# DATA 8005 Advanced Natural Language Processing

## Towards Generalist Robot:

## How to Scaling up Robotics Dataset?

Feng Chen

Fall 2024

# Outline

- 1. Introduction & Background knowledge

- 2. Methodology Generalization

- 3. Imitate from Video

- 4. Generative Simulation

# Outline

- 1. Introduction & Background knowledge

  - what's our goal?

  - what's the state now?

  - what can we learn from NLP?

  - How to Improve?

- 2. Methodology Generalization

- 3. Imitate from Video

- 4. Generative Simulation

# What is our Goal?

General Artificial Intellegence is the final goal for every AI researchers

- Is large foundation model like GPT or LLaMa general artificial intellegence?

- NO!

- When talking about general AI, what is general?

- Robot is the first figure come into your mind

- Build a universal robot to solve productivity challenges is our final goal

# What is the State Now?

Although the field of robotics has made significant progress in the past decade

- The domain of robotics research is still in special skills

- Robots can only be set up in factory settings

- So what's the reason?

# What can we Learn from NLP?

Scaling up brings something!

- We need to settle three key point for scaling up
    - 1. Good dataset (Something like Image Net)
    - 2. Good model Structure (Transformer, Next token prediction)
    - 3. Enough Compute Resource or Platform (Simulator)

# How to Improve?

- 1. Good Dataset

  - Build large scale dataset for robotics

  - Using other datasets to improve

- 2. Good Structure

  - Diffusion policy?

  - Transformer?

- 3. Simulator

  - Issac Lab, Genesis…

# What are good datasets?

- 1. Real-world data

    - Positive: Realism, easy to interact

    - Negative: Hard to scale up, no gradient, expensive

- 2. Simulation data

    - Positive: Easy to parallel train, have gradient, cheap

    - Negative: Not realism, gap to real-world, Hard to scale up

- 3. Video data

    - Positive: Realism, easy to scale up, cheap

    - Negative: Hard to interact, no gradient, no physics

# Methodology Generalization

- 1. For generalization
  - LLM/VLM for reward
  - Vision-Language-Action model
- 2. Learn from video
  - Latent Action Pretraining from Videos
  - Hand-object interaction pretraining from videos
- 3. Diffusion Polices
  - Different Model Structure

# Outline

# Basic Knowledge: Model of Robotics

The Closed-Loop Interaction Model

- How to learn this function: get observation, provide reasonable action (**POLICY**)

Observation at **t**

action at **t**

Observation at **(t+1)**

■

■

■

Other Influence factor at **t**

# Basic Methodology



- Where to learn (Data)

  - Real World

    - Control Robots to do tasks, collect sensor data for later learning

    - (For fun) "Reinforcement Learning" in real world

      - [Learning to Walk in the Real World in 1 Hour (No Simulator)www.youtube.com › watch](#)

    - Problem

      - Expensive (Buy Equipment)

      - Inefficient Data Collection

        - Build Environment; Incapability of Parallel.

# Basic Methodology



- Where to learn (Data)

  - Simulator (Chen, Feng will dig into details)

    - Build a virtual environment using software (PC games)

    - Learn the policy using the virtual environment

    - Adopt to real world (Sim-to-real)

    - Benefit:

      - Cheaper, Easy to get equipments

      - Parallel-able (Just several process in your OS…)

# Basic Methodology

- How to learn (Methods)

  - Imitation Learning

  - Reinforcement Learning

  - (https://www.researchgate.net/figure/The-framework-of-Reinforcement-Learning-Imitation-Learning-and-their-integration-The_fig4_322094035)

# Where large model can involve

Zero Shot

- Make use of LLM/VLM's

    - interpretation of web-scale knowledge

    - reasoning capability **(?)**

- Form Reward, Hierarchical Planning, …

Fine-tuning LM to input/output action.

- Start for reasonable Web-scale trained checkpoints

- How to encode/decode action

# Outline

# Text to Policy

LLM generate rewards

- Human give language instruction, then translate it into reward function for RL
- https://eureka-research.github.io/
- https://text-to-reward.github.io/

LLM generate codes

- Human give language instruction, then translate it into constraint function
- https://arxiv.org/abs/2312.06408

# Limitation of Text to Policy

**Limitation**

- Hard-to-Access Ground Truth

    - environment code, low-level state data
- Limitations of Language/Code Descriptions

    - E.g. Describe the cloth →

**Direct Vision Grounding is needed**

- **LLM →VLM**

- **Text to Policy →Vision to Policy**

# Image to Policy

What's the key intuition?

- Large Language model can generate reward function

- **VLM is stronger now!**

  - Vision feedback is more useful when generate reward

- Let's use VLM generate reward with language instruction and image

# DATA 8005 Advanced Natural Language Processing

## **RL-VLM-F**: Reinforcement Learning

## from Vision Language Foundation Model Feedback

Tutorial: Liu, Ruizhe

Fall 2024

# Overview of **RL-VLM-F**

Challenge

- **Reward engineering** in RL is labor-intensive, trial-and-error.

- CLIP Model Limitations

  - Produces noisy, high-variance signals, frequently requiring fine-tuning.

## Method: **RL-VLM-F**

- Auto-generates rewards from text goals and visual inputs via VLM feedback.

- Uses VLM to **rank observations**, learning rewards from preference labels. **(RL-[H]-F → RL-[VLM]-F)**
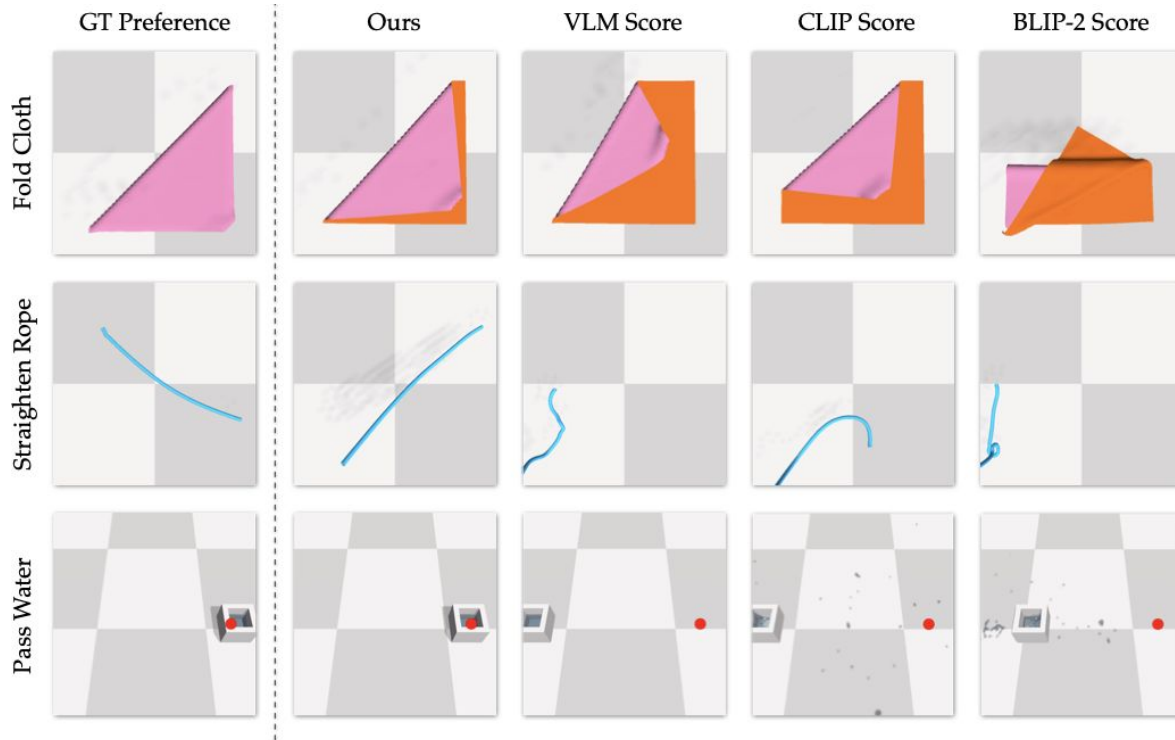
# Pipeline of **RL-VLM-F**

# VLM usage of Reward
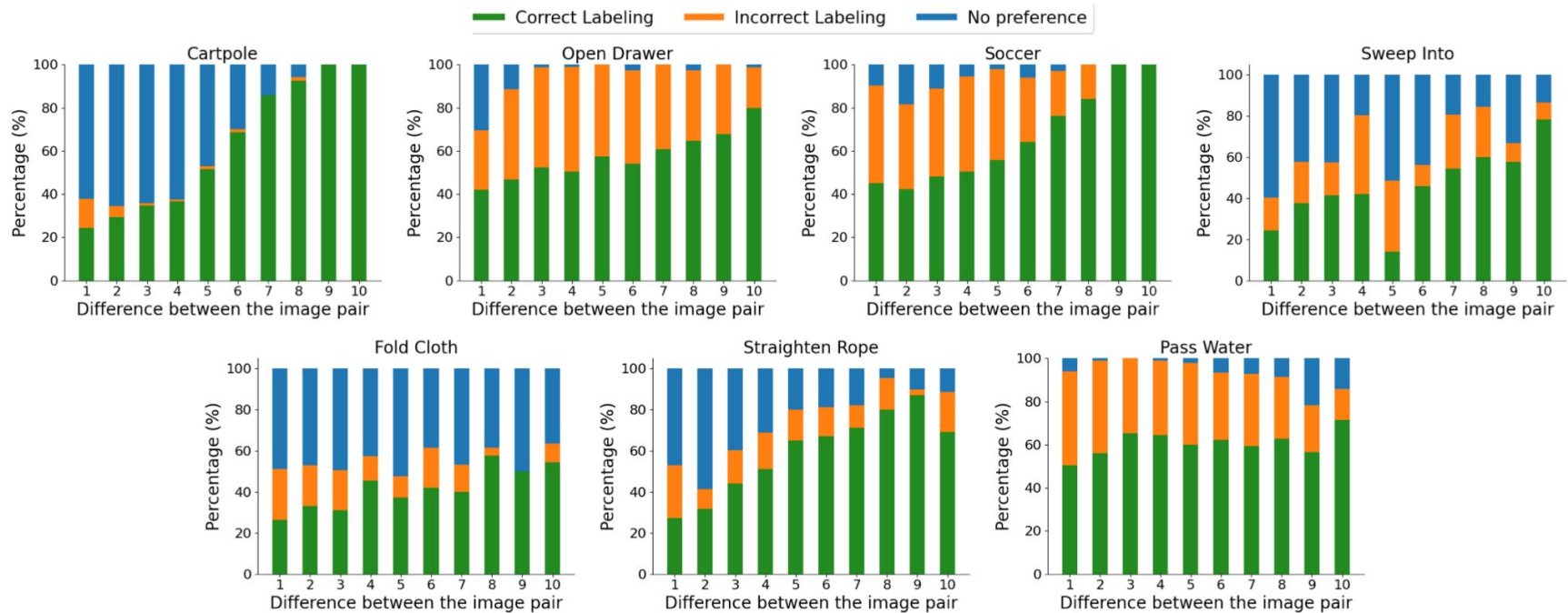
# Snap of Experiments

- Success Rates

# Snap of Experiments
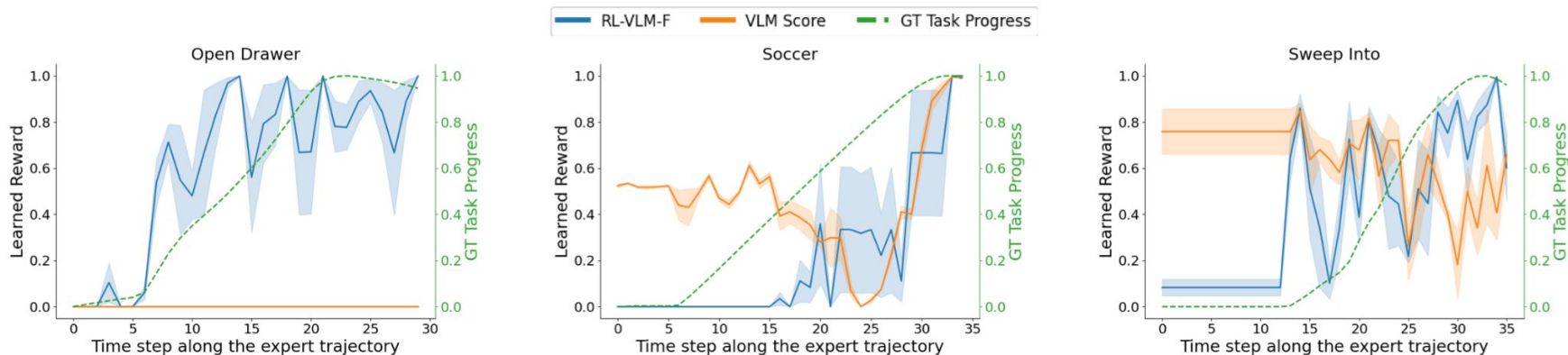
- Semantic Visualization

# Snap of Experiments

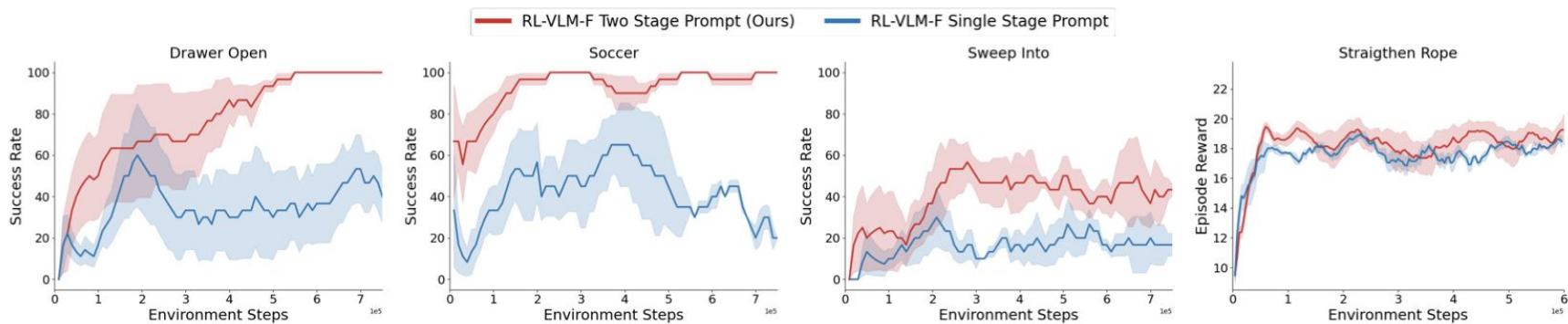- VLM labels vs. Ground Truth labels

# Snap of Experiments

- RL-VLM-F progress vs. Ground Truth progress

# Snap of Experiments

- Two Stage Prompt vs. Single Stage Prompt

# Summary of **RL-VLM-F**

Method: **RL-VLM-F**

- Auto-generates rewards from text goals and visual inputs via VLM feedback.

- Uses VLM to **rank observations**, learning rewards from preference labels. **(RL-[H]-F → RL-[VLM]-F)**

# Discussions

- Why RL-VLM-F seems to be a better structure than directly using similarity (e.g. CLIP)?

  Or more directly, why reward signal from CLIP is noisy?

- What else structure might be a good intuition?

# Limitation of (Zero-shot) Image to Policy

**Limitation**

- Limitation of "natural language"

**Dig out potential Embedding!**

# Embedding to Policy

What's the key intuition?

- Start from VLM

  - LLM has the potential of embedding latent action, but needs decoder to activate this embedding.
- Sequential Similarity between VLMs and Robotic policies

  - maximally leverage the original embedding

  - How to encode & decode.

Observation at $t$

action at $t$

Observation at $(t+1)$

Other Influence factor at $t$

"Grab the carrots"

$a_t$

Img emb          text emb          action emb

# DATA 8005 Advanced Natural Language Processing

# **OpenVLA**: An Open-Source **V**ision **L**anguage-**A**ction Model

Tutorial: Liu, Ruizhe

Fall 2024

# Background of **OpenVLA**

VLM Development

- Internet-scale vision-language data makes generalization possible

- Fine-tuning for downstream tasks adoption

  - Deep Learning / Prompt Engineering / …

Open X Embodiment



**Copy this way in Robotics :)**

- **Internet-scale** vision-language-action data makes generalization possible

  - Problem: Where is action? We lack robotic data

    - Open X Embodiment: 2,419,193

    - (Comparison) CLIP: 400,000,000

# Overview of **OpenVLA**

**Challenge**

- Lack Robotic Data (Learn from Scratch is hard)

- existing VLAs are largely closed and inaccessible to the public

- prior work fails to explore methods for efficiently **fine-tuning** VLAs for new tasks, a key component for **adoption**.
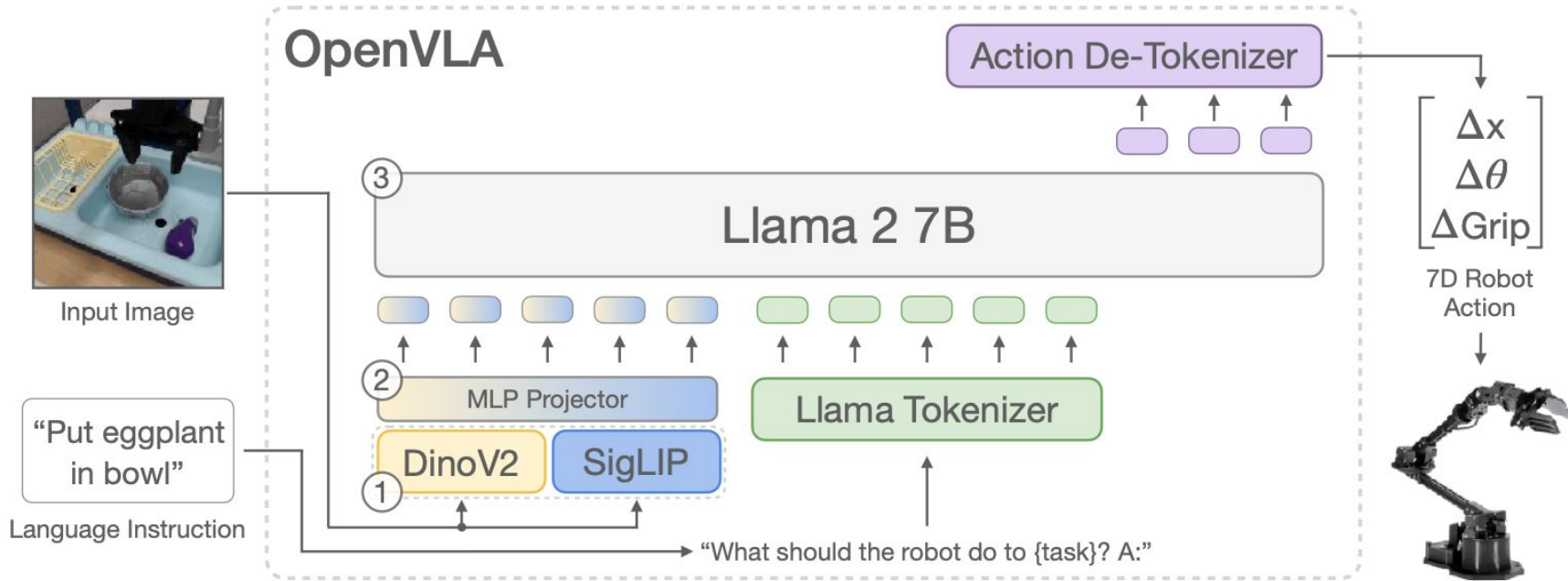
**OpenVLA**

- LLama 2 Based Transformer & **Leverage Pretrained LVM, LLM, LVLM...**

- Trained on **970K** real world robot demonstrations. (fine-tuning #1)

- Adoption Study (fine-tuning #2)

# Overview of **OpenVLA**

**How well**

- Achieves **high task success rate**, outperforming closed models (e.g., RT-2-X by 16.5% across 29 tasks, with 7x fewer parameters.)
- **Strong generalization(?)** in multi-task and multi-object environments; surpasses imitation learning methods like Diffusion Policy by 20.4%.
- **Fine-tuning** supported on consumer GPUs via low-rank adaptation; efficient deployment with quantization.
- **Open resources**: model checkpoints, fine-tuning notebooks, PyTorch codebase, and Open X-Embodiment dataset support.

# Pipeline of **OpenVLA**

# Training Data of **OpenVLA**

## Open X-Embodiment & curation

- At least one 3rd person camera
- Single-arm end-effector control.
- Data mixture weights
  - "although at a conservative mixture weight of 10%. In practice, we found that the **action token accuracy on DROID remained low throughout training, suggesting a larger mixture weight or model may be required to fit its diversity in the future**. To not jeopardize the quality of the final model, we **removed DROID** from the data mixture for the final third of training."

| OpenVLA Training Dataset Mixture | |
|---|---|
| Fractal [92] | 12.7% |
| Kuka [45] | 12.7% |
| Bridge[6, 47] | 13.3% |
| Taco Play [93, 94] | 3.0% |
| Jaco Play [95] | 0.4% |
| Berkeley Cable Routing [96] | 0.2% |
| Roboturk [97] | 2.3% |
| Viola [98] | 0.9% |
| Berkeley Autolab UR5 [99] | 1.2% |
| Toto [100] | 2.0% |
| Language Table [101] | 4.4% |
| Stanford Hydra Dataset [102] | 4.4% |
| Austin Buds Dataset [103] | 0.2% |
| NYU Franka Play Dataset [104] | 0.8% |
| Furniture Bench Dataset [105] | 2.4% |
| UCSD Kitchen Dataset [106] | <0.1% |
| Austin Sailor Dataset [107] | 2.2% |
| Austin Sirius Dataset [108] | 1.7% |
| DLR EDAN Shared Control [109] | <0.1% |
| IAMLab CMU Pickup Insert [110] | 0.9% |
| UTAustin Mutex [111] | 2.2% |
| Berkeley Fanuc Manipulation [112] | 0.7% |
| CMU Stretch [113] | 0.2% |
| BC-Z [55] | 7.5% |
| FMB Dataset [114] | 7.1% |
| DobbE [115] | 1.4% |
| DROID [11] | 10.0%[6] |

# Other Design & Feature of OpenVLA

**Start from BridgeData V2 for design decision**

- DINOv2 provides Stronger Spatial capability, making Prismatic > IDEFICS-1 and LLaVA.
- High Resolution seems provide no help (384x384 & 224x224), but needs more token… DISCARD!
- FINETUNE Vision Encoder…
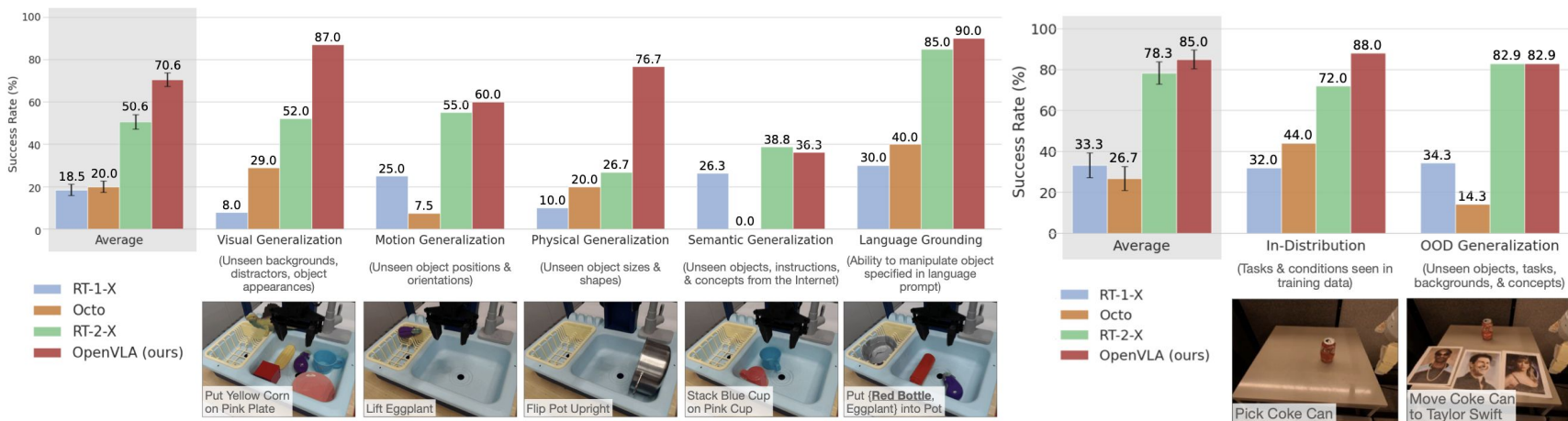- 27 epochs through training dataset (Our guess: Robotics data is not enough…)

Hardware

- 64 x A100 x 14 days x 2048
- 15GB Inference bfloat16 (without quantization), 6Hz on RTX 4090

# Snap of Experiments
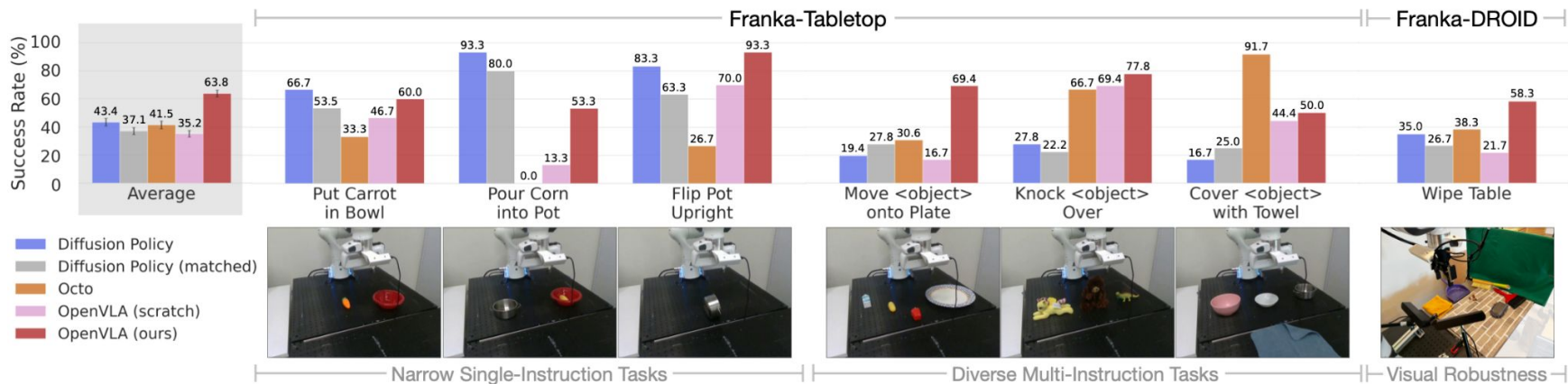
**Fine-tuning #1 Generalization**

- visual (unseen backgrounds, distractor objects, colors/appearances of objects)
- motion (unseen object positions/orientations)
- physical (unseen object sizes/shapes)
- semantic (unseen target objects, instructions, and concepts from the Internet) generalization.
- language conditioning ability of multiple objects, testing whether the policy can manipulate the correct target object, as specified in the user's prompt.
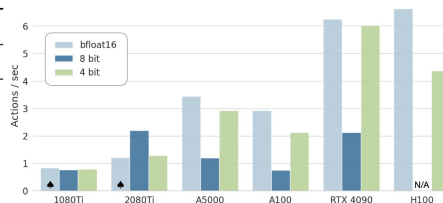
# Snap of Experiments

## Fine-tuning #2: Efficient Adoption



| Strategy | Success Rate | Train Params ($\times 10^6$) | VRAM (batch 16) |
|---|---|---|---|
| Full FT | **69.7 ± 7.2 %** | 7,188.1 | 163.3 GB* |
| Last layer only | 30.3 ± 6.1 % | 465.1 | 51.4 GB |
| Frozen vision | 47.0 ± 6.9 % | 6,760.4 | 156.2 GB* |
| Sandwich | 62.1 ± 7.9 % | 914.2 | 64.0 GB |
| LoRA, rank=32 | 68.2 ± 7.5% | **97.6** | **59.7 GB** |
| rank=64 | 68.2 ± 7.8% | 195.2 | 60.5 GB |

| Precision | Bridge Success | VRAM |
|---|---|---|
| bfloat16 | 71.3 ± 4.8% | 16.8 GB |
| int8 | 58.1 ± 5.1% | 10.2 GB |
| int4 | 71.9 ± 4.7% | 7.0 GB |

# Summary of **OpenVLA**

**OpenVLA**

- LLama 2 Based Transformer **& Leverage Pretrained LVM, LLM, LVLM…**

- Trained on **970K** real world robot demonstrations. (fine-tuning #1)
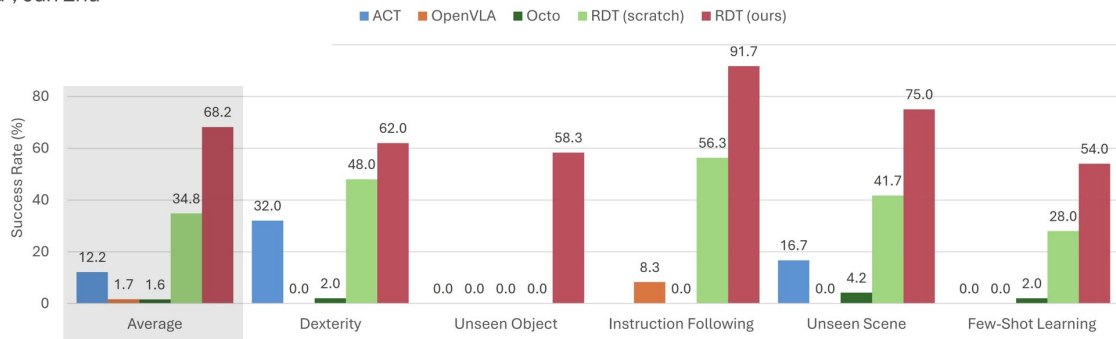
- Adoption Study (fine-tuning #2)

# Discussions

- Does open-source of OpenVLA really help? Since it is consuming (and controversial) to rigidly evaluate generalization in **real world**.

- Generalization?

## RDT-1B: a Diffusion Foundation Model for Bimanual Manipulation

Songming Liu[*,1], Lingxuan Wu[*,1], Bangguo Li[1], Hengkai Tan[1],
Huayu Chen[1], Zhengyi Wang[1], Ke Xu[1], Hang Su[1], Jun Zhu[1]
[1]Tsinghua University
[*]denotes equal contribution

# Outline

- 1. Introduction & Background knowledge

- 2. Methodology Generalization

- **3. Imitate from Video**

    - Latent Action Pretraining from Videos

    - Hand-object Interaction Pretraining from Videos

    - Discussion

- 4. Generative Simulation

# Latent Action Pretraining from Videos

Key Intuition:

- Record Observation (e.g. Video) is easy. Label Action for robots is annoying.

- How to make use of only observation.

  - **Generate Action** from it!

- Video as Vision Observation

  - Video from internet is much easier to collect than robotics dataset

  - Vision-Language-Action model could be pre-trained separately

  - Action prediction is good at generalization

# DATA 8005 Advanced Natural Language Processing

## **LAPA**: Latent Action Pretraining from Videos

Tutorial: Liu, Ruizhe

Fall 2024

# Overview of **LAPA**

Challenge

- Current VLA models rely on action labels from human teleoperators, limiting data sources and scalability.

Method: **LAPA**

- Leverages internet-scale videos without robot action labels
    - Train action quantizer (VQ-VAE) for discrete **latent actions**
- Pretrain VL[latent] A to predict **latent actions** from observations and task descriptions
- Finetune VL[latent]A on small robot manipulation data to map latent to robot actions
    - This latent action do not specify robot embodiment (One hand? Two hands? Legs? Dogs? Worms? Humanoids? Theoretically whatever in the video data is okay…)

Comment: An ambitious world model

# Why it is a World Model

Laws telling what will happen (Decoder)

$$o(T) = \int_{t=0}^{T} \frac{\mathrm{d}o(t)}{\mathrm{d}t} \mathrm{d}t = \int_{t=0}^{T} D_a(o(t), a(t)) \mathrm{d}t$$

Env. State

Env. Change ← Instant Env. Factors make changes

Action: Related to agent (robot)

Other factors make changes

# Overview of **LAPA**

How Well

- **Outperforms baselines** using actionless videos, especially in **cross-environment** and **cross-embodiment** tasks.
- LAPA effective **even with only human manipulation** video.
- Captures **environment-centric** actions (object/camera movement), aiding downstream tasks like navigation and dynamic tasks.

Comments: Action is essential; the actor is not.

$$o(T) = \int_{t=0}^{T} \frac{\mathrm{d}o(t)}{\mathrm{d}t} \mathrm{d}t = \int_{t=0}^{T} D_a(o(t), a(t)) \mathrm{d}t$$

Laws telling what will happen (Decoder)

Env. State    Env. Change    ←    Instant Env.    Factors make changes

Action: Related to agent (robot)

Other factors make changes

# Pipeline & Data



**1. Latent Action Quantization**

**2. Latent Pretraining**

Knock down the water bottle

Pick up the milk and put it in the sink

**Latent Action Pretraining** $\longrightarrow$ **Action Finetuning**

| Environment | Category | Pretraining | | Fine-tuning | |
|---|---|---|---|---|---|
| | | Dataset | # Trajs | Dataset | # Trajs |
| LangTable | In-Domain | Sim (All 5 tasks) | 181k | 5 Tasks (MT, MI) | 1k |
| | Cross-Task | Sim (All 5 tasks) | 181k | 1 Task (MI) | 7k |
| | Cross-Env | Real (All 5 tasks) | 442k | 5 tasks (MT, MI) | 1k |
| SIMPLER | In-Domain | Bridgev2 | 60k | 4 Tasks (MT) | 100 |
| | Cross-Emb | Something v2 | 220k | 4 Tasks (MT) | 100 |
| Real-World | Cross-Emb | Bridgev2 | 60k | 3 tasks (MI) | 450 |
| | Multi-Emb | Open-X | 970k | 3 tasks (MI) | 450 |
| | Cross-Emb | Open-X | 970k | 1 task (MI, Bi-manual) | 150 |
| | Cross-Emb | Something v2 | 220k | 3 tasks (MI) | 450 |

(a) LANGUAGE TABLE  (b) SIMPLER  (c) REAL
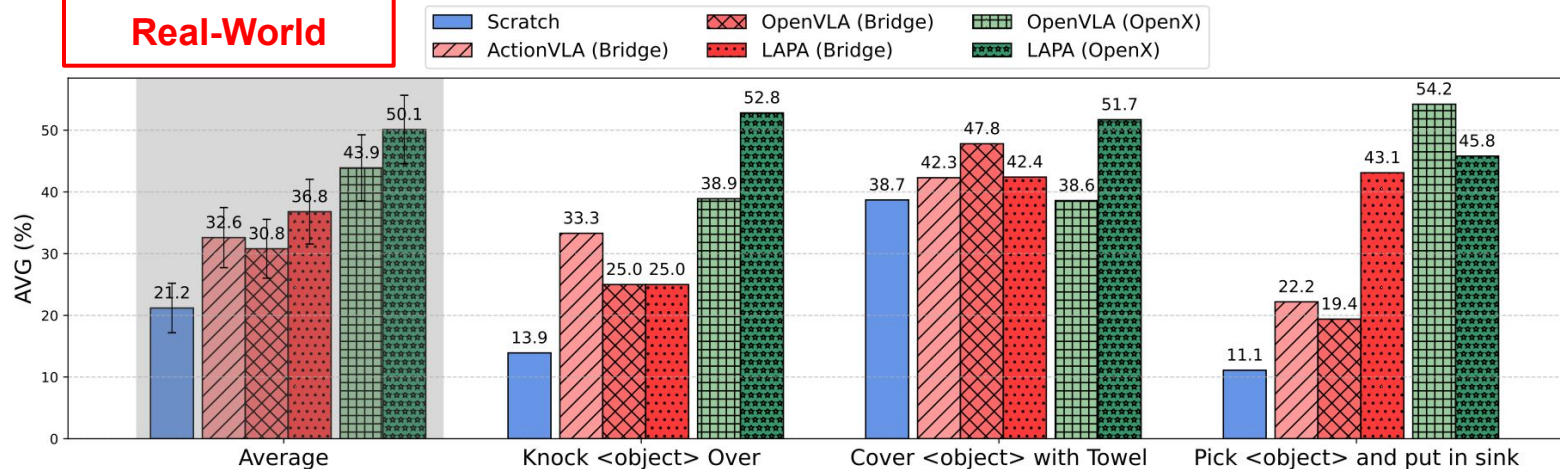
# Snap of Experiments

**LanguageTable**

| | In-domain (1k) | | Cross-task (7k) | | Cross-env (1k) | |
|---|---|---|---|---|---|---|
| | Seen | Unseen | Seen | Unseen | Seen | Unseen |
| SCRATCH | $15.6_{\pm9.2}$ | $15.2_{\pm8.3}$ | $27.2_{\pm13.6}$ | $22.4_{\pm11.0}$ | $15.6_{\pm9.2}$ | $15.2_{\pm8.3}$ |
| UNIPI | $22.0_{\pm12.5}$ | $13.2_{\pm7.7}$ | $20.8_{\pm12.0}$ | $16.0_{\pm9.1}$ | $13.6_{\pm8.6}$ | $12.0_{\pm7.5}$ |
| VPT | $44.0_{\pm7.5}$ | $32.8_{\pm4.6}$ | $72.0_{\pm6.8}$ | $\mathbf{60.8}_{\pm6.6}$ | $18.0_{\pm7.7}$ | $18.4_{\pm9.7}$ |
| LAPA | $\mathbf{62.0}_{\pm8.7}$ | $\mathbf{49.6}_{\pm9.5}$ | $73.2_{\pm6.8}$ | $54.8_{\pm9.1}$ | $\mathbf{33.6}_{\pm12.7}$ | $\mathbf{29.6}_{\pm12.0}$ |
| ACTIONVLA | $77.0_{\pm3.5}$ | $58.8_{\pm6.6}$ | $77.0_{\pm3.5}$ | $58.8_{\pm6.6}$ | $64.8_{\pm5.2}$ | $54.0_{\pm7.0}$ |

**SIMPLER**



**Real-World**

# Snap of Experiments: Human Video Only
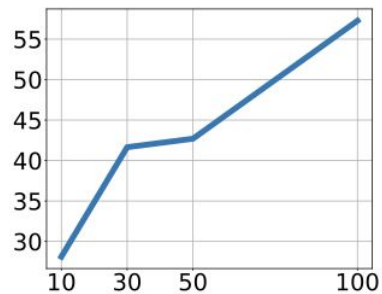


(a) SIMPLER Results
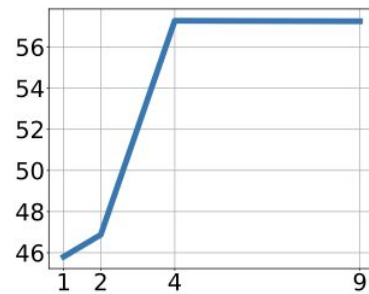
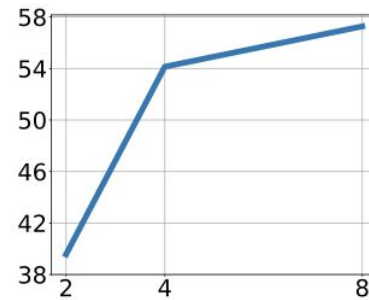(b) Real-world Tabletop Manipulation Robot Results

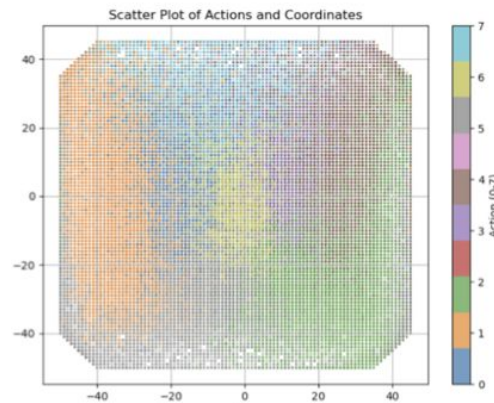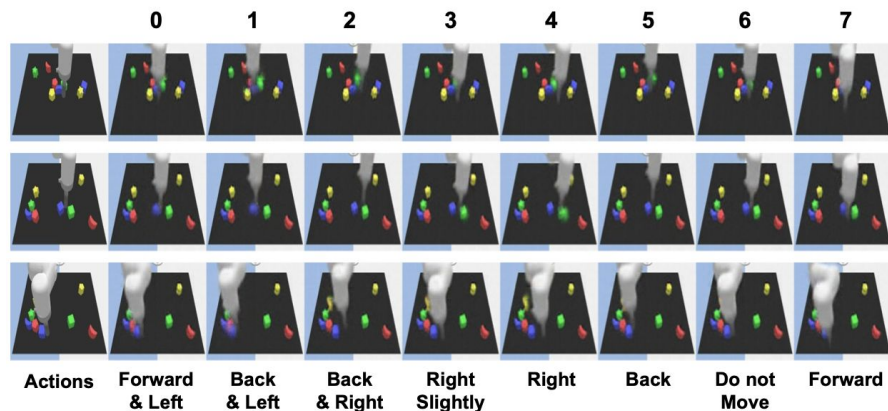# Snap of Experiments: Scaling & Beyond SR
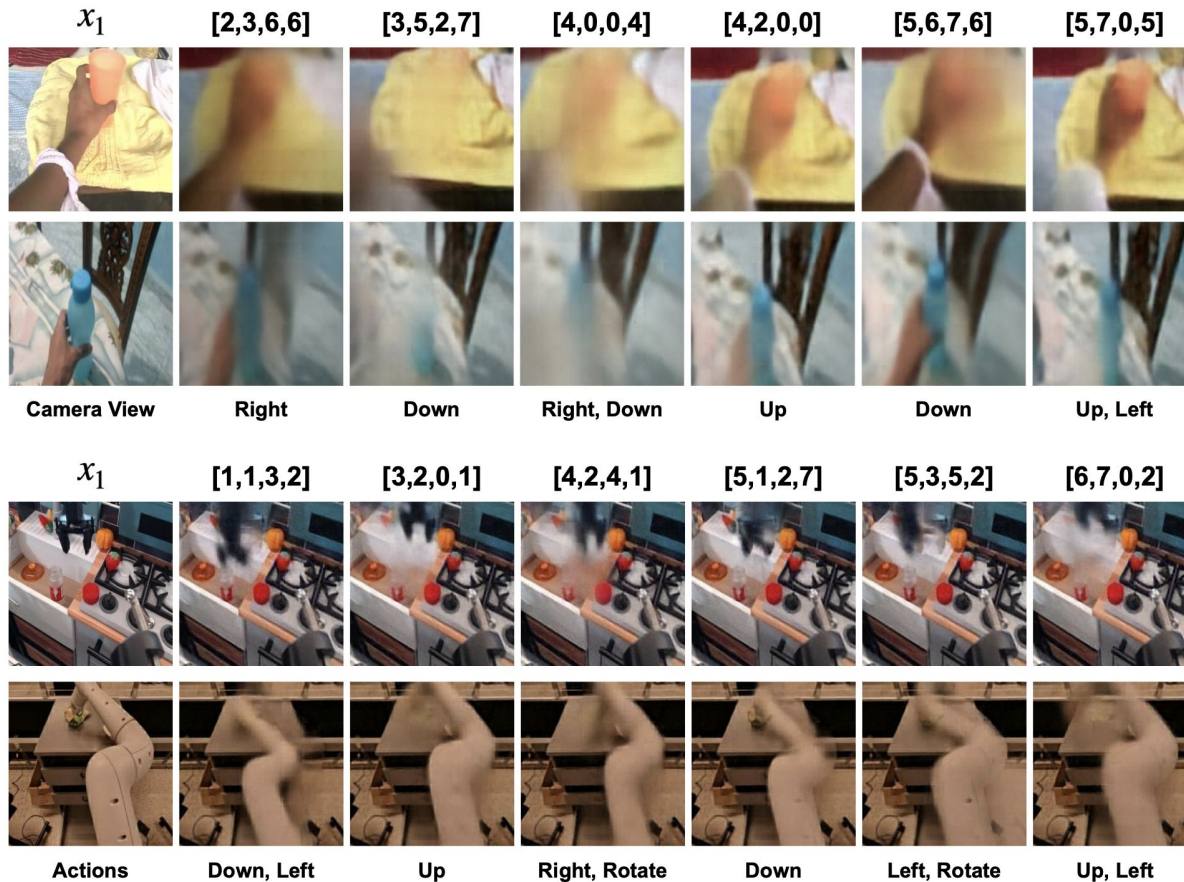


(a) Model Scaling

(b) Data Scaling (%)

(c) Latent Action Seq

(d) Latent Action Vocab

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| Actions | Forward & Left | Back & Left | Back & Right | Right Slightly | Right | Back | Do not Move | Forward |

Scatter Plot of Actions and Coordinates

Action (0-7)
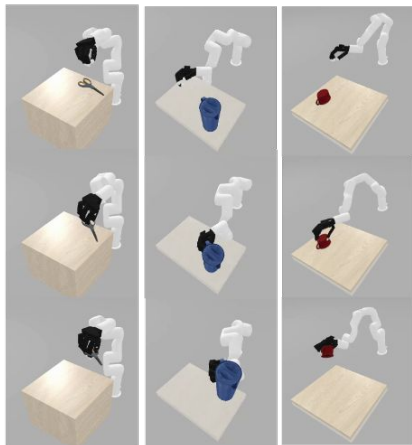
# Snap of Experiments: Latent actions & Camera Views
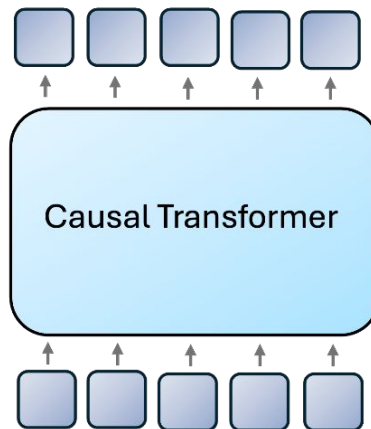
# Summary of **LAPA**

**LAPA**

- Leverages internet-scale videos without robot action labels
  - Train action quantizer (VQ-VAE) for discrete **latent actions**
- Pretrain VL[latent] A to predict **latent actions** from observations and task descriptions
- Finetune VL[latent]A on small robot manipulation data to map latent to robot actions
  - This latent action do not specify robot embodiment (One hand? Two hands? Legs? Dogs? Worms? Humanoids? Theoretically whatever in the video data is okay…)
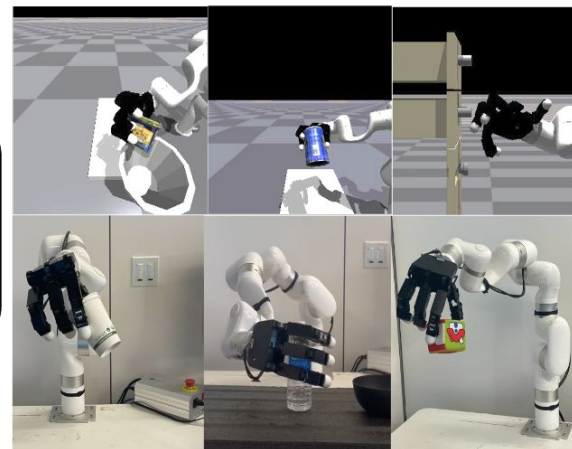
# Hand-object interaction pretraining from videos



3D hand-object trajectories    Sim-in-the-loop retargeting    Pretraining    Downstream adaptation

Causal Transformer

Another paper learn from video, detail could check on the website
https://hgaurav2k.github.io/hop/

# Discussions

- here is a statement from original paper:

  *"In the first pretraining stage, **we use a VQ-VAE-based objective to learn quantized latent actions** between raw image frames. **Analogous to Byte Pair Encoding** used for language modeling, this can be seen as learning to tokenize atomic actions without requiring predefined action priors."*

  Is there a connection between the two？

- here is a statement from original paper:

  *"By default, we freeze only the vision encoder and unfreeze the language model during training."*

  Since LAPA (claims that it) surpass OpenVLA, is OpenVLA wrong？

# Discussions

- Here is a statement from original paper:

  *"it opens the possibility of using any type of raw video paired with language instructions"*

  So why still videos that are carrying out action is necessary according to their experiment results?

# Discussions

- How to utilize video data better?

- A new idea: "Can we train two decode algin image prediction and action prediction?"

- Is latent space of video prediction same as action latent space?

# Summary of Methodology

| Methodology | Examples | Required Data |
|---|---|---|
| Zero-shot usage of VLM | Text-to-Policy family<br>RL-VLM-F | VL data (for VLM)<br>Specified robotic dataset |
| (F) Fine-tuning from VLM<br>(F) Fine-tuning adoption | OpenVLA | VL data (for VLM)<br>Various robotic dataset |
| (P) World Model<br>(F) Fine-tuning adoption | LAPA | Video data,<br>Specified robotic dataset |

# Outline

# Why is it hard to scaling up Robotics Dataset?

- Most time we require "Well-trained Ph.D. student" to generate robotics tasks

- Some papers tried to generate task through non-expert

  - https://arxiv.org/abs/2312.06408

- But it still requires manual design of tasks, training environment, algorithms, and supervision

- Can we generate task without human effort?

# Generative Simulation
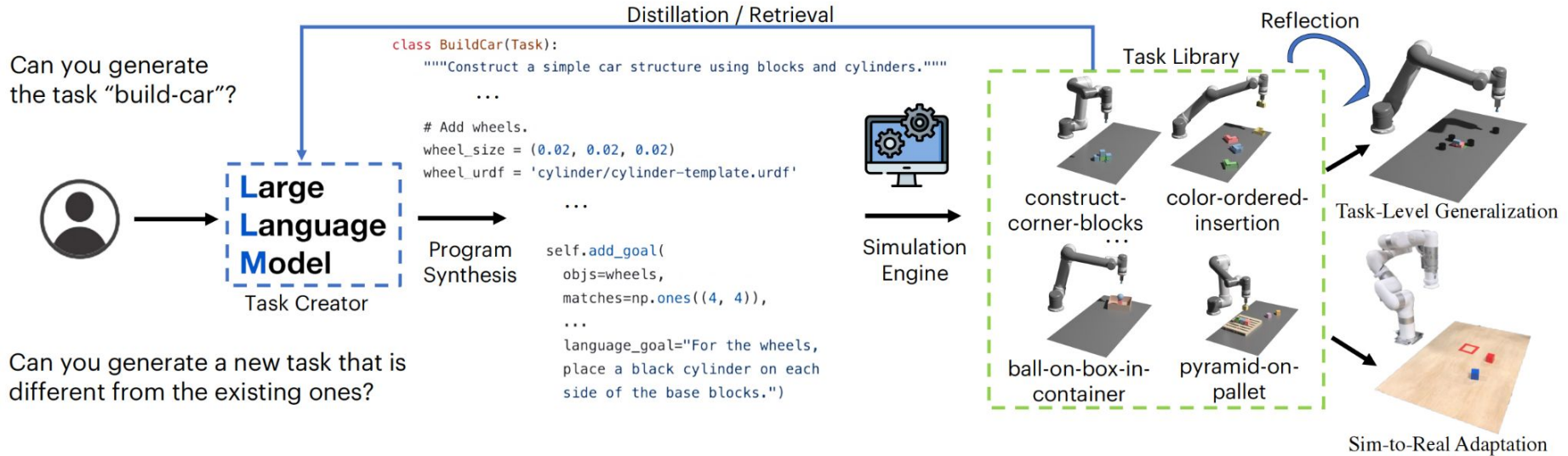
What's the key intuition?

- LLM and MLLM has shown the ability in space intelligence

- Robot tasks can usually be represented as formatted code files

- Language model could help we generate policies via different method

- https://arxiv.org/abs/2305.10455

# Gensim

# Gensim



LLM generate task file from knowledge in task library
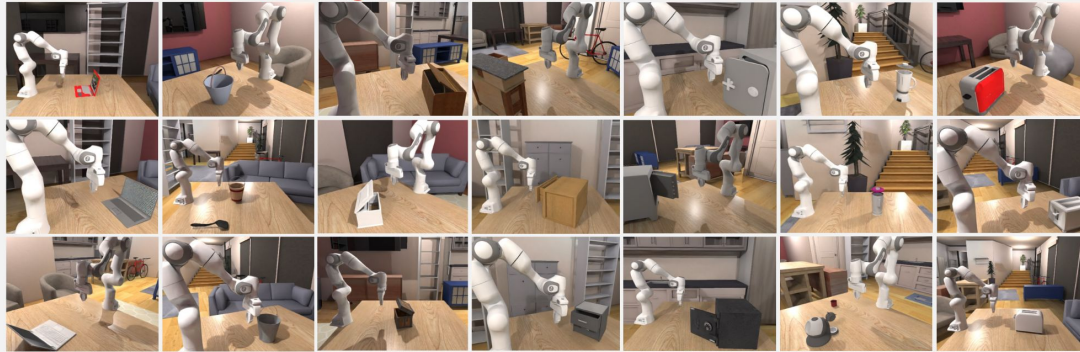Language-conditioned behavior cloning generate policy

# Gensim

Limitations:

- Tasks base on tasklib, diversity of task is kind of low

- Oracle learning sometime may not give correct policy

- All tasks are top-view table-top manipulation

- No good evaluation

# Gensim2



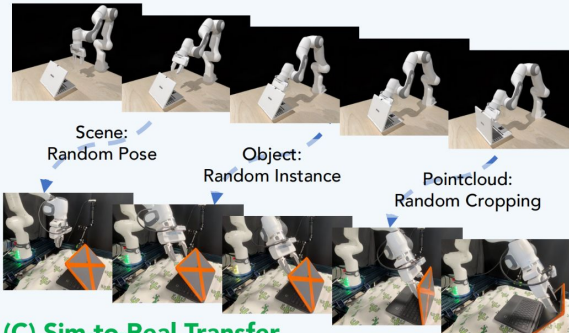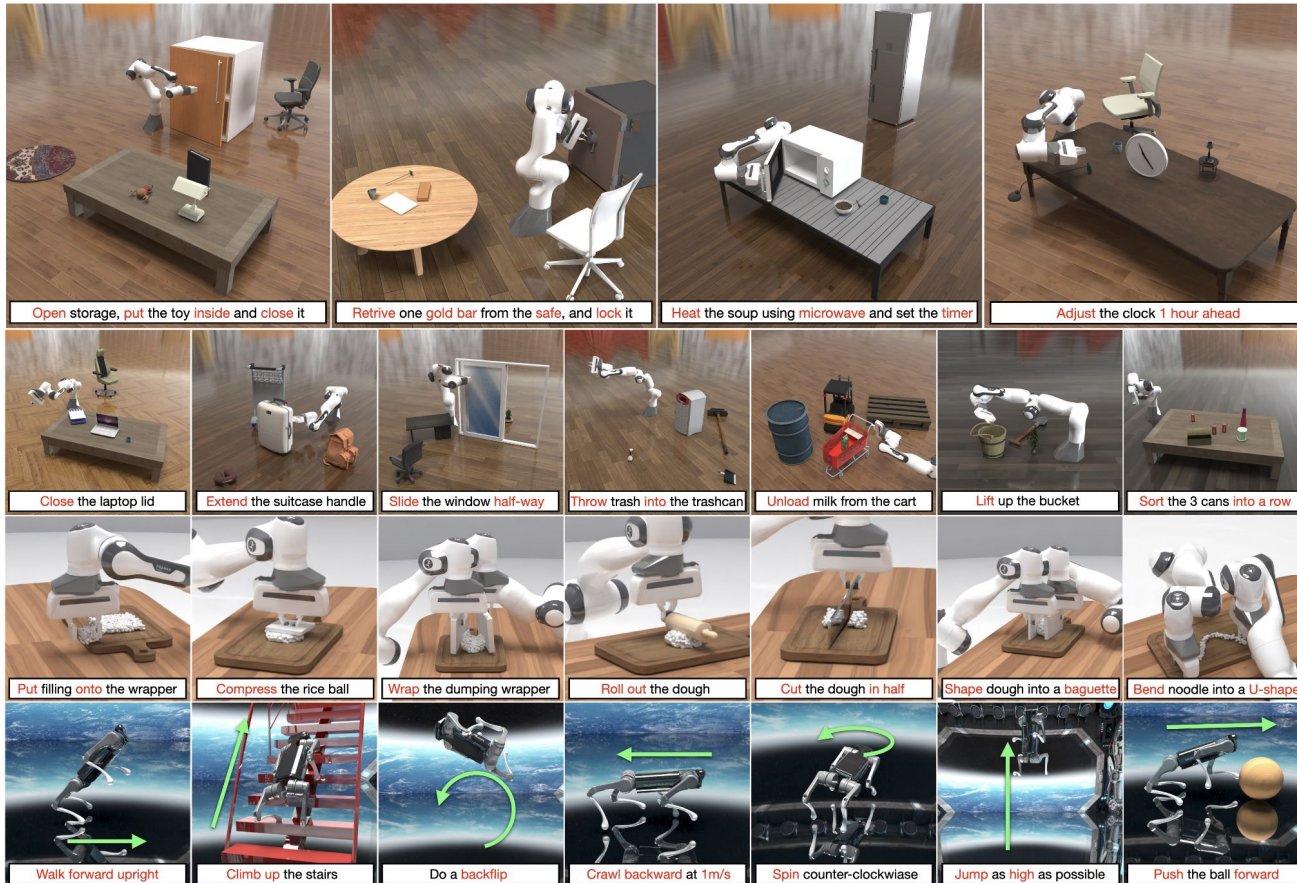**(A) Large-scale Task and Data Generation**

**(B) Multi-Task Training in Simulation**

Scene:
Random Pose

Object:
Random Instance

Pointcloud:
Random Cropping

**(C) Sim-to-Real Transfer**

https://arxiv.org/abs/2410.03645
New gensim!
No details welcome read after class!

# Robogen



Open storage, put the toy inside and close it | Retrive one gold bar from the safe, and lock it | Heat the soup using microwave and set the timer | Adjust the clock 1 hour ahead

Close the laptop lid | Extend the suitcase handle | Slide the window half-way | Throw trash into the trashcan | Unload milk from the cart | Lift up the bucket | Sort the 3 cans into a row

Put filling onto the wrapper | Compress the rice ball | Wrap the dumping wrapper | Roll out the dough | Cut the dough in half | Shape dough into a baguette | Bend noodle into a U-shape

Walk forward upright | Climb up the stairs | Do a backflip | Crawl backward at 1m/s | Spin counter-clockwiase | Jump as high as possible | Push the ball forward

# Robogen



**"Retrieve a gold bar from the safe"**
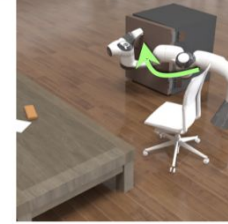
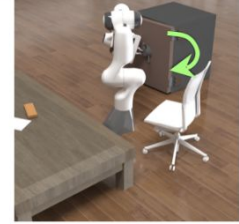Grasp the safe door — Open the safe door — Retrieve the gold bar — Move it to the table — Grasp the door again — Close the door — Rotate the knob to lock
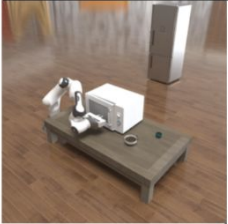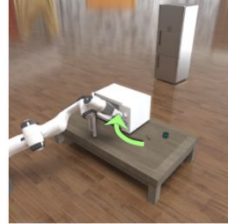
**"Heat up a bowl of soup using the microwave"**
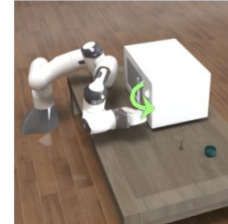
Approach the door — Open the door — Grasp the soup — Put it in the microwave — Close the door — Grasp the timer knob — Turn the knob to set timer

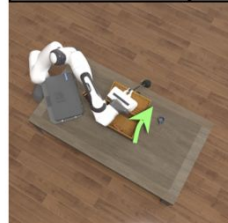**"Put the toy into the storage"**

Open the door — Grasp the toy — Move the toy inside

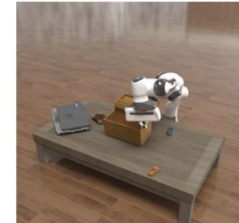**"Move the toy car out of the box"**
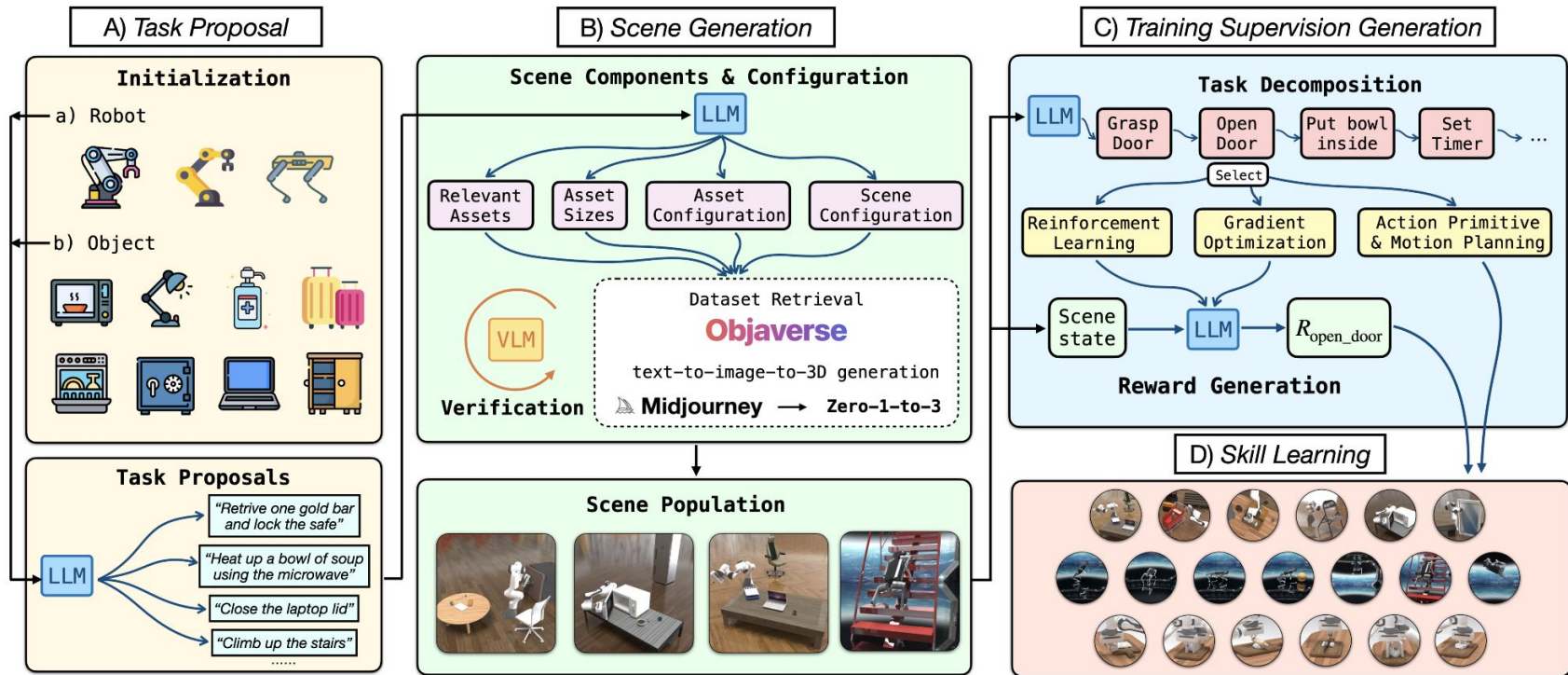
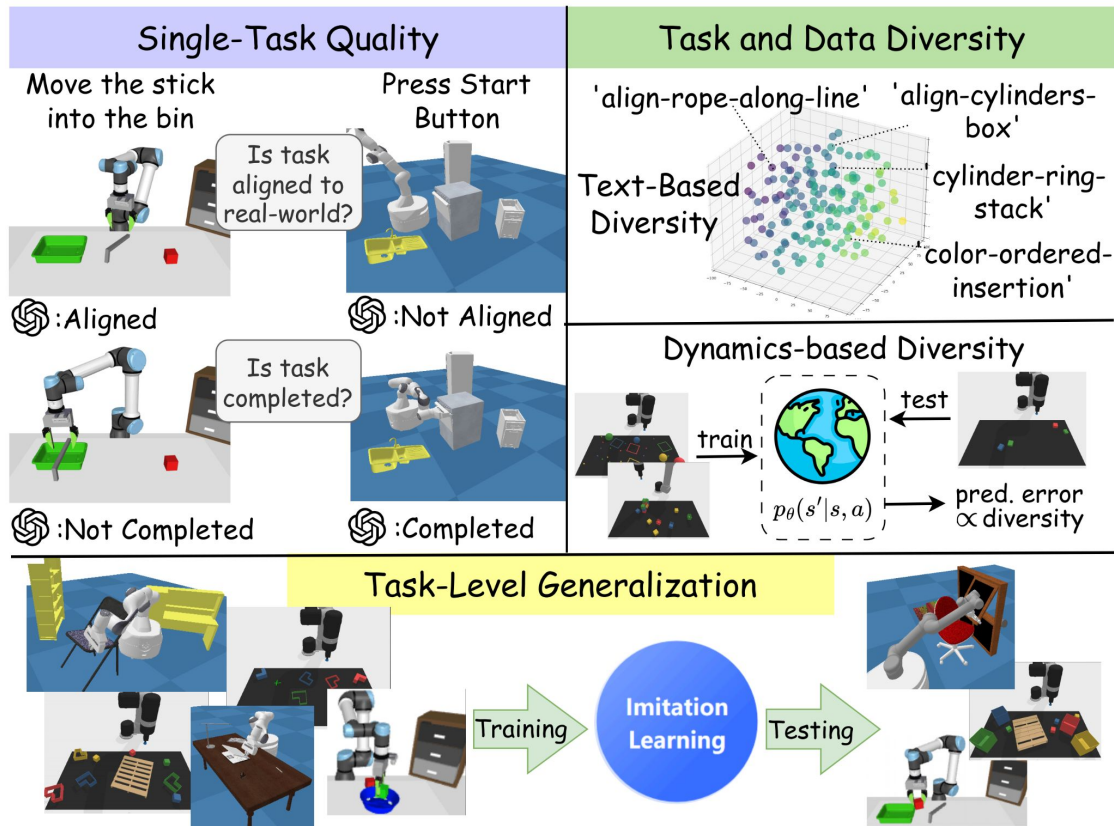Open the box — Retrieve the toy car — Move it out of the box — Release the toy car

# Robogen



A) Task proposal B) Scene generation C) Training via different methods

# Robogen

Limitations:

- Manipulation task can not generate diverse policy

- Most task can not generate good policy which can correctly solve the task

- Scene alignment sometime is not good
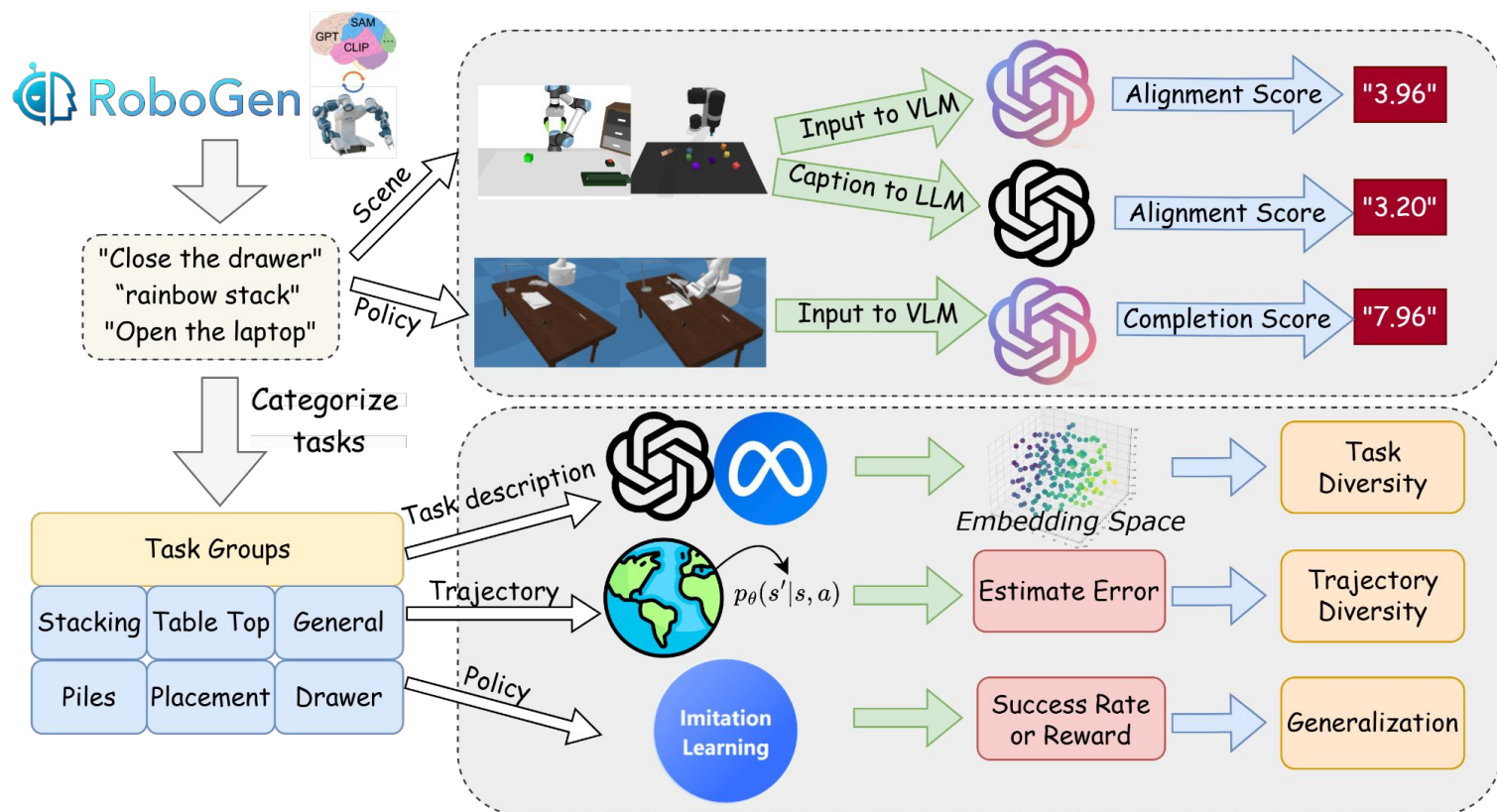
- No good evaluation

# On the Evaluation of Generative Robotic Simulations



We propose three main aspects for evaluating generative simulations: *Quality, Diversity, and Generalization*.

# On the Evaluation of Generative Robotic Simulations

# Other generative simulation works

- Auto RT: https://auto-rt.github.io/

- BBSEA: https://bbsea-embodied-ai.github.io/

- Gensim2: https://arxiv.org/abs/2410.03645

- RoboCasa: https://robocasa.ai/

# Discussions

- How can we achieve good generalization for generative simulation?

- Can generative simulation solve the thirsty of robotic data?

- What other evaluation method do you think is good metric?

- How can we achieve the final goal of embodied AI?