Visual Affordance Learning for Robot Manipulation

Yuke Zhu Toyota Research Institute, August 2021



Traditional Form of Robot Automation



Structured Environments

Fixed Set of Tasks

Pre-programmed Procedures

General-Purpose Robot Autonomy: Our North Star Goal



Natural Environments

Ever-Changing Tasks

> Human Involvement



The Perception-Action Loop

A key challenge in robot autonomy is to close the perception-action loop.



[Gibson 1979; Bajcsy 1988; Ballard 1991; Espiau et al. 1992; Hutchinson et al. 1996; Hamel & Mahony 2002; Kragic & Christensen 2002; Jonschkowski & Brock 2015; Levine et al. 2016; Agrawal et al. 2016; Bojarski et al. 2016; Finn & Levine 2017; Florence et al. 2018]

The Perception-Action Loop



[Gibson 1979; Bajcsy 1988; Ballard 1991; Espiau et al. 1992; Hutchinson et al. 1996; Hamel & Mahony 2002; Kragic & Christensen 2002; Jonschkowski & Brock 2015; Levine et al. 2016; Agrawal et al. 2016; Bojarski et al. 2016; Finn & Levine 2017; Florence et al. 2018]

A key challenge in robot autonomy is to close the perception-action loop.

The Perception-Action Loop



[Detectron - Facebook Al Research]

Conventional Computer Vision



[Zeng et al., IROS 2018]

Physically-Grounded Robot Perception





"What makes a chair a chair?"



- 1. A chair has 4 legs... although I can imagine a chair with 3 or even 5
- 2. A chair usually has arms and a back... although a stool fits in the chair category without either
- 3. It must be raised above the ground... although things that aren't chairs also do this in a similar way...

By visual attributes



Object Perception



https://www.theuncomfortable.com/



By functional properties

Object Perception

"What makes a chair a chair?"

"What makes a chair a chair is its ability to be sat on."



Aristotle's *Ontology*

Zhu Xi's *Li*(理)

"What makes a chair a chair?"

Affordances are *possibilities* for actions that the environment *affords* to the agent. (Gibson, 1977)

[Gibson 1979; Kirklik 1993; Zaff 1995; Stoytchev 2015; Amant 1999; Bousmalis 2018; Detry 2011; Zhu et al. 2017; Mahler 2017; Nagarajan 2020; Ugur 2007; Dang 2020; Fang 2018; Song 2010; Zeng 2018; Abel 2014; Abel 2015; Cruz 2016; Khetarpal 2020; Ardon 2020; Mandikal 2020]

Object Perception





Affordance Learning in Robotics





Semantic Affordance

[Zhu et al. 2014, Varadarajan & Vincze 2012]

[Katz et al. 2013; Do et al. 2018]

Prior methods focus on learning from human supervision and building staged pipelines.

[Gibson 1979; Kirklik 1993; Zaff 1995; Stoytchev 2015; Amant 1999; Bousmalis 2018; Detry 2011; Zhu et al. 2017; Mahler 2017; Nagarajan 2020; Ugur 2007; Dang 2020; Fang 2018; Song 2010; Zeng 2018; Abel 2014; Abel 2015; Cruz 2016; Khetarpal 2020; Ardon 2020; Mandikal 2020]



Detection & Segmentation

Affordance Template

[Hart et al. 2015; Pohl et al. 2020]





scene abstractions

Bridging visual perception and robot action through visual affordance

learning signals

6-DoF Robotic Grasping in Clutter

- Important modules in robot manipulation
 - Bin Picking
 - Part Assembly
 - Logistics
- Input: Partial point cloud of workspace
- Output: 6-DoF grasp pose (3D position and orientation)

How to Predict Grasps

[Bohg et al. 2011, Varley et al. 2017, Lundell et al. 2019]

Geometry Analysis

- Analytical solution
- Require full 3D model

$Reconstruction \rightarrow Grasp$

- Working with visual input
- Information bottleneck

[Mahler et al. 2017, Morrison et al. 2018, Liang et al. 2019, Breyer et al. 2020]

End-to-End Deep Learning

- High grasp performance
- Limited generalization

Affordance and geometry reasoning are not isolated

Likelihood of grasp success and grasp parameters

Affordance

Predict affordance of reconstructed part

Geometry

Reconstruct graspable region

Input TSDF

Implicit Neural Representations for 3D Shapes

- Functions that map from coordinate to quantities (SDF, occupancy).
- Functions are parametrized with neural networks.
- Shape bound is defined by level set of the parametrized functions.

DeepSDF, Park et al. CVPR 2019

Implicit Neural Representations: Advantages

Continuous & memory-efficient

• End-to-end differentiable

Structured Implicit Functions

Global implicit function

- Single global feature
- Implicit function conditioned on global feature
- Overly smooth reconstruction

Structured implicit function

ConvONets, Peng et al. ECCV 2020

- Structured feature grid
- Local features linearly sampled from the feature grid
- Fine-grained local details

GIGA Architecture

Input TSDF

3D feature grid

Projection

Aggregation

Projected 2D features

Structured feature grids

GIGA Architecture

Occupancy probability

Self-Supervised Data Collection

- Grasp affordance labels from physical trials in simulation.
- 3D geometry labels from ground-truth object meshes.

Packed scene

Pile scene

Quantitative Comparison

- Geometry learning facilitates affordance learning
- Continuity of implicit function enables higher performance

Grasp Success Rate ↑

- ■VGN (Breyer et al. 2020)

Geometry Learning Facilitates Occluded Grasps

Reconstruction Focuses on Graspable Parts

■ IoU-Grasp

GIGA-Detach Affordance Only

GIGA Affordance and Reconstruction

GIGA-Geo Reconstruction Only

- Synergies between **affordance** and **geometry** Better grasp prediction, especially in occluded regions 3D reconstruction focuses on action-relevant parts
- Structured implicit neural representation Continuous and compact representation for both affordance and geometry Combine voxel grids with neural implicit functions

Code and models are available at https://sites.google.com/view/rpl-giga2021

GIGA: Summary

Classical notion of affordance is not suitable for planning

Classical notion of affordance is not suitable for planning

feasible? **YES**

Classical affordance: whether an action is *feasible* No way to choose actions with respect to a long-horizon task goal

Classical notion of affordance is not suitable for planning

enables fetching red cube? NO

Classical affordance: whether an action is *feasible* **Our new affordance:** will an action **make future actions feasible**? Choose actions that enable the subsequent steps in the task plan

Skill Affordances

Parameterized skill plan grasp(θ) \rightarrow hook(θ) \rightarrow grasp(θ) $\rightarrow \cdots$

Affordance function:

$$\mathcal{A}_{\pi, heta}(s)$$

Whether state *S* is in the affordance set of (π, θ)

Xu et al. "Deep Affordance Foresight" ICRA 2021

Verifying that a plan is executable with affordances Parameterized skill plan grasp(θ) \rightarrow hook(θ) \rightarrow grasp(θ) $\rightarrow \cdots$ How likely is this plan executable by the robot? The plan is executable if every skill in the plan is afforded.

Probal

bility that a plan is executable:

$$p_{lan}(\{(\pi_i, \theta_i)\}_{i=1}^N) = \sum_{s \in S} Z_{N-1}(s) \mathcal{A}_{\pi_N, \theta_N}(s)$$
affordance like

Planning towards a Goal

We define a goal $g \in G$ as a condition function that checks whether a state satisfies the goal. We denote the set of states that satisfies g as S_q .

Any plan that ends with a skill whose affordance set is a *subset* of S_g are goal-directed plans for g.

We denote all plans directed at goal g as \mathcal{P}_g

Planning with Affordance

Probability that a plan is executable:

 $C_{plan}(\{(\pi_i, \theta_i)\}_{i=1}^N)$

Search for goal-directed plans that are most likely executable: $\operatorname{arg\,max}_{p\in \mathbb{C}}$

How do we plan in a partially-observed environment with unknown dynamics?

$$) = \sum_{s \in S} Z_{N-1}(s) \mathcal{A}_{\pi_N,\theta_N}(s)$$

$$\in \mathcal{P}_g C_{plan}(p)$$

Xu et al. "Deep Affordance Foresight" ICRA 2021

DAF with Model-Predictive Control (MPC)

44

Experiments: Tool-Use + Stacking

Results: Tool-Use

PlaNet

Sample Rollouts: Tool-Use + Stacking

DAF (Ours)

Experiments: Kitchen

coffee beans dispenser

coffee machine

 Affordances for fetching the mug is shared among the two tasks. Our method should be able to learn the task much faster & better

Results: Kitchen

get tea

get coffee

Rollouts: Kitchen (Coffee)

PlaNet

DAF (Ours)

DAF is reactive and can recover from mistakes

Transfer learned affordances

base task: get tea

transfer task: get coffee

Deep Affordance Foresight: Summary

• New notion of **skill affordance**

action abstraction that supports long-horizon planning estimated from visual observations

 Latent-space planning algorithm
 DAF can plan through raw perception and complex non-rigid dynamics

Bridging visual perception and robot action through visual affordance

human-centered

(human environments are shaped by affordance)

physically-grounded

- (informative representation
 - for robot actions)

task-agnostic

(inherent properties that facilitate knowledge transfer)

Architectures of Robot Manipulation

Acknowledgement

UT Robot Perception and Learning Lab https://rpl.cs.utexas.edu/

Zhenyu Jiang

Yifeng Zhu

Kuan Fang

Ajay Mandlekar

Roberto Martín-Martín

Soroush Nasiriany

Maxwell Svetlik

Danfei Xu

Fei-Fei Li

Silvio Savarese

