

Neural Fields in Robotics: A Survey

Muhammad Zubair Irshad¹, Mauro Comi², Yen-Chen Lin³, Nick Heppert⁴, Abhinav Valada⁴

Rares Ambrus¹, Zsolt Kira⁵, Jonathan Tremblay³

¹Toyota Research Institute, ²University of Bristol, ³Nvidia, ⁴University of Freiburg, ⁵Georgia Institute of Technology

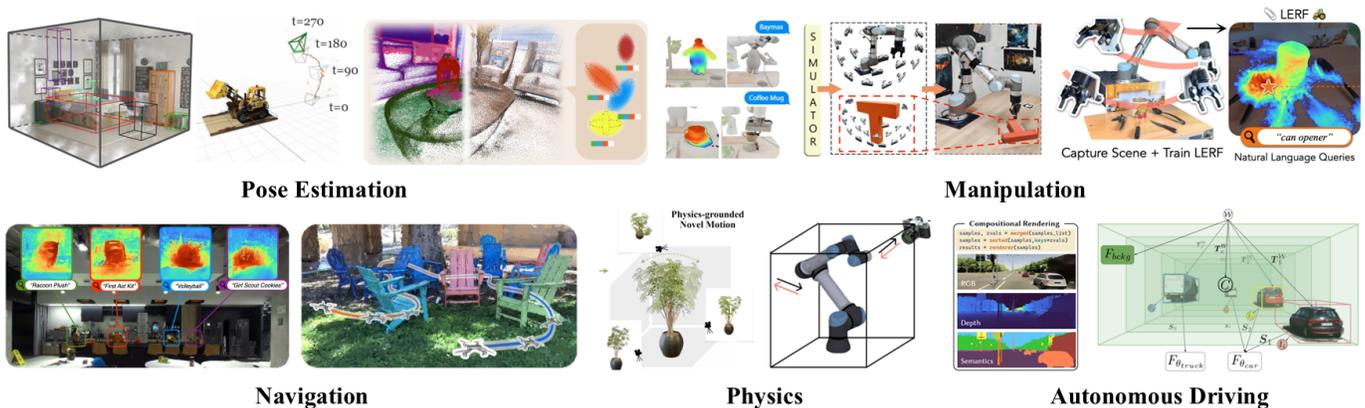


Fig. 1: **Overview:** This survey paper discusses a large variety of state-of-the-art Neural Field methods that enable robotics applications in pose estimation, manipulation, navigation, physics, and autonomous driving. Images adapted from [1–12].

Abstract—Neural Fields have emerged as a transformative approach for 3D scene representation in computer vision and robotics, enabling accurate inference of geometry, 3D semantics, and dynamics from posed 2D data. Leveraging differentiable rendering, Neural Fields encompass both continuous implicit and explicit neural representations enabling high-fidelity 3D reconstruction, integration of multi-modal sensor data, and generation of novel viewpoints. This survey explores their applications in robotics, emphasizing their potential to enhance perception, planning, and control. Their compactness, memory efficiency, and differentiability, along with seamless integration with foundation and generative models, make them ideal for real-time applications, improving robot adaptability and decision-making. This paper provides a thorough review of Neural Fields in robotics, categorizing applications across various domains and evaluating their strengths and limitations, based on over 200 papers. First, we present four key Neural Fields frameworks: Occupancy Networks, Signed Distance Fields, Neural Radiance Fields, and Gaussian Splatting. Second, we detail Neural Fields’ applications in five major robotics domains: pose estimation, manipulation, navigation, physics, and autonomous driving, highlighting key works and discussing takeaways and open challenges. Finally, we outline the current limitations of Neural Fields in robotics and propose promising directions for future research. Project page: [roboneerf.github.io](https://github.com/roboneerf)

Index Terms—Neural Radiance Field, NeRF, Neural Fields, Signed Distance Fields, 3D Gaussian Splatting, Occupancy Networks, Computer Vision, Novel View Synthesis, Neural Rendering, Volume Rendering, Pose Estimation, Robotics, Manipulation, Navigation, Autonomous Driving.

I. INTRODUCTION

Robots depend on precise and compact representations of their environment to perform a wide array of tasks, from navigating busy warehouses to organizing cluttered homes or assisting in high-stakes search-and-rescue missions. At the core of a typical robotic pipeline is the synergy between perception

and action. The perception system gathers sensory data from devices such as RGB cameras, LiDAR, and depth sensors and transforms them into a coherent model of the environment — such as a 3D map that enables the robot to maneuver through dynamic, obstacle-rich spaces. The quality of this representation directly impacts the robot’s decision-making or policy, which translates the perceived environment into actions, enabling it to avoid moving forklifts, pick up scattered objects, or plan a safe path in an emergency. Traditionally, robots have modeled their environments using data structures like point clouds [13–15], voxel grids [16], meshes [17–19], and Truncated Signed Distance Functions (TSDF) [20]. While these representations have advanced robotic capabilities, they struggle to capture fine geometric details, particularly in complex or dynamic environments, leading to suboptimal performance in adaptable scenarios.

To overcome these limitations, Neural Fields (NFs) [21] have emerged as a promising alternative, offering continuous, differentiable mappings from spatial coordinates to physical quantities like color or signed distance. Unlike traditional data structures, NFs can model 3D environments as continuous functions parameterized by neural networks or Gaussian distributions. This enables them to represent complex geometries and fine details more efficiently [22, 23]. NFs can be optimized using gradient-based methods with various types of real-world sensory data, including images and depth maps, to produce high-quality 3D reconstructions. In the realm of robotics, NFs provide several distinct advantages over traditional methods:

- **High-Quality 3D Reconstructions:** NFs generate detailed 3D representations of environments, which are crucial for tasks like navigation, manipulation, and scene understanding [24–28].
- **Multi-Sensor Integration:** NFs can seamlessly integrate data from multiple sensors, such as LiDAR and RGB

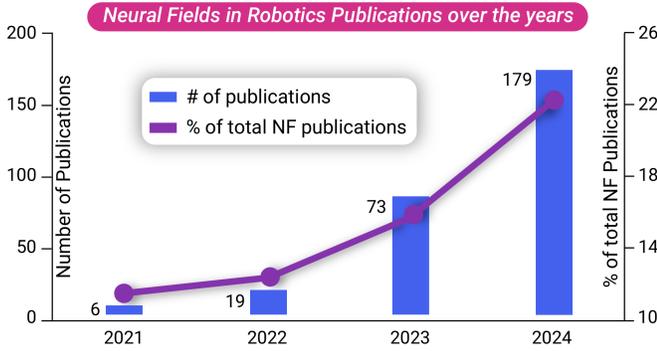


Fig. 2: **Growth of Neural Fields in Robotics:** plotted as a rough number of publications vs. % of total neural field publications per year.

cameras, providing a more robust and adaptable perception of the environment [29, 30].

- **Continuous and Compact Representations:** Unlike voxel grids or point clouds, which are inherently discrete, NFs offer continuous representations that capture fine spatial details using fewer parameters, enhancing computational efficiency [22, 31].
- **Generalization and Adaptation:** Once trained, NFs can generate novel viewpoints of a scene, even from previously unseen perspectives, which is particularly valuable for exploration or manipulation tasks. This ability is enabled by generalizable NeRF methods [32–34].
- **Integration with Foundation Models:** NFs can be combined with foundation models like CLIP [35] or DINO [36], enabling robots to interpret and respond to natural language queries or other semantic inputs [37, 38].

Recent advances in generative AI [39] have further expanded the capabilities of NFs by leveraging synthetic data as supervisory signals, thereby reducing reliance on real-world observations. This paradigm shift allows NFs to be optimized in scenarios where real-world data collection is impractical or costly. Importantly, it positions NFs as a key link between generative AI and robotics. While generative priors from 2D data are powerful, they often lack the spatial coherence needed for effective robotic decision-making. NFs integrate these priors with sparse real-world data [33], enabling them to model sensory and motor spaces in scenarios constrained by physical environments, such as limited sensor configurations and occlusions.

Given these advantages, the application of NFs in robotics is a rapidly growing area of research. Figs. 1 and 2 provide an overview of NF applications in robotics and highlight the rise in NF-related robotics publications over time. In this paper, we aim to structure and analyze their impact on the field. The manuscript is organized as follows: Sec. II covers the formulation of NFs, while Sec. III highlights their benefits across various domains, categorized into distinct themes:

- **Pose Estimation,** focuses on NF as a scene or object representation in camera pose estimation, object pose estimation, and Simultaneous Localization and Mapping (SLAM) (Sec. III-A).
- **Manipulation,** discusses how NFs’ accurate 3D

reconstruction assists robots in manipulating objects (Sec. III-B).

- **Navigation,** highlights the role of NFs in enhancing robotic navigation by enabling accurate and efficient perception of real-world environments (Sec. III-C).
- **Physics,** explores how NFs enable robots to reason about physical interactions to improve their understanding of real-world dynamics (Sec. III-D).
- **Autonomous Driving,** focuses on NFs’ role in building photorealistic simulators for the real world (Sec. III-E).

We conclude by discussing several research directions and challenges in Sec. IV. To the best of our knowledge, this survey represents one of the first comprehensive examinations of Neural Fields in the domain of robotics. We complement the closest concurrent survey [40] that focuses on NeRFs by covering a comprehensive set of fields, including 3DGS, Occupancy, Signed Distance Fields, and beyond. By integrating insights from various dimensions, this survey aims to provide a holistic understanding of the current state of NF in robotics applications, highlighting recent achievements, upcoming challenges, and unexplored areas within robotics.

II. FORMULATION OF NEURAL FIELDS

We begin by defining several types of fields that are key to the formulation of Neural Fields. In mathematics and physics, a field is a descriptor of a quantity with a value assigned to every point in space and time. Formally, this is expressed as:

- **Scalar fields:** A scalar field $f : \mathbb{R}^n \rightarrow \mathbb{R}$ assigns a scalar value (*e.g.* temperature, pressure) to every point in an n -dimensional space. Mathematically, it is defined as $f(x, y, z)$ for a 3D space, where (x, y, z) are the coordinates in space and f is the scalar value at that point. Examples of scalar fields are the Occupancy Fields (Sec. II-A) and Signed Distance Fields (Sec. II-B).
- **Vector fields:** A vector field $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is an extension of a scalar field that associates a vector (*e.g.*, velocity or force) with every point in space. For example, in three dimensions, a possible formulation is given by $F(x, y, z) = (F_0(x, y, z), \dots, F_{m-1}(x, y, z))$, where each element F_i represents the vector values corresponding to the input coordinate (x, y, z) . This function can either be represented as a neural network (Sec. II-C1), a volumetric grid [41], or the mix of both [25].

A. Occupancy Fields

Occupancy represents the binary state of whether a point \mathbf{p} in space is occupied by a surface S or not:

$$o(\mathbf{p}) = \begin{cases} 0 & \text{if } \mathbf{p} \text{ is inside } S, \\ 1 & \text{if } \mathbf{p} \text{ is outside } S. \end{cases}$$

This idea can be extended to the continuous case, where occupancy is represented by the probability $o(\mathbf{p}) \in [0, 1]$ indicating the likelihood of \mathbf{p} being inside or outside the surface. Occupancy Networks [42] leverage this idea by learning a continuous function, or NF, which maps points $\mathbf{p} \in \mathbb{R}^n$ to occupancy probabilities, conditioned on input observations \mathbf{x} (*e.g.*, point clouds or images). This function, $f_\theta(\mathbf{p} | \mathbf{x})$,

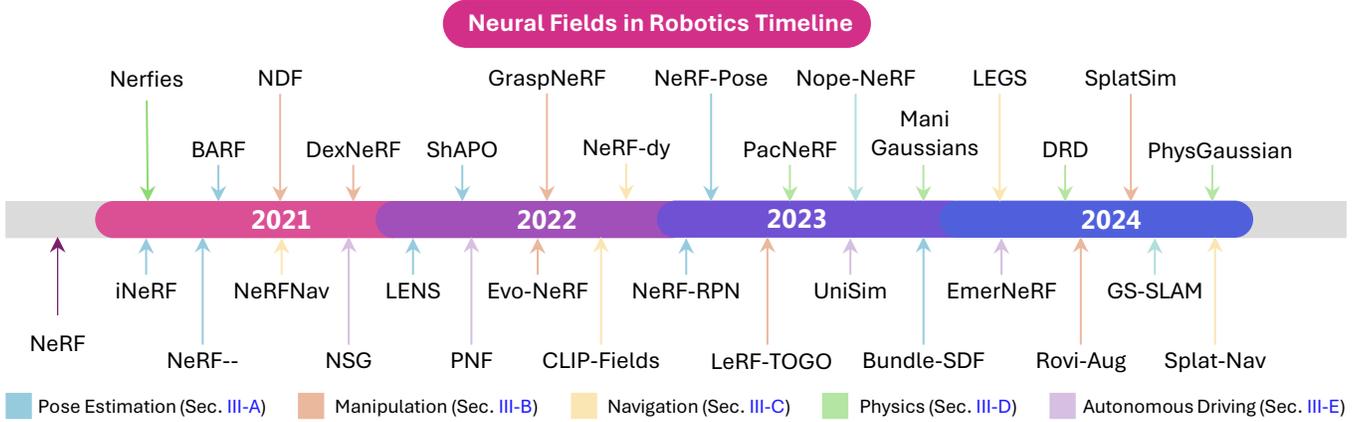


Fig. 3: **Timeline** of Neural Fields in Robotics paper showing key papers over the years divided into 5 major application areas.

is optimized during training by sampling points \mathbf{p}_i and minimizing the cross-entropy loss between predicted occupancy $f_\theta(\mathbf{p}_i | \mathbf{x})$ and ground truth occupancy $o(\mathbf{p}_i, \mathbf{x})$. This approach enables smooth, continuous surface representations, which can later be thresholded to recover discrete occupancy.

B. Signed Distance Fields (SDFs)

Signed Distance Fields assign each point in space the shortest distance to the surface boundary of an object, with the sign indicating whether the point is inside or outside. For a point \mathbf{p} and surface S , the SDF $d(\mathbf{p})$ is defined as:

$$d(\mathbf{p}) = \begin{cases} -\min_{\mathbf{q} \in S} \|\mathbf{p} - \mathbf{q}\| & \text{if } \mathbf{p} \text{ is inside } S, \\ \min_{\mathbf{q} \in S} \|\mathbf{p} - \mathbf{q}\| & \text{if } \mathbf{p} \text{ is outside } S, \end{cases}$$

where $\|\mathbf{p} - \mathbf{q}\|$ is the Euclidean distance between points \mathbf{p} and \mathbf{q} . When representing the SDF as an NF, during training, points \mathbf{p}_i are sampled, and their distance to the closest surface provides the supervision signal [23].

C. Radiance Fields

Radiance Fields represent the distribution of light in 3D space, associating each point with a radiance value (light intensity and color) in every direction. This can be described using a function $L : \mathbb{R}^3 \times S^2 \rightarrow \mathbb{R}^3$, where $L(p, \omega)$ gives the light radiance at point p in direction ω , with p in 3D space and ω on the unit sphere S^2 representing all possible directions.

1) *Neural Radiance Fields (NeRF)*: NeRFs [22] represent scenes as volumetric fields of density σ and RGB color \mathbf{c} using a neural network. The weights of NeRFs are optimized per scene using input RGB images, and their camera poses. After training, the density field captures scene geometry, while the color field models the view-dependent appearance. A multilayer perceptron (MLP) with parameters Θ predicts the density σ and RGB color \mathbf{c} for each point based on its 3D position $\mathbf{x} = (x, y, z)$ and viewing direction \mathbf{d} . To address the spectral bias of neural networks in low-dimensional spaces [43], positional encoding $\gamma(\cdot)$ is applied to inputs: $(\sigma, \mathbf{c}) \leftarrow F_\Theta(\gamma(\mathbf{x}), \gamma(\mathbf{d}))$.

To render a pixel, a camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is cast from the camera center \mathbf{o} in direction \mathbf{d} . K points $\{\mathbf{x}_k = \mathbf{r}(t_k)\}_{k=1}^K$ are sampled along the ray and passed through the MLP to generate densities and colors $\{\sigma_k, \mathbf{c}_k\}_{k=1}^K$. These are then combined to

estimate the pixel color $\hat{\mathbf{C}}(r)$ via volume rendering [44], using a numerical quadrature approximation [45]

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{k=1}^K T_k \left(1 - \exp(-\sigma_k(t_{k+1} - t_k)) \right) \mathbf{c}_k, \text{ with } (1)$$

where $T_k = \exp(-\sum_{k' < k} \sigma_{k'}(t_{k'+1} - t_{k'}))$ can be interpreted as the probability that the ray successfully transmits to point $\mathbf{r}(t_k)$. NeRF is trained by minimizing the photometric loss, $\mathcal{L}_{\text{photo}} = \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}(\mathbf{r})\|_2^2$, where $\mathbf{C}(\mathbf{r})$ is the observed RGB value for ray \mathbf{r} in a sampled set of rays \mathcal{R} .

While vanilla NeRF achieves photorealistic results, it is time-consuming to train and render from a pretrained NeRF. To reduce these costs, several improvements are proposed: a) using encodings with better speed/quality trade-offs [46], b) adopting smaller neural networks to reduce memory-bandwidth demands [41, 47], and c) skipping ray marching steps in empty space to cut down the computational cost of neural volume rendering [48]. These optimizations accelerate NeRF training and inference by several orders of magnitude, enabling real-time use in time-sensitive applications.

2) *3D Gaussian Splatting*: 3D Gaussian Splatting [49] represents a scene as a collection of anisotropic 3D Gaussians, each defined by a mean vector $\boldsymbol{\mu} \in \mathbb{R}^3$, covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^3$, color vector $\mathbf{c} \in \mathbb{R}^3$, and opacity scalar α . The influence of a Gaussian on a point $\mathbf{x} \in \mathbb{R}^3$ is given by:

$$f(\mathbf{x}; \boldsymbol{\Sigma}, \boldsymbol{\mu}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (2)$$

which quantifies each Gaussian's contribution based on spatial proximity. To render the scene onto a 2D plane, the 3D covariance $\boldsymbol{\Sigma}$ is projected into 2D as $\boldsymbol{\Sigma}' = \mathbf{J}\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T\mathbf{J}^T$ [50], where \mathbf{W} is the projective transformation and \mathbf{J} is the Jacobian. The 2D mean vector $\boldsymbol{\mu}'$ is obtained via perspective projection $\text{Proj}(\boldsymbol{\mu} | \mathbf{E}, \mathbf{K})$, using the camera's extrinsic \mathbf{E} and intrinsic \mathbf{K} matrices. To ensure valid optimization of 3D covariance $\boldsymbol{\Sigma}$, it is decomposed as $\boldsymbol{\Sigma} = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T$, where \mathbf{R} and \mathbf{S} are rotation and scaling matrices, respectively. Finally, the color of each pixel p in the image can be calculated as:

$$I(p) = \sum_{i=1}^N f(p; \boldsymbol{\mu}'_i, \boldsymbol{\Sigma}'_i) \alpha_i \mathbf{c}_i \prod_{j=1}^{i-1} 1 - f(p; \boldsymbol{\mu}'_j, \boldsymbol{\Sigma}'_j) \alpha_j. \quad (3)$$

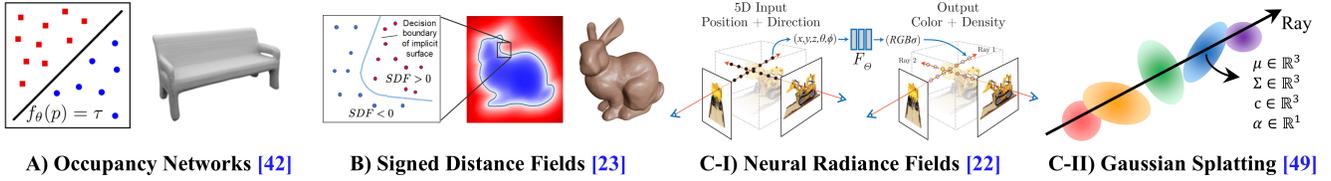


Fig. 4: **Neural Field Representations:** Section II discusses four core Neural Field representations — Occupancy Networks [42], Signed Distance Fields [23], Neural Radiance Fields [22], and 3D Gaussian Splatting [49].

Overall, scalar, vector, and radiance fields, comprising Occupancy networks, Signed Distance Fields, Neural Radiance Fields, and 3D Gaussian Splatting (as shown in Fig. 4) form the foundation of neural representations, enabling them to capture complex geometries and spatial relationships with far greater detail than traditional methods. In the next section, we explore how these mathematical tools unlock a variety of robotics applications, including pose estimation, manipulation, navigation, physical property inference, and autonomous driving in challenging environments.

III. NEURAL FIELDS FOR ROBOTICS

In this section, we delve into the application of Neural Fields across five major areas of robotics: pose estimation, manipulation, navigation, physics, and autonomous driving (see Figs. 3 and 5 for a timeline and taxonomy of selected key NFs in robotics papers). Each subsection below highlights key works within these domains, providing a comprehensive overview of the state-of-the-art methods. We conclude each subsection with a discussion of the key takeaways and the open challenges that remain in these areas, offering insights into the future directions of research in NFs for robotics.

A. Neural Fields for Pose Estimation

Neural Fields have transformed pose estimation by offering robust and efficient methods to estimate the position and orientation of cameras and objects in 3D scenes. This section explores two key areas: camera pose estimation (Sec. III-A1) and object pose estimation (Sec. III-A2). Camera pose estimation focuses on determining the viewpoint of cameras, which is crucial for tasks like mapping and reconstruction. Alternatively, object pose estimation involves localizing and orienting objects within a scene, essential for applications like manipulation and interaction. NFs optimize these tasks through gradient-based techniques, either by refining scene representations or directly providing reliable features for pose estimation.

1) *Camera Pose Estimation:* As discussed in Sec. I, NFs are differentiable, allowing gradient updates through scene representations, such as NeRF’s volumetric rendering, down to camera parameters. While differentiable rendering has been applied to meshes [51, 52], we focus here on methods applicable to NeRF-like models. This section starts by examining techniques that rely on pre-optimized NeRFs. We then explore approaches that tackle simultaneous pose estimation and geometry reconstruction, concluding with a discussion on their impact on Simultaneous Localization and Mapping (SLAM).

Pose Estimation via Optimized Neural Fields: For localization, iNeRF [2] inverts an already optimized NeRF for the task of pose estimation. Starting with an image, iNeRF finds the translation and rotation of a camera relative to a pretrained NeRF by using gradient descent to minimize the residual between pixels rendered from an optimized NeRF. Parallel iNeRF [53] parallelized the optimization processes of 6DoF poses based on fast pretrained NeRFs. Lens [54] prevented the generation of novel views in irrelevant areas by choosing virtual camera positions based on the NeRF’s internal 3D scene geometry. The rendered images were then used as synthetic data to efficiently train a camera pose regression model.

Simultaneous Pose Estimation and Reconstruction: NeRF—[55] and BARF [96] show that given RGB observations of a scene, camera poses and NeRFs can be jointly optimized, removing the need for classical Structure-from-Motion (SfM) pipelines. The former initializes cameras to the origin for forward-facing scenes, while the latter employs a coarse-to-fine reconstruction scheme that gradually introduces higher frequency position encodings [97, 98]. This coarse-to-fine approach can also be adapted to multi-resolution grids like NGP [25] by applying a weighted schedule across resolution levels [99, 100].

LocalRF [56] reconstructs long trajectories incrementally by adding images sequentially, using a subdivision approach similar to BlockNeRF [101], without relying on structure-from-motion (SfM). Notably, LocalRF uses different learning rates for translation and orientation parameters, highlighting the challenges in taming camera-optimizing NeRFs. GNeRF [102] proposes a pose-conditioned GAN to recover a NeRF. Others have explored more suitable techniques for pose optimization, such as Gaussian [103] or sinusoidal activations [104]. NoPeNeRF [105] uses monocular depth priors to constrain the scene as well as relative pose estimates. Keypoint matches or dense correspondences can also be used to constrain the relative pose estimates using ray-to-ray correspondence losses [106, 107]. DBARF [108] proposes using low-frequency feature maps to guide the bundle adjustment for generalizable NeRFs [32, 93, 109]. Nerfels [110] combines invertible neural rendering with traditional keypoint-based camera pose optimization.

Various works that apply NFs for localization and pose estimation utilize NeRF’s internal features to establish 2D-3D correspondences [111, 112], remove the need for an initial pose estimate [113], augment the training set of the pose regressor with a few-shot NeRF [114], or apply a decoupled representation of pose along with an edge-based sampling strategy to enhance the learning signal [115]. They also address dynamic scenes by integrating geometric motion and

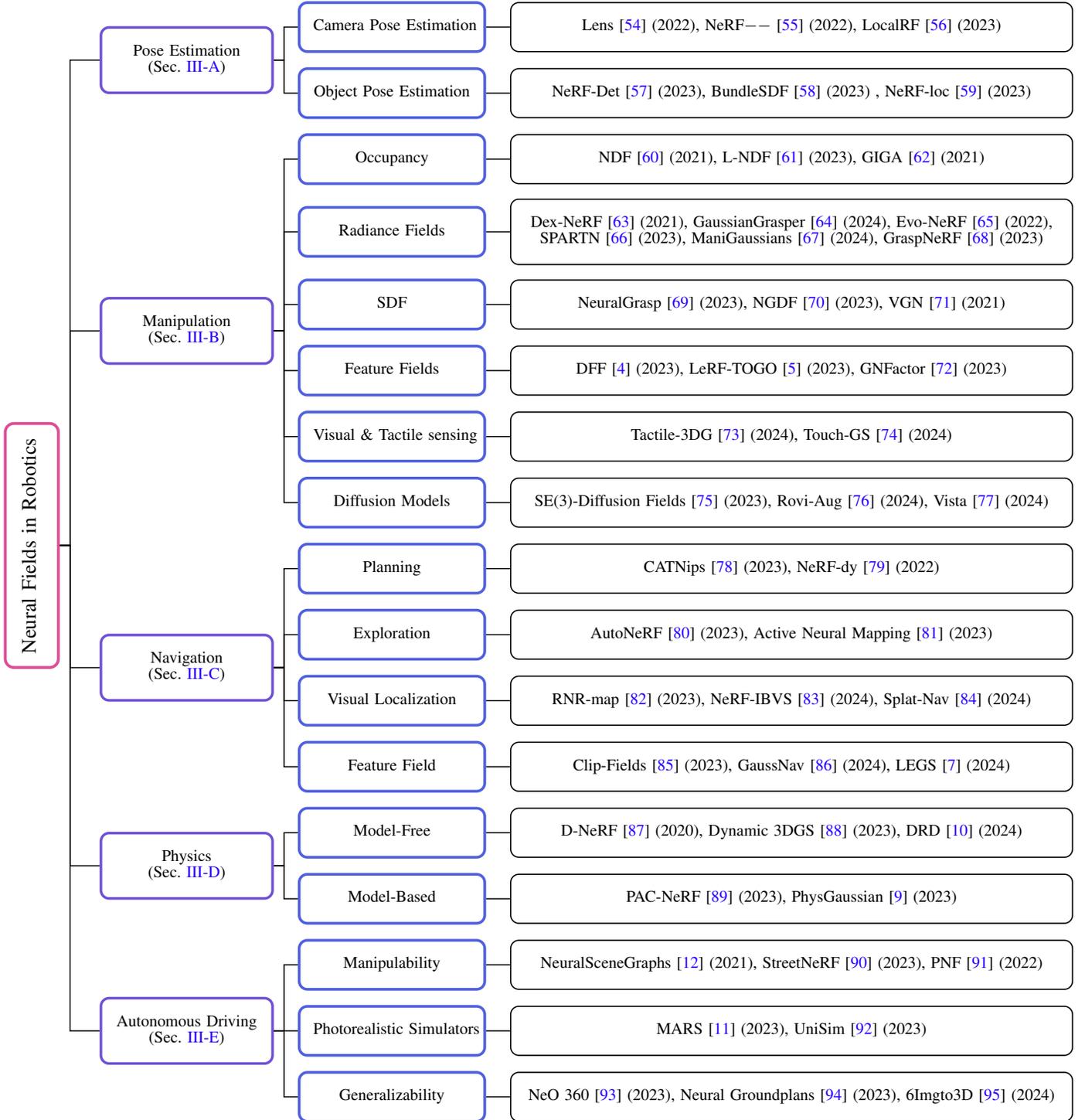


Fig. 5: Taxonomy of selected key Neural Fields papers in five major robotics application areas.

segmentation for initial pose estimation, combined with static ray sampling to speed up view synthesis [116].

Simultaneous Localization and Mapping (SLAM): Jointly optimizing camera poses and neural scene representations is a fundamental aspect of the SLAM (Simultaneous Localization and Mapping) problem. Recent advancements in NFs have reshaped SLAM by leveraging their qualities, such as the ability to model continuous surfaces, lower memory usage, and enhanced robustness against noise and outliers. For instance, iMap [117] utilizes a single MLP to predict radiance fields

as the mapping representation, employing parallel tracking and mapping threads where the tracking thread optimizes the pose of the input RGB-D frame, and the mapping thread refines both the MLP and camera pose for selected keyframes. NICE-SLAM [118] (see Fig. 7) further enhances this approach by replacing the single MLP with hierarchical feature grids, resulting in faster inference and more accurate reconstructions. Gaussian splatting-based SLAM systems exploit the advantages of 3DGS including faster runtime and improved photorealistic rendering to further enhance the performance [119–121] (see

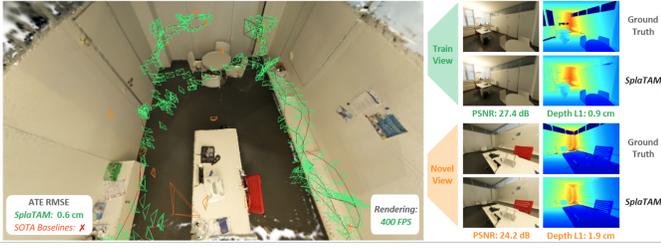


Fig. 6: Mapping and tracking results from Splatam [121].

Fig. 6). Besides better reconstruction quality, NF-based methods also provide an easier way to store various semantic information. Various semantic SLAM systems [3, 122–124] use NFs as unified representations to represent diverse information of the environment. For a more thorough survey on NFs for SLAM, we refer readers to Tosi *et al.* [125].

2) *Object Pose Estimation*: NFs have also been employed for localizing and orienting objects within a scene. Accurate pose estimation is crucial for robotics, as it enhances a robot’s ability to interact with its surroundings and perform tasks such as manipulation, navigation, and object recognition. Works in this domain use NFs’ features to establish correspondences or directly regress poses and reconstruct shapes. This facilitates the determination of bounding boxes and orientations for various objects in a 3D environment.

Neural Field features are also shown to be effective for multi-view 3D bounding box estimation of objects in the scene. NeRF-RPN [1] estimates 3D object boxes directly on NeRF’s feature grid, using a novel voxel representation and without re-rendering from a pretrained NeRF. It can be trained end-to-end to estimate high-quality 3D bounding boxes without class labels. Similarly, NeRF-Det [57] (see Fig. 9) leverages a NeRF to explicitly estimate 3D geometry and improve detection performance. It introduces geometry priors and connects detection with NeRF branches through a shared MLP, enhancing generalizability and efficiency without per-scene optimization. NeRF-RPN’s performance can be further enhanced with self-supervised representation learning directly using NeRF grids, as shown by NeRF-MAE [126]. Similarly, Gaussian splats [49] have also been effectively utilized for 3D object detection, as demonstrated by GaussianDet [127] and 3D-GSDet [128]. Furthermore, NeRF-loc [59] employs a transformer-based framework to extract labeled, oriented 3D bounding boxes of objects from NeRF scenes. It utilizes a pair of parallel transformer encoder branches to encode both the context and details of target objects, fusing these features with attention layers for accurate object localization, outperforming conventional RGB(-D) based methods. NeRF-pose [129] employs a weakly supervised 6D pose estimation pipeline that requires only 2D segmentation and known relative camera poses during training, avoiding the need for precise 6D pose annotations. It reconstructs objects from multiple views, then trains a pose regression network to predict 2D-3D correspondences with a NeRF-enabled Perspective-n-Point (PnP)+RANSAC algorithm estimating stable and accurate poses from a single input image.

Other direct pose-regression methods also simultaneously reconstruct object shapes in conjunction with object pose estimation. ShAPO [130], FSD [131] and CARTO [132] jointly

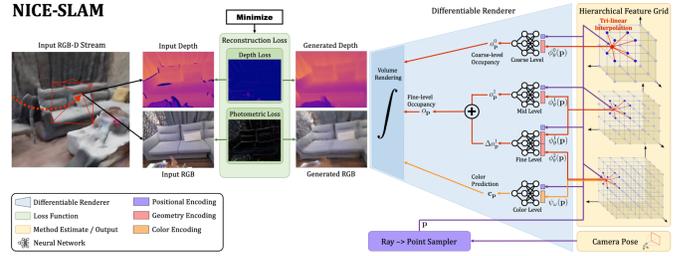


Fig. 7: Network architecture of Nice-SLAM [118].

reconstruct object shapes and regress their 6D object poses using a single-shot pipeline employing implicit representations and disentangled shape and appearance priors. UPNeRF [133] proposes a unified framework for monocular 3D reconstruction that integrates pose estimation with NeRF-based reconstruction, addressing the shortcomings of existing methods that rely on external 3D object detectors for initial poses. It decouples dimension estimation and pose refinement to resolve scale-depth ambiguity, and employs a projected-box representation for cross-domain generalization. NeRF-from-image [134] integrates NeRF with GANs to model arbitrary topologies without requiring accurate ground-truth poses or multiple views during training. It uses an unconditional 3D-aware generator and a hybrid inversion scheme to recover an SDF-parameterized 3D shape, pose, and appearance, refining initial solutions via optimization. NCF [135] estimates the 6DoF pose of a rigid object with a 3D model from a single RGB image by predicting 3D object coordinates at 3D query points sampled in the camera frustum rather than at image pixels. Bundle-SDF [58] (see Fig. 8) tracks the 6-DoF pose of an unknown object from a monocular RGBD video sequence while simultaneously performing neural 3D reconstruction. It handles arbitrary rigid objects with minimal visual texture, requiring only the object’s segmentation in the first frame. The approach uses a Neural Object Field learned alongside pose graph optimization to build a consistent 3D representation of the object’s geometry and appearance.

3) *Takeaways and Open Challenges in Neural Fields for Pose Estimation*: Despite the promising progress in NFs for pose estimation, several open challenges remain. Current methods prove effective in real-time pose estimation of cameras as well as objects. While significant progress has been made in pose estimation for static scenes, there is still room for further exploration in dynamic environments. Future work could focus on refining methods for recovering camera poses from dynamic video capture, where both cameras and objects exhibit significant movement, such as post-hoc calibration and labeling of robotic datasets [136]. This would open up the possibility of learning 3D priors from large-scale monocular videos. Another avenue for future work could explore using NFs for open-vocabulary 6D object pose estimation and solving the canonicalization problem of large-scale datasets [137].

B. Neural Fields for Robotic Manipulation

One of the key challenges in robotic manipulation is obtaining a precise geometric representation of both the objects and the environment involved in the task. An effective

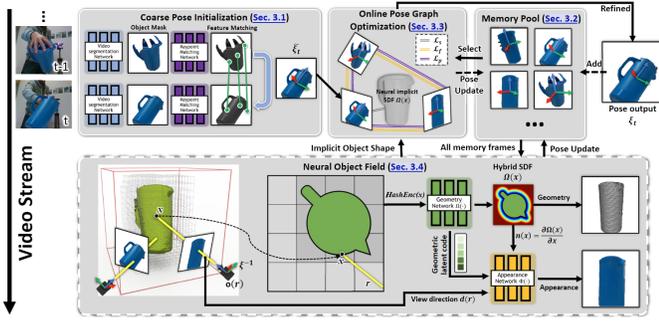


Fig. 8: BundleSDF [58] for object tracking and reconstruction.

representation must also capture the environment dynamics, offering a robust 3D understanding of the objects. In this section, we explore the application of NFs in control tasks for manipulation, with a focus on pick-and-place scenarios within 3 and 6 Degrees of Freedom (DoF). A summary of methods leveraging NFs for manipulation is provided in Table I.

Approaches that synthesize 3-DoF grasps use visual input, like RGB or depth images, to generate grasps in the image frame [71]. This means the end effector can translate horizontally and rotate around its vertical axis, with depth sensors determining vertical positioning. However, 3-DoF grasping lacks precise orientation control. In contrast, 6-DoF methods predict both position and orientation using 3D representations, enabling full control over roll, pitch, and yaw, which enhances the robot’s ability to manipulate objects in any direction.

Within the scope of 3-DoF, Dex-NeRF [63] (see Fig. 11) uses NFs to detect and infer the geometry of transparent objects, employing a transparency-aware rendering technique and additional lighting for specular reflections. Combined with Dex-Net [138], it generates 3-DoF grasping poses for transparent objects in both simulated and real environments. More recently, Evo-NeRF [65] extends this by leveraging Instant-NGP [25] to accelerate inference and adapt NeRF weights for sequential grasping tasks on transparent objects, updating the representation with each grasp.

Recently, there has been a shift toward using NFs for 6-DoF grasp pose estimation, offering an alternative to traditional point cloud-based methods. Below, we discuss these representations in detail:

1) **Occupancy Fields:** Neural Descriptor Fields (NDF) [60] propose an $SE(3)$ -equivariant object representation for manipulating novel objects in arbitrary poses, with few demonstrations. Using a Vector Neurons (VN) [146] network, 6-DoF relative poses between objects and grippers are encoded. NDF represents objects as continuous 3D descriptor fields, mapping points to descriptor vectors that capture spatial relationships to object geometry. However, NDF struggles with generalizing to new object categories, a limitation addressed by Local Neural Descriptor Fields (L-NDF) [61], which use a voxel grid of local embeddings to better capture local geometry and descriptors for new shapes. Both NDF and L-NDF rely on VN equipped with an Occupancy Network (ONet). GIGA [62] leverages a Convolutional Occupancy Network (ConvONet) to detect 6-DoF grasps in cluttered environments from a single depth image. By encoding the Truncated Signed Distance Function (TSDF), GIGA jointly predicts volumetric occupancy and 6-

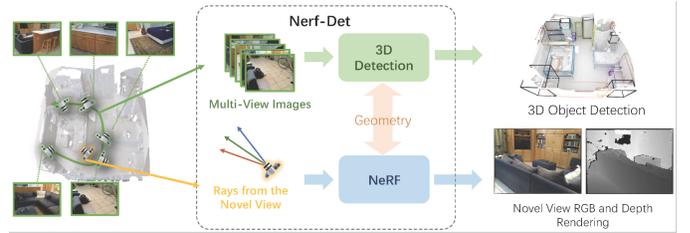


Fig. 9: NeRF-Det’s [57] 3D detection pipeline using NeRFs.

DoF grasp detection, enabling it to detect grasps on occluded objects from partial observations.

2) **Radiance Fields:** These fields are primarily modeled using two approaches: NeRFs and 3D Gaussian Splatting.

NeRF: NeRF-Supervision [140] leverages NeRFs to generate synthetic data for dense correspondence estimation by treating correspondences as depth distributions instead of pixel-wise depth. This enables 6-DoF picking tasks on thin and reflective surfaces using only RGB images, though multi-view images are required to build the NeRF representation. MIRA [141] extends this by constructing a NeRF before each action, enabling pick-conditioned placing via view synthesis. MIRA also trains a NeRF model with perspective cameras for direct orthographic rendering, which aligns well with translationally equivariant architectures like ConvNets. SPARTN [66] enhances visual grasping policies by using synthetic multi-view RGB images from eye-in-hand camera setups, significantly improving success rates in grasping tasks over standard imitation learning methods. These advancements highlight NeRF’s potential to bridge the gap between RGB-based and depth-based robotic policies.

Blukis *et al.* [144, 145] proposed an approach that jointly optimizes 3D reconstruction and grasp pose estimation by encoding objects into a unified latent representation. Encoded latents are decoded for view synthesis, 3D reconstruction, and grasp proposals. NeRFs have also been adapted for transfer learning in 6-DoF grasp pose evaluation and optimization with MVNeRF [142], which processes inputs from multiple scenes, enabling a more generalized representation and faster perception-to-action mapping. RGBGrasp [143] advances real-time applications by integrating multi-view RGB data from an eye-on-hand camera and depth maps from a pre-trained model. It further accelerates 3D reconstruction using hash-encoding [25] and a novel sampling strategy.

3D Gaussian Splatting: The introduction of 3D Gaussian Splatting (Sec. II-C2) marks a promising advancement in leveraging 3D representations for real-time robotic manipulation. GaussianGrasper [64] proposes a novel approach to 6-DoF grasping using Gaussian Splatting for open-vocabulary object grasping, thus enabling robots to understand and execute tasks based on natural language instructions. Similarly, ManiGaussian [67] builds on dynamic 3D Gaussian Splatting to capture scene-level spatiotemporal dynamics, enhancing the robot’s capability to execute tasks conditioned on natural language instructions. SplatSim [147] shows an application of 3DGS for improving Sim2Real transfer for robotic manipulation policies that rely on RGB images. This is obtained by leveraging the photorealism of 3DGS, which reduces the domain shift between synthetic and real visual information.

Fields	Input	Method	Representation	Scope	
Occupancy	Point-cloud	NDF [60]	VNN, ONet	Scene-specific	
		L-NDF [61]	VNN, ONet	General	
Signed Distance	Single-view Depth	GIGA [62]	ConvONet, TSDF	General	
	Multi-view Depth	VGN [71]	TSDF	General	
		NGDF [70]	SDF	General	
	Multi-view RGB-D	Song <i>et al.</i> [139] NeuralGrasps [69]	TSDF SDF	General Object-specific	
Radiance	Multi-view RGB	Dex-NeRF [63]	NeRF	Scene-specific	
		Evo-NeRF [65]	Instant-NGP	Scene-specific	
		NeRF-Supervision [140]	NeRF	Scene-specific	
		MIRA [141]	NeRF, Orthographic images	Scene-specific	
			SPARTN [66]	NeRF	Scene-specific
			MVNeRF [142]	NeRF	General
			RGBGrasp [143]	NeRF, Hash encoding	General
			GraspNeRF [68]	Generalizable NeRF, TSDF	General
	Multi-view RGB-D	GaussianGrasper [64] ManiGaussian [67]	3DGS 3DGS	General General	
	Single-view RGB, Annotations	Blukis <i>et al.</i> [144, 145]	NeRF	General	

TABLE I: Overview of selected methods that leverage neural fields for manipulation tasks. See Sec. III-B for more details.

3) *Signed Distance Fields*: GraspNeRF [68] (see Fig. 12) extends NeRF for 6-DoF grasp detection of transparent and specular objects. It integrates a Truncated Signed Distance Function (TSDF) with a generalizable NeRF trained on sparse RGB images for zero-shot scene reconstruction. Similarly, Volumetric Grasping Network (VGN) [71] enables real-time 6-DoF grasp detection, synthesizing collision-free grasps in cluttered environments using a 3D voxel grid representation where each voxel contains the TSDF to the nearest surface. Song *et al.* [139] employ a TSDF to map actions to rendered views, simulating future state-action pairs for 6-DoF grasping. NeuralGrasps [69] further explores neural distance fields by learning implicit representations for grasps with multiple robotic hands. This method encodes grasps into a shared latent space, with each vector corresponding to a grasp from a specific hand. The neural Grasp Distance Field (NGDF) [70] models grasp poses as the level set of an unsigned distance field, predicting the closest valid grasp for a given 6D query pose by minimizing the unsigned distance. CenterGrasp [148] extends this concept to directly predict a displacement vector, removing the need to optimize the level-set.

4) *Feature Fields*: Feature fields represent an emerging class of NFs that integrate high-dimensional features from visual data into a unified 3D representation. These fields can map 3D points to feature vectors encoding semantic information, making them useful for context-aware grasping tasks when combined with pre-trained vision-language models.

Developments in this area have focused on the creation of feature fields that enable few-shot learning and zero-shot task-oriented grasping. *Distilled Feature Fields* (DFF) [4] proposes to distill dense features from a pre-trained vision-language model (CLIP [35]) into a 3D feature field (see Fig. 10). This allows for effective generalization across diverse object categories, making DFF particularly useful for tasks that require contextual understanding. Similarly, LERF-TOGO [5] leverages

language embeddings and DINO [36] features to accurately select target objects and specific object parts for grasping. This intuition addresses the limitations of traditional learning-based grasp planners that often ignore the semantic properties of objects. The use of vision-language models for feature distillation is also proposed by GeFF [149] and GNFactor [72], where a generalizable NeRF enriched by semantic information is adopted for manipulation and navigation. In conclusion, by incorporating semantic information directly into the 3D representation, feature fields facilitate precise, context-aware, language-guided manipulation in real-world scenarios.

5) *Visual & Tactile Sensing*: The use of NFs in multimodal visual and tactile sensing is a recent and emerging area of research. Tactile data gathered from tactile sensors offers information about contact force and contact geometry. Combining visual and tactile sensing offers several advantages in situations where vision-only might be ambiguous, such as in the presence of occlusions, challenging lighting, or reflective and transparent materials [73, 74]. In the context of multimodal sensing, NFs are mainly leveraged for tactile data generation and object reconstruction. Zhong *et al.* [150] propose to use NeRF to generate realistic tactile sensory data. Similarly, TaRF [151] leverages NeRF to synthesize novel views, which are subsequently used by a conditional diffusion model to generate the corresponding tactile signal.

A challenge in using tactile data is the disparity between real and simulated tactile images. TouchSDF [152] addresses this sim-to-real gap by combining a Convolutional Neural Network (CNN) with DeepSDF for 3D shape reconstruction from tactile inputs. As a result, objects can be reconstructed using solely tactile sensing in both simulation and the real world. Suresh *et al.* [153] employ a neural signed distance field to estimate object pose and shape during in-hand manipulation. The use of NFs, in this context, allows the robot to learn and progressively refine the object’s shape online. Moreover, NFs and visuo-tactile

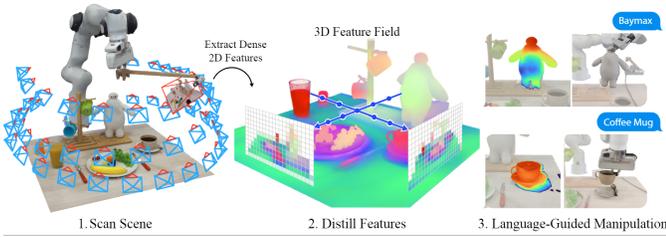


Fig. 10: Distilled feature fields [4] distill foundation model features into a feature field along with modeling a NeRF.

sensing can be employed for predicting extrinsic forces. Neural Contact Fields (NCF) [154] tracks the contact points between an object and its environment using tactile information. This method leverages NFs to generalize across different object shapes and to estimate the probability of contact at any point on an object’s surface.

Although most studies focus on implicit representations, recent techniques have begun exploring explicit approaches. Tactile-3DGS [73] and Touch-GS [74] extend 3D Gaussian Splatting by integrating visual and tactile data for 3D object reconstruction. Unlike TouchSDF, both Touch-GS and Tactile-3DGS are designed to handle the reconstruction of transparent and reflective objects. These methods have been validated in both simulated and real-world environments.

6) *Diffusion Models*: Recent work in Generative AI for robotics manipulation has explored the use of diffusion models for grasp generation, trajectory planning, and learning viewpoint and cross-embodiment invariant policies. Yoneda *et al.* [155] leverage diffusion models to predict stable object placements by learning context-dependent distributions from positive examples, eliminating the need for rejection sampling. SE(3)-Diffusion Fields [75] optimize grasp selection and trajectory generation by learning data-driven SE(3) cost functions using diffusion models. Meanwhile, VISTA [77] and RoVi-Aug [76] utilize 3D generative models for learning viewpoint-invariant policies, enabling robust generalization to new environments and unseen robots. VISTA [77] leverages Zero-NVS [34]’s zero-shot novel view synthesis capability to learn viewpoint-invariant policies, enabling robust performance in diverse environments and tasks from limited demonstration data. Similarly, RoVi-Aug [76] synthesizes augmented robot data using image-to-image generative models, allowing for zero-shot deployment on unseen robots with different embodied and largely varying camera angles. Together, these methods illustrate how 3D generative techniques can significantly improve the adaptability and effectiveness of robotic manipulation systems in real-world scenarios.

7) *Takeaways and Open Challenges in Neural Fields for Manipulation*: NFs have emerged as powerful techniques for robust 3D understanding in robotic manipulation tasks, such as grasping and pick-and-place. These representations capture detailed geometrical information and support generalization across diverse object shapes and categories. NFs have also been employed to identify optimal grasp points, improving the success rate of robotic grasps in cluttered environments. Additionally, some methods integrate these representations with language models, enabling open-vocabulary manipulation

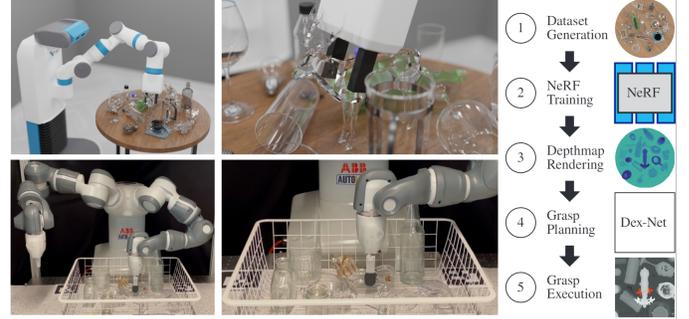


Fig. 11: Dex-NeRF [63] leverages NeRF’s depthmap rendering for transparent object grasping.

through natural language instructions.

Despite these advancements, significant challenges remain. Current approaches rely on extensive multi-view inputs or costly per-scene optimization, limiting their applicability in complex, dynamic, or unstructured environments. Furthermore, incorporating physical intuitions about object affordances and robot dynamics into the learned representations could lead to more physically grounded manipulation policies (see Sec. III-D). Finally, scaling these methods to dynamic scenes with multiple agents or articulated objects is an ongoing challenge that must be addressed for real-world deployment.

C. Neural Fields for Navigation

Autonomous navigation requires robots to perceive and model their surroundings effectively to plan collision-free paths. Traditional learning-based approaches tackle this challenge via either end-to-end [156–158] or modular systems [159–161]. Recently, the properties of NFs have proven beneficial for motion planning and navigation. NeRF’s density grid, for instance, offers a geometric approximation of the scene that aids in avoiding obstacles or learning a dynamics model. Various NeRF extensions have been proposed for navigation; some construct maps representing the visual structure of the scene, and others use autonomous agents to actively map the environment. Below, we highlight these state-of-the-art advances, structured into four key areas: Planning, Exploration, Visual Localization, and Feature Fields.

1) *Planning*: NFs’ density grid provides an approximate geometry, which is then used with a trajectory planner and state estimator in an iterative receding horizon loop for an autonomous agent to dynamically maneuver an environment with RGB camera for feedback [162]. CATNIPS [78] enables collision avoidance in a NeRF by computing collision probabilities for a robot navigating through a NeRF. It enables fast trajectory optimization using graph-based searching with spline-based trajectory optimization. SAFER-Splat [8] provides a real-time reconstruction method using Gaussian Splatting for safe robotic navigation. It operates efficiently with minimal memory, ensuring safety while maintaining high-speed performance during online mapping. NeRF’s 3D scene representation also allows learning 3D dynamical models purely from posed 2D images. Specifically, NeRF-dy [79] proposed a time-contrastive learning objective, which, when combined with Neural Radiance Fields in an auto-encoding framework,

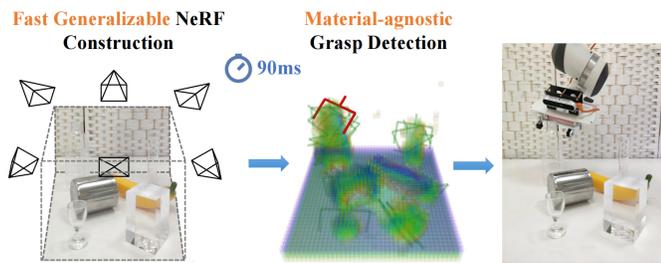


Fig. 12: Generalizable grasping with sparse multi-view images using GraspNeRF [68].

provides a viewpoint-aware neural 3D scene representation. This scene representation allows for the specification of goal points with learned predictive dynamical forward models outside the training-distribution viewpoints. CompNeRFdyn [163] extends the concepts proposed by Li *et al.* [79], introducing an auto-encoder framework in conjunction with a Graph Neural Network (GNN) [164] based dynamic model prediction in the latent space of NeRF. This forces the network to learn generalizable priors, thus aiding long-range predictions.

2) *Exploration*: A separate line of work uses modular autonomous navigation agents [165] to enable training of implicit scene representations [80]. AutoNeRF [80] (see Fig. 13) forgoes the need for carefully curated manual dataset creation with autonomously created dataset. An autonomous agent explores an unseen environment without prior access to a map and uses the experience to create implicit scene representations for novel view and semantic synthesis. DroNeRF [166] enables active reconstruction and proposed a novel optimization approach to create automated positioning of cameras for few-view reconstruction of objects in an implicit manner. Active Neural Mapping [81] (see Fig. 14) studies the problem of actively exploring an environment with a continually learned 3D scene representation, such as a NeRF. It minimizes the map uncertainty in real-time by actively selecting target spaces for exploration. By utilizing continuous geometric information encoded in the neural map, the system guides agents to find traversable paths for online scene reconstruction. DISORF [167] introduces a framework designed to facilitate real-time 3D reconstruction and visualization of scenes captured by resource-constrained mobile robots and edge devices. It addresses compute limitations and network constraints by distributing computation efficiently between the edge device and a remote server. The framework utilizes on-device SLAM systems to generate posed key frames, which are then transmitted to remote servers for high-quality 3D reconstruction and visualization using NeRF models. Finding Waldo [168] proposes baseline methods, Guided-Random Search (GRS) and Pose Interpolation-based Search (PIBS), and formulates scene exploration as an optimization problem, presenting Evolution-Guided Pose Search (EGPS) as an efficient solution.

3) *Visual Localization*: Building on the camera localization methods discussed in Sec. III-A, other works utilize implicit neural representations for top-down memory, real-time navigation, and visual localization. These approaches demonstrate the application of NFs for visual localization, which is crucial for effective navigation in dynamic environments.

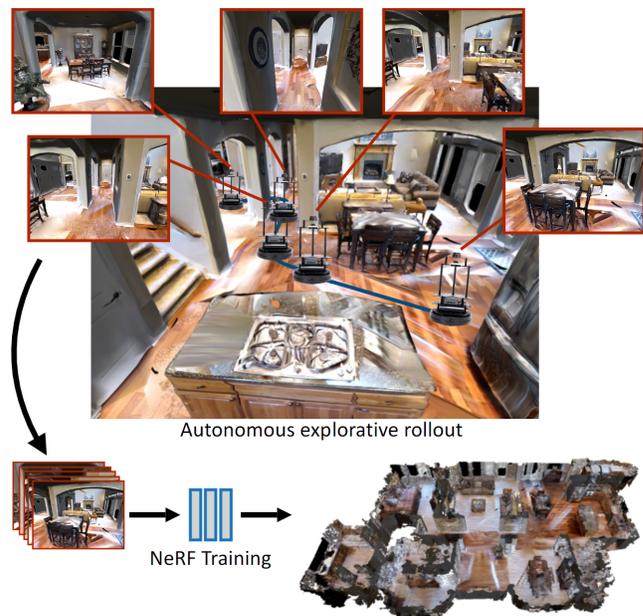


Fig. 13: AutoNeRF [80] generates 3D models of a scene by training NeRFs from data collected by autonomous agents.

RNR-Map [82] constructs a visually descriptive map of the environment, similar to Incremental Scene Synthesis [169], where a latent code at each pixel in the grid cell is embedded from an image observation and can be converted to a Neural Radiance Field. This Radiance Field can be rendered with arbitrary camera poses. A modular framework utilizing the visual information in these maps enables visual localization as well as navigation. The Le-RNRMap model [170] improves upon RNR-Map by integrating CLIP-based embedding latent codes, enabling natural language search capabilities without the need for extra-label data. Splat-Nav [84] proposes a real-time navigation system optimized for Gaussian Splat [49]-generated 3D scene representations. It consists of two key modules: Splat-Plan, which constructs collision-safe corridors and Bézier curve trajectories, and SplatLoc, facilitating robust pose estimation utilizing point cloud data and RGB images. Computational heavy-lifting, such as Bézier trajectory computation and pose optimization, is primarily handled by the CPU, allowing GPU resources to focus on tasks like online Gaussian Splat training. NeRF-IBVS [83] introduces a novel visual localization method aimed at achieving accurate localization with minimal posed images and 3D labels, addressing the challenge of acquiring such data in the real world. The method utilizes a coordinate regression network trained on a few posed images with coarse pseudo-3D labels from NeRF, followed by pose estimation with PnP and pose optimization with image-based visual servo (IBVS) leveraging scene priors from NeRF. NVINS [171] introduces a novel framework that combines NeRF with Visual-Inertial Odometry (VIO) to enhance robotic navigation in real-time. By training an absolute pose regression network using augmented image data from NeRF and quantifying uncertainty, the approach addresses positional drift and improves system reliability. Liu *et al.* [172] introduces a navigation pipeline integrating NeRF into visuomotor navigation, emphasizing the importance of memory representations for intelligent

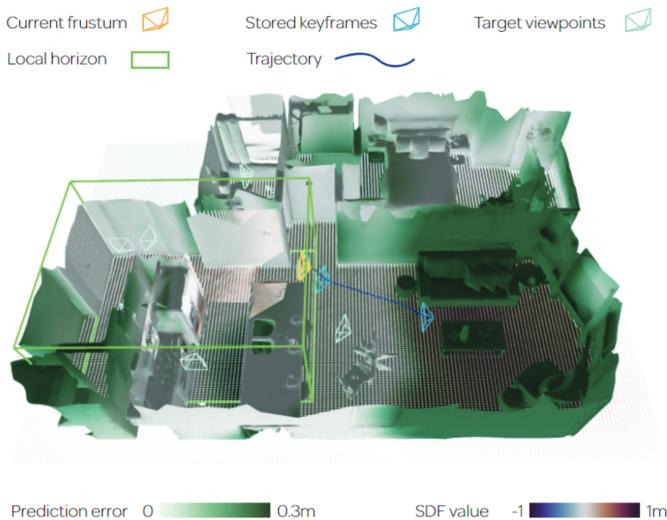


Fig. 14: Active exploration of a mobile robot to minimize prediction uncertainty [81].

agents. It presents a derivative radiance field for one-shot pose and depth estimation from a single query image, leveraging NeRF’s spatial representation for task decomposition and action generation. Other applications of NFs for navigation include inventory monitoring that enables a mobile robot to continuously update its understanding of its environment [173], visual localization [174] and exploration [175].

4) *Feature Field*: Several works have looked into lifting 2D foundation features, *i.e.*, CLIP [35], DINO [36], SAM [176] for 3D scene understanding. This 3D distillation of 2D foundation models is trained per scene and it remains to be seen their generalizability to novel scenes or environment. Owing to the world knowledge present in 2D foundation models, 3D distillation enables many real-world applications, including mobile manipulation or navigation.

For mobile manipulation, explicit 2D feature distillation into 3D using pixel-aligned open-set features can be fused into 3D maps via traditional SLAM and multi-view fusion approaches [124, 177]. CLIP-Fields [85] (see Fig. 15) implicitly utilizes a compact neural network to encode a 3D map and foundational features aligned with pixels or regions (such as LSeg [178] and Detic [179]). This specialized neural network, tailored to each scene, serves as a searchable database that aligns embeddings of images and language with 3D scene coordinates. It is designed to handle open-set queries specified in natural language. Language-embedded Radiance Fields (LEGS) [7] extended LeRF [38] to train a queryable 3D representation online as the robot traverses the environment. It enables localization of open-vocabulary object queries while training significantly faster than LeRF. GaussNav [86] creates a map representation using 3D Gaussian Splatting. This framework allows the agent to remember both the geometric and semantic details of the scene, as well as the textural features of objects using MaskRCNN distillation into the 3D domain.

Furthermore, Uni-fusion [180] proposes a universal continuous mapping framework for encoding surfaces and their properties (*e.g.*, color, infrared) without requiring training, using a Latent Implicit Map (LIM) that divides point clouds into

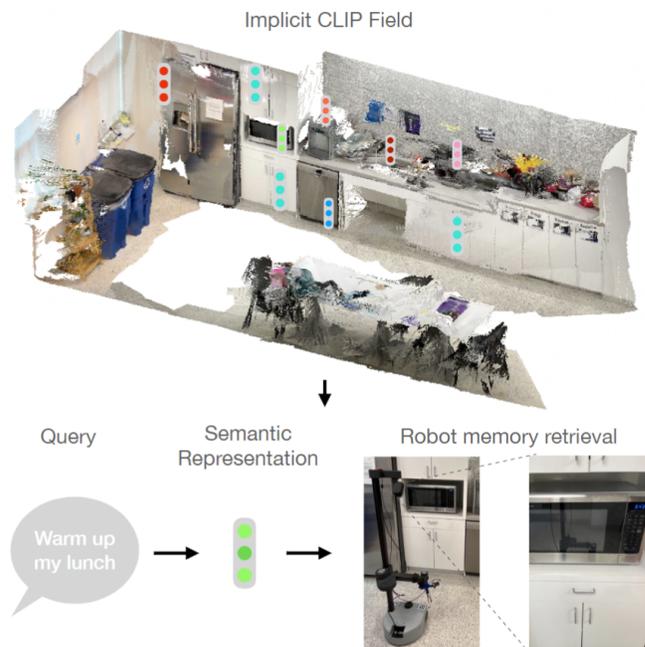


Fig. 15: Clip-Fields’s [85] semantic representation enables 3D spatial memory for mobile robots.

voxels. It supports applications such as incremental reconstruction, 2D-to-3D property transfer, and open-vocabulary scene understanding. Open-Fusion [181] proposes a real-time open-vocabulary 3D mapping and scene representation using RGB-D data, leveraging a pre-trained vision-language foundation model (VLFM) for semantic comprehension and the Truncated Signed Distance Function (TSDF) for rapid 3D reconstruction. It achieves annotation-free 3D segmentation without additional training and outperforms leading zero-shot methods.

Additionally, NF’s compact 3D representation, as well as feature distillation, makes them ideal for integration with generative models [182, 183], where the obtained 3D features from NF’s can be used directly in the projection space of the 2D Vision-Language models [184, 185] to perform a diverse set of 3D related tasks such as 3D grounding, 3D visual question answering as well as navigation.

5) *Takeaways and Open Challenges in Neural Fields for Navigation*: While Neural Fields have made significant strides in navigation, key challenges still remain. Current methods focus mainly on static environments and tasks like image-goal and vision-language navigation. Future work could extend NFs to dynamic settings, incorporating fast reconstruction techniques for real-time updates in evolving environments [186]. Another crucial direction is dynamic scene pose estimation (Sec. III-A3) to aid reconstruction and navigation in dynamic environments.

The integration of generative NFs also holds great potential. Recent diffusion model advances [39, 187] could facilitate efficient scene editing and environment creation, narrowing the sim-to-real gap. Additionally, leveraging foundation models for large-scale mobile manipulation and scene generalization could unlock further advancements. Integrating Vision-Language Models (VLMs) with implicit representations for enhanced commonsense reasoning within NFs offers another promising frontier for future exploration.

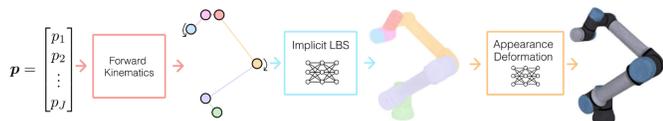


Fig. 16: Differentiable Robot rendering pipeline [10].

D. Neural Fields for Physics

Accurately simulating physics is a long-standing and challenging task traditionally marrying approaches from computer graphics and particle optimization. Linking these techniques with NFs opens up new possibilities, such as removing the need to explicitly model a scene while also imposing new challenges to researchers like balancing learned and non-learned parts.

Given the novelty of the field, NFs have seen limited use in physics-based robotics applications. One notable example is ManiGaussian [67] (see Sec. III-B2), though broader adoption remains sparse. In the following section, we discuss the possibilities and challenges introduced by the use of NFs in physical simulations for robotics. In Sec. III-D1, we review model-free approaches that do not depend on explicit physical models. In contrast, Sec. III-D2 covers physically plausible, model-based methods.

1) *Model-Free*: D-NeRF [87] was one of the first works that introduced a NeRF formulation that included a time component, allowing the representation of dynamic scenes. To decouple dynamics from structure, the authors learned an additional time-dependent MLP to map a spatial coordinate at a specific time step to a canonical space coordinate, which then serves as an input to a classical NeRF [97]. This technique is commonly known as a deformation field [188] and was extended to not just time but also arbitrary dimensions [189]. Specifically, in the case of NeRFs, deformation fields were also coined as ray bending [190, 191]. Similarly, while Li *et al.* [192] and Gao *et al.* [193] still use deformation fields to include the time component, both propose to explicitly regularize the reconstructed NeRF with an inferred scene flow.

A similar timeline of developments could be observed with 3D Gaussian Splatting. Here, Luiten *et al.* [88] was one of the first to add a time-component to each Gaussian, influencing the 3D pose while keeping the size, color, and opacity constant. Consequently, applying deformation fields to enable dynamic Gaussian Splatting was first introduced by Wu *et al.* [194]. MD-Splatting [195] extended the approach to the metric space using a rigidity and isometry regularization term adapted from Luiten *et al.* [88] and an additional momentum term. The extension to the metric space makes the reconstruction physically interpretable and thus allows applications in robotics. Yang *et al.* [196] take on a full probabilistic view and decompose the full 4D Gaussian (joint probability over space and time) into conditional Gaussian distributions. To reconstruct meshes, Liu *et al.* [197] introduce DG-Mesh, a method that utilizes a cycle-consistency loss on the predicted deformations.

Instead of conditioning their NeRF on a time component, Abou-Chakra *et al.* [198] approach the problem differently through particles that evolve over time and thus dynamically adapt to the scene in an online manner. They extended their approach also to Gaussian Splatting [199] while also

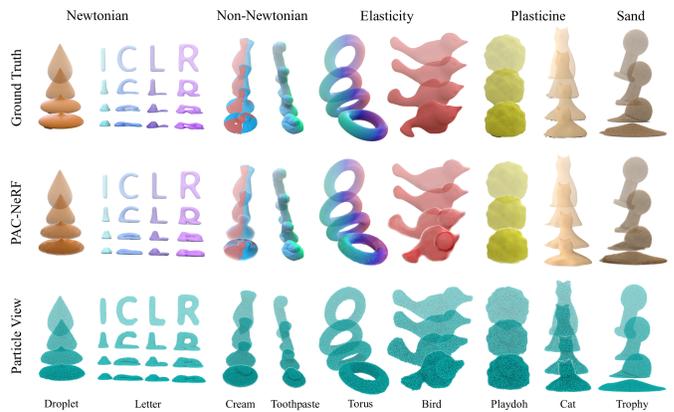


Fig. 17: An overview of the different materials model-based NFs are able to simulate [89].

introducing the concept of visual forces. Nonetheless, a robotic input to interact with the scene is missing. Opposed to that most recently, Li *et al.* [79] proposed to first learn a latent state variable which then serves as an additional input to their NeRF. Given robot demonstrations, they then also learn a latent (implicit) dynamics model, inferring the next state given the current state and a robot action. Similarly, ManiGaussian [67] does not directly condition their Gaussians on a time component but rather on a robot action and thus, learn a Gaussian world model. Using known forward kinematics of a robot to condition Gaussians in Gaussian Splatting, Liu *et al.* [10] tackle the problem of differentiable robot rendering (see Fig. 16), allowing to pass gradients from robot images down to the robot joint states, enabling various tasks such as text-to-robot pose, robot & camera pose estimation, motion retargeting with point tracker and robot control with a generative video model [200].

To conclude, one could consider the aforementioned models as model-free as no underlying physics model is explicitly used nor enforced through regularization. While most of these methods presumably have a constant density [87, 97], it can not be guaranteed that the learning procedure will vanish parts through, *e.g.*, ray bending and thus produce physically implausible results. As a result, the utility of these models for reliable and safe robotics is unclear and yet to be explored.

2) *Model-Based*: Opposed to model-free methods (see Sec. III-D1), model-based methods incorporate underlying physical principles such as the aforementioned constant density [89] and thus, could be considered as physically correct. In our analysis, we classify model-based methods based on their simulation scope: *rigid*, *articulated* and *non-rigid*.

Rigid Objects: Similar to some of the aforementioned model-free methods, Hofherr *et al.* [201] used only a photometric loss for optimization, while constraining rigid object motion through an underlying physical dynamics model. Cleac’h *et al.* [202] separated object reconstruction from parameter estimation, using static object images to train a neural representation and videos of moving objects to infer physical properties like friction and mass. Instead of videos, NeRF2Physics [203] distilled these properties from language and associated them with spatial language embeddings. MovingParts [204] removed the assumption of object separation, thus automatically detecting coherently moving rigid parts and their transformations over time.

Articulated Objects: Articulated objects can be considered a class between purely rigid objects and non-rigid objects, as they impose constraints on rigid parts and their relative movement but do not decompose them down to the particle level. Most articulated object research has focused on pure visual reconstruction [132, 205], with limited emphasis on their physical interaction properties, such as manipulation and dynamics. Notably, some recent works have begun to bridge this gap by incorporating physical reasoning and interaction capabilities into models, allowing for a more comprehensive understanding and manipulation of articulated objects in real-world scenarios. Robot See Robot Do (RSRD) [206] enables robots to imitate articulated object manipulation from a single monocular video. RSRD introduces 4D Differentiable Part Models (4D-DPM) to reconstruct 3D part motion through an analysis-by-synthesis approach that optimizes geometric regularizes from a single video. This enables the robot to replicate object part motions by planning bi-manual arm trajectories, achieving notable success rates in physical execution without task-specific training or annotation.

Non-Rigid Objects: Non-rigid objects, unlike rigid and articulated ones, consist of numerous individual moving particles, making their simulation more complex yet generalizable. Non-rigid objects can be further categorized into subtypes such as deformable objects and fluids. One of the most versatile methods for simulating these objects is the Material Point Method (MPM), a framework capable of modeling a wide range of materials [9, 89]. PAC-NeRF [89] (see Fig. 17) introduced a hybrid particle and grid-based NeRF representation that enables conversions between the two. In contrast, PhysGaussian [9] directly utilizes the particle-like nature of Gaussians in Gaussian Splatting, eliminating PAC-NeRF’s explicit conversion step. Similarly, PIE-NeRF [207] avoids PAC-NeRF’s particle-to-rest-pose conversion, reducing over-smoothing. Spring-Gaus [208] clusters Gaussians as point masses by inferring a static set of Gaussians followed by the sampling of anchor points (*i.e.*, particles).

To simulate fluids, Yu *et al.* [209] split the velocity field into a base flow and a vortex particle flow, ensuring physical accuracy through a density and projection loss. Similar to PhysGaussian [9], Gaussian Splashing [210] uses Gaussian kernel centers as particles in a physics simulation. However, unlike PhysGaussian, the authors distinguish between solids and fluids, allowing them to first reconstruct a solid scene and then synthesize fluids within it. ClimateNeRF [211] followed a similar process, first reconstructing a scene with a classical NeRF pipeline before simulating different weather effects. Additionally, Zhong *et al.* [212] combined a neural deformation field with a Kirchhoff stress field in the same latent space, enabling faster, memory-efficient simulations by operating within this latent space.

3) *Takeaways and Open Challenges in Neural Fields for Physics:* Significant progress has been made in understanding and inferring physics; however, the challenge remains to seamlessly integrate these models with robots to create a truly simulatable, general, and interactive environment. Furthermore, it is still unclear how effectively policies learned in these simulations can be transferred to the real world.

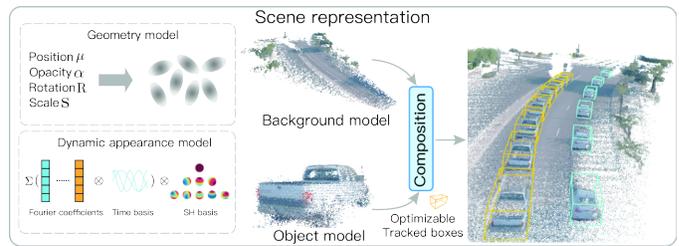


Fig. 18: The compositional pipeline for Street Gaussians [213].

E. Neural Fields in Autonomous Driving

High-quality mapping of large-scale environments is essential for autonomous driving systems. A high-fidelity map of the entire operating domain serves as a powerful prior for various tasks, including robot localization (see Sec. III-A), navigation, and collision avoidance (see Sec. III-C). Additionally, large-scale scene reconstructions facilitate closed-loop robotic simulations. Autonomous driving systems are often evaluated by re-simulating previously encountered scenarios; however, any deviation from the original encounter can alter the vehicle’s trajectory, necessitating high-fidelity novel view renderings along the adjusted path. In addition to basic view synthesis, scene-conditioned NeRFs can modify environmental lighting conditions, such as camera exposure, weather, or time of day, further enhancing simulation scenarios.

Neural Fields have become a prominent framework in autonomous driving due to their ability to generate photorealistic 3D environments from RGB images. These environments are highly valuable for constructing immersive simulation systems with several key features, as previously discussed: First, NFs offer extensive *manipulability and compositionality* (Sec. III-E1), allowing for the seamless integration and manipulation of objects within a scene. This facilitates the simulation of complex scenarios, such as collisions, which are difficult to replicate in physical settings. Second, they produce scenes with impressive *photorealism* (Sec. III-E2), enabling realistic simulations from visual data. Finally, their strong *generalizability* (Sec. III-E3) from sparse inputs allows for creating accurate, scalable environments, enhancing research in embodied AI. These traits, as discussed in the following subsections, enable the creation of simulated environments that faithfully represent real-world scenarios, thereby facilitating research in embodied AI.

1) *Manipulability and Compositionality:* The underlying principle of utilizing photorealistic simulations lies in their efficacy as proxies for real-world environments in advancing research on embodied AI. Agents operating within these simulations can strategize and execute actions, thereby enhancing their ability to handle edge cases and facilitating smoother transitions to the real world with reduced domain gaps.

One of the first works that explored this paradigm is Neural Scene Graphs [12]. It introduces a hierarchical approach to scene modeling, incorporating static and dynamic elements such as object appearance and shape. Utilizing a directed acyclic graph, scenes are uniquely defined, with nodes representing camera intrinsic, latent object codes and neural radiance fields for both static and dynamic elements and edges denoting transformations or property assignments. StreetNeRF [90] also

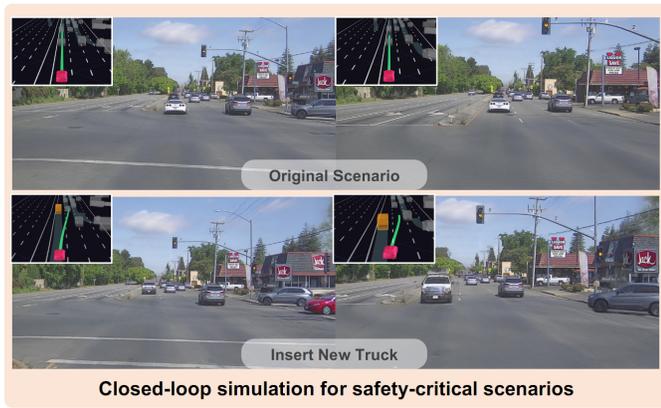


Fig. 19: Photorealistic editing results from UniSim [92].

focuses on compositional scene representation. It addresses limitations in traditional NeRFs for street-view synthesis by jointly considering large-scale background scenes and foreground moving vehicles. It improves scene parameterization and camera pose learning, leveraging noisy LiDAR points and geometry-based confidence to handle depth outliers. Experimental results demonstrate significant improvements in street-view synthesis and rendering moving vehicles compared to state-of-the-art methods. Similarly, Panoptic Neural Fields (PNF) [91] offers an object-aware neural scene representation, dividing scenes into objects and backgrounds. Leveraging compact MLPs for each object, PNF achieves faster processing while retaining category-specific priors, enabling tasks such as novel view synthesis, 2D panoptic segmentation, and 3D scene editing in real-world dynamic scenes. SUDS [214] extends NeRFs for dynamic urban scenes by efficiently encoding static, dynamic, and far-field radiance fields using separate hash table data structures. Leveraging unlabeled target signals and various reconstruction losses, SUDS decomposes scenes into static backgrounds, individual objects, and their motions, achieving state-of-the-art performance on tasks like novel-view synthesis, unsupervised 3D instance segmentation, and 3D cuboid detection while significantly reducing training time compared to previous methods. EmerNeRF [215] learns spatial-temporal representations of dynamic driving scenes by stratifying scenes into static and dynamic fields and parameterizing an induced flow field. Additionally, by lifting 2D visual foundation model features into 4D space-time, EmerNeRF improves semantic generalization and enhances 3D perception performance.

Other works leverage explicit representations to accelerate rendering. Street Gaussians [213] (see Fig. 18) introduces a dynamic urban scene model using point clouds with semantic logits and 3D Gaussians for vehicles and backgrounds, enabling efficient scene editing and fast rendering. It outperforms prior methods on benchmarks using off-the-shelf tracker poses. Driving Gaussians [216] reconstructs dynamic driving scenes with static 3D Gaussians and a dynamic Gaussian graph, using LiDAR priors for Gaussian splatting to achieve photorealistic synthesis and multi-camera consistency, surpassing existing techniques.

2) *Photorealistic Simulators*: NeRFs excel in static scenes with controlled lighting conditions, but it faces difficulties when working with image collections from unpredictable real-world environments, which include varying weather, lighting,

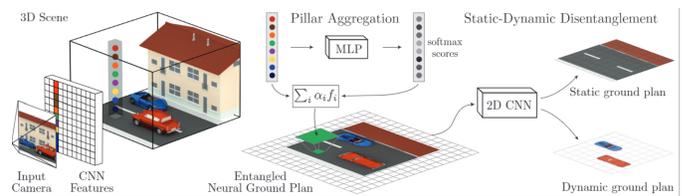


Fig. 20: An overview of Neural Groundplans approach. [94]

or temporary obstructions. NeRF-W [217] seeks to overcome these challenges by using appearance embeddings and transient networks. Subsequent neural rendering for autonomous driving applications focused on building photorealistic simulators. MARS [11] introduces an autonomous driving simulator based on NeRFs to address corner cases in autonomous vehicle driving. It features instance-aware modeling for separate foreground instances and background environments, proposing a modular design that enables flexible switching between different NeRF-related components, and achieves state-of-the-art photorealism results. Notably, MARS is open-sourced, distinguishing it from most counterparts in the field. UniSim [92] (see Fig. 19), a neural sensor simulator, enables closed-loop evaluation of self-driving vehicles by converting recorded driving logs into realistic multi-sensor simulations. Leveraging neural feature grids, UniSim reconstructs scene elements and dynamically simulates LiDAR and camera data, facilitating accurate assessment of autonomy systems on safety-critical scenarios. DriveEnv-NeRF [218] proposes to use NeRFs to create high-fidelity simulations for validating and forecasting the performance of autonomous driving agents in real-world scenes. By rendering realistic images from novel viewpoints and constructing 3D meshes to emulate collisions, it bridges the sim-to-real gap, enhancing the robustness and real-world performance of autonomous driving agents compared to those trained with standard simulators. Lindström *et al.* [219] propose methods to improve perception model robustness to NeRF artifacts, enhancing performance on both simulated and real data. Their large-scale investigation evaluates object detectors and an online mapping model on real and simulated data, demonstrating improved model robustness and, in some cases, better real-world performance.

3) *Generalizability*: A separate line of works focused on generalizable systems for outdoor scenes. NeO 360 [93] introduces a novel approach for few-view view synthesis of outdoor scenes, overcoming limitations in existing methods by reconstructing 360° scenes from a single or few posed RGB images. By capturing the distribution of complex real-world outdoor 3D scenes and using a hybrid image-conditional triplanar representation, NeO 360 offers generalizability to new views and novel scenes from as few as a single image during inference. Neural Groundplans [94] (see Fig. 20) introduces a method for mapping 2D image observations to a persistent 3D scene representation, facilitating novel view synthesis and disentangling movable and immovable scene components. Trained self-supervised from unlabeled multi-view observations, it leverages ground-aligned 2D feature grids inspired by bird's-eye-view representation, enabling efficient scene understanding tasks such as instance-level segmentation and 3D bounding box prediction. 6Imgt3D [95] uses a transformer-based encoder-

renderer method designed for efficient and scalable single-shot 3D reconstruction from six outward-facing images in large-scale, unbounded outdoor driving scenarios.

4) *Takeaways and Open Challenges in Neural Fields for Autonomous Driving*: Despite the promising progress in NFs for autonomous driving, several open challenges remain. Current methods focus on photorealistic simulators, which are dynamic, compositional, and realistic. One avenue of future work is training policies in such NF-based simulators and transferring them to the real-world. Connecting the success of NFs in autonomous driving with real-world deployment is an exciting avenue for future work. Generalizable reconstruction has seen some early signs of life with recent works but still remains largely underexplored. Future works could look at the efficiency of generalizable outdoor scene reconstruction methods, as well as advances that focus on sim2real transfer and pose-free reconstruction. This avenue of research is exciting as it opens the door for creating photorealistic simulators from a few images in the real world.

Another promising direction for autonomous driving research is integrating generative methods like diffusion models with the NFs’ paradigm. Future work could look at creating new scenarios via NF editing that are difficult to create in the real world, such as collision avoidance to train policies via reward models in NFs’ simulation. Generative asset creation through a few images from the real world is another potential avenue for NFs’ research for autonomous driving.

Furthermore, the integration of NFs into generative models such as shown in Lift3D [220] and Adv3D [221] facilitates data augmentation, addressing the challenges posed by the diversity of driving scenes. Given the high costs associated with capturing all potential scenarios, data augmentation emerges as a valuable strategy and promising future direction for expanding training datasets and improving model performance.

IV. OPEN CHALLENGES OF NEURAL FIELDS IN ROBOTICS

Despite the exciting progress in the field, there are still several open challenges for various robotic applications to adopt Neural Fields.

- **Efficiency**: NFs are computationally intensive and may not naturally operate in real-time, which is often a critical requirement for robotics applications. There is a need for significant optimization or simplification to make these models run efficiently on robotics hardware, which may have limited computational resources compared to dedicated GPUs used in data centers.
- **Dynamic environments**: Robotics often involve operating in dynamic environments where objects and scene configurations change over time. Capturing and updating NFs to reflect these changes in real-time remains a challenging task.
- **Sensor integration**: Effectively integrating data from various sensors (*e.g.*, LiDAR, RGB cameras, depth sensors) to enhance the robustness and performance of NFs is relatively under-explored. Advanced sensor fusion techniques could potentially bridge this gap.
- **Generalization**: Existing techniques often require dense input data and struggle with sensor noise or occlusions.

Developing methods that can leverage priors learned from web-scale datasets to generalize across varied scenarios offers a promising direction.

- **Physical information**: While NFs excel at representing visual aspects, they do not inherently understand physical properties like weight or friction. Extending NFs to incorporate physics simulations could enable more realistic interactions for robots.
- **Data efficiency and augmentation**: Current approaches are data-hungry, which is impractical for real-world applications. Innovations in data-efficient learning techniques and realistic data augmentation could help in overcoming these limitations.
- **Multi-modal, multi-task, and efficient scene understanding**: Developing neural field approaches that can handle multiple tasks and modalities simultaneously while maintaining efficiency in scene understanding is crucial for holistic robotic perception.
- **Performance evaluation**: Establishing standardized metrics and benchmarks for evaluating the performance of NFs in robotic applications will be essential for tracking progress and comparing different approaches.
- **Collaborative frameworks**: There is a need for frameworks that support collaboration between robots using NFs, enabling them to share learnings and improve collective understanding and decision-making in complex environments.

REFERENCES

- [1] B. Hu, J. Huang, Y. Liu, Y.-W. Tai, and C.-K. Tang, “Nerf-rpn: A general framework for object detection in nerfs,” in *Conference on Computer Vision and Pattern Recognition*, 2023, pp. 23 528–23 538.
- [2] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, “iNeRF: Inverting neural radiance fields for pose estimation,” in *International Conference on Intelligent Robots and Systems*, 2021.
- [3] M. Li, S. Liu, and H. Zhou, “Sgs-slam: Semantic gaussian splatting for neural dense slam,” *arXiv preprint arXiv:2402.03246*, 2024.
- [4] W. Shen, G. Yang, A. Yu, J. Wong, L. P. Kaelbling, and P. Isola, “Distilled feature fields enable few-shot language-guided manipulation,” in *7th Annual Conference on Robot Learning*, 2023.
- [5] A. Rashid, S. Sharma, C. M. Kim, J. Kerr, L. Y. Chen, A. Kanazawa, and K. Goldberg, “Language embedded radiance fields for zero-shot task-oriented grasping,” in *Conference on Robot Learning*, 2023.
- [6] M. N. Qureshi, S. Garg, F. Yandun, D. Held, G. Kantor, and A. Silwal, “Splatsim: Zero-shot sim2real transfer of rgb manipulation policies using gaussian splatting,” 2024.
- [7] J. Yu, K. Hari, K. Srinivas, K. El-Refai, A. Rashid, C. M. Kim, J. Kerr1, R. Cheng, M. Z. Irshad, A. Balakrishna, T. Kollar, and K. Goldberg, “Language-embedded gaussian splats (legs): Incrementally building room-scale representations with a mobile robot,” in *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2024.
- [8] T. Chen, A. Swann, J. Yu, O. Shorinwa, R. Murai, M. Kennedy III, and M. Schwager, “Safer-splat: A control barrier function for safe navigation with online gaussian splatting maps,” *arXiv preprint arXiv:2409.09868*, 2024.
- [9] T. Xie, Z. Zong, Y. Qiu, X. Li, Y. Feng, Y. Yang, and C. Jiang, “Physgaussian: Physics-integrated 3d gaussians for generative dynamics,” *arXiv preprint arXiv:2311.12198*, 2023.
- [10] R. Liu, A. Canberk, S. Song, and C. Vondrick, “Differentiable robot rendering,” in *8th Annual Conference on Robot Learning*, 2024.
- [11] Z. Wu, T. Liu, L. Luo, Z. Zhong, J. Chen, H. Xiao, C. Hou, H. Lou, Y. Chen, R. Yang, Y. Huang, X. Ye, Z. Yan, Y. Shi, Y. Liao, and H. Zhao, “Mars: An instance-aware, modular and realistic simulator for autonomous driving,” *CICAI*, 2023.
- [12] J. Ost, F. Mannan, N. Thuerey, J. Knodt, and F. Heide, “Neural scene graphs for dynamic scenes,” 2021.

- [13] H. Fan, H. Su, and L. J. Guibas, "A point set generation network for 3d object reconstruction from a single image," in *Conference on computer vision and pattern recognition*, 2017, pp. 605–613.
- [14] M. Z. Irshad, T. Kollar, M. Laskey, K. Stone, and Z. Kira, "Centersnap: Single-shot multi-object 3d shape reconstruction and categorical 6d pose and size estimation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.
- [15] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, "Learning representations and generative models for 3d point clouds," in *International conference on machine learning*. PMLR, 2018, pp. 40–49.
- [16] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4490–4499.
- [17] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang, "Pixel2mesh: Generating 3d mesh models from single rgb images," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 52–67.
- [18] Z. Chen, A. Tagliasacchi, and H. Zhang, "Bsp-net: Generating compact meshes via binary space partitioning," in *Conference on computer vision and pattern recognition*, 2020, pp. 45–54.
- [19] Y. Liao, S. Donne, and A. Geiger, "Deep marching cubes: Learning explicit surface representations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2916–2925.
- [20] M. Breyer, J. J. Chung, L. Ott, S. Roland, and N. Juan, "Volumetric grasping network: Real-time 6 dof grasp detection in clutter," in *Conference on Robot Learning*, 2020.
- [21] Y. Xie, T. Takikawa, S. Saito, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar, "Neural fields in visual computing and beyond.(2021)," 2021.
- [22] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *ECCV*, 2020.
- [23] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "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.
- [24] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, "Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction," *arXiv preprint arXiv:2106.10689*, 2021.
- [25] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [26] A. Guédon and V. Lepetit, "Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering," in *Conference on Computer Vision and Pattern Recognition*, 2024.
- [27] S. Zakharov, K. Liu, A. Gaidon, and R. Ambrus, "Refine: Recursive field networks for cross-modal multi-scene representation," 2024.
- [28] T. Takikawa, J. Litalien, K. Yin, K. Kreis, C. Loop, D. Nowrouzezahrai, A. Jacobson, M. McGuire, and S. Fidler, "Neural geometric level of detail: Real-time rendering with implicit 3D shapes," 2021.
- [29] S. Huang, Z. Gojicic, Z. Wang, F. Williams, Y. Kasten, S. Fidler, K. Schindler, and O. Litany, "Neural lidar fields for novel view synthesis," in *International Conference on Computer Vision (ICCV)*, 2023.
- [30] I. Hwang, J. Kim, and Y. M. Kim, "Ev-nerf: Event based neural radiance field," in *Winter Conference on Applications of Computer Vision (WACV)*, 2023.
- [31] M. Tancik, P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. Barron, and R. Ng, "Fourier features let networks learn high frequency functions in low dimensional domains," in *Neurips*, 2020, pp. 7537–7547.
- [32] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, "pixelnerf: Neural radiance fields from one or few images," in *CVPR*, 2021.
- [33] R. Wu, B. Mildenhall, P. Henzler, K. Park, R. Gao, D. Watson, P. P. Srinivasan, D. Verbin, J. T. Barron, B. Poole, and A. Holynski, "Reconfusion: 3d reconstruction with diffusion priors," *arXiv*, 2023.
- [34] K. Sargent, Z. Li, T. Shah, C. Herrmann, H.-X. Yu, Y. Zhang, E. R. Chan, D. Lagun, L. Fei-Fei, D. Sun, and J. Wu, "ZeroNVS: Zero-shot 360-degree view synthesis from a single real image," *CVPR*, 2024.
- [35] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [36] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby *et al.*, "Dinov2: Learning robust visual features without supervision," *arXiv preprint arXiv:2304.07193*, 2023.
- [37] S. Kobayashi, E. Matsumoto, and V. Sitzmann, "Decomposing nerf for editing via feature field distillation," *Advances in Neural Information Processing Systems*, vol. 35, pp. 23 311–23 330, 2022.
- [38] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik, "Lerf: Language embedded radiance fields," in *International Conference on Computer Vision (ICCV)*, 2023.
- [39] J. Ho, A. Jain, and P. Abbeel, "Denosing diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [40] G. Wang, L. Pan, S. Peng, S. Liu, C. Xu, Y. Miao, W. Zhan, M. Tomizuka, M. Pollefeys, and H. Wang, "Nerf in robotics: A survey," *arXiv preprint arXiv:2405.01333*, 2024.
- [41] C. Sun, M. Sun, and H.-T. Chen, "Direct voxel grid optimization: Superfast convergence for radiance fields reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5459–5469.
- [42] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3D reconstruction in function space," in *Conference on Computer Vision and Pattern Recognition*, 2019.
- [43] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng, "Fourier features let networks learn high frequency functions in low dimensional domains," *NeurIPS*, 2020.
- [44] J. T. Kajiya and B. P. Von Herzen, "Ray tracing volume densities," *ACM SIGGRAPH computer graphics*, vol. 18, no. 3, pp. 165–174, 1984.
- [45] N. Max, "Optical models for direct volume rendering," *IEEE Transactions on Visualization and Computer Graphics*, 1995.
- [46] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Trans. Graph.*, vol. 41, no. 4, pp. 102:1–102:15, Jul. 2022.
- [47] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa, "Plenoxels: Radiance fields without neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 5501–5510.
- [48] R. Li, M. Tancik, and A. Kanazawa, "Nerfacc: A general nerf acceleration toolbox," *arXiv preprint arXiv:2210.04847*, 2022.
- [49] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, 2023.
- [50] M. Zwicker, H. Pfister, J. Van Baar, and M. Gross, "Surface splatting," in *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 2001, pp. 371–378.
- [51] M. M. Loper and M. J. Black, "Opendr: An approximate differentiable renderer," in *ECCV*, 2014.
- [52] S. Liu, T. Li, W. Chen, and H. Li, "Soft rasterizer: A differentiable renderer for image-based 3d reasoning," in *ICCV*, 2019.
- [53] Y. Lin, T. Müller, J. Tremblay, B. Wen, S. Tyree, A. Evans, P. A. Vela, and S. Birchfield, "Parallel inversion of neural radiance fields for robust pose estimation," in *ICRA*, 2023.
- [54] A. Moreau, N. Piasco, D. Tsishkou, B. Stanculescu, and A. de La Fortelle, "Lens: Localization enhanced by nerf synthesis," in *Conference on Robot Learning*. PMLR, 2022, pp. 1347–1356.
- [55] Z. Wang, S. Wu, W. Xie, M. Chen, and V. A. Prisacariu, "Nerf-: Neural radiance fields without known camera parameters," *arXiv preprint arXiv:2102.07064*, 2021.
- [56] A. Meuleman, Y.-L. Liu, C. Gao, J.-B. Huang, C. Kim, M. H. Kim, and J. Kopf, "Progressively optimized local radiance fields for robust view synthesis," in *CVPR*, 2023.
- [57] C. Xu, B. Wu, J. Hou, S. Tsai, R. Li, J. Wang, W. Zhan, Z. He, P. Vajda, K. Keutzer, and M. Tomizuka, "Nerf-det: Learning geometry-aware volumetric representation for multi-view 3d object detection," in *ICCV*, 2023.
- [58] B. Wen, J. Tremblay, V. Blukis, S. Tyree, T. Muller, A. Evans, D. Fox, J. Kautz, and S. Birchfield, "Bundlesdf: Neural 6-dof tracking and 3d reconstruction of unknown objects," *CVPR*, 2023.
- [59] J. Sun, Y. Xu, M. Ding, H. Yi, C. Wang, J. Wang, L. Zhang, and M. Schwager, "Nerf-loc: Transformer-based object localization within neural radiance fields," *IEEE Robotics and Automation Letters*, 2023.
- [60] A. Simeonov, Y. Du, A. Tagliasacchi, J. B. Tenenbaum, A. Rodriguez, P. Agrawal, and V. Sitzmann, "Neural descriptor fields: Se (3)-equivariant object representations for manipulation," in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2022.
- [61] E. Chun, Y. Du, A. Simeonov, T. Lozano-Perez, and L. Kaelbling, "Local neural descriptor fields: Locally conditioned object representations for manipulation," *arXiv preprint arXiv:2302.03573*, 2023.
- [62] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu, "Synergies between affordance and geometry: 6-Dof grasp detection via implicit representations," *Robotics: science and systems*, 2021.
- [63] J. Ichnowski, Y. Avigal, J. Kerr, and K. Goldberg, "Dex-nerf: Using

- a neural radiance field to grasp transparent objects,” *arXiv preprint arXiv:2110.14217*, 2021.
- [64] Y. Zheng, X. Chen, Y. Zheng, S. Gu, R. Yang, B. Jin, P. Li, C. Zhong, Z. Wang, L. Liu *et al.*, “Gaussiangrasper: 3d language gaussian splatting for open-vocabulary robotic grasping,” *arXiv preprint arXiv:2403.09637*, 2024.
- [65] J. Kerr, L. Fu, H. Huang, Y. Avigal, M. Tancik, J. Ichnowski, A. Kanazawa, and K. Goldberg, “Evo-nerf: Evolving nerf for sequential robot grasping of transparent objects,” in *Conference on Robot Learning*, 2022.
- [66] A. Zhou, M. J. Kim, L. Wang, P. Florence, and C. Finn, “Nerf in the palm of your hand: Corrective augmentation for robotics via novel-view synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17907–17917.
- [67] G. Lu, S. Zhang, Z. Wang, C. Liu, J. Lu, and Y. Tang, “Manigaussian: Dynamic gaussian splatting for multi-task robotic manipulation,” *arXiv preprint arXiv:2403.08321*, 2024.
- [68] Q. Dai, Y. Zhu, Y. Geng, C. Ruan, J. Zhang, and H. Wang, “Grasprerf: multiview-based 6-dof grasp detection for transparent and specular objects using generalizable nerf,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 1757–1763.
- [69] N. Khargonkar, N. Song, Z. Xu, B. Prabhakaran, and Y. Xiang, “Neuralgrasps: Learning implicit representations for grasps of multiple robotic hands,” in *Conference on Robot Learning*. PMLR, 2023.
- [70] T. Weng, D. Held, F. Meier, and M. Mukadam, “Neural grasp distance fields for robot manipulation,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 1814–1821.
- [71] M. Breyer, J. J. Chung, L. Ott, R. Siegwart, and J. Nieto, “Volumetric grasping network: Real-time 6 dof grasp detection in clutter,” in *Conference on Robot Learning*. PMLR, 2021, pp. 1602–1611.
- [72] Y. Ze, G. Yan, Y.-H. Wu, A. Macaluso, Y. Ge, J. Ye, N. Hansen, L. E. Li, and X. Wang, “Gnfactor: Multi-task real robot learning with generalizable neural feature fields,” in *Conference on Robot Learning*. PMLR, 2023.
- [73] M. Comi, A. Tonioni, M. Yang, J. Tremblay, V. Blukis, Y. Lin, N. F. Lepora, and L. Aitchison, “Snap-it, tap-it, splat-it: Tactile-informed 3d gaussian splatting for reconstructing challenging surfaces,” *arXiv preprint arXiv:2403.20275*, 2024.
- [74] A. Swann, M. Strong, W. K. Do, G. S. Camps, M. Schwager, and M. Kennedy III, “Touch-gs: Visual-tactile supervised 3d gaussian splatting,” *arXiv preprint arXiv:2403.09875*, 2024.
- [75] J. Urain, N. Funk, J. Peters, and G. Chalvatzaki, “Se(3)-diffusionfields: Learning smooth cost functions for joint grasp and motion optimization through diffusion,” 2023.
- [76] L. Y. Chen, C. Xu, K. Dharmarajan, M. Z. Irshad, R. Cheng, K. Keutzer, M. Tomizuka, Q. Vuong, and K. Goldberg, “Rovi-aug: Robot and viewpoint augmentation for cross-embodiment robot learning,” in *Conference on Robot Learning (CoRL)*, 2024.
- [77] S. Tian, B. Wulfe, K. Sargent, K. Liu, S. Zakharov, V. Guizilini, and J. Wu, “View-invariant policy learning via zero-shot novel view synthesis,” *arXiv*, 2024.
- [78] T. Chen, P. Culbertson, and M. Schwager, “Catnips: Collision avoidance through neural implicit probabilistic scenes,” *arXiv preprint arXiv:2302.12931*, 2023.
- [79] Y. Li, S. Li, V. Sitzmann, P. Agrawal, and A. Torralba, “3d neural scene representations for visuomotor control,” in *Conference on Robot Learning*. PMLR, 2022, pp. 112–123.
- [80] P. Marza, L. Matignon, O. Simonin, D. Batra, C. Wolf, and D. S. Chaplot, “Autonerf: Training implicit scene representations with autonomous agents,” *arXiv preprint arXiv:2304.11241*, 2023.
- [81] Z. Yan, H. Yang, and H. Zha, “Active neural mapping,” in *Intl. Conf. on Computer Vision (ICCV)*, 2023.
- [82] O. Kwon, J. Park, and S. Oh, “Renderable neural radiance map for visual navigation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [83] Y. Wang, Y. Yan, D. Shi, W. Zhu, J. Xia, T. Jeff, S. Jin, K. Gao, X. Li, and X. Yang, “Nerf-ibvs: Visual servo based on nerf for visual localization and navigation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [84] T. Chen, O. Shorinwa, W. Zeng, J. Bruno, P. Dames, and M. Schwager, “Splat-nav: Safe real-time robot navigation in gaussian splatting maps,” *arXiv preprint arXiv:2403.02751*, 2024.
- [85] N. M. Mahi Shafiullah, C. Paxton, L. Pinto, S. Chintala, and A. Szlam, “Clip-fields: Weakly supervised semantic fields for robotic memory,” *arXiv e-prints*, pp. arXiv–2210, 2022.
- [86] X. Lei, M. Wang, W. Zhou, and H. Li, “Gaussnav: Gaussian splatting for visual navigation,” *arXiv preprint arXiv:2403.11625*, 2024.
- [87] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer, “D-NeRF: Neural Radiance Fields for Dynamic Scenes,” in *CVPR*, 2021.
- [88] J. Luiten, G. Kopanas, B. Leibe, and D. Ramanan, “Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis,” in *3DV*, 2024.
- [89] X. Li, Y.-L. Qiao, P. Y. Chen, K. M. Jatavallabhula, M. Lin, C. Jiang, and C. Gan, “Pac-nerf: Physics augmented continuum neural radiance fields for geometry-agnostic system identification,” 2023.
- [90] Z. Xie, J. Zhang, W. Li, F. Zhang, and L. Zhang, “S-nerf: Neural radiance fields for street views,” 2023.
- [91] A. Kundu, K. Genova, X. Yin, A. Fathi, C. Pantofaru, L. J. Guibas, A. Tagliasacchi, F. Dellaert, and T. Funkhouser, “Panoptic neural fields: A semantic object-aware neural scene representation,” in *Conference on Computer Vision and Pattern Recognition*, 2022.
- [92] Z. Yang, Y. Chen, J. Wang, S. Manivasagam, W.-C. Ma, A. J. Yang, and R. Urtasun, “Unisim: A neural closed-loop sensor simulator,” 2023.
- [93] M. Z. Irshad, S. Zakharov, K. Liu, V. Guizilini, T. Kollar, A. Gaidon, Z. Kira, and R. Ambrus, “Neo 360: Neural fields for sparse view synthesis of outdoor scenes,” in *International Conference on Computer Vision (ICCV)*, 2023.
- [94] P. Sharma, A. Tewari, Y. Du, S. Zakharov, R. Ambrus, A. Gaidon, W. T. Freeman, F. Durand, J. B. Tenenbaum, and V. Sitzmann, “Neural groundplans: Persistent neural scene representations from a single image,” *arXiv preprint arXiv:2207.11232*, 2022.
- [95] T. Gieruc, M. Kästingschäfer, S. Bernhard, and M. Salzmann, “6img-to-3d: Few-image large-scale outdoor driving scene reconstruction,” 2024.
- [96] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey, “Barf: Bundle-adjusting neural radiance fields,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5741–5751.
- [97] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. B. Goldman, S. M. Seitz, and R. Martin-Brualla, “Nerfies: Deformable neural radiance fields,” in *International Conference on Computer Vision*, 2021.
- [98] A. Hertz, O. Perel, R. Giryes, O. Sorkine-Hornung, and D. Cohen-Or, “Sape: Spatially-adaptive progressive encoding for neural optimization,” *Advances in Neural Information Processing Systems*, 2021.
- [99] L. Melas-Kyriazi, I. Laina, C. Rupprecht, and A. Vedaldi, “Realfusion: 360deg reconstruction of any object from a single image,” in *Conference on computer vision and pattern recognition*, 2023.
- [100] H. Heo, T. Kim, J. Lee, J. Lee, S. Kim, H. J. Kim, and J.-H. Kim, “Robust camera pose refinement for multi-resolution hash encoding,” in *International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research.
- [101] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar, “Block-nerf: Scalable large scene neural view synthesis,” in *CVPR*, 2022.
- [102] Q. Meng, A. Chen, H. Luo, M. Wu, H. Su, L. Xu, X. He, and J. Yu, “Gnerf: Gan-based neural radiance field without posed camera,” in *International Conference on Computer Vision*, 2021.
- [103] S.-F. Chng, S. Ramasinghe, J. Sherrah, and S. Lucey, “Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation,” in *European Conference on Computer Vision*, 2022.
- [104] Y. Xia, H. Tang, R. Timofte, and L. Van Gool, “Sinerf: Sinusoidal neural radiance fields for joint pose estimation and scene reconstruction,” *arXiv preprint arXiv:2210.04553*, 2022.
- [105] W. Bian, Z. Wang, K. Li, J.-W. Bian, and V. A. Prisacariu, “Nope-nerf: Optimising neural radiance field with no pose prior,” in *Conference on Computer Vision and Pattern Recognition*, 2023.
- [106] P. Truong, M.-J. Rakotosaona, F. Manhardt, and F. Tombari, “Sparf: Neural radiance fields from sparse and noisy poses,” in *Conference on Computer Vision and Pattern Recognition*, 2023.
- [107] Y. Jeong, S. Ahn, C. Choy, A. Anandkumar, M. Cho, and J. Park, “Self-calibrating neural radiance fields,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5846–5854.
- [108] Y. Chen and G. H. Lee, “Dbarf: Deep bundle-adjusting generalizable neural radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [109] Q. Wang, Z. Wang, K. Genova, P. P. Srinivasan, H. Zhou, J. T. Barron, R. Martin-Brualla, N. Snavely, and T. Funkhouser, “Ibrnet: Learning multi-view image-based rendering,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [110] G. Avraham, J. Straub, T. Shen, T.-Y. Yang, H. Germain, C. Sweeney, V. Balntas, D. Novotny, D. DeTone, and R. Newcombe, “Nerfels: renderable neural codes for improved camera pose estimation,” in *Conference on Computer Vision and Pattern Recognition*, 2022.
- [111] Q. Zhou, M. Maximov, O. Litany, and L. Leal-Taixé, “The nerfect match: Exploring nerf features for visual localization,” *arXiv preprint*

- arXiv:2403.09577*, 2024.
- [112] R. Chen, Y. Cong, and Y. Ren, “Marrying nerf with feature matching for one-step pose estimation,” *arXiv preprint arXiv:2404.00891*, 2024.
- [113] M. Bortolon, T. Tsesmelis, S. James, F. Poesi, and A. Del Bue, “Iffnerf: Initialisation free and fast 6dof pose estimation from a single image and a nerf model,” *arXiv preprint arXiv:2403.12682*, 2024.
- [114] S. Ito, H. Aizawa, and K. Kato, “Few-shot nerf-based view synthesis for viewpoint-biased camera pose estimation,” in *International Conference on Artificial Neural Networks*. Springer, 2023, pp. 308–319.
- [115] L. Claessens, F. Manhardt, R. Martin-Brualla, R. Siegwart, C. Cadena, and F. Tombari, “Robust and efficient edge-guided pose estimation with resolution-conditioned nerf,” in *BMVC*, 2023.
- [116] M. A. Karaoğlu, H. Schieber, N. Schischka, M. Görgülü, F. Grötzner, A. Ladikos, D. Roth, N. Navab, and B. Busam, “Dynamon: Motion-aware fast and robust camera localization for dynamic nerf,” *arXiv preprint arXiv:2309.08927*, 2023.
- [117] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, “imap: Implicit mapping and positioning in real-time,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6229–6238.
- [118] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, “Nice-slam: Neural implicit scalable encoding for slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 786–12 796.
- [119] H. Matsuki, R. Murai, P. H. J. Kelly, and A. J. Davison, “Gaussian Splatting SLAM,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [120] V. Yugay, Y. Li, T. Gevers, and M. R. Oswald, “Gaussian-slam: Photo-realistic dense slam with gaussian splatting,” 2023.
- [121] N. Keetha, J. Karhade, K. M. Jatavallabhula, G. Yang, S. Scherer, D. Ramanan, and J. Luiten, “Splatam: Splat, track & map 3d gaussians for dense rgb-d slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [122] S. Zhi, E. Sucar, A. Mouton, I. Haughton, T. Laidlow, and A. J. Davison, “iLabel: Interactive neural scene labelling,” in *arXiv*, 2021.
- [123] K. Mazur, E. Sucar, and A. J. Davison, “Feature-realistic neural fusion for real-time, open set scene understanding,” in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [124] K. Jatavallabhula, A. Kuwajerwala, Q. Gu, M. Omama, T. Chen, S. Li, G. Iyer, S. Saryazdi, N. Keetha, A. Tewari, J. Tenenbaum, C. de Melo, M. Krishna, L. Paull, F. Shkurti, and A. Torralba, “Conceptfusion: Open-set multimodal 3d mapping,” in *RSS*, 2023.
- [125] F. Tosi, Y. Zhang, Z. Gong, E. Sandström, S. Mattoccia, M. R. Oswald, and M. Poggi, “How nerfs and 3d gaussian splatting are reshaping slam: a survey,” *arXiv preprint arXiv:2402.13255*, 2024.
- [126] M. Z. Irshad, S. Zakharov, V. Guizilini, A. Gaidon, Z. Kira, and R. Ambrus, “Nerf-mae: Masked autoencoders for self-supervised 3d representation learning for neural radiance fields,” in *European Conference on Computer Vision (ECCV)*, 2024.
- [127] H. Yan, Y. Zheng, and Y. Duan, “Gaussian-det: Learning closed-surface gaussians for 3d object detection,” 2024.
- [128] Y. Cao, Y. Jv, and D. Xu, “3dgs-det: Empower 3d gaussian splatting with boundary guidance and box-focused sampling for 3d object detection,” 2024.
- [129] F. Li, S. R. Vutukur, H. Yu, I. Shugurov, B. Busam, S. Yang, and S. Ilic, “Nerf-pose: A first-reconstruct-then-regress approach for weakly-supervised 6d object pose estimation,” in *International Conference on Computer Vision*, 2023.
- [130] M. Z. Irshad, S. Zakharov, R. Ambrus, T. Kollar, Z. Kira, and A. Gaidon, “Shapo: Implicit representations for multi object shape appearance and pose optimization,” in *European Conference on Computer Vision (ECCV)*, 2022.
- [131] M. Lunayach, S. Zakharov, D. Chen, R. Ambrus, Z. Kira, and M. Z. Irshad, “Fsd: Fast self-supervised single rgb-d to categorical 3d objects,” in *Int. Conf. on Robotics and Automation*. IEEE, 2024.
- [132] N. Heppert, M. Z. Irshad, S. Zakharov, K. Liu, R. A. Ambrus, J. Bohg, A. Valada, and T. Kollar, “Carto: Category and joint agnostic reconstruction of articulated objects,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2023, pp. 21 201–21 210.
- [133] Y. Guo, A. Kumar, C. Zhao, R. Wang, X. Huang, and L. Ren, “Upnerf: A unified framework for monocular 3d object reconstruction and pose estimation,” *arXiv preprint arXiv:2403.15705*, 2024.
- [134] D. Pavllo, D. J. Tan, M.-J. Rakotosaona, and F. Tombari, “Shape, pose, and appearance from a single image via bootstrapped radiance field inversion,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [135] L. Huang, T. Hodan, L. Ma, L. Zhang, L. Tran, C. Twigg, P.-C. Wu, J. Yuan, C. Keskin, and R. Wang, “Neural correspondence field for object pose estimation,” in *European Conference on Computer Vision*, 2022.
- [136] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, and S. K. et al., “Droid: A large-scale in-the-wild robot manipulation dataset,” 2024.
- [137] M. Deitke, R. Liu, M. Wallingford, H. Ngo, O. Michel, A. Kusupati, A. Fan, C. Laforte, V. Voleti, S. Y. Gadre et al., “Objaverse-xl: A universe of 10m+ 3d objects,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [138] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. Goldberg, “Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics,” *arXiv preprint arXiv:1703.09312*, 2017.
- [139] S. Song, A. Zeng, J. Lee, and T. Funkhouser, “Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations,” *IEEE Robotics and Automation Letters*, 2020.
- [140] L. Yen-Chen, P. Florence, J. T. Barron, T.-Y. Lin, A. Rodriguez, and P. Isola, “Nerf-supervision: Learning dense object descriptors from neural radiance fields,” in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 6496–6503.
- [141] L. Yen-Chen, P. Florence, A. Zeng, J. T. Barron, Y. Du, W.-C. Ma, A. Simeonov, A. R. Garcia, and P. Isola, “Mira: Mental imagery for robotic affordances,” *arXiv preprint arXiv:2212.06088*, 2022.
- [142] G. Sótí, X. Huang, C. Wurrll, and B. Hein, “6-dof grasp pose evaluation and optimization via transfer learning from nerfs,” *arXiv preprint arXiv:2401.07935*, 2024.
- [143] C. Liu, K. Shi, K. Zhou, H. Wang, J. Zhang, and H. Dong, “Rgbgrasp: Image-based object grasping by capturing multiple views during robot arm movement with neural radiance fields,” *arXiv preprint arXiv:2311.16592*, 2023.
- [144] V. Blukis, T. Lee, J. Tremblay, B. Wen, I. S. Kweon, K.-J. Yoon, D. Fox, and S. Birchfield, “Neural fields for robotic object manipulation from a single image,” *arXiv preprint arXiv:2210.12126*, 2022.
- [145] V. Blukis, K.-J. Yoon, T. Lee, J. Tremblay, B. Wen, I.-S. Kweon, D. Fox, and S. Birchfield, “One-shot neural fields for 3d object understanding,” in *CVPR Workshop (CVPRW)*, 2023.
- [146] C. Deng, O. Litany, Y. Duan, A. Poulernard, A. Tagliasacchi, and L. J. Guibas, “Vector neurons: A general framework for so (3)-equivariant networks,” in *International Conference on Computer Vision*, 2021.
- [147] M. N. Qureshi, S. Garg, F. Yandun, D. Held, G. Kantor, and A. Silwal, “SplatSim: Zero-shot sim2real transfer of rgb manipulation policies using gaussian splatting,” *arXiv preprint arXiv:2409.10161*, 2024.
- [148] E. Chisari, N. Heppert, T. Welschehold, W. Burgard, and A. Valada, “Centergrasp: Object-aware implicit representation learning for simultaneous shape reconstruction and 6-dof grasp estimation,” *IEEE Robotics and Automation Letters (RA-L)*, 2024.
- [149] R.-Z. Qiu, Y. Hu, G. Yang, Y. Song, Y. Fu, J. Ye, J. Mu, R. Yang, N. Atanasov, S. Scherer et al., “Learning generalizable feature fields for mobile manipulation,” *arXiv preprint arXiv:2403.07563*, 2024.
- [150] S. Zhong, A. Albin, O. P. Jones, P. Maiolino, and I. Posner, “Touching a nerf: Leveraging neural radiance fields for tactile sensory data generation,” in *Conference on Robot Learning*. PMLR, 2023.
- [151] Y. Dou, F. Yang, Y. Liu, A. Loquercio, and A. Owens, “Tactile-augmented radiance fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [152] M. Comi, Y. Lin, A. Church, A. Tonioni, L. Aitchison, and N. F. Lepora, “Touchsd: A deepsd approach for 3d shape reconstruction using vision-based tactile sensing,” *IEEE Robotics and Automation Letters*, 2024.
- [153] S. Suresh, H. Qi, T. Wu, T. Fan, L. Pineda, M. Lambeta, J. Malik, M. Kalakrishnan, R. Calandra, M. Kaess et al., “Neural feels with neural fields: Visuo-tactile perception for in-hand manipulation,” *arXiv preprint arXiv:2312.13469*, 2023.
- [154] C. Higuera, S. Dong, B. Boots, and M. Mukadam, “Neural contact fields: Tracking extrinsic contact with tactile sensing,” in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [155] T. Yoneda, T. Jiang, G. Shakhnarovich, and M. R. Walter, “6-dof stability field via diffusion models,” *arXiv preprint arXiv:2310.17649*, 2023.
- [156] C.-Y. Ma, Z. Wu, G. AIRegib, C. Xiong, and Z. Kira, “The regretful agent: Heuristic-aided navigation through progress estimation,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [157] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Yang Wang, and L. Zhang, “Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation,” in *Conference on Computer Vision and Pattern Recognition*, 2019.
- [158] E. Wijmans, A. Kadian, A. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra, “DD-PPO: Learning near-perfect pointgoal

- navigators from 2.5 billion frames,” *International Conference on Learning Representations (ICLR)*, 2020.
- [159] M. Z. Irshad, C.-Y. Ma, and Z. Kira, “Hierarchical cross-modal agent for robotics vision-and-language navigation,” in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [160] M. Z. Irshad, N. C. Mithun, Z. Seymour, H.-P. Chiu, S. Samarasekera, and R. Kumar, “Sasra: Semantically-aware spatio-temporal reasoning agent for vision-and-language navigation in continuous environments,” in *International Conference on Pattern Recognition (ICPR)*, 2022.
- [161] D. S. Chaplot, H. Jiang, S. Gupta, and A. Gupta, “Semantic curiosity for active visual learning,” 2020.
- [162] M. Adamkiewicz, T. Chen, A. Caccavale, R. Gardner, P. Culbertson, J. Bohg, and M. Schwager, “Vision-only robot navigation in a neural radiance world,” *IEEE Robotics and Automation Letters*, 2022.
- [163] D. Driess, Z. Huang, Y. Li, R. Tedrake, and M. Toussaint, “Learning multi-object dynamics with compositional neural radiance fields,” in *Conference on Robot Learning*. PMLR, 2023, pp. 1755–1768.
- [164] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, “A comprehensive survey on graph neural networks,” *IEEE transactions on neural networks and learning systems*, 2020.
- [165] D. S. Chaplot, D. P. Gandhi, A. Gupta, and R. R. Salakhutdinov, “Object goal navigation using goal-oriented semantic exploration,” *Advances in Neural Information Processing Systems*, 2020.
- [166] D. Patel, P. Pham, and A. Bera, “Dronerf: Real-time multi-agent drone pose optimization for computing neural radiance fields,” *arXiv preprint arXiv:2303.04322*, 2023.
- [167] C. Li, R. Liang, H. Fan, Z. Zhang, S. Durvasula, and N. Vijaykumar, “Disorf: A distributed online nerf training and rendering framework for mobile robots,” *arXiv preprint arXiv:2403.00228*, 2024.
- [168] E. Skartados, M. K. Yucel, B. Manganeli, A. Drosou, and A. Saà-Garriga, “Finding waldo: Towards efficient exploration of nerf scene spaces,” in *ACM Multimedia Systems Conference*, 2024.
- [169] B. Planche, X. Rong, Z. Wu, S. Karanam, H. Kosch, Y. Tian, J. Ernst, and A. Hutter, “Incremental scene synthesis,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [170] F. Taioli, F. Cunico, F. Girella, R. Bologna, A. Farinelli, and M. Cristani, “Language-enhanced nrr-map: Querying renderable neural radiance field maps with natural language,” in *International Conference on Computer Vision (ICCV) Workshops*, 2023.
- [171] J. Han, L. L. Beyer, G. V. Cavalheiro, and S. Karaman, “Nvins: Robust visual inertial navigation fused with nerf-augmented camera pose regressor and uncertainty quantification,” *arXiv preprint arXiv:2404.01400*, 2024.
- [172] Q. Liu, N. Chen, Z. Liu, and H. Wang, “Toward learning-based visuomotor navigation with neural radiance fields,” *IEEE Transactions on Industrial Informatics*, 2024.
- [173] A. Rashid, C. M. Kim, J. Kerr, L. Fu, K. Hari, A. Ahmad, K. Chen, H. Huang, M. Gualtieri, M. Wang *et al.*, “Lifelong lrf: Local 3d semantic inventory monitoring using fogros2,” *arXiv preprint arXiv:2403.10494*, 2024.
- [174] B. Zhao, L. Yang, M. Mao, H. Bao, and Z. Cui, “Pnerfloc: Visual localization with point-based neural radiance fields,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, 2024, pp. 7450–7459.
- [175] G. S. Camps, R. Dyro, M. Pavone, and M. Schwager, “Learning deep sdf maps online for robot navigation and exploration,” *arXiv preprint arXiv:2207.10782*, 2022.
- [176] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, “Segment anything,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4015–4026.
- [177] C. Huang, O. Mees, A. Zeng, and W. Burgard, “Visual language maps for robot navigation,” in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023.
- [178] B. Li, K. Q. Weinberger, S. Belongie, V. Koltun, and R. Ranftl, “Language-driven semantic segmentation,” *arXiv preprint arXiv:2201.03546*, 2022.
- [179] X. Zhou, R. Girdhar, A. Joulin, P. Krähenbühl, and I. Misra, “Detecting twenty-thousand classes using image-level supervision,” in *European Conference on Computer Vision*. Springer, 2022, pp. 350–368.
- [180] Y. Yuan and A. Nüchter, “Uni-fusion: Universal continuous mapping,” *IEEE Transactions on Robotics*, 2024.
- [181] K. Yamazaki, T. Hanyu, K. Vo, T. Pham, M. Tran, G. Doretto, A. Nguyen, and N. Le, “Open-fusion: Real-time open-vocabulary 3d mapping and queryable scene representation,” *arXiv preprint arXiv:2310.03923*, 2023.
- [182] Y. Hong, H. Zhen, P. Chen, S. Zheng, Y. Du, Z. Chen, and C. Gan, “3d-llm: Injecting the 3d world into large language models,” 2023.
- [183] Y. Hong, C. Lin, Y. Du, Z. Chen, J. B. Tenenbaum, and C. Gan, “3d concept learning and reasoning from multi-view images,” *Conference on Computer Vision and Pattern Recognition*, 2023.
- [184] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, “Flamingo: a visual language model for few-shot learning,” *Neurips*, 2022.
- [185] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” in *NeurIPS*, 2023.
- [186] X. Puig, E. Undersander, A. Szot, M. D. Cote, T.-Y. Yang, R. Partsey, R. Desai, A. W. Clegg, M. Hlavac, S. Y. Min *et al.*, “Habitat 3.0: A co-habitat for humans, avatars and robots,” *arXiv preprint arXiv:2310.13724*, 2023.
- [187] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, “Diffusion models in vision: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [188] J.-W. Liu, Y.-P. Cao, W. Mao, W. Zhang, D. J. Zhang, J. Keppo, Y. Shan, X. Qie, and M. Z. Shou, “Devrf: Fast deformable voxel radiance fields for dynamic scenes,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 36762–36775, 2022.
- [189] K. Park, U. Sinha, P. Hedman, J. T. Barron, S. Bouaziz, D. B. Goldman, R. Martin-Brualla, and S. M. Seitz, “Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields,” *arXiv preprint arXiv:2106.13228*, 2021.
- [190] E. Tretschk, A. Tewari, V. Golyanik, M. Zollhöfer, C. Lassner, and C. Theobalt, “Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video,” in *International Conference on Computer Vision*, 2021.
- [191] Y.-L. Qiao, A. Gao, and M. Lin, “Neuphysics: Editable neural geometry and physics from monocular videos,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 12841–12854, 2022.
- [192] Z. Li, S. Niklaus, N. Snavely, and O. Wang, “Neural scene flow fields for space-time view synthesis of dynamic scenes,” in *Conference on Computer Vision and Pattern Recognition*, 2021.
- [193] C. Gao, A. Saraf, J. Kopf, and J.-B. Huang, “Dynamic view synthesis from dynamic monocular video,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5712–5721.
- [194] G. Wu, T. Yi, J. Fang, L. Xie, X. Zhang, W. Wei, W. Liu, Q. Tian, and X. Wang, “4d gaussian splatting for real-time dynamic scene rendering,” *arXiv preprint arXiv:2310.08528*, 2023.
- [195] B. P. Duisterhof, Z. Mandi, Y. Yao, J.-W. Liu, M. Z. Shou, S. Song, and J. Ichnowski, “Md-splatting: Learning metric deformation from 4d gaussians in highly deformable scenes,” *arXiv preprint arXiv:2312.00583*, 2023.
- [196] Z. Yang, H. Yang, Z. Pan, and L. Zhang, “Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting,” in *International Conference on Learning Representations (ICLR)*, 2024.
- [197] I. Liu, H. Su, and X. Wang, “Dynamic gaussians mesh: Consistent mesh reconstruction from monocular videos,” 2024.
- [198] J. Abou-Chakra, F. Dayoub, and N. Sünderhauf, “Particleclerf: A particle-based encoding for online neural radiance fields,” in *Winter Conference on Applications of Computer Vision*, 2024.
- [199] J. Abou-Chakra, K. Rana, F. Dayoub, and N. Sünderhauf, “Physically embodied gaussian splatting: A realtime correctable world model for robotics,” *arXiv preprint arXiv:2406.10788*, 2024.
- [200] A. Blattmann, T. Dockhorn, S. Kulal, D. Mendelevitch, M. Kilian, D. Lorenz, Y. Levi, Z. English, V. Voleti, A. Letts *et al.*, “Stable video diffusion: Scaling latent video diffusion models to large datasets,” *arXiv preprint arXiv:2311.15127*, 2023.
- [201] F. Hoffer, L. Koestler, F. Bernard, and D. Cremers, “Neural implicit representations for physical parameter inference from a single video,” in *Winter Conference on Applications of Computer Vision*, 2023.
- [202] S. Le Cleac’h, H.-X. Yu, M. Guo, T. Howell, R. Gao, J. Wu, Z. Manchester, and M. Schwager, “Differentiable physics simulation of dynamics-augmented neural objects,” *IEEE Robotics and Automation Letters*, vol. 8, no. 5, pp. 2780–2787, 2023.
- [203] A. J. Zhai, Y. Shen, E. Y. Chen, G. X. Wang, X. Wang, S. Wang, K. Guan, and S. Wang, “Physical property understanding from language-embedded feature fields,” in *CVPR*, 2024.
- [204] K. Yang, X. Zhang, Z. Huang, X. Chen, Z. Xu, and H. Su, “Moving-parts: Motion-based 3d part discovery in dynamic radiance field,” in *International Conference on Learning Representations*, 2024.
- [205] W.-C. Tseng, H.-J. Liao, L. Yen-Chen, and M. Sun, “CLA-NerF: Category-Level Articulated Neural Radiance Field,” Mar. 2022, arXiv:2202.00181 [cs].
- [206] J. Kerr, C. M. Kim, M. Wu, B. Yi, Q. Wang, K. Goldberg, and A. Kanazawa, “Robot see robot do: Imitating articulated object

- manipulation with monocular 4d reconstruction,” in *Conference on Robot Learning*, 2024.
- [207] Y. Feng, Y. Shang, X. Li, T. Shao, C. Jiang, and Y. Yang, “Pie-nerf: Physics-based interactive elastodynamics with nerf,” *arXiv preprint arXiv:2311.13099*, 2023.
- [208] L. Zhong, H.-X. Yu, J. Wu, and Y. Li, “Reconstruction and simulation of elastic objects with spring-mass 3d gaussians,” *arXiv preprint arXiv:2403.09434*, 2024.
- [209] H.-X. Yu, Y. Zheng, Y. Gao, Y. Deng, B. Zhu, and J. Wu, “Inferring hybrid neural fluid fields from videos,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [210] Y. Feng, X. Feng, Y. Shang, Y. Jiang, C. Yu, Z. Zong, T. Shao, H. Wu, K. Zhou, C. Jiang, and Y. Yang, “Gaussian splashing: Dynamic fluid synthesis with gaussian splatting,” 2024.
- [211] Y. Li, Z.-H. Lin, D. Forsyth, J.-B. Huang, and S. Wang, “Climatenerf: Extreme weather synthesis in neural radiance field,” in *International Conference on Computer Vision (ICCV)*, 2023.
- [212] Z. Zong, X. Li, M. Li, M. M. Chiaramonte, W. Matusik, E. Grinspun, K. Carlberg, C. Jiang, and P. Y. Chen, “Neural stress fields for reduced-order elastoplasticity and fracture,” in *SIGGRAPH Asia*, 2023.
- [213] Y. Yan, H. Lin, C. Zhou, W. Wang, H. Sun, K. Zhan, X. Lang, X. Zhou, and S. Peng, “Street gaussians for modeling dynamic urban scenes,” 2024.
- [214] H. Turki, J. Y. Zhang, F. Ferroni, and D. Ramanan, “Suds: Scalable urban dynamic scenes,” 2023.
- [215] J. Yang, B. Ivanovic, O. Litany, X. Weng, S. W. Kim, B. Li, T. Che, D. Xu, S. Fidler, M. Pavone, and Y. Wang, “Emernerf: Emergent spatial-temporal scene decomposition via self-supervision,” *arXiv preprint arXiv:2311.02077*, 2023.
- [216] X. Zhou, Z. Lin, X. Shan, Y. Wang, D. Sun, and M.-H. Yang, “Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes,” 2023.
- [217] R. Martin-Brualla, N. Radwan, M. S. M. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, “NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections,” in *CVPR*, 2021, pp. 7210–7219.
- [218] M.-Y. Shen, C.-C. Hsu, H.-Y. Hou, Y.-C. Huang, W.-F. Sun, C.-C. Chang, Y.-L. Liu, and C.-Y. Lee, “Driveenv-nerf: Exploration of a nerf-based autonomous driving environment for real-world performance validation,” *arXiv preprint arXiv:2403.15791*, 2024.
- [219] C. Lindström, G. Hess, A. Lilja, M. Fatemi, L. Hammarstrand, C. Petersson, and L. Svensson, “Are nerfs ready for autonomous driving? towards closing the real-to-simulation gap,” *arXiv preprint arXiv:2403.16092*, 2024.
- [220] L. Li, Q. Lian, L. Wang, N. Ma, and Y.-C. Chen, “Lift3d: Synthesize 3d training data by lifting 2d gan to 3d generative radiance field,” in *Conference on Computer Vision and Pattern Recognition*, 2023.
- [221] L. Li, Q. Lian, and Y.-C. Chen, “Adv3d: generating 3d adversarial examples in driving scenarios with nerf,” *arXiv preprint arXiv:2309.01351*, 2023.