

# Learning to Think Outside the Box: Wide-Baseline Light Field Depth Estimation with EPI-Shift

Titus Leistner<sup>1</sup>, Hendrik Schilling<sup>1</sup>, Radek Mackowiak<sup>2</sup>, Stefan Gumhold<sup>3</sup>, Carsten Rother<sup>1</sup>  
Visual Learning Lab, Heidelberg University<sup>1</sup>, Robert Bosch GmbH<sup>2</sup>, TU Dresden<sup>3</sup>

first.last@iwr.uni-heidelberg.de<sup>1</sup>, first.last@de.bosch.com<sup>2</sup>, first.last@tu-dresden.de<sup>3</sup>

## Abstract

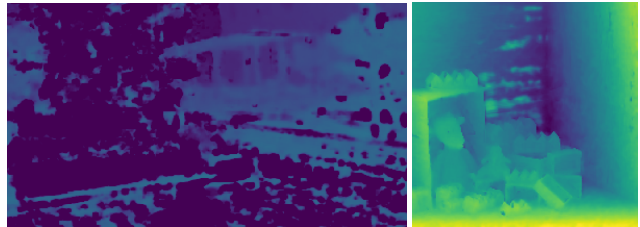
We propose a method for depth estimation from light field data, based on a fully convolutional neural network architecture. Our goal is to design a pipeline which achieves highly accurate results for small- and wide-baseline light fields. Since light field training data is scarce, all learning-based approaches use a small receptive field and operate on small disparity ranges. In order to work with wide-baseline light fields, we introduce the idea of EPI-Shift: To virtually shift the light field stack which enables to retain a small receptive field, independent of the disparity range. In this way, our approach “learns to think outside the box of the receptive field”. Our network performs joint classification of integer disparities and regression of disparity-offsets. A U-Net component provides excellent long-range smoothing. EPI-Shift considerably outperforms the state-of-the-art learning-based approaches and is on par with hand-crafted methods. We demonstrate this on a publicly available, synthetic, small-baseline benchmark and on large-baseline real-world recordings.

## 1. Introduction

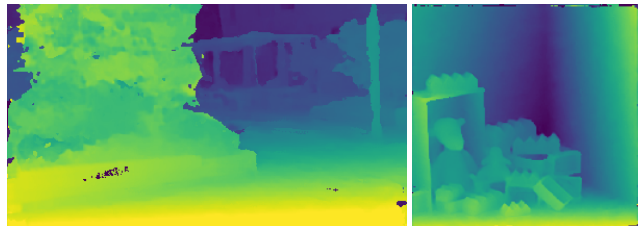
In the current deep learning era, autonomous driving and intelligent robotics are becoming more and more practical. For those applications, and many more, accurate depth perception is often a crucial component. Depth estimation can be performed with a wide variety of different hardware setups. On one end of the spectrum there is monocular depth estimation from single cameras [4][5]. In between, there is stereo-matching from rectified images captured by two cameras [21][11]. At the other end of the spectrum there is depth estimation from light field camera arrays [23] which is the focus of this work. Light field depth estimation can be seen as a well-structured multi-view stereo approach. Although the hardware setup of a light field is more involved than a stereo setup, the depth information “encoded” in light fields is considerably more precise and inherently less sensi-



(a) Center view of real (left) and synthetic (right) light field



(b) EPINET-Cross [20]



(c) Our EPI-Shift

Figure 1: **Light Field Depth Estimation.** (a) A *real* light field (left) with a large disparity range of  $[0, 12]$  and a *synthetic* light field (right) with a small disparity range of  $[-2, 2]$ . (b) The current state-of-the-art method EPINET-Cross [20] has only been trained for the small disparity range. It therefore fails at extreme disparities in the synthetic image (right, background) and outside of the trained range in the real image (left, foreground). (c) Our EPI-Shift approach performs well for both, small and large disparities. Due to better generalization, it even outperforms EPINET-Cross for small disparity ranges.



Figure 2: **Extraction of an Epipolar-Plane Image (EPI).** The vertical views captured by a **cross light field** (left) are stacked (middle). A slice through a column (right) forms an **EPI** (bottom). Note, how the slopes of lines in the EPI encode the three-dimensional structure of the scene.

tive to occlusions, due to the multiple redundant viewpoints.

There are two types of light field cameras. First, a multi-camera setup, see *e.g.* a cross-camera in Figure 2 with its corresponding EPI. Such a setup is unfortunately expensive, difficult to assemble, synchronize and calibrate, compared to *e.g.* a stereo setup. However, once these challenges are mastered, the major advantage of a multi-camera system is its accuracy due to the wide baseline, since reconstruction accuracy grows linearly with the baseline between viewpoints and hence, a small baseline yields inferior accuracy. We therefore focus on this setup. Second, a plenoptic camera based on a micro-lens array [15] which has however, a rather limited resolution and baseline.

In this work, we propose a learning-based light field depth estimation method. This is challenging due to the non-availability of real-world training data. The creation of real-world reference depth is problematic, as no other dense measurement approach is more accurate than light field depth estimation. For example, structured light scanning is problematic in the context of occlusions and Light Detection and Ranging (LIDAR) is considerably more sparse than light field data. Therefore, all training data is synthetic and the pool of publicly available datasets is small. Also, all these training images have a small baseline, emulating a micro-lens based camera rather than a camera array. Hence, current learning-based approaches fail dramatically for wide-baseline light fields, as exemplified in the real world scene in Figure 1. Interestingly, even for smaller disparities, the performance is limited by poor generalization, as demonstrated in the synthetic scene in Figure 1, where artifacts appear *within* the trained disparity range.

One major cause for this problem, the limited receptive field of previous methods, is illustrated in Figure 3. Expanding the receptive field would cause worse generalization performance. However, by applying EPI-Shifts to the input of our neural network, we circumvent this flaw.

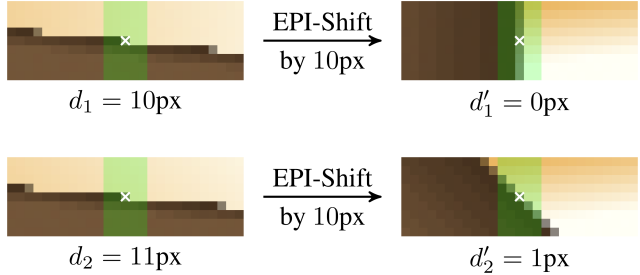


Figure 3: **Idea of EPI-Shift.** Two EPIs (left), consisting of horizontal lines from different cameras. The CNNs task is to estimate the correct disparity for the center pixel (white cross) by predicting the slope of the line going through this pixel. Note the nearly invisible difference between a disparity of 10 and 11 pixels which can only be estimated using a large receptive field. After applying an **EPI-Shift** of 10 pixels (right), the difference is clearly visible. Therefore, our network only requires a minimal receptive field (green box) to classify whether the shifted disparity lies within  $\pm 0.5$  pixels and regress a sub-pixel accurate disparity offset.

The basic idea is inspired by the technique of plane-sweep volumes. Instead of directly estimating the disparity from the light field image stacks, we utilize a plane, sweeping through space, as common for stereo and multi-view depth estimation [3][11]. Hence, we split the task into classification and regression.

The classification map states, per pixel, whether objects observed at tested plane sweep are within a refocused disparity range of  $\pm 0.5$ . It is used to merge all independent estimates from the plane sweep volume, while the disparity regression provides sub-pixel accuracy. This approach considerably improves generalization, as we are now able to infer the depth of a wide-baseline scene with a network, trained solely with small-baseline training data.

Let us summarize our main contributions:

- Applying the idea of plane-sweep volumes in the context of light fields, which we denote as EPI-Shift. This enables learning-based approaches to generalize well to large-disparity test data, even with small-disparity training data.
- A network architecture, which enables improved long-range reasoning by combining a feature extraction network [20] with a subsequent U-Net architecture [17][8] for excellent long-range smoothing with low artifacts.
- Our approach considerably outperforms the state-of-the-art learning-based approaches with same input modality and is on par with hand-crafted methods.

## 2. Related Work

Light field depth estimation methods can be categorized into optimization-based and learning-based approaches.

### 2.1. Methods based on Optimization

Heber and Pock [6] applied Robust Principal Component Analysis (RPCA) to light field depth estimation. Their method measures the quality of image alignments by projecting images to column vectors of a shared matrix, followed by a convex optimization of stereo matching and denoising. However, the method is vulnerable to occlusions.

Jeon *et al.* [12] specifically address lenslet cameras, proposing a novel distortion correction. The actual depth estimation is performed using a pixel-wise cost volume inspired by traditional stereo matching techniques which is combined with discrete Fourier transform. To ensure smoothness, the algorithm applies a sparse matching using graph cut at salient feature points based on Scale-Invariant Feature Transform (SIFT). The method shows a good performance on real images from micro-lense cameras but suffers from noise and artifacts at occlusion boundaries.

Lin *et al.* [13] utilize a focal stack composed of light field data for depth estimation. Their method is based on the local symmetry of the color distribution in a patch around the real depth. Smoothness is ensured by an optimization with graph cut. Unfortunately, the refocusing of micro lense images produces slight artifacts, making the approach not as accurate as most EPI-based methods.

Wanner and Goldluecke [22] proposed a solution that calculates the disparity directly, using linear algebra operations. To estimate the line slope, an adapted structure tensor is applied to the 4D light field. The strongest eigenvector of this tensor, composed of all partial derivations, is aligned with the EPI gradient. As this method produces inaccuracies in non-textured regions and outliers at occlusions, additional optimization methods are required.

Neri *et al.* [14] proposed another method based on  $\mathcal{L}^2$  matching costs. In order to ensure smoothness, a local optimization of discrete depth labels on a resolution pyramid is performed. The method still produces artifacts in large non-textured areas and at object boundaries.

Zhang *et al.* [24] introduced the Spinning Parallelogram Operator (SPO) fitting lines with different slopes through the center pixel of a parallelogram. Each line divides the parallelogram in two distinct areas. The SPO computes  $\chi^2$  differences of color distributions between those areas, used as confidence measure. However, the quality at occlusions still depends on the relative angle between EPI extractions and objects boundaries.

Sheng *et al.* [19] introduced a method to compensate for this effect, utilizing a bigger set of EPIs extracted at arbitrary angles. For occlusion reasoning, the algorithm utilizes the variance of depth estimates at different EPI orientations.

Using the direction of the derived occlusion boundary, the optimal non-occluded EPI is selected to determine the final depth value. The authors also utilize an operator similar to SPO but comparing the color distributions between semicircles instead of parallelograms. While this method handles occlusions very well, it depends on a large number of views that are not present in all real applications.

In contrast, Schilling *et al.* [18] introduced a method that does only depend on a cross input. It performs local optimization based on PatchMatch [1]. This model produces state-of-the-art results on the HCI Light Field Benchmark [10]. However, it contains many hand tuned parameters which could probably be improved using a learning-based method.

### 2.2. Methods based on Deep Learning

In 2016 Heber *et al.* [7] introduced the first deep learning based methods for the task, predicting 2D per-pixel hyperplane orientations. The utilized CNN processes a cross subset of the light field. A disparity map is then inferred from the hyperplane parameters by optimizing a convex energy functional. Due to the lack of training scenes, the authors also contributed a randomly generated synthetic dataset. The approach is a first step in the direction of learned depth from light fields but the results suffer from strong artifacts and blur.

In order to overcome those downsides, the same authors improved their work [8][9] by utilizing a U-Net [17]. The first paper [8] shows visual and quantitative improvement, but suffers from streaking artifacts, addressed in [9].

One of the most recent works by Shin *et al.* [20] deals with light field depth estimation using a fully-convolutional architecture. The authors introduce a multi-stream network comprised of multiple inputs for the horizontal, vertical and both diagonal light field stacks. All outputs of the individual streams are then being concatenated and fed into a second CNN. The approach reaches state-of-the-art results on the HCI 4D Light Field Benchmark [10]. However, it is limited to the disparity interval seen during training.

### 2.3. Plane Sweep Volumes

The concept of plane sweep volumes was introduced by Collins [3] in 1996. He addresses the problem of multi-image matching. Therefore, features from each input image are projected to a set of parallel planes sweeping through space at increasing depth.

Sunghoon *et al.* [11] apply this principle to multi-view stereo, using a CNN architecture. A cost volume is created by shifting features, extracted from the multiple input images using an encoder-architecture. Subsequently, the disparity is regressed by aggregating those costs. This method enables state-of-the-art occlusion handling results for depth reconstruction from stereo recordings.

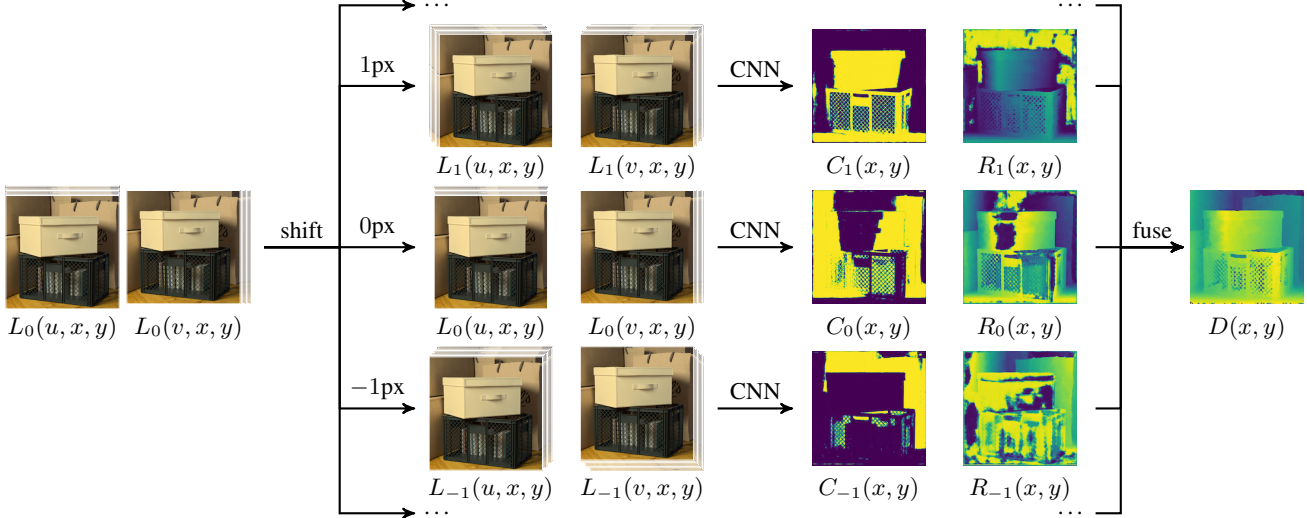


Figure 4: **Method Overview.** The **input** (left) consisting of two view stacks is shifted several times, producing stacks with different disparity ranges. Our **CNN** (middle) processes the shifted stacks, inferring a classification and regression output. Each pixel of the **final result** (right) is assigned to a discrete disparity (classification) and refined by a sub-pixel disparity offset (regression).

### 3. Method

The core idea of our method is to use a plane sweep volume [3][11] and successively apply the same neural network to each depth plane of that volume individually. The output of the network for each of these disparity ranges are two 2D maps, one for the classification (correct plane/incorrect plane) and one for the disparity offset from the plane, in the range  $-0.5$  and  $0.5$ . To generate the full disparity map of the scene, for each pixel, the shift with maximum classification activation is chosen to determine the correct plane. The corresponding per-pixel disparity offset is added to achieve sub-pixel accuracy. Because we are using a cross light field setup, the plane sweep volume can be constructed using the EPI-Shift approach, which refocuses the image stack by applying a skew transformation.

#### 3.1. Light Field Camera Setup

Goal of our method is an accurate per-pixel disparity reconstruction within the center view of a  $9 + 8$  cross-shaped light field camera setup with a large baseline. We limit ourselves to this setup due to the versatile usability in real world scenarios, compared to a star, or full 4D light field setup. Many recent submissions to the HCI 4D Light Field Benchmark [10] demonstrate that research is shifting towards using more views from the available 81 input views, *e.g.* by synthesizing a focal stack from all views. However, we argue that this is a symptom of “benchmark optimization”, because adding more views gives diminishing returns and using less views is more practical in the real world. Note that the best approach [18] is not learning-based and

requires only 17 views, compared to the inferior EPINET-Star [20] setup which requires nearly double the amount of views.

4D light fields are recorded by an array of cameras, arranged on a regular 2D grid, indexed by  $(u, v)$ . The baseline represents the distance between two adjacent cameras. Each camera captures a 2D image with pixels, indexed by  $(x, y)$ . Similar to the effect of alternatively closing the left and right eye, objects seem to move between different viewpoints. A change of viewpoints on the  $u$ -axis for example, causes movement of a projected object point along the  $x$ -axis. The straight lines in image space, which represent this depth dependent movement, are called epipolar lines. For a cross light field the 2D slices of the light field along the  $xu$ - and  $yv$ -planes represent the Epipolar-Plane Images (EPIs) (see Figure 2) [2]. The  $x$ - or  $y$ -distance between the same object point in two adjacent views is the disparity  $d$ , measured in pixels and being inversely proportional to the depth. During camera calibration and rectification, the images often get pre-shifted. Therefore, also negative disparity values occur in some light field datasets.

#### 3.2. EPI-Shift

Our EPI-Shift approach, which generates the plane sweep volume, boils down to a skew transformation on the  $xu$ -plane. Given a 3D stack of views,

$$L_0(u, x, y) \tag{1}$$

defines the color value at a given image position  $(x, y)$ , recorded by a camera  $u$ . The central camera is defined to

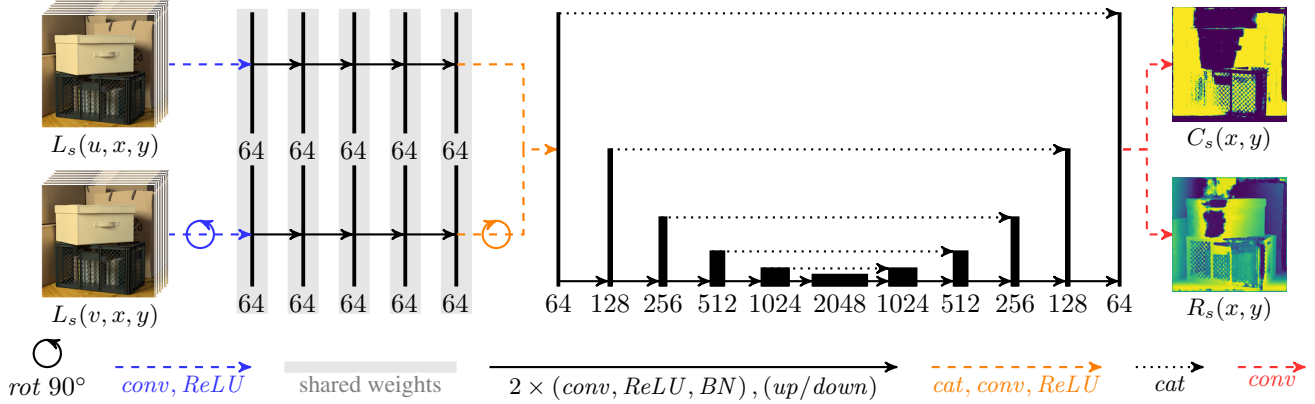


Figure 5: The **neural network architecture** of our model, consisting of two parts. First, a **siamese feature extractor** [20] (left), with four convolutional blocks for the discovery of local disparity information. Second, a **U-Net** architecture [17][8] (right) to integrate global information. The input consists of two view stacks. Our network outputs a classification of discrete depth labels and a sub-pixel accurate disparity regression. The network solely uses convolutional blocks, consisting of two consecutive  $3 \times 3$  convolutions with stride and padding of one. Numbers refer to the number of channels.

be at  $u = 0$ . Positive  $u$ -indices are assigned to cameras that are located to the right of the central camera. As visualized in Figure 3, we perform the EPI-Shift by a certain disparity  $s$  with

$$L_s(u, x, y) = L_0(u, x - us, y). \quad (2)$$

Note that this operation refers to horizontal EPIs only. However, vertical EPIs  $L_0(v, x, y)$  behave analogously after a rotation about  $90^\circ$  around the  $v$ -axis.

Because  $s$  is defined to be an integer number and we use a cross-shaped setup, no interpolation is required. We perform nearest-neighbor padding by clipping  $x - us$  to the valid pixel range. However, in order to not waste capacity of the neural network for learning to deal with image-borders, we do not apply any loss in areas effected by the padding.

To enable wide-baseline light field depth estimation, we perform three basic steps, illustrated in Figure 4. First, we generate a plane sweep volume by applying the EPI-Shift to the input light field, once per integer disparity within the disparity range of the scene. Each of those shifts can be thought of as a discrete disparity label. A pixel is assigned to a certain label if the disparity lies between  $-0.5$  and  $0.5$ , when shifted by the labels disparity. Second, we infer a classification and a regression map for each of the shifts, using the CNN architecture described in Section 3.3. Third, we compute the final result  $D$  for each pixel  $(x, y)$  by assigning an integer disparity label

$$l(x, y) = \operatorname{argmax}_s (C_s(x, y)) \quad (3)$$

according to the shift  $s$  that produced the highest classification output  $C_s$ . Using the regression map  $R_s$  of the respective shift, we add fine-grained disparity information to

achieve the sub-pixel accurate result

$$D(x, y) = l(x, y) + R_{l(x, y)}(x, y). \quad (4)$$

### 3.3. Network Architecture

Our architecture consists of two parts visualized in Figure 5. First, a siamese feature extraction network similar to [20]. The purpose of this subnetwork is the extraction of local disparity information. Second, a U-Net architecture (compare [17] and [8]) with two outputs: A classification output, assigning discrete per-pixel disparity labels and a continuous regression output, representing the sub-pixel accurate disparity relative to the label.

**Siamese Feature Extraction Network:** The cross light field introduced in Figure 2 provides a horizontal and a vertical stack of input views. Instead of concatenating the two, we chose a siamese twin architecture, consisting of one subnetwork for each stack. As both stacks contain similarly aligned EPIs after rotation of one stack by  $90^\circ$ , we share weights between the two subnetworks. This reduces the number of network parameters and therefore improves generalization.

The feature extraction network contains four fully-convolutional blocks, each consisting of two  $3 \times 3$  convolutions followed by a *ReLU* activation function and a batch normalization layer each. We chose a number of 64 channels and preserve the input dimensions with a padding and stride of one.

To facilitate the classification for the downstream U-Net, we provide it with additional data. We apply the feature extraction network to both adjacent EPI-Shifts  $L_{\pm 1}$ . The extracted features are concatenated with those, extracted

from  $L_0$  as well as the color information of the center view  $L_0(0, x, y)$ . Hence, the number of input channels for the U-Net is, for each shift: The number of channels from the feature extraction subnetwork times two (horizontal + vertical) times three for the three shifts  $(-1, 0, +1)$  plus the center view, i.e. a total of  $64 \cdot 2 \cdot 3 + 3 = 387$  channels. Our experiments showed that these additional inputs improved the distinction between foreground and background objects in ambiguous regions. This is probably due to the depth hints from adjacent shifts that provide additional information about occlusions and therefore simplify classification. The addition of the center view allows the joint model to focus on feature extraction in the first part of the network, but still use the center view to guide smoothing in the U-Net part.

**U-Net:** A U-Net [17] architecture expands the effective receptive field of the joint model without loss of generalization capability. It therefore significantly improves the smoothness of non-textured areas. We chose a depth of five down- and up-sampling layers leading to a receptive radius of 124 pixels for the U-Net part and 135 pixels for the whole network. The concatenated output of the upstream feature extractor network is reduced from 387 to 64 channels by an additional  $3 \times 3$  convolution. For the processing inside the U-Net, we chose the same convolutional blocks as for the feature network, followed by an additional  $3 \times 3$  up- or down-sampling convolution. The downsampling layers bisect the image dimensions while doubling the number of channels. Prior to upsampling, the output of the corresponding downsampling is concatenated. Therefore, the upsampling process doubles the image dimensions but divides the number of channels by four, please see [17] for more details. A final  $3 \times 3$  convolutional layer transforms the 64 output channels of the U-Net to two channels for the classification and regression output. Because the regressed disparity can be negative, no final *ReLU* activation function is applied.

### 3.4. Loss Function

Due to the drastically different outputs of our network, the choice of a well performing loss function is not trivial. For the classification output, a slight overlap between adjacent shifts might not have an effect on the final result at all. A large output at distant disparity regions however may cause a misclassification and therefore can destroy the end result. Our classification loss therefore specifically penalizes those cases. We define the loss with  $C_s(x, y)$  being the classification output for a given shift  $s$  at pixel  $(x, y)$  as

$$\mathcal{L}_{class} = \sum_{s,x,y} (C_s(x, y) - C_s^*(x, y))^2 \cdot \mathcal{W}_{disp}(x, y) \quad (5)$$

with a disparity weighting of

$$\mathcal{W}_{disp}(x, y) = (D(x, y) - D^*(x, y))^2 \quad (6)$$

computed using the final disparity output  $D$  and the ground truth disparity  $D^*$  that penalizes misclassifications during training. The classification ground truth  $C^*$  should be high for all pixels within a disparity of  $\pm 0.5$  pixels. We therefore tried two different definitions: First, the one-hot or rectangle function

$$C_s^*(x, y) = \begin{cases} 1 & \text{if } |D^*(x, y) - s| \leq 0.5 + \epsilon \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

producing hard boundaries between two labels. Second, the triangle function

$$C_s^*(x, y) = \max(0.5 + \epsilon - |D^*(x, y) - s|, 0) \quad (8)$$

which is more closely related to the regression output. It therefore should accelerate training and engage the network to share weight capacities between the two. On the other hand, it outputs lower values at boundaries, which are more vulnerable to misclassifications. In both cases, we choose a small  $\epsilon$  that produces a slight overlap at the border regions between two disparity labels in order to prevent wrong classifications. We will see in the experiments that the rectangle function (Equation 7) performs slightly better.

The regression output requires smooth surfaces but sharp edges. We see in our experiments that the  $\mathcal{L}^1$  loss function fulfills those requirements for the regression loss best. We therefore define it as

$$\mathcal{L}_{reg} = \sum_{s,x,y} |R_s(x, y) - D^*(x, y) + s| C_s^*(x, y) \quad (9)$$

with  $R$  being the regression output. We mask out all pixels outside the sub-pixel interval by weighting with the rectangle function  $C_s^*$  defined in Equation 7. In order to compensate for misclassifications at the boundaries of disparity labels, we choose a slightly higher  $\epsilon$  than for the classification ground truth. Due to the fundamentally different trend of the losses during training, a weighting between the two is also important for computation of the overall loss

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{reg} + \mathcal{L}_{class}. \quad (10)$$

The choice of a scaling factor  $\alpha$  depends on various factors such as the disparity distribution of the training data.

### 3.5. Training

We trained our model on 16 scenes of the HCI 4D Light Field Benchmark [10] that are not part of the benchmark evaluation. We implemented the model in PyTorch [16] and trained it for four days on three NVIDIA TITAN X GPUs. As optimizer we chose Adam with a learning rate of  $10^{-4}$  for the first 10000 iterations. For another 30000 iterations, we decreased the learning rate to  $10^{-5}$  and fixed the learned batch normalization parameters.

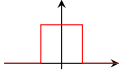
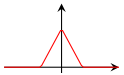
| Class GT  | $\epsilon = 0.0$<br>$\alpha = 0.25$ | $\epsilon = 0.0$<br>$\alpha = 2.5$ | $\epsilon = 0.25$<br>$\alpha = 0.25$ | $\epsilon = 0.25$<br>$\alpha = 2.5$ |
|---|-------------------------------------|------------------------------------|--------------------------------------|-------------------------------------|
|  | 37.25                               | 4.86                               | 5.74                                 | 15.47                               |
|  | 37.43                               | 52.22                              | 24.40                                | 5.27                                |

Table 1: Mean Squared Error (MSE) score for the network, trained with **different classification** ground truth functions  $C_s^*$  and values for  $\epsilon$  and  $\alpha$ .

We apply a large variety of data augmentation, comparable to [20], including random color channel re-distribution, random brightness and contrast adjustments, random rotations by multiples of  $90^\circ$ , random scales between 0.5 and 1 and random crops to a patch size of  $225 \times 225$ . This patch size leverages the utilization of global information by the U-Net. Our training batches contain seven shifts of two stacks extracted from a single RGB light field ( $7 \cdot 2 \cdot 3 = 42$  channels).

### 3.6. Refinement

When choosing the rectangle classification ground truth (Equation 7) an additional refinement step can be performed. In case

$$\max_s (C_s(x, y)) < t, \quad (11)$$

meaning that no classification exceeds some small threshold  $t = 0.01$ , we assume that the chosen disparity label at  $(x, y)$  would probably be wrong. In order to smoothly fill this pixel, we apply a median filter to each classification output first.

## 4. Experiments

In this section we present the results of our evaluations.

### 4.1. Ablation Studies

Because the choice of a classification ground truth function is highly important for our method, we evaluated different functions and parameters. Table 1 shows the MSE score of our network, trained on either the rectangle function in Equation 7 or the triangle function in Equation 8. As expected, the results show that the triangle function requires a higher  $\epsilon$  to compensate for wrong classifications at boundaries of disparity labels. Due to the slightly better performance of the rectangle function, we chose it for our subsequent evaluations, setting  $\epsilon = 0.17$  and  $\alpha = 2.5$ .

We also evaluated our CNN architecture, without EPI-Shift, similar to [20]. This model only reached an MSE

score of 31.15 compared to 0.85 for our model with EPI-Shift. This clearly indicates that our U-Net architecture requires the shifted EPIs in order to properly generalize.

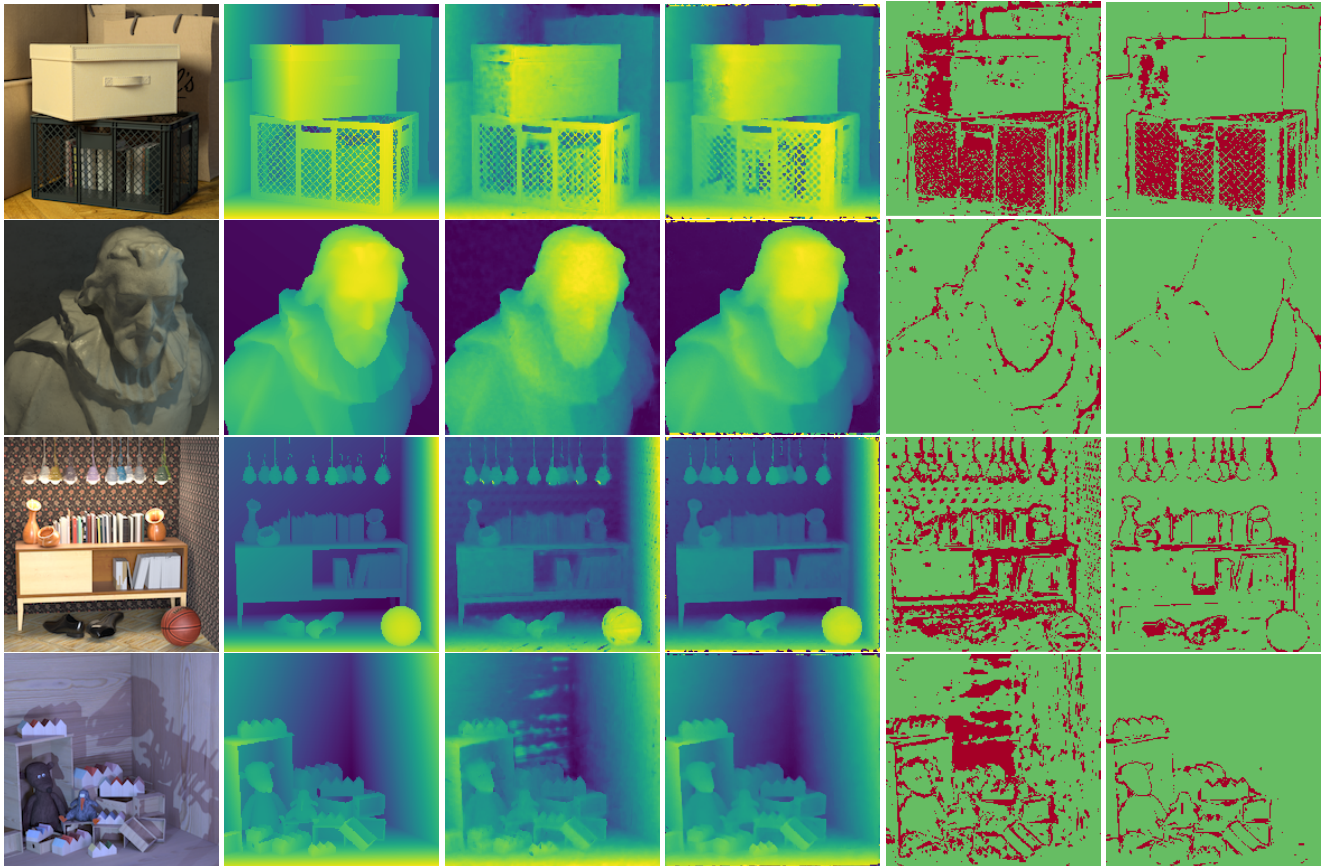
### 4.2. Results on the HCI 4D Light Field Benchmark

For the quantitative evaluation in Figure 7, we plot 13 error measures from the HCI 4D Light Field Benchmark, for details please refer to [10]. Our method outperforms EPINET-Cross [20], in 11 out of 13 metrics with a close tie of the other two. Because our network is run several times (once for each EPI-Shift), the runtime increases linearly with the disparity range and is therefore slightly higher compared to [20]. However, our approach is still significantly faster than the non-learning approaches. Our work closes the performance gap to optimization-based methods while keeping all advantages of deep learning like fewer hyper parameters and learned instead of hand-crafted heuristics. We evaluated our method on four photo-realistic scenes of the publicly available HCI 4D Light Field Benchmark [10] that were not part of the training dataset. Sadly, the benchmark uploading website is not functional anymore, hence no official submission was possible yet. In Figure 6 we show a qualitative comparison with EPINET [20]. Note, that we use the cross setup for EPINET which uses the same 17 views subset of the full light field that is used by our approach. In addition to EPINET, the quantitative evaluation in Figure 7 also includes two state-of-the-art optimization based methods (OBER [18] and SPO-MO [24]), ranking first and third in the official benchmark.

Our method provides considerably better quality at the disparity extremes (Scene 3 and 4, background) due to the improved generalization enabled by EPI-Shift. Also note the improved performance on non-textured surfaces (Scene 1, beige box) caused by the large receptive field of the U-Net which enables better smoothing and long-range reasoning. Unfortunately, the U-Net also seems to be more prone to noise at object boundaries which are not disparity label boundaries (compare Scene 1, dark box), although similar artifacts can be observed with EPINET (compare Scene 3, books).

### 4.3. Results on Real Recordings

We also evaluated our method on images recorded by a cross light field setup with a disparity range of  $[0, 12]$ , consisting of 17 cameras. As the authors of [20] did not provide us with the pre-trained parameters for the cross-version of EPINET upon request, we trained EPINET-Cross based on their implementation for four days. Figure 1 (left) shows one of the results. As expected, EPINET is only able to predict within the small training data disparity range of  $[-3.5, 3.5]$ , present in the background. In contrast, our EPI-Shift reconstructs the disparity in the whole range of  $[0, 12]$ .



(a) Center View (b) Ground Truth (c) EPINET (d) Ours (e) EPINET BadPix (f) Ours BadPix

Figure 6: **Results** on the HCI 4D Light Field Benchmark [10] compared to the best learning-based competitor EPINET [20]. The BadPix score in (e) and (f) shows all pixels (red) exceeding an absolute distance of 0.07 to the ground truth.

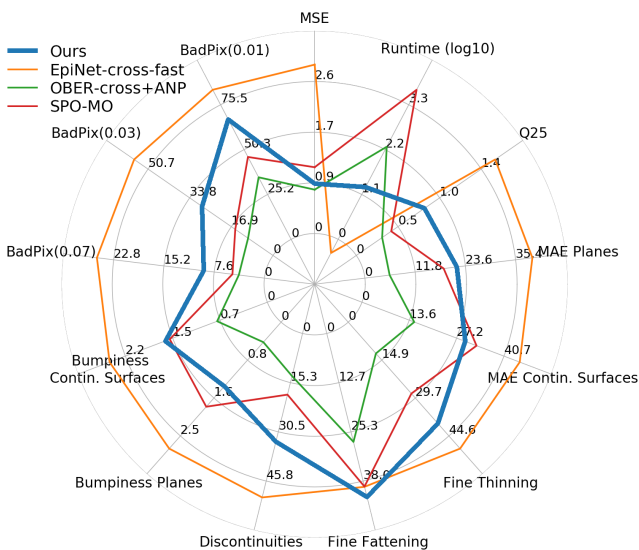


Figure 7: **Qualitative results on synthetic data.** We outperform EPINET-Cross [20] on 11 out of 13 metrics.

## 5. Conclusion

To summarize, we have introduced a new learning-based approach for depth estimation from wide-baseline light field recordings. The key idea of our approach is to use so-called EPI-Shifts, similar to plane sweep volumes for stereo depth estimation. This approach improves the generalization capability of CNN based depth estimation and enables us to increase the receptive field using a U-Net which delivers better smoothing and reduces artifacts, thanks to long range reasoning. Combining these two advantages leads to state-of-the-art performance, as demonstrated on a publicly available light field benchmark. Furthermore, the EPI-Shift concept enables depth estimation with large baseline light fields, while the training data only exhibits small disparities. We also demonstrate results on a real world recording.

## 6. Acknowledgement

We thank the Center for Information Services and High Performance Computing (ZIH) at TU Dresden for generous allocations of computer time.



## References

- [1] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACMTOG*, 28(3):24:1–24:11, 2009.
- [2] R. C. Bolles, H. H. Baker, and D. H. Marimont. Epipolar-plane image analysis: An approach to determining structure from motion. *International Journal of Computer Vision*, 1(1):7–55, Mar 1987.
- [3] R. T. Collins. A space-sweep approach to true multi-image matching. In *1996 Conference on Computer Vision and Pattern Recognition (CVPR '96), June 18-20, 1996 San Francisco, CA, USA*, pages 358–363. IEEE Computer Society, 1996.
- [4] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao. Deep ordinal regression network for monocular depth estimation. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 2002–2011. IEEE Computer Society, 2018.
- [5] X. Guo, H. Li, S. Yi, J. S. J. Ren, and X. Wang. Learning monocular depth by distilling cross-domain stereo networks. In V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, editors, *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XI*, volume 11215 of *Lecture Notes in Computer Science*, pages 506–523. Springer, 2018.
- [6] S. Heber and T. Pock. Shape from light field meets robust pca. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 751–767, Cham, 2014. Springer International Publishing.
- [7] S. Heber and T. Pock. Convolutional networks for shape from light field. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3746–3754, June 2016.
- [8] S. Heber, W. Yu, and T. Pock. U-shaped networks for shape from light field. In R. C. Wilson, E. R. Hancock, and W. A. P. Smith, editors, *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*. BMVA Press, 2016.
- [9] S. Heber, W. Yu, and T. Pock. Neural epi-volume networks for shape from light field. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2271–2279, Oct 2017.
- [10] K. Honauer, O. Johannsen, D. Kondermann, and B. Goldluecke. A dataset and evaluation methodology for depth estimation on 4d light fields. In *Asian Conference on Computer Vision*. Springer, 2016.
- [11] S. Im, H. Jeon, S. Lin, and I. S. Kweon. Dpsnet: End-to-end deep plane sweep stereo. *CoRR*, abs/1905.00538, 2019.
- [12] H.-G. Jeon, J. Park, G. Choe, J. Park, Y. Bok, Y.-W. Tai, and I. S. Kweon. Accurate depth map estimation from a lenslet light field camera. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1547–1555. IEEE, 2015.
- [13] H. Lin, C. Chen, S. B. Kang, and J. Yu. Depth recovery from light field using focal stack symmetry, 12 2015.
- [14] A. Neri, M. Carli, and F. Battisti. A multi-resolution approach to depth field estimation in dense image arrays, 09 2015.
- [15] R. Ng, M. Levoy, M. Brédif, G. Duval, M. Horowitz, and P. Hanrahan. Light field photography with a hand-held plenoptic camera. 2005.
- [16] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- [17] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In N. Navab, J. Hornegger, W. M. W. III, and A. F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III*, volume 9351 of *Lecture Notes in Computer Science*, pages 234–241. Springer, 2015.
- [18] H. Schilling, M. Diebold, C. Rother, and B. Jähne. Trust your model: Light field depth estimation with inline occlusion handling. In *CVPR*, 2018.
- [19] H. Sheng, P. Zhao, S. Zhang, J. Zhang, and D. Yang. Occlusion-aware depth estimation for light field using multi-orientation epis. 74, 09 2017.
- [20] C. Shin, H.-G. Jeon, Y. Yoon, I. S. Kweon, and S. J. Kim. Epinet: A fully-convolutional neural network using epipolar geometry for depth from light field images. In *CVPR*, 2018.
- [21] S. Tulyakov, A. Ivanov, and F. Fleuret. Practical deep stereo (PDS): toward applications-friendly deep stereo matching. In S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada.*, pages 5875–5885, 2018.
- [22] S. Wanner and B. Goldlücke. Variational light field analysis for disparity estimation and super-resolution. *IEEE Trans. Pattern Analysis Machine Intelligence*, 36:606–619, 2014.
- [23] B. S. Wilburn, M. Smulski, H.-H. K. Lee, and M. A. Horowitz. Light field video camera. In *Media Processors 2002*, volume 4674, pages 29–37. International Society for Optics and Photonics, 2001.
- [24] S. Zhang, H. Sheng, C. Li, J. Zhang, and Z. Xiong. Robust depth estimation for light field via spinning parallelogram operator. *Computer Vision and Image Understanding*, 145:148 – 159, 2016. Light Field for Computer Vision.

## A. Additional Results

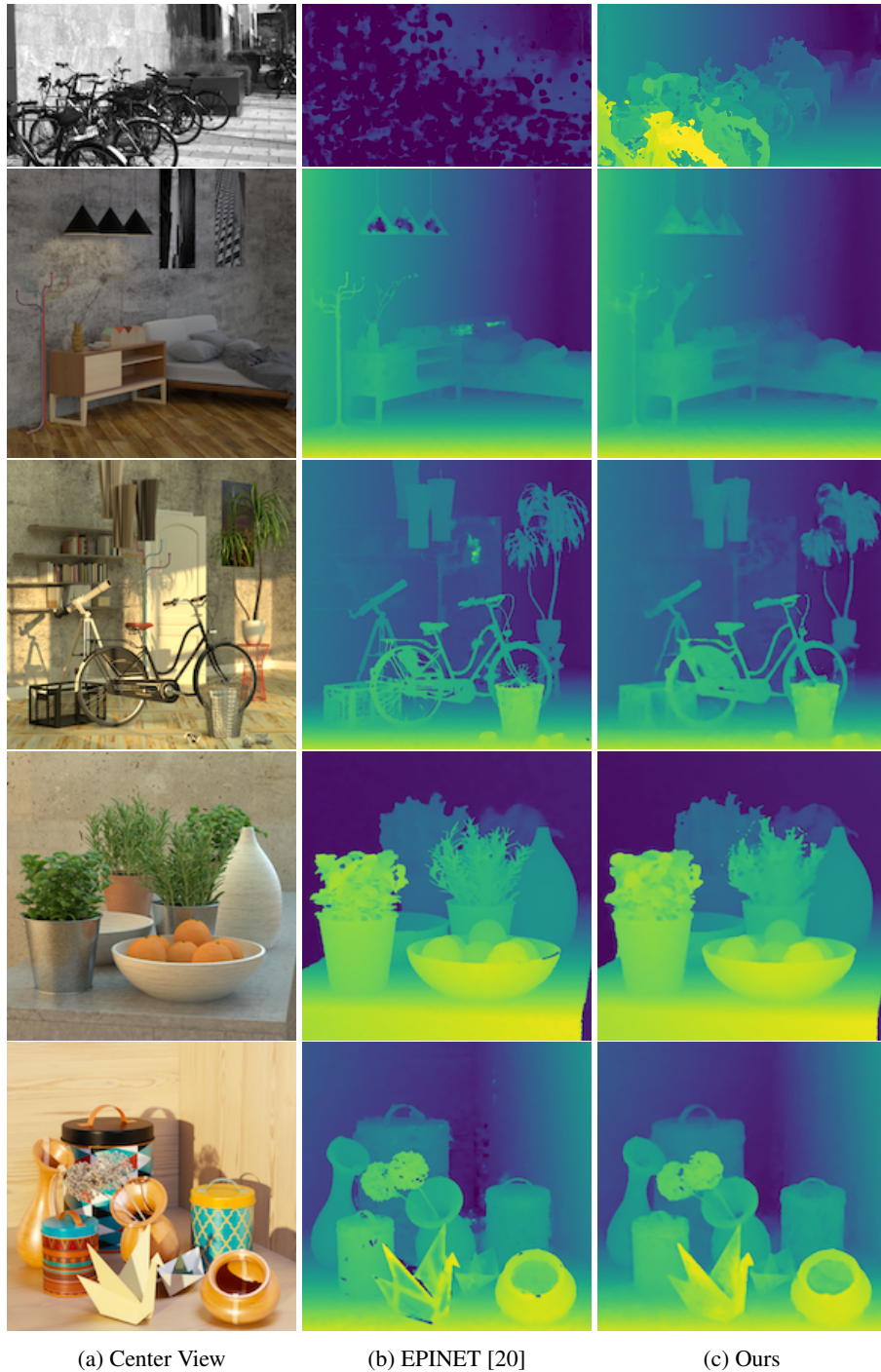


Figure A.1: **Results** on a second real scene (top) and four additional benchmark scenes [10] without publicly available ground truth. The **real recording** (top) shows the limitation of EPINET [20] to the disparity range of the training data. **Additional benchmark scenes** show an improvement in non-textured areas and at extreme disparities (compare Scene 4 (the last scene), background) but also slightly more blurry results of our method. However, blur only occurs within the small regression intervals if two objects are part of the same depth label. At extreme disparities (compare Scene 3, background), our method also performs better due to the hard transitions between adjacent labels.

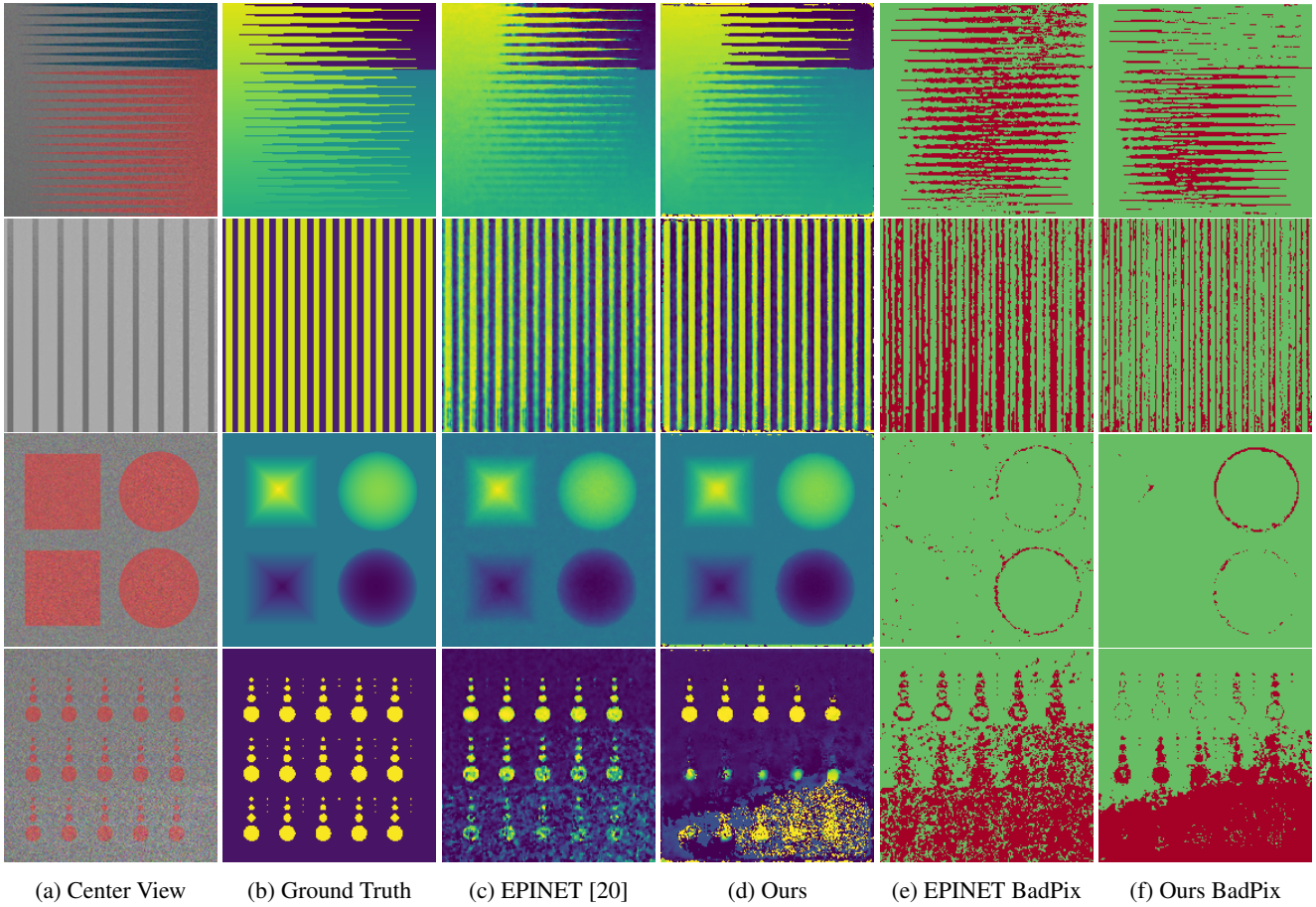


Figure A.2: **Results** on “Stratified” scenes of the HCI 4D Light Field Benchmark [10]. The BadPix score in (e) and (f) shows all pixels (red) exceeding an  $\mathcal{L}^1$ -distance of 0.07 to the ground truth. Note the **improved smoothness** on flat surfaces due to the U-Net architecture (compare Scene 3, background). Also note the **failure case** of our method in Scene 4, caused by strong noise occurring only in the bottom of the image. In those cases, EPI-Shift causes misclassifications, leading to stronger artifacts than EPINET [20].