

One-shot Learning for Task-oriented Grasping

Valerija Holomjova^{1*}, Andrew J. Starkey¹, Bruno Yun², Pascal Meißner³

Abstract—Task-oriented grasping models aim to predict a suitable grasp pose on an object to fulfill a task. These systems have limited generalization capabilities to new tasks, but have shown the ability to generalize to novel objects by recognizing affordances. This object generalization comes at the cost of being unable to recognize the object category being grasped, which could lead to unpredictable or risky behaviors. To overcome these generalization limitations, we contribute a novel system for task-oriented grasping called the One-shot Task-oriented Grasping (OS-TOG) framework. OS-TOG comprises four interchangeable neural networks that interact through dependable reasoning components, resulting in a single system that predicts multiple grasp candidates for a specific object and task from multi-object scenes. Embedded one-shot learning models leverage references within a database for OS-TOG to generalize to novel objects and tasks more efficiently than existing alternatives. Additionally, the paper presents suitable candidates for the framework’s neural components, covering essential adjustments for their integration and evaluative comparisons to state-of-the-art. In physical experiments with novel objects, OS-TOG recognizes 69.4% of detected objects correctly and predicts suitable task-oriented grasps with 82.3% accuracy, having a physical grasp success rate of 82.3%.

Index Terms—Deep Learning in Grasping and Manipulation, Grasping, Computer Vision for Automation, Recognition

I. INTRODUCTION

Task-oriented Grasping (TOG) involves finding a grasp pose on an object that enables the completion of a task [1, 2]. For instance, grasping the handle of a mug to *pour* out its contents. TOG is a vital preliminary step to accomplishing manipulations required by robotic grasping systems used for assistive robotics (e.g. doing household chores) or assembly tasks [3]. The ability to understand and interact with surrounding objects enables robotic manipulators to operate in unconstrained environments without human intervention.

Creating a dataset with sufficient coverage of the tasks and objects present in the real world to train TOG models is currently unfeasible [4], which encourages TOG solutions that can generalize to new object categories or tasks. [1, 5, 6] show capabilities of generalizing to novel objects by predicting and leveraging affordances [7, 8]. Affordances are regions of an object that represent a functional interaction (e.g. *cut*, *contain*). By mapping relationships between affordances and tasks, the robotic system can identify task-suitable grasping regions. For instance, if the robotic ma-

This research is funded by a studentship awarded by the School of Engineering at the University of Aberdeen, Scotland UK.

¹First author (*corresponding author) and second author are with the School of Engineering, University of Aberdeen, Scotland UK {v.holomjova.21, a.starkey}@abdn.ac.uk

²Third author is with the School of Natural and Computing Sciences, University of Aberdeen, Scotland UK bruno.yun@abdn.ac.uk

³Fourth author is with Würzburg-Schweinfurt Technical University of Applied Sciences (THWS), Germany pascal.meissner@thws.de

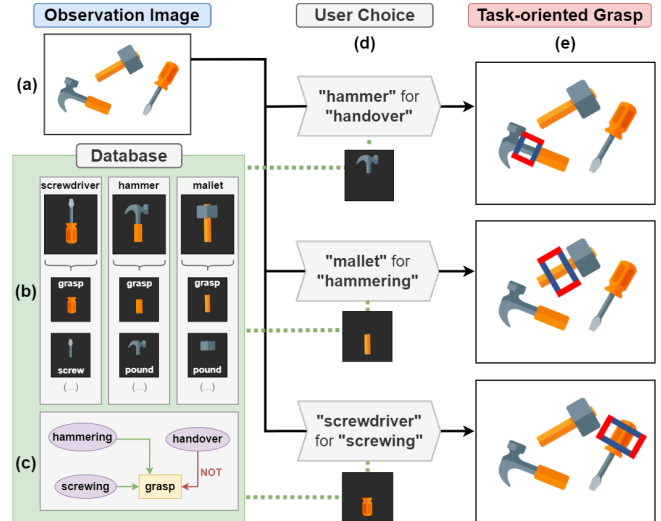


Fig. 1: Given an observation image (a), OS-TOG uses a database containing information on objects (b) and tasks (c) to predict a task-oriented grasp (e) for a user-specified object and task (d). The object database (b) contains a list of objects annotated with their labeled affordances. The task database (c) maps relationships between tasks and affordances to acquire references of task-suitable regions that should be grasped.

nipulator was given a *handover* task, it could grasp the hammer by the *pound* affordance, leaving the *grasp* region to be safely retrieved by a human. However, these solutions are unable to recognize the object category being grasped and cannot generalize to novel affordance categories. The inability to recognize objects could lead to hazardous actions when working in unconstrained environments (e.g. an assistive robot grasping the handle of a knife instead of a fork).

To address the generalization limitations of TOG systems, we present a novel system called the One-shot Task-oriented Grasping (OS-TOG) framework. OS-TOG leverages a database of objects and task-affordance relations to produce task-oriented grasps for specified objects and tasks from multi-object scenes (Fig. 1). The framework comprises four interchangeable neural networks and dependable reasoning components that enable and enhance interaction between the networks. OS-TOG can generalize to novel objects and tasks due to its embedded one-shot learning models. These models are trained to encode images into a feature representation that effectively captures the differences and similarities between them. The learned embeddings can discern the correspondence between a given image and reference images within OS-TOG’s database, even if either lies beyond its training dataset. To the best of our knowledge, this is the first TOG system that incorporates these networks.

OS-TOG is limited to recognizing objects and tasks de-

finer in its database but can generalize to more by adding a single annotated example of each new object or task to the database. This requires significantly less labeling and training effort than creating an entire dataset to re-train a standard TOG model. These generalization capabilities make OS-TOG ideal for domestic assistive robotic settings, as it can easily adapt to personal objects and user-specific tasks.

This research aims to evaluate the performance of OS-TOG on TOG and explore the extent to which it can generalize to new objects and tasks. Our contributions are three-fold; 1) We present OS-TOG, a novel framework for TOG that can generalize to new objects and tasks. 2) We propose suitable sub-models for its interchangeable neural network components, provide adjustments for their integration, and evaluate them to state-of-the-art in their respective tasks. 3) Experiments with a 7-DoF robotic arm having an RGB-D camera are carried out to demonstrate OS-TOG’s ability to perform TOG on previously unseen objects and various tasks.

II. RELATED WORK

Task-oriented grasping is a challenging research area that involves finding a suitable grasp on an object to fulfill a specified task. Machine-learning solutions have been proposed to solve TOG [1–6, 9–11], but rely on large amounts of training data. Extensive manual efforts are needed to create TOG datasets, which still have limited coverage of the objects and tasks found in the real world.

To overcome dataset limitations, most literature focuses on using alternative methods to generate datasets for training or improving the generalization capabilities of their systems. These data alternative methods include generating synthetic data [1, 5, 9], training in simulated environments [2, 3], or leveraging video footage of human-object interactions [4, 11]. However, these methods often show a drop in performance in real-world scenes or require re-training for new objects and tasks. Certain solutions have demonstrated the ability to generalize to new object categories by learning feature representations of task-relevant geometries within similar objects [1–3], leveraging semantic knowledge between tasks and objects [10] or segmenting parts of the object with a particular functionality (i.e. affordance) to assist in predictions [5, 6]. Nonetheless, some of these solutions can’t operate in multi-object scenes with all being unable to recognize the objects they are grasping. This reveals a research gap for novel TOG solutions that can recognize objects and generalize to new objects and tasks within multi-object scenes with minimal training effort required.

The task of segmenting and labeling parts of objects with a functionality is referred to as “affordance segmentation”. Machine-learning solutions that predict affordance maps are mostly semantic segmentation models consisting of object detection components. For instance, [7] proposed a CNN-based framework that detects objects and then segments their affordances from RGB images. [8] extends this solution by creating an end-to-end architecture, similar to a Mask R-CNN [12], that recognizes objects and segments affordances in parallel. [13] build upon an object detector and add domain

adaptation components to learn from synthetic data and adapt to real-world data. [14] construct an end-to-end autoencoder that learns from human-object interactions. These solutions can’t generalize to new objects and affordances without requiring re-training. To this matter, [15] and [16] demonstrate the use of one-shot learning techniques to find and segment previously unseen affordances without requiring re-training.

One-shot learning is the task of classifying objects from a single or few examples. The most popular one-shot learning technique is Siamese networks [17], composed of two sub-networks that share weights to predict the similarity between two different inputs. [18] use a Mask R-CNN followed by a Siamese network to recognize target objects from cluttered bins in order to be grasped. Alternative one-shot learning methods surpassing the performance of Siamese networks have also been introduced over the years. [19] create a novel two-branched approach for one-shot image segmentation, where one branch generates parameters from the query image which is used by the second branch to produce a segmentation mask from the query image. [20] segments objects from cluttered scenes by segmenting instances, masking their backgrounds, and computing the best match. [21] wins first place in a robotics challenge for categorizing objects in a cluttered bin. Their solution isolates each object through grasping and then matches them to the nearest object in a database using a two-stream CNN-based model. The system obtains a high recognition rate but is inefficient in multi-object scenes when only a specific object category is desired.

Inspired by the solutions of [15, 16, 18, 21] and one-shot learning techniques, we design a novel TOG framework that recognizes and retrieves objects from multi-object scenes for specified tasks. One-shot learning models in the framework compare correspondences between scene images and examples provided in a database. These correspondences are used to match object categories or determine affordance regions. This enables the framework to generalize to novel objects and tasks without requiring re-training. With the added capability of identifying object categories and functioning in multi-object scenarios, this system is more suitable for assistive robotics than existing alternatives.

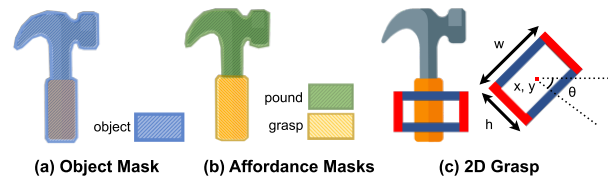


Fig. 2: An example of an object annotated with an object segmentation mask (a) and affordance segmentation masks (b), and an example of a grasp pose $g = (x^g, y^g, w^g, h^g, \theta^g)$ on an object, with center co-ordinates (x^g, y^g) , gripper opening w^g , gripper size h^g and rotation $\theta^g \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ (c).

III. PROBLEM STATEMENT

Assume a database D containing a set of user-defined tasks D_T and annotated objects D_O . Each object in D_O is represented by an RGB image which is annotated with an object segmentation mask (Fig. 2a) and suitable affordance segmentation masks (Fig. 2b) from a determined set of affordances A . These segmentation masks are binary masks

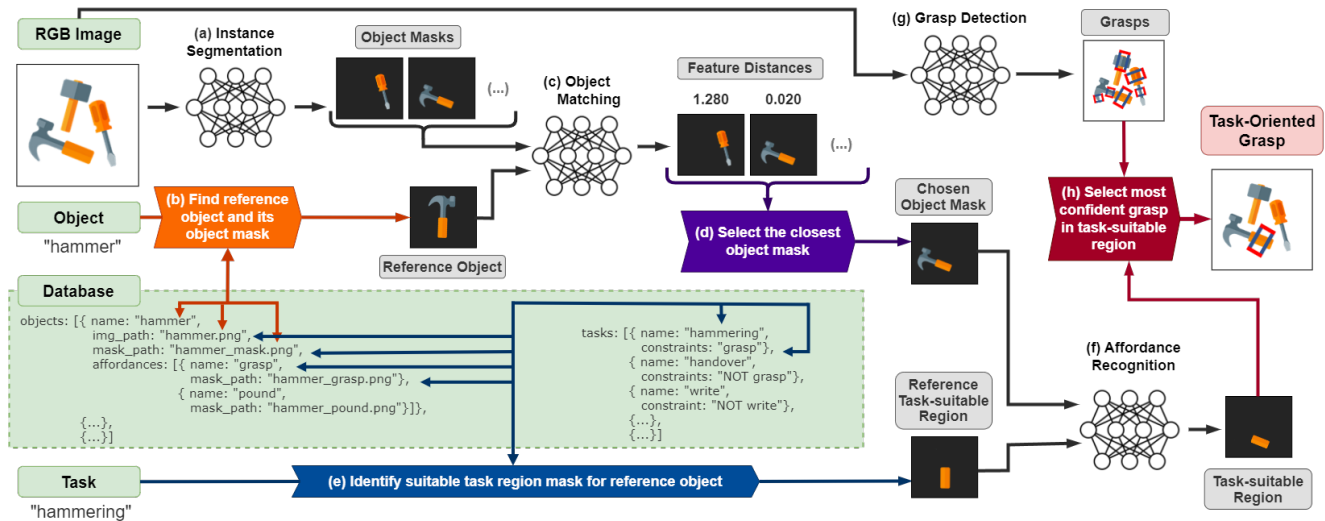


Fig. 3: An illustration of the proposed framework (OS-TOG) for task-oriented grasping. OS-TOG takes as input an RGB image, target object, target task, and database of objects and tasks (green) and outputs a 2D task-oriented grasp on the target object for the target task (red). The four sub-models are represented by grey arrows (a, c, g, f), and the reasoning components are represented by colored arrows (b, d, e, h).

of the RGB image. Each task in D_T is mapped to a suitable affordance $a \in A$ through “a” or “NOT a” relations, denoting whether grasping in the region of a will allow the task to be accomplished or not. Given an RGB image of a multi-object scene containing a set of objects N , database $D = D_T \cup D_O$, target object $o \in (N \cap D_O)$, and a target task $t \in D_T$, the objective is to localize and grasp o from the scene in a manner that satisfies the conditions of t . The target object o is localized by predicting a bounding box $b = (x^b, y^b, w^b, h^b)$ and segmentation map of an object $n \in N$ from the scene s.t. $n = o$. Given a parallel plate gripper, o can be grasped by predicting a planar grasp pose that can be parameterized as an oriented bounding box $g = (x^g, y^g, w^g, h^g, \theta^g)$ (Fig. 2c). The grasp is task-suitable if (x^g, y^g) is within the affordance regions of o that satisfy the relations defined by t in D_T .

IV. ONE-SHOT TASK-ORIENTED GRASPING

OS-TOG is a novel TOG framework with embedded one-shot learning models that leverage references from a database to generalize to new objects and tasks without needing re-training. OS-TOG comprises four neural networks and reasoning components that process predictions from the sub-models and acquire object or affordance references from the database (Fig. 3). The neural network components of OS-TOG are interchangeable allowing us to select state-of-the-art models for their respective tasks. Further details regarding the neural and reasoning components are provided below.

Instance Segmentation - The first sub-model (Fig. 3a) performs category-agnostic instance segmentation to segment and isolate all objects from N in the 640×480 RGB scene image, producing a binary mask for each object. Each binary mask is combined with the RGB scene image to create a set of object masks M . We use Mask R-CNN [12] with a ResNet-50 FPN backbone and weights pre-trained on ImageNet [22]. The heads predict an object class, a bounding box, and a binary segmentation map for each object. We replace the mask and class predictor heads to predict two

object classes; “object” or “background”, and train the model using the same multi-task loss function defined in [12].

Object Matching - Reasoning components retrieve an object mask o^m of the target object o from the database through string matching, and then combining its 256×256 RGB image $o^i \in D_O$ and binary mask $o^b \in D_O$ (Fig. 3b). Each predicted object mask $m_i \in M$ is magnified by cropping its bounding box and padding it to a size of 256×256 . The magnified predicted masks and target object mask are fed into a one-shot learning model which extracts their embedding vectors and computes the L2 distance between them to determine which m is closest and most similar s.t. $\phi = \text{argmin}_{m_i \in M} \{d(m_i, o^m)\}$ (Fig. 3d). For object matching, we re-implement N-net from [21] in PyTorch as it had the highest novel object recognition accuracy.

During training, N-net is comprised of three streams. One stream computes features for a reference object image x^a , and the other two streams compute features for two query object images (positive x^p and negative x^n). x^a shares the same object class as x^p , whereas x^n has a different object class. N-net uses embedding vectors from a frozen ResNet-50 model with pre-trained ImageNet weights for the reference image stream to improve novel object accuracy. This is further improved by using multiple product images for each reference object in training and selecting the nearest one based on L2 distances between features. We replace N-net’s original training loss function with standard triplet loss (TL) [23] and use a balanced batch sampler (BBS) after seeing improved accuracies in preliminary experiments. The BBS randomly selects p samples from k object classes in each mini-batch, generating $p \times k$ triplets in each mini-batch. Triplet loss is a metric that minimizes the L2 embedding distance between x^a and x^p , and maximizes the distance between x^a and x^n by a minimum margin α . Given that $f(x)$ is the embedding vector of an image x , triplet loss L_t for a triplet (x_i^a, x_i^p, x_i^n) can be defined as;

$$L_t = \max\{d(x_i^a, x_i^p) - d(x_i^a, x_i^n) + \alpha, 0\} \quad (1)$$

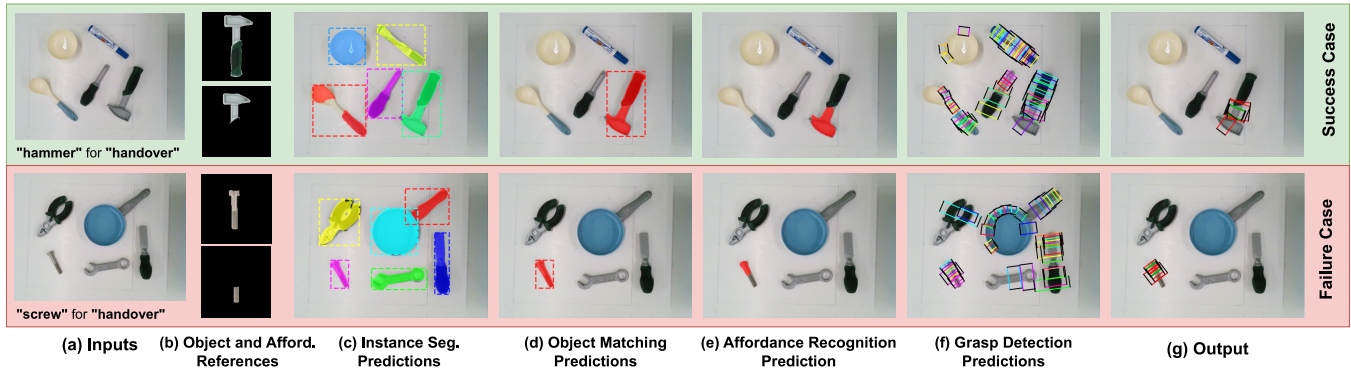


Fig. 4: Predictions made by OS-TOG in physical experiments when given an RGB image and target object and task as input (a). OS-TOG retrieves an object and suitable affordance reference of the target inputs from the database to be used by its one-shot learning models (b). This figure shows predictions of each sub-model (c-f) and the final set of task-oriented grasp candidates (g). The green candidate represents the most confident grasp which is executed.

$$d(x_i, y_i) = \|f(x_i) - f(y_i)\|_2 \quad (2)$$

One-shot Affordance Recognition - OS-TOG retrieves a reference binary affordance mask o^a of o depending on the task-affordance relation defined by t in the database (Fig. 3e). For “ a ” relations, OS-TOG retrieves a binary mask $a^b \in D_O$ of a on the target object o s.t. $o^a = a^b$. For unconstrained relations (e.g. *transport* in Table II), we take the binary mask $o^b \in D_O$ of o s.t. $o^a = o^b$. For “NOT a ” relations, we obtain a binary mask s.t. $o^a = o^b - a^b$. The one-shot affordance recognition model takes as input ϕ , o^i , and o^a to produce a binary affordance mask ϕ^a that represents the task-suitable region in ϕ (Fig. 3f). We use the AffCorrs model [16] for one-shot affordance recognition without re-implementation or training as it is unsupervised and it is the only sub-model in OS-TOG that is not re-implemented or trained. In physical experiments we found that AffCorrs performs significantly better if the orientation of the objects in ϕ , o^i are similar, hence, we rotate $\{o^i, o^a\}$ in 45° intervals and use the pair with the smallest L2 distance between o^i and ϕ .

Grasp Detection - A grasp detection model predicts grasp candidates on the image scene (Fig. 3g). We use our baseline from previous work [24] that uses a Faster R-CNN [25] model with a ResNet-50 FPN backbone and pre-trained ImageNet weights. Faster R-CNN predicts an object class and bounding box for each object in a scene, hence, we replace the object classes it predicts and an orientation class r . We discretize $\theta^g \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ orientation values into $Q = 12$ classes s.t. the set of possible orientation classes is $R = \{r_1, \dots, r_Q\}$ with an additional class denoting an invalid grasp and replace the Faster R-CNN class predictor head to predict $Q + 1$ classes. The model is trained using the same multi-task loss function defined in [25]. The grasp candidates are filtered to only consider those that have a confidence threshold > 0.5 and (x^g, y^g) lies in the predicted affordance region ϕ^a (Fig. 3h). The most confident grasp is taken giving a task-oriented grasp on object o for task t .

V. EXPERIMENTS AND EVALUATION

OS-TOG was built in PyTorch using Python 3.8. Since there is currently no publicly available gold-standard TOG dataset, the system is evaluated in three separate settings.

First, we evaluate each sub-model of OS-TOG that we implemented to the state-of-the-art in their respective tasks (Sec. V-C). Second, we evaluate the performance of OS-TOG in affordance recognition which uses all its sub-model components except for the final grasp model (Sec. V-D). Lastly, the entire framework is evaluated on TOG in physical experiments with random household objects (Sec. V-E).

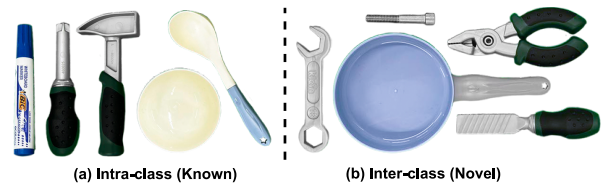


Fig. 5: Objects used for physical grasping experiments.

A. Evaluation Metrics

We adopt standard metrics from literature [16, 26] to evaluate the performance of our models; the Intersection over Union (IoU) score and F_β^w -measure [27] using $\beta = 1$ signifying equal importance to weighted recall and precision. The grasp accuracy score is calculated by classifying a predicted grasp as a success if it has an IoU score greater than 25% with a ground-truth grasp and a θ^g angle difference within 30° . We report the same metrics as [21] for one-shot learning denoting the model’s accuracy in recognizing: all object categories in the dataset (mixed), categories that were present during training (known), categories that were never present in training (novel), and the percentage it distinguishes whether an object belongs to a known or novel category (K vs N).

B. Datasets

Cornell grasp dataset [28] - contains 1,035 RGB-D images of single-object scenes covering 280 object classes hand-annotated with multiple grasps.

OCID grasp dataset [27] - has 1,763 RGB-D images of multi-object scenes covering 30 object categories annotated with multiple grasps. Each object in the scene is also labeled with a segmentation mask of its object class.

ARC image matching dataset [2] - has over 4,000 images of 61 different object categories against a green screen with matching masked product images in various orientations.

TABLE I: OBJECTS AND AFFORDANCES

Split	Object	Affordances
INTRA-CLASS	hammer	grasp, pound
	spoon	grasp, contain
	screwdriver	grasp, screw
	marker	grasp, write
	bowl	grasp, contain
INTER-CLASS	chisel	grasp, file
	frying pan	grasp, fry
	wrench	grasp, loosen, unscrew
	pliers	grasp, hold
	screw	grasp, screw

TABLE II: TASKS AND GRASP CONSTRAINTS

Task	Constraints	Affected Objects
<i>transport</i>	-	all
<i>hammering</i>	grasp	hammer
<i>handover</i>	NOT grasp	all except bowl
<i>filig</i>	NOT file	chisel
<i>loosening</i>	NOT loosen	wrench
<i>unscrewing</i>	NOT unscrew	unscrew
<i>holding</i>	grasp	pliers
<i>writing</i>	grasp	marker
<i>screwing</i>	NOT screw	screwdriver, screw
<i>scooping</i>	NOT contain	spoon
<i>frying</i>	NOT grasp	frying pan

UMD dataset [29] - has 30,000 RGB-D images of single-object scenes containing random household objects of 17 categories and 105 classes, and 7 affordance classes (grasp, cut, scoop, contain, pound, support, and wrap-grasp).

UMDⁱ dataset [16] - a subset of the UMD dataset tailored for one-shot learning containing only a single instance of each object class from UMD with original annotations kept.

C. Evaluating the sub-models of OS-TOG

Each sub-model is trained and evaluated on datasets used by state-of-the-art models in grasping literature for instance segmentation, grasp detection, and one-shot learning on their reported metrics (Tables III and IV). As mentioned in Section IV, we do not train or evaluate AffCorrs but provide reported metrics from [16] on the UMDⁱ Dataset (Table V). These metrics are not directly comparable to the other baselines which are supervised and evaluated on the full UMD Dataset.

D. Evaluation OS-TOG on Affordance Recognition

OS-TOG is trained and evaluated on the UMD dataset on two separate data splits. Intra-class signifies that all object categories are present in both data splits, whereas inter-class signifies that the test set contains exclusive object categories. A database containing a single example of each object class with all possible labeled affordances is built from the dataset. The evaluation procedure begins by iterating through each image in the test set and segmenting the object in the image, then matching it to the nearest reference object in the database. For each ground-truth affordance label in the current scene, a reference affordance mask having the same label is retrieved and used to predict an affordance mask on the scene object. Lastly, we calculate the IoU score and F_{β}^w -measure between the ground-truth and predicted affordance masks when the scene object is correctly matched (Table V).

E. Evaluation OS-TOG in Physical Experiments

For physical experiments, the grasp detection model was trained on Cornell, and the object matching model was trained on UMD. The instance segmentation model was trained on OCID, single-object UMD scenes, and then 20 multi-object scenes from UMD that we manually annotated. We conduct physical experiments using a 7-DoF robotic arm by Franka Emika equipped with a D415 Intel RealSense camera. The Frankx library [30] was used for motion planning.

Experiments are carried out on 10 random objects having at least one affordance. Half of the objects were seen by at least one of the trained sub-models in training (intra-class) (Fig. 5a), and the latter were never seen in training (inter-class) (Fig. 5b). We create a database by annotating a single instance of each object with suitable object masks and affordances (Table I) and create a list of tasks mapped to suitable affordance regions (Table II). Our approach allows us to have multiple affordances on objects even if overlapped. We carry out five trials for each object per task. Table VI shows OS-TOG’s ability to segment the object (Obj. Det.), match the detected object to the reference object correctly (Obj. Match.), detect a grasp in the correct affordance region of the correctly matched object (Grasp Det.), and physically succeed in grasping the predicted grasp (Grasp Succ.).

VI. RESULTS

Instance Segmentation and Grasp Detection - As shown by Table III, our trained grasp model and instance segmentation achieve comparable performance to Det Seg [26] on both datasets and significantly outperforms Det Seg on grasp detection in the OCID grasp dataset. This suggests that Faster R-CNN may perform better in multi-object scenes.

TABLE III: GRASP DET. AND INSTANCE SEG. RESULTS

Method	Dataset	Grasp Accuracy (%)	IOU (%)
Det Seg [26]	OCID grasp	89.0	94.1
	Cornell	98.2	-
Faster R-CNN [25] and Mask R-CNN [12] (ours)	OCID grasp	98.1	93.0
	Cornell	96.6	-

Object Matching - Our re-implementation of N-net performs better in all metrics than N-net from [21]. Our model doesn’t outperform K-net and the Two-stage model [21] in known/mixed object recognition, but excels in novel object recognition which is the most relevant metric for our system.

TABLE IV: OBJECT MATCHING RESULTS ON ARC [21]

Method	K vs N	Known	Novel	Mixed
N-net [21]	69.2	56.8	82.1	64.6
K-net [21]	93.2	99.7	29.5	78.1
Two-stage K-net + N-net [21]	93.2	93.6	77.5	88.6
N-net + TL + BBS (ours)	71.7	75.5	86.7	78.7

Affordance Recognition - OS-TOG is able to correctly detect and match objects from the scene 65.4% of the time for inter-class objects and 71.5% for intra-class. When correctly detecting and matching the object, OS-TOG is able to outperform the baseline approaches on nearly all affordance types and achieves an average IoU and F_{β}^w score of 0.77 and 0.85 for intra-class objects and 0.77 and 0.84 for inter-class objects. The similarity between OS-TOG’s inter-class and intra-class results demonstrates the generalization capabilities of the network to new objects. OS-TOG’s worse performance on inter-class splits represents the one-shot learning networks struggling to capture intricate

TABLE V: AFFORDANCE RECOGNITION RESULTS

Method	Data	Split	Grasp		Cut		Scoop		Contain		Wrap-Grasp		Pound		Support		Total Avg.	
			IoU	F_{β}^w	IoU	F_{β}^w	IoU	F_{β}^w	IoU	F_{β}^w	IoU	F_{β}^w	IoU	F_{β}^w	IoU	F_{β}^w	IoU	F_{β}^w
AffordanceNet [8]	UMD	Intra-class	-	0.73	-	0.81	-	0.76	-	0.83	-	0.82	-	0.79	-	0.84	-	0.80
ResNet [31]	UMD	Inter-class	0.33	-	0.51	-	0.69	-	0.52	-	0.85	-	0.09	-	0.51	-	0.50	-
		Intra-class	0.71	-	0.79	-	0.86	-	0.86	-	0.84	-	0.72	-	0.55	-	0.76	-
AffCorrs [16]	UMD ⁱ	Inter-class	0.39	0.41	0.51	0.50	0.62	0.65	0.71	0.75	0.83	0.87	0.72	0.73	0.82	0.79	0.66	0.68
		Intra-class	0.55	0.65	0.72	0.81	0.73	0.81	0.82	0.87	0.83	0.89	0.78	0.87	0.82	0.87	0.75	0.82
OS-TOG (ours)	UMD	Inter-class	0.58	0.69	0.65	0.76	0.78	0.86	0.85	0.91	0.80	0.89	0.86	0.93	0.87	0.93	0.77	0.84
		Intra-class	0.55	0.66	0.66	0.75	0.82	0.89	0.90	0.94	0.81	0.89	0.84	0.91	0.78	0.86	0.77	0.85

*Bolded numbers represent the top-performing scores for each affordance category and data split.

TABLE VI: PHYSICAL EXPERIMENT RESULTS

Object	Task	Success Rates (%)				
		Obj. Det.	Obj. Match.	Grasp Det.	Grasp Succ.	
INTRA-CLASS	Bowl	transport	100.0	100.0	100.0	40.0
		Screwdriver	transport	100.0	80.0	100.0
	Spoon	screwing	100.0	80.0	100.0	75.0
		handover	100.0	80.0	100.0	50.0
		transport	100.0	80.0	100.0	50.0
	Hammer	scooping	100.0	60.0	100.0	100.0
		handover	100.0	20.0	0.0	N/A
		transport	100.0	60.0	100.0	66.7
	Marker	hammering	100.0	60.0	100.0	66.7
		handover	80.0	75.0	100.0	100.0
transport		100.0	100.0	100.0	100.0	
INTER-CLASS	Chisel	writing	100.0	100.0	80.0	100.0
		handover	100.0	80.0	75.0	100.0
	Frying Pan	handover	100.0	80.0	75.0	100.0
		transport	100.0	60.0	66.7	50.0
		filing	100.0	80.0	100.0	50.0
	Pliers	handover	100.0	40.0	100.0	00.0
		transport	100.0	100.0	100.0	60.0
		frying	100.0	100.0	0.0	N/A
	Wrench	handover	100.0	100.0	100.0	100.0
		transport	100.0	40.0	100.0	100.0
holding		100.0	40.0	100.0	100.0	
Screw	handover	100.0	80.0	100.0	75.0	
	transport	100.0	80.0	100.0	100.0	
	unscrewing	100.0	80.0	100.0	100.0	
Screw	handover	100.0	40.0	50.0	100.0	
	loosening	100.0	60.0	100.0	100.0	
	transport	100.0	80.0	100.0	100.0	
Screw	screwing	80.0	80.0	75.0	100.0	
	handover	80.0	100.0	25.0	100.0	
	handover	80.0	100.0	25.0	100.0	
Avg. Total		98.5	75.0	88.8	77.4	
Avg. Total		97.5	69.4	82.3	82.3	

*N/A signifies the system failed to detect any task-suitable grasps in trials.

details within the object classes. This could be resolved by giving the networks multiple examples for better contextual understanding.

Physical Experiments - Table VI shows that OS-TOG successfully matched previously seen objects at a rate of 75.0%, and 69.4% for novel object categories when segmented. Object matching often failed due to incomplete or noisy segmentations, or object color similarity. Hammer-screwdriver mismatches occurred when the head of the hammer was poorly segmented due to shared colored properties. The frying pan was consistently segmented into two parts (Fig. 4), as they were commonly present shapes within the training datasets. Results show that task-suitable grasps were correctly predicted at a rate of 88.8% for known objects and 82.3% for novel objects. Grasp detection failures are attributed to mis-segmentations, insufficient grasps predicted on the target scene object, and affordance recognition failures. The affordance model struggled in segmenting the

correct affordance region in the screw due to the similar shapes of its thread and head (Fig. 4). Physical grasp success rates were 77.4% for known objects and 82.3% for novel objects with most failures attributed to objects slipping from the grippers or predicting w^g too small. Experimental findings suggest OS-TOG can improve by giving more intricate training examples to the instance segmentation model and transitioning to predicting 6D grasp poses which could yield more stable grasps.

VII. CONCLUSION

We present a novel framework called OS-TOG, composed of four sub-models and reasoning components that coordinate to perform task-oriented grasping. By leveraging the properties of one-shot learning models and a database of individually annotated objects and tasks, OS-TOG produces task-oriented grasps on previously unseen objects and tasks from RGB multi-object scenes. Experimentation results showed that OS-TOG is capable of generalizing substantially to new objects and tasks, which is beyond the generalization capabilities of current task-oriented grasping systems. Future work involves extending the system to 6D grasp poses.

REFERENCES

- [1] Renaud Detry et al. "Task-oriented Grasping with Semantic and Geometric Scene Understanding". In: *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 3266–3273.
- [2] Kuan Fang et al. "Learning Task-Oriented Grasping for Tool Manipulation from Simulated Self-Supervision". In: *The International Journal of Robotics Research* 39.2-3 (2020), pp. 202–216.
- [3] Bowen Wen et al. "CaTGrasp: Learning Category-Level Task-Relevant Grasping in Clutter from Simulation". In: *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 6401–6408.
- [4] Mia Kokic et al. "Learning Task-Oriented Grasping From Human Activity Datasets". In: *IEEE Robotics and Automation Letters* 5.2 (2020), pp. 3352–3359.
- [5] Yunzhi Lin et al. "Using Synthetic Data and Deep Networks to Recognize Primitive Shapes for Object Grasping". In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 10494–10501.
- [6] Weiyu Liu et al. "CAGE: Context-Aware Grasping Engine". In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 2550–2556.

IEEE Robotics and Automation Letters (RA-L) paper, presented at ICRA 2024, Yokohama, Japan. Cite as RA-L paper.

- [7] Anh Nguyen et al. "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)*. IEEE, 2017, pp. 5908–5915.
- [8] Thanh-Toan Do et al. "AffordanceNet: An End-to-End Deep Learning Approach for Object Affordance Detection". In: *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 5882–5889.
- [9] Chenjie Yang et al. "Task-oriented Grasping in Object Stacking Scenes with CRF-based Semantic Model". In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 6427–6434.
- [10] Adithyavairavan Murali et al. "Same Object, Different Grasps: Data and Semantic Knowledge for Task-Oriented Grasping". In: *Proceedings of the 2020 Conference on Robot Learning*. PMLR, 2021, pp. 1540–1557.
- [11] Hui Li et al. "Learning Task-Oriented Dexterous Grasping from Human Knowledge". In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 6192–6198.
- [12] Kaiming He et al. "Mask R-CNN". In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 2961–2969.
- [13] Fu Jen Chu et al. "Learning Affordance Segmentation for Real-World Robotic Manipulation via Synthetic Images". In: *IEEE Robotics and Automation Letters* 4.2 (2019), pp. 1140–1147.
- [14] Spyridon Thermos et al. "A Deep Learning Approach to Object Affordance Segmentation". In: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 2358–2362.
- [15] Hongchen Luo et al. "One-Shot Affordance Detection". In: *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*. 2021, pp. 895–901.
- [16] Denis Hadjivelichkov et al. "One-Shot Transfer of Affordance Regions? AffCorrs!" In: *Proceedings of The 6th Conference on Robot Learning*. PMLR, 2022, pp. 550–560.
- [17] Jane Bromley et al. "Signature Verification Using A "Siamese" Time Delay Neural Network". In: *Advances in Neural Information Processing Systems* 6 (1993), pp. 669–688.
- [18] Michael Danielczuk et al. "Mechanical Search: Multi-Step Retrieval of a Target Object Occluded by Clutter". In: *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 1614–1621.
- [19] Amirreza Shaban et al. "One-Shot Learning for Semantic Segmentation". In: *Proceedings of the British Machine Vision Conference (BMVC)*. BMVA Press, 2017, pp. 167.1–167.13.
- [20] Claudio Michaelis et al. "One-Shot Segmentation in Clutter". In: *Proceedings of the 35th International Conference on Machine Learning*. 2018, pp. 3549–3558.
- [21] Andy Zeng et al. "Robotic Pick-and-Place of Novel Objects in Clutter with Multi-Affordance Grasping and Cross-Domain Image Matching". In: *The International Journal of Robotics Research* 41.7 (2022), pp. 690–705.
- [22] Olga Russakovsky et al. "ImageNet Large Scale Visual Recognition Challenge". In: *International Journal of Computer Vision* 115.3 (2015), pp. 211–252.
- [23] Florian Schroff et al. "FaceNet: A unified embedding for face recognition and clustering". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015, pp. 815–823.
- [24] Valerija Holomjova and Pascal Meißner. "Exploring Rotated Object Detection Models for Antipodal Robotic Grasping". In: *UKRAS22 Conference "Robotics for Unconstrained Environment" Proceedings*. 2022, pp. 62–63.
- [25] Shaoqing Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". In: *Advances in Neural Information Processing Systems*. 2015, pp. 91–99.
- [26] Stefan Ainetter and Friedrich Fraundorfer. "End-to-end Trainable Deep Neural Network for Robotic Grasp Detection and Semantic Segmentation from RGB". In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 13452–13458.
- [27] Ran Margolin et al. "How to Evaluate Foreground Maps?" In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2014, pp. 248–255.
- [28] Ian Lenz et al. "Deep Learning for Detecting Robotic Grasps". In: *The International Journal of Robotics Research* 34.4-5 (2015), pp. 705–724.
- [29] Austin Myers et al. "Affordance detection of tool parts from geometric features". In: *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 1374–1381.
- [30] Lars Berscheid. *Frankx: High-Level Motion Library for the Franka Emika Robot*. <https://github.com/pantor/frankx>. 2013.
- [31] Johann Sawatzky et al. "Weakly Supervised Affordance Detection". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, pp. 2795–2804.