

# Language Embedded Radiance Fields for Zero-Shot Task-Oriented Grasping Supplementary Material

Anonymous Author(s)

Affiliation

Address

email

## 1 A LERF-TOGO

### 2 A.1 Implementation Details

3 We implement LERF-TOGO on top of the Nerfacto method from Nerfstudio [1]. For faster conver-  
4 gence and smoother optimization, we modify several parameters from the original LERF paper. We  
5 use a smaller hashgrid with 16 levels and a maximum resolution of 256, and find that using larger  
6 MLPs for the density, color, and transient output heads of NeRF results in faster convergence and  
7 better ability to handle specularities and robot shadows. We introduce weight decay of  $1e-7$  to the  
8 LERF network which smooths training. In addition, we compress the DINO embeddings into 128  
9 dimensions and supervise on these vectors rather than the original DINO outputs, but we do not  
10 normalize the resulting vectors like Tschernezki et al. [2]. While constructing the CLIP embedding  
11 pyramid for LERF, we use crops ranging from 5% to 35% of image height with 6 pyramid lev-  
12 els, biasing the pyramid to smaller crops as LERF-TOGO is primarily interested in object and part  
13 queries.

## 14 B Grasping

15 **Point Cloud Extraction** To extract a scene-wide point cloud for grasping, we use the method in  
16 Nerfstudio [1], which deprojects randomly sampled rays' depth from the train camera poses, then  
17 filters with outlier rejection. We then crop the point cloud to the workspace of the robot. For object  
18 centric point clouds, we deproject depth from views radially surrounding the object of interest.

19 **Motion Planning** A grasp is considered feasible if the robot can perform a collision-free trajectory  
20 with the following poses: the pre-grasp, grasp, and post-grasp configurations. The pre-grasp pose is  
21 positioned 5cm along the z-axis of the robot end effector, which allows the gripper to approach the  
22 target grasp pose with minimal additional motion. The post-grasp pose is located 10cm above the  
23 grasp pose, along the z-axis of the world frame. The UR5's IK solver calculates the set of viable  
24 joint configurations at these poses, and we calculate the trajectory as a linear interpolation between  
25 them. We additionally allow for a 180 degree rotation at the last wrist joint, as parallel-jaw grasps are  
26 rotationally symmetric. This facilitates the motion planning process, as the robot's camera mount is  
27 prone to colliding with the robot arm.

## 28 C Experiments

### 29 C.1 Setup Details

30 We use a UR5 arm with a Logitech BRIO webcam at 1600x896 resolution, with all camera settings  
31 frozen before each capture prevent color discrepancies among images. The camera mount points

Scene	Object Query ; Part Query
Kitchen	(black matte spoon, handle), (shiny black spoon, handle), (teapot, handle), (dish scrub brush, handle), (dust brush, handle)
Flowers	(daisy, plant stem), (rose, plant stem)
Mugs	(blue mug, handle), (pink teacup, handle), (turquoise mug, handle), (white mug, handle), (black mug, handle)
Tools	(Measuring tape, base), (screwdriver, handle), (wire cutters, handle) (soldering iron, handle), (hammer, handle)
Knives	(bread knife, handle), (steak knife, handle), (box cutter, handle)
Martinis	(red martini glass, stem), (grey martini glass, stem)
Fragile	(camera, strap), (pink sunglasses, earhooks) (blue sunglasses, earhooks), (lightbulb, screw)
Cords	(power strip, plug), (power strip, base), (ethernet dongle, usb) (ethernet dongle, ethernet)
Messy	(ice cream, cone), (green lollipop, stick), (blue lollipop, stick)
Pasta	(wine, cork), (wine, bottle neck), (saucepan, lid knob) (saucepan, handle), (corkscrew, handle)
Cleaning	(clorox, wet towel), (clorox, lid), (clorox, body) (tissue box, box), (tissue box, tissue)
Bottles	(meyer’s cleaning spray, spray trigger), (meyer’s cleaning spray, bottle neck) (meyer’s cleaning spray, body) (purple cleaning spray, body) (purple cleaning spray, spray trigger), (purple cleaning spray, bottle neck)

Table 1: **Complete list of object and part queries**

orthogonally to the gripper axis, to maximize the reachable workspace of the camera while pointing towards the workspace center. During robot capture, pre-computation of DINO, CLIP, and ZoeDepth is parallelized across 3 NVIDIA 4090 GPUs to achieve real-time performance, and all subsequent operations are carried out on a single 4090. Capturing a scene takes 30 seconds, training the LERF to 2k iterations takes 78 seconds, and finally querying LERF-TOGO takes 10 seconds.

## References

- [1] M. Tancik, E. Weber, E. Ng, R. Li, B. Yi, J. Kerr, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, et al. Nerfstudio: A modular framework for neural radiance field development. *arXiv preprint arXiv:2302.04264*, 2023.
- [2] V. Tschernezki, I. Laina, D. Larlus, and A. Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2022.