

17 **This Supplementary information file includes:**

18 Notes S1 to S17:

19 Note S1. Integration of PINS into a MAV platform.

20 Note S2. Ray-tracing verification of the focusing performance of metalenses with different
21 FoV.

22 Note S3. Wide-angle image reconstruction pipeline.

23 Note S4. Model of sensor noise during PSF-based metalens imaging simulation.

24 Note S5. Gamma Transformation for Low-Light Imaging Simulation.

25 Note S6. Depth-of-field analysis of the PINS system.

26 Note S7. Comparison of distortion-corrected experimentally captured image and simulated
27 image using the PSF model.

28 Note S8. Analysis of the effects of imaging degradation on optical-flow estimation.

29 Note S9. Structural details of MMS deep neural network.

30 Note S10. Structural details of DSMSCE module.

31 Note S11. Sequence-wise exponentially weighted L_1 loss for optical flow estimation.

32 Note S12. Characterization of the motion detection limits in PINS.

33 Note S13. Detailed formulation of the Motion Trajectory Prediction (MTP) method.

34 Note S14. Analysis of the influence of fitting parameters on trajectory prediction performance.

35 Note S15. Quantitative RMSE evaluation of trajectory prediction accuracy.

36 Note S16. Flexible trade-off between motion-information redundancy and spatial coverage by
37 tilt-phase design.

38 Note S17. Extension to stereo-vision depth measurement based on metalenses.

39

40 Table S1 to S5

41 Figures S1 to S14

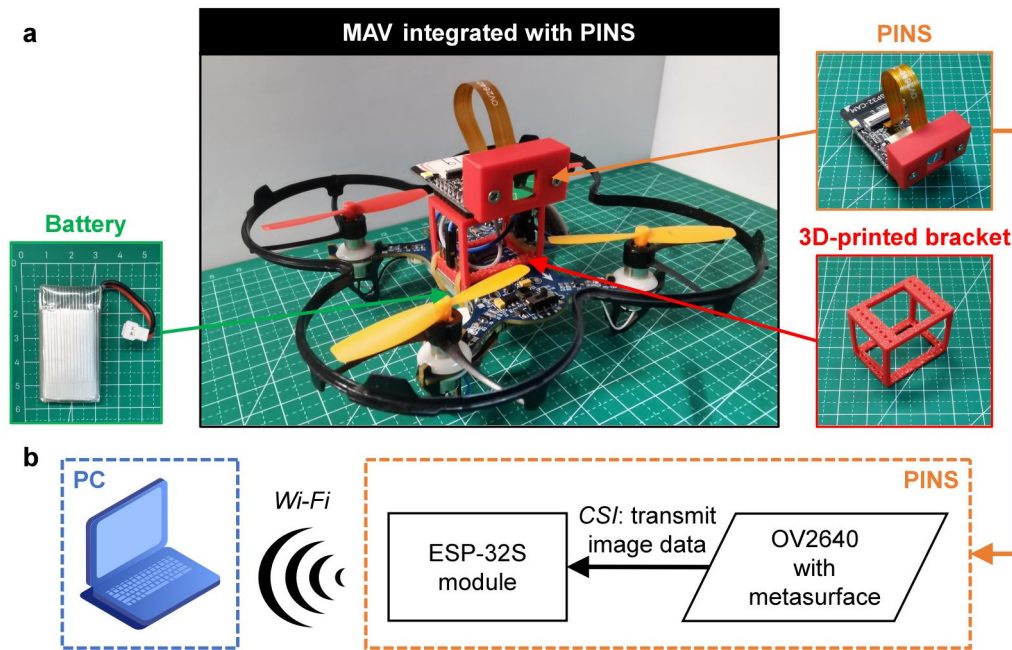
42 Legends for Movies S1 to S2

43

44 **Other Supplementary Materials for this manuscript include the following:**

45 Movies S1 to S2

46 **Note S1. Integration of PINS into a MAV platform.**



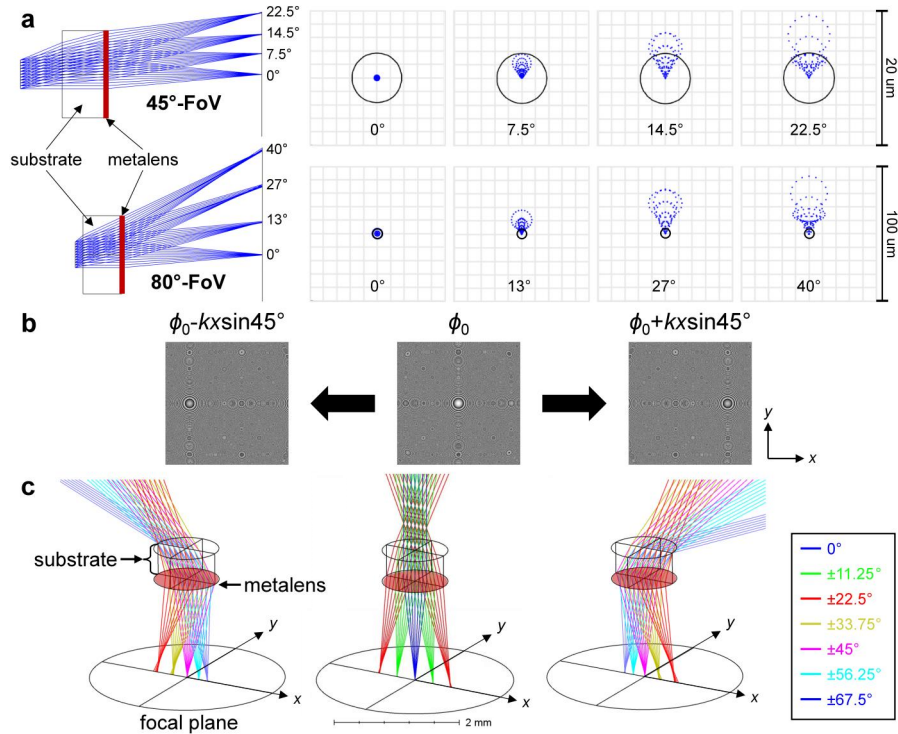
47
48 **Figure S1 | Structural components of the PINS-integrated MAV system and its working principle for**
49 **image data transmission. a**, Structural components of the PINS-integrated MAV system. The system consists
50 of three key components: the PINS module, a 5V lithium battery, and a 3D-printed mounting bracket. **b**,
51 Schematic diagram illustrating the workflow of image data transmission.

52 **Figure S1** illustrates the structural composition of the Planar Intelligent Nanophotonic Sensor
53 (PINS) integrated miniature unmanned aerial vehicle (MAV) system, along with its working
54 principle for image data transmission. As shown in **Fig. S1a**, the integrated system consists of
55 three primary components: the PINS module, a 5V lithium battery, and a custom-designed 3D-
56 printed mounting bracket. Stable and compact integration is achieved through precise mechanical
57 interlocking between the bracket, the PINS module, and the MAV. In addition, elastic bands
58 provide supplementary fixation, further enhancing the structural stability of the system.

59 As illustrated in **Fig. S1b**, the PINS module mounted on the MAV transmits real-time image
60 data to a ground-based PC through a Wi-Fi link. Specifically, the CMOS sensor (OV2640)

61 captures the image formed by the metalens and transfers the image data to the microcontroller
62 unit (MCU, ESP32-S module) through a flexible printed circuit (FPC) using the CSI image
63 transmission protocol. After initialization, the ESP32-S establishes a Wi-Fi connection and
64 launches an HTTP server, enabling the ground terminal to remotely access the real-time data
65 stream through the assigned IP address. Benefiting from its integrated 2.4 GHz Wi-Fi capability
66 and support for high-speed wireless communication (150 Mbps under the 802.11n protocol), the
67 ESP32-S provides an efficient platform for real-time wireless image backhaul in the present
68 system. Under typical operating conditions, the imaging module supports real-time transmission
69 of UXGA-resolution images (1600×1200) at a frame rate of 15 frames per second.

70 **Note S2. Ray-tracing verification of the focusing performance of metalenses with different**
 71 **FoV.**



72
 73 **Figure S2 | Ray-tracing verification of metalenses with different FoV. a,** Comparison of the ray-tracing
 74 results and focal spots for metalenses designed to cover different FoV. The black circles in the focal spot
 75 patterns are the Airy disks of the metalens. **b,** Phase profiles of the side-FoV metalenses are derived from the
 76 center-FoV metalens phase by introducing linear tilt terms of $\pm kx \sin 45^\circ$, respectively. **c,** Zemax ray-tracing
 77 simulations showing the focusing performance of the metalenses under different incident angles.

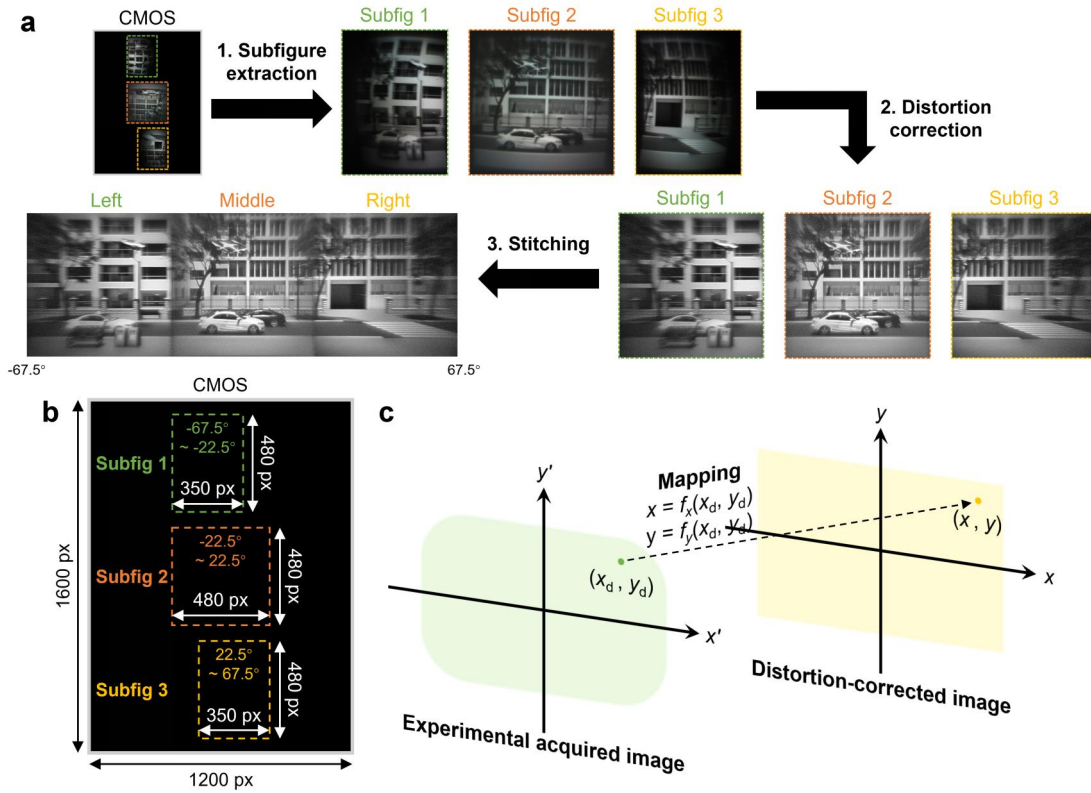
78 To determine a suitable field of view (FoV) for each individual metalens, we first carried out
 79 ray-tracing simulations in Zemax OpticStudio to evaluate and compare the focusing performance
 80 of metalenses designed for different angular ranges, as shown in **Fig. S2a**. In this comparison,
 81 we examined the focal spot distributions at several representative incident angles across the
 82 designed FoV and assessed their quality by comparing them with the corresponding Airy-disk
 83 sizes. The simulation results show that, when the FoV of a single metalens is designed to be 45°,

84 the focusing performance remains well preserved over the entire angular range considered. In
85 particular, the focal spots at representative incident angles stay close to the Airy-disk size,
86 indicating that the metalens can still provide high-quality focusing within this FoV. By contrast,
87 when the target FoV of a single metalens is further enlarged to 80° , the focusing quality
88 deteriorates noticeably, and the degradation becomes especially evident near the edge of the
89 designed angular range. The focal spots under these larger incident angles become significantly
90 broader and deviate more substantially from the Airy-disk limit. Based on this direct comparison,
91 we selected 45° as the FoV of the central metalens.

92 We next verified that the acquisition FoV of an individual metalens can be shifted from the
93 central viewing direction to a side angular sector through the introduction of an additional tilted
94 phase term. Starting from the phase profile of the center-FoV metalens, the phase profiles of the
95 two side-FoV metalenses were generated by superimposing linear tilt terms, following the design
96 principle described in the main text, as illustrated in **Fig. S2b**. From a physical point of view, this
97 additional phase term shifts the effective viewing direction of the metalens, so that incident rays
98 from an off-axis angular range are redirected toward the near-axis focusing regime of the original
99 metalens design and are then focused in a similar manner. In this way, the linear tilted phase
100 does not simply alter the outgoing wavefront direction, but effectively transfers the usable
101 focusing range of the metalens from the central FoV to a designated side FoV. The
102 corresponding ray-tracing results for the side-FoV metalenses are presented in **Fig. S2c**. These
103 simulations show that, after the tilted phase is introduced, the metalenses no longer exhibit their
104 optimal focusing performance around the central viewing direction, but instead within the shifted
105 off-axis angular range. The focal behavior therefore follows the designed change in viewing
106 direction, confirming that the FoV can indeed be reassigned through the tilted-phase modulation.

107 At the same time, the focusing performance within the shifted angular range remains acceptable
108 and comparable to that of the central-FoV metalens over its own designated range. These results
109 verify the validity of the FoV-shifting mechanism used in this work. Therefore, by combining
110 three metalenses, each individually designed to cover a 45° FoV but assigned different viewing
111 directions, the metalens array can jointly provide a total wide-angle imaging FoV of 135° .

112 **Note S3. Wide-angle image reconstruction pipeline.**



113

114 **Figure S3 | Wide-angle image reconstruction pipeline for the metalens array.** **a**, Reconstruction process of
 115 the wide-angle image, including sub-image extraction from the raw CMOS image, distortion correction of each
 116 sub-image, and final image stitching. **b**, Layout of the three sub-images on the CMOS sensor. The dashed
 117 boxes indicate the spatial regions occupied by the sub-images formed by the left, middle, and right metalenses,
 118 respectively. **c**, Schematic illustration of the distortion correction process, where the distorted pixel coordinates
 119 in each sub-image are mapped to corrected coordinates in the rectified image.

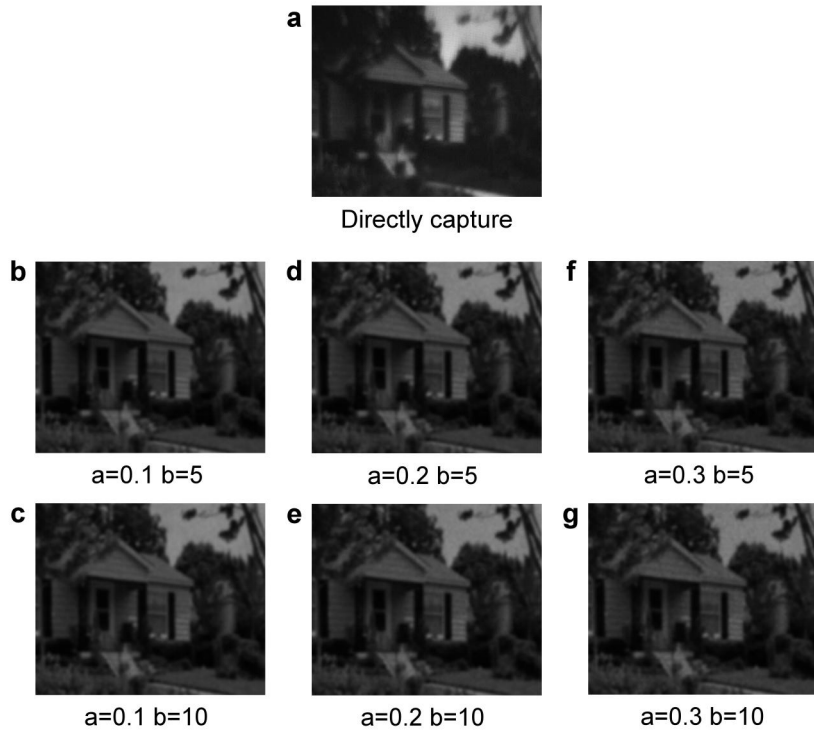
120 To reconstruct a wide-angle image from the raw output of the metalens array, we developed a
 121 pipeline consisting of three sequential steps: sub-image extraction, distortion correction, and
 122 image stitching, as illustrated in **Fig. S3a**.

123 The spatial layout of the sub-images on the CMOS sensor is determined by both the sensor
 124 geometry and the requirement to avoid crosstalk between adjacent imaging channels. The CMOS

125 sensor (OV2640, Omnivision) has a resolution of 1600×1200 pixels. To maximize utilization of
126 the available sensor area while preventing overlap between the images formed by different
127 metalenses, the three sub-images are arranged in a staggered configuration on the CMOS plane,
128 as shown in **Fig. S3b**. In this configuration, each metalens projects its image onto a distinct and
129 non-overlapping region, which suppresses inter-channel crosstalk and facilitates subsequent
130 image processing. Based on this layout, the region of interest corresponding to each metalens is
131 first extracted from the raw CMOS image. Since the three metalenses are designed for different
132 viewing directions, the extracted sub-images exhibit different imaging distortions and cannot be
133 stitched directly in their raw form. Therefore, each cropped sub-image is individually rectified
134 before stitching.

135 The distortion correction is accomplished by establishing a mapping between distorted image
136 coordinates and corrected image coordinates, as illustrated in **Fig. S3c**. To account for both the
137 asymmetric distortion for the side-FoV metalenses and symmetric distortion for the central-FoV
138 metalens, we employ a polynomial fitting approach that establishes an accurate mapping
139 between object points (x, y) and distorted image points (x_d, y_d) . For each metalens, we capture
140 images of a uniform grid of points, fit the polynomial coefficients f_x and f_y to establish this
141 mapping, and then apply it to correct every pixel in the captured image. After distortion
142 correction, the three processed sub-images represent contiguous angular segments of the full
143 scene. These rectified sub-images are then concatenated horizontally to generate the final wide-
144 angle reconstructed image (480×1440 pixels), as shown in **Fig. S3a**.

145 **Note S4. Model of sensor noise during PSF-based metalens imaging simulation.**



146

147 **Figure S4 | Comparison of sensor noise levels during PSF-based metalens imaging simulation.** **a**, The
 148 image captured directly by the CMOS sensor. **b-c**, Simulated imaging results when parameter $a = 0.1$ and b is
 149 set to 5 and 10, respectively. **d-e**, Simulated imaging results when parameter $a = 0.2$ and b is set to 5 and 10,
 150 respectively. **f-g**, Simulated imaging results when parameter $a = 0.3$ and b is set to 5 and 10, respectively.

151 In the PSF-based metalens imaging simulation process, sensor noise arises during CMOS
 152 imaging and primarily originates from two sources: photon shot noise and thermal read noise.
 153 Photon shot noise is caused by fluctuations in photon arrival due to uneven illumination and
 154 inherent temporal randomness. It follows a Poisson distribution, where the number of photons
 155 received at each pixel is modelled as a random variable:

$$P(k_{\text{ph}}, \lambda_{\text{br}}) = \frac{\lambda_{\text{br}}^{k_{\text{ph}}} \times e^{-\lambda_{\text{br}}}}{k_{\text{ph}}!} \quad (\text{S1})$$

156 Here, λ_{br} is the expected number of photon arrivals, which increases with the image brightness,
157 and k_{ph} is the actual number of photons received. Photon shot noise becomes more significant in
158 low-light images.

159 Thermal noise is the noise generated by the sensor during signal readout process. It is
160 associated to the sensor's temperature and the operating state of the electronic components.
161 Typically, thermal noise follows a zero-mean Gaussian distribution.

162 To facilitate the noise modeling, we assume that the noisy image \mathbf{Y} can be decomposed into
163 the following formula:

$$\mathbf{Y}(\mathbf{r}) = y(\mathbf{r}) + \eta_{\text{photon}}[y(\mathbf{r})] + \eta_{\text{thermal}}(\mathbf{r}) \quad (\text{S2})$$

164 where $\eta_{\text{photon}}[y(\mathbf{r})]$ represents photon shot noise, dependent on local pixel value $y(\mathbf{r})$, and
165 $\eta_{\text{thermal}}(\mathbf{r})$ represents brightness-independent thermal noise. By combining both components, the
166 total noise model can be formulated as:

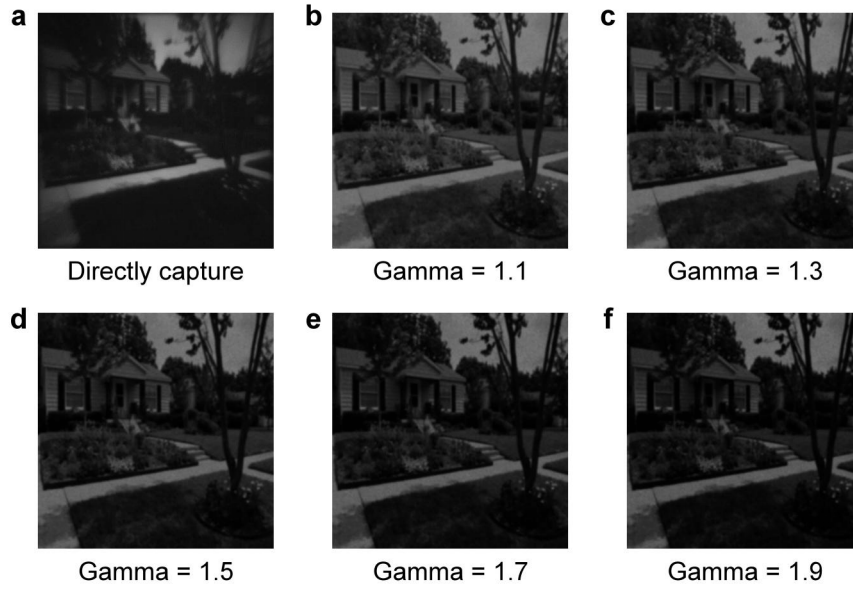
$$z(\mathbf{r}) \sim N(\mu = y(\mathbf{r}), \sigma^2 = ay(\mathbf{r}) + b) \quad (\text{S3})$$

167 where, μ represents the original pixel value of the image, σ^2 is the variance of the noise, and a, b
168 are parameters corresponding to the contributions of photon shot noise and thermal noise,
169 respectively. As image brightness decreases, the effect of photon shot noise becomes more
170 significant, while the effect of thermal noise remains relatively small.

171 **Figure S4** presents a comparison of simulated imaging results under different noise parameter
172 settings. When parameter a is fixed and b is doubled, little change is observed in the resulting
173 results. This is because the metalens imaging results are relatively dark, rendering the variation
174 in thermal noise less impactful. In contrast, photon shot noise has a greater influence on the
175 image. As shown in **Fig. S4a, S4c** and **S4e**, increasing a from 0.1 to 0.3 amplifies the noise effect

176 in the simulations. By comparing these simulated imaging results with real captured images, the
177 final noise parameters for dataset generation are determined to be $a = 0.2$, $b = 10$.

178 **Note S5. Gamma Transformation for Low-Light Imaging Simulation.**



179

180 **Figure S5 | Effect of different Gamma parameters on the simulated imaging results.** **a**, Metalens imaging
 181 captured directly by the CMOS sensor. **b-f**, Simulated imaging results with Gamma parameters set to 1.1, 1.3,
 182 1.5, 1.7, and 1.9, respectively. To better simulate the overall intensity attenuation during metalens imaging, the
 183 brightness of all images was reduced by 20%.

184 After applying PSF convolution and noise augmentation to simulate the imaging system, a
 185 nonlinear darkening transformation is required to capture the response characteristics of both the
 186 image sensor and the human visual system under low-light conditions, particularly in regions of
 187 medium to high brightness. This is achieved by applying a power-law (Gamma) transformation
 188 to the normalized pixel intensity. Let the normalized input pixel intensity after the previous
 189 processing steps be $\mathbf{I}_{\text{input}}(\mathbf{r})$ (ranging from $[0,1]$), then the output pixel $\mathbf{I}_{\text{output}}(\mathbf{r})$ is given by:

$$\mathbf{I}_{\text{output}}(\mathbf{r}) = \Gamma[\mathbf{I}_{\text{input}}(\mathbf{r})] = [\mathbf{I}_{\text{input}}(\mathbf{r})]^\gamma \quad (\text{S4})$$

190 This transformation serves two purposes in low-light simulation. First, when $\gamma > 1$, it
 191 compresses the range of brighter regions while enhances details in the darker areas. Second, the

192 power-law curve closely approximates the typical response curve of camera sensors, allowing for
193 a more realistic representation of imaging characteristics under low-light conditions. In **Fig. S5**,
194 we present the effects of different Gamma parameters on the final simulated imaging results. It
195 can be observed that when γ is close to 1, the image undergoes only mild darkening, with darker
196 regions remaining clearly visible. However, when γ approaches 2, the overall brightness
197 significantly drops, better approximating extremely low-light conditions, though this may lead to
198 the loss of certain image details. Therefore, we select $\gamma = 1.7$ as the final setting for simulating
199 metalens imaging.

200 **Note S6. Depth-of-field analysis of the PINS system.**

201 The depth of field (DoF) of the PINS system can be estimated using a geometric-optics model
202 based on image-plane shift. According to the Gaussian lens formula, when the object distance u
203 changes, the corresponding image distance v can be expressed as:

$$\frac{1}{f} = \frac{1}{u} + \frac{1}{v} \quad (\text{S5})$$

204 where f is the focal length of the metalens. When the image distance changes from v to v' , the
205 corresponding object distance u' changes to:

$$u' = \frac{f \cdot v'}{v' - f} \quad (\text{S6})$$

206 An image is considered acceptably sharp if the image plane deviation does not exceed the
207 maximum allowable shift Δv_{\max} , which is determined by the permissible circle of confusion c and
208 the F-number of the lens: $\Delta v_{\max} = c \times F/\#$. In the present metalens array, each metalens has a
209 diameter of $D = 1$ mm and a focal length of $f = 1.6$ mm, yielding a F-number of $F/\# = f/D = 1.6$.
210 The CMOS sensor has a pixel size of $p = 2.2$ μm . Considering that the permissible circle of
211 confusion c is typically taken as $2 \times p$, Δv_{\max} is then set as 7.04 μm , which means that if the image
212 plane deviation does not exceed ± 7.04 μm , the image remains acceptably sharp. For the case
213 where the metalens-CMOS distance v is set to f , the corresponding near limit of the DoF can be
214 calculated. The corresponding object distance u' satisfies:

$$u > u_{\text{near}} = \frac{f \cdot v'}{v' - f} = \frac{f \cdot (f + \Delta v_{\max})}{\Delta v_{\max}} = 0.37 \text{ m} \quad (\text{S7})$$

215 Therefore, with the image distance fixed at the focal length, the system maintains acceptable
216 sharpness for object distances from approximately 0.37 m to infinity.

217 The two adjustment screws are used for fine alignment of the metalens-to-CMOS distance, so
 218 that the actual image distance can be tuned as close as possible to the designed focal length. In
 219 practice, however, due to fabrication and assembly errors, the actual image distance may still
 220 deviate slightly from the ideal focal length. Assuming that the actual metalens-to-CMOS
 221 distance is v_0 , the corresponding DoF range is determined by the near and far limits associated
 222 with the allowable image-plane deviation:

$$u_{\text{near}} = \frac{f \cdot (v_0 + \Delta v_{\text{max}})}{v_0 + \Delta v_{\text{max}} - f}, u_{\text{far}} = \frac{f \cdot (v_0 - \Delta v_{\text{max}})}{v_0 - \Delta v_{\text{max}} - f} \quad (\text{S8})$$

223 where $u_{\text{far}} = \infty$ when $v_0 - \Delta v_{\text{max}} \leq f$. **Table S1** summarizes the resulting DoF ranges for several
 224 representative values of v_0 relative to the nominal focal length $f = 1.6$ mm.

225 **Table S1 | Depth of field ranges for different image distances.**

Image distance v_0	Near limit object distance u_{near} (m)	Far limit object distance u_{far} (m)
f	0.37	∞
$f + 10 \mu\text{m}$	0.15	0.87
$f + 15 \mu\text{m}$	0.12	0.32
$f + 20 \mu\text{m}$	0.10	0.20
$f + 25 \mu\text{m}$	0.08	0.14

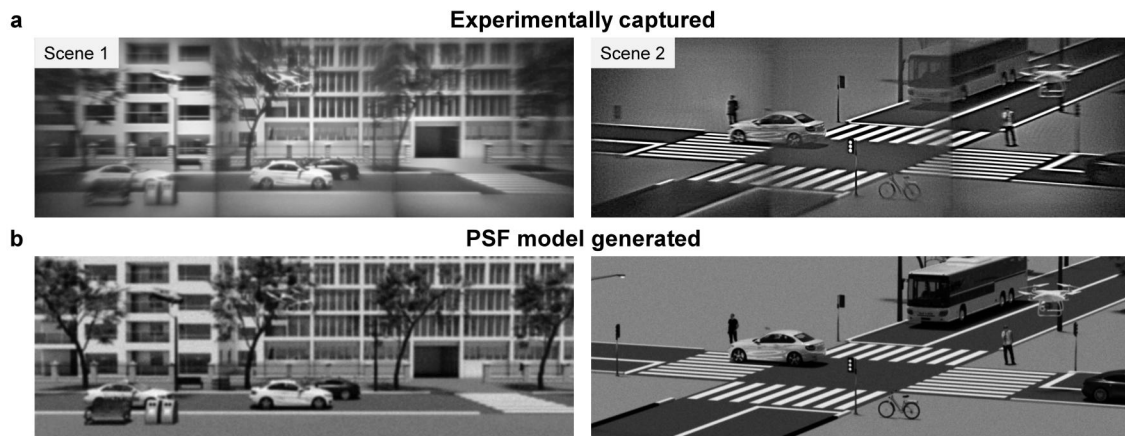
226 These results show that the depth of field is maximized when the image distance is set as close
 227 as possible to the focal length. As the image plane shifts away from the focal plane, the in-focus
 228 object-distance range rapidly narrows. The screw mechanism provides continuous depth
 229 adjustment for focus alignment, with a resolution of approximately $1.1 \mu\text{m}$ per 1° of rotation,
 230 offering sufficient positioning accuracy for practical focus tuning within the acceptable tolerance
 231 range.

232 **Note S7. Comparison of distortion-corrected experimentally captured image and simulated**
233 **image using the PSF model.**

234 First, we quantified the similarity between the corresponding image pairs in **Fig. 2e** by
235 calculating the structural similarity index measure (SSIM), mean absolute error (MAE), and
236 mean squared error (MSE), as summarized in **Table S2**.

237 **Table S2 | Quantitative similarity metrics for the corresponding image pairs shown in Fig. 2e.**

	SSIM	MAE	MSE
Image 1	0.741802	0.047315	0.007987
Image 2	0.791876	0.044247	0.008249
Image 3	0.803224	0.031747	0.005888



238
239 **Figure S6 | Comparison of distortion-corrected experimentally captured 135°-FoV image and simulated**
240 **135°-FoV image using PSF model. a,** Experimental captured image by the metalens array after reconstruction.
241 **b,** Simulated image of the same scene generated using only the central 45° PSF.

242 To further verify the validity of the PSF-based convolution imaging model for wide-FoV
243 motion sensing, we compared distortion-corrected experimentally captured wide-FoV images
244 with simulated wide-FoV images generated using only the experimentally measured central PSF.
245 The experimental wide-FoV images were reconstructed by sub-image extraction, distortion
246 correction, and image stitching. The simulated wide-FoV images were generated using the PSF
247 captured by the central-FoV metalens under normal illumination for the same scenes.

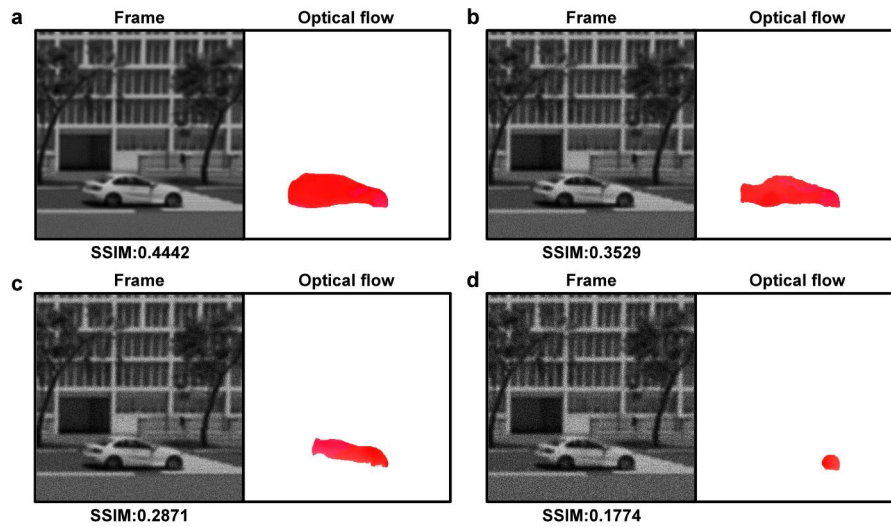
248 As shown in **Fig. S6**, the simulated images agree well with the reconstructed experimental
 249 images in overall scene structure, local blur, brightness distribution, and object contours. To
 250 further quantify the similarity, we calculated the SSIM, MAE, and MSE for two representative
 251 scenes, as summarized in **Table S3**. The results indicate that, although some residual differences
 252 remain due to experimental noise, illumination variation, and reconstruction error, the central-
 253 PSF-based model can still reproduce the main image characteristics of the wide-FoV system with
 254 reasonable accuracy. These results support the use of the central PSF as a practical
 255 approximation for training-data generation in the present wide-FoV motion-sensing framework.

256 **Table S3 | Quantitative similarity metrics for the corresponding image pairs shown in Fig. S6.**

	SSIM	MAE	MSE
Scene 1	0.6476	0.1394	0.0387
Scene 2	0.7164	0.0913	0.0172

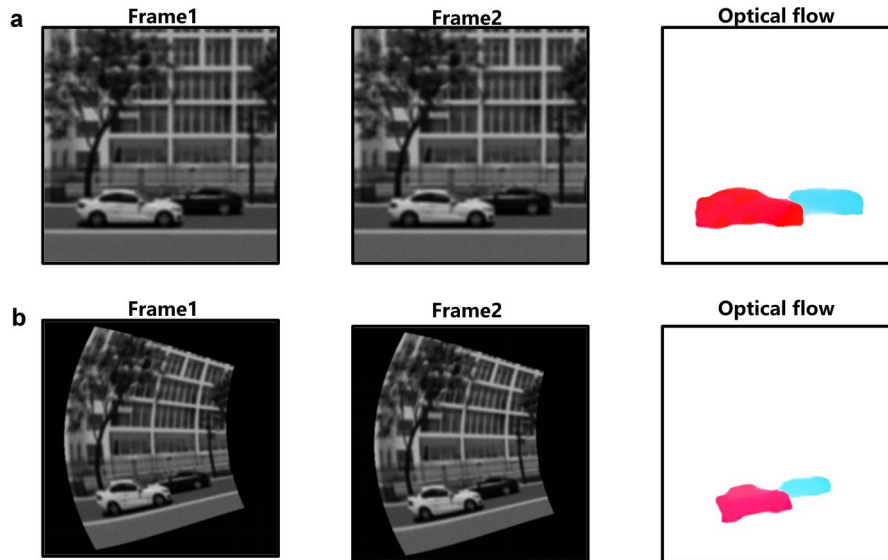
257

258 **Note S8. Analysis of the effects of imaging degradation on optical-flow estimation.**



260 **Figure S7 | Effect of imaging noise on optical-flow estimation.** Representative optical-flow results under
261 different imaging noise levels, with SSIM values of (a) 0.4442, (b) 0.3529, (c) 0.2871 and (d) 0.1774.

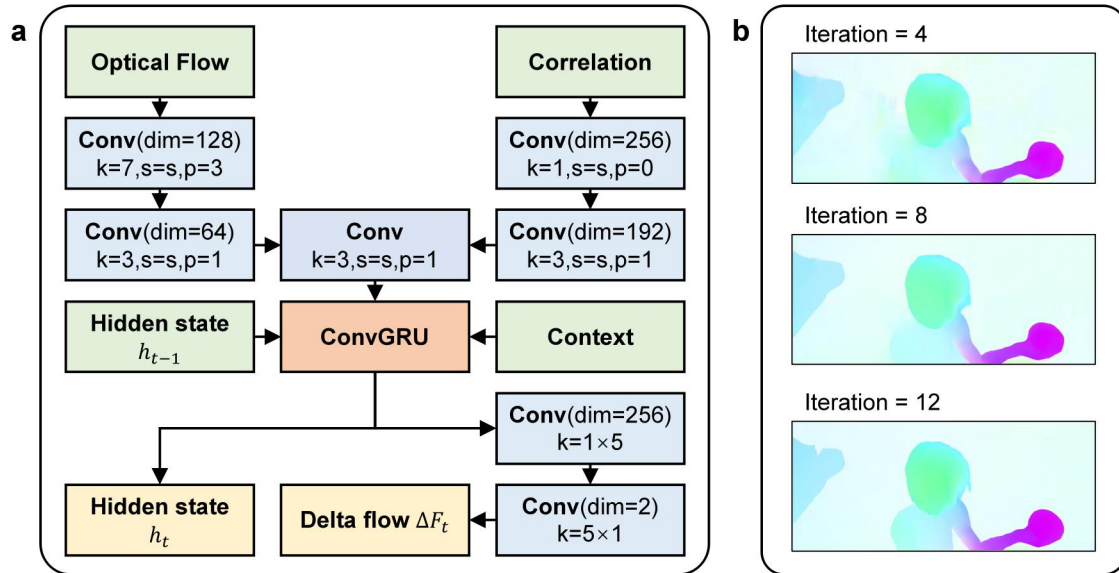
262 To evaluate the effect of imaging noise on optical-flow estimation, we introduced different levels
263 of noise into the input images and compared the corresponding motion-estimation results, as
264 shown in **Fig. S7**. The image quality is quantified by the structural similarity index (SSIM),
265 which decreases from 0.4442 in **Fig. S7a** to 0.1774 in **Fig. S7d** as the noise level increases. It
266 can be seen that increasing noise gradually degrades the completeness and accuracy of the
267 estimated motion field. When the noise becomes severe, the detected motion region shrinks
268 significantly and partial motion information is lost. These results indicate that imaging noise
269 directly affects the reliability of optical-flow estimation and thus can degrade subsequent motion
270 analysis performance.



271
 272 **Figure S8 | Effect of image distortion on optical-flow estimation.** Comparison of optical-flow results for (a)
 273 undistorted and (b) distorted image pairs.

274 To evaluate the influence of image distortion on optical-flow estimation, we compared the
 275 motion-estimation results obtained from undistorted and distorted image pairs, as shown in **Fig.**
 276 **S8**. In **Fig. S8a**, the input images are free of geometric distortion, and the corresponding optical-
 277 flow result accurately preserves the shape and relative position of the moving objects. In **Fig.**
 278 **S8b**, geometric distortion is introduced into the input image pair, leading to visible deformation
 279 of the scene structure. As a result, the estimated optical-flow vectors exhibit errors in both
 280 direction and magnitude. Although subsequent geometric rectification can correct the spatial
 281 position of the motion field, it cannot fully recover the distortion-induced errors in the optical-
 282 flow direction and magnitude. This result indicates that image distortion not only affects image
 283 geometry, but also degrades the physical accuracy of optical-flow estimation, further justifying
 284 the distortion-correction step in our wide-angle image reconstruction pipeline before motion
 285 analysis.

286 **Note S9. Structural details of MMS deep neural network.**



287
 288 **Figure S9 | Iterative update block of Meta-Motion Sense (MMS) neural network.** **a**, Schematic of the
 289 process flow and structure for a single optical flow update. The input to the module comprises four
 290 components: the current optical flow estimate, correlation features, the hidden state from the previous iteration,
 291 and context features. The module outputs the updated hidden state and the incremental optical flow correction.
 292 **b**, Impact of the number of iterations on the accuracy of optical flow estimation. The input data is generated
 293 using PSF-based simulated imaging, with specific comparison values shown in **Table S4**.

294 **Table S4 | Accuracy of optical flow estimation under different iterations.**

	EPE	APE	1-pixel	3-pixel	5-pixel
Iteration=4	5.86	6.53	0.87	0.96	0.98
Iteration=8	4.11	4.69	0.89	0.98	0.99
Iteration=12	3.22	2.69	0.92	0.99	0.99
Iteration=16	3.22	2.69	0.92	0.99	0.99

295 **Note:** EPE (End Point Error): the Euclidean distance between the predicted and ground truth flow vectors at
 296 each pixel. APE (Angular Projection Error): the angular error of the flow vectors, measuring the deviation
 297 between the predicted and ground truth flow directions. The 1-pixel, 3-pixel, and 5-pixel metrics indicate the

298 percentage of estimated flow vectors whose end-point error is within 1, 3, or 5 pixels, respectively. These
299 metrics help evaluate the performance of the motion estimation method at varying levels of accuracy.

300 From the beginning, the initial optical flow map is set to be zero at every pixel. At t -th iteration,
301 the estimated optical flow F_k is used to predict a new sampling position $\mathbf{x}' = \mathbf{x} + F_k(\mathbf{x})$,
302 representing the potential correspondence in the second frame I_2 for a given pixel position \mathbf{x} in
303 the first frame I_1 . Since the estimated optical flow value are typically sub-pixel and the
304 corresponding positions are not aligned with integer pixel coordinates, direct indexing into the
305 correlation volume is infeasible. As a result, the interpolation is required at the sampling
306 positions. Around each sampling position, a local neighbourhood with radius r is extracted from
307 every level of the 4D correlation pyramid using bilinear interpolation, yielding multi-scale
308 matching responses. These responses are concatenated with the present estimated flow features,
309 along with the contextual features extracted from the first frame I_1 via the Dual-Stage Multi-
310 Scale Composite Encoder (DSMSCE) module.

311 The resulting feature combination is subsequently fed into an update block composed of a
312 Convolutional Gated Recurrent Unit (ConvGRU) with shared weights, as detailed in **Fig. S9a**.
313 This module is designed to iteratively refine the optical flow estimation over multiple update
314 steps. By replacing the fully connected layers in a standard GRU with convolutional operations,
315 the ConvGRU preserves the spatial structure of the input feature maps. During each iteration, the
316 ConvGRU updates the hidden state and predicts an incremental flow correction, which is added
317 to the previous estimate to refine the optical flow. In each iteration, the network outputs an
318 incremental update to the flow field based on the current flow estimate, correlation features, and
319 context features, gradually approaching the final precise optical flow field. Specifically, in a
320 single iteration update, the network first accepts the current estimated flow and the correlation

321 features obtained through search. These features are then further processed and concatenated
 322 through convolutions with different dimensions. Subsequently, the concatenated features, along
 323 with context features and the hidden state from the previous iteration, are input into a ConvGRU
 324 module. Taking the t -th iteration as an example, the specific formula of the GRU module is:

$$z_t = \sigma\left(\text{Conv}_{3\times 3}\left([h_{t-1}, x_t], W_z\right)\right) \quad (\text{S9})$$

$$r_t = \sigma\left(\text{Conv}_{3\times 3}\left([h_{t-1}, x_t], W_r\right)\right) \quad (\text{S10})$$

$$\tilde{h}_t = \tanh\left(\text{Conv}_{3\times 3}\left([r_t \odot h_{t-1}, x_t], W_h\right)\right) \quad (\text{S11})$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (\text{S12})$$

325 where x_t represents the concatenated input features, h_{t-1} is the hidden state from the previous time
 326 step, σ denotes the sigmoid activation, \tanh is the hyperbolic tangent function, and \odot represents
 327 element-wise multiplication. The initial hidden state and context features are obtained from the
 328 DSMSCE based on the first frame. The hidden state output by the GRU module is then passed
 329 through two convolutional layers to predict the flow increment ΔF_t for the current iteration. The
 330 new optical flow estimate is then calculated as $F_{t+1} = F_t + \Delta F_t$. Moreover, since the output flow
 331 field resolution is 1/8 of the input image resolution, we employ a learnable convolutional
 332 weighted interpolation method to recover the original resolution. During up-sampling, the
 333 network first predicts a learnable mask of size $8 \times 8 \times 9$, and by applying SoftMax weighting to the
 334 9 neighbouring pixels, the low-resolution flow is up-sampled to full resolution. By repeating this
 335 iterative process, the optical flow estimation is progressively improved.

336 Additionally, we studied the impact of the number of iterations on the optical flow estimation
 337 accuracy. We selected two clean images from the Sintel dataset and generated simulation results
 338 close to metalens imaging using the PSF-based simulation method. As shown in **Fig. S9b**, as the

339 number of iterations increases from 4 to 12, the optical flow estimation results improve, and the
340 details gradually become clearer. The specific metrics in **Table S4** further confirm this trend.
341 However, when the number of iterations exceeds 12, the improvement in optical flow estimation
342 accuracy becomes negligible. Therefore, we set the number of iterations to 12 for metalens
343 optical flow estimation.

344 **Note S10. Structural details of DSMSCE module.**

345 Given an input feature map $X_1 \in H_1 \times W_1 \times C$, the outputs of the parallel pooling branches can be
 346 expressed as:

$$P_m = \text{Upsample}[\text{AvgPool}(X_1; \text{size} = m \times m)], m \in \{1, 2, 3\} \quad (\text{S13})$$

347 The identity mapping branch is denoted as $P_0 = X_1$. By concatenating this branch with the
 348 outputs of the parallel pooling paths, a combined feature representation $Z_1 = \text{Concat}(P_0, P_1, P_2,$
 349 $P_3) \in H_1 \times W_1 \times 4C$ is formed. A convolution operation is then applied to this concatenated feature
 350 map to reduce the dimensionality, yielding in the initial fused representation.

351 The DASPP module consists of a 1×1 standard convolution, multiple 3×3 Atrous convolutions
 352 with different dilation rates, and a global average pooling branch. Atrous convolution is a
 353 method to expand the receptive field by introducing gaps (or dilation) between the sampling
 354 points of a standard convolution operation. Given the intermediate feature input $X_2 \in H_2 \times W_2 \times C'$,
 355 the output of the 3×3 Atrous convolution branch can be expressed as:

$$A_r(i, j) = \sum_{m, n=-1}^1 W_r(m, n) X_2[i + rm, j + rn], r \in \{6, 12, 18\}, W_r \in 3 \times 3 \times C' \quad (\text{S10})$$

356 where $A_r(i, j)$ denotes the value at position (i, j) of the output feature map obtained by applying an
 357 Atrous convolution with dilation rate r . The convolutional kernel weight at position (m, n) for
 358 dilation rate r is denoted as $W_r(m, n)$, where $m, n \in \{-1, 0, 1\}$, covering all positions within the 3×3
 359 kernel. The parameter r defines the dilation rate, which controls the sampling interval, thus
 360 achieving multi-scale receptive fields. Each kernel W_r has a shape of $3 \times 3 \times C'$, indicating that
 361 each output channel is associated with a distinct 3×3 filter. In our implementation, the dilation
 362 rates are set to $r = 6, 12, 18$, to obtain the outputs from the dilated convolution branches, denoted
 363 as A_1, A_2, A_3 , respectively. Let A_0 represent the output of the standard 1×1 convolution branch,

364 and A_4 denote the output of the global average pooling branch followed by a 1×1 convolution.

365 The final output of the DASPP module is obtained by concatenating all branches:

$$Y = \text{Conv}_{1 \times 1}(\text{Concat}(A_0, A_1, A_2, A_3, A_4)), Y \in {}^{H_2 \times W_2 \times C'} \quad (\text{S11})$$

366

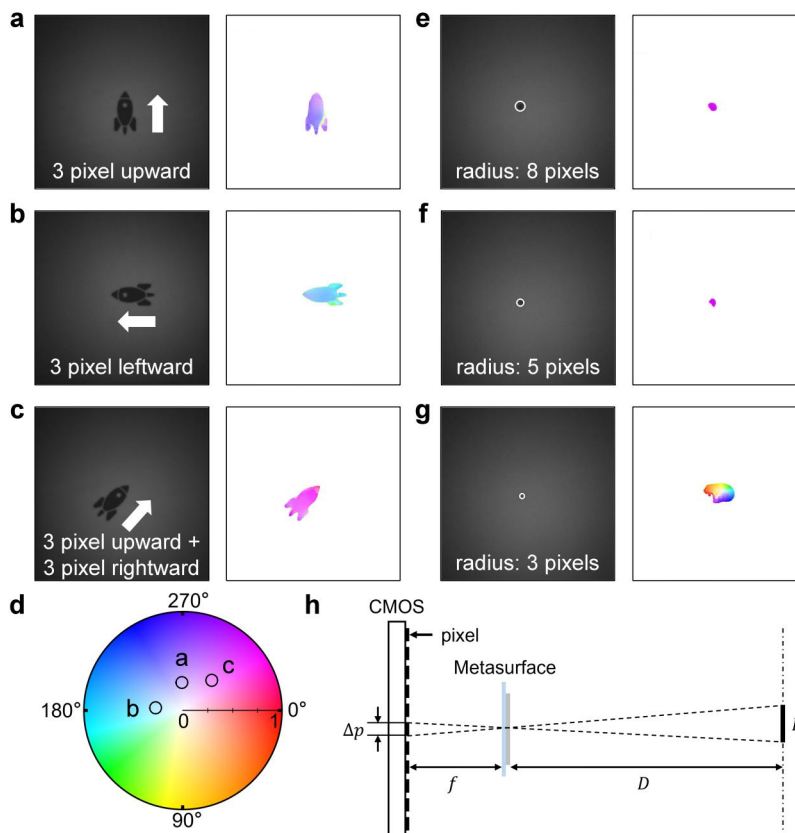
367 **Note S11. Sequence-wise exponentially weighted L₁ loss for optical flow estimation.**

368 During training, our network utilizes a sequence-wise L₁ loss function to supervise the optical
369 flow estimation process. Specifically, the network generates a sequence of N optical flow
370 predictions through iterative refinement, where the prediction at the i -th iteration is denoted as f_i ,
371 and the ground-truth flow is denoted as f_{gt} . The overall loss function is formulated as follows:

$$L = \sum_{i=1}^N \rho^{i-N} \|f_{\text{gt}} - f_i\|_1 \quad (\text{S16})$$

372 where $\|\cdot\|_1$ denotes the pixel-wise L₁ norm, which measures the absolute error between the
373 predicted and ground-truth flow fields for every pixel. The parameter N represents the total
374 number of iterations. The parameter ρ is an exponential decay factor, ranging from 0 to 1 (set to
375 0.8 here). This weighting mechanism assigns greater importance to later iterations closer to the
376 final output, while earlier predictions receive lower weights. By supervising all intermediate
377 predictions with exponentially increasing weights, this loss function accelerates network
378 convergence and improves the quality of flow estimation at each step. Furthermore, this
379 sequence-wise supervision enhances the robustness and generalization ability of the model,
380 ensuring reliable optical flow prediction even with fewer refinement steps.

381 **Note S12. Characterization of the motion detection limits in PINS.**



382
 383 **Figure S10 | Characterization of the extreme motion detection capabilities of the Meta-Motion Sense**
 384 **(MMS) deep neural network for subtle displacements and small targets.** Detected optical flow results for
 385 subtle displacements of the rocket pattern moving towards three different directions: **a**, 3 pixels upward, **b**, 3
 386 pixels leftward, and **c**, sequentially 3 pixels and 3 pixels rightward. **d**, Three movements shown in (**a**) to (**c**) are
 387 marked on the colour wheel. Detected optical flow results for small circular targets with radii of **e**, 8 pixels, **f**, 5
 388 pixels and **g**, 3 pixels, respectively, moving 15 pixels to the right. All the captured images are 701×701 pixels
 389 in size. **h**, Schematic illustration for spatial resolution estimation. Δp is the image size on CMOS sensor, f is
 390 the focal length of metalens, D is the object distance, and P is the corresponding object size.

391 To further investigate the motion detection limits of PINS, we designed test scenarios featuring
 392 minute displacements and tiny size moving targets, and analysed the resulting optical flow
 393 outputs. To assess PINS's sensitivity to minimal motion, we employed a rocket pattern, and

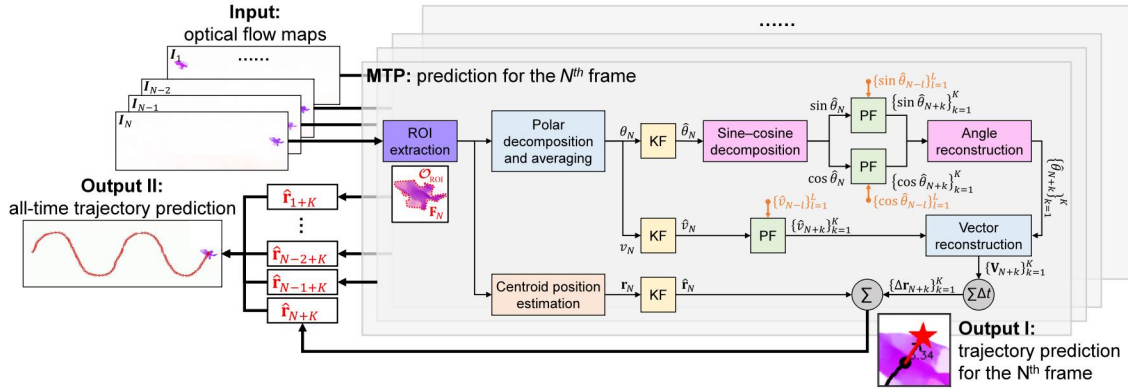
394 tested its displacement in three directions: vertical, horizontal and diagonal. **Figures S10a** and
395 **S10b** show two sequential frames and their corresponding optical flow results when the target is
396 displaced by 3 pixels upward and leftward, respectively. **Figure S10c** illustrates the results for a
397 diagonal displacement, in which the object moves 3 pixels upward and 3 pixels rightward,
398 resulting in an effective shift of approximate 4.2 pixels. These results confirm that PINS is
399 capable of accurately reconstructing the object's shape and reliably capturing both the direction
400 and magnitude of motion, even at the level of a single pixel moving on CMOS. The motion
401 information corresponding to the optical flow results in **Fig. S10a to S10c** is indicated in the
402 colour wheel shown in **Fig. S10d**.

403 Furthermore, to evaluate the detection limit for small-scale moving objects, we translated
404 small targets with radii of 8, 5, and 3 pixels, respectively, by 15 pixels to the right and analysed
405 the corresponding optical flow results. As shown in **Fig. S10e**, for the 8-pixel-radius target, the
406 MMS neural network accurately detected both the motion and its direction, indicating reliable
407 performance. When the radius was reduced to 5 pixels, minor distortions emerged around the
408 target's periphery, however, the motion remained clearly identifiable, as shown in **Fig. S10f**. In
409 contrast, for the smallest target with a radius of 3 pixels, the network failed to robustly capture
410 the motion, and the resulting optical flow exhibited considerable artifacts and noise, as illustrated
411 in **Fig. S10g**. These findings indicate that PINS can effectively detect moving objects with image
412 sizes as small as 5 pixels in diameter on CMOS.

413 The aforementioned displacements and object sizes are quantified in terms of pixel units on
414 the CMOS sensor but can be readily converted into real-world spatial dimensions. This
415 conversion depends on both the pixel size (2.2 μm) of the CMOS sensor and the focal length (1.6
416 mm) of the integrated metalens, as illustrated in **Fig. S10h**. After conversion, each pixel on the

417 CMOS sensor corresponds to an angular resolution of approximately 1.2 mrad ($\sim 0.07^\circ$).
418 Accordingly, the minimum detectable displacement of the PINS system is about 4 cm at a
419 distance of 10 meters. Furthermore, the smallest object size that can be reliably identified by
420 PINS corresponds to approximately 10 cm at the same distance.

421 **Note S13. Detailed formulation of the Motion Trajectory Prediction (MTP) method.**



422
423 **Figure S11 | Detailed processing flowchart of single-object trajectory prediction using MTP framework.**

424 The MTP framework consists of four key stages: region of interest (ROI) extraction, current centroid position
425 estimation, velocity pre-processing and curve fitting, and future centroid position estimation with trajectory
426 extrapolation. KF: Kalman filtering. PF: polynomial fitting.

427 **Figure S11** illustrates the process of predicting the future trajectory of a single object based on
428 the preceding sequence of optical flow maps $\{I_1, I_2, \dots, I_N\}$, with I_N denoting the current frame.
429 The process begins with ROI extraction, where the dominant moving region in each optical flow
430 map is identified by thresholding the flow magnitude followed by connected-component analysis,
431 as illustrated by the red dashed outline in **Fig. S11**. The largest connected region O_{ROI} is selected
432 as the target, effectively suppressing background motion and noise. Next, from the extracted
433 region O_{ROI} , the current object's centroid position r_N is estimated as the spatial average of pixel
434 coordinates:

$$\mathbf{r}_N = (x_N, y_N) = \frac{1}{M} \left(\sum_{i \in O_{ROI}} x_i, \sum_{i \in O_{ROI}} y_i \right) \quad (\text{S17})$$

435 where M denotes the total number of pixels contained within the region O_{ROI} . Subsequently,
436 velocity information is extracted from the optical flow map in the region of interest,

437 preprocessed, and fitted using polynomial regression. Given the optical flow map \mathbf{F}_N within
 438 $\mathbf{0}_{\text{ROI}}$, a polar decomposition and averaging are performed to obtain the velocity magnitude v_N
 439 and direction angle θ_N of the target object:

$$v_N = \frac{1}{M} \sum_{f \in \mathbf{F}_N} \sqrt{f_x^2 + f_y^2} \quad (\text{S18})$$

$$\theta_N = \frac{1}{M} \sum_{f \in \mathbf{F}_N} \text{atan2}(f_y, f_x) \quad (\text{S19})$$

440 where $\text{atan2}(\cdot)$ denotes the arctangent function with two arguments, which takes into account the
 441 sign of both the numerator and denominator to determine the correct quadrant of the angle. To
 442 suppress frame-to-frame noise, the current centroid position \mathbf{r}_N , velocity magnitude v_N and angle
 443 θ_N are then independently smoothed using Kalman filters, resulting in the filtered estimates $\hat{\mathbf{r}}_N$,
 444 \hat{v}_N , and $\hat{\theta}_N$, respectively. The next step involves polynomial fitting of the filtered velocity
 445 magnitude \hat{v}_N and angle $\hat{\theta}_N$. Specifically, the filtered velocity magnitude $\{\hat{v}_{N-l}\}_{l=1}^L$ and angle
 446 $\{\hat{\theta}_{N-l}\}_{l=1}^L$ over a recent window of L , from I_{N-L} to I_{N-1} , are retained and fitted using weighted
 447 polynomial regression to capture the temporal trend of motion. To emphasize recent motion and
 448 reduce the risk of overfitting to outdated motion patterns, a linearly increasing temporal
 449 weighting scheme is adopted, assigning greater weights to frames closer to the current time index.
 450 The sequence of velocity magnitudes $\{\hat{v}_{N-l}\}_{l=1}^L$ is directly modeled as a quadratic function of time
 451 index n :

$$v = a_1 n^2 + b_1 n + c_1, v \in \{\hat{v}_{N-l}\}_{l=1}^L \quad (\text{S20})$$

452 where a_1, b_1, c_1 are the polynomial coefficients estimated from the weighted least squares fitting.
 453 However, when fitting the direction angle $\{\hat{\theta}_{N-l}\}_{l=1}^L$, the periodic nature of angular data must be

454 properly addressed to avoid discontinuities near the 0 or 2π boundary. Accordingly, the velocity
 455 angle is first decomposed into its sine and cosine components:

$$\begin{cases} \sin \theta = a_2 n^2 + b_2 n + c_2 \\ \cos \theta = a_3 n^2 + b_3 n + c_3 \end{cases}, \theta \in \{\hat{\theta}_{N-l}\}_{l=1}^L \quad (\text{S21})$$

456 These components are then individually fitted using polynomial regression and subsequently
 457 recombined via inverse trigonometric reconstruction. This approach ensures angular continuity
 458 and improves the numerical stability of both the fitting and prediction processes for directional
 459 motion. The fitted curves described in **Eq. (S20)** and **Eq. (S21)** are employed to extrapolate the
 460 velocity magnitude $\{\hat{v}_{N+k}\}_{k=1}^K$ and angle $\{\hat{\theta}_{N+k}\}_{k=1}^K$ from frame $N+1$ to frame $N+K$:

$$\begin{cases} v_{N+n+1-L} = a_1 n^2 + b_1 n + c_1 \\ \sin \theta_{N+n+1-L} = a_2 n^2 + b_2 n + c_2, n \in [L, L+K-1] \\ \cos \theta_{N+n+1-L} = a_3 n^2 + b_3 n + c_3 \end{cases} \quad (\text{S22})$$

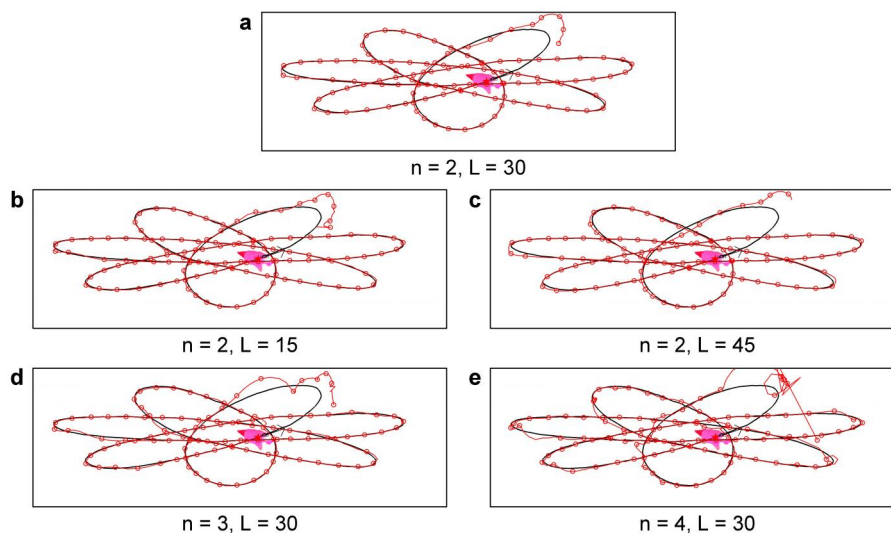
461 From these extrapolated quantities, the future velocity vectors $\{\mathbf{V}_{N+k}\}_{k=1}^K$ can be reconstructed
 462 by converting the polar representation $\{(\hat{v}_{N+k}, \hat{\theta}_{N+k})\}_{k=1}^K$. Furthermore, the predicted object's
 463 position $\hat{\mathbf{r}}_{N+k}$ at frame $N+K$ is obtained by numerically integrating the predicted velocity vectors
 464 over fine-grained temporal substep $\Delta t = T_{\text{frame}}/5$, spanning the interval $[N+1, N+k]$:

$$\hat{\mathbf{r}}_{N+k} = \hat{\mathbf{r}}_N + \Delta \mathbf{r}_{N+k} = \hat{\mathbf{r}}_N + \sum_{m=1}^{5k} \mathbf{V}_{N+m/5} \Delta t \quad (\text{S23})$$

465 Here, although the intermediate substep index $N+m/5$ is not an integer, it can still be directly
 466 substituted into **Eq. (S22)** to evaluate the extrapolated velocity using the fitted polynomial curves.
 467 This substep-based integration scheme mitigates accumulation errors and improves the numerical
 468 stability compared to direct vector summation over full frames. Finally, the predicted positions
 469 and motion trajectory of the target from frame $N+1$ to frame $N+K$ can be obtained, as illustrated

470 in Output I of **Fig. S11**. In the figure, the five-pointed star denotes the predicted position at frame
471 $N+K$; the red solid line represents the predicted trajectory over the interval $N+1$ to $N+K$; the
472 black filled circle indicates the centroid position at the current frame; the black solid line traces
473 the ground-truth trajectory of the centroid; and the black arrow and numeric label indicate the
474 current velocity direction and magnitude (in units of pixels per second). Furthermore, by
475 recording the predicted position at the furthest forecasted frame for each input frame, we can
476 construct an all-time prediction trajectory, which enables direct comparison with the ground-
477 truth motion path, as illustrated in Output II of **Fig. S11**.

478 **Note S14. Analysis of the influence of fitting parameters on trajectory prediction**
479 **performance.**



480
481 **Figure S12 | Influence of the polynomial fitting order and recent window length on trajectory prediction**
482 **performance. a**, Electron trajectory prediction results for 15 frames ahead when the polynomial fitting order is
483 $n = 2$ and the length of the recent window is $L = 30$, consistent with Fig. 5d. **b-c**, Electron trajectory prediction
484 results with fixed polynomial order $n = 2$ and window length set to $L = 15$ and $L = 45$. **d-e**, Electron trajectory
485 prediction results with fixed window length $L = 30$ and polynomial order set to $n = 3$ and $n = 4$.

486 To evaluate the sensitivity of the trajectory prediction performance to the fitting parameters, we
487 conducted a series of controlled experiments by varying the polynomial fitting order n and the
488 length of the recent window L . **Figure S12** summarizes the prediction results under different
489 parameter configurations. Among all tested configurations, the setting shown in **Fig. S12a** with
490 polynomial order $n = 2$ and recent window length $L = 30$ yields the best overall performance. It
491 achieves a good balance between smoothness and responsiveness, and is consistent with the
492 prediction result presented in **Fig. 5d** of the main text.

493 To investigate the influence of the recent window length L , we fixed the polynomial order at n
494 $= 2$ and compared prediction results under different L values. In contrast, **Fig. S12b** shows the
495 result when $L = 15$ with $n = 2$ fixed. Due to the limited number of historical data available for
496 fitting, the prediction becomes more susceptible to local noise, leading to noticeable fluctuations
497 and irregular jitter in the predicted trajectory. This indicates underfitting caused by insufficient
498 context for the polynomial model. When the window length is extended to $L = 45$ (**Fig. S12c**),
499 the prediction starts to exhibit overfitting. Although more data points are included in the fitting
500 process, the model becomes less responsive to recent directional changes, resulting in larger
501 deviations at trajectory turning points. This illustrates the trade-off between temporal resolution
502 and stability.

503 Furthermore, we analysed the impact of increasing the polynomial fitting order while keeping
504 $L = 30$ fixed. As shown in **Fig. S12d and S12e**, when the order is raised to $n = 3$ and $n = 4$, the
505 predicted trajectories become increasingly erratic and less physically plausible. This degradation
506 is likely due to Runge's phenomenon, where higher-order polynomials tend to exhibit large
507 oscillations, especially near the boundaries of the fitting interval. In the presence of noise or
508 subtle fluctuations in the input data, high-order models tend to amplify such variations, leading
509 to unstable and unreliable predictions.

510 **Note S15. Quantitative RMSE evaluation of trajectory prediction accuracy.**

511 **Table S5** summarizes the quantitative RMSE evaluation of the trajectory prediction accuracy for

512 the single-object cases in **Fig. 5** and the multi-object case in **Fig. 6**. For the single-object cases in

513 **Fig. 5**, the reference position is defined as the centroid position extracted at the corresponding

514 future frame. For the multi-object case in **Fig. 6**, the RMSE values are reported separately for

515 each tracked object ID. These results demonstrate that the proposed MTP framework achieves

516 accurate short-term trajectory prediction for both single-object and multi-object cases.

517 Meanwhile, the prediction error increases for more complex trajectories and longer forecasting

518 horizons, as expected, due to accumulated extrapolation uncertainty.

519 **Table S5 | Quantitative RMSE evaluation of the trajectory prediction accuracy in Fig. 5 and Fig. 6.**

	RMSE _X (pixels)	RMSE _Y (pixels)	RMSE _{2D} (pixels)
1 frame ahead in Fig. 5b	1.51	1.09	1.86
15 frames ahead in Fig. 5b	20.54	14.46	25.12
30 frames ahead in Fig. 5b	41.82	28.08	50.37
1 frame ahead in Fig. 5c	5.02	3.16	5.93
15 frames ahead in Fig. 5c	71.68	45.36	84.83
30 frames ahead in Fig. 5c	137.68	94.35	166.90
1 frame ahead in Fig. 5d	10.69	3.17	11.15
15 frames ahead in Fig. 5d	153.47	46.07	160.23
30 frames ahead in Fig. 5d	300.52	88.35	313.24
ID 0 in Fig. 6	10.96	3.83	11.61
ID 1 in Fig. 6	4.02	1.39	4.25
ID 2 in Fig. 6	9.91	2.46	10.21
ID 3 in Fig. 6	38.32	5.29	38.68
ID 4 in Fig. 6	7.47	2.29	7.81

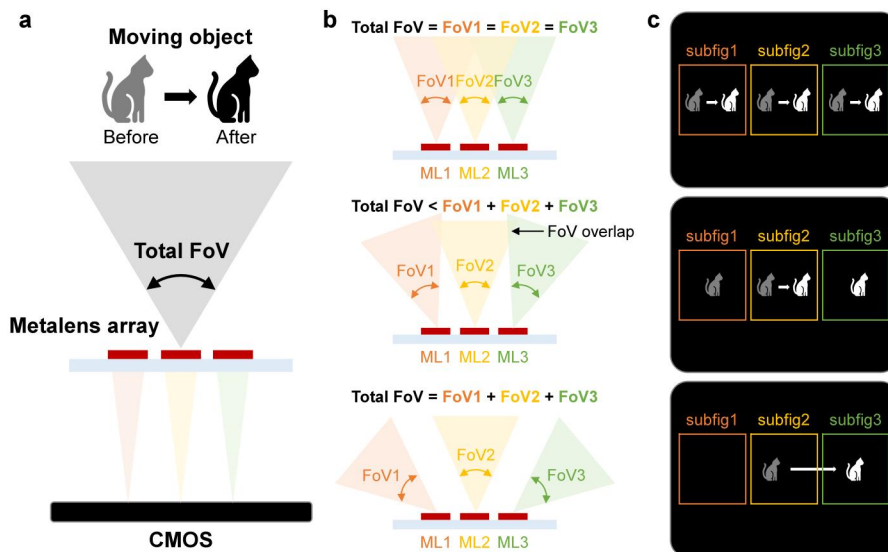
520 Note: “RMSE_X” denotes the root mean square errors of the predicted positions along the x direction, “RMSE_Y”

521 denotes the root mean square errors of the predicted positions along the y direction, “RMSE_{2D}” denotes the

522 two-dimensional root mean square error, calculated from the Euclidean distance between the predicted future

523 position and the corresponding reference position at the target frame.

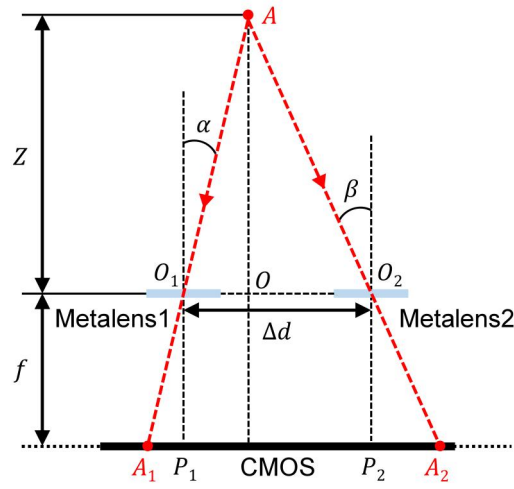
524 **Note S16. Flexible trade-off between motion-information redundancy and spatial coverage**
 525 **by tilt-phase design.**



526
 527 **Figure S13 | View-angle engineering of the metalens array by tilt-phase design.** **a**, Schematic illustration
 528 of a dynamic scene observed by the metalens array. **b**, Different view-angle allocation schemes realized by
 529 assigning different linear tilt phases to the three metalenses. **c**, Corresponding imaging results under the
 530 different phase-design schemes.

531 **Figure S13** illustrates the view-angle engineering enabled by adding different tilt phases to the
 532 metalenses in the array. By adjusting the tilt phase assigned to each metalens, the angular sectors
 533 corresponding to different sub-images can be redistributed, leading to different imaging modes.
 534 When the view angles of adjacent metalenses are more concentrated and partially overlap, the
 535 system captures the same moving target in multiple sub-images, thereby providing richer
 536 motion-related information but a smaller overall field of view. In contrast, when the view angles
 537 are more separated and the overlap between adjacent sub-images is reduced, the total field of
 538 view is enlarged and more spatial information can be covered. Therefore, this design provides a
 539 flexible trade-off between motion-information redundancy and wide-angle spatial coverage.

540 **Note S17. Extension to stereo-vision depth measurement based on metalenses.**



541

542 **Figure S14 | Principle of stereo-vision depth measurement based on metalenses.**

543 A single-lens imaging system can only provide two-dimensional (2D) projection information of a
 544 three-dimensional (3D) scene. Since depth information is inherently lost in the projection process,
 545 it is difficult to estimate the absolute distance of objects using a monocular image alone without
 546 additional assumptions or constraints. In contrast, a stereo-vision system mimics human
 547 binocular perception by using two spatially separated metalenses to simultaneously capture the
 548 same scene from different viewpoints, enabling the extraction of depth information from a 3D
 549 scene.

550 **Figure S14** shows a stereo-vision system based on two metalenses. The two optical centres O_1
 551 and O_2 of metalens1 and metalens2 are aligned along a common baseline and spaced by a
 552 distance $\Delta d = |O_1O_2|$. The imaging point A_1 is formed by metalens1 from the object point A with
 553 an incident angle α . The position of A_1 can be calculated from the projection function $A_1P_1 =$
 554 $P_f(\alpha)$. The projection function is dependent on the phase design of the metalens—for instance,
 555 $P_f(\alpha) = -f \tan \alpha$ for a hyperbolic phase profile, and $P_f(\alpha) = -f \sin \alpha$ for a quadratic phase profile.

556 Accordingly, the incident angle α can be reconstructed by inverting the projection function,
 557 yielding:

$$\alpha = P_f^{-1} [|A_1 P_1|] = \begin{cases} \arctan \left[-\frac{|A_1 P_1|}{f} \right], & \text{for hyperbolic phase profile} \\ \arcsin \left[-\frac{|A_1 P_1|}{f} \right], & \text{for quadratic phase profile} \end{cases}. \quad (\text{S24})$$

558 Here, $|A_1 P_1|$ denotes the image-plane distance between the optical axis and the projection point
 559 A_1 , and f is the focal length of the metalens. Subsequently, from the triangle $\triangle O_1 A O$, it follows
 560 that the horizontal distance from O_1 to O is given by $|O_1 O| = Z \tan \alpha$, where Z denotes the depth of
 561 point A relative to the baseline. Similarly, for metalens2, the corresponding horizontal distance
 562 can be expressed as $|O_2 O| = Z \tan \beta$, where β is the incident angle of A at metalens2. Finally, based
 563 on the geometric relationship $\Delta d = |O_1 O_2| = |O_1 O| + |O_2 O|$ the object depth Z of point A can be
 564 computed as:

$$Z = \frac{\Delta d}{\tan(\alpha) + \tan(\beta)} = \frac{\Delta d}{\tan\left(P_f^{-1} [|A_1 P_1|]\right) + \tan\left(P_f^{-1} [|A_2 P_2|]\right)}. \quad (\text{S25})$$

565

566 **Movie S1. Single-object trajectory prediction for representative motion patterns.**

567 Movie S1 provides the complete results of the single-object trajectory prediction for three
568 representative motion patterns: sinusoidal, planar spiral, and random trajectories.

569

570 **Movie S2. Multi-object tracking and trajectory forecasting in complex wide-angle scenes.**

571 Movie S2 presents the complete results of multi-object tracking and trajectory prediction in
572 complex wide-angle scenes. The system robustly segments and identifies multiple independently
573 moving targets, assigns consistent object IDs across frames, and accurately predicts their future
574 trajectories.