# Reconstructing Continuous Light Field from Single Coded Image

**YUYA ISHIKAWA[1], KEITA TAKAHASHI[1](Member, IEEE), CHIHIRO TSUTAKE[1](Member, IEEE), AND TOSHIAKI FUJII[1](Member, IEEE)**

[1]Graduate School of Engineering, Nagoya University, Nagoya 464-8603, Japan

Corresponding author: Yuya Ishikawa (e-mail: ishikawa.yuya@fujii.nuee.nagoya-u.ac.jp).

**ABSTRACT** We propose a method for reconstructing a continuous light field of a target scene from a single observed image. Our method takes the best of two worlds: joint aperture-exposure coding for compressive light-field acquisition, and a neural radiance field (NeRF) for view synthesis. Joint aperture-exposure coding implemented in a camera enables effective embedding of 3-D scene information into an observed image, but in previous works, it was used only for reconstructing discretized light-field views. NeRF-based neural rendering enables high quality view synthesis of a 3-D scene from continuous viewpoints, but when only a single image is given as the input, it struggles to achieve satisfactory quality. Our method integrates these two techniques into an efficient and end-to-end trainable pipeline. Trained on a wide variety of scenes, our method can reconstruct continuous light fields accurately and efficiently without any test time optimization. To our knowledge, this is the first work to bridge two worlds: camera design for efficiently acquiring 3-D information and neural rendering.

**INDEX TERMS** Light field, Compressed sensing, Neural representation

## I. INTRODUCTION

Light field is usually represented as a set of dense multi-view images. Thanks to the abundant information contained, light fields have been used for various applications including depth estimation [1], [2], object/material recognition [3], [4], view synthesis [5]–[7], and 3-D display [8]–[10].

Acquisition of a light field is a long-standing issue [11]–[14]. Since light-field views are highly redundant with each other, view by view sampling seems to be a waste of resources. To achieve more efficient acquisition, camera-side coding schemes have been developed [15]–[21]. Combined with learning-based reconstruction algorithms, the latest methods with joint aperture-exposure coding [22]–[25] drastically reduce the number of images required for high-quality reconstruction; even a single coded image alone is sufficient to reconstruct a light field (with, e.g., $5 \times 5$ views) with convincing quality. However, what is obtained from these methods is only a set of discretized views. Since there is no physical "viewpoint grid" in the target scene, more desirable is a continuous representation of a target 3-D scene that can be observed from continuous viewpoints.

Our goal is to develop a method for reconstructing a continuous light field for a target scene from a single observed image, as depicted in Fig. 1. To this end, we take the best of two worlds: joint aperture-exposure coding [22]–[25] for
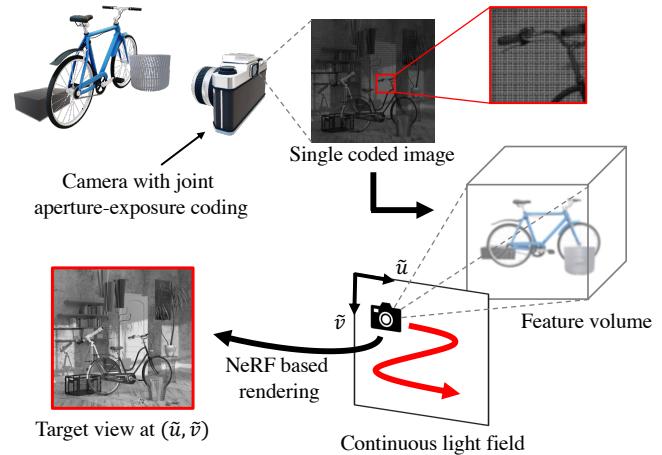


FIGURE 1: Our goal is to reconstruct continuous light field from single coded image alone. To this end, we integrate joint aperture-exposure coding and NeRF-based rendering.

compressive light-field acquisition, and a neural radiance field (NeRF) [26]–[29] for view synthesis. Joint aperture-exposure coding enables effective embedding of 3-D scene

information into an observed image. A NeRF-based representation enables high quality rendering of a target scene from continuous viewpoints. Our method integrates these two techniques into an efficient and end-to-end trainable pipeline. Trained on a wide variety of scenes, our method can reconstruct continuous light fields accurately and efficiently without any test time optimization.

To the best of our knowledge, there is no previous work that can directly obtain a continuous light field from a single coded image in the context of compressive light-field acquisition. Moreover, we are the first to use not a normal (uncoded) but a coded image as the input for NeRF-based neural rendering. We show that the integration of camera-side coding and NeRF-based representation enables accurate and efficient rendering of a 3-D scene only from a single observed image. We believe our contribution will open up a new field embracing camera design and neural rendering.

The remainder of this paper is organized as follows. Section II briefly gives some background on compressive light-field acquisition and view synthesis to clarify the position of our method. Section III describes our method including the camera-side coding scheme and reconstruction of a continuous light field. Section IV presents several experimental validations including comparisons with the state-of-the-art methods. Section V concludes the paper.

## II. BACKGROUND
### A. COMPRESSIVE LIGHT FIELD ACQUISITION
Traditionally, a light field was captured using an array of cameras [11], [13], [14] and a lenslet-based camera [12], [30]. More recently, coded aperture cameras [15]–[21] have been developed to increase the efficiency of light field acquisition. For example, two to four images, taken with different aperture-coding patterns, are sufficient to computationally reconstruct a light field with $5 \times 5$ or $8 \times 8$ views. However, the reconstruction quality obtained with a single coded image is still unsatisfactory. Joint aperture-exposure coding [22]–[25] enables more flexible coding patterns within a single exposure time. It was demonstrated that with this advanced coding scheme, only a single image is sufficient to achieve high-quality reconstruction.

A fundamental limitation of these previous works is the discontinuity of the viewpoints; what is obtained from these methods is a set of images at discretized viewpoints. To break this limitation, our method is designed to reconstruct a continuous light field. More specifically, using a single image captured with joint aperture-exposure coding as the input, our method derives a neural representation, from which arbitrary light rays at continuous coordinates can be rendered through volume rendering.

### B. VIEW SYNTHESIS
Given a set of posed images of a target scene, view synthesis aims to generate arbitrary views at continuous viewpoints. For this purpose, the traditional methods used 3-D meshes, voxels, point clouds, and depth maps as the 3-D represen-

tation [31]–[35]. Recent progress of neural networks has brought more implicit 3-D representations, such as multi-plane images [36] and neural radiance fields (NeRF) [26]. A NeRF is a coordinate-based representation of a target scene, from which high-quality images can be synthesized at continuous viewpoints through volume rendering. The original NeRF is scene-specific; the network model is optimized for each target scene, and it requires many images for training. The follow-up works of the original NeRF aimed at faster rendering [37]–[39], fewer input images [40]–[44], and generalization to unseen scenes [40]–[43], [45].

Single-view view synthesis [43], [46]–[52] refers to the extreme case of view synthesis, where only a single view is given as the input. Due to the ill-posedness of the problem, even the latest state-of-the-art methods struggle to achieve high quality rendering.

Constructed on the framework of NeRF [26], our method can synthesize arbitrary views at continuous viewpoints. As a notable difference from the previous works, our method uses only a single image captured with joint aperture-exposure coding applied on the camera. An image coded in this manner can contain richer 3-D information of a target scene than a normal (uncoded) image. Moreover, our method is designed to be light-weight and generalized over a wide variety of scenes.

## III. PROPOSED METHOD
### A. PROBLEM FORMULATION AND OVERVIEW
A light field is represented as a set of multi-view images taken from dense 2-D grid-points. It is written as $L(u, v, x, y)$, where $(u, v)$ and $(x, y)$ denote the viewpoint index and pixel positions, respectively. We also introduce a continuous light field, denoted as $L(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})$, where the coordinate $(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})$ can take arbitrary continuous values. In particular, we focus on the continuity of the viewpoints $(\tilde{u}, \tilde{v})$ because this is the main focus of view synthesis. We assume that the light field is grayscale, considering the availability of imaging hardware as mentioned in III-B.

Our goal is to reconstruct a continuous light field of a target scene from a single image alone. To fully obtain the 3-D information of the target scene, we use an image captured with joint aperture-exposure coding [22]–[25] as the input. From the coded image, we extract a feature volume that spans the target 3-D volume using a deep convolutional neural network (CNN) called FeatNet. We also construct a multi-layer perceptron (MLP) called NeRFNet to represent the radiance field of the same 3-D volume, and it is conditioned on the features extracted from the coded image. To render a view from a desired viewpoint, we make queries towards NeRFNet for the luminance and density values along the light rays that constitute the target view. Since the queries can be made at arbitrary continuous coordinates, we can reconstruct a continuous light field. Moreover, FeatNet and NeRFNet are jointly trained and can be generalized over a wide variety of scenes. Our method can reconstruct new scenes (which are unseen during training time) accurately and efficiently
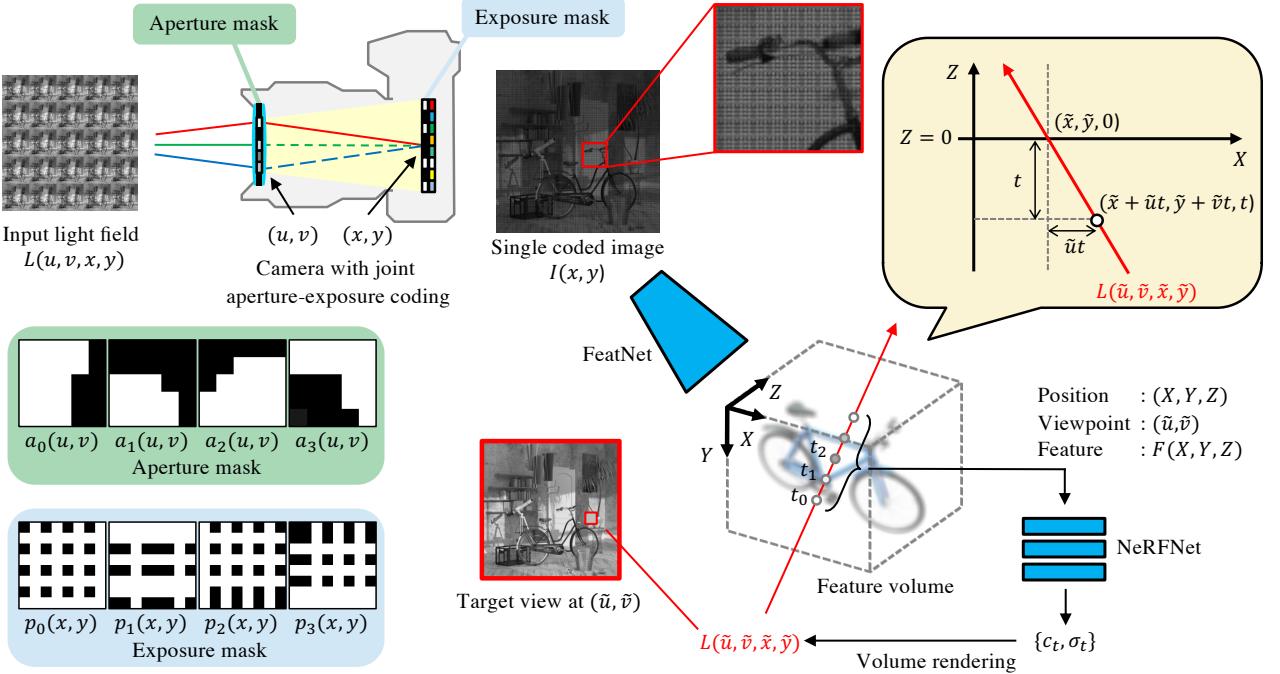
FIGURE 2: Overview of our method. Joint aperture-exposure coding is utilized to obtain single image used as input. FeatNet is responsible for extracting 3-D feature volume from single coded image. Using feature volume with NeRFNet, we perform volume rendering to synthesize views at arbitrary continuous viewpoints.

without any test time optimization. An overview of our method is illustrated in Fig. 2.

## B. JOINT APERTURE-EXPOSURE CODING

We describe the image acquisition method for capturing a coded image that is used as the input to our method. As shown in Fig. 2, all the light rays coming into a camera are considered to constitute a light field, where the intersections of each light ray with the aperture plane and sensor plane are denoted as $(\tilde{u}, \tilde{v})$ and $(\tilde{x}, \tilde{y})$, respectively. Similar to the previous works, the imaging model is constructed in a discretized domain, $(u, v, x, y)$.

If the camera has no coding mechanism for incoming light rays, all the light rays reaching a single pixel $(x, y)$ are summed together to produce a pixel value. The observed image, $I(x, y)$, is given as

$$I(x, y) = \sum_{u,v} L(u, v, x, y). \qquad (1)$$

Note that variations along $(u, v)$ dimensions are simply blurred out in $I(x, y)$, making them difficult to recover.

Coded aperture cameras [15]–[19], [21] have been used to effectively encode $(u, v)$ dimensions. Specifically, a semi-transparent mask, $a(u, v)$, is inserted at the aperture plane to encode the incoming light rays. Each pixel value is the weighted sum of light rays, described as

$$I(x, y) = \sum_{u,v} a(u, v) L(u, v, x, y). \qquad (2)$$

However, it is difficult to accurately reconstruct the light field from a single coded image. Therefore, two or more images taken with different mask patterns are used for better reconstruction quality.

To embed more abundant information into a single observed image, joint aperture-exposure coding [22]–[25] has been investigated recently. Within a single exposure time, both the aperture mask and pixel-wise exposure mask are synchronously controlled to encode the light rays. The image formation model is written as

$$I(x, y) = \sum_{u,v} \left\{ \sum_{k<K} a_k(u, v) p_k(x, y) \right\} L(u, v, x, y) \qquad (3)$$

where $a_k(u, v)$ and $p_k(x, y)$ are the $k$-th aperture and exposure coding patterns, respectively, and $K$ is the number of coding patterns along time. Thanks to the complex coding scheme, a single coded image alone is sufficient for accurate reconstruction of the light field. We adopt Eq. (3) as the image acquisition model to take the input coded image for our method.

Joint aperture-exposure coding is not easy to implement in real camera hardware. As far as we know, only Mizuno et al. [24] demonstrated a hardware implementation that can achieve Eq. (3), but there were some additional restrictions; the image sensor was grayscale, and $p_k(x, y)$ can take only limited patterns. Although their method was originally designed for time-varying scenes, it is also effective for static
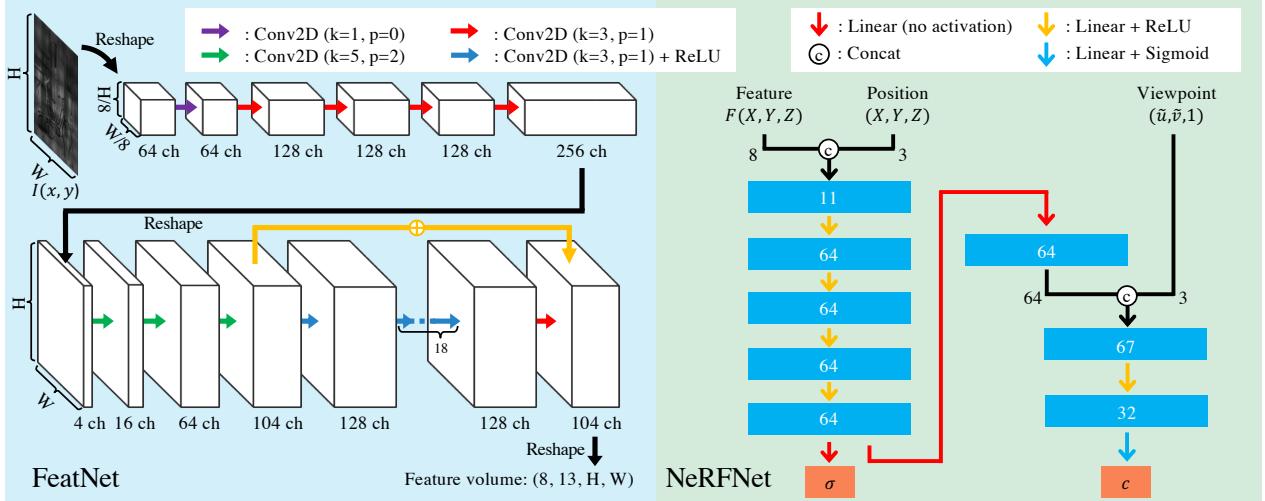
FIGURE 3: Network architectures for FeatNet (left) and NeRFNet (right); k and p mean kernel size and padding.

scenes. We strictly follow the design of Mizuno et al.'s working prototype. We consider grayscale light fields with $5\times5$ views, set $K$ to 4, and adopt the same coding patterns for $a_k(x,y)$ and $p_k(x,y)$ as Mizuno et al.'s. The mask patterns we used are shown in Fig. 2. Accordingly, our method is based on a real existing camera, rather than an idealized hypothetical camera model.

Our goal is to obtain a continuous light field of the target scene, $L(\tilde{u},\tilde{v},\tilde{x},\tilde{y})$, from a single coded image, $I(x,y)$. This is different from Mizuno et al. [24] in which the light field was reconstructed only in the discretized domain as $L(u,v,x,y)$. To achieve our goal, we integrate the idea of neural radiance fields [26], which enables rendering from continuous viewpoints, into the framework of compressive light-field acquisition, as detailed in the next subsection.

## C. RECONSTRUCTION OF CONTINUOUS LIGHT FIELD

As shown in Fig. 2, our method achieves continuous light field reconstruction using two networks: FeatNet and NeRFNet; refer to Fig. 3 for the detailed architectures.

We first mention a network for feature extraction, called FeatNet. This is a deep CNN that takes a single coded image $I \in \mathcal{R}^{H \times W}$ as the input and generates a feature volume $F \in \mathcal{R}^{H \times W \times D \times C}$ as the output.

$$F = \text{FeatNet}(I) \qquad (4)$$

The network architecture for FeatNet is almost the same as the one that Mizuno et al. [24] used for discretized light-field reconstruction. In our method, $F$ is interpreted as a 3-D volume with $D$ depth levels that spans the target scene, and each voxel takes a $C$-channel feature vector. We set $D = 13$ and $C = 8$. The 3-D volume is re-parameterized at the normalized device coordinate (NDC) with $X, Y, Z \in [-1, 1]$ and treated as being continuous. We use trilinear interpolation to enable the querying for the feature vector $\mathbf{f}$ at a

continuous 3-D coordinate $(X, Y, Z)$. We simply describe the query operation as $\mathbf{f} = F(X, Y, Z)$.

We define a neural radiance field (NeRF) for the same $(X, Y, Z)$ volume. It is implemented as a multi-layer perceptron (MLP) that takes the 3-D position $(X, Y, Z)$, angle $(\theta, \phi)$, and feature vector $\mathbf{f}$ as the input and produces the luminance $(c)$ and density $(\sigma)$ as the output.

$$c, \sigma = \text{NeRFNet}(X, Y, Z, \theta, \phi, \mathbf{f}) \qquad (5)$$

Different from the original NeRF [26], our NeRFNet takes as input the feature vector $\mathbf{f}$ extracted from the observed image. Since scene-specific information is given as the feature vector, the network weights of NeRFNet no longer need to be scene-specific but can be generalized over a wide variety of scenes. Moreover, we use a more light-weight MLP than the original NeRF. We found that, provided the feature vector as the input, a small MLP is sufficient to achieve high quality rendering.

Using the features from scene observation is not a new idea in itself; similar ideas have been used for both generalized and scene-specific NeRF-like representations [40], [42], [43], [45], [53], [54]. However, our method is different from these works in that the features are extracted from not multiple images but a single image alone, and the image is acquired with a camera-side coding process to fully capture the 3-D scene information.

We finally mention how we can obtain a light field from the feature volume $F$ and NeRFNet. As shown in Fig. 2, the continuous 3-D space $(X, Y, Z)$ is associated with the continuous light field coordinate $(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})$ via a plane+angle representation; $(\tilde{x}, \tilde{y})$ denotes the position on $Z = 0$, and $(\tilde{u}, \tilde{v})$ is considered the angle $(\theta, \phi)$. A light ray parameterized with $(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})$ is mapped to $(X, Y, Z)$ as

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \mathcal{M}_t(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y}) = \begin{bmatrix} \text{norm}(\tilde{x} + \tilde{u}t) \\ \text{norm}(\tilde{y} + \tilde{v}t) \\ \text{norm}(t) \end{bmatrix} \qquad (6)$$

where norm() denotes the normalization to the NDC, and $t \in [t_{\min}, t_{\max}]$ denotes a sampling position on the light ray. Note that $t_{\min}$ and $t_{\max}$ correspond to the minimum and maximum disparities allowable for the target scene. We set $t_{\min} = -3$ and $t_{\max} = 3$.

Using the relation of Eq. (6), we can query $\mathbf{f}_t$, $c_t$, and $\sigma_t$ for a sampling position $t$ along a light ray parameterized with $(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})$.

$$\mathbf{f}_t = F(\mathcal{M}_t(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})), \tag{7}$$

$$c_t, \sigma_t = \text{NeRFNet}(\mathcal{M}_t(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y}), \tilde{u}, \tilde{v}, \mathbf{f}_t). \tag{8}$$

For the values of $t$, we take 16 random stratified samples for the training time and 32 uniform samples for the test time. Given a set of samples along the ray, $\{c_t, \sigma_t\}$, we perform volume rendering, which is also used in the original NeRF [26], to obtain the luminance of the light ray.

$$L(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y}) = \text{VolumeRendering}(\{c_t, \sigma_t\}) \tag{9}$$

Since the coordinate $(\tilde{u}, \tilde{v}, \tilde{x}, \tilde{y})$ is continuous, we can render arbitrary light rays; thus, we can reconstruct the continuous light field of a target scene.

### D. TRAINING

We train FeatNet and NeRFNet in an end-to-end manner so that the original light field captured by the camera can be accurately reconstructed from the feature volume and NeRFNet. As the training data, we use discretized light fields with $5 \times 5$ views. Therefore, the training loss (we use MSE loss) is computed only for the original $5 \times 5$ viewpoints. However, thanks to the randomized sampling points along the rays (i.e., $t$ is drawn at random in the volume rendering process), our method is optimized over the continuous 3-D space rather than the discretized 3-D grid-points. Moreover, our method is not scene-specific (not requiring per-scene training) but trained over a wide variety of scenes to obtain a generalization capability; once the training has been completed, no test time optimization is necessary for new scenes.

## IV. EXPERIMENTS
We used the BasicLFSR dataset [55], which contains 167 light fields collected from five other datasets (EPFL, HCI new, HCI old, INRIA, and Stanford gantry); 144 light fields were assigned for training, and 23 light fields were reserved for testing. As the training samples, we extracted light field patches, each of which had $120 \times 120$ pixels and $5 \times 5$ views, from the 144 training light fields. We jointly trained FeatNet and NeRFNet over 120 epochs using Adam optimizer, which took 11.7 hours on an NVIDIA GeForce RTX 3090. To validate the effectiveness of the camera-side coding, we also trained two variants without the coding: ours (no-coding), in which Eq. (1) was used as the camera model, and ours (center-only), in which the central view of the light field was used as the input to FeatNet.

We evaluated our method from two perspectives: compressive light-field acquisition and continuous light field reconstruction. Please refer to the supplementary video because the visual quality and 3-D consistency of a light field cannot be well presented on paper.

### A. COMPRESSIVE LIGHT-FIELD ACQUISITION
We compared our method with several state-of-the-art methods for compressive light-field acquisition [19], [21], [24]. Inagaki et al.'s method [19] and Guo et al.'s method [21] are based on coded aperture imaging (Eq. (2)). These methods were retrained on the BasicLFSR dataset until convergence in the same configuration as ours: $5 \times 5$ views were reconstructed from a single coded image. Mizuno et al.'s method [24] adopted joint aperture-exposure coding (Eq. (3)). Since Mizuno et al.'s method was originally designed for a moving scene, it produced $5 \times 5$ views over 4 temporal sub-frames. To acquire a light field for a static scene, we took the average over the temporal domain to reduce them into $5 \times 5$ views. We used two models for this method; one was the pre-trained model provided by Mizuno et al. [24], and the other was re-trained on the BasicLFSR dataset by ourselves. Note that our method can reconstruct **continuous** viewpoints, whereas the others [19], [21], [24] can reconstruct only **discretized** $5 \times 5$ views. In other words, these methods [19], [21], [24] devoted all resources to reconstructing individual $5 \times 5$ views rather than a continuous 3-D representation of the target scene.

Table 1 summarizes the quantitative scores (PSNR and SSIM; greater is better) for the 23 light fields reserved for testing. These 23 light fields were divided into 5 groups corresponding to the source datasets, and we report the average scores for each group and all the groups. Note that the scores were evaluated only for discrete $5 \times 5$ viewpoints, whereas our method can reconstruct continuous viewpoints. As seen from the table, the camera-side coding had a significant impact on the reconstruction quality. Joint aperture-exposure coding (Mizuno et al.'s [24]) significantly outperformed aperture coding (Inagaki et al.'s [19] and Guo et al.'s [21]). Using joint aperture-exposure coding for image acquisition, our method achieved a reconstruction quality comparable to Mizuno et al.'s [24] for the discretized $5 \times 5$ viewpoints. The quantitative scores for Mizuno et al.'s [24] were slightly better than those for ours; this is understandable, because Mizuno et al.'s [24] devoted all the resources to the discrete $5 \times 5$ views without considering other viewpoints. Meanwhile, ours (no-coding) and ours (center-only) produced poor results, showing the difficulty of light field reconstruction from a single image without camera-side coding.

### B. CONTINUOUS LIGHT-FIELD RECONSTRUCTION
We evaluated the capability of our method to reconstruct a continuous light field. For a quantitative evaluation, we used a computer generated scene, *Planets*, provided by Sakai et al. [20]. We used $5 \times 5$ views to compute a coded image by using Eq. (3), from which we reconstructed a light field with $13 \times 13$ views. As shown in Fig. 5, the $13 \times 13$ views include both the original $5 \times 5$ views and interpolated/extrapolated views.

TABLE 1: Quantitative evaluation for compressive light-field acquisition on BasicLFSR test dataset. Reported scores are PSNR/SSIM; greater is better for both of them. Note that our method can reconstruct **continuous** light fields, whereas others [19], [21], [24] can obtain only **discretized** $5 \times 5$ views. Our method achieved comparable quality to Mizuno et al.'s [24] for these discretized viewpoints.

| Method | EPFL | HCI (new) | HCI (old) | INRIA | Stanford | ALL |
|---|---|---|---|---|---|---|
| Inagaki [19] | 29.95/0.913 | 27.85/0.784 | 34.20/0.906 | 31.25/0.918 | 25.12/0.815 | 29.67/0.867 |
| Guo [21] | 30.69/0.920 | 28.49/0.809 | 33.14/0.891 | 33.45/0.938 | 25.78/0.847 | 30.31/0.881 |
| Mizuno (pre-trained) [24] | 32.86/0.932 | 30.86/0.841 | 37.38/0.941 | 35.02/0.939 | **29.27/0.899** | 33.08/0.910 |
| Mizuno (re-trained) [24] | **33.50/0.939** | **31.08/0.844** | **38.02/0.947** | **35.53/0.942** | 29.18/0.891 | **33.46/0.913** |
| Ours | 32.23/0.926 | 30.82/0.837 | 37.49/0.941 | 34.03/0.932 | 29.10/0.888 | 32.73/0.905 |
| Ours (no-coding) | 29.02/0.892 | 27.02/0.760 | 32.69/0.879 | 30.45/0.911 | 23.16/0.773 | 28.47/0.843 |
| Ours (center-only) | 27.14/0.860 | 25.98/0.742 | 31.49/0.856 | 28.47/0.881 | 21.65/0.718 | 26.95/0.812 |



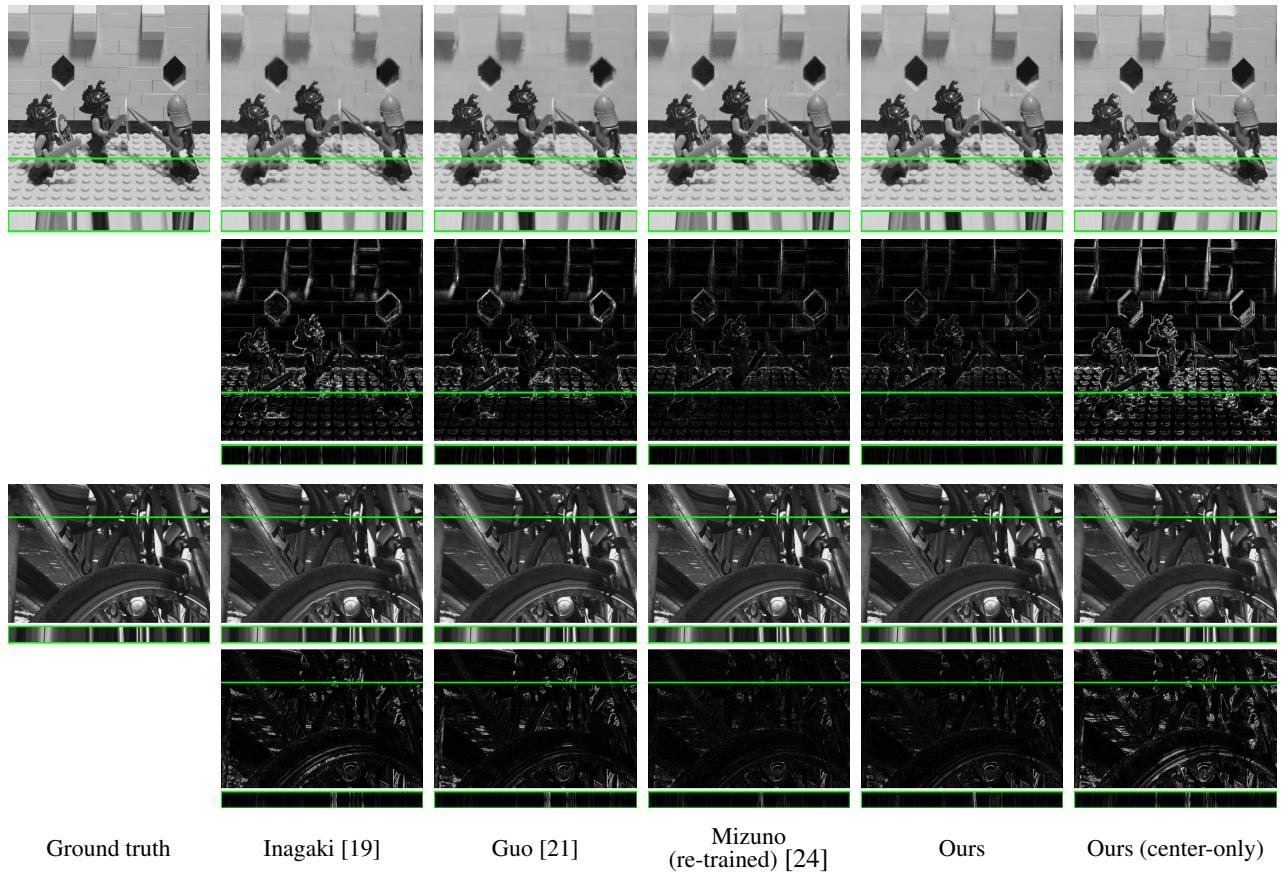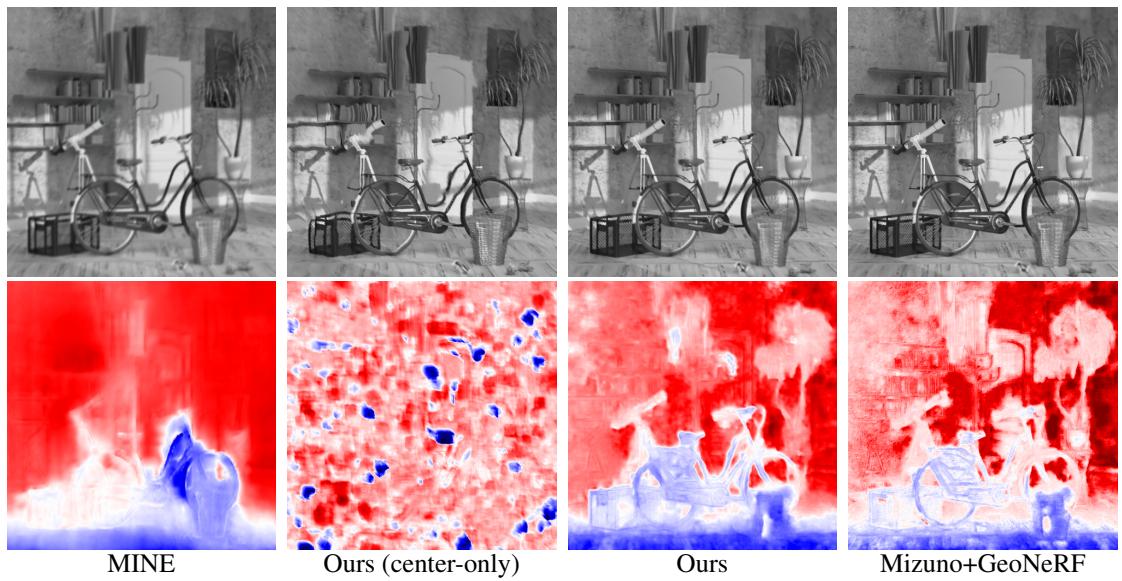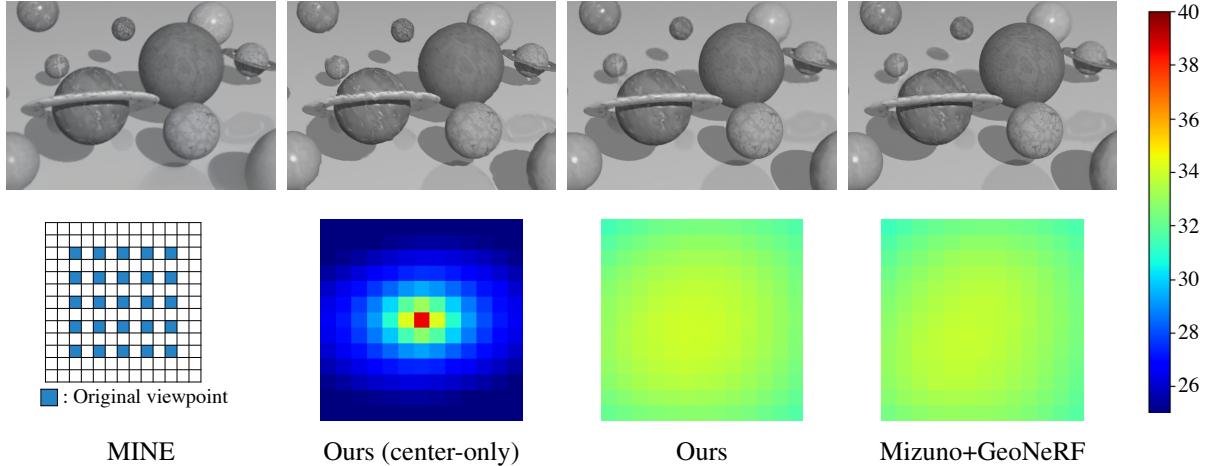| Ground truth | Inagaki [19] | Guo [21] | Mizuno (re-trained) [24] | Ours | Ours (center-only) |

FIGURE 4: Visual results of compressive light-field acquisition for *Lego* (top) and *Bike* (bottom) scenes. For each method, we present top-left view, epipolar plane image (EPI) corresponding to green line, and difference from ground truth ($\times 3$ for better visualization).

For comparison, we tested two other methods for continuous light-field reconstruction. The first is **MINE** [51], a state-of-the-art single-view view synthesis method without test time training. It was reported in [51] that MINE was superior to several other methods [47], [48], [56], [57]. As the input to MINE, we used the central view of a light field without camera-side coding. The other is **Mizuno+GeoNeRF**, which is a brute-force concatenation of two state-of-the-art methods; first, by using Mizuno et al.'s method [24] (the re-trained model), a light field with $5 \times 5$ views was reconstructed

from a single coded image; then, these reconstructed views along with the corresponding camera parameters were fed to GeoNeRF [45], a state-of-the-art NeRF-based rendering method without test time training. For both MINE and GeoNeRF, we used pre-trained network models [1] provided by the corresponding authors. Similar to the case with our method, grayscale images were used as input to these methods.

---

[1] For MINE, we used a model trained on *RealEstate10K* because the ones trained on *Flowers* and *Kitti* were overfitted to a specific category of scenes and performed very poorly for other scenes.

FIGURE 5: Continuous light field reconstruction for *Planets* scene. PSNR values for $13 \times 13$ views are presented as heat maps. PSNR values for MINE cannot be computed, and thus, we present viewpoint configuration in column of MINE.



FIGURE 6: Rendered views (top) and pseudo depth maps (bottom) for *Bicycle* scene. Please refer to supplementary video for better visualization with viewpoint movement and other results.

It should be noted with MINE that single-view view synthesis is scale ambiguous; since only a single view is given as the input, the method cannot know the absolute scale of the scene in principle. Accordingly, we cannot exactly align the rendered viewpoints with those of the target light field.[2] For this reason, we did not evaluate the quantitative scores for MINE. For visual comparison, we manually configured the viewpoints for MINE to obtain a similarly-looking viewpoint arrangement to the target light field.

The quantitative scores (PSNRs) for the $13 \times 13$ views are visualized as heat maps (each grid corresponds to each viewpoint) in Fig. 5. With ours (center-only), the reconstruction quality degraded rapidly as the viewpoint diverged from the central viewpoint, resulting in very poor quality for the marginal viewpoints. This indicates that continuous light-field reconstruction is difficult to achieve without the camera-side coding. In contrast, our method could maintain the quality of all the viewpoints; the reconstruction quality degraded only gradually as the viewpoint diverged from the central viewpoint. Moreover, the reconstruction quality was consistent (not changing abruptly) among the original and the other viewpoints, supporting our claim that our method can reconstruct a light field at continuous viewpoints. The

---

[2]We provided MINE with the ground truth camera parameters of the *Planets* scene, but due to the scale difference, the rendered viewpoints were significantly different from those of the target light field.

TABLE 2: Comparison of computation times (in seconds) between our method and Mizuno+GeoNeRF measured on NVIDIA GeForce RTX 3090. Although average PSNR scores are almost same, our method is significantly more efficient than Mizuro+GeoNeRF.

| | Average PSNR [dB] (13 × 13 views) | 5 × 5 view reconstruction | Feature/geometry computation | Rendering (per view) |
|---|---|---|---|---|
| Ours | **33.01** | - | **0.14** | **0.40** |
| Mizuno+GeoNeRF | 32.93 | 0.12 | 1.58 | 79 |

reconstruction quality of our method was almost the same as that of Mizuno+GeoNeRF, which can be regarded as the best achievable quality with the current state-of-the-art.

Besides the results for the *Planets* scene in Fig. 5, we present visual results for the *Knight* scene in Fig. 6. Along with the rendered views, we present pseudo depth maps, which can be obtained during the process of volume rendering. Although the depth values were not that accurate in particular for textureless regions, they indicated how each method understood the 3-D structure of the target scene. The reconstruction quality of MINE was obviously poor. The pseudo depth map indicates that MINE seriously failed to reconstruct the 3-D shape of the target scene. As the viewpoint moved, we observed that the object shapes were distorted significantly. These observations indicate that light field reconstruction from a normal (uncoded) image alone is an ill-posed problem. For the same reason, ours (center-only) also resulted in significantly poor visual quality. Meanwhile, thanks to the camera-side coding, our method achieved visually-convincing and 3-D-consistent reconstruction. The reconstruction quality of our method was comparable to that of Mizuno+GeoNeRF.

We finally mention the computational advantage of our method over Mizuno+GeoNeRF. We measured the computational times for these methods using the *Planets* scene on an NVIDIA GeForce RTX 3090. As shown in Table 2, our method was substantially more efficient than Mizuno+GeoNeRF. Our method derived a feature volume directly from the coded image with a small amount of computational time (0.14 sec), which was comparable to that for the discretized light field reconstruction in Mizuno+GeoNeRF (0.12 sec). Using the feature volume, our method can perform rendering from continuous viewpoints. Meanwhile, Mizuno+GeoNeRF required another 1.58 sec to compute the geometry before the rendering process. Moreover, the neural rendering process of our method was much more light-weight than that for GeoNeRF; 0.40 sec/view with our method against 79 sec/view with Mizuno+GeoNeRF. Despite being very efficient, our method achieved a rendering quality comparable to (slightly better than, for this scene) Mizuno+GeoNeRF. This can be attributed to our unified design of the feature extraction and neural rendering processes, which enabled end-to-end optimization over the entire pipeline.

## V. CONCLUSION

We proposed a method for reconstructing a continuous light field of a target scene from a single observed image. To this end, we integrated two state-of-the-art techniques into a unified and end-to-end trainable pipeline: joint aperture-exposure coding for compressive light-field acquisition, and a NeRF-based representation for high-quality rendering from continuous viewpoints. Experimental results showed that our method can reconstruct continuous light fields accurately and efficiently without any test time optimization. To our knowledge, this is the first work to bridge two worlds: camera design for efficiently acquiring 3-D information and neural rendering for high-quality view synthesis from continuous viewpoints.

**Limitations and future work**. Although joint aperture-exposure coding is quite effective at obtaining the 3-D information of a target scene, it requires an elaborate hardware implementation. Considering the existence of imaging hardware, we assumed that the target light fields were grayscale. Extension to RGB colors remains as future work, which should involve further hardware development. Moreover, we kept our networks simple for computational efficiency, which put an upper-bound on the reconstruction quality. Using deeper networks would improve the reconstruction quality at the cost of increased computational cost. Finding the balance between quality and efficiency is also an important future direction.

## REFERENCES

[1] 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, 2016.

[2] 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 IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 4748–4757.

[3] K. Maeno, H. Nagahara, A. Shimada, and R.-I. Taniguchi, "Light field distortion feature for transparent object recognition," in IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 2786–2793.

[4] T.-C. Wang, J.-Y. Zhu, E. Hiroaki, M. Chandraker, A. A. Efros, and R. Ramamoorthi, "A 4D light-field dataset and CNN architectures for material recognition," in European Conference on Computer Vision, 2016.

[5] N. K. Kalantari, T.-C. Wang, and R. Ramamoorthi, "Learning-based view synthesis for light field cameras," ACM Transactions on Graphics, vol. 35, no. 6, 2016.

[6] B. Mildenhall, P. P. Srinivasan, R. Ortiz-Cayon, N. K. Kalantari, R. Ramamoorthi, R. Ng, and A. Kar, "Local light field fusion: Practical view synthesis with prescriptive sampling guidelines," ACM Transactions on Graphics, vol. 38, pp. 1–14, 2019.

[7] M. Broxton, J. Flynn, R. Overbeck, D. Erickson, P. Hedman, M. DuVall, J. Dourgarian, J. Busch, M. Whalen, and P. Debevec, "Immersive light

field video with a layered mesh representation," in ACM Transactions on Graphics (Proc. SIGGRAPH), 2020.

[8] G. Wetzstein, D. Lanman, M. Hirsch, and R. Raskar, "Tensor displays: Compressive light field synthesis using multilayer displays with directional backlighting," ACM Transactions on Graphics, vol. 31, no. 4, pp. 1–11, 2012.

[9] F. Huang, K. Chen, and G. Wetzstein, "The Light Field Stereoscope: Immersive Computer Graphics via Factored Near-Eye Light Field Displays with Focus Cues," ACM Transactions on Graphics (SIGGRAPH), no. 4, 2015.

[10] S. Lee, C. Jang, S. Moon, J. Cho, and B. Lee, "Additive light field displays: realization of augmented reality with holographic optical elements," ACM Transactions on Graphics, vol. 35, no. 4, pp. 1–13, 2016.

[11] B. Wilburn, N. Joshi, V. Vaish, E.-V. Talvala, E. Antunez, A. Barth, A. Adams, M. Horowitz, and M. Levoy, "High performance imaging using large camera arrays," ACM Transactions on Graphics, vol. 24, no. 3, pp. 765–776, 2005.

[12] R. Ng, "Digital light field photography," Ph.D. dissertation, Stanford University, 2006.

[13] T. Fujii, K. Mori, K. Takeda, K. Mase, M. Tanimoto, and Y. Suenaga, "Multipoint measuring system for video and sound - 100-camera and microphone system," in IEEE International Conference on Multimedia and Expo, 2006, pp. 437–440.

[14] Y. Taguchi, T. Koike, K. Takahashi, and T. Naemura, "TransCAIP: A live 3D TV system using a camera array and an integral photography display with interactive control of viewing parameters," IEEE Transactions on Visualization and Computer Graphics, vol. 15, no. 5, pp. 841–852, 2009.

[15] A. Veeraraghavan, R. Raskar, A. Agrawal, A. Mohan, and J. Tumblin, "Dappled photography: Mask enhanced cameras for heterodyned light fields and coded aperture refocusing," ACM Transactions on Graphics, vol. 26, no. 3, p. 69, 2007.

[16] C.-K. Liang, T.-H. Lin, B.-Y. Wong, C. Liu, and H. H. Chen, "Programmable aperture photography: multiplexed light field acquisition," ACM Transactions on Graphics, vol. 27, no. 3, pp. 1–10, 2008.

[17] H. Nagahara, C. Zhou, T. Watanabe, H. Ishiguro, and S. K. Nayar, "Programmable aperture camera using LCoS," in European Conference on Computer Vision, 2010, pp. 337–350.

[18] K. Marwah, G. Wetzstein, Y. Bando, and R. Raskar, "Compressive light field photography using overcomplete dictionaries and optimized projections," ACM Transactions on Graphics, vol. 32, no. 4, pp. 1–12, 2013.

[19] Y. Inagaki, Y. Kobayashi, K. Takahashi, T. Fujii, and H. Nagahara, "Learning to capture light fields through a coded aperture camera," in European Conference on Computer Vision, 2018, pp. 418–434.

[20] K. Sakai, K. Takahashi, T. Fujii, and H. Nagahara, "Acquiring dynamic light fields through coded aperture camera," in European Conference on Computer Vision, 2020, pp. 368–385.

[21] M. Guo, J. Hou, J. Jin, J. Chen, and L.-P. Chau, "Deep spatial-angular regularization for light field imaging, denoising, and super-resolution," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 44, no. 10, pp. 6094–6110, 2022.

[22] K. Tateishi, K. Sakai, C. Tsutake, K. Takahashi, and T. Fujii, "Factorized modulation for single-shot light-feild acquisition," in IEEE International Conference on Image Processing, 2021.

[23] E. Vargas, J. N. P. Martel, G. Wetzstein, and H. Arguello, "Time-multiplexed coded aperture imaging: Learned coded aperture and pixel exposures for compressive imaging systems," in International Conference on Computer Vision, 2021.

[24] R. Mizuno, K. Takahashi, M. Yoshida, C. Tsutake, T. Fujii, and H. Nagahara, "Acquiring a dynamic light field through a single-shot coded image," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022.

[25] K. Tateishi, C. Tsutake, K. Takahashi, and T. Fujii, "Time-multiplexed coded aperture and coded focal stack -comparative study on snapshot compressive light field imaging," IEICE Transactions on Information and Systems, vol. E105.D, no. 10, pp. 1679–1690, 2022.

[26] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "NeRF: Representing scenes as neural radiance fields for view synthesis," in European Conference on Computer Vision, 2020.

[27] K. Zhang, G. Riegler, N. Snavely, and V. Koltun, "NeRF++: Analyzing and improving neural radiance fields," arXiv:2010.07492, 2020.

[28] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, "Mip-Nerf: A multiscale representation for anti-aliasing neural radiance fields," International Conference on Computer Vision, 2021.

[29] B. Y. Feng and A. Varshney, "SIGNET: Efficient neural representations for light fields," in International Conference on Computer Vision, 2021.

[30] E. H. Adelson and J. Y. Wang, "Single lens stereo with a plenoptic camera," IEEE transactions on pattern analysis and machine intelligence, vol. 14, no. 2, pp. 99–106, 1992.

[31] S. J. Gortler, R. Grzeszczuk, R. Szeliski, and M. F. Cohen, "The lumigraph," in ACM SIGGRAPH, 1996, pp. 43–54.

[32] W. Matusik, C. Buehler, R. Raskar, S. J. Gortler, and L. McMillan, "Image-based visual hulls," in Proceedings of the 27th annual conference on Computer graphics and interactive techniques, 2000.

[33] C. Buehler, M. Bosse, L. McMillan, S. Gortler, and M. Cohen, "Unstructured lumigraph rendering," in Proceedings of the 28th annual conference on Computer graphics and interactive techniques, 2001.

[34] E. Penner and L. Zhang, "Soft 3D reconstruction for view synthesis," vol. 36, no. 6, 2017.

[35] G. Chaurasia, S. Duchêne, O. Sorkine-Hornung, and G. Drettakis, "Depth synthesis and local warps for plausible image-based navigation," ACM Transactions on Graphics, vol. 32, pp. 30:1–30:12, 2013.

[36] T. Zhou, R. Tucker, J. Flynn, G. Fyffe, and N. Snavely, "Stereo magnification: Learning view synthesis using multiplane images," ACM Transactions on Graphics, vol. 37, pp. 1–12, 2018.

[37] Sara Fridovich-Keil and Alex Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa, "Plenoxels: Radiance fields without neural networks," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022.

[38] S. J. Garbin, M. Kowalski, M. Johnson, J. Shotton, and J. Valentin, "FastNeRF: High-fidelity neural rendering at 200FPS," arXiv preprint arXiv:2103.10380, 2021.

[39] C. Reiser, S. Peng, Y. Liao, and A. Geiger, "KiloNeRF: Speeding up neural radiance fields with thousands of tiny MLPs," in International Conference on Computer Vision, 2021.

[40] A. Chen, Z. Xu, F. Zhao, X. Zhang, F. Xiang, J. Yu, and H. Su, "MVS-NeRF: Fast generalizable radiance field reconstruction from multi-view stereo," in International Conference on Computer Vision, 2021.

[41] J. Chibane, A. Bansal, V. Lazova, and G. Pons-Moll, "Stereo radiance fields (SRF): Learning view synthesis from sparse views of novel scenes," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021.

[42] Q. Wang, Z. Wang, K. Genova, P. Srinivasan, H. Zhou, J. T. Barron, R. Martin-Brualla, N. Snavely, and T. Funkhouser, "IBRNet: Learning multi-view image-based rendering," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021.

[43] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, "pixelNeRF: Neural radiance fields from one or few images," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021.

[44] A. Jain, M. Tancik, and P. Abbeel, "Putting NeRF on a diet: Semantically consistent few-shot view synthesis," in International Conference on Computer Vision, 2021.

[45] M. Johari, Y. Lepoittevin, and F. Fleuret, "GeoNeRF: Generalizing NeRF with geometry priors," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022.

[46] S. Niklaus, L. Mai, J. Yang, and F. Liu, "3d ken burns effect from a single image," ACM Transactions on Graphics, vol. 38, no. 6, pp. 184:1–184:15, 2019.

[47] O. Wiles, G. Gkioxari, R. Szeliski, and J. Johnson, "SynSin: End-to-end view synthesis from a single image," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[48] R. Tucker and N. Snavely, "Single-view view synthesis with multiplane images," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[49] Q. Li and N. Khademi Kalantari, "Synthesizing light field from a single image with variable MPI and two network fusion," ACM Transactions on Graphics, 2020.

[50] A. Trevithick and B. Yang, "GRF: Learning a general radiance field for 3D representation and rendering," in International Conference on Computer Vision, 2021.

[51] J. Li, Z. Feng, Q. She, H. Ding, C. Wang, and G. H. Lee, "MINE: Towards continuous depth MPI with NeRF for novel view synthesis," in International Conference on Computer Vision, 2021.

[52] D. Xu, Y. Jiang, P. Wang, Z. Fan, H. Shi, and Z. Wang, "SinNeRF: Training neural radiance fields on complex scenes from a single image," in European Conference on Computer Vision, 2022.

[53] A. R. Kosiorek, H. Strathmann, D. Zoran, P. Moreno, R. Schneider, S. Mokr'a, and D. J. Rezende, "NeRF-VAE: A geometry aware 3D scene generative model," ArXiv, vol. abs/2104.00587, 2021.

[54] E. R. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein, "Pi-GAN: Periodic implicit generative adversarial networks for 3D-Aware image synthesis," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021.

[55] ZhengyuLiang24, "BasicLFSR," https://github.com/ZhengyuLiang24/BasicLFSR.

[56] P. P. Srinivasan, T. Wang, A. Sreelal, R. Ramamoorthi, and R. Ng, "Learning to synthesize a 4D RGBD light field from a single image," in IEEE International Conference on Computer Vision, 2017, pp. 2262–2270.

[57] S. Tulsiani, R. Tucker, and N. Snavely, "Layer-structured 3D scene inference via view synthesis," in European Conference on Computer Vision, 2018.

TOSHIAKI FUJII received B.E., M.E., and Dr.E. degrees in electrical engineering from the University of Tokyo, Japan, in 1990, 1992, and 1995. In 1995, he joined the Graduate School of Engineering, Nagoya University, where he is currently a professor. From 2008 to 2010, he was with the Graduate School of Science and Engineering, Tokyo Institute of Technology. From 2019 to 2021, he also served as a senior science and technology policy fellow of the Cabinet Office, Government of Japan. His current research interests include multidimensional signal processing, multi-camera systems, multi-view video coding and transmission, free-viewpoint video, and their applications. He is a member of the IEEE Signal Processing Society, the ISO/IEC JTC1/SC29/WG4, WG1 (MPEG-I Visual, JPEG) standardization committee of Japan, the Information Processing Society of Japan, and the Institute of Image Information and Television Engineers of Japan.

· · ·

YUYA ISHIKAWA received his B.E. in electrical engineering from Nagoya University, Japan, in 2021. He is currently a graduate student at the Graduate School of Engineering, Nagoya University, Japan. His research topics are light-field and neural representation.

KEITA TAKAHASHI received B.E., M.S., and Ph.D. degrees in information and communication engineering from the University of Tokyo, Japan, in 2001, 2003, and 2006. He was a project assistant professor at the University of Tokyo from 2006 to 2011 and an assistant professor at the University of Electro-Communications from 2011 to 2013. Since 2013, he has been with the Graduate School of Engineering, Nagoya University, as an associate professor. His research interests include image processing, computational photography, and 3D displays. He is a member of the IEEE Computer Society.

CHIHIRO TSUTAKE received B.E., M.E., and Ph.D. degrees from the University of Fukui, Japan, in 2015, 2017, and 2020. Since 2020, he has been with the Graduate School of Engineering, Nagoya University, as an assistant professor. His research interests include 3D image processing, image coding, optical systems, and applied mathematics.