All
GMFlow [41]	10.56	0.65	1.74	15.80	1.32	2.90
GMFlowNet [46]	8.49	0.52	1.39	13.88	1.27	2.65
SKFlow [38]	7.24	0.55	1.28	11.51	1.46	2.28
FlowFormer [33]	7.16	0.42	1.16	11.30	0.96	2.09
TransFlow [21]	6.77	0.36	1.06	10.96	0.99	2.08
Baseline	7.60	0.45	1.23	11.70	0.93	2.11
Ours	6.42	0.38	1.04	10.44	0.82	1.87

Table 3. Occlusion analysis on Sintel test set. Unm. and Mat. denote performance on unmatched and matched areas, respectively.

end, we introduce the additional parameters by widening the baseline network. Namely, we extract higher-dimension features along the spatial dimension and concatenate them with the original motion feature. All models in this section are equipped with the ISC module. “w/o” denotes the baseline Twins-SKFlow network. “w” means adding additional parameters. “Ours” denotes the method equipped with our temporal modeling modules. Results show the improvement achieved by simply adding more parameters is minor, and the performance gain is primarily attributed to the effectiveness of StreamFlow modules.

GTR module. We also examined whether the GTR module could enhance flow predictions. “w/o” means applying vanilla SKFlow decoder while “w” denotes using GTR. All models in this part utilize the ISC module in the encoder. Tab. 2 demonstrates the necessity of incorporating the GTR. With GTR, StreamFlow could further achieve stable improvement on multiple benchmarks. We could also learn that GTR especially helps the flow estimation on the challenging final passes, with the performance gain of 0.14.

ISC module. In this part, we verify the effectiveness of the proposed ISC module. All models in this part adopt GTR as the flow decoder. From Tab. 2, we could learn that the ISC module is efficient and effective in temporal modeling and makes a significant contribution to the improvement of the multiple-frame pipeline. It introduces no additional parameters and a modest increase in runtime, while significantly boosting the performance.

Number of frames. We delve into the influence of different numbers of input frames, as illustrated in Tab. 2. We set the number of frames to 4 due to limitations in GPU memory. From an efficiency standpoint, augmenting the number of input frames results in a higher proportion of redundant computations eliminated by StreamFlow within the to-

tal computational workload, consequently leading to a more substantial improvement in processing time. Although there is an increase in the parameter count for temporal modeling, the efficiency of StreamFlow is further enhanced in the context of four input frames due to a reduced proportion of redundant computations, resulting in a shorter average prediction time per frame compared to the three-frame setting.

4.4. Qualitative results

In this section, we demonstrate visualization results on both synthetic and real-world scenes. We test the models on the challenging Sintel [3] and KITTI [25], as shown in Fig. 5. In the appendix, we also demonstrate the qualitative performance on the real-world dataset DAVIS [28]. Our models are pre-trained using the T+H+S+K schedule. We could learn that StreamFlow could still achieve remarkable qualitative results when generalized to real-world scenes.

4.5. Efficiency analysis

In this section, we evaluate the efficiency of the StreamFlow method in terms of runtime and parameter counts. Our experiments were conducted on an NVIDIA A100 GPU. Models are trained using the (C+)T schedule and evaluated on the Sintel dataset. The runtime is measured as the average inference time per frame of five runs on the Sintel training set. Figure Fig. 1 depicts the results, where the size of the bubble corresponds to the number of parameters, the horizontal axis represents time, and the vertical axis represents end-point-error. We could learn StreamFlow achieves nearly comparable efficiency with state-of-the-art two-frame methods while achieving superior performance. The key to maintaining high efficiency is its non-overlapping SIM pipeline. StreamFlow does not perform pairwise redundant computation and predicts all flows simultaneously. Another reason for the high speed is the CNN-based decoder of StreamFlow. We could learn that StreamFlow is much faster than the pure two-frame transformer architecture FlowFormer. Besides, the specially designed lightweight temporal-modeling modules also contribute to the efficiency, simultaneously aiding in better results compared to the 2-frame baseline Twins-SKFlow.

5. Conclusion

In this work, we proposed StreamFlow, a multi-frame optical flow estimation approach proficient in identifying optical flow across multiple video frames using efficient Spatio-temporal relationship mining. StreamFlow proposes to estimate multi-frame optical flow via an in-batch method (SIM pipeline) and explores the design of temporal modeling modules under such constraints. In specific, StreamFlow introduces a parameter-efficient Integrative Spatio-temporal Coherence (ISC) module that is seamlessly equipped with the encoder, and designs an efficient and effective Global

Temporal Regressor (GTR) module in the decoder. Extensive experiments demonstrate the efficiency and effectiveness of StreamFlow. With the proposed SIM pipeline, ISC, and GTR module, StreamFlow showed comparable efficiency with two-frame methods while achieving remarkable accuracy, especially in occluded regions.

StreamFlow: Streamlined Multi-Frame Optical Flow Estimation for Video Sequences

Appendix

Qualitative analysis on real-world scenes In this section, we facilitate our visualizations and evaluations using two prominent real-world datasets, namely DAVIS [28]. The DAVIS dataset, short for Densely Annotated Video Segmentation, is a widely recognized benchmark in the field of computer vision. It comprises high-quality video sequences captured in diverse scenarios, encompassing a broad range of challenging visual conditions such as occlusions, motion blur, and dynamic object interactions. The dataset provides pixel-level annotations for every frame, facilitating precise evaluation and comparison of various video segmentation methods. The visualizations on the DAVIS dataset is shown in Fig. 6. Our model is pretrained using the “T” and “T+S+H+K” schedule and then fine-tuned on KITTI [25]. “T” denotes the FlyingThings [24] dataset and “T+S+H+K” refers to the combination of the FlyingThings, Sintel [3], HD1K [17], and KITTI datasets. Then we infer our models on the DAVIS dataset. The number of refinements is set to 12. The number of input frames for each non-overlapping group is 3. We could learn that StreamFlow demonstrates remarkable adaptability across real-world datasets, showing its robust performance in challenging scenes for optical flow estimation. This is particularly evident in scenarios such as the occlusion of the bear’s hind legs in the first row, first column and the small motion of the small tennis ball in the last column. Additionally, it can be observed that in the motion captured in the first row, second and third columns, the hind legs of the camel and the leg movements of the dancer are also vividly delineated. These instances reaffirm its efficacy in diverse and demanding environments for optical flow estimation.

6. Initialization of GTR

In this section, we investigate the impact of different GTR initialization methods. Previous works in spatio-temporal modeling such as [2] have suggested initializing the temporal modules with zero values. We employed two distinct initialization approaches, namely zero initialization and PyTorch’s default initialization, and the corresponding results are presented in Tab. 4. Following training on the FlyingThings dataset, the model was tested on the Sintel and KITTI datasets. It is evident from the results that the zero initialization could contribute to a better overall performance.

References

[1] Elaine Angelino, Daniel Yamins, and Margo Seltzer. Starflow: A script-centric data analysis environment. In

Provenance and Annotation of Data and Processes: Third International Provenance and Annotation Workshop, IPAW 2010, Troy, NY, USA, June 15-16, 2010. Revised Selected Papers 3, pages 236–250. Springer, 2010. 3, 6

- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, page 4, 2021. 1
- [3] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *European conference on computer vision*, pages 611–625. Springer, 2012. 5, 8, 1
- [4] Linda Capito, Umit Ozguner, and Keith Redmill. Optical flow based visual potential field for autonomous driving. In *2020 IEEE Intelligent Vehicles Symposium (IV)*, pages 885–891. IEEE, 2020. 1
- [5] Yonghu Chen, Dongchen Zhu, Wenjun Shi, Guanghui Zhang, Tianyu Zhang, Xiaolin Zhang, and Jiamao Li. Mfcflow: A motion feature compensated multi-frame recurrent network for optical flow estimation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5068–5077, 2023. 1, 3, 6
- [6] Hosik Choi, Byungmun Kang, and DaeEun Kim. Moving object tracking based on sparse optical flow with moving window and target estimator. *Sensors*, 22(8):2878, 2022. 1
- [7] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. *Advances in Neural Information Processing Systems*, 34:9355–9366, 2021. 3, 5
- [8] 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 *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015. 2, 5
- [9] 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. *arXiv preprint arXiv:2203.16194*, 2022. 1, 3, 5, 6
- [10] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-flownet: A lightweight convolutional neural network for optical flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8981–8989, 2018. 2
- [11] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. A lightweight optical flow cnn—revisiting data fidelity and regularization. *IEEE transactions on pattern analysis and machine intelligence*, 43(8):2555–2569, 2020. 6
- [12] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *Proceed-*

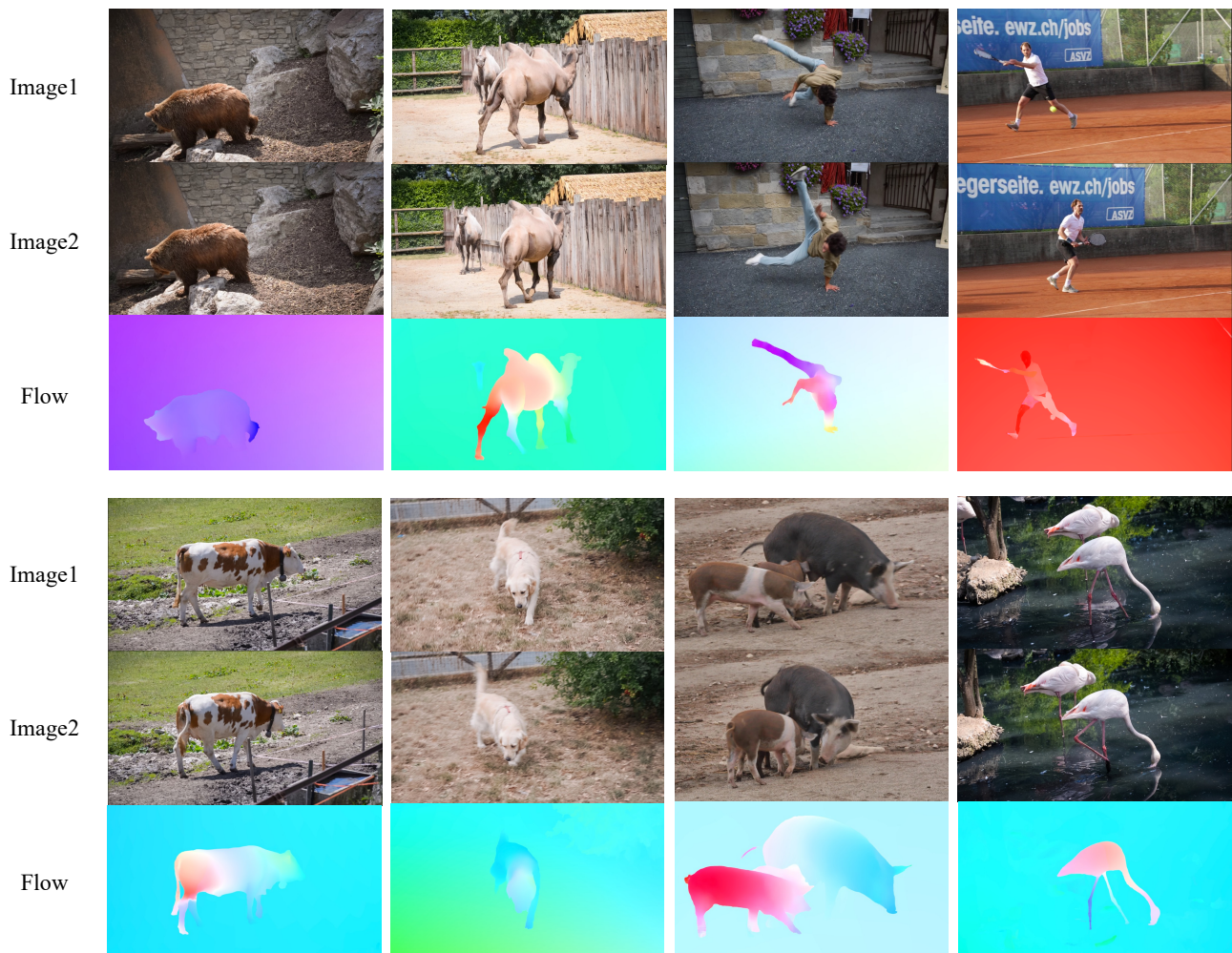


Figure 6. Visualizations of predicted flows on DAVIS [28]. StreamFlow demonstrates robust generalization to other real-world datasets, performing well in challenging scenarios for optical flow estimation, as evidenced by instances such as the occluded hind legs of the bear in the first column and the small tennis ball in the last column.

Method	Sintel (Clean)	Sintel (Final)	KITTI (EPE)	KITTI (Fl-all)
Default	0.91	2.20	4.05	13.44
Zero-init	0.93	2.15	3.92	12.36

Table 4. Comparison of different ways of initialization. All models are trained under the FlyingThings.

- ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5754–5763, 2019. 6
- [13] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 6
- [14] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9772–9781, 2021. 1, 2, 3, 4, 5, 6
- [15] Yang Jiao, Guangming Shi, and Trac D Tran. Optical flow estimation via motion feature recovery. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 2558–2562. IEEE, 2021. 1, 3, 6
- [16] Kiran Kale, Sushant Pawar, and Pravin Dhulekar. Moving object tracking using optical flow and motion vector estimation. In *2015 4th international conference on reliability, infocom technologies and optimization (ICRITO)(trends and*

- future directions), pages 1–6. IEEE, 2015. 1
- [17] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 19–28, 2016. 5, 1
- [18] Jiahao Li, Bin Li, and Yan Lu. Deep contextual video compression. *Advances in Neural Information Processing Systems*, 34:18114–18125, 2021. 1
- [19] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 5
- [20] Guo Lu, Wanli Ouyang, Dong Xu, Xiaoyun Zhang, Chunlei Cai, and Zhiyong Gao. Dvc: An end-to-end deep video compression framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11006–11015, 2019. 1
- [21] Yawen Lu, Qifan Wang, Siqi Ma, Tong Geng, Yingjie Victor Chen, Huaijin Chen, and Dongfang Liu. Transflow: Transformer as flow learner. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18063–18073, 2023. 1, 3, 4, 6, 7, 8
- [22] Ao Luo, Fan Yang, Kunming Luo, Xin Li, Haoqiang Fan, and Shuaicheng Liu. Learning optical flow with adaptive graph reasoning. *arXiv preprint arXiv:2202.03857*, 2022. 2, 4, 6
- [23] Ao Luo, Fan Yang, Xin Li, Lang Nie, Chunyu Lin, Haoqiang Fan, and Shuaicheng Liu. Gaflow: Incorporating gaussian attention into optical flow. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9642–9651, 2023. 2, 4, 6
- [24] N. Mayer, E. Ilg, P. Häusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. arXiv:1512.02134. 5, 1
- [25] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015. 5, 8, 1
- [26] Michal Neoral, Jan Šochman, and Jiří Matas. Continual occlusion and optical flow estimation. In *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part IV 14*, pages 159–174. Springer, 2019. 3
- [27] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5
- [28] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017. 8, 1, 2
- [29] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *proceedings of the IEEE International Conference on Computer Vision*, pages 5533–5541, 2017. 7
- [30] Zhile Ren, Orazio Gallo, Deqing Sun, Ming-Hsuan Yang, Erik B Sudderth, and Jan Kautz. A fusion approach for multi-frame optical flow estimation. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 2077–2086. IEEE, 2019. 3
- [31] Hao Shi, Yifan Zhou, Kailun Yang, Xiaoting Yin, and Kaiwei Wang. Csflo: Learning optical flow via cross strip correlation for autonomous driving. In *2022 IEEE Intelligent Vehicles Symposium (IV)*, pages 1851–1858. IEEE, 2022. 1
- [32] Xiaoyu Shi, Zhaoyang Huang, Weikang Bian, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Videoflow: Exploiting temporal cues for multi-frame optical flow estimation. *arXiv preprint arXiv:2303.08340*, 2023. 1, 2, 3, 4, 5, 6
- [33] Xiaoyu Shi, Zhaoyang Huang, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Flowformer++: Masked cost volume autoencoding for pretraining optical flow estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1599–1610, 2023. 3, 8
- [34] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, page 1100612. International Society for Optics and Photonics, 2019. 5
- [35] 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 *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 17602–17611, 2022. 3, 6
- [36] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8934–8943, 2018. 2, 3
- [37] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Models matter, so does training: An empirical study of cnns for optical flow estimation. *IEEE transactions on pattern analysis and machine intelligence*, 42(6):1408–1423, 2019. 2, 6
- [38] Shangkun Sun, Yuanqi Chen, Yu Zhu, Guodong Guo, and Ge Li. Skflow: Learning optical flow with super kernels. *Advances in Neural Information Processing Systems*, 35: 11313–11326, 2022. 1, 2, 3, 4, 5, 6, 8
- [39] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020. 1, 2, 3, 4, 5, 6
- [40] Bo Wang, Yifan Zhang, Jian Li, Yang Yu, Zhenping Sun, Li Liu, and Dewen Hu. Splatflow: Learning multi-frame optical flow via splatting. *arXiv preprint arXiv:2306.08887*, 2023. 3
- [41] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching. In *Proceedings of the IEEE/CVF conference on*

- computer vision and pattern recognition*, pages 8121–8130, 2022. [1](#), [2](#), [3](#), [6](#), [8](#)
- [42] Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow. *Advances in neural information processing systems*, 32, 2019. [2](#), [6](#)
- [43] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6044–6053, 2019. [6](#)
- [44] Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable flow: Learning motion cost volumes for optical flow estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10807–10817, 2021. [4](#), [6](#)
- [45] Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. Maskflownet: Asymmetric feature matching with learnable occlusion mask. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6278–6287, 2020. [2](#), [6](#)
- [46] Shiyu Zhao, Long Zhao, Zhixing Zhang, Enyu Zhou, and Dimitris Metaxas. Global matching with overlapping attention for optical flow estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17592–17601, 2022. [1](#), [2](#), [3](#), [6](#), [8](#)