

Generative Embodied AI

Guest Lecture

Ruoshi Liu Columbia University

10-16-2024

Generative Embodied AI

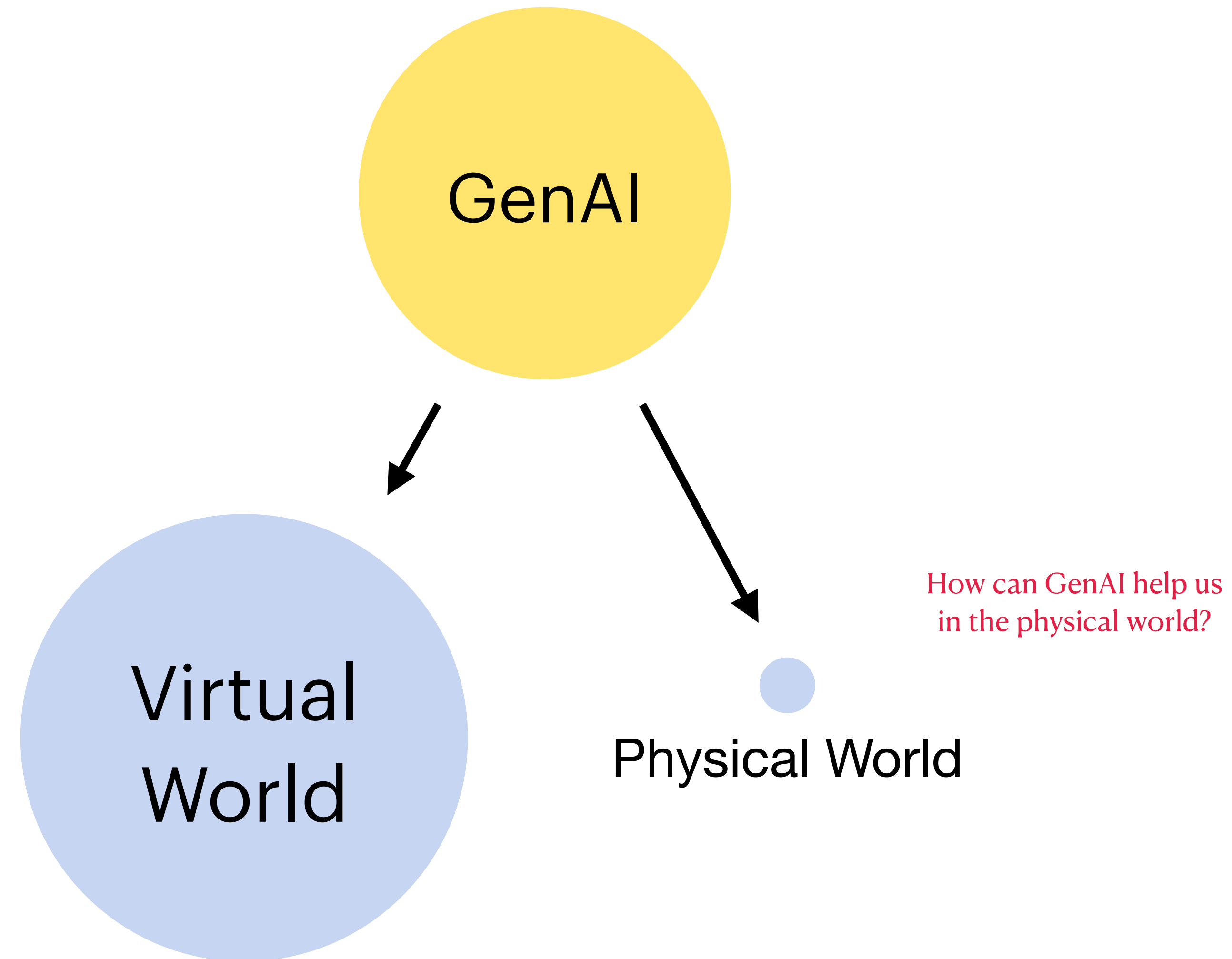
(With the generous help of neural rendering)

Guest Lecture

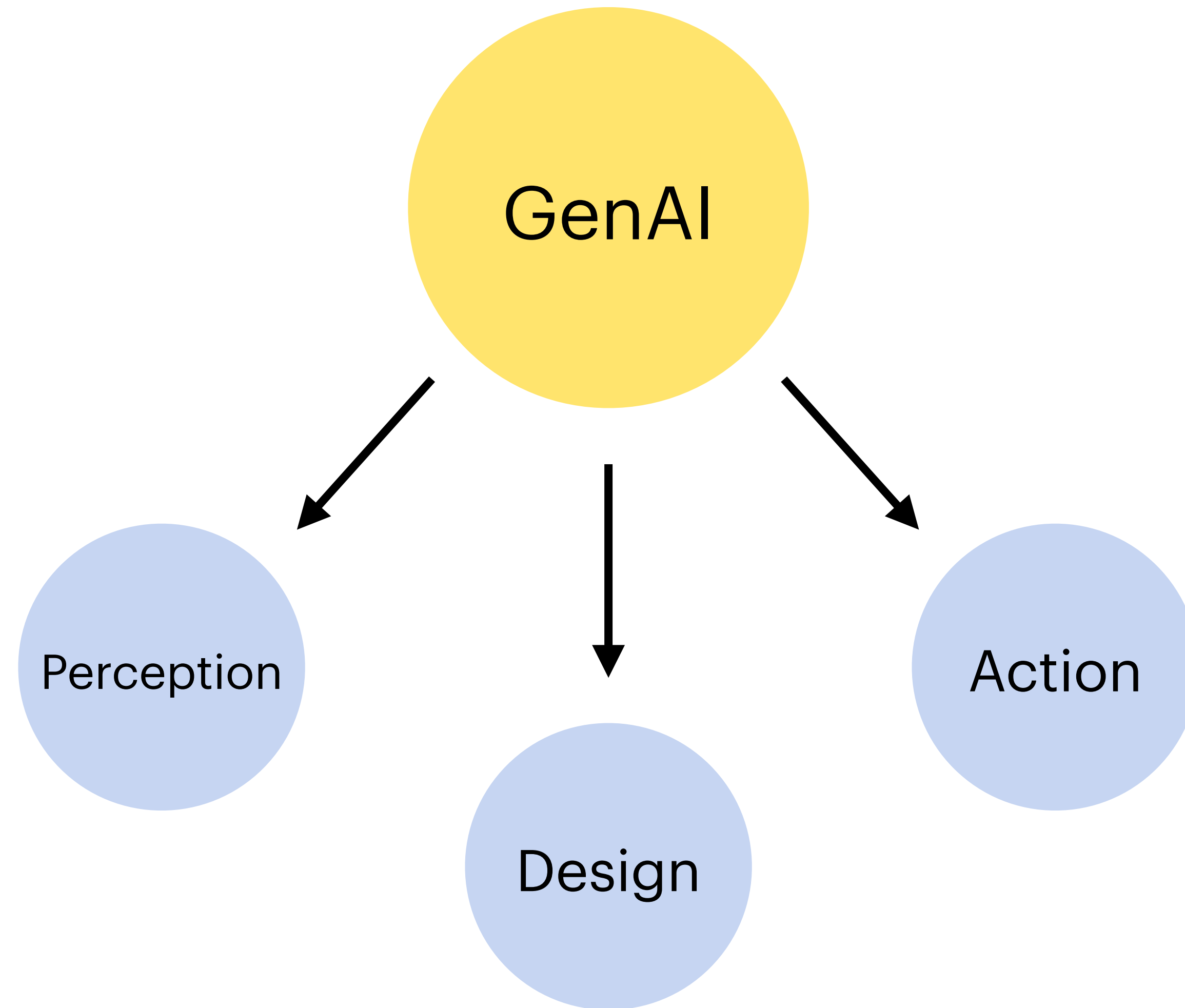
Ruoshi Liu Columbia University

10-16-2024

Generative Models Today



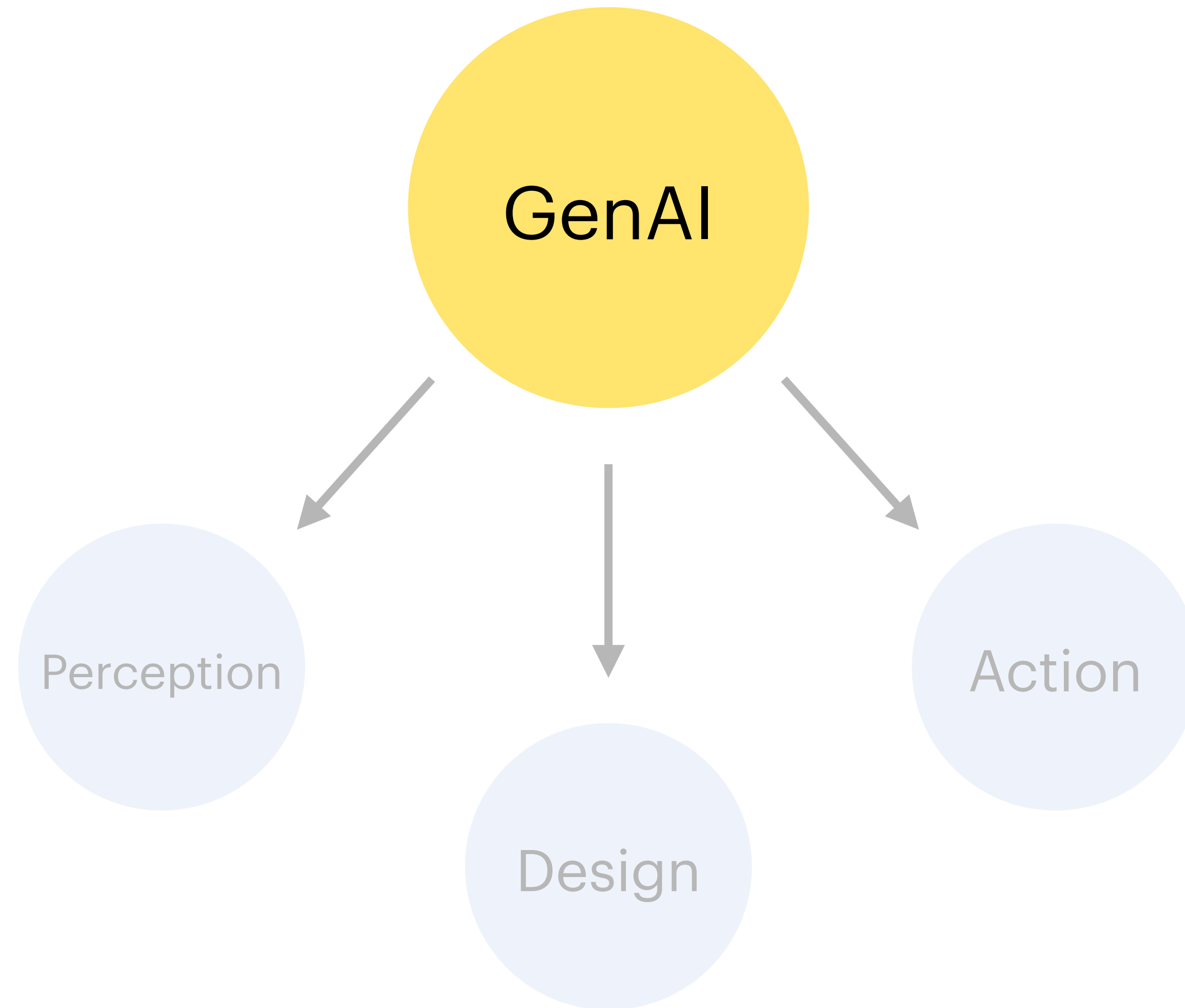
Generative Embodied AI



How do we build better 3D generative models?

How can 3d Gen help us in the physical world?

Generative Embodied AI



How do we build better 3D generative models?

Vision Generative Model Timeline

PCA
(Tu et al.)

2006

⋮



Vision Generative Model Timeline



Image Parsing: Unifying Segmentation, Detection, and Recognition

ZHUOWEN TU AND XIANGRONG CHEN

Departments of Statistics, University of California, Los Angeles, Los Angeles, CA 90095, USA

ztu@stat.ucla.edu

xrchen@stat.ucla.edu

ALAN L. YUILLE

Departments of 'Statistics' and 'Psychology', University of California, Los Angeles, Los Angeles, CA 90095, USA

yuille@stat.ucla.edu

SONG-CHUN ZHU

Departments of 'Statistics' and 'Computer Science' University of California, Los Angeles, Los Angeles, CA 90095, USA

sczhu@stat.ucla.edu

Figure 6. Random samples drawn from the PCA face model.

Vision Generative Model Timeline



Vision Generative Model Timeline

PCA (Tu et al.)
2006



NeurIPS 2010 DL Workshop: <https://twitter.com/ethanCaballero/status/1544400983261954048?s=20&t=6ocUrsX7sOCYkONn6Zas3w>

Vision Generative Model Timeline



Vision Generative Model Timeline

Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com



Vision Generative Model Timeline

[Submitted on 16 Jan 2014 (v1), last revised 30 May 2014 (this version, v3)]

Stochastic Backpropagation and Approximate Inference in Deep Generative Models

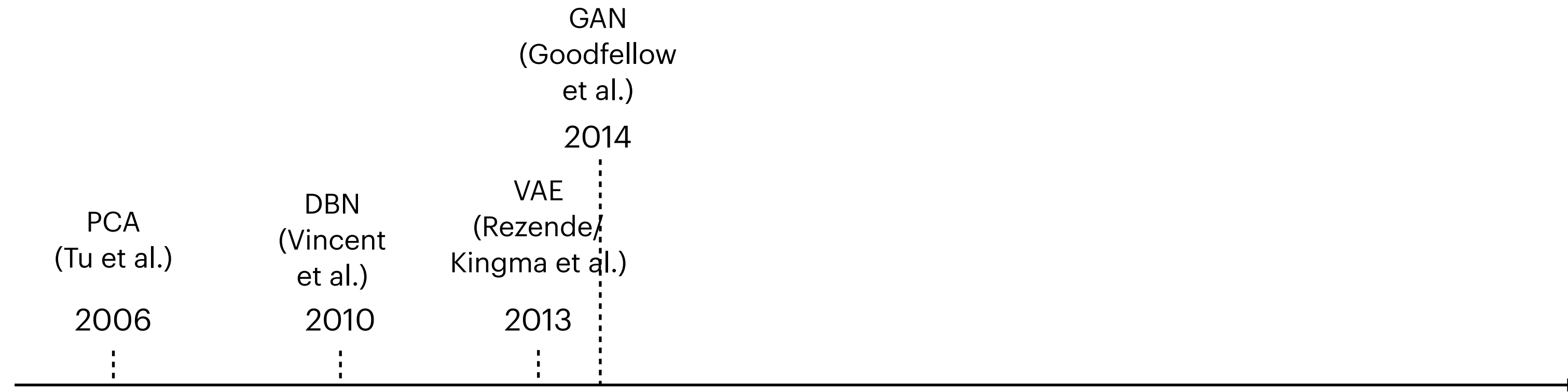
[Danilo Jimenez Rezende](#), [Shakir Mohamed](#), [Daan Wierstra](#)

[Submitted on 20 Dec 2013 (v1), last revised 10 Dec 2022 (this version, v11)]

Auto-Encoding Variational Bayes

[Diederik P Kingma](#), [Max Welling](#)

Vision Generative Model Timeline



Generative Adversarial Nets

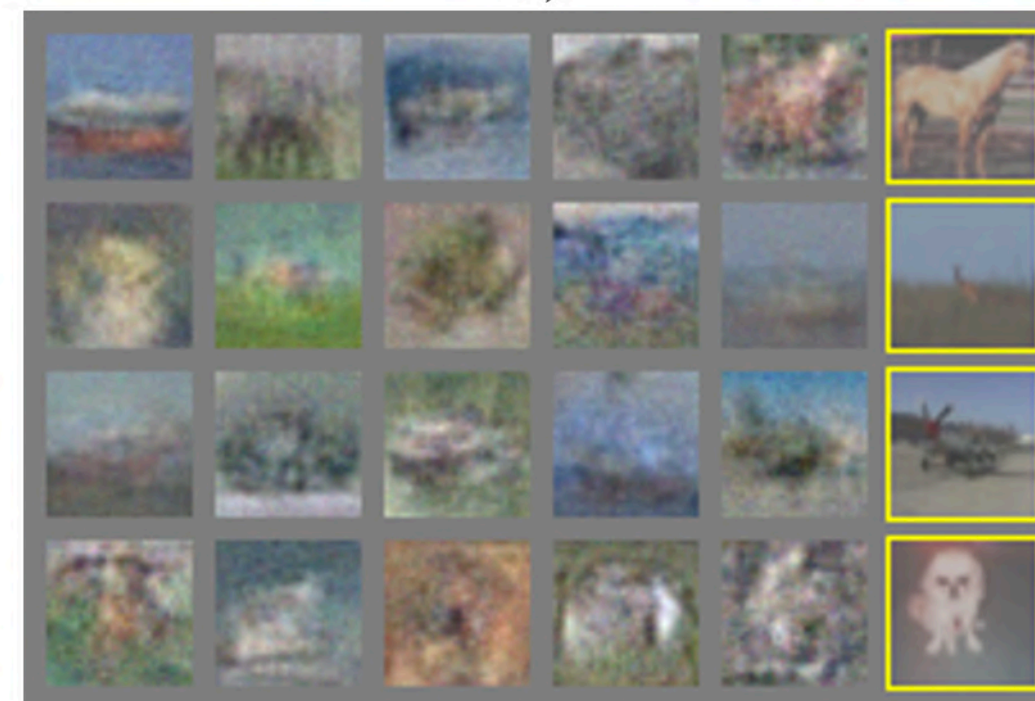
Ian J. Goodfellow, Jean Pouget-Abadie*, Mehdi Mirza, Bing Xu, David Warde-Farley,
Sherjil Ozair,† Aaron Courville, Yoshua Bengio‡
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7



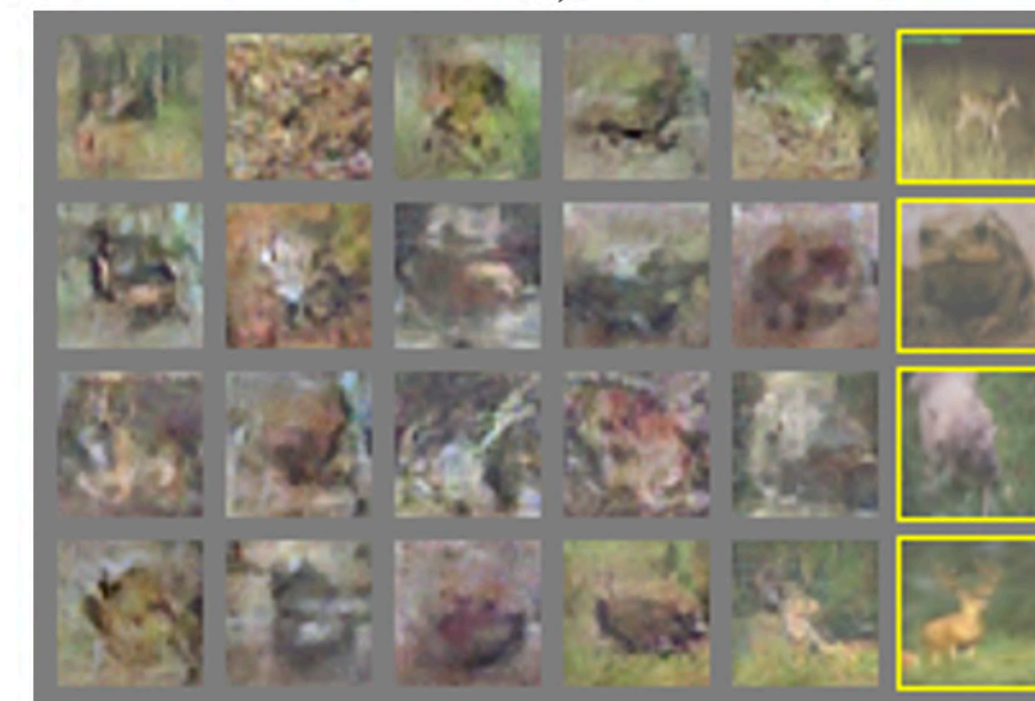
a)



b)

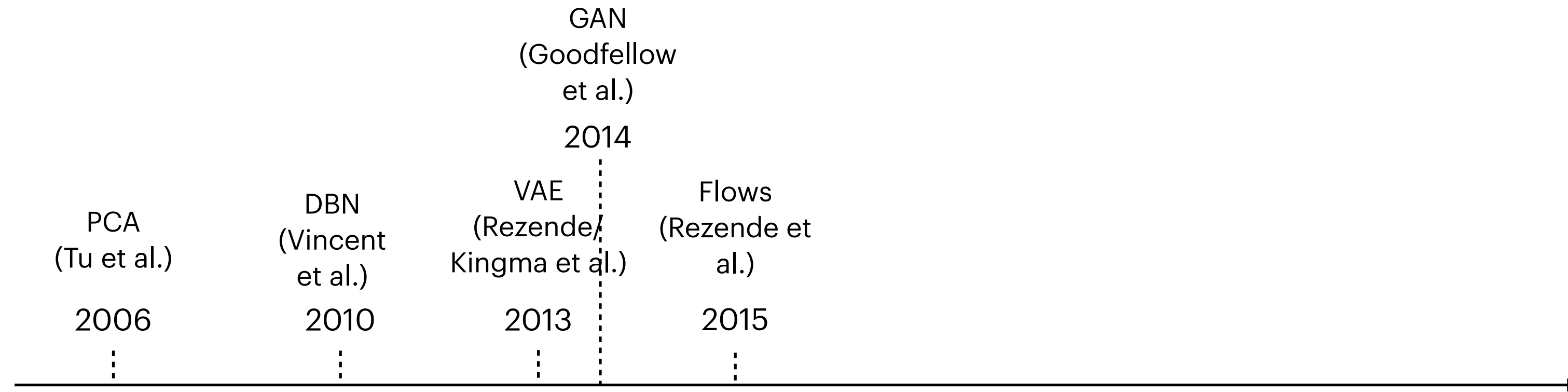


c)

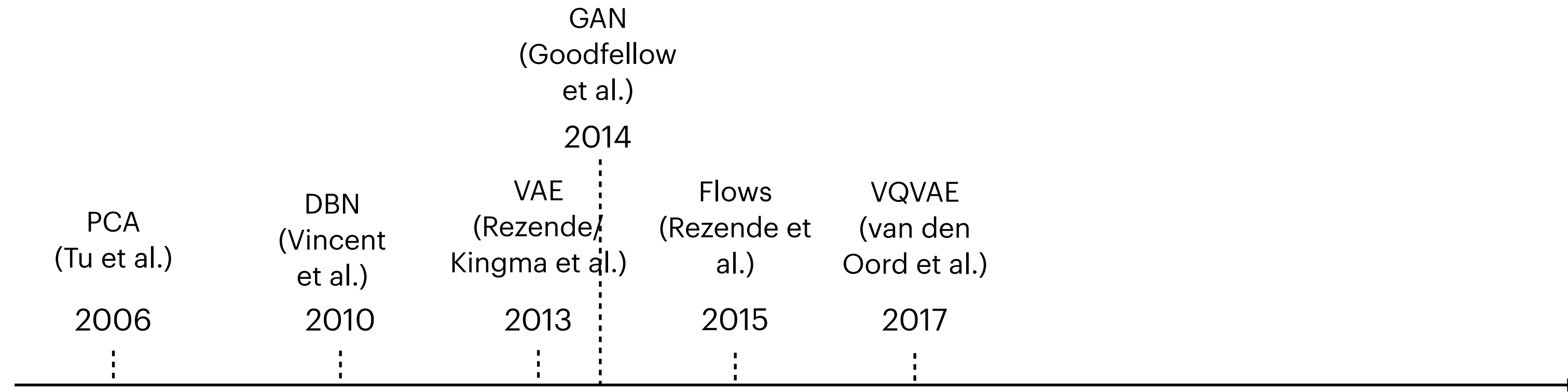


d)

Vision Generative Model Timeline



Vision Generative Model Timeline



Vision Generative Model Timeline

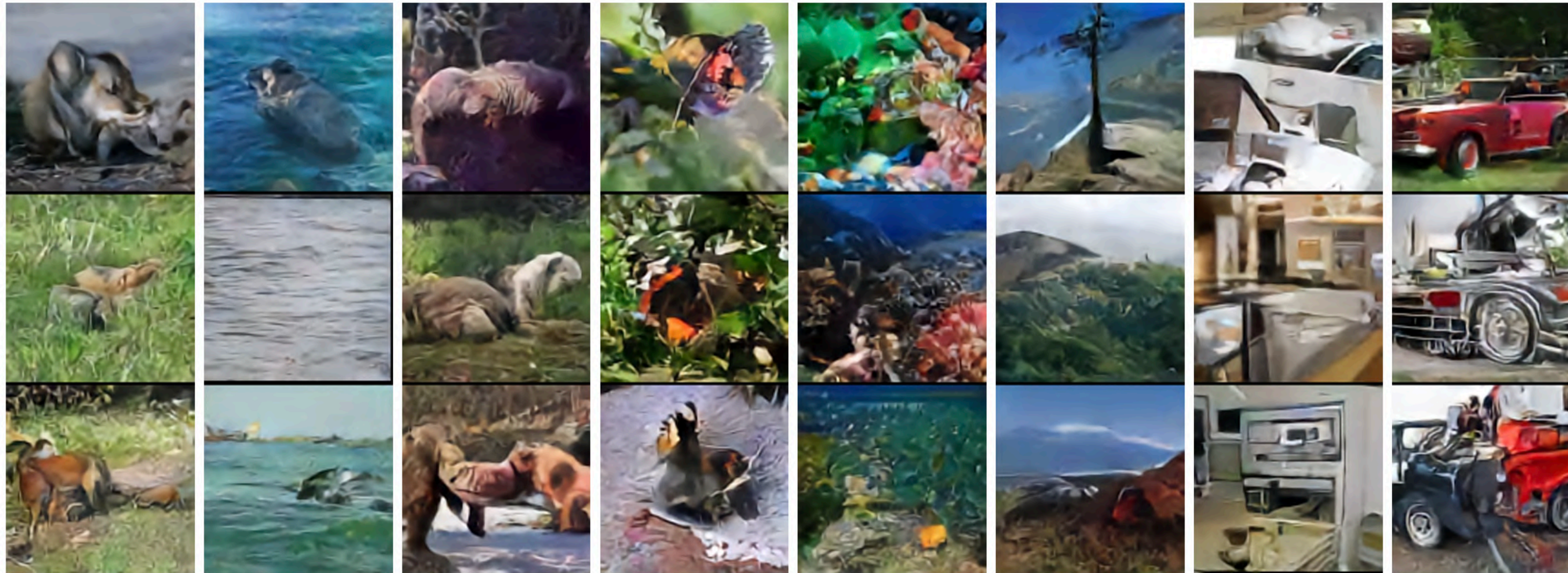
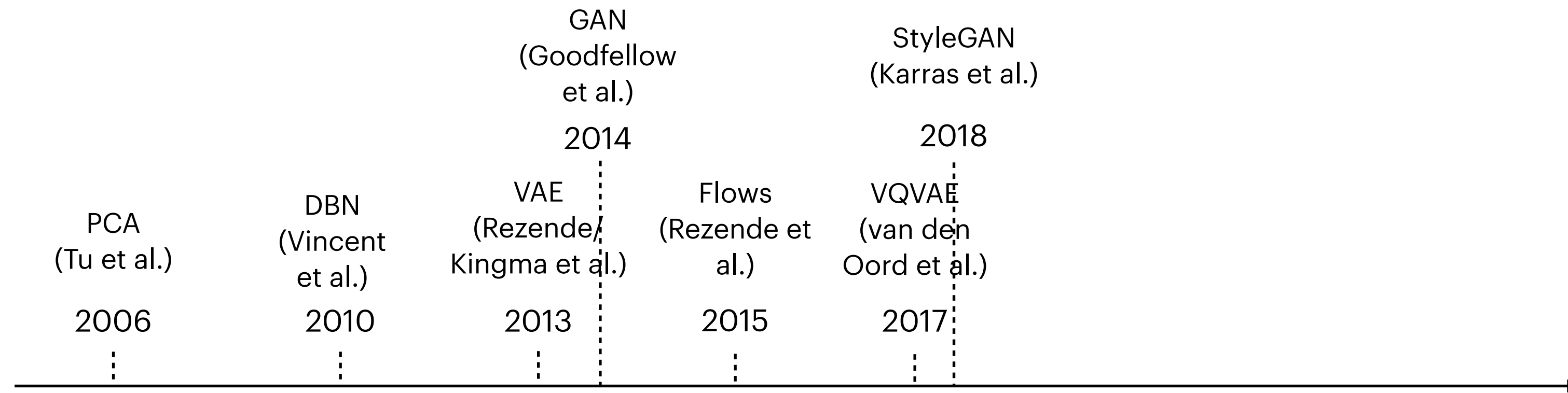


Figure 3: Samples (128x128) from a VQ-VAE with a PixelCNN prior trained on ImageNet images. From left to right: kit fox, gray whale, brown bear, admiral (butterfly), coral reef, alp, microwave, pickup.

Vision Generative Model Timeline



Vision Generative Model Timeline

A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras
NVIDIA

tkarras@nvidia.com

Samuli Laine
NVIDIA

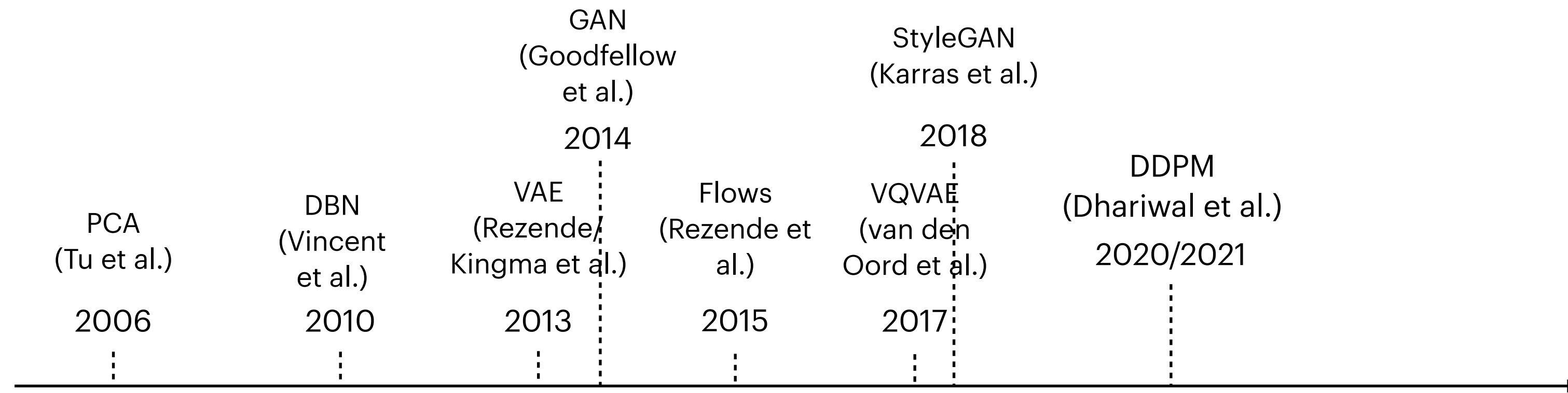
slaine@nvidia.com

Timo Aila
NVIDIA

taila@nvidia.com



Vision Generative Model Timeline

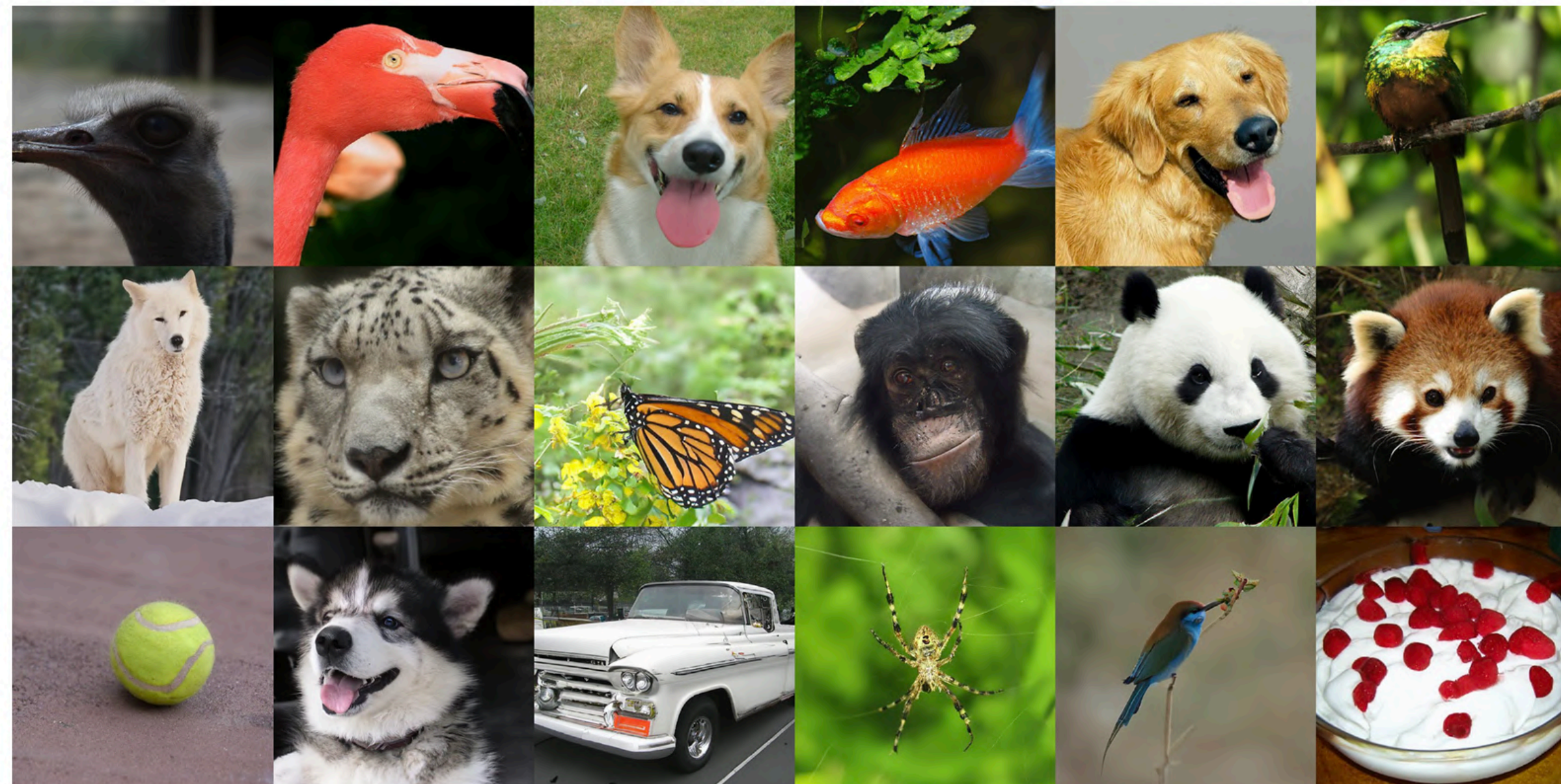


Vision Generative Model Timeline

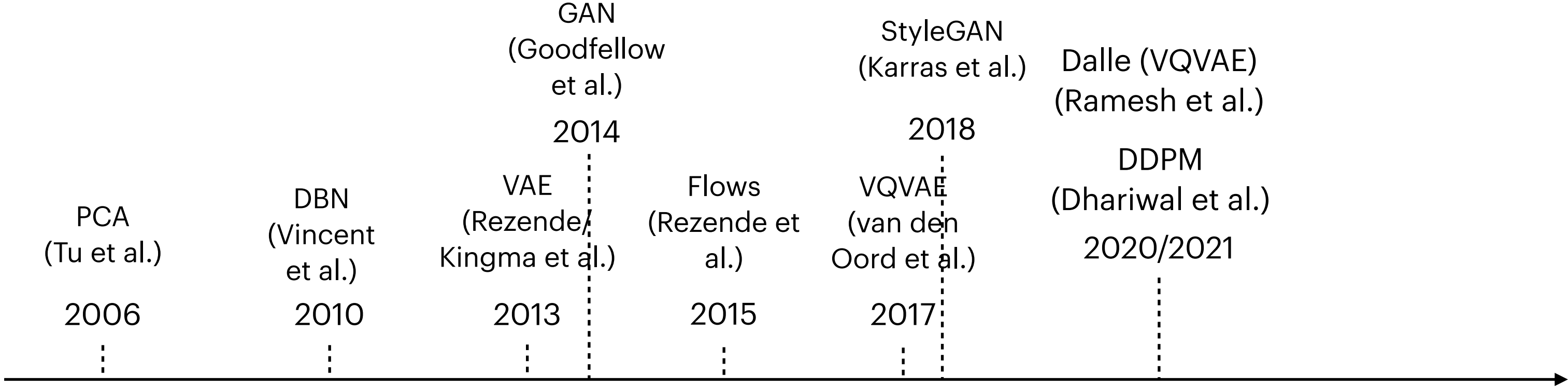
Diffusion Models Beat GANs on Image Synthesis

Prafulla Dhariwal*
OpenAI
prafulla@openai.com

Alex Nichol*
OpenAI
alex@openai.com



Vision Generative Model Timeline



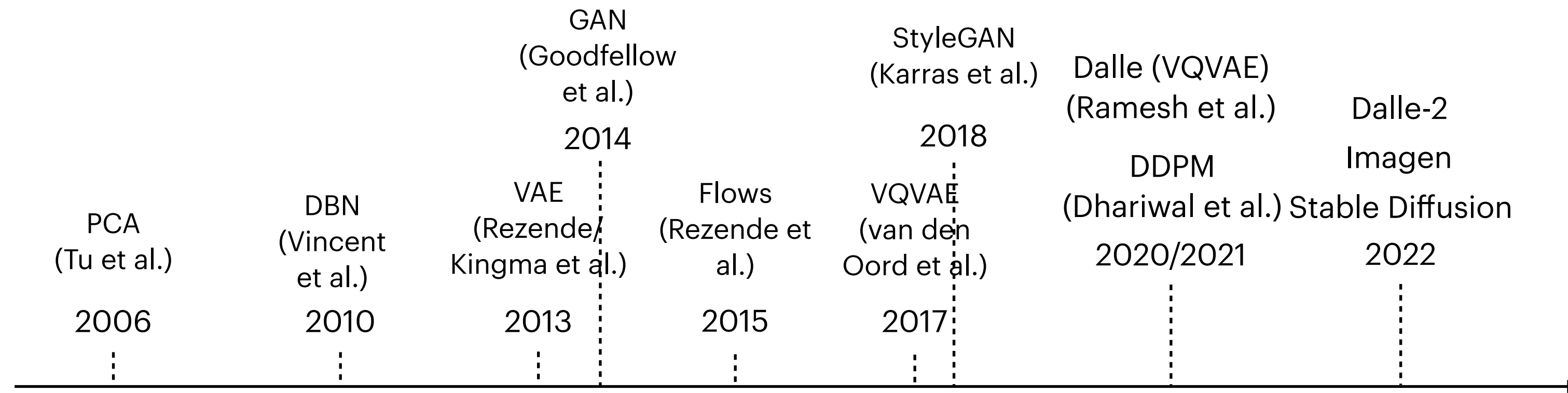
Vision Generative Model Timeline

Zero-Shot Text-to-Image Generation

Aditya Ramesh¹ Mikhail Pavlov¹ Gabriel Goh¹ Scott Gray¹
Chelsea Voss¹ Alec Radford¹ Mark Chen¹ Ilya Sutskever¹



Vision Generative Model Timeline



Vision Generative Model Timeline



Dalle-2

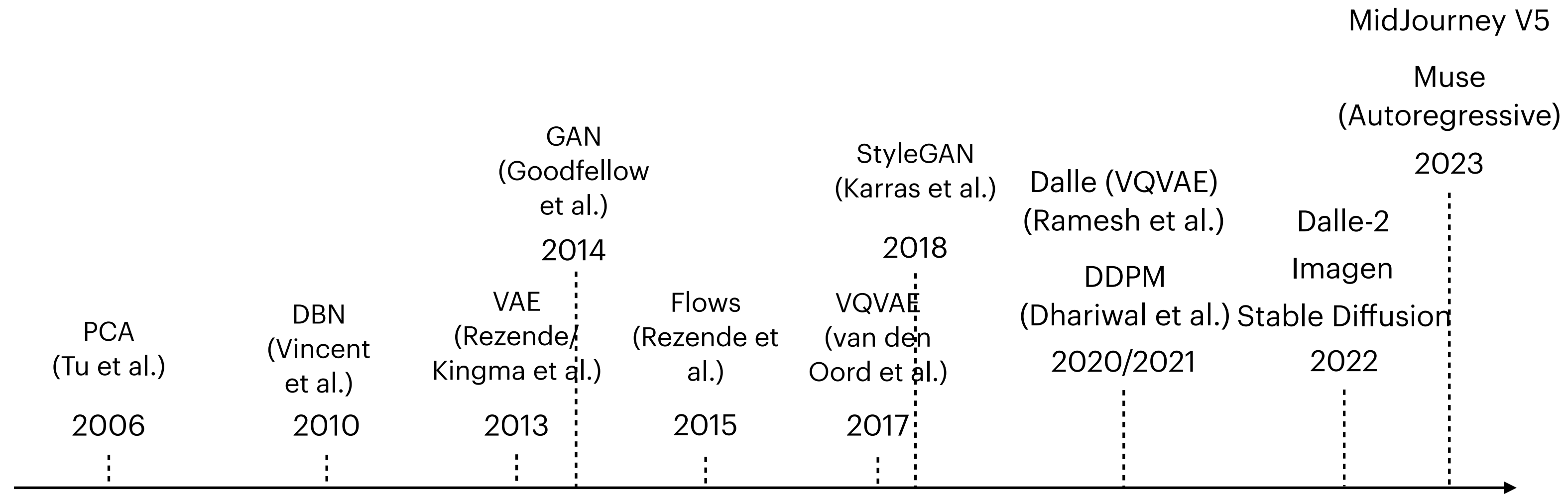


Stable Diffusion



Imagen

Vision Generative Model Timeline



Vision Generative Model Timeline



2006

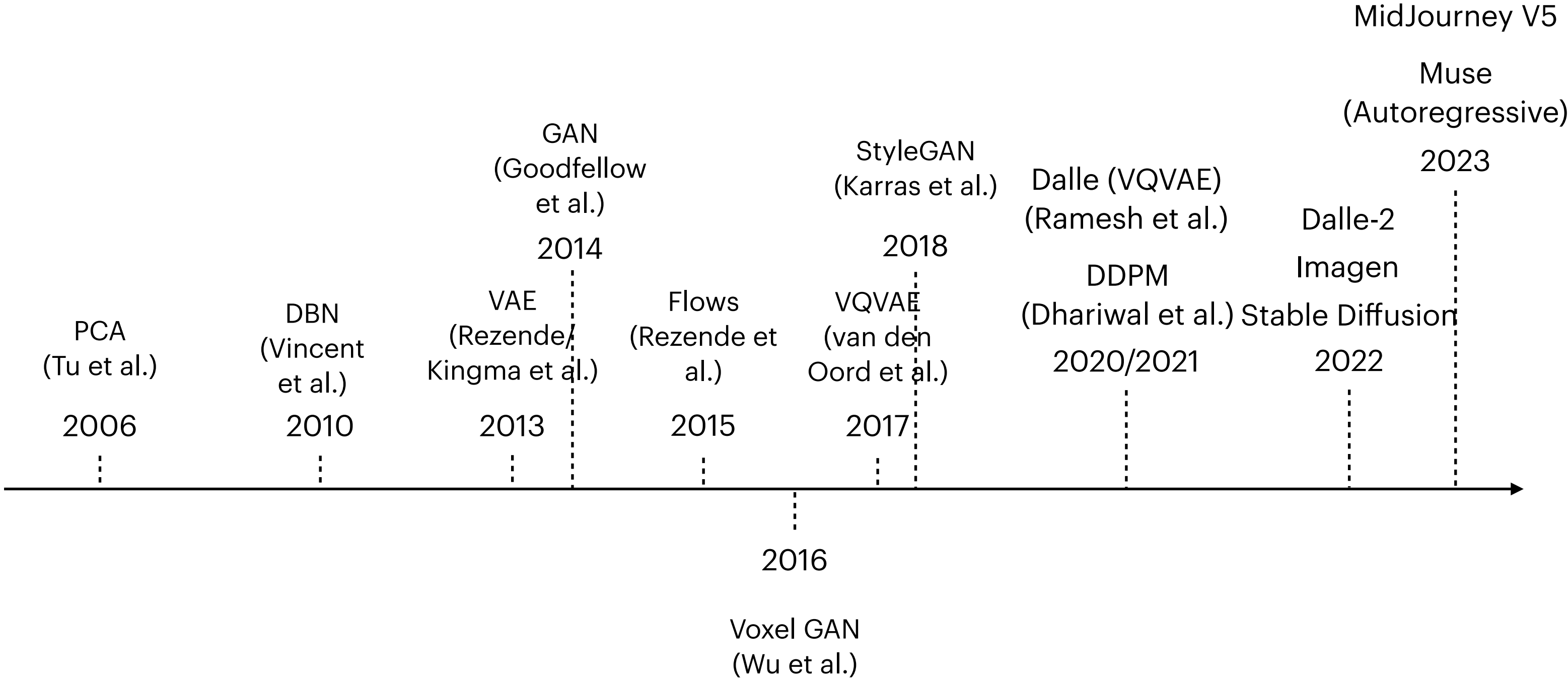


2018



2023

Vision Generative Model Timeline



Vision Generative Model Timeline

Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling

Jiajun Wu*
MIT CSAIL

Chengkai Zhang*
MIT CSAIL

Tianfan Xue
MIT CSAIL

William T. Freeman
MIT CSAIL, Google Research

Joshua B. Tenenbaum
MIT CSAIL

Our results ($64 \times 64 \times 64$)

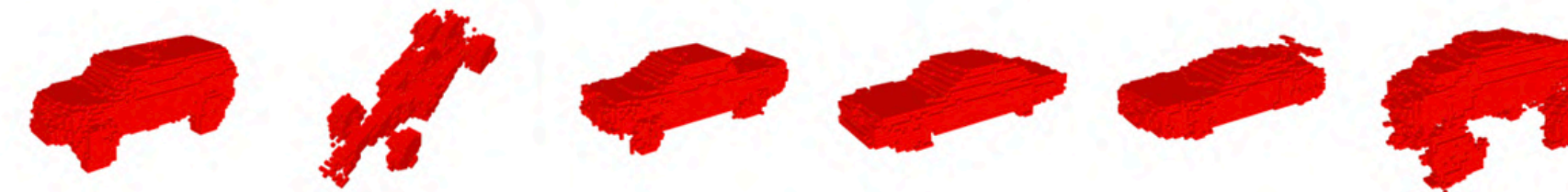
Gun



Chair



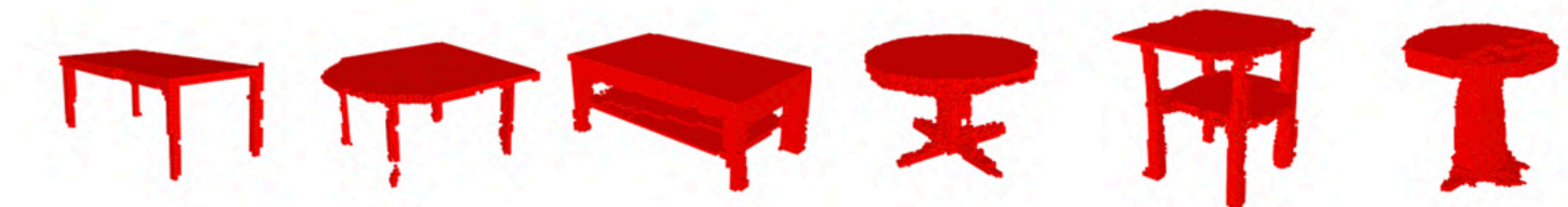
Car



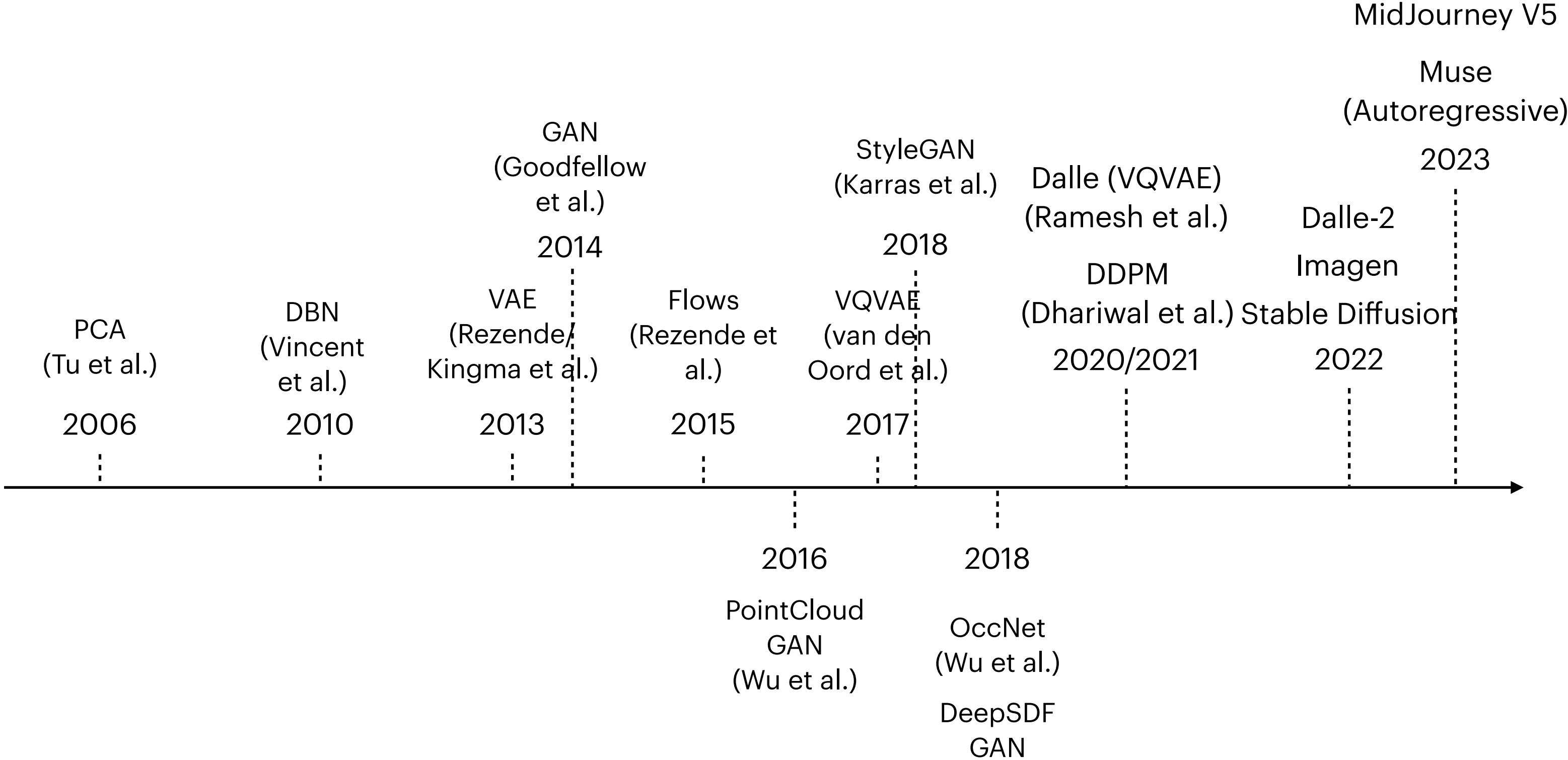
Sofa



Table



Vision Generative Model Timeline



Vision Generative Model Timeline

Occupancy Networks: Learning 3D Reconstruction in Function Space

Lars Mescheder¹ Michael Oechsle^{1,2} Michael Niemeyer¹ Sebastian Nowozin^{3†} Andreas Geiger¹

¹Autonomous Vision Group, MPI for Intelligent Systems and University of Tübingen

²ETAS GmbH, Stuttgart

³Google AI Berlin

{firstname.lastname}@tue.mpg.de nowozin@gmail.com

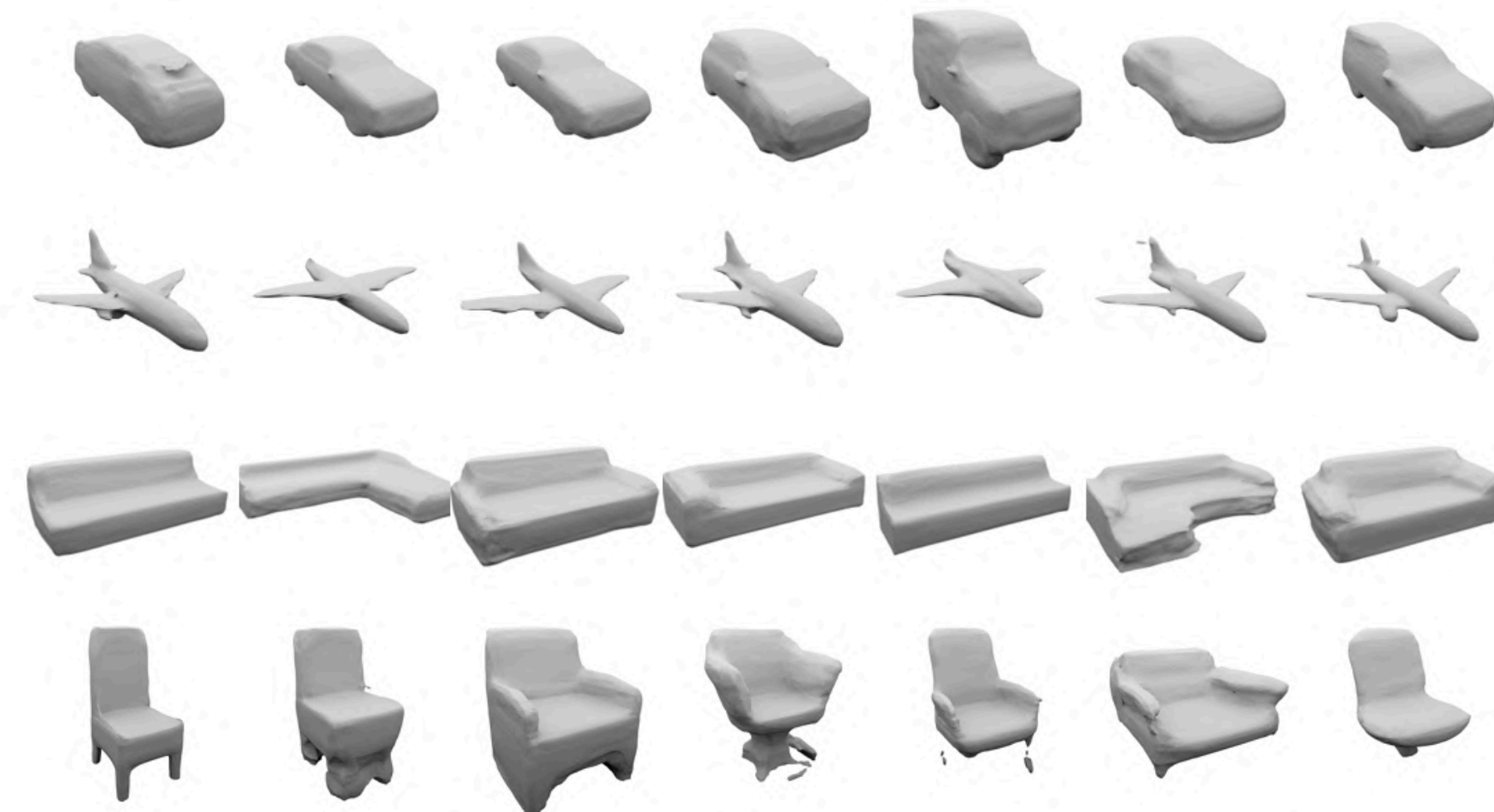
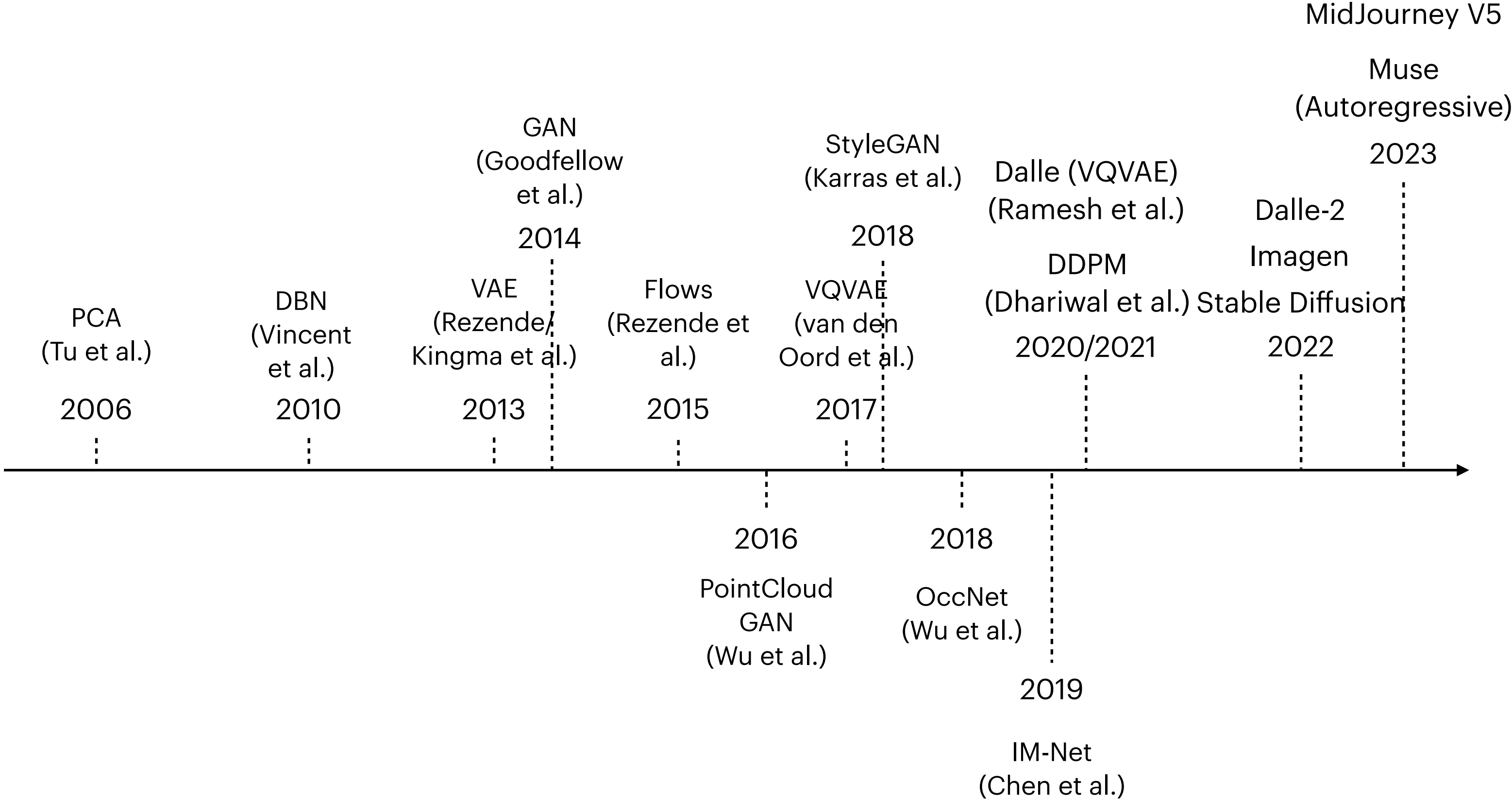


Figure 7: **Unconditional 3D Samples.** Random samples of our unsupervised models trained on the categories “car“, “airplane“, “sofa“ and “chair“ of the ShapeNet dataset. We see that our models are able to capture the distribution of 3D objects and produce compelling new samples.



Vision Generative Model Timeline



Learning Implicit Fields for Generative Shape Modeling

Zhiqin Chen
Simon Fraser University

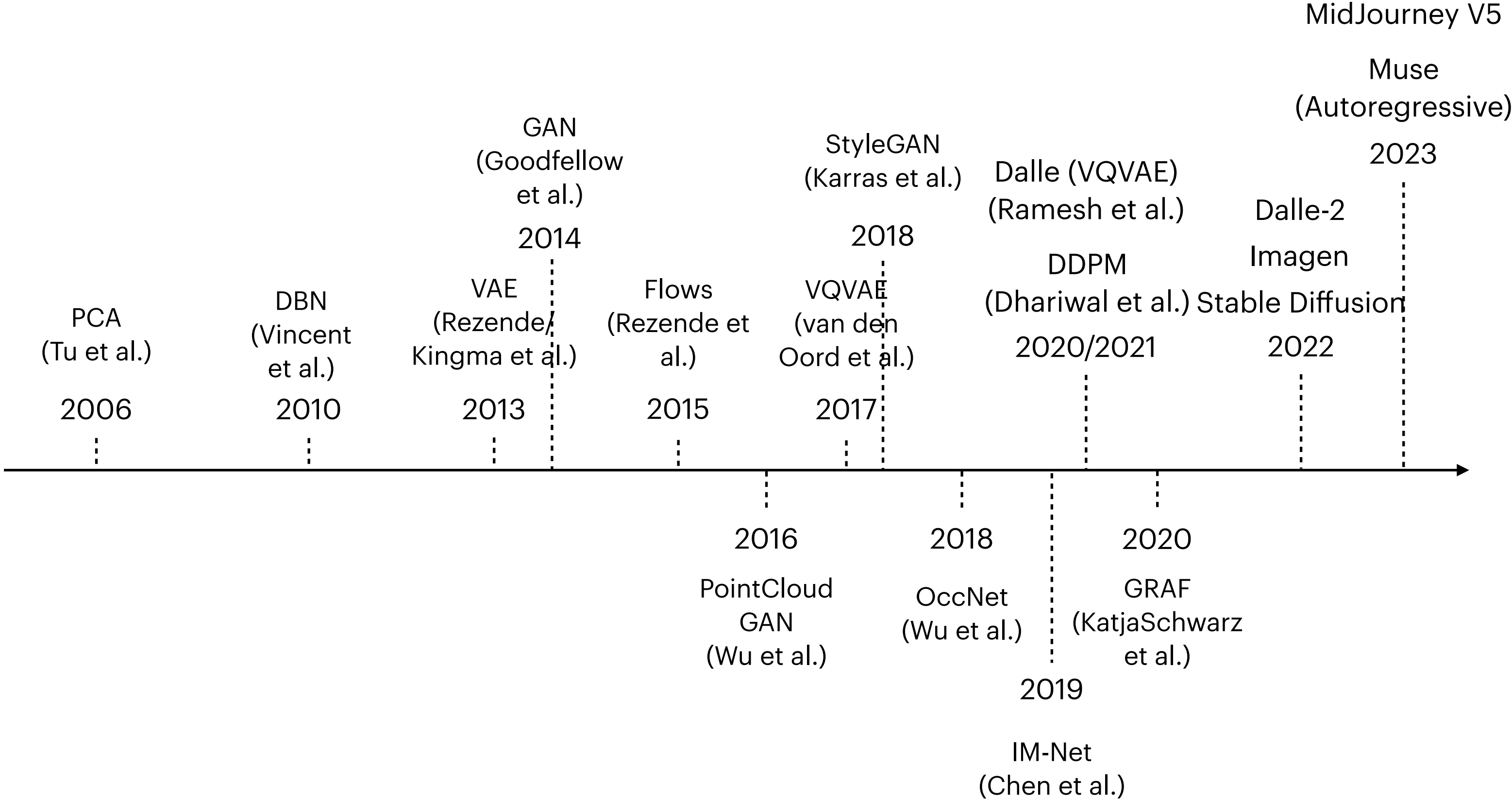
zhiqinc@sfu.ca

Hao Zhang
Simon Fraser University

haoz@sfu.ca

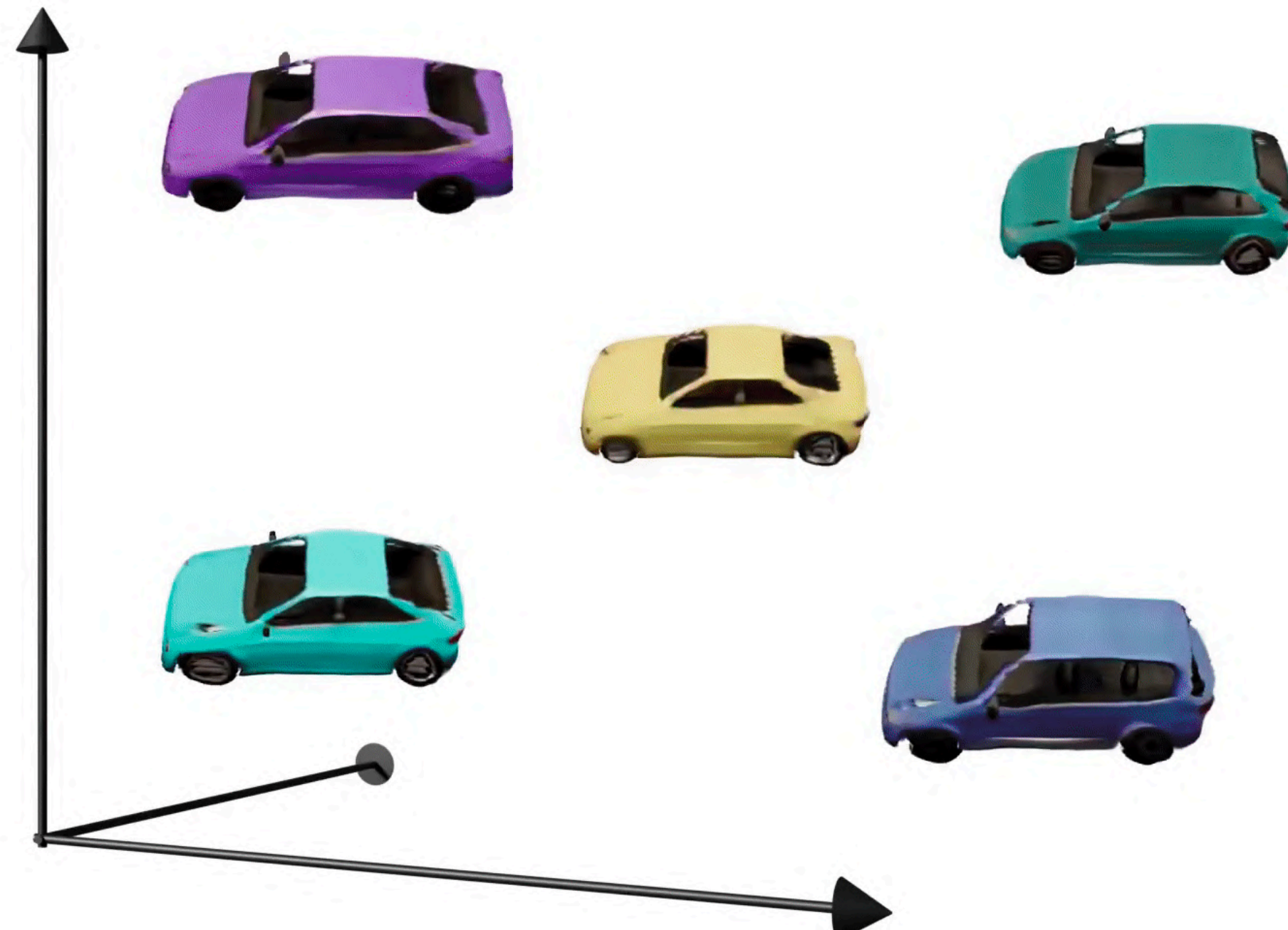


Vision Generative Model Timeline

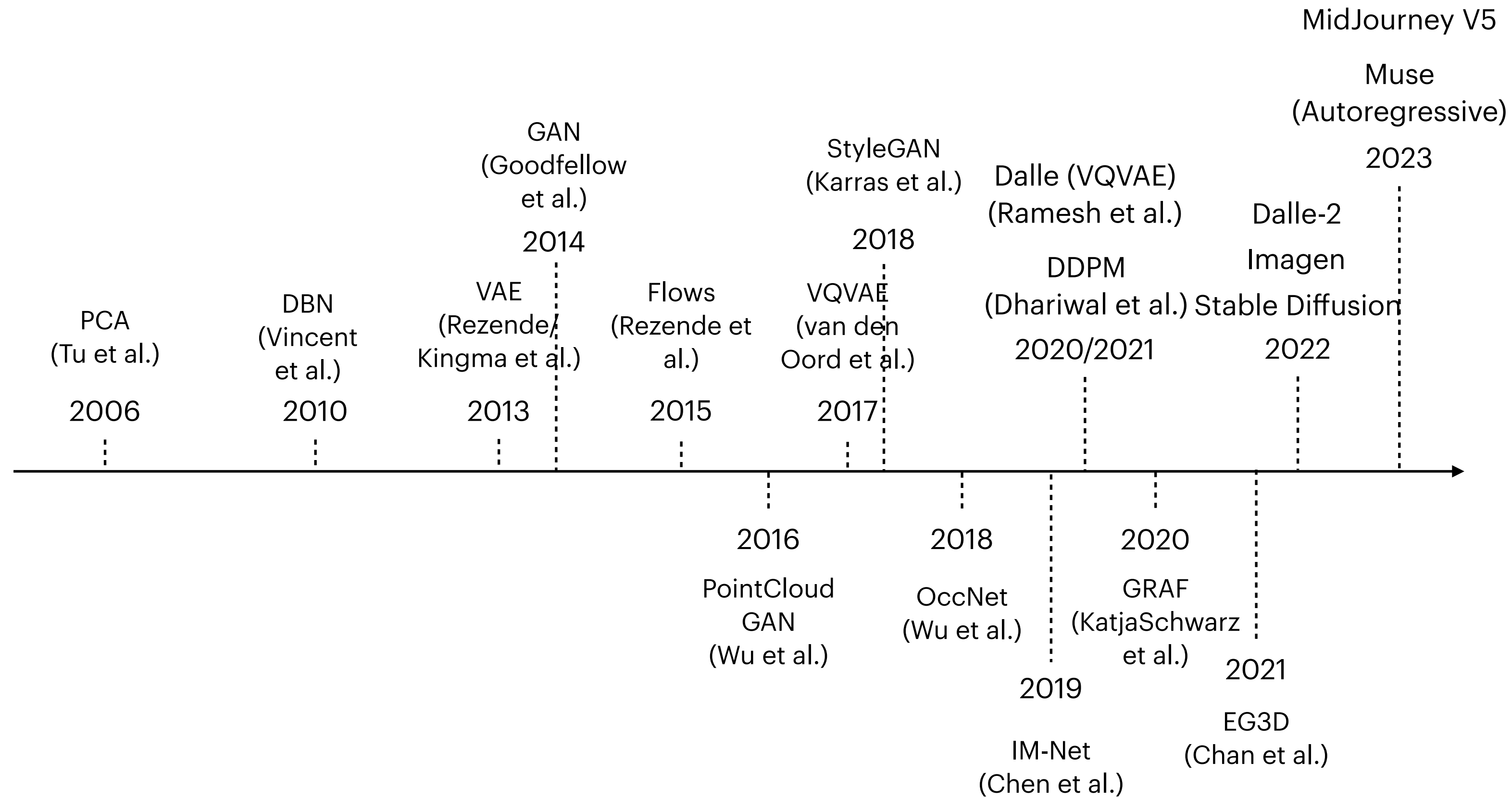


GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis

Katja Schwarz* **Yiyi Liao*** **Michael Niemeyer** **Andreas Geiger**
Autonomous Vision Group
MPI for Intelligent Systems and University of Tübingen
`{firstname.lastname}@tue.mpg.de`



Vision Generative Model Timeline

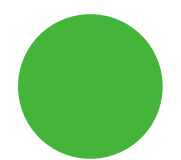
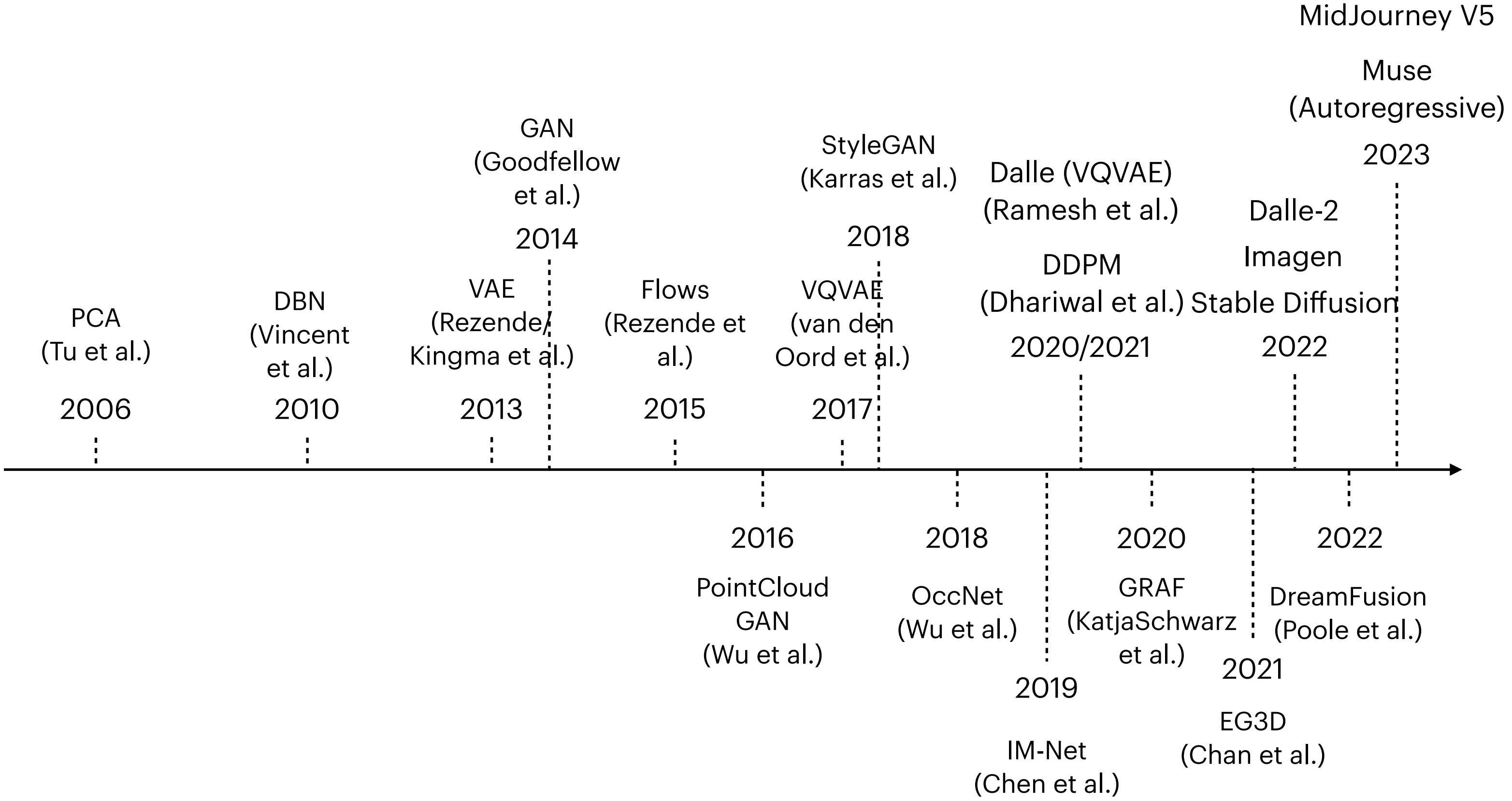


Efficient Geometry-aware 3D Generative Adversarial Networks

Eric R. Chan ^{*†1,2}, Connor Z. Lin^{*1}, Matthew A. Chan^{*1}, Koki Nagano^{*2}, Boxiao Pan¹, Shalini De Mello², Orazio Gallo², Leonidas Guibas¹, Jonathan Tremblay², Sameh Khamis², Tero Karras², and Gordon Wetzstein¹



Vision Generative Model Timeline



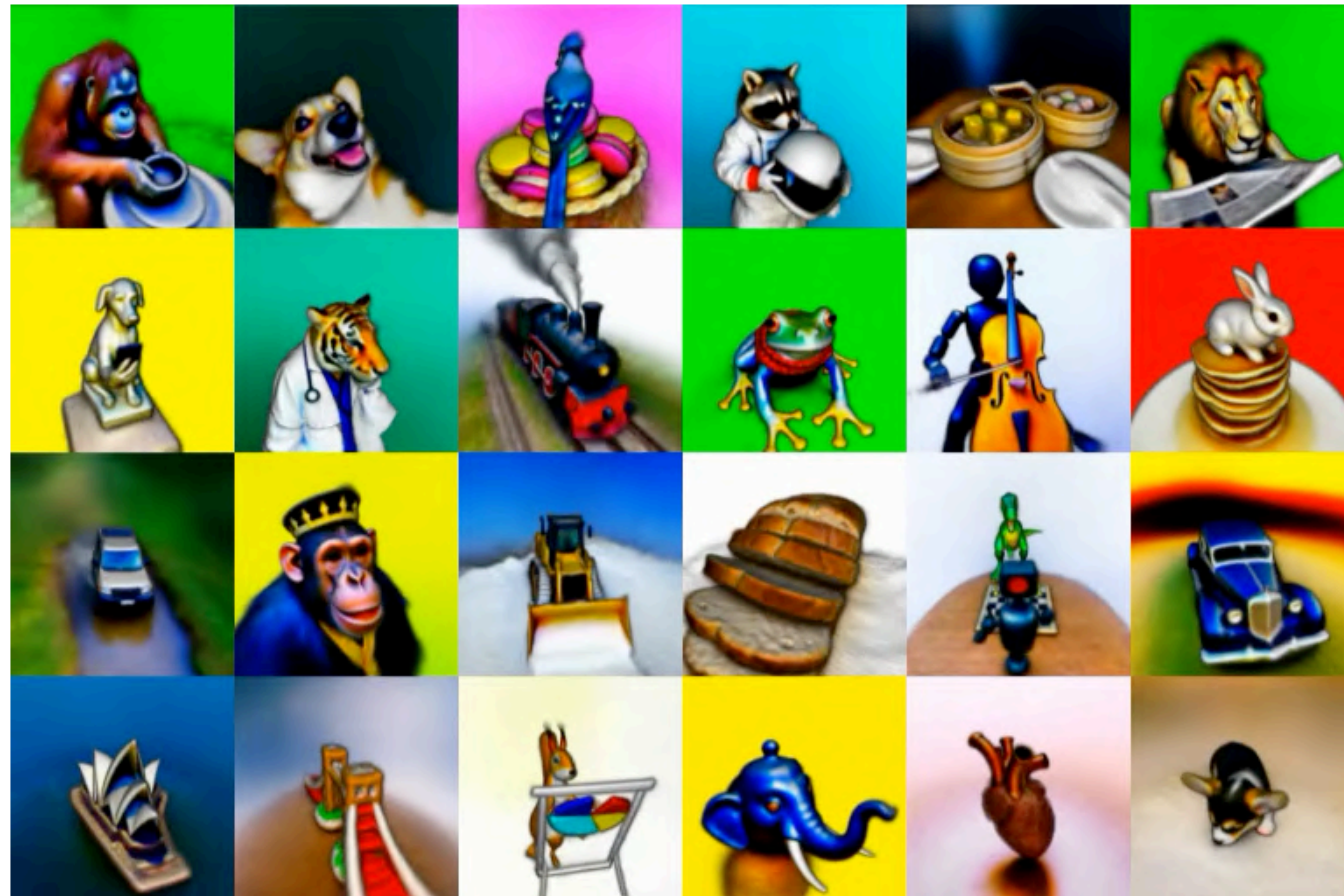
DreamFusion: Text-to-3D using 2D Diffusion

Ben Poole
Google Research

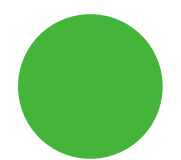
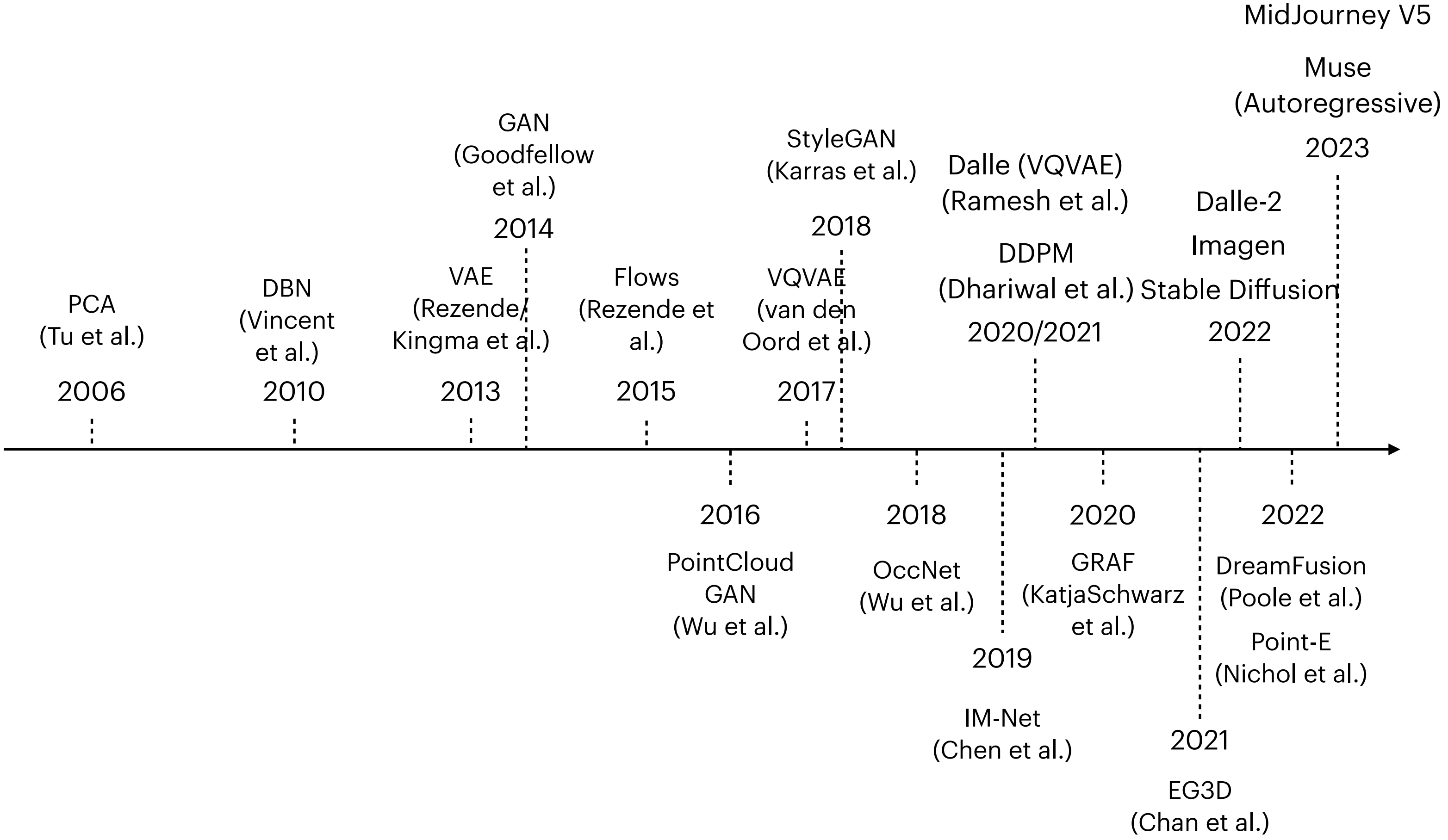
Ajay Jain
UC Berkeley

Jonathan T. Barron
Google Research

Ben Mildenhall
Google Research



Vision Generative Model Timeline



Point-E: A System for Generating 3D Point Clouds from Complex Prompts

Alex Nichol^{*1} Heewoo Jun^{*1} Prafulla Dhariwal¹ Pamela Mishkin¹ Mark Chen¹



“a corgi wearing a red santa hat”



“a multicolored rainbow pumpkin”



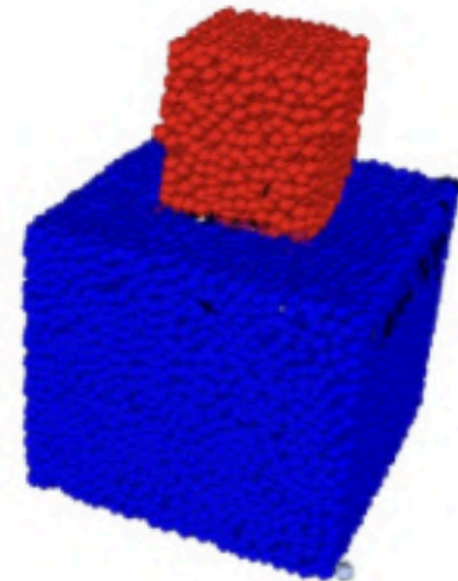
“an elaborate fountain”



“a traffic cone”



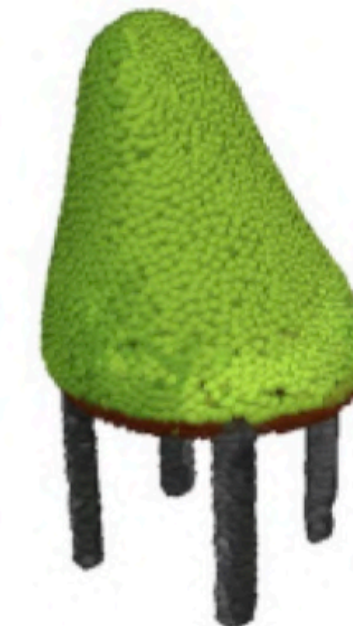
“a vase of purple flowers”



“a small red cube is sitting on top of a large blue cube. red on top, blue on bottom”



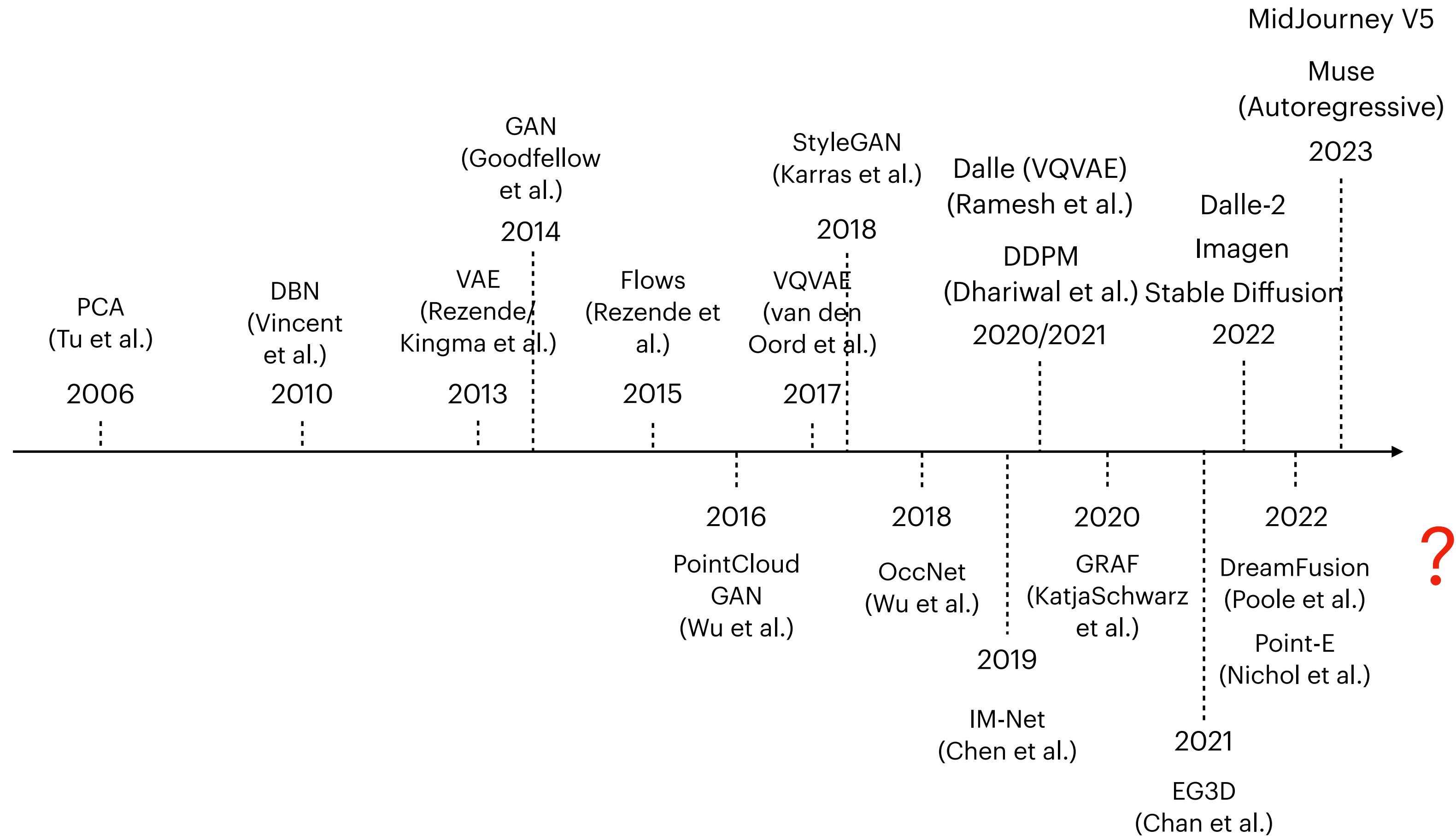
“a pair of 3d glasses, left lens is red right is blue”



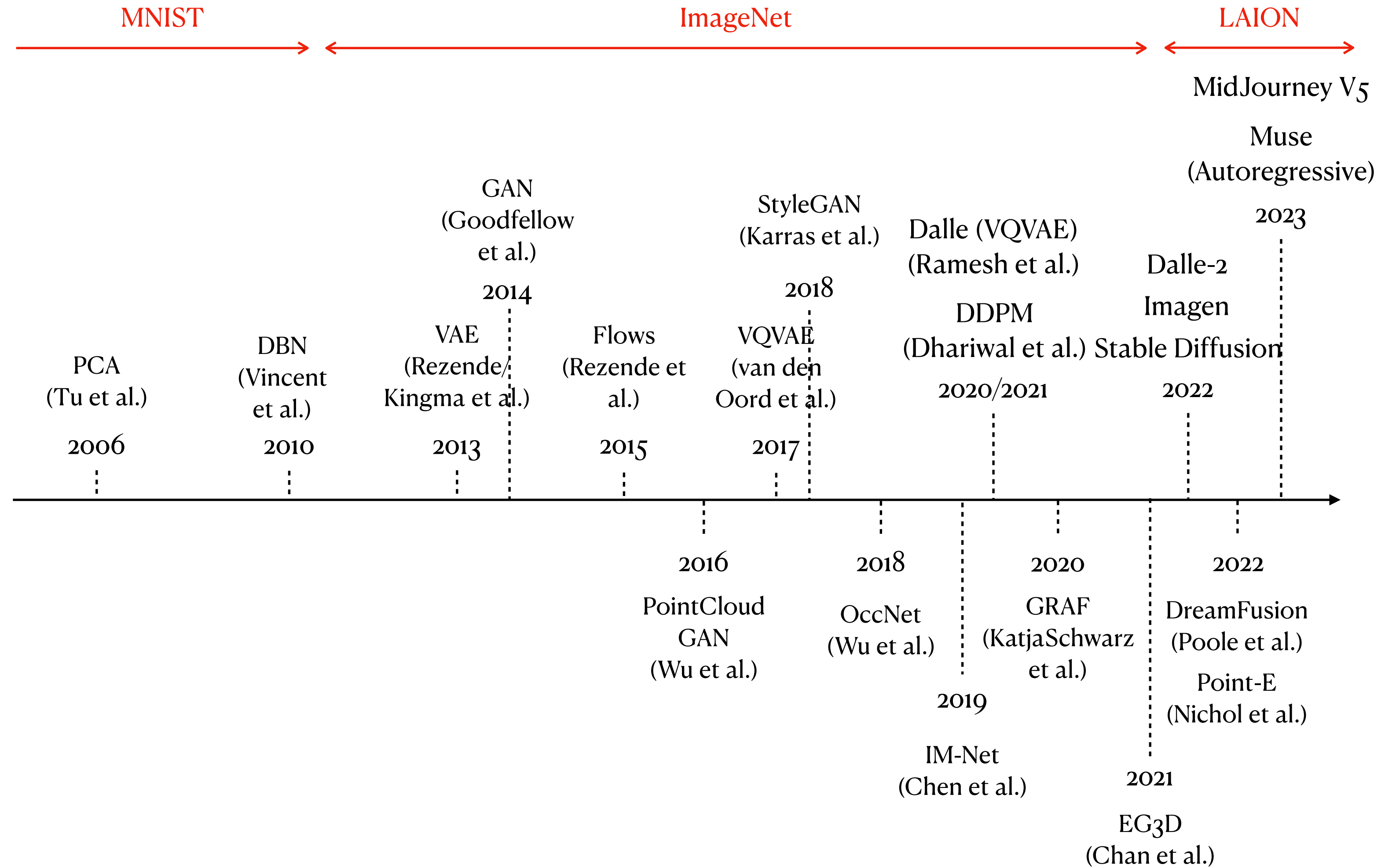
“an avocado chair, a chair imitating an avocado”



What's next?



High-Quality Dataset is the Unseen Hero



Objaverse-XL: 10M+ high-quality 3D assets

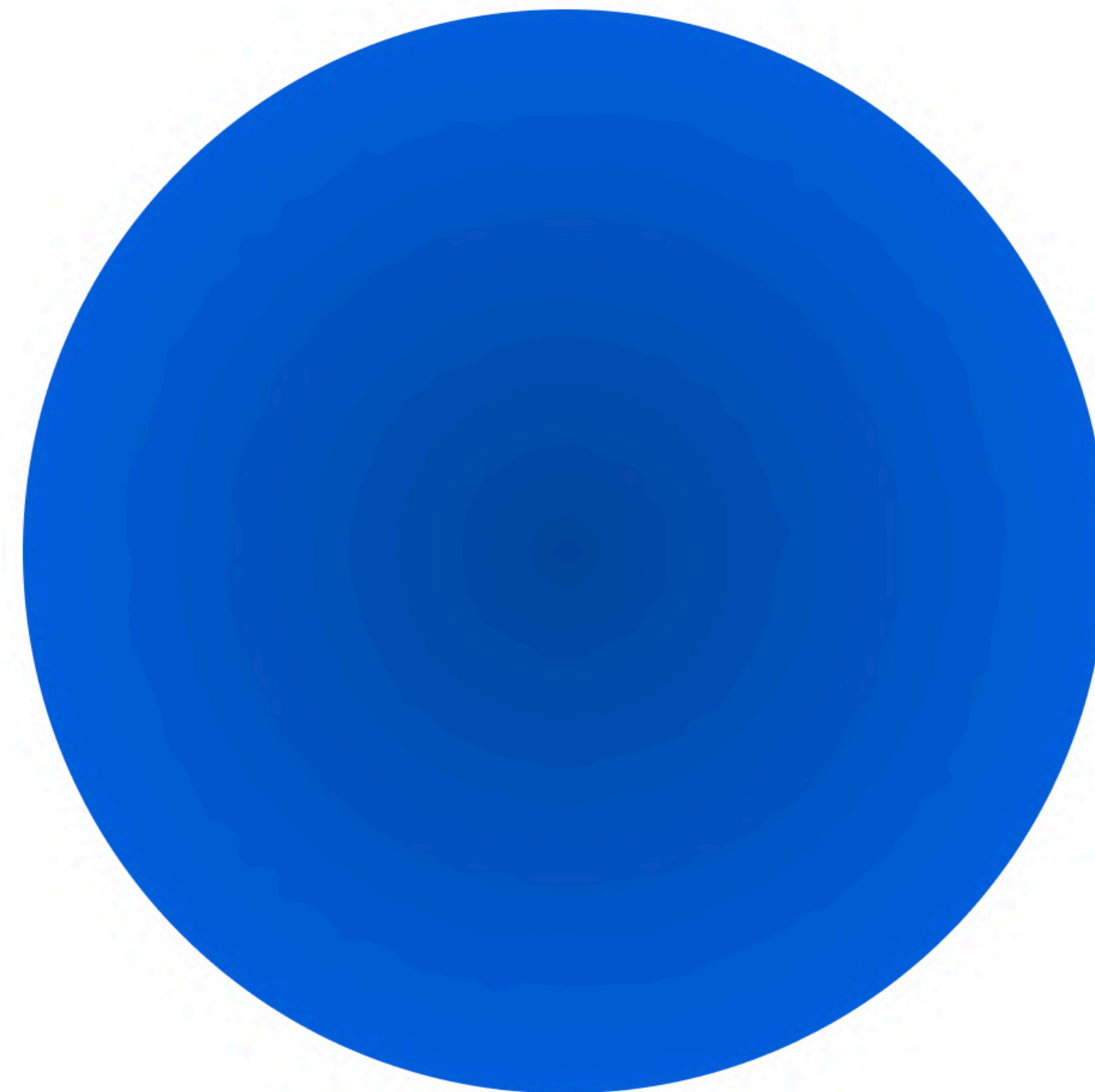


Objaverse-XL: 10M+ high-quality 3D assets

Everything Else
(Combined)



Objaverse-XL



Objaverse-XL: 10M+ high-quality 3D assets

GitHub



Thingiverse



Polycam



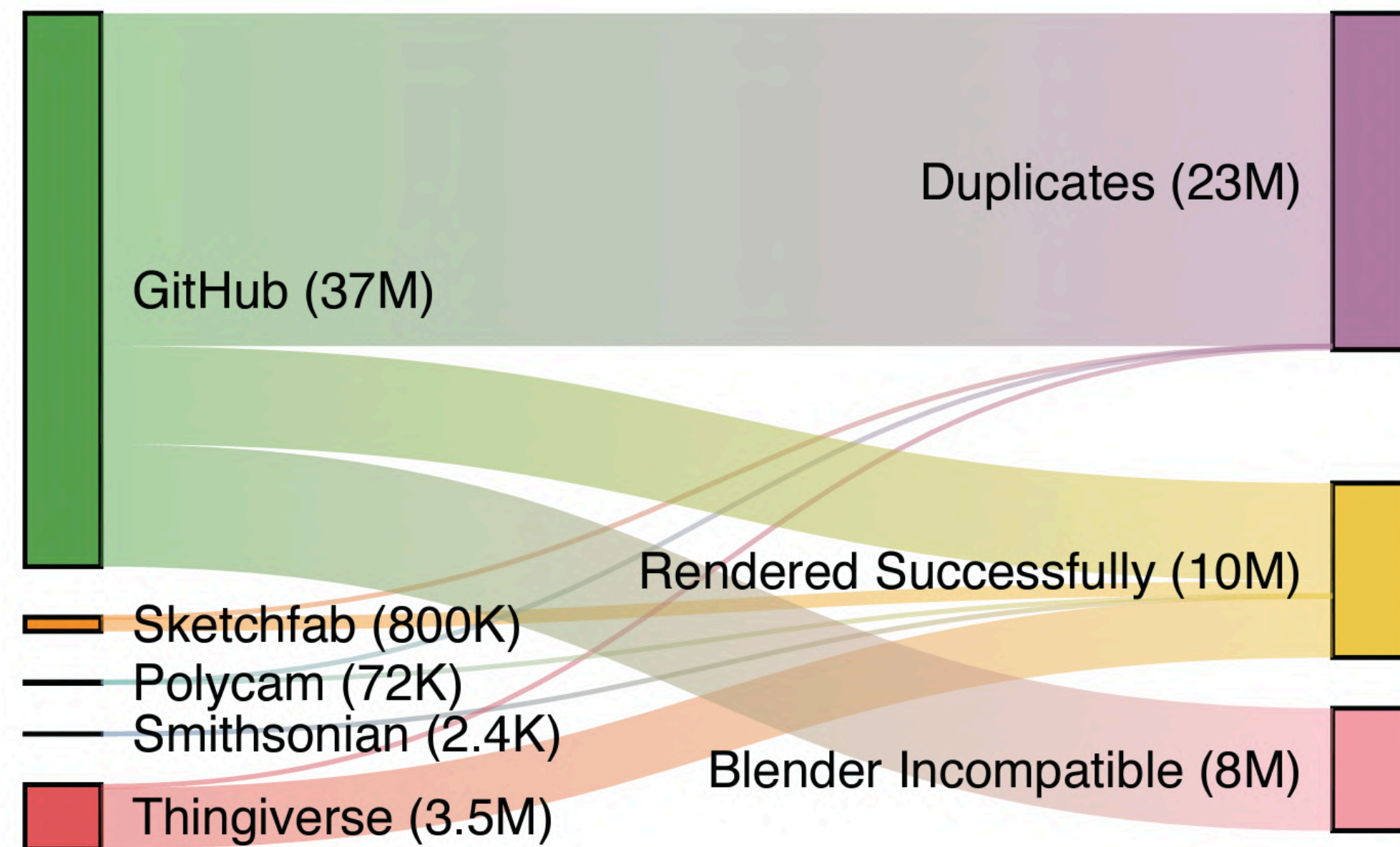
Smithsonian Institution



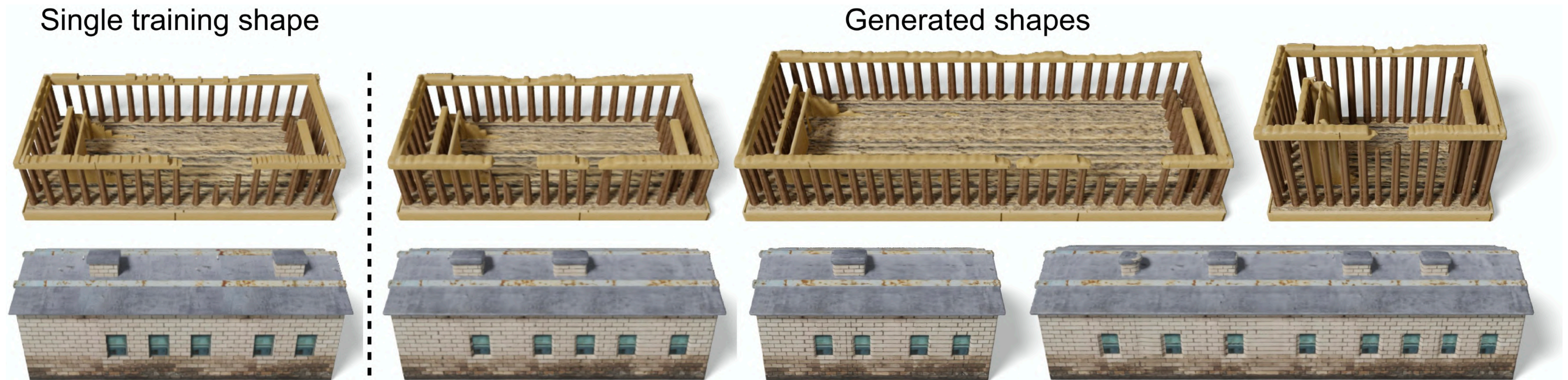
Sketchfab



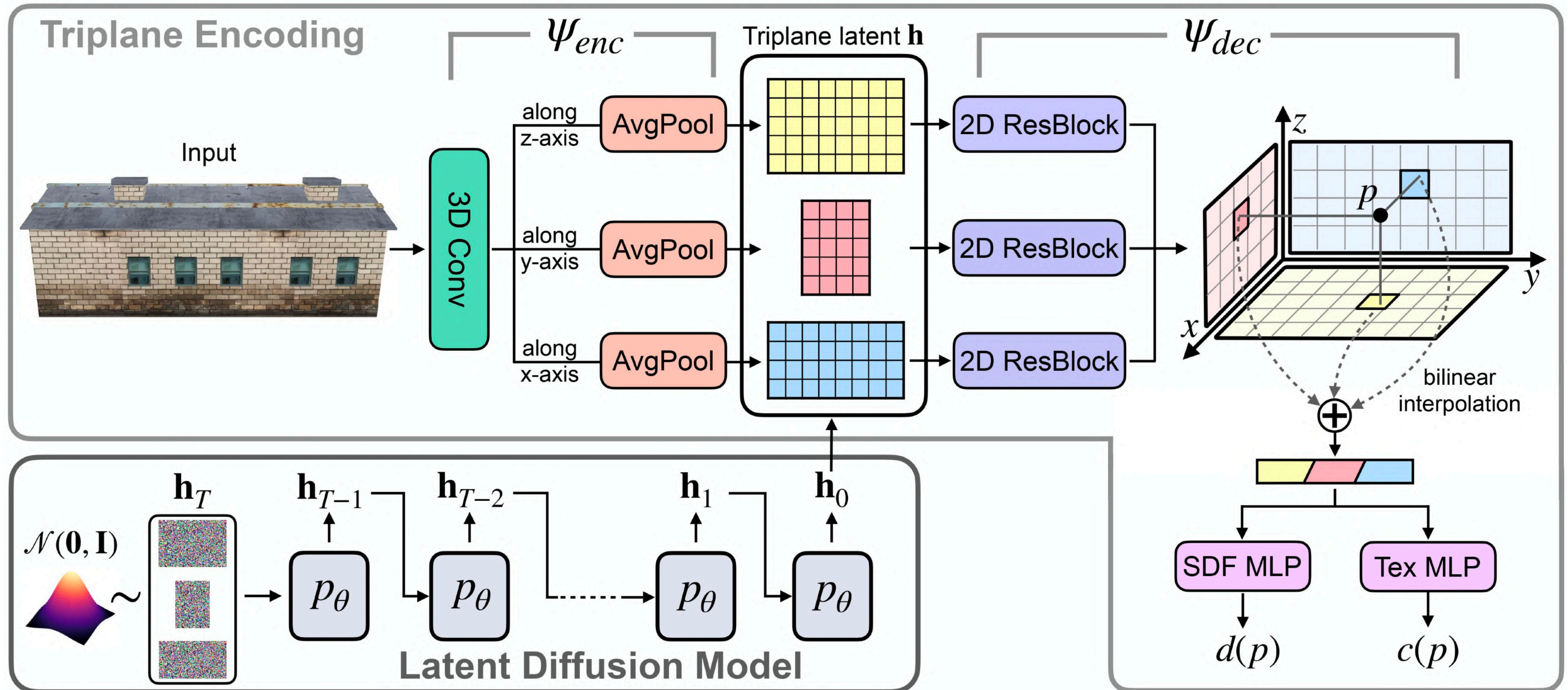
Objaverse-XL: 10M+ high-quality 3D assets



Sin3DM: a method for 3D data augmentation

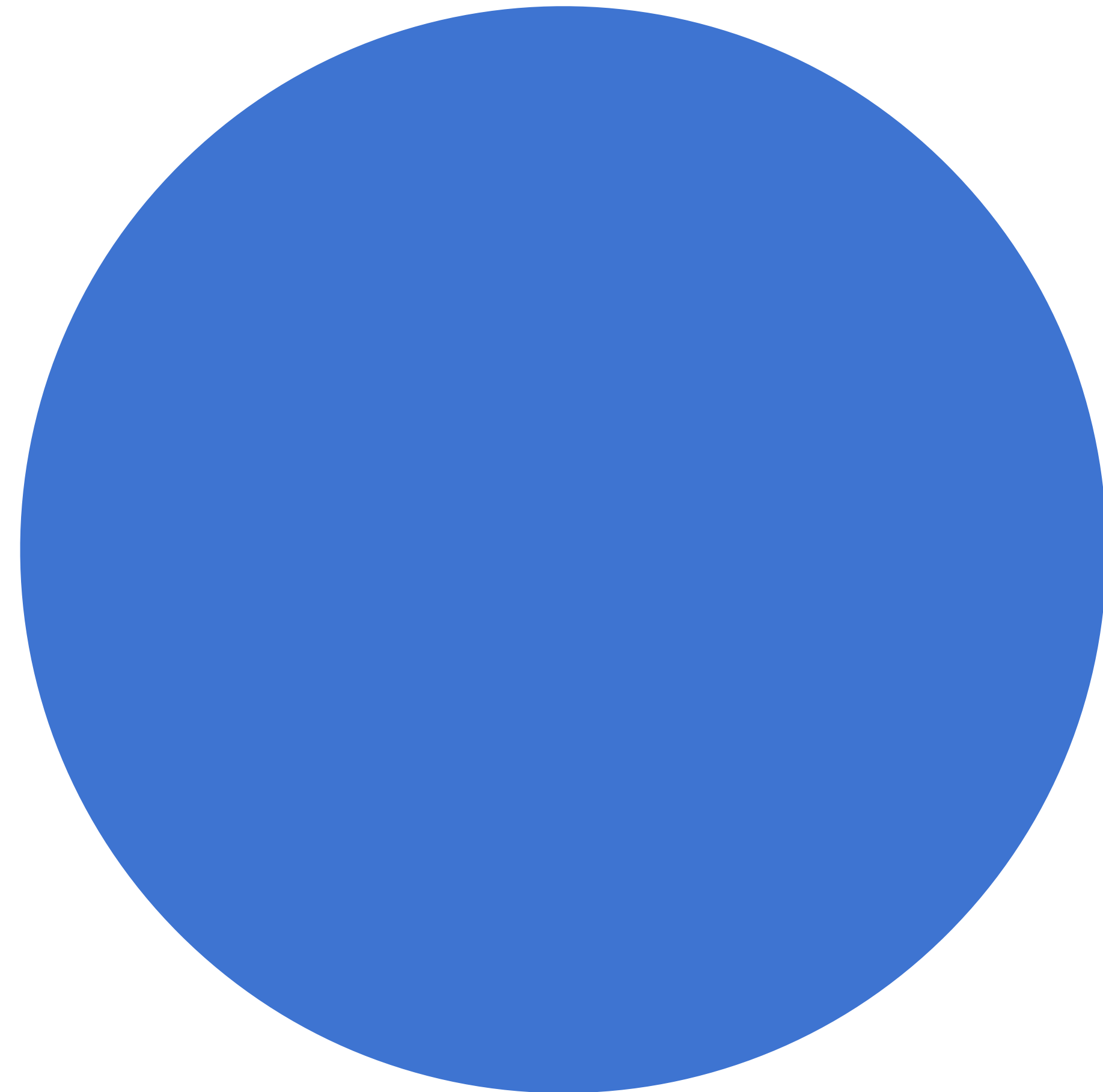


Learning a Diffusion Model from a Single 3D Textured Shape



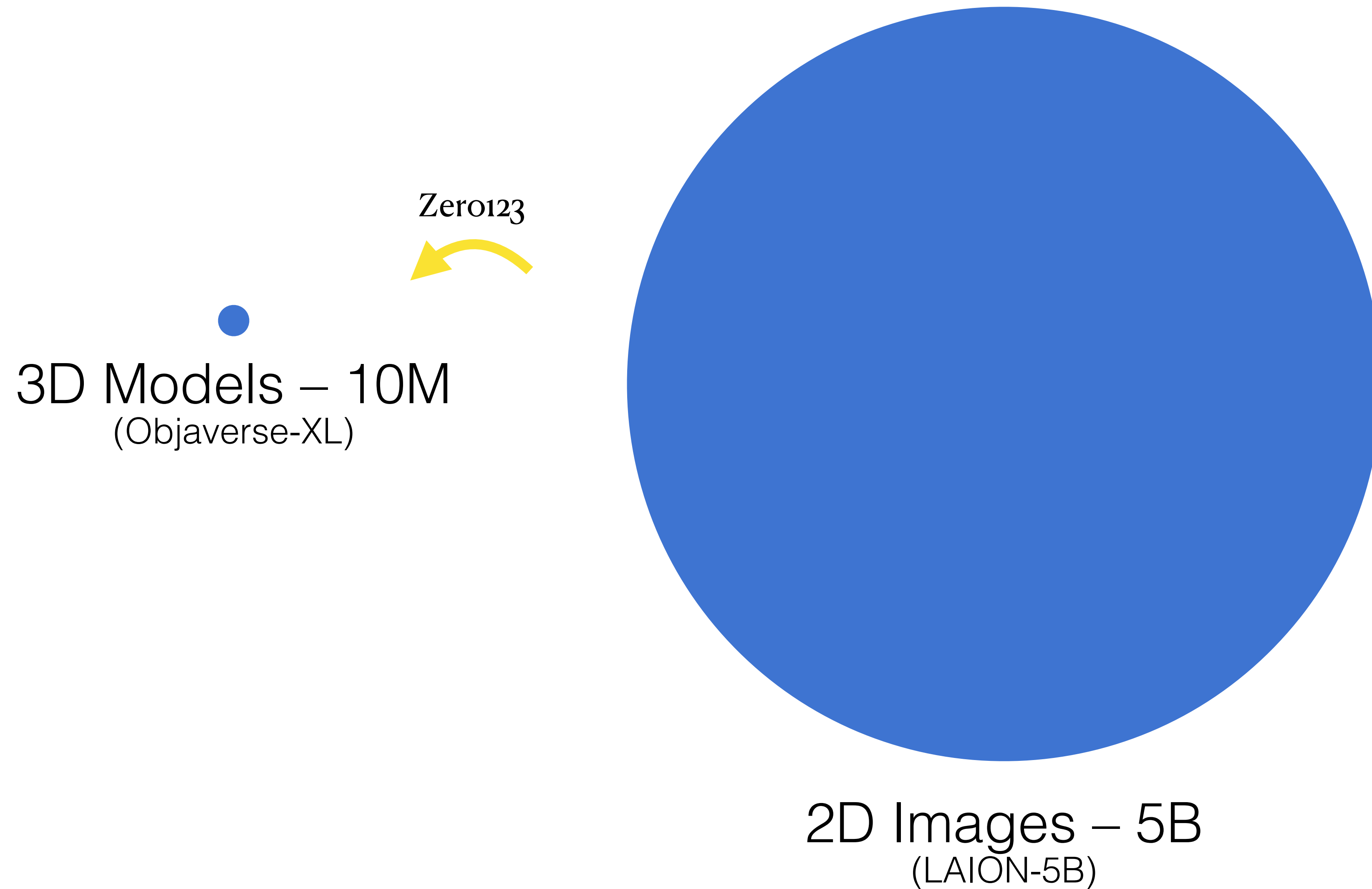
Is Dataset Scaling Enough for 3D Gen?

3D Models – 800K
(Objaverse)



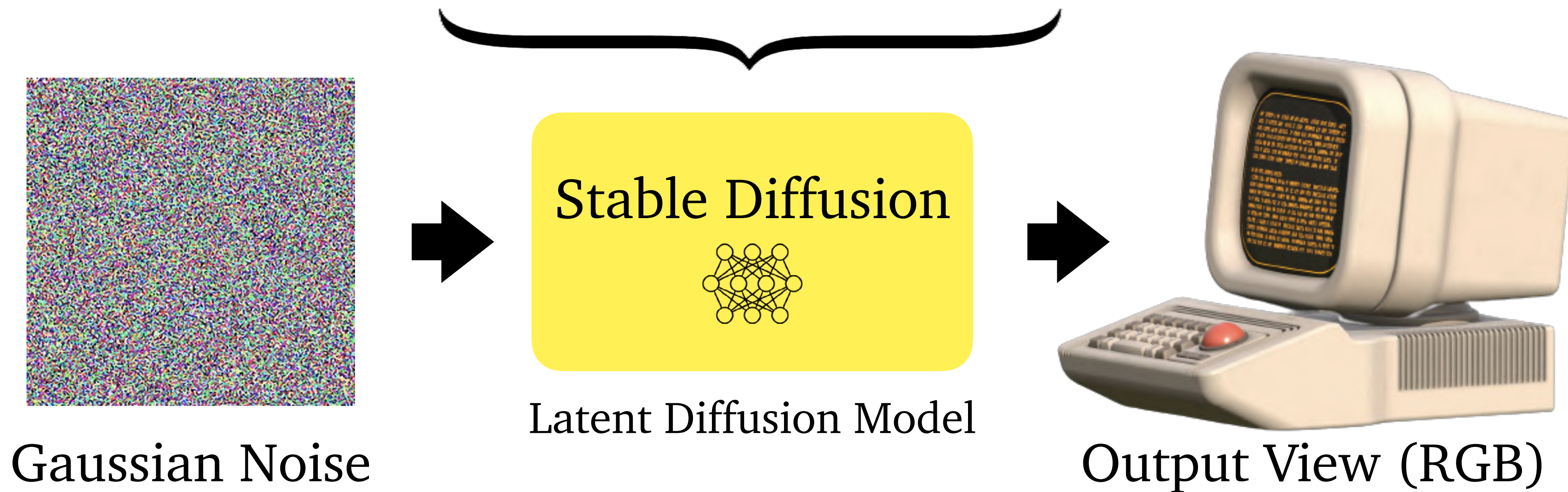
2D Images – 5B
(LAION-5B)

Is Dataset Scaling Enough for 3D Gen?

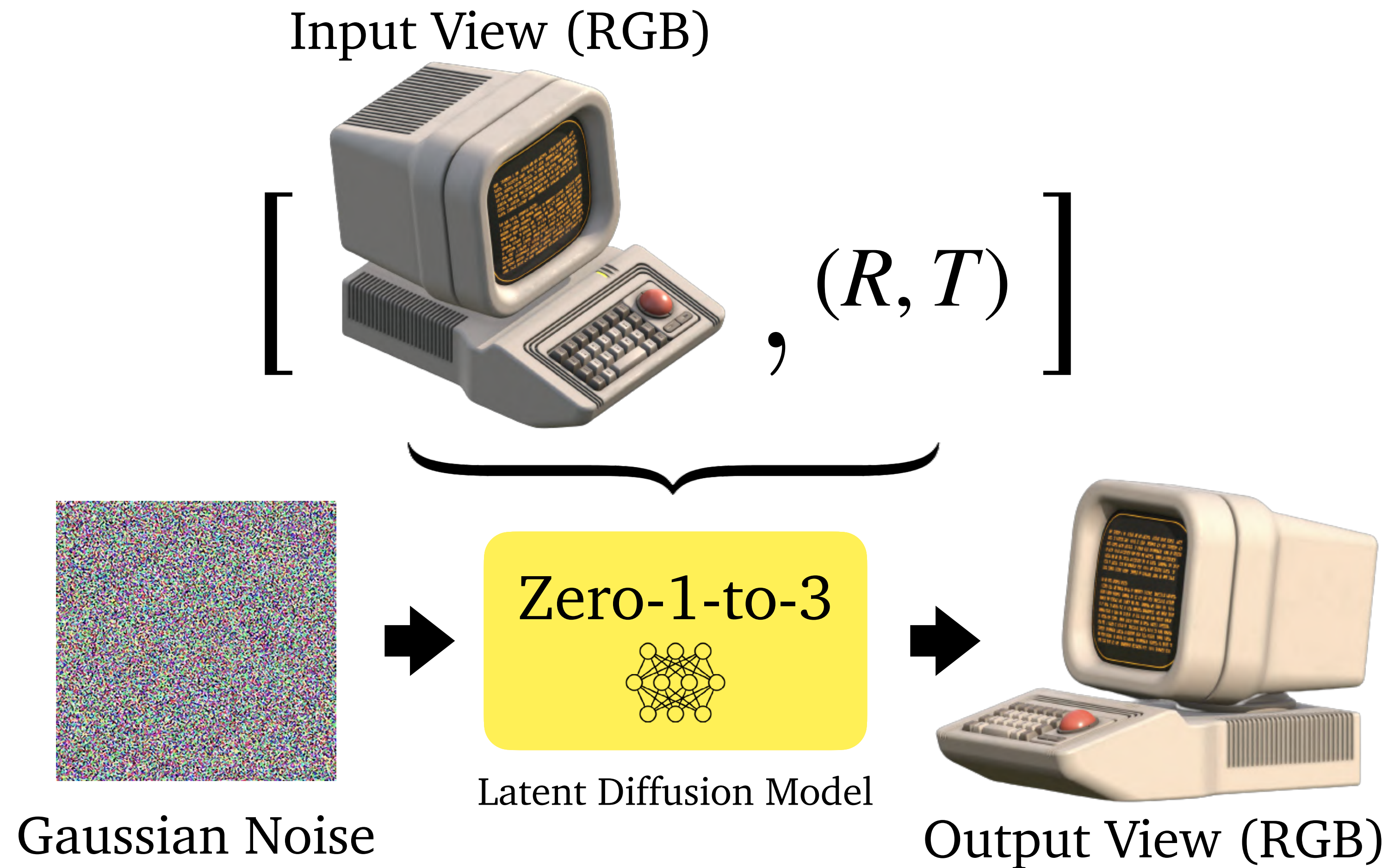


Stable Diffusion

“A retrospective computer”



Zero-1-to-3: Zero-Shot One Image to 3D Objects



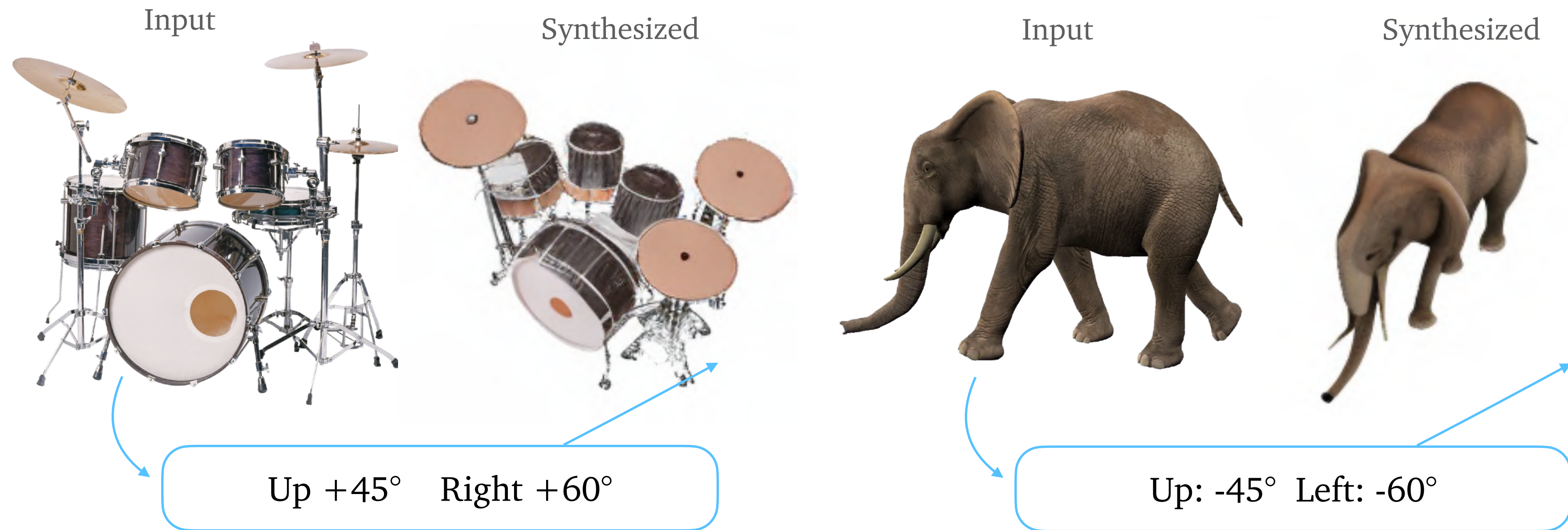
Dataset of Novel Views



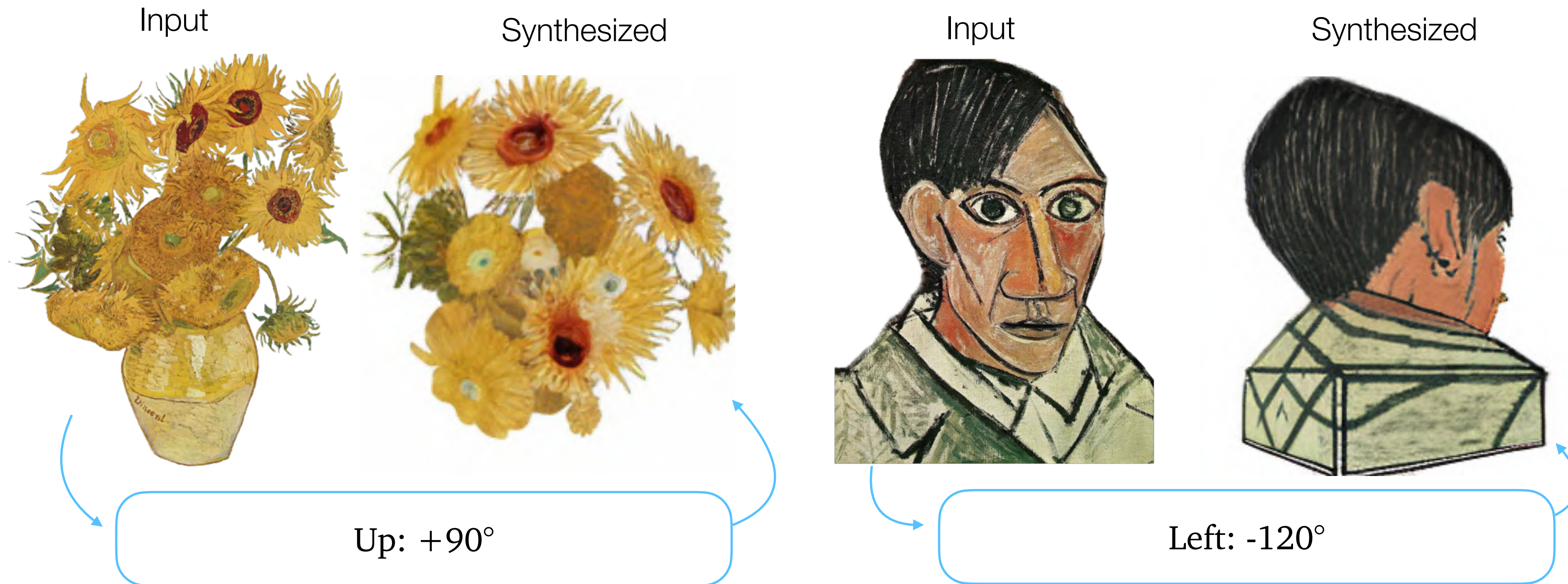
Everyday Objects



Complex/Deformable Objects



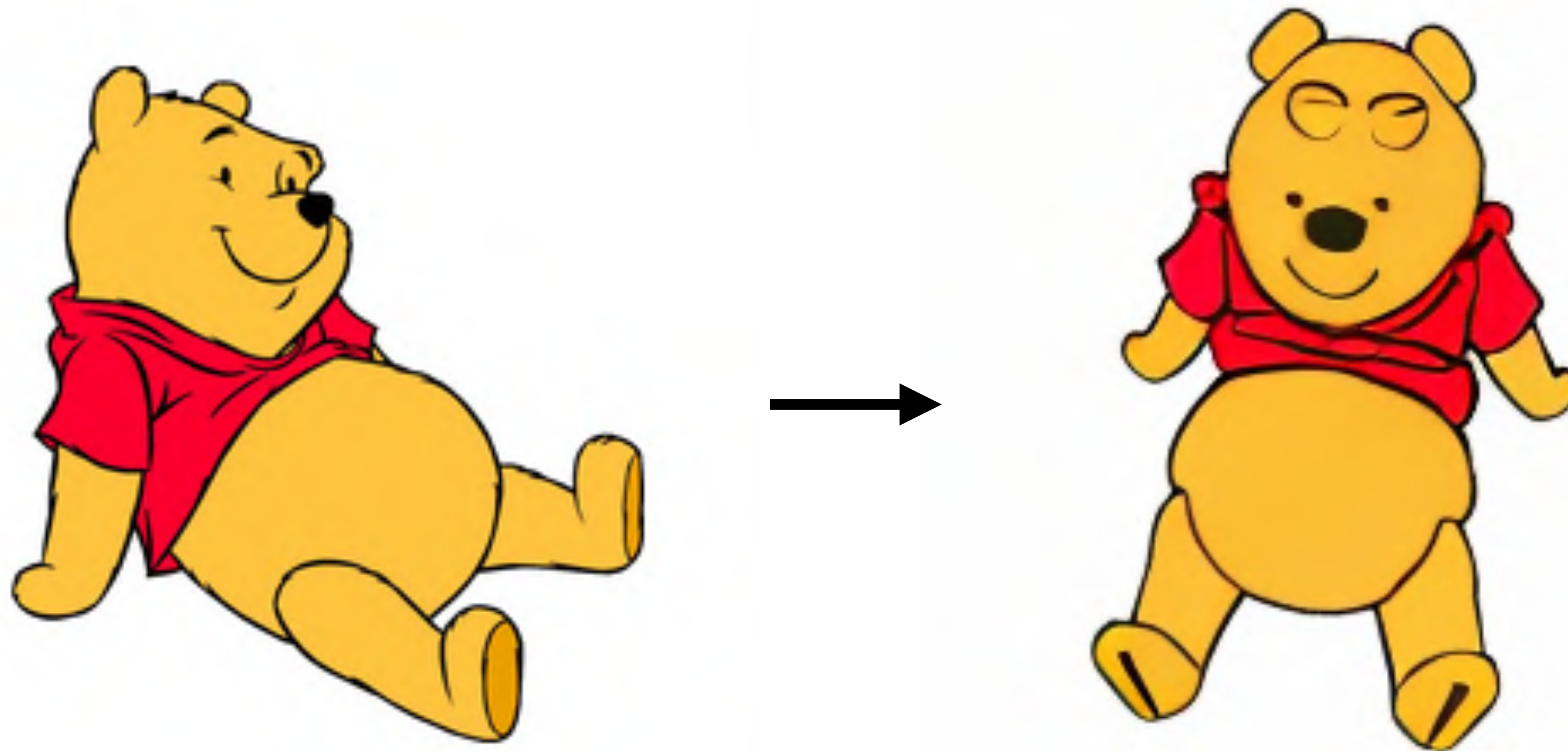
Abstract Paintings



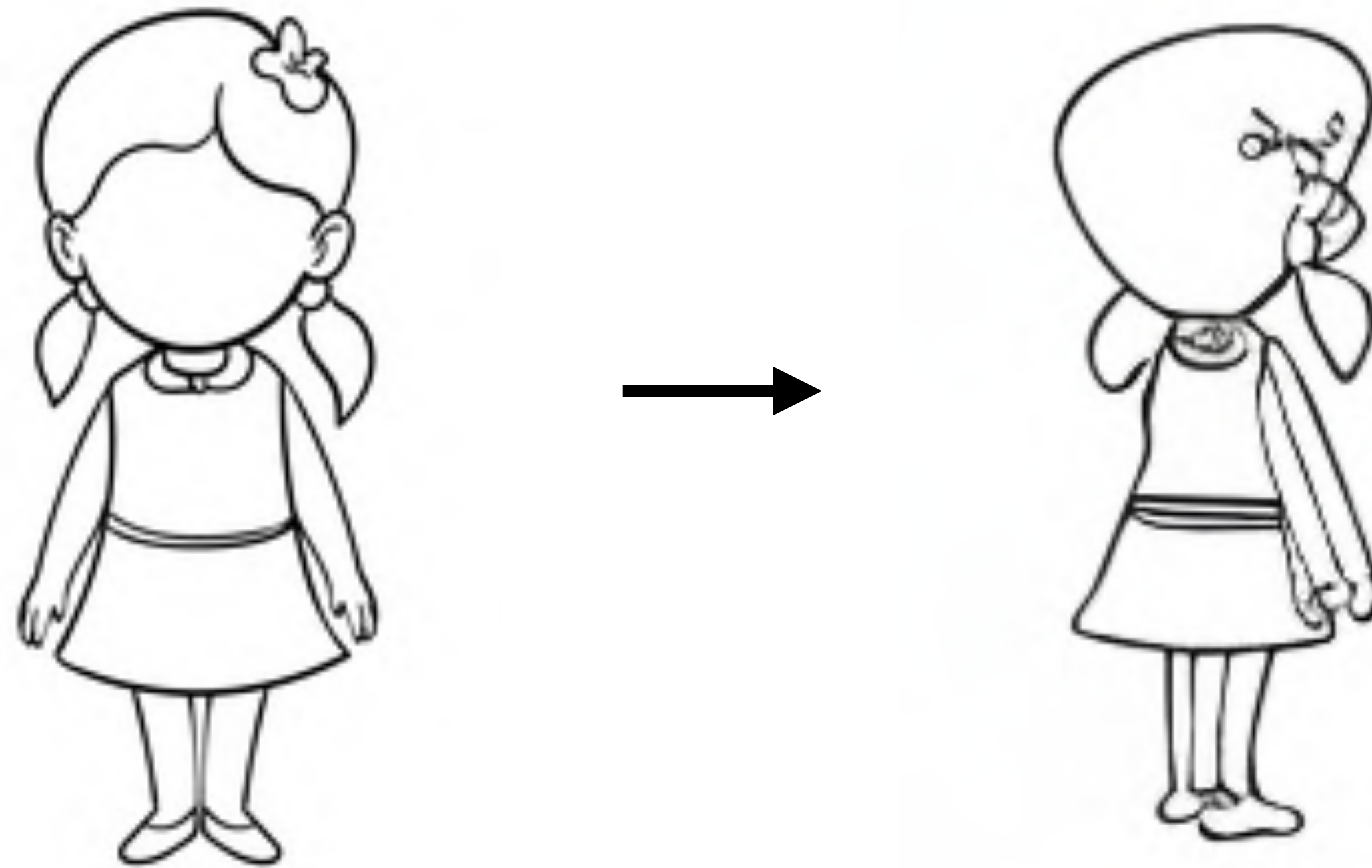
Oil Paintings



Cartoons



Line Drawings



Sketches

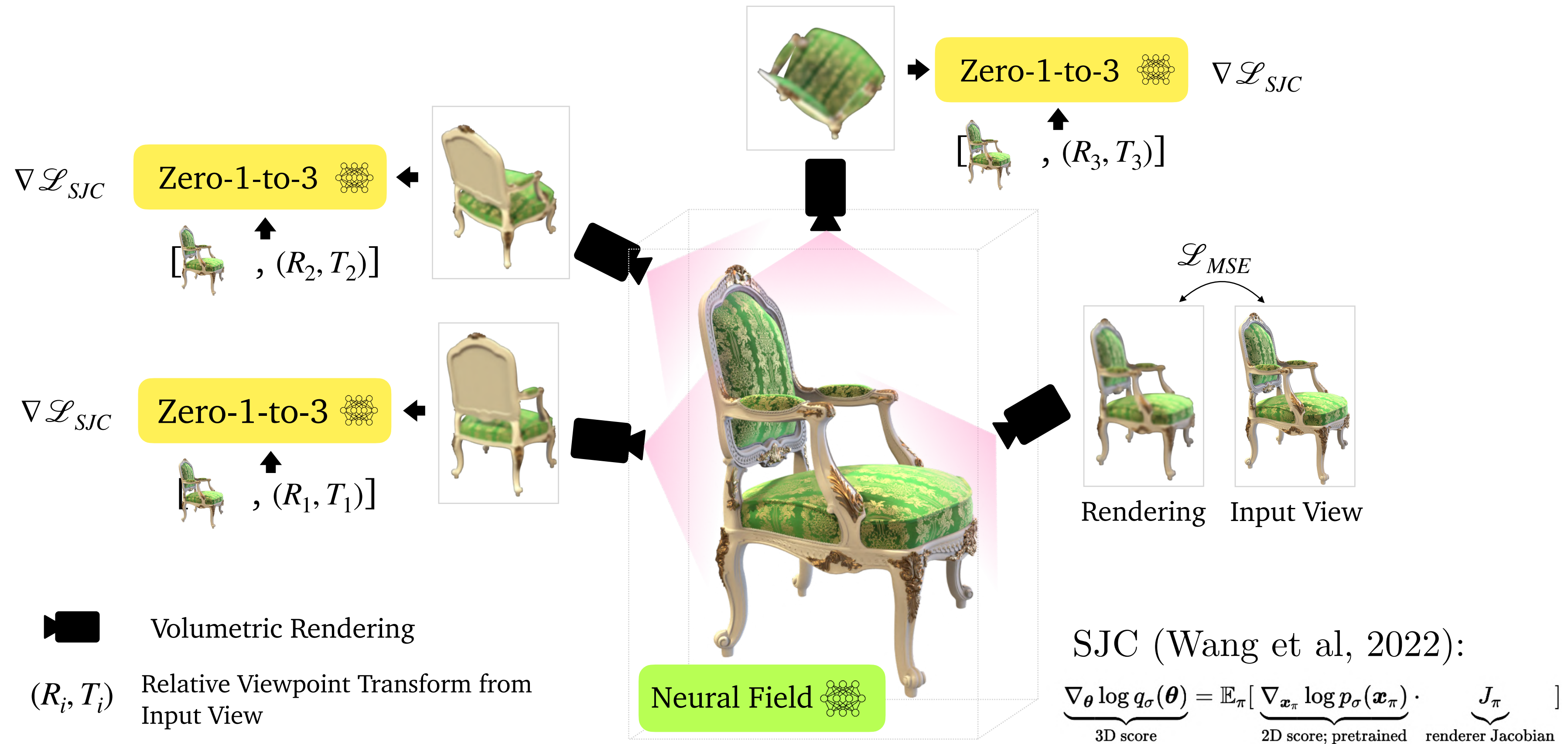


Ambiguity from Occlusion



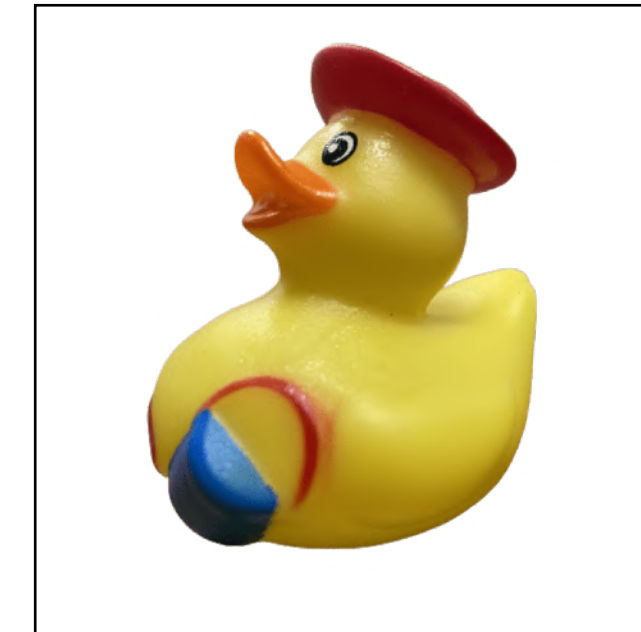
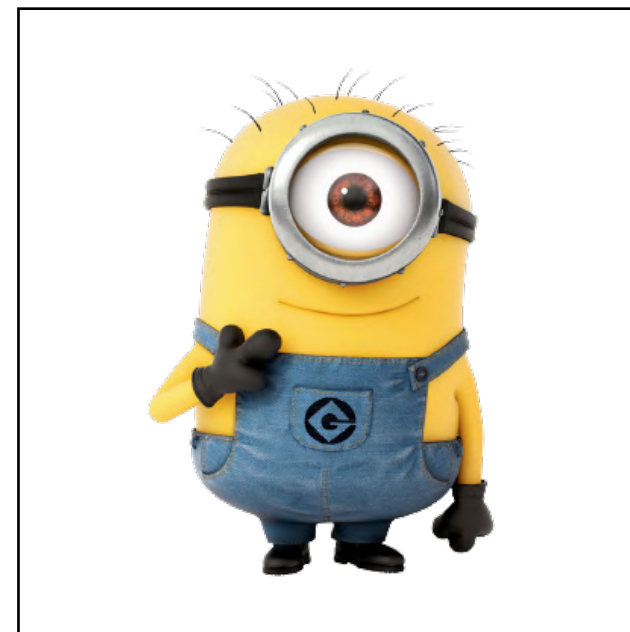
Input View

Method: 3D Reconstruction

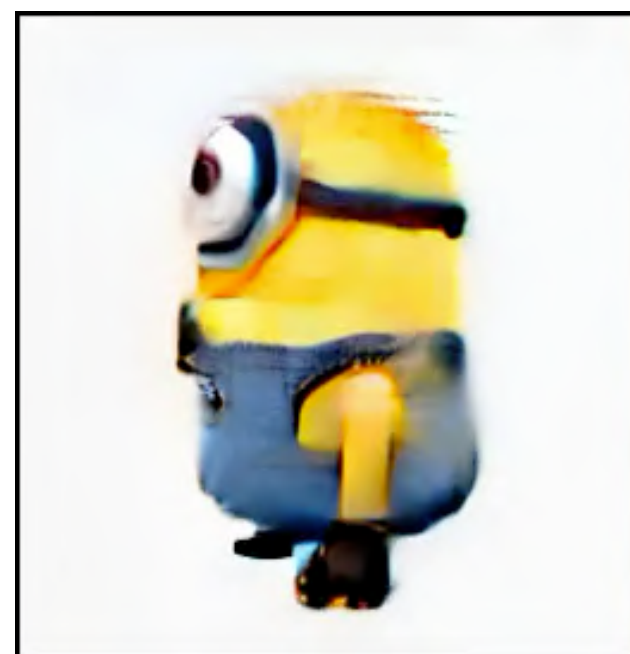
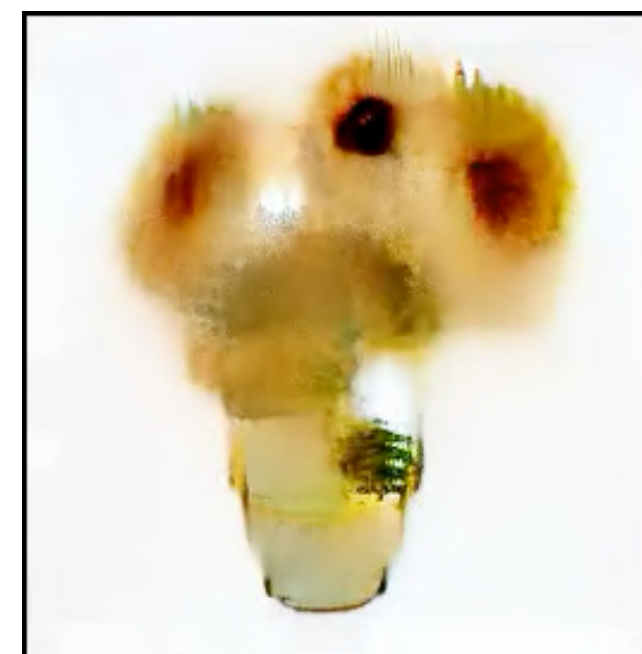


Results: 3D Reconstruction

Input Image



Reconstruction



Voxel NeRF + SJC



Results: 3D Reconstruction

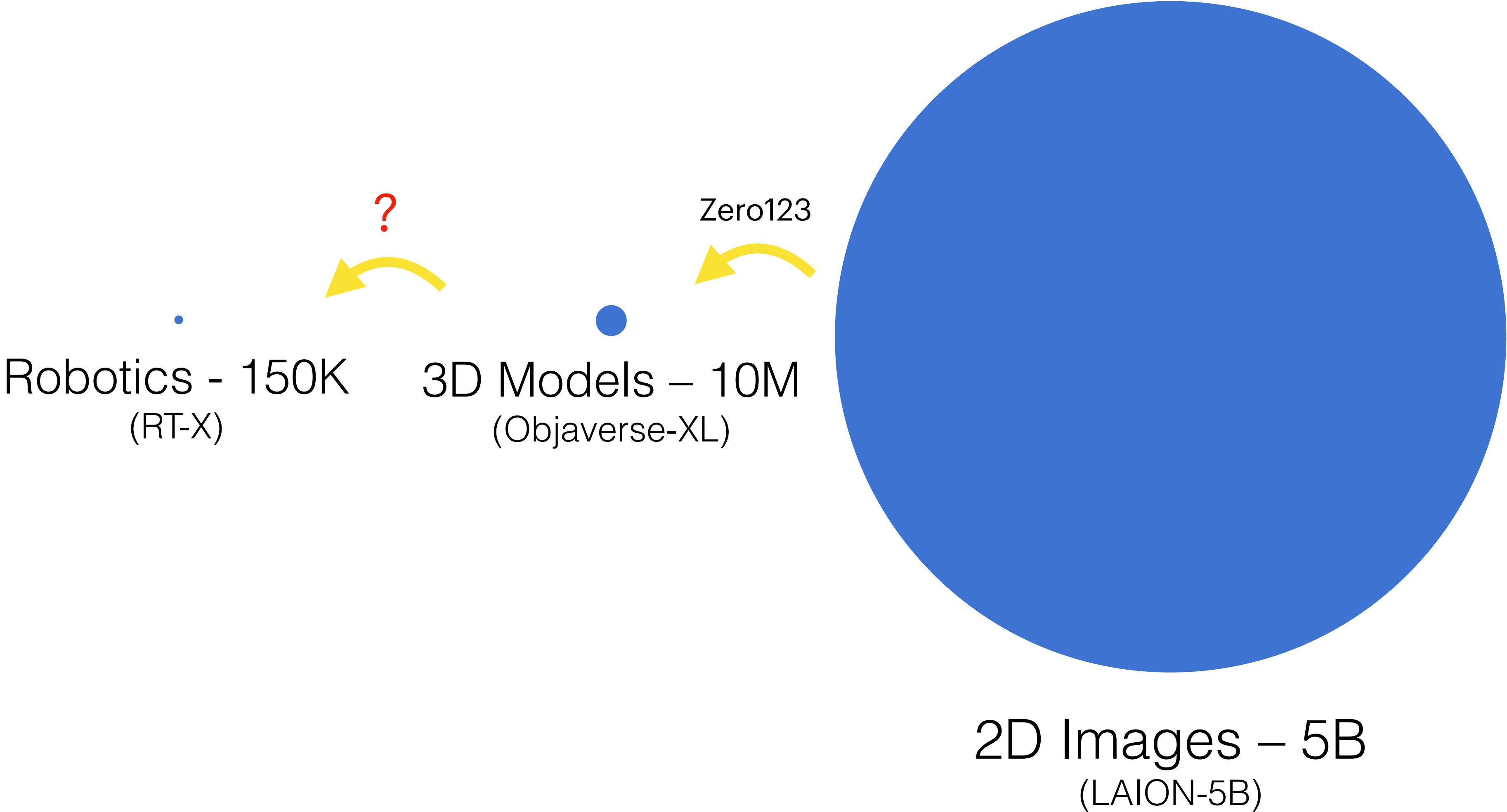


Zero123 + Instant-NGP + SDS + DMTet

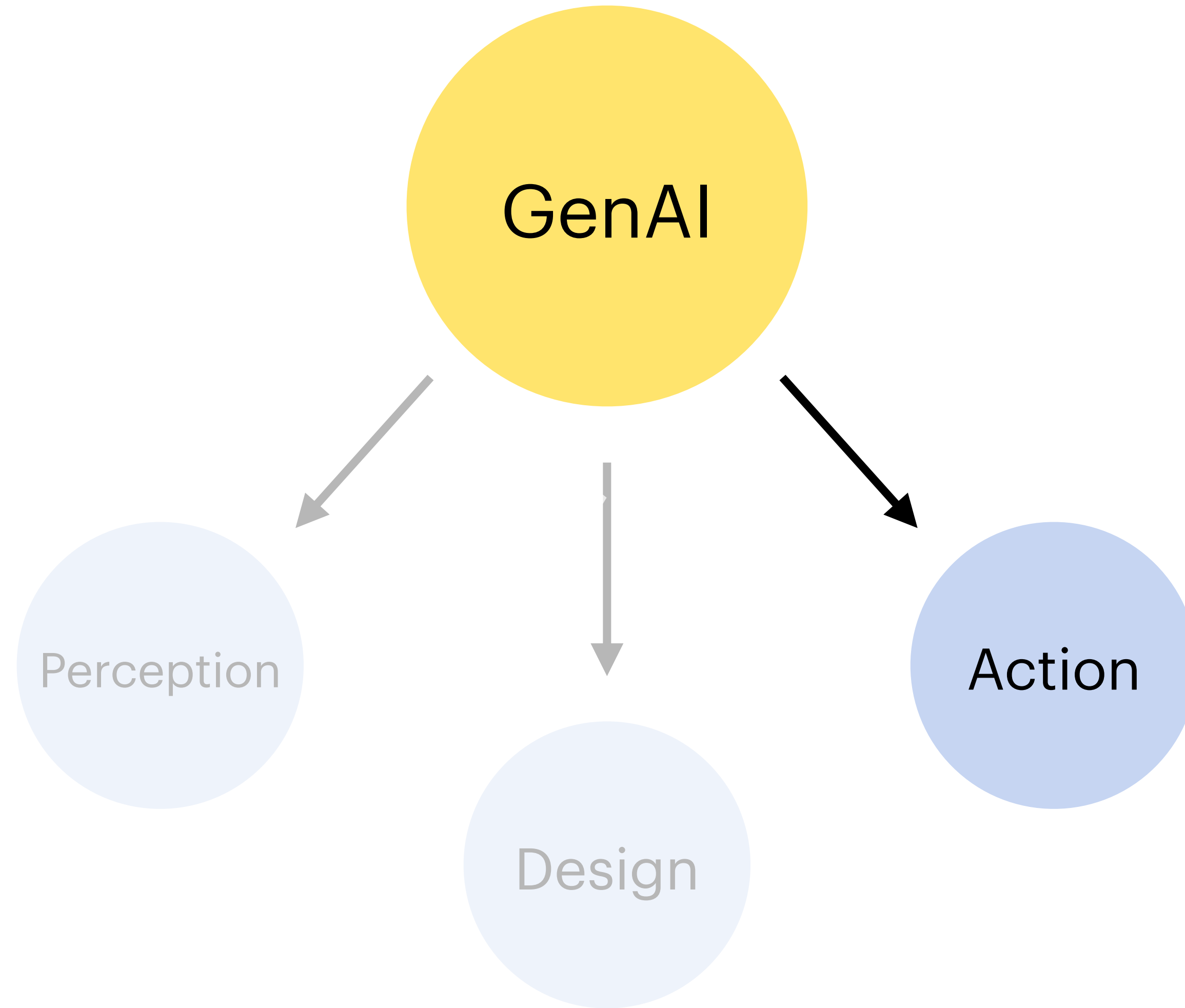
From Stability AI: <https://github.com/threestudio-project/threestudio#zero-1-to-3->



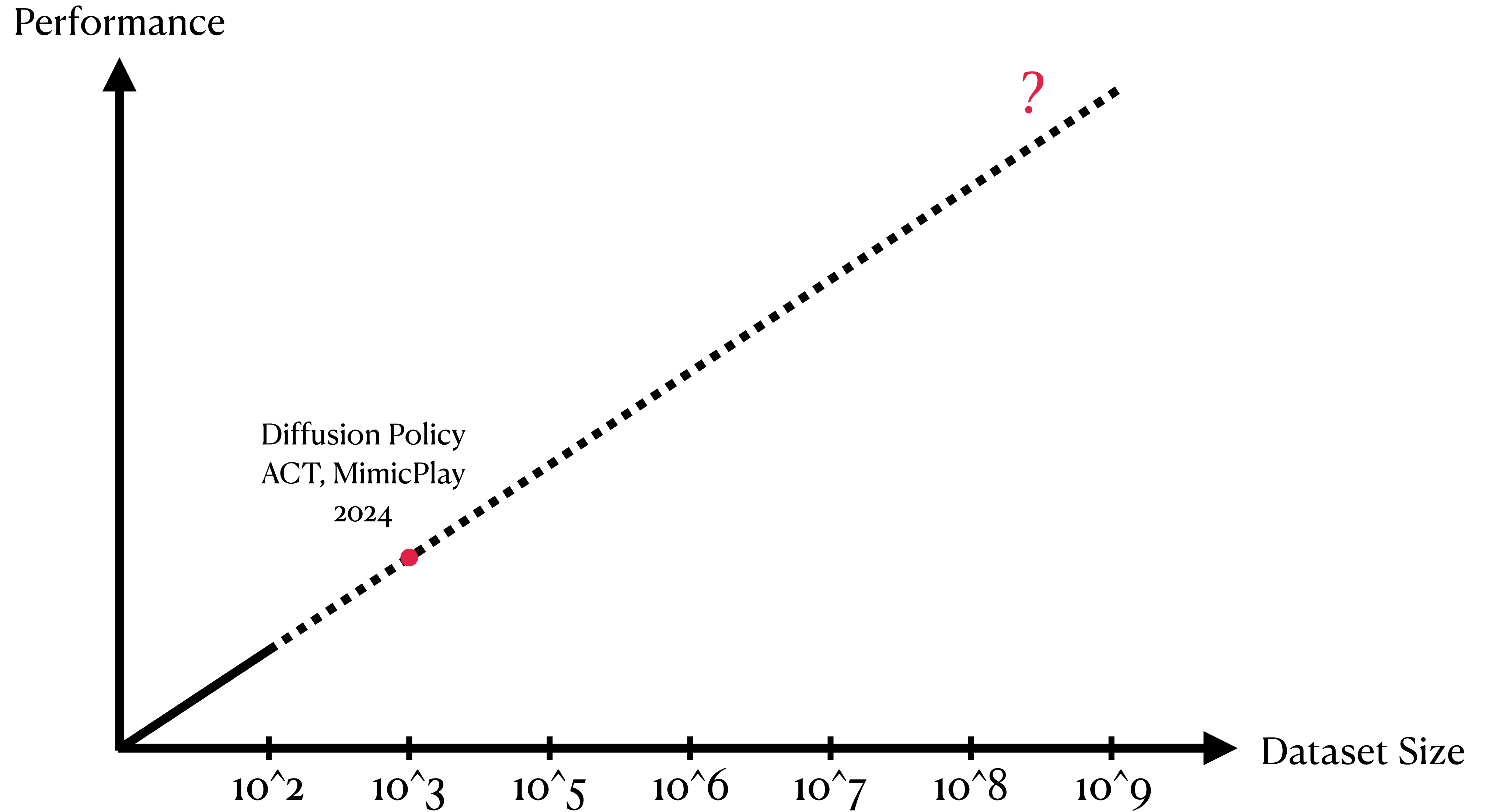
What about robotics?



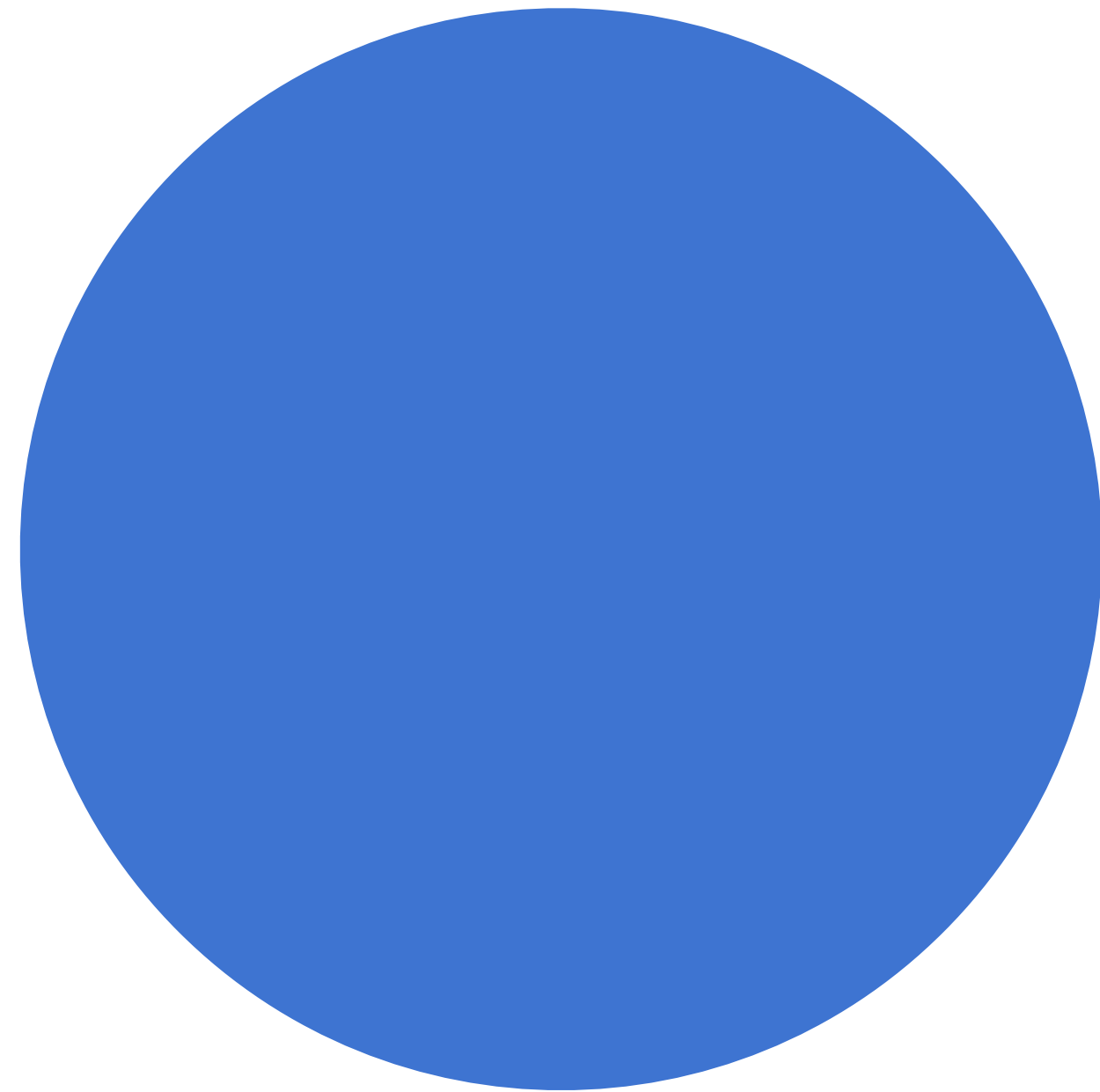
Generative Embodied AI



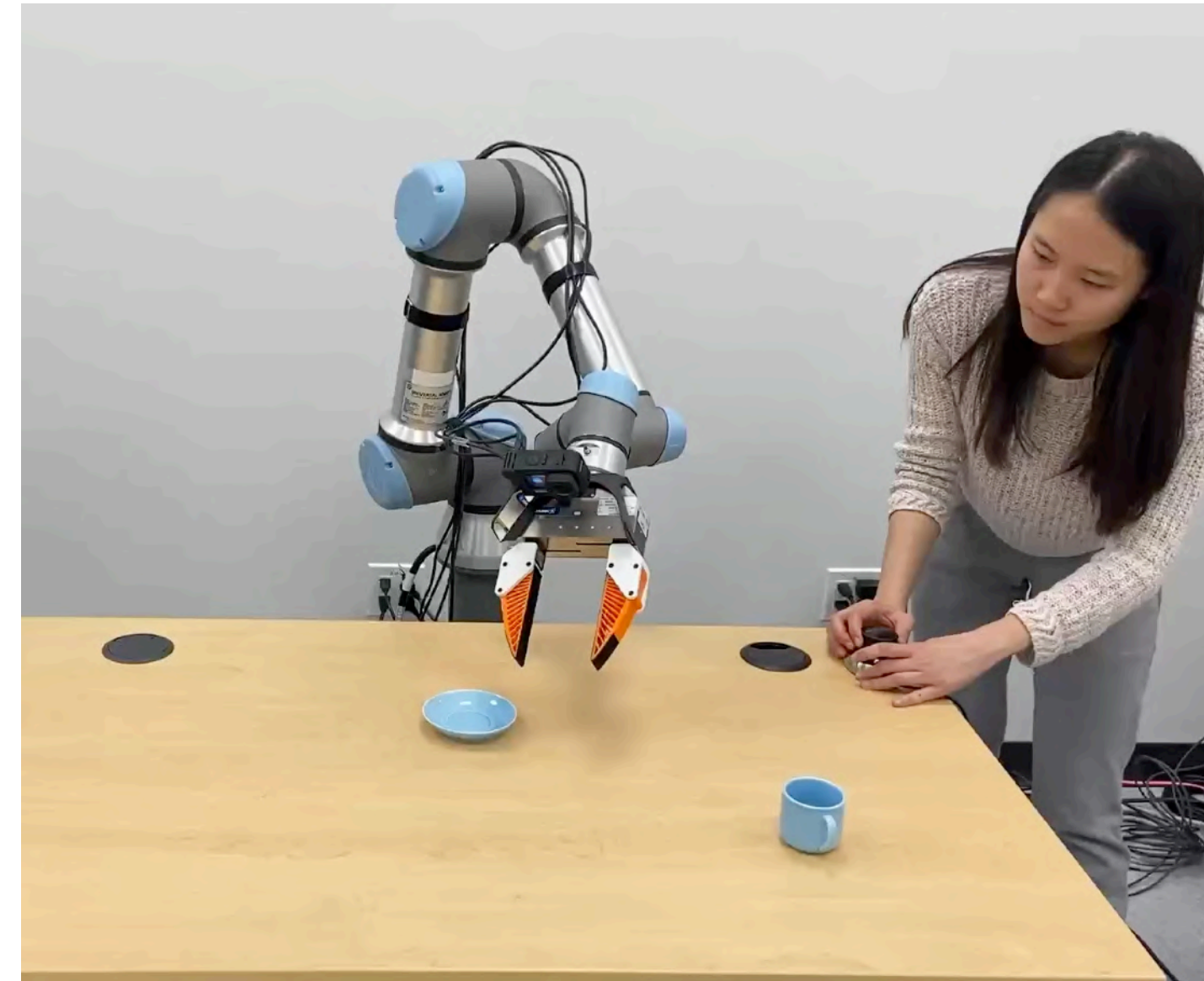
Behavior Cloning Puts Robot Learning on a Scaling Curve



Robot Learning Has a Data Problem



Robot Data



The Dilemma of Robot vs. Visual Data

Limited
but Robot Complete



Robot Data

Diverse
but Embodiment Gap

Internet Image / Video



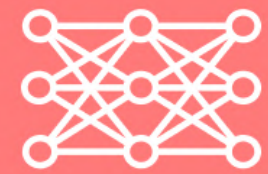
We need an interface between visual data and robot control!

Robot Complete
But Not Diverse



Robot Data

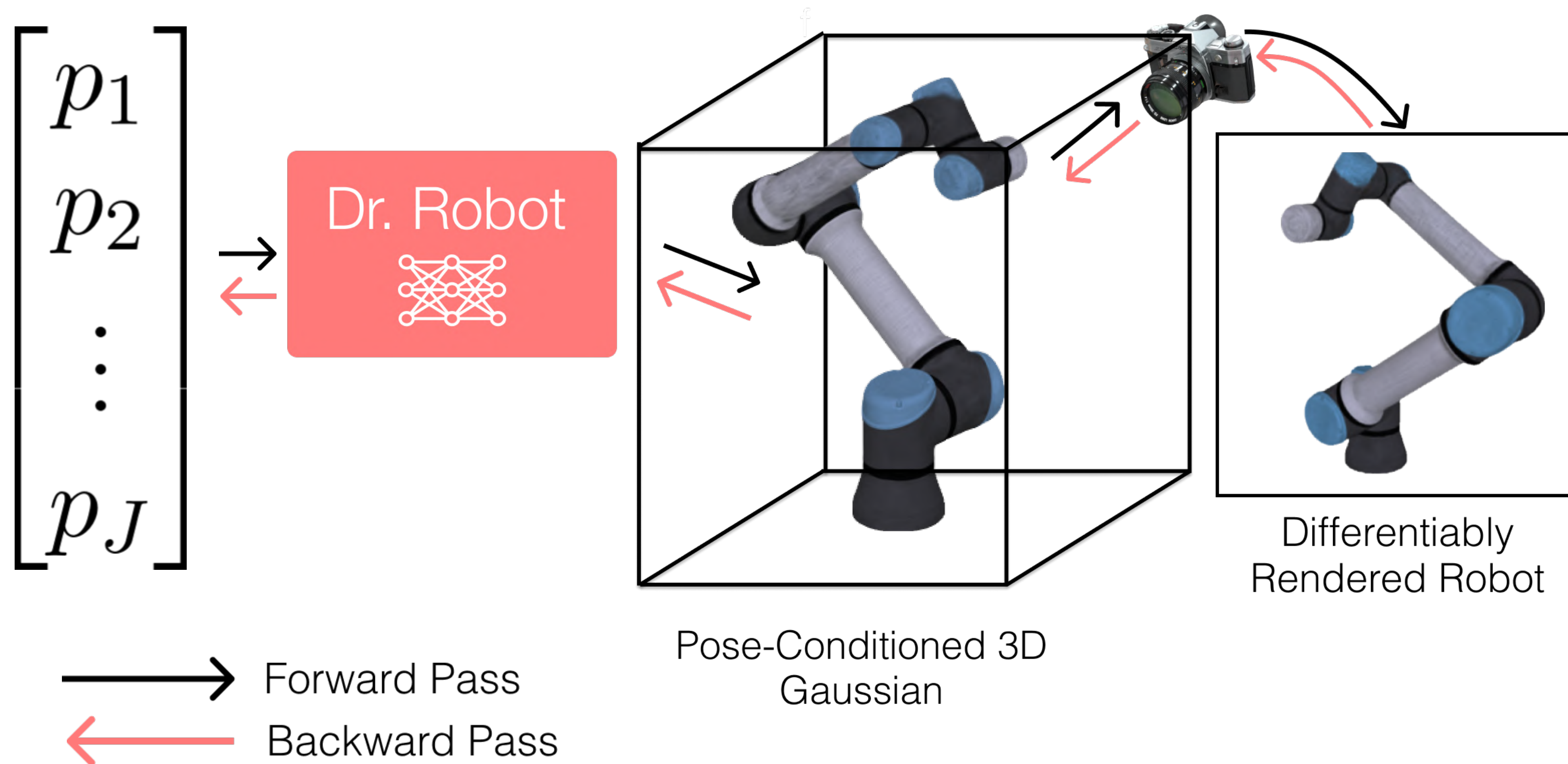
Dr. Robot



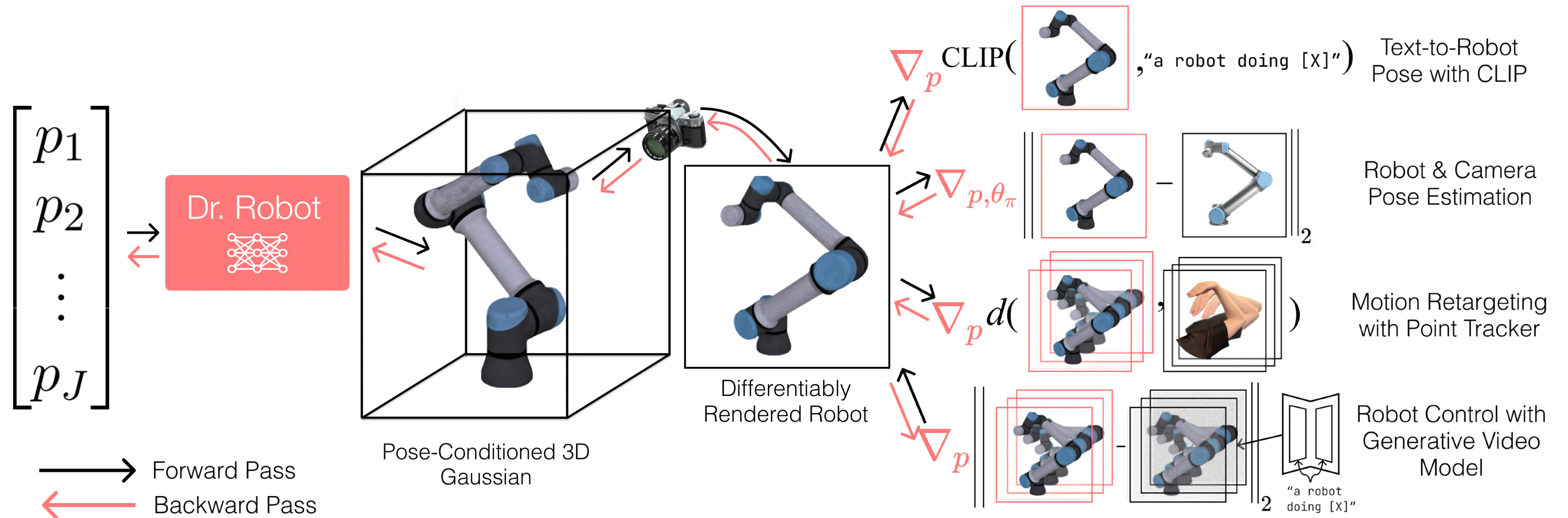
Diverse
But Embodiment Gap

Visual Data

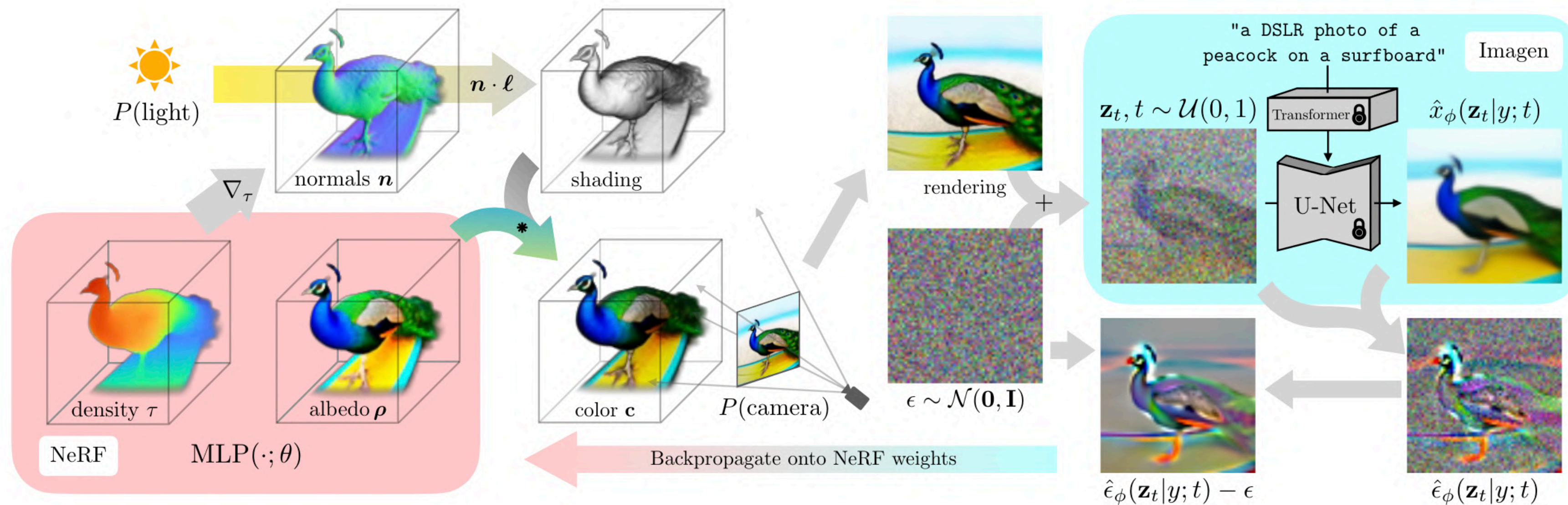
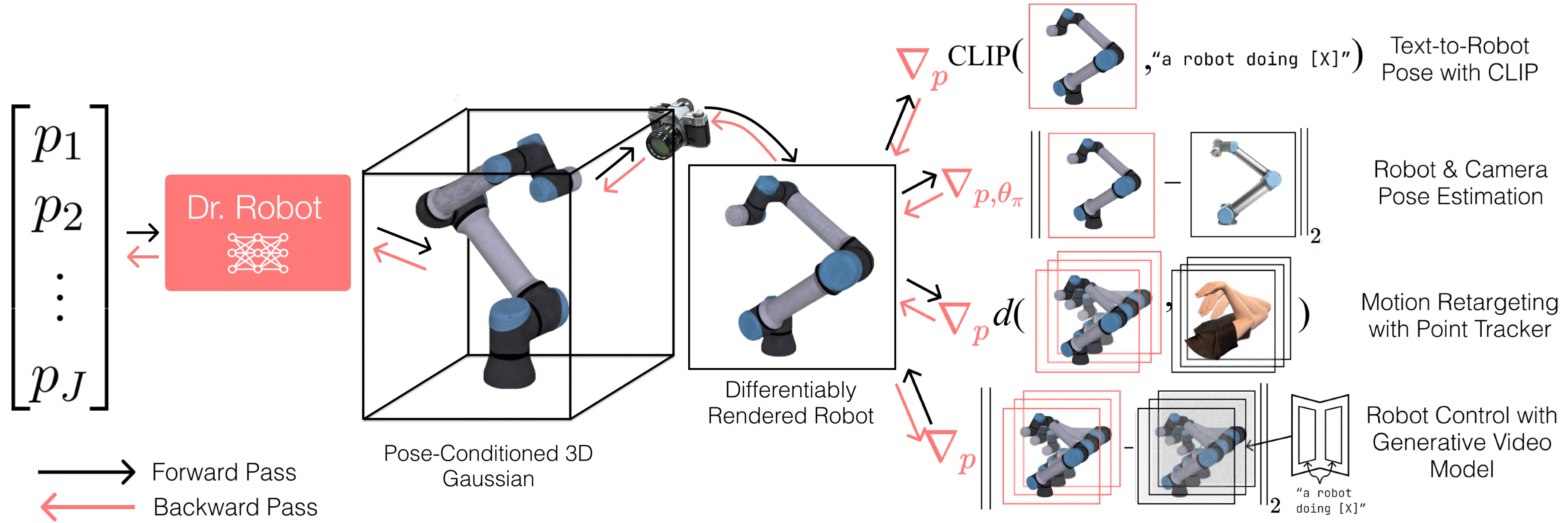
Differentiable Robot Rendering



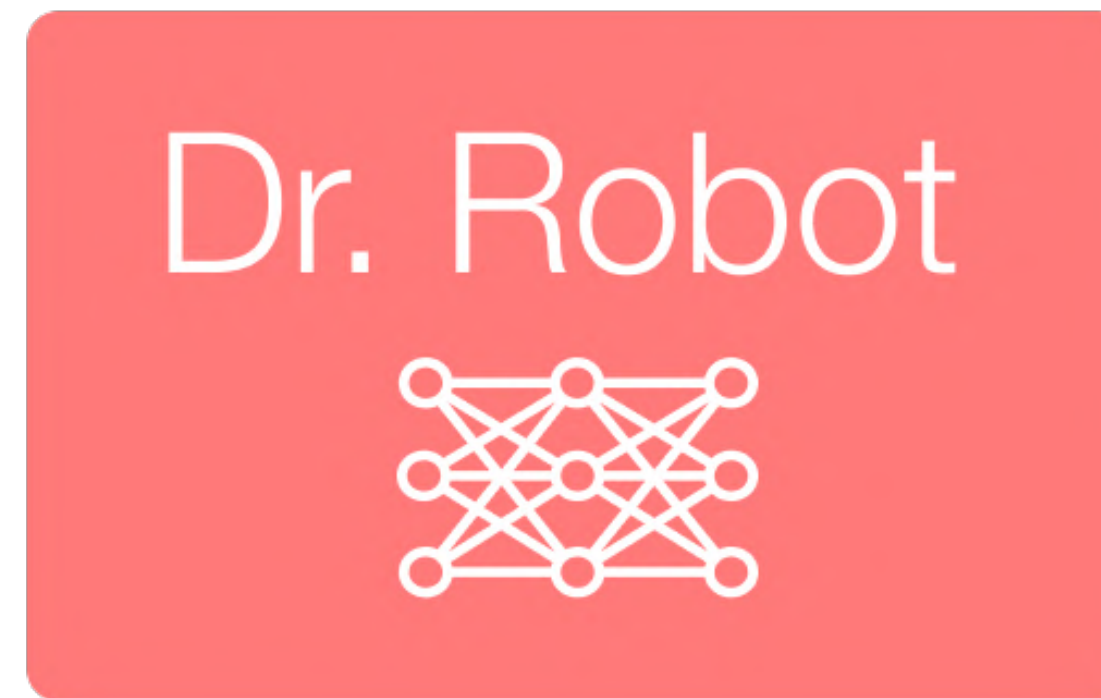
Dr. Robot Allows Visual Appearance of Robots to be Differentiable w.r.t. to Control Parameters



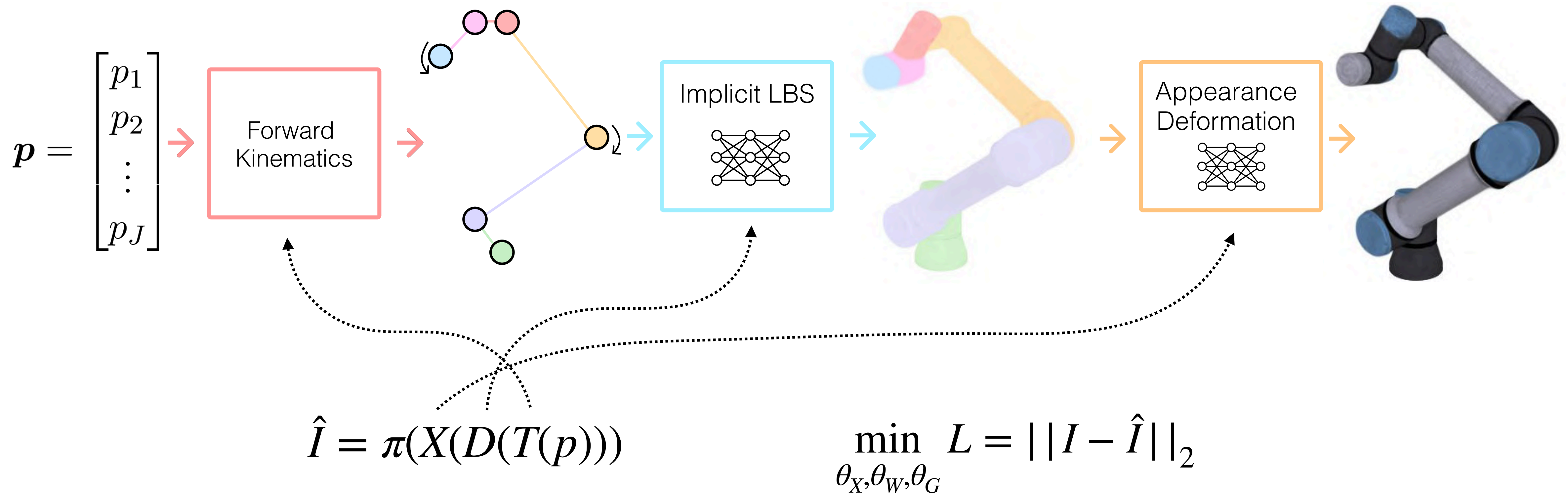
3D Gen, but for robot



How did we achieve differentiability?



Differentiable Robot Rendering



Built from any URDF



Rendering



GT



Rendering



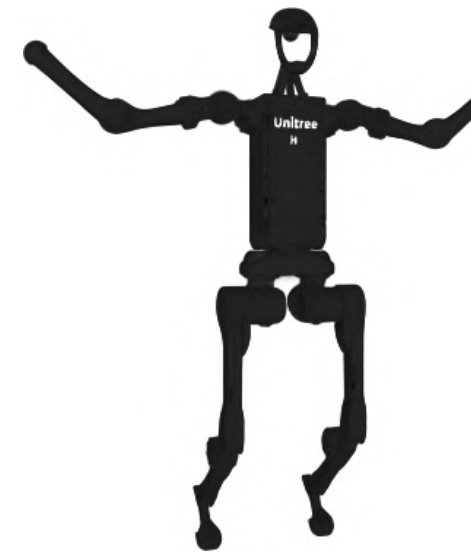
GT



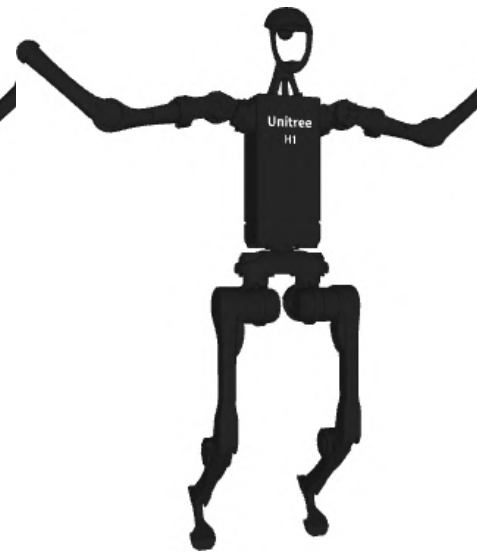
Rendering



GT



Rendering



GT



Rendering



GT



Rendering



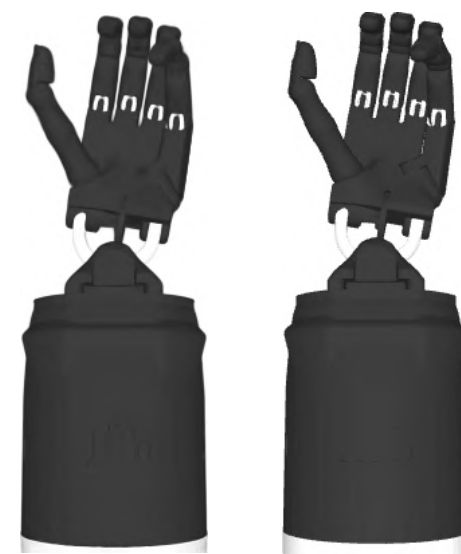
GT



Rendering



GT



Rendering



GT



Rendering



GT



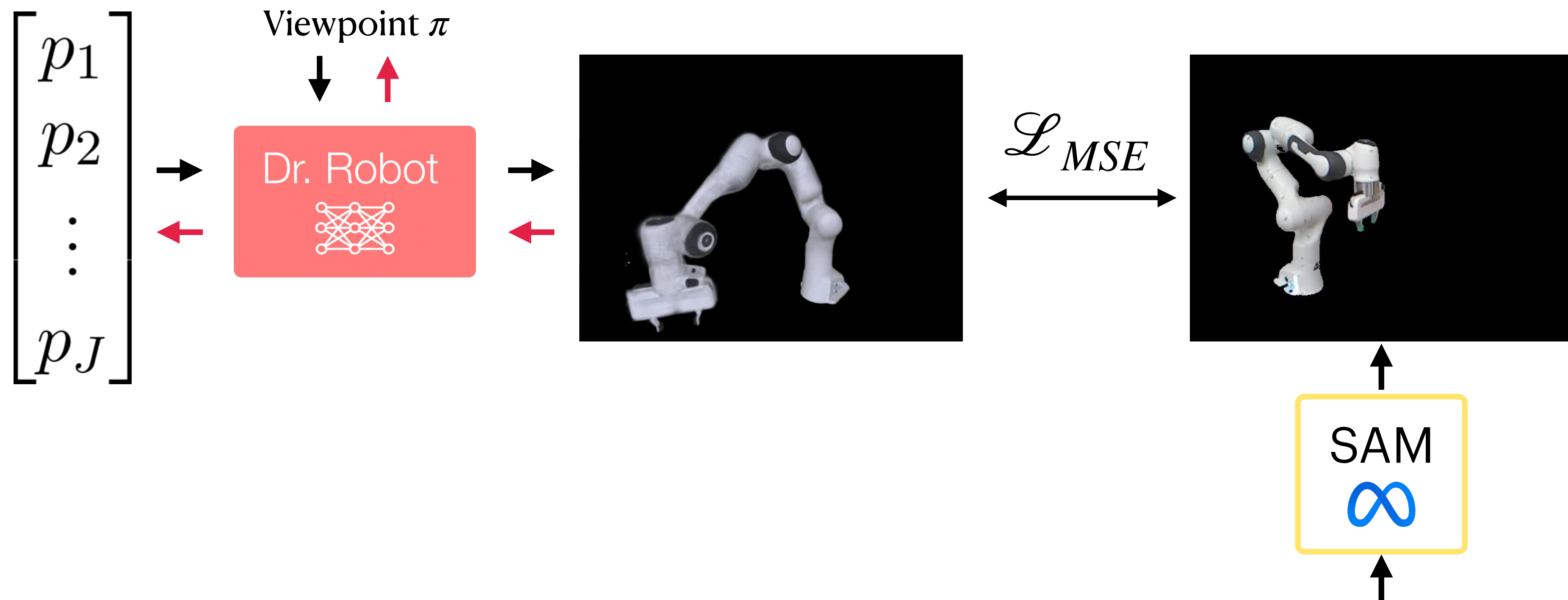
Rendering



GT



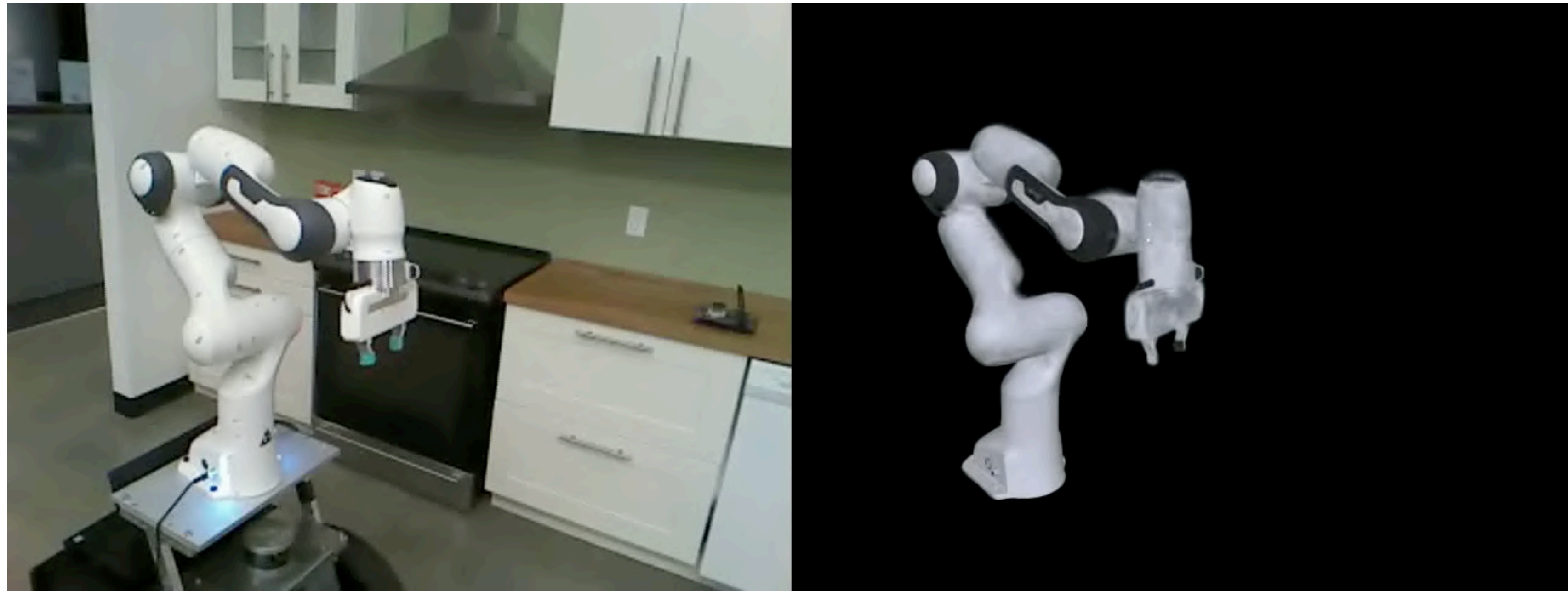
Robot Poses Reconstruction from Single Image Through Analysis-by-Synthesis



Optimization Objective: $\min_{p, \theta_\pi} \mathcal{L}_{MSE}(p, I) = ||I - \pi(f(p))||_2$



Robot Poses Reconstruction from Single Image Through Analysis-by-Synthesis

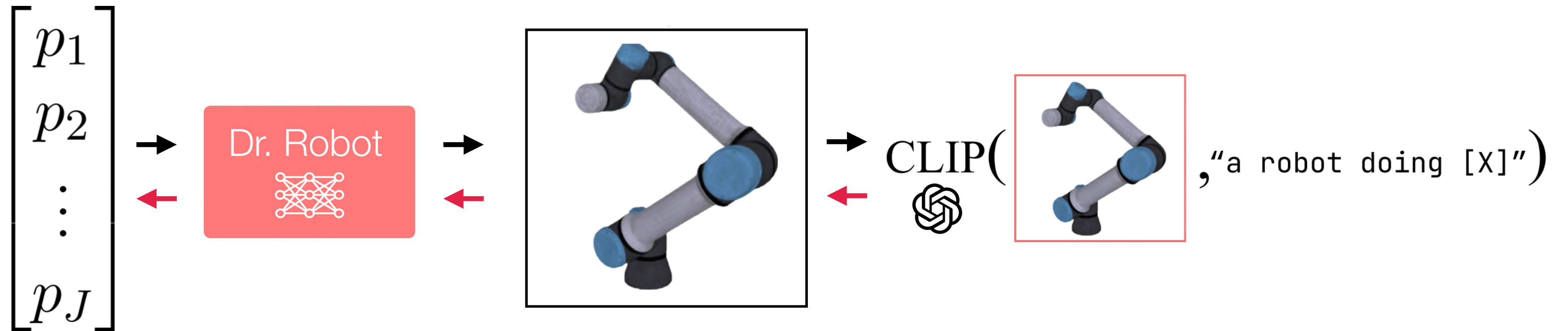


Original Video

Reconstructed Robot



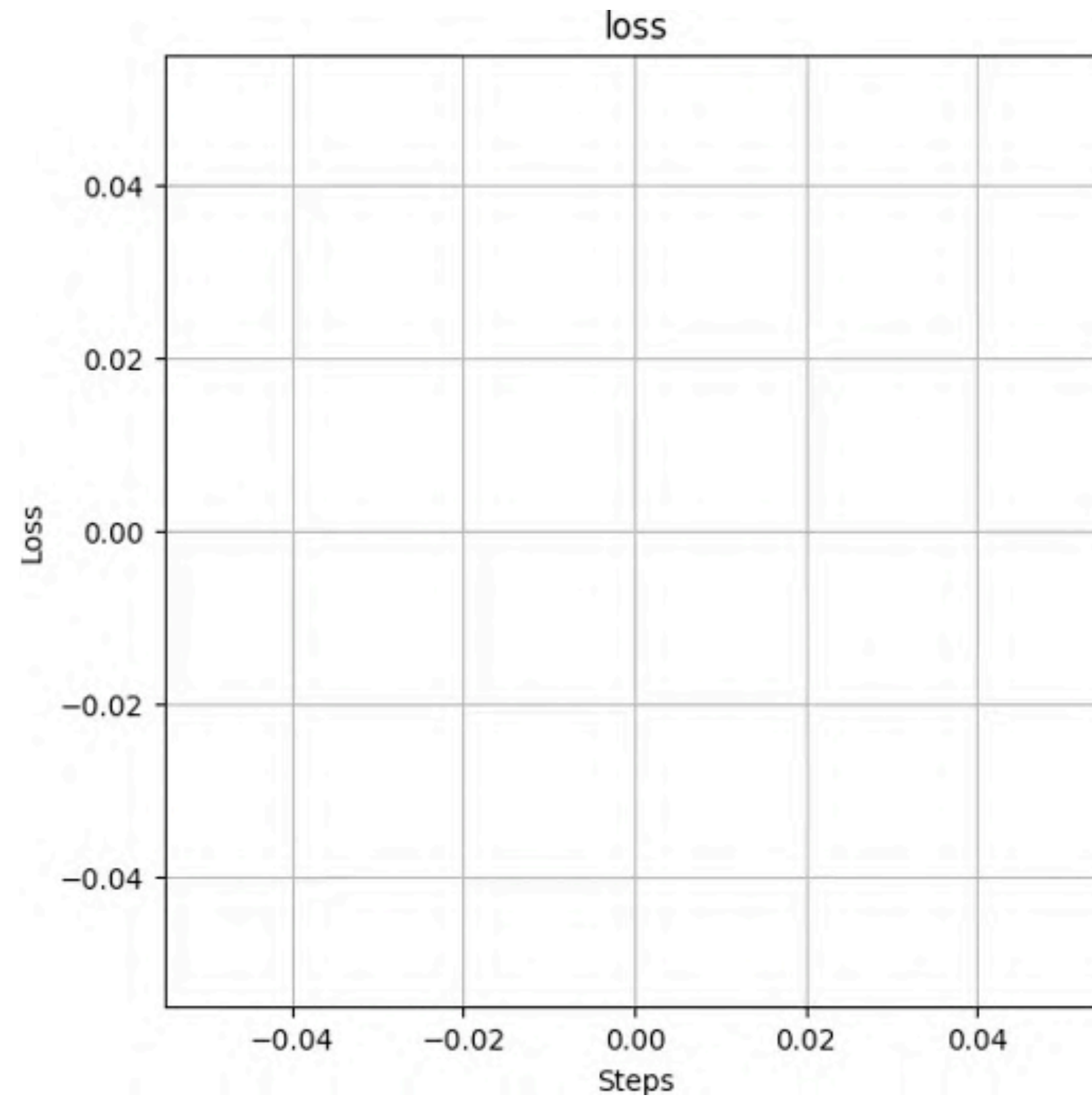
Visual MPC with CLIP



Optimization Objective: $\min_p \mathcal{L}_{\text{CLIP}}(p; \text{prompt}) = \text{CLIP}(\text{prompt}, \pi(f(p)))$



Visual MPC with CLIP

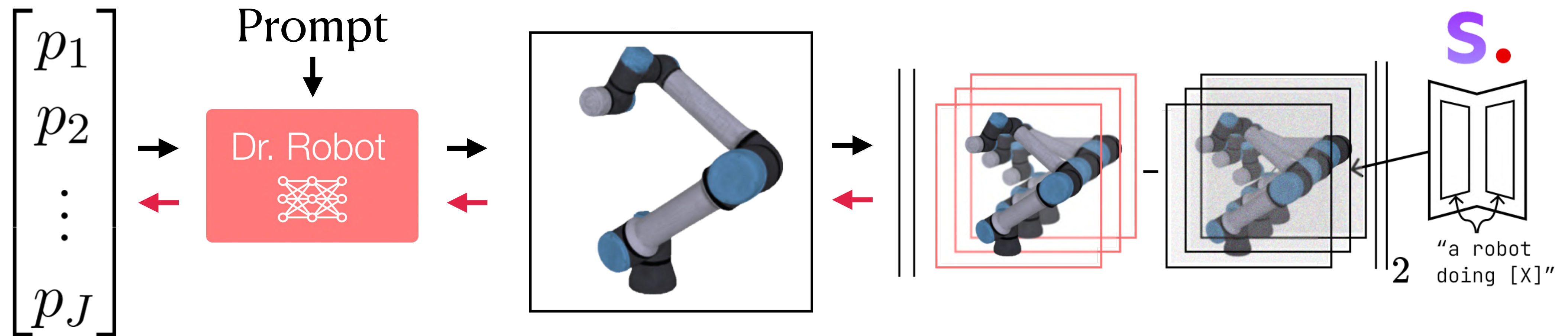


$\mathcal{L}_{\text{CLIP}}(p; \text{prompt})$

“Victory Sign”



Robot Control with Text2Video Model



Optimization Objective:
$$\min_{p, \theta_\pi} \sum_i^T \mathcal{L}_{MSE}(p_i, I_i) = ||I_i - \pi(f(p_i))||_2$$

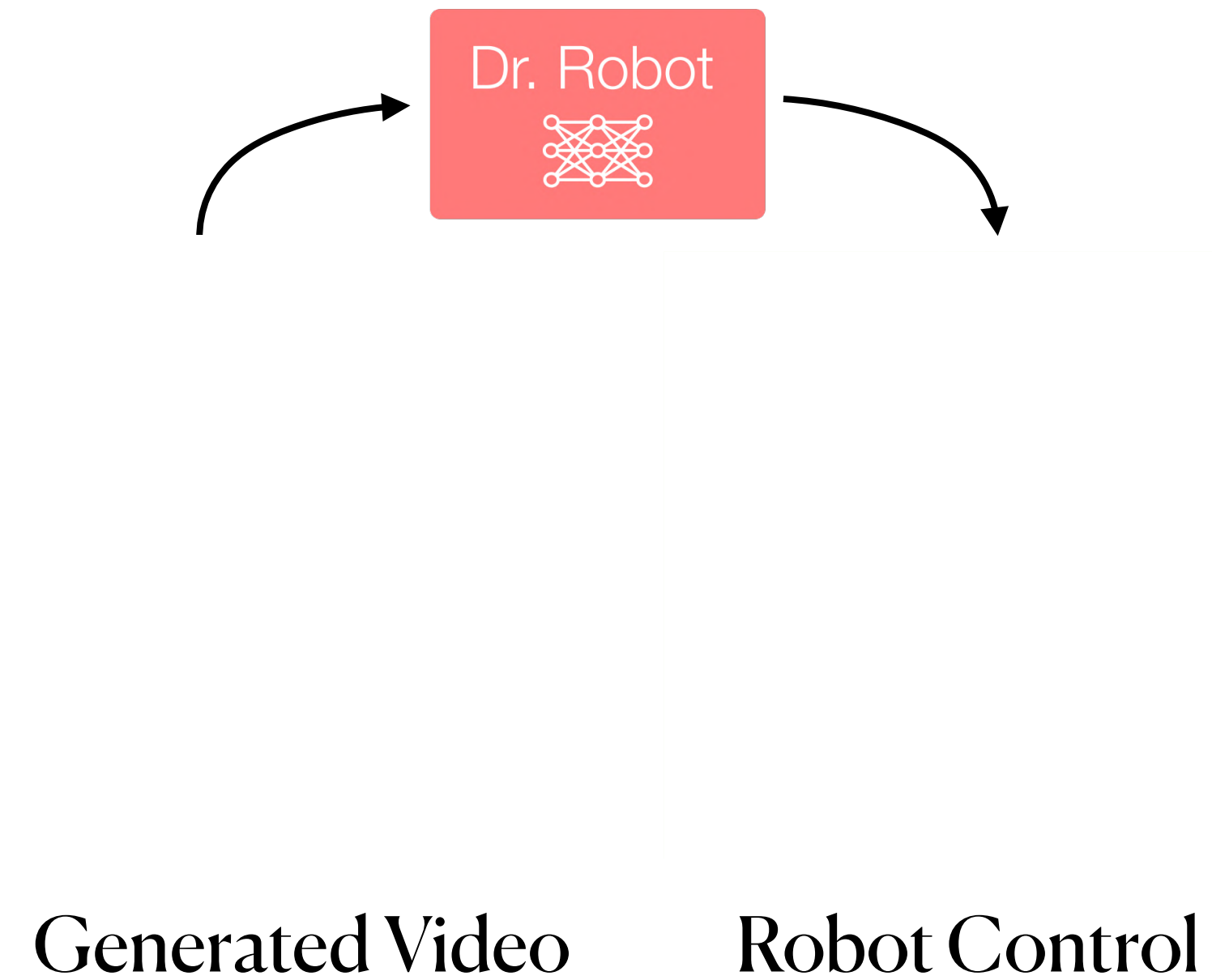
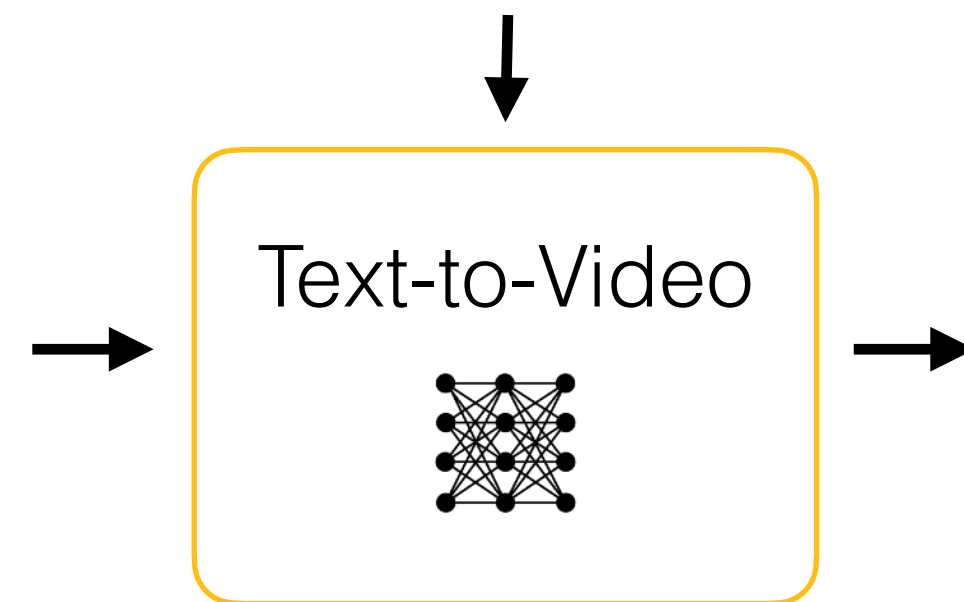


Robot Control with Text2Video Model



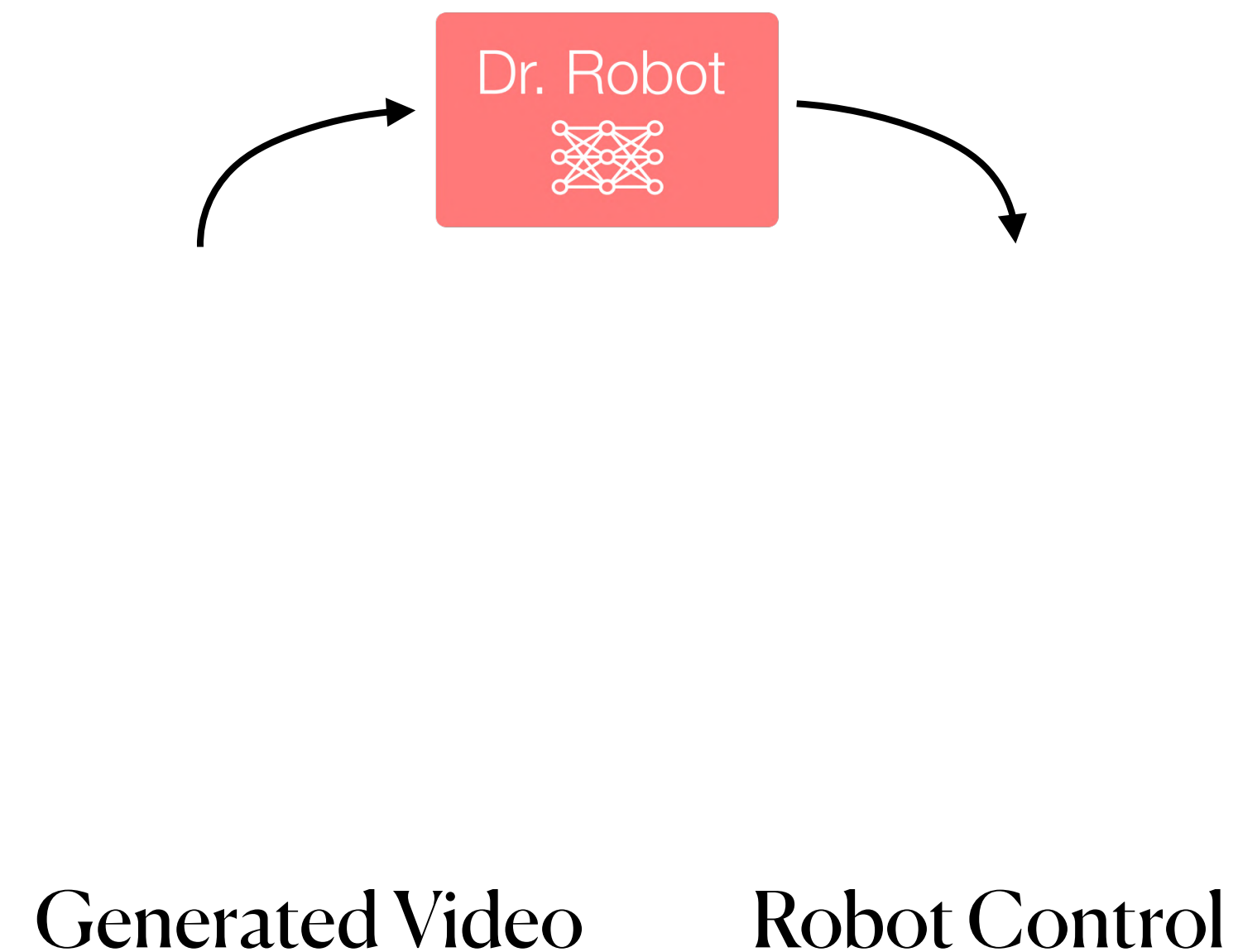
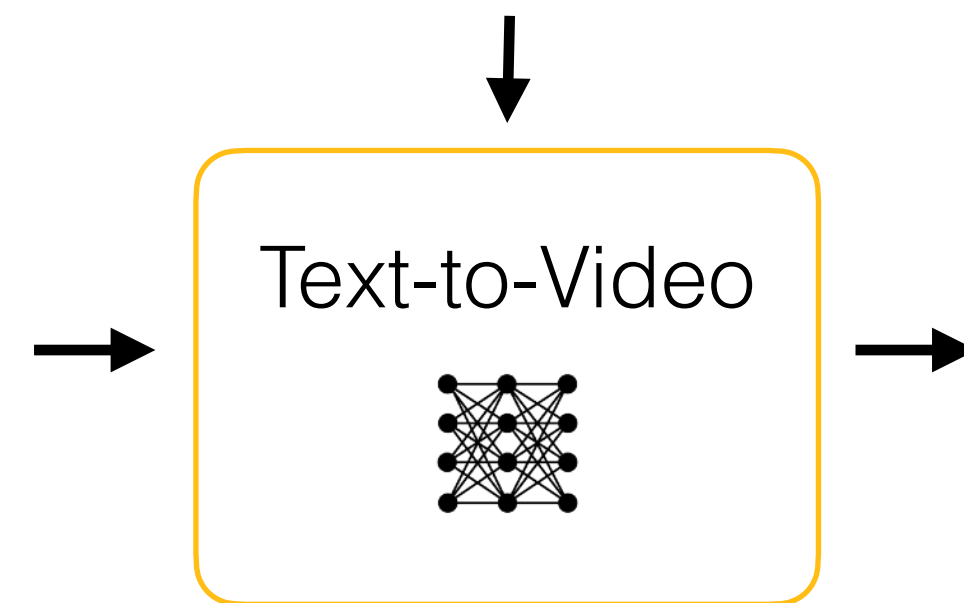
Current Observation

“Pick up the blue can.”

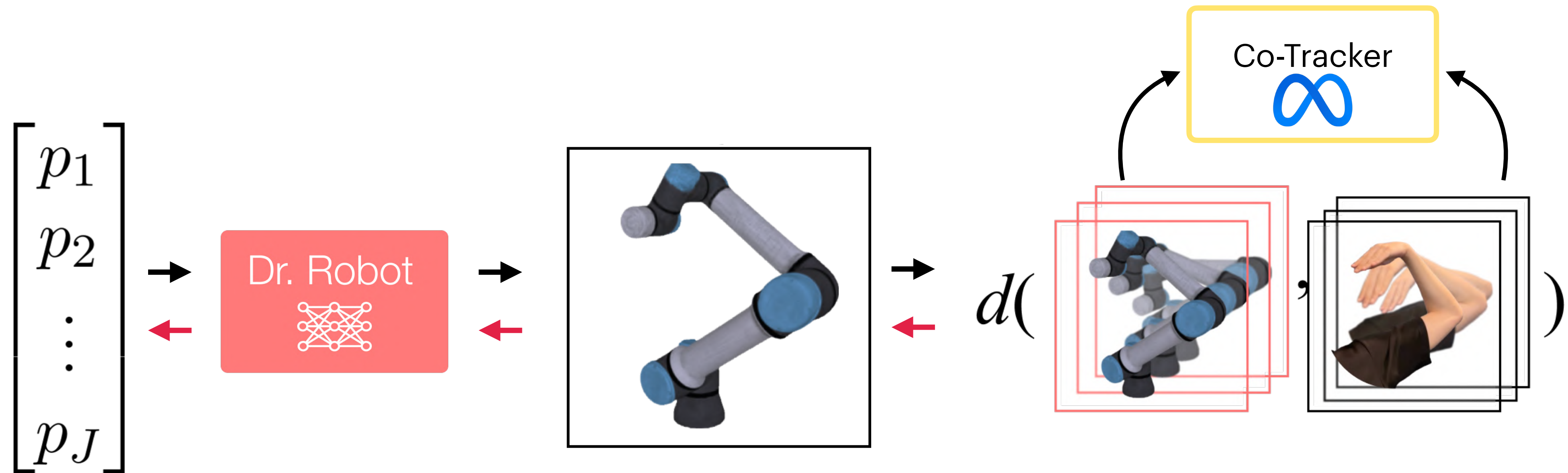


Current Observation

“Pick up the red block.”



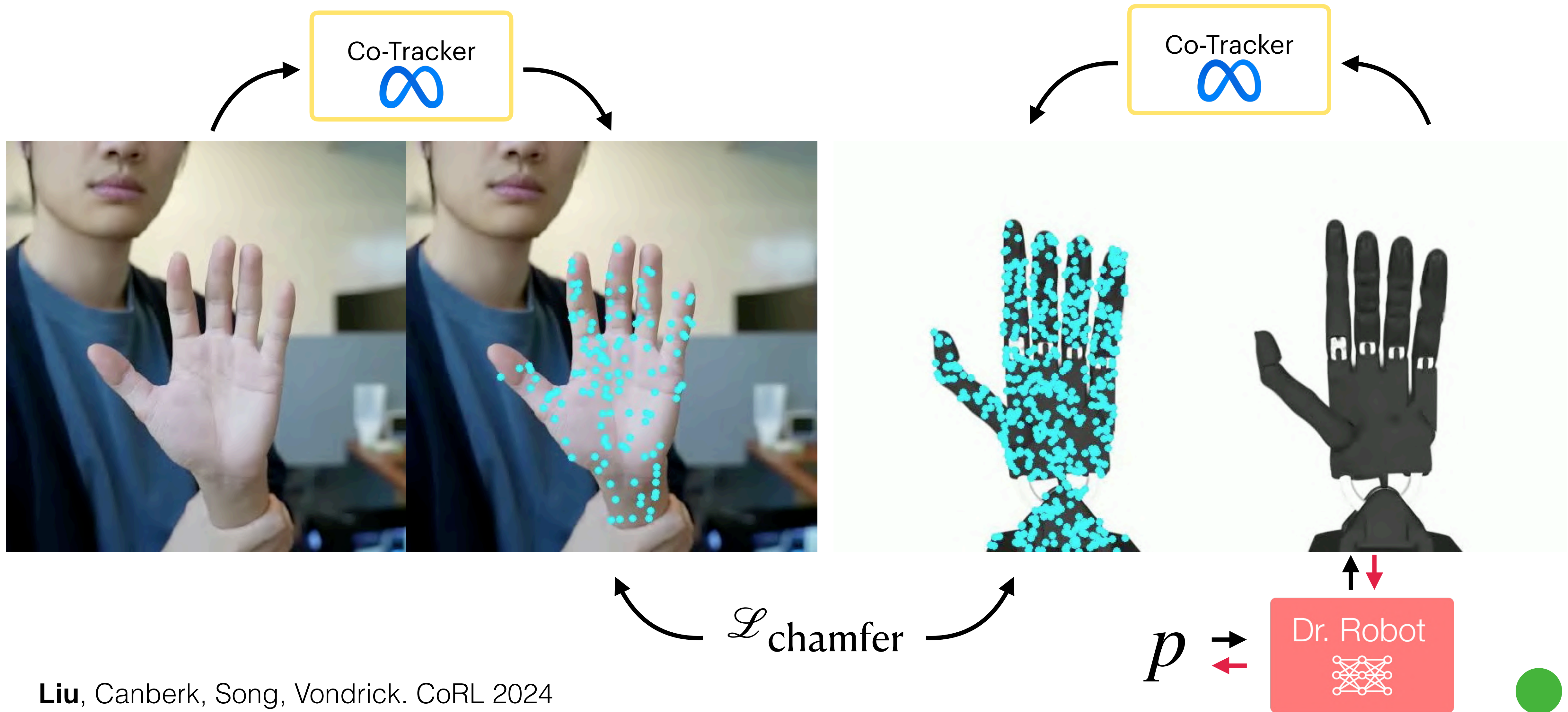
Motion Retargetting with Point Tracker



Optimization Objective:
$$\mathcal{L}_{\text{Track}}(p_{1:T}; c_0^r, c_{1:T}) = \sum_{t=1}^T d(c_t, c_t^r) = \sum_{t=1}^T d(c_t, \mathcal{P}(c_0^r, \pi(p_{1:T}))_t)$$



Motion Retargetting with Point Tracker

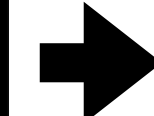
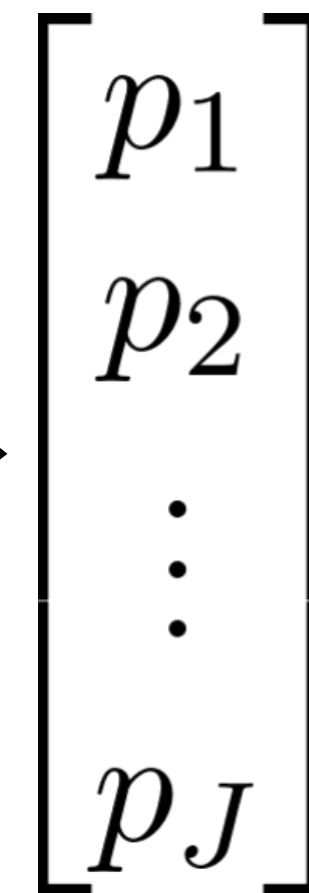
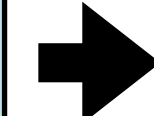
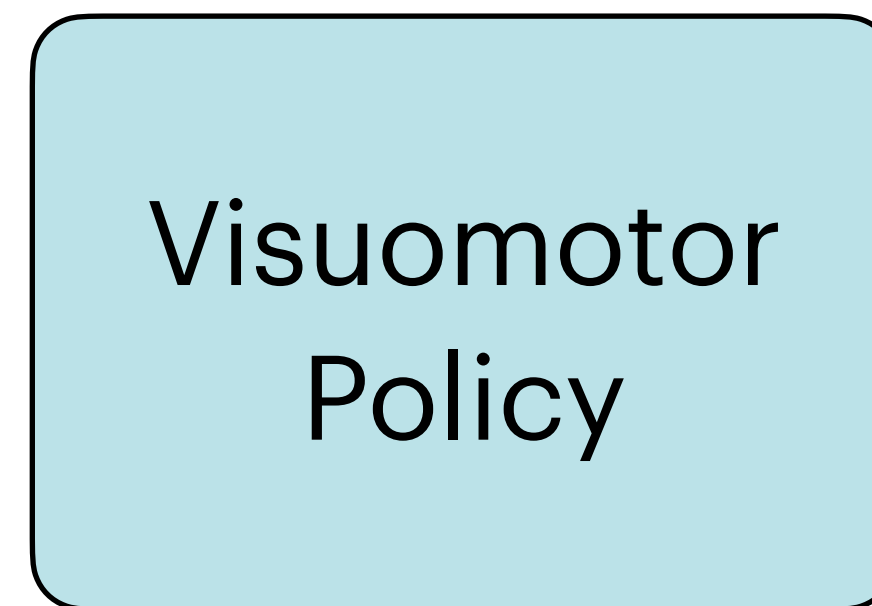
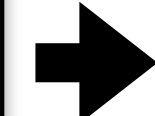


What about policy learning?

Visuomotor Policy

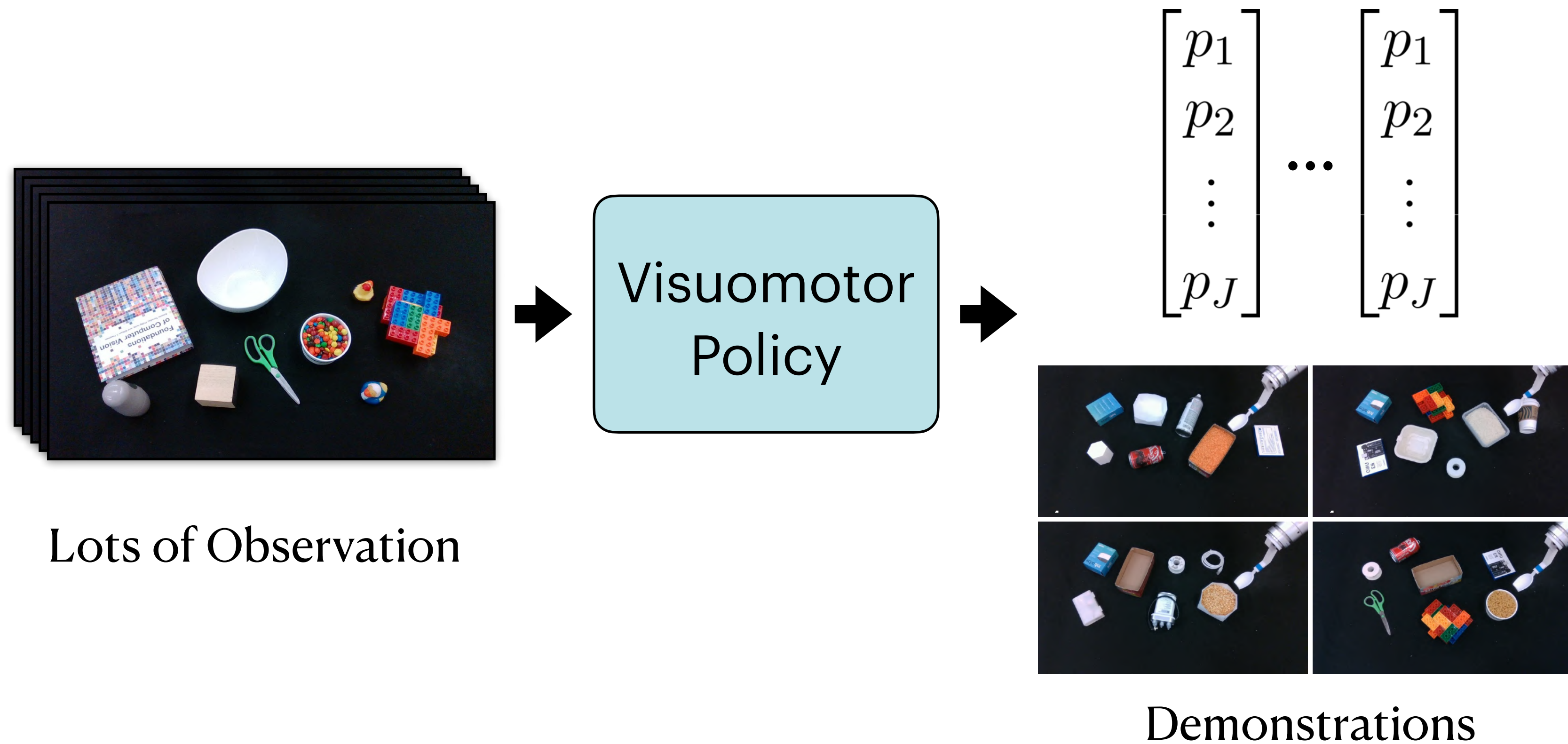


Observation

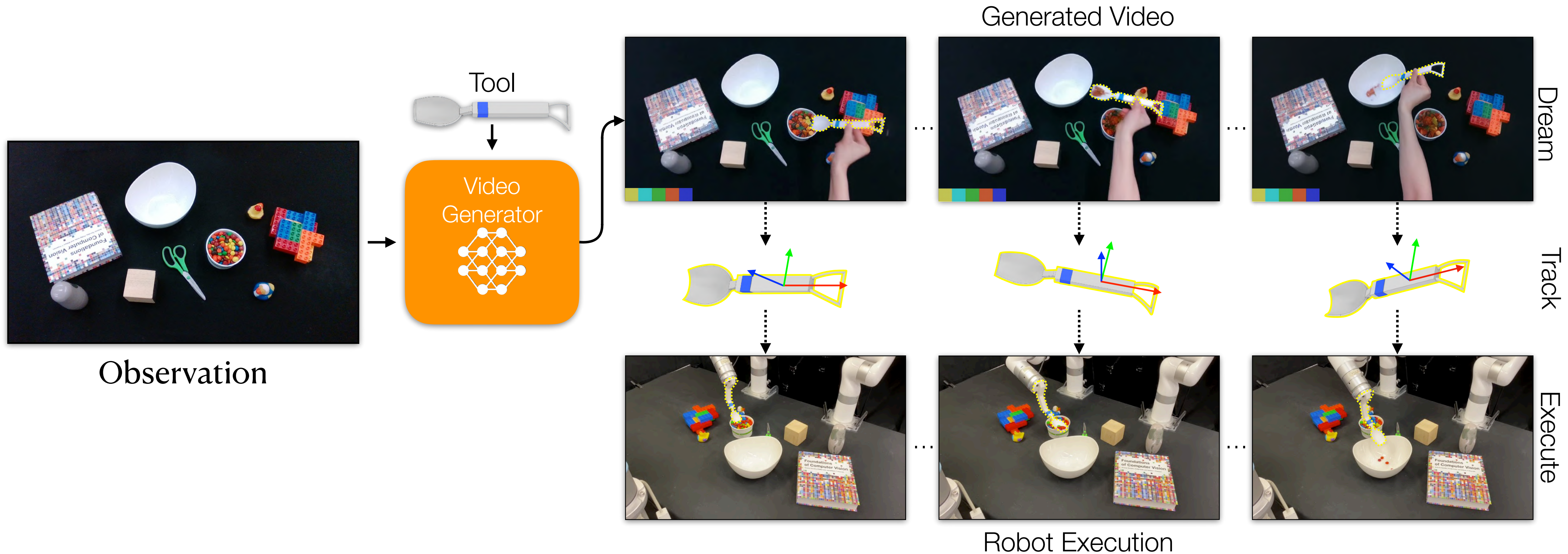


Motor Actions

Behavior Cloning is Supervised Learning of Human Behavior




Video Generation as Intermediate Action Representation



Dreamitate

Visuomotor
Policy

$$= \text{Dream} + \text{Imitate}$$

Generated Video (indicated by )

Robot Execution

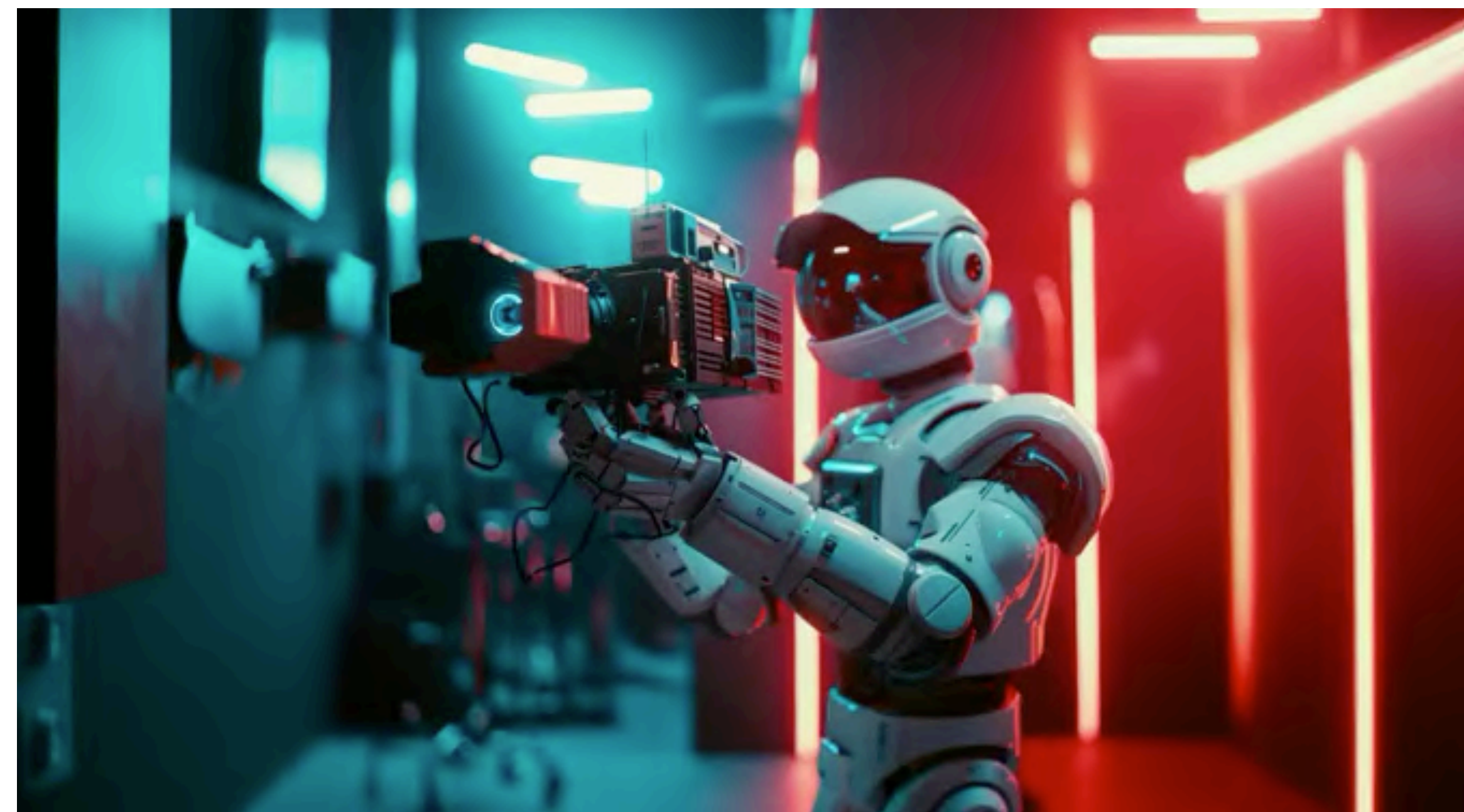
Input Image



Aligning Video Model to Robot

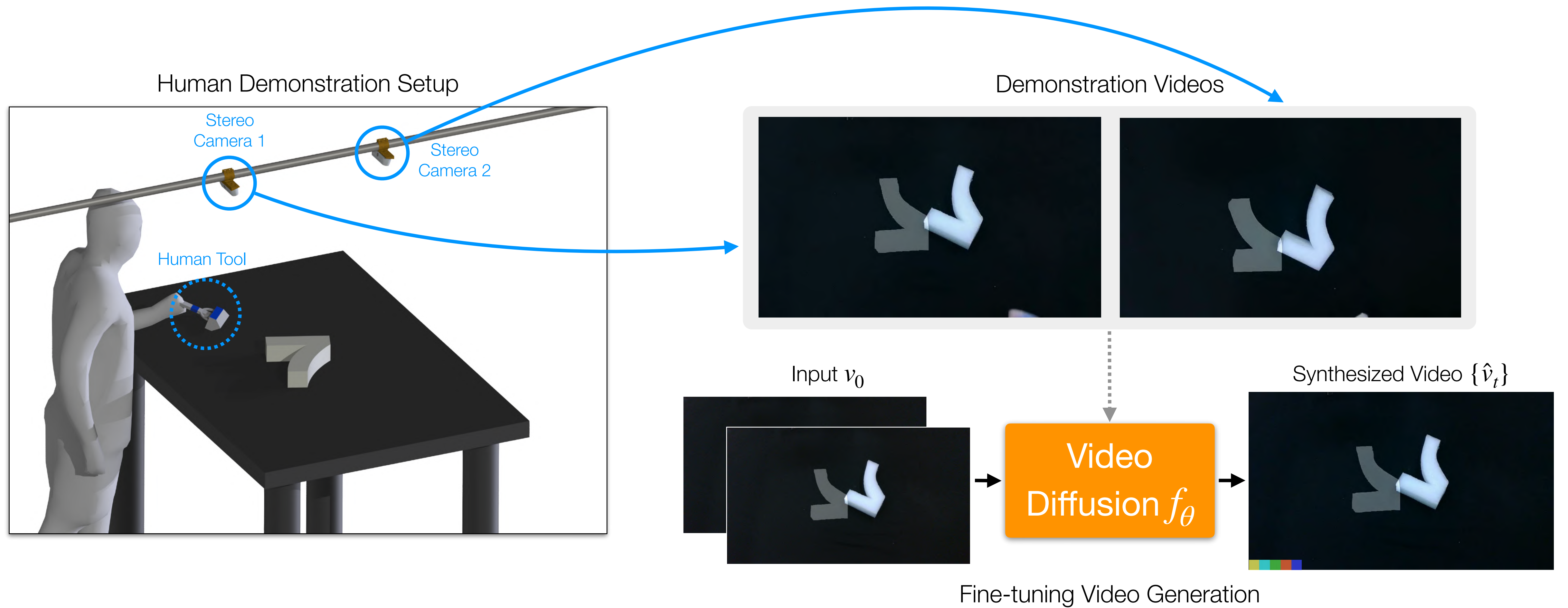


OpenAI SORA

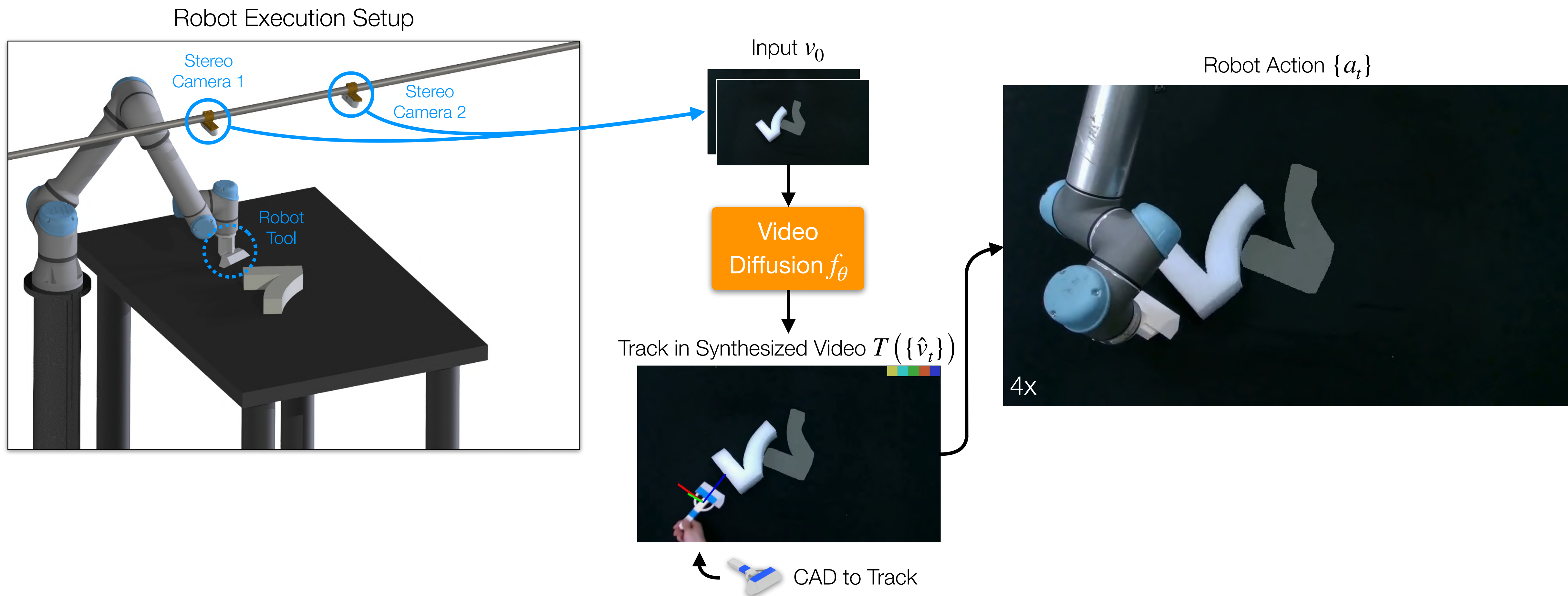


Luma Dream Machine

Data Collection



Robot Execution



Evaluation Tasks

Demonstration Videos



Rotation Task

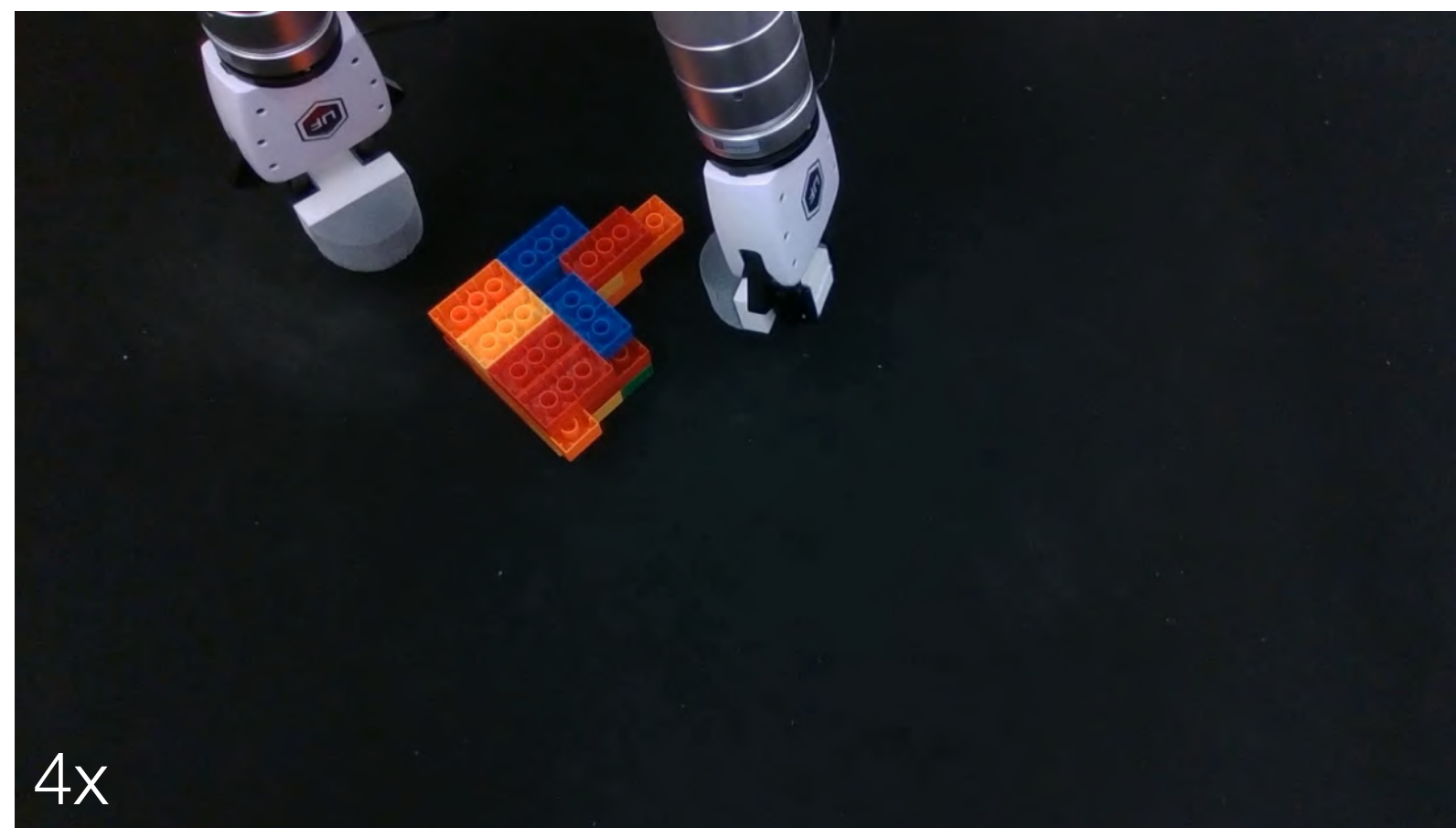
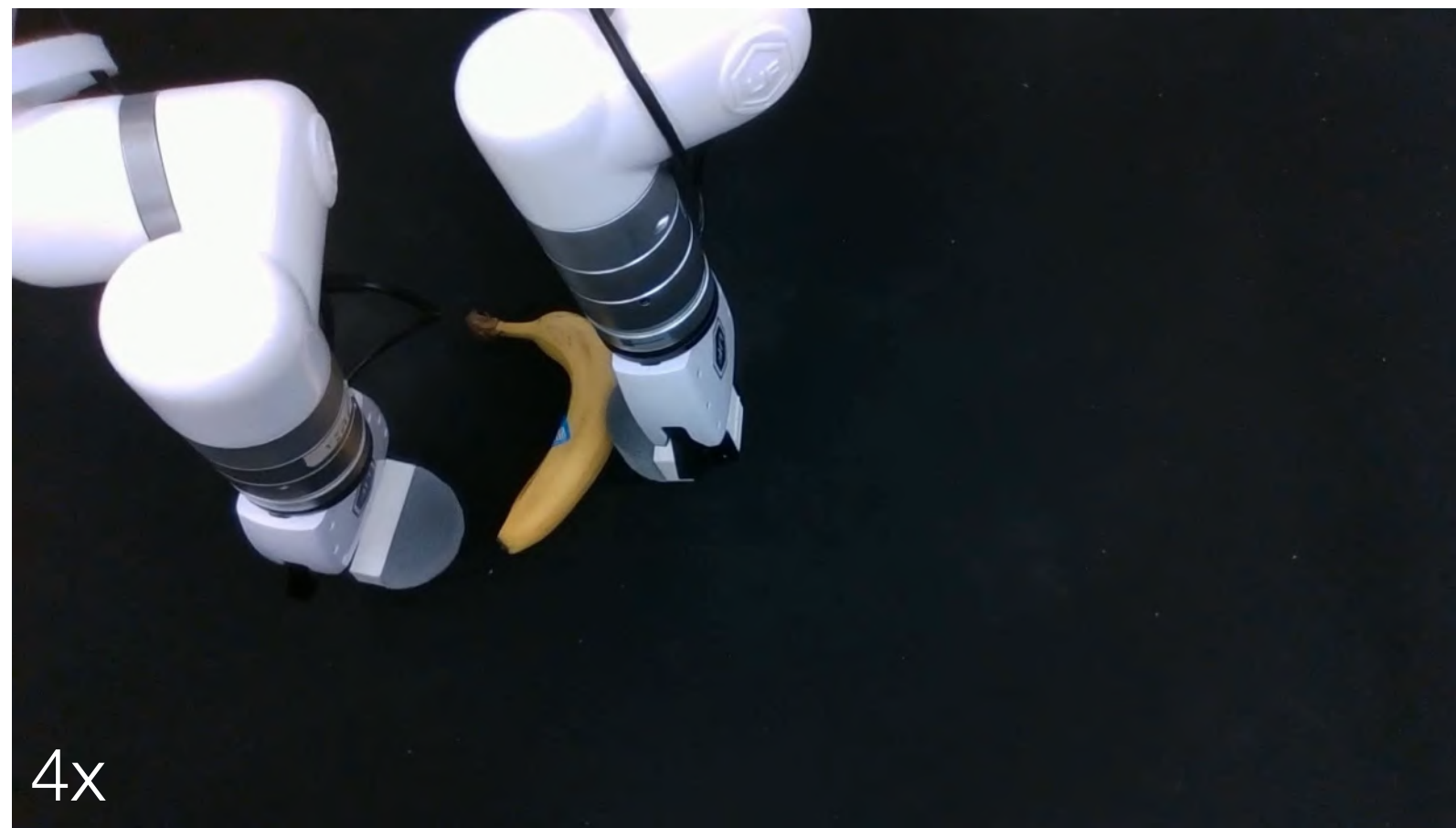
- Diffusion Policy: 22/40
- Ours: **37/40**

Test Generations



Rotation Task

Execution



Scooping Task

- Diffusion Policy: 22/40
- Ours: **36/40**

Test Generations



Scooping Task

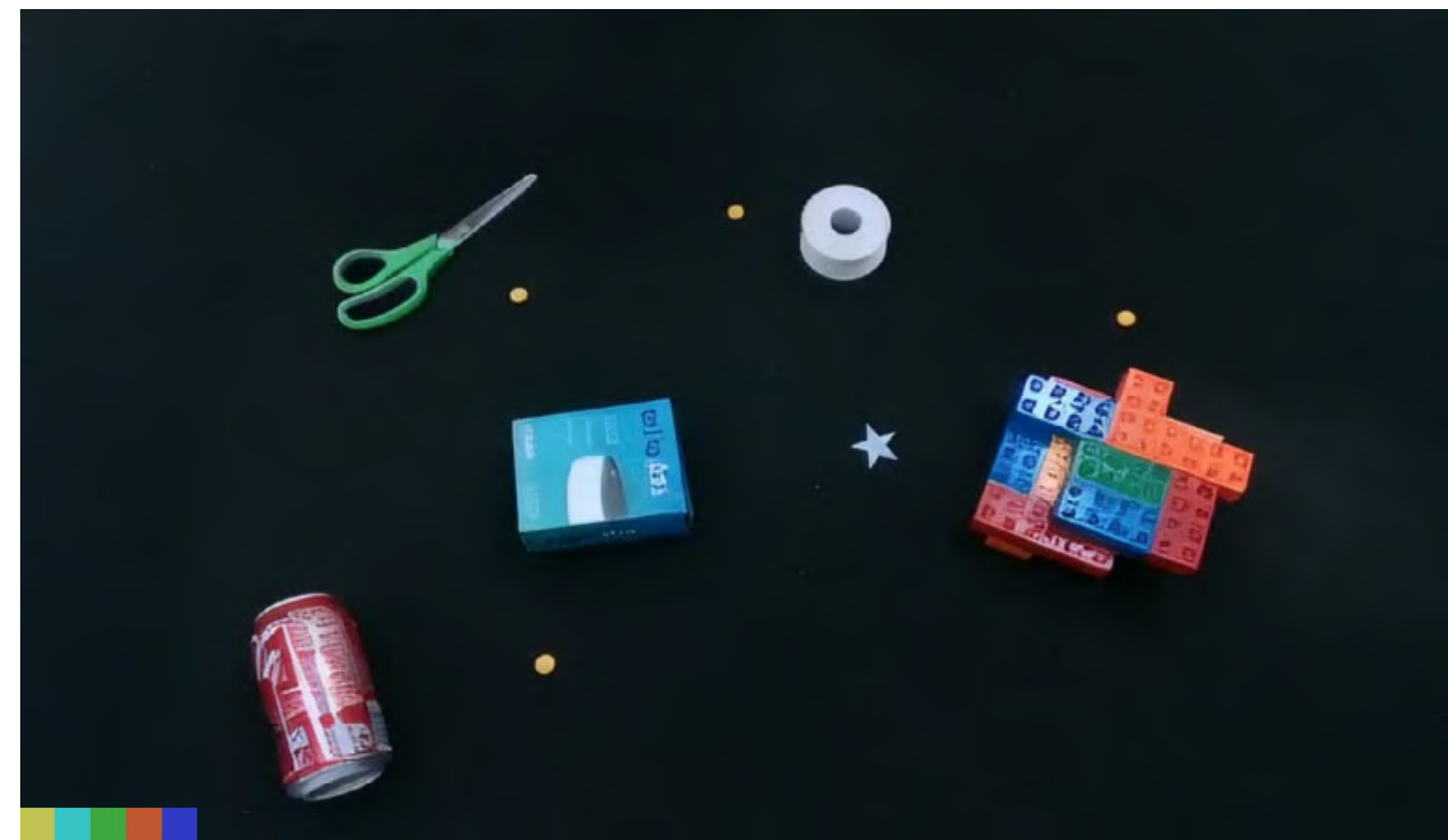
Robot Execution



Sweeping Task

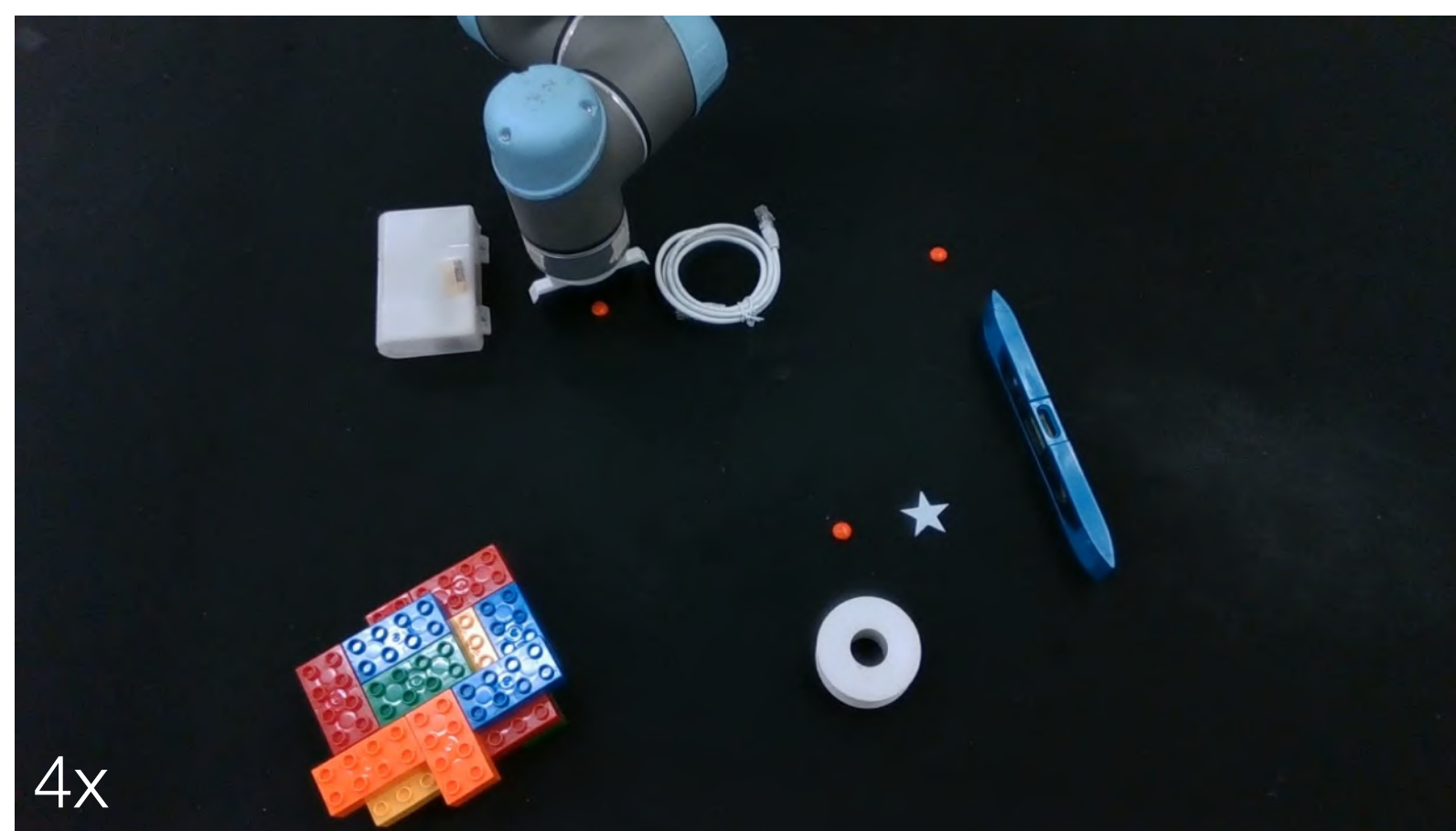
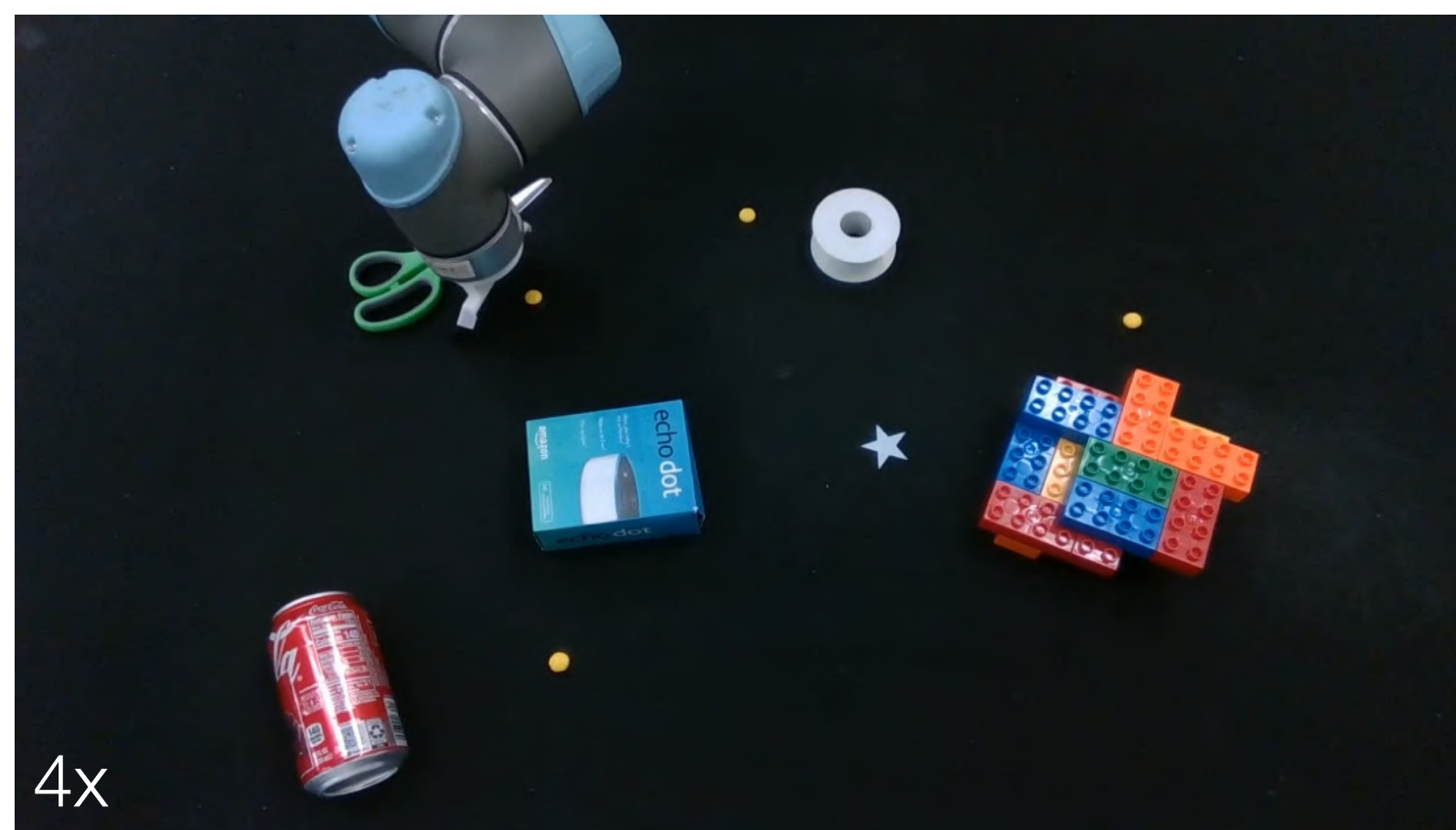
- Diffusion Policy: 5/40
- Ours: **37/40**

Test Generations



Sweeping Task

Robot Execution



Push-Shape Task

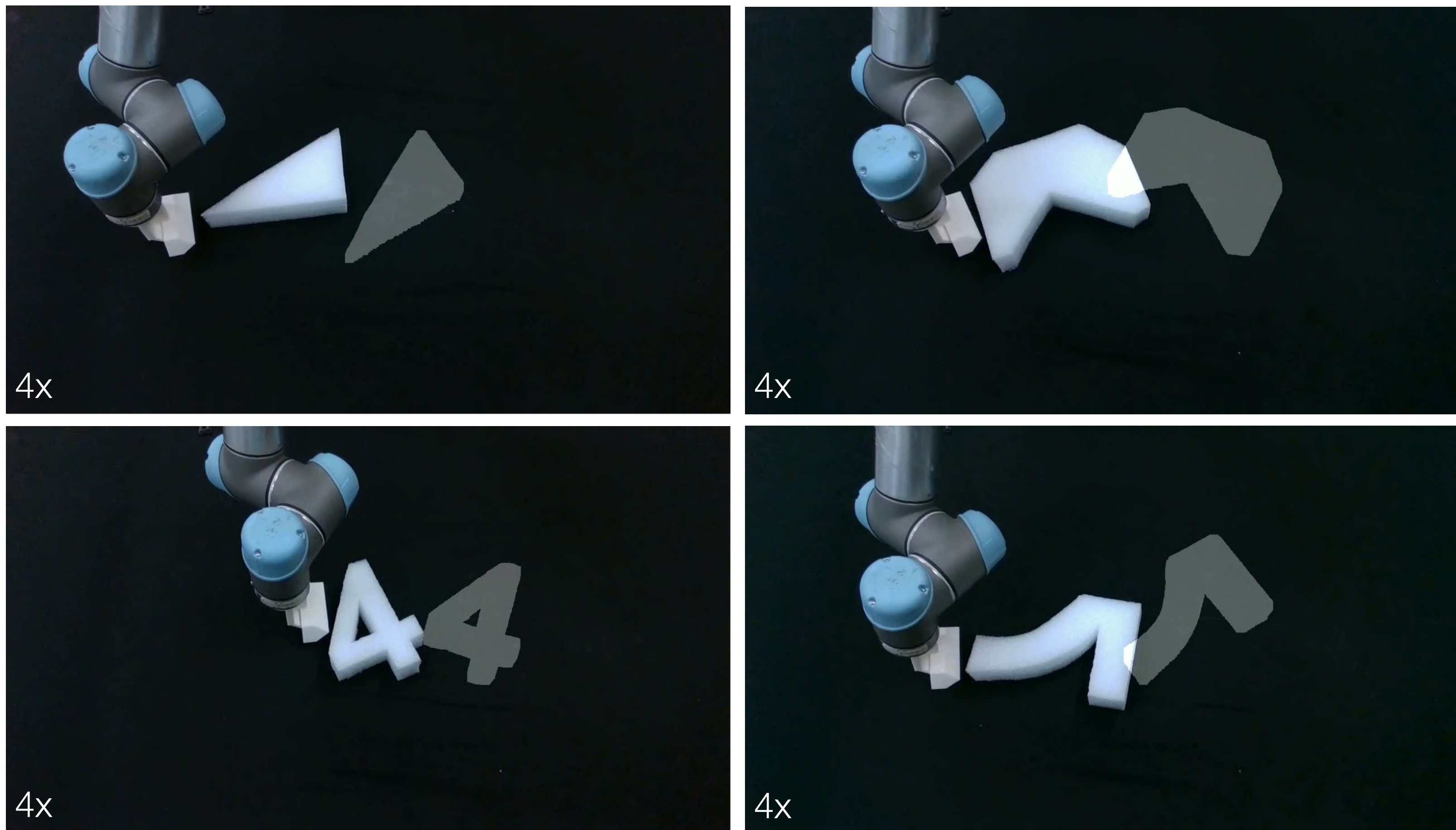
- Diffusion Policy: 0.550 mIoU, 48.2° avg. rotation error
- Ours: **0.731 mIoU, 8.0° avg. rotation error**

Test Generations



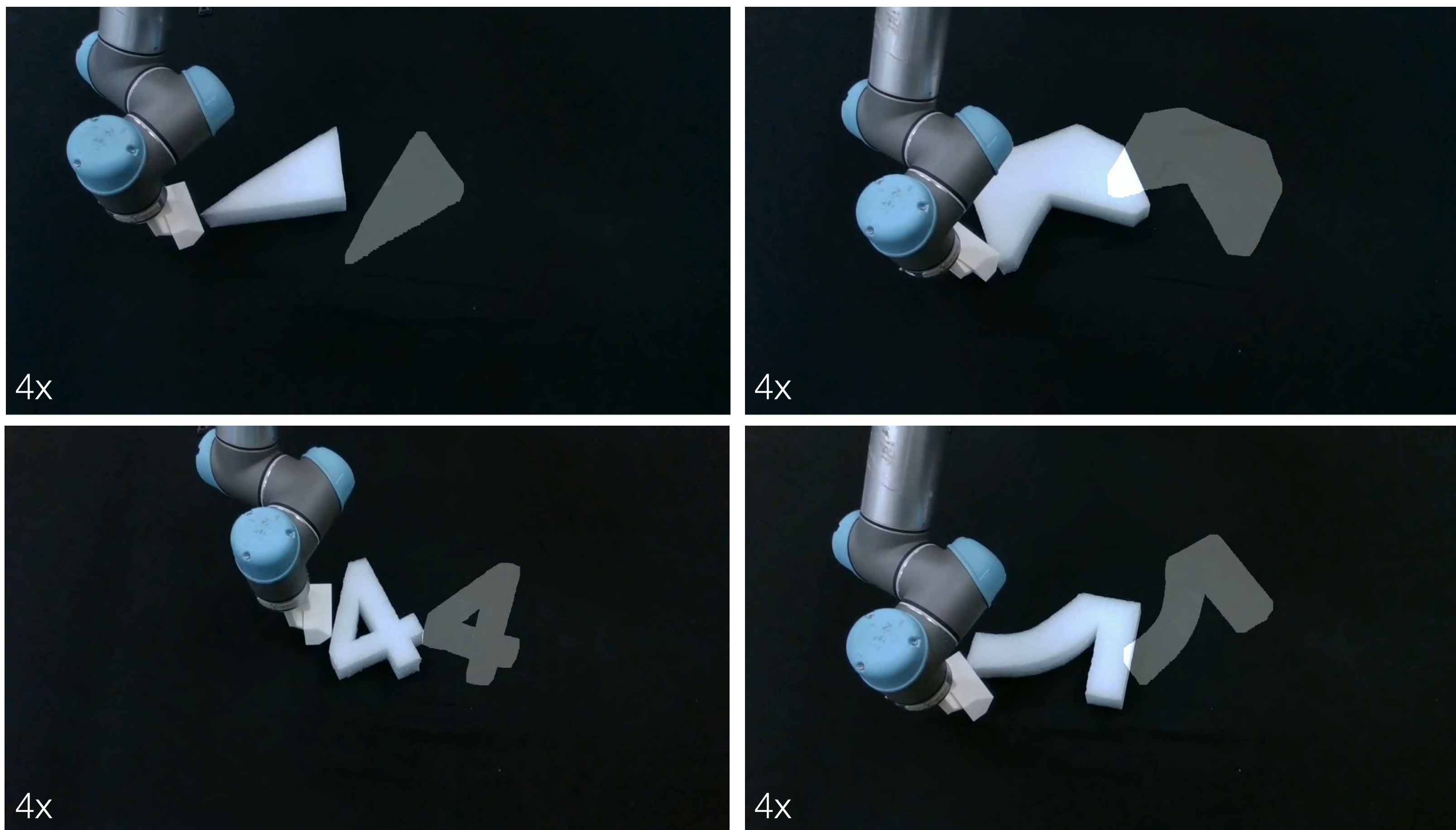
Push-Shape Task

Robot Execution

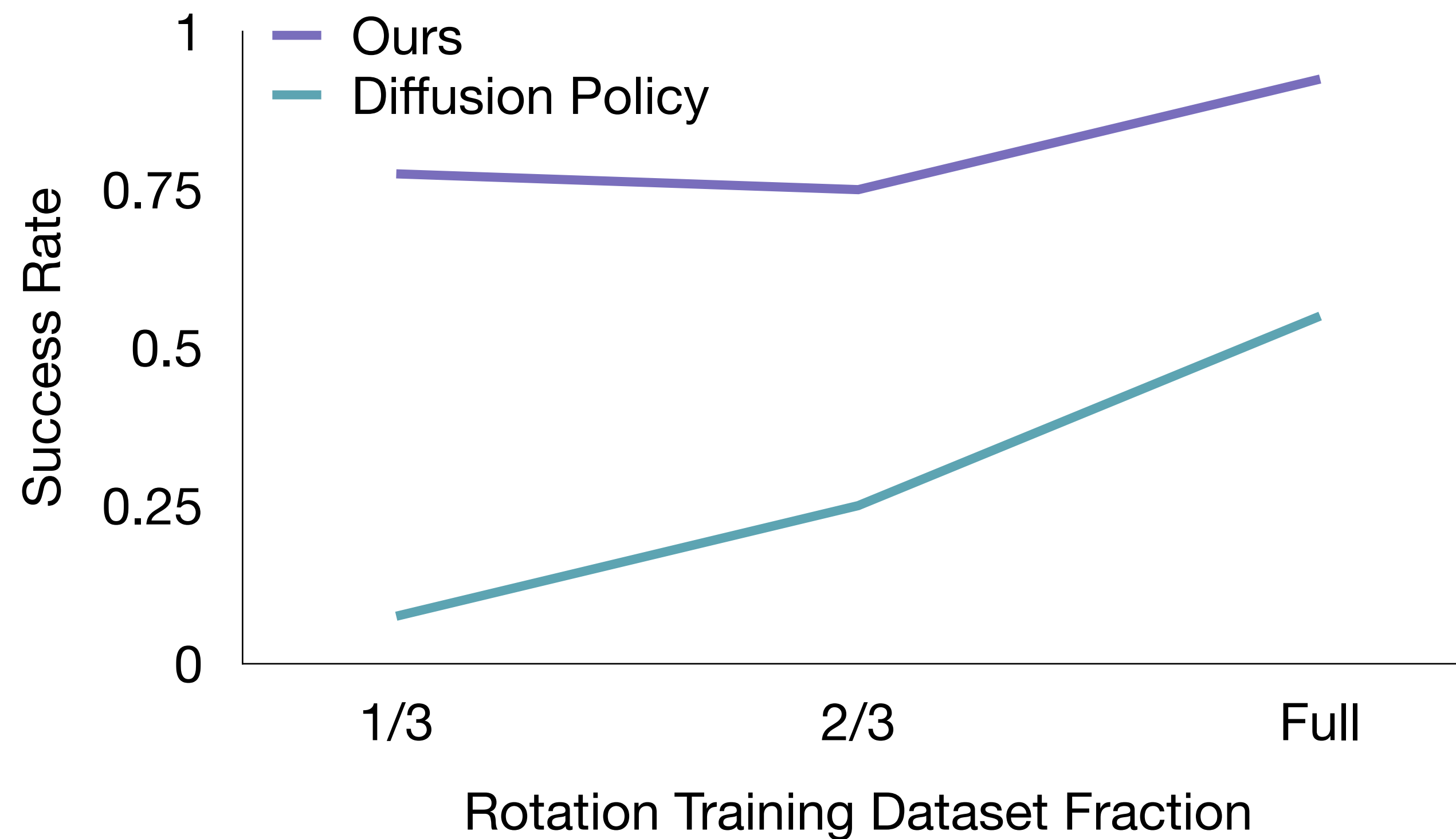


Push-Shape Task

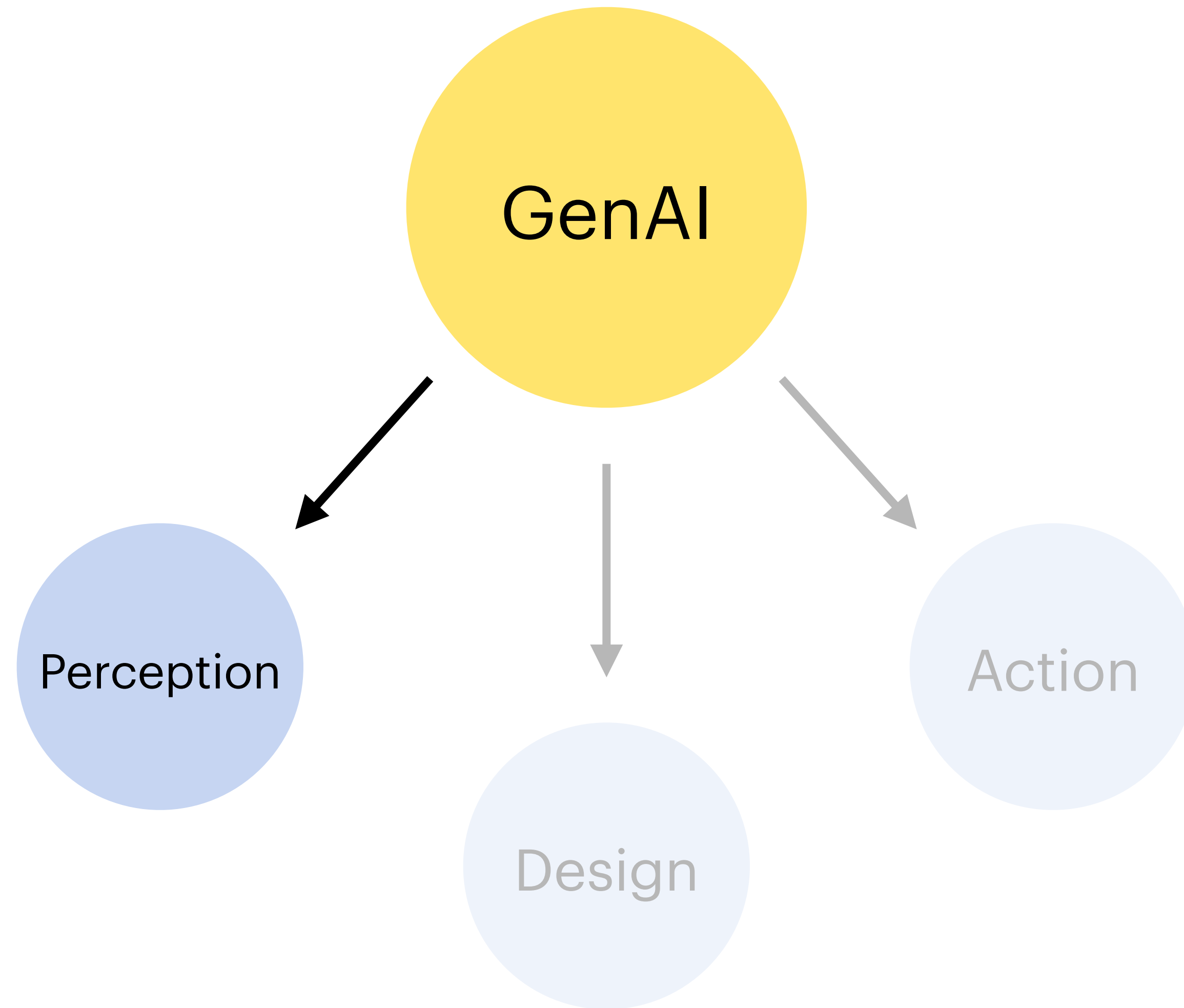
Diffusion Policy Execution



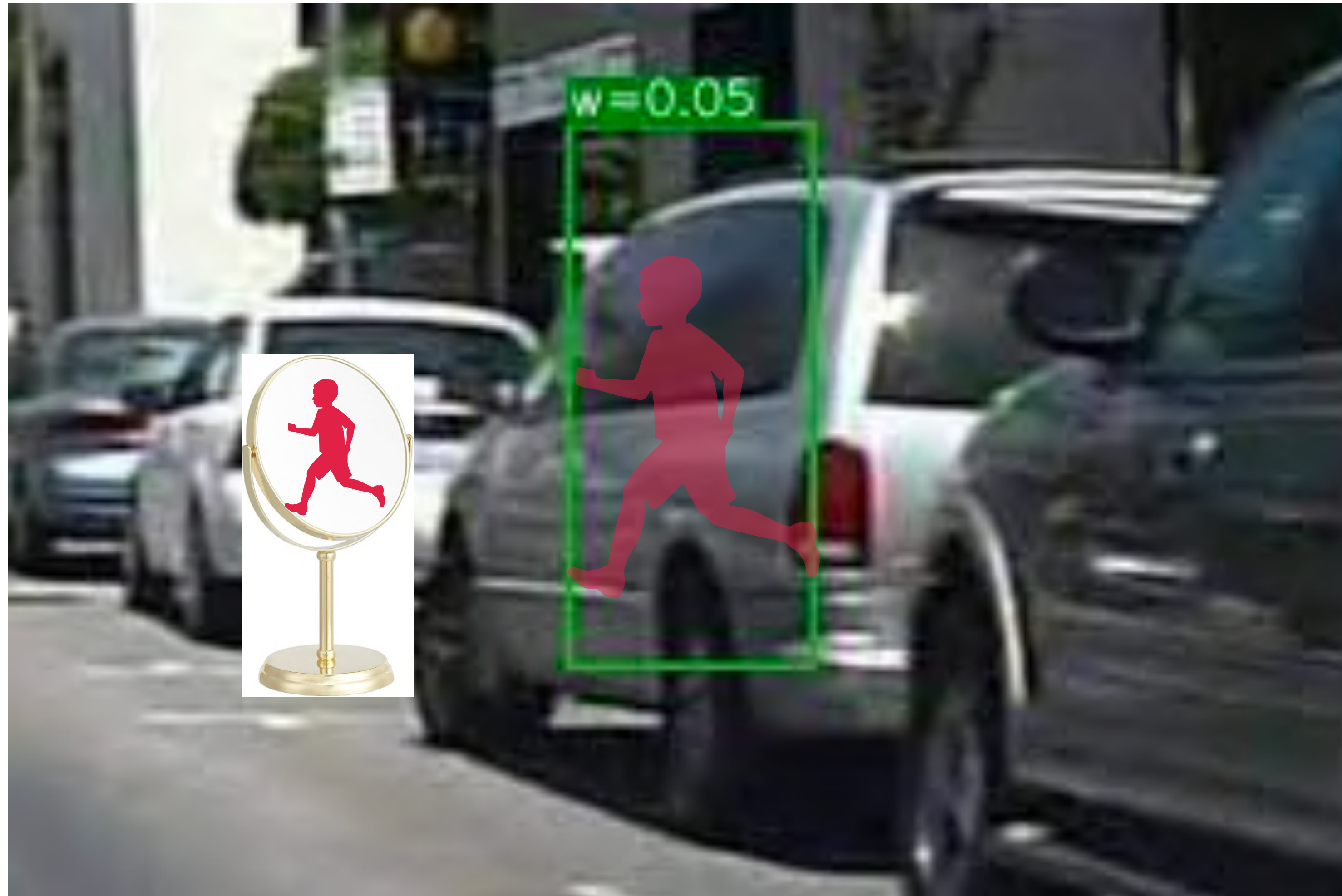
Performance Scaling Curve



Generative Embodied AI



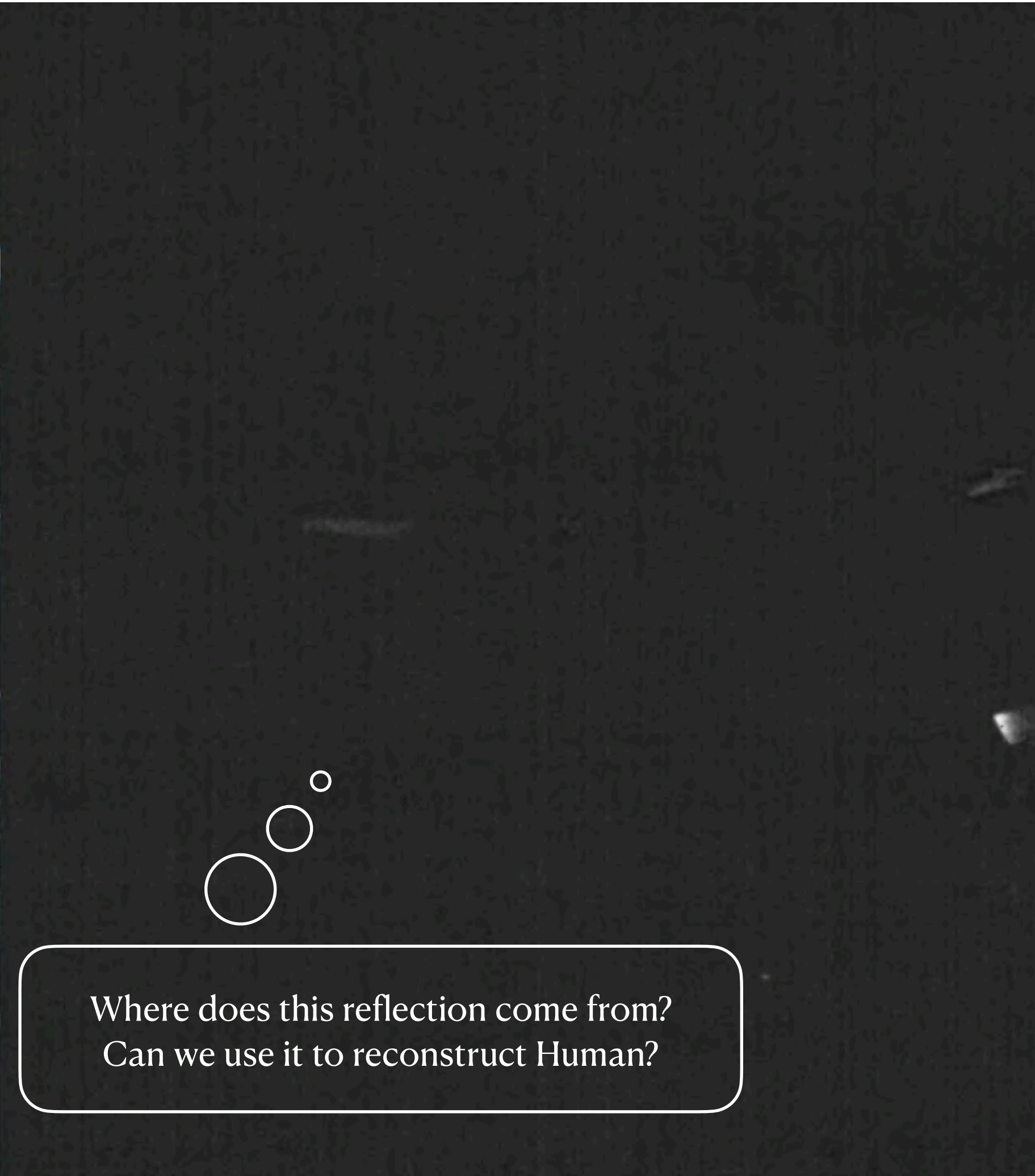
3D Reconstruction Occluded Human



Infrared Turns Cars into Mirrors!



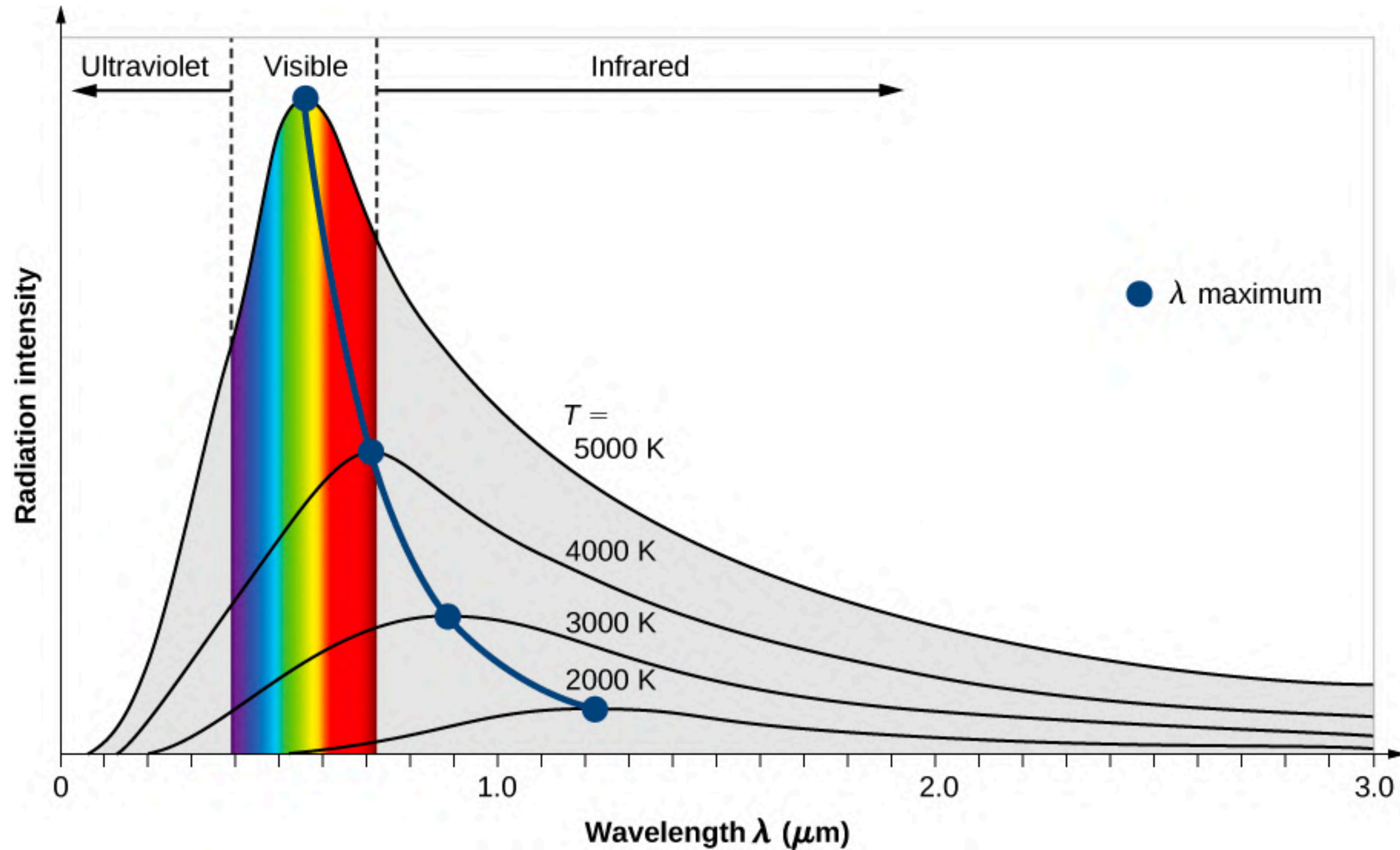
Normal Camera (0.4-0.7 μm)



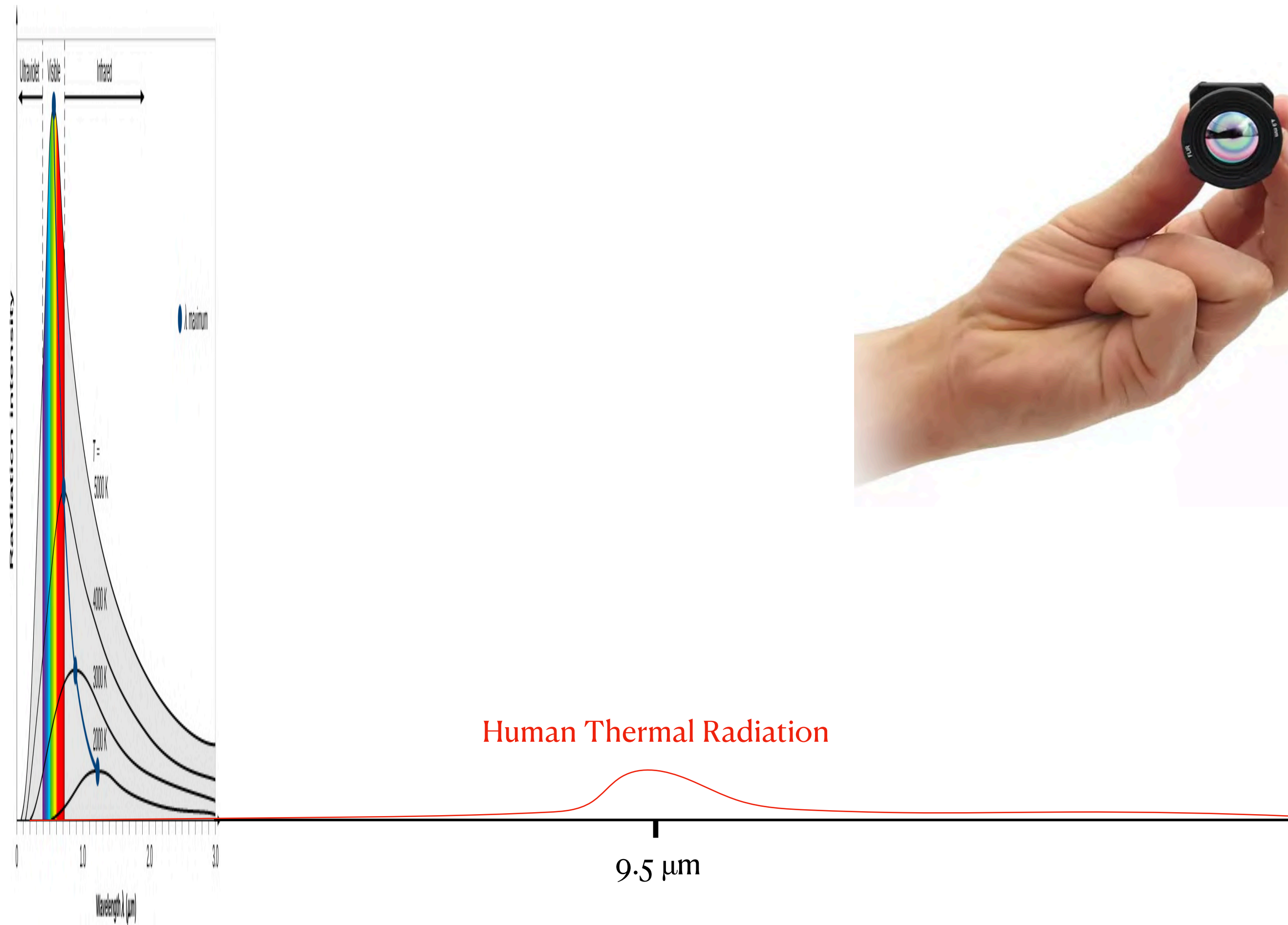
Where does this reflection come from?
Can we use it to reconstruct Human?

Thermal Camera (7-14 μm)

Black-body Radiation



Human bodies are long wavelength infrared light bulbs



Humans are infrared light bulbs



Normal Camera (0.4-0.7 μm)

Thermal Camera (7-14 μm)

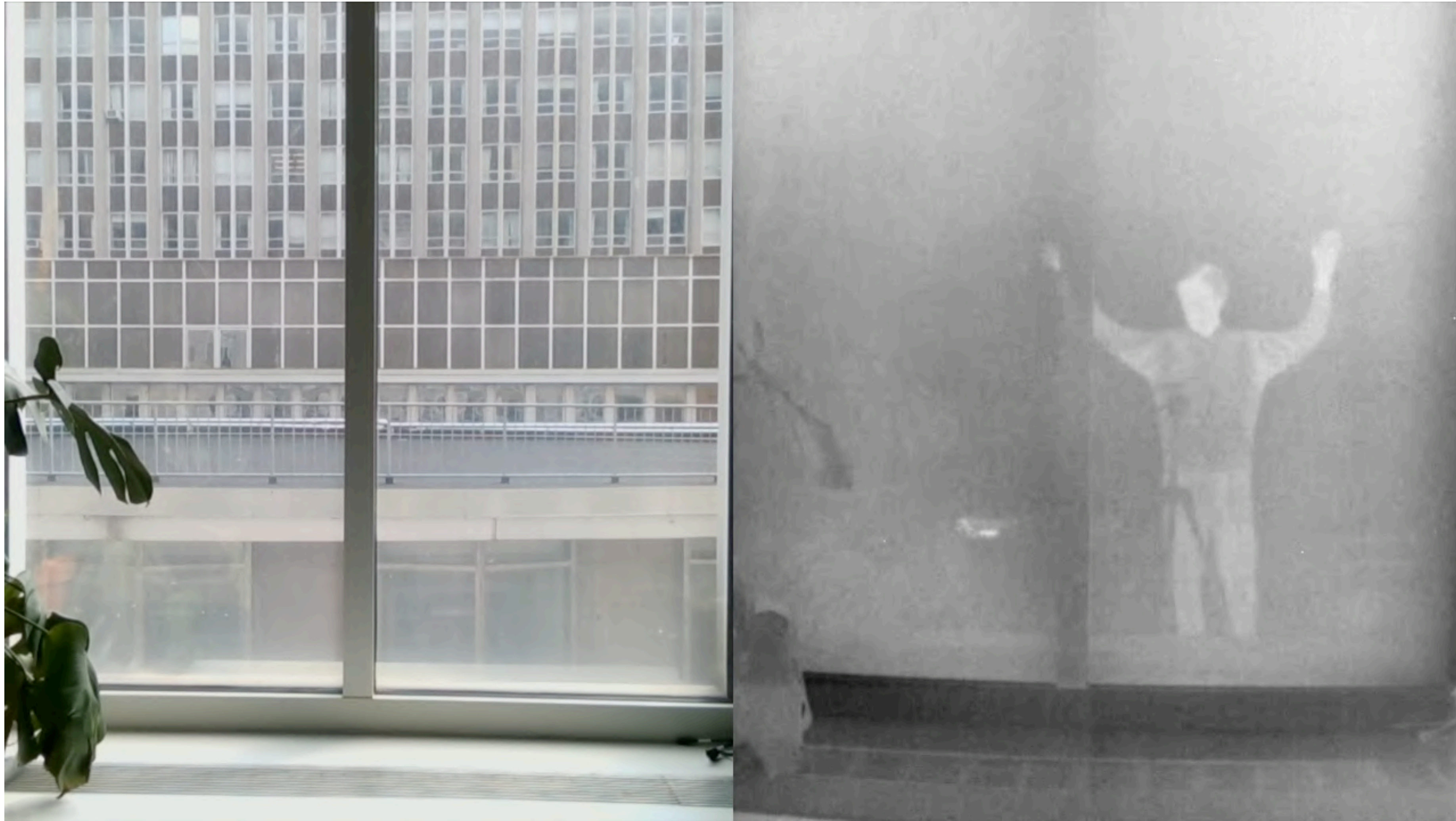
Cats are infrared light bulbs



Normal Camera (0.4-0.7 μm)

Thermal Camera (7-14 μm)

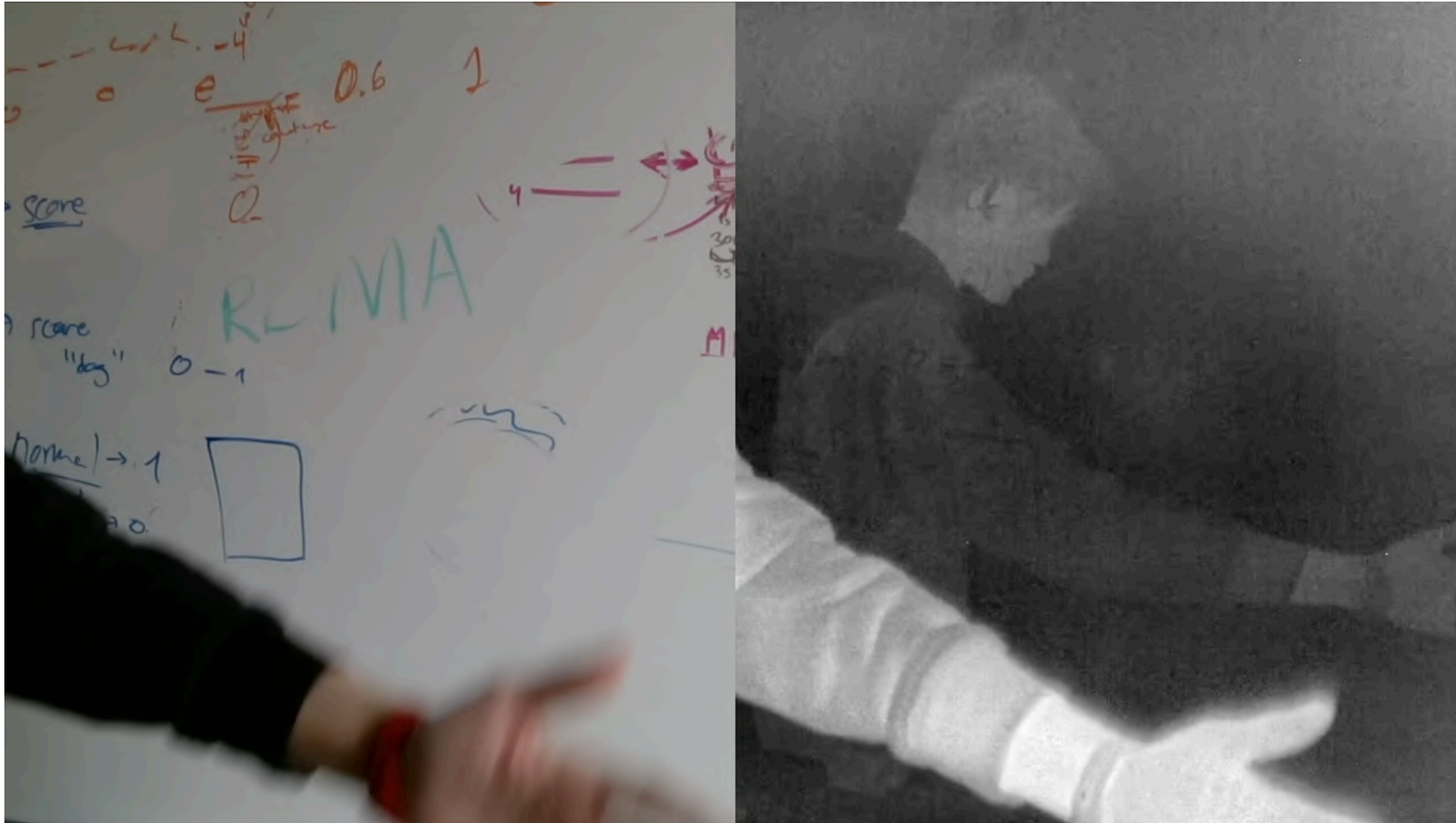
Many objects are more reflective in infrared



Normal Camera (0.4-0.7 μm)

Thermal Camera (7-14 μm)

Many objects are more reflective in infrared



Normal Camera (0.4-0.7 μm)

Thermal Camera (7-14 μm)

Many objects are more reflective in infrared



Normal Camera (0.4-0.7 μm)



Thermal Camera (7-14 μm)

Many objects are more reflective in infrared



Normal Camera (0.4-0.7 μm)

Thermal Camera (7-14 μm)

Many objects are more reflective in infrared

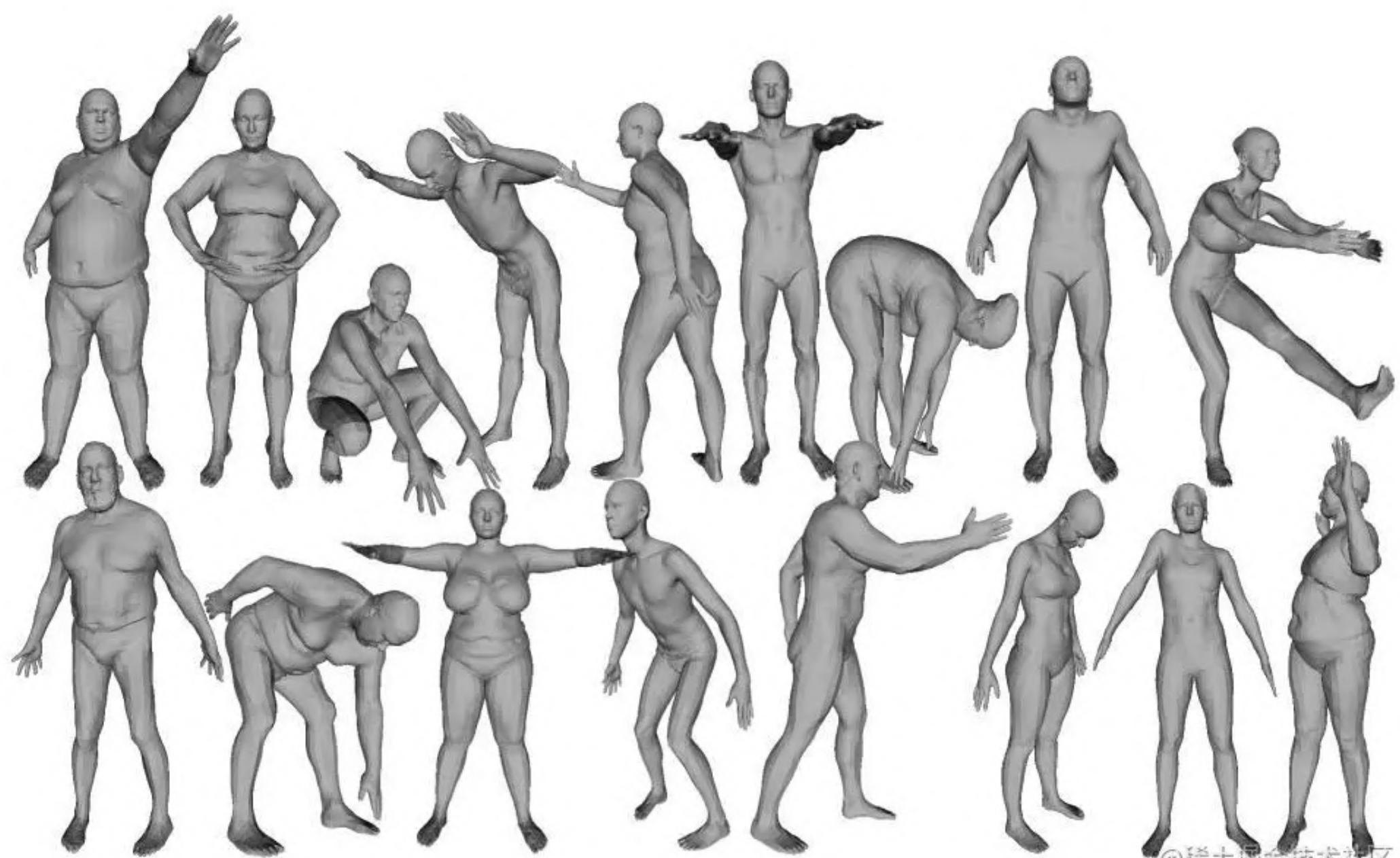


Normal Camera (0.4-0.7 μm)



Thermal Camera (7-14 μm)

3D Generative Models



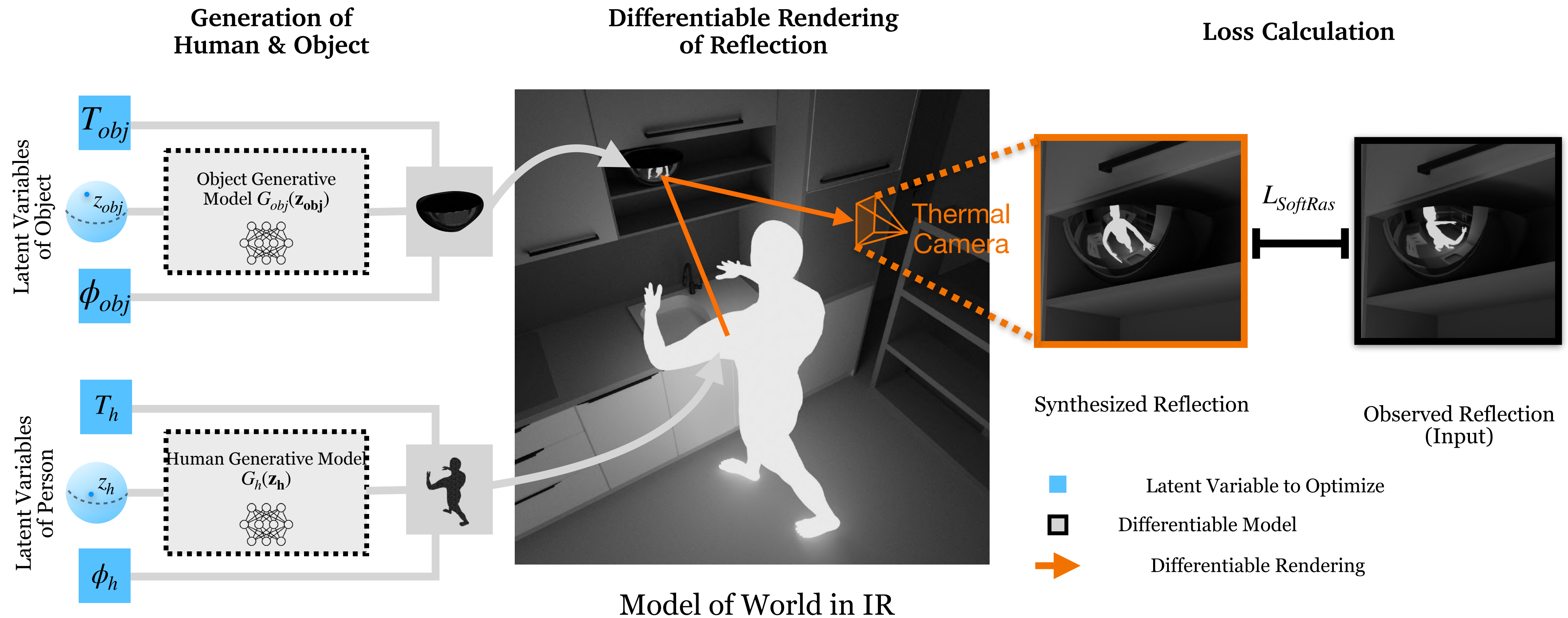
3D Human



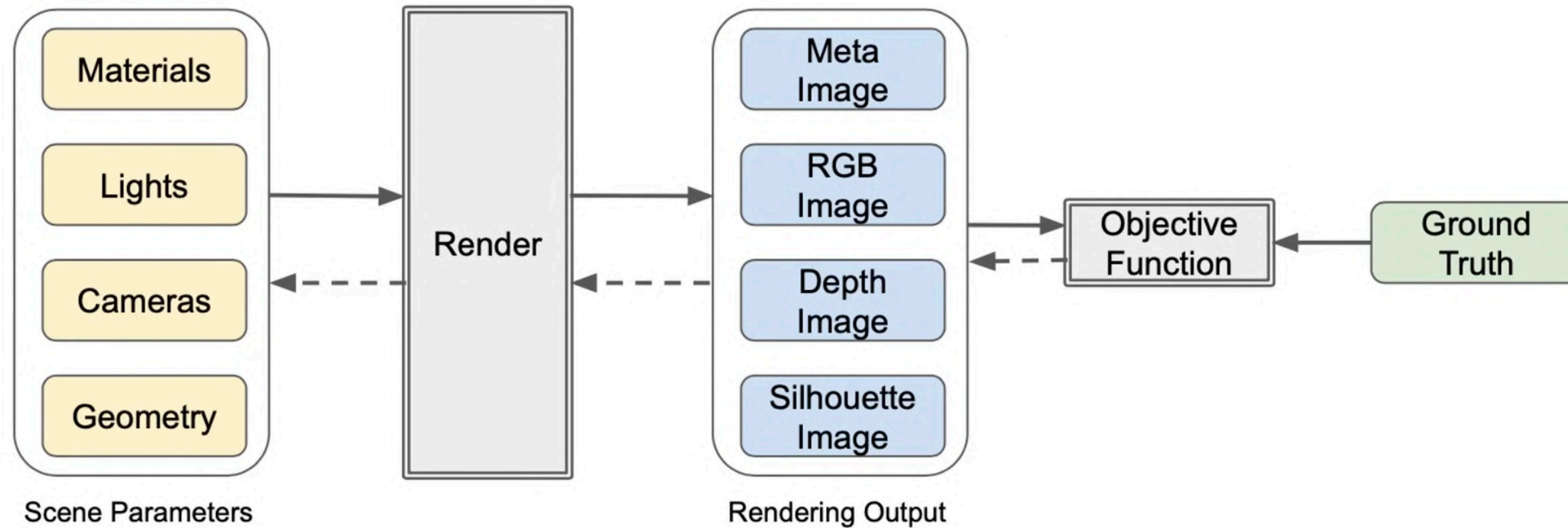
3D Objects



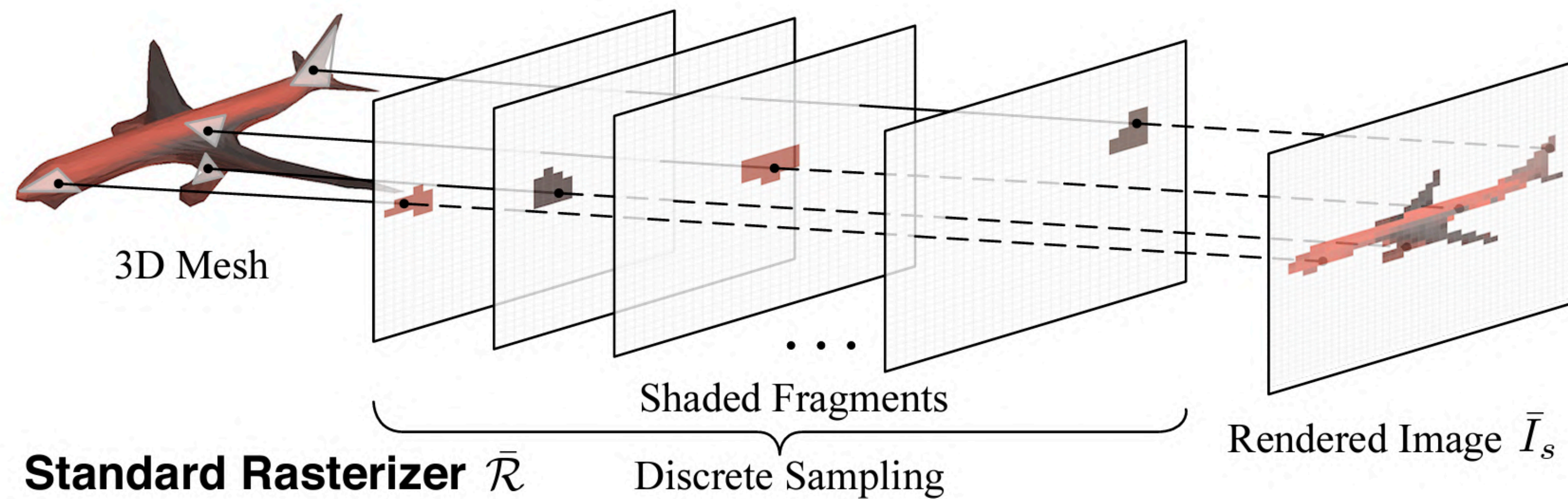
Method



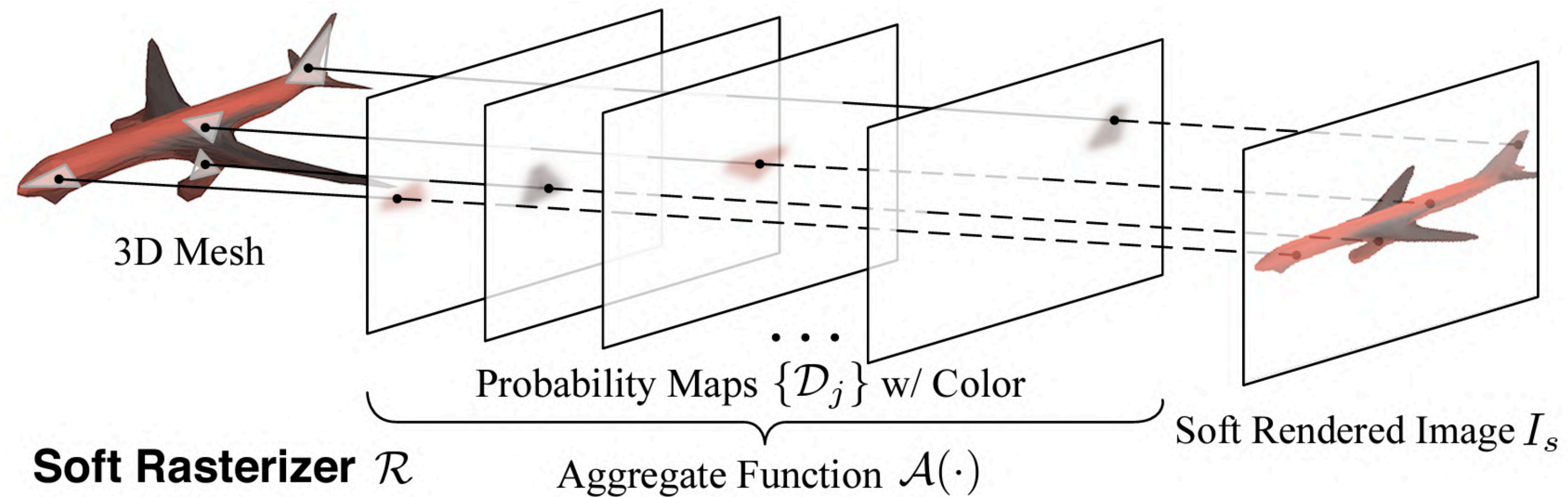
Differentiable Rendering



Non-differentiable Rendering

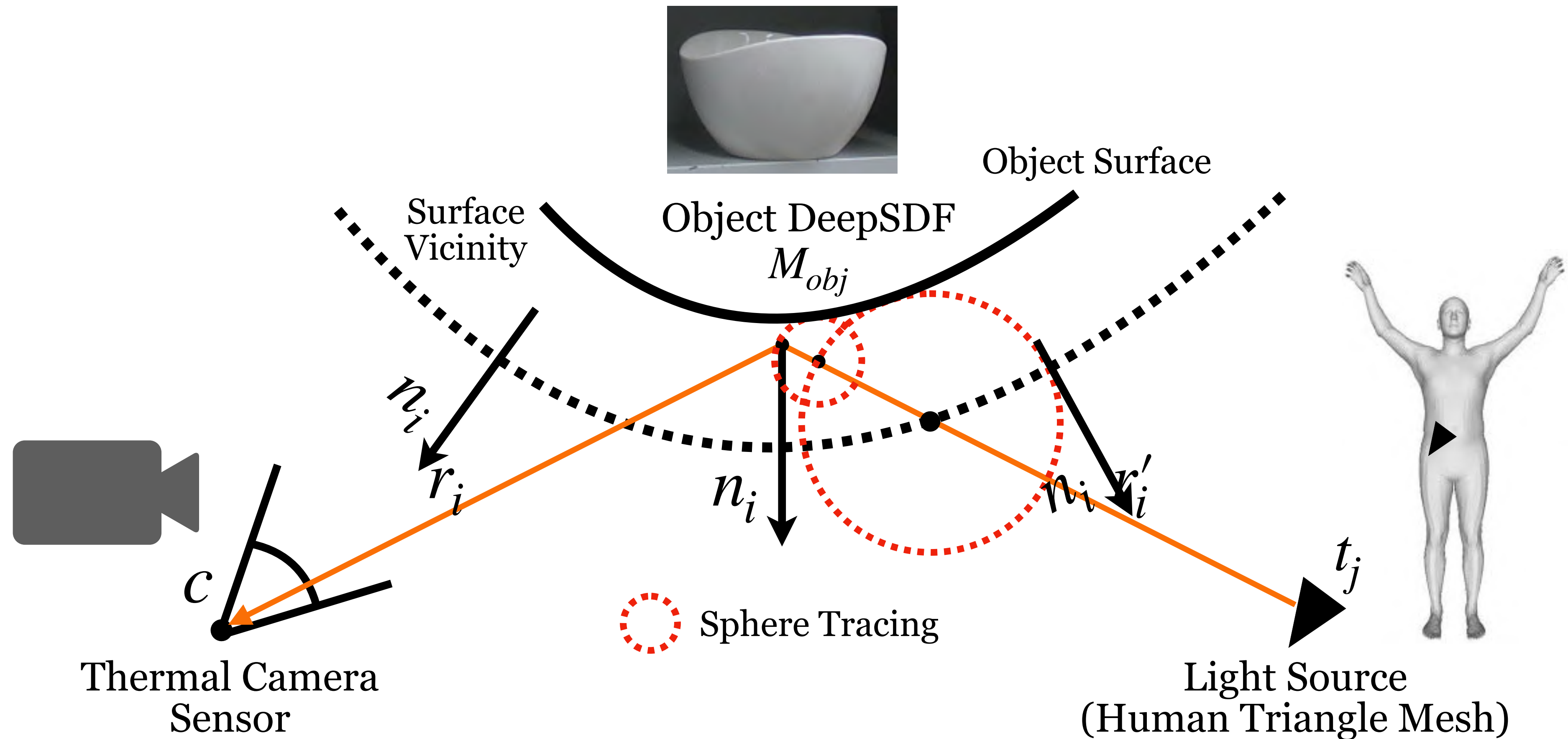


Differentiable Rendering



What about reflection?

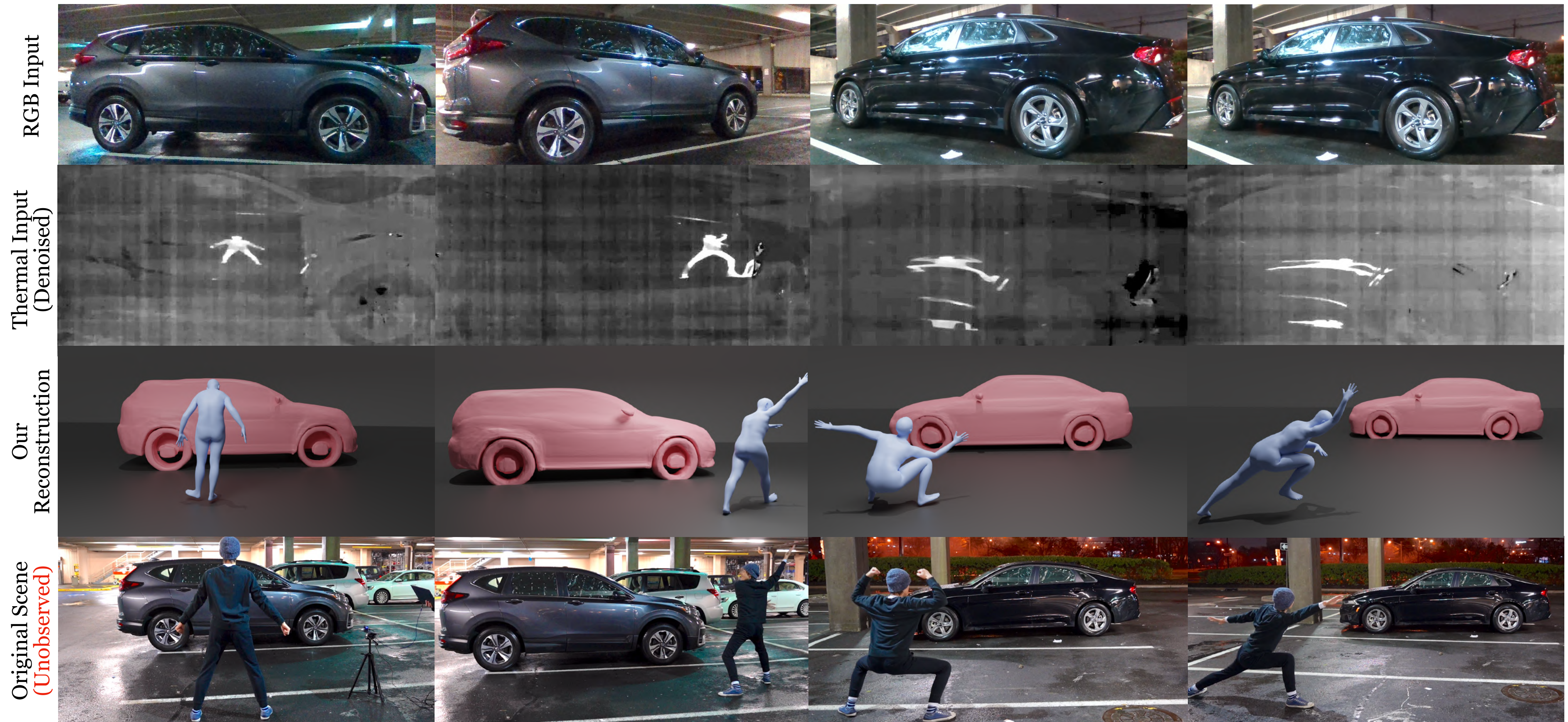
Differentiable Rendering of Reflection



Results



Results



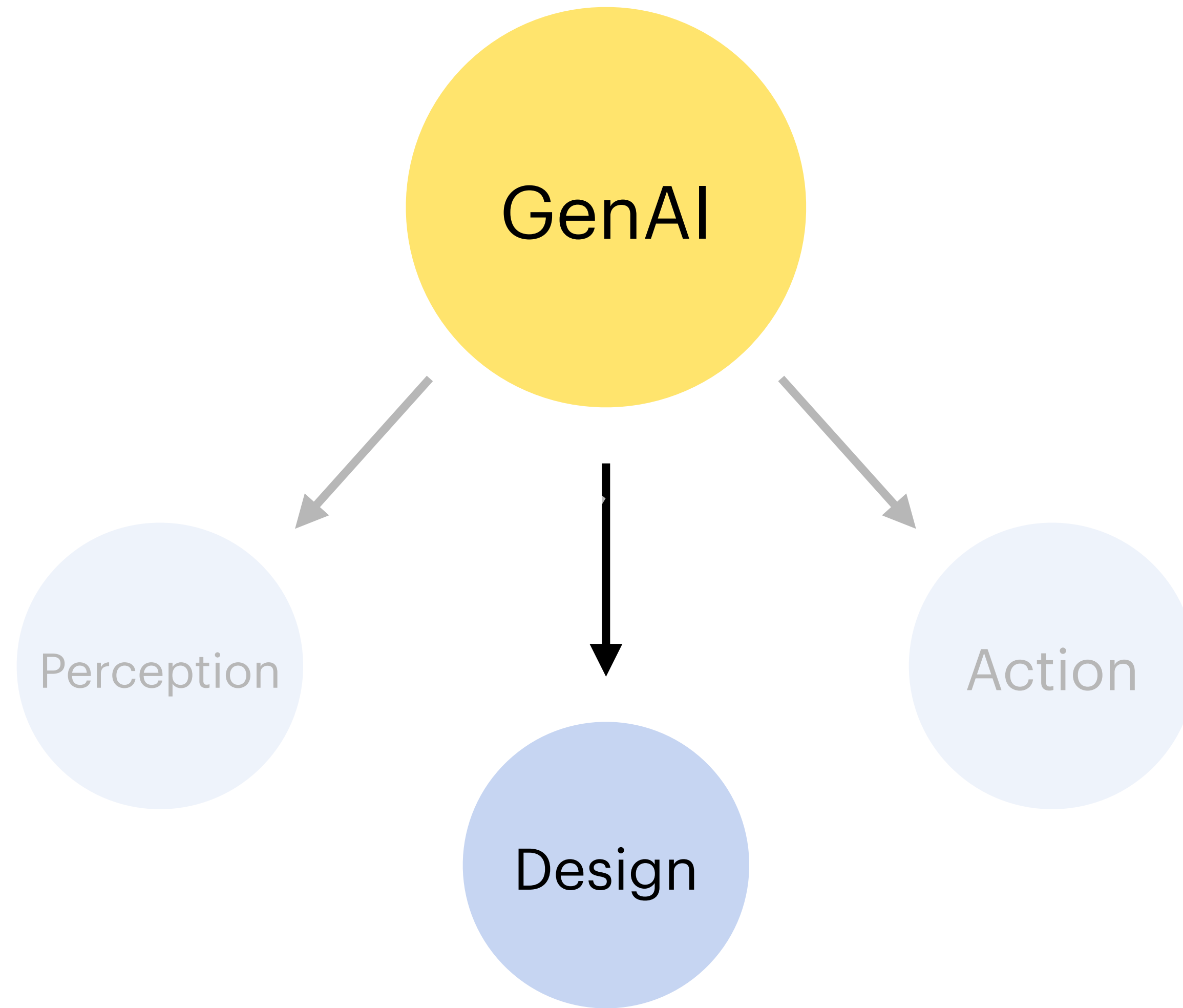
RGB Input



RGB Input



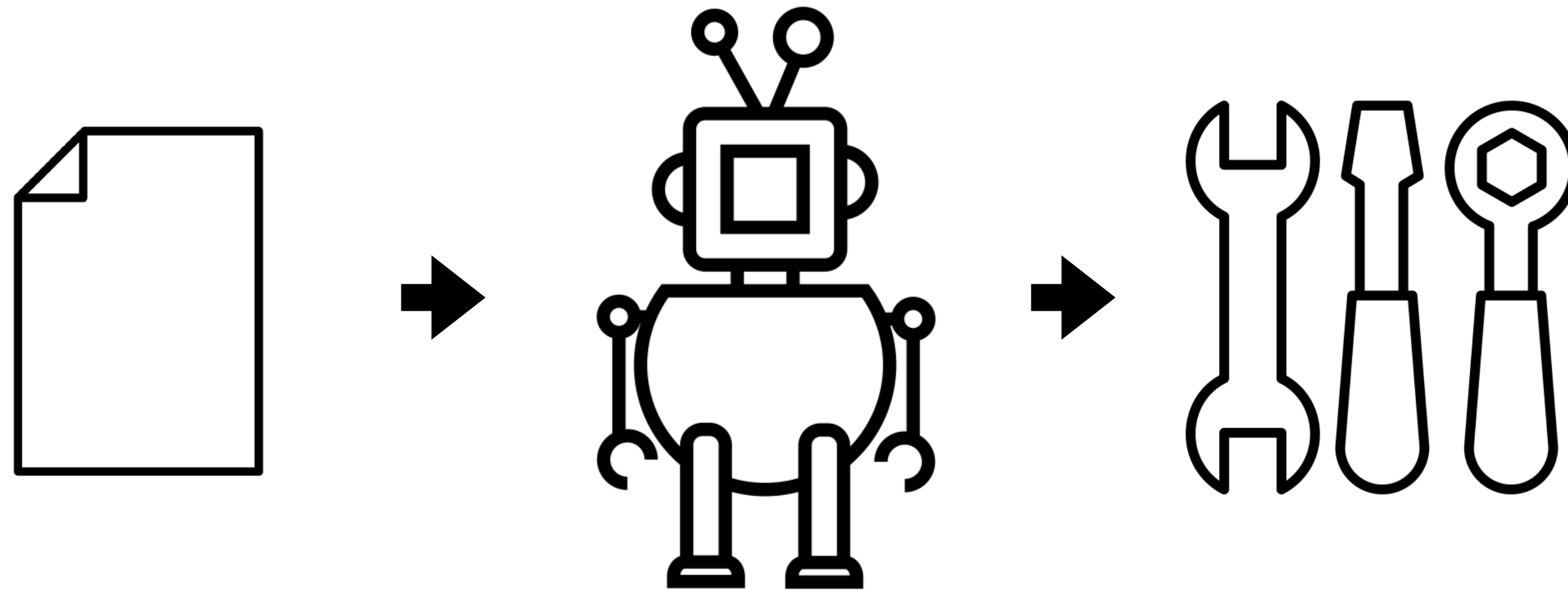
Generative Embodied AI



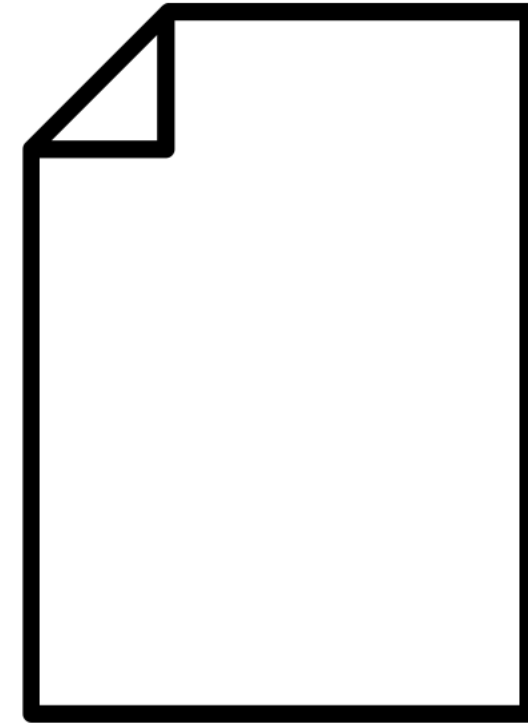
Robotic Tool Design



Designing Paper Tools



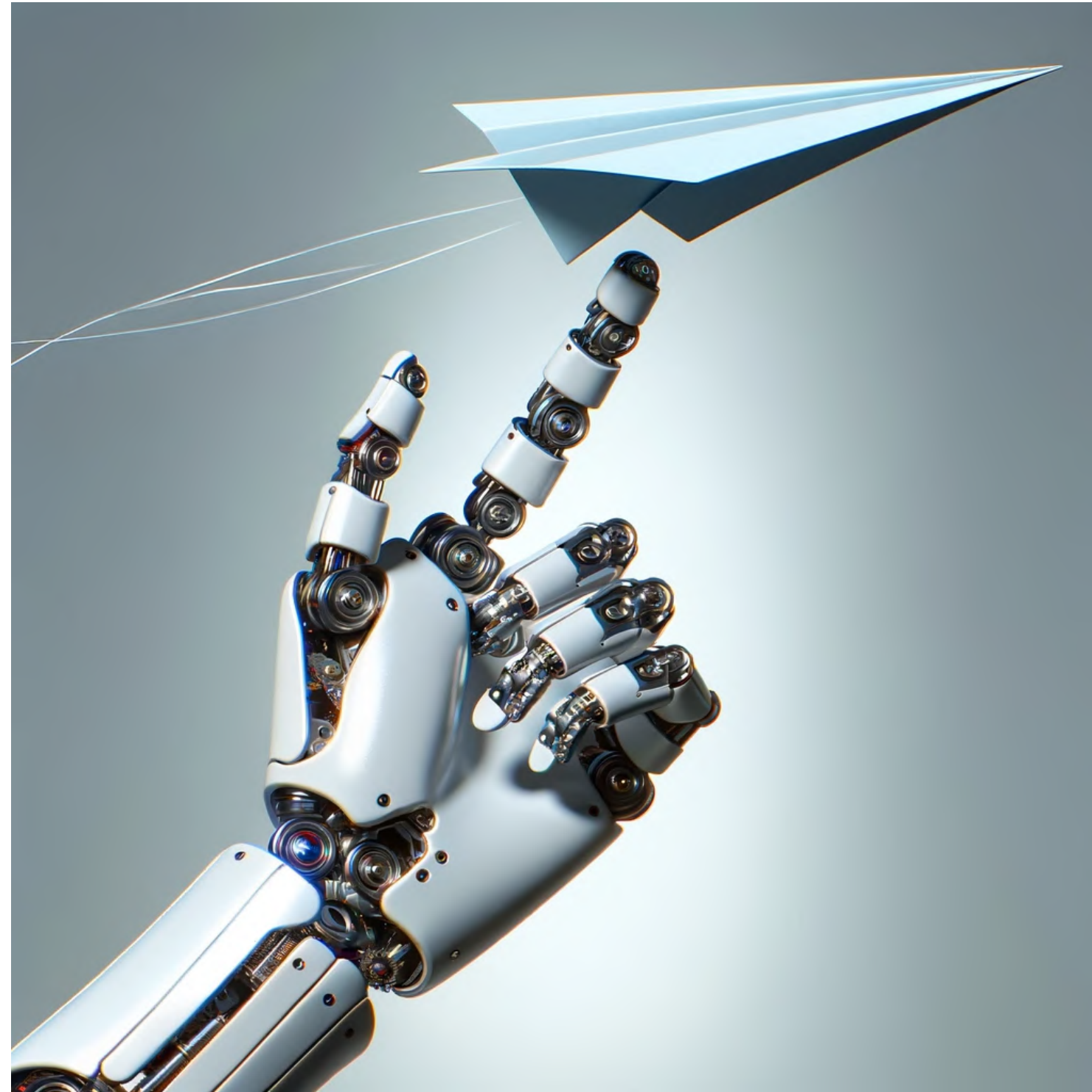
Why Paper?



Paper Makes Practical Recyclable Tools!



Case Study: Paper Airplane





Design



Design



Build



Design



Throw



Design



Measure



Design



Build



Throw



Measure



Design



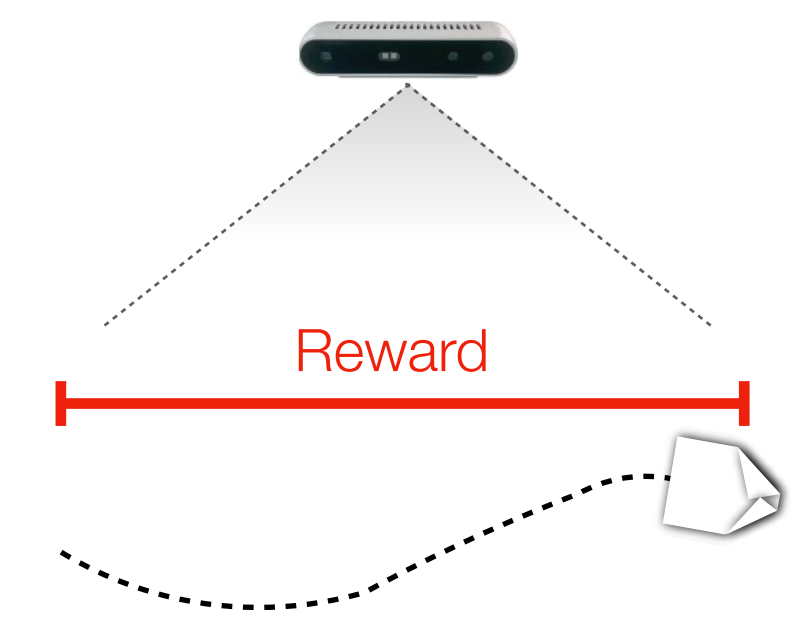
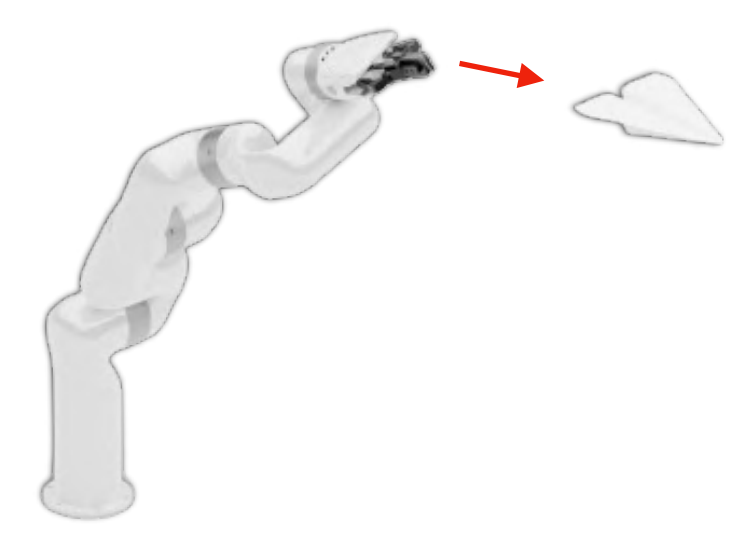
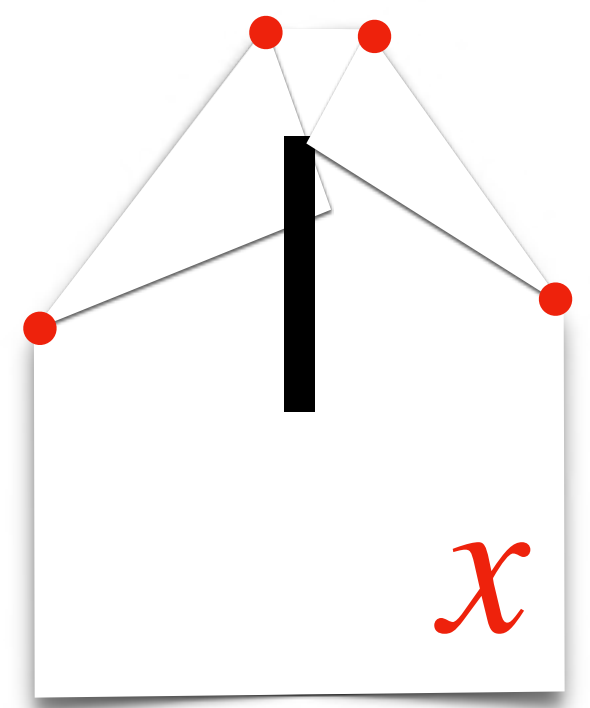
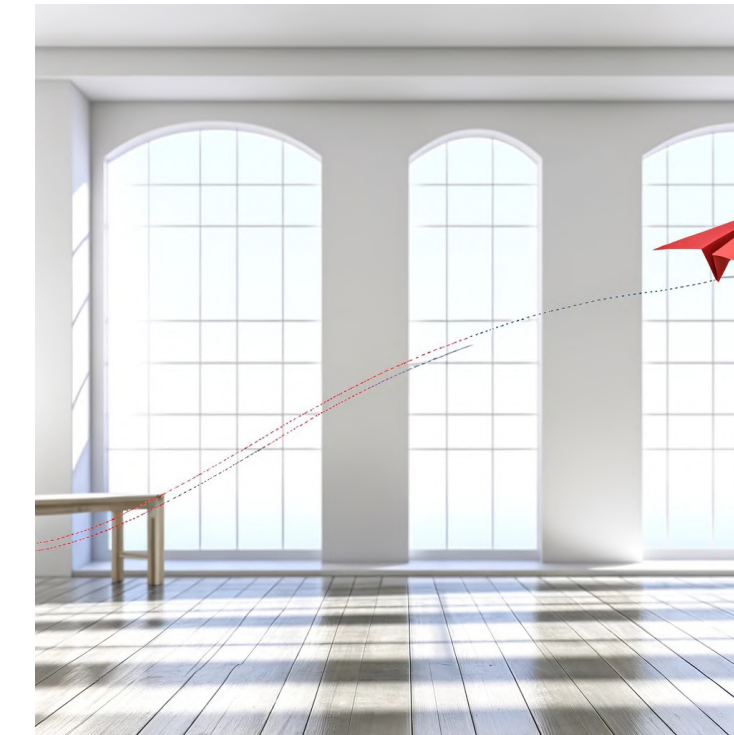
Build



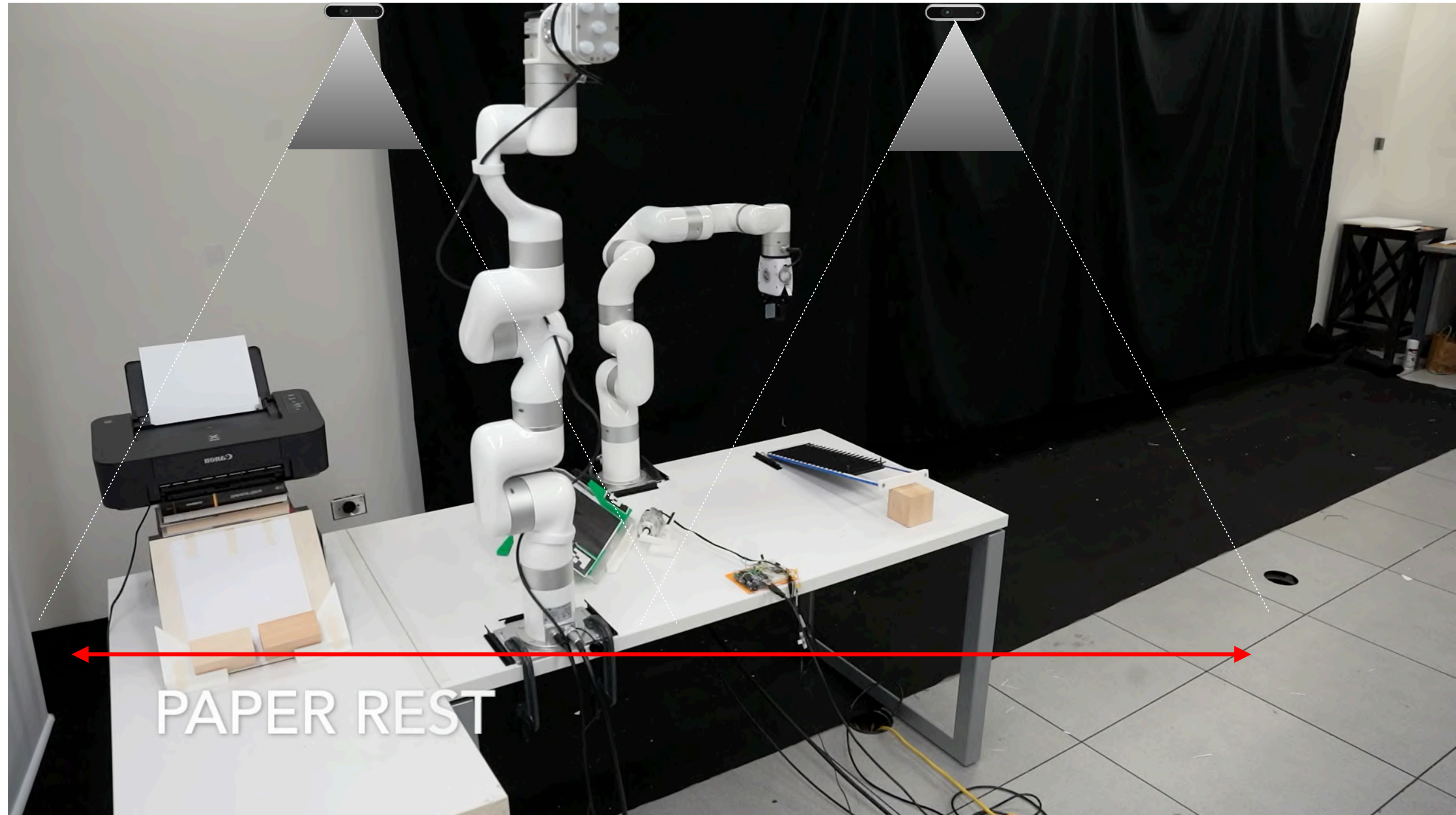
Throw



Measure



Automation



PAPER REST

Design



Build



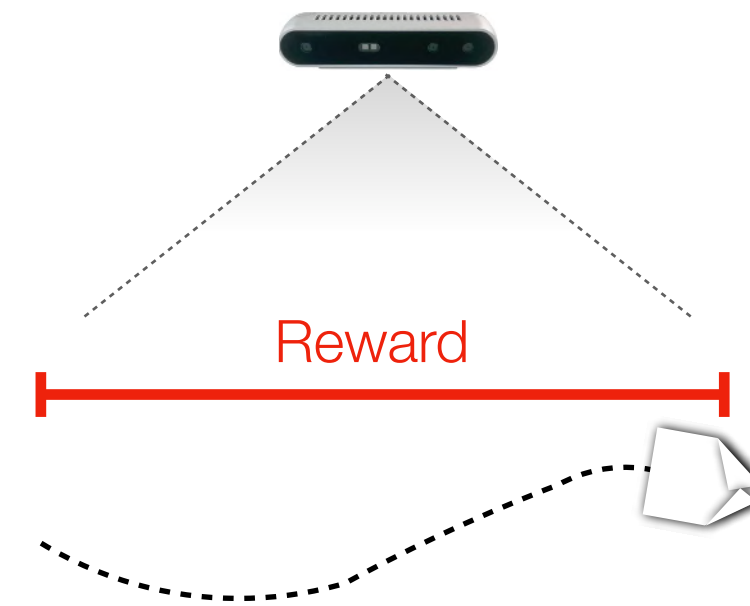
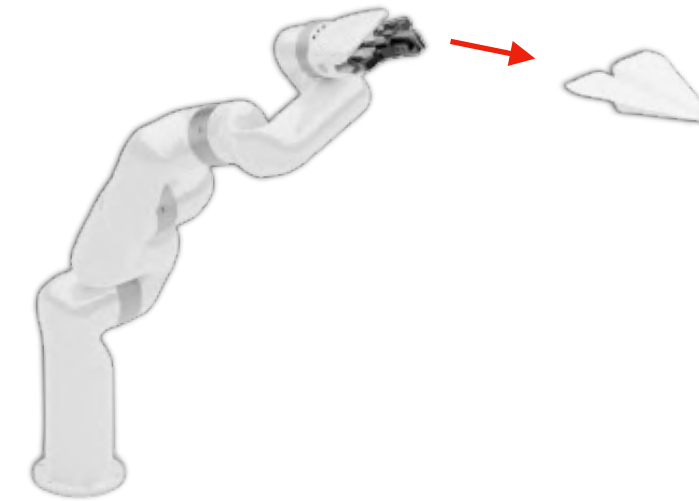
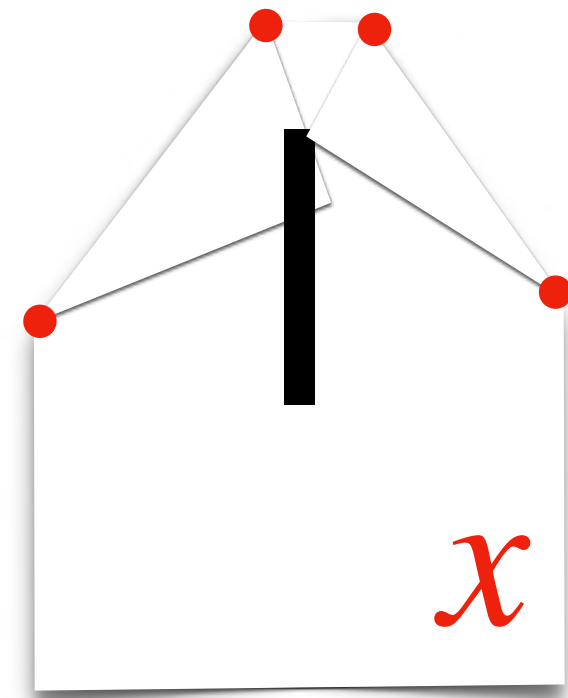
Throw



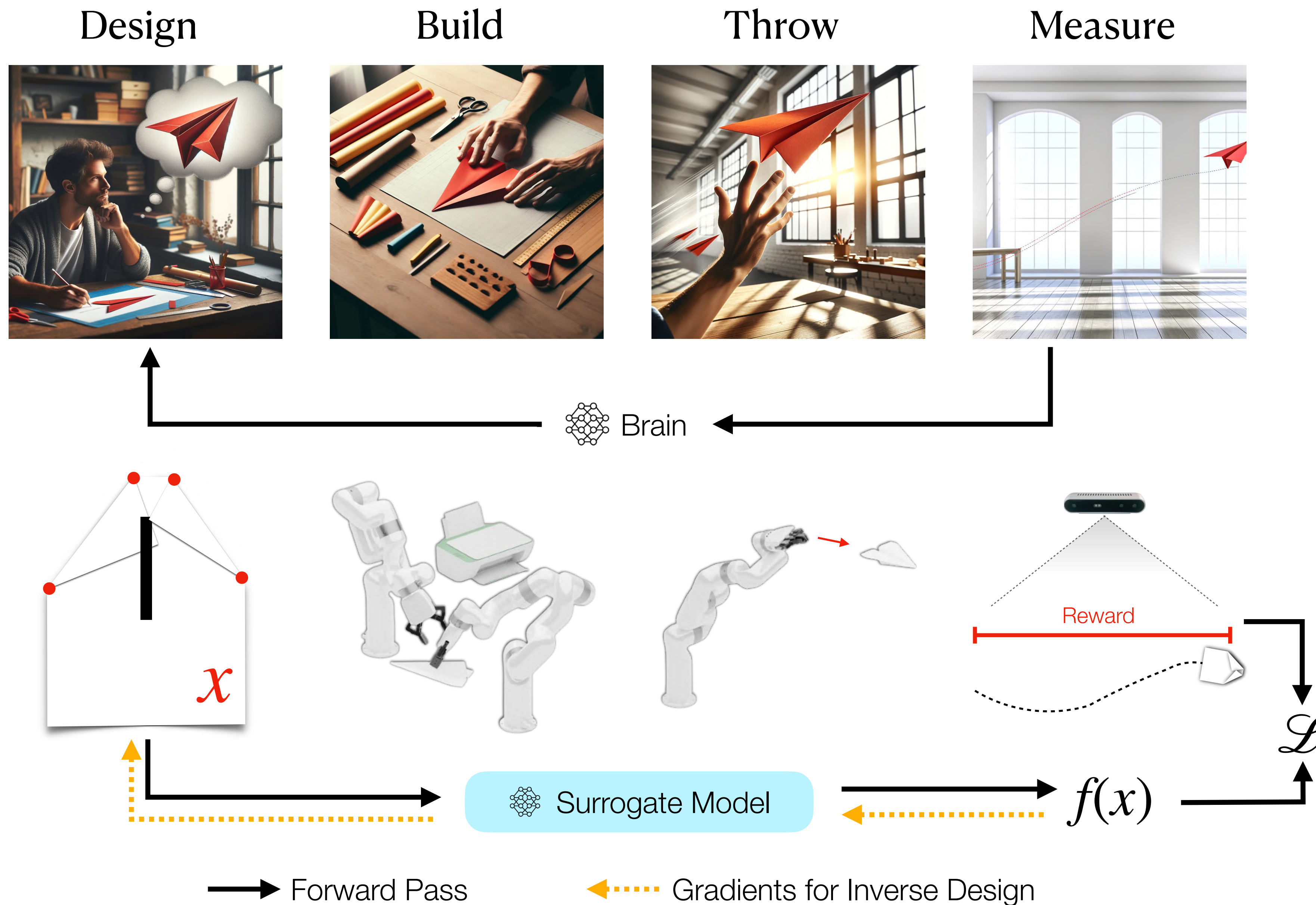
Measure



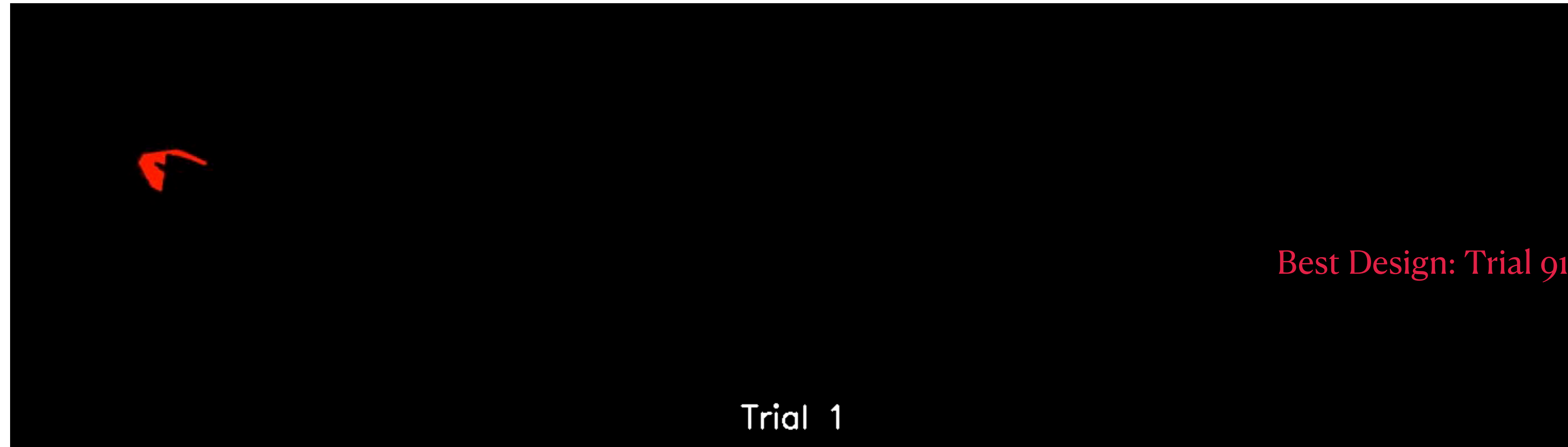
Brain



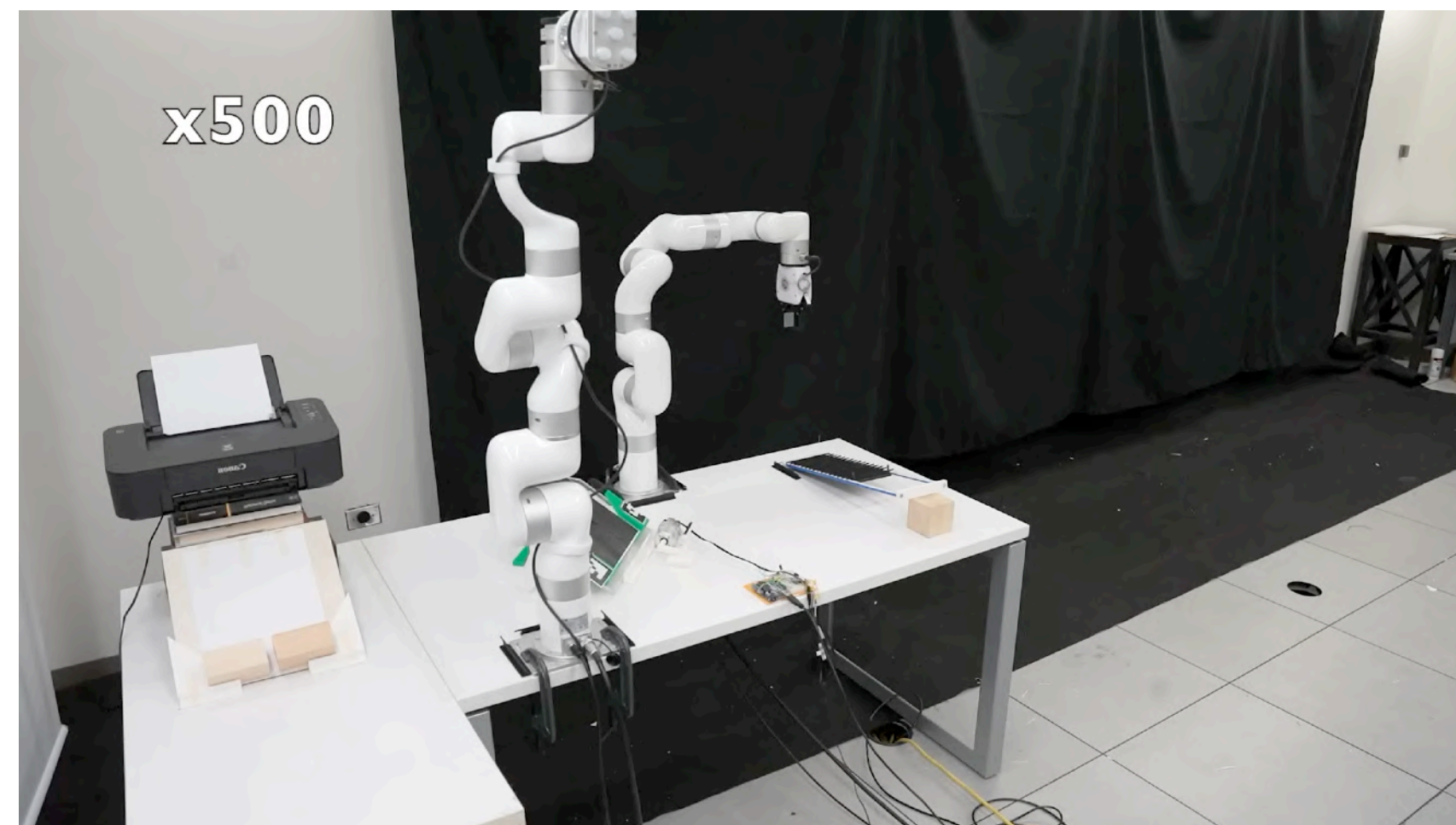
?



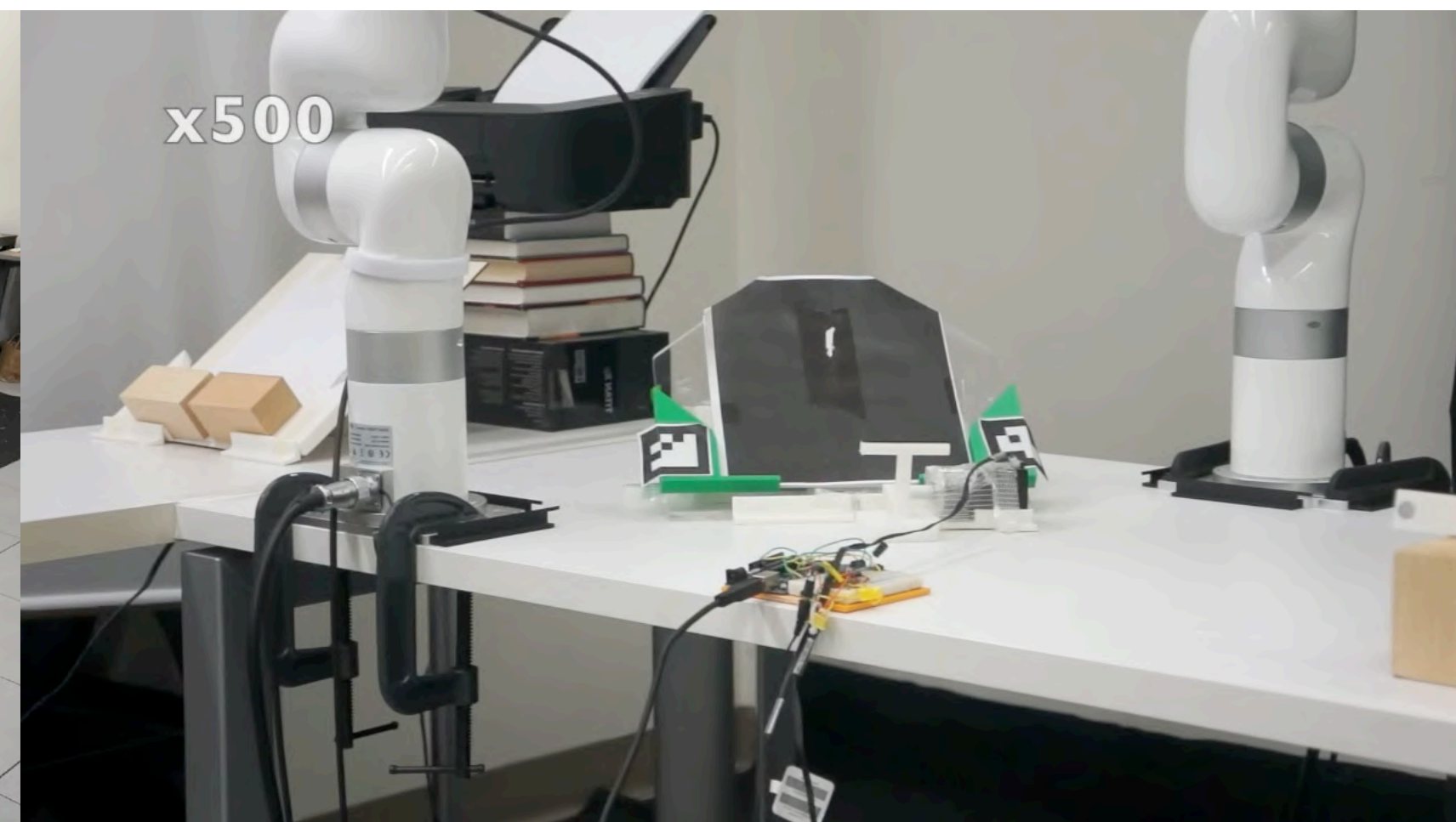
Learning



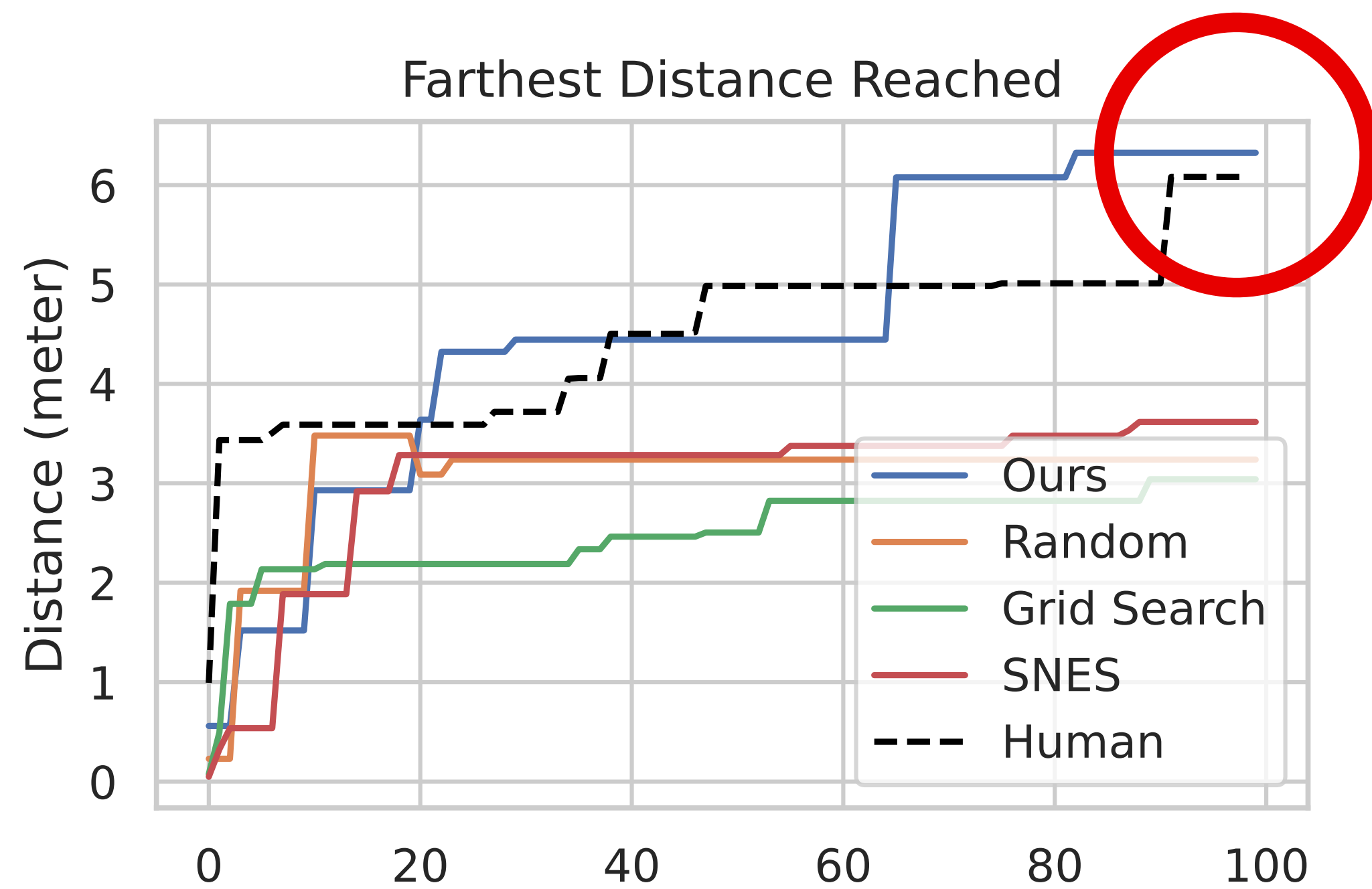
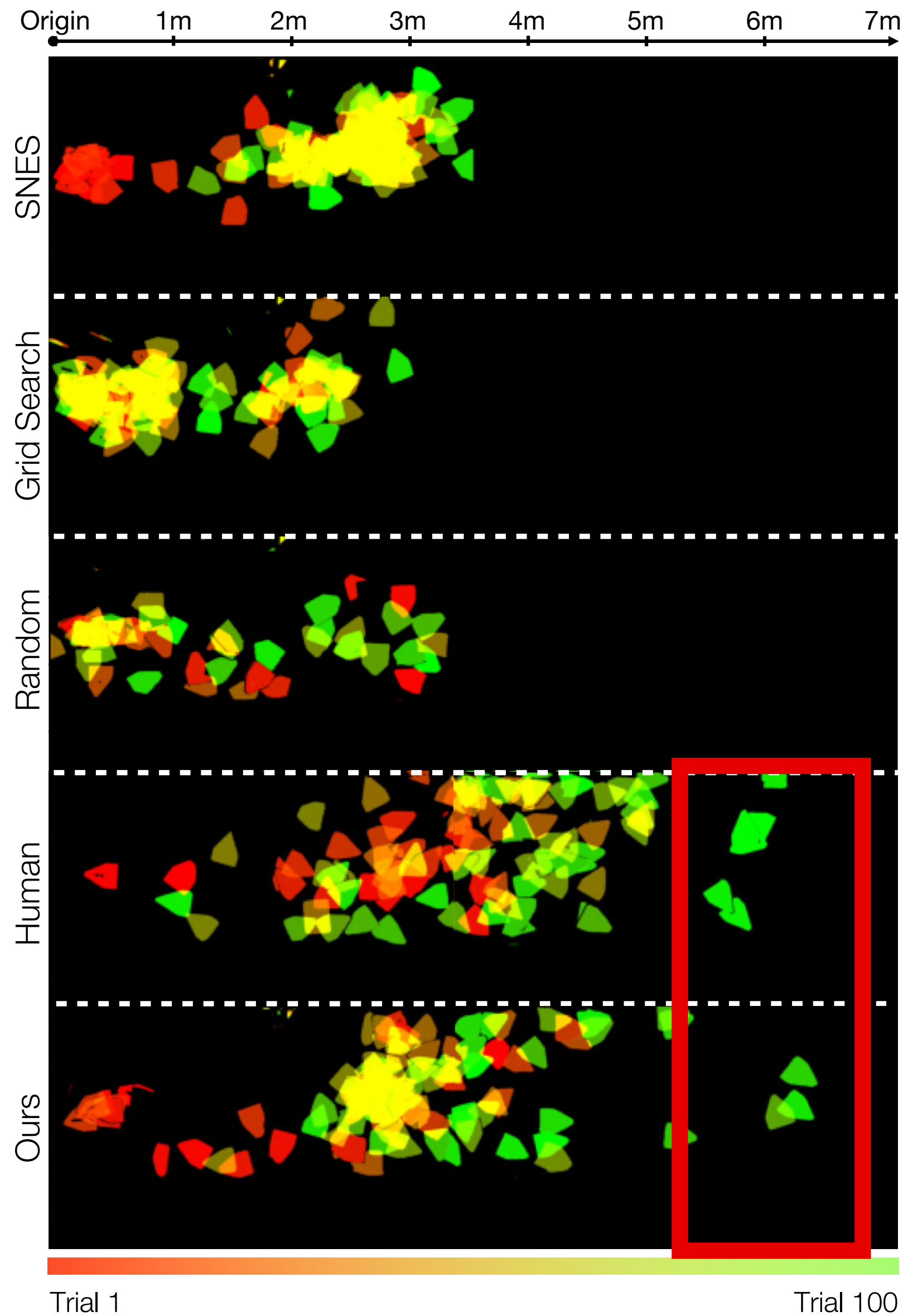
Trial 1



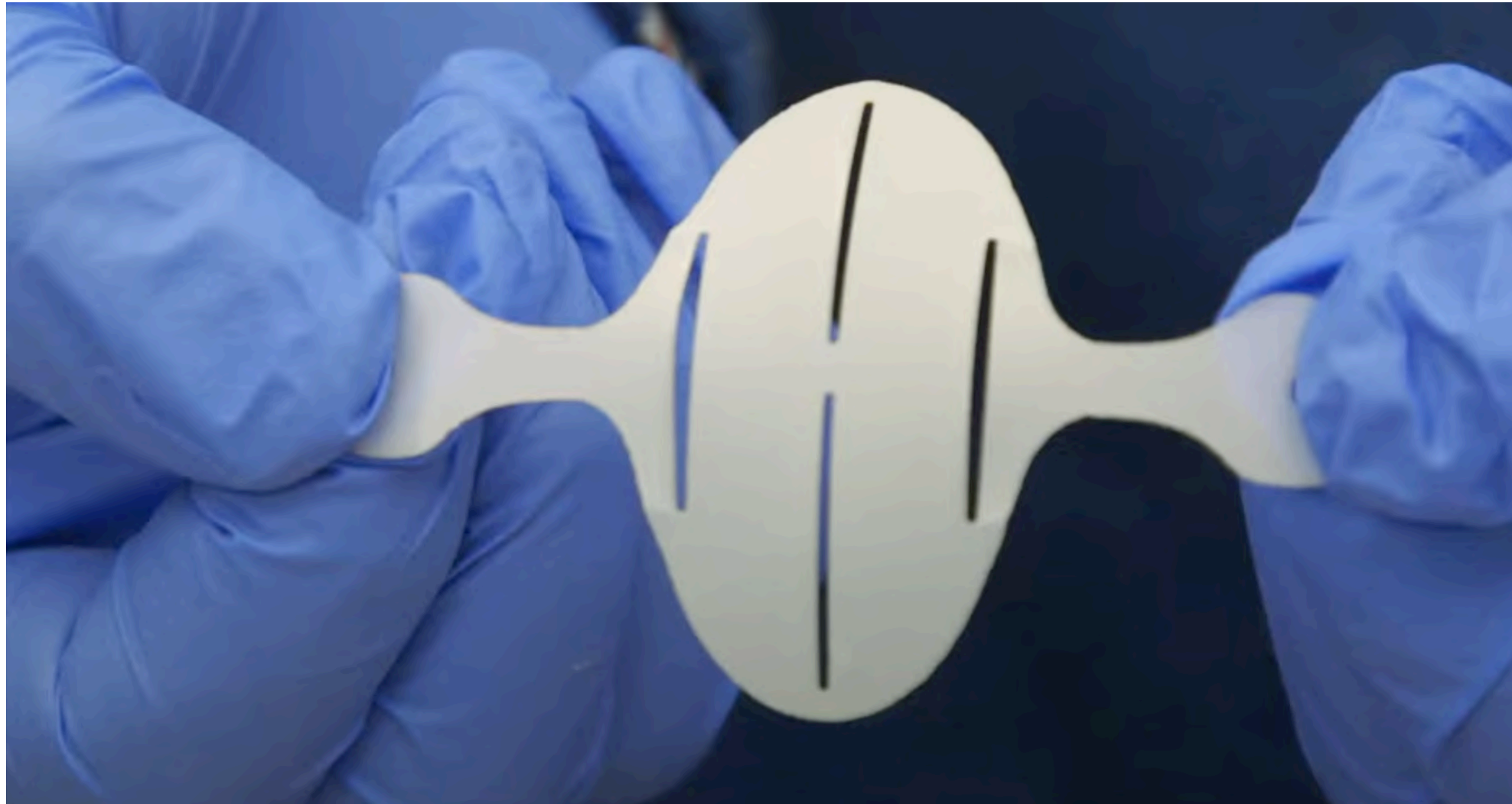
Trial 100



Comparison Against Baselines

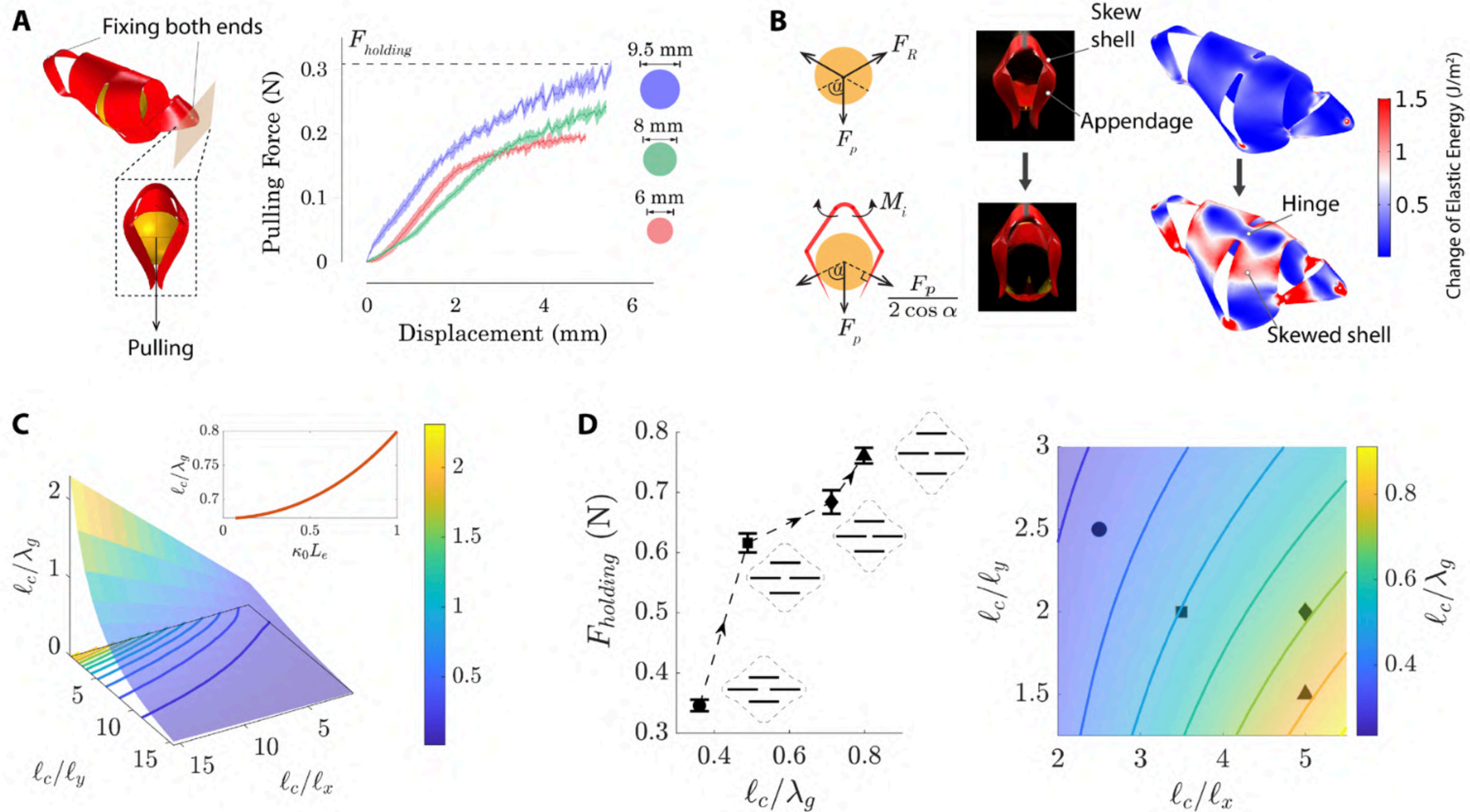


Case Study 2: Kirigami Gripper



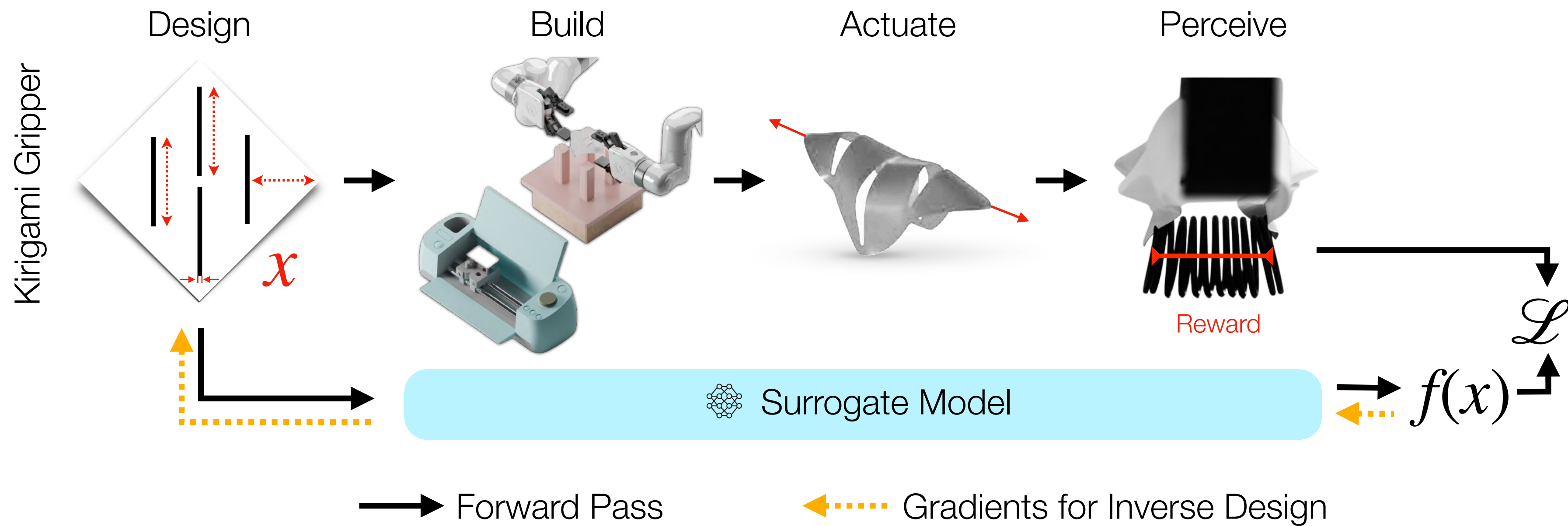
Grasping with kirigami shells. Yi Yang, Katherine Vella, Douglas P. Holmes
Science Robotics 2021

Prior Work Relies on Physical Analysis



Task 2: Kirigami Gripper

Learning to cut paper into grippers that exert maximum gripping force



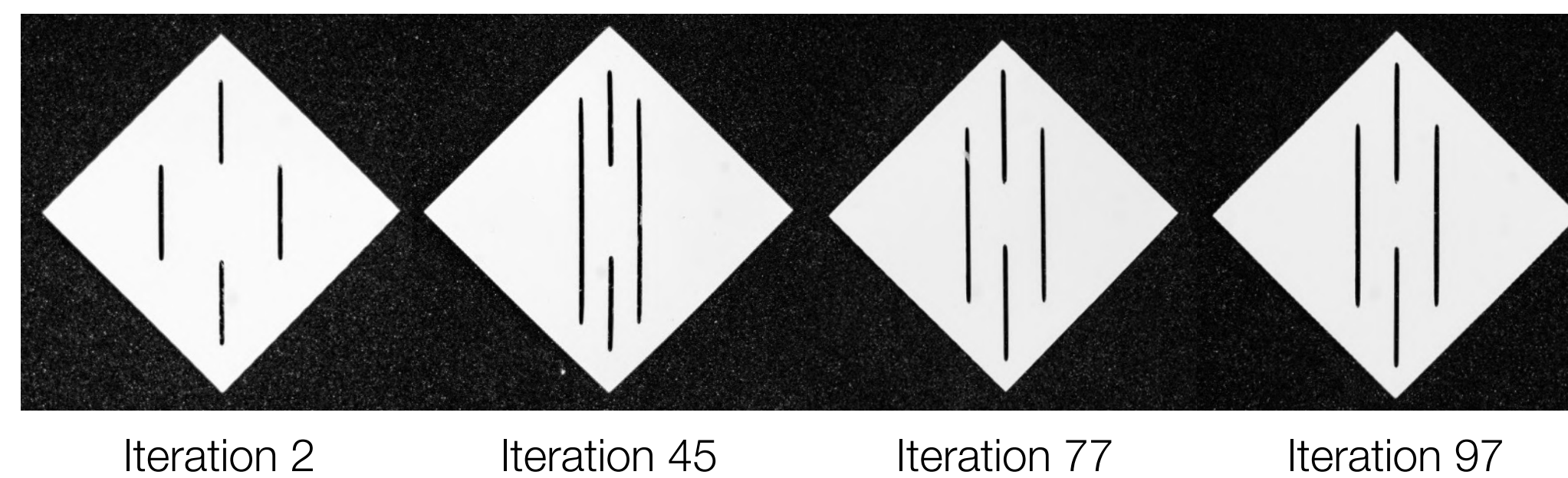
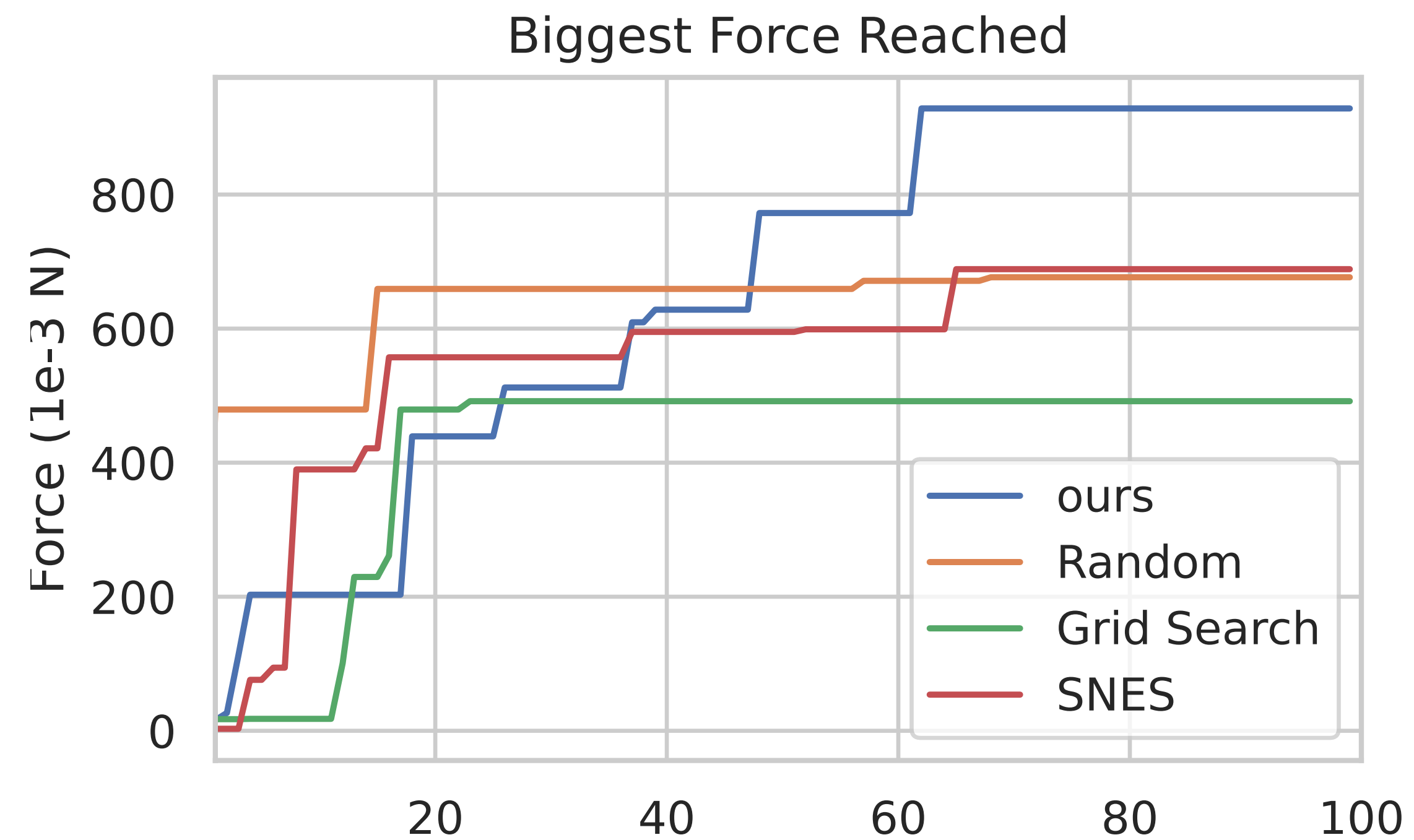
System Setup



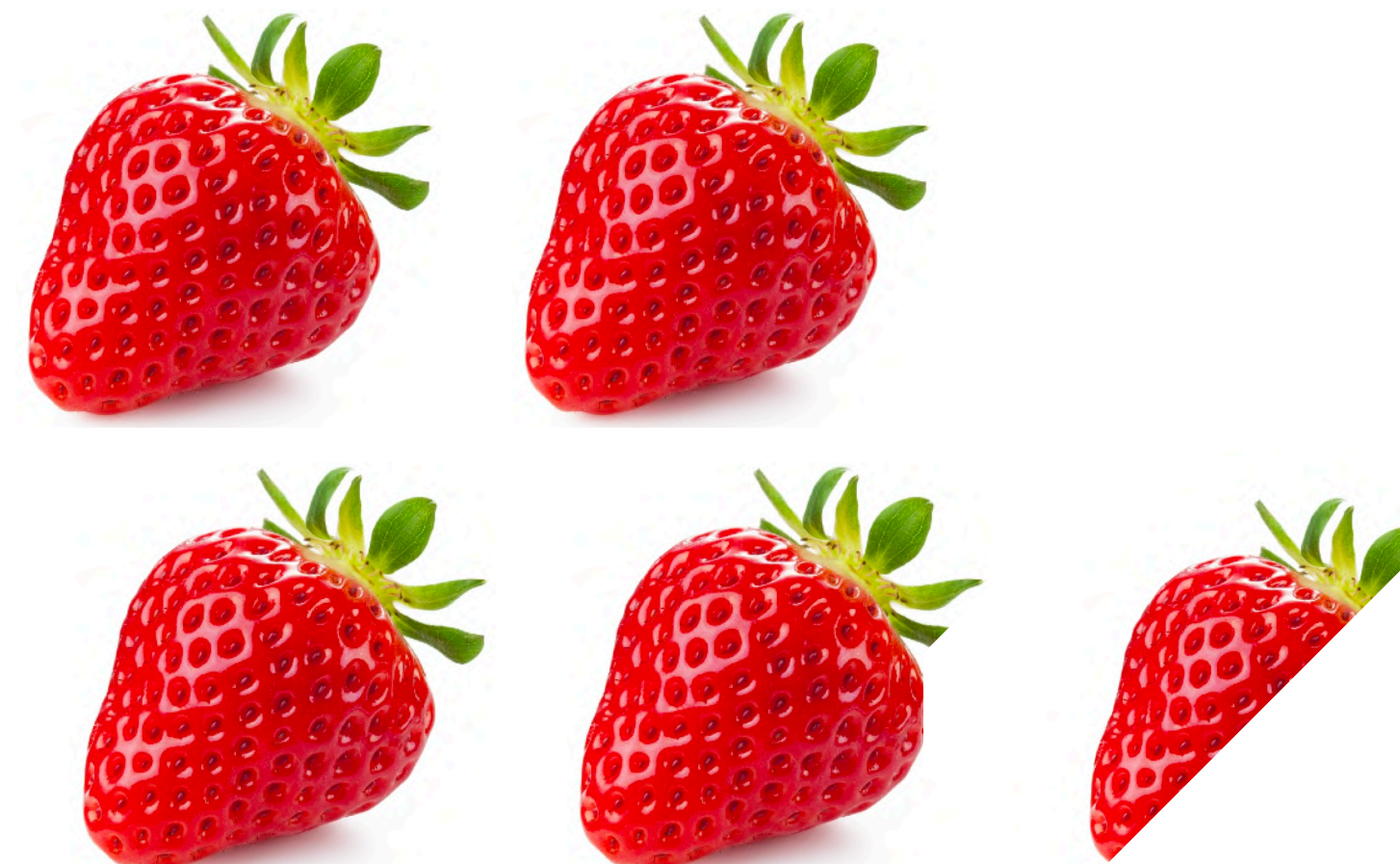
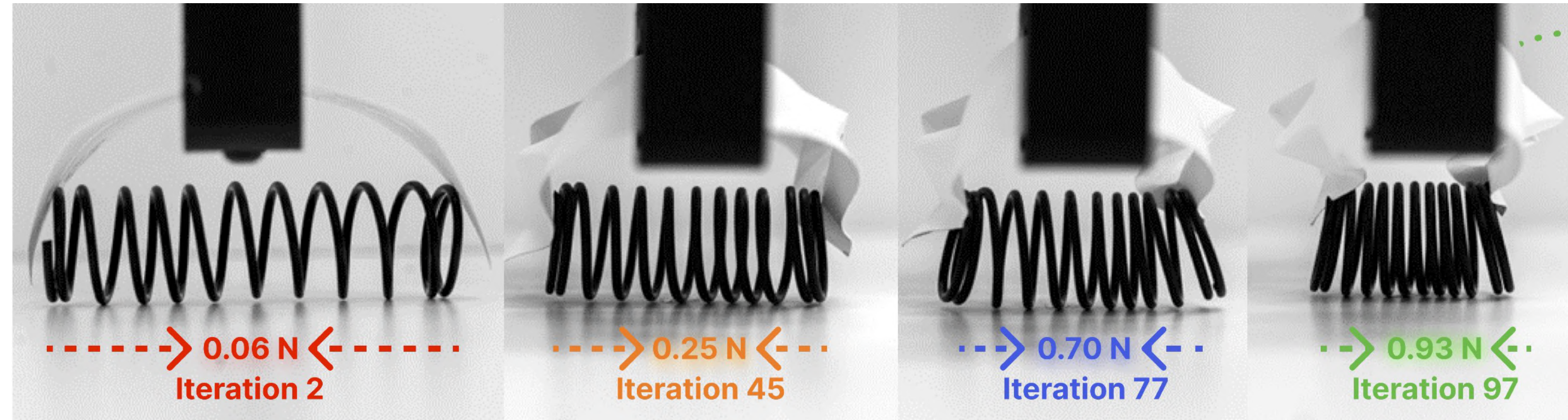
Cricut Maker 3

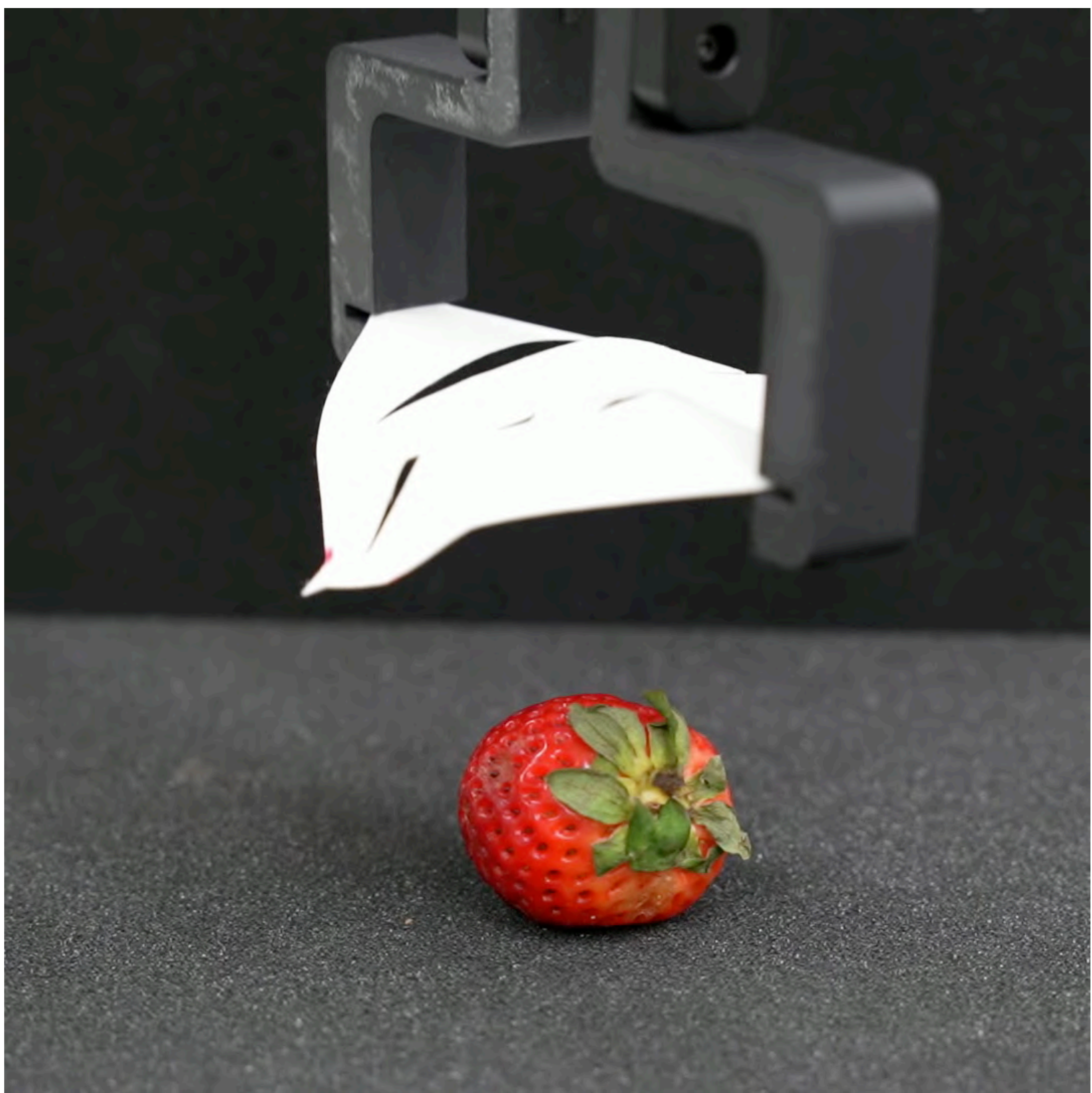
Load Cells

Experiments



Results

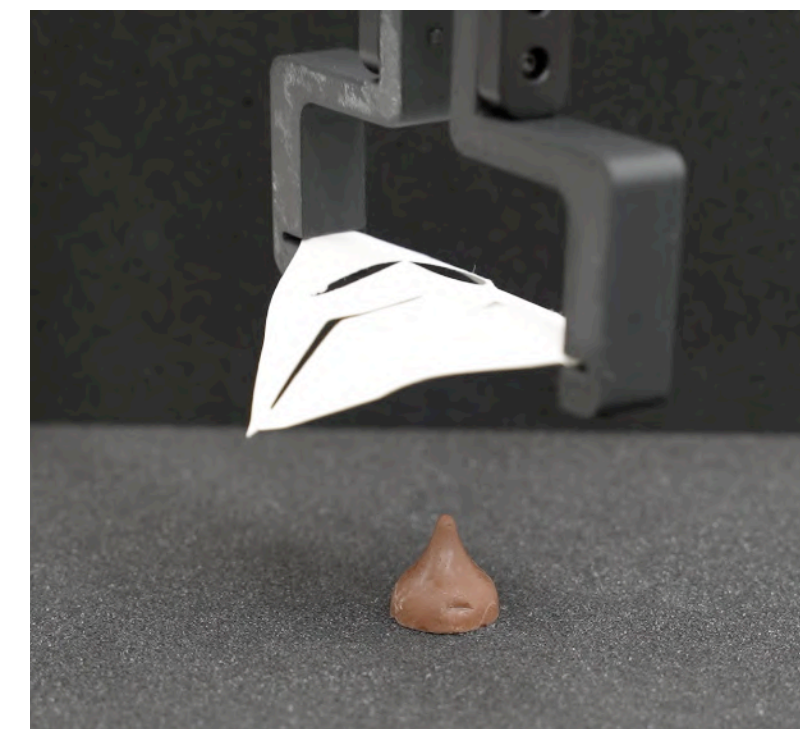
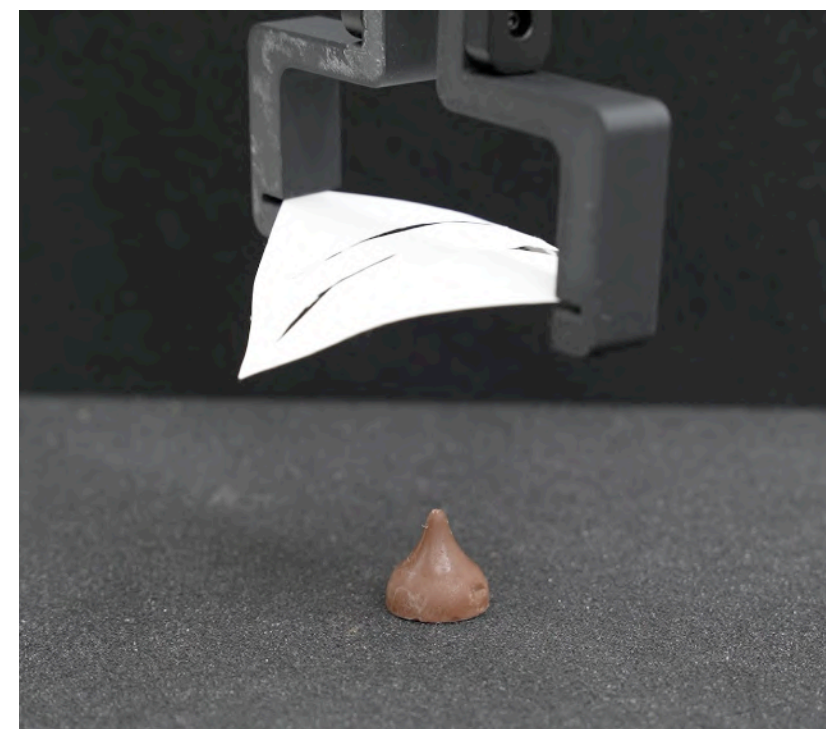
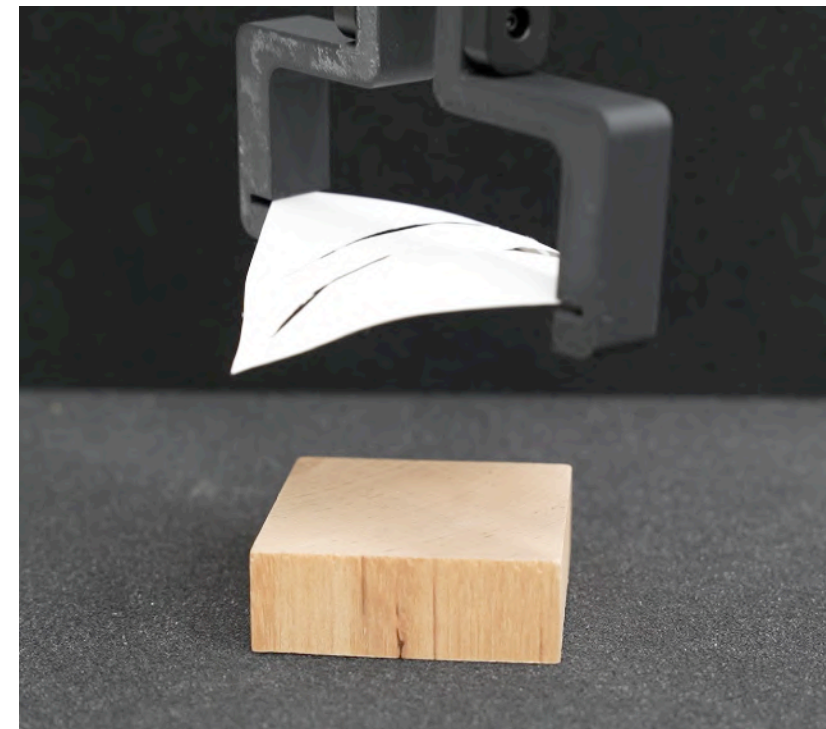






Adaptation

| | Small (1.5cm) | Large (8cm) |
|-------------------------------------|-------------------|-------------------|
| Original (optimized for 5cm) | 0.302 ± 0.012 | 0.055 ± 0.027 |
| Adapted Gripper | 0.442 ± 0.080 | 1.131 ± 0.235 |



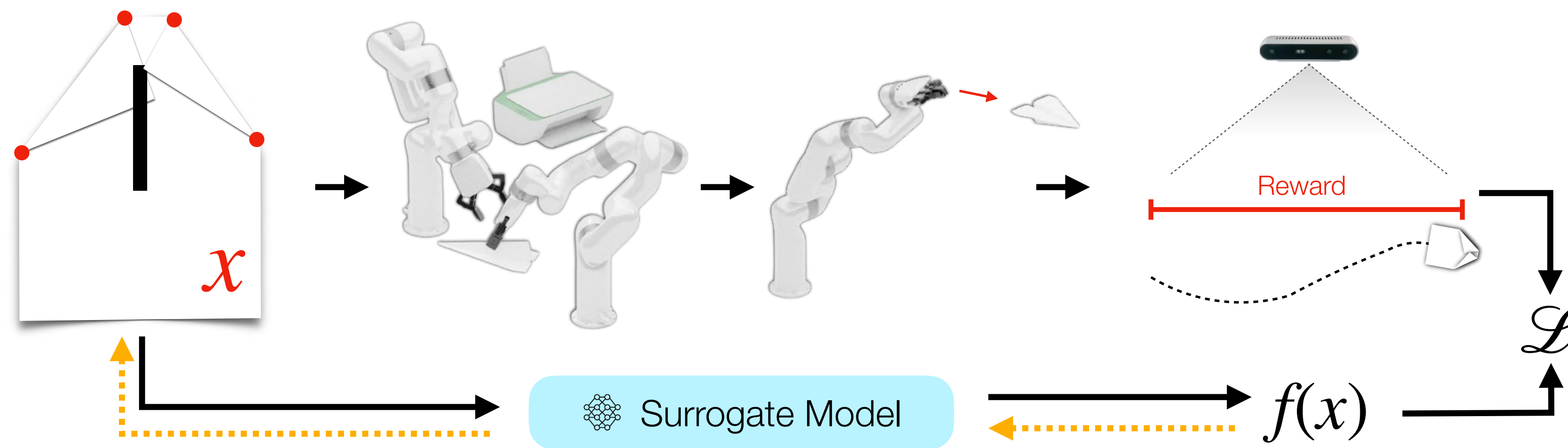
Before Adaptation

After Adaptation



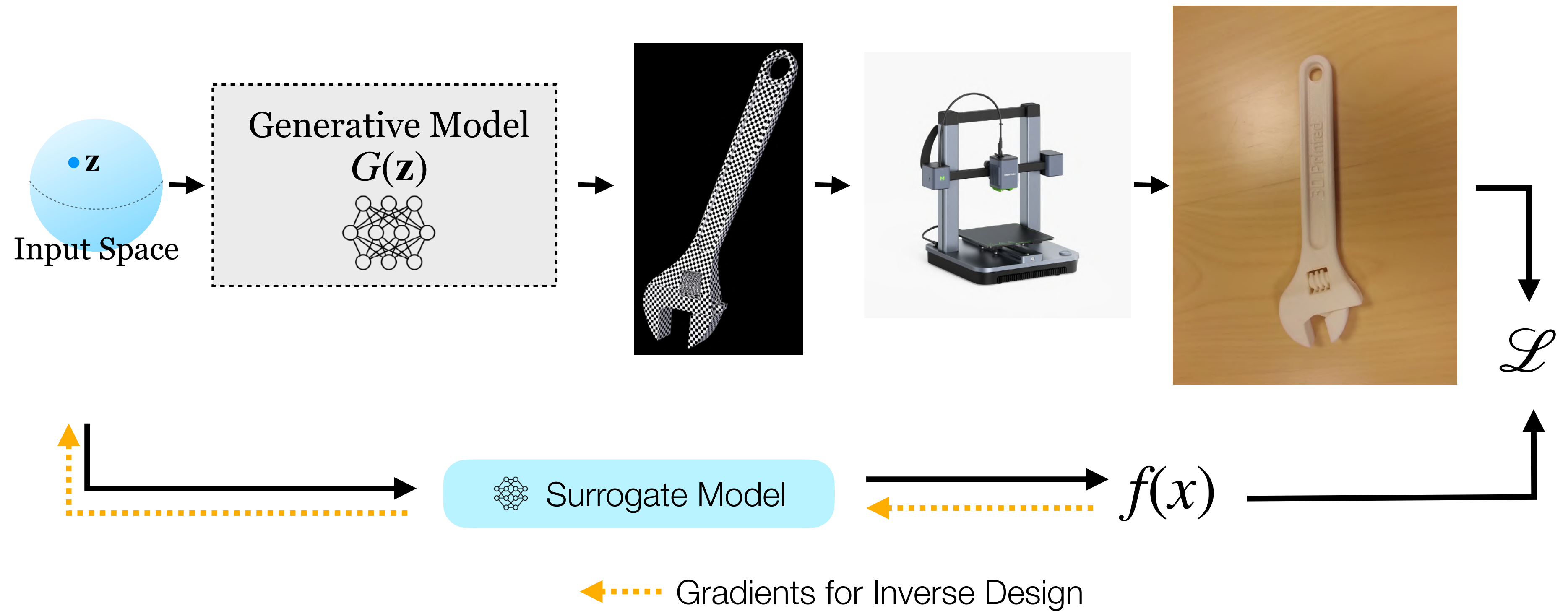
paperbot.cs.columbia.edu



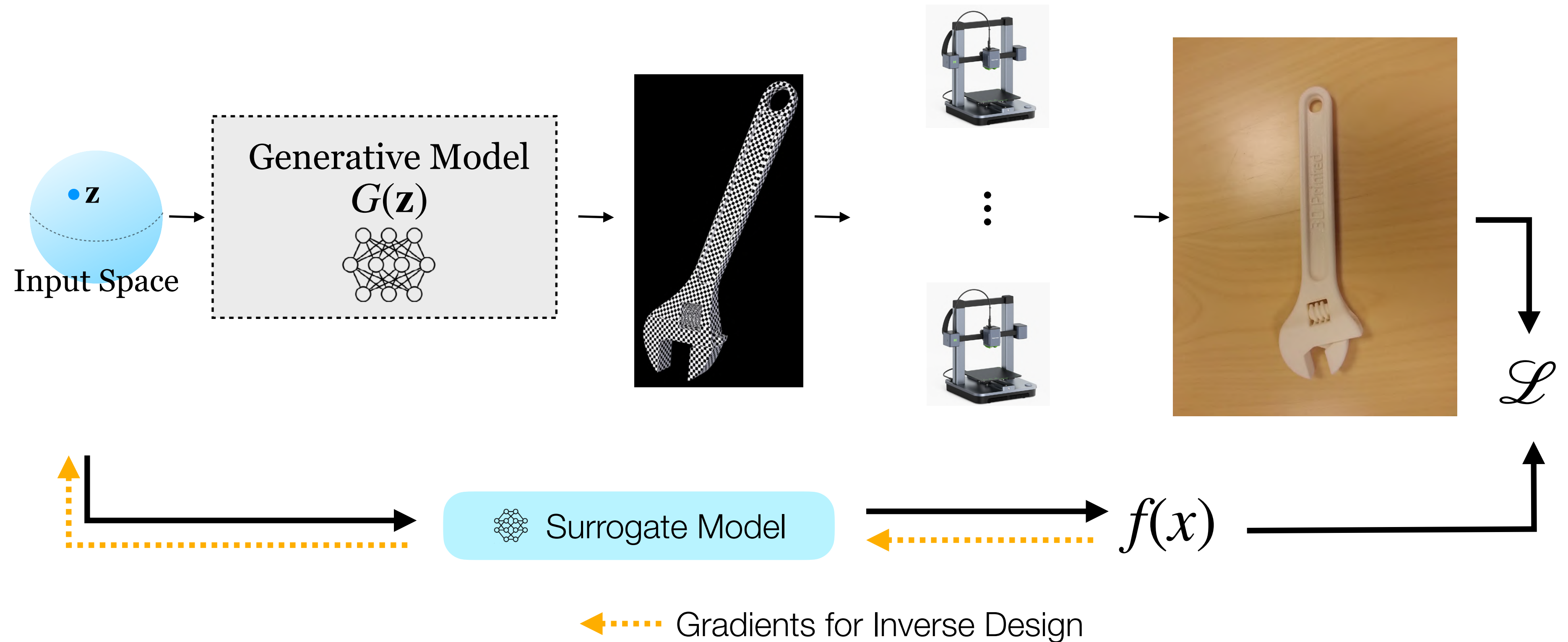


← Gradients for Inverse Design

3D Generation for Tool Design




3D Generation for Tool Design



Generative Embodied AI

3D Generation



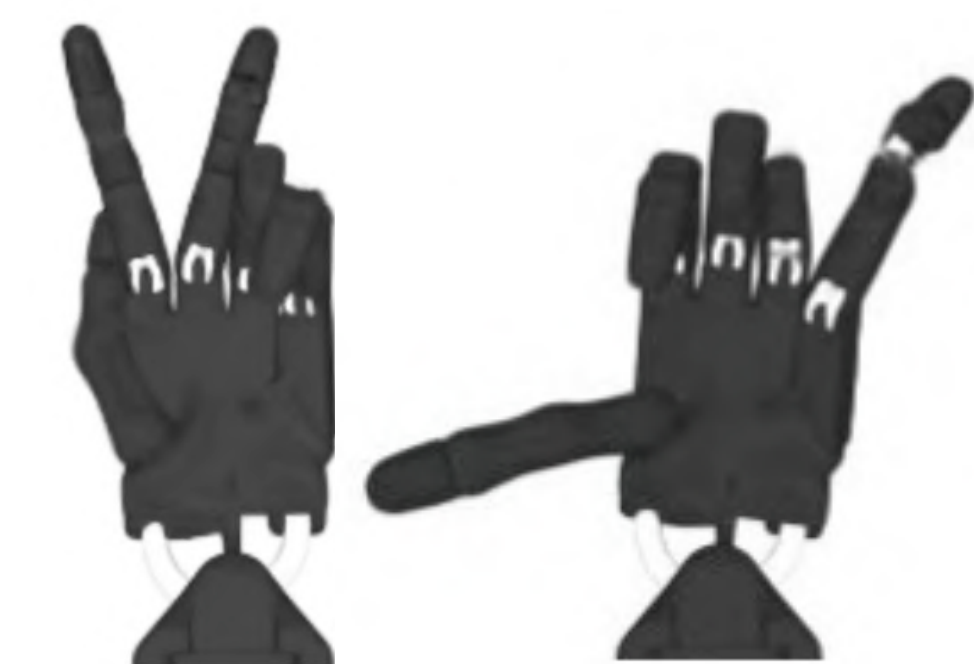
Input Synthesized

Up: 90°

Zero123
ICCV 2023
Liu et al.


Objaverse-XL
NeurIPS 2023
Deitke, Liu et al.

This block illustrates 3D generation. On the left, an 'Input' image of a vase of sunflowers is shown next to its 'Synthesized' 3D model. A blue arrow labeled 'Up: 90°' indicates the camera angle change. On the right, a large collection of diverse 3D models is displayed, including a satellite, a pagoda, a hot air balloon, a dinosaur, a rocket, a car, a house, and various tools and objects.




Dr. Robot
CoRL 2024
Liu*, Canberk* et al.

This block shows a 3D model of a hand, specifically a 'Dr. Robot' hand, in two different poses. The hand is black with white markings on the fingers. Below the images, the text identifies the model and its associated conference and authors.



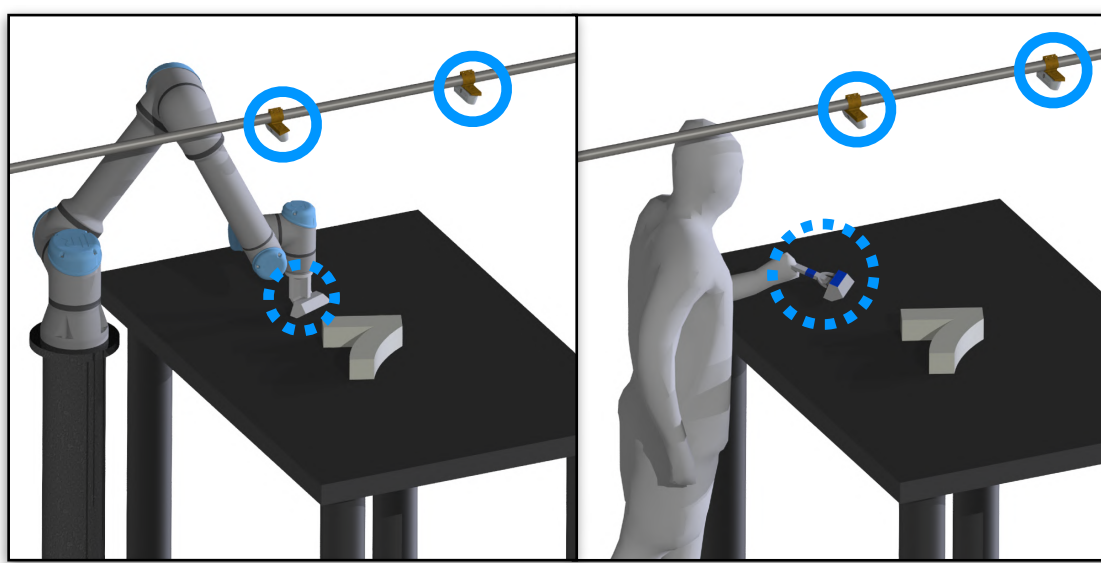
Humans as Light Bulbs
CVPR 2023
Liu et al.

This block shows a person standing in a room with their arms raised, next to a 3D model of a human figure with arms raised. The person is in a room with a camera on a tripod and a shelf with red objects. The 3D model is a blue, featureless human figure. Below the images, the text identifies the model and its associated conference and authors.



PaperBot
RSS 2024
Liu et al.

This block shows a 3D model of a robot arm, specifically a 'PaperBot' arm, in two different poses. The arm is white and has a green sensor on its end. Below the images, the text identifies the model and its associated conference and authors.



Dreamitate
CoRL 2024
Liang*, Liu* et al.

This block shows a 3D model of a robot arm, specifically a 'Dreamitate' arm, in two different poses. The arm is white and is shown interacting with a table and a box. Blue circles highlight the interaction points. Below the images, the text identifies the model and its associated conference and authors.

Physical Reconstruction

Physical Design

Physical Interaction