# Continual Occlusions and Optical Flow Estimation

Michal Neoral, Jan Šochman, and Jiří Matas

Center for Machine Perception, Faculty of Electrical Engineering
Czech Technical University in Prague, Czech Republic
{neoramic,jan.sochman,matas}@fel.cvut.cz

**Abstract.** Two optical flow estimation problems are addressed: i) occlusion estimation and handling, and ii) estimation from image sequences longer than two frames. The proposed ContinualFlow method estimates occlusions before flow, avoiding the use of flow corrupted by occlusions for their estimation. We show that providing occlusion masks as an additional input to flow estimation improves the standard performance metric by more than 25% on both KITTI and Sintel. As a second contribution, a novel method for incorporating information from past frames into flow estimation is introduced. The previous frame flow serves as an input to occlusion estimation and as a prior in occluded regions, i.e. those without visual correspondences. By continually using the previous frame flow, ContinualFlow performance improves further by 18% on KITTI and 7% on Sintel, achieving top performance on KITTI and Sintel.

## 1 Introduction

Optical flow is a two-dimensional displacement field describing the projection of scene motion between two images. Occlusions caused by scene motion contribute to the ill-posedness of optical flow estimation – at occluded pixels no visual correspondences exist. Classical non-CNN methods address this problem by using regularisation which extrapolates the flow from the surrounding non-occluded area. Current state-of-the-art CNN algorithms for optical flow use the correlation cost volume [9,18,31,35,26,16] to estimate the most likely correspondences. Their regularisation is only implicit and the network has to learn when to rely on the cost volume and when to extrapolate. In both cases, the occluded areas are processed the same way as non-occluded ones which leads to errors in the occluded areas as well as in the nearby non-occluded regions.

Approaches dealing with occlusions [26,1] usually first estimate initial forward and backward optical flows. Occlusions are found by a forward-backward consistency check and occlusion maps are then used for estimating of the final optical flow. The problem here is that occlusions affect the initial flow and thus the final output.

As our first contribution, we extend a current state-of-the-art CNN optical flow method [35] by estimating the occluded areas first, *without estimating the flow*, and then passing the occlusion maps to the optical flow estimation network.

The correlation cost volume for flow estimation is re-used for occlusion estimation. Intuitively, the cost will be low in non-occluded areas with good correspondences and high in occluded regions. While preserving end-to-end trainability, we accurately estimate occlusions and significantly improve the estimated flow.

Optical flow estimation over more than two frames is a problem whose difficulty stems from the need for the pixels to be mapped to a reference coordinate system before loss evaluation. The mapping is defined by the unknown optical flow itself. Hence, it is difficult to apply temporal regularisation before the flow is known. A typical solution over three frames is to use the middle one as the reference defining the coordinate system and to compute the forward flow to the future frame and the backward flow to the past frame and to apply regularisation to these two flows. Published multi-frame approaches assume various motion constraints: constant rigid motion for three images [41], adaptive trajectory regularisation over five images [38], multi-frame subspace constrains [19] and other complex motion models [12] over the whole sequence.

We avoid modelling the motion regularity explicitly and let a CNN model learn the relations of the current and previous optical flows. The CNN is fed pairs of consecutive images together with the flow computed between the penultimate and last images. We solve the coordinate system mapping by bilinear warp [20]. The proposed method is not limited to a fixed temporal horizon, the network uses previously estimated flows and thus, by recursion, all prior frames.

The two above-mentioned problems – occlusion estimation and the use of multiple frames – are related. Since there are no correspondences in occluded areas, optical flow cannot be estimated from the cost volume and the CNN is forced to use regularisation. Knowing the occlusions and given the previous flow, the network has prior information about the motion to be used when no correspondences are available. So, the last estimated flow is also fed into the occlusions estimation as it is a source of information about possible occlusions.

Finally, we add a specialised refinement network [18,29] to the proposed architecture. It has been shown to improve fine detail accuracy of the flow, which is confirmed by our experiments. We integrate this network with both occlusion estimation and temporal processing.

**Contributions.** We introduce integrated occlusion estimation, i.e. the algorithm does not operate on an occlusion-ignorant flow estimate, to the state-of-the-art PWC-net [35]. Second, we propose a novel method that implicitly uses all previous frames for optical flow estimation. Finally, we add refinement blocks with additional feature map inputs leading to improved spatial resolution of the final flow. ContinualFlow is state-of-the-art on several public benchmarks[1]: 1st place in Sintel [6][2] and 1st place in the KITTI'15 [27] optical flow benchmark among Robust Vision Challenge (ROB) participants and 3rd over all optical flow methods[3] with a large margin in precision in occluded areas. Continual flow ranked 3rd in ROB [32] for the optical flow category.

---

[1] As of the submission date, July 7, 2018.

[2] The "Final pass" category.

[3] Excluding scene flow methods.

## 2   Related Work

**Occlusion estimation and occlusion handling.** Most optical flow methods detect occlusions as outliers of the correspondence field [1,13,2] or by a consistency check on the estimated forward and backward optical flows [36,8]. The optical flow is then extrapolated into the occluded areas. The shortcoming of such approaches is that the initial flow is already adversely affected by the occlusions. Other methods incorporate occlusion estimation directly into the energy minimisation [42,37,34] by truncating the data term, avoiding the problematic post-processing of already affected optical flow. The current best performing non-CNN method [17] formulates optical flow estimation symmetrically - estimating the forward and backward flows, occlusions and dis-occlusions in a single joint optimisation.

Most of the current state-of-the-art CNN networks [9,18,31,35] do not explicitly deal with occlusions. The network in [26] estimates the forward and backward flows independently and uses the forward-backward consistency check to estimate the occlusions. The estimated occlusions are then used for network training only. In LiteFlowNet [16] an occlusion probability map is a function of brightness inconsistency between the reference frame and warped target frame. The occlusion probability map is used in a flow regularisation module.

To our best knowledge, no published CNN method estimates occlusions prior to optical flow estimation to improve the flow in the test phase.

**Using multiple frames.** Most methods that process more than two frames impose some kind of regularisation on the flow. Murray and Buxton [28] introduced an approach that uses spatio-temporal smoothness term which regularises optical flow trajectory over multiple frames. However, the algorithm does not work well for large displacements. Black et al. [5] extrapolate the flow from the previous frame as a starting point for the optimisation in the current frame. In Garg et al. [11], the motion regularisation was relaxed from several rigid motions into multi-frame subspace constraints allowing non-rigid motions. Multi-frame subspace constraints were used in [19] over long trajectories. Its extension [12] allows more complex motions using soft constraints between frames. An adaptive trajectory regularisation over five consecutive frames was used in [38], where optical flow was parametrised w.r.t. the central reference frame. Wulff *et. al.* [41] use super-pixel segmentation and a rigid motion assumption over triplets of images. ProFlow [23] uses three consecutive frames, a CNN regularises non-CNN-estimated forward ($I_t \rightarrow I_{t+1}$) and backward ($I_{t-1} \rightarrow I_t$) optical flows.

While many non-CNN algorithms use more than two frames in some form, to our best knowledge, no CNN-based method using more frames has been published. Unlike the above-mentioned approaches, the proposed method trains the regularisation from data and does not need any hand-crafted approximations.

**The refinement network.** The last important component added to the proposed architecture is a specialised refinement network [18,29]. We confirm it

improves accuracy of fine details of the flow. We integrate the network with both occlusion estimation and temporal processing.

The refinement network was introduced in [18] for optical flow estimation as a part of an architecture specialised on optical flow fine detail refinement. The inputs to the network are the optical flow estimated by previous blocks, the brightness error of the warped image and the input images themselves. In [18,29], it was shown that training the first flow estimation block and the refinement network sequentially leads to improvements in estimated optical flow.

## 3   ContinualFlow

The proposed ContinualFlow method builds on the state-of-the-art PWC-Net architecture [35]. We extend the architecture by adding i) occlusion estimation blocks and use the estimated occlusions for flow estimation, ii) an refinement network to improve fine detail accuracy, and iii) temporal connections for utilising the previous flow for estimation of both the flow and the occlusions. Fig 1 shows a schematic of the PWC-Net with both the occlusions estimation blocks and temporal connections. Another diagram containing also the refinement network is shown in Fig 2.

The original PWC-Net [35] is composed of two networks: a *feature pyramid extractor* and a coarse-to-fine *optical flow decoder*. The feature pyramid extractor takes as input two images $I_t$ and $I_{t+1}$ and encodes them into a pyramid of feature vectors $\mathcal{F}_t^s$ and $\mathcal{F}_{t+1}^s$ with gradually decreasing spatial resolution (indexed by $s$) and with increasing channel dimension. The decoder, in a coarse-to-fine manner, takes features from the corresponding resolution $s$, warps features $\mathcal{F}_{t+1}^s$ using the up-sampled flow $F_{t+1}^{s-1}$ estimated at a coarser iteration $s-1$ (if not at the coarsest resolution) and builds a correlation cost volume - a volume of feature correlations over a limited displacement range. The cost volume is then fed to the optical flow estimator, which produces the current scale optical flow $F^s$ and the process is repeated for higher resolution. We refer the reader to the original paper for further details. We are using the version with DenseNet [15] and a context network as described in the original paper.

### 3.1   Occlusion Estimation

PWC-Net and many other state-of-the-art approaches rely on the correlation cost volume for estimation of the optical flow [9,18,31,35,26]. Apart from being useful for the flow estimation, it is also indicative of possible occlusions. Intuitively, when the cost for all displacements for some pixel is high, the pixel is likely occluded in the next frame. In order to utilise this information, we propose to connect the occlusions estimator directly after the cost volume computation, even before any flow is estimated as shown in Fig 1. The output of the occlusions estimator is then sent to the optical flow estimator together with the cost volume itself. This way the occlusion estimation does not rely on the imprecise flow estimation and the flow estimator benefits from the additional input. Same
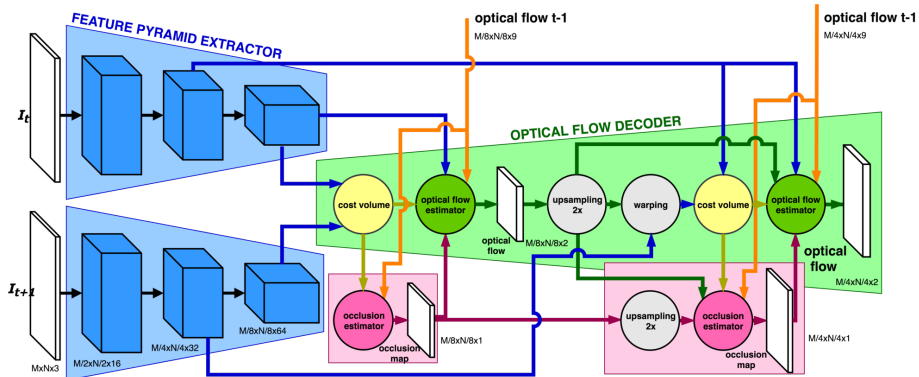
**Fig. 1:** ContinualFlow - optical flow and occlusion decoder, which extends the PWC-Net [35] flow decoder for occlusion estimation. The feature pyramid extractor (in blue) is a convolutional network which produces a feature pyramid given an input image. A correlation cost volume is computed on each scale from warped features from the second frame using up-sampled flow estimated at a coarser level of decoder. The cost volume is used to estimate occlusions in occlusion estimator (in magenta). The cost volume and the occlusion map are inputs to the optical flow estimator. For clarity, the diagram shows only three of the six levels of the ContinualFlow pyramid extractor. The output resolution is quarter of the input reference frame. Please, refer to the text for additional network details and inputs explanation.

as the flow estimator, the occlusions estimator works in a coarse-to-fine manner with higher resolution estimators receiving also up-sampled flow estimate from the lower resolution.

In experiments, we use an occlusion estimator with five convolutional layers with $D$, $\lfloor \frac{D}{2} \rfloor$, $\lfloor \frac{D}{4} \rfloor$, $\lfloor \frac{D}{8} \rfloor$ and two output channels (occluded/not occluded maps), where $D = 89$ in our case (the number of correlation cost volume layers + 8). All layers use ReLU activation except for the last one, which uses soft-max.

### 3.2   Refinement Network

It was shown that a specialised refinement network which processes the output of the initial network boosts the precision of the flow estimate, especially the fine details recovery [18,29]. The refinement network takes several extra inputs, like the current estimate of the optical flow, image $I_{t+1}$ warped back to time step $t$ and brightness error between $I_t$ and the warped $I_{t+1}$, and produces a refined optical flow [18].

The refinement network used in ContinualFlow has the same architecture as the optical flow decoder, but without the DenseNet connections. The main difference is in the network inputs. Instead of using the input images and their warps
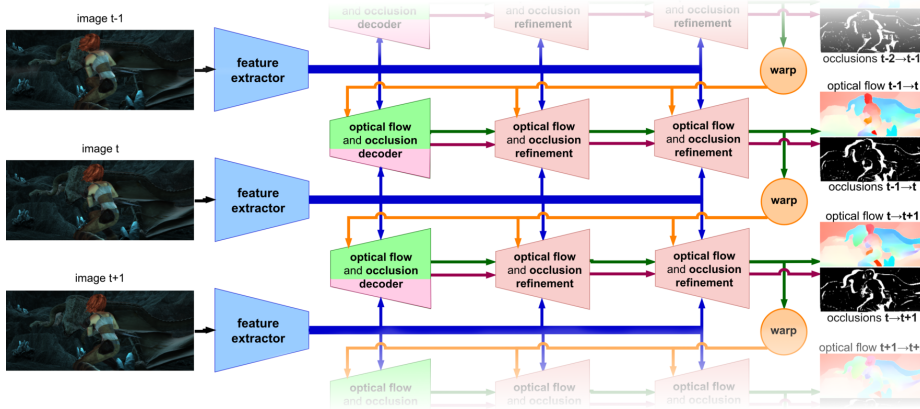
**Fig. 2:** Block diagram of ContinualFlow. Feature extractors with shared weights compute a feature pyramid from the input images. Features are input to the optical flow and occlusion decoder and the refinement blocks. The decoder estimates the optical flow and the occlusion map from the input features and from the temporal connection – the warped optical flow from the previous time step. Optical flow and occlusion maps are finalised by the refinement blocks.

as in [18], we use the features from the feature pyramid on the corresponding scale and their warps as a richer input representation. The input error channel for these features is computed as a sum of the $L_1$ distance and structure similarity (SSIM) [39]. We applied the refinement two times, additional refinements did not improve the accuracy in our experiments.

### 3.3   ContinualFlow Estimation over Image Sequence

We use temporal connections, which give the optical flow decoder, the occlusions decoder and the refinement network an additional input: the flow estimated in the previous time step (see the orange arrows in Fig 1 and Fig 2). When processing sequences longer than two frames these connections allow the network to learn typical relations between the previous and current flows and use them in the current frame flow estimation.

However, as discussed in Sec 1, the coordinate systems in which the two flows are expressed differ and need to be transformed onto each other in order to apply the previous flow to the correct pixels in the current time step. Here we describe two such transformations, forward and backward warping, and we test them independently as well as in combination (concatenation of both) in Sec 4.

**Forward warping transformation.** Forward warping transforms the coordinate system from time step $t-1$ using the optical flow $F_{t-1}$ itself. The warped flow $\hat{F}_{t-1}$ is computed as

$$\hat{F}_{t-1}\left(\mathbf{x} + \mathrm{round}(F_{t-1}\left(\mathbf{x}\right))\right) = F_{t-1}\left(\mathbf{x}\right), \tag{1}$$

for all pixel positions $\mathbf{x}$. For positions to which the flow $F_{t-1}$ maps more than once we preserve the larger of the mapped flows. This prioritises larger motions, thus faster moving objects. Although the experiments show usefulness of this warping, the main disadvantage of this approach is that the transformation is not differentiable. Thus, the training cannot propagate gradients through this step and relies on the shared weights only.

**Backward warping transformation.** Alternatively, the coordinate system could be transformed using the backward flow $B_t$ from frame $t$ to frame $t-1$. This requires an extra evaluation of the network, but then the warping is a direct application of the differentiable spatial transformer [20]. Thus, in this case the gradients are propagated through the temporal connections during training. A disadvantage of this approach is the computationally expensive computation of the backward flow.

**Combining forward and backward warping.** It is possible to use both warpings at the same time. In ContinualFlow we combine forward warped previous flow, backward warped previous flow and backward flow by simply concatenating their outputs. The only difference becomes that the previous flow input has nine channels: three times two for the flow warps and a validity masks for each warp (set to zero if the measurement is not available, e.g. at the beginning of the sequence).

**Multi-frame sequence initialisation.** The network is fed a pair of input images and the previously estimated flow. For the first frame in the sequence, no previous flow estimation is available. We estimate the initial optical flow between the first and second frame twice. First, we mask out the temporal connection and, in the second estimation, we use the first estimate as a temporal input.

### 3.4   Training Loss

The network is trained end-to-end with a weighted multi-task loss over the flow and occlusions estimators at all scales,

$$\mathcal{L} = \sum_{s=1}^{S} \alpha^s \mathcal{L}_F^s + \alpha_O \sum_{s=1}^{S} \alpha^s \mathcal{L}_O^s \,, \tag{2}$$

where $\alpha^s$ is the weight of individual scale $s$ losses and $\alpha_O$ is the occlusion estimation weight. The sums are over all $S$ spatial resolutions. The flow estimator loss $\mathcal{L}_F$ is the same as in PWC-Net, i.e. the end-point error

$$\mathcal{L}_F^s = \sum_{\mathbf{x}} \gamma(\mathbf{x}) ||F^s(\mathbf{x}) - F_{gt}^s(\mathbf{x})||_2 \,, \tag{3}$$

where $F^s$ is the estimated optical flow at scale $s$, $F_{gt}^s$ the corresponding ground-truth optical flow and $\gamma$ is the valid ground-truth flow mask (one for valid flow

and zero otherwise). The sum is over all pixel positions. As in [43,7] we adopted the weighted pixel-wise cross-entropy loss for occlusion map estimation

$$
\begin{aligned}
\mathcal{L}_O^s = - \, w_{noc} \sum_{\mathbf{x}:\, O_{gt}(\mathbf{x})=1} \rho(\mathbf{x}) \log \Pr(O(\mathbf{x})=1|X) \\
- \, w_{occ} \sum_{\mathbf{x}:\, O_{gt}(\mathbf{x})=\mathbf{0}} \rho(\mathbf{x}) \log \Pr(O(\mathbf{x})=0|X) \,,
\end{aligned}
\tag{4}
$$

where $\Pr(O(\mathbf{x})=1|X)$ is computed using soft-max $\sigma(\cdot)$ function on the occlusion estimator output $O$, $O_{gt}$ is the ground truth occlusion map, $\rho$ the valid ground-truth occlusion mask used for masking out images without ground-truth occlusions, and $w_{occ}$ and $w_{noc}$ are the fractions of occluded and non-occluded ground truth pixels respectively.

As suggested by [35], we modify this loss for the final fine-tuning on the most complex evaluation benchmark datasets. Here we change the $\mathcal{L}_F^s$ loss to the generalised Charbonnier loss (with $q = 0.4$, $\epsilon = 0.01$ as in [35]):

$$
\mathcal{L}_F^s = \sum_{\mathbf{x}} \gamma(\mathbf{x}) \left( |\hat{F}^s(\mathbf{x}) - F_{\mathbf{gt}}^s(\mathbf{x})| + \epsilon \right)^q .
\tag{5}
$$

## 4    Experiments

**Training details.** The ContinualFlow network is trained using a curriculum learning approach [4] starting from a dataset with less complex motions and increasing gradually the task complexity [18,35]. First, we train on FlyingChairs dataset [9] using the training parameters introduced in [35] and following the learning rate schedule from [18]. We do not use rotation, scaling and translation augmentations. Since the FlyingChairs dataset contains only two frames sequences and no occlusion ground truth, we cannot train the full ContinualFlow model with temporal connections and the occlusion map estimation. Instead, we use it for pre-training the PWC-Net part of the ContinualFlow network. The network is trained for 1200k iteration and the learning rate 1e-4 is divided by 2 each 200k iteration, starting from 400k. Images in a batch of size eight are randomly cropped to $448 \times 384$ px.

Next, the all parts of the ContinualFlow network are trained on the FlyingThings dataset [25]. Since occlusion maps were not available for this dataset, we computed them using the available backward and forward ground truth flows and the object segmentation masks. The mask $O_t(\mathbf{x})$ is set to "occluded" for pixel $x$ when the object labels $L_t(\mathbf{x})$ and $L_{t+1}(F_{gt}(\mathbf{x}))$ differ or the bi-directional consistency between backward and forward flows differs by more than one pixel. The network is trained for 500k iteration and the learning rate, set to 1e-4 for the first 200k iterations, is divided by 2 at that point and after 100k iterations. First, we train the network without the refinement. Then, only the refinement is trained while all other weights are fixed. Images in the batch of size four are randomly cropped to $768 \times 384$ px. After cropping, optical flow pointing out of the frame is labelled as occluded.

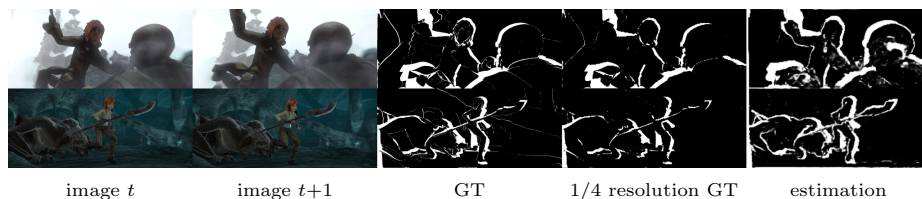| image $t$ | image $t+1$ | GT | 1/4 resolution GT | estimation |

**Fig. 3:** Example estimated occlusion maps on the Sintel (final) dataset, our validation split. ContinualFlow estimates occlusions up to quarter resolution.

Finally, the ContinualFlow is trained on data from six datasets: Driving [25], KITTI'15 [27], VirtualKITTI [10], Sintel [6], HD1K [22] and the FlyingChairs small motions dataset [18]. These datasets, except for FlyingChairs, contain sequences longer than two frames and are suitable for the training of temporal connections. We used the first image twice for the FlyingChairs dataset to obtain the same batch size for all input data, the loss on the estimate of the (zero) flow $F^{0,1}$ is not used. Dense occlusion maps are available only for the Sintel and Driving datasets. We set occlusion estimation loss to zero on the rest. The network is fine-tuned for 500k iteration and the learning rate, set to 1e-5 for the first 200k interactions, is divided by 2 at that point and after 100k iterations. Images in batches of size four are randomly cropped to $768 \times 320$ px. We sample images from all datasets uniformly.

We set weights for individual scales as in [35]. Maximal displacement in the cost volume is set to four. The same scale weights are set to train the refinement network and for the occlusion map estimation. The occlusion estimation weight $\alpha_O$ is set to 0.1. All experiments are trained with the ADAM optimiser [21] and 0.0004 weight decay. All parts of the network are implemented in TensorFlow.

The ContinualFlow training has the same three phases as training of PWC-Net. Only when training the refinement network separately, there is an additional phase which updates only the refinement parameters as mentioned above. ContinualFlow without the refinement network has 9.6M parameters, 0.8M more than the PWC-Net. The refinement network adds 5.0M parameters, it is based on the PWC-Net-small architecture. ContinualFlow runs at 8 FPS on KITTI-resolution of 1240x375 px.

In the following, we focus on the Robust Vision Challenge [32], where one trained model with the same parameters has to be evaluated on four individual benchmarks [27,6,22,3] instead of fine-tuning for each particular dataset independently.

### 4.1 Ablation Study

In this section, we experimentally evaluate the individual contributions and design choices for the ContinualFlow network trained on FlyingChairs [9] and fine-tuned on FlyingThings [25] as described above. Below, the term *baseline* refers to our TensorFlow implementation of PWC-Net. Unlike the PWC-Net settings [35], we trained the network without rotation, scaling and translation augmentation of input frames.

**Occlusion map learning.** Table 1 shows the results of optical flow estimation with and without occlusion learning. Temporal connections are not used. Application of the occlusion map estimator improves performance on all tested datasets not only in occluded regions but also in all non-occluded regions. Fig 3 shows example estimated occlusion maps.

**The specialised refinement block** improves results of the estimated optical flow as is shown in [18]. Table 1 compares the optical flow estimation with and without the refinement block. No temporal connections are used. The refinement block improves the estimated optical flow, especially in occluded areas.

**Influence of the coordinate warping methods.** We evaluated the three approaches for warping the previous flow estimate introduced in section 3.3. Results for individual datasets are shown in Table 1. Forward warping $W_f$ is beneficial for the KITTI dataset [27] and the Sintel Clean dataset [6], while backward warping $W_b$ is more suitable for the complex Sintel Final sequences. The combination of both, $W_{bf}$, is the most accurate on FlyingThings sequences [25]. All evaluated variants use the occlusion estimator in the decoder and no refinement.

**Temporal connection placement.** We experimented with passing the warped optical flow from previous frame to different network components, thus creating different temporal connections. In one variant, only the refinement network received the previous frame flow estimates. In another variant, all temporal connections as depicted in Fig 2 were used. Table 1 shows how feeding these connections with different warpings influences the estimated flow. The best results were obtained with temporal connections leading into both the decoder and refinement networks and the combination of forward and backward warpings.

**Number of refinement blocks.** Table 1 shows results for 1, 2, 3 and 5 stacked refinement networks. Stacking more than two refinement networks is not beneficial. Thus the final network architecture contains only two refinements. All evaluated variants use the occlusion estimator and warps the previously estimated optical flow using both warping methods in the first part of the network and the refinement.

**Multi-frame sequences initialisation** For the first frame in the sequence, no previous frame flow estimate is available to be passed to the temporal connections. Unsurprisingly, the estimation on the first frame is usually slightly worse than at the consecutive frames. We tested two initialisations of the first frame flow estimation: (i) no flow (zero displacements) instead of the previously estimated optical flow and, (ii) a two-pass initial estimation of the currently estimated optical flow as described in Section 3.3. We evaluated both approaches for an increased length of the sequence on different datasets. As shown in Table 1, the two-pass initialisation leads to quicker convergence and is most beneficial for the first optical flow estimation in the sequence.

**Table 1: Ablation study of ContinualFlow**. The leftmost column codes the experiment configurations: occlusion estimator (+OC); refinement network (+R); temporal connection with forward warping ($W_f$), backward warping ($W_b$) and both warping methods ($W_{bf}$); previous flow input in the refinement ($RW_x$); and two pass (2pass) initialisation of the first frame of the sequence N frames long. Performance measure are the KITTI 3-pixel error metric (column Fl) and the end-point error (in pixels, all other columns) for background (bg), foreground (fg), occluded (occ), non-occluded (noc) and all (all) pixels. The best performance in bold. All models trained on FlyingChairs and fine-tuned on FlyingThings. See section 3 for details.

| | FlyingThings | | | | | KITTI'15 noc | | KITTI'15 occ | | Sintel Clean | | | Sintel Final | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | all | occ-bg | occ-fg | noc-bg | noc-fg | Fl-all | all | Fl-all | all | all | occ | noc | all | occ | noc |
| common: baseline | **Occlusion map learning** | | | | | | | | | | | | | | |
| | 22.79 | 25.31 | 53.88 | 10.64 | 26.78 | 37.73 | 7.82 | 43.56 | 14.16 | 3.45 | 9.29 | 2.38 | 5.36 | 12.03 | 4.17 |
| +OC | **18.01** | **18.27** | **47.53** | **7.10** | **20.13** | **23.98** | **5.22** | **31.12** | **10.60** | **2.45** | **7.46** | **1.53** | **4.02** | **9.99** | **2.91** |
| common: baseline+OC | **The specialised refinement block** | | | | | | | | | | | | | | |
| | 18.01 | 18.27 | 47.53 | **7.10** | **20.13** | 23.98 | 5.22 | 31.12 | 10.60 | 2.45 | 7.46 | 1.53 | 4.02 | 9.99 | 2.91 |
| +R | **17.80** | **17.49** | **45.90** | 7.31 | 21.46 | **21.14** | **4.78** | **28.61** | **9.83** | **2.30** | **7.11** | **1.42** | **3.87** | **9.68** | **2.76** |
| common: baseline+OC | **Influence of coordinate warping methods** | | | | | | | | | | | | | | |
| | 18.01 | 18.27 | 47.53 | 7.10 | 20.13 | 23.98 | 5.22 | 31.12 | 10.60 | 2.45 | 7.46 | 1.53 | 4.02 | 9.99 | 2.91 |
| +$W_f$ | 14.90 | 14.89 | 38.75 | 6.55 | **16.69** | **20.78** | **4.13** | **27.85** | **8.28** | **2.18** | **6.67** | **1.37** | 4.04 | 9.48 | 3.03 |
| +$W_b$ | 16.33 | 17.10 | 39.68 | 6.49 | 20.72 | 26.52 | 4.56 | 33.80 | 10.64 | 2.58 | 7.49 | 1.70 | **3.79** | 9.27 | **2.80** |
| +$W_{bf}$ | **14.64** | **14.84** | **36.05** | **6.10** | 17.65 | 23.64 | 4.56 | 30.92 | 9.46 | 2.36 | 6.79 | 1.59 | 3.81 | **8.97** | 2.87 |
| common: baseline+OC | **Temporal connection placement** | | | | | | | | | | | | | | |
| +$RW_f$ | 16.10 | 15.87 | 38.71 | 6.76 | 19.01 | 23.11 | 4.88 | 30.35 | 9.82 | 2.27 | 6.89 | **1.45** | 3.92 | 9.34 | 2.90 |
| +$RW_{bf}$ | 14.90 | 15.35 | 37.55 | **5.78** | 17.69 | 24.54 | 4.84 | 32.18 | 10.12 | 2.35 | 6.93 | 1.54 | **3.55** | **8.62** | **2.65** |
| +$W_{bf}$ | 14.64 | 14.84 | 36.05 | 6.10 | 17.65 | 23.64 | 4.56 | 30.92 | 9.46 | 2.36 | 6.79 | 1.59 | 3.81 | 8.97 | 2.87 |
| +$W_{bf}$+ $RW_{bf}$ | **14.28** | **14.24** | **35.58** | 5.82 | **17.56** | **21.72** | **4.41** | **29.48** | **9.33** | **2.26** | **6.66** | 1.49 | 3.70 | 8.81 | 2.76 |
| common: baseline+OC+$W_{bf}$ | **Number of refinement blocks** | | | | | | | | | | | | | | |
| +1x$RW_{bf}$ | 14.28 | 14.24 | 35.58 | 5.82 | **17.56** | **21.72** | **4.41** | **29.48** | **9.33** | **2.26** | **6.71** | **1.47** | **3.76** | **8.93** | 2.80 |
| +2x$RW_{bf}$ | **14.26** | **14.13** | **35.60** | 5.78 | 17.62 | 21.77 | 4.45 | 29.62 | 9.35 | **2.26** | 6.72 | **1.47** | **3.76** | 8.96 | **2.79** |
| +3x$RW_{bf}$ | 14.30 | **14.13** | 35.71 | **5.75** | 17.77 | 21.98 | 4.50 | 29.86 | 9.40 | **2.26** | 6.74 | **1.47** | 3.77 | 8.99 | 2.80 |
| +5x$RW_{bf}$ | 14.43 | 14.24 | 36.16 | **5.75** | 17.93 | 22.48 | 4.58 | 30.35 | 9.49 | 2.28 | 6.80 | 1.48 | 3.80 | 9.03 | 2.83 |
| common: baseline+OC+$W_{bf}$+$RW_{bf}$ | **Multi-frame sequence initialisation** | | | | | | | | | | | | | | |
| 2 frames | - | - | - | - | - | 25.08 | 5.50 | 32.59 | 11.56 | 2.48 | 7.72 | 1.48 | 3.84 | 9.64 | 2.75 |
| 2 frames+2pass | - | - | - | - | - | **23.06** | **5.03** | **30.92** | **11.00** | **2.41** | **7.60** | **1.41** | **3.74** | **9.48** | **2.66** |
| 3 frames | - | - | - | - | - | 21.72 | 4.41 | 29.48 | 9.33 | 2.26 | 6.71 | **1.47** | 3.76 | 8.93 | 2.80 |
| 3 frames+2pass | - | - | - | - | - | **21.65** | **4.36** | **29.42** | **9.23** | **2.26** | **6.71** | 1.48 | **3.73** | **8.92** | **2.76** |
| 4 frames | - | - | - | - | - | **21.53** | 4.30 | **29.32** | 9.05 | **2.23** | 6.59 | 1.46 | 3.75 | 8.83 | 2.82 |
| 4 frames+2pass | - | - | - | - | - | 21.54 | **4.30** | 29.33 | **9.02** | 2.24 | **6.59** | **1.46** | **3.73** | **8.80** | **2.80** |
| 5 frames | - | - | - | - | - | **21.48** | 4.25 | **29.27** | **8.92** | **2.21** | **6.51** | **1.45** | 3.80 | 8.85 | 2.87 |
| 5 frames+2pass | - | - | - | - | - | **21.48** | **4.25** | 29.28 | **8.92** | **2.21** | 6.52 | 1.46 | **3.79** | **8.83** | **2.86** |
| 10 frames | - | - | - | - | - | 21.49 | **4.24** | 29.28 | 8.90 | - | - | - | - | - | - |
| 10 frames+2pass | - | - | - | - | - | **21.48** | **4.24** | **29.27** | **8.89** | - | - | - | - | - | - |

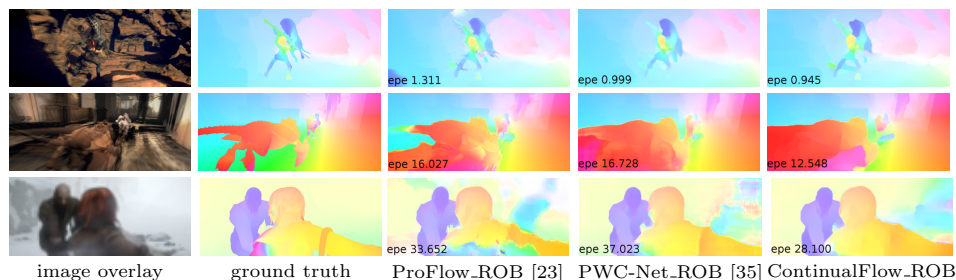| image overlay | ground truth | ProFlow_ROB [23] | PWC-Net_ROB [35] | ContinualFlow_ROB |

**Fig. 4:** Example results on Sintel Final-pass for ContinualFlow closest competitors in the Robust Vision Challenge. End-point-error for each method is shown for particular scenes.

### 4.2   Comparison with State of the Art

We start by noting that a single model was used for all benchmarks without further fine-tuning to individual datasets. We were not able to evaluate occlusions on public benchmarks since there is no benchmark available for occlusion map estimation. ContinualFlow achieves recall 0.87 and F1-score 0.83 for the validation split of FlyingThings [25] and recall 0.72 and F1-score 0.48 for Sintel [6]. Examples of estimated occlusion maps are shown in Fig 3.

**KITTI'15** optical flow benchmark [27] results are reported in Table 2. Fl refers to the KITTI evaluation metric – the percentage of pixels with end-point-error greater than 3 px. Our method ranked first among methods participating in the Robust Vision Challenge (ROB) and third for all optical flow estimation methods with score 10.03% on all evaluated pixels. We are interested in ROB Challenge since methods outside ROB fine-tune on each particular dataset, resulting in over-fitting, which we wanted to avoid.

**Sintel.** Fig 4 shows visual comparison with the closest competitors. Results of ROB participants on the Sintel dataset are reported in Table 3. ContinualFlow ranked first on Sintel Final for the all pixels end-point-error evaluation. As we are focused on occlusion estimation and handling, we point out that ContinualFlow achieves best results for estimation in occluded areas with significant margin.

**Robust Vision Challenge.** A snapshot of the leaderboard[4] of optical flow Robust Vision Challenge [32] is shown in Table 4. ContinualFlow is built on our implementation of PWC-Net [35]. While ContinualFlow did not achieve a better results in the ROB than the original PWC-Net, the experiments show that our contributions outperform the results of our baseline.

The source code for PWC-Net was released by the authors just days before the ACCV submission deadline, so a direct comparison was possible only

[4] As of July 7, 2018.

**Table 2:** KITTI'15 optical flow benchmark results of Robust Vision Challenge participants as of June 7, 2018. Performance measured by the KITTI 3-pixel error metric (column Fl) and the end-point error (in pixels, all other columns) for background (bg), foreground (fg), occluded (occ), non-occluded (noc) and all (all) pixels. The best results in bold. Anonymous entries in time of paper submission are marked [anon]. Methods are sorted according to Fl-all, the default ranking for KITTI.

| Fl (%) | KITTI'15 occ (%) | | | KITTI'15 noc (%) | | |
|---|---|---|---|---|---|---|
| | bg | fg | all | bg | fg | all |
| **ContinualFlow_ROB** | **8.54** | 17.48 | **10.03** | **5.90** | 14.99 | 7.55 |
| LFNet_ROB [anon] | 11.18 | 10.20 | 11.01 | 6.14 | 6.87 | **6.27** |
| PWC-Net_ROB [35] | 11.22 | 13.69 | 11.63 | 7.12 | 10.29 | 7.69 |
| ProFlow_ROB [23] | 14.15 | 21.82 | 15.42 | 8.44 | 17.90 | 10.15 |
| FF++_ROB [33] | 15.32 | 19.27 | 15.97 | 7.82 | 15.33 | 9.18 |
| ResPWCR_ROB [anon] | 16.63 | 16.18 | 16.55 | 10.10 | 12.23 | 10.49 |
| AugFNG_ROB [anon] | 19.77 | **9.95** | 18.14 | 13.75 | **6.71** | 12.47 |
| DMF_ROB [40] | 30.74 | 30.07 | 30.63 | 19.32 | 25.60 | 20.46 |

**Table 3:** Sintel benchmark results for Robust Vision Challenge participants. Performance measured the end-point error (EPE, in pixels) for matched (noc), unmatched (occ) and all (all) pixels. The best results in bold. Anonymous entries marked [anon]. Methods are sorted by EPE all, the default ranking for Sintel.

| | Sintel Final | | | Sintel Clean | | |
|---|---|---|---|---|---|---|
| | all | noc | occ | all | noc | occ |
| **ContinualFlow_ROB** | 4.528 | 2.723 | **19.248** | 3.341 | 1.752 | **16.292** |
| PWC-Net_ROB [35] | 4.903 | **2.454** | 24.878 | 3.897 | 1.726 | 21.637 |
| ProFlow_ROB [23] | 5.015 | 2.659 | 24.192 | **2.709** | **1.013** | 16.549 |
| AugFNG_ROB [anon] | 5.500 | 2.978 | 26.052 | 3.606 | 1.603 | 19.939 |
| LFNet_ROB [anon] | 5.966 | 3.278 | 27.893 | 4.815 | 2.333 | 25.065 |
| FF++_ROB [33] | 6.496 | 2.990 | 35.057 | 3.953 | 1.148 | 26.836 |
| ResPWCR_ROB [anon] | 6.530 | 3.849 | 28.371 | 5.674 | 3.138 | 26.380 |
| DMF_ROB [40] | 7.475 | 3.575 | 39.245 | 5.368 | 1.742 | 34.899 |

**Table 4:** Robust Vision Challenge. Performance measured by ranking of all metrics in individual datasets. The best results in bold. Anonymous entries marked [anon]. Methods are sorted by the Robust Vision Challenge rank.

|  | Middlebury | KITTI | MPI Sintel | HD1K |
|---|---|---|---|---|
| PWC-Net_ROB [35] | 2 | 4 | 2 | 1 |
| ProFlow_ROB [23] | 1 | 6 | 1 | 4 |
| **ContinualFlow_ROB** | 5 | 2 | 3 | 3 |
| LFNet_ROB [anon] | 7 | 1 | 6 | 5 |
| AugFNG_ROB [anon] | 9 | 3 | 4 | 2 |
| FF++_ROB [33] | 3 | 5 | 5 | 6 |
| DMF_ROB [40] | 4 | 8 | 7 | 8 |
| ResPWCR_ROB [anon] | 6 | 7 | 8 | 7 |
| WOLF_ROB [anon] | 8 | 9 | 9 | 9 |
| TVL1_ROB [30] | 10 | 10 | 10 | 10 |
| H+S_ROB [14] | 11 | 11 | 11 | 11 |

through ROB vision challenge submissions, which are limited in number by the challenge rules. We did our best to follow the paper regarding the architecture, parameters and training. Later, when analysing the results, we found two main differences: i) Due to implementation issues, we omitted rotation and scaling data augmentations, which in retrospect could harm the performance significantly as suggested in [24]. ii) Our implementation is in Tensorflow whereas the original implementation is in Caffe, so some of the suggested training parameter values may need to be fine-tuned for this framework. Still, the ablation study clearly shows the impact and significance of the novelties (occlusion estimation, feeding the previous estimate of optical flow as input).

## 5   Conclusion

The ContinualFlow network for optical flow estimation was introduced, with two novelties - occlusion estimation integrated in the optic flow computation and the use of the optic flow from the previous time instant, and, through recursion, of all prior flows. We showed that the two contributions improve performance, especially in occluded areas or areas close to motion discontinuities. In evaluation on standard dataset ContinualFlow is top ranked in Sintel and 3rd in KITTI.

## Acknowledgements

# References

1. Bailer, C., Taetz, B., Stricker, D.: Flow fields: Dense correspondence fields for highly accurate large displacement optical flow estimation. In: ICCV (2015)
2. Bailer, C., Varanasi, K., Stricker, D.: Cnn-based patch matching for optical flow with thresholded hinge embedding loss. In: CVPR. pp. 3250–3259 (2017)
3. Baker, S., Scharstein, D., Lewis, J.P., Roth, S., Black, M.J., Szeliski, R.: A database and evaluation methodology for optical flow. ICCV (2011). https://doi.org/10.1007/s11263-010-0390-2
4. Bengio, Y., Louradour, J., Collobert, R., Weston, J.: Curriculum learning. In: International conference on machine learning. pp. 41–48. ACM (2009)
5. Black, M.J., Anandan, P.: Robust dynamic motion estimation over time. In: CVPR (1991)
6. Butler, D.J., Wulff, J., Stanley, G.B., Black, M.J.: A naturalistic open source movie for optical flow evaluation. In: ECCV (2012)
7. Caelles, S., Maninis, K.K., Pont-Tuset, J., Leal-Taixé, L., Cremers, D., Van Gool, L.: One-shot video object segmentation. In: CVPR 2017. IEEE (2017)
8. Chen, Q., Koltun, V.: Full flow: Optical flow estimation by global optimization over regular grids. In: CVPR (2016)
9. Dosovitskiy, A., Fischer, P., Ilg, E., Häusser, P., Hazırbaş, C., Golkov, V., van der Smagt, P., Cremers, D., Brox, T.: Flownet: Learning optical flow with convolutional networks. In: ICCV. pp. 2758–2766 (Dec 2015)
10. Gaidon, A., Wang, Q., Cabon, Y., Vig, E.: Virtual worlds as proxy for multi-object tracking analysis. In: CVPR (2016)
11. Garg, R., Pizarro, L., Rueckert, D., Agapito, L.: Dense multi-frame optic flow for non-rigid objects using subspace constraints. In: ACCV (2010)
12. Garg, R., Roussos, A., Agapito, L.: A variational approach to video registration with subspace constraints. IJCV (3) (2013)
13. Güney, F., Geiger, A.: Deep discrete flow. In: ACCV (2016)
14. Horn, B.K., Schunck, B.G.: Determining optical flow. Artificial intelligence (1981)
15. Huang, G., Liu, Z., Weinberger, K.Q., van der Maaten, L.: Densely connected convolutional networks. CVPR (2017)
16. Hui, T.W., Tang, X., Loy, C.C.: LiteFlowNet: A Lightweight Convolutional Neural Network for Optical Flow Estimation. In: CVPR (June 2018)
17. Hur, J., Roth, S.: Mirrorflow: Exploiting symmetries in joint optical flow and occlusion estimation. ICCV (2017)
18. Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., Brox, T.: Flownet 2.0: Evolution of optical flow estimation with deep networks. CVPR (2017)
19. Irani, M.: Multi-frame correspondence estimation using subspace constraints. IJCV (2002)
20. Jaderberg, M., Simonyan, K., Zisserman, A., et al.: Spatial transformer networks. In: ANIPS (2015)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. ICLR (2015)
22. Kondermann, D., Nair, R., Honauer, K., Krispin, K., Andrulis, J., Brock, A., Gussefeld, B., Rahimimoghaddam, M., Hofmann, S., Brenner, C., Jahne, B.: The hci benchmark suite: Stereo and flow ground truth with uncertainties for urban autonomous driving. In: CVPR (2016)
23. Maurer, D., Bruhn, A.: Proflow: Learning to predict optical flow. BMVC (2018)
24. Mayer, N., Ilg, E., Fischer, P., Hazirbas, C., Cremers, D., Dosovitskiy, A., Brox, T.: What makes good synthetic training data for learning disparity and optical flow estimation? IJCV pp. 1–19 (2018)

25. Mayer, N., Ilg, E., Hausser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T.: A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In: CVPR (2016)
26. Meister, S., Hur, J., Roth, S.: Unflow: Unsupervised learning of optical flow with a bidirectional census loss. AAAI (2018)
27. Menze, M., Geiger, A.: Object scene flow for autonomous vehicles. In: CVPR (2015)
28. Murray, D.W., Buxton, B.F.: Scene segmentation from visual motion using global optimization. PAMI (1987)
29. Pang, J., Sun, W., Ren, J.S., Yang, C., Yan, Q.: Cascade residual learning: A two-stage convolutional neural network for stereo matching. CVPR (2017)
30. Pérez, J.S., Meinhardt-Llopis, E., Facciolo, G.: Tv-l1 optical flow estimation. Image Processing On Line (2013)
31. Ranjan, A., Black, M.J.: Optical flow estimation using a spatial pyramid network. CVPR (2018)
32. Robust Vision Challenge team: Robust vision challenge. `http://www.robustvision.net` (2018), [Online; accessed 8-July-2018]
33. Schuster, R., Bailer, C., Wasenmüller, O., Stricker, D.: Flowfields++: Accurate optical flow correspondences meet robust interpolation. ICIP (2018)
34. Sun, D., Liu, C., Pfister, H.: Local layering for joint motion estimation and occlusion detection. In: CVPR. pp. 1098–1105 (2014)
35. Sun, D., Yang, X., Liu, M.Y., Kautz, J.: PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. CVPR (2018)
36. Sundaram, N., Brox, T., Keutzer, K.: Dense point trajectories by gpu-accelerated large displacement optical flow. In: ECCV (2010)
37. Unger, M., Werlberger, M., Pock, T., Bischof, H.: Joint motion estimation and segmentation of complex scenes with label costs and occlusion modeling. In: CVPR. pp. 1878–1885. IEEE (2012)
38. Volz, S., Bruhn, A., Valgaerts, L., Zimmer, H.: Modeling temporal coherence for optical flow. In: ICCV (2011)
39. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing (2004)
40. Weinzaepfel, P., Revaud, J., Harchaoui, Z., Schmid, C.: Deepflow: Large displacement optical flow with deep matching. In: ICCV (2013)
41. Wulff, J., Sevilla-Lara, L., Black, M.J.: Optical flow in mostly rigid scenes. CVPR (2017)
42. Xiao, J., Cheng, H., Sawhney, H., Rao, C., Isnardi, M.: Bilateral filtering-based optical flow estimation with occlusion detection. In: ECCV (2006)
43. Xie, S., Tu, Z.: Holistically-nested edge detection. In: ICCV. pp. 1395–1403 (2015)