# Geometric Deep Learning
## for Camera Pose Prediction, Registration, Depth Estimation, and 3D Reconstruction

Xueyang Kang

ORCID: 0000-0001-7159-676X

**Doctor of Philosophy**
**September 2025**

Faculty of Engineering and Information Technology

Submitted in fulfillment of the requirements for the awarded degree of Doctor of Philosophy at the University of Melbourne

# Geometric Deep Learning For Camera Pose Prediction, Registration, Depth Estimation, and 3D Reconstruction

Xueyang Kang[1]

KU Leuven, Belgium

The University of Melbourne, Australia

September 3, 2025

[1]Email: kangxueyang@126.com

# Preface
# -Acknowledgements

## Preface

I sincerely thank my collaborators and mentors for their invaluable support, guidance, and encouragement throughout this journey. I am especially grateful to my parents for their unwavering belief in me, both spiritually and financially. Their love and support have been a constant source of strength.

A special thanks to Ariel, Henry, Alex, Cipher, and Iacopo, the undergraduate and postgraduate students who collaborated with me during my Ph.D. Their enthusiasm and dedication played a crucial role in advancing my research.

I extend my deepest gratitude to Prof. Patrick Vandewalle (KU Leuven) and Prof. Kourosh Khoshelham (the University of Melbourne) for shaping and supporting my joint Ph.D. program. Their efforts in bridging research resources between Leuven and Melbourne provided me with an invaluable cross-disciplinary research experience, integrating theory with practical applications.

I am also grateful to Dr. Gong Dong, Dr. Jin Wanxin, and Dr. Wang Bing for their invaluable guidance in academic writing. Their mentorship helped me meet the quality standards required for top-tier venues in computer vision and robotics, and I have learned tremendously from these talented researchers.

Finally, I express my profound appreciation to Prof. Matthias Niessner (Technical University of Munich), who mentored me during my first Ph.D. year. His guidance was instrumental in shaping my early research on 3D deep learning, and his pioneering work inspired my academic idea growth. I am deeply grateful for his mentorship and for the knowledge I gained from his exceptional research contributions.

Thank you all for your support in my joint Ph.D. journey.

## Abstract

Modern deep learning developments create new opportunities for 3D mapping technology, scene reconstruction pipelines, and virtual reality development. Despite advances in 3D deep learning technology, direct training of deep learning models on 3D data faces challenges due to the high dimensionality inherent in 3D data and the scarcity of labeled datasets. Structure-from-motion (SfM) and Simultaneous Localization and Mapping (SLAM) exhibit robust performance when applied to structured indoor environments but often struggle with ambiguous features in unstructured environments. These techniques often struggle to generate detailed geometric representations effective for downstream tasks such as rendering and semantic analysis. Current limitations require the development of 3D representation methods that combine traditional geometric techniques with deep learning capabilities to generate robust geometry-aware deep learning models.

The dissertation provides solutions to the fundamental challenges in 3D vision by developing geometric deep learning methods tailored for essential tasks such as camera pose estimation, point cloud registration, depth prediction, and 3D reconstruction. The integration of geometric priors or constraints, such as including depth information, surface normals, and equivariance into deep learning models, enhances both the accuracy and robustness of geometric representations. This study systematically investigates key components of 3D vision, including camera pose estimation, point cloud registration, depth estimation, and high-fidelity 3D reconstruction, demonstrating their effectiveness across real-world applications such as digital cultural heritage preservation and immersive VR/AR environments.

The first research project of this dissertation introduces a vision-based camera pose tracking system for robust camera orientation estimation in natural environments, particularly for UAV-based imaging and data collection. The proposed method uses geometric cues of skylines and ground planes from nature to enhance orientation stability and mitigate motion drift of predicted orientation. By integrating a lightweight ResNet backbone and an adaptive particle filter, the system achieves real-time performance on embedded hardware, outperforming IMU-based solutions in stability and robustness against drift.

The second part focuses on point cloud registration, a fundamental problem in 3D vision. To address the limitations of conventional feature-based registration methods, this work introduces a 2D surfel-based $\mathbf{SE(3)}$-equivariant deep learning framework. The model learns robust point cloud alignment

by leveraging surfel representations extracted from RGB-D data or LiDAR scans. The proposed method demonstrates superior and robust registration accuracy through extensive evaluations of indoor and outdoor datasets, particularly in low inlier-to-outlier ratio input point clouds, making it highly applicable for robotics and mobile cases.

The third part of this dissertation explores depth estimation, a key component for generating dense 3D reconstructions. Specifically, it investigates depth prediction from focal stack images, which infers depth information based on focus/defocus cues. A novel Transformer-based network, FocDepth-Former, is proposed, integrating self-attention mechanisms with an LSTM-based recurrent module to handle focal stacks of arbitrary lengths. Compared to traditional CNN-based methods, the proposed approach achieves state-of-the-art performance across benchmark datasets, significantly improving depth prediction accuracy while reducing reliance on large-scale focal stack training data.

The final part presents a high-fidelity 3D reconstruction framework using implicit Signed Distance Fields (SDF). The wavelet-transformed depth feature is used to condition the implicit SDF model to preserve fine-grained geometric details in implicit geometry representations. By integrating wavelet-transformed depth features with a fusion in triplane latent space, the model achieves superior accuracy and detail preservation in reconstructed 3D surfaces. Extensive experiments on DTU, Tanks, Temples, and cultural heritage datasets validate the effectiveness of this approach, demonstrating its applicability to both small-scale objects and large architectural buildings.

Through extensive experimentation on public datasets, the four proposed methods consistently outperform state-of-the-art techniques across individual 3D vision tasks. This dissertation highlights how geometric deep learning bridges the gap between traditional geometry-based methods and data-driven deep learning approaches, driving advancements in 3D vision research and its practical applications in real-world scenarios, such as digital cultural heritage.

# List of Figures

iv

vi

vii

# List of Tables

# Chapter 1

# Introduction

3D vision is a rapidly evolving field that addresses real-world challenges such as pose estimation for localization, point cloud registration for global mapping, depth estimation from images to generate dense point clouds, and 3D reconstruction from casually captured photos. The primary goal of 3D vision is to infer structure and recover geometry from raw images and laser scan data. However, the reconstruction process is complex, typically involving multiple modular components, including pose estimation between frames, dense point cloud generation from depth maps, and establishing accurate correspondences between frames.

## 1.1 Motivation

Improving 3D vision techniques for real-world applications can create great profits for society, helping automate tasks that used to be expensive, error-prone, inefficient, or just too computationally heavy. For instance, predicting camera position and orientation, usually known as camera pose estimation is very important for applications like self-driving cars, augmented reality, and robot vision, where you need to know exactly where you are in dynamic, challenging environments for navigation. To tackle this, we have built a vision-based system that tracks camera orientation and stabilizes camera motion by using geometric hints from the geometric primitives of the skyline and ground plane, then further fused by IMU prediction by the adaptive particle filter constrained on the manifold. This makes pose estimation much more reliable and accurate in natural settings, like drone surveys, outdoor mapping, or autonomous driving.

Another key challenge in 3D vision is point cloud registration, which is crucial for tasks like mapping cities or digital twins of real properties, helping

robots perceive or inspect industrial components We have come up with a surfel-based registration framework that uses $\mathbf{SE}(3)$-equivariant features to handle input data noise and improve the learning efficiency and generalization for both translation and rotation. This improves alignment accuracy in real-world 3D scans, which is a key technique for creating digital twins, inspecting infrastructure, or automating 3D modeling through autonomous alignments of local scan data.

Depth estimation from a focal stack is a great solution for situations where LiDAR is too costly to use or camera motion is constrained in some workspace settings, like robotic grasping of small objects, or medical and biological imaging. We have developed a Transformer-based network called FocDepthFormer, which can predict depth from focal stacks of any length, making it much more flexible for different stack imaging setups. This opens up exciting possibilities for 3D photography on smartphones, microscopy, and quality control in industries where depth sensing is key but depth sensor hardware is limited.

3D reconstruction is also a critical technique for things like virtual reality, self-driving, preserving cultural heritage, and urban planning. To help boost the creation efficiency of digital 3D assets, we have created a wavelet-feature-conditioned implicit SDF model that uses multi-scale wavelet-transformed depth features to capture those geometry details more accurately for the reconstruction of meshes with fine-grained details. This leads to better, more detailed 3D models, which is a huge help for digitizing historical sites, building virtual worlds, or digital manufacturing.

By introducing various kinds of geometric constraints or geometry prior to deep learning models, our research aims to create 3D vision solutions that are more reliable, scalable, efficient, and generalizable. We are working to boost automation, cut down on manual efforts, and promote the use of 3D deep learning models with robust and good performance, to open up new possibilities across fields as diverse as robotics and the creative industries.

## 1.2 Research Problems & Objectives

This thesis explores how to develop geometric deep learning models for camera pose estimation, point cloud registration, depth estimation, and 3D reconstruction, handling the key shortcomings of both traditional and deep learning-based approaches in complex real-world environments, which may be unstructured and natural. While deep learning offers a promising path by directly learning 3D geometric representations, existing approaches remain limited by high computational costs, the need for large-scale 3D data,

and heavy reliance on data augmentation. These challenges stem from three main factors: the absence of structured geometric priors, the difficulty of training in high-dimensional 3D feature spaces, and the lack of robustness to real-world uncertainties. To address these gaps, this thesis proposes hybrid frameworks that integrate traditional geometric constraints with modern deep learning architectures, aiming to improve accuracy, generalization, and efficiency across several core 3D vision tasks.

Given that multi-view 3D reconstruction is a complex system composed of interconnected modules—such as camera pose estimation, depth prediction, and feature correspondence matching—this thesis decomposes the overall problem into a series of subtasks. Each subtask is examined in depth to develop targeted solutions, which are then progressively combined to optimize the reconstruction process. Rather than relying solely on an end-to-end black-box deep learning model, the proposed approach emphasizes subtask-level optimization, ensuring both interpretability and robustness. Finally, the concluding chapter presents the integration of these techniques into a complete 3D reconstruction framework, demonstrating how the modular solutions contribute to the performance of the overall system.

**Camera pose estimation for images using natural geometry cues and manifold constraints.** Traditional camera pose estimation depends on point feature matching and iterative optimization techniques, such as Random Sample Consensus (RANSAC) [57]. However, these approaches perform poorly in situations where visual ambiguity and noisy input lead to a low ratio of inliers to outliers, which makes pose-solving fail.

Deep learning-based pose regression models have been proposed as an alternative, but they still face significant challenges, particularly in accurately predicting rotation due to the complex nature of geometric transformations. One major limitation is that these models lack explicit constraints on the underlying geometric structure and reliable reference cues, making it difficult for them to predict accurate poses in the long-term run due to pose drift. How to leverage some new cues from nature and how to use the geometric constraints along with deep learning techniques to create a hybrid system for robust and accurate pose estimation.

Thus, the goal of Chapter 3 is to design a robust and lightweight camera pose prediction model capable of running on mobile devices, leveraging natural reference cues extracted from images and robust orientation constraints for improved stability and accuracy of camera pose prediction.

**Point cloud registration using 2D Surfel-based equivariance constraint.** Standard point cloud registration techniques align 3D point clouds by matching keypoint features. However, these methods struggle when the features are sparse, noisy, or ambiguous, making alignment unreliable and

challenging. While recent learning-based approaches attempt to learn geometric representation and spatial transformations directly, they often overlook the impact of uncertainty factors and rotation feature representations, resulting in poor generalization to noisy and unseen data, particularly for the input point cloud scans with small overlaps.

To address these challenges, Chapter 4 presents a surfel-based registration method that incorporates $\mathbf{SE}(3)$-equivariant surfel constraints into deep learning architectures, ensuring robust and generalizable point cloud alignment, particularly for input scans with small overlaps or a high ratio of outliers.

**Depth prediction from the focal stack using focal geometry Constraint.** Depth estimation from a focal stack takes advantage of the way focus distance changes with scene depth. However, traditional depth from focal stack methods are limited by the requirement for a fixed number of input images, making them impractical for many real-world applications where image sequences vary in length, and both the training and testing are very inefficient. Furthermore, most existing models are built on convolutional architectures, which have limited receptive fields to capture global multi-scale features.

Chapter 5 proposes a Transformer-based adaptive depth estimation model, which overcomes this limitation by leveraging self-attention mechanisms to extract long-range focus cues in the focal stack and perform latent fusion by LSTM to handle varying input lengths. By embedding focal geometry information into the training loss, this approach allows the model to efficiently process an arbitrary number of input images, significantly improving the performance of depth from the focal stack.

**3D reconstruction using implicit SDF with wavelet feature-based prior.** 3D reconstruction from multi-view images is progressing very fast. The key techniques for 3D reconstruction can be categorized into explicit and implicit models, where the explicit Gaussian Splatting model although fast yet struggles with continuous geometry presentation due to the discretized Gaussian blobs, and the Implicit Signed Distance Function-based methods have proven to be effective for 3D reconstruction, as they preserve continuous surface representations. However, current implicit approaches struggle to capture fine-grained geometric details, especially when dealing with multi-scale geometric structures and sharp edges. This limitation arises because existing implicit methods have difficulty encoding high-frequency geometric features due to the limited learning capability of implicit models via MLPs, leading to over-smoothed reconstructions.

To tackle this, Chapter 6 introduces a wavelet-transformed-depth-feature-conditioned implicit SDF model, which incorporates multi-resolution geo-

metric priors using wavelet-based feature decomposition and triplane feature fusion. By leveraging the features encoded in different frequency bands, the proposed model enhances multi-scale feature representation, preserving sharp geometric details of the input images and improving the accuracy and details of reconstructed surfaces, to overcome the limitation of the implicit model learning capacity.

## 1.3    Contributions

My thesis contribution focuses on combining deep learning with geometric constraints or priors for a range of 3D vision tasks, introducing new ways to represent geometric features and guide the learning process more efficiently for 3D challenges. The primary focus areas of this thesis are camera pose estimation, point cloud registration, depth prediction, and 3D reconstruction.

The first three tasks—camera pose estimation, point cloud registration, and depth prediction—form the essential building blocks of a robust reconstruction pipeline. Camera pose estimation provides accurate geometric alignment across views, point cloud registration ensures consistency in merging multiple scans of data, and depth prediction supplies dense geometric information to recover detailed scene structure. By addressing each of these subtasks individually with hybrid deep learning–geometric approaches, the thesis establishes strong foundations that directly feed into the final 3D reconstruction system.

The contributions of the early chapters, therefore, not only advance the state of the art within their respective domains but also collectively support the integration of a modular, optimized 3D reconstruction framework presented in the concluding chapter. This modular-to-system perspective highlights how solving core subtasks with principled designs can lead to a more accurate, generalizable, modular, and efficient solution for the broader reconstruction challenge. In the following sections, we break down the key methodological contributions of each chapter.

**Camera pose estimation for images using natural geometry cues and manifold constraints.** Chapter 3 presents a new vision-based orientation tracking and fusion algorithm for robust camera pose estimation in natural outdoor environments, particularly for UAV-based wild investigation applications. The approach leverages a lightweight ResNet-18 backbone, deployed on an embedded Jetson Nano, to perform real-time binary segmentation of ground and sky regions. By utilizing the skyline and ground plane as reference cues, a natural geometric primitive-based camera pose estimation framework is developed to enhance visual tracking under challenging

natural conditions. Additionally, an adaptive particle filter operates on a multi-resolution manifold surface, enabling the flexible fusion of orientation estimates from both vision-based cues and IMU data, to make the camera pose tracking more robust.

**Point cloud registration using 2D surfel-based equivariance constraint.** Chapter 4 introduces a novel point cloud registration framework that leverages **SE**(3)-equivariant 2D Gaussian surfel features to improve point cloud registration accuracy and robustness. By representing local surface geometry as learned 2D Gaussian-based surfels, the method preserves explicitly **SE**(3) equivariance for learning, ensuring robustness to rigid transformations. The encoder is an adapted E2PN encoder [287] to take in both point position and orientation, improving registration performance by learning both point position and orientation representations in noisy and partial point cloud scenarios. Additionally, structured huber loss refines pose regression, enhancing robustness to outliers.

**Depth prediction from the focal stack using focal geometry constraint.** In this part, a novel Transformer-LSTM-based network for depth estimation from focal stacks, addresses limitations in conventional CNN-based depth from focal stack approaches. The model integrates a vision Transformer encoder to capture local spatial feature cues and an LSTM-based recurrent module to aggregate stack cues along the main stack dimension, which can be arbitrary lengths. This design overcomes the constraints of fixed-stack-size processing of conventional methods, improving flexibility and generalization to diverse focal stacks. Additionally, a multi-scale convolutional encoder extracts fine-grained focus/defocus features, enhancing depth prediction accuracy by learning more accurate feature representations.

**3D reconstruction using implicit SDF with wavelet feature-based prior.** This chapter introduces a wavelet-transformed depth feature to condition the implicit SDF model for high-fidelity 3D reconstruction, addressing the problem of fine-grained geometric detail loss in deep implicit representations. By leveraging a pre-trained wavelet autoencoder trained on sharp depth maps, our approach extracts multi-scale wavelet-transformed features that efficiently capture high-frequency geometric details like edges. These features are aligned via triplane projection and fused with implicit triplane features, enhancing surface accuracy while preserving complete structures. A hybrid UNet-based fusion module further refines SDF predictions, leading to more precise isosurface extraction.

Overall, this thesis pushes geometric learning in 3D vision applications forward, tackling current challenges in various 3D tasks. Combining traditional geometric constraints or prior knowledge with deep learning helps to boost 3D vision models that are efficient, accurate, generalizable, and robust

Figure 1.1: Chapter relationships and outline structures of dependencies, focusing on the chapters with novel contributions.

for real-world applications like robotics, augmented reality, and self-driving tech.

## 1.4 Thesis Outline

This thesis is structured into eight chapters, including the four main projects from Chapter 3 to Chapter 6.

In the beginning, we lay out the overall research problem and define the objectives of each chapter, giving readers a big-picture view of the thesis and the limitations of existing methods. This is followed by a summary of our contributions to addressing these challenges. Chapter 2 provides a high-level overview of related work across the core problems, complementing the literature reviews in the individual chapters and helping readers contextualize the research questions within current trends.

The subsequent chapters are organized around three foundational subtasks—camera pose estimation, point cloud registration, and depth prediction—which represent the key components of a 3D reconstruction pipeline. Each subtask is studied in depth with tailored frameworks that integrate deep learning and geometric priors, yielding more accurate, robust, and efficient solutions. These contributions are not isolated but are designed to build toward the final system: the last chapter presents the complete 3D reconstruction framework, where the techniques developed in the earlier chapters are integrated and evaluated as part of a unified solution. This modular-

to-system progression illustrates how targeted advances in the core subtasks collectively address the broader challenge of efficient and reliable 3D reconstruction.

From Chapters 3 to 6, each chapter builds upon the previous one, progressively increasing complexity from basic pose estimation to full 3D reconstruction. The step-by-step presentation demonstrates the strength of uniting geometric constraints with deep learning for 3D vision tasks. A typical 3D reconstruction pipeline consists of several modular stages: starting with camera pose estimation to align input images in a global map frame, followed by dense depth prediction and point cloud unprojection. The resulting point clouds, initially in local frames, are then registered to form a global point cloud. Given this modular structure, we decompose the reconstruction system into distinct components for focused analysis and discussion. This organization also highlights that each module can be optimized independently, rather than relying solely on end-to-end training. Modular optimization is often more practical from an engineering perspective and can effectively improve the overall quality of the final 3D reconstruction when integrated into the complete pipeline. Consequently, the overall structure of this thesis is organized according to this modular pipeline optimization strategy. The structure outlined in Figure 1.1 reflects this design: each chapter builds upon the previous one to form a coherent research progression, with inter-block dependencies marked by arrows to indicate how each component supports the next. Specifically, Chapter 3, which treats the subject of image pose estimation, establishes the relative spatial relationship among the frames and provides 3D reconstruction from camera poses. Chapter 4, in the case of point cloud registration, initializes 2D surfels from the 3D point cloud and predicts the relative transformation from source to target scans. Chapter 5 for depth estimation facilitates dense point cloud unprojection in the local camera frame. Finally, Chapter 6 for 3D reconstruction integrates image poses, registered point clouds to a globally consistent map, and predicted depth to generate high-fidelity 3D reconstructions. Every chapter serves as a building block upon which the later pipeline uses the earlier module to create a more complex framework for 3D vision tasks.

Chapter 7 moves from the theoretical research side of things to real-world applications, demonstrating how the methods I have developed contribute to digital cultural heritage preservation and robotic perception or other fields, and it evaluates the expected social benefits, economic values, and target business customers. This chapter explores practical deployment examples and evaluates the potential effectiveness of these methods when applied to real-world scenarios. Finally, Chapter 8 concludes the thesis by summarizing the key findings of the thesis and each project chapter, highlighting the main

contributions again, and reflecting on the limitations of the approaches. It also outlines a plan for future research, suggesting potential improvements, like adapting these methods to other fields, or further enhancing model performance through novel contribution ideas.

# Chapter 2

# Related Work

Geometric Deep Learning (GDL) has emerged as a powerful technique for integrating deep learning with geometric representation or constraints, enabling more effective and efficient solutions for 3D vision tasks such as registration, depth estimation, and reconstruction. This chapter reviews prior work relevant to these topics, progressively linking traditional and modern approaches to provide a structured understanding of the 3D field.

We begin by discussing the foundations of Geometric Deep Learning, clarifying its evolution from classical deep learning applied to non-Euclidean domains to the broader incorporation of geometric constraints and representations in neural network models. Following this, we present the development of camera pose estimation, point cloud registration, depth prediction from images, and 3D reconstruction methods in the same order as the content chapters in the dissertation. This provides readers with a high-level understanding of the progression of each research problem, while each content chapter also includes a detailed review of related work and state-of-the-art methods for a more in-depth understanding of each research problem (Chapter 3 to Chapter 6). As 3D reconstruction is the most complex system, relying on the preceding three techniques, we place a lot of emphasis on its technique development timeline. For the 3D reconstruction problem, we examine traditional 3D reconstruction techniques, focusing on Structure-from-Motion (SfM) and Simultaneous Localization and Mapping (SLAM), which have been widely used in vision-based 3D applications for a long time. These methods, while robust, struggle with large-scale, unstructured, or ambiguous data, driving the transition to data-driven learning approaches.

We explore advancements in 3D Deep Learning-based reconstruction, categorizing models into explicit and implicit representations. Explicit models, such as point cloud and voxel-based networks, directly process 3D data but often face scalability challenges and artifacts from raw data. In contrast,

implicit models represent 3D structures as continuous functions, enabling high-fidelity reconstruction—exemplified by neural radiance fields (NeRF) [168] and signed distance functions (SDF). More recently, Gaussian Splatting [117] has significantly accelerated learning and rendering but at the cost of certain geometric representation continuity. This discussion follows the evolution from classical geometric methods to modern data-driven approaches, providing context for contributions of the thesis in leveraging Geometric Deep Learning for robust and efficient 3D vision applications.

## 2.1   Geometric Deep Learning

Geometric Deep Learning (GDL), introduced by Bronstein *et al.* [23], originally refers to deep learning in non-Euclidean spaces. As shown in Figure 2.1, Geometric Deep Learning emerges from combining deep learning models and non-Euclidean space distributions. Here, I adopt a broader definition, including integrating deep learning models with traditional geometric representations or constraints. This work will use Geometric Deep Learning in this generalized meaning. In essence, GDL leverages the prior of geometric primitive representations, such as points, voxels, normals, and mesh surfaces, to train neural networks to learn the general geometric representation. This approach is very promising and evolves very fast in the field of 3D vision, where high-dimensional feature representation of 3D data can be encoded into a general geometry before handling challenging problems. Unlike conventional deep learning methods that operate primarily on data commonly aligned in regular grids (e.g., images and sequences), GDL excels in handling more complex non-Euclidean structures like graphs and manifolds and facilitates the model to learn topological information underlying the input data, to handle the various down-stream tasks.

The conventional deep learning models are trained on Euclidean data structures, such as images and voxels, using convolutional neural networks (CNNs). However, many real-world problems involve data that exist on more complex structures, such as social media networks, biological or chemical molecular structures, and 3D point clouds of deformable objects. Geometric deep learning provides theoretical and practical frameworks for extending neural networks to these domains, enabling new applications and improving existing ones.

Figure 2.1: Geometric deep learning diagram combines various deep learning model structures and geometry constraints.

## 2.2 Camera Pose Estimation

Pose regression using deep learning models has advanced significantly, evolving from early CNN-based or ResNet-based architectures to transformer-based models. The objective is to directly predict the 6-DoF pose of an object while keeping the camera static [256, 156, 6], or to estimate the camera location while the scene remains static [35, 208, 22, 154]. As shown in Figure 2.2, deep learning regression models are applied to both object pose estimation and visual localization. Various model structures have been explored for these tasks.

**CNN-based pose regression.** Early learning-based pose regression approaches leverage convolutional neural networks (CNNs) as encoders to extract spatial features from images and then regress object or camera poses using MLPs. Early works such as PoseNet [116] employed PoseCNN [256] to predict the absolute pose of the camera, demonstrating robustness in challenging environments. CNNs can also be extended to 6D camera regression via 3D surface-to-image alignment regression. However, the inherent limitations of CNNs restrict feature representation accuracy due to their local receptive fields. For further analysis of the limitations of CNNs in pose regression, please refer to the work by Torsten *et al.* [208]. Consequently, more powerful models have been explored to enhance pose estimation by expanding feature receptive fields and improving feature representation after encoding.

**Transformer-based approaches.** Recent advancements in vision transformers (ViTs) [48] have demonstrated strong performance in pose regression by capturing long-range spatial dependencies and global features within the encoder. For instance, an AutoEncoder-based transformer backbone can be used for pose regression [219], or a cascaded transformer can learn feature representations at multiple scales [219]. Transformer-based encoders generalize well to pose regression across single or multiple scenes due to their supe-

Figure 2.2: Deep learning methods for pose regression. (a) PoseNet: a CNN-based model for 6D camera relocalization [116]. (b) PoseCNN: a deep learning approach for object pose prediction [256]. (c) Map-relative pose regression: a model estimating camera pose relative to a pre-built map representation by using Transformer attention of camera intrinsic and feature maps after CNN encoder [35].

rior feature representation capabilities [218]. However, in highly challenging environments, such as natural landscapes where distinctive appearance features are difficult to identify, *e.g.*, mountains, and meadows with repetitive patterns, most deep learning methods struggle to perform effectively.

## 2.3 Point Cloud Registration

Point cloud registration is a fundamental technique in 3D vision, aiming to align two or more point cloud scans within a global mapping frame. Traditional point cloud registration methods rely on iterative optimization to identify sufficient correspondence inliers between input points and subsequently estimate the transformation pose. Deep learning-based approaches, on the other hand, are data-driven and learn to extract geometric features from input points, which can be leveraged for correspondence matching and pose regression. Recent advancements in point cloud registration have focused on leveraging equivariant feature learning to enhance learning efficiency and generalization. Representative state-of-the-art methods are illustrated in Figure 2.3, including the traditional ICP-based approach, multi-level feature correspondence establishment via learned superpoints, and the equivariant

transformer with specialized attention mechanisms designed to learn both equivariant and invariant features.



Figure 2.3: Deep learning methods for point cloud registration. (a) [240]: an iterative approach for point cloud registration across diverse scenes, requiring only a few shared parameters. (b) Geometric Transformer with multi-level feature matching [194]: a hierarchical method that aligns point clouds from coarse to fine to determine the final pose. (c) Equivariant Transformer SE3ET [136]: a registration model incorporating equivariant cross-attention and invariant self-attention for robust alignment.

**Iterative closest point (ICP) and KISS-ICP.** Traditional registration methods, such as Iterative Closest Point (ICP) and its many variants like Generalized-ICP [214], remain widely used due to their simplicity and efficiency, requiring no parameter training. ICP iteratively refines the transformation between two point clouds by minimizing the point-wise distances of matched point pairs, given an initial transformation. KISS-ICP [240] improves registration robustness across diverse scenes by enforcing stability constraints in optimization and using a single system configuration. Despite their efficiency, these methods struggle with noisy data, high outlier ratios, symmetric structures in point clouds, and sparsity, often requiring a good initial transformation guess, which limits their utility in complex scenarios.

In contrast, geometric feature learning-based methods for point cloud registration offer improved generalization capabilities.

**Registration through geometric feature learning.** Deep learning methods [248, 71] have significantly advanced point cloud registration by learning high-dimensional geometric features directly from raw 3D point coordinates.

Early approaches, such as PointNet [191] and its variants, have been applied to registration tasks by extracting global and local point features, improving correspondence matching by searching for matches in high-dimensional feature space to mitigate ambiguity, as demonstrated in PCRNet [206]. Subsequent works further enhance CNN-based geometric feature descriptor extraction, such as FCGF [42], which strengthens learned feature representations for robust registration through a convolutional U-Net combined with hardest-contrastive and hardest-triplet losses.

Additionally, methods like Maximal Cliques [282] improve correspondence establishment by identifying maximal cliques in the point set. Jaesim *et al.* [182] introduce color information into point cloud registration to enhance robustness, while NICP [215] incorporates normal information in addition to point coordinates for improved alignment. Geometric Transformer [194] further utilizes invariant self-attention and equivariant cross-attention to learn geometric feature tokens and match point tokens, enabling robust registration even in small overlap settings.

These methods demonstrate strong performance in handling noisy points and partial data scans; however, they often rely on extensive data augmentation to generalize across different scenes, making the learning process inefficient.

**Equivariant point cloud registration.** Recent works leverage equivariant feature learning to enforce transformation consistency directly within network models, ensuring that the learned features transform correspondingly with the input points. SE(3)-equivariant models, such as E2PN [287], Equi-GSPR [108] and Equivariant Transformers SE3ET [136], guarantee that learned features remain invariant under rigid transformations while maintaining equivariant edge features, leading to more accurate and robust registration. By incorporating specifically designed equivariant convolution kernels or self-attention mechanisms, these models achieve superior alignment accuracy with fewer training samples and exhibit strong generalization even on unseen data.

The evolution from iterative optimization to deep learning and equivariant models highlights the growing research interest in improving accuracy, generalization, robustness, and efficiency in point cloud registration.

## 2.4   Depth Estimation from Images

Depth prediction from images has progressed rapidly with the emergence of image foundation models [58, 115]. Deep learning models have advanced from CNN-based networks to state-of-the-art diffusion models. These approaches

estimate scene depth from a single image by leveraging relative feature scale distribution or stereo input priors constructed from a virtual camera frame, facilitating applications in 3D reconstruction and robotic mapping.



Figure 2.4: Deep learning architectures for depth prediction. (a) Deep Virtual Stereo Odometry: a method leveraging stereo displacements to generate a virtual stereo camera from a single image for depth estimation [269]. (b) Vision Transformer: a transformer-based model for dense depth prediction [197]. (c) Marigold: a diffusion-based approach for depth prediction using a generative model [115].

**CNN-based depth prediction.** Early deep learning methods for depth estimation leveraged convolutional neural networks (CNNs) to extract multi-scale features and infer depth maps from a single image [126, 151, 65] or video [54]. ResNet-based models, such as the pioneering work MonoDepth2 [70], demonstrated that deep networks could learn depth directly from monocular images. Later research, including monocular depth estimation via unsupervised learning [29, 65], improved prediction performance by incorporating geometric or structural constraints. Despite their success, CNN-based models struggled with predicting fine-grained details and generalizing across diverse environments.

**Virtual stereo with left-to-right displacement.** To mitigate the scale ambiguity in monocular depth estimation, researchers introduced virtual stereo methods that synthesize a right viewpoint by learning left-to-right image transformations through network models like DVSO [269]. These methods, often based on disparity estimation, leverage epipolar constraints between left and right views to refine depth predictions. Works like StereoNet [119] improve generalization by utilizing edges as an objective during training

17

refinement. While effective, these models remain sensitive to view occlusions and textureless regions, limiting their robustness in cluttered scenes.

**Vision Transformers for depth estimation.** Recent advancements in vision transformers (ViTs) have significantly improved depth prediction by leveraging Transformer attention mechanisms to capture global features, as seen in ViT-based dense depth prediction [197]. Transformer models outperform traditional CNNs in depth estimation accuracy and robustness while preserving fine-grained details in depth maps. However, their computational cost remains high, posing challenges for deployment on resource-constrained platforms like mobile hardware.

**Diffusion models for depth prediction.** Recent breakthroughs, such as GeoWizard [58], LOTUS [80], and Repurposing Diffusion [115], leverage diffusion-based foundation models, which iteratively refine depth predictions through a generative denoising process during training. Inspired by the denoising diffusion process, these models predict depth by progressively removing pixel-wise noise. Unlike deterministic depth prediction approaches, diffusion-based depth estimation iteratively captures high-frequency details, better reflecting the feature learning process from low to high frequencies. This results in sharper and more accurate depth maps. Recent advancements in diffusion models have enabled the generation of multi-view images that maintain consistency for 3D representation [109].

The progression from CNNs to diffusion models highlights continuous improvements in depth prediction accuracy and robustness by leveraging the power of generative models. This has enabled the development of more generalizable depth estimation models for practical applications in robotics, AR/VR, and autonomous systems when dealing with unseen data.

## 2.5   3D Reconstruction

### 2.5.1   Traditional 3D Reconstruction Methods

Geometric deep learning often depends on large-scale 3D data obtained either through LiDAR scanning [110], which is efficient but costly and less scalable, or multi-view image reconstruction, which offers greater scalability. This dependence highlights the importance of understanding traditional 3D vision techniques for reconstruction, particularly Structure from Motion (SfM) and Simultaneous Localization and Mapping (SLAM). These methods have long served as the processing tool for pose estimation and 3D mapping in computer vision, and they offer robust frameworks for reconstructing 3D scenes from 2D images captured from multiple viewpoints. Even today, these traditional

methods remain the foundational tool for generating camera view poses in NeRF [168] and sparse 3D point cloud for Gaussian Splatting [117].

**Structure from Motion (SfM)**

Structure from Motion (SfM) ([180, 211]) is a photogrammetric technique that estimates 3D structures from 2D image sequences captured at various viewpoints. Such a 3D reconstruction pipeline recovers the 3D structure of the scene visible in the multi-view images. Taking the common 3D reconstruction tool COLMAP ([210, 212]) as an example, the reconstruction process is composed of the following steps:

- **Feature detection and extraction.** Distinctive feature points are first detected in each image, followed by the extraction of their feature descriptors to enable reliable matching across views.

- **Correspondence establishment.** Feature matching is performed by finding correspondences between detected features across multi-view images. A geometric verification step filters out incorrect matches. The verification ensures only geometrically consistent correspondences are retained for accurate relative pose calculation.

- **Camera pose estimation and 3D point reconstruction.** With the established correspondences, the camera poses are estimated using a linear SVD transformation, and 3D points are reconstructed by triangulating matched 2D feature points. This is achieved through incremental structure-from-motion, where new 3D points are continuously added while refining camera intrinsic and extrinsic parameters.

- **Global bundle adjustment.** To enhance pose and 3D map accuracy, a global bundle optimization is performed over all camera poses and 3D points. This optimization process minimizes the projective function of reprojection errors. The resulting sparse 3D point cloud can be further densified through post-processing techniques, to create a point cloud with uniform density.

Figure 2.5 depicts a fundamental Structure from the Motion approach, inferring 3D point locations through the triangulation of multi-view images. The framework can reconstruct multi-scale 3D assets, ranging from tabletop objects to large-scale buildings.

Figure 2.5: Illustration of structure from motion by multi-view images, image source from GLOMAP [181].

**Simultaneous Localization and Mapping (SLAM)**

Simultaneous Localization and Mapping (SLAM) ([236, 61, 112, 263]) is a method employed by robotics and computer vision communities for building the map of an unknown environment incrementally, while at the same time localizing the ego-agent in the constructed map. SLAM algorithms find broad use in autonomous vehicles, drones, and augmented reality systems.

The key components of SLAM include a front-end and a back-end:

- **Front-end.** Collecting raw data from various sensors, such as cameras, LiDAR, and IMU, then feature extraction and matching are performed over the raw sensor data to establish correspondences. Consequently, the position and orientation of the ego-agent are estimated based on the feature correspondences of frames.

- **Back-end.** A mapping module updates the map by fusing new information from the selected feature points. Once the agent has returned to a previously visited location, the triggered loop closure will correct the map misalignment accordingly by global map optimization.

In the context of SLAM classification, the classification relies mostly on the type of front-end sensor employed: (1) Vision SLAM, which utilizes camera

images to aid in feature detection, motion estimation, and map construction; (2) LiDAR SLAM, which employs LiDAR sensors to scan 3D point clouds for scan matching, odometry estimation, and mapping; and (3) Sensor fusion-based SLAM, which fuses measurements from various sensors (*e.g.*, cameras, LiDAR, IMU, GPS, and wheel odometry) for improving accuracy and robustness in front-end and back-end procedures. The two most representative SLAM methods are illustrated in Figure 2.6. Visual SLAM tracks image features (green patch in the left part of 2.6a) to reconstruct sparse 3D points (right part of 2.6a), while LiDAR-based SLAM generates a global grid map along the motion trajectory (Figure 2.6b).

## 2.5.2 3D Deep Learning

While traditional methods like Structure-from-Motion (SfM) ([211]) and Simultaneous Localization and Mapping (SLAM) have given rise to opportunities for 3D reconstruction and mapping, they are subject to severe limitations. These methods fall short in terms of computational feature matching, thus making them ineffective in processing dense 3D data. They are also prone to robustness breakdown under high outlier ratios or feature ambiguities within the input data.

In contrast, 3D deep learning models provide significant advancements by employing a large number of parameters to learn strong and precise geometric representations in parallel computation. These representations are learned through end-to-end loss backpropagation, offering general and strong geometric priors for many 3D tasks.

To provide a better understanding, let us start with a quick review of the recent 3D deep learning architectures.

### Explicit Model for 3D Geometry Learning

Explicit 3D deep learning models such as PointNet [191] and VoxelNet [160] directly process 3D data such as point clouds and voxel grids for object detection, semantic segmentation, detection [111], and 3D reconstruction. Such models learn geometric representations automatically through neural networks under the guidance of explicit supervision signals. They have the limitation of high computational complexity, high memory, and low scalability for processing dense or large data. The effectiveness of their performance is prone to degradation due to sparse or noisy inputs and is sensitive to biases in data representation. It is necessary to overcome these difficulties for the improvement of both the efficiency and robustness of 3D deep learning models in practical applications.

(a) ORB vision SLAM demonstration. Source image from work by Mur-Artal *et al.* [174].



(b) Lidar-based Cartographer SLAM by Google. Source image from work by Hess *et al.* [83].

Figure 2.6: Common SLAM framework: (a) visual SLAM, (b) Lidar-based SLAM.

Explicit deep learning models for 3D reconstruction are illustrated in Figure 2.7. The top row shows a classical recurrent model by [43] that encodes multi-view images using LSTMs in the latent feature space. Another common approach directly takes in the point cloud to predict a shape occupancy map, as shown in Figure 2.7b.

(a) 3D recurrent neural network (3D-R2N2) for 3D shape reconstruction from sequence images. Source from work by [43].



(b) 3D Unet convolution for reconstruction.

Figure 2.7: 3D reconstruction by using a 3D recurrent method (a), and explicit 3D U-Net model (b).

## Implicit Model for 3D Geometry Learning

Implicit deep learning models, such as Neural Radiance Fields (NeRF) by [167] and DeepSDF by [183], represent 3D geometry and appearance through neural networks—typically Multi-Layer Perceptrons (MLPs)—instead of explicit data structures like point clouds or voxel grids. These models learn continuous volumetric representations from sparse, noisy input data by capturing the underlying geometry distribution. This enables them to generate high-fidelity 3D reconstructions and novel view synthesis by optimizing over a latent space encoding geometric information.

A key advantage of implicit models is their ability to handle complex topologies and generate continuous geometric details, and they overcome the resolution limitations of grid-based input representations. However, they also pose challenges, such as high computational demands during training with a large volume of pixel rays to query the radiance field, and the generalization ability of implicit models on diverse datasets is quite weak. With increasing computational power, these challenges are gradually being ignored through scaling law. Overall, implicit deep learning models still have their merits in 3D reconstruction and rendering while maintaining a compact model size.

23

Next, we discuss recent developments in implicit 3D modeling.

**Implicit SDF.** Implicit deep learning models have changed 3D shape representation by modeling shapes as continuous functions. Unlike explicit models that rely on discrete representations such as point clouds, meshes, or voxel grids, this implicit continuous representation, learns 3D features in MLP-based weights, enabling smooth surface reconstructions. I can be scaled to complex environments with a small memory footprint. Implicit models also offer robust invariance to incomplete or noisy data, making them well-suited for robust applications in 3D shape representation and reconstruction.

A pivotal shape reconstruction work in this domain is DeepSDF by [183], which introduced Signed Distance Functions (SDFs) for high-fidelity 3D shape reconstruction. This method represents shape by an implicit SdF model, allowing for smooth feature learning and robust reconstruction given partial inputs.

Building on this foundation, Convolutional Occupancy Networks ([185]) integrated convolutional neural networks to learn occupancy probabilities. Such an occupancy prediction-based model improves the details and accuracy of 3D geometry representation by using multi-scale triplane or volumetric features. This approach efficiently processes volumetric data, enabling high-fidelity reconstructions of complex scenes. This is achieved through grid-based latent feature interpolation in volumes or triplanes.

Next, IF-Net ([38]) introduced implicit function constraints in feature space, which jointly performs 3D shape reconstruction and completion. IF-Net enhanced the ability to reconstruct complex and incomplete shapes, such as partial human body scans, offering a more flexible 3D processing model.

More recently, Neural Shape Deformation Priors ([234]) incorporated learned deformable shape priors into neural networks. Such deformable prior-guided implicit model improves the accuracy of both 3D reconstruction and shape manipulation. This method leverages a hierarchical structure, including global shape templates and local deformations encoded by different MLPs. Such a design can capture inherent shape distributions and dynamic variations more effectively.

In a nutshell, these advancements highlight the rapid evolution of implicit 3D modeling techniques. From the shape DeepSDF to the advanced Neural Shape Deformation Priors, each pushing the limits of 3D geometry reconstruction with increasing complexity and showcasing the rising influence of implicit deep learning in computer vision and graphics.

**Implicit neural radiance field.** The transition from Neural Radiance Fields (NeRF) by [167] to 3D Gaussian Splatting by [117] and the latest 2D Gaussian Splatting by [92] marks a significant evolution in 3D scene representation and rendering, and this development is demonstrated in Figure

Figure 2.8: Implicit SDF for a wide range of reconstruction problems, from shape to room scan, human body, and deformable animation. Source images rearranged from [183, 185, 38, 234].

2.9.

NeRF ([167]) introduced a novel approach to three-dimensional scene representation by using implicit neural networks to represent volumetric geometry and color appearance. This approach enables photorealistic rendering of new viewpoints from a small set of input images by mapping three-dimensional spatial positions and viewing directions to their corresponding color and density values. Unlike implicit Signed Distance Function (SDF) models that are limited to geometry representation, NeRF combines geometric and appearance information in real-time, hence producing high-quality, continuous three-dimensional representations without the need for explicit structures like meshes or point clouds. Extensions of NeRF [144, 166] have improved the rendering of occluded regions and performance under varying lighting conditions, hence expanding its use in virtual reality, augmented reality, and 3D content creation. The compact scene representation provided by NeRF saves memory, hence improving computational efficiency and making it possible to deploy on mobile and edge devices.

Implicit NeRF suffers from complex and time-consuming training, and it often requires extensive GPU resources for large-scale scene reconstruction. To overcome this computing limitation, 3D Gaussian Splatting by [117] is proposed to enhance real-time rendering capabilities without sacrificing rendering quality. This approach represents scenes as a collection of 3D Gaussian blobs, which are optimized efficiently and explicitly through a Gaussian feature representation. A key advantage of 3D Gaussian Splatting is its ability to achieve real-time performance through a specialized CUDA kernel imple-

25

Figure 2.9: Neural Radiance Field (at the left) alongside its extensions: 3D and 2D Gaussian Splatting (at the right). Source images from [167, 117, 92].

mentation, so it is suitable for applications such as localization and mapping. However, Gaussian Splatting represents scenes with discrete Gaussian primitives, enabling fast rendering but lacking the continuity of implicit SDF models. Unlike SDFs, which enforce smooth and coherent surface geometry, Gaussian Splatting may produce disconnected or floating artifacts. Gaussian Splatting struggles with fine-grained topology reconstruction, making it less ideal for precise 3D reconstruction.

The latest advancement, 2D Gaussian Splatting by [92], further refines the 3D Gaussian splatting approach by improving geometric accuracy and rendering efficiency. This method projects 3D Gaussian splats into a 2D representation, and it reduces computational overhead while preserving high visual fidelity. Notably, initializing 3D Gaussians from sparse point clouds remains computationally demanding, making this refinement particularly impactful.

Overall, the transition from NeRF to 3D and even 2D Gaussian Splatting outlines a clear roadmap toward more efficient and accurate scene representations. These advancements pave the way for real-time applications in large-scale environments, enabling more complex tasks in autonomous driving and robotic perception.

**Implicit Simultaneous Localization and Mapping.**

Implicit deep learning SLAM systems achieve good performance by leveraging continuous scene representations that are obtained directly from raw color and depth images, as shown in Figure 2.10. Such systems use neural networks to map spatial coordinate information to continuous functions, thereby enabling real-time localization and mapping and producing high-quality mesh reconstructions without the need for further processing.

A major advantage of implicit models is their ability to learn autonomously

26

Figure 2.10: The NICE-SLAM framework [288] tracks ego-body pose and constructs a 3D mesh in real time (top row). vMap [124] simultaneously reconstructs scene objects and renders the environment, enabling a more interactive 3D scene representation. Image from work by [288, 124]

robust geometric representations through end-to-end training. This allows the implicit model to better handle noisy and incomplete data, thus outperforming traditional SLAM methods under challenging conditions. Additionally, implicit models can fuse multimodal sensor data, such as images and depth maps, in a unified framework, thus the implicit model has good robustness and flexibility to use.

NICE-SLAM ([288]) boosts SLAM (simultaneous localization and mapping) technology by leveraging neural implicit functions for real-time 3D reconstruction and pose tracking given live-streaming images. It supports scalable and robust 3D mapping by learning a continuous indoor scene representation, which easily captures complex details while being robust to noise perturbations. Additionally, its hierarchical feature interpolation improves multi-scale geometry learning, so that it effectively alleviates the oversmoothing issues of one global scene representation by leveraging multilayer perceptrons (MLP) at different feature resolution levels.

Gaussian Splatting SLAM [158] uses Gaussian splats for 3D representation, significantly improving real-time rendering and computational efficiency compared to NICE-SLAM. This approach provides an approximate geometry representation of the environment while simultaneously keeping a fast reconstruction and rendering performance.

Vectorized Object Mapping for Neural Field SLAM ([124]) extends neural field SLAM by integrating vectorized object representations. vMap detects and registers object instances on the fly, enabling complex scene and object-level reconstruction. Combining semantic understanding with interactive scene representations enhances robot manipulation and navigation capabilities [255]. In short, the fusion of traditional geometric techniques with modern deep learning architectures has become an active field of research in recent 3D vision research. Traditional approaches offer interpretability and build upon proven mathematical models but they often struggle with scalability issues and the treatment of complex, unstructured environments. Deep learning, on the other hand, enables data-driven generalization and achieves high-quality reconstruction but often lacks geometric consistency, in addition to high training complexity. This thesis builds on these recent developments, and it aims to leverage Geometric Deep Learning, in combination with the strengths of classical geometric representations, aiming at the development of robust and efficient deep learning models for 3D vision applications.

# Chapter 3

# Camera Pose Estimation for Images Using Natural Geometry Cues and Manifold Constraint

## Abstract

Accurate camera pose estimation from image frames is essential for 3D reconstruction. This chapter [1] presents a vision-based orientation tracking and fusion algorithm for drones operating in natural environments. The system can be used to mitigate motion blur and stabilize camera orientation by taking advantage of the skyline and ground plane as reference signals, ensuring high-quality image capture for downstream 3D tasks. The key contributions

**Author Contributions:**

- **Xueyang Kang**: Idea Design, Methodology, Software, Experiment Validation, Formal Analysis, Data Curation, Writing, Review, and Editing.

- Ariel Herrera: Data Curation and Experiment Validation.

- Henry Lema: Data Curation.

- Esteban Valencia: Review.

- Patrick Vandewalle: Review, Editing, Supervision, and Funding.

include: a) A lightweight ResNet-18 backbone, trained from scratch and deployed on a Jetson Nano, for real-time segmentation of images into binary ground and sky regions. b) A geometry-based framework that utilizes skyline and ground cues for robust visual tracking in challenging outdoor conditions. c) An adaptive particle filter sampling technique on a multi-resolution manifold surface, enabling flexible fusion of orientation estimates from multiple frames. The proposed method was implemented and tested in real-world environments, including rooftop and drone-mounted experiments. The final experimental results demonstrate its robustness in frame-to-frame pose estimation under challenging natural conditions.

## 3.1 Introduction

Pose estimation from image frames is a fundamental step in 3D reconstruction. Typically, a camera follows an orbital trajectory around the target, incorporating pan rotation and elevation changes to ensure sufficient overlap between consecutive frames for reliable feature matching. This structured motion enhances pose estimation accuracy. However, abrupt tilting with minimal height change or pure rotation can severely degrade performance due to the lack of overlap between consecutive frames. Additionally, unintended jitter during motion can introduce blur, making camera stabilization essential for obtaining sharp, high-quality frames for precise pose estimation.

Motion blur degrades image quality, making feature extraction and correspondence matching significantly more challenging, often leading to incorrect matches. Reliable frame-to-frame correspondence is crucial for accurate pose estimation and consistent 3D alignment.

In UAV-based image capture, especially in mountainous regions with unpredictable rotations and motion jitter, conventional IMU-based camera pose tracking methods struggle with long-term drift and noise. To address these challenges, we propose a novel approach to estimate camera orientation by leveraging natural scene cues and a manifold-based fusion strategy. The key contributions of this chapter are as follows.

- A lightweight binary segmentation model trained to classify ground and sky pixels in real-time on an embedded device.

- Geometry-based camera pose estimation method that uses the skyline and ground plane as reference cues to infer rotation angles.

- A nonlinear particle filter with adaptive sampling resolution on the manifold surface to fuse rotation estimates from multiple sensor modal-

ities (IMU and vision-based cues), ensuring robust orientation tracking in challenging outdoor conditions.

Our approach integrates manifold-based constraints to enforce geometric consistency in rotation estimation over a spherical domain. Traditional methods relying solely on high-precision IMUs often suffer from noise drift and instability. In contrast, our fusion framework mitigates these issues through adaptive sampling on the manifold surface, ensuring robust multi-modal integration. Using natural scene features, our method provides stable and accurate pose estimation, even in dynamic and unpredictable environments, demonstrating a significant advancement in geometry-aware deep learning for UAV-based vision tasks.

## 3.2   Related Work

Pose estimation from image frames is a fundamental problem in computer vision, often approached through feature matching, motion estimation, or direct pose regression. Traditional methods rely on handcrafted feature descriptors such as SIFT [17], SURF [188], or optical flow [275] to establish correspondences between consecutive frames. However, these techniques are highly sensitive to motion blur and textureless regions, which degrade tracking performance [159].

To improve robustness, inertial sensors have been integrated into pose estimation pipelines. Some works compensate for motion jitter by fusing IMU data with vision-based tracking, as seen in humanoid robotics [5] and complex motion prediction tasks [9]. A review of motion estimation techniques for video stabilization, many of which are applicable to pose estimation, is provided by Rawat et al. [199].

Recent advances in deep learning have significantly improved pose estimation by leveraging learned priors. Yu et al. [274] introduced a scene representation approach that extracts structured cues such as background feature points, foreground contours, and 3D facial meshes to aid tracking. Li et al. [130] trained a model to estimate visual odometry from raw IMU data, demonstrating the effectiveness of learning-based fusion for motion prediction.

Self-supervised methods have also been explored, where pose estimation is embedded within larger mapping and localization frameworks. Lee et al. [128] proposed a graph-based self-supervised method, while Liu et al. [145] introduced a dense warping field for motion compensation, synthesizing stabilized frames from sequential observations. Choi et al. [39] further refined

pose estimation by integrating optical flow-based motion prediction. Despite their accuracy, these approaches involve significant computational overhead, making them impractical for real-time inference on edge devices.

Pose estimation is crucial in UAV-based applications, particularly for navigation and localization in challenging environments [241, 113]. Several works use skyline tracking [165, 27, 51] to estimate camera orientation, but these methods require careful tuning and perform poorly in non-ideal conditions. Other approaches rely on geometric constraints, such as the five-point algorithm [177] for epipolar geometry estimation or curvature alignment techniques for motion tracking [217]. Moving object detection has been tackled using tracking filters [242] and SIFT-based correspondence matching [139]. However, feature-based methods struggle in natural environments due to spurious correspondences and occlusions, limiting their reliability for pose estimation.

## 3.3    Method

As shown in Figure 3.1, the system is mounted beneath the airplane body and consists of a camera, an IMU, and a barometer. Ensuring stable and accurate camera orientation during flight is crucial, particularly in dynamic and unstructured environments where feature-based tracking often fails due to motion blur and inconsistent correspondences. To overcome these challenges, our approach leverages natural visual cues—such as the skyline and ground plane—to aid pose estimation. The lower part of Figure 3.1 illustrates a binary mask where the skyline and ground region are distinctly segmented, providing reliable geometric references.

A key contribution of our method is the manifold-based adaptive particle filter, which fuses IMU-based orientation estimates with visual cues extracted from the segmented scene. Unlike traditional filtering approaches, which may suffer from noise drift or require precise sensor calibration, our method operates directly on the rotation manifold, ensuring smooth and consistent pose estimation. This fusion not only enhances robustness against sensor noise but also improves long-term orientation estimation stability, enabling accurate camera orientation in challenging aerial conditions.

Rays passing through the red dots, in conjunction with height measurements obtained from the barometer, enable us to accurately determine the 3D position of the ground plane using trigonometric calculations. This geometry relationship is fundamental to our approach, as it allows for precise estimation of the camera's orientation relative to the ground, thereby improving the stability of the captured imagery. The hardware diagram details

the integration of the camera, IMU, and barometer, while the software diagram outlines the algorithmic pipeline that processes the visual and inertial data.

Overall, the capability of such an orientation estimation system to leverage natural environmental features, coupled with advanced geometry deep learning techniques and efficient sensor fusion, represents a significant advancement in UAV-based camera pose tracking in the natural environment. This robust design is particularly well suited for long-term operations in remote and challenging terrains, such as volcanic regions, where traditional pose estimation systems might fail due to noise drift or the lack of high-precision sensors. The integration of manifold-based constraints in our pose estimation approach ensures high precision and stability, making it an ideal solution to improve the quality of captured video and the reliability of surveillance missions.



Figure 3.1: Illustration of pose orientation estimation from image frames on the fixed-wing airplane.

**Hardware.** The hardware design is based on open-source hardware components. The main processing unit is a Jetson Nano, equipped with 2GB of GPU memory and Quad-core ARM A57, connected to IMU, camera, and barometer sensors. An OpenCR 2.0 driver board maps the driving command to control commands for two servo motors.

**ROS Nodes.** The software is composed mainly of three parts: the preprocessing part including the network model and geometry primitive extraction, followed by a tracking module to align the skyline and normal of the ground

Figure 3.2: Open source hardware setup.

plane in the current frame with those in the reference frame. The compensation angles from various pipelines are then fed into the proposed particle filter presented in the Section below to obtain fusion orientations, further as input for the controller to stabilize the camera.



Figure 3.3: Block diagram of the presented orientation estimation algorithm. Circular nodes in pink are signals. Rectangular boxes in yellow are ROS nodes for the algorithm, and the dashed region is the front-end perception part, including tracking of the skyline and ground plane.

### 3.3.1 Perception

The perception part is structured into preprocessing and rotation estimation, where the rotation estimation can be further separated into two pipelines:

34

roll and pitch prediction from skyline tracking and rotation estimation from ground plane tracking. The following subsections follow this structure.



(a) Correspondences based on "SIFT" [150] feature points of neighboring image frames.



(b) HSV image converted from a raw RGB image.         (c) Boundary curve-fitting on the detected Canny edges.

Figure 3.4: Failure case demo using traditional OpenCV pipeline.

We first tried the general computer vision processing pipelines, but they all failed due to spurious feature candidates. As illustrated in the bottom left of Figure 3.4, the similar appearances in the grassland region and cloudy sky all pose a great challenge to correspondence search. In the top of Figure 3.4, many false correspondences are found. Canny edge detection [26] is applied to find the boundary between the ground and sky region on the HSV image; nevertheless, some brightness of the sky is cast onto the grass ground, generating the wrong boundary in the bottom right of Figure 3.4. To tackle this challenging segmentation task, we finally choose the data-driven approach by training a Resnet-18 [259] network on the "Skyfinder" dataset [125] first, followed by a fine-tuning on the self-collected dataset with one hundred images. Binary cross entropy loss is applied for pre-training through 100 epochs and fine-tuning with 30 epochs, respectively. The total training time is less than two hours. Some training data samples, along with ground truth masks, are presented in Figure 3.5. The model is exported into "ONNX" and optimized by "TensorRT" to convert to "FP16" precision

for Jetson Nano deployment. We found the model can achieve above 90% success rate for the segmentation on average; only under some extreme cases like overexposure does the failure happen. Additionally, our use case is for mountainous terrain. The scenario with mirror effects and reflections by water on the ground is not advisable due to the similar color distributions of the sky and the ground.



(a) Sample training images.



(b) Corresponding ground truth masks.

Figure 3.5: Sample images and ground truth masks for training.

### 3.3.2 Skyline Tracking

Starting from the binary mask results gained from the network model, the skyline can be extracted along the boundary direction. The extracted skyline can be further considered as a cue to estimate the roll and pitch of the camera, as shown in the top of Figure 3.6. The boundary points, represented by yellow dots at the bottom of Figure 3.6, can be implemented firstly over the whole reference image. A straight line is fitted as illustrated in red (the bottom of Figure 3.6), using two parameters, slope $m'$ and intercept $b'$ in Equation 1. For the subsequent images, we use a constant angular velocity model derived from skyline tracking of the previous two frames to predict the skyline position in the current frame, like the green line presented at the bottom of Figure 3.6. The assumption is that the camera remains not upside down, due to the airplane maneuverability constraints. The boundary can always be searched along the vertical direction of the predicted skyline (green). In practice, the skyline points are further down-sampled to speed up the search. Once the current sample points of the skyline are attained, the estimated skyline (red) in the current frame can be derived via least-square

at the bottom of Figure 3.6.

$$m'x + b' = y' \tag{3.1}$$

$$mx + b = y \tag{3.2}$$

$$\alpha = \arctan(m) - \arctan(m') \tag{3.3}$$

$$\beta = \arctan\left(\frac{h_1 - c_y}{f_y}\right) - \arctan\left(\frac{h_2 - c_y}{f_y}\right) \tag{3.4}$$

The roll $\alpha$ and pitch $\beta$ angles can be derived from Equations 3, and 4 respectively by subtracting the angles. $c_y$ is half of the image height, and $h_1$ and $h_2$ are the height of the center point of the skyline in the current image frame and reference frame respectively. $f_y$ is the focal length $y$ of the camera.



(a) Reference image (left) and current image frame (right).



(b) Reference mask with skyline (left) and current mask with skyline (right).

Figure 3.6: Segmented images for skyline search.

Roll and pitch have a specific tolerance range to avoid unnecessary operations, so processing is only triggered when the movement is out of this range. Roll angle can be predicted from the slope $m$ of the skyline, followed by pitch estimation, which is on top of the image result after roll compensation. There is a total of three cases in our system, that is pure roll, pure pitch, or both happening simultaneously. Here the height shift resulting from translation is subtracted before rotation processing by barometer readings.

The orientation estimation algorithm is designed to handle roll and pitch angle changes separately. It first checks for significant roll angle changes

37

and applies roll compensation if necessary. After that, it checks for significant pitch angle changes and applies pitch estimation, either on the roll-compensated image or directly on the raw image, depending on whether roll compensation was applied or not. This decoupled approach can be useful in applications where roll and pitch corrections need to be applied independently, such as in image stabilization systems or camera orientation tracking. Although this compensation mechanism is designed for roll and pitch only, it can be easily adapted to the full 3D case with roll, pitch, and yaw angles.

### 3.3.3 Ground Plane Tracking

Ground plane tracking relies on the normal vector of the ground plane, as demonstrated in Figure 3.1. A set of points in the ground region of the binary mask are sampled evenly, followed by a back projection to the camera frame corresponding to Equations 5-11.

$$\mathbf{K} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \tag{3.5}$$

$$\rho_{\mathbf{i}} = [u, v, 1]^T \tag{3.6}$$

$$\mathbf{P_i} = \mathbf{K}^{-1}\rho_{\mathbf{i}}, i \in (1...N) \tag{3.7}$$

$$\mathbf{N_G} = [0, 0, g_z]^T \tag{3.8}$$

$$\cos\theta = \frac{\mathbf{P_i} \cdot \mathbf{N_G}}{||\mathbf{P_i}|| \cdot ||\mathbf{N_G}||} \tag{3.9}$$

$$l_i = \frac{h}{\cos\theta} \tag{3.10}$$

$$\mathbf{P'_i} = l_i\mathbf{P_i} \tag{3.11}$$

The height $h$ is measured from the barometer, and the ray direction passing through the pixel position is $\mathbf{P_i}$ on the left side of Equation 7. $\mathbf{K}$ is the intrinsic matrix obtained from calibration [283]. $\theta$ is derived from the dot product of the gravitational vector and the ray direction. Length scale $l_i$ can be calculated by trigonometry in Equation 10. Finally, the current normal vector $\mathbf{m}$ of the ground plane is shaped by the cross product of points as below:

$$\mathbf{m} = (\mathbf{P'_i} - \mathbf{P'_j}) \times (\mathbf{P'_i} - \mathbf{P'_k}) \tag{3.12}$$

The ground plane tracing mode is only triggered when the camera is over 300 meters above ground so that the variance of uneven grassland can be approximated by a flat plane compared to the height. Next, the rotation

matrix to align the normal vector $\mathbf{m}$ in the current frame and the reference normal $\mathbf{n}$ at the start can be derived as follows.

$$s = \frac{\mathbf{m}}{||\mathbf{m}||} \cdot \frac{\mathbf{n}}{||\mathbf{n}||}, \tag{3.13}$$

$$\mathbf{k} = \frac{\mathbf{m}}{||\mathbf{m}||} \times \frac{\mathbf{n}}{||\mathbf{n}||}, \tag{3.14}$$

$$\mathbf{k}_\times = \begin{bmatrix} 0 & -k_3 & k_2 \\ k_3 & 0 & -k_1 \\ -k_2 & k_1 & 0 \end{bmatrix} \tag{3.15}$$

$$\mathbf{R} = \mathbf{I} + \mathbf{k}_\times + \mathbf{k}_\times^2 \frac{1}{1+s} \tag{3.16}$$

$\mathbf{k}_\times$ is the skew matrix, where non-zero elements are in off-diagonal positions, corresponding to the components of the cross product of $\mathbf{m}$ and $\mathbf{n}$. $s$ is a scale derived from the dot product of two vectors. The rotation matrix is calculated following Rodrigues' rotation formula in Equation 16.

In the end, the 3D Euler angles are retrieved from the rotation matrix according to Equations 17-20, in a right-hand order, "yaw←pitch←roll". Only roll and pitch are used for later fusion.

$$\mathbf{R} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \tag{3.17}$$

$$\alpha = \arctan(r_{32}, r_{33}) \tag{3.18}$$

$$\beta = \arctan(-r_{31}, \sqrt{r_{32}^2 + r_{33}^2}) \tag{3.19}$$

$$\gamma = \arctan(r_{21}, r_{11}) \tag{3.20}$$

### 3.3.4 Adaptive Particle Filter for Pose Fusion

The Particle Filter is easy to implement and applicable in non-linear problems, in particular for positioning [77]. Here, a variant of the vanilla particle filter [265], sampling on the spherical surface adaptively is proposed in the pseudo-code below. In a real configuration, the roll and pitch are virtually constrained by a limit range due to mechanical kinematics, e.g., roll in a range from -45° to 45°.

The general idea behind this adaptive particle filter in Figure 20 is straightforward. The filter mainly comprises three steps: sampling from orientation measurements of the IMU $s_I$ (line 2); sampling according to observation from the Computer Vision (CV) pipeline (lines 8 and 15); resampling proportionally to the updated weight of each particle (line 19). The weight in line 6 is

**Algorithm 1:** Particle Filter on Spherical Surface

---

**Data:** $S_I = (\alpha_I, \beta_I), S_{C_{1,2}} = (\alpha_{C_{1,2}}, \beta_{C_{1,2}}), \Omega_{1,2,3}$.
**Result: Output** $\bar{S}_k = [\alpha, \beta]^T$.

1: **if** new $S^I$ available **then**
2:      Initialize the particle $s_0^{(j)}$ or sampling $s_k^{(j)}$ (on $\Omega_1$) from
       $\mathcal{N}(\tilde{\mu}_I | \mu_I + \omega \delta t, (\sigma_I + b)), j = 1...N$;
3: **end**
4: **if** new $S_{C_1}$ and $S_{C_2}$ both are available **then**
5:      **for** $s_k^{(j)}$ in the particle set of $N$ samples **do**
6:          $\omega_k^{(j)} = \omega_{k-1}^{(j)} \exp{(s_k^{(j)} - S_{C_{1,2}})}$;
7:          **if** $||s_k^{(j)} \text{-} S_{C_{1,2}}||_2 \leqslant \epsilon$ **then**
8:             Sampling $\hat{s}_k^{(1...m)}$ from $\mathcal{N}(\tilde{\mu}_f | \mu_f, \delta_f)$ (on $\Omega_3$)
9:             initialize new weights: $\omega_k^{(1...m)} = \omega_k^{(j)} \exp{(\hat{s}_k^{(1...m)} - \mu_f)}$;
10:          **end**
11:      **end**
12: **end**
13: **if** new $S_{C_1}$ or $S_{C_2}$ is available **then**
14:      **for** $s_k^{(j)}$ in $N$ samples (on $\Omega_2$ resolution) **do**
15:          Sampling from $\mathcal{N}(\tilde{S}_{C_1} | S_{C_1}, \delta_{C_1})$ or $\mathcal{N}(\tilde{S}_{C_2} | S_{C_2}, \delta_{C_2})$, same as
            line 6-9;
16:      **end**
17: **end**
18: $\hat{s}_k^{(1...m)} \cup S_\Omega$;
19: Resampling $s_k^{(j)}$ according to $\omega_k^{(j)}$;
20: $\bar{S}_k = \sum_{j=1}^{N+m} \omega_k^{(j)} \hat{s}_k^{(j)}$;

---

derived from a normal distribution, as a function of the square root of the angular distance, which is between sampled cell position and sensor observation on the manifold surface. $S_{C_1}$ and $S_{C_2}$ indicate the orientation estimation from the skyline and the ground plane respectively.

It is noteworthy that all the particle samples are generated on the discretized cells, spreading over the manifold space formed by the roll and pitch angles. Each particle is a 2D vector represented by a cell position on the spherical surface. There are three levels of cell resolution ($\Omega_{1,2,3}$ in Algorithm 1) from coarse to fine, where the longitudinal direction of the spherical surface represents the pitch, while the latitudinal direction is the roll. In line 2, particles are sampled from a Gaussian distribution $\mathcal{N}$, with a mean value at the IMU measurements plus a shift by a constant angular velocity propagating through a certain interval plus an offset $b$.

When both observations from the skyline and ground plane are available (line 4), more samples will be created around those cells close to sensor measurements (line 8), and their weights are initialized by multiplication of parent weight and local weight as a function of angular distance (line 9), whereas the other particles' weights will be down-weighed by an aforementioned normal distribution as a function of angular distance. Line 7 manifests the angular distance criteria for neighboring particle cells close to observation. The sample cells meeting the criteria are used as parents to create more children particles around $\mu_f$ at line 8 of Algorithm 1, and $\mu_f$ is derived from Equation 3.21, as a weighted sum of results from two Computer Vision pipelines, with each result used as a mean value of the normal distribution. A simple inverse of corresponding variance $\delta_{C_1}, \delta_{C_2}$ respectively is considered as weight for mean sum.

$$\mu_f = \frac{\mu_{C_1}}{\delta_{C_1}} + \frac{\mu_{C_2}}{\delta_{C_2}} \tag{3.21}$$

$$\delta_f = \left( \frac{1}{\delta_{C_1}} + \frac{1}{\delta_{C_2}} \right)^{-1} \tag{3.22}$$

IMU and CV observation variances are set as a constant according to practical tests. $\mu_f$ and $\delta_f$ are fused results from two CV pipelines in the same weighted sum form of Equations 3.21. The same strategy repeats when a single CV observation pipeline is present (line 15), but sampling rather on a middle-level resolution.

Each particle's life cycle can be represented in four phases, as shown in Figure 3.7. Arrows between them stand for transition conditions. In our setting, the Computer Vision based orientation observation has a smaller variance compared to IMU. As aforementioned, three levels of cell resolution

Figure 3.7: Lifetime phases of particle filter sampling on a spherical surface.

are employed. The top left part of Figure 3.7 represents the initial state, sampled from a normal distribution centered at the measurement of timestamp $t$. The bottom left is the prediction based on the propagation of the previous particle states by extrapolation in time. When the orientation from a single pipeline, either skyline or ground plane, is available, the new samples in red are generated on finer resolution neighboring the green dots corresponding to $\mu_f$ in Equation 3.21. Each dot is located at the center of a cell. If both skyline and ground plane pipelines are available, the highest resolution $\Omega_3$ is employed to generate more new samples. Each particle's timestamp of creation is kept as well. If the lapsed time exceeds a certain interval $\delta t$, the particles will be placed back at a coarse resolution, like the arrow direction from the bottom right to the top right. At a certain time point, the particles in the set have various precision. A lifetime check will be called periodically to eliminate the particles that have existed for a quite long time.

## 3.4    Experiments and Results

In practice, the video is scaled down to 640×480 resolution to achieve a raw frame rate of 20, while the overall frame rate scales down to 12-15 after the fusion on Jetson Nano. Intrinsics of the camera are acquired following the calibration guide of [283]. Extrinsic calibration between low-cost IMU (BNO055) and Raspi-camera is established using an open source tool "Kalibr" [200], [64]. Figure 3.8 shows the simulation setup on top of the



Figure 3.8: Simulation test setup on top of the building.

building in the landscape. There are three parts, a motor driver board, Jetson Nano, and sensors. All 3D-printed cases have enhanced connections, taking aerodynamics into account for flying efficiency. The camera on the pose estimation system is placed forward facing the landscape. Then the system cases along with the sensors are attached to a pole end (not in the view of Figure 3.8). The other end of the pole is controlled manually to simulate random rotation. Here, the ground truth roll and pitch angles are read from the servo motors, and the protractor in the figure is only adopted for verification of the test. In the demo test, we use the fused estimation from our particle filter to steer the motors. Closed-loop PID controller is leveraged

43

for actuation. All the following test sequences were recorded from a static position at the start. Furthermore, it is guaranteed the ground plane should be orthogonal to the gravitational vector at the start. This configuration remains unchanged for the real UAV test on the drones.

Open-source datasets for orientation estimation research mostly overlapped with SLAM research, and the SLAM datasets were often captured indoors or within urban regions, rarely including unpopulated areas, viewed from the top. We thus recorded sequences by using the aforementioned setup in a real mountain landscape on a tall building roof nearby, which should be quite similar to the camera view on the airplane. During recording. Each pose configuration of the system is kept on par with angle readings from motors containing hall sensors as ground truth. For comparisons, the SOTA visual-inertial frameworks "ORBSLAM3" [25], "R-VIO" [91], and "DM-VIO" [192] were selected at first, but we found these algorithms are dedicated to 6-DoF visual odometry and are all relying on feature points extracted from the images, which are not suitable for the challenging featureless scenes of our datasets, so we compare our fusion algorithm against the IMU filtered by quaternion-based Madgwick, CV only pipelines, like skyline or ground respectively.

Table 3.1: Average RMSE of roll and pitch angles (in radians). The framewise error bigger than radians 0.3 occurring over the half sequence length is considered a failure. The sequence test period is indicated in the brackets next to the test type.

| | Sequence | IMU (Madgwick Filter) | Skyline only | Ground plane only | Fusion |
|---|---|---|---|---|---|
| | test01 | 0.0121 | 0.0182 | 0.0344 | **0.0090** |
| Roll test (125s) | test02 | 0.0147 | 0.0236 | 0.0457 | **0.0126** |
| | test03 | 0.0162 | 0.0325 | 0.0593 | **0.0152** |
| | test01 | 0.0147 | 0.0196 | 0.0325 | **0.0118** |
| Pitch test (127s) | test02 | 0.0174 | 0.0214 | 0.0291 | **0.0139** |
| | test03 | 0.0208 | 0.0241 | —— | **0.0165** |
| | test01 | 0.386 | 0.0713 | 0.0674 | **0.0451** |
| Mixed test (960s) | test02 | 0.415 | 0.0651 | —— | **0.0584** |
| | test03 | 0.491 | 0.0742 | —— | **0.0617** |

Here, the RMSE results of my model are compared to other popular SLAM baseline models, which are based on visual-inertial odometry. Regarding the roll-only test, the proposed method (Ours) generally performs better or on par with the other methods. It shows lower error values compared to ORB3 and R-VIO and is comparable to or slightly better than DM-VIO. As for pitch only test, the proposed method again shows competitive performance. In some cases, it outperforms other methods, especially in the test02 sequence where R-VIO does not provide results (indicated by

Table 3.2: Comparison results of our method and baseline method of average RMSE of roll and pitch angles (in radians).

|  | Sequence | ORB3[25] | R-VIO[91] | DM-VIO[192] | Ours |
|---|---|---|---|---|---|
| | test01 | 0.01011 | 0.02071 | 0.00942 | **0.00862** |
| Roll test (125s) | test02 | 0.00994 | 0.03904 | **0.00300** | 0.00824 |
| | test03 | 0.01309 | 0.06566 | **0.00324** | 0.00945 |
| | test01 | **0.01207** | 0.02828 | 0.04666 | 0.01657 |
| Pitch test (127s) | test02 | 0.04254 | —— | —— | **0.03446** |
| | test03 | 0.00914 | —— | —— | **0.00875** |
| | test01 | **0.01845** | 0.03975 | —— | 0.01972 |
| Mixed test (9200s) | test02 | 0.01828 | 0.05635 | —— | **0.01804** |
| | test03 | 0.01617 | 0.03982 | —— | **0.01615** |

dashed lines). In the mixed test, the proposed method maintains its competitive margins over the baseline methods, showing lower or comparable error values compared to the other methods. It performs consistently across different sequences. It should be noted that dashed lines indicate that R-VIO failed to produce results in several test sequences, particularly in the pitch and mixed tests. This indicates robustness issues with the R-VIO method in such challenging scenarios. In general, the proposed method demonstrates consistent performance across all test sequences, often showing lower error rates compared to the other methods. This suggests a more robust and reliable performance. Although DM-VIO and ORB3 also show good performance, the proposed method often achieves the lowest error values, indicating that it might be the best-performing method overall in this comparison. This summary highlights the strengths of the proposed method in terms of performance consistency and robustness across various test conditions compared to existing methods.

We can conjecture from Table 3.1 that our fusion approach consistently outperforms the other baseline approaches without fusion by a considerable margin on all sequences. The sequences cover three movement patterns: sequences with pure roll, pure pitch movement, and random rotation on both axes, and each case is implemented at different angular speeds, ordered from three levels, 3, 9, and 15 degrees per second. The good performance with the lowest RMSE in most tests can be attributed to our good assumptions of the environment, skyline, and ground plane in the wild. The IMU results in the table are from the 6-DoF Madgwick quaternion filter [153], without the use of a magnetometer. This is because in our case, the mountain region with active volcanoes is affected by the disturbances from the earth's magnetism field change. The filtered IMU results are prone to drift, as presented in the mixed test of 16 minutes, and the errors are nearly one order of magnitude bigger than roll and pitch test sequences. Either the use of skyline or the

ground plane as a tracking cue can guarantee the error remaining at a lower level compared to the Madgwick filter over IMU measurements, but in some fast rotation cases, the ground planes are partially or not present in the image, which may fail. All of the results in the table justify the merits of using an adaptive particle filter over the manifold, improving the robustness, sensor redundancy, and accuracy.

Figure 3.9 further validates the consistency of our method. The test is repeated 10 times per sequence, and then an average error of all trials is taken over the mean error of the whole sequence. The slowest angular speed of the sequence (test01) for pure roll or pitch in Table 3.1 is employed. The variances of our fusion results are always the smallest compared to the results of the single sensor modality.



Figure 3.9: The Green arrow is a mean error, the orange line is a median error, and the box bounds represent the min/max errors. Roll and pitch results are in pink and blue respectively.

## 3.5 Conclusion

In this chapter, we proposed a stand-alone gimbal system for pose estimation from image frames, leveraging natural geometric primitives such as skylines and ground plane approximations. By using two lines and their normal vectors, we derived the current frame's rotation relative to a reference frame. A

particle filter with adaptive resolution-based sampling was introduced to fuse orientation estimates from both computer vision (CV) and IMU pipelines, adapting to different phases of particle lifetime. The system was implemented on a 3D-printed gimbal platform and tested in real-time on a Jetson Nano. Comparisons with IMU-only solutions demonstrated superior accuracy and robustness against drift and disturbances.

While our approach benefits from a simplified geometric assumption, it faces challenges in extreme weather conditions, abrupt illumination changes, and scenarios where the skyline is not visible. To improve robustness in diverse environments, a hybrid system incorporating feature points alongside the skyline could be explored. Currently, the skyline is approximated as a straight line for real-time processing, but complex terrains with curved mountains challenge this assumption. Image-level matching techniques, such as iterative closest point (ICP), could help refine alignment. Additionally, integrating a fisheye camera or multiple cameras from different viewpoints would enhance the system's robustness and adaptability across varied scenes.

# Chapter 4

# Point Cloud Registration Using 2D Surfel-Based Equivariance Constraint

## Abstract

Point cloud registration is a critical task in 3D reconstruction. Traditional methods extract and match 2D or 3D feature points, but these matches are often sparse and unreliable due to a low inlier-to-outlier ratio, making registration prone to failure. Recent point cloud-based registration methods, both learning-based and non-learning-based, focus primarily on spatial alignment while often ignoring point orientations and uncertainties. Additionally, these methods require extensive transformation enhancements in training data, making these methods sensitive to noise and large rotations. To address these limitations, we propose a surfel-based point cloud registration framework [1] that leverages $\mathbf{SE(3)}$-equivariant features for robust alignment. Our

**Author Contributions:**

- **Xueyang Kang**: Idea Design, Methodology, Software, Experiment Validation, Formal Analysis, Data Curation, Writing, Review, and Editing.

- Hang Zhao: Data Curation and Review.

- Zhaoliang Luan: Review.

- Patrick Vandewalle: Review, and Supervision.

- Kourosh Khoshelham: Review, and Supervision.

method initializes surfels from clustered superpixels aligned with depth maps using camera parameters or directly from LiDAR scans, learning the equivariant position and orientation representations through $\mathbf{SE(3)}$-equivariant convolutions. The model integrates an equivariant convolutional encoder, a cross-attention module for similarity estimation, a fully connected decoder, and a non-linear Huber loss for robust optimization. Experimental results on indoor and outdoor datasets demonstrate the effectiveness and robustness of our approach, outperforming state-of-the-art methods in real-world point cloud registration.

## 4.1    Introduction

Registration is a critical challenge in 3D reconstruction, shape alignment, VR/AR, and various applications, as it involves estimating relative transformations between source and target frames. Traditional 2D registration relies on finding correspondences between distinctive feature points extracted from image frames. However, repetitive or ambiguous textures, along with high outlier ratios in correspondence candidates, often make these methods unreliable. Furthermore, 2D view registration frequently results in under-constrained solutions for full $\mathbf{SE}(3)$ pose estimation due to these limitations.

Recent progress in monocular foundation models, such as DepthAnything [266], DepthAnythingV2 [268], and LOTUS [80], predict consistent depth and normal maps from color images. This enables unprojecting pixels into 3D point clouds in the camera frame, and this approach effectively lifts 2D image registration into 3D point cloud registration to avoid correspondence ambiguity of 2D views. However, these predicted depth maps inherently have uncertainty introduced by view projection geometry, making the point cloud noisy. Current point cloud registration methods, including ICP [149], DGR [41], and PointDSC [11], neglect such uncertainties and associated point cloud normals. Additionally, while deep learning models process point clouds, they are often translation-equivariant but not rotation-equivariant, so they require extensive data augmentation during training to handle registration under large rotation, making training very inefficient.

Traditional registration techniques typically use raw colored point clouds initialized from RGB-D images [67, 223, 49]. Rigid methods like ICP optimize point-wise match errors but are sensitive to noisy inputs and high outlier ratios, while non-rigid approaches address registration of deformable objects in dynamic environments [230, 170, 175] with temporal information. Learning-based methods, such as D3feat [12], SpinNet [8], and RoReg [243], aim to learn $\mathbf{SO(3)}$-equivariant features for rotation and position estimation.

However, they struggle in challenging scenarios involving noisy inputs, large rotations, symmetry ambiguity existing in the input, or near-planar camera motion, often failing to resolve orientation ambiguities, in particular for transformations that are orthogonal or antiparallel. We provide a demonstration Figure 4.1 to clarify the model equivariance, invariance, and its general relationship diagram,

A model encoder function $f$ is claimed to be equivariant if applying a transformation $g$ to the input $x$, followed by the function $f$, results in the same transformation $g'$ (an approximation for group transformation in feature space) being applied to the output. Mathematically, the equivariance can be formulated as below:

$$f(g(x)) = g'(f(x)). \tag{4.1}$$

In Eq. 4.1, the input transformation causes a corresponding transformation in the output (set of points), preserving relationships. In contrast, a function $f$ is invariant if applying a transformation $g$ to the input $x$ does not change the output. Mathematically, the invariance can be stated as below:

$$f(g(x)) = f(x). \tag{4.2}$$

Although the input undergoes the same rotation, the output remains the same, demonstrating invariance, as shown in the right side of Eq. 4.2.

Various equivariant model structures exist, including rotation-equivariant CNNs with steerable kernels [143], equivariant graph models with steerable features [30], and attention-based equivariant models [60] where queries, keys, and values are mapped from edge features. Each equi-model structure has its own merits, such equi-CNN is efficient, while equi-GNN and equi-Transformer have a large receptive field to learn global features.

To address these aforementioned registration challenges, we introduce a surfel-based pose regression model. Surfels, small oriented disks initialized from reprojected view point clouds or LiDAR data, serve as 2D Gaussian primitives encoding point position, normal, and uncertainty. Compared to point clouds, surfels offer superior robustness by leveraging data uncertainties, as demonstrated by recent advancements in Gaussian-based methods such as Gaussian Splatting [118] for 3D rendering tasks. Building on previous work [46, 18, 187], we propose a novel **SE(3)**-equivariant deep learning model for surfel-based registration.

Our method integrates a pipeline for initializing point clouds from 2D depth maps using camera parameters or directly from LiDAR scans [101, 90] and the depth can be predicted by a 2D monocular foundation model given a single image. The model employs a rotation-equivariant convolutional

(a) Equi-CNN Kernel

(b) Equi-Graph

$$\mathbf{v}_{ij} = \mathbf{W}_V \left( \mathbf{x}_j - \mathbf{x}_i \right) \mathbf{f}_j$$
$$\mathbf{k}_{ij} = \mathbf{W}_K \left( \mathbf{x}_j - \mathbf{x}_i \right) \mathbf{f}_j$$

$$\mathbf{q}_i = \mathbf{W}_Q \mathbf{f}_i$$

(c) Equi-Transformer

Figure 4.1: Different Equivariant model designs. (a) is equivariant $\mathbf{SO}(3)$ kernel [143] of CNN, and (b) is the equivariant graph model [30]. Lastly, (c) is the equivariant Transformer attention [59].

encoder based on E2PN [34] to extract $\mathbf{SE(3)}$-equivariant features, explicitly capturing both position and orientation. Pairwise equivariant features undergo cross-attention to compute similarity maps for establishing correspondences. These attention-enhanced features are concatenated and passed through fully connected layers to predict the relative transformation. Additionally, we introduce a differentiable $\mathbf{SE(3)}$ Huber loss function to supervise the soft correspondences, inspired by the node-wise supervision approach of PointDSC [11].

Experimental results on real-world indoor and outdoor datasets demonstrate the effectiveness and robustness of our model, particularly in large-scale scenarios. Compared to state-of-the-art methods, our approach addresses key limitations in traditional and learning-based registration methods, offering improved accuracy and scalability.

In summary, the main contribution of this work can be briefly outlined as follows:

- A surfel-based initialization pipeline from depth maps or LiDAR scans for pose regression tasks.

- A $\mathbf{SE(3)}$-equivariant deep learning model leveraging surfel features for robust and efficient registration.

- A differentiable $\mathbf{SE(3)}$ Huber loss function for supervision using soft correspondences.

## 4.2   Related Work

**Applications of 3D Surfels.** SurfelMeshing [213] has been applied to 3D mapping in indoor scenes, while other approaches have used surfels for large-scale outdoor mapping [278, 245, 18]. Surfels model the probability of a point by considering sensor accuracy and observation uncertainty, such as errors in pixel location due to camera projection. This data structure of surfels explicitly encodes $\mathbf{SE(3)}$ information, unlike point cloud which relies on neighboring points to derive rotation information implicitly.

**3D Registration.** Registration techniques are widely used in shape alignment ([41, 286]) and deformable target scanning [21]. Traditional non-learning methods like vanilla Iterative Closest Point (ICP), kiss-ICP [149], and point-to-plane ICP [184] struggled with nonlinear optimization errors. In contrast, state-of-the-art deep learning models like PointDSC [11] and Deep Global Registration (DGR) [41] explore correspondences in high-dimensional feature spaces. Max-Clique [282] and GeoTransformer [193] leverage the latest

graph or attention learning to create the registration backbone. Additionally, research focuses on enhancing feature descriptor learning to capture neighboring geometry information, exemplified by Fully Convolutional Geometric Features [42].

**Equivariant Feature Representation.** Equivariant models have emerged as robust solutions for 3D applications such as point cloud registration and pose estimation. These models maintain **SE(3)** rotational and translational equivariance, essential for consistent alignment performance. Examples include Spherical CNNs [45], group CNNs [56], and SE(3)-Transformers proposed [59]. Unlike the state-of-the-art models that are solely translation-equivariant [42, 41], incorporating rotation-equivariant and invariant feature characteristics can enhance learning efficiency and robustness. D3feat [12], SpinNet [8], and RoReg [243] leverage rotation-equivariant features for point cloud registration. Other approaches integrate equivariance learning into model structure design, such as vector neurons [47] by extending neurons from 1D scalars to 3D vectors, using steerable kernel filters [45, 253, 252, 229, 251, 102], projecting features onto spherical harmonics [235], employing rotation-guided attention mechanisms [100, 59, 60], and incorporating equivariant high-dimensional features into graph models for message propagation [108, 50, 207].

## 4.3 Method

For perspective image input, we initialize surfels from the unprojected depth map and its associated normal map, where the normals are derived from the depth gradients. This process also involves determining the surfel uncertainty 1D radius based on a camera perspective projection model, accounting for both inverse depth uncertainty and camera view angle to the image center. The whole initialization process serves as a pre-processing step.

For LiDAR point clouds, surfels can be created from the neighboring non-coplanar triplet points for normal vector estimation, and 1D uncertainty can be derived proportionally from the point density.

To ensure registration efficiency, surfels from source and target frames are downsampled before feeding into the neural network model. The model architecture, as illustrated in 4.4, consists of three components: an **SE(3)** equivariant convolution kernel-based encoder, a cross-attention, and a decoder for predicting relative transformation.

53

### 4.3.1  Surfel Initialization

The general goal of sampling is to find an optimal surfel representation of the geometry with minimal redundancy. Most sampling methods perform object discretization based on geometry parameters of the surface, such as curvature or silhouettes. This object-space discretization often results in either excessive or insufficient primitives for rendering.

While traditional surfel-based methods map the environment by assigning per-pixel surfels, our approach extracts surfels based on superpixels from RGB images using the SLIC algorithm [2]. This method generates 3D primitives for lower-resolution surface reconstruction compared to raw input, enabling faster processing with reduced memory consumption and noise. At the same time, it preserves critical information, such as texture and depth discontinuities. The pixels are clustered by intensity, the main parameter being the minimum number of pixels $k$ per superpixel to ensure approximately equal sizes (Figure 4.2). This superpixel-based approach minimizes memory overload for large-scale scenarios while reducing outliers and noise from low-quality depth maps.



Figure 4.2: (Left) Color image decomposed into superpixels of large size. (Right) Color image decomposed into superpixels of small size.

Each surfel is initialized from a single superpixel. The position $\mathbf{p_i}$, normal $\mathbf{n_i}$, and color $\mathbf{c_i}$ are computed as the average values of the pixels within the superpixel:

$$\mathbf{p_i} = \frac{1}{k} \sum_k \mathbf{p_k^{s_i}}, \quad \mathbf{n_i} = \frac{1}{k} \sum_k \mathbf{n_k^{s_i}}, \quad \mathbf{c_i} = \frac{1}{k} \sum_k \mathbf{c_k^{s_i}}, \tag{4.3}$$

where $\mathbf{p_k^{s_i}}$, $\mathbf{n_k^{s_i}}$, and $\mathbf{c_k^{s_i}}$ represent the position, normal, and color values of the $k^{th}$ pixel in superpixel $s_i$. The surfel radius $r_i$, which represents uncer-

tainty, is calculated as the dot product between the normalized pixel direction $\mathbf{v}_i$ and the camera principal axis facing forward:

$$\rho_i = |\tan\left(\mathbf{v_i} \cdot \mathbf{z}\right)| \tag{4.4}$$

This approach ensures smaller uncertainty at the center and larger uncertainty near the edges, where the viewing angle is greater. The resulting surfels have corresponding radii, as illustrated in Figure 4.3, where brighter points indicate larger radii and darker points indicate smaller radii.

To visualize the results, we separated the properties into three different point clouds due to the limitations of available open-source libraries. These visualizations are shown in Figure 4.3.

This way, the uncertainty is smaller in the center and larger on the sides where the angle given by the point of view is also large. The resulting surfels will have corresponding radii as shown in Figure 4.3, where brighter points correspond to large radii as opposed to dark-colored points having small radius values.

To visualize the 3D surfel results, due to the limitations in the open source libraries ready to use, we separated the different properties into three different point clouds visible in Figure 4.3.



Figure 4.3: (Left) Color points cast from the depth map aligned with the color pixels. (Middle) Point normals cast from the normals map aligned with the color pixels. (Right) Point uncertainties are cast from the pixel-wise uncertainties aligned with the color pixels.

In particular, each surfel in the source frame, indexed by $i$, consists of three main components: the 3D position $\mathbf{x}_i \in \mathbb{R}^3$, the normal vector $\mathbf{n}_i \in \mathbb{R}^3$, and a scalar radius $\epsilon_i \in \mathbb{R}$. The surfel $\mathbf{y}_j$ is in the target frame. Here, we mainly focus on clarifying the complex process of surfels created from a perspective view. The surfel center position is determined by the depth map and the camera intrinsic matrix through unprojection ($\pi$). To calculate the normal map from the depth gradient, a CUDA-based Sobel operator is

applied to the depth map in parallel. Subsequently, the resulting normal map is converted into 3D normals using the camera intrinsic matrix, providing information about the orientation of the disk associated with the surfel center. The surfel radius $\epsilon_i$ is then derived by the following expression:

$$\epsilon_i = C \frac{e^{-\hat{\rho}}}{1 + e^{-\tan(\theta)}}, \tag{4.5}$$

where $C$ is the normalization factor, and $\theta$ represents the view angle between the camera principal axis $\vec{o}$ and the ray $\vec{r}$ emitted from the camera center through the pixel location, as expressed in the following equation,

$$\theta = \arccos\left(\frac{\vec{r} \cdot \vec{o}}{\|\vec{r}\|\|\vec{o}\|}\right). \tag{4.6}$$

The value of $\rho$ represents the inverse of the depth, which is then truncated to $\hat{\rho}$ within the inverse depth range $(\rho_{min}, \rho_{max})$ of the sensor,

$$\hat{\rho} = \min\left(\max\left(\rho, \rho_{min}\right), \rho_{max}\right). \tag{4.7}$$

## 4.3.2 Network Structure

As exhibited in Figure 4.4, given initialized surfels of source frame $\mathbf{s}_i \in \mathbf{s}_1, ..., \mathbf{s}_N$ and target frame $\mathbf{s}_j \in \mathbf{s}_1, ..., \mathbf{s}_N$, all surfels are encoded by the same encoder $f_\theta(\cdot)$. Notably, the position and normal vectors $\mathbf{n}_{(\cdot)}, \mathbf{p}_{(\cdot)}$ of each surfel are weighted by a factor of $(1 - \|\epsilon(\cdot)\|)$ to reduce the influence of highly uncertain surfels. The encoder architecture is augmented on E2PN [34], with doubled feature dimensions compared to the original point cloud input. Equivariance is maintained through a symmetric conv-kernel $\kappa$ arranged in an icosahedral shape solid.

In the following convention, the symbol $'$ next to a symbol indicates discretization operation. E2PN encoder features are aligned in the spherical space $\mathbf{S}^{2'} \times \mathcal{R}^3$, where coordinates of each feature vertex in $\mathcal{R}^3$, associated with the 128-dimension feature descriptor, and $\mathbf{S}^{2'}$ signifies the discretized sphere surface. This discretized feature representation is determined by $\mathbf{SO}(3)'/\mathbf{SO}(2)'$, where $\mathbf{SO}(2)'$ is a subgroup of $\mathbf{SO}(3)$. The quotient space is defined as a group of rotations $\mathbf{R}_{i/j}$ with the same endpoint, such as the sphere's north pole after rotation. This discretization of $\mathbf{SO}(3)$ facilitates more efficient and accelerated learning. The $\mathbf{SE}(3)$ feature is constructed by extending the $\mathbf{SO}(3)$ rotation feature, incorporating translation through the concatenation of point coordinates.

Figure 4.4: The network structure features a shared encoder for surfels (6 dimensions plus uncertainty radius) from both source and target frames in **SE(3)** space. This encoder maps 1024 surfels with 6 dimensions (position and normal), weighted by confidence value $(1 - \epsilon(\cdot))$, into 12 feature descriptors in 128-dim with 60 group rotation orders. Then each descriptor undergoes linear embedding to produce triplet token embeddings $\mathbf{Q}$, $\mathbf{K}$, and $\mathbf{V}$. Cross-attention $g_\theta(\cdot)$ is applied to feature descriptors from source and target frames in 12-channel dimension. The resulting tokens are in the shape of $(12 \times 128 \times 60)$, where each token is in 128-dim $\times$ 60 order groups, and the attention map is $12 \times 12$, where each element token in the attention map is formulated via descriptor dot product. Then the features are flattened and processed through Fully-Connected (FC) layers, mapping features to relative position $\mathbf{t}$ and relative quaternion $\mathbf{q}$ rotation. A close-up of the E2PN module is provided below.

Figure 4.5: Recovering discreteized $\mathbf{SO(3)}'$ from the quotient feature $S^{2'}$ by permutation order.

The surfel undergoes convolution with two distinct symmetric kernels, $\kappa_1$ and $\kappa_2$, as shown in Figure 4.5, relating to point position and normal respectively. These are concatenated post-convolution. The icosahedron comprises 60 rotations, each denoted by various permutation orders $\mathbf{R}_{(\cdot)}$. The E2PN symmetric kernel selects one rotation from the 60 options to generate the output equivariant feature. This is achieved by choosing the maximum sum of each rotation feature along the 12-channel dimension. The Equivariant features after E2PN encoder of source and target frame can be stated as the $\mathbf{D}_i, i \in (1, ..., 12)$ and $\mathbf{D}_j, j \in (1, ..., 12)$, where each has 128 dimensions. Next, we use the linear layer to project each descriptor into a triplet composed of $\mathbf{Q}_{(\cdot)}, \mathbf{K}_{(\cdot)}, \mathbf{V}_{(\cdot)}$. We use the same index convention as the input in the following. $i$ states the feature of the source frame, and $j$ indicates the index of the target frame. The cross-attention $g_\theta(\cdot)$ is then applied to calculate the attention-weighted equivariant features from the pairwise frames. Different from the normal attention-based multiplication, The feature embedding vectors (denoted as $\mathbf{Q}_i, \mathbf{K}_j$ respectively, $1 \leq i, j \leq N$), and each output token

feature $\hat{\mathbf{V}}_i$ is calculated as formulations below,

$$\alpha_{ij} = \frac{\exp(\mathbf{Q}_i^T \mathbf{K}_j)}{\sum_{j=1}^{N} \exp(\mathbf{Q}_i^T \mathbf{K}_j)}, \tag{4.8}$$

$$\hat{\mathbf{V}}_i = \sum_{j=1}^{N} (\alpha_{ij} \mathbf{V}_j), \tag{4.9}$$

The transformation applied to the input point cloud is denoted by group rotation $g'(\cdot)$. The key idea of equivariant feature learning is that the output features of the encoder are transformed accordingly to preserve $\mathbf{SE(3)}$ equivariance $f_\theta(g'(x)) = g'(f_\theta(x))$. This ensures that the features are equivariant to the input transformations. The equ-features after the E2PN encoder of the source and target frames are denoted as $\mathbf{D}_i, i \in (1, ..., 12)$ and $\mathbf{D}_j, j \in (1, ..., 12)$, where each has 128 dimensions. Next, we use the linear layer to project each descriptor into a triplet composed of $\mathbf{Q}_{(\cdot)}, \mathbf{K}_{(\cdot)}, \mathbf{V}_{(\cdot)}$ tokens. We use the same index convention throughout the paper, where $i$ denotes the index of the source frame, and $j$ refers to the target frame. The triplet tokens, derived from feature descriptors at the 12 corners of icosahedral planoids (see Figure 4.5), are aggregated from neighboring surfel coordinates. This effectively fuses the input surfels into 12 distinct regions, which are then used in the subsequent cross-attention module. The cross-attention $g_\theta(\cdot)$ is then applied to calculate the attention-weighted equivariant features (attention map in $12 \times 12$) from the pairwise frames, finally to be decoded by the fully connected layers into the transformation estimation.

### 4.3.3 Loss Function

Inspired by the node-wise supervision in PointDSC [11], we adapt the original binary cross-entropy loss to a non-linear Huber loss. This maps the transformed point error into $\mathcal{L}_2$ norm when it is small and into $\mathcal{L}_1$ normal when the error is large. The point position from the source frame is transformed by the predicted rotation $\mathbf{R}$ and translation $\mathbf{t}$ into $\mathbf{y}_{j*} = \mathbf{R}\hat{\mathbf{x}}_{i*} + \mathbf{t}$. The rotation matrix $\mathbf{R}$ is derived from the predicted quaternion $\mathbf{q}$. The Huber loss is defined as below,

$$\mathcal{L}_{\text{Huber}} = \begin{cases} \frac{1}{2}(\hat{\mathbf{x}}_{i*} - \mathbf{y}_{j*})^2, & \text{if } |\hat{\mathbf{x}}_{i*} - \mathbf{y}_{j*}| \leq \delta \\ \delta(|\hat{\mathbf{x}}_{i*} - \mathbf{y}_{j*}| - \frac{1}{2}\delta), & \text{otherwise} \end{cases} \tag{4.10}$$

The threshold is set to 0.6m. The correspondence index pair $(i, j)$ is established using the nearest neighboring point search.

## 4.4 Experiments and Results

To evaluate the model performance, we utilized two indoor scan datasets, the outdoor dataset KITTI [66] and 3DMatch [280]. These datasets include RGB-D sequence frames and extrinsic poses. Our model was trained on each dataset separately for comparison fairness against other baseline models. For ARKitScenes, we selected 10 different scenes. We further selected 100 pairs of frames out from each scene sequence by choosing the depth frames close to each other temporally, to guarantee enough overlap between the source and target scans created from the depth frame images. We employ the 3D voxel-based downsampling to generate 1024 points un-projected from the depth map of each frame for surfel initialization.

**Evaluation metrics.** We use the Rotation Error (RE) and Translation Error (TE) to evaluate the accuracy of rotation and translation separately. Furthermore, we incorporate the Registration Recall (RR) and *F1 score* as registration success evaluation metrics. Our model was trained on each dataset separately for fair comparison against other baseline models. We employed 3D voxel-based downsampling to generate 1024 points unprojected from the depth map of each frame for the surfel initialization.

$$\delta = \sqrt{\frac{1}{\mathcal{N}(\Omega)} \sum_{(\mathbf{x}_i, \mathbf{y}_j) \in \Omega} \mathbb{1}[\|\mathbf{R}\mathbf{x}_i + \mathbf{t} - \mathbf{y}_j\|^2 < \tau]}, \tag{4.11}$$

where $\mathcal{N}(\Omega)$ represents the total number of ground truth correspondences in the set $\Omega$. The symbol $[\cdot]$ is an indicator of whether the condition is satisfied. Based on the calculation of recall rate, the *F1 score* is defined as $2 \cdot \frac{Precision \times Recall}{Precision + Recall}$.

**Baseline models** We compare our model with popular deep learning-based models, including PointDSC [11], Deep Global Registration [41] (DGR), Maximal Clique [282] (MAC). Additionally, we choose equivariant methods, like D3Feat [12], SpinNet [8], RoReg [243] and MAC+GeoTransformer [282, 193] (MAC+GeoTrans) for the two datasets, For all the point feature descriptor dependent approaches, like DGR, FCGF [42] descriptor is applied for 3DMatch [279], while FPHF [204] descriptor is used instead for KITTI [66]. We provide quantitative evaluation results in Table 4.1, and qualitative comparison results including three top performance baseline models in Figure 4.6. All these results exhibit the superior performance of our model over baselines.

## 4.4.1 Baseline Comparisons

Our proposed model demonstrates superior performance compared to the baseline models across all test scenes, as shown in Table 4.1, with consistent performance superiority over other learning models on both datasets in terms of all metrics, in particular, the proposed model has a smaller rotation error (around 8-11% reduction) compared to the second best model of each dataset respectively. This can be attributed to the explicit orientation signal as the input for learning, along with the good equivariant feature learning through the symmetric kernel.

We employed 3D voxel-based downsampling to generate 1024 points unprojected from the depth map of each frame for surfel initialization.

Table 4.1: Evaluation results of registration approaches on 3DMatch [279] (left) and KITTI [66] (right).

| Method | 3DMatch [279] | | | | KITTI [66] | | | |
|---|---|---|---|---|---|---|---|---|
| | RE(°) ↓ | TE(cm) ↓ | RR(%) ↑ | F1(%) ↑ | RE(°) ↓ | TE(cm) ↓ | RR(%) ↑ | F1(%) ↑ |
| DGR [41] ↑ | 2.40 | 7.48 | 91.30 | 89.76 | 1.45 | 14.60 | 76.62 | 73.84 |
| D3Feat [12] | 2.57 | 8.16 | 89.70 | 87.40 | 2.07 | 18.92 | 70.06 | 65.31 |
| RoReg [243] ↓ | 1.84 | 6.28 | 93.70 | 91.60 | | | | |
| SpinNet [8] | 1.93 | 6.24 | 93.74 | 92.07 | 1.08 | 10.75 | 82.83 | 80.91 |
| PointDSC [11] | 2.06 | 6.55 | 93.28 | 89.35 | 1.63 | 12.31 | 74.41 | 70.08 |
| MAC [282] | 1.89 | 6.03 | 93.72 | 91.46 | 1.42 | 8.46 | 91.37 | 89.25 |
| MAC+GeoTF [282, 193] | 1.74 | 6.01 | 95.02 | 91.80 | 1.37 | 8.01 | 90.59 | 88.45 |
| Ours | **1.34** | **5.72** | **95.08** | **93.32** | **1.57** | **6.09** | **92.05** | **90.61** |

In addition to the quantitative comparison results, we further provide the visual comparison results in Figure 4.6 of the top three performance baseline models of each dataset. For better illustration purposes, we use the dense point cloud scan of source and target frames, instead of using the sparse input for training/testing to show each model's performance intuitively. While PointDSC [11], DGR [41] and RoReg [243] exhibit good registration accuracy in some test scenes, yet they suffer from orientation ambiguity, *e.g.*, in the second row, where only flat planes are dominant in the scans. PointDSC performs the worst among all the models in terms of both position and orientation accuracy. In challenging scenarios, such as the cluttered environment in the first row, DGR may even fail, resulting in an obvious misalignment, by seeing penetration into walls.

To verify the proposed model performance under varying numbers of sparsely sampled input points, we also implement the sparse tests as table below, by comparing with FCGF registration, D3Feat with or without PointDSC combination, and SpinNet. Our model has a consistent plausible performance on 3DMatch over the comparison models. In addition, more use

MAX Clique [282]   GeoTF [193]   Equi-(GSPR[108])   Ours   Ground Truth

PointDSC [11]   MAX Clique [282]   GeoTF [193]   Ours   Ground Truth

Figure 4.6: Comparison results on KITTI [66]. For each dataset, the top three models with good performance are presented.

points can boost the proposed model registration performance. In our setting. 1024 is chosen as the input point number to initialize the surfels in all the tests to get a good trade-off between accuracy and real-time performance. The surfel representation shows strong robustness even on extremely sparse points (256), with only a tiny performance drop of 2.5% compared to the 4096 number of use points, while all the other baseline models demonstrate a remarkable performance difference between 4096 and 256 points.

Table 4.2: RR results on 3DMatch with a different number of sampled points.

| #Sampled Points | 4096 | 2048 | 1024 | 512 | 256 | Average |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| FCGF [42] | 91.7 | 90.3 | 89.5 | 85.7 | 80.5 | 87.5 |
| D3Feat [12] | 91.9 | 90.4 | 89.8 | 86.0 | 82.5 | 88.1 |
| PointDSC [11] | 92.1 | 92.5 | 90.8 | 87.4 | 83.6 | 89.3 |
| SpinNet [8] | 93.8 | 93.6 | 93.7 | 89.5 | 85.7 | 91.3 |
| Ours | **95.4** | **94.8** | **94.1** | **92.4** | **87.8** | **92.9** |

## 4.4.2 Ablation Study

1) We perform eight different types of ablation tests to verify the contribution of each design to our model performance, as shown in Table 4.3. Point cloud position as input fed into various SOTA point cloud encoders, or vanilla E2PN encoders, cannot achieve performance on par with our surfel-based equi-model design. The uncertainties and Huber loss are all beneficial to improving the model performance. 2) We provide the robustness analysis of input scans perturbed by various levels of rotation and translation in Table 4.4. 3) Finally, we provide model complexity comparisons for the top five baseline models and our model in Table 4.5, as shown below, to showcase the low latency and small model size complexity of our model compared to other baselines.

Table 4.3: Ablation study on 3DMatch dataset.

| Method | RE(°) ↓ | TE(cm) ↓ | RR(%) ↑ |
|---|---|---|---|
| 1. Surfel w/o uncertainty weight | 1.64 | 6.75 | 91.73 |
| 2. Point cloud + vanilla E2PN | 2.24 | 7.16 | 89.85 |
| 3. Point Cloud + 3D CNN | 12.08 | 14.75 | 51.47 |
| 4. Point Cloud + PointNet++ | 8.73 | 10.82 | 68.91 |
| 5. W/o attention module | 4.62 | 18.24 | 56.31 |
| 6. $\mathcal{L}_1$ loss only | 1.96 | 6.75 | 87.08 |
| 7. $\mathcal{L}_2$ loss only | 2.49 | 7.93 | 84.80 |
| 8. Full model + Huber Loss | **1.34** | **5.72** | **95.08** |

Our model estimates the relative transformation between the source and target frame scans. We report the average metric errors across various translation and rotation intervals, applied to the same input scan for generating the source and target input pair.

Table 4.4: Robustness analysis of average errors on 3DMatch.

| | $5°, 10cm$ | $25°, 50cm$ | $50°, 100cm$ | $75°, 150cm$ | $100°, 200cm$ |
|---|---|---|---|---|---|
| RE(°) ↓ | 0.18 | 0.71 | 2.37 | 3.12 | 4.60 |
| TE(cm) ↓ | 0.51 | 1.75 | 2.86 | 6.73 | 12.59 |
| RR(%) ↑ | 94.02 | 93.26 | 90.15 | 88.50 | 87.41 |

The input consists of 100 scan pairs from the 3DMatch subset, with the source scan perturbed by the translations and rotations shown in the first row. As shown in the table, the registration error remains within 4-5% of the input translation, and the rotation error within 3-4%, demonstrating that the error does not scale with increased perturbations.

63

Table 4.5: Ablation study of model complexity on 3DMatch.

|  | DGR | PointDSC | Spinnet | RoReg | MAC+Trans | Ours |
|---|---|---|---|---|---|---|
| Latency(s) ↓ | 1.26 | **0.08** | 2.84 | 22226 | 0.31 | 0.09 |
| Params(Mb) ↓ | 1.94 | 1.07 | 3.16 | 83.25 | 28.64 | **0.98** |

We also provide the training loss curves within 300 epochs under the different threshold $\delta$ in Eq. 4.10. The threshold $\delta$ 0.5 can achieve a good trade-off between training convergence speed and accuracy (curve in blue). The Huber loss curve of the threshold bigger than 0.5 (curve in yellow and magenta), even suffers from the overfitting problem after 150 epochs. The training converges relatively slowly under the threshold $\delta$ 0.02 (in red) and 0.2 (in green) compared to 0.5, although the final loss reaches the accuracy of the same level as threshold 0.5. The small threshold $\delta$ makes the $\mathcal{L}_2$ loss range very limited, and the $\mathcal{L}_1$ loss is dominant in the learning process then.



Figure 4.7: Huber loss learning curve under different thresholds.

## 4.5 Conclusion

We propose a surfel-based **SE(3)** equivariant network model, incorporating surfel initialization from raw RGB-D depth maps or LiDAR point clouds. Our framework consists of a shared E2PN encoder, a cross-attention module, and an MLP-based decoder. Extensive experiments on two datasets demonstrate the model's robustness and state-of-the-art accuracy. Furthermore,

64

the model's modular surfel representation enables generalization across diverse 3D scenes. Future work could explore leveraging surfel primitives for 3D mapping, reconstruction, and mesh conversion to enhance downstream tasks, while also improving model robustness across varying levels of point cloud sparsity and even for noisier data, including dynamic objects and uncertain information. Lastly, a more efficient and robust correspondence structure can be investigated. Lastly, more efficient equivariant deep learning methods can be explored. Equivariance can be achieved through structural design, such as extending neurons from 1D scalars to 3D vectors, or by integrating modules that enforce rotation-equivariant properties. However, these approaches often increase the complexity of the model. Therefore, the trade-off between implementation complexity and efficient rotation generalization should be further investigated. In addition, we acknowledge that certain objects exhibit strong view-dependent appearance characteristics, which challenge conventional point cloud registration methods that typically assume view-independence. As the main focus of this chapter and paper is on geometric modeling rather than appearance, we have chosen to leave this issue as a noted limitation and highlight it as a promising direction for future work.

# Chapter 5

# Depth Prediction from Focal Stack Using Focal Geometry Constraint

## Abstract

Depth estimation from images is crucial for generating dense point clouds in the camera frame, enabling advanced 3D reconstruction tasks. While depth sensors and LiDAR provide accurate measurements, they remain expensive and limited in range. Alternatively, triangulated point clouds from sparse 2D feature points require extensive image coverage, increasing data collection challenges. Monocular depth priors in deep learning models offer a cost-effective and scalable alternative for dense point cloud generation. Depth from a focal stack is a specific challenge within depth estimation, where we use the focus differences across a sweep of images to figure out the depth of a scene for our proposed FocDepthFormer[1]. Unlike broader depth estimation techniques that might use stereo vision, structure-from-motion, or deep

---

[1]The majority of this chapter was published as a peer-reviewed conference paper [106] (Kang, Xueyang, et al. "FocDepthFormer: Transformer with Latent LSTM for Depth Estimation from the Focal Stack." Australasian Joint Conference on Artificial Intelligence. Singapore: Springer Nature Singapore, 2024).

**Author Contributions:**

- **Xueyang Kang**: Idea Design, Methodology, Software, Experiment Validation, Formal Analysis, Data Curation, Writing, Review, and Editing.
- Fengze Han: Data Curation and Review.
- Abdur R. Fayjie: Data Curation and Review.
- Patrick Vandewalle: Review, and Supervision.
- Kourosh Khoshelham: Review, and Supervision.

learning, focal stack methods tap into defocus clues to recover depth. This makes them especially handy in situations where other depth hints are hard to detect. Most existing methods for depth estimation from a focal stack of images employ convolutional neural networks (CNNs) using 2D or 3D convolutions over a fixed set of images. However, their effectiveness is constrained by the local properties of CNN kernels, which restricts them from processing only focal stacks of a fixed number of images during both training and inference. This limitation hampers their ability to generalize to stacks of arbitrary lengths. To overcome these limitations, we present a novel Transformer-based network, FocDepthFormer, which integrates a Transformer with an LSTM module and a CNN decoder. The Transformer's self-attention mechanism allows for the learning of more informative spatial features by implicitly performing non-local cross-referencing. The LSTM module is designed to integrate representations across image stacks of varying lengths. Additionally, we employ multi-scale convolutional kernels in an early-stage encoder to capture low-level features at different degrees of focus/defocus. By incorporating the LSTM, FocDepthFormer can be pre-trained on large-scale monocular RGB depth estimation datasets, improving visual pattern learning and reducing reliance on difficult-to-obtain focal stack data. Extensive experiments on diverse focal stack benchmark datasets demonstrate that our model outperforms state-of-the-art approaches across multiple evaluation metrics.

## 5.1   Introduction

In scenarios where random camera motion is constrained, depth can be inferred by capturing a focal stack, multiple images taken at varying aperture sizes and focal distances, where each image encodes depth-dependent defocus information. However, existing depth-from-focus methods typically rely on convolutional neural networks (CNNs) with fixed 2D or 3D kernels, which struggle to generalize across focal stacks of varying lengths and fail to capture long-range dependencies effectively.

To overcome these limitations, we propose FocDepthFormer, a novel transformer-based network designed to estimate the depth of focal stacks. Our model integrates a Transformer module to capture global spatial relationships, an LSTM to aggregate information across focal stacks of arbitrary lengths, and a CNN decoder to refine depth predictions. A multi-scale convolutional encoder extracting fine-grained focus/defocus features is also used

- Dong Gong: Review and Supervision.

before the Transformer encoder. By leveraging pre-training before monocular RGB-depth datasets, FocDepthFormer can achieve good generalization ability without using a large volume of focal stack datasets for training.

Extensive experiments on diverse datasets demonstrate that FocDepthFormer outperforms state-of-the-art methods across multiple focal stack benchmarks. Furthermore, its strong generalization ability enables accurate depth prediction on unseen focal stack datasets, and this shows the great potential of using such a technique for direct inferences on open-world dataset.

Traditional 3D reconstruction pipelines, such as COLMAP [212], rely on triangulating matched 2D feature points extracted from images using descriptors like SIFT [138]. These methods require dense image coverage to ensure sufficient correspondences for reconstructing 3D points. For large-scale scenes, such as outdoor buildings, generating dense point clouds often requires collecting thousands of images. As a result, the 3D points reconstructed from sparse feature points can vary significantly in density across spatial regions, depending on the distribution of input images.



Figure 5.1: Point cloud reconstructed from the extracted feature points of multiview image (Kings College Scene from Cambridge dataset [116]) by COLMAP [212, 210] tool.

With the advent of deep neural networks (DNNs) and large-scale image datasets, depth can now be predicted directly from images, enabling the generation of dense point clouds from color pixels. When the point clouds in the local frame are aligned consistently into a global and dense map, these dense point clouds can reduce the reliance on extensive image collections. Monocular depth estimation methods [273, 75] have demonstrated success on benchmarks [66], but their reliance on prior knowledge of scene textures, perspectives, and contexts often limits generalization to unseen data. Multi-

view consistent cues [228, 10] have been employed to improve generalization, but challenges persist in handling diverse scenes.

Depth estimation from RGB images remains a core paradigm for its convenience and flexibility. Researchers have explored various cues to infer depth, including image context [273, 75], geometry cues [228, 281, 68, 258], and focus/defocus cues [264, 72, 203]. Among these, focus/defocus cues, which leverage sharpness variations across focal distances, offer a promising yet underexplored approach to depth estimation [13, 28, 233, 161]. Despite being a fundamental depth estimation problem, focal stack distinctively varies from monocular depth estimation [52, 137, 70, 69, 54, 88, 162, 198, 164, 196, 86], stereo depth or disparity estimation [68, 258], and multi-frame depth estimation [211, 72]. These models are trained on RGB-D image pairs or disparity maps. However, the focal stack task can not directly utilize these models as focal stack images contain out-of-focus regions and no disparity cues from motion are available. Depth estimation from focus and defocus involves predicting the depth map from a captured *focal stack* of the scene, which comprises images captured by a camera focused on different focal planes [260], also referred to as the depth of field control problem [186], where focal stack images are obtained using a light field camera [141]. Conventional methods [232, 231, 171] address this task using handcrafted features based on sharpness. However, these methods often fail in textureless scenes. To enhance feature extraction for this task, Convolutional Neural Networks (CNNs) have been employed to learn depth map prediction from focal stack [79, 246, 264, 233, 28, 76, 7, 81, 63]. Specifically, DDFFNet [79], AiFDepthNet [246], and DFVNet [264] leverage in-focus cues for depth estimation. DefocusNet [161] aims to learn permutation invariant defocus cues, also known as Circle-of-Confusion (CoC). These methods utilize 2D or 3D convolutions to represent visual and focal features across spatial domains and stack channels. They either fuse stacked network 2D depth outputs [161] or predict a single 2D depth map from the 3D feature volume [264]. Despite the potential to enlarge receptive fields with increased network depth, CNN-based models are confined to capturing features in local areas. Additionally, existing methods are limited to focal stacks with a constant number of images during both training and testing, making it challenging to generalize to stacks with an arbitrary number of images.

In this chapter, we propose a novel LSTM + Transformer-based network for depth estimation from focal stacks, which is referred to as FocDepthFormer. The core component of FocDepthFormer is a module of *Transformer with latent LSTM*, which consists of a Transformer encoder [48], an LSTM-based recurrent module [84] applied to the latent tokens, and a CNN decoder. The Transformer and LSTM are used to separately model the spa-

tial and stack information. Different from the CNNs [264] restricted to local representation, the Transformer encoder captures the visual features with a larger receptive field. The self-attention mechanism in Vision Transformer (ViT) [48] facilitates the *cross reference* among non-local patterns, allowing the Transformer encoder to capture more informative features to represent the sharpness and blur characteristics. Considering that the focal stacks may have arbitrary and unknown numbers of images, we utilize the LSTM in latent feature space to fuse focusing information across the entire stack for depth prediction. This differs from existing focal stack depth estimation methods [264, 246, 161], as well as monocular depth estimation methods based on CNNs or Transformers [3, 197], which typically handle inputs with a constant number of images. Specifically, we compactly fuse activated token features via the recurrent LSTM module after the Transformer encoder. This design allows the proposed model to handle focal stacks of arbitrary lengths with a predefined order during both training and testing, offering greater flexibility in practice.

Before inputting data into the Transformer, we employ an early-stage convolutional encoder with multi-scale kernels [257] to directly capture low-level focus/defocus features at different scales. In light of the limited availability of focal stack data, our model exhibits the capability to enhance its scene feature representation through pre-training on monocular depth estimation datasets. This is facilitated by the use of the recurrent LSTM module in our model, allowing it to take varying numbers of input images. The main contributions of this work can be summarized as:

- We introduce a novel Transformer-based network model designed for depth estimation from focal stack images. The model utilizes a vision Transformer encoder with self-attention to capture non-local spatial visual features, effectively representing sharpness and blur patterns. To handle an arbitrary number of input images, we incorporate an LSTM-based recurrent module. This structural flexibility enables us to pre-train the model with a monocular depth estimation dataset, mitigating the demands for focal stack data, which is both limited and expensive to collect.

- To fuse the stack features, the LSTM is applied, and before the fusion, we employ a grouping operation to manage recurrent complexity over tokens without scaling complexity, as the token count increases due to the larger stack size. This is achieved by applying the LSTM solely to a subset of activated embedding tokens while preserving information on other non-activated tokens through averaging.

70

- Specifically, we also propose the use of multi-scale kernels in an early-stage convolutional encoder to enhance the capture of low-level focus/defocus cues at various scales.

## 5.2 Related Work

**Depth from Focus/Defocus.** Depth estimation from focal stacks relies on discerning relative sharpness within the stack of images for predicting depth. Traditional machine learning methods [260, 232, 231] treat this problem as an image filtering and stitching process. Johannsen *et al.* [105] provide a comprehensive overview of methods addressing the challenges posed by light field cameras, laying a foundation for research in this direction. More recently, CNN-based approaches have emerged in the context of focal stacks. DDFFNet [79] introduces the first end-to-end learning model trained on the DDFF 12-Scene dataset. DFVNet [264] utilizes the first-order derivative of volume features within the stack. AiFNet [246] aims to bridge the gap between supervised and unsupervised methods, accommodating ground truth depth or its absence. Barratt *et al.* [13] formulate the problem as an inverse optimization task, utilizing gradient descent search to simultaneously recover an all-in-focus image and depth map. DefocusNet [161] exploits the Circle-of-Confusion, a defocus cue determined by focal plane depth, for generating intermediate defocus maps in the final depth estimation. Anwar *et al.* [7] leverage defocus cues to recover all-in-focus images by eliminating blur in a single image. Recently, the DEReD model [224] learns to estimate both depth and all-in-focus (AIF) images from focal stack images in a self-supervised way, by taking the optical model into the loop to reconstruct the defocus effects. Gur and Wolf [76] present depth estimation from a single image by leveraging defocus cues to infer disparity from varying viewpoints.

**Attention-Based Models.** The success of attention-based models [239] in sequential tasks has led to the rise of Vision Transformer for computer vision tasks. The Vision Transformer represents input images as a series of patches ($16 \times 16$). While this model performs well in image recognition compared to CNN-based models, a recent study [257] demonstrates that injecting a small convolutional inductive bias in early kernels significantly enhances the performance and stability of the Transformer encoder. In the context of depth estimation, Ranftl *et al.* [197] utilize a Transformer-based model as the backbone to generate tokens from images, and these tokens are assembled into an image-like representation at multiple scales. DepthFormer [3] merges tokens at different layer levels to improve depth estimation performance. The latest advancement in this domain, Swin Transformer [146], achieves a larger recep-

tive field by shifting the attention window, revealing the promising potential of the Transformer model.

**Recurrent Networks.** Recurrent networks, specifically LSTM [84], have found success in modeling temporal distributions for video tasks such as tracking [178] and segmentation [261]. The use of LSTM introduces minimal computational overhead, as demonstrated in SliceNet [189], where multi-scale features are fused for depth estimation from panoramic images. Recent works [99, 97] combine LSTM with Transformer for language understanding via long-range temporal attention.

## 5.2.1 Preliminaries

In the paper, we present the features of our network model, a primary factor to explain our favorable results is the good identification ability of pixel sharpness, as mentioned in our paper. Usually, the sharpness of pixels can be evaluated according to the Circle-of-Confusion, $\mathbf{C}$ (in Figure 5.2), defined as follows,

$$C = \frac{f^2}{N(z - d_f)}\left|1 - \frac{d_f}{z}\right|, \tag{5.1}$$

where, $N$ denotes $f$-number, as a ratio of focal length to the valid aperture diameter. $C$ is the Circle of Confusion diameter (CoC). $d_f$ is the focus distance of the lens. $z$ represents the distance from the lens to the target object. In general, the range of $z$ is $[0, \infty]$. However, in reality, the range is always constrained by lower and upper bounds.

In Eq. 5.9, $C$ is zero ($C^*$) when the image pixel is in focus, and it is a signed value, where $C > 0$ indicates the camera focused in front of the sensor plane, while $C < 0$ is the reverse case, with a camera focused behind the sensor plane. It can be inferred from the denominator, that CoC $C$ has two divergence points, at 0 and $d_f$ respectively. Our model learns the mapping relationship from CoC to depth autonomously.

Given a focal stack, $\mathbf{S}$, comprising $N$ images ordered from near to far by focus distance, denoted as $\mathbf{S} = (\mathbf{x}_i)_{i=1}^{N}$, where each image $\mathbf{x}_i \in \mathbb{R}^{H \times W \times 3}$, our objective is to generate a single depth map $\mathbf{D} \in \mathbb{R}^{H \times W \times 1}$ for a stack of images. In contrast to the vanilla Transformer [48], we initially encode each image $\mathbf{x}$ using an *early-stage multi-scale kernel-based convolution* $\mathcal{F}(\cdot)$. This convolution ensures a multi-scale feature representation $\mathbf{x}'$ for the focal stack images. Subsequently, the *transformer encoder* $g(\cdot)$ processes the feature maps, transforming them into a series of ordered tokens that share information through self-attention. The self-attention weights between the in-focus features and blur features, encode spatial information from each input image. The *recurrent LSTM module* sequentially processes cached latent

Figure 5.2: The rays emitted from an object placed at an axial distance $z$ to the lens, converging at a distance $d_f$ behind the lens. The sensor is situated at a distance $d$ from the focal lens. The pixels are sharply imaged when the sensor is placed right at focus distance, $C$ is the CoC, which grows as the



Figure 5.3: The overview of our proposed network, FocDepthFormer, is presented with its core components: the Transformer encoder, the recurrent LSTM module, and the CNN decoder. Preceding the Transformer encoder, early-stage multi-scale convolutional kernels are depicted within the dashed line. The resulting multi-scale feature maps are concatenated and subjected to spatial and depth-wise convolution. Subsequently, the fused feature map of an image stack is divided into patches, which are then individually projected by a linear embedding layer into tokens. A red token represents a global embedding token mapped from the entire image and is summed with each patch embedding token.

tokens from different frames of a focal stack and fuses them along the stack dimension. This stack feature fusion process is learned in the latent space by the LSTM module. Our attention design with LSTMs enhances the model's capability to handle an arbitrary number of input images. The final disparity map is decoded (denoted as $d(\cdot)$) from the fusion feature, utilizing the aggregated tokens from all images in the stack.

## 5.2.2  Early-stage Encoding with Multi-scale Kernels

To capture low-level focus and defocus features at different scales, we employ an early-stage convolutional encoder with multi-scale kernels, which is different from methods using fixed-size kernel convolution stem before the Transformer [257]. As illustrated in Figure 5.3, the early-stage encoder utilizes three convolutional kernels to generate multi-scale feature maps $f_m(\mathbf{x}), \{m = 1, 2, 3\}$. All feature maps are concatenated and merged into the feature map $\mathbf{x}' \in \mathbb{R}^{H' \times W' \times 1}$ through spatial convolution, followed by $3 \times 3$ and $1 \times 1$ convolution on the feature map depth channel:

$$\mathbf{x}' = \mathcal{F}(\mathbf{x}) = \text{Conv}(\text{Concat}(f_m(\mathbf{x}))), \tag{5.2}$$

where $m$ ranges from 1 to 3. Feature concatenation after convolutions with multiple kernel sizes preserves fine-grained details of features across varying depth scales. The first module from the left in Figure 5.3, comprising parallel multi-scale kernel convolutions followed by depth-wise convolution, ensures the model has a large receptive field beyond the $7 \times 7$ kernel size. This facilitates capturing more defocus features while preserving intricate details.

## 5.2.3  Transformer with LSTM

**Transformer encoder.** The Transformer in Figure 5.3, denoted as $g(\cdot)$, processes the feature maps $\mathbf{x}'$ from the preceding early-stage multi-scale convolutions to generate a series of tokens $(\mathbf{t}'_p)_{p=1}^k$:

$$\mathbf{t}'_1, \mathbf{t}'_2, ..., \mathbf{t}'_k = g(\mathbf{x}'). \tag{5.3}$$

Specifically, the early kernel CNNs and Transformer encoder take the focal stack images sequentially, cache and concatenate the feature maps of a certain stack of images into $\mathbf{x}'$. Transformer uses first a linear embedding layer, which divides the feature maps $\mathbf{x}'$ into $k$ patches of size $16 \times 16$. That is, $\mathbf{x}'_p \in \mathbf{x}', p = 1, 2, ..., k$, is projected by a linear embedding layer (MLP) into corresponding embedding tokens $(\mathbf{l}_p)_{p=1}^k$, each token with a dimension of 768 (576 in total), and all the tokens of a full stack $N$ are cached into $\{(\mathbf{l}_p)_{p=1}^k\} \times N$

before LSTM to be fused at once. The Transformer's *Position Embedding* encodes the positional information of the image patches of a single frame in an iterative order from the top-left of the image. An MLP layer generates the Global Embedding Token (Figure 5.3) by mapping the entire image into a global token and then adding each patch embedding token. Each linear embedding token is projected via a weight matrix $\mathbf{W}^{(\cdot)}$ from dimension $d_m$ into three vectors: the Query, $\mathbf{l}_Q$, the Key, $\mathbf{l}_K$, and the Value, $\mathbf{l}_V$, with dimensions $d_Q$, $d_K$, and $d_V$ respectively. The Queries, Keys, and Values pass through the Multi-Head Attention (MHA) units in parallel.



Figure 5.4: To illustrate the LSTM module in our network, the initial step involves grouping all cached output tokens from the Transformer encoder into activated and non-activated tokens. These two groups are then individually processed, with activated tokens undergoing LSTMs followed by max pooling and non-activated tokens undergoing average pooling. Following this, the output tokens undergo reshaping and concatenation before being fed into the CNN decoder for predicting the depth map.

$$\text{MHA}(\mathbf{l}_Q, \mathbf{l}_K, \mathbf{l}_V) = (\text{head}_1 \oplus ... \oplus \text{head}_N)\mathbf{W}^O, \tag{5.4}$$

$$\text{head}_i = \text{softmax}\left(\frac{\mathbf{l}_Q \mathbf{W}^{\mathbf{l}_Q} \mathbf{l}_K \mathbf{W}^{\mathbf{l}_K}}{\sqrt{d_k}}\right) \mathbf{l}_V \mathbf{W}^{\mathbf{l}_V}, \tag{5.5}$$

where $\mathbf{W}^{\mathbf{l}_Q} \in \mathbb{R}^{d_m \times d_V}$, $\mathbf{W}_i^{\mathbf{l}_K} \in \mathbb{R}^{d_m \times d_K}$, $\mathbf{W}^{\mathbf{l}_V} \in \mathbb{R}^{d_m \times d_V}$, and $\mathbf{W}^O \in \mathbb{R}^{d_m \times d_V}$. Following the Multi-Head-Attention modules contained in the encoder $g(\mathbf{x}')$, the resulting tokens $(\mathbf{t}'_p)_{p=1}^k$ encode features that differentiate between focus and defocus cues from different stack image patches at the same spatial location within the image. This capability is illustrated in Figure 5.5. Consequently, the more in-focus and sharp features of the image patches receive increased attention in the embedding space.

Figure 5.5: Comparison of Transformer attention on the two left column images. Cropped Image patches within green and orange boxes in (a) and (c) are used as the query input to calculate the self-attention map over the whole input image, respectively. In (b) and (d), the attention map on the left and right of the green line represents the attention output of the green and orange boxes, respectively. This demonstrates that the patches can selectively attend to both foreground and background areas, exhibiting the ability to differentiate the focus and defocus cues.

**LSTM module.** To attain the flexibility of our model in handling stacks with arbitrary lengths, as opposed to the fixed length in existing methods [264, 246, 79], we employ an LSTM to progressively fuse sharp features along the stack. The LSTM treats patch embedding tokens at the same image position as the sequential features along the stack dimension. The corresponding feature tokens $(\mathbf{t'}_p)_{p=1}^k$ from stack images at stack number $N$ are

spatially ordered and fed into LSTM modules, which are arranged in the original spatial image order. At each position, each LSTM module incrementally fuses the latent token $\mathbf{t}'_p$ from a stack. This approach differs from existing models constrained to a 3D volume stack with a pre-defined and fixed size [264, 246, 79]. Importantly, the sequential processing in the latent space after a shared encoder for the stack images incurs limited complexity in practice.

Preceding the LSTM modules, tokens associated with image patches from a single frame are categorized into activated and non-activated tokens, reflecting the informative level of the features. The $L_2$ norm, denoted as $\| \cdot \|$, of each embedding token is compared with a threshold of 0.4, as shown in Figure 5.4. Specifically, for tokens within a single frame, only the activated tokens, identified through this threshold comparison, are forwarded to the LSTM, with the number of $k_1$ tokens. This operation reduces the computational complexity of the LSTM by applying the operation solely to a fraction of the latent tokens:

$$\mathbf{t}_p^n, \mathbf{h}_p^n = LSTM(\mathbf{t}'^n_p, \mathbf{h}_p^{n-1}, c^n), \tag{5.6}$$

this is a single LSTM layer expression, where $p = 1, 2, ..., k_1$, and $n$ is the frame index number of a stack. We set the number of hidden layers of the LSTM module to be the same as the stack size $N$. The memory cell $c$ undergoes continuous updates at each step, influenced by the input $\mathbf{t}'^n_p$ and the hidden state $\mathbf{h}$. Then all the LSTM layer outputs are merged by max pooling $\max\{\mathbf{t}'^1_p, ..., \mathbf{t}'^n_p\}$ into $\mathbf{t}'_p$. For the non-activated tokens $(\mathbf{t}'_p)_{p=k-k_1}^k$, an averaging operation is performed with the corresponding cached tokens from the previous step at the same embedding position. Finally, the two groups of output tokens are arranged as the original input embedding order, yielding the final fused tokens $(\mathbf{t}_p)_{p=1}^k$. **CNN decoder.** Our decoder $d(\cdot)$ follows the approach presented by Ranftl *et al.* [197], employing Transpose-convolutions to integrate feature maps after LSTMs. The decoder also incorporates feature maps of $i - th$ layer $g_i(\mathbf{x}'))$ from the encoder through skip connections, as depicted in Figure 5.3. Finally, the decoder $d(\cdot)$ predicts the depth.

$$\hat{D} = d((\mathbf{t}_p)_{p=1}^k, g_i(\mathbf{x}')), \quad i = \{1, 2, 3\}. \tag{5.7}$$

### 5.2.4  Training Loss

Our training loss comprises the Mean Squared Error (MSE) loss, denoted as $\mathcal{L}_{MSE}$, and a sharpness regularizer $\mathcal{L}_{log}$ weighted by $\alpha$:

$$\mathcal{L}_{total} = \mathcal{L}_{MSE}(\hat{D}, D) + \alpha\mathcal{L}_{\log}(\delta_{\mathbf{\Delta}\hat{D}}, \delta_{\mathbf{\Delta}D}), \tag{5.8}$$

where $D$ represents the ground truth depth, and $\hat{D}$ indicates the predicted depth. The $\boldsymbol{\Delta}$ is the laplacian operator, applied to predicted and ground truth depth images respectively. $\delta$ is the variance of the depth image. The regularizer item is formulated as $\log(\delta_{\boldsymbol{\Delta}\hat{D}}/\delta_{\Delta D})$. The pixel blurriness due to out-of-focus can be described by the Circle-of-Confusion (CoC),

$$\sigma = \frac{C}{2 \cdot r} = \frac{1}{2r} \frac{f^2}{N(z - d_f)} \Big| 1 - \frac{d_f}{z} \Big|, \tag{5.9}$$

where, $N$ denotes $f$-number, as a ratio of focal length to the valid aperture diameter, and $r$ is the CMOS pixel size. $C$ is the Circle of Confusion (CoC) diameter. $d_f$ is the focus distance of the lens. $z$ represents the distance from the lens to the target object. In general, the range of $z$ is $[0, \infty]$. However, in reality, the range is always constrained by lower and upper bounds. The goal of the model is to learn to shape a depth map from the focus/defocus features.

### 5.2.5  Pre-training with Monocular Depth Prior

Focal stack datasets are typically limited in size due to the high cost and challenges associated with data collection. To address the scarcity of data and fully exploit the potential of the Transformer, we can optionally pre-train the Transformer encoder on widely available monocular depth estimation datasets, such as NYUv2 [226], to enhance spatial representation learning. This can be achieved by bypassing the latent LSTM module while back-propagating the gradients of the encoder $g(\cdot)$ and decoder $d(\cdot)$ weights only. The LSTM-based fusion design after Transformer allows it to process an arbitrary number of input images, making it leverage pre-training before monocular datasets. While monocular depth datasets may differ from focal stack datasets, pre-training facilitates the learning of a versatile visual representation for depth prediction, thanks to our separate spatial image and focal stack learning structure design. Notably, our model demonstrates plausible performance even in the absence of pre-training, as evidenced by the experiments presented below. For detailed parameters of each network block module, please refer to the figure below,

## 5.3  Experiments and Results

The experiment is performed on four different datasets, and compared with the four baseline models for depth estimation from focal stack. The comparisons include both metric and visual results. After that, the detailed ablation

study analysis is performed to evaluate each module block contribution, including data plots, visual results, and metric results.

## 5.3.1 Experimental settings

Experiment results are composed of the **Datasets.** We conducted extensive evaluations of our model using four benchmark focal stack datasets: DDFF 12-Scene [79], Mobile Depth [19], LightField4D [85], and FOD500 [264]. Additionally, our model offers the flexibility of pre-training on the monocular RGB-D dataset NYUv2 [226]. Specifically, we conducted training on DDFF 12-Scene and FOD500 separately for the subsequent experiments, while Mobile Depth and LightField4D were employed for evaluating the model's generalizability with pre-trained on DDFF 12-Scene only. A comprehensive summary of the evaluation datasets, including their properties and captured sensor types, is presented in Table 5.1.

**Evaluation metrics.** The metrics used in our work for quantitative evaluations, are defined as follows,

$$RMSE : \sqrt{\frac{1}{|M|} \sum_{p \in \mathbf{x}} \|f(\mathbf{x}) - \mathbf{D}\|^2}, \tag{5.10}$$

$$logRMSE : \sqrt{\frac{1}{|M|} \sum_{p \in \mathbf{x}} \|log f(\mathbf{x}) - \mathbf{D}\|^2}, \tag{5.11}$$

$$absRel : \frac{1}{|M|} \sum_{p \in \mathbf{x}} \frac{|f(\mathbf{x}) - \mathbf{D}|}{\mathbf{D}}, \tag{5.12}$$

$$sqrRel \frac{1}{|M|} \sum_{p \in M} \frac{\|f(\mathbf{x}) - \mathbf{D}\|}{\mathbf{D}}, \tag{5.13}$$

$$Bump : \frac{1}{|M|} \sum_{p \in \mathbf{x}} \min(0.05, \|\mathbf{H}_\Delta(p)\|) \times 100, \tag{5.14}$$

$$Accuracy(\delta) : \max\left(\frac{f(\mathbf{x})}{\mathbf{D}}, \frac{\mathbf{D}}{f(\mathbf{x})}\right) = \delta < threshold,$$
$$\%\ of\ \mathbf{D}, \tag{5.15}$$

where $\Delta = f(\mathbf{x}) - \mathbf{D}$ and $\mathbf{H}$ is the Hessian matrix. The accuracy threshold of bumpiness (bump) is set at three levels ($1.25$, $1.25^2$, and $1.25^3$).

**Implementation details.** For the evaluation on DDFF 12-Scene, we conducted training experiments using our model with and without pre-training on NYUv2 [226], presenting results for both scenarios. We employed a patch size of $16 \times 16$ and an image size of $384 \times 384$ for the Transformer. Our

Table 5.1: Summary of evaluation datasets.

| Dataset | Image source | GT type | Cause of defocus |
|---|---|---|---|
| DDFF 12-Scene [79] | Real | Depth | Light-field settings |
| Mobile Depth [19] | Real | — | Real |
| LightField4D [85] | Real | Disparity | Light-field settings |
| FOD500 [264] | Synthetic | Depth | Synthesis blendering |

network utilizes the Adam optimizer with a learning rate of $1 \times 10^{-4}$ and a momentum of 0.9. The regularization scalar $\alpha$ in Eq. (5.8) is set to 0.2. In terms of hardware configuration, all training and tests below were conducted on a single Nvidia RTX 2070 GPU with 8GB of VRAM.

Table 5.2: Evaluation results on DDFF 12-Scene. The best results are denoted in **Red** while <u>Blue</u> indicates the second-best. $\delta = 1.25$.

| Model | RMSE↓ | logRMSE↓ | absRel↓ | sqrRel↓ | Bump↓ | $\delta$ ↑ | $\delta^2$ ↑ | $\delta^3$ ↑ |
|---|---|---|---|---|---|---|---|---|
| DDFFNet [79] | 2.91e-2 | 0.320 | 0.293 | 1.2e-2 | 0.59 | 61.95 | 85.14 | 92.98 |
| DefocusNet [161] | 2.55e-2 | 0.230 | 0.180 | 6.0e-3 | 0.46 | 72.56 | 94.15 | 97.92 |
| DFVNet [264] | 2.13e-2 | 0.210 | <u>0.171</u> | 6.2e-3 | 0.32 | 76.74 | 94.23 | 98.14 |
| AiFNet [246] | 2.32e-2 | 0.290 | 0.251 | 8.3e-3 | 0.63 | 68.33 | 87.40 | 93.96 |
| Ours (w/o Pre-training) | <u>2.01e-2</u> | <u>0.206</u> | 0.173 | <u>5.7e-3</u> | <u>0.26</u> | <u>78.01</u> | <u>95.04</u> | <u>98.32</u> |
| Ours (w/ Pre-training) | **1. 96e-2** | **0.197** | **0.161** | **5.4e-3** | **0.23** | **79.06** | **96.08** | **98.57** |

Table 5.3: Evaluation results on FOD500 test dataset. Here the first 400 FOD500 focal stacks are used for training, following the standard setting from DFVNet [264]. The best results are denoted in **Red**, while <u>Blue</u> indicates the second-best. $\delta = 1.25$.

| Model | RMSE↓ | logRMSE↓ | absRel↓ | sqrRel↓ | Bump↓ | $\delta$ ↑ | $\delta^2$ ↑ | $\delta^3$ ↑ |
|---|---|---|---|---|---|---|---|---|
| DDFFNet [79] | 0.167 | 0.271 | 0.172 | 3.56e-2 | 1.74 | 72.82 | 89.96 | 96.26 |
| DefocusNet [161] | 0.134 | 0.243 | 0.150 | 3.59e-2 | 1.57 | 81.14 | 93.31 | 96.62 |
| DFVNet [264] | <u>0.129</u> | <u>0.210</u> | <u>0.131</u> | <u>2.39e-2</u> | <u>1.44</u> | 81.90 | <u>94.68</u> | <u>98.05</u> |
| AiFNet [246] | 0.265 | 0.451 | 0.400 | 4.32e-1 | 2.13 | <u>85.12</u> | 91.11 | 93.12 |
| Ours (w/o Pre-training) | **0.121** | **0.203** | **0.129** | **2.36e-2** | **1.38** | **85.47** | **94.75** | **98.13** |

**Runtime.** We evaluated the runtime of the proposed method and baseline approaches by executing them on focal stacks from DDFF 12-Scene. Our FocDepthFormer processes a stack with 10 images sequentially in 15ms, an average of 2ms per image. In comparison, DDFFNet [79] requires 200ms for each stack under the same conditions, and DFVNet [264] performs in the range of 20-30ms.

## 5.3.2 Baseline Comparisons

DDFFNet [79] and DefocusNet [161] lacked pre-trained weights; therefore, we utilized their open-source codebases to train the networks from scratch. No-

tably, DefocusNet [161] offers two architectures, and we chose the "PoolAE" architecture due to its consistently good performance for comparison. Conversely, for AiFNet [246] and DFVNet [264], we employed the pre-trained weights provided by the authors to conduct the following evaluations.

**Results on DDFF 12-Scene.** Table 5.2 presents the quantitative evaluation results of our model on the DDFF 12-Scene dataset. As the ground truth for the "test set" is not publicly available, we adhere to the standard evaluation protocol used in other comparative works, assessing the models on the "validation set" as per the split provided by DDFFNet [79]. Moreover, we demonstrate that our model, trained on DDFF-12, performs robustly on other completely unseen datasets, highlighting its generalization ability and mitigating concerns of over-fitting. The results in the table indicate that our model without pre-training outperforms prior models on all metrics, except for absRel compared to DFVNet. Specifically, our model achieves an accuracy of 78.01% ($\delta = 1.25$), marking a notable improvement of 1.27% over DFVNet. The Bumpiness metric also reflects this superiority with a value of 0.26, one-third less than DFVNet. This can be attributed to the compact design of the Transformer and LSTM, efficiently learning spatial features and stack features separately. Examining the table, pre-training brings considerable advantages across all metrics, particularly for absRel and Bump, with improvements of around 7% and 12%, respectively, over pure training. This underscores the potential of our model. Figure 5.6 illustrates the qualitative performance of our model on DDFF 12-Scene. The visualizations showcase depth estimation results preserving fine-grained details such as the thin wire over the sofa (first row) and the cup handle on the shelf (second row).

**Results on FOD500.** Table 5.3 presents the quantitative evaluation of our model on the synthetic FOD500 dataset. For testing, we use the last 100 image stacks from the dataset, while the initial 400 image stacks are reserved for training. DDFFNet and DefocusNet are re-trained on FOD500 from scratch. The results highlight the consistent superiority of our model across all metrics when compared to the baseline methods. Notably, in our experiments, we observed that pre-training on NYUv2 did not provide significant benefits, likely due to the gap between synthetic and real data. Interestingly, we also noted that the proposed method can achieve satisfactory and competitive results with only a few training epochs.

### 5.3.3 Cross Dataset Evaluation

To evaluate the generalizability of our model, it is initially trained on DDFF 12-Scene and subsequently evaluated on the Mobile Depth and LightField4D datasets. The Mobile Depth dataset poses a challenge as it comprises 11

Input    GT    DDFF    DefocusNet    AiFNet    DFVNet    Ours



Figure 5.6: Qualitative evaluation of our model on DDFF 12-Scene dataset.

Input    DDFF    DefocusNet    AiFNet    DFVNet    Ours



Figure 5.7: Qualitative evaluation of our model on Mobile Depth dataset.

Input    GT    DDFF    DefocusNet    AiFNet    DFVNet    Ours



Figure 5.8: Qualitative evaluation of our model on LightField4D dataset.

aligned focal stacks captured by a mobile phone camera, each with varying numbers of focal planes and lacking ground truth. Results on the Mobile Depth dataset are illustrated in Figure 5.7, showcasing the model's ability to preserve sharp information for depth prediction in complex scenes, such as the ball with many holes in the first row. Notably, the model excels in recognizing fine details in complex topological structures, such as grape granules. Additionally, our model effectively fuses depth information across a diverse range of scenes, including backgrounds like bananas in the second row. An advantage of our model is its capability to handle varying numbers of input images, a feature not supported by baseline methods that are restricted to fixed training settings. For generalization tests on the LightField4D dataset, our model, along with other baseline models, is pre-trained exclusively on DDFF 12-Scene, except for AiFNet [246], which is also pre-trained on LightField4D all-in-focus color images in an unsupervised manner. Visual results are presented in Figure 5.8, illustrating that our model produces more accurate depth maps for complex objects compared to AiFNet [246]. Overall, these test results affirm that our model demonstrates comparable generalization performance to previous models without specific pre-training on specific focal stack datasets. We present the quantitative results for cross-dataset evaluation of our model on the LightField4D dataset in Table 5.4. Our model achieves a comparable performance in terms of accuracy (58.90%) on this completely unseen dataset. Although the AiFNet can attain the least RMSE (0.231) and logRMSE(0.407) error, our model generates smoother boundaries and fine-grained details, indicated by a lower bumpiness value (2.53). We use the available pre-trained AiFNet on LightField3D which uses the all-in-focus color image as supervision for evaluation. It also explains why AiFNet achieves satisfactory performance, and it further validates our model's good generability without using any supervision signal from this new focal stack dataset, while our model can achieve plausible performance only by learning a good stack distribution and spatial sharp feature representation from DDFF 12-Scene dataset.

Table 5.4: Metric evaluation results on "additional" set of LightField4D dataset. The best results are denoted in **Red**, while Blue indicates the second-best.

| Model | RMSE↓ | logRMSE↓ | absRel↓ | Bump↓ | $\delta(1.25)$ ↑ |
|---|---|---|---|---|---|
| DDFFNet | 0.431 | 0.790 | 0.761 | 2.93 | 44.39 |
| DefocusNet | 0.273 | 0.471 | 0.435 | 2.84 | 48.73 |
| DFVNet | 0.352 | 0.647 | 0.594 | 2.97 | 43.54 |
| AiFNet | **0.231** | **0.407** | 0.374 | 2.53 | 55.04 |
| Ours | 0.237 | 0.416 | **0.364** | **1.54** | **58.90** |

### 5.3.4 Ablation Study

The following tests all use the DDFF 12-Scene validation set and we omit the data name for brevity.

**Transformer encoder:** Table 5.5 provides a comparative analysis between Transformer and CNN-based encoders, focusing on the vanilla Transformer, Swin Transformer, and the CNN-based DDFF-Net model. To ensure fairness in the comparison, we conducted experiments by excluding the early kernels from our model. The results reveal that the ViT-based model combined with LSTM has the highest accuracy in this experimental setup, showcasing an approximately 33% improvement over the CNN-based DDFF-Net encoder in terms of RMSE.

Table 5.5: Metric evaluation of various encoders.

|  | RMSE↓ | absRel↓ | Bump↓ |
|---|---|---|---|
| ViT-base encoder | **2.06e-2** | **0.197** | **0.29** |
| Swin Transformer encoder | 2.21e-2 | 0.205 | 0.32 |
| CNN encoder | 3.12e-2 | 0.268 | 0.46 |

We present the attention heat map of different focus or out-of-focus images patched over the whole image in Figure 5.5. The attention of the Transformer module can attend to more in-focus feature information, *e.g.*, the toy on the right side of Figure (d) is with higher attention compared to Figure (b) at the same location. It shows that the patches can attend to the fore and background regions with related focus and defocus cues of the corresponding depth field quite well. It further manifests that the proposed model based on Transformer attention can differentiate the pixel sharpness variance on the image input. Although some attention is put wrongly, as in Figure (c), where the attention on the monitor and toy is incorrect, most of the attention is distributed consistently with the input patch appearance. Additionally, the attention of the chosen channel of the Transformer encoder for visualization is not scattered around the whole image, which further discloses that the attention is mainly focused on the similar semantic, sharpness, and appearance information of the input patch.

**Multi-scale early-stage kernels:** Table 5.6 presents a comparison of different early CNN kernel design configurations in conjunction with our Transformer encoder (ViT). The effectiveness of the proposed *early-stage multi-scale kernels* encoder is evident, demonstrating robust performance. Conversely, omitting multi-scale kernels or forgoing subsequent convolutions after in-parallel convolutions leads to a degradation in model performance.

**Fusion by LSTM.** The LSTM module facilitates our network to fuse each image from the focal stack incrementally, which extends the model capabil-

Table 5.6: Results of different designs of the early-stage conv.

|  | RMSE↓ | absRel↓ | Bump↓ |
|---|---|---|---|
| *multi-scale kernels* w/ depth convos + ViT | **2.01e-2** | **0.173** | **0.26** |
| Constant kernel size at $3 \times 3$ + ViT | 2.18e-2 | 0.216 | 0.31 |
| *multi-scale kernels* w/o depth convos + ViT | 2.27e-2 | 0.229 | 0.29 |

ity to varying focal stack lengths. Figure 5.9 depicts the fusion process of ordered input images of one stack. As the images from the stack are given sequentially, starting from the in-focus plane close to the camera, the model can fuse the sharp in-focus features from various frames to attain a final all-in-focus prediction depth map at the bottom right figure.



(a) 1st image.          (b) 5th image.          (c) 10th image.

(d) Depth prediction of (e) Depth prediction of (f) Fusion prediction of
1st input frame.        two input frames.       3 input frames.

Figure 5.9: The top row is the input, and the bottom is the output disparity map. The final disparity map is at the bottom right. The red rectangle highlights the incremental fusion results of depth information in the background.

In table 5.7, we compare our model to its base model without the LSTM module. For the performance without LSTM, the encoder, and decoder are connected directly, and all the depth maps of each image in the stack are averaged to get the metric results of a whole stack. The results indicate the necessity and importance of LSTM in-depth estimation from the focal stack problem. It further validates that the modeling stack information sep-

arately from the image spatial features can help to improve depth prediction accuracy.

Table 5.7: Metric evaluation on different settings for LSTM module on DDFF 12-Scene validation dataset.

| Structure design | RMSE↓ | absRel↓ | Bump↓ |
|---|---|---|---|
| Model w/o LSTM module | 3.68e-2 | 0.324 | 0.37 |
| Full model | **1.92e-2** | **0.161** | **0.19** |

To justify the model structure design, we further implement the tests by replacing The Transformer encoder with the CNN encoder, and then the features after convolution are fused by LSTM. The pure Transformer-based encoder without LSTM fusion takes the image stack as a whole, and then concatenates the feature maps to be decoded into the depth map, The Trans-



(a) Ours  (b) LSTM+CNN  (c) Transformer

Figure 5.10: Different model structures' comparison.

former + LSTM design proposed in the paper can predict a more detailed feature map with fine-grained details, while the naive concatenation of feature maps after the Transformer can not achieve the equivalent performance (as reflected in the middle image) compared to our proposed structure, by combining the focus/defocus cues directly. The CNN encode even with the help of LSTM can not learn a good image feature representation after the encoder for depth map prediction, due to the local receptive field limits, resulting in the blur effect in the final prediction. It is noteworthy that for pure Transformer without LSTM, we downsample the focal stack size to three images for these comparison experiments, as the original size of the image results in out-of-memory in pure Transformer-based encoder implementation, which takes the whole stack of images at once and processes them with multiple encoders in parallel.

**LSTM for handling arbitrary stack length.** Our proposed LSTM-based method exhibits flexibility in processing focal stacks of arbitrary lengths, a feature distinguishing it from designs that limit inputs to fixed lengths. To illustrate the advantages of our LSTM-based model, we conducted experiments using DDFF 12-scene data. The results, including the RMSE comparison between our model and DFVNet [264], are presented in Table 5.8. Initially, we trained our model and DFVNet using 10-frame (10F) stacks (*i.e.,* Ours-10F and DFVNet-10F). During testing, DFVNet-10F is constrained and cannot process stacks with fewer than 10 frames. In contrast, due to the fixed stack size requirements during training and testing, DFVNet must be retrained for different stack sizes (DFVNet-#F). In comparison, our model is trained once on 10-frame stacks and can be directly tested with varying numbers of stack images. The focal stack images are ordered based on focus distances. Despite our model's initial performance being inferior to DFVNet, the learning curve of LSTM indicates rapid convergence.

Table 5.8: RMSE for evaluation of LSTM compared to DFVNet.

| Model | 2 Frames | 4 Frames | 6 Frames | 8 Frames | 10 Frames |
|---|---|---|---|---|---|
| Ours-10F | 3.2e-2 | **2. 61e-2** | **2.18e-2** | **2.16e-2** | **2.04e-2** |
| DFVNet-10F | —- | —-- | —- | —- | 2.43e-2 |
| DFVNet-#F | **2.97e-2** | 2.70e-2 | 2.52e-2 | 2.47e-2 | 2.43e-2 |

We perform experiments to observe the results of our model on various focal stack sizes. Figure 5.11 shows the findings of our experiments on a focal stack from the DDFF 12-Scene validation set, with multiple objects placed at varying depths. In the figure, we plot RMSE and accuracy (in percentage) ($\delta = 1.25$) w.r.t. the number of focal stack images. We observe that the model accuracy increases with more input images from the focal stack in use, while the RMSE decreases correspondingly. Our model achieves a decent performance around six frames with a notable increase, then followed by a marginal increase after six frames.

To evaluate the impact of the token $L_2$ norm threshold, we further tested the RMSE under the various thresholds. A lower threshold smaller than 0.4 like 0.3 can improve the RMSE by a marginal increase at around 0.6%, yet scaling the model parameter complexity a lot, because more and more tokens are activated, and they are thrown into the LSTM module for processing. The threshold lower than 0.3 can even lead to out-of-memory issues. Conversely, the large threshold can reduce the LSTM number in use but also sacrifice the accuracy raging from 0.8% to 1.2% as the threshold increased from 0.5 to 0.8. **Loss Function:** Table 5.9 presents the results obtained by employing three distinct loss functions: Mean Squared Error (MSE), Mean

Figure 5.11: Our model performance w.r.t. the frame size of one focal stack sample from DDFF 12-Scene test.

Absolute Error (MAE), and Regularized MSE (MSE with a gradient regularizer). The findings indicate that MSE loss consistently outperforms MAE in terms of overall performance. Notably, the inclusion of the gradient regularizer contributes to achieving the highest accuracy, as evidenced by the bumpiness metric (0.26).

Table 5.9: Evaluation for our model with different losses.

|  | RMSE↓ | absRel↓ | Bump↓ |
|---|---|---|---|
| MSE loss | 2.94e-2 | 0.280 | 0.50 |
| MAE loss | 3.76e-2 | 0.372 | 0.62 |
| MSE + Gradient loss | **2.01e-2** | **0.173** | **0.26** |

**Pre-training:** Table 5.10 illustrates the impact of pre-training on the performance of the proposed method. The results showcase an enhancement in the model's capabilities through pre-training, leveraging the advantages of the compact design with Transformer and LSTM. Even in the absence of pre-training, our model demonstrates competitive results. Notably, attempts were made to apply pre-training to DFVNet by creating a stack from

repeated monocular images of NYUv2. However, the pre-training did not yield any improvement for DFVNet and, in some instances, led to performance degradation due to the data modality gap.

Table 5.10: Pre-training contribution comparisons.

|        | Pre-training | RMSE↓   | logRMSE↓ | absRel↓ | Bump↓ |
|--------|--------------|---------|----------|---------|-------|
| Ours   | ✗            | 2.01e-2 | 0.206    | 0.173   | 0.26  |
|        | ✓            | **1.96e-2** | **0.197** | **0.161** | **0.19** |
| DFVNet | ✗            | 2.13e-2 | 0.210    | 0.171   | 0.32  |
|        | ✓            | 2.57e-2 | 0.233    | 0.184   | 0.49  |

# 5.4 Conclusion

We introduced FocDepthFormer, a novel model for depth estimation from focal stacks. The key innovation lies in its hybrid architecture, which combines a Transformer encoder for capturing global spatial relationships with an LSTM module in latent space to effectively aggregate information across focal stacks of varying lengths. This design enhances flexibility and generalization compared to traditional CNN-based approaches.

While FocDepthFormer achieves strong performance, its use of Transformers leads to a higher memory footprint than simpler CNN-based models. Future work will explore more efficient Transformer architectures to mitigate this limitation. Additionally, we aim to extend our framework to image synthesis with defocus effects, further leveraging focal stacks for advanced 3D vision applications. Lastly, depth prediction from focal stacks faces challenges due to the limited depth range and the need for a sufficient number of focal stack images to estimate depth with adequate resolution. Successful training relies on distinct focus and defocus signals between images, which becomes difficult in textureless scenes, such as outdoor environments or plain walls, where such signals are absent. In these cases, inferring depth is highly ambiguous. A promising direction may involve combining focal stack techniques with state-of-the-art deep prior models trained on monocular or video inputs to recover fine-grained details more robustly.

# Chapter 6

# 3D Reconstruction Using Implicit SDF with Wavelet Feature-Based Prior

## Abstract

3D reconstruction aims to recover the underlying structure of a scene from sensor data, primarily LiDAR point clouds or multi-view images. Point clouds provide a direct geometric representation but are limited by sampling density and resolution. Multi-view images, on the other hand, leverage photometric cues for dense surface reconstruction. Recent advancements in deep learning have enabled implicit representations, such as Signed Distance Fields (SDF), which model continuous surfaces and can be converted into explicit 3D geometry using techniques like Marching Cubes. Implicit SDF models process both point cloud and image inputs. For point clouds, they estimate the signed distance of query points to the nearest surface, defining the zero-crossing isosurface. For multi-view images, SDF values are predicted along sample rays, optimizing the implicit representation for view-consistent rendering. While these models offer superior topological accuracy, they struggle to capture fine-grained geometric details due to high-frequency information loss during feature extraction, leading to suboptimal multi-scale representation. To address this limitation, we propose a wavelet-conditioned implicit SDF model that enhances geometric fidelity by integrating a pre-trained wavelet autoencoder optimized with sharp depth maps[1]. This autoencoder extracts multi-scale wavelet-transformed features, which are fused

---

[1]The majority of this chapter is based on the paper submitted to Transactions on Machine Learning Research (under review).
**Author Contributions:**

with implicit 3D triplane representations, preserving structural details more effectively. Our approach serves as a plug-and-play module that seamlessly integrates with existing implicit SDF frameworks. Extensive evaluations on DTU, Tanks, and Temples, and a cultural heritage dataset demonstrate that our model outperforms state-of-the-art implicit and explicit methods, producing more complete and detailed 3D reconstructions across various scene scales—from small objects to large architectural structures.

## 6.1 Introduction

3D reconstruction is closely tied to the underlying 3D representation format. Various representations exist, each with advantages and limitations. For instance, voxel grids discretize 3D space into uniform cubic elements, making them compatible with neural networks but suffering from cubic memory complexity $O(n^3)$, limiting resolution and introducing a "Manhattan world bias" [160]. Meshes represent surfaces efficiently using vertices and faces as an approximation of continuous geometry representation, but mesh also has self-intersection issues for complex shapes [183]. Point clouds provide a simple, lightweight representation by recoding the 3D point coordinates [53], but such discrete point cloud representation lacks connectivity information, and the underlying topology is ignored.

A more recent alternative is implicit functions, which model shapes continuously rather than as discrete elements. These functions map 3D coordinates (and optional conditions like images) to either occupancy probabilities $[0, 1]$ [185] or signed distance values $[-D, D]$ [262], capturing fine details while being memory efficient.

Before the advent of deep learning, Truncated Signed Distance Functions (TSDFs) were already widely used for 3D reconstruction. At first, a ground truth (GT) mesh is voxelized accordingly, with each voxel storing the signed distance value to the nearest mesh surface. The positive TSDF values indicate points outside the surface, while negative TSDF values indicate interior regions. The distance values are usually truncated to a fixed distance range centered on the surface, ensuring a finite update to avoid numeric instability.

- **Xueyang Kang**: Idea Design, Methodology, Software, Experiment Validation, Formal Analysis, Data Curation, Writing, Review, and Editing.

- Hang Zhao: Data Curation and Review.

- Kourosh Khoshelham: Review and Supervision.

- Patrick Vandewalle: Review and Supervision.

This volumetric representation is integrated into software product kits like KinectFusion [101] for real-time 3D reconstruction from RGB and Depth views. However, TSDF-based methods are constrained by cubic memory complexity, limiting its scalability for large-scale scenes.

To address these limitations, implicit SDFs have been used as a powerful solution for 3D reconstruction. Instead of using explicit voxel grids with TSDF values, an implicit SDF represents the watertight surface as a neural function $\mathbf{f}(\mathbf{x})$, where MLP layers are trained to predict signed distance values for any query point $\mathbf{x}_i$. This approach enables high-fidelity surface reconstruction, moreover, it has a lower memory footprint and is unlimited to the input resolution.

Implicit SDFs can take in two modalities usually: point clouds and multi-view images. For point clouds, the model learns a continuous surface representation by fitting the implicit function to the sparse input point samples. For multi-view images, it leverages photometric consistency loss, along with some other geometric constraints like depth or normal losses, to learn 3D structure. The 2D cross-section snapshot in Figure 6.1 illustrates how TSDF values are distributed in the proximity of the surface, from far (blue) to near (red) for the surface.



Figure 6.1: 2D cross-section of the implicit SDF volume space, where red indicates regions near the surface, and blue represents areas farther from the surface.

To efficiently generate Ground Truth (GT) SDF values for training an implicit SDF model, we primarily sample points near the surface. Points in free space, far from the surface, have large distance values, which can lead

to gradient instability during training. By concentrating on points near the surface, the model learns a more stable and continuous function. Figure 5.1 illustrates this sampling strategy, with colors indicating point positions from left to right.



Figure 6.2: Sampling points near the surface are generated to compute Signed Distance Values, providing training input for GT SDF calculation.

Sampling methods in the surface region can be categorized into two approaches, as illustrated in Figure 6.3. When the GT mesh is available, points are sampled by shifting along the surface normal (red arrow) and its opposite direction, scaled by a specified distance to create an envelope around the surface. The shift magnitude determines the SDF value.

Alternatively, points are randomly sampled within a predefined range around the surface, covering both interior and exterior regions. Their closest distance to the surface is then computed to obtain the corresponding SDF value.

Image-based 3D reconstruction methods, such as Structure from Motion (SfM), recover 3D structures from multi-view 2D images [210], yet they are struggling to preserve high-fidelity details. Alternatively, reconstruction from structured light scans [89, 244] or a fusion of images and LiDAR scans [110, 172] has seen active progress, but these point-based methods are prone to noise in scans, making it difficult to obtain a plausible reconstruction mesh. High-fidelity reconstruction with fine structural details remains a core challenge in computer vision. Recent advances in implicit representations, such as neural radiance fields (NeRF) [168], and explicit methods such as Gaussian splatting (GS) proposed [117], have significantly advanced 3D applications.

Implicit models leverage photometric consistency loss to learn Signed Distance Fields (SDFs) from multi-view images [78]. Unisurf [179] unifies surface

Normal vector-based SDF          Sampling-based SDF

Figure 6.3: GT SDF value generation can be achieved using two main approaches: surface normal vector-based sampling (left), where points are shifted along the normal direction and random sampling with Gaussian deviation within a narrow region around the surface.

and volume rendering to improve generalization, while hybrid volume-surface representations can be converted into high-quality meshes for real-time rendering, like the work [271]. Multi-resolution hash grids further enable coarse-to-fine optimization for detailed neural surface reconstruction [134], making implicit SDF models effective for complex topologies and continuous geometry fields. Image-based implicit reconstruction models recover geometry using multi-view photometric consistency under the Lambertian assumption of uniform light reflection. However, for reflective surfaces such as glass or transparent materials, the reconstruction pipeline should also explicitly model external glass planes or mirror surfaces [254, 135, 195, 237] to recover geometry from appearance features better while maintaining the consistency assumption.

Explicit Gaussian splatting (GS) represents scenes by anisotropic 3D Gaussians [117], enabling efficient training and real-time rendering. However, while GS offers speed, it often sacrifices geometric quality. To address this, AGS-Mesh [201] incorporates meshing priors, PGSR proposed [33] enforces planar constraints, DN-Splatter [238] utilizes depth and normal priors for Gaussian-based reconstruction, and 2D GS [93] simplifies 3D Gaussian parameters to improve surface alignment. Gaussian Splatting-based methods can be further enhanced to render reflective surfaces using material shaders [104], or to handle complex mirror reflections [272, 142].

Most prior work emphasizes global shape reconstruction and coarse geometric structures. While some methods incorporate geometric priors to enhance shape representation, they often struggle with fine-grained details due to high-frequency feature loss, as current network architectures have limited band representation capacity, often requiring complex 3D prior integration.

94

To overcome these challenges, we propose a multi-scale wavelet-based feature approach utilizing a pre-trained depth image autoencoder trained on monocular depth priors. Wavelets efficiently capture high-frequency geometric details while preserving spatial localization, unlike Fourier transforms, which lose spatial information. This property is crucial for retaining fine surface details that deep learning models often neglect due to the lack of specialized multi-scale representation. The autoencoder is trained on wavelet-transformed depth images generated by a state-of-the-art monocular depth diffusion model [80]. The extracted wavelet features are aligned with implicit 3D triplane features via triplane projection and fused to enhance SDF predictions. Our method outperforms state-of-the-art reconstruction models across diverse scenes. The main contributions of our work can be summarized as follows:

- **Wavelet-Transformed Depth Feature Conditioning**: We introduce a pre-trained multi-scale wavelet autoencoder for depth image reconstruction. During implicit SDF training, wavelet features extracted from depth maps condition the network, enhancing geometric detail preservation.

- **Triplane-Aligned Wavelet Feature Projection**: A triplane projection strategy aligns 2D wavelet features with 3D implicit representations, ensuring seamless fusion and improved geometric consistency.

- **Hybrid Feature Fusion for SDF Prediction**: A UNet-based fusion mechanism integrates implicit 3D features with wavelet-transformed depth representations, yielding more structured and accurate SDF predictions for isosurface mesh extraction.

## 6.2  Related Work

**Geometry Representation.** 3D geometry representation follows two main paradigms: implicit and explicit. Implicit methods model surfaces via neural radiance fields (NeRF) [168], surface reconstruction like Unisurf [179], or signed distance functions like BakedSDF [271]. Explicit methods, such as Structure from Motion (SfM) [210] and Multi-View Stereo (MVS) [221], reconstruct 3D geometry from multi-view images. Recent advances, like 3D Gaussian Splatting [117], enable real-time rendering while maintaining high fidelity. Each approach balances reconstruction accuracy, efficiency, and rendering quality.

Further refinements address aliasing artifacts, such as Mip-NeRF created [15], and Mip-NeRF 360 created [16], which extends NeRF-based models to large-scale unconstrained environments akin to NeRF in the Wild [163]. Implicit SDF models reconstruct shapes from single images akin to DISN [262] and enhance local geometry with SDF priors [31].

Recent work integrates SDFs with diffusion models for high-fidelity shape generation from text or single image input [222, 284, 40, 132, 37]. Explicit Gaussian-based methods [117] continue evolving: AGS-Mesh [201] incorporates meshing priors, PGSR [33] enforces planar constraints for structured Gaussian point clouds, and DN-Splatter [238] integrates depth and normal supervision for improved reconstruction.

**Spectrum Techniques.** Spectral methods have long played a crucial role in computer vision. The frequency analysis of a Fourier Transform has inspired spectral convolution kernels in CNNs like the work [127] and enabled Fourier Convolutional Neural Networks (FCNN) [190]. Similarly, the Fourier transform has also been integrated into multi-head attention mechanisms for Fourier Transformers [82, 176, 24], enhancing image feature learning. However, Fourier-based features often face training challenges due to the broad frequency distribution and the loss of locality.

Wavelet transforms provide a more localized spectral representation, preserving spatial details lost in the standard Fourier transform. They have been widely applied in image denoising [173, 32], super-resolution [74, 95], and restoration, as well as compression [220, 202, 152] and inpainting [96, 276, 55]. Wavelet autoencoders [62, 36, 169, 205, 209] efficiently represent image features while reducing parameters for lightweight models.

Recent advances extend spectral methods to 3D tasks. Periodic activation functions can be leveraged in implicit MLP to capture repeating geometric patterns [227], while FINER [147] adapts spectral variables for feature learning. Fourier bases have been explored for implicit representations [131], with models like Bacon [140] and BANF [216] make use of progressive learning for band-limited feature capture in 3D reconstruction.

Wavelets have also been integrated into multi-scale triplane radiance fields [120] and SDF diffusion models [87, 285, 98], enhancing shape generation with fine-grained local details. Despite these advances, spectral models still face convergence challenges, and implementing wavelet decomposition in 3D feature spaces remains computationally demanding.

## 6.3  Method

### 6.3.1  Point Cloud Input for Reconstruction

The standard backbone for 3D reconstruction is a 3D CNN-based U-Net, comprising an encoder and a transposed decoder. Figure 6.4 illustrates a 3D U-Net architecture for shape completion and reconstruction from point cloud input. The process begins with an incomplete point cloud, voxelized into a sparse volumetric representation.

The 3D U-Net follows an encoder-decoder structure with 3D convolutional layers (Conv), transposed convolutions (ConvTr), and skip connections. The encoder progressively downsamples the input through multiple convolutional layers (Conv1 to Conv6), capturing hierarchical spatial features and compressing the data into a latent representation. The decoder then upscales the features via transposed convolutions (ConvTr1 to ConvTr5), incorporating skip connections to preserve fine-grained details. Pruning at specific layers maintains sparsity, enabling efficient processing of large-scale data.

Finally, the network reconstructs a dense, completed 3D shape. While this multi-scale feature learning improves computational efficiency, the model struggles to represent complex geometries due to the CNN kernel's limited receptive field.



Figure 6.4: 3D CNN-based UNet for shape point cloud completion, image from Minkowski Sparse CNN demo.

In contrast, the implicit SDF model offers greater flexibility in capturing complex topology and geometry. By leveraging implicit SDF representations, it generates continuous, watertight meshes from sparse point cloud inputs, enabling robust shape completion and reconstruction. Figure 6.5 illustrates the training phase of an implicit Signed Distance Function (SDF) model for mesh reconstruction from point cloud input.

The process begins with a sparse point cloud representing the target shape as raw 3D data. It is then processed through a series of 3D convolutional layers (Conv1–Conv6) to extract hierarchical features, capturing both global and local geometric details. These features are used to predict SDF values, which encode the signed distances of points to the object's surface. The

blue cube in Figure 6.5 represents intermediate occupancy or SDF predictions, while the red cube denotes the fused 3D features from the encoder and voxelized shape points.

Finally, training is guided by two loss functions: binary cross-entropy loss on SDF values and absolute Distance L1 loss between predicted and ground truth SDF values. This dual-head learning strategy enhances reconstruction accuracy.



Figure 6.5: The implicit SDF model for shape reconstruction in the training phase.

During the test/inference phase, the model processes the input sparse point cloud and predicts signed distance values and occupancy probabilities. The Marching Cubes algorithm is then applied to extract a triangular mesh surface from the predicted SDF values of the query points.



Figure 6.6: The implicit SDF model for shape reconstruction in the test phase.

Unlike the recent 3D Gaussian Splatting approach [117], which employs explicit Gaussians to learn 3D representations from multiview images, Gaussian Splatting offers fast training but relies on a discretized representation. While it can render plausible views, its geometric accuracy remains limited. To improve mesh reconstruction, methods such as 2D Gaussian-based [92] and Geometry-aware Gaussian-based approaches [133] introduce surface-aligned geometry constraints. However, these techniques still fall short of the

reconstruction performance achieved by implicit SDF, which benefits from its continuous representation. Exploring hybrid methods that integrate explicit Gaussians with implicit SDF could further enhance reconstruction quality.

## 6.3.2  Multiview Image Input for 3D Reconstruction

The implicit SDF model can also process image inputs by generating sample pixel rays from multiview posed images. Along these rays, sampled points are queried through the 3D implicit MLP layers for representation learning. However, MLP layers excel at capturing low-frequency features but struggle to preserve high-frequency details due to their intrinsic averaging effect across multiple layers.

To enhance reconstruction quality, feature augmentation before the decoder is crucial. One promising approach is conditioning the implicit SDF with a "codebook" representation—a pre-trained set of embedding features encoding the structural characteristics of building point clouds.

Despite its potential, this approach faces several limitations. The pre-trained codebook is specialized for building point cloud data, which may restrict its generalizability to other 3D shapes or environments beyond its feature distribution. Additionally, the limited availability of large-scale 3D building scans makes acquiring diverse and representative training data challenging.

To overcome these challenges, future work could focus on creating more generic and adaptive codebook representations by leveraging 3D synthetic data priors or foundational 3D generation models to develop a universal geometry codebook. Incorporating additional cues, such as semantic segmentation or part-based geometry priors, could further enhance reconstruction capabilities. While this approach illustrates the value of pre-trained priors for improving 3D reconstruction from limited data, continued research is essential to make these methods more robust and broadly applicable.

Recent advancements in 2D foundational models, such as DepthAnythingV1 [266], DepthAnythingV2 [268], and LOTUS [80], have achieved remarkable success in predicting sharp depth and normal maps from single image inputs. These models provide a strong geometry prior, which can be further enhanced using a frequency-based encoder to capture high-frequency features from the input images. These enhanced features are then aligned with 3D triplane features for final SDF value decoding.

As illustrated in Figure 6.7, the process begins with input images, which are processed through a foreground mask to isolate the region of interest, such as a building gate. The masked images, denoted as $\mathbf{X}_i$, are then passed through a multi-scale wavelet encoder. This encoder utilizes wavelet features'

spatial and multi-scale properties to capture the geometry structure and visual details of the 2D images. The wavelet-encoded features from multiple views are fused into a unified 2D feature representation, $\mathbf{C}_i$, which is used as a conditioning input to a 2D UNet architecture. The fused features are then passed through an MLP-based decoder to predict 3D signed distance values.

This approach offers several advantages. Wavelet-encoded 2D features effectively capture geometry and visual information, enabling the 2D UNet to generate more accurate 3D reconstructions. By aligning wavelet features across multiple views and using them as a conditioning input, the model leverages complementary perspectives to produce high-quality results. Compared to the codebook-based approach, this pipeline demonstrates improved generalization, as wavelet features adapt more effectively to diverse visual patterns beyond specific building structures. Moreover, the 2D UNet architecture provides robust and flexible latent feature fusion, better accommodating variations in the input images.

In summary, this wavelet encoder-based 2D-3D reconstruction framework integrates the strengths of wavelet-encoded features and deep learning-based 3D reconstruction. This results in a more versatile and generalizable system capable of handling a wide range of 3D shapes and scenes.

The proposed reconstruction method leverages implicit triplane features $\mathbf{F}_{xy}, \mathbf{F}_{xz}, \mathbf{F}_{yz}$, while there are learned 2D feature grids aligned with three orthogonal planes to encode both geometric and appearance information of a 3D scene. As shown in Figure 6.7, our framework utilizes these triplane representations for efficient 3D reconstruction from images with known poses. For any 3D query point along a sampled ray, features are retrieved from the three orthogonal planes and aggregated to predict the Signed Distance Function (SDF) value at that location.

To enhance this representation, we introduce a pipeline that enriches triplane features with wavelet-encoded geometric details extracted from input images $\mathbf{X}_i$. These input images from the subset undergo monocular depth estimation and multiresolution wavelet transforms before being aggregated and fused into refined triplane characteristics $\{\mathbf{F}_{xy}^{fused}, \mathbf{F}_{xz}^{fused}, \mathbf{F}_{yz}^{fused}\}$ for high-quality SDF prediction. In essence, our method improves surface reconstruction by integrating implicit triplane features with multi-scale wavelet features. The following sections detail each stage of the method along with its mathematical formulation.

### 6.3.3 Preliminaries

**Implicit Neural Rendering.** Implicit NeRF encodes a 3D scene by representing its volume density and color field, leveraging multi-view posed

Figure 6.7: Our model is based on implicit triplane feature fusion for Signed Distance Function (SDF) prediction. Given an input view image $\mathbf{X}_i$, a foreground mask $\mathbf{X}_i^{fg}$ is extracted to focus on the target region for SDF queries. For each pixel, its ray is traced from the camera view $\mathbf{C}_i$ to query the implicit triplane features $\{\mathbf{F}_{xy}, \mathbf{F}_{xz}, \mathbf{F}_{yz}\}$. Images with close-up details are processed via a monocular depth before predicting depth maps, followed by wavelet transforms in three resolutions. The transformed features $\mathbf{W}_*$ are encoded through a multi-scale wavelet feature encoder ($\Phi_1, \Phi_2, \Phi_3$) and aggregated into a fused wavelet feature map. This map is projected onto three orthogonal planes, producing $\mathbf{Z}_{proj} = \{\mathbf{Z}_{xy}, \mathbf{Z}_{xz}, \mathbf{Z}_{yz}\}$. The triplane features [185] and wavelet features are concatenated and further fused using a 2D U-Net $\psi$. Finally, MLPs $\mathbf{g}(\cdot)$ decode the fused features to predict SDF values. During inference, the isosurface is extracted via marching cubes to generate the mesh.

images through volume rendering. A pixel ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is defined, starting from the camera position $\mathbf{o} \in \mathbb{R}^3$ and traversing along the view direction $\mathbf{d} \in \mathbb{R}^3$. Radiance integration along the ray accumulates color contributions from the sampling points of each ray to generate the final pixel color. For each sampling point, volume density $\sigma$ and radiance $\mathbf{c}$ are predicted using separate MLPs. The rendered pixel color $\hat{\mathbf{C}}$ is calculated by $\mathbf{T}(t) = \exp(-\int_{t_n}^{t} \sigma(\mathbf{r}(u))\, du)$, density $\sigma(t)$, and color $\mathbf{c}(t)$ over the ray, bounded by $t_n$ (near) and $t_f$ (far):

$$\hat{\mathbf{C}} = \int_{t_n}^{t_f} \mathbf{T}(t)\sigma(\mathbf{r}(t))\mathbf{c}(t)\,dt. \tag{6.1}$$

For practical computation, the numerical quadrature-based integration [4] is used to approximate continuous integral calculation.

**SDF-Based Neural Implicit Surface.** A 3D surface $\mathcal{S}$ can be implicitly represented using the zero-level-set of its signed distance function $f(\mathbf{x})$ : $\mathbb{R}^3 \to \mathbb{R}$, with a 3D point initialized from the depth map of color image $\mathbf{X}$ as input. Here, $\mathcal{S} = \{\mathbf{x} \in \mathbb{R}^3 \mid f(\mathbf{x}) = 0\}$, can be seen as the zero-crossing of the signed distance function. NeuS [247] reformulates volume density rendering in NeRF into a signed distance field (SDF) representation by employing a logistic function to optimize for neural volume rendering,

$$\sigma(\mathbf{x}) = \phi_s(f(\mathbf{x})), \tag{6.2}$$

where $\phi_s(x) = se^{-sx}/(1 + e^{-sx})^2$ is a logistic density function. It can be derived as the derivative of the sigmoid function $\Phi_s(x) = (1 + e^{-sx})^{-1}$, and is parameterized by the slope $s$. The final opaque density $\sigma(t)$ along the ray is thus given by:

$$\sigma(t) = \max\left(-\frac{d\Phi_s}{dt}(f(\mathbf{r}(t)))/\Phi_s(f(\mathbf{r}(t))), 0\right). \tag{6.3}$$

## 6.3.4   Preprocessing

To get rid of clutter pixels like humans and animals existing in the random online images of landmarks collected in the wild, we further utilize a preprocessing pipeline to create cleaner images and masks for training high-fidelity reconstruction meshes.

The preprocessing pipeline including distractor detection, distractor mask, background mask, effectively filters out non-architectural elements to focus the training only on the relevant structural components. The end result provides clean input data where query rays are only generated for the actual building geometry, improving the quality of the learned implicit representation.

Furthermore, our model can also be directly used to render a clean color image by introducing a color rendering head. The images in Figure 6.9 show a comparison of removing unwanted pixels (like people) from a photo of the Brandenburg Gate in Berlin. The input to our framework is a color image (a), which contains pedestrians in front of the Berlin Gate. This image serves as the initial scene for further cleanup. In the next step, distractors

Figure 6.8: A three-stage preprocessing pipeline for distractor removal: (a) Initial detection identifies unwanted elements like people and objects in the foreground using object detection, (b) Segment Anything Model (SAM) created [122] converts these detections into precise segmentation masks shown in black silhouettes, and (c) The final masked result isolates the architectural structure by removing both the detected distractors and sky background, leaving only the foreground building pixels that will be used for training the implicit model given query rays.



Figure 6.9: (a) Raw image with distractor on the ground. (b) Inpainted image without distractor as training input. (c) Rendered color image predicted by conditioning on trained implicit SDF model. The whole distractor removal process on the raw image is followed by the diffusion model. In the end, our implicit model after training can render the full image without the distractor pixels.

are detected and removed, followed by an inpainting process to recover the missing pixels, resulting in the processed image (b). Finally, (c) presents the rendered output generated by the pre-trained implicit model, demonstrating the scene reconstruction without distractors and validating the effectiveness of our approach. Such color rendering result is implemented by introducing an additional head for color prediction conditioned on the SDF value prediction of the original implicit 3D model to justify the clean 3D representation of the implicit SDF model. This paper is still focused on the results of the

3D reconstruction instead of the results of the rendering.

Wavelet transforms are applied selectively to high-quality close-up images to optimize training efficiency for wavelet feature fusion.

The overall pipeline effectively removes the transient elements (people) while preserving and reconstructing the underlying static architecture through a combination of detection, masking, and inpainting for a cleaner 3D reconstruction.

### 6.3.5 Model Structure

We adopt the same implicit volumetric rendering expression as clarified in the previous section for the following model introduction. The whole model is composed of five parts, including a multi-scale wavelet feature encoder, a triplane feature query, a wavelet encoder with output feature projection onto a triplane, a triplane feature fusion, and an implicit SDF decoder. In particular, the input of the wavelet encoder is monocular depth, while all multi-view color images are used as input to the triplane feature encoder.

**Wavelet encoder for multi-scale features.** Given a selected input image $\mathbf{X}_i \in \mathbb{R}^{H \times W \times 3}$ from the original multiview images, the monocular depth before selected images predicts a depth map $\mathbf{D}_i \in \mathbb{R}^{H \times W}$. The selection of particular close-up images is based on the quality and details that exist in the input view image, and most views are quite distant with blurry pixels, thus making it hard for the monocular depth estimation to provide accurate depth details. The depth map $\mathbf{D}_i$ undergoes a wavelet transform in three resolutions to produce multi-scale wavelet features $\{\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3\}$. These are generated using a wavelet transform in three resolutions. These features are then processed by the wavelet encoder $\Phi$ with three different sizes, resulting in a final fused wavelet feature map $\mathbf{Z}_{wave} \in \mathbb{R}^{H' \times W' \times C}$ through feature aggregation:

$$\mathbf{Z}_{wave} = \Phi_1(\mathbf{W}_1) + \Phi_2(\mathbf{W}_2) + \Phi_3(\mathbf{W}_3), \tag{6.4}$$

where $\Phi_{1,2,3}$ are scale-specific encoding functions to extract various sized features. Usually, $C$ is four channels, representing 2D signals through four filters, defined as $\mathbf{LL}$, $\mathbf{LH}$, $\mathbf{HL}$, and $\mathbf{HH}$. Given an input image $\mathbf{X}$, the 2D wavelet transform with specific scale decomposes the image into a low-frequency component $\mathbf{x_L}$ and three high-frequency components $\{\mathbf{x_H}, \mathbf{x_V}, \mathbf{x_D}\}$, corresponding to horizontal, vertical, and diagonal details respectively.

We train the wavelet feature encoder using an autoencoder similar to the design of LiteVAE ([205]), aiming to reconstruct the original depth map from its wavelet-transformed representation. We apply a single-level wavelet

decomposition independently at multiple scales of the input depth map, generating multi-scale wavelet feature maps. This allows the encoder to capture fine-to-coarse spatial details efficiently. After pretraining three separate wavelet encoders, we obtain their extracted multi-scale feature representations. To ensure alignment, We downsample the feature maps from the two higher-resolution encoders by factors of 1/2 and 1/4, respectively, to align with the smallest-scale feature map. This downsampling and aggregation follow the same process as LiteVAE [205]. To mitigate the loss of fine details, we retain multi-scale information by aggregating features across different resolutions. Furthermore, since each wavelet decomposition produces four sub-bands (LL, LH, HL, HH), we stack these sub-maps along the channel dimension before passing them to the subsequent processing pipeline.



Figure 6.10: Wavelet transform of the depth map in finest resolution, (a) is the original depth map, and (b) is the wavelet transform of the depth map, composed of four parts: Low-Low (LL), Low-High (LH), High-Low (HL), and High-High (HH). The wavelet transformed depth is used as input for the AutoVAE Encoder.

We provide example results of wavelet transformed features in Figure 6.10, which demonstrates the effectiveness of wavelet transforms in preserving geometric information from depth maps. The visualization compares original depth maps $\mathbf{D}_j$ (top row) with their corresponding wavelet decompositions $\mathbf{W}_j$ (bottom row). The input depth maps, predicted by the state-of-the-art diffusion-based monocular depth estimation network of LOTUS [80], capture detailed geometric structures and continuous depth variations. Our wavelet transform decomposes these depth maps in three resolutions into multi-scale

feature representations $\mathbf{W}_j, j = 1, 2, 3$ with three levels, where each level $j$ preserves both spatial and frequency information critical for geometric detail reconstruction. This multi-resolution representation enables the model to effectively encode both fine-grained surface details and global shape features. The wavelet transform decomposes each depth map with a specific resolution into four sub-bands (LL, LH, HL, HH), effectively capturing different frequency components. While LL retains a global structure, LH and HL emphasize horizontal and vertical details, and HH captures diagonal features. This highlights the ability to preserve and distinguish depth-specific geometry, and such spatial can also be easily aligned with the image feature map.

**Pixel ray query for implicit triplane features.** Each pixel of the foreground masked input image $\mathbf{X}_i^{fg}$ is associated with a ray cast from the camera view $\mathbf{o} \in \mathbb{SE}(3)$. All the sampled points along query rays of each image retrieve implicit triplane features $\{\mathbf{F}_{xy}, \mathbf{F}_{xz}, \mathbf{F}_{yz}\}$ from three orthogonal planes $\{xy, xz, yz\}$ of the 3D space via ray projection, where each plane has a feature resolution of $\mathbb{R}^{H' \times W' \times C'}$, with feature channel dimension $C'$:

$$\{\mathbf{F}_{xy}, \mathbf{F}_{xz}, \mathbf{F}_{yz}\} = \text{Query}_{\{xy, xz, yz\}}(\mathbf{r}(t)). \tag{6.5}$$

**Wavelet feature projection onto triplane.** Meanwhile, the wavelet feature map $\mathbf{Z}_{wave}$ of Equation 6.4 is projected onto the three orthogonal 2D planes to match the implicit triplane feature resolution. This cosine projection generates three projected wavelet feature maps $\{\mathbf{Z}_{xy}, \mathbf{Z}_{xz}, \mathbf{Z}_{yz}\}$ respectively.

To incorporate wavelet features into the implicit Signed Distance Field (SDF) model, we first generate a structured 3D representation of the scene by leveraging aligned depth maps. Specifically, we reconstruct a dense unprojected point cloud in the camera frame followed by a camera-to-world transform. This transformation involves back-projecting depth pixels into 3D space using the known intrinsic and extrinsic camera parameters. The resulting 3D points are then associated with wavelet-based features aligned with 2D pixels.

Once the wavelet-enhanced feature map is obtained, it is projected onto the three feature-aligned triplane representations corresponding to the orthogonal planes defined by the normal vectors $(1, 0, 0), (0, 1, 0), (0, 0, 1)$. This projection ensures that the 3D wavelet features are properly integrated into the implicit triplane feature space. Mathematically, this process is formulated as follows:

$$\{\mathbf{Z}_{xy}, \mathbf{Z}_{xz}, \mathbf{Z}_{yz}\} = \mathbf{P}\mathbf{Z}_{wave} \cdot \cos(\{\alpha, \beta, \gamma\}), \tag{6.6}$$

Figure 6.11: Sampling points along the pixel ray $\mathbf{r}(t)$ starting from $\mathbf{o}$ for implicit triplane feature learning via projection. For Wavelet feature projection onto triplane. The ray $\mathbf{r}(t)$ starts from the camera origin $\mathbf{o}$, then passes through a single unprojected point $\mathbf{S}$. Dashed lines represent the orthogonal projections onto the $xy$, $xz$, and $yz$ planes to obtain triplane features $s_{xy}$, $s_{xz}$, and $s_{yz}$ for a 3D point.

Where $\{\mathbf{Z}_{xy}, \mathbf{Z}_{xz}, \mathbf{Z}_{yz}\}$ represent the projected wavelet-enhanced features on the three orthogonal feature planes $xy, xz, yz$. The projection angles $\alpha, \beta, \gamma$ correspond to each feature plane, ensuring an optimal alignment between the wavelet features and the triplane encoding. Such an angle is the dot product between the ray direction and the axis direction. The transformation matrix $\mathbf{P}$, derived from the camera intrinsic parameters and the camera-to-world extrinsic pose, maps unprojected 3D points from the camera frame to the world coordinate system. These mapped points are then projected onto the triplane feature planes, where they serve as inputs to our method. Each pixel ray maps 3D spatial information onto a set of triplane feature representations. Given a camera pixel ray $\mathbf{r}(t)$, we analyze the sampled point $\mathbf{S}$ along the ray scaled by the predicted depth value and compute its orthogonal projections onto the three principal planes: $xy$, $xz$, and $yz$. These projections provide the corresponding triplane feature locations $\mathbf{s}_{xy}$, $\mathbf{s}_{xz}$, and $\mathbf{s}_{yz}$.

Figure 6.11 illustrates the feature extraction process along a pixel ray $\mathbf{r}(t)$. The ray originates from the camera at $\mathbf{o}$, extends through the sampled point $\mathbf{S}$, and continues along its trajectory. Dashed lines indicate the orthogonal

projections of $\mathbf{S}$ onto the three triplane feature planes ($xy$, $xz$, and $yz$), which are used for feature representation.

In Figure 6.11, implicit triplane features are obtained by sampling multiple points along the pixel ray uniformly, capturing continuous spatial information. Wavelet triplane features are projected in the same way as the implicit features. A key distinction is that each feature plane in the implicit approach consists of 16 channels, whereas the wavelet-based plane features are compressed into 4 channels, reducing redundancy while preserving essential spatial details.

This structured feature projection enables seamless integration of 2D wavelet-transformed depth features with 3D implicit features, leading to more accurate SDF predictions and higher-fidelity 3D reconstructions.

**Feature concatenation and fusion.** The implicit triplane features and the projected wavelet features of each feature plane are concatenated along the channel dimension and fused using a 2D U-Net. This results in three fused triplane features $\{\mathbf{F}_{xy}^{fused}, \mathbf{F}_{xz}^{fused}, \mathbf{F}_{yz}^{fused}\}$, where $\mathbf{F}_{*}^{fused} \in \mathbb{R}^{H' \times W' \times 2C}$:

$$\{\mathbf{F}_{xy}^{fused}, \mathbf{F}_{xz}^{fused}, \mathbf{F}_{yz}^{fused}\} = \{\psi_{\text{fusion}}([\mathbf{F}_{xy}; \mathbf{Z}_{xy}]),$$
$$\psi_{\text{fusion}}([\mathbf{F}_{xz}; \mathbf{Z}_{xz}]), \psi_{\text{fusion}}([\mathbf{F}_{yz}; \mathbf{Z}_{yz}])\}, \tag{6.7}$$

where $[;]$ denotes the concatenation of feature maps along the channel dimension.

**Encoder for SDF prediction.** The fused features are finally decoded by a neural network $g(\cdot)$, which predicts the SDF value $v \in \mathbb{R}$ for the given pixel ray query:

$$v = g(\mathbf{F}_{xy}^{fused}, \mathbf{F}_{xz}^{fused}, \mathbf{F}_{yz}^{fused}). \tag{6.8}$$

The predicted SDF values are used to extract the isosurface via marching cubes [148], producing a reconstructed 3D mesh.

**Loss function.** The total training loss $\mathcal{L}_{\text{total}}$ for the implicit model is defined as a combination of three components: the mean cross-entropy loss $\mathcal{L}_{\text{RGB}}$, the Eikonal regularizer $\mathcal{L}_{\text{Eik}}$, The total loss is defined as:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{RGB}} + \lambda_{\text{Eik}} \mathcal{L}_{\text{Eik}}, \tag{6.9}$$

The photometric loss is defined as the $L_2$ loss between rendered image $\hat{C}(\mathbf{r}_i)$ and ground truth image $C(\mathbf{r}_i)$,

$$\mathcal{L}_{\text{rgb}} = \frac{1}{N} \sum_{i=1}^{N} \|\hat{C}(\mathbf{r}_i) - C(\mathbf{r}_i)\|_2^2, \tag{6.10}$$

where the Eikon loss is applied to $M$ neighboring sampled points to regularize the smoothness of local SDF prediction.

$$\mathcal{L}_{\text{Eik}} = \frac{1}{M} \sum_{i=1}^{M} \left| \|\nabla \hat{v}_i\| - 1 \right|^2, \tag{6.11}$$

In the formulation 6.9, $\mathcal{L}_{\text{Eik}}$ regularizes the gradients to enforce the local smoothness of the signed distance field. For color image rendering, we need to add another mean-squared photometric loss to leverage the supervision of RGB pixels.

# 6.4 Experiments and Results

## 6.4.1 Point Cloud-based 3D Reconstruction

At first, we present 3D point cloud completion results using a 3D U-Net backbone with partial shape point cloud inputs. The network features a CNN-based encoder and a decoder with transposed convolutional layers that mirror the encoder structure.

As shown in Figure 6.12, the model effectively completes shapes with symmetrical input, such as the chair in the top-left subfigure. However, it struggles with complex, asymmetrical structures, as indicated by the red circles. Thin elements, like chair legs, are not reconstructed accurately, highlighting the UNet model's limitations in capturing fine geometric details.

Next, we present mesh reconstruction results using a 3D CNN encoder paired with an implicit MLP decoder, as illustrated in Figure 6.6. In the last column of Figure 6.13, integrating occupancy probability and SDF loss during training enables the model to reconstruct fine details and generate smooth meshes from coarse voxels. While occupancy probability loss alone produces reasonable shapes, it fails to fully capture surface details, leading to gaps and holes, such as those in the back of the chair. This underscores the importance of combining both loss functions to enhance surface reconstruction quality.

## 6.4.2 Image-based 3D Reconstruction

**Datasets.** To evaluate the general performance of our 3D reconstruction approach, we make use of a wide variety of datasets, including the DTU ([103]) dataset, which is collected from a turntable; the Tanks and Temples dataset ([123]), captured as video scans of sculptures and buildings; and the Cultural Heritage dataset ([157]), which features large-scale historical sites.

Figure 6.12: The results of 3D UNet model for point cloud completion, where blue is the partial input, green is the GT point cloud, and red is the predicting output of 3D UNet model.

|       |     |          |           |
|-------|-----|----------|-----------|
| Input | GT  | Occ only | Occ + SDF |

Figure 6.13: The shape reconstruction results from 3D voxel input via implicit SDF model.

111

For DTU ([103]), we use the Chamfer distance metric calculated between the reconstructed model and the ground truth. In contrast, for Tanks and Temples ([123]), we evaluate reconstruction accuracy using the F1 score ($F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$), as Chamfer distance makes it difficult to differentiate the performance for some scenes. Due to the large scene scale of the Cultural Heritage dataset, obtaining ground truth (GT) meshes or pseudo-GT is challenging, so we primarily provide qualitative results. For DTU and Tanks and Temples, we uniformly sample 1,000 and 10,000 points, respectively, and compare them with the nearest 3D points from the GT mesh, which is obtained through the COMLAP tool, followed by delicate post-processing to make the GT mesh complete and smooth enough without artifacts like holes.

**Baseline models.** For baseline evaluation, we compare our method against several state-of-the-art implicit and explicit 3D reconstruction models. The implicit SDF baselines include VolSDF by [270], NeuS by [247], Neuralangelo by [134], and BakedSDF by [271]. The explicit reconstruction baselines include SuGaR by [73], GOF by [94], and 2DGS by [93]. We provide quantitative comparisons across the DTU and Tanks and Temples datasets, while qualitative visual comparisons highlight the top three performing models. The input for all baselines is images and camera poses.

**Implementation details.** For monocular depth estimation, we use the diffusion-based LOTUS model by [80] to predict depth maps for selected views of the heritage dataset. In contrast, for the Tanks and Temples, and DTU datasets, we process all training images. Wavelet decomposition is performed using the Fast Wavelet Transform (FWT) ([155]) with Haar basis filters. Please note that depth is used only for training the wavelet features and not for supervising the 3D reconstruction. Furthermore, for the cultural heritage dataset, all input images undergo the same preprocessing step to remove transient pedestrians, ensuring that all baseline models are trained with the same masked image inputs. This guarantees fairness in the baseline comparisons.

Our reconstruction pipeline is trained per scene, consistent with standard neural implicit surface reconstruction settings. Regarding the wavelet feature encoder, we emphasize that it is not trained in a scene-specific way. Rather, this module is pre-trained offline using an autoencoding objective (inspired by LiteVAE) on a collection of monocular depth predictions across a variety of generic scenes. These depth maps, along with their wavelet-transformed representations, are used to learn a compact latent prior over geometric features, analogous to learned image priors used in low-level vision tasks. This prior is then fixed and reused across all scenes during reconstruction. Our approach is thus fair for comparison, as it does not rely on specific geometry information supervision, and instead builds on learned geometric priors from

112

the common depth prediction foundation model prior, which are estimated from monocular RGB images in a self-supervised manner. This setup maintains the generalizability of the wavelet-depth encoder and adheres to the same per-scene optimization assumptions used in all baselines by using the generalized depth wavelet feature prior.

The autoencoder for wavelet-transformed depth features consists of a ResNet encoder followed by a fully convolutional decoder, similar to LiteAutoVAE ([205]). We apply a Gaussian blurring loss to low-frequency sub-bands and a Charbonnier loss ([14]) to high-frequency sub-bands. During implicit SDF training, the AutoVAE encoder remains frozen. The triplane feature representation is structured as $3 \times 64 \times 64 \times 16$, with an SDF decoder composed of fully connected layers. Wavelet-triplane fusion is achieved via a 2D U-Net with four downsampling and upsampling blocks, followed by a $1 \times 1$ convolution along the depth channel. The fused representation consists of three orthogonal triplane feature planes ($64 \times 64 \times 16$ each), combined with a projected wavelet feature map ($64 \times 64 \times 4$), and refined through the 2D U-Net.

The wavelet autoencoder processes four spectral channels—low-frequency, vertical high-frequency, horizontal high-frequency, and diagonal high-frequency—using ResNet blocks. Wavelet transforms are applied to Lotus-generated depth maps at three resolutions, with extracted features used to train the autoencoder. To balance fine-grained details and global features, we incorporate self-modulated convolutional layers ([205]). The loss function includes reconstruction, regularization, and adversarial terms ([205]). Features are extracted at $256 \times 256$, $128 \times 128$, and $64 \times 64$ resolutions, with higher-resolution features downsampled by 1/4 and 1/2 for alignment. For the cultural heritage dataset, we manually selected 100 close-up images to enhance the implicit SDF model with wavelet-transformed features.

For the implicit SDF model, we use the Facto-SDF implementation from SDFStudio by [277], integrating it with the triplane feature representation as the encoder backbone.

**Training Complexity.** Our training pipeline consists of two stages: training the wavelet encoder and training the implicit SDF conditioned on the frozen wavelet encoder. The wavelet encoder training takes approximately 8 hours on an RTX 3090. For the implicit SDF training of the DTU model, training is completed in 1-2 hours. As for Tanks and Temples, and Cultural Heritage), initial implicit SDF training on color images takes 6-8 hours due to data diversity, followed by 1-2 hours of fine-tuning with wavelet-triplane features.

All experiments were conducted on an NVIDIA RTX 3090 GPU, ensuring efficient training and inference. This modular approach enables scalable

learning across datasets of varying sizes and complexities.

## Baseline Comparisons

We first provide the qualitative comparison results of the sample targets or scenes in the three datasets, including the qualitative results of the DTU, Tank, and Temple, and the Cultural Heritage dataset in Figure 6.14. These data sets are collected through cameras that point towards a target. Furthermore, the quantitative results on the DTU and Tank and Temple datasets are also provided in Table 6.1 and Table 6.2, respectively.

As seen in Figure 6.14, our model can reconstruct fine-grained details on the mesh surface, such as feature details on birds, details of clothes of happy buddha, owl, texts on the Berlin gate. We recommend that readers have a close look at the red highlighted circles. The 2D Gaussian Splatting seems to struggle to preserve details and also has obvious artifacts and holes on the mesh surface. The 2D GS also fails to reconstruct a mesh of a large-scale Berlin gate. Neuralangelo is very good at preserving some details, but still has some artifacts or obstructions, as shown on the bottom of the Berlin gate with unexpected blockings, although normals are consistent along with details of texts. BakedSDF has the worst performance, smoothing the results, with a smooth surface and loss of details, particularly obvious on the Berlin gate, which may even incur some unexpected reconstruction mesh regions in front of the house.

Finally, we provide quantitative evaluation results in Table 6.1 and Table 6.2 using the Chamfer distance. On DTU, our model is the best, while Neuralangelo gets the second-best performance. On the Tank and Temple dataset, our model still scores best on three out of six scans, while Neuralangelo follows next.

Table 6.1: Chamfer Distance (CD, ↓) comparison across scan IDs. Green sates the best, Orange is second best, Yellow represents third best. Lower is better.

| Method | 24 | 37 | 40 | 55 | 63 | 65 | 69 | 83 | 97 | 105 | 106 | 110 | 114 | 118 | 122 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| VolSDF [270] | 1.14 | 1.26 | 0.81 | 0.49 | 1.25 | 0.70 | 0.72 | 1.29 | 1.18 | 0.70 | 0.66 | 1.08 | 0.42 | 0.61 | 0.55 |
| NeuS [247] | 1.00 | 1.37 | 0.93 | 0.43 | 1.10 | 0.65 | 0.57 | 1.48 | 1.09 | 0.83 | 0.52 | 1.20 | 0.35 | 0.49 | 0.54 |
| Neuralangelo [134] | 0.41 | 0.36 | 0.35 | 0.35 | 1.29 | 0.54 | 0.73 | 0.52 | 0.97 | 0.56 | 0.48 | 0.73 | 0.32 | 0.40 | 0.36 |
| BakedSDF [271] | 0.63 | 0.58 | 0.40 | 0.52 | 1.37 | 0.63 | 0.81 | 0.56 | 1.02 | 0.81 | 0.50 | 0.82 | 0.39 | 0.45 | 0.38 |
| SuGaR [73] | 1.47 | 1.33 | 1.13 | 0.61 | 2.25 | 1.71 | 1.15 | 1.63 | 1.62 | 1.07 | 0.79 | 2.45 | 0.98 | 0.88 | 0.79 |
| GOF [94] | 0.50 | 0.37 | 0.38 | 0.74 | 1.18 | 0.76 | 0.90 | 0.47 | 1.29 | 0.68 | 0.77 | 0.90 | 0.42 | 0.41 | 0.42 |
| 2DGS [93] | 0.48 | 0.39 | 0.41 | 0.83 | 1.36 | 0.83 | 1.04 | 0.70 | 1.27 | 0.76 | 0.70 | 1.40 | 0.40 | 0.43 | 0.40 |
| Ours | 0.45 | 0.34 | 0.32 | 0.34 | 0.97 | 0.52 | 0.54 | 0.50 | 0.82 | 0.53 | 0.45 | 0.68 | 0.30 | 0.36 | 0.34 |

Figure 6.14: Baseline comparison results on five targets from DTU ([103]), and Cultural Heritage ([157]) dataset. Each model result (split by a dashed line) contains mesh and normals.

Table 6.2: Chamfer Distance (CD, ↓) results on the Tanks and Temples dataset. Green is best, Orange is second best, and Yellow is third best. Lower values indicate better reconstruction accuracy.

| Method | Barn | Caterpillar | Courthouse | Ignatius | Meetingroom | Truck |
|---|---|---|---|---|---|---|
| NeuS | 0.69 | 0.68 | 0.54 | 1.62 | 0.63 | 0.90 |
| Geo-NeuS | 0.74 | 0.61 | 0.49 | 1.50 | 0.52 | 0.88 |
| Neuralangelo | 1.21 | 0.89 | 0.82 | 1.72 | 0.79 | 0.95 |
| SuGaR | 0.42 | 0.53 | 0.38 | 0.83 | 0.59 | 0.66 |
| 2DGS | 0.40 | 0.58 | 0.41 | 0.74 | 0.53 | 0.58 |
| Ours | 0.34 | 0.51 | 0.36 | 0.76 | 0.49 | 0.52 |

### 6.4.3 Ablation Study

We first perform ablation studies to assess the contribution of each key component in our model structure, with qualitative results in Figure 6.15 and quantitative metrics in Table 6.3.



Figure 6.15: Ablation study on 3D reconstruction. From left to right: (1) Removing the triplane leads to fragmented geometry. (2) Without the wavelet encoder, fine details are distorted. (3) Omitting 2D U-Net fusion results in less sharp features. (4) The full model achieves the best quality.

The wavelet feature encoder can effectively capture fine-grained details from input views, such as edge features, as shown in Figure 6.16. The output feature maps of the learned wavelet autoencoder provide a rich representation of the input views to encode more geometric details, like the carved letters on the right side of the gate. The final reconstruction mesh details can be enhanced by the decoder conditioning on the wavelet features. Figure 6.16 showcases the feature output of the largest wavelet encoder in an autoencoder

pre-trained separately on the wavelet-transformed depth input in three resolutions, highlighting the progressive decomposition of image features across multiple frequency bands. The wavelet encoder outputs are visualized across four columns in Figure 6.16. Column (a) shows the encoded depth features, preserving the overall geometric structure. Column (b) displays vertical gradient features $\Phi$ that highlight edge transitions along the y-axis. Column (c) presents horizontal gradient features, capturing edge variations along the x-axis. Column (d) shows diagonal gradient features that encode diagonal directional geometric variations. This learned decomposition through our wavelet encoder enables comprehensive feature extraction at multiple orientations, crucial for accurate 3D surface reconstruction. Each component contributes specific directional information, allowing the model to capture both directional surface variations and overall geometric structure.

Table 6.3: Ablation study of our model on DTU [103], reporting Chamfer Distance (CD, $\downarrow$) to evaluate geometric accuracy. Lower values indicate better reconstruction quality.

| Metric (CD$\downarrow$) | w/o Triplane Feature | w/o Multi-scale Wavelet Feature | w/ Single Scale Wavelet | w/o UNet Channelwise Fusion |
|---|---|---|---|---|
| CD $\downarrow$ | 0.87 | 0.65 | 0.60 | 0.56 |
| Metric (CD$\downarrow$) | w/o Wavelet Autoencoder | w/o Photometric Loss | w/o Eikonal Loss | Full Model |
| CD $\downarrow$ | 0.60 | 0.94 | 0.55 | **0.51** |

First, we evaluate the impact of removing the triplane representation, reverting to a purely MLP-based implicit SDF modeling. This results in severe surface fragmentation, geometric noise, and disconnected structures, reflected in a significantly worse Chamfer Distance (CD) of 0.87, indicating the critical role of the triplane representation in preserving global spatial coherence. Removing the multi-scale wavelet feature while retaining the triplane and decoder degrades surface quality, especially around fine structures (e.g., object edges and thin regions), increasing the CD to 0.65. Incorporating only a single-scale wavelet feature improves reconstruction to 0.60, whereas our full multi-scale wavelet pipeline further reduces CD to 0.51, confirming the importance of multi-resolution geometric encoding.

In addition, removing the 2D U-Net channel-wise fusion module and instead directly concatenating features leads to slight oversmoothing and loss of fine detail, reflected in a CD of 0.56, underscoring the benefit of learned feature fusion. On the loss design side, excluding the photometric loss produces the most degraded result (0.94 CD), demonstrating its necessity for accurate geometry recovery. Removing the wavelet autoencoder or the Eikonal regularization also leads to increased CD scores of 0.60 and 0.55, respectively, showing their contribution to feature compactness and surface smoothness.

Overall, the complete model, which combines triplane encoding, multi-scale wavelet features, U-Net fusion, and all loss terms, achieves the best reconstruction accuracy with a CD of 0.51, effectively balancing global consistency with local surface precision across the DTU dataset.



Figure 6.16: Visualization of learned wavelet encoder feature maps at the highest resolution level. The four columns demonstrate different components of the encoded representation: (a) depth features preserving the overall geometric structure, (b) vertical gradient features capturing y-axis surface variations, (c) horizontal gradient features encoding x-axis transitions, and (d) diagonal gradient features representing cross-directional geometric patterns. Each component is processed through our wavelet encoder $\Phi$ to extract orientation-specific geometric information.

To further validate the design of our depth-guided wavelet feature extractor, we perform an extended ablation study isolating the roles of depth input, wavelet transform, and the autoencoder. Removing the entire module leads to a CD of 0.60, while using only depth maps without wavelet or autoencoder yields a poorer result of 0.72, indicating insufficient geometric encoding from

Table 6.4: Extended ablation study on the DTU dataset [103]. This table evaluates the contribution of depth, wavelet, and autoencoder components. Lower Chamfer Distance (CD) indicates better reconstruction.

| Ablation Setting | Chamfer Distance (CD ↓) |
|---|---|
| w/o Depth Wavelet Autoencoder module | 0.60 |
| w/ Depth only (no Wavelet, no Autoencoder) | 0.72 |
| w/ Depth + Wavelet (no Autoencoder) | 0.58 |
| w/ Depth + Autoencoder (no Wavelet) | 0.65 |
| w/ Color + Wavelet + Autoencoder | 0.84 |
| w/ Full wavelet pipeline (Depth + Wavelet + Autoencoder) | **0.51** |

depth alone. Incorporating depth and wavelet features without the autoencoder improves performance to 0.58, showing that multi-scale spatial cues from wavelet features significantly enhance reconstruction. However, using depth with an autoencoder but without wavelet results in a higher CD of 0.65, highlighting the essential role of wavelet transforms in the encoding structure beyond the raw depth. Notably, relying on color wavelet features alone, even with the autoencoder, produces the weakest result (0.84 CD), confirming that geometric cues from depth are critical. Our complete wavelet pipeline, which combines depth, wavelet transform, and auto-encoder, achieves the best reconstruction with a CD of 0.51, demonstrating the complementary benefits of multiscale frequency representation and learned feature compression.

We further present visualization results from the ablation study on the wavelet autoencoder module design.

As shown in 6.17, the raw depth contains noticeable local noise. Removing either the wavelet transform or the autoencoder results in feature maps that remain noisy and slightly blurred. In contrast, the second row, second column—representing the full wavelet-autoencoder pipeline—produces clean, fine-grained, and geometrically detailed feature maps. On the other hand, using color images for wavelet-based autoencoding yields a completely different feature distribution, failing to capture meaningful geometric information, as seen in the bottom-right corner of the color input example.

## 6.5 Conclusion

We propose an implicit SDF model that integrates wavelet-transformed depth features into a latent triplane feature space. By combining spatially decomposed wavelet representations with triplane embeddings, our approach enhances the preservation of geometric details. During inference, fused features are sampled along query rays and decoded into SDF values, enabling high-fidelity mesh reconstruction. Our model requires only monocular priors

Figure 6.17: Ablation study visualizations of different design combinations for the wavelet-transformed depth and autoencoder module.

from state-of-the-art diffusion-based depth estimation models or a subset of selected heritage dataset images. Compared to existing implicit SDF and explicit Gaussian Splatting methods, our approach achieves superior shape completeness while retaining intricate geometric details. Despite these advances, opportunities remain for further improvement. Future work could explore optimized sampling strategies to enhance computational efficiency. Additionally, integrating discrete Gaussian representations may accelerate training while maintaining high reconstruction fidelity. These extensions could expand our method's applicability to large-scale scenarios and real-time applications.

A key assumption in our work is the photometric consistency of features across multi-view color images, which generally holds for Lambertian surfaces where appearance remains constant across viewpoints. However, this assumption breaks down for reflective or transparent materials due to non-Lambertian effects. To reconstruct such objects, alternative approaches are needed that explicitly model or mitigate view-dependent appearance changes. For example, incorporating uncertainty prediction to filter out reflective regions or leveraging shape priors to guide geometry inference—rather than relying solely on pixel feature matching and classical structure-from-motion—can improve robustness. While Gaussian Splatting has recently emerged as a promising prior art for handling appearance and

view-dependent effects, our focus in this chapter has been on geometry rather than appearance. We therefore leave the exploration of such appearance-oriented methods and their potential integration with geometric representations as an important direction for future work.

# Chapter 7

# Conclusion

This dissertation presents a set of 3D geometric deep learning frameworks that integrate geometric constraints with deep learning models to address key computer vision tasks. Specifically, we explore geometric deep learning for camera pose estimation, point cloud registration, focal stack depth estimation, and implicit SDF-based 3D reconstruction. These techniques enhance mobile digital applications, providing robust, scalable, and high-quality solutions for the broader 3D vision pipeline. Beyond the academic contribution of making 3D learning representation more robust, accurate, and efficient, these models significantly affect Virtual Reality (VR), Augmented Reality (AR), and digital twin generation. Such high-fidelity 3D reconstructions and 3D asset generation enable virtual museums, digital twins, and interactive learning for a better education experience, creating immersive experiences that connect virtual with modern reality. With companies like Ubisoft, the gaming industry has already integrated reconstructed historical heritage into the entertainment industry, demonstrating great commercial value.

This dissertation presents a clear roadmap for advancing 3D vision by integrating geometric priors and constraints (such as surfels, manifolds, and wavelet features) to condition and guide deep learning models. Many of these hybrid strategies and combinations can be adapted to solve other challenges in 3D vision. In the end, this thesis demonstrates how hybrid approaches can introduce physics-based methods to data-driven learning models, which can boost the learning model's performance. Regarding evaluation metrics: in Chapter 3, the real-time camera pose estimation system uses algorithms involving random sampling. As a result, repeated runs on the same video sequence yield different outputs, so we performed multiple runs and reported the mean and standard deviation in the evaluation. In contrast, Chapters 4 and 5 focus on point cloud registration and depth estimation using deep learning models with fixed weights and random seeds in an offline setting,

where the same input consistently produces the same output, making statistical analysis unnecessary. For 3D reconstruction, we performed multiple runs to extract meshes and selected the visually best quality mesh for metric evaluation.

## 7.1  Camera pose Estimation

For pose estimation, our manifold-based adaptive particle sampling method taps into natural geometric primitives, such as skylines and ground planes, as reliable reference points alongside IMU signals. This approach delivers outstanding accuracy and keeps drift impressively low over long-term periods, outperforming traditional camera pose estimation techniques. By providing the method with these intuitive geometric cues, it handles the correspondence uncertainty of appearance features in real-world natural settings, proving its superiority over traditional pose-tracking methods that often struggle with natural environments. Additionally, as it is quite challenging to run these models on mobile hardware with constrained computing resources and a limited power supply, we cannot run some heavy deep learning models on the hardware. In such cases, a more efficient algorithm should be preferred, or, as in our case, we directly use deep learning only for image segmentation, while the subsequent feature tracking, pose fusion, and pose estimation are non-learning-based. We implemented it in real-time on a 3D-printed gimbal platform mounted onto the polar stick to conduct the simulation in the lab with the camera pointing toward a landscape picture. Next, we put the whole system onto a UAV flying through natural outdoor environments to validate the performance in practice. The results confirm its practical use in challenging wild settings, where precision and reliability are very satisfying. This not only shows the robustness of the method but also inspires us to borrow similar ideas by leveraging natural cues or other permanent environmental features for vision tasks. Moreover, such manifold-based particle filters can better approximate the Lie group-based rotation representation and improve the pose estimation accuracy and efficiency of deep learning or partial learning models. This technique is important for aerial mapping or autonomous navigation, where steady pose estimation in natural environments is a foundational ability for such systems.

## 7.2 Point Cloud Registration

In point cloud registration, we introduced a surfel-based **SE**(3)-equivariant network that achieves state-of-the-art accuracy on both indoor and outdoor benchmark datasets. This approach leverages a 2D surfel representation, incorporating surfel initialization from raw RGB-D depth maps or LiDAR point clouds, and consists of a shared E2PN encoder, a cross-attention module, and an MLP-based decoder. Extensive experiments across two datasets highlight its robustness and plausible accuracy. By explicitly embedding **SE**(3)-equivariance into the framework, it effectively handles the geometric ambiguity in 3D data, proving its superiority over point cloud registration techniques that often fail with noisy data and input points with small overlaps. The modular surfel design also enables strong generalization across varied 3D environments, laying a solid foundation for consistent 3D global mapping.

This method not only showcases good performance and robustness but also opens doors for broader real applications, like 3D mapping for robotics, 3D reconstruction, and even for downstream tasks like visual investigation and visual analysis in robotics or infrastructure engineering. It inspires fresh ideas too, encouraging us to explore these 2D Gaussian surfel primitives further for tackling sparse or noisy point clouds with small overlaps, including those with dynamic objects or uncertain data. Furthermore, refining correspondence and feature learning representation could boost the efficiency and performance of registration, especially under unstructured settings. Such equivariance constraints also improve the learning efficiency of geometric feature representation. This technique matters for applications like autonomous navigation or large-scale scene modeling, where reliable point cloud registration in noisy, real-world conditions is mandatory for next-level systems.

## 7.3 Depth from Focal Stack Images

We came up with a new model to estimate depth from focal stacks that overcomes the limitations of previous CNN-based methods. It uses a Transformer encoder to capture 2D spatial features and an LSTM module to learn depth cues across different stack images. Before the model, the multi-scale convolutional encoding further enhances detailed feature extraction, ensuring accurate depth prediction. Then we further use the lens distance constraint to create the loss as supervision for focal stack depth estimation learning.

Through extensive evaluations, FocDepthFormer demonstrated state-of-the-art performance, outperforming prior methods across multiple bench-

marks. More importantly, its ability to handle arbitrarily sized focal stacks offers practical advantages in real-world applications, where fixed stack sizes are often not flexible for training and testing. The pre-training strategy on monocular depth estimation datasets also proved effective in mitigating data scarcity, reinforcing the model's adaptability.

Lastly, this work highlights the potential of hybrid deep learning architectures for depth estimation by merging attention and recurrence for latent feature fusion, along with focal stack constraints to differentiate focus and defocus cues better for depth feature learning. FocDepthFormer contributes to advancing 3D reconstruction techniques and lays the foundation for future research in computational photography and depth prediction for dense local viewpoint cloud generation.

## 7.4    3D reconstruction

Our proposed implicit SDF learning framework incorporates wavelet-transformed depth features into a triplane representation through triplane projection, achieving high-fidelity 3D reconstructions from limited image views. Such depth input for wavelet feature learning is generated by the latest monocular depth prior model, which is based on the diffusion model. By integrating spatial decomposition with neural representations, the model can help preserve intricate geometric details while maintaining global shape completeness.

In a nutshell, these contributions advance the progress of 3D reconstruction by demonstrating how such geometric constraints in high-frequency bands can be effectively integrated into deep learning frameworks. While challenges remain, including computational efficiency, extreme environmental conditions, sparse input view, and scalability to complex and large-scale scenes. This wavelet feature prior can be used as a plug-and-play module for other model backbones to boost the overall performance of the reconstruction. This research establishes a foundation for more robust, accurate, and geometrically consistent reconstruction results suitable for real-world VR/AR and robotic applications. It also inspires us to look back into the traditional spectral technique to combine the spectral features in high frequencies for reconstruction.

## 7.5    Future Work

While this dissertation advances 3D vision across multiple tasks, several key areas remain open for exploration. Future research can build upon these

findings to enhance accuracy, efficiency, and adaptability in real-world applications.

**Robust Multi-Modal Fusion for Pose Estimation.** Current pose estimation models rely on geometric priors and learning features, but robustness in more challenging environments with dynamics remains an issue. Future work could integrate feature tracking with lightweight neural implicit representations to improve performance under low illumination, occlusion, or dynamic scenes. We can also extend the framework to a multi-modal sensor fusion system, such as combining Lidar, event cameras to further enhance stability and precision.

**Scalable and Efficient 3D Registration.** The surfel-based registration method in this thesis provides accurate point cloud alignment but can struggle with large-scale scenes or cluttered environments with extremely low overlap. Incorporating Transformer-style global attention mechanisms or other probabilistic matching strategies could improve feature aggregation and robustness against sparse, noisy, or incomplete point clouds. Additionally, exploring fast, robust, memory-efficient correspondence learning could enhance registration performance for real-time SLAM and large-scale scene reconstruction. Lastly, integrating more efficient equivariant deep learning approaches, such as using 3D vector neurons [47] or modules that enforce rotation-equivariant properties, may improve the model's ability to generalize across varying orientations. However, as these designs often increase implementation complexity, a careful trade-off between rotation generalization and model design simplicity should also be considered.

**Memory-Efficient and More Generalizable Depth Estimation.** FocDepth-Former demonstrates strong depth estimation performance from focal stacks, but its high computational cost remains a bottleneck. Future work could explore lightweight Transformer architectures, knowledge distillation, or pruning techniques to enable deployment on edge devices or mobile platforms. Additionally, depth prediction from focal stacks faces inherent challenges, such as limited depth range and the need for numerous focal stack images to infer accurate disparity. The method also relies on clear focus and defocus cues, which can be difficult to extract in textureless or ambiguous scenes, like natural outdoor environments or plain walls, where such cues are absent. To address this, integrating focal stack techniques with state-of-the-art deep prior models trained on monocular or video input may help recover fine-grained depth details more robustly. Furthermore, adopting self-supervised learning paradigms could reduce dependence on large-scale, labeled focal stack datasets and improve generalization to unseen domains.

**Real-Time and More Robust Implicit Reconstruction for Reflective Surface.** Implicit SDF-based reconstruction offers high-quality 3D model-

126

ing, but real-time applications still require significant improvement in training and inference speed. Future work could focus on accelerating SDF optimization with diffusion models, real-time differentiable rendering, or even explore introducing the wavelet feature into the Gaussian Splatting method to leverage multi-band frequency features in the learning process. Furthermore, addressing challenges in reconstructing reflective surfaces, such as glass or mirrors, by explicitly modeling these materials or using learned uncertainty masks could improve robustness in real-world scenes. Lastly, extending the method to dynamic and deformable objects could enable real-time tracking of complex shapes and adaptive reconstruction, unleashing applications in robotics, AR/VR, and interactive digital modeling.

Beyond these individual research directions, this work also opens up broader opportunities for the fields:

Interactive AI for the virtual museum and efficient AI for a mobile robotic agent, connected to the Vision Language Model for real deployment on complex tasks, enable smarter robotics, where adaptive reconstruction enhances robotic interaction with complex dynamic environments.

Ultimately, this thesis demonstrates how geometric priors and deep learning models can be seamlessly integrated to advance 3D vision. Future research should focus on making these methods more efficient, generalizable, robust, and scalable, pushing the boundaries of robotics, AR/VR, digital twins, and other 3D vision fields beyond.

# Chapter 8

# Valorization Plan

This chapter outlines a concrete valorization plan for the novel geometric deep learning techniques developed in this research, focusing on four key areas: pose estimation, point cloud registration, depth estimation from focal stacks, and 3D reconstruction using implicit SDF models. While these advancements have broad applications in autonomous systems, robotics, and AR/VR, this work emphasizes their significant impact on digital cultural heritage, exploring practical use cases and deployment strategies.

First, we present a pose estimation system leveraging geometric primitives such as skylines and ground planes to derive camera orientation from image frames. By integrating an adaptive particle filter, the approach enhances robustness against sensor drift and environmental disturbances, demonstrating real-time feasibility on embedded hardware.

Next, we introduce a surfel-based $\mathbf{SE(3)}$-equivariant model for point cloud registration. Using a modular surfel representation, our method enables state-of-the-art alignment performance across diverse 3D scenes, offering potential extensions for real-time mapping and reconstruction.

For depth estimation, we propose FocDepthFormer, a hybrid Transformer-LSTM model that effectively aggregates focal stack images to recover fine depth details. While achieving high accuracy, the model's computational efficiency can be further optimized for large-scale applications such as defocus-based image synthesis.

Finally, we present an implicit SDF model that fuses wavelet-transformed depth features with triplane embeddings to improve shape reconstruction. Our approach surpasses explicit Gaussian Splatting in preserving fine geometric details, with future directions including optimized sampling and hybrid implicit-explicit representations for real-time performance.

Together, these contributions provide a robust foundation for applying geometric deep learning in cultural heritage digitization, offering scalable

solutions for high-fidelity 3D reconstruction, localization, and scene understanding.

## 8.1 Introduction

This chapter presents a comprehensive valorization plan for the geometric deep learning techniques developed in the previous chapters, with a primary focus on applications in digital cultural heritage. The plan outlines this research's social, practical, and academic impacts, structured around four core aspects: pose estimation from image frames, point cloud registration, focal stack depth estimation, and 3D reconstruction using an implicit SDF model. These techniques address critical challenges in 3D computer vision while advancing Virtual Reality (VR) and Augmented Reality (AR) technologies, with significant contributions to cultural preservation and education.

By bridging cutting-edge geometric deep learning research with practical applications, this work enhances digital preservation, interactive education, virtual tourism, archaeological analysis, and digital heritage restoration. High-fidelity 3D reconstructions enable the creation of digital twins for historic artifacts and buildings, while immersive learning tools enhance user interactive experience with cultural heritage.

Before discussing the detailed valorization plan, let us first recap the key techniques developed in this thesis:

**Pose estimation for the camera via natural cues.** Accurate pose estimation is crucial for 3D data capture and reconstruction, especially in challenging outdoor environments where traditional feature matching is unreliable when the inlier/outlier ratio is low. This research utilizes natural geometric cues, such as skylines and ground plane approximations, to estimate camera orientation. By extracting and analyzing these features, the system integrates inertial sensor data with vision-based estimation for robust orientation fusion. This method enhances stability in outdoor settings and significantly improves pose accuracy, making it particularly valuable for VR/AR applications in uncontrolled wild environments.

**Point cloud registration.** Aligning multiple point cloud frames is essential for constructing complete 3D scans. This research introduces novel algorithms leveraging data equivariance and uncertainty modeling to enhance registration accuracy. These techniques facilitate large-scale reconstructions of sculptures and buildings by merging partial scans into consistent models.

**Focal stack depth estimation.** Depth estimation from focal stacks captures fine details by analyzing images taken at varying focal distances. This approach enables high-precision 3D reconstruction of paintings, textiles, and

calligraphy, offering a flexible method for digitizing cultural artifacts and heritage.

**3D reconstruction of buildings.** Chapter 6 explores 3D reconstruction techniques for architectural structures, enabling detailed virtual models for historical preservation and visualization. These reconstructions support VR/AR applications for virtual tourism, interactive education, and heritage conservation, providing an accessible platform for studying architectural and cultural history.

Each of these contributions advances VR/AR technologies, mobile robotics, and self-driving, but here we will focus on the demo case for cultural heritage applications, facilitating digital preservation and interactive education, to give readers a more concrete understanding of the usability of these techniques. By integrating geometric deep learning with the digital cultural industry, this research paves the way for a cross-disciplinary field to use Geometric Deep Learning for history, culture, arts, and the gaming industry.

## 8.2 Use case of Pose Estimation for Camera Frame

Capturing stable, high-quality images and videos with UAV-mounted cameras, while simultaneously estimating the camera pose in real-time, is crucial for 3D applications such as reconstruction and 3D mapping in dynamic outdoor environments where wind and other disturbances can cause significant jitter. Chapter 3 presents an advanced camera pose estimation technique that can stabilize the camera by leveraging sensor fusion and natural geometric cues like skylines as reference. By accurately determining the camera orientation in real-time, the system compensates for unwanted movements, ensuring smooth, jitter-free video footage for downstream computer vision tasks. This improves data quality, enhances image quality and pose precision for mapping and surveying, and increases operational efficiency and stability by reducing the need for post-processing, ultimately making data collection in the wild more reliable for 3D applications like reconstruction.

Figure 8.1 illustrates pose estimation from camera images based on tracking features of images. These 2D features are projected into the 3D space to determine the camera pose.

### 8.2.1 Validation

A real demo of using a UAV equipped with a gimbal system, including our camera pose estimation framework for a flight to collect scan images of a

Figure 8.1: Camera pose is calculated from the live-streaming images, demo image used from [44].

building,g is provided in Figure 8.2,



Figure 8.2: Experimental Hardware Test Setup: Our camera pose estimation system is mounted on top of a UAV for real-world testing in a natural environment.

We conducted two tests to demonstrate practical applicability: one on the ground, where the hardware was mounted on a wooden stick facing the landscape from a rooftop for easy Ground Truth collection (Figure 8.3a), and another in the air using a UAV. In the aerial test, the onboard camera transmitted live images remotely while a real-time deep learning model generated skyline boundary segmentation output (Figure 8.3b).



(a) Camera pose estimation system tested on the ground by mounting it on a stick.

(b) Camera pose estimation system tested in the air by mounting it on a UAV

Figure 8.3: Camera pose estimation system tested both on the ground and in the air.

The camera pose tracking system can estimate pose robustly and quickly over long periods, without orientation drift. This demonstrates the reliable performance of our proposed tracking system, even when using a gimbal-based controller, to stabilize the camera orientation robustly.

## 8.2.2   Target Customers

This research work targets customers who require precise camera pose estimation for UAV-based 3D reconstruction and image capture in natural environments like mountains. For example, companies in aerial surveying and 3D mapping, such as those providing global mapping data for environmental monitoring and analysis, can leverage our efficient geometry-based orientation tracking system to improve the quality of camera pose for 3D reconstructions in challenging outdoor settings. Although the demo case in this chapter focuses on digital cultural heritage, it has applications beyond this field and can be extended to many other fields. For instance, Drone manufacturers and robotics companies, particularly those focused on autonomous navigation in outdoor terrains, can integrate our method to enhance motion stabilization and reduce orientation drift, ensuring high-quality image capture for downstream tasks. Additionally, entertainment industries

like filmmaking and natural life documentation, which often need to operate drones in unpredictable natural conditions, can benefit from the system's ability to mitigate motion blur and stabilize camera orientation using natural cues like skylines and ground planes. The proposed approach, with its real-time implementation on cheap embedded devices like the Jetson Nano, offers a practical and affordable solution for the software or as a complementary tool to the sensor fusion system for pose estimation. This can also be applied to cross-disciplinary research, such as for archaeology researchers at KU Leuven.

### 8.2.3 Economic Value

The primary economic value of our camera pose estimation system lies in its ability to boost the efficiency and quality of UAV-based imaging, creating opportunities in outdoor environments. By licensing this technology or partnering with drone manufacturers and aerial imaging companies in regions like Asia and Europe, we can expect the fast-growing drone market projected to reach billion-dollar valuations by 2030, according to the drone market report [121]. This research work offers a software toolkit solution bundled with drone manufacturing firms, driven by demand in surveying, environmental monitoring, and autonomous navigation in the wild. Our method design and real-time capabilities on embedded hardware cut the need for costly computational resources, offering a cost-effective solution for small- and medium-sized companies in these sectors. Additionally, its applications in filmmaking, natural life documentation, and virtual tourism, where reliable, high-quality 3D reconstructions are key for immersive experiences to support industry economic goals through sustainable data collection practices. By tackling the natural challenges of pose estimation, such as drift and noise in outdoor settings, this research delivers scalable, robust solutions that industries can adopt to enhance operational performance, reliability, and the overall value of drone applications.

## 8.3 Use case of Point Cloud Registration

Archaeology involves intricate engineering to excavate artifacts. Fragile materials such as porcelain can degrade over time due to many environmental factors. Reconstructing these artifacts from numerous fragments is labor-intensive, akin to solving a complex puzzle. To address this, each fragment is scanned and advanced registration techniques, as described in Chapter 4, are used. These techniques, known for their robust performance under

uncertainty, enable the precise and robust assembly of fragments into a complete 3D model, significantly enhancing the efficiency and accuracy of artifact recovery.

### 8.3.1 Validation

Figure 8.4 shows the registration of head scans in the body scan of terracotta warriors based on the overlapping regions, based on the proposed registration technique. This task is crucial for archaeologists and historians aiming to reconstruct and understand cultural artifacts from fragmented remains for digital recovery. The intricate process involves identifying and matching numerous scanned parts to recreate a complete and accurate representation of the original artifact, thereby preserving and interpreting historical and cultural heritage.



Figure 8.4: Registration of source and target scans of Terracotta Warriors [249].

### 8.3.2 Target Customers

This point cloud registration research targets customers who require a robust and accurate 3D alignment for large-scale reconstruction from Lidar scan tasks in diverse environments. Archaeologists can use this technique to assist in assembling cultural fragments by aligning scanned pieces within computer simulation software. Robotics companies developing autonomous

navigation systems can leverage our surfel-based $\mathbf{SE(3)}$-equivariant framework to improve spatial alignment performance in real-world indoor and outdoor settings, ensuring reliable mapping for applications like building global 3D maps for city blocks. Industries involved in 3D mapping and surveying, such as those in architecture or urban planning, can benefit from our framework's ability to handle noise and large rotations, producing consistent point cloud registrations for detailed digital twins of buildings or objects. Additionally, augmented reality (AR) developers creating immersive experiences can use our approach to align point clouds from LiDAR scans or depth maps, enabling the integration of virtual objects into real-world scenes for interaction. By addressing the limitations of traditional and learning-based methods, our framework offers a scalable solution for these fields, ensuring high accuracy in challenging registration scenarios.

### 8.3.3 Economic Value

The economic value of our surfel-based point cloud registration system lies in its ability to achieve robust and efficient 3D alignment, unlocking opportunities in high-growth industries. Licensing this technology to robotics and AR companies presents a promising billion-dollar market, particularly in the entertainment and education sectors. Our model design enhances robustness by reducing sensitivity to noise and large rotations through explicit $\mathbf{SE(3)}$-equivariant features, making registration more reliable even when the inlier-to-outlier ratio is low. Additionally, this technique mitigates the need for extensive scan transformation augmentations during training, offering a cost-effective solution for small- and medium-sized enterprises.

Beyond these applications, our system supports 3D mapping for architecture and surveying, enabling the creation of digital twins that reduce project costs and promote sustainable digital development. By addressing key challenges in point cloud registration, this research provides a high-precision tool that industries can leverage to enhance operational efficiency and drive innovation in 3D reconstruction applications.

For assembling complex cultural heritage fragments, advanced algorithms such as PuzzleFusion++ [250] can be used. These methods utilize iterative techniques to accurately search, match, and align fragments, particularly valuable when dealing with pieces that have minimal overlapping areas and irregular shapes that make candidate matching spurious and challenging. The algorithm begins by denoising and verifying the scanned fragments to have a clean input for geometry representation learning. In each iteration, the method incrementally refines the fit of the pieces, leveraging geometric and topological features to ensure consistency and accuracy between the pieces.

135

After refinement, the pieces are fitted more precisely, forming a coherent and accurate reconstruction of the original artifact. The ability to reconstruct artifacts digitally not only aids in preservation and study but also allows for virtual restoration and exhibition, making cultural heritage more accessible to the public, furthermore, the complete assembly digital scan can help the archaeologist to recover the full artifact by indicating the positions of each piece in the complete scan.

## 8.4 Use case of Depth Estimation from Focal Stack

Depth estimation from the focal stack technique can be used to extract 3D information from cultural paintings or sculptures, as artists have created these 2D paintings according to the rules of perspective projections.



Figure 8.5: Depth estimation for the planar sculpture of Marigold. The original image resource is from the public website link.

Figure 8.5 presents an example to demonstrate the application of depth estimation from monocular image techniques in the 3D modeling of art and sculpture. Depth estimation from the image is particularly well-suited for creating 3D models of intricate details and sculptures with planar constraints or primarily one-sided feature distributions. Unlike monocular depth estimation, the focal stack depth estimation technique captures a sweep of images of the target at varying focus distances and then uses algorithms to analyze the areas of sharp focus in each image to reconstruct depth information without requiring camera motion. Focal stack imaging offers several

key advantages for 3D digitization when camera motion is prohibited: as a passive perception technique, it enables non-invasive capture without physical contact, which is essential for preserving fragile artwork compared to active methods like LiDAR scanning. The approach creates high-resolution detail by utilizing multiple focused images with large stack sizes, capturing fine features such as intricate decorations and facial expressions. It excels at handling complex geometries with varying depth planes by differentiating subtle depth changes along stack dimensions, as evidenced by the accurate 3D feature representation, which is critical for the precise digital reproduction of sculptures without the need for camera motion, simply by adjusting the focal lens geometry distances continuously to reconstruct fine details.

### 8.4.1  Validation



Focal Stack Input          GT          DDFF    DefocusNet    Ours

Figure 8.6: Depth prediction from focal stack technique applied to heritage sculpture scanning.

The sweep of focal stack images of a sculpture in Figure 8.6 demonstrates how modern 3D computer vision techniques can extract 3D information from 2D focal stack images, especially those created with strong perspective principles by artists. The 3D reconstruction of historic artwork relies on visual analysis. Perspective analysis inverts spatial relationships by examining how artists employed pinhole camera models to create depth; the depth information can be inferred by our method through the focus/defocus cues in focal stack images, as the target scan details may be at the millimeter level, such as painting strokes, as fine-grained details. Thus, traditional viewpoint change-based methods for 3D reconstruction are not applicable, as the pose error is usually several orders of magnitude larger.

This application of depth estimation to paintings not only provides interesting insights into the use of perspective and camera focal geometry but also opens up possibilities for creating 3D visualizations or augmented reality experiences based on classical artworks. It represents an intersection of

art history and modern 3D computer vision technology. For instance, recent depth estimation techniques like "Depth Anything" [267] and "Marigold" [115] already show the great potential of the latest deep learning models applied to such challenging problems through the powerful learning representation ability of these models.

### 8.4.2 Target Customers

This depth estimation research is aimed at customers who need precise 3D models for cultural heritage preservation with small-scale level change, and academic study. Museums and cultural heritage organizations can use FocDepthFormer to create highly accurate 3D models of sculptures and artifacts like paintings, enabling detailed digital replicas for exhibitions or digital preservation without physical contact, which is crucial for fragile ancient artifacts at risk of being damaged. Researchers and art historians worldwide can benefit from these digital models for in-depth analysis of these historical artworks, studying intricate details like surface textures or structural patterns that might otherwise be lost to time. Additionally, educational institutions, museums, and virtual exhibition platforms can utilize these 3D reconstructions to provide immersive experiences, allowing global audiences to inspect cultural artifacts in detail. By addressing the challenge of handling arbitrary focal stack sizes, our method offers a flexible solution for these fields, ensuring high-quality depth estimation for diverse cultural business and research needs.

### 8.4.3 Economic Value

The economic value of FocDepthFormer lies in its ability to produce detailed 3D models through focal stack depth estimation, creating 3D assets in cultural heritage and related industries with just a sweep of images. By licensing this technology to museums and cultural preservation organizations, we hope to find opportunities in the growing market for affordable digital heritage solutions. The technique can expand to research institutes to document artifacts digitally. The creation of digital replicas for exhibitions or online platforms enhances global access to more cultural resources. The method is a cost-effective approach, using pre-training on monocular depth datasets and focal stack datasets to reduce reliance on expensive 3D data or multi-view data, making it an affordable option for smaller institutions, while its ability to preserve the digital twin of degrading and historic sculptures ensures long-term value for society through digitalization. Furthermore, applications in virtual exhibitions and education can generate revenue through subscription-

based access to the use of copyright of such digital 3D models, aligning with sustainable preservation practices. By overcoming limitations in traditional depth estimation, this research provides a scalable, high-precision tool that industries can use to advance cultural preservation and reproduction.

## 8.5 Use case of 3D Reconstruction via implicit SDF with Wavelet Feature Prior

Geometric deep learning techniques can help transform various aspects of the industry, such as robotics [255], self-driving, VR/AR, 3D printing or the gaming industry; here, we will focus on introducing the latter three fields in this part, which are more closely tied to the usage of digital cultural heritage. **VR/AR for Virtual museum.** In this digital era, 3D reconstruction techniques are essential for generating high-quality 3D assets. Cultural heritage artifacts can be digitally preserved eternally through the creation of digital 3D reconstruction, typically reconstructed using non-learning or deep learning models applied to multi-frame point clouds or images. These geometric models serve as the foundation for downstream tasks such as photorealistic rendering and relighting, enabling a more immersive and interactive experience of cultural heritage assets in high-quality virtual environments.

Once the digital twin generation becomes affordable and scalable, cultural heritage artifacts can establish virtual museum databases offering multiple benefits: VR/AR spatial computing creates immersive experiences that surpass traditional 2D displays by enabling visitors to navigate freely through the historical environments; interactive exploration allows users to manipulate virtual artifacts and navigate reconstructed ancient sites without access limits, promoting both educational and entertainment purposes; the mixing of reality with virtual properties offers new analytical perspectives on cultural heritage; global accessibility eliminates expensive travel cost while aiding digital preservation of fragile artifacts for the public to enjoy; educational value is enhanced through powerful interactive tools that place students within historical contexts; and heritage research capabilities are also expanded as VR environments provide platforms for visual analysis and efficient sharing of digital twins, potentially generating more great business insights.

By harnessing VR and AR technologies, virtual museums offer visitors immersive experiences that seamlessly mix technology, education, entertainment, and cultural art appreciation in unprecedented ways. Specifically, in these virtual museums, visitors can examine artifacts with close-up checks, view them from any angle, and even manipulate virtual artifacts. The in-

tegration of AR technology further enhances the experience by overlaying digital information onto physical spatial spaces or objects in real-time, providing additional guide introduction for tourism, as illustrated in Figure 8.7, it shows a traditional Chinese temple complex with architectural labels, illustrating the potential for visualizing cultural heritage sites through 3D digital reconstruction and augmented reality (AR). By overlaying digital information onto the physical structure, AR technology can provide visitors with a more immersive and educational experience, so that it allows for a vivid presentation of the temple components without physically altering the site.



Figure 8.7: 3D AR Guides overlaid with the real temple. The image was retrieved from public website link.

**Gaming industry.** The famous game design company, Ubisoft, has been a long-lasting pioneer in integrating historical 3D reconstructions of buildings and environments into their game products, particularly in the Assassin's Creed series. This approach has evolved to incorporate VR and AR technologies, enhancing the immersive experience for game players and the educational value of their games to learn history and culture.

Ubisoft development teams collaborate closely with historians and archaeologists to create detailed, historically accurate reconstructions of real ancient cities and landmark sites. These range from Ancient Egypt in Assassin's Creed Origins to Renaissance Italy. The company has also experimented with VR adaptations of their reconstructed sites for a better immersive experience. The level of detail and accuracy in these reconstructions has been so impressive that they are now used for educational purposes beyond gaming [20]. As such, high-quality digital twins of real cultural heritage sites have great potential for the gaming industry.

Figure 8.8 demonstrates the beautiful interior of a historic cathedral with

Figure 8.8: The demo image is from public website link.

intricate decorations, arched windows, and a huge domed ceiling. The player can be immersed in exploring this digital architectural heritage. Such collaborations between game developers and cultural institutions are opening up new possibilities for education, conservation, and public engagement with historical architecture, making centuries-old buildings accessible to people around the world through the power of technology. Lastly, this collaborative project can be scaled to other cultural relics in global cities.

Another example from the gaming industry is the video game "Black Myth Wukong" from China, which has generated considerable attention in the gaming world of 2024. The game story is based on an old Chinese mythology novel published in the Ming dynasty, and the video game has provided quite an immersive experience based on real cultural building scans across China, including traditional Chinese towers, pagodas, tower, and Buddha sculptures, as demonstrated in Figure 8.9.

These cases demonstrate the strong potential of 3D reconstruction techniques for generating assets in virtual gaming, a rapidly growing consumer market within the broader entertainment industry.

**3D Printing.**

Figure 8.10 exhibits another critical application of how 3D digital reconstruction technology is revolutionizing the reproduction of cultural heritage sites like the Dunhuang Mogao Grottoes. By creating detailed digital models of these ancient cave temples, preservationists can now offer immersive ex-

Figure 8.9: Chinese 3A Video Game, Black-myth Wukong, with many cultural heritage environments created from real scans of cultural buildings in Chinese cultural sites. The images depict the game's scene and character.

periences to visitors far from the original location through 3D printing. This technology not only allows tourists to enjoy but also enables the creation of precise 3D-printed replicas for other business purposes. Such reproductions can be displayed in museums worldwide, bringing the intricate art details and historical significance of the Dunhuang Mogao Grottoes to a global audience. This approach balances the need for the preservation of delicate original sites, which often have limited access due to conservation concerns or fragile preservation environments, with the desire to share these cultural treasures through 3D printing. The 3D printing results of such a digital twin of heritage paintings in caves presented here, with its vivid colors and detailed paintings, demonstrate how technology can bridge the gap between remote cultural sites and tourists and art enthusiasts, making these fragile and inaccessible heritage sites available for study and appreciation in new physical ways.

## 8.5.1  Validation

Our proposed reconstruction technique in Chapter 6 enables high-fidelity 3D reconstruction with fine-grained geometric details being preserved from multi-view images. Once the reconstruction is converted into a mesh, we can further apply post-processing steps to refine it by completing missing regions and removing noisy artifacts or unwanted background regions, ensuring a clean and accurate mesh representation. This post-processed mesh of cultural heritage sites can then be utilized for various downstream applications, such as high-quality rendering, recoloring (the left branch in Figure 8.11), or 3D printing of digital models (the right branch in Figure 8.11), as illustrated in Figure 8.11. This can prove the great usage potential of our reconstruction method.

Figure 8.10: 3D printed copy of Dunhuang Mogao Grottes inside the mountain caves. The image was retrieved from public website link.

The digitally reconstructed geometry can serve as a 3D representation primitive, such as a point cloud or mesh surface, for downstream applications like semantic segmentation [1], colorization [129], and texturing [225]. As a fundamental shape representation, it provides shared structural knowledge for various object and target representations, enabling the reuse of geometry for diverse post-processing tasks such as color and texture editing of 3D assets.

## 8.5.2 Target Customers

This 3D reconstruction research targets customers who need high-fidelity 3D models for diverse applications in virtual reality, the gaming industry, and the 3D printing industry. Educational or, museums, and training providers using VR/AR can leverage our wavelet-feature-conditioned implicit SDF model to create detailed, accurate historic architectures or cultural artifacts, enabling students or users to explore reconstructed cultural heritage sites or artifacts in the digital world. Gaming companies, especially those developing historically themed storytelling games, can use our method to generate detailed 3D assets from realistic data like images for digital worlds, improving players' immersive experience with virtual environments in games. Additionally, museums and cultural heritage organizations can benefit from 3D printing the reconstructed digital meshes to produce precise replicas of artifacts or

Figure 8.11: Practical application of my 3D reconstruction model: The heritage site mesh is first reconstructed from online images, followed by post-processing and specific downstream tasks. The left part shows recoloring to create a colored mesh, while the right part illustrates preparation for 3D printing.

cultural architectures, supporting preservation to make a copy of original heritage properties and public exhibitions to reduce the protection cost and risk for the valuable heritage collections. By addressing the challenge of capturing and learning fine-grained geometric details, our approach offers a comprehensive solution for these three fields, ensuring high-quality reconstructions across various object or scene scales.

### 8.5.3 Economic Value

The economic value of our 3D reconstruction system lies mainly in its ability to create detailed, accurate 3D models, opening opportunities in high-demand industries like virtual reality, the gaming industry, and 3D printing. By licensing this technology to VR/AR developers, we can penetrate the growing market for immersive training or exhibition programs, reaching a vast market, particularly in Asia, Europe, and the USA, covering a population of over multiple billion people. Such reconstructions of historical sites or cultural artifacts enhance training and educational experiences, which can be offered as a digital subscription service, allowing users to access our model via web or mobile app for a valid period. In the gaming industry, companies can easily integrate our high-fidelity 3D assets into their pipelines, reducing development costs for realistic environments and saving efforts in 3D asset creation, particularly for small or medium game development teams, driven by the growing market, game players, and investments. Furthermore, the ability to 3D print accurate replicas of reconstructed meshes offers economic potential for museums and cultural heritage institutions to broadcast or guide the protection of original, old, and fragile heritage properties, enabling them to sell or exhibit artifact replicas while supporting sustainable preservation practices. By overcoming limitations in capturing fine geometric details, this research provides a cheap, flexible, and high-precision solution that industries can adopt to drive innovation in digital 3D applications and create great value for the mass market in entertainment or cultural research fields.

# Bibliography

[1] ABDELREHEEM, A., SKOROKHODOV, I., OVSJANIKOV, M., AND WONKA, P. Satr: Zero-shot semantic segmentation of 3d shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 15166–15179.

[2] ACHANTA, R., SHAJI, A., SMITH, K., LUCCHI, A., FUA, P., AND SUSSTRUNK, S. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence 34* (05 2012).

[3] AGARWAL, A., AND ARORA, C. Depthformer: Multiscale vision transformer for monocular depth estimation with global local information fusion. In *2022 IEEE International Conference on Image Processing (ICIP)* (2022), IEEE, pp. 3873–3877.

[4] ALPERT, B. K. Hybrid gauss-trapezoidal quadrature rules. *SIAM Journal on Scientific Computing 20*, 5 (1999), 1551–1584.

[5] ALQUISIRIS-QUECHA, O., AND MARTINEZ-CARRANZA, J. Video stabilization of the nao robot using imu data. In *Robot Operating System (ROS)*. Springer, Cham, 2020, pp. 147–162.

[6] AMINI, A., SELVAM PERIYASAMY, A., AND BEHNKE, S. Yolopose: Transformer-based multi-object 6d pose estimation using keypoint regression. In *International Conference on Intelligent Autonomous Systems* (2022), Springer, pp. 392–406.

[7] ANWAR, S., HAYDER, Z., AND PORIKLI, F. Deblur and deep depth from single defocus image. *Machine vision and applications 32*, 1 (2021), 1–13.

[8] AO, S., HU, Q., YANG, B., MARKHAM, A., AND GUO, Y. Spinnet: Learning a general surface descriptor for 3d point cloud registration.

In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2021), pp. 11753–11762.

[9] Auysakul, J., Xu, H., and Pooneeth, V. A hybrid motion estimation for video stabilization based on an imu sensor. *Sensors 18*, 8 (2018), 2708.

[10] Bae, G., Budvytis, I., and Cipolla, R. Multi-view depth estimation by fusing single-view depth probability with multi-view geometry. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 2842–2851.

[11] Bai, X., Luo, Z., Zhou, L., Chen, H., Li, L., Hu, Z., Fu, H., and Tai, C.-L. Pointdsc: Robust point cloud registration using deep spatial consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 15859–15869.

[12] Bai, X., Luo, Z., Zhou, L., Fu, H., Quan, L., and Tai, C.-L. D3feat: Joint learning of dense detection and description of 3d local features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2020), pp. 6359–6367.

[13] Barratt, S., and Hannel, B. Extracting the depth and all-in-focus image from a focal stack. In *Proceedings of the IEEE International Conference on Computer Vision* (2015), pp. 3451–3459.

[14] Barron, J. T. A general and adaptive robust loss function. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 4331–4339.

[15] Barron, J. T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., and Srinivasan, P. P. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 5855–5864.

[16] Barron, J. T., Mildenhall, B., Verbin, D., Srinivasan, P. P., and Hedman, P. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 5470–5479.

[17] Battiato, S., et al. Sift features tracking for video stabilization. In *14th International Conference on Image Analysis and Processing (ICIAP 2007)* (2007), IEEE.

[18] BEHLEY, J., AND STACHNISS, C. Efficient surfel-based slam using 3d laser range data in urban environments. In *Robotics: Science and Systems* (2018), vol. 2018, p. 59.

[19] BENAVIDES, F. T., IGNATOV, A., AND TIMOFTE, R. Phonedepth: A dataset for monocular depth estimation on mobile devices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 3049–3056.

[20] BEVILACQUA, M. G., RUSSO, M., GIORDANO, A., AND SPALLONE, R. 3d reconstruction, digital twinning, and virtual reality: Architectural heritage applications. In *2022 IEEE Conference on Virtual Reality and 3D User Interfaces Abstracts and Workshops (VRW)* (2022), IEEE, pp. 92–96.

[21] BHATNAGAR, B. L., SMINCHISESCU, C., THEOBALT, C., AND PONS-MOLL, G. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. *Advances in Neural Information Processing Systems 33* (2020), 12909–12922.

[22] BRACHMANN, E., AND ROTHER, C. Learning less is more-6d camera localization via 3d surface regression. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 4654–4662.

[23] BRONSTEIN, M. M., BRUNA, J., LECUN, Y., SZLAM, A., AND VANDERGHEYNST, P. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine 34*, 4 (2017), 18–42.

[24] BUCHHOLZ, T.-O., AND JUG, F. Fourier image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 1846–1854.

[25] CAMPOS, C., ELVIRA, R., GOMEZ, J. J., MONTIEL, J. M. M., AND TARDOS, J. D. ORB-SLAM3: An accurate open-source library for visual, visual-inertial and multi-map SLAM. *IEEE Transactions on Robotics 37*, 6 (2021), 1874–1890.

[26] CANNY, J., ET AL. A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence 6* (1986), 679–698.

[27] CARRIO, A., BAVLE, H., AND CAMPOY, P. Attitude estimation using horizon detection in thermal images. *International Journal of Micro Air Vehicles 10*, 4 (2018), 352–361.

[28] CARVALHO, M., LE SAUX, B., TROUVÉ-PELOUX, P., ALMANSA, A., AND CHAMPAGNAT, F. Deep depth from defocus: how can defocus blur improve 3d estimation using dense neural networks? In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops* (2018), pp. 0–0.

[29] CASSER, V., PIRK, S., MAHJOURIAN, R., AND ANGELOVA, A. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI conference on artificial intelligence* (2019), vol. 33, pp. 8001–8008.

[30] CESA, G., LANG, L., AND WEILER, M. A program to build e (n)-equivariant steerable cnns. In *International conference on learning representations* (2022).

[31] CHABRA, R., LENSSEN, J. E., ILG, E., SCHMIDT, T., STRAUB, J., LOVEGROVE, S., AND NEWCOMBE, R. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *Computer Vision– ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16* (2020), Springer, pp. 608–625.

[32] CHANG, S. G., YU, B., AND VETTERLI, M. Adaptive wavelet thresholding for image denoising and compression. *IEEE transactions on image processing 9*, 9 (2000), 1532–1546.

[33] CHEN, D., LI, H., YE, W., WANG, Y., XIE, W., ZHAI, S., WANG, N., LIU, H., BAO, H., AND ZHANG, G. Pgsr: Planar-based gaussian splatting for efficient and high-fidelity surface reconstruction. *arXiv preprint arXiv:2406.06521* (2024).

[34] CHEN, H., LIU, S., CHEN, W., LI, H., AND HILL, R. Equivariant point network for 3d point cloud analysis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2021), pp. 14514–14523.

[35] CHEN, S., CAVALLARI, T., PRISACARIU, V. A., AND BRACHMANN, E. Map-relative pose regression for visual re-localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 20665–20674.

[36] CHEN, T., LIN, L., ZUO, W., LUO, X., AND ZHANG, L. Learning a wavelet-like auto-encoder to accelerate deep neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2018), vol. 32.

[37] CHENG, Y.-C., LEE, H.-Y., TULYAKOV, S., SCHWING, A. G., AND GUI, L.-Y. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 4456–4465.

[38] CHIBANE, J., ALLDIECK, T., AND PONS-MOLL, G. Implicit functions in feature space for 3d shape reconstruction and completion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (jun 2020), IEEE.

[39] CHOI, J., PARK, J., S, I., AND KWEON. Self-supervised real-time video stabilization, 2021. arXiv preprint arXiv:2111.05980.

[40] CHOU, G., BAHAT, Y., AND HEIDE, F. Diffusion-sdf: Conditional generative modeling of signed distance functions. In *Proceedings of the IEEE/CVF international conference on computer vision* (2023), pp. 2262–2272.

[41] CHOY, C., DONG, W., AND KOLTUN, V. Deep global registration. In *CVPR* (2020).

[42] CHOY, C., PARK, J., AND KOLTUN, V. Fully convolutional geometric features. In *Proceedings of the IEEE/CVF international conference on computer vision* (2019), pp. 8958–8966.

[43] CHOY, C. B., XU, D., GWAK, J., CHEN, K., AND SAVARESE, S. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proceedings of the European Conference on Computer Vision (ECCV)* (2016).

[44] CITRARO, L., MÁRQUEZ-NEILA, P., SAVARE, S., JAYARAM, V., DUBOUT, C., RENAUT, F., HASFURA, A., BEN SHITRIT, H., AND FUA, P. Real-time camera pose estimation for sports fields. *Machine Vision and Applications 31*, 3 (2020), 16.

[45] COHEN, T. S., GEIGER, M., KÖHLER, J., AND WELLING, M. Spherical cnns. *arXiv preprint arXiv:1801.10130* (2018).

[46] DAHL, V., AANÆS, H., AND BÆRENTZEN, J. Surfel based geometry reconstruction. pp. 39–44.

[47] DENG, C., LITANY, O., DUAN, Y., POULENARD, A., TAGLIASACCHI, A., AND GUIBAS, L. J. Vector neurons: A general framework for so (3)-equivariant networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 12200–12209.

[48] DOSOVITSKIY, A., BEYER, L., KOLESNIKOV, A., WEISSENBORN, D., ZHAI, X., UNTERTHINER, T., DEHGHANI, M., MINDERER, M., HEIGOLD, G., GELLY, S., ET AL. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[49] DU, J., WANG, R., AND CREMERS, D. Dh3d: Deep hierarchical 3d descriptors for robust large-scale 6dof relocalization. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16* (2020), Springer, pp. 744–762.

[50] DU, W., ZHANG, H., DU, Y., MENG, Q., CHEN, W., ZHENG, N., SHAO, B., AND LIU, T.-Y. Se(3) equivariant graph neural networks with complete local frames. In *International Conference on Machine Learning* (2022), PMLR, pp. 5583–5608.

[51] DUSHA, D., BOLES, W., AND WALKER, R. Fixed-wing attitude estimation using computer vision based horizon detection. In *Proceedings of AIAC12: 2nd Australasian Unmanned Air Vehicles Conference* (2007), Waldron Smith Management.

[52] EIGEN, D., PUHRSCH, C., AND FERGUS, R. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems 27* (2014).

[53] FAN, H., SU, H., AND GUIBAS, L. J. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 605–613.

[54] FENG, Y., WU, S., KÖPÜKLÜ, O., KANG, X., AND TOMBARI, F. Unsupervised monocular depth prediction for indoor continuous video streams. *arXiv preprint arXiv:1911.08995* (2019).

[55] FIGUEIREDO, M. A., AND NOWAK, R. D. An em algorithm for wavelet-based image restoration. *IEEE Transactions on Image Processing 12*, 8 (2003), 906–916.

[56] FINZI, M., STANTON, S., IZMAILOV, P., AND WILSON, A. G. Generalizing convolutional neural networks for equivariance to lie groups on arbitrary continuous data. In *International Conference on Machine Learning* (2020), PMLR, pp. 3165–3176.

[57] FISCHLER, M. A., AND BOLLES, R. C. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM 24*, 6 (1981), 381–395.

[58] FU, X., YIN, W., HU, M., WANG, K., MA, Y., TAN, P., SHEN, S., LIN, D., AND LONG, X. Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a single image. In *ECCV* (2024).

[59] FUCHS, F., WORRALL, D., FISCHER, V., AND WELLING, M. Se (3)-transformers: 3d roto-translation equivariant attention networks. *Advances in neural information processing systems 33* (2020), 1970–1981.

[60] FUCHS, F. B., WAGSTAFF, E., DAUPARAS, J., AND POSNER, I. Iterative se (3)-transformers. In *Geometric Science of Information: 5th International Conference, GSI 2021, Paris, France, July 21–23, 2021, Proceedings 5* (2021), Springer, pp. 585–595.

[61] FUENTES-PACHECO, J., RUIZ-ASCENCIO, J., AND RENDÓN-MANCHA, J. M. Visual simultaneous localization and mapping: a survey. *Artificial intelligence review 43* (2015), 55–81.

[62] FUJIEDA, S., TAKAYAMA, K., AND HACHISUKA, T. Wavelet convolutional neural networks. *arXiv preprint arXiv:1805.08620* (2018).

[63] FUJIMURA, Y., IIYAMA, M., FUNATOMI, T., AND MUKAIGAWA, Y. Deep depth from focal stack with defocus model for camera-setting invariance. *International Journal of Computer Vision* (2023), 1–16.

[64] FURGALE, P., REHDER, J., AND SIEGWART, R. Unified temporal and spatial calibration for multi-sensor systems. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems* (2013), IEEE, pp. 1280–1286.

[65] GARG, R., BG, V. K., CARNEIRO, G., AND REID, I. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII 14* (2016), Springer, pp. 740–756.

[66] GEIGER, A., LENZ, P., AND URTASUN, R. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition* (2012), IEEE, pp. 3354–3361.

[67] GLOCKER, B., IZADI, S., SHOTTON, J., AND CRIMINISI, A. Real-time rgb-d camera relocalization. In *2013 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)* (2013), IEEE, pp. 173–179.

[68] GODARD, C., MAC AODHA, O., AND BROSTOW, G. J. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 270–279.

[69] GODARD, C., MAC AODHA, O., FIRMAN, M., AND BROSTOW, G. J. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 3828–3838.

[70] GODARD, C., MAC AODHA, O., FIRMAN, M., AND BROSTOW, G. J. Digging into self-supervised monocular depth prediction.

[71] GOJCIC, Z., ZHOU, C., WEGNER, J. D., GUIBAS, L. J., AND BIRDAL, T. Learning multiview 3d point cloud registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2020), pp. 1759–1769.

[72] GORDON, A., LI, H., JONSCHKOWSKI, R., AND ANGELOVA, A. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 8977–8986.

[73] GUÉDON, A., AND LEPETIT, V. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 5354–5363.

[74] GUO, T., SEYED MOUSAVI, H., HUU VU, T., AND MONGA, V. Deep wavelet prediction for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops* (2017), pp. 104–113.

[75] GUO, X., LI, H., YI, S., REN, J., AND WANG, X. Learning monocular depth by distilling cross-domain stereo networks. In *Proceedings of the European Conference on Computer Vision (ECCV)* (2018), pp. 484–500.

[76] GUR, S., AND WOLF, L. Single image depth estimation trained via depth from defocus cues. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 7683–7692.

[77] GUSTAFSSON, F., ET AL. Particle filter theory and practice with positioning applications. *IEEE Aerospace and Electronic Systems Magazine 25*, 7 (2010), 53–82.

[78] HASSON, Y., TEKIN, B., BOGO, F., LAPTEV, I., POLLEFEYS, M., AND SCHMID, C. Leveraging photometric consistency over time for sparsely supervised hand-object reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2020), pp. 571–580.

[79] HAZIRBAS, C., SOYER, S. G., STAAB, M. C., LEAL-TAIXÉ, L., AND CREMERS, D. Deep depth from focus. In *Asian conference on computer vision* (2018), Springer, pp. 525–541.

[80] HE, J., LI, H., YIN, W., LIANG, Y., LI, L., ZHOU, K., LIU, H., LIU, B., AND CHEN, Y.-C. Lotus: Diffusion-based visual foundation model for high-quality dense prediction. *arXiv preprint arXiv:2409.18124* (2024).

[81] HE, R., HONG, H., FU, B., AND LIU, F. Multi-task learning for monocular depth and defocus estimations with real images. *arXiv preprint arXiv:2208.09848* (2022).

[82] HE, Z., YANG, M., FENG, M., YIN, J., WANG, X., LENG, J., AND LIN, Z. Fourier transformer: Fast long range modeling by removing sequence redundancy with fft operator. *arXiv preprint arXiv:2305.15099* (2023).

154

[83] HESS, W., KOHLER, D., RAPP, H., AND ANDOR, D. Real-time loop closure in 2d lidar slam. In *2016 IEEE international conference on robotics and automation (ICRA)* (2016), IEEE, pp. 1271–1278.

[84] HOCHREITER, S., AND SCHMIDHUBER, J. Long short-term memory. *Neural computation 9*, 8 (1997), 1735–1780.

[85] HONAUER, K., JOHANNSEN, O., KONDERMANN, D., AND GOLD-LUECKE, B. A dataset and evaluation methodology for depth estimation on 4d light fields. In *Asian conference on computer vision* (2016), Springer, pp. 19–34.

[86] HORNAUER, J., AND BELAGIANNIS, V. Gradient-based uncertainty for monocular depth estimation. In *European Conference on Computer Vision* (2022), Springer, pp. 613–630.

[87] HU, J., HUI, K.-H., LIU, Z., LI, R., AND FU, C.-W. Neural wavelet-domain diffusion for 3d shape generation, inversion, and manipulation. *ACM Transactions on Graphics 43*, 2 (2024), 1–18.

[88] HU, J., OZAY, M., ZHANG, Y., AND OKATANI, T. Revisiting single image depth estimation: Toward higher resolution maps with accurate object boundaries. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)* (2019), IEEE, pp. 1043–1051.

[89] HU, Y., RAO, W., QI, L., DONG, J., CAI, J., AND FAN, H. A refractive stereo structured-light 3-d measurement system for immersed object. *IEEE Transactions on Instrumentation and Measurement 72* (2022), 1–13.

[90] HUAI, J., ZHANG, Y., AND YILMAZ, A. Real-time large scale 3d reconstruction by fusing kinect and imu data. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences 2* (2015), 491–496.

[91] HUAI, Z., AND HUANG, G. Robocentric visual-inertial odometry. *The International Journal of Robotics Research 41*, 7 (2022), 667–689.

[92] HUANG, B., YU, Z., CHEN, A., GEIGER, A., AND GAO, S. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 conference papers* (2024), pp. 1–11.

[93] HUANG, B., YU, Z., CHEN, A., GEIGER, A., AND GAO, S. 2d gaussian splatting for geometrically accurate radiance fields. In *SIGGRAPH 2024 Conference Papers* (2024), Association for Computing Machinery.

[94] HUANG, B., YU, Z., CHEN, A., GEIGER, A., AND GAO, S. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 conference papers* (2024), pp. 1–11.

[95] HUANG, H., HE, R., SUN, Z., AND TAN, T. Wavelet-srnet: A wavelet-based cnn for multi-scale face super resolution. In *Proceedings of the IEEE international conference on computer vision* (2017), pp. 1689–1697.

[96] HUANG, Y., HUANG, J., LIU, J., YAN, M., DONG, Y., LYU, J., CHEN, C., AND CHEN, S. Wavedm: Wavelet-based diffusion models for image restoration. *IEEE Transactions on Multimedia* (2024).

[97] HUANG, Z., XU, P., LIANG, D., MISHRA, A., AND XIANG, B. Trans-blstm: Transformer with bidirectional lstm for language understanding. *arXiv preprint arXiv:2003.07000* (2020).

[98] HUI, K.-H., LI, R., HU, J., AND FU, C.-W. Neural wavelet-domain diffusion for 3d shape generation. In *SIGGRAPH Asia 2022 Conference Papers* (2022), pp. 1–9.

[99] HUTCHINS, D., SCHLAG, I., WU, Y., DYER, E., AND NEYSHABUR, B. Block-recurrent transformers. *arXiv preprint arXiv:2203.07852* (2022).

[100] HUTCHINSON, M., LAN, C. L., ZAIDI, S., DUPONT, E., TEH, Y. W., AND KIM, H. Lietransformer: Equivariant self-attention for lie groups, 2020.

[101] IZADI, S., KIM, D., HILLIGES, O., MOLYNEAUX, D., NEWCOMBE, R., KOHLI, P., SHOTTON, J., HODGES, S., FREEMAN, D., DAVISON, A., ET AL. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology* (2011), pp. 559–568.

[102] JENNER, E., AND WEILER, M. Steerable partial differential operators for equivariant neural networks. In *International Conference on Learning Representations* (2022).

[103] JENSEN, R., DAHL, A., VOGIATZIS, G., TOLA, E., AND AANÆS, H. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2014), pp. 406–413.

[104] JIANG, Y., TU, J., LIU, Y., GAO, X., LONG, X., WANG, W., AND MA, Y. Gaussianshader: 3d gaussian splatting with shading functions for reflective surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 5322–5332.

[105] JOHANNSEN, O., HONAUER, K., GOLDLUECKE, B., ALPEROVICH, A., BATTISTI, F., BOK, Y., BRIZZI, M., CARLI, M., CHOE, G., DIEBOLD, M., ET AL. A taxonomy and evaluation of dense light field depth estimation algorithms. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops* (2017), pp. 82–99.

[106] KANG, X., HAN, F., FAYJIE, A. R., VANDEWALLE, P., KHOSHELHAM, K., AND GONG, D. Focdepthformer: Transformer with latent lstm for depth estimation from focal stack. In *Australasian Joint Conference on Artificial Intelligence* (2024), Springer, pp. 273–290.

[107] KANG, X., HERRERA, A., LEMA, H., VALENCIA, E., AND VANDEWALLE, P. Adaptive sampling-based particle filter for visual-inertial gimbal in the wild. In *2023 IEEE International Conference on Robotics and Automation (ICRA)* (2023), IEEE, pp. 2738–2744.

[108] KANG, X., LUAN, Z., KHOSHELHAM, K., AND WANG, B. Equi-gspr: Equivariant se (3) graph network model for sparse point cloud registration. In *European Conference on Computer Vision* (2024), Springer, pp. 149–167.

[109] KANG, X., XIANG, Z., ZHANG, Z., AND KHOSHELHAM, K. Multiview geometry-aware diffusion transformer for indoor novel view synthesis. In *ICLR 2025 Workshop on Deep Generative Model in Machine Learning: Theory, Principle and Efficacy* (2025).

[110] KANG, X., YIN, S., AND FEN, Y. 3d reconstruction & assessment framework based on affordable 2d lidar. In *2018 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM)* (2018), IEEE, pp. 292–297.

[111] KANG, X., AND YUAN, S. Robust data association for object-level semantic slam. *arXiv preprint arXiv:1909.13493* (2019).

[112] KANG, X., AND YUAN, S. Integrated visual-inertial odometry and image stabilization for image processing. *Google Patents, US Patent App 18*, 035,479 (2023).

[113] KANG, X., AND YUAN, S. Integrated visual-inertial odometry and image stabilization for image processing, December 28 2023.

[114] KANG, X., ZHAO, H., KHOSHELHAM, K., AND PATRICK, V. 2d surfel-based 3d point cloud registration with robust equivariant se (3) features. In *The conference proceedings and published in IEEE Xplore of 2025 IEEE International Geoscience and Remote Sensing Symposium* (2025).

[115] KE, B., OBUKHOV, A., HUANG, S., METZGER, N., DAUDT, R. C., AND SCHINDLER, K. Repurposing diffusion-based image generators for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 9492–9502.

[116] KENDALL, A., GRIMES, M., AND CIPOLLA, R. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision* (2015), pp. 2938–2946.

[117] KERBL, B., KOPANAS, G., LEIMKÜHLER, T., AND DRETTAKIS, G. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics 42*, 4 (July 2023).

[118] KERBL, B., KOPANAS, G., LEIMKÜHLER, T., AND DRETTAKIS, G. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics 42*, 4 (2023), 1–14.

[119] KHAMIS, S., FANELLO, S., RHEMANN, C., KOWDLE, A., VALENTIN, J., AND IZADI, S. Stereonet: Guided hierarchical refinement for real-time edge-aware depth prediction. In *Proceedings of the European conference on computer vision (ECCV)* (2018), pp. 573–590.

[120] KHATIB, R., AND GIRYES, R. Trinerflet: A wavelet based multiscale triplane nerf representation. *arXiv preprint arXiv:2401.06191* (2024).

[121] KIPPONEN, S. Defining the market potential of industry specific drone software: Case: Metropolia innovation hub of smart mobility.

[122] Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A. C., Lo, W.-Y., Dollár, P., and Girshick, R. Segment anything. *arXiv:2304.02643* (2023).

[123] Knapitsch, A., Park, J., Zhou, Q.-Y., and Koltun, V. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG) 36*, 4 (2017), 1–13.

[124] Kong, X., Liu, S., Taher, M., and Davison, A. J. vmap: Vectorised object mapping for neural field slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 952–961.

[125] La Place, C., Urooj, A., and Borji, A. Segmenting sky pixels in images: Analysis and comparison. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)* (2019), IEEE.

[126] Laina, I., Rupprecht, C., Belagiannis, V., Tombari, F., and Navab, N. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth international conference on 3D vision (3DV)* (2016), IEEE, pp. 239–248.

[127] Lavin, A., and Gray, S. Fast algorithms for convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 4013–4021.

[128] Lee, Y.-C., et al. 3d video stabilization with depth estimation by cnn-based optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021).

[129] Leifman, G., and Tal, A. Pattern-driven colorization of 3d surfaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2013), pp. 241–248.

[130] Li, C., et al. Deep online video stabilization using imu sensors. *IEEE Transactions on Multimedia* (2022).

[131] Li, J. C. L., Liu, C., Huang, B., and Wong, N. Learning spatially collaged fourier bases for implicit neural representation. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2024), vol. 38, pp. 13492–13499.

[132] LI, M., DUAN, Y., ZHOU, J., AND LU, J. Diffusion-sdf: Text-to-shape via voxelized diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 12642–12651.

[133] LI, Y., LYU, C., DI, Y., ZHAI, G., LEE, G. H., AND TOMBARI, F. Geogaussian: Geometry-aware gaussian splatting for scene rendering. In *European Conference on Computer Vision* (2025), Springer, pp. 441–457.

[134] LI, Z., MÜLLER, T., EVANS, A., TAYLOR, R. H., UNBERATH, M., LIU, M.-Y., AND LIN, C.-H. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 8456–8465.

[135] LI, Z., YEH, Y.-Y., AND CHANDRAKER, M. Through the looking glass: Neural 3d reconstruction of transparent shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), pp. 1262–1271.

[136] LIN, C. E., ZHU, M., AND GHAFFARI, M. Se3et: Se (3)-equivariant transformer for low-overlap point cloud registration. *IEEE Robotics and Automation Letters* (2024).

[137] LIN, H., CHEN, C., KANG, S. B., AND YU, J. Depth recovery from light field using focal stack symmetry. In *Proceedings of the IEEE International Conference on Computer Vision* (2015), pp. 3451–3459.

[138] LINDEBERG, T. Scale invariant feature transform.

[139] LINDEBERG, T., ET AL. Scale invariant feature transform. In *Proceedings of the International Conference on Computer Vision* (2012), p. 10491.

[140] LINDELL, D. B., VAN VEEN, D., PARK, J. J., AND WETZSTEIN, G. Bacon: Band-limited coordinate networks for multiscale scene representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 16252–16262.

[141] LIU, C., QIU, J., AND JIANG, M. Light field reconstruction from focal stack based on landweber iterative scheme. In *Mathematics in Imaging* (2017), Optica Publishing Group, pp. MM2C–3.

[142] LIU, J., TANG, X., CHENG, F., YANG, R., LI, Z., LIU, J., HUANG, Y., LIN, J., LIU, S., WU, X., ET AL. Mirrorgaussian: Reflecting 3d gaussians for reconstructing mirror reflections. In *European Conference on Computer Vision* (2024), Springer, pp. 377–393.

[143] LIU, R., LAUZE, F., BEKKERS, E., ERLEBEN, K., AND DARKNER, S. Se (3) group convolutional neural networks and a study on group convolutions and equivariance for dwi segmentation.

[144] LIU, Y., PENG, S., LIU, L., WANG, Q., WANG, P., THEOBALT, C., ZHOU, X., AND WANG, W. Neural rays for occlusion-aware image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 7824–7833.

[145] LIU, Y.-L., ET AL. Hybrid neural fusion for full-frame video stabilization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021).

[146] LIU, Z., LIN, Y., CAO, Y., HU, H., WEI, Y., ZHANG, Z., LIN, S., AND GUO, B. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 10012–10022.

[147] LIU, Z., ZHU, H., ZHANG, Q., FU, J., DENG, W., MA, Z., GUO, Y., AND CAO, X. Finer: Flexible spectral-bias tuning in implicit neural representation by variable-periodic activation functions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 2713–2722.

[148] LORENSEN, W. E., AND CLINE, H. E. Marching cubes: A high resolution 3d surface construction algorithm. In *Seminal graphics: pioneering efforts that shaped the field.* 1998, pp. 347–353.

[149] LOW, K.-L. Linear least-squares optimization for point-to-plane icp surface registration. *Chapel Hill, University of North Carolina 4*, 10 (2004), 1–3.

[150] LOWE, G. Sift-the scale invariant feature transform. *Int. J 2*, 91-110 (2004), 2.

[151] MA, F., AND KARAMAN, S. Sparse-to-dense: Depth prediction from sparse depth samples and a single image. In *2018 IEEE international conference on robotics and automation (ICRA)* (2018), IEEE, pp. 4796–4803.

[152] MA, H., LIU, D., YAN, N., LI, H., AND WU, F. End-to-end optimized versatile image compression with wavelet-like transform. *IEEE Transactions on Pattern Analysis and Machine Intelligence 44*, 3 (2020), 1247–1263.

[153] MADGWICK, S. An efficient orientation filter for inertial and inertial/magnetic sensor arrays. Tech. Rep. 25, x-io and University of Bristol (UK), 2010.

[154] MAHENDRAN, S., ALI, H., AND VIDAL, R. 3d pose regression using convolutional neural networks. In *Proceedings of the IEEE international conference on computer vision workshops* (2017), pp. 2174–2182.

[155] MALLAT, S. G. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE transactions on pattern analysis and machine intelligence 11*, 7 (1989), 674–693.

[156] MAO, W., GE, Y., SHEN, C., TIAN, Z., WANG, X., AND WANG, Z. Tfpose: Direct human pose estimation with transformers. *arXiv preprint arXiv:2103.15320* (2021).

[157] MARTIN-BRUALLA, R., RADWAN, N., SAJJADI, M. S. M., BARRON, J. T., DOSOVITSKIY, A., AND DUCKWORTH, D. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *CVPR* (2021).

[158] MATSUKI, H., MURAI, R., KELLY, P. H., AND DAVISON, A. J. Gaussian splatting slam. *arXiv preprint arXiv:2312.06741* (2023).

[159] MATSUSHITA, Y., OFEK, E., TANG, X., AND SHUM, H.-Y. Full-frame video stabilization. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)* (2005), vol. 1, IEEE, pp. 50–57.

[160] MATURANA, D., AND SCHERER, S. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)* (2015), IEEE, pp. 922–928.

[161] MAXIMOV, M., GALIM, K., AND LEAL-TAIXÉ, L. Focus on defocus: bridging the synthetic to real domain gap for depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), pp. 1071–1080.

[162] MENG, X., FAN, C., MING, Y., AND YU, H. Cornet: Context-based ordinal regression network for monocular depth estimation. *IEEE Transactions on Circuits and Systems for Video Technology* (2021).

[163] MESHRY, M., GOLDMAN, D. B., KHAMIS, S., HOPPE, H., PANDEY, R., SNAVELY, N., AND MARTIN-BRUALLA, R. Neural rerendering in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 6878–6887.

[164] MIANGOLEH, S. M. H., DILLE, S., MAI, L., PARIS, S., AND AKSOY, Y. Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 9685–9694.

[165] MIHAIL, R. P. W., BESSINGER, S., JACOBS, Z., AND JACOBS, N. Sky segmentation in the wild: An empirical study. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)* (March 2016), IEEE, pp. 1–6.

[166] MILDENHALL, B., HEDMAN, P., MARTIN-BRUALLA, R., SRINIVASAN, P. P., AND BARRON, J. T. Nerf in the dark: High dynamic range view synthesis from noisy raw images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 16190–16199.

[167] MILDENHALL, B., SRINIVASAN, P. P., TANCIK, M., BARRON, J. T., RAMAMOORTHI, R., AND NG, R. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV* (2020).

[168] MILDENHALL, B., SRINIVASAN, P. P., TANCIK, M., BARRON, J. T., RAMAMOORTHI, R., AND NG, R. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM 65*, 1 (2021), 99–106.

[169] MISHRA, D., SINGH, S. K., AND SINGH, R. K. Wavelet-based deep auto encoder-decoder (wdaed)-based image compression. *IEEE Transactions on Circuits and Systems for Video Technology 31*, 4 (2020), 1452–1462.

[170] MITRA, N. J., FLÖRY, S., OVSJANIKOV, M., GELFAND, N., GUIBAS, L. J., AND POTTMANN, H. Dynamic geometry registration. In *Symposium on geometry processing* (2007), pp. 173–182.

[171] MOELLER, M., BENNING, M., SCHÖNLIEB, C., AND CREMERS, D. Variational depth from focus reconstruction. *IEEE Transactions on Image Processing 24*, 12 (2015), 5369–5378.

[172] MOEMEN, M. Y., ELGHAMRAWY, H., GIVIGI, S. N., AND NOURELDIN, A. 3-d reconstruction and measurement system based on multimobile robot machine vision. *IEEE Transactions on Instrumentation and Measurement 70* (2020), 1–9.

[173] MOHIDEEN, S. K., PERUMAL, S. A., AND SATHIK, M. M. Image de-noising using discrete wavelet transform. *International Journal of Computer Science and Network Security 8*, 1 (2008), 213–216.

[174] MUR-ARTAL, R., MONTIEL, J. M. M., AND TARDOS, J. D. Orbslam: a versatile and accurate monocular slam system. *IEEE transactions on robotics 31*, 5 (2015), 1147–1163.

[175] MYRONENKO, A., AND SONG, X. Point set registration: Coherent point drift. *IEEE transactions on pattern analysis and machine intelligence 32*, 12 (2010), 2262–2275.

[176] NGUYEN, T., PHAM, M., NGUYEN, T., NGUYEN, K., OSHER, S., AND HO, N. Fourierformer: Transformer meets generalized fourier integral theorem. *Advances in Neural Information Processing Systems 35* (2022), 29319–29335.

[177] NISTÉR, D. An efficient solution to the five-point relative pose problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence 26*, 6 (2004), 756–770.

[178] NWOYE, C. I., MUTTER, D., MARESCAUX, J., AND PADOY, N. Weakly supervised convolutional lstm approach for tool tracking in laparoscopic videos. *International journal of computer assisted radiology and surgery 14*, 6 (2019), 1059–1067.

[179] OECHSLE, M., PENG, S., AND GEIGER, A. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 5589–5599.

[180] ÖZYEŞIL, O., VORONINSKI, V., BASRI, R., AND SINGER, A. A survey of structure from motion*. *Acta Numerica 26* (2017), 305–364.

[181] PAN, L., BARÁTH, D., POLLEFEYS, M., AND SCHÖNBERGER, J. L. Global structure-from-motion revisited. In *European Conference on Computer Vision* (2024), Springer, pp. 58–77.

[182] PARK, J., ZHOU, Q.-Y., AND KOLTUN, V. Colored point cloud registration revisited. In *Proceedings of the IEEE international conference on computer vision* (2017), pp. 143–152.

[183] PARK, J. J., FLORENCE, P., STRAUB, J., NEWCOMBE, R., AND LOVEGROVE, S. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 165–174.

[184] PARK, S.-Y., AND SUBBARAO, M. An accurate and fast point-to-plane registration technique. *Pattern Recognition Letters 24*, 16 (2003), 2967–2976.

[185] PENG, S., NIEMEYER, M., MESCHEDER, L., POLLEFEYS, M., AND GEIGER, A. Convolutional occupancy networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16* (2020), Springer, pp. 523–540.

[186] PENTLAND, A. P. A new sense for depth of field. *IEEE transactions on pattern analysis and machine intelligence* (1987), 523–531.

[187] PFISTER, H., ZWICKER, M., BAAR, J., AND GROSS, M. Surfels: Surface elements as rendering primitives. *Proceedings of the ACM SIGGRAPH Conference on Computer Graphics* (05 2000).

[188] PINTO, B., AND ANURENJAN, P. R. Video stabilization using speeded up robust features. In *2011 International Conference on Communications and Signal Processing* (2011), IEEE.

[189] PINTORE, G., AGUS, M., ALMANSA, E., SCHNEIDER, J., AND GOBBETTI, E. Slicenet: deep dense depth estimation from a single indoor panorama using a slice-based representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 11536–11545.

[190] PRATT, H., WILLIAMS, B., COENEN, F., AND ZHENG, Y. Fcnn: Fourier convolutional neural networks. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2017, Skopje, Macedonia, September 18–22, 2017, Proceedings, Part I 17* (2017), Springer, pp. 786–798.

[191] Qi, C. R., Su, H., Mo, K., and Guibas, L. J. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 652–660.

[192] Qin, T., Li, P., and Shen, S. Vins-mono: A robust and versatile monocular visual-inertial state estimator. *IEEE Transactions on Robotics 34*, 4 (2018), 1004–1020.

[193] Qin, Z., Yu, H., Wang, C., Guo, Y., Peng, Y., Ilic, S., Hu, D., and Xu, K. Geotransformer: Fast and robust point cloud registration with geometric transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence 45*, 8 (2023), 9806–9821.

[194] Qin, Z., Yu, H., Wang, C., Guo, Y., Peng, Y., and Xu, K. Geometric transformer for fast and robust point cloud registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 11143–11152.

[195] Qiu, J., Jiang, P.-T., Zhu, Y., Yin, Z.-X., Cheng, M.-M., and Ren, B. Looking through the glass: Neural surface reconstruction against high specular reflections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 20823–20833.

[196] Ramamonjisoa, M., Firman, M., Watson, J., Lepetit, V., and Turmukhambetov, D. Single image depth prediction with wavelet decomposition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2021), pp. 11089–11098.

[197] Ranftl, R., Bochkovskiy, A., and Koltun, V. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 12179–12188.

[198] Ranftl, R., Lasinger, K., Hafner, D., Schindler, K., and Koltun, V. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2020).

[199] Rawat, P., and Singhai, J. Review of motion estimation and video stabilization techniques for hand held mobile video. *Signal & Image Processing: An International Journal (SIPIJ) 2* (2011).

[200] Rehder, J., Nikolic, J., Schneider, T., Hinzmann, T., and Siegwart, R. Extending kalibr: Calibrating the extrinsics of multiple imus and of individual axes. In *2016 IEEE International Conference on Robotics and Automation (ICRA)* (2016), IEEE, pp. 4304–4311.

[201] Ren, X., Turkulainen, M., Wang, J., Seiskari, O., Melekhov, I., Kannala, J., and Rahtu, E. Ags-mesh: Adaptive gaussian splatting and meshing with geometric priors for indoor room reconstruction using smartphones. In *International Conference on 3D Vision (3DV)* (2025).

[202] Rippel, O., and Bourdev, L. Real-time adaptive image compression. In *International Conference on Machine Learning* (2017), PMLR, pp. 2922–2930.

[203] Ruan, L., Chen, B., Li, J., and Lam, M.-L. Aifnet: All-in-focus image restoration network using a light field-based dataset. *IEEE Transactions on Computational Imaging 7* (2021), 675–688.

[204] Rusu, R. B., Blodow, N., and Beetz, M. Fast point feature histograms (fpfh) for 3d registration. In *2009 IEEE international conference on robotics and automation* (2009), IEEE, pp. 3212–3217.

[205] Sadat, S., Buhmann, J., Bradley, D., Hilliges, O., and Weber, R. M. Litevae: Lightweight and efficient variational autoencoders for latent diffusion models. *arXiv preprint arXiv:2405.14477* (2024).

[206] Sarode, V., Li, X., Goforth, H., Aoki, Y., Srivatsan, R. A., Lucey, S., and Choset, H. Pcrnet: Point cloud registration network using pointnet encoding. *arXiv preprint arXiv:1908.07906* (2019).

[207] Satorras, V. G., Hoogeboom, E., and Welling, M. E(n) equivariant graph neural networks, 2021.

[208] Sattler, T., Zhou, Q., Pollefeys, M., and Leal-Taixe, L. Understanding the limitations of cnn-based absolute camera pose regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 3302–3312.

[209] Schelkens, P., Munteanu, A., Barbarien, J., Galca, M., Giro-Nieto, X., and Cornelis, J. Wavelet coding of volumetric medical datasets. *IEEE Transactions on medical Imaging 22*, 3 (2003), 441–458.

167

[210] Schönberger, J. L., and Frahm, J.-M. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)* (2016).

[211] Schonberger, J. L., and Frahm, J.-M. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 4104–4113.

[212] Schönberger, J. L., Zheng, E., Pollefeys, M., and Frahm, J.-M. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)* (2016).

[213] Schöps, T., Sattler, T., and Pollefeys, M. Surfelmeshing: Online surfel-based mesh reconstruction. *IEEE transactions on pattern analysis and machine intelligence 42*, 10 (2019), 2494–2507.

[214] Segal, A., Haehnel, D., and Thrun, S. Generalized-icp. In *Robotics: science and systems* (2009), vol. 2, Seattle, WA, p. 435.

[215] Serafin, J., and Grisetti, G. Nicp: Dense normal based point cloud registration. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2015), IEEE, pp. 742–749.

[216] Shabanov, A., Govindarajan, S., Reading, C., Goli, L., Rebain, D., Yi, K. M., and Tagliasacchi, A. Banf: Band-limited neural fields for levels of detail reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 20571–20580.

[217] Shabayek, A. E. R., et al. Vision-based uav attitude estimation: Progress and insights. *Journal of Intelligent & Robotic Systems 65*, 1 (2012), 295–308.

[218] Shavit, Y., Ferens, R., and Keller, Y. Learning single and multi-scene camera pose regression with transformer encoders. *Computer Vision and Image Understanding 243* (2024), 103982.

[219] Shavit, Y., and Keller, Y. Camera pose auto-encoders for improving pose regression. In *European Conference on Computer Vision* (2022), Springer, pp. 140–157.

[220] Shen, K., and Delp, E. J. Wavelet based rate scalable video compression. *IEEE transactions on circuits and systems for video technology 9*, 1 (1999), 109–122.

[221] SHEN, S. Accurate multiple view 3d reconstruction using patch-based stereo for large-scale scenes. *IEEE transactions on image processing 22*, 5 (2013), 1901–1914.

[222] SHIM, J., KANG, C., AND JOO, K. Diffusion-based signed distance fields for 3d shape generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 20887–20897.

[223] SHOTTON, J., GLOCKER, B., ZACH, C., IZADI, S., CRIMINISI, A., AND FITZGIBBON, A. Scene coordinate regression forests for camera relocalization in rgb-d images. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2013), pp. 2930–2937.

[224] SI, H., ZHAO, B., WANG, D., GAO, Y., CHEN, M., WANG, Z., AND LI, X. Fully self-supervised depth estimation from defocus clue. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 9140–9149.

[225] SIDDIQUI, Y., THIES, J., MA, F., SHAN, Q., NIESSNER, M., AND DAI, A. Texturify: Generating textures on 3d shape surfaces. In *European Conference on Computer Vision* (2022), Springer, pp. 72–88.

[226] SILBERMAN, N., HOIEM, D., KOHLI, P., AND FERGUS, R. Indoor segmentation and support inference from rgbd images. In *European conference on computer vision* (2012), Springer, pp. 746–760.

[227] SITZMANN, V., MARTEL, J., BERGMAN, A., LINDELL, D., AND WETZSTEIN, G. Implicit neural representations with periodic activation functions. *Advances in neural information processing systems 33* (2020), 7462–7473.

[228] SMOLYANSKIY, N., KAMENEV, A., AND BIRCHFIELD, S. On the importance of stereo for accurate depth estimation: An efficient semi-supervised deep neural network approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops* (2018), pp. 1007–1015.

[229] SOSNOVIK, I., SZMAJA, M., AND SMEULDERS, A. Scale-tooltool steerable networks. *arXiv preprint arXiv:1910.11093* (2019).

[230] STÜCKLER, J., AND BEHNKE, S. Multi-resolution surfel maps for efficient dense 3d modeling and tracking. *Journal of Visual Communication and Image Representation 25*, 1 (2014), 137–147.

[231] SURH, J., JEON, H.-G., PARK, Y., IM, S., HA, H., AND SO KWEON, I. Noise robust depth from focus using a ring difference filter. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017), pp. 6328–6337.

[232] SUWAJANAKORN, S., HERNANDEZ, C., AND SEITZ, S. M. Depth from focus with your mobile phone. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015), pp. 3497–3506.

[233] TANG, H., COHEN, S., PRICE, B., SCHILLER, S., AND KUTULAKOS, K. N. Depth from defocus in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 2740–2748.

[234] TANG, J., MARKHASIN, L., WANG, B., THIES, J., AND NIESSNER, M. Neural shape deformation priors. *Advances in Neural Information Processing Systems 35* (2022), 17117–17132.

[235] THOMAS, N., SMIDT, T., KEARNES, S., YANG, L., LI, L., KOHLHOFF, K., AND RILEY, P. Tensor field networks: Rotation- and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219* (2018).

[236] THRUN, S. Simultaneous localization and mapping. In *Robotics and cognitive approaches to spatial mapping*. Springer, 2008, pp. 13–41.

[237] TONG, J., MUTHU, S., MAKEN, F. A., NGUYEN, C., AND LI, H. Seeing through the glass: Neural 3d reconstruction of object inside a transparent container. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 12555–12564.

[238] TURKULAINEN, M., REN, X., MELEKHOV, I., SEISKARI, O., RAHTU, E., AND KANNALA, J. Dn-splatter: Depth and normal priors for gaussian splatting and meshing. *arXiv preprint arXiv:2403.17822* (2024).

[239] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, Ł., AND POLOSUKHIN, I. Attention

170

is all you need. *Advances in neural information processing systems 30* (2017).

[240] VIZZO, I., GUADAGNINO, T., MERSCH, B., WIESMANN, L., BEHLEY, J., AND STACHNISS, C. Kiss-icp: In defense of point-to-point icp–simple, accurate, and robust registration if done the right way. *IEEE Robotics and Automation Letters 8*, 2 (2023), 1029–1036.

[241] WALHA, A., WALI, A., AND ALIMI, A. M. Video stabilization for aerial video surveillance. *Aasri Procedia 4* (2013), 72–77.

[242] WALHA, A., WALI, A., AND ALIMI, A. M. Video stabilization with moving object detecting and tracking for aerial video surveillance. *Multimedia Tools and Applications 74* (2015), 6745–6767.

[243] WANG, H., LIU, Y., HU, Q., WANG, B., CHEN, J., DONG, Z., GUO, Y., WANG, W., AND YANG, B. Roreg: Pairwise point cloud registration with oriented descriptors and local rotations. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).

[244] WANG, J., GONG, Z., TAO, B., AND YIN, Z. A 3-d reconstruction method for large freeform surfaces based on mobile robotic measurement and global optimization. *IEEE Transactions on Instrumentation and Measurement 71* (2022), 1–9.

[245] WANG, K., GAO, F., AND SHEN, S. Real-time scalable dense surfel mapping. In *2019 International conference on robotics and automation (ICRA)* (2019), IEEE, pp. 6919–6925.

[246] WANG, N.-H., WANG, R., LIU, Y.-L., HUANG, Y.-H., CHANG, Y.-L., CHEN, C.-P., AND JOU, K. Bridging unsupervised and supervised depth from focus via all-in-focus supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 12621–12631.

[247] WANG, P., LIU, L., LIU, Y., THEOBALT, C., KOMURA, T., AND WANG, W. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689* (2021).

[248] WANG, Y., AND SOLOMON, J. M. Deep closest point: Learning representations for point cloud registration. In *Proceedings of the IEEE/CVF international conference on computer vision* (2019), pp. 3523–3532.

[249] WANG, Y., ZHOU, P., GENG, G., AN, L., AND ZHOU, M. Enhancing point cloud registration with transformer: cultural heritage protection of the terracotta warriors. *Heritage Science 12*, 1 (2024), 314.

[250] WANG, Z., CHEN, J., AND FURUKAWA, Y. Puzzlefusion++: Auto-agglomerative 3d fracture assembly by denoise and verify. *arXiv preprint arXiv:2406.00259* (2024).

[251] WEILER, M., AND CESA, G. General E(2)-Equivariant Steerable CNNs. In *Conference on Neural Information Processing Systems (NeurIPS)* (2019).

[252] WEILER, M., GEIGER, M., WELLING, M., BOOMSMA, W., AND COHEN, T. S. 3d steerable cnns: Learning rotationally equivariant features in volumetric data. *Advances in Neural Information Processing Systems 31* (2018).

[253] WEILER, M., HAMPRECHT, F. A., AND STORATH, M. Learning steerable filters for rotation equivariant cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018), pp. 849–858.

[254] WHELAN, T., GOESELE, M., LOVEGROVE, S. J., STRAUB, J., GREEN, S., SZELISKI, R., BUTTERFIELD, S., VERMA, S., NEWCOMBE, R. A., GOESELE, M., ET AL. Reconstructing scenes with mirror and glass surfaces. *ACM Trans. Graph. 37*, 4 (2018), 102.

[255] WONG, L. H. K., KANG, X., BAI, K., AND ZHANG, J. A survey of robotic navigation and manipulation with physics simulators in the era of embodied ai. *arXiv preprint arXiv:2505.01458* (2025).

[256] XIANG, Y., SCHMIDT, T., NARAYANAN, V., AND FOX, D. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *arXiv preprint arXiv:1711.00199* (2017).

[257] XIAO, T., SINGH, M., MINTUN, E., DARRELL, T., DOLLÁR, P., AND GIRSHICK, R. Early convolutions help transformers see better. *Advances in Neural Information Processing Systems 34* (2021), 30392–30400.

[258] XIE, J., GIRSHICK, R., AND FARHADI, A. Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks. In *European conference on computer vision* (2016), Springer, pp. 842–857.

[259] XIE, S., ET AL. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017).

[260] XIONG, Y., AND SHAFER, S. A. Depth from focusing and defocusing. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition* (1993), IEEE, pp. 68–73.

[261] XU, N., YANG, L., FAN, Y., YANG, J., YUE, D., LIANG, Y., PRICE, B., COHEN, S., AND HUANG, T. Youtube-vos: Sequence-to-sequence video object segmentation. In *Proceedings of the European conference on computer vision (ECCV)* (2018), pp. 585–601.

[262] XU, Q., WANG, W., CEYLAN, D., MECH, R., AND NEUMANN, U. Disn: Deep implicit surface network for high-quality single-view 3d reconstruction. *Advances in neural information processing systems 32* (2019).

[263] XUEYANG, K., XU, L., ZOU, Y., XU, H., AND MA, L. Simultaneous localization and mapping using cameras capturing multiple spectra of light, June 8 2023. US Patent App. 18/004,795.

[264] YANG, F., HUANG, X., AND ZHOU, Z. Deep depth from focus with differential focus volume. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 12642–12651.

[265] YANG, J., SCHONFELD, D., AND MOHAMED, M. Robust video stabilization based on particle filter tracking of projected camera motion. *IEEE Transactions on Circuits and Systems for Video Technology 19*, 7 (2009), 945–954.

[266] YANG, L., KANG, B., HUANG, Z., XU, X., FENG, J., AND ZHAO, H. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR* (2024).

[267] YANG, L., KANG, B., HUANG, Z., XU, X., FENG, J., AND ZHAO, H. Depth anything: Unleashing the power of large-scale unlabeled data. *arXiv preprint arXiv:2401.10891* (2024).

[268] YANG, L., KANG, B., HUANG, Z., ZHAO, Z., XU, X., FENG, J., AND ZHAO, H. Depth anything v2. *arXiv:2406.09414* (2024).

[269] YANG, N., WANG, R., STUCKLER, J., AND CREMERS, D. Deep virtual stereo odometry: Leveraging deep depth prediction for monocular

direct sparse odometry. In *Proceedings of the European conference on computer vision (ECCV)* (2018), pp. 817–833.

[270] YARIV, L., GU, J., KASTEN, Y., AND LIPMAN, Y. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems 34* (2021), 4805–4815.

[271] YARIV, L., HEDMAN, P., REISER, C., VERBIN, D., SRINIVASAN, P. P., SZELISKI, R., BARRON, J. T., AND MILDENHALL, B. Bakedsdf: Meshing neural sdfs for real-time view synthesis. In *ACM SIGGRAPH 2023 Conference Proceedings* (2023), pp. 1–9.

[272] YE, K., HOU, Q., AND ZHOU, K. 3d gaussian splatting with deferred reflection. In *ACM SIGGRAPH 2024 Conference Papers* (2024), pp. 1–10.

[273] YIN, W., LIU, Y., SHEN, C., AND YAN, Y. Enforcing geometric constraints of virtual normal for depth prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 5684–5693.

[274] YU, J., ET AL. Real-time selfie video stabilization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021).

[275] YU, J., AND RAMAMOORTHI, R. Learning video stabilization using optical flow. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020).

[276] YU, Y., ZHAN, F., LU, S., PAN, J., MA, F., XIE, X., AND MIAO, C. Wavefill: A wavelet-based generation network for image inpainting. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 14114–14123.

[277] YU, Z., CHEN, A., ANTIC, B., PENG, S., BHATTACHARYYA, A., NIEMEYER, M., TANG, S., SATTLER, T., AND GEIGER, A. Sdfstudio: A unified framework for surface reconstruction, 2022.

[278] YUAN, C., XU, W., LIU, X., HONG, X., AND ZHANG, F. Efficient and probabilistic adaptive voxel mapping for accurate online lidar odometry. *IEEE Robotics and Automation Letters 7*, 3 (2022), 8518–8525.

[279] ZENG, A., SONG, S., NIESSNER, M., FISHER, M., XIAO, J., AND FUNKHOUSER, T. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 1802–1811.

[280] ZENG, A., SONG, S., NIESSNER, M., FISHER, M., XIAO, J., AND FUNKHOUSER, T. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *CVPR* (2017).

[281] ZHANG, F., PRISACARIU, V., YANG, R., AND TORR, P. H. Ga-net: Guided aggregation net for end-to-end stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 185–194.

[282] ZHANG, X., YANG, J., ZHANG, S., AND ZHANG, Y. 3d registration with maximal cliques. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 17745–17754.

[283] ZHANG, Z. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence 22*, 11 (2000), 1330–1334.

[284] ZHENG, X.-Y., PAN, H., WANG, P.-S., TONG, X., LIU, Y., AND SHUM, H.-Y. Locally attentional sdf diffusion for controllable 3d shape generation. *ACM Transactions on Graphics (ToG) 42*, 4 (2023), 1–13.

[285] ZHOU, J., ZHANG, W., MA, B., SHI, K., LIU, Y.-S., AND HAN, Z. Udiff: Generating conditional unsigned distance fields with optimal wavelet diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 21496–21506.

[286] ZHOU, Q.-Y., PARK, J., AND KOLTUN, V. Fast global registration. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14* (2016), Springer, pp. 766–782.

[287] ZHU, M., GHAFFARI, M., CLARK, W. A., AND PENG, H. E2pn: Efficient se (3)-equivariant point network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 1223–1232.

[288] ZHU, Z., PENG, S., LARSSON, V., XU, W., BAO, H., CUI, Z., OSWALD, M. R., AND POLLEFEYS, M. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 12786–12796.

175