

Semantically-aware Neural Radiance Fields for Visual Scene Understanding: A Comprehensive Review

Thang-Anh-Quan Nguyen^{1,*}, Amine Bourki^{2,*}, Mátyás Macudzinski³,
Anthony Brunel⁴, Mohammed Bennamoun⁵

¹Noah's Ark, Huawei Paris Research Center, France.

²Inception Lab, Paris, France.

³Centrale Lille, University of Lille, France.

⁴Nice, France.

⁵The University of Western Australia, Perth, Australia.

*: Denotes co-primary authorship.

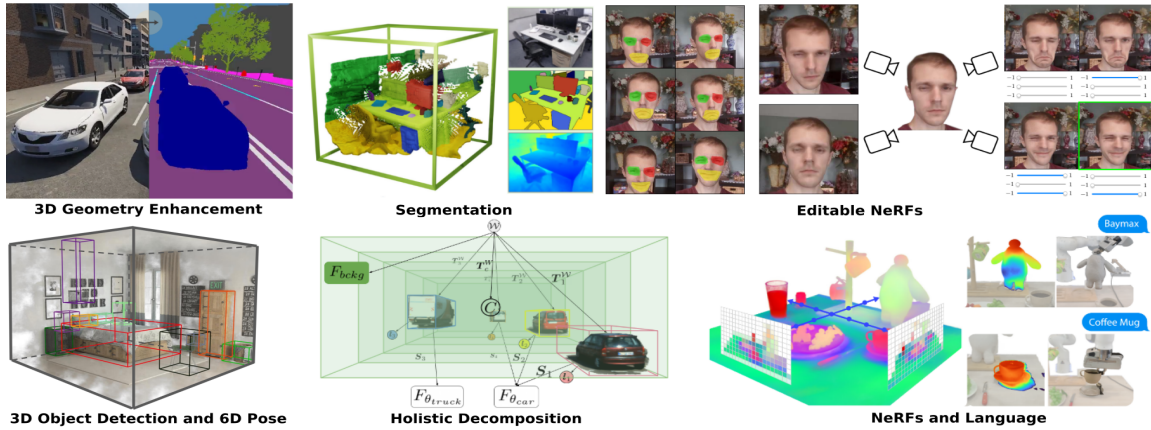


Fig. 1: Our survey covers 250+ papers on Neural Radiance Fields with semantic scene understanding capabilities, spanning the 6 main categories depicted above. Illustrations from [81, 195, 94, 72, 160, 192].

Abstract

This review thoroughly examines the role of semantically-aware Neural Radiance Fields (NeRFs) in visual scene understanding, covering an analysis of over 250 scholarly papers. It explores how NeRFs adeptly infer 3D representations for both stationary and dynamic objects in a scene. This capability is pivotal for generating high-quality new viewpoints, completing missing scene details (inpainting), conducting comprehensive scene segmentation (panoptic segmentation), predicting 3D bounding boxes, editing 3D scenes, and extracting object-centric 3D models. A significant aspect of this study is the application of semantic labels as viewpoint-invariant functions, which effectively map spatial coordinates to a spectrum of semantic labels, thus facilitating the recognition of distinct objects within the scene. Overall, this survey highlights the progression and diverse applications of semantically-aware neural radiance fields in the context of visual scene interpretation.

Keywords: Neural Radiance Fields, NeRFs, Visual Scene Understanding, 3D Scene Representation, Generative AI, Literature Survey.

1 Introduction

Neural Radiance Fields (NeRFs) have marked a significant development since their inception [145], offering unprecedented capabilities in synthesizing photorealistic unseen views from a set of 2D images through a new 3D scene representation. The core strength of NeRFs lies in their ability to intricately model the complex interactions of light within a scene, thereby generating 3D representations that are both detailed and realistic. Traditional NeRFs, however, primarily focus on geometric and photometric accuracy, and often overlook the underlying semantics of the observed scenes.

The advent of semantically-aware NeRFs (SRFs) marks a significant advancement in this domain. These models not only capture the physical characteristics of a scene, but also incorporate an understanding of semantic and contextual information. This leap in technology facilitates a range of sophisticated applications, such as scene editing, improved object recognition, and more interactive and realistic virtual environments.

Recent developments in implicit neural rendering have also been pivotal. These methods demonstrate the possibility of learning accurate view synthesis for complex scenes by predicting their volumetric density and color, using only a set of RGB images as supervision. Despite these advances, most of the existing methods are limited to static scenes. They tend to encode all scene objects into a single neural network, thus falling short of representing dynamic scenes or breaking down scenes into individual objects. This limitation is a significant hurdle on the path towards creating more dynamic and responsive 3D environments.

Visual Scene Understanding, which is often categorized into the three R's of Computer Vision [139], namely Reconstruction, Recognition, and Re-organization (i.e., bottom-up segmentation), has received a massive attention in the field, both as individual problems as well as joint multitask approaches aiming to take advantage of their inherent mutually-informative nature, leading to improved performance and efficiency [2, 224, 203, 42, 272]. Similarly, the conventional methods which sequentially solve for NeRF then perception not only introduce additional computational costs and inefficiencies for perception tasks, but

also fail to fully leverage the mutually-beneficial potential of volumetric renderings w.r.t perception during the training phase. This gap represents a missed opportunity to maximize the synergies between 3D scene reconstruction and semantic perception.

Our comprehensive survey explores deep into these aspects, exploring the most recent advancements in semantically-aware NeRFs. We examine how the integration of semantic information can substantially enhance the capabilities of NeRFs, especially in complex and dynamic environments. Our discussion covers various methodologies for integrating semantic data into radiance fields, the challenges inherent in these processes, and the vast potential applications of these enriched models across diverse domains.

Positioning and Impact. The ultimate goal of this paper is to provide a thorough understanding of the current state and potential of semantically-aware NeRFs. We aim to identify gaps in existing methodologies, highlight the challenges yet to be overcome, and offer a vision for future research directions. To our knowledge, *this survey is the first in the field to specifically concentrate on semantic coupling in Neural Radiance Fields*. This is significant considering the growing interest in this area of study.

1.1 Prior Related Surveys

This section focuses on previously conducted surveys that have explored Neural Radiance Fields (NeRFs), with a particular emphasis on semantic scene understanding from different perspectives, including 2D, 2.5D, and multi-view imaging techniques. These surveys have laid the groundwork in the field and provide insight into the development, capabilities, and limitations of NeRFs in processing and interpreting complex visual data.

Our survey expands upon this existing knowledge base by considering a wide range of venues and studies within a specific timeframe, offering a contemporary snapshot of the advancements in the field of semantically aware NeRFs. We not only review the findings and methodologies of these prior surveys but also highlight how our survey stands out in its approach and focus.

In particular, our survey explores how recent advancements in NeRF technology have been tailored to enhance semantic scene understanding.

Survey	Venue	Sem. Tasks	Sem. Focus	Strengths	Limitations
Xia <i>et al.</i> [244]	CS 2023	..E...	✗	Extensively covers 3D-aware image synthesis, implicit scene representations, differentiable neural rendering in NeRF-based methods.	Broad scope in image synthesis with little attention to semantics . Includes generation and editability. Study limited to December 2022 , before most of competitive SRF methods were published.
Zhu <i>et al.</i> [288]	APSIPA 2023	G.E...	✗	Provides synthetic overview of the field on a wide variety of topics with emphasis of current limitations of generic approaches.	Broad scope with very superficial attention to semantic aspects . Study limited to December 2022 , before most of competitive SRF methods were published.
Tewari <i>et al.</i> [216]	CGF 2022	G.E.H.	✗	Comprehensively explains neural rendering and scene representation techniques, challenges and improvement strategies. Discusses editable and NeRFs and compositionality of its time. Extensive discussion on challenges and perspectives.	Broad scope with very superficial attention to semantic aspects (handful of papers), focuses on graphics and rendering related geometric aspects. Study limited to 2021 papers before most of competitive SRF methods were published.
Xie <i>et al.</i> [249]	CGF 2022	G.E.H.	✗	In depth review with a strong tutorial on theoretical aspects of differentiable rendering, NeRFs, and scene representations.	Broad scope focusing on appearance, textures and relighting applications with very superficial attention to semantic aspects . Study limited to very early 2022 , before most competitive SRF methods were published.
Mittal <i>et al.</i> [151]*	arXiv 2023	GSEOHL	✗	Emphasizes on the basics. Lists out short abstracts of 500 papers and pre-prints. Organized by loss functions, applications, publication year. Up-to-date and gets incremental updates for each new NeRF publication.	It is a 400+ page unpublished tech report with comprehensive tutorials, not a journal publication . Mainly a list of abstract summaries, lacks analysis, discussions, and global insights . It is a constantly updated work in progress .
Li <i>et al.</i> [112]	arXiv 2023	..E..L	✗	In-depth review of text-guided and text-controlled strategies, with applications on the generation of avatars, textures, scenes, and shapes.	Narrow scope restricted to TextTo3D generative methods and partially covers editability. Study limited to May 2023. Only addresses a small fraction of our scope in time and applications .
Rabby <i>et al.</i> [169]	arXiv 2023	G.E.H.	✗	Includes discussions on compositionality, scene editing, and public datasets. Provides a centralized benchmark.	Broad scope with very superficial attention to semantics . Study limited to 2022 papers, before most of competitive SRF methods were published .
Slapak <i>et al.</i> [196]	arXiv 2023	G...H.	✗	Good general overview on geometric approaches relevant to industrial and robotics fields. Discusses efficiency and effectiveness improvements of traditional NeRFs.	Relatively short survey. Narrow scope restricted to industrial and robotic applications . Only partially and superficially covers semantic aspects G and H . Study limited to early 2023, missing many references from CVPR 2023 onwards.
Gao <i>et al.</i> [51]	arXiv 2022	GSE...	✗	Recent 2023 scope. Good generalist overview. Covers public datasets and evaluation.	Broad scope with very superficial overview of a few SRFs (handful of papers) . Focus is diluted among many geometric and efficiency oriented non-semantic discussions, leaving modest text real estate for semantics.
Ours	-	GSEOHL	✓	Comprehensive analysis on semantics, up to submission date in January 2024.	Purposefully focuses on semantics.

Table 1: Comparative overview of previously existing NeRF surveys w.r.t semantics (SRFs). Semantic Tasks include: G: 3D Geometry Enhancement, S: Segmentation, E: Editable NeRFs, O: Object Detection and 6D Pose, H: Holistic Decomposition, L: NeRFs and Language, .: denotes missing task. Semantic Focus refers to whether the primary focus of the study is on semantics. *: Interesting reference, but not a journal paper.

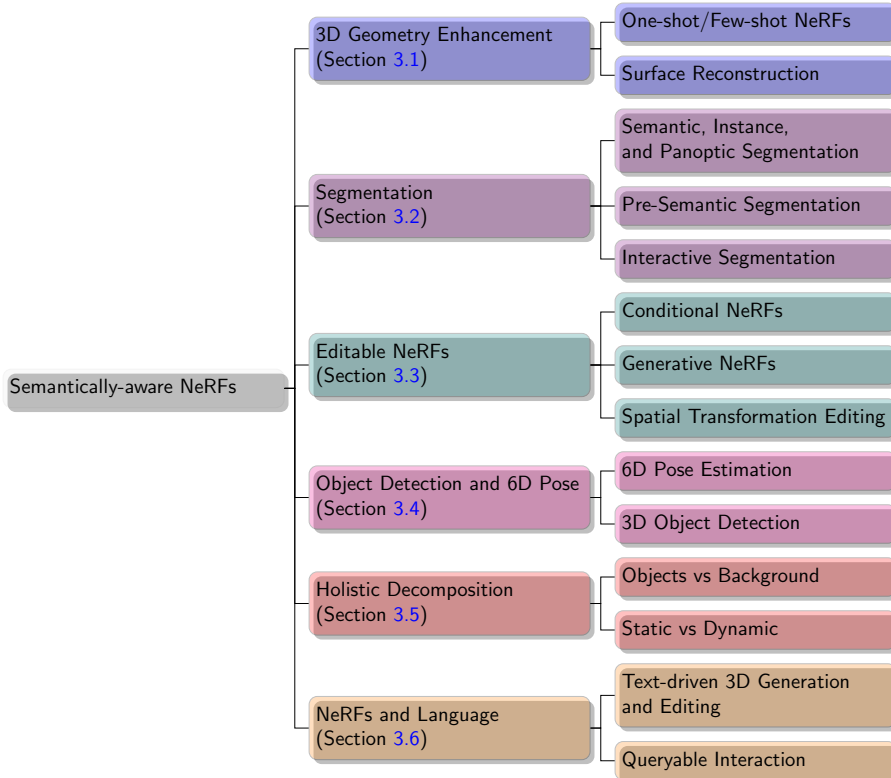


Fig. 2: Taxonomy of our study on Semantically-aware Neural Radiance Fields (SRFs).

This includes exploring how these advanced NeRF models interpret and interact with complex visual scenes, pushing the boundaries of what is possible in terms of visual perception and scene interpretation. We also discuss the methodological approaches these studies have adopted, providing an analytical framework that contrasts our approach with those of previous surveys.

More specifically, in the survey published in November 2023, Xia et al. [244] extensively cover the field of 3D-aware image synthesis, including detailed discussions on implicit scene representations such as occupancy fields, signed distance fields, and radiance fields, with a particular emphasis on NeRFs. It also examines differentiable neural rendering, underscoring its crucial role in fine-tuning neural networks for 3D rendering and highlighting the importance of volume rendering in NeRF-based methods. However, compared to our survey, certain limitations become apparent.

The broader scope of [244] contrasts with the more focused approach of our survey, which

explores deeply the integration of semantic understanding into NeRFs. Our narrower focus allows for a more comprehensive exploration of how semantic integration can enhance or extend NeRFs, especially in complex and dynamic environments, which is an aspect that may not be covered as thoroughly in [244].

Furthermore, our survey potentially offers richer insights into the practical applications and challenges associated with implementing semantically enhanced NeRFs. These practical considerations are underrepresented in [244]. In terms of future research directions, our survey provides more targeted guidance specific to semantic understanding in NeRFs, whereas [244] may present a wider range of future trends and research areas across the broader field of 3D-aware image synthesis.

In summary, while both surveys make significant contributions to the field of Computer Vision and 3D Image Synthesis, our survey stands out for its specialized and in-depth focus on the semantic aspects of Neural Radiance Fields, offering

nanced perspectives and insights that are particularly relevant to the advancement of semantic integration in this area.

In a tech report also from November 2023, Gao et al. [51] offer a broad overview of Neural Radiance Fields, discussing NeRF models, training requirements, various datasets used in research, and quality assessment metrics such as PSNR, SSIM, and LPIPS. However, compared to our survey, theirs shows limitations in its focus and depth. Our survey specifically concentrates on the integration of semantic understanding into NeRFs, providing detailed insights into semantic enhancement in complex and dynamic environments, practical applications, and specific future research directions. In contrast, [51] covers a wider range of topics in NeRF but lacks the specialized focus on semantic aspects, making our survey more comprehensive and targeted toward advancing semantic integration in Neural Radiance Fields.

In contrast to published surveys which typically cover NeRFs, as outlined in Table 1, our goal is to provide readers with a thorough understanding of semantically-aware NeRFs research. By comparing our approach with previous surveys, we aim to emphasize the distinct contributions and insights our study provides, especially regarding the incorporation and interpretation of semantic information within the Neural Radiance Fields framework.

1.2 Scope and Methodology

The papers referenced in this survey are predominantly published in the top venues for Computer Vision, Computer Graphics, Machine Learning, and Robotics. They cover the period from 2020 (first NeRF paper [145]) to the submission of the present paper in January 2024.

Our study primarily focuses on a taxonomy 6 main categories of semantically-aware NeRFs that define the underlying notion of semantics considered in this work, as summarized in Figure 2. First, we consider 3D geometry approaches that primarily use semantic information to improve performance in geometry-oriented tasks such as novel view synthesis and surface reconstruction. Specifically, in the very challenging setups of one-to-few shot scenarios, i.e., with a very limited number of input views, NeRF-based methods

can cope with the underconstrained nature of such challenging settings by leveraging higher-level information. In addition to ‘reconstruction’ applications, our study also includes segmentation which considers both the ‘recognition’ and ‘re-organization’ R’s of Visual Scene Understanding (resp. semantic and pre-semantic segmentations) [139]. Editable NeRFs allow to manipulate scenes through various priors and strategies. We also discuss works that enrich a radiance field formulation with 3D Object Detection or 6D Pose considerations. Holistic decomposition which aims at encoding the exhaustive structure of an input scene in a top-down manner. Lastly, we study language-rich NeRFs that enable new multi-modal applications for human interaction or effective scene manipulation.

Semantics in 3D visual computing applications has been thoroughly explored with often vastly differing definitions and considerations. In the context of this study, our intended definition of ‘semantics’ can be dissected into three main categories. **Initially**, we consider semantics as an explicit higher level construct to designate object- and / or instance-level labels [195], 3D object bounding boxes [72], 3D object 6D poses [184], or scene-wide decomposition [160, 119]. **Secondly**, as for Editable NeRFs and certain 3D Geometry Enhancement methods, we consider the use of semantics through the lens of language representation learning (e.g., [227]) which aim to describe scene objects with compact, controllable codes [85]. This is typically used in order to efficiently improve multi-view consistency, cope with missing views [83], or to enable object- or scene level manipulation [155]. **Finally**, we consider SRF strategies that explicitly interconnect vision and language, in order to generate novel 3D contents through text prompts [97], or enable higher-level scene interactions and user-guided manipulation [192].

Our work is also intended to help Computer Vision researchers unfamiliar with the topic step into semantically-aware NeRFs. Therefore, we cover the key concepts of the original NeRF architecture [145] and a classic extension for jointly considering semantic segmentation [285]. Additionally, we present a comprehensive overview of the relevant public datasets and evaluation tools. This additionally includes a centralized view of leading methods on these public benchmarks,

by grouping results that are initially scattered across dozens of referenced papers, and presenting original discussions and insights.

1.3 Organization of the Paper

The remainder of this paper is organized as follows. Section 2 discusses the key principles of the standard NeRF formulation and its extension to a basic semantic task, i.e., semantic segmentation. Section 3 provides an extensive literature review of SRFs, while Section 4 reviews the main public datasets, metrics, and evaluation tools commonly used in this field. Section 5 explores current challenges and perspectives, highlighting potential improvements in understanding semantic scenes and exploring real-world applications. Finally, Section 6 concludes the paper, giving higher-level perspectives for the field. We will also maintain an in-depth project repository on GitHub at github.com/aboutki/SoTA-Semantically-aware-NeRFs, including a comprehensive list of references, datasets, and performance evaluations, with regular updates to provide the latest state-of-the-art developments.

2 Fundamentals of Neural Radiance Fields

This section presents the core principles and terminology involved in the initial NeRF paper, as well as one of its simpler extensions to incorporate semantic reasoning capabilities. To do so, we cover how the 3D scene is represented along the formal definition of NeRFs, how they are employed to generate novel views. For more general or geometry-oriented details, we refer the interested reader to other existing surveys that emphasize more on such aspects than our study, which focuses on semantic considerations, e.g., [244, 216, 249, 51].

2.1 Scene Representation and Problem Statement

Neural Radiance Fields (NeRFs), introduced by Mildenhall et al. [145], have revolutionized the field of novel view synthesis. A NeRF model encapsulates a 3D scene through a radiance field, which is essentially a 5D function that

describes the light intensity traversing every direction within the scene. This is achieved by specifying both the color (as RGB values) and the volume density at each point in space. The core of a NeRF model lies in its ability to approximate this radiance function using Multi-Layer Perceptrons (MLPs). In the standard NeRF framework [145], a single MLP, denoted F_{Θ} , is used for this purpose, as follows:

$$(\mathbf{c}, \sigma) = F_{\Theta}(\mathbf{x}, \mathbf{d}) \quad (1)$$

where $\mathbf{x} = (x, y, z)$ is a given 3D point with x, y, z coordinates, $\mathbf{d} = (\theta, \phi)$ represents the viewing direction in Euler angles, $\mathbf{c} = (r, g, b)$ the color, and σ the corresponding volume density.

2.1.1 3D Scene Representation

Radiance fields are typically encoded using either of two different approaches to representing 3D scenes: implicit and explicit representations, making them implicit or explicit radiance fields respectively. When using an implicit scene representation, e.g., Signed Distance Functions (SDFs) or Deep Neural Networks (DNNs) in the case of NeRFs, the underlying geometry of the scene is not explicitly defined nor stored. It has to be retrieved using a post-processing or querying step, making then much more memory-efficient to the expense of additional computation.

Explicit Radiance Fields on the other hand rely on a data structure that explicitly defines the scene geometry as, e.g., point clouds [78], voxel grids [207], or permutohedral lattices [183] that allow to store radiance information with faster access rates but often with scene resolution constraints linked to their superior memory complexity.

2.1.2 Volumetric Rendering

Volume rendering [92] is a technique used to compute the color $C(r)$ of any camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ where \mathbf{o} represents the camera position and \mathbf{d} is the viewing direction, given the volume density and color functions of the scene being rendered. The color $C(r)$ is given by:

$$C(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i \quad (2)$$

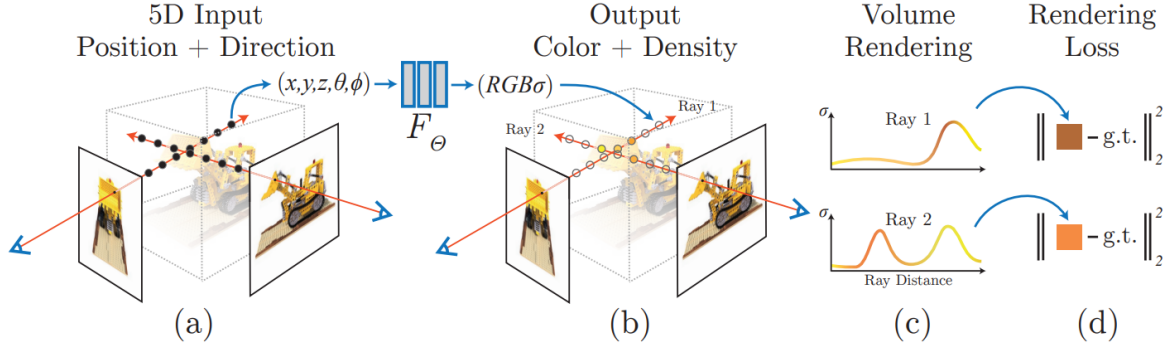


Fig. 3: Overview of NeRF [145] scene representation and differentiable rendering. (a) Images are synthesized by sampling 5D coordinates (location and viewing direction) along camera rays, (b) an MLP produces a color and volume density from those sampled points, and (c) volume rendering allow to reconstruct the final image using those values, all of which is end-to-end differentiable (d).

where $T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)$ and $\delta_i = t_{i+1} - t_i$ is the distance between adjacent samples. This function is trivially differentiable and reduces to traditional alpha compositing with alpha values $\alpha_i = 1 - \exp(-\sigma_i \delta_i)$.

2.1.3 Training for Novel View Synthesis

During training, for each pixel, a square error photometric loss \mathcal{L}_{color} is used to optimize the MLP parameters, as follows:

$$\mathcal{L}_{color} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left\| \hat{C}(r) - C(r) \right\|_2^2 \quad (3)$$

where \mathcal{R} is the set of rays in each batch, and $C(r)$, $\hat{C}(r)$ are the ground truth, volume predicted RGB colors for ray r respectively. The training procedure is typically scene-specific and requires dense images along their 3D poses and intrinsic parameters, and scene bounds which can be estimated using Structure-from-Motion (SfM) end-to-end frameworks, e.g., COLMAP [186], OpenMVG [153], or PixelPerfect [125].

Here is an outline of the novel view synthesis procedure using NeRFs (cf. (a-c) in Figure 3).

- i Send camera rays throughout the image pixels across the scene to produce sampling points.
- ii Use the MLP(s) to compute local color and density data for those sampling points with corresponding viewing direction.

- iii Compute the volume rendering to reconstruct the output image by integrating color and density information, throughout.

2.2 Positional Encoding

By processing the scene with the standard method we have described so far, experiments show that small displacements in input spatial coordinates to the MLPs F_Θ may result in sometimes severe outcomes in synthesized images, in particular in high-frequency textured areas. To mitigate this problem, Mildenhall et al. [145] considered positional encoding, which is the mapping of the coordinate inputs to a higher dimensional space using non-linearities prior to passing them to the neural network. This enables better fitting of data that contain high-frequency variations. The encoding function takes the following form:

$$\gamma(\mathbf{x}) = [\sin(\mathbf{x}), \cos(\mathbf{x}), \dots, \sin(2^{L-1}\mathbf{x}), \cos(2^{L-1}\mathbf{x})] \quad (4)$$

where $\gamma(\cdot)$ is separately applied to each normalized coordinate value in \mathbf{x} and to the three components of viewing direction unit vector \mathbf{d} where L is the encoding dimensionality parameter (typically $L = 10$ for \mathbf{x} and $L = 4$ for \mathbf{d} [145]).

2.3 Depth Rendering

Depth is a valuable source of data for view synthesis and 3D representations. Depth values from a particular pose are calculated in a similar fashion

to rendering RGB pixels:

$$\hat{D}(\mathbf{r}) = \sum_{i=1}^N T_i(1 - \exp(-\sigma_i \delta_i)) t_i \quad (5)$$

where $\hat{D}(\mathbf{r})$ is the expected depth along the camera axis of ray \mathbf{r} and the weighting term is the ray-termination probability of sample i along the ray, defined earlier when rendering color. Following this idea, alternative sources of depth supervision are often utilized to enhance the training and enforce consistency between photometric and geometric constraints. These sources include LiDAR depth [176]; depth cameras; projected point clouds from Structure from Motion packages [38] or pre-trained depth estimation/-completion models [179].

Depth loss \mathcal{L}_{depth} is defined in various formulations with the most widely used approach being the MSE between the predicted depth values $\hat{D}(\mathbf{r})$ and the ground truth depth values $D(\mathbf{r})$.

$$\mathcal{L}_{depth} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \left\| \hat{D}(\mathbf{r}) - D(\mathbf{r}) \right\|_2^2 \quad (6)$$

Furthermore, some methods [159, 251] incorporate depth smoothness constraints in addition to the depth loss. This is based on the observation that real-world geometry often exhibits piece-wise smooth characteristics, where flat surfaces are more common than high-frequency structures. To enforce depth smoothness, these methods typically introduce penalties that encourage neighboring pixels of a rendered patch to have similar depth values.

$$\mathcal{L}_{smooth}(d_i) = e^{-\nabla^2 \mathcal{I}(x_i)} (|\partial_{xx} d_i| + |\partial_{xy} d_i| + |\partial_{yy} d_i|) \quad (7)$$

where d_i is the depth map, $-\nabla^2 \mathcal{I}(x_i)$ refers to the Laplacian of pixel value at location x_i .

2.4 Empowering NeRFs with Semantic Reasoning

Research on neural field representations has shown that MLP networks can be trained from scratch for complex scenes by predicting their volumetric density and color supervised solely by a set of RGB images. However, radiance fields only provide low-level representations of geometry and

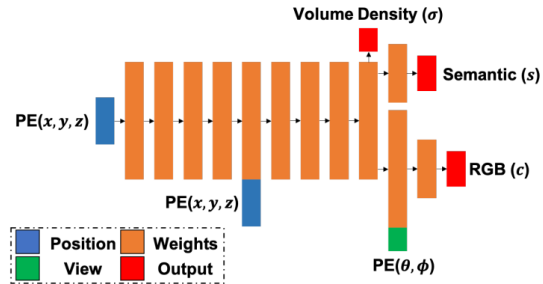


Fig. 4: Semantic NeRFs [285]. 3D positions (x, y, z) and viewing directions (θ, ϕ) are fed into the network after positional encoding (PE). Volume densities σ and semantic logits \mathbf{s} are functions of (x, y, z) while \mathbf{c} additionally depend on (θ, ϕ) .

radiance and lack a higher-level (e.g., semantic or object-centric) understanding of the scene. The standard NeRF approach typically suffers from slow training and fails to recover reliable geometry in some cases when the number of input views is sparse and the depth range is infinite [83, 45]. They are also restricted to learning efficient representations of static scenes that encode all objects of the scene and lack the ability to represent complex scenes and the decomposition into individual objects [160] that populate the scene.

2.4.1 Semantic Radiance Fields

Semantic labels can also be formalized as an inherently view-invariant function that maps only a world coordinate \mathbf{x} to a distribution over semantic labels via pre-softmax semantic logits $\mathbf{s}(\mathbf{x})$, while for instance, it is a one-hot encoding of the object instance identifier. This was done by appending additional branches before injecting the viewing direction \mathbf{d} into the rendering function:

$$\hat{S}(\mathbf{r}) = \sum_{i=1}^N T_i(1 - \exp(-\sigma_i \delta_i)) \mathbf{s}_i \quad (8)$$

Semantic logits can then be transformed into multi-class probabilities through a softmax normalization layer. During inference, the semantic label is determined as the class of the maximum probability in $\hat{S}(\mathbf{r})$.

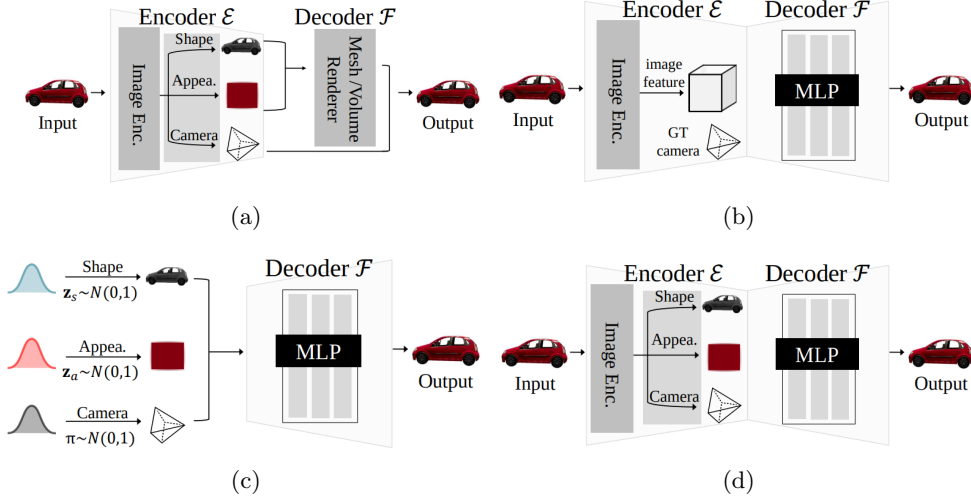


Fig. 5: Different approaches for conditional 3D representation, which can be effectively used for 3D-aware object manipulation: (a) conditional surface or volumetric representation [93, 7], (b) image-conditional NeRFs [268, 32, 217, 191] that train the feature encoder and NeRF as decoder (c) generative NeRFs [187, 18, 255] that render images from randomly sampled disentangled 3D attributes, and (d) auto-encoding NeRFs [85, 133, 100] that extract the disentangled 3D latent codes from input and renders images from these attributes.

The semantic loss \mathcal{L}_{sem} is usually chosen as a multi-class cross-entropy loss to encourage the rendered semantic/instance labels to be consistent with the provided labels, whether these are ground-truth, noisy, or partial observations:

$$\mathcal{L}_{sem}^{2D} = -\frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left[\sum_{l=1}^L p^l(\mathbf{r}) \log \hat{p}^l(\mathbf{r}) \right] \quad (9)$$

where $1 \leq l \leq L$ denotes the class index, $p^l(\mathbf{r})$, $\hat{p}^l(\mathbf{r})$ are the multi-class semantic probability (by forwarding the logits into a softmax normalization layer) of the camera ray \mathbf{r} at class l of the ground truth and predicted map.

2.4.2 Prior Learning and Conditional NeRFs

A conditional neural field introduces the ability to alter the characteristics of a radiance field through the manipulation of latent variables \mathbf{z} . These latent variables can encompass diverse aspects, ranging from random samples drawn from any distribution to geometric/semantic attributes such as shape, type, size, color, and more. Alternatively, they could be derived from the encoding

of other data types, including embedded text or audio data. Instance-specific details can be encoded within the conditioning latent variable \mathbf{z} , whereas information shared across instances is stored in the parameters of the neural field. When these latent variables are mapped to a semantic or smoothly varying space, it allows for their interpolation or editing.

The conditioning latent code $\mathbf{z} = \mathcal{E}(I)$ is generated by an encoder or embedding mechanism \mathcal{E} usually implemented as a neural network (as shown in Figure 5). The parameters within \mathcal{E} are capable of encoding prior knowledge, which can be learned from pre-training on datasets or through auxiliary tasks. The decoder is the neural field that is conditioned by the latent code:

$$(\mathbf{c}, \sigma) = F_{\Theta}(\mathbf{x}, \mathbf{d}, \mathbf{z}) \quad (10)$$

This adaptation can be achieved by conditioning the field on latent variables \mathbf{z} that encapsulate specific higher-level, semantic characteristics from the scene. When these latent variables are edited, the corresponding neural field can be modified accordingly.

2.4.3 High-level Feature Consistency

While the field of 3D scenes has unique challenges and complexities, the image domain stands out with its abundance of extensive, high-quality datasets and a wealth of established techniques for effective feature extraction. The semantic richness captured in image feature spaces can be harnessed to establish correspondences and enhance understanding through text, image queries, or clustering. Although there exists pixel-wise misalignment between the views, it is observed that the extracted representation of pre-trained deep neural networks as feature extractors is robust to this misalignment and provides supervision at the semantic level [83, 251, 150]. Intuitively, this occurs naturally because the content and style of the two views are alike, allowing a deep network to learn a representation that remains consistent across them. Perceptual Loss \mathcal{L}_{feat} , also known as feature loss or content loss, is a measure of the discrepancy between the high-level features of the predicted image and the ground truth image, both extracted from the pre-trained network:

$$\mathcal{L}_{feat} = \frac{1}{N} \sum_{i=1}^N \left\| \Phi(\hat{I}_i) - \Phi(I_i) \right\|_2^2 \quad (11)$$

where $\Phi(\cdot)$ refers to the extracted features/tokens. I and \hat{I} are patches or images from the reference view and rendered view, respectively. This commonly appears in tasks where preserving high-level global features is important.

It is shown [104, 221] that the knowledge distilled from the teacher model aligns with the scene’s geometry, thereby enhancing feature quality across viewpoints and occlusion awareness; and the infusion of features pre-trained on diverse external datasets, bringing a broader open-world perspective to the 3D representation without collecting annotations for them.

3 Semantically-aware NeRFs for Visual Scene Understanding

In this section, we review the most prominent NeRF-based approaches and strategies that either use semantic-level reasoning as leverage to enhance 3D geometry or that aim to achieve a

higher level of scene understanding through either of the tasks and applications considered in our considered taxonomy (Fig. 2 and Fig. 6).

3.1 3D Geometry Enhancement

In this category, several notable approaches incorporate semantic reasoning to improve performance in Novel View Synthesis (NVS), to make up for limited amounts of input views, to generalize to unseen environments, or to address 3D surface reconstruction.

3.1.1 One-shot/Few-shot NeRFs

PixelNeRF [268] (Fig. 7) and **S-RF** [32] use image-level CNN features, whereas **MVS-NeRF** [22] builds a 3D cost volume via image warping which is then processed by a 3D CNN. This fully convolutional strategy allows the network to be trained across multiple scenes to learn scene-level priors and, thus, generalize to unseen environments and object categories. Building on this concept, **MINE** [113], **Behind the Scenes** [238], and **SceneRF** [14] reduce the scene representation complexity by leveraging monocular depth estimation and redefine feature extraction and rays and color sampling accordingly to account for the self-supervised depth network. **DietNeRF** [83] and **SinNeRF** [251] match high-level and global semantic attributes to semantic pseudo-labels with texture guidance across different views, allowing us to supervise the training process from random poses. This improves the perceptual quality of NVS in the few-shot setting in particular.

Single-view (i.e., one-shot) reconstruction can also be formulated as a conditioned 3D generation problem for a single-image NVS task without explicit 3D supervision. **RealFusion** [142] and **Zero-1-to-3** [132] extract a neural field from the original image input and a internet-level pre-trained diffusion models, thus achieving a comprehensive reconstruction of the object from unseen viewpoints, or in a prompt-constrained zero-shot setting. This process captures both appearance and geometry. Additionally, image-level text embeddings can be extracted through textual inversion, which captures additional high-level visual cues. However, such a strategy yields

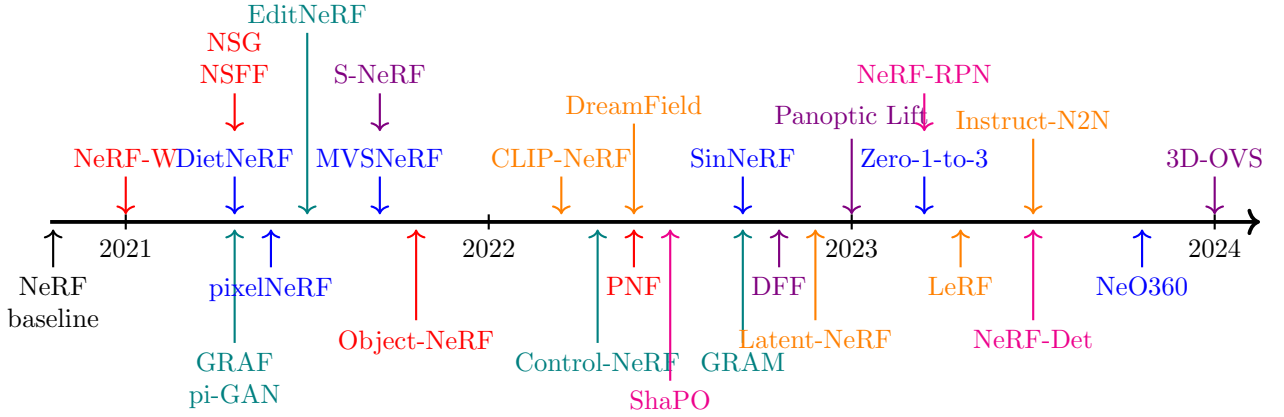


Fig. 6: Chronological overview of the most relevant semantically-aware NeRFs spanning all 6 categories covered by our study: 3D geometry enhancement, segmentation, editable NeRFs, object detection and 6D pose, holistic decomposition, NeRFs and language.

substantially ambiguous representations in unobserved areas, and they are mostly object-centric assuming a plain background.

NeRDi [37] also uses diffusion priors trained on large image datasets. It utilizes a two-section semantic guidance to refine the general prior knowledge conditioned on the input image. This ensures that synthesized novel views are both semantically and visually consistent. Despite the training of the model being carried out on a synthetic dataset, it shows robust zero-shot generalization capabilities. It effectively extends to both out-of-distribution datasets and real-world, in-the-wild images. **SegNeRF** [273] and **S4C** [66] address generalization and learn a semantic field, performing reconstruction and segmentation in a self-supervised fashion from a single view, while also allowing for semantic object/scene completion. **Neural groundplans** [191] conditions a self-supervised NeRF on ground-aligned 2D feature grids trained from multi-view videos. **NeO 360** [81] leverages a hybrid conditional triplanar representation which combines the strengths of voxel and bird’s eye view (BEV) representation. These hybrid discrete-continuous representations allow to learn from a large collection of 360 unbounded scenes while addressing different downstream tasks including NVS, object localization, and scene editing from as few as a single image during inference.

3.1.2 Surface Reconstruction

The piecewise planarity assumption, i.e., which assumes that a given scene can be mostly explained by piecewise planar surfaces has been a stable prior in the traditional 3D reconstruction literature and has also proven effective in the context of implicit neural representations. Guo et al. [60] formulate geometric constraints of floors and walls within the normal loss function adhering to the Manhattan World Assumption [34], assuming three mutually orthogonal surface orientations. These regions were obtained by a 2D semantic segmentation networks. To address inaccurate segmentations, they encode the semantics of 3D points with another MLP that jointly optimizes the scene geometry and semantics. **PlaNeRF** [230] also performs a planar regularization based on Singular Value Decomposition (SVD). This improves the underlying geometry in that correspond to image regions with low texture, without any additional geometric prior. **S3PRecon** [262] introduces an iterative training scheme for grouping pixels and optimizing the reconstruction network via a superplane constraint. This in particular yields better performance than using explicit 3D plane supervision which is costly to obtain.

SS-NeRF [279] and **MuvieNeRF** [283] are versatile multi-task frameworks. They can render images from novel viewpoints and manage various scene properties such as appearance, geometry, and semantic segmentation. Both utilize a shared

scene encoding network that allows for cross-view and cross-task attention modules to ensure view consistency. They also examine the relationships among different scene properties to enhance performance. This approach highlights the potential of multi-task learning and knowledge transfer within a synthesis paradigm in benefiting from the mutually-informative relationships between different tasks and properties, e.g., semantic labels, surface normal, shading, keypoints, and edges.

3.2 Segmentation

The most common approach for scene understanding typically focuses on 2D reasoning in image space, using classic image-to-image architectures that are trained on extensive sets of semantically annotated images. Although these techniques are easy to implement, they generate only pixel-by-pixel annotations and mostly overlook the underlying 3D structure of the scene. In contrast, our objective is to use a set of RGB images with established poses to produce a 3D semantic/instance field. This involves devising a function that assigns probability distributions over semantic and/or instance-level categories to specific 3D positions and viewpoints.

3.2.1 Semantic, Instance, and Panoptic Segmentation

NeSF [226] uses a pre-trained NeRF to generate a volumetric density grid. Following this, a 3D UNet is used to produce a feature grid that maintains the same spatial resolution. This process enables high-level reasoning within the 3D space. Semantic maps are generated through the application of the volumetric rendering equation, using camera poses on the semantic field. Consequently, NeSF is trained comprehensively on various scenes, eliminating the need for segmentation input when making inferences about new scenes.

Semantic-NeRF [285] is a groundbreaking work that extends NeRF to include both semantics along with appearance and geometry. By adding semantic class predictions to radiance and density within a scene-specific implicit MLP model, it can ensure multi-view consistency between semantic labels. Consequently, the experiments demonstrate its ability to perform multi-view semantic label fusion in various scenarios: pixel-wise label noise, region-wise label noise,

low-resolution dense or sparse labeling, partial labeling, and using the output from an imperfect segmentation model. In this respect, several studies leverage 3D geometry together with semantic predictions to resolve label uncertainties. For example, **Panoptic NeRF** [45, 46] introduces an optimization process guided by semantics to enhance the underlying geometry. This technique uses a dual of semantic fields: a fixed semantic field which focuses on guiding the underlying density, defined by 3D bounding primitives, and a learned semantic field designed to capture the semantic distribution.

Another work, **Semantic Ray** [126], fully exploits semantic information along the ray direction from its multi-view re-projections. The authors tackle the limitations of prior methods that depend on positional encoding and scene-specific models for semantic learning. Unlike these approaches, they harness insights from multiple views using a new module called Cross-Reprojection Attention. This module efficiently captures contextual information along the reprojected ray paths, enriching the understanding from various views.

JacobiNeRF [253] introduces a regularization of learning processes to align the Jacobians of highly correlated entities, effectively maximizing their mutual information amid random perturbations in the scene. This approach of mutual information modeling is key in configuring NeRF to perform sparse label propagation for semantic and instance segmentation. For a given target view of a scene that is unlabeled, one can produce labels by selecting the argmax of the perturbation responses from the source view annotations.

Liu et al. [137] propose the training of a Semantic-NeRF network for each scene by fusing the predictions of a segmentation model and using the view-consistent rendered semantic labels as pseudo-labels for model adaptation. Their method simultaneously trains the frame-level semantic network and the scene-level NeRF, ensuring that the semantic forecasts and NeRF renderings are in alignment. This transfer strategy not only boosts the performance of both models but also reflects a real-world deployment scenario that accounts for the covariate shift across different scenes and the possibility of revisiting previously observed scenes.

Traditional methods depend on accurately labeled ground truth data to train models for

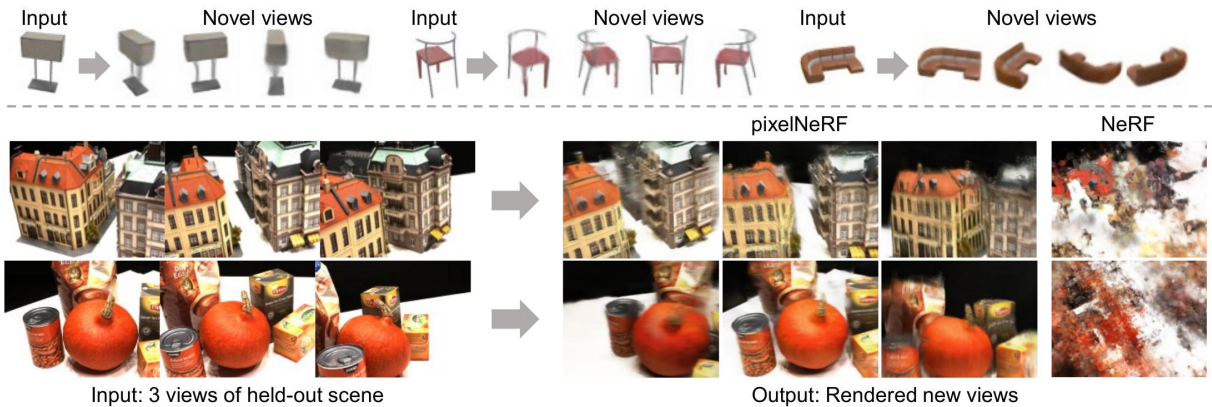


Fig. 7: Generalization capabilities of PixelNeRF [268] – The method can be trained from multiple datasets in order to synthesize plausible novel views from very few input views, without test-time optimization. In contrast, standard NeRF can not generalize to previously unseen environments.

object-compositional scene representations. It’s important to recognize that these manual annotations are designed to be 3D consistent, ensuring that identifiers for specific objects remain constant across different viewpoints. However, a major challenge presents itself when using pseudo-labels generated by off-the-shelf networks. These labels, inferred from individual views, often fail to maintain the 3D alignment of instance indices, leading to inconsistency. Several studies have focused on addressing the discrepancies and maintaining consistency across different viewpoints within the same scene, particularly when employing an off-the-shelf 2D panoptic segmentation network. These efforts strive to preserve the object instance identities from machine-generated panoptic labels within the implicit 3D volumetric representation. For instance, **Panoptic Lifting** [195] assigns 3D surrogate identifiers to machine-generated instances by solving linear assignment problems, using these associations to guide the training of the instance field through an NCE loss.

Contrastive Lift [6] introduces a change to the labeling process by using a low-dimensional Euclidean space, which simplifies the model by reducing the dimensions needed to calculate pairwise distances. This slow-fast clustering objective function is scalable and suitable if there are a large number of objects (up to 500 per scene). On the other hand, **PCFF** [31] proposes an Instance Quadruplet loss which leads to a discriminating feature space for the scene decomposition at instance levels. The model is further

refined with strategies that are added to the architecture, like semantic-appearance hierarchical learning and semantic-guided regional refinement. Finally, **Instance-NeRF** [136] seeks to match 3D object masks projected from a proposal-based NeRF-RCNN with inconsistent segmentation maps in image space, thereby refining the initial instance segmentation results.

3.2.2 Pre-Semantic Segmentation

DFE [104], **N3F** [221] and **FeatureNeRF** [263] adopt a 2D-teacher-3D-student framework. In this setup, pre-trained 2D image feature extractors like LSeg [111], SAM [102], and DINO [16] act as ‘teachers’ that guide the learning process of a NeRF ‘student’ network. The loss function in this context is designed by imposing penalties on the discrepancies between rendered features and the outputs generated by the feature descriptor. These methods pave the way for applications in language-guided editing, 3D spatial rearrangements, and targeted scene removal. **NeRF-SOS** [41] integrates a self-supervised pre-trained framework to generate feature tensors from color patches rendered by the model. This approach then uses these features to create volumes for appearance-segmentation, applying contrastive losses to correlate both appearance-segmentation and geometry-segmentation. During inference, the model perform a clustering process

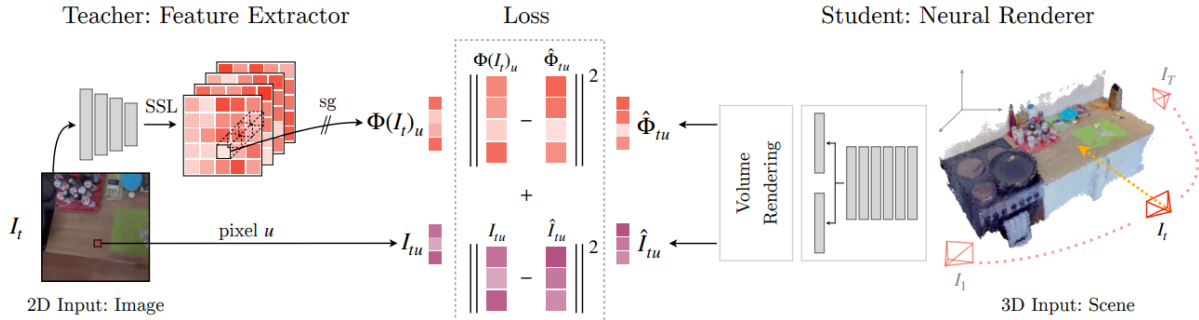


Fig. 8: N3F paradigm [221] - A student-teacher framework distills predicted 2D features from images into a 3D student NeRF-like model. The student optimizes image and feature reconstruction, while the teacher remains untrained. The resulting representation can operate in 2D or 3D contexts.

on the rendered feature field to produce segmentation masks. Similarly, **3D-OVS** [128] demonstrates that aligning the class relevancy distribution with these pre-trained foundation models in a weakly supervised way can achieve precise, annotation-free segmentation, as shown in Figure 9. **Feature-Realistic Fusion** [141] fuses general features learned from EfficientNet into NeRF representation. With a SLAM backend, this system operates incrementally in real-time, effectively managing the exploration of new, unobserved regions of a scene.

RFP [135] introduces an innovative propagation method that uses a bidirectional photometric loss. This approach allows for unsupervised partitioning of a scene into distinct, salient regions that correspond to individual object instances, effectively performing object segmentation within the scene. **IntrinsicNeRF** [264] goes further by producing outputs like reflectance, shading, and residual terms. The model is trained using unsupervised prior and reflectance clustering as constraints in the loss function. These terms are particularly useful for real-time augmented applications such as recoloring, illumination variations, and, importantly, semantic segmentation.

SNeRL [193] integrates NeRF with semantic and distilled feature fields specifically for reinforcement learning applications. It employs a NeRF-based autoencoder, trained to act as a feature extractor, for fine-tuning in multi-view reinforcement learning (RL) tasks. This method has shown to outperform current representation

learning techniques in both model-free and model-based RL algorithms across various 3D environments.

3.2.3 Interactive Segmentation

For a practical scene-annotation tool, simple user annotations like sparse clicks can be extended and propagated to achieve dense and accurate labeling of the scene. This process allows for the creation of complete and accurate 2D semantic labels with minimal in-scene annotations specific to the scene. **iLabel** [284] takes this concept further by integrating semantic label-propagation into an online, interactive 3D scene-capturing system, enabling segmentation of coherent 3D entities with minimal user click annotations. The authors also introduced a novel hierarchical semantic representation using a binary tree, facilitating the prediction of semantics at different levels. **Baking in the Feature** [10] and **ISRF** [56] merge distilled features with a bilateral search in a unified spatial-semantic space for an interactive segment user interface. **NVOS** [177] trains a 3D segmentation network to classify each voxel as foreground or background, using partial user scribbles as supervision. This is followed by applying the learned classifier and further refining the segmentation with a 3D graph-cut, leveraging the 3D distance field of the scribble. Other methods [17, 236] aim to generalize the Segment Anything Model (SAM) [102] for 3D object extraction. These alternate between mask inverse rendering and cross-view self-prompting across different views to iteratively complete the 3D object mask from a single

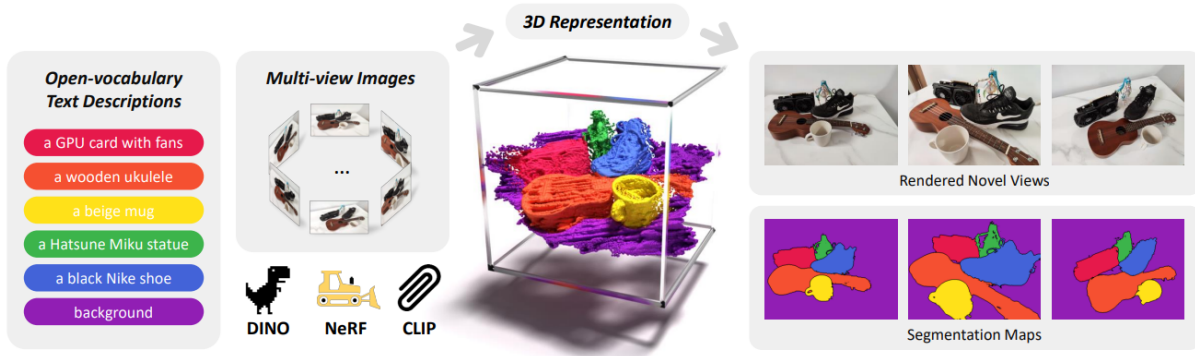


Fig. 9: The method developed by Liu et al. [128] merges multimodal open-vocabulary knowledge from CLIP [170] with the object reasoning capabilities of DINO [16] resulting in accurate delineation of 3D objects without relying on segmentation annotations during training. This showcases the model’s ability to render novel views with corresponding segmentation maps.

view. Users can annotate frames in an RGB video sequence with brush strokes, while the system concurrently fits a model to the scene and annotations. These strategies surpass the labeling accuracy of conventional pre-trained semantic segmentation methods. **SGISRF** [215] takes this even further, requiring fewer user interactions for interactive segmentation by using Cross-Dimension Guidance Propagation and Concealment-Revealed Learning schemes. Another key interactive fea-



Fig. 10: iLabel [284] demonstrates the capability to create high-quality segmentations of various entities within a scene using only a minimal number of user-provided clicks.

ture in 3D scene manipulation is the removal of undesired objects in a way that the resulting area blends seamlessly and logically with its surroundings, a process commonly known as image inpainting. This technique starts with a pre-trained NeRF model and its associated image dataset. In the first stage, known as mask generation, an initial mask is created from a single-view annotation using one-shot segmentation methods like Mask

R-CNN [67], SAM [102], or GLIP [114, 229]. Following this, **NeRF-In** [127] uses an inpainting network [15, 211] to produce a guiding image and a depth image, based on the user-selected area to be removed. This process updates the NeRF model by optimizing both the color-guiding and depth-guiding losses. While NeRF-In doesn’t fully resolve the 3D inconsistencies in the output of the inpainters, and only minimizes the number of views used, there are proposals to overcome blurring and ensure consistency between views. These include using a relaxation approach based on perceptual loss [150], applying bilateral solvers, and incorporating estimated depth to introduce view-dependent effects in the inpainted regions [149]. Another technique involves selectively excluding views using an uncertainty mechanism and a pixel-wise loss [235]. Weder et al. [235]’s method updates the set of images used for optimization iteratively, based on confidence scores, to maintain consistency during the inpainting process. This enables the generation of realistic novel views of the scene without the removed objects.

3.3 Editable NeRFs

3.3.1 Conditional NeRFs

CodeNeRF [85] implements the learning of separate embeddings, whereas **EditNeRF** [133] incorporates a shared shape branch within the

conditional radiance field, aiming to better reconstruct shape instances. Both approaches encourage the network to develop a common representation across different object instances, resulting in enhanced shape editing and consistency. **ShaRF** [175] employs a shape network that maps the shape latent code into a 3D shape in the form of a voxel grid. The NeRF network then conditions on two additional factors: the occupancy value estimated from the voxel grid and the appearance latent code that dictates the object’s appearance. **AE-NeRF** [100] introduces two specific losses - global-local attribute consistency loss and swapped-attribute classification loss - to enhance disentanglement capabilities. Furthermore, this conditional model benefits from a GAN-based stage-wise training approach, significantly elevating its performance. **AutoRF** [155] and **Car-NeRF** [134] develop an object-level radiance field specifically for cars, effectively disentangling shape and appearance within their image encoders. For each car instance, they use a panoptic segmentation mask and a 3D bounding box, which describe the pose and size of the object. These models transform each ray from the camera space into the Normalized Object Coordinate Space (NOCS), creating an object-centric ray that allows for the generation of high-quality car images from any single-view input. Thanks to their ability to perform shape and color edits within the network layers, these models facilitate a hybrid network update strategy. This approach enables the formulation of optimization problems for color and shape editing that meet specific user requirements while maintaining the integrity of the original object’s structure. Such features are key in preserving the overall visual coherence of the edited object and reduce the number of images needed during testing.

3.3.2 Generative NeRFs

Recent advancements in NeRF-based generative models, including VAEs, GANs, and diffusion models, have significantly progressed in creating 3D-aware generators. These models have the capability to disentangle the underlying 3D aspects of the objects they represent, enabling precise manipulation of camera poses while still producing high-fidelity object renderings. Additionally, these models are designed to generate view-consistent

and varied images that accurately reflect specified conditions. The versatility of these models is further enhanced by their ability to incorporate a range of user-defined conditions, such as text and images, into their generation process. **GRAF** [187] and **pi-GAN** [18] introduce a generative model that employs implicit radiance fields for the synthesis of novel scenes. These models are trained on unposed images and focus on simple objects. Building on GRAF, **GIRAFFE** [157] enhances this approach by representing scenes through compositional generative neural feature fields. This advancement allows for the disentanglement of individual objects’ shapes and appearances from their backgrounds without the need for explicit supervision. As a result, users are provided with greater flexibility to compose more complex scenes. While this method is less demanding on memory when scaling up to higher resolutions compared to voxel-based techniques, it still requires considerable computational power to train and render images at high resolutions. In response, **StyleNeRF** [59], **GIRAFFE HD** [255], and work by Chan et al. [19] all aim to retain the 3D controllability characteristic of GIRAFFE while generating images of much higher quality and resolution (exceeding 512×512), using the architecture of StyleGAN2 [95]. Following in these footsteps, **UrbanGIRAFFE** [261] extends this concept by using coarse panoptic priors in the form of semantic voxel grids and object layouts. This approach enhances controllability even more, particularly for substantial changes in camera viewpoint and semantic layouts.

Reconstructing a dynamic human face poses unique challenges due to the complexity of facial geometry and the varying appearances caused by diverse expressions. Facial expressions, involving a mix of local deformations, can be represented through controllable attributes defined as latent variables. These attributes can be flexibly applied to different types of conditions such as landmarks [171], sketches, low-resolution images, and text [90] as input conditions. Methods like **CoNeRF** [94] and **FaceCLIPNeRF** [79], which build upon HyperNeRF [163], are capable of being trained on dynamic scenes to control facial deformations using only sparse input views. Users can manipulate facial attributes effectively by providing simple expression codes [47, 277], mask



Fig. 11: Illustration of view-controllable images generated by GAN-based NeRFs, showcasing the high quality and 3D consistency of the output. The left set (a) features images from the GRAM [40] study, while the right set (b) includes images from the work of Chan et al. [19].

annotations of facial regions [208, 94] (such as eyes being open/closed or mouths smiling/frowning), or textual descriptions [79] (like “happy”, “surprised”, “fearful”, “angry”, and “sad”). These methods allow for precise control over facial expressions and attributes.

Recent progress in this field has led to more fine-grained applications, particularly in avatar generation [39, 40, 19, 289, 267, 27, 286, 245, 49] and human pose generation [130, 206, 276, 282, 237, 86, 28, 73, 154, 25]. A significant achievement of these technologies is their ability to produce high-fidelity animations of real subjects using only a limited number of input images. This breakthrough not only conserves resources but also opens up exciting research prospects, especially in fields like video gaming, augmented and virtual reality (AR/VR), and human-computer interaction.

3.3.3 Spatial Transformation Editing

ST-NeRF [275] presents a layered representation approach for each dynamic entity within scenes, where every entity is represented as a separate continuous function that spans both space and time. The MLP network of the model is composed of two key modules: a space-time deform module and a neural radiance module. In this setup, the frame number is directly encoded in the model. This approach to disentanglement of space and time facilitates various spatial editing techniques, such as affine transformation, insertion, and removal, along with temporal editing

capabilities such as re-timing, as demonstrated in Figure 12.

Approaches like **AutoRF** [155], **Neural Scene Graph** [160], **PNF** [107] or **Dis-CoScene** [254] build a full 3D radiance field for each object contained within a bounding box. By treating objects’ radiance fields as independent entities, we can render scenes more efficiently by focusing only on the relevant points where the rays intersect with these bounding boxes (ray-box intersection). This allows for image editing through the manipulation of bounding boxes, enabling the repositioning (rotation and translation) of objects in a scene without altering their visual appearance. For operations like removal or replication, users can adjust the scene’s layout by deleting or cloning bounding boxes. In scenarios without bounding boxes [195], object removal is accomplished by reducing the density of points associated with the target instance to zero. Meanwhile, duplicating the weights of an object’s MLPs or its latent codes can result in the cloning of that instance in the scene.

Control-NeRF [109] learns volumetric representations for multiple scenes by using a single shared rendering model. During testing, because the feature volumes are separate from the rendering model, the authors can perform spatial adjustments to these volumes or combine them. This process allows for the editing of the scene content without altering the fixed parameters of the rendering network.



Fig. 12: Illustration of ST-NeRF’s [275] capabilities to perform complex editing tasks in dynamic scenes. Here, we see the application of spatial affine transformations, temporal retiming, and transparency adjustments to selected objects within a scene. The transformations are applied across different timestamps and 3D bounding boxes around the target objects.

3.4 Object Detection and 6D Pose Estimation

The task of 3D object detection is essential for a variety of applications, as it provides a detailed understanding of objects’ sizes and positions in three dimensions. This task is more complex than 2D object detection due to the challenges in obtaining precise 3D data and the additional degrees of freedom (DoF). Methods based on point cloud representations depend heavily on accurate data from specialized sensors. Therefore, innovative techniques are necessary to leverage the capabilities of NeRFs while addressing the complexities of accurate 3D object detection from 2D images.

3.4.1 3D Object Detection

NeRF-RPN [72] is designed to identify all bounding boxes in a scene. The process begins by sampling a grid of points, from which RGB and density values are extracted using a pre-trained NeRF model. These volumetric features are then processed through a 3D Feature Pyramid Network (FPN) [123] backbone, yielding deep, multi-scale 3D features. These features are inputted into a 3D Region Proposal Network (RPN) head, generating region proposals. A key innovation in

NeRF-RPN is its use of a novel voxel representation, integrating multi-scale 3D neural volumetric features. This allows for the direct regression of 3D bounding boxes within NeRF without needing to render from any viewpoint. In contrast, **NeRF-Det** [250], a joint NeRF-and-Det method, links the NeRF branch with the detection branch using a shared geometry-based MLP. This setup enables the detection branch to use the gradient flow from NeRF in estimating the opacity fields. Consequently, it effectively masks out free space and reduces ambiguity in the feature volume, offering improvements over the NeRF-to-Det approach.

MonoNeRD [252] approaches the concept of monocular 3D detection with NeRFs by considering intermediate frustum representations as SDF-based (Signed Distance Function-based) NeRFs. These are then optimized using volume rendering techniques. The process involves grid sampling on these frustum features to construct regular 3D voxel features along with corresponding densities. These voxel features are subsequently inputted into detection modules. This methodology establishes a new standard in monocular 3D detection using NeRFs.

On another front, techniques like **Neural groundplans** [191] and **SUDS** [222] use feature field clustering to derive object-centric 3D representations in an unsupervised manner. These

methods start with a dynamic field and apply traditional connected-component labeling in the feature space, considering the cumulative density values. This process aids in identifying individual objects. The smallest box enclosing each connected component is then computed, resulting in a 3D bounding box for every detected object.

3.4.2 6D Pose Estimation

ShAPO [80] extracts comprehensive 3D details of multiple objects from a single RGB-D observation. This includes the objects’ shape, 6D pose, scale, and appearance. The technique uses an octree-based differentiable optimization, drawing on pose, texture, and masks derived from an FPN [123] backbone. **NCF** [75] is a method that estimates the 6D pose of a rigid object using a single RGB image. It maps from the camera space to the object model space. NCF predicts the corresponding 3D point in the model space and its signed distance. This facilitates the creation of 3D-3D correspondences, crucial for determining the object’s pose. **NeurOCS** [147] focuses on predicting the object mask and NOCS (Normalized Object Coordinate Space) map, which are then used in PnP (Perspective-n-Point) algorithms to estimate object pose. Additionally, a separate detector is applied to the NOCS and predicted depth data, aiding in precise 3D object localization.

3.5 Holistic Decomposition

3.5.1 Objects vs Background

NeRF-W [140] incorporates per-frame embeddings and a transient branch to model non-photometric consistent effects in unconstrained photo collections. Although it wasn’t specifically designed to distinctly separate objects from their surroundings, it offers an innovative method for foreground element capture within various environments.

Subsequent research, including works [204, 248, 257, 8], has led to the division of NeRF into a dual-pathway architecture. This structure comprises a scene (background) branch for encoding scene geometry and appearance, and an object (foreground) branch for individual object encoding. These models learn to encode multiple objects simultaneously by assigning activation codes to

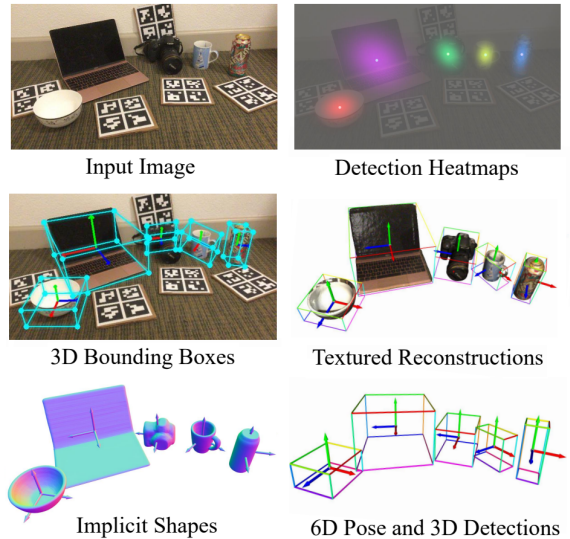


Fig. 13: ShAPO’s [80] multifaceted 3D object detection and pose estimation capabilities illustrated through an input image’s transformation into detection heatmaps, 3D bounding boxes, textured reconstructions, and implicit shape representations.

training rays for each object, eliminating the need for separate training per object. For view generation with object manipulation, they render the transformed objects using the conditioned object branch and the background from the scene branch together. An added feature includes an object manipulator for precise radiance and object field editing, taking into account challenges such as object collisions and occlusions. Meanwhile, works like **uORF**’s [270, 197], aim to deduce latent object-centric representations into distinct slots through an attention mechanism, facilitating unsupervised segmentation.

In their panoramic room capture study, Yang et al. [258] initially predict object metadata and infer object-to-object and object-to-room relationships, leveraging object-level predictions and geometric cues. They also incorporate pre-convolved HDR maps and surface normals into a global optimization, enabling the synthesis of novel lighting conditions and scene touring. Zhu et al. [287] employ one MLP to accurately model a scene with occlusions and another MLP for the background. They train the background MLP and remove occlusions from aggregated information

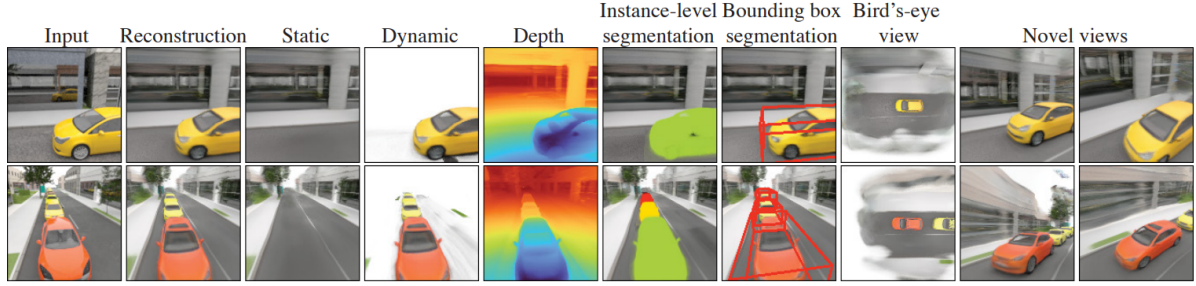


Fig. 14: Semantically-trained Neural Groundplans [191] for Urban Scene Analysis – The model discerns static/dynamic entities and object instances, while enhancing novel view synthesis. Key features include synthesis of novel viewpoints, filling in missing information, accurate 3D bounding box predictions, 3D scene editing capabilities, and extraction of object-centric 3D representations.

from scene MLP to determine if the output of the background NeRF matches the observed color, learning a mask from the ray’s weights. This approach includes a depth constraint for probing occluded areas, by comparing the depth of occlusion and background, based on the assumption that occlusions are foreground objects at a closer distance.

vMAP [105] designs a vectorized object-level mapping, where each object is detected through its 3D point cloud and instance segmentation map, which is then represented by a separate MLP. The 3D bounds are continually updated, via data association across frames, leading to improvements in object-level reconstruction quality and runtime efficiency compared to traditional SLAM systems.

Zhang et al. [281] focus on representing scenes using small local radiance fields, termed “nerflets”. Each nerflet covers a specific scene portion, determined by its influence function. These nerflets can collectively represent complex object instances, providing a more efficient and compact representation for outdoor environments that can be rendered, decomposed, and edited.

AssetField [247] presents a natural visualization of scenes in Bird’s Eye View (BEV) using informative ground feature planes aligned with the physical ground. This approach extracts and categorizes neural representations of scene objects, enabling users to manipulate and compose assets directly on the ground feature plane using feature patches from multiple scenes.

Haughton et al. [65] and Chen et al. [26] demonstrate how robots can identify objects and

build composed 3D representations through physical interactions like pushing, grasping, or poking. The coherence of their model allows for the efficient propagation of measured physical properties (e.g., poses, rigidity, material) throughout the scene. Their experiments highlight the potential for automated sorting and grasping tasks.

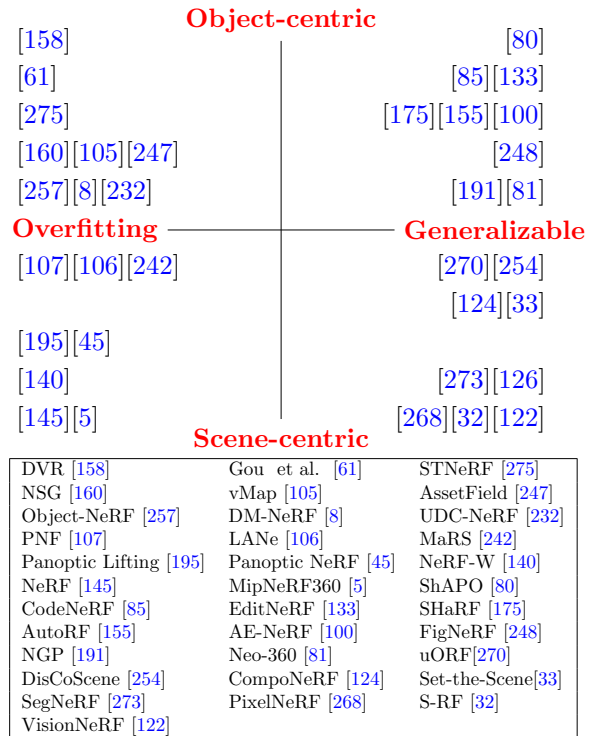


Fig. 15: Categorization of methods that jointly infer appearance and shapes.

3.5.2 Static vs Dynamic Objects

Perceiving and representing dynamic environments is essential for autonomous agents to understand and interact with their surroundings. The key challenge involves disentangling camera and object motion while simultaneously reconstructing the dynamic scene. Such representations permit the synthesis of novel views within dynamic settings or the distinction between moving and stationary elements, offering flexibility in perspective and timing, i.e., in a free-view and time-varying manner. Incorporating a temporal element into an MLP could serve as a viable approach. This would involve encoding the time variable t , either by mapping it to a higher-dimensional space using frequency encoding or a 4D-hash grid, in a similar manner to the spatial coordinates \mathbf{x} and \mathbf{d} , or through learnable, time-dependent latent codes as suggested in several studies [116, 164].

For scenes that are predominantly static, optimizing a single model could lead to blurry outputs and inconsistencies. Solutions like **Dyn-NeRF** [50] and **STaR** [271] have been developed to segregate moving objects from the static background. They use two separate branches for static landscapes and dynamic objects: a static branch containing non-moving topography consistent across videos and a dynamic branch that handles dynamic objects. The training of these branches is often directed using pre-existing semantic and motion segmentation methods, creating masks that exclude “dynamic” pixels from the static training process. This approach ensures the background is reconstructed accurately without conflicting the losses, avoiding errors caused by moving objects. Additionally, temporal variations can be accounted for in a self-supervised manner through regularization [271, 220, 191, 240], which enables the dynamic field to learn as necessary. **D²NeRF** [240], an extension of **HyperNeRF** [163] to dynamic scenes, can handle complex scenes involving multiple non-rigid and topologically varying objects. This method is able to decouple dynamically moving shadows with a separated field that reduces the static radiance output as well. The features encoded from both branches can be regularized during training and can be interpolated using MLPs or 4D hash grids in both short- and long-term space-time ranges [164]. This technique not only delivers

high-quality, smooth rendering performance but also enhances the efficiency and stability of the training process.

On the contrary, the dynamic model proposed by Li in **NSFF** [117] goes a step further by directly predicting forward/backward scene flows along with disocclusion weights from a multilayer perceptron (MLP). These disocclusion weights act as an unsupervised confidence, determining the locations and intensity at which to apply the temporal photoconsistency loss. The model uses a pre-trained 2D optical flow model to supervise the predicted 3D flows, which are also refined using a cycle consistency term for regularization. Building on this work, **SAFF** [119] enhances the model by also generating semantic and saliency features, which are instrumental in refining the segmentation of static and dynamic elements within the scene. In a similar vein, **Factored-NeRF** [239] leverages annotations from keyframes, propagating them to adjacent frames to deduce scene flows, map object trajectories, and determine rigidities. Through comprehensive end-to-end optimization, this model gains the ability to modify object placements, trajectories, and even adapt to non-rigid movements. By computing static and dynamic fields independently, these methods facilitate the separate rendering of stationary and moving parts within a scene. In addition, a significant body of research is dedicated to testing these approaches in dynamic environments, notably those with complex movements, such as vehicles and pedestrians in urban settings. **Neural Scene Graphs** [160] introduces a learned scene graph representation that encodes the transformations and radiance of objects. This method uses tracking data and video frames to learn distinct representations for each object within the scene graph, thereby streamlining the process of synthesizing and decomposing views across various object arrangements and dynamic conditions. This progress enables not only the realistic rendering of new scenes and objects, but also the potential for 3D object detection through the technique of inverse rendering.

PNF [107], **LANe** [106], and **MARS** [242] also break down scenes into distinct objects and backgrounds, employing panoptic segmentation and bounding boxes for each object. Each object is represented by an oriented 3D bounding box and is characterized by a dedicated MLP that computes density and radiance from inputs like position,

direction, and pose. These MLPs are tailored to individual instances and refined through a meta-learning initialization process [107]. **LANe** [106] trains on a single scene under varying lighting conditions, learning to adapt by creating a light field and using a corresponding shader to modulate the appearances of objects for coherent integration into scenes with different lighting. **SUDS** [222] and the subsequent **EmerNeRF** [260] handle scalability by employing a multi-resolution hash table for scene partitioning, enabling dynamic management of vast numbers of objects over extensive areas (hundreds of kilometers), using implicit scene flows and DINO [16] features for enhancement. **Neural groundplans** [191] process their ground-aligned 2D feature grids through a 2D CNN, effectively disentangling the representation into two distinct groundplans for static and dynamic features, thus achieving a clear disentanglement.

3.6 NeRFs and Language

3.6.1 Text-driven 3D Generation and Editing

Text-guided image generation has seen tremendous success in recent years, primarily due to the breathtaking progress in language image and diffusion models. These have also inspired major breakthroughs in text-guided shape generation. This progress has influenced research, linking NeRFs with textual input descriptions.

CLIP-NeRF [227] extends the work on conditional Neural Radiance Fields (NeRF) by promoting similarity between the CLIP [170] embeddings of scenes, facilitating user-friendly manipulation of NeRF through short text prompts or example images. This approach disentangles the latent representation, allowing for separate control over the shape and appearance of objects. Consequently, it enables the creation of code mappers that modify latent codes based on user-specified edits via text prompts or images, demonstrating improved editing capabilities and narrowing the gap between textual and visual editing cues. **DreamField** [84] employs a CLIP model pre-trained on large datasets of captioned images from the web. It guides the generation process so that the rendered images achieve high scores with a target caption according to the CLIP model, even

without access to 3D shape or multi-view data. This method facilitates the zero-shot generation of diverse 3D objects from captions. Additionally, Lee et al. [110] explored the performance of different CLIP model architectures in voxel grid representations, finding that an ensemble of models for guidance can prevent adversarial generations and improve geometrical structure, memory, and training speed.

In parallel, the application of 2D diffusion models for similar purposes, as discussed in [181], is introduced. Since NeRFs operate in image space, guiding a NeRF scene with the diffusion model involves practical solutions like deriving a Score Distillation loss or leveraging the training process in latent space, as seen in **DreamFusion** [168] and **Latent-NeRF** [143]. However, these approaches often result in unsatisfactory outputs and low diversity in objects generated from the same input text, coupled with lengthy synthesis times. Addressing these challenges, **DITTO-NeRF** [188] introduces progressive reconstruction schemes focusing on scales (from low to high resolution), angles (from inner to outer boundaries initially), and masks (from object to background boundary). This methodology achieves significant improvements in diversity and quality, as well as speed and fidelity of the generated objects, marking a significant progress in the field.

LaTeRF [148] enhances the NeRF framework by incorporating an “objectness” probability for each point, allowing the extraction of objects from scenes using pixel annotations. **SINE** [3] enhances semantic editing by introducing advanced methods: cyclic constraints alongside a proxy mesh for accurate geometric modifications, a color compositing system for better texture editing, and feature cluster-based regularization to manage the edited areas while maintaining the integrity of content that is not being edited. These enhancements facilitate compatibility with off-the-shelf text-prompt editing methods, enabling modifications to an object’s appearance and geometry, and the inpainting of missing parts of an object based on textual cues. **NeRF-Art** [228] and **Blending-NeRF** [199] integrate a pre-trained NeRF with an editable NeRF. The editable NeRF is trained to render a blended image that aligns with a target text, allowing precise editing of 3D object regions while preserving their original appearance.

Instruct-NeRF2NeRF [64] iteratively updates dataset images during NeRF model training with global text instructions from InstructPix2Pix [12]. This process, involving a loss that combines rays from various viewpoints, leads to higher quality results and more stable optimization.

Existing methods, however, face limitations in controlling individual objects within a scene. Modifying specific scene aspects without affecting others remains a challenge, and scene-level editing with long text prompts can lead to guidance collapse, preventing specific scene component edits. **CompoNeRF** [124] and **Set-the-Scene** [33] address these issues by employing a composition module to adjust text guidance levels and ensure the distinctiveness of entities while maintaining overall scene coherence. They represent the scene as a composition of multiple NeRFs, each optimized to represent specific objects “locally” and integrate seamlessly into the broader scene “globally”, thus eliminating guidance ambiguity. Through proxy manipulation, scenes can be decomposed and reassembled for editing without the need for additional fine-tuning.

3.6.2 Queryable Interaction

CLIP-Fields [190] integrates the strengths of the CLIP [170] image encoder, Sentence BERT [172], and NeRF to create a 3D scene representation that is queryable for mobile robots. This architecture is equipped with heads that output vectors corresponding to natural language descriptions, the visual appearance of objects, and the instance identification of every specific point in space. It uses two contrastive losses: one for the label token and another for the visual language embedding. CLIP-Fields demonstrate robustness in low-shot scenarios and label errors, capable of answering queries with varying degrees of real-life complexity. **VL-Fields** [219] aims to overcome the limitations of CLIP-Fields, which are restricted to a subset of scene points with known object classes. It proposes an open-set visual-language model that operates without prior knowledge of the object classes present in the scene.

LERF [97] uses a multi-scale feature pyramid combining 3D CLIP field and DINO [16] features to refine object boundaries for language query interactions. It allows for pixel-aligned queries of distilled 3D CLIP embeddings, bypassing the need

for region proposals, masks, or fine-tuning. LERF supports hierarchical, long-tail, open-vocabulary queries across the scene volume. **F3RM** [192] conducts few-shot learning experiments for grasping and placing tasks, drawing on Deep Fusion Field [104] (DFF) methodologies. This enables robots to perform 6-DoF object manipulation in response to natural language commands, exhibiting open-set generalization capabilities for handling unseen objects with significant differences. **GNFactor** [274] optimizes a generalizable NeRF for reconstruction alongside a Perceiver Transformer [82] for decision-making. This transformer integrates the robot’s proprioception and language features to execute decisions based on a Q-function [152], facilitating advanced decision-making processes in robotic applications.

4 Datasets and Evaluation

4.1 Core Metrics and Principles

4.1.1 Reconstruction and Novel View Synthesis

Image reconstruction and novel view synthesis in the standard setting use visual quality assessment metrics for benchmarks. The following metrics are the common standards in the NeRF literature:

Peak Signal-to-Noise Ratio (PSNR) quantifies the ratio of the maximum possible signal power, represented by the highest pixel intensity value, to the power of the noise corrupting the signal. A higher PSNR value indicates superior image quality. However, PSNR may not reliably reflect perceptual similarity since it fails to precisely represent how humans perceive image quality.

The **Structural Similarity Index Metric (SSIM)** [233] offers a perceptually more relevant evaluation by comparing two images through aspects such as luminance, contrast, and structural integrity. It considers variations in pixel intensities, spatial relationships, and texture contrasts. SSIM values span from -1 to 1, with 1 signifying an exact correspondence between the original and the reconstructed images. In terms of aligning with human visual perception, SSIM delivers a more accurate measure of image quality than PSNR.

The **Learned Perceptual Image Patch Similarity (LPIPS)** [280] metric assesses perceptual similarity between rendered views/poses and their corresponding ground truth images from specific viewing directions. Using deep learning, this perceptual metric measures the similarity between two images based on features extracted from a pre-trained Convolutional Neural Network (CNN), such as AlexNet or VGG, trained on the ImageNet dataset. Designed to more closely mirror human perception of image similarity, lower LPIPS scores indicate a greater perceptual similarity between the compared images. LPIPS proves to be especially effective in identifying subtle geometric and textural differences, making it particularly valuable for evaluating generative models and tasks related to image synthesis.

To enable easier comparison, an “average” error metric that summarizes all three above metrics is supplementarily presented [4]:

$$\text{Average} = \sqrt[3]{10^{-\text{PSNR}/10} \cdot \sqrt{1 - \text{SSIM}} \cdot \text{LPIPS}}$$

Fréchet Inception Distance (FID) [70] is a metric used to measure the similarity between the distribution of real images and the distribution of generated images in feature space. It uses the Inception-v3 [212] model to extract features from real and generated images. The FID score is calculated by computing the Fréchet distance between the multivariate Gaussian distributions of the feature representations of real and generated images. A lower FID score indicates that the generated images are more similar to the real images in terms of visual appearance and diversity.

Kernel Inception Distance (KID) [9] is an extension of the FID that aims to address some limitations of the FID. It measures the Maximum Mean Discrepancy between the feature distributions of real and generated images using kernel functions. KID focuses on a more robust and informative evaluation of image quality and diversity by considering the distributional properties of the features. It provides an unbiased estimation of the true distance between distributions of real and generated images, ensuring a more accurate representation of their similarity in feature space. Moreover, KID’s robustness to the choice of sample size minimizes the variability stemming from different sample sizes and requires fewer samples for calculation compared to alternative metrics.

4.1.2 Segmentation

Various evaluation metrics are employed to assess the performance of segmentation algorithms, quantifying the accuracy and reliability of the delineation between different regions in an image. Here are some of the common metrics used:

Pixel Accuracy computes the proportion of correctly classified pixels over the total number of pixels. It’s a simple and intuitive measure but might not capture the overall performance accurately, especially when dealing with imbalanced classes.

Mean Intersection over Union (mIoU), also referred to as the Jaccard index, is essentially a method to quantify the percent overlap between the target mask and our prediction output. mIoU is calculated by taking the IoU of each class and averaging them.

$$\text{IoU}(p, g) = \frac{|p \cap g|}{|p \cup g|} \quad (12)$$

Panoptic segmentation combines both semantic segmentation and instance segmentation. As a result, evaluation metrics for panoptic segmentation are crucial for quantitatively assessing the performance of algorithms that classify each pixel in an image into predefined classes or instance IDs and need to consider both aspects. **Panoptic Quality (PQ)** [101] is defined as the average IoU of the matched segments, while the denominator (see Equation 13) is designed to penalize segments without matches. PQ treats the quality of segmentation masks for all classes in an interpretable and unified manner, capturing all aspects of the task.

$$\text{PQ} = \frac{\sum_{(p,g) \in TP} \text{IoU}(p, g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|} \quad (13)$$

4.2 Public Datasets for Semantically-aware NeRFs

Existing datasets for novel view synthesis in the classical NeRF literature can be grouped into the following major categories:

- a Hemispherical 360° inward-facing views around an object of interest, which is mostly set against a plain white background (these include ShapeNet [21], CO3D [174], OmniObject3d [241], and Realistic Synthetic [145]).

- b Forward-facing scenes which aim the camera in a single direction and move in the vicinity facing the object (these include DTU [87, 1] and LLFF [144]).
- c Unbounded 360° real-scenes that provide full surrounding coverage with detailed backgrounds (these include Tanks and Templates [103] and MipNeRF360 [5] dataset).

Although fine-grained reconstruction is possible with provided camera intrinsics and poses, these datasets often lack compositional annotations (such as 3D bounding boxes or multi-object masks) and usually include a limited number of scenes. Efforts have been made to minimize photometric variation and avoid introducing multiple complex objects during capture. However, radiance field scene representations trained on these datasets typically focus on individual, per-scene optimization without additional semantic annotations or learning generalized priors. This makes it challenging to evaluate the performance of most semantically aware NeRFs.

Certain methods use hand-crafted annotations or pre-trained models to extract regions of interest from the scenes, but these approaches lack reliability as official benchmarks for comparing different methods. Therefore, in this section, we will discuss publicly available datasets that contain high-quality data with semantic annotations and which are the most relevant and most widely used in the literature.

4.2.1 Indoor Scenes

Scannet [35] is an RGB-D video dataset containing 2.5M views stemming from more than 1,500 scans, annotated with 3D camera poses, surface reconstructions, and instance-level semantic segmentations. It includes both 2D and 3D data and supports several 3D scene understanding tasks, including 3D object classification, semantic voxel labeling, and CAD model retrieval.

Replica [205] consists of 18 highly photorealistic 3D indoor scene reconstructions at room and building scale. Each scene consists of a dense mesh, high-resolution HDR textures, per-primitive semantic class and instance information, and planar mirror and glass reflectors.

Hypersim [178] is a photorealistic synthetic dataset for holistic indoor scene understanding. It contains 77,400 images of 461 indoor scenes

with detailed per-pixel labels and corresponding ground truth geometry, including complete scene geometry, material information, lighting information for every scene, dense per-pixel semantic instance segmentations, and complete camera information for every image.

HM3DSem [256] consists of 142,6K annotations of object instances in 216 spaces and 3,100 rooms within those spaces, built on top of Matterport 3D [20] for Embodied AI applications. A key difference setting from other datasets is the use of texture information to annotate pixel-accurate object boundaries.

The most recently referenced dataset, **ScanNet++** [266], offers high-resolution and high-quality RGBD captures, supporting the novel view synthesis task along with dense semantic annotations. It encompasses 460 scenes, featuring 280K DSLR images and more than 3.7M iPhone RGBD frames.

PeRFception [89] uses radiance fields (Plenoxels [269, 43]) as another representation of data that effectively conveys the same information for both 2D and 3D in a unified and compressed model, eliminating the need to store different data formats separately. At the moment, PeRFception-CO3D and PeRFception-ScanNet are created, which have covered object-centric and scene-centric environments respectively.

To tackle the problems of collecting, processing, and annotating datasets for 3D scene understanding at scale, **Kubric** [57] is introduced as a framework for generating synthetic datasets with fine-grain control over data complexity and rich ground truth annotations. The pipeline is linked with an open-source Python framework and Blender, allowing facilitating reuse of data-generation code, across multiple scales. Furthermore, Kubric can provide various randomization options for custom use cases. There have been many papers that applied this framework to create their own datasets [226, 240, 6, 52]. However, most of the collected datasets are new in the field and limited within the approaches of the articles without proper benchmarks, they still remain as important parts and are waiting to be tested by the community.

Datasets	Venue	#Scenes	#Imgs	Centricity	Type	Data Modalities	Annotations	URL
3DMV-VQA [71]	CVPR 2023	5000	600K	S+O	Indoor	RGB	Visual question & answer	🔗
NeRDS 360 [81]	ICCV 2023	75	15k	S+O	Urban	Synthetic	3D object boxes 2D panoptic segmentation	🔗
ScanNet++ [266]	ICCV 2023	460	3.7M	S	Indoor	RGB-D	2D/3D panoptic segmentation	🔗
KITTI-360 [120]	PAMI 2022	10	150K	S+O	Urban	RGB & LiDAR	2D/3D object boxes 2D panoptic segmentation	🔗
SHIFT [210]	CVPR 2022	4850	2.5M	S+O	Urban	Synthetic	2D/3D object boxes 2D panoptic segmentation	🔗
HM3D Sem [256]	arXiv 2022	216	-	S	Indoor	Mesh	3D semantic segmentation	🔗
3D-FRONT [44]	ICCV 2021	18968	-	S+O	Indoor	Synthetic	3D semantic segmentation	🔗
HyperSim [178]	ICCV 2021	461	77.4K	S+O	Indoor	Synthetic	2D/3D object boxes 2D/3D panoptic segmentation	🔗
Waymo [209]	CVPR 2020	1150	1M	S+O	Urban	RGB & LiDAR	2D/3D object boxes 2D panoptic segmentation	🔗
nuScenes [13]	CVPR 2020	1000	1.4M	S+O	Urban	RGB & LiDAR	3D object boxes 2D semantic segmentation	🔗
Replica [205]	arXiv 2019	18	-	S	Indoor	Mesh	2D/3D panoptic segmentation	🔗
Matterport 3D [20]	3DV 2017	90	194.4K	S	Indoor	RGB-D	2D/3D panoptic segmentation	🔗
CLEVR [91]	CVPR 2017	-	100K	O	Indoor	Synthetic	Visual question & answer	🔗
ScanNet [35]	CVPR 2017	1513	2.5M	S+O	Indoor	RGB-D	3D object boxes 2D/3D panoptic segmentation	🔗
Virtual KITTI [48]	CVPR 2016	5	17K	S+O	Urban	Synthetic	2D/3D object boxes 2D panoptic segmentation	🔗
SUN RGB-D [202]	CVPR 2015	47	10.3K	S+O	Indoor	RGB-D	2D/3D object boxes 2D panoptic segmentation	🔗
Shapenet [21]	arXiv 2015	-	-	O	Objects	CAD model	3D part segmentation	🔗
KITTI [54, 55]	CVPR 2012	22	15K	S+O	Urban	RGB & LiDAR	2D/3D object boxes 2D panoptic segmentation	🔗

Table 2: Overview of existing datasets for SRF-based multi-view scene understanding. ‘Centricity’ refers to scene and/or object-centric datasets, respectively denoted with S and O above.

4.2.2 Outdoor Urban Scenes

The **KITTI** [54, 55] dataset is a renowned collection tailored for computer vision research in the context of urban-scale 2D-3D environments, specifically designed to train and evaluate algorithms aimed at autonomous driving technologies. This dataset was compiled using raw LiDAR and video data collected in Karlsruhe, Germany, employing a vehicle-mounted system equipped with GPS and an inertial measurement unit. To accommodate various research objectives, parts of the dataset have been manually annotated by researchers, making KITTI a comprehensive resource that includes labeled data for a range of tasks such as stereo 2D-3D segmentation, optical

flow, odometry, 2D-3D object detection, tracking, lane detection, and depth prediction/completion. However, the absence of complete semantic labeling limits its use for tasks like synthesizing new view images or constructing large-scale semantic maps, as these activities require fully labeled datasets for accurate evaluation.

Other datasets such as **nuScenes** [13], **Waymo** [209] try to address this shortcoming by providing more comprehensive data with semantic/instance labels in 2D and 3D, and richer 360° sensory information corresponding to longer driving logs with more accurate and geolocalized vehicle poses. Especially, **KITTI-360** [120] with

its 3D-to-2D label transfer that opens more interesting tasks, e.g., semantic SLAM or novel view semantic synthesis.

Adapting to a continuously evolving environment is a safety-critical challenge inevitably faced by all autonomous driving systems. Existing image and video-driving datasets, however, fall short of capturing the mutable nature of the real world. In other words, they are captured under approximately stationary conditions. **Virtual KITTI** [48] and **SHIFT** [210] captures these driving scenarios in various environmental directions: time of day, cloudiness, rain, fog strength and vehicle and pedestrian density with more detailed object class annotations (persons, cars, license plates, ...) in separate discrete variations [48] or continuously shifting conditions [210].

NERDS 360 [81] is a large-scale dataset for 3D urban scene understanding. This dataset consists of 75 outdoor urban scenes with diverse backgrounds in 360° hemispherical views, encompassing over 15K images. The dataset and the corresponding tasks are extremely challenging due to occlusions, diversity of backgrounds, and rendered objects with various lightning and shadows.

4.2.3 Vision and Language

CLEVR [91] is a diagnostic dataset for studying the ability of Visual Question Answering (VQA) systems. It contains 100K rendered images and 853K generated unique questions for visual reasoning abilities such as counting, comparing, logical reasoning, and storing information. Each object present in the scene, aside from position, is characterized by a set of four attributes: 2 sizes: large, and small, 3 shapes: square, cylinder, and sphere, 2 material types: rubber, and metal, 8 color types: gray, blue, brown, yellow, red, green, purple, and cyan, resulting in 96 unique combinations.

3DMV-VQA [71] consists of approximately 5K scenes, and 600K images, paired with 50K questions in total. This dataset is built on top of the HM3DSem dataset [256] with four types of questions: conceptual, counting, relational, and comparison questions. The authors further propose a 3D concept learning and reasoning framework that is grounded on open-vocabulary semantic concepts on 3D representation.

5 Challenges and Perspectives

To progress in the field of semantically-aware NeRFs, targeted research efforts are essential. This section outlines the primary challenges and opportunities for enhancement that we have identified as critical focus areas.

i Scene Generalizability. Current Semantically-aware NeRF (SRF) methods, capable of processing datasets without the need for scene-specific training or optimization, mark a significant progress over the original NeRF methodology [145], which lacked any capability for cross-dataset generalization. Despite these improvements, there are still clear limitations. Current approaches might necessitate expensive, dense semantic annotations [81], require substantial data volumes [185], predominantly operate within synthetic settings, or produce blurriness in novel view synthesis, often attributed to L2 loss training [185, 268]. These challenges are further amplified by the variation in viewpoint densities during test-time, affecting traditional performance and efficiency. Additionally, while some strategies employ pre-trained, sparse detectors and segmentation networks, they tend to achieve only object-centric generalization [155, 61]. Addressing these challenges to improve cross-dataset generalization, or combining their respective benefits, would represent a substantial leap forward, enabling true real-time applications from acquisition to rendering.

ii Camera Calibration. While some NeRF-based methods are designed to take unposed images and simultaneously recover their extrinsic matrices [88, 234], most NeRF derivatives work under the assumption that the RGB views provided as input are already posed. As a result, even minor calibration errors can lead to significant semantic misalignments across different views. Such misalignments can cause early failures in the process, which are often irreversible during the subsequent training or scene optimization phases. Therefore, there is a clear necessity within the NeRF domain as a whole to not only improve calibration techniques but also to develop specific mechanisms for pose refinement during the training process.

iii **Data Efficiency and Augmentation.**

Addressing the data efficiency challenge is essential for making NeRF more practical in real-world settings. Future work may involve exploring methods to train accurate semantically-aware models with less training data, and fewer annotations, making them more accessible for a broader range of applications, in particular in one/few-shot settings in real-world environments. The successful integration of semantic understanding into NeRFs has the potential to dramatically enhance applications in augmented reality, autonomous navigation, and beyond, by providing more meaningful and context-aware interpretations of 3D environments. These functionalities make NeRF well-suited to serve as the foundation for various components. By understanding the decomposition of the scene and allowing dynamic adjustments to the simulation environment, NeRF becomes a valuable asset in creating realistic scenarios for closed-loop simulations, providing a crucial element in training and testing scenarios for many systems. Additionally, its adaptability for data augmentation enhances its utility in improving the robustness and generalization of machine learning models.

iv **Multi-modal, Multi-task, and Efficient Scene Understanding.**

Currently, the majority of multi-modal approaches investigated within the NeRF domain are centered around textTo3D. Despite the wide range of potential multi-task combinations available [272, 2], many remain largely unexplored, representing missed opportunities to discover new, mutually informative tasks specifically within the area of Radiance Fields. For instance, areas such as sound processing or other types of inputs [138] have yet to be fully explored.

v **Real-Time and Mobile Performance.**

NeRF encounters challenges in terms of computational efficiency, particularly due to specialized volumetric rendering algorithms mismatched to widely deployed graphics hardware. These computationally intensive methods often necessitate extended rendering times and substantial resources, hindering real-time applications. To address this, exploring alternative data structures or rendering techniques, especially those suitable for low computational mobile devices, presents a promising avenue.

e.g., 3D Gaussian Splatting [96] (3DGS) with unprecedented rendering efficiency and other non-semantically aware strategies could similarly serve as backbone models to SRFs.

vi **Ethical Concerns and Societal Impact.**

The generative capabilities of editable NeRFs, allowing for the creation of photorealistic 3D objects, humans, and scenes not previously seen, may lead to challenges akin to those seen with DeepFakes [194] in 2D image generation. These potential issues require similar scrutiny and efforts to mitigate. Conversely, the generative and editable nature of these methods could offer substantial opportunities for 3D enthusiasts and content creators, attributed to their user-friendly design.

vii **Performance Evaluation.**

Current metrics, such as those used for novel view synthesis, are well-established in the field. However, these metrics are decoupled from human perception, meaning that quantitative evaluations cannot guarantee objective optimality in a way that aligns with human assessment. Multi-task models also face issues due to the absence of composite metrics, and instead rely on disjoint, linear combination metrics. Existing learned perceptual metrics, e.g., LPIPS, are restricted to the evaluation of static image frames [280]. They do not consider video or 3D consistency [118] in terms of shape, appearance, and semantics. This is a necessary research opportunity to evaluate complex 3D environments that are typically dynamic.

viii **Collaborative Frameworks.**

Recognizing the difficulties arising from scattered codebases and the lack of consolidated support, Nerfstudio [214] is an end-to-end framework which aggregates modular plug-and-play components, e.g., viewers, algorithms, datasets, and benchmarking tools. This facilitates the integration of features across diverse implementations, simplifies the collaborative process for researchers and practitioners, and enhances accessibility through real-time visualization tools that natively support semantic information. Providing a cohesive and extensible platform, such frameworks can foster collaboration, accelerate progress, and contribute to the advancement of NeRF research more efficiently and cohesively.

6 Conclusion

We have conducted the first survey on Neural Radiance Fields (NeRFs), specifically focusing on semantically-aware NeRFs. Our comprehensive review has shed light on the state-of-the-art methodologies, challenges, and a wide array of applications. It also highlights the need for further advancements in this field, to enable more sophisticated, efficient, and context-aware interpretations of 3D scenes by achieving the full potential of NeRFs. This will pave the way for truly real-time end-to-end applications from acquisition to rendering, on commodity hardware.

References

- [1] Henrik Aanæs et al. “Large-Scale Data for Multiple-View Stereopsis”. In: *International Journal of Computer Vision* (2016), pp. 1–16.
- [2] Andrei Atanov et al. “Task discovery: Finding the tasks that neural networks generalize on”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 15702–15717.
- [3] Chong Bao et al. “SINE: Semantic-driven Image-based NeRF Editing with Prior-guided Editing Field”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 20919–20929.
- [4] Jonathan T Barron et al. “Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5855–5864.
- [5] Jonathan T Barron et al. “Mip-nerf 360: Unbounded anti-aliased neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 5470–5479.
- [6] Yash Bhalgat et al. “Contrastive Lift: 3D Object Instance Segmentation by Slow-Fast Contrastive Fusion”. In: *arXiv preprint arXiv:2306.04633* (2023).
- [7] Anand Bhattad et al. “View generalization for single image textured 3d models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 6081–6090.
- [8] WANG Bing, Lu Chen, and Bo Yang. “DM-NeRF: 3D Scene Geometry Decomposition and Manipulation from 2D Images”. In: *The Eleventh International Conference on Learning Representations*. 2022.
- [9] Mikołaj Bińkowski et al. “Demystifying MMD GANs”. In: *International Conference on Learning Representations*. 2018.
- [10] Kenneth Blomqvist et al. “Baking in the feature: Accelerating volumetric segmentation by rendering feature maps”. In: *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2023, pp. 7629–7634.
- [11] Kenneth Tor Blomqvist et al. “Grounding Pretrained Features in 3D Representations”. In: *ICRA2023 Workshop on Pretraining for Robotics (PT4R)*. 2023.
- [12] Tim Brooks, Aleksander Holynski, and Alexei A Efros. “Instructpix2pix: Learning to follow image editing instructions”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 18392–18402.
- [13] Holger Caesar et al. “nuscenes: A multimodal dataset for autonomous driving”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020, pp. 11621–11631.
- [14] Anh-Quan Cao and Raoul de Charette. “Scenerf: Self-supervised monocular 3d scene reconstruction with radiance fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 9387–9398.
- [15] Chenjie Cao and Yanwei Fu. “Learning a Sketch Tensor Space for Image Inpainting of Man-made Scenes. 2021 IEEE”. In: *CVF International Conference on Computer Vision, ICCV*. 2021, pp. 10–17.
- [16] Mathilde Caron et al. “Emerging properties in self-supervised vision transformers”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 9650–9660.
- [17] Jiazhong Cen et al. “Segment anything in 3d with nerfs”. In: *arXiv preprint arXiv:2304.12308* (2023).
- [18] Eric R Chan et al. “pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2021, pp. 5799–5809.
- [19] Eric R Chan et al. “Efficient geometry-aware 3D generative adversarial networks”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 16123–16133.
- [20] Angel Chang et al. “Matterport3D: Learning from RGB-D Data in Indoor Environments”.

- In: *International Conference on 3D Vision (3DV)* (2017).
- [21] Angel X Chang et al. “Shapenet: An information-rich 3d model repository”. In: *arXiv preprint arXiv:1512.03012* (2015).
- [22] Anpei Chen et al. “Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 14124–14133.
- [23] Di Chen et al. “Geoaug: Data augmentation for few-shot nerf with geometry constraints”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 322–337.
- [24] Haoran Chen et al. “Panoptic Vision-Language Feature Fields”. In: *IEEE Robotics and Automation Letters* (2024).
- [25] Jianchuan Chen et al. “GM-NeRF: Learning Generalizable Model-based Neural Radiance Fields from Multi-view Images”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 20648–20658.
- [26] Linghao Chen et al. “Perceiving Unseen 3D Objects by Poking the Objects”. In: *ICRA*. 2023.
- [27] Shuhong Chen et al. “PAniC-3D: Stylized Single-view 3D Reconstruction from Portraits of Anime Characters”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 21068–21077.
- [28] Xinya Chen et al. “VeRi3D: Generative Vertex-based Radiance Fields for 3D Controllable Human Image Synthesis”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 8986–8997.
- [29] Yurui Chen et al. “Single-view Neural Radiance Fields with Depth Teacher”. In: *arXiv preprint arXiv:2303.09952* (2023).
- [30] Zhiqin Chen et al. “Mobilerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 16569–16578.
- [31] Xinhua Cheng et al. “Panoptic Compositional Feature Field for Editable Scene Rendering With Network-Inferred Labels via Metric Learning”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 4947–4957.
- [32] Julian Chibane et al. “Stereo radiance fields (srf): Learning view synthesis for sparse views of novel scenes”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 7911–7920.
- [33] Dana Cohen-Bar et al. “Set-the-Scene: Global-Local Training for Generating Controllable NeRF Scenes”. In: *arXiv preprint arXiv:2303.13450* (2023).
- [34] James Coughlan and Alan L Yuille. “The manhattan world assumption: Regularities in scene statistics which enable bayesian inference”. In: *Advances in Neural Information Processing Systems* 13 (2000).
- [35] Angela Dai et al. “Scannet: Richly-annotated 3d reconstructions of indoor scenes”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 5828–5839.
- [36] Matt Deitke et al. “ProcTHOR: Large-Scale Embodied AI Using Procedural Generation”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 5982–5994.
- [37] Congyue Deng et al. “Nerdi: Single-view nerf synthesis with language-guided diffusion as general image priors”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 20637–20647.
- [38] Kangle Deng et al. “Depth-supervised nerf: Fewer views and faster training for free”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 12882–12891.
- [39] Yu Deng, Baoyuan Wang, and Heung-Yeung Shum. “Learning Detailed Radiance Manifolds for High-Fidelity and 3D-Consistent Portrait Synthesis from Monocular Image”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.
- [40] Yu Deng et al. “Gram: Generative radiance manifolds for 3d-aware image generation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10673–10683.
- [41] Zhiwen Fan et al. “Nerf-sos: Any-view self-supervised object segmentation on complex scenes”. In: *arXiv preprint arXiv:2209.08776* (2022).
- [42] Chris Fifty et al. “Efficiently identifying task groupings for multi-task learning”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 27503–27516.

- [43] Sara Fridovich-Keil et al. “Plenoxels: Radiance fields without neural networks”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 5501–5510.
- [44] Huan Fu et al. “3d-front: 3d furnished rooms with layouts and semantics”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 10933–10942.
- [45] Xiao Fu et al. “Panoptic NeRF: 3D-to-2D Label Transfer for Panoptic Urban Scene Segmentation”. In: *International Conference on 3D Vision (3DV)*. 2022.
- [46] Xiao Fu et al. “PanopticNeRF-360: Panoramic 3D-to-2D Label Transfer in Urban Scenes”. In: *arXiv preprint arXiv:2309.10815* (2023).
- [47] Guy Gafni et al. “Dynamic neural radiance fields for monocular 4d facial avatar reconstruction”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 8649–8658.
- [48] Adrien Gaidon et al. “Virtual worlds as proxy for multi-object tracking analysis”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 4340–4349.
- [49] Stathis Galanakis et al. “3DMM-RF: Convolutional Radiance Fields for 3D Face Modeling”. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, pp. 3536–3547.
- [50] Chen Gao et al. “Dynamic view synthesis from dynamic monocular video”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5712–5721.
- [51] Kyle Gao et al. “Nerf: Neural radiance field in 3d vision, a comprehensive review”. In: *arXiv preprint arXiv:2210.00379* (2022).
- [52] Siyu Gao et al. “Object-Centric Voxelization of Dynamic Scenes via Inverse Neural Rendering”. In: *arXiv preprint arXiv:2305.00393* (2023).
- [53] Stephan J Garbin et al. “Fastnerf: High-fidelity neural rendering at 200fps”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 14346–14355.
- [54] Andreas Geiger, Philip Lenz, and Raquel Urtasun. “Are we ready for autonomous driving? the kitti vision benchmark suite”. In: *2012 IEEE conference on computer vision and pattern recognition*. IEEE. 2012, pp. 3354–3361.
- [55] Andreas Geiger et al. “Vision meets robotics: The kitti dataset”. In: *The International Journal of Robotics Research* 32.11 (2013), pp. 1231–1237.
- [56] Rahul Goel et al. “Interactive segmentation of radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 4201–4211.
- [57] Klaus Greff et al. “Kubric: A scalable dataset generator”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 3749–3761.
- [58] Jiaming Gu et al. “UE4-NeRF: Neural Radiance Field for Real-Time Rendering of Large-Scale Scene”. In: *Thirty-seventh Conference on Neural Information Processing Systems*. 2023.
- [59] Jiatao Gu et al. “Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis”. In: *arXiv preprint arXiv:2110.08985* (2021).
- [60] Haoyu Guo et al. “Neural 3d scene reconstruction with the manhattan-world assumption”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 5511–5520.
- [61] Michelle Guo et al. “Object-centric neural scene rendering”. In: *arXiv preprint arXiv:2012.08503* (2020).
- [62] Yulan Guo et al. “Deep learning for 3d point clouds: A survey”. In: *IEEE transactions on pattern analysis and machine intelligence* 43.12 (2020), pp. 4338–4364.
- [63] Yasaman Haghighi et al. “Neural Implicit Dense Semantic SLAM”. In: *arXiv preprint arXiv:2304.14560* (2023).
- [64] Ayaan Haque et al. “Instruct-nerf2nerf: Editing 3d scenes with instructions”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023.
- [65] Iain Haughton et al. “Real-time Mapping of Physical Scene Properties with an Autonomous Robot Experimenter”. In: *Conference on Robot Learning*. PMLR. 2023, pp. 118–127.
- [66] Adrian Hayler et al. “S4C: Self-Supervised Semantic Scene Completion with Neural Fields”. In: *arXiv preprint arXiv:2310.07522* (2023).
- [67] Kaiming He et al. “Mask r-cnn”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2961–2969.
- [68] Peter Hedman et al. “Baking neural radiance fields for real-time view synthesis”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5875–5884.

- [69] Quentin Herau et al. “MOISST: Multi-modal Optimization of Implicit Scene for SpatioTemporal calibration”. In: *arXiv preprint arXiv:2303.03056* (2023).
- [70] Martin Heusel et al. “Gans trained by a two time-scale update rule converge to a local nash equilibrium”. In: *Advances in neural information processing systems* 30 (2017).
- [71] Yining Hong et al. “3D Concept Learning and Reasoning from Multi-View Images”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 9202–9212.
- [72] Benran Hu et al. “NeRF-RPN: A general framework for object detection in NeRFs”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 23528–23538.
- [73] Shoukang Hu et al. “SHERF: Generalizable Human NeRF from a Single Image”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023.
- [74] Tao Hu et al. “Efficientnerf efficient neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 12902–12911.
- [75] Lin Huang et al. “Neural correspondence field for object pose estimation”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 585–603.
- [76] Shi-Sheng Huang et al. “Real-time globally consistent 3D reconstruction with semantic priors”. In: *IEEE transactions on visualization and computer graphics* (2021).
- [77] Shi-Sheng Huang et al. “Supervoxel convolution for online 3d semantic segmentation”. In: *ACM Transactions on Graphics (TOG)* 40.3 (2021), pp. 1–15.
- [78] Xiaoyang Huang et al. “Boosting point clouds rendering via radiance mapping”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 37. 2023, pp. 953–961.
- [79] Sungwon Hwang et al. “FaceCLIPNeRF: Text-driven 3D Face Manipulation using Deformable Neural Radiance Fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 3469–3479.
- [80] Muhammad Zubair Irshad et al. “Shapo: Implicit representations for multi-object shape, appearance, and pose optimization”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 275–292.
- [81] Muhammad Zubair Irshad et al. “NeO 360: Neural Fields for Sparse View Synthesis of Outdoor Scenes”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 9187–9198.
- [82] Andrew Jaegle et al. “Perceiver: General perception with iterative attention”. In: *International conference on machine learning*. PMLR. 2021, pp. 4651–4664.
- [83] Ajay Jain, Matthew Tancik, and Pieter Abbeel. “Putting nerf on a diet: Semantically consistent few-shot view synthesis”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5885–5894.
- [84] Ajay Jain et al. “Zero-shot text-guided object generation with dream fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 867–876.
- [85] Wonbong Jang and Lourdes Agapito. “Codenerf: Disentangled neural radiance fields for object categories”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 12949–12958.
- [86] Vinoj Jayasundara et al. “FlexNeRF: Photorealistic Free-viewpoint Rendering of Moving Humans from Sparse Views”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.
- [87] Rasmus Jensen et al. “Large scale multi-view stereopsis evaluation”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014, pp. 406–413.
- [88] Yoonwoo Jeong et al. “Self-calibrating neural radiance fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5846–5854.
- [89] Yoonwoo Jeong et al. “PeRFception: Perception using radiance fields”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 26105–26121.
- [90] Kyungmin Jo et al. “CG-NeRF: Conditional Generative Neural Radiance Fields for 3D-aware Image Synthesis”. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, pp. 724–733.
- [91] Justin Johnson et al. “Clevr: A diagnostic dataset for compositional language and elementary visual reasoning”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 2901–2910.
- [92] James T Kajiya and Brian P Von Herzen. “Ray tracing volume densities”. In: *ACM*

- SIGGRAPH computer graphics* 18.3 (1984), pp. 165–174.
- [93] Angjoo Kanazawa et al. “Learning category-specific mesh reconstruction from image collections”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018, pp. 371–386.
- [94] Kacper Kania et al. “Conerf: Controllable neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 18623–18632.
- [95] Tero Karras et al. “Analyzing and improving the image quality of stylegan”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020, pp. 8110–8119.
- [96] Bernhard Kerbl et al. “3D Gaussian Splatting for Real-Time Radiance Field Rendering”. In: *ACM Transactions on Graphics* 42.4 (2023).
- [97] Justin Kerr et al. “Lerf: Language embedded radiance fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 19729–19739.
- [98] Salman H Khan et al. “Integrating geometrical context for semantic labeling of indoor scenes using rgbd images”. In: *International Journal of Computer Vision* 117 (2016), pp. 1–20.
- [99] Mijeong Kim, Seonguk Seo, and Bohyung Han. “Infonerf: Ray entropy minimization for few-shot neural volume rendering”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 12912–12921.
- [100] Mira Kim et al. “Ae-nerf: Auto-encoding neural radiance fields for 3d-aware object manipulation”. In: *arXiv preprint arXiv:2204.13426* (2022).
- [101] Alexander Kirillov et al. “Panoptic segmentation”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019, pp. 9404–9413.
- [102] Alexander Kirillov et al. “Segment anything”. In: *arXiv preprint arXiv:2304.02643* (2023).
- [103] Arno Knapitsch et al. “Tanks and temples: Benchmarking large-scale scene reconstruction”. In: *ACM Transactions on Graphics (ToG)* 36.4 (2017), pp. 1–13.
- [104] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. “Decomposing NeRF for Editing via Feature Field Distillation”. In: *Advances in Neural Information Processing Systems*. 2022.
- [105] Xin Kong et al. “vmap: Vectorised object mapping for neural field slam”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 952–961.
- [106] Akshay Krishnan et al. “LANe: Lighting-Aware Neural Fields for Compositional Scene Synthesis”. In: *arXiv preprint arXiv:2304.03280* (2023).
- [107] Abhijit Kundu et al. “Panoptic neural fields: A semantic object-aware neural scene representation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 12871–12881.
- [108] Minseop Kwak, Jiuhn Song, and Seung Wook Kim. “GeCoNeRF: Few-shot Neural Radiance Fields via Geometric Consistency”. In: *ICML abs/2301.10941* (2023).
- [109] Verica Lazova et al. “Control-nerf: Editable feature volumes for scene rendering and manipulation”. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, pp. 4340–4350.
- [110] Han-Hung Lee and Angel X Chang. “Understanding pure clip guidance for voxel grid nerf models”. In: *arXiv preprint arXiv:2209.15172* (2022).
- [111] Boyi Li et al. “Language-driven Semantic Segmentation”. In: *International Conference on Learning Representations*. 2022.
- [112] Chenghao Li et al. “Generative AI meets 3D: A Survey on Text-to-3D in AIGC Era”. In: *arXiv preprint arXiv:2305.06131* (2023).
- [113] Jiaxin Li et al. “Mine: Towards continuous depth mpi with nerf for novel view synthesis”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 12578–12588.
- [114] Liunian Harold Li et al. “Grounded language-image pre-training”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10965–10975.
- [115] Ruilong Li, Matthew Tancik, and Angjoo Kanazawa. “Nerfacc: A general nerf acceleration toolbox”. In: *arXiv preprint arXiv:2210.04847* (2022).
- [116] Tianye Li et al. “Neural 3d video synthesis from multi-view video”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 5521–5531.
- [117] Zhengqi Li et al. “Neural scene flow fields for space-time view synthesis of dynamic scenes”. In: *Proceedings of the IEEE/CVF Conference*

- on *Computer Vision and Pattern Recognition*. 2021, pp. 6498–6508.
- [118] Hanxue Liang et al. “Perceptual Quality Assessment of NeRF and Neural View Synthesis Methods for Front-Facing Views”. In: *arXiv preprint arXiv:2303.15206* (2023).
- [119] Yiqing Liang et al. “Semantic Attention Flow Fields for Monocular Dynamic Scene Decomposition”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 21797–21806.
- [120] Yiyi Liao, Jun Xie, and Andreas Geiger. “KITTI-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45.3 (2022), pp. 3292–3310.
- [121] Chen-Hsuan Lin et al. “Barf: Bundle-adjusting neural radiance fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5741–5751.
- [122] Kai-En Lin et al. “Vision transformer for nerf-based view synthesis from a single input image”. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, pp. 806–815.
- [123] Tsung-Yi Lin et al. “Feature pyramid networks for object detection”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 2117–2125.
- [124] Yiqi Lin et al. “Componerf: Text-guided multi-object compositional nerf with editable 3d scene layout”. In: *arXiv preprint arXiv:2303.13843* (2023).
- [125] Philipp Lindenberger et al. “Pixel-perfect structure-from-motion with featuremetric refinement”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 5987–5997.
- [126] Fangfu Liu et al. “Semantic Ray: Learning a Generalizable Semantic Field with Cross-Reprojection Attention”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 17386–17396.
- [127] Hao-Kang Liu, I Shen, Bing-Yu Chen, et al. “NeRF-In: Free-form NeRF inpainting with RGB-D priors”. In: *arXiv preprint arXiv:2206.04901* (2022).
- [128] Kunhao Liu et al. “Weakly Supervised 3D Open-vocabulary Segmentation”. In: *Thirty-seventh Conference on Neural Information Processing Systems*. 2023.
- [129] Lingjie Liu et al. “Neural sparse voxel fields”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 15651–15663.
- [130] Lingjie Liu et al. “Neural actor: Neural free-view synthesis of human actors with pose control”. In: *ACM transactions on graphics (TOG)* 40.6 (2021), pp. 1–16.
- [131] Yu-Lun Liu et al. “Robust dynamic radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 13–23.
- [132] Ruoshi Liu et al. “Zero-1-to-3: Zero-shot one image to 3d object”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 9298–9309.
- [133] Steven Liu et al. “Editing conditional radiance fields”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 5773–5783.
- [134] Tianyu Liu et al. “Car-Studio: Learning Car Radiance Fields from Single-View and Endless In-the-wild Images”. In: *arXiv preprint arXiv:2307.14009* (2023).
- [135] Xinhang Liu et al. “Unsupervised multi-view object segmentation using radiance field propagation”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 17730–17743.
- [136] Yichen Liu et al. “Instance neural radiance field”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 787–796.
- [137] Zhizheng Liu et al. “Unsupervised Continual Semantic Adaptation through Neural Rendering”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 3031–3040.
- [138] Andrew Luo et al. “Learning neural acoustic fields”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 3165–3177.
- [139] Jitendra Malik et al. “The three R’s of computer vision: Recognition, reconstruction and reorganization”. In: *Pattern Recognition Letters* 72 (2016), pp. 4–14.
- [140] Ricardo Martin-Brualla et al. “Nerf in the wild: Neural radiance fields for unconstrained photo collections”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 7210–7219.
- [141] Kirill Mazur, Edgar Sucar, and Andrew J Davison. “Feature-realistic neural fusion for real-time, open set scene understanding”. In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2023, pp. 8201–8207.

- [142] Luke Melas-Kyriazi et al. “Realfusion: 360deg reconstruction of any object from a single image”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 8446–8455.
- [143] Gal Metzer et al. “Latent-nerf for shape-guided generation of 3d shapes and textures”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 12663–12673.
- [144] Ben Mildenhall et al. “Local Light Field Fusion: Practical View Synthesis with Prescriptive Sampling Guidelines”. In: *ACM Transactions on Graphics (TOG)* (2019).
- [145] Ben Mildenhall et al. “NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis”. In: *European Conference on Computer Vision*. Springer. 2020, pp. 405–421.
- [146] Ben Mildenhall et al. “Nerf in the dark: High dynamic range view synthesis from noisy raw images”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 16190–16199.
- [147] Zhixiang Min et al. “NeurOCS: Neural NOCS Supervision for Monocular 3D Object Localization”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 21404–21414.
- [148] Ashkan Mirzaei et al. “Laterf: Label and text driven object radiance fields”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 20–36.
- [149] Ashkan Mirzaei et al. “Reference-guided Controllable Inpainting of Neural Radiance Fields”. In: *arXiv preprint arXiv:2304.09677* (2023).
- [150] Ashkan Mirzaei et al. “SPIn-NeRF: Multiview segmentation and perceptual inpainting with neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 20669–20679.
- [151] Ansh Mittal. “Neural Radiance Fields: Past, Present, and Future”. In: *arXiv preprint arXiv:2304.10050* (2023).
- [152] Volodymyr Mnih et al. “Playing atari with deep reinforcement learning”. In: *arXiv preprint arXiv:1312.5602* (2013).
- [153] Pierre Moulon et al. “Openmvg: Open multiple view geometry”. In: *Reproducible Research in Pattern Recognition: First International Workshop, RRPR 2016, Cancún, Mexico, December 4, 2016, Revised Selected Papers 1*. Springer. 2017, pp. 60–74.
- [154] Jiteng Mu et al. “ActorsNeRF: Animatable Few-shot Human Rendering with Generalizable NeRFs”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023.
- [155] Norman Müller et al. “Autorf: Learning 3d object radiance fields from single view observations”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 3971–3980.
- [156] Thomas Müller et al. “Instant neural graphics primitives with a multiresolution hash encoding”. In: *ACM Transactions on Graphics (ToG)* 41.4 (2022), pp. 1–15.
- [157] Michael Niemeyer and Andreas Geiger. “Giraffe: Representing scenes as compositional generative neural feature fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 11453–11464.
- [158] Michael Niemeyer et al. “Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020, pp. 3504–3515.
- [159] Michael Niemeyer et al. “Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 5480–5490.
- [160] Julian Ost et al. “Neural scene graphs for dynamic scenes”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 2856–2865.
- [161] Xuran Pan et al. “Activenerf: Learning where to see with uncertainty estimation”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 230–246.
- [162] Michael Pantic et al. “Sampling-free obstacle gradients and reactive planning in Neural Radiance Fields”. In: *Workshop on Motion Planning with Implicit Neural Representations of Geometry” at 2022 IEEE International Conference on Robotics and Automation (ICRA 2022)*. 2022.
- [163] Keunhong Park et al. “Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields”. In: *arXiv preprint arXiv:2106.13228* (2021).
- [164] Sungheon Park et al. “Temporal Interpolation Is All You Need for Dynamic Neural Radiance Fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 4212–4221.

- [165] Dario Pavlo et al. “Shape, Pose, and Appearance from a Single Image via Bootstrapped Radiance Field Inversion”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.
- [166] Martin Píala and Ronald Clark. “Terminerf: Ray termination prediction for efficient neural rendering”. In: *2021 International Conference on 3D Vision (3DV)*. IEEE. 2021, pp. 1106–1114.
- [167] Giovanni Pintore et al. “State-of-the-art in automatic 3D reconstruction of structured indoor environments”. In: *Computer Graphics Forum*. Vol. 39. 2. Wiley Online Library. 2020, pp. 667–699.
- [168] Ben Poole et al. “Dreamfusion: Text-to-3d using 2d diffusion”. In: *arXiv preprint arXiv:2209.14988* (2022).
- [169] AKM Rabby and Chengcui Zhang. “Beyond-Pixels: A Comprehensive Review of the Evolution of Neural Radiance Fields”. In: *arXiv preprint arXiv:2306.03000* (2023).
- [170] Alec Radford et al. “Learning transferable visual models from natural language supervision”. In: *International conference on machine learning*. PMLR. 2021, pp. 8748–8763.
- [171] Daniel Rebaín et al. “LOLNeRF: Learn from One Look”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 1558–1567.
- [172] Nils Reimers and Iryna Gurevych. “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics. 2019.
- [173] Christian Reiser et al. “Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 14335–14345.
- [174] Jeremy Reizenstein et al. “Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 10901–10911.
- [175] Konstantinos Rematas, Ricardo Martín-Brualla, and Vittorio Ferrari. “Sharf: Shape-conditioned radiance fields from a single view”. In: *International Conference on Machine Learning*. 2021.
- [176] Konstantinos Rematas et al. “Urban radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 12932–12942.
- [177] Zhongzheng Ren et al. “Neural volumetric object selection”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 6133–6142.
- [178] Mike Roberts et al. “Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 10912–10922.
- [179] Barbara Roessle et al. “Dense depth priors for neural radiance fields from sparse input views”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 12892–12901.
- [180] Luis Roldao, Raoul De Charette, and Anne Verroust-Blondet. “3D semantic scene completion: A survey”. In: *International Journal of Computer Vision* 130.8 (2022), pp. 1978–2005.
- [181] Robin Rombach et al. “High-resolution image synthesis with latent diffusion models”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022, pp. 10684–10695.
- [182] Antoni Rosinol, John J Leonard, and Luca Carlone. “Nerf-slam: Real-time dense monocular slam with neural radiance fields”. In: *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2023, pp. 3437–3444.
- [183] Radu Alexandru Rosu and Sven Behnke. “Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 8466–8475.
- [184] Denys Rozumnyi et al. “Tracking by 3D Model Estimation of Unknown Objects in Videos”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 14086–14096.
- [185] Mehdi SM Sajjadi et al. “Scene representation transformer: Geometry-free novel view synthesis through set-latent scene representations”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 6229–6238.
- [186] Johannes L Schonberger and Jan-Michael Frahm. “Structure-from-motion revisited”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 4104–4113.

- [187] Katja Schwarz et al. “Graf: Generative radiance fields for 3d-aware image synthesis”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 20154–20166.
- [188] Hoigi Seo et al. “DITTO-NeRF: Diffusion-based Iterative Text To Omni-directional 3D Model”. In: *arXiv preprint arXiv:2304.02827* (2023).
- [189] Seunghyeon Seo et al. “MixNeRF: Modeling a Ray with Mixture Density for Novel View Synthesis from Sparse Inputs”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 20659–20668.
- [190] Nur Muhammad Mahi Shafiullah et al. “CLIP-Fields: Weakly Supervised Semantic Fields for Robotic Memory”. In: *ICRA2023 Workshop on Pretraining for Robotics (PT4R)*. 2023.
- [191] Prafull Sharma et al. “Neural Groundplans: Persistent Neural Scene Representations from a Single Image”. In: *The Eleventh International Conference on Learning Representations*. 2022.
- [192] William Shen et al. “Distilled Feature Fields Enable Few-Shot Language-Guided Manipulation”. In: *Conference on Robot Learning*. PMLR. 2023, pp. 405–424.
- [193] Dongseok Shim, Seungjae Lee, and H Jin Kim. “SNeRL: Semantic-aware Neural Radiance Fields for Reinforcement Learning”. In: *International Conference on Machine Learning*. PMLR. 2023.
- [194] Kaede Shiohara and Toshihiko Yamasaki. “Detecting deepfakes with self-blended images”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 18720–18729.
- [195] Yawar Siddiqui et al. “Panoptic lifting for 3d scene understanding with neural fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 9043–9052.
- [196] Eugen Šlapak et al. “Neural radiance fields in the industrial and robotics domain: applications, research opportunities and use cases”. In: *arXiv preprint arXiv:2308.07118* (2023).
- [197] Cameron Omid Smith et al. “Unsupervised Discovery and Composition of Object Light Fields”. In: *Transactions on Machine Learning Research* (2023).
- [198] Nagabhushan Somraj, Adithyan Karanayil, and Rajiv Soundararajan. “SimpleNeRF: Regularizing Sparse Input Neural Radiance Fields with Simpler Solutions”. In: *SIGGRAPH Asia 2023 Conference Papers*. 2023, pp. 1–11.
- [199] Hyeonseop Song et al. “Blending-NeRF: Text-Driven Localized Editing in Neural Radiance Fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 14383–14393.
- [200] Jiuhn Song et al. “DàRF: Boosting Radiance Fields from Sparse Input Views with Monocular Depth Adaptation”. In: *Thirty-seventh Conference on Neural Information Processing Systems*. 2023.
- [201] Liang Song et al. “SC-NeRF: Self-Correcting Neural Radiance Field with Sparse Views”. In: *arXiv preprint arXiv:2309.05028* (2023).
- [202] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. “Sun rgb-d: A rgb-d scene understanding benchmark suite”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 567–576.
- [203] Trevor Standley et al. “Which tasks should be learned together in multi-task learning?” In: *International Conference on Machine Learning*. PMLR. 2020, pp. 9120–9132.
- [204] Karl Stelzner, Kristian Kersting, and Adam R Kosiorok. “Decomposing 3d scenes into objects via unsupervised volume segmentation”. In: *arXiv preprint arXiv:2104.01148* (2021).
- [205] Julian Straub et al. “The Replica dataset: A digital replica of indoor spaces”. In: *arXiv preprint arXiv:1906.05797* (2019).
- [206] Shih-Yang Su et al. “A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 12278–12291.
- [207] Cheng Sun, Min Sun, and Hwann-Tzong 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.
- [208] Jingxiang Sun et al. “Fenerf: Face editing in neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 7672–7682.
- [209] Pei Sun et al. “Scalability in perception for autonomous driving: Waymo open dataset”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020, pp. 2446–2454.
- [210] Tao Sun et al. “SHIFT: a synthetic driving dataset for continuous multi-task domain adaptation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 21371–21382.

- [211] Roman Suvorov et al. “Resolution-robust large mask inpainting with fourier convolutions”. In: *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 2022, pp. 2149–2159.
- [212] Christian Szegedy et al. “Rethinking the inception architecture for computer vision”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 2818–2826.
- [213] Matthew Tancik et al. “Block-nerf: Scalable large scene neural view synthesis”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 8248–8258.
- [214] Matthew Tancik et al. “Nerfstudio: A modular framework for neural radiance field development”. In: *ACM SIGGRAPH 2023 Conference Proceedings*. 2023, pp. 1–12.
- [215] Songlin Tang et al. “Scene-Generalizable Interactive Segmentation of Radiance Fields”. In: *Proceedings of the 31st ACM International Conference on Multimedia*. 2023, pp. 6744–6755.
- [216] Ayush Tewari et al. “Advances in neural rendering”. In: *Computer Graphics Forum*. Vol. 41. 2. Wiley Online Library. 2022, pp. 703–735.
- [217] Alex Trevithick and Bo Yang. “Grf: Learning a general radiance field for 3d representation and rendering”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 15182–15192.
- [218] Prune Truong et al. “SPARF: Neural Radiance Fields from Sparse and Noisy Poses”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.
- [219] Nikolaos Tsagkas, Oisín Mac Aodha, and Chris Xiaoquan Lu. “VL-Fields: Towards Language-Grounded Neural Implicit Spatial Representations”. In: *arXiv preprint arXiv:2305.12427* (2023).
- [220] Vadim Tschernezki, Diane Larlus, and Andrea Vedaldi. “NeuralDiff: Segmenting 3D objects that move in egocentric videos”. In: *2021 International Conference on 3D Vision (3DV)*. IEEE. 2021, pp. 910–919.
- [221] Vadim Tschernezki et al. “Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations”. In: *2022 International Conference on 3D Vision (3DV)*. IEEE. 2022, pp. 443–453.
- [222] Haithem Turki et al. “SUDS: Scalable Urban Dynamic Scenes”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 12375–12385.
- [223] Mikaela Angelina Uy et al. “SCADE: NeRFs from Space Carving with Ambiguity-Aware Depth Estimates”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 16518–16527.
- [224] Simon Vandenhende et al. “Multi-task learning for dense prediction tasks: A survey”. In: *IEEE transactions on pattern analysis and machine intelligence* 44.7 (2021), pp. 3614–3633.
- [225] Dor Verbin et al. “Ref-nerf: Structured view-dependent appearance for neural radiance fields”. In: *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE. 2022, pp. 5481–5490.
- [226] Suhani Vora et al. “NeSF: Neural Semantic Fields for Generalizable Semantic Segmentation of 3D Scenes”. In: *Transactions on Machine Learning Research* (2022). ISSN: 2835-8856.
- [227] Can Wang et al. “Clip-nerf: Text-and-image driven manipulation of neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 3835–3844.
- [228] Can Wang et al. “Nerf-art: Text-driven neural radiance fields stylization”. In: *IEEE Transactions on Visualization and Computer Graphics* (2023).
- [229] Dongqing Wang et al. “InpaintNeRF360: Text-Guided 3D Inpainting on Unbounded Neural Radiance Fields”. In: *arXiv preprint arXiv:2305.15094* (2023).
- [230] Fusang Wang et al. “PlaNeRF: SVD Unsupervised 3D Plane Regularization for NeRF Large-Scale Scene Reconstruction”. In: (2024).
- [231] Guangcong Wang et al. “Sparsenerf: Distilling depth ranking for few-shot novel view synthesis”. In: *arXiv preprint arXiv:2303.16196* (2023).
- [232] Yuxin Wang, Wayne Wu, and Dan Xu. “Learning unified decompositional and compositional nerf for editable novel view synthesis”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 18247–18256.
- [233] Zhou Wang et al. “Image quality assessment: from error visibility to structural similarity”. In: *IEEE transactions on image processing* 13.4 (2004), pp. 600–612.

- [234] Zirui Wang et al. “NeRF-: Neural radiance fields without known camera parameters”. In: *arXiv preprint arXiv:2102.07064* (2021).
- [235] Silvan Weder et al. “Removing objects from neural radiance fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 16528–16538.
- [236] Xiaobao Wei et al. “NOC: High-Quality Neural Object Cloning with 3D Lifting of Segment Anything”. In: *arXiv preprint arXiv:2309.12790* (2023).
- [237] Chung-Yi Weng et al. “Humannerf: Free-viewpoint rendering of moving people from monocular video”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern Recognition*. 2022, pp. 16210–16220.
- [238] Felix Wimbauer et al. “Behind the Scenes: Density Fields for Single View Reconstruction”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 9076–9086.
- [239] Yu-Shiang Wong and Niloy J Mitra. “Factored Neural Representation for Scene Understanding”. In: *Computer Graphics Forum*. Wiley Online Library. 2023, e14911.
- [240] Tianhao Walter Wu et al. “D²NeRF: Self-Supervised Decoupling of Dynamic and Static Objects from a Monocular Video”. In: *Advances in Neural Information Processing Systems*. 2022.
- [241] Tong Wu et al. “Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 803–814.
- [242] Zirui Wu et al. “Mars: An instance-aware, modular and realistic simulator for autonomous driving”. In: *arXiv preprint arXiv:2307.15058* (2023).
- [243] Jamie Wynn and Daniyar Turmukhambetov. “Diffusionerf: Regularizing neural radiance fields with denoising diffusion models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 4180–4189.
- [244] Weihao Xia and Jing-Hao Xue. “A Survey on Deep Generative 3D-aware Image Synthesis”. In: *ACM Computing Surveys* 56.4 (2023), pp. 1–34.
- [245] Jianfeng Xiang et al. “Gram-hd: 3d-consistent image generation at high resolution with generative radiance manifolds”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 2195–2205.
- [246] Yuanbo Xiangli et al. “Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering”. In: *European conference on computer vision*. Springer. 2022, pp. 106–122.
- [247] Yuanbo Xiangli et al. “AssetField: Assets Mining and Reconfiguration in Ground Feature Plane Representation”. In: *arXiv preprint arXiv:2303.13953* (2023).
- [248] Christopher Xie et al. “Fig-nerf: Figure-ground neural radiance fields for 3d object category modelling”. In: *2021 International Conference on 3D Vision (3DV)*. IEEE. 2021, pp. 962–971.
- [249] Yiheng Xie et al. “Neural fields in visual computing and beyond”. In: *Computer Graphics Forum*. Vol. 41. 2. Wiley Online Library. 2022, pp. 641–676.
- [250] Chenfeng Xu et al. “NeRF-Det: Learning Geometry-Aware Volumetric Representation for Multi-View 3D Object Detection”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 23320–23330.
- [251] Dejie Xu et al. “Sinnerf: Training neural radiance fields on complex scenes from a single image”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 736–753.
- [252] Junkai Xu et al. “MonoNeRD: NeRF-like Representations for Monocular 3D Object Detection”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 6814–6824.
- [253] Xiaomeng Xu et al. “JacobiNeRF: NeRF Shaping with Mutual Information Gradients”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 16498–16507.
- [254] Yinghao Xu et al. “DisCoScene: Spatially Disentangled Generative Radiance Fields for Controllable 3D-aware Scene Synthesis”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 4402–4412.
- [255] Yang Xue et al. “Giraffe hd: A high-resolution 3d-aware generative model”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 18440–18449.
- [256] Karmesh Yadav et al. “Habitat-matterport 3d semantics dataset”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 4927–4936.

- [257] Bangbang Yang et al. “Learning object-compositional neural radiance field for editable scene rendering”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 13779–13788.
- [258] Bangbang Yang et al. “Neural rendering in a room: amodal 3d understanding and free-viewpoint rendering for the closed scene composed of pre-captured objects”. In: *ACM Transactions on Graphics (TOG)* 41.4 (2022), pp. 1–10.
- [259] Jiawei Yang, Marco Pavone, and Yue Wang. “FreeNeRF: Improving Few-shot Neural Rendering with Free Frequency Regularization”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.
- [260] Jiawei Yang et al. “EmerNeRF: Emergent Spatial-Temporal Scene Decomposition via Self-Supervision”. In: *arXiv preprint arXiv:2311.02077* (2023).
- [261] Yuanbo Yang et al. “UrbanGIRAFFE: Representing Urban Scenes as Compositional Generative Neural Feature Fields”. In: *arXiv preprint arXiv:2303.14167* (2023).
- [262] Botao Ye et al. “Self-Supervised Super-Plane for Neural 3D Reconstruction”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 21415–21424.
- [263] Jianglong Ye, Naiyan Wang, and Xiaolong Wang. “FeatureNeRF: Learning Generalizable NeRFs by Distilling Pre-trained Vision Foundation Models”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023.
- [264] Weicai Ye et al. “Intrinsicnerf: Learning intrinsic neural radiance fields for editable novel view synthesis”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 339–351.
- [265] Lin Yen-Chen et al. “inerf: Inverting neural radiance fields for pose estimation”. In: *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2021, pp. 1323–1330.
- [266] Chandan Yeshwanth et al. “ScanNet++: A High-Fidelity Dataset of 3D Indoor Scenes”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 12–22.
- [267] Yu Yin et al. “NeRFInvertor: High Fidelity NeRF-GAN Inversion for Single-shot Real Image Animation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.
- [268] Alex Yu et al. “pixelnerf: Neural radiance fields from one or few images”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 4578–4587.
- [269] Alex Yu et al. “Plenotrees for real-time rendering of neural radiance fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 5752–5761.
- [270] Hong-Xing Yu, Leonidas J. Guibas, and Jiajun Wu. “Unsupervised Discovery of Object Radiance Fields”. In: *International Conference on Learning Representations*. 2022.
- [271] Wentao Yuan et al. “Star: Self-supervised tracking and reconstruction of rigid objects in motion with neural rendering”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021, pp. 13144–13152.
- [272] Amir R Zamir et al. “Taskonomy: Disentangling task transfer learning”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 3712–3722.
- [273] Jesus Zarzar et al. “SegNeRF: 3D Part Segmentation with Neural Radiance Fields”. In: *arXiv preprint arXiv:2211.11215* (2022).
- [274] Yanjie Ze et al. “Gnfactor: Multi-task real robot learning with generalizable neural feature fields”. In: *Conference on Robot Learning*. PMLR. 2023, pp. 284–301.
- [275] Jiakai Zhang et al. “Editable free-viewpoint video using a layered neural representation”. In: *ACM Transactions on Graphics (TOG)* 40.4 (2021), pp. 1–18.
- [276] Jichao Zhang et al. “3D-aware semantic-guided generative model for human synthesis”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 339–356.
- [277] Jingbo Zhang et al. “Fdnerf: Few-shot dynamic neural radiance fields for face reconstruction and expression editing”. In: *SIGGRAPH Asia 2022 Conference Papers*. 2022, pp. 1–9.
- [278] Kai Zhang et al. “Nerf++: Analyzing and improving neural radiance fields”. In: *arXiv preprint arXiv:2010.07492* (2020).
- [279] Mingtong Zhang et al. “Beyond RGB: Scene-Property Synthesis with Neural Radiance Fields”. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, pp. 795–805.

- [280] Richard Zhang et al. “The unreasonable effectiveness of deep features as a perceptual metric”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 586–595.
- [281] Xiaoshuai Zhang et al. “Nerflets: Local radiance fields for efficient structure-aware 3d scene representation from 2d supervision”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 8274–8284.
- [282] Fuqiang Zhao et al. “Humannerf: Efficiently generated human radiance field from sparse inputs”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 7743–7753.
- [283] Shuhong Zheng et al. “Multi-task View Synthesis with Neural Radiance Fields”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 21538–21549.
- [284] Shuaifeng Zhi et al. “ilabel: Interactive neural scene labelling”. In: *arXiv preprint arXiv:2111.14637* (2021).
- [285] Shuaifeng Zhi et al. “In-place scene labelling and understanding with implicit scene representation”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 15838–15847.
- [286] Peng Zhou et al. “CIPS-3D++: End-to-End Real-Time High-Resolution 3D-Aware GANs for GAN Inversion and Stylization”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [287] Chengxuan Zhu et al. “Occlusion-Free Scene Recovery via Neural Radiance Fields”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 20722–20731.
- [288] Fang Zhu et al. “Deep Review and Analysis of Recent NeRFs”. In: *APSIPA Transactions on Signal and Information Processing* 12.1 (2023).
- [289] Yiyu Zhuang et al. “Mofanerf: Morphable facial neural radiance field”. In: *European Conference on Computer Vision*. Springer. 2022, pp. 268–285.