# ReAgent:
# Point Cloud Registration using Imitation and Reinforcement Learning

Dominik bauer, Timothy Patten and Markus Vincze
TU Wien
Vienna, Austria

**Presenter : 김진용**

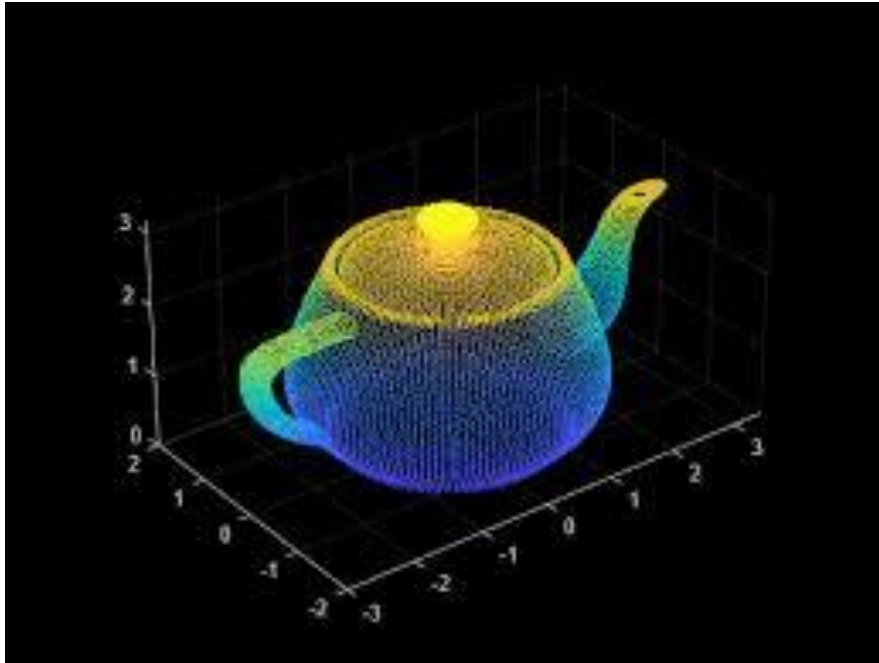# Content
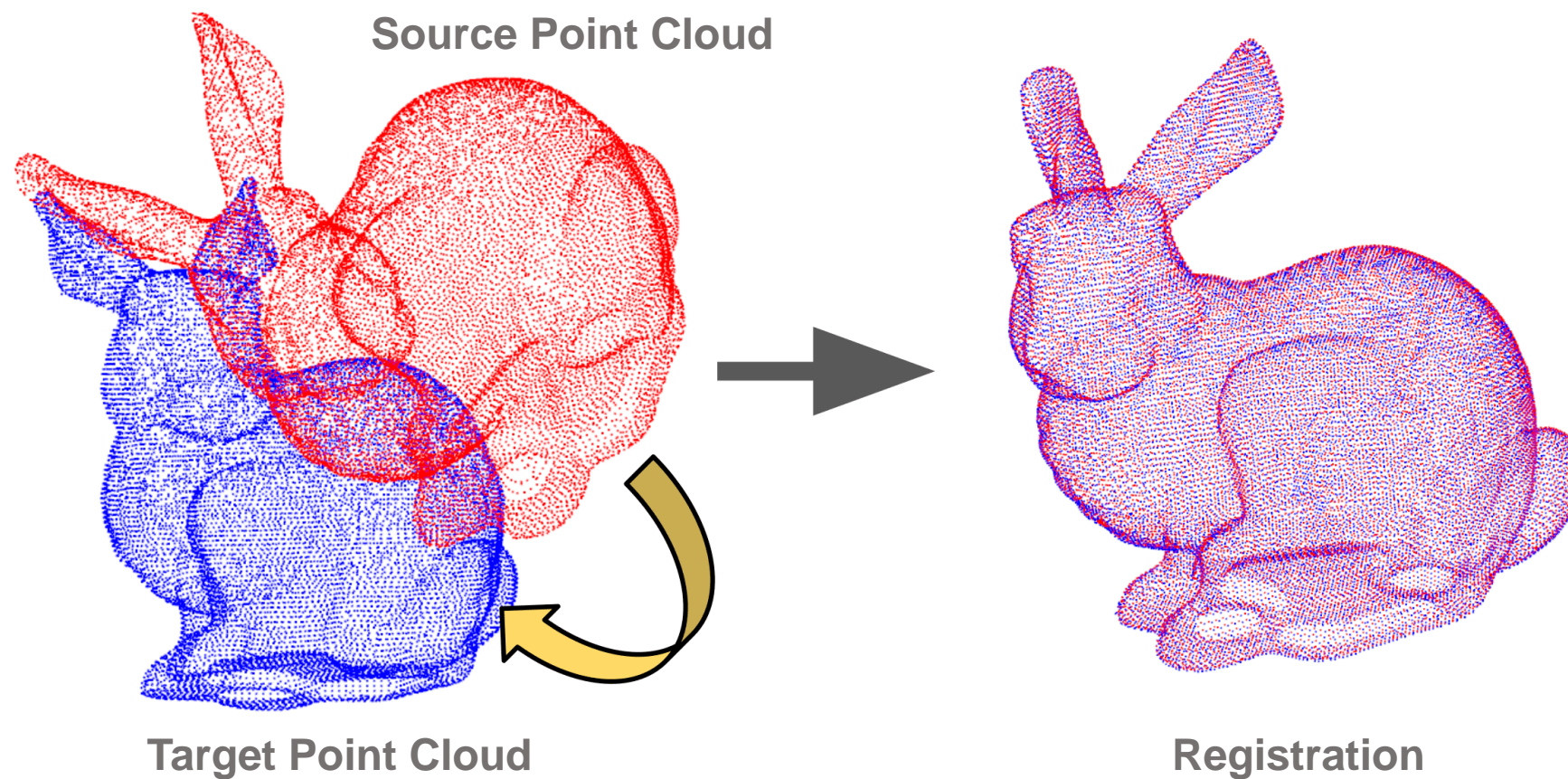
# Background

## What is Point Cloud?



**Point Cloud**

(X, Y, Z) **n x 3** set

If number of point cloud is **n**, it has below structure.

[[X1, Y1, Z1],
[X2, Y3, Z2],

…
[Xn, Yn, Zn]]

# Background

**What is Point Cloud Registration?**

**Source Point Cloud**

**Target Point Cloud**

**Registration**
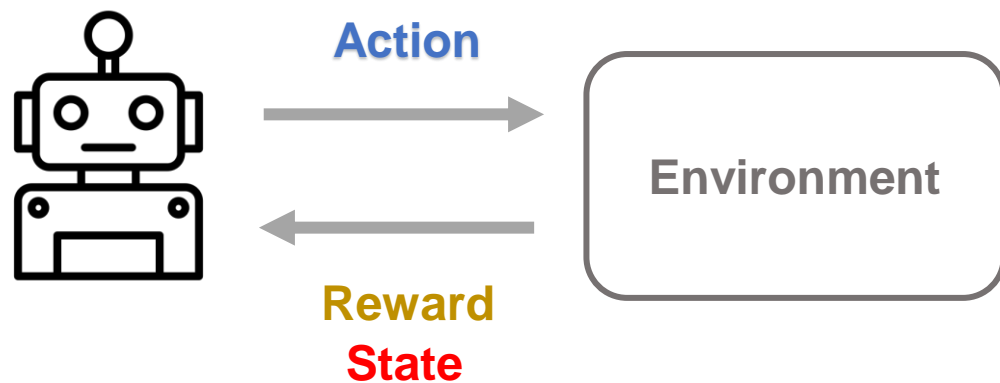
# Background

What is Reinforcement Learning and Imitation Learning?

**Reinforcement Learning**

Action

Environment
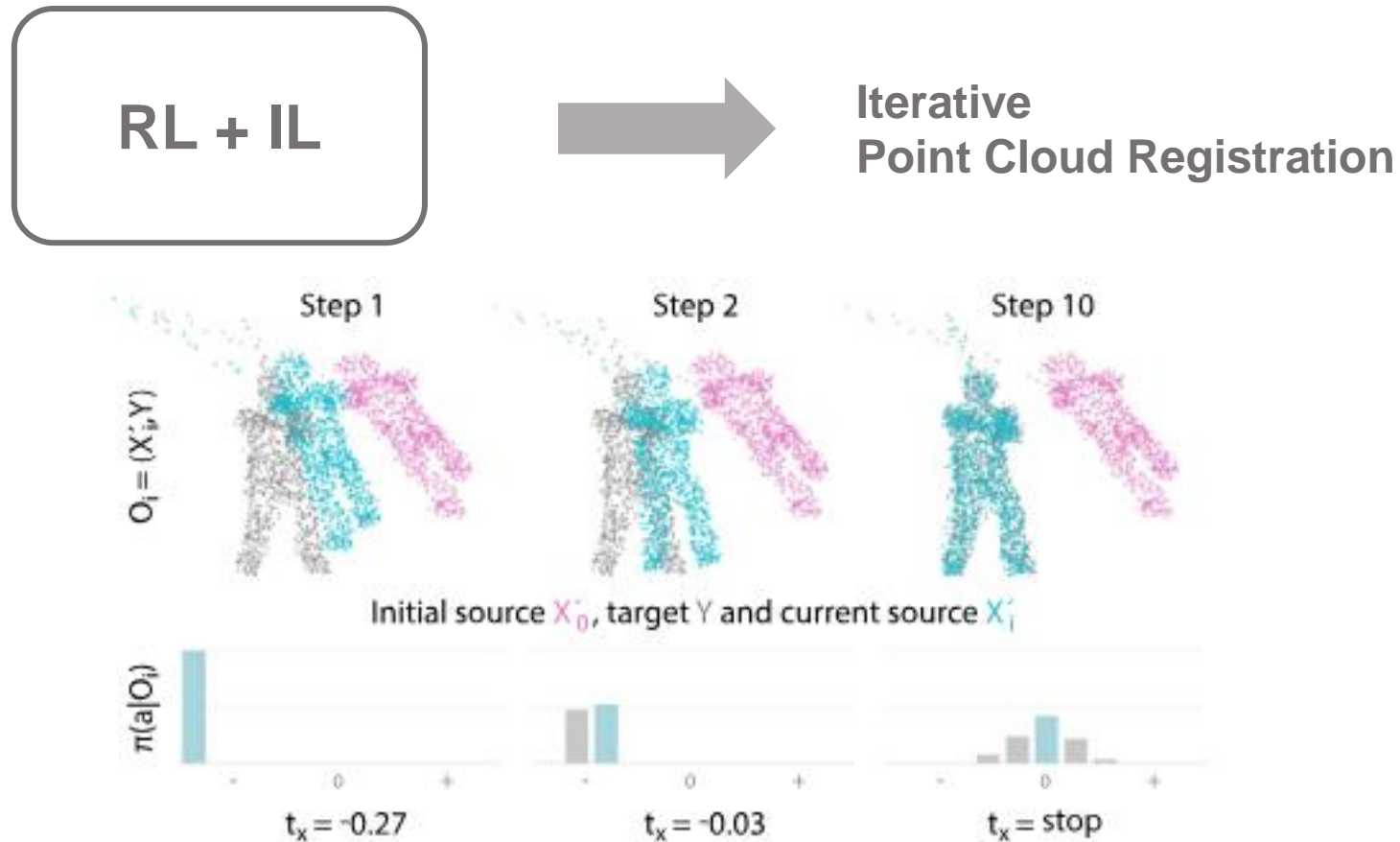
Reward
State

**Imitation Learning**

Action

Environment

Imitation

State

# Define Problem

- **Classical Registration Methods generalize well to novel domains but fail when given a noisy observation or a bad initialization.**
- **Learning-based methods, in contrast, are more robust but lack in generalization capacity.**

RL + IL → Iterative Point Cloud Registration



$O_i = (X'_i, Y)$

Step 1   Step 2   Step 10

Initial source $X'_0$, target $Y$ and current source $X'_i$

$\pi(a|O_i)$

$t_x = -0.27$   $t_x = -0.03$   $t_x = stop$

# Point Cloud Registration Agent

Reagent Architecture



**State Vector**
The number of point cloud is 1024.
We have to extract feature to reduce computational cost.
We may use MLP to get **state vector** $S\_i$.

**Action vector**
From state vector, actor-critic predicts **rotation and translation action vector**.

**Step**
We **transform previous observation to new observation** by using transformation matrix. And environment offers reward comparing chamfer distance between two observations.

# Point Cloud Registration Agent



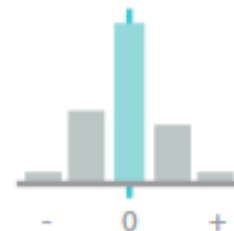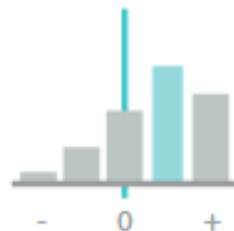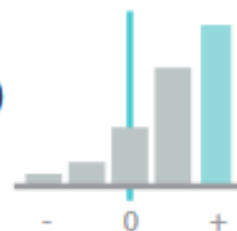Discrete action space

Translation : cartesian
Rotation : radian
x,y,z,no_ops : 33 dim

P(a|s)

-0.27  -0.09  -0.03  -0.01  -0.0033  0  0.0033  0.01  0.03  0.09  0.27

O

π(a|O)

Large step to right
Small step to right
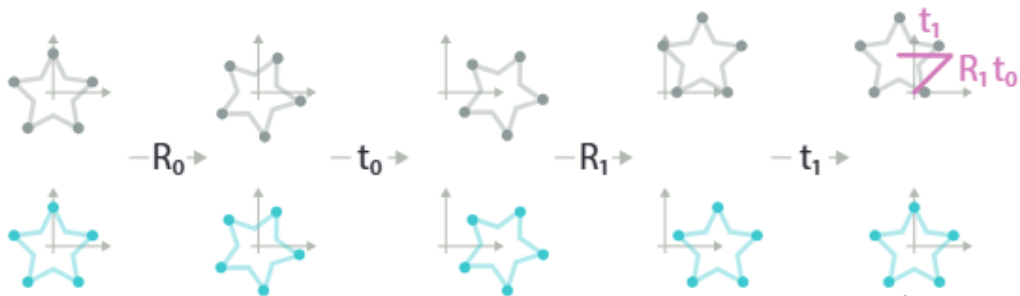No operation
Small step to left

# Point Cloud Registration Agent

**Disentangled transformation**



## Global Transformation

$$\begin{bmatrix} R_1 & t_1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} R_0 & t_0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} R_1 R_0 & R_1 t_0 + t_1 \\ 0 & 1 \end{bmatrix}$$

$$R_1(R_0 X + t_0) + t_1 = R_1 R_0 X + R_1 t_0 + t_1.$$

The center point of rotation matrix is global origin coordinate. It entangles rotation and translation matrices. It is hard to know global translation whether it is lead by rotation matrix or translation matrix.

## Disentangled Transformation

$$X_i = (\prod^{i} R_i) X + \sum^{i} t_i.$$

$$R_i = \hat{R}_i R_{i-1}, \quad t_i = \hat{t}_i + t_{i-1},$$

$$X'_i = R_i(X' - \mu_{X'}) + \mu_{X'} + t_i.$$

We may use rotation matrix by object center point. It improves **interpretability** how model acts by using object's local rotation.
It means that translation and rotation matrices are separated.

# Imitating an Expert Policy

Expert policy

$$Target = T_{gt} \times Source$$

State1 $\qquad Pred_1 = T_{pred_1} \times Source \qquad\qquad Target = \boxed{T_{gt} T_{pred_1}^{-1}} \times Pred_1$

State2 $\qquad Pred_2 = T_{pred_2} \times Pred_1 \qquad\qquad Target = \boxed{T_{gt} T_{pred_1}^{-1} T_{pred_2}^{-1}} \times Pred_2$

State3 $\qquad Pred_3 = T_{pred_3} \times Pred_2 \qquad\qquad Target = \boxed{T_{gt} T_{pred_1}^{-1} T_{pred_2}^{-1} T_{pred_3}^{-1}} \times Pred_3$

$\bullet \bullet \bullet$

$\bullet \bullet \bullet$

State10 $\qquad Pred_{10} = T_{pred_{10}} \times Pred_9$

# Imitating an Expert Policy

- If expert data of translation x is 0.3, since maximum size of expert policy is 0.27, 0.27 is chosen as gt.
- If expert data of translation x is 0.0001, 0 is chosen as gt.
- If expert data of translation x is -0.09, -0.09 is chosen as gt.

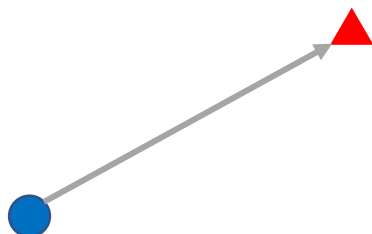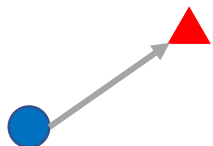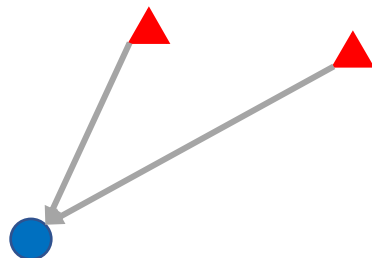# Improving through Reinforcement

**Reward function**

$$r = \begin{cases} -\varepsilon^-, & CD(X_i', X) > CD(X_{i-1}', X) \\ -\varepsilon^0, & CD(X_i', X) = CD(X_{i-1}', X) \\ \varepsilon^+, & CD(X_i', X) < CD(X_{i-1}', X). \end{cases}$$

$$d_{\text{ch}}(A, B) = \frac{1}{|A|} \sum_{a \in A} \min_{b \in B} \|a - b\|^2 + \frac{1}{|B|} \sum_{b \in B} \min_{a \in A} \|b - a\|^2$$
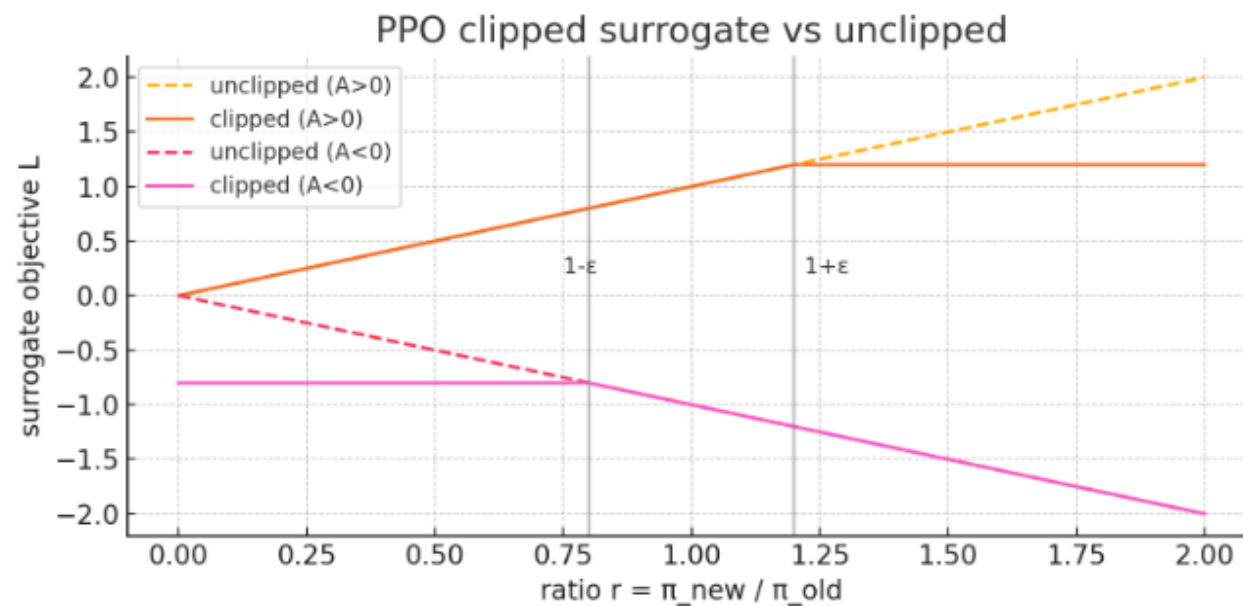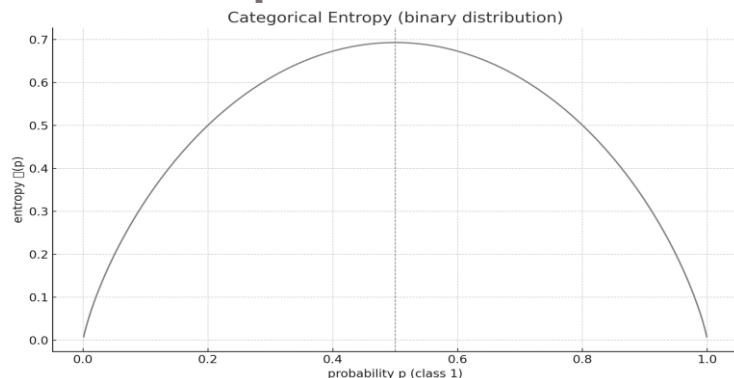
# Improving through Reinforcement

**PPO**

$$-L^{\text{clip}} + c_v \underbrace{\frac{1}{2}(V_\theta(s) - G_t)^2}_{\text{value}} - c_e \underbrace{\mathcal{H}[\pi_\theta]}_{\text{entropy}}$$

<span style="color:blue">**Policy Gradient**</span> **: If the action is good behavior, log probability of action is encouraged.**

<span style="color:blue">**Value loss**</span> **: High advantage means that value function doesn't predict value well. It is encouraged to reduce advantage to predict value precisely.**

<span style="color:blue">**Entropy loss**</span> **: Low probability is encouraged to be increased and High probability is encouraged to be reduced. Exploration is maximize.**



Categorical Entropy (binary distribution)



PPO clipped surrogate vs unclipped

<span style="color:blue">**PPO Clip**</span>
- **If action is good and ratio is over 1+E, surrogate objective is limited not to be increased for training stability.**
- **If action is good and ratio is under 1-E, loss is encouraged to be increased.**

# Improving through Reinforcement

---

**Algorithm 1** Combined Imitation and Reinforcement Learning using a Replay Buffer

---

1: **for all** observations $O$ in $\mathcal{O}$ **do**
2:    *% Gather replay buffer*
3:      **for** $N$ trajectories **do**
4:        **for** $n$ refinement steps **do**
5:          agent predicts policy $\pi(O)$ and value $\hat{v}$
6:          action $a$ is sampled from policy $\pi(O)$
7:          take action $a$, receive reward $r$ and next $O'$
8:          add sample to buffer $b$, step observation $O = O'$
9:        **end for**
10:      **end for**
11:    *% Process replay buffer*
12:    compute return $R$, shuffle buffer $b$
13:    **for all** samples in buffer $b$ **do**
14:      agent predicts new policy $\pi'(O)$ and value $\hat{v}'$
15:      *% Imitate expert*
16:      expert predicts action $a^*$
17:      compute cross-entropy loss $l_{IL}$ of $\pi'(O)$ and $a^*$
18:      *% Reinforce*
19:      compute PPO loss $l_{RL}$ of $\pi'(O)$ and $\pi(O)$
20:      *% Update agent*
21:      $l = l_{IL} + l_{RL} \cdot \alpha$
22:      backpropagate combined loss $l$
23:    **end for**
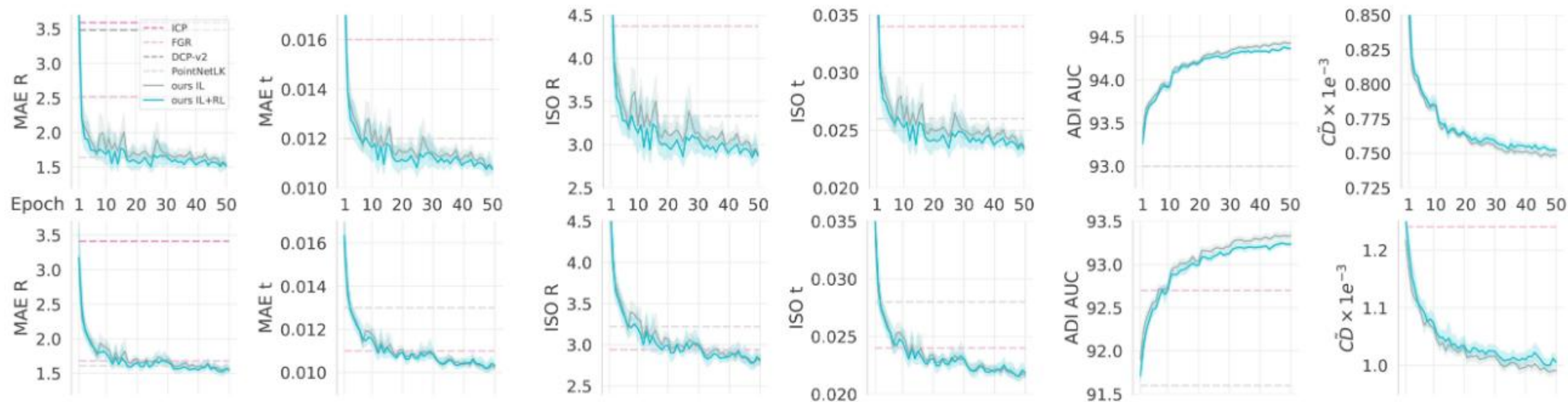24:    clear buffer $b$
25: **end for**

---

Fig. 5: Convergence of ReAgent with 10 random seeds on held-out models (top) and categories (bottom) of ModelNet40. The lines show the mean and the shaded areas indicate the 95%-confidence intervals. Best viewed digitally.

# Experiment

| | Segmented Objects | | | | | | |
|---|---|---|---|---|---|---|---|
| | MAE (↓) | | ISO (↓) | | ADI (↑) | $\tilde{CD}$ (↓) | T (↓) |
| | R | t | R | t | AUC | $\times 1e^{-3}$ | [ms] |
| ICP | 5.34 | 0.036 | 10.47 | 0.076 | 88.1 | 2.99 | **19** |
| FGR$^+$ | **0.11** | **0.001** | **0.19** | **0.001** | **99.7** | **0.16** | 131 |
| DCP-v2 | 7.42 | 0.050 | 14.93 | 0.102 | 72.4 | 4.93 | 54 |
| PointNetLK | 0.90 | 0.010 | 1.74 | 0.020 | 92.5 | 1.09 | 45 |
| ours IL | 0.77 | 0.006 | 1.33 | 0.012 | 95.7 | 0.30 | 21 |
| ours IL+RL | 0.93 | 0.007 | 1.66 | 0.014 | 95.4 | 0.34 | |

TABLE IV: Results on ScanObjectNN with the object segmented from the observation. Learning-based methods use the model trained on ModelNet40. Note that ↓ indicates that smaller values are better. Runtimes are for a single registration and 2048 points per cloud. $^+$ indicates that FGR additionally uses normals.

# Thank you