

UT Robot Perception & Learning Lab



- Objects, Skills, and the Quest for
- **Compositional Robot Autonomy**

Yuke Zhu February 18, 2022

TEXASRobotics



Two Worlds of Building Robot Autonomy

James Webb Space Telescope



robotic systems engineering

- special-purpose robotic systems
- high engineering cost: over 9,000,000,000 dollars
- decades of teamwork: 20+ years, originally planned for 2013, launched in the final week of 2021
- extremely complex yet reliable system: over 300 are

"single points of failure"

"The James Webb Space Telescope — making 300 points of failure reliable" Robert Barron



Two Worlds of Building Robot Autonomy

- general-purpose robotic systems
- intractable to be manually engineered
- months of research effort: small team of graduate

students and researchers

• one-off and unreliable system: publish a paper with

~70% success rate and move on to the next one

General-Purpose Collaborative Robots



robot learning research

Two Worlds of Building Robot Autonomy

James Webb Space Telescope



robotic systems engineering

General-Purpose Collaborative Robots



robot learning research

Key to deployable robot autonomy in the wild!

James Webb Space Telescope



robotic systems engineering

General-Purpose Collaborative Robots



robot learning research

Lessons from Systems Engineering Principles

Abstraction



"division of labor"

Composition



"harmony of labor"

"Types, abstraction, and parametric polymorphism." John Reynolds, 1983 "What Are Abstractions in Software Engineering with Examples." Matthieu Cneude, 2019 "Between abstraction and composition..." Jonathan Sterling, 2021



 π :



the "pixels to torques" approach





 π :



compositional robot autonomy stack





[Source: Detectron2, FAIR]

composition of objects

 π :



compositional robot autonomy stack





[Source: Daniel M. Wolpert]

composition of skills

 π :



compositional robot autonomy stack



Neural Task Programming (NTP): Hierarchical Policy Learning as Neural Program Induction



Objects: color cubes

Skills: pick and place

 \rightarrow Strong generalization!

"Neural Task Programming: Learning to Generalize Across Hierarchical Tasks." Xu et al. ICRA 2018 "Neural Task Graphs: Generalizing to Unseen Tasks from a Single Video Demonstration." Huang et al. CVPR 2019





Neural Task Programming (NTP): Hierarchical Policy Learning as Neural Program Induction

Objects: color b

Skills: pick and

q dener Stror

can be clearly defined, like the block world.

- **Observation:** Compositionality works well in
- domains where a finite set of objects and skills
- **Question:** Can we "unblock the block world"?





v. Observation	Ir	nput Task Spec.
: block_stacking	EOP: False	
: pick_and_place	Output Task Spec.	
Env. Observation		Input Task Spec
P _{in} : place		EOP: True
P _{out} : release		Args: N/A
release()		













Compositional task modeling with skill abstractions



BUDS [Zhu et al., RA-L 2022]

Library of Behavior



MAPLE [Nasiriany et al., ICRA 2022]





Object Representations for Robotic Grasping



- Input: Partial point cloud of workspace
- Output: 6-DoF grasp pose (3D position and orientation)





ace tion

Affordance and geometry reasoning are connected.

Likelihood of grasp success and grasp parameters

Affordance

Predict affordance of reconstructed part



Geometry

Reconstruct graspable region

GIGA: Grasping via Implicit Geometry and Affordance

Input TSDF





Grasp affordance



Neural Fields for 3D Scenes

- Neural networks that map from coordinate to quantities (SDF, occupancy)
- Smooth, continuous encoding of 3D scenes
- Fully differentiable models that can be trained with rich supervisions



DeepSDF, Park et al. CVPR 2019



GIGA: Grasping via Implicit Geometry and Affordance



Self-Supervised Data Collection

- Affordance labels from grasp trials in simulation.



Packed scene

• 3D geometry labels from ground-truth object meshes.



Pile scene



Geometry Learning Facilitates Occluded Grasps.



failure





Reconstruction Focuses on Graspable Parts



■ IoU-Grasp



GIGA-Detach Affordance Only

GIGA Affordance and Reconstruction

GIGA-Geo Reconstruction Only





- Task-dependent: optimized for downstream grasping tasks
- Multimodal: captures synergies between affordance and geometry for robotic grasping.
- Structured: combines 3D voxel grids with neural field methods
- Self-supervised: trained with self-supervised interactions



Zhenyu Jiang

"Synergies Between Affordance and Geometry: 6-DoF Grasp Detection via Implicit Representations." Jiang, Zhu, Svetlik, Fang, and Z. Code and models are available at https://sites.google.com/view/rpl-giga2021

GIGA: Learning Object Representations





Affordance

High Affordance





Can we build object models from interaction?



Interaction creates novel sensory stimuli for object learning.



Can we build object models from interaction?



physical object in universe

digital twin creation



virtual object in "metaverse"

Ditto: Digital Twin of Articulated Objects

Before interaction

After interaction



point cloud

part-level segmentation & 3d geometry

articulation parameters

Ditto Architecture





Recreated Digital Twins of Articulated Objects







laptop







oven







faucet







drawer



From Real World to Simulation and Back



https://robosuite.ai





Ditto: Building Interactive Object Models

- **Embodied interactions:** Emit useful sensory information for understanding an object
- Structured implicit neural representations: Jointly encode geometry and articulation
- **Digital twins:** Bridge simulation and the real world



Zhenyu Jiang





"Building Digital Twins of Articulated Objects from Interaction." Jiang, Hsu, and Z. Code and models are available at https://ut-austin-rpl.github.io/Ditto





Learning object abstractions from embodied interactions











sensorimotor skills: perceptually grounded, temporally extended behaviors





The Context Principle: words have meaning only as constituents of (hence, presumably, only in virtue of their use in) sentences

The Compositionality Principle: the meaning of the whole sentence is a function of the meanings of its parts

"How to Stop Worrying About Compositionality", Aurelie Herbelot



BUDS: Bottom-Up Discovery of sensorimotor Skills



Unstructured demonstrations

No temporal annotations

Sensorimotor skills

Composition of skills

Operating on raw sensory data





BUDS: Bottom-Up Discovery of sensorimotor Skills

Hierarchical task structure



Unstructured demonstration

Bottom-up approach

- Building hierarchical task structures
- Agglomerative clustering on temporal segments w/ multi-sensory cues
- Identify reusable skills from recurring temporal patterns

ures nporal



Compose Skills to Solve Tasks

Cluster segments into skills

meta controller

Skils kodex

Subgoal

imageompute reaction or each input temporal segment

arn

Temporal segments from set of demonstrations



Learned skills improve task success.

Behavioral Cloning (Zhang et al. 2018)







24.4% success

23.4% success

Trained on 30min demonstrations for each task \bullet

CHAMP (Niekum et al. 2015)





Multi-sensory cues improve skill segmentation.



Task Success Rate (%)

Kitchen



"Event Structure in Perception and Conception." Zacks and Tversky, 2001



Multi-task learning improves quality of skills.



Multi-task learning in Kitchen leads to an 8% increase in success rate.

Discovered skills are reusable in novel tasks.



Through reusing the skills while re-training the meta-controller.

Real Robot: Autonomous Execution

Workspace View





- Discover reusable sensorimotor skills from unstructured lacksquaretask demonstrations
- Multi-sensory, multi-task learning improves skill quality
- Sample-efficient learning of complex behaviors using 30min of demonstrations



BUDS: Bottom-Up Sensorimotor Skill Discovery





"Bottom-Up Skill Discovery from Unsegmented Demonstrations for Long-Horizon Robot Manipulation." Zhu, Stone, and Z. Code and models are available at https://github.com/UT-Austin-RPL/BUDS











The Context Principle: words have meaning only as constituents of (hence, presumably, only in virtue of their use in) sentences

The Compositionality Principle: the meaning of the whole sentence is a function of the meanings of its parts



The Context Principle: words have meaning only as constituents of (hence, presumably, only in virtue of their use in) sentences

The Compositionality Principle: the meaning of the whole sentence is a function of the meanings of its parts

BUDS



The Context Principle: words have meaning only as constituents of (hence, presumably, only in virtue of their use in) sentences

The Compositionality Principle: the meaning of the whole sentence is a function of the meanings of its parts



- Heterogenous library of primitives
- Primitives are parameterized with inputs. For example, Grasp(x, y, z, ψ): reach location (x, y, z) at angle ψ

- Heterogenous library of primitives
- Primitives are parameterized with inputs. For example, Grasp(x, y, z, ψ): reach location (x, y, z) at angle ψ
- Atomic primitive dedicated to lowlevel motor actions

Selected Primitive

Atomic (0% Success)

MAPLE (Non-Atomic) (0% Success)

Evaluation: Peg Insertion

DAC (16% Success)

MAPLE (ours) (100% Success)

Evaluation: Compositionality

Evaluation: Compositionality

From Digital Twin Training to Real Robot Deployment

MAPLE Policy Trained in Digital Twin

"Fast Uncertainty Quantification for Deep Object Pose Estimation." Shi, Zhu, Tremblay, Birchfield, Ramos, Anandkumar, and Z.

Zero-Shot Policy Transfer to Real Robot

- Heterogeneous behavior primitives to scaffold longhorizon manipulation tasks
- Hierarchical policy to invoke behavior primitives as modular APIs
- Emergent compositional structures from primitive composition

MAPLE: Primitive-Augmented Reinforcement Learning

"Augmenting Reinforcement Learning with Behavior Primitives for Diverse Manipulation Tasks." Nasiriany, Liu, and Z. Code and models are available at https://github.com/UT-Austin-RPL/maple

What is an **Object**?

- Context-, task- dependent
- Multi-modal, multi-faceted
- Learned from embodied interaction

What is a Skill?

- Context vs Compositionality principle
- Versatility and reusability
- Temporal abstraction and perceptual grounding

Al Architectures of Robot Autonomy

Acknowledgement UT Robot Perception and Learning (RPL) Lab

https://rpl.cs.utexas.edu/

Google

