
SKFlow: Learning Optical Flow with Super Kernels

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Optical flow estimation is a classical yet challenging task in computer vision. One
2 of the essential factors to accurately predict optical flow is to alleviate occlusions
3 between frames. However, it is still a thorny problem for current top-performing op-
4 tical flow estimation methods due to insufficient local evidence to model occluded
5 areas. In this paper, we propose Super Kernel Flow Network (SKFlow), a CNN
6 architecture to ameliorate the impacts of occlusions on optical flow estimation.
7 SKFlow benefits from the super kernels which bring enlarged receptive fields to
8 complement the absent matching information and recover the occluded motions.
9 We present efficient super kernel designs by utilizing conical connections and hy-
10 brid depth-wise convolutions. Extensive experiments demonstrate the effectiveness
11 of SKFlow on multiple benchmarks, especially on the occluded areas. Without
12 pre-trained backbones on ImageNet and with modest increase in computation,
13 SKFlow achieves compelling performance and ranks **1st** among current published
14 methods on Sintel benchmark. On the challenging Sintel final pass test set, SKFlow
15 attains the average end-point error of 2.23, which surpasses the best published
16 result 2.47 by 9.72%.

17 1 Introduction

18 Optical flow is the task of modeling per-pixel motion across a pair of frames with various downstream
19 applications, e.g., pedestrian re-identification, video segmentation, and scene reconstruction, etc.
20 Nowadays, optical flow estimation still faces tough challenges such as occlusions and motion blurs,
21 etc. Among those challenges, occlusion is considered as one of the most difficult problems which is
22 still under-explored. Particularly, the term *occlusion* in optical flow estimation can be extended to a
23 broader definition [16]: *regions where pixels appear in the current frame while disappear in the next*
24 *frame*. Occlusions pose apparent difficulties in predicting optical flow since it directly violates the
25 brightness constancy constraint [10], where the intensities of pixels are regarded as the same between
26 consecutive frames. Therefore, the correspondence matching of occluded areas between frames could
27 be extremely hard.

28 One of the potential solutions is to utilize the neighbouring pixels to recover motions of the occluded
29 regions, such as learning the neighbouring relationship by CNNs [31, 32, 36], or interpolating the
30 hidden motions via smoothness terms in Markov Random Fields [4]. However, both convolution
31 and interpolation are restricted to the local information within the small operation window, which
32 only focus on learning the local evidence. As occlusions get worse, those local evidence would
33 be insufficient to recover the hidden motion and thus severely degrade the performance. Recent
34 works [16, 42] propose to model long-range dependencies between local descriptors via non-local
35 methods to make up the missing local evidence. These methods alleviate the issue to some extent, but
36 still tend to fail since the representative capabilities of local descriptors have been largely weakened
37 when facing severe occlusions.

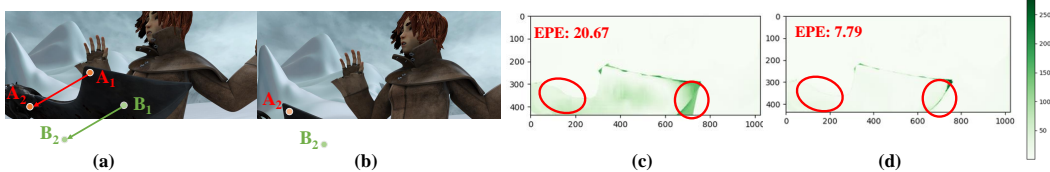


Figure 1: Occlusions between frames and flow error map of different methods. From (a) to (b), the most of the blade moves out-of-frame (e.g. pixel B_1 moves to pixel B_2). But the motion of the occluded B_1 could still be recovered via the non-occluded motion of A_1 over long distances due to the large receptive fields. (c) and (d) denote the error between the predicted flow and ground-truth from method with normal convolutions (e.g. GMA) and our SKFlow, respectively. The darker the color, the greater the error. SKFlow achieves lower end-point-error (EPE) and performs better in the occluded blade area.

38 Motivated by [21, 7] that larger kernels presents an effective way compared to deeper layers for
 39 larger receptive fields, using large kernels in optical flow estimation network is arguably a possible
 40 solution to handle the occlusion problem, as illustrated in Figure 1. However, directly applying large
 41 kernels is not practical due to: **(1)** The quadratic growth in computation. As the kernel size grows
 42 from 3×3 to 15×15 , the model would expand 25 times of its original size. **(2)** The optimization
 43 issue. As shown in [7, 24, 11], even with a careful design and training on a huge dataset, large kernel
 44 networks take effort to optimize and even prone to performance drop. It is more challenging on
 45 optical flow networks, where the training data is often much smaller. For example, Sintel [3] dataset
 46 has only 1,041 training samples.

47 In this work, we propose Super Kernel Flow Network (SKFlow), where we introduce a new architec-
 48 ture design which efficiently utilize the conical connections and hybrid depth-wise convolutions, and
 49 accordingly develop an effective optical flow network to handle the occlusions. Detailed description
 50 of the approach is described in Section 3. SKFlow is better at resolving the ambiguity caused by
 51 severe occlusions and obtains compelling performance on standard benchmarks (see Section 4).
 52 Besides, SKFlow attains a good trade-off between accuracy and computation cost. On the challenging
 53 Sintel final pass test set, SKFlow ranks **1st** among all published methods, improving the GMA by
 54 **9.72%**, and the increase in MACs is less than **8.42%**.

55 Our contributions are summarized as follows: **(1)** We introduce the super kernel schemes to optical
 56 flow task for the first time. **(2)** We explore three new architecture designs for super kernel designs
 57 in optical flow network and proposed a new network which we named SKFlow. **(3)** Our proposed
 58 SKFlow achieves the state-of-the-art performance on standard benchmarks and obtains strong cross-
 59 dataset generalization.

60 2 Related work

61 **Optical flow estimation.** Traditionally, optical flow is formulated as an energy minimization
 62 problem. Based on brightness constancy and spatial smoothness, Horn and Schunck [10] pioneered the
 63 variational approach to compute optical flow. On that basis, there emerged subsequent improvements
 64 for visual similarity designs and regularization terms [10, 1, 2, 30]. In the deep learning era, CNNs
 65 have emerged as a powerful technique for optical flow estimation since FlowNet [8]. Then the coarse-
 66 to-fine strategy is widely adopted [31, 32, 12, 13, 40, 14, 43]. Coarse-to-fine methods are influential
 67 but tend to miss small motions due to the inaccurate flow guidance in the coarse level. Therefore,
 68 RAFT [36] presents an iterative refinement method with the all-pairs field transform. It maintains a
 69 single fixed flow field at high resolutions and achieves significant improvements, inspiring lots of
 70 works such as Flow1D [39], FM-RAFT [17], Separable flow [42], AGFlow [20], and GMA [16], etc.

71 **Kernel sizes in convolution layers.** Recently, there have emerged some explorations of applying
 72 large kernels to CNNs, which have not been widely used since early designs [34, 35, 33, 27]. The
 73 ERF theory [21] indicates that larger kernels are more effective to obtain larger receptive fields than
 74 deeper layers. LRNet [11] proposes the local relation layer where 7×7 convolutions are widely
 75 used. However there is a performance drop when kernel size is increased to 9×9 . GCN [24]
 76 further enlarges the kernel size in the segmentation task by dense-connected symmetric separable

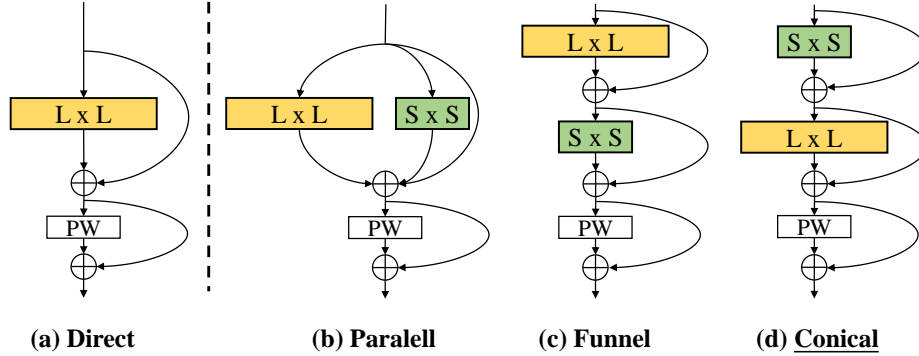


Figure 2: Architecture designs for super kernel blocks. L and S denote the kernel size of large and the auxiliary depth-wise convolutions respectively. PW is the point-wise convolution. The underlined design is used in our SKFlow model.

77 convolutions. It is observed that the performance tend to decrease when transferred to other tasks.
 78 CKCNN [26] and FlexConv [25] adopt larger kernel sizes, while model size and computation cost
 79 grow correspondingly. RepLKNet [7] develops the re-parameterization technique to help optimize
 80 large kernels and scales the kernel size to 31×31 , achieving the state-of-the-art performance on
 81 classification and downstream tasks such as detection and segmentation. Despite its competitive
 82 results on these high-level vision tasks, our experiments show that it could not be directly applied for
 83 dense prediction tasks such as optical flow estimation (see Section 4.4). Therefore, we explore the
 84 large kernel designs in optical flow network, aiming at alleviating the occlusions and improving the
 85 estimation performance.

86 3 Approach

87 We propose to utilize the large receptive field to handle the local ambiguity caused by occlusions.
 88 Since the motion of the occluded pixels is hard to estimate with local ambiguity, our SKFlow method
 89 introduces super kernels to help complement local evidence, and an efficient architecture design
 90 is proposed to reduce the computation burden. We describe the details of super kernel block in
 91 Section 3.1. The overall network architecture is then presented in Section 3.2. Finally, the learning
 92 process of SKFlow is demonstrated in Section 3.3.

93 3.1 Super kernel block designs

94 We enlarge the kernel size instead of deepening layers to attain large receptive field given the following
 95 considerations. **(1)** Deep layers lead to the optimization issue. Although ResNet [9] has resolved the
 96 optimization problems for deeper networks [7, 9] to some extent, recent works [6, 37] still suggest
 97 that it might behave like shallow networks and obtain the limited receptive fields with deeper layers.
 98 **(2)** The Effective Receptive Field (ERF) [21] grows sub-linearly with the number of layers and
 99 linearly with the kernel size, which demonstrates that the large kernel size is more effective.

100 **Super kernel block components.** In order to reduce the massive computation cost caused by large
 101 kernel sizes, we propose the super kernel block that contain three components: **(1)** Hybrid depth-wise
 102 convolution kernels. Inspired by separable convolutions [5], we first split a convolution into a couple
 103 of depth-wise convolutions. Namely, a large depth-wise kernel with size $L \times L$ and an auxiliary small
 104 depth-wise kernel with size $S \times S$. With the input of shape $N \times C_{in} \times H \times W$, the computation cost
 105 for this hybrid depth-wise convolutions is $N \times C_{in} \times H \times W \times (L^2 + S^2) \times 1$. The auxiliary kernel is
 106 designed to help capture small-scale patterns in the frame. **(2)** Residual connections. Residuals [9] are
 107 used (a) to combine the large kernel and the auxiliary kernel; (b) connect depth-wise and point-wise
 108 convolutions. **(3)** Point-wise convolution. An 1×1 point-wise convolution is applied after hybrid
 109 depth-wise convolutions to help the information flow across channels. The point-wise convolution
 110 does not change the input dimension and the computation cost is: $N \times C_{in} \times H \times W \times C_{in}$ with
 111 the input shape $N \times C_{in} \times H \times W$.

112 **Super kernel block architecture.** We propose three architecture designs for our super kernel
 113 blocks. As shown in Figure 2, our designs (*b*, *c*, and *d* columns) are derived from the Direct design in
 114 column *a*. Specifically, **(1) Parallel** (Figure 2 (b)), is a block of layers that adopting a parallel small
 115 kernel in the large depth-wise convolution layer. An point-wise convolution is followed by the parallel
 116 VGG-style convolutions combined with skip connections. **(2) Funnel** (Figure 2 (c)), is a ResNet-style
 117 block of layers where the kernel sizes are decreased from large kernel L to auxiliary kernel S . **(3)**
 118 **Conical** (Figure 2 (c)), is similar to the Funnel design but applies the hybrid convolutions in the
 119 opposite order. Experimental results in Section 4.4 that all three large kernel designs outperform
 120 normal convolution layers with normal small kernels and the computation increase is very limited.
 121 Among all these architectures, the conical block obtains the best performance and is adopted as our
 122 final design, which can be formulated as:

$$\mathbf{h} = \sigma(\mathbf{x} + Conv_{S \times S}^{dw}(\mathbf{x})) \quad (1)$$

$$\mathbf{d} = \sigma(\mathbf{h} + Conv_{L \times L}^{dw}(\mathbf{h})) \quad (2)$$

$$\mathbf{p} = Conv_{1 \times 1}^{pw}(\mathbf{d}) \quad (3)$$

$$\mathbf{o} = \mathbf{d} + \sigma(\mathbf{p}) \quad (4)$$

123 where \mathbf{x} and \mathbf{o} denote the input and output feature map, respectively. $Conv_{S \times S}^{dw}$, $Conv_{L \times L}^{dw}$ refer to
 124 the depth-wise convolution with large and small kernels, respectively. $Conv_{1 \times 1}^{pw}$ is the point-wise
 125 convolution and σ is the activation function.

126 For each block, there are a set of 1×1 convolution layers to match the dimension and increase more
 127 non-linear transforms. The computation cost is $(N \times H \times W)[D \times \alpha C_{in} \times (C_{in} + C_{out})]$, where
 128 D is the depth of convolution layer set and α is the channel scaling factor. The total computation cost
 129 of our super kernel block is $N \times H \times W \times C_{in} \times [C_{in} + L^2 + S^2 + \alpha D(C_{in} + C_{out})]$. We can
 130 compute the computation cost ratio of our super kernel block and the normal convolution as:

$$\begin{aligned} \frac{Cost_{ours}}{Cost_{normal}} &= \frac{N \times H \times W \times C_{in} \times [C_{in} + L^2 + S^2 + \alpha D(C_{in} + C_{out})]}{N \times H \times W \times C_{in} \times C_{out} \times L^2} \\ &= \frac{1}{L^2} \left[\left(1 + \frac{C_{in}}{C_{out}}\right) \alpha D + \frac{S^2 + C_{in}}{C_{out}} \right] + \frac{1}{C_{out}} \end{aligned} \quad (5)$$

131 where α , D , S , C_{in} and C_{out} are constants. Notably, the computation cost of our super kernel is
 132 reduced to $O(1/L^2)$ compared with normal large kernels. Besides, our super kernel also achieves
 133 compelling performance (see Section 4 for results).

134 3.2 Super kernel flow network

135 Our Super Kernel Flow (SKFlow) network follows the similar framework to GMA [16] except the
 136 super kernel modules, which are made up of super kernel blocks. We present the overall network
 137 architecture in the following.

138 **All-pairs correlation cost volume.** The all-pairs correlation cost volume proposed by [36] is used
 139 to model correlations for all possible displacements. Features from two frames perform a dot-product
 140 and then the matching correlations \mathbf{c} across different levels can be built, namely,

$$\mathbf{c}^l(i, j, m, n) = \frac{1}{2^{2l}} \sum_u^{2^l} \sum_v^{2^l} \langle \mathbf{x}_1(i, j), \mathbf{x}_2(2^l m + u, 2^l n + v) \rangle \quad (6)$$

141 where \mathbf{x}_1 and \mathbf{x}_2 are feature maps extracted from input frames, and l refers to the correlation level.
 142 $\langle \cdot, \cdot \rangle$ denotes the inner product function.

143 **Global motion aggregation (GMA) module.** GMA module adopts the self-attention mechanism
 144 to model long-range connections between local descriptors. Following GMA [16], we apply the
 145 module in the flow decoder. Given the input motion feature \mathbf{x} with height H and width W , the output
 146 \mathbf{o} is defined as:

$$\mathbf{o}(i, j) = \mathbf{x}(i, j) + \gamma \sum_u^H \sum_v^W f(\mathbf{q}_y(i, j), \mathbf{k}_y(u, v)) \mathbf{v}_x(u, v) \quad (7)$$

147 where \mathbf{q}_y and \mathbf{k}_y denote the query and key vector derived from context feature map. V_m is the value
 148 vector derived from \mathbf{x} , f is the dot-product attention function and γ is a learnable coefficient.

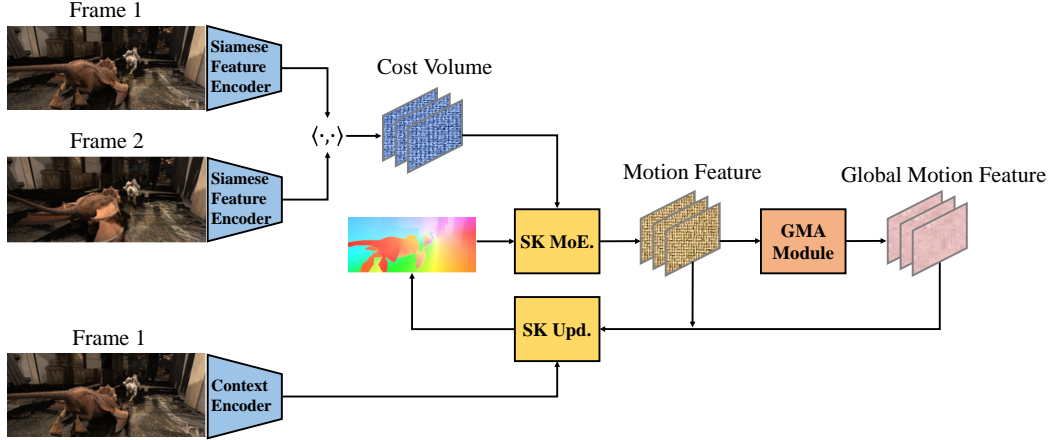


Figure 3: Overview of our proposed SKFlow. SK MoE. and SK Upd. denote our proposed super kernel motion encoder and updater respectively. GMA module is the non-local method proposed in [16].

149 **Super kernel modules.** Super kernel modules consist of **(1) Super kernel motion encoder**, which
 150 shares similar architectures with GMA but applies super kernel blocks. It derives motion features
 151 from cost volume vectors \mathbf{c} and the predicted flow \mathbf{f} , namely,

$$\mathbf{c}' = SKBlock(SKBlock(\mathbf{c})) \quad (8)$$

$$\mathbf{f}' = SKBlock(SKBlock(\mathbf{f})) \quad (9)$$

$$\mathbf{o} = Concat(SKBlock(Concat(\mathbf{c}', \mathbf{f}')), \mathbf{f}) \quad (10)$$

152 where $SKBlock$ denotes our proposed super kernel block, and $Concat$ is the concatenation
 153 operation. **(2) Super kernel updater.** Different from the ConvGRU used in RAFT and GMA, we
 154 directly adopt our super kernel block as the updater to refine the predicted residual flow $\Delta\mathbf{f}$, which
 155 could be formulated as:

$$\Delta\mathbf{f} = SKBlock(Concat(\mathbf{x}_m, \mathbf{x}_c, \mathbf{x}_g)) \quad (11)$$

156 where $\mathbf{x}_m, \mathbf{x}_c$ denote the motion feature and context feature. \mathbf{x}_g refers to the global motion feature
 157 from GMA module.

158 The whole architecture of our SKFlow is shown in Figure 3. It is notable that the super kernel blocks
 159 are applied in the decoder only based on the following consideration: Applying in the decoder is
 160 more efficient due to its relatively small size. Given the lack of efficient implementation like normal
 161 3×3 kernel for large depth-wise kernels in most libraries, super kernels in those two large encoders
 162 can take lots of additional time in inference and training. Although it may be solved when faster
 163 implementations occur in the future, currently applying super kernels in the smaller decoder part
 164 brings modest time increase and achieves significant performance improvements.

165 3.3 Supervision

166 Following previous works [36, 16], we adopt the following loss function to supervise the parameter
 167 update. The l_1 loss of the predicted flow after every refinement are weighted with the exponentially
 168 increasing coefficients:

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

169 where $\|\cdot\|_1$ denotes the l_1 distance between flow ground-truth f_{gt} and output flow, and f_i represents
 170 the predicted flow field at the i th refinement. λ denotes weights on different predicted flows and is
 171 set to 0.8 in our experiments.

Table 1: Quantitative results on Sintel and KITTI. We report average End-Point Error(EPE) unless specified. C+T refers to the FlyingChairs \rightarrow FlyingThings schedule. + S + K + H denotes that Sintel, KITTI and HD1K training sets are combined when finetuning on Sintel. * denotes using warm-start strategy in RAFT [36]. SKFlow achieves the best generalization performance on C+T and C+T+S+K+H.

Training Data	Method	Sintel (train)		KITTI-15 (train)		Sintel(test)		KITTI-15 (test)
		Clean	Final	Fl-epe	Fl-all	Clean	Final	Fl-all
C+T	HD3 [41]	3.84	8.77	13.17	24.0	-	-	-
	PWC-Net [31]	2.55	3.93	10.35	33.7	-	-	-
	VCN [40]	2.21	3.68	8.36	25.1	-	-	-
	MaskFlowNet [43]	2.25	3.61	-	23.1	-	-	-
	FlowNet2 [15]	2.02	3.54	10.08	30.0	3.96	6.02	-
	DICL-Flow [38]	1.94	3.77	8.70	23.6	-	-	-
	RAFT [36]	1.43	2.71	5.04	17.4	-	-	-
	AGFlow [20]	1.31	2.69	4.82	17.0	-	-	-
	Separable Flow [42]	1.30	2.59	4.60	15.9	-	-	-
	GMA [16]	1.30	2.74	4.69	17.1	-	-	-
	SKFlow (Ours)	1.22	2.46	4.27	15.5	-	-	-
C+T+S+K+H	LiteFlowNet2 [13]	(1.30)	(1.62)	(1.47)	(4.8)	3.48	4.69	7.74
	PWC-Net+ [32]	(1.71)	(2.34)	(1.50)	(5.3)	3.45	4.60	7.72
	MaskFlowNet [43]	-	-	-	-	2.52	4.17	6.10
	RAFT [36]	(0.76)	(1.22)	(0.63)	(1.5)	1.61*	2.86*	5.10
	AGFlow [20]	(0.65)	(1.07)	(0.58)	(1.2)	1.43	2.47	4.89
	Separable Flow [42]	(0.69)	(1.10)	(0.69)	(1.6)	1.50	2.67	4.64
	GMA [16]	(0.62)	(1.06)	(0.57)	(1.2)	1.39*	2.47*	5.15
		SKFlow (Ours)	(0.52)	(0.78)	(0.51)	(0.94)	1.28*	2.23*

172 4 Experiments

173 **Experimental setup.** Our model is evaluated on Sintel [3] and KITTI [22] datasets. We adopt the
174 average End-Point-Error (EPE) as the evaluation metric, which denotes the mean flow error over
175 all pixels. KITTI also adopts the Fl-All (%) metric, which computes the percentage of pixels with
176 EPE larger than 3 pixels or over 5% of ground truth. Sintel consists of different passes rendered with
177 different levels of difficulty. Specifically, the Albedo pass is rendered without illumination effects
178 and has approximately piecewise constant colors. The clean pass introduces reflection properties
179 and various illumination. The final pass adds motion blur, atmospheric effects and other artistic
180 embellishments to lighting, which is the most challenging. In our experiments, the clean and final
181 passes are used for training and the albedo pass is for testing on occluded areas.

182 **Implementation details.** Following previous works [36, 16, 42, 20], we first pre-train our SKFlow
183 using the FlyingChairs \rightarrow FlyingThings schedule. On FlyingChairs, our model is trained for 120K
184 iterations with a batch size of 8 and the learning rate of 2.5×10^{-4} . On FlyingThings, we train the
185 model for 150K iterations with batch size 6 and the learning rate of 1.75×10^{-4} . We then fine-tune
186 for Sintel with the combined dataset from Sintel, FlyingThings, KITTI and HD1K [18] for 180K
187 iterations. The batch size is 6 and learning rate is 1.75×10^{-4} . Our model is further fine-tuned on
188 KITTI for another 50K iterations. We use the AdamW [19] optimizer and one-cycle policy [28]. Our
189 SKFlow is built with PyTorch [23] library and trained using two Tesla V100 GPUs.

190 4.1 Quantitative results

191 Following previous works [36, 16, 42, 20], we evaluate our SKFlow on Sintel and KITTI-2015
192 dataset. We first evaluate the generalization ability of models by pretraining on FlyingChairs and
193 FlyingThings datasets and evaluate on Sintel and KITTI datasets. Next, we compare the performance
194 on each dataset. Models are finetuned on the training set combined from Sintel, HD1K and KITTI
195 dataset, and evaluate on the Sintel and KITTI dataset. The experimental results are shown in Table 1.

196 **Results for cross-dataset generalization.** The cross-dataset results can be seen in Tables 1 "C+T"
197 rows. From the results we can see the generalization capability of different methods. Specifically,
198 on the challenging Sintel final pass, our SKFlow obtains the best performance among all published

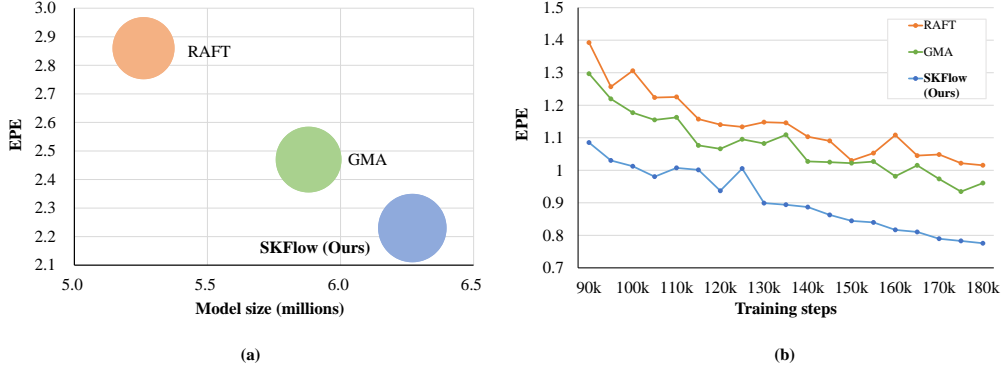


Figure 4: (a) Comparisons between MACs, parameters and performance. A larger bubble represents higher MACs. (b) EPE curve on Sintel during training. All models are trained under the same settings.

199 methods, outperforming GMA by 10.21%. On KITTI and clean pass, our method achieves the Fl-epe
 200 of 4.27 and Fl-all of 15.5, and the EPE of 1.22, outperforming GMA by 8.96%, 9.36% and 6.15%,
 201 respectively.

202 **Results on Sintel benchmark.** The experimental results are shown in Table 1 "C+T+S+K+H" rows.
 203 From the results we can see that our SKFlow achieves significant improvements in performance on
 204 both Sintel training and test set among listed methods. Compared with GMA, our SKFlow achieves
 205 the epe reduction of 7.91% and 9.72% on the clean and final pass of Sintel test set. On the training
 206 clean and final pass, SKFlow outperforms GMA by 16.13% and 26.42%.

207 **Results on KITTI benchmark.** The results are shown in Table 1 "C+T+S+K+H" rows. Our
 208 SKFlow achieves 4.84, achieving competitive results. On KITTI training set, our method achieves
 209 with the improvement on GMA with the EPE of 0.51% and of Fl-epe of 0.94%. On KITTI test set,
 210 our method outperforms GMA by 6.02%.

211 4.2 Model efficiency

212 We then compare the efficiency of our SKFlow with current state-of-the-art methods. We use
 213 MACs and model size to measure the model efficiency. Specifically, MACs are measured using
 214 PTFLOPS [29] library. Model size is measured in number of parameters. All the models are fine-
 215 tuned using the combined dataset from FlyingThings3D, HD1K, KITTI and Sintel training sets. The
 216 evaluation is conducted on Sintel test set. As illustrated in Figure 4 (a), our proposed SKFlow shows
 217 better accuracy (9.72%) compared with RAFT and GMA models with similar MACs (8.42%) and
 218 limited increasing in model size (6.63%).

219 4.3 Occlusion analysis

220 We split pixels into occluded (Occ) and non-
 221 occluded (Noc) with the help of the occlusion
 222 maps provided in Sintel dataset. Following the
 223 previous work [16], the model is pre-trained us-
 224 ing the C+T schedule and then fine-tuned using
 225 the C+T+S+H+K schedule for 150K iterations.
 226 The Albedo pass of Sintel dataset is used to
 227 verify the effectiveness of estimating occluded
 228 motions and not used in training. We choose
 229 Albedo pass as the test set because it is rendered
 230 without illumination effects and thus adheres
 231 to the brightness constancy constraint except
 232 the occluded areas, which highlights the perfor-
 233 mance on occlusions. The results are shown in

Table 2: Performance on occluded and non-occluded areas on the Sintel dataset.

Data	Type	RAFT	GMA	Ours
Clean (train)	Noc	0.29	0.27	0.26
	Occ	5.00	4.44	3.96
	All	0.63	0.58	0.53
Final (train)	Noc	0.60	0.52	0.46
	Occ	6.90	6.32	5.39
	All	1.06	0.95	0.82
Albedo (test)	Noc	0.31	0.31	0.27
	Occ	6.18	5.96	4.75
	All	0.74	0.72	0.60

234 table 2, where "Noc" denotes non-occluded pixels, "Occ" denotes the occluded pixels and "All"
235 denotes the overall pixels. We can see that our SKFlow obtains significant improvements on both
236 occluded and non-occluded areas, compared with RAFT and GMA models. This result demonstrates
237 that our SKFlow method can better resolve the ambiguity caused by occlusion areas.

238 4.4 Ablation study

239 We perform a set of ablation experiments to show the importance of each component in the network
240 design. The results are shown in Figure 4 (b) and Table 3 respectively. All ablated versions are
241 trained using the C+T schedule. The settings which are used in our final model is underlined. In the
242 following, we describe each of the experiments in more detail.

243 **Training schedules.** In particular, we find that the performance of our proposed SKFlow keeps
244 improving when training steps is increasing. Specifically, we train the model by additional 30K
245 steps on FlyingThings dataset and additional 60K steps for fine-tuning on Sintel dataset. For a fair
246 comparison, we then train GMA and RAFT using the same settings. As illustrated in Figure 4
247 (b), the EPE of SKFlow keeps decreasing while the the training steps is increased, and an obvious
248 improvement of performance can be obtained comparing to the RAFT and GMA. We argue that large
249 receptive fields are inherently helpful to capture more motion information and raise the upper bound
250 of the learning ability of networks, although it tends to be hard to optimize.

251 **Kernel sizes.** We explore the impact of different kernel sizes. Specifically, the kernel size is
252 increased from 1 to 31, growing $2\times$ each time. As shown in Table 3, we can see that kernel
253 size 15×15 achieves the best overall performance. Although it obtains 0.89% and 0.92% lower
254 performance than 7×7 kernel on KITTI, it achieves 6.98% and 5.56% higher in performance on
255 Sintel clean and final pass. We also observe much lower performance when the kernel size is 31,
256 especially on a small dataset such as KITTI, which indicates the optimization of the extremely large
257 kernel is still a challenge given the limited number of samples in optical flow datasets.

258 **Super kernel block components.** We explore the effect of each component in Section 3.1). Specif-
259 ically, "Res" denotes the residual connections and "Aux" denotes the Auxiliary small kernel. From
260 Table 3 "Components" we can see that by using both the residual connections and the auxiliary small
261 kernel, the overall performance improves significantly. (comparable result 2.46 vs 2.45 on Sintel
262 Final). Nevertheless, without the auxiliary kernel, there will be a significant drop for small dataset
263 like KITTI, indicating the importance of the auxiliary kernel on optimization.

264 **Super kernel block architecture.** We explore the performance of three hybrid architecture designs
265 in Section 3.1, as shown in Table 3 "Architecture". We can see that the conical connection design
266 achieves the best performance, showing a strong capability to optimize optical flow networks with
267 large kernels. It is different from the existing work [7] showing that the parallel re-parameterization
268 structure perform well in high-level tasks such as image classification. On optical flow estimation
269 our conical structure obtains better performance. There might be two reasons: (1) Differences in the
270 distribution and size of datasets. Compared with large datasets such as ImageNet where samples are
271 from real life, most optical flow datasets are synthetic and contains much smaller samples, which
272 makes it harder for large kernels to optimize. (2) Different focus of tasks. Compared with optical
273 flow tasks, high-level vision tasks pay more attention to semantic information instead of per-pixel
274 correspondence between two images.

275 **Auxiliary kernel size.** We also explore different sizes for the auxiliary kernel. According to the
276 results shown in Table 3 "Aux kernel size", the most effective auxiliary kernel size is 1×1 . Although
277 3×3 achieves a slightly better performance on Sintel, 1×1 kernel outperforms 3×3 kernel on
278 Fl-epe and Fl-all of KITTI dataset with more improvement and less computation.

279 **Super kernel modules.** We compare the performance of using super kernels in update block and
280 motion encoder and the results are shown in Table 3 "SK modules". Upd. refers to the update block
281 and MoE. denotes the motion encoder. We can see that by using super kernels in the update block the
282 performance is 2.65 on Sintel final pass, while by using the super kernels in the motion encoder, the
283 performance is 2.56. When super kernels are applied on both motion encoder and update block in
284 SKFlow network, the performance on both Sintel and KITTI datasets is futher improved significantly.

Table 3: Ablations on super kernel block designs. Models are trained using the C+T schedule.

Experiment	Method	Sintel(train)		KITTI-15(train)		Parameters
		Clean	Final	Fl-epe	Fl-all	
Large kernel size	31	1.27	2.65	4.73	17.40	7.08M
	<u>15</u>	1.20	2.55	4.51	16.26	6.27M
	7	1.29	2.70	4.47	16.11	6.08M
	3	1.30	2.77	4.83	17.31	6.04M
	1	1.41	2.83	7.12	21.58	6.03M
Components	No Res	1.35	2.56	4.76	17.33	6.27M
	No Aux	1.25	2.45	4.53	16.55	6.27M
	No Res & Aux	1.33	2.61	5.04	17.38	6.27M
	<u>With Res & Aux</u>	1.22	2.46	4.27	15.47	6.27M
Architecture	Parallel	1.24	2.53	4.71	17.82	6.27M
	Funnel	1.19	2.54	4.55	16.97	6.27M
	<u>Conical</u>	1.22	2.46	4.27	15.47	6.27M
Aux kernel size	3	1.21	2.44	4.33	15.56	6.29M
	<u>1</u>	1.22	2.46	4.27	15.47	6.27M
SK modules	Upd.	1.37	2.65	4.89	16.89	5.30M
	MoE.	1.28	2.56	4.60	16.86	6.67M
	<u>Upd.+MoE.</u>	1.22	2.46	4.27	15.47	6.27M

285 4.5 Visualizations

286 As shown in Figure 5, we visualize the predicted optical flows on sample image from Sintel test set.
 287 Results are highlighted using red circle in the figure, which demonstrates that our proposed SKFlow
 288 predicts a more accurate flow and captures more details compared with other methods.



Figure 5: Visualizations of predicted optical flow among different methods on Sintel test set.

289 5 Conclusions

290 In this work, we propose SKFlow, which utilizes large receptive fields brought by super kernels to
 291 recover the occluded motion in optical flow estimation. To the best of our knowledge, SKFlow is the
 292 first to apply the kernel size up to 15×15 to optical flow networks with significant performance gain
 293 and modest computation increase. To leverage the computation cost for the large kernels and attains
 294 improved accuracy, SKFlow explores and provides effective and efficient design schemes for applying
 295 super kernels, such as conical connections and hybrid depth-wise convolutions. Comprehensive
 296 experiments have demonstrated the effectiveness of SKFlow. With less than 8.42% increase in MACs,
 297 SKFlow outperform previous state-of-the-art method by 9.72% on Sintel final pass, ranking 1st
 298 among all published methods.

299 **References**

- 300 [1] Michael J Black and Padmanabhan Anandan. A framework for the robust estimation of optical
301 flow. In *1993 (4th) International Conference on Computer Vision*, pages 231–236. IEEE, 1993.
- 302 [2] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Lucas/kanade meets horn/schunck:
303 Combining local and global optic flow methods. *International journal of computer vision*,
304 61(3):211–231, 2005.
- 305 [3] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source
306 movie for optical flow evaluation. In *European conference on computer vision*, pages 611–625.
307 Springer, 2012.
- 308 [4] Qifeng Chen and Vladlen Koltun. Full flow: Optical flow estimation by global optimization
309 over regular grids. In *Proceedings of the IEEE conference on computer vision and pattern
310 recognition*, pages 4706–4714, 2016.
- 311 [5] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceed-*
312 *ings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258,
313 2017.
- 314 [6] Soham De and Sam Smith. Batch normalization biases residual blocks towards the identity
315 function in deep networks. *Advances in Neural Information Processing Systems*, 33:19964–
316 19975, 2020.
- 317 [7] Xiaohan Ding, Xiangyu Zhang, Yizhuang Zhou, Jungong Han, Guiguang Ding, and Jian Sun.
318 Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. *arXiv preprint
319 arXiv:2203.06717*, 2022.
- 320 [8] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov,
321 Patrick van der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with
322 convolutional networks. In *Proceedings of the IEEE International Conference on Computer
323 Vision (ICCV)*, December 2015.
- 324 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
325 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
326 pages 770–778, 2016.
- 327 [10] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*,
328 17(1-3):185–203, 1981.
- 329 [11] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image
330 recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
331 pages 3464–3473, 2019.
- 332 [12] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Liteflownet: A lightweight convolutional
333 neural network for optical flow estimation. In *Proceedings of the IEEE conference on computer
334 vision and pattern recognition*, pages 8981–8989, 2018.
- 335 [13] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. A lightweight optical flow cnn—revisiting
336 data fidelity and regularization. *IEEE transactions on pattern analysis and machine intelligence*,
337 43(8):2555–2569, 2020.
- 338 [14] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion
339 estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
340 Recognition*, pages 5754–5763, 2019.
- 341 [15] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas
342 Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of
343 the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- 344 [16] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate
345 hidden motions with global motion aggregation. In *Proceedings of the IEEE/CVF International
346 Conference on Computer Vision*, pages 9772–9781, 2021.

- 347 [17] Shihao Jiang, Yao Lu, Hongdong Li, and Richard Hartley. Learning optical flow from a
348 few matches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
349 Recognition*, pages 16592–16600, 2021.
- 350 [18] Daniel Kondermann, Rahul Nair, Katrin Honauer, Karsten Krispin, Jonas Andrusis, Alexander
351 Brock, Burkhard Gusefeld, Mohsen Rahimimoghaddam, Sabine Hofmann, Claus Brenner, et al.
352 The hci benchmark suite: Stereo and flow ground truth with uncertainties for urban autonomous
353 driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition
354 Workshops*, pages 19–28, 2016.
- 355 [19] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International
356 Conference on Learning Representations*, 2018.
- 357 [20] Ao Luo, Fan Yang, Kunming Luo, Xin Li, Haoqiang Fan, and Shuaicheng Liu. Learning optical
358 flow with adaptive graph reasoning. *arXiv preprint arXiv:2202.03857*, 2022.
- 359 [21] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive
360 field in deep convolutional neural networks. *Advances in neural information processing systems*,
361 29, 2016.
- 362 [22] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene
363 flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015.
- 364 [23] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
365 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
366 pytorch. 2017.
- 367 [24] Chao Peng, Xiangyu Zhang, Gang Yu, Guiming Luo, and Jian Sun. Large kernel matters—
368 improve semantic segmentation by global convolutional network. In *Proceedings of the IEEE
369 conference on computer vision and pattern recognition*, pages 4353–4361, 2017.
- 370 [25] David W Romero, Robert-Jan Bruintjes, Jakub M Tomczak, Erik J Bekkers, Mark Hoogendoorn,
371 and Jan C van Gemert. Flexconv: Continuous kernel convolutions with differentiable kernel
372 sizes. *arXiv preprint arXiv:2110.08059*, 2021.
- 373 [26] David W Romero, Anna Kuzina, Erik J Bekkers, Jakub M Tomczak, and Mark Hoogendoorn.
374 Ckconv: Continuous kernel convolution for sequential data. *arXiv preprint arXiv:2102.02611*,
375 2021.
- 376 [27] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
377 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 378 [28] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks
379 using large learning rates. In *Artificial intelligence and machine learning for multi-domain
380 operations applications*, volume 11006, page 1100612. International Society for Optics and
381 Photonics, 2019.
- 382 [29] Vladislav Sovrasov. Flops counter for convolutional networks in pytorch framework, 2019.
- 383 [30] Deqing Sun, Stefan Roth, and Michael J Black. A quantitative analysis of current practices
384 in optical flow estimation and the principles behind them. *International Journal of Computer
385 Vision*, 106(2):115–137, 2014.
- 386 [31] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow
387 using pyramid, warping, and cost volume. In *Proceedings of the IEEE conference on computer
388 vision and pattern recognition*, pages 8934–8943, 2018.
- 389 [32] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Models matter, so does training:
390 An empirical study of cnns for optical flow estimation. *IEEE transactions on pattern analysis
391 and machine intelligence*, 42(6):1408–1423, 2019.
- 392 [33] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4,
393 inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI
394 conference on artificial intelligence*, 2017.

- 395 [34] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov,
396 Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions.
397 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9,
398 2015.
- 399 [35] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Re-
400 thinking the inception architecture for computer vision. In *Proceedings of the IEEE conference*
401 *on computer vision and pattern recognition*, pages 2818–2826, 2016.
- 402 [36] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In
403 *European conference on computer vision*, pages 402–419. Springer, 2020.
- 404 [37] Andreas Veit, Michael J Wilber, and Serge Belongie. Residual networks behave like ensembles
405 of relatively shallow networks. *Advances in neural information processing systems*, 29, 2016.
- 406 [38] Jianyuan Wang, Yiran Zhong, Yuchao Dai, Kaihao Zhang, Pan Ji, and Hongdong Li.
407 Displacement-invariant matching cost learning for accurate optical flow estimation. *Advances*
408 *in Neural Information Processing Systems*, 33:15220–15231, 2020.
- 409 [39] Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-resolution opti-
410 cal flow from 1d attention and correlation. In *Proceedings of the IEEE/CVF International*
411 *Conference on Computer Vision*, pages 10498–10507, 2021.
- 412 [40] Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow.
413 *Advances in neural information processing systems*, 32, 2019.
- 414 [41] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition
415 for match density estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
416 *and Pattern Recognition*, pages 6044–6053, 2019.
- 417 [42] Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable flow:
418 Learning motion cost volumes for optical flow estimation. In *Proceedings of the IEEE/CVF*
419 *International Conference on Computer Vision*, pages 10807–10817, 2021.
- 420 [43] Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. Maskflownet: Asymmetric
421 feature matching with learnable occlusion mask. In *Proceedings of the IEEE/CVF Conference*
422 *on Computer Vision and Pattern Recognition*, pages 6278–6287, 2020.

423 Checklist

- 424 1. For all authors...
- 425 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
426 contributions and scope? [Yes]
- 427 (b) Did you describe the limitations of your work? [No]
- 428 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 429 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
430 them? [Yes]
- 431 2. If you are including theoretical results...
- 432 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 433 (b) Did you include complete proofs of all theoretical results? [N/A]
- 434 3. If you ran experiments...
- 435 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
436 mental results (either in the supplemental material or as a URL)? [Yes] See supplemen-
437 tal materials.
- 438 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
439 were chosen)? [Yes] See Section 4
- 440 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
441 ments multiple times)? [No]

- 442 (d) Did you include the total amount of compute and the type of resources used (e.g., type
443 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4
- 444 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 445 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 446 (b) Did you mention the license of the assets? [No]
- 447 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 448 (d) Did you discuss whether and how consent was obtained from people whose data you're
449 using/curating? [No]
- 450 (e) Did you discuss whether the data you are using/curating contains personally identifiable
451 information or offensive content? [No]
- 452 5. If you used crowdsourcing or conducted research with human subjects...
- 453 (a) Did you include the full text of instructions given to participants and screenshots, if
454 applicable? [N/A]
- 455 (b) Did you describe any potential participant risks, with links to Institutional Review
456 Board (IRB) approvals, if applicable? [N/A]
- 457 (c) Did you include the estimated hourly wage paid to participants and the total amount
458 spent on participant compensation? [N/A]