

MemFlow: Optical Flow Estimation and Prediction with Memory

Qiaole Dong and Yanwei Fu✉
 School of Data Science, Fudan University
 {qldong18, yanweifu}@fudan.edu.cn

Abstract

Optical flow is a classical task that is important to the vision community. Classical optical flow estimation uses two frames as input, whilst some recent methods consider multiple frames to explicitly model long-range information. The former ones limit their ability to fully leverage temporal coherence along the video sequence; and the latter ones incur heavy computational overhead, typically not possible for real-time flow estimation. Some multi-frame-based approaches even necessitate unseen future frames for current estimation, compromising real-time applicability in safety-critical scenarios. To this end, we present MemFlow, a real-time method for optical flow estimation and prediction with memory. Our method enables memory read-out and update modules for aggregating historical motion information in real-time. Furthermore, we integrate resolution-adaptive re-scaling to accommodate diverse video resolutions. Besides, our approach seamlessly extends to the future prediction of optical flow based on past observations. Leveraging effective historical motion aggregation, our method outperforms VideoFlow with fewer parameters and faster inference speed on Sintel and KITTI-15 datasets in terms of generalization performance. At the time of submission, MemFlow also leads in performance on the 1080p Spring dataset. Codes and models will be available at: <https://dqiaole.github.io/MemFlow/>.

1. Introduction

Optical flow, a critical area in computer vision, plays a key role in various real-world applications like video inpainting [23], action recognition [60], and video prediction [25, 69]. In essence, it captures the displacement vector field for each pixel between successive video frames. Recent advances in optical flow estimation, as highlighted by works such as FlowNet [30], PWC-Net [59], RAFT [62], SKFlow [61], FlowFormer [28], and a rethinking training approach by MatchFlow[19], have been successful. This success is attributed to advancements in model architectures [28, 59, 62] and dedicated datasets [19, 21, 44].

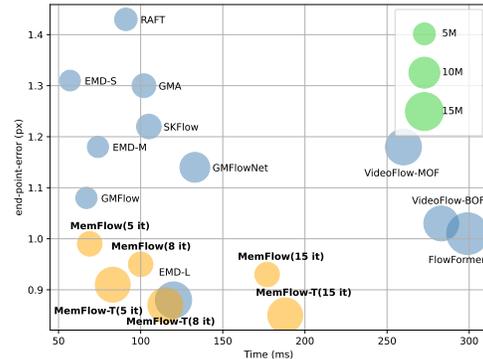


Figure 1. End-point-error on Sintel (clean) vs. inference time (ms) and model size (M). All models are trained on FlyingChairs and FlyingThings3D, and tested with one NVIDIA A100 GPU. MemFlow(-T) (x it) indicates running our network with only x iterations of GRU. Our MemFlow(-T) achieves significant reductions in computational overhead as well as substantial performance boosts over the state-of-the-art methods.

The classical optical flow works use two frames as input, potentially limiting their ability to fully leverage temporal coherence along a video sequence. This limitation results in a bottleneck, prompting an increasing reliance on computationally intensive vision transformer encoders [13] for improved performance, as noted in [28].

Conversely, some recent approaches [42, 51, 55] explore the use of multi-frame videos as input. Typically, these methods either employ simple fusion modules with modest improvements or explicitly model long-range information, incurring heavy computational overhead. For instance, PWC-Fusion [51] straightforwardly fuses backward warped past flow with current flow, resulting in a modest improvement of 0.65% over the baseline PWC-Net [59]. TransFlow [42] and VideoFlow [55] explicitly model long-range motion within a 5-frame context, leading to a significant computational overhead. Importantly, these methods [42, 55] operate in an offline mode, demanding access to unseen future frames in advance for current estimation. Additionally, VideoFlow runs considerably slower than its 2-frame baseline, SKFlow [61], as depicted

in Fig. 1. The substantial number of parameters (13.5M) also poses a significant burden on model deployment, causing out-of-memory issues when tested on the 1080p Spring dataset [45] with a single NVIDIA A100 80 GB GPU.

To this end, we present **MemFlow**, an innovative architecture by proposing the memory module [9, 10, 27, 37, 38, 47, 48, 65, 68, 71] for effective optical flow estimation. It operates in real-time (online mode) efficiently with the following notable strengths: (1) *Strong Cross-dataset Generalization Performance*. On both clean and final pass of the Sintel [6] dataset, our model achieves an end-point-error (EPE) of 0.93 and 2.08 pixels. This represents a substantial 23.8% and 15.4% error reduction compared to our 2-frame baseline, SKFlow [61] (1.22 and 2.46 pixels). When evaluated on the KITTI dataset [24], our method demonstrates an error rate of 13.7%, showcasing an 11.6% improvement over SKFlow (15.5%). (2) *High Inference Efficiency*. Our MemFlow achieves an impressive inference speed of 5.6 frames per second (fps) on A100 GPU for processing 1024x436 videos. Even a faster variant of our model, with 5 iterations of GRU, can run at 14.5 fps. Importantly, this accelerated version maintains the near-best generalization performance, as illustrated in Fig. 1.

Technically, our MemFlow is a real-time approach designed for optical flow estimation and prediction, incorporating a memory component. Specifically, MemFlow maintains a memory buffer that stores both historical motion information and context features from the input video stream. As new frame pairs are inputted, the memory buffer is continually updated with the latest context and motion features. We employ an attention mechanism to query the memory buffer using the context feature, extracting useful motion information as the aggregated motion feature. By combining this aggregated motion feature with the current motion and context features, we can regress the residual flow.

Additionally, we introduce a resolution-adaptive re-scaling for similarity computation within the attention mechanism to enhance cross-resolution generalization during inference. Furthermore, we provide the option to replace our feature encoder with a more robust vision transformer [13], referred to as **MemFlow-T**, resulting in improved outcomes. As depicted in Fig. 1, our MemFlow (-T) demonstrates significant reductions in computational overhead and substantial performance enhancements compared to state-of-the-art methods. Notably, even with only 5 iterations, our MemFlow surpasses the performance of heavyweight state-of-the-art approaches such as VideoFlow [55] and FlowFormer [28], while running faster than RAFT [62].

In addition to estimating optical flow, we’re investigating future flow prediction using our versatile memory module. This capability is crucial for intelligent systems like autonomous vehicles and robots, enabling effective planning and response in dynamic environments. Our adapted

network for flow prediction is named **MemFlow-P**. In MemFlow-P, we maintain a memory buffer for decoding residual flow from context and aggregated motion features, eliminating the need for 2D motion features between the current and next frame. Upon the arrival of the next frame, we update the memory module with a new motion feature calculated between the current and next frames. To showcase our flow prediction quality, we combine MemFlow-P with Softmax Splatting [46] and image inpainting [7, 8, 18] for video prediction, involving the synthesis of future video frames based on past ones. Despite not being specifically trained for video prediction, our method achieves comparable results with two competitive flow-based methods [25, 69] in SSIM [66] and LPIPS [73].

In summary, we make four significant contributions: 1) *Innovative Real-Time Optical Flow Estimation*: We introduce a novel architecture that effectively employs a memory module, allowing for real-time optical flow estimation. 2) *Enhanced Generalization with Resolution-Adaptive Re-scaling*: We propose the use of resolution-adaptive re-scaling in attention computation, enhancing cross-resolution generalization performance. 3) *Superior Optical Flow Estimation*: Our MemFlow(-T) achieves state-of-the-art or near-SOTA performance on various standard optical flow estimation benchmarks, demonstrating exceptional performance with minimal computational overhead. 4) *Future Prediction Capability without Explicit Training*: Repurposing MemFlow for optical flow future prediction, we achieve competitive results in video prediction without the need for specific training for this downstream task, highlighting the adaptability of our approach.

2. Related Work

Optical Flow Estimation by Two Frames. Traditionally, it is solved by optimizing energy function [2–5, 26, 52, 72] to maximize visual similarity between images. Recent efforts resort to deep neural networks for directly regressing [17, 21, 29, 30, 50, 59] or generating [20, 54] optical flow from two frames. Specifically, FlowNet [21] first proposed optical flow estimation with end-to-end trainable CNN. PWC-Net [59] and Lite-FlowNet [29] modified the network following the strategy of coarse-to-fine. RAFT [62] further developed a convolutional GRU block upon multi-scale 4D correlation volume for iteratively updating. The following works [19, 28, 34, 42, 55, 61, 75] continually improve the performance of optical flow estimation based on this recurrent architecture. Our MemFlow is also built upon and benefited from the most recently developed GRU-based network, while we employ memory to maintain past motion information and attention for gathering temporal cues for optical flow estimation.

Optical Flow Estimation by Multiple Frames. Traditional works [12, 22] employed Kalman filter for opti-

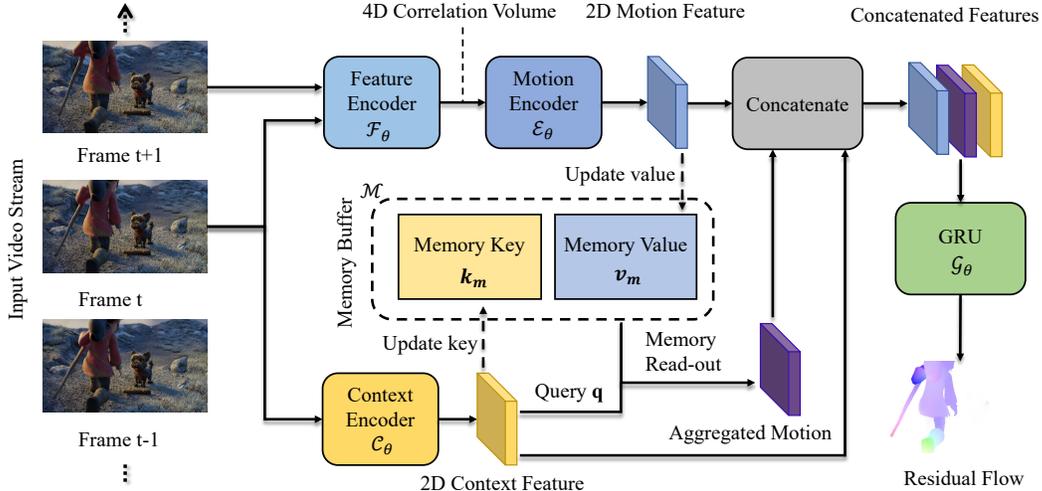


Figure 2. Overview of our MemFlow. MemFlow maintains a memory buffer to store historical motion states of video, together with an efficient update and read-out process that retrieves useful motion information for the current frame’s optical flow estimation. It has three key components: 1) *Feature Extractors*. Feature and motion encoder extract and construct the motion feature for the current frame. Another context encoder produces the context feature. 2) *Memory buffer*. Memory buffer stores historical context and motion features and read-out the aggregated motion feature. 3) *Update Modules*. GRU updates the optical flow with a series of residual flows. And the Memory buffer is kept updating when a new frame comes.

cal flow estimation with temporal coherence. Some recent un/self-supervised deep models [33, 39, 40] take three frames to estimate the optical flow of the current frame. On the other hand, there are also several supervised efforts. Particularly, PWC-Fusion [51] fused the backward warped past flow with current flow through a fusion module. RAFT [62] implicitly utilized the historical frames as the initialization to “warm-start” the optical flow estimation of current frame. TransFlow [42] and VideoFlow [55] take five frames (including both past and future frames) to better model the long-range temporal information, while this demands prohibitive computational cost and memory footprint to store and process these frames. To make a balance of performance and cost, our MemFlow integrates a memory module using past frames for optical flow, and thus is capable of being running in an online mode. Moreover, MemFlow outperforms these competitors with better generalization performance while being much more computationally friendly.

Future Prediction by Flows. It aims to predict the optical flow into the future based on past frames [35, 43], motion history [14] or even single image [1, 63, 64]. Luo *et al.* [43] firstly proposed to predict future 3D flow through a convolutional LSTM architecture. And OFNet [14] employed a UNet and ConvLSTM to predict the optical flow autoregressively based on past flows. However, Walker *et al.* [63, 64] predicted the optical flow from a single image and utilized variational autoencoders to model the uncertainty. In contrast, we repurpose our MemFlow for one time step ahead of optical flow prediction with minimal changes and achieve better prediction performance compared to recent OFNet [14] and several strongly competitive baselines.

3. MemFlow

3.1. Definition and Overview

Problem Setup. Optical flow is a per-pixel displacement vector field: $\mathbf{f}_{t \rightarrow t+1} = (f^1, f^2)$, mapping the location (u, v) in current frame \mathbf{I}_t to next frame \mathbf{I}_{t+1} as $(u + f^1(u), v + f^2(v))$. In the setting of our online multi-frame optical flow estimation, we have access to memory of past history and current frame pair: $\mathbf{I}_t, \mathbf{I}_{t+1}$. We aim to output an optical flow $\mathbf{f}_{t \rightarrow t+1}$ and update the memory accordingly. Note that the same MemFlow framework can be directly utilized to estimate optical flow for future prediction. Thus such a task is also evaluated here: we only get video frame \mathbf{I}_t and memory of past, while predicting optical flow $\mathbf{f}_{t \rightarrow t+1}$ into future.

Overview. We present an overview of our Memory module for optical Flow estimation (MemFlow) as in Fig. 2. Specifically, our MemFlow consists of three key components: 1) *Feature Extractors*. We have a feature encoder \mathcal{F}_θ extracting features from the frames to construct the 4D correlation volume subsequently. By correlation lookup operation [62], the following motion encoder \mathcal{E}_θ can produce the motion feature. To provide context features of the current frame, we also employ the context encoder \mathcal{C}_θ . 2) *Memory buffer*. We store historical context and motion features in the buffer \mathcal{M} , while only the aggregated motion feature can be read-out through an attention mechanism. 3) *Update Modules*. We iterate the modules of GRU and Memory for flow and feature updating, respectively. Typically, GRU \mathcal{G}_θ outputs a series of residual flows. And we update the memory buffer with the final optical flow.

In the next sections, we’ll elaborate on our proposed

memory module for optical flow estimation (MemFlow(-T)) in Secs. 3.2 and 3.3. Subsequently, we’ll demonstrate its application in future prediction through minimal modifications, resulting in our optical flow prediction model, MemFlow-P, as outlined in Sec. 3.4.

3.2. Memory Read-out

We will first introduce the necessary feature extraction module, then present our novel memory read-out and resolution-adaptive re-scaling for the aggregated motion feature here.

Feature Extraction. Given current input image pairs: $\mathbf{I}_t, \mathbf{I}_{t+1}$, we first extract the feature of images at 1/8 resolution by feature encoder \mathcal{F}_θ : $\mathcal{F}_\theta(\mathbf{I}_t), \mathcal{F}_\theta(\mathbf{I}_{t+1}) \in \mathbb{R}^{H \times W \times D}$, where D is the number of channel; H, W indicate the 1/8 height and width of original images. We then construct the 4D correlation volume C through the dot product between all pairs of features:

$$C = \mathcal{F}_\theta(\mathbf{I}_t) \times \mathcal{F}_\theta(\mathbf{I}_{t+1})^T \in \mathbb{R}^{H \times W \times H \times W}. \quad (1)$$

With the current estimation of optical flow f_i , which is initialized as an all zeros tensor, we can lookup correlation values from C as in [62]. Combined with the current flow, we can get the motion feature $f_m = \mathcal{E}_\theta(f_i, \text{LookUp}(C, f_i))$. Finally, we extract the context feature f_c from \mathbf{I}_t with our context encoder \mathcal{C}_θ , which is trained with the same network architecture of feature encoder \mathcal{F}_θ . Please refer to the implementation and supplementary for the network details.

Memory Read-out. We present a novel module for memory read-out. Our memory buffer $\mathcal{M} = \{k_m \in \mathbb{R}^{L \times D_k}, v_m \in \mathbb{R}^{L \times D_v}\}$, initialized from an empty set, consists of memory keys and values, where $L = l \times H \times W$ is the number of keys and values. D_k, D_v are the feature dimension. We further define l as the length of memory buffers. With current context feature f_c and motion feature f_m , we read-out the aggregated motion feature through an attention mechanism. Specifically, we first linear project f_c, f_m and get the corresponding query, key, and value by concatenation with memory buffer,

$$q = f_c W_q, \quad k = [f_c W_k; k_m], \quad v = [f_m W_v; v_m], \quad (2)$$

where W_q, W_k, W_v are the learnable projection parameters, and $[\cdot]$ is the concatenation operation along first dimension. The aggregated motion feature can be read-out by

$$f_{am} = f_m + \alpha \cdot \text{Softmax}(1/\sqrt{D_k} \times q \times k^T) \times v, \quad (3)$$

where α is a learnable scalar initialized from 0. And we omit the necessary reshape operation here for simplicity. Note that, GMA [34] utilizes attention to aggregate spatial information, while we employ the attention for gathering additional temporal information, as illustrated in Eq. (2). Furthermore, we enhance the attention for resolution adaptivity through a re-scaling technique, as explained later.

Resolution-adaptive Re-scaling. We also introduce a novel strategy for adapting resolution here. Specifically, if the model is trained using sequences up to length N , it struggles to generalize attention effectively to sequences longer than N . Pioneer work [11] found that the dilution of similarity score accounts for this. So Chiang and Cholak [11] proposed to fix this problem by scaling similarity with $\log n$, where n is the sequence length. In contrast to them, we further update the scaling with average training sequence length n_{avg} as the logarithmic base, and use the length of key k as the sequence length in our cross-attention. So the softmax function in Eq. (3) is updated as

$$\text{Softmax}\left(\frac{\log_{n_{avg}}(L + H * W)}{\sqrt{D_k}} \times q \times k^T\right). \quad (4)$$

After incorporating this novel scaling coefficient into memory read-out, it can work for various resolutions and even generalize well to 1080p video as verified in the experiment.

3.3. Memory Update and Flow Estimation

In this section, we introduce a novel memory update strategy and flow estimation with our new memory module.

Particularly, with the context, motion, and aggregated motion features, we can now output a residual flow through a GRU unit: $\Delta f_i = \text{GRU}(f_c, f_m, f_{am})$. After N iterations of GRU, we can get the final optical flow and corresponding motion feature f_m . We then update the memory buffer by inserting the transformed context and motion feature into the key and value tensors of memory,

$$k_m = [f_c W_k; k_m], \quad v_m = [f_m W_v; v_m]. \quad (5)$$

When the memory buffer length l exceeds a pre-defined maximum of l_{max} , we simply discard the obsolete features. Though we try to distill these obsolete features into long-term memory and model the long-range motion information, we find it has no effect on the final performance.

Loss Functions. Our loss functions are inherited from the classical works - RAFT [62]. Generally, we supervise our network with l_1 distance between our partially summed residual flow $\{f_1, \dots, f_N\}$ and groundtruth f_{gt} with exponentially increasing weights,

$$\mathcal{L} = \sum_{i=1}^N 0.85^{N-i} \|f_{gt} - f_i\|_1, \quad N = 12. \quad (6)$$

3.4. Beyond Flow Estimation: Future Prediction

The framework in Fig. 2 can be directly utilized for future prediction with minimal changes, and we present the modifications here. More details are in the supplementary.

As we do not have access to frame I_{t+1} , we are not able to calculate the correlation volume and encode the motion feature for the current frame. So we extract the context feature f_c from the current frame and read-out the aggregated motion feature f_{am} from the memory buffer as

in Sec. 3.2. Then we predict the optical flow by a small convolutional network based on f_c , f_{am} , and past flow f_p : $f = \text{Convs}(f_c, f_{am}, f_p)$. After the next frame comes, we can now calculate motion feature based on our predicted flow or flow estimated by MemFlow. Finally, we could update the memory buffer of MemFlow-P as in Sec. 3.3 and are ready for optical flow prediction of the next frame. We also use l_1 distance as our loss function.

We further utilize the flow prediction for video prediction. Generally, we forward warp the last video frame by Softmax Splatting [46] with monocular depth from DPT [49] and our predicted optical flow. And we will fill the holes due to splatting with image inpainting method [18].

4. Experiments

Dataset and Implementation. We adopt SKFlow [61] as the network architecture of MemFlow(-P). And we further replace the feature encoder of SKFlow with Twins-SVT [13] as MemFlow-T. The maximum length of memory buffer l_{max} is set to 1. The iteration number of GRU is set to 15 by default during inference. In order to learn better correlation and motion features, we first pre-train our network with 2-frame input on FlyingChair [21] and FlyingThings3D [44] following SKFlow. Subsequently, with 3-frame video as input, we train MemFlow(-T) and MemFlow-P on FlyingThings3D for generalization evaluation and then finetune for Sintel [6] submission with the combination of Sintel, KITTI [24], HD1K [36], FlyingThings3D. Finally, we finetune the model with KITTI and newly proposed Spring [45] for KITTI and Spring submission, respectively. Our network is trained with AdamW [41] optimizer with one-cycle [57] learning rate on two NVIDIA A100 GPUs. Further details are provided in the supplementary material.

Evaluation Metric. We utilize end-point-error (EPE) and Fl-all for Sintel and KITTI evaluation. EPE denotes the l_2 distance between estimated flow and groundtruth. And Fl-all refers to the percentage of outliers whose EPE is larger than 3 pixels and 5% of groundtruth flow magnitude. For the Spring benchmark, we also adopt 1-pixel outlier rate (1px) and WAUC [53], which is the weighted average of the inlier rates for a range of thresholds from 0 to 5 pixels. In the following tables, the best results are in bold, while the second-best ones are underlined.

4.1. Optical Flow Estimation

Generalization Performance. Following previous works, we first show the generalization performance as in Tab. 1. Our MemFlow(-T) achieves state-of-the-art zero-shot performance on both challenging datasets, even with fewer iterations and inference time as shown in Fig. 1. MemFlow with 5 iterations of GRU can even run in real-time while still keeping near-SOTA in terms of generalization. Particularly, though share the same model architecture, our MemFlow

Table 1. Generalization performance of optical flow estimation on Sintel and KITTI-15 after trained on FlyingChairs and FlyingThings3D. *MF* indicates methods using multi frames for optical flow.

Model	Sintel		KITTI-15	
	Clean	Final	Fl-epe	Fl-all
RAFT [62]	1.43	2.71	5.04	17.4
GMA [34]	1.30	2.74	4.69	17.1
GMFlow [70]	1.08	2.48	7.77	23.4
GMFlowNet [74]	1.14	2.71	4.24	15.4
SKFlow [61]	1.22	2.46	4.27	15.5
MatchFlow [19]	1.03	2.45	4.08	15.6
FlowFormer++ [56]	0.90	2.30	3.93	14.1
EMD-L [17]	0.88	2.55	4.12	13.5
TransFlow ^(MF) [42]	0.93	2.33	3.98	14.4
VideoFlow-BOF ^(MF) [55]	1.03	2.19	3.96	15.3
VideoFlow-MOF ^(MF) [55]	1.18	2.56	3.89	14.2
MemFlow (Ours) ^(MF)	0.93	<u>2.08</u>	<u>3.88</u>	13.7
MemFlow-T (Ours) ^(MF)	0.85	2.06	3.38	12.8

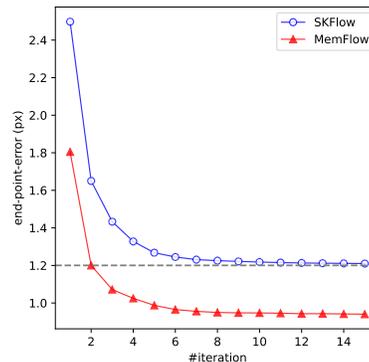


Figure 3. End-point-error of optical flow vs. number of iterations during inference. This figure provides the generalization performance on Sintel (clean) training set. Our method outperforms 15-iteration SKFlow’s performance, after using only 2 iterations.

reduces EPE by 0.38 and 0.39 from SKFlow [61] on Sintel final pass and KITTI-15, respectively. Besides, we visualize the EPE evolution over iterations in Fig. 3. With only 2 iterations, our MemFlow can beat the previous SKFlow, showing superior efficiency. Fig. 4 further provides a qualitative comparison on Sintel final pass with SKFlow and VideoFlow-MOF [55]. Our MemFlow(-T) not only exhibits more accurate details on large motion regions of hands and head but also are good at fine detail of small object.

Besides, we find that recent multi-frame based VideoFlow is typically not good at estimating optical flow for the first frame of a video. Because VideoFlow needs strict 3-frame or 5-frame input, and output optical flow for the center frame only. So they need to copy the first frame and insert it into the video as a pseudo previous frame, which accounts for why it tends to output some zero flow for the first frame as shown in the first and fourth row of Fig. 4. Yet our MemFlow can work normally without regard to the position of the input frame within the video.

Table 2. Optical flow finetuning evaluation on the public benchmark. *MF* indicates methods using multi frames for optical flow. * uses RAFT’s multi-frame "warm-start" strategy on Sintel.

Model	Sintel		KITTI-15
	Clean	Final	Fl-all
RAFT* [62]	1.61	2.86	5.10
GMA* [34]	1.39	2.47	5.15
GMFlow [70]	1.74	2.90	9.32
GMFlowNet [74]	1.39	2.65	4.79
SKFlow* [61]	1.28	2.23	4.84
MatchFlow* [19]	1.16	2.37	4.63
FlowFormer++ [56]	1.07	1.94	4.52
EMD-L [17]	1.32	2.51	4.51
PWC-Fusion ^(MF) [51]	3.43	4.57	7.17
TransFlow ^(MF) [42]	1.06	2.08	4.32
VideoFlow-BOF ^(MF) [55]	1.01	1.71	4.44
VideoFlow-MOF ^(MF) [55]	0.99	1.65	3.65
VideoFlow-MOF ^(MF) (online) [55]	-	-	4.08
MemFlow (Ours) ^(MF)	1.05	1.91	4.10
MemFlow-T (Ours) ^(MF)	1.08	1.84	<u>3.88</u>

Finetuning Evaluation. We further report the finetuning results on public benchmark datasets, Sintel and KITTI, in Tab. 2. Our MemFlow(-T) improves SKFlow by a large margin and outperforms most previous methods, e.g. SOTA 2-frame based methods FlowFormer++ [56] and multi-frame based TransFlow [42], except recent VideoFlow. However, the online version of 5-frame VideoFlow-MOF degrades Fl-all on the KITTI-15 test set from 3.65 to 4.08, which is worse than our 3.88. We suspect that using multi-frame in an offline mode explicitly can lead to better dataset-specific performance, while with limited cross-dataset generalization performance as in Tab. 1, where 5-frame VideoFlow-MOF performs much worse than 3-frame VideoFlow-BOF on Sintel. Finally, a qualitative result in Fig. 5 on the KITTI test set shows our MemFlow(-T) outperforms others in distinguishing between different vehicles and between the foreground and the sky.

Evaluation on Full-HD Spring Dataset. We also report the generalization and finetuning results on the newly proposed Full-HD (1080p) Spring benchmark in Tab. 3. We first test MemFlow trained on Sintel for evaluation of generalization. Tab. 3 shows that our MemFlow achieves the best EPE and Fl-all while being competitive with MS-RAFT+ [32] in terms of 1px within different regions. After finetuning on the Spring training set, MemFlow outperforms previous SOTA CroCo-Flow [67] by a large margin, though CroCo-Flow is pretrained with additional 5.3M real-world image pairs. Qualitative comparison in Fig. 6 also shows that MemFlow performs better on fine details, while CroCo-Flow employs the much slower tile-technique [31] for high-resolution testing and leads to block-like artifacts. Besides, we should point out that VideoFlow encounters

out-of-memory when tested on the Spring dataset with a single NVIDIA A100 80 GB GPU. This further shows that our MemFlow is much more computationally friendly.

4.2. Future Prediction of Optical Flow

For future prediction of optical flow, we compare with following three baselines: (1) *MemFlow*, we use MemFlow to estimate $\mathbf{f}_{t-1 \rightarrow t}$ with available frames and forward warp the flow to next time step as $\hat{\mathbf{f}}_{t \rightarrow t+1}$. (2) *Warped Oracle*, in contrast to (1), we forward warp the available optical flow groundtruth in dataset as a performance upper bound of the warping-based method. (3) *OFNet* [14], a learning-based method that trains a UNet and ConvLSTM with 6 past optical flows as input for future prediction. Note that all trainable models are trained on FlyingThings3D for comparison. **Flow Prediction Results.** Left part of Tab. 4 reports the EPE on test split of FlyingThings3D, training set of Sintel and KITTI-15. MemFlow-P outperforms other competitors on all three datasets by a large margin, showing great dataset-specific and cross-dataset performance. More results can be found in supplementary material.

Downstream Task: Video Prediction. We show here quantitatively that our MemFlow-P generalizes well to video prediction. We compare our method with recent two flow-based video prediction models [25, 69] on KITTI. Though without training for video prediction specifically, our method can indeed achieve comparable or even better SSIM [66] and LPIPS [73] as in the right part of Tab. 4.

4.3. Ablation study

In this section, we provide ablation studies on MemFlow by evaluating the generalization performance. We also finetune SKFlow with the same data volume as ours, denoted as SKFlow* in Tab. 5. Note that a T -frame video has $T - 1$ flow labels. So the training step varies for fair comparison.

2-Frame Pretraining. We first assess the effect of 2-frame pre-training. As in Tab. 5, pretraining substantially improves the generalization performance on Sintel and KITTI-15, except for the slightly worse EPE of Sintel final pass.

Inference Memory Length. In this experiment, MemFlow is trained using 5-frame video on FlyingThings3D, and the maximum length of memory buffer l_{max} is set to 2 at training. However, Tab. 5 shows that during inference, set l_{max} to 1 can achieve the best performance on three metrics. Compared to not using memory module during inference, enabling the memory module can improve the performance a lot, which shows the benefit of our designed memory module. But increasing l_{max} from 1 to 3 can degrade the results. We think that it’s because motion state a few frames ago is less relevant to current motion.

Training Video Length. We are also interested in how long the video we need for training the memory module. Fortunately, Tab. 5 shows that video of 3-frame is good enough

Table 3. Optical flow generalization and finetuning results on Spring [45]. We provide the 1px outlier rate for low/high-detail, (un)matched, (non-)rigid, and (not) sky regions. We also show the EPE, Fl error [24], and WAUC [53]. Important metrics are highlighted in blue.

Dataset	Model	1px												EPE	Fl	WAUC
		total	low-det.	high-det.	matched	unmat.	rigid	non-rig.	not sky	sky	s0-10	s10-40	s40+			
C+T+S+K+H	RAFT [62]	6.79	6.43	64.09	6.00	39.48	4.11	27.09	5.25	30.18	3.13	5.30	41.40	1.476	3.20	90.92
	GMA [34]	7.07	6.70	66.20	6.28	39.89	4.28	28.25	5.61	29.26	3.65	5.39	40.33	0.914	3.08	90.72
	GMFlow [70]	10.36	9.93	76.61	9.06	63.95	6.80	37.26	8.95	31.68	5.41	9.90	52.94	0.945	2.95	82.34
	FlowFormer [28]	6.51	6.14	64.22	5.77	37.29	3.53	29.08	5.50	<u>21.86</u>	3.38	5.53	35.34	0.723	2.38	91.68
	MS-RAFT+ [32]	5.72	5.37	61.50	5.04	<u>33.95</u>	3.05	25.97	4.84	19.15	2.06	5.02	<u>38.32</u>	<u>0.643</u>	2.19	92.89
	MemFlow	5.76	5.39	<u>63.35</u>	5.11	32.76	3.29	24.42	4.49	24.99	<u>2.92</u>	4.82	32.07	0.627	2.11	<u>92.25</u>
+Spring	CroCo-Flow [67]	4.57	4.21	60.59	3.85	34.20	2.19	22.50	4.48	5.87	1.23	4.33	33.13	0.498	1.51	93.66
	MemFlow	4.48	4.12	61.70	3.74	35.12	2.39	20.31	3.93	12.81	1.31	4.44	31.18	0.471	1.42	93.86

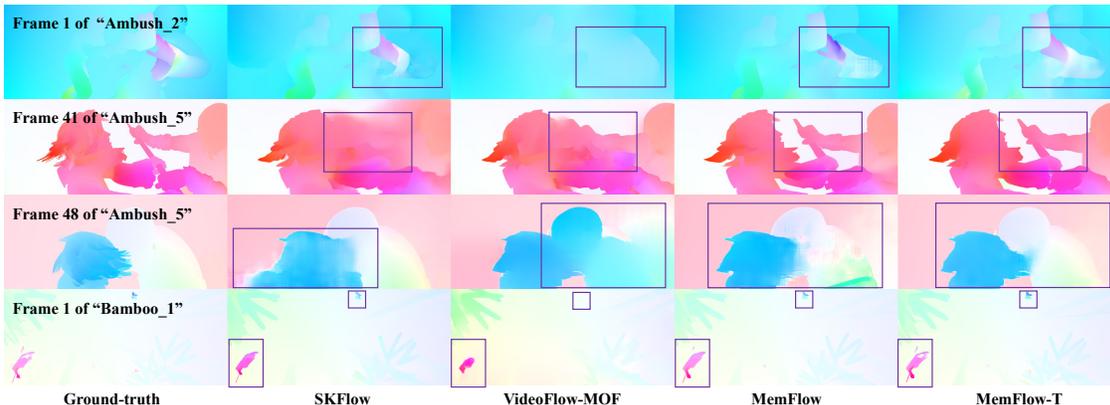


Figure 4. Qualitative comparison on the training set of Sintel final pass after pre-training on FlyingChair and FlyingThings3D. Notable areas are marked by a bounding box. Please zoom in for details.

Table 4. *Left*: End-point-error of flow prediction on FlyingThings3D (Final), Sintel (Final), and KITTI-15. *Right*: Comparison of next frame prediction on KITTI test set (256x832). Note that our method is not trained for video prediction specifically.

Method	Things	Sintel	KITTI	Method	SSIM \uparrow	LPIPS \downarrow
Warped Oracle	14.76	5.76	-	VPVFI [69]	0.827	0.123
MemFlow	15.70	6.23	12.95	VPCL [25]	0.820	0.172
OFNet [14]	<u>13.76</u>	<u>6.03</u>	<u>12.43</u>	Ours	<u>0.825</u>	0.138
MemFlow-P	7.56	5.38	8.82			

to train MemFlow. Training on longer video only results in comparable results but with much more computational overhead, e.g. GPU memory and training time.

Warm-start. Warm-start was originally proposed by RAFT for better initialization with the previous flow. It’s also compatible with MemFlow. However, equipped with our memory module, warm-start has little effect on the results as shown in Tab. 5. So we don’t use warm-start by default.

Resolution-adaptive Re-scaling. We ablate the proposed resolution-adaptive re-scaling on the newly proposed 1080p Spring dataset, which has optical flow groundtruth at a resolution of 1920x1080 and therefore an ideal testbed for research of cross-resolution generalization. MemFlow is trained at a resolution of 368x768 by default, while we also train another model at a much higher resolution of 432x960

Table 5. Ablation studies. Parameters used in our final model are underlined. * means finetuning SKFlow [61] with 1200k steps.

Experiment	Method	Sintel		KITTI-15	
		Clean	Final	Fl-epe	Fl-all
Baseline	SKFlow*	1.13	2.39	4.03	14.63
<i>Reference Model, Training: 300k on FlyingThings, max Mem is 2</i>					
2-Frame Pretraining	No	1.16	2.18	4.27	16.55
	<u>Yes</u>	1.00	2.19	3.86	14.76
Inference Mem Length	0	1.16	2.41	4.20	15.44
	<u>1</u>	1.00	2.19	3.86	14.76
	2	1.02	2.36	3.86	14.65
	3	1.07	2.25	3.87	14.64
<i>Reference Model, Training: 600k on FlyingThings, max Mem is 2</i>					
Training Video Length	3	0.93	2.08	3.88	13.71
	5	0.97	2.11	3.80	14.14
	8	0.95	2.08	3.64	14.65
Warm-start	With	1.02	2.14	3.92	13.70
	<u>Without</u>	0.93	2.08	3.88	13.71

for comparison. As shown in Tab. 6, our proposed method can substantially improve the cross-resolution generalization performance. Compared to high-resolution finetuning, our method also performs better with minimal training cost.

Discussion: why set max memory length to 1? We’ve

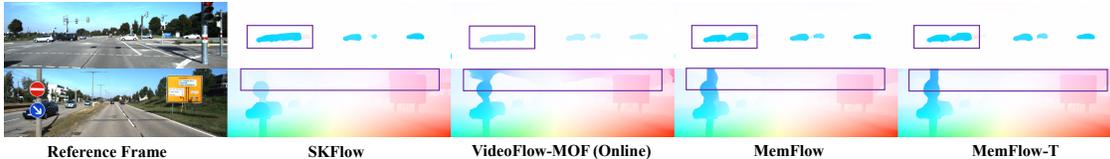


Figure 5. Qualitative comparison on test set of KITTI-15 after finetuning. Ours do much better at distinguishing between different vehicles (first row) and between the foreground and the sky (second row). Please zoom in for details.

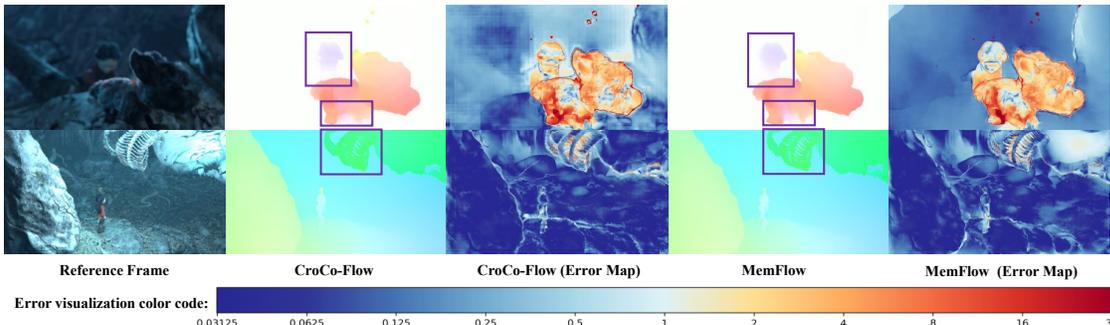


Figure 6. Qualitative comparison with CroCo-Flow [67] on Spring test set. MemFlow performs better on fine details. Notable areas are also marked by a bounding box. Please zoom in for details.

Table 6. Resolution-adaptive re-scaling substantially improves the generalization to Full-HD (*i.e.*, 1920x1080) Spring training set.

Method	1px				EPE
	total	s0-10	s10-40	s40+	
None	4.245	2.535	12.216	42.630	0.448
High-reso. ft	4.456	2.741	12.360	43.174	0.436
Ada. Re-scaling	4.211	2.512	12.113	42.404	0.433

demonstrated that optimizing the memory buffer’s maximum length to 1 achieves the best performance. Essentially, this implies that our model operates in a mode similar to a 3-frame setup. The underlying intuition behind this is that a 3-frame video represents the minimal length needed for effective temporal modeling, ensuring the most consistent motion along the time axis for the same object. Moreover, a 3-frame configuration already encompasses all the matching information required for the center frame [33]. In simpler terms, if there’s a pixel in the center frame, it’s usually visible in at least one other frame. This is probably one reason why various methods [33, 39, 40, 51, 55] opt for a 3-frame approach in flow estimation. However, it’s important to note that our MemFlow operates differently. Our memory module efficiently accumulates and updates the motion state without redundant computations over time and doesn’t explicitly extract a 3-frame sequence for estimation.

Discussion: why not a longer range? We have taken an initial step in capturing long-range motion cues for flow estimation. Specifically, we train MemFlow on an 8-frame video and enhance motion features in the memory buffer by adding relative position encoding [58] along the time axis. Despite these efforts, setting the memory length to 1 continues to yield the best results, as indicated in Tab. 7. Furthermore, we have also experimented with distilling outdated memory features into long-term memory based on historical

Table 7. More ablations about memory module. All models are trained on 8 frames video with relative position encoding here.

Experiment	Method	Sintel		KITTI-15	
		Clean	Final	Fl-epe	Fl-all
Inference Mem Length	<u>1</u>	1.05	2.12	3.65	13.78
	2	1.06	2.27	3.65	13.81
	3	1.03	2.18	3.67	13.83
Long Term Mem	With	1.05	2.12	3.65	13.78
	<u>Without</u>	1.03	2.10	3.64	13.80

attention scores, following a similar approach to XMem [9]. However, introducing long-term memory doesn’t significantly impact performance, as evidenced in Tab. 7. We believe a potential avenue for future research involves exploring long-range motion history for optical flow estimation while ensuring efficiency for real-time applications.

5. Conclusion

We introduced MemFlow, a novel online approach for video-based optical flow estimation and prediction. What sets MemFlow apart is its use of a memory module to store historical motion states. Additionally, MemFlow incorporates resolution-adaptive re-scaling, enhancing cross-resolution performance at minimal training cost. Notably, MemFlow stands out with state-of-the-art cross-dataset generalization and high inference efficiency. Besides, with minimal adjustments, MemFlow can be repurposed for flow prediction, achieving top-notch prediction performance.

Acknowledgements: Yanwei Fu is the corresponding author. Yanwei Fu is also with Shanghai Key Lab of Intelligent Information Processing, Fudan University, and Fudan ISTBI-ZJNU Algorithm Centre for Brain-inspired Intelligence, Zhejiang Normal University. The computations in this research were performed using the CFFF platform of Fudan University.

References

- [1] Dawit Mureja Argaw, Junsik Kim, Francois Rameau, Jae Won Cho, and In So Kweon. Optical flow estimation from a single motion-blurred image. In *Proceedings of the AAAI conference on Artificial Intelligence*, pages 891–900, 2021. 3
- [2] Michael J Black and Padmanabhan Anandan. A framework for the robust estimation of optical flow. In *1993 (4th) International Conference on Computer Vision*, pages 231–236. IEEE, 1993. 2
- [3] Michael J Black and Paul Anandan. The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. *Computer vision and image understanding*, 63(1):75–104, 1996.
- [4] Thomas Brox, Andrés Bruhn, Nils Papenberg, and Joachim Weickert. High accuracy optical flow estimation based on a theory for warping. In *European conference on computer vision*, pages 25–36. Springer, 2004.
- [5] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *International journal of computer vision*, 61(3):211–231, 2005. 2
- [6] 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. 2, 5
- [7] Chenjie Cao, Qiaole Dong, and Yanwei Fu. Learning prior feature and attention enhanced image inpainting. In *European Conference on Computer Vision*, pages 306–322. Springer, 2022. 2
- [8] Chenjie Cao, Qiaole Dong, and Yanwei Fu. Zits++: Image inpainting by improving the incremental transformer on structural priors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2
- [9] Ho Kei Cheng and Alexander G Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *European Conference on Computer Vision*, pages 640–658. Springer, 2022. 2, 8
- [10] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. *Advances in Neural Information Processing Systems*, 34:11781–11794, 2021. 2
- [11] David Chiang and Peter Cholák. Overcoming a theoretical limitation of self-attention. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7654–7664, Dublin, Ireland, 2022. Association for Computational Linguistics. 4
- [12] Toshio M Chin, William Clement Karl, and Alan S Willsky. Probabilistic and sequential computation of optical flow using temporal coherence. *IEEE Transactions on Image Processing*, 3(6):773–788, 1994. 2
- [13] 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. In *Advances in Neural Information Processing Systems*, pages 9355–9366. Curran Associates, Inc., 2021. 1, 2, 5
- [14] Andrea Ciamarra, Federico Becattini, Lorenzo Seidenari, and Alberto Del Bimbo. Forecasting future instance segmentation with learned optical flow and warping. In *International Conference on Image Analysis and Processing*, pages 349–361. Springer, 2022. 3, 6, 7, 2
- [15] Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. 2023. 1
- [16] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*, 2022. 1
- [17] Changxing Deng, Ao Luo, Haibin Huang, Shaodan Ma, Jiangyu Liu, and Shuaicheng Liu. Explicit motion disentangling for efficient optical flow estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9521–9530, 2023. 2, 5, 6
- [18] Qiaole Dong, Chenjie Cao, and Yanwei Fu. Incremental transformer structure enhanced image inpainting with masking positional encoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11358–11368, 2022. 2, 5, 1
- [19] Qiaole Dong, Chenjie Cao, and Yanwei Fu. Rethinking optical flow from geometric matching consistent perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1337–1347, 2023. 1, 2, 5, 6
- [20] Qiaole Dong, Bo Zhao, and Yanwei Fu. Open-ddvm: A reproduction and extension of diffusion model for optical flow estimation. *arXiv preprint arXiv:2312.01746*, 2023. 2
- [21] 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*, pages 2758–2766, 2015. 1, 2, 5
- [22] Michael Elad and Arie Feuer. Recursive optical flow estimation-adaptive filtering approach. *Journal of Visual Communication and image representation*, 9(2):119–138, 1998. 2
- [23] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf. Flow-edge guided video completion. In *European Conference on Computer Vision*, pages 713–729. Springer, 2020. 1
- [24] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 2, 5, 7
- [25] Daniel Geng, Max Hamilton, and Andrew Owens. Comparing correspondences: Video prediction with correspondence-wise losses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3365–3376, 2022. 1, 2, 6, 7
- [26] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 2
- [27] Li Hu, Peng Zhang, Bang Zhang, Pan Pan, Yinghui Xu, and Rong Jin. Learning position and target consistency for memory-based video object segmentation. In *Proceedings of*

- the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4144–4154, 2021. [2](#)
- [28] 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](#), [2](#), [7](#)
- [29] 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](#)
- [30] 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*, pages 2462–2470, 2017. [1](#), [2](#)
- [31] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppala, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021. [6](#)
- [32] Azin Jahedi, Maximilian Luz, Lukas Mehl, Marc Rivinius, and Andrés Bruhn. High resolution multi-scale raft (robust vision challenge 2022). *arXiv preprint arXiv:2210.16900*, 2022. [6](#), [7](#)
- [33] Joel Janai, Fatma Guney, Anurag Ranjan, Michael Black, and Andreas Geiger. Unsupervised learning of multi-frame optical flow with occlusions. In *Proceedings of the European conference on computer vision (ECCV)*, pages 690–706, 2018. [3](#), [8](#)
- [34] 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. [2](#), [4](#), [5](#), [6](#), [7](#)
- [35] Xiaojie Jin, Huaxin Xiao, Xiaohui Shen, Jimei Yang, Zhe Lin, Yunpeng Chen, Zequn Jie, Jiashi Feng, and Shuicheng Yan. Predicting scene parsing and motion dynamics in the future. *Advances in neural information processing systems*, 30, 2017. [3](#)
- [36] Daniel Kondermann, Rahul Nair, Katrin Honauer, Karsten Krispin, Jonas Andrulis, Alexander Brock, Burkhard Gusefeld, Mohsen Rahimimoghadam, 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](#)
- [37] Sangmin Lee, Hak Gu Kim, Dae Hwi Choi, Hyung-Il Kim, and Yong Man Ro. Video prediction recalling long-term motion context via memory alignment learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3054–3063, 2021. [2](#)
- [38] Mingxing Li, Li Hu, Zhiwei Xiong, Bang Zhang, Pan Pan, and Dong Liu. Recurrent dynamic embedding for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1332–1341, 2022. [2](#)
- [39] Liang Liu, Jiangning Zhang, Ruifei He, Yong Liu, Yabiao Wang, Ying Tai, Donghao Luo, Chengjie Wang, Jilin Li, and Feiyue Huang. Learning by analogy: Reliable supervision from transformations for unsupervised optical flow estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6489–6498, 2020. [3](#), [8](#)
- [40] Pengpeng Liu, Michael Lyu, Irwin King, and Jia Xu. Self-low: Self-supervised learning of optical flow. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4571–4580, 2019. [3](#), [8](#)
- [41] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. [5](#)
- [42] 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](#), [2](#), [3](#), [5](#), [6](#)
- [43] Zelun Luo, Boya Peng, De-An Huang, Alexandre Alahi, and Li Fei-Fei. Unsupervised learning of long-term motion dynamics for videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2203–2212, 2017. [3](#)
- [44] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. [1](#), [5](#)
- [45] Lukas Mehl, Jenny Schmalfluss, Azin Jahedi, Yaroslava Nalivayko, and Andrés Bruhn. Spring: A high-resolution high-detail dataset and benchmark for scene flow, optical flow and stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4981–4991, 2023. [2](#), [5](#), [7](#)
- [46] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020. [2](#), [5](#), [1](#)
- [47] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9226–9235, 2019. [2](#)
- [48] Matthieu Paul, Martin Danelljan, Luc Van Gool, and Radu Timofte. Local memory attention for fast video semantic segmentation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1102–1109. IEEE, 2021. [2](#)
- [49] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. [5](#), [1](#)
- [50] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4161–4170, 2017. [2](#)
- [51] 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. [1](#), [3](#), [6](#), [8](#)
- [52] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1164–1172, 2015. [2](#)
- [53] Stephan R Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2213–2222, 2017. [5](#), [7](#)
- [54] Saurabh Saxena, Charles Herrmann, Junhwa Hur, Abhishek Kar, Mohammad Norouzi, Deqing Sun, and David J Fleet. The surprising effectiveness of diffusion models for optical flow and monocular depth estimation. *Advances in Neural Information Processing Systems*, 36, 2024. [2](#)
- [55] 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](#), [5](#), [6](#), [8](#)
- [56] 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. [5](#), [6](#)
- [57] 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*, pages 369–386. SPIE, 2019. [5](#)
- [58] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021. [8](#)
- [59] 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. [1](#), [2](#)
- [60] Shuyang Sun, Zhanghui Kuang, Lu Sheng, Wanli Ouyang, and Wei Zhang. Optical flow guided feature: A fast and robust motion representation for video action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1390–1399, 2018. [1](#)
- [61] Shangkun Sun, Yuanqi Chen, Yu Zhu, Guodong Guo, and Ge Li. Skflow: Learning optical flow with super kernels. *arXiv preprint arXiv:2205.14623*, 2022. [1](#), [2](#), [5](#), [6](#), [7](#)
- [62] 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](#), [7](#)
- [63] Jacob Walker, Abhinav Gupta, and Martial Hebert. Dense optical flow prediction from a static image. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2443–2451, 2015. [3](#)
- [64] Jacob Walker, Carl Doersch, Abhinav Gupta, and Martial Hebert. An uncertain future: Forecasting from static images using variational autoencoders. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, pages 835–851. Springer, 2016. [3](#)
- [65] Hao Wang, Weining Wang, and Jing Liu. Temporal memory attention for video semantic segmentation. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 2254–2258. IEEE, 2021. [2](#)
- [66] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [2](#), [6](#)
- [67] Philippe Weinzaepfel, Thomas Lucas, Vincent Leroy, Johann Cabon, Vaibhav Arora, Romain Brégier, Gabriela Csurka, Leonid Antsfeld, Boris Chidlovskii, and Jérôme Revaud. Croco v2: Improved cross-view completion pre-training for stereo matching and optical flow. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17969–17980, 2023. [6](#), [7](#), [8](#)
- [68] Chao-Yuan Wu, Yanghao Li, Karttikeya Mangalam, Haoqi Fan, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Memvit: Memory-augmented multiscale vision transformer for efficient long-term video recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13587–13597, 2022. [2](#)
- [69] Yue Wu, Qiang Wen, and Qifeng Chen. Optimizing video prediction via video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17814–17823, 2022. [1](#), [2](#), [6](#), [7](#)
- [70] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofghi, 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. [5](#), [6](#), [7](#)
- [71] Jiyang Yu, Jingen Liu, Liefeng Bo, and Tao Mei. Memory-augmented non-local attention for video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17834–17843, 2022. [2](#)
- [72] Christopher Zach, Thomas Pock, and Horst Bischof. A duality based approach for realtime tv-l 1 optical flow. In *Joint pattern recognition symposium*, pages 214–223. Springer, 2007. [2](#)
- [73] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. [2](#), [6](#)
- [74] 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. [5](#), [6](#)
- [75] Zihua Zheng, Ni Nie, Zhi Ling, Pengfei Xiong, Jiangyu Liu, Hao Wang, and Jiankun Li. Dip: Deep inverse patch-match for high-resolution optical flow. In *Proceedings of*

the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 8925–8934, 2022. [2](#)

MemFlow: Optical Flow Estimation and Prediction with Memory

Supplementary Material

6. Details of Future Prediction

MemFlow-P. We present an overview of our Memory module for future Prediction of optical Flow (MemFlow-P) as in Fig. 7. Specifically, given current frame \mathbf{I}_t , we should first calculate the 2D motion feature f_m with previous frame \mathbf{I}_{t-1} . We are now able to update the memory buffer with f_m and the context feature $C_\theta(\mathbf{I}_{t-1})$ from \mathbf{I}_{t-1} . Then we extract the context feature $C_\theta(\mathbf{I}_t)$ from the current frame \mathbf{I}_t , which also serves as a query and reads out the aggregated motion feature f_{am} from the memory buffer. Besides, we also forward warp the previous flow $f_{t-1 \rightarrow t}$ as a base f_p for flow prediction. Finally, we concatenate the aggregated motion feature from history, context feature from the current frame, and forward warped flow f_p for flow prediction with a simple CNN: $f = \text{Convs}(f_c, f_{am}, f_p)$. Our CNN has similar convolutional layers as the original GRU. It consists of two SKBlocks as introduced by SKFlow [61]. Each SKBlock consists of two Feed Forward Networks (FFN), two depth-wise convolutional layers, and one point-wise convolutional layer. The total parameter of our MemFlow-P is 5.1 M. Our loss function is the l_1 distance between our predicted flow and the groundtruth:

$$\mathcal{L} = \|f_{gt} - f\|_1. \quad (7)$$

MemFlow-P for Video Prediction. As shown in Fig. 8, we first predict the optical flow $f_{t \rightarrow t+1}$ for the last video frame \mathbf{I}_t . Besides, we also estimate the monocular depth from DPT [49] for the last video frame. We then utilize the Softmax Splatting [46] for forward warping the last video frame. As shown in the right part of Fig. 8, we get the splatted frame and a disocclusion mask indicating the blank regions. We finally inpaint the disocclusion region with image inpainting method ZITS [18] and get the synthesised frame $\hat{\mathbf{I}}_{t+1}$.

7. Implementation Details

Network Details. Our MemFlow shares the same network architecture with SKFlow [61]. Specifically, our feature encoder and context encoder consist of 6 residual blocks, 2 at 1/2 resolution, 1/4 resolution, and 1/8 resolution, respectively. Besides, our motion encoder and GRU are based on 6 and 2 SKBlocks as in SKFlow [61], respectively. And our MemFlow-P only replaces the GRU with a small CNN as illustrated in Sec. 6. As for our MemFlow-T, we utilize the first two stages of ImageNet-pretrained Twins-SVT [13] as our feature and context encoder.

Training Details. During training, we employ FlashAttention-2 [15, 16] for faster memory read-out.

Training Schedule. We first pre-train our networks with 2-frame in FlyingChair and FlyingThings3D for 120k (batch size 8) and 150k (batch size 6) iterations, respectively. Then, we train our networks with 3-frame and batch size 8 on the following datasets, for

- **MemFlow**, we train on FlyingThings3D for additional 600k iterations for generalization evaluation. Then, we finetune our model for 600k iterations on Sintel, KITTI, HD1K, and FlyingThings3D for Sintel submission. Finally, we finetune on KITTI for 40k and on Spring for 400k iterations, respectively.
- **MemFlow-T**, we train on FlyingThings3D for additional 600k iterations for generalization evaluation. Then, we finetune our model for 300k iterations on Sintel, KITTI, HD1K, and FlyingThings3D for Sintel submission. Finally, we finetune on KITTI for 40k iterations.
- **MemFlow-P**, we randomly initialized the newly added CNN. We then train MemFlow-P on FlyingThings3D for an additional 40k iterations for generalization evaluation. For the experiment of video prediction, we train our models on Sintel, KITTI, HD1K, and FlyingThings3D with 300k iterations.

Evaluation Protocol of Video Prediction. We evaluate the performance of video prediction on four sequences from the KITTI test set following previous works [25, 69]. The four sequences we employed are:

- "2011_09_26_drive_0060_sync",
- "2011_09_26_drive_0084_sync",
- "2011_09_26_drive_0093_sync", and
- "2011_09_26_drive_0096_sync".

Besides, as in prior works [25, 69], we use a context of T=4 past frames as input. All algorithms synthesize the next frame based on past frames.

8. More Qualitative Comparison

More qualitative results on Sintel training set and KITTI training set after pre-training on FlyingChair and FlyingThings3D are given in Figs. 9 and 10. We highlight the areas where our MemFlow(-T) achieves substantial improvements with bounding boxes, compared to previous state-of-the-art VideoFlow-MOF [55] and our baseline SKFlow [61]. Please zoom in for more details.

We also provide more qualitative results on the 1080p Spring test set as shown in Fig. 11. The qualitative results show superior cross-resolution generalization performance of our MemFlow, which is trained with the image resolution of 368x768.

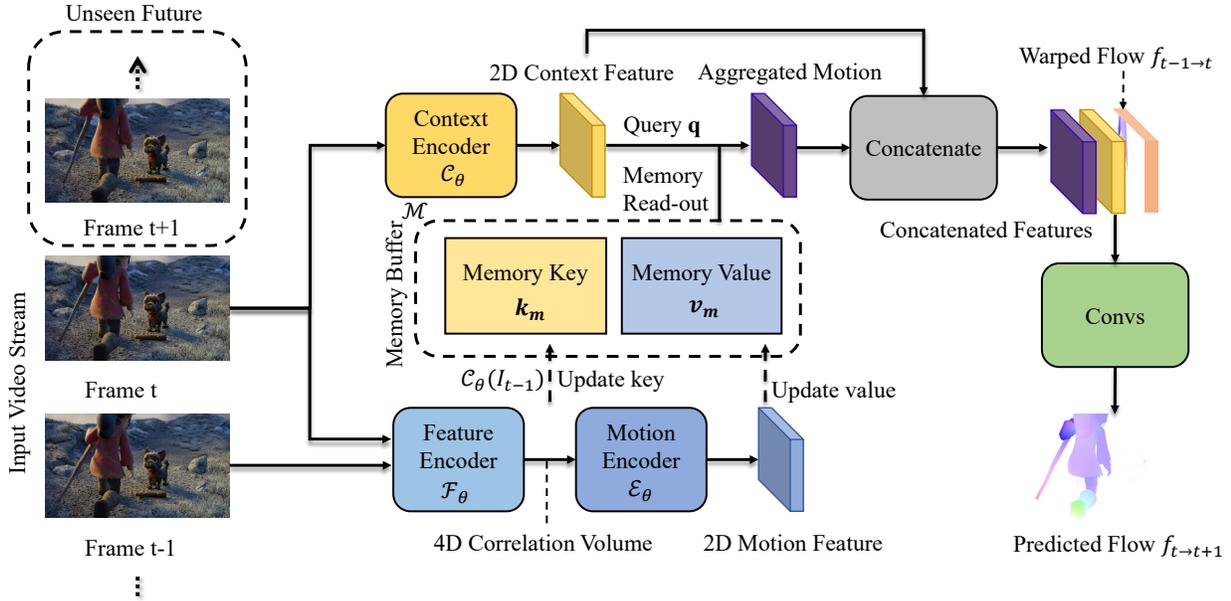


Figure 7. Overview of our MemFlow-P for future prediction of optical flow.

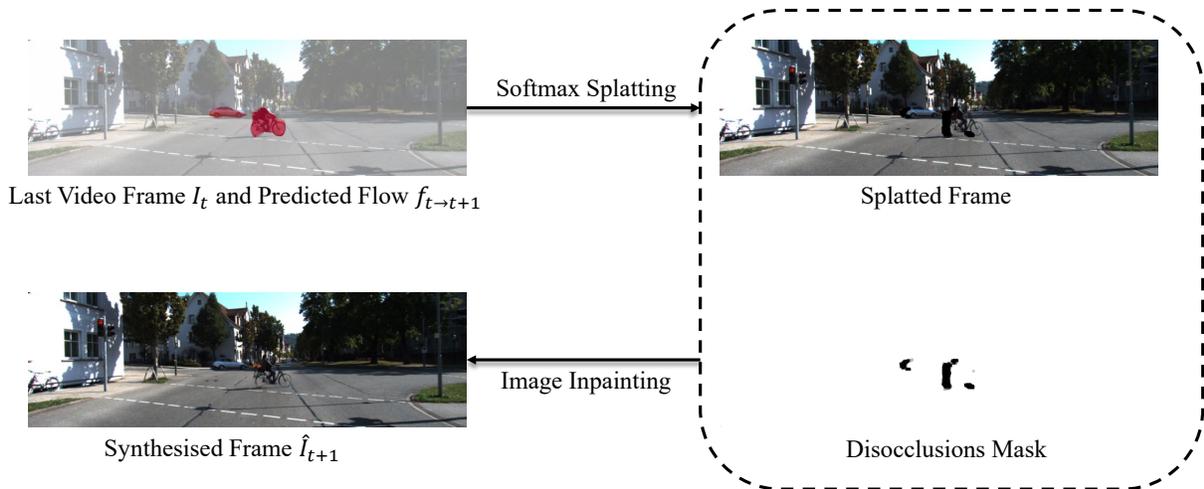


Figure 8. Overview of our MemFlow-P for video prediction.

9. More Results on Future Prediction of Optical Flow

Flow Prediction Results. We further show the full results of generalization performance evaluation for flow prediction in Tab. 8. MemFlow-P still outperforms other competitors in terms of the EPE on the clean pass of datasets and the Fl-all on KITTI-15 by a large margin, showing great dataset-specific and cross-dataset performance.

Ablation Studies. In this section, we report the ablation studies of flow prediction. First, we train a baseline model

Table 8. Generalization evaluation of flow prediction on FlyingThings3D, Sintel, and KITTI-15.

Method	Things		Sintel		KITTI-15	
	Clean	Final	Clean	Final	Fl-epe	Fl-all
Warped Oracle	14.76	14.76	5.76	5.76	-	-
MemFlow	15.55	15.70	5.92	6.23	12.95	54.48
OFNet [14]	13.73	13.76	5.78	6.03	12.43	59.17
MemFlow-P	7.81	7.56	4.97	5.38	8.82	43.93

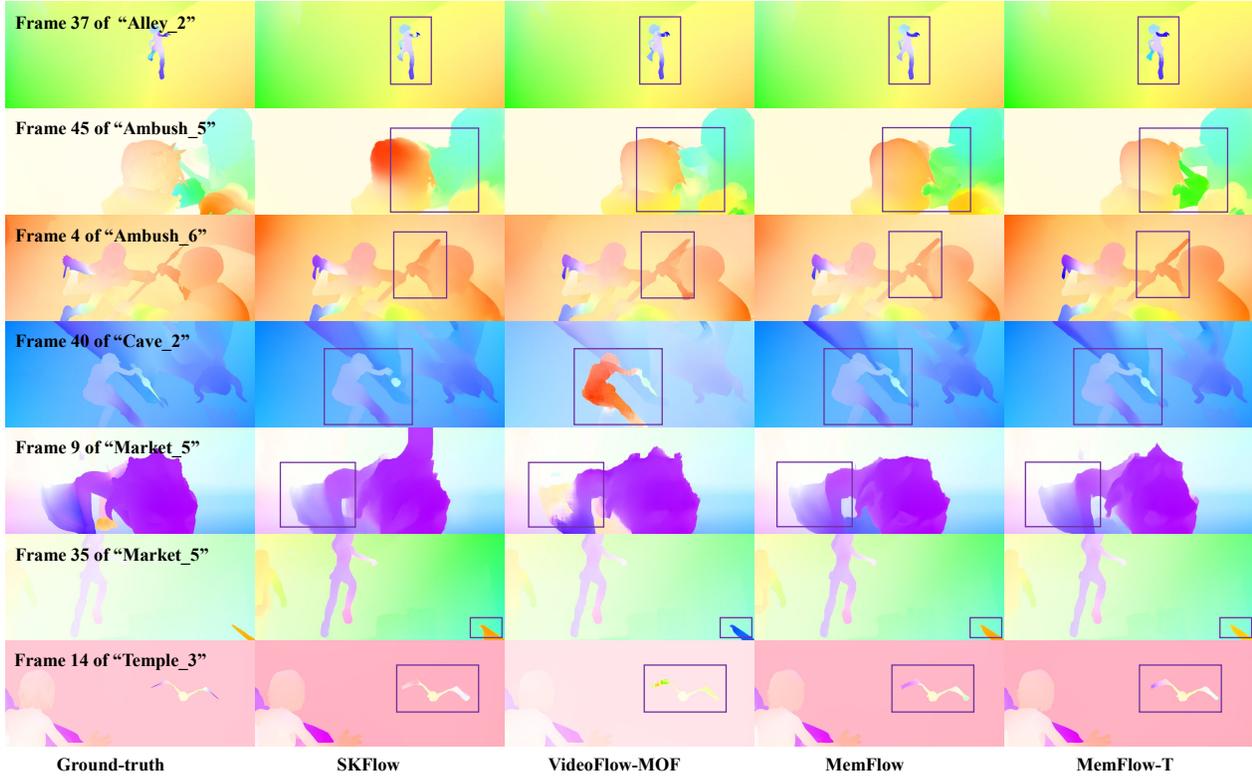


Figure 9. More qualitative results on Sintel training set final pass after pre-training on FlyingChair and FlyingThings3D. Bounding boxes mark the regions of substantial improvements. Please zoom in for details.

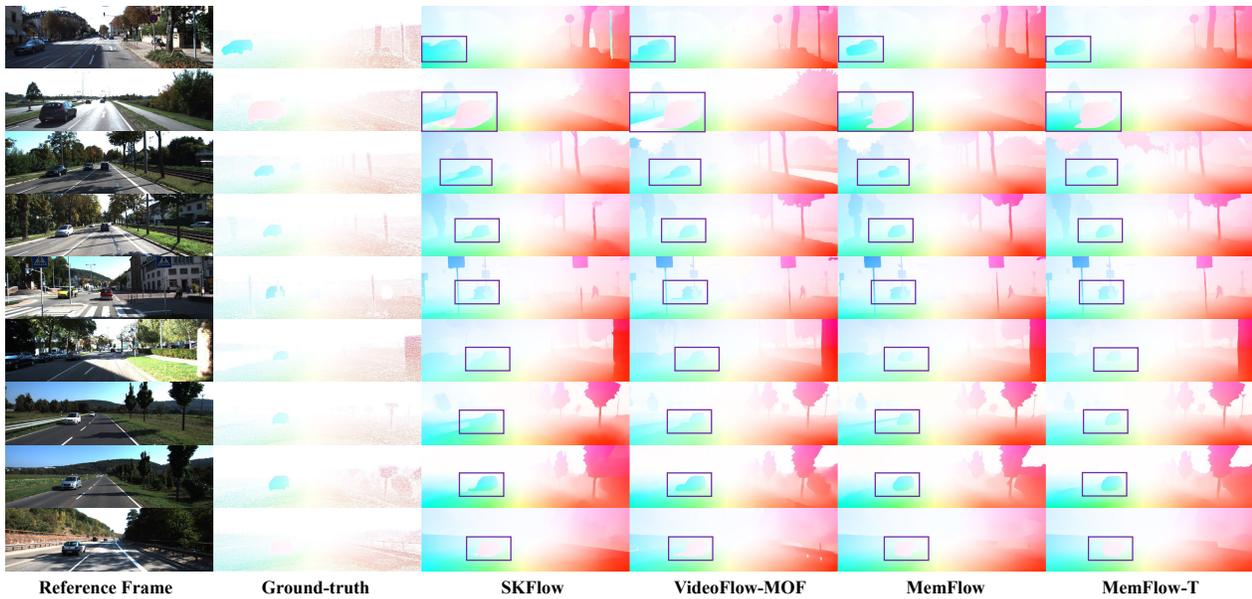


Figure 10. Qualitative results on KITTI training set after pre-training on FlyingChair and FlyingThings3D. Bounding boxes mark the regions of substantial improvements. Please zoom in for details.

for flow prediction without the forward warped past flow as input of CNN. The model is trained with 6-frame videos

sampled from FlyingThings3D. As shown in Tab. 9, concatenating the forward warped flow can improve the cross-

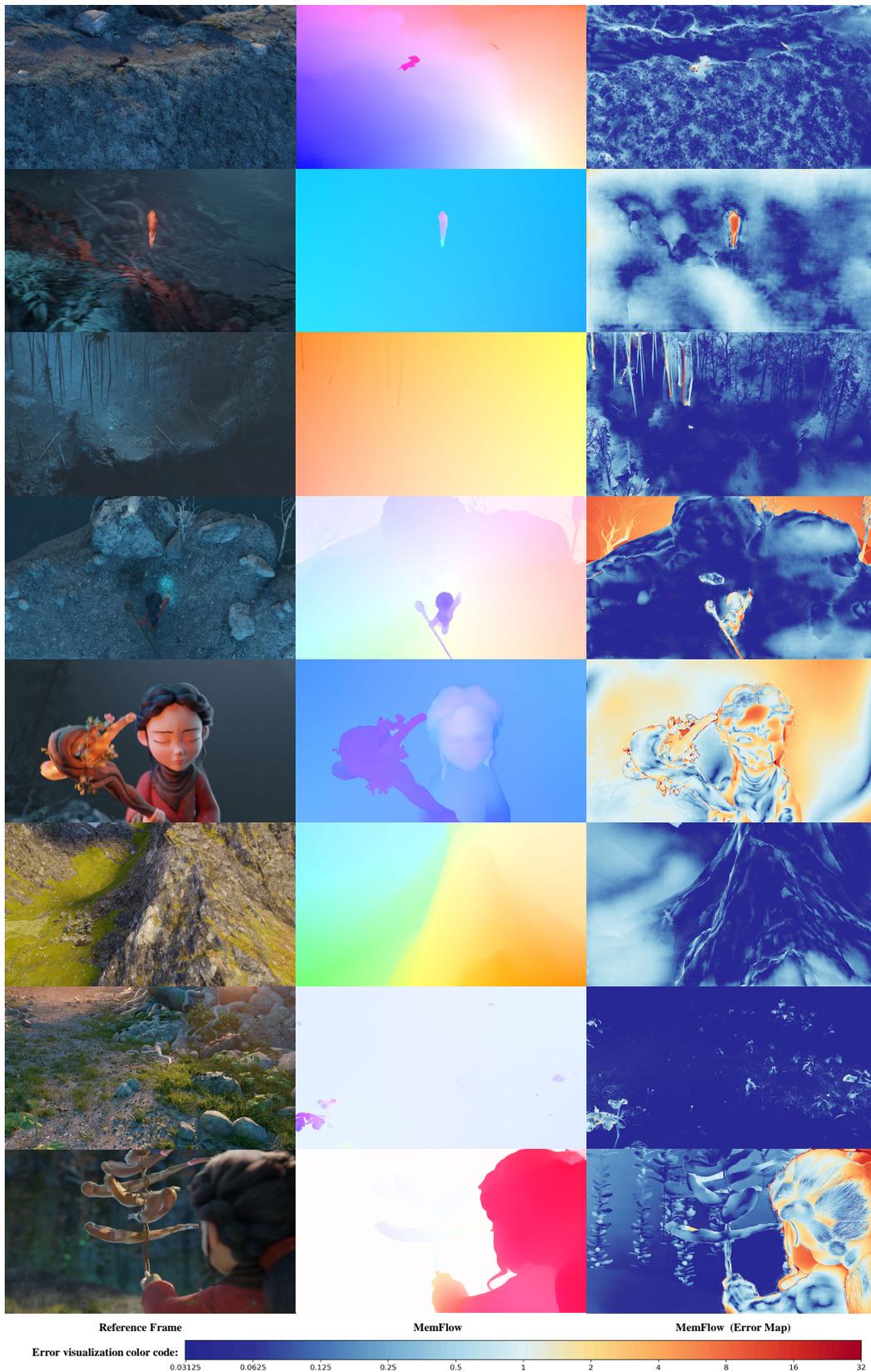


Figure 11. Qualitative results on 1080p Spring test set after finetuning on Spring. Error maps are downloaded from the official website. Please zoom in for details.

Table 9. Ablation studies on optical flow prediction.

Experiment	Things		Sintel		KITTI-15	
	Clean	Final	Clean	Final	Fl-epe	Fl-all
Baseline	7.62	6.58	5.25	5.79	8.84	56.63
+Forward Warped Flow	7.76	7.57	4.96	5.47	8.57	53.15
+Training with 3-frame	7.81	7.56	4.97	5.38	8.82	43.93

dataset generalization performance a lot, though with little worse results on the FlyingThings3D test split. Moreover, we find that training MemFlow-P with 3-frame videos can achieve similar results as the one trained with 6-frame. Therefore, we choose to train our MemFlow-P with 3-frame videos and forward warped flow.

Qualitative Results of Future Prediction by Optical Flow. We further provide several qualitative results of future prediction by optical flow as shown in Fig. 12. Our MemFlow-P can predict credible flow for the last video frame, and successfully synthesize the next frame.

Limitations of Long-term Future Prediction by Optical Flow. Our approach can generate nice results for short-term (one time step) future prediction as shown in Fig. 12. However, in the long term, as the predicted frame deviates from the distribution of training images, performance will drop quickly due to error accumulation like other video prediction methods. We further provide the quantitative and qualitative results of long-term future prediction in Figs. 13 and 14.

10. Screenshots of 1080p Spring, Sintel, and KITTI Results

We further provide anonymous screenshots of Spring, Sintel, and KITTI results on the test server as in Figs. 15 to 17. Our MemFlow ranks first on the 1080p Spring benchmark. The one without finetuning on Spring also performs well in terms of cross-dataset generalization performance. On Sintel, our MemFlow-T and MemFlow take the third and fourth places on the final pass, which improves the performance of SKFlow a lot. We also achieved great improvement on the KITTI benchmark compared to the baseline SKFlow.

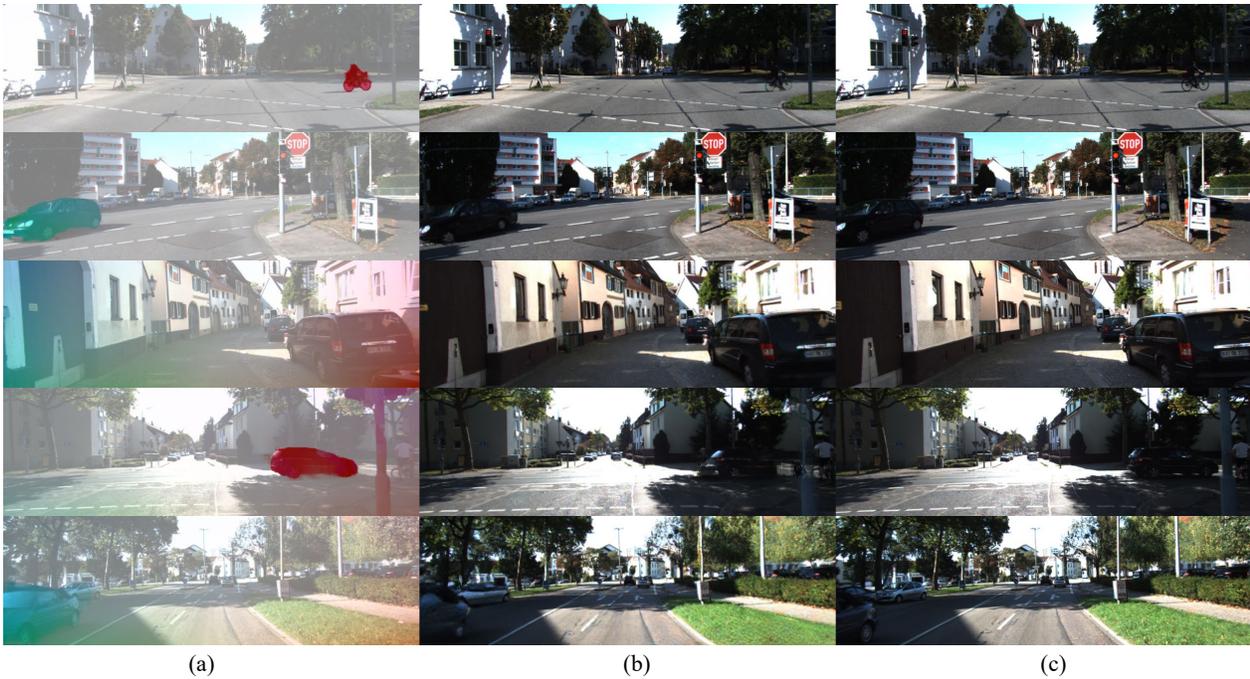


Figure 12. Qualitative Results of Future Prediction by Optical Flow. (a) Predicted optical flow superimposed on the last video frame. (b) Synthesized video frame based on our predicted flow. (c) Groundtruth next frame.

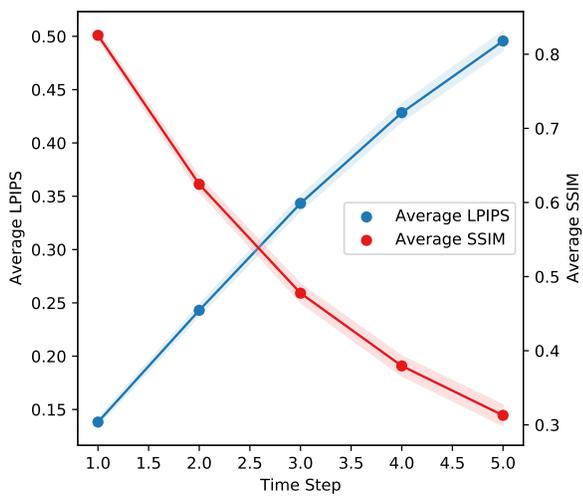


Figure 13. Quantitative results of long-term future prediction by optical flow. The plot shows the average LPIPS and SSIM-time step chart over KITTI test videos (256x832) and shadow is the 95% confidence interval. We calculate the metric with predicted frames for up to time step $T + 5$ from a context of $T=4$ past frames.

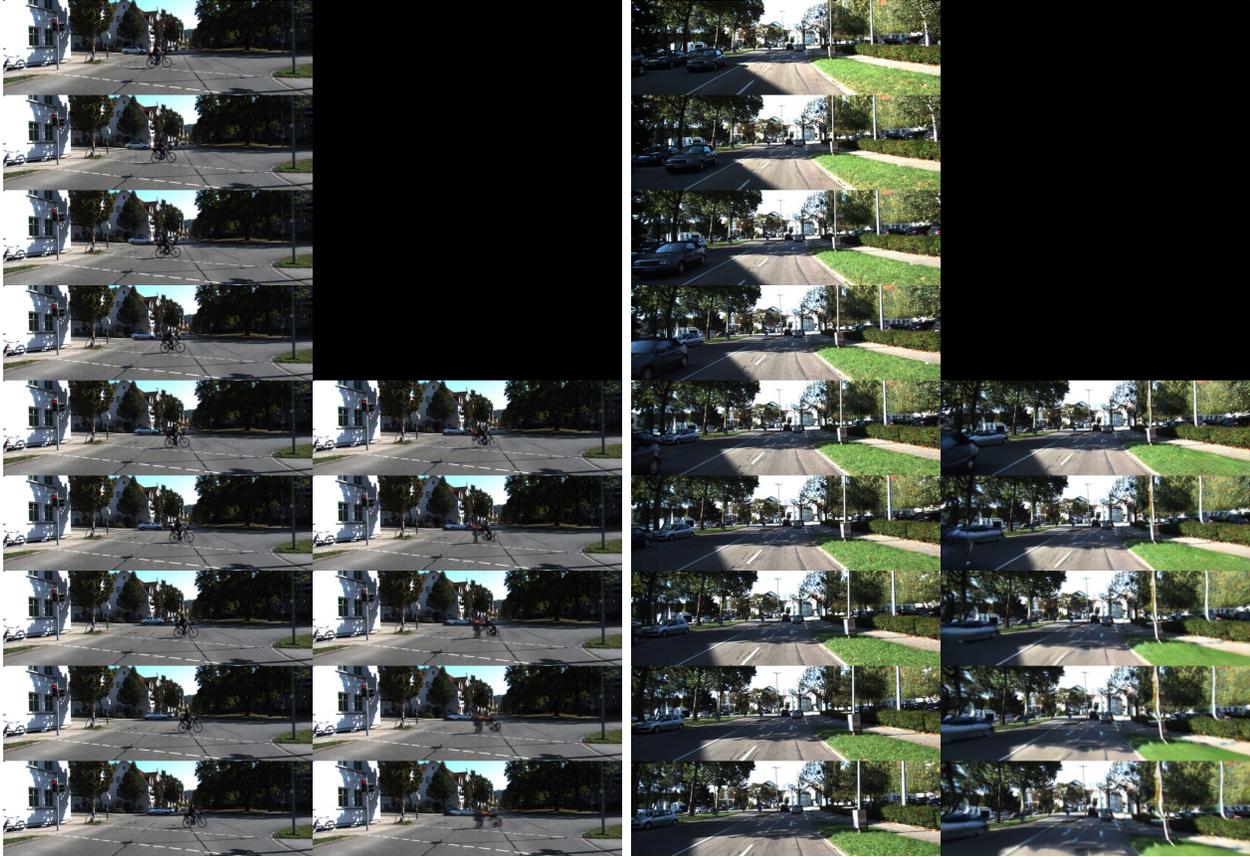


Figure 14. Our approach may fail to generate high-quality frames many steps into the future autoregressively due to error accumulation: Given 4 conditioning frames (top left), we show 5 predicted future frames in column 2 (bottom right) of two videos. Groundtruth frames are shown in the bottom left.

SPRING

Dataset & Benchmark

L. Mehl, J. Schmaljuss, A. Jahedi, Y. Nalivayko, A. Bruhn — University of Stuttgart

Download
Stereo
Optical Flow
Scene Flow
Submit
FAQ

Not logged in | [Login](#) | [Create Account](#)

Name	1px Δ total	1px low-det.	1px high-det.	1px matched	1px unmat.	1px rigid	1px non-rig.	1px not sky	1px sky	1px s0-10	1px s10-40	1px s40+	EPE	FI	WAUC
1 MemFlow <small>Anonymous.</small>	4.482	4.119	61.703	3.742	35.115	2.391	20.306	3.934	12.809	1.305	4.437	31.184	0.471	1.416	93.855
2 XCAFlow <small>Anonymous.</small>	4.493	4.145	59.236	3.853	30.966	2.105	22.559	4.291	7.570	1.727	4.605	27.333	0.506	1.566	93.163
3 CroCo-Flow <small>code</small> <small>CroCo v2: Improved Cross-view Completion Pre-training for Stereo Matching and Optical Flow. Weizsäpfel et al. ICCV 2023.</small>	4.565	4.209	60.594	3.848	34.200	2.194	22.501	4.479	5.868	1.225	4.332	33.134	0.498	1.508	93.660
4 RPKNet <small>Anonymous.</small>	4.809	4.460	59.716	4.171	31.198	2.298	23.802	4.478	9.834	1.665	4.757	31.249	0.657	1.756	92.638
5 Win-Win <small>Anonymous.</small>	5.371	5.003	63.211	4.624	36.274	2.706	25.531	4.965	11.535	1.318	4.854	40.679	0.475	1.621	92.720
6 MS-RAFT+ <small>code</small> <small>submitted by spring team A. Jahedi, M. Lutz, L. Mehl, M. Rivinius, and A. Bruhn, "High Resolution Multi-Scale RAFT." in Robust Vision Challenge, 2022.</small>	5.724	5.370	61.497	5.041	33.954	3.047	25.973	4.840	19.150	2.055	5.022	38.315	0.643	2.189	92.888
7 MemFlow(w/o ft) <small>Anonymous.</small>	5.759	5.394	63.348	5.107	32.755	3.293	24.422	4.494	24.990	2.918	4.820	32.071	0.627	2.114	92.253
8 FlowFormer <small>code</small> <small>submitted by spring team Z. Huang, X. Shi, C. Zhang, Q. Wang, K. C. Cheung, H. Qin, J. Dai, and H. Li, "FlowFormer: A Transformer Architecture for Optical Flow." in European Conference on Computer Vision (ECCV), 2022.</small>	6.510	6.144	64.219	5.766	37.294	3.527	29.084	5.500	21.858	3.381	5.530	35.344	0.723	2.384	91.679
9 FlowNet2 <small>code</small> <small>submitted by spring team E. Ilg, N. Mayer, T. Sakia, M. Keuper, A. Dosovitskiy, and T. Brox, "FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks." in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.</small>	6.710	6.346	64.061	5.691	48.892	3.711	29.404	6.039	16.908	1.862	5.816	49.693	1.040	2.823	90.907
10 RAFT <small>code</small> <small>submitted by spring team Z. Teed, and J. Deng, "RAFT: Recurrent All-Pairs Field Transforms for Optical Flow." in European Conference on Computer Vision (ECCV), 2020.</small>	6.790	6.426	64.087	5.999	39.481	4.107	27.088	5.250	30.183	3.134	5.301	41.403	1.476	3.198	90.920

Figure 15. Screenshots for 1080p Spring optical flow evaluation on the official website.

Final Clean

	EPE all	EPE matched	EPE unmatched	d0-10	d10-60	d60-140	s0-10	s10-40	s40+	
GroundTruth ^[1]	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	Visualize Results
VideoFlow-MOF ^[2]	1.649	0.788	8.660	2.090	0.609	0.334	0.403	1.243	8.804	Visualize Results
VideoFlow-BOF ^[3]	1.713	0.812	9.054	2.056	0.636	0.387	0.387	1.242	9.422	Visualize Results
MemFlow-T ^[4]	1.840	0.874	9.710	2.233	0.671	0.370	0.467	1.351	9.828	Visualize Results
MemFlow ^[5]	1.914	0.931	9.928	2.332	0.736	0.419	0.430	1.382	10.556	Visualize Results
FlowFormer++ ^[6]	1.943	0.878	10.627	2.302	0.720	0.384	0.438	1.404	10.712	Visualize Results

(a) Screenshot of Sintel Final results

Final Clean

	EPE all	EPE matched	EPE unmatched	d0-10	d10-60	d60-140	s0-10	s10-40	s40+	
GroundTruth ^[1]	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	Visualize Results
VideoFlow-MOF ^[2]	0.991	0.397	5.832	1.028	0.317	0.218	0.229	0.694	5.484	Visualize Results
SAMFlow ^[3]	0.995	0.384	5.966	1.012	0.293	0.191	0.252	0.760	5.245	Visualize Results
VideoFlow-BOF ^[4]	1.005	0.389	6.023	1.029	0.310	0.189	0.229	0.695	5.605	Visualize Results
GMFlow+ ^[5]	1.028	0.335	6.680	0.868	0.264	0.183	0.227	0.689	5.826	Visualize Results
MemFlow ^[6]	1.046	0.426	6.091	1.169	0.308	0.206	0.253	0.778	5.623	Visualize Results
GMFlow_RVC ^[7]	1.055	0.420	6.227	1.084	0.326	0.227	0.302	0.754	5.513	Visualize Results
XCAFlow ^[8]	1.057	0.340	6.908	0.904	0.248	0.194	0.186	0.623	6.432	Visualize Results
TransFlow ^[9]	1.058	0.357	6.770	0.876	0.285	0.194	0.246	0.706	5.943	Visualize Results
CCMR+ ^[10]	1.067	0.311	7.235	0.832	0.262	0.143	0.148	0.560	6.864	Visualize Results
FlowFormer++ ^[11]	1.073	0.390	6.635	1.099	0.296	0.179	0.252	0.796	5.810	Visualize Results
NA ^[12]	1.077	0.398	6.610	1.142	0.297	0.179	0.250	0.792	5.865	Visualize Results
MemFlow-T ^[13]	1.081	0.430	6.384	1.171	0.351	0.184	0.246	0.750	6.024	Visualize Results

(b) Screenshot of Sintel Clean results

Figure 16. Screenshots for Sintel optical flow evaluation on the official website.

15	MemFlow-T		3.44 %	6.09 %	3.88 %	100.00 %				<input type="checkbox"/>
16	RAFT-it+ RVC	code	3.62 %	5.33 %	3.90 %	100.00 %	0.14 s	1 core @ 2.5 Ghz (Python)		<input type="checkbox"/>
D. Sun, C. Herrmann, F. Reda, M. Rubinstein, D. Fleet and W. Freeman: Disentangling Architecture and Training for Optical Flow . ECCV 2022.										
17	RRTC		3.77 %	4.70 %	3.93 %	100.00 %	0.3 s	1 core @ 2.5 Ghz (Python)		<input type="checkbox"/>
18	RAFT-OCIC		3.72 %	5.39 %	4.00 %	100.00 %	0.2 s	GPU @ 2.5 Ghz (Python)		<input type="checkbox"/>
J. Jeong, J. Lin, F. Porikli and N. Kwak: Imposing Consistency for Optical Flow Estimation (Qualcomm AI Research) . CVPR 2022.										
19	RCA-Flow		3.67 %	6.25 %	4.10 %	100.00 %	0.16 s	1 core @ 2.5 Ghz (Python)		<input type="checkbox"/>
20	MemFlow		3.67 %	6.27 %	4.10 %	100.00 %				<input type="checkbox"/>

Figure 17. Screenshots for KITTI-15 optical flow evaluation on the official website.