

Probabilistic Pixel-Adaptive Refinement Networks

Anne S. Wannewetsch^{1,2*} Stefan Roth²
¹Amazon, Germany ²TU Darmstadt, Germany

Abstract

Encoder-decoder networks have found widespread use in various dense prediction tasks. However, the strong reduction of spatial resolution in the encoder leads to a loss of location information as well as boundary artifacts. To address this, image-adaptive post-processing methods have shown beneficial by leveraging the high-resolution input image(s) as guidance data. We extend such approaches by considering an important orthogonal source of information: the network’s confidence in its own predictions. We introduce probabilistic pixel-adaptive convolutions (PPACs), which not only depend on image guidance data for filtering, but also respect the reliability of per-pixel predictions. As such, PPACs allow for image-adaptive smoothing and simultaneously propagating pixels of high confidence into less reliable regions, while respecting object boundaries. We demonstrate their utility in refinement networks for optical flow and semantic segmentation, where PPACs lead to a clear reduction in boundary artifacts. Moreover, our proposed refinement step is able to substantially improve the accuracy on various widely used benchmarks.

1. Introduction

Convolutional neural networks (CNNs) have become a standard tool in computer vision. Especially in dense prediction tasks [7, 41, 49, 60], encoder-decoder or pyramid-structured CNNs are a common choice. While originating in unsupervised learning [21], such architectures have become popular also in supervised settings. The encoder builds a powerful feature representation, reducing the spatial resolution of the inputs to aggregate global information [41]. The decoder takes the feature representation from the bottleneck, enlarges its size, and transforms it into the desired output, *e.g.* a segmentation map or optical flow field.

While downsampling in the encoder increases the receptive field and allows to deal with large image sizes, it also leads to a drastic loss in spatial resolution. As such, valuable location information is lost and boundary artifacts can

*This work was done at TU Darmstadt prior to Anne S. Wannewetsch joining Amazon.

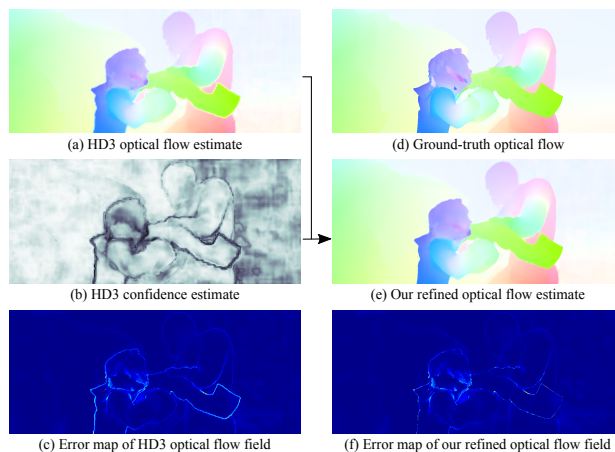


Figure 1. Our PPAC refinement method leverages the close relationship between estimated confidences and prediction errors to refine and improve the prediction itself, here the optical flow field.

arise [18, 41], *e.g.* segmentation maps that are misaligned w.r.t. the underlying objects. Moreover, the decoder typically yields low-resolution outputs and simple components are used to upscale predictions to the input size. This often results in blurry outputs since estimates of different objects are combined, *e.g.* motion from background and foreground objects is mixed as can be seen in Fig. 1.

Several approaches have been proposed to reduce these disadvantages, *e.g.* skip connections [41, 49] or densely connected blocks [24, 32]. Additionally, different types of generalized convolutions, taking into account a high-resolution RGB guidance image, have shown beneficial as part of the decoder [25, 26] or in a separate upsampling and/or refinement step [31, 48, 56, 58]. Many of the above approaches require a large number of additional parameters, are computationally expensive, or restrictive w.r.t. to the filtering method and the applicable guidance data. Most recently, pixel-adaptive convolutions (PACs) were introduced by Su *et al.* [51]. PACs combine spatially-invariant convolution weights with a content-adaptive kernel that depends on guidance data. In [51], PACs are shown to yield state-of-the-art results in joint upsampling tasks.

In this paper, we argue that we can leverage another source of information for refinement that is complementary

to the input image: *the uncertainty of each pixel’s estimate*. For dense classification networks, this quantity is generally provided implicitly. For instance, segmentation networks usually output log-probabilities of the classes, which are passed through an $\arg\max$ operation at test time. Even though these uncertainties might not be well calibrated [16], we argue that they still contain valuable information for refining the predictions. Beyond classification, explicit uncertainty estimates are also gaining increased attention for dense regression problems such as geometry [47] or motion estimation [12, 27, 60]. They are, for example, helpful for applications in which the reliability of network estimates is crucial, *e.g.* in autonomous driving. Here, we show that we can also leverage them to refine the regression output itself.

Fig. 1 shows an optical flow field estimated by the probabilistic HD3 method [60], as well as the corresponding confidence map and endpoint error per pixel. We observe that regions of high uncertainty (*b*, dark gray) correspond quite well with large errors (*c*, dark red). When applying post-processing to the network output, it seems desirable to take the available pixel uncertainty into account. As such, only reliable pixels should be spatially propagated while uncertain pixels can be replaced. To allow for probability-aware¹ filtering, we propose *probabilistic pixel-adaptive convolutions (PPACs)*, and therefore extend the adaptive convolution operation of [51]. The kernels of PPACs and thus the filtering output vary dependent on two properties: guidance data, *e.g.* the input image, as well as a probability map estimated by the deep network, either inherently or explicitly.

This paper focuses on the application of PPACs for the refinement of outputs from dense prediction networks, illustrated with the tasks of optical flow estimation and semantic segmentation. Therefore, we introduce a *PPAC refinement network*, which leverages RGB guidance data and probabilities for several content- and probability-adaptive convolutions. For both tasks, PPACs not only allow to improve estimates at boundaries but also remove outliers of low reliability. As shown in Fig. 1, blurry edges in the flow field (*a*) are transformed into crisp boundaries and the overall prediction is smoothed (*e*). Along with the visual improvements, PPAC refinement leads to a clear accuracy gain in optical flow and semantic segmentation. For instance, PPACs substantially improve state-of-the-art HD3 [60] optical flow estimates on the widely used KITTI 2012 and 2015 datasets. Our proposed PPAC-HD3 method ranks 1st among published optical flow approaches² on both benchmarks, improving the outlier rate by $\sim 11.1\%$ and $\sim 7.5\%$ over the underlying baseline method.

¹By *probability* we refer to a measure that approximates or summarizes the marginal posterior over the network’s estimates, *e.g.* by taking the marginal posterior of the chosen prediction value. To ease readability, the terms probability, confidence, (un)certainty, and reliability will be used interchangeably in the paper.

²All rankings at the time of publication.

2. Related Work

Probabilistic deep networks. Combining probabilistic approaches with deep networks is an active field of research, which is pursued to cope with model and/or input uncertainty [35]. As such, we can only provide a rough summary and refer to the cited references for a broader overview.

Bayesian neural networks [5, 15, 20, 37, 42, 57] often learn parametric distributions over the network weights to capture model uncertainty. The predictive distribution over the outputs is obtained by taking an expectation over the weights through approximate inference [5, 20]. However, Bayesian neural networks introduce many additional parameters and are not always easy to handle in practice [39].

Sampling-based approaches, *e.g.* [2, 11, 35], are often simpler to apply and include a random component, such as dropout, in the network structure. At test time, the predictive uncertainty is computed as Monte Carlo estimates from several network passes. Similarly, an ensemble of networks can be trained and combined at test time [23, 39]. A major drawback of both avenues is the increased runtime as they require multiple forward passes.

Another line of research uses deep networks to output the parameters of an assumed predictive distribution, either directly [35, 46] or by propagation of input uncertainty [12, 54]. There, the difficulty is to find a parametric distribution that is sufficiently easy to handle in practice and appropriately describes the quantity of interest.

Beyond such general purpose probabilistic treatments, probabilistic networks have also been developed in the context of specific vision problems. Yin *et al.* [60] propose a method to aggregate correspondence uncertainty in the context of optical flow and stereo matching through various spatial scales. In [27], a multi-hypothesis network for optical flow estimation is trained to output an ensemble at once. Novotny *et al.* [47] use uncertainty estimates obtained with probabilistic losses to predict the reliability of descriptors.

Content-adaptive convolutions. One category of content-adaptive convolutions adjusts the sampling location of neighbor pixels [8, 33]. Deformable convolutions [8] predict data-dependent offsets to determine at which locations neighboring pixels should be sampled for a spatially-invariant convolution. Another line of research adjusts the convolution weights [34, 59] of standard convolutions. Dynamic filter networks [34] use a subnetwork to predict location-specific weight kernels, which have already shown benefits for optical flow estimation, *e.g.* [25, 26]. A common drawback is the significant amount of additional parameters, which increases the risk of overfitting – especially if only a limited amount of training data is available.

Several works, therefore, approach content-adaptive convolutions in a more constrained setting. Spatial transformers [30] as well as CARAFE [55] rearrange features

in a content-adaptive way with a global or local transformation before performing the convolution itself. However, they remain restricted to a 2D grid structure. [31, 56] perform image-adaptive convolutions with predefined or learned features in the high-dimensional permutohedral lattice [1]. Such convolutions are computationally expensive and thus not suitable for fast processing. [18] incorporates semantics by learning input-dependent attention masks but requires object classes for (pre-)training. Deep guided filters [58] extend the classical guided filter [19] to learned guidance data. In [48], the parameters for spatially-variant linear representation models of the guided filter are learned with a CNN. In both approaches only 1D guidance data can be used for each output channel. Deep joint image filtering [40] applies standard convolutions to the concatenation of pre-processed guidance data and estimates, but misses explicit knowledge on the relation of both components.

We base our approach on pixel-adaptive convolutions (PACs) [51]. Here, the convolution kernels are split into a fixed weight as well as a location-specific component that depends upon a feature embedding learned from guidance data. PACs have a small computational overhead and are easy to train in practice. Nevertheless, like most adaptive filtering approaches, PACs do not allow to explicitly leverage knowledge about the reliability of filter inputs. While such information can be included as guidance data for the content-adaptive part of the weight kernel, our explicit probabilistic formulation leads to a significant performance gain.

Probabilistic joint filtering. There are only few filtering approaches that jointly consider guidance data as well as probabilities. Different filtering methods have been applied to refine semantic segmentations [50], optical flow [43, 53], and especially depth [17, 38, 45]. However, these approaches are task-specific, tailored to certain filtering methods, and/or rely on time-intensive iterative approaches. [14] replaces unreliable pixels using a network that takes uncertainties and guidance data as input, but uncertainties are not explicitly leveraged to improve estimates. The Fast Bilateral Solver [4] and the Domain Transform Solver [3] allow to perform fast, edge-aware optimization on different estimates. They leverage uncertainty by requiring a closer connection between inputs and outputs for more reliable pixels. As these approaches optimize a predefined objective, their flexibility is restricted. Moreover, [3, 4] cannot backpropagate into the features used to determine the pixel similarity. Closest to ours with regard to probabilistic joint filtering are [22, 29], where confidences are used to extend the bilateral and the guided filter, respectively, by weighing a pixel’s importance with its confidence. However, both methods rely on predefined features for pixel similarity, hand-crafted reliability measures, and fixed filter kernels.

In comparison to previous work, our approach is very general as pixel feature embeddings, the filter weights, as

well as a pre-processing of the probabilities are learned from data. Moreover, the proposed approach is fast and easily integrable into different task-specific neural networks.

3. Probabilistic Pixel-Adaptive Convolutions

We begin by presenting pixel-adaptive convolutions (PACs) as introduced in [51]. We then propose an advanced normalization approach for pixel-adaptive convolutions and finally extend PACs to allow for probabilistic filtering.

3.1. Pixel-adaptive convolutions

Assume, we aim to perform a convolution with neighborhood size s that transforms features $\mathbf{v} \in \mathbb{R}^d$ into features $\tilde{\mathbf{v}} \in \mathbb{R}^{d'}$. We denote the corresponding convolution weights as tensor $\mathbf{W} \in \mathbb{R}^{d' \times d \times s \times s}$ and the bias term as $\mathbf{b} \in \mathbb{R}^{d'}$. Following the notation of [51], the output of a standard convolution at pixel i is then given as

$$\tilde{\mathbf{v}}_i = \sum_{j \in \mathcal{N}(i)} \mathbf{W}[\mathbf{p}_i - \mathbf{p}_j] \mathbf{v}_j + \mathbf{b}, \quad (1)$$

where $\mathcal{N}(i)$ denotes the $s \times s$ neighborhood of the pixel. The vectors \mathbf{p}_i and \mathbf{p}_j represent the 2D pixel positions. $\mathbf{W}[\mathbf{p}_i - \mathbf{p}_j] \in \mathbb{R}^{d' \times d}$ corresponds to the 2D slice from weight tensor \mathbf{W} , evaluated at position $\mathbf{p}_i - \mathbf{p}_j$.

PACs generalize such spatially-invariant convolutions by augmenting the convolution weight with an additional location-adaptive component $K(\mathbf{f}_i, \mathbf{f}_j)$. In [51], the vectors \mathbf{f}_i and \mathbf{f}_j are denoted as *pixel features* and characterize the pixels i and j , respectively. For instance, one could use the RGB components of a guidance image as feature $\mathbf{f}_{(\cdot)}$, as is done in the bilateral filter [52]. However, more advanced features learned from data have shown to be advantageous [51]. The function $K(\mathbf{f}_i, \mathbf{f}_j) = K(\mathbf{f}_i - \mathbf{f}_j)$ is a (fixed) kernel, which evaluates the difference between \mathbf{f}_i and \mathbf{f}_j . If pixels i and j show similar characteristics, $K(\mathbf{f}_i, \mathbf{f}_j)$ weighs the corresponding values \mathbf{v}_j more than the ones of a more deviating pixel. Various choices for $K(\cdot, \cdot)$ are possible; we will apply a Gaussian RBF kernel in the following, *i.e.*

$$K(\mathbf{f}_i, \mathbf{f}_j) = e^{-\frac{1}{2}(\mathbf{f}_i - \mathbf{f}_j)^T(\mathbf{f}_i - \mathbf{f}_j)}. \quad (2)$$

The PAC convolution [51] is then defined as

$$\tilde{\mathbf{v}}_i = \sum_{j \in \mathcal{N}(i)} K(\mathbf{f}_i, \mathbf{f}_j) \cdot \mathbf{W}[\mathbf{p}_i - \mathbf{p}_j] \mathbf{v}_j + \mathbf{b}. \quad (3)$$

The same feature kernel K is used for all input channels, while the weight \mathbf{W} differs dependent on the spatial location within the mask and for each feature channel.

3.2. Advanced normalization step

PACs are a powerful tool for deep dense prediction architectures, but they are not without challenges. One major issue is the fact that the number of closely related pixels, *i.e.* pixels j that show a high value of $K(\mathbf{f}_i, \mathbf{f}_j)$, varies

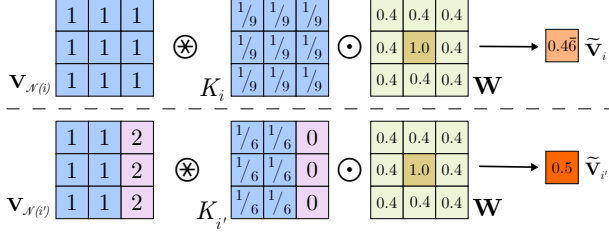


Figure 2. PAC normalization w.r.t. kernel only (see text).

across the image. This is natural and even desirable, since only a restricted neighborhood region should be taken into account at object boundaries or similar. Nevertheless, applying PACs in different neighborhoods should lead to results with the same output magnitude as long as the input variables are equal. Otherwise, learning of convolution weights might be difficult since the output values are inevitably smaller at boundaries or in highly structured areas.

Basic scheme. The implementation of PACs provides an option to normalize kernels such that $\sum_j K(\mathbf{f}_i, \mathbf{f}_j) = 1$ for all pixels i . However, we argue that such a *kernel normalization* is not sufficient. Consider the illustration of two pixels i and i' and their neighborhoods $\mathcal{N}(i)$ and $\mathcal{N}(i')$ in Fig. 2. For simplicity, we assume equal input values for pixels of the same object. Pixel i is part of a homogeneous area and the kernel function leads to an equal distribution of the kernel weights. In contrast, the neighborhood of pixel i' contains also elements from a different object. Here, the kernel weights are only distributed over the elements from the same object as pixel i' . Even though both kernels sum to one, their convolution with an exemplary weight \mathbf{W} leads to clearly different results. As mentioned above, this complicates the learning of PACs.

Our advanced scheme. We address this issue with an advanced normalization scheme. To that end, we adapt the normalization from [56] to the context of PACs. An auxiliary array $\mathbf{v}^{aux} = \mathbf{1}$ with constant value 1 is defined and passed through the filtering step in Eq. (3):

$$\tilde{\mathbf{v}}_i^{aux} = \sum_{j \in \mathcal{N}(i)} K(\mathbf{f}_i, \mathbf{f}_j) \cdot \mathbf{W} [\mathbf{p}_i - \mathbf{p}_j] \cdot \mathbf{1}. \quad (4)$$

In slight abuse of notation, we now denote by $\tilde{\mathbf{v}}_i$ the PAC output in Eq. (3) *before* adding the bias term. Then, the normalized convolution output $\tilde{\mathbf{v}}_{i, \text{norm}}$ is given by

$$\tilde{\mathbf{v}}_{i, \text{norm}} = \tilde{\mathbf{v}}_i / \tilde{\mathbf{v}}_i^{aux} + \mathbf{b}. \quad (5)$$

With the novel normalization, not only the number of similar pixels is taken into account but also their weighting by \mathbf{W} . However, the normalization becomes invalid as soon as the weight \mathbf{W} becomes negative [56]. We thus follow [56] to introduce an additional *normalization weight* \mathbf{W}' , which

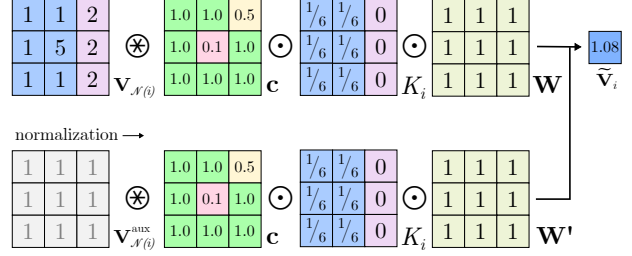


Figure 3. Illustration of our proposed probabilistic PAC approach. The weighting of estimates according to their pixel feature similarity K as well as their confidence \mathbf{c} allows to remove outliers.

replaces \mathbf{W} in the normalization convolution in Eq. (4). \mathbf{W}' is ensured to remain positive and set to the same initialization as \mathbf{W} . During training, we independently update \mathbf{W} as well as \mathbf{W}' using regular gradient-based optimization and thus omit to explicitly enforce their similarity.

Reconsidering the example in Fig. 2, both output values remain close when using our advanced normalization with the appropriate $\mathbf{W}' \approx \mathbf{W}$. We will show in Sec. 6.1 that the proposed scheme leads to superior results in practice.

3.3. Probabilistic pixel-adaptive convolutions

While the definition of PACs allows for more advanced filtering than a standard convolution, the approach is still restricted. The kernel function $K(\cdot, \cdot)$ in Eq. (2) only takes differences of pixel features as input, thus excluding properties that cannot be reasonably expressed as such. Consider, for instance, that we have a per-pixel probability alongside the estimate. In this case, it seems beneficial to consider such information during filtering to improve unreliable estimates. As $K(\cdot, \cdot)$ rewards similar pixel features, neighbors with a similar level of reliability are more closely connected. However, this seems counterintuitive, given that *the values of reliable pixels should particularly propagate to neighbors with a very different, i.e. low, confidence*.

Therefore, we extend PACs such that we are able to perform convolutions that also take unary properties – especially probabilities – into account. Let c_j describe the confidence assigned to a certain pixel location j . For consistency, we assume that only one confidence estimate is given per spatial location.³ Similar to [29], we then propose to define a *probabilistic pixel-adaptive convolution (PPAC)* as:

$$\tilde{\mathbf{v}}_i = \sum_{j \in \mathcal{N}(i)} c_j \cdot K(\mathbf{f}_i, \mathbf{f}_j) \cdot \mathbf{W} [\mathbf{p}_i - \mathbf{p}_j] \mathbf{v}_j + \mathbf{b}. \quad (6)$$

Here, each pixel value is not only weighted by its distance to the center pixel but also with its individual confidence.

To illustrate our proposed approach, consider Fig. 3: An outlier pixel is surrounded by more reliable estimates,

³An extension to individual confidences per channel is straightforward.

which belong to the same and a different object. For simplicity, we assume that \mathbf{W} performs an averaging operation. If the pixel confidence was not taken into account, the outlier value would spread to the surrounding pixels due to its higher magnitude. In contrast, the proposed PPAC allows to propagate the reliable pixel values from the same object and thus almost completely replaces the outlier in the center.

The normalization as proposed in Sec. 3.2 can be easily extended to PPACs. To that end, $\tilde{\mathbf{v}}_i^{aux}$ is obtained as

$$\tilde{\mathbf{v}}_i^{aux} = \sum_{j \in \mathcal{N}(i)} c_j \cdot K(\mathbf{f}_i, \mathbf{f}_j) \cdot \mathbf{W}[\mathbf{p}_i - \mathbf{p}_j] \cdot \mathbf{1}. \quad (7)$$

The normalization step is performed as before and the PPAC outputs are divided per pixel by $\tilde{\mathbf{v}}_i^{aux}$ before the bias (Eq. 5)

4. Refinement Networks with PPACs

Deep dense predictors often output predictions at a scale lower than the input resolution to save time and parameters, e.g. [7, 60]. The low-resolution estimates are then upsampled by simple methods such as bilinear interpolation. In [51, 56, 58], image-adaptive convolutions have proven very helpful in this upsampling step. Following that, we propose a *PPAC refinement network*, which takes image and reliability data into account to upscale and refine network outputs.

A straightforward approach is to upscale results with transposed PPACs. However, we found that this can lead to difficulties, especially for optical flow. This is due to the fact that flow networks often assume the input sizes to be divisible by a certain power of 2, e.g. 2^6 [60]. As this is mostly not the case, e.g. for the Sintel benchmark [6], input images are resized and the output flow is afterwards rescaled with a non-integer factor. To apply transposed convolutions, which can only upscale by integers, one has to pad the inputs and crop the output after upscaling. We observed that this leads to severe artifacts, which clearly reduce the accuracy.

Instead, we propose to first upscale the estimates by the default method of the original network. A lightweight network with PPACs is then applied at full resolution. As we only use a small number of PPACs, the computational expense of the approach remains low and the prediction accuracy does not decrease due to padding artifacts or similar.

Our proposed refinement networks consist of three branches as illustrated in Fig. 4. In addition to the upsampled estimates, the network takes corresponding probability data, e.g. a full marginal posterior per pixel or other probability measures, and the high-resolution images as inputs. The first subnetwork transforms the probability data into scalar confidence values for the individual PPACs. The second branch processes the guidance images to generate meaningful pixel features. Both intermediate outputs as well as the underlying network predictions are then fed into the PPACs of the combination branch to create a refined estimate.

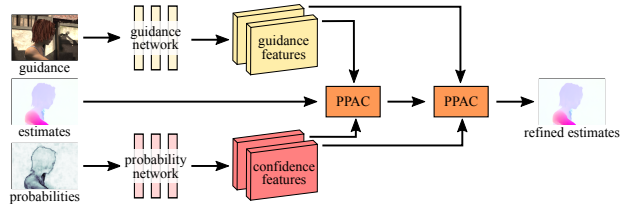


Figure 4. Exemplary architecture of our proposed PPAC refinement networks.

5. Implementation

5.1. Network architecture

To illustrate their capabilities, we experiment with PPAC refinement networks for the tasks of optical flow estimation and semantic segmentation. All networks are fully convolutional and include two consecutive PPACs. We use log values as inputs to the probability branch and add a sigmoid function at its end for normalization. Please refer to the supplement for a detailed description of the networks' setup.

Here, we only highlight architectural choices that differ significantly from the ones in [51]. First, we found no benefit from increasing the number of channels in the combination branch. Moreover, using group convolutions with the number of groups being equal to the number of inputs did not decrease the performance, but significantly lowers the number of parameters. We even go further and share the convolution weights across all channels, which proved especially beneficial for semantic segmentation, possibly due to the large reduction in parameters. Following the upscaling setup in [51], we also experimented with standard convolution layers to pre- or post-process the data itself. However, we found no benefit from such convolutions and thus stick with a combination branch that only includes PPACs.

5.2. Additional baselines

To assess the benefits of our PPAC refinement networks, we introduce two baseline networks. The *simple refinement network* takes estimates, log-probabilities, and input images and processes them jointly with standard convolutions. Here, we set the channel depth such that the number of parameters used for the PPAC network is approximately equal to the simple baseline. Additionally, we use a PAC baseline that replaces all PPACs with its non-probabilistic PAC counterpart. For this network, the probability branch is removed and the probabilities are instead concatenated with the guidance images and fed to the guidance subnetwork. We again ensure that the number of parameters remains comparable for PAC and PPAC refinement networks.

5.3. Training procedure

For fair comparison, we only train the refinement networks and do not backpropagate into the original networks

Table 1. Average end-point error (*AEE*) and 3-pixel outlier rate (*out*) on our Sintel and KITTI test splits for different normalizations of a PAC refinement network.

| | Sintel (AAE) | | KITTI | |
|-------------------------------------|--------------|--------------|--------------|---------------|
| | clean | final | AEE | out |
| HD3 [60] | 1.672 | 1.357 | 1.990 | 6.14 % |
| PAC w/o normal. | 1.665 | 1.352 | 1.924 | 6.34 % |
| PAC w/ kernel normal. | 1.622 | 1.323 | 1.921 | 6.15 % |
| PAC w/ adv. normal. (<i>ours</i>) | 1.594 | 1.302 | 1.868 | 5.81 % |

themselves. However, our proposed probabilistic refinement is also easily applicable in fully end-to-end training. For networks with PACs or PPACs, we apply our advanced normalization from Sec. 3.2. Here, we found it important to initialize weights \mathbf{W} and \mathbf{W}' with the same values. We thus initialize both with positive, random numbers and ensure that \mathbf{W}' remains positive by learning in log-space.

As we aim to compare several different approaches, carrying out ablations on the official test datasets is not feasible. Thus, we split the data with available ground truth randomly into custom training, validation, and test sets. Learning rates for the optimization with Adam [36] are determined for all networks individually on the validation set. Please see Sec. 6 as well as the supplement for further training specifics. Our PyTorch code is publicly available.⁴

6. Experiments

6.1. Optical flow

We first apply our probabilistic refinement networks to the task of optical flow estimation. As underlying network, we use the state-of-the-art HD3 method [60], which yields competitive results on the major benchmarks. HD3 predicts flow in a residual fashion and estimates a discrete probability distribution at each scale. In [60], the full (discrete) matching distribution of the flow is composed, which is time- and memory-consuming. We instead upsample all probability maps via bilinear interpolation and provide the network with the probability value of the respective flow residual at all five scales. Since the residuals are subpixel-accurate and mostly fall outside of the discrete grid, we use a nearest neighbor interpolation of the probabilities.

We evaluate the proposed refinement networks on two widely used optical flow benchmarks: Sintel [6] and KITTI [13, 44]. The Sintel data is split into 862 images for training, 80 for validation, and 99 for test. Moreover, we merge the KITTI 2012 and 2015 images to obtain 319 training, 31 validation, and 44 test images. Using the procedure described in [60], we fine-tune two individual HD3 models on our own Sintel and KITTI training splits. We initialize the

Table 2. Average-end-point error (*AEE*) and 3-pixel outlier rate (*out*) on our Sintel and KITTI test splits for different refinements.

| | Sintel (AAE) | | KITTI | |
|------------------------------|--------------|--------------|--------------|---------------|
| | clean | final | AEE | out |
| HD3 [60] | 1.672 | 1.357 | 1.990 | 6.14 % |
| Simple refinement | 1.638 | 1.334 | 1.872 | 5.92 % |
| PAC network [51] | 1.594 | 1.302 | 1.868 | 5.81 % |
| PPAC network (<i>ours</i>) | 1.562 | 1.283 | 1.848 | 5.50 % |
| Oracle network | <i>1.430</i> | <i>1.149</i> | <i>1.500</i> | <i>4.48 %</i> |
| SL [56] | 1.634 | 1.340 | 1.953 | 6.41 % |
| FBS [4] | 1.643 | 1.354 | – | – |

networks from the author-provided checkpoint pre-trained on FlyingChairs [9] as well as FlyingThings3D [28] and determine the best models on our validation images. Following [60], only images from Sintel final are used for fine-tuning of the Sintel model.

All refinement networks are trained for 500 epochs with the average end-point error (AEE) as loss function and a batch size of 8. The base learning rates are cut by a factor of two every 100 epochs. As augmentations, we only apply random cropping to sizes (384, 768) and (320, 896) for Sintel and KITTI data, respectively. On Sintel, individual networks are trained on the final and clean subsets.

Comparison of normalization approaches. We first demonstrate the benefit of our proposed advanced normalization scheme. Therefore, we train PAC refinement networks without normalization, similar to [51], and with kernel normalization, *c.f.* Sec. 3.2. Table 1 shows results evaluated on our test sets of Sintel and KITTI. Note that the results on Sintel clean tend to be worse than on final as this data has not been seen during HD3 fine-tuning. We observe that a PAC network with kernel normalization is able to improve the accuracy over the HD3 baseline as well as the PAC approach without normalization. However, the same network leveraging our proposed advanced normalization shows clearly the best accuracy on all test sets. We attribute this to the fact that kernel normalization does not sufficiently compensate different neighborhood conditions.

Evaluation of PPACs. We now evaluate refinement networks including PPACs in comparison to the simple and PAC baselines as introduced in Sec. 5.2. The results of flow refinement on our test splits are summarized in Table 2. The simple refinement approach based only on standard convolutions (with the same inputs) improves the flow predictions only slightly on all datasets. Applying guidance data, including the estimated probabilities, explicitly in a PAC refinement network already leads to a clear improvement. However, our proposed PPAC refinement network outperforms the PAC approach by a large margin, improving the AEE by 6.6%, 5.5% and 7.1% on Sintel clean, Sintel final

⁴https://github.com/visinf/ppac_refinement



Figure 5. Examples of ground truth (*top*), HD3 optical flow (*middle*), and PPAC-refined optical flow (*bottom*) on KITTI. *Best viewed on screen.*

and KITTI, respectively. We also observe that the improvement of using content-adaptive, probabilistic convolutions over a simple setup with standard convolutions is more significant on Sintel than on KITTI. We attribute this to the fact that guidance data has shown to be most effective in boundary regions [56], which play a lesser role in the KITTI dataset as the ground truth is sparse.

Even more striking than these significant accuracy gains, are the improvements in visual quality. Figs. 1 and 5 show example flow fields on Sintel final and KITTI. PPACs clearly improve the underlying HD3 estimates and lead to substantial improvements especially near motion boundaries. Additionally, our proposed approach is able to correctly propagate flow estimates into outlier regions.

To further understand the results of PPAC refinement, we evaluate an additional *oracle network*, which takes oracle confidences as input to the probability branch, which we take to be the inverse of the AEE. As this correctly assesses the reliability of each pixel, this network provides an upper bound on the possible accuracy improvement from probabilistic refinement. Comparing the results in Table 2, we observe that an even more significant improvement would be possible if more precise probability estimates were available. This observation holds especially for the evaluation on KITTI, where a rather small amount of ground truth is available during fine-tuning. This suggests that future work on improved probability estimates in deep network has the potential of improving the accuracy in difficult areas further.

We finally compare our proposed PPAC refinement to other approaches from the literature. As PACs have shown to be the state-of-the-art method for joint upsampling, we restrict further comparison to the Semantic Lattice (SL) [56], which appeared concurrently to [51], and the Fast Bilateral Solver (FBS) [4] as a representative method that explicitly considers confidences for post-processing. The SL is trained as described in [56]. For the FBS, we use the probabilities of the last output layer of HD3 [60]. All FBS weight parameters and the trade-off parameter λ are determined via grid search on the validation set. Note that we were not able to find stable parameters for FBS on KITTI.

Table 3. Average-end-point error (AEE) and AEE of regions ≤ 10 pixels from occlusion boundaries (*d0-10*) evaluated on Sintel test.

| | clean | | final | |
|--------------------------|--------------|--------------|--------------|--------------|
| | AEE | d0-10 | AEE | d0-10 |
| HD3 [60] | 4.788 | 3.225 | 4.666 | 3.786 |
| PPAC-HD3 (<i>ours</i>) | 4.589 | 2.788 | 4.599 | 3.521 |

Table 4. 3-pixel outlier rate of non-occluded/all pixels (*Out-Noc/all*) and runtimes evaluated on KITTI 2012 and 2015 test.

| | KITTI 2012 | | KITTI 2015 | |
|--------------------------|---------------|---------------|---------------|---------------|
| | Out-Noc | Out-all | Out-all | Runtime |
| HD3 [60] | 2.26 % | 5.41 % | 6.55 % | 0.11 s |
| PPAC-HD3 (<i>ours</i>) | 2.01 % | 5.09 % | 6.06 % | 0.19 s |

Table 2 shows that SL and FBS are both able to improve the HD3 baseline accuracy. However, PPAC refinement outperforms both previous methods by a large margin.

Evaluation on benchmarks. We finish our flow experiments by evaluating PPAC-HD3, *i.e.* HD3 optical flow [60] with PPAC refinement, on the official benchmarks. For Sintel, we initialize HD3 with the same checkpoint as used in [60]. On KITTI, we use the fine-tuned checkpoint with context module as provided online. To train our PPAC network, we leverage the entire training data provided for the Sintel and KITTI benchmarks, respectively. Again, we train separate networks for Sintel clean and final, as well as a joint refinement network for both KITTI benchmarks. In comparison to our previous experiment, we adjust the number of epochs in the learning scheme such that the total number of iterations remains approximately the same despite the larger number of images. Note that we do not use other augmentations than random cropping, since we found our refinement network to be robust w.r.t. overfitting. We attribute this to the lightweight structure of our PPAC network, which only adds approximately 12k parameters.

Results on Sintel test are summarized in Table 3. PPAC-HD3 clearly improves the accuracy of the underlying HD3 method on clean and final splits by $\sim 4.2\%$ and $\sim 1.4\%$, respectively. The larger improvement on clean might be partly due to the fact that no data from this pass was used during fine-tuning of the HD3 baseline. Moreover, we observe a very significant improvement of 7.0% on the final pass and of $\sim 13.6\%$ on Sintel clean when considering the AEE of regions closer than 10 pixels to motion boundaries (*d0-10*). Overall, PPAC-HD3 ranks 4th on Sintel final among published two-frame methods and is thus highly competitive.

In Table 4, we show results on the official test sets of KITTI 2012 and 2015. PPAC-HD3 outperforms its underlying method by a large margin, leading to a substantial relative improvement of $\sim 11.1\%$ for the outlier rate of non-

Table 5. Mean intersection over union (*mIoU*) evaluated on our Pascal VOC 2012 test split. †Results taken from [56].

| | mIoU |
|---------------------------------|----------------|
| DeepLabv3+ [7] | 82.20 % |
| PAC refinement [51] | 82.39 % |
| PPAC refinement (<i>ours</i>) | 82.62 % |
| SL [56]† | 82.25 % |
| FBS [4] | 82.28 % |

occluded pixels on KITTI 2012. On KITTI 2015, PPAC refinement improves the outlier rate of all pixels from 6.55% to 6.06% (relative improvement of $\sim 7.5\%$). Our proposed method clearly ranks 1st among published optical flow approaches on both datasets. PPAC-HD3 is even able to outperform several strong scene flow methods, which leverage additional stereo data as input. Table 4 also shows computational times measured on a single GTX 1080 Ti GPU. Adding PPACs has a computational overhead of $\sim 1.7\times$, which seems justified by the strongly improved results.

6.2. Semantic segmentation

In a second set of experiments, we apply our probabilistic refinement networks to the task of semantic segmentation. We choose DeepLabv3+ [7] as a baseline using the Xception65 variant of the model. Before feeding the logits of the segmentation network into the probability branch of our PPAC network, we normalize them with a log-softmax operation. The logits are equally used as input to the combination branch. Here, we left the values unnormalized as we found a normalization to lead to inferior results.

We evaluate our probabilistic refinement network on Pascal VOC 2012 [10] and use the training, validation, and test split of [56]. When training PAC and PPAC networks, we found it important to use different learning rates for the pre-processing branches as well as the weights in the combination branch. To determine appropriate values, we first fixed the PPAC or PAC weights to a Gaussian kernel with standard deviation $\sigma = 1$ and searched for the optimal learning rate of the guidance and probability subnetworks on the validation set. In a second step, all weights are randomly initialized and a second learning rate is determined for the convolution weights in the combination branch. Furthermore, we observed improved accuracy when initializing the bias term of PACs and PPACs to zero. All networks are trained for 500 epochs with constant learning rate, a batch-size of 16, and random image crops with size 200×272 . We use the cross entropy as loss function and evaluate the mean intersection over union for validation. Please see the supplement for details, *e.g.* the network architectures.

Table 5 summarizes the results on our test split of Pascal VOC 2012. In this setting, our PPAC network requires on average 0.055s per image on a single GTX 1080 Ti GPU and

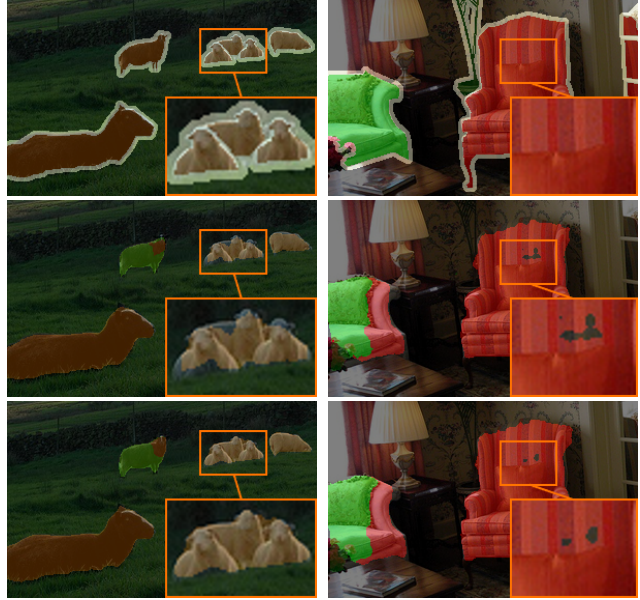


Figure 6. Cropped examples of ground truth (*top*), DeepLabv3+ (*middle*), and PPAC refined segmentation maps (*bottom*) on Pascal VOC 2012. *Best viewed on screen.*

is able to improve the segmentation accuracy even though DeepLabv3+ takes already special care of the decoder. In contrast, PAC refinement shows a considerably smaller benefit. We were not able to find a simple configuration that improved the original results. As for optical flow, we also compare to SL [56] and FBS [4], and use the probability of the most likely class as confidence input. Both previous methods are clearly outperformed by our proposed PPACs.

Fig. 6 shows examples of segmentation maps from Pascal VOC 2012. PPAC refinement leads to a better alignment with the underlying objects especially at object intersections or for smaller objects. Moreover, PPACs are able to successfully close smaller holes in the segmentation maps.

7. Conclusion

We introduced probabilistic pixel-adaptive convolutions (PPACs), which allow for filtering operations that respect guidance data as well as per-pixel confidences. Building on the work of [51], we first proposed an advanced normalization scheme, which we show to clearly improve the results in practice. Subsequently, we extend PACs to include confidence information during the filtering step to especially improve regions of low reliability. We proposed to use PPACs for the refinement of dense prediction networks and demonstrated their benefits for optical flow estimation and semantic segmentation. Here, PPAC refinement resulted in significant accuracy gains; our PPAC-HD3 clearly leads both KITTI benchmarks for optical flow. Moreover, refined estimates show fewer boundary artifacts and are smoother overall while correctly respecting object boundaries.

References

- [1] Andrew Adams, Jongmin Baek, and Myers Abraham Davis. Fast high-dimensional filtering using the permutohedral lattice. *Comput. Graph. Forum*, 29(2):753–762, 2010.
- [2] Murat Seçkin Ayhan and Philipp Berens. Test-time data augmentation for estimation of heteroscedastic aleatoric uncertainty in deep neural networks. In *International Conference on Medical Imaging with Deep Learning*, 2018.
- [3] Akash Bapat and Jan-Michael Frahm. The domain transform solver. In *CVPR*, pages 6014–6023, 2019.
- [4] Jonathan T. Barron and Ben Poole. The fast bilateral solver. In *ECCV*, volume 3, pages 617–632, 2016.
- [5] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. In *ICML*, pages 1613–1622, 2015.
- [6] Daniel J. Butler, Jonas Wulff, Garrett B. Stanley, and Michael J. Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, volume 4, pages 611–625, 2012.
- [7] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, volume 7, pages 833–851, 2018.
- [8] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017.
- [9] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Häusser, Caner Hazırbaş, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, pages 2758–2766, 2015.
- [10] Mark Everingham, S. M. Ali Eslami, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL visual object classes challenge: A retrospective. *Int. J. Comput. Vision*, 111(1):98–136, 2015.
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *ICML*, pages 1050–1059, 2016.
- [12] Jochen Gast and Stefan Roth. Lightweight probabilistic deep networks. In *CVPR*, pages 3369–3378, 2018.
- [13] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *CVPR*, pages 3354–3361, 2012.
- [14] Spyros Gidaris and Nikos Komodakis. Detect, replace, refine: Deep structured prediction for pixel wise labeling. In *CVPR*, pages 5248–5257, 2017.
- [15] Alex Graves. Practical variational inference for neural networks. In *NIPS*2011*, pages 2348–2356.
- [16] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *ICML*, pages 1321–1330, 2017.
- [17] Bumsu Ham, Minsu Cho, and Jean Ponce. Robust guided image filtering using nonconvex potentials. *IEEE T. Pattern Anal. Mach. Intell.*, 40(1):192–207, 2018.
- [18] Adam W. Harley, Konstantinos G. Derpanis, and Iasonas Kokkinos. Segmentation-aware convolutional networks using local attention masks. In *ICCV*, pages 5038–5047, 2017.
- [19] Kaiming He, Jian Sun, and Xiaoou Tang. Guided image filtering. In *ECCV*, volume 1, pages 1–14, 2010.
- [20] José Miguel Hernández-Lobato and Ryan P. Adams. Probabilistic backpropagation for scalable learning of Bayesian neural networks. In *ICML*, pages 1861–1869, 2015.
- [21] Geoffrey E. Hinton and Richard S. Zemel. Autoencoders, minimum description length and Helmholtz free energy. In *NIPS*1993*, pages 3–10.
- [22] Jobst Hörentrup and Markus Schlosser. Confidence-aware guided image filter. In *ICIP*, pages 3243–3247, 2014.
- [23] Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E. Hopcroft, and Kilian Q. Weinberger. Snapshot ensembles: Train 1, get M for free. In *ICLR*, 2017.
- [24] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *CVPR*, pages 770–778, 2017.
- [25] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. LiteFlowNet: A lightweight convolutional neural network for optical flow estimation. In *CVPR*, pages 8981–8989, 2018.
- [26] Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *CVPR*, pages 5754–5763, 2019.
- [27] Eddy Ilg, Ozgun Cicek, Silvio Galesso, Aaron Klein, Osama Makansi, Frank Hutter, and Thomas Brox. Uncertainty estimates and multi-hypotheses networks for optical flow. In *ECCV*, pages 652–667, 2018.
- [28] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, pages 1647–1655, 2017.
- [29] Jörn Jachalsky, Markus Schlosser, and Dirk Gandolph. Confidence evaluation for robust, fast-converging disparity map refinement. In *ICME*, pages 1399–1404, 2010.
- [30] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *NIPS*2015*, pages 2017–2025.
- [31] Varun Jampani, Martin Kiefel, and Peter V. Gehler. Learning sparse high dimensional filters: Image filtering, dense CRFs and bilateral neural networks. In *CVPR*, pages 4452–4461, 2016.
- [32] Simon Jégou, Michal Drozdal, David Vazquez, Adriana Romero, and Yoshua Bengio. The one hundred layers Tiramisu: Fully convolutional DenseNets for semantic segmentation. In *CVPR Workshops*, pages 1175–1183, 2017.
- [33] Yunho Jeon and Junmo Kim. Active convolution: Learning the shape of convolution for image classification. In *CVPR*, pages 4201–4209, 2017.
- [34] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc Van Gool. Dynamic filter networks. In *NIPS*2016*, pages 667–675.
- [35] Alex Kendall and Yarin Gal. What uncertainties do we need in Bayesian deep learning for computer vision? In *NIPS*2017*, pages 5574–5584.
- [36] Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [37] Diederik P. Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *NIPS*2015*, pages 2575–2583.

- [38] Patrick Knöbelreiter and Thomas Pock. Learned collaborative stereo refinement. In *GCPR*, pages 3–17, 2019.
- [39] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NIPS*2017*, pages 6402–6413.
- [40] Yijun Li, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Joint image filtering with deep convolutional networks. *IEEE T. Pattern Anal. Mach. Intell.*, 41(8):1909–1923, 2019.
- [41] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015.
- [42] Christos Louizos and Max Welling. Multiplicative normalizing flows for variational Bayesian neural networks. In *ICML*, pages 2218–2227, 2017.
- [43] Tan Khoa Mai, Michèle Gouiffès, and Samia Bouchafa. Optical flow refinement using reliable flow propagation. In *VIS-APP*, pages 451–458, 2017.
- [44] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *CVPR*, pages 3061–3070, 2015.
- [45] Marcus Müller, Frederik Zilly, and Peter Kauff. Adaptive cross-trilateral depth map filtering. In *3DTV*, 2010.
- [46] David A. Nix and Andreas S. Weigend. Estimating the mean and variance of the target probability distribution. In *ICNN*, pages 55–60, 1994.
- [47] David Novotny, Samuel Albanie, Diane Larlus, and Andrea Vedaldi. Self-supervised learning of geometrically stable features through probabilistic introspection. In *CVPR*, pages 3637–3645, 2018.
- [48] Jinshan Pan, Jiangxin Dong, Jimmy S. Ren, Liang Lin, Jinhui Tang, and Ming-Hsuan Yang. Spatially variant linear representation models for joint filtering. In *CVPR*, pages 1702–1711, 2019.
- [49] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, 2015.
- [50] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Guided curriculum model adaptation and uncertainty-aware evaluation for semantic nighttime image segmentation. In *ICCV*, pages 7373–7382, 2019.
- [51] Hang Su, Varun Jampani, Deqing Sun, Orazio Gallo, Erik Learned-Miller, and Jan Kautz. Pixel-adaptive convolutional neural networks. In *CVPR*, pages 11158–11167, 2019.
- [52] Carlo Tomasi and Roberto Manduchi. Bilateral filtering for gray and color images. In *ICCV*, pages 839–846, 1998.
- [53] Christoph Vogel, Patrick Knöbelreiter, and Thomas Pock. Learning energy based inpainting for optical flow. In *ACCV*, volume 6, pages 340–356, 2018.
- [54] Hao Wang, Xingjian Shi, and Dit-Yan Yeung. Natural-parameter networks: A class of probabilistic neural networks. In *NIPS*2016*, pages 118–126.
- [55] Jiaqi Wang, Kai Chen, Rui Xu, Ziwei Liu, Chen Change Loy, and Dahua Lin. CARAFE: Content-Aware ReAssembly of FEatures. In *ICCV*, pages 3007–3016, 2019.
- [56] Anne S. Wannenwetsch, Martin Kiefel, Peter V. Gehler, and Stefan Roth. Learning task-specific generalized convolutions in the permutohedral lattice. In *GCPR*, pages 345–359, 2019.
- [57] Anqi Wu, Sebastian Nowozin, Edward Meeds, Richard E. Turner, José Miguel Hernández-Lobato, and Alexander L. Gaunt. Deterministic variational inference for robust Bayesian neural networks. In *ICLR*, 2019.
- [58] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. Fast end-to-end trainable guided filter. In *CVPR*, pages 1838–1847, 2018.
- [59] Jialin Wu, Dai Li, Yu Yang, Chandrajit Bajaj, and Xiangyang Ji. Dynamic filtering with large sampling field for convnets. In *ECCV*, volume 10, pages 188–203, 2018.
- [60] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *CVPR*, pages 6044–6053, 2019.

Probabilistic Pixel-Adaptive Refinement Networks

– Supplemental Material –

Anne S. Wannewetsch^{1,2*} Stefan Roth²
¹Amazon, Germany ²TU Darmstadt, Germany

In this supplemental material, we give further implementation details of the different types of refinement networks and provide results for a more comprehensive comparison on optical flow benchmarks. Moreover, we present an analysis considering the PPAC improvements on unreliable pixels as well as additional visualizations of PPAC-refined optical flow fields and segmentation maps.

A. Additional Implementation Details

A.1. Learning procedure

To train our networks, we use the Adam optimizer [36] with default parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$ and without weight decay. PPAC refinement networks are trained with a learning rate of 1×10^{-3} for networks on Sintel and a learning rate of 5×10^{-3} on KITTI. For semantic segmentation on Pascal VOC 2012, we use a learning rate of 1×10^{-4} for guidance and probability branches and 1×10^{-5} for the remaining PPAC parameters. The image inputs to all networks are normalized while estimates and log-probabilities remain unchanged. For faster training of all refinement networks, we save the outputs of the underlying backbone networks (*i.e.* HD3 or DeepLabv3+), and only propagate through the refinement step.

A.2. Network architectures

Tables 6 – 8 show the network structures used for PPAC, PAC, and our baseline simple refinement network, respectively. Here, ‘C’ represents standard convolution layers, ‘P’ layers with non-probabilistic PACs, and ‘PP’ layers with our PPACs. The networks for optical flow and semantic segmentation differ mainly by the number of input and output channels (2 or 21, respectively). For optical flow, the guidance branch uses only the first image as input since the flow fields should be aligned w.r.t. the objects in this image. All standard as well as PAC and PPAC-convolutions pad the inputs with zeros to preserve the feature size and use a stride of one. Moreover, a bias term is learned for all types of convolutions. The output of guidance and, if applicable,

*This work was done at TU Darmstadt prior to Anne S. Wannewetsch joining Amazon.

Table 6. Network structure of our PPAC networks for optical flow/semantic segmentation with $\sim 12.3\text{k}/14.3\text{k}$ parameters.

| | layer type | kernel size | non-linearity | shared weights |
|-------------|-----------------------------|--------------|---------------|----------------|
| guidance | C: 3 \rightarrow 15 | 5×5 | ReLU | \times |
| branch | C: 15 \rightarrow 15 | 5×5 | ReLU | \times |
| | C: 15 \rightarrow 10 | 5×5 | \times | \times |
| probability | C: 5/21 \rightarrow 5 | 5×5 | ReLU | \times |
| branch | C: 5 \rightarrow 5 | 5×5 | ReLU | \times |
| | C: 5 \rightarrow 2 | 5×5 | Sigmoid | \times |
| combination | PP: 2/21 \rightarrow 2/21 | 7×7 | \times | \checkmark |
| branch | PP: 2/21 \rightarrow 2/21 | 7×7 | \times | \checkmark |

Table 7. Network structure of our PAC baseline networks for optical flow/semantic segmentation with $\sim 12.6\text{k}/15.5\text{k}$ parameters.

| | layer type | kernel size | non-linearity | shared weights |
|-------------|------------------------------|--------------|---------------|----------------|
| guidance | C: 8/24 \rightarrow 15/13 | 5×5 | ReLU | \times |
| branch | C: 15/13 \rightarrow 15/13 | 5×5 | ReLU | \times |
| | C: 15/13 \rightarrow 10 | 5×5 | \times | \times |
| combination | P: 2/21 \rightarrow 2/21 | 7×7 | \times | \checkmark |
| branch | P: 2/21 \rightarrow 2/21 | 7×7 | \times | \checkmark |

Table 8. Network structure of our simple baseline network for optical flow with a total of $\sim 12.4\text{k}$ parameters.

| | layer type | kernel size | non-linearity | shared weights |
|--------|------------------------|--------------|---------------|----------------|
| simple | C: 10 \rightarrow 11 | 7×7 | ReLU | \times |
| branch | C: 11 \rightarrow 11 | 7×7 | ReLU | \times |
| | C: 11 \rightarrow 2 | 7×7 | \times | \times |

probability branches is split equally by the number of PAC or PPACs such that individual guidance is used for the components of the combination branch. For the simple setup, we equally experimented with networks with two convolutions and thus more channels but found the given one with three convolutions to perform better.

Table 9. Average end-point error (*AEE*) of top-ranked two-frame optical flow methods on Sintel train and test. *Re-evaluated for comparability.

| | train | | test | |
|--------------------------|---------------|---------------|-------------|-------------|
| | clean | final | clean | final |
| VCN [67] | (1.66) | (2.24) | 2.81 | 4.40 |
| IRR-PWC [26] | (1.92) | (2.51) | 3.84 | 4.58 |
| PWC-Net+ [64] | (1.71) | (2.34) | 3.45 | 4.60 |
| PPAC-HD3 (<i>ours</i>) | (1.54) | (1.05) | 4.59 | 4.60 |
| HD3 [60] | (1.68)* | (1.15)* | 4.79 | 4.67 |

Table 10. Average end-point error (*AEE*) and 3-pixel outlier rate on non-occluded/all pixels (*Out-Noc/all*) of top-ranked optical flow methods on KITTI 2012 train and test. Results in parentheses indicate that data was used in training. †Methods use left and right stereo images. *Re-evaluated for comparability.

| | train | | test | |
|--------------------------|---------------|------------|---------------|---------------|
| | AEE | AEE | Out-Noc | Out-all |
| PPAC-HD3 (<i>ours</i>) | (0.71) | 1.2 | 2.01 % | 5.09 % |
| HD3 [60] | (0.81)* | 1.4 | 2.26 % | 5.41 % |
| PRSM† [65] | – | 1.0 | 2.46 % | 4.23 % |
| LiteFlowNet2 [62] | – | 1.4 | 2.63 % | 6.16 % |
| SPS-StFl† [66] | – | 1.3 | 2.82 % | 5.61 % |

Table 11. Average end-point error (*AEE*), 3-pixel outlier rate on all pixels (*Out-all*), and runtimes (*time*) of top-ranked methods on KITTI 2015 train and test. Results in parentheses indicate that data was used in training. †Methods use left and right stereo images. *Re-evaluated for comparability.

| | train | | test | |
|--------------------------|---------------|-----------------|---------------|-----------------|
| | AEE | Out-all | Out-all | time |
| UberATG-DRISF† [63] | – | – | 4.73 % | 0.75 s |
| PPAC-HD3 (<i>ours</i>) | (1.20) | (3.56 %) | 6.06 % | 0.19 s |
| ISF† [61] | – | – | 6.22 % | 600 s |
| VCN [67] | (1.16) | (4.1 %) | 6.30 % | 0.18 s |
| HD3 [60] | (1.40)* | (4.39 %)* | 6.55 % | 0.11 s * |

B. Detailed Comparison on Optical Flow Benchmarks

For completeness, we give a more detailed comparison on benchmarks for optical flow, including the training results of PPAC-HD3 as well as the results of related work.

Table 9 shows results on Sintel clean and final. For comparability, we re-evaluated the flow fields of HD3 on the training splits, taking into account the available invalid masks. Our proposed PPAC-HD3 ranks 4th w.r.t. to the AEE on Sintel final.

Tables 10 and 11 summarize results for the best-ranked published methods on KITTI 2012 and 2015. For com-

Table 12. Relative improvement of average end-point error (*AEE*), evaluated on the 10% most unreliable and the remaining pixels of our Sintel and KITTI test splits.

| | Sintel | | KITTI |
|------------------------|--------|--------|--------|
| | clean | final | |
| Most unreliable pixels | 9.86 % | 8.93 % | 4.28 % |
| Remaining pixels | 4.11 % | 3.01 % | 9.42 % |

pleteness, we also include scene flow methods. Note, however, that such approaches are not fully comparable as they leverage additional stereo images to compute flow. On both datasets, PPAC-HD3 ranks 1st among optical flow methods and 2nd over all published approaches on KITTI 2015. As we used the publically available checkpoint for HD3, which differs slightly from the one used in [60], we report re-evaluated results on the training sets. Moreover, we provide HD3 runtimes evaluated on the same GTX 1080 Ti GPU as PPAC-HD3 for fair comparison.

C. Improvement of Unreliable Pixels

In the main paper, we argue that probabilistic pixel-adaptive refinement allows to propagate correct estimates into unreliable regions. Here, we examine the influence of PPACs on unreliable pixels in more detail. We evaluate the refinement of optical flow by computing the AEE on the 10% most unreliable pixels of each flow field and comparing it to the AEE of the remaining pixels. To assess the reliability of a pixel estimate, we upsample the probabilities of the last output scale and use nearest neighbor interpolation if the estimated residuals fall outside the probability grid. We found these reliabilities to correlate better with the optical flow errors than the ones obtained from the composed full matching probability distribution proposed in [60]. Moreover, we use the same PPAC refinement networks as trained for the experiments in Table 2.

Table 12 shows the relative improvement on unreliable and remaining pixels evaluated on our test splits of Sintel and KITTI. On Sintel clean and final, we clearly observe a more significant improvement on the unreliable pixels, justifying the conclusion that PPACs allow to replace pixels of low reliability. In contrast, our evaluation on KITTI shows a larger improvement for the remaining pixels. This correlates well with the fact that we found the output probabilities of [60] to be less well calibrated on KITTI, judging by the comparatively larger benefit of oracle confidences, *c.f.* Sec. 6.1. However, when comparing the relative improvements of PPACs to the ones obtained by PACs (8.87% for more reliable and 2.72% for uncertain pixels), we observe that PPACs nevertheless allow for better handling of unreliable regions even if the reliability estimates are not completely accurate themselves.



Figure 7. Examples of ground truth (*top*), HD3 optical flow [60] (*middle*), and our PPAC-refined optical flow (*bottom*) on Sintel final. *Best viewed on screen.*

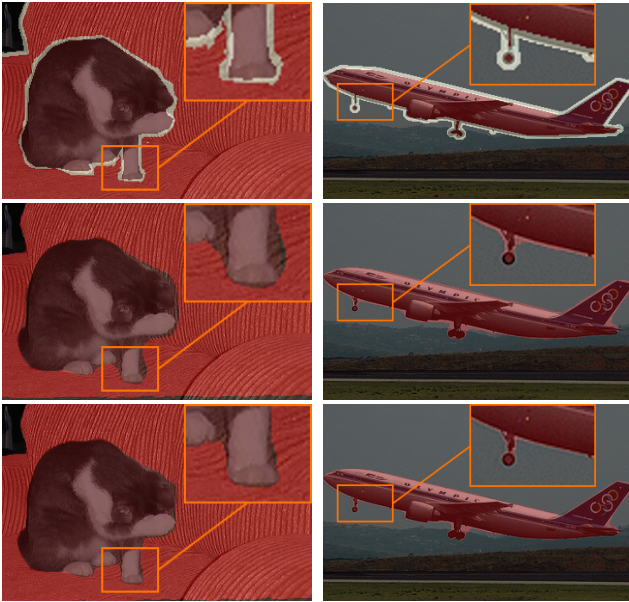


Figure 8. Additional examples of cropped ground truth (*top*), DeepLabv3+ [7] (*middle*), and PPAC-refined segmentation maps (*bottom*) on Pascal VOC 2012. *Best viewed on screen.*

D. Additional Visualizations

Fig. 7 shows additional visualizations of refined optical flow fields on our own validation and test splits of Sintel final. As such, none of these flow fields was presented to the PPAC refinement network during training. We clearly observe improved motion boundaries but also the ability of our approach to correctly propagate estimates into erroneous regions, *e.g.* the bird wings on the leftmost example.

In Fig. 8, we provide additional visualizations of refined segmentation maps on Pascal VOC 2012. PPAC refinement leads to a clear reduction of errors near object boundaries, *e.g.* by considerably minimizing the segmentation margin visible at the cat paw.

References

- [61] Aseem Behl, Omid Hosseini Jafari, Siva Karthik Mustikovela, Hassan Abu Alhaija, Carsten Rother, and Andreas Geiger. Bounding boxes, segmentations and object coordinates: How important is recognition for 3D scene flow estimation in autonomous driving scenarios? In *ICCV*, pages 2574–2583, 2017.
- [62] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. A lightweight optical flow CNN – Revisiting data fidelity and regularization. *arXiv:1903.07414 [cs.CV]*, 2019.
- [63] Wei-Chiu Ma, Shenlong Wang, Rui Hu, Yuwen Xiong, and Raquel Urtasun. Deep rigid instance scene flow. In *CVPR*, pages 3614–3622, 2019.
- [64] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Models matter, so does training: An empirical study of CNNs for optical flow estimation. *arXiv:1809.05571 [cs.CV]*, 2018.
- [65] Christoph Vogel, Konrad Schindler, and Stefan Roth. 3D scene flow estimation with a piecewise rigid scene model. *Int. J. Comput. Vision*, 115(1):1–28, 2015.
- [66] Koichiro Yamaguchi, David McAllester, and Raquel Urtasun. Efficient joint segmentation, occlusion labeling, stereo and flow estimation. In *ECCV*, volume 5, pages 756–771, 2014.
- [67] Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow. In *NeurIPS*2019*, pages 784–805.