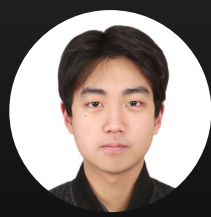

Fine-tuning Robotic Foundation Model Using Dynamic Scene Generation

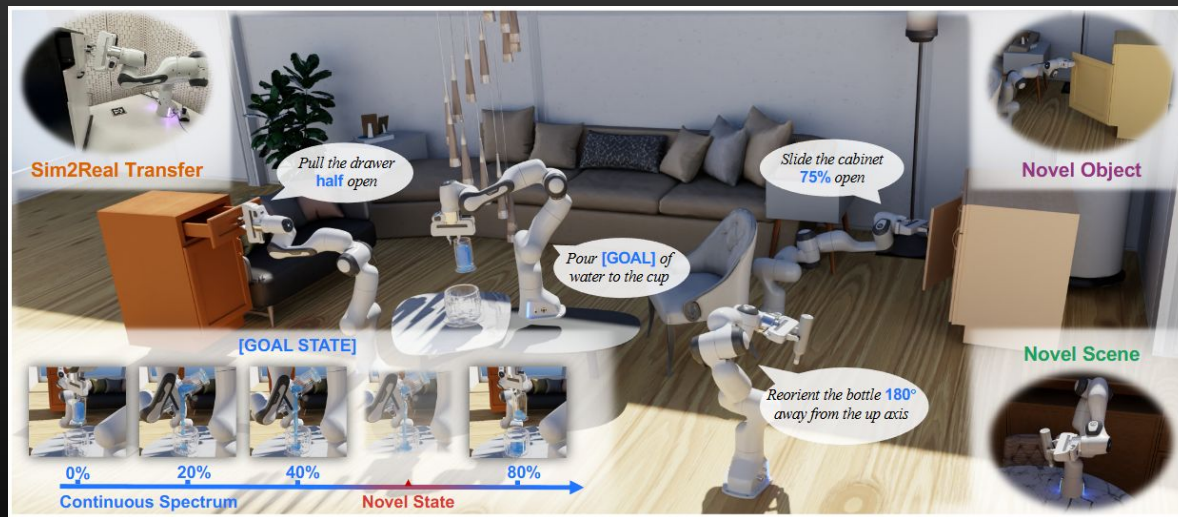
Challenge Team: **Reality Lab**



Dowon Kim, Chaewoo Lim, Sungyong Park, Sangmin Lee, Heewon Kim
Soongsil University

Robotic Manipulation in Embodied AI

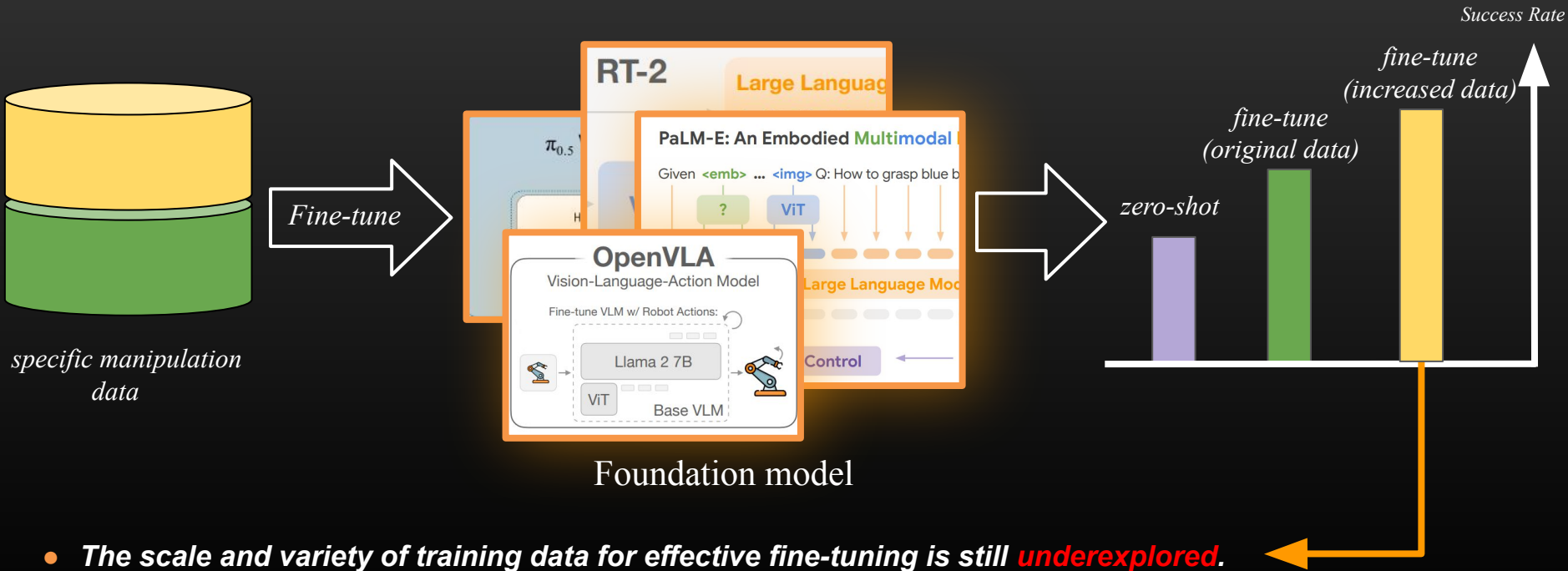
- Goal: **Perceive** and **manipulate** objects to complete designated tasks.
- Challenge: Real-world manipulation requires **generalization** to unseen scenarios.



→ To support this, ARNOLD provides **novel objects**, **scenes**, and **goal states** in evaluation.

Foundation Model for Robotic Manipulation

- Foundation models are a promising approach to **generalization** in robotic manipulation.
- However, even these models **require fine-tuning** data to optimize for specific tasks.



- The scale and variety of training data for effective fine-tuning is still **underexplored**.

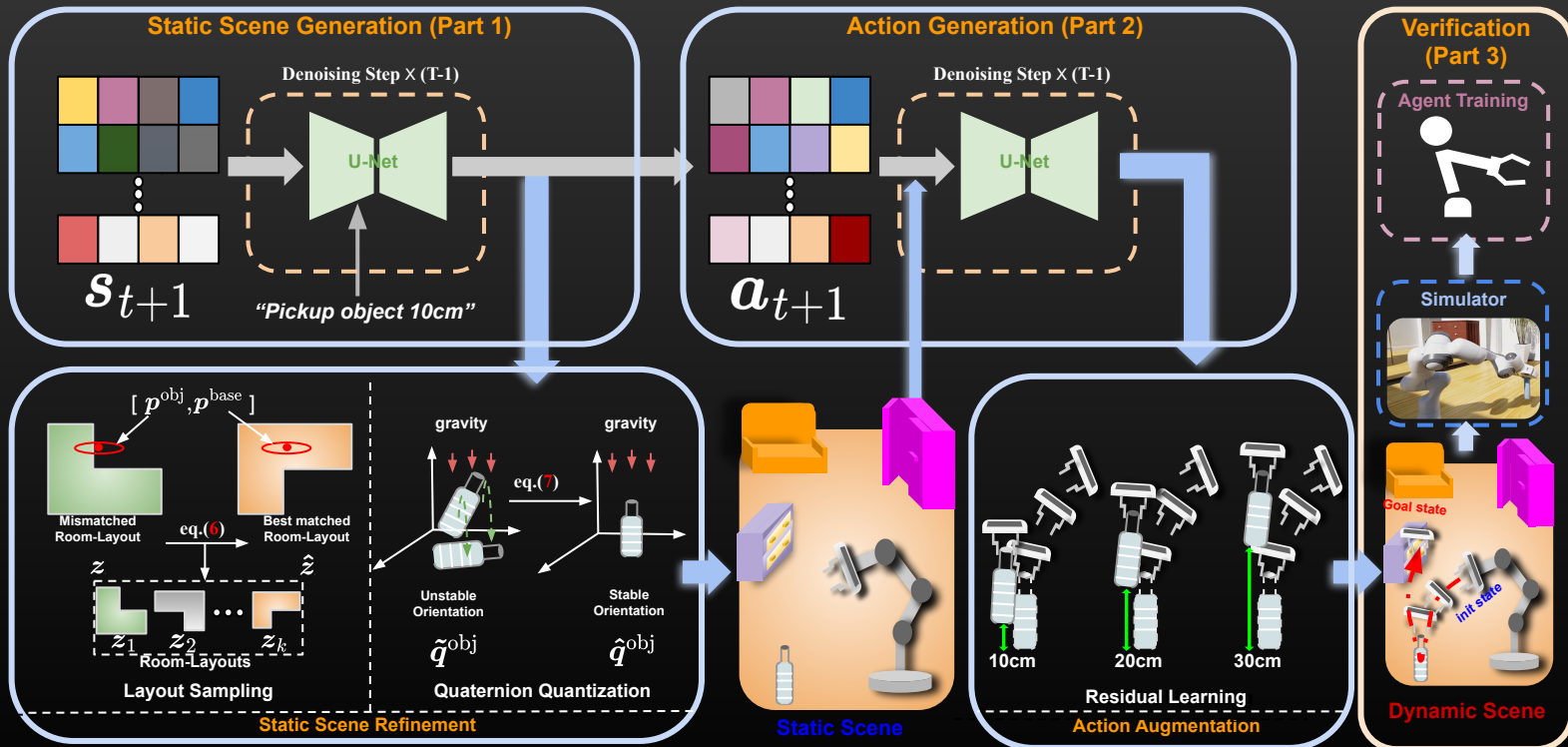
- In this challenge, we investigate:

We explore strategies to generate training data **effective** for agent learning.

- Our solution: **DynScene** [1] - robotic manipulation scene generation method enabling **scalable** dataset creation.

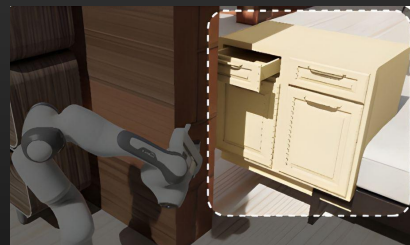
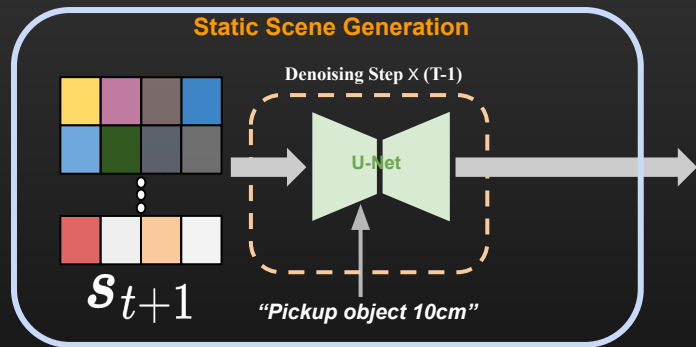
Proposed Method : DynScene

Framework Overview



Part 1: Static Scene Generation

- Diffusion model effectively generates **realistic** static scenes aligned with instructions.
- However, requires refinement for **physical plausibility** and **task accessibility**.



Unable to reach the target



Robot-Layout collision



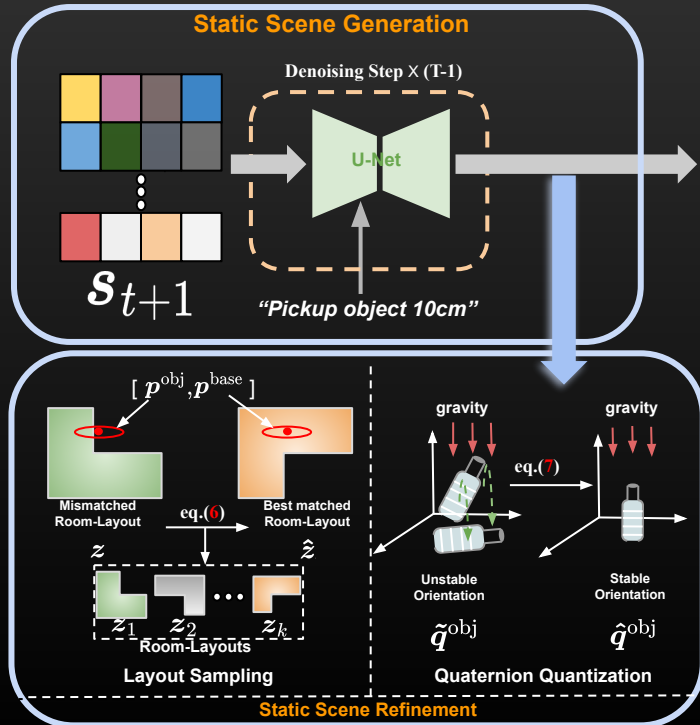
Unstable Orientation



Unintended Opening

Part 1: Static Scene Generation

- Diffusion model effectively generates **realistic** static scenes aligned with instructions.
- However, requires refinement for **physical plausibility** and **task accessibility**.



- **Layout Sampling** (eq.6)

$$\hat{r} = \arg \min_r (\|p_r^{\text{obj}} - \tilde{p}^{\text{obj}}\|^2 + \|p_r^{\text{base}} - \tilde{p}^{\text{base}}\|^2)$$

Unable to reach the target

Robot-Layout collision

Gravity

- **Quaternion Quantization** (eq.7)

$$\hat{q}^{\text{obj}} = \text{round} \left(\frac{\tilde{q}^{\text{obj}}}{\delta} \right) \cdot \delta$$

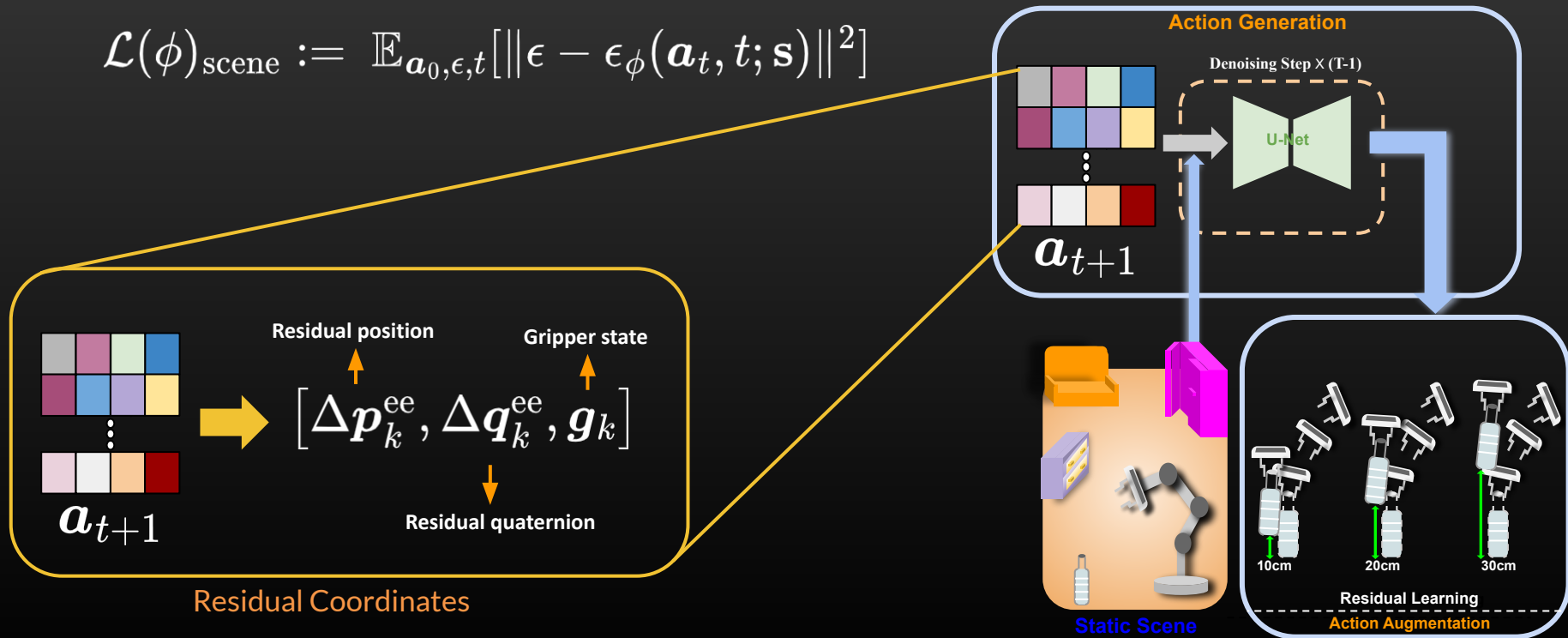
Unstable Orientation

Unintended Opening

Part 2: Action Generation and Augmentation

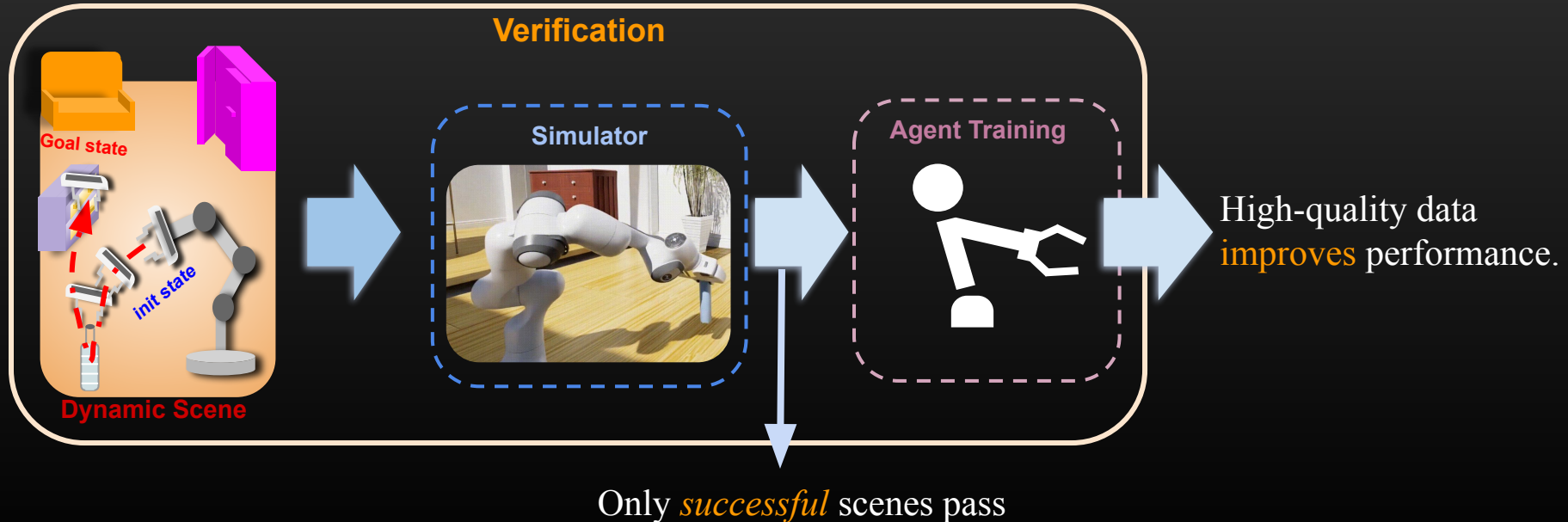
- Action generation uses diffusion model conditioned on static scenes.

$$\mathcal{L}(\phi)_{\text{scene}} := \mathbb{E}_{\mathbf{a}_0, \epsilon, t} [\|\epsilon - \epsilon_\phi(\mathbf{a}_t, t; \mathbf{s})\|^2]$$

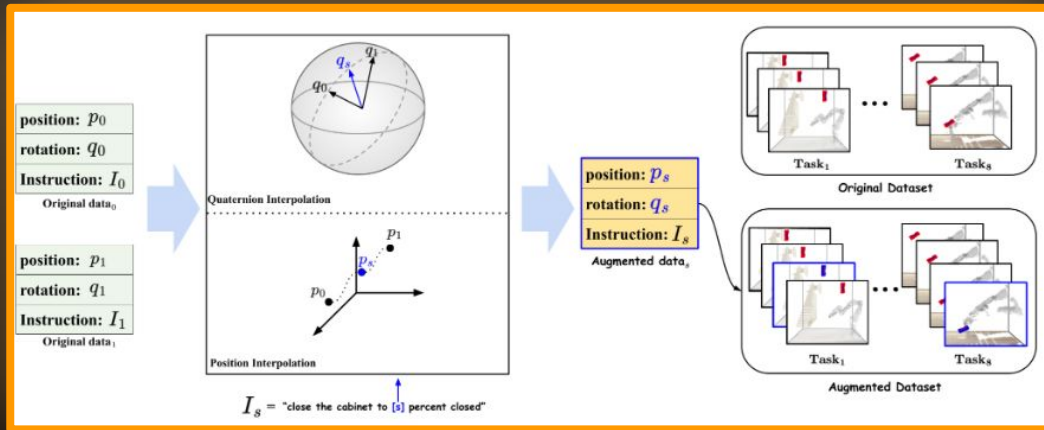


Part3: Filtering Invalid Dynamic Scenes

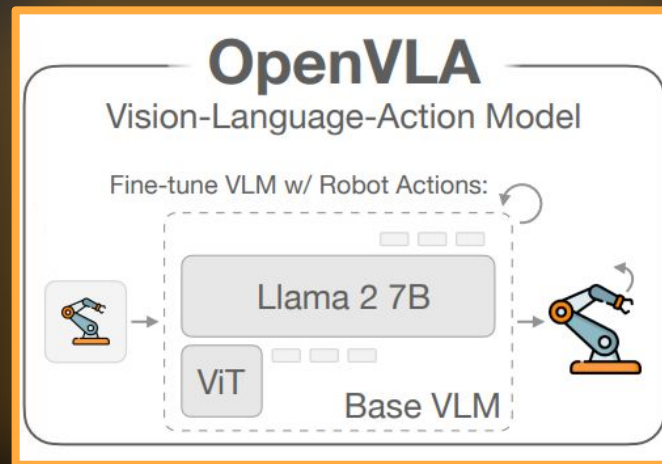
- Inaccurate actions can degrade agent performance and cause task **failures**.
 - Generated dynamic scenes undergo physics simulation verification in *NVIDIA Isaac Sim*.



Training OpenVLA with state interpolation augmentation

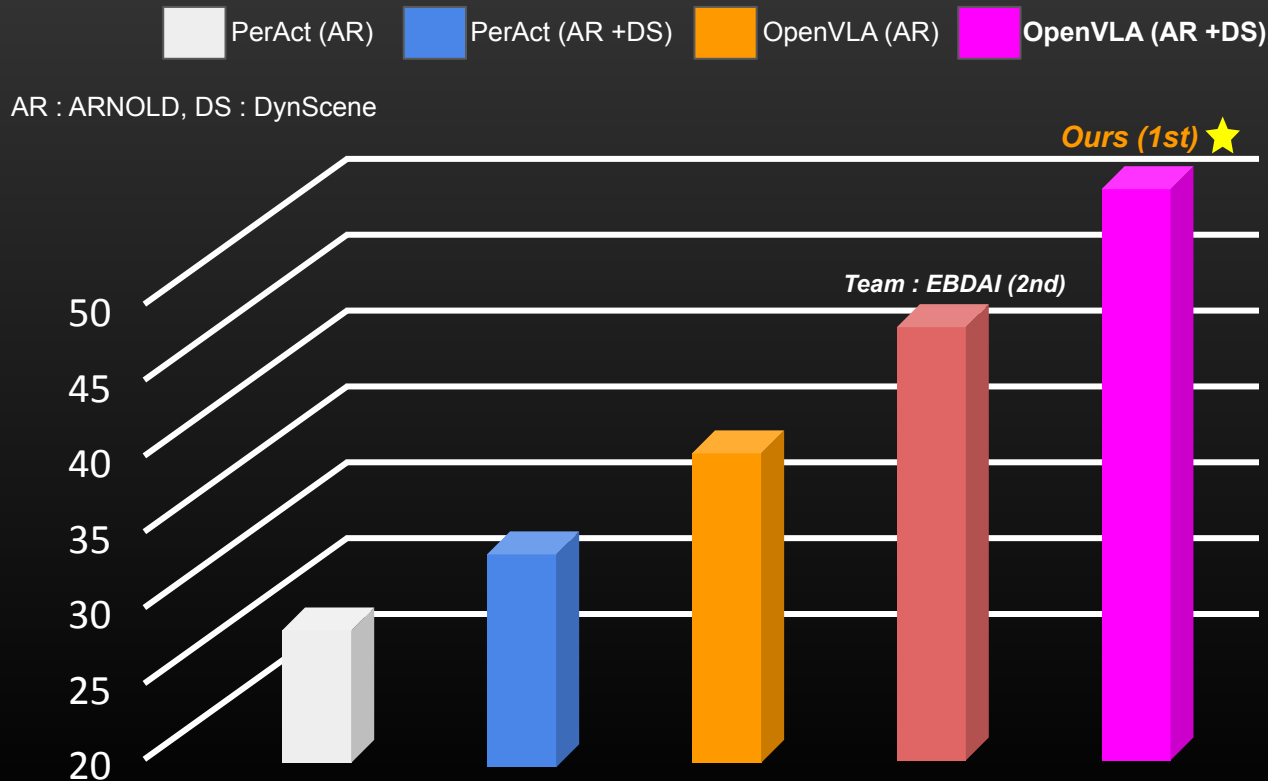


State Interpolation
(3rd place in 2024)



OpenVLA [2]

- We trained foundation models using DynScene-generated data.



Thank you

<https://reality.ssu.ac.kr/>