

Configurable Holography: Towards Display and Scene Adaptation

YICHENG ZHAN, University College London, UK

LIANG SHI, Massachusetts Institute of Technology, USA

WOJCIECH MATUSIK, Massachusetts Institute of Technology, USA

QI SUN, New York University, USA

KAAN AKŞIT, University College London, UK

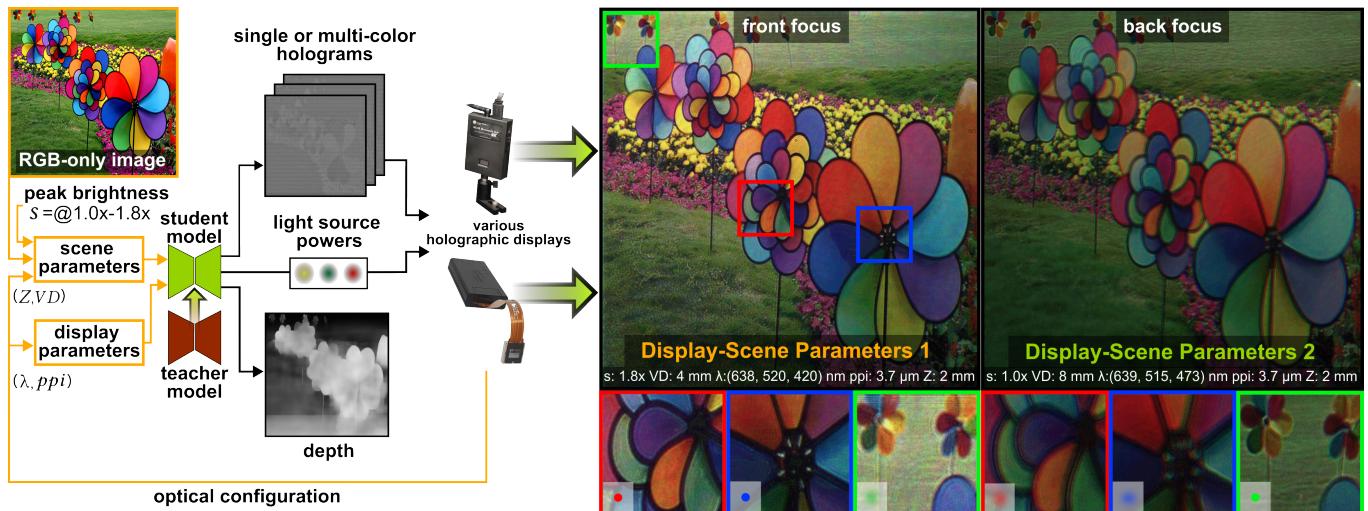


Fig. 1. Our display-scene parameter adaptive learned holography model can generate Three-Dimensional (3D) holograms for various display-scene parameters without requiring retraining. These display-scene parameters include working wavelengths (λ), pixel pitch (ppi), propagation distance (Z), volume depth (VD), and peak brightness (s). Our model also predicts 3D holograms from given RGB-only 2D images, eliminating the need for depth inputs. With the help of a knowledge-distillation technique, our student model estimates holograms with 2x speed up w.r.t. literature [Shi et al. 2022]. We verify our findings experimentally using two different holographic display prototypes and various display-scene parameters (Image Source: [Spragg 2009]).

Emerging learned holography approaches have enabled faster and high-quality hologram synthesis, setting a new milestone toward practical holographic displays. However, these learned models require training a dedicated model for each set of display-scene parameters. To address this shortcoming, our work introduces a highly configurable learned model structure, synthesizing 3D holograms interactively while supporting diverse display-scene parameters. Our family of models relying on this structure can be conditioned continuously for varying novel scene parameters including input images, propagation distances, volume depths, peak brightnesses, and novel display parameters of pixel pitches and wavelengths. Uniquely, our findings

unearth a correlation between depth estimation and hologram synthesis tasks in the learning domain, leading to a learned model unlocking accurate 3D hologram generation from 2D images across varied display-scene parameters. We validate our models by synthesizing high-quality 3D holograms in simulations and also verify our findings with two different holographic display prototypes. Moreover, our family of models can synthesize holograms with a 2x speed-up compared to the state-of-the-art learned holography approaches in the literature.

CCS Concepts: • Computing methodologies → Mixed / augmented reality; 3D imaging; Computer graphics; Rendering.

Additional Key Words and Phrases: deep learning, neural networks, depth estimation, hard-parameter sharing, computer-generated holography

ACM Reference Format:

Yicheng Zhan, Liang Shi, Wojciech Matusik, Qi Sun, and Kaan Akşit. 2025. Configurable Holography: Towards Display and Scene Adaptation. 1, 1 (April 2025), 11 pages. <https://doi.org/XXXXXXX.XXXXXXX>

Authors' addresses: Yicheng Zhan, UCABY83@ucl.ac.uk, University College London, Gower Street, London, UK, WC1E 6BT; Liang Shi, Massachusetts Institute of Technology, MA 02139, Massachusetts, USA, liangs@mit.edu; Wojciech Matusik, Massachusetts Institute of Technology, MA 02139, Massachusetts, USA, wojciech@mit.edu; Qi Sun, New York University, , New York, USA, qisun@nyu.edu; Kaan Akşit, University College London, Gower Street, London, UK, k.aksit@ucl.ac.uk.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Association for Computing Machinery.

XXXX-XXXX/2025/4-ART \$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Holographic displays [Kim et al. 2024] support optical focus cues and multiple perspectives, promising authentic immersive 3D visual experiences as a potential future display technology. The demanding computational requirements of rendering visuals for these displays [Blinder et al. 2019] necessitate the emerging research field of

	Input	PA	Speed	Hologram			SLM Refresh Rate	Stage	3D	Depth Accuracy	Learned	Max Z (mm)	VD (mm)
				Type	PD	DP							
Our Method	RGB-only	Yes	Fast	M+S	8	No	60 Hz	Single	True	Moderate	Full	~10.0	~8.0
HoloBeam[Akşit and Itoh 2023]	G-only	No	Fast	S	8	Yes	60 Hz	Single	True	Inaccurate	Full	~0.0	~6.0
Multi-DNN [Ishii et al. 2023]	RGB-only	No	Slow	S	8	Yes	60 Hz	Three	True	Moderate	Full	~50.0	~2.0
NH [Peng et al. 2020]	RGB-only	No	Fast	S	8	No	60 Hz	Two	False	Accurate	Full	~100.0	~0.0
NH3D [Choi et al. 2021a]	RGB-D	No	Slow	S	8	No	60 Hz	Two	True	Accurate	Semi	~8.2	~4.4
TensorV2 [Shi et al. 2022]	RGB-D	No	Fast	S	8	Yes	60 Hz	Two	True	Accurate	Full	~12.0	~6.0
DGE-CNN [Liu et al. 2023a]	RGB-D	No	Slow	S	8	No	58 Hz	Two	True	Moderate	Semi	~10.0	~30.0
4K-DMDNet [Liu et al. 2023b]	RGB-D	No	Slow	S	8	No	58 Hz	Two	False	Accurate	Full	~300.0	~0.0
Time-multiplexed [Choi et al. 2022]	RGB-D	No	Slow	S	4	No	480 Hz	Two	True	Accurate	Semi	~79.0	~12.0

Table 1. Comparison of hologram synthesizing methods. Our method can generate both single (S) and multi-color holograms (M) from an RGB-only input for a given preferred set of display-scene parameters. In *Input* column, G-only refers to green-channel-only. In *Speed* column, Fast refers to > 10 FPS, and Slow refers to < 10 FPS, at 1920 × 1080 resolution. In *Learned* column, Full refers to the full learning-based method, and Semi refers to a combination of learning and optimization. NH [Peng et al. 2020] is fully learning-based. Conversely, NH3D [Choi et al. 2021b] is a hybrid, using learning only for the refinement of holograms. The label PA, PD, DP, Z, VD, and SLM refers to parameter adaptive, pixel depth, double phase encoding [Hsueh and Sawchuk 1978], maximum propagation distance, volume depth of 3D scenes, and Spatial Light Modulator, respectively.

learned holography methods [Shi et al. 2022] in recent years. However, existing learned holography models require training a dedicated model for each set of display-scene parameters, leading to challenges in model management and generalization.

Our work focuses on synthesizing 3D holographic images at interactive rates for a wide variety of display-scene parameters in holographic displays without requiring training a dedicated model for each configuration. To achieve this goal, we introduce a highly configurable model structure to efficiently synthesize conventional single-color [Kavaklı et al. 2023a] and emerging multi-color [Kavaklı et al. 2023b; Markley et al. 2023] 3D holograms. We derived a family of models from this structure and conditioned them for a variety of display-scene parameters including working wavelengths, propagation distance, volume depth, pixel pitch, and peak brightness. Our work also unearths the correlation between depth estimation and 3D hologram synthesis tasks in learned methods, a new contribution to the relevant literature. Although estimating depth is not our aim in this new model, we utilize depth estimation task in parallel to hologram synthesis task to improve 3D accuracy of our predicted holograms when working with RGB-only inputs without depth, the most common media. This way, our family of models can also support a variant with RGB-only inputs, the most common media type, without requiring depth in the input. Furthermore, we distill our model using student-teacher learning strategy [Wang and Yoon 2021] to achieve hologram synthesis at interactive rates. Our contributions are as follows:

- Display-scene parameter adaptive Learned Holography. We introduce a new learned holography model structure, from which we derive a family of models, supporting a range of display-scene parameters, single and multi-color holograms. With various training examples, we show that our family of models can continuously support propagation distance, volume depth, pixel pitch, peak brightness, and wavelength.
- Accurate 3D Hologram Synthesis from RGB-only inputs. We unearth the correlation between depth estimation and 3D hologram synthesis tasks in learned methods. Our learned model leverages

this depth-hologram correlation and adopts multitask learning with hard-parameter sharing [Caruana 1993] to convert RGB-only Two-Dimensional (2D) images to accurate 3D holograms. Applications where no depth information is provided, such as viewing photos and videos or live streaming, can benefit from our model as depth is no longer required in hologram synthesis.

- Advancing 3D Holograms towards Interactive Rates. We distill our model using a student-teacher learning strategy, achieving up to 2x faster hologram synthesis compared with the literature [Shi et al. 2022] under 32-bit precision (fp32).

We evaluate the image quality and speed of our model by comparing it with existing learned holography methods and their potential future derivatives. We verify our findings experimentally using two holographic display prototypes, demonstrating comparable image quality with the state-of-the-art (SOTA). Our model is designed to cater to conventional holographic displays, sharing their common FoV and Eyebox constraints [Shi et al. 2017]. Further investigations are needed to support emerging holographic display architectures [Chae et al. 2023; Jang et al. 2024; Kuo et al. 2023] and novel hologram types [Choi et al. 2022; Kim et al. 2024]. Our model and weights are publicly available at [REVIEW].

2 Related Work

Our work represents a family of learned 3D Computer-Generated Holography (CGH) methods that support working with RGB-only and RGB-D inputs while leveraging the concepts of multi-task learning and knowledge distillation. Tbl. 1 provides a summary comparison of our method against the relevant literature.

Learned Computer-Generated Holography. Early works on conventional hologram (single-color) synthesis utilize U-Net-based Convolutional Neural Network (CNN) structures [Siddique et al. 2021]. A recent study by Shi et al. [2022] proposes a stack of residual blocks without resolution reduction to speed up the synthesis process for single-color holograms. Both of these works [Shi et al. 2022; Siddique et al. 2021] require RGB-D inputs, supporting a specific set of display-scene parameters. Liu et al. [2023a] propose a two-stage

solution for removing the dependency on depth information in the input. The first stage involves depth estimation, while the second stage optimizes single-color holograms using RGB-D inputs, leading to a slow hologram synthesis process. Akşit and Itoh [2023] merge these two-stages solutions [Liu et al. 2023a] into a single stage CNN to accelerate the process. However, their synthesized single-color holograms provide inaccurate and limited depth in their reconstructed images. Ishii et al. [2023a] propose a triple-stage solution where the first estimates depth, the second generates holograms using RGB-D input and the third optimizes the holograms' quality, leading to a slower process than both works [Akşit and Itoh 2023; Liu et al. 2023a]. *Our work offers the first learned CGH model to support single-color and multi-color hologram types, RGB-only inputs, and various display-scene parameters at interactive rates.*

Multi-task Learning. The work by Zhang and Yang [2021] introduces a comprehensive review of Multi-task Learning (MTL). Our work leverages MTL [Caruana 1993] to solve depth inaccuracies in previous works with RGB-only inputs [Akşit and Itoh 2023], jointly learning input depths, holograms, and light source intensities using shared weights and biases. A common strategy in MTL is to share layers among all tasks (hard-parameter sharing) while keeping output layers or several input layers (heads) task-specific [Ruder et al. 2019; Sarwar et al. 2019]. Similarly, our work follows hard-parameter sharing while keeping only output heads task-specific. This approach is separate from soft-parameter sharing [Duong et al. 2015], where each task has an independent CNN, while other tasks constrain the weights of each CNN during training. *Uniquely, our work identifies that jointly learning subtasks of Monocular Depth Estimation (MDE) and hologram synthesis can help generate depth-accurate holograms when working with RGB-only inputs.*

Knowledge Distillation. Knowledge Distillation (KD) [Hinton et al. 2015] transfers knowledge from large models to introduce compact and lightweight models for various applications. As stated in a survey by Gou et al. [2021], KD methods follow the teacher-student architecture. A teacher model can transfer knowledge to a student model either by regularizing the output responses (response-based) [Micaelli and Storkey 2019] or layer responses (feature-based) [Chen et al. 2020] between the two models. In the past, KD has also been applied in MTL [Ignatov et al. 2021; Li and Bilen 2020]. *Our work adopts these findings to holography first time and transfers knowledge from a teacher model based on MTL to a student model through pixel-wise distillation (response-based) to speed up hologram synthesis.*

3 Display-scene Parameter Adaptive Learned Holography

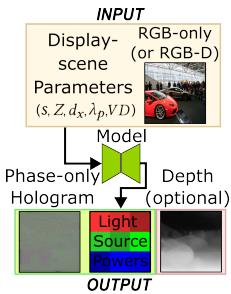


Fig. 2. Model overview.
(Source: [C 2014])

depth estimation task as a beneficial parallel task to help generate accurate 3D holograms from 2D images. Thus, our model does not aim to compete and advance the MDE task alone.

3.1 Synthesizing 3D Holograms

Holographic displays play successive holograms rapidly to generate full-color images. The Human Visual System (HVS) helps to fuse the light from these illuminated holograms into full-color images. Suppose only a single-color primary illuminates a specific hologram, which is calculated for the particular wavelength of that single-color primary. In that case, this hologram type is called a single-color hologram. In contrast, if multiple color primaries illuminate a hologram, and if that hologram is calculated for all the wavelengths provided by these multiple color primaries, then this hologram is a multi-color hologram. Regardless of the type, hologram synthesis could be illustrated by Fig. 3 and modeled by the following optimization,

$$Z_0 = Z - \frac{VD}{2}, Z_n = Z + \frac{VD}{2},$$

$$I_r(p, t, z) = \left| l_{(p,t)} e^{i \frac{\lambda_p}{\lambda_{p,\text{anchor}}} u_t} * h_p(\lambda_p, z, dx) \right|^2, \quad (1)$$

$$\hat{u}_t, \hat{l}_{(p,t)} \leftarrow \underset{u_t, l_{(p,t)}}{\operatorname{argmin}} \mathcal{L}_{\text{img}} \left(\sum_{z=Z_0}^{Z_n} \sum_{p=1}^P \sum_{t=1}^T I_r(p, t, z), sI_{(p,z)} \right),$$

where $Z \in \mathbb{R}$ denotes the light propagation distance, $VD \in \mathbb{R}$ denotes the volume depth, $P \in \mathbb{Z}$ denotes the total number of color primaries (i.e. typically three in existing display technology), $p \in \mathbb{Z}$ denotes the index of a color primary, $T \in \mathbb{Z}$ denotes the total number of subframes reproducing a full-color image (i.e. typically three in holographic displays), $t \in \mathbb{Z}$ denotes the index of a subframe, $l_{(p,t)} \in \mathbb{R}^{P \times T}$ represents the light source power for the p -th primary at the t -th subframe, $\lambda_p \in \{400 - 700\text{nm}\}$ denotes the wavelength of the active primary, $\lambda_{p,\text{anchor}}$ denotes the wavelength of the anchor primary, for which the nominal value of the SLM phase is calibrated against, $u_t \in \mathbb{C}^{H \times W}$ represents the phase-only hologram (a subframe), dx represents the pixel pitch in a hologram frame, $I_{(p,z)} \in \mathbb{R}^{H \times W}$ is a target image intensity, $s \in \mathbb{R}$ represents a parameter to scale the brightness of a target image intensity, $h_p \in \mathbb{C}^{H \times W}$ denotes the wavelength, λ_p , and location, z , dependent light transport kernel [Choi et al. 2021b; Kavakli et al. 2022; Matsushima and Shimobaba 2009], and \mathcal{L}_{img} represents the visual difference between the intended target image for the given plane and the reconstructed image by a hologram. Note that l is an identity matrix in the single-color holograms, where only one

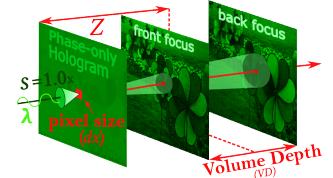


Fig. 3. A collimated light illuminates a phase-only hologram, reconstructing the images at distance Z with a certain depth VD .

color primary illuminates a subframe. Recent studies [Chao et al. 2023] introduce improvements in brightness levels in single-color holograms. Multi-color holograms extend these improvements over a broader range of s , supporting brighter images in holographic displays [Kavaklı et al. 2023b; Markley et al. 2023] without requiring higher power light sources. Please note that a typical hologram optimization pipeline [Kavaklı et al. 2023b] optimizes multi-color holograms in minutes for requested s while requiring RGB-D inputs.

3.2 Learning To Synthesis 3D Holograms

Our model converts the hologram synthesis process in Eq. (1) into a single-stage configurable learned model. The final form of our solution is a student model refined from a teacher model using a KD technique. Thus, we have to introduce our teacher model first and follow up with our student model. *Note that the distillation process plays a crucial role in deriving faster and smaller student models as these student models could not be trained for our tasks from scratch without the guidance of a teacher model.*

3.2.1 Teacher Model. We provide a complete layout of our teacher model in Fig. 4. Our teacher model takes I_{input} , λ_p , d_x , s , VD , Z , as inputs in our training. During our training, we vary these inputs so that our model can condition the hologram synthesis process for preferred display-scene parameters at the test time. Specifically, our training draws the input variables from the following set,

$$\begin{aligned} \lambda_p &\subseteq \{(640, 515, 470)\} \text{ nm}, s \subseteq \{1.0, 1.4, 1.8\}, \\ VD &\subseteq \{4.0, 8.0\} \text{ mm}, Z \subseteq \{2, 4, 7, 10\} \text{ mm}, d_x \subseteq \{3.74\} \mu\text{m}. \end{aligned} \quad (2)$$

Our choice in Z follow the most recent literature [Akşit and Itoh 2023; Shi et al. 2021, 2022] and our choice of VD supports the common focal length of VR headsets (40–75 mm), which roughly translates to placing virtual images from our VD between 5 Diopter to infinity. We generate permutations of our set in Eq. (2), resulting in total 24 cases in training. We train our teacher model for $I_{input} \in \mathbb{R}^{3 \times H \times W}$, using an entire set of our permutations in one training session. Our teacher model, shown in Fig. 4, processes I_{input} using a U-Net structure [Ronneberger et al. 2015]. The encoder of our U-Net, Enc , is the EfficientNet 1B¹ by Tan and Le [2019]. In the decoder of our U-Net, Dec , the input from the previous layer goes into a transpose convolutional layer, and the processed input is then concatenated with the input from the skip connection. The later stage in our decoder involves a CBAM [Woo et al. 2018]². Note that CBAM is an attention block that processes local channel and spatial information, which we found necessary in our trials as it allowed us to estimate depth more accurately. The last stage of our decoder conditions our U-Net with display-scene parameters: λ , s , d_x , VD , Z , denoted as $Param_{cond}$. More specifically, we embed $Param_{cond}$ using a combination of floating-point numbers and One-Dimensional (1D) Point-Spread Function (PSF), which encapsulate the conditions defined by these parameters. The parameters in floating-point numbers are concatenated and processed using sinusoidal embedding. Then, to obtain the 1D PSF, we first generate the 2D PSF of amplitude and phase using λ , d_x , and Z . Given the circular symmetry and separable characteristics of 2D PSF, we extract the central x-axis to create

¹Implemented from “https://github.com/qubvel-org/segmentation_models.pytorch”.

²Source code adopted from “<https://github.com/Jongchan/attention-module>”.

the 1D PSF, followed by two consecutive linear layers to reduce its dimension. The sinusoidal embedded numbers and 1D PSF are then concatenated, and processed by another two consecutive linear layers. Finally, we sum the linear layer and the CBAM output together. The 2D PSF has been used for Z conditioning in CGH [Asano et al. 2024]; We choose the 1D PSF due to its computational efficiency. The conditioning process, as depicted in Fig. 5, enables our teacher model to synthesize holograms for various parameters.

We use features harvested from every level of our U-Net’s decoder and process it following the work by Lin et al. [2017]. The features from various levels of our U-Net are processed with convolutional layers with batch normalization and a non-linear activation function of ReLU. We upsample these features using bilinear filters to match

the dimensions with the rest of our teacher model. Finally, we process these upsampled features further using another set of convolutional layers with batch normalization and sum them to form our latent code.

We provide teacher model’s latent code to our “heads” dedicated to three separate tasks. These tasks are predicting phase-only holograms, light source powers, and depth of a given RGB-only image. The outputs from all heads are regularized using the following penalization term,

$$\mathcal{L}_{train} = \alpha_0 \mathcal{L}_{recon} + \alpha_1 \mathcal{L}_{light} + \alpha_2 \mathcal{L}_{depth}, \quad (3)$$

where α s weight our penalization terms ($\alpha_0 = 1$, $\alpha_1 = 1$, $\alpha_2 = 30$), \mathcal{L}_{recon} represents multiplane image reconstruction loss introduced by Kavaklı et al. [2023a], \mathcal{L}_{light} represents light power constrain as introduced by Kavaklı et al. [2023b], and \mathcal{L}_{depth} represents the depth estimation. In our implementation, we also compare the color performance of reconstructed images using regularizations from Andersson et al. [2020]. This modification to \mathcal{L}_{recon} is an addition to the original work by Kavaklı et al. [2023a] and improves color consistency in our reconstructions. Please refer to Supplementary Sec.C for details of the practical aspects of our \mathcal{L}_{train} .

Phase-only hologram synthesis task. The phase head consists of one convolutional layer. The output of the phase head, depicted in Fig. 4, provides pseudo amplitude and phase values, which are later converted into a complex-valued field. This complex-valued field is propagated using a band-limited Angular Spectrum Method (ASM) method with the same Z and λ over all subframes. We then retrieve the imaginary part from the complex field as our phase-only hologram. Additionally, to support long Z (e.g. 10 mm), we introduce

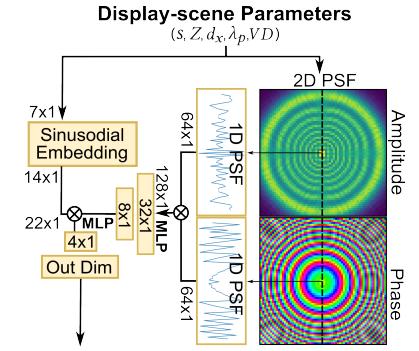


Fig. 5. Display-scene parameters and 1D PSF are used as the conditions of our model.

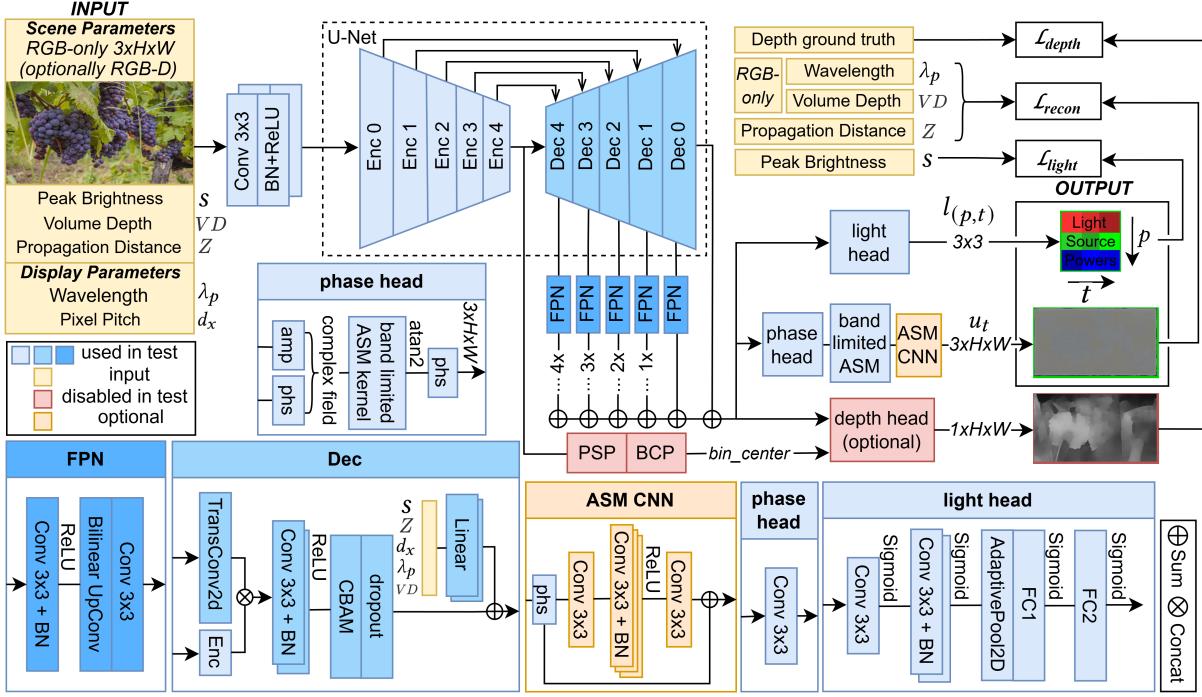


Fig. 4. Our teacher model. A Feature Pyramid Network (FPN) [Lin et al. 2017] is connected to every stage of our U-Net’s decoder to leverage spatial information at varying scales from an input (RGB-only or RGB-D). Convolutional Block Attention Module (CBAM) [Woo et al. 2018] in our model introduce both channel and spatial attention mechanisms, while each decoder is conditioned with peak brightness, wavelength, volume depth, propagation distance, and pixel pitch. A Pyramid Spatial Pooling (PSP) layer [He et al. 2015] enhances our model to aggregate spatial context for depth estimation. Our light head layer predicts the required light source powers. Here, BCP represents the Bin Center Predictor (BCP) [Agarwal and Arora 2023] (RGB-only source: [Gashi 2012]).

the ASM CNN block, depicted in Fig. 4, which consists of a convolutional layer, followed by three consecutive convolutional layers, batch normalization, and ReLU (See Supplementary Sec. A.2.4 for the details of ASM CNN). The output from ASM CNN is added with the previous phase from the propagated complex field, culminating in the final phase output shown in Fig. 4.

Light power estimation task. The light head in Fig. 4 takes the latent code from our teacher model as the input. The latent code is first processed by a convolutional layer and consecutive convolutional layers with batch normalizations and sigmoids. Next, the output from the previous layer is fed into an adaptive pooling layer to aggregate the spatial information. Finally, the pooled output is processed by two linear layers followed by a sigmoid to ensure the shape of the prediction is $t \times p$ (e.g. 3×3) ranging between 0 and 1.

Depth estimation task. Our teacher model includes a depth head as depicted in Fig. 4. We feed the latent output of our U-Net and the bin center predictor’s output to the depth head as shown in Fig. 4. The bin center predictor cascades two modules: PSP [2015] and BCP³ [2023], which estimates the distribution of depth values across different bins in the depth map. We feed our U-Net’s last encoder layer output (Enc 4) into the PSP within the bin center

predictor. In the depth head, the model’s latent code passes through a depthwise convolutional layer and softmax. Then, we multiply the softmax layer’s output with the BCP prediction to generate the final depth prediction. We apply regularization to the depth map from the depth head using a loss function, $\mathcal{L}_{depth} = \mathcal{L}_{silog} + \mathcal{L}_{gm} + \mathcal{L}_{tv}$. The Scale Invariant Log (SILog) loss, \mathcal{L}_{silog} , is introduced by Eigen et al. [2014]. The Gradient Matching (GM) loss, \mathcal{L}_{gm} , compares the edges of estimated depths with the ground truth. The Total Variation (TV) loss, \mathcal{L}_{tv} , smoothes the edges of the depth objects. See Supplementary Sec. C for the implementation details of \mathcal{L}_{depth} .

3.2.2 Student Model. Our teacher model can be accelerated by converting it into a simpler student model, depicted in Fig. 6, through the KD technique [Wang and Yoon 2021], achieving up to 2x speed compared with Shi et al. [2022] under fp32 (See Supplementary Fig. 1 for inference times). *Without KD, successful hologram synthesis wouldn’t be possible in smaller models when they are trained individually (See Supplementary Sec. M).* Further engineering efforts are required to explore acceleration through quantization optimization (TensorRT). Similar to the teacher model, the student model also takes I_{input} , λ_p , s , d_x , VD , and Z as inputs. Furthermore, our student model employs the U-Net structure to process I_{input} , but

³Source code adopted from “<https://github.com/ashutosh1807/PixelFormer>”.

with a different encoder - MobileNetV3⁴ [Howard et al. 2019]. In our student decoders, similar to the teacher model, the input from the previous layer also goes into a transpose convolutional layer, and the processed input is then concatenated with the input from the skip connection. Concatenated information is processed by a consecutive depthwise separable convolutional layer with batch normalization and ReLU. The last stage of our student decoder conditions our U-Net with λ , s , d_x , VD , and Z . The structure of the linear condition layers remains the same as the teacher model. The structures of the prediction heads in the student model are also the same as in the teacher model. In the training of our student model,

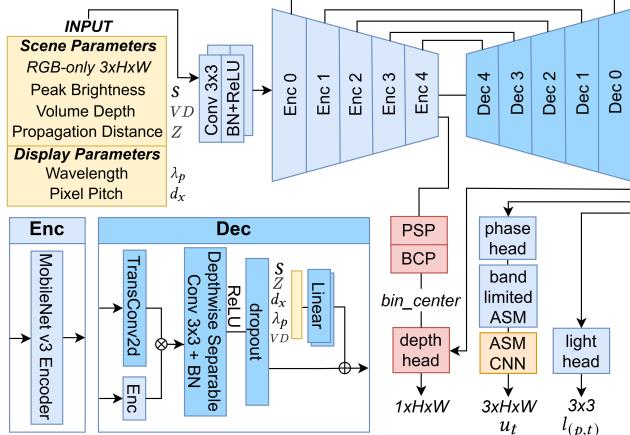


Fig. 6. The overview of the student model.

the logit-based KD [Hinton et al. 2015] is applied to balance the model performance and the inference speed. The original KD loss \mathcal{L}_{KD} , is defined as:

$$\mathcal{L}_{KD} = T^2 \cdot D_{KL}(S(y_{student}/T), S(y_{teacher}/T)), \quad (4)$$

where D_{KL} represents the Kullback-Leibler divergence [Kullback 1951], $y_{student}$ and $y_{teacher}$ are the phase and depth predictions of the student and teacher models, S denotes the softmax function, and the hyperparameter, T , also known as temperature, controls the softness of probabilities. The scalar factor T^2 is used to scale the loss to the correct magnitude. To have a robust and accurate depth distillation, we combined the $\mathcal{L}_{KD_{depth}}$ with the depth estimation loss, \mathcal{L}_{depth} . The target of the \mathcal{L}_{depth} will be depth prediction from the teacher model instead of the ground truth. To maintain a smooth phase profile, we involved Charboninner loss [Barron 2017] as our pixel-wise distillation loss of phase prediction between our student and teacher model. The updated distillation loss $\mathcal{L}_{distill}$ is defined as $\mathcal{L}_{distill} = \mathcal{L}_{KD_{phase}} + \mathcal{L}_{charboninner} + \mathcal{L}_{KD_{depth}} + \mathcal{L}_{depth}$. The final loss that combines the original training loss and the KD loss is defined as $\mathcal{L} = \mathcal{L}_{distill} + \mathcal{L}_{train}$.

4 Implementation

We utilize the Segment Anything (SA) dataset SA-1B to create a dataset to train our models. We used MiDaS [Lasinger et al. 2019]

⁴Source code adopted from “https://github.com/qubvel-org/segmentation_models.pytorch”.

to generate ground truth depth in the SA-1B dataset. We utilize a subset of SA-1B dataset for creating our dataset, specifically 44,000 images (0.4% of the entire dataset). We employ the ground truth images and depth information to generate target reconstruction images using the work by Kavaklı et al. [2023b]. We implemented our solutions all using Pytorch [2019] and Odak [2024]. We choose Adam [Kingma and Ba 2015] as the optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.99$) and CosineAnnealingLR [2017] as the Learning Rate (LR) scheduler in our training. A single batch in our training contains 16 images. Due to memory constraints, the teacher model was trained with 44K images at the resolution of 800×800 , while the student model was trained with the same images at the resolution of 1024×1024 . Both models were trained for 20 epochs with an initial LR of 0.001. All experiments were conducted on eight NVIDIA RTX 3090s. The total time to train the teacher model and distill the student model was approximately four days. Note that training our teacher model for a wide variety of display-scene parameters could easily exhaust computational resources. Therefore, to demonstrate broader generalization capabilities of our model structure with reduced computational cost, we also derived an RGB-D version of our teacher model where the depth estimation task is not included, called the RGB-D condition model. We trained the RGB-D condition model with more permutations and a wider range of λ , d_x , Z , and VD , see Supplementary Sec. F, Sec. J, Sec. K, Sec. L, Supplementary’s Fig. 6-7, 12-17 and Tbl. 4, 7, 8, 9 for results of novel cases (display-scene parameters outside of the training set). Moreover, we provide results from two different holographic displays, see Supplementary Sec. E, Supplementary’s Fig. 2-4 for our experiment results.

5 Evaluation

Quantitative Analysis. Tbl. 2 provides the image metric values of PSNR, SSIM, LPIPS [2018], FLIP [2020] and FVVDP [2021] for our implementations. We calculate these metrics using 100 test images from the DIV2K dataset [2017] that are not used in our training. Our student model maintained comparable image quality with our teacher model (Δ PSNR<0.1%). Compared with the other models, the student model can generate 3D holograms with various display-scene parameters without compromising the image quality. This is evidenced by Δ PSNR and Δ SSIM differences of 7.0% and 6.3% between the best and worst set of display-scene parameters, and improvements of +4.9% and +2.9% between student model and the other display-scene specific models. For LPIPS, though Tensor V2 has the highest value, overall the student model is as good as the other models. For FLIP, all the models have similar performance. For FVVDP, our student model achieves the highest value. In addition, we conduct a pair-wise ANOVA test and observed no significant effect, see Supplementary Sec. F for details. In line with existing iterative methods [Kavaklı et al. 2023a,b], increasing the s or Z tends to lower image quality for all models across different configurations. For instance, $s = 1.8\times$ induces a decrease of 1.9%, 1.1%, 2.6%, 8.3% and 3.6% in PSNR, SSIM, LPIPS, FLIP, FVVDP, respectively. To avoid overcrowding, we only present a subset of our model results in Tbl. 2. See Supplementary Tbl. 3 for the extended evaluation table.

Qualitative Analysis. We evaluate our student model, using simulation results. We also provide an extensive set of experimental

Method	Input	Display-scene Parameters	PSNR↑ (dB)		SSIM↑		LPIPS↓	FVVD↑	Parameters	Speed (fp32)
			Mean	Std	Mean	Std	Mean	Mean		
Our Method (teacher)	RGB-only	$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.0x$	29.33	2.73	0.95	0.03	0.35	0.10	8.30	
		$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.4x$	28.78	2.84	0.94	0.03	0.36	0.11	8.23	
		$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.8x$	27.65	2.80	0.94	0.04	0.36	0.11	8.20	
		$d_x: 3.74 \mu m, Z: 4 mm, VD: 4 mm, 1.0x$	28.16	3.07	0.94	0.03	0.38	0.12	8.18	
		$d_x: 3.74 \mu m, Z: 7 mm, VD: 4 mm, 1.0x$	27.80	3.06	0.93	0.03	0.40	0.12	8.11	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.0x$	26.92	2.71	0.91	0.04	0.46	0.14	7.89	10.74 M
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.4x$	26.97	2.66	0.91	0.04	0.47	0.14	7.93	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.8x$	26.54	2.56	0.90	0.04	0.48	0.15	7.83	
		$d_x: 3.74 \mu m, Z: 2 mm, VD: 8 mm, 1.0x$	28.23	2.87	0.93	0.05	0.39	0.12	8.00	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 8 mm, 1.0x$	26.31	2.69	0.88	0.05	0.49	0.15	7.70	
Our Method (student)	RGB-only	$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.0x$	28.55	2.88	0.95	0.03	0.35	0.10	8.48	
		$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.4x$	28.32	2.79	0.94	0.03	0.35	0.10	8.33	
		$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.8x$	27.69	2.78	0.94	0.04	0.36	0.11	8.18	
		$d_x: 3.74 \mu m, Z: 4 mm, VD: 4 mm, 1.0x$	28.29	3.01	0.94	0.03	0.36	0.11	8.22	
		$d_x: 3.74 \mu m, Z: 7 mm, VD: 4 mm, 1.0x$	28.28	3.15	0.93	0.03	0.39	0.12	8.15	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.0x$	27.15	2.81	0.91	0.03	0.44	0.15	7.94	2.19 M
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.4x$	27.07	2.72	0.91	0.04	0.44	0.15	7.90	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.8x$	26.92	2.43	0.90	0.04	0.45	0.16	7.64	
		$d_x: 3.74 \mu m, Z: 2 mm, VD: 8 mm, 1.0x$	27.89	2.98	0.93	0.04	0.39	0.12	8.02	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 8 mm, 1.0x$	26.56	2.71	0.89	0.05	0.46	0.15	7.78	
Our Method (RGB-D condition)	RGB-D	$d_x: 6.4 \mu m, Z: 4.88 mm, VD: 6.75 mm, 1.0x^1$	27.73	1.98	0.93	0.02	0.37	0.11	8.64	
		$d_x: 7.19 \mu m, Z: 2 mm, VD: 4.23 mm, 1.0x^1$	29.83	2.81	0.96	0.02	0.31	0.09	8.79	6.84 M
		$d_x: 4.57 \mu m, Z: 10 mm, VD: 7.61 mm, 1.0x^1$	26.93	2.19	0.91	0.03	0.42	0.13	8.07	
HoloBeam	RGB-only ²	$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.0x$	29.12	2.70	0.93	0.04	0.37	0.10	8.37	
		$d_x: 3.74 \mu m, Z: 4 mm, VD: 4 mm, 1.0x$	27.15	2.62	0.90	0.03	0.43	0.11	8.24	1.94 M
		$d_x: 3.74 \mu m, Z: 7 mm, VD: 4 mm, 1.0x$	25.31	2.24	0.87	0.03	0.47	0.13	8.11	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.0x$	20.62	2.51	0.82	0.03	0.49	0.14	7.99	
Tensor V2	RGB-D	$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.0x$	29.23	2.26	0.97	0.03	0.34	0.10	8.47	
		$d_x: 3.74 \mu m, Z: 4 mm, VD: 4 mm, 1.0x$	28.01	2.01	0.95	0.03	0.37	0.10	8.32	0.04 M
		$d_x: 3.74 \mu m, Z: 7 mm, VD: 4 mm, 1.0x$	25.89	2.09	0.93	0.03	0.39	0.11	8.29	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.0x$	23.04	1.98	0.91	0.03	0.41	0.11	8.12	
modified 3D NH ³	RGB-D	$d_x: 3.74 \mu m, Z: 2 mm, VD: 4 mm, 1.0x$	29.01	2.45	0.94	0.03	0.40	0.10	8.41	
		$d_x: 3.74 \mu m, Z: 4 mm, VD: 4 mm, 1.0x$	29.03	2.31	0.92	0.04	0.41	0.11	8.25	3.87 M
		$d_x: 3.74 \mu m, Z: 7 mm, VD: 4 mm, 1.0x^4$	28.92	2.18	0.90	0.04	0.43	0.13	8.17	
		$d_x: 3.74 \mu m, Z: 10 mm, VD: 4 mm, 1.0x^4$	28.98	2.39	0.89	0.04	0.45	0.13	8.05	

Table 2. Evaluation of various models using different display-scene parameters. *Parameters* refers to the size of the model. *Speed* refers to the inference time of the model under fp32. Note that *RGB teacher*, *RGB student*, *HoloBeam*, *Tensor V2* and *modified 3D NH* were trained with 800×800 , 1024×1024 , 1024×1024 , 384×384 and 1024×1024 data, respectively. The test resolution is 1920×1080 . ¹Novel cases are in italic and *not* included in metrics comparison. ²We improve the *HoloBeam* from G-only to RGB-only input. ³Since NH only supports 2D holograms, we modified NH’s *HoloNet* [Peng et al. 2020] to generate 3D holograms. ⁴In our experiment, training *modified 3D NH* with 1024×1024 resolution for long propagation distance will severely decrease image quality. Hence, we used 1920×1080 resolution to train the *modified 3D NH*. See Supplementary Tbl. 3 for the complete evaluation table, Sec. J, Sec. K and Sec. L for the full experiments on novel cases, and Supplementary Sec. H for the details of modified 3D NH.

results from two actual holographic display prototypes in our Supplementary Sec. E. Fig. 7 shows the simulated and captured results of the student model (RGB-only) at a short Z (2 mm). Our method with short Z retained the reconstruction of fine details, color accuracy and texture preservation in different s . Fig. 8 shows the simulated and captured results of the student model (RGB-only) at a long Z (10 mm). Compared with short Z , the results show moderate color deviations. In long propagation, the larger diffraction cone requires higher resolution to capture the beam propagation accurately. We believe training at lower resolutions results in artifacts that manifest as color deviations for long Z . Due to demanding computation ($8 \times A100$ GPUs) and time (> 20 days) requirements, training our RGB-only models with larger images was not feasible. To validate our argument, we used full HD data to train our RGB-D condition model for long Z and compared it with modified 3D NH and Tensor V2, achieving competitive image quality with no color deviations, see Supplementary Sec. H for details.

Configurability Evaluation. Our RGB-only models continuously support s between $1.0 - 1.8 \times$, VD between $4-8 mm$ and support a discrete set of $Z \subseteq \{2, 4, 7, 10\} mm$, while maintaining fixed d_x and

Ours	s	Z	VD	λ	d_x
RGB-only	con	dis	con	fixed	fixed
RGB-D	fixed	con (Sec. K*)	con	con (Sec. L*)	con (Sec. J*)

Table 3. Configurability of our models. *See the corresponding Supplementary sections.

We benefit from the reduced computational complexity of RGB-D condition model and conduct three training examples to verify our model structure, supporting continuous ranges for Z , d_x , and λ . These examples include: (1) Z experiment: Continuous ranges of 2-3, 4-5, 6-7, 8-9, and 10-11 mm (total of 5 mm distance) with three discrete d_x values ($3.74, 6.4, 8.0 \mu m$); (2) d_x experiment: Continuous range of $3.7-8.0 \mu m$ with two discrete Z values ($2.0, 10.0 mm$); (3) λ experiment: Continuous ranges of $425-480, 510-565, 625-680 nm$ with two discrete d_x values ($3.74, 8.0 \mu m$) and two discrete Z values ($2.0, 10.0 mm$). All our RGB-D training supports continuous VD between 4-8 mm with fixed s . Based on our various experiments, Tbl. 3 shows the respective flexibility of our models in generalizing various display-scene parameters. See our Supplementary’s Sec. J, Sec. K and Sec. L for the details of our three experiments.

Limitations. Extended Training Time: Covering a wide range of display-scene parameters increases the training set. Future variants need to invent a training strategy with fewer display-scene permutations in training. Improving Defocus: The defocus in our holograms is strongly correlated with the quality of predicted depth (See Supplementary Sec. I). To improve defocus accuracy, we plan to borrow strategies from SOTA MDE works [Bhat et al. 2023; Yang et al. 2024]. Parameter Sensitivity: In Supplementary’s Fig. 13, 15–17, we observe image quality differences with long range of changes in display-scene parameters (e.g., Z), requiring further investigation.

Discussion. Holographic displays have shown promising improvements, especially in applications akin to glasses [Gilles et al. 2023; Jang et al. 2024]. Our research represents the first attempt to resolve the challenging case of 3D hologram synthesis for various display-scene parameters in a single model. Our ultimate goal is to make the model continuously configurable, adapting to any novel configuration outside of the variable set without training densely.

Acknowledgments

In review.

References

- Ashutosh Agarwal and Chetan Arora. 2023. Attention Attention Everywhere: Monocular Depth Prediction with Skip Attention. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2023, Waikoloa, HI, USA, January 2–7, 2023*. IEEE, 5850–5859. doi:10.1109/WACV56688.2023.00581
- Eirikur Agustsson and Radu Timofte. 2017. NTIRE 2017 Challenge on Single Image Super-Resolution: Dataset and Study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- Kaan Akşit and Yuta Itoh. 2023. HoloBeam: Paper-Thin Near-Eye Displays. In *2023 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*. IEEE, 581–591.
- Kaan Akşit, Jeanne Beyazian, Praneeth Chakravarthula, Ziyang Chen, Mustafa Doğa Doğan, Ahmet Hamdi Güzel, Yuta Itoh, Henry Kam, Ahmet Serdar Karadeniz, Koray Kavaklı, Liang Shi, Josef Spjut, David Robert Walton, Doğa Yılmaz, Wang Yujie, Runze Zhu, and Yicheng Zhan. 2024. *Odak*. doi:10.5281/zenodo.10407993
- Pontus Anderson, Jim Nilsson, Tomas Akenine Moller, Magnus Oskarsson, Kalle AAstrom, and Mark D. Fairchild. 2020. FLIP: A Difference Evaluator for Alternating Images. *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 3, 2 (2020), 15:1–15:23. doi:10.1145/3406183
- Yuto Asano, Kenta Yamamoto, Tatsuki Fushimi, and Yoichi Ochiai. 2024. Distance-adaptive unsupervised CNN model for computer-generated holography. In *ACM SIGGRAPH 2024 Posters* (Denver, CO, USA) (*SIGGRAPH’24*). Association for Computing Machinery, New York, NY, USA, Article 64, 2 pages. doi:10.1145/3641234.3671051
- Jonathan T. Barron. 2017. A More General Robust Loss Function. *CoRR* abs/1701.03077 (2017). CVPR 2019:1701.03077 <http://arxiv.org/abs/1701.03077>
- Shariq Farooq Bhat, Reiner Birk, Diana Wofk, Peter Wonka, and Matthias Müller. 2023. ZoeDepth: Zero-shot Transfer by Combining Relative and Metric Depth. *CoRR* abs/2302.12288 (2023). doi:10.48550/ARXIV.2302.12288 arXiv:2302.12288
- David Blinder, Ayyoub Ahar, Stijn Bettens, Tobias Birnbaum, Athanasia Symeonidou, Heidi Ottevaere, Colas Schretter, and Peter Schelkens. 2019. Signal processing challenges for digital holographic video display systems. *Signal Processing: Image Communication* 70 (2019), 114–130.
- Maurizio C. 2014. Lambos. <https://openverse.org/image/1fe39884-2382-4857-b119-981380a14aac?l=lambo>
- R Caruana. 1993. Multitask learning: A knowledge-based source of inductive bias1. In *Proceedings of the Tenth International Conference on Machine Learning*. Citeseer, 41–48.
- Minseok Chae, Kiseung Bang, Dongheon Yoo, and Yoonchan Jeong. 2023. Étendue Expansion in Holographic Near Eye Displays through Sparse Eye-box Generation Using Lens Array Eyepiece. *ACM Trans. Graph.* 42, 4 (2023), 58:1–58:13. doi:10.1145/3592441
- Brian Chao, Manu Gopakumar, Suyeon Choi, and Gordon Wetzstein. 2023. High-brightness holographic projection. *Optics Letters* 48, 15 (2023), 4041–4044.
- Hanting Chen, Yunhe Wang, Chang Xu, Chao Xu, and Dacheng Tao. 2020. Learning student networks via feature embedding. *IEEE Transactions on Neural Networks and Learning Systems* 32, 1 (2020), 25–35.
- Suyeon Choi, Manu Gopakumar, Yifan Peng, Jonghyun Kim, Matthew O’Toole, and Gordon Wetzstein. 2022. Time-multiplexed Neural Holography: A Flexible Framework for Holographic Near-eye Displays with Fast Heavily-quantized Spatial Light Modulators. In *SIGGRAPH ’22: Special Interest Group on Computer Graphics and Interactive Techniques Conference, Vancouver, BC, Canada, August 7 – 11, 2022*. Munkhtsetseg Nandigav, Niloy J. Mitra, and Aaron Hertzmann (Eds.). ACM, 32:1–32:9. doi:10.1145/3528233.3530734
- Suyeon Choi, Manu Gopakumar, Yifan Peng, Jonghyun Kim, and Gordon Wetzstein. 2021a. Neural 3D Holography: Learning Accurate Wave Propagation Models for 3D Holographic Virtual and Augmented Reality Displays. *ACM Trans. Graph. (SIGGRAPH Asia)* (2021).
- Suyeon Choi, Manu Gopakumar, Yifan Peng, Jonghyun Kim, and Gordon Wetzstein. 2021b. Neural 3D holography: learning accurate wave propagation models for 3D holographic virtual and augmented reality displays. *ACM Transactions on Graphics (TOG)* 40, 6 (2021), 1–12.
- Long Duong, Trevor Cohn, Steven Bird, and Paul Cook. 2015. Low resource dependency parsing: Cross-lingual parameter sharing in a neural network parser. In *Proceedings of the 53rd annual meeting of the Association for Computational Linguistics and the 7th international joint conference on natural language processing (volume 2: short papers)*. 845–850.
- David Eigen, Christian Puhrsch, and Rob Fergus. 2014. Depth Map Prediction from a Single Image using a Multi-Scale Deep Network. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8–13 2014, Montreal, Quebec, Canada*. Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger (Eds.). 2366–2374. <https://proceedings.neurips.cc/paper/2014/hash/7bcfde7714a1ebad06c5f4cea752c1-Abstract.html>
- Bujar I Gashi. 2012. Rahovec Grapes and Wine. <https://openverse.org/image/4f2d2421-f354-449a-971a-cb6340080e23?q=grapes>
- Antonin Gilles, Pierre Le Gargasson, Grégory Hocquet, and Patrick Gioia. 2023. Holographic Near-eye Display with Real-time Embedded Rendering. In *SIGGRAPH Asia 2023 Conference Papers, SA 2023, Sydney, NSW, Australia, December 12–15, 2023*. June Kim, Ming C. Lin, and Bernd Bickel (Eds.). ACM, 21:1–21:10. doi:10.1145/3610548.3618179
- Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision* 129 (2021), 1789–1819.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 37, 9 (2015), 1904–1916. doi:10.1109/TPAMI.2015.2389824
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *NIPS 2014 Deep Learning Workshop* (2015).
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. 2019. Searching for MobileNetV3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Chung-Kai Hsueh and Alexander A Sawchuk. 1978. Computer-generated double-phase holograms. *Applied optics* 17, 24 (1978), 3874–3883.
- Andrey Ignatov, Grigory Malivenko, David Plowman, Samarth Shukla, and Radu Timofte. 2021. Fast and accurate single-image depth estimation on mobile devices, mobile ai 2021 challenge: Report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2545–2557.
- Yoshiyuki Ishii, Fan Wang, Harutaka Shiomi, Takashi Kakue, Tomoyoshi Ito, and Tomoyoshi Shimobaba. 2023. Multi-depth hologram generation from two-dimensional images by deep learning. *Optics and Lasers in Engineering* 170 (2023), 107758. doi:10.1016/j.optlastec.2023.107758
- Changwon Jang, Kiseung Bang, Minseok Chae, Byoungho Lee, and Douglas Lanman. 2024. Waveguide holography for 3D augmented reality glasses. *Nature Communications* 15, 1 (02 Jan 2024), 66. doi:10.1038/s41467-023-44032-1
- Koray Kavaklı, Yuta Itoh, Hakan Urey, and Kaan Akşit. 2023a. Realistic defocus blur for multiplane computer-generated holography. In *2023 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*. IEEE, 418–426.
- Koray Kavaklı, Liang Shi, Hakan Urey, Wojciech Matusik, and Kaan Akşit. 2023b. Multi-color Holograms Improve Brightness in Holographic Displays. In *ACM SIGGRAPH ASIA 2023 Conference Proceedings*. –. doi:10.1145/3610548.3618135
- Koray Kavaklı, Hakan Urey, and Kaan Akşit. 2022. Learned holographic light transport. *Applied Optics* 61, 5 (2022), B50–B55.
- Dongyeon Kim, Seung-Woo Nam, Suyeon Choi, Jong-Mo Seo, Gordon Wetzstein, and Yoonchan Jeong. 2024. Holographic Parallax Improves 3D Perceptual Realism. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–13.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings*. Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6980>
- Solomon Kullback. 1951. Kullback-leibler divergence.

- Grace Kuo, Florian Schiffers, Douglas Lanman, Oliver Cossairt, and Nathan Matsuda. 2023. Multisource Holography. *ACM Trans. Graph.* 42, 6, Article 203 (dec 2023), 14 pages. doi:10.1145/3618380
- Katrin Lasinger, René Ranftl, Konrad Schindler, and Vladlen Koltun. 2019. Towards Robust Monocular Depth Estimation: Mixing Datasets for Zero-Shot Cross-Dataset Transfer. *CoRR* abs/1907.01341 (2019). TPAMI:1907.01341 <http://arxiv.org/abs/1907.01341>
- Wei-Hong Li and Hakan Bilen. 2020. Knowledge distillation for multi-task learning. In *Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*. Springer, 163–176.
- Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature Pyramid Networks for Object Detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 936–944. doi:10.1109/CVPR.2017.106
- Kexuan Liu, Jiachen Wu, Zehao He, and Liangcai Cao. 2023b. 4K-DMDNet: diffraction model-driven network for 4K computer-generated holography. *OptoElectron Adv*, 5 (2023), 220135–1–220135–13. doi:10.29026/oea.2023.220135
- Ninghe Liu, Zhengzhong Huang, Zehao He, and Liangcai Cao. 2023a. DGE-CNN: 2D-to-3D holographic display based on depth gradient extracting module and CNN network. *Optics Express* 31, 15 (2023), 23867–23876.
- LogicalRaifan. 2023. City Hall north platform NB. <https://openverse.org/image/6483e8d9-e965-4e98-9bdc-db3f776b6495?q=city>
- Ilya Loshchilov and Frank Hutter. 2017. SGD: Stochastic Gradient Descent with Warm Restarts. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference Track Proceedings*. OpenReview.net. <https://openreview.net/forum?id=Skq89Scxx>
- Rafal K. Mantiuk, Gyorgy Denes, Alexandre Chapiro, Anton Kaplyanyan, Gizem Rufo, Romain Bachy, Trisha Lian, and Anjul Patney. 2021. FovVideoVDP: a visible difference predictor for wide field-of-view video. *ACM Trans. Graph.* 40, 4 (2021), 49:1–49:19. doi:10.1145/3450626.3459831
- Eric Markley, Nathan Matsuda, Florian Schiffers, Oliver Cossairt, and Grace Kuo. 2023. Simultaneous Color Computer Generated Holography. In *SIGGRAPH Asia 2023 Conference Papers, SA 2023, Sydney, NSW, Australia, December 12–15, 2023*, June Kim, Ming C. Lin, and Bernd Bickel (Eds.). ACM, 22:1–22:11. doi:10.1145/3610548.3618250
- Kyoji Matsushima and Tomoyoshi Shimobaba. 2009. Band-limited angular spectrum method for numerical simulation of free-space propagation in far and near fields. *Optics express* 17, 22 (2009), 19662–19673.
- Paul Micaelli and Amos J Storkey. 2019. Zero-shot knowledge transfer via adversarial belief matching. *Advances in Neural Information Processing Systems* 32 (2019).
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *CoRR* abs/1912.01703 (2019). arXiv:1912.01703 <http://arxiv.org/abs/1912.01703>
- Yifan Peng, Suyeon Choi, Nitish Padmanaban, and Gordon Wetzstein. 2020. Neural holography with camera-in-the-loop training. *ACM Transactions on Graphics (TOG)* 39, 6 (2020), 1–14.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*. Springer, 234–241.
- Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. 2019. Latent multi-task architecture learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 4822–4829.
- Syed Shakib Sarwar, Aayush Ankit, and Kaushik Roy. 2019. Incremental learning in deep convolutional neural networks using partial network sharing. *IEEE Access* 8 (2019), 4615–4628.
- Liang Shi, Fu-Chung Huang, Ward Lopes, Wojciech Matusik, and David Luebke. 2017. Near-eye light field holographic rendering with spherical waves for wide field of view interactive 3D computer graphics. *ACM Trans. Graph.* 36, 6 (2017), 236:1–236:17. doi:10.1145/3130800.3130832
- Liang Shi, Beichen Li, Changil Kim, Petr Kellnhofer, and Wojciech Matusik. 2021. Towards real-time photorealistic 3D holography with deep neural networks. *Nature* 591, 7849 (2021), 234–239.
- Liang Shi, Beichen Li, and Wojciech Matusik. 2022. End-to-end learning of 3d phase-only holograms for holographic display. *Light: Science & Applications* 11, 1 (2022), 247.
- Nahian Siddique, Sidike Paheding, Colin P Elkin, and Vijay Devabhaktuni. 2021. U-net and its variants for medical image segmentation: A review of theory and applications. *Ieee Access* 9 (2021), 82031–82057.
- Bernard Spragg. 2009. Colorful windmills Shenzhen China. <https://openverse.org/image/50c510ac-be7d-4485-b397-d0bdd2dec723?q=colorfulChina>
- Mingxing Tan and Quoc Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*. PMLR, 6105–6114.
- Lin Wang and Kuk-Jin Yoon. 2021. Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks. *IEEE transactions on pattern analysis and machine intelligence* 44, 6 (2021), 3048–3068.
- Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. 2018. CBAM: Convolutional Block Attention Module. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. 2024. Depth Anything: Unleashing the Power of Large-Scale Unlabeled Data. In *CVPR*.
- Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18–22, 2018*. Computer Vision Foundation / IEEE Computer Society, 586–595. doi:10.1109/CVPR.2018.00068
- Yu Zhang and Qiang Yang. 2021. A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering* 34, 12 (2021), 5586–5609.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

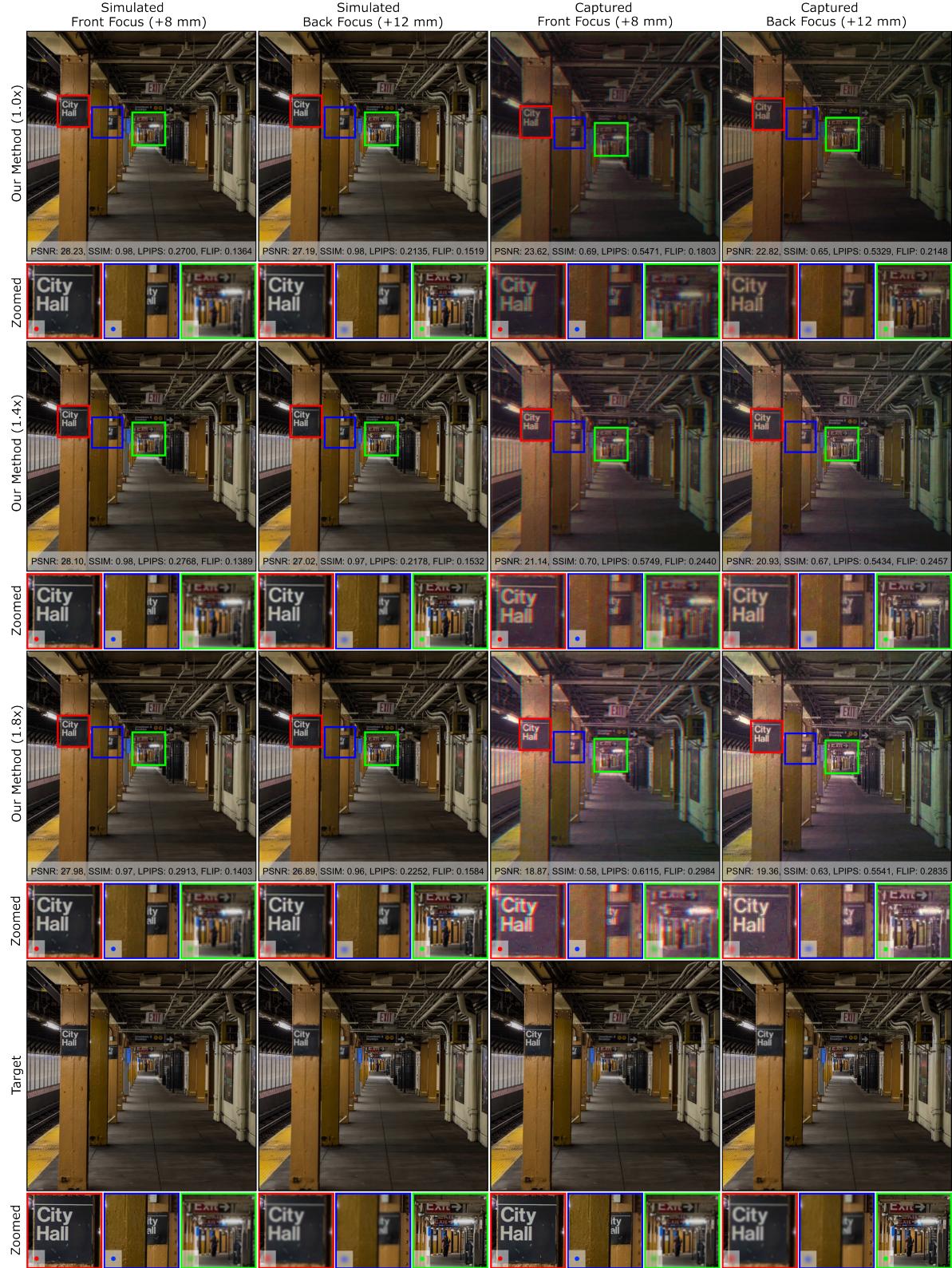


Fig. 7. The simulated and captured short Z distance reconstructions comparison of student model (RGB-only) when peak brightnesses are 1.0, 1.4, and 1.8. The volume depth of the results is 4 mm and the propagation distance is 2 mm (Image Source: [LogicalRailfan 2023]). The test resolution is 2816 × 2816.



Fig. 8. The simulated and captured long Z distance reconstructions comparison of student model (RGB-only) when peak brightnesses are 1.0, 1.4, and 1.8. The volume depth of the results is 4 mm and the propagation distance is 10 mm (Image Source: [LogicalRailfan 2023]). The test resolution is 2816 × 2816.