# Unsupervised Learning of Visual 3D Keypoints for Control

Boyuan Chen

UC Berkeley

Pieter Abbeel

UC Berkeley

Deepak Pathak
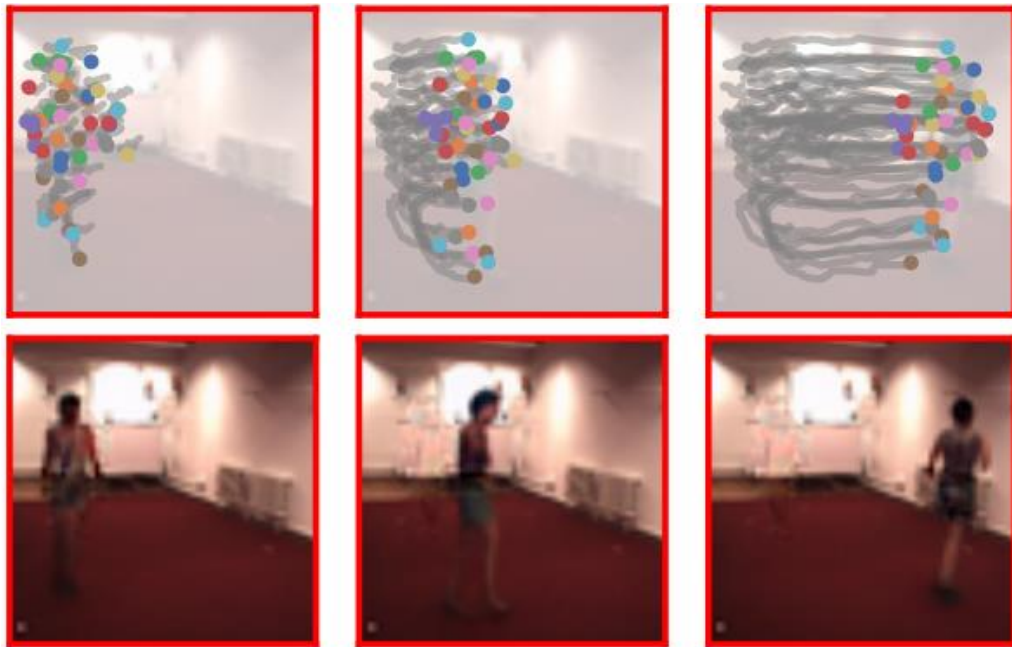
CMU

Learning keypoints from pixels
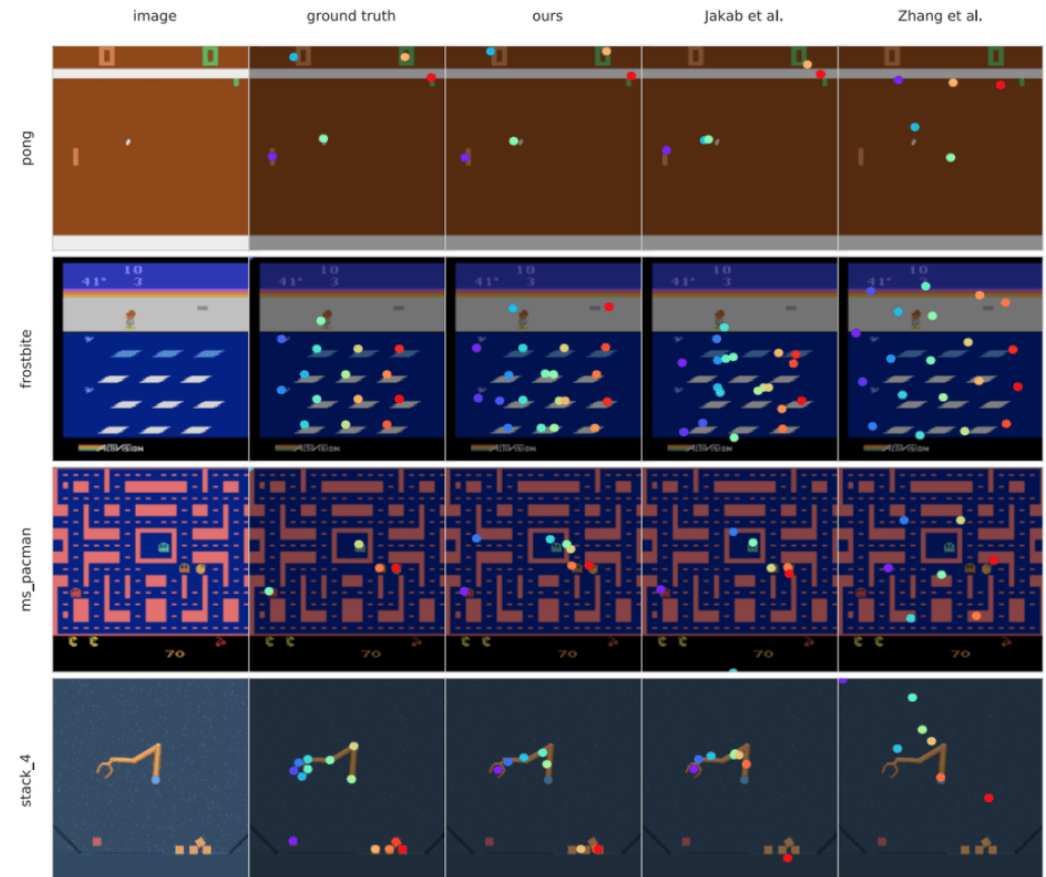


OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, Cao et al., 2018

# Unsupervised 2D Keypoints learning



Unsupervised Learning of Object Structure and Dynamics from Videos, Minderer et al., 2019
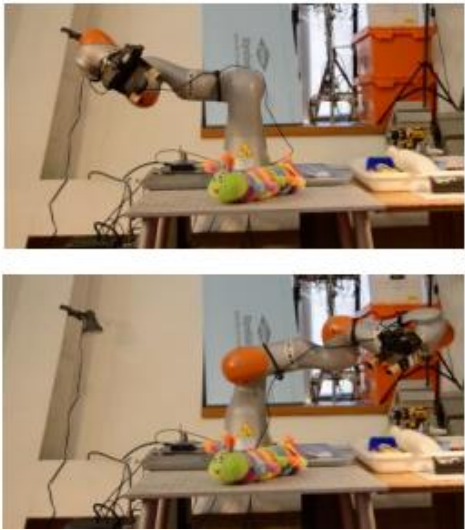
Unsupervised Learning of Object Keypoints for Perception and Control, Kulkarni, Gupta et al., 2019

# Unsupervised 3D Keypoints



Discovery of Latent 3D Keypoints via End-to-end Geometric Reasoning,  Suwajanakorn et al., 2018
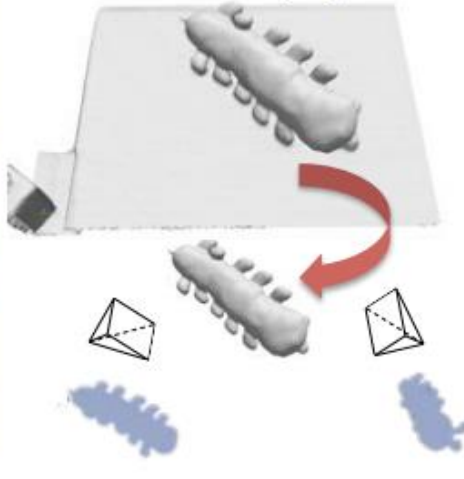
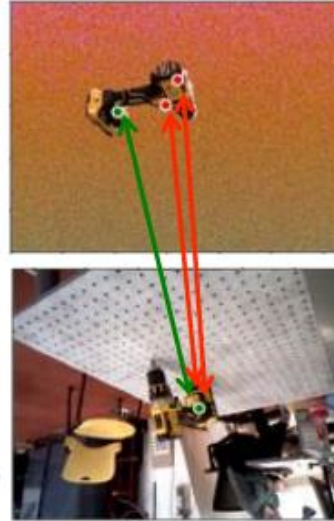# Self-supervised 3D structure learning



(a) Robot-Automated Data Collection
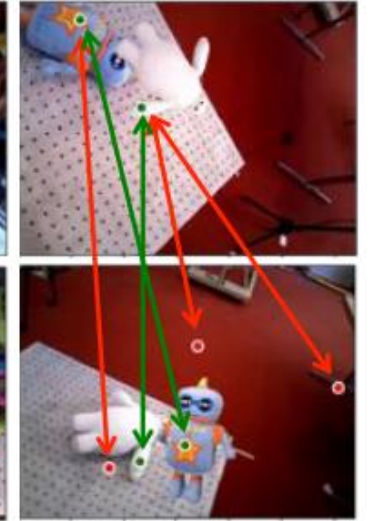(b) 3D Reconstruction based Change Detection and Masked Sampling
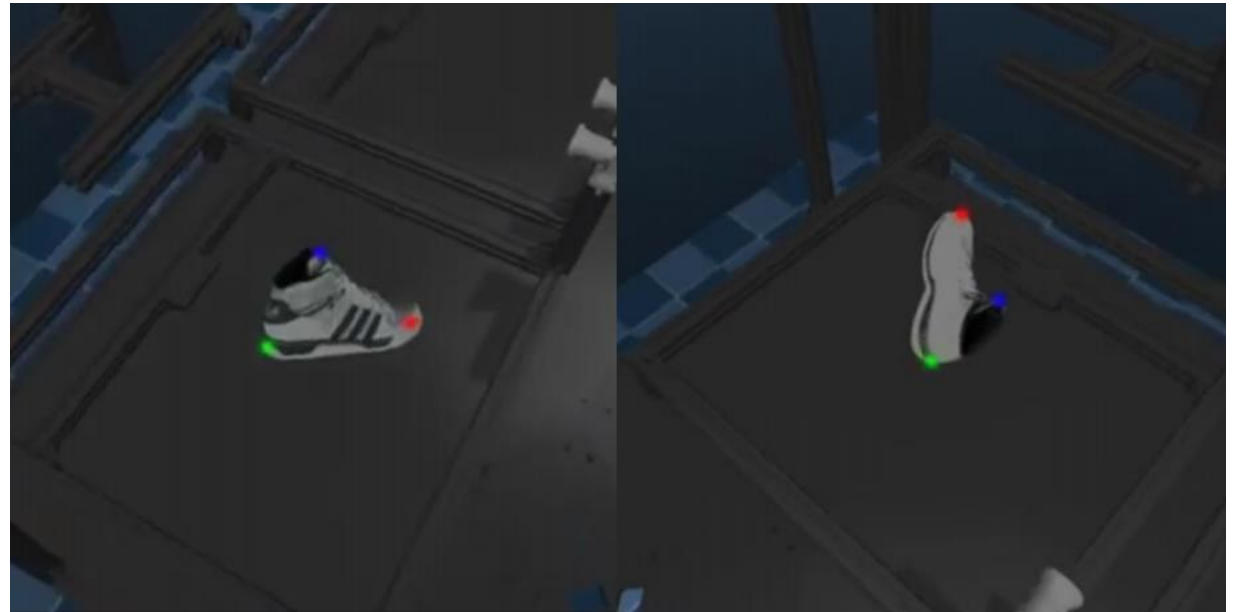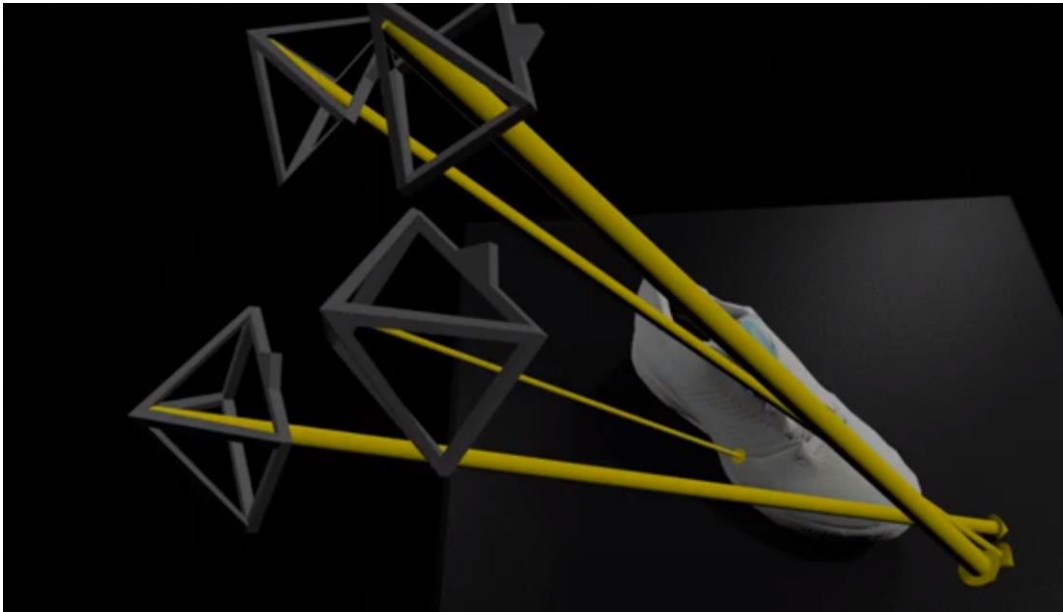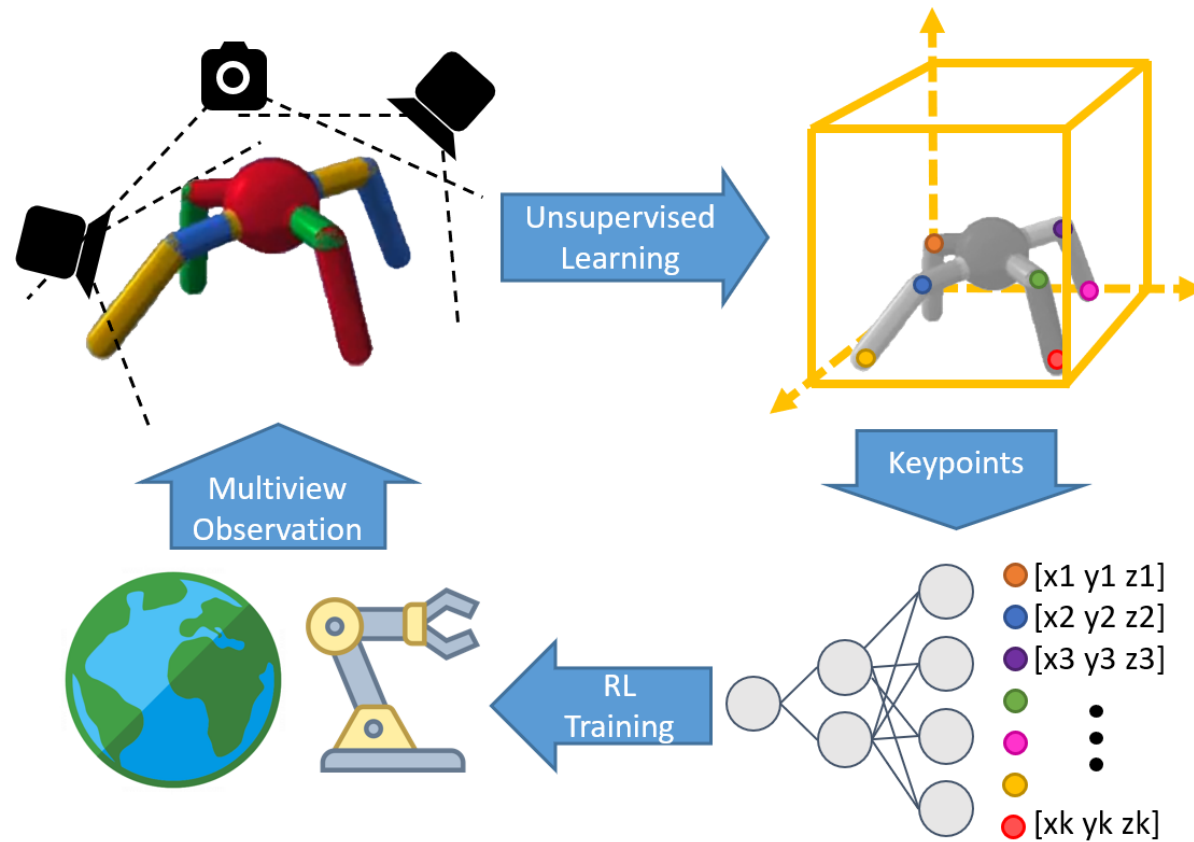(c) Background Randomization
(d) Cross Object Loss
(e) Direct Multi Object
(f) Synthetic Multi Object

Dense Object Nets: Learning Dense Visual Object Descriptors By and For Robotic Manipulation, Florence et al., 2018

# Semi-supervised 3D Keypoints



S3K: Self-Supervised Semantic Keypoints for Robotic Manipulation via Multi-View Consistency, Sushkov et al., 2020

Our work: Keypoint 3D

Unsupervised Learning

Multiview Observation

Keypoints

RL Training

[x1 y1 z1]
[x2 y2 z2]
[x3 y3 z3]

[xk yk zk]
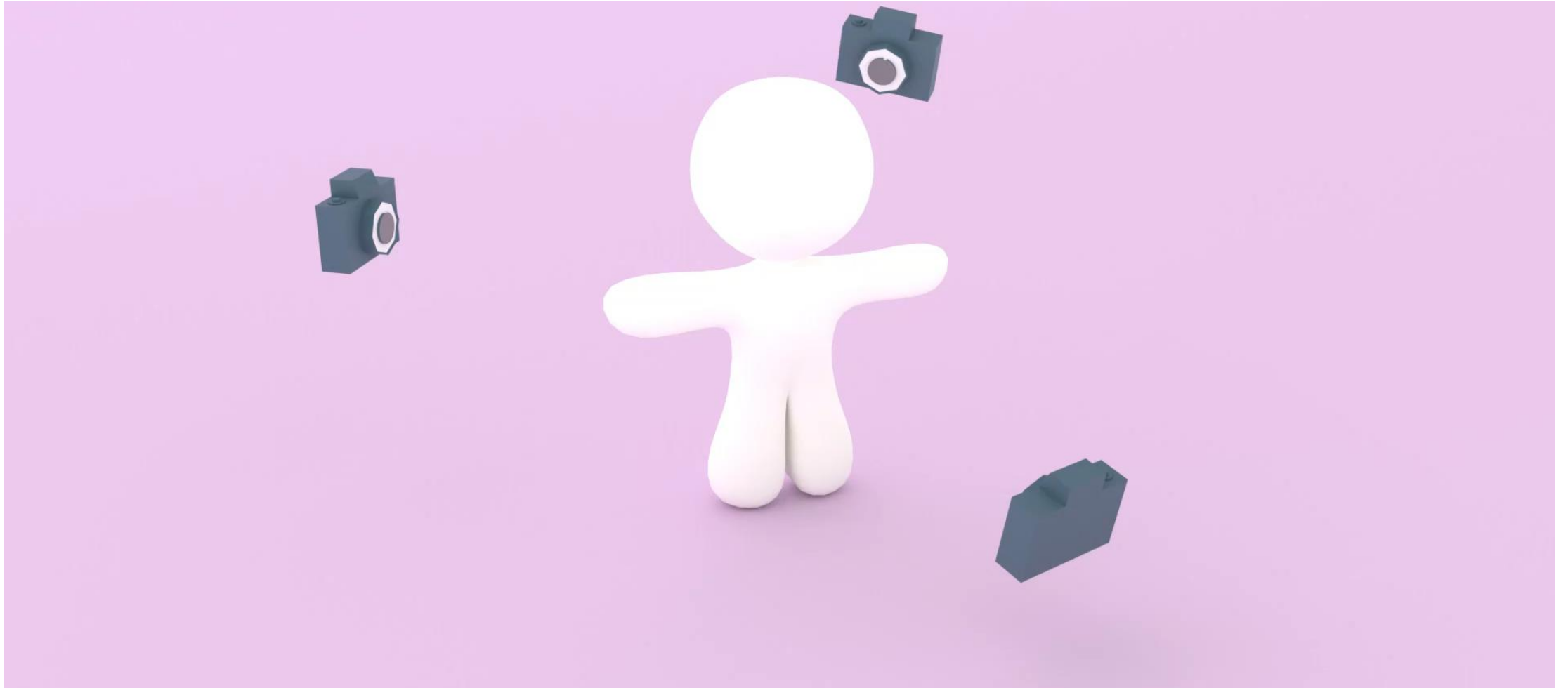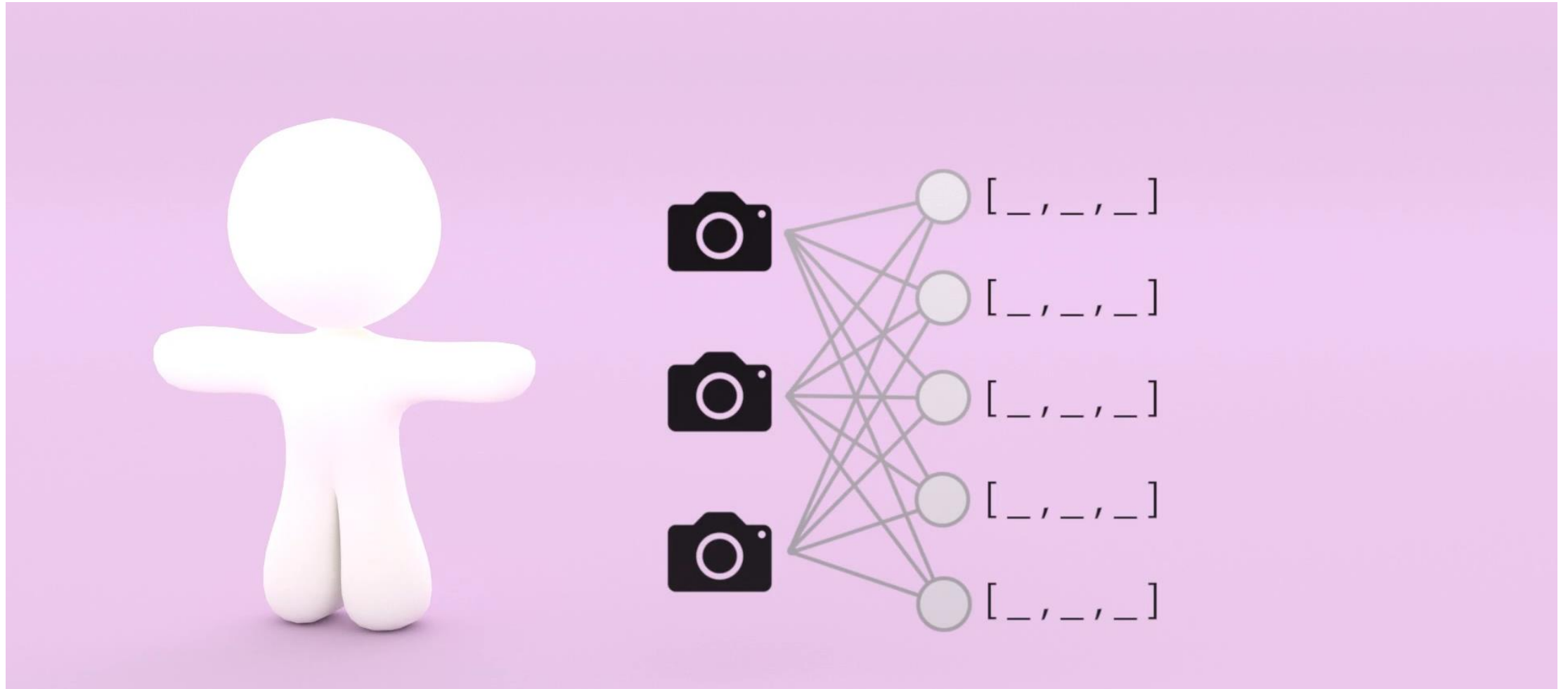
We hope to learn keypoints:
- in 3D world coordinates
- without supervision
- are good representation for control

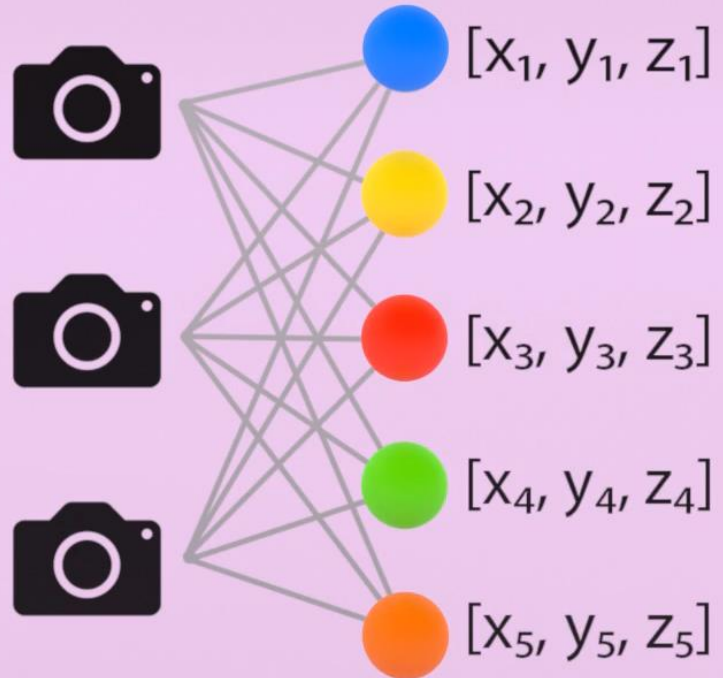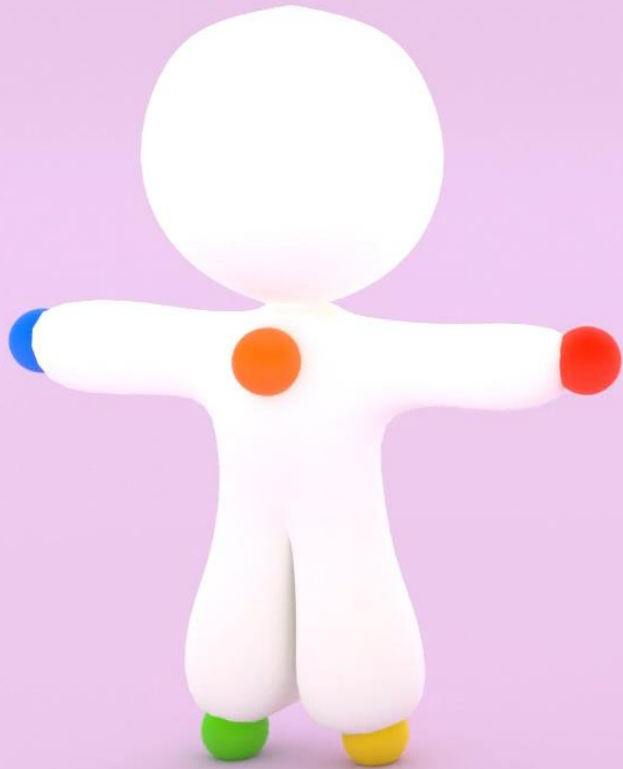# Our 3D Keypoint: Setup

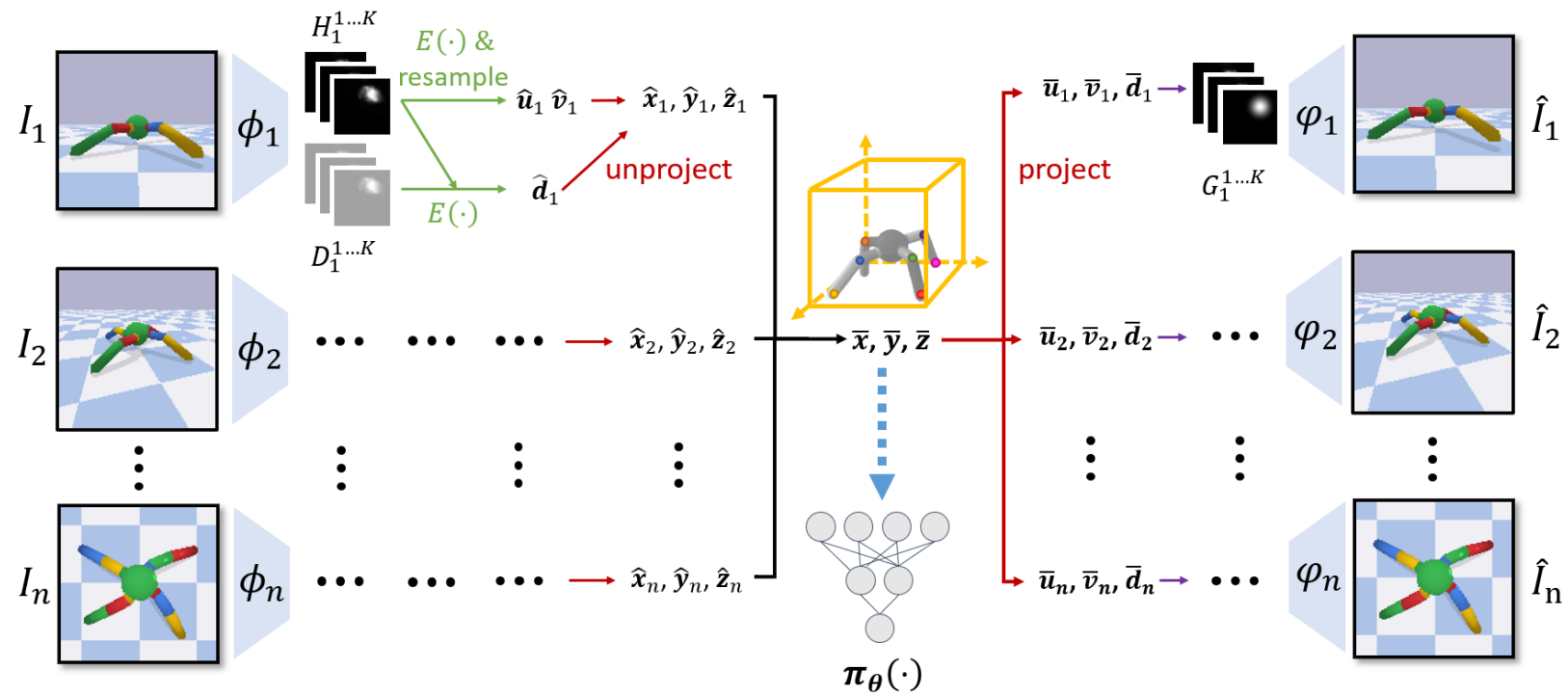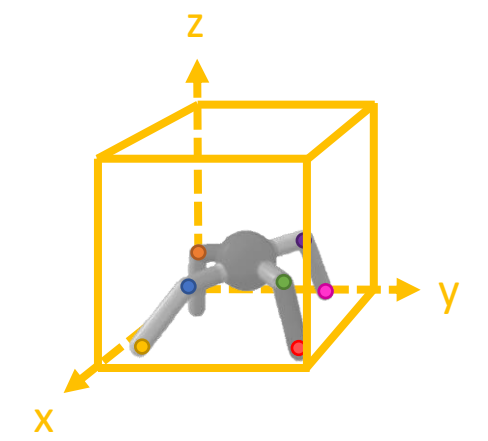# Our 3D Keypoint: Keypoint learning

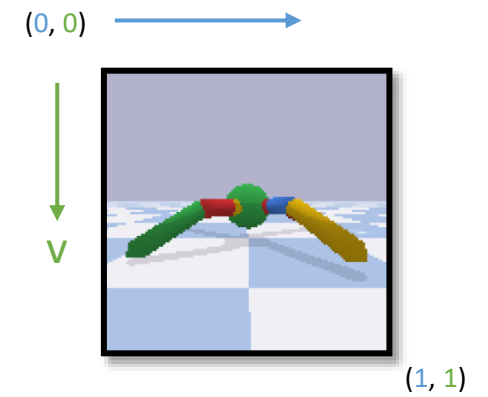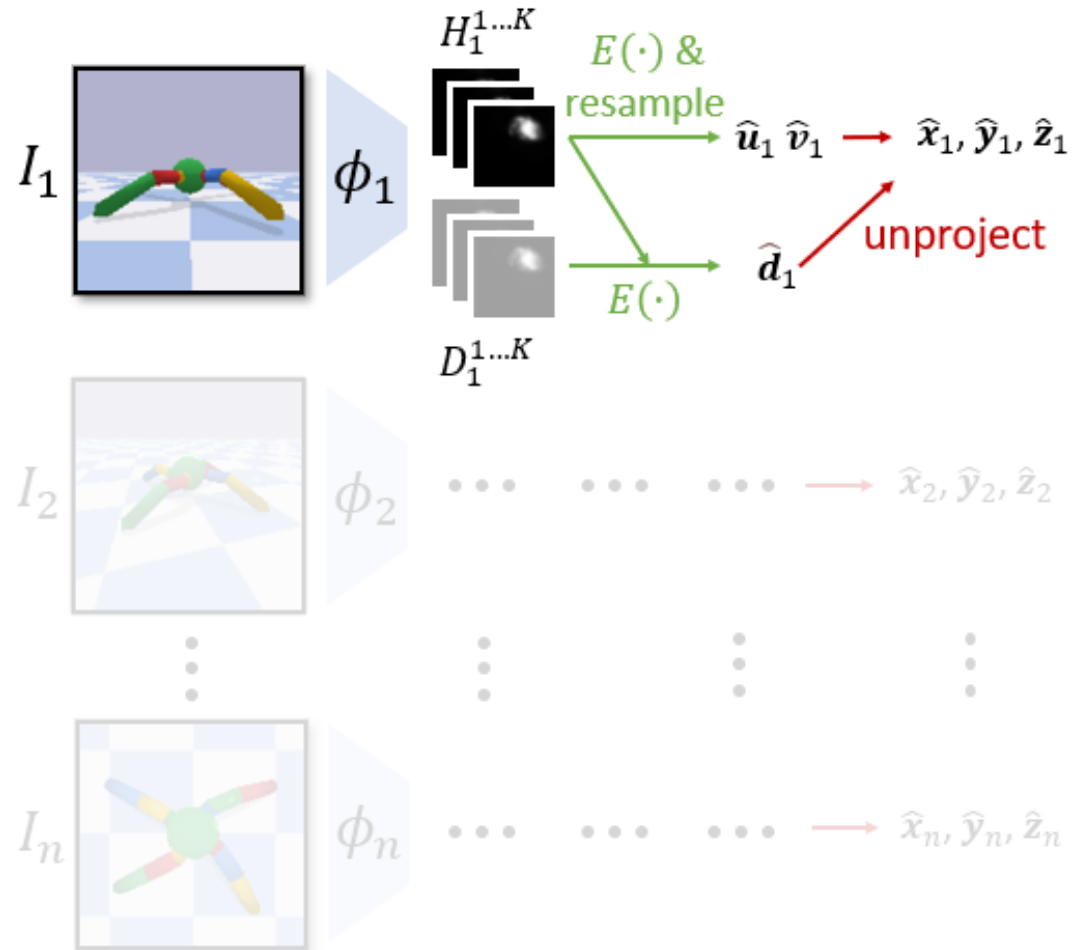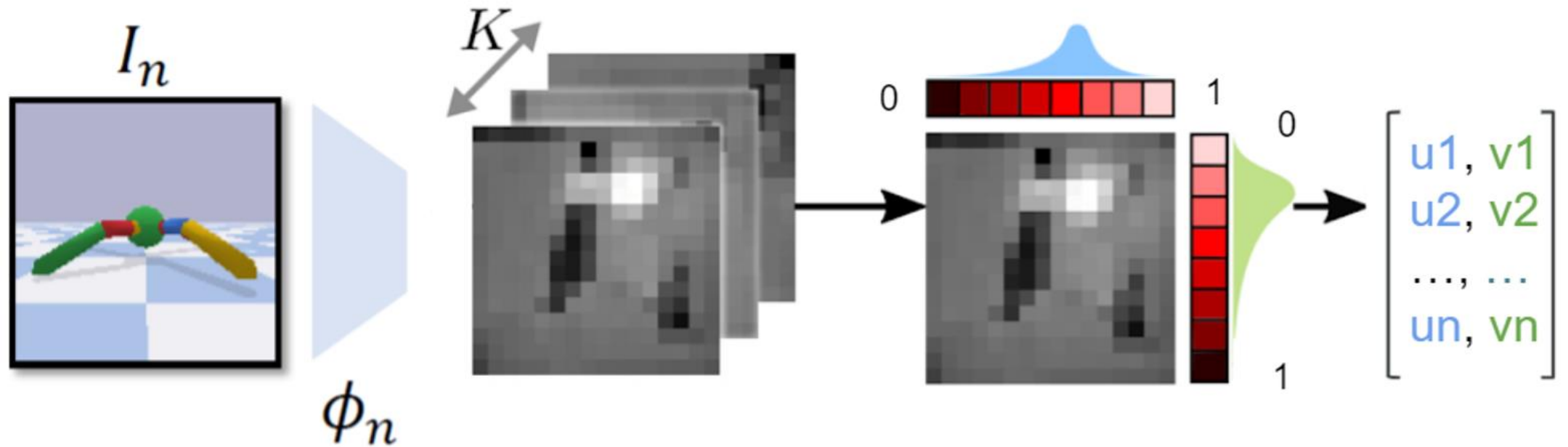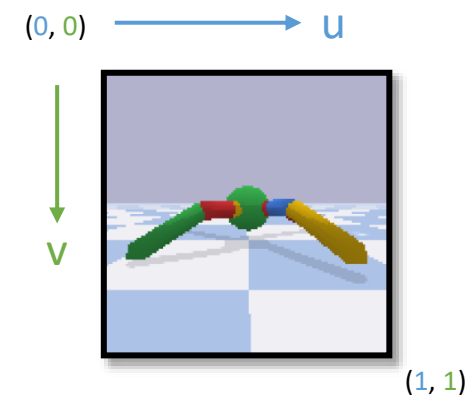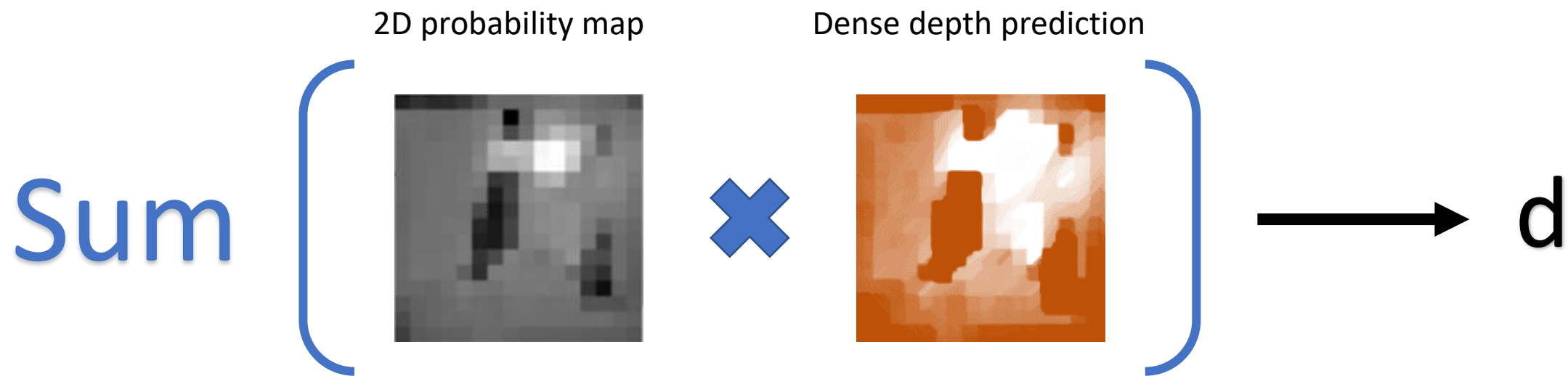# Our 3D Keypoint: policy

# Method: encoder

Unsupervised learning of object landmarks through conditional image generation, Jakab et al., 2018

# Fully differentiable keypoint bottleneck

Sum $\left[\phantom{xx}\times\phantom{xx}\right]$ $\longrightarrow$ d

2D probability map    Dense depth prediction

Depth parameterization

# Method: encoder

# Method: decoder

$$\overline{x}, \overline{y}, \overline{z} \longrightarrow \overline{u}, \overline{v}, \overline{d} \longrightarrow$$



z

x      y

(0, 0)      u

v

(1, 1)

2D Gaussian with
Mean at $(\overline{u}, \overline{v})$
Std $\propto 1 / \overline{d}$

xyz coordinate regains 2D structure in a fully differentiable way!

# Method: decoder

# Method: auto-encoding loss



Core intuition:
To best decode to original image, the 2D gaussians have centers aligned with meaningful points

# Method: multi-view consistency loss



Core intuition:
Some point movements are visible from camera A but not camera B,
B must learn to "hallucinate" these points to minimize disagreement

# Method: policy

# Method: attention

N

$I_1$ $\phi_1$ $H_1^{1...K}$ $E(\cdot)$ & resample $\widehat{u}_1\,\widehat{v}_1$ $\rightarrow$ $\widehat{x}_1, \widehat{y}_1, \widehat{z}_1$

unproject

$\widehat{d}_1$

$E(\cdot)$

$D_1^{1...K}$

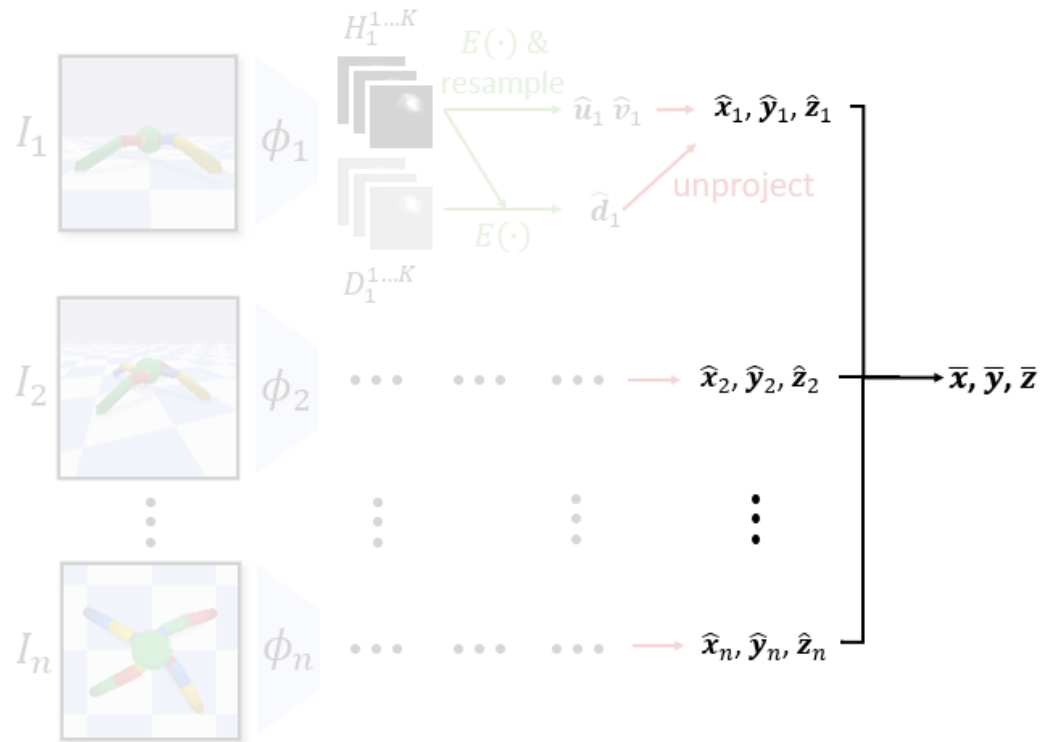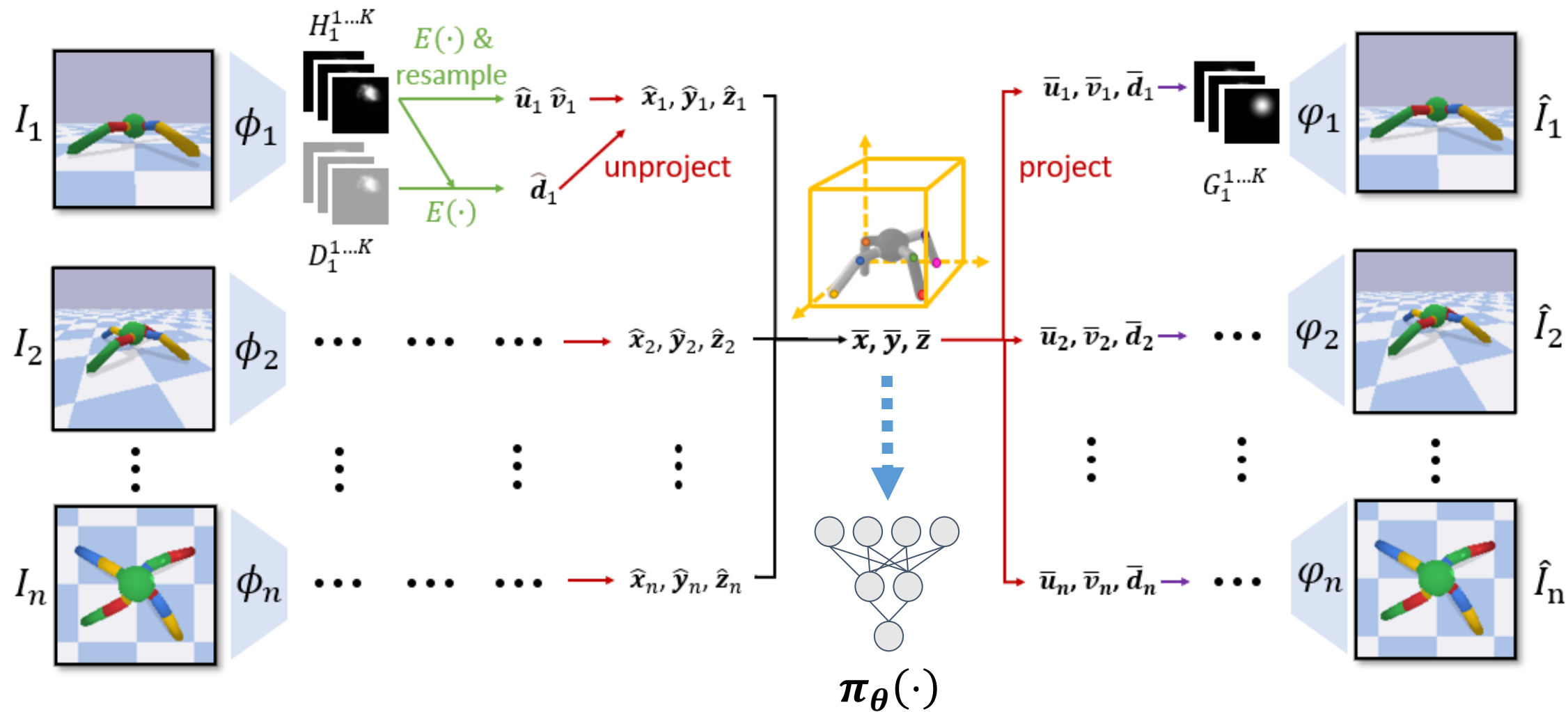Softmax along # camera dimension
Use mean logits of each map as attention logit

$I_2$ $\phi_2$ $\cdots$ $\cdots$ $\cdots$ $\rightarrow$ $\widehat{x}_2, \widehat{y}_2, \widehat{z}_2$ $\longrightarrow$ $\overline{x}, \overline{y}, \overline{z}$

$I_n$ $\phi_n$ $\cdots$ $\cdots$ $\cdots$ $\rightarrow$ $\widehat{x}_n, \widehat{y}_n, \widehat{z}_n$

Allow model to ignore unconfident estimations!

# Random crop as self-supervision

No crop



$(0, 0) \longrightarrow u$

$v$

$(1, 1)$

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \\ ... & ... \\ u_n & v_n \end{bmatrix}$$

Crop



$(0, 0) \longrightarrow u'$

$v'$

$(1, 1)$

$$\begin{bmatrix} u'_1 & v'_1 \\ u'_2 & v'_2 \\ ... & ... \\ u'_n & v'_n \end{bmatrix}$$

$$f \left( \begin{bmatrix} u'_1 & v'_1 \\ u'_2 & v'_2 \\ ... & ... \\ u'_n & v'_n \end{bmatrix} \right)$$

$f(\cdot)$ maps u'v' coordinate to uv coordinate

Coordinates must align with the random cropping to predict well

# Temporal variant



$I^t$

$[\,\overline{x}, \overline{y}, \overline{z}\,]^t$

Policy

$I^{t+1}$

$[\,\overline{x}, \overline{y}, \overline{z}\,]^{t+1}$

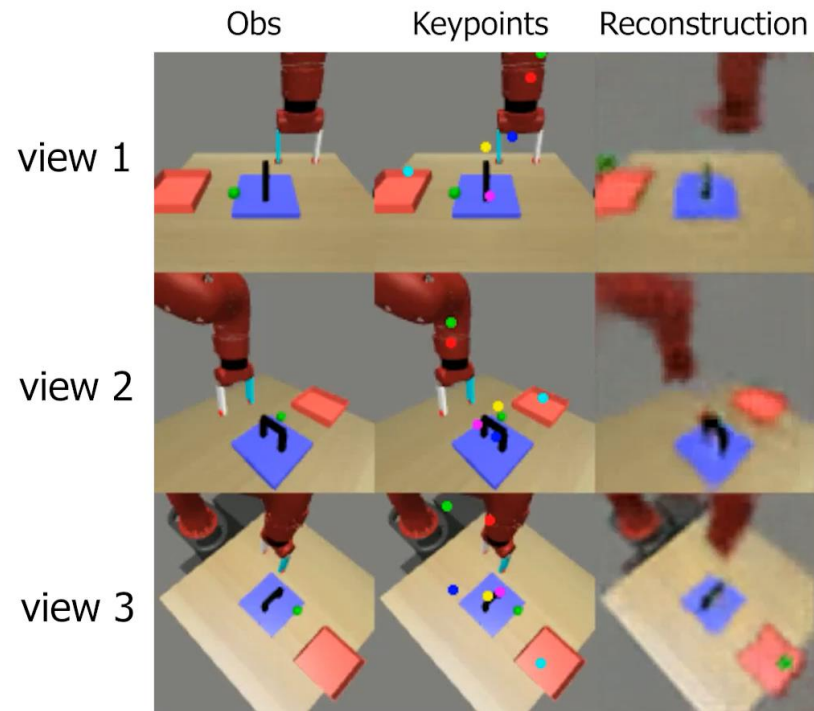$normalized([\,\overline{x}, \overline{y}, \overline{z}\,]^{t+1} - [\,\overline{x}, \overline{y}, \overline{z}\,]^t)$
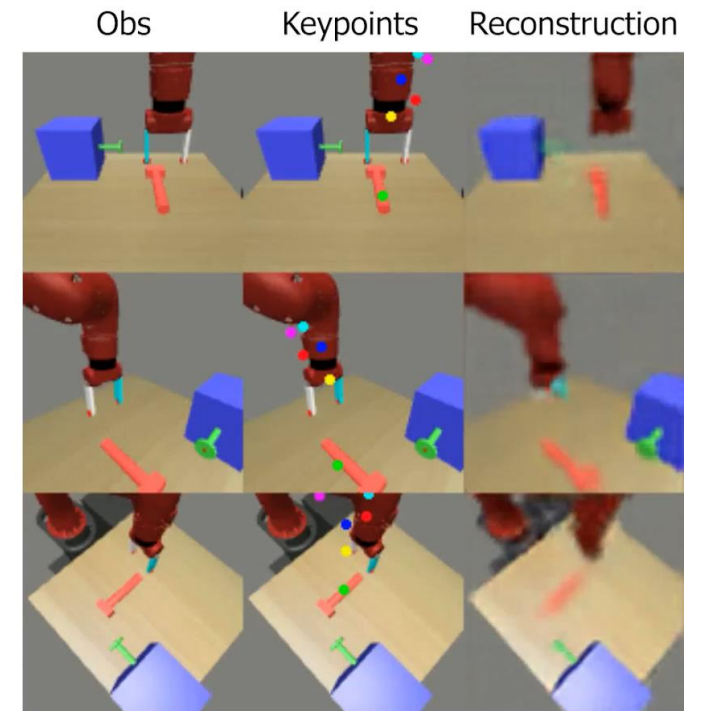
Differences between keypoint prediction is velocity vector! Explicitly normalize as movement feature

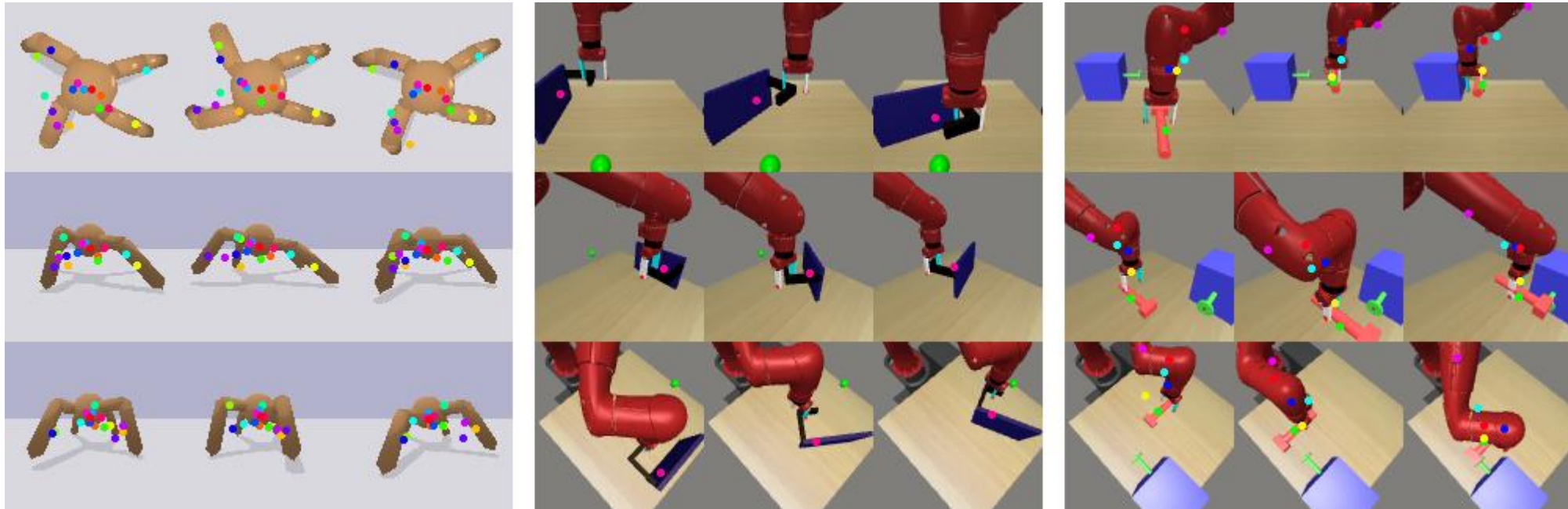# Visualization of Learned 3D Keypoints
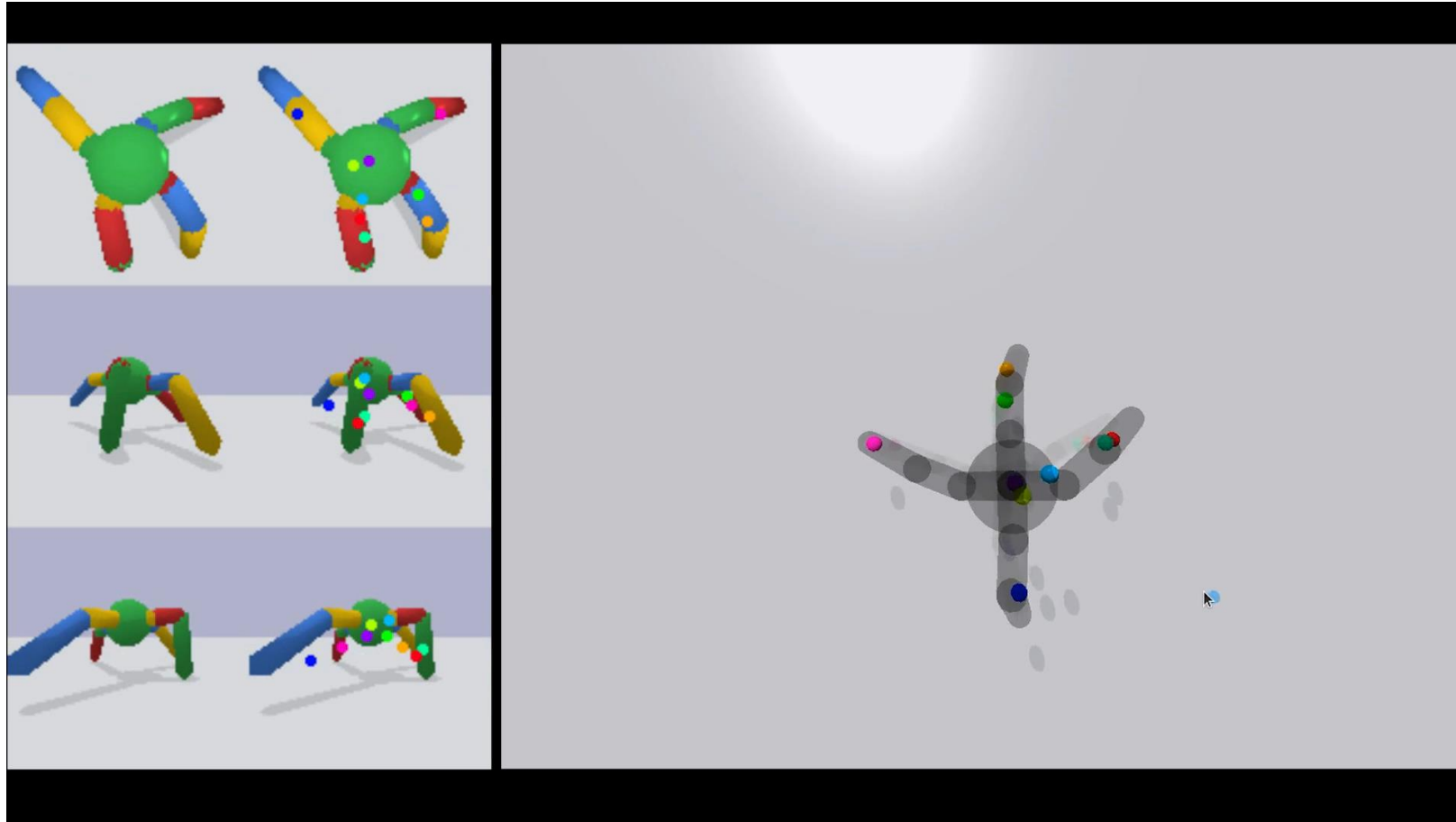


Close Box

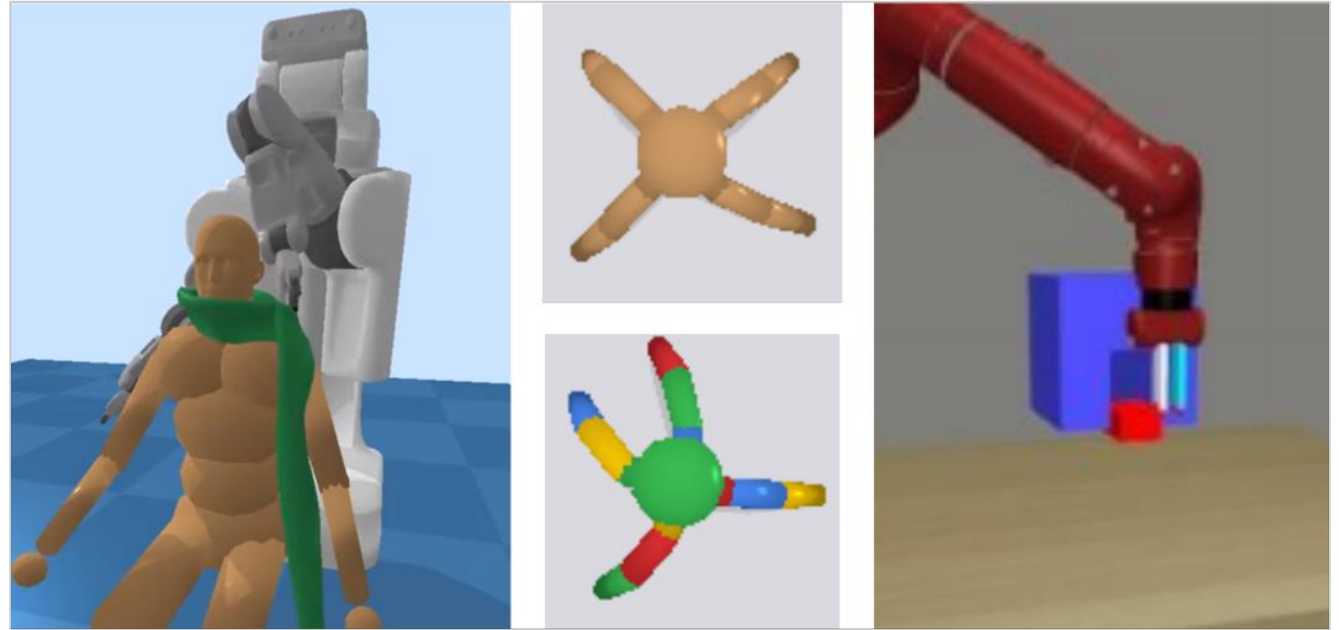Hammer Nail

# Visualization of Learned 3D Keypoints

*attention can be used to filter out unconfident predictions with a threshold!

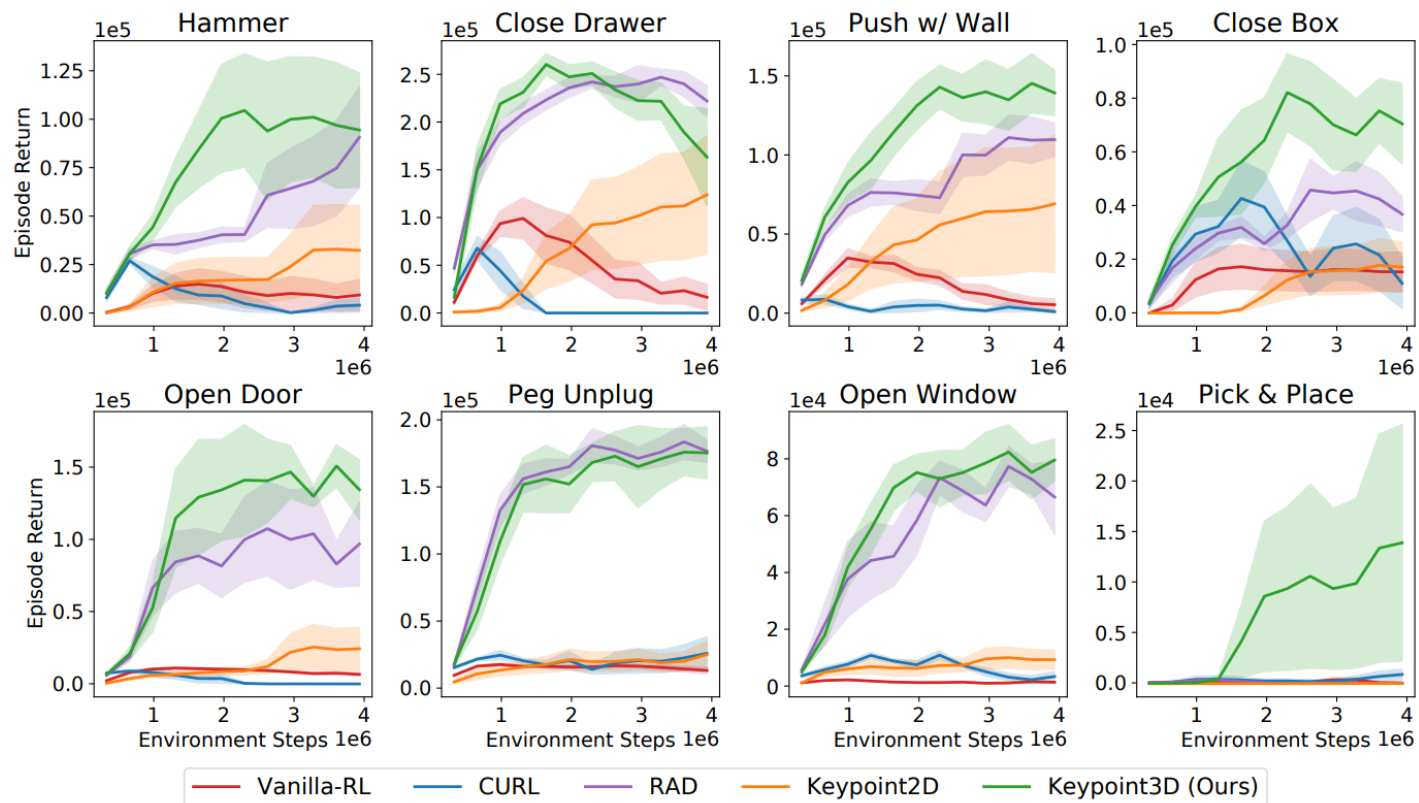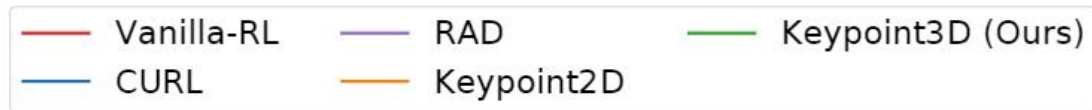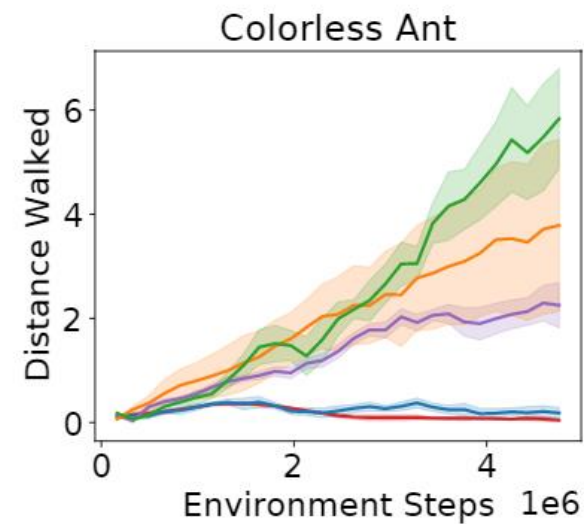# Visualization of Learned 3D Keypoints
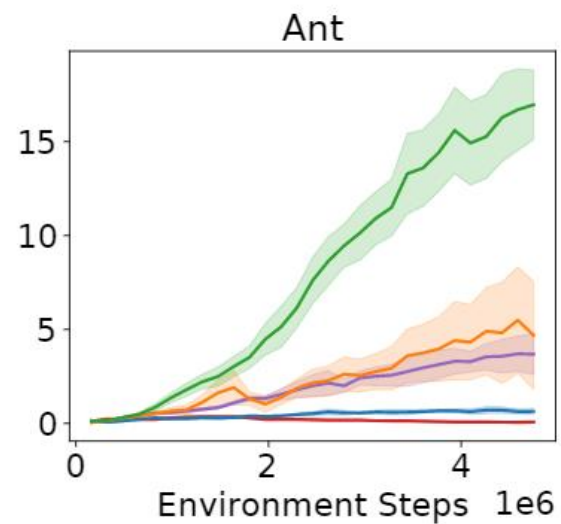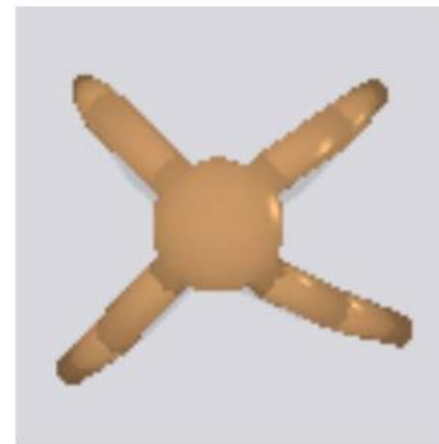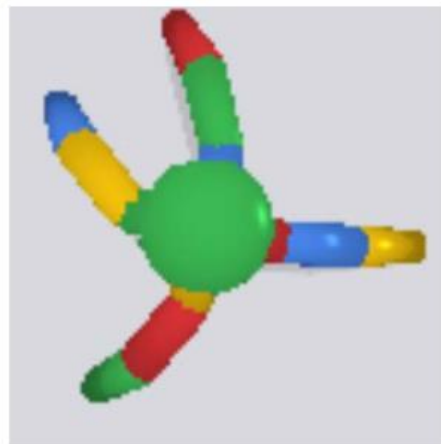
# Experiments Overview



Effectiveness of 3d keypoints for control

- Sample efficiency compared to other representation

- Scalability to higher dimensional control problems(pybullet ant)

- Effectiveness on low-textured objects

- Ability to adapt to deformable objects (scarf manipulation)
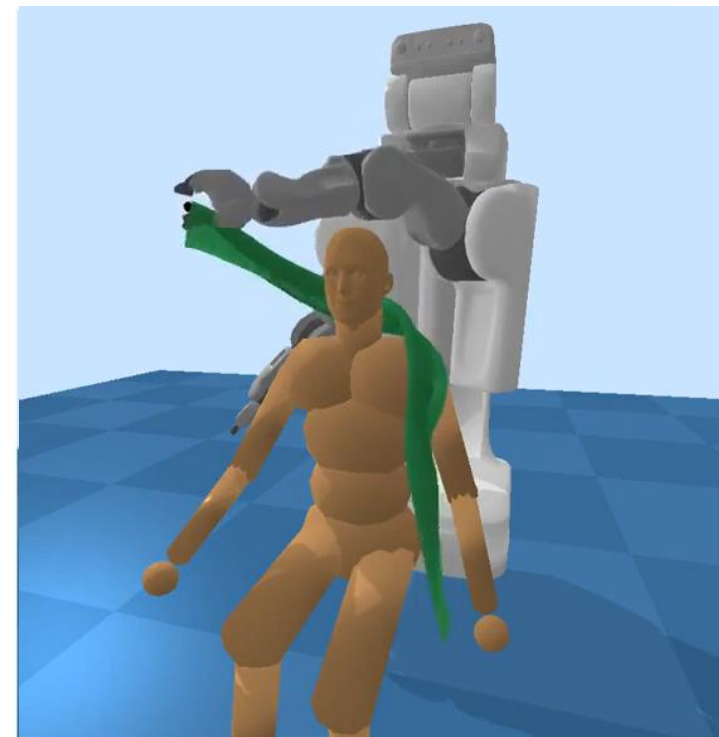
# Sample efficiency in manipulation

High dimensional control and low textured variant



Ant

Colorless Ant

Distance Walked

Environment Steps 1e6

Environment Steps 1e6

Vanilla-RL    RAD    Keypoint3D (Ours)
CURL    Keypoint2D

Deformable manipulation

Deformable Manipulation

Episode Return vs Environment Steps 1e6

Vanilla-RL
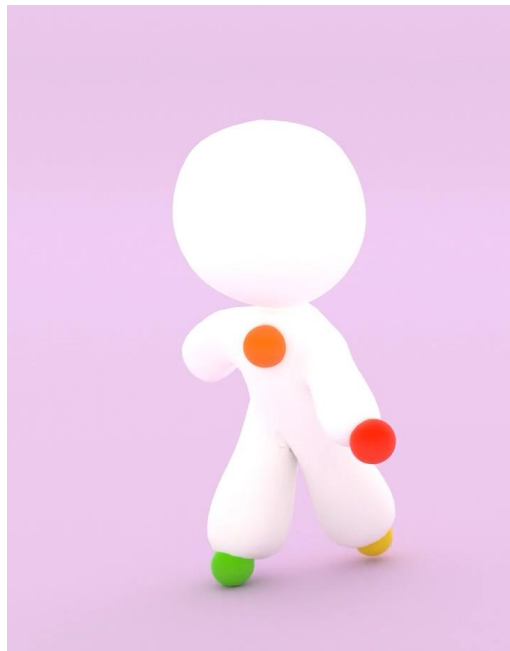CURL
RAD
Keypoint2D
Keypoint3D (Ours)

# Summary

- We propose a framework to learn 3D keypoints without supervision for continuous control

- We leverage multi-view auto-encoding with a 3D keypoint bottleneck to learn meaningful 3d keypoints; We jointly train policy learning in conjunction with keypoint learning

- Our method achieves significant sample efficiency improvement in a variety of 3D environments.

- The 3D keypoints learned by our algorithm are consistent across space and time.

We hope our method serves as a bridge between pixel domain and 3D control tasks.

# More Details:



- [Website]  https://buoyancy99.github.io/unsup-3d-keypoints
- [Code]  https://github.com/buoyancy99/unsup-3d-keypoints
- [Paper]  https://arxiv.org/pdf/2106.07643.pdf