

Comprehensive Machine Learning Benchmarking for Fringe Projection Profilometry with Photorealistic Synthetic Data

Anush Lakshman S^{a,*}, Adam Haroon^{a,*}, and Beiwen Li^b

^aDepartment of Mechanical Engineering, Iowa State University, Ames, Iowa, USA

^bCollege of Engineering, University of Georgia, Athens, Georgia, USA

*These authors contributed equally.

ABSTRACT

Machine learning approaches for fringe projection profilometry (FPP) are hindered by the lack of large, diverse datasets and comprehensive benchmarking protocols. This paper introduces the first open-source, photorealistic synthetic dataset for FPP, generated using NVIDIA Isaac Sim with 15,600 fringe images and 300 depth reconstructions across 50 diverse objects. We benchmark four neural network architectures (UNet, Hformer, ResUNet, Pix2Pix) on single-shot depth reconstruction, revealing that all models achieve similar performance (58-77 mm RMSE) despite substantial architectural differences. Our results demonstrate fundamental limitations of direct fringe-to-depth mapping without explicit phase information, with reconstruction errors approaching 75-95% of the typical object depth range. This resource provides standardized evaluation protocols enabling systematic comparison and development of learning-based FPP approaches.

Keywords: Fringe projection profilometry, machine learning, synthetic data, deep learning, 3D reconstruction, structured light, NVIDIA Isaac Sim, benchmarking

1. INTRODUCTION AND RELATED WORK

Fringe projection profilometry (FPP) has emerged as a critical non-destructive technology in robotic scanning,¹ manufacturing inspection,² 3D printing optimization,³ offering high-precision surface measurements with submillimeter accuracy.^{4,5} While traditional FPP uses multi-step phase-shifting algorithms requiring sequential pattern capture, deep learning offers possibilities for single-shot reconstruction enabling real-time applications.⁶⁻¹²

Moreover, machine learning (ML) has shown promise in phase unwrapping,¹³ fringe denoising,¹⁴ and depth regression,⁶ but studies rely on limited datasets that don't generalize well. To solve the issue of limited datasets, synthetic data generation through virtual twins has proven powerful for optical metrology,^{15,16} using Blender,⁸ Unity,¹⁷ or MATLAB.¹⁸ However, these systems either require pre-calibrated physical systems or provide simplified optical models, thus limiting the ability to create a diverse dataset with different camera-projector configurations.

On the other hand, the current approaches to the formulation and evaluation of the single-shot problem acts as a major impediment to progress. First, unlike computer vision benchmarks such as ImageNet¹⁹ or COCO,²⁰ FPP lacks large-scale datasets with absolute ground truth and standardized evaluation protocols. Second, it is economically and physically infeasible to generate large-scale training data across diverse object geometries and lighting conditions. Third, obtaining perfect ground truth 3D geometry remains challenging as measurement systems introduce their own errors.

To address these issues, we build on VIRTUS-FPP,²¹ which introduced the first physics-based virtual FPP system with end-to-end camera-projector modeling in NVIDIA Isaac Sim, to present a systematic machine learning benchmarking framework. Our contributions include:

- First open-source synthetic FPP dataset: 15,600 fringe images and 300 depth maps for 50 diverse objects with perfect ground truth

Send correspondence to Adam Haroon: E-mail: aharoon@iastate.edu

- Comprehensive data acquisition methodology leveraging VIRTUS-FPP’s physics-based rendering and virtual calibration
- Benchmarking protocols showing UNet, Hformer, and Pix2Pix achieve nearly identical performance (58.89-60.26 mm RMSE) while ResUNet underperforms at 76.55 mm RMSE
- Demonstration that reconstruction errors (58-77 mm) approach the typical 80 mm depth range, revealing networks learn coarse shape priors rather than accurate geometry

2. VIRTUAL FRINGE PROJECTION PROFILOMETRY

The VIRTUS-FPP²¹ used for benchmarking, is built in NVIDIA Isaac Sim, integrating OptiX ray tracing for or photorealistic rendering, PhysX for physics, and Universal Scene Description (USD) for 3D composition. This section is structured as follows: Section 2.1 discusses the configuration of the FPP system used for data acquisition and Section 2.2 elaborates the calibration process of the constructed virtual system.

2.1 System Configuration

The virtual system consists of a calibrated camera-projector pair (Table 1). The camera uses Isaac Sim’s pinhole primitive (960×960 resolution, 50 cm focal length), while the projector is modeled using a rectangular light source (0.625 m × 0.5 m, 40 nits) with texture projection. The projector is positioned 0.1 m below and 0.125 m left of the camera for optimal triangulation geometry.

Table 1. Virtual Camera and Projector System Parameters

Camera Parameters	Value
Focal Length	50 cm
Horizontal Aperture	20.9995 cm
Vertical Aperture	15.2908 cm
Resolution	960 × 960 pixels
Projector Parameters	Value
Intensity	40 nits
Height	0.625 m
Width	0.5 m
Pattern Resolution	912 × 1140 pixels

VIRTUS-FPP’s key innovation is projector modeling through the inverse camera model:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = (M_{ext})^{-1}(M_{int})^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \quad (1)$$

enabling accurate dimensional correspondence of projected fringe patterns at any distance without hardware constraints. All objects in our dataset use consistent matte material properties (roughness=0.95, specular=0.15, AO-to-diffuse=0.95) representative of typical structured light scanning.^{22, 23}

The rendering pipeline uses OptiX path tracing with specific configurations: disabled sampled direct lighting mode to prevent phase map artifacts, and disabled shadows for clean fringe patterns. This physics-based approach captures complex light transport including multi-bounce illumination, surface reflectivity variations, and ambient occlusion.

2.2 Virtual Calibration

VIRTUS-FPP performs complete virtual calibration using procedurally generated 5×9 asymmetric circular boards (10 mm diameter, 20 mm spacing). The system captures 18 calibration poses yielding 936 calibration images in 5 minutes (10,530 images/hour). The calibrated system achieves sub-pixel accuracy (stereo reprojection error: 0.055506 pixels, projector error: 0.048609 pixels).²¹ Our VIRTUS-FPP simulation setup is illustrated in Figure 1.

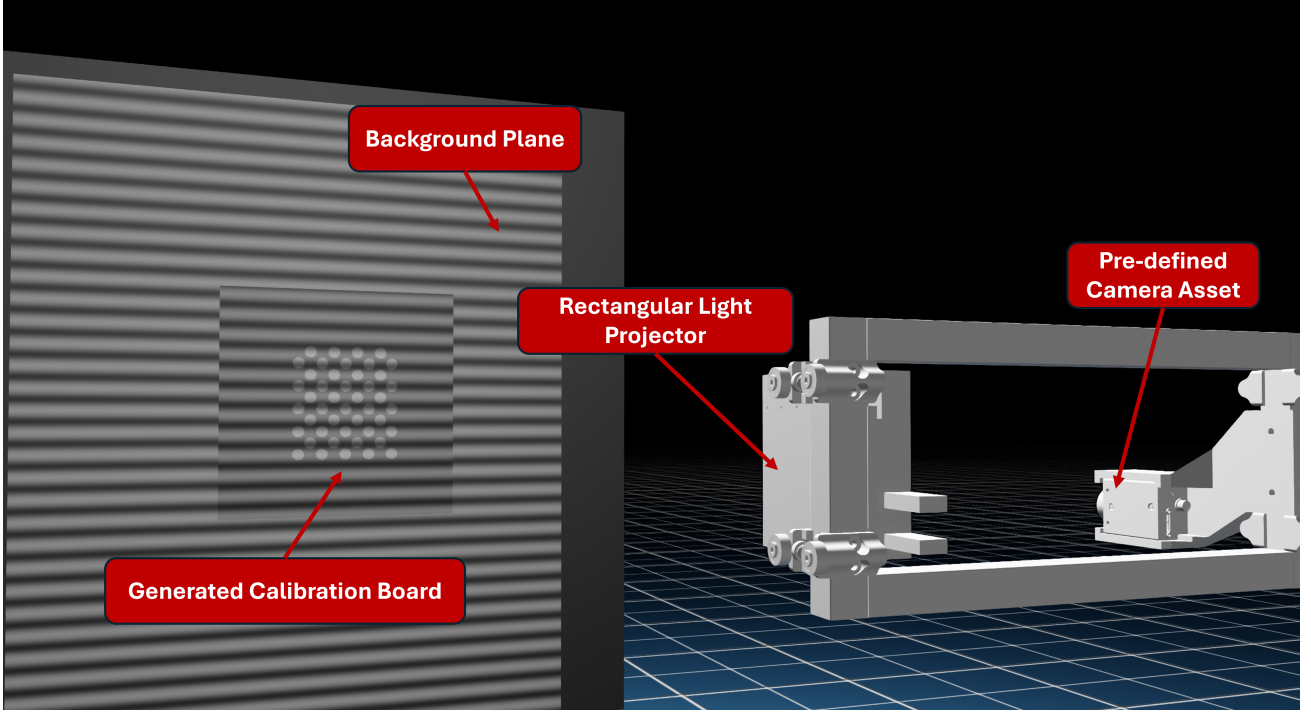


Figure 1. Virtual camera-projector calibration setup with a pinhole camera model, rectangular light-source projector, calibration board, and matte background plane.

3. DATA ACQUISITION METHODOLOGY

3.1 Dataset Composition

We collected data for 50 USD objects from YCB datasets²⁴ and NVIDIA Physical AI Warehouse²⁵ spanning cylindrical containers, rectangular boxes, complex shapes (power drills, sprayguns), and industrial components. This diversity evaluates robustness across varying surface characteristics and morphological complexity from simple geometric primitives to intricate shapes with concavities and fine-scale features.

Objects are positioned on a background plane with identical matte properties to provide consistent lighting and minimize reflections. Multi-view acquisition rotates each object about the vertical axis with 60° increments, yielding 6 viewpoints per object with 50% overlap between adjacent views:

$$R_z(\theta_i) = \begin{bmatrix} \cos \theta_i & -\sin \theta_i & 0 \\ \sin \theta_i & \cos \theta_i & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (2)$$

for $\theta_i = i \cdot 60$ where $i = 0, 1, \dots, 5$.

3.2 Fringe Acquisition and Ground Truth Generation

At each viewpoint, an 18-step phase-shifting sequence ($\delta_n = 2\pi n/18$, $n = 0, \dots, 17$) is captured at 960×960 resolution:⁴

$$I_n(u, v) = I'(u, v) + I''(u, v) \cos \left[\phi(u, v) + \frac{2\pi n}{18} \right] \quad (3)$$

where (u, v) are pixel coordinates, $I'(u, v)$ is background intensity, $I''(u, v)$ is modulation amplitude, and $\phi(u, v)$ is the phase. The GPU-accelerated pipeline achieves 3 fps, over twice the speed of previous approaches.⁸

Captured patterns are processed using standard N-step phase-shifting, Gray-code temporal unwrapping,²⁶ and triangulation to generate depth maps $D(u, v)$. Per-object normalization maps depths to $[0, 1]$:

$$D_{norm}(u, v) = \frac{D(u, v) - D_{min}}{D_{max} - D_{min}} \quad (4)$$

where D_{min} and D_{max} are object-specific. Normalized maps are stored as 16-bit PNG images (1.2×10^{-3} mm precision for 80 mm range) with normalization parameters stored separately for metric reconstruction during evaluation.

3.3 Dataset Summary

The dataset comprises 15,600 fringe images ($50 \text{ objects} \times 6 \text{ viewpoints} \times 18 \text{ patterns}$), 300 normalized depth maps, 300 normalization parameter files, and 50 ground truth mesh geometries. Data are partitioned 80/10/10 at object level: 240 training samples ($40 \text{ objects} \times 6 \text{ viewpoints}$), 30 validation samples, and 30 test samples, ensuring evaluation on completely unseen geometries.

4. SINGLE-SHOT RECONSTRUCTION BENCHMARKING

4.1 Problem Formulation

Single-shot reconstruction predicts depth $\hat{D}_{norm} = f_{\theta}(I)$ from a single fringe image I , where f_{θ} is a neural network. We use the first fringe from each 18-step sequence as input, ensuring identical conditions across models.

This task is inherently challenging, as single fringe images contain ambiguity in depth estimation due to the periodic nature of sinusoidal patterns. In the absence of temporal dependency or spatial unwrapping, each fringe cycle spans a 2π phase range, making it difficult to uniquely associate a specific cycle with a surface point. As a result, learning-based approaches rely on inferring depth from learned shape priors and statistical regularities, rather than from fully explicit geometric cues alone.

Table 2. Quantitative evaluation on 30 test samples (mm). Results show comparable performance across architectures with errors approaching the 80 mm depth range.

Model	RMSE (mm)	MAE (mm)	Median RMSE (mm)	Std RMSE (mm)	Min RMSE (mm)	Max RMSE (mm)
Pix2Pix	58.89	52.93	62.13	26.71	14.19	108.68
Hformer	59.92	52.92	66.67	29.65	9.51	106.62
UNet	60.26	53.23	67.88	27.72	10.08	102.46
ResUNet	76.55	66.92	75.64	29.31	29.61	124.34

4.2 Network Architectures

We benchmark four architectures representing different paradigms:

UNet:²⁷ Encoder-decoder with skip connections, four stages (960×960 to 60×60), channel depth 64 to 1024. Dropout 0.5 at bottleneck, Adam optimizer with RMSE loss.

Hformer:⁹ Hybrid CNN-transformer with HRNet-W18 backbone for multi-scale features [18,36,72,144], transformer encoder-decoder with window-based attention (size 8), patch expansion upsampling. Dropout 0.5, Adam optimizer with RMSE loss.

ResUNet:¹¹ UNet with residual blocks replacing convolutional blocks, four levels (960×960 to 120×120), identity skip connections for gradient flow. Dropout 0.5, RMSProp optimizer ($\alpha = 0.99$, $\text{lr}=10^{-4}$ with ReduceLROnPlateau).

Pix2Pix:²⁸ Conditional GAN with U-Net generator and PatchGAN discriminator, adapted from NVIDIA Pix2Pix-HD.²⁹ LeakyReLU, instance norm, adversarial + L1 loss, Adam optimizer ($\text{lr}=2 \times 10^{-4}$, $\beta_1 = 0.5$).

4.3 Training and Evaluation

All networks except Pix2Pix use Adam optimizer ($\text{lr}=10^{-4}$, $\beta=(0.9,0.999)$, weight decay= 10^{-5}), with learning rate reduction by 0.1 after 10 plateau epochs. Training continues for max 1000 epochs with early stopping after 50 non-improving epochs.

Loss function is RMSE over normalized depth values:

$$\mathcal{L}_{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (D_{norm}^{(i)} - \hat{D}_{norm}^{(i)})^2} \quad (5)$$

where N is the number of valid pixels (excluding background). Training uses batch size 4 (UNet) or 1 (Hformer) on NVIDIA A100 GPUs with mixed precision.

Models are evaluated after denormalizing to metric units (mm):

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (D^{(i)} - \hat{D}^{(i)})^2} \quad (6)$$

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |D^{(i)} - \hat{D}^{(i)}| \quad (7)$$

Results are reported as aggregate statistics across all 30 test samples.

4.4 Experimental Results

4.4.1 Quantitative Comparison

Table 2 shows three models achieve nearly identical performance (Pix2Pix: 58.89 mm, Hformer: 59.92 mm, UNet: 60.26 mm RMSE) despite architectural differences, while ResUNet underperforms (76.55 mm, 30% worse). Mean errors represent 74-96% of the typical 80 mm depth range. High standard deviations (26.71-29.65 mm, 50% of mean) reveal geometry-dependent brittleness: simple objects yield 9-15 mm errors while complex shapes produce 100-120 mm errors.

Figure 2 shows error distributions and per-sample errors. The curves track closely for Pix2Pix, Hformer, and UNet, while ResUNet consistently exhibits higher errors. Certain test samples present extreme difficulty for all models (80-120 mm errors), while others yield relatively low errors (10-30 mm), indicating reconstruction accuracy is dominated by geometric factors rather than model-specific capabilities.

4.4.2 Qualitative Analysis

Figure 3 shows representative results for one test object. Models capture coarse shape and approximate depth ordering but fail on fine details with exception of UNet and ResUNet. Error concentrates at boundaries and discontinuities, suggesting networks learn semantic shape completion rather than geometric reconstruction from fringe deformation. Predictions for Hformer in particular resemble smooth, regularized versions of true geometry, as if networks project ambiguous observations onto learned shape manifolds. Meanwhile, Pix2Pix is unable to predict the rough overall geometry of the object, while every other model can.

4.5 Discussion

Results demonstrate that direct single-shot fringe-to-depth regression without explicit phase extraction yields limited reconstruction accuracy. Mean errors of 59-77 mm (74-96% of 80 mm depth range) indicate networks learn coarse shape approximations rather than accurate geometry.

First, single fringe images lack sufficient information for accurate reconstruction. The periodic nature of sinusoidal patterns creates fundamental ambiguities unresolved by training on diverse geometries. Traditional FPP uses temporal redundancy or spatial unwrapping; our experiments show networks instead learn to map fringe patterns to plausible depth distributions based on statistical regularities.

Second, high variance across test objects (standard deviations 26.71-29.65 mm, 50% of mean) reveals geometry-dependent brittleness. Simple objects yield 9-15 mm errors while complex shapes produce 100-120 mm errors, suggesting single-shot approaches may only work for constrained domains with limited geometric variation.

Third, network architecture plays surprisingly limited role. Pix2Pix, Hformer, and UNet achieve nearly identical performance (58.89-60.26 mm RMSE) despite substantial design differences: adversarial vs supervised, hybrid CNN-transformer vs pure CNN. Although in the case of Hformer and UNet the overall shape of the object is captured, Pix2Pix fails to do that, which can be attributed the requirement of larger training data for Pix2Pix.³⁰ This suggests architectural innovations provide minimal benefit for this problem without addressing the fundamental information deficit. ResUNet’s underperformance (76.55 mm) despite increased depth suggests overfitting on the limited 240-sample training set.

These findings motivate incorporating explicit phase information as intermediate representations. Rather than end-to-end fringe-to-depth regression, future work should train networks to refine phase unwrapping, denoise wrapped phase, or predict depth from coarse phase estimates, leveraging geometric structure while benefiting from learned priors.

5. CONCLUSION AND FUTURE WORK

This paper presents the first comprehensive machine learning benchmarking framework for FPP, creating a large-scale synthetic dataset (15,600 images, 300 reconstructions, 50 objects) with perfect ground truth using VIRTUS-FPP. Our benchmarking reveals four architectures achieve similar performance (58-77 mm RMSE) with errors approaching 75-95% of the depth range, demonstrating fundamental limitations of single-shot fringe-to-depth mapping without explicit phase information.

The near-identical performance across diverse architectures indicates the information deficit, not model design, limits reconstruction quality. Networks learn semantic shape priors rather than accurate geometry. These results strongly motivate hybrid approaches combining traditional phase-based FPP with learned refinement.

Future directions include: (1) Phase-guided learning using wrapped/unwrapped phase maps, (2) Sim-to-real transfer via domain adaptation leveraging VIRTUS-FPP’s digital twin capability, (3) Multi-view fusion using 6 viewpoints per object, (4) Task reformulation for post-processing traditional reconstructions, (5) Dataset expansion to challenging materials and lighting with domain randomization,³¹ and (6) Uncertainty quantification through probabilistic deep learning.

By providing comprehensive synthetic data and standardized evaluation protocols, this work establishes a foundation for systematic, data-driven FPP research, enabling development of robust systems for manufacturing, biomedical imaging, and automated inspection.

ACKNOWLEDGMENTS

We thank Iowa State University for access to computational resources.

REFERENCES

- [1] Haroon, A., Lakshman, A., Mundy, M., and Li, B., “Autonomous robotic 3d scanning for smart factory planning,” in [*Dimensional Optical Metrology and Inspection for Practical Applications XIII*], **13038**, 110–118, SPIE (2024).
- [2] Lakshman, A., Delzendehrooy, F., Balasubramaniam, B., Kremer, G. E., Liao, Y., and Li, B., “Corrosion characterization of engine connecting rods using fringe projection profilometry and unsupervised machine learning,” *Measurement Science and Technology* **35**(8), 085021 (2024).
- [3] Lakshman, A., Huang, Y., Bussey, W., Liu, L., and Li, B., “Characterizing the 3-dimensional printability of alginate–gelatin and nanocellulose gels via fringe projection,” *Advanced Devices & Instrumentation* **6**, 0116 (2025).
- [4] Zhang, S., [*High-Speed 3D Imaging with Digital Fringe Projection Techniques*], CRC Press, 1st ed. (2016).
- [5] Geng, J., “Structured-light 3d surface imaging: a tutorial,” *Advances in optics and photonics* **3**(2), 128–160 (2011).
- [6] Zuo, C., Qian, J., Feng, S., Yin, W., Li, Y., Fan, P., Han, J., Qian, K., and Chen, Q., “Deep learning in optical metrology: a review,” *Light: Science & Applications* **11**(1), 39 (2022).
- [7] Zhang, S., “Recent progresses on real-time 3d shape measurement using digital fringe projection techniques,” *Optics and lasers in engineering* **48**(2), 149–158 (2010).
- [8] Zheng, Y., Wang, S., Li, Q., and Li, B., “Fringe projection profilometry by conducting deep learning from its digital twin,” *Optics Express* **28**(24), 36568–36583 (2020).
- [9] Zhu, X., Han, Z., Yuan, M., Guo, Q., Wang, H., and Song, L., “Hformer: Hybrid convolutional neural network transformer network for fringe order prediction in phase unwrapping of fringe projection,” *Optical Engineering* **61**(9), 093107–093107 (2022).
- [10] Wang, F., Wang, C., and Guan, Q., “Single-shot fringe projection profilometry based on deep learning and computer graphics,” *Optics Express* **29**(6), 8024–8040 (2021).
- [11] Ikeda, K., Usuki, T., Kurita, Y., Matsueda, Y., Koyama, O., and Yamada, M., “Deep-learning-assisted single-shot 3d shape and color measurement using color fringe projection profilometry,” *Optical Review*, 1–12 (2025).
- [12] Balasubramaniam, B. and Li, B., “Single shot 3d shape measurement of non-volatile data storage devices,” in [*International Manufacturing Science and Engineering Conference*], **87240**, V002T06A010, American Society of Mechanical Engineers (2023).
- [13] Wang, K., Kemao, Q., Di, J., and Zhao, J., “Deep learning spatial phase unwrapping: a comparative review,” *Advanced Photonics Nexus* **1**(1), 014001 (2022).
- [14] Yan, K., Yu, Y., Huang, C., Sui, L., Qian, K., and Asundi, A., “Fringe pattern denoising based on deep learning,” *Optics Communications* **437**, 148–152 (2019).
- [15] Nikolenko, S. I. et al., [*Synthetic data for deep learning*], vol. 174, Springer (2021).
- [16] De Melo, C. M., Torralba, A., Guibas, L., DiCarlo, J., Chellappa, R., and Hodgins, J., “Next-generation deep learning based on simulators and synthetic data,” *Trends in cognitive sciences* **26**(2), 174–187 (2022).
- [17] Ueda, K., Ikeda, K., Koyama, O., and Yamada, M., “Fringe projection profilometry system verification for 3d shape measurement using virtual space of game engine,” *Optical Review* **28**(6), 723–729 (2021).
- [18] Zhang, Q., Xing, M., Li, H., Li, X., and Wang, T., “Measurement simulation system of fringe projection profilometry based on ray tracing,” *IEEE Access* **11**, 89616–89624 (2023).
- [19] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L., “ImageNet: A large-scale hierarchical image database,” in [*2009 IEEE conference on computer vision and pattern recognition*], 248–255, IEEE (2009).
- [20] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L., “Microsoft COCO: Common objects in context,” in [*Computer Vision–ECCV 2014: 13th European Conference*], 740–755, Springer (2014).

- [21] Haroon, A., Lakshman, A., Balasubramaniam, B., and Li, B., “Virtus-fpp: Virtual sensor modeling for fringe projection profilometry in nvidia isaac sim,” (2025).
- [22] Zapico, P., Meana, V., Cuesta, E., and Mateos, S., “Optical characterization of materials for precision reference spheres for use with structured light sensors,” *Materials* **16**(15), 5443 (2023).
- [23] Ou, J., Xu, T., Gan, X., He, X., Li, Y., Qu, J., Zhang, W., and Cai, C., “Comparative analysis on the effect of surface reflectance for laser 3D scanner calibrator,” *Micromachines* **13**(10), 1607 (2022).
- [24] Calli, B., Singh, A., Bruce, J., Walsman, A., Konolige, K., Srinivasa, S., Abbeel, P., and Dollar, A. M., “Yale-cmu-berkeley dataset for robotic manipulation research,” *The International Journal of Robotics Research* **36**(3), 261–268 (2017).
- [25] NVIDIA, “Physical AI Spatial Intelligence Warehouse Dataset.” <https://huggingface.co/datasets/nvidia/PhysicalAI-Spatial-Intelligence-Warehouse> (2025). Accessed: 2026-01-12.
- [26] Sansoni, G., Trebeschi, M., and Docchio, F., “Three-dimensional vision based on a combination of gray-code and phase-shift light projection: analysis and compensation of the systematic errors,” *Applied optics* **38**(31), 6565–6573 (1999).
- [27] Ronneberger, O., Fischer, P., and Brox, T., “U-Net: Convolutional networks for biomedical image segmentation,” in [*Medical Image Computing and Computer-Assisted Intervention (MICCAI)*], 234–241, Springer (2015).
- [28] Isola, P., Zhu, J.-Y., Zhou, T., and Efros, A. A., “Image-to-image translation with conditional adversarial networks,” in [*Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*], 1125–1134 (2017).
- [29] Wang, T.-C., Liu, M.-Y., Zhu, J.-Y., Tao, A., Kautz, J., and Catanzaro, B., “High-resolution image synthesis and semantic manipulation with conditional GANs,” in [*Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*], 8798–8807 (2018).
- [30] Guo, Y., Fang, T., Cui, Z., and Stouffs, R., “A dual-aspect evaluation framework for architectural-like plan generation via pix2pix series algorithms,” *Frontiers of Architectural Research* (2025).
- [31] Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., and Abbeel, P., “Domain randomization for transferring deep neural networks from simulation to the real world,” in [*2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*], 23–30, IEEE (2017).

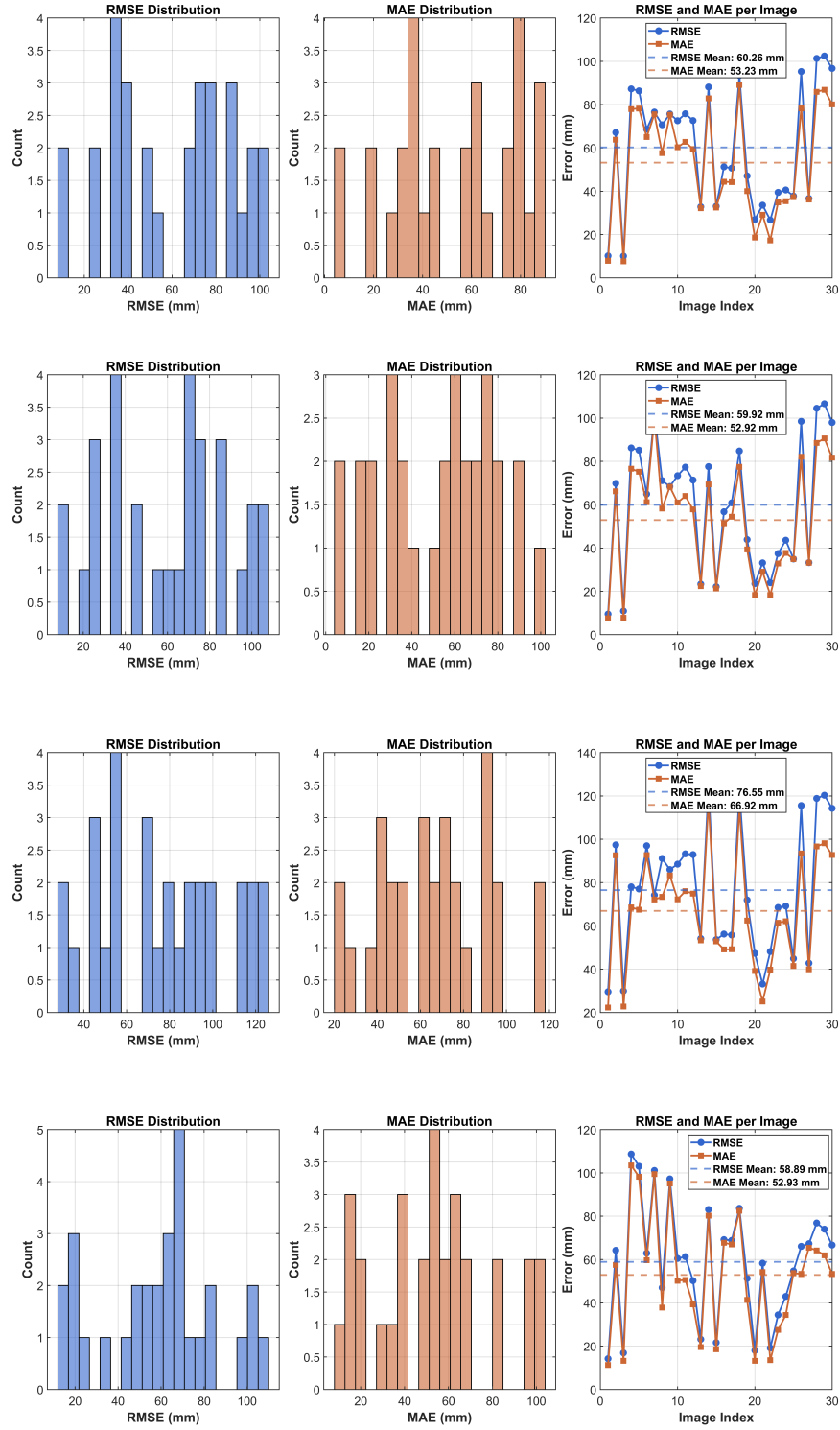


Figure 2. RMSE and MAE distributions and per-sample errors for all four models. Left: error distributions, center: RMSE/MAE distributions, right: per-sample errors with mean lines. Rows: UNet, Hformer, ResUNet, Pix2Pix. Error curves track closely for Pix2Pix, Hformer, and UNet.

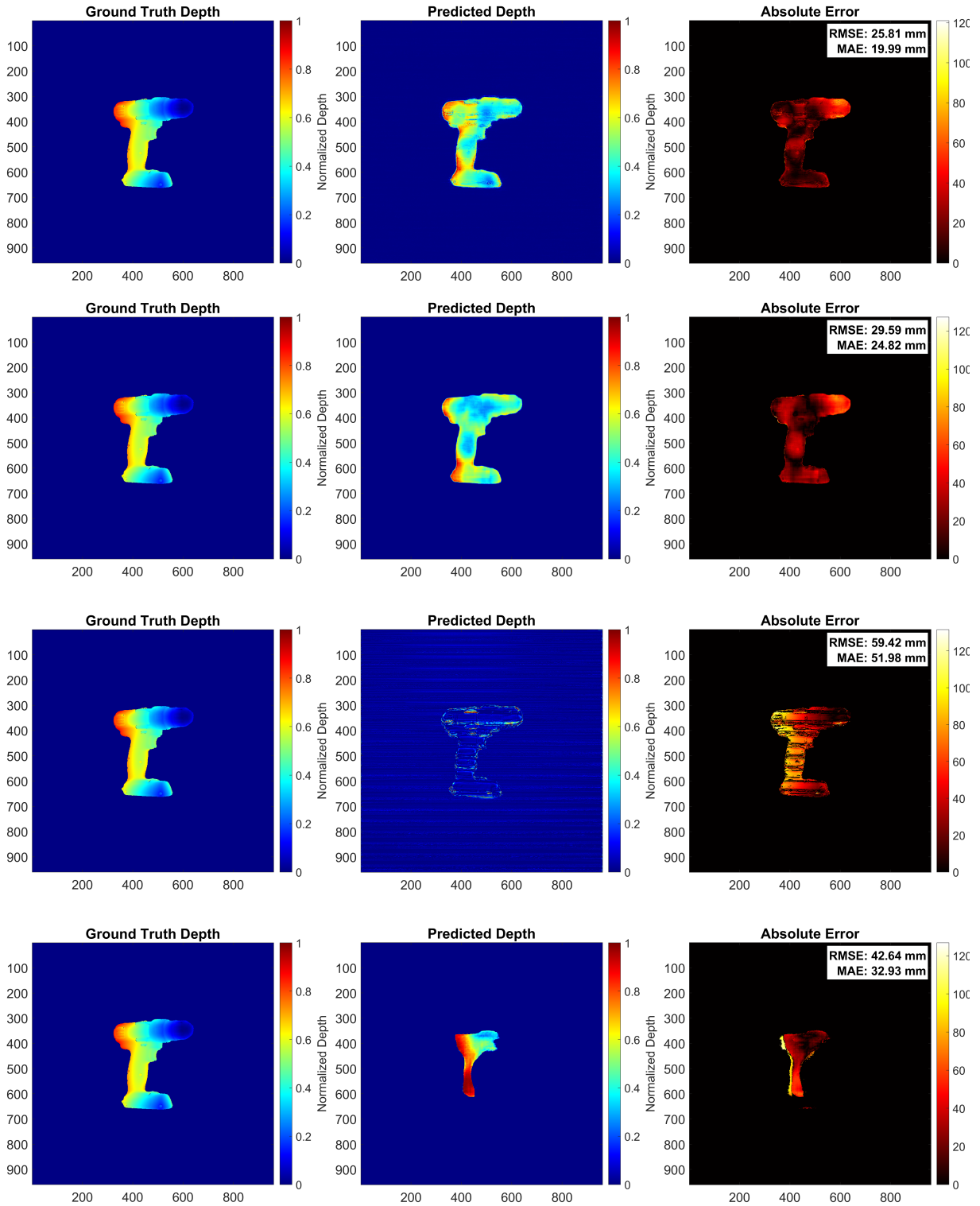


Figure 3. Qualitative results for one test object: From left to right we have the ground truth normalized depth, prediction normalized depth, and absolute error. From top to bottom, the models are UNet, Hformer, ResUNet, and Pix2Pix. Models overall capture coarse shape but fail on fine details and accurate depth prediction.