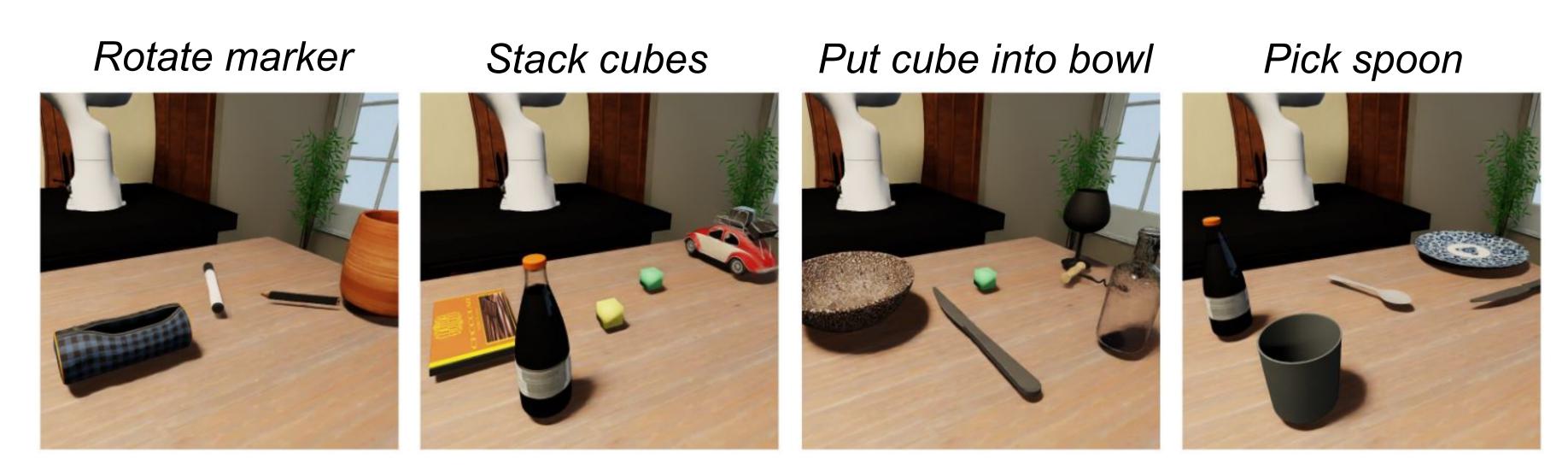




UNIVERSITY OF AMSTERDAM

Goal

- > Reliable and scalable evaluation of generalist policies for robot manipulation.
- Understanding failure modes to guide model development, data collection, and evaluations.



Motivation

- > Real world performance is the only true benchmark for robotics.
- > Many simulated environments are disconnected from reality and fail to serve as a good proxy.
- > Answering "how does my policy perform on task X if Y is changed" without spending hours on a real robot.

Challenges

Realistic robot control and protocol for large-scale that reflects real-world VLA performance.

Contributions

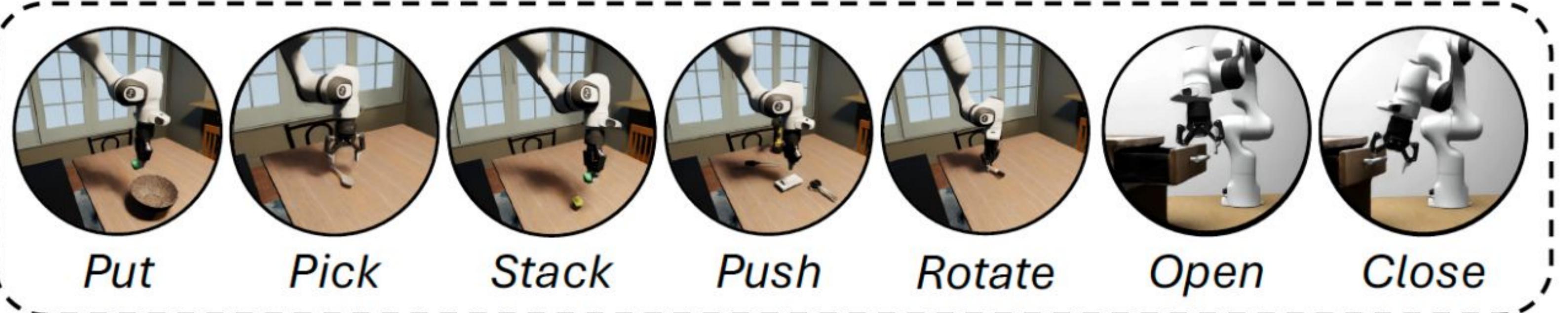
- > A real-to-sim aligned simulation for DROID.
- > REALM: A generalization benchmark supporting 15 perturbations from 3 categories.
- > Results for 3 state-of-the-art VLA models.

REALM: A Real-to-Sim Aligned Benchmark for Generalization in Robotic Manipulation

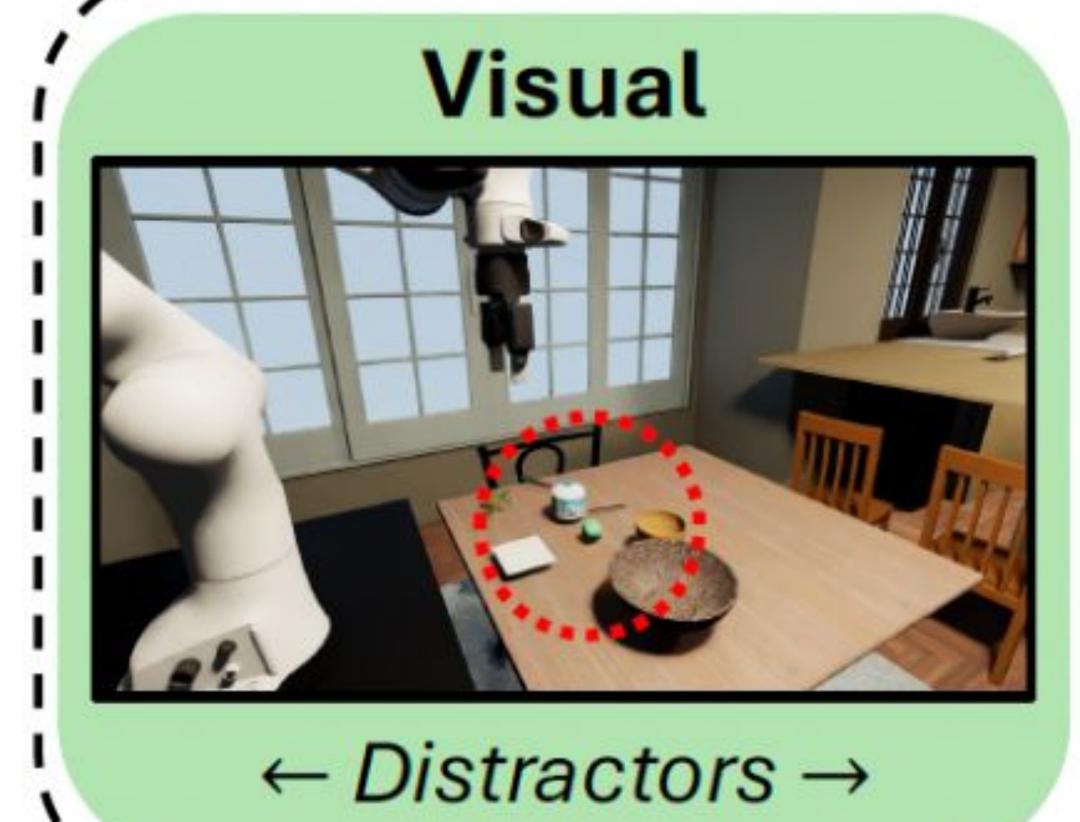
Martin Sedlacek¹, Pavlo Yefanov¹, Georgy Ponimatkin¹, Simon Pilc¹, Mederic Fourmy¹, Evangelos Kazakos¹, Cees Snoek², Josef Sivic¹, Vladimir Petrik¹ ¹CIIRC, Czech Technical University ²VIS Lab, University of Amsterdam



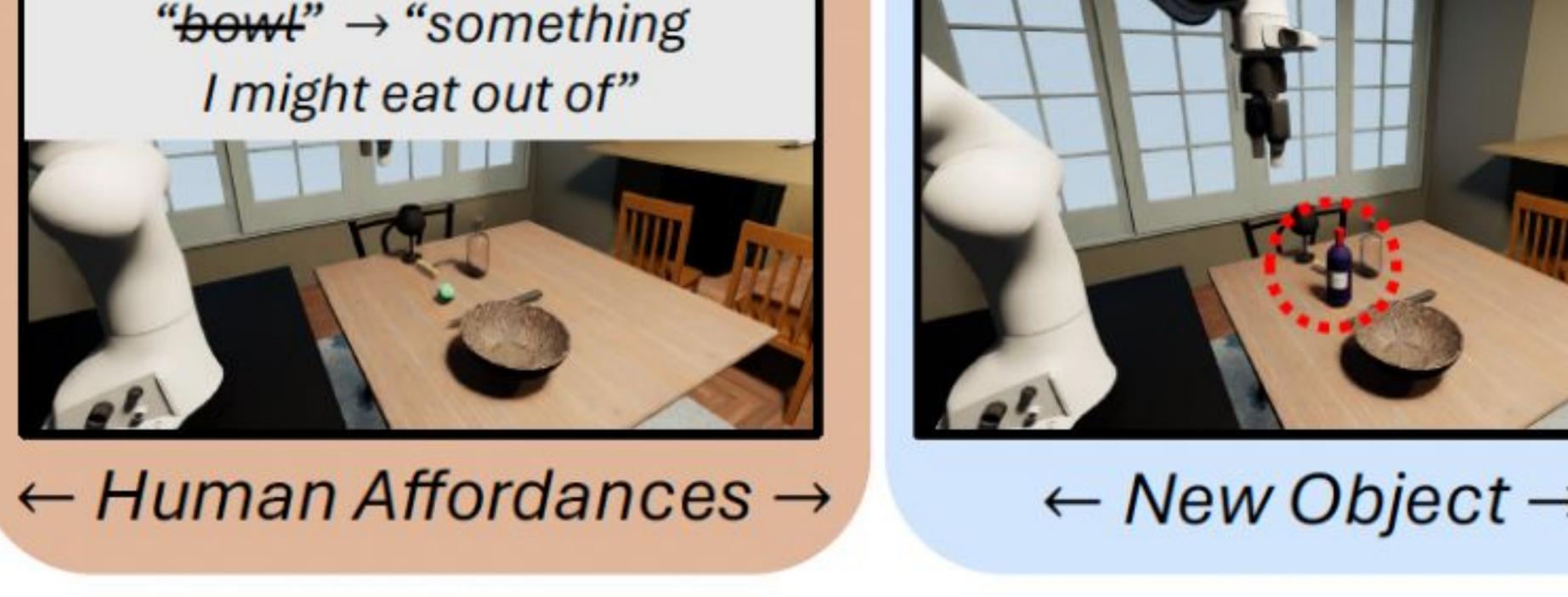
https://martin-sedlacek.com/realm/



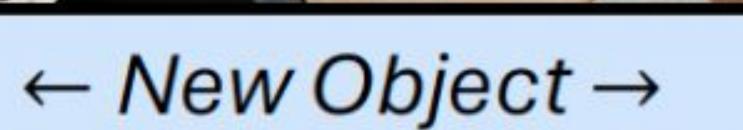
7 manipulation skills



Semantic "bowl" → "something I might eat out of"

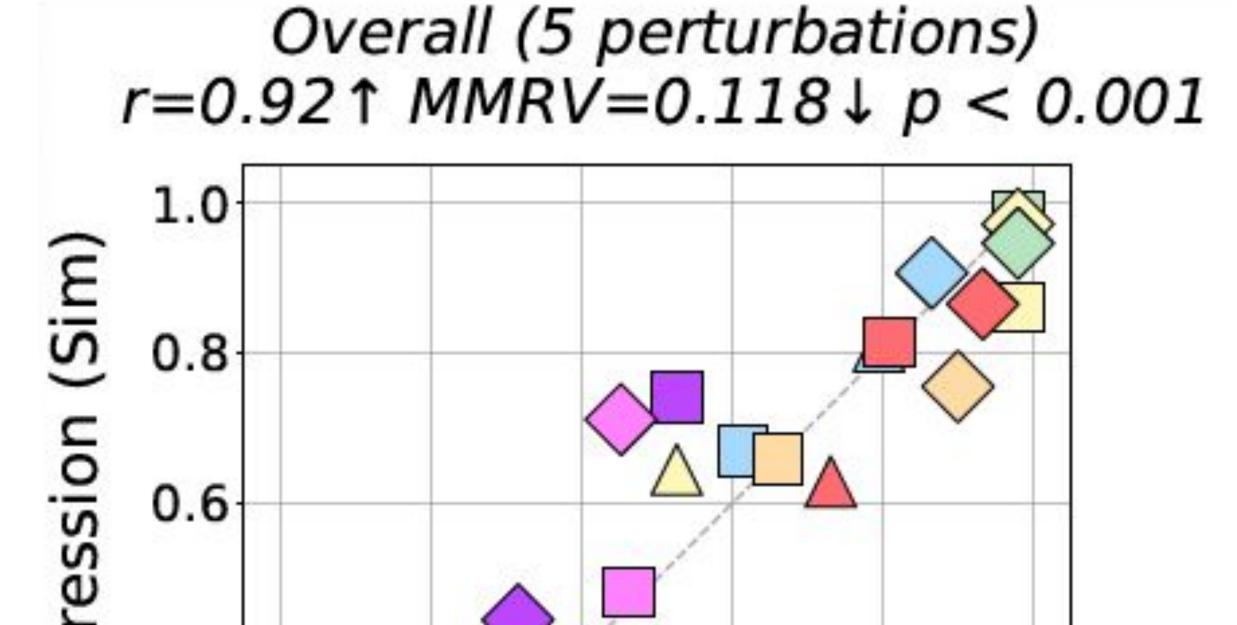


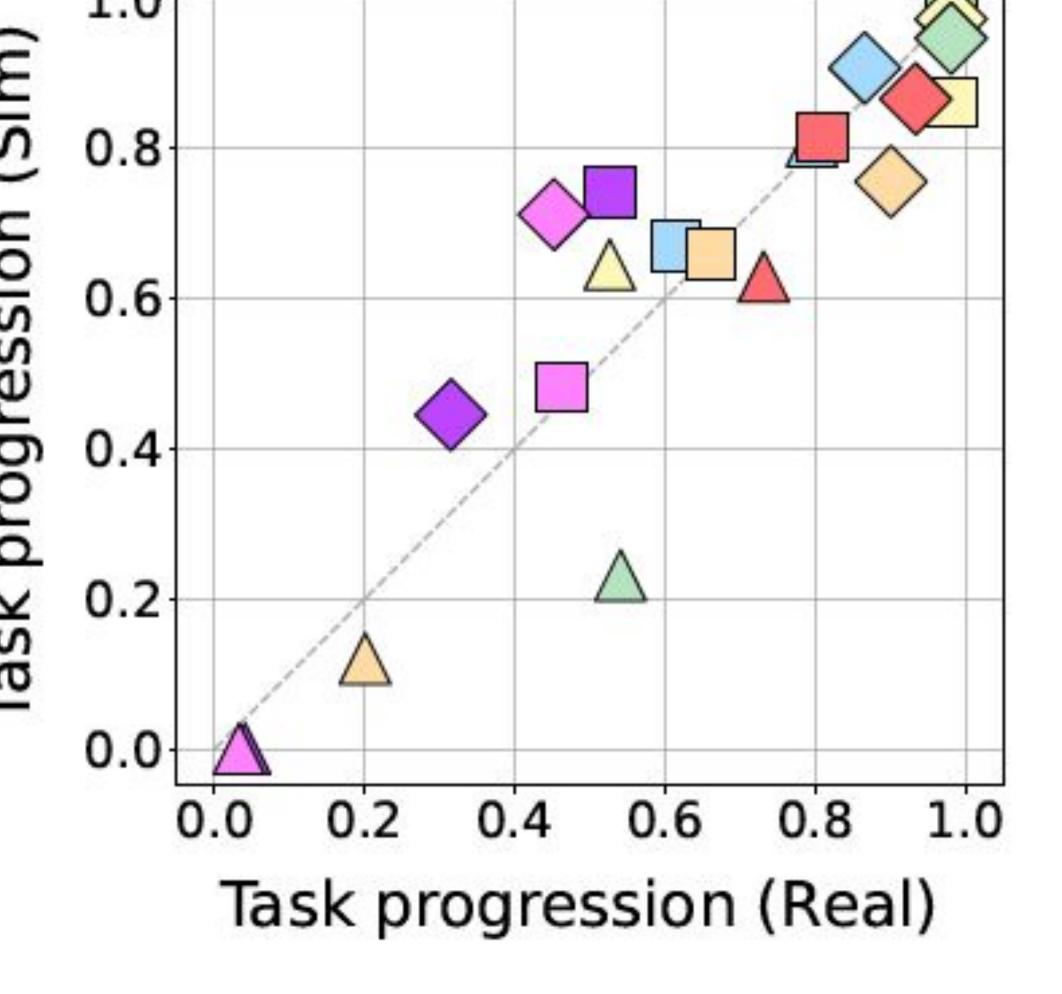
Behavioral



real-to-sim aligned

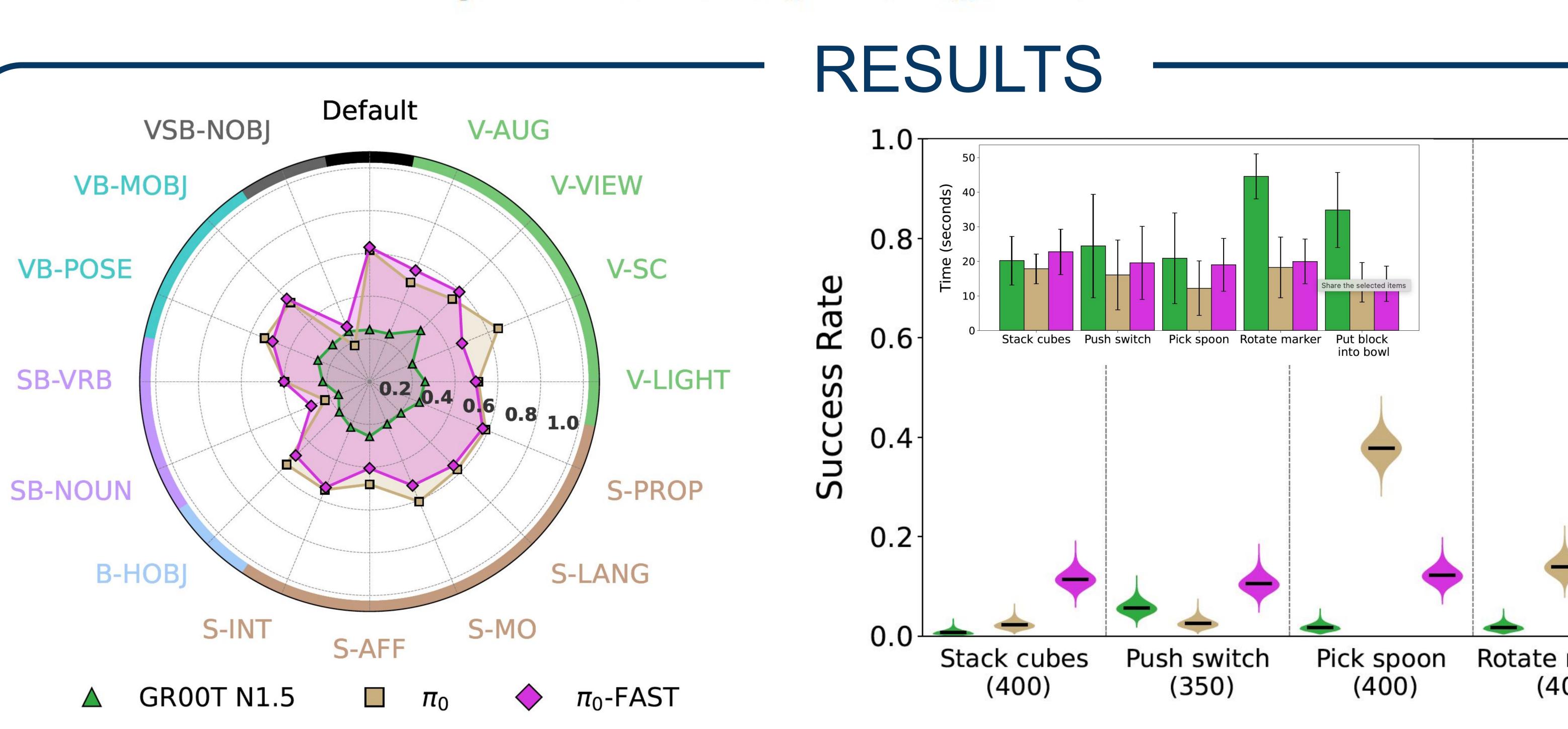
into bowl





- > High correlation between sim and real under diverse perturbations / tasks.
- > Realistic robot control and high-fidelity simulated visuals.

15 perturbations | 3 categories



TAKEAWAYS

- > High-fidelity simulation with aligned robot control provides a valuable proxy for real-world performance, helping mitigate the issue of saturated and disconnected simulation benchmarks.
- > Despite VLM backbones pre-trained on Internet-scale data, there is a noticeable drop in performance under most perturbations.
- > But, performance for most models does not go to zero, indicating they do possess some generalization capabilities.
- > Behavioral generalization across skills and objects is still the most challenging, even for cases that are over-represented in the data.