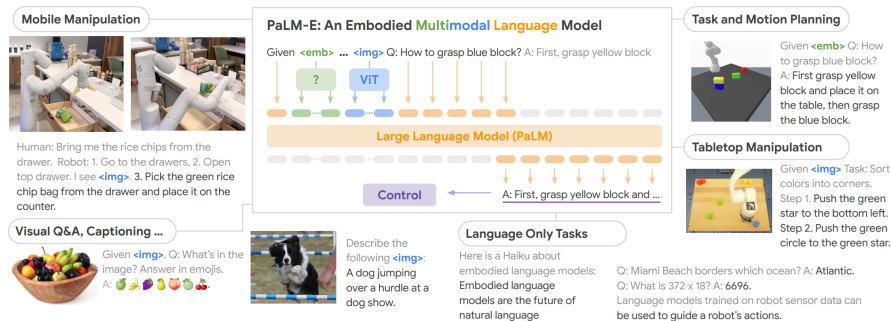
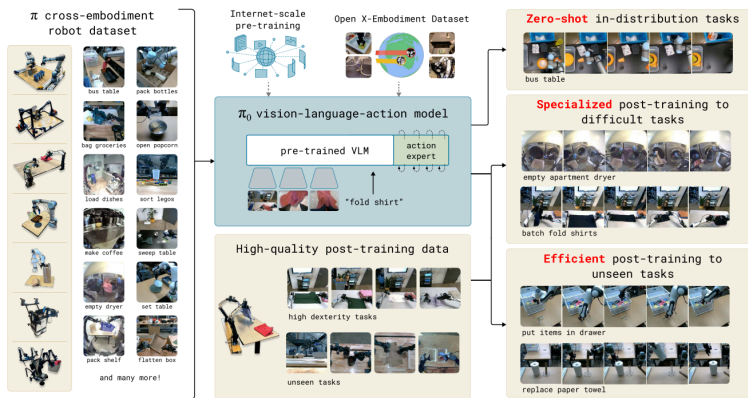


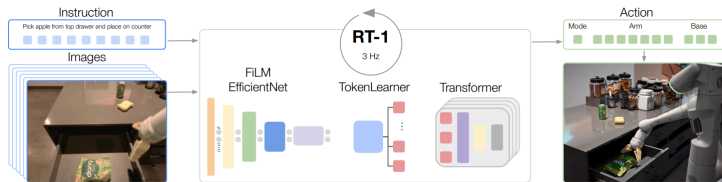


RobotArena ∞ : Scalable Robot Benchmarking via
Real-to-Sim Translation

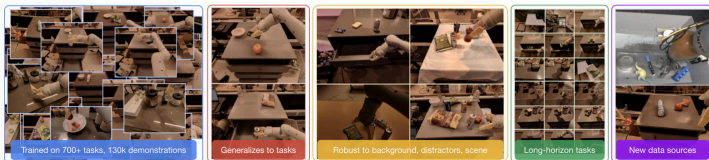
Robot Foundational Policies Need General Scalable Benchmarks



$\pi 0$: A Vision-Language-Action Flow Model for General Robot Control

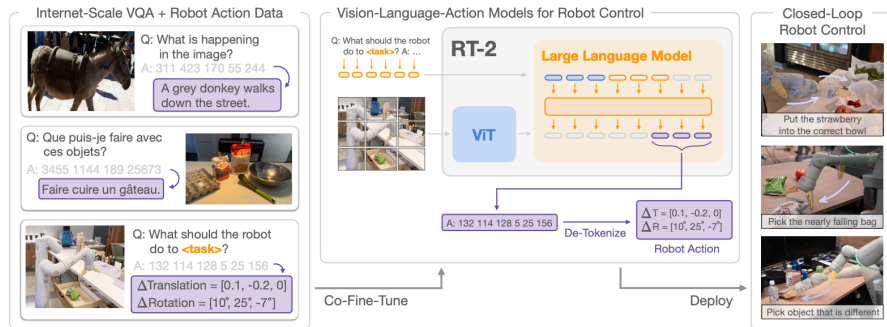


(a) RT-1 takes images and natural language instructions and outputs discretized base and arm actions. Despite its size (35M parameters), it does this at 3 Hz, due to its efficient yet high-capacity architecture: a FiLM (Perez et al., 2018) conditioned EfficientNet (Tan & Le, 2019), a TokenLearner (Ryoo et al., 2021), and a Transformer (Vaswani et al., 2017).



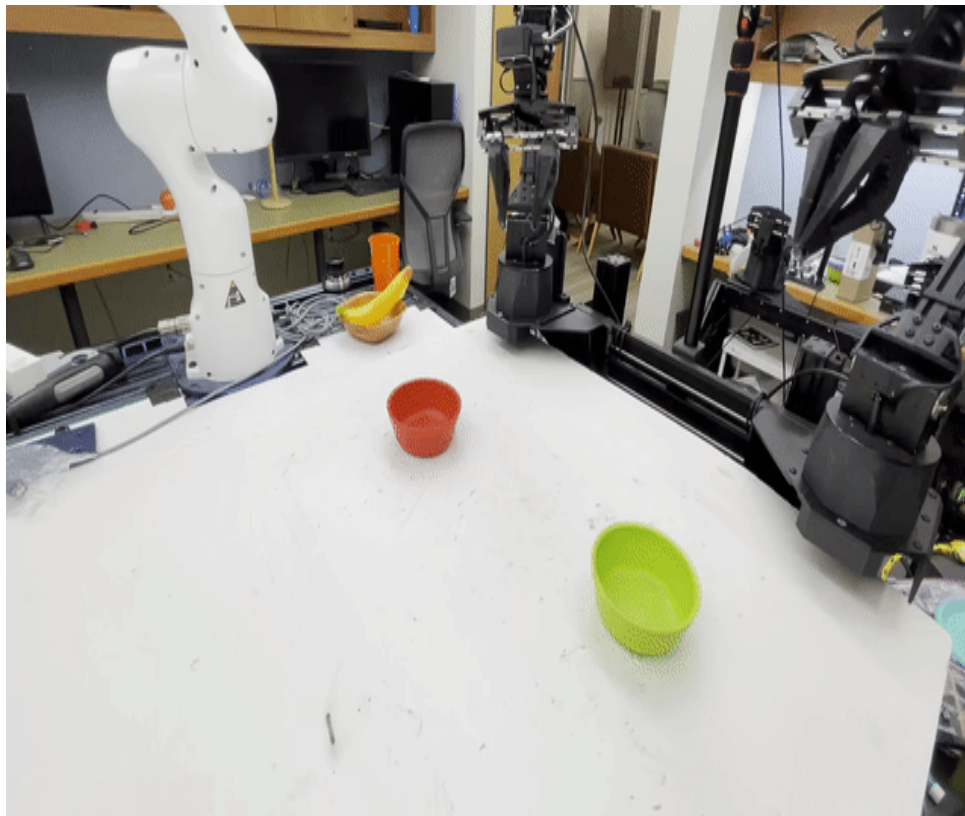
RT-1, Brohan et al., Dec. 2022

PALM-E, Driess et al., March 2023



RT-2, Brohan et al., July. 2023

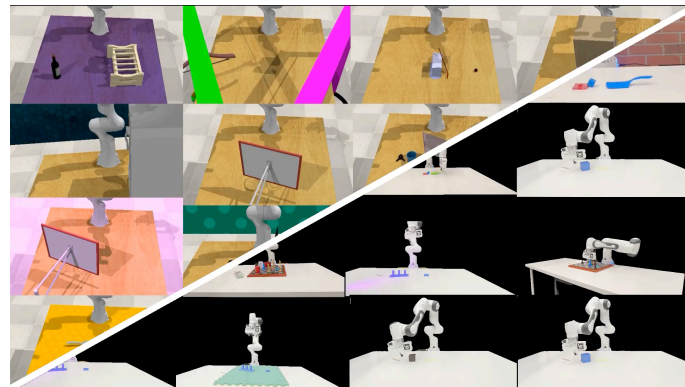
Evaluating robot policies in the real world does not scale



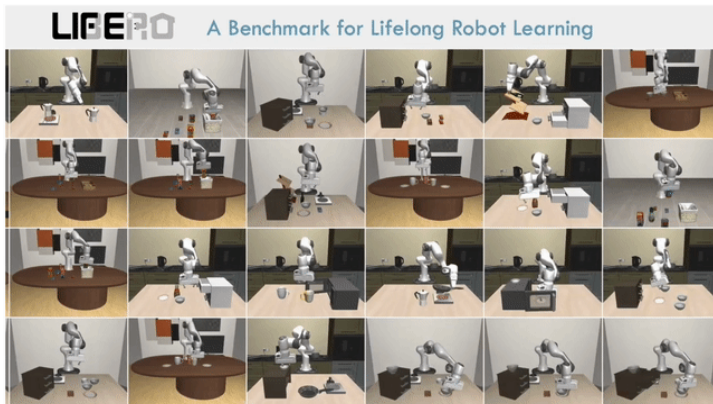
Simulation Benchmarks: Train and Test in Simulation



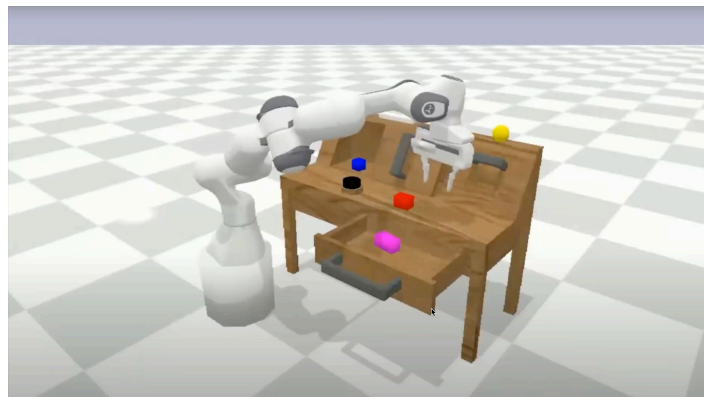
RL Bench



Colosseum

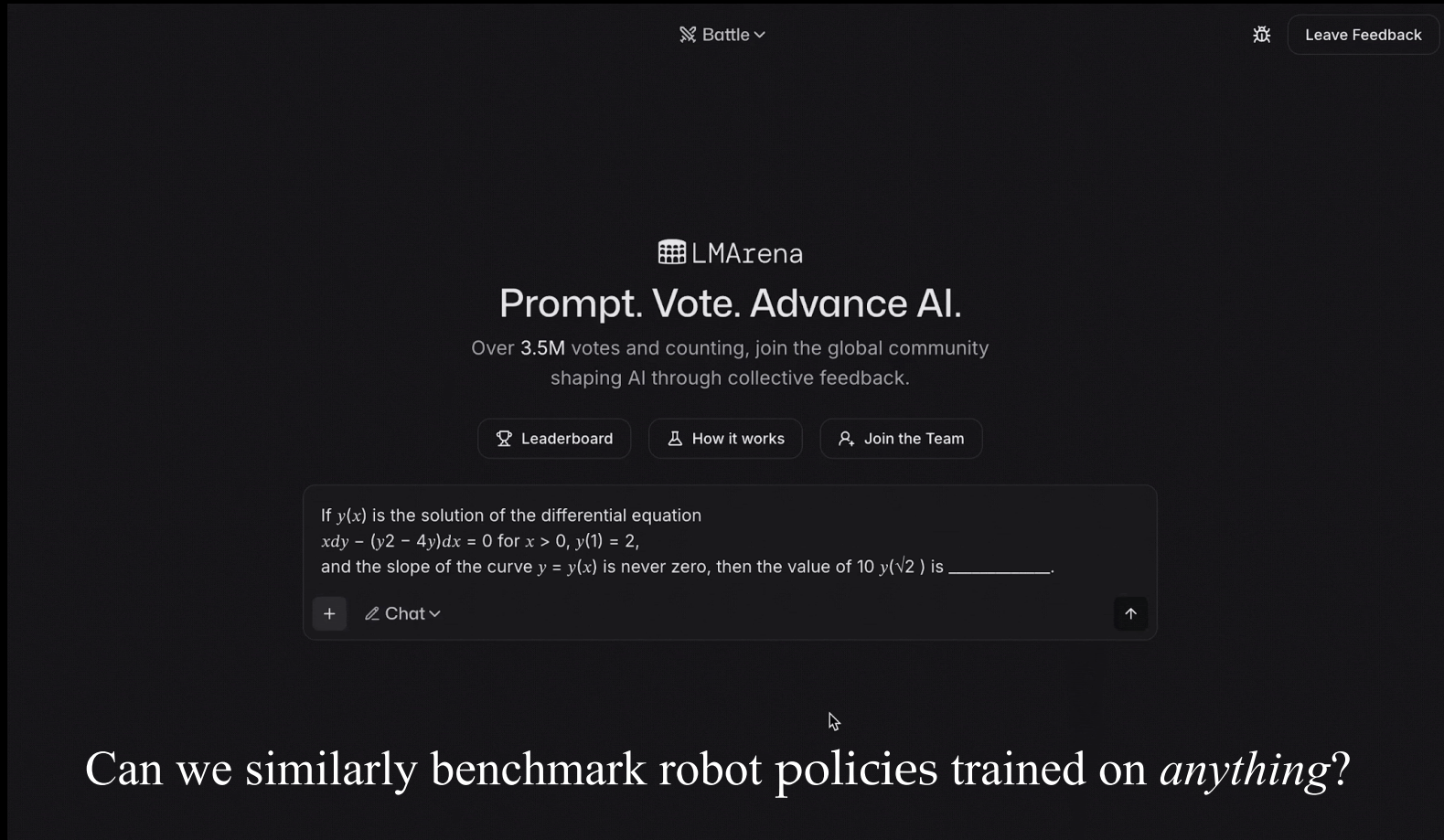


LIBERO



Calvin

LMarena: Evaluating LLMs trained on ANYTHING



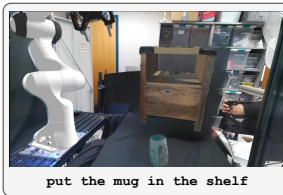
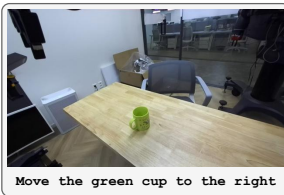
The screenshot shows the LMarena website interface. At the top, there is a "Battle" dropdown menu and a "Leave Feedback" button. The main heading is "LMarena" with the tagline "Prompt. Vote. Advance AI." Below this, it states "Over 3.5M votes and counting, join the global community shaping AI through collective feedback." There are three navigation buttons: "Leaderboard", "How it works", and "Join the Team". The central part of the image shows a chat prompt: "If $y(x)$ is the solution of the differential equation $xdy - (y^2 - 4y)dx = 0$ for $x > 0$, $y(1) = 2$, and the slope of the curve $y = y(x)$ is never zero, then the value of $10 y(\sqrt{2})$ is _____." At the bottom of the chat area, there is a "+ Chat" button and an upward arrow.

Can we similarly benchmark robot policies trained on *anything*?

RobotArena ∞ : Unlimited Robot Benchmarking via Real-to-Sim Translation

Large Scale Simulated Environments

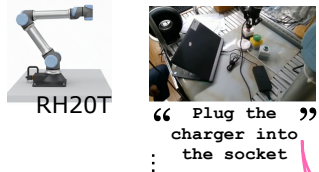
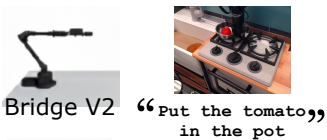
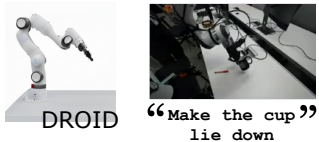
Automated Environment Creation from Real Video for VLA Evaluation



ROBOTARENA



Robotic Datasets



User-Recorded Video

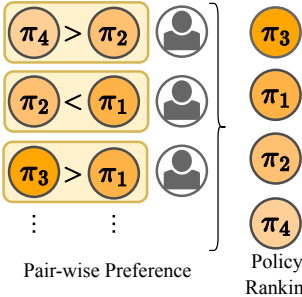


Real2Sim

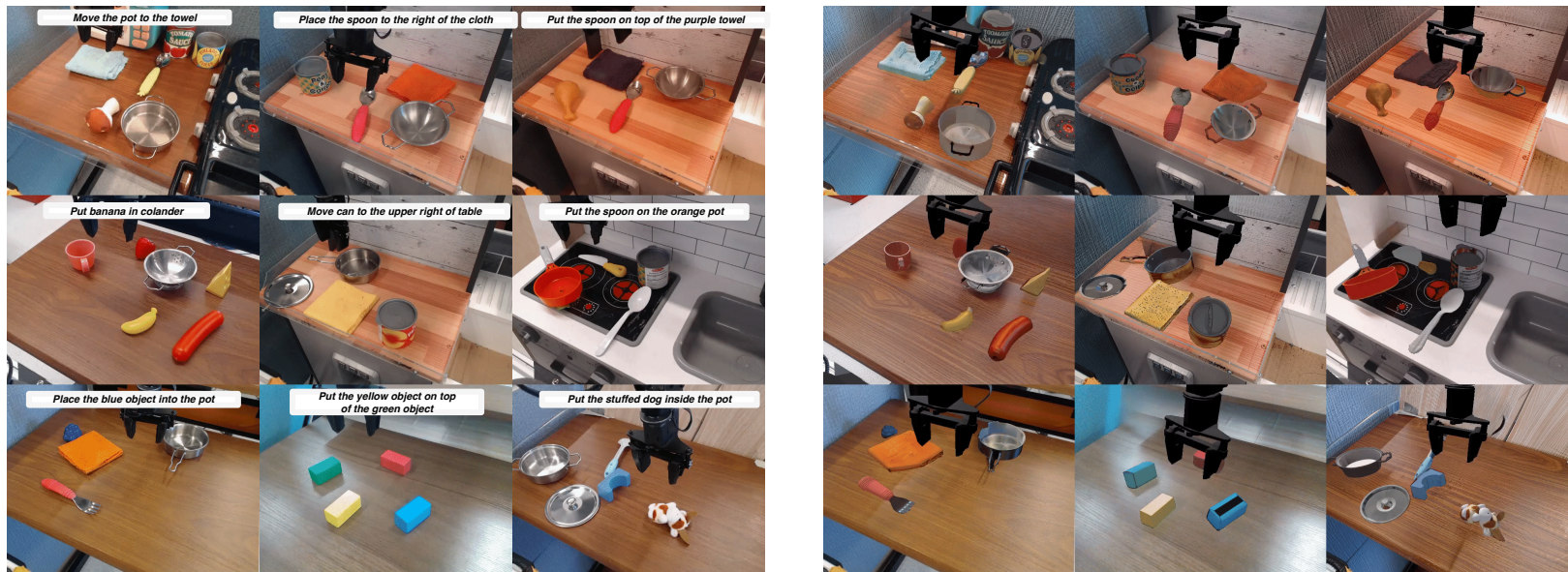
Automatic Policy Evaluation

Policy	VLM Scores		
π_1	79.32	77.03	...
π_2	76.38	77.50	...
π_3	83.77	77.18	...
⋮	⏟	⏟	⏟
	Base Env	Under Controlled	Perturbation

Human Preference Evaluation

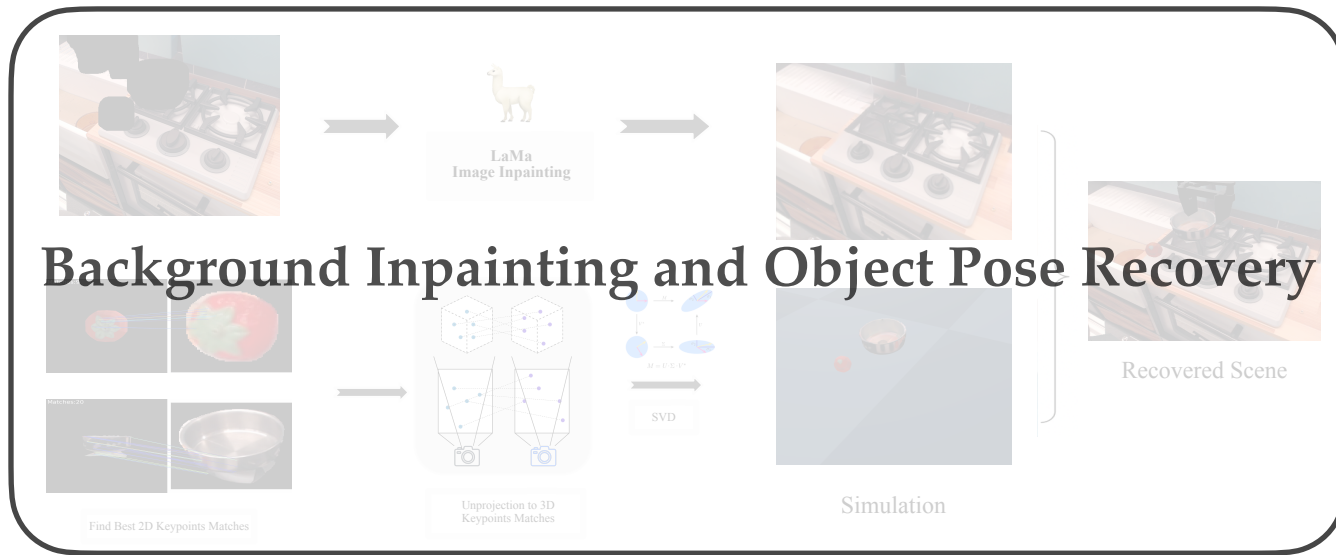
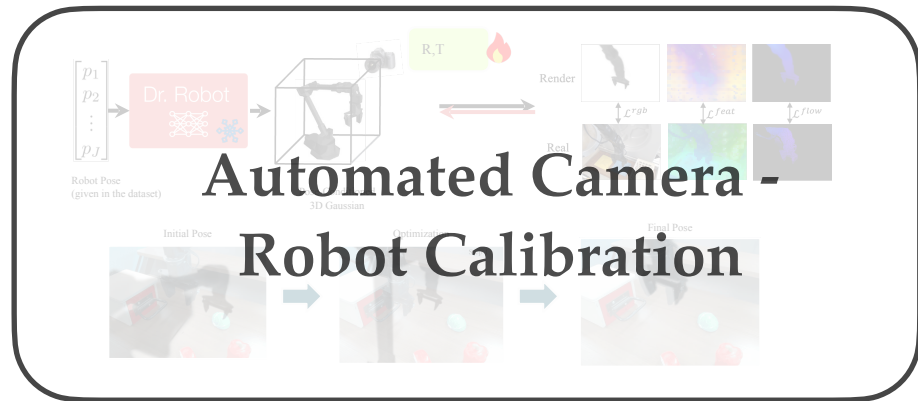


RobotArena ∞ : **Unlimited** Robot Benchmarking via Real-to-Sim Translation

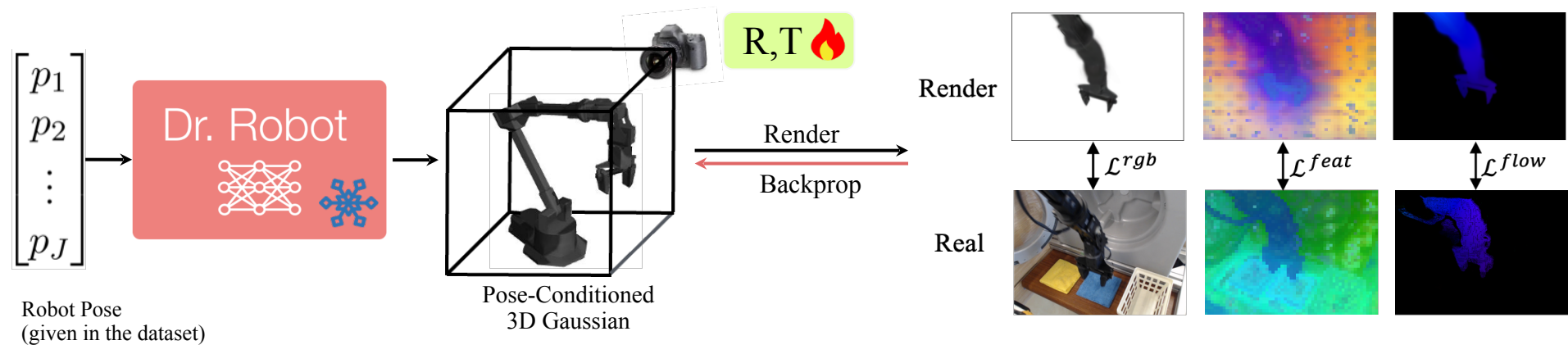


No human involvement

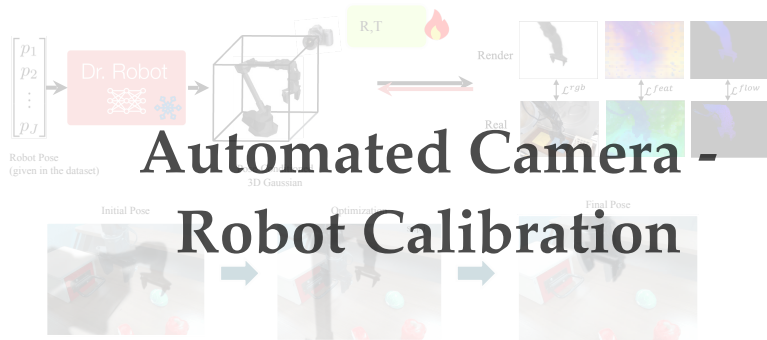
Real2Sim Translation



Automated Camera - Robot Calibration



Real2Sim Translation

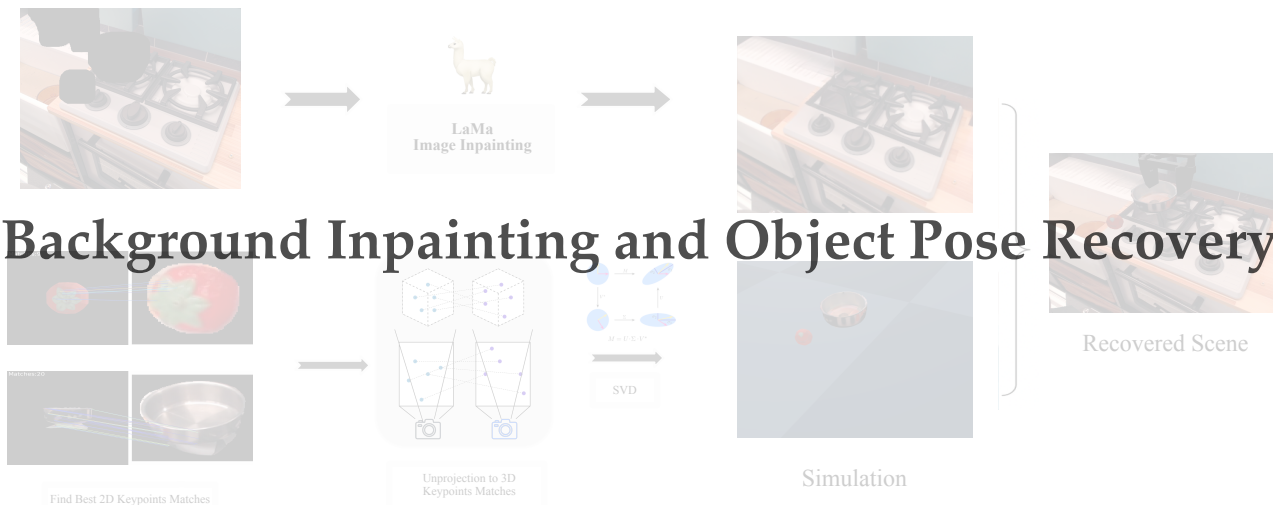


Relevant Object Reconstruction

Prompt: "...Output a JSON list of segmentation masks where each entry contains the 2D

Red Tomato Toy Metal Pot

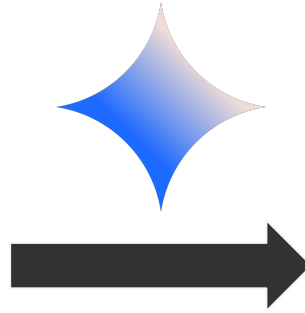
Background Inpainting and Object Pose Recovery



Segmenting and Labelling Objects



Prompt: "...Output a JSON list of segmentation masks where each entry contains the 2D bounding box in "box_2d", a descriptive text label in "label", and the mask in "mask" "".

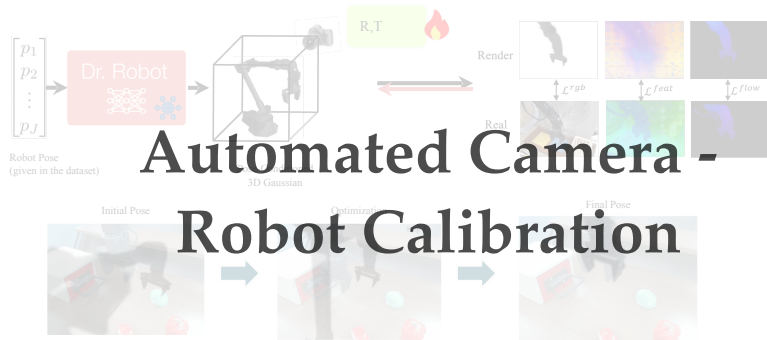


Red Tomato Toy



Metal Pot

Real2Sim Translation

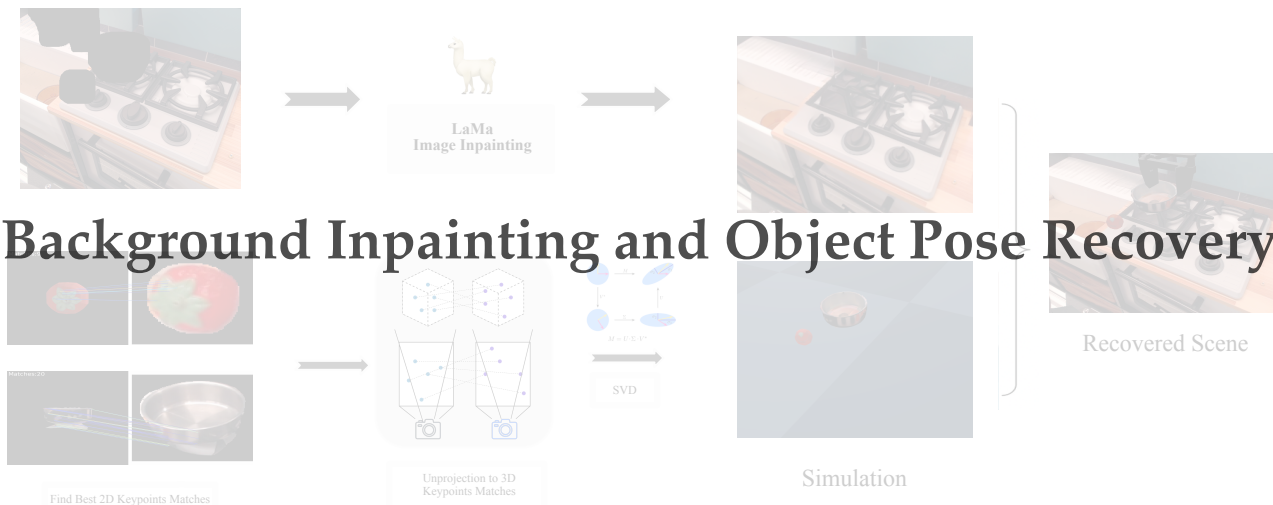


Relevant Object Reconstruction

Prompt: "...Output a JSON list of segmentation masks where each entry contains the 2D

Red Tomato Toy Metal Pot

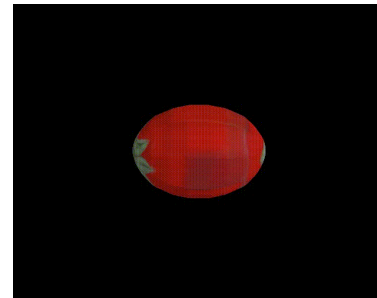
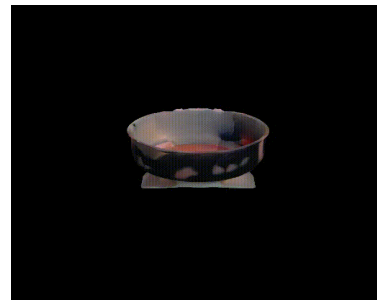
Background Inpainting and Object Pose Recovery



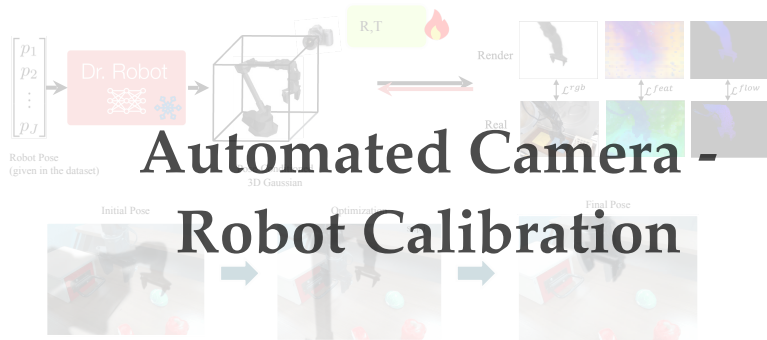
Reconstructing 3D Objects



**Image to 3D Mesh
Generation**



Real2Sim Translation

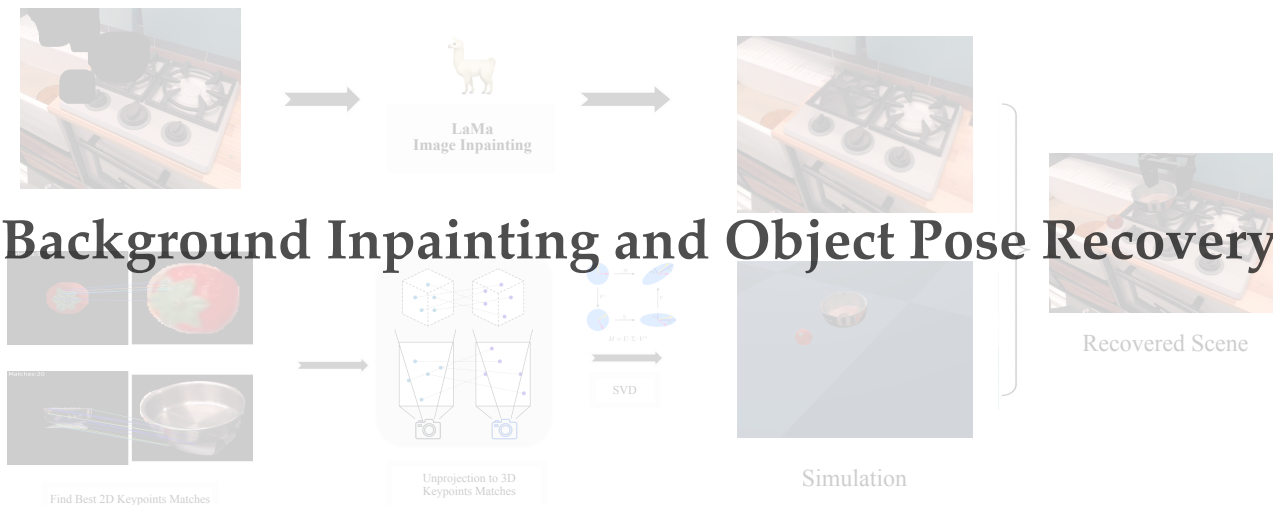


Relevant Object Reconstruction

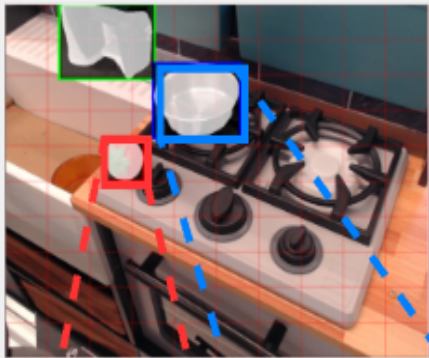
Prompt: "...Output a JSON list of segmentation masks where each entry contains the 2D

Red Tomato Toy Metal Pot

Background Inpainting and Object Pose Recovery



Querying VLMs for Object Physics Properties



Red Tomato Toy

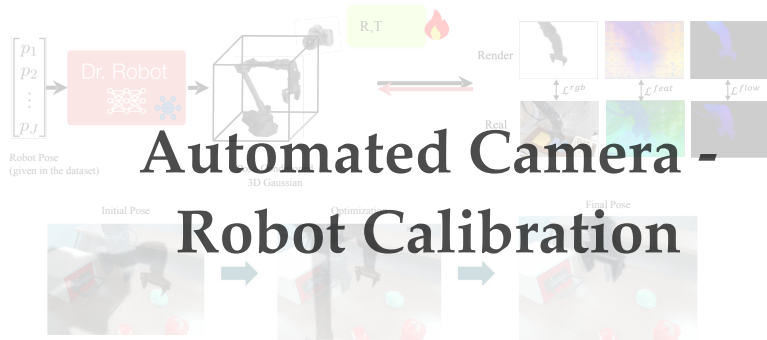


Metal Pot

Prompt: I have a {class_label}. Estimate and return a JSON object describing its possible physical properties. For each entry, include one plausible material the object is most likely made of,, along with their estimated physics attributes. Do not hallucinate incorrect material properties or non-existent material types. Do not include duplicate materials. Please return it as a JSON object that can be parsed automatically in the following format. The format is:

```
json_format = f'{{ \
  '{class_label}': [ \
    'Mass (kg)': ..., \
    'Friction Coefficient': ..., \
    'Material': one of ['Glass', 'Water', 'Emission', 'Plastic', 'Rough', \
    'Smooth', 'Reflective', 'Metal', 'Iron', 'Aluminium', 'Copper', 'Gold'] \
  ] \
}
```

Real2Sim Translation

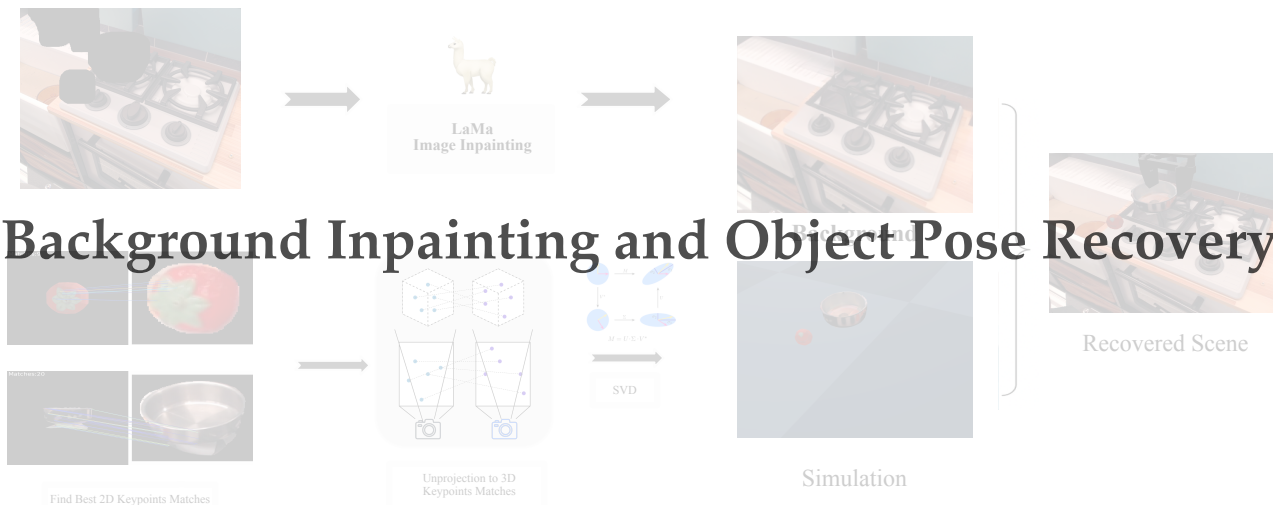


Relevant Object Reconstruction

Prompt: "...Output a JSON list of segmentation masks where each entry contains the 2D

Red Tomato Toy Metal Pot

Background Inpainting and Object Pose Recovery



Background Inpainting

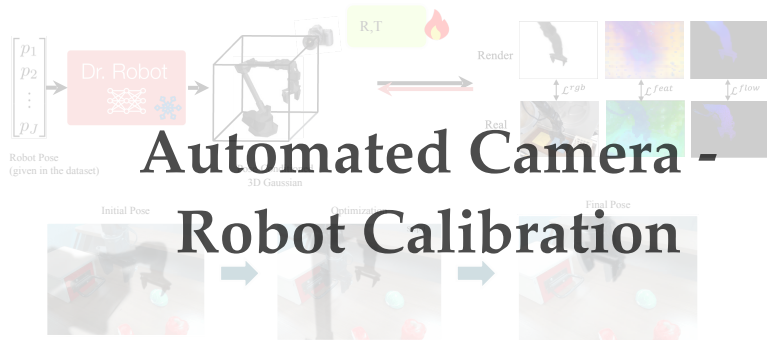


**LaMa
Image Inpainting**



Background

Real2Sim Translation

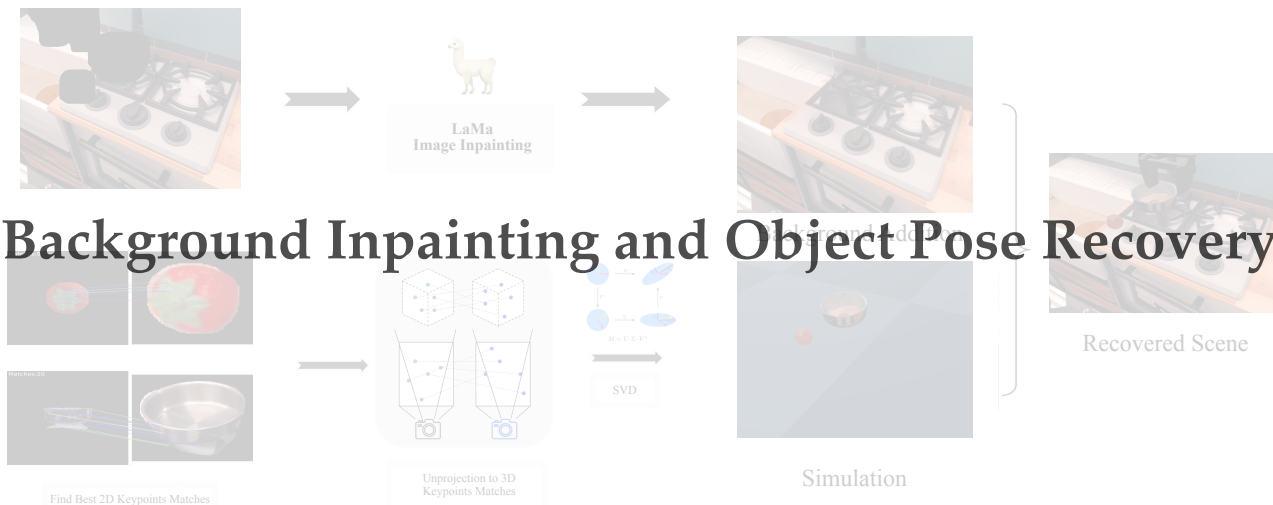


Relevant Object Reconstruction

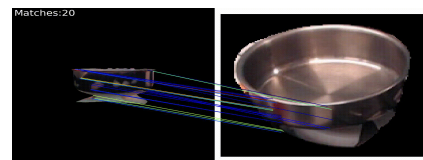
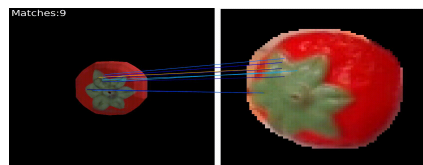
Prompt: "...Output a JSON list of segmentation masks where each entry contains the 2D

Red Tomato Toy Metal Pot

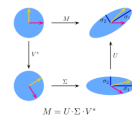
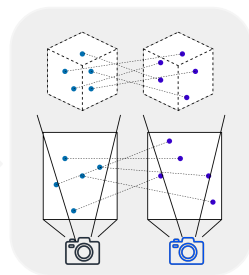
Background Inpainting and Object Pose Recovery



Estimating Object Scale, Orientation and Translation with Render-and-Compare



Find Best 2D Keypoints Matches



SVD



Background Addition

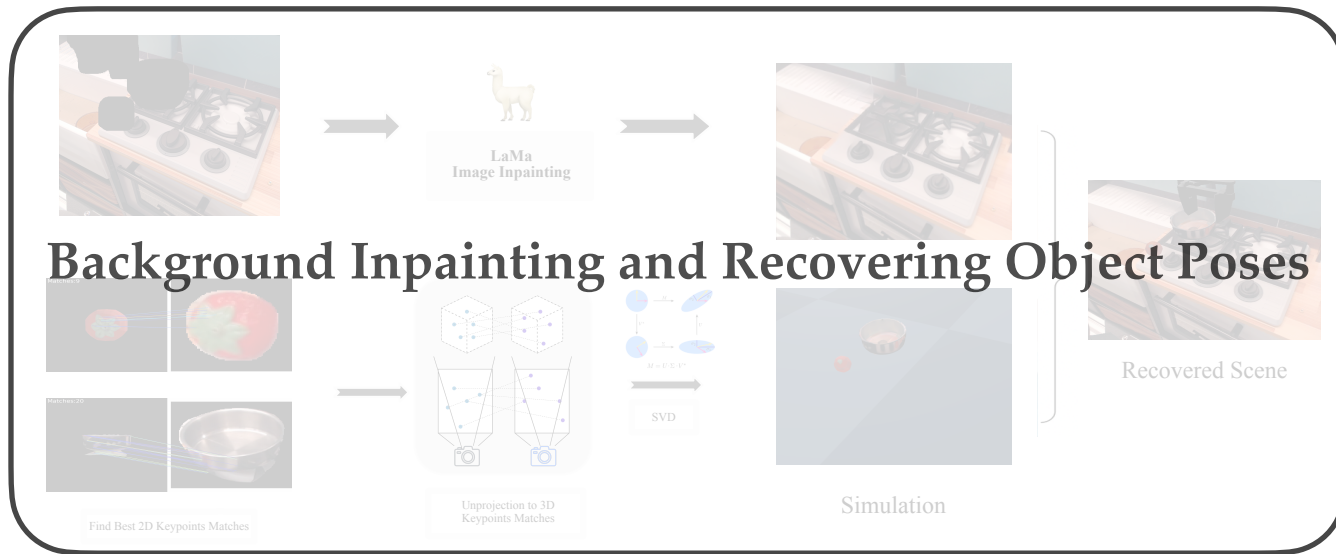
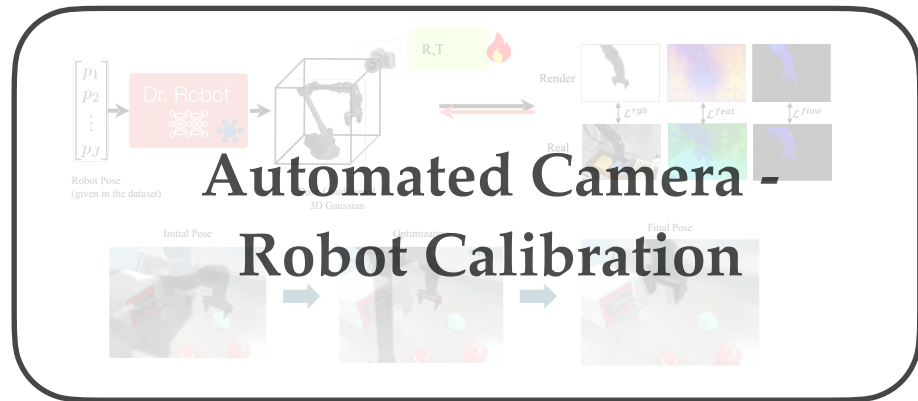


Simulation



Recovered Scene

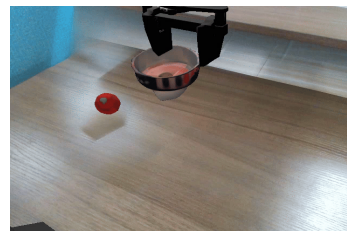
Real2Sim Translation



Scene Perturbations



Base Scene



Background Change



Color Shift



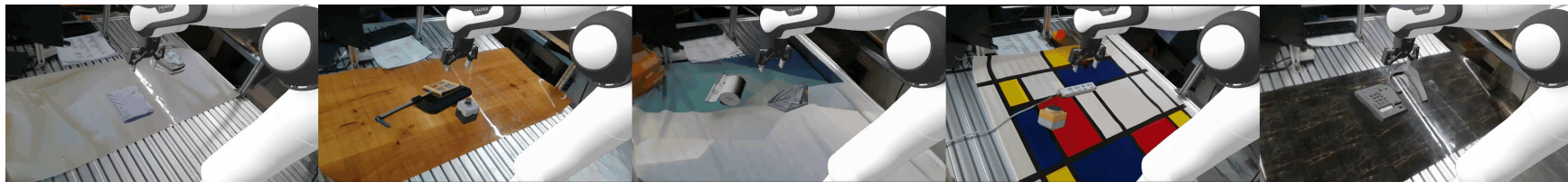
Object Pose Change



BridgeSim

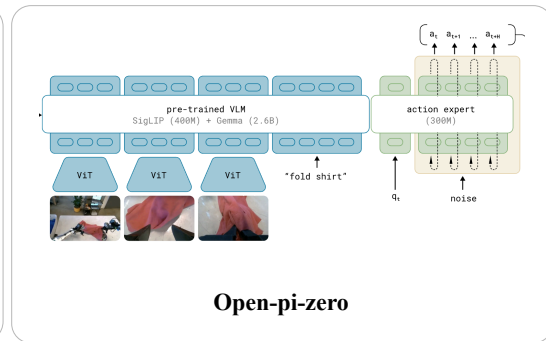
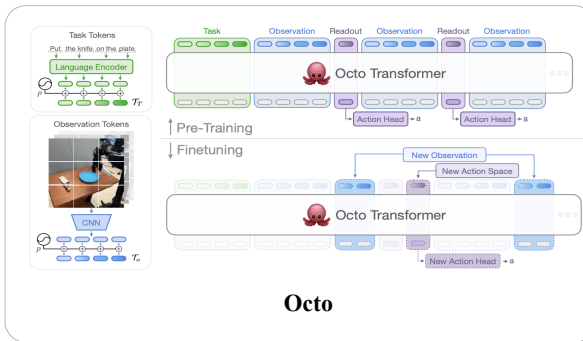
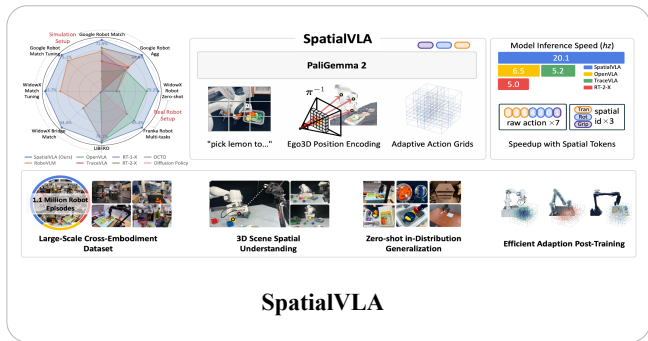
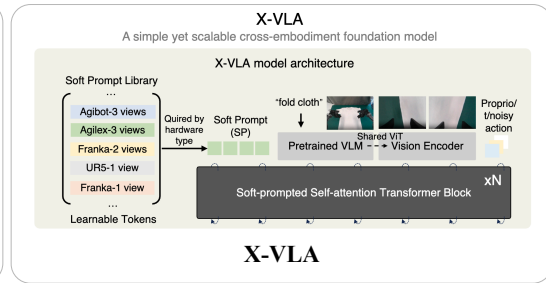
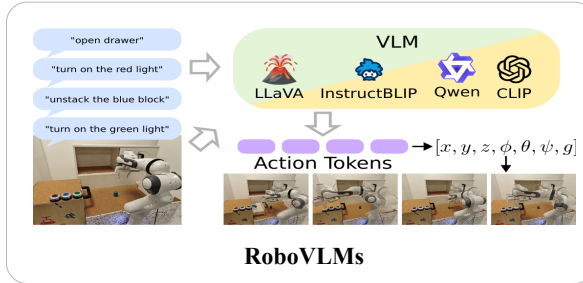
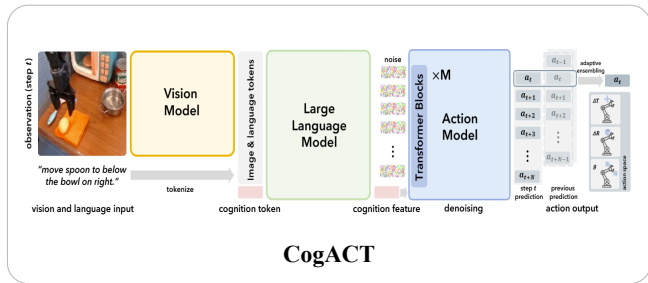


DROIDSim



Rh20TSim

RobotArena ∞ : Deploying VLAs



Li et al. (2024), CogACT: A Foundational Vision-Language-Action Model for Synergizing Cognition and Action in Robotic Manipulation. arXiv:2411.19650
 Li et al. (2024), Towards Generalist Robot Policies: What Matters in Building Vision-Language-Action Models. arXiv:2412.14058
 Qu et al. (2024), SpatialVLA: Exploring Spatial Representations for Visual-Language-Action Models. RSS 2025
 Octo Model Team and Ghosh et al. (2023), Octo: An Open-Source Generalist Robot Policy. RSS 2024
 The X-VLA Team and Chen et al. (2025), X-VLA: Soft-Prompted Transformer as Scalable Cross-Embodiment Vision-Language-Action Model. arXiv:2510.10247
 Physical Intelligence and Black et al. (2024), S_{pi}-OS: A Vision-Language-Action Flow Model for General Robot Control. RSS 2025
 Ren et al. (2024), Open-Pi-Zero: An Open-Source Re-implementation of Pi-0. GitHub Repository:allenzren/open-pi-zero

RobotArena ∞ : Deploying VLAs

“Put the red object in the pot”



CogACT



RoboVLM



X-VLA



SpatialVLA



Octo



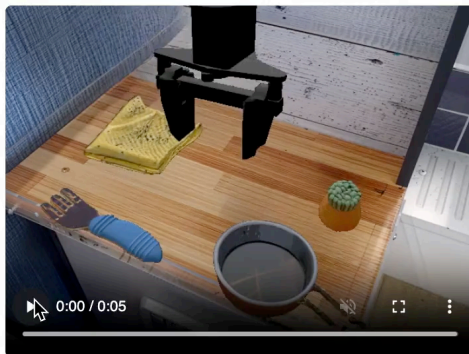
Open-pi-zero

Which policy is better?

Task: Place the fork near the yellow cloth

Video 7 of 150

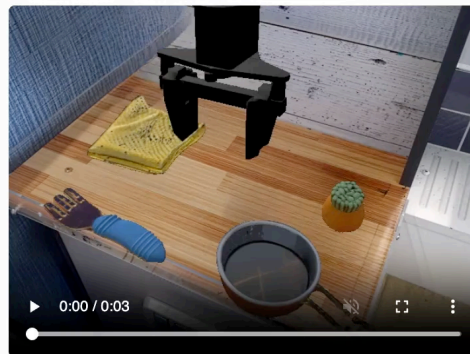
Policy A



Describe the robot's attempt to complete the task:

e.g., The robot reaches to the cup, picks it up and...

Policy B



Describe the robot's attempt to complete the task:

e.g., The robot made no meaningful move...

Submit Descriptions

Please provide descriptions for both videos to continue.

Your evaluation:

👉 Left is Better

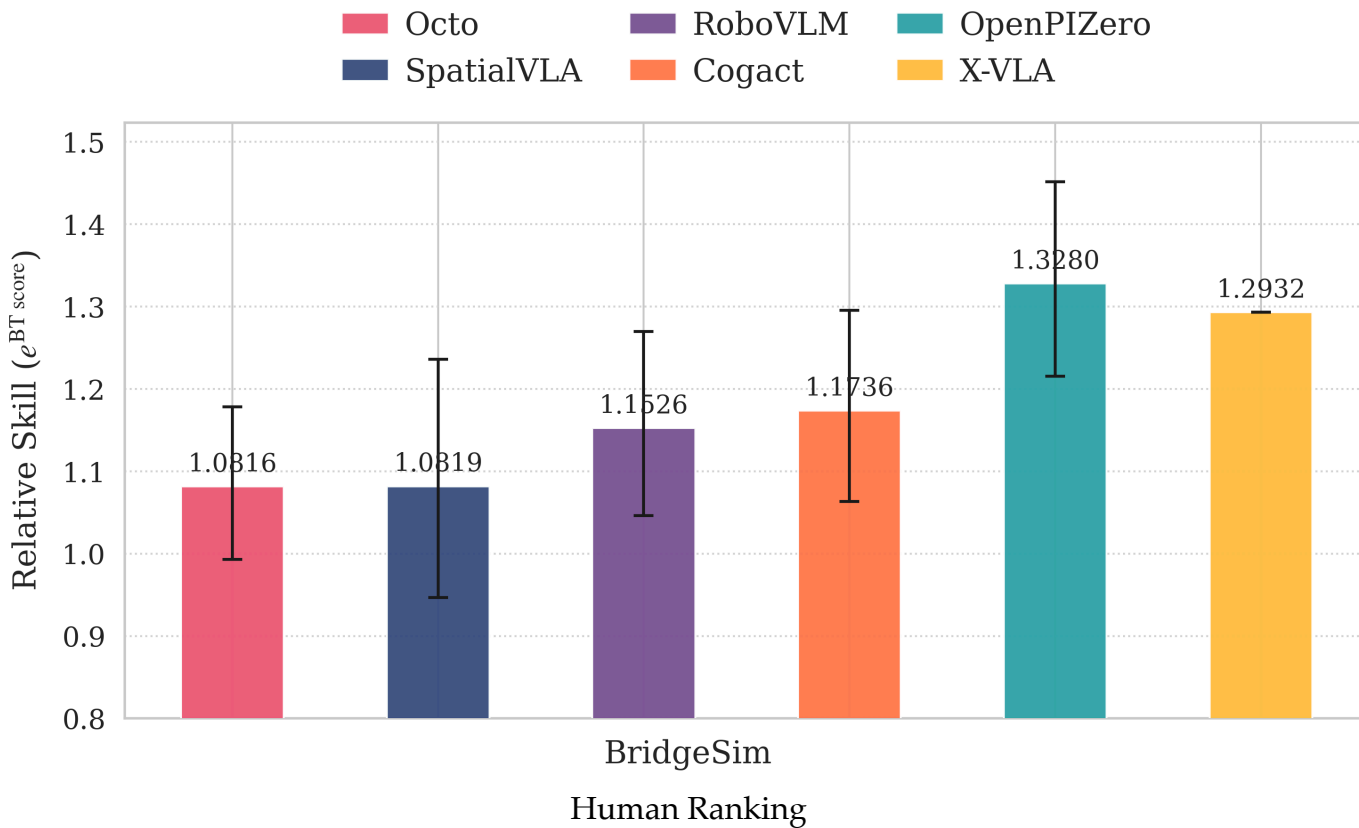
👉 Tie

Right is Better 👉

(equally good or equally bad)

Next

RobotArena ∞ : Policy Rankings at Scale



Task Progress Estimation with VLMs

“ Put the tomato
in the pot ”

Task Prompt



Frames



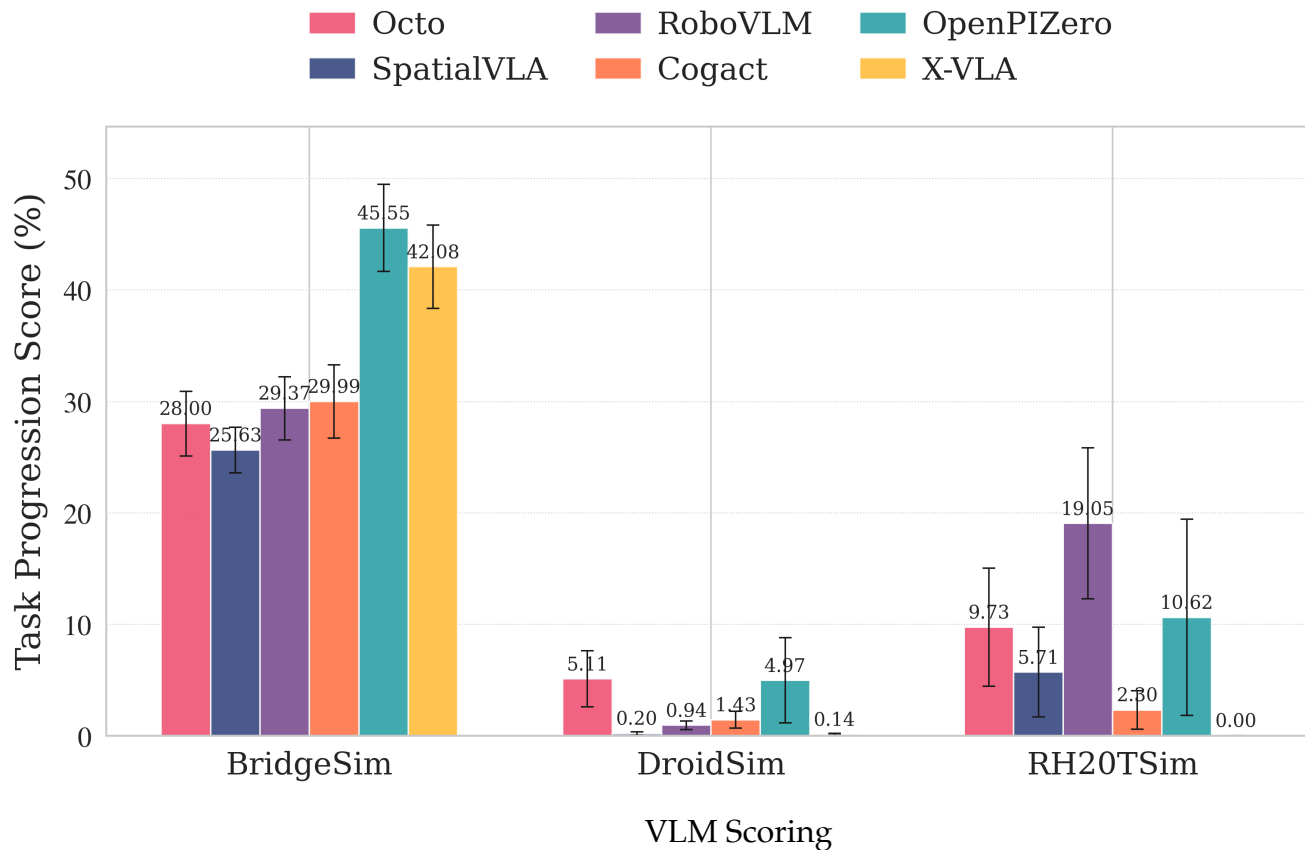
Robot and
Object State





Vision Language Model



RobotArena ∞ : Policy Rankings at Scale



Limitations and Ongoing Work

-  **Static Camera Assumption**
The image background superimposition assume a static camera
-  **Going beyond rigid objects:** toward fluids, gases, and deformable materials

Thank you!

