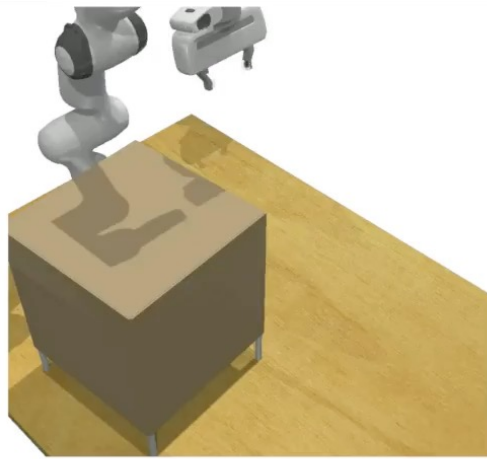
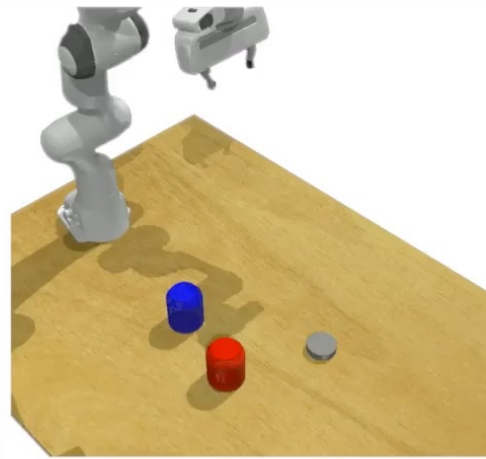


Acting with Perception and Language

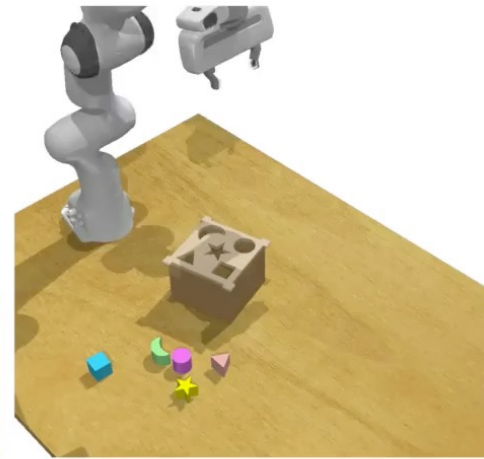
Mohit Shridhar



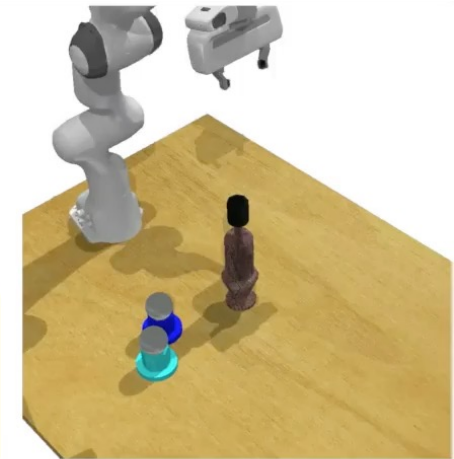
"open the middle drawer"



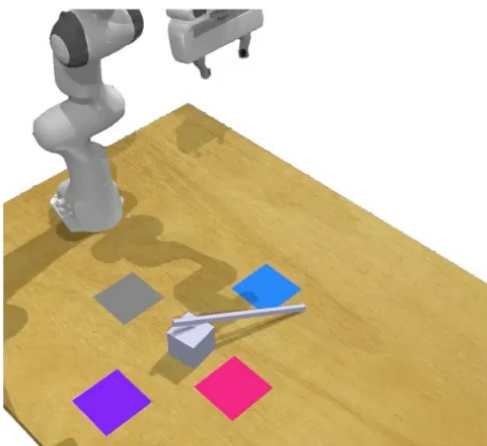
"close the blue jar"



"put the cylinder
in the shape sorter"



"screw in the
blue light bulb"



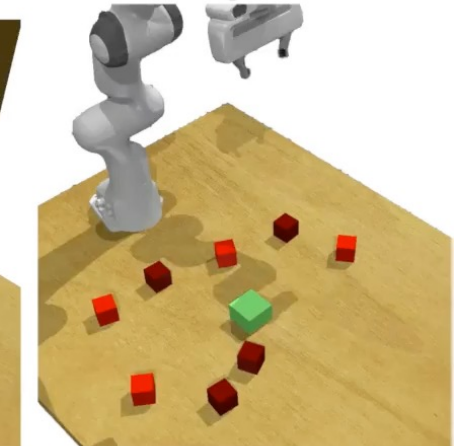
"use the stick to drag the cube
onto the gray target"



"take the steak off the grill"



"put the mustard
in the cupboard"



"stack 2 maroon blocks"

So far ...



Object Detection



Object Pose Estimation

What are the **objects** here?



Credit: Tasty (Youtube)

There is something weird
about **objects** 🤔

David Marr



“Vision is a computational process that transforms the retinal image into an objective representation of 3D shape.”

Other Perspectives?



Prof. Lana Lazebnik
UIUC



Computer Vision: Looking Back to Look Forward

Svetlana Lazebnik
IRIM Short Course
Spring 2020

<https://slazebni.cs.illinois.edu/spring20>

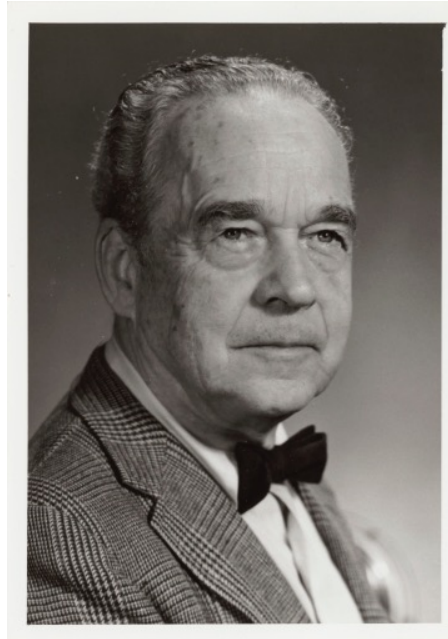
Three Perspectives on Vision

David Marr



“Vision is a computational process that transforms the retinal image into an objective representation of 3D shape.”

James Gibson



“There is no computation. There is no retinal image. There are no representations. There is no 3D shape. There is only direct pickup of ecologically relevant variants and invariants. Vision is in the world, not the observer.”

Affordances



Can I **stretch** this?

Can I **tear** this?

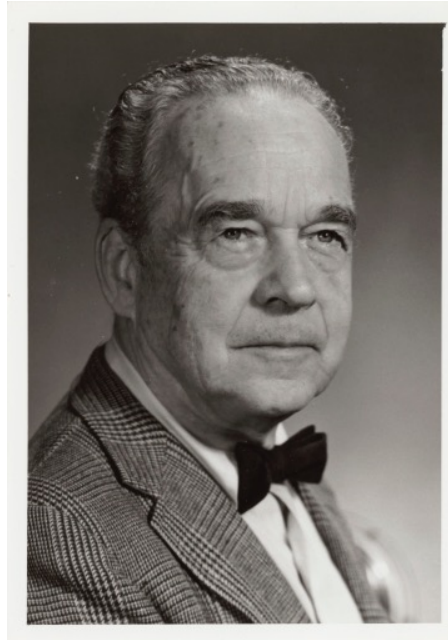
Three Perspectives on Vision

David Marr



“Vision is a computational process that transforms the retinal image into an objective representation of 3D shape.”

James Gibson



“There is no computation. There is no retinal image. There are no representations. There is no 3D shape. There is only direct pickup of ecologically relevant variants and invariants. Vision is in the world, not the observer.”

Jan Koenderink



“There is no objective world, only the observer’s *umwelt*. Thus, vision cannot be in the world but is a creative act of the observer.”

Umwelt

*An ant has difference
experience and goals*



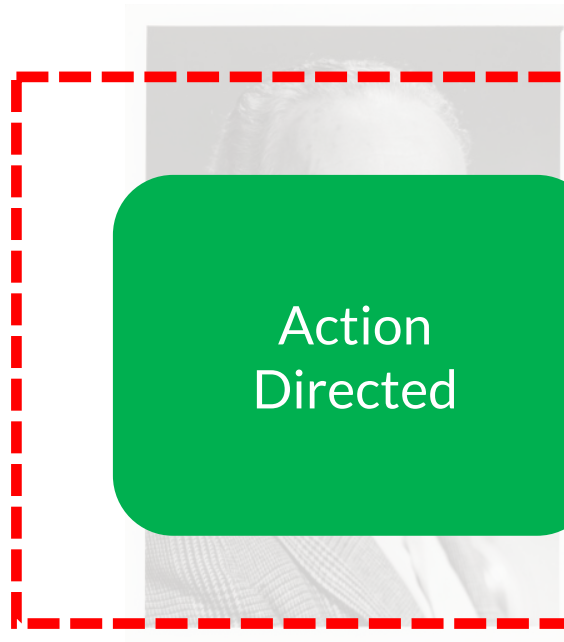
Three Perspectives on Vision

David Marr



“Vision is a computational process that transforms the retinal image into an objective representation of 3D shape.”

James Gibson



“There is no computation. There is no retinal image. There are no representations. There is no 3D shape. There is only direct pickup of ecologically relevant variants and invariants. Vision is in the world, not the observer.”

Jan Koenderink



“There is no objective world, only the observer’s *umwelt*. Thus, vision cannot be in the world but is a creative act of the observer.”

This Lecture

CASE STUDY

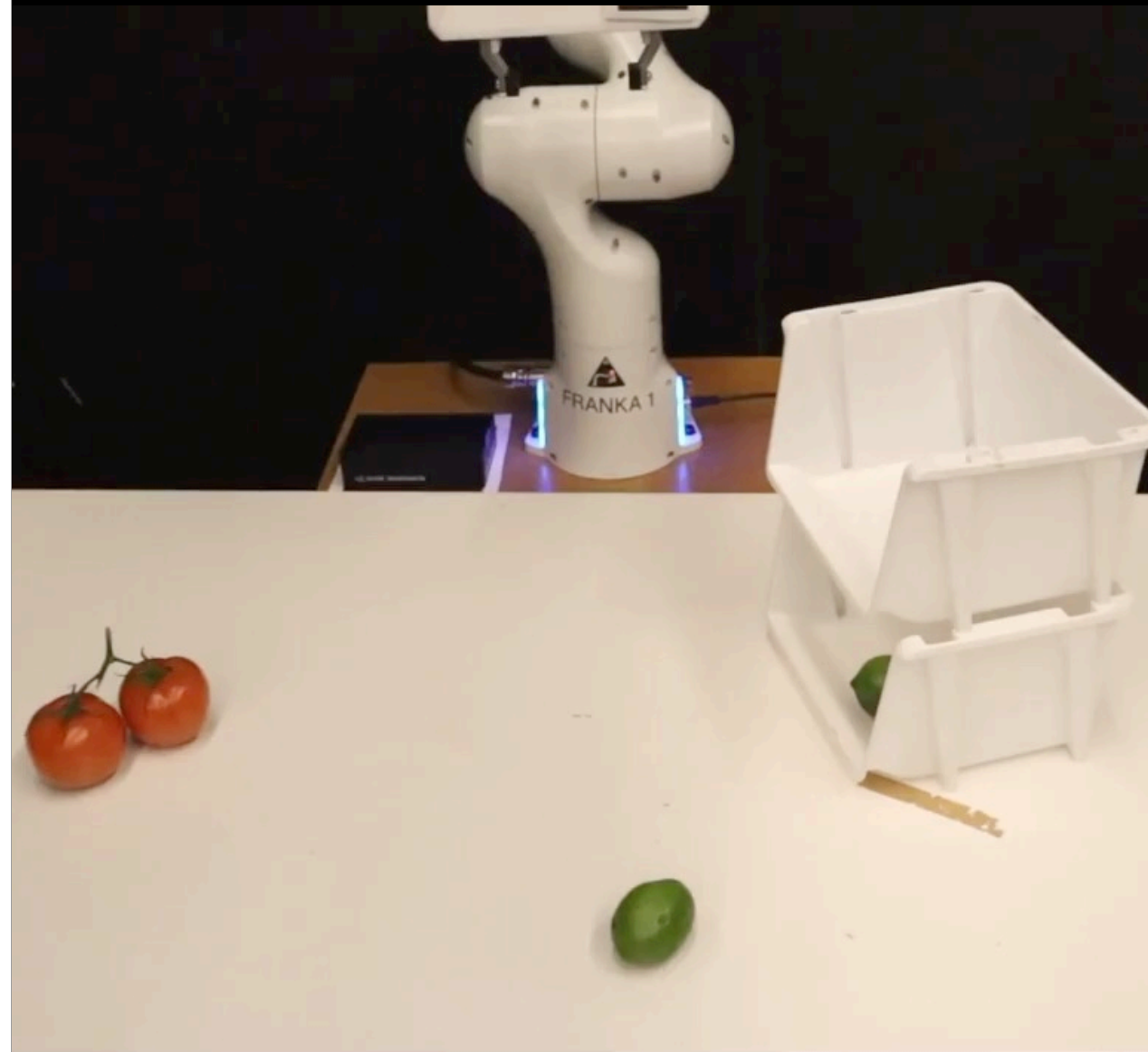
Perceiver-Actor: A Multi-Task Transformer for Robotic Manipulation

Mohit Shridhar¹, Lucas Manuelli², Dieter Fox^{1, 2}

¹University of Washington, ²NVIDIA

Perceiver-Actor

Multi-task 6-DoF manipulation agent

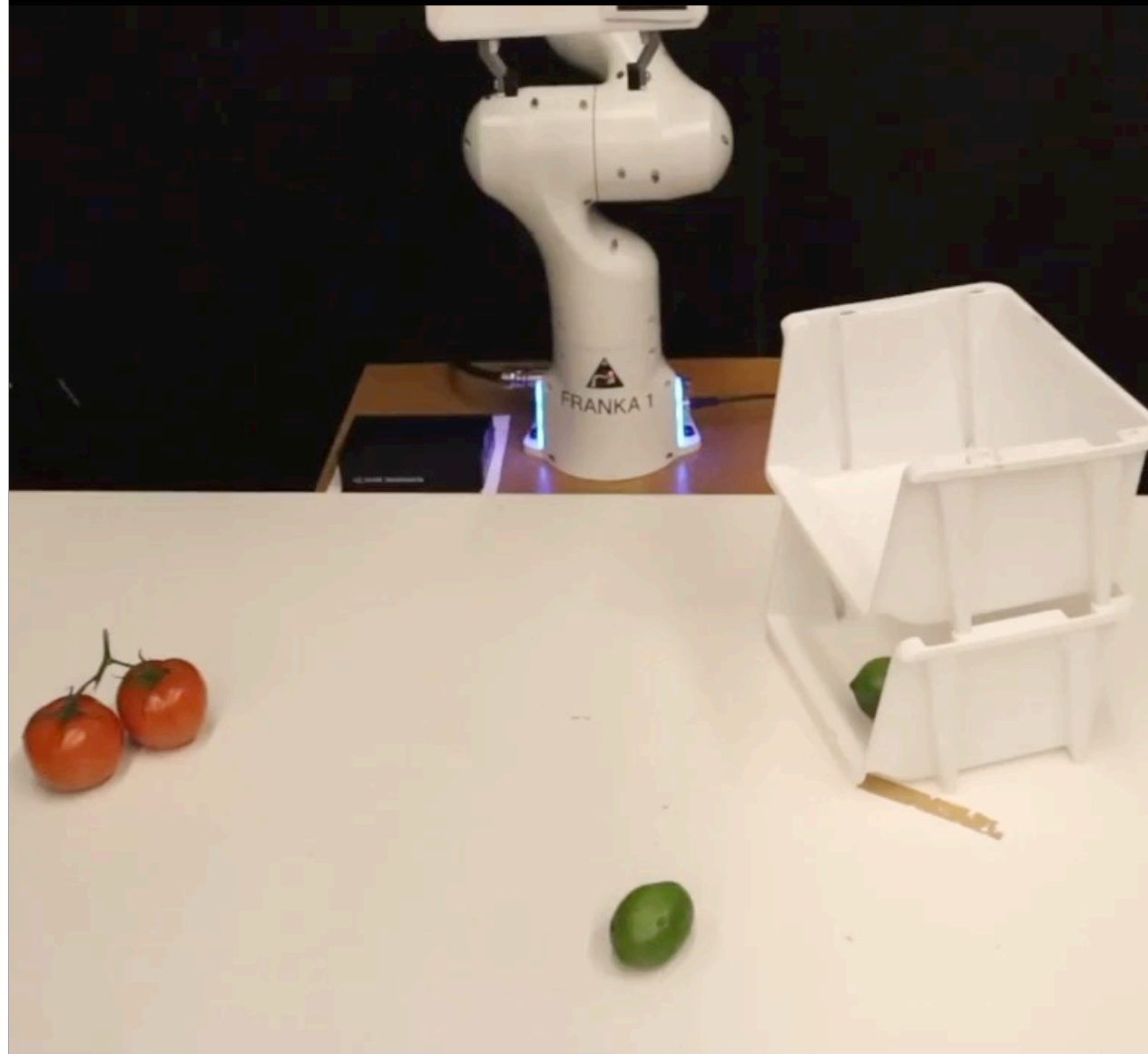


“put the tomatoes in the top bin”

Perceiver-Actor

Multi-task 6-DoF manipulation agent

End-to-end few-shot imitation learning



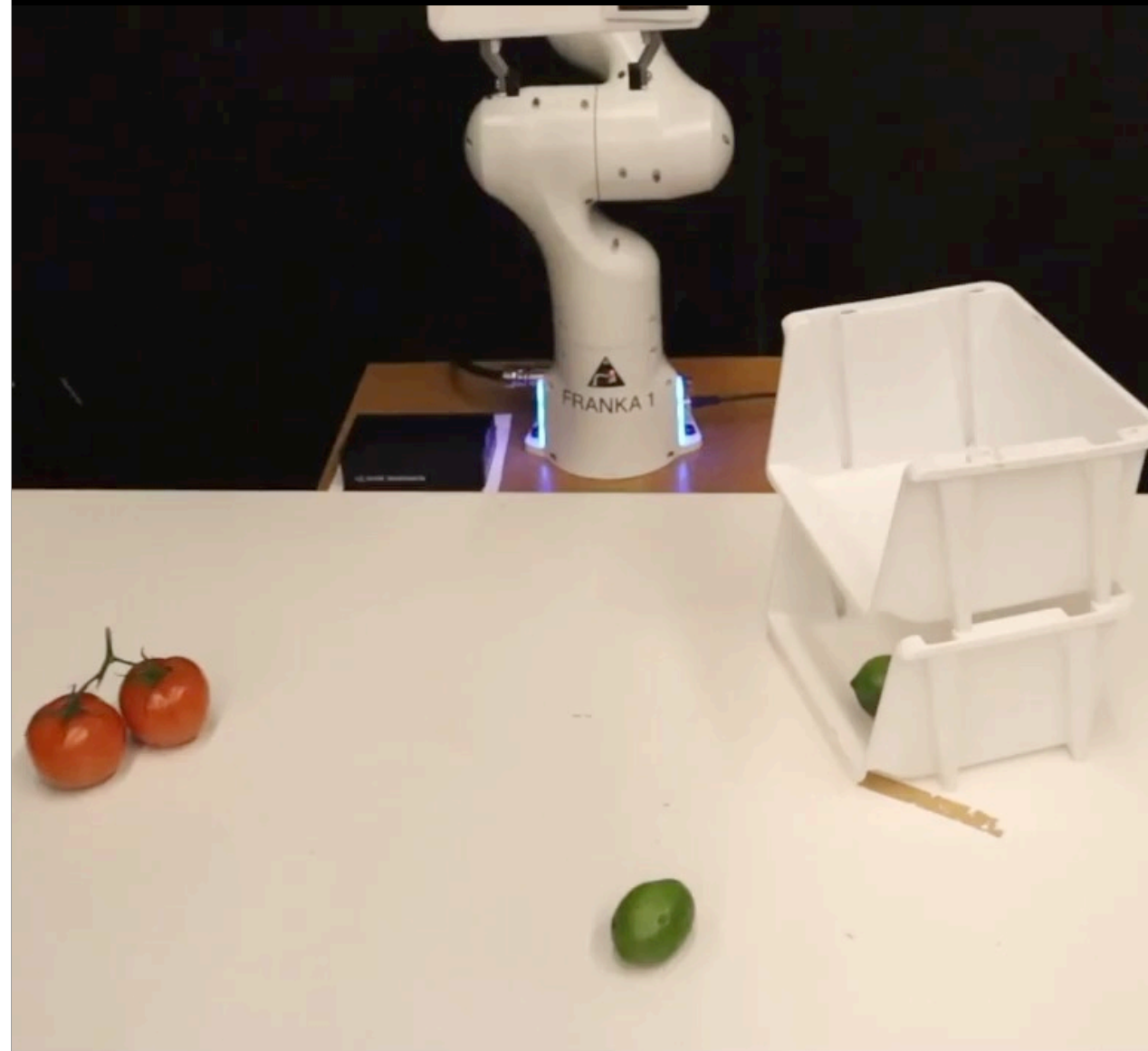
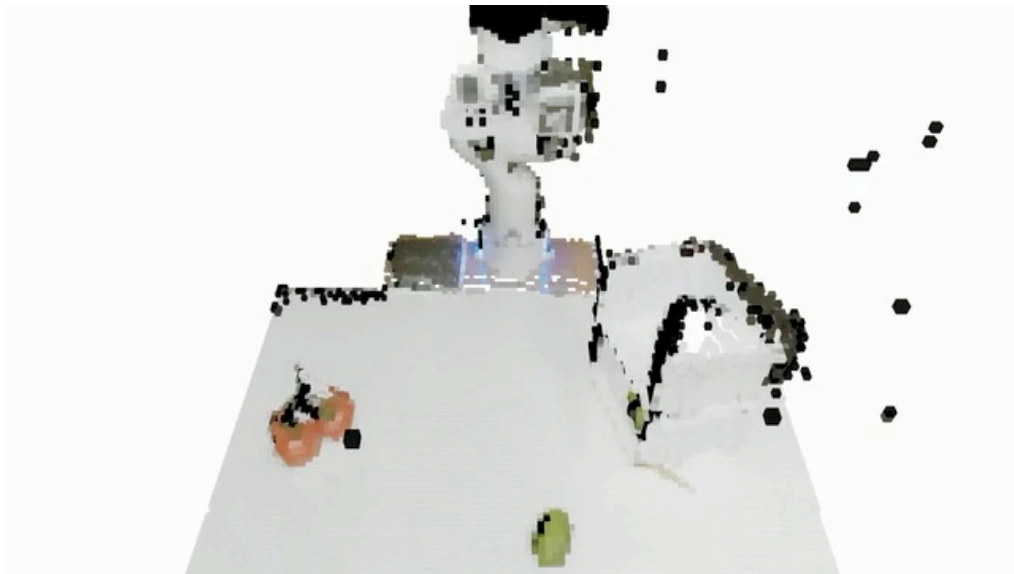
“put the tomatoes in the top bin”

Perceiver-Actor

Multi-task 6-DoF manipulation agent

End-to-end few-shot imitation learning

Input: RGB-D Voxels & [Language Goal](#)



“put the tomatoes in the top bin”

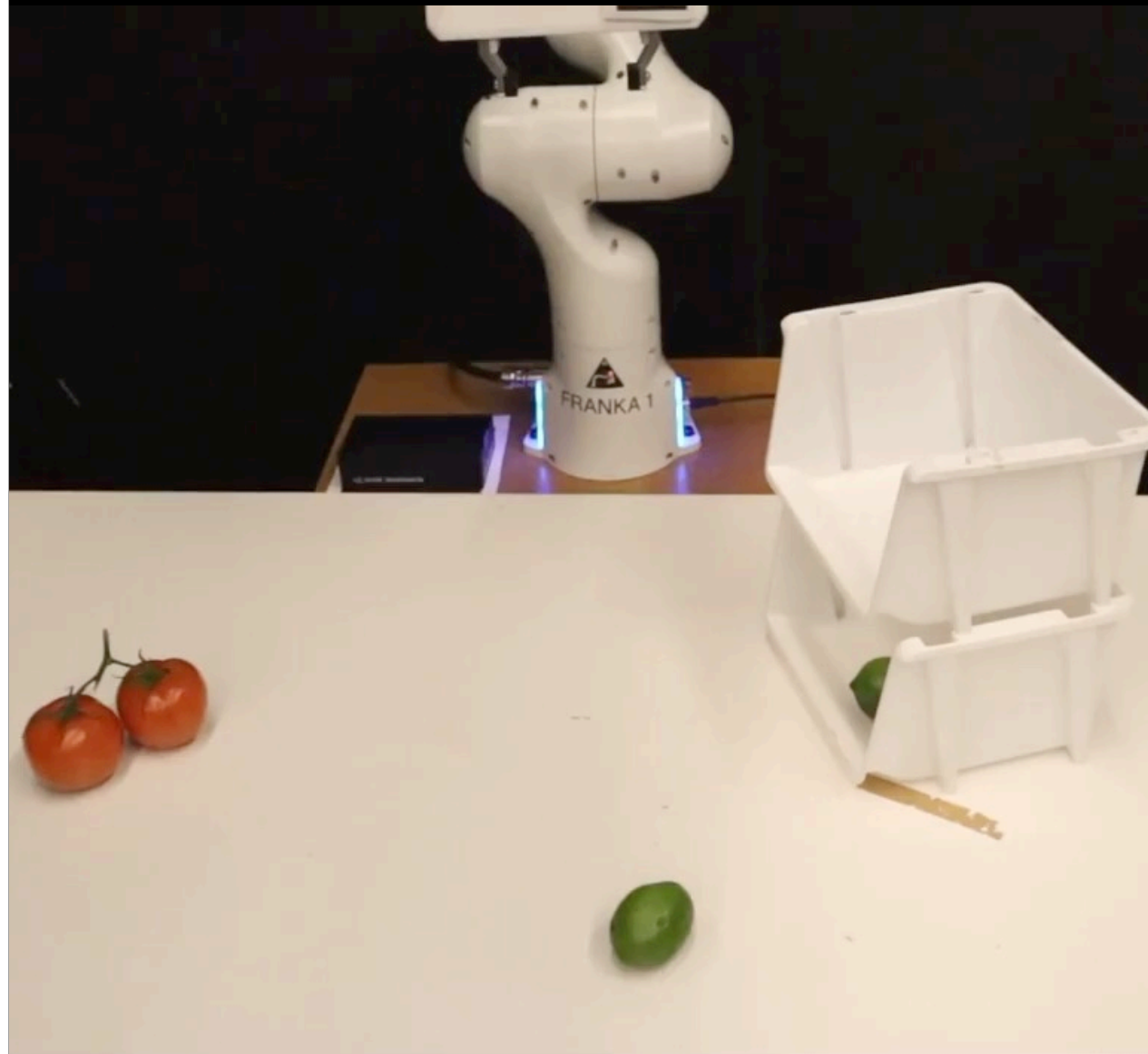
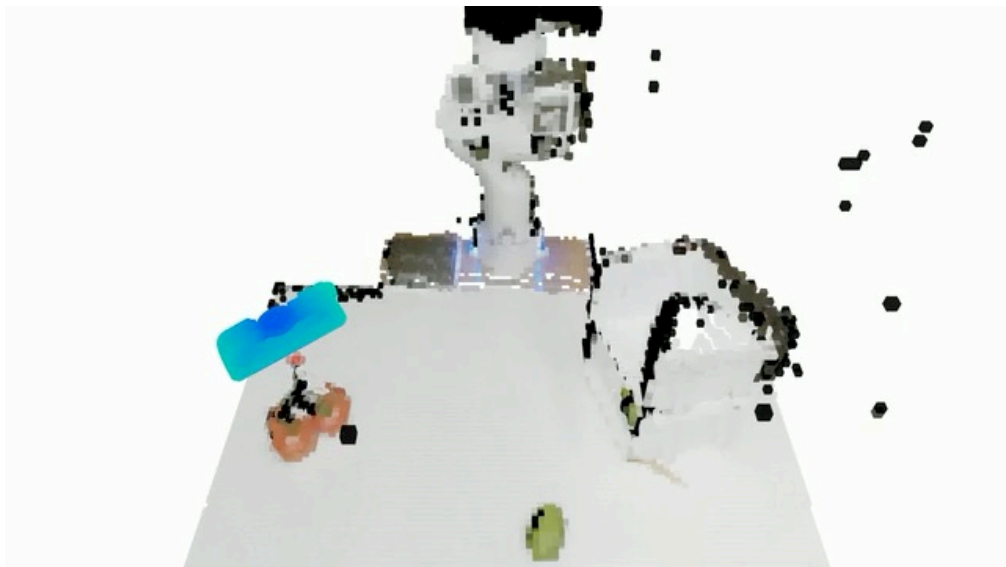
Perceiver-Actor

Multi-task 6-DoF manipulation agent

End-to-end few-shot imitation learning

Input: RGB-D Voxels & Language Goal

Output: Discretized 6-DoF action + open/close



“put the tomatoes in the top bin”

Perceiver-Actor

Multi-task 6-DoF manipulation agent

End-to-end few-shot imitation learning

Input: RGB-D Voxels & Language Goal

Output: Discretized 6-DoF action + c

obs

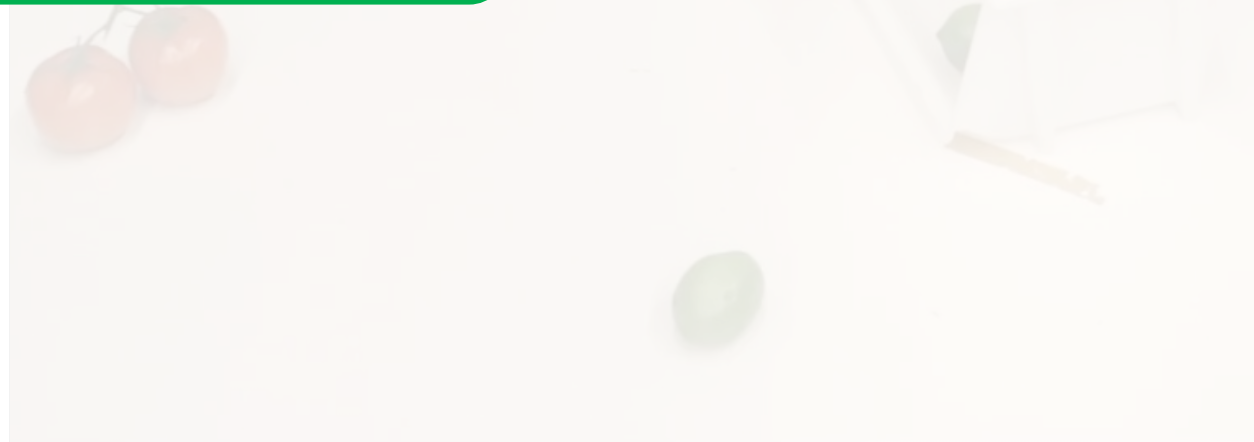
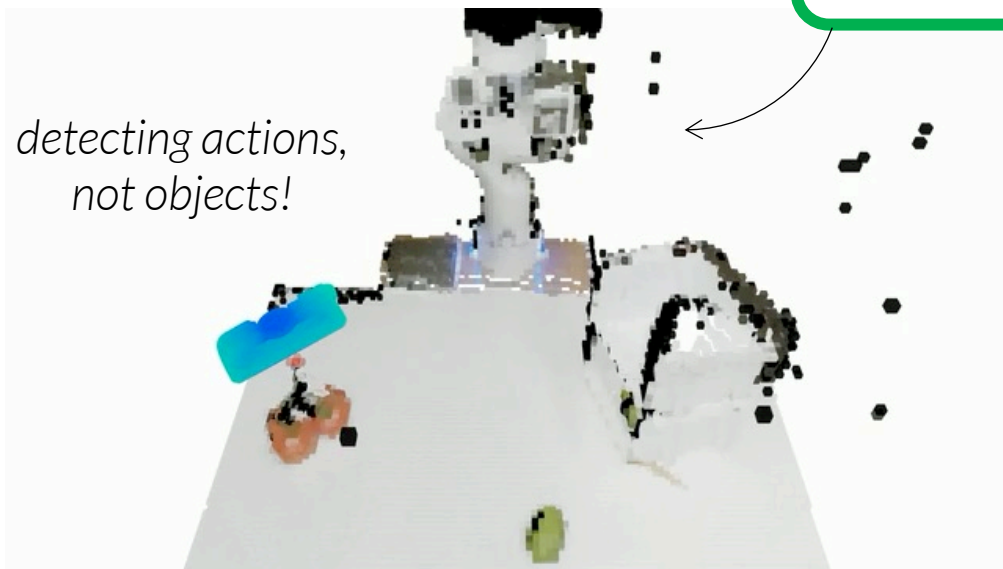


action



Obs Space = Action Space

*detecting actions,
not objects!*



“put the tomatoes in the top bin”

Perceiver-Actor

Multi-task 6-DoF manipulation agent

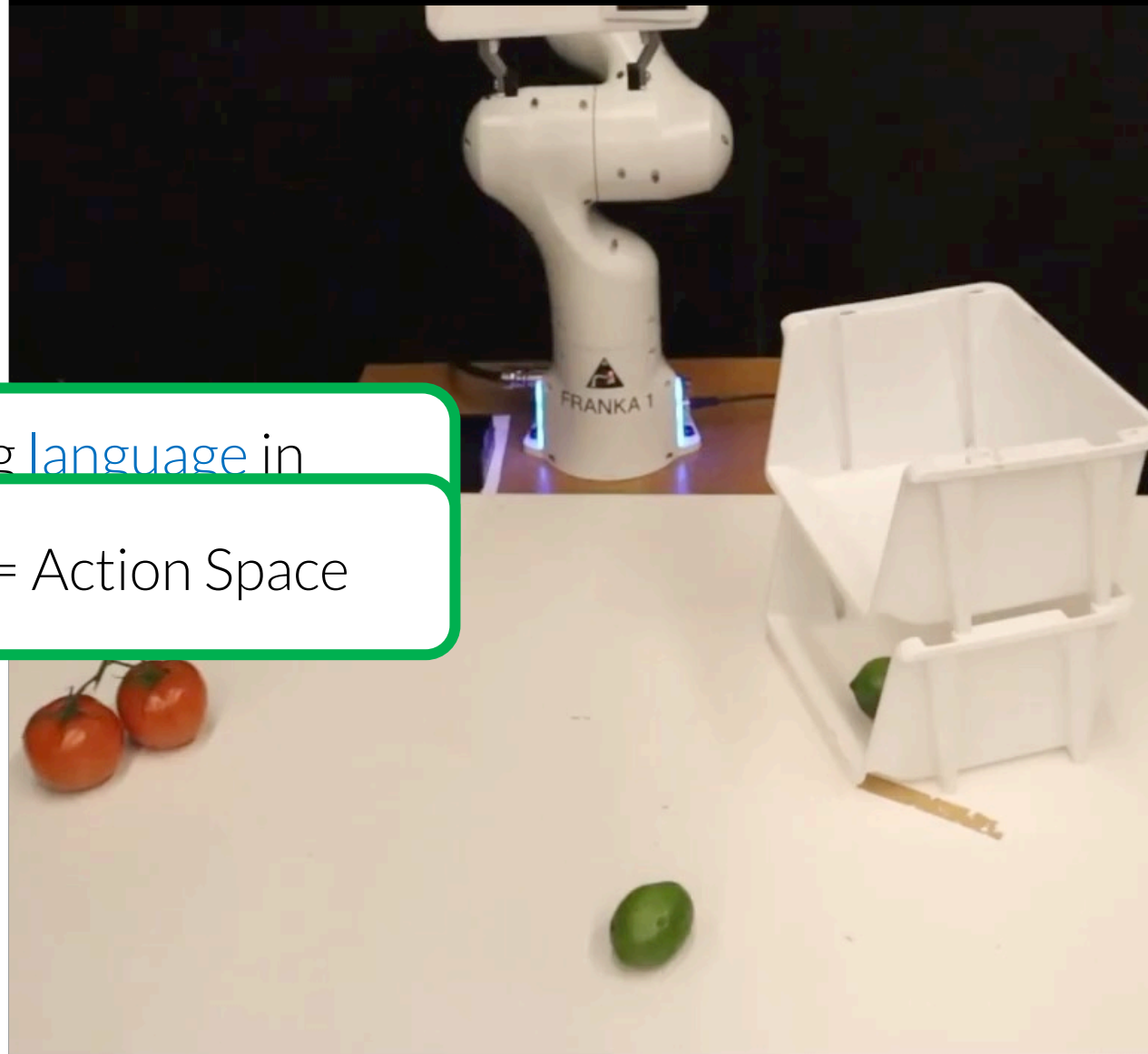
End-to-end few-shot imitation learning

Input: RGB-D Voxels & Language Goals

Output: Discretized 6-DoF action + c


Grounding **language** in

Obs Space = Action Space

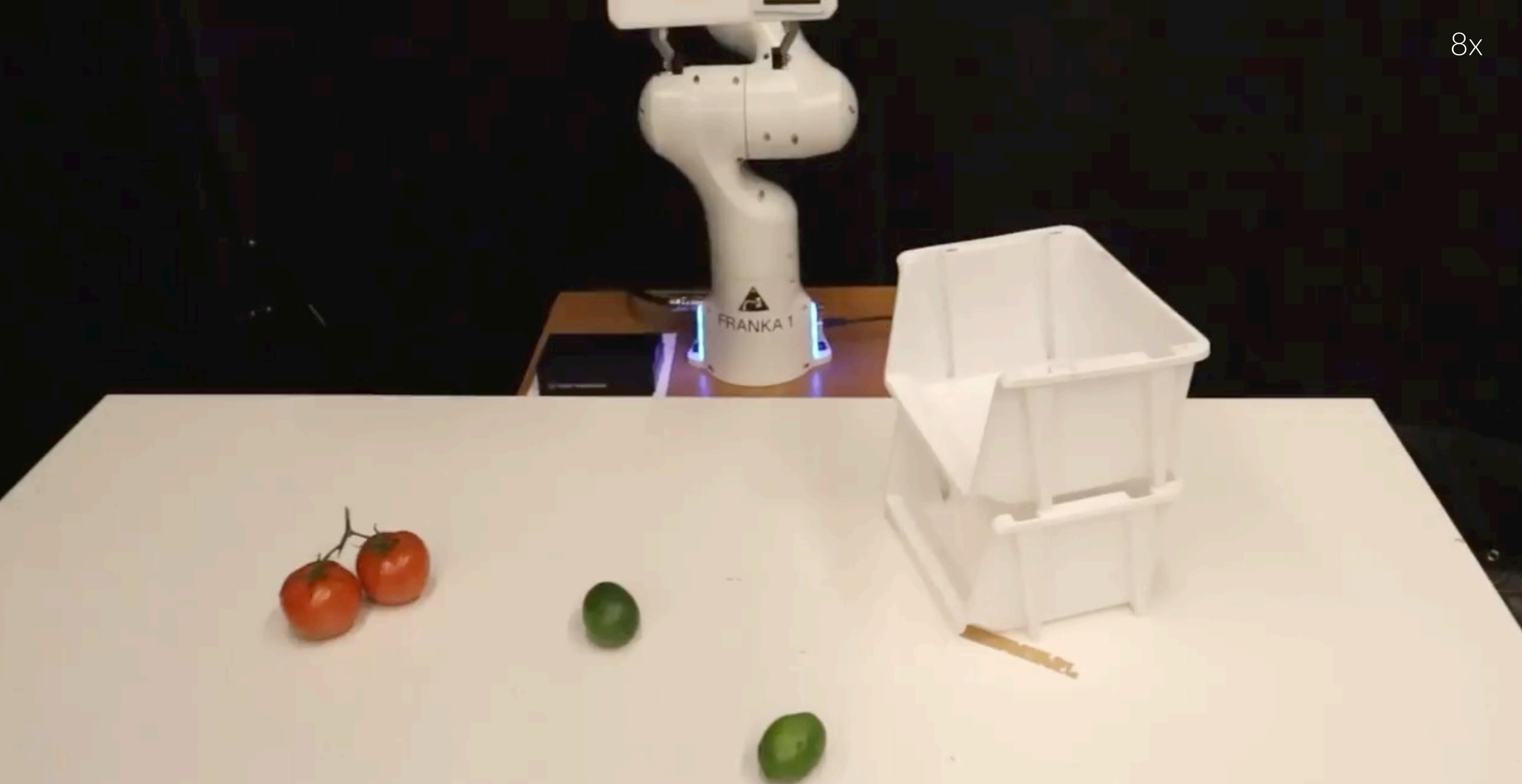


“put the tomatoes in the top bin”

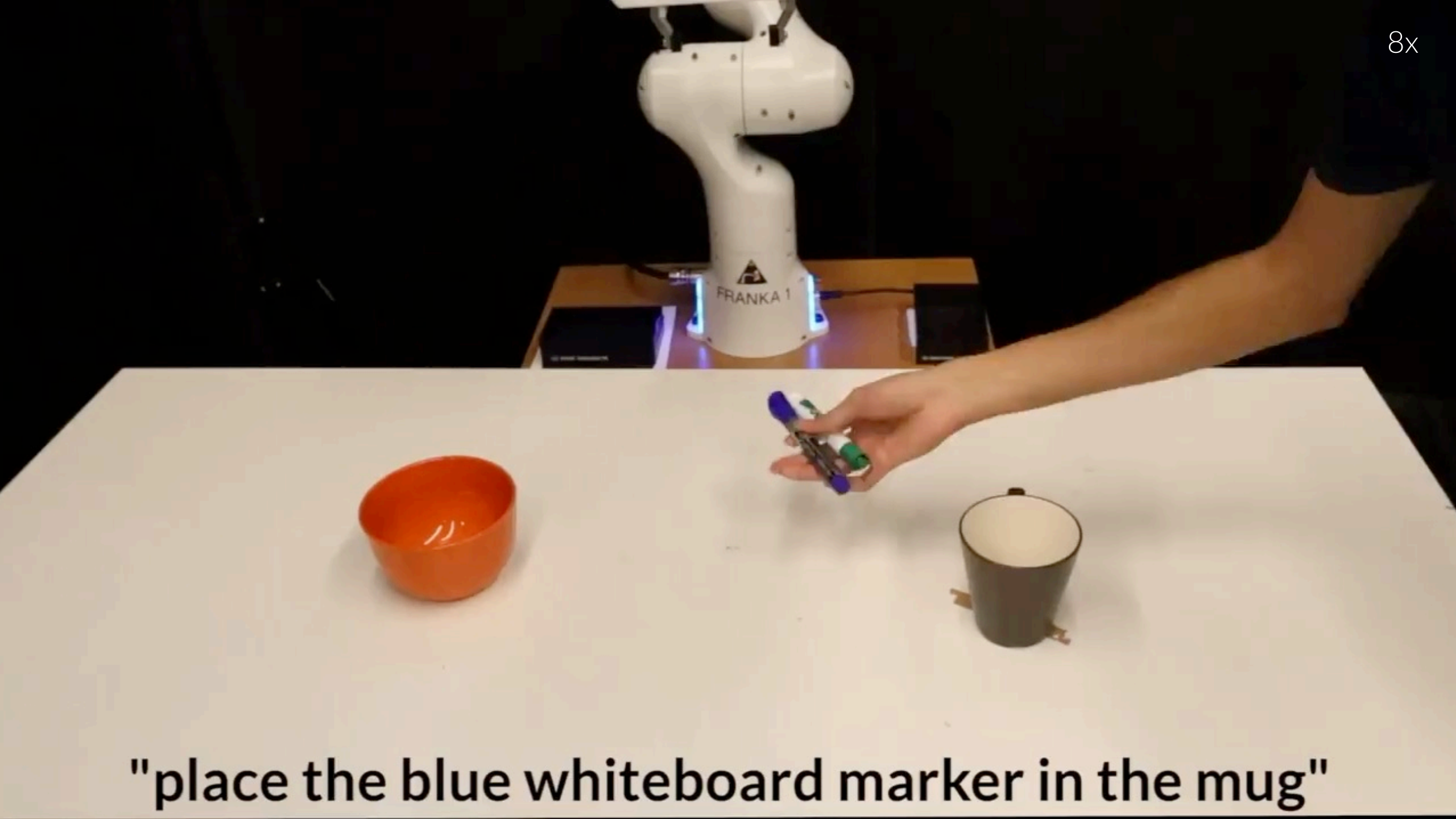
These results are from a one multi-task Transformer
trained *from scratch* with just **53 demos**



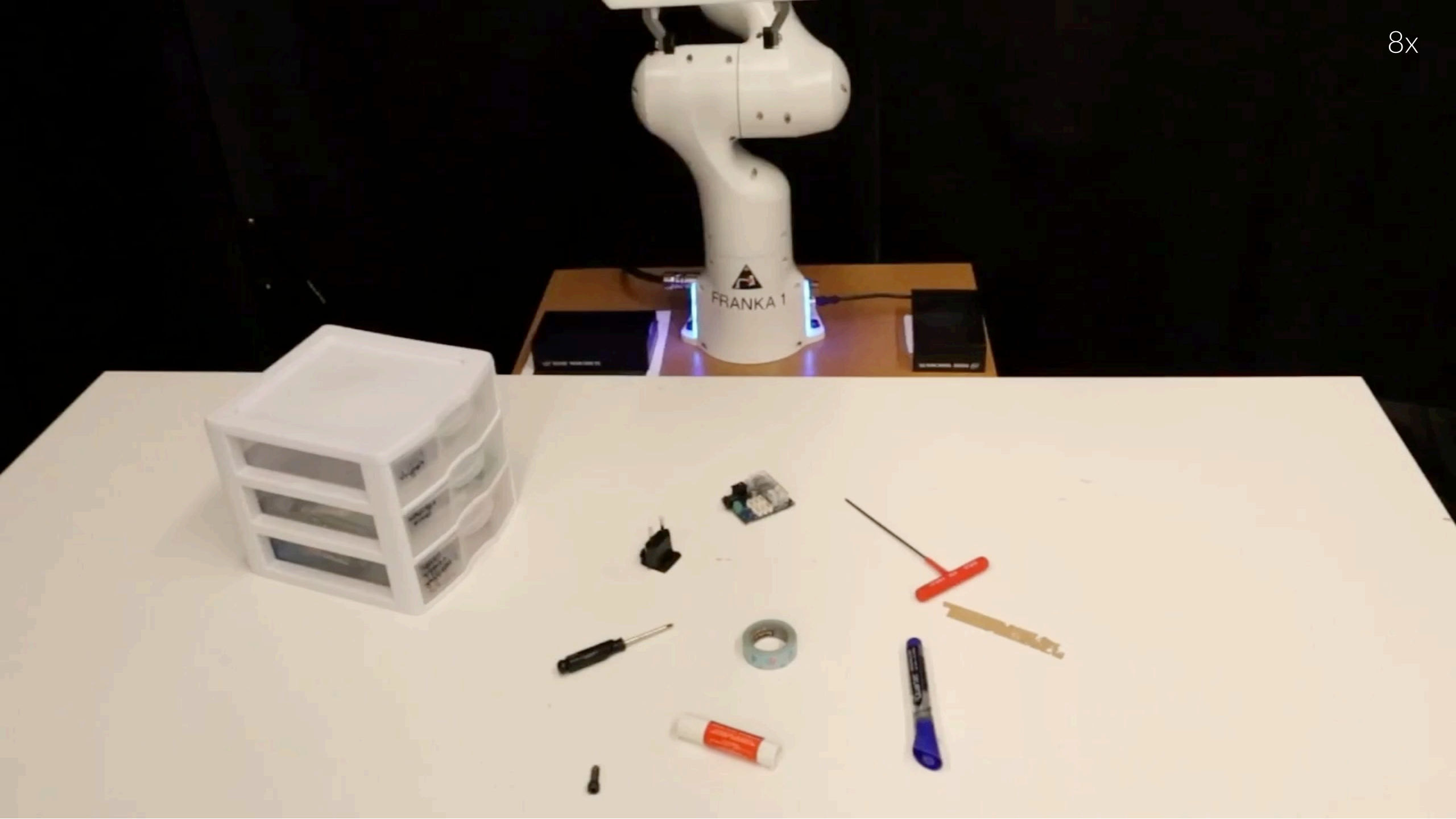
"press the hand san"

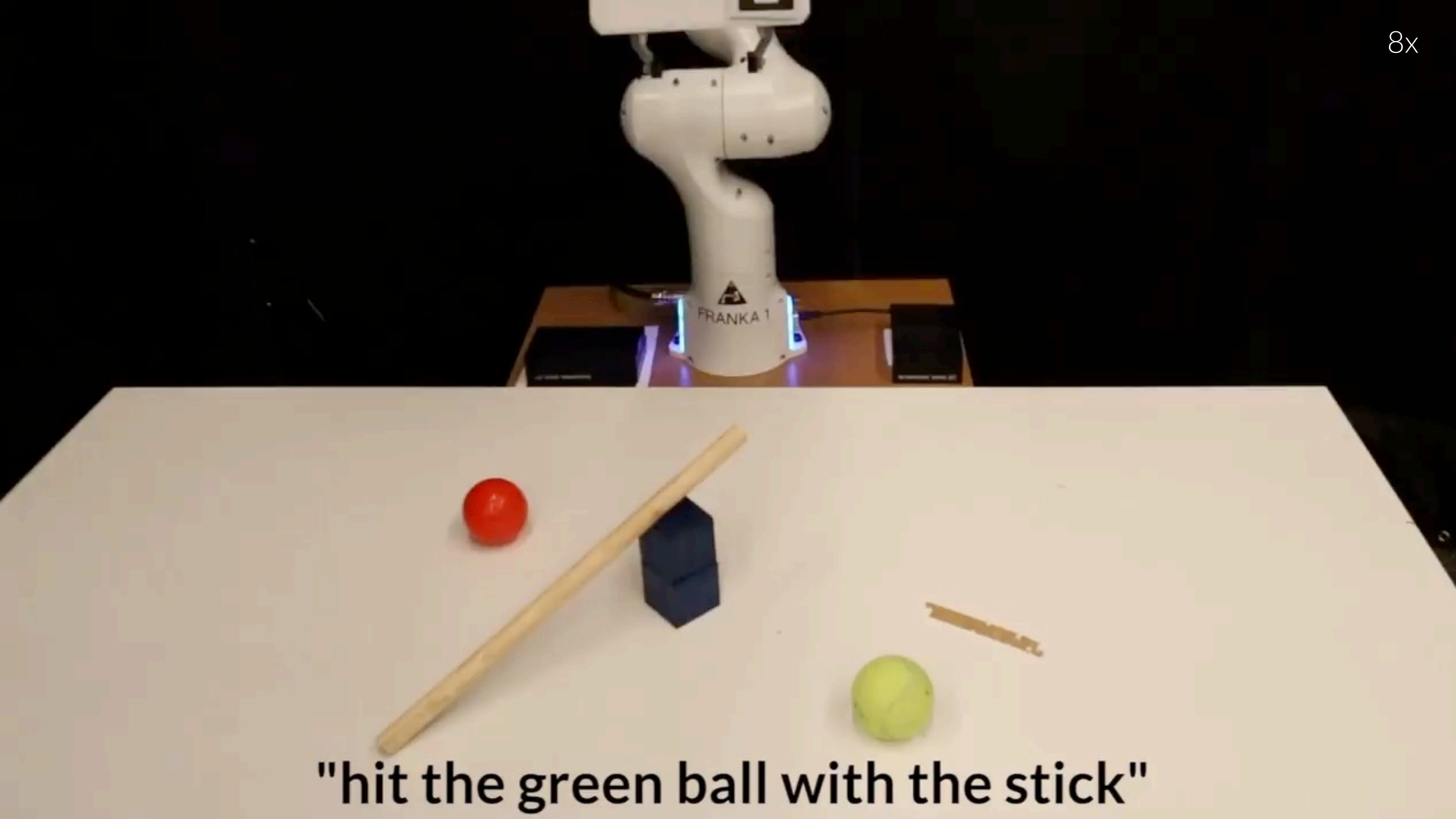


"put one lime in the bottom bin"

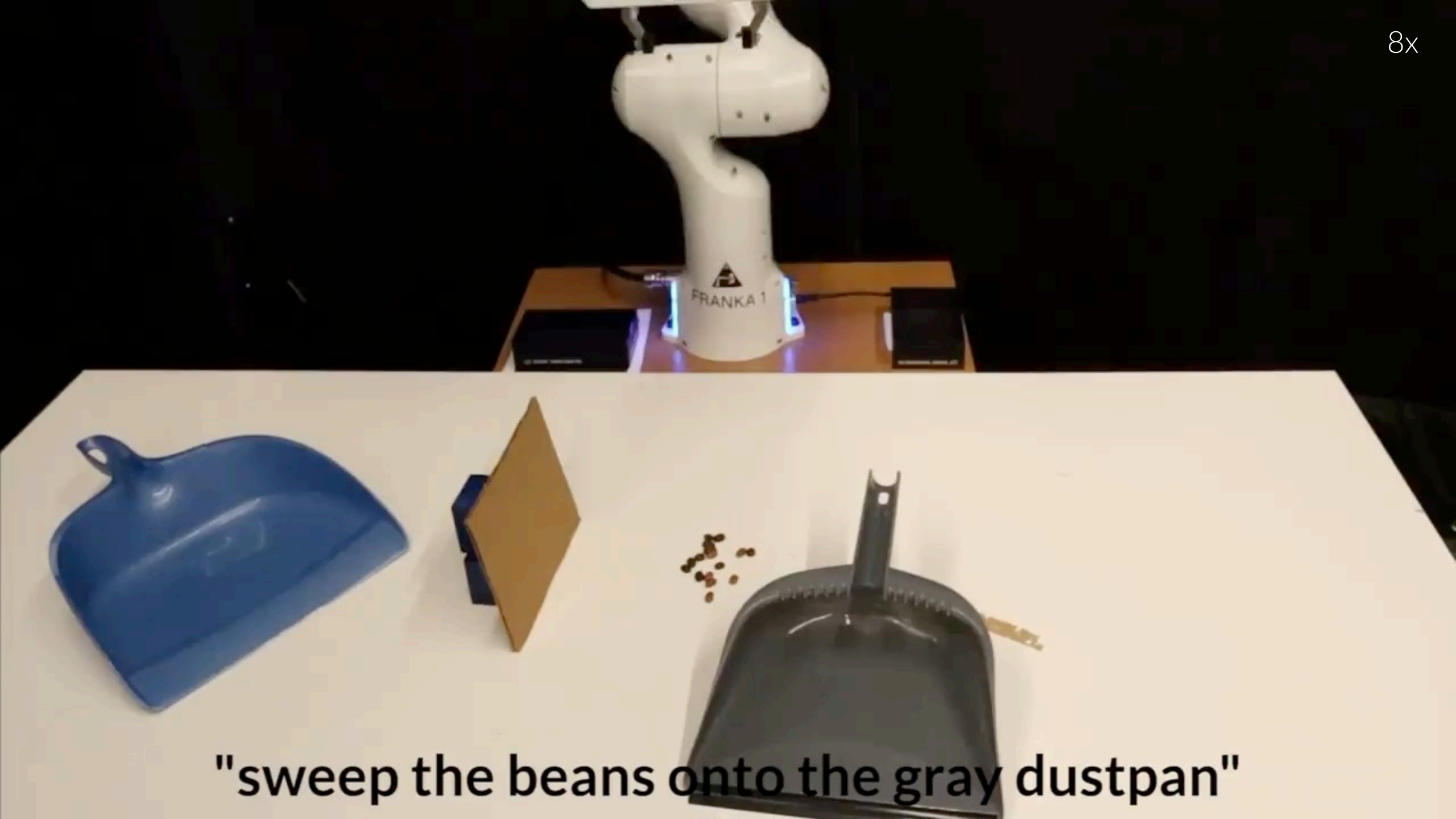


"place the blue whiteboard marker in the mug"



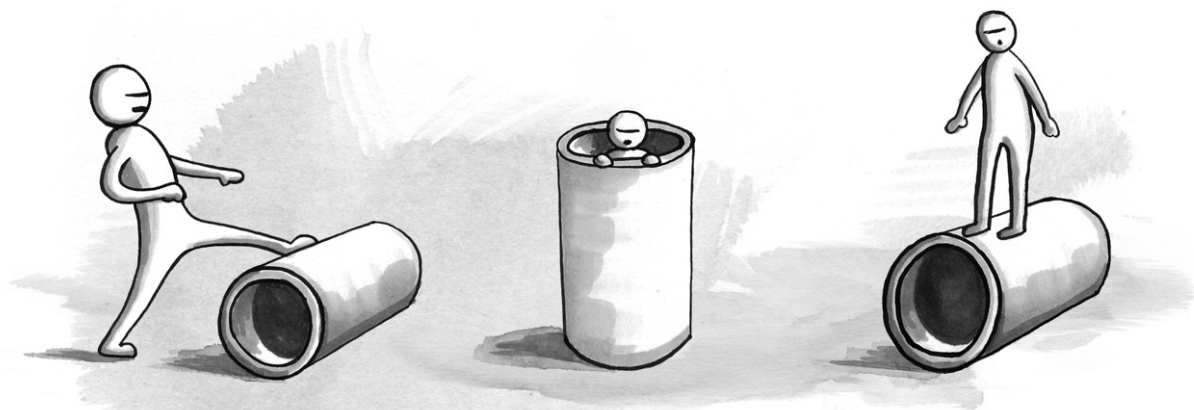


"hit the green ball with the stick"

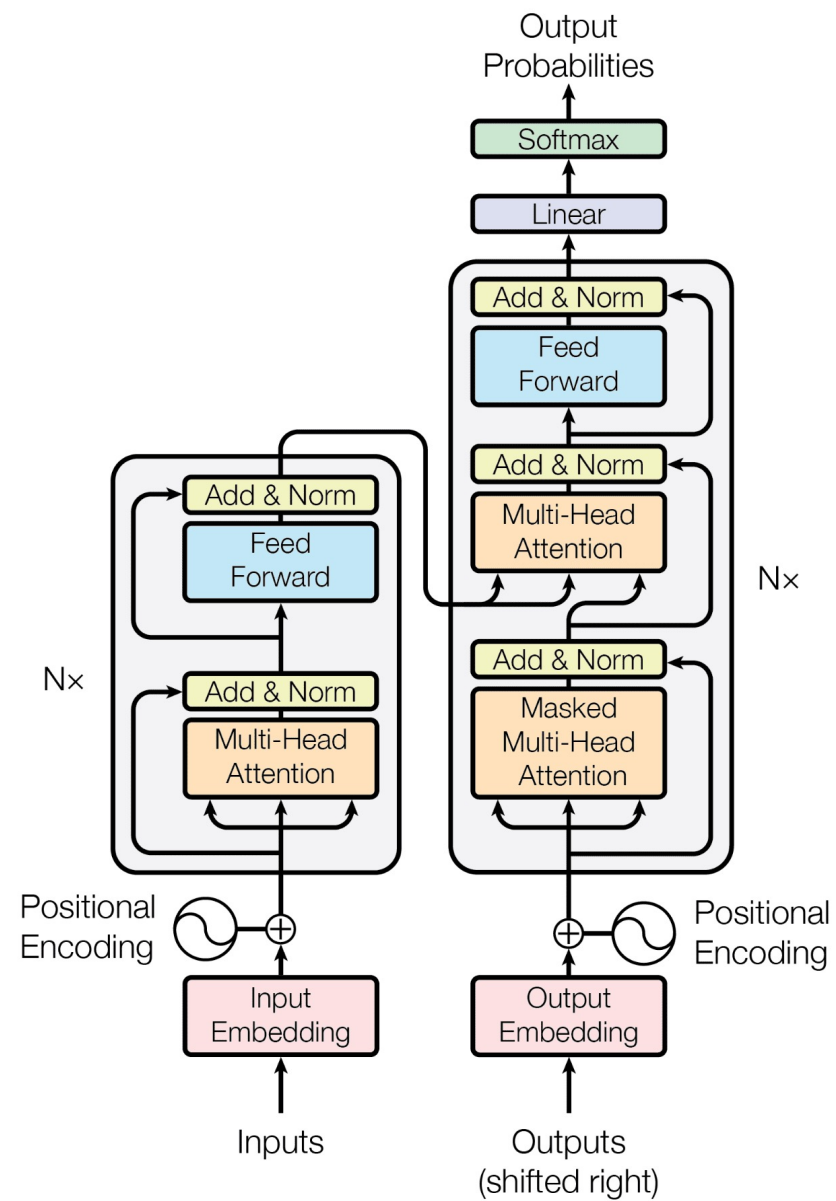


"sweep the beans onto the gray dustpan"

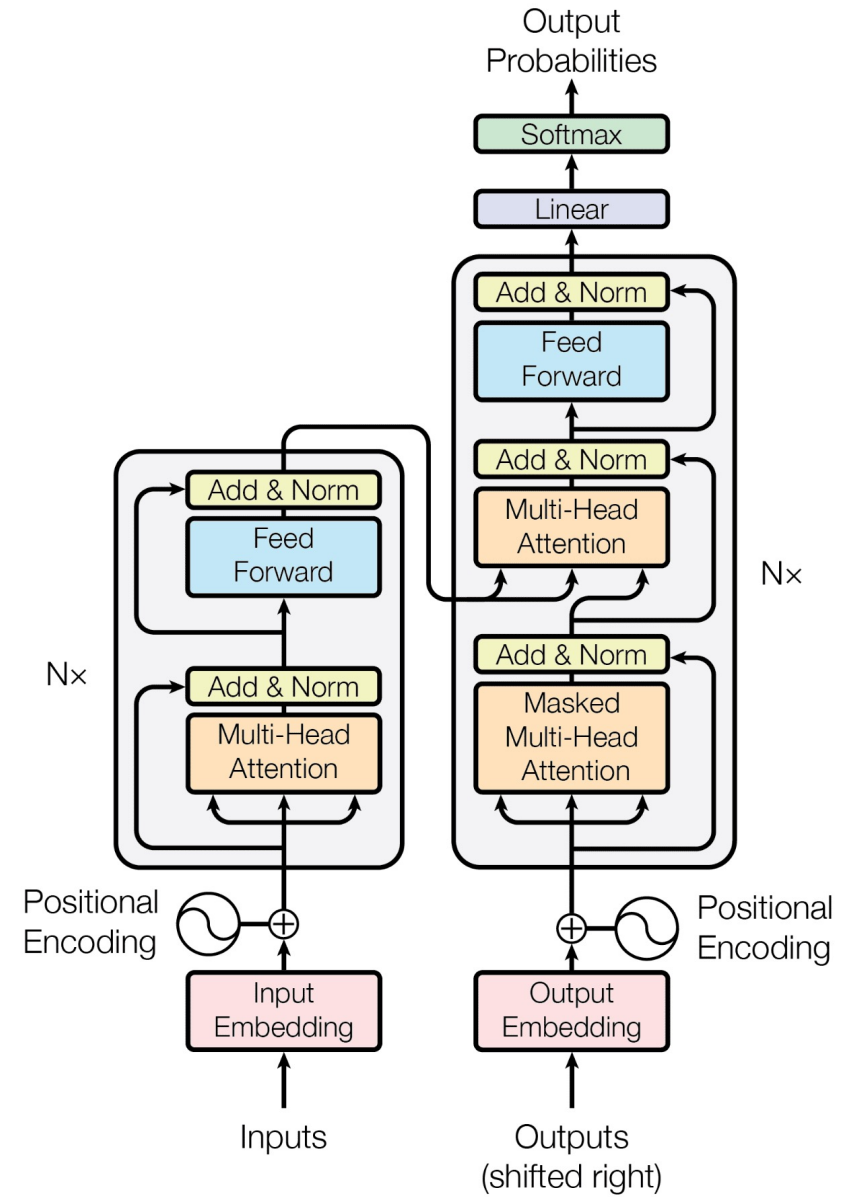
How does it work?



Credit: <https://uxdesign.cc/>



 **Transformer**
Vaswani et al, 2017



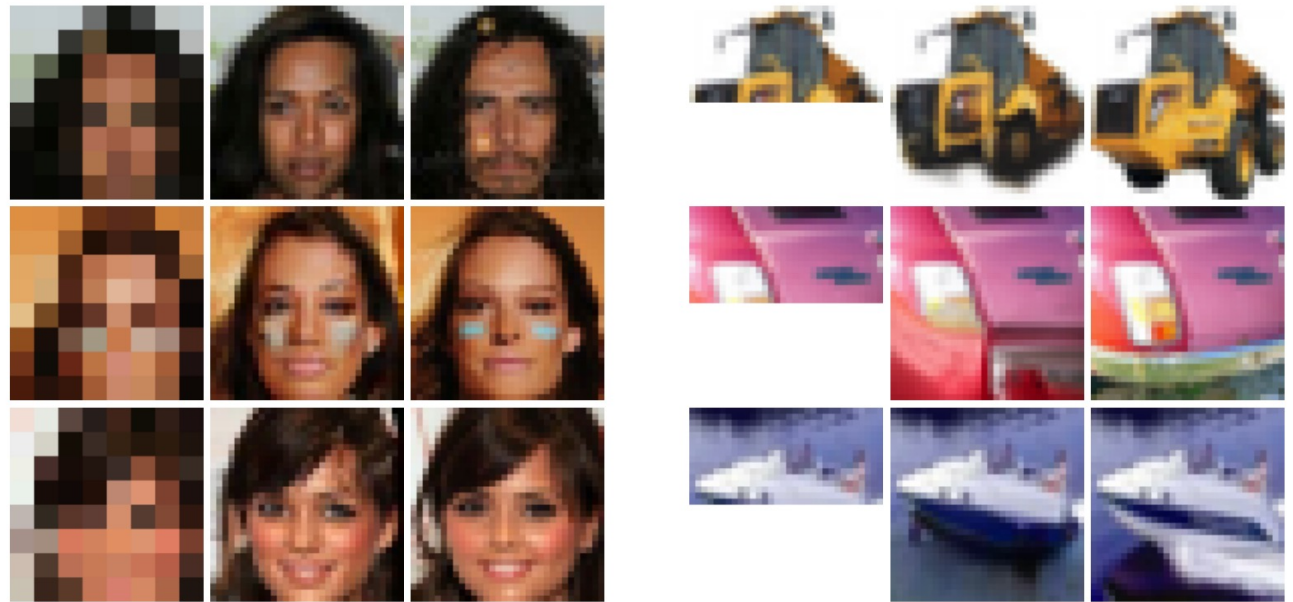


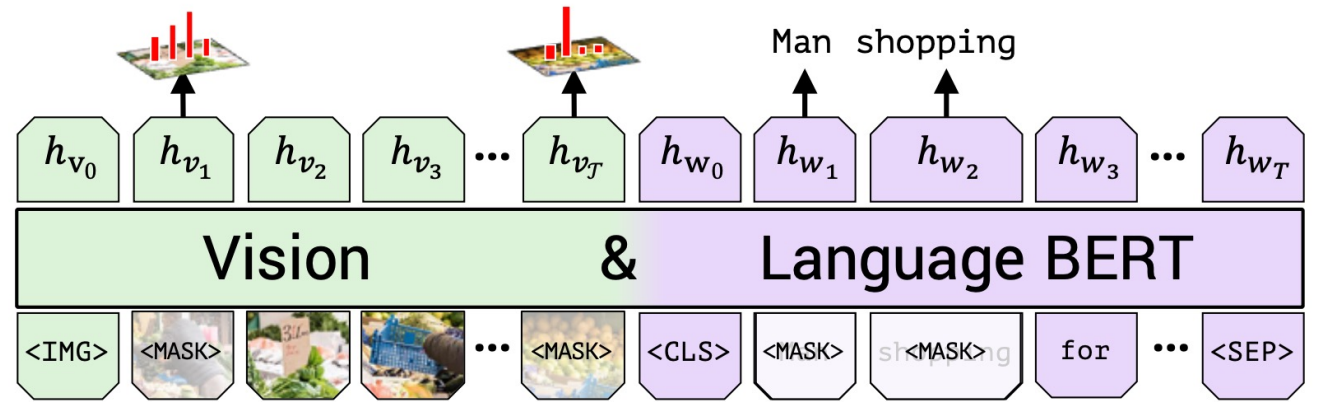
Image Transformer
Parmar et al, 2018



Transformer
Vaswani et al, 2017

Problem:

■ vs. word



Problem:
Object Detections

Vision-Lang BERTs

Lu et al, 2019
Tan et al, 2019

Image Transformer
Parmar et al, 2018

Transformer
Vaswani et al, 2017

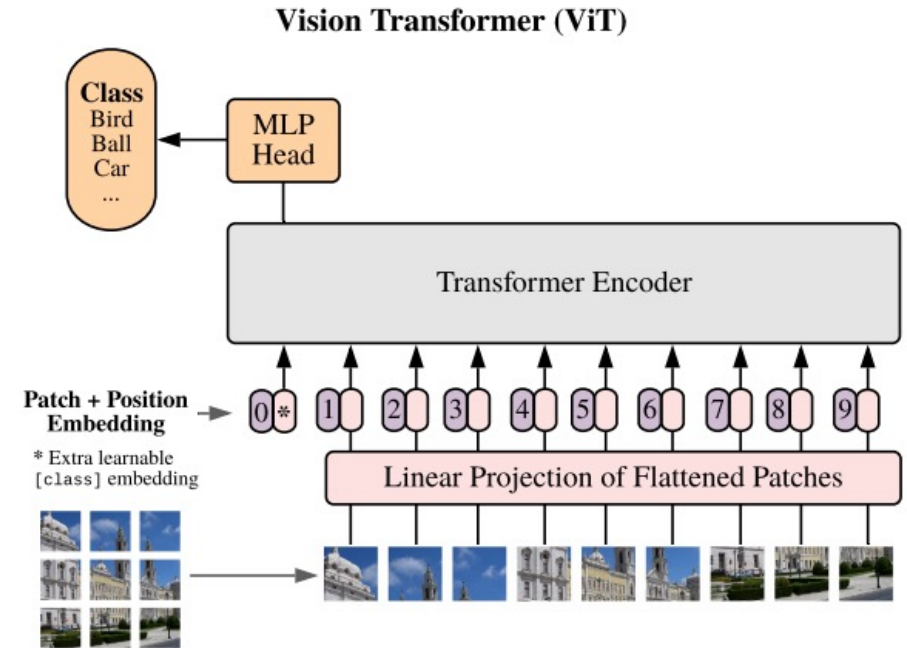
ViT
Dosovitskiy et al, 2020



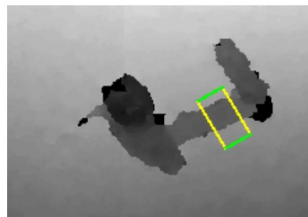
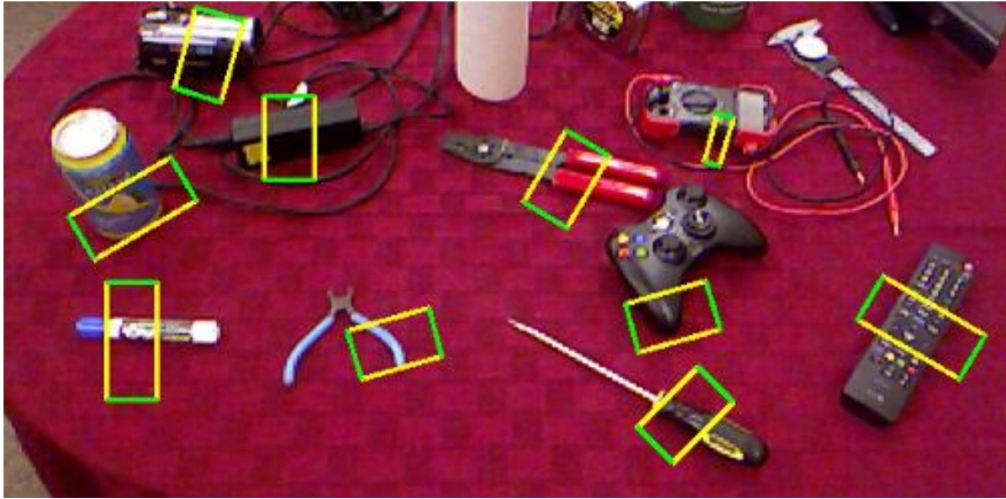
Vision-Lang BERTs
Lu et al, 2019
Tan et al, 2019

Image Transformer
Parmar et al, 2018

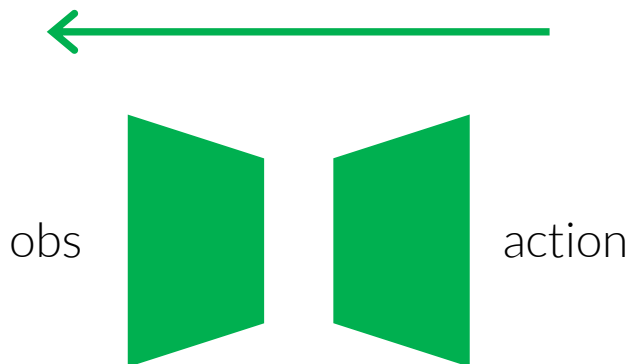
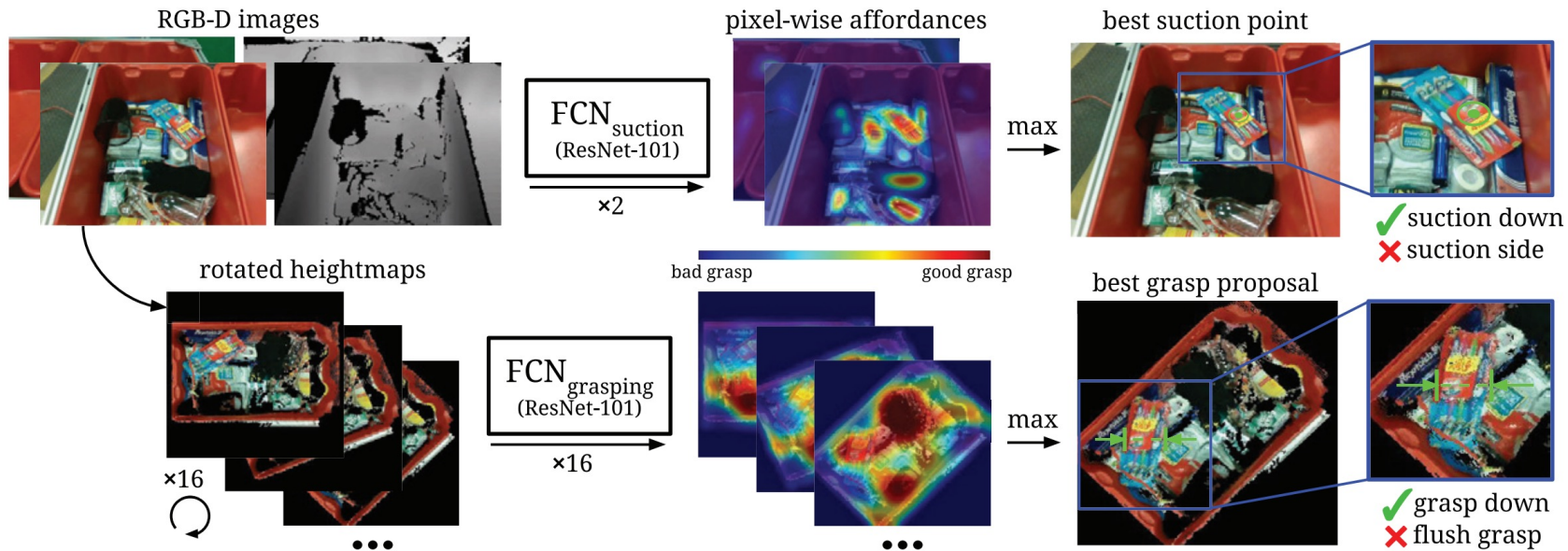
Transformer
Vaswani et al, 2017



Solution:
2D patches!



Deep Grasping
Lenz et al, 2014



Visual Affordances

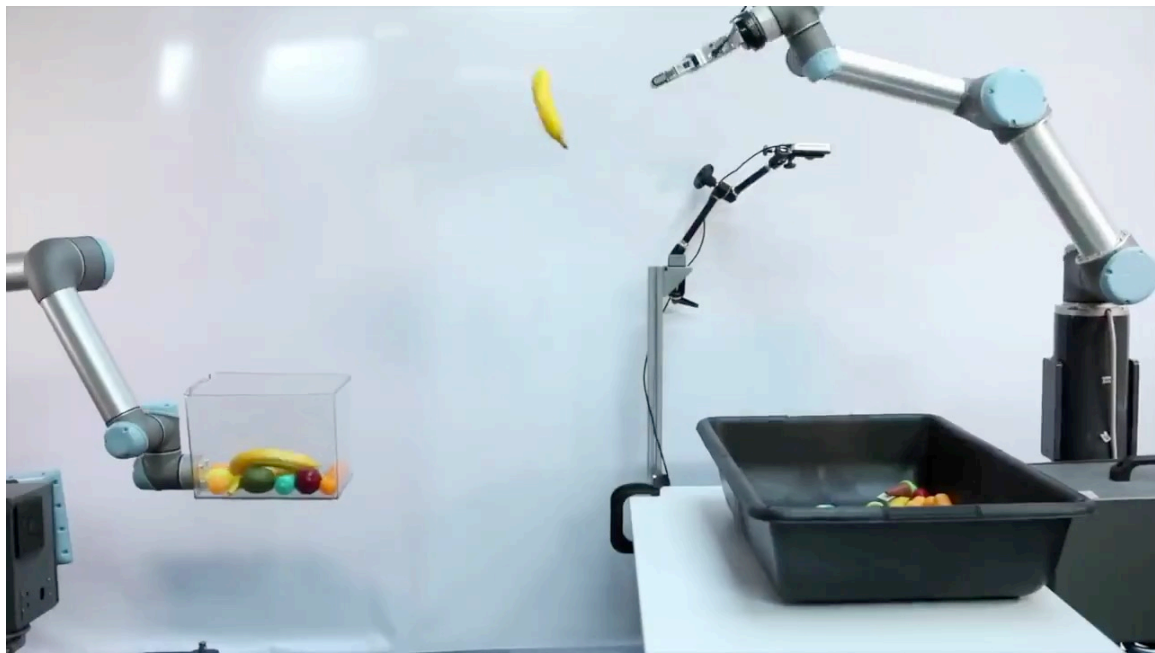
Zeng et al, 2017

Zeng et al, 2019



Deep Grasping

Lenz et al, 2014



Visual Affordances

Zeng et al, 2017

Zeng et al, 2019



Deep Grasping

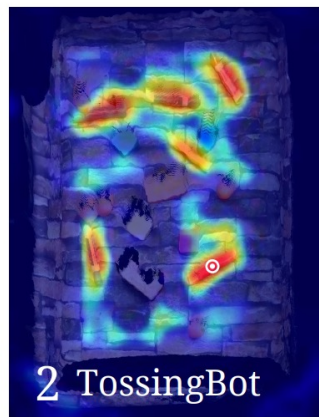
Lenz et al, 2014



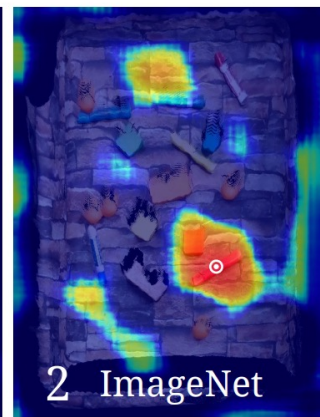
(a)



(b)



(e)



(f)

Problem:

Missing Natural Interface
for Goal-Conditioning

Visual Affordances

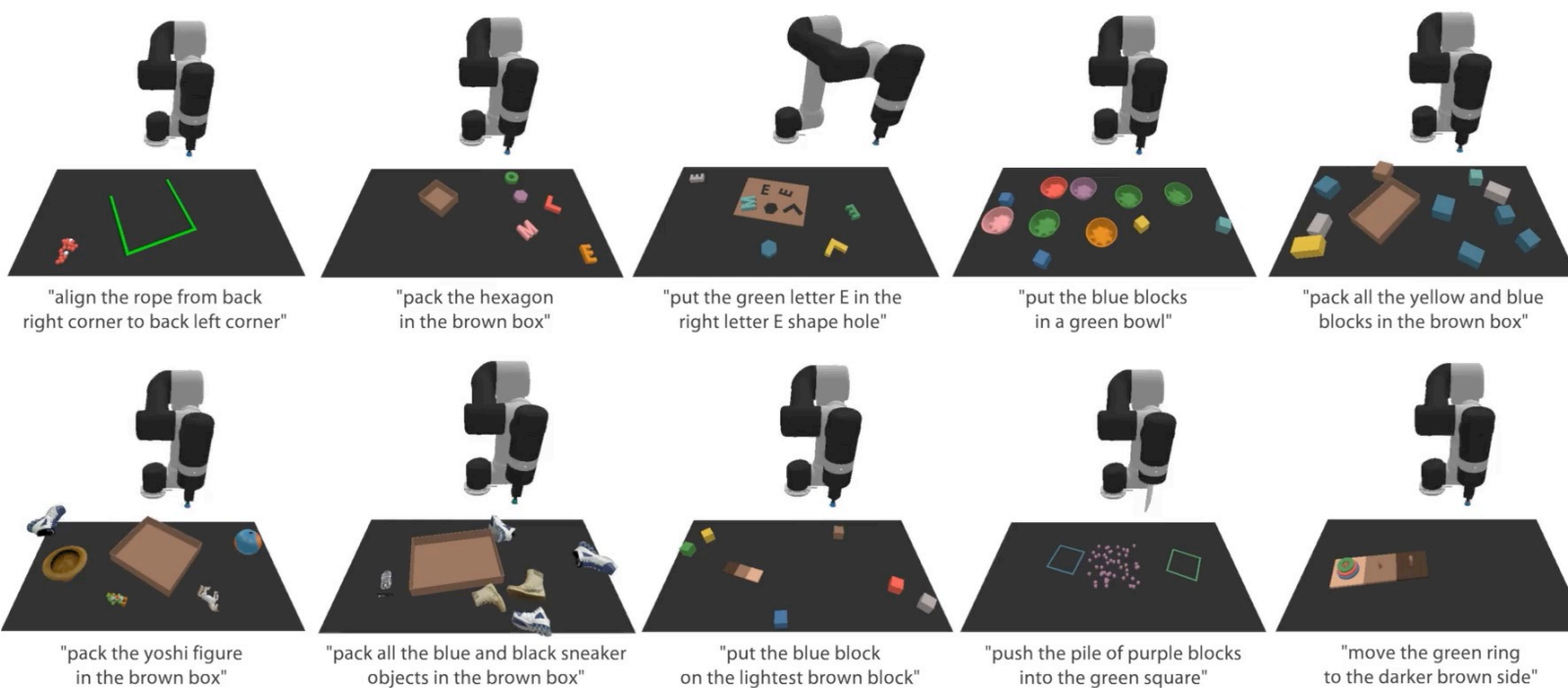
Zeng et al, 2017

Zeng et al, 2019



Deep Grasping

Lenz et al, 2014



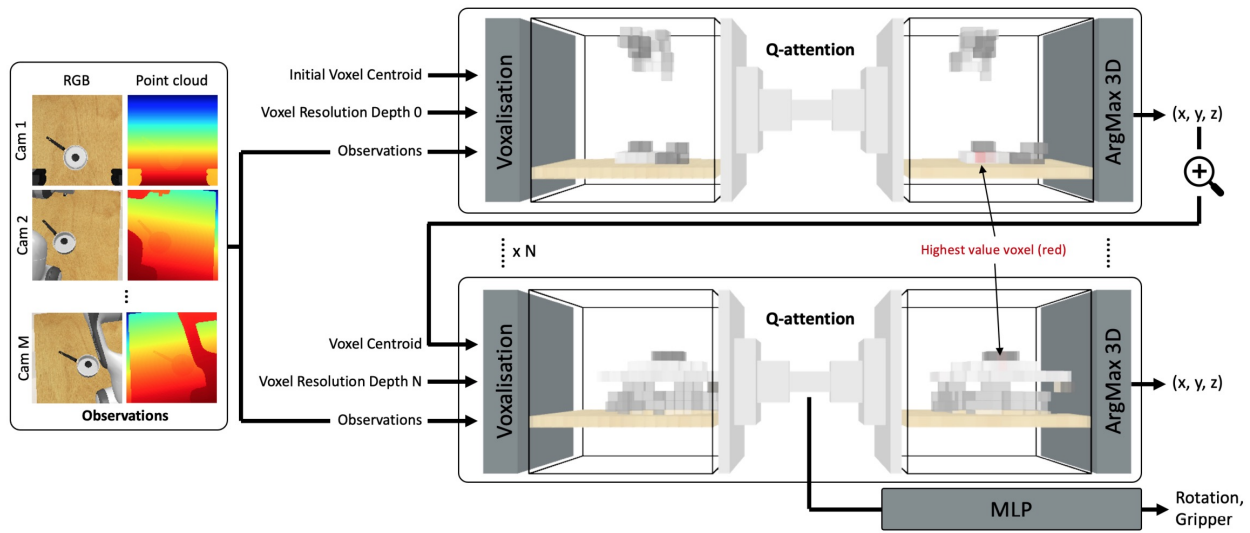
Problem:
2D Obs & Action Space

Visual Affordances
Zeng et al, 2017
Zeng et al, 2019



CLIPort
Shridhar et al, 2021

Deep Grasping
Lenz et al, 2014



Problem:
High Dimensional Input

C2FARM
James et al, 2021

Visual Affordances
Zeng et al, 2017
Zeng et al, 2019

CLIPort
Shridhar et al, 2021

Deep Grasping
Lenz et al, 2014



ViT
Dosovitskiy et al, 2020



Vision-Lang BERTs
Lu et al, 2019
Tan et al, 2019

Image Transformer
Parmar et al, 2018



Transformer
Vaswani et al, 2017



C2FARM
James et al, 2021



CLIPort
Shridhar et al, 2021



Visual Affordances
Zeng et al, 2017
Zeng et al, 2019



Deep Grasping
Lenz et al, 2014



ViT
Dosovitskiy et al, 2020



Vision-Lang BERTs
Lu et al, 2019
Tan et al, 2019



Image Transformer
Parmar et al, 2018



Transformer
Vaswani et al, 2017



What are the right
tokens for manipulation?

3D Voxel Patches

*Problems with
2D Static Monocular RGB*



hand-eye coordination
depth cues
camera perturbations
distractors
spatial data augmentation?

What are the right
tokens for manipulation?

3D Voxel Patches

How to deal with
high dimensional input?

Latent-Space Transformer

C2FARM
James et al, 2021

Visual Affordances
Zeng et al, 2017
Zeng et al, 2019

CLIPort
Shridhar et al, 2021

Deep Grasping
Lenz et al, 2014



ViT
Dosovitskiy et al, 2020

How to deal with
high dimensional input?

Image
Latent-Space Transformer
Parmar et al, 2018

Transformer
Vaswani et al, 2017

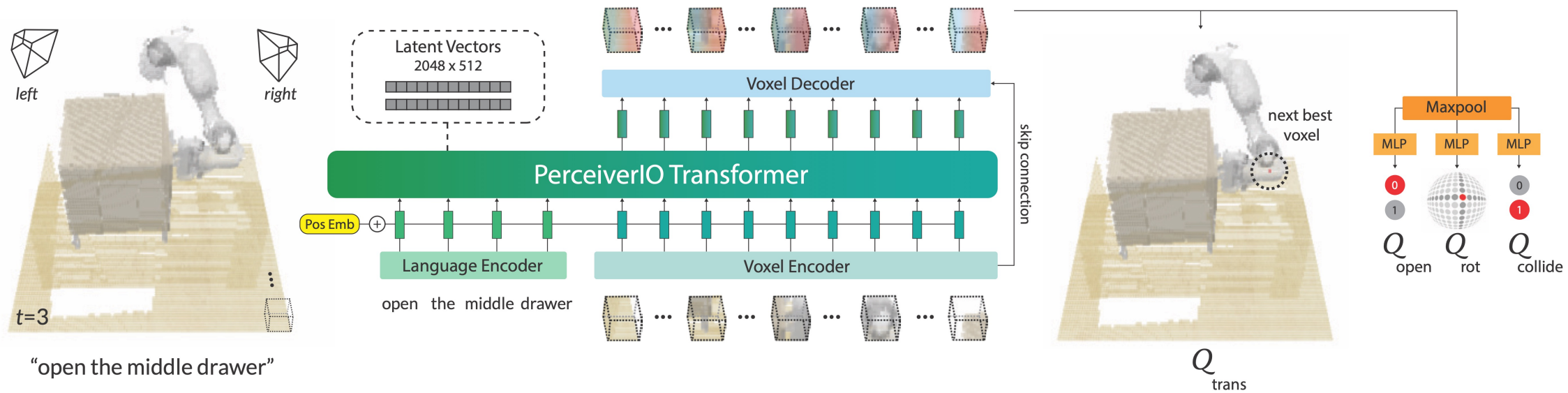
C2FARM
James et al, 2021

What are the right
tokens for manipulation?

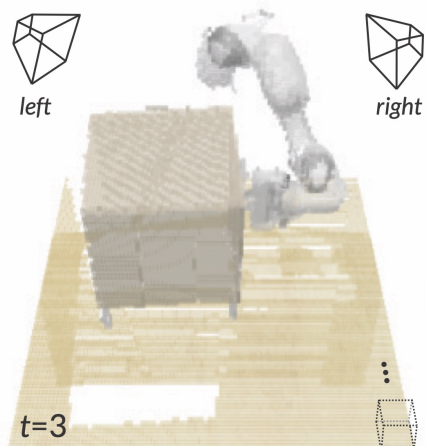
Visual Affordance
Zeng et al, 2017
3D Voxel Patches
Zeng et al, 2019

Deep Grasping
Lenz et al, 2014

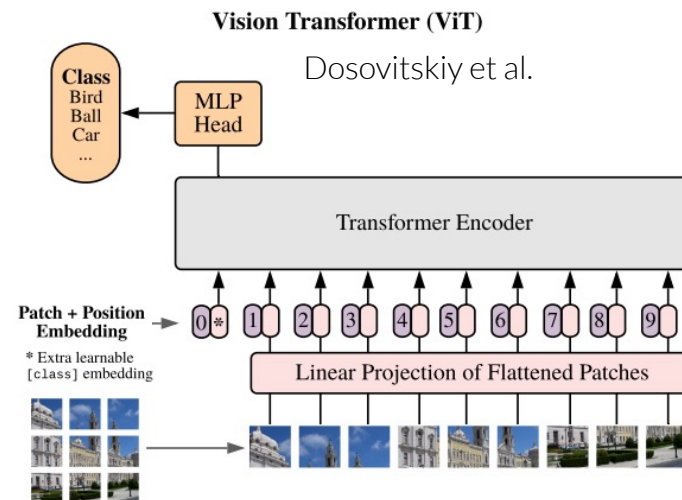
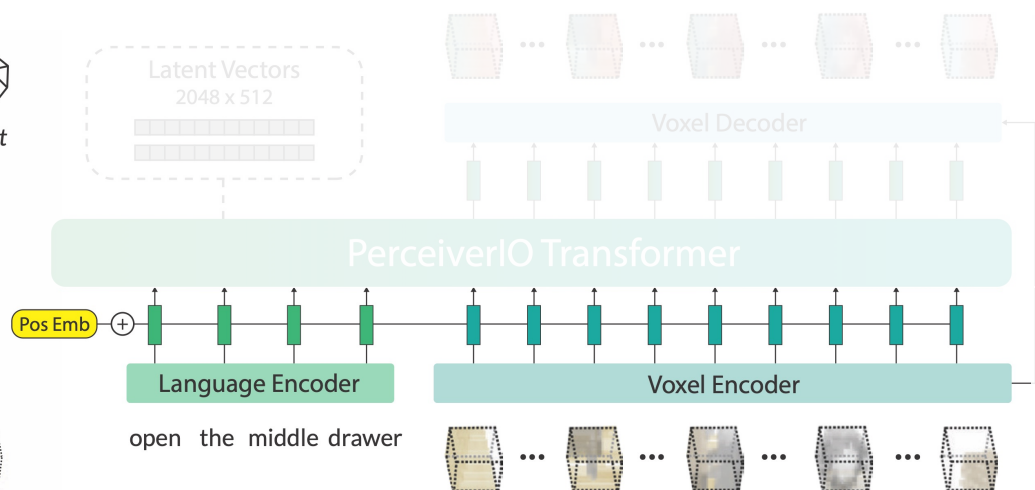
Perceiver-Actor



Perceiver-Actor



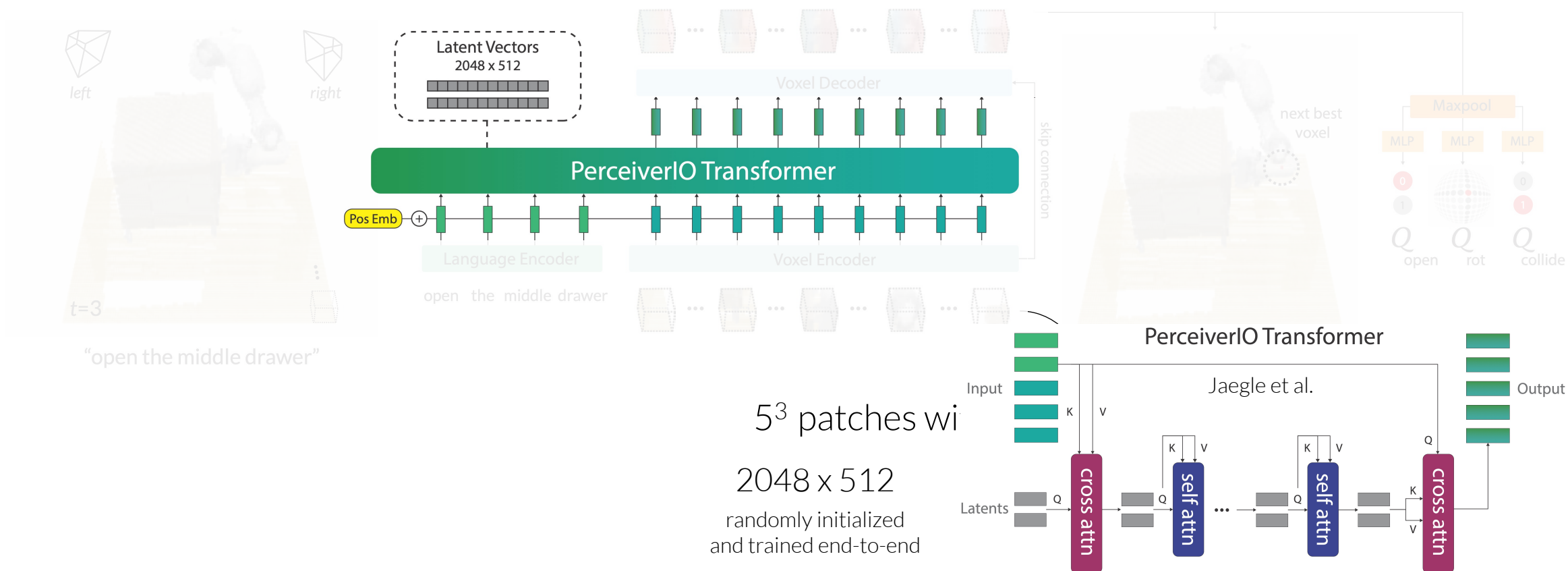
“open the middle drawer”



ViT : 2D image patches
PerAct : 3D voxel patches

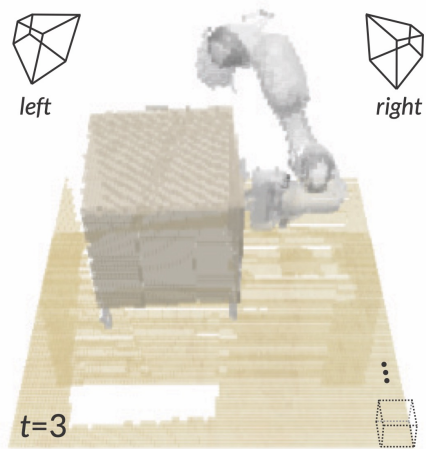


Perceiver-Actor

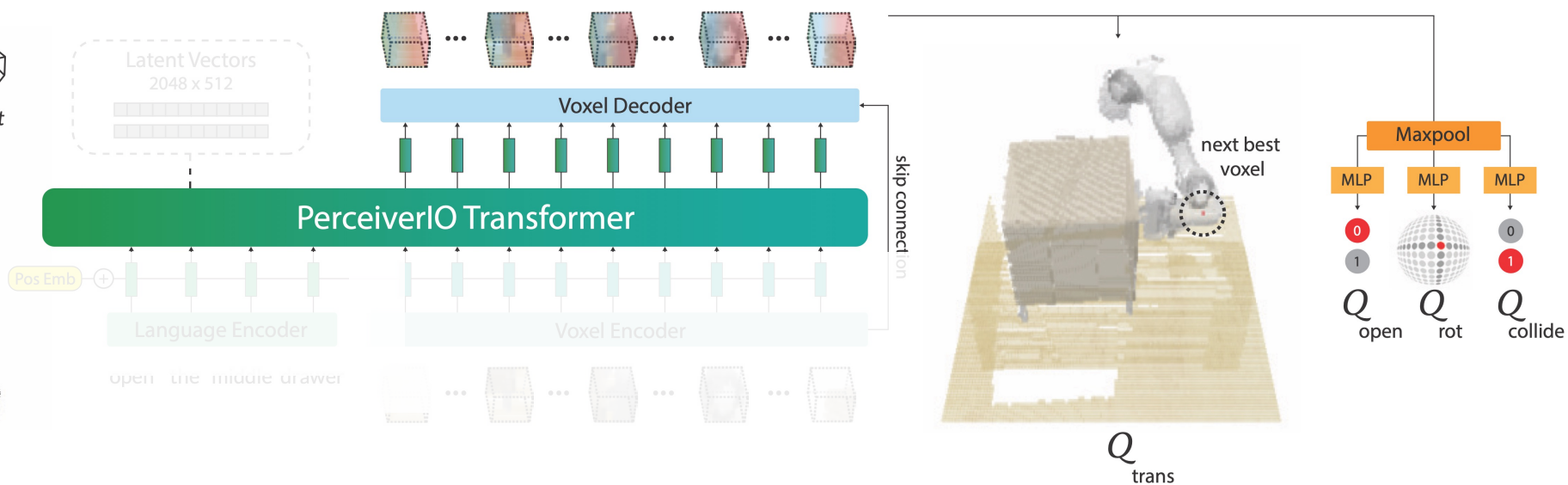


Perceiver-Actor

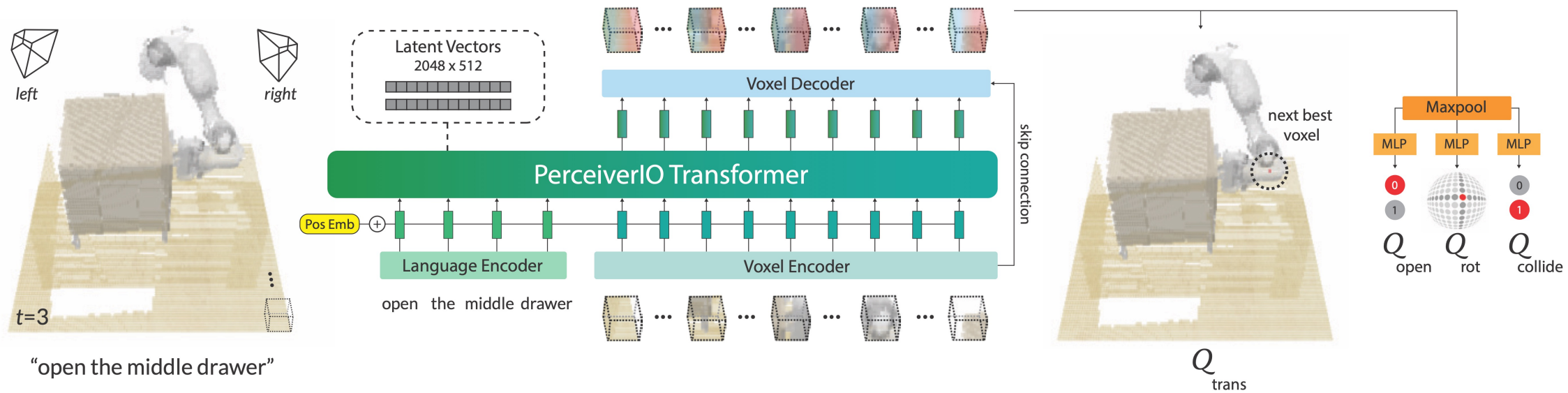
100 x 100 x 100 x 64 features



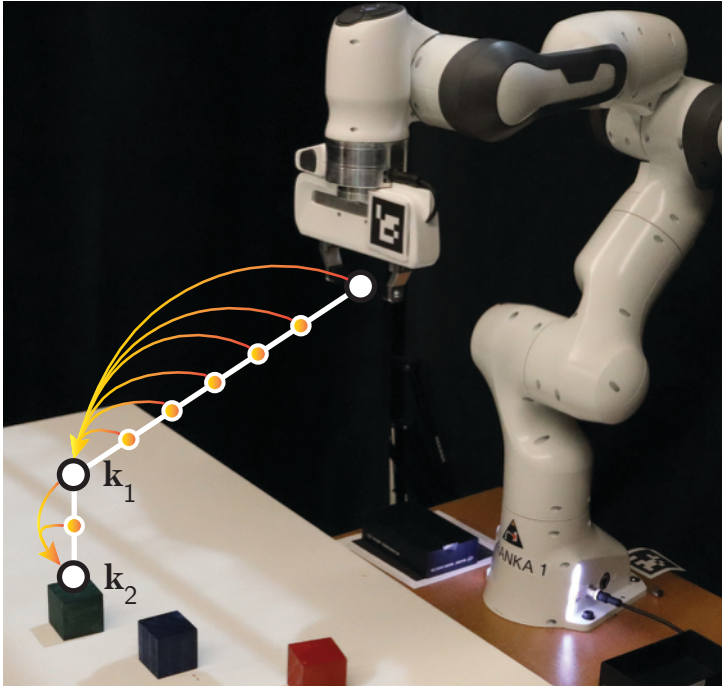
"open the middle drawer"



Perceiver-Actor



Dataset Setup



Heuristic for Keyframe Extraction:

- (1) Joint velocities are near zero &
- (2) Gripper open state has not changed

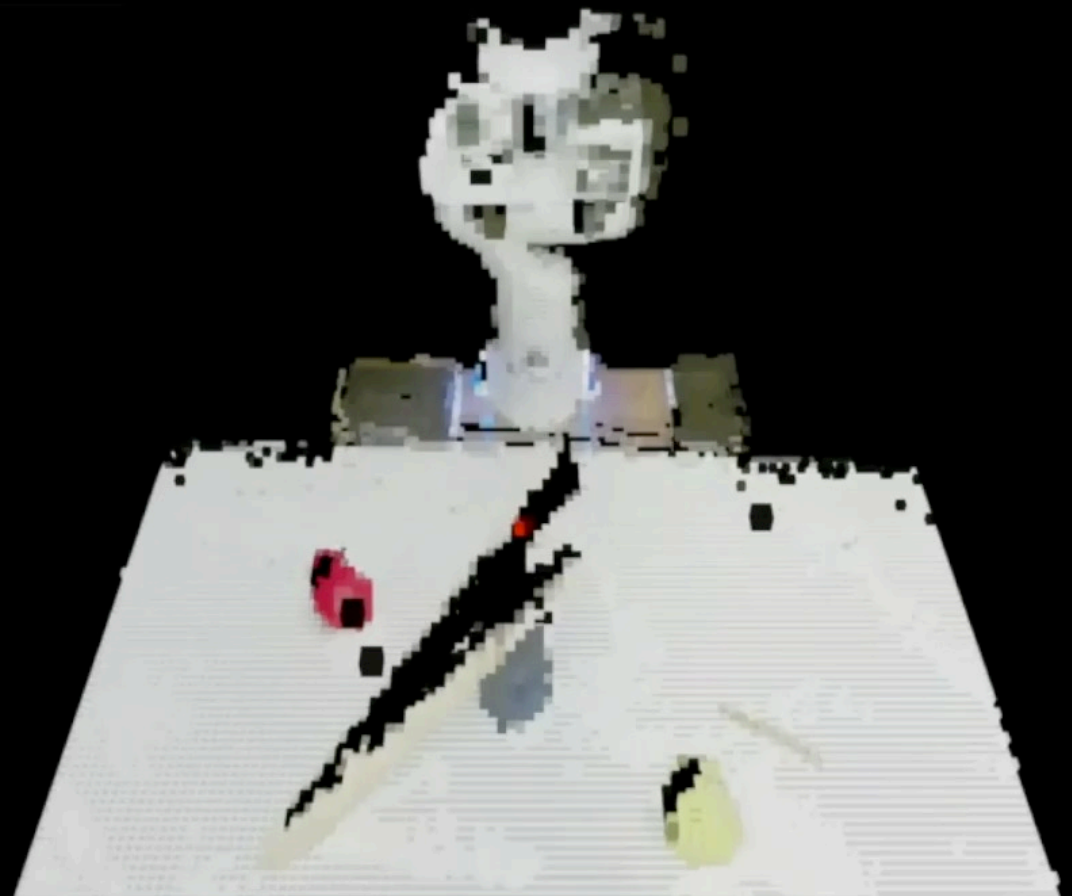
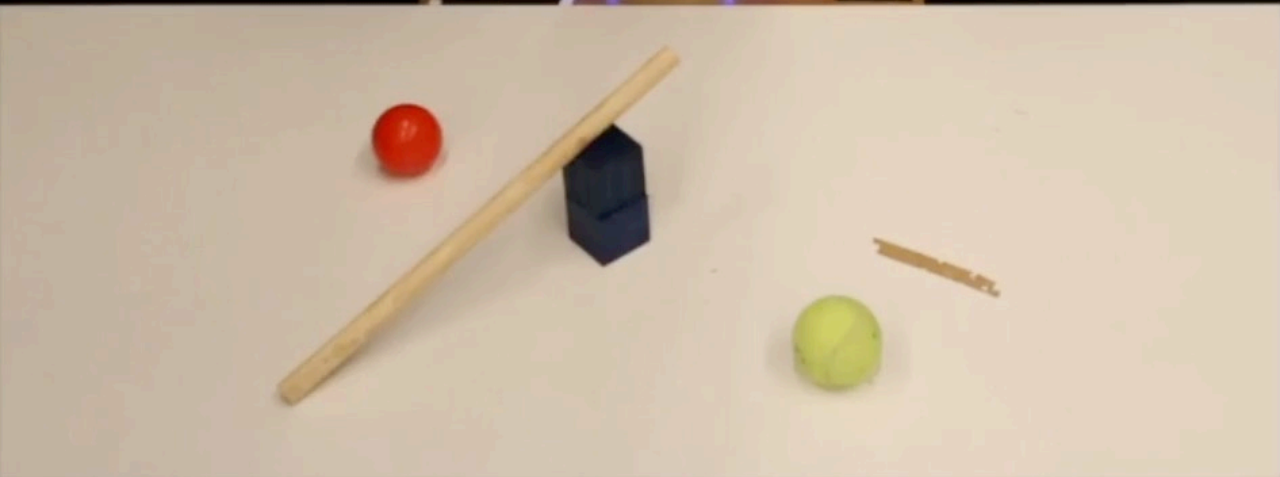
Predict the “next best keyframe action”
classification task

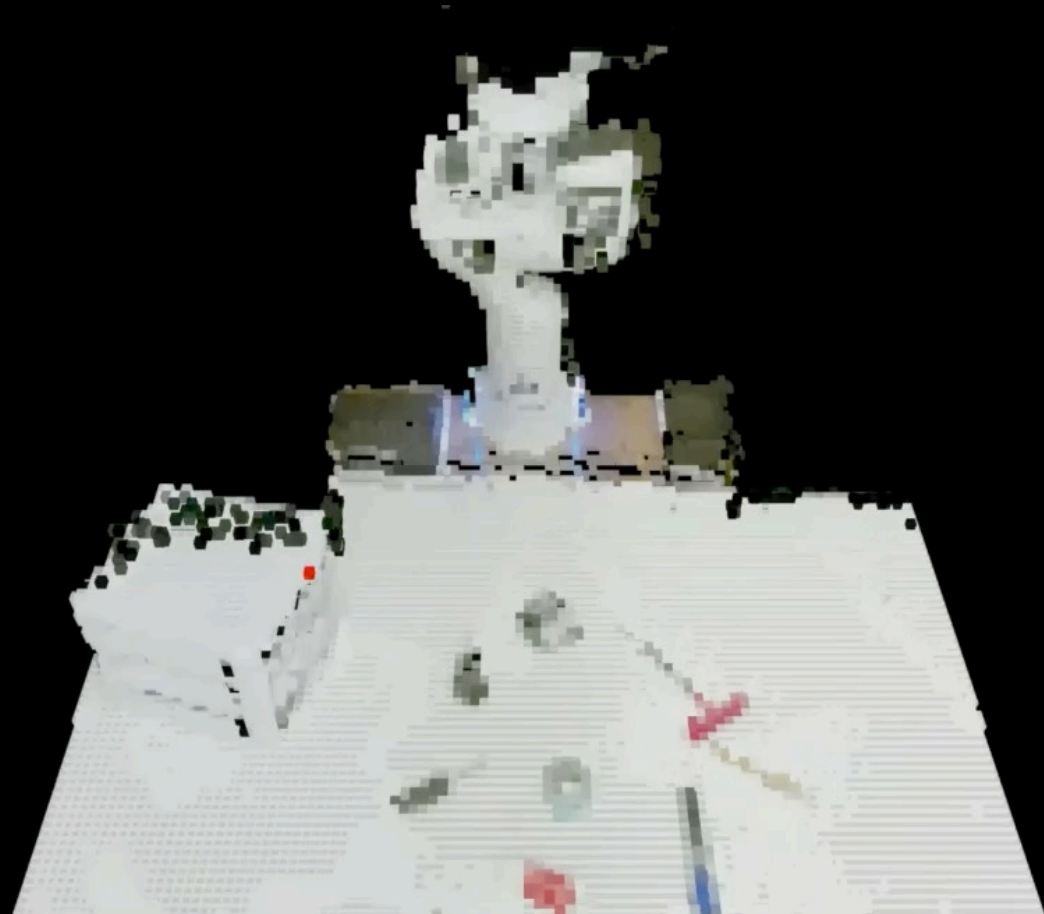
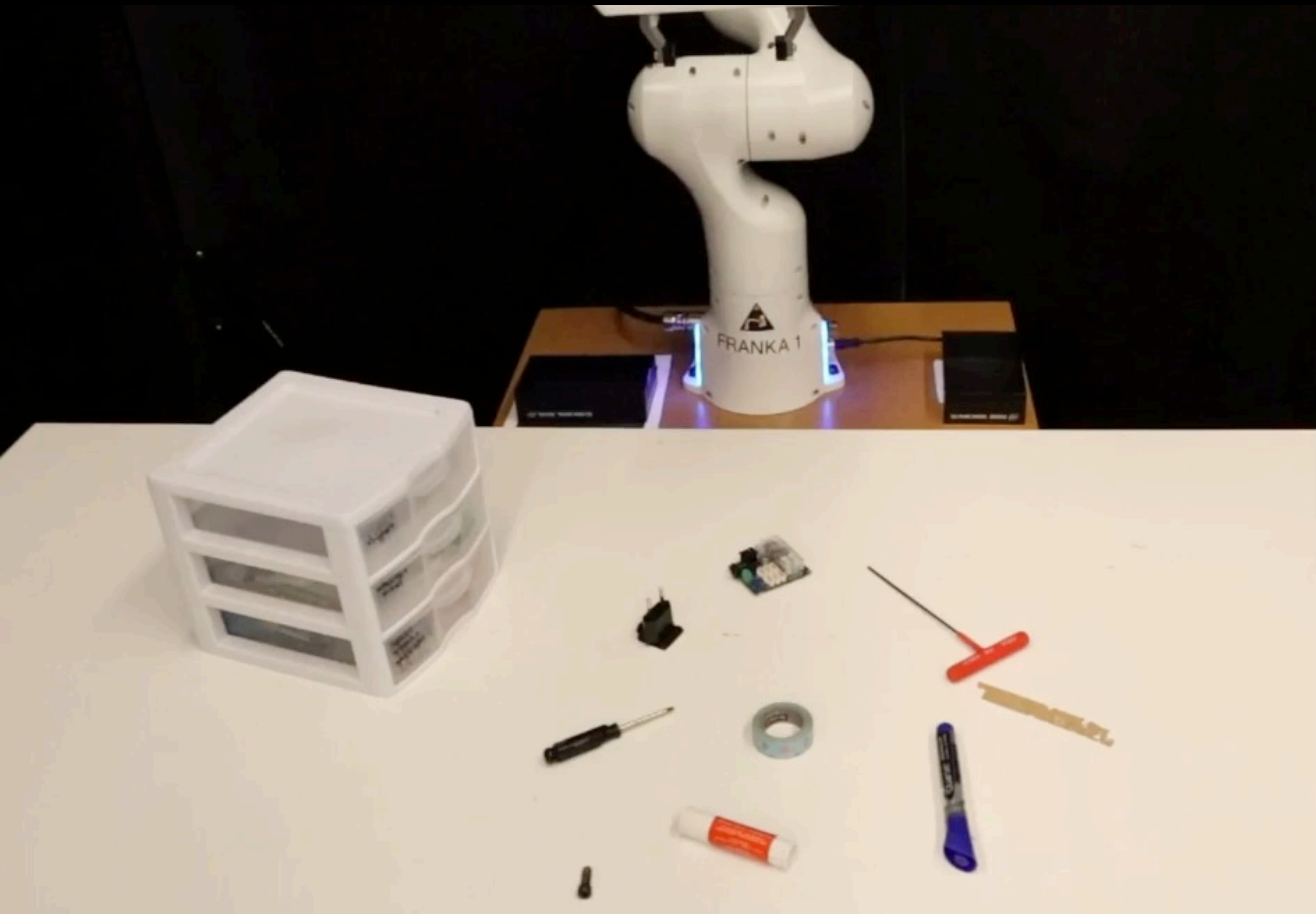
James et al

2x









PerAct Takeaways

Multi-task paradigm (from vision & NLP) might also work for robotics

PerAct Takeaways

Multi-task paradigm (from vision & NLP) might also work for robotics

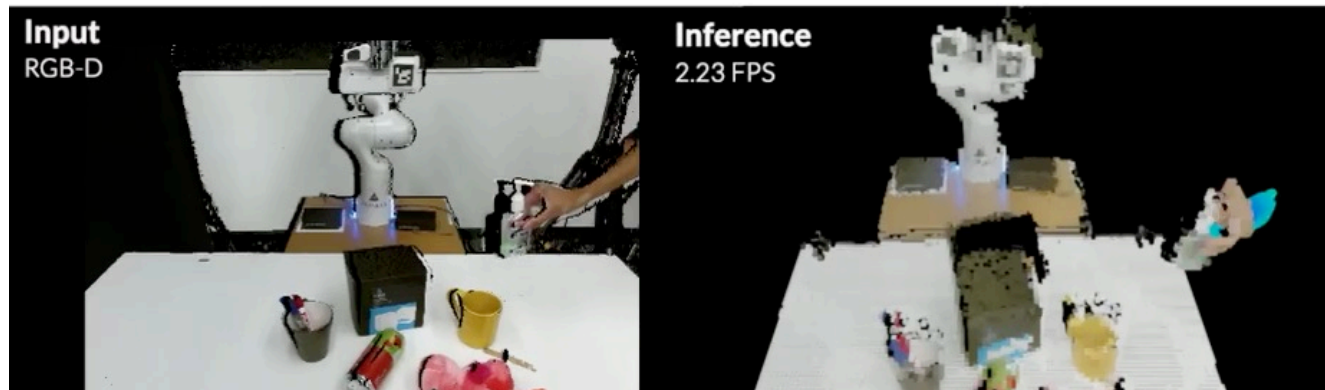
The right problem formulation can make a huge difference for scaling Transformers for robotics

PerAct Takeaways

Multi-task paradigm (from vision & NLP) might also work for robotics

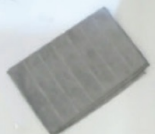
The right problem formulation can make a huge difference for scaling Transformers for robotics

Traditional robotic perception (detecting, pose estimating, grasping) might fall out of action-centric models



More 2D action detections
from **CLIPort**

t=1



Input



Pick



Place



"fold the cloth in half"

t=3



Input



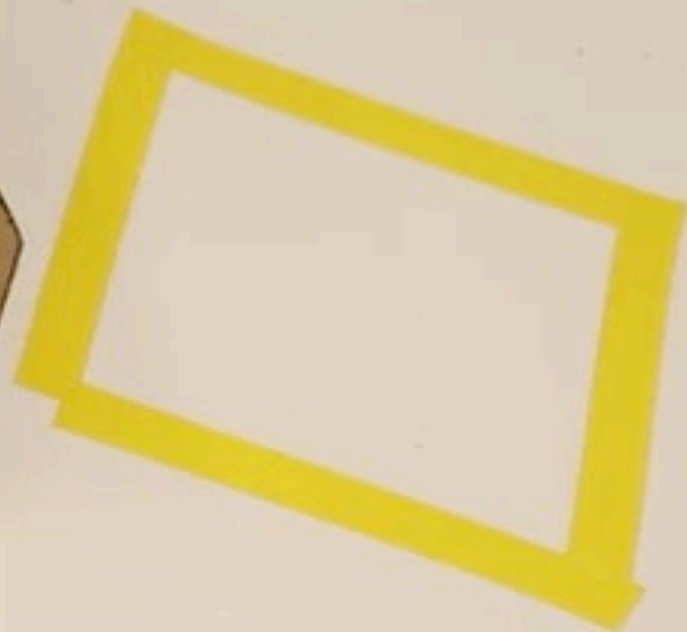
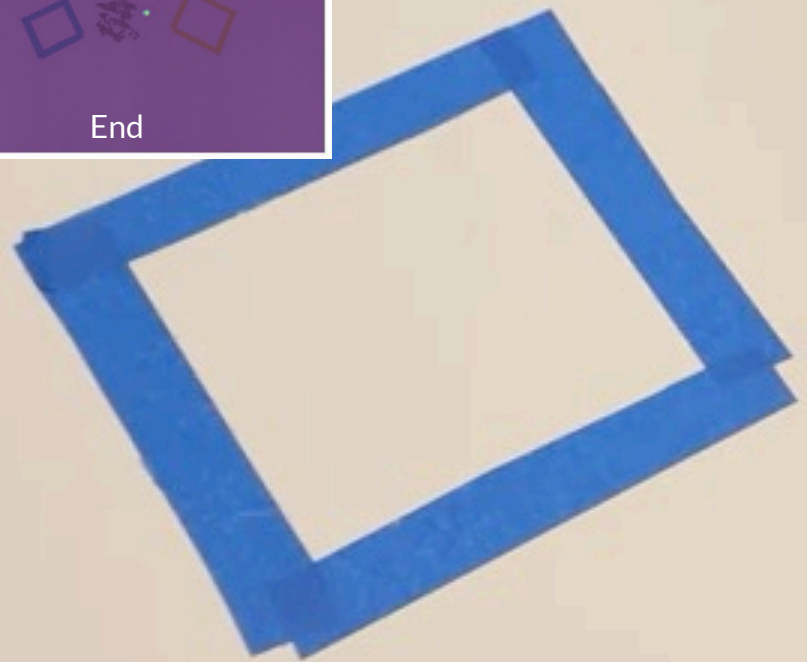
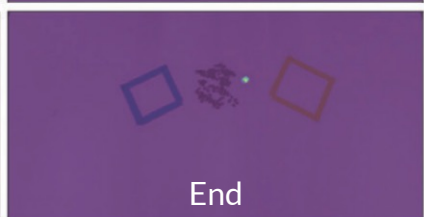
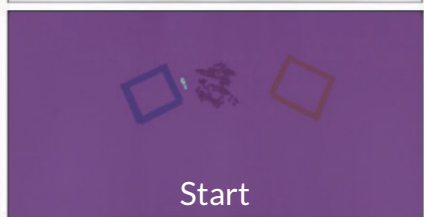
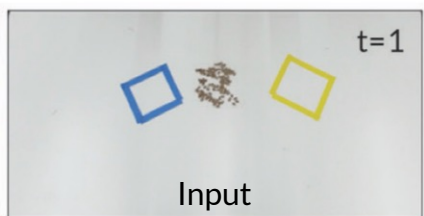
Pick



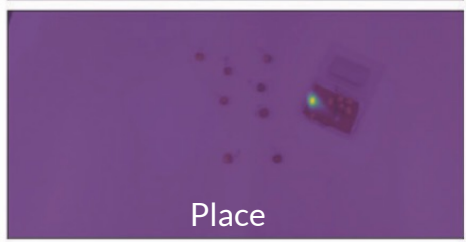
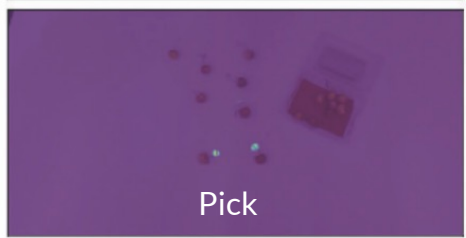
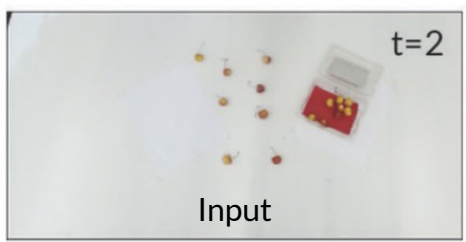
Place



"move the rook one block right"



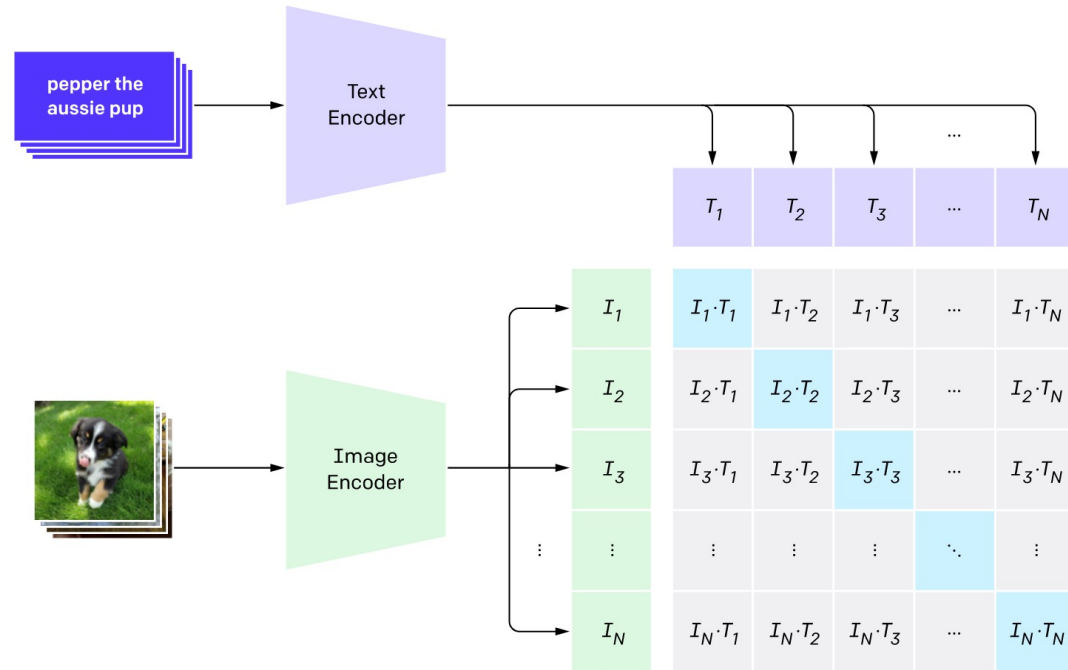
"sweep the beans into the blue zone"



David Marr

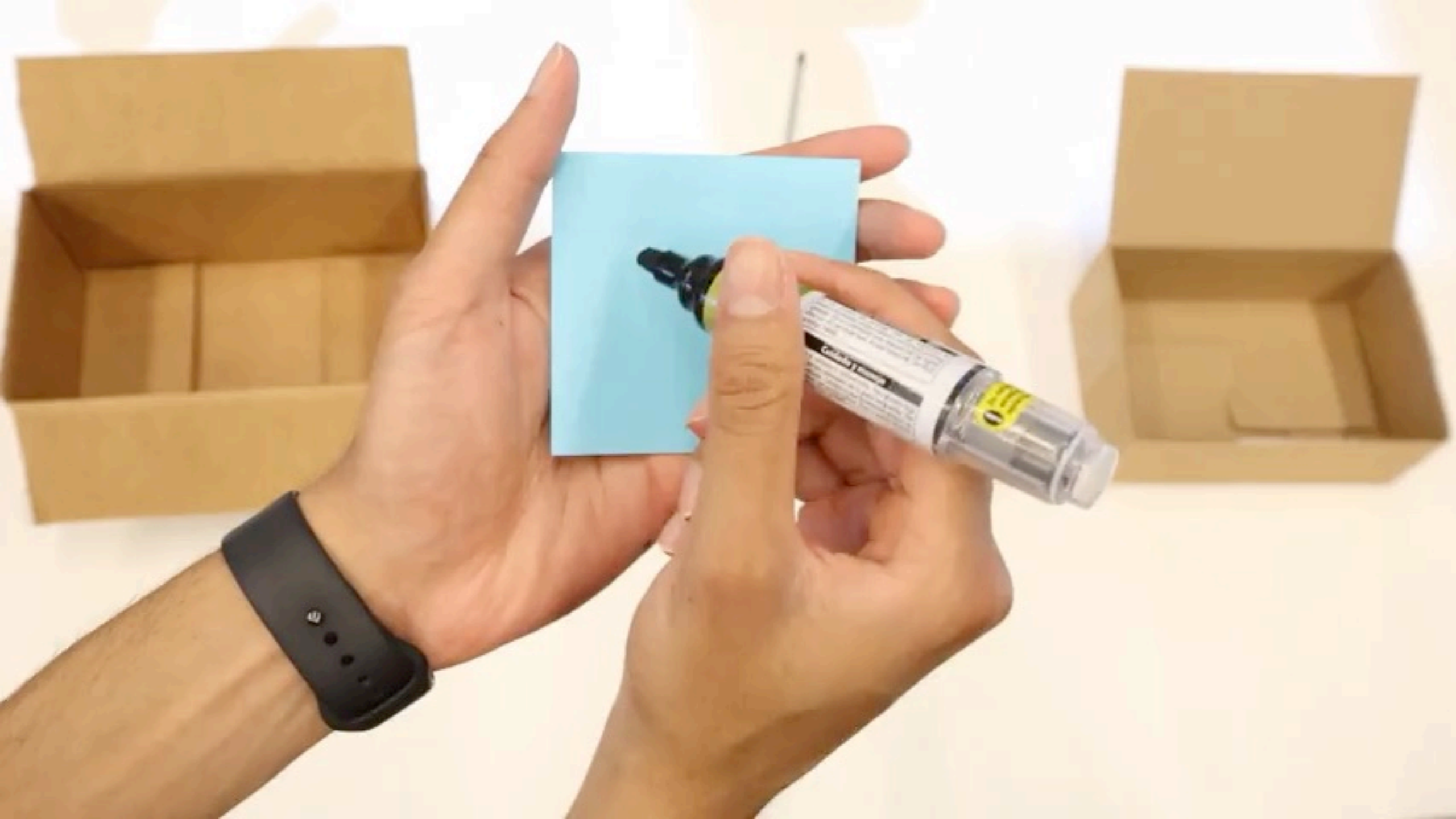


“Vision is a computational process that transforms the retinal image into an objective representation of 3D shape.”



Actionless priors (like pretrained CLIP)

can speed up learning



Three Perspectives on Vision

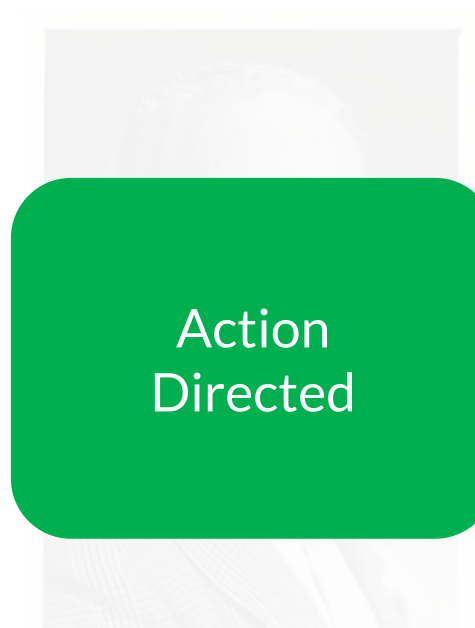
David Marr



Image
Processing

Vision Priors
like CLIP
process the
retinal image into an objective
representation of 3D shape.

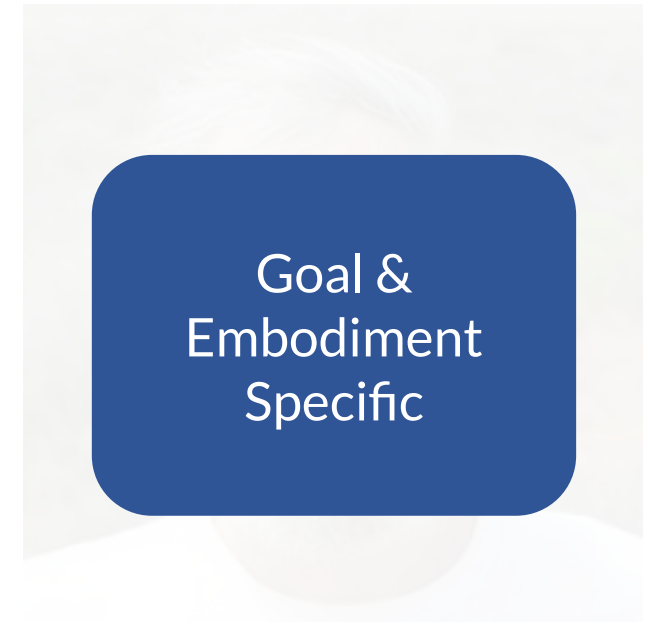
James Gibson



Action
Directed

Action
Affordances
"There is no copy in the world. There is no retinal
image. There are no representations. There is no
3D shape. There is only direct pickup of
ecologically relevant variants and invariants.
Vision is in the world, not the observer."

Jan Koenderink



Goal &
Embodiment
Specific

Language
Goals
"There is no world, only
the observer. Thus, vision
cannot be in the world but is a
creative act of the observer."

Summary

Summary



Learning **visual** representations

Summary



Learning **visual** representations of **actions**



Summary



Learning **visual** representations of **actions**
conditioned on **language**



PerAct

Paper, videos, Colab, code:
peract.github.io

CLIPort

Paper, videos, code, models:
cliport.github.io



"align the rope from back right corner to back left corner"



"pack the hexagon in the brown box"



"put the green letter E in the right letter E shape hole"



"put the blue blocks in a green bowl"



"pack all the yellow and blue blocks in the brown box"

Limitations

Hard to extend to dynamic and dexterous manipulation

Struggles with unseen objects

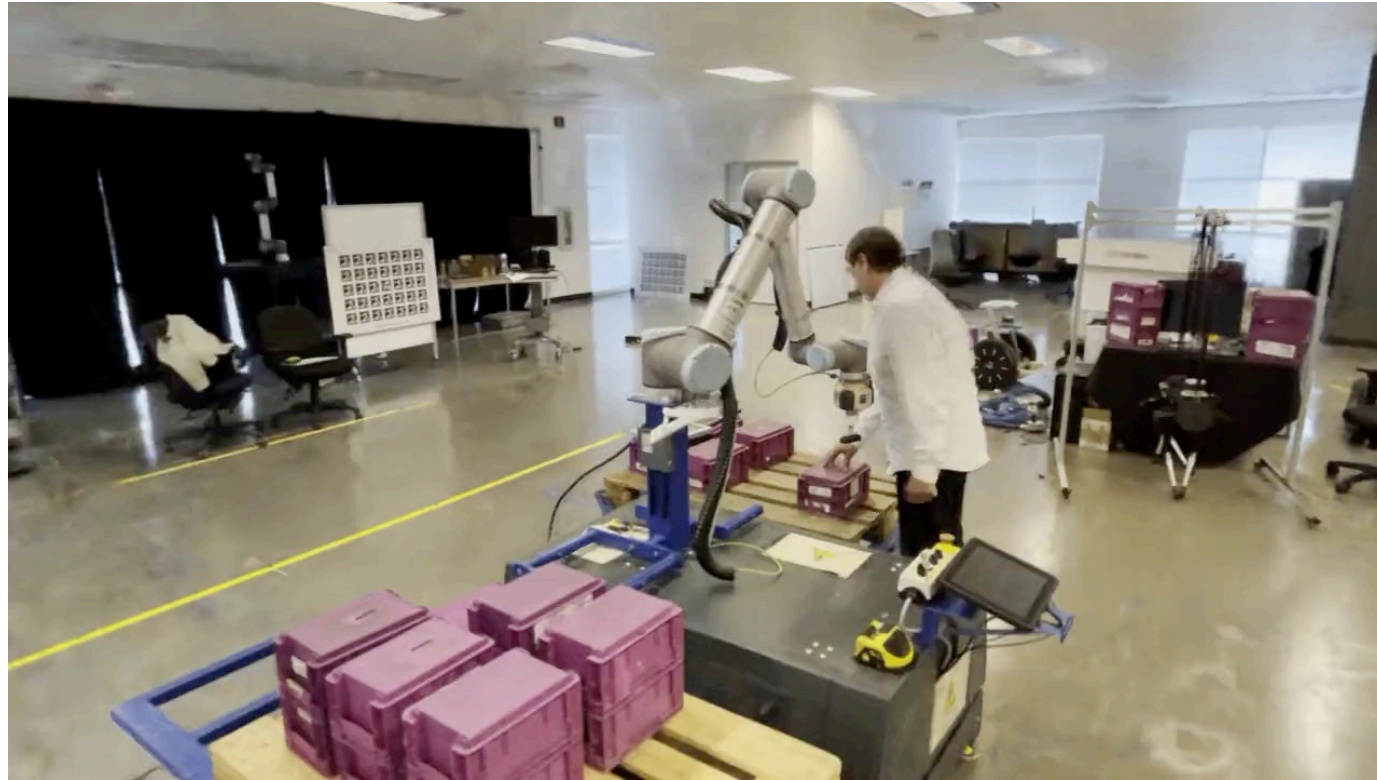
Does not predict task-completion

Struggles with complex spatial relationships

Needs good hand-eye

Scope of language (especially verbs) is mostly limited to the training distribution

What next?



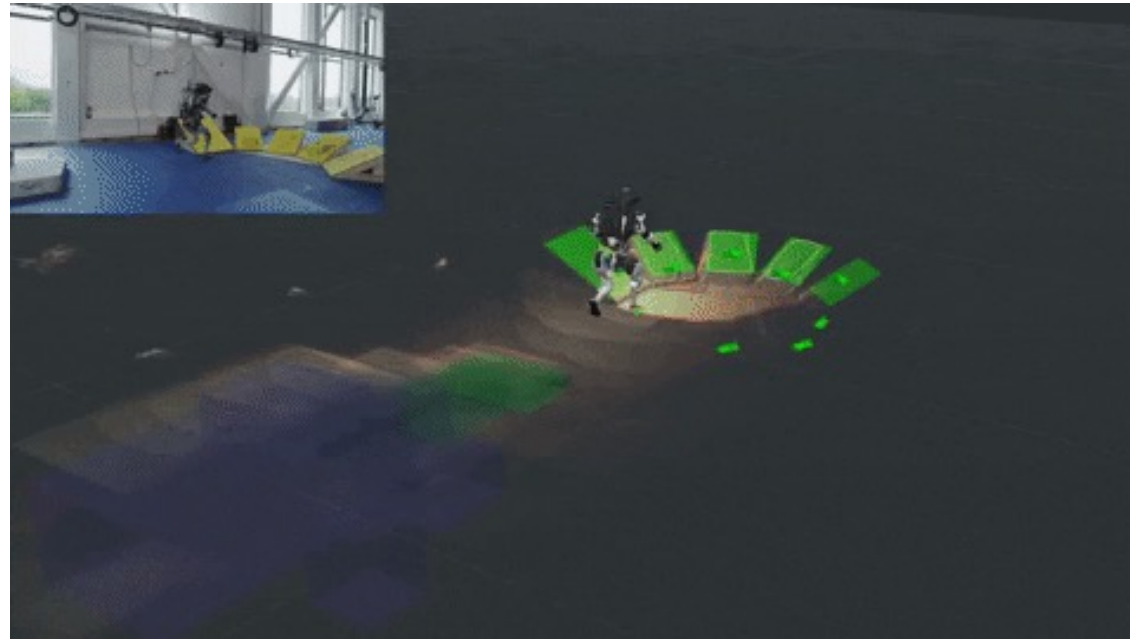
NeRF voxel features?

What next?



NeRF voxel features?

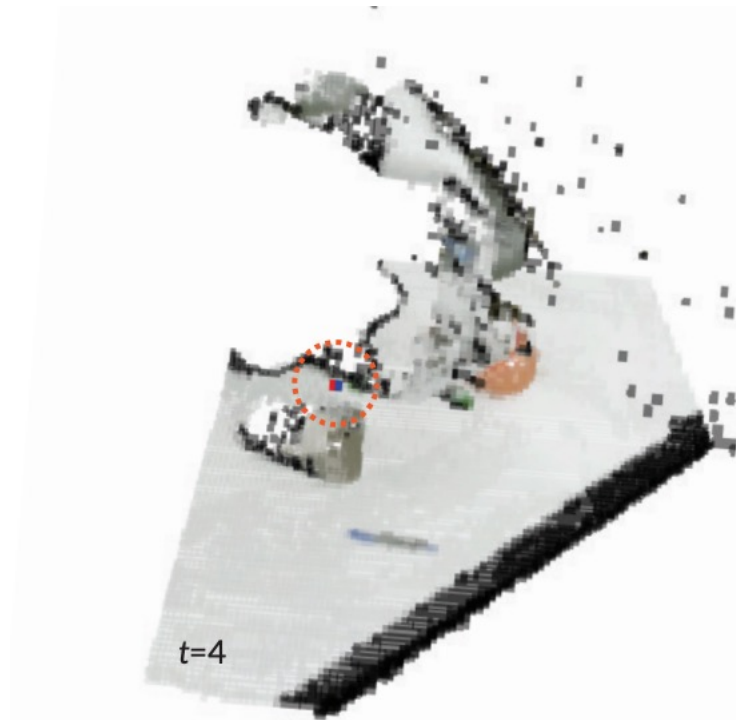
What next?



Detect footprint poses? bimanual grippers poses? finger tips?

Appendix

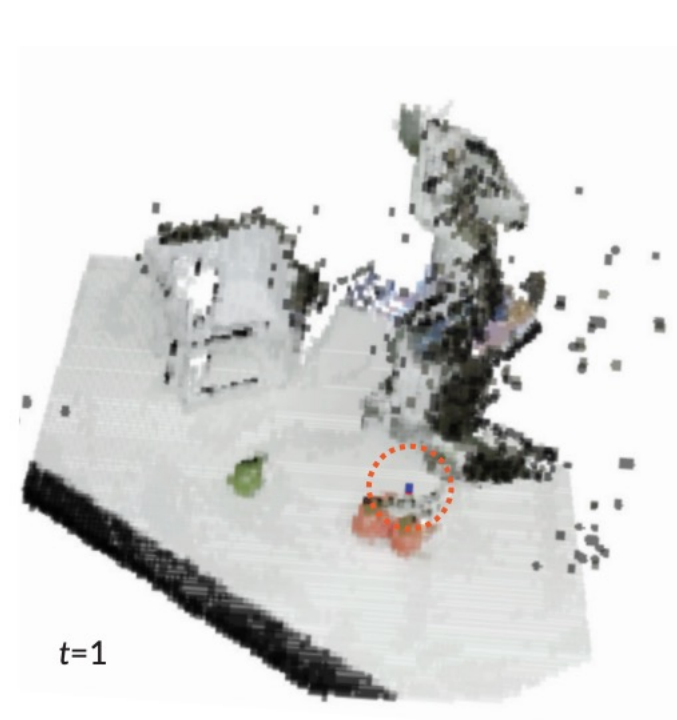
Data Augmentation



"put the green whiteboard marker in the mug"



"sweep the beans into the gray dustpan"



"put the tomatoes in the top bin"

Data Augmentation

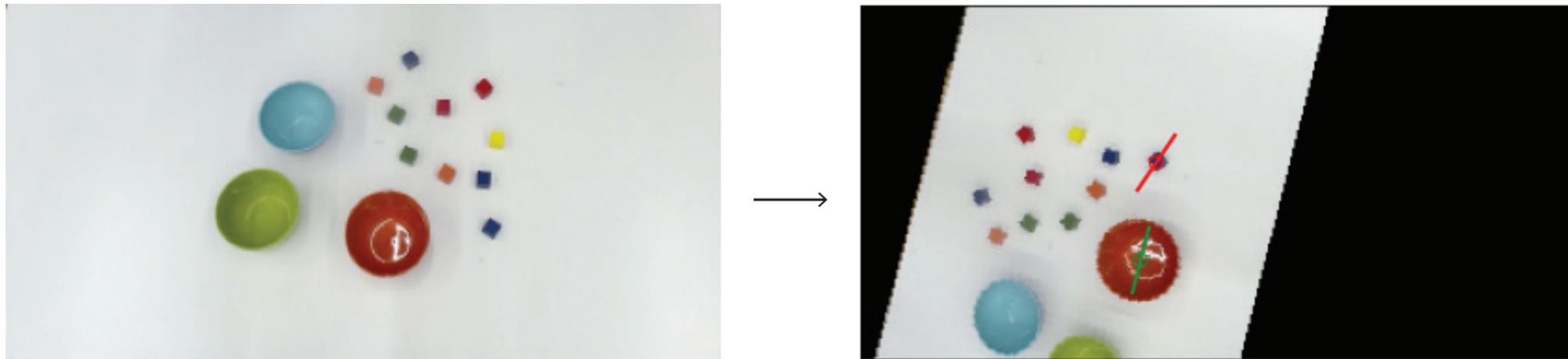


Figure 9. Data Augmentation: $SE(2)$ transform applied to RGB-D input. The left image shows the original input, and the right image shows the transformed input along with expert \mathcal{T}_{pick} (red) and \mathcal{T}_{place} (green) actions.

Perturbation Tests

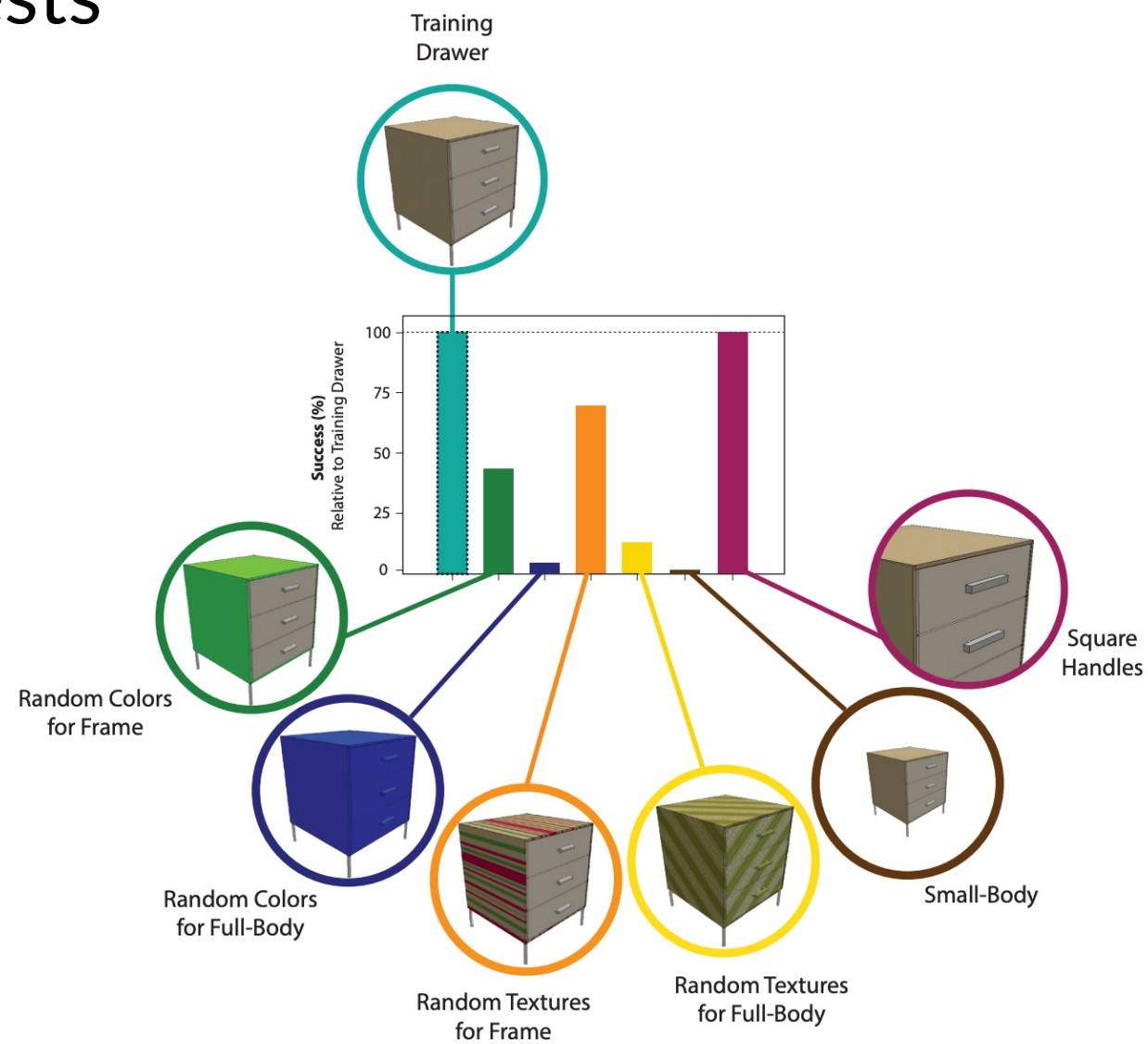
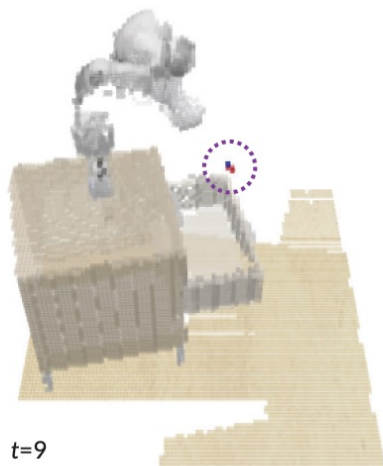


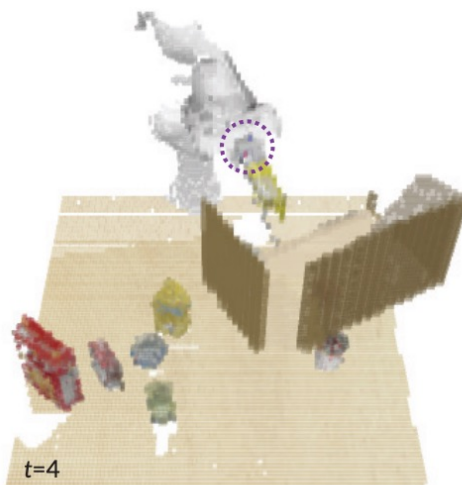
Figure 11. Perturbation Tests. Results from a multi-task PERACT agent trained on a single drawer and evaluated on several instances perturbed drawers. Each perturbation consists of 25 evaluation episodes, and reported successes are relative to the training drawer.

More Affordances

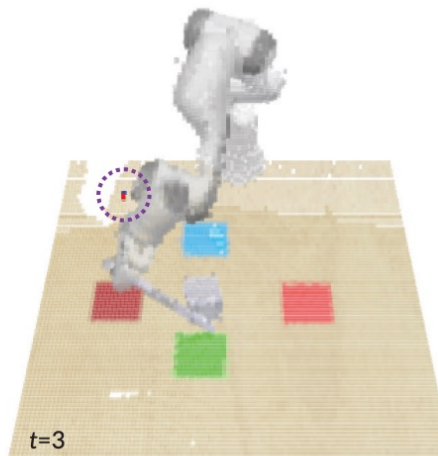
● Q-Prediction ● Expert Action



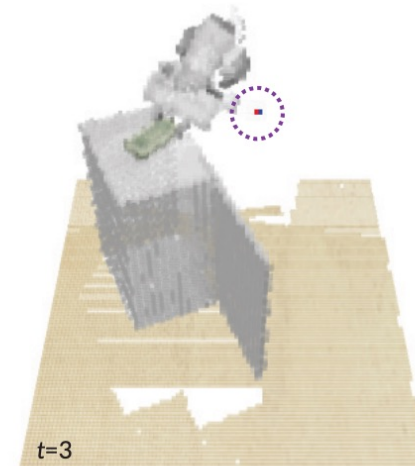
t=9
"put the item
in the middle drawer"



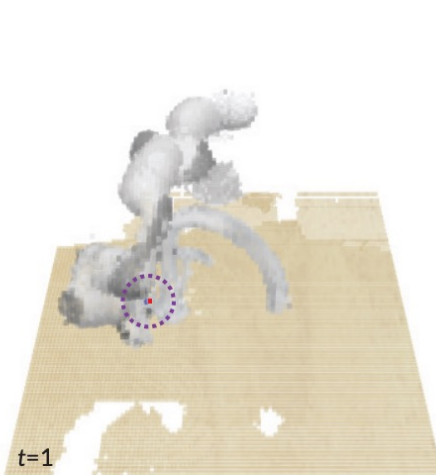
t=4
"put the sugar
in the cupboard"



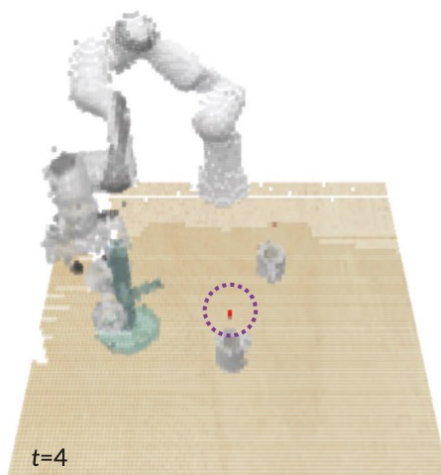
t=3
"use the stick to drag the cube
onto the blue target"



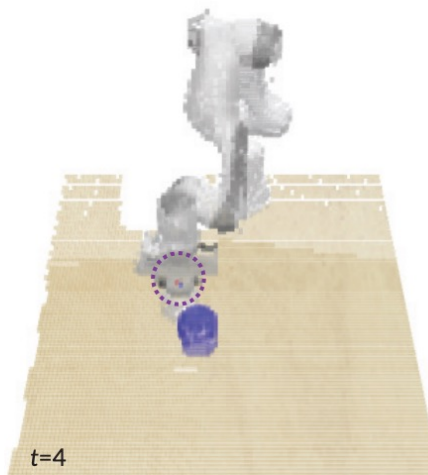
t=3
"put the money away in the safe
on the top shelf"



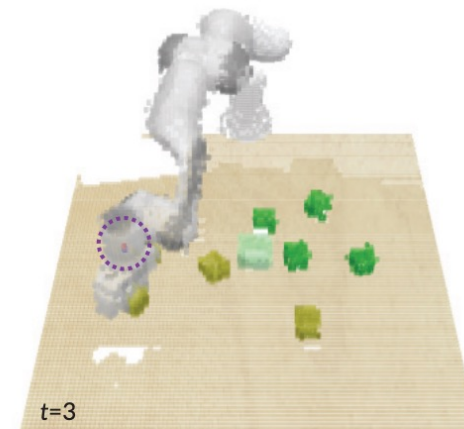
t=1
"turn right tap"



t=4
"place 3 mugs
on the cup holder"



t=4
"close the gray jar"



t=3
"stack 2 olive blocks"

More Affordances

"pack all the yellow and blue blocks into the brown box"

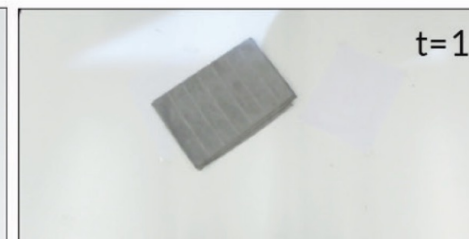
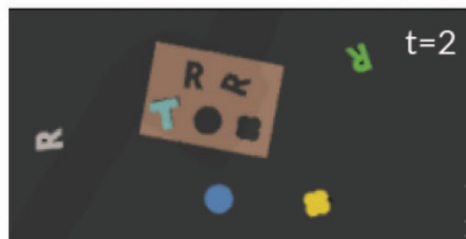
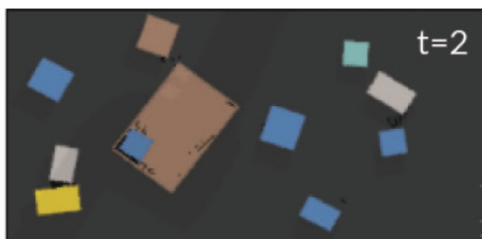
"put the green letter R shape in the right R shape hole"

"pack the white tape in the brown box"

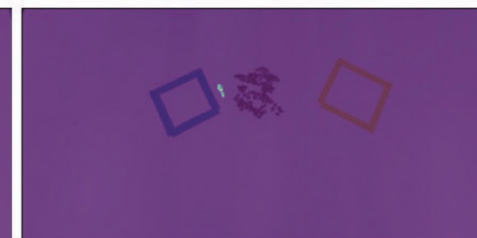
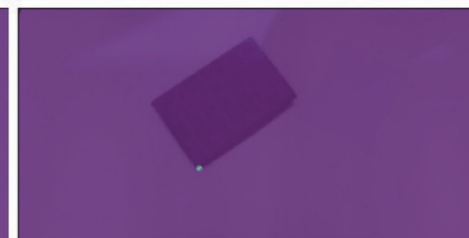
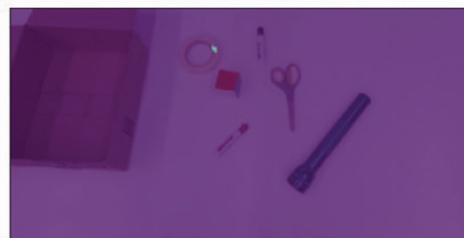
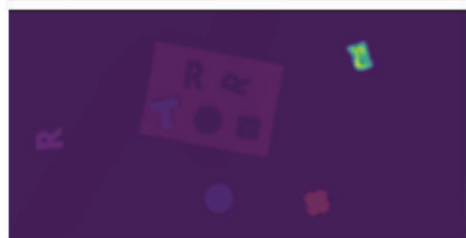
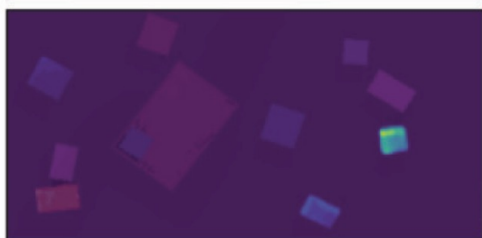
"unfold the cloth"

"sweep the beans into the yellow zone"

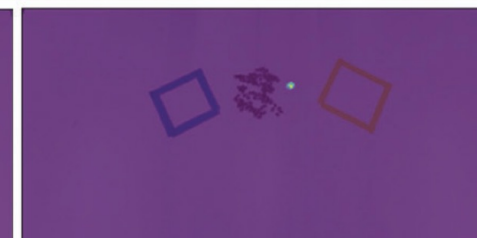
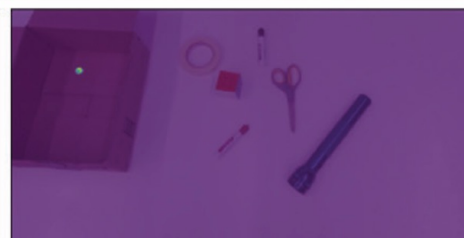
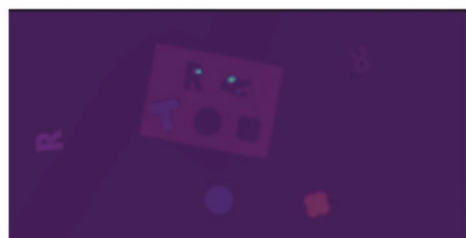
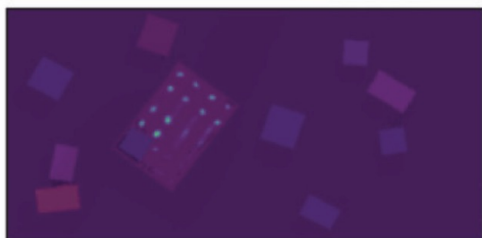
Input



Pick



Place



More Affordances

"align the rope from front right tip to back right tip to back right"

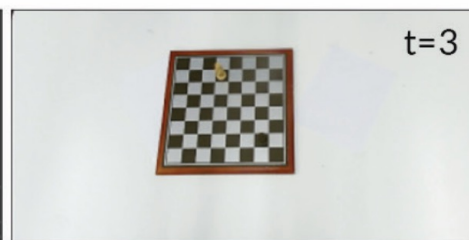
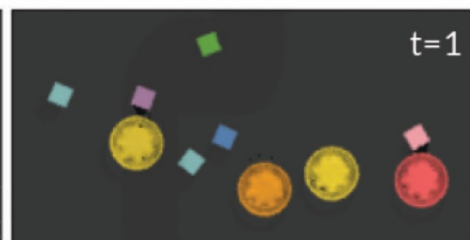
"pack the black shoe with orange stripes in the brown box"

"push the pile of yellow blocks into the brown square"

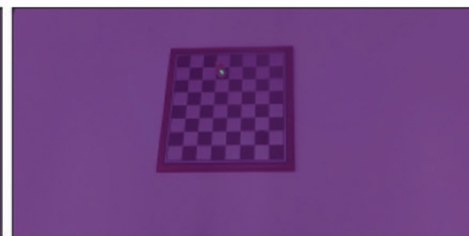
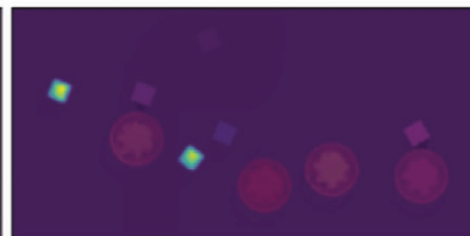
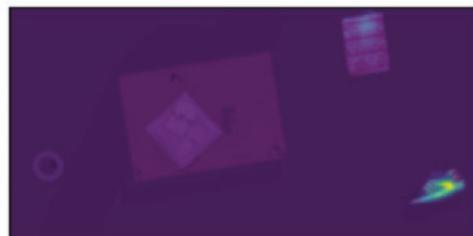
"put the cyan blocks in the yellow bowl"

"move the rook one block forward"

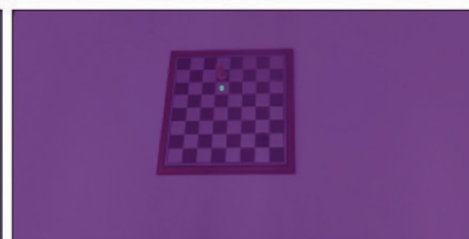
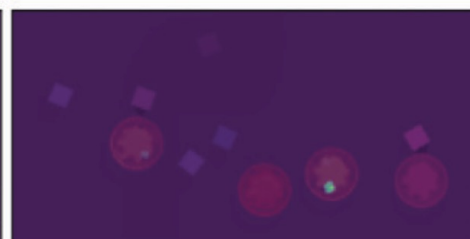
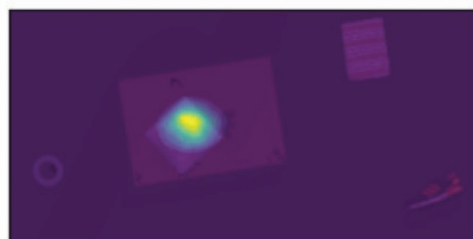
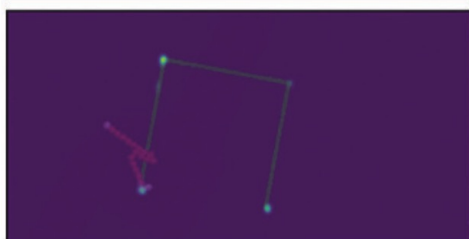
Input



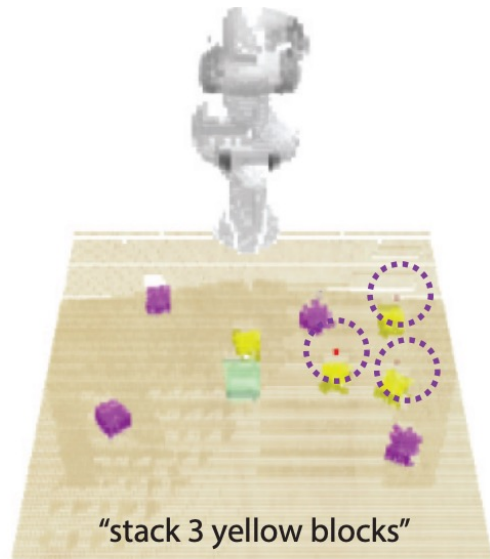
Pick



Place



Multi Modal Actions



"pack all the yellow and blue blocks into the brown box"

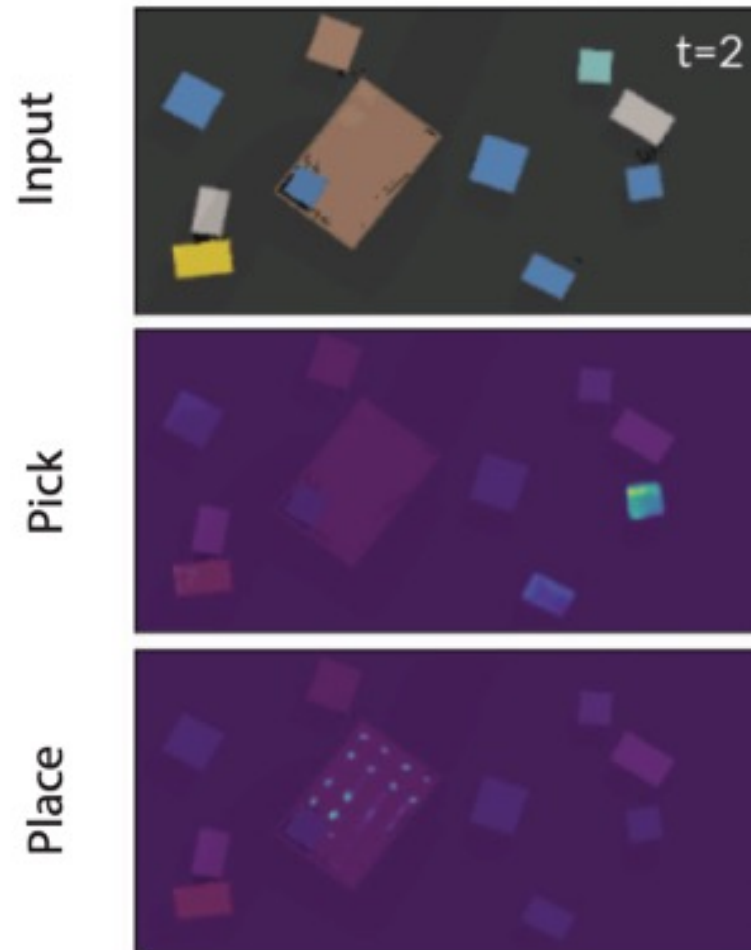


Figure 13. Examples of Multi-Modal Predictions.