Weakly-Supervised Learning of Dense Functional Correspondences

Stefan Stojanov* Linan Zhao* Yunzhi Zhang Daniel L. K. Yamins Jiajun Wu Stanford University

Abstract

Establishing dense correspondences across image pairs is essential for tasks such as shape reconstruction and robot manipulation. In the challenging setting of matching across different categories, the function of an object, i.e., the effect that an object can cause on other objects, can guide how correspondences should be established. This is because object parts that enable specific functions often share similarities in shape and appearance. We derive the definition of dense functional correspondence based on this observation and propose a weakly-supervised learning paradigm to tackle the prediction task. The main insight behind our approach is that we can leverage vision-language models to pseudo-label multi-view images to obtain functional parts. We then integrate this with dense contrastive learning from pixel correspondences to distill both functional and spatial knowledge into a new model that can establish dense functional correspondence. Further, we curate synthetic and real evaluation datasets as task benchmarks. Our results demonstrate the advantages of our approach over baseline solutions consisting of off-the-shelf self-supervised image representations and grounded vision language models.

1. Introduction

Finding pixel correspondence across image pairs is fundamental for object understanding and is critical for applications like shape reconstruction [40, 45, 48, 76], editing [19], and object manipulation in robotics [17, 29, 30, 58]. This task requires reasoning beyond visual similarity in local appearance, geometry, and texture across images. It also involves structural similarity, *e.g.*, the part-whole relationships of objects and their part components, and semantic similarity, *e.g.*, the functional properties of parts of objects.

These aspects of similarity are essential for learning efficient generalizable systems for downstream applications. For example, in imitation learning in robotics, human demonstrations are a scarce and valuable data source.

Given a demonstration with an object, such as pouring with a kettle, establishing dense functional correspondence with another object that supports this function, *e.g.*, a bottle, enables the efficient transfer of the demonstration.

It becomes harder to find dense correspondence when the input images shift from being two views of the same object to different objects from the same category, and finally to objects from distinct categories, as the visual similarity becomes less apparent. This work focuses on the most challenging scenario with objects from different categories. We aim to establish dense pixel-level correspondence between pairs of images containing objects with parts whose shape enables the execution of similar functions. Specifically, by "function", we refer to the effect one object can have on another object or substance, *e.g.*, the function "cut-with" for a knife and a spatula or "hang-onto" for objects with hooks.

Practically, training and evaluation for this task are challenging due to the lack of labeled data. Supervised training at scale is infeasible because manual dense correspondence labeling is intractable, emphasizing the need for a self- or weakly-supervised approach. For evaluation, while datasets exist for dense within-category correspondence [35, 72, 81] and sparse functional keypoint correspondence across categories [38], there is still no established task or dataset for dense correspondence across categories. In this work, we make progress toward addressing both the challenges of training and evaluation.

The key insight behind our training approach is that the capabilities of self-supervised image representations like DINOv2 [56] or Stable Diffusion [66] and vision language models (VLMs) [25, 82] are complementary but individually insufficient for solving this task. On the one hand, surprisingly accurate dense correspondences can be established using image features from pre-trained self-supervised models. This works well when the input images contain visually similar object instances from the same category, *e.g.*, two cats or two cars [91]. However, the accuracy decreases for the more generic scenario when objects come from distinct categories. On the other hand, VLMs can detect the bounding boxes of object parts with similar functions in a zero-shot manner [25, 82] but cannot perform fine-grained reasoning about correspondences across objects.

^{*}Equal contribution.

 $^{^{1}}Project\ website:$ https://dense-functional-correspondence.github.io/

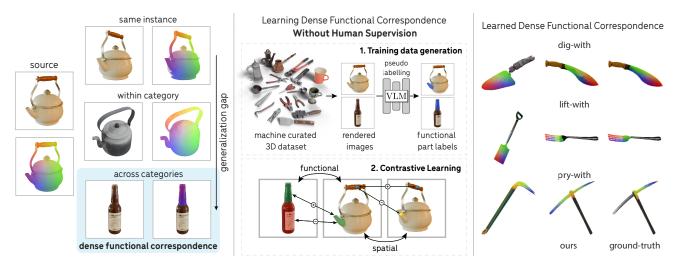


Figure 1. **Dense Functional Correspondence** refers to establishing dense correspondences across object instances based on function similarity (*e.g.*, "pour-with"). This task is especially challenging when objects have visually different but functionally similar parts, requiring both semantic understanding, *i.e.*, identifying which parts can perform the same function, and structural understanding, *i.e.*, establishing dense correspondence across the parts at a surface-point-level based on functionally equivalent alignment. We propose a method to learn such correspondences with little human supervision, leveraging automated data curation and annotation, and dense contrastive learning.

We distill the strengths of each approach into a new model using a scalable technique that requires minimal human supervision. Specifically, we first obtain multi-view-consistent pseudo labels of functionally relevant regions of 3D object assets [7] using an off-the-shelf grounded VLM [82]. We then combine these labels with multi-view correspondences [17, 68] in a contrastive learning framework building on pre-trained DINOv2 [56] feature extractor. For evaluation, we define the dense functional 2D correspondence task and develop an annotation procedure based on aligning 3D object pairs in functionally equivalent poses.

In sum, we define the task of dense functional correspondence as a means for investigating cross-category dense correspondence. We then curate synthetic and real-world evaluation datasets for this task. We further propose a scalable, weakly-supervised method leveraging vision foundation models, which empirically outperforms baselines.

2. Related Work

Learning representations to establish dense functional correspondence requires fine-grained structural and semantic visual reasoning about objects. The most relevant prior works come from the object-level correspondences and affordance learning research domains. We also review recent work on vision foundation models and VLMs, focusing on works relating to fine-grained object understanding.

Learning Correspondences. For this work, it is relevant to categorize correspondence learning methods based on their degrees of generalization. For generalizing across geometric scene transforms, works on multi-view correspondences aim to match different views of the same scene [28, 67, 69], whereas optical flow techniques match consecu-

tive video frames [26, 27, 73]. For generalizing within categories, NOCS-style representations [35, 80, 81] enable dense matching across instances of a category, whereas learning sparse keypoints [49, 70] enables sparse matching based on pre-defined semantic keypoint taxonomies. For generalizing across categories, Lai et al. [38] propose matching based on object function by learning five keypoints per object function category. The main drawback of keypoint-based correspondences is the requirement for a keypoint taxonomy, which by definition limits such techniques' capability to capture nuanced similarities across highly dissimilar objects (*e.g.*, a bottle and a kettle). Through our dense functional correspondence formulation, we overcome the limitation of keypoint definitions and enable higher precision in downstream applications.

Learning Affordances. In his seminal work [18], James J. Gibson defines affordances as objects' "opportunities for interaction." Various object affordance formalisms have been developed in computer vision and robotics, such as estimating grasps [2, 16, 50, 51], and localizing affordance regions in 2D [5, 14, 46, 52-54] and 3D [10, 20, 87] through bounding boxes and segments [14, 52–54], heatmaps [10, 46, 52] or keypoints [62, 77, 85]. Early works adopt a fully supervised learning paradigm [2, 14, 52], while more recent works aim to use less supervision by learning from human object interaction videos [53], egocentric videos [42] or unlabeled exocentric images [41, 46]. Our work has two key distinctions: First, affordance heatmaps or segments identify object regions or parts in individual images. They do not allow for fine-grained spatial correspondence across object parts in different images (e.g., can identify the blades of two knives but cannot find correspondences for pixels between the tips or edges of the blades). Second, our focus is on object function – the effect an object can cause on something else, rather than the broader concept of affordance, which emphasizes potential interactions with a specific object instance (*e.g.*, striking with a hammer vs. holding). Last, our goal is to learn dense functional correspondence in a weakly-supervised manner, without relying on human annotations of ground-truth correspondences.

Vision Foundation Models. Recent developments in large-scale language [11, 64, 75] and image [63, 90] pretraining have led to the development of vision-language models (VLMs) capable of strong zero-shot performance through vision-question answering [44, 82], which have been adapted to reasoning about functional affordances and grasping in robotics [13, 25, 61, 89]. Powerful correspondence representations have been found to emerge [1, 71] in DINO [3, 56] and Stable Diffusion [66], which have led to direct applications in low-shot affordances [30] and object manipulation [12, 30, 36, 58, 59]. In this work, we leverage the complementary characteristics of VLMs and self-supervised image models to go beyond their individual capabilities for dense functional correspondence.

3. Dense Functional Correspondence

Distinct object categories with similar functionality, *e.g.*, a "kettle" and a "bottle" which can both pour liquid, may have different visual shapes and appearances as well as distinct part organizations. However, individual parts that serve the specific functionality of interest, *e.g.*, the spout of a kettle and the mouth of a bottle in this example, have a higher resemblance with each other than at the overall object level. Such consistency is a consequence of how form follows function – object parts that fulfill a specific function tend to remain consistent across objects, even if other parts vary greatly in shapes and arrangements. The partlevel consistency provides a crucial ground from which we can derive the definition of functional correspondence (Section 3.1) and develop a corresponding evaluation data curation pipeline to benchmark this task (Section 3.2).

3.1. Problem Definition

We refer to the effect that an object causes on other objects or substances as an "object function." This concept has been widely studied in model generalization in visual computing [37, 38, 74, 93] and the development of categorization in humans [32, 39, 83]. Examples are shown in Figure 1, e.g., "pour with." When executing a function with an object, such as pouring with a kettle, the functional part (the spout) follows a specific 3D trajectory. To replicate this function with a different object, e.g., a bottle, the neck of the bottle would be aligned with the spout and follow the same trajectory. This illustrates how the *same* object function is fulfilled with different objects via aligning function-

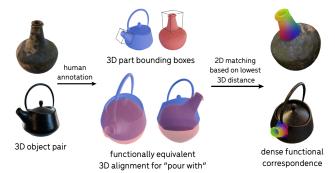


Figure 2. Annotation Pipeline (Evaluation Only). Given a 3D object pair (left) and a function ("pour-with"), we annotate the functional alignment of two objects by aligning the functional parts in 3D (middle). Afterward, we derive dense 2D correspondences (right) based on 3D distances of corresponding object surface points, with matching pixels shown in the same color.

ally equivalent parts, which is central to robotic applications with imitation learning approaches [23, 58, 94].

The above observations lead us to define dense functional correspondence through 3D object alignment based on functionally equivalent parts. Specifically, given two objects (e.g., a kettle and a bottle) and an object function (e.g., "pour with"), the objects are aligned if and only if the parts that fulfill this function (e.g., the kettle spout and the bottle neck) are spatially close to each other. The alignment induces an image-space distance: for any pair of pixels on the functional parts of two objects, the pixels are in functional correspondence if their respective surface points are close in 3D when the objects are aligned. Since this distance is defined at the pixel level, it is inherently dense.

Formally, the input consists of an object function $\mathcal F$ and an image pair (I_1,I_2) , where each image is a view of a 3D object O_1 and O_2 . Let $\pi^{-1}:I\to O$ represent the backprojection function that maps an image pixel to the corresponding 3D object surface point. We define $M(O;\mathcal F)$ as the functional part of object O responsible for executing $\mathcal F$, and let $M(I;\mathcal F)$ be its projected 2D mask in the image. In our setup, the functional parts of both objects, $M(O_1;\mathcal F)$ and $M(O_2;\mathcal F)$, are assumed to be aligned in 3D such that they follow the same trajectory when performing $\mathcal F$. We therefore define dense functional correspondence as a mapping $f(I_1,I_2;\mathcal F):M(I_1;\mathcal F)\to M(I_2;\mathcal F)$ that minimizes $\sum_{p\in M(I_1;\mathcal F)}||\pi^{-1}(p)-\pi^{-1}(f(p))||_2$. This ensures that pixel pairs in functional correspondence are from spatially close locations in 3D when the objects are aligned.

3.2. Evaluation Dataset Curation

The problem definition in Section 3.1 provides a guiding principle to obtain ground truth annotations for dense functional correspondence in image pairs by *aligning objects in 3D*. We introduce the annotation procedure and use it to construct both synthetic and real-world evaluation datasets for quantitative evaluation in Section 5.

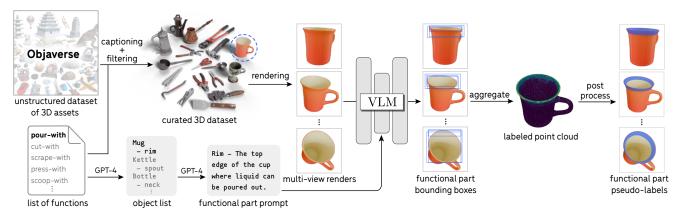


Figure 3. **Training Data Curation via VLM Pseudo Labeling.** Given a large unstructured dataset like Objaverse [7], we leverage off-the-shelf VLMs to curate and label the functional parts. Specifically, GPT-4 [55] generates category-specific functional part prompts, and CogVLM [82] produces bounding box proposals for multi-view image renderings, which are aggregated onto a 3D point cloud. The point cloud is post-processed to produce pixel-level functional part labels for training.

Annotation Procedure. To obtain ground-truth functional correspondence for an image pair, we assume each object is rendered from a known 3D asset. By aligning the two assets in 3D, we derive dense pixel correspondences between the images. This procedure eliminates the need for manual dense 2D labeling, enabling large-scale evaluation. An overview is shown in Figure 2.

Specifically, given two 3D meshes of objects supporting the same function, we first align them based on their functional parts and annotate a 3D bounding box around each functional part. Then, for a pair of rendered images, we unproject pixels from the functional parts onto the object surfaces and compute 3D distances between these points to perform minimum-cost matching. Pixels corresponding to visible surface points that are spatially close in 3D are matched. A detailed description of the annotation procedure is provided in the Supplement.

Synthetic Evaluation Dataset. We use the 3D assets from Objaverse [7], a large collection of diverse 3D models, to obtain a synthetic evaluation dataset. We hand-label 3D annotations for 950 pairs of assets from Objaverse spanning 24 functions, selected for quality and diversity. See Section 4.1 for how the assets and functions are selected. As such, 85% of the ground-truth pairs contain across-category correspondences. From these 3D annotations, we derive 1,800+ unique 2D image pairs rendered from the 3D assets, with ground truth dense functional correspondences.

Real Evaluation Dataset. Setting up a real-world benchmark is crucial for measuring model performance on real images. Thus, we utilize the HANDAL dataset [20], which contains images and 3D reconstructions of real-world objects. After manually fixing the geometry of the 3D scanned assets, *e.g.*, the missing concavities of mugs, pots and pans, we label 190 asset pairs spanning 13 functions. This results in a real evaluation dataset of 500+ unique 2D real image pairs with ground-truth functional correspondence labels.

4. Approach

Our goal is to develop a scalable learning framework for dense functional correspondences without relying on human-labeled ground truth. Since this task requires both semantic and structural knowledge, we distill from off-the-shelf VLMs to obtain pseudo-labeled training data (Sec. 4.1), which is further combined with dense spatial correspondences from synthetic data in a contrastive learning framework (Sec. 4.2). This approach enables the model to generalize to real-world data, as we will show in Sec. 5.

4.1. Dataset of Pseudo-labeled Functional Parts

A dataset for learning dense functional correspondences at scale requires a diverse source of object images, a diverse taxonomy of functions and associated functional parts, and a low-cost, reliable means for part labeling.

Image Data. Our approach requires a large and diverse multi-view image dataset where functional parts are visible. Existing multi-view object datasets [65, 84, 88, 92] are suboptimal because they have few desired objects like tools and utensils, the objects are in canonical poses that may not reveal functional parts, or are placed in cluttered contexts where occlusions often occur. To overcome this, we render high-quality images from the Objaverse [7] dataset using ray-tracing and HDRI environments [21] in Blender [6], obtaining arbitrary amounts of diverse multi-view data.

Object and Function Taxonomy. To curate relevant object assets for our training dataset, we prompt GPT-4 [55] for common functions and refer to object functions studied in [38, 51]. Then, we prompt GPT-4 to generate a comprehensive list of object categories for each function. After deduplication and manual filtering, our taxonomy has 24 functions and 160 object categories.

Object Asset Selection. To retrieve relevant assets from the noisy-labeled Objaverse dataset based on the list of object

categories, we utilize asset captions from Caption3D [47]. We use Llama 3.1 [15] to summarize the captions into category names and use Llama word embeddings to match the summaries to our category list. Finally, we prompt Llama to verify these matches. To ensure diversity, we cap each category at 200 assets. To ensure quality, we manually filter the retrieved assets to obtain 8,285 assets in total, 80% of which are used for training. Details about prompting, filtering, and the taxonomy are included in the Supplement.

Functional Part Pseudo-Labeling. Labeling data at scale using large pre-trained models has been shown as an effective approach for achieving high performance with minimal human effort [79, 86]. The key elements for success are a sufficiently accurate pre-trained model and a low-cost and reliable procedure for rejecting low-quality labels. Grounded VLMs [8, 44, 82] have shown remarkable capabilities for zero-shot prompt-based object detection. We, therefore, use the 17B grounded CogVLM [82] model, which has state-of-the-art referring expression detection performance. For an overview of the pseudolabeling pipeline, see Figure 3. Given our list of object categories and functions, we prompt GPT-4 to obtain the names and appearance descriptions of functional parts to serve as prompts for CogVLM, which we then manually filtered and deduplicated. Because functional part names can be different across categories (e.g., the spout of a kettle vs. the neck of a bottle), we generate these functional part lists separately for each category. We empirically found that prompting CogVLM with part names and appearance descriptions significantly improves the bounding box predictions.

Given a set of rendered views for an object and a functional part text prompt, we generate bounding box predictions with CogVLM [82], which vary due to sampling in VLM inference. The accuracy of the bounding boxes also depends on viewpoint because of part pose and visibility. To aggregate these possibly noisy labels and obtain a final part label, we sample a dense point cloud on the surface of the object, and accumulate the 2D labels across views onto the 3D points. We post-process these labeled point clouds to generate 2D masks for views rendered for training.

This dataset curation and pseudo-labeling procedure allows us to generate a large dataset of functional part segmentation labels with relatively little human effort, which was mostly necessary for prompt engineering and quality control. In this work, we apply this approach on the $\approx 600 \mathrm{K}$ labeled meshes from Caption3D, but it is straightforward to scale up to the millions of meshes in Objaverse-XL [9].

4.2. Learning Dense Functional Correspondence

To learn dense functional correspondence, we train a feature embedding that captures both the high-level function semantics and the structural similarity between functional parts. For instance, given a bottle and a kettle for the func-

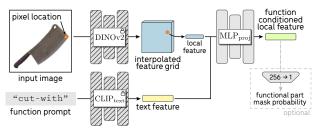


Figure 4. **Local Functional Feature Extraction.** To obtain dense functionally conditioned features, we apply an MLP on top of a function text embedding and the spatial DINO features. The MLP is trained with both functional and spatial contrastive losses.

tion "pour-with," the features for the neck of the bottle and the spout of the kettle should be similar. Moreover, the mouth of the bottle and the tip of the kettle spout should be in correspondence, as well as the bottom of the bottle's neck and the bottom of the kettle's spout. To achieve this, we train a function-conditioned network on top of frozen DINOv2 [56] and CLIP [63] (illustrated in Figure 4), that is applied at the local feature level. Because of significant developments in object segmentation [31, 34] and our focus on object-level understanding, we assume that the input images consist of segmented objects.

Function-Conditioned MLP. Given an image and a function, we first extract the image features from the last three blocks of DINOv2 and the function conditioning from CLIP text embeddings. We average the DINOv2 features from each block using learned weights into a single feature grid, and use bilinear interpolation to obtain a feature vector for each pixel location. Then, we concatenate the image feature with the CLIP embedding of the function and pass it through a 3-layer MLP, which produces the final feature at each pixel location. This network can be thought of as a function-conditioned version of the final projection layer used in contrastive learning [4, 22]. We parameterize our model as $g_{\theta}(p|I, \mathcal{F})$, which outputs the normalized feature of pixel p on image I conditioned on the function \mathcal{F} .

We also investigate the option of adding an extra fully connected layer that maps the output feature vector to a prediction for the functional part mask. This allows us to obtain a binary functional part mask at inference time.

Functional Part Contrastive Learning. To distill the knowledge of functional part semantics from the VLM, we use contrastive learning based on the pseudo-labeled functional part masks. The parts from two objects that can be used to perform the same function should share a more similar embedding space. Specifically, given two images, I_1 and I_2 of objects that can perform the same function \mathcal{F} , let the functional part segments be P_1^+ and P_2^+ . Then, define the rest of the objects' pixels as P_1^- and P_2^+ . Learning correspondence requires the pixels in P_1^+ to be similar to the ones in P_2^+ but different from the ones in P_2^- . In addition, to encourage the model to focus on the functionally relevant

regions of objects, we add a term that pushes the features of P_1^- away from that of P_2^- .

Let $\sin(x,y|I_1,I_2,\mathcal{F})=g(x|I_1,\mathcal{F})\cdot g(y|I_2,\mathcal{F})$ represent the feature similarity between pixel x on image I_1 and pixel y on image I_2 when conditioned on function \mathcal{F} . For brevity, we short-hand it to $\sin(x,y)$ below. The infoNCE loss [78] for the function-part contrastive learning given $(p_1^+,p_1^-,p_2^+,p_2^-)\in (P_1^+,P_1^-,P_2^+,P_2^-)$ is thus

$$\mathcal{L}_{\text{func}} = -\log \frac{e^{\sin(p_1^+, p_2^+)/\tau}}{e^{\sin(p_1^+, p_2^+)/\tau} + e^{\sin(p_1^+, p_2^-)/\tau} + e^{\sin(p_1^-, p_2^-)/\tau}}$$
(1)

for temperature τ .

When the model predicts functional part masks, we add a binary cross-entropy loss \mathcal{L}_{mask} to compare the predicted mask with the pseudo-labeled functional part segment.

Part Structure via Multi-view Contrastive Learning. If we train the embedding with only the functional part contrastive loss, we inevitably run into mode collapse issues. That is, the whole spout of the kettle would have the same features regardless of the pixel's spatial location. To preserve the structural information, we apply dense contrastive learning from multi-view correspondences.

Given two views of an object, we can find corresponding pixels that project to the same location in 3D space. We require a view-invariant feature embedding – a pixel should have high similarity with its corresponding pixel on the other image but remain different from all the other pixels. This encourages the model to learn the structural information of the object, to not collapse the embedding space, and to encode the object part consistently across different views. This multi-view contrastive objective only applies to two images of the same asset. However, because the underlying DINOv2 embedding space enables generalization for visually similar regions, the trained feature embedding can retain information about the structural similarities between functional parts *across* categories.

Formally, let q be a pixel in the first view I, q'_+ be a pixel in the second view I' that corresponds to the same location in 3D as q, and any other pixel on I' be denoted as q'_- . The multiview contrastive objective is

$$\mathcal{L}_{\text{spatial}} = -\log \frac{e^{\sin(q, q'_+)/\tau}}{e^{\sin(q, q'_+)/\tau} + e^{\sin(q, q'_-)/\tau}}.$$
 (2)

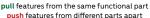
Combining the terms, we obtain the final loss

$$\mathcal{L} = \mathcal{L}_{\text{func}} + \lambda_{\text{spatial}} \mathcal{L}_{\text{spatial}} + \lambda_{\text{mask}} \mathcal{L}_{\text{mask}}.$$
 (3)

4.3. Implementation Details

We use DINOv2-B as the backbone and an image size of 224. The MLP projector has 3 layers with 1024 hidden dimensions each. We use the Adam [33] optimizer with







pull features from the same 3D point close push features from different points apart

Figure 5. **Training Objectives.** To ensure functional part similarity in the learned feature space, we use a part-level contrastive objective to distill *functional part semantics* from VLMs (left). The *spatial* contrastive loss (right) serves a complementary role and prevents the model from collapsing predictions for different regions of a part, *e.g.*, the top and bottom of a kettle spout.

default hyperparameters, a batch size of 50 image pairs, 128 positive and negative sampled points on each image, and a learning rate of 1×10^{-4} . In addition, we use a weight of $\lambda_{\rm spatial}=10$ for the spatial loss and a weight of $\lambda_{\rm mask}=1$ for the mask loss. We use random-color background augmentation during training following [17]. A sensitivity analysis of loss weights and a breakdown of computational costs are provided in the Supplement.

5. Experiments

In this section, we benchmark our approach in Sec. 4 and several baseline solutions on the dense functional correspondence task. Since our problem formulation in Sec. 3.1 requires a function as input and focuses on matches within functional parts, it differs significantly from existing benchmarks on semantic correspondence [49, 70]. As such, we leverage the evaluation datasets from Sec. 3.2.

5.1. Metrics

We evaluate dense functional correspondence from two different aspects: *correspondence label transfer*, which assesses the precision with which the model can transfer one functional part to another, and *correspondence discovery*, which assesses the model's ability to identify relevant functional correspondences without any reference input labels.

Correspondence Label Transfer. To evaluate the precision of the correspondences that can be found using the learned features, we use normalized pixel distance (Normalized Dist) and percentage of correct keypoints (PCK).

Specifically, let the ground-truth correspondences between images I_1, I_2 given the function $\mathcal F$ be $\{p_1^1, p_1^2, \cdots, p_1^k\}, \{p_2^1, p_2^2, \cdots, p_2^k\}$. For each pixel p_1^i on image I_1 , we can find its most similar match $p_2^{j(i)}$ on I_2 using feature similarity. The normalized distance metric is simply the mean of $||p_2^{j(i)}-p_2^i||_2$ normalized by the image size, and PCK@k pixels is the mean of $\mathbbm{1}_{||p_2^{j(i)}-p_2^i||_2 < k}$.

Correspondence Discovery. In addition to label transfer, models should discover the relevant set of functional correspondences on its own, without assuming a priori that the

Model	Correspond	lence Label Tr	ansfer	Correspondence Discovery					
Wide	Normalized Dist (\(\psi \)	PCK@23p (†)	PCK@10p (†)	Best F1@23p (†)	Best F1@10p (†)	AP@23p (†)	AP@10p (†)		
Synthetic Evaluation Dataset									
Chance	0.310	0.165	0.046	0.416	0.176	0.256	0.093		
DINO [56]	0.212	0.381	0.148	0.578	0.281	0.381	0.130		
SD [91]	0.268	0.298	0.126	0.479	0.231	0.267	0.097		
SD-DINO [91]	0.227	0.376	0.161	0.563	0.301	0.341	0.144		
CogVLM [82] + DINO	0.180	0.416	0.158	0.678	0.333	0.556	0.188		
ManipVQA-P [25] + DINO	0.223	0.346	0.130	0.575	0.269	0.418	0.134		
ManipVQA-F [25] + DINO	0.272	0.259	0.093	0.528	0.244	0.320	0.097		
Ours (functional only)	0.228	0.287	0.094	0.575	0.233	0.441	0.112		
Ours (spatial only)	0.204	0.470	0.227	0.610	0.369	0.412	0.211		
Ours (full without mask loss)	0.170	0.486	0.227	0.768	0.470	0.685	0.338		
Ours (full with mask loss)	<u>0.172</u>	0.480	0.223	0.774	0.471	<u>0.684</u>	0.330		
Real Evaluation Dataset									
Chance	0.313	0.170	0.045	0.417	0.167	0.248	0.087		
DINO [56]	0.206	0.408	0.159	0.589	0.294	0.382	0.138		
SD [91]	0.259	0.309	0.127	0.503	0.238	0.285	0.101		
SD-DINO [91]	0.220	0.385	0.163	0.577	0.301	0.343	0.142		
CogVLM [82] + DINO	0.172	0.440	0.169	0.695	0.350	0.561	0.198		
ManipVQA-P [25] + DINO	0.204	0.398	0.153	0.600	0.295	0.420	0.148		
ManipVQA-F [25] + DINO	0.256	0.309	0.114	0.575	0.281	0.368	0.126		
Ours (functional only)	0.200	0.336	0.115	0.652	0.283	0.532	0.148		
Ours (spatial only)	0.203	0.472	0.228	0.708	0.353	0.382	0.182		
Ours (full without mask loss)	0.152	0.516	0.249	0.775	0.476	0.691	0.344		
Ours (full with mask loss)	<u>0.153</u>	0.501	0.235	0.808	0.502	0.730	0.360		

Table 1. **Quantitative Evaluation** on the synthetic and real evaluation datasets. The simplest baselines, self-supervised features from Stable Diffusion and DINOv2, perform relatively poorly. Adding semantic knowledge from predicted functional part labels from VLMs can offer slight improvement. Our approach, combining the strengths of both self-supervised features and VLMs, achieves the best performance.

relevant pixels on one image have been given. This capability is essential for potential downstream applications such as object alignment in robot object manipulation.

First, since we assume that the input images are segmented, let M_1, M_2 be the object masks for images I_1, I_2 . For every pixel $p_1^i \in M_1$, we find its most similar match $p_2^{j(i)}$ on I_2 and find the backward match of $p_2^{j(i)}$ on I_1 , denoted as q_1^i . As such, $||p_1^i-q_1^i||_2$ captures the level of cycleconsistency of the match. We therefore construct a score $s=(1-||p_1^i-q_1^i||_2)\cdot \sin(p_1^i,p_2^{j(i)})$ to rank each pair of $(p_1^i,p_2^{j(i)})$, using both similarity and cycle consistency.

Then, we consider the top t% of all pairs as "discovered" and compare them with the ground-truth. A discovered pair (x_1,x_2) is equivalent to a ground-truth pair (y_1,y_2) if both end points are within k pixels of the ground truth. Increasing t results in higher recall but potentially lower precision: the number of discovered ground-truth correspondences monotonically increases while the percentage of correct correspondence tends to decrease. Sweeping t produces a precision-recall curve, from which we can calculate the best F1 score (at k pixels) and the average precision (AP) (at k pixels). Formally, Best F1 = $\max_t \frac{2\times \operatorname{Precision}_t \times \operatorname{Recall}_t}{\operatorname{Precision}_t + \operatorname{Recall}_t}$ and $\operatorname{AP} = \sum_t (\operatorname{Recall}_t - \operatorname{Recall}_{t-1}) \operatorname{Precision}_t$.

5.2. Baselines

We describe several baseline methods below.

Self-Supervised Features. Powerful correspondences emerge in the feature space of large pre-trained vision foundation models, as reviewed in Sec. 2. We use features extracted from DINOv2 [56], Stable Diffusion [66, 91], and fused features of the two [91] as baselines. We use feature-level similarity between pixel pairs to find correspondences.

Self-Supervised Features and VLM Grounding. Since our task requires both semantic and structural reasoning based on the function prompt, these baselines chain a VLM that grounds functional parts with a pre-trained model that provides structural priors. Given an image pair, we use functional part bounding boxes generated by the VLM for each image, and then use self-supervised features to find correspondences within these part labels. This approach can benefit both label transfer and discovery because the functional part prediction adds a constraint on the space of possible matches, making it easier to find accurate matches. We consider two VLMs as the functional part grounding modules to be combined with off-the-shelf DINOv2 features:

 CogVLM [82], which outputs bounding boxes based on prompts of the functional part.

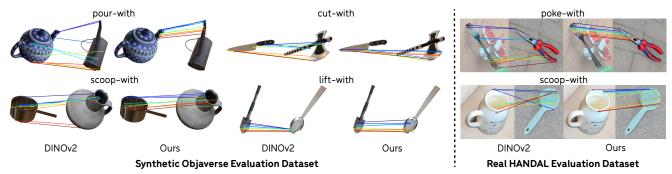


Figure 6. **Correspondence Discovery Comparisons.** We observe that our approach more reliably retrieves the functionally relevant correspondences than off-the-shelf DINOv2. The top 10 highest-ranked matches are shown.

 ManipVQA [25], an affordance-grounding model that outputs bounding boxes conditioned on actions. We use the 7B model in our experiments. We also prompt ManipVQA in two ways, one with the functional part name and the other with the function itself because the model is finetuned for robotic tasks. We refer to these as ManipVQA-P and ManipVQA-F, respectively.

5.3. Quantitative Comparisons

Results in Table 1 evaluate the performance of our method and baseline solutions on the synthetic and real evaluation datasets introduced in Sec. 3.2. Results show that our model trained on fully synthetic data can generalize to real images.

Compared to baseline solutions that solely use self-supervised features, our full model – trained with both functional and spatial contrastive loss – consistently outperforms. These metrics demonstrate that the pseudolabel quality is sufficient for learning meaningful functional correspondences. Additionally, given that the evaluation dataset predominantly includes cross-category pairs, Table 1 illustrates that self-supervised features struggle with cross-category generalization. Further evidence is provided in the Supplement, where we present a detailed breakdown of metrics for both within- and across-category pairs.

Compared to baselines using VLM grounding, even with CogVLM bounding boxes as additional functional part information, off-the-shelf DINOv2 features underperform relative to our full model. The margin is generally smaller, which highlights the importance of understanding the context of the function. On the other hand, ManipVQA outputs less accurate bounding boxes, which is reflected in the metrics. In particular, prompting with the part instead of the function is significantly better, which shows the difficulty of zero-shot affordance grounding given a function name. Note also that running CogVLM inference is roughly 50 times slower than our model and running ManipVQA inference is roughly 1000 times slower than our model.

Ablations. We ablate the role of the functional and spatial contrastive loss in Table 1. The model trained solely with functional loss performs poorly in both label transfer and

correspondence discovery. The model trained solely with spatial loss is better but still falls short compared to the full model due to its lack of functional information. Finally, models with and without mask loss share similar performances. The model with mask loss does outperform the model without it in all metrics for correspondence discovery on the real evaluation dataset, which represents the least constrained and most realistic case. This optional mask prediction module can learn functional part masks with minimal additional cost.

5.4. Qualitative Results

We present results for correspondence label transfer in Figure 1 and correspondence discovery in Figure 6. Our model predictions not only capture object parts specific to the input function, but also preserve the structural relation among parts. Figure 6 shows top 10 matches according to the score from Section 5.1 separated by 5 pixels each. DINOv2 features are not function-aware and result in inaccurate matching, especially in cross-category image pairs. In comparison, our model produces dense matches between functional parts from different object categories with high spatial precision, *e.g.*, matching the rim of a saucepan with the rim of a jug. Overall, our model demonstrates a deep understanding of functional and structural information of objects, which produces high-quality dense functional correspondences.

6. Conclusion

We have introduced the problem of dense functional correspondence, where input images contain objects with similar functionality but possibly come from distinct object categories. We have proposed a principled approach to obtain dense 2D functional correspondences from 3D object alignments and curated datasets for comprehensive evaluations. To tackle the task, we have presented a weakly-supervised framework that distills semantic information from vision-language models, while learning structural information through tuning self-supervised features with a multi-view contrastive loss. Our model outperforms a set of baselines in both synthetic and real-world benchmarks.

Acknowledgments. This work is in part supported by NSF RI #2211258 and #2338203, NSF CAREER #1844724, NSF NCS-FR #2123963, ONR YIP N00014-24-1-2117, ONR MURI N00014-22-1-2740, ONR N00014-20-1-2589, ONR MURI N00014-21-1-2801, ONR MURI N00014-24-1-2748, and Simons Foundation grant SFI-AN-NC-GB-Culmination-00002986-05. We also thank the Stanford HAI for their support with computing resources.

References

- [1] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. *arXiv preprint arXiv:2112.05814*, 2021. 3
- [2] Samarth Brahmbhatt, Cusuh Ham, Charles C Kemp, and James Hays. Contactdb: Analyzing and predicting grasp contact via thermal imaging. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8709–8719, 2019. 2
- [3] 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. 3
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 5
- [5] Ching-Yao Chuang, Jiaman Li, Antonio Torralba, and Sanja Fidler. Learning to act properly: Predicting and explaining affordances from images. In *Proceedings of the IEEE Con*ference on Computer Vision and Pattern Recognition, pages 975–983, 2018. 2
- [6] Blender Online Community. Blender a 3D modelling and rendering package. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 4
- [7] 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. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13142–13153, 2023. 2, 4, 13, 16
- [8] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. arXiv preprint arXiv:2409.17146, 2024.
- [9] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. Advances in Neural Information Processing Systems, 36, 2024. 5
- [10] Shengheng Deng, Xun Xu, Chaozheng Wu, Ke Chen, and Kui Jia. 3d affordancenet: A benchmark for visual object affordance understanding. In proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 1778–1787, 2021. 2

- [11] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- [12] Norman Di Palo and Edward Johns. On the effectiveness of retrieval, alignment, and replay in manipulation. *IEEE Robotics and Automation Letters*, 2024. 3
- [13] Kairui Ding, Boyuan Chen, Ruihai Wu, Yuyang Li, Zongzheng Zhang, Huan-ang Gao, Siqi Li, Guyue Zhou, Yixin Zhu, Hao Dong, et al. Preafford: Universal affordance-based pre-grasping for diverse objects and environments. arXiv preprint arXiv:2404.03634, 2024. 3
- [14] Thanh-Toan Do, Anh Nguyen, and Ian Reid. Affordancenet: An end-to-end deep learning approach for object affordance detection. In 2018 IEEE international conference on robotics and automation (ICRA), pages 5882–5889. IEEE, 2018. 2
- [15] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 5, 14, 21
- [16] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 11444–11453, 2020. 2
- [17] Peter R Florence, Lucas Manuelli, and Russ Tedrake. Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. In *Conference on Robot Learning*, pages 373–385. PMLR, 2018. 1, 2, 6
- [18] James J Gibson. The ecological approach to visual perception: Classic edition. *Taylor & Francis*, 1979.
- [19] Yuchao Gu, Yipin Zhou, Bichen Wu, Licheng Yu, Jia-Wei Liu, Rui Zhao, Jay Zhangjie Wu, David Junhao Zhang, Mike Zheng Shou, and Kevin Tang. Videoswap: Customized video subject swapping with interactive semantic point correspondence. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 7621– 7630, 2024. 1
- [20] Andrew Guo, Bowen Wen, Jianhe Yuan, Jonathan Tremblay, Stephen Tyree, Jeffrey Smith, and Stan Birchfield. Handal: A dataset of real-world manipulable object categories with pose annotations, affordances, and reconstructions. In 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 11428–11435. IEEE, 2023. 2, 4, 16
- [21] HDRI Haven. Hdri haven. https://hdri-haven.com, 2024. Accessed: 2024-9. 4
- [22] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 5
- [23] Nick Heppert, Max Argus, Tim Welschehold, Thomas Brox, and Abhinav Valada. Ditto: Demonstration imitation by trajectory transformation. In 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2024. 3

- [24] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- [25] Siyuan Huang, Iaroslav Ponomarenko, Zhengkai Jiang, Xiaoqi Li, Xiaobin Hu, Peng Gao, Hongsheng Li, and Hao Dong. Manipvqa: Injecting robotic affordance and physically grounded information into multi-modal large language models. arXiv preprint arXiv:2403.11289, 2024. 1, 3, 7, 8, 15, 16, 17
- [26] Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Flowformer: A transformer architecture for optical flow. In *European conference on computer vision*, pages 668–685. Springer, 2022. 2
- [27] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Pro*ceedings of the IEEE conference on computer vision and pattern recognition, pages 2462–2470, 2017. 2
- [28] Wei Jiang, Eduard Trulls, Jan Hosang, Andrea Tagliasacchi, and Kwang Moo Yi. Cotr: Correspondence transformer for matching across images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6207–6217, 2021. 2
- [29] Zhenyu Jiang, Hanwen Jiang, and Yuke Zhu. Doduo: Learning dense visual correspondence from unsupervised semantic-aware flow. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 12420– 12427. IEEE, 2024. 1
- [30] Yuanchen Ju, Kaizhe Hu, Guowei Zhang, Gu Zhang, Mingrun Jiang, and Huazhe Xu. Robo-abc: Affordance generalization beyond categories via semantic correspondence for robot manipulation. *arXiv preprint arXiv:2401.07487*, 2024.
- [31] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36, 2024. 5
- [32] Deborah G Kemler Nelson, Rachel Russell, Nell Duke, and Kate Jones. Two-year-olds will name artifacts by their functions. *Child development*, 71(5):1271–1288, 2000. 3
- [33] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 6
- [34] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Con*ference on Computer Vision, pages 4015–4026, 2023. 5
- [35] Akshay Krishnan, Abhijit Kundu, Kevis-Kokitsi Maninis, James Hays, and Matthew Brown. Omninocs: A unified nocs dataset and model for 3d lifting of 2d objects. *arXiv preprint arXiv:2407.08711*, 2024. 1, 2
- [36] Yuxuan Kuang, Junjie Ye, Haoran Geng, Jiageng Mao, Congyue Deng, Leonidas Guibas, He Wang, and Yue Wang. Ram: Retrieval-based affordance transfer for generalizable zero-shot robotic manipulation. *arXiv preprint arXiv:2407.04689*, 2024. 3

- [37] Hamid Laga, Michela Mortara, and Michela Spagnuolo. Geometry and context for semantic correspondences and functionality recognition in man-made 3d shapes. ACM Transactions on Graphics (TOG), 32(5):1–16, 2013. 3
- [38] Zihang Lai, Senthil Purushwalkam, and Abhinav Gupta. The functional correspondence problem. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15772–15781, 2021. 1, 2, 3, 4, 13, 18, 19
- [39] Barbara Landau, Linda Smith, and Susan Jones. Object shape, object function, and object name. *Journal of mem*ory and language, 38(1):1–27, 1998. 3
- [40] Yixing Lao, Xiaogang Xu, Xihui Liu, Hengshuang Zhao, et al. Corresnerf: Image correspondence priors for neural radiance fields. Advances in Neural Information Processing Systems, 36:40504–40520, 2023. 1
- [41] Gen Li, Varun Jampani, Deqing Sun, and Laura Sevilla-Lara. Locate: Localize and transfer object parts for weakly supervised affordance grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10922–10931, 2023. 2
- [42] Gen Li, Nikolaos Tsagkas, Jifei Song, Ruaridh Mon-Williams, Sethu Vijayakumar, Kun Shao, and Laura Sevilla-Lara. Learning precise affordances from egocentric videos for robotic manipulation. arXiv preprint arXiv:2408.10123, 2024. 2
- [43] Yuanqi Li, Shun Liu, Xinran Yang, Jianwei Guo, Jie Guo, and Yanwen Guo. Surface and edge detection for primitive fitting of point clouds. In ACM SIGGRAPH 2023 conference proceedings, pages 1–10, 2023. 14
- [44] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv* preprint arXiv:2311.07575, 2023. 3, 5
- [45] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-perfect structure-frommotion with featuremetric refinement. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5987–5997, 2021. 1
- [46] Hongchen Luo, Wei Zhai, Jing Zhang, Yang Cao, and Dacheng Tao. Learning affordance grounding from exocentric images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2252–2261, 2022. 2
- [47] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. *Advances in Neural Information Processing Systems*, 36, 2024. 5, 13, 21
- [48] Jinjie Mai, Wenxuan Zhu, Sara Rojas, Jesus Zarzar, Abdullah Hamdi, Guocheng Qian, Bing Li, Silvio Giancola, and Bernard Ghanem. Tracknerf: Bundle adjusting nerf from sparse and noisy views via feature tracks. *arXiv preprint arXiv:2408.10739*, 2024. 1
- [49] Juhong Min, Jongmin Lee, Jean Ponce, and Minsu Cho. Spair-71k: A large-scale benchmark for semantic correspondence. *arXiv preprint arXiv:1908.10543*, 2019. 2, 6

- [50] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspnet: Variational grasp generation for object manipulation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2901–2910, 2019. 2
- [51] Adithyavairavan Murali, Weiyu Liu, Kenneth Marino, Sonia Chernova, and Abhinav Gupta. Same object, different grasps: Data and semantic knowledge for task-oriented grasping. In *Conference on robot learning*, pages 1540–1557. PMLR, 2021. 2, 4, 13
- [52] Austin Myers, Ching L Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection of tool parts from geometric features. In 2015 IEEE International Conference on Robotics and Automation (ICRA), pages 1374–1381. IEEE, 2015. 2
- [53] Tushar Nagarajan, Christoph Feichtenhofer, and Kristen Grauman. Grounded human-object interaction hotspots from video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8688–8697, 2019. 2
- [54] Anh Nguyen, Dimitrios Kanoulas, Darwin G. Caldwell, and Nikos G. Tsagarakis. Object-based affordances detection with convolutional neural networks and dense conditional random fields. In 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 5908–5915, 2017. 2
- [55] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen
- Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob Mc-Grew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. 4, 13, 21
- [56] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193, 2023. 1, 2, 3, 5, 7, 15, 16, 18
- [57] Nobuyuki Otsu et al. A threshold selection method from gray-level histograms. Automatica, 11(285-296):23–27, 1975. 14
- [58] Norman Di Palo and Edward Johns. Dinobot: Robot manipulation via retrieval and alignment with vision foundation models. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024. 1, 3
- [59] Georgios Papagiannis and Edward Johns. Miles: Making

- imitation learning easy with self-supervision. In *Proceedings* of the Conference on Robot Learning (CoRL), 2024. 3
- [60] Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer, 2017. 14
- [61] Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhuowen Tu, and Li Erran Li. Affordancellm: Grounding affordance from vision language models. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7587–7597, 2024. 3
- [62] Zengyi Qin, Kuan Fang, Yuke Zhu, Li Fei-Fei, and Silvio Savarese. Keto: Learning keypoint representations for tool manipulation. In 2020 IEEE International Conference on Robotics and Automation (ICRA), pages 7278–7285. IEEE, 2020. 2
- [63] 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. 3, 5
- [64] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learn*ing Research, 21(140):1–67, 2020. 3
- [65] 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. 4
- [66] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of* the IEEE/CVF conference on computer vision and pattern recognition, pages 10684–10695, 2022. 1, 3, 7
- [67] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of* the IEEE/CVF conference on computer vision and pattern recognition, pages 4938–4947, 2020. 2
- [68] Stefan Stojanov, Anh Thai, Zixuan Huang, and James M Rehg. Learning dense object descriptors from multiple views for low-shot category generalization. Advances in Neural Information Processing Systems, 35:12566–12580, 2022. 2
- [69] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF con*ference on computer vision and pattern recognition, pages 8922–8931, 2021. 2
- [70] Yixuan Sun, Yiwen Huang, Haijing Guo, Yuzhou Zhao, Runmin Wu, Yizhou Yu, Weifeng Ge, and Wenqiang Zhang. Misc210k: A large-scale dataset for multi-instance semantic correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7121–7130, 2023. 2, 6

- [71] Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent correspondence from image diffusion. Advances in Neural Information Processing Systems, 36:1363–1389, 2023. 3
- [72] Tatsunori Taniai, Sudipta N Sinha, and Yoichi Sato. Joint recovery of dense correspondence and cosegmentation in two images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4246–4255, 2016. 1
- [73] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23– 28, 2020, Proceedings, Part II 16, pages 402–419. Springer, 2020. 2
- [74] Skye Thompson, Leslie Pack Kaelbling, and Tomas Lozano-Perez. Shape-based transfer of generic skills. In 2021 IEEE International Conference on Robotics and Automation (ICRA), pages 5996–6002. IEEE, 2021. 3
- [75] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 3
- [76] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. Sparf: Neural radiance fields from sparse and noisy poses. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pages 4190–4200, 2023. 1
- [77] Dylan Turpin, Liquan Wang, Stavros Tsogkas, Sven Dickinson, and Animesh Garg. Gift: Generalizable interaction-aware functional tool affordances without labels. arXiv preprint arXiv:2106.14973, 2021. 2
- [78] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019.
- [79] Raviteja Vemulapalli, Hadi Pouransari, Fartash Faghri, Sachin Mehta, Mehrdad Farajtabar, Mohammad Rastegari, and Oncel Tuzel. Knowledge transfer from vision foundation models for efficient training of small task-specific models. In Forty-first International Conference on Machine Learning, 2024. 5
- [80] Boyan Wan, Yifei Shi, and Kai Xu. Socs: Semantically-aware object coordinate space for category-level 6d object pose estimation under large shape variations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14065–14074, 2023. 2
- [81] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 2642– 2651, 2019. 1, 2
- [82] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. 1, 2, 3, 4, 5, 7, 14, 15, 16, 17, 18
- [83] Elizabeth A Ware and Amy E Booth. Form follows function:

- Learning about function helps children learn about shape. *Cognitive Development*, 25(2):124–137, 2010. 3
- [84] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Liang Pan Jiawei Ren, Wayne Wu, Lei Yang, Jiaqi Wang, Chen Qian, Dahua Lin, and Ziwei Liu. Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 4
- [85] Ruinian Xu, Fu-Jen Chu, Chao Tang, Weiyu Liu, and Patricio A Vela. An affordance keypoint detection network for robot manipulation. *IEEE Robotics and Automation Letters*, 6(2):2870–2877, 2021. 2
- [86] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10371–10381, 2024. 5
- [87] Yuhang Yang, Wei Zhai, Hongchen Luo, Yang Cao, Jiebo Luo, and Zheng-Jun Zha. Grounding 3d object affordance from 2d interactions in images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10905–10915, 2023. 2
- [88] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimgnet: A large-scale dataset of multi-view images. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 9150–9161, 2023. 4
- [89] Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan Murali, Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial affordance prediction for robotics. *arXiv preprint arXiv:2406.10721*, 2024. 3
- [90] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference* on Computer Vision, pages 11975–11986, 2023. 3
- [91] Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence, 2023. 1, 7, 15
- [92] Chenchen Zhu, Fanyi Xiao, Andrés Alvarado, Yasmine Babaei, Jiabo Hu, Hichem El-Mohri, Sean Culatana, Roshan Sumbaly, and Zhicheng Yan. Egoobjects: A large-scale egocentric dataset for fine-grained object understanding. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 20110–20120, 2023. 4
- [93] Yixin Zhu, Yibiao Zhao, and Song Chun Zhu. Understanding tools: Task-oriented object modeling, learning and recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2855–2864, 2015. 3
- [94] Yifeng Zhu, Arisrei Lim, Peter Stone, and Yuke Zhu. Vision-based manipulation from single human video with open-world object graphs. arXiv preprint arXiv:2405.20321, 2024. 3

Appendix

7. Training Data Generation

7.1. Function and Object Taxonomy

List of functions. To obtain our taxonomy of functions, we first take the function lists defined by [38, 51] and ask GPT-4 [55] to expand them. Our prompt is a simple "Given this list of functions, generate more options for object functions." We manually process this list by simplifying synonymous functions into the most generic function to reduce redundancy, *e.g.*, "slice-with" and "chop-with" get absorbed into "cut-with," or "skewer-with" and "bore-with" get absorbed into "pierce-with." The final list of 24 functions is shown in Table 7.

Finding object categories given a function. To generate a list of object categories suitable for the chosen functions, we use the prompting strategy shown in Figure 13. We combine the common and uncommon lists and remove object names that are synonyms or that would require significant improvisation to achieve a certain function well.

7.2. Functional Part Description Generation

Using LLMs, we have created a list of functions and a list of object categories that can carry out those functions. Given an object category and a function, we now require a means to generate part names and descriptions to prompt the grounded VLM. To obtain a list of functional part names, we use the prompting strategy shown in Figure 14. This produces a list of parts for each (object category, function) pair, which we manually filter based on the most precise part. For instance, if GPT generates "blade" and "point" for (knife, pierce-with), we will choose "point." Querying different functions for the same object may result in the same functional part description being output multiple times with small variability. To combine these descriptions, we simply prompt GPT-4 to summarize them into one.

7.3. Objaverse Dataset Filtering

The Objaverse [7] dataset does not come with high-quality labels, making it challenging to use as a training dataset for tasks that require semantic object understanding. There is the Objaverse-LVIS split, but it is a small subset of the complete Objaverse, and the labels are noisy. To address this, Caption3D [47] proposes a technique for generating captions for $\approx 600K$ of the assets in Objaverse based on a combination of VLMs and LLMs. However, these captions are still insufficient for our purpose because they do not contain explicit category labels.

For each caption from Caption3D, we propose to filter it by comparing it with our list of object categories from Section 7.1. However, doing this naively using a large language

model like Llama [15] would require about 100M model inferences, making this intractable. To resolve this and make the procedure more efficient, we propose summarizing each caption into a single noun using Llama3 with the prompting strategy described in Figure 15. After converting the list of captions into a list of nouns, we use Llama3 word embedding distance to determine whether the noun belongs to the list of categories we generated in Section 7.1. Last, we ask Llama to verify the matches from word embedding as a final pass.

7.4. CogVLM Prompting and Aggregation

We use the descriptions generated in Section 7.2 to prompt the cogvlm-grounding-generalist-v1.1 variant of CogVLM [82], which has been tuned for referring expression comprehension. Specifically, given a prompt like "What are the exact bounding boxes of <expr> in the provided picture?", where <expr> can be a noun or a descriptive phrase, the model is tuned to produce a text sequence describing a bounding box. Because of the sampling inherent to language transformer model inference, the bounding boxes vary across trials. Our procedure to label functional parts using CogVLM outputs consists of the following steps:

- 1. Render 19 views per object that shows it from various angles, including from above and below.
- For each functional part description and each view, query CogVLM for four trials to obtain the bounding box pseudo-labels. For small parts like points or tips, we do a second iteration that zooms into the initial bounding boxes to improve precision.
- 3. Aggregate all trials and views onto a point cloud of 100K randomly sampled points on the object's surface. Every time a given point in the point cloud gets labeled by a bounding box in a different view, we increment its score. The final numbers are normalized to be in 0-1. For prompts that specifically ask for the labeling of edges, we multiply the point cloud with the per-point edge probabilities from SED-net [43], a method for decomposing point clouds into primitives.
- 4. Given this point cloud, for any rendered image of the object, we can project the point cloud into 2D and produce a binary mask with Otsu's method [57] and a series of binary dilation/erosion steps to close holes in the mask. Example outputs of this procedure are shown in Figure 7.

8. Additional Training and Evaluation Details

8.1. Ground-Truth Generation

In this section, we provide additional details for deriving ground-truth 2D dense correspondences from 3D alignment. Given two object meshes that can perform the same function, we obtain their 3D functional alignment and the

3D bounding boxes for the functional parts using the procedure in Section 3.2. Given rendered images I_1, I_2 of the two assets, we first find 2D pixels P_1, P_2 that would back-project to 3D points within the labeled functional part bounding boxes. The set of pixels P_1 and P_2 represent the functional part segmentation on the two images. Then, we perform minimum cost matching where the cost between two pixels $p_1 \in P_1$ and $p_2 \in P_2$ is measured by the distance between their back-projected 3D locations. In particular, we use the Hungarian algorithm. However, since the Hungarian algorithm requires one-to-one matches, we subsample the set between P_1 and P_2 that has more pixels using furthest point sampling. The output of the Hungarian algorithm constitutes the ground-truth 2D dense functional correspondences.

Practically, we randomly sample rendered images from the top 5 out of 30 views where the functional part is most visible. We do so for six trials and repeat the procedure above to obtain 2D ground-truth annotations for the six view pairs for each pair of assets. Among these trials, ambiguity in the correspondence definition may arise due to 3D symmetries. We disambiguate this based on the objects' orientation when projected in the 2D images. For instance, for two rims in 2D, the top (in 2D, relative to the sides of the image) of one rim should align with the top of the other rim. We believe this is appropriate as it is the first investigation of this problem setting. In future work, we aim to refine the task and model to capture such ambiguity. As such, we manually filter the derived 2D annotations based on the ground-truth dense visualizations to disambiguate and ensure high quality.

8.2. Additional Technical Details

Model training. For training our full models, we found that sampling points solely on the functional part for the spatial contrastive loss helped performance. However, when training the model with spatial loss only, we found that sampling points on the whole object helped more.

Feature representation complexity. We experimented with LoRA [24] finetuning of DINOv2 and FiLM [60] layers for text conditioning. Despite the increased training cost of LoRA, we did not observe consistent improvements (e.g. normalized distance increased from 0.172 to 0.181).

Evaluation details. Metrics are computed at fixed pixels because the input images are center-cropped (all objects have similar sizes), making it equivalent to normalizing with respect to a percentage of bounding box sizes as in prior work. For the SD-DINO baseline, we follow their standard resolutions and scale the input images accordingly.

Evaluation with predicted functional part masks. Below, we explain the evaluation protocol for models that involve a functional part mask prediction (*e.g.*, CogVLM [82]



Figure 7. Examples of pseudo-labeled functional parts in point clouds and images using CogVLM [82]. Using the procedure outlined in Section 7.4, we pseudo-label images with masks for the object functional parts. Notably, this pipeline has the ability to generate part labels for non-convex object parts, such as a mug's rim, and for parts that lack clear edge boundaries, such as a teapot's spout. Point clouds are shown in views that best capture the aggregated functional part labels.

Model	Correspon	dence Label Tr	ansfer	Correspondence Discovery				
1170001	Normalized Dist (↓)	PCK@23p (†)	PCK@10p (†)	Best F1@23p (†)	Best F1@10p (†)	AP@23p (†)	AP@10p(†)	
	within / across	within / across	within / across	within / across	within / across	within / across	within / across	
Synthetic Evaluation Dataset								
Chance	0.317 / 0.309	0.162 / 0.166	0.047 / 0.046	0.382 / 0.421	0.163 / 0.178	0.234 / 0.260	0.085 / 0.095	
DINO [56]	0.132 / 0.225	0.589 / 0.347	0.283 / 0.126	0.708 / 0.557	0.425 / 0.257	0.555 / 0.352	0.265 / 0.108	
SD [91]	0.221 / 0.275	0.423 / 0.278	0.210 / 0.112	0.528 / 0.471	0.295 / 0.220	0.322 / 0.258	0.153 / 0.087	
SD-DINO [91]	0.154 / 0.240	0.553 / 0.347	0.284 / 0.141	0.642 / 0.550	0.406 / 0.284	0.443 / 0.324	0.239 / 0.129	
CogVLM [82] + DINO	0.126 / 0.188	0.596 / 0.387	0.281 / 0.138	0.840 / 0.651	0.519 / 0.303	0.749 / 0.525	0.362 / 0.160	
CogVLM [82] + SD-DINO	0.135 / 0.188	0.578 / 0.404	0.292 / 0.161	0.825 / 0.683	0.554 / 0.368	0.717 / 0.551	0.400 / 0.216	
ManipVQA-P [25] + DINO	0.181 / 0.230	0.493 / 0.323	0.232 / 0.113	0.737 / 0.548	0.437 / 0.242	0.608 / 0.387	0.284 / 0.110	
ManipVQA-F [25] + DINO	0.234 / 0.278	0.352 / 0.244	0.152 / 0.084	0.650 / 0.508	0.374 / 0.222	0.444 / 0.300	0.193 / 0.081	
Ours (functional only)	0.187 / 0.235	0.412 / 0.266	0.154 / 0.084	0.723 / 0.551	0.358 / 0.212	0.617 / 0.412	0.220 / 0.094	
Ours (spatial only)	0.128 / 0.217	<u>0.674</u> / 0.436	0.385 / <u>0.201</u>	0.686 / 0.597	0.469 / 0.353	0.493 / 0.398	0.295 / 0.198	
Ours (full without mask loss)	0.112 / 0.180	0.680 / 0.454	<u>0.377</u> / 0.203	<u>0.878</u> / <u>0.750</u>	0.643 / 0.442	0.823 / 0.662	0.537 / 0.306	
Ours (full with mask loss)	<u>0.122</u> / <u>0.180</u>	0.655 / <u>0.451</u>	0.367 / 0.199	0.879 / 0.757	0.645 / 0.443	<u>0.820</u> / <u>0.661</u>	<u>0.528</u> / <u>0.297</u>	
Real Evaluation Dataset								
Chance	0.311 / 0.313	0.170 / 0.170	0.044 / 0.046	0.431 / 0.413	0.171 / 0.165	0.262 / 0.243	0.090 / 0.086	
DINO [56]	0.130 / 0.230	0.570 / 0.356	0.252 / 0.129	0.734 / 0.542	0.434 / 0.250	0.577 / 0.320	0.275 / 0.095	
SD [91]	0.204 / 0.277	0.411 / 0.276	0.192 / 0.106	0.587 / 0.477	0.308 / 0.215	0.355 / 0.263	0.148 / 0.086	
SD-DINO [91]	0.151 / 0.243	0.514 / 0.344	0.244 / 0.137	0.679 / 0.544	0.400 / 0.270	0.468 / 0.303	0.224 / 0.116	
CogVLM [82] + DINO	0.142 / 0.182	0.544 / 0.407	0.239 / 0.147	0.782 / 0.667	0.465 / 0.314	0.686 / 0.521	0.312 / 0.161	
CogVLM [82] + SD-DINO	0.154 / 0.186	0.506 / 0.402	0.234 / 0.158	0.762 / 0.683	0.462 / 0.360	0.618 / 0.540	0.295 / 0.219	
ManipVQA-P [25] + DINO	0.148 / 0.222	0.534 / 0.354	0.234 / 0.127	0.719 / 0.563	0.415 / 0.256	0.577 / 0.370	0.260 / 0.112	
ManipVQA-F [25] + DINO	0.236 / 0.263	0.405 / 0.279	0.174 / 0.095	0.714 / 0.531	0.412 / 0.239	0.509 / 0.323	0.231 / 0.093	
Ours (functional only)	0.179 / 0.206	0.405 / 0.313	0.152 / 0.103	0.730 / 0.627	0.356 / 0.260	0.599 / 0.511	0.199 / 0.132	
Ours (spatial only)	<u>0.129</u> / 0.227	0.631 / 0.421	0.343 / 0.192	0.708 / 0.564	0.470 / 0.316	0.501 / 0.344	0.295 / 0.145	
Ours (full without mask loss)	0.122 / <u>0.161</u>	0.639 / 0.477	0.352 / 0.216	<u>0.835</u> / <u>0.756</u>	0.589 / 0.441	<u>0.741</u> / <u>0.675</u>	<u>0.469</u> / <u>0.304</u>	
Ours (full with mask loss)	0.132 / 0.160	0.603 / <u>0.469</u>	0.321 / <u>0.208</u>	0.857 / 0.792	0.611 / 0.467	0.773 / 0.716	0.485 / 0.321	

Table 2. **Quantitative evaluation by within- and across-category pairs.** We further break down Table 1 by within- and across-category performance for all the metrics. Additional result for CogVLM + SD-DINO is also included. Off-the-shelf self-supervised features tend to perform worse at cross-category generalization compared to our full model.

+ DINO, ManipVQA [25] + DINO, and our full model with mask loss). In label transfer, for each pixel p_1^i on image I_1 , we restrict its most similar match $p_2^{j(i)}$ on I_2 to be within the predicted functional part mask of I_2 . In correspondence

discovery, predicted part masks are produced for both images. We restrict matches to only happen between the two predicted masks and between their complements. Matches that fall within the two predicted part masks are prioritized

in the ranking explained in Section 5.1.

Dense correspondence visualization. The dense label transfer visualizations use the ground-truth mask for the source image but the predicted mask for the target image. For each pixel on the target image's functional part mask, we find its most similar match in the source image's functional part mask to produce the label transfer color map.

8.3. Computational Costs

Rendering multi-view images on selected Objaverse [7] assets takes one day with four 2080 Ti GPUs. Functional part pseudo-labeling takes one week on eight A6000 GPUs, as CogVLM [82] inference is slow and memory-intensive. We emphasize that rendering and pseudo-labeling are only done once and scale significantly better than human annotation. Our model can be trained on a single NVIDIA GeForce RTX 3090 in ≈ 2 days for 100 epochs. These computational demands are fairly standard and are justified by the capability to trade off compute for expensive and time-consuming human annotation.

9. Additional Quantitative Results

9.1. Within- and Cross-Category Comparison

Since the evaluation dataset contains both within-category pairs and across-category pairs, we further separate the metrics in Table 1 into within-category results and acrosscategory results in Table 2. In general, all the models and baselines perform better on within-category cases than on across-category cases. This illustrates the inherent difficulty of cross-category generalization. In addition, the performance margin between off-the-shelf self-supervised features and our model is often larger on the across-category pairs. On average, DINOv2 performs 46.3% worse on cross-category pairs, while ours is 33.3% worse. This serves as evidence that off-the-shelf self-supervised features struggle more with cross-category generalization. Last, without any functional part information, our spatial-only model performs competitively on within-category pairs on label transfer metrics but is worse on across-category pairs.

9.2. Scaling Experiments

In this section, we show scaling experiments where we replace the backbone in our full model with mask loss with DINOv2 [56] of different ViT sizes. The results are shown in Table 3. As the ViT size increases, we generally observe an improvement in the evaluation metrics. In addition, when we reduce the stride size from 14 pixels to 7, the model performance also improves, especially in correspondence discovery. This demonstrates that both higher spatial resolution and higher backbone capacity can improve the performance of our approach.

Note that due to computational resource constraints, DI-NOv2 with ViT-G was only trained for 30 epochs, and ViT-B with half stride was trained for 80 epochs, while other models were trained for 100 epochs. Compared to ViT-B, using ViT-S is ≈ 1.6 times faster, using half stride is ≈ 2.6 times slower, using ViT-L is ≈ 2.1 times slower, and using ViT-G is ≈ 5.8 times slower.

9.3. Sensitivity Analysis of Loss Weights

We further ablate the spatial and mask loss weights in Table 4. Varying $\lambda_{\rm spatial}$ has an effect, but the model does not appear to be highly sensitive, making it easy to converge on $\lambda_{\rm spatial}=10$ to achieve the best result. On the other hand, we observe low variance when increasing $\lambda_{\rm mask}$. The benefits of the mask loss are illustrated in Figure 9.

9.4. Functional Part Prediction Accuracy

Some of the methods we evaluate generate functional part segmentation predictions. Accordingly, we compare their segmentation accuracies in Table 5. Specifically, ManipVQA-P and ManipVQA-F [25] refer to segmentations produced by ManipVQA using part label prompts and function name prompts, respectively. For CogVLM [82] on 2D images, predictions are generated from single-image inputs into CogVLM, aggregated across four trials via K-Means clustering. These three methods produce bounding boxes, which are further multiplied with the object masks. CogVLM [82] with 3D aggregation follows the pipeline illustrated in Figure 3. Since our full model with mask loss incorporates a functional part mask prediction module, we also evaluate its segmentation performance as part of this comparison.

To evaluate these methods, we use ground-truth part masks generated by our evaluation pipeline on both the synthetic Objaverse [7] data and the real HANDAL [20] data. Specifically, for each (object, function) pair, we label a 3D bounding box, and any pixel that projects to a 3D point within this bounding box is classified as part of the functional region. As shown in Table 5, both CogVLM methods and our learned model have good accuracy. Note that the pseudo-labeling pipeline can produce very fine-grained parts like small tips or edges that do not necessarily align with the human annotations. As such, the main advantage of the 3D aggregation pipeline is illustrated in Figure 7. In addition, on the real HANDAL data, our model's predictions perform better than the CogVLM model, which has state-of-the-art referring expression detection capabilities.

9.5. Ranking Scheme.

We designed our feature similarity and cycle consistencybased ranking scheme in Section 5.1 to enable strong performance across all methods. To show its impact, we include results from a simpler version using only feature similarity

Model	Correspond	lence Label Tr	ansfer	Correspondence Discovery				
Model	Normalized Dist (↓)	PCK@23p (†)	PCK@10p (†)	Best F1@23p (†)	Best F1@10p (†)	AP@23p (†)	AP@10p (†)	
Synthetic Evaluation Dataset								
DINOv2 ViT-S	0.171	0.476	0.218	0.768	0.466	0.676	0.325	
DINOv2 ViT-B	0.172	0.480	0.223	0.774	0.471	0.684	0.330	
DINOv2 ViT-B w/ half stride	0.166	0.494	0.229	0.799	0.508	0.721	0.373	
DINOv2 ViT-L	0.164	0.493	0.233	0.789	0.490	0.705	0.351	
DINOv2 ViT-G	$\overline{0.161}$	0.505	0.239	0.792	0.498	0.711	0.361	
Real Evaluation Dataset								
DINOv2 ViT-S	0.162	0.494	0.229	0.788	0.481	0.697	0.335	
DINOv2 ViT-B	0.153	0.501	0.235	0.808	0.502	0.730	0.360	
DINOv2 ViT-B w/ half stride	0.150	0.519	0.247	0.821	0.525	0.751	0.403	
DINOv2 ViT-L	0.152	0.515	$\overline{0.244}$	0.809	0.514	0.730	0.377	
DINOv2 ViT-G	0.146	0.523	0.252	0.808	0.507	0.729	0.370	

Table 3. **Quantitative evaluation of our model trained with different backbones.** In general, performance increases when the vision transformer backbone becomes larger or when the stride size is reduced.

Loss Weights	Correspone	dence Label Tr	ansfer	Correspondence Discovery				
2055 Weights	Normalized Dist (↓) PCK@23p (↑)		PCK@10p (†)	Best F1@23p (†)	Best F1@10p (†)	AP@23p (†)	AP@10p (↑)	
Synthetic Evaluation Data	set							
$\lambda_{\text{spatial}} = 1, \lambda_{\text{mask}} = 1$	0.193	0.402	0.161	0.707	0.367	0.601	0.224	
$\lambda_{\text{spatial}} = 5, \lambda_{\text{mask}} = 1$	0.177	0.458	0.207	0.761	0.445	0.664	0.304	
$\lambda_{\text{spatial}} = 10, \lambda_{\text{mask}} = 1$	0.172	0.480	0.223	0.774	0.471	0.684	0.330	
$\lambda_{\text{spatial}} = 10, \lambda_{\text{mask}} = 5$	0.173	0.477	0.222	0.775	0.471	0.685	0.330	
$\lambda_{\mathrm{spatial}}^{\mathrm{r}} = 10, \lambda_{\mathrm{mask}} = 10$	0.170	0.478	0.221	0.778	0.470	0.687	0.329	
Real Evaluation Dataset								
$\lambda_{\text{spatial}} = 1, \lambda_{\text{mask}} = 1$	0.169	0.443	0.175	0.759	0.405	0.671	0.260	
$\lambda_{\text{spatial}} = 5, \lambda_{\text{mask}} = 1$	0.158	0.492	0.223	0.793	0.481	0.713	0.345	
$\lambda_{\text{spatial}} = 10, \lambda_{\text{mask}} = 1$	0.153	0.501	0.235	0.808	0.502	0.730	0.360	
$\lambda_{\text{spatial}} = 10, \lambda_{\text{mask}} = 5$	0.155	0.499	0.232	0.804	0.497	0.729	0.359	
$\lambda_{\mathrm{spatial}} = 10, \lambda_{\mathrm{mask}} = 10$	<u>0.154</u>	0.501	0.231	0.800	0.495	0.719	0.353	

Table 4. **Quantitative evaluation of varying loss weights.** Model performance improves with increasing spatial loss weight (up to 10) and remains stable with different mask loss weights.

Method	IoU on Objaverse	IoU on HANDAL
ManipVQA-P [25]	0.453	0.276
ManipVQA-F [25]	0.240	0.146
CogVLM [82] on 2D images	0.656	0.597
CogVLM [82] w/ 3D aggregation	0.635	N/A
Our model prediction	0.628	0.636

Table 5. Quantitative evaluation of functional part segmentation accuracy. This table compares the accuracy of functional part segmentation produced by different methods. Both CogVLM [82] and the predictions of the distilled model demonstrate strong performance in this task. Note that the pipeline described in Figure 3 in the main paper was applied only to the Objaverse dataset; therefore, results for CogVLM [82] with 3D aggregation are omitted for the HANDAL dataset.

in Table 6. The ordering is consistent with the main text, but all methods perform worse. This confirms that all methods benefit from the improved ranking scheme and that our findings are not sensitive to this.

Model	Best F1@23p (†)	Best F1@10p (\uparrow)	AP@23p (↑)	AP@10p (†)
DINO	0.573	0.277	0.376	0.128
CogVLM + DINO	0.672	0.329	0.551	0.184
Ours (full with mask loss)	0.767	0.465	0.679	0.325

Table 6. Correspondence discovery evaluation using only feature similarity. Compared with Table 1, using only feature similarity in the ranking scheme achieves worse performance overall but preserves relative performance among methods.

10. Additional Qualitative Results

Additional dense label transfer results on the synthetic Objaverse dataset, which further validate the effectiveness of our approach, are presented in Figure 8. These results highlight the strong performance of our model in transferring functional part labels across diverse object categories.

More qualitative discovery results on the synthetic Objaverse dataset are shown in Figure 11, and more qualitative discovery results on the real HANDAL dataset are shown in



Figure 8. Additional Label Transfer Dense Visualizations. For each target image (right), our model predicts the functional part mask. To generate the transferred color map, each pixel in the predicted mask is matched to its best corresponding pixel within the ground-truth mask of the source image (left) in terms of feature similarity.

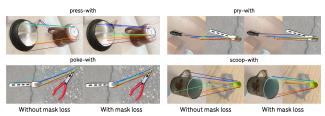


Figure 9. **Qualitative examples for the impact of mask loss.** Functional part predictions can help avoid incorrect matches outside the functional parts in correspondence discovery.

Figure 12. We compare our model with the DINO [56] and CogVLM [82] + DINO baselines. In line with the conclusion in Section 5.4, our model can focus on the functionally relevant part and produce more spatially precise correspondences.

Lastly, we show qualitative evidence for the potential benefits of the optional mask loss in Figure 9. In cases where the predicted functional part mask is accurate, the mask loss can prevent incorrect matches outside functionally relevant regions.

11. Discussion

Differences with FunkPoint [38]. The concept of functional correspondence was previously introduced by [38]. However, our formulation is different in three key aspects.

First, our problem requires dense functional correspondences to be established, whereas [38] defines five sparse keypoints. The manual definition of semantic keypoints

at the function type level is not guaranteed to be well-defined across all object categories. Consequently, we observe inconsistencies and labeling ambiguities in the sparse keypoint annotations. In addition, establishing dense correspondences requires fine-grained and precise reasoning about the structure of object parts, which may make it more useful for downstream applications like transferring demonstrations in robotics.

Second, keypoint matches from [38] include both the object's functional part and where the human interacts with the object. In many cases, like a "bottle" and a "kettle," the functionally irrelevant parts cannot be well aligned. As a result, the key points outside of the functional parts are highly ambiguous. In addition, where an agent interacts with the object depends on the end-effector design, introducing complexity and confounding. On the other hand, our formulation introduces a more precise definition based on the concept of functionally equivalent 3D alignments (discussed in Section 3).

Last, the model proposed by Lai et al. [38] relies on human annotations of sparse keypoints, which inherently limits scalability. In contrast, our approach leverages selfsupervised features and pseudo-labeling, requiring minimal human input, and offers a significantly more scalable solution.

Because of these fundamental differences, our method is not directly applicable to the dataset in Lai et al. [38]. While the feature maps from Lai et al. [38] could be used



Figure 10. **Comparison with [38]** A visual comparison of dense correspondence between [38] (left) and our method (right).

for dense correspondence, the method is not designed for this and it qualitatively appears to be relatively coarse. A visual comparison is provided in Figure 10.

Limitations. A limitation of our work is the existence of ambiguity in some cross-category cases. Ambiguity can arise when an object has multiple parts that can be used for the same function. For instance, both the tip and side rim of a spoon can be used for the function "scrape-with." On the other hand, ambiguity can also arise due to radial symmetry: the rim of a cup and the rim of a bowl can be matched in infinitely many ways. As such, a compelling direction for future work can be developing a probabilistic model to handle the multi-modal nature of the problem and use additional conditioning to resolve such ambiguities.



Figure 11. **Additional Correspondence Discovery Results on Objaverse Evaluation Dataset.** We show more qualitative examples of correspondence discovery on the synthetic Objaverse evaluation dataset, comparing our model against baselines.

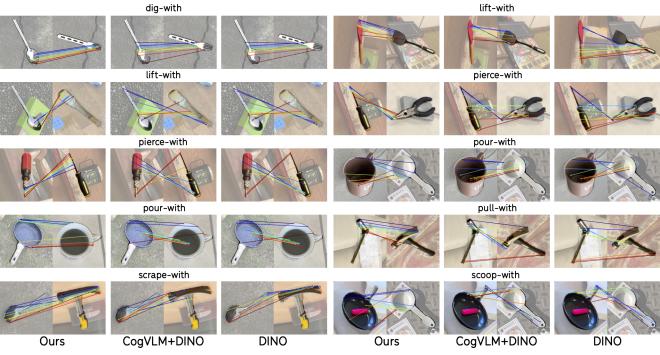


Figure 12. Additional Correspondence Discovery Results on HANDAL Evaluation Dataset.

System Prompt: You are an assistant trying to help a robot figure out what objects can be used to perform certain actions. The objects should be hand-held and unpowered, rigid, unarticulated. The object shape should be what allows it to perform its function. Given an input that is an (action, preposition) pair, output two lists: First output at least 20 object names that are commonly used to perform that action, and at least 20 object names that are not very commonly used to perform an action but the action could be done with some improvisation. Example inputs and outputs. Format your output as a ison Example 1: input: "pour, with" output:{ "common": ["bottle", "cup", "wine glass", "watering can", "teapot", "pitcher", "jug", "saucepan", "measuring cup", "drinking flask"] "uncommon": ["erlenmeyer flask", "lab flask", "bota bag", "bailer shell", "drinking horn"]} Example 2: input: "pound, with" output:{ "common":["hammer", "mallet", "meat tenderizer", "pestle", "sledgehammer", "rubber paver mallet"] "uncommon":["axe", "knife", "boot", "rock", "wrench"]} Input: cut-with Output: { "common": ["knife", "scissors", "box cutter", "razor blade", "chef's knife", "paring knife", "utility knife", "carving knife", "bread knife", "cleaver", "hacksaw", "chisel", "guillotine", "paper cutter", "exacto knife", "shears", "pruning shears", "scalpel", "clipper", "lawn "uncommon": ["glass cutter", "wire cutter", "pizza cutter", "cookie cutter", "plane blade", "credit card", "sharp stone", "broken glass", "metal spatula", "saw blade", "ceramic shard", "guitar string", "fishing line", "shard of mirror", "dental floss", "key", "i skate blade", "shovel edge", "hoe", "axe edge"]} Input:"dig-with" Output: "dig-with": { "common": ["shovel", "trowel", "spade", "pickaxe", "" "" "" "" "hood", "bood", "pickaxe", "hood", "pickaxe", "hood", "pickaxe", "hood", "pickaxe", "hood", "pickaxe", "p "hoe", "mattock", "posthole digger", "scoop", "hand shovel", "garden fork", "edger", "excavator scoop", "dibble", "ice auger", "clam shovel", "snow shovel", "coal shovel", "drain spade", "trenching shovel", "folding shovel"], "uncommon": ["spoon", "stick", "piece of glass", "fork", "knife", "old credit card", "can lid", "metal rod", "plastic cup", "ruler", "pen cap", "toy shovel", "saucepan", "bottle cap", "wooden dowel", "scalpel", "chopstick", "brush handle", "car key", "phone case"]

Figure 13. Our GPT-4 [55] prompting procedure for creating object categories given a list of functions.

System Prompt: Here are your instructions for the rest of the chat: Respond as if you are a human expert giving simplifying instructions to a robot learning to interact with the world by identifying object parts that correspond to verbs. We want to know what area of the object can be used to perform this action. We do not want to know the part that needs to be held or interacted by a human to do this action. Respond with a list of part names, each with a sentence describing the part appearance in ``name - description" format. When answering user questions, carefully consider the following 4 examples. Each example contains a question, a good answer, and a bad answer. The bad answers generally contain parts that the human explicitly interact with. Be sure to avoid bad answers. Question: What are the names of object parts of a "knife" that can be used to perform the action "cut-with"? Respond with only a bulleted list of single word responses paired with short descriptions. Good Answer: - Blade - the flat, sharp part used for cutting. - Edge - The sharpened side of the blade that slices through materials. Bad Answer: - Handle - the part where you grip the knife. **Prompt:** What are the names of object parts of a "dagger" that can be used to perform the action "pierce-with"? Output: - Point - The sharp, tapered end of the dagger used for piercing. - Blade - The flat, sharp part used for slicing or Prompt: What are the names of object parts of a "trowel" that can be used to perform the action "dig-with"? Output:

Figure 14. Overview of our strategy for prompting GPT-4 [55] to obtain functional part names.

- Blade - The flat, pointed part used for digging

- Tip - The sharp end of the blade which helps

into soil.

penetrate the ground.

System Prompt: You will be provided with a brief caption or description of a 3D asset. Your task is to generate the most concise, accurate, and contextually appropriate object name based on the given description. The object name should reflect the core identity of the asset, avoiding overly specific labels. Output only the object Input Caption: "a screwdriver with a blue wooden handle" Prompt: The caption is 'a screwdriver with a blue wooden handle'. Based on this description, provide the most fitting and concise object name. Output: screwdriver Input Caption: "a white and blue coffee mug with a label, featuring a blue lid and a yellow and white design, resembling a honey jar and a plastic container with the word 'Ulma' on it." Prompt: The caption is 'a white and blue coffee mug with a label, featuring a blue lid and a yellow and white design, resembling a honey jar and a plastic container with the word 'Ulma' on it.'. Based on this description, provide the most fitting and concise object name. Output: Coffee mug

Figure 15. Summarizing Caption3D [47] captions into nouns with Llama3 [15]. The LLM is capable of finding the noun that is the main subject of the caption.

Function	Objects
scrape with	knife, screw, card, dagger, pen, coin, pencil, screwdriver, shovel, key, spoon, needle, scissors, pickaxe, fork
	spatula, CD, hook, ruler, credit card, pitchfork, lid, pin, comb, awl, cleaver, trowel, razor blade, nail, toothpick
	hockey stick, machete, rake, paddle, paper clip, license plate, hoe, corkscrew, box cutter, chisel, brush, grater
	stylus, scalpel, file, letter opener, squeegee, peeler
press with	smartphone, bottle, shoe, stone, bowl, mug, water bottle, jug, teapot, hammer, bucket, cup, jar, book, plate, candle
	holder, tray, brick, pot, coffee pot, boot, flask, spoon, cutting board, pan, mallet, spatula, glass, kettle, plank, tablet
	credit card, can, lid, ladle, CD case, saucepan, stamp, clipboard, paddle, pestle, hoe, meat tenderizer
pound with	axe, bottle, shoe, bowl, water bottle, hammer, jar, pipe, candle holder, flashlight, wrench, brick, pot, basebal
	bat, screwdriver, dumbbell, shovel, boot, remote control, spoon, pan, mallet, pickaxe, spatula, bowling pin, log
	crowbar, can, ladle, rolling pin, gavel, cleaver, hockey stick, baton, saucepan, cricket bat, club, hairbrush, pestle
	meat tenderizer
pierce with	screw, knife, sword, dagger, pen, pencil, drill, screwdriver, needle, scissors, pickaxe, stilettos, fork, hook, pitchfork
	pin, spear, fish hook, dart, awl, chopsticks, harpoon, nail, toothpick, machete, skewer, golf tee, corkscrew, box
	cutter, chisel, stylus, scalpel, safety pin, letter opener
poke with	pipe, pencil, stick, pliers, screwdriver, key, rod, spoon, needle, pickaxe, fork, toothbrush, branch, pin, paintbrush
	awl, chopsticks, nail, coat hanger, dowel, baton, antenna, toothpick, skewer, crayon, matchstick, tweezers, tongs
	drumstick, stylus, stirrer, letter opener
mix with	knife, pen, pencil, screwdriver, rod, spoon, fork, spatula, toothbrush, branch, ruler, ladle, awl, chopsticks, baton
	straw, marker, skewer, paddle, brush, tongs, whisk, scalpel, stylus, stirrer, letter opener
pour with	mug, bottle, shoe, bowl, jug, water bottle, teapot, bucket, cup, jar, hat, pot, coffee pot, oil can, flask, pan, watering
	can, hard hat, kettle, glass, can, ladle, saucepan, decanter, coconut shell
cut with	axe, knife, sword, dagger, key, scissors, spatula, CD, ruler, credit card, saw, cleaver, razor blade, machete, box
*.1	cutter, ice skate, chisel, scalpel, pizza cutter, letter opener
scoop with	mug, shoe, bowl, jug, seashell, bucket, cup, hat, pot, shovel, flask, spoon, pan, hard hat, glass, ladle, trowel
roll onto	saucepan, dustpan, coconut shell
TOII OIIIO	cylinder, mug, bottle, water bottle, cup, jar, pen, pipe, flashlight, glass, bowling pin, log, can, rolling pin, battery lipstick, dowel, marker, roller
dig with	knife, stick, screwdriver, shovel, key, spoon, pickaxe, fork, ruler, awl, trowel, chopsticks, nail, paddle, hoe, dustpan
dig with	chisel, plow
sweep with	card, shovel, fork, spatula, broom, credit card, pitchfork, trowel, hockey stick, feather, rake, paddle, hairbrush
sceptal	hoe, mop, brush, squeegee
pry with	knife, dagger, wrench, screwdriver, shovel, key, spoon, pickaxe, fork, spatula, ruler, crowbar, chopsticks, car
1 7	opener, bottle opener, corkscrew, chisel
lift with	knife, seashell, plate, tray, shovel, cutting board, spoon, fork, spatula, ruler, lid, cleaver, trowel, paddle, clipboard
	dustpan
pull with	hammer, L-bracket, hook, crowbar, fish hook, coat hanger, harpoon, carabiner, hoe, grappling hook, grabber
spread with	knife, card, spoon, spatula, ruler, credit card, cleaver, trowel, brush
brush with	broom, toothbrush, paintbrush, feather, hairbrush, mop, brush
write with	pen, pencil, paintbrush, lipstick, marker, crayon, stylus
hang onto	mug, curtain ring, hook, fish hook, coat hanger, carabiner, paper clip
peel with	knife, dagger, box cutter, chisel, scalpel, peeler
pick up with	pliers, chopsticks, tweezers, tongs, grabber
wedge with	axe, wedge, spatula, chisel
apply torque with	wrench, pliers, grabber
sift with	basket, strainer, colander

Table 7. The list of functions and the associated objects in our dataset's function/object taxonomy.

axe	213	pencil	91	broom	26	razor blade	12	paper clip	8
knife	200	brick	85	CD	25	nail	12	license plate	8
screw	200	pot	80	glass	24	chopsticks	12	hoe	8
card	200	stick	78	kettle	24	harpoon	12	corkscrew	8
smartphone	200	drill	73	bowling pin	24	coat hanger	12	box cutter	8
bottle	200	pliers	71	toothbrush	24	dowel	12	pestle	8
shoe	200	baseball bat	66	plank	23	lipstick	12	matchstick	8
stone	200	screwdriver	56	hook	22	toothpick	11	tweezers	8
bowl	200	dumbbell	54	tablet	22	hockey stick	11	ice skate	8
mug	200	coffee pot	53	log	22	machete	11	dustpan	8
water bottle	200	shovel	51	branch	22	saucepan	11	mop	8
jug	200	L-bracket	51	ruler	21	baton	11	colander	8
cylinder	200	oil can	50	credit card	20	antenna	11	chisel	7
teapot	199	key	49	crowbar	19	carabiner	11	brush	7
hammer	199	wedge	48	pitchfork	18	stamp	10	tongs	7
sword	199	boot	47	can	18	cricket bat	10	grater	6
seashell	192	remote control	47	lid	17	skewer	10	stylus	6
bucket	187	flask	41	ladle	16	golf tee	10	scalpel	6
dagger	183	rod	41	rolling pin	16	crayon	10	drumstick	6
cup	181	spoon	39	saw	16	straw	10	whisk	6
jar	175	cutting board	39	pin	15	marker	10	grappling hook	6
book	173	pan	35	CD case	15	roller	10	safety pin	5
basket	163	watering can	34	spear	15	feather	10	grabber	5
pen	160	needle	33	gavel	14	strainer	10	file	4
plate	158	scissors	33	fish hook	14	rake	9	stirrer	4
coin	140	hard hat	33	battery	14	paddle	9	pizza cutter	4
candle holder	126	mallet	32	comb	13	clipboard	9	letter opener	3
pipe	126	curtain ring	32	awl	13	club	9	squeegee	2
hat	118	pickaxe	29	cleaver	13	hairbrush	9	peeler	2
tray	114	stilettos	29	dart	13	decanter	9	meat tenderizer	2
flashlight	114	fork	27	paintbrush	13	can opener	9	coconut shell	2
wrench	112	spatula	27	trowel	12	bottle opener	9	plow	1
				-					

Table 8. Categories in our curated dataset and the number of assets in each category.