

InterpoNet, A brain inspired neural network for optical flow dense interpolation

Shay Zweig¹ and Lior Wolf^{2,3}

¹The Gonda Multidisciplinary Brain Research Center, Bar Ilan University, Israel

²The Blavatnik School of Computer Science, Tel Aviv University, Israel

³Facebook AI Research

shayzweig@gmail.com, wolf@cs.tau.ac.il

Abstract

Sparse-to-dense interpolation for optical flow is a fundamental phase in the pipeline of most of the leading optical flow estimation algorithms. The current state-of-the-art method for interpolation, EpicFlow, is a local average method based on an edge aware geodesic distance. We propose a new data-driven sparse-to-dense interpolation algorithm based on a fully convolutional network. We draw inspiration from the filling-in process in the visual cortex and introduce lateral dependencies between neurons and multi-layer supervision into our learning process. We also show the importance of the image contour to the learning process. Our method is robust and outperforms EpicFlow on competitive optical flow benchmarks with several underlying matching algorithms. This leads to state-of-the-art performance on the Sintel and KITTI 2012 benchmarks.

1. Introduction

The leading optical flow algorithms to date, with few exceptions, are not end-to-end deep learning. While some of them employ deep matching scores for estimating the best match in image I' for every location in image I , almost all methods employ multiple steps that do not involve learning. With the current affinity toward end-to-end deep learning solutions, the existence of large training datasets, and many concurrent contributions in the field of deep optical flow and related fields, one may wonder why this is the case.

Out of the four steps of modern optical flow pipelines: matching, filtering, interpolation and variational refinement, we focus on the third. In this step, a sparse list of matches is transformed into dense optical flow maps. It is one of the most crucial steps and without the availability of the EpicFlow method [33], which currently dominates this step, a large number of sparse matching techniques would not have been competitive enough to gain attention.

EpicFlow is an extremely effective method that is based on solid computer vision foundations. However, despite using sophisticated heuristics for improved runtime, it is still

rather slow and as a non-learning method, it is bounded in the performance it can deliver. Replacing EpicFlow by a deep learning method is harder than it initially seems. Feedforward neural networks excel in analyzing image information, but neuroscience tells us that in biological networks, lateral and top-down feedback loops are involved in solving cases where the information is missing or corrupted at random locations.

Artificial feedback networks are slower than feedforward networks, harder to train, and have not proven themselves in the practice of computer vision. We note that feedback networks with a predefined number of feedback iterations can be unrolled into deep feedforward networks with one major caveat – while in most feedforward networks, the supervision flows from the top down, in feedback networks, the supervision occurs at each iteration. To resolve this, we equip our network with supervision at every layer.

Inspired by neuroscience, we also suggest a loss involving lateral dependencies. Here, too, we replace the process of lateral feedback during run-time with additional supervision during training. In this way, the feedforward network learns how to mimic a network with lateral feedback loops by utilizing the training labels.

Taken together, our contributions are: (a) We propose, for the first time, to the best of our knowledge, a neural network based sparse-to-dense interpolation for optical flow. Our network performs better than the current state of the art, it is robust and can be adjusted to different matching algorithms and serve as the new default interpolation method in optical flow pipelines. (b) We introduce a new lateral dependency loss, embedding the correlations between neighbors into the learning process. (c) We define a novel architecture involving detour networks in each layer of the network. The new architecture provides a substantial increase in performance. (d) We solidify the importance of motion boundaries in learning dense interpolation for optical flow.

2. Related Work

Interpolation In The Visual Cortex. The visual system often receives a noisy and missing input. However, it is

known to robustly denoise and fill-in the gaps in the input image. This phenomenon termed - perceptual filling-in [20], was reported to occur for occlusions [17], illusory contours and surfaces [29], in the "blind spot"[31] and in visual scotomas [32]. Different features in the visual stimulus are filled in, including brightness[28], color[10], texture and motion[32]. The neurophysiological mechanism underlying perceptual filling-in is still under debate. However, many have found evidence of the existence of a neuronal filling-in mechanism [28, 30, 39, 15, 43]. In this mechanism, neurons that are retinotopically mapped to visible or salient parts of an image (such as the edges) are activated first. This initial activation is followed by a later spread to neurons that are mapped to the missing parts, resulting in a complete representation of the image [7, 42, 21]. This activation spread is mediated by both lateral connections within areas in the cortex as well as top down connections [15, 30, 43]. It was also shown to be very sensitive to edges in the image, usually originating in edges and stops when encountered with edges [38, 43]. Finally, neuronal filling-in was found to take place in multiple areas in the visual cortex hierarchy, from V1 and V2 [15, 34] via V4 [30, 35] and in higher areas [25, 24].

We designed our interpolation network to incorporate three concepts inspired by neuronal filling-in: the interactions between neighbor neurons, multi-layer supervision and the importance of edges. Neighbor neurons' interactions can be modeled by recurrent connections within a layer, such as the model suggested by Liang and Hu [23]. While the anatomic resemblance of such models to the cortex is appealing, in reality, they are unfolded to a feedforward network with shared weights. We, therefore, preferred to utilize the loss to force the interaction between neighbor neurons while using simpler, strictly feedforward networks, which were shown to perform extremely well for vision tasks while excelling in training time and simplicity.

Interpolation For Optical Flow. Most current optical flow approaches are based on a four phase pipeline. The first phase matches pixels between the images in the image pair, based on nearest neighbor fields or feature matching techniques (hand engineered or learned) [4, 14, 27]. The second phase filters matches with low confidence, producing a noisy and missing flow map[2]. The missing pixels usually undergo large displacements, a significant shift in appearance or are occluded in one of the images. Therefore, a third phase is needed to interpolate the missing parts and reduce the noise. A fourth and final phase applies refinement to the interpolated dense map from Phase 3.

The best and most used algorithm for optical flow interpolation (the third phase) is currently EpicFlow [33]. EpicFlow computes the flow of each pixel using a weighted sum of the pixel's local environment. Locality is defined by

a geodesic distance function based on the image edges that correspond to the motion boundaries. This edges aware approach yields good interpolation results for occluded pixels and large displacement. EpicFlow excels in interpolation. However, it is less robust to noisy matches, especially in the vicinity of large missing regions, as displayed in their Figure 8. This sensitivity to noise is increased by the fact that the noise produced by each matching algorithm displays slightly different patterns. To overcome these difficulties, a trained algorithm like ours that learns the noise patterns is more suitable. We suggest a new interpolation method based on a deep convolutional neural network. The method is applied in a feedforward manner and leads to an improvement in both accuracy and speed over the EpicFlow method.

Finally, it is noteworthy that some of the new optical flow methods do not rely on the aforementioned pipeline [36, 16]. One interesting example is presented by Dosovitskiy et al. [9] in their FlowNet model. They present an end to end convolutional neural network for optical flow that outputs a dense flow map. While their method does not reach the state of the art performance, it runs in real-time and demonstrates the power of feedforward deep learning in optical flow estimation.

3. Network Architecture

The optical flow dense interpolation problem is defined in the following way: given a sparse and noisy set of matches between pixels $M = \{(p_m, p'_m)\}$, we want to approximate the dense flow field $F : I \rightarrow I'$ between a source image I and a target image I' . To solve this problem, we use a fully convolutional network with no pooling. The main branch of the network consists of ten layers, each applying a 7×7 convolution filter followed by an Elu [5] non-linearity (Fig. 1). We use zero-padding to maintain the same image dimensions at each layer of the network.

3.1. Network Input

The input to our algorithm is a set of sparse and noisy matches M . These matches can be produced by any third party matching algorithm. In our experiments, we used several of the leading matching algorithms: FlowFields (FF) [2], CPM-Flow (CPM) [14], DiscreteFlow (DF) [27], and finally DeepMatching (DM) [41]. From the matches, we produce a sparse flow map of size $h \times w \times 2$ where h and w are the height and width of the image pair. Each pixel is initialized with the displacement to its match in the x and y axis. Missing pixels are filled with zeros. Apart from the sparse flow map, we add two additional matrices as guiding inputs to the networks: A binary mask of the missing pixels, and the edges map (Fig. 1).

We create a binary mask of all the missing pixels to indicate their position to the network (since zero can be a valid

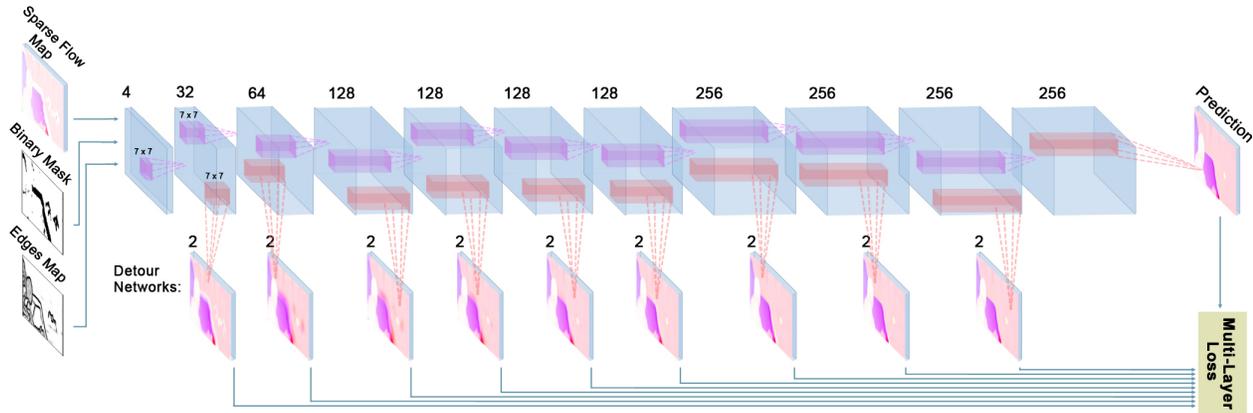


Figure 1: The architecture of InterpoNet.

displacement value). It was shown by others [19] to enhance performance in deep neural networks for inpainting.

The last input to the network is an edges map of one of the images in the image pair for which the flow is computed. The contours of an image were shown to be a key feature in image processing in the early visual cortex [11, 42, 43, 6, 30]. EpicFlow [33] already showed the benefit of the image edges as motion boundaries for optical flow estimation. In our work, we show evidence that a learning system also benefits from receiving the edges as input (see Fig. 4). We used an off-the-shelf edges detector - the "structured edges detector" (SED) [8] - the same one used by EpicFlow.

All of the inputs are stacked together and downsampled by 8 to form an $h/8 \times w/8 \times 4$ matrix. Rather than a simple stacking, we also considered different ways of introducing the edges map into the network. Among others, we have tried feeding the edges to all layers in the deep network, feeding the map to a different network and combining its output with the main branch in a deeper layer as well as constructing different networks to deal with pixels around the edges and far from the edges. However, we found that the simplest approach used here produced the best results.

3.2. The lateral dependency loss

To optimize the network results, we used the EPE (End Point Error) loss function, which is one of the standard error measures for optical flow. It is defined as the Euclidean distance between two flow pixels:

$$EPE(p_1, p_2) = \|p_1 - p_2\|_2 \quad (1)$$

The loss for an image pair was the average EPE over pixels:

$$L_{epe} = \frac{1}{n} \sum_{i,j} EPE(\hat{Y}_{i,j}, Y_{i,j}) \quad (2)$$

Where \hat{Y} is the network prediction, Y is the ground truth flow map and n is the number of pixels in the flow map.

This standard loss by itself does not yield good enough results (see Sec. 4.2). We, therefore, resort to cortical neuronal filling-in processes in our search for better losses.

Neuronal filling-in is characterized by spatial spread of activation. There is evidence that the activation spread is mediated by both lateral and top-down connections. To imitate the lateral dependency between neighbors in the network, we define a new lateral dependency loss. This loss pushes the distance between neighboring pixels to be similar to the distance in the ground truth flow. It is defined in the following way:

$$L_{ld} = \frac{1}{n} \sum_{i,j} |EPE(\hat{Y}_{i,j}, \hat{Y}_{i-1,j}) - EPE(Y_{i,j}, Y_{i-1,j})| + |EPE(\hat{Y}_{i,j}, \hat{Y}_{i,j-1}) - EPE(Y_{i,j}, Y_{i,j-1})|$$

The proposed loss term directly includes the local spatial dependencies within the training process, similar to what happens in the early stages of the visual cortex [1, 15].

3.3. Multi-layer loss using detour networks

Top-down connections are tricky to implement in artificial neural networks. We, therefore, use the loss function, which is the main feedback to the network, to imitate top-down connections. Also inspired by the evidence that neuronal filling-in takes place in many layers in the visual system hierarchy [30, 25, 24], we used detour networks connecting each and every layer directly to the loss function.

During training, the loss function served as top down information pushing each layer to perform interpolation in the best possible manner. The detour networks were kept simple: aside from the main branch of the network, each of the layer's activations was transformed into a two channels flow map using a single convolution layer with linear activations

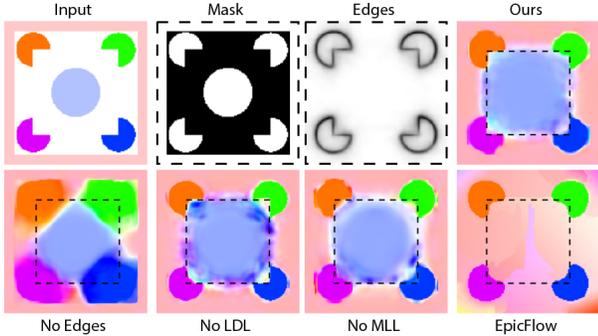


Figure 2: The network prediction for the Kanizsa illusion.

(no nonlinearity, see Fig. 1). Each of the flow maps produced by the detour networks was compared to the ground truth flow map using the EPE and LD losses. The final network loss was the weighted sum of all the losses:

$$L_{net} = \sum_l w^l (L_{epe}^l + L_{ld}^l) \quad (3)$$

Where w^l , L_{epe}^l and L_{ld}^l are the weight, EPE loss and LD loss of layer l . We found that weights of 0.5 for each of the middle layers and 1 for the last yielded the best results. For inference, we use only the last detour layer output - the one connected to the last layer of the network’s main branch.

Our approach has some similarities to the one used in the inception model introduced by Szegedi et al. [37], which employs auxiliary networks with independent losses during training. They found it to provide regularization and combat the vanishing gradients problem. However, in their network, the first auxiliary network was added in the tenth layer. We found that adding a detour network for each layer gave the best results. Szegedi et al.’s auxiliary networks were also built of several layers and performed some computation within them. We found that the simplest linear convolution was the best architecture. Additional layers or nonlinearities did not improve the performance of the network.

Taken together, our network was equipped with mechanisms with which it could imitate interpolation in the visual cortex. Interestingly, not only did it learn to perform interpolation of regular motion, it also performed strikingly similar to the visual cortex, when presented with an illusion. Figure 2 shows the interpolation applied by different variants of our network and EpicFlow on a given Kanizsa like motion pattern. The network never saw such a pattern in training time. When masking parts of the image, our network interpolates the motion pattern from the background and the interior. The propagation from the background stops in the borders of the imaginary square contour (marked by a dashed line), much like our visual perception. Importantly, only the real edges, not those of the imaginary contour, were fed to the network. Other networks that were not equipped

with all the tools we presented as well as EpicFlow, performed different levels of a simpler interpolation.

3.4. Post-processing

Our fully convolution with zero padding and no pooling network produces an output in the same size of the input. We, therefore, upsample the output by a factor of 8 using bi-linear interpolation. Like others before us [9], we found that using the variational energy minimization post-processing used in EpicFlow [33] slightly improved our final prediction (0.25px. gain in mean EPE). We employ the same parameters as EpicFlow, as appears in their Section 4.

4. Experiments

We report the results of our network on the Sintel [3], KITTI 2012 [12] and KITTI 2015 [26] datasets. We also show the effectiveness of different features in the network: the lateral dependency loss, the multi-layer loss and the edges input.

4.1. Training details

Preprocessing. As described in Section 3, the network receives a four channel input composed of two sparse flow channels given as the output of a matching algorithm, a binary mask and the edges map. To reduce training time, we downsample all the inputs by 8 (some matching algorithms output a downsampled version by default [41, 27, 14]). To reduce the number of missing pixels in training time, we apply bi-directional averaging (see supplementary).

We apply flipping as our only data augmentation method. Other transformations such as scaling, shearing, rotating and zooming did not improve the network performance, probably due to the interpolations that accompany them and drastically change the flow map.

Datasets. We evaluated our network on the three main optical flow benchmarks: MPI Sintel [3] is a collection of several scenes taken from a graphical animation movie. Each scene consists of several consecutive frames for which a dense ground truth flow map is given (a total of 1041 training image pairs). The scenes are diverse and include battle scenes with challenging large displacements. KITTI 2012 [12] is composed of real world images taken from a moving vehicle (194 training images). And KITTI 2015, [26] is similar to the KITTI 2012 dataset but with more challenging scenes (200 training images).

Since convolutional networks demand a large set of training data, we use the same approach used by Dosovitskiy et al. [9]. For initial pre-training, we use the Flying Chairs dataset that they introduced. This is a relatively large synthetic dataset (22,875 image pairs) composed of chair objects flying over different backgrounds. We train on all the dataset and use a sub-sample of the Sintel dataset as validation. Due to time constraints, we could not apply all of

Loss	EPE
EPE only	6.104
EPE + LDL	5.833
EPE + MLL	5.656
EPE + LDL + MLL	5.470

Table 1: Comparing losses for the Sintel 'final' pass validation set, trained on the output of FlowFields.

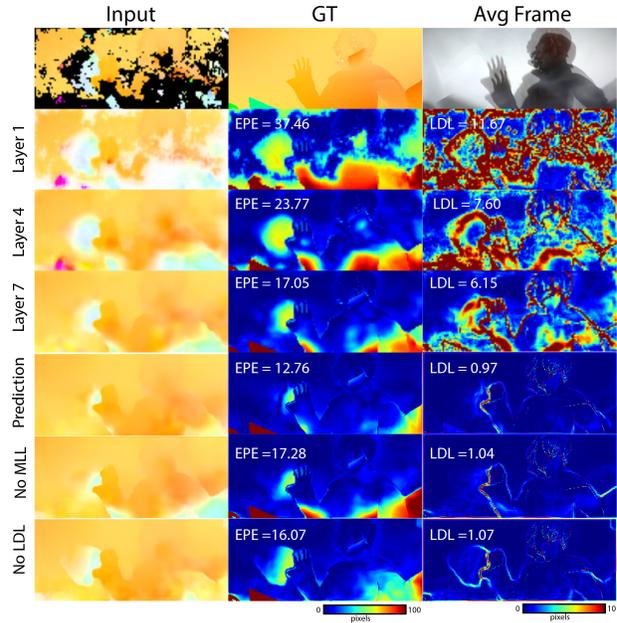
the matching algorithms to the big flying chairs dataset. We, therefore, used only FlowFields [2] for this initial training on Flying Chairs. Additional fine tuning was applied using the training sets of specific benchmarks and for the specific matching algorithm (see supplementary). In all presented experiments to follow, we pre-train the networks on Flying Chairs and fine tune on Sintel using the FlowFields matching algorithm - unless stated otherwise. All the analysis, results and visualizations are done without the variational post-processing, except for the benchmark results.

Optimization. We use Adam [18] with standard parameters ($\beta_1 = 0.9, \beta_2 = 0.999$). A learning rate of 5×10^{-5} for the pre-training and 5×10^{-6} for the fine tuning is used.

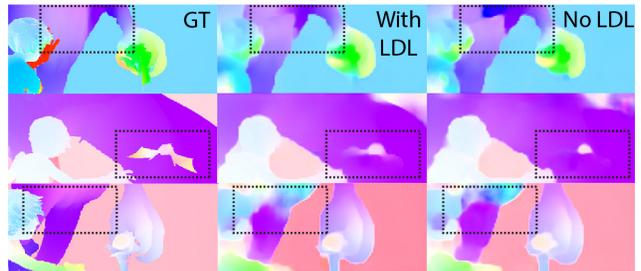
4.2. Comparison of loss variants

To ensure the efficiency of the different losses and the new architecture we introduce, we trained several variants of our network - with only the EPE loss, with the EPE + LD loss and with the EPE + ML loss. As Table 1 shows, each of our introduced losses yields a performance boost. Figure 3a shows the output of the different detour networks in different layers as well as the error maps for the two losses we used - EPE and LD loss, for an example image in our Sintel validation set. Notice how both the EPE and LD loss improves as the network deepens - this is consistent over all of the images in the validation set (Fig. 3c). At the first layers of the network, it seems that it is focused on performing a simple interpolation to mainly fill the missing parts. This initial interpolation is less aware of the motion boundaries. As the network deepens, it mainly polishes the details and reduces noise according to the segmentation introduced by the edges (for example. the green patches in Fig. 3a left column). The prediction of a network trained without multi-layer loss is noisier (6th row in Fig. 3a, lower right corner). It seems that the added supervision in all the layers helps to extinguish errors in the interpolation and adjust it according to the motion boundaries.

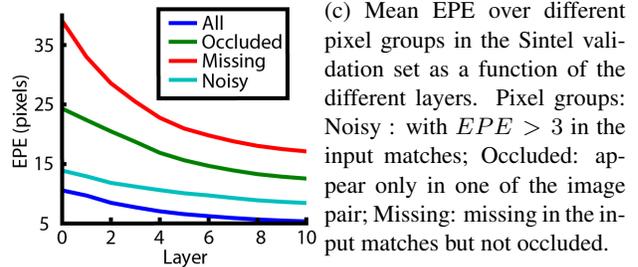
The LD loss is introduced to enforce a certain dependency between neighboring pixels. Much like in neuronal filling-in, lateral dependency plays a role in propagation, especially in terms of uncertainty. For example, in Figure 3a, there is a big missing part in the center-left with some false matches (light green in the flow map) to its right. The network can choose to either propagate the background or the false matches. To avoid the wrong local dependencies, the network with LD loss uses the background to fill most



(a) The progression of the prediction process throughout the different layers in the network, as shown by the detour networks outputs. Starting from the second row, the second and third columns are the EPE and LD loss maps respectively. The last three rows are the final predictions of networks with different losses. They are presented after upsampling which contributes to the decrease in LDL. Missing pixels in the input are marked in black.



(b) Comparison of the network performance with and without the LD loss. Left column is the ground truth, center and right columns are the predictions with and without LDL respectively.



(c) Mean EPE over different pixel groups in the Sintel validation set as a function of the different layers. Pixel groups: Noisy : with $EPE > 3$ in the input matches; Occluded: appear only in one of the image pair; Missing: missing in the input matches but not occluded.

Figure 3: The contribution of different losses.

of the area, almost extinguishing the false matches (notice the shrinking "bubble" in the LDL maps in Fig. 3a rows

Input	EPE
Sparse map + mask	6.240
Sparse map + mask + edges map	5.470

Table 2: Comparison of the network results with and without the edges as input. The reported EPE is for the Sintel 'final' pass validation set. The network is trained on the FlowFields algorithm output.

2-5). The network without the LD loss does not use this information and leaves a high contrast where it should not appear (Fig. 3a last row). LD loss mostly enforces smoothness in the outcome (Fig. 3b first row), but it also encourages high contrasts where they should appear, as shown in the example in Figure 3b (second row) for the wings of the small dragon. In rear cases, the LD loss combined with a poor input could decrease performance like in the third row of Figure 3b where the smooth transition introduced by the LD loss decreased the performance. Overall, the LD loss improves the EPE in 60% of predictions in our validation set, and 80% out of the noisy examples in the set (those with over one percent of noisy and missing pixels).

4.3. The importance of the edges

To validate the importance of the edges as an input to the network, we perform an experiment in which the edges are not fed into the network (in both train and test time). Table 2 shows the significant impact the edges input has on the performance of the network. Much like neuronal filling-in in the visual cortex, our network uses the edges as a boundary for local spread in missing or occluded areas. Figure 4a shows a comparison between the prediction of the two networks (with and without the edges input) for two examples from the Sintel validation set. Notice the spread into the missing pixels, while the network without the edges input performs what seems like a simple interpolation from all of the surroundings, the network with the edges input uses this information and stops the spread at the edges.

To quantify the effect of the edges on the missing and occluded pixels, we define an improvement index (II):

$$II_p = \frac{EPE_{p-noedges} - EPE_{p-edges}}{EPE_{p-noedges} + EPE_{p-edges}} \quad (4)$$

where $EPE_{p-noedges}$ and $EPE_{p-edges}$ are the EPE in pixel p between the prediction of the edges network and non-edges network respectively. Positive values of this index indicate improvement as a result of the edges input, while negative values indicate a decrease in performance. Mean II over occluded and missing pixels is significantly higher than the mean II over the non-missing pixels in the Sintel validation set ($Mean \pm SEM$ II difference = $0.0235 \pm 5.83 \times 10^{-3}$; paired t-test $p < 1 \times 10^{-4}$; $n=167$).

pre-training	Evaluated on	Fine tuned on		
		None	FF	Self
fc FF	Sintel FF	5.802	5.470	-
fc FF	Sintel CPM	6.165	5.782	5.851
fc FF	Sintel DM	6.498	6.075	5.971
fc FF	Sintel DF	6.543	6.35	6.142
fc DM	Sintel DM	6.665	-	5.934

Table 3: The network performance (EPE) without fine tuning, with FlowFields fine tuning and with fine tuning for the specific matching algorithm used for evaluation. EPE is reported for the Sintel 'final' pass validation set. Notations: fc - Flying Chairs, FF -FlowFields [2], CPM - CPM Flow [14], DM - DeepMatching [41], DF - DiscreteFlow [27]

Interestingly, as demonstrated in Figure 4b, the contribution of the edges input to the performance in the occluded and missing areas is not influenced by the distance from the edges. This is expected, since the decision about the spread is dependent on the segmentation by the edges and, therefore, even far away from the edges the effect is considerable (see top right corner in our prediction in the first row of Fig. 5 as an example). For non-missing pixels, however, the performance gains decrease almost monotonically with the distance from the edges (green line in Fig. 4b). These pixels are processed differently in the network, since they have initial values. They are more affected by their immediate surroundings. Therefore, a nearby edge can improve their prediction but less so far from edges. In all distances, the II values are significantly higher for the missing and occluded pixels (Wilcoxon signed rank test $p < 0.05$).

4.4. Fine tuning

Our network is trained in two phases. First, it is pre-trained on the flying chairs dataset using the FlowFields matching algorithm followed by fine tuning to the specific dataset and matching algorithm at hand. Table 3 shows the performance of the networks trained only on the flying chairs dataset, compared to the networks fine tuned on the Sintel training set with either the FlowFields matching algorithm or the same matching algorithms used for evaluation. The network performance is quite good even without fine tuning. However, the fine tuning phase still improves the performance by a considerable margin. Fine tuning on FlowFields applied to Sintel yields results comparable to fine tuning on the evaluation algorithm. Finally, using a different matching algorithm for pre-training (DeepMatching, last line in table 3) does not improve the results. We, therefore, suggest the best practice for incorporating new matching algorithms with our method as follows: For most cases using the network trained and fine tuned on FlowFields as an out-of-the-box solution should be sufficient.

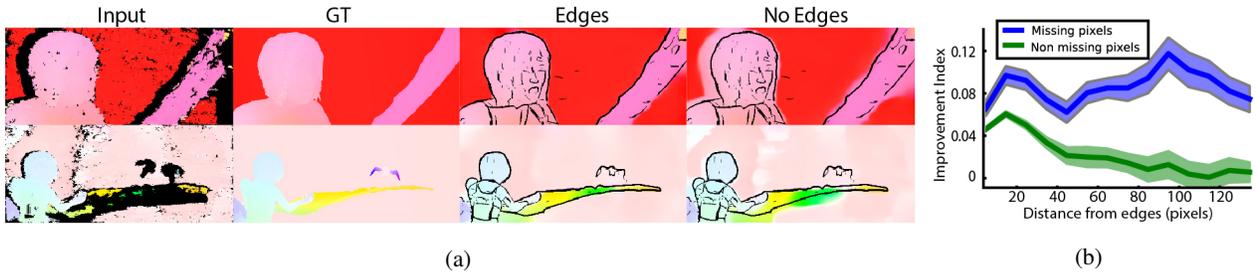


Figure 4: The contribution of the edges input to the network. (a) The predictions of the networks with and without the edges input. Edges are marked with black lines. (b) Mean improvement index over missing (blue) and non-missing (green) pixels. Shaded areas marks \pm SEM over image pairs in the Sintel validation set.

Method	EPE	EPE-noc	EPE-occ
FF [2]+Ours	5.535	2.372	31.296
PGM-C (anon.)	5.591	2.672	29.389
RicFlow (anon.)	5.620	2.765	28.907
CPM [14]+Ours	5.627	2.594	30.344
FF+ [2]	5.707	2.684	30.356
DM [41]+Ours	5.711	2.650	30.642
Deep DF [13]	5.728	2.623	31.042
FN2-ft-s(anon.)	5.739	2.752	30.108
SBFlow (anon.)	5.734	2.676	30.654
FF [2]+Epic	5.810	2.621	31.799
SPM-BPv2 [22]	5.812	2.754	30.743
FullFlow [4]	5.895	2.838	30.793
CPM+Epic [14]	5.960	2.990	35.14
FN2(anon.)	6.016	2.977	30.807
GPC [40]	6.040	2.938	31.309
DF [27]+Ours	6.044	2.788	32.581
DF [27]+Epic	6.077	2.937	31.685
DM+Epic [33]	6.285	3.060	32.564

Table 4: Leading results for the Sintel benchmark using the 'final' rendering pass. EPE-noc and EPE-occ are the EPE in non-occluded and occluded pixels respectively.

For improved results, we suggest fine tuning on the specific dataset and matching algorithm. Pre-training on the specific matching algorithm applied to the flying chairs dataset is not necessary, although it could be beneficial in some cases.

4.5. Benchmarks results

We applied our method to the output of several of the leading matching algorithms for Sintel and KITTI. The chosen matching algorithms are the highest on the leaderboards that have an available code and a reasonable running time. We used FlowField [2], DiscreteFlow [27] and CPM-Flow [14]. We also used DeepMatching, since it was used in the original EpicFlow paper [33].

For Sintel (Table 4), we achieve state of the art results using FlowFields as the matching algorithm. For all the matching algorithms used, we achieve better results com-

Method	2012 - EPE		2012 - %Out		2015
	Noc	All	Noc	All	%Out-All
DF [27]+Ours	1.0	2.4	4.94	14.13	23.55
DF+Epic	1.3	3.6	6.23	16.63	22.38
CPM+Ours	1.0	2.5	5.28	14.57	23.84
CPM [14]+Epic	1.3	3.2	5.79	13.70	23.23
FF+Ours	1.1	2.6	5.57	14.76	–
FF [2]+Epic	1.4	3.5	5.77	14.01	–
DM+Ours	1.1	2.7	5.85	15.03	24.65
DM+Epic [33]	1.5	3.8	7.88	17.08	27.10

Table 5: KITTI 2012 and KITTI 2015 benchmarks results. The %Out is the percentage of outlier pixels as defined by the benchmarks. FF does not have results on KITTI2015.

pared to EpicFlow improving the EPE by an average of 0.3px. Our performance is better in most areas including occluded, non-occluded and pixels in different distances from occlusion boundaries (with the exception of occluded pixels in CPM-flow and discrete flow). Figure 5 shows a comparison of EpicFlow's and our outputs on several sparse flow maps produced by FlowFields for the Sintel validation set. Notice the performance difference in missing areas with noise (top right corner in the first row, the hand in the third and bottom right in the last row). Due to its non-learning nature, EpicFlow is clinging to any information that it finds within a segmented area and is, therefore, prone to fail in such regions. The flexibility of a data-driven algorithm, like ours, is more suitable here. Further analysis demonstrated the superior performance of our method over EpicFlow in different regions (see supplementary). Based on our results, we believe that applying our method to a matching algorithm ranked higher than FlowFields (ranked 7 before our contribution) should yield even better results.

For KITTI 2012 [12], using DiscreteFlow [27] as the baseline matching algorithm, we achieve state-of-the-art results out of the published, pure optical flow methods, excluding semantic segmentation methods. We have the best performance, both in terms of EPE and the percentage outlier for non-occluded pixels (Table 5). Compared to

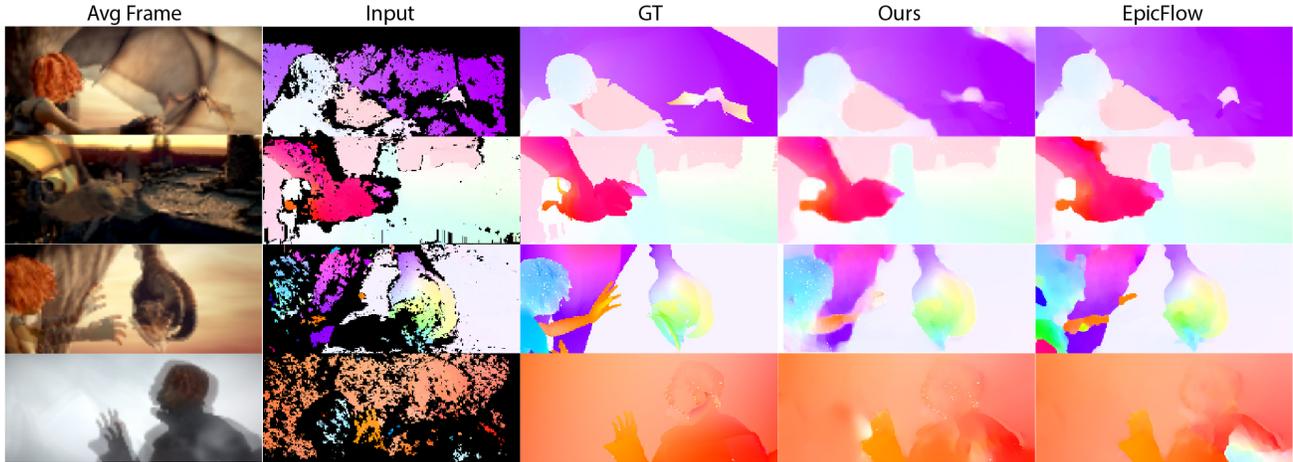


Figure 5: A comparison of the predictions of our network to EpicFlow.

EpicFlow, the EPE is improved by a margin (21%–33%), using all matching algorithms. The %Out measurement used in the KITTI datasets, calculates the percentage of pixels with $EPE > 3$. It is not linearly correlated with the EPE which we use as the target measurement, as reflected from the network’s loss function. Consequently, while this measurement was improved for non-occluded pixels (3%–25%) we achieved mixed results for all pixels (-6%–+15%; Table 5). Our results for KITTI 2015 [26], which uses only the %Out as the evaluation system, were also mixed (Table 5). The EPE measurement is not available in this benchmark, but our KITTI 2012 results support the possibility of an improvement in the EPE that is not reflected in the %Out. The results for our validation set were better than EpicFlow using all the matching algorithms (see supplementary).

4.6. Runtime analysis

Table 6 shows the runtime of the different components of our algorithm computed for one Sintel image pair (1024x436 pixels). The network inference ran on one NVIDIA GTX Titan black GPU (6GB RAM) while the other steps were performed on a single 3.4GHz CPU core. The run time of the edges detection and variational post-processing is as reported in [33]. The entire runtime was 1.333 seconds. This is slightly better than the reported runtime for EpicFlow (1.4 seconds). Notably, several parts in the pipeline could be dropped for better runtime without a big decrease in performance. The bi-directional average can be dropped in inference time (which will also reduce the downsampling by half), as well as the variational post processing, leaving the edges detection as the biggest bottleneck. Therefore, without much performance loss, our method can be as fast as 5 fps. Combined with a fast edge detection and matching algorithm, future work can produce a real-time optical flow algorithm.

Step	runtime (sec)
Downsampling	0.058
Bi-directional average	0.091
Edges detection	0.150
Network inference	0.025
Upsampling	0.009
Variational post proc.	1
Total	1.333

Table 6: Runtime of various steps of our solution for an image pair in the Sintel dataset.

5. Conclusions

Using a fully convolutional neural network, we have presented a data-driven solution for sparse-to-dense interpolation for optical flow producing state-of-the-art results. Our solution was inspired by ideas taken from interpolation processes in the visual cortex. We embedded anatomical features, like lateral dependency and multi-layer processing, by using the loss function, thereby applying supervision rather than using the architecture of the network which contributes to the simplicity of our solution. We also showed that the edge information is crucial for learning to interpolate. The network learns to use the edges as stoppers for the spread of interpolation, much like in the visual cortex.

Our solution is robust and can be applied to the output of different matching algorithms and our code and models will be made completely public. We encourage new solutions to use our method as part of their pipeline.

6. Acknowledgments

This research is supported by the Intel Collaborative Research Institute for Computational Intelligence (ICRI-CI) and by the Israeli Ministry of Science, Technology, and Space.

References

- [1] J. M. Alonso. Neural connections and receptive field properties in the primary visual cortex. *Neuroscientist*, 8(5):443–456, 2002. 3
- [2] C. Bailer, B. Taetz, and D. Stricker. Flow fields: Dense correspondence fields for highly accurate large displacement optical flow estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 11-18-Dece, pages 4015–4023, 2016. 2, 5, 6, 7
- [3] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, 2012. 4
- [4] Q. Chen and V. Koltun. Full Flow: Optical Flow Estimation By Global Optimization over Regular Grids. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 7
- [5] D.-A. Clevert, T. Unterthiner, and S. Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015. 2
- [6] J. Dai and Y. Wang. Representation of surface luminance and contrast in primary visual cortex. *Cerebral Cortex*, 22:776–787, 2012. 3
- [7] P. De Weerd, R. Gattass, R. Desimone, and L. G. Ungerleider. Responses of cells in monkey visual cortex during perceptual filling-in of an artificial scotoma. *Nature*, 377:731–734, 1995. 2
- [8] P. Dollar and C. L. Zitnick. Structured Forests for Fast Edge Detection. *2013 IEEE International Conference on Computer Vision*, pages 1841–1848, 2013. 3
- [9] A. Dosovitskiy, P. Fischery, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. V. D. Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 11-18-Dece, pages 2758–2766, 2016. 2, 4
- [10] H. S. Friedman, H. Zhou, and R. V. D. Heydt. Color filling-in under steady fixation: Behavioral demonstration in monkeys and humans. *Perception*, 28(11):1383–1395, 1999. 2
- [11] H. S. Friedman, H. Zhou, and R. von der Heydt. The coding of uniform colour figures in monkey visual cortex. *The Journal of physiology*, 548(2003):593–613, 2003. 3
- [12] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 4, 7
- [13] F. Güney and A. Geiger. Deep Discrete Flow. In *Asian Conference on Computer Vision (ACCV)*, 2016. 7
- [14] Y. Hu, R. Song, and Y. Li. Efficient Coarse-to-Fine Patch-Match for Large Displacement Optical Flow. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, (June):5704–5712, 2016. 2, 4, 6, 7
- [15] X. Huang and M. a. Paradiso. V1 response timing and surface filling-in. *Journal of neurophysiology*, 100(1):539–47, 7 2008. 2, 3
- [16] J. Hur and S. Roth. Joint Optical Flow and Temporally Consistent Semantic Segmentation. *CVRSUAD workshop at ECCV*, 2016. 2
- [17] P. J. Kellman, C. Yin, and T. F. Shipley. A common mechanism for illusory and occluded object completion. *Journal of experimental psychology. Human perception and performance*, 24(3):859–869, 1998. 2
- [18] D. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations*, pages 1–13, 2014. 5
- [19] R. Köhler, C. Schuler, B. Schölkopf, and S. Harmeling. Mask-specific inpainting with deep neural networks. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 8753, pages 523–534, 2014. 3
- [20] H. Komatsu. The neural mechanisms of perceptual filling-in. *Nature reviews. Neuroscience*, 7(March):220–231, 2006. 2
- [21] H. Komatsu, M. Kinoshita, and I. Murakami. Neural responses in the retinotopic representation of the blind spot in the macaque V1 to stimuli for perceptual filling-in. *The Journal of Neuroscience*, 20(24):9310–9319, 2000. 2
- [22] Y. Li, D. Min, M. S. Brown, M. N. Do, and J. Lu. SPM-BP: Sped-up patchmatch belief propagation for continuous MRFs. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 11-18-Dece, pages 4006–4014, 2016. 7
- [23] M. Liang and X. Hu. Recurrent convolutional neural network for object recognition. *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*, (Figure 1):3367–3375, 2015. 2
- [24] J. D. Mendola, A. M. Dale, B. Fischl, A. K. Liu, and R. B. Tootell. The representation of illusory and real contours in human cortical visual areas revealed by functional magnetic resonance imaging. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 19(19):8560–8572, 1999. 2, 3
- [25] M. Meng, D. A. Remus, and F. Tong. Filling-in of visual phantoms in the human brain. *Nature neuroscience*, 8(9):1248–54, 9 2005. 2, 3
- [26] M. Menze and A. Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June, pages 3061–3070, 2015. 4, 8
- [27] M. Menze, C. Heipke, and A. Geiger. Discrete Optimization for Optical Flow. In *37th German Conference, GCPR 2015*, pages 16–28, 2015. 2, 4, 6, 7
- [28] M. Paradiso and K. Nakayama. Brightness perception and filling-in. *Vision research*, 31(7):1221–1236, 1991. 2
- [29] B. Pinna, G. Brelstaff, and L. Spillmann. Surface color from boundaries: A new 'watercolor' illusion. *Vision Research*, 41:2669–2676, 2001. 2
- [30] J. Poort, F. Raudies, A. Wannig, V. a. F. Lamme, H. Neumann, and P. R. Roelfsema. The role of attention in figure-ground segregation in areas V1 and V4 of the visual cortex. *Neuron*, 75(1):143–56, 7 2012. 2, 3
- [31] V. S. Ramachandran*. Blind spots. *Scientific American*, 266(5):86 ST – Blind spots., 1992. 2

- [32] V. S. Ramachandran and R. L. Gregory. Perceptual filling in of artificially induced scotomas in human vision. *Nature*, 350(6320):699–702, 4 1991. [2](#)
- [33] J. Revaud, P. Weinzaepfel, Z. Harchaoui, and C. Schmid. EpicFlow: Edge-preserving interpolation of correspondences for optical flow. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June, pages 1164–1172, 2015. [1](#), [2](#), [3](#), [4](#), [7](#), [8](#)
- [34] A. W. Roe, H. D. Lu, and C. P. Hung. Cortical processing of a brightness illusion. *Proceedings of the National Academy of Sciences of the United States of America*, 102(10):3869–3874, 2005. [2](#)
- [35] Y. Sasaki and T. Watanabe. The primary visual cortex fills in color. *Proceedings of the National Academy of Sciences of the United States of America*, 101(52):18251–18256, 2004. [2](#)
- [36] L. Sevilla-Lara, D. Sun, V. Jampani, and M. J. Black. Optical Flow with Semantic Segmentation and Localized Layers. *Cvpr*, pages 1–10, 2016. [2](#)
- [37] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June, pages 1–9, 2015. [4](#)
- [38] R. von der Heydt, H. S. Friedman, and H. Zhou. Searching for the neural mechanisms of color filling-in. In *Filling-in: From perceptual completion to cortical reorganization*, pages 106–127. Oxford University Press, Oxford, 2003. [2](#)
- [39] T. Wachtler, T. J. Sejnowski, and T. D. Albright. Representation of color stimuli in awake macaque primary visual cortex. *Neuron*, 37:681–691, 2003. [2](#)
- [40] S. Wang, S. R. Fanello, C. Rhemann, S. Izadi, and P. Kohli. The Global Patch Collider. In *cvpr*, 2016. [7](#)
- [41] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid. DeepFlow: Large displacement optical flow with deep matching. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1385–1392, 2013. [2](#), [4](#), [6](#), [7](#)
- [42] G. Zurawel, I. Ayzenshtat, S. Zweig, R. Shapley, and H. Slovin. A contrast and surface code explains complex responses to black and white stimuli in V1. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 34(43):14388–402, 10 2014. [2](#), [3](#)
- [43] S. Zweig, G. Zurawel, R. Shapley, and H. Slovin. Representation of Color Surfaces in V1: Edge Enhancement and Unfilled Holes. *Journal of Neuroscience*, 35(35):12103–12115, 2015. [2](#), [3](#)