

DOI: [10.29026/oes.2024.230033](https://doi.org/10.29026/oes.2024.230033)

Edge enhanced depth perception with binocular meta-lens

Xiaoyuan Liu^{1,2,3}, Jingcheng Zhang¹, Borui Leng¹, Yin Zhou¹,
Jialuo Cheng¹, Takeshi Yamaguchi^{4,5,6}, Takuo Tanaka^{4,5,6*} and
Mu Ku Chen^{1,2,3*}

The increasing popularity of the metaverse has led to a growing interest and market size in spatial computing from both academia and industry. Developing portable and accurate imaging and depth sensing systems is crucial for advancing next-generation virtual reality devices. This work demonstrates an intelligent, lightweight, and compact edge-enhanced depth perception system that utilizes a binocular meta-lens for spatial computing. The miniaturized system comprises a binocular meta-lens, a 532 nm filter, and a CMOS sensor. For disparity computation, we propose a stereo-matching neural network with a novel H-Module. The H-Module incorporates an attention mechanism into the Siamese network. The symmetric architecture, with cross-pixel interaction and cross-view interaction, enables a more comprehensive analysis of contextual information in stereo images. Based on spatial intensity discontinuity, the edge enhancement eliminates ill-posed regions in the image where ambiguous depth predictions may occur due to a lack of texture. With the assistance of deep learning, our edge-enhanced system provides prompt responses in less than 0.15 seconds. This edge-enhanced depth perception meta-lens imaging system will significantly contribute to accurate 3D scene modeling, machine vision, autonomous driving, and robotics development.

Keywords: metasurfaces; meta-lenses; deep learning; depth perception; edge detection

Liu XY, Zhang JC, Leng BR et al. Edge enhanced depth perception with binocular meta-lens. *Opto-Electron Sci* **3**, 230033 (2024).

Introduction

Spatial computing¹ and the emerging meta-verse represent a paradigm shift in how humans interact with a machine. Spatial computing refers to integrating digital information and virtual objects into the physical world, creating a mixed reality where the boundaries between the digital and physical realms are blurred. Common augmented reality devices rely on spatial computing to

perceive the depth of the real physical world while embedding virtual objects into real scenes three-dimensionally². One of the key technologies of spatial computing is its depth perception capability, which bridges the gap between the physical and digital realms. This promises intuitive and natural interaction with virtual objects. Therefore, digital information can be correctly placed and manipulated in the scene following physical laws. However, the weight and volume of traditional depth

¹Department of Electrical Engineering, City University of Hong Kong, Hong Kong SAR 999077, China; ²Centre for Biosystems, Neuroscience, and Nanotechnology, City University of Hong Kong, Hong Kong SAR 999077, China; ³The State Key Laboratory of Terahertz and Millimeter Waves, and Nanotechnology, City University of Hong Kong, Hong Kong SAR 999077, China; ⁴Innovative Photon Manipulation Research Team, RIKEN Center for Advanced Photonics, 351-0198, Japan; ⁵Metamaterial Laboratory, RIKEN Cluster for Pioneering Research, 351-0198, Japan; ⁶Institute of Post-LED Photonics, Tokushima University, 770-8506, Japan.

*Correspondence: T Tanaka, E-mail: t-tanaka@riken.jp; MK Chen, E-mail: mkchen@cityu.edu.hk

Received: 26 September 2023; Accepted: 18 December 2023; Published online: 2 April 2024



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License.

To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024. Published by Institute of Optics and Electronics, Chinese Academy of Sciences.

sensing systems result in a lack of comfort in human-computer interaction wearable devices, which contain many sensors (mainly cameras and LiDAR). At the same time, the space occupied by bulky sensors also limits battery life, causing the device to need to be recharged frequently. Advancements in portable and accurate imaging and depth sensing systems are crucial for next-generation human-computer interaction wearable devices.

Complementing spatial computing, binocular meta-lens³ offers a breakthrough approach to depth sensing and imaging with the advantages of being lightweight^{4,5}, thin, and compact. Meta-lens create advanced optical functionalities that surpass the limitations of traditional optics^{6,7}, such as wavefront shaping⁸, polarization control^{9–11}, and spectral manipulation^{12,13}. Meta-lens utilizes nanoantennas to manipulate light¹⁴, offering an opportunity for engineering optical properties such as thinness, flatness, broadband capability¹⁵, high diffraction efficiency¹⁶, extreme depth-of-field¹⁷, and compatibility with complementary metal-oxide-semiconductor (CMOS) technology. By leveraging the unique properties of meta-optics, this compact and miniaturized optical meta-device allows for capturing three-dimensional information from the surrounding environment. Binocular meta-lens enable precise and accurate depth perception, similar to human binocular vision. In recent years, the support of artificial intelligence has increasingly promoted the development of meta-devices in terms of inverse design^{18,19}, prompt data analysis^{20,21}, optical computation^{22,23}, and intelligent reconfigurable meta-devices^{24,25}. These advancements pave the way for compact, lightweight, and highly efficient optical systems seamlessly integrated into spatial computing devices, enhancing their performance and enabling novel applications.

The principle underlying depth acquisition in binocular imaging relies on presenting a stereo-image pair exhibiting discernible disparities²⁶. Disparity denotes the horizontal displacement between corresponding pixels in the left and right images. Traditional binocular disparity computation pipeline often entails the utilization of block matching algorithms for calculating matching losses²⁷. The combination of deep learning and photonics has been widely researched in recent years, encompassing applications such as orbital angular momentum communication²⁸, optical neural networks²⁹, optical encryption³⁰, enhancing holographic data storage (HDS)³¹, photonic inverse design³² and hyperspectral imaging³³.

Nonetheless, convolutional neural networks (CNNs) have garnered greater preference owing to their inherent advantages of rapidity, precision, and operational simplicity in processing. Despite significant advancements in accuracy and speed achieved by various binocular stereo systems, finding accurate corresponding points within inherently ill-posed regions for depth computation remains challenging, such as textureless areas and reflective surfaces³⁴. Ambiguous depth prediction has a serious impact on subsequent machine decision-making. Edge is the typical representation of texture. There must be texture feature points in the edge area for stereo matching. Numerous studies have explored edge detection techniques utilizing meta-lenses, each with distinct characteristics. For instance, the Green function^{35,36}, and spiral phase³⁷ have been employed to enable edge detection using a single meta-lens. Another approach involves utilizing meta-lens arrays for three-dimensional (3D) edge detection³⁸. Polarization control has been leveraged for switchable bright field imaging and edge detection capabilities^{39,40}. Edge detection by the Pancharatnam–Berry phase⁴¹ has emerged as a noteworthy technique, demonstrating potential in quantum applications⁴². Edge-based depth perception offers superior fidelity in the estimation of depth. Within the framework of depth edge views, non-textured regions that lack prominent edges or transitions are efficiently discarded. This filtering process reduces the impact of unreliable or ambiguous depth information originating from textureless regions, thereby enhancing the overall accuracy and reliability of depth estimation. By focusing on edges that signify depth discontinuities, edge-based depth perception provides a more resilient and accurate depiction of depth.

We develop an edge-enhanced depth perception based on binocular meta-lens for spatial computing. The whole system is miniaturized, intelligent, lightweight and compact. Its physical working mechanism consists of a binocular lens, a 532 nm filter, and a CMOS sensor. Each meta-lens, measuring 2.6 mm in diameter, weighs 2.45×10^{-5} g and occupies a volume of 3.98×10^{-6} cm³. The weight of the Sapphire substrate is 0.115 g with a volume of 0.0288 cm³. Thin and flat nature make it simple in both physical system configuration and image processing pipeline. Without preprocessing, the raw captured image is processed directly by our proposed pyramid stereo-matching neural network, H-Net, to obtain the disparity. A novel symmetric H-module with an attention mechanism allows the H-Net to dynamically

allocate resources based on the significance of contextual features of each view and the correlation between the left and right views. With depth-sensing results, an edge enhancement is performed to filter the feature information that detects the 3D space gradients.

Figure 1 demonstrates the edge-enhanced depth perception system schematic with our binocular meta-lens. There are two letter objects in front of the binocular stereo-vision meta-lens. The application scenario shown in Fig. 1 has ill-posed regions, such as the letter objects' unpatterned backgrounds and untextured surfaces. But with the support of a proposed neural network for comprehensive context analysis and a Canny edge detector for filtering, an edge-enhanced depth perception view is realized, perceiving both intensity and depth discontinuities simultaneously.

The convergence of spatial computing and meta-optics holds immense potential for transforming our daily lives. From augmented and virtual reality experiences that blend seamlessly with our physical surroundings to smart glasses that provide personalized information overlays, edge-enhanced spatial computing powered by meta-optics promises to revolutionize how we perceive and interact with the world around us. This integration can lead to breakthroughs in robotics, autonomous systems, underwater exploration, and medical imaging, where accurate depth perception is crucial for navigation, object recognition, and scene reconstruction.

Methods

Simulation and fabrication

We utilize the commercial simulation software COMSOL Multiphysics® to design and analyze the unit cells of the meta-lens. We set periodic boundary conditions for the x and y directions and a perfect match layer (PML) boundary condition for the z -direction. The meta-lens consists of unit cells of gallium nitride (GaN) cylindrical nanopillars on a sapphire substrate. The diameter of the nanopillars varies across the meta-lens. The refractive index of the sapphire substrate is set to 1.77, while the refractive index of GaN at the working wavelength is 2.42. Using this configuration, we calculate the cylindrical nanopillars' simulated transmission spectra and phase shift, as shown in Supplementary information Fig. S1. The meta-atom arrangement layout for fabrication is designed according to the focusing phase distribution

$$\varphi(x, y, \lambda) = - \left[\frac{2\pi}{\lambda} \left(\sqrt{x^2 + y^2 + f^2} - f \right) \right], \quad (1)$$

in which $\varphi(x, y, \lambda)$ is the phase compensation requirement at the (x, y) position under the illumination of wavelength $\lambda = 532$ nm, f is the desired focal length of 10.0 mm. The target diameter of each meta-lens is 2.6 mm.

The proposed binocular meta-lens is fabricated by adopting the following process (see details in Supporting

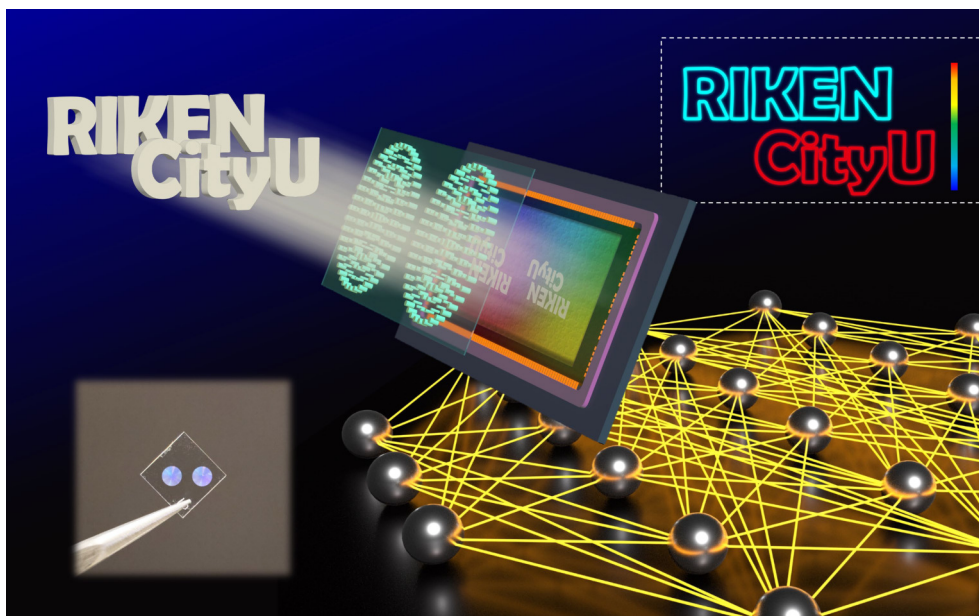


Fig. 1 | Schematic of the edge-enhanced spatial computing with binocular meta-lens. There are two letter objects in front of the binocular meta-lens, which are texture-less and have no background. A binocular meta-lens is designed and fabricated to develop the stereo vision system for texture-less spatial computing scenarios. An edge-enhanced depth perception is realized with the support of a proposed neural network.

Information Fig. S2): A 750-nm-thick GaN is firstly deposited on a sapphire substrate via metalorganic chemical vapor deposition (MOCVD). A 200-nm-thick SiO₂ film, which serves as the hard mask for pattern transfer to the GaN layer with a high aspect ratio, is subsequently deposited using an E-gun evaporator. A PMMA layer is spin-coated on the SiO₂ film, followed by pre-baking at 180 °C for 3 min. A layer of conductive polymer is then spin-coated on the PMMA to avoid charge accumulation. The PMMA layer is exposed under EBL (ELS-HS50, ELIONIX INC.) for pattern definition. After being immersed in DI water to remove the conductive polymer layer, the patterned sample is developed with methyl isobutyl ketone (MIBK)/ isopropyl alcohol (IPA) for 75 s and is rinsed in IPA for 20 s. An additional Cr layer with 40 nm thickness is deposited on the patterned sample using an E-gun evaporator. Followed by the lift-off process in Acetone, the pattern is transferred into the Cr layer. Taking the Cr layer as the hard mask, the SiO₂ layer is etched by inductively coupled plasma reactive ion etching (ICP-RIE) with CF₄ gas. Chromium etchant is adopted to remove the remaining Cr. A second ICP-RIE with a mixture of Ar and Cl₂ is applied for pattern transfer from the patterned SiO₂ film to the GaN film. After removing the residual SiO₂ using a buffered oxide etch (BOE) solution, the desired GaN nanostructure on the sapphire substrate is finally realized.

Figure 2(a) demonstrates the optical image of fabricated binocular meta-lens. The fabrication process of the well structure was characterized based on scanning electron microscope (SEM) images. There is no cracks or pores on the fabricated nanopillars, as shown in the top-view SEM image of Fig. 2(b). From the zoomed-in tilted view of the nanopillar SEM image in Fig. 2(c), the good collimation of the 750-nm high nanopillars can be observed with precise etching. The physical dimension analysis of the binocular sample is divided into two parts: the sapphire substrate and two GaN meta-lens. Each

meta-lens, measuring 2.6 mm in diameter with a volume of 4.25×10^{-6} cm³, weighs 2.61×10^{-5} g, which is lighter than one percent of the weight of a hair. The weight of the sapphire substrate is 0.115 g and occupies a volume of 0.0288 cm³. Even though the sapphire substrate brings much more occupation, the overall weight and volume are still tiny and ignorable.

For disparity computation, we propose a pyramid stereo-matching neural network (named H-Net) with a novel "H"-shaped attention module (H-Module), as shown in Fig. 3(a). The H-Net follows an end-to-end learning framework from stereo input images to disparity map prediction without any other pre- or post-processing. The global context aggregation is vital to derive the disparity information from stereo image pairs. Besides the conventional encoder-decoder architecture and pyramid pooling, H-Net adopts cross-pixel interaction and cross-view interaction to enable the utilization of contextual information and the integration of diverse perspectives (see details in Supplementary information Section 4). Compared with the conventional block matching method⁴³ and two advanced neural networks^{34,44}, H-Net demonstrates significant performance improvements and more comprehensive analysis. (see details in Supplementary information Section 5) With the backbone of PSMNet³⁴, the head of H-Net is a Siamese network⁴⁵, whose two branch networks are weight sharing. These head Siamese CNNs utilize residual blocks⁴⁶ to extract features and weight-sharing spatial pyramid pooling (SPP) modules³⁴ to aggregate context information. The output left and right feature maps from the head backbone (Siamese CNNs) are integrated by the proposed H-Module. The introduction of H-Module with attention mechanism^{47,48} allows the model to dynamically allocate its attention or resources based on the relevance or significance of specific features or contexts. H-Module is a symmetric processing pipeline composed of four cross-pixel interaction blocks and one

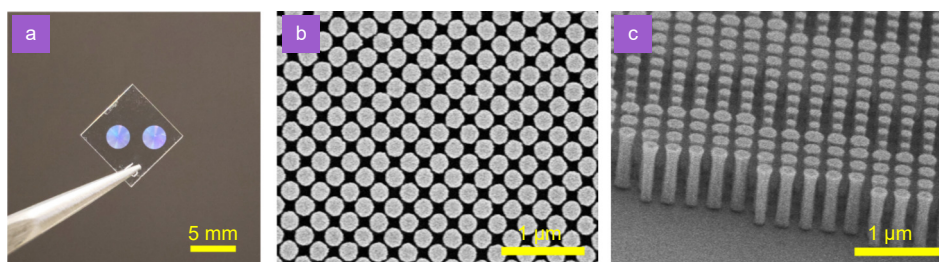


Fig. 2 | Optical and SEM images of fabricated binocular meta-lens. (a) Optical image of the binocular meta-lens. (b) The zoomed-in top-view SEM image of the meta-lens. (c) The zoomed-in tilted-view SEM image at the edge of the meta-lens.

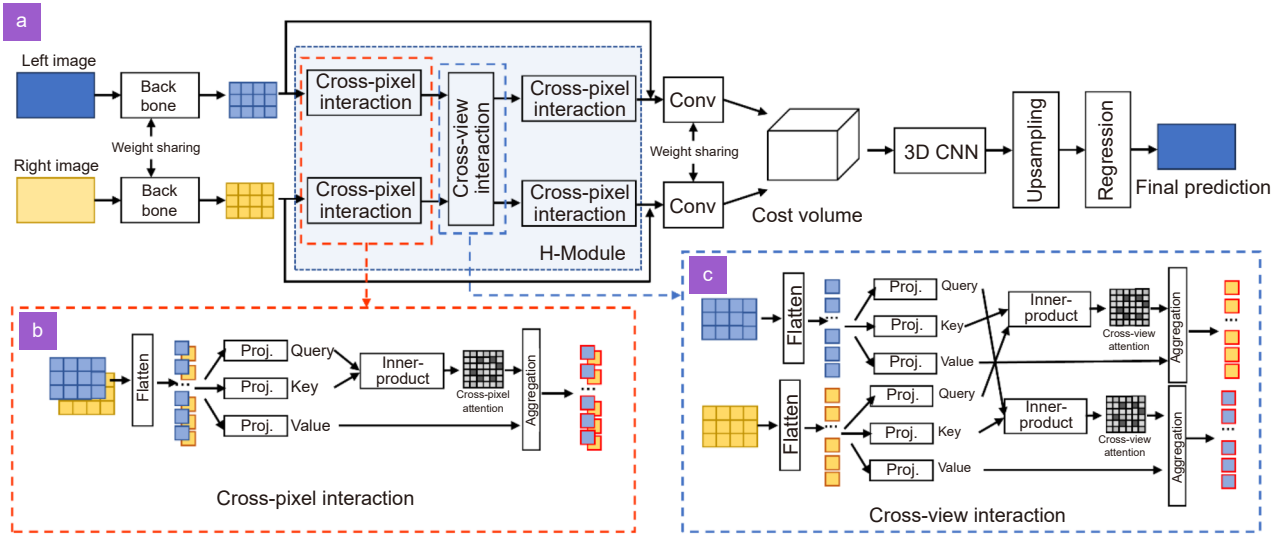


Fig. 3 | Disparity computation with neural network. (a) Architecture overview of proposed neural network H-Net with H-Module. The stereo images are processed by weight-sharing backbones to extract features. These features are then combined using cross-pixel interaction and cross-view interaction in an H-Module. A 4D cost volume is created from the left and right image features, which is then used in a 3D CNN for depth estimation. A disparity regression module is performed before the final disparity map prediction. (b) Detailed pipelines of the cross-pixel interaction. The left and right feature maps are flattened and processed through separate fully connected layers to generate Query, Key, and Value vectors. The inner product is utilized to compute the similarity between Query and Key, resulting in weight coefficients for each Key. These coefficients are used for cross-pixel attention, associating each Key with its corresponding Value. The weighted Values are aggregated to produce enhanced features. (c) Detailed pipelines of the cross-view interaction. The difference from the cross-pixel interaction is the inner product of Key and Query vector comes from different stereo views.

cross-view interaction block. Cross-pixel interaction is the mutual interaction or influence between pixels in an image or visual representation. It involves considering the relationships and dependencies between neighboring pixels to capture contextual information and improve the understanding or analysis of the image. As illustrated in Fig. 3(b), the left and right feature maps are flattened and projected through separate fully connected layers into three essential vectors: Query, Key, and Value. The similarity or correlation between Query and Key is computed using the inner product, yielding weight coefficients for each Key corresponding to its associated Value, known as cross-pixel attention. The Value is then weighted and aggregated based on attention coefficients to obtain enhanced features. Corresponding attention calculation equation⁴⁹ is

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V}, \quad (2)$$

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}, \quad (3)$$

where \mathbf{Q} is the Query vector, \mathbf{K} is the Key vector, \mathbf{V} is the Value vector, $\sqrt{d_k}$ serves as a scale to control the result range, d_k is the dimension of Query vector and Key

vector, and *softmax* is a normalization function utilized to transform a vector of numerical values into a vector of probability distributions. This transformation ensures that the probability associated with each value is directly proportional to its relative proportion within the original vector.

Cross-view interaction refers to the interaction or integration of information from multiple views or perspectives. In multi-view analysis, cross-view interaction aims to leverage information from different viewpoints or modalities to enhance the overall understanding or interpretation of the scene. Detailed processing steps are depicted in Fig. 3(c), which is similar to cross-pixel interaction. The difference is that the calculation of cross-view attention is based on the Query and Key from different features. Specifically, the Query of the left feature map is computed with the Key of the right feature map through inner product and vice versa. This interaction involves feature matching and data fusion, allowing the alignment and combination of information from different views. The attention mechanism enhances the model's ability to capture dependencies, focus on relevant information, and leverage contextual relationships within the visual data (see the ablation study details in Supplementary information Section 5.4 Ablation study).

The enhanced left and right feature maps from H-Module are concatenated as a 4D cost volume. Three repeated encoder-decoder architectures are utilized in the 3D CNN module to further comprehensively understand the contextual information. Before the final prediction of the disparity map, a disparity regression⁵⁰ is performed with a soft attention mechanism. For the disparity map $\mathcal{D} = \{d_a\}_{a=0}^{A_{\max}}$, each final disparity value \hat{d}_a is the original depth value d_a weighted by its probability. The disparity regression is performed as the equation below

$$\hat{d}_a = \sum_{a=0}^{A_{\max}} d_a \cdot \text{softmax}(-c_a), \quad (4)$$

where \hat{d}_a is the final predicted disparity, c_a is the corresponding cost from cost volume for each disparity d_a , a is the annotation number associated with each disparity value d_a in disparity map \mathcal{D} , A_{\max} is the maximum value of a within the range of annotations, softmax function is discussed in Eq. (3). We adopt the smooth L1 loss as the loss function for its fast convergence and robustness to outliers. The final loss is averaged over the N -pixel disparity map, as shown in Eq. (5).

$$\text{Loss}(D, \hat{D}) = \frac{1}{N} \sum_n \text{smooth}_{L1}(d_n - \hat{d}_n), \quad (5)$$

in which

$$\text{smooth}_{L1}(x) = \begin{cases} 0.5x^2, & \text{if } |x| < 1 \\ |x| - 0.5, & \text{otherwise} \end{cases}, \quad (6)$$

where D is the ground truth disparity map, \hat{D} is the predicted disparity map, N is the number of pixels in the disparity map, d_n is the ground truth disparity data for pixel n , and \hat{d}_n is the predicted disparity data for pixel n . H-net was trained on the stereo vision dataset KITTI 2012⁵¹. We employed the Adam Optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$). The learning rate was 0.001 for the first 10 epochs and 0.0001 for the rest. The batch size was 3 on a Nvidia GeForce RTX 3090 GPU. After 800 epochs (64,000 iterations) of training, the final model converged with a training loss of approximately 0.3 (see details in Supporting Information Fig. S8).

The depth map is calculated based on the predicted disparity map. The depth calculation formula³ is

$$\text{depth} = \frac{fb}{ps \cdot \left| \hat{D} + U_{\text{offs}} + O_{\text{offs}} \right|}, \quad (7)$$

in which

$$O_{\text{offs}} = \frac{b}{ps} - |x_1 - x_0|, \quad (8)$$

where focal length f is 10 mm, baseline b is measured 4.056 mm, the side length of the physical pixel on CMOS sensor ps is 3.45 μm , misalignment of lens and sensor on the x -axis U_{offs} is 0, the principal point offset along the x -axis O_{offs} is calculated as -396.6 pixels with x coordinate of left image center $x_0 = 1232$, and the x coordinate of the right image center $x_1 = 2789$ (see more details in Supplementary information Fig. S3). The edge image is derived from the raw captured stereo image with a Canny edge detector⁵², which approximates the first derivative of a Gaussian operator. Through the lower bound cut-off suppression and edge tracking by hysteresis, the detected edges are constrained to be one pixel wide and located at the center discontinuous area without false noise edge points. There are no edges in the non-textured regions in images with uniform intensity distribution. These ill-posed regions will cause ambiguous depth prediction because of the feature-matching calculation mechanism. Under the guidance of the edge image, these ill-posed regions on the depth map are discarded. The edge-enhanced depth perception is the depth map filtered by logical conjunction (AND) operations on edge images. Both the discontinuity of intensity and depth are preserved with high fidelity and accuracy.

Results and discussion

The optical performance of the fabricated meta-lens is derived under 532 nm illumination. The measured intensity profile of left and right meta-lenses along the propagation direction is presented in Fig. 4(a). The corresponding measured focal lengths of left and right meta-lenses are 10.048 mm and 10.046 mm, which matches the designed focal length of 10.0 mm. The diameter of a single meta-lens is 2.6 mm, and the metalens' numerical aperture (NA) is about 0.13. The measured full-width at half-maximum (FWHM) of the focal spots of both meta-lens along X - and Y -axes range from 2.21 to 2.36 μm , with the minimum measurement accuracy of 0.2809 μm per division. Therefore, the averaged FWHM is $2.26 \pm 0.14 \mu\text{m}$, which is close to the diffraction-limited system with an FWHM of 2.1 μm (FWHM = $0.514\lambda/\text{NA}$). The modulation transfer function (MTF), the Fourier transform of the point spread function (PSF), was also calculated, which further confirms that the fabricated meta-lens is a diffraction-limited lens (see more details in Supporting Information Fig. S4). The

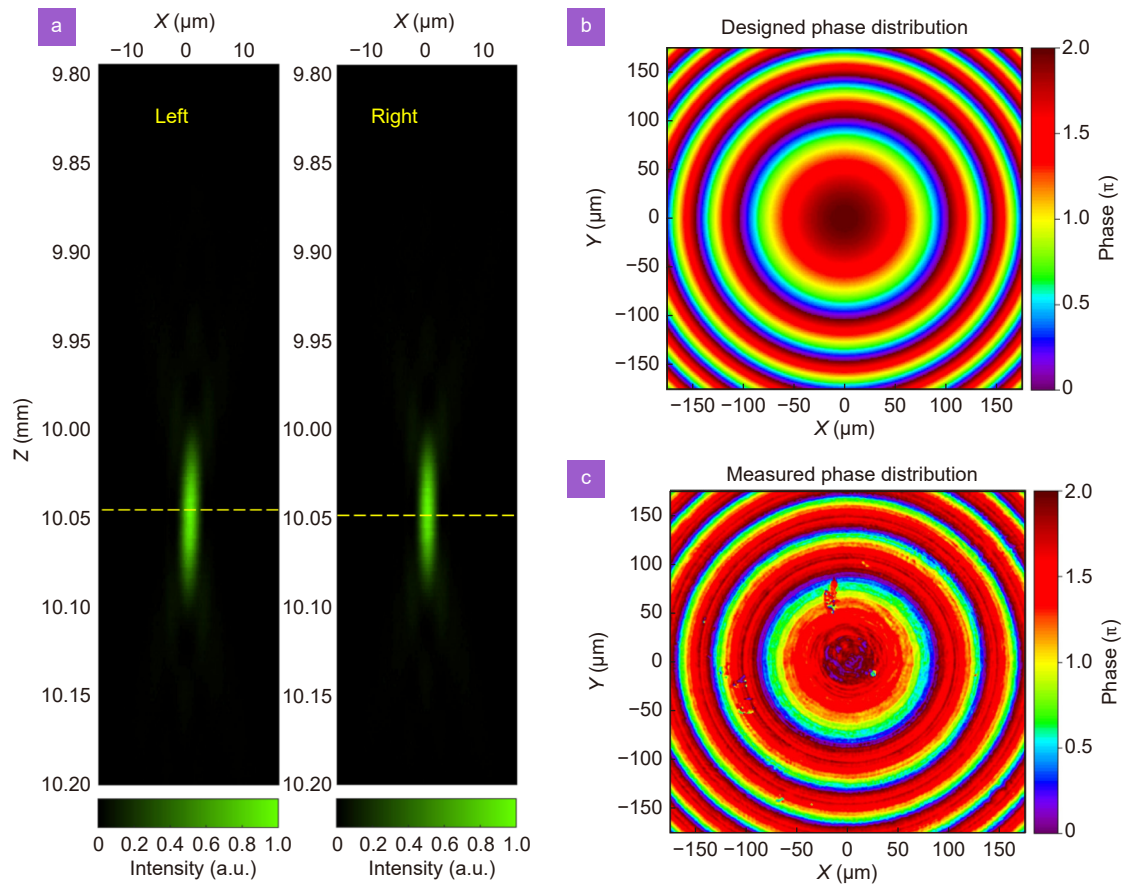


Fig. 4 | Characterization of binocular meta-lens. (a) X-Z plane focusing profiles of left and right meta-lens under 532 nm of wavelength. The measured focal lengths of left and right meta-lenses are 10.048 mm and 10.046 mm, respectively, which are denoted by yellow dashed lines. (b) Designed phase distribution of the meta-lens. (c) Corresponding measured phase distribution of the meta-lens in (b).

measured focusing efficiency is 73.86% at the working wavelength of 532 nm. The focusing efficiency is calculated by dividing the total light power of the focal point area at the focal plane by the total input light power of the bare substrate surface (the selected area is equal to the size of the meta-lens). Several experiments were performed to characterize the 2.6 mm meta-lens using a commercial measurement system (AR-Meta-P, IDEAOPTICS INC.). The phase profile of the fabricated meta-lens was measured to check the agreement between the calculated phase profile and the fabricated phase profile. The detailed experimental setup for meta-lens phase measurement was demonstrated in our previous work⁵³. The simulated and experimental phase distribution maps at the central region of the meta-lens are depicted in Fig. 4(b) and 4(c), respectively, which are in good agreement with each other. The small disparities can be attributed to the fabrication defects and the spherical aberrations in the measurement system. More theoretical and the measured phase profile comparison results are depicted in Supporting Information Fig. S5.

Various imaging and depth sensing experiments are conducted to test the performance of edge-enhanced depth perception of our binocular meta-lens. The configuration of the binocular meta-lens camera for imaging is shown in Fig. S7 in the Supporting Information. Figure 5 demonstrates the raw captured images, depth sensing results, edge-enhanced depth maps, and the integration results of raw images and 3D edges. The raw captured image I_{raw} is cropped from the common stereoscopic region of the left image. Proposed H-Net outputs corresponding disparity map of the stereo images. Through Eq. (5), the depth map \hat{D}_{epth} is calculated accordingly and illustrated in pseudocolor, as shown in the second column of Fig. 5. The 2D edge images that represent the spatial intensity discontinuity are derived from the raw captured image (first column) with the Canny operator. The 2D edge image is converted into a binary matrix E_b , in which the edge pixel is 1, otherwise it is 0. The edge-enhanced depth map DE is calculated by the Hadamard product of the depth map \hat{D}_{epth} and the binary edge matrix E_b , which is similar to a logical conjunction (AND)

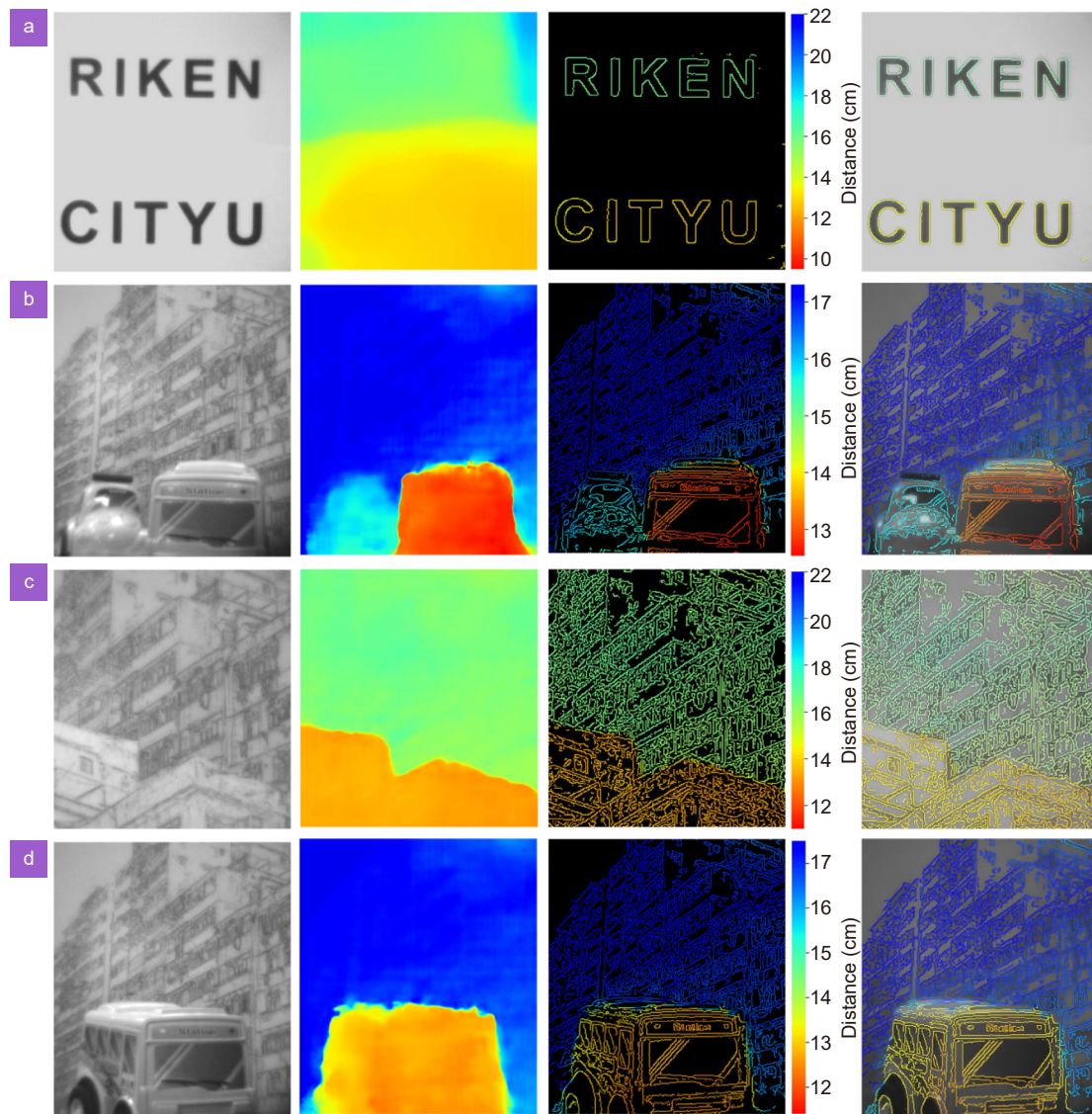


Fig. 5 | Edge-enhanced depth perception of various objects. The first column is the raw left image. The second column is the corresponding depth map. The third column is the edge-enhanced depth map. The second and third columns use the same color bar on the right of the third column. The fourth column is the integration image of the raw image and edge-enhanced depth map. (a) Two pieces of transparent plastic paper printed with "RIKEN" and "CITYU" in black letters are placed at 16.0 cm and 12.8 cm, respectively. (b) A piece of sketch paper printed with a tilted three-dimensional building is located at 17.3 cm as the background. The front ends of the two toy cars are approximately 12.9 cm and 15.7 cm, respectively. (c) The two architectural sketches are at 13.5 cm and 16.5 cm, respectively. (d) The background architecture sketch is positioned at 17.3 cm. The depth of a toy car's body ranges from 12.5 cm to 15.5 cm.

operation. The specific calculation equation is

$$DE = \widehat{D}_{\text{epth}} \odot E_b. \quad (9)$$

The edge-enhanced depth maps DE are displayed in the third column of Fig. 5 in pseudocolor. The non-edge regions with 0 values are set to be black. The integration images I_{integ} in the fourth column of Fig. 5 are merged using the following expression:

$$I_{\text{integ}} = 1.2DE + 0.8I_{\text{raw}}. \quad (10)$$

The integration images aim to demonstrate the fidel-

ity of edge-enhanced depth perception in spatial intensity and depth discontinuity detection.

Figure 5(a) depicts a scenario with ill-posed regions. Two black letter objects, "RIKEN" and "CITYU," printed on transparent plastic papers, are positioned at 16.0 cm and 12.8 cm, respectively. The letter carrier is transparent plastic paper. The background is a white wall without any texture. The absence of texture makes it difficult to establish reliable correspondences between image points in the left and right views, leading to unreliable or ambiguous depth estimates (see the middle region of the

depth map in Fig. 5(a). Such unreliable and ambiguous depth estimates will cause severe trouble for decision-making tasks. In edge-enhanced depth perception, the 3D edge data agree well with the ground truth with the completed preservation of essential details of the scene. Figure 5(b) demonstrates a multi-object traffic scene with two toy cars located at about 12.9 cm and 15.7 cm. An architecture sketch background providing is placed at 17.3 cm. Figure 5(c) shows two architecture sketches with false 3D feelings positioned at 13.5 cm and 16.5 cm, respectively. With edge-enhanced depth perception, the planar false 3D objects do not deceive the system. Figure 5(d) displays a toy car with a continuous depth change, ranging from 12.5 cm to 15.5 cm. All depth sensing results are correct, demonstrating the accuracy capability of our H-Net. The edge-enhanced depth results discard all uniform regions and amplify the 3D feature details with high confidence.

Conclusions

Spatial computing has attracted growing attention from both academia and industry, driven by the rising popularity of the metaverse. A portable and accurate imaging and depth sensing system is of vital importance for next-generation virtual reality devices. In this work, we demonstrate an edge-enhanced depth perception system based on binocular meta-lens, which is intelligent, lightweight, and compact for spatial computing. The miniaturized system contains a binocular meta-lens, a 532 nm filter, and a CMOS sensor. The binocular meta-lens only weighs about 0.115 g with 0.0288 cm³ volume consumption. The imaging system based on our meta-lens minimizes the discomfort caused by the weight and volume of wearable devices to users. We propose a stereo-matching neural network with a novel H-Module for the disparity computation. The H-Module introduces the attention mechanism into the Siamese network. The symmetric architecture with cross-pixel interaction and cross-view interaction enables a more comprehensive analysis of the contextual information in stereo images. The edge enhancement based on the spatial intensity discontinuity discards the ill-posed regions in the image, where ambiguous depth prediction will be generated due to the lack of texture information. With the support of deep learning, our edge-enhanced provides a prompt, intelligent response in less than 0.15 seconds. This edge-enhanced depth perception system will facilitate accurate 3D scene modeling to promote the development of ma-

chine vision, autonomous driving, and robotics.

References

- Greenwald S. Spatial computing (Massachusetts Institute of Technology, Cambridge, 2003).
- Pangilinan E, Lukas S, Mohan V. *Creating Augmented and Virtual Realities: Theory and Practice for Next-Generation Spatial Computing* (O'Reilly Media, Inc. , Sebastopol, 2019).
- Liu XY, Chen MK, Chu CH et al. Underwater binocular metalens. *ACS Photonics* **10**, 2382–2389 (2023).
- Chen MK, Chu CH, Liu XY et al. Meta-lens in the sky. *IEEE Access* **10**, 46552–46557 (2022).
- Jeon D, Shin K, Moon SW et al. Recent advancements of metalenses for functional imaging. *Nano Convergence* **10**, 24 (2023).
- Li T, Chen C, Xiao XJ et al. Revolutionary meta-imaging: from superlens to metalens. *Photon Insights* **2**, R01 (2023).
- Moon SW, Lee C, Yang Y et al. Tutorial on metalenses for advanced flat optics: design, fabrication, and critical considerations. *J Appl Phys* **131**, 091101 (2022).
- Pu MB, Li X, Ma XL et al. Catenary optics for achromatic generation of perfect optical angular momentum. *Sci Adv* **1**, e1500396 (2015).
- Hu YQ, Li L, Wang YJ et al. Trichromatic and tripolarization-channel holography with noninterleaved dielectric metasurface. *Nano Lett* **20**, 994–1002 (2020).
- Wu PC, Sokhoyan R, Shirmanesh GK et al. Near - infrared active metasurface for dynamic polarization conversion. *Adv Opt Mater* **9**, 2100230 (2021).
- Song QH, Baroni A, Wu PC et al. Broadband decoupling of intensity and polarization with vectorial Fourier metasurfaces. *Nat Commun* **12**, 3631 (2021).
- Song MW, Feng L, Huo PC et al. Versatile full-colour nanopainting enabled by a pixelated plasmonic metasurface. *Nat Nanotechnol* **18**, 71–78 (2023).
- Li X, Chen QM, Zhang X et al. Time-sequential color code division multiplexing holographic display with metasurface. *Opto-Electron Adv* **6**, 220060 (2023).
- Guo YH, Pu MB, Zhang F et al. Classical and generalized geometric phase in electromagnetic metasurfaces. *Photon Insights* **1**, R03 (2022).
- Wang SM, Wu PC, Su VC et al. A broadband achromatic metalens in the visible. *Nat Nanotechnol* **13**, 227–232 (2018).
- Zhang F, Pu MB, Li X et al. Extreme-angle silicon infrared optics enabled by streamlined surfaces. *Adv Mater* **33**, 2008157 (2021).
- Fan QB, Xu WZ, Hu XM et al. Trilobite-inspired neural nanophotonic light-field camera with extreme depth-of-field. *Nat Commun* **13**, 2130 (2022).
- Chen MK, Liu XY, Sun YN et al. Artificial Intelligence in Meta-optics. *Chem Rev* **122**, 15356–15413 (2022).
- Krasikov S, Tranter A, Bogdanov A et al. Intelligent metaphotonics empowered by machine learning. *Opto-Electron Adv* **5**, 210147 (2022).
- Chen MK, Liu XY, Wu YF et al. A meta-device for intelligent depth perception. *Adv Mater* **35**, 2107465 (2023).
- Li ZS, Sun JS, Fan Y et al. Deep learning assisted variational Hilbert quantitative phase imaging. *Opto-Electron Sci* **2**, 220023 (2023).
- Liu C, Ma Q, Luo ZJ et al. A programmable diffractive deep neural network based on a digital-coding metasurface array. *Nat Electron* **5**, 113–122 (2022).
- Gao XX, Ma Q, Gu Z et al. Programmable surface plasmonic

- neural networks for microwave detection and processing. *Nat Electron* 6, 319–328 (2023).
24. Li LL, Ruan HX, Liu C et al. Machine-learning reprogrammable metasurface imager. *Nat Commun* 10, 1082 (2019).
 25. Li LL, Zhao HT, Liu C et al. Intelligent metasurfaces: control, communication and computing. *eLight* 2, 7 (2022).
 26. Blake R, Wilson H. Binocular vision. *Vision Res* 51, 754–770 (2011).
 27. Hirschmuller H. Accurate and efficient stereo processing by semi-global matching and mutual information. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)* 807–814 (IEEE, 2005); <http://doi.org/10.1109/CVPR.2005.56>.
 28. Zhou HQ, Wang YT, Li X et al. A deep learning approach for trustworthy high-fidelity computational holographic orbital angular momentum communication. *Appl Phys Lett* 119, 044104 (2021).
 29. He C, Zhao D, Fan F et al. Pluggable multitask diffractive neural networks based on cascaded metasurfaces. *Opto-Electron Adv* 7, 230005 (2024).
 30. Liao MH, Zheng SS, Pan SX et al. Deep-learning-based ciphertext-only attack on optical double random phase encryption. *Opto-Electron Adv* 4, 200016 (2021).
 31. Hao JY, Lin X, Lin YK et al. Lensless complex amplitude demodulation based on deep learning in holographic data storage. *Opto-Electron Adv* 6, 220157 (2023).
 32. Ma TG, Tobah M, Wang HZ et al. Benchmarking deep learning-based models on nanophotonic inverse design problems. *Opto-Electron Sci* 1, 210012 (2022).
 33. Lin CH, Huang SH, Lin TH et al. Metasurface-empowered snapshot hyperspectral imaging with convex/deep (CODE) small-data learning theory. *Nat Commun* 14, 6979 (2023).
 34. Chang JR, Chen YS. Pyramid stereo matching network. In *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* 5410–5418 (IEEE, 2018); <http://doi.org/10.1109/CVPR.2018.00567>.
 35. Zhou Y, Zheng HY, Kravchenko II et al. Flat optics for image differentiation. *Nat Photonics* 14, 316–323 (2020).
 36. Guo C, Xiao M, Minkov M et al. Photonic crystal slab Laplace operator for image differentiation. *Optica* 5, 251–256 (2018).
 37. Kim Y, Lee GY, Sung J et al. Spiral metalens for phase contrast imaging. *Adv Funct Mater* 32, 2106050 (2022).
 38. Chen MK, Yan Y, Liu XY et al. Edge detection with meta-lens: from one dimension to three dimensions. *Nanophotonics* 10, 3709–3715 (2021).
 39. Badloe T, Kim Y, Kim J et al. Bright-field and edge-enhanced imaging using an electrically tunable dual-mode metalens. *ACS Nano* 17, 14678–14685 (2023).
 40. Huo PC, Zhang C, Zhu WQ et al. Photonic spin-multiplexing metasurface for switchable spiral phase contrast imaging. *Nano Lett* 20, 2791–2798 (2020).
 41. Zhou JX, Qian HL, Chen CF et al. Optical edge detection based on high-efficiency dielectric metasurface. *Proc Natl Acad Sci USA* 116, 11137–11140 (2019).
 42. Zhou JX, Liu SK, Qian HL et al. Metasurface enabled quantum edge detection. *Sci Adv* 6, eabc4385 (2020).
 43. Hamid MS, Manap NA, Hamzah RA et al. Stereo matching algorithm based on deep learning: A survey. *J King Saud Univ - Comput Inf Sci* 34, 1663–1673 (2022).
 44. Xu HF, Zhang J, Cai JF et al. Unifying flow, stereo and depth estimation. *IEEE Trans Pattern Anal Mach Intell* 45, 13941–13958 (2023).
 45. Taigman Y, Yang M, Ranzato MA et al. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition* 1701–1708 (IEEE, 2014); <http://doi.org/10.1109/CVPR.2014.220>.
 46. He KM, Zhang XY, Ren SQ et al. Deep residual learning for image recognition. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition* 770–778 (IEEE, 2016); <http://doi.org/10.1109/CVPR.2016.90>.
 47. Li WY, Liu XY, Yuan YX. SIGMA: Semantic-complete graph matching for domain adaptive object detection. In *Proceedings of the 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition* 5281–5290 (IEEE, 2022); <http://doi.org/10.1109/CVPR52688.2022.00522>.
 48. Li WY, Liu XY, Yuan YX. SIGMA++: Improved semantic-complete graph matching for domain adaptive object detection. *IEEE Trans Pattern Anal Mach Intell* 45, 9022–9040 (2023).
 49. Vaswani A, Shazeer N, Parmar N et al. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* 6000–6010 (Curran Associates Inc., 2017).
 50. Kendall A, Martirosyan H, Dasgupta S et al. End-to-end learning of geometry and context for deep stereo regression. In *Proceedings of the 2017 IEEE International Conference on Computer Vision* 66–75 (IEEE, 2017); <http://doi.org/10.1109/ICCV.2017.17>.
 51. Geiger A, Lenz P, Urtasun R. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition* 3354–3361 (IEEE, 2012); <http://doi.org/10.1109/CVPR.2012.6248074>.
 52. Canny J. A computational approach to edge detection. *IEEE Trans Pattern Anal Mach Intell PAMI-8*, 679–698 (1986).
 53. Zhao MX, Chen MK, Zhuang ZP et al. Phase characterisation of metalenses. *Light Sci Appl* 10, 52 (2021).

Acknowledgements

We are grateful for financial supports from the Research Grants Council of the Hong Kong Special Administrative Region, China [Project No. C5031-22G; CityU11310522; CityU11300123], the Department of Science and Technology of Guangdong Province [Project No. 2020B1515120073], City University of Hong Kong [Project No. 9610628] and, JST CREST (Grant No. JPMJCR1904).

Author contributions

XY Liu, T Tanaka, and MK Chen. organized the project. XY Liu, BR Leng, JC Zhang, Y Zhou, JL Cheng, and MK Chen conceived the principle, numerical design, and characterization of the metasurface and meta-system. T Yamaguchi and T Tanaka conceived the fabrication of metasurface. XY Liu, and MK Chen built up the optical system for measurement and the deep learning model and collected the data with experiments for analysis. All authors discussed the results, prepared the manuscripts, and commented on the manuscript.

Competing interests

The authors declare no competing financial interests.

Supplementary information

Supplementary information for this paper is available at <https://doi.org/10.29026/oes.2024.230033>