# TEXAS Robotics

**RPL** UT Robot Perception & Learning Lab

**LARG** Learning Agents Research Group — The University of Texas at Austin

# GROOT saves you from retraining visuomotor policies whenever you change backgrounds, move cameras, or use new objects!

GROOT: Learning Generalizable Manipulation Policies with Object-Centric 3D Representations

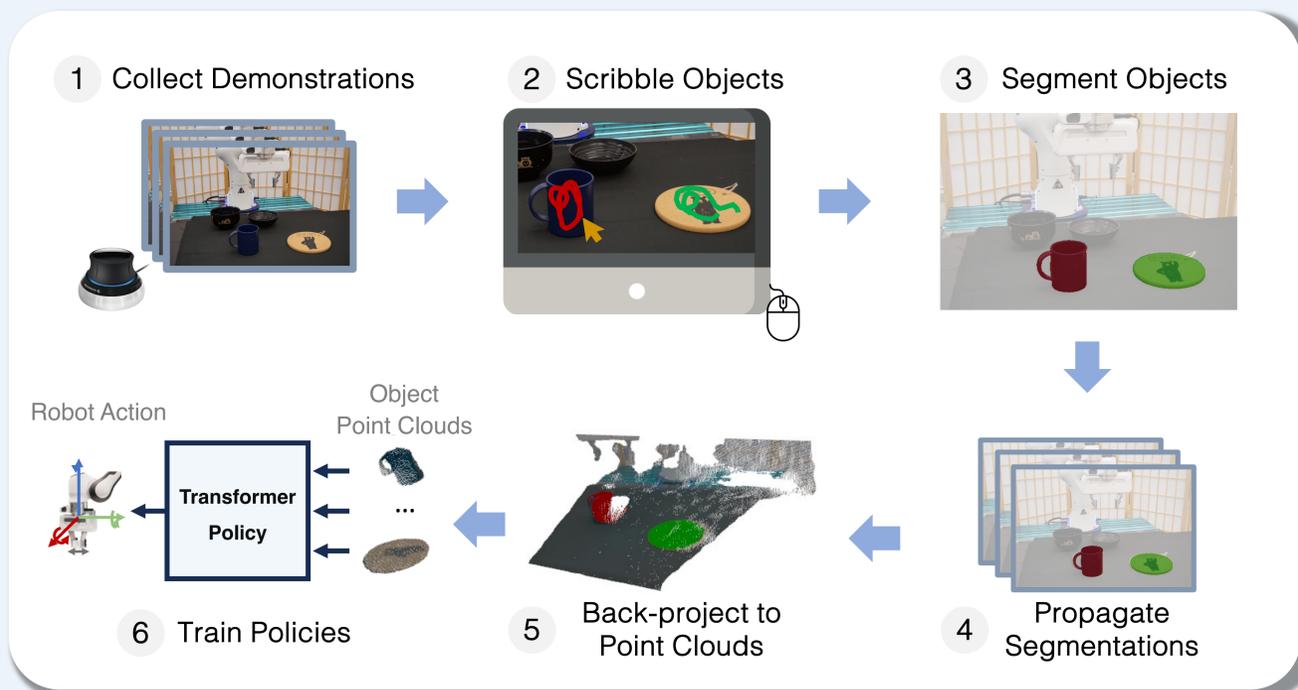Yifeng Zhu, Zhenyu Jiang, Peter Stone, Yuke Zhu

## Motivation

- Prior imitation learning methods fall short in generalization beyond the training conditions
- We want to avoid the time-consuming data re-collection and model re-training in every new setting

## Insights

- Object-centric:
  - Exploit compositional structure of visual scenes in objects and entities
  - Attend to task-relevant objects, minimize visual distractions
- 3D-aware:
  - Lift the spatial reasoning from the 2D plane to a unified reference frame of 3D coordinates
- New object generalization:
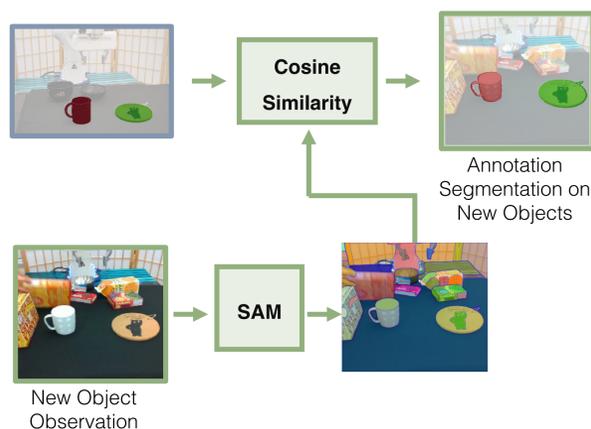  - Use the open-vocabulary visual understanding of vision foundation models

## GROOT Overview



1. Collect Demonstrations
2. Scribble Objects
3. Segment Objects
4. Propagate Segmentations
5. Back-project to Point Clouds
6. Train Policies

Robot Action · Transformer Policy · Object Point Clouds

## Generalization of GROOT



Various Visual Backgrounds
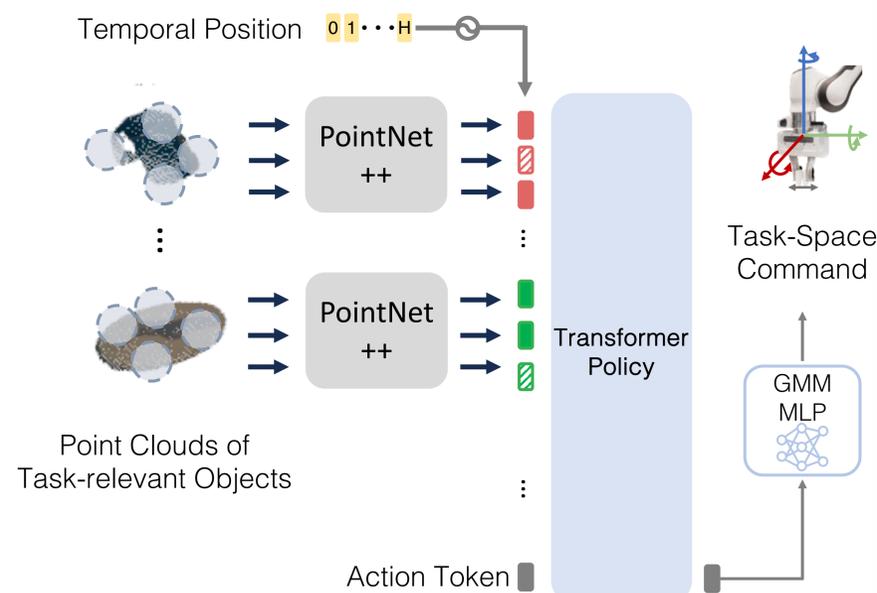
Different Camera Angles

New Object Instances

## Segmentation Correspondence Model

*Only used during deployment

- SCM helps identify new object
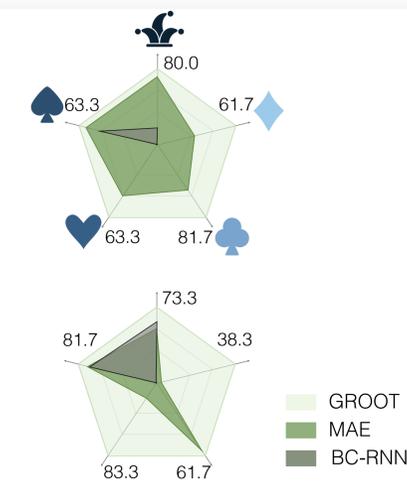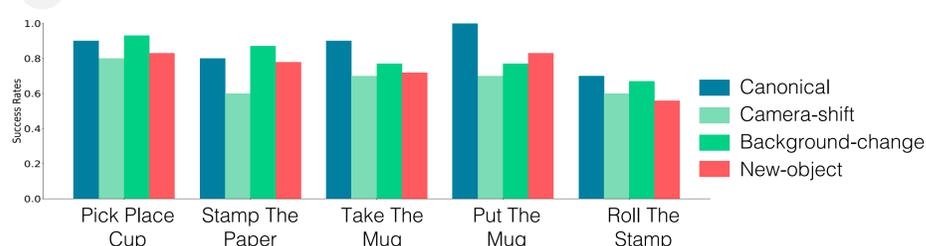- The new annotation segmentation then propagates to new objects



Cosine Similarity · Annotation Segmentation on New Objects

New Object Observation · SAM

## Transformer Policy



Temporal Position 0 1 ··· H

PointNet ++

Transformer Policy

Task-Space Command

GMM MLP

Point Clouds of Task-relevant Objects

Action Token

## Experiments

### A Simulation



Put the moka pot on the stove

Reposition the yellow and white mug

Radar charts: 80.0 · 63.3 · 61.7 · 63.3 · 81.7

73.3 · 81.7 · 38.3 · 83.3 · 61.7

Legend: GROOT · MAE · BC-RNN

### C New Object Generalization



### B Real Robot



Success Rates — Pick Place Cup, Stamp The Paper, Take The Mug, Put The Mug, Roll The Stamp

Legend: Canonical · Camera-shift · Background-change · New-object