

# Visual Affordance and Function Understanding: A Survey

Mohammed Hassanin, Salman Khan, Murat Tahtali

**Abstract**—Nowadays, robots are dominating the manufacturing, entertainment and healthcare industries. Robot vision aims to equip robots with the ability to discover information, understand it and interact with the environment. These capabilities require an agent to effectively understand object affordances and functionalities in complex visual domains. In this literature survey, we first focus on ‘Visual affordances’ and summarize the state of the art as well as open problems and research gaps. Specifically, we discuss sub-problems such as affordance detection, categorization, segmentation and high-level reasoning. Furthermore, we cover functional scene understanding and the prevalent functional descriptors used in the literature. The survey also provides necessary background to the problem, sheds light on its significance and highlights the existing challenges for affordance and functionality learning.

**Index Terms**—affordance prediction, functional scene understanding, deep learning, object detection



## 1 INTRODUCTION

Affordance understanding is concerned with the possible set of actions that an environment allows to an actor. In other words, this area of study aims to answer the question of how an object can be used by an agent? Ecological psychologist James Gibson was the first to introduce the concept of affordances in 1966 [1]. Since then, the theory of affordances has been extensively used in the design of better and robust robotic systems capable of operating in complex and dynamic environments [2]. In contrast to affordances which are directly dependent on the actor, function understanding relates to identifying the possible set of tasks which can be performed with an object. Object function is therefore a permanent property of an object independent of the characteristics of the user. Affordance and function understanding not only allow humans or AI agents to better interact with the world, but also provide valuable feedback to the product designers who need to consider possible interactions between users and products. As a result, this research topic is highly important for domestic robotics, content analysis and context-aware scene understanding.

Despite being an indispensable step towards the design of intelligent machines, affordance learning is a complex and highly integrated task. First, to understand how an object can be used by an agent requires reasoning about what it is and where it is located? Furthermore, it is necessary to know about the object’s geometry and pose e.g., an inverted cup cannot afford a ‘pouring’ action. Unlike traditional classification and detection tasks where each object takes a single label represented as one-hot encoding, a single object can simultaneously take multiple affordances, e.g., a bed is

both ‘sittable’ and ‘layable’. The object affordances are also dynamic and an intelligent agent should be able to consider both the prior knowledge and past experiences e.g., a cup may be first ‘graspable’, then ‘liftable’ and finally ‘pourable’. These challenges offer room to novel ideas and innovative solutions for visual affordance learning and functionality understanding.

Affordance learning has been reviewed from different perspectives in the literature; Bohg et al. [3] reviewed data-driven grasping tasks particularly for manipulation of objects and grasping, Yamanobe et al. [4] reviewed the affordance tasks to cover the grasping and manipulation of objects, Min et al. [5] introduced a general overview about the affordances and existing techniques; however, it is devoted to developmental robots and related tasks such as formalization of affordances. Up to the best of our knowledge, this survey is the first effort to review the literature from the perspective of ‘visual’ affordance and functionality understanding. Notably, other literature reviews cover affordance learning from the aspect of robotics perception, sensory-motor coordination or psychology and neuroscience [6], [5], [7], [4], [3]. However, affordance based reasoning is equally important for machine vision and visual scene understanding, as demonstrated by a growing activity in this area (see Figure 1).

Figure 2 shows the taxonomy of this survey and our scope which we are going to introduce in the following lines. It has two main parts: affordance-based techniques and functionality understanding methods. The rest of this survey is organized as follows: First, a comprehensive background to the area is provided along with the definition of specific terms frequently used in the visual affordance literature in Sec. 2. Afterwards, we provide the significance and challenges in Secs. 3 and 4 respectively. We then cover the research in visual affordance learning in Sec. 5 and categorize the methods according to specific sub-problems such as affordance detection, categorization, semantic labeling. Efforts to understand functions of different objects and tools are summarized in Sec.6. The main computer vision datasets

- M. Hassanin and M. Tahtali are with University of New South Wales (UNSW), Canberra, AU.  
E-mail: —
- S. Khan is with Data61-CSIRO and Australian National University (ANU), Canberra, AU.

Manuscript received –; revised –.



Fig. 1: Growth in the number of papers on visual affordance in computer and robot vision literature in the recent years (from 2014 to 2017).

with affordance and function annotations are listed in Sec. 7. Finally, we summarize open research problems in this area and mention new research directions in Sec. 8.

## 2 TERMINOLOGY AND BACKGROUND

The term affordance was coined by the psychologist Gibson to define the interactions between an actor and its environment [1], [8]. In his own words:

**Affordance:** “The affordances of the environment are what it offers the animal, what it provides or furnishes, either for good or ill. The word affordance implies the complementarity of the animal and the environment.”

—Gibson, 1979

He advocated that the target of computer vision should be to estimate possible interactions between a human, animal and environment in a scene rather than merely detecting the contents of a scene [9]. It was claimed at the time that perception is only a property of agent, but Gibson argued that the perception’s meaning is inherited from the environment [1]. Meanwhile, many researchers tried to find out the best definition of affordance, such as Turvey [10] who defined it as: “An affordance is a particular kind of disposition, one whose complement is a dispositional property of an organism”. He used the term ‘effectivity’ to denote the dispositional property of an organism (or an AI agent). Based on the agent’s effectivity and the environment affordance properties, an action is realized. Stoffregen [11] sought to clarify the relation between the affordances and actions and concluded that the affordance and actions are not identical in many aspects. In another study of Stoffregen, he critiqued the definition of Turvey that it neglected the relationship between animal and environment, therefore it is sub-optimal from the view of perception-action affordance [12]. Instead of agreeing with their predecessors on the definition of affordance, Sahin et al. [13] introduced a new formalization for affordances. They argued that the affordance has three main constituents: the agent, the environ-

ment and the observer. This concept allowed affordances to cover every aspect of robot control ranging from perception to planning i.e., first perceiving the object, then applying behavior to that object and finally the generation of desired effect.

Next, we define several key terms frequently used in the affordance learning literature, that we will consider for the purpose of this survey:

–**Visual Affordance:** The term visual affordance means extracting information related to affordance from an image or a video. Similar to other machine vision domains such as object and human activity recognition, it uses computer-vision techniques to perceive the affordance characteristics in visual media.

–**Functionality understanding:** It transcends the traditional tasks of object detection and segmentation in visual scene understanding and aims to understand the function of objects in a scene. For example, detecting the electricity plug that is required to charge the laptop or mobile phone as shown in Figure 3 (a).

–**Affordance learning:** In the context of complex and dynamic environments, robots need to learn what can be done and what cannot? In other words, affordance learning involves teaching the robot to learn the possible set of actions that can be performed in a given environment. It addresses the possible effects that arise through object-agent actions. It circulates (overlaps) three main factors, that are object, action and effect as shown in Figure 13.

–**Affordance detection:** Similar to object detection, it localizes and labels the affordance for scene objects. Different from conventional detection tasks, it targets only the salient objects that are most relevant for actions (instead of all objects in the scene). The goal of detecting these affordances is to predict the next action or recognize the function of some objects, therefore it selectively targets only significant objects. Formally, let  $X = \{x_1, x_2, \dots, x_n\}$  denotes the set of candidate bounding boxes with  $n$  instances and  $Y = \{y_1, y_2, \dots, y_n\}$  denotes the output label space where

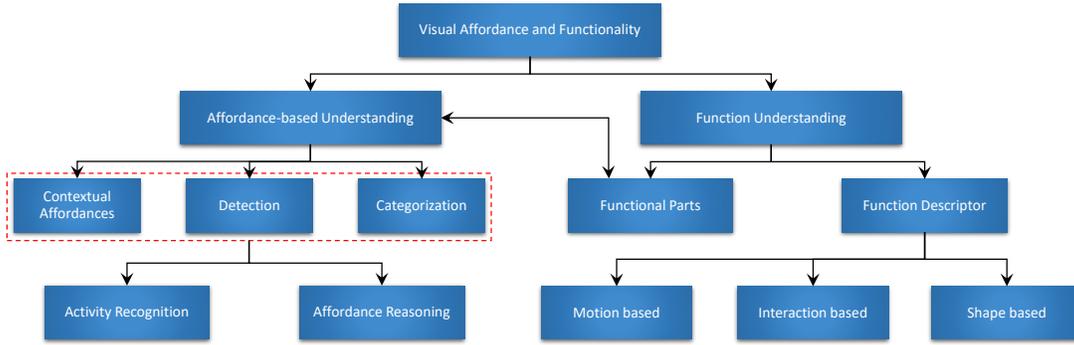


Fig. 2: Survey taxonomy shows the structure of methods which have been used to solve affordance issues.

$y_i = \{(r_i, l_i) : i \in [1, n]\}$  denotes a tuple consisting of instance location and affordance label respectively. Then, the affordance detection aims to find the function  $f : X \rightarrow Y$ .

–**Affordance categorization:** It means classifying the input images into a possible set of affordance categories. Often, this step is used as a precursor to affordance detection to make the process of localization and recognition more easier. Formally, let  $I$  denotes an input image and  $Y$  denotes its ground-truth affordance labels. The categorization task seeks to learn the optimal function  $f : I \rightarrow Y$ .

–**Affordance reasoning:** Affordance reasoning refers to more complex understanding of affordances which requires higher-order contextual modeling. The main purpose of such reasoning is to infer hidden variables like the amount of liquid inside a bottle or how much water can be pured into an empty bottle. Formally, let  $X = \{x_1, x_2, \dots, x_n\}$  denotes the input space where  $n$  is the number of object instances,  $C = \{c_1, c_2, \dots, c_n\}$  denotes the contextual space (physical attributes, material properties, neighborhood and semantic relationships) and  $Y = \{y_1, y_2, \dots, y_n\}$  denotes the output label space. The reasoning task is to find a function to map a relation such that  $f : (C, X) \rightarrow Y$ .

–**Affordance semantic labeling:** This task involves segmenting an image into a set of regions which are labeled with a semantically meaningful affordance category. Remarkably, this task assigns a category label to each pixel in a region of interest. Formally, let  $P = \{p_1, p_2, \dots, p_n\}$  denotes the set of  $n$  image pixels and  $Y = \{y_1, y_2, \dots, y_n\}$  denotes the output label space. The segmentation requires learning a function  $f : P \rightarrow Y$  to assign a label  $y$  to each image pixel.

–**Affordance-based activity recognition:** Objects bear possible actions of an agent which are represented as their affordances. Thus, recognizing the affordances is a crucial step towards complete activity recognition. This task aims to represent the activities in terms of affordances. Note that we refer to atomic single-person operation as an ‘action’, while actions performed by multiple people in a complex environment as an ‘activity’, e.g., a moving robot is performing an action while a group of marching robots are performing an activity. Formally, let  $X = \{x_1, x_2, \dots, x_n\}$  denotes the input space where  $n$  is the number of object instances,  $A = \{a_1, a_2, \dots, a_n\}$  denotes affordance space for each instance in the input space  $X$  and  $Y = \{y_1, y_2, \dots, y_n\}$  denotes the activities label space. Then, the activity recognition task is to learn optimal function  $f : (X, A) \rightarrow Y$ .

–**Social affordances:** Social affordance are a type of affordances that offer possible object-social actions. The term social imposes the human interaction through affordance learning. These social affordances may be positive such as sit on empty chair or negative (socially forbidden) to open a lady bag lying on the next chair. Formally, let  $X = \{x_1, x_2, \dots, x_n\}$  denotes the input space where  $n$  is the number of object instances,  $E = \{e_1, e_2, \dots, e_n\}$  denotes objects interactions space for each instance in the input space  $X$  and  $Y = \{y_1, y_2, \dots, y_n\}$  denotes the activities label space. The activity recognition task is to find the function  $f : (X, E) \rightarrow Y$ .

### 3 SIGNIFICANCE

Affordance learning is a crucial task in computer vision and human machine interaction. Affordance relates to objects, actions and effects, therefore addressing it (i mean affordance) benefits all these associated fields. The significance of affordance and functionality understanding is summarized below:

–**Anticipating and predicting future actions:** Because affordances represent the possible actions that can be performed on or with objects, learning helps in action anticipation in a given environment. The examples in the literature that used affordances to predict future actions are [15], [16], [17].

–**Recognizing agent’s activities:** The presence of an agent (e.g., a robot or a human) in the scene opens the possibility of its interactions with the surrounding environment. Recognizing its activities becomes an indispensable task to develop a complete understanding of scenes. Agent activities are highly dependent on the types of actions that are possible or more likely in a given environment. Several studies have targeted the problem of affordance-based human activity recognition [18], [19], [20].

–**Provides valid functionality of the objects:** Conventional object recognition task does not offer knowledge about its function and affordance, for example an occupied or a broken chair will still be recognized as a chair but a person cannot sit on it. In other words, recognizing what functions an object offers is a compulsory task to use it, particularly for the case of interactive robots [21].

–**Understanding social scene situations:** In social situations, forbidden actions such as a healthy person occupying a disabled seat, happen daily in our life. However, it needs

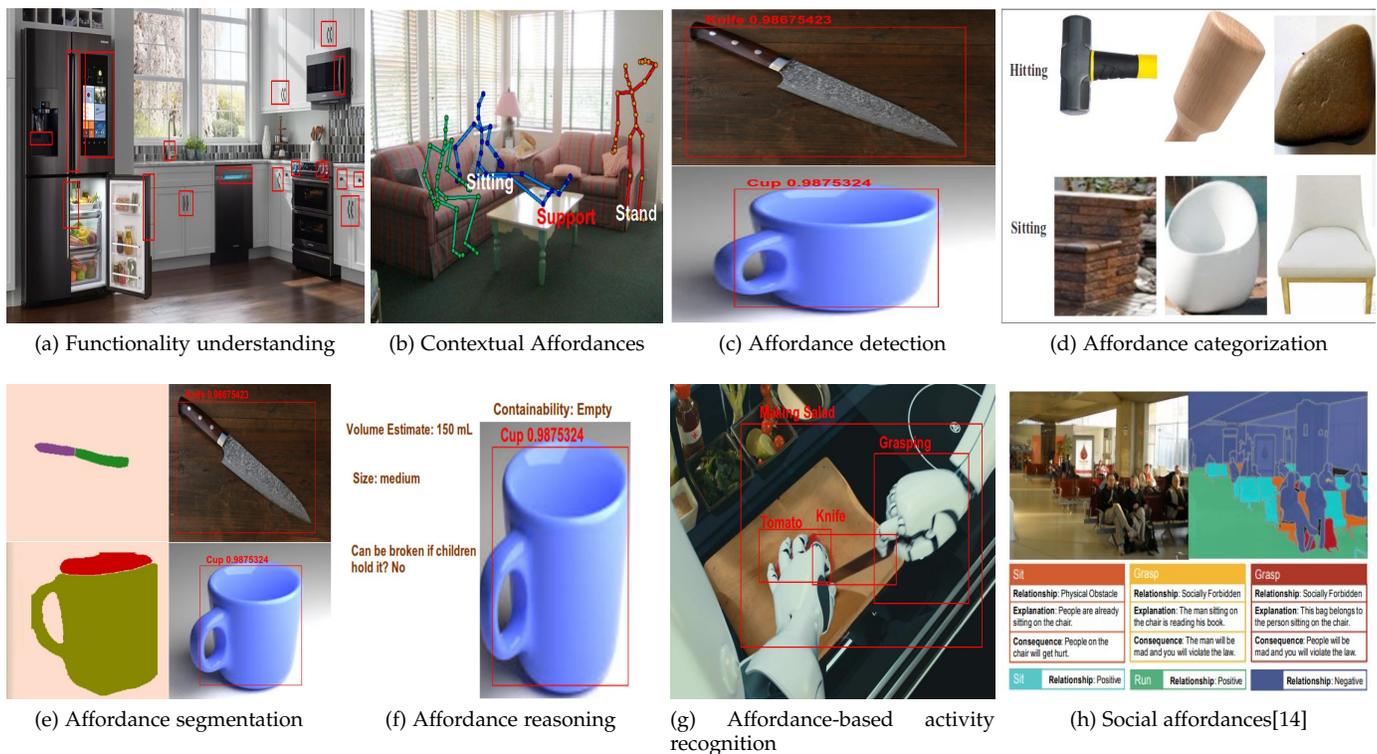


Fig. 3: (a),(b) Detecting functional parts to understand the scene usage and learn the affordances e.g. learning by observation. (c) Detecting the affordance objects and labels together using bounding boxes. (d),(e) affordance-based categorization and semantic labeling. (f) reasoning affordances from visual inputs such as how much water in the glass. (g), (h) given affordances, recognize the activities and social situations.

to be learned in case of visual recognition. Learning social affordances helps in understanding the whole scene [14].

–**Understanding the hidden values of the objects** The recognition of an object’s category does not provide the intuition about its value or significance. For instance, the task of knocking a nail would be normally done with a hammer. However, what if the hammer is not available? Learning affordances tells us that appropriate sized stones can also be used for the same task [22].

–**Detailed scene understanding** In traditional classification tasks, different kinds of pots are labeled with same label "pots". However, the usage of each pot could be different e.g., some could store rice while other could be suitable for vegetables. Therefore, categorizing objects according to their functionality is extremely important for detailed (or fine-grained) scene understanding [23], [21].

–**Affordance cues benefit object recognition** Extracting affordance cues from a scene provide the model with contextual information about objects that make the task of recognition easier. Many researchers used affordance values as contextual information to classify objects [24], [25]

#### 4 CHALLENGES

It is worth noting that detection is the first step in the process of affordance learning. It is a crucial task and it should be accomplished carefully. It also involves several challenges, such as scale, illumination, appearance and viewpoint variations. Furthermore, because affordance detection is a kind of

object detection, it inherits all the object detection problems as well. Some specific challenges are explained in detail below:

–**Illumination Conditions:** Illumination changes affect the final results of affordance learning because they change the quality of the image being processed. Robot will be stuck if the electricity switches off suddenly particularly indoors building. Likewise, outdoor robots may face the same situation if the sun light changes suddenly from sunny to cloudy or even slowly from day to night. It affects the image at the pixel level, therefore any changes in the illumination will change the final accuracy.

–**Occlusion and Clutter:** Occlusions often occur in dense images where one object obstructs parts of other objects. Background clutter adds difficulty to image recognition because of unordered objects. This sort of conditions leads the algorithm to incorrect predictions. Generally, this problem exists in multi-object images such as kitchens in indoor scenes and crowded scenes in the outdoor case. To allow the robot to perform its tasks easily, this problem needs to be addressed properly [26], [27], [28].

–**Viewpoint Variations:** Images acquired from different viewpoints can affect the performance of recognition algorithms. Therefore, the orientation of an object with respect to the camera defines the accuracy of the method. For robots, robustness against viewpoint and pose changes is one of the main effective factors in the process of affordance learning [29].

–**Scale Variations:** Scale of the object in terms of the size

is an important factor in the process of affordance learning especially with tools. For example, a fruit knife is generally small in comparison to a meat knife. Affordance detection algorithms should therefore be scale invariant to generalize well to unseen examples.

–**Deformation and intra-class variations:** Deformation or the different shapes of the same object is another criterion that needs to be treated to ensure reliable actions [30].

–**Single Object Multiple Labels (SOML):** The characteristics of affordance learning problems is different from traditional problems which have a single label for every instance or scene [31], [32]. For example, a knife in the kitchen has a grasping label in the hand and a cutting label in the cutter. Similarly, a cup of tea has a grasping label outside and a pouring label inside. Hence, visual affordance inherit all the problems of multi-label learning paradigm such as ranking, correlation, dependency and multi-label scheme representation.

–**Multiple Objects Multiple Labels (MOML):** In contrast to multi label problems, affordance cases have multiple labels at the object level rather than the complete scene level. It is closer to multi-instance multi-label (MIML) problems if we consider the objects as instances [33], [34], [35]. Intuitively, MOML inherits MIML challenges as well as single object multiple label difficulties, it has to address higher-order correlations between instance, dependency as well as the exponential size bottleneck.

–**Multi-source features:** The features required to perform affordance learning are diverse in nature and should come from complementary sources [36] (see Figure 4):

- 1) **Visual cues** such as color and texture.
- 2) **Physical properties** such as weight and volume, for instance, the chair’s visual feature is movable, but the weight, which is a physical property, does not allow sitting.
- 3) **Material properties** e.g., to assess the comfort of a chair.



Fig. 4: Multiple sources of features that can be fused together to create the affordance context.

–**Inter-object dependency:** The affordance of an object sometimes relies on other objects, particularly in service robots. For example, for the task of preparing a cup of tea, the robot has to boil water which depends on the electricity plug and pour hot water in a cup. This task will be complicated if one of those objects is unseen or occluded. [37].

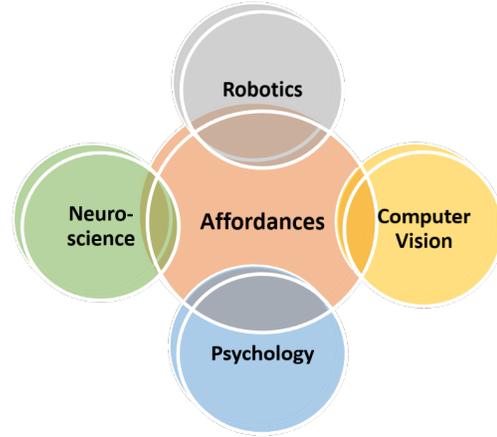


Fig. 5: Affordances and other fields

–**Human-Objects dependency:** The existence of a human in a scene makes it more complex in terms of actions, events and affordance learning in particular because the affordance is then dependent on the attributes of the human. For example, putting some fruit in the fridge depends on the height of the human.

## 5 WHAT IS A VISUAL AFFORDANCE?

The field of scene understanding aims to allow computers to be able to understand the environment and its contents. It has been studied from different perspectives: object recognition [38], [39], object detection [40], [41], [42], [43], [44], scene classification [45], [46], indoor scene understanding [47], [48] and so on. Much research has been conducted to address these problems. However, understanding possible interactions and developing a higher level reasoning about scenes has been less investigated. In other words, merely detecting the scene elements is not enough to make intelligent decisions but inferring complex interactions and dynamics in a scene: What the possible actions that can me made? For example, learn how to use the vacuum cleaner inside the kitchen? Which button should be pressed? Where is the electricity plug? What is required to name the chair a chair? It may not be sittable, occupied or broken [21]. To sum up, maximizing the benefits of scene understanding requires other associated factors such as affordance detection and reasoning to make best use of it.

The concept of affordance has initially been engaged strongly in the fields of perceptual psychology, cognitive psychology and environmental psychology. Affordances are often used for testing the capability of objects’ interactions because it has tight relations with different environment types, it addresses various research areas and it is linked to the object where ever it is. In detail, affordances differ according to the environment. For instance, the robot inside the kitchen should understand the affordance of tools [49], CAD objects such as chairs [50] and the functionalities that can be done such as the functionality of electricity plug to run the vacuum cleaner [51]. Moreover, if the human is involved in the scene, a new aspect, human-robot interaction and action recognition, should be addressed. After that it has been introduced to the fields of computer vision,

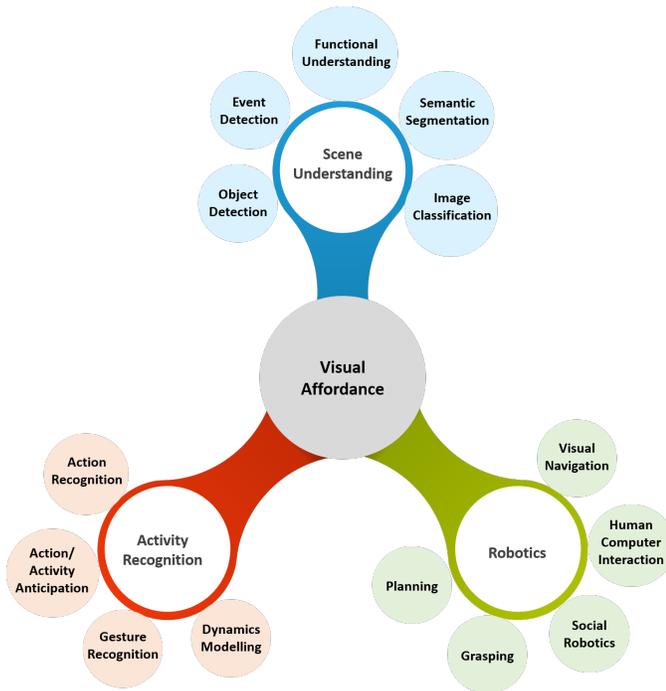


Fig. 6: Relation between affordances and visual fields

human-robot interaction and robotics [52]. The Affordance field is inter-related to different fields such developmental learning, robot manipulation and psychology. In this survey, the visual affordance will be summarized because of the lack of any existing literature review in this area.

Visual affordance is a branch of research that deals with the affordances as computer vision problem based on images and videos; and uses the machine learning methods, such as deep learning, to solve these challenges. It is tightly connected to various fields: action recognition [23], [16], scene understanding, grasping, human-robot interaction [53], [54] and function recognition [55], [51]. Therefore it overlaps with important computer vision problems such as image classification and action recognition; and affected by identical set of challenges such as illumination changes, occluded scenes, pose of objects and dynamic scenes [56]. Apart from that, the problem of affordance learning is more complicated because on top of all inherited difficulties from computer vision, it encounters additional challenges e.g. each object may have different labels/functionalities.

### 5.1 Affordance Detection

Recent advances in robotics and computer vision have paved the way for autonomous agents to impact every aspect of our life. To allow service robots do tasks, they need to understand the environment. For example, the robot cannot use the hammer without localizing the handle. Therefore, detecting the affordance of the scene, tool or a given image is a mandatory task for effective interaction and manipulation (see Figure 7). This task is, however, quite challenging, as explained below.

#### 5.1.1 Feature-engineering approaches

The affordance detection task deals with labeling and localization of particular parts that can afford certain actions.

It also involves reasoning about the functionality of a tool. Early work in affordance detection sought to recognize 3D CAD objects e.g., chairs based on the object functions. Stark and Bowyer [50] built the Generic Recognition Using Form and Function (GRUFF) system to recognize different objects according to the functionalities rather than the shapes. They used several functional primitives which were assigned to perform function-based recognition. Aldoma *et al.* [57] proposed a visual cue method to find affordances in the scene depending on the pose of the object. Their method depends on first recognizing the objects in a scene and then estimating its 3D pose. Finally, they learn the so-called 0-Order affordances which refers to hidden and unhidden affordances (see Table 1 for the investigated affordances).

Myers *et al.* [49] were the first to treat images from the perspective of geometry as traditional vision tasks. They detected affordance of tool parts based on its importance for robot vision. They introduced RGB-D data set with ground truth annotations whereas SVM was the learning algorithm. This study used a pixel-wise method to generate geometrical features to learn through it, but they used hand-crafted features. The idea behind using pixel-wise methods was brand new in this sense, however, it is complicated because the same pixel may share different affordance labels. They used two methods to train this model. Firstly, S-HMP (Superpixel-based Hierarchical Matching Pursuit) [66] to extract geometric features (depth, surface normals, principle curvatures, and shape-index and curvedness) and SVM as the main classifier. In addition, they used S-RF (Structured Random Forest) [67] to infer the affordance labels based on extracted decision trees particularly in real-time basis. They introduced seven affordance labels as shown in Table 1.

Herman *et al.* [58] sought to introduce a new method which depends on physical and visual features such as material, shape, size and weight to learn the affordances labels (as shown in Table 1). They collected their own data, which belongs to six categories: balls, books, boxes containers (mugs, bottles, and pitchers) shoes and towels, using a mobile robot Pioneer 3 DX. Based on these features, they used SVM and k-nn classifiers to test their method. Their work emphasized the concept that combining physical and visual attributes together enhances the affordance learning.

Grabner *et al.* [21] used the concept of functionality to recognize whether an object (e.g. a chair) is sittable or not. In other words, the chair allows the sitting affordance, but it may be used by another object. In order to detect whether the chair affords sitting or is "sittable", they used a human skeleton 3D model to test sitting in 3D models of chairs. Through different sitting human poses; and interaction between an actor and object, they detect not only a chair's affordances but also how to use it. Despite its effect in detection and reasoning, it required additional cues to improve the performance.

Moldovan *et al.* [26] proposed a novel method to estimate objects affordances in an occluded environment. They used the relational affordances concept to search for objects which can afford certain actions [68]. Additionally, they used Statistical Relational Learning (SRL) [69], [70] methods to model probability distributions that encode the relationships between objects. However, they generate these probabilities through an external knowledge base i.e., web

Affordance Label	References	Description	Examples	Method
rollable	[58], [59], [57]	indicate whether the object is rollable or not	roads, trolley	detection
containment	[60], [49], [57], [61], [62]	indicates contain-ability of an object	pots	detection
liquid-containment	[57]	indicates liquid contain-ability of an object	glasses, cups, mug	detection
unstable	[57]	indicates whether the object pose is stable after pushing or not	glass cups will be broken in case of pushing	detection
stackable-onto	[57]	indicates that the object can be stacked	mugs, pots	detection
sittable	[57], [59], [63], [64]	indicates whether the object can be used to sit or not	chairs, desks	detection, segmentation
grasp	[49], [58], [61], [59], [64], [65]	indicates the location of manipulation of flat tools	hammer, cups	detection
cut	[49], [60], [61], [65], [62]	indicates cutting	knife, penknife, key	detection
scoop	[49]	indicates curved surfaces tools	trowels, cookie scoop, gutter scoop	detection
pound	[49], [61]	indicates striking tools	hammer head	detection
support,place-on	[49], [60], [61]	indicates flat tools, helpers or support an agent	flat tools (turners,spatulas), place-on (tables, desks), agent support (walls) as walls	detection
wrap-grasp	[49], [61]	indicates the location of grasping of rounded tools like cups	the outside of a cup)	detection
pushable, pushable forward, pushable right, pushable left	[58], [64], [65]	indicates whether the object is push-able of an object	trolley, bike	detection,segmentation
liftable	[58], [64], [65]	indicates whether the object can be lifted or no	liftable chairs	detection, segmentation
dragable, pushable backward	[58], [64]	indicates whether the object can be dragged	desk, table	detection, segmentation
carryable	[58]	indicates whether the object can be carried	light-weight pots, balls	detection
traversable	[58]	indicates whether the object can be traversed	road, grass	detection
openable	[60], [65]	indicate whether the object can be opened	fridge, room, microwave, book, box	detection
pourable	[60]	indicates whether the object is pour-able	mug	detection
holdable	[60]	indicates whether the object can be hold	the outside of mug	detection
display, observe	[61], [59]	refers to display sources	TV, monitor screen	detection
engine	[61]	refers to tool's engine parts	drill engine	detection
hit	[61], [62]	refers to tools could be used to strike other objects.	racket head	detection
obstruct	[59]	indicates the locations of obstructer	wall	detection
break	[59]	indicates break-sensitive objects	glass cups	detection
pinch-pull	[59]	indicates objects that should be pulled with punch	knob	detection
hook-pull	[59]	indicates objects that should be pulled with hooking up	handle	detection
tip-push	[59]	indicates objects that perform actions after pushing	electricity buttons	detection
warmth	[59]	indicates warmth objects	fireplaces	detection
illumination	[59]	indicates light objects	lamps	detection
dry	[59]	indicates objects that absorb water	towels	detection
walk	[59], [63]	indicates places that allow walking	gardens	detection, segmentation
lyable	[63]	refers to long free space that allow person to lie down	bed	segmentation
reachable	[63]	refers to object in a scene that is reachable for a person to pick it	water bottle from the fridge	segmentation
movable	[64]	refers to objects that can be moved around	small objects like balls, mugs	segmentation

TABLE 1: Indoor affordance labels used to detect objects as in studies [57], [49], [58], [60], [61], [59], [63], [64], [65], [62]

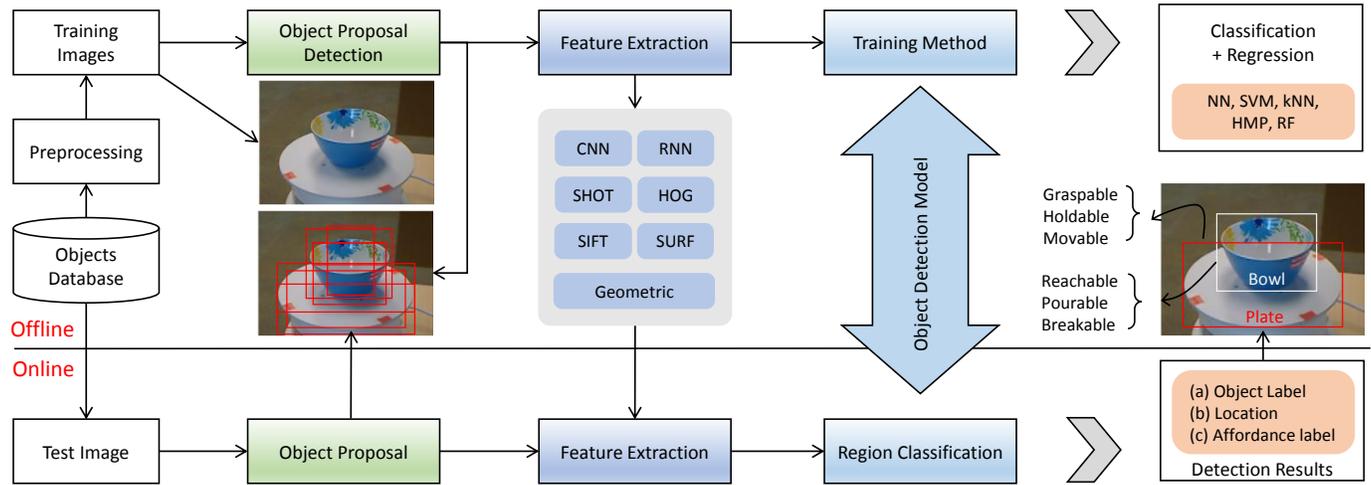


Fig. 7: Detection process through machine learning techniques.

images. Hassan and Dharmaratne [37] proposed an affordance detection method based on the object, human and the ambient environment. They used the objects attributes (physical, material, shape, etc), human attributes (poses) and object-to-object to train their scheme. Local feature have been extracted and used as inputs for classification based on Bayesian networks.

**5.1.2 Feature-learning approaches:**

Inspired by [49], Nguyen et al. [71] built their model to extract geometric deep features using Convolutional Neural Network (CNN) to detect affordances. They used an encoder-decoder architecture based on deep CNNs with multi-modal features (horizontal disparity, height and angle between pixel normals and inferred gravity) [72], [73]. It was demonstrated that the automatic feature learning performed on top of geometric features resulted in better performance compared to [49]. In contrast, semantic segmentation [74] is introduced to treat the affordance pixels. They tested it by real robot for grasping and they conducted their study on the UMD dataset [49]. Despite their significant enhancements, the data set images are simple i.e. it has no occlusion nor clutter. Sawatzky *et al.* [60], [75] proposed weakly supervised method to learn affordance detection using deep CNN based expectation maximization framework. This framework adequately handles weakly-labeled data at the image-level or key-point level annotations. They sought to fix the problem of affordance segmentation which needs special care because every pixel may be assigned to multiple labels. They learned deep features from the training data, however, they used human pose to represent the context. They introduced affordance RGBD datasets with rich contextual information. The annotated labels are shown in Table 1. However, the authors employed an additional step to update the parameters of CNN and estimate required masks for segmentation [75]. For this reason, they proposed the adaptive binarization threshold approach to get rid of that step and enhance the results. Nguyen *et al.* [61] conducted another study to detect the affordance labels in the images using a deep CNN architecture. Similar to Ye *et al.* [51], their work was inspired by popular CNN based object

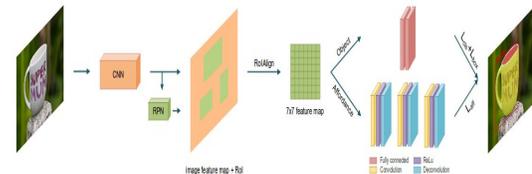


Fig. 8: The deep learning architecture that is proposed by [83]. The concept of RoIs have been used [84] to share the weights with the main CNN layers whereas the VGG is the feature extractor. In addition, they used deconvolution layer to refine the results of affordances.

detectors [55], [40], [44], Resultantly they treated affordance learning as an object detection problem. This approach starts with a candidate set of bounding box proposals for objects which is generated using [76], [77]. Although They tested through Faster R-CNN [77] and R-FCN [76] with various popular network architectures like VGG-16 [78], ResNet-51 and ResNet-101 [79], R-FCN outperformed the Faster R-CNN by a slight margin. This detection stage is followed by atrous convolution technique [80] to extract deep features which are finally followed by a Conditional Random Fields (CRF) model. The CRF, as a post-processing mechanism, provides substantial improvements [80], [81]. The authors published a new RGB-D dataset called IIT-AFF (with nine affordance labels as given in Table 1) to prove the efficacy of their method. The collected images in the dataset have good quality i.e., there do not exist many occlusions, cluttered regions and low resolution images. The authors in [61] used most recent architectures in deep learning, however, it resulted in high computational complexity due to high number of parameters in these models [82]. Do *et al.* [83] followed the two previous studies of the same group [61], [71] with a new effort to build an end-to-end deep learning architecture. Notably, the end-to-end model learning concept has recently predominated the recognition techniques, which has led to algorithms that can perform model training

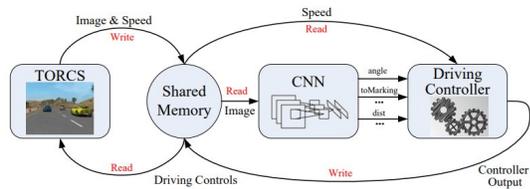


Fig. 9: **The deep driving architecture** [86]. Given an image from TORCS [88], the CNN extracts 13 indicators that will be fused along with the car speed to reason the proper command for this case

Affordance Indicator	Description or Function
angle	angle between the car and road's tangent
<b>When driving in the lane called "in lane system"</b>	
to_Marking_LL	marking distance between left & left lane
to_Marking_ML	marking distance between it & left lane
to_Marking_MR	marking distance between lane & right lane
to_Marking_RR	marking distance between right & right lane
dist_LL	distance between left lane & preceding car
dist_MM	distance between its lane & preceding car
dist_RR	distance from right lane to the preceding car
<b>When driving in the lane called "in marking system"</b>	
to_Marking_L	distance to the left lane marking
to_Marking_M	distance to the central lane marking
to_Marking_R	distance to the right lane marking
dist_L	distance to the preceding if it is in the left
dist_R	distance to the preceding if it is in the right

TABLE 2: The outdoor affordance labels that are proposed by Chen et al. [86]

in a single framework [41], [40]. To elaborate further, these systems detect the objects and their affordances in a single stage instead of multiple isolated and disintegrated steps. Hence, they reduce the training time and lead to better performances. Although they mainly followed the same strategy of [84], [85] as shown in Figure 8, they added new components like deconvolutional layers and robust resizing strategies to handle the multiple affordance classes problem. They relied on the affordance labels of [61] and tested their approach on two datasets: UMD dataset [49] and IIT-AFF [61]. In contrast to the above mentioned methods, Chen *et al.* [86] proposed a new out of the box idea to utilize affordance learning to reason about autonomous driving actions. They trained deep Convolutional Neural Network (CNN) on the KITTI dataset [87] and twelve recorded hours of video game [88] (see Figure 9). The authors introduced 13 affordance indicators as shown in Table 2 while, learning these indicators depend on the lane of that car and its perception. These learned affordances are tested to predict the right action. Eventually, a detailed comparison between affordance detection techniques is presented in Table 3.

## 5.2 Affordance-Semantic Labeling

The affordance semantic labeling task involves assigning pixel-level affordance category labels to relevant regions in an image. This problem combines segmentation and detection tasks and is relatively more challenging. It requires local and global contextual modeling for accurate pixel-level predictions. Affordance labeling is highly useful for precisely

locating where appropriate actions can be performed in a given scene. Figure 10 shows the most important steps to do segmentation.

Inspired by the semantic segmentation framework proposed by Eigen and Fergus [89] for generic object categories, Roy *et al.* proposed a multi-scale method to segment semantically meaningful affordances through CNN [63]. Given a RGB image, their model predicted three types of information: 1- the depth map, 2- surface normals, 3- labels for semantic segmentation. Thereafter, the outputs were merged together in a CNN network to predict the affordance maps. The experiments were performed on the NYU Depth dataset which consists of real-world indoor images [90]. This dataset was extended with affordance ground-truth annotations. The authors suggested five affordance labels as shown in Table 1. Although this method introduced a new feature encoding hierarchy based on intermediate semantic segmentation, it does not address the problem of multi-label affordances which is conflicting with the concept of segmentation. This is due to the reason that semantic segmentation aims to assign a single label to all pixels belonging to an object, while each object usually has multiple affordance labels e.g. a knife has both *cut* and *handle* affordances. Likewise, Kim *et al.* [64] performed affordance segmentation by using surface geometry features (e.g. linearity, normal and occupancy) of RGB-D images. The authors suggested six affordances as shown in Table 1.

Most recently, Luddecke and Florentin [59] proposed a new method to label affordances in RGB images using a refined version of Residual CNN [79] which was inspired by the work of Piheiro *et al.* [91]. As a major novelty, they developed a new cost function to handle multiple affordances in case of incomplete data. In a similar work, Roy and Todorovic [63] used the concept of action maps which predict the ability of users to do actions at various locations [92], [93]. This results in pixel-wise affordance segmentations given RGB images. Different from previous works, the authors introduce two original concepts:

- **Object Parts** – These are used to detect relevant segments of an object e.g., the surface of a table is important for placement while table legs are useless as shown in Figure 11. As a consequence, they can train on the only data set that supports object parts i.e. ADE20K [94].
- **Transfer Table** – It is a manual look-up table to map between object labels and affordance labels. In order to cover the affordance parts, the authors specifically suggested fifteen labels as shown in Table 1.

Apart from that, Table 4 compares between the recent methods in the literature.

## 5.3 Affordance as a Context

Affordances are inter-linked with both physical and semantic characteristics of an object and a scene. Object affordances can provide useful clues about object properties such as their category, location and function. In this section, we describe research efforts that aim to use affordance relationships as a context for other associated tasks such as action recognition, object detection and gesture recognition.

	Object Features			Features Extraction		Evaluation			Training			Model	
	2D	3D	Multimodal	Feature Learning	Hand-crafted	Real Robot	Simulation	Benchmark	Supervised	Unsupervised	Weakly Supervised	Statistical/ Mathematical	Neural Net
Stark et al. [50]		✓			✓	✓			✓			✓	
Aldoma et al. [57]		✓			✓				✓			✓	
Myers et al. [49]		✓			✓					✓		✓	
Herman et al. [58]	✓				✓				✓			✓	
Moldovan et al. [26]	✓				✓				✓			✓	
Nguyen et al. [71]			✓	✓	✓	✓	✓		✓			✓	✓
Sawatzky et al. [60]		✓		✓	✓	✓	✓		✓		✓	✓	✓
Nguyen et al. [61]		✓		✓	✓	✓	✓		✓			✓	✓
Do et al. [83]		✓		✓	✓	✓	✓		✓			✓	✓
Chen et al. [86]		✓		✓	✓		✓		✓			✓	✓

TABLE 3: Comparison between affordance detection methods

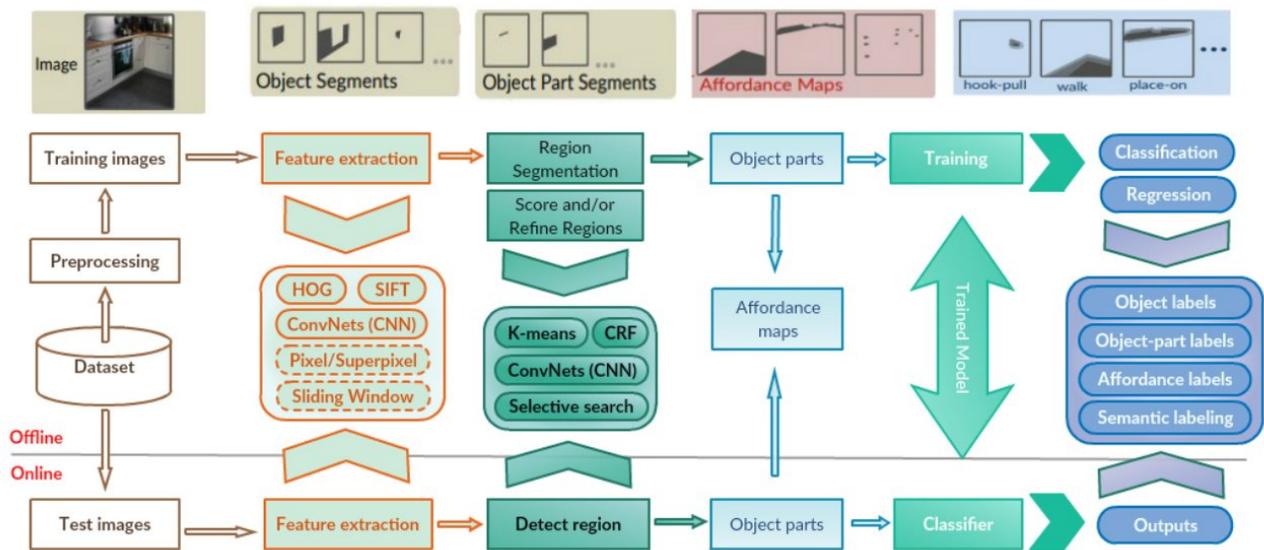


Fig. 10: affordance-based segmentation process diagram.

	Input		Features		Evaluation		Training			Model	
	2D	3D	Feature learning	Hand crafted	Real robot	Benchmark	Supervised	Unsupervised	Weakly Supervised	Mathematical	Neural
Roy et al. [63]	✓			✓	✓	✓	✓			✓	
Kim et al. [64]		✓		✓	✓		✓			✓	
Luddecke et al. [59]	✓		✓			✓	✓				✓
Nguyen et al. [61]		✓	✓			✓	✓				✓
Do et al. [83]		✓	✓		✓	✓	✓				✓

TABLE 4: Comparison between affordance-semantic labeling methods

Fig. 11: The part-segmentation methodology that is proposed by [59]. They employed segmentation for objects parts and train the the model to predict the affordances. In the same essence, they use manually look up (transfer) table to map between affordances labels and object labels

In an early work, Fitzpatrick et al. [95] proposed a method that allows a robot to learn how to segment the ob-

jects through imitation. The goal was to allow the humanoid robots to understand through acting. Similarly, through imitating the human actions on objects, a robot can learn object affordances e.g., whether a spherical shape is rollable or a cubic shape is slide-able. In addition to predicting affordances, it could interpret other’s actions. The intertwining of objects and actions, which is formally known by motor actions [96], to learn about objects affordances or segment objects was a promising at that time. Montesano *et al.* [97] sought to learn affordances through robot-environment interactions. Additionally, they used affordance labels such as eatable, movable and graspable as sensing capabilities. Local features e.g., color characteristics were used to detect the shape of the object which was eventually used to detect these affordances. Then the robot learned the model of grasping and used affordances through imitation and self observation. In the same context of motor actions, They proposed a probabilistic model based on Markov Chain Monte Carlo sampling. Inspired by Montesano *et al.* [98], Lopes *et al.* [99] proposed a probabilistic technique based on Bayesian Networks to learn affordances, and thereafter these learned affordances were used to recognize the demonstrations of an agent and learn the given task. Based on the affordances (e.g. tappable or graspable) and through self observation, the robot could relate the action with the resulting effects. However, learning affordances was employed for only single objects. Likewise, Ugur *et al.* [100] used self-interaction and self-observation to build their model. However, they provided behavioral parameters to enhance the accuracy. In contrast to previous methods [99], [97], the authors used unsupervised clustering to segment grasping through 300 trials whereas SVM was used to learn affordance labels. Varadarajan *et al.* [101], [102] developed a dataset to build knowledge ontologies similar to MIT ConceptNet [103] and KnowRob Semantic Map [104], but for household RGB-D images. They presented various affordance features such as grasp, material and structural. The authors built affordance filtrations which started by localizing the affordances, then they identified the entities related to that object and named it affordance duals. After that, they looked for all the entities that share the same affordance with that object. Moreover, they used a semantic part segmentation algorithm [105] for segmentation whereas modified the Levelberg Marquardt Algorithm (LMA) [106] using swarm PSO in order to recognize the objects.

Gupta *et al.* [24] modeled affordances in 3D indoor images to detect the workspace based on human poses while they used upright, lay down, reach sit as the human poses. Inspired by the idea of Gibson which says that the recognition of objects based on the function is better than the visual appearance [8], Castellini *et al.* [107] proposed using affordances as visual features and motor features, which are defined by kinematic features of the hand when grasping (e.g. time and instance of contact), to enhance the accuracy of recognition. Additionally, these features were defined as human-hand poses while grasping an object. The authors introduced the CONTACT VMGdb dataset which has visual features and kinematics content of grasping in various illumination conditions. They focused their learning on five affordance labels (cylindric power, flat, pinch, spherical and tripodal) along with various objects for grasping. Moldovan

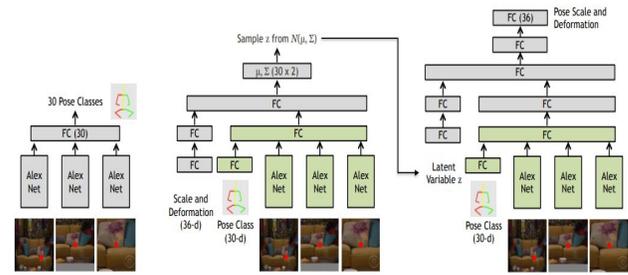


Fig. 12: The deep learning architecture that is proposed by [110]. The three parts of images from left to right as following: 1- Classification network 2- Variational AutoEncoder (VAE) 3- VAE decoder. VAE encoders and decoders share the weights in layers which have green color

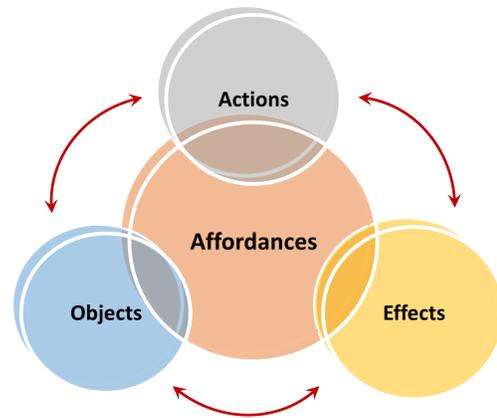


Fig. 13: The relations between objects, actions and effects

*et al.* [108] extended the model of [109], [97] that considers three related concepts: actions, objects and effects as shown in figure 13. They used spatial relations, which were defined by distance between objects and affordances to tackle the problem of multiple-objects manipulation. They used probabilistic programming with logical rules and probabilities to build their model. On the contrary, Lopes *et al.* [109] used object affordances as priori information to enhance gesture recognition and reduce ambiguities based on motor terms. Recently, Wang *et al.* [110] used Variational-Auto Encoders (VAE) [111] to build their model to predict the affordance poses. Based on the location, the algorithm classifies it into one of the 30 pose classes. Thereafter, it uses a VAE to extract the deformation of this pose. They proved that the nonexistent poses can be predicted using ConvNets and VAE. In addition, the authors claimed that deep learning based efforts are quite few in this area due to the non-existence of a big dataset for learning affordances. For this reason, they build a large-scale dataset from sitcoms to learn affordances.

Sun *et al.* [112] presented a method to learn affordances from object’s interactions. Given a video with labeled sequenced frames, their approach produces interactive motion models between pairs of objects and then represents it in a

Bayesian network with human actions. Ruiz and Cuevas [113], [114] proposed using bisector surfaces and developed them to include rational weighting, provenance vectors and affordance key points. They aimed to study the affordance locations for objects over one another (e.g. man riding bike, placing pot in the kitchen shelf or where the kids ride bike in a flat) and simultaneously detect unseen views of the object. Thermos *et al.* [65] presented a deep learning paradigm to investigate the problem of sensorimotor 3D object recognition. They employed a biological neural network architecture (VGG-16 network [78]) to fuse multiple evidence sources to learn affordances. They used fifteen affordance types as shown in Table 1. Some of these types describe complex affordances like "squeeze" and continuous like "write". Moreover, they introduced a new RGB-D dataset which has human-interaction and affordance types to test their method. All of the above mentioned methods have been summarized in Table 5 to show what has been done and what should be solved.

	Input		Features		Evaluation		Training			Model	
	2D	3D	Feature learning	Handcrafted	Real Robot	Benchmark	Supervised	Unsupervised	Semi-supervised	Mathematical	Neural
Fitzpatrick <i>et al.</i> [95]	✓			✓	✓			✓		✓	
Lopes <i>et al.</i> [99]	✓			✓	✓				✓	✓	
Montesano <i>et al.</i> [97]	✓			✓	✓		✓	✓		✓	
Ugur <i>et al.</i> [100]	✓			✓	✓			✓		✓	
Varadarajan <i>et al.</i> [101], [102]		✓		✓		✓				✓	
Gupta <i>et al.</i> [24]		✓		✓		✓				✓	
Castellini <i>et al.</i> [107]		✓		✓		✓				✓	
Grabner <i>et al.</i> [21]		✓		✓		✓				✓	
Lopes <i>et al.</i> [109]	✓			✓	✓				✓	✓	
Wang <i>et al.</i> [110]	✓		✓			✓		✓			✓
Sun <i>et al.</i> [112]	✓			✓		✓				✓	
Thermos <i>et al.</i> [65]		✓	✓			✓					✓

TABLE 5: Comparison between contextual affordance recognition methods

### 5.4 Affordance Categorization

The affordance categorization task aims to tag an image with the relevant set of affordance labels. To this end, a general approach is to represent an image in the form of discriminative features and employ a classifier to assign affordance labels. This task is relatively simple compared to affordance detection and segmentation, which also localize affordance categories (see Figure ??).

Varadarajan and Vincze [115] proposed hybrid parallel architecture of deep learning and suggestive activation (PDLA) to overcome the problems of deep learning like uni-modality and serialization in order to categorize the objects based on the affordance features. The authors extracted the semantic features of affordances as proposed in [102] as well as the structural and material features to enhance the recognition results. The Washington RGB-D dataset was used to test the efficacy of their model. Sun *et al.* employed object categorization as an intermediate step to infer the affordances more correctly [116]. They developed a

visual category-based affordance model, which encoded the relationships among visual features, learned categories and affordances in probabilistic form. Such a probabilistic modeling allows knowledge transfer and enhance accuracy, especially when limited annotated data is available. Additionally, they addressed the problem of incremental learning of affordances. The authors suggested seven object categories and six affordance labels as shown in Table 1. However, this study has been devoted to indoor buildings and it was not enough to treat all the cases such as defining that the table is movable requires knowing its physical attributes. In [21], the authors have categorized the images based on functions to enhance the performance of detection. The concept of bootstrapping, which uses the past knowledge to accelerate the learning process, has been applied to affordance learning [62], [117]. Schoeler *et al.* [62] proposed to infer any possible usage of a tool even if that usage is possible through another main tool e.g., using stone as hammer or using helmet instead of water cup. They sought to divide the tools into six functional categories (contain, cut, hit, hook, poke and sieve). Afterwards, an ontology of tool functions was created to allow a deeper understanding by exploring their usage in absence of main tool. The authors developed their main algorithm in three steps: **(1)** Part Segmentation through Constrained Planar Cuts (CPC) algorithm [118], **(2)** Extraction of part-based visual features via a Signature of Histograms of Orientations (SHOT) descriptor [119], [120], **(3)** The pose of individual parts with respect to each other is encoded via a Pose Signature that models the alignment and attachment between parts. Although the authors introduced a new idea that was thoroughly tested on a 3D synthetic dataset. However, their method gets confused when many tools exist to perform the same action. Also, the selected tool for some action may not be well-suited due to its size or shape. To avoid these problems they need to find the best correspondences between tools and objects.

Given basic affordances, Ugur *et al.* [117] proposed a bootstrapping method to learn the complex affordances through the relational affordances or so-called paired-object affordances. They evaluated their approach using a real robot on various object shapes like boxes, spheres and cylinders. Additionally they trained the robot to perform actions such as side-poke and stack. Similarly, Fichtl *et al.* [62] addressed the problem of "related affordances", which denotes the cases where affordances are related together e.g. open kitchen door to fetch the glass from inside. They used pose and size as the visual features to build their model. Abelha *et al.* [121], [122], [123] proposed methods to identify tools and their substitutes based on matching superquadrics [124] of these tools in point cloud data. Mar *et al.* [125], [126] used learning by exploration as the methodology to train the iCub robot. In the latter work [126], they categorized the objects before using parallel Self Organizing Maps (SOM). Following their study [127], Pieropan *et al.* used the same method to learn affordances rather than functionalities [128]. Kjellström *et al.* [23] proposed functional categorization to learn object affordances through human demonstration. They used a Conditional Random Field (CRF) and factorial conditional random field to train the model and infer the object's affordances and actions even when they are novel. Overall, Table 6 presents detailed

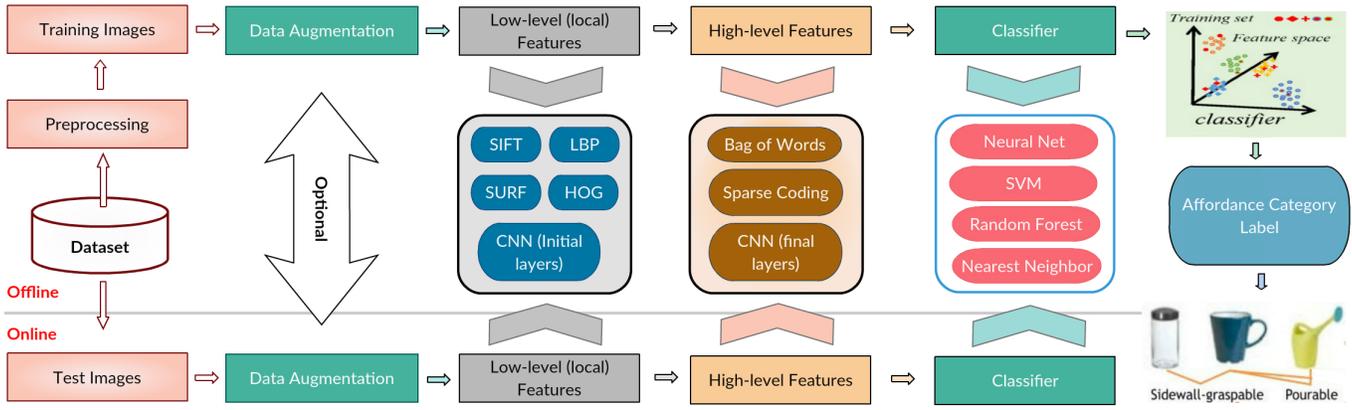


Fig. 14: Affordance classification process through machine learning techniques.

comparison between affordance classification methods.

### 5.5 Affordance-based activity recognition

Enabling seamless human-robot interaction is a crucial step towards ubiquitous use of personal robots. It is a multi-disciplinary research area which overlaps with Robotics, Human-Computer Interaction, Cognitive Science, Artificial Intelligence, Action Recognition and Affordance Prediction. To understand the interaction between robot and its environment, the affordance of objects should be predicted which can be a very useful cue for activity recognition as shown in Figure 15.

Due to the close relationship between affordances, actions and humans, its use in recognizing human action and hence activities has been investigated. In this section, the relation between affordances and actions will be reviewed. Koppula *et al.* [16] learned the human activities from RGB-D videos considering object affordances. They structured their scheme using a graphical representation where the nodes represented the sub-activities and objects; and the edges referred to the object’s affordances and the relations between human and objects. Additionally, they learned their model through a structural support vector machine algorithm. In a recent study, Koppula *et al.* [15] presented a modified version of their previous work [16] where they used a Conditional Random Field (CRF) to represent the model. They merged CRF structure with object affordances and sub-activities to form so-called Anticipatory Temporal CRF (ATCRF). Following these studies [16], [15], Jain *et al.* merged spatial-temporal graph with a Recurrent Neural Network (RNN) to address problems in graphical models [17]. They trained and tested different kinds of spatial-temporal cases like motion modeling, human activity prediction and action anticipation. Qi *et al.* [18] represented affordances, human actions and interacting objects in a Spatial-Temporal And-Or graph (ST-AOG) to predict human activities in RGB-D videos. They built the model in two main stages: video parsing and activity prediction. The parsing is done using segmentation by a dynamic programming approach and later label refinement using Gibbs sampling. For activity prediction, it depended on an Earley parser [19] to predict sub-activities and all the learned cues (parsed graph and sub-activities) to estimate human activity. Vu

*et al.* [20] challenged that various scenes under the same category have similar functional features. They described scenes in terms of functionalities to predict actions from static images. Dutta and Zielinska [129] presented a novel method to predict next action based on object affordances and human interaction. They represented the model in a spatio-temporal based probabilistic state automaton. The generated motion trajectory was used to build action heat maps that led to infer next actions. Shu *et al.* [54] proposed learning social affordances from human to human interactions. They represented their model in a graphical scheme that has nodes as subevents/subgoals. They provided a RGB-D video dataset (HHOI) to describe human to human interactions. Given a RGB-D video, Shu *et al.* [130] learned social affordance grammars and then represented them as a ST-AOG to perform motion modeling. To sum up, Table 7 compares between affordance-based activity recognition techniques.

### 5.6 High-level Affordance Reasoning

Affordances can be used as a tool to perform reasoning about more complex object properties and events in a scene. As an example, affordances have been used to infer the hidden properties e.g., What is inside a container? What are the intricate relationships between objects? Or to answer complex questions about a scene. In this section, we cover research works which perform high-level reasoning based on affordances (see Figure 17).

Zhu *et al.* [131] proposed the first study to discuss the visual reasoning of affordances. They developed Knowledge base (KB) that represents the object along with other nodes which describe attributes (Visual, Physical and Categorical) or affordances to infer the affordance labels, human poses or relative locations. They learned their model using Markov Logic Network (MLN) [132] whereas zero-shot learning [133] has been used to predict affordances for novel objects. However, their approach assumed that the affordance semantics and attributes are given in advance and use static models. Chao *et al.* proposed this study [134]. They modeled the affordance semantics problems in the form of action-object pairs as connected verb-noun nodes in WordNet [135] or encoding the plausibility of a matrix such as this study [136]. They used novel statistical methods

	Input		Features		Evaluation			Training				Abstraction	
	2D	3D	Feature learning	Handcrafted	Real Robot	Simulation	Benchmark	Supervised	Unsupervised	Self-supervised	Semi-supervised	Mathematical	Neural
Varadarajan <i>et al.</i> [115]		✓	✓				✓	✓					✓
Sun <i>et al.</i> [116]	✓			✓	✓			✓				✓	
Schoeler <i>et al.</i> [62]		✓		✓					✓			✓	
Ugur <i>et al.</i> [117]		✓		✓	✓				✓			✓	
Fichtl <i>et al.</i> [62]		✓		✓		✓					✓	✓	
Abelha <i>et al.</i> [121], [122], [123]		✓		✓			✓	✓				✓	
Mar <i>et al.</i> [125]			✓	✓	✓				✓			✓	
Mar <i>et al.</i> [126]		✓		✓		✓				✓		✓	
Pieropan <i>et al.</i> [127]		✓		✓				✓	✓			✓	✓
Kjellström <i>et al.</i> [23]	✓			✓			✓	✓				✓	

TABLE 6: Comparison between affordance classification methods.

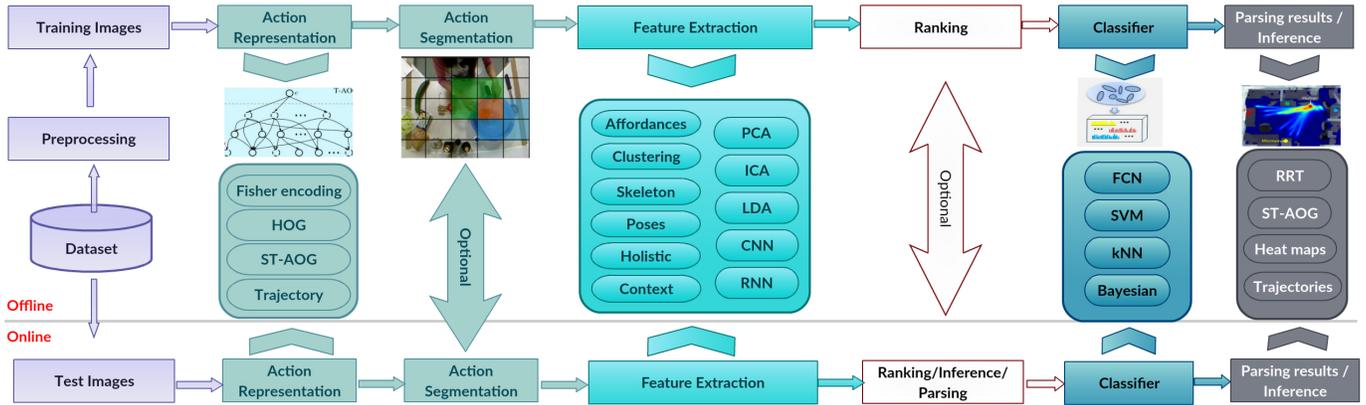


Fig. 15: Affordance-based activity recognition process.

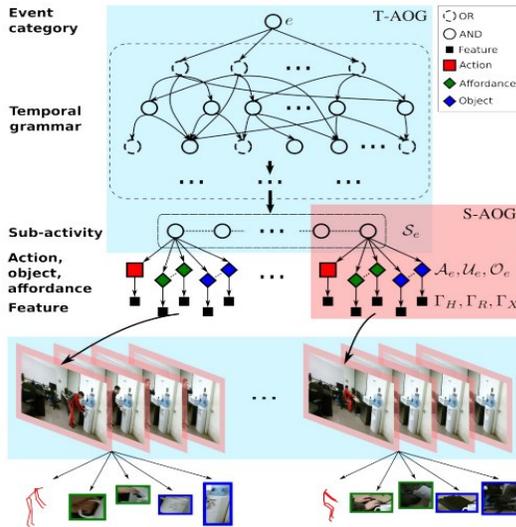


Fig. 16: ST-AOG model has two main parts: T-AOG on top represents the activity as the root, S-AOG represents the subactivities nodes which encode actions, object affordances and interactions as the state context.[18]

	Input		Features		Evaluation		Training		Model	
	2D	3D	Feature learning	Handcrafted	Real Robot	Benchmark	Supervised	Weakly supervised	Mathematical	Neural
Koppula <i>et al.</i> [16]		✓		✓	✓	✓	✓		✓	
Koppula <i>et al.</i> [15]		✓		✓		✓	✓		✓	
Jain <i>et al.</i> [17]		✓	✓	✓	✓	✓	✓			✓
Qi <i>et al.</i> [18]		✓	✓	✓	✓	✓	✓			✓
Vu <i>et al.</i> [20]	✓		✓	✓	✓	✓	✓		✓	
Dutta and Zielinska [129]		✓		✓	✓	✓	✓		✓	
Shu <i>et al.</i> [54]		✓		✓	✓	✓	✓		✓	
Shu <i>et al.</i> [130]		✓		✓	✓	✓	✓		✓	

TABLE 7: Comparison between affordance-based activity recognition methods

like co-occurrences to infer about the affordances and give the best description (verb) for this object(noun). Regarding reasoning about the containers' contents, GÄijler *et al.* introduced the visual-tactile method to infer what is inside the container. A kinect camera was used to capture the object and then deformation model was detected before using tactile signals for reasoning. They used a three-fingered Schunk Dextrous Hand to grasp and squeeze the container to check whether it is full or empty. PCA was the extractor and kNN

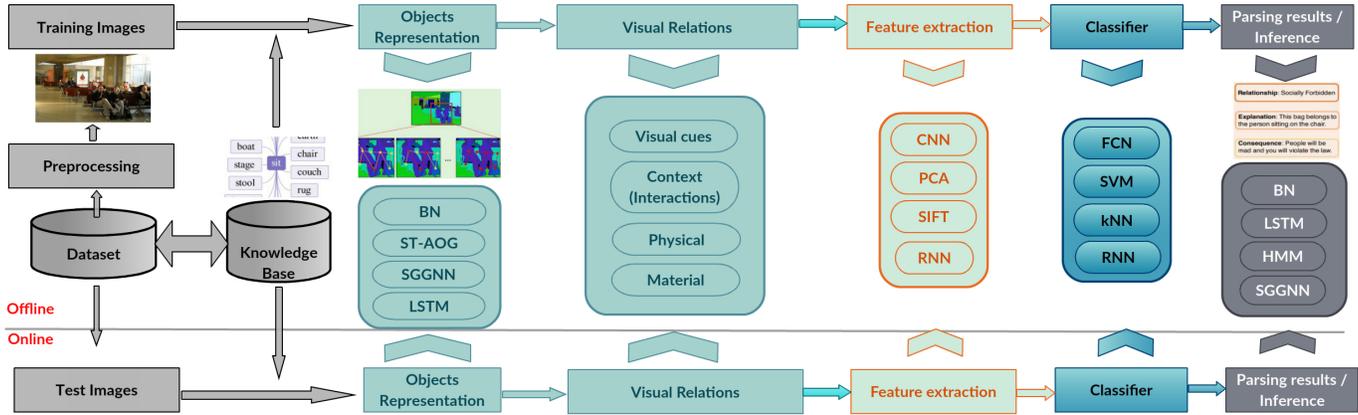


Fig. 17: Affordance reasoning process.

and SVM were the classifiers of the model. However, the authors assumed the presence of the 3D object model to perform the learning.

Yu et al. [137] proposed a physics-based method to reason liquid containability affordances through two steps: 1) the best filling direction 2) considering the direction, transfer the liquid outside to estimate the containability of an object. They used fluid dynamics in 3D space to simulate human motion of liquid transfer to test their approach. Although this study tried to predict containability based on the visual images, it did not use the human factors in the scene which makes the reasoning easier. Further, they did not relate the material attributes to the reasoning e.g. any curved tool can contain liquid even if it is made of this paper. Seeking to overcome the human factors in the scene problem, Wang et al. [138] developed their approach as an enhancement of this paper [137]. Given a 3D scene and considering the compatibility among container, containee and pose, they learn their model to reason the best pose and container to do the task of transferring the liquid. They provided a RGB-D dataset and they used SVM to train the approach. The process of filling the containers is the same process of this study [137]; they voxelize the object and simulate the filling inside its space. In the same context, Mottaghi et al. [139] developed an approach to reason about the affordance of liquids inside the container (the volume, amount of liquid) and predict its behavior. They introduced Containers Of liQuid contEnt (CODE) dataset of RGB images along with 3D CAD models. They used deep learning in the form of CNN and RNN to learn the containability affordance based on contextual cues. This method depended on visual features. Phillips et al. [140] introduced a method to detect, localize and segment pose estimation of transparent objects like glasses; and they provided an annotated dataset of transparent objects. Liang et al. studied the human cognition of the containers to infer its affordances (object’s containment and number of objects that can be contained inside) [141]. In a recent study, Liang et al. [142] inferred containment affordances and relations in RGB-D videos over time. For example, the fridge contains the eggs carton which contains the eggs. They used the human actions (move-in, move-out, no-change and paranormal-change) to draw containment graphs based on spatial-temporal relations; that is, the ac-

tion used to detect containment objects. They introduced RGB-D videos dataset and they developed probabilistic dynamic programming to optimize the containment graphs. Zhu et al. [143] used physics-based simulation to infer the forces and pressures for the different body parts while sitting on a chair. They predicted the object affordances through human utilities while sitting e.g. comfortable and lazy. Krunic et al. introduced a model to include verbal information to link between utterances and objects through inferring the context between words, actions and its outcomes. Zhu et al. [144] built a knowledge base (KB) to reason answers for image questions. They represented nodes in their model as attributes, affordance labels, scene categories or image features. In their most recent study, Chuang et al. presented a promising study to reason about the action-object affordances based on physical and social norms [14]. They annotated the ADE20k [94] dataset with affordance features, detailed explanation and potential consequences for every object. For instance, pouring water into a cup has explanation that it is improper to pour because the cup is full and consequence that you will make a mess in that place. They built the model using a Gated Graph Neural Network (GGNN) while Spatial version of GGNN has been employed to reason affordances. This study combined the social norms, physical features, visual attributes and situation parameters to infer the affordances and its relations whether positive or negative. To summarize, Table 8 gives more details about the used methods in the literature of reasoning.

## 6 FUNCTIONAL SCENE UNDERSTANDING

Zhu et al. [22] made a distinction between functionality and affordances. The problem of affordance particularly depends on detecting and classifying objects; learning their affordances; and performing more detailed understanding and reasoning based on the learned affordance properties. The problem of understanding the functions is related to affordance learning, but specifically aims to identify the tasks that can be performed with an object. In contrast, affordance learning reasons about object functions in the context of agent (animal or robot). Further more, some objects have both affordance parts and function parts e.g. the hammer

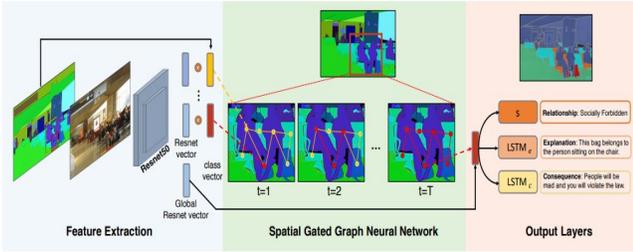


Fig. 18: The deep learning architecture that is proposed by [14]. ResNet is used as the feature extraction layer and the results are fed into a customized Spatial Gated Graph Neural Network (SGGNN) to represent the visual affordances through graphs. Hence, use this graph to reason the affordances, explanations and consequences

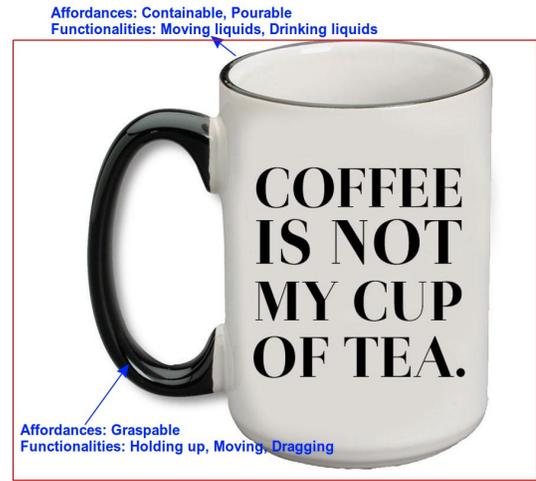


Fig. 19: The relation between functions and affordances is many to many.

Methods	Input		Features	Evaluation		Training		Model
	2D	3D		Simulation	Benchmark	Supervised	Unsupervised	
Zhu <i>et al.</i> [131]	✓		✓		✓		✓	
Chao <i>et al.</i> [134]	✓		✓		✓		✓	
Yu <i>et al.</i> [137]		✓	✓		✓		✓	
Zhu <i>et al.</i> [143]		✓	✓	✓		✓	✓	
Wang <i>et al.</i> [138]		✓	✓	✓		✓	✓	
Liang <i>et al.</i> [142]		✓	✓	✓		✓	✓	
Mottaghi <i>et al.</i> [139]	✓		✓	✓		✓	✓	
Phillips <i>et al.</i> [140]	✓		✓	✓		✓	✓	
Zhu <i>et al.</i> [144]		✓	✓	✓		✓	✓	
Chuang <i>et al.</i> [14]	✓		✓	✓		✓	✓	

TABLE 8: Comparison between affordance reasoning methods

has a head which is suitable for striking objects (function) and a handle from which it can be grasped (affordance). Some of these object parts have multiple affordances but a single function e.g. the *hammer handle* can be used push or pull objects so it has pushable and pullable affordances while the function is head support. In contrast, the *hammer claws* have multiple functions and a multiple affordances, i.e., claws can function as a lever and also can be used to pull out nails from timber but has the affordance of hooking, grasping and pushing. Thus, the relation between the affordance and the function is many to many as Figure 19 shows.

As scene understanding is an old problem, it has been studied from many perspectives. However, functional scene understanding has not been thoroughly investigated before. Although this problem is of high significance for robotics, a few research efforts have focused on functional scene understanding. For example, a robot can not clean the kitchen dynamically without understanding how to use taps or electricity plugs to run the vacuum. So that the functional scene understanding is crucial to the robots particularly cognitive robots. As Figure 20 shows, categorizing the images according to its purpose is meaningful in many cases than categorizing according to its appearance.

### 6.1 Function-parts recognition

Inspired by Gibson [8], the authors in [145], [146] looked for the relation between recognizing objects and its functions. In other words, they aim to recognize the objects according to their usage features instead of their visual properties. Rivlin *et al.* used the relations of objects (e.g. size and orientation) and parts to reason about the functionality [147]. Desai and Ramanan performed an interesting study to predict functional regions and functional landmarks in an image [148]. They used deformable part models (DPM) [43] and [149] as pose detector to extract spatial relations of objects followed by kNN to predict functional regions in the scene. For instance, the functional part of the vacuum cleaner is the power buttons. They targeted to detect 3D objects based on their functional parts, spatial relations and functional landmarks and tested their approach on the attributes of people dataset [150]. Zhao and Zhu proposed a stochastic method, Function-Geometry Appearance (FGA), to parse 3D scenes by combining features of the functionality, geometry and appearance [151]. They modeled the FGA through top-down/bottom-up hierarchy and used MCMC to build the algorithm and infer the parse tree (see Figure 21). The functional descriptor was composed of functional scene categories, groups, objects and parts and use AND-OR rules to understand the affordance and therefore the full scene.

Shiraki *et al.* [152] introduced the differentiation between main parts and subordinate parts (like hammer head and hammer handle) which is important to reason about object functions. Zhu *et al.* [22] used a simulation engine to perform affordance and functionality analysis. They differentiated between functional learning which was defined as the right location to do a task on the target object; and affordance learning which was defined as the best location to grasp the object depending on the tool type. Additionally, they fused physical features (e.g. forces, pressure and volume), human pose from imagined action, affordance features and functional features together to understand and infer tool’s usage. Ye *et al.* [51] addressed object function detection by introducing a novel approach through three main steps:



Fig. 20: Categories according to the functionality. The left side shows the traditional categorizing or object-labels categorization whereas the right side categorizes the objects in functional manner. For example, handle-graspable category includes the first four objects from the right side

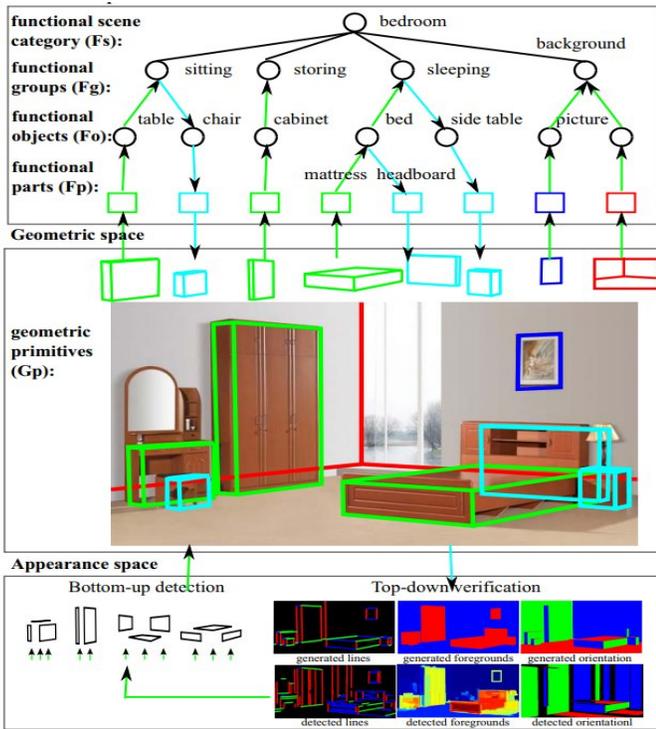


Fig. 21: Function-Geometry Appearance (FGA) use functional cues (categories, objects and parts), geometrical features (e.g. object 3D size and corners values ) and appearance parameters to enhance the prediction accuracy [151].

selective search to detect object proposals; feature extraction by VGG-F and VGG-S network and training the neural network with stochastic gradient descent to perform functional object detection. The authors inspired that problem, but it has the following limitations: (a) the accuracy has not exceeded 25%, (b) it requires long time to train because of using three different stages i.e. detecting, feature extraction and training. Since this problem is relatively new, it needs many further investigations to achieve a practical and successful baseline. In a different way of processing, Pechuk *et al.* structured their problem into function classification model in 3D images, where the category labels were defined to be e.g., “to sit” and “to use” [153]. Their scheme operates on a multi-level hierarchy of parts that have various functionality signatures. They named these parts as primitives and classified each part or group of them as functional part if they offer the same functionality. To link between

primitive and functional parts, they defined association, connection and mapping relationships. Kitano *et al.* [154] build a CNN to estimate the object functions based on their appearance. As [21], the usage of CAD models is repeated in [155] while Hinkle and Olson proposed learning method for robots that face novel (unknown) objects. Based on physical simulations of a falling sphere onto objects, the robot can classify the objects according to functionality or potential capabilities. They learned three functions: drinking vessel, table and sittable. Awaad *et al.* [156] proposed using functional affordances to serve social robots. They challenge that using functional affordances of objects represent the start point for socializing robots to achieve the goal. For example, using a mug instead of the glass, which is not available, can achieve the goal of drinking water. To learn functional affordances, they built an ontology for the objects usage and ranked them according to whether it is the primary function or a secondary function. Madry *et al.* proposed categorizing the 2D-3D images according to its functionality to reason about grasp planning [157]. They modeled their approach as a Bayesian Network and extended this probabilistic method [158] for reasoning. They used five functional categories to perform tasks as follows: hand-over, pouring, dish washing, playing and tool-use. Xie *et al.* [159] named functional objects as “dark matter” and all the functional objects in a scene as “dark energy”; and proposed a Bayesian framework to represent their model whereas Markov Chain Monte Carlo (MCMC) is used for inference. In contrast to the previous works, this method used the human trajectories in order to localize the functional objects from videos. In the context of functional categorization, Gall *et al.* proposed a novel method to solve the problem of unsupervised categorization of objects based on functionality. They inferred the functionality from motion of human-object interactions. Furthermore, they extracted the human poses through interaction, encoded a string of poses and measured the similarity between the actions with Levenshtein distance method [160]. Saponaro *et al.* [161] presented a geometric transformation method to the hand actions (push, grasp and tap) to predict physical actions and reason intentions of the human. They built their model based on the affordance relationships mentioned in Figure 13.

## 6.2 Function from Motion

Gupta *et al.* [162] used spatial and functional features to understand human actions in static images and videos. They argued that the functionality, which has been pointed out as

motion trajectories and human poses, and pose can interpret the scene in terms of action recognition. They built Bayesian Network model which fused functional and spatial data for object, action and object effect for recognition. Oh *et al.* [163] detected the functionality of moving objects in order to understand the scene context. They claimed that appearance features are not enough for some objects in traffic scenes (e.g. parks and way lanes) because all of them are similar without much difference. Turek *et al.* [164] presented a novel method to categorize a scene according to the motion of objects. The categories are function-based such as sidewalks or parking areas. Local descriptors to model object properties, behavioral relationships and temporal features were used to discriminate between functional regions. Inspired by [164], Zen *et al.* used the same steps but additionally included the semantics of the moving object to get a more informative descriptor [165]. Additionally, they tested their model on a traffic dataset [166]. Rhinehart and Kitani [167] analyzed egocentric videos to extract functional descriptors to learn "Action Maps" (six actions) [92] which can tell a user about how to perform activities in novel scenes. Pirk *et al.* [168] used scanners to track interaction parts and built the functionality descriptor by analyzing motion of these parts interaction with an object e.g. analyzing an agent tries to pour water into a mug.

### 6.3 Function from Interactions

Gupta *et al.* [24] modeled human-scene interaction with a functional descriptor in terms of 3D human poses to predict various poses of the human and they named it as the human workspace. They used a 3D cuboidal model to fit objects and predicted the subjective affordances. In [169], Delaitre *et al.* functional descriptors were proposed to allow a better scene understanding. The main idea is to extract functional descriptors from human-object interaction aiming to achieve better object recognition results. They extended an existing dataset [170] to evaluate their proposed model. Similarly, Fouhey *et al.* used the idea of [162] to predict the affordance labels from videos by observing person-object interactions [170]. For instance, the knife has the "cut" function because many persons used it for cutting. Additionally, they provided a video dataset for indoor scenes with function labels. In [171], Pieropan *et al.* proposed using object-to-object spatio-temporal relationships to create a so called "object context" along with functional descriptor to predict the human activities. As an example, only the presence of a mug does not confirm if a drinking action could take place, but the presence of a water bottle beside it increases the likelihood of a drinking action. They used the kitchen images of the dataset proposed in [127] and trained a probabilistic model, Conditional Random Field (CRF), to classify functional classes into four types: tools, ingredients, support and containers. Yao *et al.* represented function understanding problem as weakly supervised to discover all possible functionalities for each object [55]. They used unsupervised clustering in an iterative manner to categorize human-object interactions, then used these updates as input for detection and pose estimation, and finally discovered the functionalities. The authors tested their model through musical instrument dataset [172] which con-

tains images for human-object interactions. Given human-object interaction observations, Stark *et al.* [173] proposed a learning mechanism to categorize grasping affordances as object functionalities as shown in Figure 20. Through observing the interaction with some object, the affordance cues are defined. In other words, the affordance has been defined as relation between robot hand and an object. The implicit shape model (ISM) [174] was the main algorithm used by [173] to categorize the objects. To this end, these cues have been used to predict the grasping points of a 2D image. Pieropan *et al.* [127], applied functional descriptors/cues, which have been learned from hand-object interaction, to understand human activities from RGB-D videos. They represented the objects by their interaction with human hands as well as they encoded these objects as strings through which the string kernel [175] measured their similarity. Likewise, they used spatial location and temporal trajectory to estimate object position relative to the hand. Hence, the estimated object position produced a so called functional descriptor. Thereafter, they fused the similarity measures with functional descriptors to recognize human activities. Mar *et al.* [125] proposed a method to learn tool affordances based on its function and the way of interaction (grasping). To find functional descriptors, they learned geometrical features through Self-Organized Maps (SOM) and K-means. Hu *et al.* introduced what called Interaction CONtext (ICON) to describe the functionality of 3D object through geometric features [176] focused on other aspects of functionality usage in vision research. The main idea was to define what is called the contextual descriptor to describe the shape functionality in the presence of other objects using Interaction Bisector Surface (IBS) [114]. In other words, they used object-to-object interaction to build a geometrical descriptor for functionality analysis and hence they recognized the correspondences between similar parts on various shapes of 3D images. They used the Trimble 3D Warehouse<sup>1</sup> to test their experiment. Savva *et al.* [168] built action maps for potential actions through scanning geometry of the captured 3D scenes, reconstructing depth meshes, tracking human interactions to define the functionality descriptors and therefore predicting the affordances of unseen objects.

### 6.4 Functional descriptor shape and correspondence

Jain *et al.* [177] proposed a probabilistic model based on geometric features, which is related to the functionalities such as those based on material and shape, along with probabilistic dependencies between the effects, tool's actions and tool's functional features. A Bayesian Network (BN) is used in their scheme because of its ability to handle the probabilistic dependencies between nodes. Laga *et al.* proposed a model that extracted the pair-wise semantics of shapes through combining structural and geometric features [178]. Additionally, they recognized the functionality of the shapes (e.g. graspable and container) using supervised learning methods. Kim *et al.* used affordances as priori information to predict the correspondence in human poses through which they predicted the functional descriptors [179]. Apart from above-mentioned studies which used abstract functionality

1. <https://3dwarehouse.sketchup.com/>

	Input		Features		Evaluation			Training		Abstraction		
	2D	3D	Feature learning	Handcrafted	Real Robot	Simulation	Benchmark	Supervised	Weakly supervised	Unsupervised	Mathematical	Neutral
Zhu <i>et al.</i> [22]		✓		✓			✓	✓			✓	
Shiraki <i>et al.</i> [152]		✓		✓			✓				✓	
Turek <i>et al.</i> [164]	✓			✓			✓			✓		
Yao <i>et al.</i> [55]	✓			✓			✓		✓		✓	
Pechuk <i>et al.</i> [153]		✓		✓			✓					✓
Hinkle and Olson [155]		✓		✓			✓		✓		✓	
Jain <i>et al.</i> [177]	✓			✓		✓				✓		
Awaad <i>et al.</i> [156]	✓			✓		✓				✓		
Madry <i>et al.</i> [157]		✓		✓			✓				✓	
Xie <i>et al.</i> [159]	✓			✓			✓		✓		✓	
Oh <i>et al.</i> [163]	✓			✓			✓		✓		✓	
Zen [165] <i>et al.</i>	✓			✓			✓		✓		✓	
Gupta <i>et al.</i> [162]	✓			✓			✓		✓		✓	
Delaitre <i>et al.</i> [169]	✓			✓			✓		✓		✓	
Fouhey <i>et al.</i> [170]	✓			✓			✓		✓		✓	
Rhinehart and Kitani [167]	✓			✓			✓		✓		✓	
Zhao and Zhu [151]	✓			✓			✓		✓		✓	
Kim <i>et al.</i> [179]		✓		✓			✓		✓		✓	
Hu <i>et al.</i> [176]		✓		✓			✓		✓		✓	
Laga <i>et al.</i> [178]		✓		✓			✓		✓		✓	
Lun <i>et al.</i> [180]		✓		✓		✓			✓		✓	
Savva <i>et al.</i> [183]		✓		✓			✓		✓		✓	
Lun <i>et al.</i> [180]		✓		✓			✓		✓		✓	
Saponaro <i>et al.</i> [161]	✓			✓		✓			✓		✓	
Stark <i>et al.</i> [173]	✓			✓			✓		✓		✓	

TABLE 9: Comparison between function-scene understanding methods

such as "to pour" and "to move". In recent study, Lun *et al.* [180] designed a unified model to detect a human pose according to human-object affordances (leaning, holding, sitting and treating) along with object parts. The functionality descriptor has been employed to recover mechanical assemblies or parts from raw scans [181]. They used segmentation and joint optimization to learn their scheme. Hu *et al.* in recent study [182] proposed a method to analyze inter-object relations and intra-object relation aiming to categorize the objects based on their functionalities. They used objects' parts contexts, semantics and functionalities to recognize their shapes.

Various methods with different ways of representation have been introduced to address the functionality issues. Table 9 summarizes and compares the most important properties to make it more understandable.

## 7 DATASETS

In this section, we investigate the available datasets provided with affordance annotations. As the following (Table 10) shows, the distribution of them range from RGB, RGB-D for images and videos. For visual cues, many datasets have been proposed such as UMD and IIT-AFF [49], [61] to facilitate detecting affordance objects from the scene i.e. detecting objects that bear affordances or functionalities from the input image. In other words, these datasets enable researchers to treat affordances as traditional image detection tasks like pedestrian detection or face detection. Since physical and material attributes are important for defining the functionality and affordances, it has been provided in these methods [58], [94], [139], [142]. Regarding human activities recognition, these methods [54], [107], [60]

advanced annotated datasets that contain human subjects. In the essence of objects parts, this dataset [94] has part's annotations.

## 8 OPEN PROBLEMS AND RESEARCH QUESTIONS

**Lack of Consensus on Affordance Definition:** Nearly 5 decades after the introduction of the affordance concept by Gibson, no formal definition of affordances has been agreed upon by AI researchers and ecological psychologists. Gibson's own description of the concept in his seminal work "The Ecological Approach to Visual Perception" (1979) shows the complex nature of the concept, e.g., he said "Affordance is equally a fact of the environment and a fact of behavior. It is both physical and psychological, yet neither. An affordance points both ways, to the environment and to the observer" [9]. Previous efforts have tried to describe affordance as the property of the environment, the mutual phenomenon between an agent and its environment or an observer's perception of the relationships between an agent and its surroundings [185], [13]. However, there does not exist a unified definition for affordances so far.

**Function vs Affordance:** The terms, function and affordance, are sometimes used with the same meaning in the literature even though they are completely different. The relation between these two terms is complicated where one object may have different affordances and functions such as the case of a hammer which has affordances (grasping, striking, dragging) and functions (drive nails, fit parts, forge metal, and break apart objects). Differently, some objects have affordances and functions with the same meaning such as knife which has affordance (cut) and function (cut). However, the agreed concept (also emphasized in this survey) is that affordance always relate to the object itself rather than function which relates to only another object. In other words, affordance relates to the possible actions whereas the function relates to the effect. Much effort has to be devoted in this issue to distinguish between affordance and functionality in the right manner.

**Affordance and Attributes:** Recently, much research have been introduced to describe objects with their attributes [186], [187], [188], [189]. Along with the object label, semantically meaningful attributes are needed to understand the object characteristics. For instance, a cup of tea may be described with some attributes like glass, white and has handle. Therefore, describing the cup by its attributes will help in deciding the best way of grasping. Furthermore, the affordance learning problem needs attributes to address some difficulties such as getting some fruits from the fridge requires prior knowledge about the height of the fridge and the agent as well. Despite its significance, the use of attributes has not been addressed in the context of visual affordances before.

**Multi-class Labeling:** Considering affordances and functions together gives rise to an advanced set of labels for every object. Therefore, it is a multi-label problem in its nature. For example, the Figure of cup 19 shows nine labels. If more detailed analysis is required, these labels need to be ordered according to the object, scene and the situation. Remarkably, these labels should be identified in terms of the parts rather than the whole object.

Reference	Year	Properties	Affordances	Format	Subjects Number	Categories	Image/Video	Subjects
UMD[49]	2015	visual	Table 1	RGB-D Image		17	30,000	Indoor
[58]	2011	visual, physical and material	Table 1	RGB-D Image		7	375	Indoor
CAD120 [60]	2016	visual, human-interactions	Table 1	RGB-D Video	35		3090 / 215	Indoor
IIT-AFF [61]	2017	visual	Table 1	RGB-D Image		9	8,835	Tools
ADE- Affordance [94]	2016	visual, physical, social, object-action pairs, exceptions, explanations	Table 1	RGB Images		7	10,000	Indoors
[184]	2016	visual	Table 1	RGB Image		8	10,360	Indoor tools
Extended NYUv2 [63]	2016	visual	Table 1	RGB Image		5	1449	Indoor
CONTACT VMGdB [107]	2011	visual, human interactions	Table 1	RGB-D Video	20	5	5200	Indoor
HHOI[54]	2016	visual, human interactions	Table1	RGB-D video	14	5	-	Indoor
Binge Watching [110]	2017	visual, human poses	30 poses	RGB-D		30	11449	indoor
CERTH-SOR3D [65]	2017	visual, human interactions	Table 1	RGB-D		14	20,800	indoor, tools
COQE [139]	2017	visual, physical	-	RGB		10	5000	Containers
[142]	2016	visual, physical	-	RGB-D video		4	1326	Containers
Tool & Tool-Use (TTU) [22]	2015	visual, physical, human demonstrations	Table 1	RGB-D Image		10	452	Tools

TABLE 10: The datasets that have ground-truth annotation for affordances and functionalities.

**Deep Learning for Affordance Learning:** Nowadays, deep learning is dominating the field of vision and it achieved remarkable improvements in this context. However, it has not received much attention in addressing the challenges particular to affordance and function understanding. By looking at comparison tables, the number of feature learning methods proposed in the literature is too few. Another related factor to deep learning is the size of datasets. In other words, modeling the affordance algorithms with deep learning needs large annotated datasets which are currently unavailable.

**Complex Affordances:** Affordances of an object can be classified to basic and higher order affordances [49], [51] e.g. the "pouring" hot water into a cup is basic affordance while making tea from this water is higher-order because it depends on the basic one. These higher-order affordances have not been explored in the existing works.

**Outdoor Affordances:** Mostly all of the previous studies focused on indoor scenes. Only one method has been introduced to address the outdoor affordances [86]. Despite of

less focus on outdoor affordances, this direction of research is highly valuable due to the relationship between affordances and silent actions which is important in self-driving cars, autonomous driving, traffic monitoring applications.

**Affordances for Developmental Robotics:** Since developmental robotics and affordances are tightly related to each other, visually understanding the environment can enhance the learning process. Developmental robotics seek to enhance robot understanding through the environmental interactions while affordances are emergent environmental variables. Thus, learning affordances visually would shorten the time required for interactive robots to build their knowledge accurately. Yet, the conducted studies in this paradigm are still not sufficient and need more investigations to get efficient baselines. In case of extrinsic motivation, visual affordances can be used as reward evaluators and indicators to measure the progress of that learning scheme. For intrinsic motivations, the visual affordances can be utilized to point out the most important features in the environment which will increase the robots curiosity to learn more.

**Visual Questions Answering (VQA) & Learning-by-Asking (LBA):** VQA deals with answering intelligent questions about a visual element. Because affordances are constant environment variables that provide highly valuable information about scene content, they can assist in developing a deep insight. Similar to affordance-based recognition, fusing affordances and VQA will improve the accuracy of these answers. Unlike VQA, LBA seeks to understand the environment by asking questions and requesting supervision. Assuming that affordances are a precursor for an object interaction, merging affordances and LBA will reduce the agent time to build its knowledge.

We believe that the application of visual scene understanding algorithms including those for semantic/instance segmentation, physics based reasoning and 3D volumetric analysis to affordance will help in resolving several underlying challenges.

## 9 CONCLUSION

In this survey, visual affordance and functional scene understanding has been reviewed. We introduce a hierarchical taxonomy and cover the progress according to each sub task e.g., classification, segmentation and detection. We begin with formal definition of each sub-task and provide significance and applications to motivate the readers. The paper succinctly compares best approaches in every sub-category in a tabular form to help researchers identify research gaps and open questions. The paper also covers datasets proposed in the field and provides detailed comparisons between them. We discussed the open problems and challenges in this field. Finally, we hope this survey will be helpful for researchers particularly as it is the first survey to review visual affordances and function understanding.

## REFERENCES

- [1] J. J. Gibson, "The senses considered as perceptual systems." 1966.
- [2] T. E. Horton, A. Chakraborty, and R. S. Amant, "Affordances for robots: a brief survey." *Avant*, vol. 3, no. 2, 2012.
- [3] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesis; a survey," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 289–309, April 2014.
- [4] N. Yamanobe, W. Wan, I. G. Ramirez-Alpizar, D. Petit, T. Tsuji, S. Akizuki, M. Hashimoto, K. Nagata, and K. Harada, "A brief review of affordance in robotic manipulation research," *Advanced Robotics*, vol. 31, no. 19-20, pp. 1086–1101, 2017.
- [5] H. Min, C. Yi, R. Luo, J. Zhu, and S. Bi, "Affordance research in developmental robotics: A survey," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 8, no. 4, pp. 237–255, 2016.
- [6] L. Jamone, E. Ugur, A. Cangelosi, L. Fadiga, A. Bernardino, J. Piater, and J. Santos-Victor, "Affordances in psychology, neuroscience and robotics: a survey," *IEEE Transactions on Cognitive and Developmental Systems*, 2016.
- [7] L. Montesano, M. Lopes, A. Bernardino, and J. Santos-Victor, "Learning object affordances: from sensory-motor coordination to imitation."
- [8] J. J. Gibson, "The theory of affordances," *The people, place, and space reader*, pp. 56–60, 1979.
- [9] G. James, "The ecological approach to visual perception," *Dallas: Houghton Mifflin*, 1979.
- [10] M. T. Turvey, "Affordances and prospective control: An outline of the ontology," *Ecological psychology*, vol. 4, no. 3, pp. 173–187, 1992.
- [11] T. A. Stoffregen, "Affordances and events," *Ecological psychology*, vol. 12, no. 1, pp. 1–28, 2000.
- [12] —, "Affordances as properties of the animal-environment system," *Ecological psychology*, vol. 15, no. 2, pp. 115–134, 2003.
- [13] E. Şahin, M. Çakmak, M. R. Doğar, E. Uğur, and G. Üçoluk, "To afford or not to afford: A new formalization of affordances toward affordance-based robot control," *Adaptive Behavior*, vol. 15, no. 4, pp. 447–472, 2007.
- [14] C.-Y. Chuang, J. Li, A. Torralba, and S. Fidler, "Learning to act properly: Predicting and explaining affordances from images," *arXiv preprint arXiv:1712.07576*, 2017.
- [15] H. S. Koppula and A. Saxena, "Anticipating human activities using object affordances for reactive robotic response," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 1, pp. 14–29, 2016.
- [16] H. S. Koppula, R. Gupta, and A. Saxena, "Learning human activities and object affordances from rgb-d videos," *The International Journal of Robotics Research*, vol. 32, no. 8, pp. 951–970, 2013.
- [17] A. Jain, A. R. Zamir, S. Savarese, and A. Saxena, "Structural-rnn: Deep learning on spatio-temporal graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5308–5317.
- [18] S. Qi, S. Huang, P. Wei, and S.-C. Zhu, "Predicting human activities using stochastic grammar," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [19] J. Earley, "An efficient context-free parsing algorithm," *Communications of the ACM*, vol. 13, no. 2, pp. 94–102, 1970.
- [20] T.-H. Vu, C. Olsson, I. Laptev, A. Oliva, and J. Sivic, "Predicting actions from static scenes," in *European Conference on Computer Vision*. Springer, 2014, pp. 421–436.
- [21] H. Grabner, J. Gall, and L. Van Gool, "What makes a chair a chair?" in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 1529–1536.
- [22] Y. Zhu, Y. Zhao, and S. Chun Zhu, "Understanding tools: Task-oriented object modeling, learning and recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2855–2864.
- [23] H. Kjellström, J. Romero, and D. Kragić, "Visual object-action recognition: Inferring object affordances from human demonstration," *Computer Vision and Image Understanding*, vol. 115, no. 1, pp. 81–90, 2011.
- [24] A. Gupta, S. Satkin, A. A. Efros, and M. Hebert, "From 3d scene geometry to human workspace," in *CVPR 2011*, June 2011, pp. 1961–1968.
- [25] D. Parikh and K. Grauman, "Relative attributes," in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 503–510.
- [26] B. Moldovan and L. D. Raedt, "Occluded object search by relational affordances," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, May 2014, pp. 169–174.
- [27] L. L. S. Wong, L. P. Kaelbling, and T. Lozano-Pérez, "Manipulation-based active search for occluded objects," in *2013 IEEE International Conference on Robotics and Automation*, May 2013, pp. 2814–2819.
- [28] M. Gupta, J. Mājller, and G. S. Sukhatme, "Using manipulation primitives for object sorting in cluttered environments," *IEEE Transactions on Automation Science and Engineering*, vol. 12, no. 2, pp. 608–614, April 2015.
- [29] H. O. Song, M. Fritz, C. Gu, and T. Darrell, "Visual grasp affordances from appearance-based cues," in *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, Nov 2011, pp. 998–1005.
- [30] A. Saxena, L. L. Wong, and A. Y. Ng, "Learning grasp strategies with partial shape information." in *AAAI*, vol. 3, no. 2, 2008, pp. 1491–1494.
- [31] E. Gibaja and S. Ventura, "A tutorial on multilabel learning," *ACM Computing Surveys (CSUR)*, vol. 47, no. 3, p. 52, 2015.
- [32] —, "Multi-label learning: a review of the state of the art and ongoing research," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 4, no. 6, pp. 411–444, 2014.
- [33] Z.-H. Zhou, M.-L. Zhang, S.-J. Huang, and Y.-F. Li, "Multi-instance multi-label learning," *Artificial Intelligence*, vol. 176, no. 1, pp. 2291–2320, 2012.
- [34] Z.-H. Zhou and M.-L. Zhang, "Multi-instance multi-label learning with application to scene classification," in *Advances in neural information processing systems*, 2007, pp. 1609–1616.
- [35] J. Wu, X. Bai, M. Loog, F. Roli, and Z.-H. Zhou, "Multi-instance learning in pattern recognition and vision," 2017.
- [36] A. Gupta and L. S. Davis, "Objects in action: An approach for combining action understanding and object perception," in *2007*

- IEEE Conference on Computer Vision and Pattern Recognition*, June 2007, pp. 1–8.
- [37] M. Hassan and A. Dharmaratne, *Attribute Based Affordance Detection from Human-Object Interaction Images*. Cham: Springer International Publishing, 2016, pp. 220–232.
- [38] M. Liang and X. Hu, “Recurrent convolutional neural network for object recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3367–3375.
- [39] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [40] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.
- [41] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single shot multibox detector,” in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [42] L. Bourdev and J. Malik, “Poselets: Body part detectors trained using 3d human pose annotations,” in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 1365–1372.
- [43] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part-based models,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [44] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders, “Selective search for object recognition,” *International journal of computer vision*, vol. 104, no. 2, pp. 154–171, 2013.
- [45] M. Dixit, S. Chen, D. Gao, N. Rasiwasia, and N. Vasconcelos, “Scene classification with semantic fisher vectors,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [46] L. Zhang, X. Zhen, and L. Shao, “Learning object-to-class kernels for scene classification,” *IEEE Transactions on Image Processing*, vol. 23, no. 8, pp. 3241–3253, Aug 2014.
- [47] S. Gupta, P. Arbeláez, R. Girshick, and J. Malik, “Indoor scene understanding with rgb-d images: Bottom-up segmentation, object detection and semantic segmentation,” *International Journal of Computer Vision*, vol. 112, no. 2, pp. 133–149, 2015.
- [48] W. Choi, Y.-W. Chao, C. Pantofaru, and S. Savarese, “Indoor scene understanding with geometric and semantic contexts,” *International Journal of Computer Vision*, vol. 112, no. 2, pp. 204–220, 2015.
- [49] A. Myers, C. L. Teo, C. Fermajller, and Y. Aloimonos, “Affordance detection of tool parts from geometric features,” in *2015 IEEE International Conference on Robotics and Automation (ICRA)*, May 2015, pp. 1374–1381.
- [50] L. Stark and K. Bowyer, “Function-based generic recognition for multiple object categories,” *CVGIP: Image Understanding*, vol. 59, no. 1, pp. 1–21, 1994.
- [51] C. Ye, Y. Yang, R. Mao, C. Fermajller, and Y. Aloimonos, “What can i do around here? deep functional scene understanding for cognitive robots,” in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, May 2017, pp. 4604–4611.
- [52] E. Ugur, M. R. Dogar, M. Cakmak, and E. Sahin, “The learning and use of traversability affordance using range images on a mobile robot,” in *Robotics and Automation, 2007 IEEE International Conference on*. IEEE, 2007, pp. 1721–1726.
- [53] Y. W. Chao, Z. Wang, R. Mihalcea, and J. Deng, “Mining semantic affordances of visual object categories,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 4259–4267.
- [54] T. Shu, M. S. Ryoo, and S.-C. Zhu, “Learning social affordance for human-robot interaction,” in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, ser. IJCAI’16. AAAI Press, 2016, pp. 3454–3461. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3061053.3061104>
- [55] B. Yao, J. Ma, and L. Fei-Fei, “Discovering object functionality,” in *2013 IEEE International Conference on Computer Vision*, Dec 2013, pp. 2512–2519.
- [56] S. Aarathi and S. Chitrakala, “Scene understanding; a survey,” in *2017 International Conference on Computer, Communication and Signal Processing (ICCCSP)*, Jan 2017, pp. 1–4.
- [57] A. Aldoma, F. Tombari, and M. Vincze, “Supervised learning of hidden and non-hidden 0-order affordances and detection in real scenes,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 1732–1739.
- [58] T. Hermans, J. M. Rehg, and A. Bobick, “Affordance prediction via learned object attributes,” in *IEEE International Conference on Robotics and Automation (ICRA): Workshop on Semantic Perception, Mapping, and Exploration*, 2011, pp. 181–184.
- [59] T. Luddecke and F. Worgotter, “Learning to segment affordances,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 769–776.
- [60] J. Sawatzky, A. Srikantha, and J. Gall, “Weakly supervised affordance detection,” July 2017.
- [61] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis, “Object-based affordances detection with convolutional neural networks and dense conditional random fields,” in *International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- [62] M. Schoeler and F. Worgotter, “Bootstrapping the semantics of tools: Affordance analysis of real world objects on a per-part basis,” *IEEE Transactions on Cognitive and Developmental Systems*, vol. 8, no. 2, pp. 84–98, 2016.
- [63] A. Roy and S. Todorovic, “A multi-scale cnn for affordance segmentation in rgb images,” in *European Conference on Computer Vision*. Springer, 2016, pp. 186–201.
- [64] D. I. Kim and G. S. Sukhatme, “Semantic labeling of 3d point clouds with object affordance for robot manipulation,” in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 5578–5584.
- [65] S. Themos, G. T. Papadopoulos, P. Daras, and G. Potamianos, “Deep affordance-grounded sensorimotor object recognition,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [66] L. Bo, X. Ren, and D. Fox, “Unsupervised feature learning for rgb-d based object recognition,” in *Experimental Robotics*. Springer, 2013, pp. 387–402.
- [67] P. Dollar and C. L. Zitnick, “Structured forests for fast edge detection,” in *Computer Vision (ICCV), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1841–1848.
- [68] B. Moldovan, P. Moreno, M. van Otterlo, J. Santos-Victor, and L. De Raedt, “Learning relational affordance models for robots in multi-object manipulation tasks,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 4373–4378.
- [69] L. De Raedt, P. Frasconi, K. Kersting, and S. H. Muggleton, *Probabilistic inductive logic programming*. Springer, 2008, vol. 4911.
- [70] L. De Raedt, A. Kimmig, and H. Toivonen, “Problog: A probabilistic prolog and its application in link discovery,” 2007.
- [71] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis, “Detecting object affordances with convolutional neural networks,” in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct 2016, pp. 2765–2770.
- [72] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, “Multimodal deep learning,” in *Proceedings of the 28th international conference on machine learning (ICML-11)*, 2011, pp. 689–696.
- [73] S. Gupta, R. Girshick, P. Arbelaez, and J. Malik, “Learning rich features from rgb-d images for object detection and segmentation,” in *European Conference on Computer Vision*. Springer, 2014, pp. 345–360.
- [74] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [75] J. Sawatzky and J. Gall, “Adaptive binarization for weakly supervised affordance segmentation,” in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [76] J. Dai, Y. Li, K. He, and J. Sun, “R-fcn: Object detection via region-based fully convolutional networks,” in *Advances in neural information processing systems*, 2016, pp. 379–387.
- [77] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 91–99. [Online]. Available: <http://papers.nips.cc/paper/5638-faster-r-cnn-towards-real-time-object-detection-with-region-proposal-net.pdf>
- [78] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.

- [79] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [80] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," in *Proceedings of the International Conference on Learning Representations*, 2015.
- [81] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr, "Conditional random fields as recurrent neural networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1529–1537.
- [82] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama et al., "Speed/accuracy trade-offs for modern convolutional object detectors," *arXiv preprint arXiv:1611.10012*, 2016.
- [83] T.-T. Do, A. Nguyen, I. Reid, D. G. Caldwell, and N. G. Tsagarakis, "Affordancenet: An end-to-end deep learning approach for object affordance detection," *arXiv preprint arXiv:1709.07326*, 2017.
- [84] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask r-cnn," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [85] Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei, "Fully convolutional instance-aware semantic segmentation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [86] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec 2015, pp. 2722–2730.
- [87] A. Geiger, P. Lenz, C. Stillner, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [88] B. Wymann, E. Espi e, C. Guionneau, C. Dimitrakakis, R. Coulom, and A. Sumner, "Torcs, the open racing car simulator," *Software available at <http://www.torcs.org>*, 2014.
- [89] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2650–2658.
- [90] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgb-d images," in *European Conference on Computer Vision*. Springer, 2012, pp. 746–760.
- [91] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Doll ar, "Learning to refine object segments," in *European Conference on Computer Vision*. Springer, 2016, pp. 75–91.
- [92] M. Savva, A. X. Chang, P. Hanrahan, M. Fisher, and M. Nie sner, "Scenegrok: Inferring action maps in 3d environments," *ACM transactions on graphics (TOG)*, vol. 33, no. 6, p. 212, 2014.
- [93] N. Rhinehart and K. M. Kitani, "Learning action maps of large environments via first-person vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 580–588.
- [94] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Semantic understanding of scenes through the ade20k dataset," *arXiv preprint arXiv:1608.05442*, 2016.
- [95] P. Fitzpatrick, G. Metta, L. Natale, S. Rao, and G. Sandini, "Learning about objects through action-initial steps towards artificial cognition," in *Robotics and Automation, 2003. Proceedings. ICRA'03. IEEE International Conference on*, vol. 3. IEEE, 2003, pp. 3140–3145.
- [96] H. Kozima, C. Nakagawa, and H. Yano, "Emergence of imitation mediated by objects," 2002.
- [97] L. Montesano, M. Lopes, A. Bernardino, and J. Santos-Victor, "Learning object affordances: From sensory–motor coordination to imitation," *IEEE Transactions on Robotics*, vol. 24, no. 1, pp. 15–26, Feb 2008.
- [98] —, "Modeling affordances using bayesian networks," in *Intelligent Robots and Systems, 2007. IROS 2007. IEEE/RSJ International Conference on*. IEEE, 2007, pp. 4102–4107.
- [99] M. Lopes, F. S. Melo, and L. Montesano, "Affordance-based imitation learning in robots," in *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Oct 2007, pp. 1015–1021.
- [100] E. Ugur, E. Oztop, and E. Sahin, "Going beyond the perception of affordances: Learning how to actualize them through behavioral parameters," in *2011 IEEE International Conference on Robotics and Automation*, May 2011, pp. 4768–4773.
- [101] K. M. Varadarajan and M. Vincze, "Afnet: The affordance network," in *Asian Conference on Computer Vision*. Springer, 2012, pp. 512–523.
- [102] —, "Afrob: The affordance network ontology for robots," in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*. IEEE, 2012, pp. 1343–1350.
- [103] C. Havasi, R. Speer, and J. Alonso, "Conceptnet 3: a flexible, multilingual semantic network for common sense knowledge," in *Recent advances in natural language processing*. John Benjamins Philadelphia, PA, 2007, pp. 27–29.
- [104] M. Tenorth, L. Kunze, D. Jain, and M. Beetz, "Knowrob-map - knowledge-linked semantic object maps," in *2010 10th IEEE-RAS International Conference on Humanoid Robots*, Dec 2010, pp. 430–435.
- [105] K. M. Varadarajan and M. Vincze, "Object part segmentation and classification in range images for grasping," in *Advanced Robotics (ICAR), 2011 15th International Conference on*. IEEE, 2011, pp. 21–27.
- [106] J. J. Mor e, "The levenberg-marquardt algorithm: implementation and theory," in *Numerical analysis*. Springer, 1978, pp. 105–116.
- [107] C. Castellini, T. Tommasi, N. Noceti, F. Odone, and B. Caputo, "Using object affordances to improve object recognition," *IEEE Transactions on Autonomous Mental Development*, vol. 3, no. 3, pp. 207–215, Sept 2011.
- [108] B. Moldovan, P. Moreno, D. Nitti, J. Santos-Victor, and L. De Raedt, "Relational affordances for multiple-object manipulation," *Autonomous Robots*, vol. 42, no. 1, pp. 19–44, 2018.
- [109] M. Lopes and J. Santos-Victor, "Visual learning by imitation with motor representations," *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 35, no. 3, pp. 438–449, June 2005.
- [110] X. Wang, R. Girdhar, and A. Gupta, "Binge watching: Scaling affordance learning from sitcoms," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [111] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [112] Y. Sun, S. Ren, and Y. Lin, "Object–object interaction affordance learning," *Robotics and Autonomous Systems*, vol. 62, no. 4, pp. 487–496, 2014.
- [113] M. Peternell, "Geometric properties of bisector surfaces," *Graphical Models*, vol. 62, no. 3, pp. 202–236, 2000.
- [114] X. Zhao, H. Wang, and T. Komura, "Indexing 3d scenes using the interaction bisector surface," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 3, p. 22, 2014.
- [115] K. M. Varadarajan and M. Vincze, "Parallel deep learning with suggestive activation for object category recognition," in *International Conference on Computer Vision Systems*. Springer, 2013, pp. 354–363.
- [116] J. Sun, J. L. Moore, A. Bobick, and J. M. Rehg, "Learning visual object categories for robot affordance prediction," *The International Journal of Robotics Research*, vol. 29, no. 2-3, pp. 174–197, 2010.
- [117] E. Ugur, S. Szedmak, and J. Piater, "Bootstrapping paired-object affordance learning with learned single-affordance features," in *Development and Learning and Epigenetic Robotics (ICDL-Epirob), 2014 joint IEEE International Conferences on*. IEEE, 2014, pp. 476–481.
- [118] M. Schoeler, J. Papon, and F. Worgotter, "Constrained planar cuts-object partitioning for point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5207–5215.
- [119] F. Tombari, S. Salti, and L. Di Stefano, "Unique signatures of histograms for local surface description," in *European conference on computer vision*. Springer, 2010, pp. 356–369.
- [120] W. Mustafa, N. Pugeault, and N. Kr uger, "Multi-view object recognition using view-point invariant shape relations and appearance information," in *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. IEEE, 2013, pp. 4230–4237.
- [121] P. Abelha, F. Guerin, and M. Schoeler, "A model-based approach to finding substitute tools in 3d vision data," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, May 2016, pp. 2471–2478.
- [122] P. Abelha Ferreira and F. Guerin, "Learning how a tool affords by simulating 3d models from the web," in *Proceedings of IEEE International Conference on Intelligent Robots and Systems (IROS 2017)*. IEEE Press, 2017.

- [123] P. Abelho and F. Guerin, "Transfer of tool affordance and manipulation cues with 3d vision data," *arXiv preprint arXiv:1710.04970*, 2017.
- [124] A. Jaklic, A. Leonardis, and F. Solina, *Segmentation and recovery of superquadrics*. Springer Science & Business Media, 2013, vol. 20.
- [125] T. Mar, V. Tikhonoff, G. Metta, and L. Natale, "Multi-model approach based on 3d functional features for tool affordance learning in robotics," in *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, Nov 2015, pp. 482–489.
- [126] —, "Self-supervised learning of tool affordances from 3d tool representation through parallel self mapping," in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 894–901.
- [127] A. Pieropan, C. H. Ek, and H. Kjellström, "Functional object descriptors for human activity modeling," in *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1282–1289.
- [128] —, "Functional descriptors for object affordances," in *IEEE International Conference on Intelligent Robots and Systems Workshop*. IEEE, 2015.
- [129] V. Dutta and T. Zielinska, "Action prediction based on physically grounded object affordances in human-object interactions," in *Robot Motion and Control (RoMoCo), 2017 11th International Workshop on*. IEEE, 2017, pp. 47–52.
- [130] T. Shu, X. Gao, M. S. Ryoo, and S. C. Zhu, "Learning social affordance grammar from videos: Transferring human interactions to human-robot interactions," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, May 2017, pp. 1669–1676.
- [131] Y. Zhu, A. Fathi, and L. Fei-Fei, "Reasoning about object affordances in a knowledge base representation," in *European conference on computer vision*. Springer, 2014, pp. 408–424.
- [132] M. Richardson and P. Domingos, "Markov logic networks," *Machine learning*, vol. 62, no. 1-2, pp. 107–136, 2006.
- [133] M. Rohrbach, M. Stark, and B. Schiele, "Evaluating knowledge transfer and zero-shot learning in a large-scale setting," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 1641–1648.
- [134] Y. W. Chao, Z. Wang, R. Mihalcea, and J. Deng, "Mining semantic affordances of visual object categories," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 4259–4267.
- [135] G. A. Miller, "Wordnet: a lexical database for english," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [136] T. Zhou, H. Shan, A. Banerjee, and G. Sapiro, "Kernelized probabilistic matrix factorization: Exploiting graphs and side information," in *Proceedings of the 2012 SIAM International Conference on Data Mining*. SIAM, 2012, pp. 403–414.
- [137] L.-F. Yu, N. Duncan, and S.-K. Yeung, "Fill and transfer: A simple physics-based approach for containability reasoning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 711–719.
- [138] H. Wang, W. Liang, and L.-F. Yu, "Transferring objects: Joint inference of container and human pose," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2933–2941.
- [139] R. Mottaghi, C. Schenck, D. Fox, and A. Farhadi, "See the glass half full: Reasoning about liquid containers, their volume and content," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [140] C. J. Phillips, M. Lecce, and K. Daniilidis, "Seeing glassware: from edge detection to pose estimation and shape recovery." in *In Robotics: Science and Systems*, vol. 3, 2016.
- [141] W. Liang, Y. Zhao, Y. Zhu, and S.-C. Zhu, "Evaluating human cognition of containing relations with physical simulation." in *Proceedings of the 37th Annual Meeting of the Cognitive Science Society*, 2015.
- [142] —, "What is where: Inferring containment relations from videos."
- [143] Y. Zhu, C. Jiang, Y. Zhao, D. Terzopoulos, and S.-C. Zhu, "Inferring forces and learning human utilities from videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3823–3833.
- [144] Y. Zhu, C. Zhang, C. Ré, and L. Fei-Fei, "Building a large-scale multimodal knowledge base system for answering visual queries," *arXiv preprint arXiv:1507.05670*, 2015.
- [145] P. H. Winston, T. O. Binford, B. Katz, and M. Lowry, *Learning physical descriptions from functional definitions, examples, and prece-*
- dents*. Department of Computer Science, Stanford University, 1983.
- [146] L. Stark and K. Bowyer, "Achieving generalized object recognition through reasoning about association of function to structure," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 13, no. 10, pp. 1097–1104, Oct 1991.
- [147] E. Rivlin, S. J. Dickinson, and A. Rosenfeld, "Recognition by functional parts," *Computer Vision and Image Understanding*, vol. 62, no. 2, pp. 164–176, 1995.
- [148] C. Desai and D. Ramanan, "Predicting functional regions on objects," in *2013 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, June 2013, pp. 968–975.
- [149] —, "Detecting actions, poses, and objects with relational phraselets," in *Computer Vision – ECCV 2012*, A. Fitzgibbon, S. Lazebnik, P. Perona, Y. Sato, and C. Schmid, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 158–172.
- [150] L. Bourdev, S. Maji, and J. Malik, "Detection, attribute classification and action recognition of people using poselets (in submission)," *IEEE PAMI*, 2013.
- [151] Y. Zhao and S.-C. Zhu, "Scene parsing by integrating function, geometry and appearance models," in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*. IEEE, 2013, pp. 3119–3126.
- [152] Y. Shiraki, K. Nagata, N. Yamanobe, A. Nakamura, K. Harada, D. Sato, and D. N. Nenchev, "Modeling of everyday objects for semantic grasp," in *The 23rd IEEE International Symposium on Robot and Human Interactive Communication*, Aug 2014, pp. 750–755.
- [153] M. Pechuk, O. Soldea, and E. Rivlin, "Learning function-based object classification from 3d imagery," *Computer Vision and Image Understanding*, vol. 110, no. 2, pp. 173–191, 2008.
- [154] Y. Kitano, T. Takiguchi, and Y. Ariki, "Estimation of object functions using convolutional neural network," in *The Korea-Japan Joint Workshop on Frontiers of Computer Vision*, 2016.
- [155] L. Hinkle and E. Olson, "Predicting object functionality using physical simulations," in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Nov 2013, pp. 2784–2790.
- [156] I. Awaad, G. K. Kraetzschmar, and J. Hertzberg, "The role of functional affordances in socializing robots," *International Journal of Social Robotics*, vol. 7, no. 4, pp. 421–438, 2015.
- [157] M. Madry, D. Song, and D. Kragic, "From object categories to grasp transfer using probabilistic reasoning," in *2012 IEEE International Conference on Robotics and Automation*, May 2012, pp. 1716–1723.
- [158] D. Song, K. Huebner, V. Kyrki, and D. Kragic, "Learning task constraints for robot grasping using graphical models," in *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*. IEEE, 2010, pp. 1579–1585.
- [159] D. Xie, S. Todorovic, and S.-C. Zhu, "Inferring "dark matter" and "dark energy" from videos," in *The IEEE International Conference on Computer Vision (ICCV)*, December 2013.
- [160] V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions, and reversals," in *Soviet physics doklady*, vol. 10, no. 8, 1966, pp. 707–710.
- [161] G. Saponaro, G. Salvi, and A. Bernardino, "Robot anticipation of human intentions through continuous gesture recognition," in *2013 International Conference on Collaboration Technologies and Systems (CTS)*, May 2013, pp. 218–225.
- [162] A. Gupta, A. Kembhavi, and L. S. Davis, "Observing human-object interactions: Using spatial and functional compatibility for recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 10, pp. 1775–1789, 2009.
- [163] S. Oh, A. Hoogs, M. Turek, and R. Collins, "Content-based retrieval of functional objects in video using scene context," in *European Conference on Computer Vision*. Springer, 2010, pp. 549–562.
- [164] M. W. Turek, A. Hoogs, and R. Collins, "Unsupervised learning of functional categories in video scenes," in *European Conference on Computer Vision*. Springer, 2010, pp. 664–677.
- [165] G. Zen, N. Rostamzadeh, J. Staiano, E. Ricci, and N. Sebe, "Enhanced semantic descriptors for functional scene categorization," in *Pattern Recognition (ICPR), 2012 21st International Conference on*. IEEE, 2012, pp. 1985–1988.
- [166] M. Wang and X. Wang, "Automatic adaptation of a generic pedestrian detector to a specific traffic scene," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 3401–3408.

- [167] N. Rhinehart and K. M. Kitani, "Learning action maps of large environments via first-person vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 580–588.
- [168] S. Pirk, V. Krs, K. Hu, S. D. Rajasekaran, H. Kang, Y. Yoshiyasu, B. Benes, and L. J. Guibas, "Understanding and exploiting object interaction landscapes," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 3, p. 31, 2017.
- [169] V. Delaitre, D. F. Fouhey, I. Laptev, J. Sivic, A. Gupta, and A. A. Efros, "Scene semantics from long-term observation of people," in *European conference on computer vision*. Springer, 2012, pp. 284–298.
- [170] D. F. Fouhey, V. Delaitre, A. Gupta, A. A. Efros, I. Laptev, and J. Sivic, "People watching: Human actions as a cue for single view geometry," *International journal of computer vision*, vol. 110, no. 3, pp. 259–274, 2014.
- [171] A. Pieropan, C. H. Ek, and H. Kjellstr  m, "Recognizing object affordances in terms of spatio-temporal object-object relationships," in *2014 IEEE-RAS International Conference on Humanoid Robots*, Nov 2014, pp. 52–58.
- [172] B. Yao and L. Fei-Fei, "Grouplet: A structured image representation for recognizing human and object interactions," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 9–16.
- [173] M. Stark, P. Lies, M. Zillich, J. Wyatt, and B. Schiele, "Functional object class detection based on learned affordance cues," *Computer Vision Systems*, pp. 435–444, 2008.
- [174] B. Leibe, A. Leonardis, and B. Schiele, "An implicit shape model for combined object categorization and segmentation," in *Toward category-level object recognition*. Springer, 2006, pp. 508–524.
- [175] H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins, "Text classification using string kernels," *Journal of Machine Learning Research*, vol. 2, no. Feb, pp. 419–444, 2002.
- [176] R. Hu, C. Zhu, O. van Kaick, L. Liu, A. Shamir, and H. Zhang, "Interaction context (icon): Towards a geometric functionality descriptor," *ACM Transactions on Graphics (TOG)*, vol. 34, no. 4, p. 83, 2015.
- [177] R. Jain and T. Inamura, "Bayesian learning of tool affordances based on generalization of functional feature to estimate effects of unseen tools," *Artificial Life and Robotics*, vol. 18, no. 1-2, pp. 95–103, 2013.
- [178] H. Laga, M. Mortara, and M. Spagnuolo, "Geometry and context for semantic correspondences and functionality recognition in man-made 3d shapes," *ACM Transactions on Graphics (TOG)*, vol. 32, no. 5, p. 150, 2013.
- [179] V. G. Kim, S. Chaudhuri, L. Guibas, and T. Funkhouser, "Shape2pose: Human-centric shape analysis," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 4, p. 120, 2014.
- [180] Z. Lun, E. Kalogerakis, R. Wang, and A. Sheffer, "Functionality preserving shape style transfer," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 6, p. 209, 2016.
- [181] M. Lin, T. Shao, Y. Zheng, N. J. Mitra, and K. Zhou, "Recovering functional mechanical assemblies from raw scans," *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 3, pp. 1354–1367, March 2018.
- [182] R. Hu, O. van Kaick, B. Wu, H. Huang, A. Shamir, and H. Zhang, "Learning how objects function via co-analysis of interactions," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 4, p. 47, 2016.
- [183] M. Savva, A. X. Chang, P. Hanrahan, M. Fisher, and M. Nie  sner, "Scenegrok: Inferring action maps in 3d environments," *ACM transactions on graphics (TOG)*, vol. 33, no. 6, p. 212, 2014.
- [184] H. O. Song, M. Fritz, D. Goehring, and T. Darrell, "Learning to detect visual grasp affordance," *IEEE Transactions on Automation Science and Engineering*, vol. 13, no. 2, pp. 798–809, 2016.
- [185] A. Chemero, "An outline of a theory of affordances," *ECOLOGICAL PSYCHOLOGY*, vol. 15, no. 2, pp. 181–195, 2003.
- [186] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth, "Describing objects by their attributes," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, June 2009.
- [187] G. Patterson and J. Hays, "The sun attribute database: Organizing scenes by affordances, materials, and layout," in *Visual Attributes*. Springer, 2017, pp. 269–297.
- [188] J. Wang, X. Zhu, S. Gong, and W. Li, "Attribute recognition by joint recurrent learning of context and correlation," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [189] Y. Li, C. Huang, C. C. Loy, and X. Tang, "Human attribute recognition by deep hierarchical contexts," in *European Conference on Computer Vision*. Springer, 2016, pp. 684–700.
- [190] T. Bagautdinov, A. Alahi, F. Fleuret, P. Fua, and S. Savarese, "Social scene understanding: End-to-end multi-person action localization and collective activity recognition," *arXiv preprint arXiv:1611.09078*, 2016.
- [191] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," in *Advances in neural information processing systems*, 2016, pp. 379–387.
- [192] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [193] —, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [194] A. Gupta, A. Kembhavi, and L. S. Davis, "Observing human-object interactions: Using spatial and functional compatibility for recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 10, pp. 1775–1789, Oct 2009.
- [195] —, "Observing human-object interactions: Using spatial and functional compatibility for recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 10, pp. 1775–1789, Oct 2009.
- [196] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [197] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy, "Speed/accuracy trade-offs for modern convolutional object detectors," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 3296–3297.
- [198] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, "Spatial transformer networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 2017–2025.
- [199] J. Johnson, L. Ballan, and L. Fei-Fei, "Love thy neighbors: Image annotation by exploiting image metadata," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4624–4632.
- [200] D. I. Kim and G. S. Sukhatme, "Semantic labeling of 3d point clouds with object affordance for robot manipulation," in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 5578–5584.
- [201] N. Krishnamoorthy, G. Malkarnenkar, R. J. Mooney, K. Saenko, and S. Guadarrama, "Generating natural-language video descriptions using text-mined knowledge." in *AAAI*, vol. 1, 2013, p. 2.
- [202] J. fran  ois Lalonde, N. V. D. F. Huber, and M. Hebert, "Natural terrain classification using three-dimensional lidar data for ground robot mobility," *Journal of Field Robotics*, vol. 23, pp. 839–861, 2006.
- [203] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Doll  r, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [204] G. Medioni, M.-S. Lee, and C.-K. Tang, *A computational framework for segmentation and grouping*. Elsevier, 2000.
- [205] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. Tsagarakis, "Object-based affordances detection with convolutional neural networks and dense conditional random fields," in *ResearchGate*, 2017.
- [206] F. Pittino, M. Driusso, A. D. Torre, and C. Marshall, "Outdoor and indoor experiments with localization using lte signals," in *2017 European Navigation Conference (ENC)*, May 2017, pp. 311–321.
- [207] J. Thomason, S. Venugopalan, S. Guadarrama, K. Saenko, and R. J. Mooney, "Integrating language and vision to generate natural language descriptions of videos in the wild." in *Coling*, vol. 2, no. 5, 2014, p. 9.
- [208] Y. Yang, C. L. Teo, H. Daum   III, and Y. Aloimonos, "Corpus-guided sentence generation of natural images," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2011, pp. 444–454.
- [209] A. Saxena, J. Driemeyer, and A. Y. Ng, "Robotic grasping of novel objects using vision," *The International Journal of Robotics Research*, vol. 27, no. 2, pp. 157–173, 2008.
- [210] J. Bohg and D. Kragic, "Grasping familiar objects using shape context," in *Advanced Robotics, 2009. ICAR 2009. International Conference on*. IEEE, 2009, pp. 1–6.

- [211] I. Lenz, H. Lee, and A. Saxena, "Deep learning for detecting robotic grasps," *The International Journal of Robotics Research*, vol. 34, no. 4-5, pp. 705–724, 2015.
- [212] C. C. Kemp and A. Edsinger, "Robot manipulation of human tools: Autonomous detection and control of task relevant features," in *Proc. of the Fifth Intl. Conference on Development and Learning*, 2006.
- [213] P. Fitzpatrick, G. Metta, L. Natale, S. Rao, and G. Sandini, "Learning about objects through action-initial steps towards artificial cognition," in *Robotics and Automation, 2003. Proceedings. ICRA'03. IEEE International Conference on*, vol. 3. IEEE, 2003, pp. 3140–3145.
- [214] T. Mar, V. Tikhonoff, G. Metta, and L. Natale, "Self-supervised learning of grasp dependent tool affordances on the icub humanoid robot," in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. IEEE, 2015, pp. 3200–3206.
- [215] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," 2013.
- [216] J. G. Wang, P. S. Mahendran, and E. K. Teoh, "Deep affordance learning for single- and multiple-instance object detection," in *TENCON 2017 - 2017 IEEE Region 10 Conference*, Nov 2017, pp. 321–326.
- [217] L. Montesano, M. Lopes, A. Bernardino, and J. Santos-Victor, "Modeling affordances using bayesian networks," in *Intelligent Robots and Systems, 2007. IROS 2007. IEEE/RSJ International Conference on*. IEEE, 2007, pp. 4102–4107.
- [218] K. Lai, L. Bo, X. Ren, and D. Fox, "A large-scale hierarchical multi-view rgb-d object dataset," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 1817–1824.
- [219] S. Fichtl, D. Kraft, N. Kruger, and F. Guerin, "Bootstrapping relational affordances of object pairs using transfer," *IEEE Transactions on Cognitive and Developmental Systems*, 2016.
- [220] H. Kjellström, J. Romero, and D. Kragić, "Visual object-action recognition: Inferring object affordances from human demonstration," *Computer Vision and Image Understanding*, vol. 115, no. 1, pp. 81–90, 2011.
- [221] V. Krunić, G. Salvi, A. Bernardino, L. Montesano, and J. Santos-Victor, "Affordance based word-to-meaning association," in *2009 IEEE International Conference on Robotics and Automation*, May 2009, pp. 4138–4143.
- [222] J. Gall, A. Fossati, and L. van Gool, "Functional categorization of objects using real-time markerless motion capture," in *CVPR 2011*, June 2011, pp. 1969–1976.
- [223] E. Ruiz and W. Mayol-Cuevas, "Geometric affordances from a single example via the interaction tensor," *arXiv preprint arXiv:1703.10584*, 2017.
- [224] P. Güler, Y. Bekiroglu, X. Gratal, K. Pauwels, and D. Kragić, "What's in the container? classifying object contents from vision and touch," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*. IEEE, 2014, pp. 3961–3968.
- [225] E. Ugur, E. Oztop, and E. Sahin, "Goal emulation and planning in perceptual space using learned affordances," *Robotics and Autonomous Systems*, vol. 59, no. 7-8, pp. 580–595, 2011.