

Defocus Map Estimation and Deblurring from a Single Dual-Pixel Image: Supplementary Material

Shumian Xin¹ Neal Wadhwa² Tianfan Xue² Jonathan T. Barron²
Pratul P. Srinivasan² Jiawen Chen³ Ioannis Gkioulekas¹ Rahul Garg²
¹Carnegie Mellon University ²Google Research ³Adobe Inc.

1. Introduction

In this supplementary material, we cover the following topics:

1. In Sec. 2, we describe blur kernel calibration in more detail, and explore how blur kernels change with respect to scene depth and focus distance.
2. In Sec. 3, we provide more technical details about our method. More specifically, we explain how we render defocus maps from the multi-plane image (MPI) representation, provide the derivation of the bias correction term, and define the total variation function $V(\cdot)$ and the edge map E used in the regularization terms.
3. In Sec. 4, we provide additional implementation details, and show comparison results and ablation studies on more data in our collected Google Pixel 4 dataset. **To facilitate comparisons, we also provide an interactive HTML viewer [2] at the project website [11].**

2. Blur Kernel Calibration

We provide more information about our calibration procedure for the left and right blur kernels used as input to our method. We use a method similar to the one proposed by Mannan and Langer [6], and calibrate blur kernels for left and right dual-pixel (DP) images independently (Fig. 1) for a specific focus distance. Specifically, we image a regular grid of circular discs on a monitor screen at a distance of ~ 45 cm from the camera. We apply global thresholding and binarize the captured image, perform connected component analysis to identify the individual discs and their centers, and generate and align the binary sharp image M with the known calibration pattern by solving for a homography between the calibration target disc centers and the detected centers. In order to apply radiometric correction, we also capture all-white and all-black images displayed on the same screen, and generate the grayscale latent sharp image as $I_l = M \odot I_w + (1 - M) \odot I_b$, where \odot represents pixel-wise multiplication, and I_w and I_b are captured all-white and all-black images. Once we have the aligned latent image and the captured image, we can solve for spatially-varying blur kernels using the optimization proposed by Mannan and Langer [6]. Specifically, we solve for a 8×6 grid of kernels corresponding to 1344×1008 central field of view.

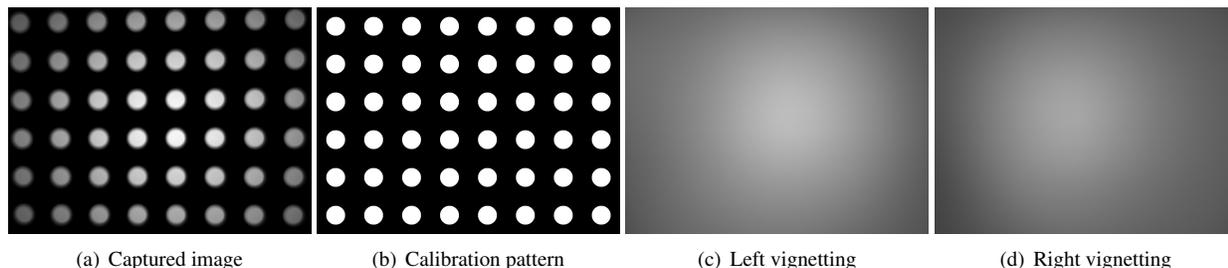


Figure 1: Captured image (a) of the calibration pattern (b) that is used to calibrate the blur kernels. Left DP image (c) and right DP image (d) of a white sheet shot through a diffuser that is used to correct for vignetting.

In addition to the blur kernels, we calibrate for different vignetting in left and right DP images. Specifically, for the same focus distance as above, we capture six images of a white sheet through a diffuser. We then average all left and right images individually to obtain the left and right vignetting patterns W_l and W_r , respectively.

Next, we explore how DP blur kernels change with respect to scene depth and focus distance (Fig. 2). As observed by Tang and Kutulakos [8], we find that kernels behave differently on the opposite sides of the focus plane. Therefore we choose focus settings such that all scene contents are at or behind the focus plane for all of our experiments, including this kernel analysis. We observe that DP blur kernels are approximately resized versions of each other as the scene depth or focus distance changes, similar to the expected behavior for blur kernels in a regular image sensor.

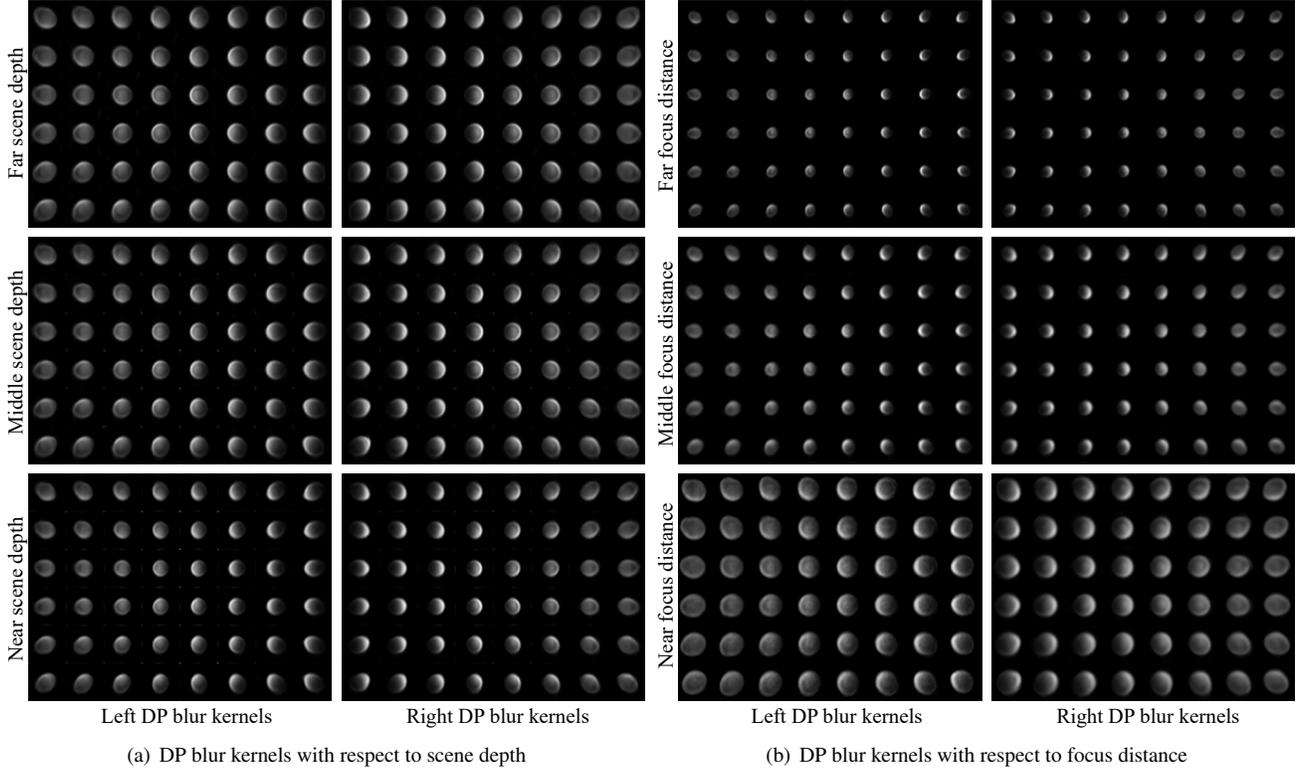


Figure 2: DP blur kernels with respect to scene depth and focus distance. We choose focus settings such that all scene contents are at or behind the focus plane, and calibrate for blur kernels either with the same focus settings but at different depths (a), or at the same depth but with various focus distances (b).

3. Additional Method Details

In this section, we provide more technical details about our method. We explain how we render defocus maps from the MPI representation in Sec. 3.1, provide the derivation of the bias correction term in Sec. 3.2, and finally define the total variation function $V(\cdot)$ and the edge map E used in the regularization terms in Sec. 3.3.

3.1. Defocus Map from MPI

We have shown in the main paper that an all-in-focus image can be composited from an MPI representation as:

$$\hat{\mathbf{I}}_s = \sum_{i=1}^N t_i c_i = \sum_{i=1}^N \left[c_i \alpha_i \prod_{j=i+1}^N (1 - \alpha_j) \right]. \quad (1)$$

We can synthesize a continuous-valued defocus map $\hat{\mathbf{D}}$ in a similar way as discussed by Tucker and Snavely [9], by replacing all pixel intensities in Eq. (1) with the defocus blur size d_i of that layer:

$$\hat{\mathbf{D}} = \sum_{i=1}^N \left[d_i \alpha_i \prod_{j=i+1}^N (1 - \alpha_j) \right]. \quad (2)$$

3.2. Proof of Eq. (4) of the Main Paper

In this section, we provide a detailed derivation of the bias correction term. To be self-contained, we restate our assumed image formation model. Given an MPI representation, its corresponding DP images can be expressed as:

$$\mathbf{I}_o^{\{l,r\}} = \mathbf{I}_b^{\{l,r\}} + \mathbf{N}^{\{l,r\}}, \quad (3)$$

where $\mathbf{I}_b^{\{l,r\}}$ are the latent noise-free left and right defocused images, and $\mathbf{N}^{\{l,r\}}$ is additive white Gaussian noise with entries independent identically distributed with distribution $\mathcal{N}(0, \sigma^2)$. Our goal is to optimize for an MPI with intensity-alpha layers $(\hat{c}_i, \hat{\alpha}_i)$, with defocus sizes $d_i, i \in [1, \dots, N]$, such that the L_2 loss $\|\hat{\mathbf{I}}_b^{\{l,r\}} - \mathbf{I}_o^{\{l,r\}}\|_2^2$ is minimized. We show that, in the presence of image noise, minimizing the above loss biases the estimated defocus map towards smaller blur values. Specifically, we quantify this bias and then correct for it in our optimization.

For simplicity, we assume for now that all scene contents lie on a single fronto-parallel plane with ground truth defocus size d^* , and our scene representation is an MPI with a single opaque layer (i.e., $\hat{\alpha}_i = \mathbf{1}$) with a defocus size hypothesis d_i . Under this assumption, the defocused image rendering equation (Eq. (2) of the main paper)

$$\hat{\mathbf{I}}_b^{\{l,r\}} = \sum_{i=1}^N \left[\left(\mathbf{k}_{d_i}^{\{l,r\}} * (\mathbf{c}_i \alpha_i) \right) \odot \prod_{j=i+1}^N \left(\mathbf{1} - \mathbf{k}_{d_j}^{\{l,r\}} * \alpha_j \right) \right] \quad (4)$$

reduces to

$$\hat{\mathbf{I}}_b^{\{l,r\}} = \mathbf{k}_{d_i}^{\{l,r\}} * \hat{\mathbf{c}}_i. \quad (5)$$

Similarly, Eq. (3) becomes:

$$\mathbf{I}_o^{\{l,r\}} = \mathbf{k}_{d^*}^{\{l,r\}} * \mathbf{c}_i + \mathbf{N}^{\{l,r\}}. \quad (6)$$

We can express the above equations in the frequency domain as follows:

$$\mathcal{I}_o^{\{l,r\}} = \mathbf{K}_{d^*}^{\{l,r\}} \mathbf{C}_i + \mathcal{N}^{\{l,r\}}, \quad (7)$$

where $\mathcal{I}_o^{\{l,r\}}, \mathbf{K}_{d^*}^{\{l,r\}}, \mathbf{C}_i$, and $\mathcal{N}^{\{l,r\}}$ are the Fourier transforms of $\mathbf{I}_o^{\{l,r\}}, \mathbf{k}_{d^*}^{\{l,r\}}, \mathbf{c}_i$, and $\mathbf{N}^{\{l,r\}}$, respectively. Note that the entries of $\mathcal{N}^{\{l,r\}}$ are also independent identically distributed with the same Gaussian distribution $\mathcal{N}(0, \sigma^2)$ as the entries of $\mathbf{N}^{\{l,r\}}$.

We can obtain a maximum a posteriori (MAP) estimate of $\hat{\mathbf{C}}_i$ and d_i by solving the following optimization problem [12]:

$$\begin{aligned} & \arg \max P \left(\mathcal{I}_o^l, \mathcal{I}_o^r | \hat{\mathbf{C}}_i, d_i, \sigma \right) P \left(\hat{\mathbf{C}}_i, d_i \right) \\ & = \arg \max P \left(\mathcal{I}_o^l, \mathcal{I}_o^r | \hat{\mathbf{C}}_i, d_i, \sigma \right) P \left(\hat{\mathbf{C}}_i \right). \end{aligned} \quad (8)$$

According to Eq. (7), we have

$$P \left(\mathcal{I}_o^l, \mathcal{I}_o^r | \hat{\mathbf{C}}_i, d_i, \sigma \right) \propto \exp \left(-\frac{1}{2\sigma^2} \sum_{v=\{l,r\}} \|\mathbf{K}_{d_i}^v \hat{\mathbf{C}}_i - \mathcal{I}_o^v\|^2 \right). \quad (9)$$

We also follow Zhou et al. [12] in assuming a prior for the latent all-in-focus image such that:

$$P \left(\hat{\mathbf{C}}_i \right) \propto \exp \left(-\frac{1}{2} \|\Phi \hat{\mathbf{C}}_i\|^2 \right), \quad (10)$$

where we define Φ such that

$$|\Phi(f)|^2 = \frac{1}{|\hat{C}_i(f)|^2}, \quad (11)$$

and f is the frequency. As \hat{C}_i is the unknown variable, we approximate Eq. (11) by averaging the power spectrum over a set of natural images $\{C_i\}$:

$$|\Phi(f)|^2 = \frac{1}{\int_{C_i} |C_i(f)|^2 \mu(C_i)}, \quad (12)$$

where $\mu(C_i)$ represents the probability distribution of C_i in image domain.

Maximizing the log-likelihood of Eq. (8) is equivalent to minimizing the following loss:

$$E(d_i | \mathcal{I}_o^l, \mathcal{I}_o^r, \sigma) = \min_{\hat{C}_i} \left(\sum_{v=\{l,r\}} \|K_{d_i}^v \hat{C}_i - \mathcal{I}_o^v\|^2 \right) + \|\sigma \Phi \hat{C}_i\|^2. \quad (13)$$

d_i can be estimated as the minimizer of the above energy function. Then given d_i , setting $\partial E / \partial \hat{C}_i = 0$ yields the following solution of \hat{C}_i , known as a *generalized Wiener deconvolution with two observations*:

$$\hat{C}_i = \frac{\mathcal{I}_o^l \overline{K_{d_i}^l} + \mathcal{I}_o^r \overline{K_{d_i}^r}}{|K_{d_i}^l|^2 + |K_{d_i}^r|^2 + \sigma^2 |\Phi|^2}, \quad (14)$$

where $\overline{K_{d_i}^{\{l,r\}}}$ is the complex conjugate of $K_{d_i}^{\{l,r\}}$, and $|K_{d_i}^{\{l,r\}}|^2 = K_{d_i}^{\{l,r\}} \overline{K_{d_i}^{\{l,r\}}}$.

We then evaluate the defocus size hypothesis d_i by computing the minimization loss given the latent ground truth depth d^* , and the noise level σ , that is,

$$E(d_i | K_{d^*}^l, K_{d^*}^r, \sigma) = \mathbb{E}_{C_i, \mathcal{I}_o^l, \mathcal{I}_o^r} E(d_i | K_{d^*}^l, K_{d^*}^r, \sigma, C_i, \mathcal{I}_o^l, \mathcal{I}_o^r) \quad (15)$$

$$= \mathbb{E}_{C_i, \mathcal{I}_o^l, \mathcal{I}_o^r} \left[\left(\sum_{v=\{l,r\}} \|K_{d_i}^v \hat{C}_i - \mathcal{I}_o^v\|^2 \right) + \|\sigma \Phi \hat{C}_i\|^2 \right], \quad (16)$$

where $\mathbb{E}(\cdot)$ is the expectation. Substituting \hat{C}_i with Eq. (14) gives us:

$$\begin{aligned} & E(d_i | K_{d^*}^l, K_{d^*}^r, \sigma) \\ &= \mathbb{E}_{C_i, \mathcal{I}_o^l, \mathcal{I}_o^r} \left[\left(\sum_{v=\{l,r\}} \|K_{d_i}^v \frac{\mathcal{I}_o^l \overline{K_{d_i}^l} + \mathcal{I}_o^r \overline{K_{d_i}^r}}{|K_{d_i}^l|^2 + |K_{d_i}^r|^2 + \sigma^2 |\Phi|^2} - \mathcal{I}_o^v\|^2 \right) + \|\sigma \Phi \frac{\mathcal{I}_o^l \overline{K_{d_i}^l} + \mathcal{I}_o^r \overline{K_{d_i}^r}}{|K_{d_i}^l|^2 + |K_{d_i}^r|^2 + \sigma^2 |\Phi|^2}\|^2 \right]. \end{aligned} \quad (17)$$

Then substituting \mathcal{I}_o^v with Eq. (7), we get:

$$\begin{aligned} & E(d_i | K_{d^*}^l, K_{d^*}^r, \sigma) \\ &= \mathbb{E}_{C_i, \mathcal{N}^l, \mathcal{N}^r} \left[\left(\sum_{v=\{l,r\}} \|K_{d_i}^v \frac{(K_{d^*}^v C_i + \mathcal{N}^v) \overline{K_{d_i}^v} + (K_{d^*}^v C_i + \mathcal{N}^v) \overline{K_{d_i}^v}}{|K_{d_i}^l|^2 + |K_{d_i}^r|^2 + \sigma^2 |\Phi|^2} - (K_{d^*}^v C_i + \mathcal{N}^v)\|^2 \right) + \right. \\ & \quad \left. \|\sigma \Phi \frac{(K_{d^*}^l C_i + \mathcal{N}^l) \overline{K_{d_i}^l} + (K_{d^*}^r C_i + \mathcal{N}^r) \overline{K_{d_i}^r}}{|K_{d_i}^l|^2 + |K_{d_i}^r|^2 + \sigma^2 |\Phi|^2}\|^2 \right]. \end{aligned} \quad (18)$$

We now define $B = |\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2 + \sigma^2 |\Phi|^2$. We can rearrange the above equation as:

$$\begin{aligned}
& E(d_i | \mathbf{K}_{d_i}^l, \mathbf{K}_{d_i}^r, \sigma) \\
&= \mathbb{E}_{C_i, \mathcal{N}^l, \mathcal{N}^r} \left[\left(\sum_{v=\{l,r\}} \left\| \frac{C_i \left[\mathbf{K}_{d_i}^v \left(\mathbf{K}_{d_i}^l \overline{\mathbf{K}_{d_i}^l} + \mathbf{K}_{d_i}^r \overline{\mathbf{K}_{d_i}^r} \right) - \mathbf{K}_{d_i}^v B \right]}{B} + \frac{\mathbf{K}_{d_i}^v \left(\mathcal{N}^l \overline{\mathbf{K}_{d_i}^l} + \mathcal{N}^r \overline{\mathbf{K}_{d_i}^r} \right)}{B} - \mathcal{N}^v \right\|^2 \right) + \right. \\
& \quad \left. \left\| \sigma \Phi \frac{C_i \left(\mathbf{K}_{d_i}^l \overline{\mathbf{K}_{d_i}^l} + \mathbf{K}_{d_i}^r \overline{\mathbf{K}_{d_i}^r} \right)}{B} + \sigma \Phi \frac{\mathcal{N}^l \overline{\mathbf{K}_{d_i}^l} + \mathcal{N}^r \overline{\mathbf{K}_{d_i}^r}}{B} \right\|^2 \right]. \tag{19}
\end{aligned}$$

Given that the entries of $\mathcal{N}^{\{l,r\}}$ are independent identically distributed with distribution $\mathcal{N}(0, \sigma^2)$, we have $\mathbb{E}(\mathcal{N}^v) = \mathbf{0}$, $\mathbb{E}(\mathcal{N}^{v^2}) = \sigma^2$ and $\mathbb{E}(\mathcal{N}^l \mathcal{N}^r) = \mathbf{0}$, and we can simplify the above equation as:

$$\begin{aligned}
& E(d_i | \mathbf{K}_{d_i}^l, \mathbf{K}_{d_i}^r, \sigma) \\
&= \mathbb{E}_{C_i, \mathcal{N}^l, \mathcal{N}^r} \left[\left(\sum_{v=\{l,r\}} \left\| \frac{C_i \left[\mathbf{K}_{d_i}^v \left(\mathbf{K}_{d_i}^l \overline{\mathbf{K}_{d_i}^l} + \mathbf{K}_{d_i}^r \overline{\mathbf{K}_{d_i}^r} \right) - \mathbf{K}_{d_i}^v B \right]}{B} \right\|^2 + \left\| \frac{\mathbf{K}_{d_i}^v \left(\mathcal{N}^l \overline{\mathbf{K}_{d_i}^l} + \mathcal{N}^r \overline{\mathbf{K}_{d_i}^r} \right)}{B} - \mathcal{N}^v \right\|^2 \right) + \right. \\
& \quad \left. \left\| \sigma \Phi \frac{C_i \left(\mathbf{K}_{d_i}^l \overline{\mathbf{K}_{d_i}^l} + \mathbf{K}_{d_i}^r \overline{\mathbf{K}_{d_i}^r} \right)}{B} \right\|^2 + \left\| \sigma \Phi \frac{\mathcal{N}^l \overline{\mathbf{K}_{d_i}^l} + \mathcal{N}^r \overline{\mathbf{K}_{d_i}^r}}{B} \right\|^2 \right] \tag{20}
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{C_i} \left\{ \left[\sum_{v=\{l,r\}} \left\| \frac{C_i \left[\mathbf{K}_{d_i}^v \left(\mathbf{K}_{d_i}^l \overline{\mathbf{K}_{d_i}^l} + \mathbf{K}_{d_i}^r \overline{\mathbf{K}_{d_i}^r} \right) - \mathbf{K}_{d_i}^v B \right]}{B} \right\|^2 + \sigma^2 \left(\left\| \frac{\mathbf{K}_{d_i}^v}{B} \right\|^2 + \left\| \frac{\mathbf{K}_{d_i}^l \mathbf{K}_{d_i}^r}{B} \right\|^2 \right) \right] + \right. \\
& \quad \left. \left\| \sigma \Phi \frac{C_i \left(\mathbf{K}_{d_i}^l \overline{\mathbf{K}_{d_i}^l} + \mathbf{K}_{d_i}^r \overline{\mathbf{K}_{d_i}^r} \right)}{B} \right\|^2 + \sigma^2 \left(\left\| \sigma \Phi \frac{\mathbf{K}_{d_i}^l}{B} \right\|^2 + \left\| \sigma \Phi \frac{\mathbf{K}_{d_i}^r}{B} \right\|^2 \right) \right\}. \tag{21}
\end{aligned}$$

Recall that, in Eq. (12), we defined $\Phi(f)$ such that $\frac{1}{|\Phi(f)|^2} = \int_{C_i} |C_i(f)|^2 \mu(C_i)$. Then we can further simplify $E(d_i | \mathbf{K}_{d_i}^l, \mathbf{K}_{d_i}^r, \sigma)$ as:

$$\begin{aligned}
& E(d_i | \mathbf{K}_{d_i}^l, \mathbf{K}_{d_i}^r, \sigma) \\
&= \sum_f \left[\frac{\frac{1}{|\Phi|^2} |\mathbf{K}_{d_i}^l \mathbf{K}_{d_i}^r - \mathbf{K}_{d_i}^r \mathbf{K}_{d_i}^l|^2}{B} \right] + \sum_f \left[\frac{\frac{1}{|\Phi|^2} \sigma^2 |\Phi|^2 \left(|\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2 \right)}{B} \right] + \\
& \quad \sum_f \left[\sigma^2 \left(\left\| \frac{\mathbf{K}_{d_i}^l}{B} \right\|^2 + \left\| \frac{\mathbf{K}_{d_i}^r}{B} \right\|^2 + 2 \left\| \frac{\mathbf{K}_{d_i}^l \mathbf{K}_{d_i}^r}{B} \right\|^2 + \left\| \sigma \Phi \frac{\mathbf{K}_{d_i}^l}{B} \right\|^2 + \left\| \sigma \Phi \frac{\mathbf{K}_{d_i}^r}{B} \right\|^2 \right) \right] \tag{22}
\end{aligned}$$

$$= \sum_f \left[\frac{\frac{1}{|\Phi|^2} |\mathbf{K}_{d_i}^l \mathbf{K}_{d_i}^r - \mathbf{K}_{d_i}^r \mathbf{K}_{d_i}^l|^2}{B} \right] + \sigma^2 \sum_f \left[\frac{|\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2}{B} + \frac{\sigma^2 |\Phi|^2}{B} + 1 \right] \tag{23}$$

$$= \sum_f \left[\frac{\frac{1}{|\Phi|^2} |\mathbf{K}_{d_i}^l \mathbf{K}_{d_i}^r - \mathbf{K}_{d_i}^r \mathbf{K}_{d_i}^l|^2}{|\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2 + \sigma^2 |\Phi|^2} \right] + \sigma^2 \sum_f \left[\frac{|\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2 + \sigma^2 |\Phi|^2}{|\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2 + \sigma^2 |\Phi|^2} + 1 \right]. \tag{24}$$

If we define $C_1(\mathbf{K}_{d_i}^{\{l,r\}}, \sigma, \Phi) = \frac{1}{|\mathbf{K}_{d_i}^l|^2 + |\mathbf{K}_{d_i}^r|^2 + \sigma^2 |\Phi|^2}$, and $C_2(\sigma) = \sigma^2 \sum_f 1$, then Eq. (24) boils down to Eq. (4) of the main paper.

3.3. Edge-aware Total Variation Function

We first define a pixel-wise total variation function of a single-layer image \mathbf{I} that is used in both the intensity smoothness prior $\mathcal{L}_{\text{intensity}}$ and the alpha and transmittance smoothness prior $\mathcal{L}_{\text{alpha}}$:

$$V(\mathbf{I}) = \sqrt{\mathbf{I}^2 * g - (\mathbf{I} * g)^2}, \quad (25)$$

where g is a two-dimensional Gaussian blur kernel:

$$g = \begin{bmatrix} 1/16 & 1/8 & 1/16 \\ 1/8 & 1/4 & 1/8 \\ 1/16 & 1/8 & 1/16 \end{bmatrix}. \quad (26)$$

Each ‘‘pixel’’ in $V(\mathbf{I})(x, y)$ is equivalent to, for the 3×3 window surrounding pixel (x, y) in \mathbf{I} , computing the sample standard deviation (weighted by a Gaussian kernel) of the pixel intensities in that window. This follows easily from two facts: 1) as g sums to 1 by construction, $\mathbf{I} * g$ produces an image whose pixel intensities can be viewed as expectations of their surrounding 3×3 input patch; and 2) the standard deviation $\sqrt{\mathbb{E}[(X - \mathbb{E}[X])^2]}$ can be written equivalently as $\sqrt{\mathbb{E}[X^2] - \mathbb{E}[X]^2}$.

As done in prior work [9], we would like to encourage edge-aware smoothness in addition to minimizing total variation, so a bilateral edge mask is computed using this total variation:

$$\mathbf{E}(\mathbf{I}) = \mathbf{1} - \exp\left(-\frac{\mathbf{I}^2 * g - (\mathbf{I} * g)^2}{2\beta^2}\right). \quad (27)$$

In this equation, β is set to $1/32$ (assuming pixel intensities are in $[0, 1]$). A joint total variation function that takes into account both the original and the edge-aware total variation is then defined as:

$$V_{\mathbf{E}}(\mathbf{I}, \mathbf{E}) = \ell(V(\mathbf{I})) + (1 - \mathbf{E}) \odot \ell(V(\mathbf{I})). \quad (28)$$

4. Additional Implementation Details and Experimental Results

We first discuss more implementation details about our method, then show qualitative results on more data in our collected Google Pixel 4 dataset in Fig. 4-6. We also provide an interactive HTML viewer [2] at the project website [11] to facilitate the comparisons.

Data normalization. Before running the optimization, we first compute an intensity scaling factor $s = 0.5 / \text{mean}(\mathbf{I}_o^{\{l, r\}})$, and normalize the inputs $\bar{\mathbf{I}}_o^{\{l, r\}} = s\mathbf{I}_o^{\{l, r\}}$ to account for global intensity variations. After optimization, we undo the normalization by dividing the all-in-focus image by s .

Scaling factors of each loss term. Recall that our optimization loss is

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{data}} + \lambda_2 \mathcal{L}_{\text{aux}} + \lambda_3 \mathcal{L}_{\text{intensity}} + \lambda_4 \mathcal{L}_{\text{alpha}} + \lambda_5 \mathcal{L}_{\text{entropy}}. \quad (29)$$

$\mathcal{L}_{\text{data}}$ and \mathcal{L}_{aux} have the same weight: $\lambda_1 = \lambda_2 = 2.5 \cdot 10^4$. For most scenes, $\lambda_3 = 30$, $\lambda_4 = 7.5 \cdot 10^4$, and $\lambda_5 = 12$. We set higher weights on the regularization terms $\mathcal{L}_{\text{intensity}}$, $\mathcal{L}_{\text{alpha}}$, and $\mathcal{L}_{\text{entropy}}$ for scenes with less texture, e.g., data from Abuolaim and Brown [1].

Kernel size of each MPI layer. We manually determine the kernel sizes of the front and back layers, and evenly distribute MPI layers in disparity space. As mentioned in Sec. 2, we choose focus settings such that all scene contents are at or behind the focus plane. Therefore the kernel size of the front MPI layer is usually set to a small positive number, e.g., in the range of 1×1 to 3×3 , to mimic a 2D delta function, while the kernel size corresponding to the back MPI layer is set to a large enough value, e.g., in the range of 57×57 to 61×61 , to represent blur kernels at infinity.

Quantitative Metrics. We use the commercial software, Helicon Focus [4] to compute the ground truth all-in-focus images and the defocus maps using focus stacking. There may be a small *shift* between the ground truth all-in-focus image and the all-in-focus image from the deblurring algorithms we evaluate. This is because one can apply an arbitrary transform to the blur kernels and an inverse transform to the recovered all-in-focus image to yield the same blurred image. We determine this shift for each algorithm by using OpenCV to align the ground-truth all-in-focus image with the all-in-focus image from the algorithm via an affine transform for a single specific scene, and then using that transform to align all images before

computing the metrics for all-in-focus images. We also crop a small border of 8 pixels before computing the metrics as it may contain invalid pixels after alignment.

Additional results for bias correction. Fig. 3 shows additional results for our ablation study. Specifically, it shows that without bias correction term \mathcal{B} , the estimated defocus size is smaller on average as predicted by our analysis.

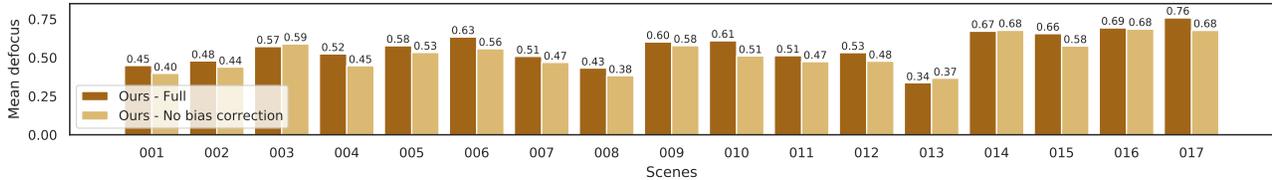
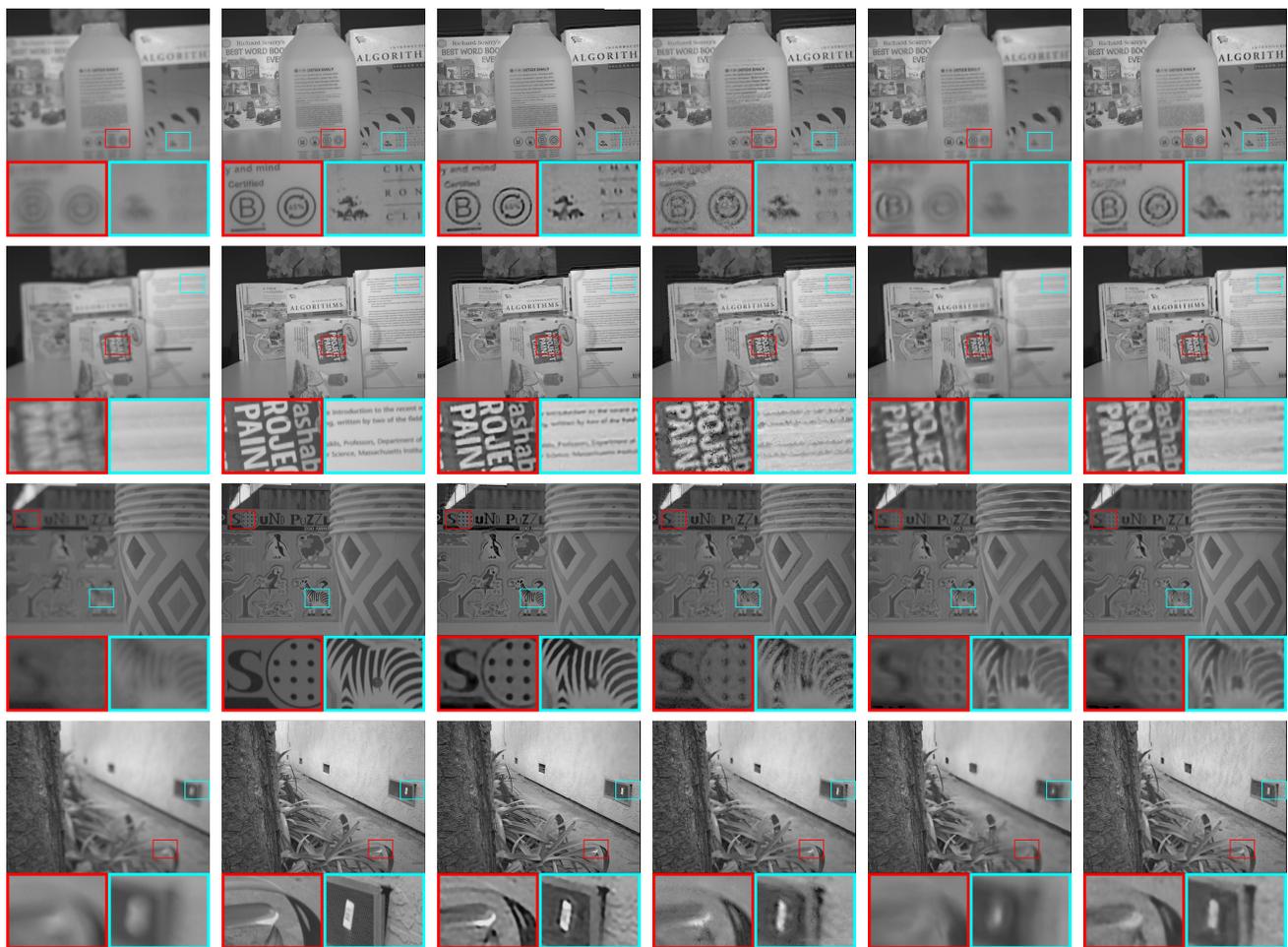


Figure 3: Mean of the predicted defocus map for our full pipeline vs an ablation where bias correction term is not applied. Defocus is measured as the relative scaling applied to the calibrated kernels. Without bias correction, the mean defocus is lower in 14 of the 17 scenes, i.e., the prediction is biased towards smaller defocus size.

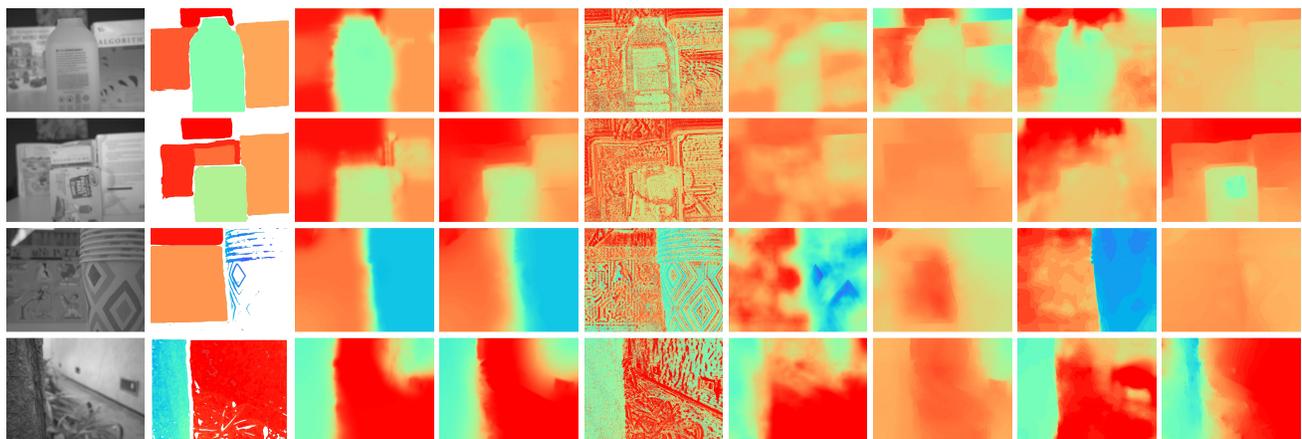
References

- [1] Abdullah Abuolaim and Michael S. Brown. Defocus deblurring using dual-pixel data. *European Conference on Computer Vision*, 2020. 6, 8
- [2] Benedikt Bitterli, Wenzel Jakob, Jan Novák, and Wojciech Jarosz. Reversible jump metropolis light transport using inverse mappings. *ACM Transactions on Graphics*, 2017. 1, 6
- [3] Rahul Garg, Neal Wadhwa, Sameer Ansari, and Jonathan T. Barron. Learning single camera depth estimation using dual-pixels. *IEEE/CVF International Conference on Computer Vision*, 2019. 8
- [4] Helicon focus. <https://www.heliconsoft.com/>. 6
- [5] Junyong Lee, Sungkil Lee, Sunghyun Cho, and Seungyong Lee. Deep defocus map estimation using domain adaptation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 8
- [6] Fahim Mannan and Michael S. Langer. Blur calibration for depth from defocus. *Conference on Computer and Robot Vision*, 2016. 1
- [7] Abhijith Punnappurath, Abdullah Abuolaim, Mahmoud Afifi, and Michael S. Brown. Modeling defocus-disparity in dual-pixel sensors. *IEEE International Conference on Computational Photography*, 2020. 8
- [8] Huixuan Tang and Kiriakos N. Kutulakos. Utilizing optical aberrations for extended-depth-of-field panoramas. *Asian Conference on Computer Vision*, 2012. 2
- [9] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 3, 6
- [10] Neal Wadhwa, Rahul Garg, David E. Jacobs, Bryan E. Feldman, Nori Kanazawa, Robert Carroll, Yair Movshovitz-Attias, Jonathan T. Barron, Yael Pritch, and Marc Levoy. Synthetic depth-of-field with a single-camera mobile phone. *ACM Transactions on Graphics*, 2018. 8
- [11] Shumian Xin, Neal Wadhwa, Tianfan Xue, Jonathan T. Barron, Pratul P. Srinivasan, Jianwen Chen, Ioannis Gkioulekas, and Rahul Garg. Project website, 2021. https://imaging.cs.cmu.edu/dual_pixels. 1, 6
- [12] Changyin Zhou, Stephen Lin, and Shree K. Nayar. Coded aperture pairs for depth from defocus. *IEEE/CVF International Conference on Computer Vision*, 2009. 3, 8



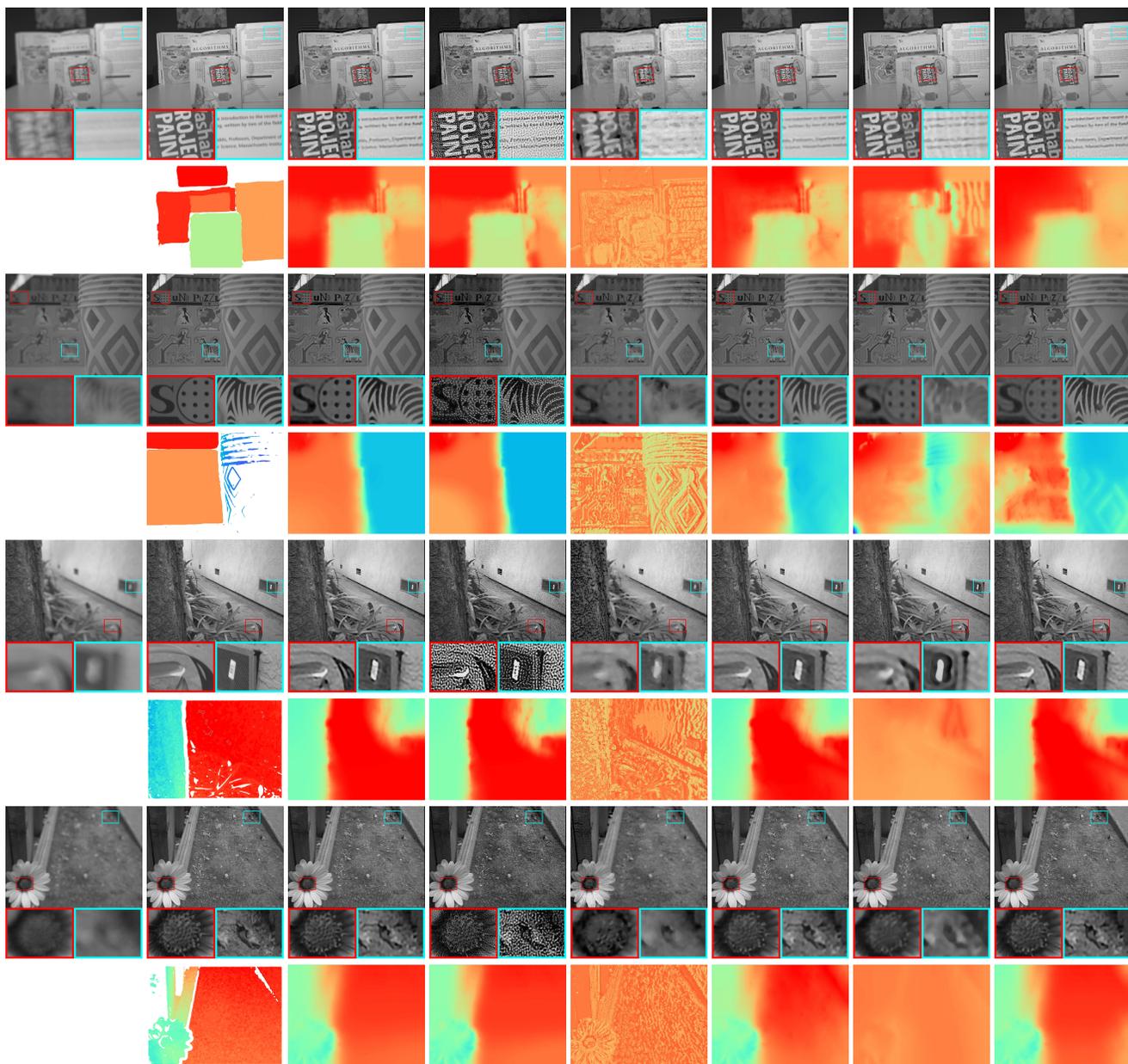
(a) Input image (b) GT all-in-focus image (c) Ours (d) Wiener deconv. [12] (e) DPDNet [1] (f) DPDNet w/ calib [1]

Figure 4: More qualitative comparisons of various defocus deblurring methods.



(a) Input image (b) Ground Truth (c) Ours (d) Ours w/ GF (e) Wiener [12] (f) DMENet [5] (g) [7] (h) Garg [3] (i) Wadhwa [10]

Figure 5: More qualitative comparisons of defocus map estimation methods.



(a) Input image (b) Ground truth (c) Ours full (d) No $\mathcal{L}_{intensity}$ (e) No \mathcal{L}_{alpha} (f) No $\mathcal{L}_{entropy}$ (g) No \mathcal{L}_{aux} (h) No \mathcal{B}

Figure 6: More qualitative results on ablation study.