



University of Hong Kong

Data Collection in DexHand Imitation Learning

Qian Luo

05/12/2024

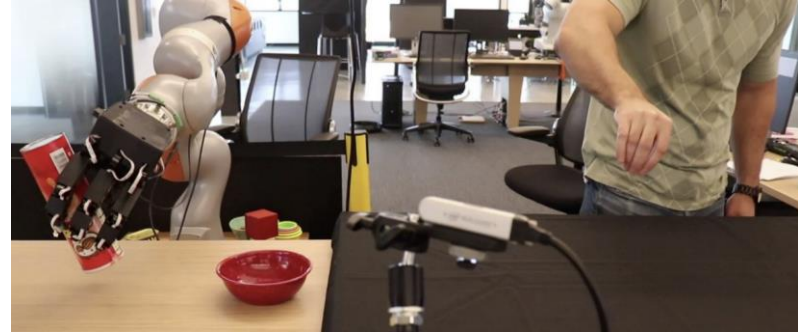
Contents

- Data Collection
 - Vision Teleportation
 - VR Teleportation
 - Motion Capture System
- Imitation Learning
 - 3D diffusion policy

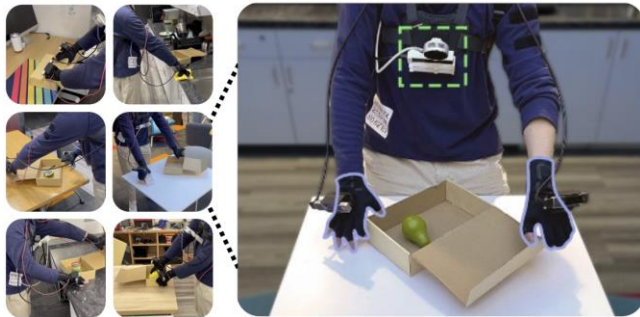
Data Collection

- Three Types
- Difference: The way capturing hand motions

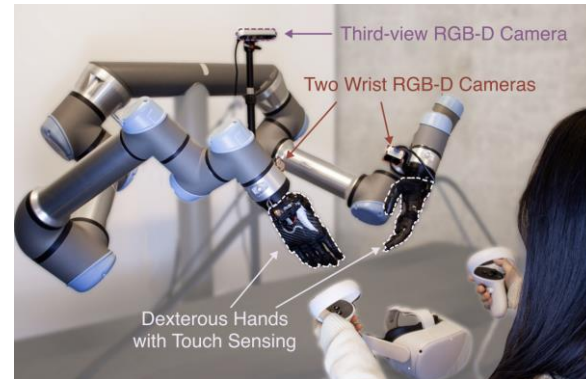
Vision Teleporation



Motion Capture System

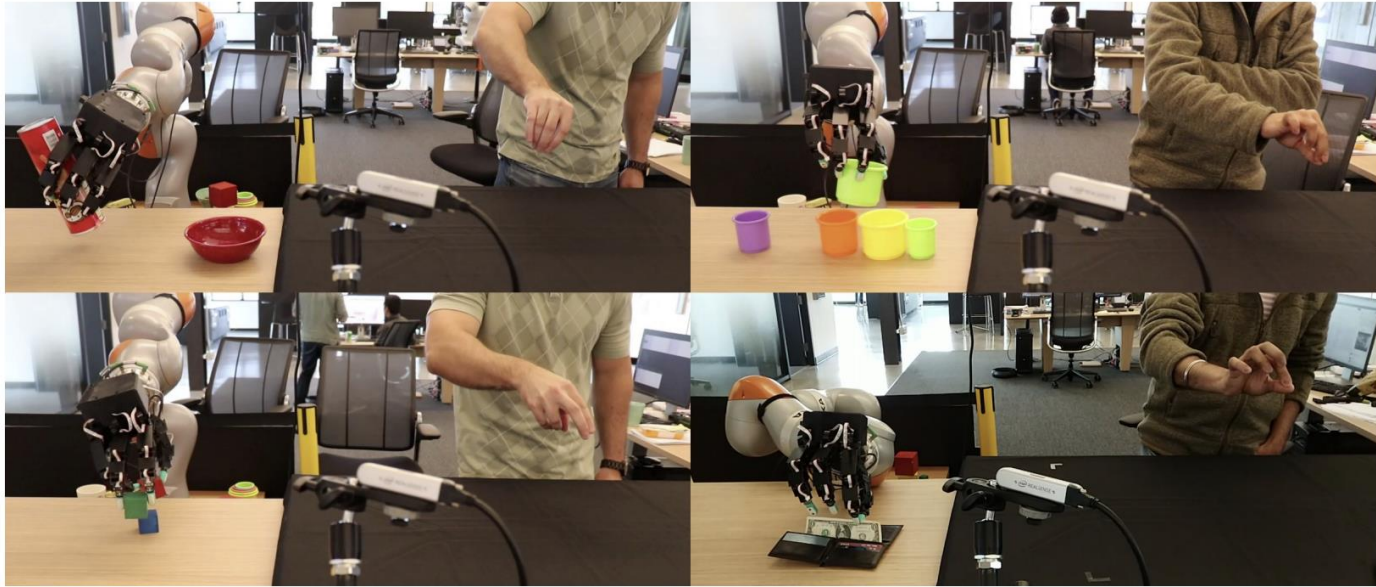


VR Teleporation



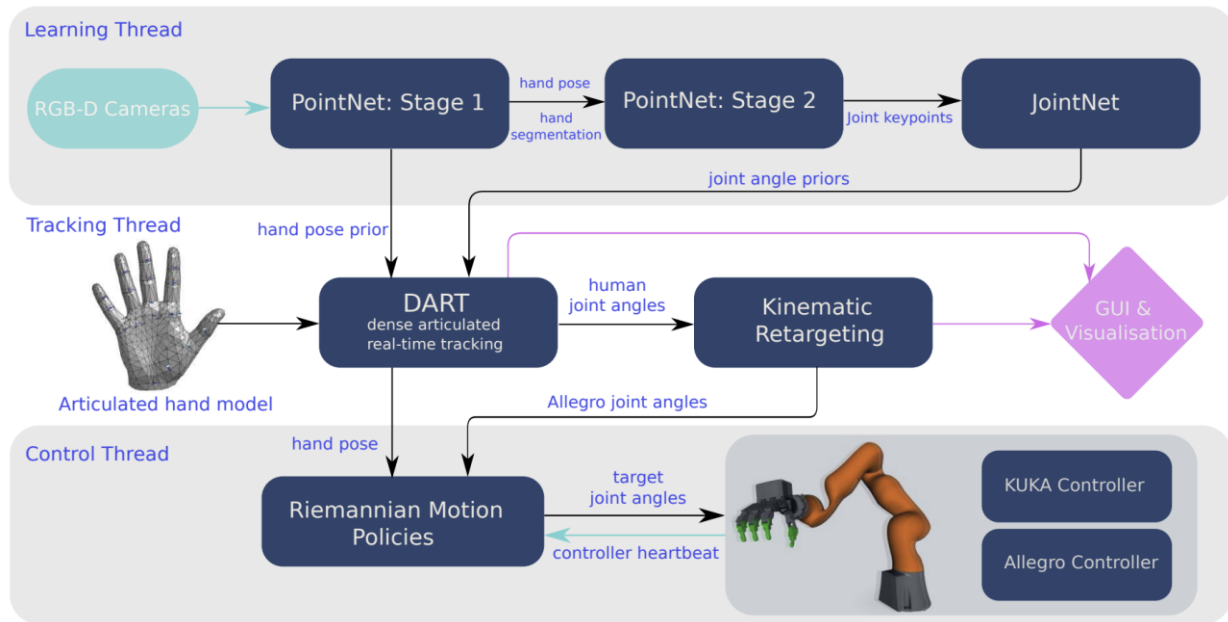
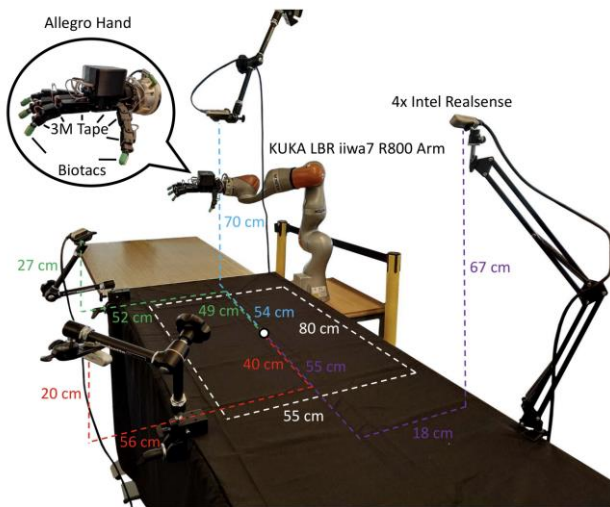
Vision Teleoperation

- Capturing visual information of the human hand using a camera and mapping the hand poses to a robotic hand for teleoperation.
- DexPilot (UW, 2019)



Vision Teleoperation

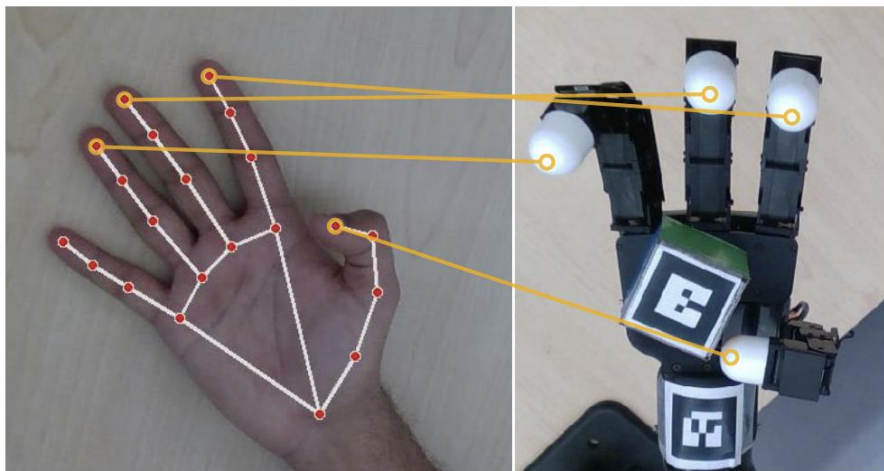
- DexPilot (UW, 2019)



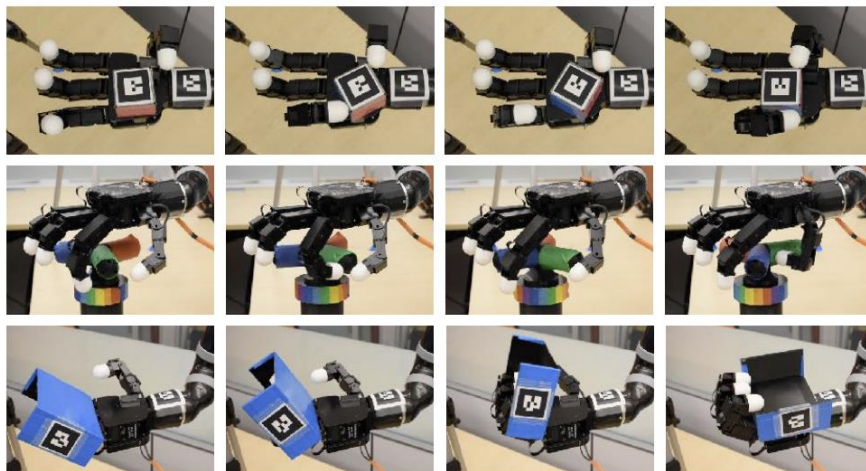
- Derive human hand poses from the 3D point cloud of the hand and re-target them to the robotic hand.

Vision Teleportation

- DIME (NYU, 2022)



(a) Teleoperation through a single RGB camera.



(b) Learned policies for dexterous manipulation.

- Obtain human hand poses from a single RGB image and re-target them to the robotic hand.

Vision Teleoperation

- DIME (NYU, 2022)

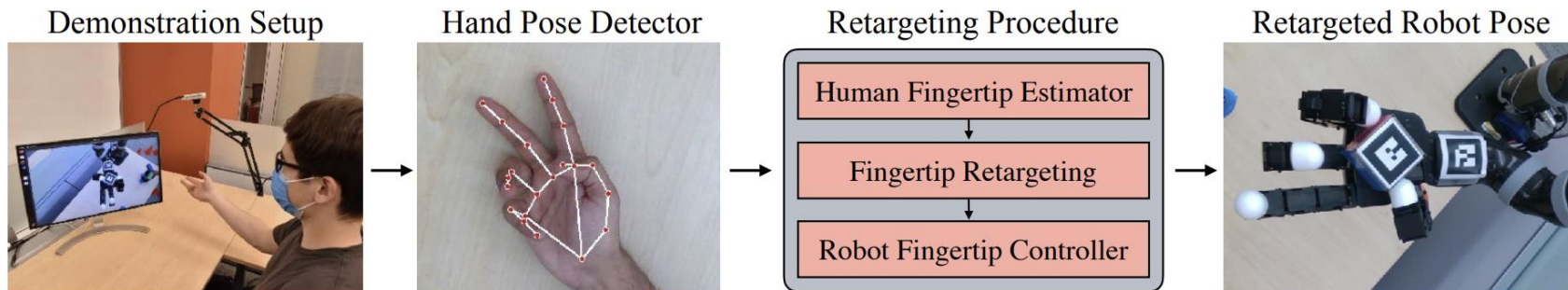
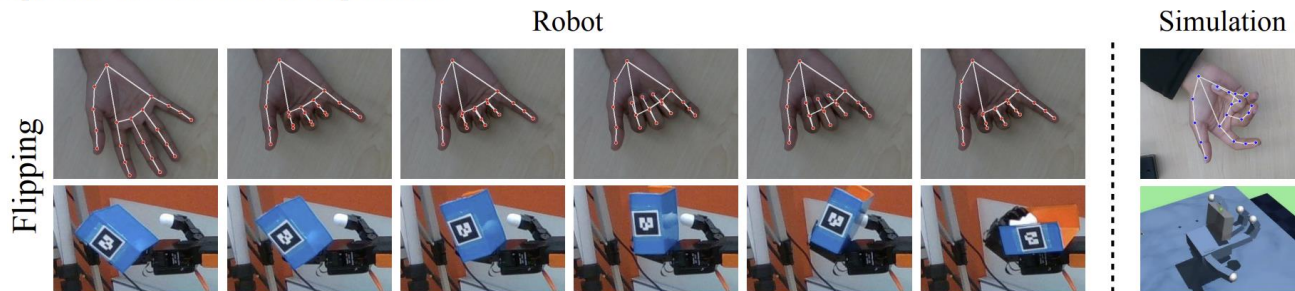
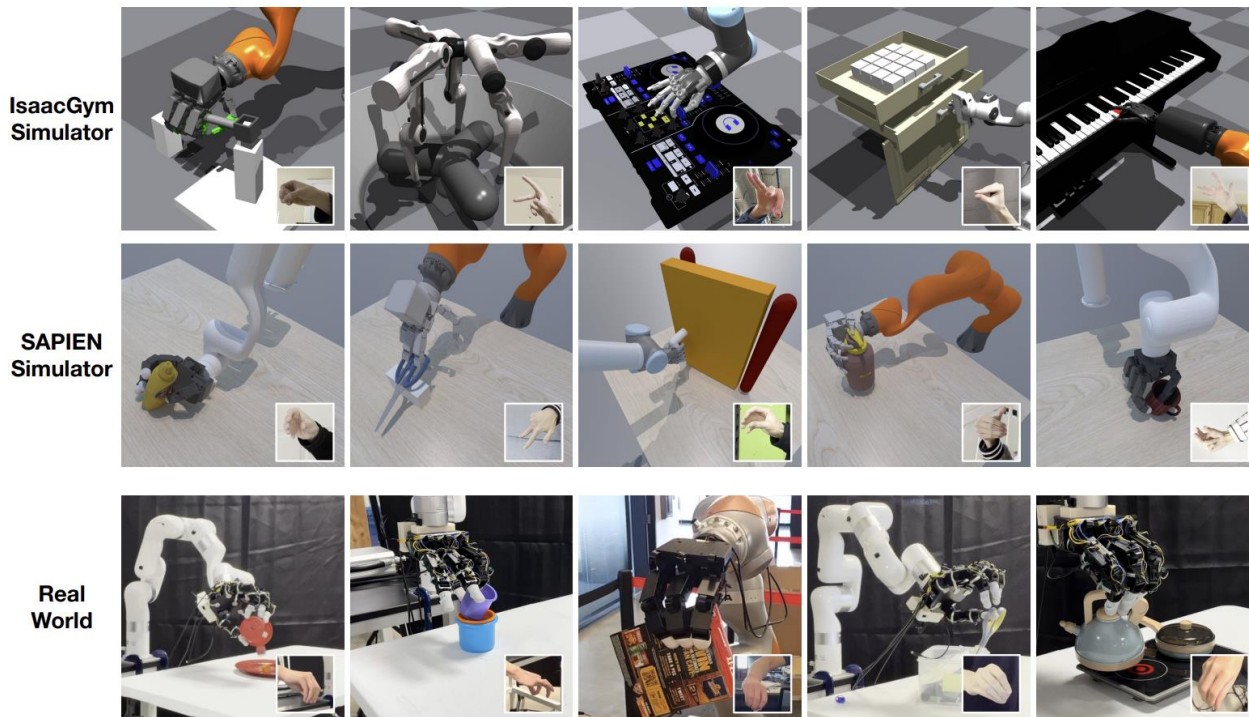


Fig. 2: Overview of the teleoperation framework in DIME. Given RGB streams of a human operator’s hand, a hand pose detector followed by a retargeting procedure is used to control the fingertips of the robot’s hand. Visual feedback of the teleoperated actions is then provided back to the operator for real-time teleoperation.



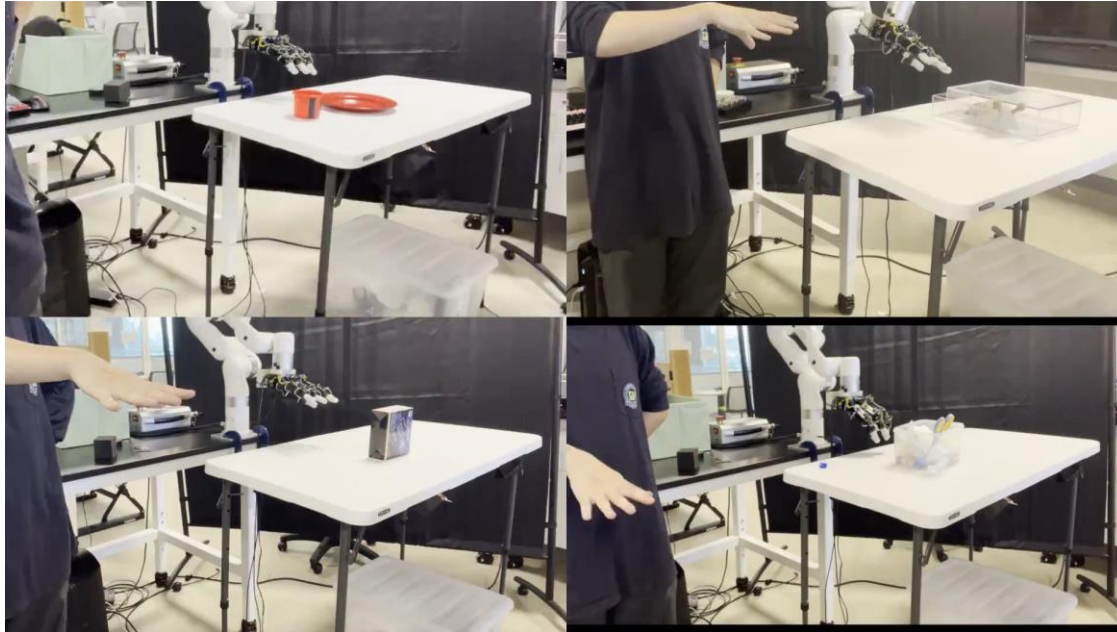
Vision Teleoperation

- AnyTeleop (UCSD, 2023)



Vision Teleporation

- AnyTeleop (UCSD, 2023)



Vision Teleoperation

- AnyTeleop (UCSD, 2023)



Vision Teleoperation

- AnyTeleop (UCSD, 2023)

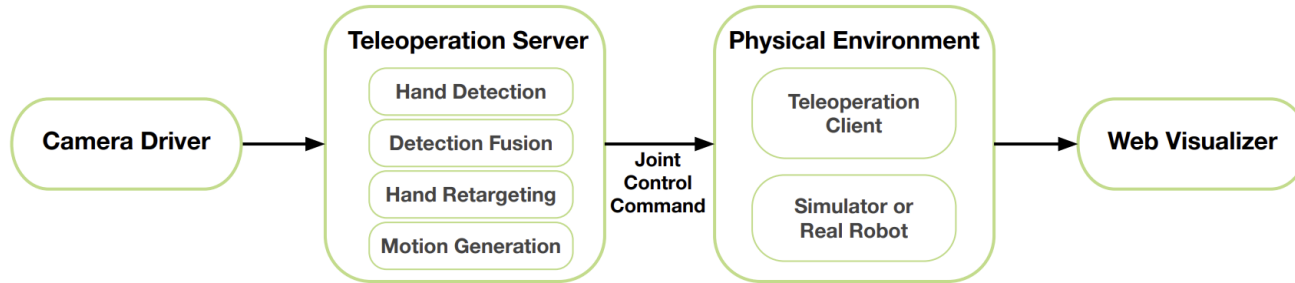


Fig. 3: **System Architecture.** AnyTeleop is composed of four components: (i) camera driver, which captures the human hand pose in RGB or RGB-D format; (ii) teleoperation server, the core component in our system, which performs hand pose detection and converts detection results to robot control commands; (iii) teleoperated robot, which is either a real robot or a simulated robot in a virtual environment; (iv) web visualizer, which enables remote visualization across the internet.

- Human hand pose detection + re-targeting + motion generation.

Vision Teleporation

Pros :

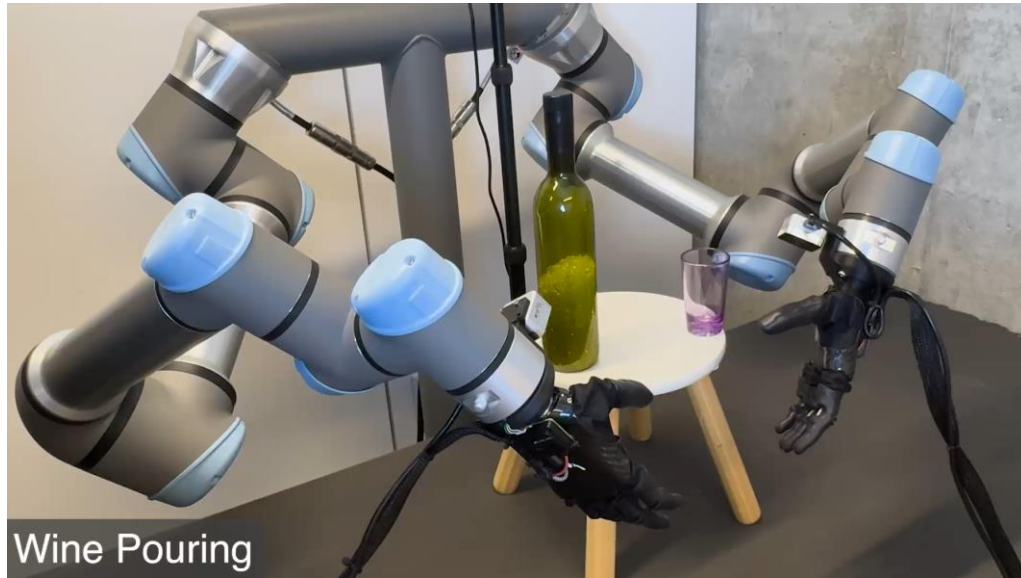
- Portable
- Real-time, low latency
- Realistic(real robot data, easy for policy training)

Cons :

- Relatively low accuracy (the occlusion problem is challenging to resolve).

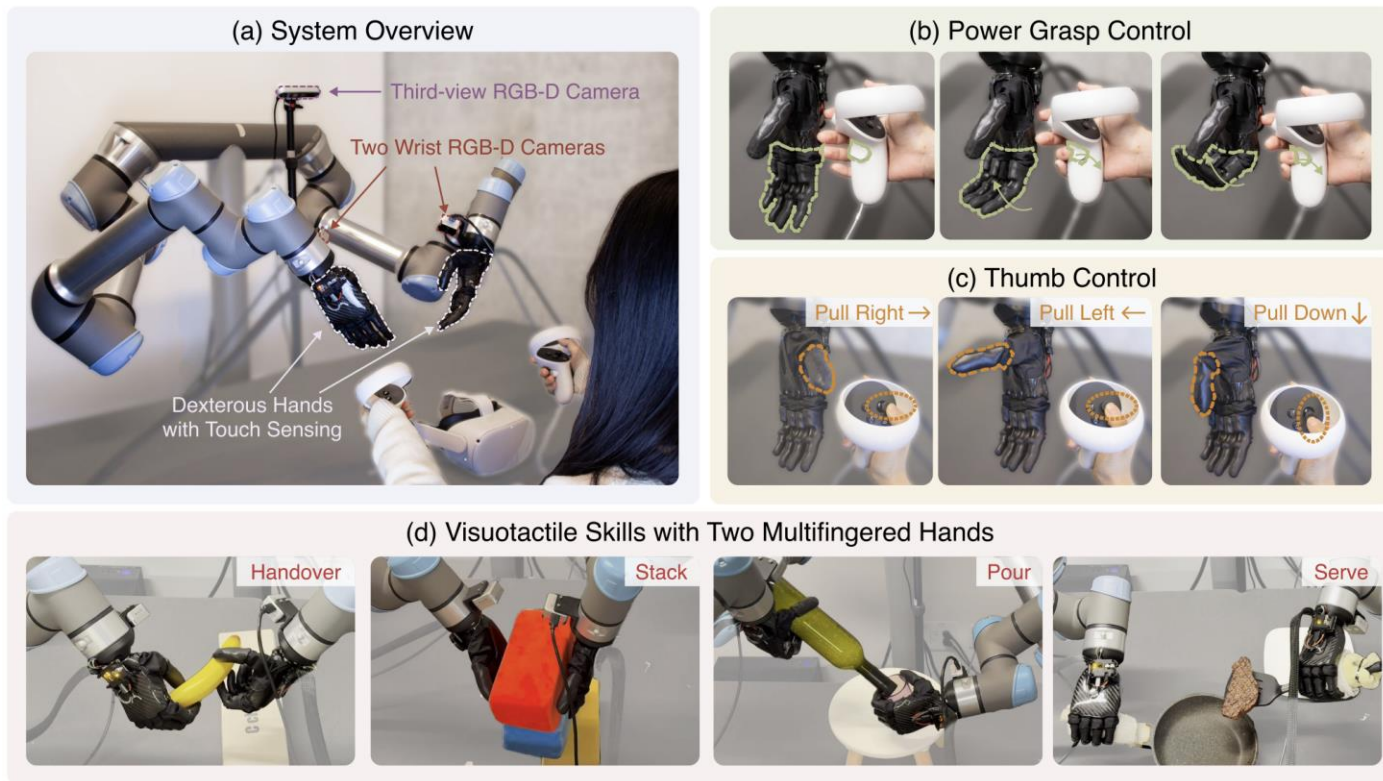
VR Teleporation

- Human teleoperate the robotic hand's movements through a (VR) headset.
- HATO (UCB, 2024)



VR Teleoperation

- HATO (UCB, 2024)



VR Teleportation

Pros :

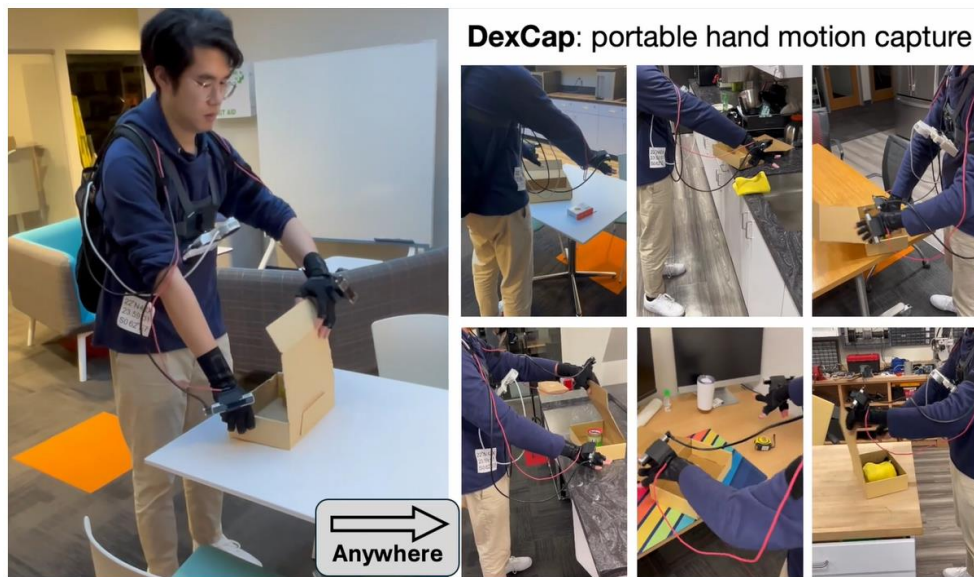
- Immersive teleportation
- Real-time, low latency
- Realistic(real robot data, easy for policy training)

Cons :

- Relatively low accuracy (the occlusion problem is challenging to resolve).

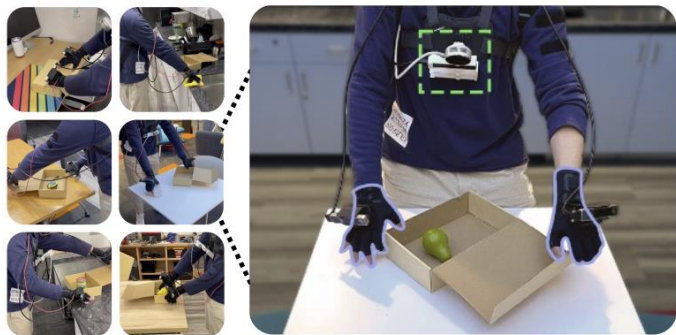
Motion Capture System

- Human hand directly performs the task, capturing hand poses during the process.
- DexCap (Stanford, 2024)

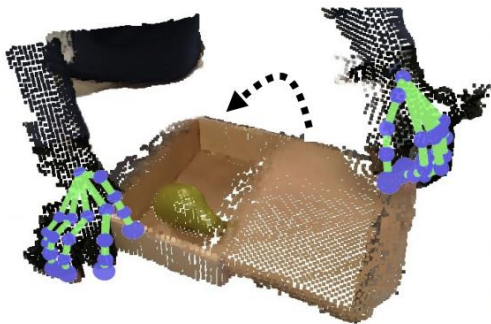


Motion Capture System

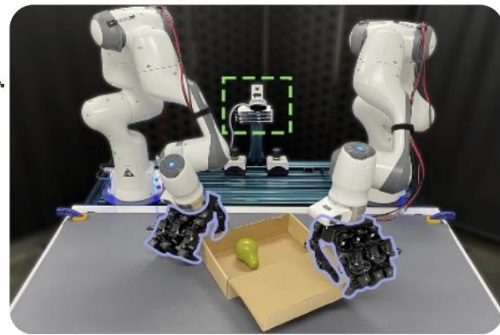
- DexCap (Stanford, 2024)



(a) DexCap: Portable motion capture system



(b) Mocap data and 3D scene



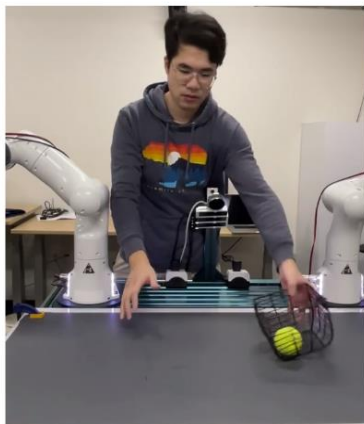
(c) DexIL: Dexterous imitation learning

Fig. 1: **DEXCAP** facilitates the in-the-wild collection of high-quality human hand motion capture data and 3D observations. Leveraging this data, **DEXIL** adapts it to the robot embodiment and trains control policy to perform the same task.

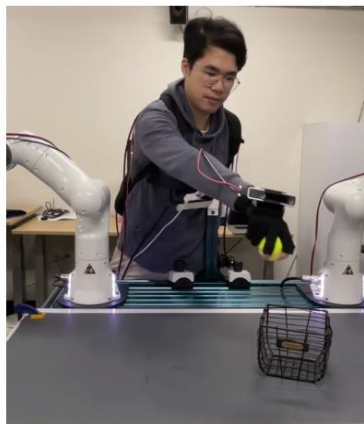
- Unlike teleoperation, DexCap directly uses data collected from the human hand for training and executes tasks on the robotic hand.

Motion Capture System

- DexCap (Stanford, 2024)



Human motion
5.17 s/demo



DexCap
7.75 s/demo



Teleoperation
21.50 s/demo

- Compared to teleoperation, motion capture is more efficient.

Motion Capture System

- DexCap (Stanford, 2024)

Vision-based (VR headset)



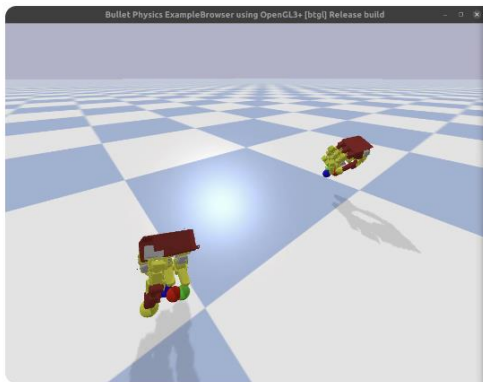
DexCap (ours)



- Compared to vision-based capture, (IMU-based) motion capture offers higher accuracy.

Motion Capture System

- DexCap (Stanford, 2024)



Visual gap: To further bridge the visual gap between human hand and robot hand, we use forward kinematics to generate a point cloud mesh of the robot hand and add it to the pointcloud observation as is shown in this video.

- A visual gap exists. It is necessary to align the human hand and the robotic hand visually.

Motion Capture System

Pros :

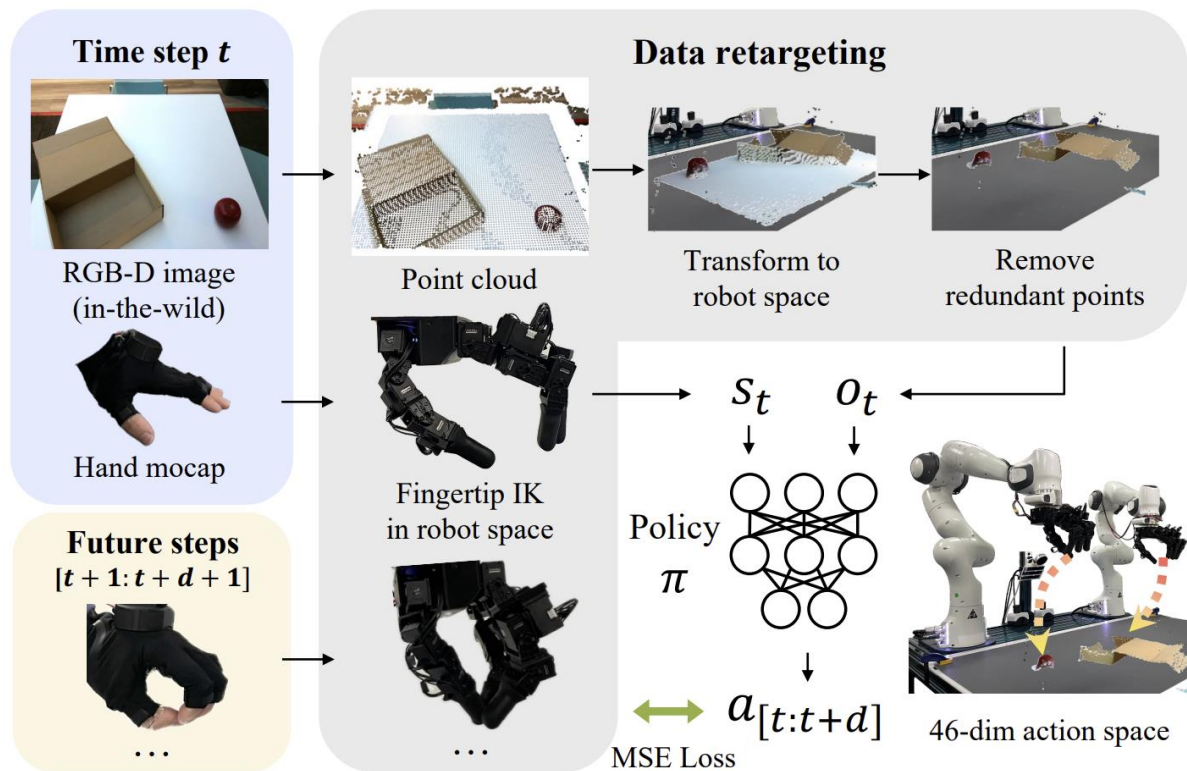
- Highly portable
- No latency
- High accuracy (IMU data)

Cons :

- It is necessary to compensate for the error between the human hand and the robotic hand.

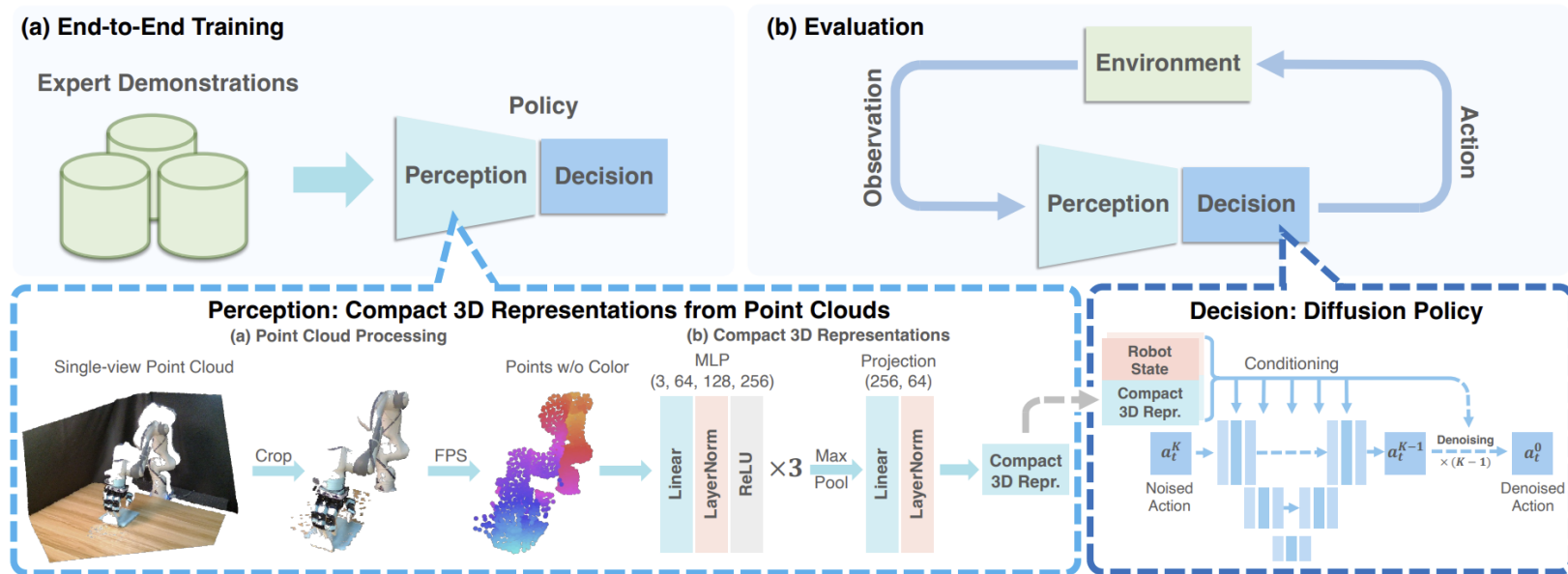
Imitation Learning

- Pipeline



Imitation Learning

- 3D diffusion policy (THU, 2024)



- Perception: 3D point cloud Decision-making: Diffusion policy

Thank you for listening!