

Metameric Varifocal Holograms

David R. Walton*
UCL

Koray Kavaklı†
Koç University

Rafael Kuffner dos Anjos‡
University of Leeds

David Swapp§
UCL

Tim Weyrich||
UCL

Hakan Urey||
Koç University

Anthony Steed**
UCL

Tobias Ritschel††
UCL

Kaan Aksit‡‡
UCL

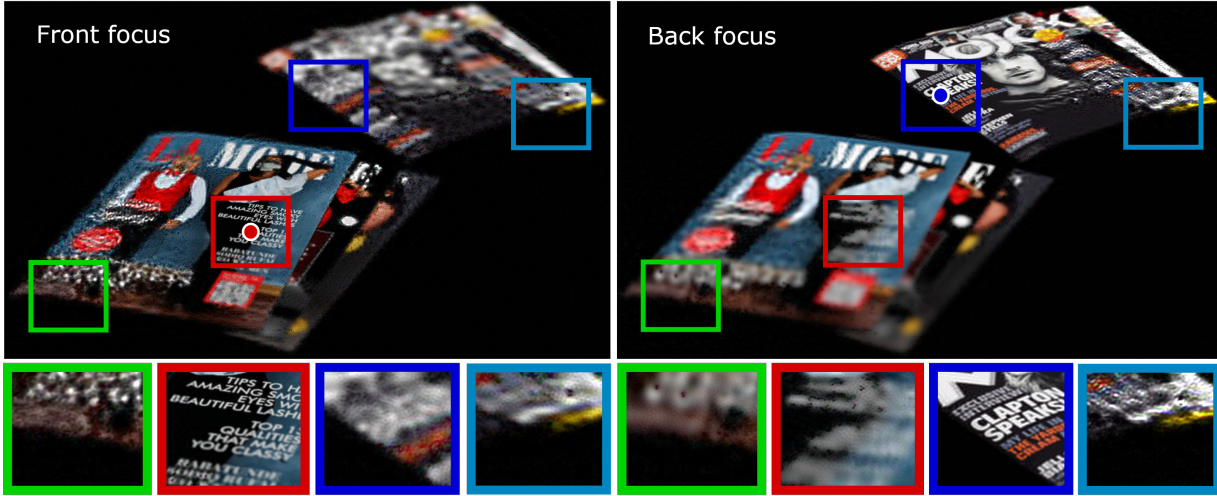


Figure 1: Simulated reconstructions of metameric varifocal holograms. Our holograms reconstruct single-plane images at the correct focus levels, reconstructing high-resolution visuals at a user’s fovea while displaying statistically correct content across their peripheral vision indistinguishable from the target images (metamers). Top row: simulated image reconstructions at two different focus levels (gaze location marked with a dot). Bottom row: zoomed-in insets from these two reconstructions. All foveated images in this paper are best viewed at a 60 cm wide display from a distance of 80 cm. (Three-dimensional assets from Vilém Duha ©2021)

ABSTRACT

Computer-Generated Holography (CGH) offers the potential for genuine, high-quality three-dimensional visuals. However, fulfilling this potential remains a practical challenge due to computational complexity and visual quality issues. We propose a new CGH method that exploits gaze-contingency and perceptual graphics to accelerate the development of practical holographic display systems. Firstly, our method infers the user’s focal depth and generates images only at their focus plane without using any moving parts. Second, the images displayed are metamers; in the user’s peripheral vision, they need only be statistically correct and blend with the fovea seamlessly. Unlike previous methods, our method prioritises and improves foveal visual quality without causing perceptually visible distortions at the periphery. To enable our method, we introduce a novel metameric loss function that robustly compares the statistics of two given images for a known gaze location. In parallel, we implement a model representing the relation between holograms and their image recon-

structions. We couple our differentiable loss function and model to metameric varifocal holograms using a stochastic gradient descent solver. We evaluate our method with an actual proof-of-concept holographic display, and we show that our CGH method leads to practical and perceptually three-dimensional image reconstructions.

Keywords: Computer-Generated Holography, Foveated Rendering, Metamerisation, Varifocal Near-Eye Displays, Virtual Reality, Augmented Reality

Index Terms: Computing methodologies—Computer graphics—Graphics systems and interfaces—Perception; Hardware—Communication hardware, interfaces and storage—Displays and imagers

1 INTRODUCTION

In recent years, improving display technology to enable lifelike three-dimensional visuals has attracted much attention from industry and academia as displays are crucial for future Human-Computer Interaction (HCI) [33]. An emerging trend, Computer-Generated Holography (CGH) [47], promises such realistic visuals in the next-generation of displays [29]. Unlike conventional displays, pixelated images are not sent to holographic displays directly.

In a typical phase-only Spatial Light Modulator (SLM)-based holographic display, laser light illuminates an array of pixels which modulate the phase of the light. The reflected light interferes to produce the image. Finding the correct phase values to send to the SLM is challenging due to the complexity of light transport. Also, SLMs have limited resolution. As a result, real CGH displays suffer from noise and other artefacts.

Gaze-contingent approaches [1, 27] are often used to reduce the hardware and computational requirements of displays. In this work,

*e-mail: david.walton.13@ucl.ac.uk

†e-mail: kkavakli@ku.edu.tr

‡e-mail: r.kuffnerdosanjos@leeds.ac.uk

§e-mail: d.swapp@ucl.ac.uk

||e-mail: t.weyrich@cs.ucl.ac.uk

||e-mail: hurey@ku.edu.tr

**e-mail: a.steed@ucl.ac.uk

††e-mail: t.ritschel@ucl.ac.uk

‡‡e-mail: k.aksit@ucl.ac.uk

we explore whether gaze-contingency for CGH can help meet the demands of the Human Visual System (HVS) in practice. Knowing the user’s gaze gives us two critical pieces of information we exploit.

First, it tells us which parts of the image fall in the user’s periphery, rather than their fovea. To exploit this, we draw inspiration from the state of the art in foveated graphics literature [53]. This work focuses on generating visuals which are not pixel-accurate to a target image in the periphery of the user’s vision, but are still perceived as identical to the target. We exploit their work to dedicate more of the expressive power of the SLM to generating high-quality visuals at the fovea as described in work by Chakravarthula et al. [11]. In contrast, visuals at the periphery need only be statistically correct (in a sense precisely described in Sec. 3.2) and will still be perceived as accurate. As highly accurate simulation models become available in the future, such a method can pave the way towards distributing the speckle noise at a holographic display [13] in a statistically correct way, enabling indistinguishable images at the periphery in the future.

Second, given the depths of each pixel in the displayed image, it allows us to infer the user’s current focal depth. We can use this information to only enforce our reconstruction to be correct at the user’s current focus. Whilst CGH is certainly capable of displaying multi-plane images, this often leads to image quality issues as the hologram pixels are used to deliver images at multiple planes at once. For that purpose, we draw inspiration from existing literature on varifocal near-eye displays [2, 34] and varifocal holograms [37]. We argue that generating images at a single plane instead of multiple planes will help assure quality in visuals generated by CGH. We combine these arguments to enable CGH computation pipelines that are perceptually accurate and offer high visual quality.

Specifically, this work introduces the following contributions:

(1) Metameric loss function. We introduce a fast metameric loss that can help us quantify image quality within the peripheral field of view by comparing the statistics of images. We believe this loss couples well with a gaze-contingent display and graphics application, specifically holographic displays, as they are often proposed as the next-generation display technology.

(2) Metameric varifocal holograms. We introduce a complete optimisation pipeline for metameric varifocal holograms using our metameric loss function. Note that our holograms change focus in a gaze-contingent manner, avoiding the complexity of representing light fields or multiplane images using CGH;

(3) Proof-of-concept prototype. We build a single colour holographic display to experiment with our metameric varifocal holograms. We assess the results of our CGH method using this proof-of-concept display.

2 RELATED WORK

Our work combines the state of the art in visual perception and CGH while relying on gaze contingency. Hence, we review the relevant work in visual perception, gaze-contingent displays and CGH fields.

2.1 Gaze-contingent displays

Eye-gaze tracking [4, 30] is of great interest to AR and VR research. A major reason for this is that visual [17] and depth acuity [55] of the HVS drops sharply with increasing eccentricity towards far peripheral vision. Combined with the visual and depth acuity of HVS, eye-gaze information opens up opportunities towards reducing computational and hardware complexity of displays. Such displays that take advantage of eye-gaze information are known as gaze-contingent displays.

A form of gaze-contingent displays – foveated displays [48] – present images at high resolution at the fovea and lower resolution at the periphery. A foveated display tracks a user’s gaze and can either actively move a foveal inset display [51], move both foveal and peripheral insets [27], or change the distribution of resolution

by distorting optical fields [1] to generate images with fewer pixels but with no perceptual difference.

Our work falls into the category of foveated displays. It merges the ideas of distorting [1] or shifting [2] optical fields while taking advantage of CGH in a foveated manner following the spirit of the work by Chakravarthula et al. [11]. A major advantage of CGH in this setting is that it can facilitate foveated rendering without the need for moving parts in the form of displays or lenses.

2.2 Metamers

Most popular foveated rendering approaches focus on decreasing resolution with increasing eccentricity [19, 56]. However, traditional literature on human vision [5] refers to the objects in the periphery as difficult to see and different, but not particularly blurry. Objects are not only less sharp, but the size of stimuli [60], visual crowding [21] and texture content [52] also play an important part in how things are perceived. For a comprehensive review of the behavior and the limitations of peripheral vision, we recommend the review from Rozenholtz [41].

With this in mind, Freeman and Simoncelli [18] showed that it is possible to devise a process to generate “ventral metamers”; pairs of images perceived as identical for a given fixation point (see Fig. 2). Their work [18] models the correlation between the size of pooling regions and different eccentricities, describing visuals in the periphery with local image statistics rather than individual pixel values. Unfortunately, their process is computationally expensive as it depends on complex image statistics and iterative optimisation processes. Deza et al. [14] proposed an approach to approximate this effect using techniques from style transfer [20], blending Visual Geometry Group (VGG) network [46] features of the target image with those from a noise image using Adaptive Instance Normalisation (AdaIN)[23] over foveated pooling regions. Similarly, recent work from Surace et al. [49] uses a texture synthesis approach combined with Generative Adversarial Networks (GANs) to generate ventral metamers. There is no guarantee that the generated visuals will be statistically correct for any known human vision model in both cases. These techniques are also unable to operate at interactive framerates, although they are significantly faster than that of Freeman and Simoncelli [18].

To our knowledge, the first work that achieves generation of ventral metamers at interactive rates is the work by Walton et al. [53], which uses a simplified statistical model focused on using fast calculated means and variances of a steerable pyramid [38]. Their simplified model allows fast synthesis of metamers by scaling and biasing bands of a steerable pyramid constructed from a noise image.

Our work reformulates their model [53] as a general-purpose differentiable loss function. This opens it up to a range of other applications, including but not limited to CGH as described in this paper.

2.3 Computer-Generated Holography

CGH has garnered much interest from the research community in recent years. This interest primarily stems from the widespread availability of powerful, highly parallel processors coupled with modern machine-learning frameworks that automatically differentiate given models [36]. These advancements accelerate and improve the accuracy of hologram generation (phase retrieval), particularly when taking advantage of advances in deep learning [59]. As a result, modern phase retrieval techniques offer dramatically improved image quality over classical hologram calculation methods such as the Gerchberg-Saxton method [57]. From the recent past, the work by Chakravarthula et al. [9] revisits Wirtinger complex derivatives and shows that the visual quality of two-dimensional image reconstructions in CGH can be improved in common Gerchberg-Saxton and Double-phase coding [22] approaches. The works by Peng et al. [37] and Chakravarthula et al. [10] help to bridge the gap between

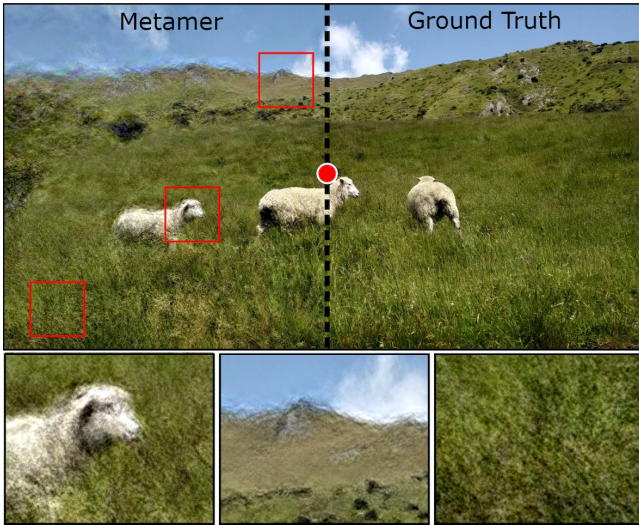


Figure 2: A sample metamer display the top row is generated by following the work by Walton et al. [53] using a gaze location at the center of the image (red dot at the center). For comparison purposes, we also show the ground truth image at the right portion of the same image. Red highlighted regions from the top image are zoomed-in and provided as insets at the bottom row.

CGH simulations and actual image reconstructions in a physical display by learning a model of display hardware using a camera and convolutional neural networks. Their findings have drastically improved the quality of two-dimensional image reconstructions in actual holographic displays. Closest to our work, Chakravarthula et al. [11] show that naïve foveation using an importance map can help to combat speckle in two-dimensional CGH image reconstruction. Their technique does not guarantee peripheral visuals that are indistinguishable from a target image, however.

Meanwhile, three-dimensional image reconstructions in holographic displays have also seen dramatic improvement. The work by Maimone et al. [31] represents each point in a three-dimensional scene by adding a relevant sub-hologram to a final hologram. Their work exploits separable functions and introduces superior image quality with their fast pointwise method. Meanwhile, alternative approaches [43] that treat a three-dimensional scene as a multi-view image stack also prove themselves in terms of image quality. However, all of these approaches are computationally expensive, do not yet run at interactive rates and require a high level of sophistication in data representations. The work by Shi et al. [44] shows that pointwise approaches can potentially run at interactive rates with a learning methodology that cleverly stitches occluded sub-holograms to a final hologram. Thus, their work paves the way towards three-dimensional CGH at interactive rates in the future.

Our work deals with two-dimensional image reconstructions in CGH. However, unlike existing work, our holograms actively change the depth plane with a user’s focus and use state-of-the-art perceptual graphics to maintain the highest visual quality possible. In addition, our work does not require sophisticated data representations (we operate on images rather than three-dimensional data). To our knowledge this is a unique combination, and we believe it can pave a path to practical CGH.

3 METHOD

Our simulation pipeline is composed of two primary blocks, a holographic display model (Sec. 3.1) and a perceptual model (Sec. 3.2) which is used to define our metameric loss. Both blocks are differentiable and have no tunable or learned components. Our work

aims to optimise the input to the display (phase values, ϕ) such that the difference between the resulting percept and the percept of a reference image I_t is minimised (Fig. 3). To this end, we rely on a display model H and a model of perception P .

The display model $H(z)$ will map phase values to image intensities at a certain distance z . Note that in a varifocal display [2, 16, 40], the required focus z is assumed to be known at every frame. How light is propagated will depend on that focus distance.

The perception model P maps image intensities into a perceptual space where distances between points are perceptually uniform [18, 53]. This means that images which are perceived as similar should map to nearby points in the perceptual space, and images perceived as dissimilar should map to more distant points.

Putting both display and perception model together, we optimise

$$\operatorname{argmin}_{\phi} \mathcal{M}(H(z) \cdot \phi, I_t). \quad (1)$$

where \mathcal{M} is our metameric loss function, defined as:

$$\mathcal{M}(A, B) := |P(A) - P(B)|_2. \quad (2)$$

We provide details of our display and perceptual model in the following Sec. 3.1 and Sec. 3.2.

3.1 Display system and model

We use a combination of an actual holographic display and a commonly accepted differentiable model of that display hardware.

System: We design phase-only diffractive components that can either be represented with a static diffractive optical element [50] or a programmable phase-only SLM [42]. While amplitude modulation physically blocks light, phase-only holograms are a light-efficient form of optical beam shaping. More sophisticated versions of these components such as cascaded [32] or volumetric [24] holograms do not fall into the scope of this work. Importantly, such a system can be reliably modeled using a differentiable operator, explained next.

Model: The relation between a complex light wave with unit amplitude and phase ϕ leaving our phase-only SLM and the light wave $u(z)$ at the image plane depth z is described by the Fresnel diffraction operator H as $u(z) = H(z) \cdot \phi$. For derivation of this operator see [7] and section 1 of the supplementary material for implementation details. Notably, H is a linear operator and hence differentiable and can be used to optimise holograms [58].

In our case, z varies when the user changes gaze \mathbf{g} but is fixed for any point in time where there is a unique gaze and a unique focus. The HVS perceives intensity (wave amplitude-square) of light; therefore, the perceived reconstructed image, I_p is

$$I_p = |H(z) \cdot \phi|^2. \quad (3)$$

A differentiable PyTorch implementation of this model is readily available in the odak library [3]. Using this library and work by Kavaklı et al. [26], it is also possible to optimise this H operator in a camera-in-the-loop fashion to produce the best possible output on an actual holographic display. In our pipeline, we use an H operator optimised in this way to best suit our display hardware.

3.2 Perception model

Our perceptual model maps images to a perceptual space as outlined above. This mapping forms the core of our metameric loss, which is computed by measuring the distance between two images after transforming both to this perceptual space. In this work, we strive to make our model efficient (following [53]) and also differentiable, allowing it to be used effectively for any optimisation or machine learning task requiring foveated output.

Our perceptual model is inspired by the analysis step of [53], with some alterations to make it more suitable for backpropagation

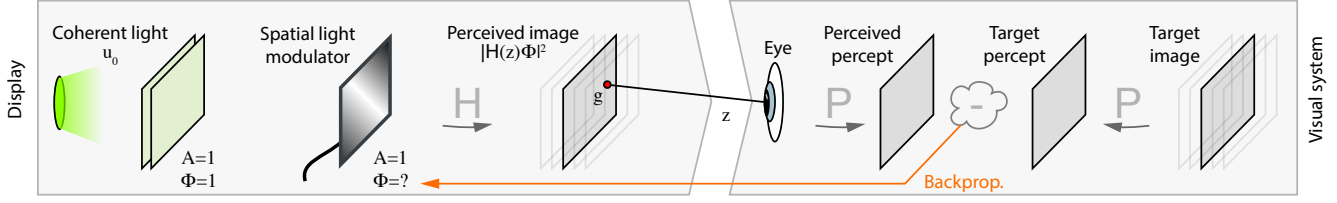


Figure 3: Overview of our system comprised of a display model (left) and a perception model (right). This scheme shows a specific fixation g and hence a specific propagation $H(z)$ for a specific focus z out of many possible changes over time as we gaze and focus.

(see supplementary material for further details). When processing colour images, as in [53] these are first converted to a YCbCr colourspace [39]. We then compute the (real-valued) steerable pyramid [38] of an input image. From each level, i of this pyramid, local statistics s_i are then computed. These s_i consist of means and variances computed over local pooling regions around each pixel which grow with eccentricity. We first describe how these local statistics are computed, then describe how the size of the pooling regions is determined.

Local means of an image can be found by convolving the image with a normalised low-pass filter F . To determine local variances, we use the identity:

$$\mathbb{V}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (4)$$

Thus, we identify the local variances by applying another lowpass filter F , squaring and subtracting: $\mathbb{V}[I] = F * (I^2) - (F * I)^2$.

The bandwidth of the lowpass filter decreases with increasing eccentricity. We set the angular pooling size to be proportional to the square of the eccentricity, as we found this gave the best results with the method of [53]. The constant of proportionality α is the parameter of the approach that controls the foveation effect's aggressiveness. In practice, to accelerate this spatially-varying lowpass filter, as in [53], we compute a MIP map of the input image that we sample using trilinear interpolation to achieve the correct pooling size at each pixel.

At each pixel, we determine the Level of Detail (LoD) value to use in this sampling so the pooling region covers the appropriate angular size. Full details of this specific process are provided in the supplementary material.

Since pooling size is proportional to eccentricity squared, some pixels near the fixation point will have pooling sizes less than or equal to one pixel. For these pixels, rather than measuring loss over all pyramid levels we calculate a direct \mathcal{L}_2 loss against the target image. In principle this does not change the loss function, but in practice we found it helped the most critical part of the image near the fovea to converge to a noise-free result more quickly.

Our loss function is differentiable and accelerated through a GPU implementation using a modern machine learning library with automatic differentiation [36]. It is straightforward to use it directly in any desired optimisation/training (a PyTorch implementation is available at [3]).

4 IMPLEMENTATION

The implementation of our proposal contains two building blocks. These are the actual hologram optimisation pipeline and a proof-of-concept holographic display prototype. We provide details of each in the next sections.

4.1 Hologram optimisation pipeline

We implement a differentiable metameric varifocal hologram optimisation pipeline using a modern machine learning library with automatic differentiation [36]. Source code of our implementation is publicly available at [54].

Target Image Capture/Generation: We use both real captured images and virtual rendered images as target images for our optimisation. In both cases, we intentionally limit the Depth of Field (DoF) of the images to mimic the optics of the human eye. We note that images formed on the retina already contain DoF blur due to the optics of the eye. Metamerisation relies upon the processing in the HVS that takes place right after the optics of the eye. We mimic this approach by applying a metameric loss to a target with limited DoF, to replicate the physical process more accurately. If the target images were sharp everywhere, our output reconstructions would be sharp in regions of the image which would have DoF blur if the real scene were viewed by the user, giving unrealistic output.

When capturing real images, we adjust the depth of field of our camera by opening the aperture to the degree that best approximates the DoF we observed when viewing the real scene (f -number 2 to 8). Virtual scenes were rendered using the Cycles raytracer in Blender [12], where we can render realistic scenes with Global Illumination (GI), and also simulating DoF blur (rendering scripts are available at [54]).

```

1 def holographic_metamer(I_t, g, z, opt, steps):
2     tp = percept(I_t, g) # See Listing 3
3
4     phi = define_initial_phase(type='random')
5     phi.requires_grad = True
6
7     for i in range(steps):
8         optimiser.zero_grad()
9         for lambda in range(lambda_r, lambda_g, lambda_b):
10            I_p = propagate(phi, lambda, z) # See Listing 2
11            pp = percept(I_p, g) # See Listing 3
12            loss = l1(pp - tp)
13            loss.backward()
14            optimiser.step()
15
16     return phi

```

Listing 1: Metameric varifocal phase-only hologram optimisation.

Main loop: Our implementation (Listing 1) follows the design of recent hologram optimisation methods [9, 11, 58]. The variable to optimise is a three-channel grid of phase values ϕ , initialised randomly. In an optimisation loop, the holographic image formation is simulated using `propagate()`, followed by a mapping into a perceptual space `percept()` (we will discuss implementation of both below). Comparing this percept to the reference percept results in a scalar loss that is back-propagated to the phase values ϕ .

All our results were produced using Stochastic Gradient Descent with ADAM [28] as the optimiser on a computer with Intel[®] i9 CPU and NVIDIA[®] RTX 2080 Ti with 200 steps. Note that our pipeline can run both on CPU and GPU. In our case, optimizing a hologram for all colour channels takes 90 seconds in total. Often, multiple holograms of a similar scene are desired. In this case, the initialisation of the next rounds of hologram optimisation with a previously calculated hologram can generally decrease this time down to 4 seconds with only five iterations. The number of itera-

tions required will be somewhat higher if the viewpoint, focus or gaze location change drastically between frames. Initializing with the previous hologram in this way also improves temporal consistency between holograms, avoiding the flicker that can result from changing between very different metamers at each frame.

Display model: We model the propagation of light from the SLM to the human eye as Fresnel diffraction [7], which is typically approximated as a convolution with a non-compact and dense kernel, best implemented using a Fourier transform. The kernel changes depending on the distance z and the wavelength λ . Hence both are taken into account by `fresnel_kernel`.

```

1 def propagate( $\phi, \lambda, z$ ):
2     for  $\lambda$  in  $\lambda_r, \lambda_g, \lambda_b$ :
3          $u_0 = \text{generate\_complex\_field}(1., \phi_\lambda)$ 
4          $H = \text{fresnel\_kernel}(z, \lambda)$ 
5          $I_{p,\lambda} = \text{ifft}(H * \text{fft}(u_0))$ 
6     return norm( $I_p$ , axis = -1)**2

```

Listing 2: Propagation from RGB phase to perceived RGB images.

Perception model: Our key contribution is a practical, efficient and differentiable mapping, `percept()`, from an image to a perceptual space that accounts for peripheral perception and is suitable for hologram optimisation (Listing 3). This method is called twice, once on the target and in each iteration on the image simulated from the current phase state. Internally, it uses a function pool for computing local means and variances over regions of varying size using convolution. This function relies on two components which we briefly describe here. Firstly `make_steerable_pyramid()` recursively and in constant time computes a steerable pyramid following the method described in [45], in our case with two orientations (vertical and horizontal) using 5×5 kernels. Secondly, `make_lod_map()` computes a map, that for every pixel holds the level at which a MIP map needs to be read to achieve the correct pooling region size for a specific eccentricity. The full details of computing this map are given in the supplemental material.

```

1 def pool(image, gaze):
2     lod_map = make_lod_map(image, gaze)
3     mipmap = make_mipmap(image)
4     return trilinear_sample(mipmap, lod_map)
5
6 def percept(I, g):
7     p = make_steerable_pyramid(I)
8     for b in p:
9         m = pool(b, g)
10        s = sqrt(pool(b*b, g) - m*m)
11        features.append(m)
12        features.append(s)
13    return features

```

Listing 3: Mapping an image and a gaze to a perceptual space.

4.2 Holographic display prototype

Having established the theoretical basis of our hologram optimisation pipeline and how to simulate outcomes from this pipeline, we tested using a physical holographic display. At the time of this manuscript, there was no commodity holographic display that we could purchase off-the-shelf. Therefore, we constructed a single colour phase-only holographic display on an optical bench (see Fig. 4).

For this purpose, we use a fibre-coupled laser diode (OSI Laser Diode, Inc - TCW RGBS-400R) with an operating wavelength of 520 nm. We collimate and polarise the laser light source with a Thorlabs LB1945-A bi-convex lens with a 200 mm focal length and Thorlabs LPVISE100-A polariser. The linearly polarised collimated beam bounces off the beamsplitter, towards our 0.93 degrees tilted

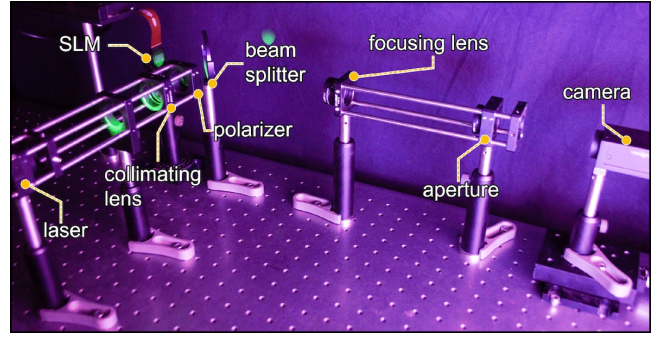


Figure 4: Proof-of-concept holographic display prototype. We use a green laser as a light source. We collimate the light from that green laser, and illuminate a phase-only SLM. The pattern displayed on the SLM modulates the phase of the light. We observe images reconstructed at various depths using a focusing lens, an aperture and a bare image sensor.

phase-only SLM, Holoeye Pluto 2.0 (tilted half order). To avoid undiffracted light, we add a horizontal grating to the displayed holograms on our SLM. The horizontally grating hologram, u'_0 , is

$$u'_0(x, y) = \begin{cases} e^{-j(\phi(x, y) + \pi)} & \text{if } [x] \text{ is even} \\ e^{-j\phi(x, y)} & \text{if } [x] \text{ is odd} \end{cases}, \quad (5)$$

where ϕ , the original phase of u_0 , is modified. This grating ensures that the image reconstructions are not visually affected by the undiffracted beam (0th order reflections). We capture the reconstructed images from these modulated beams using a Point Grey GS3-U3-23S6M-C USB 3.0 camera and a cascade of beam focusing lenses Thorlabs LA1908-A and LB1056-A. We also added a pinhole aperture, Thorlabs SM1D12, in between the camera and these lenses to avoid the undiffracted beam interfering with the result. In our system, the target image plane for our holograms is about 15cm away from the optical setup. We note that our current prototype display is not a complete near-eye display or projection display. However, it does allow us a way to practically test our method in a safe way, and to verify that the optical reconstructions appear correct on real hardware.

The resolution of holographic displays is only affected by the SLM resolution. Therefore, the image resolution of the holographic display can be calculated as $8 \mu\text{m}$ as lateral and 1 mm as axial spot sizes that are defined by Abbe's law [35].

The full realization of a complete near-eye display would require appropriate eye-piece optics. A standard eye-piece optic that can be used for such a system can be a lens with a focal length of 50 mm. This architecture's field of view (FOV) can be calculated as 17.5° horizontal x 9.9° vertical.

5 EVALUATION

In this section, we will compare our method to different alternatives in terms of fidelity (Sec. 5.1) and demonstrate results on our display prototype (Sec. 5.2).

5.1 Comparison

We compare different methods qualitatively, using the same iteration count on a set of natural images.

Methods: We study four methods that differ only in their loss. The first is the pixel-value \mathcal{L}_2 loss in RGB. The second and third are naïve baseline foveated losses inspired by [49], which account only for the acuity of the visual system in the periphery. We study two variants. The first loss \mathcal{B}_m convolves the target T according to the

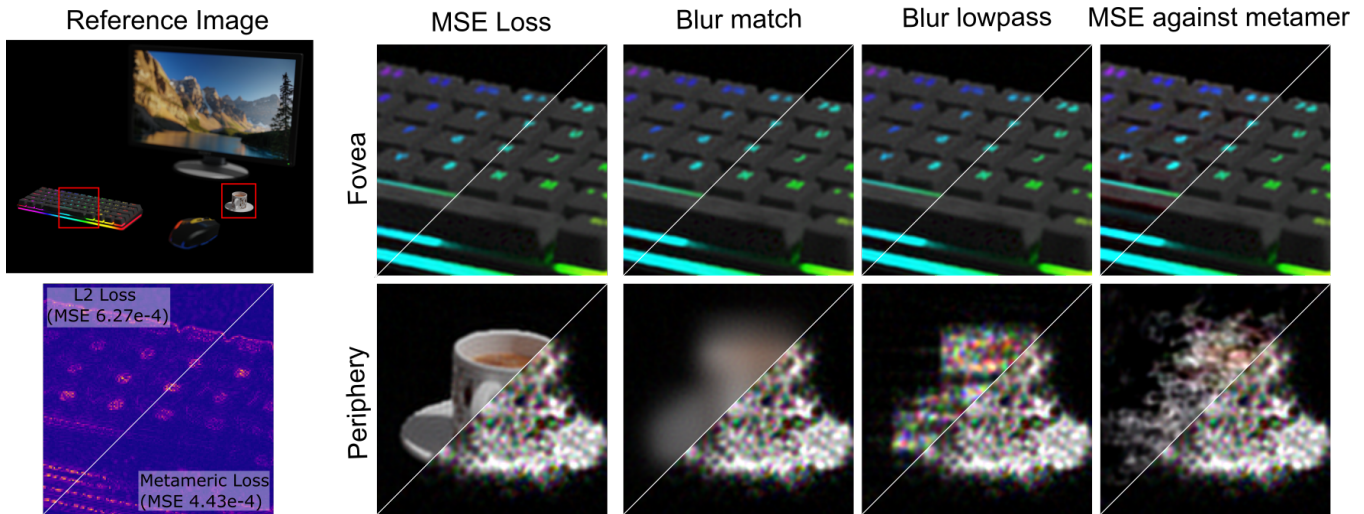


Figure 5: Comparison of different methods in a simulated environment. Bottom-left compares the output using metamer loss and MSE loss to the reference, and gives overall MSE losses for this part of the fovea. Right: output of each method in the fovea and periphery in the upper left triangle, and our method (\mathcal{M}) in the lower-right triangle for comparison. The competing methods are: reconstructing images with the correct depth of field (MSE Loss), reconstructing using MSE against a blurred image (\mathcal{B}_m) Blur match, reconstructing using MSE between blurred source and target (\mathcal{B}_1) Blur lowpass, reconstructing a target metamer image (MSE against metamer).

acuity-dependent blur kernel B and compares this to the reference I using \mathcal{L}_2 :

$$\mathcal{B}_m := \mathcal{L}_2(I, B(T)) \quad (6)$$

The second \mathcal{B}_1 blurs the result of the optimisation accordingly and compares this to the target.

$$\mathcal{B}_1 := \mathcal{L}_2(B(I), B(T)) \quad (7)$$

We note that although the definitions of these two losses are very similar, they behave in very different ways. The \mathcal{B}_m loss enforces the image I to exactly match a blurred target $B(T)$ and have no high frequency content in the periphery. In contrast \mathcal{B}_1 only enforces the low frequency content of I in the periphery to match $B(T)$, and does not constrain the higher frequencies of I away from the fovea.

The fourth method, \mathcal{M}_{L2} , is \mathcal{L}_2 loss against a metamer generated using the method of [53]. We note this is not the same as metamer loss. Our metamer loss only constrains the output to be any metamer of the target. This loss constrains the output to be identical to one particular metamer of the target.

The final method is our metamer loss \mathcal{M} (see eqn. 2).

Data: We study results on our dataset, which consists of both natural images and rendered virtual scenes. This was produced as outlined in Sec. 4.

Metric: As there is no reliable metric to capture the perceptual effects of focus and fixation, the evaluation has to rely on qualitative examples. To judge the quality, we show insets from the fovea as well as insets from the periphery. All comparisons are made after the same number of gradient descent steps (200). We note that the foveated results from the metamer and blur losses are intended to be viewed whilst fixated on the intended gaze location. As such, whilst the foveal insets can be compared directly to the reference to assess quality, the peripheral insets should be compared by fixating at a different location and keeping them in the periphery of one’s vision.

Results: Results of all five methods on a natural scene are shown in Fig. 5.

The standard \mathcal{L}_2 loss distributes error uniformly across the image. This approach would likely be the best if the user’s gaze were not known. However, if gaze information is available, it does not

prioritise the fovea in any way, meaning visible artefacts will still be uniformly distributed and presented there in the case of an actual holographic display. In this case, the fovea is noticeably blurry and some noise can be seen.

\mathcal{B}_m produces images with a similar quality in the fovea, although the periphery is naturally very blurry. As noted in [53] this is acceptable for low levels of blur. However, as blur increases, the lack of high frequencies becomes increasingly noticeable, even when fixating in the correct location. As with \mathcal{L}_2 it does not prioritise the fovea.

The \mathcal{B}_1 also has similar quality to \mathcal{L}_2 in the fovea. However, this loss does not restrict the higher frequency content in the periphery. As a result, the appearance of the periphery can vary greatly depending on the optimisation task. This application produced disturbing grid-like noise that we found to be visible even when in the periphery of vision.

The \mathcal{M}_{L2} loss approximated a metamer of the target image. However, it in no way prioritises the foveal region. Consequently, the result will not necessarily appear superior to standard \mathcal{L}_2 loss against the original target image, even when fixating in the correct location. In fact in this case the fovea appears noticeably worse than the previous approaches, possibly because the metamer generated using the approach of [53] is harder to approximate with a holographic reconstruction than the original target image.

The metamer loss \mathcal{M} , like \mathcal{B}_1 does not enforce pixel-level accuracy to the target image. Unlike \mathcal{B}_1 however, it requires that the orientation statistics in the periphery match the target. This extra constraint results in a periphery that appears similar to the metamers of [53]. The metamer loss tolerates some degree of noise, ringing or other artefacts in the periphery, making it more flexible than \mathcal{L}_2 whilst still forcing the output to be close to the target perceptually. In this case, the extra flexibility and the fact the loss prioritises the fovea have allowed it to produce a sharper result with fewer artefacts in the fovea, and perceptually correct content elsewhere.

5.2 Prototype results

Now that we have established our metamer loss (not MSE against a metamer) and hologram optimisation methods, we assess the outcome of our entire pipeline in an actual holographic display.

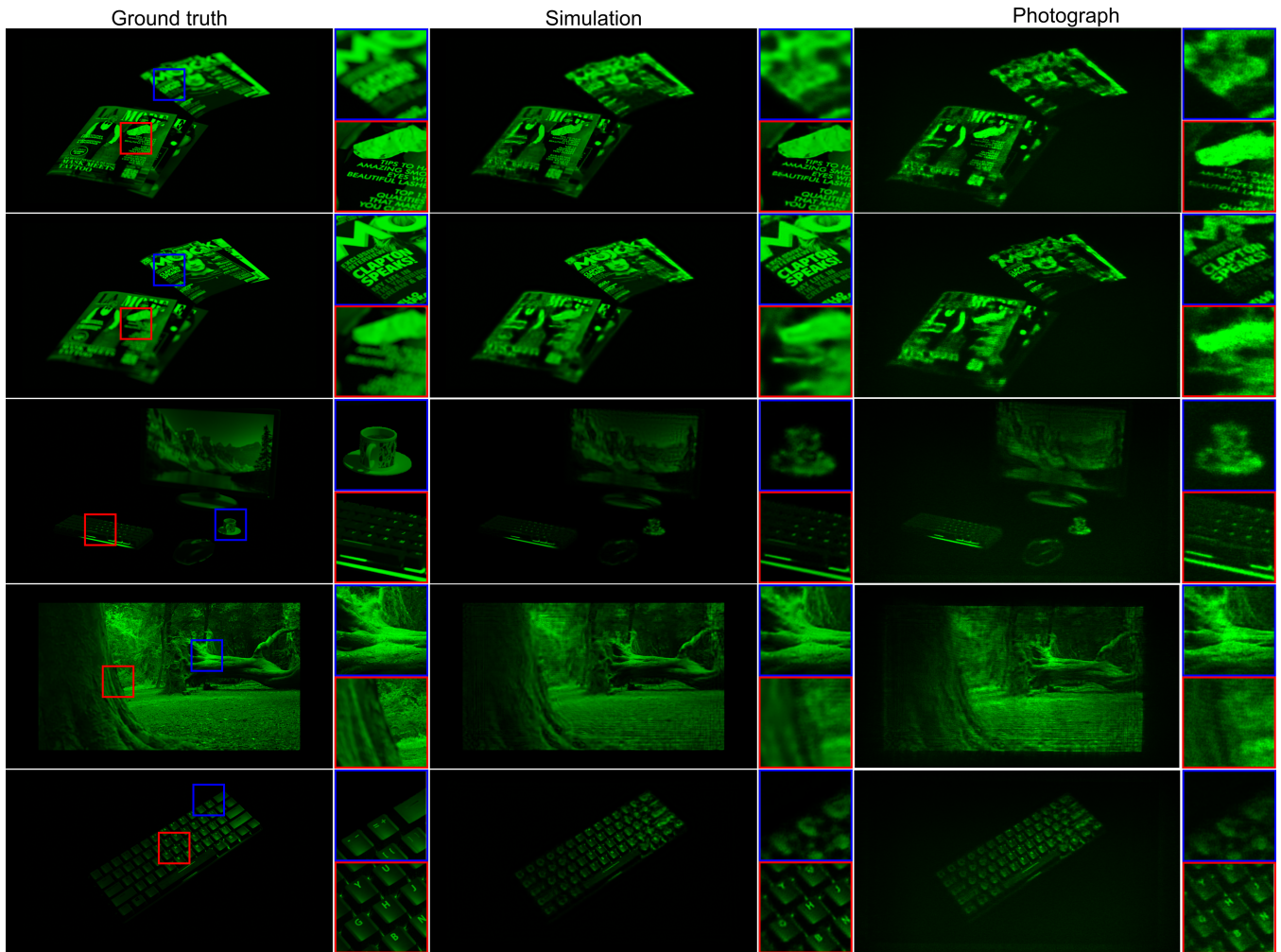


Figure 6: Metameric varifocal holographic reconstructions for holographic displays. The first column provides ground truth images that are going to be metamerised. The second column shows simulated image reconstructions of optimised holograms derived using our metameric varifocal hologram optimisation pipeline. The last column shows captured photographs of image reconstructions of our holograms when displayed in our proof-of-concept display prototype. For all columns, insets that shows zoomed-in versions of foveal and peripheral regions are also provided.

A temporal averaging approach is commonly used in holographic displays to avoid speckle noise and improve the results’ perceived quality when showing static images [13]. Rather than displaying a single, static hologram to a user, holograms optimised with slightly modified target images are rapidly displayed in sequence. These images have different realisations of the high-frequency noise and, if displayed at a rate above the user’s Critical Flicker Fusion (CFF) threshold, will be “averaged” by the HVS.

There is one issue to address when using this averaging approach with the metameric loss. We note that in general, the mean of two distinct metamers to an image is not a metamer. Taking the mean of multiple, very different metamers to an image tends to produce a blurry result in the periphery. For this reason, we need to ensure that we use the same metamer for a particular target image in our temporal averaging approach, just adding different frequencies of high-frequency holographic noise in each case. We achieve this by initializing with a noise phase image to create the first hologram, and for the subsequent holograms we initialise with the previous optimised phase, optimising for just 5 iterations with a lower learning rate. This has the added benefit of reducing the number of iterations needed to generate the subsequent holograms.

Fig. 7 shows an example. Here we show a single hologram and

the average of five consecutive holograms with a different noise, simulating the image perceived by the user. Our current prototype display operates at 60Hz, meaning that with 5-frame temporal averaging the effective framerate is reduced to 12Hz. Note, however, that the noise is greatly reduced and the image closer to the ground truth on the right. One of the goals in CGH is to avoid using multiple holograms so that holographic displays can use the total frame rate of an SLM. Our metameric loss improves the precision of a single hologram, as later shown in Fig. 6, helping towards meeting the goals of CGH.

We provide sample image reconstructions of our holograms optimised using our metameric varifocal hologram generation pipeline as in Fig. 6. These samples contain a set of ground truth images that are targeted to be metamerised. Gaze locations and focal planes ($\sim 0.15m$) are also provided for each sample during an optimisation session. We intentionally choose images varying from synthetic sparsely populated cases to densely populated natural photographs. In the figure, the results of our simulated image reconstructions are provided for the holograms optimised using our pipeline. Finally, we also provide photographs of image reconstructions from our proof-of-concept display. Though it is not pixel-perfect, note that simulated cases and actual photographs very closely resemble each



Figure 7: Simulated Temporal averaging. Rather than using a single hologram based image reconstruction (images on the left column), we rely on showing multiple images (center column). This way, we can generate images that suffer less from noise. Note that we use five holograms for this temporal averaging example.

other, and the noise patterns at both images are close in terms of spatial distribution. How well the actual photographic results match the simulation is heavily dependent on the model used in optimisation. We relied on a state-of-the-art model from the literature to achieve the best results possible at the present time. As models continue to improve, our pipeline can take advantage of improvements to combat noise in CGH in the future.

Focal planes used in our captures varied between 0.14m to 0.16m. These distances are suitable for virtual reality and augmented reality applications. For example, in a typical virtual reality display, a magnifier glass typically of focal length from 35 to 60mm is used between a flat image modulator and the eye. An image volume in the order of a few millimetres translates to covering virtual image distances from very close distances (10cm) to far away (6m and beyond). Thus, our CGH pipeline can present images at a wide range of focal distances.

6 DISCUSSION

Both the simulations and the results from our physical prototype further our understanding of the technical challenges in deriving a metameric varifocal hologram pipeline. These efforts also lay the foundation of future holographic displays that are perceptually guided. Nevertheless, more work remains before we can achieve fully practical CGH.

Interactive rates: Our current CGH optimisation does not run at interactive rates, as is typical for most CGH optimisation pipelines. In the meantime, there has been a push towards taking advantage of learning approaches in CGH in recent years [10, 37, 44]. While learning-based acceleration is outside of the scope of our work, our metameric loss appears eminently suited for learning frameworks.

Metameric loss: Our readers may argue that setting a metamer or a full resolution image as a target image and optimizing a hologram using MSE loss can lead to a faster optimisation routine. This method is indeed computationally efficient but, as shown in Sec. 5 above, produces results with similar or worse quality than regular \mathcal{L}_2 loss in the fovea. Also, in an actual or simulated holographic display, perfectly matching target images may be physically impossible, and visual imperfections such as noise are very likely to occur, as can be observed in Fig. 7. Rather than using MSE loss against a metamer target or a complete resolution target, optimizing holograms with our metameric loss guarantees metamer solutions that play well with a given holographic light transport model in a

simulation. This advantage clearly shows in our simulated results; however, at the time being, the advantages of our method and other foveation methods for CGH [11] do not always translate to the real world due to the lack of accurate-enough holographic light transport and display models in the literature – which is beyond the scope of our focus of this work. We argue that the importance of metameric loss will become evident as the holographic light transport models match physical hardware accurately in the future.

Varifocal visuals: There has been a long-standing debate about the qualities of a display when it comes to supporting optical depth cues. A recent survey on displays [29] captures a detailed background of this debate. The varifocal approach in displays has recently proven to improve visual comfort [25]. Fortunately, the latency requirements for a varifocal display are not demanding as the accommodation duration of human eyes has been measured as between 500 ms and 1 s in various studies [6, 8]. One disadvantage of varifocal displays is that they cannot handle rare cases where gaze does not uniquely determine the focal depth. Nevertheless, we argue that varifocal displays are strong candidates for supporting optical depth cues in displays.

Gaze-contingent displays: A complete gaze-contingent display requires eye-gaze tracking systems to be involved in the process of generating visuals. The accuracy requirements of a gaze-contingent varifocal display have recently been systematically studied [15]. Unfortunately, our proof-of-concept prototype is not equipped with eye-gaze tracking hardware [4, 30]. Our work assumes that such eye-gaze information is readily available. We calculate holograms with this assumption.

Real-world gaze trackers suffer from varying degrees of inaccuracy and latency. Our method could be adapted to better tolerate these issues, by modifying the foveation and increasing the size of the foveal region. There is however a trade-off between foveal size and improvement in image quality over MSE loss.

Though some hurdles remain in our implementation, our work resembles a reliable blueprint for a practical CGH pipeline that can deliver perceptually accurate visuals to users.

7 CONCLUSION

The versatility in generating high-resolution and three-dimensional visuals makes Computer-Generated Holography a powerful technique suitable for next-generation displays. However, among many technical issues, achieving three-dimensional visuals with CGH still poses a significant challenge in real holographic displays. We argue that gaze-contingent CGH can be key to achieving practical holographic displays with perceptually accurate three-dimensional visuals. For this purpose, we build upon state-of-the-art perceptual graphics. We formulate a new differentiable hologram optimisation pipeline with a perceptually guided loss function. Rather than reconstructing imperfect three-dimensional scenes, our CGH method reconstructs visuals at the user’s focus. It offers improved image quality at the fovea, while displaying true metamers of target images in the periphery. Using gaze-contingency, we formulate our phase optimisation as a two-dimensional problem, removing the need to match a light field or multiplane image. In this way, our method paves the way towards a practical display that provides perceptually accurate three-dimensional visuals more efficiently.

ACKNOWLEDGMENTS

The authors thank the reviewers for their useful feedback. The authors also thank Duygu Ceylan for the fruitful and inspiring discussions improving the outcome of this research, and Selim Ölçer for helping with the fibre alignment of laser light source in the proof-of-concept display prototype. This work was partially funded by the EPSRC/UKRI project EP/T01346X/1 and Royal Society’s RGS\R2\212229 - Research Grants 2021 Round 2.

REFERENCES

- [1] K. Akşit, P. Chakravarthula, K. Rathinavel, Y. Jeong, R. Albert, H. Fuchs, and D. Luebke. Manufacturing application-driven foveated near-eye displays. *IEEE TVCG*, 25(5):1928–1939, 2019.
- [2] K. Akşit, W. Lopes, J. Kim, P. Shirley, and D. Luebke. Near-eye varifocal augmented reality display using see-through screens. *ACM Trans Graph*, 36(6):1–13, 2017.
- [3] K. Akşit, A. S. Karadeniz, P. Chakravarthula, W. Yujie, K. Kavakli, Y. Itoh, and D. R. Walton. kunguz/odak: Odak 0.1.9. <https://doi.org/10.5281/zenodo.5526684>, Sept. 2021. doi: 10.5281/zenodo.5526684
- [4] A. N. Angelopoulos, J. N. Martel, A. P. Kohli, J. Conrath, and G. Wetzstein. Event based, near eye gaze tracking beyond 10,000 hz. *arXiv preprint arXiv:2004.03577*, 2020.
- [5] H. Aubert and R. Förster. Beiträge zur kenntniss des indirecten sehens.(i). untersuchungen über den raumsinn der retina. *Archiv für Ophthalmologie*, 3(2):1–37, 1857.
- [6] S. R. Bharadwaj and C. M. Schor. Acceleration characteristics of human ocular accommodation. *Vis Res*, 45(1):17–28, 2005.
- [7] M. Born and E. Wolf. *Principles of optics: electromagnetic theory of propagation, interference and diffraction of light*. 2013.
- [8] F. Campbell and G. Westheimer. Dynamics of accommodation responses of the human eye. *J Physiology*, 151(2):285–295, 1960.
- [9] P. Chakravarthula, Y. Peng, J. Kollin, H. Fuchs, and F. Heide. Wirtinger holography for near-eye displays. *ACM Trans Graph*, 38(6):1–13, 2019.
- [10] P. Chakravarthula, E. Tseng, T. Srivastava, H. Fuchs, and F. Heide. Learned hardware-in-the-loop phase retrieval for holographic near-eye displays. *ACM Trans Graph*, 39(6):1–18, 2020.
- [11] P. Chakravarthula, Z. Zhang, O. Tursun, P. Didyk, Q. Sun, and H. Fuchs. Gaze-contingent retinal speckle suppression for perceptually-matched foveated holographic displays. *IEEE TVCG*, 2021.
- [12] B. O. Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2021.
- [13] V. R. Curtis, N. W. Caira, J. Xu, A. G. Sata, and N. C. Pégard. Dcgh: Dynamic computer generated holography for speckle-free, high fidelity 3d displays. In *2021 IEEE Virtual Reality and 3D User Interfaces (VR)*, pp. 1–9. IEEE, 2021.
- [14] A. Deza, A. Jonnalagadda, and M. P. Eckstein. Towards metamerism via foveated style transfer. In *ICLR*, 2019.
- [15] D. Dunn. Required accuracy of gaze tracking for varifocal displays. In *VRST*, pp. 1838–1842, 2019.
- [16] D. Dunn, C. Tippets, K. Torell, P. Kellnhofer, K. Akşit, P. Didyk, K. Myszkowski, D. Luebke, and H. Fuchs. Wide field of view varifocal near-eye display using see-through deformable membrane mirrors. *IEEE TVCG*, 23(4):1322–1331, 2017.
- [17] D. B. Elliott, K. Yang, and D. Whitaker. Visual acuity changes throughout adulthood in normal, healthy eyes: seeing beyond 6/6. *Optometry and vision science*, 72(3):186–191, 1995.
- [18] J. Freeman and E. P. Simoncelli. Metamers of the ventral stream. *Nature neuroscience*, 14(9):1195–1201, 2011.
- [19] S. Friston, T. Ritschel, and A. Steed. Perceptual rasterization for head-mounted display image synthesis. *ACM Trans Graph*, 38(4):1–14, 2019.
- [20] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *CVPR*, pp. 2414–2423, 2016.
- [21] M. Gong, Y. Xuan, L. J. Smart, and L. A. Olzak. The extraction of natural scene gist in visual crowding. *Scientific Reports*, 8(1):1–13, 2018.
- [22] C.-K. Hsueh and A. A. Sawchuk. Computer-generated double-phase holograms. *Applied optics*, 17(24):3874–3883, 1978.
- [23] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pp. 1501–1510, 2017.
- [24] C. Jang, O. Mercier, K. Bang, G. Li, Y. Zhao, and D. Lanman. Design and fabrication of freeform holographic optical elements. *ACM Trans Graph*, 39(6):1–15, 2020.
- [25] P. V. Johnson, J. A. Parnell, J. Kim, C. D. Saunter, G. D. Love, and M. S. Banks. Dynamic lens and monovision 3d displays to improve viewer comfort. *Optics express*, 24(11):11808–11827, 2016.
- [26] K. Kavaklı, H. Urey, and K. Akşit. Learned holographic light transport. *arXiv preprint arXiv:2108.08253*, 2021.
- [27] J. Kim, Y. Jeong, M. Stengel, K. Akşit, R. Albert, B. Boudaoud, T. Greer, J. Kim, W. Lopes, Z. Majercik, et al. Foveated ar: dynamically-foveated augmented reality display. *ACM Trans Graph*, 38(4):1–15, 2019.
- [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [29] G. A. Koulieris, K. Akşit, M. Stengel, R. K. Mantiuk, K. Mania, and C. Richardt. Near-eye display and tracking technologies for virtual and augmented reality. In *Comp Graph Forum*, vol. 38, pp. 493–519, 2019.
- [30] R. Li, E. Whitmire, M. Stengel, B. Boudaoud, J. Kautz, D. Luebke, S. Patel, and K. Akşit. Optical gaze tracking with spatially-sparse single-pixel detectors. In *ISMAR*, pp. 117–126, 2020.
- [31] A. Maimone, A. Georgiou, and J. S. Kollin. Holographic near-eye displays for virtual and augmented reality. *ACM Trans Graph*, 36(4):1–16, 2017.
- [32] D. Mengü, Y. Luo, Y. Rivenson, and A. Ozcan. Analysis of diffractive optical neural networks and their integration with electronic neural networks. *IEEE Journal of Selected Topics in Quantum Electronics*, 26(1):1–14, 2019.
- [33] J. Orlosky, M. Sra, K. Bektaş, H. Peng, J. Kim, N. Kos’ myna, T. Hollerer, A. Steed, K. Kiyokawa, and K. Akşit. Telelife: The future of remote living. *arXiv preprint arXiv:2107.02965*, 2021.
- [34] N. Padmanaban, R. Konrad, T. Stramer, E. A. Cooper, and G. Wetzstein. Optimizing virtual reality for all users through gaze-contingent and adaptive focus displays. *PNAS*, 114(9):2183–2188, 2017.
- [35] N. Padmanaban, Y. Peng, and G. Wetzstein. Holographic near-eye displays based on overlap-add stereograms. *ACM Trans. Graph.*, 38(6), nov 2019. doi: 10.1145/3355089.3356517
- [36] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.
- [37] Y. Peng, S. Choi, N. Padmanaban, and G. Wetzstein. Neural holography with camera-in-the-loop training. *ACM Trans Graph*, 39(6):1–14, 2020.
- [38] J. Portilla and E. P. Simoncelli. A parametric texture model based on joint statistics of complex wavelet coefficients. *J Computer Vision*, 40(1):49–71, 2000.
- [39] C. Poynton. *Digital Video and HD: Algorithms and Interfaces*. 2012.
- [40] K. Rathinavel, G. Wetzstein, and H. Fuchs. Varifocal occlusion-capable optical see-through augmented reality display based on focus-tunable optics. *IEEE TVCG*, 25(11):3125–3134, 2019.
- [41] R. Rosenholtz. Capabilities and limitations of peripheral vision. *Annual review of vision science*, 2:437, 2016.
- [42] N. Savage. Digital spatial light modulators. *Nature Photonics*, 3(3):170–172, 2009.

- [43] L. Shi, F.-C. Huang, W. Lopes, W. Matusik, and D. Luebke. Near-eye light field holographic rendering with spherical waves for wide field of view interactive 3d computer graphics. *ACM Trans Graph*, 36(6):1–17, 2017.
- [44] L. Shi, B. Li, C. Kim, P. Kellnhofer, and W. Matusik. Towards real-time photorealistic 3d holography with deep neural networks. *Nature*, 591(7849):234–239, 2021.
- [45] E. P. Simoncelli and W. T. Freeman. The steerable pyramid: A flexible architecture for multi-scale derivative computation. In *ICIP*, vol. 3, pp. 444–447, 1995.
- [46] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [47] C. Slinger, C. Cameron, and M. Stanley. Computer-generated holography as a generic display technology. *Computer*, 38(8):46–53, 2005.
- [48] J. Spjut, B. Boudaoud, J. Kim, T. Greer, R. Albert, M. Stengel, K. Akşit, and D. Luebke. Toward standardized classification of foveated displays. *IEEE TVCG*, 26(5):2126–2134, 2020.
- [49] L. Surace, M. Wernikowski, O. Tursun, K. Myszkowski, R. Mantiuk, and P. Didyk. Learning foveated reconstruction to preserve perceived image statistics, 2021.
- [50] G. J. Swanson. Binary optics technology: the theory and design of multi-level diffractive optical elements. Technical report, 1989.
- [51] G. Tan, Y.-H. Lee, T. Zhan, J. Yang, S. Liu, D. Zhao, and S.-T. Wu. Foveated imaging for near-eye displays. *Optics express*, 26(19):25076–25085, 2018.
- [52] T. S. A. Wallis, C. M. Funke, A. S. Ecker, L. A. Gatys, F. A. Wichmann, and M. Bethge. Image content is more important than Bouma’s Law for scene metamers. *bioRxiv*, 2018.
- [53] D. R. Walton, R. K. Dos Anjos, S. Friston, D. Swapp, K. Akşit, A. Steed, and T. Ritschel. Beyond blur: Ventral metamers for foveated rendering. *ACM Trans. Graph.*, 40(4), 2021.
- [54] D. R. Walton, K. Kavakli, R. K. D. Anjos, D. Swapp, T. Weyrich, H. Urey, A. Steed, T. Ritschel, and K. Akşit. Metameric holography repository. https://github.com/complight/metameric_holography.
- [55] B. Wang and K. J. Ciuffreda. Depth-of-focus of the human eye in the near retinal periphery. *Vis Res*, 44(11):1115–1125, 2004.
- [56] L. Wang, M. Hajiesmaili, and R. K. Sitaraman. Focas: Practical video super resolution using foveated rendering. In *ACM Multimedia*, 2021.
- [57] G.-z. Yang, B.-z. Dong, B.-y. Gu, J.-y. Zhuang, and O. K. Ersoy. Gerchberg–saxton and yang–gu algorithms for phase retrieval in a nonunitary transform system: a comparison. *Applied optics*, 33(2):209–218, 1994.
- [58] J. Zhang, N. Pégard, J. Zhong, H. Adesnik, and L. Waller. 3d computer-generated holography by non-convex optimization. *Optica*, 4(10):1306–1313, 2017.
- [59] Y. Zhang, M. A. Noack, P. Vagovic, K. Fezzaa, F. Garcia-Moreno, T. Ritschel, and P. Villanueva-Perez. Phasegan: a deep-learning phase-retrieval approach for unpaired datasets. *Optics Express*, 29(13):19593–19604, 2021.
- [60] C. M. Ziemba and E. P. Simoncelli. Opposing effects of selectivity and invariance in peripheral vision. *Nature Communications*, 12(1):1–11, 2021.