# Cross-Camera Deep Colorization

Yaping Zhao[1,2], Haitian Zheng[3], Mengqi Ji[4], and Ruqi Huang[1,*]

[1]Tsinghua-Berkeley Shenzhen Institute and Tsinghua Shenzhen International Graduate School,
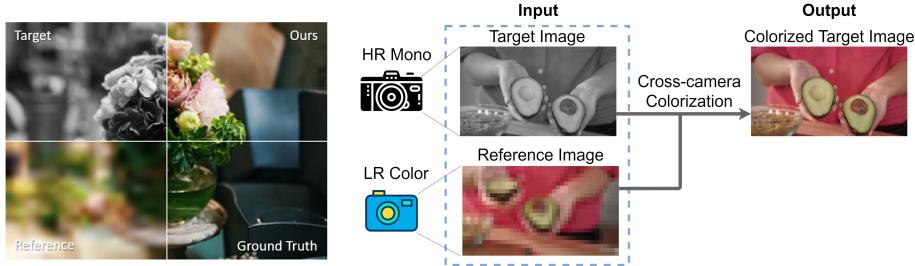[2]Zhejiang Future Technology Institute (Jiaxing), [3]University of Rochester, [4]Beihang University

Fig. 1: Aiming at pursuing high-resolution and high color depth simultaneously, we proposed a general binocular imaging framework that performs cross-camera colorization.

**Abstract.** In this paper, we consider the color-plus-mono dual-camera system and propose an end-to-end convolutional neural network to align and fuse images from it in an efficient and cost-effective way. Our method takes cross-domain and cross-scale images as input, and consequently synthesizes HR colorization results to facilitate the trade-off between spatial-temporal resolution and color depth in the single-camera imaging system. In contrast to the previous colorization methods, ours can adapt to color and monochrome cameras with distinctive spatial-temporal resolutions, rendering the flexibility and robustness in practical applications. The key ingredient of our method is a cross-camera alignment module that generates multi-scale correspondences for cross-domain image alignment. Through extensive experiments on various datasets and multiple settings, we validate the flexibility and effectiveness of our approach. Remarkably, our method consistently achieves substantial improvements, *i.e.*, around 10dB PSNR gain, upon the state-of-the-art methods. Code is at: github.com/IndigoPurple/CCDC.

**Keywords:** image colorization · image fusion · computational imaging.

## 1 Introduction

Nowadays, it has become a common practice of leveraging fusion of multi-sensory data from camera arrays to improve imaging quality [41,10,1,33,4,24,46,55,26,59]. In this paper, we consider the color-plus-mono dual-camera system and propose a novel learning-based framework for data fusion. Such a setting enjoys the advantages of photosensibility and high resolution of the monochrome camera, and the color sensibility of the color camera simultaneously. More specifically, our pipeline takes as inputs a

high-resolution (HR) grayscale image and a low-resolution (LR) color image, which are taken of the same scene by the respective cameras from similar viewpoints, as shown on the left of Fig. 1. After fusion, we obtain a HR color image illustrated on the right of Fig. 1, which to some extent facilitates the trade-off between spatial-temporal resolution and color depth in the single camera due to the space-bandwidth-product (SBP) [31].

In fact, the color-plus-mono dual camera has attracted an increasing amount of attention [23,11,29,34,38] recently. Most of the existing methods focus on addressing the cross-domain image color transfer. For example, traditional methods [49,35,20,40,14,9,30,3] employ global color statistics or low-level feature correspondences while emerging learning-based approaches [27,17] find high-level feature correspondences between images. Nevertheless, they commonly assume similar spatial resolution between data, and therefore can only deal with minor resolution gap, which is typically less than $2\times$.

On the other hand, it is naturally desirable to retain both high spatial-temporal resolution and high color depth. Thus methods assuming low-resolution gap have to either sacrifice resolution or require high-resolution input from *both* cameras, which can be costly and computationally heavy. Such discrepancy limits their practical applicability, for example, in using cameras with huge-resolution gaps to capture gigapixel video [54,56], or in reducing the budget of a camera system [47].

In contrast, our method enables the imaging system to flexibly employ various cameras with different resolution gaps. The key ingredient of our method is a cross-camera alignment module that generates multi-scale correspondence for cross-domain image alignment. Without resorting to hand-crafted design on image registration or fusion [23,11], we propose a novel neural network that leverages joint image alignment and fusion for cross-camera colorization. To improve the correspondence, we design visibility maps computation that explicitly computes and compensates warping errors. Finally, we utilize a warping regularization [42] to further improve the alignment quality.

Extensive experiments are performed to evaluate the proposed method under various settings, i.e., combinations of different resolution gaps, viewpoints, temporal steps, dynamic/static scenes. We test our method on various datasets, including the video dataset Vimeo90k [51], the light field dataset Flower [39], and the light field video dataset LFVideo [45]. Both the quantitative and qualitative experiments show the substantial improvements of our method over the state-of-the-art methods – remarkably, we achieve around 10dB gain in terms of PSNR in most of the test cases upon the best baseline.

Our main contributions are summarized as follows:
- A flexible and cost-effective imaging framework that is applicable to various color-plus-mono camera settings with multiple resolution gaps. In particular, our method can adapt to both spatial and temporal resolution gap more than $8\times$.
- A novel network design for cross-camera colorization: the cross-camera alignment generates dense correspondence for multi-scale feature alignment; the fusion module compensates alignment error via visibility map computation and performs synthesis; the warping regularization further improves the alignment.
- Extensive evaluation on a wide range of settings, *i.e.* different resolution gaps, combinations of viewpoints and temporal steps show the substantial improvements of our method, *i.e.*, around 10dB PSNR gain over the state-of-the-art ones.

## 2    Related Work

### 2.1    Automatic Image Colorization

Most traditional methods perform colorization by optimization, *e.g.*, Regression Tree Fields (RTFs) [22] and graph cuts [7]. With deep learning, some approaches [8,58] leverage large-scale color image data [63,37] to colorize grayscale images automatically. Yoo *et al.* [53] propose a colorization network augmented by external neural memory networks. Another line of works uses GANs [13,21,5,44] to colorize grayscale images by learning probability distributions.

However, automatic colorization is ill-conditioned since many potential colors can be assigned to the gray pixels of an input image. As a result, these methods tend to generate unnatural colorized images.

### 2.2    Reference-based Image Colorization

According to what is used as a reference, reference-based methods can be divided into strokes/palette-based ones and example-based ones.

*Strokes/Palette-based Image Colorization.* Strokes/palette-based image colorization methods [25,18,52,32,6] seek to reconstruct color images from the users-provided sparse stroke. However, these methods require intensive manual works. More importantly, scribbles and palettes provide insufficient color information for colorization, leading to unsatisfying results.

*Example-based Image Colorization.* To avoid the manual labor of scribbling or palette selection while facilitating controllable colorization, example-based methods are proposed to transfer the color of a reference image to the target image. The initial works [35,49] attempt to match color statistics globally. Targeting at more accurate color transfer, later works [20,40,14,9,30,3] utilize hand-crafted feature correspondences, *e.g.*, SIFT, Gabor wavelet to enforce local color consistency. Recently, He *et al.* [16] and Liao *et al.* [27] perform dense patch matching for color transfer. However, the patch-based correspondence [16,27] are inherently inefficient to compute. Furthermore, the inter-patch misalignment may hurt the image synthesis. Xiao *et al.* [50] propose a self-supervised approach that reconstructs a color image from its grayscale version and global color distribution coding. Although obtaining impressive results, the proposed approach does not utilize the spatial correspondence between two images for local colorization.

### 2.3    Flow-based or non-rigid correspondences

The cross-camera image colorization is also related to flow-based and non-rigid correspondence estimation. Specifically, back-warping reference images using estimated optical flow fields [28,19,60,61] can serve to align images from two viewpoints. In addition, HaCohen *et al.* [15] showcase the usage of estimating non-rigid dense correspondence for image enhancement. However, due to the domain and resolution gap between the two cameras, the above approaches fail in finding correct matching and result in poor performance under our cross-domain and cross-resolution setting.
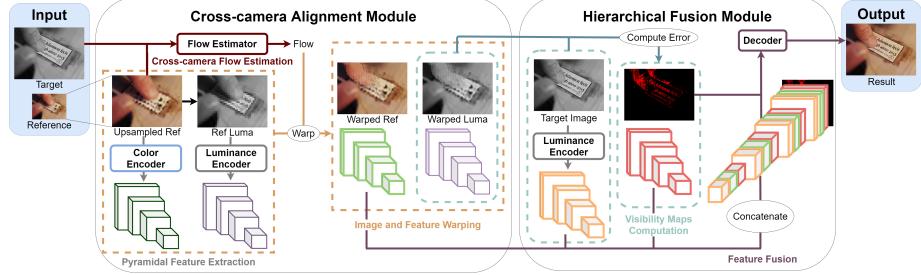
Fig. 2: Our network contains two modules: 1) the alignment module, which performs non-rigid transformation on the cross-camera image inputs and extracted features; 2) the fusion module, which performs features fusion and colorization synthesis.

## 3   Approach

Assume two images from different cameras capturing the same scene at similar viewpoints are given – an LR color image as the reference and an HR gray image as the target. We denote the target single-channel image as $I_1 \in \mathbb{R}^{sH \times sW \times 1}$ and the three-channel reference as $I_2 \in \mathbb{R}^{H \times W \times 3}$, where $s \geq 1$ is the scale factor representing the spatial-resolution gap in the dual camera, $H$ and $W$ are the horizontal and vertical spatial-resolution, respectively. Similarly, the ground-truth color image is denoted as $I_g \in \mathbb{R}^{sH \times sW \times 3}$, where $I_g$ shares the same viewpoint with the target image $I_1$. Our goal is to generate an HR color image $I_c \in \mathbb{R}^{sH \times sW \times 3}$ as similar as possible to $I_g$.

To achieve this goal, we propose an end-to-end and fully convolutional deep neural network (see illustration in Fig. 2). Our network contains a cross-camera alignment module and a hierarchical fusion module: the former (Sec. 3.1) consists of a color encoder $Net_E^c$, a luminance encoder $Net_E^l$, and a flow estimator $Net_{flow}$; the latter (Sec. 3.2) consists of a decoder $Net_D$. On top of that, our network trains jointly the warping flow estimator and colorization neural networks by combining warping loss $\mathcal{L}_w$ and colorization loss $\mathcal{L}_c$ (Sec. 3.3).

### 3.1   Cross-camera Alignment Module

This module is designed to perform the temporal and spatial alignment. As shown on the left of Fig. 2, we first extract feature maps of input images and then utilize a flow estimator to generate the cross-camera correspondence at multiple scales. After flow estimation, we use cross-camera warping to perform the non-rigid transformation. In the following, we provide details of each part.

*Pyramidal Feature Extraction.* Considering the different resolutions of the input images, we first upsample the reference $I_2$ to the same resolution as $I_1$ via bicubic upsampling, denoted by $I_2^{\uparrow}$. $I_1$ and $I_2^{\uparrow}$ are of the same resolution but belong to different image modalities, thus we design two encoders to extract their pyramidal features respectively:

$$\{F_i^l\} = Net_E^l(I_1), \quad \{F_i^c\} = Net_E^c(I_2^{\uparrow}), \quad i = 1, 2, 3, 4, \qquad (1)$$

where $F_i^l$ (resp. $F_i^c$) is the feature map of target image $I_1$ (resp. $I_2$) at scale $i$.

To measure the image alignment across different modalities, we convert the upsampled reference image $I_2^{\uparrow}$ from RGB color space to YUV [43] color space. Thus, $I_2^{\uparrow}$ can be separated into a luminance component and two chrominance components. Discarding chrominance components, we retain the luminance components denoted as $I_2^Y$, which is a grayscale counterpart of the upsampled reference image $I_2^{\uparrow}$.

Then the luminance encoder is also used to extract multi-scale feature maps of $I_2^Y$, which is utilized latter in Sec. 3.2 to calculate warping errors on feature domain:

$$\{F_i^Y\} = Net_E^l(I_2^Y), \quad i = 1, 2, 3, 4. \tag{2}$$

*Cross-camera Flow Estimation.* For image alignment, we adopt the widely used FlowNetS [12] as our flow estimator, $Net_{flow}$, to generate the dense cross-camera correspondence at multiple scales. Tailored for our setting, we change the input channel number of the first convolutional layer of FlowNetS from 6 to 4, and obtain the following flow fields:

$$\{f_i\} = Net_{flow}(I_1, I_2^{\uparrow}), \quad i = 0, 1, 2, 3, 4, \tag{3}$$

where $f_i$ is the estimated flow field at scale $i$.

*Image and Features Warping.* To perform the temporal and spatial alignment with the estimated flow fields, we utilize a warping operation similar to [62]. More specifically, our warping operation considers the cross-camera flow field $f$:

$$\widetilde{I} = \mathcal{W}(I, f), \tag{4}$$

where $\mathcal{W}(I, f)$ denotes the result of warping the input $I$ using the flow field $f$.

After flow estimation, we perform the warping operation on the reference image features $F_i^c$ and corresponding luminance features $F_i^Y$. Using the multi-scale flow $f_i$ in Eq. 3, we generate the temporally and spatially aligned features $\{\widetilde{F}_i^c\}$ and $\{\widetilde{F}_i^Y\}$:

$$\widetilde{F}_i^c = \mathcal{W}(F_i^c, f_i), \quad \widetilde{F}_i^Y = \mathcal{W}(F_i^Y, f_i), \quad i = 1, 2, 3, 4. \tag{5}$$

To measure image alignment latter in Sec. 3.2, we also perform the warping operation on the image domain. According to Eq. 4, we have:

$$\widetilde{I}_2^Y = \mathcal{W}(I_2^Y, f_0), \tag{6}$$

where $f_0$ is the estimated flow field at scale 0 in Eq. 3, $\widetilde{I}_2^Y$ is the warping result utilized latter in Sec. 3.2 to calculate warping errors on the image domain.

### 3.2   Hierarchical Fusion Module

*Visibility Maps Computation.* On the image domain, the warping error indicates the different light intensity between the target and the reference, *i.e.*, optical visibility. While

on the feature domain, warping error represents the different activation value of the feature maps, *i.e.*, feature recognition. We combine the errors from both perspectives and define the multi-scale visibility maps $\{V_i\}$ as warping errors on both domains:

$$V_0 = \widetilde{I}_2^Y - I_1, \quad V_i = \widetilde{F}_i^Y - F_i^l, \quad i = 1, 2, 3, 4, \tag{7}$$

where $V_i$ is the warping error at scale $i$.

To give an intuition of the visibility maps, we visualize it on the image domain in Fig. 3. There are lost details caused by the low-quality reference (*e.g.*, blurry words shown on the top of Fig. 3), motion blur caused by fast-moving objects (*e.g.*, waving hand shown on the bottom of Fig. 3), local occlusion caused by motion or parallax (*e.g.*, the garment occluded by the hand shown on the bottom of Fig. 3). The pixel-wise positive and negative values of visibility maps represent the invisible regions in the reference and the target image, respectively.



Fig. 3: Visualizing the visibility map on image domain. Red: invisible regions in the reference image; green: invisible regions in the target image. Zoom in to see details.

*Feature Fusion.* In the end, we design a U-Net [36] like decoder to fuse feature maps and visibility maps, and synthesize the colorization result. As shown on the right of Fig. 2, the target features $\{F_i^l\}$, warped reference features $\{\widetilde{F}_i^c\}$ and visibility maps $\{V_i\}$ are concatenated as the input of fusion decoder. Finally we obtain the result $I_c$ by:

$$I_c = Net_D(\{F_i^l\}, \{\widetilde{F}_i^c\}, \{V_j\}), \quad i = 1, 2, 3, 4, j = 0, 1, 2, 3, 4. \tag{8}$$

### 3.3   Loss Function

We use two loss functions: warping loss and colorization loss. The former encourages the flow estimator to generate precise cross-camera correspondence for image alignment. The latter is responsible for the final synthesized image.

*Warping Loss.* Since the ground-truth flow is unavailable, it is difficult to train the flow estimator in an unsupervised fashion. To solve this problem, we adopt the warping loss from [42]. Specifically, since the input images capture the same scene, it is reasonable to require the warped-upsampled reference image owning a intensity distribution similar to the ground truth $I_g$ as much as possible. Thus our warping loss is defined as:

$$\mathcal{L}_w = \frac{1}{2N} \sum_{i=1}^{N} \sum_j ||\widetilde{I}_2^{\uparrow(i)} - I_g^{(i)}||_2^2, \quad \widetilde{I}_2^{\uparrow} = \mathcal{W}(I_2^{\uparrow}, f_0), \tag{9}$$

where $N$ is the sample number, $i$ iterates over training samples, $f_0$ is the estimated flow field at scale $0$ in Eq. 3.

*Colorization Loss.* Given the network prediction $\widetilde{I}_c$ and the ground truth $I_c$, the colorization loss is defined as:

$$\mathcal{L}_c = \frac{1}{N} \sum_{i=1}^{N} \sum_{j} \rho(I_c^{(i)} - I_g^{(i)}), \tag{10}$$

where $\rho(x) = \sqrt{x^2 + 0.001^2}$ is the Charbonnier penalty function [2], $\widetilde{I}_c$ is obtained from Eq. 8, $N$ is the sample number, $i$ iterates over training samples.



Fig. 4: Colorization comparisons on Vimeo90K dataset under cross-scale $8\times$ settings.

Fig. 5: Comparison under cross-scale $8\times$ settings on Flower (the 1st, 2nd, 3rd rows) and LFVideo (the last row) dataset, respectively.
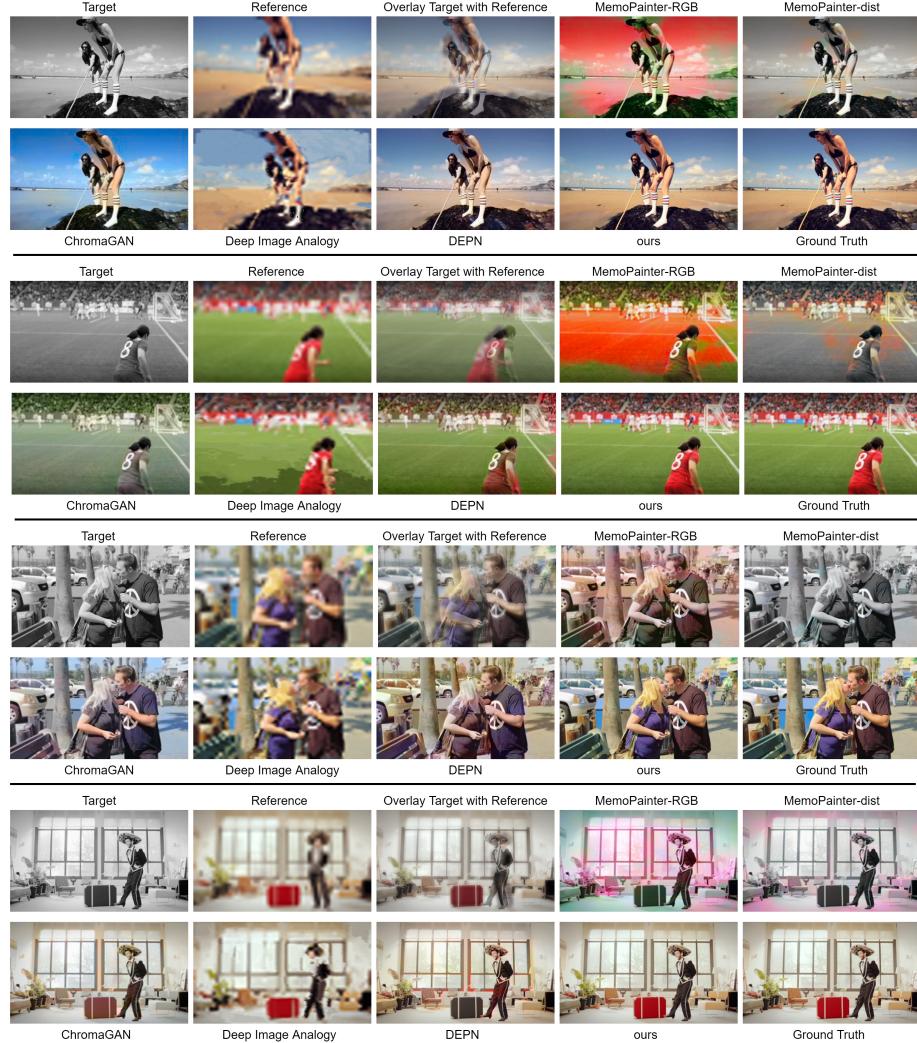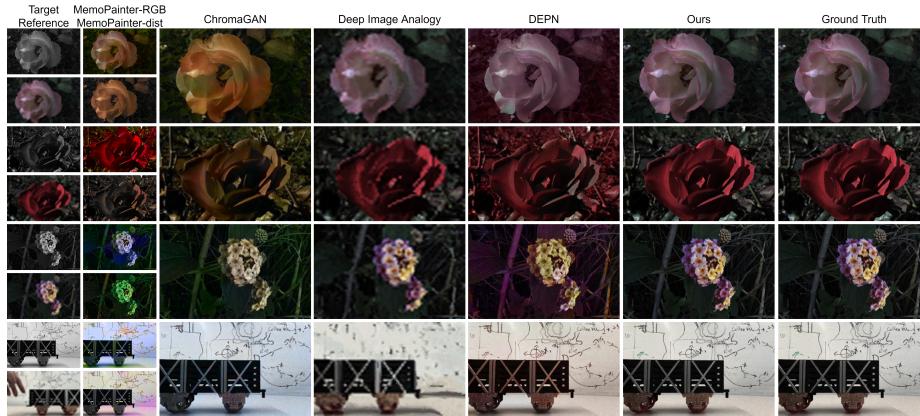
## 4    Experiments

### 4.1    Dataset

*Video Dataset.* The Vimeo90K [51] dataset consists of videos that cover a large variety of scenes and actions. Following [51], we selected $66,178$ video sequences, each of which contains 7 frames with resolution of $448 \times 256$. To construct the training and testing datasets, we randomly divide it into $60,000$ sequences for training and $6,178$ sequences for testing. For training, we downsample the first frame at each video sequence as a reference image and randomly select the $t$th$(1 < t < 8)$ frame converted to the grayscale image as a target. For testing, reference images are sampled from the first video frame, while the target images are from the second and the last frame. See the supplementary material for training details.

*Light-field Dataset.* The Flower dataset [39] contains flowers and plants light-field images with the $376 \times 541$ spatial resolution, and $14 \times 14$ angular samples. Following [39], we extract the central $8 \times 8$ grid of angular sample, and randomly selected 343 samples for evaluations. The images at viewpoints $(1, 1)$ and $(7, 7)$ are converted to grayscale as target, and images at viewpoint $(0, 0)$ are downsampled as reference.

*Light-field Video Dataset.* The LFVideo dataset [45] contains real-scene light-field videos with the spatial resolution as $376 \times 541$, while the angular samples are $8 \times 8$. For evaluations, we randomly selected 270 video frames. The images of the $t$th$(t = 2, 9)$ frame at viewpoints $(i, i), i = 1, 7$ are converted to monochrome as target, and images of the first frame at viewpoint $(0, 0)$ are downsampled as reference.

### 4.2    Comparison to State-of-the-Art Methods

Our method is compared against the state-of-the-art example-based colorization methods, namely Deep Image Analogy [27] and DEPN [50], and the recent automatic colorization approaches, including MemoPainter [53] and ChromaGAN [44]. Following [53], we evaluated MemoPainter using RGB information and color distribution for color features, respectively. Both the quantitative evaluation in Table 1 and the qualitative results in Figs. 4 and 5 suggest that our performance is far beyond that of the baselines.

Table 1: Quantitative evaluations of the state-of-the-art automatic and example-based algorithms on different datasets, in terms of NRMSE/PSNR/SSIM/LPIPS for different scale factors, frame gaps and parallax settings respectively.

| Dataset | Target Position | | Scale | Methods | Quantitative evaluations | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Frame | View | | | NRMSE | PSNR | SSIM | LPIPS | Runtime |
| | 2 | - | N/A | MemoPainter-RGB [53] | 0.3463 | 19.9920 | 0.7989 | 0.4206 | 0.5124 |
| | 2 | - | N/A | MemoPainter-dist [53] | 0.2352 | 22.5768 | 0.8751 | 0.3015 | 0.5620 |
| | 2 | - | N/A | ChromaGAN [44] | 0.2097 | 23.5276 | 0.8697 | 0.2917 | 2.6915 |
| | 2 | - | 4× | Deep Image Analogy [27] | 0.1544 | 25.6775 | 0.7741 | 0.3184 | 219.5779 |
| | 2 | - | 4× | DEPN [50] | 0.1313 | 27.5916 | 0.9313 | 0.1585 | 0.3425 |
| | 2 | - | 4× | Ours | **0.0227** | **43.2263** | **0.9884** | **0.0157** | **0.0838** |
| | 2 | - | 8× | DEPN [50] | 0.1316 | 27.5724 | 0.9310 | 0.1588 | 0.3525 |
| | 2 | - | 8× | Deep Image Analogy [27] | 0.1951 | 23.5943 | 0.6954 | 0.4142 | 203.4619 |
| Vimeo | 2 | - | 8× | Ours | **0.0275** | **41.6039** | **0.9845** | **0.0241** | **0.0847** |
| | 7 | - | N/A | MemoPainter-RGB [53] | 0.3150 | 20.0014 | 0.7995 | 0.4201 | 0.5301 |
| | 7 | - | N/A | MemoPainter-dist [53] | 0.2326 | 22.4321 | 0.8738 | 0.3028 | 0.5483 |
| | 7 | - | N/A | ChromaGAN [44] | 0.2064 | 23.5428 | 0.8696 | 0.2909 | 2.7384 |
| | 7 | - | 4× | Deep Image Analogy [27] | 0.2063 | 23.0956 | 0.7227 | 0.3529 | 214.3759 |
| | 7 | - | 4× | DEPN [50] | 0.1306 | 27.5174 | 0.9313 | 0.1593 | 0.3548 |
| | 7 | - | 4× | Ours | **0.0380** | **39.5223** | **0.9823** | **0.0321** | **0.0846** |
| | 7 | - | 8× | Deep Image Analogy [27] | 0.2356 | 21.8671 | 0.66771 | 0.4260 | 220.8749 |
| | 7 | - | 8× | DEPN [50] | 0.1309 | 27.4989 | 0.9311 | 0.1596 | 0.3478 |
| | 7 | - | 8× | Ours | **0.0408** | **38.6736** | **0.9796** | **0.0382** | **0.0891** |
| | - | (1,1) | N/A | MemoPainter-RGB [53] | 0.4172 | 22.0304 | 0.7508 | 0.3389 | 0.5623 |
| | - | (1,1) | N/A | MemoPainter-dist [53] | 0.3237 | 24.2114 | 0.8822 | 0.2668 | 0.6184 |
| | - | (1,1) | N/A | ChromaGAN [44] | 0.3046 | 24.4863 | 0.8407 | 0.3081 | 1.1562 |
| | - | (1,1) | 4× | Deep Image Analogy [27] | 0.1874 | 28.6252 | 0.8065 | 0.2581 | 411.6000 |
| | - | (1,1) | 4× | DEPN [50] | 0.1797 | 29.0692 | 0.9286 | 0.1711 | 0.4899 |
| | - | (1,1) | 4× | Ours | **0.0205** | **45.9354** | **0.9938** | **0.0041** | **0.1173** |
| | - | (1,1) | 8× | Deep Image Analogy [27] | 0.2620 | 25.6807 | 0.6976 | 0.4029 | 404.4650 |
| | - | (1,1) | 8× | DEPN [50] | 0.1797 | 29.0730 | 0.9285 | 0.1710 | 0.4729 |
| Flower | - | (1,1) | 8× | Ours | **0.0255** | **43.8623** | **0.9923** | **0.0077** | **0.1171** |
| | - | (7,7) | N/A | MemoPainter-RGB [53] | 0.4216 | 21.9233 | 0.7525 | 0.3424 | 0.5912 |
| | - | (7,7) | N/A | MemoPainter-dist [53] | 0.3158 | 24.4801 | 0.8858 | 0.2617 | 0.6034 |
| | - | (7,7) | N/A | ChromaGAN [44] | 0.3065 | 24.4616 | 0.8404 | 0.3081 | 1.5492 |
| | - | (7,7) | 4× | Deep Image Analogy [27] | 0.2532 | 25.9906 | 0.7100 | 0.2954 | 409.2839 |
| | - | (7,7) | 4× | DEPN [50] | 0.1775 | 29.2156 | 0.9300 | 0.1663 | 0.4725 |
| | - | (7,7) | 4× | Ours | **0.0237** | **44.8138** | **0.9931** | **0.0060** | **0.1170** |
| | - | (7,7) | 8× | Deep Image Analogy [27] | 0.2903 | 24.7787 | 0.6703 | 0.4006 | 405.3829 |
| | - | (7,7) | 8× | DEPN [50] | 0.1774 | 29.2186 | 0.9299 | 0.1663 | 0.4793 |
| | - | (7,7) | 8× | Ours | **0.0276** | **43.2567** | **0.9920** | **0.0092** | **0.1175** |
| | 2 | (1,1) | N/A | MemoPainter-RGB [53] | 0.4190 | 21.4680 | 0.7227 | 0.3470 | 0.6633 |
| | 2 | (1,1) | N/A | MemoPainter-dist [53] | 0.3317 | 23.7711 | 0.8565 | 0.2181 | 0.6927 |
| | 2 | (1,1) | N/A | ChromaGAN [44] | 0.1982 | 27.6557 | 0.8592 | 0.2120 | 0.5773 |
| | 2 | (1,1) | 4× | Deep Image Analogy [27] | 0.3304 | 23.2553 | 0.6374 | 0.4040 | 755.4143 |
| | 2 | (1,1) | 4× | DEPN [50] | 0.1129 | 32.7250 | 0.9668 | 0.0933 | 0.4824 |
| | 2 | (1,1) | 4× | Ours | **0.0459** | **41.0111** | **0.9857** | **0.0288** | **0.1173** |
| | 2 | (1,1) | 8× | Deep Image Analogy [27] | 0.3695 | 22.3074 | 0.5970 | 0.4944 | 748.9276 |
| | 2 | (1,1) | 8× | DEPN [50] | 0.1129 | 32.7220 | 0.9669 | 0.0933 | 0.4793 |
| | 2 | (1,1) | 8× | Ours | **0.0492** | **40.3612** | **0.9849** | **0.0312** | **0.1167** |
| | 2 | (7,7) | N/A | MemoPainter-RGB [53] | 0.4261 | 21.3552 | 0.7125 | 0.3527 | 0.6325 |
| | 2 | (7,7) | N/A | MemoPainter-dist [53] | 0.3058 | 24.4653 | 0.8662 | 0.2037 | 0.6429 |
| | 2 | (7,7) | N/A | ChromaGAN [44] | 0.1959 | 27.8429 | 0.8609 | 0.2097 | 0.5826 |
| | 2 | (7,7) | 4× | Deep Image Analogy [27] | 0.3394 | 22.9054 | 0.6315 | 0.4041 | 751.4360 |
| | 2 | (7,7) | 4× | DEPN [50] | 0.1104 | 32.9745 | 0.9679 | 0.0880 | 0.4529 |
| | 2 | (7,7) | 4× | Ours | **0.0448** | **41.0335** | **0.9859** | **0.0274** | **0.1185** |
| | 2 | (7,7) | 8× | Deep Image Analogy [27] | 0.3718 | 22.0710 | 0.5977 | 0.4920 | 753.5917 |
| | 2 | (7,7) | 8× | DEPN [50] | 0.11070 | 32.9560 | 0.9680 | 0.0883 | 0.4826 |
| LFVideo | 2 | (7,7) | 8× | Ours | **0.0480** | **40.4257** | **0.9853** | **0.0301** | **0.1179** |
| | 9 | (1,1) | N/A | MemoPainter-RGB [53] | 0.4501 | 21.1598 | 0.6991 | 0.3608 | 0.6530 |
| | 9 | (1,1) | N/A | MemoPainter-dist [53] | 0.3290 | 23.7217 | 0.8582 | 0.2204 | 0.6498 |
| | 9 | (1,1) | N/A | ChromaGAN [44] | 0.2117 | 27.5192 | 0.8512 | 0.2244 | 0.5728 |
| | 9 | (1,1) | 4× | Deep Image Analogy [27] | 0.3270 | 23.5479 | 0.6404 | 0.4102 | 749.7418 |
| | 9 | (1,1) | 4× | DEPN [50] | 0.1165 | 32.5500 | 0.9670 | 0.0917 | 0.4692 |
| | 9 | (1,1) | 4× | Ours | **0.0508** | **40.7632** | **0.9853** | **0.0323** | **0.1169** |
| | 9 | (1,1) | 8× | Deep Image Analogy [27] | 0.3709 | 22.4716 | 0.5979 | 0.4954 | 747.7591 |
| | 9 | (1,1) | 8× | DEPN [50] | 0.1166 | 32.5409 | 0.9671 | 0.0919 | 0.4792 |
| | 9 | (1,1) | 8× | Ours | **0.0538** | **40.1106** | **0.9845** | **0.0349** | **0.1167** |
| | 9 | (7,7) | N/A | MemoPainter-RGB [53] | 0.4104 | 21.9786 | 0.7211 | 0.3464 | 0.6009 |
| | 9 | (7,7) | N/A | MemoPainter-dist [53] | 0.3271 | 23.9009 | 0.8628 | 0.2170 | 0.6947 |
| | 9 | (7,7) | N/A | ChromaGAN [44] | 0.2112 | 27.5598 | 0.8542 | 0.2248 | 0.5927 |
| | 9 | (7,7) | 4× | Deep Image Analogy [27] | 0.3221 | 23.7166 | 0.6447 | 0.4058 | 747.9275 |
| | 9 | (7,7) | 4× | DEPN [50] | 0.1130 | 32.7364 | 0.9676 | 0.0882 | 0.4865 |
| | 9 | (7,7) | 4× | Ours | **0.0492** | **40.8788** | **0.9858** | **0.0304** | **0.1420** |
| | 9 | (7,7) | 8× | Deep Image Analogy [27] | 0.3564 | 22.8959 | 0.6035 | 0.4921 | 750.8465 |
| | 9 | (7,7) | 8× | DEPN [50] | 0.1132 | 32.7194 | 0.9678 | 0.0881 | 0.4539 |
| | 9 | (7,7) | 8× | Ours | **0.0518** | **40.3270** | **0.9851** | **0.0328** | **0.1453** |

Our quantitative evaluation involves four image quality metrics: NRMSE, PSNR, SSIM [48], and LPIPS [57]. Table 1 shows quantitative comparisons, and it is evident that our method outperforms the baselines by a large margin in most of the settings, including combinations of different scale factors, frame gaps, and parallax. In general, our method achieves approximately 10dB gain of PSNR upon the baselines.

We provide qualitative comparisons in Figs. 4 and 5 respectively on the Vimeo90k, Flower, LFVideo datasets under the challenging scale 8×, largest parallax and largest frame gap setting. Firstly, benefiting from the reference image, the example-based approaches show better results than the automatic colorization approaches. Moreover, among the example-based approaches, our method generates the most coherent images with less color bleeding effects, showing the clear advantage of our algorithm.

### 4.3   Ablation Study

This section investigates the role of proposed visibility maps and warping loss by two variants of our pipeline, which respectively turn off one of the components while remaining others. In particular, we run our test on the Vimeo90k dataset, with resolution gaps being 4× and 8×. According to Fig. 6, disabling either component leads to a performance drop, suggesting the necessities of both.



*Visibility Maps*  When visibility maps are disabled, slight degradation in PSNR is observed in Fig. 6. Also, as shown in Fig. 3, artifacts are more evident without visibility maps, suggesting that visibility maps can reduce the wrong colorization due to inconsistent visibility between input images.

Fig. 6: The PSNR performance of the different variants with some of the components disabled during training on Vimeo90K dataset.

*Warping Loss*  From Fig. 6, the performance of our network drops if warping loss is remove. We speculate that it is because the colorization loss is for image synthesis and does not explicitly define terms for flow estimation. In contrast, the warping loss is necessary to regularize the flow estimator training and enable better convergence.

## 5   Conclusion

We propose a novel convolutional neural network to facilitate the trade-off between spatial-temporal resolution and color depth in the imaging system. Our method fuses data to obtain an HR color image given cross-domain and cross-scale input pairs captured by the color-plus-mono dual camera. In contrast to previous works, our method can adapt to color and monochrome cameras with the various spatial-temporal resolution; thus it enables more flexible and generalized imaging systems. The key ingredient of our method is a cross-camera alignment module that generates multi-scale correspondence for cross-domain image alignment. Experiments on several datasets demonstrate the superior performance of our method (around 10dB in PSNR) compared to state-of-the-art methods.
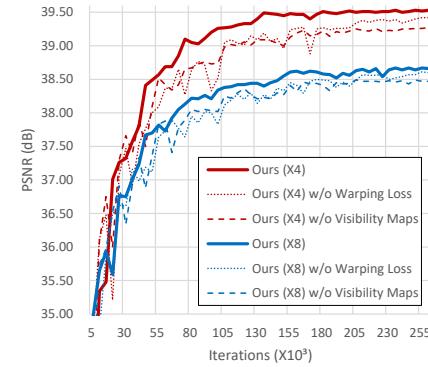
# References

1. Brady, D.J., et al.: Multiscale gigapixel photography. Nature **486**(7403), 386–389 (2012)
2. Bruhn, A., Weickert, J., Schnörr, C.: Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. International journal of computer vision **61**(3), 211–231 (2005)
3. Bugeau, A., et al.: Variational exemplar-based image colorization. IEEE TIP (2013)
4. Cao, X., Tong, X., Dai, Q., Lin, S.: High resolution multispectral video capture with a hybrid camera system. In: CVPR 2011. pp. 297–304. IEEE (2011)
5. Cao, Y., Zhou, Z., Zhang, W., Yu, Y.: Unsupervised diverse colorization via generative adversarial networks. In: ECML PKDD. pp. 151–166. Springer (2017)
6. Chang, H., et al.: Palette-based photo recoloring. ACM Trans. Graph. (2015)
7. Charpiat, G., et al.: Automatic image colorization via multimodal predictions. In: ECCV. Springer (2008)
8. Cheng, Z., Yang, Q., Sheng, B.: Deep colorization. In: ICCV. pp. 415–423 (2015)
9. Chia, A.Y.S., et al.: Semantic colorization with internet images. ACM TOG (2011)
10. Cossairt, O.S., et al.: Gigapixel computational imaging. In: ICCP. IEEE (2011)
11. Dong, X., Li, W.: Shoot high-quality color images using dual-lens system with monochrome and color cameras. Neurocomputing **352**, 22–32 (2019)
12. Dosovitskiy, A., et al.: Flownet: Learning optical flow with conv networks. In: ICCV (2015)
13. Goodfellow, I., et al.: Generative adversarial nets. In: NeuIPS. pp. 2672–2680 (2014)
14. Gupta, R.K., et al.: Image colorization using similar images. In: ACM MM (2012)
15. HaCohen, Y., Shechtman, E., Goldman, D.B., Lischinski, D.: Non-rigid dense correspondence with applications for image enhancement. ACM TOG **30**(4), 1–10 (2011)
16. He, M., Liao, J., Yuan, L., Sander, P.V.: Neural color transfer between images. arXiv (2017)
17. He, M., et al.: Deep exemplar-based colorization. ACM TOG (2018)
18. Huang, Y.C., et al.: An adaptive edge detection based colorization algorithm and its applications. In: ACM MM (2005)
19. Ilg, E., et al.: Flownet 2.0: Evolution of optical flow estimation with deep networks. In: CVPR. pp. 2462–2470 (2017)
20. Ironi, R., et al.: Colorization by example. In: Rendering Techniques. Citeseer (2005)
21. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: CVPR. pp. 1125–1134 (2017)
22. Jancsary, J., Nowozin, S., Sharp, T., Rother, C.: Regression tree fields—an efficient, non-parametric approach to image labeling problems. In: CVPR. pp. 2376–2383. IEEE (2012)
23. Jeon, H.G., et al.: Stereo matching with color and monochrome cameras in low-light conditions. In: CVPR (2016)
24. Jin, D., et al.: All-in-depth via cross-baseline light field camera. In: ACM MM (2020)
25. Levin, A., et al.: Colorization using optimization. In: ACM SIGGRAPH 2004 Papers (2004)
26. Li, G., et al.: Zoom in to the details of human-centric videos. In: ICIP. IEEE (2020)
27. Liao, J., et al.: Visual attribute transfer through deep image analogy. arXiv (2017)
28. Liu, C., et al.: Sift flow: Dense correspondence across scenes and its applications. IEEE TPAMI (2010)
29. Liu, C., Shan, J., Liu, G.: High resolution array camera (Apr 19 2016), uS Patent 9,319,585
30. Liu, X., et al.: Intrinsic colorization. In: ACM SIGGRAPH Asia 2008 papers, pp. 1–9 (2008)
31. Lohmann, et al.: Space–bandwidth product of optical signals and systems. JOSA A (1996)
32. Luan, Q., Wen, F., Cohen-Or, D., Liang, L., Xu, Y.Q., Shum, H.Y.: Natural image colorization. In: Eurographics conference on Rendering Techniques. pp. 309–320 (2007)
33. Ma, C., Cao, X., Tong, X., Dai, Q., Lin, S.: Acquisition of high spatial and spectral resolution video with a hybrid camera system. IJCV **110**(2), 141–155 (2014)
34. Mantzel, W., et al.: Shift-and-match fusion of color and mono images (2017), uS Patent

35. Reinhard, E., et al.: Color transfer between images. IEEE Computer graphics and applications **21**(5), 34–41 (2001)
36. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241. Springer (2015)
37. Russakovsky, O., et al.: Imagenet large scale visual recognition challenge. IJCV (2015)
38. Sharif, S., Jung, Y.J.: Deep color reconstruction for a sparse color sensor. Optics express **27**(17), 23661–23681 (2019)
39. Srinivasan, P.P., et al.: Learning to synthesize a 4d rgbd light field from a single image. In: ICCV. pp. 2243–2251 (2017)
40. Tai, Y.W., Jia, J., Tang, C.K.: Local color transfer via probabilistic segmentation by expectation-maximization. In: CVPR. vol. 1, pp. 747–754. IEEE (2005)
41. Tai, Y.W., et al.: Image/video deblurring using a hybrid camera. In: CVPR. IEEE (2008)
42. Tan, Y., Zheng, H., Zhu, Y., Yuan, X., Lin, X., Brady, D., Fang, L.: Crossnet++: Cross-scale large-parallax warping for reference-based super-resolution. IEEE TPAMI (2020)
43. Union, I.T.: Encoding parameters of digital television for studios. CCIR Recommend. (1992)
44. Vitoria, P., Raad, L., Ballester, C.: Chromagan: An adversarial approach for picture colorization. arXiv preprint arXiv:1907.09837 (2019)
45. Wang, T.C., et al.: Light field video capture using a learning-based hybrid imaging system. ACM TOG (2017)
46. Wang, X., et al.: Panda: A gigapixel-level human-centric video dataset. In: CVPR (2020)
47. Wang, Y., Liu, Y., Heidrich, W., Dai, Q.: The light field attachment: Turning a dslr into a light field camera using a low budget camera ring. IEEE TVCG **23**(10), 2357–2364 (2016)
48. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE TIP **13**(4), 600–612 (2004)
49. Welsh, T., Ashikhmin, M., Mueller, K.: Transferring color to greyscale images. In: annual conference on Computer graphics and interactive techniques. pp. 277–280 (2002)
50. Xiao, C., Han, C., Zhang, Z., others, G., He, S.: Example-based colourization via dense encoding pyramids. In: Computer Graphics Forum. Wiley Online Library (2020)
51. Xue, T., et al.: Video enhancement with task-oriented flow. IJCV (2019)
52. Yatziv, L., Sapiro, G.: Fast image and video colorization using chrominance blending. IEEE transactions on image processing **15**(5), 1120–1129 (2006)
53. Yoo, S., Bahng, H., Chung, S., Lee, J., Chang, J., Choo, J.: Coloring with limited data: Few-shot colorization via memory augmented networks. In: CVPR. pp. 11283–11292 (2019)
54. Yuan, X., Fang, L., Dai, Q., Brady, D.J., Liu, Y.: Multiscale gigapixel video: A cross resolution image matching and warping approach. In: ICCP. pp. 1–9. IEEE (2017)
55. Yuan, X., et al.: A modular hierarchical array camera. Light: Science & Applications (2021)
56. Zhang, J., Zhu, T., Zhang, A., et al.: Multiscale-vr: Multiscale gigapixel 3d panoramic videography for virtual reality. In: ICCP. pp. 1–12. IEEE (2020)
57. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: CVPR. pp. 586–595 (2018)
58. Zhang, R., et al.: Colorful image colorization. In: ECCV. pp. 649–666. Springer (2016)
59. Zhao, Y., Li, G., Wang, Z., Lam, E.Y.: Cross-camera human motion transfer by time series analysis. arXiv preprint arXiv:2109.14174 (2021)
60. Zhao, Y., et al.: Efenet: Reference-based video super-resolution with enhanced flow estimation. In: CICAI. pp. 371–383. Springer (2021)
61. Zhao, Y., et al.: Manet: Improving video denoising with a multi-alignment network. arXiv preprint arXiv:2202.09704 (2022)
62. Zheng, H., Ji, M., Wang, H., Liu, Y., Fang, L.: Crossnet: An end-to-end reference-based super resolution network using cross-scale warping. In: ECCV. pp. 88–104 (2018)
63. Zhou, B., et al.: Learning deep features for scene recognition using places database. In: NeuIPS. pp. 487–495 (2014)