

**OCT 13 2020**

# **NLP + ACTION**

---

**PRESENTERS: PRIYANKA MANDIKAL, BINGYE LI**

# NLP + Action

- ❖ Trends in language understanding
  - ❖ Plain Text Corpora → NLP + Vision → NLP + Vision + Action

> 2015

> 2018

no grounding

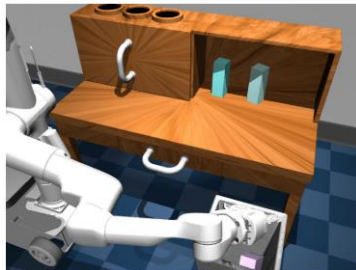
static datasets

active perception

Source	An admitting privilege i to carry out a diagnosis
Reference	Le privilège d'admissio d'un hôpital, d'admettr diagnostic ou un traiten
RNNenc-50	Un privilège d'admissio centre médical d'un dia
RNNsearch-50	Un privilège d'admissio centre médical pour effe soins de santé à l'hôpit
Google Translate	Un privilège admettre centre médical pour eff que travailleur de soins



\*construction worker in orange safety vest is working on road.\*



now: put the object in the trash  
next:

- ❖ Language + Embodied AI
- ❖ Instruction following in realistic situated environments – egocentric RGB cameras, agent actions

# NLP + Action

1. ALFRED, A Benchmark for Interpreting Grounded Instructions for Everyday Tasks, CVPR 2020
  - ❖ Egocentric vision + Natural language instructions → Action sequences for household tasks
2. Grounding Language in Play, arXiv 2020
  - ❖ Teleoperated play + Natural language instructions → Continuous robotic control

**OCT 13 2020**

# **ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks**

---

Mohit Shridhar, Jesse Thomason,  
Daniel Gordon, Yonatan Bisk, Winson  
Han, Roozbeh Mottaghi, Luke  
Zettlemoyer, Dieter Fox


**PRESENTER: BINGYE LI**

# Introduction

ALFRED is a benchmark for learning a mapping from natural language instructions and egocentric vision to sequences of actions for household tasks.

**Goal:** "Rinse off a mug and place it in the coffee maker"



	— Language —		— Virtual Environment —			— Inference —		
	# Human Annotations	Granularity	Visual Quality	Movable Objects	State Changes	Vis. Obs.	Navigation	Interaction
TACoS [43]	17k+	High&Low	Photos	✗	✗	–	–	–
R2R [3]; Touchdown [14]	21k+; 9.3k+	Low	Photos	✗	✗	Ego	Graph	✗
EQA [15]	✗	High	Low	✗	✗	Ego	Discrete	✗
Matterport EQA [55]	✗	High	Photos	✗	✗	Ego	Discrete	✗
IQA [20]	✗	High	High	✗	✓	Ego	Discrete	Discrete
VirtualHome [42]	2.7k+	High&Low	High	✓	✓	3 <sup>rd</sup> Person	✗	Discrete
VSP [58]	✗	High	High	✓	✓	Ego	✗	Discrete
ALFRED 	25k+	High&Low	High	✓	✓	Ego	Discrete	Discrete + Mask

- Include both high-level goal and low-level natural language instructions.
- Include object and state interactions.
- Enable discretized, grid-based movement rather than topological graph navigation.
- Require spatially located interaction masks instead of choosing from a set of object classes.

## Related Work

- Vision & Language Navigation:
  - Navigation in static environment
  - No object interactions and state changes
- Vision & Language Task Completion
  - Based on simpler block worlds and fully observable scenes
  - AI2-THOR, an interactive 3D environment for visual AI, where AI agents can navigate in the scenes and interact with objects to perform tasks.

## Related Work

- Embodied Question Answering
  - Question answering using templated language or static scenes
  - No task completion
- Instruction Alignment
  - Learning visual correspondence from recorded videos
  - Not in an interactive setting
- Robotics Instruction Following
  - Consider different tasks individually
  - Limited to fewer scenes and objects



# ALFRED Dataset

	Pick & Place	Stack & Place	Pick Two & Place	Clean & Place	Heat & Place	Cool & Place	Examine in Light
item(s)	Book	Fork (in) Cup	Spray Bottle	Dish Sponge	Potato Slice	Egg	Credit Card
receptacle	Desk	Counter Top	Toilet Tank	Cart	Counter Top	Side Table	Desk Lamp
scene #	Bedroom 14	Kitchen 10	Bathroom 2	Bathroom 1	Kitchen 8	Kitchen 21	Bedroom 24
expert demonstration							



	Annotation # 1	Annotation # 2	Annotation # 3
Goals	Put a clean sponge on a metal rack.	Place a clean sponge on the drying rack	Put a rinsed out sponge on the drying rack
Instructions	Go to the left and face the faucet side of the bath tub. Pick up left most green sponge from the bath tub. Turn around and go to the sink. Put the sponge in the sink. Turn on then turn off the water. Take the sponge from the sink. Go to the metal bar rack to the left. Put the sponge on the top rack to the left of the lotion bottle.	Turn around and walk over to the bathtub on the left. Grab the sponge out of the bathtub. Turn around and walk to the sink ahead. Rinse the sponge out in the sink. Move to the left a bit and face the drying rack in the corner of the room. Place the sponge on the drying rack.	Walk forwards a bit and turn left to face the bathtub. Grab a sponge out of the bathtub. Turn around and walk forwards to the sink. Rinse the sponge out in the sink and pick it up again. Turn left to walk a bit, then face the drying rack. Put the sponge on the drying rack.

ALFRED includes 25,743 English language directives describing 8,055 expert demonstrations averaging 50 steps each, resulting in 428,322 image-action pairs.



# Expert demonstrations

- Expert demonstrations are composed of an agent's egocentric visual observations of the environment, actions taken at each timestep, and ground-truth interaction masks.
- Navigation actions: MoveAhead, RotateRight, RotateLeft, LookUp, and LookDown.
- Manipulation actions: Pickup, Put, Open, Close, ToggleOn, ToggleOff, and Slice.

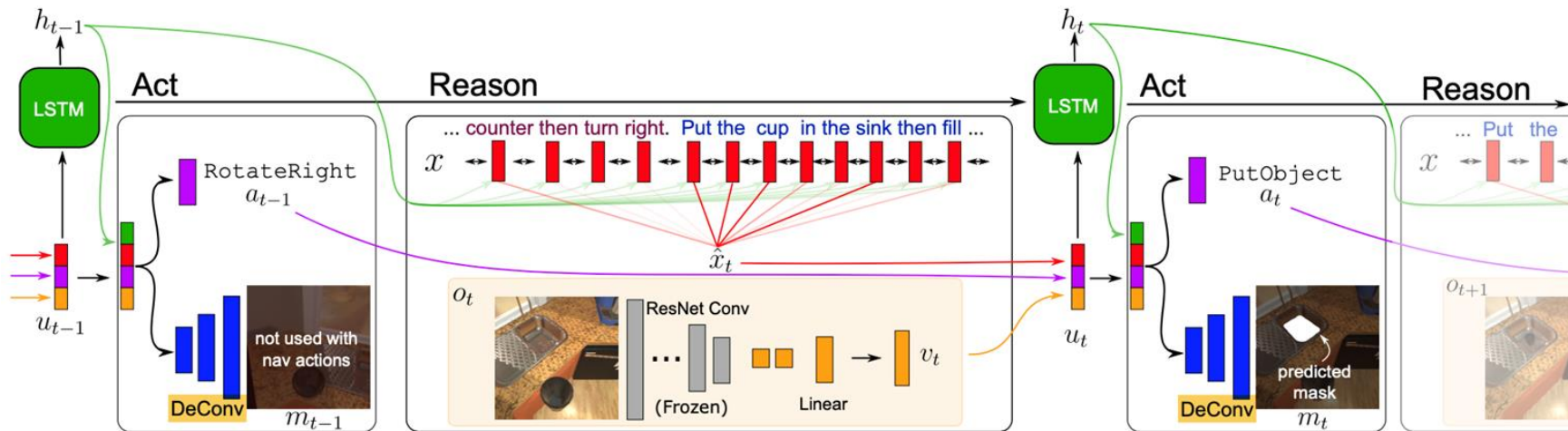
# Language directives

- Language directives include a high-level goal together with low-level instructions.
- AMT workers write low-level, step-by-step instructions for each highlighted sub-goal segment.
- Ex: *“Walk to the coffee maker on the right.”*
- They also write a high-level goal that summarizes what the robot should accomplish during the expert demonstration.
- Ex: *“Rinse off a mug and place it in the coffee maker.”*

# Sequence-to-Sequence Model

- A bidirectional-LSTM generates a representation of the language input
- A CNN encodes the visual input
- A decoder LSTM infers a sequence of low-level actions while attending over the encoded language
- At each timestep, the model produces the expert action and associated interaction mask for manipulation actions.

# Sequence-to-Sequence Model

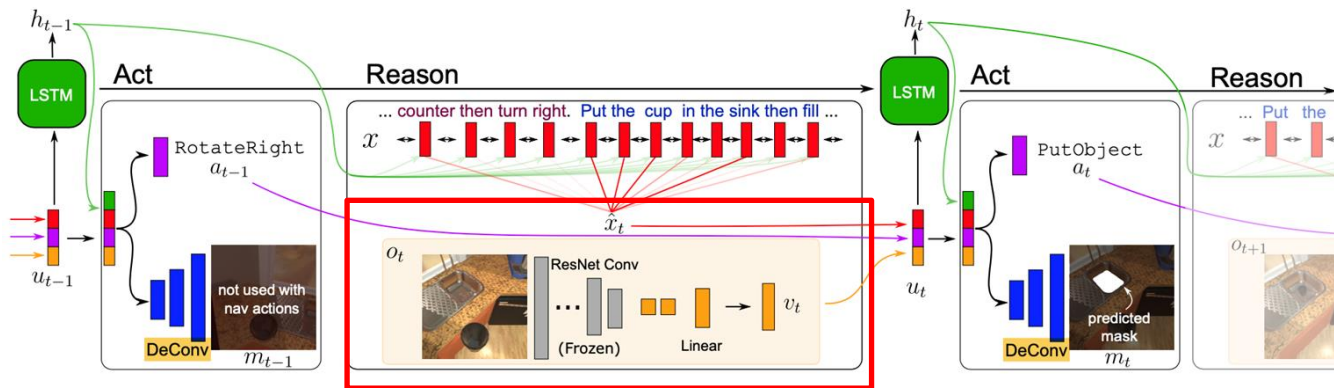


# Language encoding

- a natural language goal  $\bar{G} = \langle g_1, g_2, \dots, g_{L_g} \rangle$
- step-by-step instructions  $\bar{S} = \langle s_1, s_2 \dots s_{L_s} \rangle$
- Construct a single input sequence  $\bar{X} = \langle g_1, g_2, \dots, g_{L_g}, \langle \text{SEP} \rangle, s_1, s_2 \dots s_{L_s} \rangle$
- Fed the sequence into a bidirectional LSTM encoder to produce an encoding

$$x = \{x_1, x_2, \dots, x_{L_g+L_s}\}$$

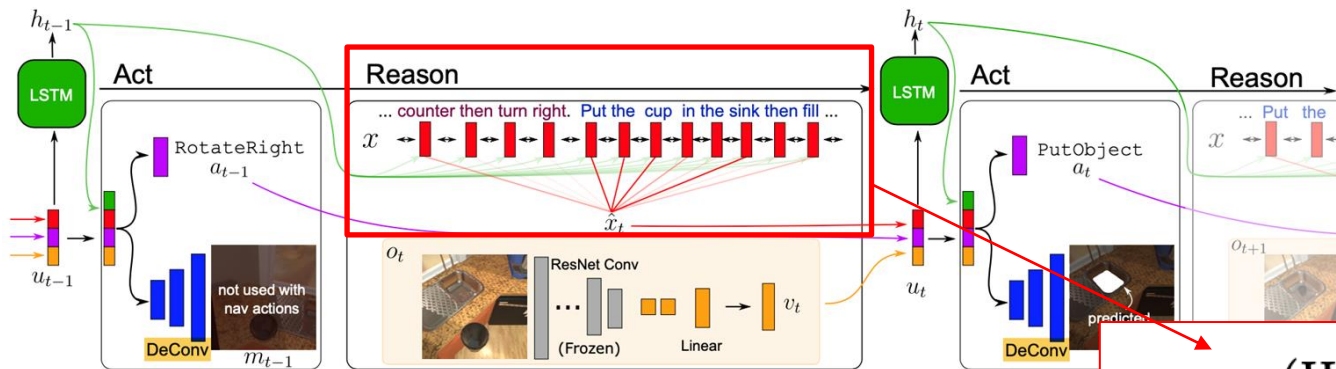
# Visual Encoding



Each visual observation  $o_t$  is encoded with a frozen ResNet-18 CNN followed by two more  $1 \times 1$  convolution layers and a fully-connected layer.



# Attention over language



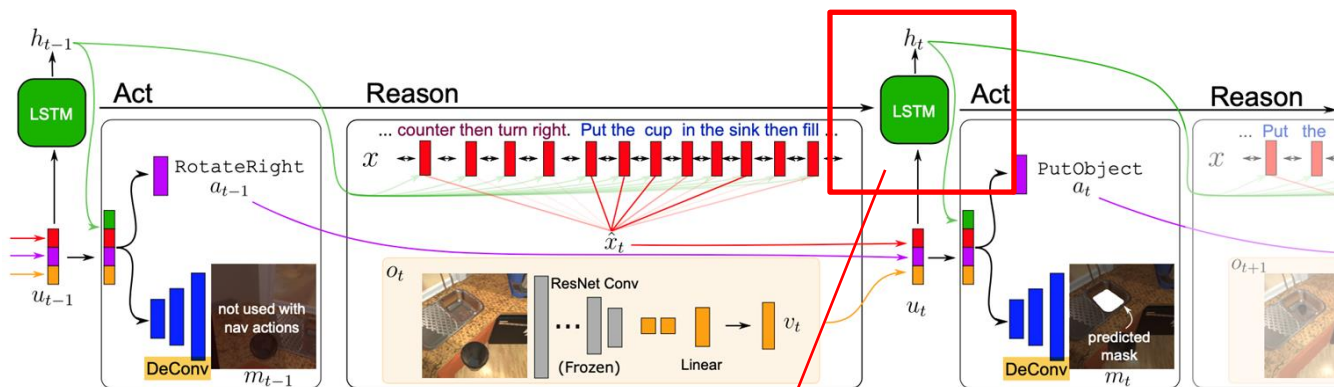
perform soft-attention on the language features  $x$

$$z_t = (W_x h_{t-1})^\top x,$$

$$\alpha_t = \text{Softmax}(z_t),$$

$$\hat{x}_t = \alpha_t^\top x$$

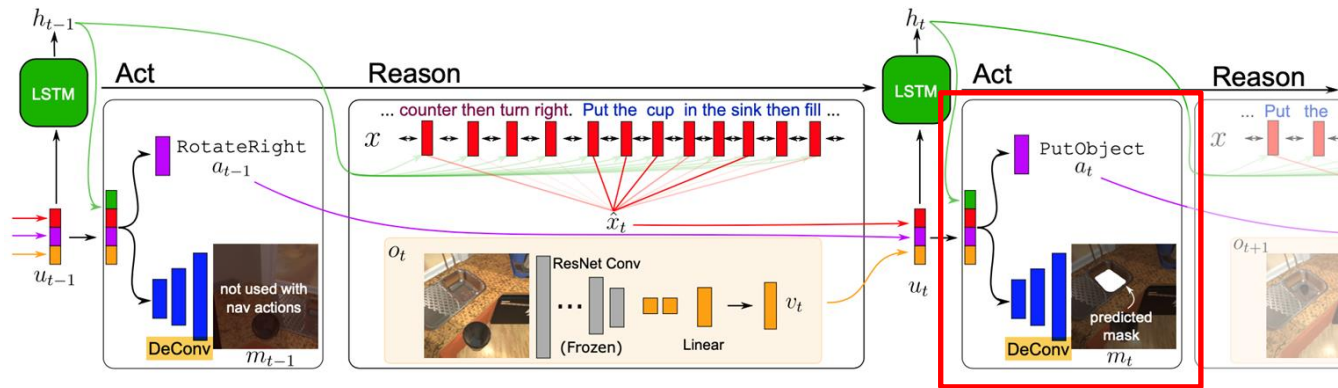
# Action decoding



$$u_t = [v_t; \hat{x}_t; a_{t-1}],$$

$$h_t = \text{LSTM}(u_t, h_{t-1})$$

# Action and mask prediction



Action loss: softmax cross entropy

Mask loss: binary cross entropy

rebalanced for sparsity

$$a_t = \operatorname{argmax} (W_a [h_t; u_t])$$

$$m_t = \sigma (\mathbf{deconv} [h_t; u_t])$$

# Progress Monitors

Progress prediction helps learn the utility of each state in the process of achieving the overall task.

$p_t = \sigma (W_p [h_t; u_t]) \in [0, 1]$  normalized time-stamp value

Sub-goal prediction encourages the agent to coarsely track its progress through the language directive.

$c_t = \sigma (W_c [h_t; u_t]) \in [0, 1]$  normalized number of completed sub-goals

# Evaluation Metrics

- Task Success
- Goal-Condition Success
- Path Weighted Metrics
- Sub-Goal Evaluation

# Task Success

1 if the object positions and state changes correspond correctly to the task goal-conditions at the end of the action sequence, and 0 otherwise.

*“Put a hot potato slice on the counter”*

*succeed if any potato slice object has changed to the heated state and is resting on any counter top surface.*

# Goal-Condition Success

The goal-condition success of a model is the ratio of goal-conditions completed at the end of an episode to those necessary to have finished a task.

*“Put a hot potato slice on the counter”*

- *a potato must be sliced*
- *a potato slice should become heated*
- *a potato slice should come to rest on a counter top.*
- *the same potato slice that is heated should be on the counter top.*

## Path Weighted Metrics

Both two metrics have a path weighted version, that considers the length of the expert demonstration.

The path weighted score  $P_s$  for metric  $s$  is given as  $p_s = s \times \frac{L^*}{\max(L^*, \hat{L})}$

where  $\hat{L}$  is the number of actions the model took in the episode, and  $L^*$  is the number of actions in the expert demonstration.



# Goal-Condition Success

The ability of a model to accomplish the next sub-goal conditioned on the preceding expert sequence.

# Analysis

Model	Validation				Test			
	<i>Seen</i>		<i>Unseen</i>		<i>Seen</i>		<i>Unseen</i>	
	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond
NO LANGUAGE	0.0 (0.0)	5.9 (3.4)	0.0 (0.0)	6.5 (4.7)	0.2 (0.0)	5.0 (3.2)	0.2 (0.0)	6.6 (4.0)
NO VISION	0.0 (0.0)	5.7 (4.7)	0.0 (0.0)	6.8 (6.0)	0.0 (0.0)	3.9 (3.2)	0.2 (0.1)	6.6 ( <b>4.6</b> )
GOAL-ONLY	0.1 (0.0)	6.5 (4.3)	0.0 (0.0)	6.8 (5.0)	0.1 (0.1)	5.0 (3.7)	0.2 (0.0)	6.9 (4.4)
INSTRUCTIONS-ONLY	2.3 (1.1)	9.4 (6.1)	0.0 (0.0)	<b>7.0</b> (4.9)	2.7 (1.4)	8.2 (5.5)	<b>0.5 (0.2)</b>	7.2 ( <b>4.6</b> )
SEQ2SEQ	2.4 (1.1)	9.4 (5.7)	<b>0.1</b> (0.0)	6.8 (4.7)	2.1 (1.0)	7.4 (4.7)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM PROGRESS-ONLY	2.1 (1.1)	8.7 (5.6)	0.0 (0.0)	6.9 (5.0)	3.0 (1.7)	8.0 (5.5)	0.3 (0.1)	<b>7.3</b> (4.5)
+ PM SUBGOAL-ONLY	2.1 (1.2)	9.6 (5.5)	0.0 (0.0)	6.6 (4.6)	3.8 (1.7)	8.9 (5.6)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM Both	<b>3.7 (2.1)</b>	<b>10.0 (7.0)</b>	0.0 (0.0)	6.9 ( <b>5.1</b> )	<b>4.0 (2.0)</b>	<b>9.4 (6.3)</b>	0.4 (0.1)	7.0 (4.3)
HUMAN	-	-	-	-	-	-	91.0 (85.8)	94.5 (87.6)

~8% goal-condition success rate (partially complete tasks)

Model	Validation				Test			
	Seen		Unseen		Seen		Unseen	
	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond
NO LANGUAGE	0.0 (0.0)	5.9 (3.4)	0.0 (0.0)	6.5 (4.7)	0.2 (0.0)	5.0 (3.2)	0.2 (0.0)	6.6 (4.0)
NO VISION	0.0 (0.0)	5.7 (4.7)	0.0 (0.0)	6.8 (6.0)	0.0 (0.0)	3.9 (3.2)	0.2 (0.1)	6.6 ( <b>4.6</b> )
GOAL-ONLY	0.1 (0.0)	6.5 (4.3)	0.0 (0.0)	6.8 (5.0)	0.1 (0.1)	5.0 (3.7)	0.2 (0.0)	6.9 (4.4)
INSTRUCTIONS-ONLY	2.3 (1.1)	9.4 (6.1)	0.0 (0.0)	<b>7.0</b> (4.9)	2.7 (1.4)	8.2 (5.5)	<b>0.5 (0.2)</b>	7.2 ( <b>4.6</b> )
SEQ2SEQ	2.4 (1.1)	9.4 (5.7)	<b>0.1</b> (0.0)	6.8 (4.7)	2.1 (1.0)	7.4 (4.7)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM PROGRESS-ONLY	2.1 (1.1)	8.7 (5.6)	0.0 (0.0)	6.9 (5.0)	3.0 (1.7)	8.0 (5.5)	0.3 (0.1)	<b>7.3</b> (4.5)
+ PM SUBGOAL-ONLY	2.1 (1.2)	9.6 (5.5)	0.0 (0.0)	6.6 (4.6)	3.8 (1.7)	8.9 (5.6)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM Both	<b>3.7 (2.1)</b>	<b>10.0 (7.0)</b>	0.0 (0.0)	6.9 ( <b>5.1</b> )	<b>4.0 (2.0)</b>	<b>9.4 (6.3)</b>	0.4 (0.1)	7.0 (4.3)
HUMAN	-	-	-	-	-	-	91.0 (85.8)	94.5 (87.6)

- Vision and language modalities are necessary to accomplish the tasks.
- The NO LANGUAGE model finishes some goal-conditions by interacting with familiar objects seen during training.
- The NO VISION model similarly finishes some goal-conditions by following low-level language instructions for navigation and memorizing interaction masks for common objects.

Model	Validation				Test			
	Seen		Unseen		Seen		Unseen	
	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond
NO LANGUAGE	0.0 (0.0)	5.9 (3.4)	0.0 (0.0)	6.5 (4.7)	0.2 (0.0)	5.0 (3.2)	0.2 (0.0)	6.6 (4.0)
NO VISION	0.0 (0.0)	5.7 (4.7)	0.0 (0.0)	6.8 (6.0)	0.0 (0.0)	3.9 (3.2)	0.2 (0.1)	6.6 ( <b>4.6</b> )
GOAL-ONLY	0.1 (0.0)	6.5 (4.3)	0.0 (0.0)	6.8 (5.0)	0.1 (0.1)	5.0 (3.7)	0.2 (0.0)	6.9 (4.4)
INSTRUCTIONS-ONLY	2.3 (1.1)	9.4 (6.1)	0.0 (0.0)	<b>7.0</b> (4.9)	2.7 (1.4)	8.2 (5.5)	<b>0.5 (0.2)</b>	7.2 ( <b>4.6</b> )
SEQ2SEQ	2.4 (1.1)	9.4 (5.7)	<b>0.1</b> (0.0)	6.8 (4.7)	2.1 (1.0)	7.4 (4.7)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM PROGRESS-ONLY	2.1 (1.1)	8.7 (5.6)	0.0 (0.0)	6.9 (5.0)	3.0 (1.7)	8.0 (5.5)	0.3 (0.1)	<b>7.3</b> (4.5)
+ PM SUBGOAL-ONLY	2.1 (1.2)	9.6 (5.5)	0.0 (0.0)	6.6 (4.6)	3.8 (1.7)	8.9 (5.6)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM Both	<b>3.7 (2.1)</b>	<b>10.0 (7.0)</b>	0.0 (0.0)	6.9 ( <b>5.1</b> )	<b>4.0 (2.0)</b>	<b>9.4 (6.3)</b>	0.4 (0.1)	7.0 (4.3)
HUMAN	-	-	-	-	-	-	91.0 (85.8)	94.5 (87.6)

- Providing only high-level, underspecified goal language is insufficient to complete the tasks but is enough to complete some goal-conditions.
- Using just low-level, step-by-step instructions, performs similarly to using both high- and low-levels.

Model	Validation				Test			
	Seen		Unseen		Seen		Unseen	
	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond	Task	Goal-Cond
NO LANGUAGE	0.0 (0.0)	5.9 (3.4)	0.0 (0.0)	6.5 (4.7)	0.2 (0.0)	5.0 (3.2)	0.2 (0.0)	6.6 (4.0)
NO VISION	0.0 (0.0)	5.7 (4.7)	0.0 (0.0)	6.8 (6.0)	0.0 (0.0)	3.9 (3.2)	0.2 (0.1)	6.6 ( <b>4.6</b> )
GOAL-ONLY	0.1 (0.0)	6.5 (4.3)	0.0 (0.0)	6.8 (5.0)	0.1 (0.1)	5.0 (3.7)	0.2 (0.0)	6.9 (4.4)
INSTRUCTIONS-ONLY	2.3 (1.1)	9.4 (6.1)	0.0 (0.0)	<b>7.0</b> (4.9)	2.7 (1.4)	8.2 (5.5)	<b>0.5 (0.2)</b>	7.2 ( <b>4.6</b> )
SEQ2SEQ	2.4 (1.1)	9.4 (5.7)	<b>0.1</b> (0.0)	6.8 (4.7)	2.1 (1.0)	7.4 (4.7)	<b>0.5 (0.2)</b>	7.1 (4.5)
+ PM PROGRESS-ONLY	2.1 (1.1)	8.7 (5.6)	0.0 (0.0)	6.9 (5.0)	3.0 (1.7)	8.0 (5.5)	0.3 (0.1)	<b>7.3</b> (4.5)
+ PM SUBGOAL-ONLY	2.1 (1.2)	9.6 (5.5)	0.0 (0.0)	6.6 (4.6)	3.8 (1.7)	8.9 (5.6)	<b>0.5 (0.2)</b>	7.1 (4.5)
<b>+ PM Both</b>	<b>3.7 (2.1)</b>	<b>10.0 (7.0)</b>	0.0 (0.0)	6.9 ( <b>5.1</b> )	<b>4.0 (2.0)</b>	<b>9.4 (6.3)</b>	0.4 (0.1)	7.0 (4.3)
HUMAN	-	-	-	-	-	-	91.0 (85.8)	94.5 (87.6)

- The two progress monitoring signals are marginally helpful, increasing the success rate by ~1% to ~2%.
- They also lead to more efficient task completion, as indicated by the consistently higher path weighted scores.

Sub-Goal Ablations - **Validation**

Model		<i>Goto</i>	<i>Pickup</i>	<i>Put</i>	<i>Cool</i>	<i>Heat</i>	<i>Clean</i>	<i>Slice</i>	<i>Toggle</i>	Avg.
<i>Seen</i>	No Lang	28	22	71	<b>89</b>	<b>87</b>	64	19	90	59
	S2S	49	32	80	87	85	<b>82</b>	23	97	67
	S2S + PM	<b>51</b>	<b>32</b>	<b>81</b>	88	85	81	<b>25</b>	<b>100</b>	<b>68</b>
<i>Unseen</i>	No Lang	17	9	31	75	86	13	8	4	30
	S2S	21	20	<b>51</b>	<b>94</b>	88	21	<b>14</b>	<b>54</b>	45
	S2S + PM	<b>22</b>	<b>21</b>	46	92	<b>89</b>	<b>57</b>	12	32	<b>46</b>

- Visual semantic navigation (Goto, Pickup) is considerably harder in unseen environments.
- Simple sub-goals like Cool, and Heat are achieved at a high success rate of ~90% because these tasks are mostly object-agnostic.

Rank ↕	Submission	Created ↕	Unseen Success Rate ↕	Seen Success Rate ↕	Seen PLWSR ↕	Unseen PLWSR ↕	Seen GC ↕	Unseen GC ↕	Seen PLW GC Success Rate ↕	Unseen PLW GC Success Rate ↕
1	<b>Hierarchical Attention Model</b> <i>Van-Quang Nguyen and Takayuki...</i>	08/01/2020	0.0445	0.1239	0.0820	0.0224	0.2068	0.1234	0.1879	0.0944
2	<b>Baseline Seq2Seq + Progress M...</b> <i>Singh, Bhambri, Kim, Choi (Gl...</i>	07/22/2020	0.0150	0.0541	0.0251	0.0070	0.1232	0.0808	0.0827	0.0520
3	<b>Baseline + ImprovedMask personal</b>	07/29/2020	0.0066	0.0385	0.0207	0.0032	0.1018	0.0798	0.0757	0.0551
4	<b>baseline</b> <i>DeepblueAI</i>	07/10/2020	0.0059	0.0365	0.0202	0.0019	0.0851	0.0714	0.0533	0.0410
5	<b>Baseline Seq2Seq+PM (both)</b> <i>Shridhar et. al (UW)</i>	03/28/2020	0.0039	0.0398	0.0202	0.0008	0.0942	0.0703	0.0627	0.0426
5	<b>baseline personal</b>	07/05/2020	0.0039	0.0359	0.0161	0.0006	0.0889	0.0690	0.0618	0.0421
7	<b>baseline v2</b> <i>DeepblueAI</i>	07/17/2020	0.0033	0.0150	0.0071	0.0008	0.0605	0.0695	0.0383	0.0401
8	<b>Baseline_v2 Personal</b>	08/04/2020	0.0026	0.0033	0.0007	0.0005	0.0453	0.0690	0.0316	0.0413
9	<b>test personal</b>	08/20/2020	0.0020	0.0027	0.0009	0.0006	0.0463	0.0690	0.0333	0.0458
9	<b>Baseline personal</b>	08/04/2020	0.0020	0.0027	0.0009	0.0006	0.0463	0.0690	0.0333	0.0458

# Conclusion

- Result: A baseline model based on recent embodied vision-and-language tasks performs poorly on ALFRED
- Challenges: long-horizon task planning, visual semantic navigation, object detection, referring expression grounding, and action grounding
- Goal: Shrink the gap between research benchmarks and real-world applications



OCT 13 2020



# GROUNDING LANGUAGE IN PLAY

---

Pierre Sermanet, Corey Lynch

Robotics at Google

**PRESENTER: PRIYANKA MANDIKAL**

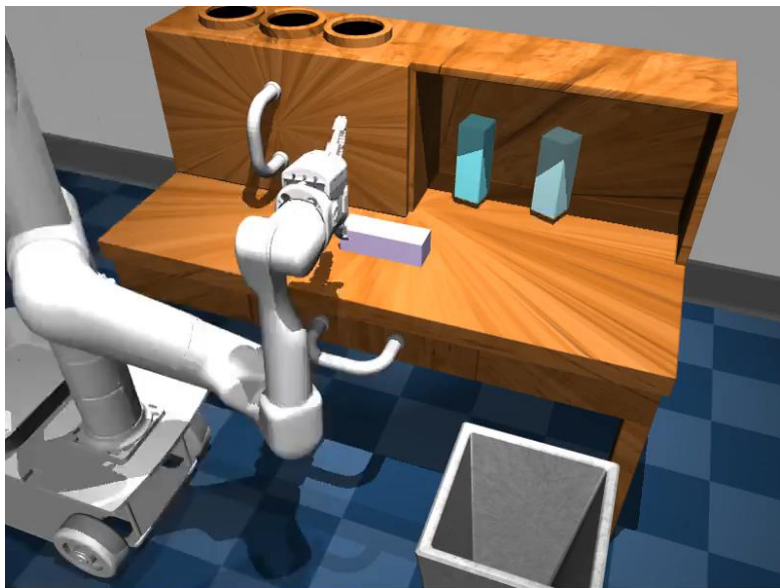
# Motivation

- ❖ Children learn language in the context of rich, sensorimotor experience
  - ➔ language acquisition is embodied
  
- ❖ Infants contribute actions while care-takers contribute relevant words
  - ➔ language acquisition is highly-social

# Main Problem

- ❖ Assuming real humans play a critical role in robot language acquisition, what is the most efficient way to go about it?
- ❖ How can we scalably pair robot experience with relevant human language to bootstrap instruction following?

# Problem Setting



now: **do not do anything**

next:

- ❖ Control a robotic arm within a physics simulator
- ❖ Manipulate objects in the environment
- ❖ Conditioned on an external natural language instruction

# Challenges

- ❖ High-dimensional continuous sensory inputs and actuators
- ❖ Even simple instruction following is notoriously hard
  - E.g. “Sweep the block into the drawer”
    - ❖ Relate language to low-level perception (What does a block look like? What is a drawer?)
    - ❖ Perform visual reasoning (What does it mean for block to be in drawer?)
    - ❖ Solve a complex sequential decision problem (What commands do I send to my arm to “sweep”)
- ❖ Complex task specification → long-horizon robotic object manipulation from natural language instructions

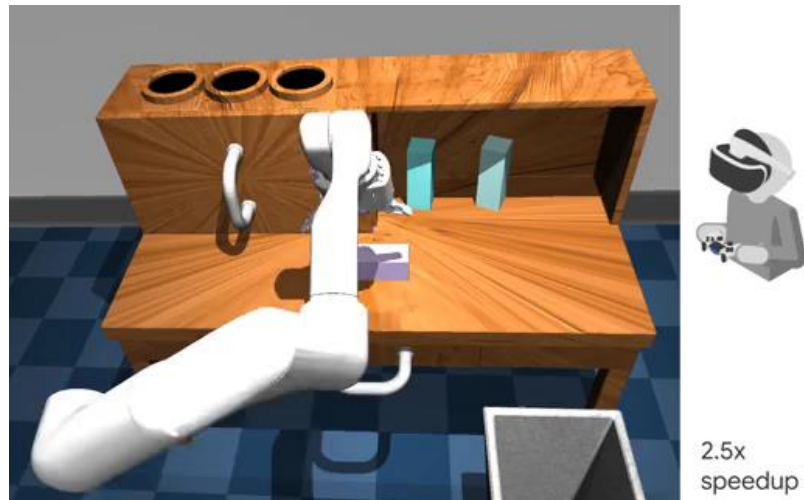
# Related Work

- ❖ Robot learning from general sensors
  - ❖ Imitation learning: requires many human demonstrations
  - ❖ Reinforcement Learning: hand-designed reward functions
- ❖ Task-agnostic control: Single agent must reach any goal on command
  - ❖ Model-based control: Learn model through interactions and then plan; exploration issues
  - ❖ Goal-relabeling: used in both IL and RL ; this paper
- ❖ Covering state space: Exploration vs Tele-operated play
- ❖ Instruction following: restricted env and simplified actuators  
“learning to follow *natural* language is still not the standard in instruction following research” → restricted vocab and grammar

# Prior Work

Learning from Play (LfP): Lynch et. al.  
CORL 2019

- ❖ Learning general-purpose skills from onboard sensors
- ❖ Tele-operated “play” data → relabeled imitation learning → goal-directed policy
- ❖ Limitation: tasks need to be specified using a *goal image* → impractical in real-world environments

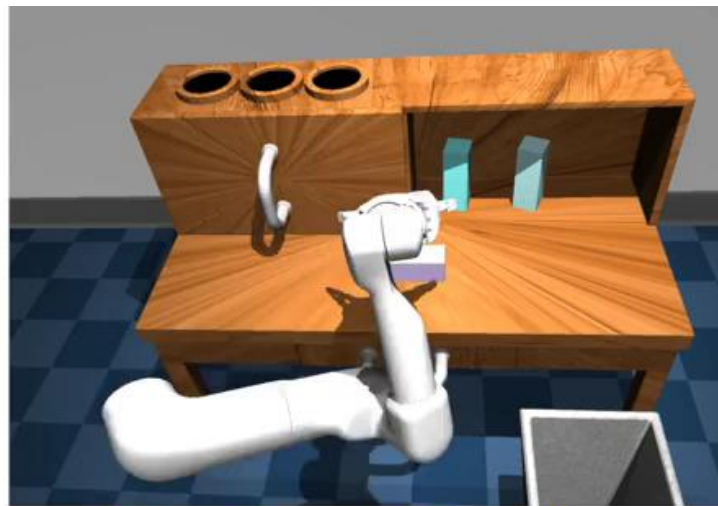


# Prior Work

## Learning from Play



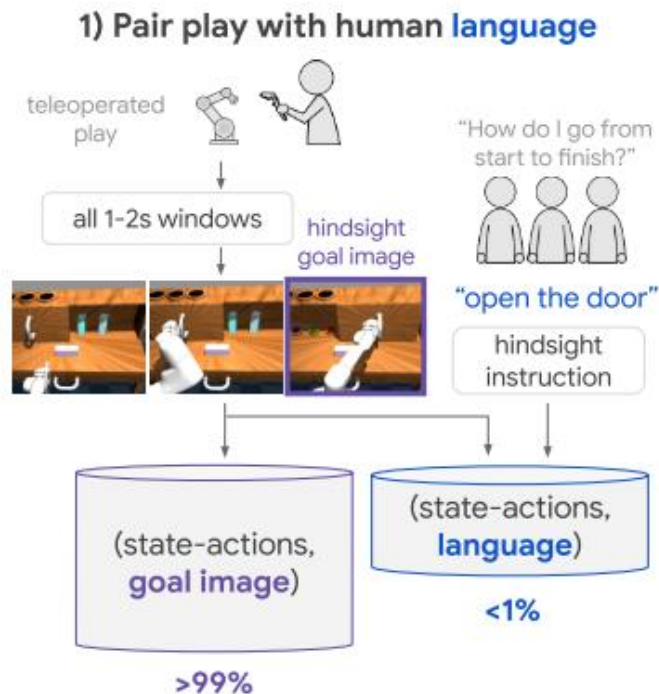
Goal



Single Play-LMP policy



# Overview

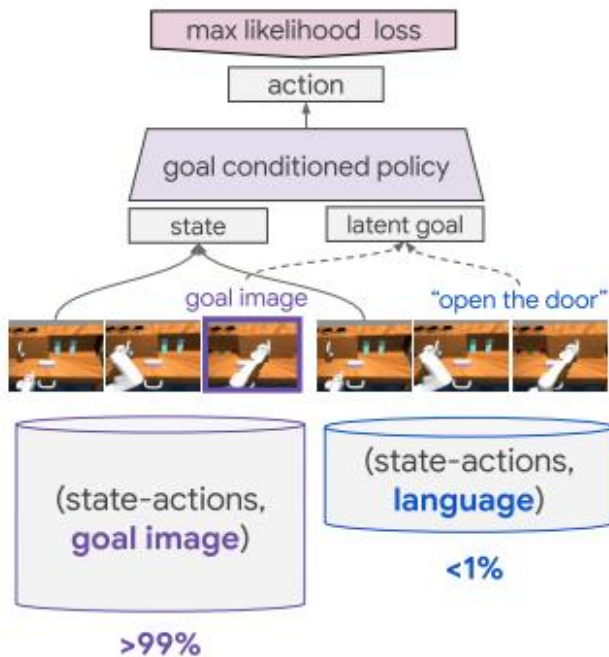


1a. Cover the space with teleoperated play

1b. Pair play with human language (Hindsight Instruction Pairing)

# Overview

## 2) Train on **image** and **language** goals



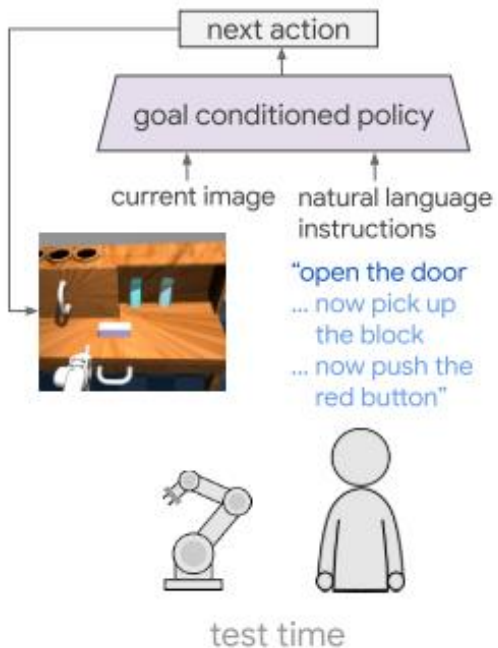
## 2. Multicontext imitation learning

- train a single policy to solve image or language goals

- highly data efficient

# Overview

## 3) Follow human language



3. Condition on human language at test time

# Preliminaries

## ❖ Relabeled Imitation Learning

- Goal conditioned learning – train a single agent to reach any goal
- Goal conditioned behavior cloning → relabel collected data

$$\mathcal{L}_{\text{GCBC}} = \mathbb{E}_{(\tau, s_g) \sim \mathcal{D}_R} \left[ \sum_{t=0}^{|\tau|} \log \pi_{\theta}(a_t | s_t, s_g) \right]$$

## ❖ Teleoperated Play

- adds diversity to the dataset (fully cover state space)

## ❖ Learning from Play (LfP)

- combines relabeled imitation learning with teleoperated play

$$\mathcal{L}_{\text{LfP}} = \mathbb{E}_{(\tau, s_g) \sim D_{\text{play}}} \left[ \sum_{t=0}^{|\tau|} \log \pi_{\theta}(a_t | s_t, s_g) \right]$$

# APPROACH

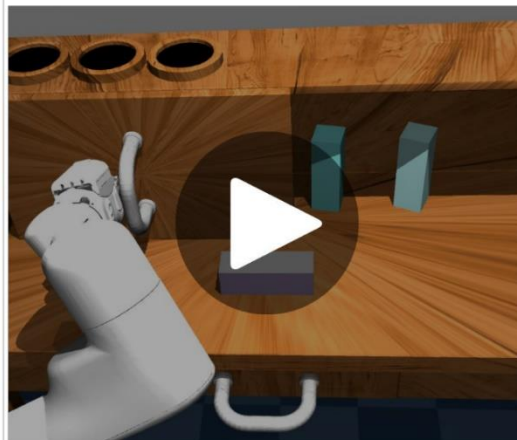
# Hindsight Instruction Pairing

- ❖ Sample any robot behavior from play, then collect an optimal instruction
- ❖ After-the-fact natural lang instructions, operator actions not affected by instructions → more diverse play and instruction dataset
- ❖ No restrictions on vocabulary or grammar

“pick the object and then lift it up.”  
 “pull the drawer.”  
 “drag the object into the drawer”  
 “drop the object, and again pickup the object high”  
 “close the drawer”  
 “do nothing.”

Task: type the instruction that answers:  
 “How do I go from start to finish?”

<type instruction>



video (start to finish)



start frame

finish frame

# Hindsight Instruction Pairing

- Assumes access to  $D_{\text{play}}$  consisting of hindsight goal image samples
- From  $D_{\text{play}} \rightarrow D_{(\text{play}, \text{lang})}$  consisting of hindsight instruction samples

---

**Algorithm 3** Pairing robot sensor data with natural language instructions.

---

- 1: **Input:**  $D_{\text{play}}$ , a relabeled play dataset holding  $(\tau, s_g)$  pairs.
  - 2: **Input:**  $D_{(\text{play}, \text{lang})} \leftarrow \{\}$ .
  - 3: **Input:** `get_hindsight_instruction()`: human overseer, providing after-the-fact natural language instructions for a given  $\tau$ .
  - 4: **Input:**  $K$ , number of pairs to generate,  $K \ll |D_{\text{play}}|$ .
  - 5: **for**  $0 \dots K$  **do**
  - 6:   # Sample random trajectory from play.
  - 7:    $(\tau, -) \sim D_{\text{play}}$
  - 8:   # Ask human for instruction making  $\tau$  optimal.
  - 9:    $l = \text{get\_hindsight\_instruction}(\tau)$
  - 10:   Add  $(\tau, l)$  to  $D_{(\text{play}, \text{lang})}$
  - 11: **end for**
-

# Instruction Samples

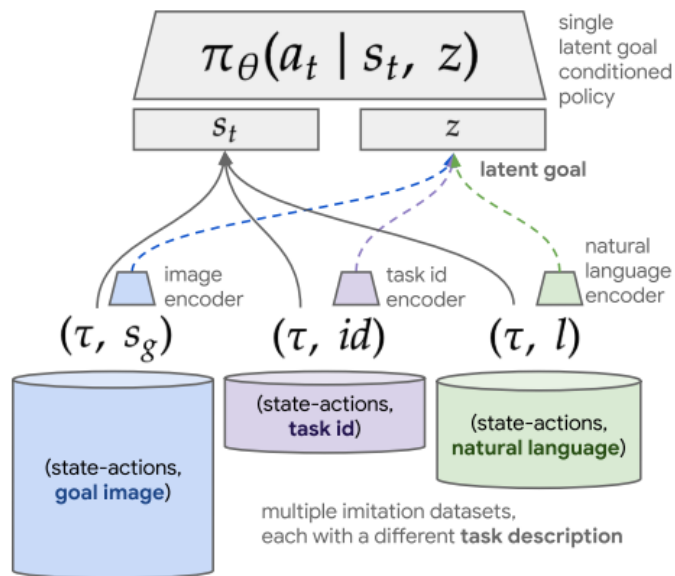
Task	Natural language instructions
open sliding door	“move the door all the way to the right” “slide the door to the right” “move the sliding door all the way to the right and let go”
close sliding door	“Grasp the door handle, then slide the door to the left” “move the door all the way to the left”
open drawer	“open the cabinet drawer” “open the drawer and let go”
close drawer	“close the drawer and let go” “close the drawer”
grasp flat	“Pick up the block” “grasp the object and lift it up” “grasp the object and move your hand up”
grasp lift	“Pick up the object from the drawer and drop it on the table.” “hold the block and place it on top of the table”
grasp upright	“Pick the object and lift it up” “grasp the object and lift”
knock	“push the block forward” “push the object towards the door”
pull out shelf	“Drag the block from the shelf towards the drawer” “pick up the object from the shelf and drop it on the table”

Task	Natural language instructions
put in shelf	“grasp the object and place it inside the cabinet shelf” “Pick the object, move it into the shelf and then drop it.”
push red	“go press the red button” “Press the red button”
push green	“press the green button” “push the green button”
push blue	“push the blue button” “press down on the blue button.”
rotate left	“Pick up the object, rotate it 90 degrees to the left and drop it on the table” “rotate the object 90 degrees to the left”
rotate right	“turn the object to the right” “rotate the block 90 degrees to the right”
sweep	“roll the object into the drawer” “drag the block into the drawer”
sweep left	“Roll the block to the left” “close your fingers and roll the object to the left”
sweep right	“roll the object to the right” “Push the block to the right.”



# Multicontext Imitation Learning (MCIL)

- ❖ Generalization of contextual imitation to multiple heterogenous contexts
- ❖ Multiple imitation learning datasets, each with a different way of describing tasks and different cost of collection
  - ❖ E.g. goal image, task id, natural language, video demonstration, etc
- ❖ Trains
  - ❖ A single *latent goal* conditioned policy  $\pi_{\theta}(a_t | s_t, z)$  over *all* datasets simultaneously
  - ❖ A set of encoders, one per dataset; each maps task description  $\rightarrow$  *shared* latent space



# Multicontext Imitation Learning (MCIL)

## Training Procedure:

- ❖ At each training step, for each dataset:
  - ❖ sample a minibatch of trajectory-context pairs
  - ❖ encode the contexts in the latent space

- ❖ Contextual imitation objective (per dataset)

$$\mathcal{L}_{\text{context}}(D, h) = \mathbb{E}_{(\tau, c) \sim D} \left[ \sum_{t=0}^{|\tau|} \log \pi_{\theta}(a_t | s_t, f_{\theta}(c)) \right]$$

- ❖ Full MCIL objective  $\rightarrow$  averaged over all datasets

$$\mathcal{L}_{\text{MCIL}} = \frac{1}{|\mathcal{D}|} \sum_k^{|\mathcal{D}|} \mathcal{L}_{\text{context}}(D_k, h_k)$$

---

### Algorithm 1 Multicontext imitation learning

---

- 1: **Input:**  $\mathcal{D} = \{D^0, \dots, D^K\}$ ,  $D^k = \{(\tau_i^k, c_i^k)\}_{i=0}^{D^k}$ , One dataset per context type (e.g. goal image, language instruction, task id), each holding pairs of (demonstration, context).
  - 2: **Input:**  $\mathcal{F} = \{f_{\theta}^0, \dots, f_{\theta}^K\}$ , One encoder per context type, mapping context to shared latent goal space, e.g.  $z = f_{\theta}^k(c^k)$ .
  - 3: **Input:**  $\pi_{\theta}(a_t | s_t, z)$ , Single latent goal conditioned policy.
  - 4: **Input:** Randomly initialize parameters  $\theta = \{\theta_{\pi}, \theta_{f^0}, \dots, \theta_{f^K}\}$
  - 5: **while** True **do**
  - 6:    $\mathcal{L}_{\text{MCIL}} \leftarrow 0$
  - 7:   # Loop over datasets.
  - 8:   **for**  $k = 0 \dots K$  **do**
  - 9:     # Sample a (demonstration, context) batch from this dataset.
  - 10:      $(\tau^k, c^k) \sim D^k$
  - 11:     # Encode context in shared latent goal space.
  - 12:      $z = f_{\theta}^k(c^k)$
  - 13:     # Accumulate imitation loss.
  - 14:      $\mathcal{L}_{\text{MCIL}} += \sum_{t=0}^{|\tau^k|} \log \pi_{\theta}(a_t | s_t, z)$
  - 15:   **end for**
  - 16:   # Average gradients over context types.
  - 17:    $\mathcal{L}_{\text{MCIL}} *= \frac{1}{|\mathcal{D}|}$
  - 18:   # Train policy and all encoders end-to-end.
  - 19:   Update  $\theta$  by taking a gradient step w.r.t.  $\mathcal{L}_{\text{MCIL}}$
  - 20: **end while**
-

# Multicontext Imitation Learning (MCIL)

## ❖ Advantages

- ❖ Being context-agnostic → enables highly efficient training
  - ❖ Learn the majority of control from the cheapest data source
  - ❖ Learn general task conditioning from a small number of labelled examples  
E.g. Natural language instructions < 1% of collected robot experience!
- ❖ Broadly useful beyond this paper

# LangLfP

- ❖ Special case of MCIL

- ❖ Two datasets

$$\mathcal{D} = \{D_{\text{play}}, D_{(\text{play}, \text{lang})}\}$$

- ❖ Tasks

- ❖ hindsight goal image
- ❖ hindsight instructions

- ❖ Encoders:  $\mathcal{F} = \{g_{\text{enc}}, s_{\text{enc}}\}$

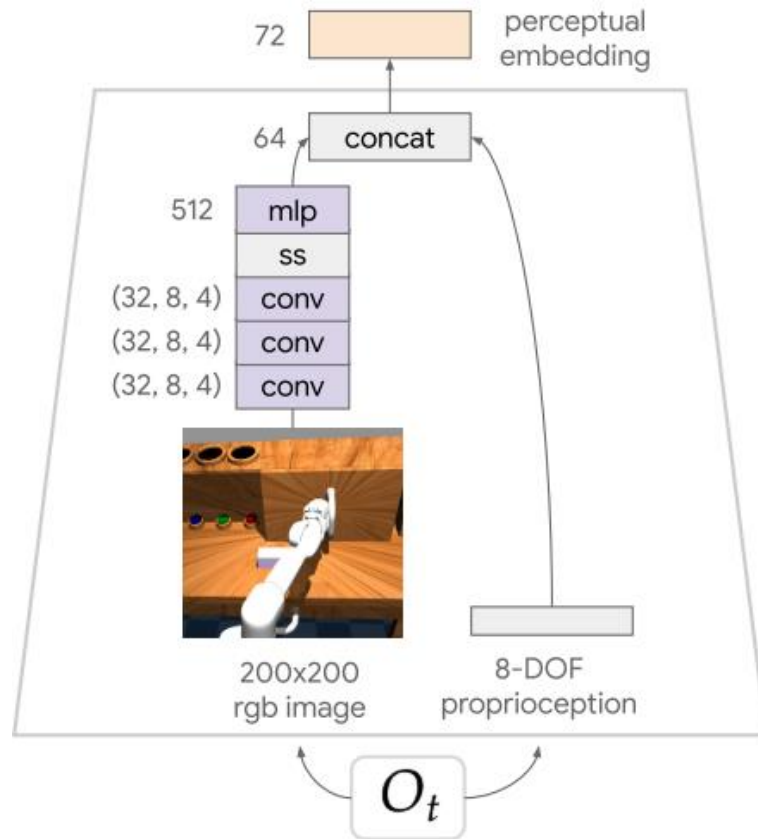
- ❖ Maps image goal and instructions  $\rightarrow$  shared *visuo-lingual* goal space

- ❖ *Learns perception, language understanding and control end-to-end*

# LangLfP: Perception Module

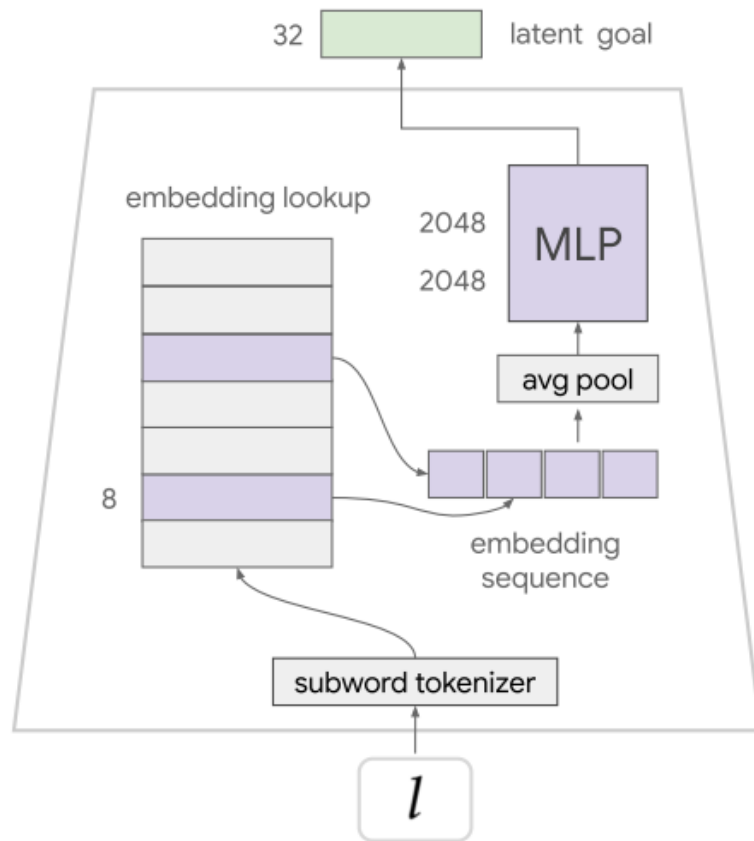
## ❖ Observations:

- ❖ High-dim image (200x200x3)
- ❖ Proprioceptive sensor readings i.e. robot joint angles and locations in Cartesian coordinate space



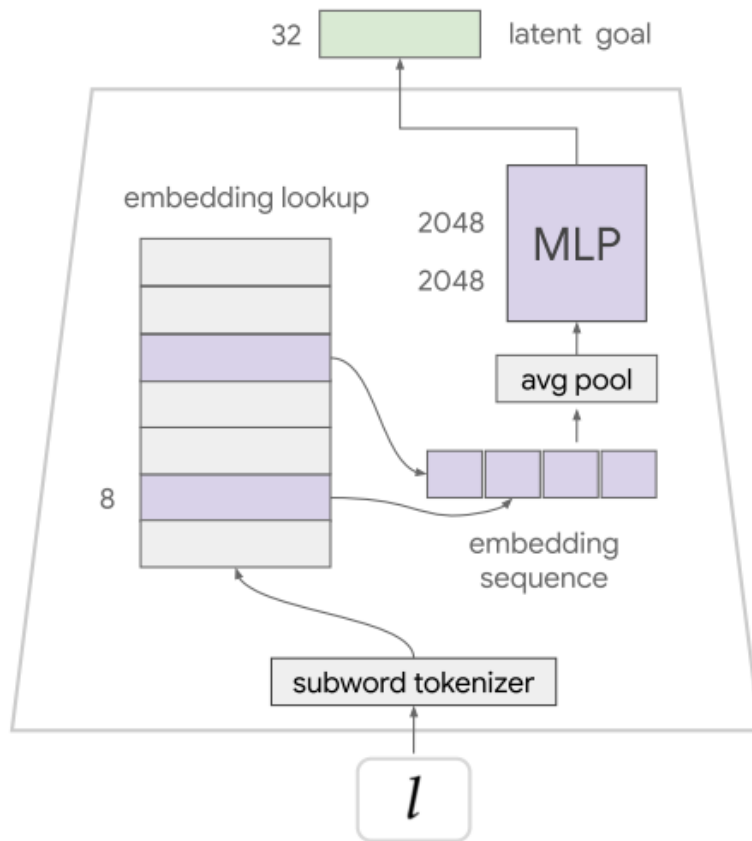
# LangLfP: Language Module

- ❖ Two approaches:
  - ❖ From scratch (LangLfP)
  - ❖ Transfer Learning (TransferLangLfP)



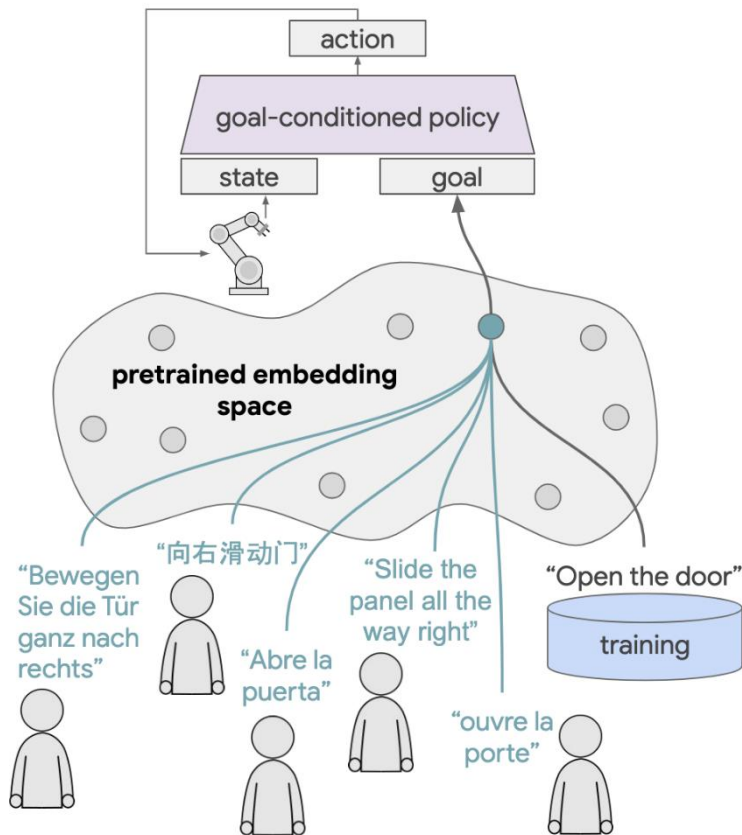
## LangLfP: Language Module from Scratch

- ❖ Tokenize raw text into subwords
- ❖ Retrieve subword embeddings from a lookup table
- ❖ Summarize embeddings into a point in  $z$  space
- ❖ Embedding fed to 2-layer MLP



# TransferLangLfP: Language Module via Transfer Learning

- ❖ Pretrained embeddings from Multilingual Universal Sentence Encoder (MUSE)
- ❖ Maps sentences  $\rightarrow$  512-D vector
- ❖ Benefits
  - ❖ Serves as a strong prior if there is a semantic match between source and target domains
  - ❖ Encodes word similarity  $\rightarrow$  follow out-of-distribution instructions in zero shot

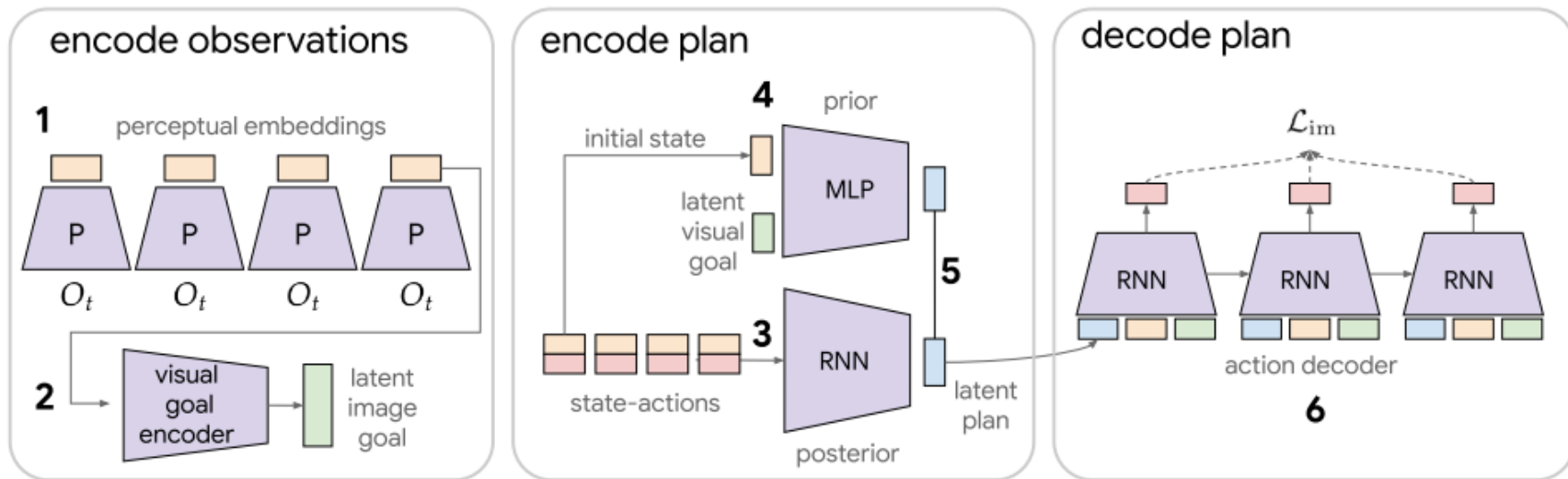




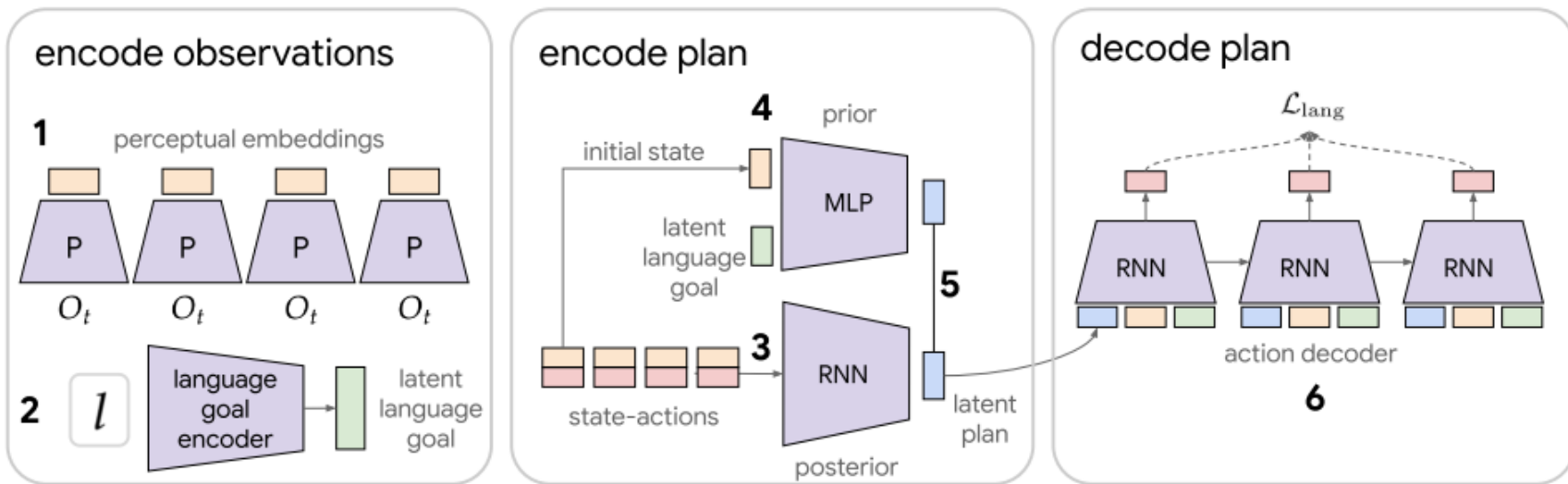
# LangLfP: Control Module

- ❖ Implement multicontext control policy:  $\pi_{\theta}(a_t|s_t, z)$
- ❖ Use Latent Motor Plans (from LfP paper)
  - ❖ Goal directed imitation architecture
  - ❖ Uses latent variables to model multimodality
  - ❖ Seq2seq CVAE that auto-encodes contextual demos through a latent “plan” space
  - ❖ Decoder: goal-conditioned policy
  - ❖ Refer to LfP for more details

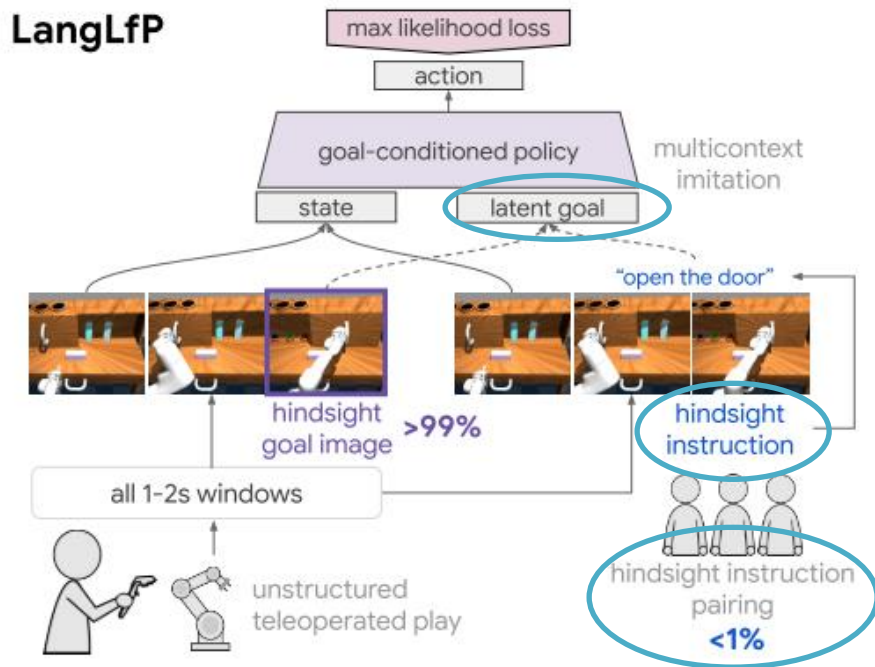
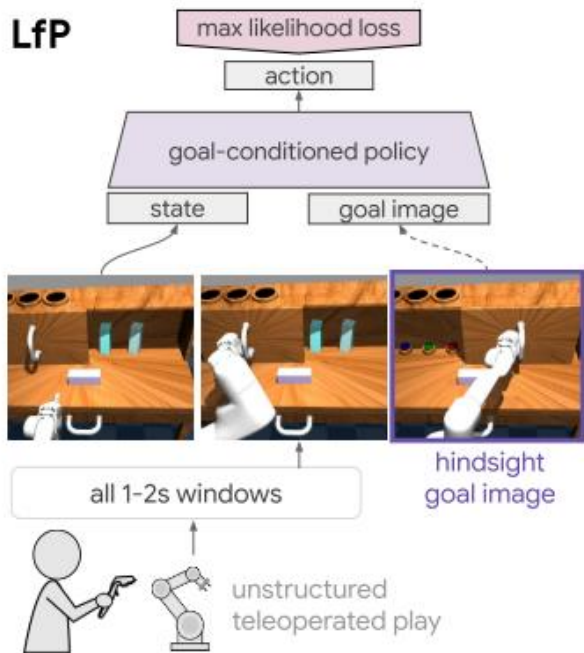
# Multicontext LMP: Goal Image



# Multicontext LMP: Language

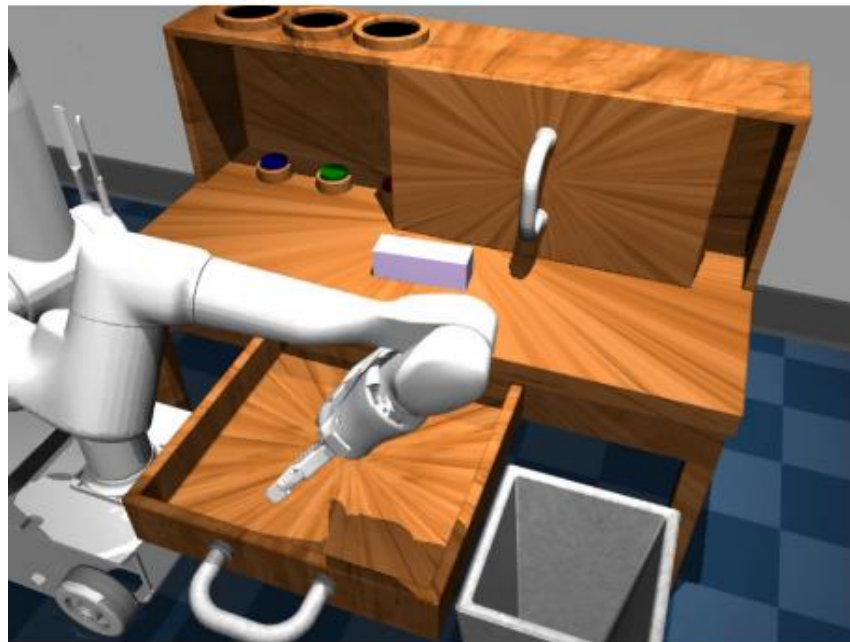


# LfP vs LangLfP



# Experimental Setup: Environment

- ❖ Situated robot in a 3D environment
- ❖ 8-DOF robot arm and parallel gripper
- ❖ RGB video sensors
- ❖ Proprioceptive sensors
- ❖ Goal: Agent must perform high-frequency, closed-loop continuous control to solve user-described manipulation tasks



# Experimental Setup: Methods

- ❖ LangBC – language, but no play – multi-task demos –  $D_{(\text{demo}, \text{lang})}$
  - ❖ LfP – play, but no language –  $D_{\text{play}}$
  - ❖ LangLfP – play and language –  $D_{\text{play}}$  and  $D_{(\text{play}, \text{lang})}$
  - ❖ Restricted LangLfP – LangLfP restricted to size of  $D_{(\text{demo}, \text{lang})}$
  - ❖ TransferLangLfP – LangLfP using MUSE embeddings -  $D_{(\text{play}, \text{lang})}$
- 2 sets of experiments – pixel and state

# Experiments: Ask-Me-Anything (AMA)

- ❖ Multi-stage instruction following
- ❖ Derived from Multi-18 → 18 evaluation tasks described in LfP. E.g. open sliding door, sweep, close sliding door, etc
- ❖ Consider all valid N-stage transitions between the 18 tasks → Chain-2, Chain-3, Chain-4 manipulation benchmarks
- ❖ Multi-18 can be seen as a subset of this extended set

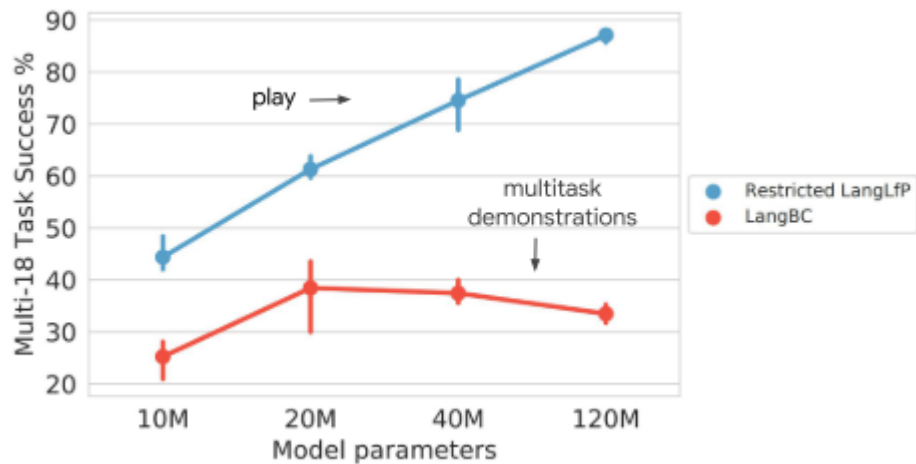
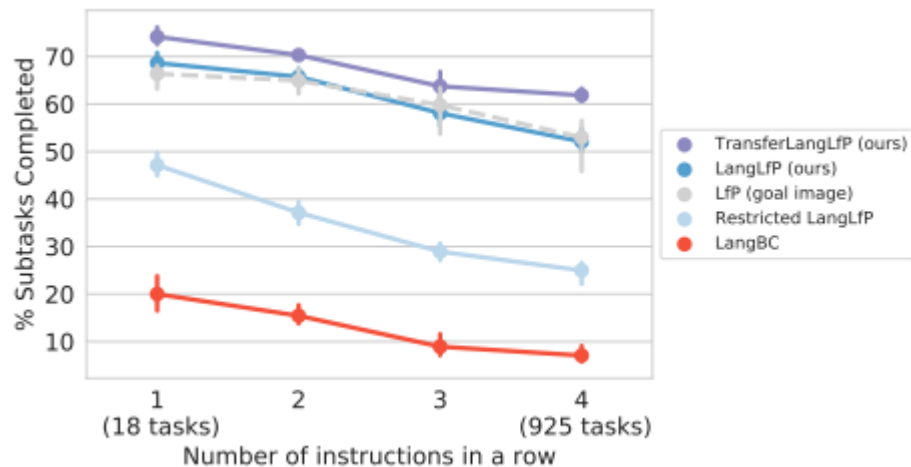
# Experiments: Ask-Me-Anything (AMA)

Method	Input	Training source	Task conditioning	Multi-18 Success (18 tasks)	Chain-4 Success (925 long-horizon tasks)
LangBC	pixels	predefined demos	text	20.0% $\pm$ 3.0	7.1% $\pm$ 1.5
Restricted LangLfP	pixels	play	text	47.1% $\pm$ 2.0	25.0% $\pm$ 2.0
LfP	pixels	play	goal image	66.4% $\pm$ 2.2	53.0% $\pm$ 5.0
LangLfP (ours)	pixels	play	text	68.6% $\pm$ 1.7	52.1% $\pm$ 2.0
TransferLangLfP (ours)	pixels	<b>play</b>	<b>text</b>	<b>74.1%</b> $\pm$ 1.5	<b>68.6%</b> $\pm$ 1.6
LangBC	states	predefined demos	text	38.5% $\pm$ 6.3	13.9% $\pm$ 1.4
Restricted LangLfP	states	play	text	88.0% $\pm$ 1.4	64.2% $\pm$ 1.5
LangLfP (ours)	states	play	text	88.5% $\pm$ 2.9	63.2% $\pm$ 0.9
TransferLangLfP (ours)	states	<b>play</b>	<b>text</b>	<b>90.5%</b> $\pm$ 0.8	<b>71.8%</b> $\pm$ 1.6

- ❖ LangLfP ~ LfP, but is more scalable in terms of task conditioning
- ❖ TransferLangLfP > LangLfP and original LfP
- ❖ RestrictedLangLfP > LangBC; Restricted LangLfP can transition well between tasks; LangBC fails to recover from compounding errors



# Experiments: Ask-Me-Anything (AMA)



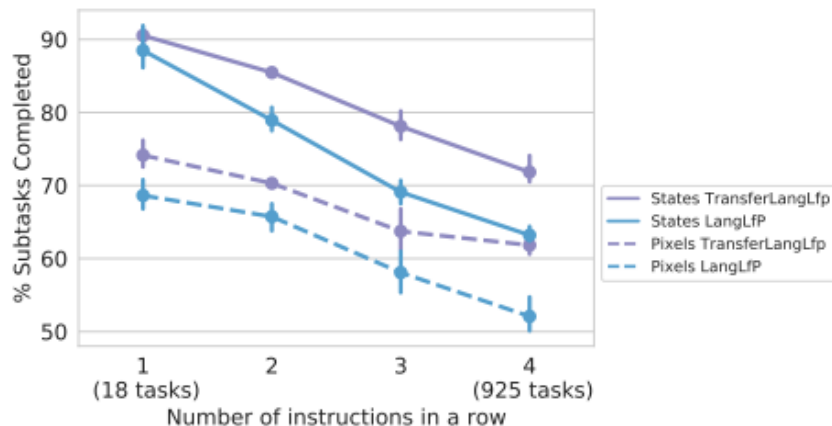
- ❖ As model capacity increases, play model can capitalize on increased strength because of diversity in dataset
- ❖ LangBC constrained to predefined behaviors

# Results: Ask-Me-Anything (AMA)



# Experiments: Knowledge Transfer

- ❖ TransferLangLfP outperforms LangLfP
- ❖ → evidence that world knowledge in large corpora is beneficial for downstream robotic manipulation tasks



# Experiments: Knowledge Transfer

Out-of-distribution instructions:

## ❖ Synonyms

- ❖ “*Drag the block from the shelf*” → “*Retrieve the brick from the cupboard*”
- ❖ OOD-syn eval set: 14k OOD samples across 18 tasks
- ❖ TransferLangLfP generalizes *substantially!*

## ❖ 16 different languages

- ❖ OOD-16-lang eval set: Translate (Multi-18 + OOD-syn) 240k samples across 18 tasks

Method	OOD-syn (~15k tasks)	OOD-16-lang (~240k tasks)
Random Policy	0.0% ± 0.0	0.0% ± 0.0
LangLfP	37.6% ± 2.3	27.94% ± 3.5
TransferLangLfP	<b>60.2% ± 3.2</b>	<b>56.0% ± 1.4</b>

- ❖ TransferLangLfP > LangLfP
- ❖ LangLfP resorts to producing max likelihood play actions

# Results: Knowledge Transfer



now: **trascina fuori il blocco**

**(it → en: drag the block out)**

next: **tirare la maniglia del cassetto fino in fondo**

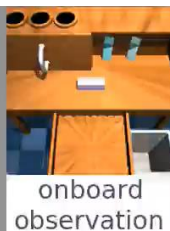
# Results: TransferLangLfP vs LangLfP

## TransferLangLfP



now: **ferme le tiroir entièrement**  
(fr → en: close the drawer all the way)  
next: **tire la poignée du tiroir jusqu'au bout**

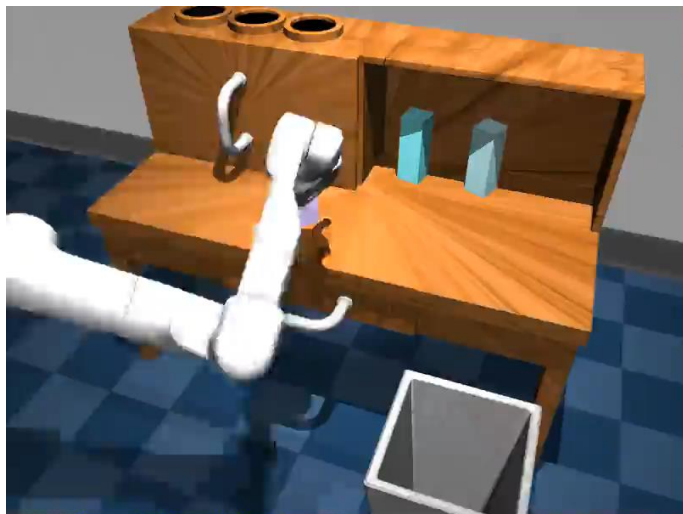
## LangLfP



now: **ferme le tiroir entièrement**  
(fr → en: close the drawer all the way)  
next: **tire la poignée du tiroir jusqu'au bout**

# Limitations

Agent times out before task completion



now: **knock the object**

next: pick up the object and hold it up high



Compounding error → awkward arm configurations



now: **push the object into the cabinet**

next: push the door to the right



# Future Work

❖ Current → goal-directed imitation, lacks autonomous policy improvement

Future → Imitation + RL for autonomous policy improvement, not restricted to human actions

❖ Current → Single env

Future → Large play corpora, generalization to new rooms and objects



# Summary

- ❖ Introduced LangLfP, an extension of LfP trained both on relabeled goal image play and play paired with human language instructions
- ❖ Multicontext Imitation Learning → reduce the cost of language pairing
- ❖ Single policy trained with LangLfP can solve many 3D robotic manipulation tasks over a long horizon from onboard sensors via human language
- ❖ Simple technique for knowledge transfer; 16 different languages

# Discussion

- ❖ Transfer learning on 16 languages
  - ❖ For LangLfP, could have translated instructions to English before feeding into the model
- ❖ How do ALFRED and LangLfP compare with each other?
  - ❖ ALFRED
    - ❖ Pros: mobile robot, larger env diversity, large # of obj state changes (door open/close, lights on/off, bread whole/sliced, tap on/off, vase intact/broken ...) → rich lang vocabulary
    - ❖ Cons: No physical realism, no fine motor control
  - ❖ LangLfP
    - ❖ Pros: contact-rich physics env, 8-DOF motor control
    - ❖ Cons: Fixed robot, limited state space, limited obj state changes (restricted to pick-up, door open/close, lights on/off) → limited action vocabulary