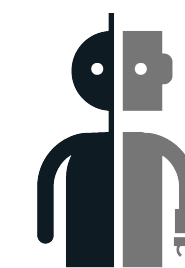


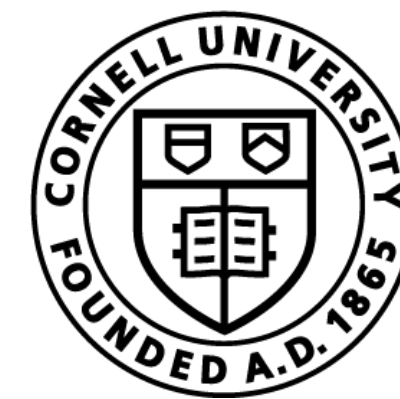


# ICML

International Conference  
On Machine Learning



PORTAL



# Imitation Learning from a Single Temporally Misaligned Video



William Huey\*, Huaxiaoyue (Yuki) Wang\*, Anne Wu, Yoav Artzi, Sanjiban Choudhury

# Designing reward function is tedious



Robot state space?

Environment  
constraints?

Unintended behaviors  
to avoid?

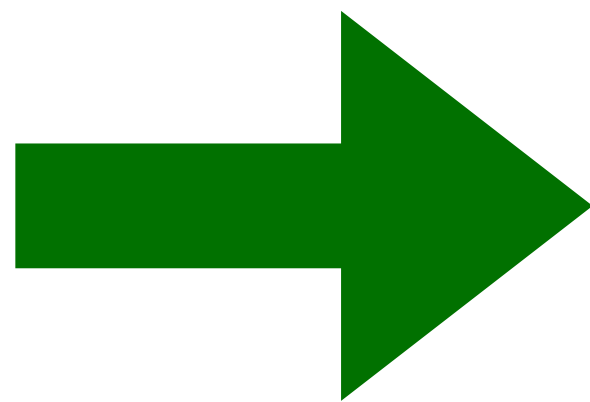
...

# Designing reward function is tedious

## Instead, easier to provide a video demonstrating the task



Robot state space?  
Environment  
constraints?  
Unintended behaviors  
to avoid?  
...





# Demonstrations are often **temporally misaligned**

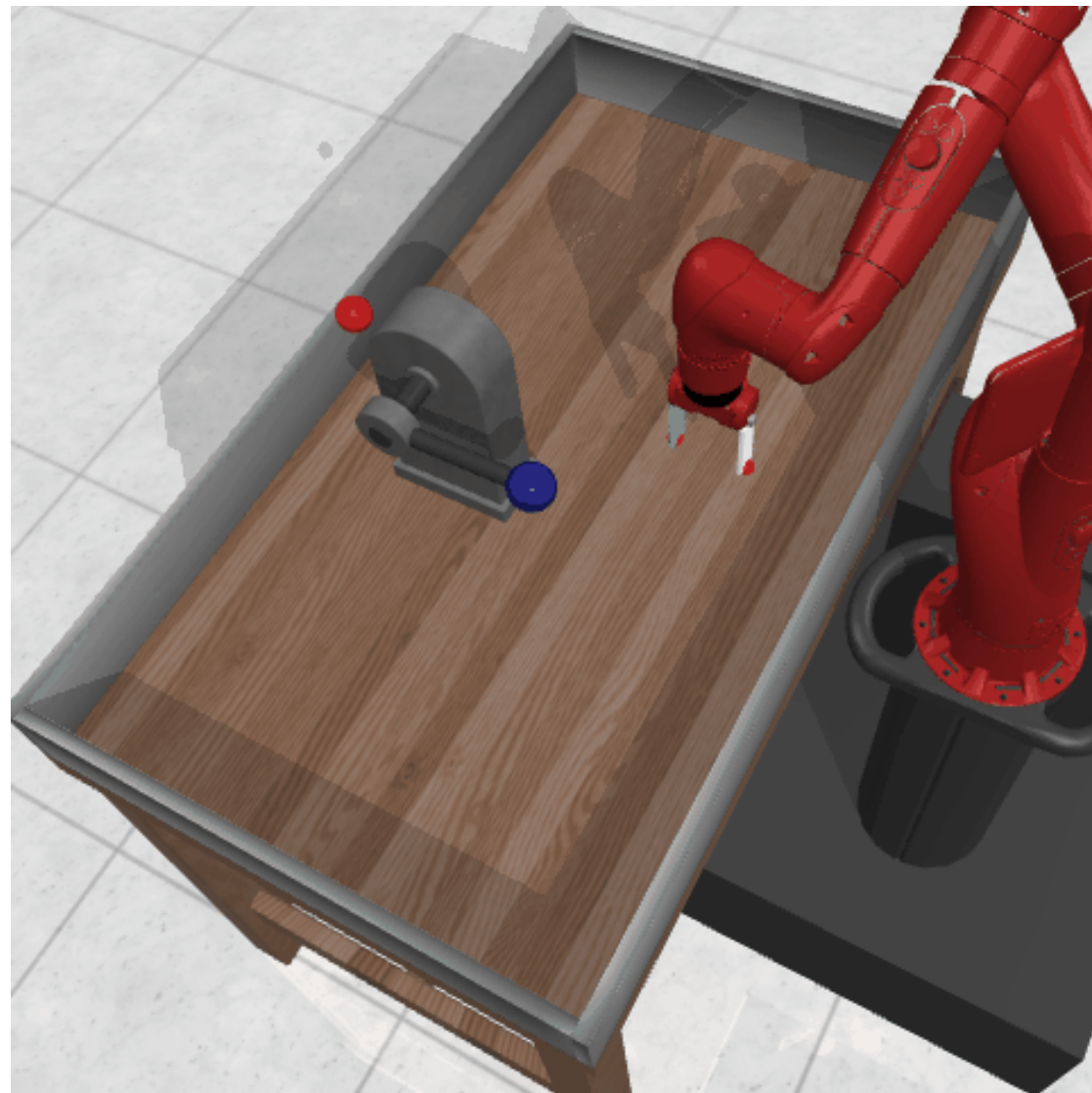


**Different Speeds**

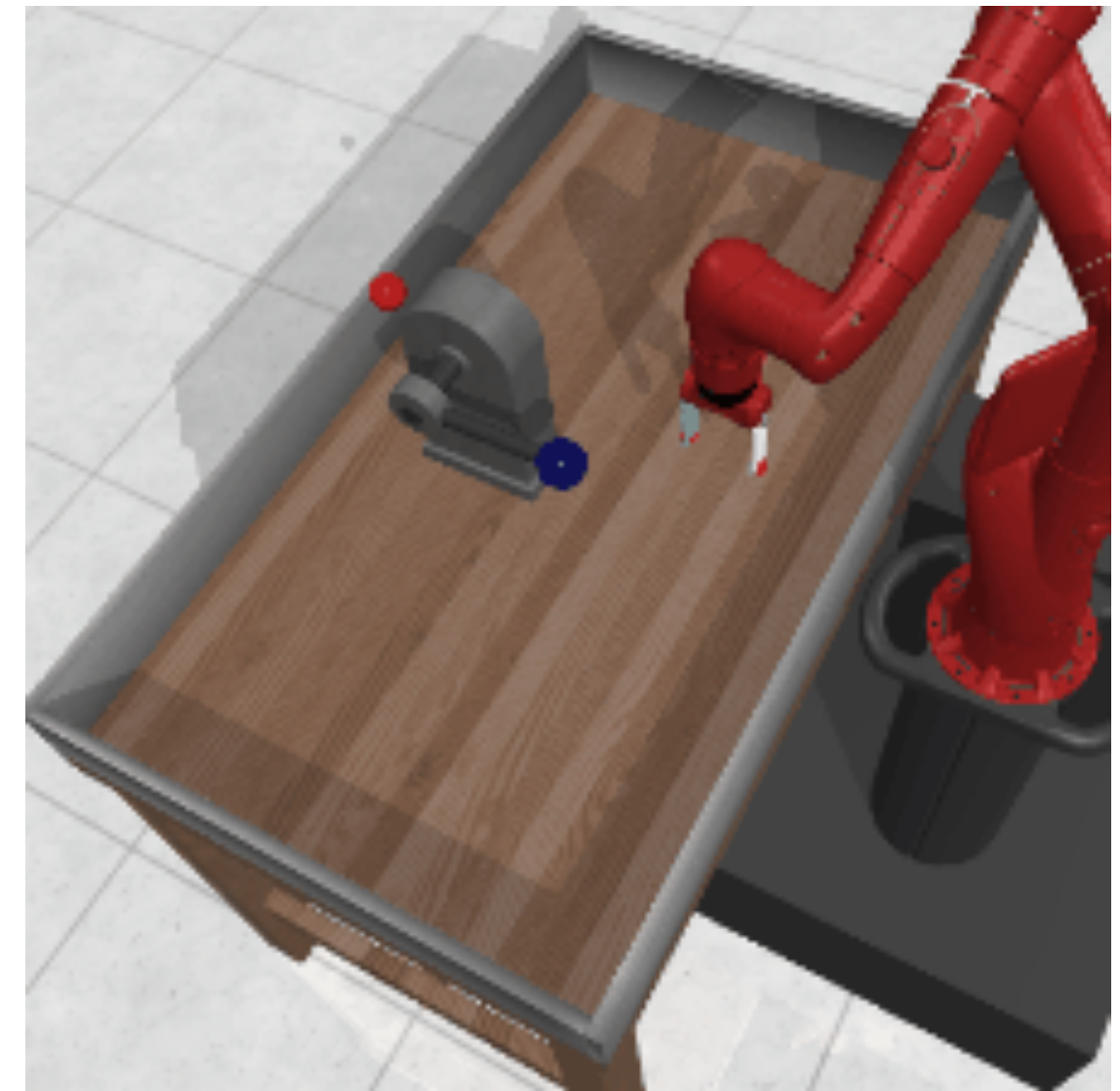




# Demonstrations are often **temporally misaligned**

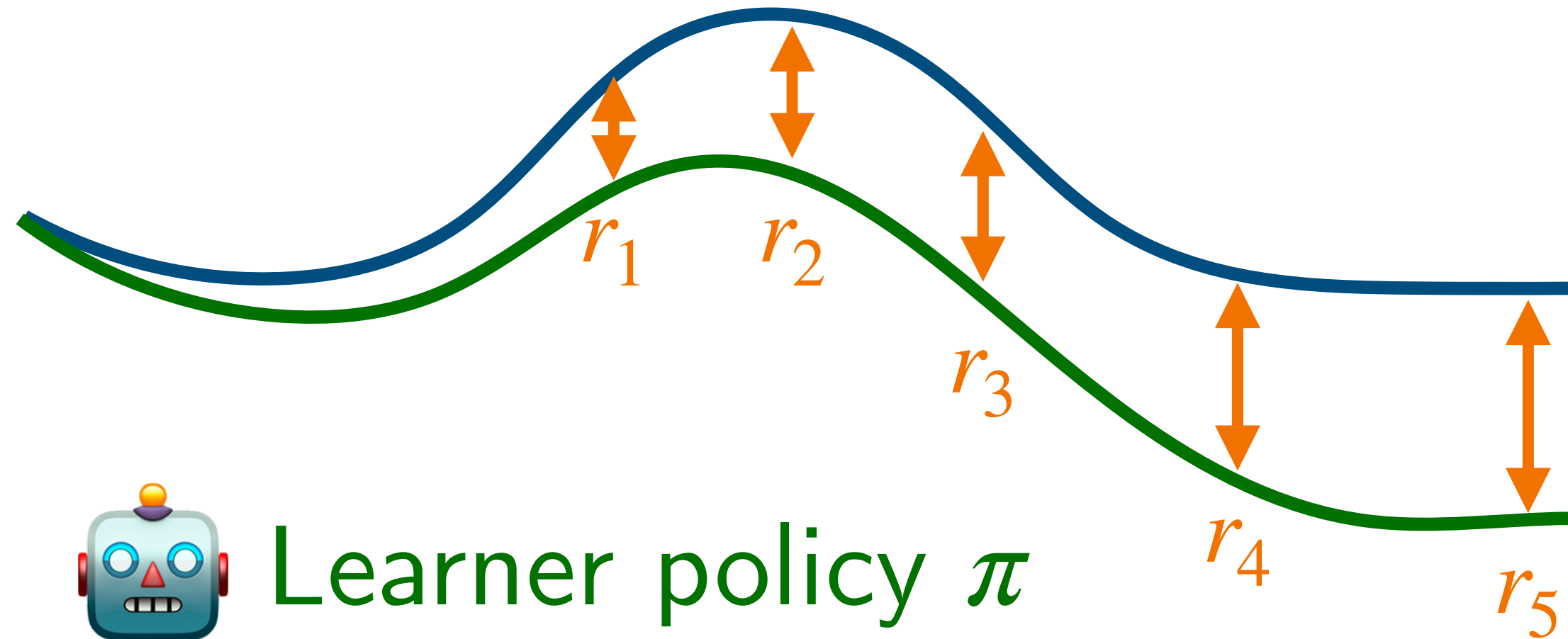


**Different Speeds**



Existing approaches to learning from videos  
use **inverse reinforcement learning** (IRL)

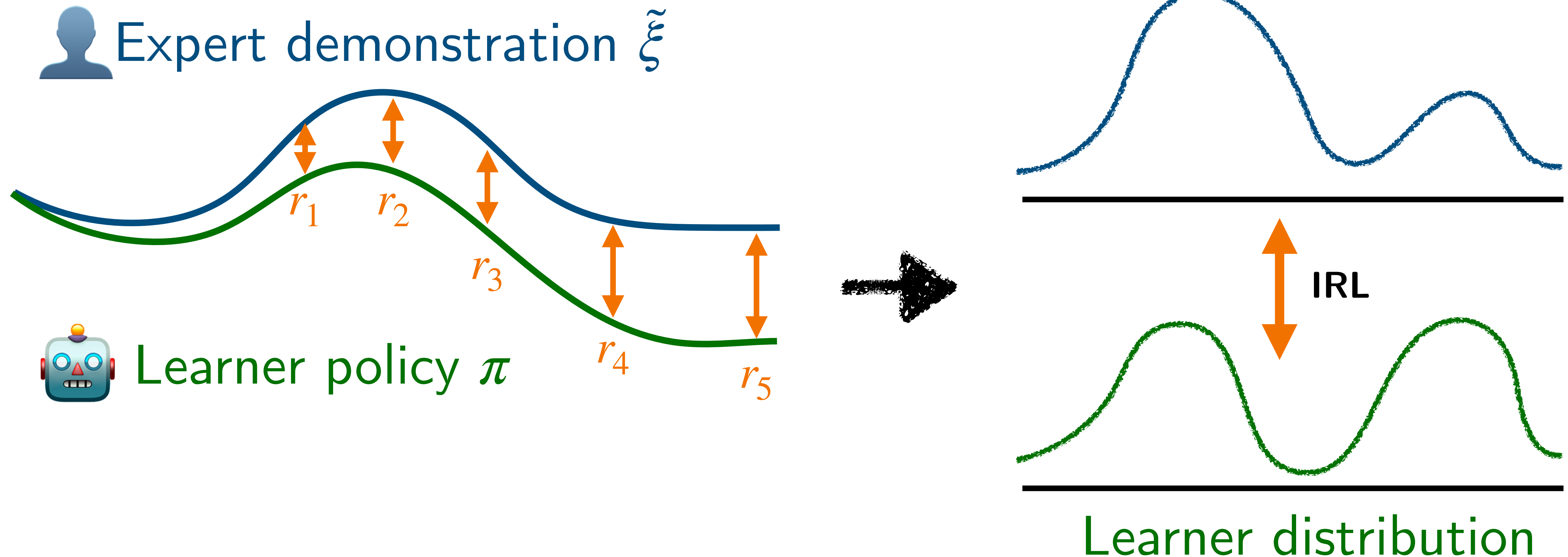
 Expert demonstration  $\xi$



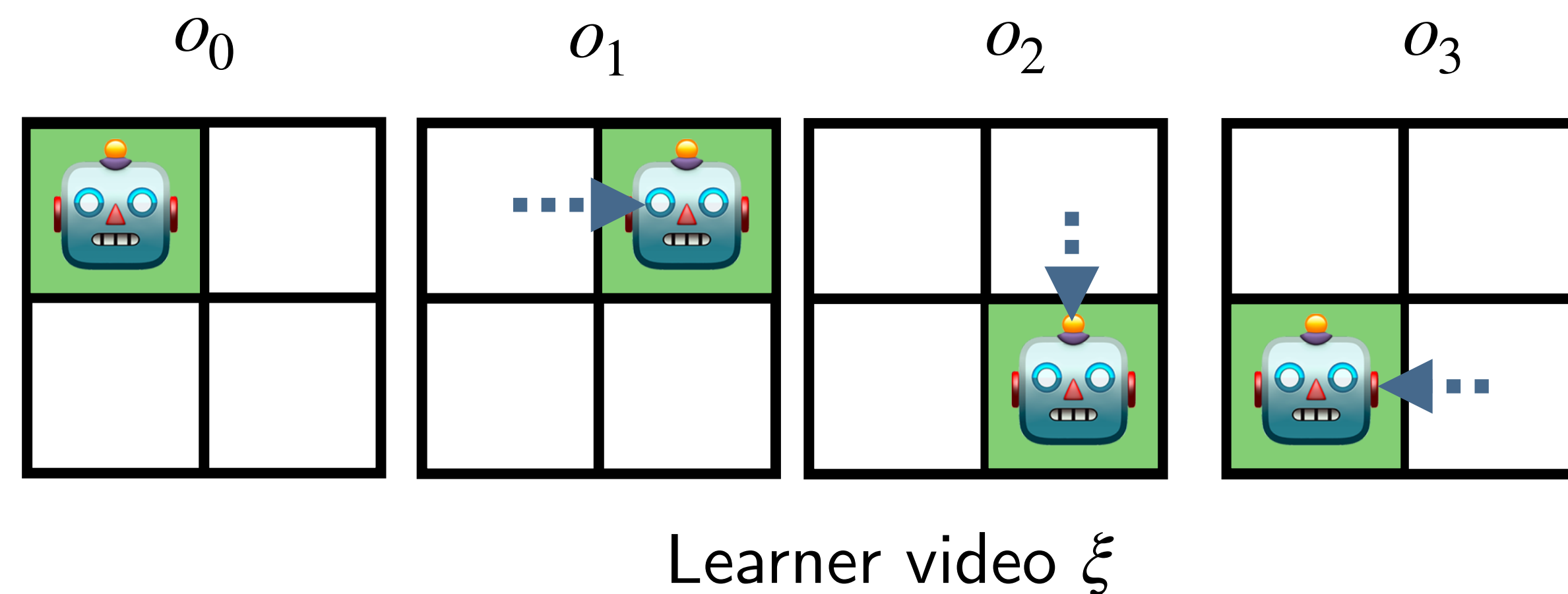
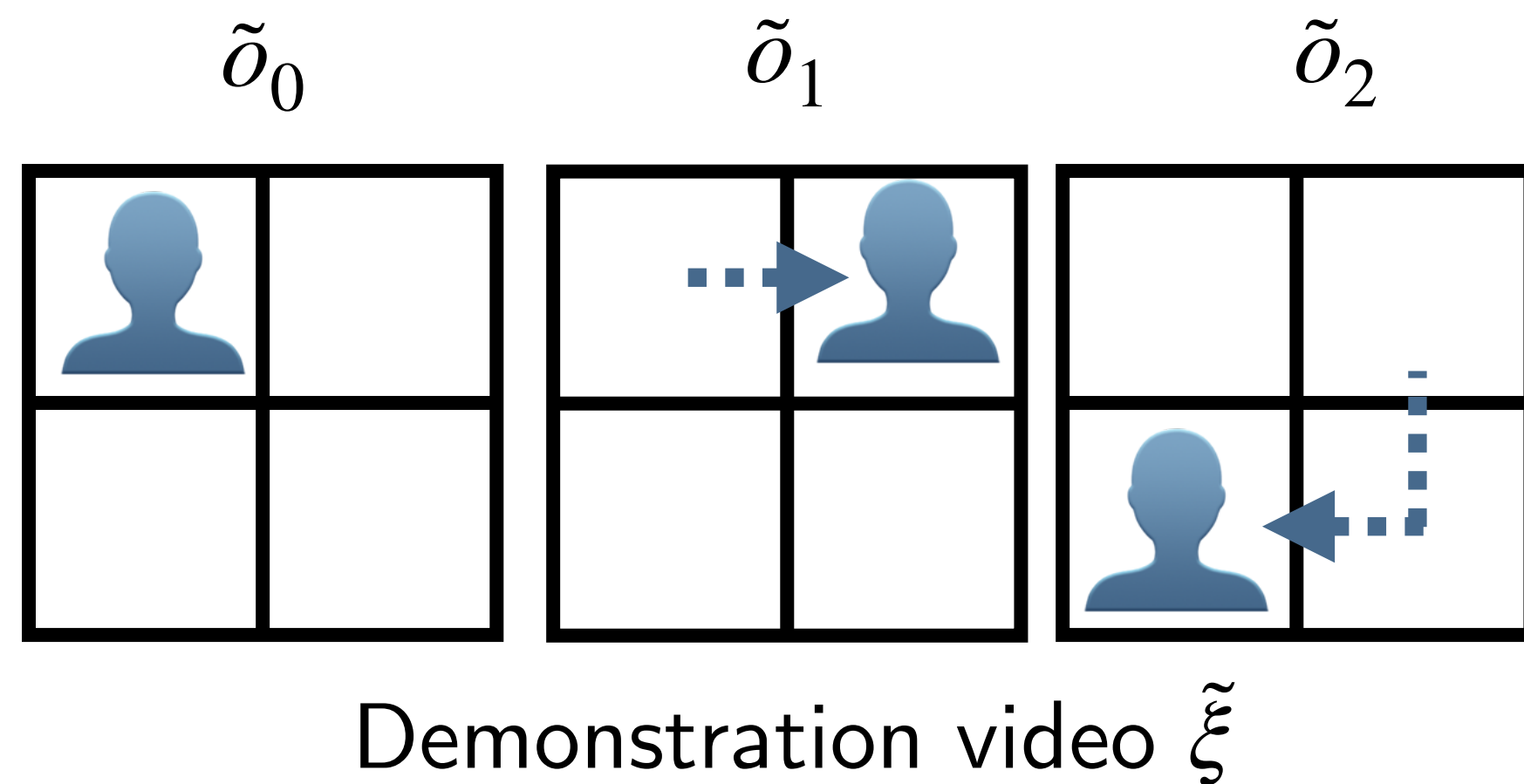
 Learner policy  $\pi$



# IRL matches learner and expert distributions in expectation

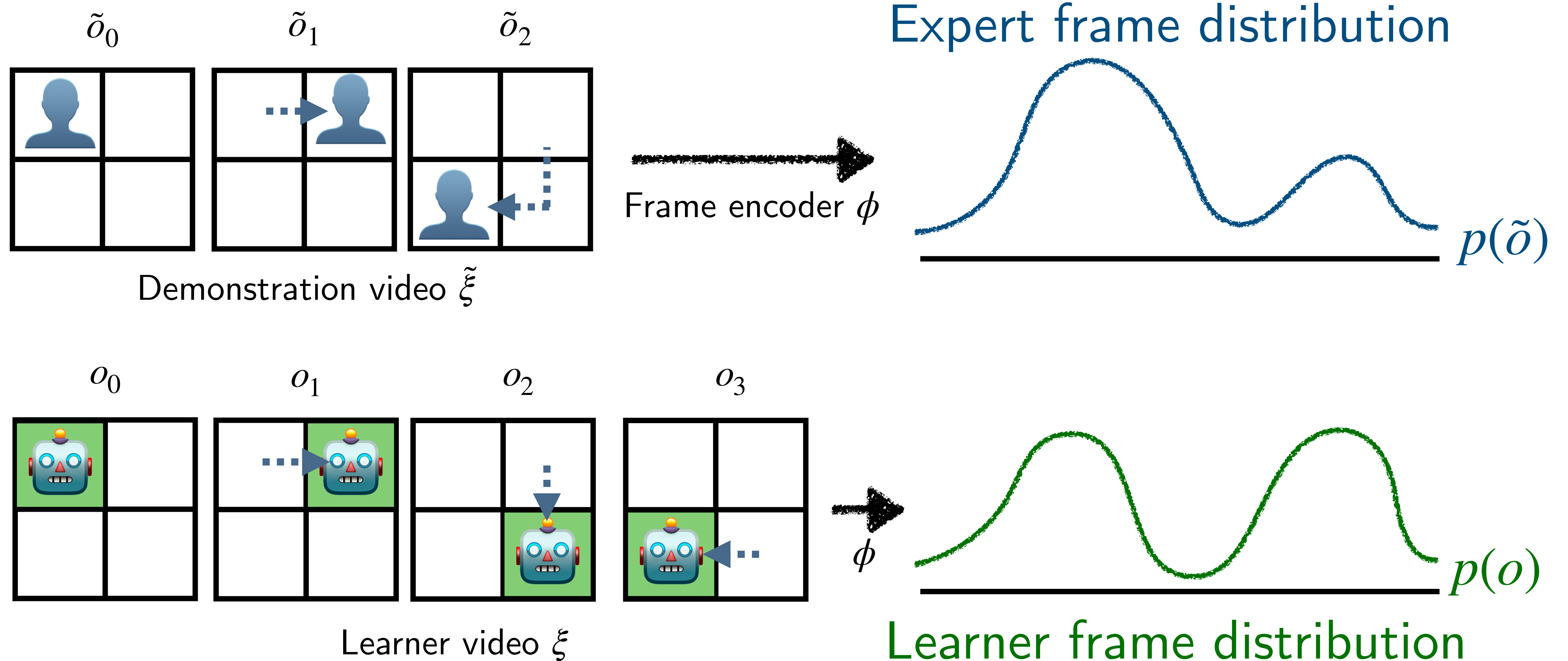


# Existing approaches to learning from videos match **distributions over video frames**

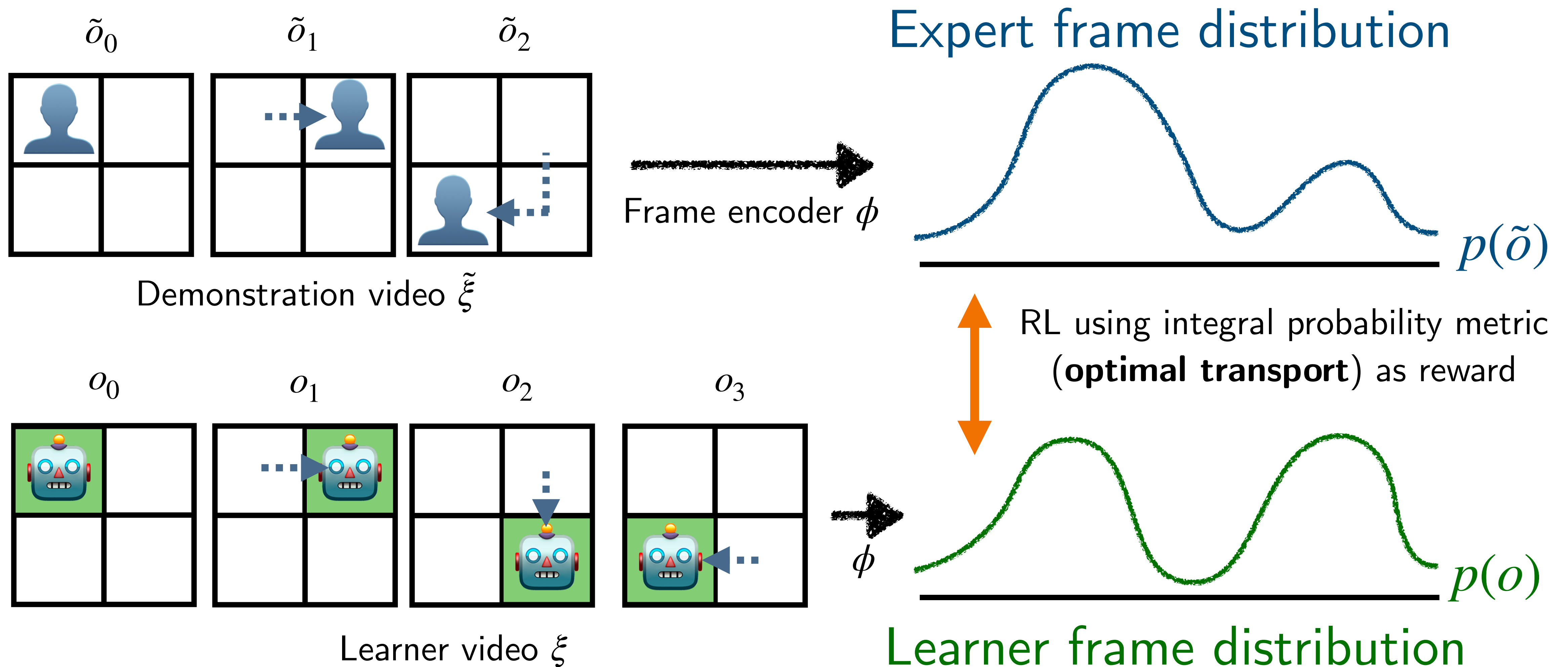




# Existing approaches to learning from videos match **distributions over video frames**

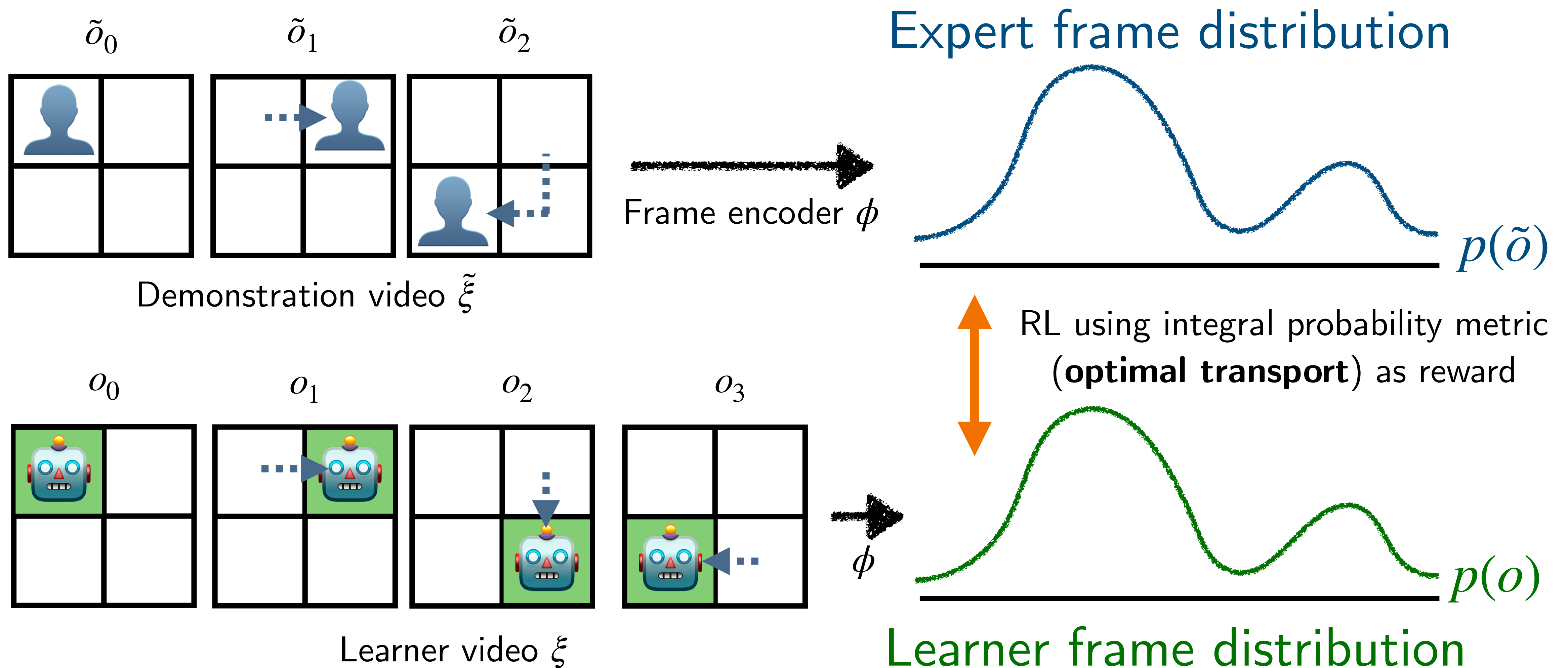


# Existing approaches to learning from videos match **distributions over video frames**

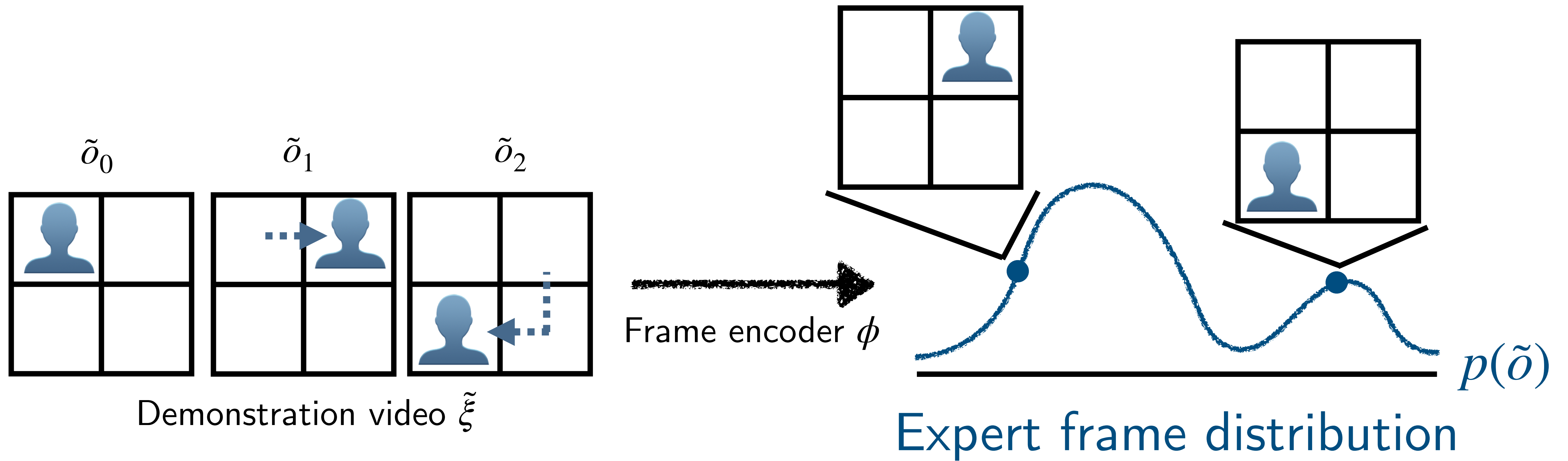




# Existing approaches to learning from videos match **distributions over video frames**



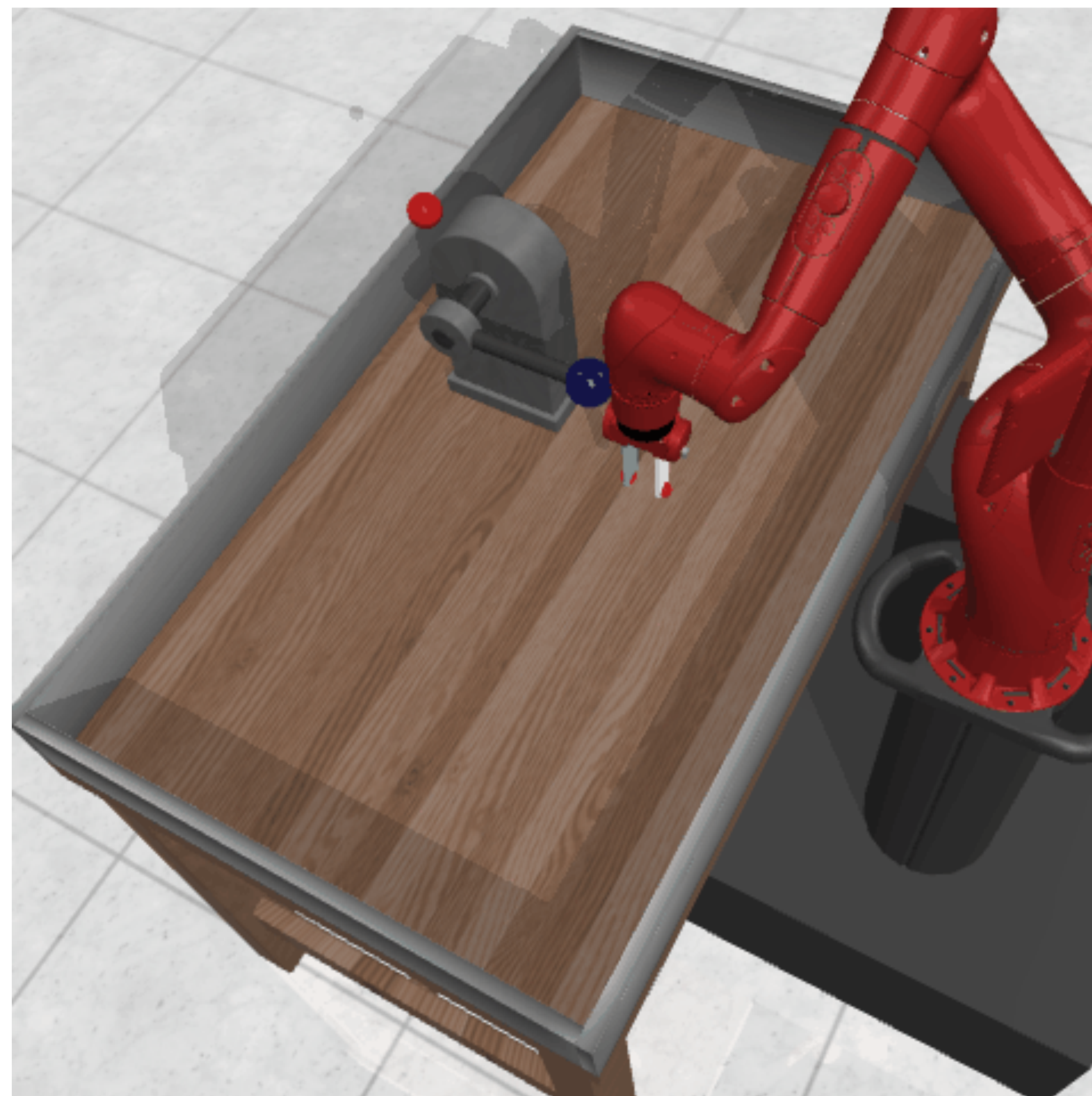
# Problem! Individual frames lack temporal information



$p(\tilde{o})$  represents a collection of frames, not an ordered sequence

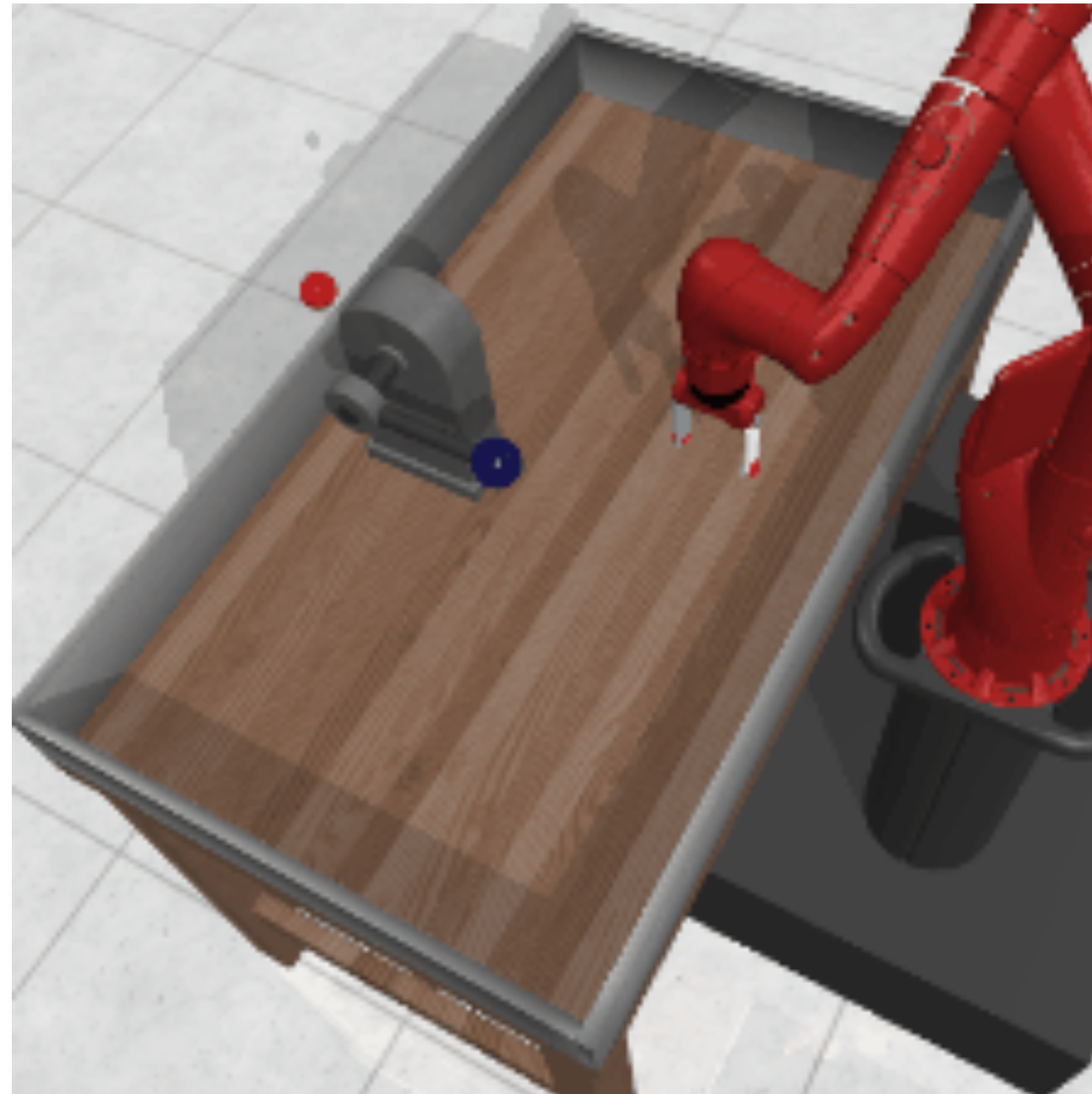


# Experimental Example in Meta-world



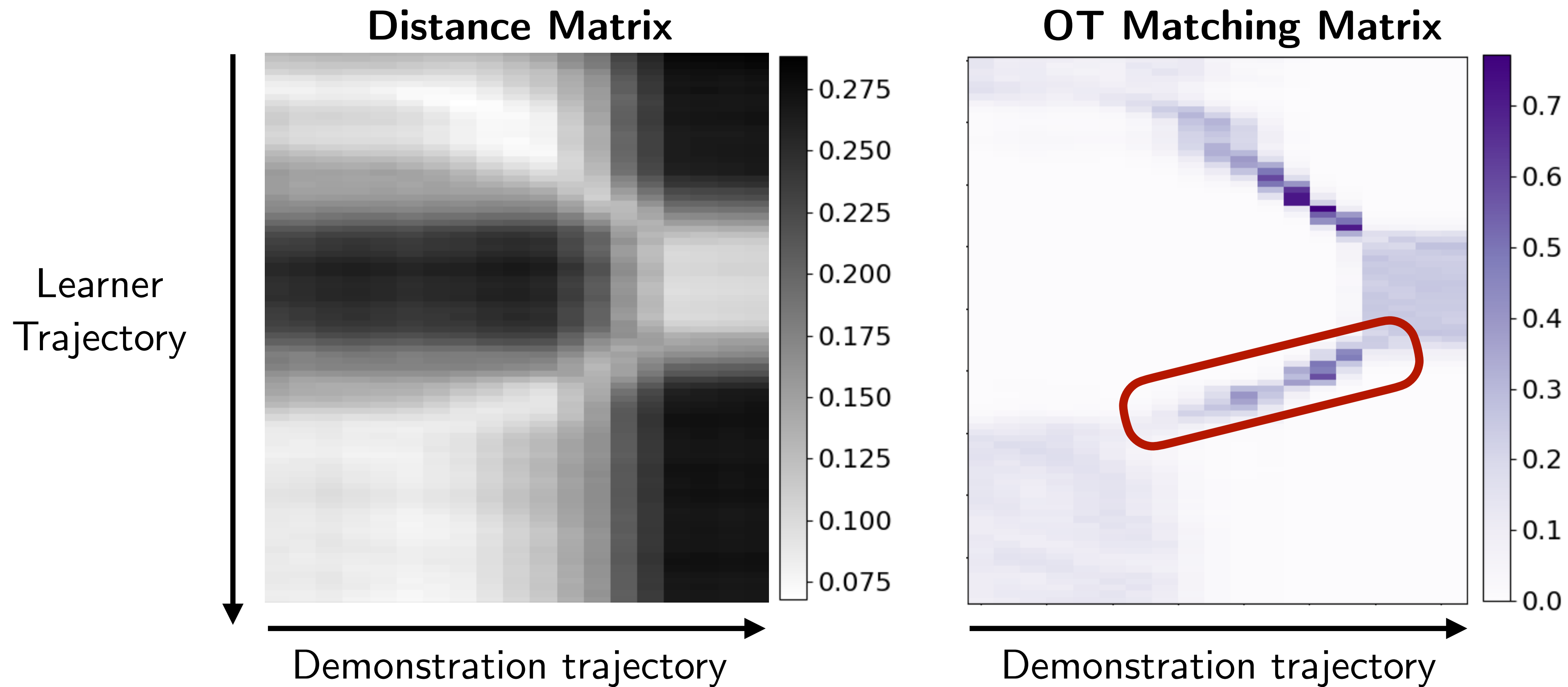
Demo

# OT fails to enforce **temporal ordering**



OT

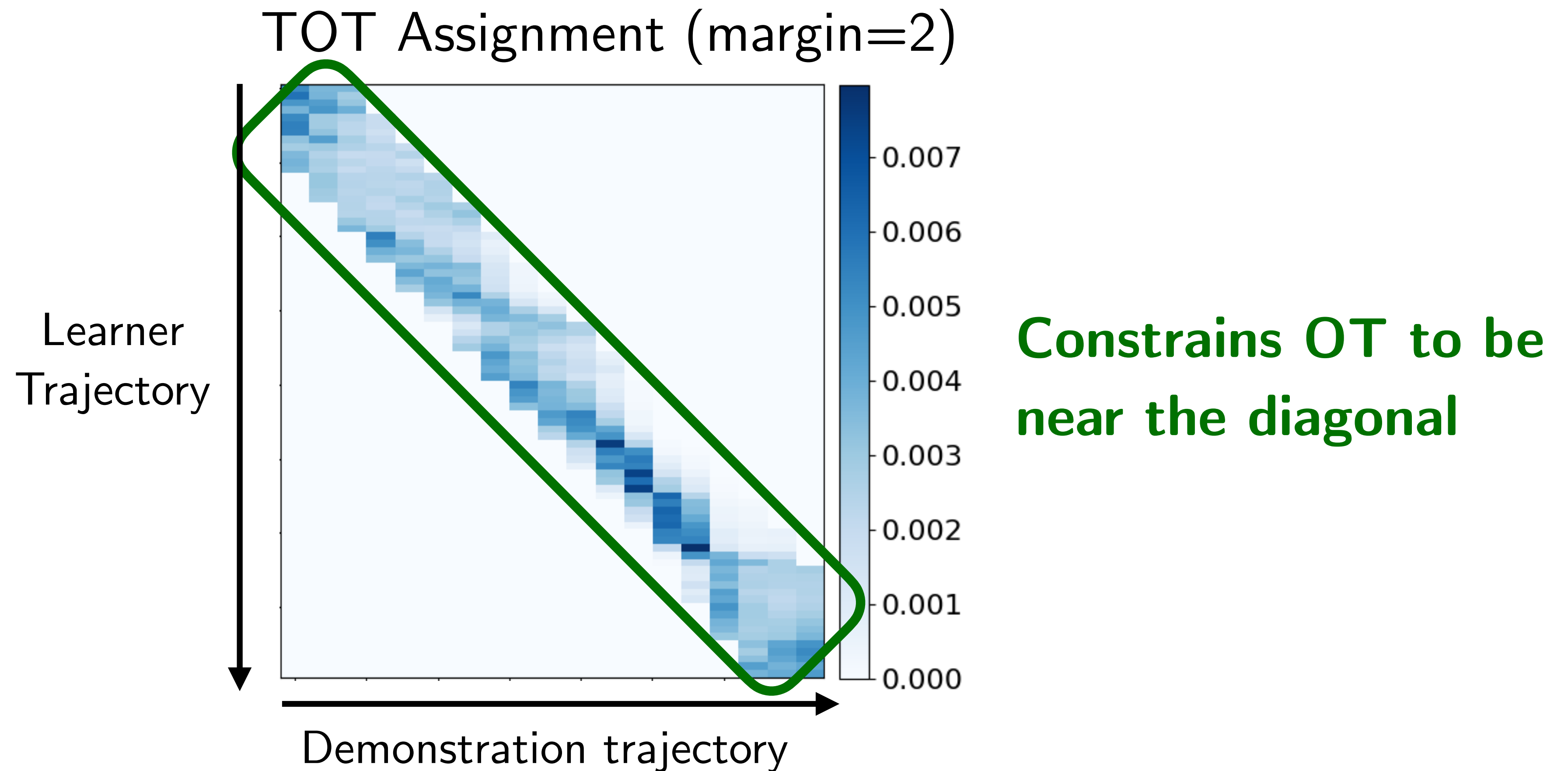
# OT fails to enforce **temporal ordering**



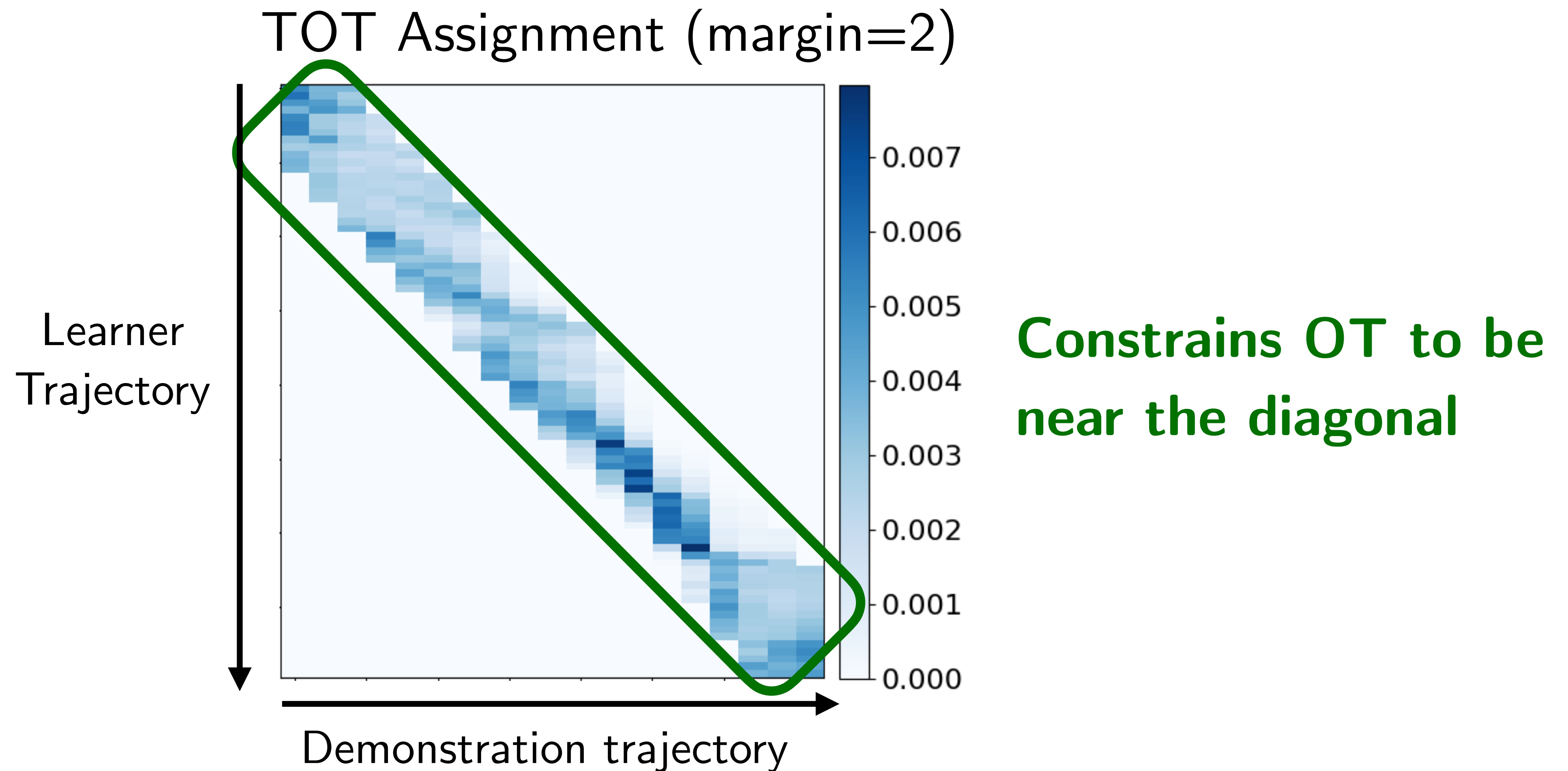
**OT matches later learner frames to earlier subgoals.**



# TemporalOT solves the ordering problem



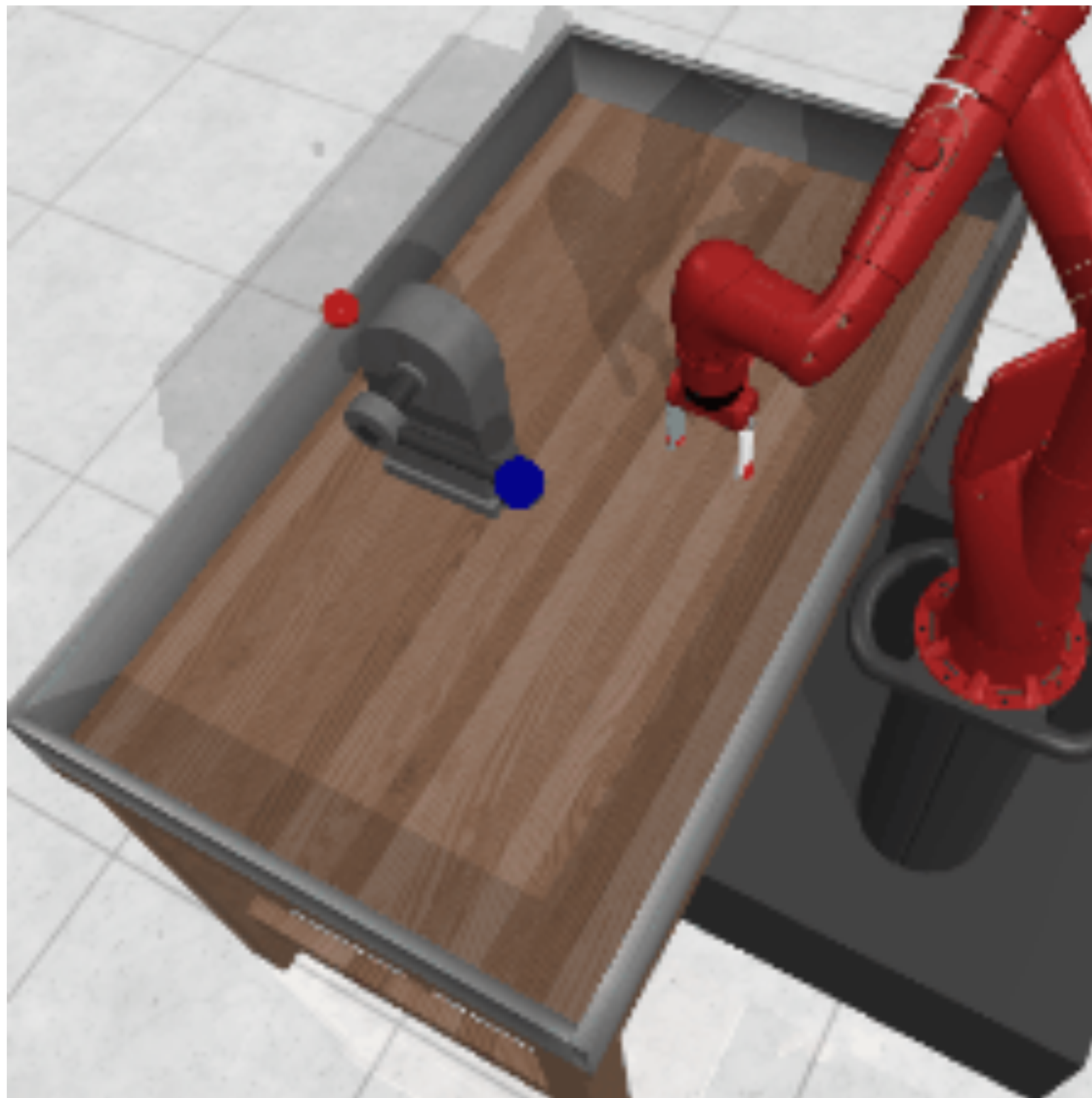
# TemporalOT solves the ordering problem



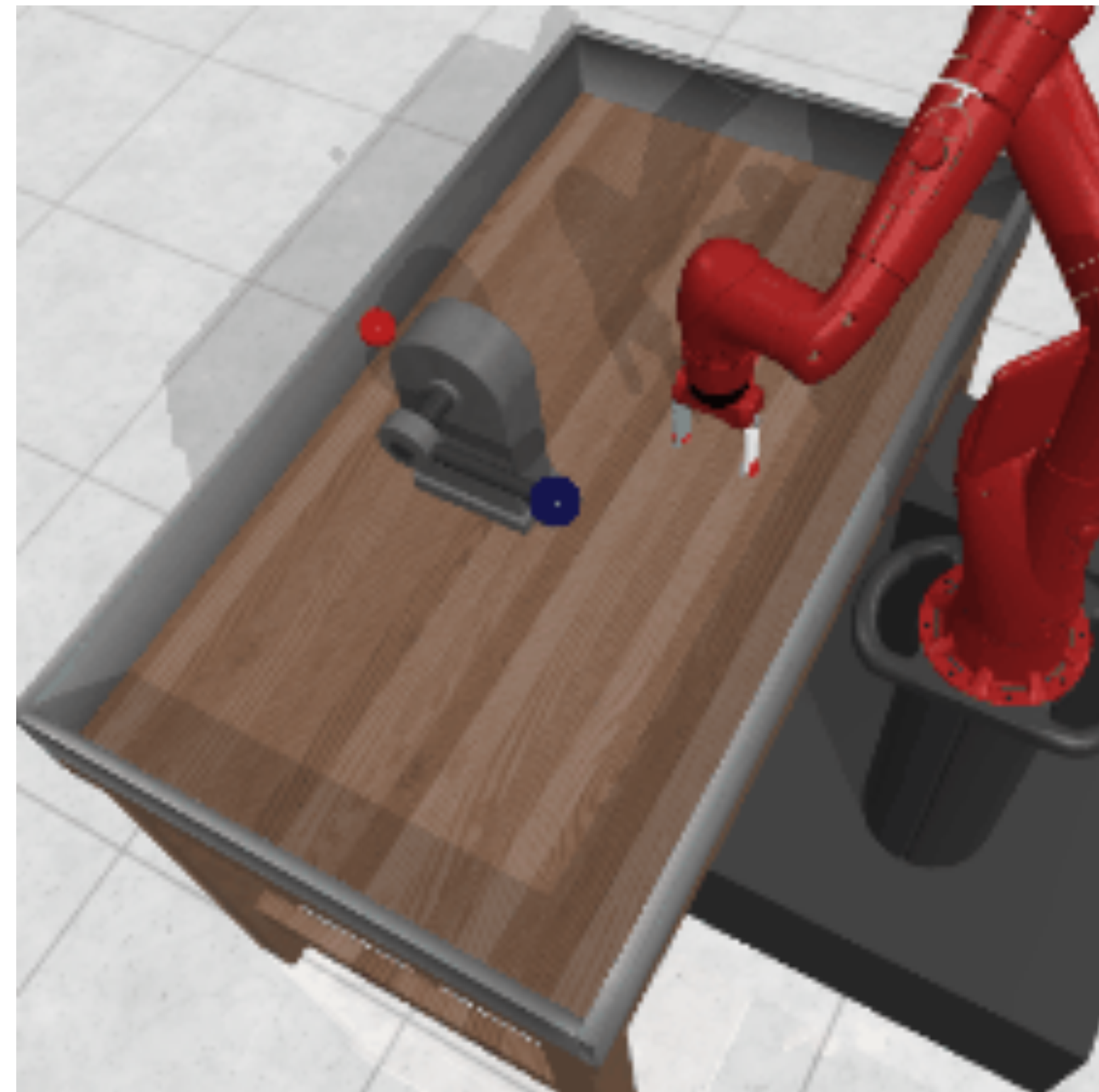
**Problem:** assumes learner and expert are *temporally aligned*

# TemporalOT fails to enforce **subgoal coverage** under temporal misalignment

TOT (margin=10)

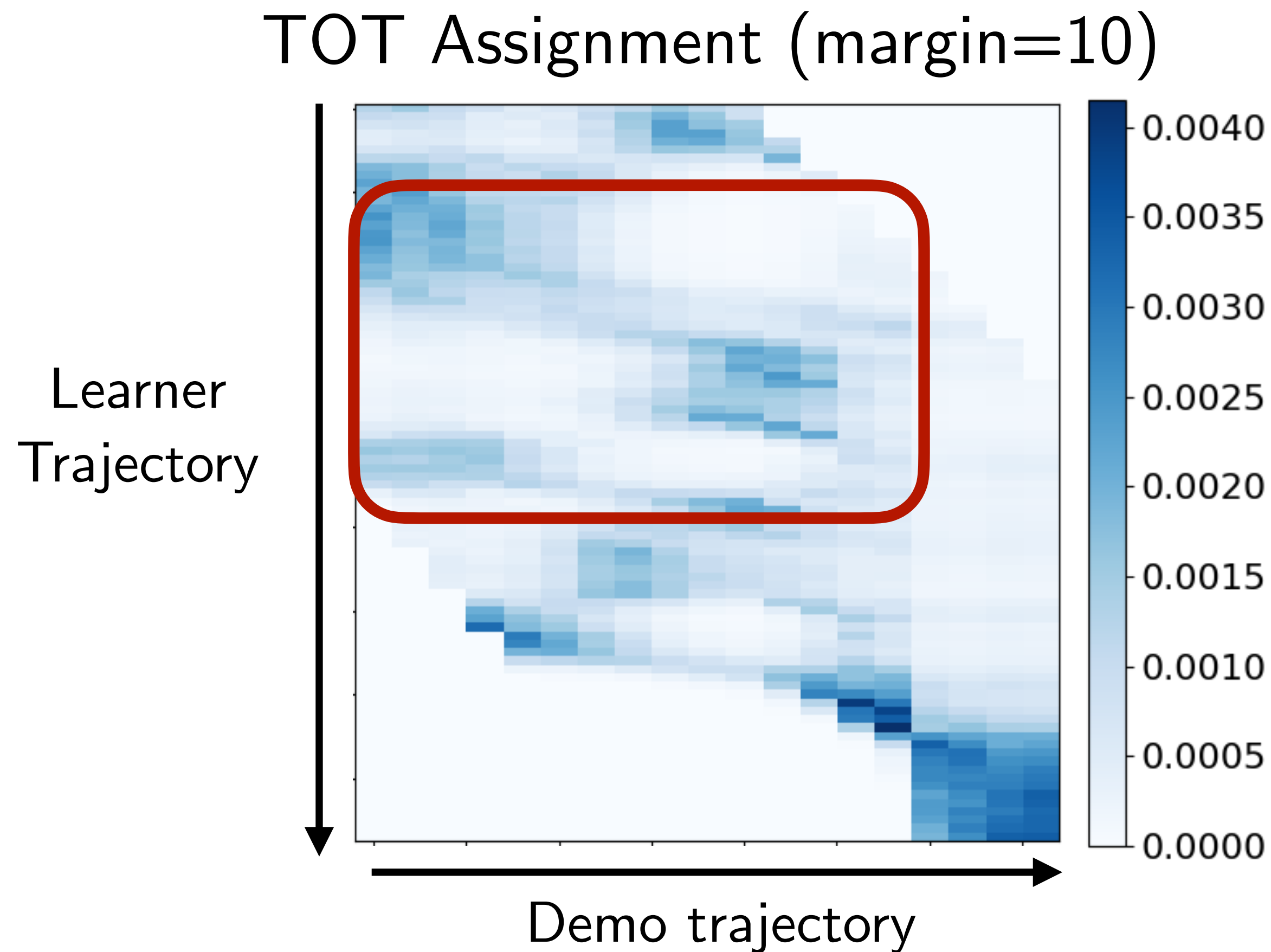


TOT (margin=2)

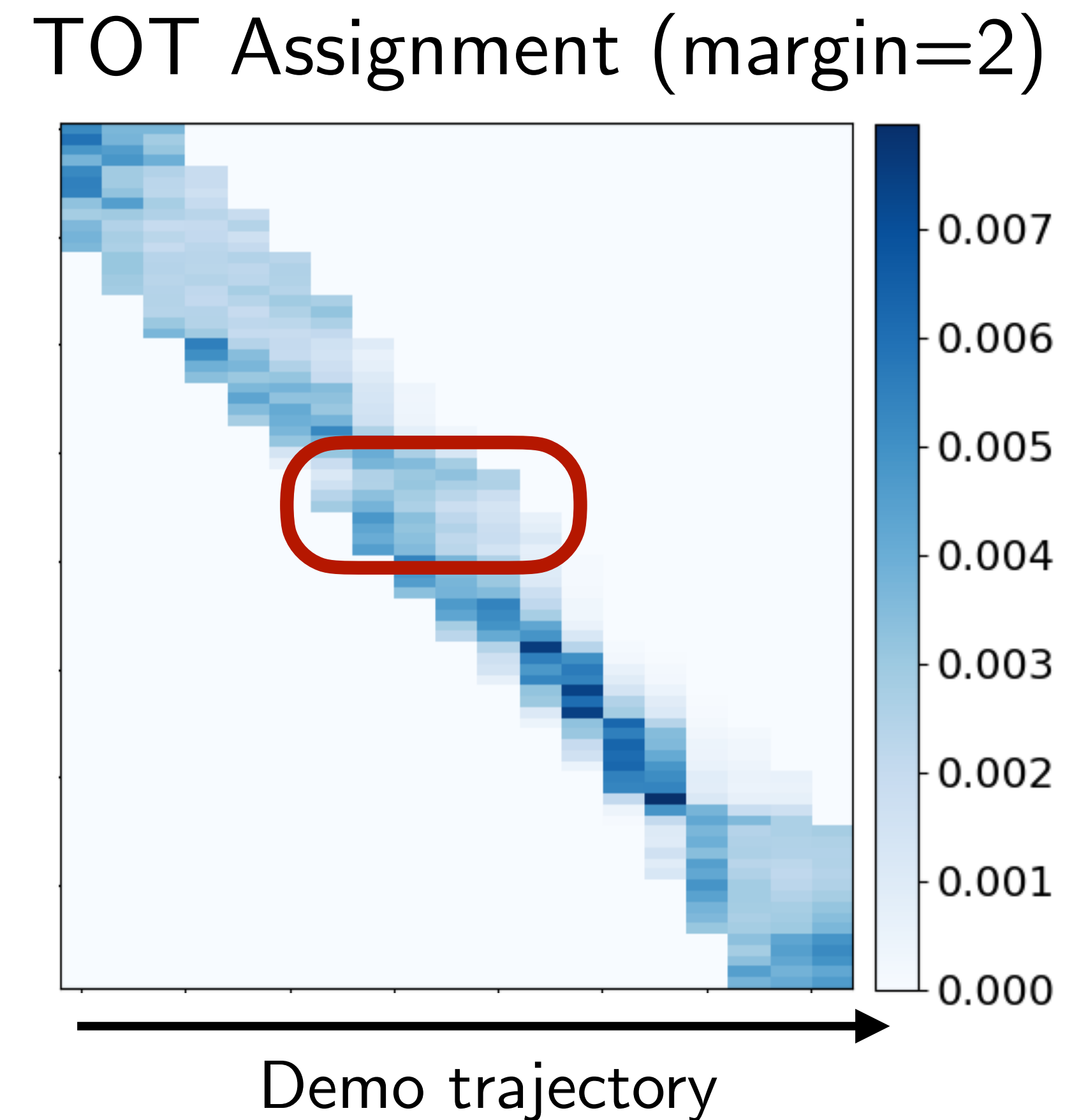




# TemporalOT fails to enforce **subgoal coverage** under temporal misalignment



**Fails to enforce temporal ordering  
due to the large mask window**



**Over-constrains assignment,  
leading to slow and shaky motion**



Instead of matching at  
the frame-level,  
we should match at the  
**sequence-level**

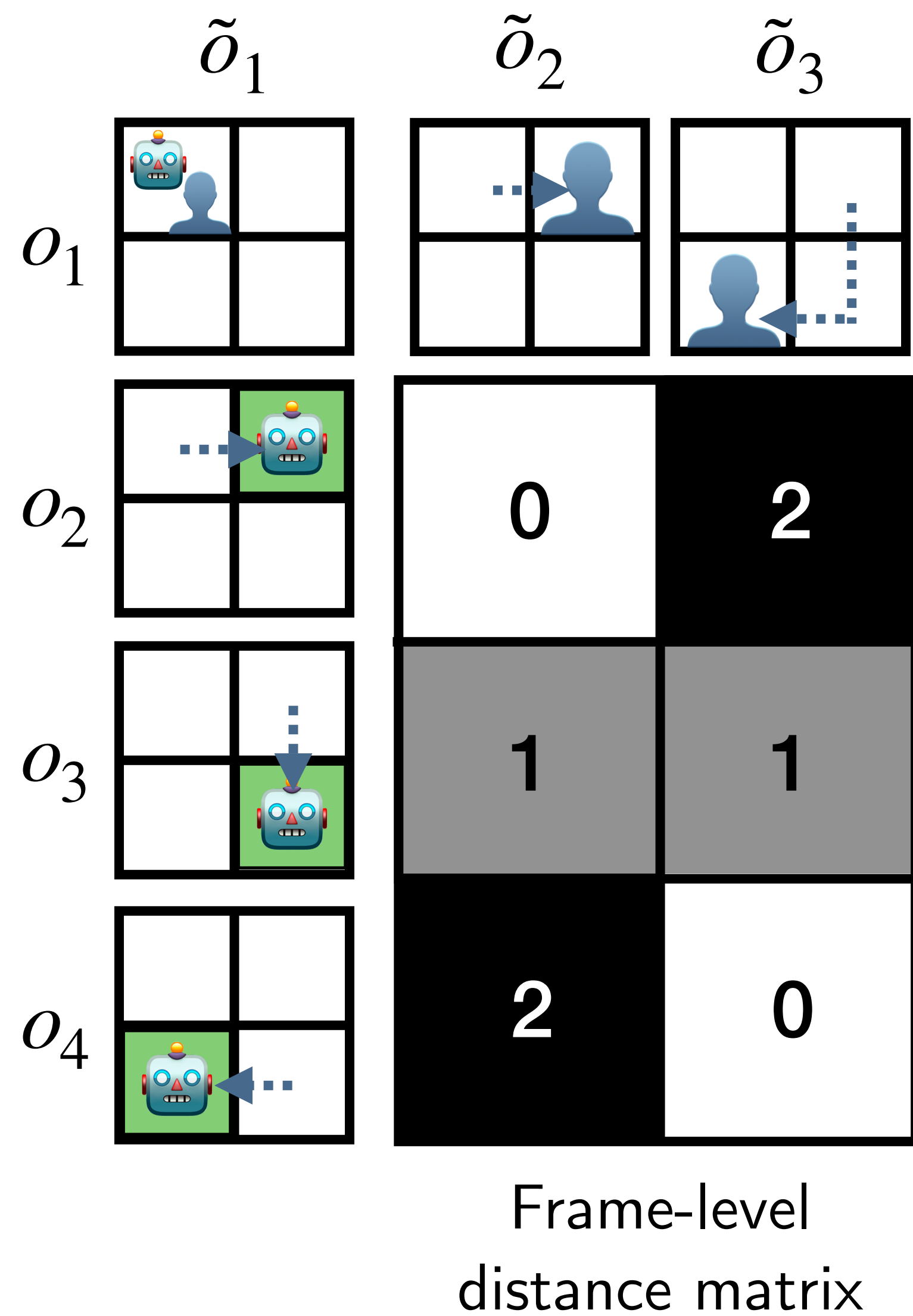


Sequence-level matching  
should enforce:

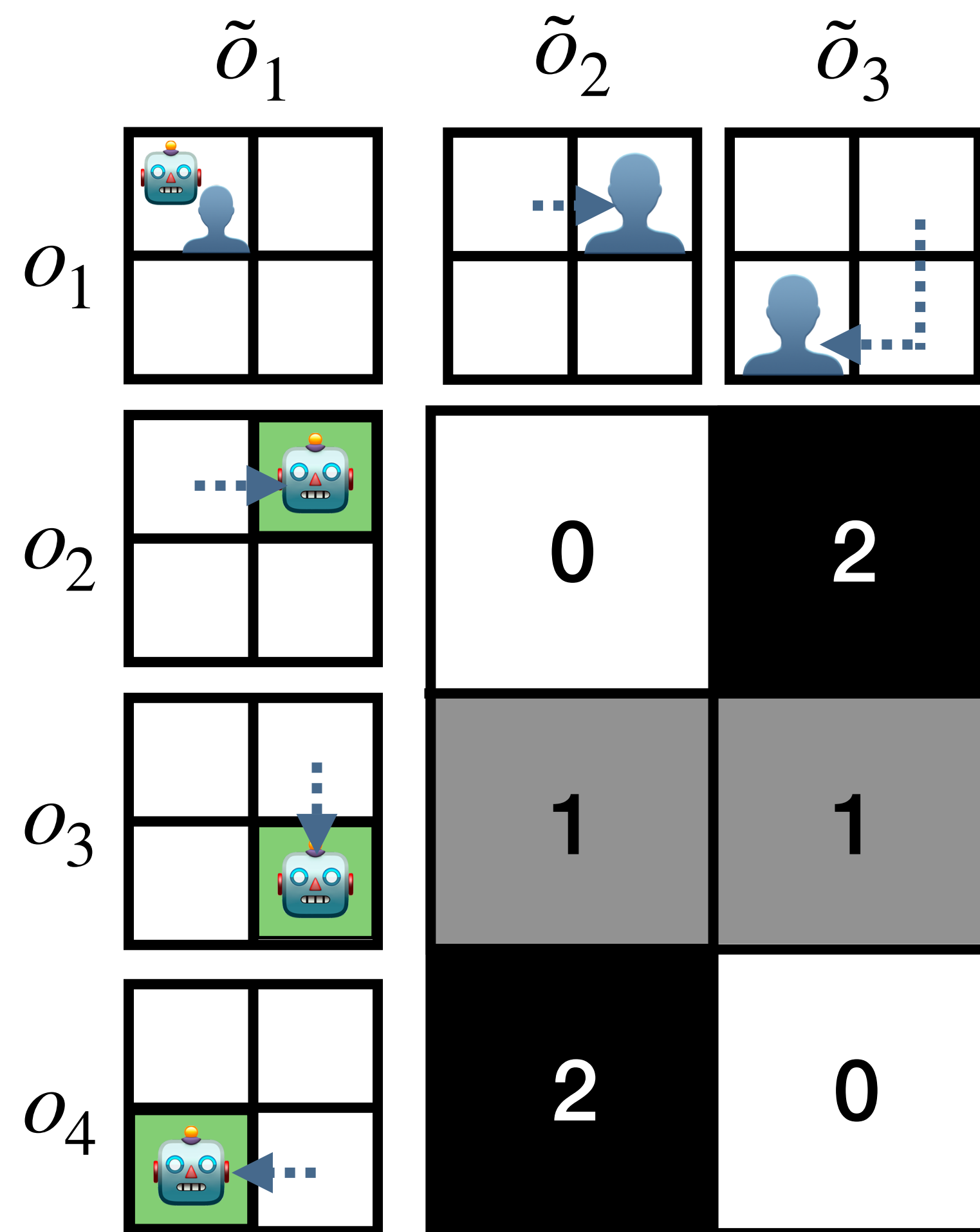
1. Subgoal **ordering**
2. Subgoal **coverage**



# ORdered Coverage Alignment (**ORCA**)



# ORdered Coverage Alignment (**ORCA**)



$$P_{t,j} = \exp(-\lambda d(o_t, \tilde{o}_j))$$

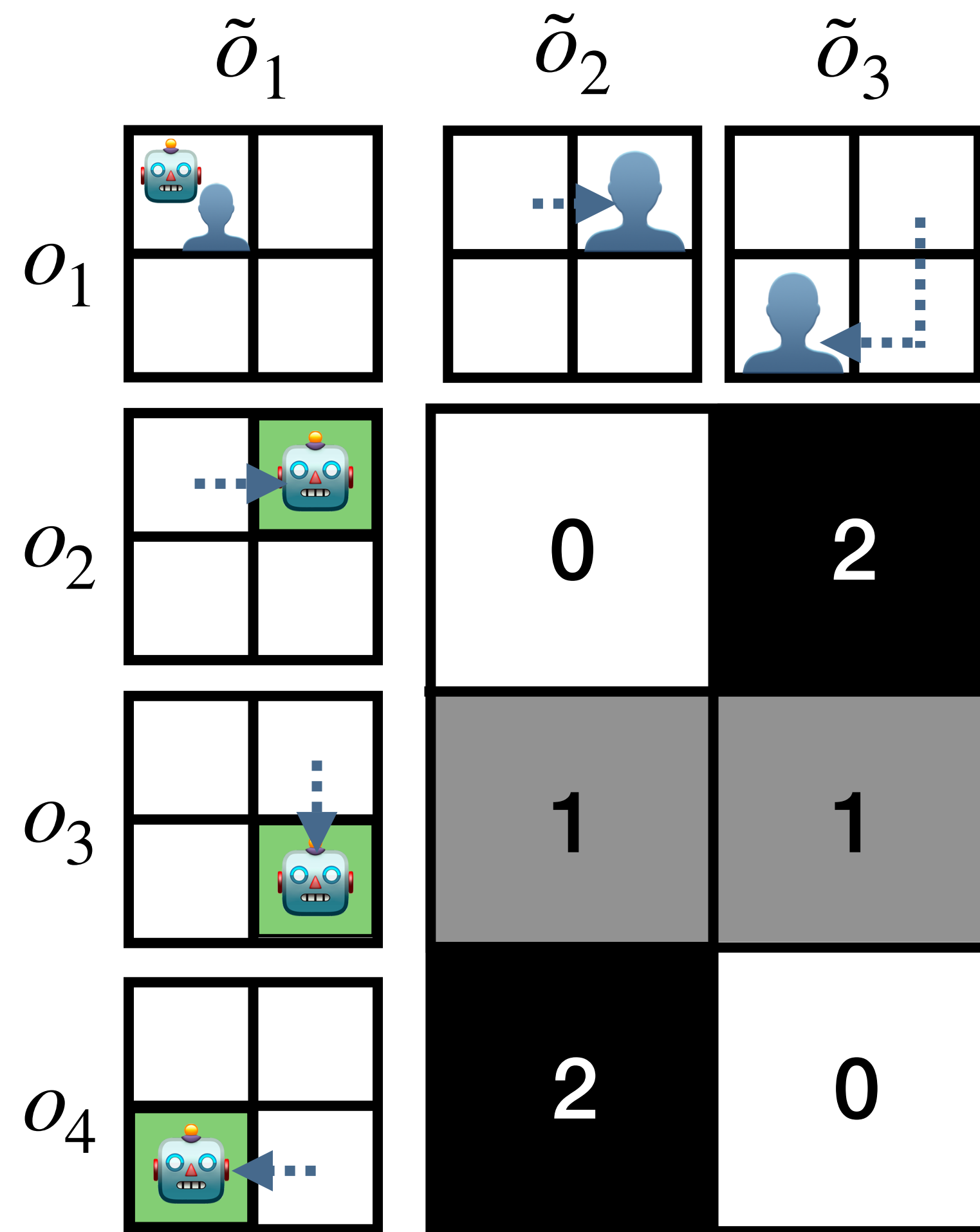


1.0	0.14
.37	.37
0.14	1.0

Frame-level  
distance matrix

Soft occupancy  
probability **P**

# ORdered Coverage Alignment (**ORCA**)



Frame-level  
distance matrix

$$P_{t,j} = \exp(-\lambda d(o_t, \tilde{o}_j))$$

$$C_{t,j} = \max\{C_{t-1,j}, C_{t,j-1}P_{t,j}\}$$

1.0	0.14
.37	.37
0.14	1.0

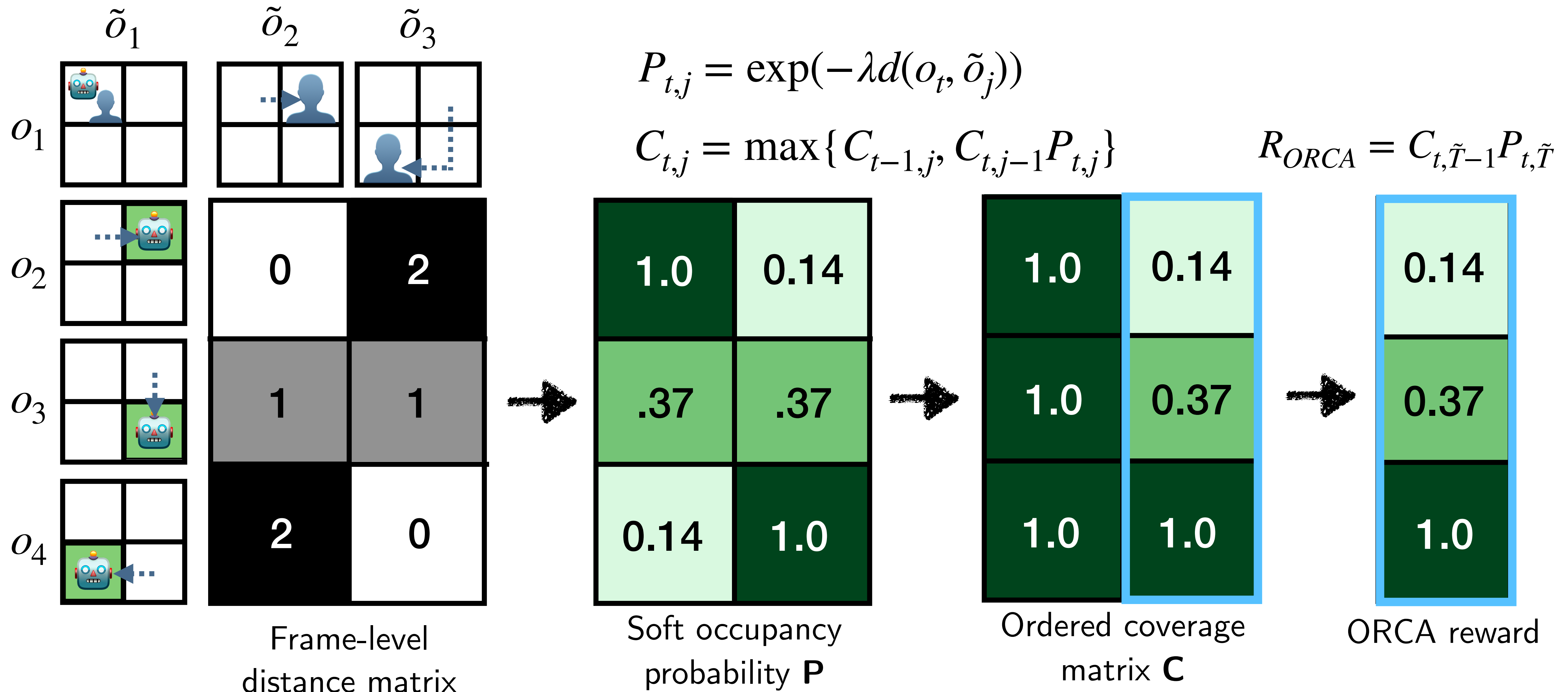
Soft occupancy  
probability **P**

1.0	0.14
1.0	0.37
1.0	1.0

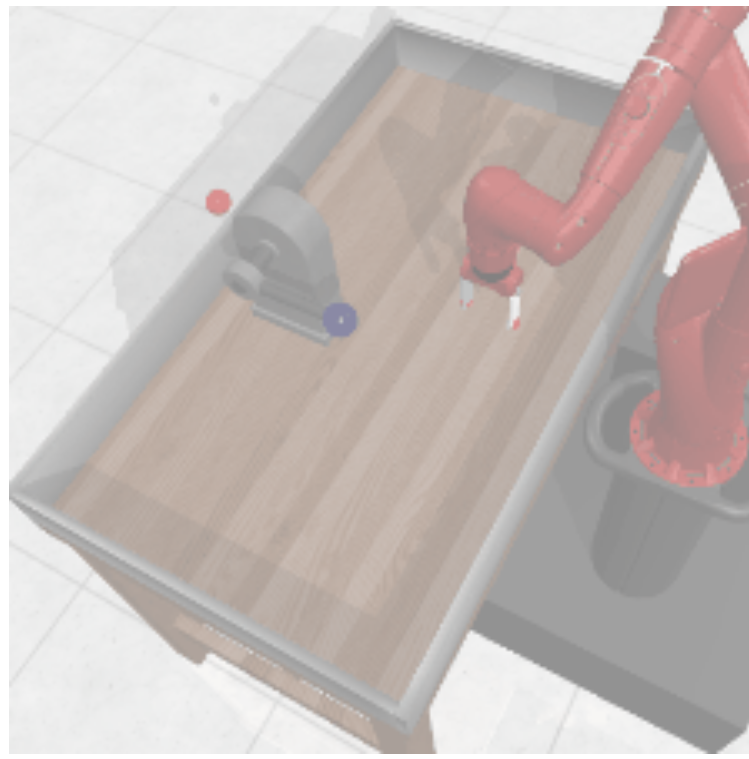
Ordered coverage  
matrix **C**



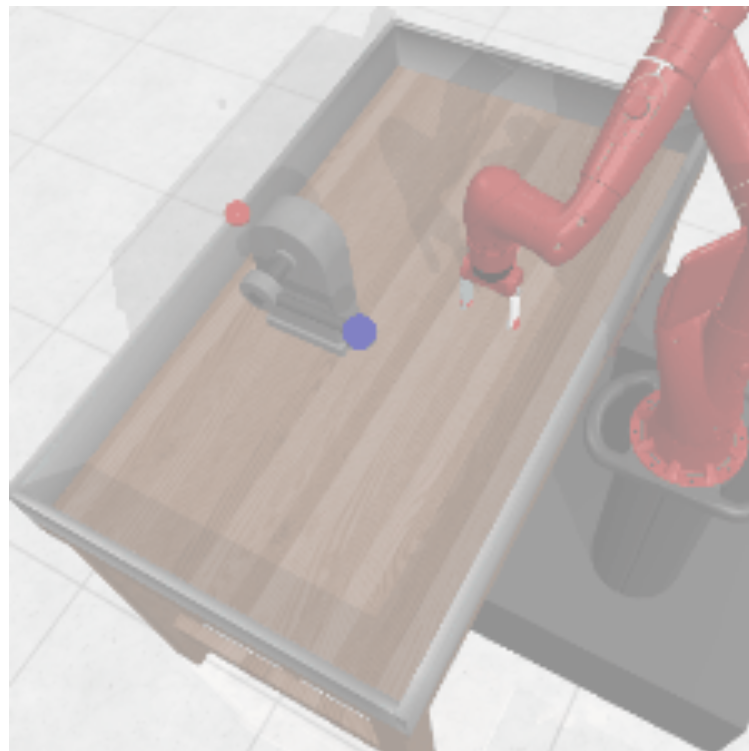
# ORdered Coverage Alignment (**ORCA**)



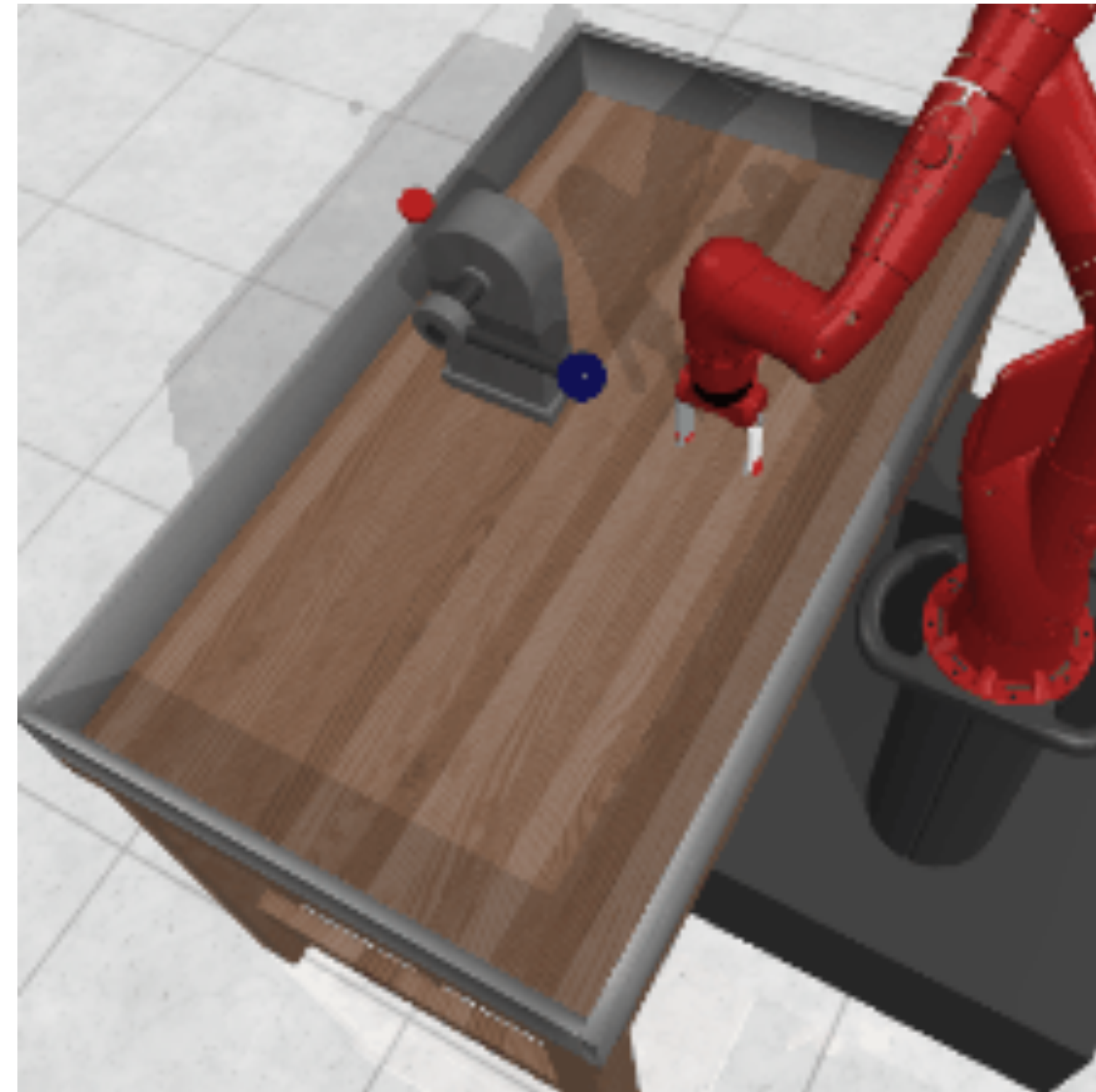
# ORCA completes tasks efficiently and effectively



OT: **✗** subgoal order



TOT: **✗** subgoal coverage



**ORCA: ✓ subgoal order ✓ subgoal coverage**

\*See our paper for proofs that ORCA enforces ordering and coverage



# ORCA achieves better performance on temporally misaligned demonstrations

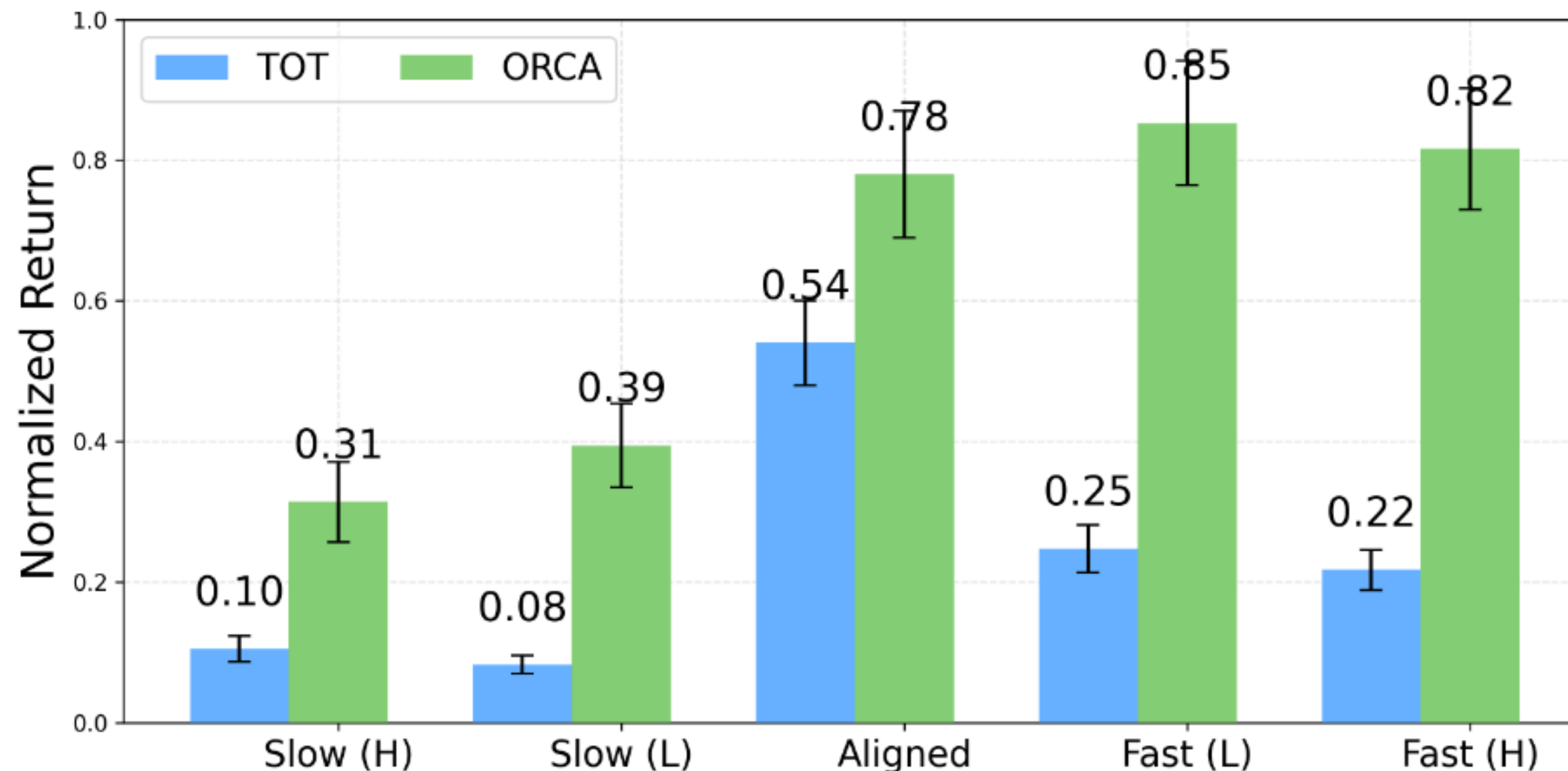
Category	Environment	Threshold	DTW	OT	TemporalOT	ORCA (NP)	<b>ORCA</b>
Easy	Button-press	0.30 (0.10)	0.00 (0.00)	0.00 (0.00)	0.10 (0.02)	0.45 (0.11)	<b>0.62 (0.11)</b>
	Door-close	0.34 (0.07)	0.00 (0.00)	0.00 (0.00)	0.19 (0.01)	0.86 (0.01)	<b>0.88 (0.01)</b>
Medium	Door-open	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.08 (0.01)	<b>1.60 (0.09)</b>	0.89 (0.13)
	Window-open	0.72 (0.14)	0.00 (0.00)	0.19 (0.06)	0.26 (0.05)	<b>0.86 (0.17)</b>	<b>0.85 (0.16)</b>
	Lever-pull	0.07 (0.02)	0.00 (0.00)	0.00 (0.00)	0.07 (0.03)	<b>0.27 (0.08)</b>	<b>0.28 (0.09)</b>
	Hand-insert	0.00 (0.00)	0.00 (0.00)	0.03 (0.02)	0.00 (0.00)	<b>0.08 (0.08)</b>	<b>0.04 (0.04)</b>
	Push	0.07 (0.05)	0.00 (0.00)	0.03 (0.01)	0.01 (0.01)	0.02 (0.02)	0.00 (0.00)
Hard	Basketball	0.00 (0.00)	0.00 (0.00)	0.00 (0.00)	0.01 (0.01)	<b>0.07 (0.03)</b>	0.01 (0.00)
	Stick-push	0.12 (0.04)	0.00 (0.00)	0.07 (0.02)	0.36 (0.00)	0.46 (0.13)	<b>1.25 (0.04)</b>
	Door-lock	0.00 (0.00)	0.05 (0.02)	0.04 (0.02)	0.00 (0.00)	<b>0.23 (0.09)</b>	<b>0.19 (0.08)</b>
<b>Average</b>		0.16 (0.02)	0.01 (0.00)	0.04 (0.01)	0.11 (0.01)	<b>0.49 (0.04)</b>	<b>0.50 (0.04)</b>

(Top) Meta-world tasks

(Right) Humanoid tasks

Task	TOT	ORCA (NP)	<b>ORCA</b>
Arm up (L)	5.29 (2.22)	65.9 (8.25)	<b>81.6 (3.65)</b>
Arm up (R)	7.67 (2.88)	<b>92.5 (4.71)</b>	49.6 (5.00)
Arms out	1.62 (0.75)	<b>72.7 (10.1)</b>	8.50 (2.60)
Arms down	11.6 (3.56)	19.7 (5.03)	<b>33.4 (7.20)</b>
Average	6.55 (2.35)	<b>62.9 (7.02)</b>	43.3 (4.61)

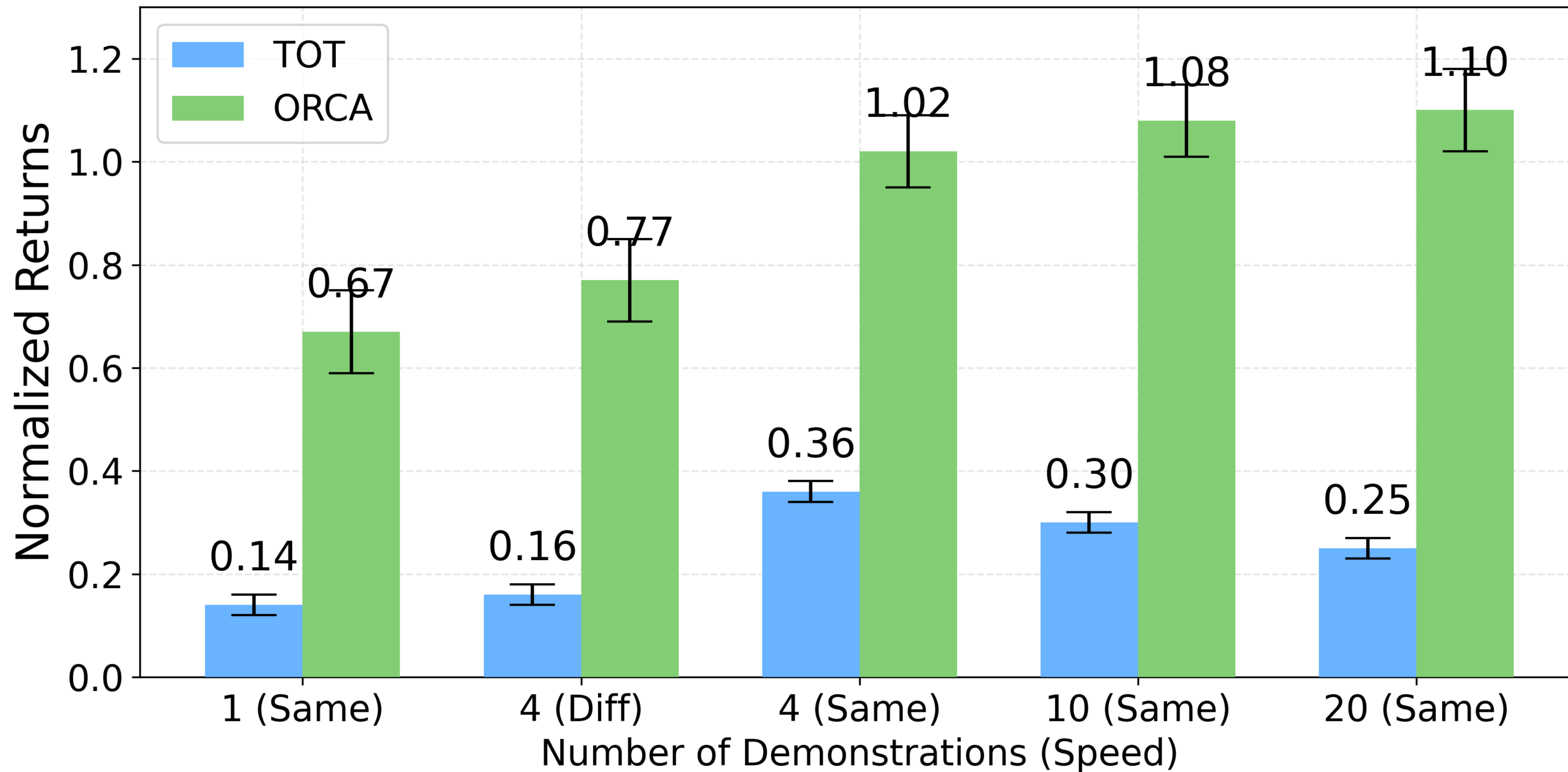
# ORCA achieves better performance on temporally misaligned demonstrations



ORCA beats baselines given faster, slower, and same-speed demonstrations

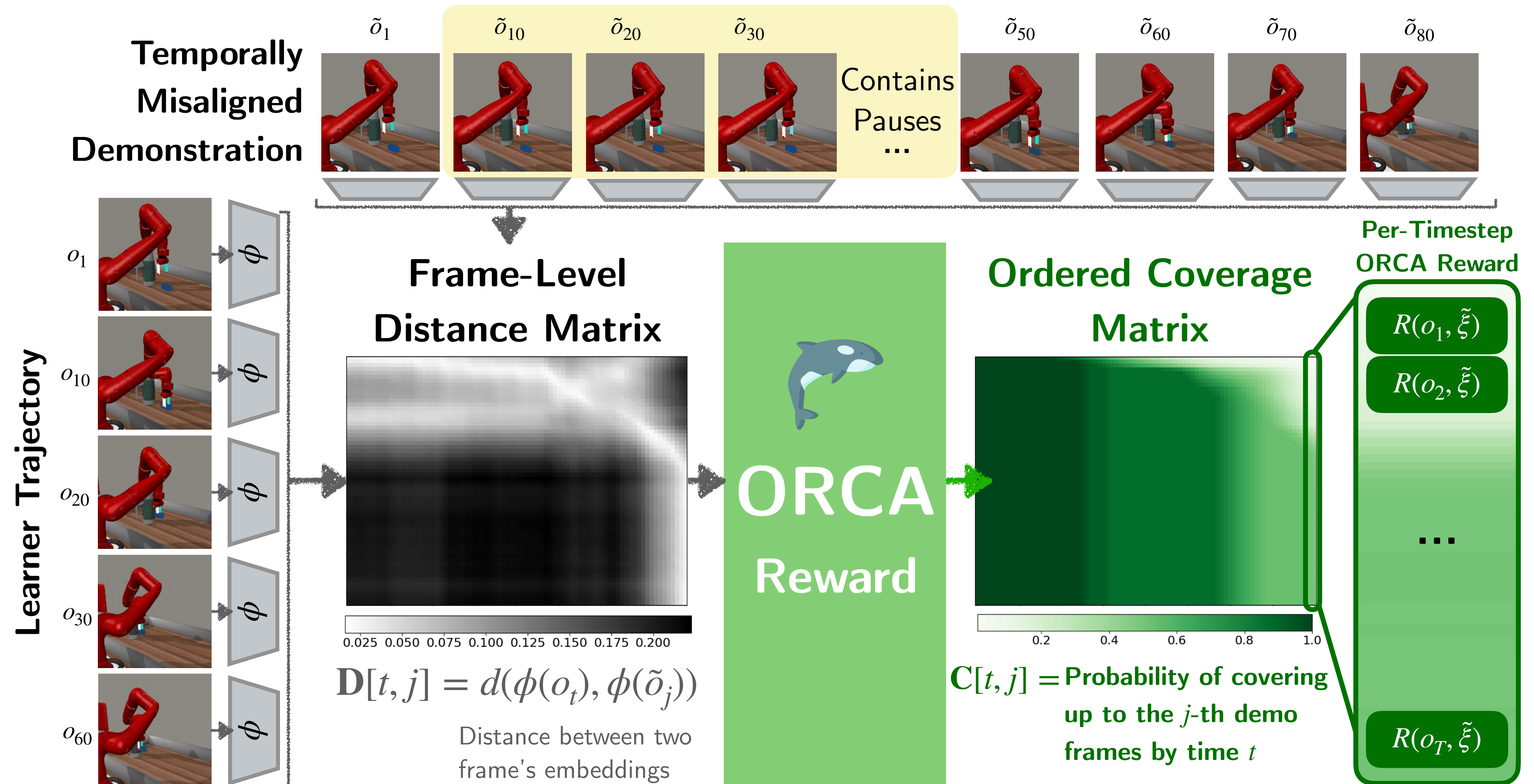


# ORCA scales with more demonstrations



ORCA beats baselines given multiple demonstrations,  
whether they are the same speed or different speeds

# ORCA is a principled reward function for imitation learning from a **temporally misaligned video**





# Thank You!



William Huey\*, Huaxiaoyue (Yuki) Wang\*, Anne Wu, Yoav Artzi, Sanjiban Choudhury  
{wph52, hw575}@cornell.edu, [willhuey.com](http://willhuey.com), [lunay0yuki.github.io](http://lunay0yuki.github.io)