

Article

Holoscopic Elemental-Image-Based Disparity Estimation Using Multi-Scale, Multi-Window Semi-Global Block Matching

Bodor Almatrouk , Hongying Meng * and Mohammad Rafiq Swash 

Department of Electronic and Electrical Engineering, Brunel University London, London UB8 3PH, UK; bodor.almatrouk@brunel.ac.uk (B.A.); rafiq.swash@brunel.ac.uk (M.R.S.)

* Correspondence: hongying.meng@brunel.ac.uk

Abstract: In Holoscopic imaging, a single aperture is used to acquire full-colour spatial images like a fly's eye by gently altering angles between nearby lenses with a micro-lens array. Due to its simple data collection and visualisation methods, which provide robust and scalable spatial information, and its motion parallax, binocular disparity, and convergence, this technique may be able to overcome traditional 2D imaging issues like depth, scalability, and multi-perspective problems. A novel disparity-map-generating method uses angular information from a single Holoscopic image's micro-images, or Elemental Images (EIs), to create a scene's disparity map. Not much research has used EIs instead of Viewpoint Images (VPIs) for disparity estimation. This study investigates whether angular perspective data may replace spatial orthographic data. Using noise reduction and contrast enhancement, EIs with a low resolution and lack of texture are pre-processed to calculate the disparity. The Semi-Global Block Matching (SGBM) technique is used to calculate the disparity between EI pixels. A multi-resolution approach overcomes EIs' resolution constraints, and a content-aware analysis dynamically modifies the SGBM window size settings to generate disparities across different texture and complexity levels. A background mask and nearby EIs with accurate backgrounds detect and rectify EIs with erroneous backgrounds. Our method generates disparity maps that outperform two state-of-the-art deep learning algorithms and VPIs in real images.

Keywords: holoscopic; elemental images; viewpoint images; micro-lenses; disparity; SGBM



Citation: Almatrouk, B.; Meng, H.; Swash, M.R. Holoscopic Elemental-Image-Based Disparity Estimation Using Multi-Scale, Multi-Window Semi-Global Block Matching. *Appl. Sci.* **2024**, *14*, 3335. <https://doi.org/10.3390/app14083335>

Academic Editor: Sungho Kim

Received: 13 March 2024

Revised: 6 April 2024

Accepted: 8 April 2024

Published: 15 April 2024



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Depth estimation from Holoscopic images is a promising technique that has gained attention recently due to its advantage of calculating depth using a single-aperture camera. Holoscopic cameras are based on the same fundamental principles as conventional cameras but with an additional array of micro-lenses (MLA) in front of the image sensor. In traditional cameras, the main lens translates the object plane into the camera's image plane. The micro-lenses focus light beams from various directions onto a single pixel, thereby capturing the scene in three dimensions.

The pixels behind each micro-lens record the same data as traditional cameras but with a greater precision by measuring information from different angles, as shown in Figure 1. The images formed behind each micro-lens, known as the Elemental Images (EIs), represent unique angles of light incidence. Thus, by analysing the EIs, the location and orientation of each light beam can be determined on a pixel-by-pixel basis. A sub-aperture image of a scene, or a Viewpoint Image (VPI), is created by re-sampling pixels from the same locations across the EIs. The EIs provide angular information, whereas VPIs provide spatial information.

Traditionally, disparity estimation is performed on VPIs, which encompass the entire scene from a certain perspective, whereas EIs only include a portion of it. VPIs share visual characteristics with 2D orthographic stereo images, allowing existing stereo-image-based disparity estimation methods to be applied with few adjustments. Additionally, VPIs can be up-sampled using super-resolution methods [1].

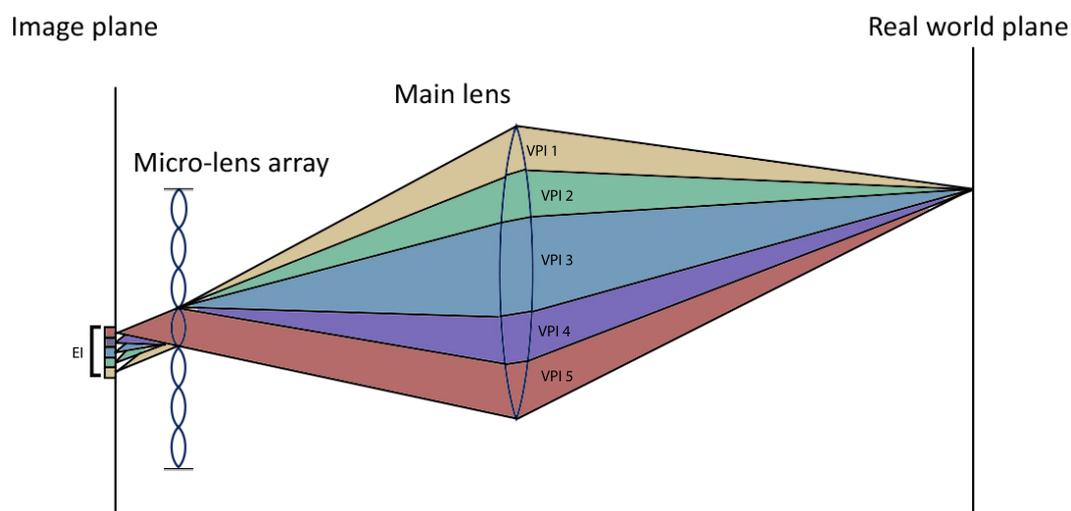


Figure 1. Light beams from various perspectives (VPIs) hitting the same EI in a Holographic sensor.

However, extracting VPIs requires mapping the information gathered from the sensor to reconstruct the scene, which is not always straightforward and can sometimes lead to strong aliasing artefacts [2–4]. Lens error correction and camera calibration must firstly be performed to avoid artefacts. The geometry of the scene must also be considered during VPI creation to avoid image artefacts in areas not ‘in focus’ [5]. Additionally, some micro-lens array designs feature multiple micro-lens sizes with different focal lengths, making it impractical to extract pixels from the same location across all EIs. The convergence of light rays from multiple VPIs might result in overlapping on the image sensor, complicating the separation and extraction of individual rays. Therefore, selecting a ‘patch’ of pixels from each EI might be more effective in increasing the resolution and reducing the artefacts. However, it is more challenging than choosing a single pixel as these patches depend on the depth level within the scene; thus, using the same patch size throughout the entire scene could result in a distorted VPI. For instance, the ideal patch size for displaying the foreground can be excessively large for the background, leading to the occurrence of artefacts in the background.

Extracting VPIs from Holographic images is time-consuming and requires significant storage due to the large number of VPIs generated. Estimating depth from video frames or in real-time scenarios is particularly challenging due to the large number of frames [6]. For these reasons, EIs provide a more straightforward method for estimating disparities, requiring only lens correction as a pre-processing step. In this paper, disparity estimation using perspective EIs is employed, contrary to conventional methods that use extracted, corrected, and up-sampled VPIs.

Perspective projection and orthographic projection are two types of 3D projection. As seen in Figure 2, perspective projection is comparable to the human visual system, in which parallel lines in an image appear to converge at a single point; the closer the object is to the point of convergence, the smaller it appears (change in scale). Orthographic (orthogonal to the scene) projection assumes parallel lines will continue to be parallel and disregards the scaling impact.

Understanding depth via perspective projection is far more precise than using an orthographic technique because, in perspective depth, every light ray is tracked to the precise pixel of its source, unlike in orthographic depth, where light is considered to be emanating from infinity [7]. Although perspective projection has shown more accurate disparity estimation results [8–10], most depth estimation algorithms are performed on orthographic images due to the simplicity of the capturing mechanisms.

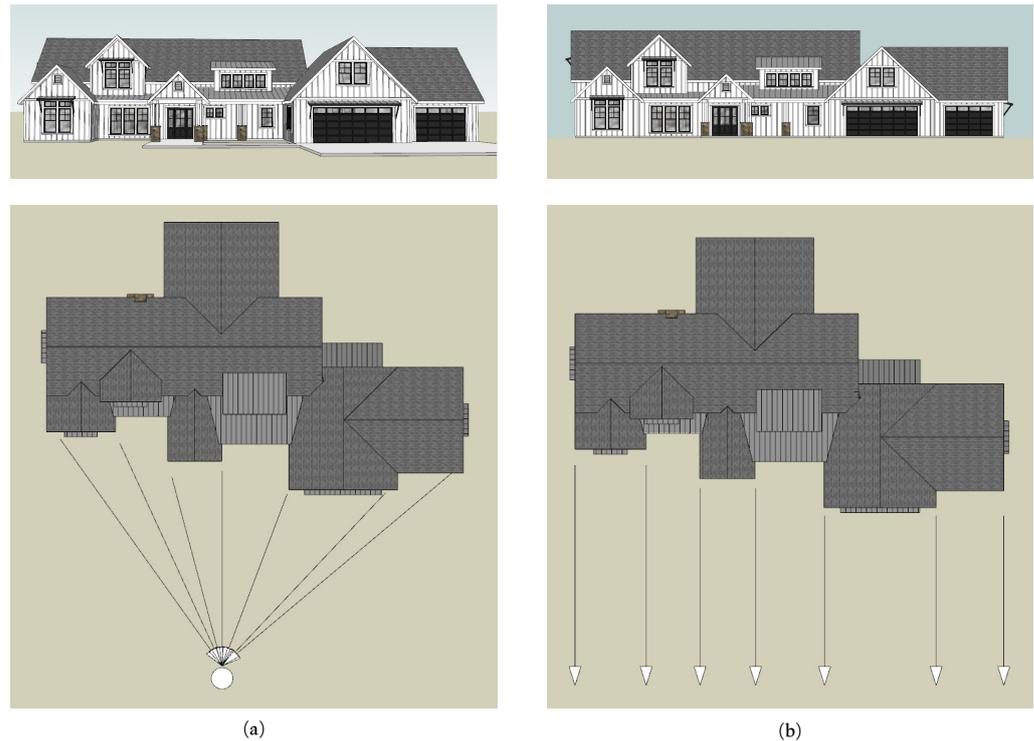


Figure 2. (a) Perspective projection. (b) Orthographic projection.

As seen in Figure 1, the EIs in the Holographic setup record light from different angles, resulting in perspective images that contain angular information. Conversely, VPIs are obtained from various locations on the primary lens, replicating different viewpoints. These images typically exhibit orthographic projection, predominantly capturing surface characteristics. The differentiation here between EIs and VPIs is linked to their ways of spatial representation [11].

2. Methodology

A single Holographic image records the scene's spatial and angular details. Hence, it is possible to compute the scene's depth map from a single shot. In our proposed method, as seen in Figure 3, pre-processing is carried out on the EIs, which is crucial before computing the disparity to improve their quality, as they inherently have a low resolution and lack texture. This procedure consists of two primary stages: noise reduction by bilateral filtering and contrast enhancement via histogram equalisation.

The disparity among EI pixels is computed via the Semi-Global Block Matching (SGBM) algorithm [12], which is favoured due to its flexibility to adapt to the unique features of EIs and its optimal balance between precision and computing efficiency. The SGBM algorithm is enhanced through a multi-resolution approach to address the limitations of EIs in terms of resolution. This involves creating an upscaled pyramid of EIs to capture details at different scales and performing a content-aware analysis to adaptively adjust the SGBM window size parameters. This ensures optimal estimation of disparities across various texture and complexity levels within the EIs. Ultimately, a weighted least squares (WLS) filter is employed to further enhance the optimisation process.

EIs are known to be low in resolution, lack texture, and only capture a portion of the scene. Several deep learning models have been designed to estimate disparity maps, including many specifically designed for VPIs. Deep learning necessitates a substantial and comprehensive dataset specifically designed for Holographic imagery. Pre-existing, pre-trained deep learning stereo-matching solutions would not be compatible with EIs due to differences in training data properties. These solutions are mostly learned using high-quality images, while EIs have a low resolution and lack texture. Deep learning

algorithms have the potential to be highly effective in stereo vision tasks, but their effectiveness is contingent upon the quality and range of the training data. If the training data lack sufficient representation of scenarios, including low resolutions, limited textures, and narrow disparity ranges, the model may exhibit poor generalisation in these settings. Deep learning algorithms may encounter difficulties in generating intricate details in such situations, resulting in unclear outcomes.

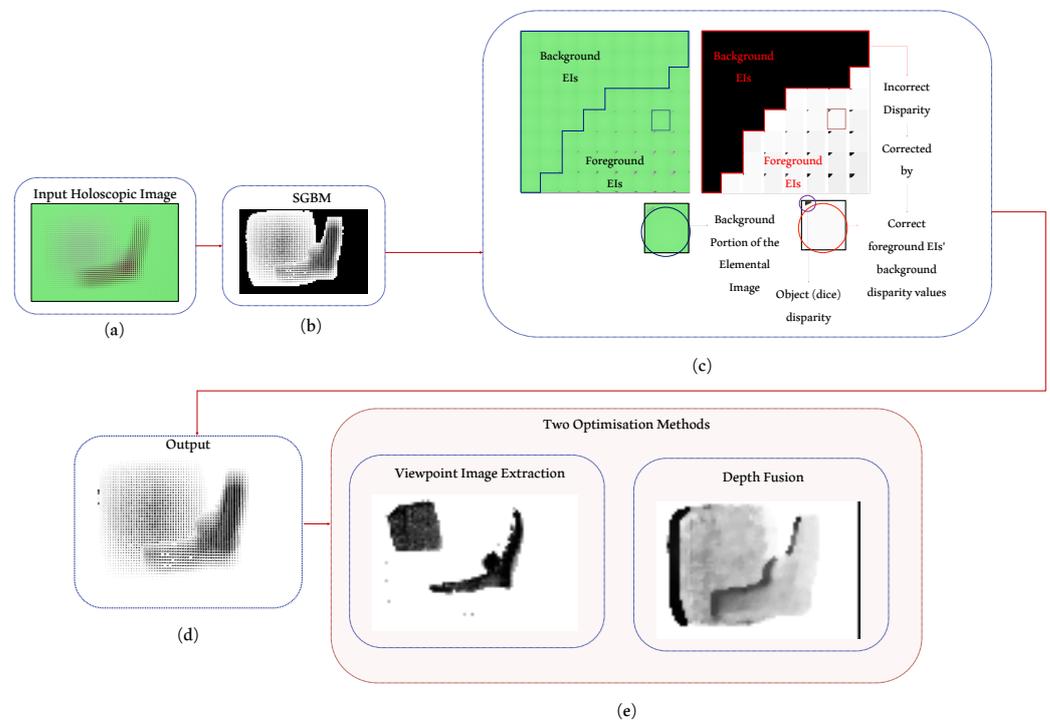


Figure 3. The disparity estimation from the EI pipeline. (a) Input raw Holoscopic image. (b) Disparity map using multi-resolution content-aware SGBM matching. (c) Background correction using background/foreground masks. (d) Output disparity image after SGBM and background correction. (e) Output: two optimisation results. Top: extracted central VPI. Bottom: fusing depth from multiple EIs.

2.1. Pre-Processing of Elemental Images

Before initiating a disparity estimate on the EIs, it is crucial to carry out pre-processing on the EIs to adequately prepare them to achieve an improved outcome. EIs exhibit a low resolution and limited texture. Therefore, while implementing pre-processing techniques, it is crucial to eliminate noise while preserving the critical features.

2.1.1. Noise Reduction through Bilateral Filtering

Applying image blurring is a conventional technique for diminishing image noise. Particularly with images that have minimal texture, such as the background, there may be instances where a stepping effect occurs. This effect is caused by discontinuous disparity levels, resulting in noticeable “steps” in areas with reduced changes in depth, as seen in the textured map in Figure 4. The limited resolution, subtle variations in lighting, limited bit depth, and lack of texture can cause seamless transitions to look like sudden shifts. However, using image blurring will inevitably cause a loss of fine information such as edges, hence reducing the accuracy of disparity estimation. To address this issue, bilateral filtering [13] is employed. This technique, known for its ability to preserve edges, is considered as an advanced way of blurring.

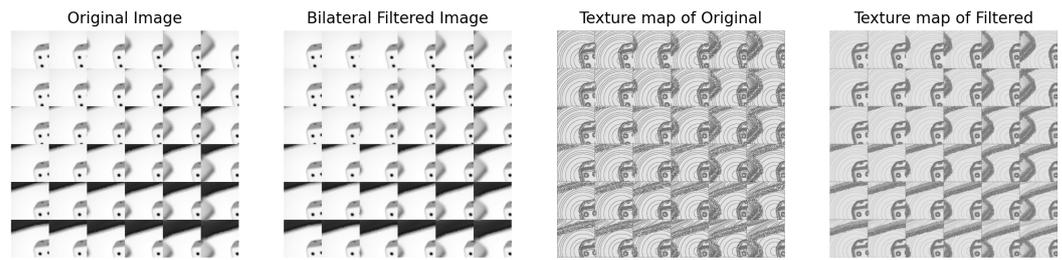


Figure 4. EI before and after applying bilateral filtering. As seen in the texture map of the original image, there is a noticeable stepping effect in the background. Although the filter did not eliminate it, it did assist in reducing the impact while maintaining edge information.

Bilateral blurring is applied to each EI to reduce the noise:

$$I_{\text{filtered}}(p) = \frac{1}{W_p} \sum_{q \in S} I(q) \cdot f_r(\|I(p) - I(q)\|) \cdot f_s(\|p - q\|) \quad (1)$$

Let $I_{\text{filtered}}(p)$ represent the filtered intensity of pixel p , $I(q)$ denote the intensity of the next pixel, S be the set of pixels surrounding p , and W be the normalised factor. The variable f_r represents the spatial range of the kernel, which corresponds to the dimensions of the neighbouring region. On the other hand, f_s denotes the minimum magnitude required for an edge to be detected. This procedure ensures that only pixels with similar intensity levels to the core pixel are considered for blurring while maintaining distinct intensity fluctuations. A lower value of f_r leads to a more distinct edge. As the value of f_s tends towards infinity, the equation approaches convergence to a Gaussian blur.

2.1.2. Contrast Enhancement via Histogram Equalisation

Due to their low resolution and the settings under which they are captured (tiny micro-lenses), EIs often exhibit a lack of contrast. Histogram equalisation is commonly employed to enhance the image contrast by spreading the intensity levels, hence boosting feature visibility by:

$$I_{\text{equalised}}(p) = CDF(I(p)) \quad (2)$$

where $I(p)$ is the intensity of pixel p in the original image, $I_{\text{equalised}}(p)$ is the intensity of pixel p in the equalised image, and CDF is the cumulative distribution function of the original image's histogram. The function $CDF(I(p))$ assigns a new intensity value to each pixel by utilising the cumulative distribution, hence improving the contrast of the image.

2.2. Content-Aware Multi-Resolution Disparity Estimation Using Semi-Global Block Matching

2.2.1. Overview

Deriving disparity from EIs using SGBM presents challenges, mostly attributed to the lack of texture and low resolution. Upsampling the EIs would lead to data loss, resulting in the introduction of noise and a decrease in image quality. Many studies [14,15] have investigated the computation of disparity at various resolutions to enhance the accuracy of disparity maps, particularly in the context of developing deep learning models. While rescaling EIs may result in a loss of image quality, calculating the disparity at multiple resolutions instead of merely one upsampled resolution still is an effective approach for handling varying levels of details and textures. Lower resolutions may result in the loss of some details in the scene, while higher resolutions may exhibit an inconsistent overall structure. This indicates a trade-off between maintaining structural consistency and capturing high-frequency details based on the input resolution [15].

Content information can vary across different regions within EIs, particularly at varying resolutions. By utilising a content-aware approach, the disparity window size can be dynamically modified according to the characteristics of each location, resulting in more accurate disparity estimations. When working with areas that have a high level of

texture, using a smaller window size would be advantageous in capturing intricate details. Conversely, in areas that lack texture, a larger window size can be used to minimise noise.

2.2.2. Multi-Resolution Elemental Image Pyramid

Typically, when constructing a pyramid with multiple resolutions for any objective, the procedure commences by taking the original image and reducing its size. However, when it comes to EIs, the images are already of low resolution. Creating a pyramid by progressively down-sampling them will result in extremely small images that lack significant information. Thus, in the instance of the EIs, the pyramid is formed by enlarging the EIs into two additional layers and down-sampling the image by one layer, as seen in Figure 5, enabling the algorithm to encompass characteristics that span from large-scale structures at lower levels to intricate details at higher resolutions. Start with the EI of the original size as the base level L_0 , with each level increased by a factor of two using bicubic interpolation [16]. A minimum resolution threshold is implemented to prevent further down-sampling of EIs with an extremely low resolution. If the value of the EI is less than 40×40 , the down-sampling step is omitted.

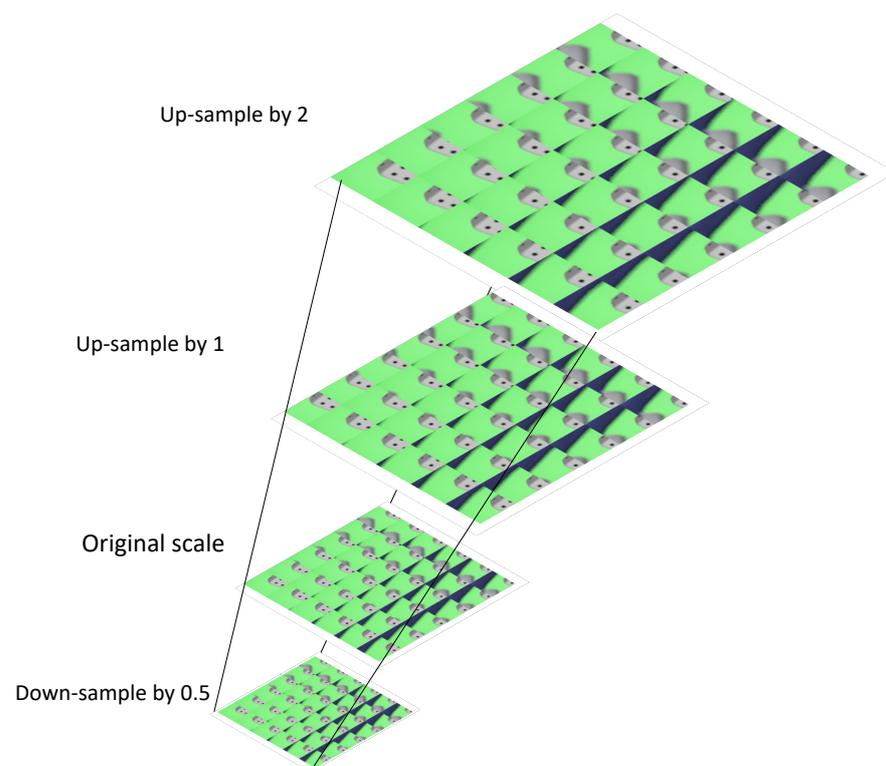


Figure 5. Multi-resolution pyramid of EIs.

2.2.3. Multi-Resolution Content Analysis

Content-aware analysis is an essential process for evaluating the visual attributes in the EIs. Its purpose is to optimise the window size parameters used in disparity estimation based on the complexity and textures present at different scales. This analysis is particularly valuable for adjusting window size parameters at both single and multiple scales to enhance the precision and resilience of the disparity.

EIs possess a high degree of sensitivity. Consequently, a simpler approach involving edge segmentation and texture analysis is employed. The Sobel filter is utilised for accomplishing edge detection. The filter's sensitivity is contingent upon the resolution of the images. Low-resolution images necessitate a higher threshold for detecting significant structures, while high-resolution images require a lower threshold to identify finer details as depicted in Figure 6.

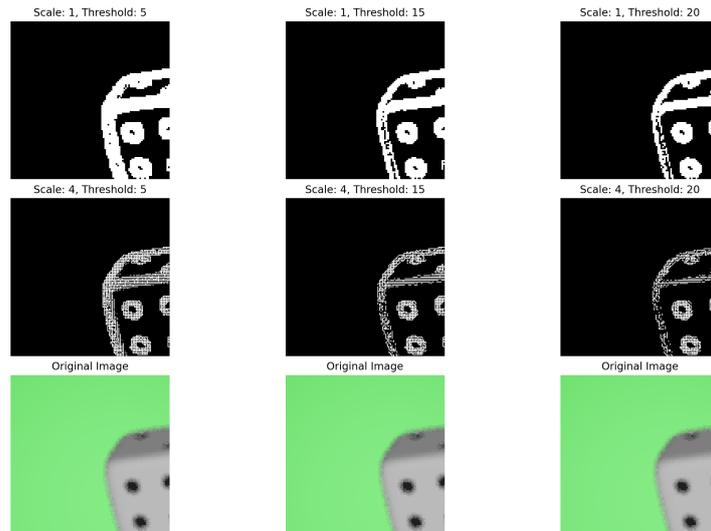


Figure 6. Examples of extreme edge thresholds show that the sensitivity of the filter depends on the resolution of the image. High-resolution images need a lower threshold to identify finer structures, while low-resolution images need a higher threshold to identify significant features.

Textures are ideal to identify intensity patterns which is great for identifying regions for disparity estimation. Local Binary Patterns (LBPs) are used in this case to identify the textures in the EIs. Here, the focus is on larger patterns at lower resolutions and finer textural details at higher resolutions.

$$LBP_n(p) = \sum_{k=0}^{P-1} 2^k \cdot \mathbf{1}(I_n(p_k) \geq I_n(p)) \tag{3}$$

where LBP is computed for pixel p located at location n within the image used to classify the texture. P represents the total number of pixels neighbouring to p , with the summation ranging from $k = 0$ to $P - 1$ and 2^k represents the weighting factor assigned to each neighbouring element, which is determined by its location. The neighbouring pixel's (p_k) intensity is compared with the central pixel $I_n(p)$. $\mathbf{1}(I_n(p_k) \geq I_n(p))$ returns 1 if it is true and 0 if false. Texture maps across different scales are shown in Figure 7.

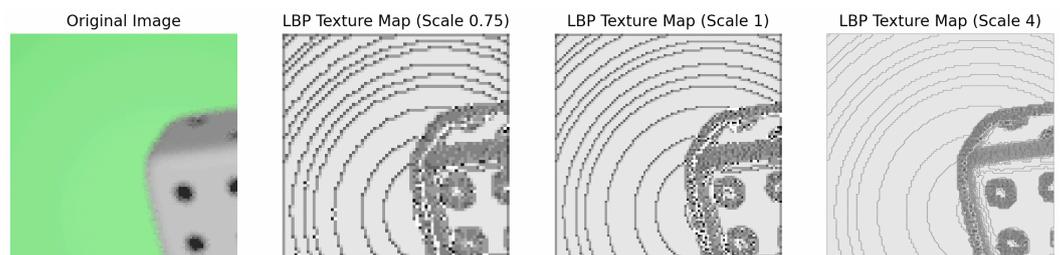


Figure 7. LBP texture maps across different scales before filtering to show the effect.

To enhance simplicity and preserve time, the edge map E and texture map T are combined to then form a dynamic adaptive window for the computation of disparities.

$$F(x, y) = \alpha \cdot E(x, y) + (1 - \alpha) \cdot T(x, y) \tag{4}$$

The combined feature at pixel (x, y) is denoted as F , and it is influenced by a weighted factor, α , which ranges from 0 to 1. This factor determines the appropriate ratio between edge and texture data. The value of α has been cautiously adjusted to achieve the ideal result for each image.

2.2.4. Multi-Resolution Multi-Window Disparity Estimation Using SGBM

Dynamic Window: Given that the EI’s resolution varies from 50×50 to around 400×400 , it is necessary to select a range of window sizes. The value of W_{\min} is selected to be around 5% of the minimum resolution, resulting in an amount of 5. Similarly, the value of W_{\max} is chosen to be roughly 20% of the resolution, resulting in a value of 80. The size of the window adjusts according to the value of the feature map F . The dynamic window size W at pixel (x, y) can be calculated by:

$$W(i, j) = W_{\min} + (W_{\max} - W_{\min}) \cdot (1 - F_{\text{norm}}(i, j)) \tag{5}$$

Greater values of F_{norm} , representing the normalised F values, will result in the selection of smaller windows for more complicated regions, and vice versa.

Semi-Global Block Matching Disparity: The disparity is calculated by comparing blocks of pixels along the epipolar line and obtaining the associated vertical displacement, as demonstrated in our previous work [17] (see Figure 8). This problem can be represented by a comprehensive cost function:

$$E(D) = \sum_{d \in D} \left(C(d) + \sum_{d' \in N(d)} P_1 I_{\{|d-d'|=1\}} + \sum_{d'' \in N(d)} P_2 I_{\{|d-d''|>1\}} \right) \tag{6}$$

where I is a function that indicates whether an input is true or false and returns 1 or 0 accordingly. (d) is the chosen disparity’s data term similarity metric. A 3D cost structure is used to hold each similarity cost, and this process is repeated for each pixel block, with a cost of d . The 3D structure stack’s minimal costs stand for possible disparity estimates.

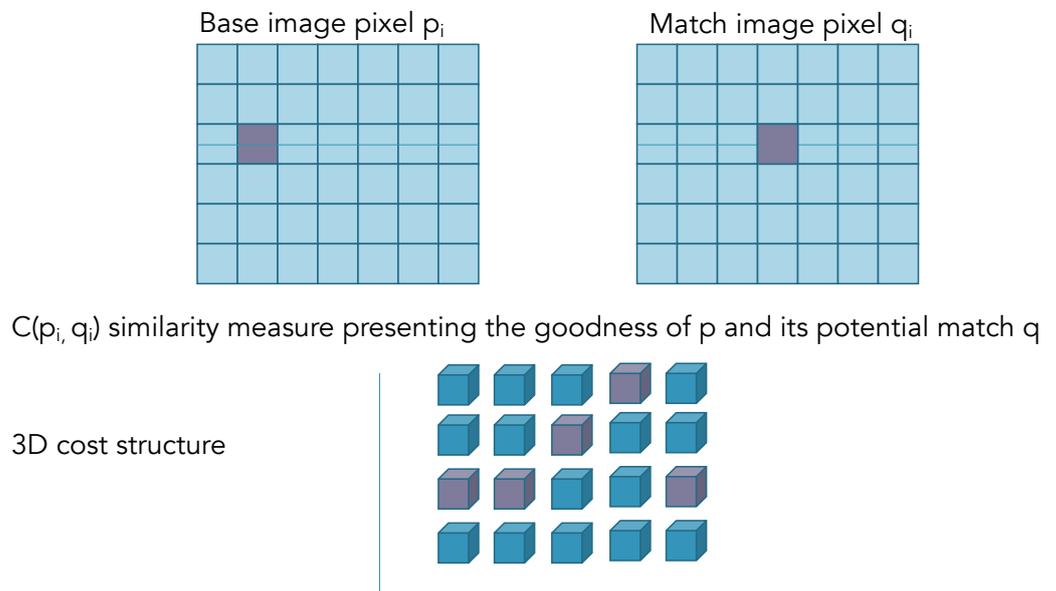


Figure 8. The produced minimal costs are not highly distinctive, which could result in incorrect disparity estimation [17].

Disparity Aggregation: The resulting minimum costs may lack significant distinctiveness, thus resulting in an incorrect assessment of disparity. This issue is addressed by employing cost aggregation within these 3D cost structures. The total cost is determined by aggregating the lowest costs across various image paths. A total of eight paths were utilised in this paper. Potential cost values were pooled, and a weighted summing of these cost possibilities was conducted. The weights were obtained from the normalised feature map $F_{\text{norm}}(x, y)$, which characterises the contents (texture and edges) at each scale level. The feature map underwent normalisation:

$$F_{\text{norm}}(x, y) = \frac{F(x, y)}{\sum_L F_L(x, y)} \tag{7}$$

The function F_{norm} represents the normalised feature for pixels (x, y) , whereas L is the scale level. Normalising the feature map guarantees that, throughout the content-aware analysis, disparities from all resolutions contribute proportionally. Thus, the final content-aware disparity map D_{final} can be represented as:

$$D_{\text{final}}(x, y) = \sum_i \left(\frac{F_i(x, y)}{\sum_j F_j(x, y)} \right) \cdot D_L(x, y) \tag{8}$$

where D_L is the disparity optimised at each level of resolution. To achieve greater accuracy, a higher weight is assigned to the original scale, since this method is still sensitive to multiple scales:

$$D_{\text{final}}(x, y) = \left(\frac{\alpha \cdot F_{L_0}(x, y)}{\sum_j F_j(x, y) + (\alpha - 1) \cdot F_{L_0}(x, y)} \right) \cdot D_{L_0}(x, y) + \sum_{i \neq L_0} \left(\frac{F_i(x, y)}{\sum_j F_j(x, y) + (\alpha - 1) \cdot F_{L_0}(x, y)} \right) \cdot D_L(x, y) \tag{9}$$

The expression $\sum_j F_j(x, y) + (\alpha - 1) \cdot F_{L_0}(x, y)$ ensures normalisation for assigning a larger weight to F_{L_0} , where F_{L_0} represents the feature map at the original level and α is the weighting factor.

Penalty terms P_1 and P_2 are introduced, which are based on the difference in neighbourhood disparities, where $N(d)$ is the neighbour of d . Accordingly, for each pixel, all its neighbouring pixels along the routes are analysed; the greater the difference between the lateral parallax axes of the pixel and its neighbours, the greater the penalty, resulting in a considerable increase in the source value of the matching costs (Figure 9). This procedure ensures a smooth surface by forcing the strings along the path to be somewhat continuous. This process is repeated for each path and each correspondence in the image to obtain the final cost.

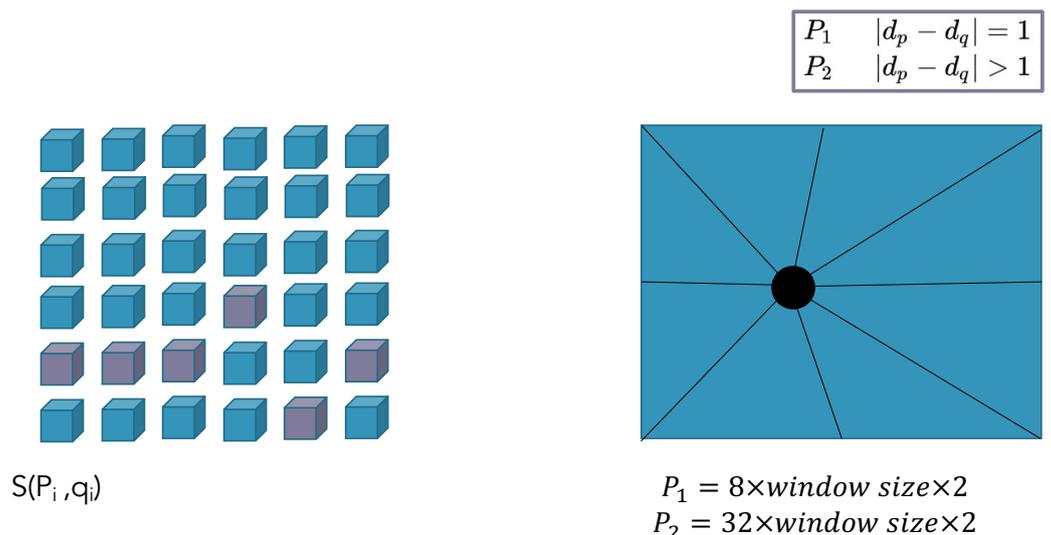


Figure 9. The ultimate cost is the sum of the least costs along picture routes. Eight pathways were used. Cost possibilities are pooled and weighted. P_1 and P_2 are based on neighbourhood disparities, where $N(d)$ is d 's neighbour [17].

To minimise the noise in the computed disparity image, a weighted least squares (WLS) filter [18] is applied [17]. The WLS filter, a well-known edge-preserving smoothing

technique, has weights that highly depend on the image gradients. The final disparity image can be seen in Figure 10.

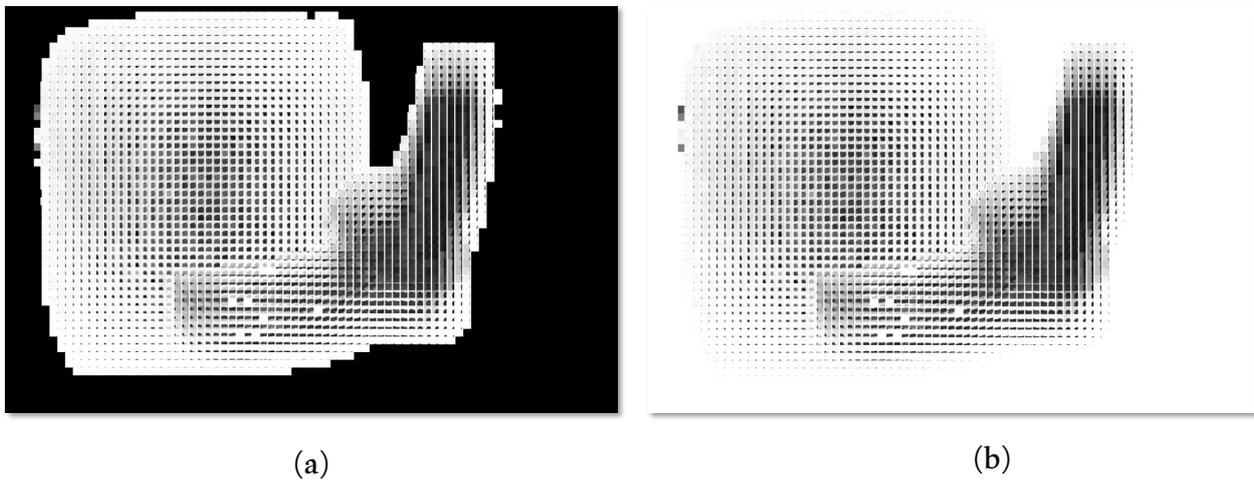


Figure 10. (a) The final disparity image using EIs before background correction. (b) The final disparity image using EIs after background correction. In this disparity map, the darker the pixel, the closer it is to the camera for clarity.

2.3. Background's Disparity Correction

EI-based disparity estimation allows for the recovery of angular information. However, as can be seen in Figure 10a, incorrect disparity may emerge from large texture-less areas such as the background because the EIs only represent segments of the whole scene. Thus, a solution is implemented in which background extraction is first performed to create a background mask, and then the disparity is corrected.

Initially, a disparity map D of the same size as the Holographic image is filled with zeros. Then, the Holographic image's EIs are iterated over to select the left and right pairs:

$$EI_L = EI(i, j), EI_R = EI(i, j + 1); \text{ where } i \in [0, n], j \in [0, m) \quad (10)$$

where i and j are the EIs' locations in the Holographic image of size (n, m) . The disparity for each left and right pair of EIs is computed and the resulting disparity is filtered. $D(i, j)$ is filled with the computed disparity map.

To separate background EIs from foreground EIs, the background threshold value bg_{th} , which is in the range $[0, 1]$ based on the disparity map, is defined. The ratio between non-zero disparity values and the total number of values in the disparity map is computed as r . If this ratio is greater than bg_{th} , $EI(i, j)$ is labelled as a foreground EI; otherwise, it is labelled as a background EI. Increasing the value of bg_{th} will add more images to background EIs, and vice versa.

$$EI(i, j) \begin{cases} M_{bg}, r \leq bg_{th} \\ M_{fr}, r > bg_{th} \end{cases} \quad (11)$$

Background EIs' disparity values are corrected using the correct background disparity values in the foreground EIs, as seen in Figure 11 using colour descriptors [19,20]. To obtain a mask for background regions within foreground EIs ($bg_{r_{fg}}$), the mean and standard deviation of each channel (RGB) of foreground EIs are generated. A pixel in the foreground of an EI is considered to be part of the background if its value is less than one standard deviation from the mean (across all three RGB channels). This presupposes that the majority of foreground image pixels are part of the background region. This implies that the average pixel value should be within one standard deviation of the intensity of the background pixels at the very least. Finally, calculate the mode of the disparity values for $bg_{r_{fg}}$.

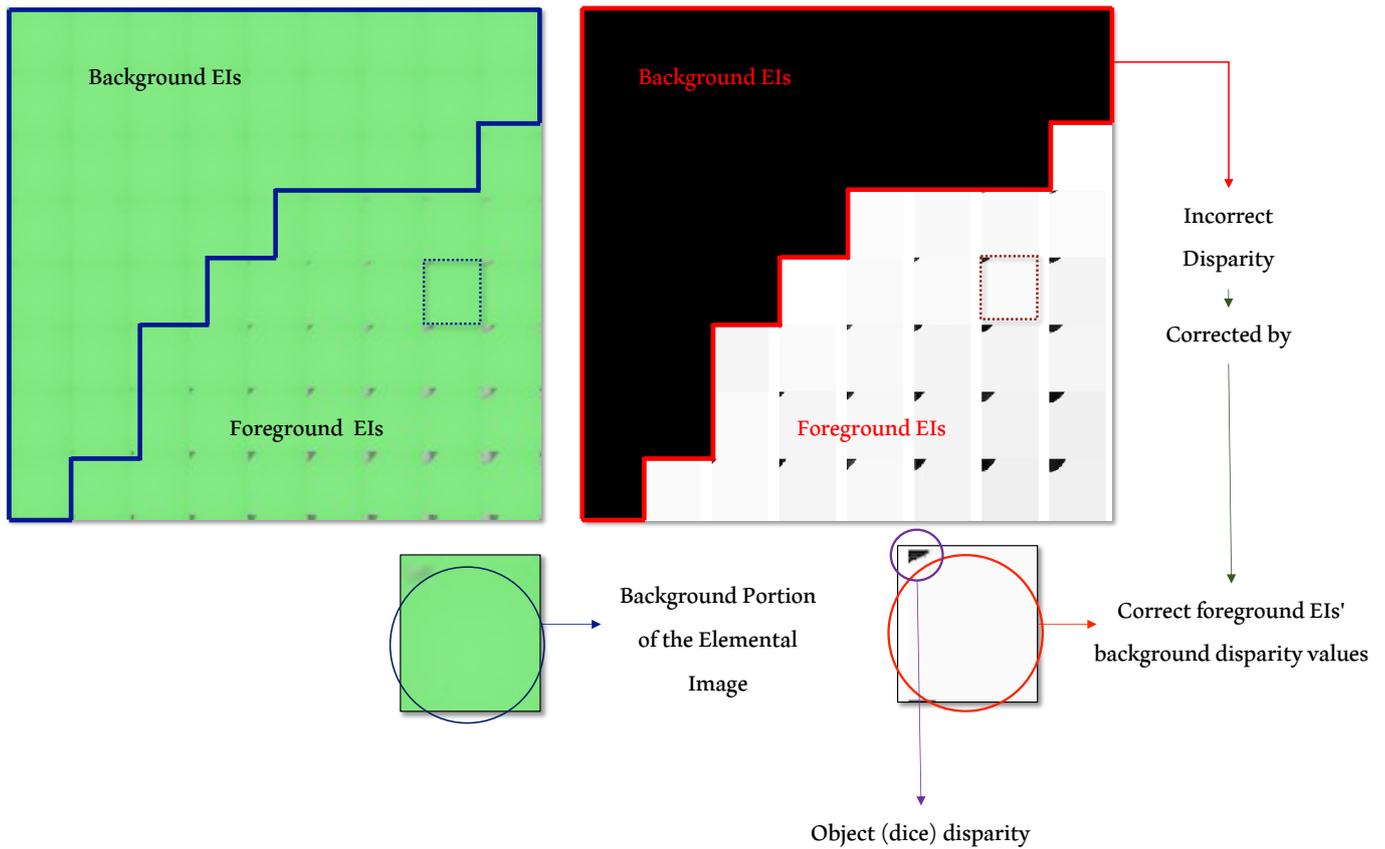


Figure 11. Background EIs' disparity values are corrected using the foreground EIs' background disparity values. **Left:** Holoscopic image showing background and foreground EIs. **Right:** Holoscopic disparity map showing incorrect background EIs being corrected by the background information of the correct disparity of the foreground EIs.

Finally, the mode (mean or median) of the disparity values for bgr_{fg} is substituted for the disparity values (in D) of all background EIs. This acts as a disparity correction step for EIs that only contain background images since stereo SGBM will fail to work for such pairs. Instead, the background disparity is corrected by replacing it with disparity from background regions in foreground EIs. The output result can be seen in Figure 10b, where the background disparity information is fixed.

3. Evaluation

3.1. Dataset

The methodology underwent three evaluations: one comparing the method on VPIs against EIs, another evaluating the method across multiple resolutions, and a third evaluating the method on the same dataset but against two other deep learning methods. This study utilised two Holoscopic datasets to determine the effectiveness of the methodology. The first is a synthetic dataset [21] that was specifically created to replicate the features of Brunel's Holoscopic full-frame camera sensor (Figure 12), which has a sensor size of 35×24 mm and a resolution of 40 megapixels, resulting in image dimensions of 7900×5300 pixels. This dataset has five EI resolutions: 20×20 , 40×40 , 60×60 , 80×80 , and 100×100 pixels. The simulated images were used to evaluate different resolutions and compare them with deep learning techniques. The second dataset was acquired using Brunel's Holoscopic camera. This dataset is utilised because the synthetic one provides flawless VPI and EI pixel mapping, resulting in perfect VPIs that are free from lens effects, distortion, and artefacts. Hence, it is not feasible to directly compare the disparity outcomes between EIs and VPIs derived from the synthetic images.

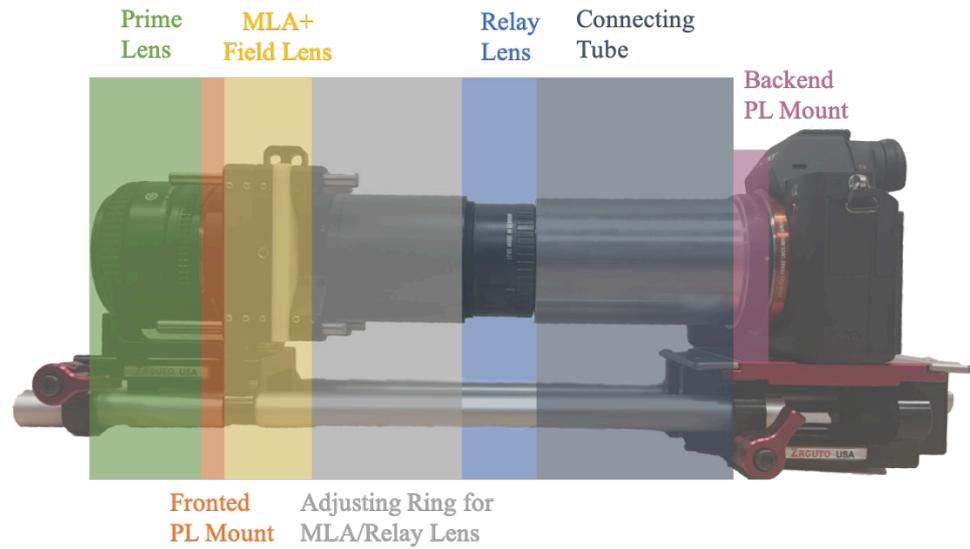


Figure 12. Brunel Holoscopic camera that includes a prime lens, a microlens array, a relay lens to focus light beams onto the sensor, and a CMOS imaging sensor.

3.2. Metrics

Two types of metrics were used to assess the accuracy of the disparity estimation methodology: non-ground truth metrics and ground truth metrics. Ground-truth-based metrics provide dependable evaluation outcomes, but real images from the Brunel camera lack ground truth disparity, necessitating alternative measurements.

Non-ground truth metrics: The consistency check metric, or left–right disparity consistency, evaluates disparity uniformity between left and right images, ensuring pixel correspondence. It is used for refining disparities by scanning both disparities to identify errors at the pixel level, with the error value indicating precision in the disparity map:

$$E = |d_l(x, y) - d_r(x - d_l(x, y), y)| \leq \theta \quad (12)$$

The average error, E_{avg} , calculates the mean disparity error for each pixel:

$$E_{\text{avg}} = \frac{1}{N} \sum_{x,y} |d_l(x, y) - d_r(x - d_l(x, y), y)| \quad (13)$$

Edge alignment evaluates disparity near the edges using the Sobel operator for edge detection. The Mean Absolute Error (MAE) and its normalised version assess the disparity accuracy:

$$\text{MAE} = \frac{1}{N} \sum_{(x,y)} |I_e(x, y) - d_e(x, y)| \quad (14)$$

$$\text{MAE}_{\text{norm}} = 1 - \frac{\text{MAE}}{\text{MAE}_{\text{max}}} \quad (15)$$

Non-ground truth metrics, though less reliable, provide insight into disparity errors.

Ground truth metrics: For synthetic datasets, ground truth metrics include the Mean Absolute Error (MAE) for the average absolute difference between predicted and actual disparities, and the Percentage of Bad Pixels (PBP) for recognising significantly incorrect disparity pixels:

$$\text{MAE} = \frac{1}{N} \sum_{x=1}^W \sum_{y=1}^H |d_e(x, y) - d_{\text{gr}}(x, y)| \quad (16)$$

$$\text{PBP} = \frac{1}{N_p} \sum_{(x,y)} (|d(x, y) - d_T(x, y)| > \delta) \cdot 100 \quad (17)$$

Both MAE and PBP metrics are utilised for evaluation, with values normalised for simplicity.

3.3. Elemental Image Compared to Viewpoint Image

VPIs and EIs are two image structures that can be obtained from Holoscopic images. Previous research has shown significant results in estimating disparity maps utilising VPIs. VPIs can be created by extracting a single pixel from each EI and arranging them in a tiled manner. However, the process of extracting VPIs does not consistently provide ideal images, unlike the VPIs found in synthetic datasets and those obtained from Lytro (the camera's performance was hindered by extensive pre-processing, resulting in a slow performance). The production of these images involves a significant amount of pre-processing. Occasionally, these procedures may require choosing a group of pixels instead of just one, employing shift and integration techniques, and utilising other methodologies to remove artefacts.

As depicted in Figure 13, the Holoscopic images used have undergone calibration and rectification, ensuring that the grid of the EIs aligns perfectly to extract the VPI images accurately. VPI images are extracted using traditional methods, obtaining one pixel per EI. Figure 14 displays three extracted VPIs from different locations. These images exhibit a lower clarity, higher noise, and a reduced resolution when compared to the images typically obtained from publicly available VPI datasets that have undergone extensive pre-processing. The difference in clarity between the EIs and VPIs can be seen in Figures 13 and 14.

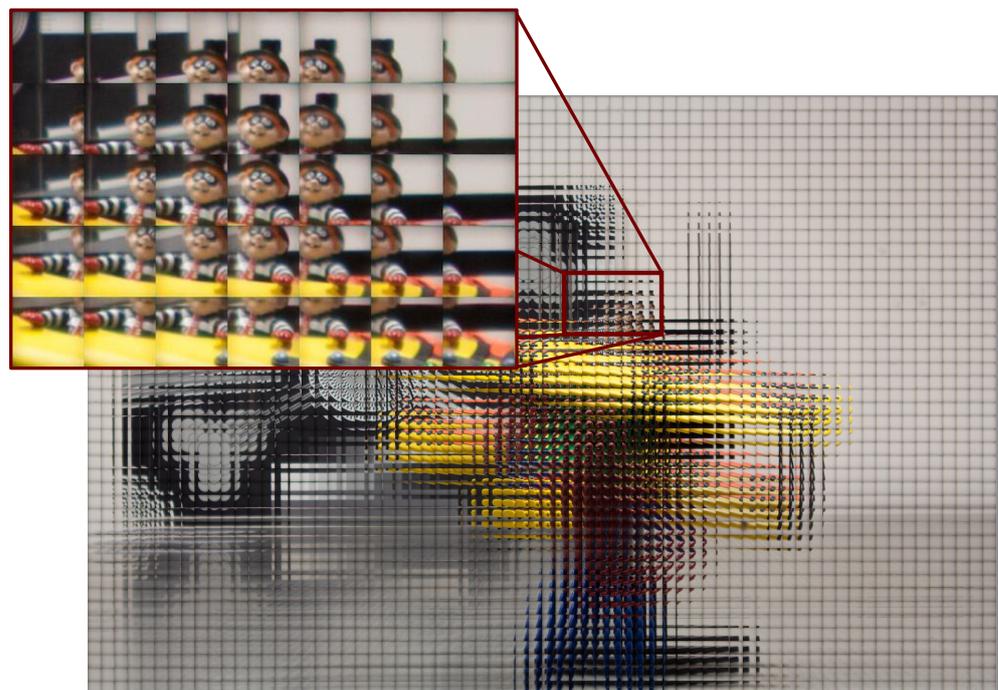


Figure 13. The Holoscopic image was calibrated and rectified, resulting in a total of 68×45 EIs, with each EI measuring 74×74 in size.



Figure 14. VPIs extracted from three different positions (0, 0), (50, 20), (68, 45).

Utilising pixel patches instead of single pixels to extract VPIs during pre-processing might lead to better outcomes, as demonstrated in Figure 15. However, increasing the size of the extracted patch leads to a decrease in angular information, as the number of VPIs obtained is dramatically reduced. The number of VPIs is directly related to the resolution of the EI, which represents the amount of angular information captured. Moreover, while examining Figure 15, it is apparent that the images require additional pre-processing to enhance the outcome. The process of obtaining VPIs also results in a substantial rise in the image generation time. This process can become particularly cumbersome when dealing with Holoscopic videos.

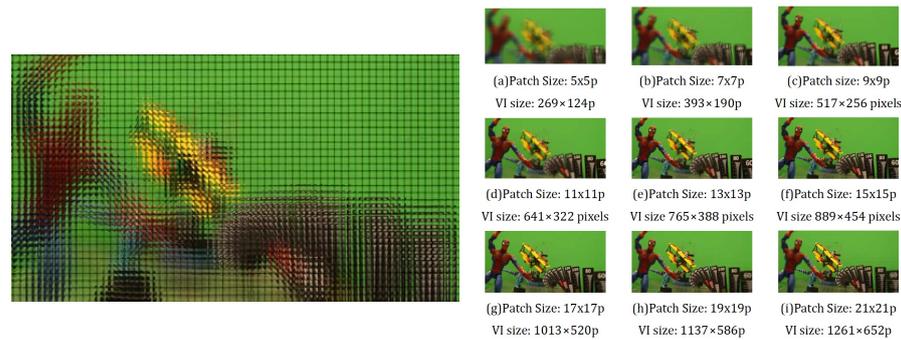


Figure 15. Holoscopic image, Spiderman: (64 × 34 MLA) 5160 × 2743 and sample images from different VPIs retrieved using patch sizes ranging from 5 × 5 pixels to 21 × 21 pixels (p). As seen in the extracted VPIs, they still exhibit some artefacts.

The disparity map was obtained from the EIs of the dataset captured by the Brunel Holoscopic camera using our approach, and subsequently obtained from the extracted VPIs. The disparity maps obtained from VPIs are then transformed to generate EIs, enabling a comparison between EIs with direct disparity estimates and EIs with disparity estimations derived from VPIs, as seen in Figure 16. The closeup crops from the Holoscopic image reveal that the disparity calculated by the EIs is distinct and clear, whereas the EIs obtained from the disparity generated by the VPIs are distorted and ambiguous, as depicted in Figure 17.

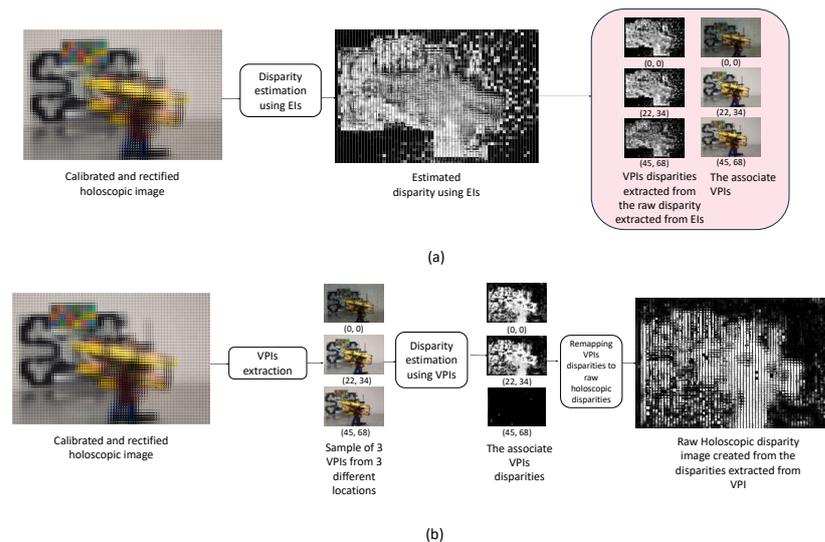


Figure 16. Disparity map derived from EIs and VPIs. (a) The disparity is calculated directly from the EIs using the raw Holoscopic image. Within the red-coloured box, there are a few extracted VPI disparities from the EI-based disparity. Their clarity is compromised by the low resolution. (b) VPIs are extracted from calibrated and rectified Holoscopic images, and the disparity map is obtained from them. These VPI disparities are then mapped back to EIs, allowing for a comparison between VPI-based and EI-based disparity maps.

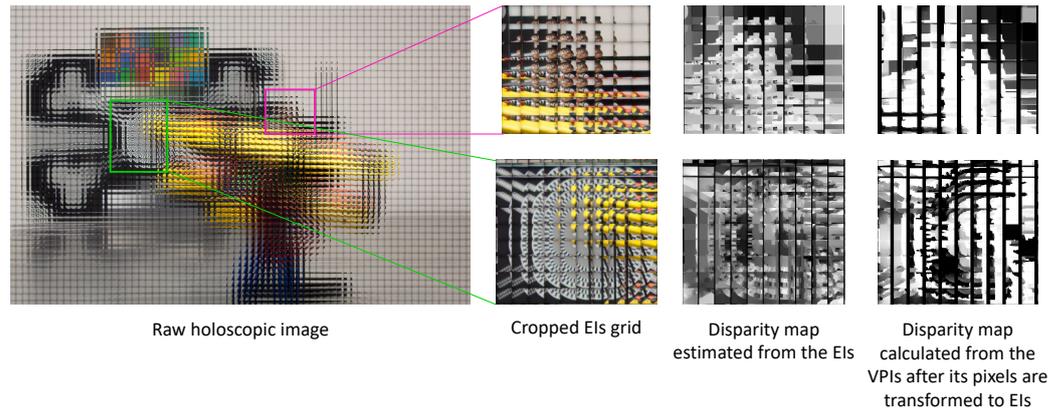


Figure 17. This is a close-up view of a raw Holoscopic image, along with the disparity maps derived from EIs and VPIs. The disparity map created from EIs has greater clarity compared to the one derived from VPIs.

The disparity map, evaluated by the consistency check metric, utilises the entire raw Holoscopic image to optimise the efficiency and minimise the amount of time and effort required. However, the evaluation of disparity using an edge-preserving approach is conducted between individual EIs. This metric is capable of detecting both the grid of EIs and the edges of the features within them. By utilising individual EIs, more reliable results can be obtained.

The edge alignment bar graph in Figure 18(top) illustrates the MAE values for 12 distinct raw real Holoscopic images, which range between approximately 0.352 and 0.781. The changes seen can be attributed to disparities in the scene, texture, colour, and complexity throughout the images. The results generally show lower values (better) in comparison to the edge-alignment metric results derived using VPI disparity where the range of values for different images is approximately 0.498 to 0.797. Overall, EIs demonstrate better results in comparison to VPIs, with an average MAE of approximately 0.523, whereas VPIs have an average MAE of approximately 0.681, which can be viewed in the averaged bar “All”.

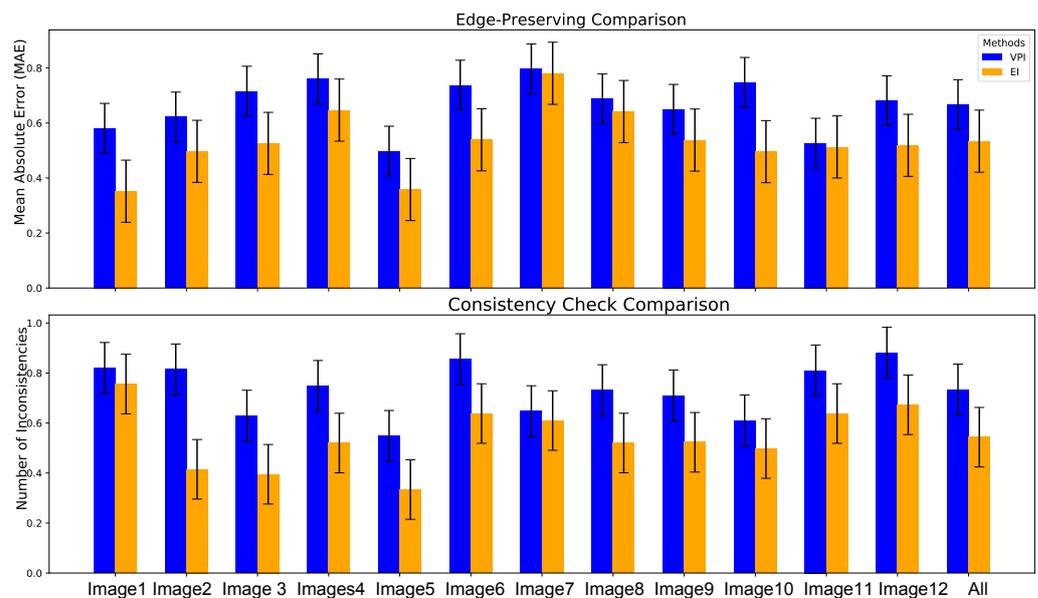


Figure 18. The bar graphs display the edge-alignment matrices (top) and consistency check matrices (bottom) calculated from 12 raw Holoscopic images captured by the Brunel Holoscopic camera. The averaged result is labelled as “All”. EIs generally outperform VPIs, as seen by their lower average MAE and consistency check metric.

The bar graph depicted in Figure 18(bottom) illustrates the range of values for the consistency check metric derived from the disparity of EIs and VPIs. The values vary between approximately 0.334 and 0.756. The results exhibit lower values when compared to the consistency check metric results generated using VPI disparity, where the values range between approximately 0.548 and 0.881 for different images. Overall, EIs yield better results compared to VPIs, exhibiting average values of roughly 0.520, while VPIs demonstrate an average value of around 0.731. A selection of four raw Holoscopic images is shown in Figure 19. Simple scenes were captured to compare the disparity of EIs vs. VPIs directly, rather than assessing the algorithm in a complicated scene configuration.

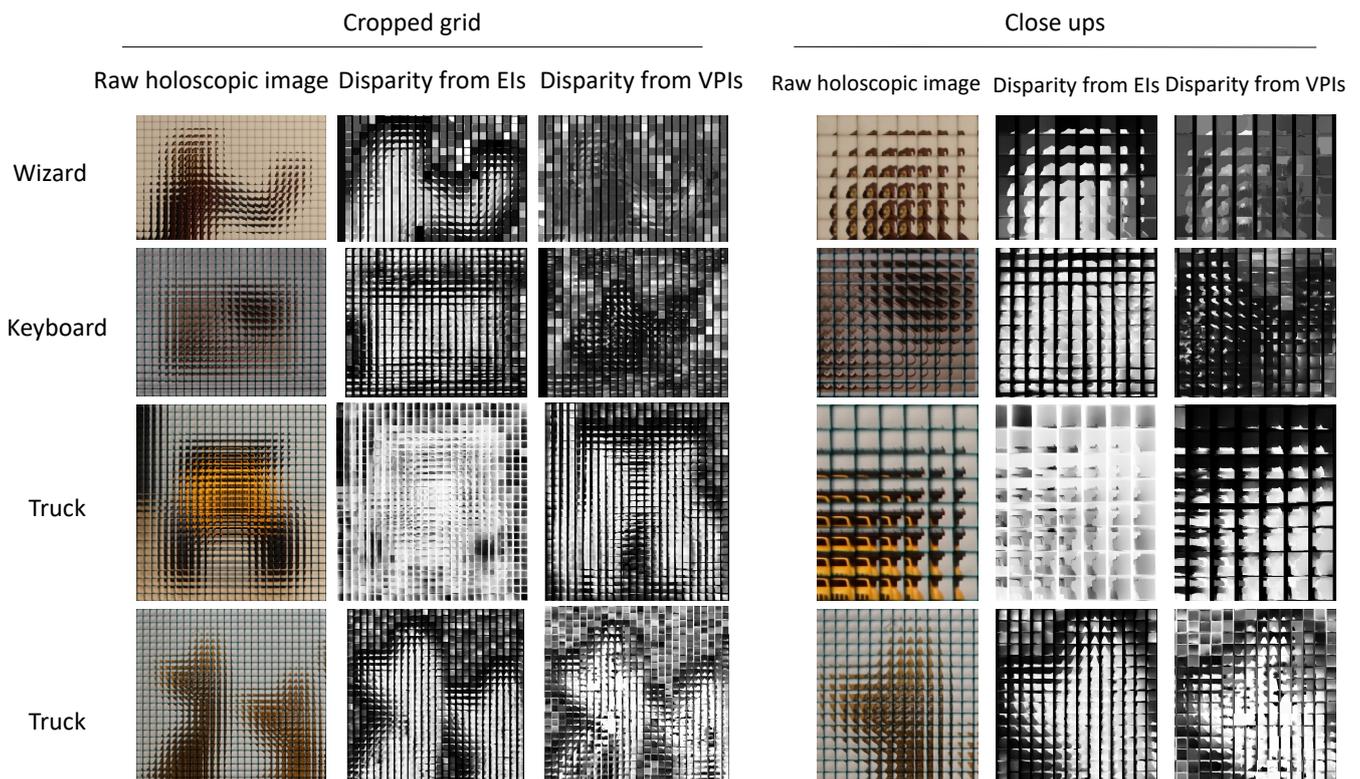
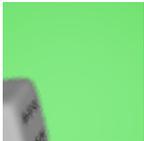


Figure 19. The algorithm was tested on a total of 12 real Holoscopic images. This is a collection of four images showcasing close-up sections to illustrate the disparity between EIs and VPIs.

3.4. Elemental Image Compared to Viewpoint Image Resolution

The algorithm's performance was assessed by utilising 24 synthetic Holoscopic images with five distinct EI resolutions: 20×20 , 40×40 , 60×60 , 80×80 , and 100×100 , as shown in Figure 20. The MAE and PBP were calculated for all resolutions. As depicted in Figure 21, an increase in the EI's resolution does not consistently result in an improved accuracy. EIs with a high resolution are expected to lead to a high score. Yet, the clarity of the EIs relies on the clarity of the produced VPIs. Smaller EIs typically originate from VPIs with higher resolutions compared to those larger EIs (trade-off in resolution), which allows for more information to be presented in the EIs. This ultimately leads to a sharper image, as seen in Table 1. This table displays a single EI on various scales. Although the EI with a resolution of 100×100 is larger, it is noticeable that the circles on the dice in the EI with a resolution of 60×60 are more defined and sharper. As the scale increases, the EI loses more information, resulting in the presence of noisy features. Future research can employ this dataset with many resolutions to construct a multi-resolution pyramid, thereby capturing all the accessible information at each level of resolution.

Table 1. EIs of three different scales, original, down-sampled, and up-sampled.

EI Slice Resolution	20 × 20	40 × 40	60 × 60	80 × 80	100 × 100
Original Resolution					
Scaled-Down (20 × 20)					
Scaled-Up (100 × 100)					

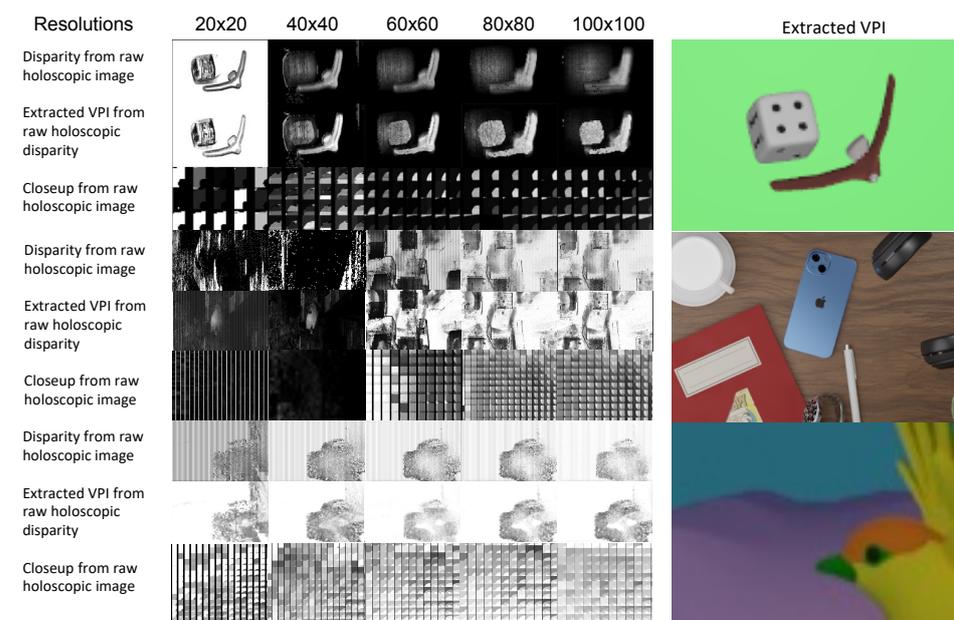


Figure 20. Example of three Holoscopic images alongside their calculated disparities at various resolutions. Observing the disparity from the low-resolution images is difficult. Consequently, close-up views are offered.

As seen in Figure 21, the MAE values for the methodology across different resolutions reveal varying degrees of accuracy. The MAE for the 20 × 20 resolution ranges between 0.673 and 0.804, suggesting a significant amount of errors. For the 40 × 40 resolution, the MAE ranges between 0.613 and 0.755, indicating a significantly enhanced performance in comparison to the 20 × 20 resolution. The 60 × 60 resolution’s MAE ranges from 0.430 to 0.650, demonstrating a significant improvement in accuracy compared to the lower resolutions. The MAE for the 80 × 80 resolution varies between 0.419 and 0.625, indicating a higher level of precision. At a resolution of 100 × 100, the MAE varies between 0.462 and 0.640, suggesting a somewhat lower level of precision compared to the 80 × 80 resolution.

The PBP for the 20 × 20 resolution ranges from 68.8% to 86.1%, suggesting a significant presence of bad pixels. The PBP of the 40 × 40 resolution falls within the range of 72.7%

to 87.0%, indicating comparable performance to that of the 20×20 resolution. Moving to 60×60 resolution, the range is from 49.4% to 66.9%, suggesting a significant reduction in bad pixels compared to the lower levels. With an 80×80 resolution, the PBP falls between 39.6% and 64.1%, indicating a significant enhancement in performance and a reduction in the number of bad pixels. Finally, at 100×100 resolution, the range is from 44.7% to 64.3%, exhibiting an accuracy that is slightly lower than 80×80 .

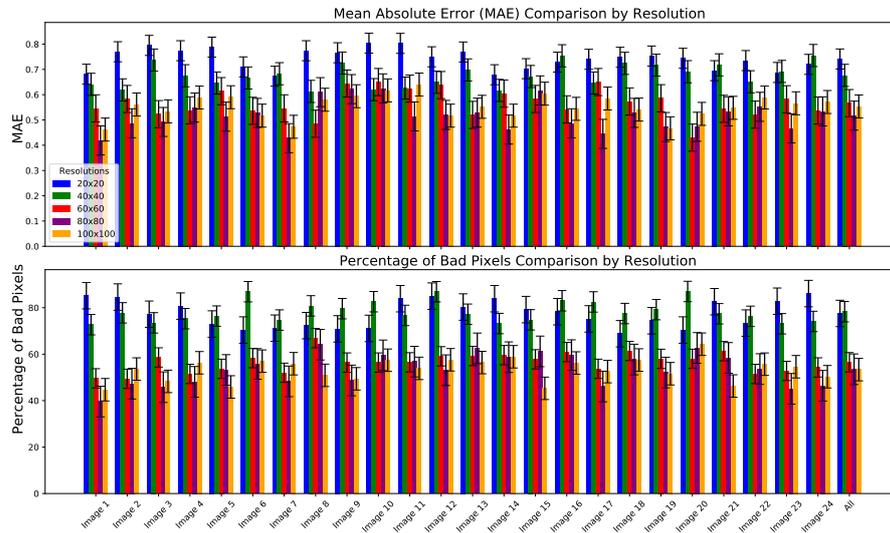


Figure 21. The bar graph illustrates the performance of disparity calculated at five different resolutions. The graph shows that the EIs achieve the highest level of precision at a resolution of 80×80 , followed by 100×100 .

Images with a low resolution, such as 20×20 and 40×40 images, still have a noticeably reduced accuracy. This is because achieving accurate disparity typically requires a combination of a wide baseline and a high-resolution image. Since larger EIs demonstrate a greater baseline and a reduced number of texture-less EIs, as depicted in Figure 22, the results for high-resolution EIs are better than the accuracy in low-resolution images.

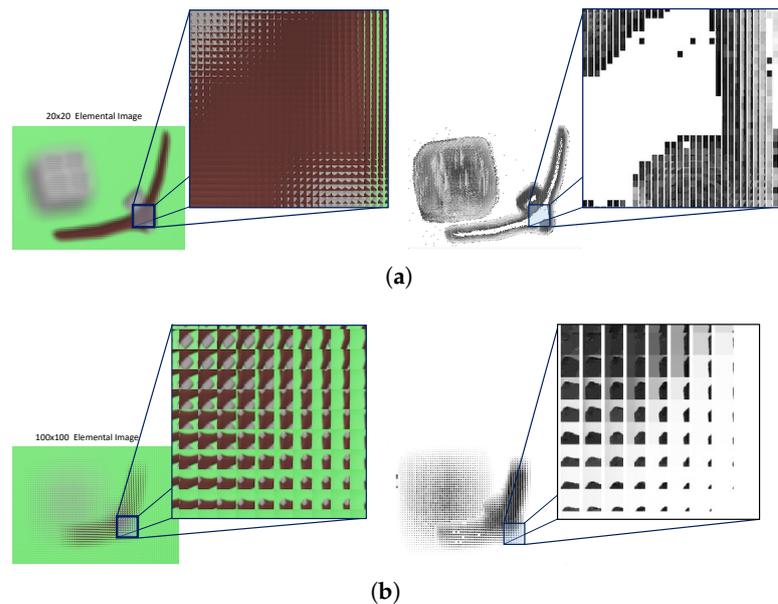


Figure 22. (a) shows the output of 20×20 pixel EIs with a significant texture-less area, leading to an incorrect disparity computation. In (b), the outcome of 100×100 EIs taken from the same point with a wider baseline and a larger portion of the objects presented is a more accurate disparity.

3.5. Comparative Analysis of Stereo-Matching Networks

The results from raw Holoscopic images (EIs) were also compared against two state-of-the-art deep learning stereo-matching algorithms: the methods by Zhang et al. [22] and Chang and Chen [23]. Zhang et al. [22] proposed a technique to enhance the generalisation abilities of stereo-matching networks. Their main objective was to maintain the consistency of features between corresponding pixels. Their methodology combines pixel-level contrastive learning with a stereo-selective whitening loss to enhance the consistency of features across various domains. This technique is highly versatile and may be easily integrated into pre-existing networks without any disruptions.

Chang and Chen [23] employed supervised learning and convolutional neural networks (CNNs) to address the task of estimating disparities from stereo image pairs. They proposed a Pyramid Stereo-Matching Network (PSMNet) as an alternative to the patch-based Siamese networks commonly employed in current architectures. The PSMNet overcomes the limitation of incorporating contextual information in uncertain regions by incorporating spatial pyramid pooling and a 3D CNN.

Both of the pre-trained models were used to extract disparity from all 24 raw Holoscopic images in the dataset, and an 80×80 resolution was chosen based on the accuracy level from the previous section. These results were then compared with those obtained from this paper's method applied to the same dataset. The disparity outcome of a basic EI of both deep learning models was highly blurred and undefined, as seen in Figure 23. These outcomes can be attributed to various factors, including the dissimilar characteristics of the higher-resolution stereo images utilised for training the models developed by Zhang et al. [22] and Chang and Chen [23] compared to the low-resolution and low-texture EIs. Therefore, when these models are employed on the EIs, they struggle with accurately capturing intricate details. Furthermore, the efficacy of these models is greatly influenced by their specific architecture, particularly Chang and Chen [23]'s PSMNet, which further reduces the resolution of low-resolution EIs, resulting in unsatisfactory outcomes.

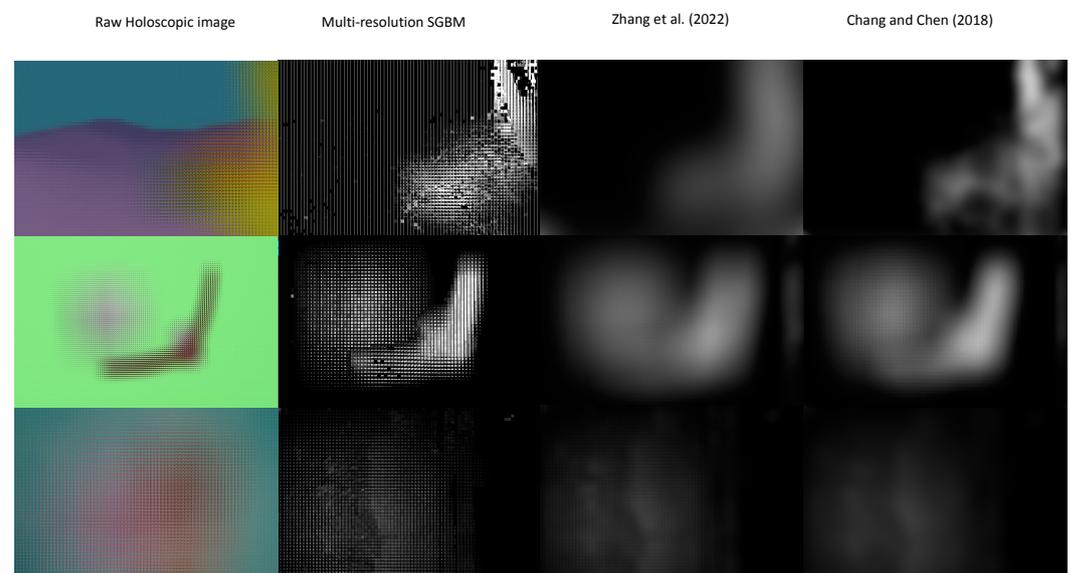


Figure 23. Both deep learning algorithms result in blurry and undefined EIs compared to our result [22,23].

As depicted in Figure 24, MAE values for this paper's method ranged from 0.419 to 0.625, which were considerably lower than the MAE values reported by Zhang et al. [22], ranging from 0.637 to 0.801, and Chang and Chen [23], ranging from 0.686 to 0.798. Regarding the PBP, this paper's method achieved percentages ranging from 39.6% to 64.1%, which indicates a superior performance. In comparison, Zhang et al. [22] and Chang and Chen [23] obtained greater percentages, with ranges of 63.3% to 78.7% and 65.7% to 77.9%, respectively.

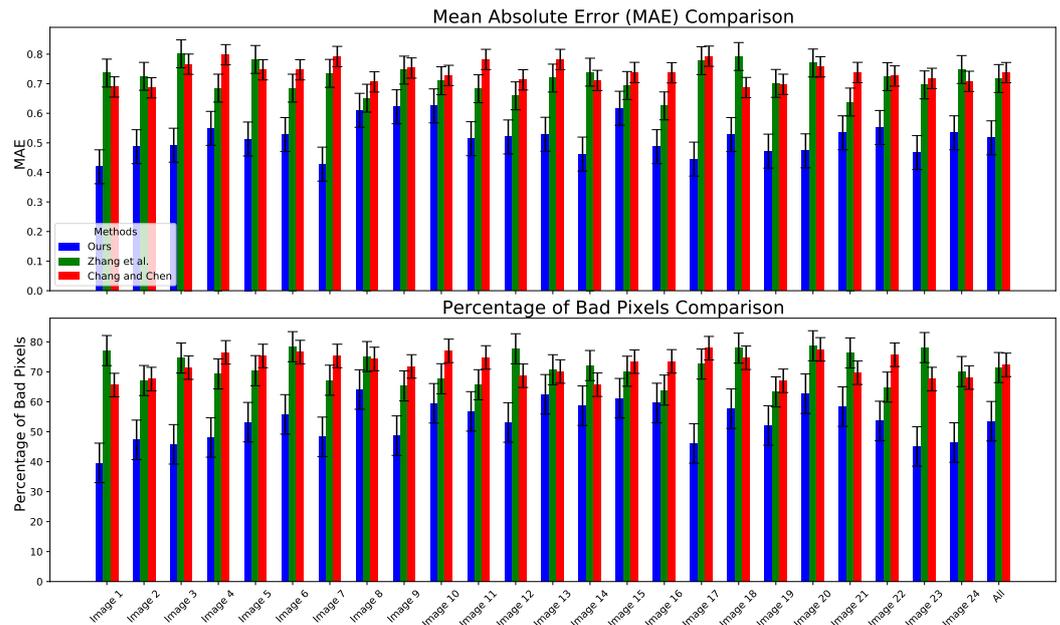


Figure 24. The bar graph depicts the comparative performance of the extracted disparity to that of the methods of Zhang et al. [22] and Chang and Chen [23], demonstrating that our method outperforms both methods' algorithms.

4. Conclusions

This study introduces a novel method for disparity estimation in Holographic 3D imaging, leveraging angular information from Elemental Images (EIs) over traditional spatial data from Viewpoint Images (VPs). Our goal was to evaluate if EIs could serve as a more accurate foundation for disparity estimation, offering an alternative to conventional methods.

Through detailed experimentation, we developed an approach that not only generates accurate disparity maps from EIs but also surpasses both traditional strategies and advanced deep learning algorithms. This achievement stems from an innovative application of the Semi-Global Block Matching (SGBM) method, enhanced by multi-resolution techniques and content-aware analysis, optimizing the use of EIs' angular data.

Our findings indicate that EIs, even without extensive manipulation, provide a more reliable basis for disparity estimation than VPs. This suggests that the inherent angular information in EIs is better suited for precise disparity assessments, paving the way for advancements in Holographic 3D imaging technology.

Furthermore, our method outperformed comparable deep learning models, a result attributable to EIs' unique characteristics, notably their lower texture and resolution. This highlights a critical gap in current deep learning approaches: the lack of training on datasets specifically designed for the nuances of EIs. Our study underscores the urgent need to create comprehensive EI datasets for training deep learning models for EI-based applications, promising significant progress in automated disparity estimation.

Moreover, our investigation into EI resolutions revealed that higher resolutions do not always equate to a better disparity estimation accuracy. We identified an optimal resolution range for EIs, challenging the assumption that higher is always better and offering insights into how resolution influences angular information extraction for disparity calculations.

Additionally, the potential of EIs for developing compact depth-sensing devices opens new possibilities for their use in microscopes, mobile devices, medical instruments like endoscopes, and real-time depth estimation applications such as autonomous vehicles and medical diagnostics. This is due to EIs' ability to be captured directly by image sensors, simplifying the depth estimation process.

In conclusion, our research confirms the viability of using angular perspective data from EIs for disparity estimation in Holographic 3D imaging, marking a significant advance-

ment over traditional methods. The implications for future 3D imaging technologies are vast, necessitating continued research to unlock the full potential of EIs in enhancing the depth estimation accuracy and efficiency across various applications.

Author Contributions: Conceptualisation, M.R.S.; methodology, B.A.; software, B.A.; validation, B.A., H.M. and M.R.S.; formal analysis, B.A.; investigation, B.A.; resources, B.A.; data curation, B.A.; writing—original draft preparation, B.A.; writing—review and editing, B.A. and H.M.; visualisation, B.A.; supervision, H.M.; project administration, B.A. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: The data presented in this study are available on request from the corresponding author, Bodor Almatrouk, at bodor.almatrouk@brunel.ac.uk. The data are not publicly available due to commercial privacy.

Conflicts of Interest: The authors declare no conflicts of interest.

References

- Cheng, Z.; Xiong, Z.; Chen, C.; Liu, D. Light field super-resolution: A benchmark. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, Long Beach, CA, USA, 16–17 June 2019.
- Lumsdaine, A.; Georgiev, T. Full resolution lightfield rendering. *Indiana Univ. Adobe Syst. Tech. Rep.* **2008**, *91*, 92.
- Lumsdaine, A.; Georgiev, T. The focused plenoptic camera. In Proceedings of the 2009 IEEE International Conference on Computational Photography (ICCP), San Francisco, CA, USA, 16–17 April 2009; IEEE: Piscataway, NJ, USA, 2009; pp. 1–8.
- Georgiev, T.; Lumsdaine, A. Reducing plenoptic camera artifacts. *Comput. Graph. Forum* **2010**, *29*, 1955–1968. [[CrossRef](#)]
- Ng, R.; Levoy, M.; Brédif, M.; Duval, G.; Horowitz, M.; Hanrahan, P. Light field photography with a hand-held plenoptic camera. *Comput. Sci. Tech. Rep. CSTR* **2005**, *2*, 1–11.
- Kinoshita, T.; Ono, S. Depth estimation from 4D light field videos. In Proceedings of the International Workshop on Advanced Imaging Technology (IWAIT) 2021, Online, 5–6 January 2021; SPIE: Paris, France, 2021; Volume 11766, pp. 56–61.
- Mousavi, M.; Khanal, A.; Estrada, R. Ai playground: Unreal engine-based data ablation tool for deep learning. In Proceedings of the International Symposium on Visual Computing, San Diego, CA, USA, 5–7 October 2020; Springer: Berlin/Heidelberg, Germany, 2020; pp. 518–532.
- Tankus, A.; Kiryati, N. Photometric stereo under perspective projection. In Proceedings of the Tenth IEEE International Conference on Computer Vision (ICCV'05), Beijing, China, 17–21 October 2005; IEEE: Piscataway, NJ, USA, 2005; Volume 1, pp. 611–616.
- Park, J.H.; Baasantseren, G.; Kim, N.; Park, G.; Kang, J.M.; Lee, B. View image generation in perspective and orthographic projection geometry based on integral imaging. *Opt. Express* **2008**, *16*, 8800–8813. [[CrossRef](#)] [[PubMed](#)]
- Eagle, R.; Hogervorst, M. The role of perspective information in the recovery of 3D structure-from-motion. *Vis. Res.* **1999**, *39*, 1713–1722. [[CrossRef](#)] [[PubMed](#)]
- Thomas, G.A.; Stevens, R.F. Processing of Images for 3D Display. US Patent 6,798,409, 28 September 2004.
- Hirschmuller, H. Stereo processing by semiglobal matching and mutual information. *IEEE Trans. Pattern Anal. Mach. Intell.* **2007**, *30*, 328–341. [[CrossRef](#)] [[PubMed](#)]
- Elad, M. On the origin of the bilateral filter and ways to improve it. *IEEE Trans. Image Process.* **2002**, *11*, 1141–1151. [[CrossRef](#)]
- Buades, A.; Facciolo, G. Reliable multiscale and multiwindow stereo matching. *SIAM J. Imaging Sci.* **2015**, *8*, 888–915. [[CrossRef](#)]
- Miangoleh, S.M.H.; Dille, S.; Mai, L.; Paris, S.; Aksoy, Y. Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Nashville, TN, USA, 20–25 June 2021; pp. 9685–9694.
- Keys, R. Cubic convolution interpolation for digital image processing. *IEEE Trans. Acoust. Speech Signal Process.* **1981**, *29*, 1153–1160. [[CrossRef](#)]
- Almatrouk, B.; Meng, H.; Swash, M.R. Disparity estimation from holoscopic elemental images. In Proceedings of the International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery, Xi'an, China, 1–3 August 2020; Springer: Berlin/Heidelberg, Germany, 2020; pp. 1106–1113.
- Liu, W.; Chen, X.; Shen, C.; Liu, Z.; Yang, J. Semi-global weighted least squares in image filtering. In Proceedings of the Proceedings of the IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 5861–5869.
- Alzu'bi, A.; Amira, A.; Ramzan, N. Semantic content-based image retrieval: A comprehensive study. *J. Vis. Commun. Image Represent.* **2015**, *32*, 20–54. [[CrossRef](#)]
- Rublee, E.; Rabaud, V.; Konolige, K.; Bradski, G.R. ORB: An efficient alternative to SIFT or SURF. In Proceedings of the ICCV, Barcelona, Spain, 6–13 November 2011; Volume 11, p. 2.

21. Almatrouk, B.; Meng, H.; Aondoakaa, A.; Swash, R. A New Raw Holographic Image Simulator and Data Generation. In Proceedings of the 2023 8th International Conference on Image, Vision and Computing (ICIVC), Dalian, China, 27–29 July 2023; IEEE: Piscataway, NJ, USA, 2023; pp. 489–494.
22. Zhang, J.; Wang, X.; Bai, X.; Wang, C.; Huang, L.; Chen, Y.; Gu, L.; Zhou, J.; Harada, T.; Hancock, E.R. Revisiting domain generalized stereo matching networks from a feature consistency perspective. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, New Orleans, LA, USA, 18–24 June 2022; pp. 13001–13011.
23. Chang, J.R.; Chen, Y.S. Pyramid stereo matching network. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–23 June 2018; pp. 5410–5418.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.