

Global Matching with Overlapping Attention for Optical Flow Estimation

Shiyu Zhao^{1,*} Long Zhao² Zhixing Zhang¹ Enyu Zhou³ Dimitris Metaxas¹

¹Rutgers University

²Google Research

³SenseTime Research

Abstract

Optical flow estimation is a fundamental task in computer vision. Recent direct-regression methods using deep neural networks achieve remarkable performance improvement. However, they do not explicitly capture long-term motion correspondences and thus cannot handle large motions effectively. In this paper, inspired by the traditional matching-optimization methods where matching is introduced to handle large displacements before energy-based optimizations, we introduce a simple but effective global matching step before the direct regression and develop a learning-based matching-optimization framework, namely GMFlowNet. In GMFlowNet, global matching is efficiently calculated by applying argmax on 4D cost volumes. Additionally, to improve the matching quality, we propose patch-based overlapping attention to extract large context features. Extensive experiments demonstrate that GMFlowNet outperforms RAFT, the most popular optimization-only method, by a large margin and achieves state-of-the-art performance on standard benchmarks. Thanks to the matching and overlapping attention, GMFlowNet obtains major improvements on the predictions for textureless regions and large motions. Our code is made publicly available at <https://github.com/xiaofeng94/GMFlowNet>.

1. Introduction

Optical flow estimation is a key computer vision task, which benefits various applications, including video interpolation [25], deblurring [53], video segmentation [44] and action recognition [38]. Prevalent work in this area has been largely dominated by either matching-optimization or direct-regression methods. Previous energy-based optimization methods [6, 18, 32] usually fail to handle large displacements due to their inability to capture long-term motion correspondences. To remedy this, matching-optimization methods [2, 7, 48] introduce a matching step

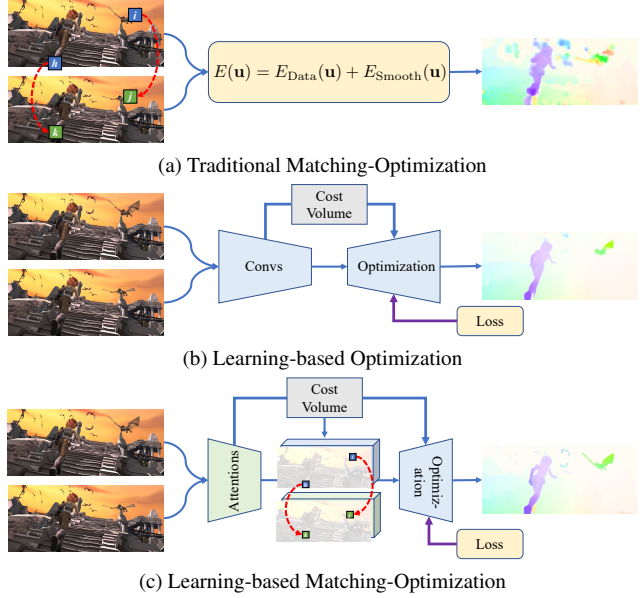


Figure 1. **Main frameworks for optical flow estimation.** (a) Traditional matching-optimization methods first build a sparse matching to get a coarse flow and then exploit energy-based optimization to refine the flow. (b) Direct-regression methods mimic the energy-based optimization with learned parameters. They can be regarded as learning-based optimizations without matching. (c) Our framework introduces matching before the learning-based optimization and further improves the performance.

before the optimization, which aims to find correspondences between pixels or patches across frames. However, their matching process depends on complicated hand-crafted features and is time-consuming and inaccurate.

Recent direct-regression methods [21, 39, 42, 55] regard optical flow estimation as a regression task and achieve considerable improvements especially in predicting small changes in optical flow. These methods typically calculate 4D cost volumes representing the similarity between pixels and then directly regress flows from cost volumes by neural networks. Similar to energy-based optimization, direct-regression methods cannot capture long-term motion correspondences in an explicit way and thus suffer from a perfor-

*Correspondence to: Shiyu Zhao (sz553@rutgers.edu).

mance drop in areas with large motions.

In this paper, we incorporate a matching step to explicitly handle large displacements for direct-regression methods, inspired by the improvement matching-optimization methods brought to energy-based optimization approaches. Based on this idea, we develop a novel framework for optical flow estimation, namely Global Matching Flow Network (GMFlowNet), where global matching is introduced before the direct regression. Unlike traditional methods, GMFlowNet provides an efficient and accurate matching step. For efficiency, we apply argmax to the typical 4D cost volume to build the global matching since it results in minor computational overhead. For accuracy, we propose a Patch-based OverLapping Attention (POLA) block to extract large context features to diminish regional ambiguities in matching, e.g., repeated patterns and textureless regions. Specifically, POLA divides input feature maps into patches and attends each patch with itself and its neighboring patches. Since direct-regression methods mimic the traditional energy-based optimizations in a data-driven manner [42], they can be interpreted as learning-based optimizations. Thus, our method can be regarded as a learning-based matching-optimization framework. Fig. 1 illustrates differences between previous related frameworks and ours.

We evaluate GMFlowNet on standard datasets for optical flow estimation. Extensive experiments demonstrate that GMFlowNet significantly outperforms the most popular optimization-only model RAFT [42] and achieves state-of-the-art performance. As expected, GMFlowNet provides better flow estimations especially for large motion areas and textureless regions. Besides, we thoroughly investigate our global matching and POLA, showing that they are both effective and efficient.

Our contributions are summarized as follows: 1) We introduce a global matching step to explicitly handle large displacement optical flow estimations for direct-regression methods. With typical 4D cost volumes, our global matching is effective and efficient. 2) We propose a well-designed Patch-based OverLapping Attention (POLA) to address local ambiguities in matching and demonstrate its effectiveness via extensive experiments. 3) Following traditional matching-optimization frameworks, we propose a learning-based matching-optimization framework named GMFlowNet that achieves state of the art performance on standard benchmarks.

2. Related Work

Optical flow as energy optimization. Previous methods formulated the optical flow as a continuous global energy function optimization problem [18]. Black and Anandan [5] introduced a robust estimation framework to address outliers caused by occlusions or significant brightness variations. Later research made further improvements by using

better regularization terms [6, 35, 54] or additional robust optimization terms [5, 8]. However, these approaches lack the ability to compute long-term dependencies and, thus, only work well for small displacements. To handle large displacements, later methods [6, 9] introduced the coarse-to-fine strategy where large and small displacements are handled at different levels of an image pyramid.

However, coarse-to-fine approaches can neither handle small and fast-moving objects that disappear at coarse levels nor remedy mistakes made in the early stages. To address those issues, Brox and Malik [7] introduced feature matching to the energy-based optimization framework, which was further improved in later works [36, 48, 50]. Following studies [1, 3, 12, 20, 49] widely adopted this approach. However, all these studies consider global matching highly time-consuming, so they only conducted local matching for computation efficiency, e.g., EpicFlow [36]. Contrary to previous methods, we calculate global matching efficiently by applying the argmax operator on widely adopted 4D cost volumes and achieve better performance.

Optical flow as network regression. More recently, the community has been motivated by the success of CNNs on high-level vision tasks [28] to exploit learning-based solutions for optical flow estimation. Relevant studies [4, 15, 17, 23, 39, 42, 55, 56] typically formulate optical flow estimation as regression instead of matching. In regression, cost volumes are the critical component that represents the similarity between pixels. For example, Sun et al. [39] designed a network using stacked image pyramids, feature warping, and cost volumes. Hofinger et al. [17] employed a sampling-based strategy to improve the calculation of cost volumes. Teed and Deng [42] built 4D cost volumes for all pairs of pixels. However, due to the high cost of memory and time, they did not aggregate the cost volumes to involve the global information. Separable Flow [55] proposed a separable cost volume module for efficient aggregations. In this work, we sidestep the high-cost global aggregation and leverage global information by constructing global matching using existing 4D cost volumes.

Attention mechanism in vision. This work extracts large context information for matching via leveraging recent advances in Vision Transformers [11, 14, 29]. Methods leveraging Transformers’ ability of modeling long-term dependencies have outperformed convolutional neural networks in various high-level computer vision tasks [14, 43]. Inspired by these, Jiang et al. [26] introduced an attention-based module to resolve occlusions for optical flow estimation. Furthermore, LoFTR [41] adopted the self- and cross-attention to extract better descriptors for feature matching. Prevailing Vision Transformer architectures, e.g., Swin Transformer [29], conduct indirect inter-patch information exchange with shifted windows. We propose POLA to exchange information across patches directly.

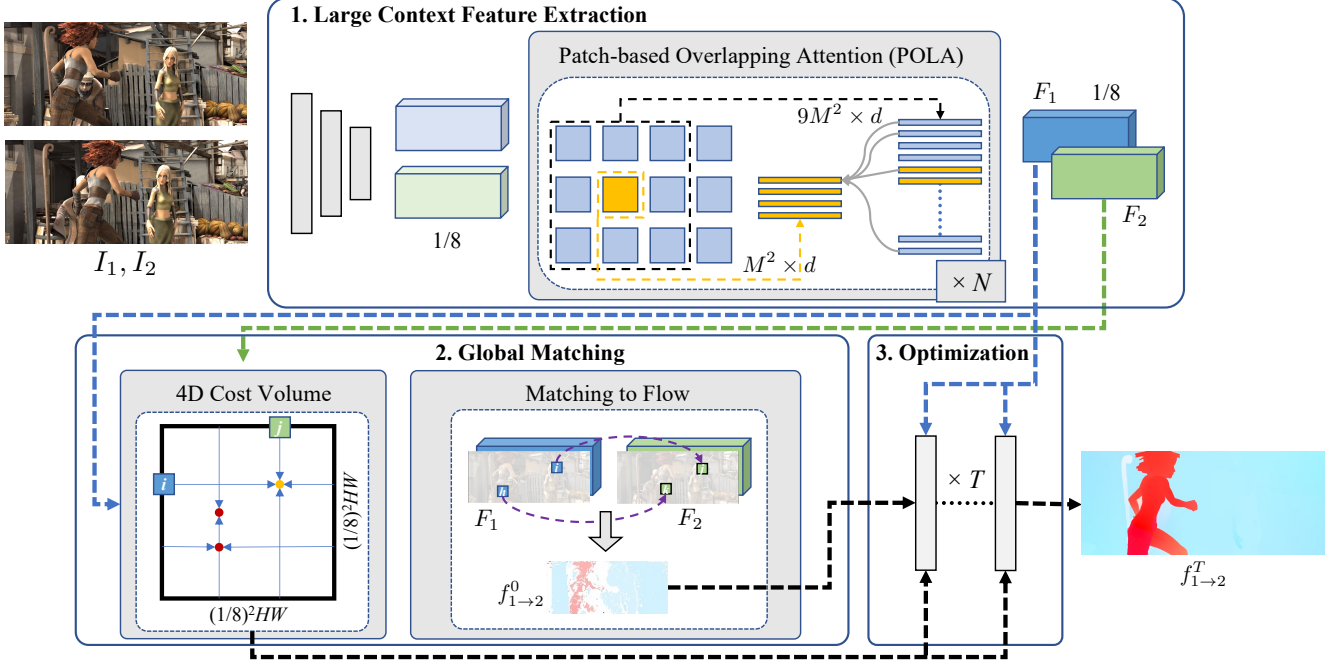


Figure 2. **Overview of GMFlowNet.** GMFlowNet has three components: **1)** The large context feature extraction module generates initial features from 3 convolutional layers and adopts the proposed POLA to extract large context information. N refers to the number of attention blocks. **2)** The global matching module adopts large context features and constructs a 4D cost volume. Then, a global matching is built by applying argmax on the cost volume and refined by mutual matching. A coarse flow $f_{1 \rightarrow 2}^0$ is generated from the matching. **3)** The optimization module takes $f_{1 \rightarrow 2}^0$ as the initial state and updates the flow estimation iteratively. T refers to the number of iterations. We employ the off-the-shelf optimization from RAFT [42].

3. Approach

We propose a novel framework GMFlowNet where a simple and effective global matching is introduced before the learning-based optimization. Our GMFlowNet consists of three modules, namely, large context feature extraction, global matching, and learning-based optimization. Fig. 2 provides an overview of GMFlowNet, and each module is elaborated in the following sections.

3.1. Large Context Feature Extraction

Large context information is the key to handle matching in locally ambiguous locations, e.g. repeated patterns and textureless regions. GMFlowNet first employs 3 convolutional layers (3-Convs) to extract initial features and then adopts Transformer blocks to include long-term dependency information. Due to the large dimension of image features, it’s computationally prohibitive to apply vanilla self-attention [45] on whole feature maps. To reduce the computation cost, we propose a well-designed local attention module POLA for optical flow estimation. In this section, we first describe attention in Transformer and then we introduce POLA. In the end, we compare POLA with other feature extractors and discuss why ours is better for our task.

Attention in Transformer. Given query vectors $Q \in \mathbb{R}^{N_q \times d}$, key vectors $K \in \mathbb{R}^{N_k \times d}$, and value vectors $V \in \mathbb{R}^{N_v \times d}$, where d is the feature dimension, attention module attends Q with V by the similarity between Q and K . Additionally, Ramachandran *et al.* [34] suggest a learned relative position bias $B \in \mathbb{R}^{N_q \times N_k}$ for better performance, and the attention is calculated as,

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d} + B) \cdot V. \quad (1)$$

For more details about Transformers, please refer to [45].

Patch-based overlapping attention. Our POLA divides features into $M \times M$ non-overlapping patches and attends every patch with itself and its eight neighboring patches. Fig. 3a illustrates our POLA with $M = 2$. Following prior work [29, 45], we adopt multi-head attentions in our attention block, as well. Given a patch vectorized as $P \in \mathbb{R}^{M^2 \times d}$ and its surrounding 3×3 patches vectorized as $S \in \mathbb{R}^{9M^2 \times d}$, for the i -th head of our attention, we first project P and S into d_k dimensions by learned linear projections and denote the projected results as P_i and S_i , respectively. Then, we perform attention with P_i and S_i and get the output h_i . Finally, we concatenate h_i from all heads as H and project H to d dimensions as the final result $O \in \mathbb{R}^{M^2 \times d}$. Our multi-head patch-based overlapping

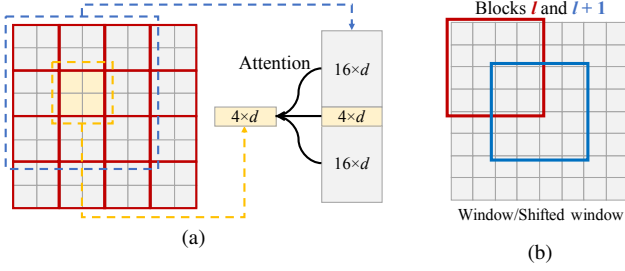


Figure 3. **Local attention.** (a) The proposed POLA. (b) Window partitions in Swin Transformer [29]. Red boxes highlight windows, and blue boxes highlight shifted windows. Two blocks are required to propagate information between windows.

attention can be formulated as,

$$\begin{aligned} h_i &= \text{Attention}(L_i^Q(P), L_i^K(S), L_i^V(S)). \\ H &= \text{Concat}([h_1, h_2, \dots, h_n]) \\ O &= L^O(H). \end{aligned} \quad (2)$$

Here n is the number of heads, L_i^Q , L_i^K , L_i^V and L^O are linear projection functions. In the experiments, we set $n = 8$ and $d_k = d/n$.

Why POLA is an improved attention method. Swin Transformer [29] provides a general local attention mechanism for vision tasks with windows and shifted windows as shown in Fig. 3b. However, the shifted window scheme requires two individual attention blocks to propagate inter-patch features, leading to information loss. Such loss is especially detrimental to matching because matching heavily depends on context information to reduce local ambiguities. By contrast, our POLA involves inter-patch features within one block and propagates information directly with less information loss. Moreover, POLA can be viewed as a generalization of per pixel overlapping attention that has been explored in [19, 34]. Compared with the per-pixel one, POLA enjoys at least three advantages: 1) consumes less memory, 2) can be efficiently implemented in existing deep learning platforms, and 3) arranges features by patch, which may provide better performance as suggested in recent research [14, 29, 43].

3.2. Global Matching

We extract the context features F_1 and F_2 for the first input image I_1 and the second input image I_2 , respectively. Then, a 4D cost volume is constructed on F_1 and F_2 . After that, a global matching is computed from the cost volume and outputs a coarse flow $f_{1 \rightarrow 2}^0$ for I_1 and I_2 , which is taken as the initial state of the later optimization.

4D cost volume calculation. We follow prior work [26, 42] to construct the 4D cost volume on $1/8$ of the input

resolution. The cost volume C is calculated as,

$$C(i, j, u, v) = F_1(i, j) \cdot F_2(u, v), \quad (3)$$

where (i, j) and (u, v) refer to locations in F_1 and F_2 .

Matching confidence calculation. We adopt a dual-softmax operator [37] to convert the cost volume into matching confidence. This operator is efficient and enables the supervision of matching. In our case, the matching confidence P_c is computed by,

$$P_c(i, j, u, v) = \text{softmax}(C(i, j, \cdot)) \odot \text{softmax}(C(\cdot, u, v)), \quad (4)$$

where $C(i, j, \cdot)$ means all (u, v) for given (i, j) . $C(\cdot, u, v)$ is similar. Pixel-wise production is denoted as \odot .

Matching selection and flow generation. Based on P_c , we obtain the matching for I_1 at (i, j) as

$$M_{1 \rightarrow 2}(i, j) = \arg \max_{u, v} P_c(i, j, u, v). \quad (5)$$

The matching for I_2 , $M_{2 \rightarrow 1}(u, v)$, is attained similarly. Then, we pick robust matches that satisfy both $M_{1 \rightarrow 2}(i, j)$ and $M_{2 \rightarrow 1}(u, v)$ and define the matching set M_c as,

$$M_c = \{(\hat{i}, \hat{j}) | (\hat{i}, \hat{j}) = M_{2 \rightarrow 1}(M_{1 \rightarrow 2}(\hat{i}, \hat{j}))\}. \quad (6)$$

The coarse flow is computed as,

$$f_{1 \rightarrow 2}^0 = \begin{cases} M_{1 \rightarrow 2}(i, j) - (i, j) & (i, j) \in M_c \\ (0, 0) & \text{Otherwise} \end{cases}. \quad (7)$$

3.3. Optimization

We use the off-the-shelf update operator from RAFT [42] as our optimization. This optimization predicts a delta flow and adds it to the current flow estimation. It iterates on such additions and outputs a series of flow predictions $\{f_{1 \rightarrow 2}^1, f_{1 \rightarrow 2}^2, \dots, f_{1 \rightarrow 2}^T\}$, where T is the total number of iterations and $f_{1 \rightarrow 2}^T$ is used as the final prediction. We initialize the optimization with our coarse flow $f_{1 \rightarrow 2}^0$ instead of the zero flow used in [42]. The optimization part in GMFlowNet is replaceable. We adopt RAFT's because it achieves the best performance. Any future optimization may be applied here for further improvements.

3.4. Supervision

Matching loss. We round the ground truth optical flow $f_{1 \rightarrow 2}^{gt}$ to the pixel level and collect the ground truth matching set M_c^{gt} . We consider regions as matched if they appear in both frames and set occlusion areas as unmatched. As the supervision in feature matching [41], we minimize the negative log-likelihood of P_c in matched regions as,

$$L_M = -\frac{1}{|M_c^{gt}|} \sum_{(\hat{i}, \hat{j}) \in M_c^{gt}} \log P_c(\hat{i}, \hat{j}) \quad (8)$$

Optimization loss. We follow RAFT [42] and supervise the optimization with ℓ_1 distance between the predicted flow and f_{gt} . The optimization loss is defined as,

$$L_O = \sum_{i=1}^T \gamma^{(i-T)} \|f_{1 \rightarrow 2}^{gt} - f_{1 \rightarrow 2}^i\|_1. \quad (9)$$

The overall loss function of GMFlowNet is,

$$L = L_O + \lambda L_M \quad (10)$$

where λ balances different loss terms.

4. Experiments

This section elaborates on the experimental results to demonstrate the effectiveness of GMFlowNet. We show that GMFlowNet improves optical flow estimation when large motions and textureless regions are present based on both quantitative and qualitative evaluations. We also discuss the improvements in the results. An ablation study and an efficiency evaluation finalize the evaluation.

We implemented GMFlowNet in PyTorch [33] and followed the training setting of RAFT [42]. We first train our model on FlyingChairs [24] (C) for 120k iterations (batch size of 10) and then finetune it on FlyingThings [30] (T) for 160k iterations (batch size of 6). After that, our model is further finetuned on a combination of data from FlyingThings (T), Sintel with both clean and final passes [10] (S), KITTI [31] (K), and/or HD1K [27] (H). In the following sections, C+T refers to FlyingChairs and FlyingThings. C+T+S/K means C+T with either Sintel or KITTI. C+T+S+K+H refers to all training datasets. We set the patch size to $M = 7$ for POLA and the feature dimension to $d = 256$. When evaluating on Sintel, we improve the model by replacing 4 heads of our POLA blocks with 2 vertical and 2 horizontal axial-attention heads that are proposed by Wang *et al.* [46].

4.1. Quantitative Evaluations

Evaluations on different displacements. Our global matching aims at addressing large motions explicitly. To evaluate its performance, we divide all regions of the Sintel training set (both clean and final passes) into different subsets, i.e., $s10$, $s10-40$, $s40+$, based on displacements. $s10$ refers to regions with displacements between 0 and 10, $s10-40$ for 10 and 40, and $s40+$ for larger than 40. Then, we train the optimization-only baseline model RAFT [42] and GMFlowNet on C+T and evaluate them on the different subsets. Table 1 provides the evaluation results in terms of average end-point-error (AEPE). As shown, for the clean pass, GMFlowNet improves RAFT by 22.4% (from 8.80 from 6.83) on $s40+$ and 18.3% (from 1.38 from 1.69) on $s10-40$. For the final pass, GMFlowNet is close to RAFT

Sintel Dataset	Type	RAFT [42] (AEPE)	Ours (AEPE)	Rel. Impr. (%)
Clean (train)	$s0-10$	0.37	0.28	24.3
	$s10-40$	1.69	1.38	18.3
	$s40+$	8.80	6.83	22.4
	All	1.47	1.14	22.4
Final (train)	$s0-10$	0.53	0.54	-1.9
	$s10-40$	3.11	3.09	0.6
	$s40+$	18.11	17.25	4.7
	All	2.78	2.71	2.5

Table 1. **Quantitative results on different displacements.** Models are trained on C+T. *Rel. Impr.* refers to relative improvement. Our method improves more on regions with extremely large motions ($s40+$) than on $s10-40$.

on $s10$ and $s10-40$ but outperforms RAFT on $s40$ by 4.7%. Those results indicate that GMFlowNet enjoys great improvements on regions with extremely large displacements, which demonstrates that the global matching with large context information is beneficial to handle large motions.

Cross-domain evaluations. Following previous studies [26, 42, 55], we trained the proposed GMFlowNet on C+T and evaluated it on the training sets of Sintel and KITTI as cross-domain evaluations. Table 2 displays the results of GMFlowNet and other competitive approaches. As a common practice, AEPE is reported for Sintel. Fl-epe and Fl-all are reported for KITTI.

As shown, GMFlowNet is close to the best method Separable Flow [55] on Sintel Final and achieves better performance on the other datasets. Our method achieves an AEPE of 1.14 on Sintel Clean, a Fl-all of 15.4 on KITTI, which are 19.6% and 11.5% better than the optimization-only baseline, RAFT. Those results demonstrate that GMFlowNet boasts a better generalization ability than RAFT as well as other methods. Considering that GMFlowNet and RAFT share the same optimization stage, we attribute the huge improvement in generalization to our global matching. We believe this is a fair claim because RAFT exploits regression but GMFlowNet considers both matching and regression. Since regression is more likely to overfit specific datasets than matching, GMFlowNet generalizes better.

Evaluations on standard benchmarks. We evaluate GMFlowNet on standard online benchmarks, i.e., Sintel [10] and KITTI [31]. For a fair comparison, we follow previous methods [26, 42, 55] and train GMFlowNet on C+T+S/K and C+T+S+K+H, respectively. Table 2 exhibits the evaluation results. GMFlowNet adopts the optimization process of RAFT, but outperforms RAFT by a large margin. Moreover, GMFlowNet outperforms the state-of-the-art method Separable Flow [55] on Sintel, but achieves slightly lower performance on KITTI. This is probably be-

Training	Method	Sintel (train)		KITTI-15 (train)		Sintel (test)		KITTI-15 (test)
		Clean	Final	F1-epe	F1-all	Clean	Final	F1-all
C+T	HD3 [52]	3.84	8.77	13.17	24.0	-	-	-
	PWC-Net [39]	2.55	3.93	10.35	33.7	-	-	-
	LiteFlowNet2 [22]	2.24	3.78	8.97	25.9	-	-	-
	VCN [51]	2.21	3.68	8.36	25.1	-	-	-
	MaskFlowNet [56]	2.25	3.61	-	23.1	-	-	-
	FlowNet2 [24]	2.02	3.54	10.08	30.0	3.96	6.02	-
	DICL-Flow [47]	1.94	3.77	8.70	23.6	-	-	-
	RAFT [42]	1.43	<u>2.71</u>	5.04	17.4	-	-	-
	GMA [26]	<u>1.30</u>	2.74	4.69	17.1	-	-	-
	Separable Flow [55]	<u>1.30</u>	2.59	<u>4.60</u>	<u>15.9</u>	-	-	-
	GMFlowNet (Ours)	1.14	<u>2.71</u>	4.24	15.4	-	-	-
C+T+S/K	FlowNet2 [24]	(1.45)	(2.01)	(2.30)	(6.8)	4.16	5.74	11.48
	HD3 [52]	(1.87)	(1.17)	(1.31)	(4.1)	4.79	4.67	6.55
	PWC-Net [39]	-	-	-	-	4.39	5.04	9.60
	LiteFlowNet [21]	(1.35)	(1.78)	(1.62)	(5.58)	4.54	5.38	9.38
	ScopeFlow [4]	-	-	-	-	3.59	4.10	6.82
	VCN [51]	(1.66)	(2.24)	(1.16)	(4.1)	2.81	4.40	6.30
	DICL-Flow [47]	(1.11)	(1.60)	(1.02)	(3.60)	2.12	3.44	6.31
	RAFT* [42]	(0.77)	(1.20)	(0.64)	(1.5)	2.08	3.41	5.27
	Separable Flow [55]	(0.71)	(1.14)	(0.68)	(1.57)	<u>1.99</u>	<u>3.27</u>	4.89
	GMFlowNet (Ours)	(0.65)	(1.06)	(0.63)	(1.49)	1.59	2.91	4.89
C+T+S+K+H	LiteFlowNet2 [22]	(1.30)	(1.62)	(1.47)	(4.8)	3.48	4.69	7.74
	PWC-Net+ [40]	(1.71)	(2.34)	(1.50)	(5.3)	3.45	4.60	7.72
	MaskFlowNet [56]	-	-	-	-	2.52	4.17	6.10
	RAFT* [42]	(0.76)	(1.22)	(0.63)	(1.5)	1.94	3.18	5.10
	GMA* [26]	(0.62)	(1.06)	(0.57)	(1.2)	<u>1.40</u>	2.88	5.15
	Separable Flow [55]	(0.69)	(1.10)	(0.69)	(1.60)	1.50	<u>2.67</u>	4.64
	GMFlowNet (Ours)	(0.59)	(0.91)	(0.64)	(1.51)	1.39	2.65	<u>4.79</u>

Table 2. **Quantitative results on Sintel and KITTI datasets.** “C+T”: We test the generalization ability on Sintel and KITTI training sets after training on FlyingChairs (C) and FlyingThing (T). “C+T+S/K”: We train models on C+T and finetune them on either Sintel (S) or KITTI (T) and evaluate on the test set of S or T. “C+T+S+K+H”: Our training set contains training samples from C, T, S, K and HD1K (H). Parentheses denote results on the training set. The best and runner up results are highlighted in bold and underlined, respectively. *We report results of the 2-view setting that is adopted by other methods.

cause GMFlowNet adopts attention blocks to extract large context features. However, KITTI only provides 200 training images that are far from enough to train high quality attention blocks. We assume that with more training data, GMFlowNet may result in larger improvements compared to CNN-based approaches.

4.2. Qualitative Evaluations

We visualize the estimated flows and cost volumes to illustrate the exact aspects that GMFlowNet improves. The supplementary document provides additional visualizations for ours coarse flow from matching.

Visualizations of estimated flows. Fig. 4 provides several test samples from KITTI and the corresponding flow estimations of RAFT and GMFlowNet. As we can see,

compared with RAFT, GMFlowNet provides better predictions on locally ambiguous regions like textureless regions. For example, there are two white cars moving forward side by side in the last row of Fig. 4. Since the two cars share similar colors and shapes, RAFT interprets them as one car. In contrast, our method succeeds in estimating the difference between the two cars and predicts the flow correctly. For more results, please refer to The supplementary Sect 7. These improvements are strong evidence of the effectiveness of the introduced global matching and POLA.

Visualizations of cost volumes. Fig. 5 visualizes the average and normalized cost volumes of both RAFT and GMFlowNet for large displacement regions (> 20 pixels). The supplementary Sect 7 provides more details about the visualization. For a fair comparison, we trained RAFT and

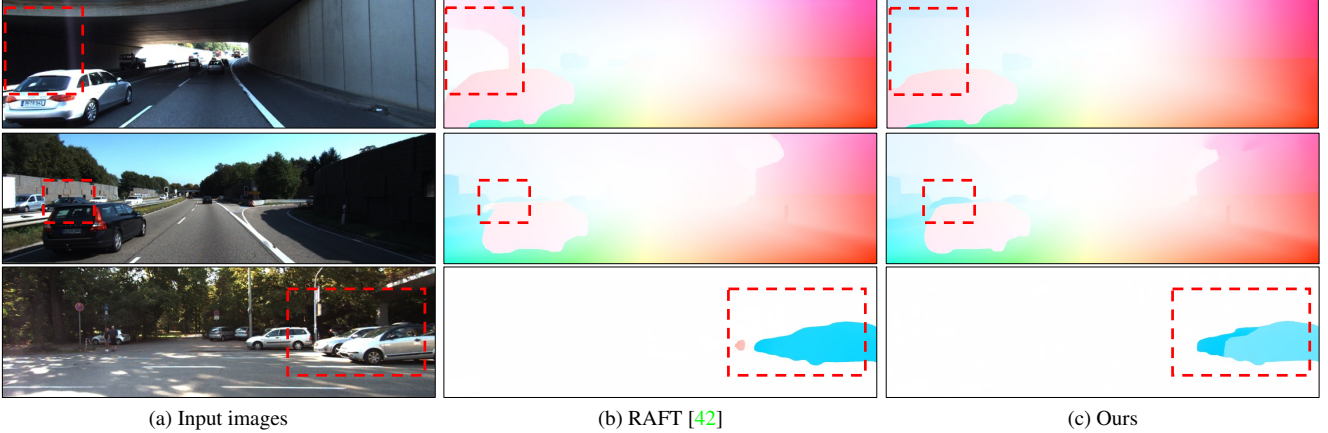


Figure 4. **Qualitative evaluations** for four samples from KITTI test set. (b) Results of the widely adopted optimization-only baseline model RAFT [42]. (c) Results of our GMFlowNet. Regions with significant improvements are highlighted by red dash boxes. GMFlowNet works better especially in textureless regions, because our overlapping attention provides more context information to diminish regional ambiguities.

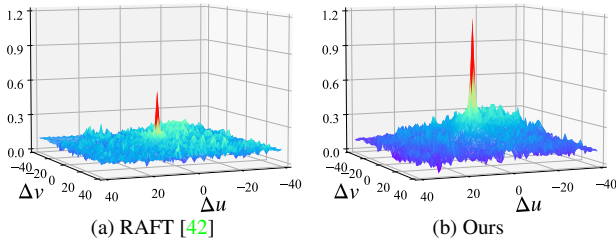


Figure 5. **Visualizations of cost volumes for large motions.** The peak of our cost volume is twice higher than that of RAFT’s, which demonstrates that our method handles large displacements better.

GMFlowNet on C+T and drew the figure using the training set of Sintel. As shown, the peak of our cost volume is much higher than that of RAFT, which clearly demonstrates that GMFlowNet is better at handling large displacements. This is plausible because our matching is designed to handle large motions and the proposed POLA extracts large context information that is crucial to overcome regional ambiguities for matching.

4.3. Ablation Study

We perform a set of ablation studies to show the importance and effectiveness of each component in GMFlowNet. All models in the experiments are trained on C+T and tested on Sintel and KITTI training sets. Table 3 provides the results for various ablation experiments. In each section of the table, we study a specific component of our approach in isolation and underline the settings used in our final model.

Initial feature extraction. We tried three different modules, i.e., None, ResNet [16], and 3-Convs, to extract initial features for the following POLA blocks. None means no

initial features. For this setting, we use the Swin Transformer architecture as the overall feature extractor, but we replace its attention blocks with POLA blocks. As shown in Table 3, 3-Convs achieve the best performance. This is likely because, on the one hand, attention blocks have more difficulties than CNNs to learn rich features from raw images. On the other hand, ResNet is much deeper than 3-Convs and may extract more high level features that are less useful for matching.

Large context feature extraction. To verify the effectiveness of our POLA, we compare it with ResNet, Swin Transformer [29], and ViT [14]. For ViT, we further reduce the feature maps by 4x. Otherwise, ViT will run out of memory because it takes global attentions instead of local attentions used in POLA. As shown in Table 3, POLA outperforms others by a large margin.

The number of our attention blocks. A simple way to expand GMFlowNet is to increase the number of attention blocks. Table 3 shows that more attention blocks achieve better performance, which is probably because more blocks provide larger receptive fields and better context information. However, more blocks increases the computation and memory costs. As a trade-off, we take 6 blocks finally.

Overlapping type. Our POLA can be viewed as a generalization of per pixel overlapping attention proposed in [34]. As shown in Table 3, POLA shares the same amount of parameters with the per pixel attention and outperforms it.

Global matching. Our key motivation is to introduce global matching into direct-regression methods. We remove the global matching in GMFlowNet and observe a significant performance drop on Sintel Final and KITTI shown in Table 3. Those results clearly demonstrate the effectiveness of the global matching.

Experiment	Method	Sintel (train)		KITTI-15 (train)		Parameters
		Clean	Final	F1-epe	F1-all	
Baseline [42]	-	1.43	2.71	5.04	17.4	5.3M
Initial Features	None	1.31	2.83	4.89	17.4	11.1M
	ResNet	1.27	2.86	4.44	16.3	10.0M
	<u>3-Convs</u>	1.14	2.71	4.24	15.4	9.3M
Large Context Features	ResNet	1.26	2.95	4.74	17.1	5.3M
	Swin Transformer	1.33	2.90	5.65	17.3	9.3M
	ViT	1.33	2.94	4.86	16.4	12.5M
	<u>POLA</u>	1.14	2.71	4.24	15.4	9.3M
Number of Attention Blocks	3	1.27	2.87	4.76	16.9	6.9M
	<u>6</u>	1.14	2.71	4.24	15.4	9.3M
	12	1.12	2.63	4.04	15.6	14.1M
Overlapping Type	Per pixel	1.32	2.88	5.11	16.8	9.3M
	<u>Patch-based</u>	1.14	2.71	4.24	15.4	9.3M
Global Matching	No	1.24	2.82	4.58	16.4	9.3M
	<u>Yes</u>	1.14	2.71	4.24	15.4	9.3M

Table 3. **Ablation experiments.** Settings used in the final model are underlined. See Sec. 4.3 for details.

Method	Param	Speed	Sintel Clean	KITTI F1-epe
RAFT [42]	5.3M	0.382s	1.43	5.04
RAFT+GM	5.3M	0.384s	1.26	4.74
+SWIN [29]	9.3M	0.422s	1.33	5.65
Ours	9.3M	0.500s	1.14	4.24

Table 4. **Comparisons of parameters and inference time.** All models are trained on C+T and tested on S and K. Speed measurements are evaluated on Sintel with the same platform.

4.4. Efficiency

Running time cost of our global matching. Running time cost is a major concern to adopt global matching. To address this concern, we compare the running time of the widely adopted RAFT and RAFT+GM. RAFT+GM refers to RAFT with our global matching step. As shown in Table 4, the global matching is very efficient and only takes 0.002s or 0.52% of extra time. Moreover, compared with RAFT, GMFlowNet runs slightly slower with 4M more parameters but significantly improves the performance. Therefore, the main benefit of our method is the performance improvement.

Running time cost of our overlapping attention. Our overlapping attention introduces more calculations but is not necessarily inefficient. To demonstrate this, we com-

pare our model with +SWIN in Table 4. +SWIN is a variant model where the POLA blocks are replaced with local attention blocks from Swin Transformer [29]. As shown, compared with +Swin, GMFlowNet requires 0.078s of extra time and improves the performance by 13.5% on Sintel clean pass and by 24.9% on KITTI. We believe that the overhead is acceptable given the performance improvement.

5. Conclusion

We have shown that matching improves the performance of direct-regression optical flow estimation methods in handling large displacements. We proposed a novel framework, GMFlowNet, where a global matching step is introduced before learning-based optimization. To improve the matching, we proposed a patch-based overlapping attention that extracts large context features to diminish regional ambiguities. GMFlowNet significantly improves predictions for large motions and textureless regions and achieves state-of-art performance on standard benchmark datasets. Future work may focus on addressing GMFlowNet’s limitations on running time cost and number of parameters.

Acknowledgments. We thank Samuel Schuster from NEC Laboratories America for helpful discussions. This research has been partially funded by the following grants, NSF IUCRC CARTA, ARO MURI 805491, NSF IIS-1793883, NSF CNS-1747778, NSF IIS 1763523, DOD-ARO ACC-W911NF, NSF OIA-2040638 to Dimitris Metaxas.

References

- [1] Min Bai, Wenjie Luo, Kaustav Kundu, and Raquel Urtasun. Exploiting semantic information and deep matching for optical flow. In *ECCV*, pages 154–170. Springer, 2016. [2](#)
- [2] Christian Bailer, Bertram Taetz, and Didier Stricker. Flow Fields: Dense correspondence fields for highly accurate large displacement optical flow estimation. In *ICCV*, pages 4015–4023, 2015. [1](#)
- [3] Christian Bailer, Kiran Varanasi, and Didier Stricker. CNN-based patch matching for optical flow with thresholded hinge embedding loss. In *CVPR*, pages 3250–3259, 2017. [2](#)
- [4] Aviram Bar-Haim and Lior Wolf. ScopeFlow: Dynamic scene scoping for optical flow. In *CVPR*, pages 7998–8007, 2020. [2](#), [6](#)
- [5] Michael J Black and Paul Anandan. The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. *CVIU*, 63(1):75–104, 1996. [2](#)
- [6] Thomas Brox, Andrés Bruhn, Nils Papenberg, and Joachim Weickert. High accuracy optical flow estimation based on a theory for warping. In *ECCV*, pages 25–36. Springer, 2004. [1](#), [2](#)
- [7] Thomas Brox and Jitendra Malik. Large displacement optical flow: descriptor matching in variational motion estimation. *IEEE TPAMI*, 33(3):500–513, 2010. [1](#), [2](#)
- [8] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Combining the advantages of local and global optic flow methods. In *Joint Pattern Recognition Symposium*, pages 454–462. Springer, 2002. [2](#)
- [9] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *IJCV*, 61(3):211–231, 2005. [2](#)
- [10] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, pages 611–625. Springer, 2012. [5](#), [11](#), [13](#), [14](#)
- [11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229. Springer, 2020. [2](#)
- [12] Qifeng Chen and Vladlen Koltun. Full flow: Optical flow estimation by global optimization over regular grids. In *CVPR*, pages 4706–4714, 2016. [2](#)
- [13] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014. [11](#)
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. [2](#), [4](#), [7](#)
- [15] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, pages 2758–2766, 2015. [2](#)
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. [7](#)
- [17] Markus Hofinger, Samuel Rota Bulò, Lorenzo Porzi, Arno Knapitsch, Thomas Pock, and Peter Kontschieder. Improving optical flow on a pyramid level. In *ECCV*, pages 770–786. Springer, 2020. [2](#)
- [18] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial Intelligence*, 17(1-3):185–203, 1981. [1](#), [2](#)
- [19] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *ICCV*, pages 3464–3473, 2019. [4](#)
- [20] Yinlin Hu, Rui Song, and Yunsong Li. Efficient coarse-to-fine patchmatch for large displacement optical flow. In *CVPR*, pages 5704–5712, 2016. [2](#)
- [21] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. LiteFlowNet: A lightweight convolutional neural network for optical flow estimation. In *CVPR*, pages 8981–8989, 2018. [1](#), [6](#)
- [22] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. A lightweight optical flow CNN — revisiting data fidelity and regularization. *IEEE TPAMI*, 43(8):2555–2569, 2020. [6](#)
- [23] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *CVPR*, pages 5754–5763, 2019. [2](#)
- [24] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, pages 2462–2470, 2017. [5](#), [6](#)
- [25] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super SloMo: High quality estimation of multiple intermediate frames for video interpolation. In *CVPR*, pages 9000–9008, 2018. [1](#)
- [26] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *ICCV*, pages 9772–9781, 2021. [2](#), [4](#), [5](#), [6](#)
- [27] Daniel Kondermann, Rahul Nair, Katrin Honauer, Karsten Krispin, Jonas Andrulis, Alexander Brock, Burkhard Gussefeld, Mohsen Rahimimoghaddam, Sabine Hofmann, Claus Brenner, et al. The HCI benchmark suite: Stereo and flow ground truth with uncertainties for urban autonomous driving. In *CVPRW*, pages 19–28, 2016. [5](#)
- [28] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, pages 1097–1105, 2012. [2](#)
- [29] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10019, 2021. [2](#), [3](#), [4](#), [7](#), [8](#)
- [30] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, pages 4040–4048, 2016. [5](#)

- [31] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *CVPR*, pages 3061–3070, 2015. 5, 11, 13, 15
- [32] Nils Papenberg, Andrés Bruhn, Thomas Brox, Stephan Didas, and Joachim Weickert. Highly accurate optic flow computation with theoretically justified warping. *IJCV*, 67(2):141–158, 2006. 1
- [33] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8026–8037, 2019. 5
- [34] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. In *NeurIPS*. IEEE, 2019. 3, 4, 7
- [35] René Ranftl, Kristian Bredies, and Thomas Pock. Non-local total generalized variation for optical flow estimation. In *ECCV*, pages 439–454. Springer, 2014. 2
- [36] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. EpicFlow: Edge-preserving interpolation of correspondences for optical flow. In *CVPR*, pages 1164–1172, 2015. 2
- [37] I. Rocco, M. Cimpoi, R. Arandjelović, A. Torii, T. Pajdla, and J. Sivic. Neighbourhood consensus networks. In *NeurIPS*, 2018. 4
- [38] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NeurIPS*, 2014. 1
- [39] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *CVPR*, pages 8934–8943, 2018. 1, 2, 6
- [40] 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 TPAMI*, 42(6):1408–1423, 2019. 6
- [41] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. LoFTR: Detector-free local feature matching with transformers. In *CVPR*, pages 8922–8931, 2021. 2, 4
- [42] Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field transforms for optical flow. In *ECCV*, pages 402–419. Springer, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 11, 12, 13, 15
- [43] Asher Trockman and J Zico Kolter. Patches are all you need? *arXiv preprint arXiv:2201.09792*, 2022. 2, 4
- [44] Yi-Hsuan Tsai, Ming-Hsuan Yang, and Michael J Black. Video segmentation via object flow. In *CVPR*, pages 3899–3908, 2016. 1
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 3, 11
- [46] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-DeepLab: Stand-alone axial-attention for panoptic segmentation. In *ECCV*, pages 108–126. Springer, 2020. 5
- [47] Jianyuan Wang, Yiran Zhong, Yuchao Dai, Kaihao Zhang, Pan Ji, and Hongdong Li. Displacement-invariant matching cost learning for accurate optical flow estimation. In *NeurIPS*, 2020. 6
- [48] Philippe Weinzaepfel, Jerome Revaud, Zaid Harchaoui, and Cordelia Schmid. DeepFlow: Large displacement optical flow with deep matching. In *ICCV*, pages 1385–1392, 2013. 1, 2
- [49] Jia Xu, René Ranftl, and Vladlen Koltun. Accurate optical flow via direct cost volume processing. In *CVPR*, pages 1289–1297, 2017. 2
- [50] Li Xu, Jiaya Jia, and Yasuyuki Matsushita. Motion detail preserving optical flow estimation. *IEEE TPAMI*, 34(9):1744–1757, 2011. 2
- [51] Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow. In *NeurIPS*, pages 794–805, 2019. 6
- [52] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *CVPR*, pages 6044–6053, 2019. 6
- [53] Yuan Yuan, Wei Su, and Dandan Ma. Efficient dynamic scene deblurring using spatially variant deconvolution network with optical flow guided training. In *CVPR*, pages 3555–3564, 2020. 1
- [54] Christopher Zach, Thomas Pock, and Horst Bischof. A duality based approach for realtime tv-l1 optical flow. In *Joint Pattern Recognition Symposium*, pages 214–223. Springer, 2007. 2
- [55] Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable Flow: Learning motion cost volumes for optical flow estimation. In *ICCV*, pages 10807–10817, 2021. 1, 2, 5, 6
- [56] Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. MaskFlowNet: Asymmetric feature matching with learnable occlusion mask. In *CVPR*, pages 6278–6287, 2020. 2, 6

6. Architecture Details

6.1. Large Context Feature Extraction

In the paper, we exploit Transformer blocks to extract large context features to improve the matching step in GMFlowNet. In the original Transformer block [45], input features are updated by a Multi-head Self-Attention (MSA) followed by a Multilayer perceptron (MLP). MSA is able to extract the long-term dependency, and MLP projects the features to the required dimension. Both MSA and MLP calculate residuals that are added to the input features as the output features. The update in a transformer block can be formulated as,

$$\begin{aligned}\hat{x}^l &= \text{MSA}(\text{LN}(x^{l-1})) + \hat{x}^l \\ x^l &= \text{MLP}(\text{LN}(\hat{x}^l)) + \hat{x}^l,\end{aligned}\quad (11)$$

where LN refers to layer norm, and x^{l-1} and x^l represent output features of the previous block and the current block, respectively. The MSA is originally designed for language tasks and takes the whole 1D features as input, but it is computationally prohibitive to apply it on 2D feature maps for optical flow estimation. To extract the long-term dependency with an acceptable computation cost, we propose the patch-based overlapping attention (POLA) to replace MSA of the original attention block and call our attention block as multi-head POLA (M-POLA).

In our large context feature extraction module (Section 3.1), we take 3 convolutional layers (3-Convs) to extract initial features and 6 M-POLA blocks to extract large context information based on initial features. The detailed structure of this module is listed in Table 5.

6.2. Optimization Network

We adopt the iterative update operator proposed in RAFT [42] as the optimization step of GMFlowNet. As stated in [42], this operator mimics the steps of an optimization algorithm and iteratively outputs a series of flow predictions $\{f_{1 \rightarrow 2}^{(1)}, f_{1 \rightarrow 2}^{(2)}, \dots, f_{1 \rightarrow 2}^{(T)}\}$. For the t -th iteration, the flow prediction $f_{1 \rightarrow 2}^{(t)}$ is calculated by a Convolutional GRU [13] (ConvGRU) as,

$$\begin{aligned}x^{(t)} &= [f_{1 \rightarrow 2}^{(t-1)}, F_1, \text{lookup}(C, f_{1 \rightarrow 2}^{(t-1)}, r)], \\ r^{(t)} &= \sigma(\text{Conv}([h^{(t-1)}, x^{(t)}])), \\ \tilde{h}^{(t)} &= \sigma(\text{Conv}([r^{(t)} \odot h^{(t-1)}, x^{(t)}])), \\ z^{(t)} &= \mu(\text{Conv}([h^{(t-1)}, x^{(t)}])), \\ h^{(t)} &= (1 - z^{(t)}) \odot h^{(t-1)} + z^{(t)} \odot \tilde{h}^{(t)}, \\ \Delta f_{1 \rightarrow 2}^{(t)} &= \text{Conv}(h^{(t)}), \\ f_{1 \rightarrow 2}^{(t)} &= f_{1 \rightarrow 2}^{(t-1)} + \Delta f_{1 \rightarrow 2}^{(t)}\end{aligned}\quad (12)$$

where F_1 is the context features, C is the 4D cost volume (See Section 3.2 of the paper), $\text{Conv}(\cdot)$ refers to a con-

	Layer name(s)
3-Convs	Conv(3, 64, 7, 2), ReLU
	Conv(64, 128, 3, 2), ReLU
	Conv(128, 256, 3, 2), ReLU
6 M-POLA	M-POLA (dim=256, head=8, win_size=7)
	M-POLA (dim=256, head=8, win_size=7)
	M-POLA (dim=256, head=8, win_size=7)
	M-POLA (dim=256, head=8, win_size=7)
	M-POLA (dim=256, head=8, win_size=7)
	M-POLA (dim=256, head=8, win_size=7)

Table 5. Large context feature extraction. The arguments in $\text{Conv}(\cdot)$ are the input channel number, the output channel number, the kernel size, and the convolution stride, respectively.

volution layer, $\sigma(\cdot)$ means sigmoid, and $\mu(\cdot)$ means tanh. $\text{lookup}(\cdot)$ represents the cost volume within the range of r . For each location \mathbf{x} in I_1 , $\text{lookup}(\cdot)$ is defined as,

$$\text{lookup}(\cdot) = \{C(\mathbf{x}, f_{1 \rightarrow 2}^{(t-1)}(\mathbf{x}) + \delta \mathbf{x}) \mid r > \|\delta \mathbf{x}\|_1\}. \quad (13)$$

Different iterations share the weights in the ConvGRU.

7. More Visualizations

7.1. Attention maps

Fig. 6 visualizes full attention score maps of the first POLA for three pixels highlighted in white. The more red a pixel is, the higher the score is. Yellow dash boxes indicate the local regions that are used in POLA. As shown, a pixel is more likely to attend to those that are visually similar to the pixel.

7.2. Coarse Flows

Figure 7 displays the coarse flows from our matching step as well as the final flow estimation for samples from Sintel [10] and KITTI [31] datasets. We compare our GMFlowNet with RAFT [42] because they share the same optimization architecture. For Sintel, both models are trained on C+T. For KITTI, they are trained on all the training data. As shown, the coarse flow results in better predictions especially in large motion areas and textureless regions. For example, the hand of the character in Fig. 7b moves fast, leading to failures of RAFT. On the contrary, our matching step finds the optical flow for the hand and improves the final prediction.

7.3. More Visual Results

Figure 8 provides the qualitative evaluation of GMFlowNet and RAFT on the Sintel test set. We highlight with white arrows and red dash boxes the regions where our method outperforms RAFT. Fig. 9 exhibits the visualization of more samples from the KITTI test set. Red dash

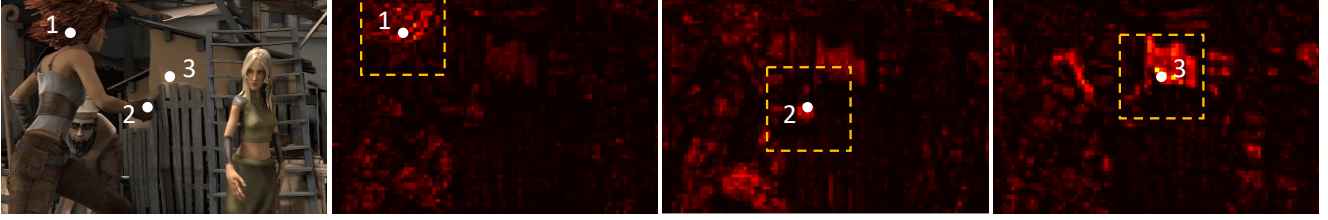


Figure 6. Visualization of attention scores. The more red a pixel is, the higher the score is.

boxes highlights the regions where our method outperforms RAFT.

8. How We Visualize Cost Volumes

In order to compare the 4D cost volumes C of RAFT [42] and our method, we extract the matrix $F_{x,y}$ as the matching matrix for the point (x, y) ,

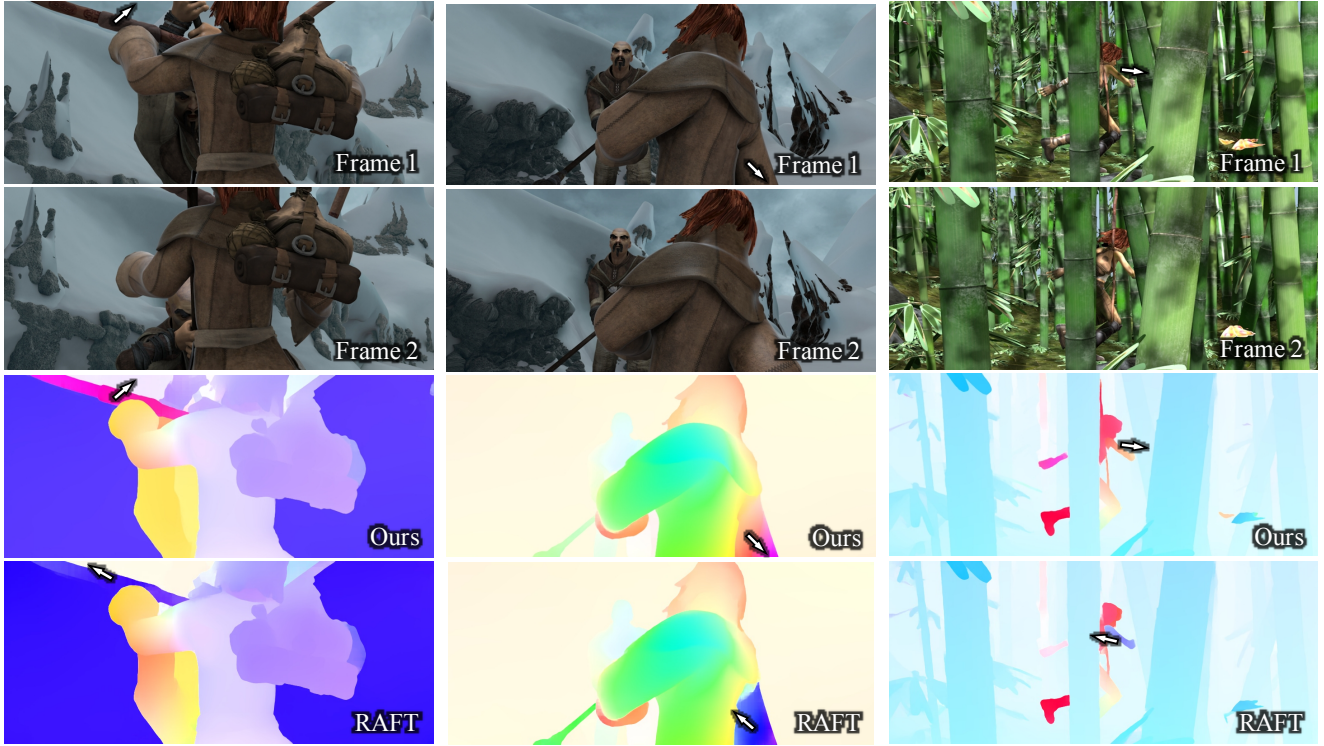
$$F_{x,y} = \text{softmax}(C[x, y, (x + \delta x - 40) : (x + \delta x + 40), (y + \delta y - 40) : (y + \delta y + 40)]) \quad (14)$$

where δx and δy are indicated by the ground truth flow at (x, y) . The symbol $C[\cdot]$ means to fetch values from C within a given range. Then, we average $F_{x,y}$ on all points within a specific displacement range for all images in Sintel and visualize the averaged matching matrix.

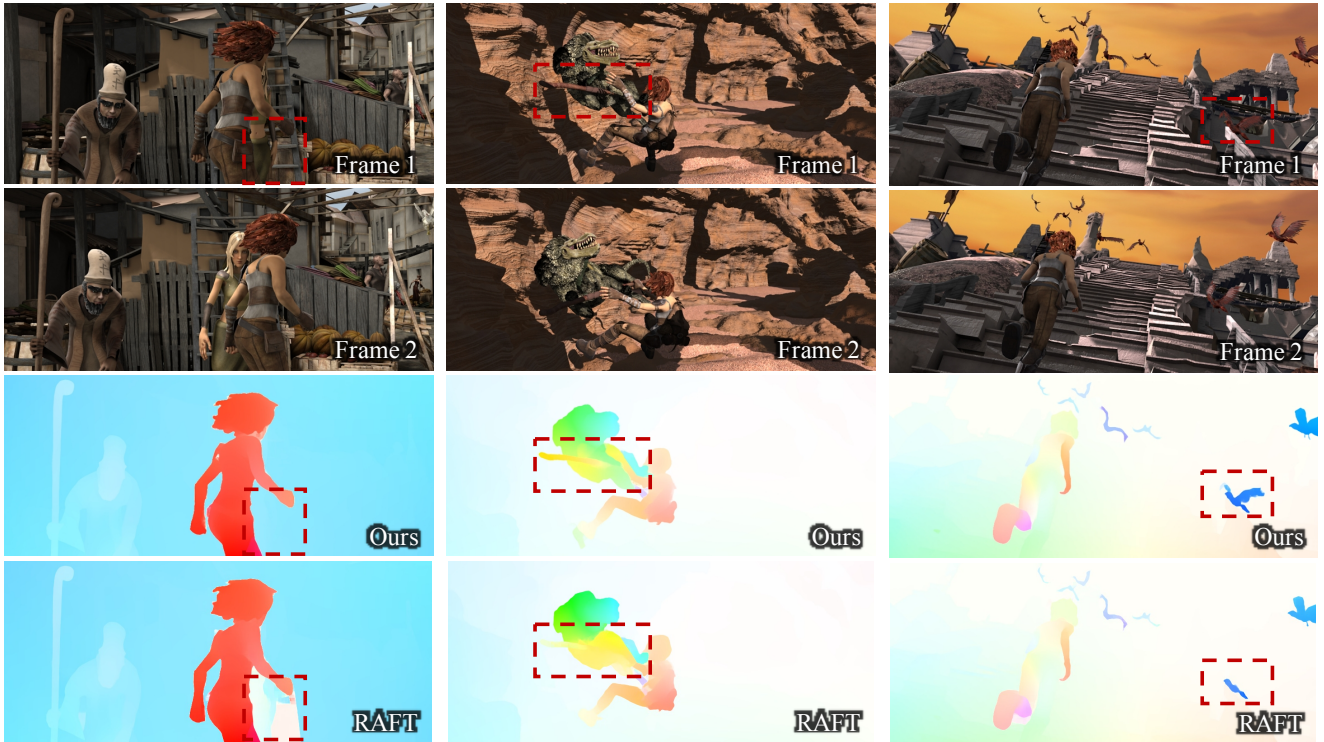
We visualize the cost volume for different ranges of displacements in Fig. 10. The larger the value at the center of the averaged matching matrix is, the higher quality the cost volume has. As shown, GMFlowNet outperforms RAFT in all displacement ranges, which indicates that our approach provides better cost volumes not only for small displacements but also for large ones.



Figure 7. **Visualizations of coarse flow.** For (a) Sintel, models are trained on C+T. For (b) KITTI, models are trained on C+T+S+K+H. Ground-truth flows for KITTI are unavailable and thus are not shown. With the coarse flow, our method outperforms the most popular optimization-only method RAFT [42]. Red dash boxes highlight the main differences between RAFT’s predictions and ours.

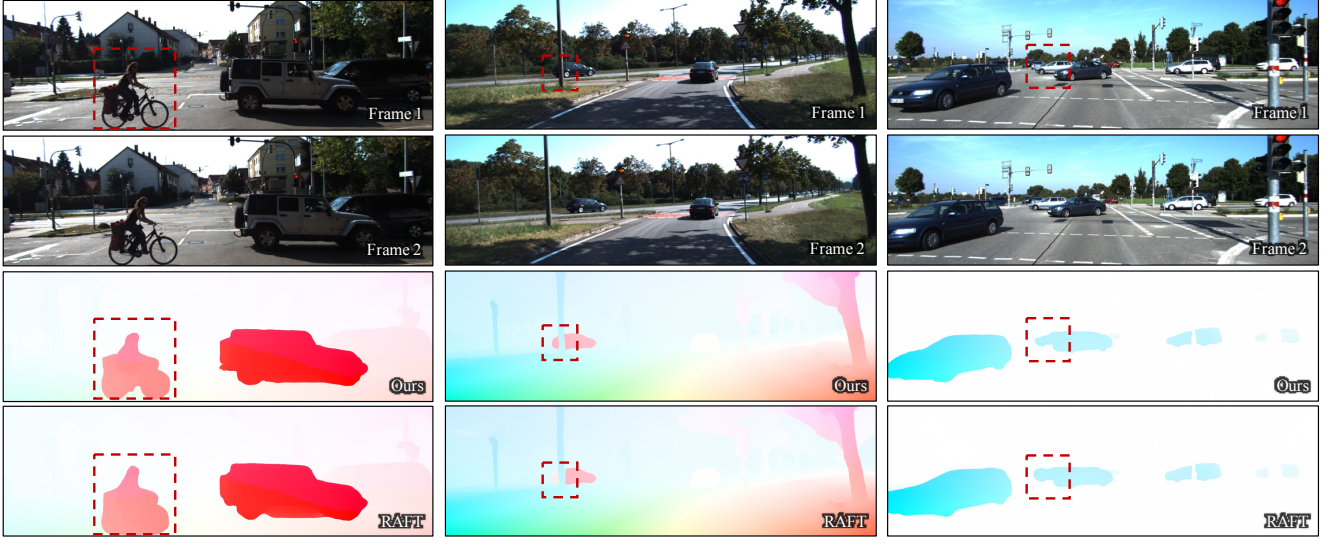


(a)



(b)

Figure 8. **Qualitative evaluation** on the Sintel test set [10]. White arrows in (a) and red dash boxes in (b) highlight the differences between our method and RAFT. Ground-truth optical flows are not available and are not shown. Models are trained on the same training data.



(a)

Figure 9. **Qualitative evaluation** on the KITTI test set [31]. Red dash boxes highlight the differences between our method and RAFT. Models are trained on the same training data.

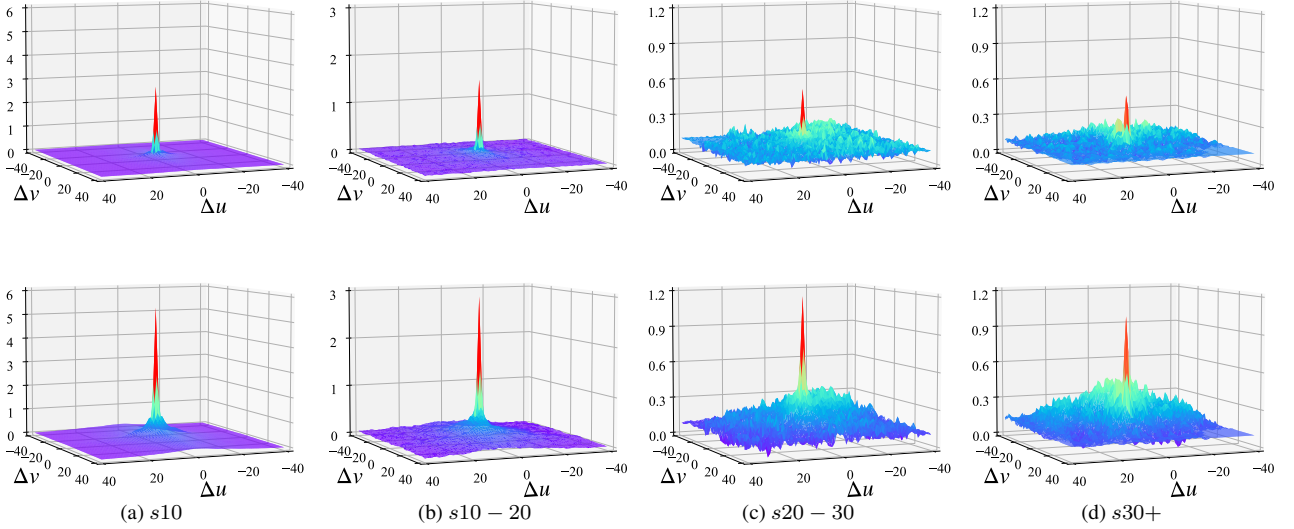


Figure 10. **Visualization of cost volumes in different range of displacements.** The first row is for RAFT [42], and the second row is ours. s_{10} refers to regions with displacements below 10 pixels, s_{10-20} for displacements between 10 and 20 pixels, s_{20-30} for displacements between 20 and 30 pixels, and s_{30+} for displacements larger than 30 pixels.