

# LERF: Language Embedded Radiance Fields

Justin Kerr\*, Chung Min Kim\*, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik

UC Berkeley

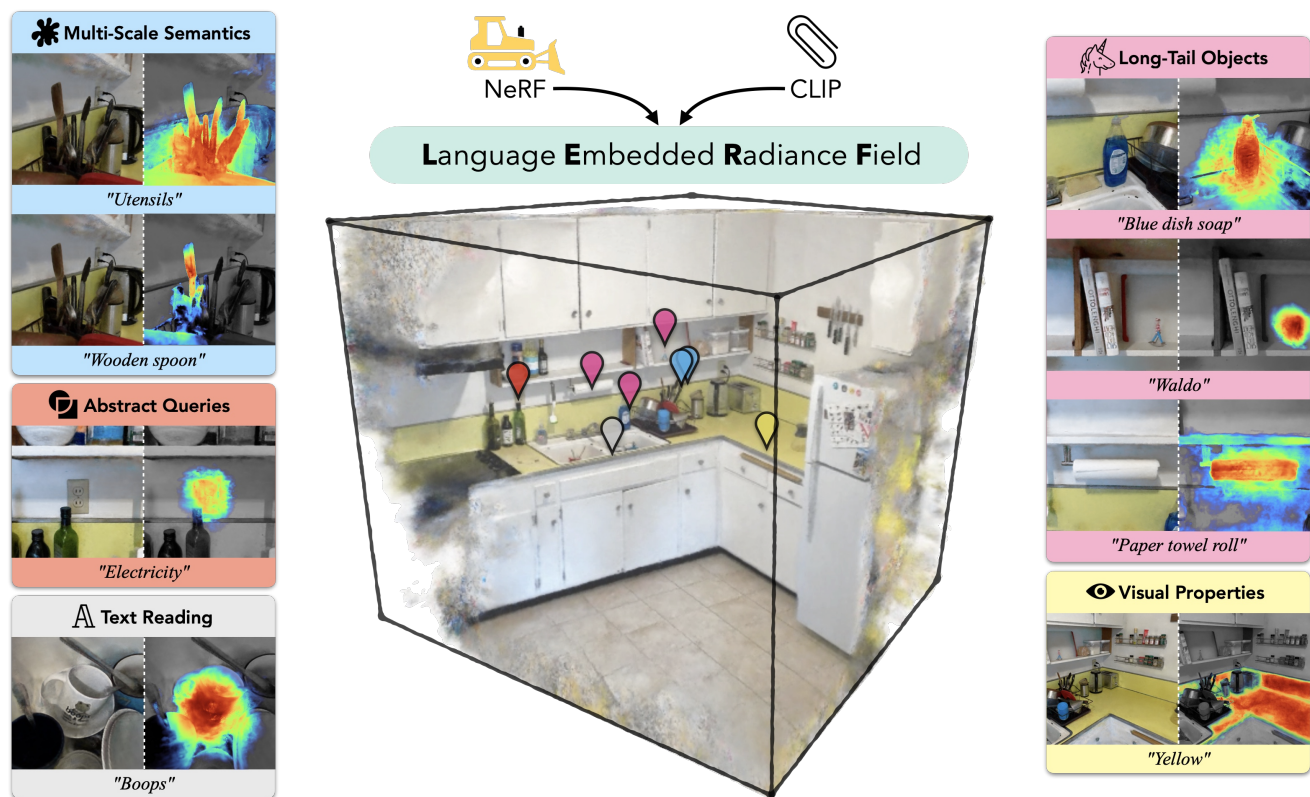


Figure 1: **Language Embedded Radiance Fields (LERF)**. LERF grounds CLIP representations in a dense, multi-scale 3D field. A LERF can be reconstructed from a hand-held phone capture within 45 minutes, then can render dense relevancy maps given textual queries interactively in real-time. LERF enables a broad range of concepts to be queried via natural language, from abstract queries like “Electricity”, visual properties like “Yellow”, long-tail objects such as “Waldo”, and even reading text like “Boops” on the mug. For each prompt, an RGB image and relevancy map are rendered focusing on the location with maximum relevancy activation.

## Abstract

Humans describe the physical world using natural language to refer to specific 3D locations based on a vast range of properties: visual appearance, semantics, abstract associations, or actionable affordances. In this work we propose Language Embedded Radiance Fields (LERFs), a method for grounding language embeddings from off-the-shelf models like CLIP into NeRF, which enable these types of open-

ended language queries in 3D. LERF learns a dense, multi-scale language field inside NeRF by volume rendering CLIP embeddings along training rays, supervising these embeddings across training views to provide multi-view consistency and smooth the underlying language field. After optimization, LERF can extract 3D relevancy maps for a broad range of language prompts interactively in real-time, which has potential use cases in robotics, understanding vision-language models, and interacting with 3D scenes. LERF enables pixel-aligned, zero-shot queries on the distilled 3D

\*Equal contribution, corresponding authors.

*CLIP embeddings without relying on region proposals or masks, supporting long-tail open-vocabulary queries hierarchically across the volume. See the project website at <https://lerf.io>.*

## 1. Introduction

Neural Radiance Fields (NeRFs) [25] have emerged as a powerful technique for capturing photorealistic digital representations of intricate real-world 3D scenes. However, the immediate output of NeRFs is nothing but a colorful density field, devoid of meaning or context, which inhibits building interfaces for interacting with the resulting 3D scenes.

Natural language is an intuitive interface for interacting with a 3D scene. Consider the capture of a kitchen in Figure 1. Imagine being able to navigate this kitchen by asking where the “*utensils*” are, or more specifically for a tool that you could use for “*stirring*”, and even for your favorite mug with a specific logo on it — all through the comfort and familiarity of everyday conversation. This requires not only the capacity to handle natural language input queries but also the ability to incorporate semantics at multiple scales and relate to long-tail and abstract concepts.

In this work, we propose Language Embedded Radiance Fields (LERF), a novel approach that grounds language within NeRF by optimizing embeddings from an off-the-shelf vision-language model like CLIP into 3D scenes. Notably, LERF utilizes CLIP directly without the need for finetuning through datasets like COCO or reliance on mask region proposals, which limits the ability to capture a wide range of semantics. Because LERF preserves the integrity of CLIP embeddings at multiple scales, it is able to handle a broad range of language queries, including visual properties (“*yellow*”), abstract concepts (“*electricity*”), text (“*boops*”), and long-tail objects (“*waldo*”) as illustrated in Figure 1.

We construct a LERF by optimizing a language field jointly with NeRF, which takes both position and physical scale as input and outputs a single CLIP vector. During training, the field is supervised using a multi-scale feature pyramid that contains CLIP embeddings generated from image crops of training views. This allows the CLIP encoder to capture different scales of image context, thus associating the same 3D location with distinct language embeddings at different scales (*e.g.* “*utensils*” vs. “*wooden spoon*”). The language field can be queried at arbitrary scales during test time to obtain 3D relevancy maps. To regularize the optimized language field, self-supervised DINO [5] features are also incorporated through a shared bottleneck.

LERF offers an added benefit: since we extract CLIP embeddings from multiple views over multiple scales, the relevancy maps of text queries obtained through our 3D CLIP embedding are more localized compared to those ob-

tained via 2D CLIP embeddings. By definition, they are also 3D consistent, enabling queries directly in the 3D fields without having to render to multiple views.

LERF can be trained without significantly slowing down the base NeRF implementation. Upon completion of the training process, LERF allows for the generation of 3D relevancy maps for a wide range of language prompts in real-time. We evaluate the capabilities of LERF on a set of hand-held captured in-the-wild scenes and find it can localize both fine-grained queries relating to highly specific parts of geometry (“*fingers*”), or abstract queries relating to multiple objects (“*cartoon*”). LERF produces view-consistent relevancy maps in 3D across a wide range of queries and scenes, which are best viewed in videos on our website. We also provide quantitative evaluations against popular open-vocab detectors LSeg [22] and OWL-ViT [26], by distilling LSeg features into 3D [21] and querying OWL-ViT from rendered novel views. Our results suggest that features in 3D from LERF can localize a wide variety of queries across in-the-wild scenes. The zero-shot capabilities of LERF leads to potential use cases in robotics, analyzing vision-language models, and interacting with 3D scenes. Code and data is available at <https://lerf.io>.

## 2. Related Work

**Open-Vocabulary Object Detection** A number of approaches study detecting objects in 2D images given natural language prompts. These lie on a spectrum from purely zero-shot to fully trained on segmentation datasets. LSeg [22] trains a 2D image encoder on labeled segmentation datasets, which outputs pixel-wise embeddings that best match the CLIP text embedding of the segmentation label at that given pixel. CRIS [40] and CLIPSeg [24] train a 2D image decoder to output a relevancy map based on the query CLIP embedding and the intermediate outputs of the pretrained CLIP image encoder. However, such fine-tuning approaches tend to lose significant language capabilities by training on a smaller dataset.

Another common approach for 2D images is a two-stage framework wherein class-agnostic region or mask proposals direct where to query open-vocabulary classification models. OpenSeg [15] simultaneously learns a mask prediction model while predicting the text embeddings for each mask, while ViLD [17] directly uses CLIP [30] to classify 2D regions from class-agnostic mask proposal networks. Detic [42] builds on existing two-stage object detector approaches, but demonstrates a greater generalization ability by allowing detector classifiers to train with image classification data. OWL-ViT [26] attaches lightweight object classification and localization heads after a pre-trained 2D image encoder. Although region proposal-based approaches can leverage more detection data, these proposal generators still tend to output within the training set distribution. Con-

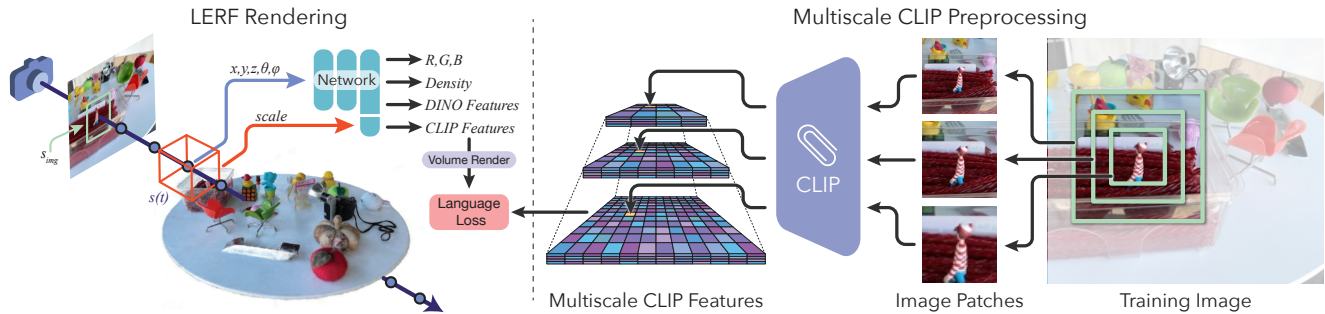


Figure 2: **LERF Optimization:** *Left:* LERF represents a field of 3D volumes, parameterized by position  $x, y, z$  and scale  $s$  (orange cube). To render a CLIP embedding along a ray, the field is sampled and averaged according to NeRF’s volume rendering weights. Physical scale corresponds to an image scale  $s_{\text{img}}$  via projective geometry. *Right:* We pre-compute a multi-scale feature pyramid of CLIP embeddings over training views, and during training interpolate this pyramid with  $s_{\text{img}}$  and the ray’s pixel location to obtain CLIP supervision. The CLIP loss maximizes cosine similarity, and other outputs are supervised with mean squared-error using standard per-pixel rendering.

sequently, as noted by the authors of OV-Seg [23], such networks often face difficulties in accurately segmenting unlabeled hierarchical components of the original masks, such as object parts. LERF strives to avoid region proposals by incorporating language embeddings in a dense, 3D, multi-scale field which allows hierarchical text queries.

Grad-CAM [33], attention-based methods [7], or patch-aligned contrastive learning [27] provide a relevancy mapping between 2D images and text in vision-language models (e.g., CLIP). Works such as Semantic Abstraction [18] have shown that these frameworks can be used to detect long-tail objects for scene understanding. Outputs from LERF are most similar in spirit to these methods, outputting a 3D relevancy score given a query. However, LERF builds a 3D representation that can be queried with different text prompts without reconstructing the underlying representation each time, and in addition fuses multiple viewpoints into a single shared scene representation, rather than operating per-image.

**Distilling 2D Features into NeRF** NeRF has an attractive property of averaging information across multiple views. Several prior works leverage this to improve the quality of 2D semantics, segmentations, or feature vectors by distilling them into 3D. Semantic NeRF [41], Panoptic Lifting [35], and NeRF-SOS [14] embed semantic information from semantic segmentation networks into 3D, showing that combining noisy or sparse labels in 3D can result in crisp 3D segmentations. This concept has been applied to segmenting objects with extremely sparse user input scribbles of foreground-background masks [32]. Our approach draws inspiration from these works by averaging multiple potentially noisy language embeddings over input views.

Distilled Feature Fields [21] and Neural Feature Fusion Fields [39] explore embedding pixel-aligned feature vectors from LSeg or DINO [5] into a NeRF, and show they can be used for 3D manipulations of the underlying geometry.

LERF similarly embeds feature vectors into NeRF, but also demonstrates a method to distill non pixel-aligned embeddings (e.g., from CLIP) into 3D without fine-tuning.

**3D Language Grounding** Incorporating language into 3D has been explored in a wide range of contexts: 3D visual question answering [16, 2, 6] leverage 3D information to extract answers to queries about the environment. In addition, incorporating language with shape information can improve object recognition via text [11, 38].

LERF is more similar to 3D scene representations in robotics which fuse vision-language embeddings to support natural language interaction. VL-Maps [19] and OpenScene [29] build a 3D volume of language features which can be queried for navigation tasks, by using pre-trained pixel-aligned language encoders [15, 22]. In LERF, we compare against one such encoder, LSeg, in 3D and find it loses significant expressive capability compared to CLIP.

CLIP-Fields [34] and NLMaps-SayCan [8] fuse CLIP embeddings of crops into pointclouds, using a contrastively supervised field and classical pointcloud fusion respectively. In CLIP-Fields, the crop locations are guided by Detic [42]. On the other hand, NLMaps-SayCan relies on region proposal networks. These maps are sparser than LERF as they primarily query CLIP on detected objects rather than densely throughout views of the scene. Concurrent work ConceptFusion [20] fuses CLIP features more densely in RGBD pointclouds, using Mask2Former [9] to predict regions of interest, meaning it can lose objects which are out of distribution to Mask2Former’s training set. In contrast, LERF does not use region or mask proposals.

LERF contributes a new dense, volumetric interface for 3D text queries which can integrate with a broad range of downstream applications of 3D language, improving the resolution and fidelity at which these methods can query the environment when multi-view inputs are available. This is enabled by the smoothing behavior of embedding po-

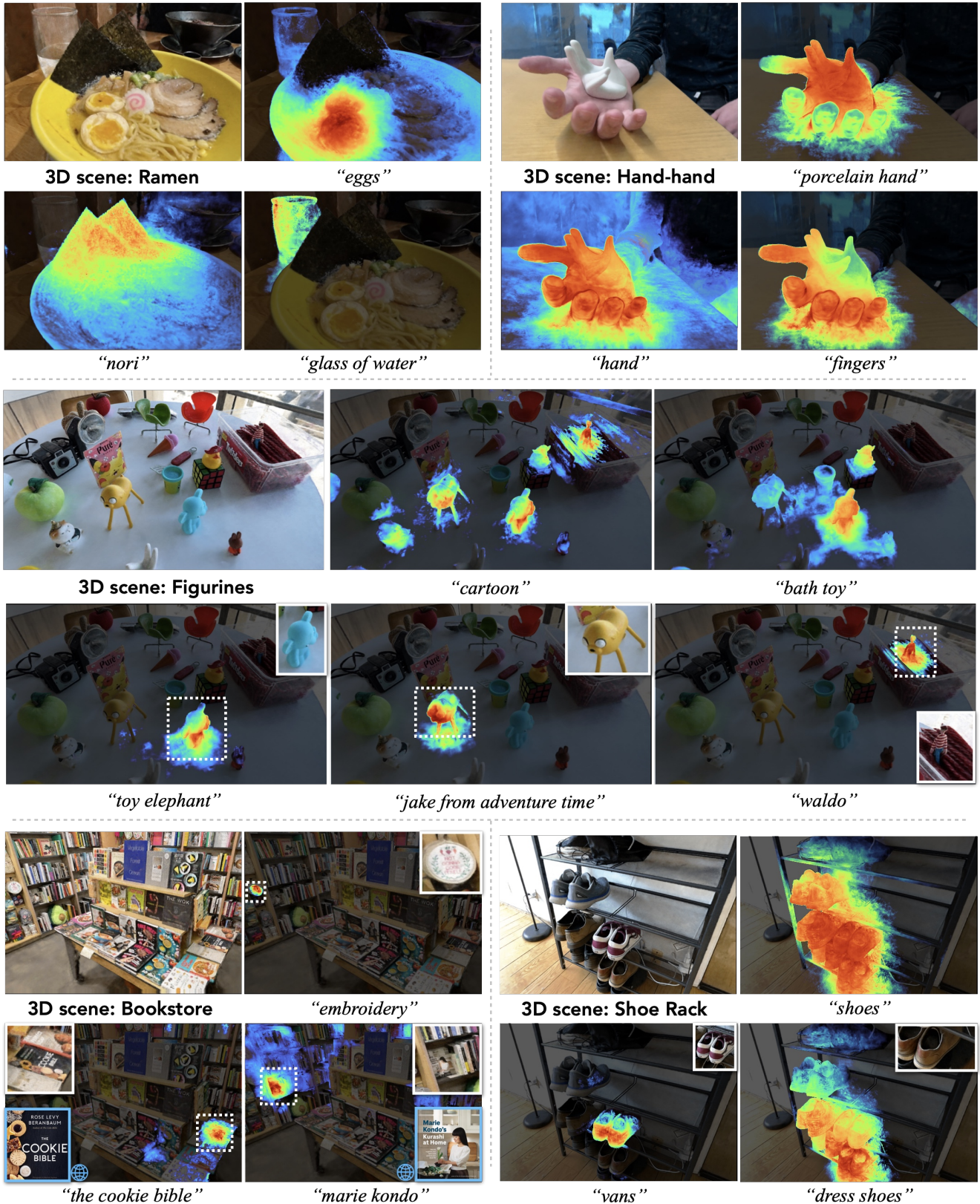


Figure 3: **Results with LERF for 5 in-the-wild scenes.** Each image shows a visual rendering of the LERF (Sec. 3), along with relevancy renderings (Sec. 3.5) for each text query and a cropped view of the activated region. For the bookstore scene, the original book cover images are shown in blue with a globe icon. See Sec. 4.1 for discussion and details on relevancy visualization.



Figure 4: **Multi-scale Semantics**: Each 3D crop visualizes the  $\text{argmax}$  position and scale of each query. Queries for smaller objects correspond to smaller crops, since those scales produce higher relevancy activations. This behavior highlights the importance of scale inside LERF’s parameterization (Sec. 3). Note that for “Table”, the crop is centered on the edge of the table and only includes a section of it, due to the phenomena discussed in Fig. 11.

tentially noisy feature vectors from multiple views into the dense LERF structure.

### 3. Language Embedded Radiance Fields

Given a set of calibrated input images, we ground CLIP embeddings into a 3D field within NeRF. However, querying a CLIP embedding for a single 3D point is ambiguous, as CLIP is inherently a global image embedding and not conducive to pixel-aligned feature extraction. To account for this property, we propose a novel approach that involves learning a field of language embeddings over *volumes* centered at the sample point (Fig. 2). Specifically, the output of this field is the average CLIP embedding across all training views of image crops containing the specified volume. By reframing the query from points to volumes, we can effectively supervise a dense field from coarse crops of input images, which can be rendered in a pixel-aligned manner by conditioning on a given volume scale.

#### 3.1. LERF Volumetric Rendering

NeRF takes in a position  $\vec{x}$  and view direction  $\vec{d}$  and outputs color  $\vec{c}$  and density  $\sigma$ . Samples of these values can be composited along a ray to produce a pixel’s color. To create LERF, we augment NeRF’s outputs with a language embedding  $F_{\text{lang}}(\vec{x}, s) \in \mathbb{R}^d$ , which takes an input position  $\vec{x}$  and physical scale  $s$ , and outputs a  $d$ -dimensional language embedding. We choose this output to be *view-independent*, since the semantic meaning of a location should be invariant to viewing angle. This allows multiple views to contribute to the same field input, averaging their embeddings together. The scale  $s$  represents the side length in world coordinates of a cube centered at  $\vec{x}$ , and is analogous to how Mip-NeRF[3, 4] incorporates different scales via integrated positional encodings.

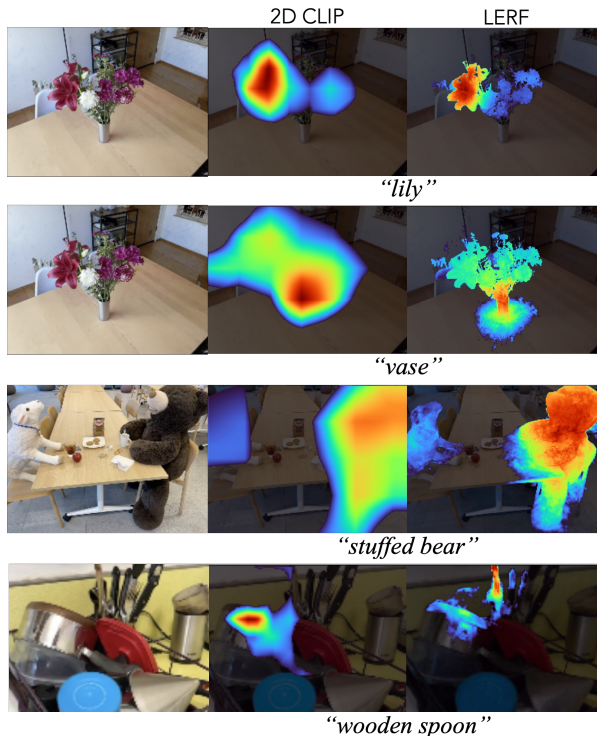


Figure 5: **2D CLIP vs LERF**: The left visualizes similarity interpolated over patchwise CLIP embeddings, and the right rendered from LERF. Because volumetric language rendering incorporates information from multiple views, 3D relevancy activation maps have better alignment with the underlying scene geometry.

Rendering color and density from LERF remains exactly the same as NeRF. To render language embeddings into an image, we adopt a similar technique as prior work[21, 41] to render language embeddings along a ray  $\vec{r}(t) = \vec{o}_t + t\vec{d}$ . However, since LERF is a field over *volumes*, not points, we must also define a scale parameter for each position along the ray. We achieve this by fixing an initial scale in the image plane  $s_{\text{img}}$  and define  $s(t)$  to increase proportionally with focal length and sample distance from the ray origin:  $s(t) = s_{\text{img}} * t / f_{xy}$  (Fig. 2, left). Geometrically, this represents a frustrum along the ray. We calculate rendering weights as in NeRF:  $T(t) = \int_t \exp(-\sigma(s)ds)$ ,  $w(t) = \int_t T(t)\sigma(t)dt$ , then integrate the LERF to obtain raw language outputs:  $\hat{\phi}_{\text{lang}} = \int_t w(t)F_{\text{lang}}(r(t), s(t)) dt$ , and finally normalize each embedding to the unit sphere as in CLIP:  $\phi_{\text{lang}} = \hat{\phi}_{\text{lang}} / \|\hat{\phi}_{\text{lang}}\|$ . We find that spherical interpolation (slerp) between samples on a ray is unnecessary because non-zero weight samples along a ray tend to be spatially close, and instead opt for a weighted euclidean average followed by normalization for implementation simplicity.

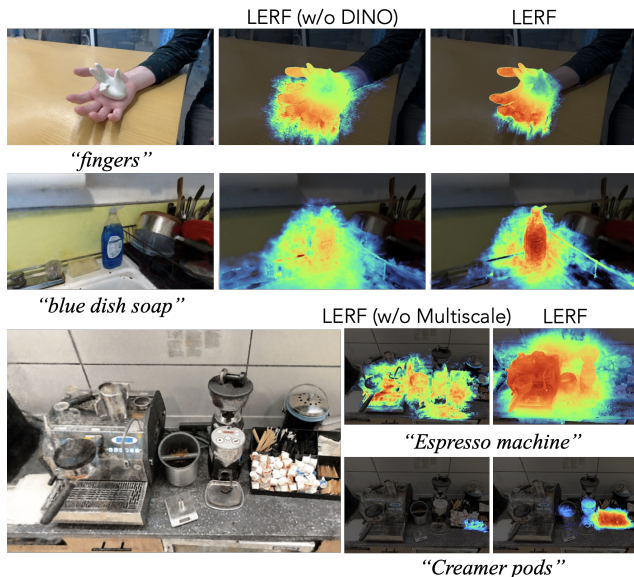


Figure 6: **Ablations:** We ablate DINO regularization and multi-scale training (Sec. 4.4), and highlight qualitative degradation in relevancy maps here.

### 3.2. Multi-Scale Supervision

To supervise language field outputs  $F_{\text{lang}}$ , recall that we can only query language embeddings over image patches, not pixels. Therefore, to supervise the multi-scale LERF, we supervise each rendered frustrum with an image crop of size  $s_{\text{img}}$  centered at the image pixel where the ray originated. In practice, re-computing a CLIP embedding for each ray during LERF optimization would be prohibitively expensive, so instead we pre-compute an image pyramid over multiple image crop scales and store the CLIP embeddings of each crop (Fig 2, right). This pyramid has  $n$  layers sampled between  $s_{\text{min}}$  and  $s_{\text{max}}$ , with each crop arranged in a grid with 50% overlap between crops.

During training, we randomly sample ray origins uniformly throughout input views, and uniformly randomly select  $s_{\text{img}} \in (s_{\text{min}}, s_{\text{max}})$  for each. Since these samples don't necessarily fall in the center of a crop in this image pyramid, we perform trilinear interpolation between the embeddings from the 4 nearest crops for the scale above and below to produce the final ground truth embedding  $\phi_{\text{lang}}^{\text{gt}}$ . We minimize a loss between rendered and ground truth embeddings which maximizes cosine similarity between the two, scaling by a constant  $\lambda_{\text{lang}}$  (Sec. 3.3):  $L_{\text{lang}} = -\lambda_{\text{lang}} \phi_{\text{lang}} \cdot \phi_{\text{lang}}^{\text{gt}}$ .

### 3.3. DINO Regularization

Naïvely implementing LERF as described produces cohesive results, but the resulting relevancy maps can sometimes be patchy and contain outliers in regions with few views or little foreground-background separation (Fig. 6).

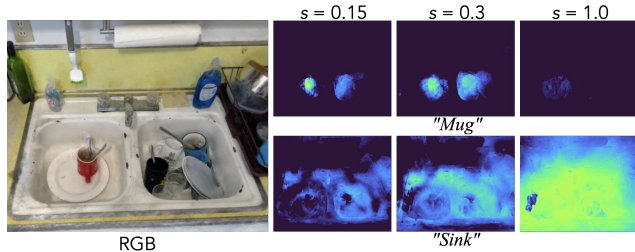


Figure 7: **Effect of scale:** We visualize the same queries at varying scale to highlight the effect of different context to the LERF activations. Physically smaller objects correspond to higher relevancy activations at smaller scale.

To mitigate this, we additionally train a field  $F_{\text{dino}}(\vec{x})$  which outputs a DINO [5] feature at each point. DINO has been shown to exhibit emergent object decomposition properties despite training on no labels [1], and additionally distills well into 3D fields [21], making it a good candidate for grouping language in 3D without relying on labeled data or imparting too strict a prior. Because DINO outputs pixel-aligned features,  $F_{\text{dino}}$  does not take in scale as an input, and is directly supervised for each ray with the DINO feature it corresponds to. We render  $\phi_{\text{dino}}$  identically to  $\phi_{\text{lang}}$  except without normalizing to a unit sphere, and supervise it with MSE loss on ground-truth DINO features. DINO is not used explicitly during inference, and only serves as an extra regularizer during training since CLIP and DINO output heads share an architectural backbone.

### 3.4. Field Architecture

Intuitively, optimizing a language embedding in 3D should not influence the distribution of density in the underlying scene representation. We capture this inductive bias in LERF by training two separate networks: one for feature vectors (DINO, CLIP), and the other for standard NeRF outputs (color, density). Gradients from  $L_{\text{lang}}$  and  $L_{\text{dino}}$  *do not affect* the NeRF outputs, and can be viewed as jointly optimizing a language field in conjunction with a radiance field.

We represent both fields with a multi-resolution hashgrid [28]. The language hashgrid has two output MLPs for CLIP and DINO respectively. Scale  $s$  is passed into the CLIP MLP as an extra input in addition to the concatenated hash-grid features. We adopt the Nerfacto method from Nerfstudio [37] as the backbone for our approach, leveraging the same proposal sampling, scene contraction, and appearance embeddings.

### 3.5. Querying LERF

Often, language models like CLIP are evaluated on zero-shot classification, where a category is selected from a list guaranteed to include the correct category [30]. However, in practical usage of LERF on in-the-wild scenes, an exhaus-

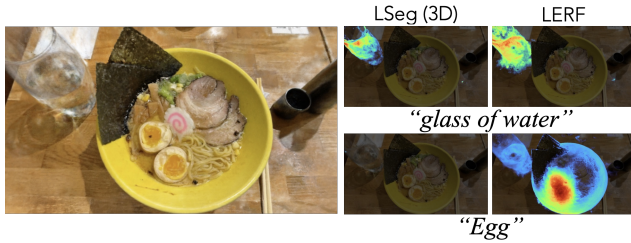


Figure 8: **Comparison to LSeg in 3D:** LSeg performs well on “glass of water” since cups are in the COCO dataset, but cannot locate an out-of-distribution query like an egg.

tive list of categories is not available. We view the open-endedness and ambiguity of natural language as a benefit, and propose a method to query 3D relevancy maps from the LERF given an arbitrary text query. Querying LERF has two parts: 1) obtaining a relevancy score for a rendered embedding and 2) automatically selecting a scale  $s$  given the prompt. **Relevancy Score:** To assign each rendered language embedding  $\phi_{\text{lang}}$  a score, we compute the CLIP embedding of the text query  $\phi_{\text{quer}}$ , along with a set of canonical phrases  $\phi_{\text{canon}}^i$ . We compute cosine similarity between the rendered embedding and canonical phrase embeddings, then compute the pairwise softmax between the rendered embedding text prompts. The relevancy score is then  $\min_i \frac{\exp(\phi_{\text{lang}} \cdot \phi_{\text{quer}})}{\exp(\phi_{\text{lang}} \cdot \phi_{\text{canon}}^i) + \exp(\phi_{\text{lang}} \cdot \phi_{\text{quer}})}$ . Intuitively, this score represents how much closer the rendered embedding is towards the query embedding compared to the canonical embeddings. All renderings use the same canonical phrases: “object”, “things”, “stuff”, and “texture”. We chose these as qualitatively “average” words for queries users might make, and found them to be surprisingly robust to queries ranging from incredibly specific to visual or abstract. We acknowledge that choosing these phrases is susceptible to prompt-engineering, and think fine-tuning them could be an interesting future work, perhaps incorporating feedback from user interaction about negative prompts they *do not* consider relevant.

**Scale Selection:** For each query, we compute a scale  $s$  to evaluate  $F_{\text{lang}}$ . To accomplish this, we generate relevancy maps across a range of scales 0 to 2 meters with 30 increments, and select the scale that yields the highest relevancy score. This scale is used for all pixels in the output relevancy map. We found this heuristic to be robust across a broad range of queries and is used for all the images and videos rendered in this paper. This assumes relevant parts of a scene are the same scale, see Limitations Sec. 5.

**Visibility Filtering:** Regions of the scene that lack sufficient views, such as those in the background or near floaters, may generate noisy embeddings. To address this issue, during querying we discard samples that were observed by fewer than five training views (approximately 5% of the

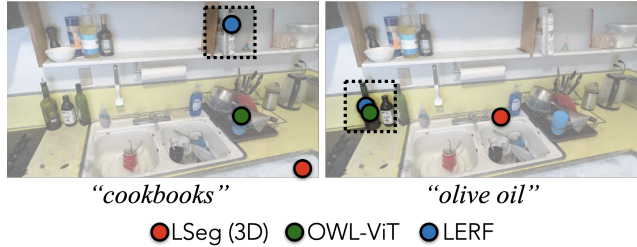


Figure 9: **Localization comparison** Qualitative comparison on localizing long-tail objects from novel views with LSeg in 3D (DFE) and OWL-ViT (Tab. 1)

views in our datasets).

### 3.6. Implementation Details

We implement LERF in Nerfstudio [37], on top of the Nerfacto method. Proposal sampling is the same except we reduce the number of LERF samples from 48 to 24 to increase training speed. We use the OpenClip [10] ViT-B/16 model trained on the LAION-2B dataset, with an image pyramid varying from  $s_{\text{min}} = .05$  to  $s_{\text{min}} = .5$  in 7 steps. The hashgrid used for representing language features is much larger than a typical RGB hashgrid: it has 32 layers from a resolution of 16 to 512, with a hash table size of  $2^{21}$  and feature dimension of 8. The CLIP MLP used for  $F_{\text{lang}}$  has 3 hidden layers with width 256 before the final 512 dimension CLIP output. The DINO MLP for  $F_{\text{DINO}}$  has 1 hidden layer of dimension 256.

We use the Adam optimizer for proposal networks and fields with weight decay  $10^{-9}$ , with an exponential learning rate scheduler from  $10^{-2}$  to  $10^{-3}$  over the first 5000 training steps. All models are trained to 30,000 steps (45 minutes), although good results can be obtained in as few as 6,000(8 minutes) as presented in the Appendix. We train on an NVIDIA A100, which takes roughly 20GB of memory total. One can interactively query in real-time within the Nerfstudio viewer. The  $\lambda$  used in weighting CLIP loss is 0.01, chosen empirically and ablated in Sec 4.4. When computing relevancy score, we multiply similarity by 10 as a temperature parameter within the softmax.

## 4. Experiments

We examine the capabilities of LERF and find that it can effectively process a wide variety of input text queries, encompassing various aspects of natural language specifications that current open-vocab detection frameworks encounter difficulty with. Though existing 3D scan datasets exist, they tend to be either of singulated objects [31, 13], or are RGB-D scans without enough views to optimize high quality NeRFs [12], and such simulated or scanned scenes contain few long-tail objects [36]. Emphasizing the capability of LERF to handle real-world data, we collect 13 scenes containing a mixture of in-the-wild (grocery store,

Test Scene	LSeg (3D)	OWL-ViT	LERF
waldo_kitchen	13.0%	42.6%	<b>81.5%</b>
bouquet	50.0%	66.7%	<b>91.7%</b>
ramen	15.0%	<b>92.5%</b>	62.5%
teatime	28.1%	75.0%	<b>93.8%</b>
figurines	8.9%	38.5%	<b>79.5%</b>
Overall	18.0%	54.8%	<b>80.3%</b>

Table 1: **Localization accuracy** comparison between Distilled Feature Fields using LSeg, OWL-ViT, and LERF. Overall performance is calculated by aggregating scene results. Refer to supplements for more details on scenes and text queries.

kitchen, bookstore) and posed long-tail (teatime, figurines, hand) scenes. We capture scenes using the iPhone app Polycam, which runs on-board SLAM to find camera poses, and use images of resolution  $994 \times 738$ .

#### 4.1. Qualitative Results

We visualize relevancy score by normalizing the colormap for each query from 50% (less relevant than canonical phrases) to the maximum relevancy. Extensive visualizations of all scenes can be found in the Appendix, and in Fig. 3 we select 5 representative scenes which demonstrate LERF’s ability to handle natural language. Visual comparison with LSeg in 3D are presented in Fig 8.

LERF captures language features of a scene at different levels of detail, supporting queries of properties like “yellow”, as well as highly specific queries like names of books and specific characters from TV shows (“jake from adventure time”). Because of the lack of discrete categories, objects can be relevant to multiple queries: in the figurine scene, abstract text queries can create semantically meaningful groupings. “Cartoon” selects the cat, Jake, rubber duck, miffy, waldo, toy elephant. “Bath toy” selects rubber-like objects, such as rubber duck, Jake, and toy elephant (made of rubber). Toy elephant is strongly highlighted for three different queries, demonstrating the ability of LERF to support different semantic tags for the same object.

#### 4.2. Existence Determination

We evaluate if LERF can detect whether an object exists within a scene. We label ground truth existence for 5 scenes, collecting two sets of labels: 1) labels from COCO to represent in-distribution objects to LSeg and 2) our own long-tail labels, which consist of queries of 15-20 objects in each scene concatenated together, for 81 total queries. See the Appendix for all queries. LERF determines whether an object exists in the scene by rendering a dense pointcloud over visible geometry, and returns “True” if any point has a relevancy score over a threshold.

We compare against distilling LSeg features into 3D as in DFF [21], but implemented in our own codebase for

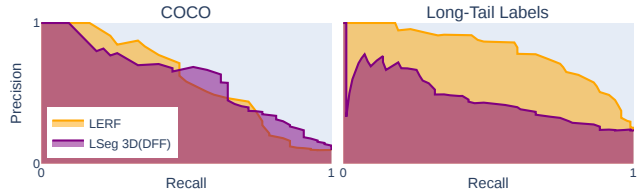


Figure 10: Precision recall curves for 3D existence experiments, Sec 4.2. LSeg performs similarly to LERF on in-distribution labels, but significantly suffers on long-tail labels of wild scenes.

an apples-to-apples comparison. We remove scale as a parameter to  $F_{\text{lang}}$  for LSeg since it outputs pixel-aligned features. We report precision-recall curves over relevancy score thresholds in Fig. 10. This experiment reveals that LSeg, trained on limited segmentation datasets, lacks the ability to represent natural language effectively. Instead, it only performs well on common objects that are within the distribution of its training set, as demonstrated in Fig. 8.

#### 4.3. Localization

To evaluate how well LERF can localize text prompts in a scene we render *novel* views and label bounding boxes for 72 objects across 5 scenes. For 3D methods we consider a label a success if the highest relevancy pixel lands inside the box, or for Owl-ViT if the center of the predicted box does. Results are presented in Table 1 and Fig. 9, and suggest that language embeddings embedded in LERF strongly outperform LSeg in 3D for localizing relevant parts of a scene. We also compare against the 2D open-vocab detector OWL-ViT by rendering full-HD NeRF views and selecting the bounding box with the highest confidence score for the text query. Owl-ViT outperforms LSeg in 3D, but suffers compared to LERF on long-tail queries.

#### 4.4. Ablations

**No DINO:** Removing DINO results in a qualitative deterioration in the smoothness and boundaries of relevancy maps, especially in regions with few surrounding views or little geometric separation between foreground and background. We show two illustrative examples where DINO improves the quality of relevancy maps in Fig. 6.

**Single-Scale Training:** We ablate multi-scale CLIP supervision from the pipeline by only training on a fixed  $s_0 = 15\%$  image scale. Doing so significantly impairs LERF’s ability to handle queries of *all* scales, failing on both large (“espresso machine”) queries it doesn’t have enough context for, as well as queries for which it does (“creamer pods”). These results imply that multi-scale training regularizes the language field at all scales, not just ones with relevant context for a given query.



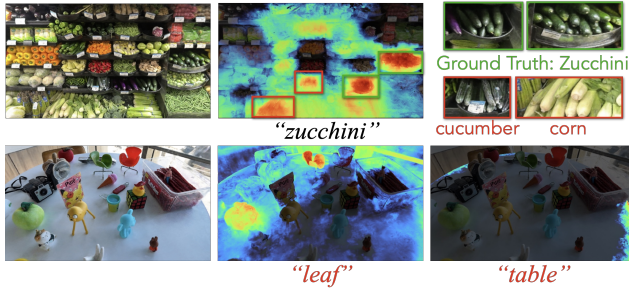


Figure 11: **Failure cases:** LERF struggles with identifying objects that appear visually similar to the query: “Zucchini” also activates on other long, green-ish vegetables, and “leaf” activates on the green plastic chair. LERF also struggles with global/spatial reasoning, where “table” only activate on the edges of the table.

## 5. Limitations

LERF has limitations associated with both CLIP and NeRF; some are visualized in Fig. 11. Like CLIP, language queries from LERF often exhibit “bag-of-words” behavior (i.e., “not red” is similar to “red”) and struggles to capture spatial relationships between objects. LERF can be prone to false positives with queries that appear visually or semantically similar: “zucchinis” activate on other similarly-shaped vegetables, though zucchinis are more relevant than the distractors (Fig. 11).

LERF requires known calibrated camera matrices and NeRF-quality multi-view captures, which aren’t always available or easy to capture. The quality of language fields is bottlenecked by the quality of the NeRF reconstruction. In addition, because of the volumetric input to  $F_{\text{lang}}$ , objects which are near other surfaces without side views can result in their embeddings being blurred to their surroundings since no views see the background without the object. This results in similar blurry relevancy maps to single-view CLIP (Fig. 5). In addition, we only render language embeddings from a single scale for a given query. Some queries could benefit from or even require incorporating context from multiple scales (eg “table”). See additional figures on limitations in the Appendix.

## 6. Conclusions

We present LERF, a novel method of fusing raw CLIP embeddings into a NeRF in a dense, multi-scale fashion without requiring region proposals or fine-tuning. We find that it can support a broad range of natural language queries across diverse real-world scenes, strongly outperforming pixel-aligned LSeg in supporting natural language queries. LERF is a general framework that supports any aligned multi-modal encoders, meaning it can naturally support improvements to vision-language models. Code and datasets will be released after the submission process.

## 7. Acknowledgements

This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE 2146752. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation. This project is funded in part by Bakar Fellows. The authors were supported in part by equipment grants from NVIDIA. We thank our colleagues who provided helpful feedback and suggestions, in particular Jessy Lin, Simeon Adebola, Satvik Sharma, Abhik Ahuja, Rohan Mathur, and Alexander Kristoffersen for their feedback on paper drafts, and Suzie Petryk, Boyi Li, Eric Wallace, and members of the Nerfstudio team for their conversations about the project.

## References

- [1] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. *arXiv preprint arXiv:2112.05814*, 2(3):4, 2021.
- [2] Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19129–19139, 2022.
- [3] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.
- [4] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- [6] Paola Cascante-Bonilla, Hui Wu, Letao Wang, Rogerio S Feris, and Vicente Ordonez. Simvqa: Exploring simulated environments for visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5056–5066, 2022.
- [7] Hila Chefer, Shir Gur, and Lior Wolf. Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In *ICCV*, pages 397–406, 2021.
- [8] Boyuan Chen, Fei Xia, Brian Ichter, Kanishka Rao, Keerthana Gopalakrishnan, Michael S Ryoo, Austin Stone, and Daniel Kappler. Open-vocabulary queryable scene representations for real world planning. *arXiv preprint arXiv:2209.09874*, 2022.
- [9] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask

- transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022.
- [10] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. *arXiv preprint arXiv:2212.07143*, 2022.
- [11] Rodolfo Corona, Shizhan Zhu, Dan Klein, and Trevor Darrell. Voxel-informed language grounding. *arXiv preprint arXiv:2205.09710*, 2022.
- [12] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017.
- [13] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022.
- [14] Zhiwen Fan, Peihao Wang, Yifan Jiang, Xinyu Gong, De-jia Xu, and Zhangyang Wang. Nerf-sos: Any-view self-supervised object segmentation on complex scenes. *arXiv preprint arXiv:2209.08776*, 2022.
- [15] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Open-vocabulary image segmentation. *arXiv preprint arXiv:2112.12143*, 2021.
- [16] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4089–4098, 2018.
- [17] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021.
- [18] Huy Ha and Shuran Song. Semantic abstraction: Open-world 3d scene understanding from 2d vision-language models. In *Conference on Robot Learning*, 2022.
- [19] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. *arXiv preprint arXiv:2210.05714*, 2022.
- [20] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, Ayush Tewari, Joshua B. Tenenbaum, Celso Miguel de Melo, Madhava Krishna, Liam Paull, Florian Shkurti, and Antonio Torralba. Conceptfusion: Open-set multimodal 3d mapping, 2023.
- [21] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. In *NeurIPS*, volume 35, 2022.
- [22] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022.
- [23] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yanan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. *arXiv preprint arXiv:2210.04150*, 2022.
- [24] Timo Lüddecke and Alexander Ecker. Image segmentation using text and image prompts. In *CVPR*, pages 7086–7096, 2022.
- [25] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [26] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022.
- [27] Jishnu Mukhoti, Tsung-Yu Lin, Omid Poursaeed, Rui Wang, Ashish Shah, Philip HS Torr, and Ser-Nam Lim. Open vocabulary semantic segmentation with patch aligned contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19413–19423, 2023.
- [28] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [29] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. *arXiv preprint arXiv:2211.15654*, 2022.
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [31] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. 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*, pages 10901–10911, 2021.
- [32] Zhongzheng Ren, Aseem Agarwala<sup>†</sup>, Bryan Russell<sup>†</sup>, Alexander G. Schwing<sup>†</sup>, and Oliver Wang<sup>†</sup>. Neural volumetric object selection. In *CVPR*, 2022. (<sup>†</sup> alphabetic ordering).
- [33] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, pages 618–626, 2017.
- [34] Nur Muhammad Mahi Shafullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv:2210.05663*, 2022.
- [35] Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Buló, Norman Müller, Matthias Nießner, Angela Dai, and Peter Kotschieder. Panoptic lifting for 3d scene understanding with neural fields. *arXiv preprint arXiv:2212.09802*, 2022.
- [36] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijnmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl

- Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.
- [37] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, et al. Nerfstudio: A modular framework for neural radiance field development. *arXiv preprint arXiv:2302.04264*, 2023.
- [38] Jesse Thomason, Mohit Shridhar, Yonatan Bisk, Chris Paxton, and Luke Zettlemoyer. Language grounding with 3d objects. In *Conference on Robot Learning*, pages 1691–1701. PMLR, 2022.
- [39] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations. In *3DV*, 2022.
- [40] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. Cris: Clip-driven referring image segmentation. In *CVPR*, pages 11686–11695, 2022.
- [41] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew Davison. In-place scene labelling and understanding with implicit scene representation. In *ICCV*, 2021.
- [42] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, pages 350–368. Springer, 2022.