



DR

DeepRob

[Group 5] Lecture 5
Multisensory Learning + Manipulation
by Mason Hawver, Ryan Diaz, and Hanchen Cui
University of Minnesota



Image Source: D. F. Gomes, P. Paoletti, and S. Luo, "Generation of gelsight tactile images for sim2real learning," IEEE Robotics and Automation Letters, vol. 6, no. 2, pp. 4177–4184, 2021.



Outline

Each person has one theme and take away - think about that when adding slides to the deck

mason - datagen

hanchen- multimae

ryan - learning with visuotactile data

cite the image/vid src
slide for questions
practice slides

. intro to multimodal sensors

. Different sensors (camera, depth, audio, contact, etc, force and torque, temporal)

. Structures (Dimensions and examples) [one each]

. Maps

. Image, Depth, Segmentation, visual tactile (GelSight / bubble grippers!!!!) [Mason]

. Multi Spectral images (all frequencies of light <- great for farming!) [Mason]

. Clouds

. Point cloud (Lidars, multi camera), gaussian splate [Hanchen]

. Time series:

. Force Torque, audio [Ryan]

. This is not the only set of sensors...

. Intro to Multi Modal Learning wrt to Foundation models

. **Multimodal** (vision+language) vs **multisensory** (raw sensor data)

. Intro pretraining (foundation models) (quick) [mason]

. individual [one each]

. CNNs, ViTs for maps

. Point++ for point clouds

. Time series models????? -> ask ryan (FFT, MLP)

. Individual encoders (images (ResNet, ViT), force-torque (MLP, CausalConv, etc.), audio)

. pretraining with multimodal data (add timeline stuff here)

— ~~LLM vs VLM stuff [Hanchen]~~

. Multimae [hanchen]

. MSVT, VTT, SVFL, AugInsert [Ryan] <- Using visuotactile data, [with audio] Hearing Touch, See Hear Feel





How do we sense and perceive the world?

How can **robots** sense and perceive the world?

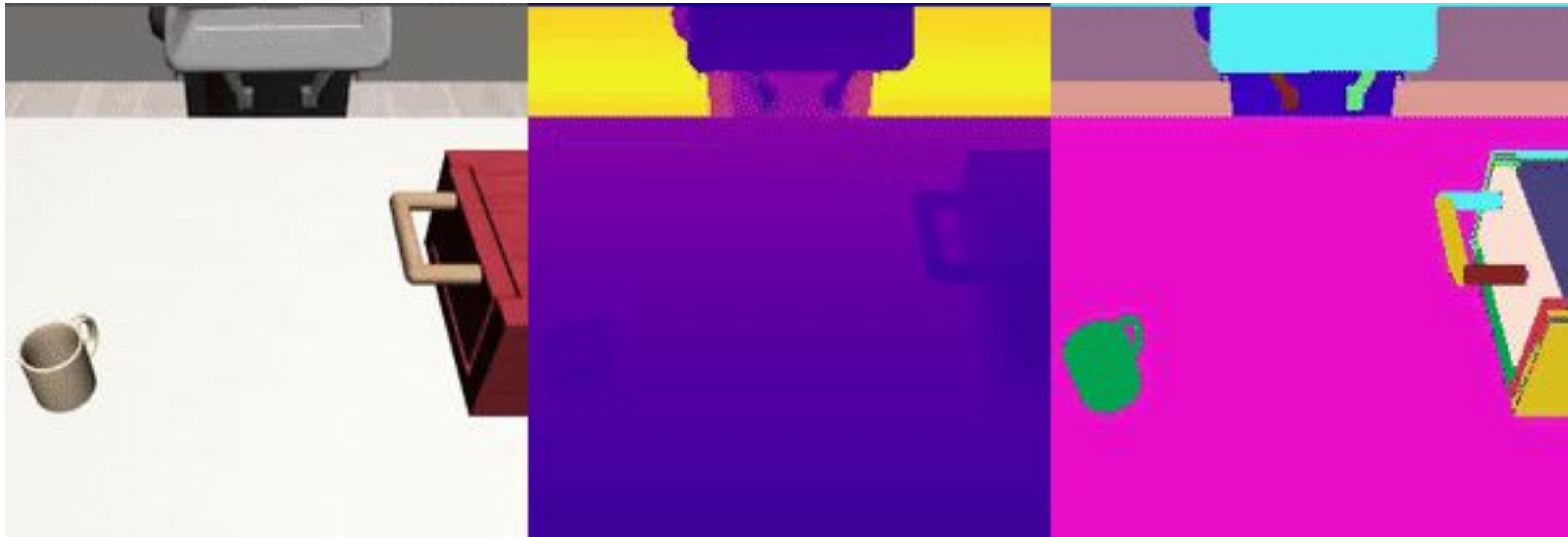


2D Maps

RGB Color

Depth

Segmentation



3xWxH, uint8

1xWxH, float

1xWxH, uint8

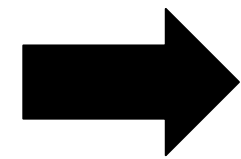
2D Maps - Visuotactile

Visuotactile: Represent Tactile Information with Vision



2D Maps - Visuotactile

Visuotactile: Represent Tactile Information with Vision



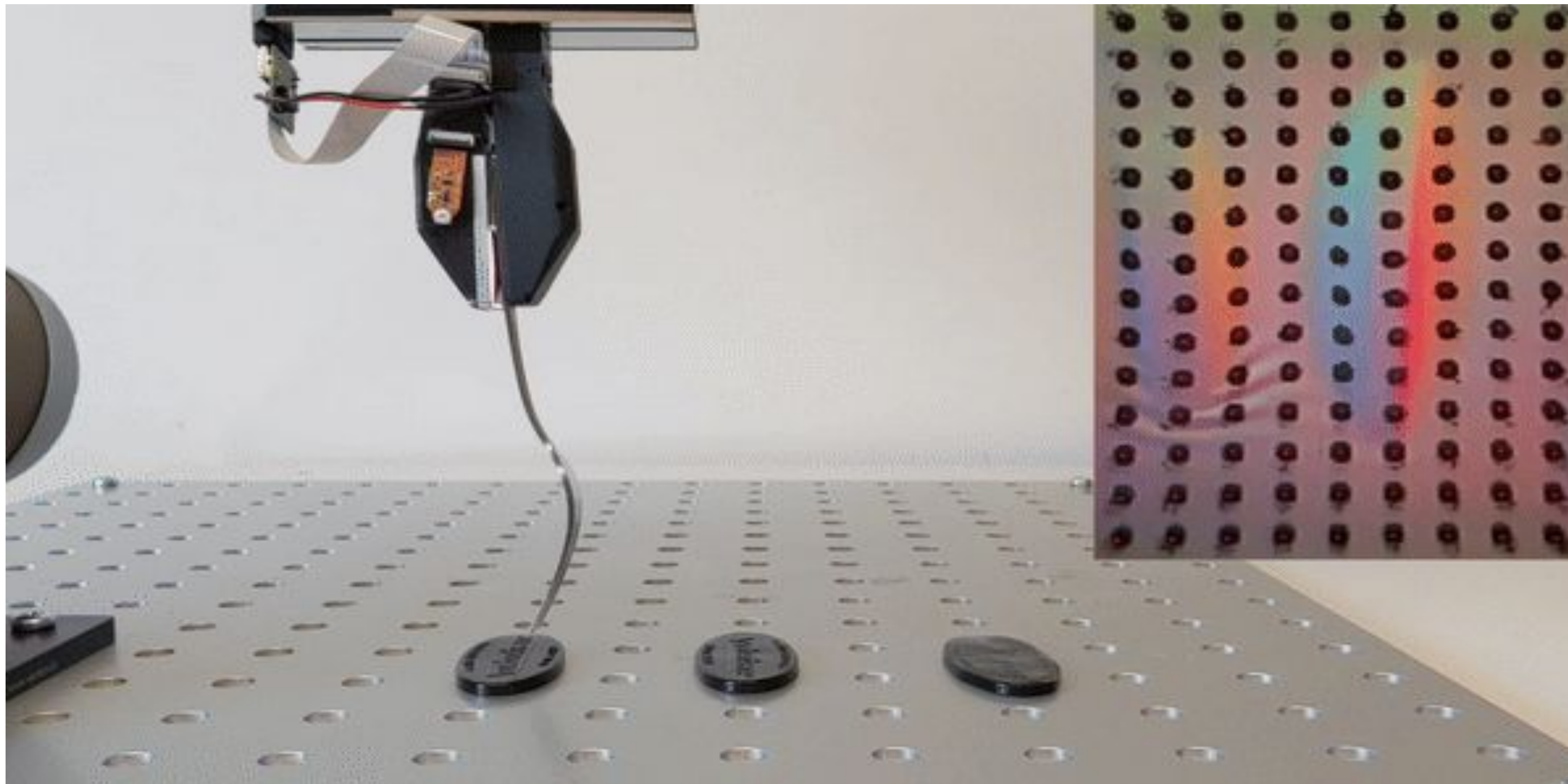
Two Grayscale Images
2x1xWxH, unit8

Credit: RPM Miles and Aaron, and TRI



2D Maps - Visuotactile

Visuotactile: Represent Tactile Information with Vision



Vector Field, Map of Vectors:
 $2 \times W \times H$, floats

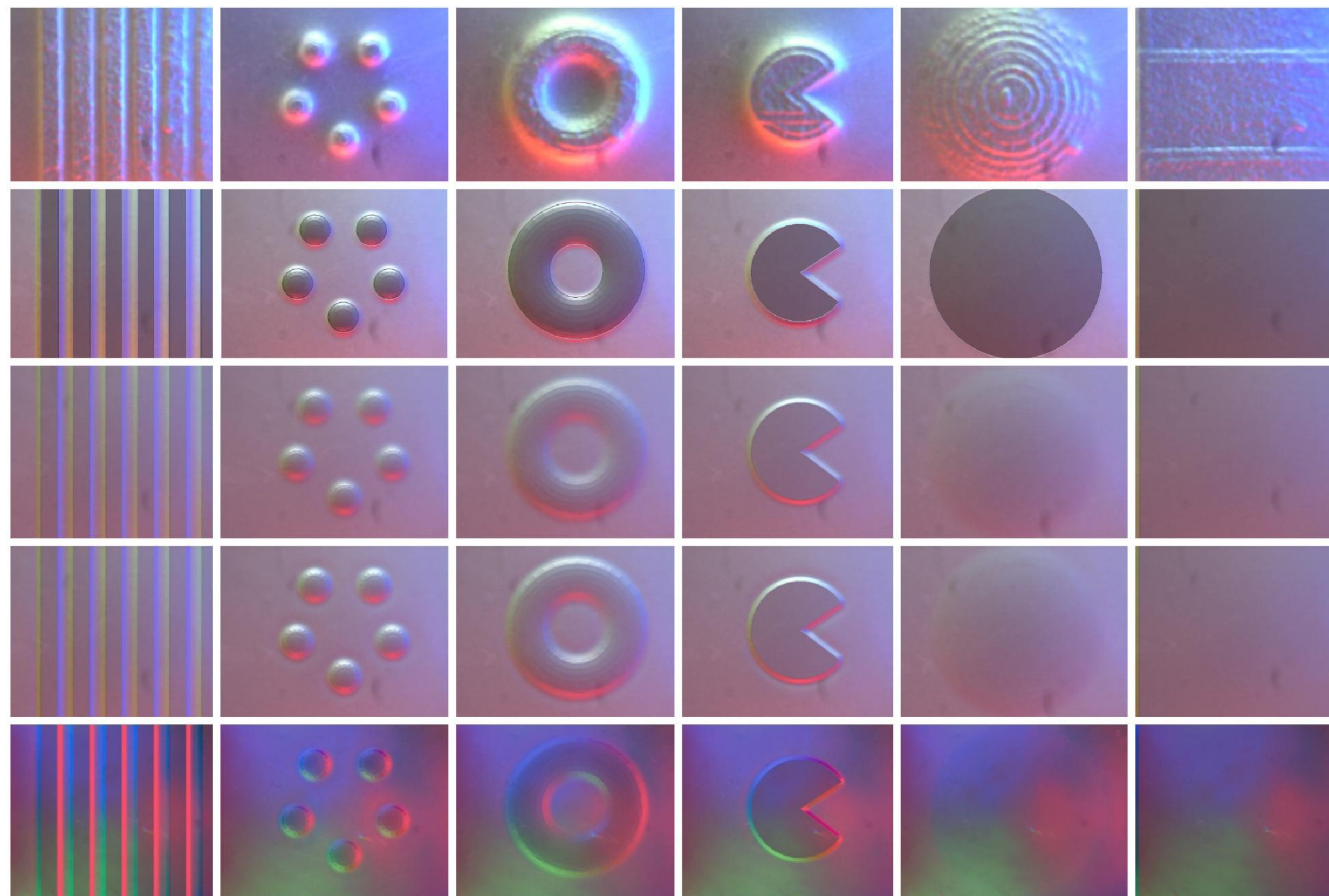
Credit: Rui Li and GelSight





2D Maps - Visuotactile

Visuotactile: Represent Tactile Information with Vision

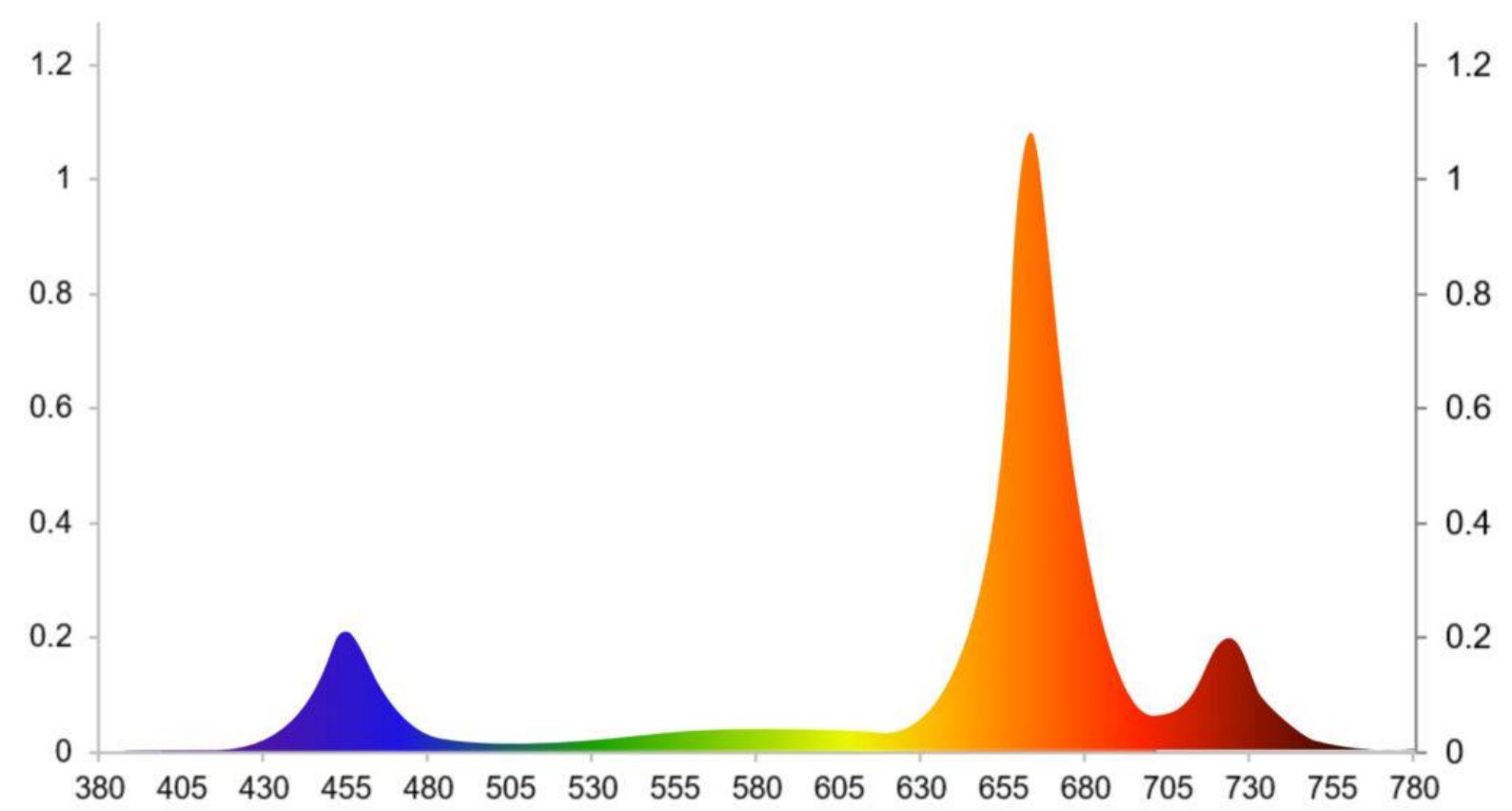
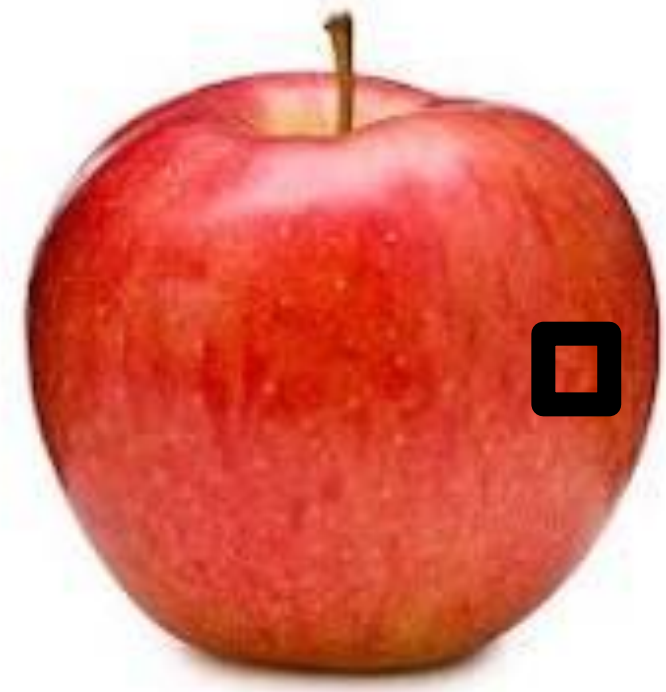


Gel Sight, RGB Image:
3xWxH, unit8





2D Maps - MultiSpectral



2D Map of Histograms
BinsxWxH, unit8

Credit: IdeaForge and Mathworks





2D Maps - MultiSpectral



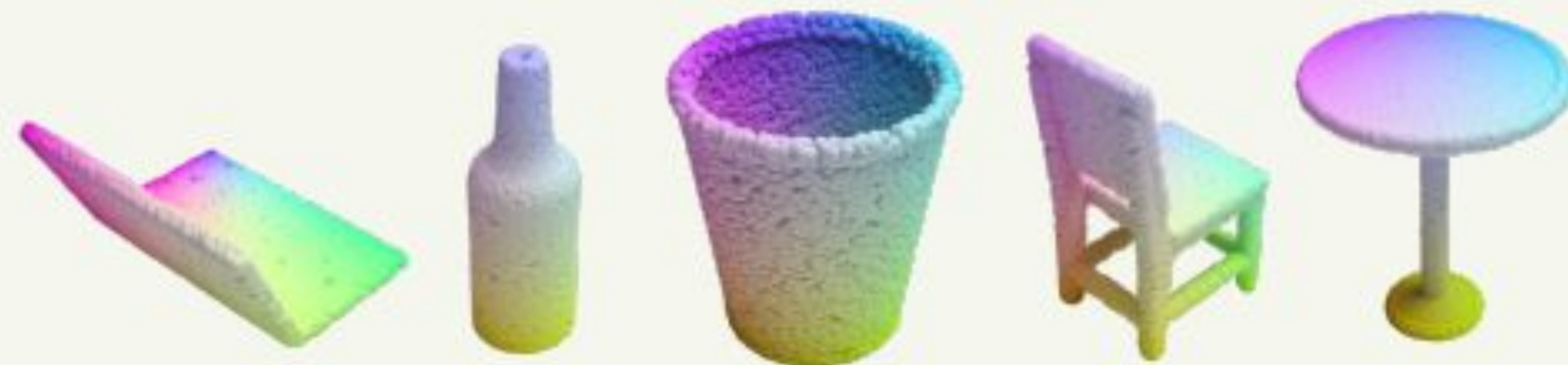
Credit: IdeaForge and Mathworks

DR

Perceive the 3D world



3D Point Clouds
Accurate geometric information

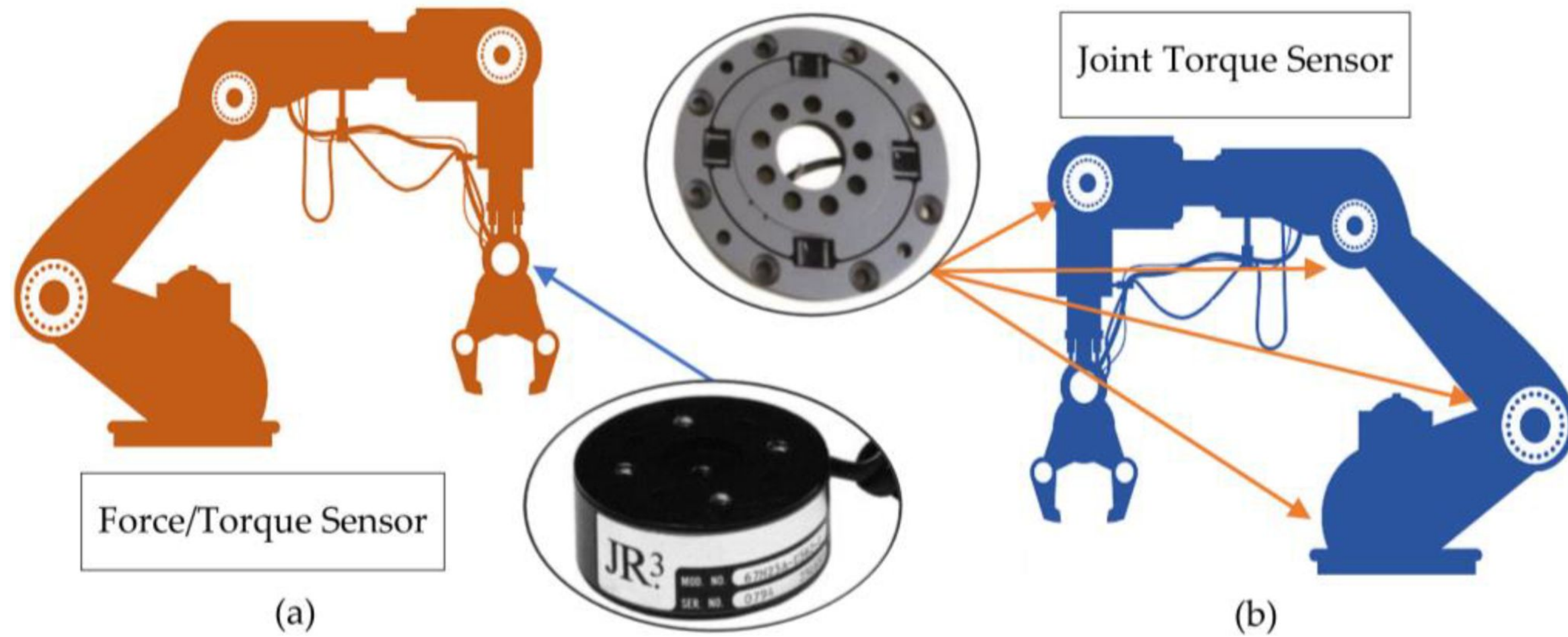


data structure: (N, X, Y, Z)

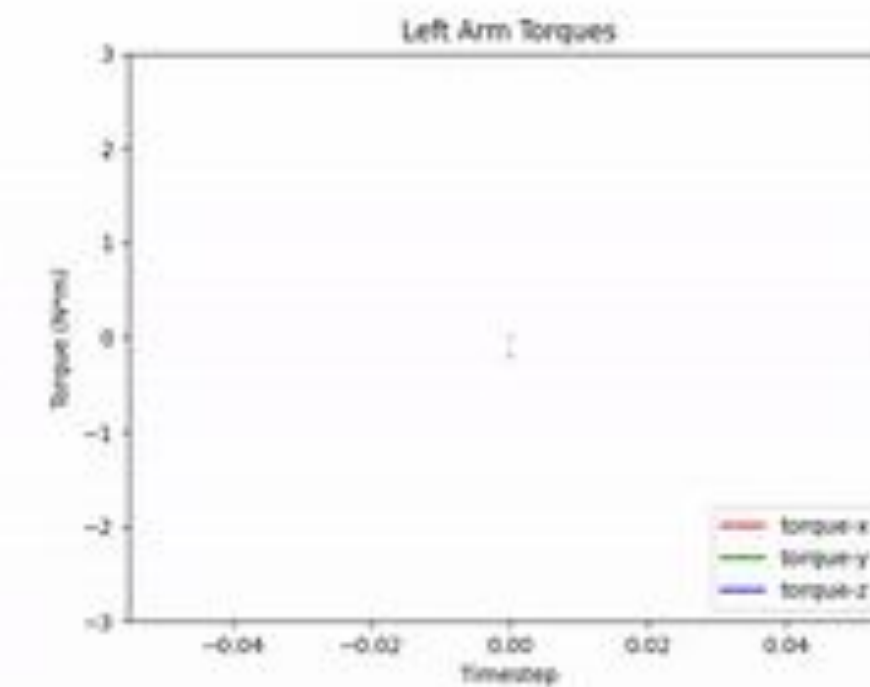
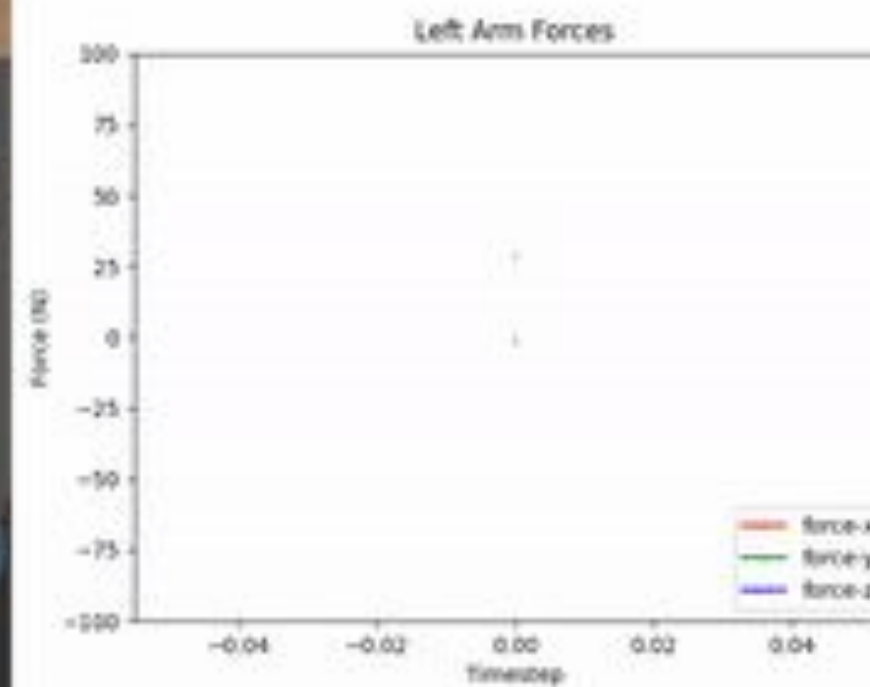
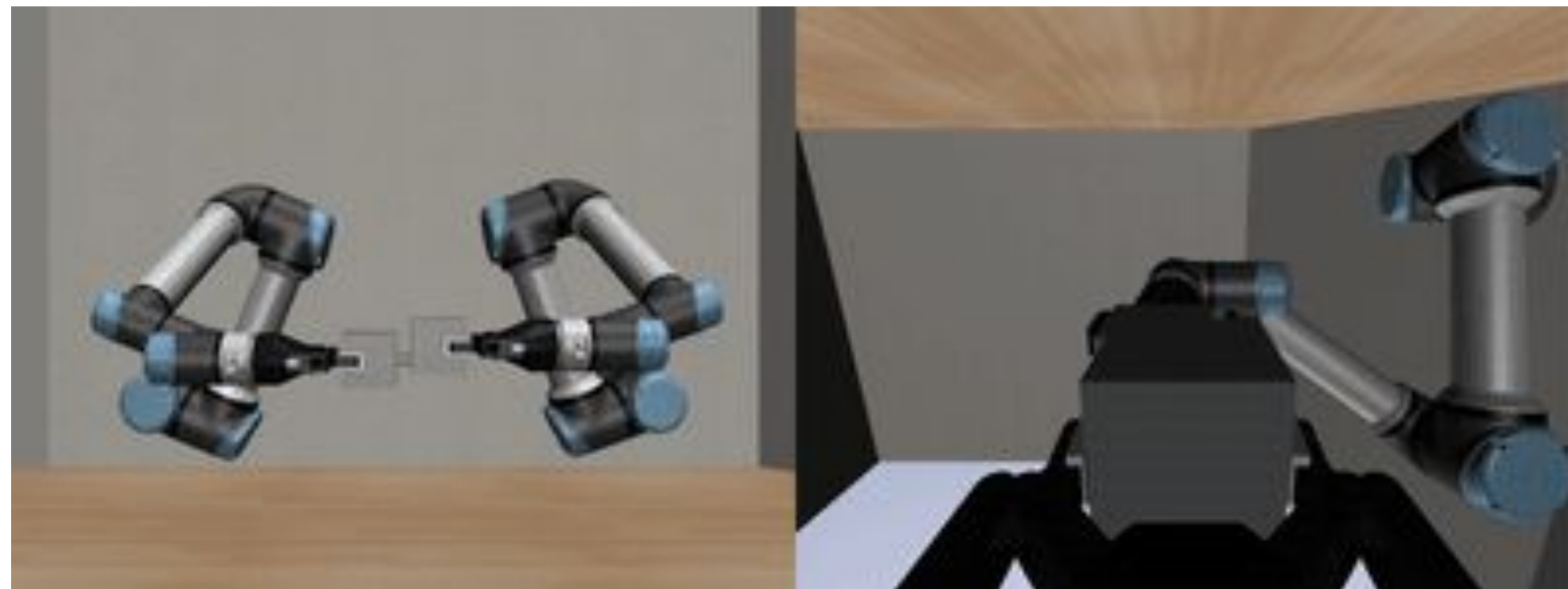
unordered: permutation invariant



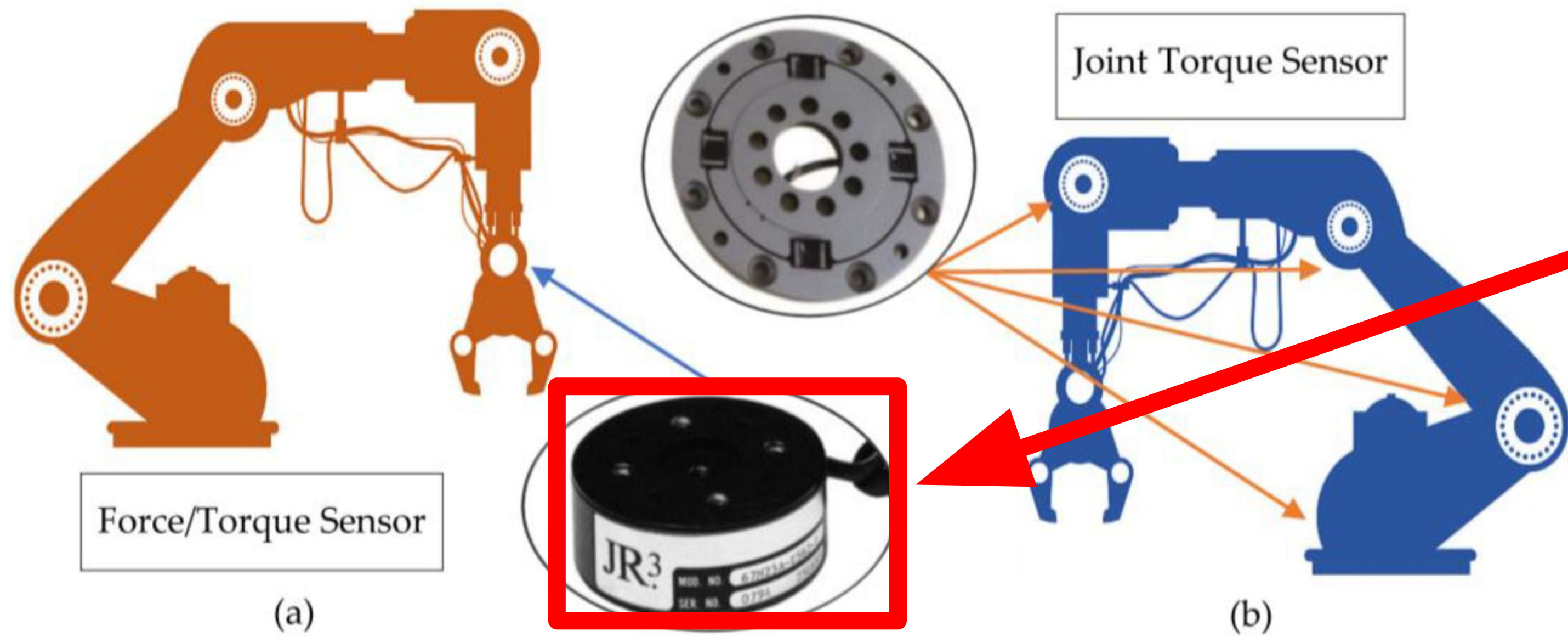
Force-Torque Sensing



R. P. Ubeda, S. C. Guti´errez Rubert, R. Zotovic Staniscic, and ´A. Perles Ivars, “Design and manufacturing of an ultra-low-cost custom torque sensor for robotics,” *Sensors*, vol. 18, no. 6, 2018.



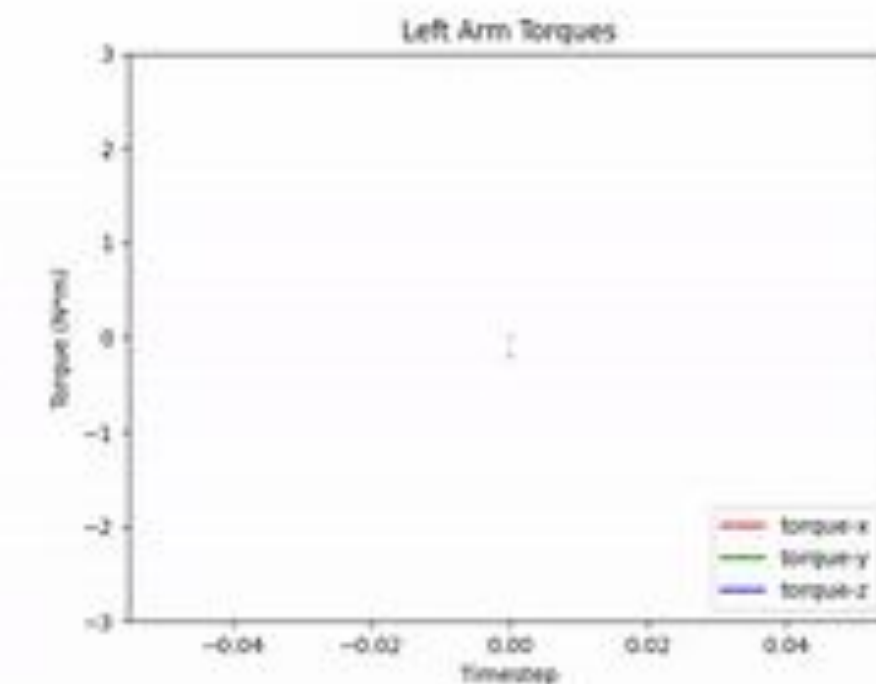
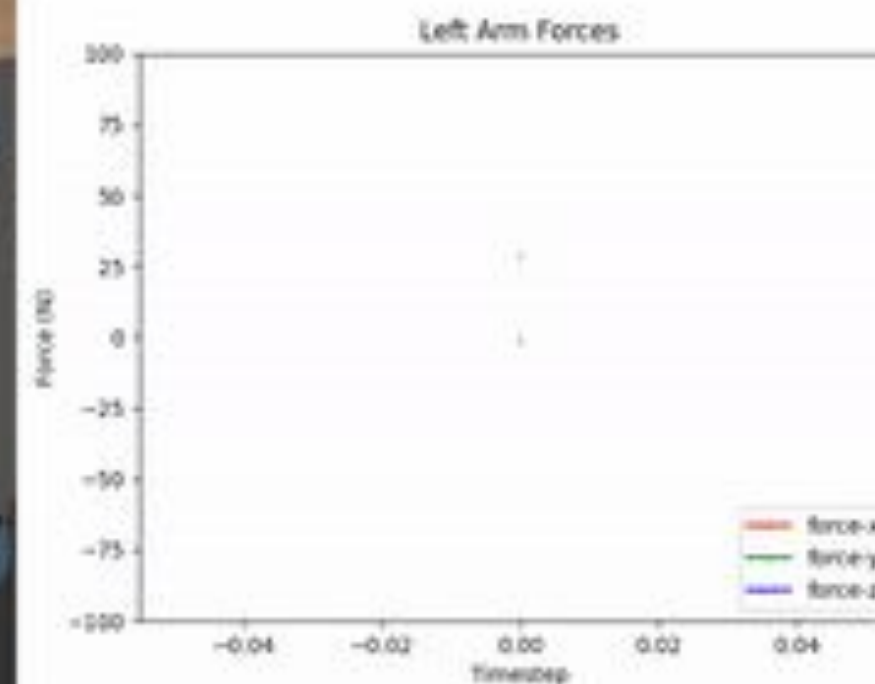
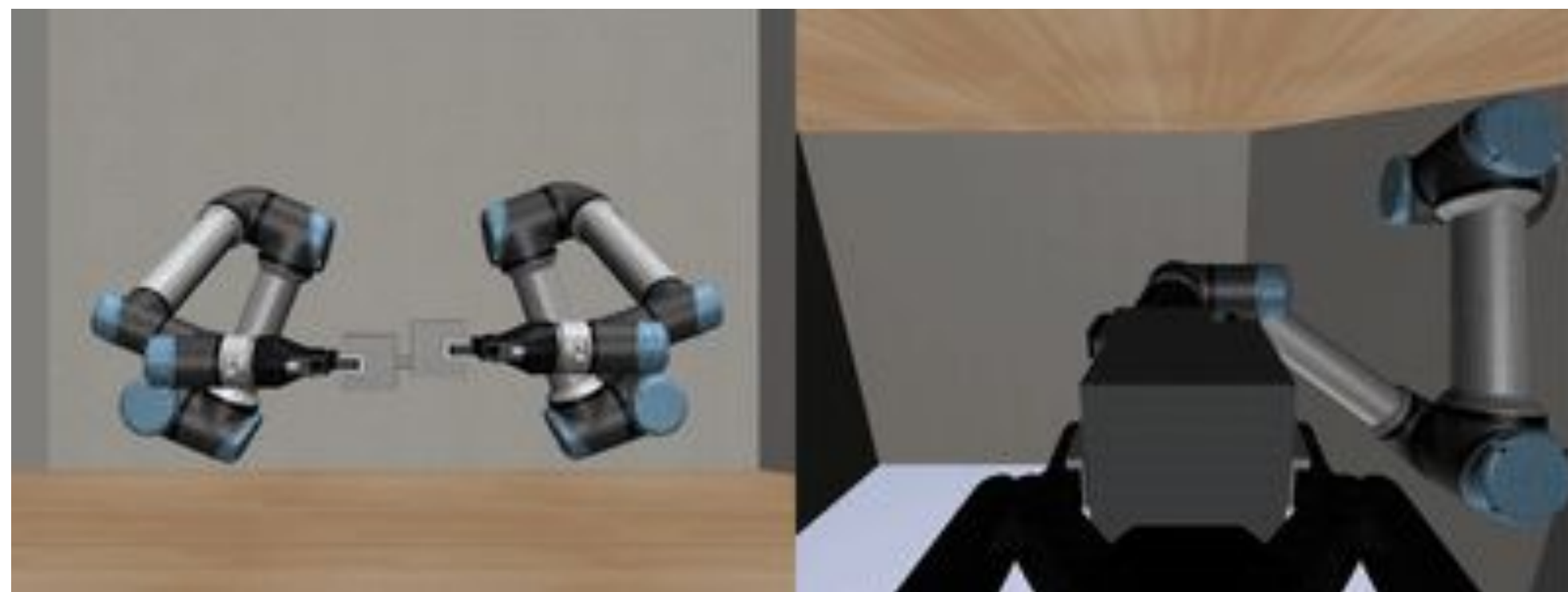
Force-Torque Sensing



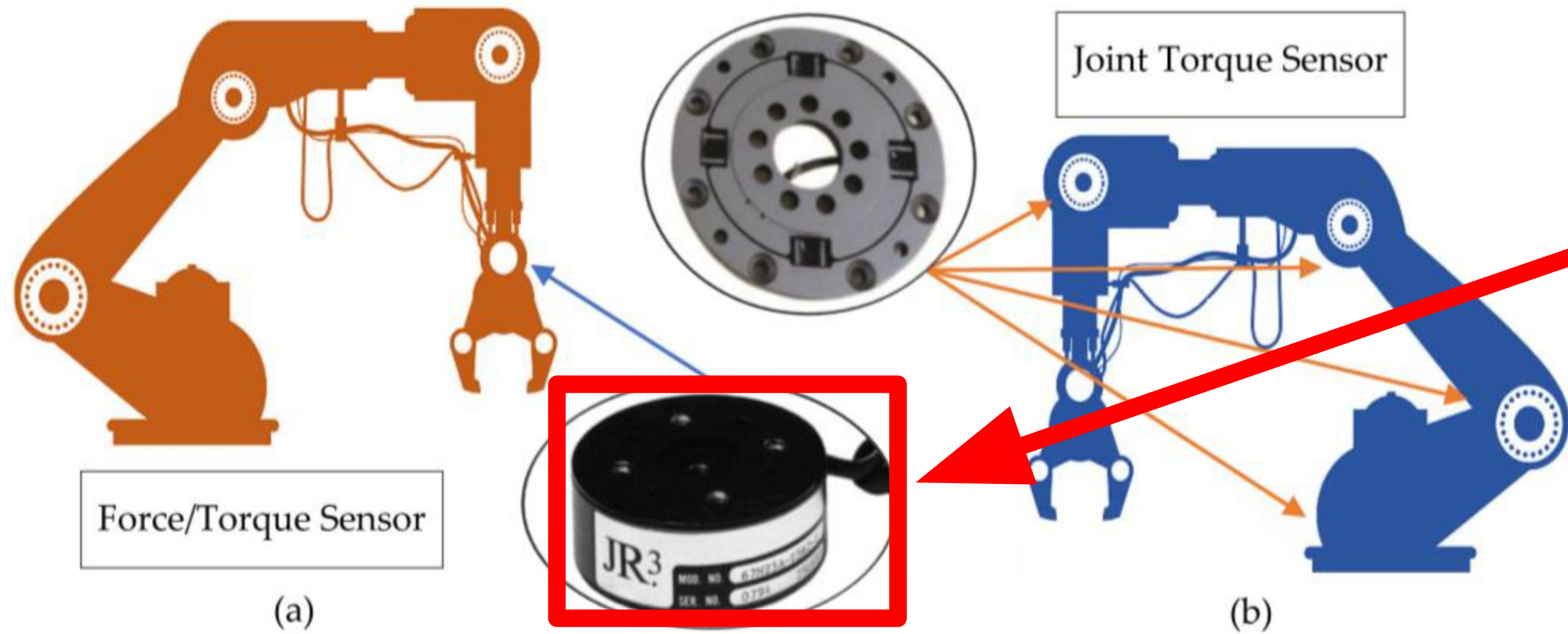
6-axis force-torque sensor along X,Y,Z axes relative to EEF frame

F/T measured using electrical signals

R. P. Ubeda, S. C. Guti´errez Rubert, R. Zotovic Stanisic, and ´A. Perles Ivars, “Design and manufacturing of an ultra-low-cost custom torque sensor for robotics,” *Sensors*, vol. 18, no. 6, 2018.



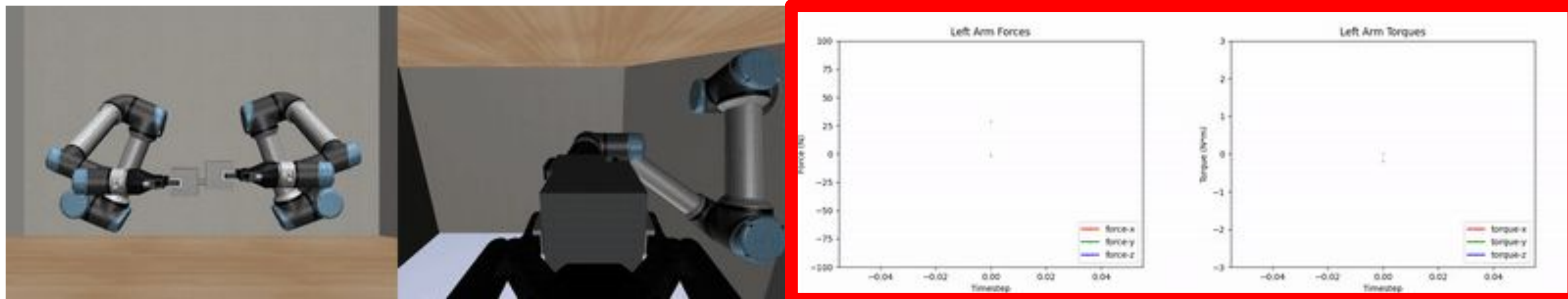
Force-Torque Sensing



6-axis force-torque sensor along X,Y,Z axes relative to EEF frame

F/T measured using electrical signals

R. P. Ubeda, S. C. Guti´errez Rubert, R. Zotovic Staniscic, and ´A. Perles Ivars, “Design and manufacturing of an ultra-low-cost custom torque sensor for robotics,” Sensors, vol. 18, no. 6, 2018.



Hx6, float





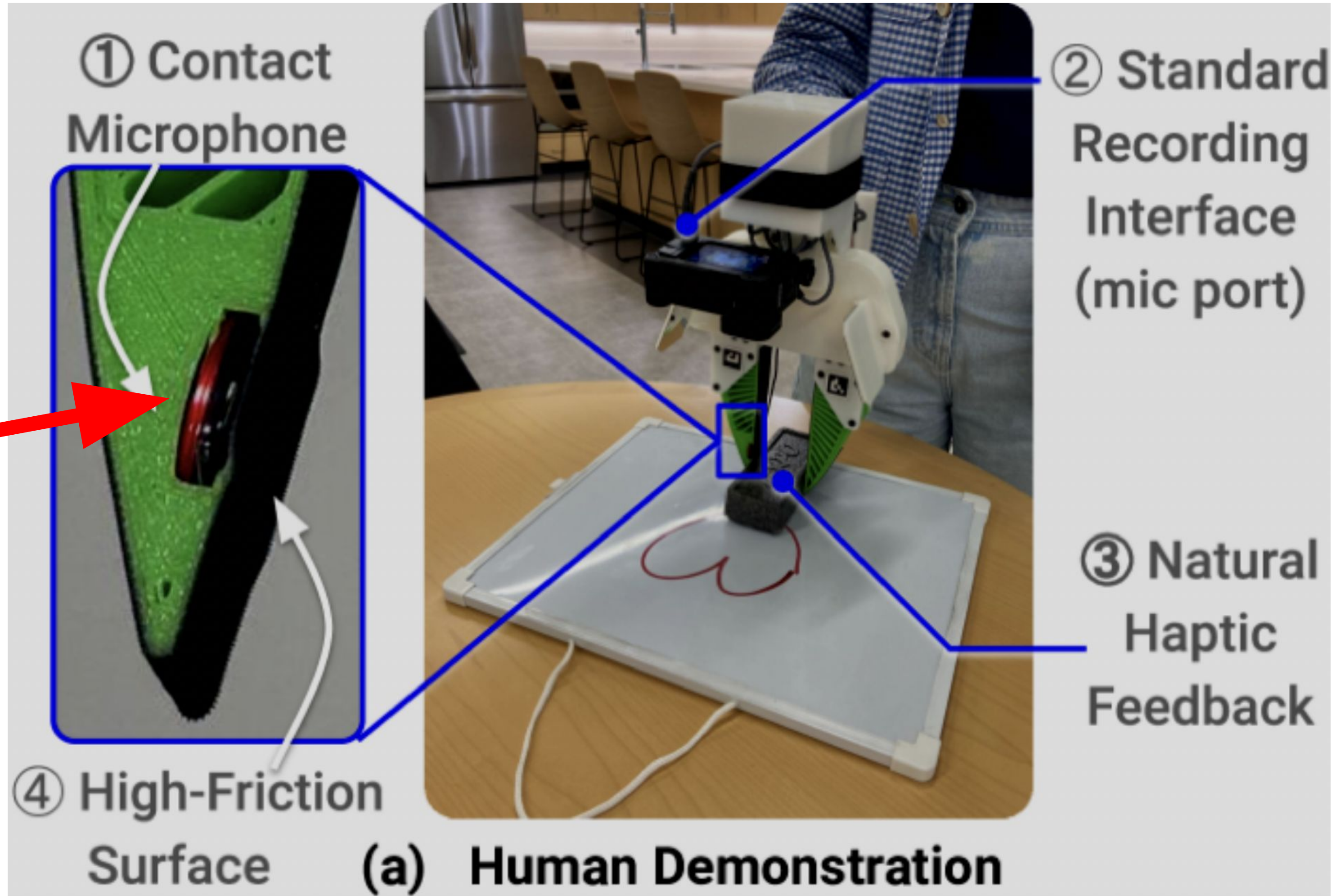
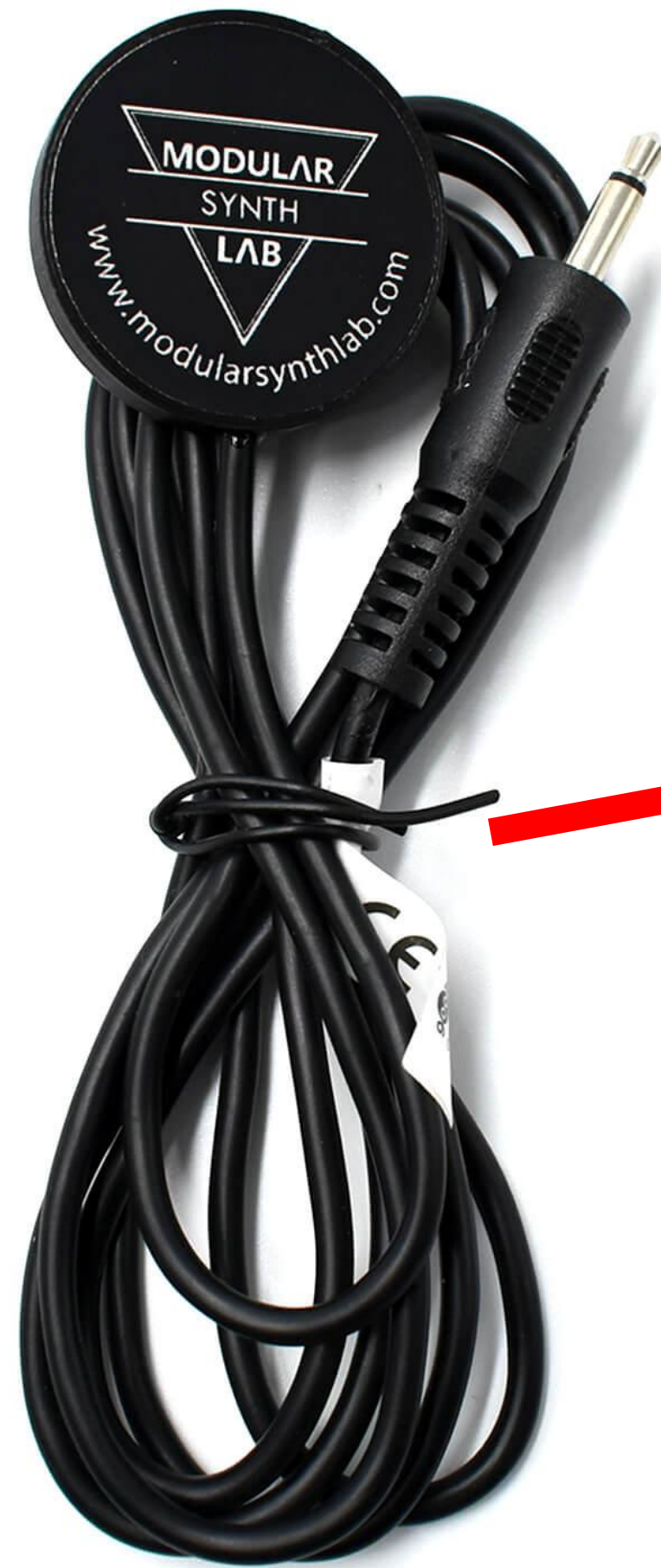
Contact Audio Sensing



<https://modularsynthlab.com/product/high-quality-piezo-contact-microphone-piezo-transducer-27mm-120cm-cable-mono-jack-3-5mm/?v=0b3b97fa6688>



Contact Audio Sensing



<https://modularsynthlab.com/product/high-quality-piezo-contact-microphone-piezo-transducer-27mm-120cm-cable-mono-jack-3-5mm/?v=0b3b97fa6688>

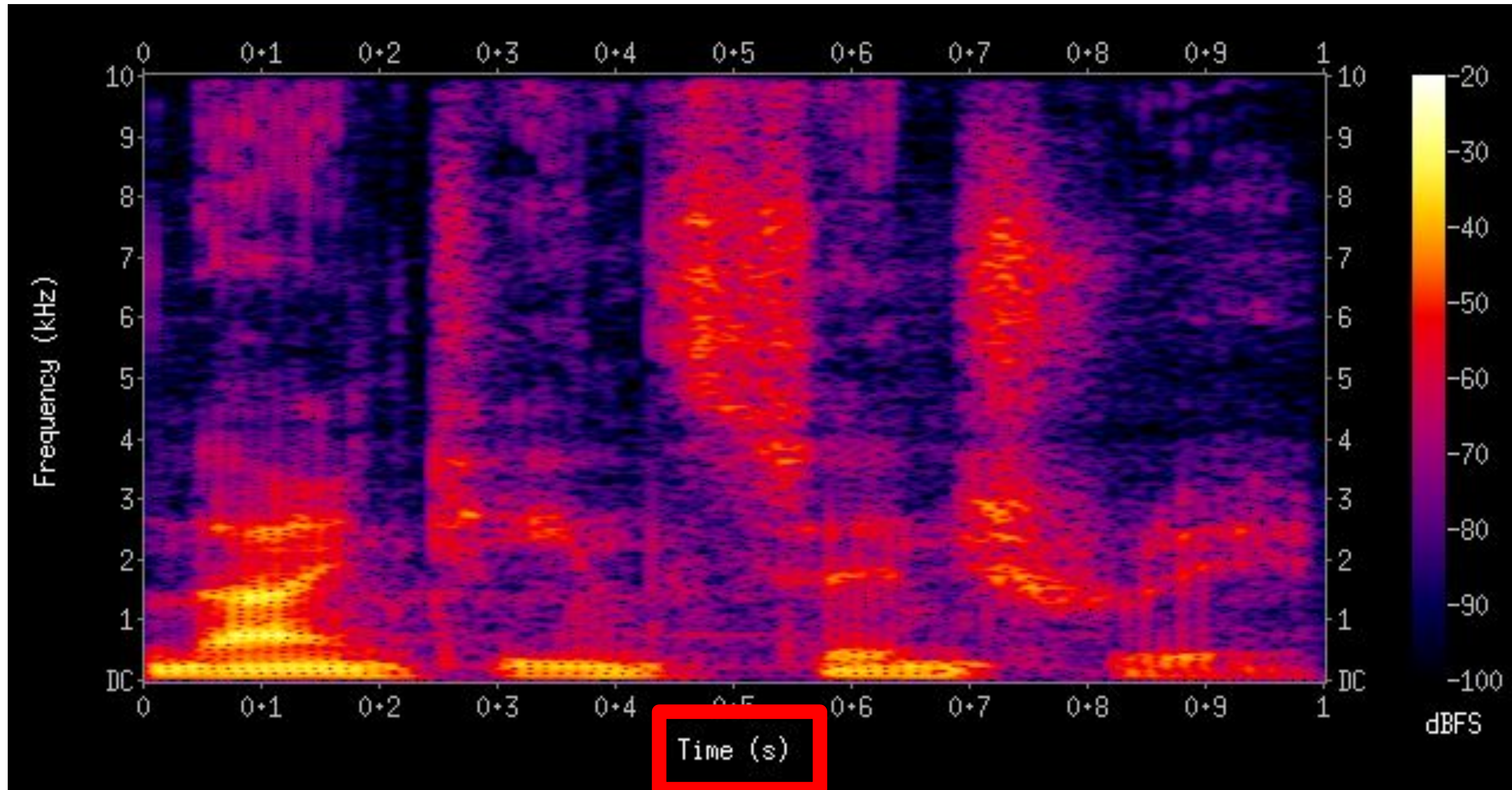
Z. Liu, C. Chi, E. Cousineau, N. Kuppaswamy, B. Burchfiel, and S. Song, "Maniway: Learning robot manipulation from in-the-wild audio-visual data," arXiv preprint arXiv:2406.19464, 2024.





Contact Audio Sensing

Spectrograms: Transferring audio into the vision domain



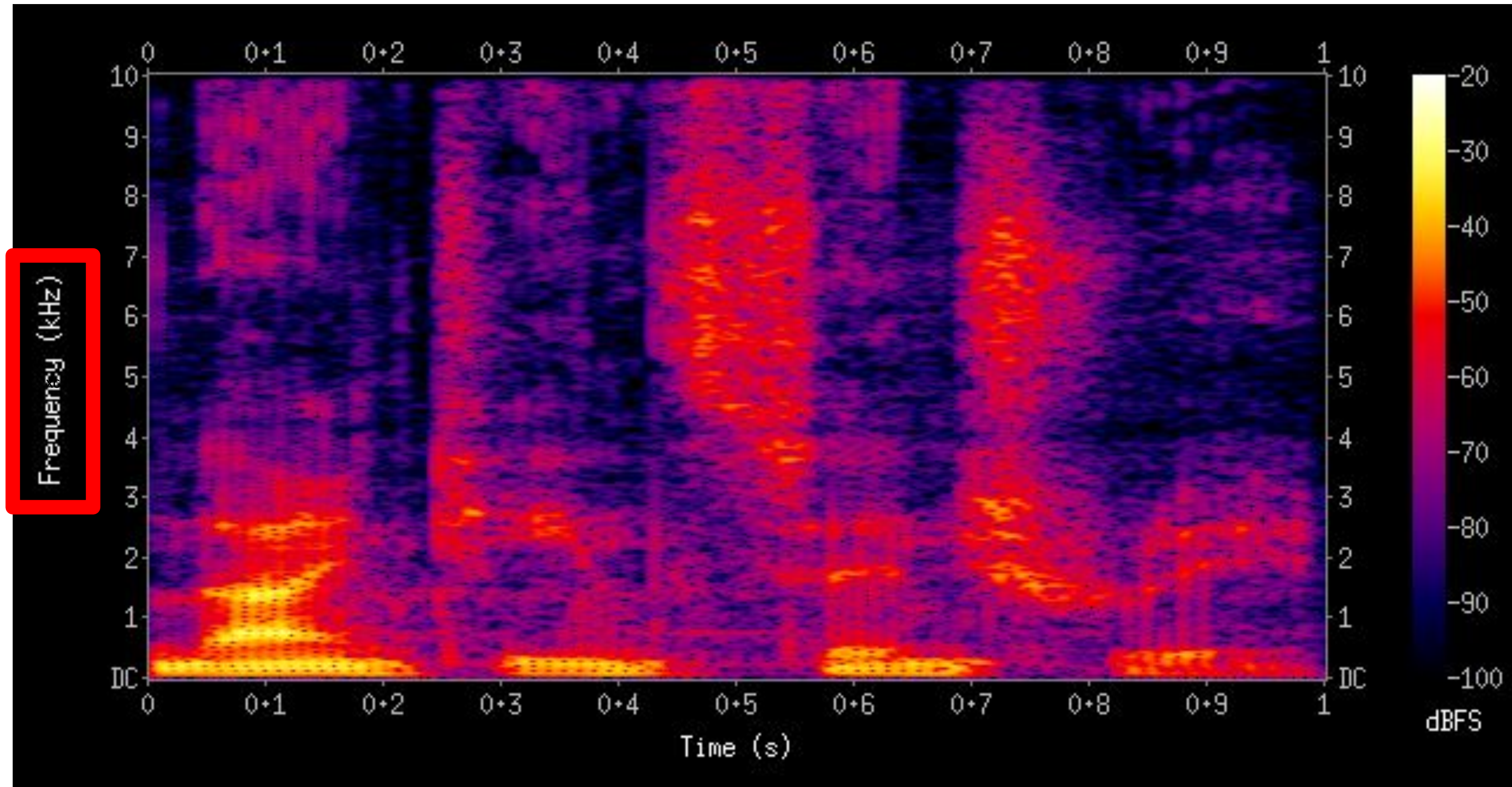
1. Record an audio sample over time (i.e. last few seconds from contact microphone)





Contact Audio Sensing

Spectrograms: Transferring audio into the vision domain



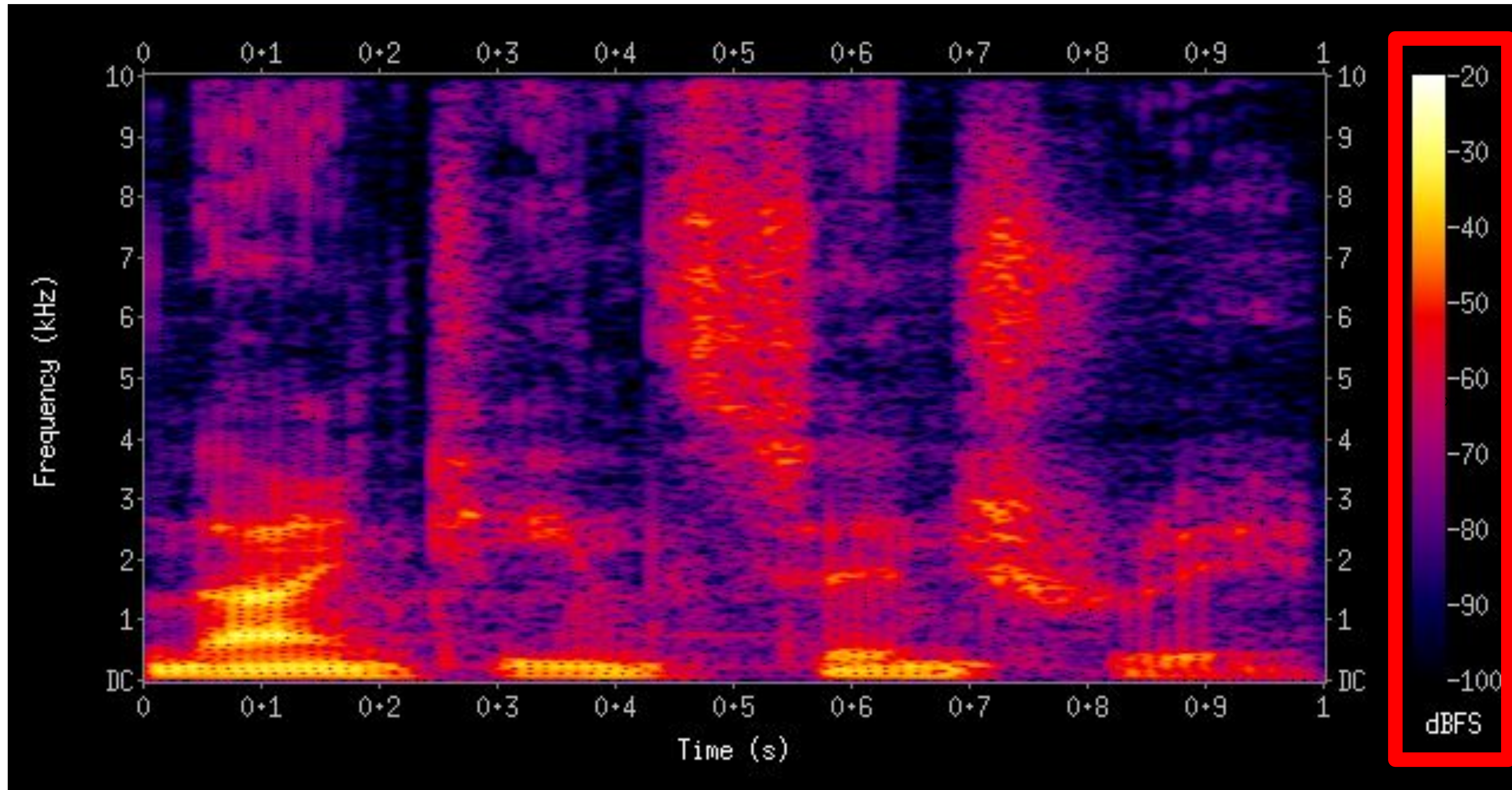
2. Extract frequencies for each time step





Contact Audio Sensing

Spectrograms: Transferring audio into the vision domain



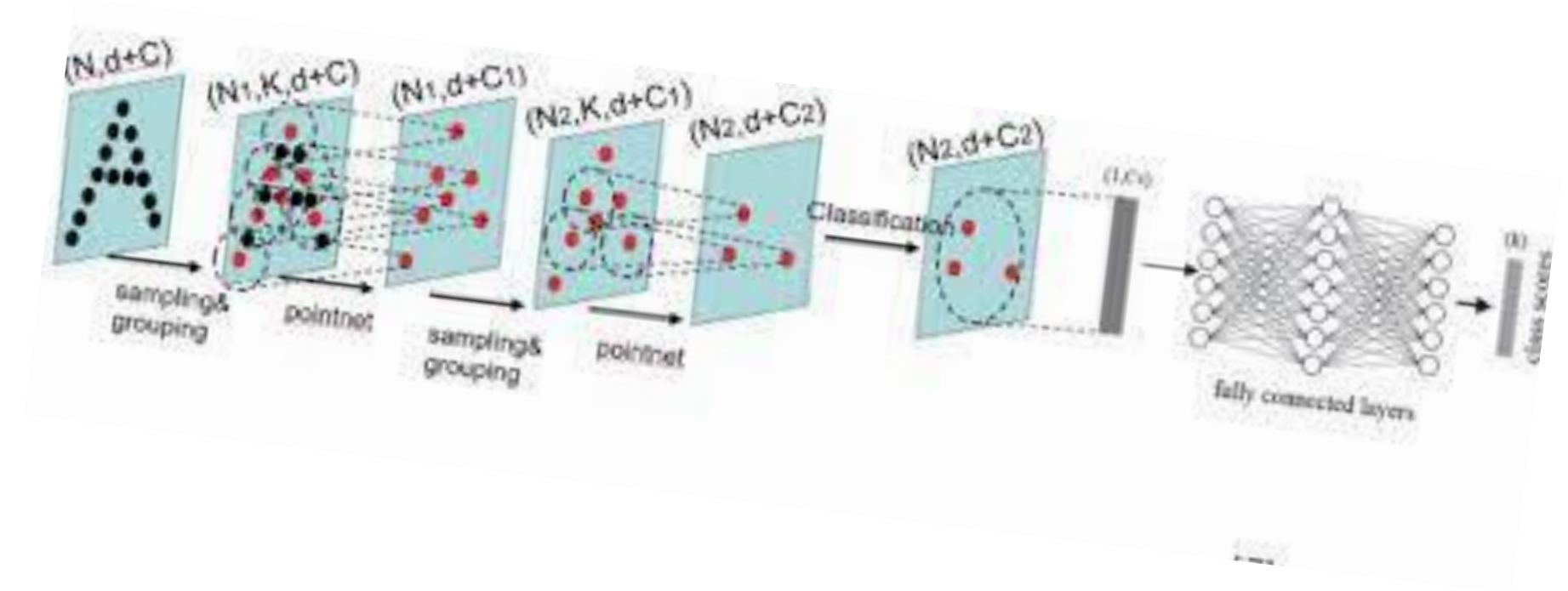
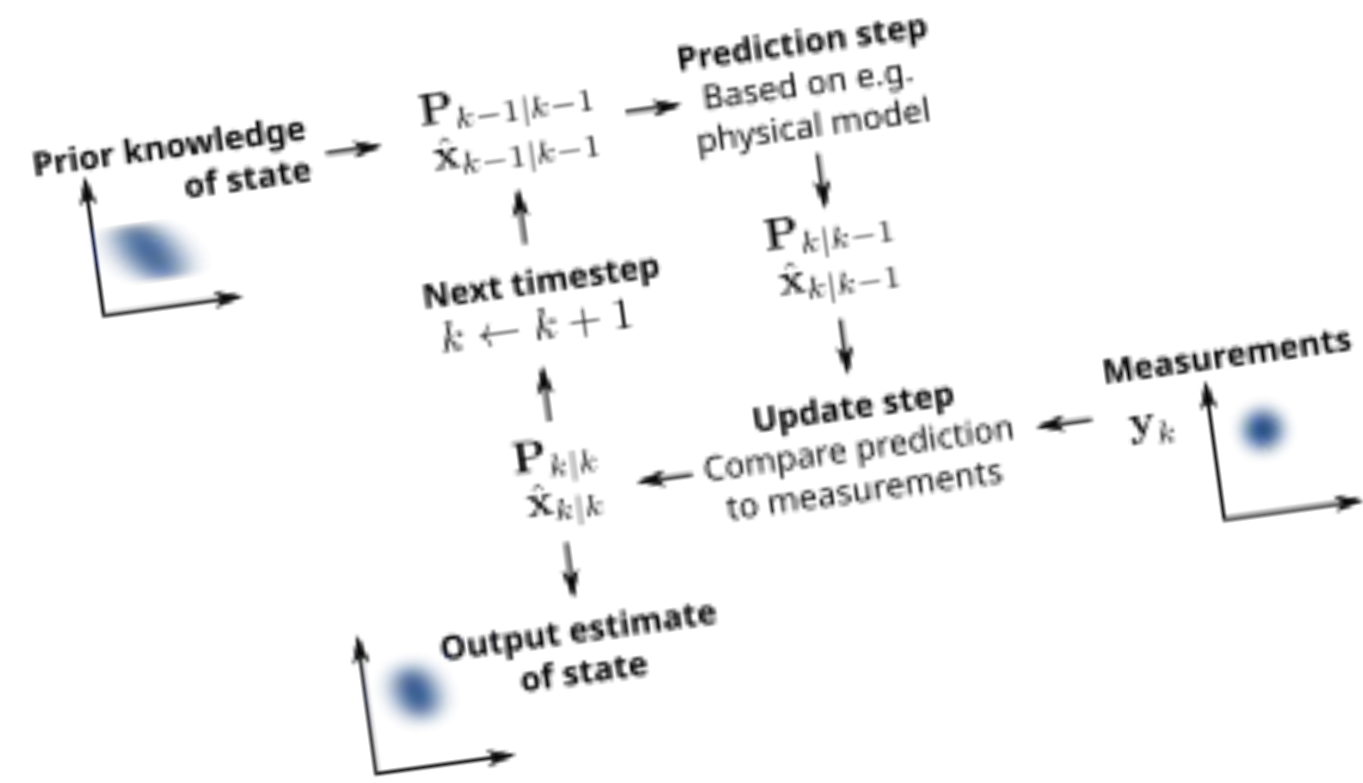
3. Plot amplitudes of each frequency for each timestep



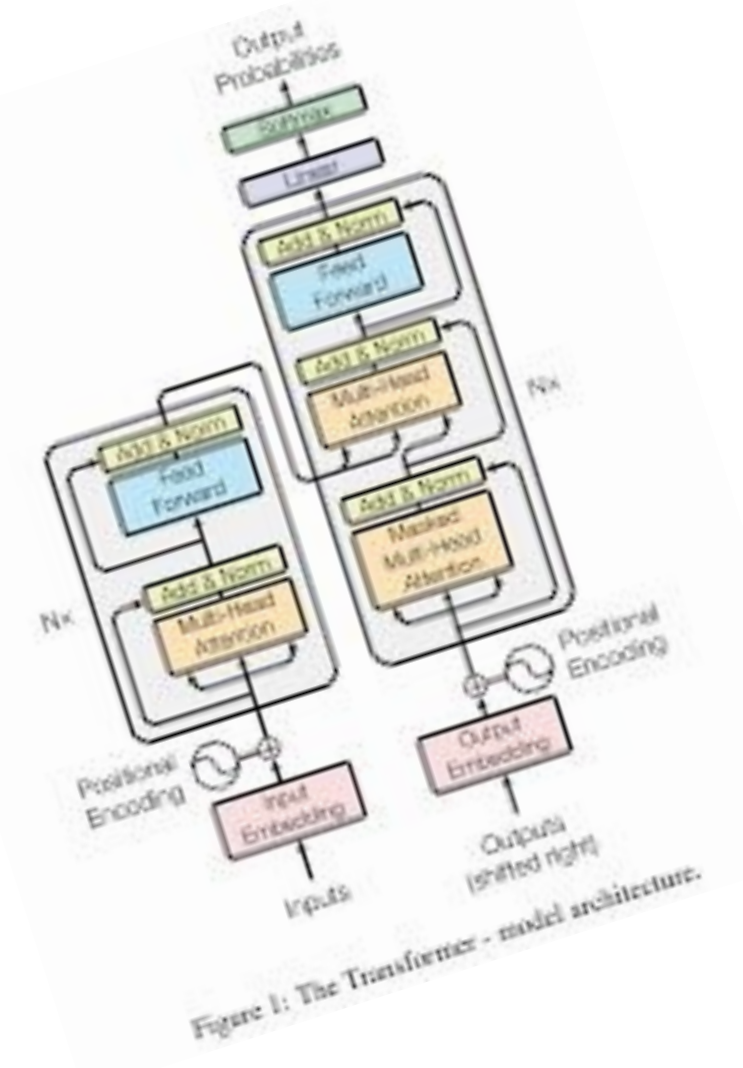
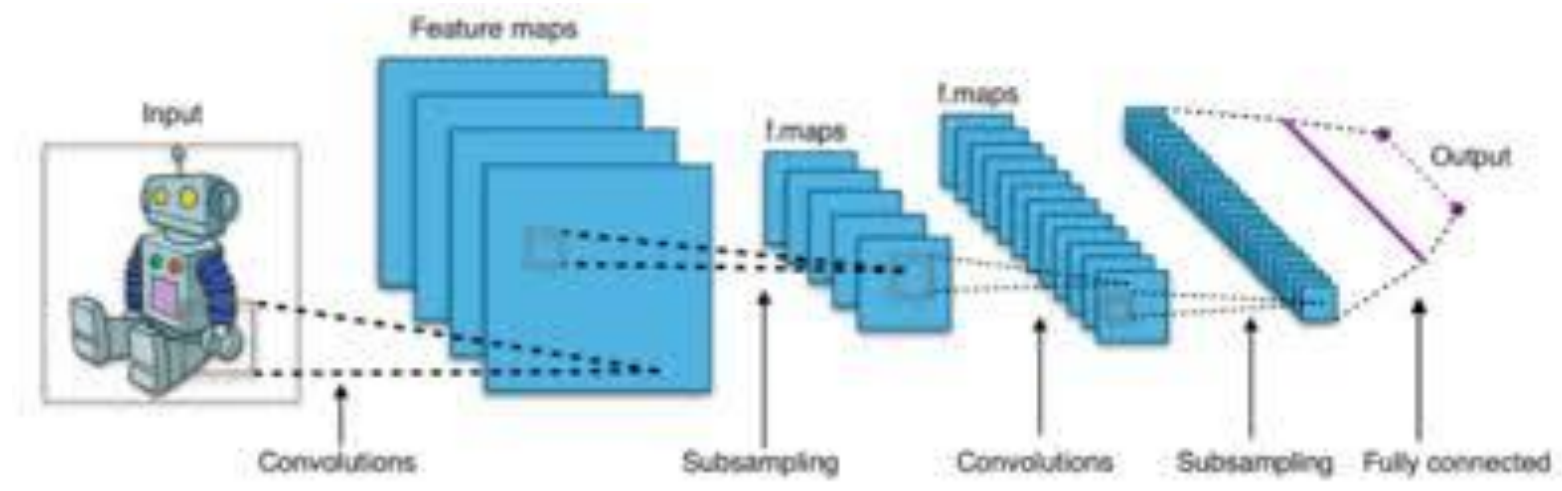
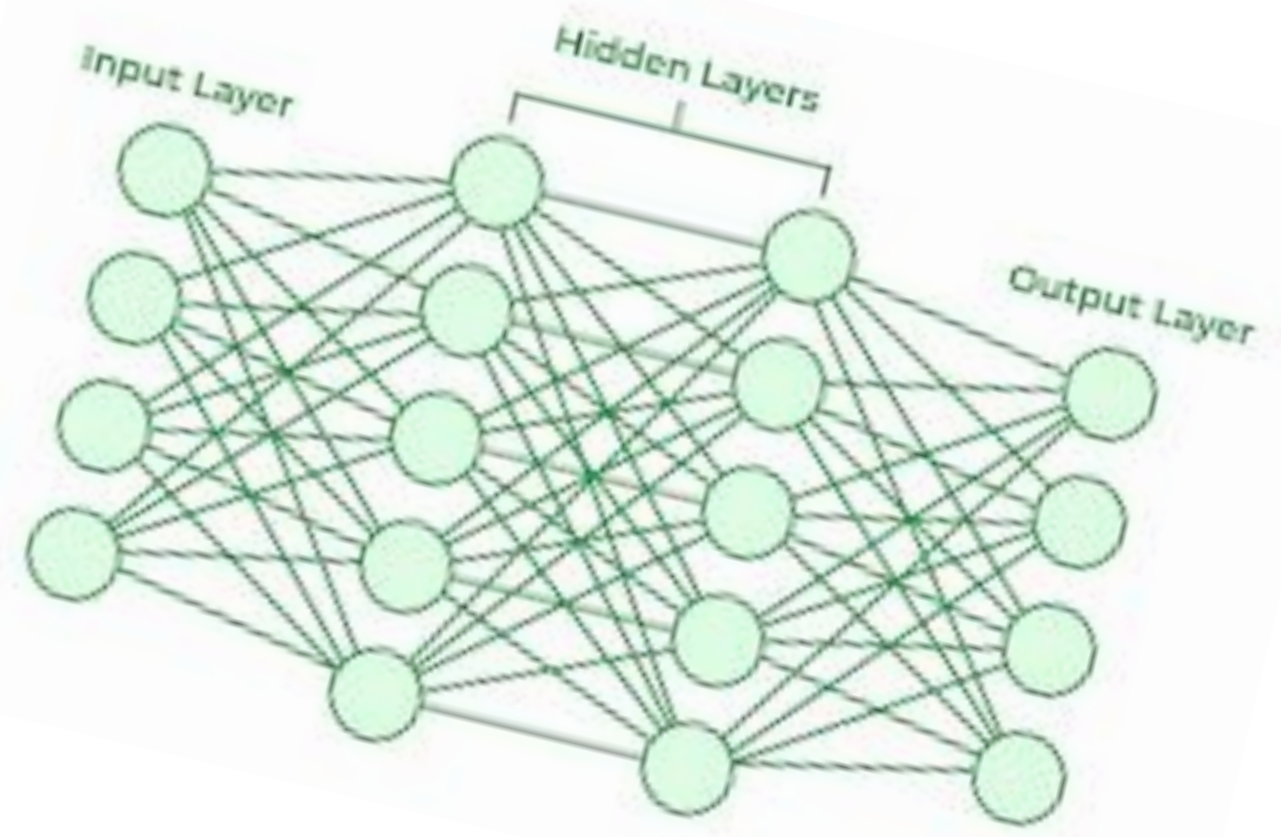


These are not all the types of sensors...

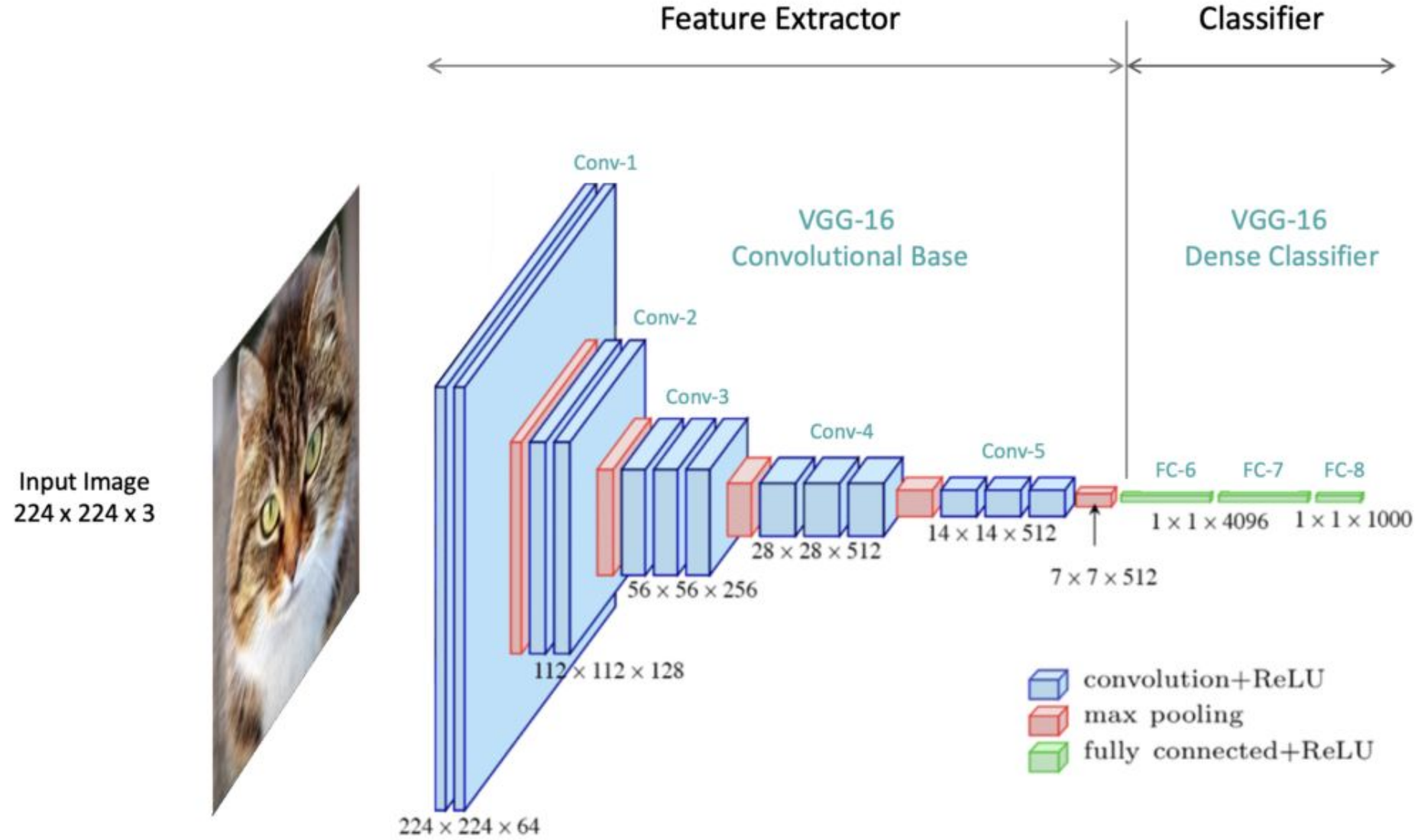




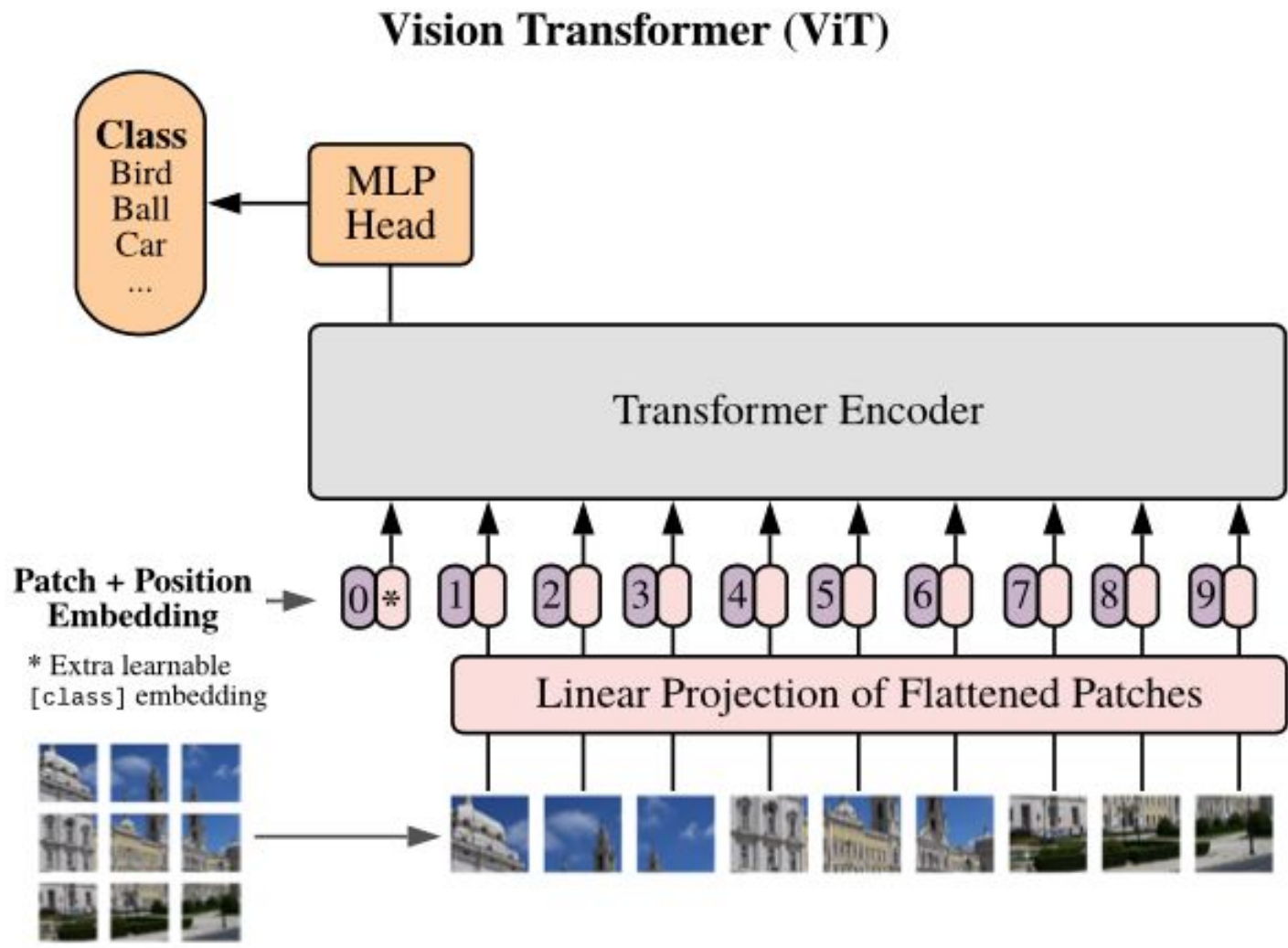
How can we encode each sensor data type?



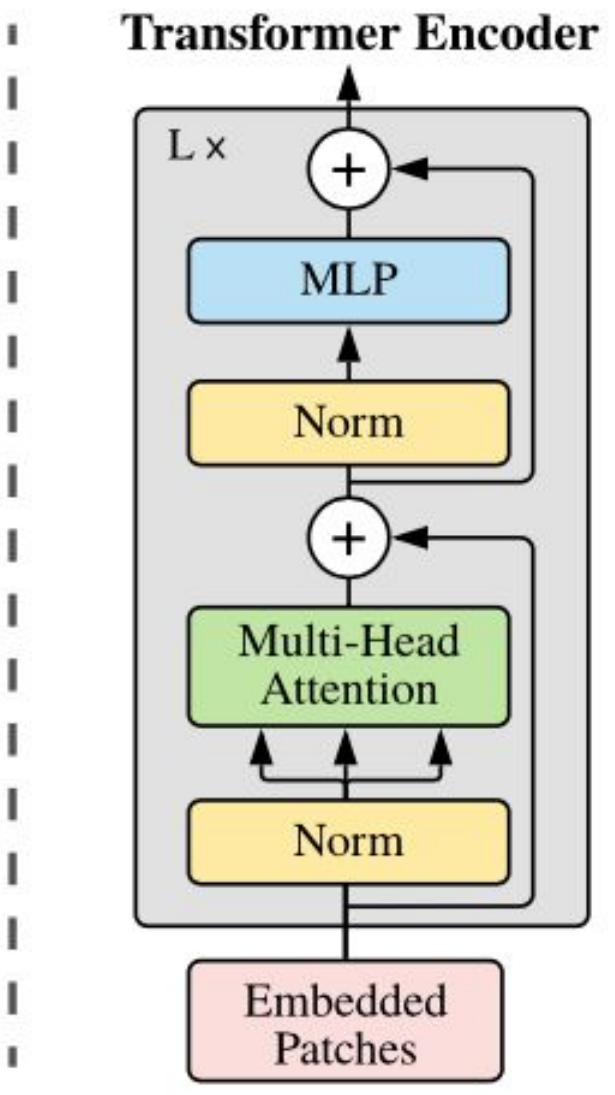
Classifying 2D Maps



CNNs



ViTs

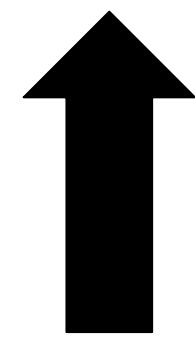
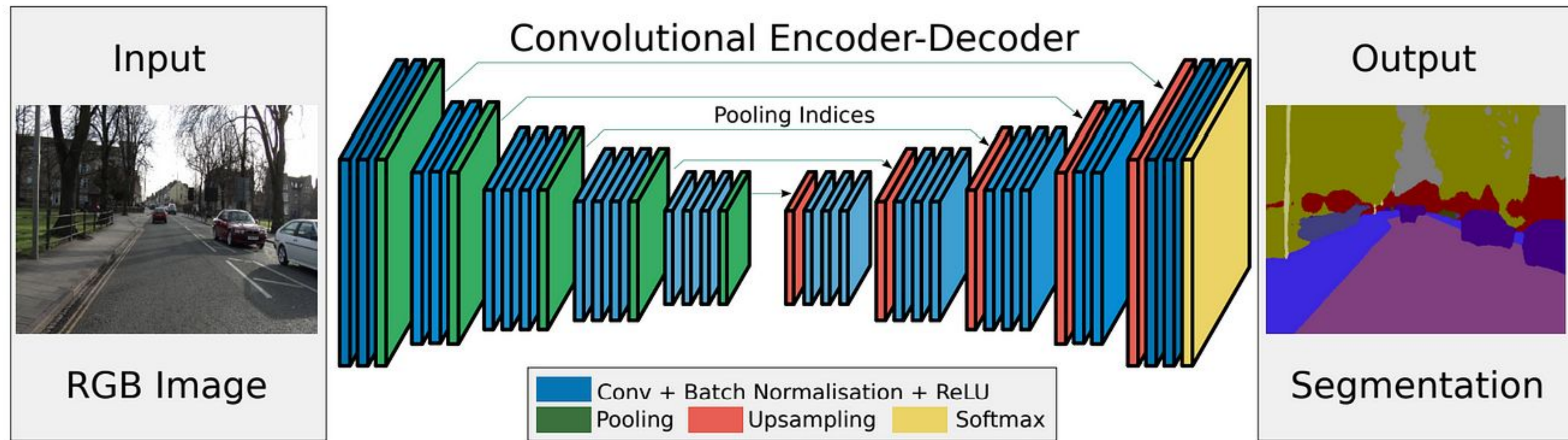


Credit: Learn OpenCV and Unite.AI



Encoding 2D Maps

Task Encoding



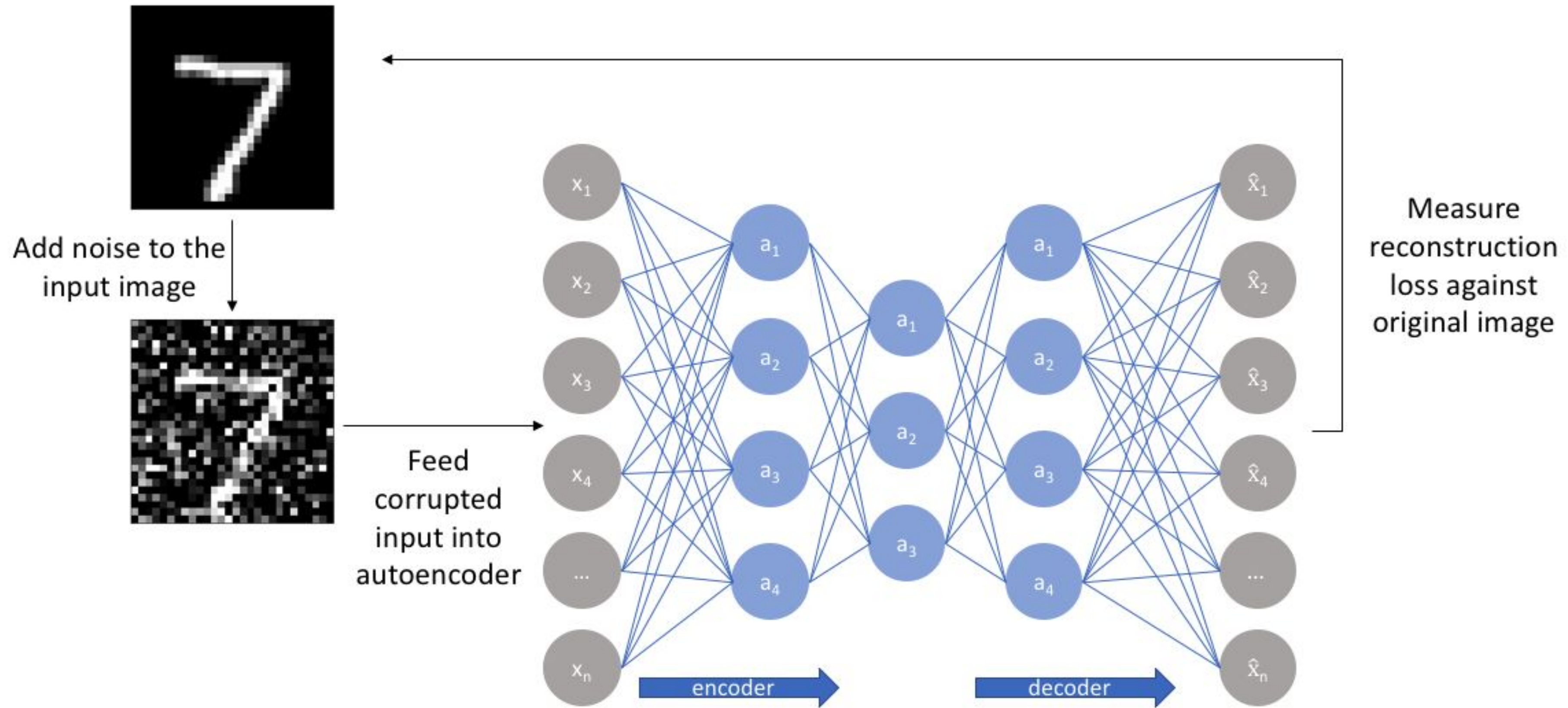
Latent Vector

Credit: Jeremy Jordan



Encoding 2D Maps

Autoregressive Encoding

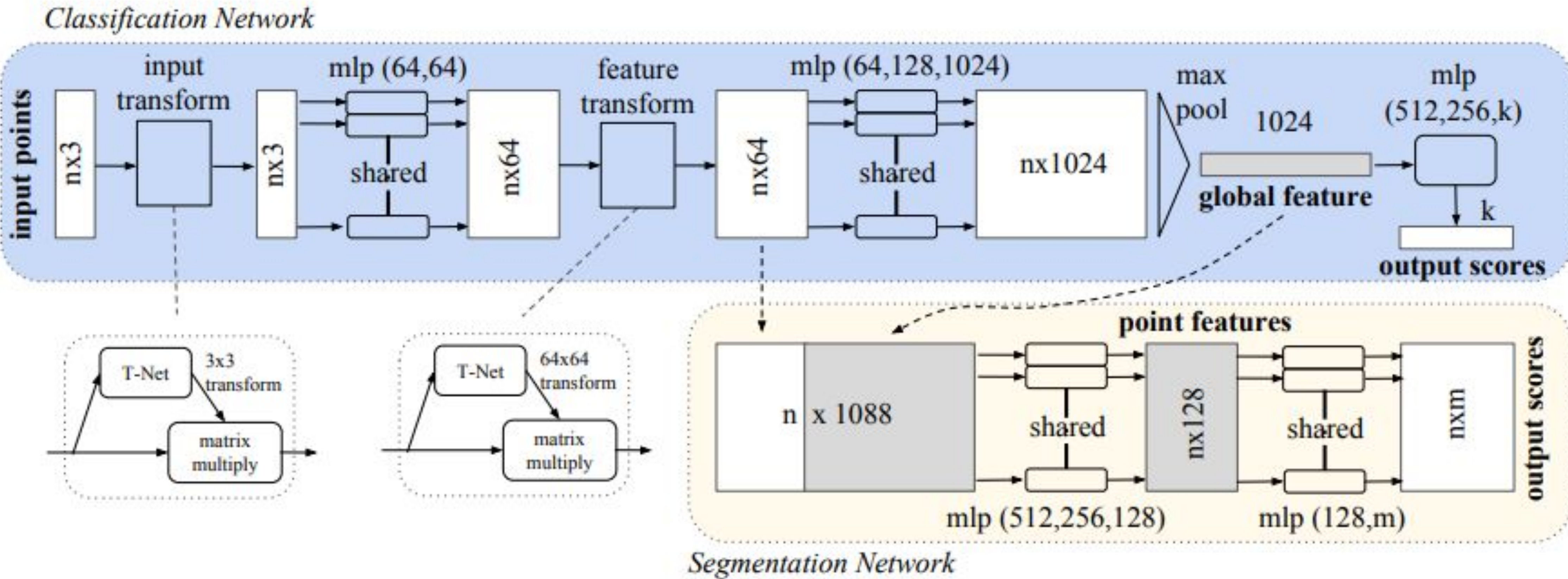


Credit: Jeremy Jordan



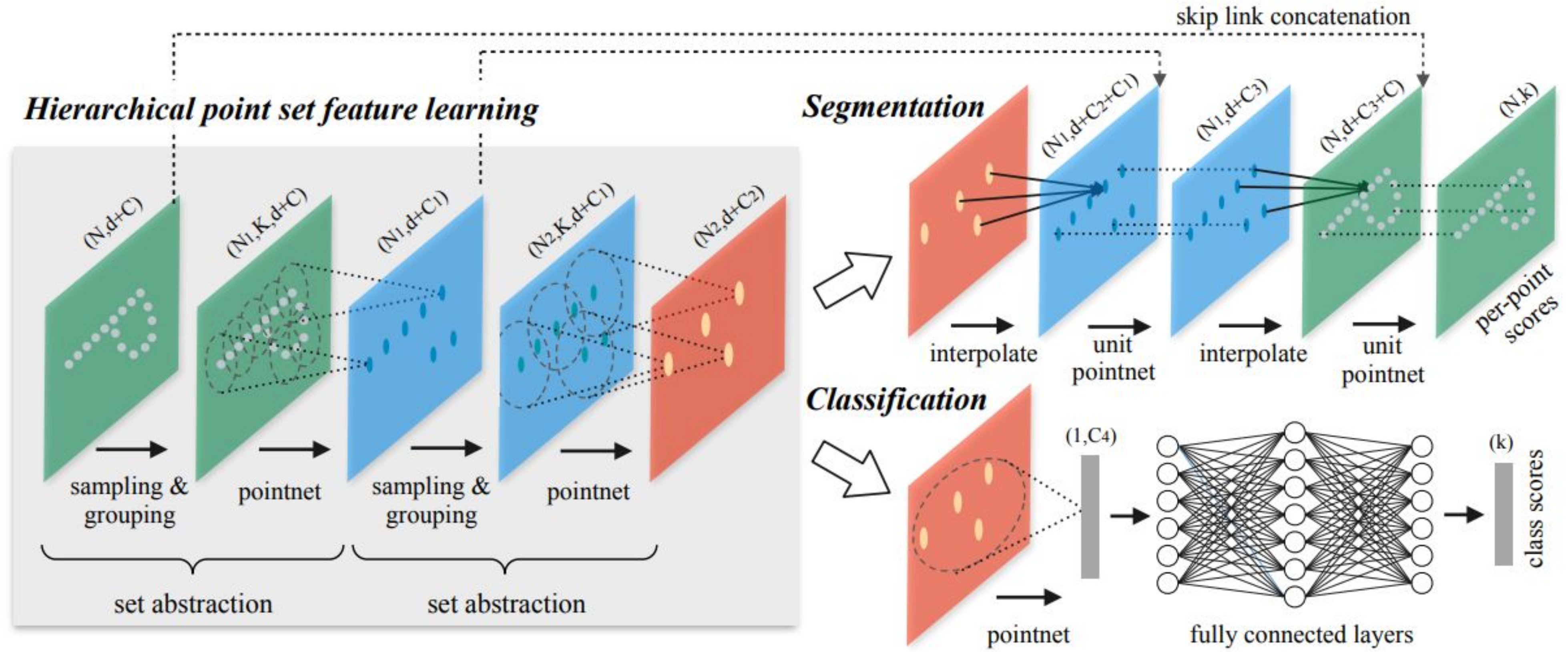


How to process point cloud data—pointnet





Pointnet++



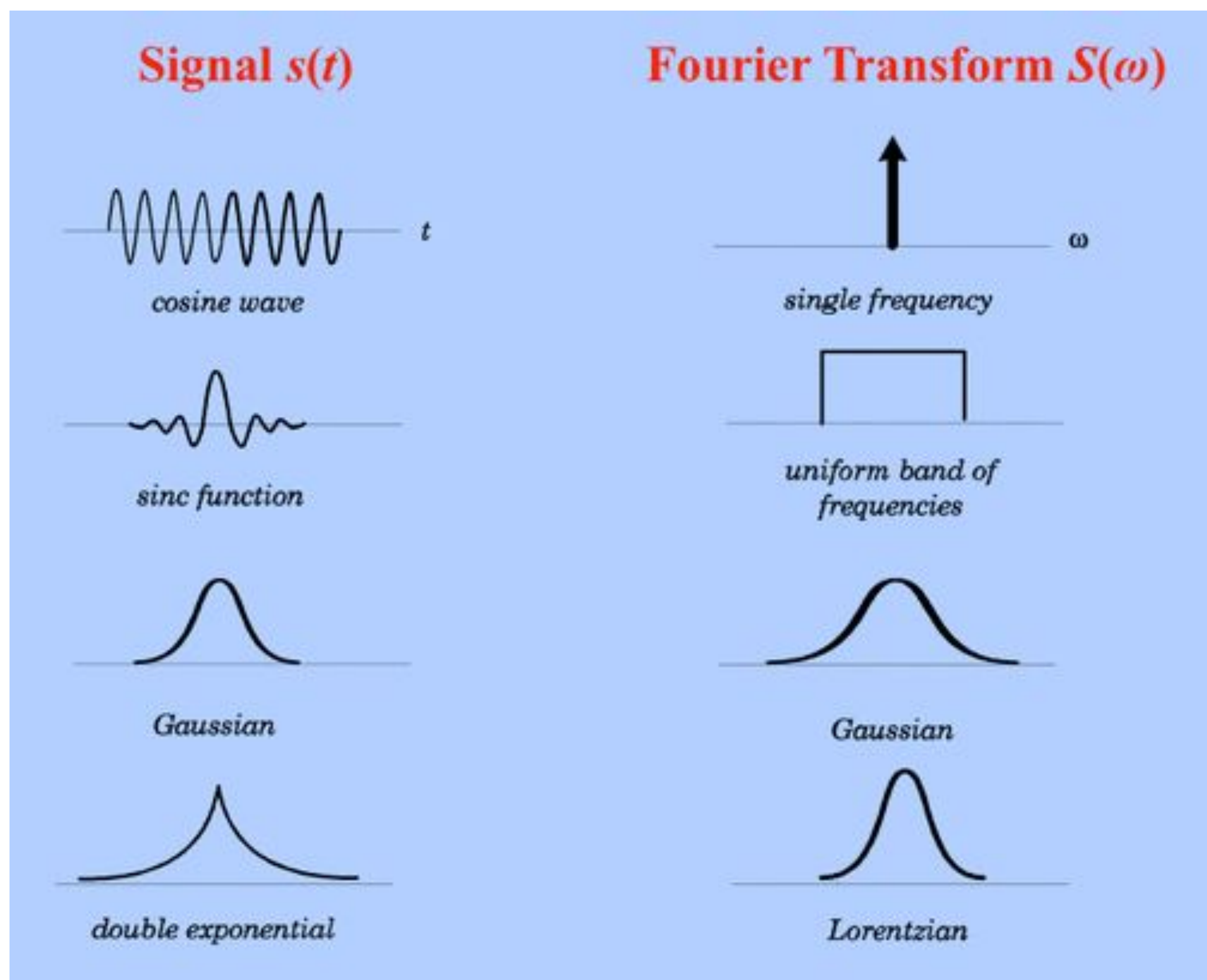
capture local geometric features

1 .hierarchical structure

2. multi-scale feature aggregation

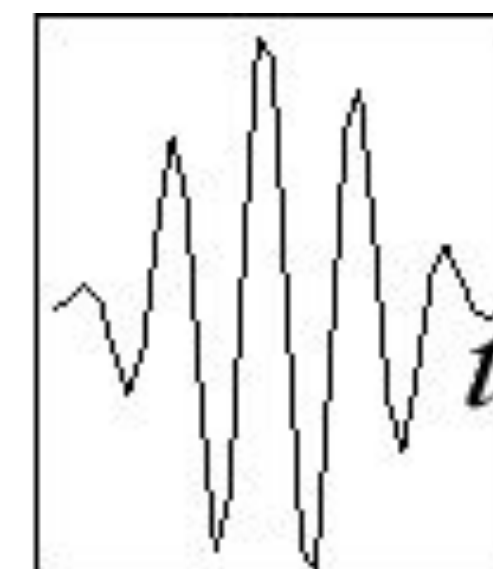


Processing Time Series Data

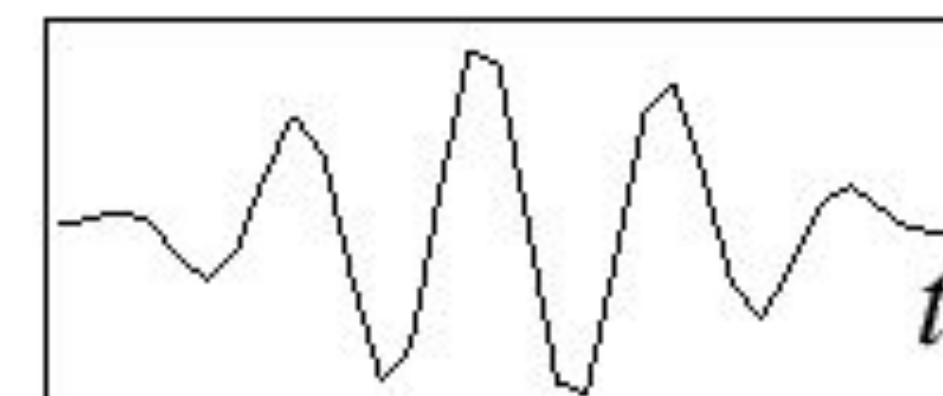


Fourier Transform

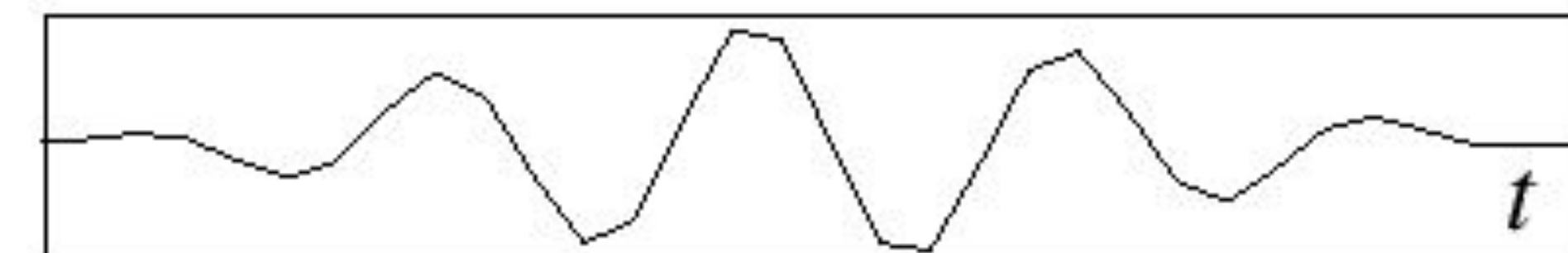
$f(t)$
(Mother Wavelet)



$\frac{1}{\sqrt{2}} f\left(\frac{t}{2}\right)$



$\frac{1}{2} f\left(\frac{t}{4}\right)$

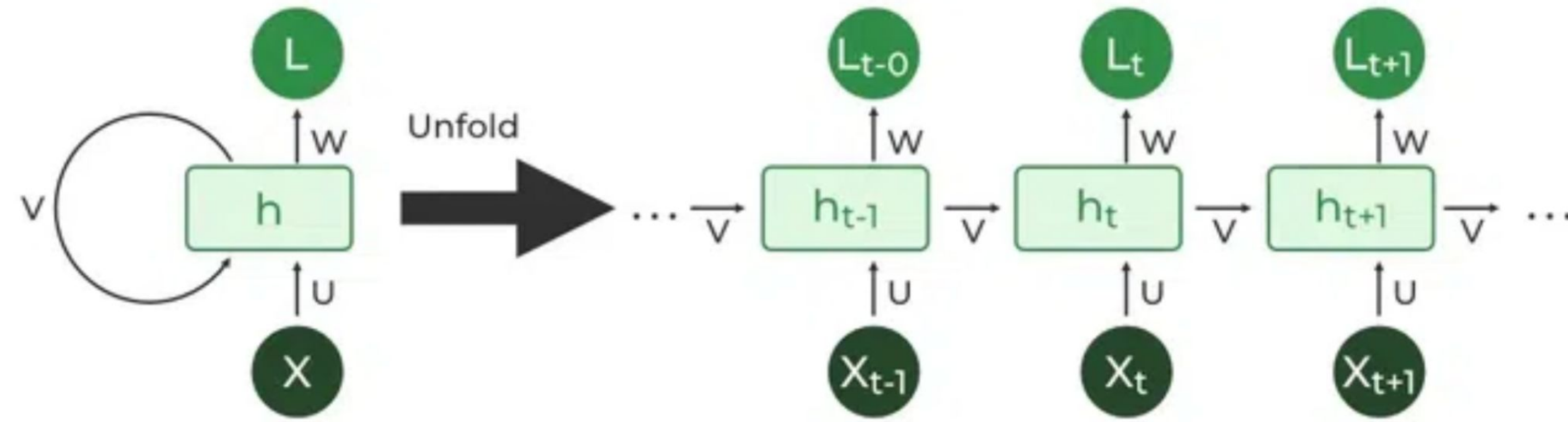


Wavelet Transform





Processing Time Series Data



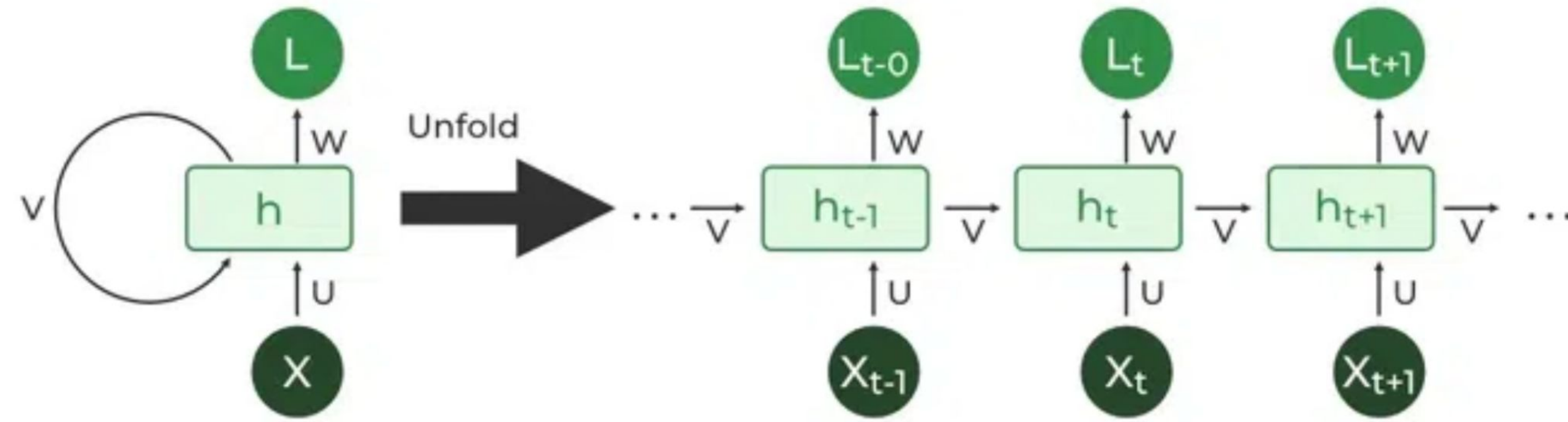
RNNs / Attention

<https://www.geeksforgeeks.org/introduction-to-recurrent-neural-network/#>



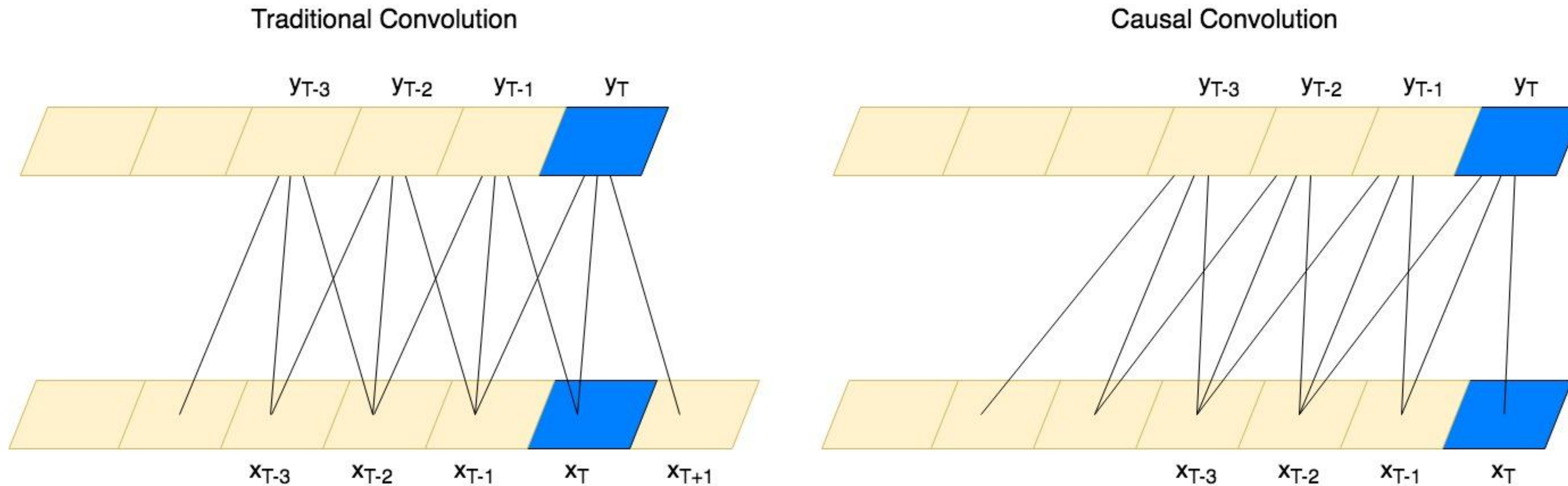


Processing Time Series Data



RNNs / Attention

<https://www.geeksforgeeks.org/introduction-to-recurrent-neural-network/#>



Causal Convolution

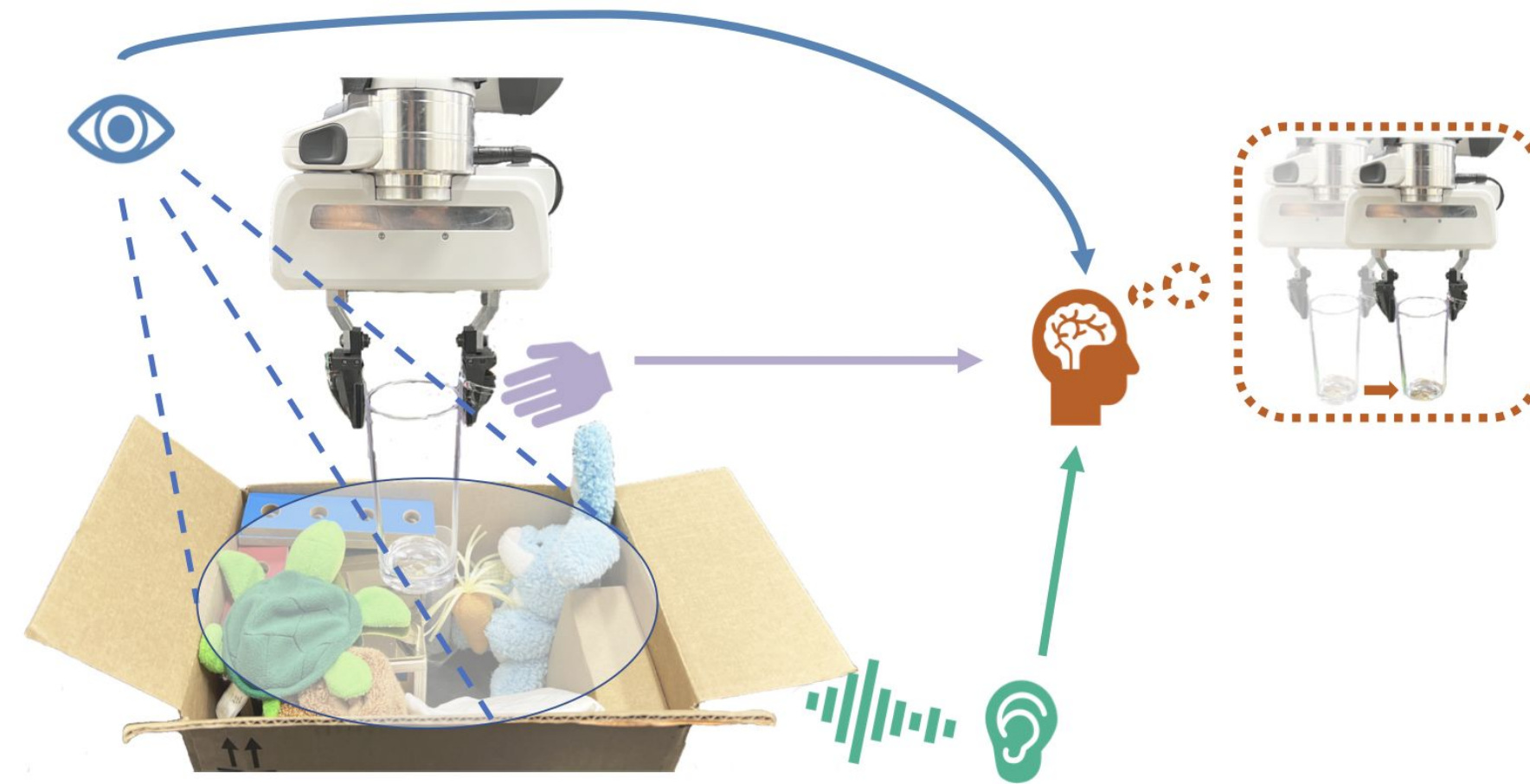
<https://blog.fastforwardlabs.com/2018/05/31/convolve-all-the-things.html>



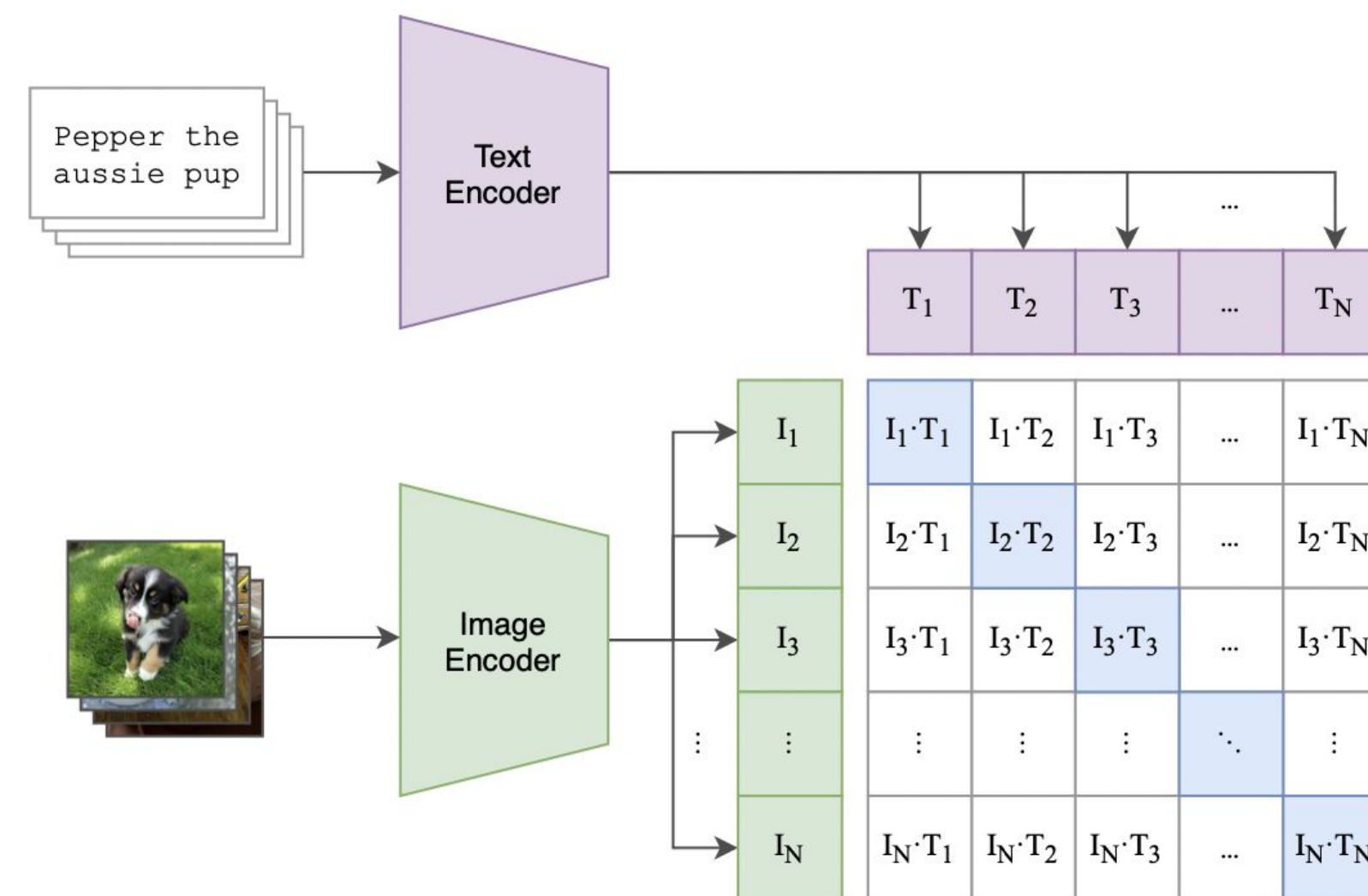


Questions?





What is the difference between *Multisensory* vs. *Multimodal* ?





Putting it all together



Making Sense of Vision and Touch: Self-Supervised Learning of Multimodal Representations for Contact-Rich Tasks

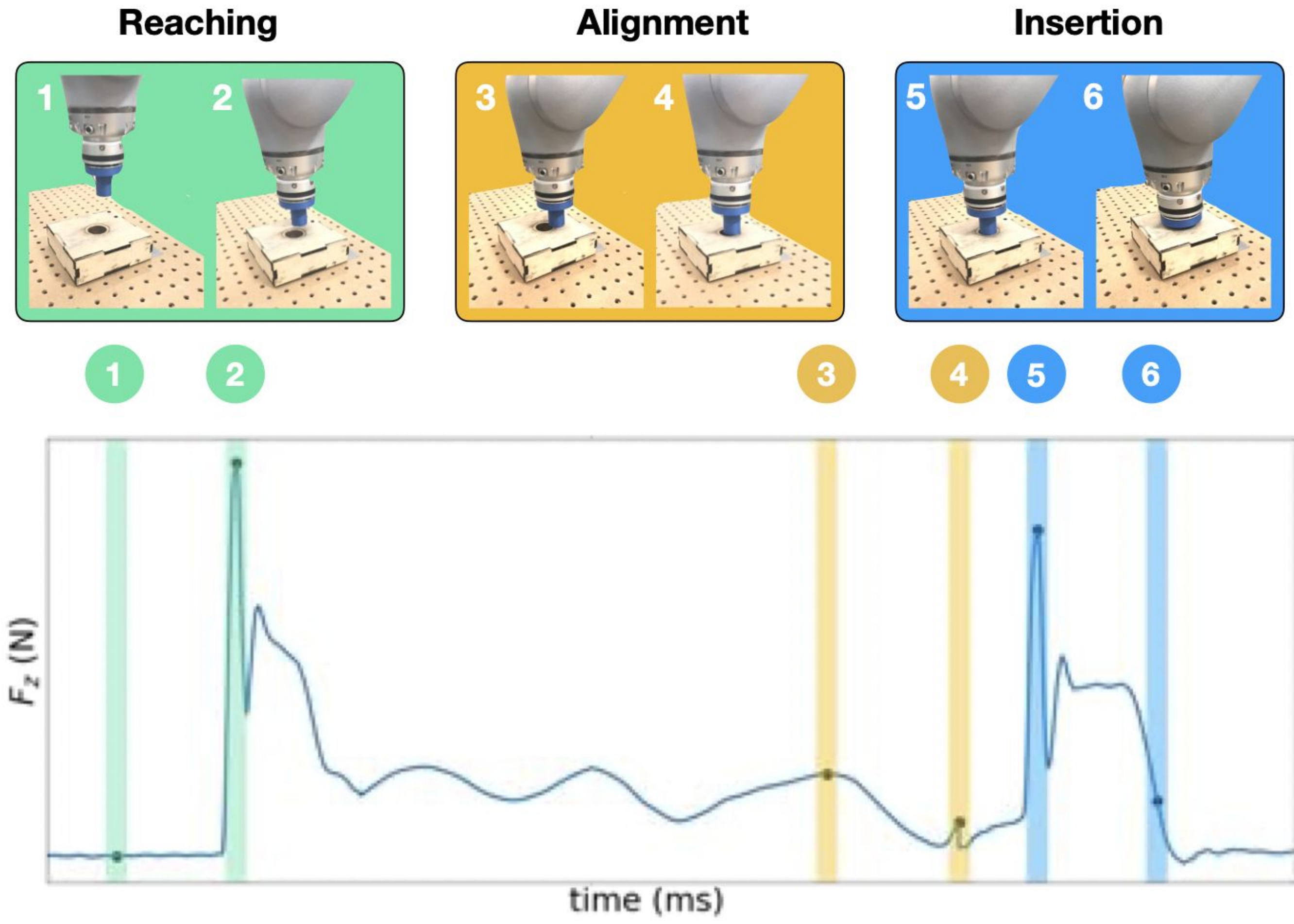
Michelle A. Lee, Yuke Zhu, Krishnan Srinivasan, Parth Shah, Silvio
Savarese, Li Fei-Fei, Animesh Garg, Jeannette Bohg

ICRA 2019 [Best Paper Award]

T-RO 2020 [Extended Version]

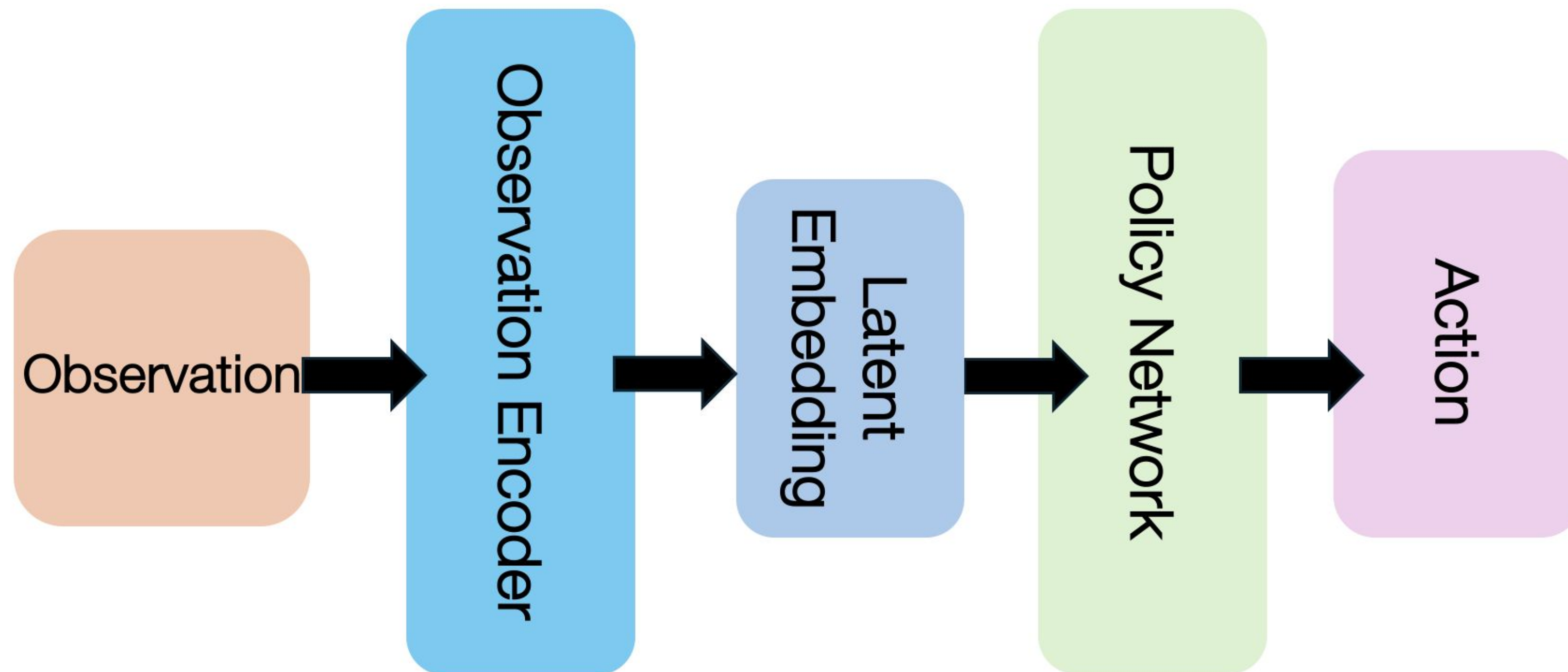


How can we learn good latent representations for contact-rich tasks?



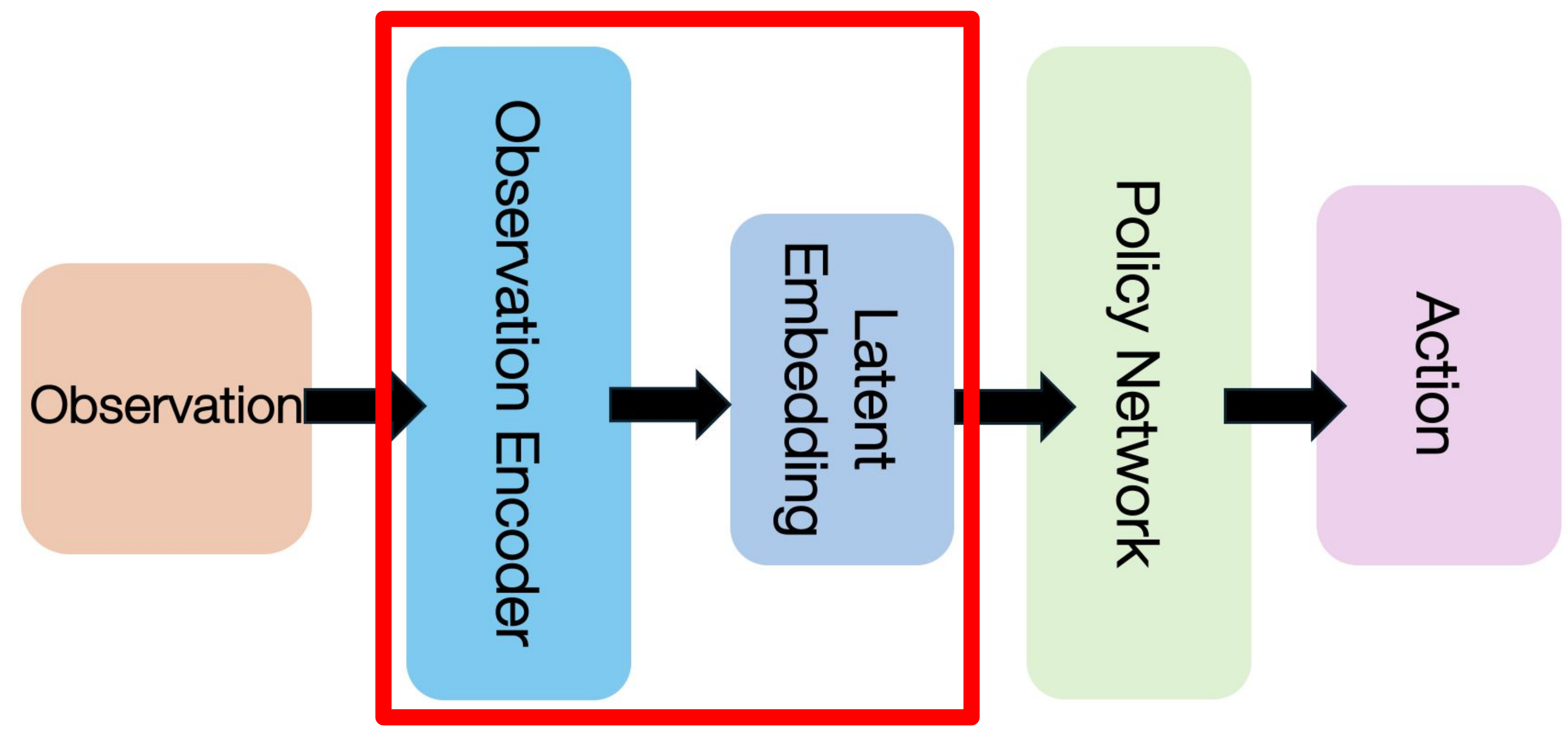


Idea: Decouple representation and policy learning



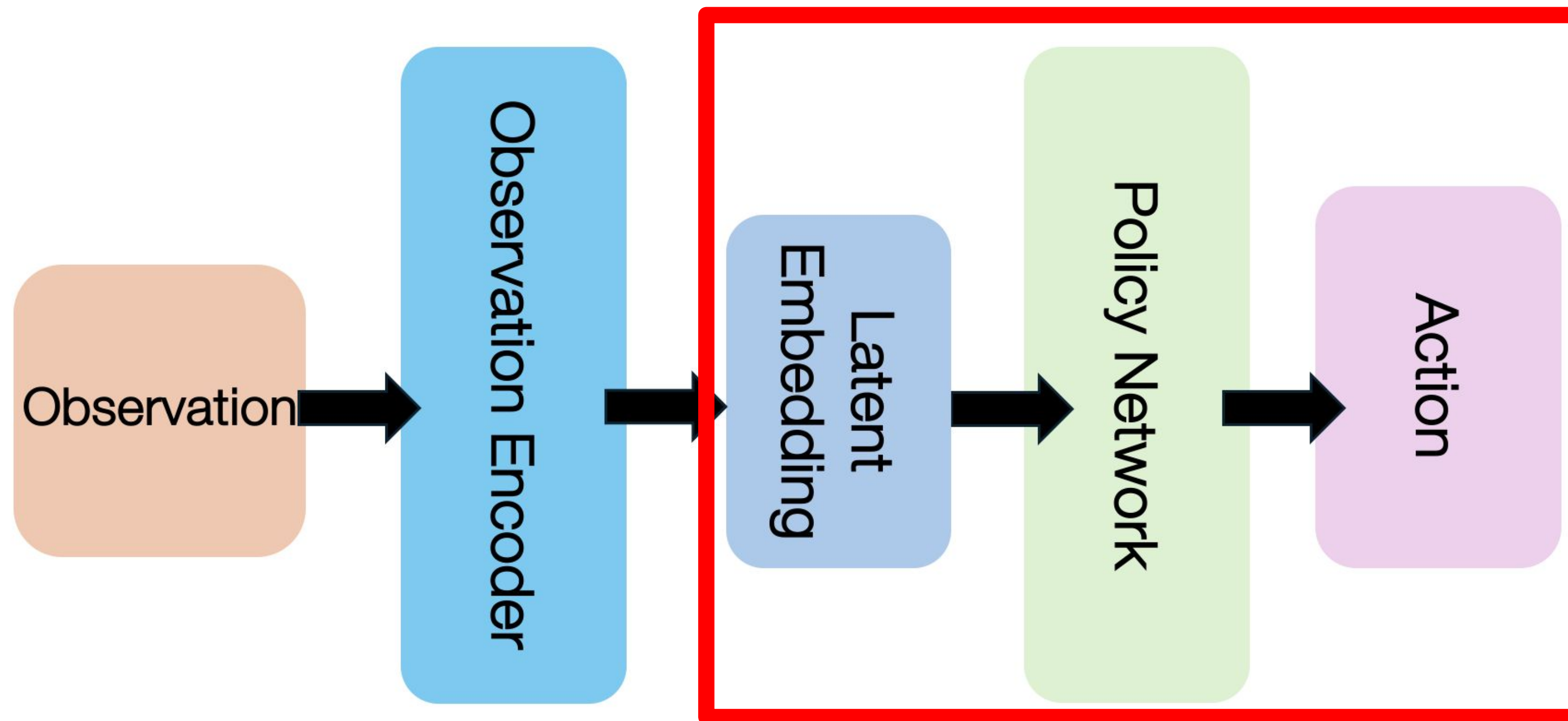
Idea: Decouple representation and policy learning

- 1. Learn latent embedding space through self-supervised learning



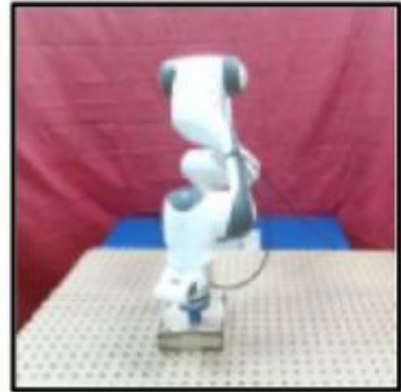
Idea: Decouple representation and policy learning

1. Learn latent embedding space through self-supervised learning
2. Use pretrained representation for policy learning

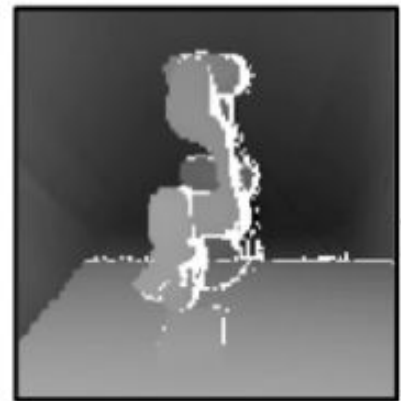




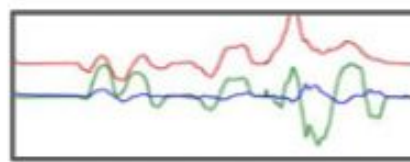
1. Learn latent embedding space through self-supervised learning



RGB Camera



Depth



Force-Torque Sensor



Proprioception

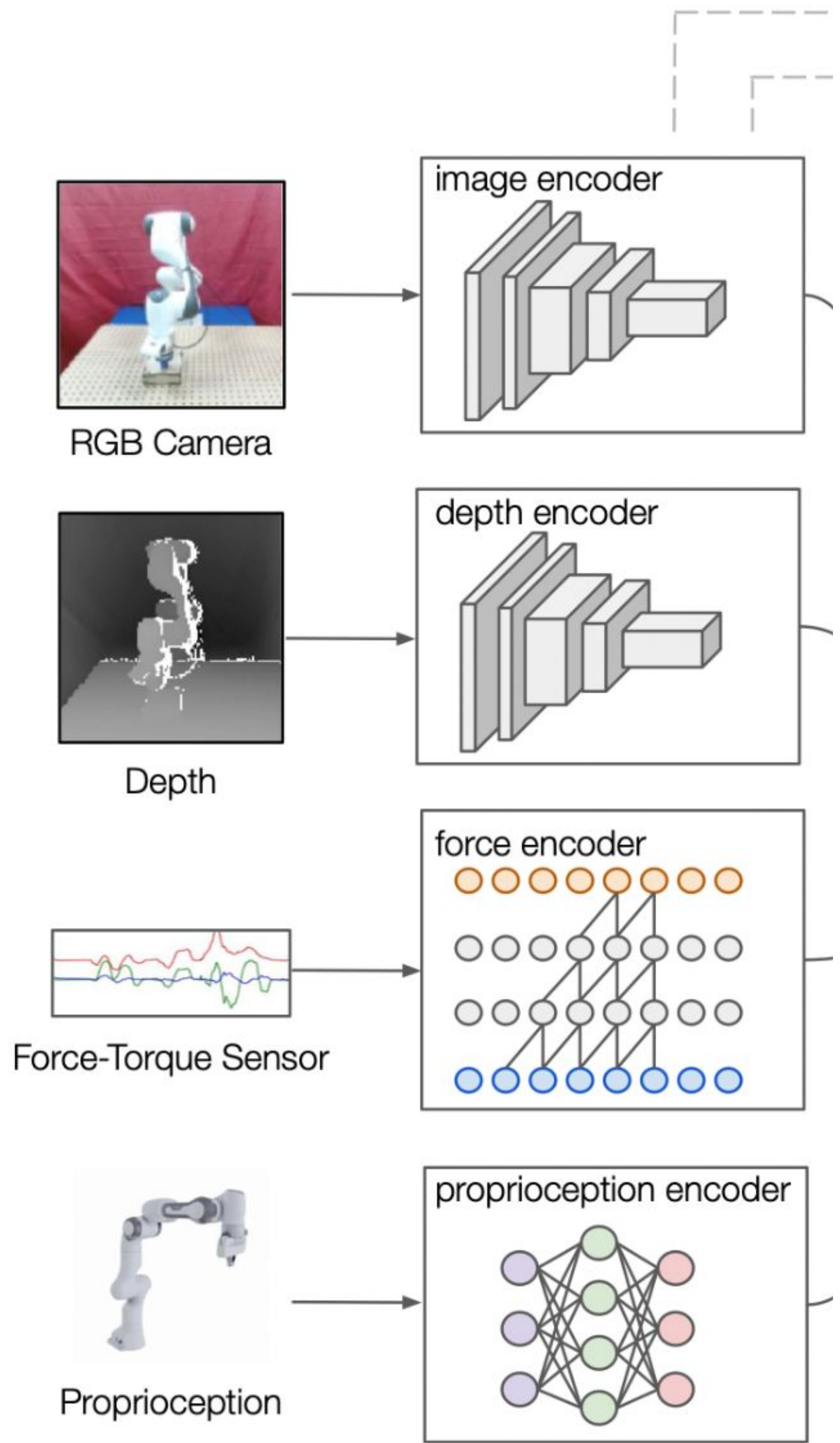
Sensing modes





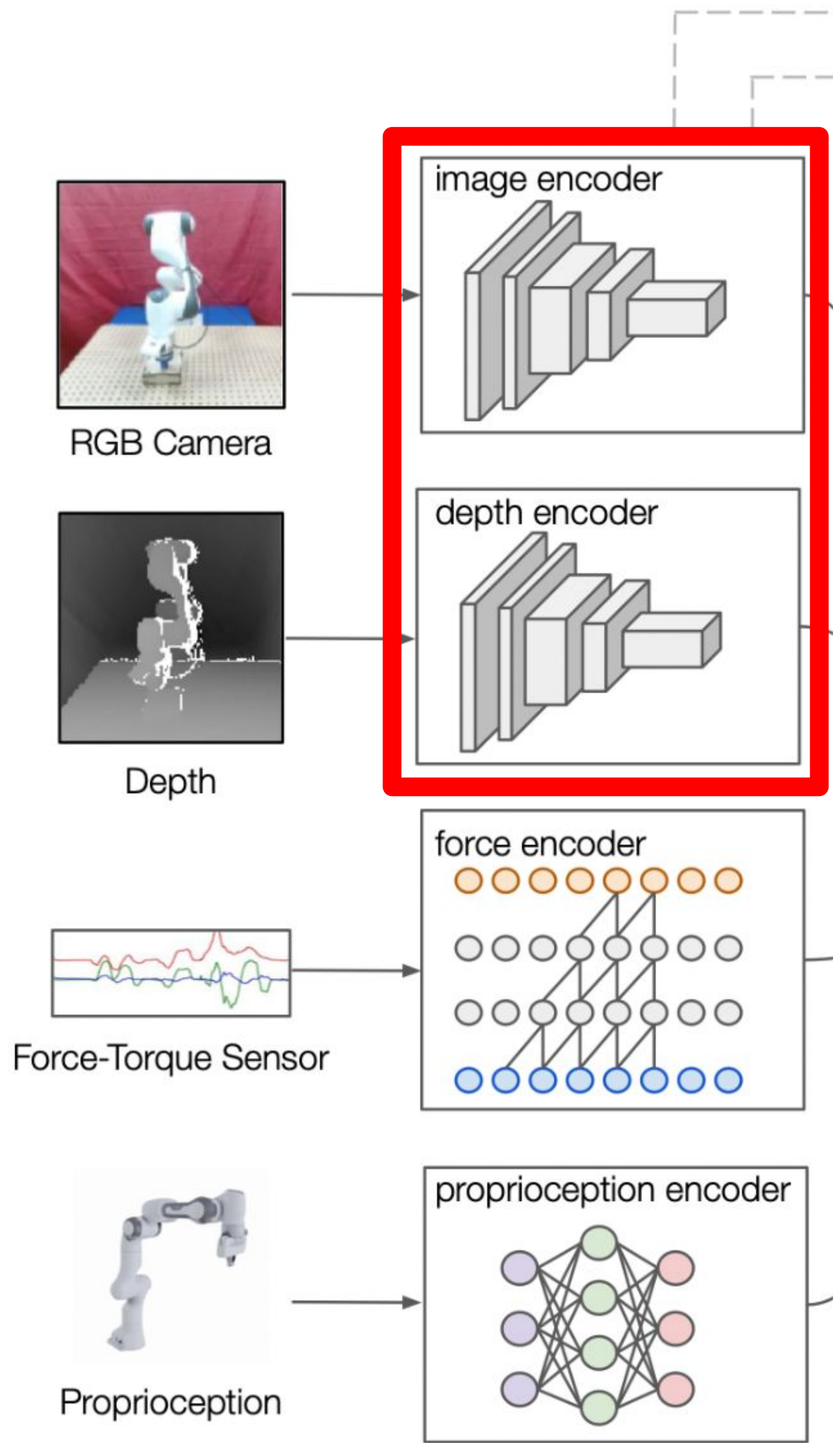
1. Learn latent embedding space through self-supervised learning

Domain-specific encoders





1. Learn latent embedding space through self-supervised learning



Domain-specific encoders

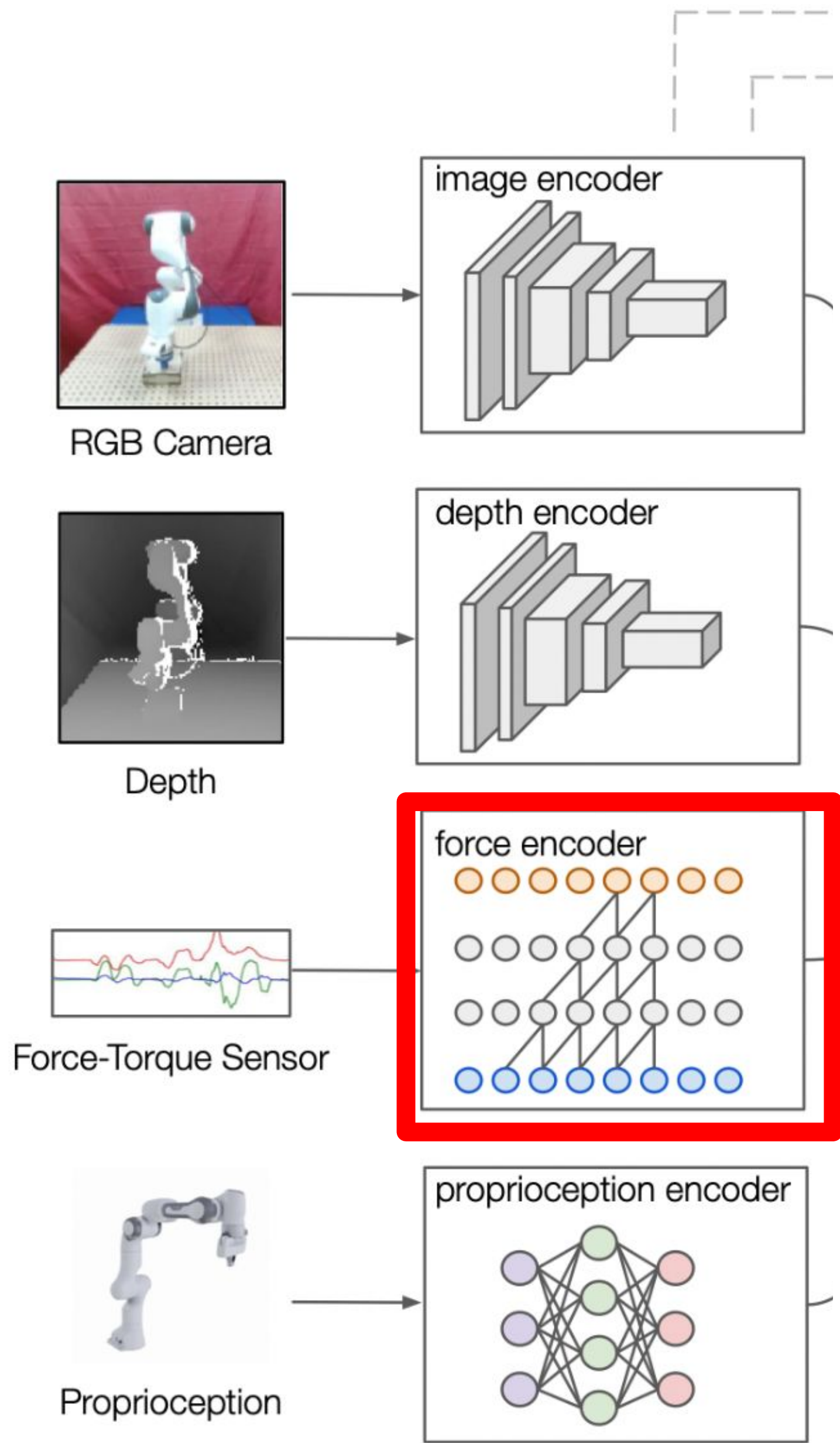
CNN (FlowNet)

CNN (VGG-16)





1. Learn latent embedding space through self-supervised learning



Domain-specific encoders

CNN (FlowNet)

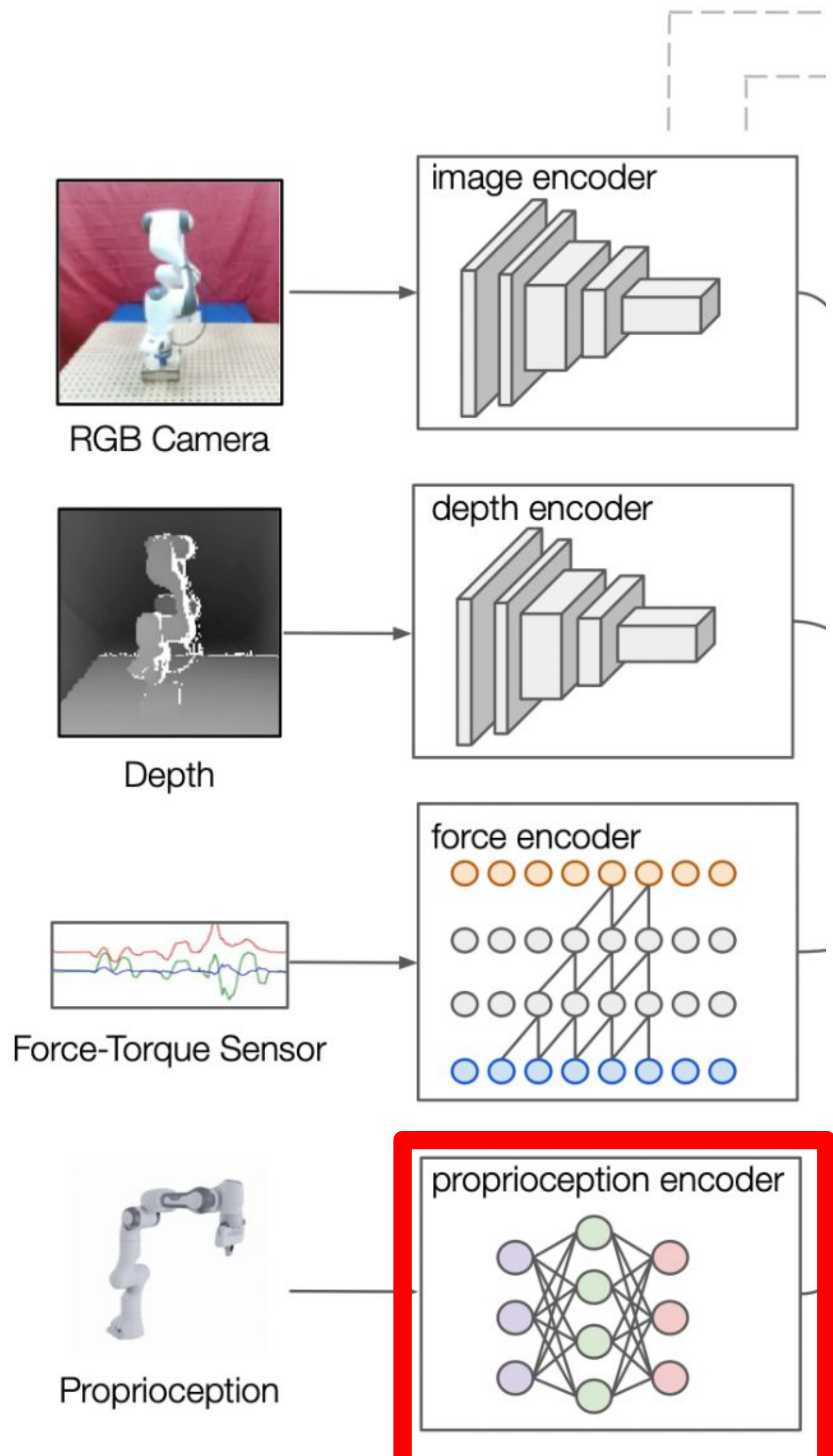
CNN (VGG-16)

Causal Convolution Layer





1. Learn latent embedding space through self-supervised learning



Domain-specific encoders

CNN (FlowNet)

CNN (VGG-16)

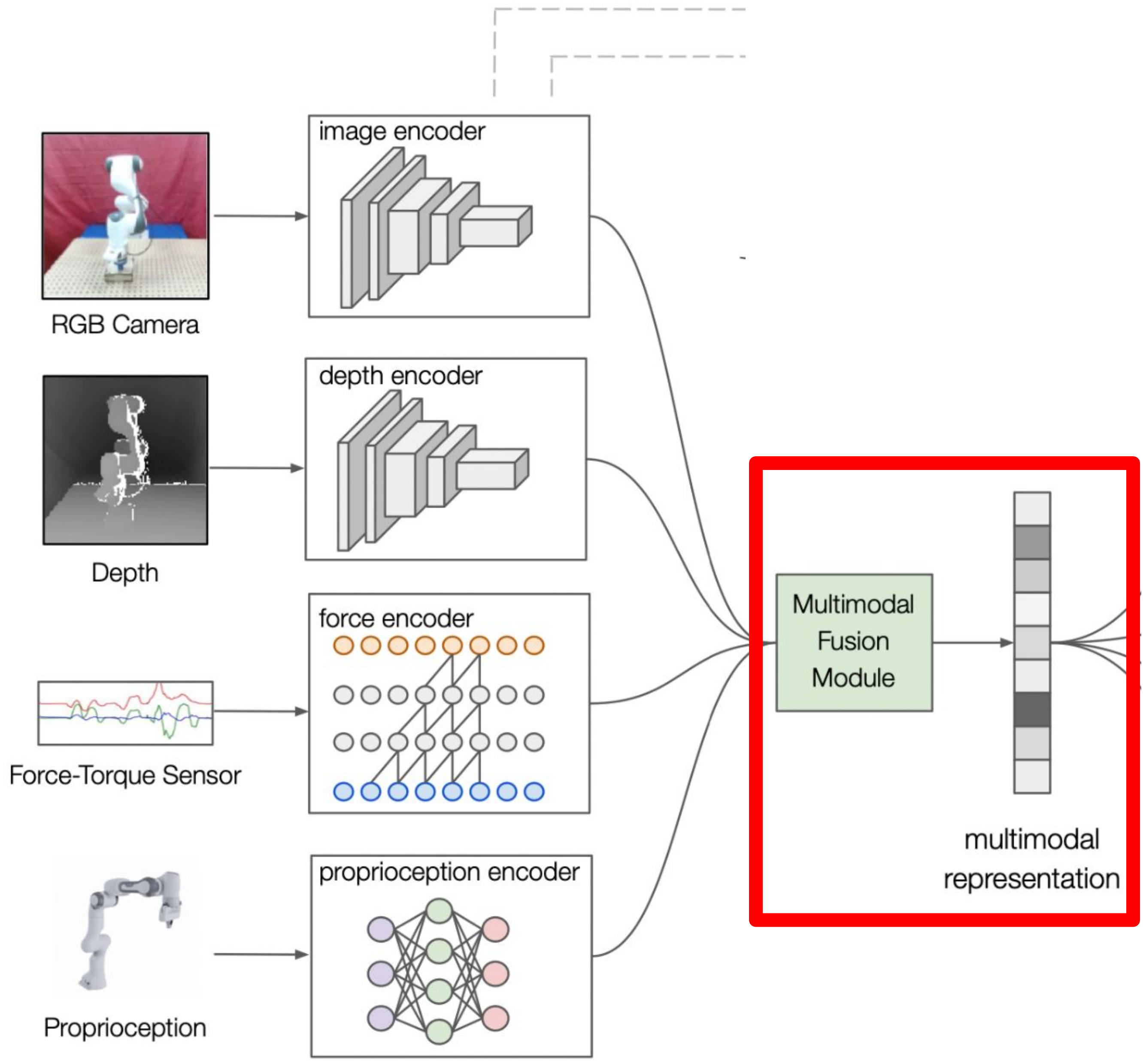
Causal Convolution Layer

Fully-Connected Network (MLP)





1. Learn latent embedding space through self-supervised learning



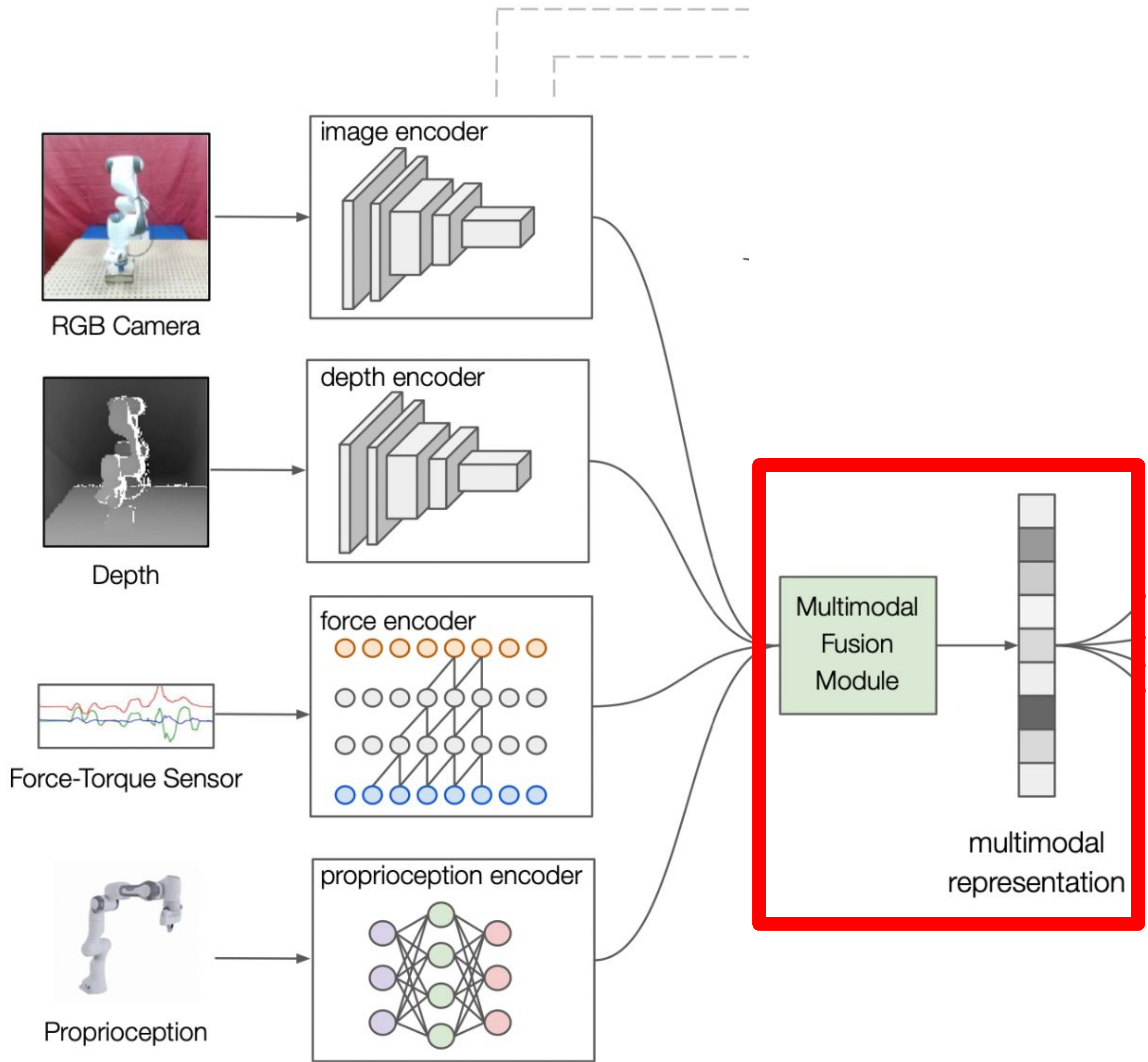
Multimodal Fusion

1. Simple concatenation





1. Learn latent embedding space through self-supervised learning



Multimodal Fusion

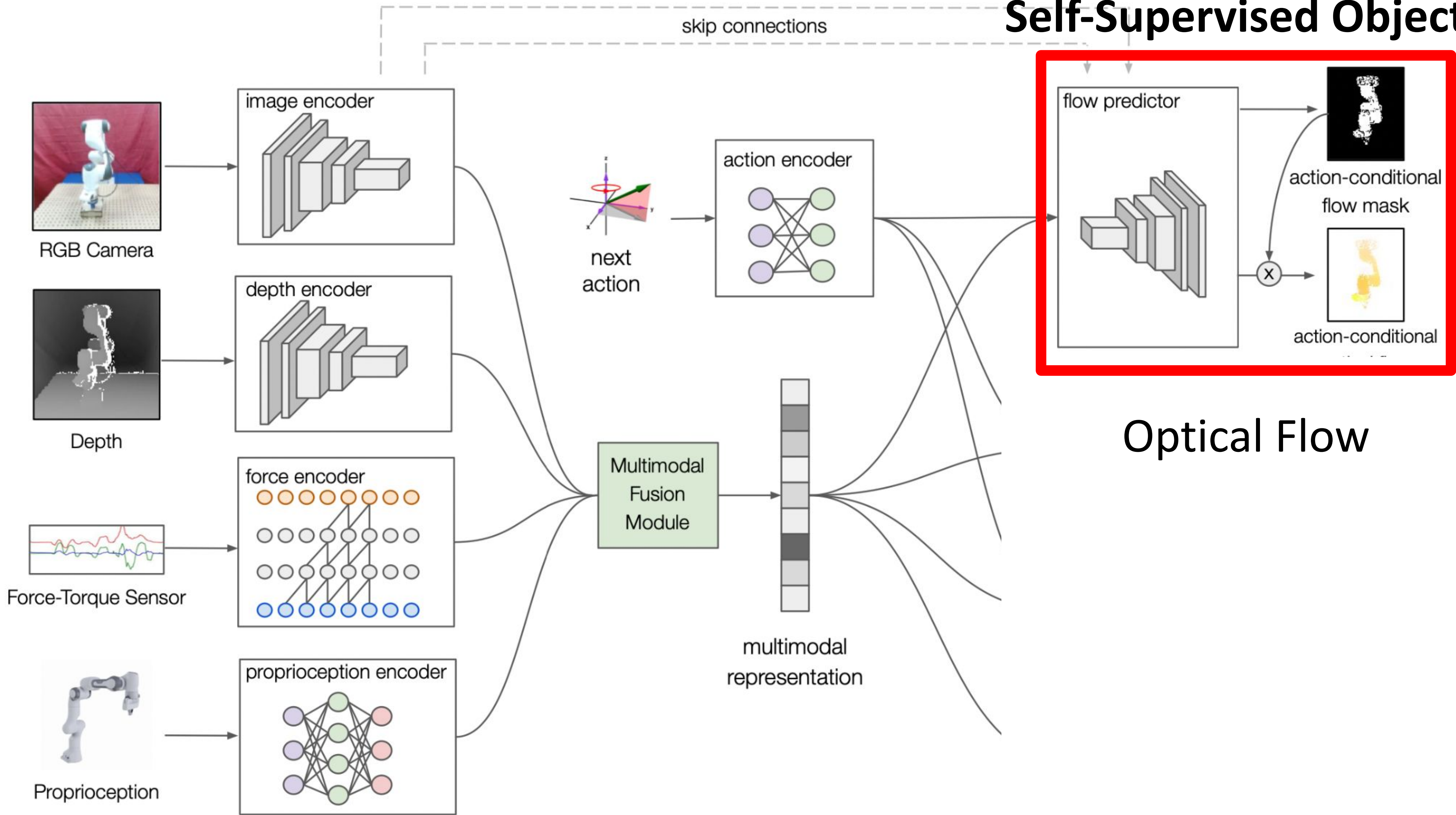
1. Simple concatenation
 2. Variational inference + product of experts
- (very similar to VAEs)





1. Learn latent embedding space through self-supervised learning

Self-Supervised Objectives



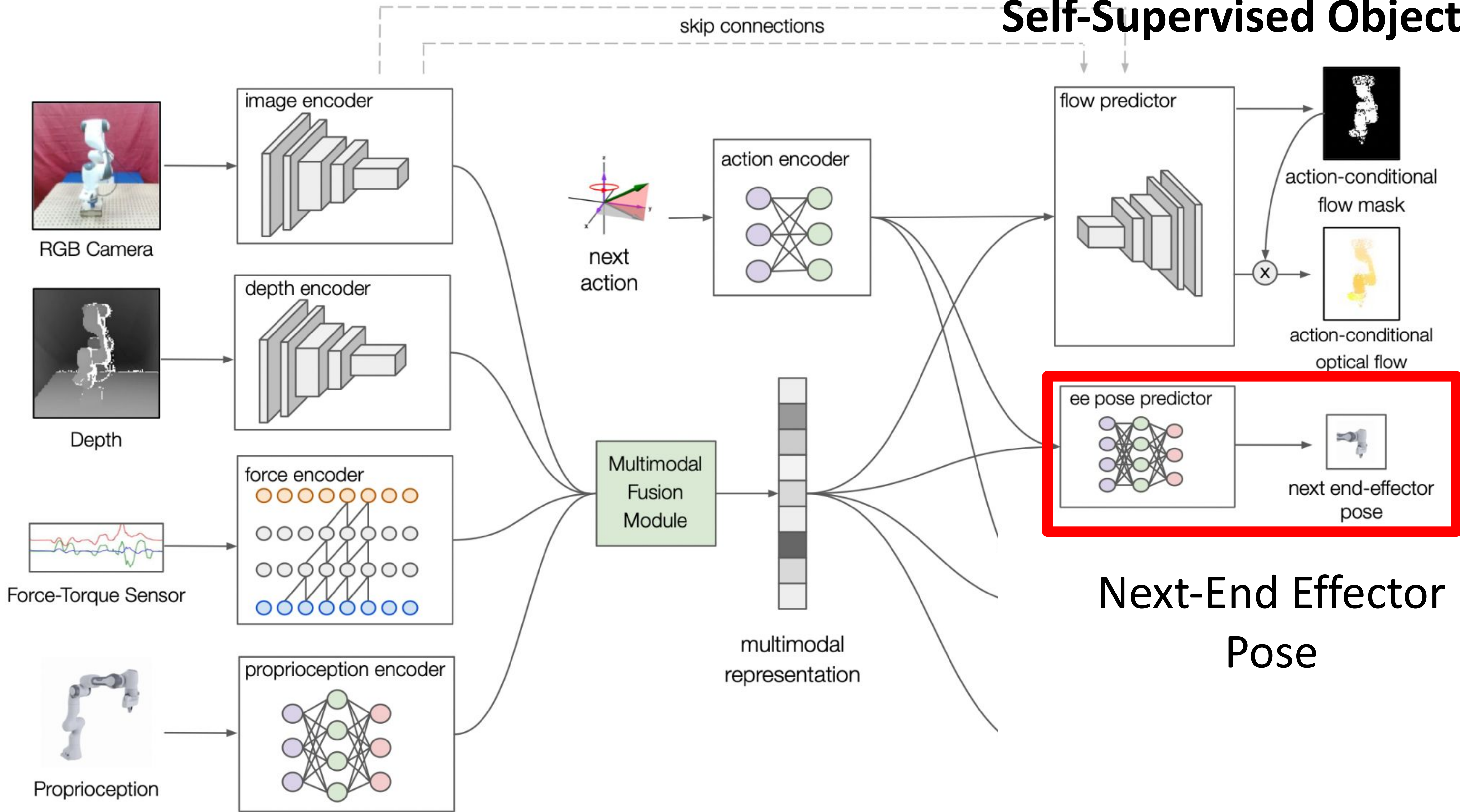
Optical Flow





1. Learn latent embedding space through self-supervised learning

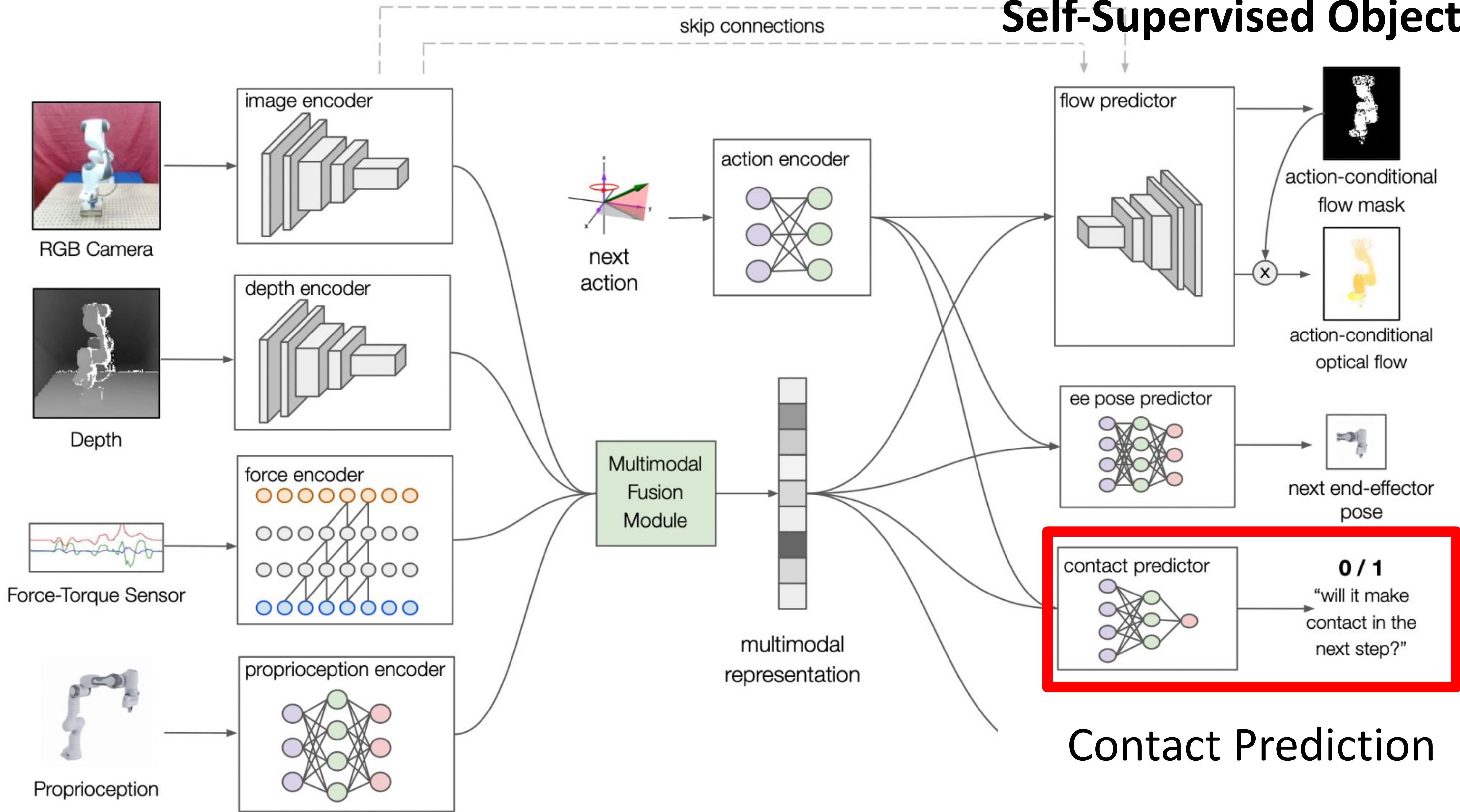
Self-Supervised Objectives





1. Learn latent embedding space through self-supervised learning

Self-Supervised Objectives



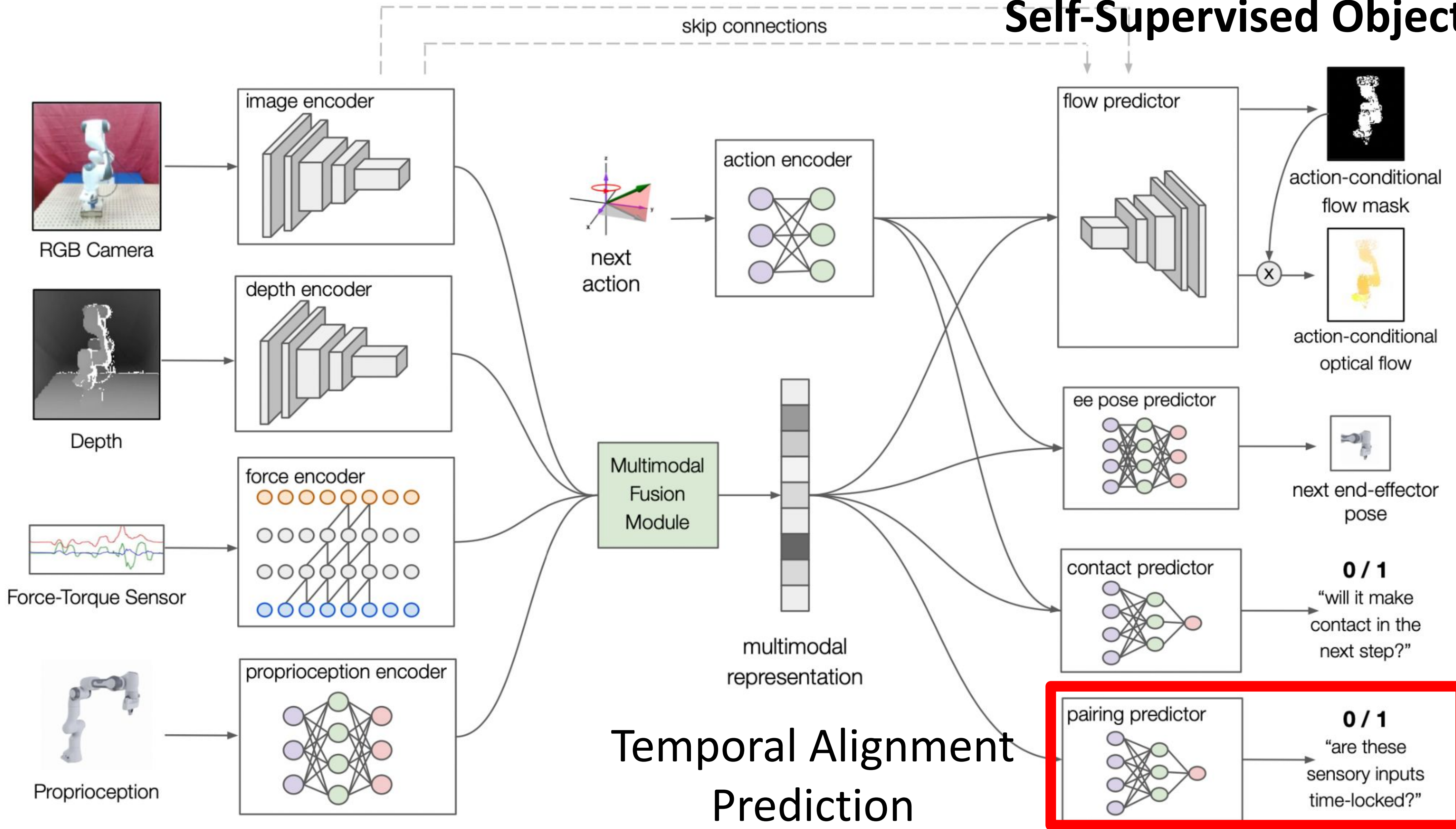
Contact Prediction





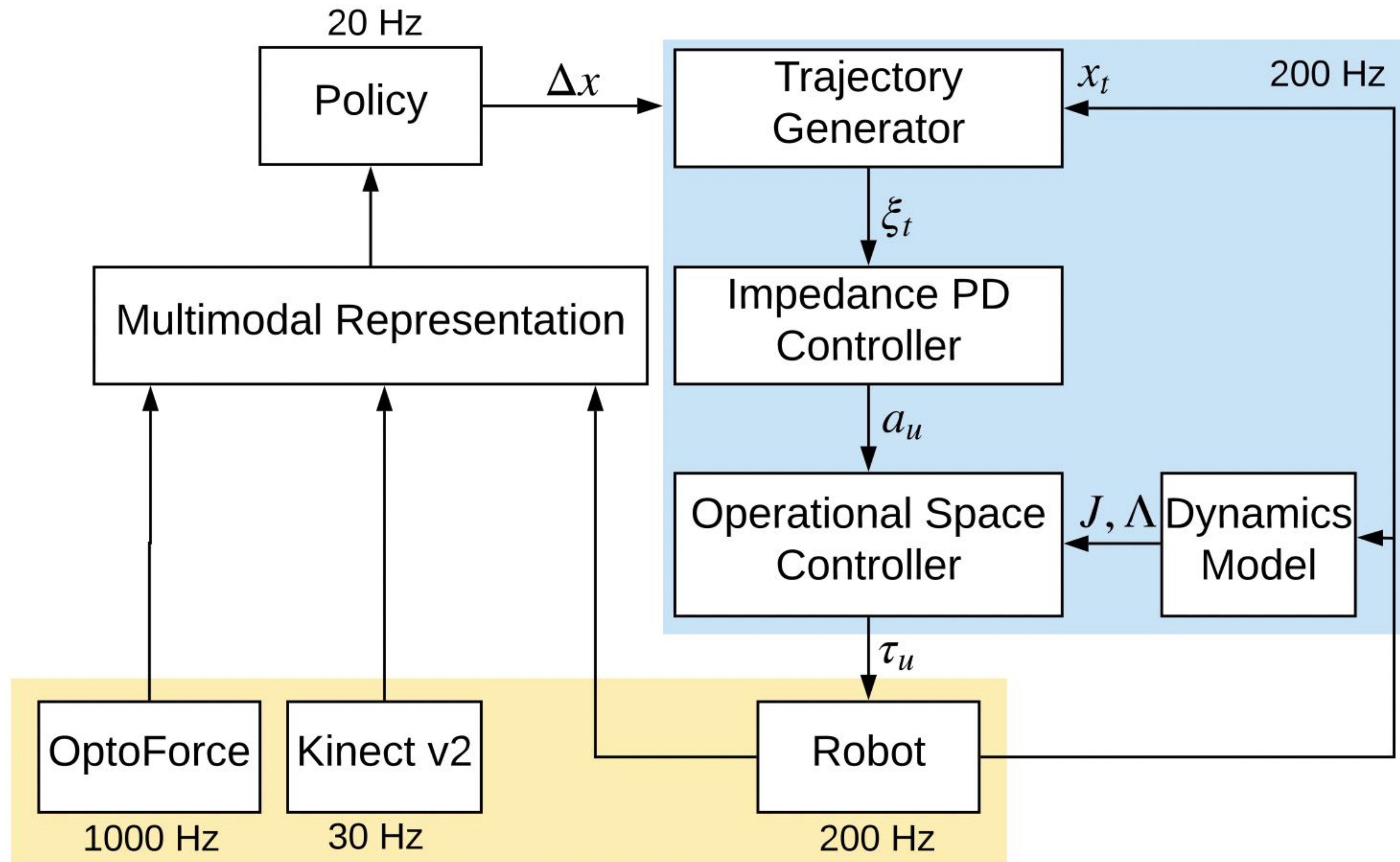
1. Learn latent embedding space through self-supervised learning

Self-Supervised Objectives



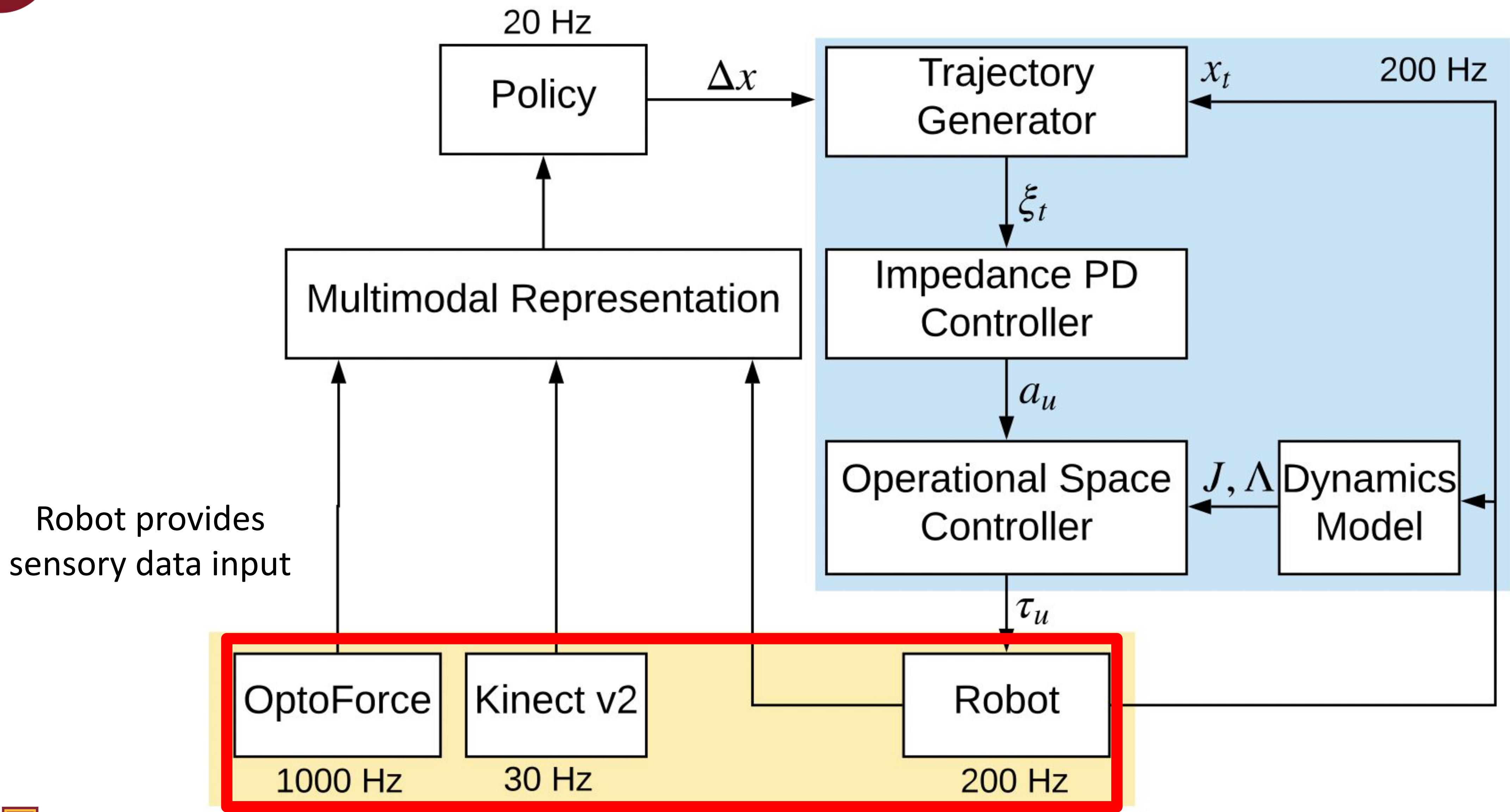


2. Use pretrained representation for policy learning



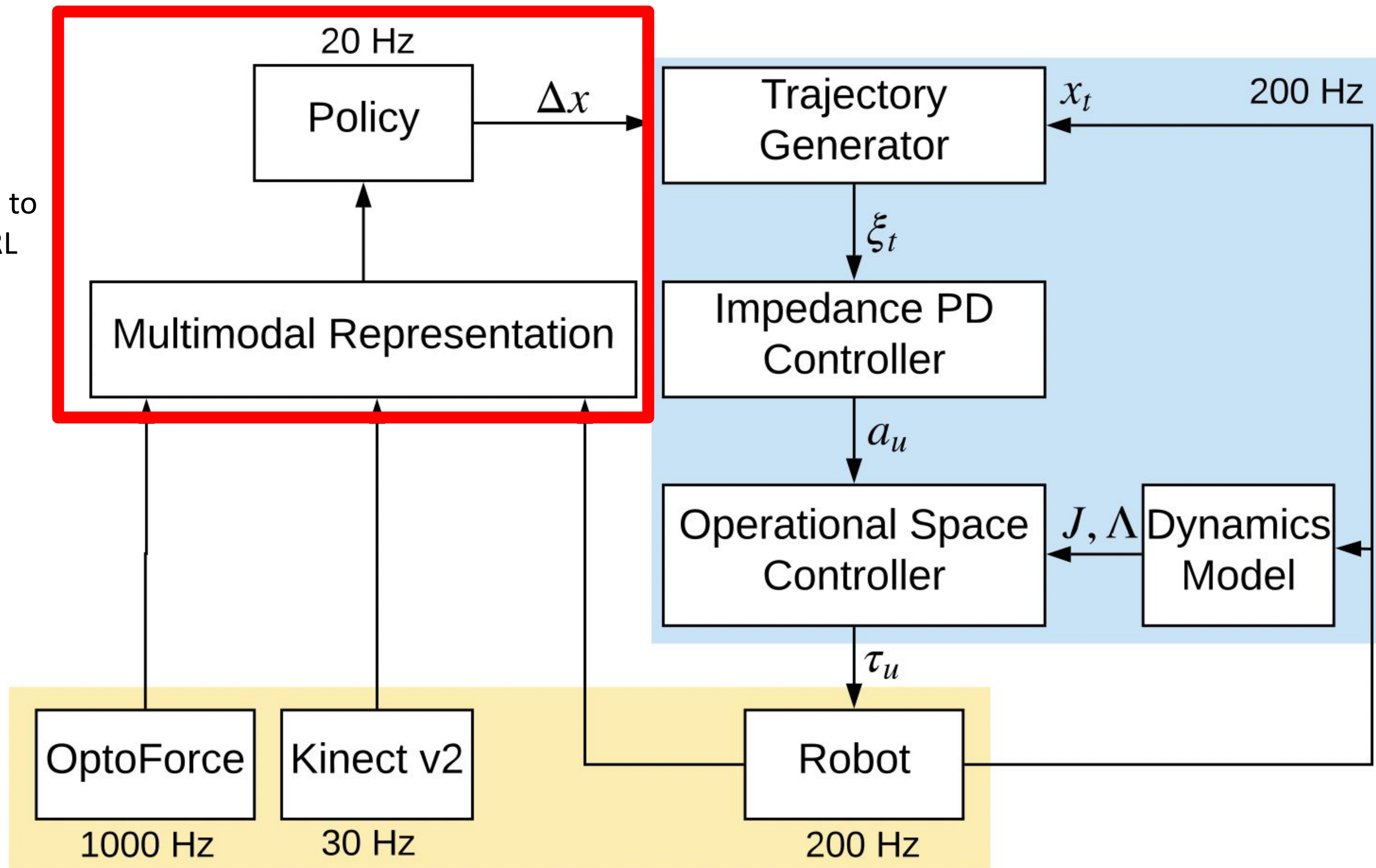


2. Use pretrained representation for policy learning



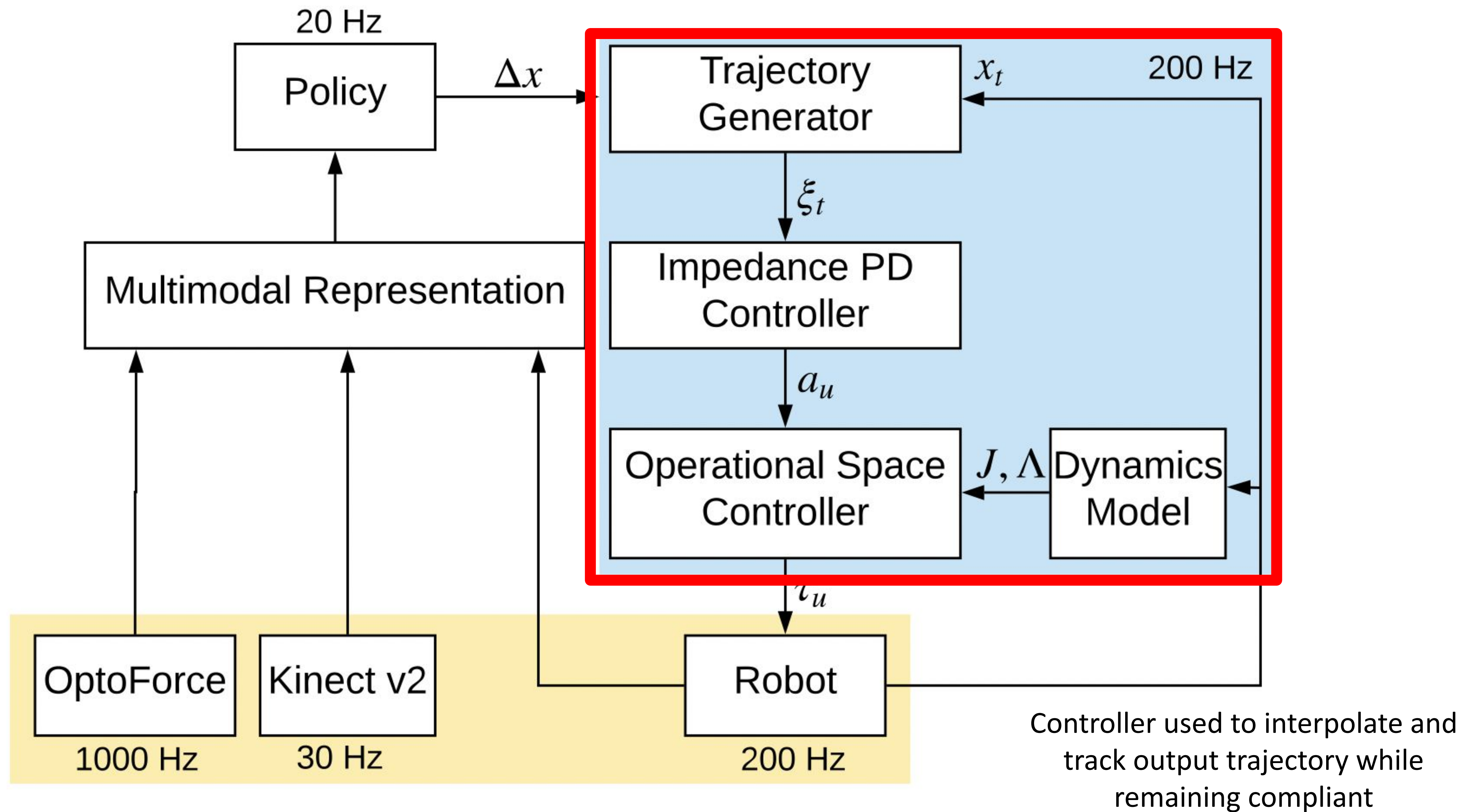
2. Use pretrained representation for policy learning

Pretrained representation used to train policy using RL





2. Use pretrained representation for policy learning





Visuo-Tactile Transformers for Manipulation

Yizhou Chen, Andrea Sipos, Mark Van der Merwe, Nima Fazeli

CoRL 2022



Focusing on contact-rich tasks

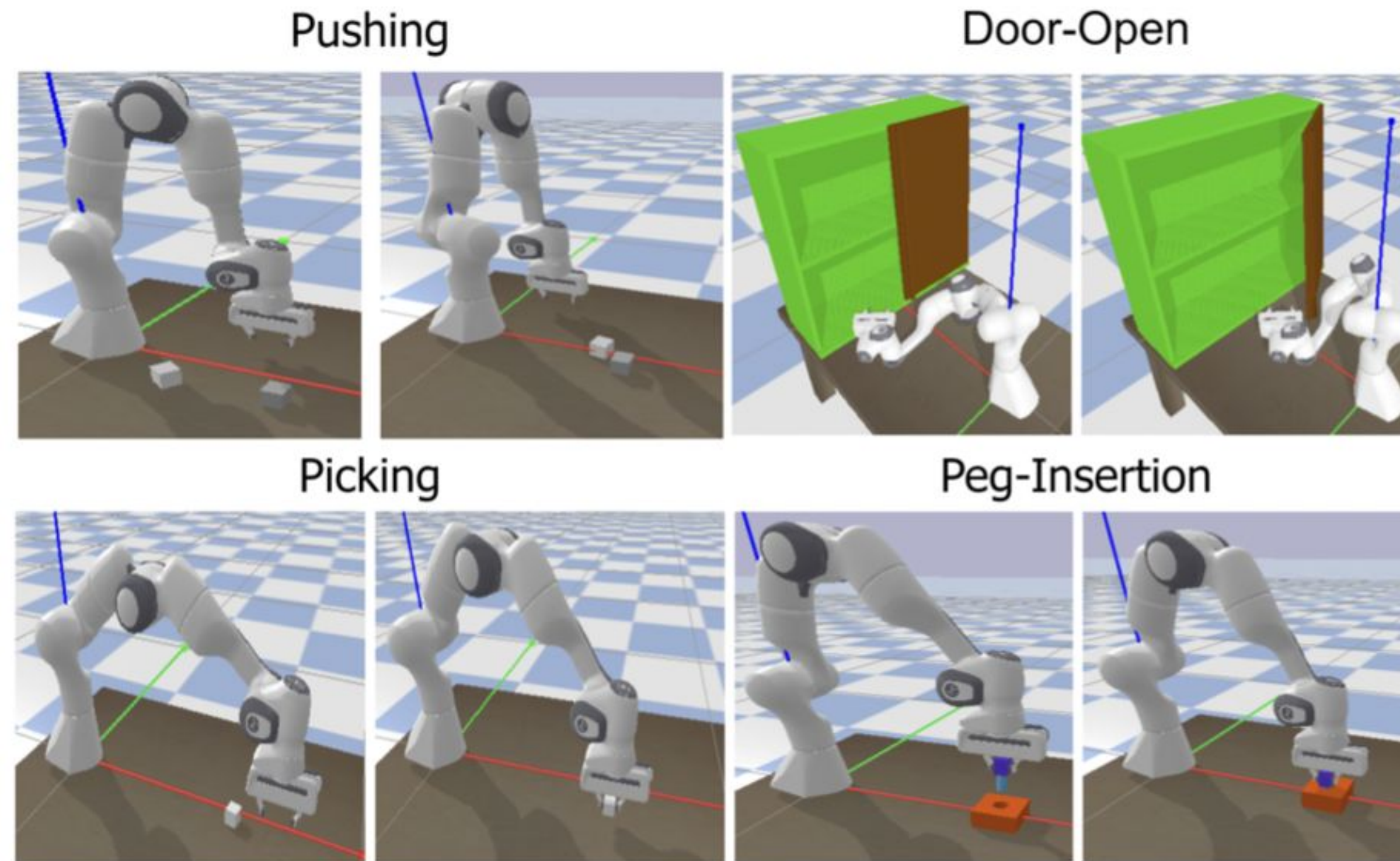
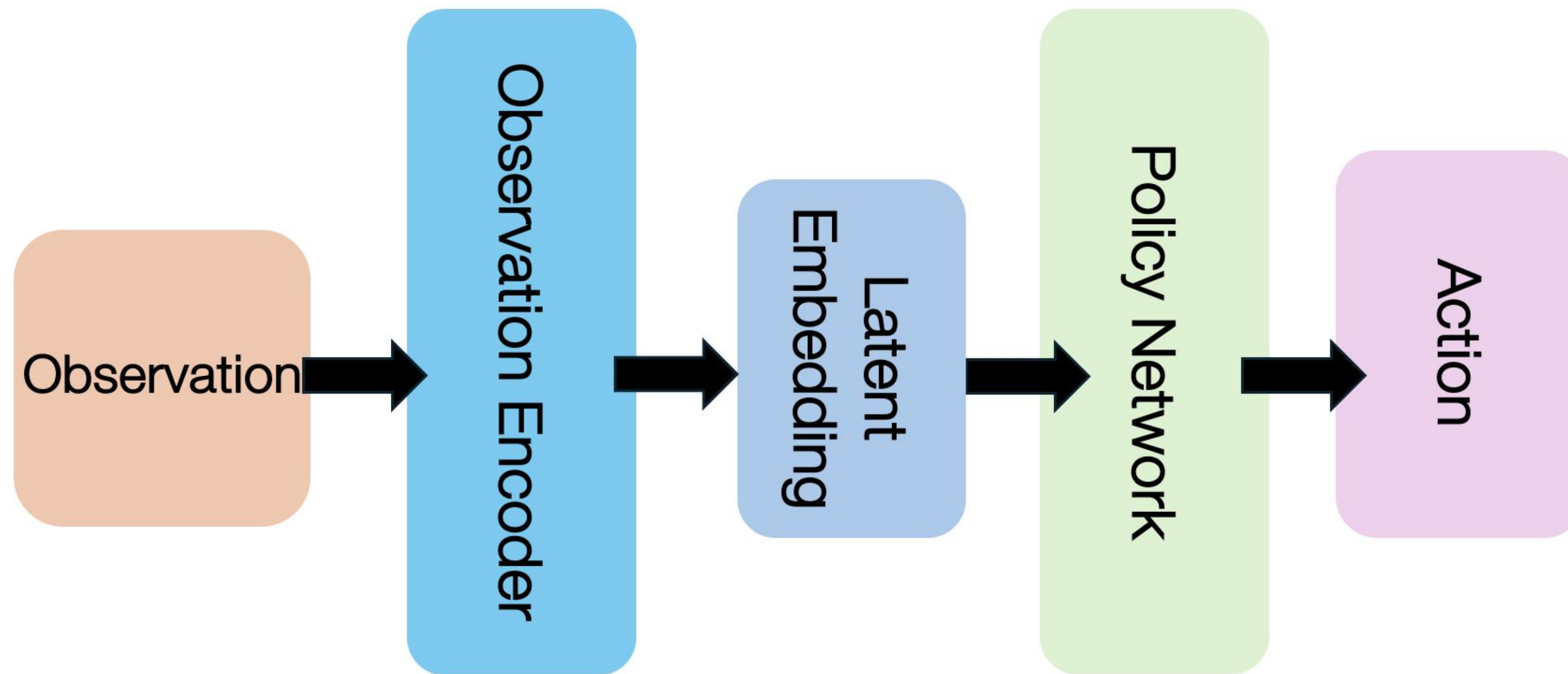


Figure 3: Manipulation Tasks: We evaluate VTT on four tasks in Pybullet. We vary visual and physical parameters in each task.

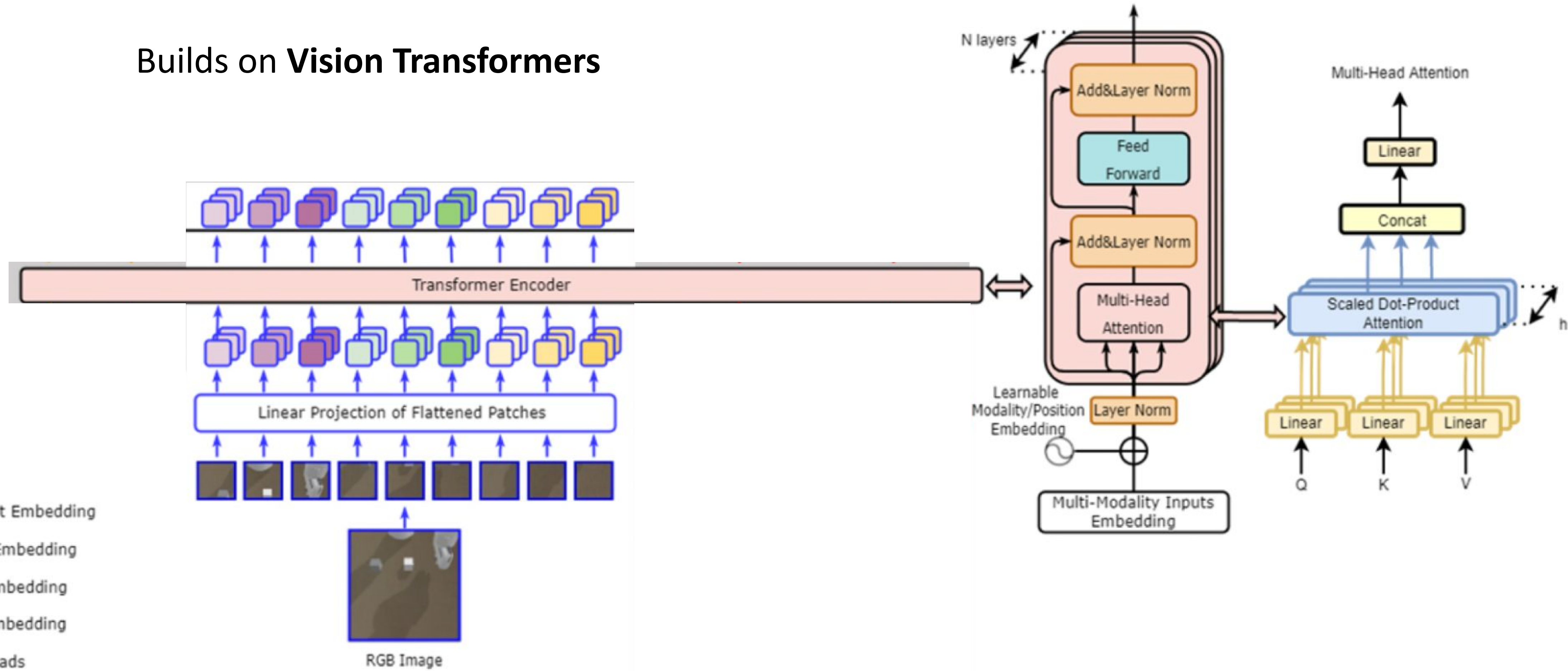
Same decoupling with a new encoder!





A shared (not separate) visuotactile encoder

Builds on **Vision Transformers**

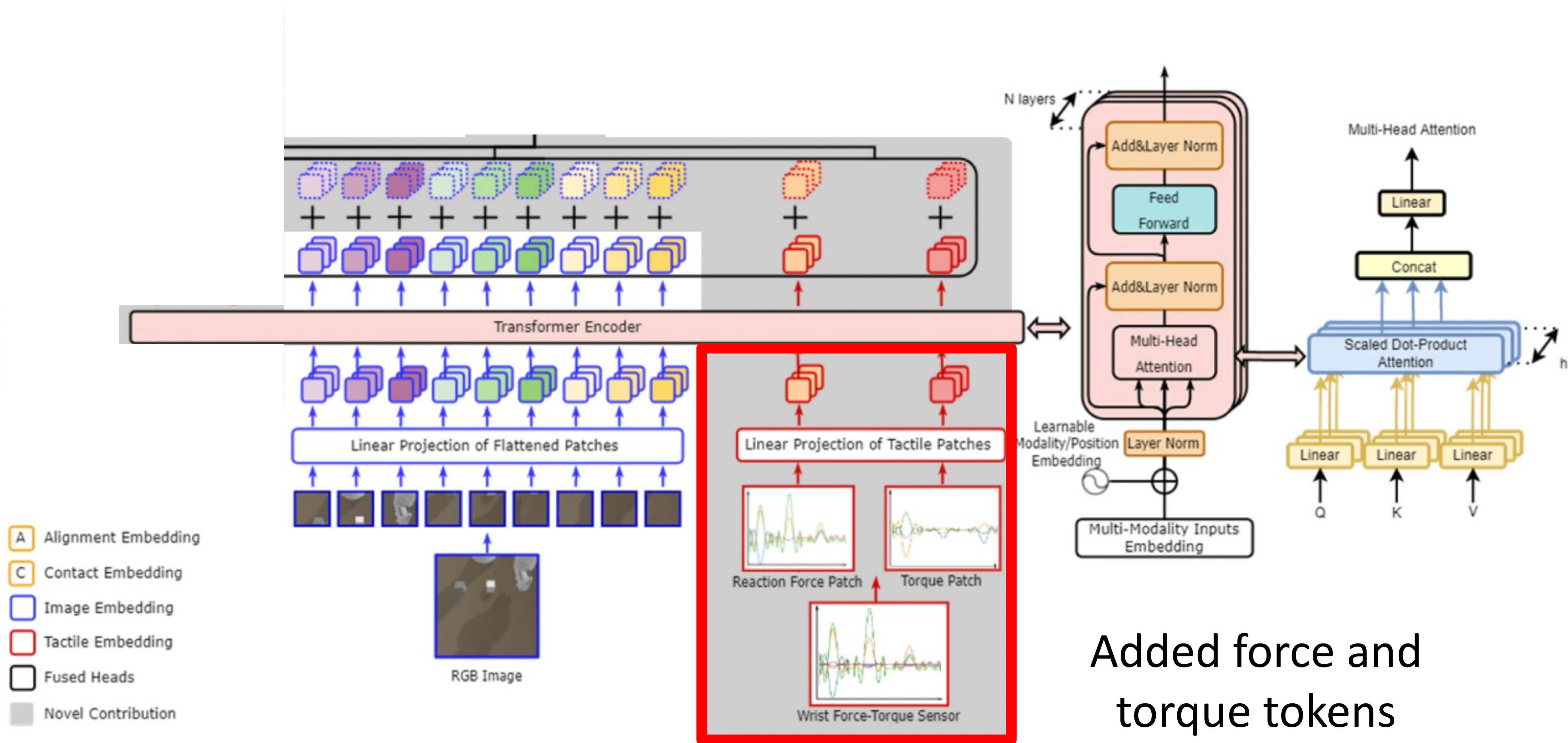


- A Alignment Embedding
- C Contact Embedding
- I Image Embedding
- T Tactile Embedding
- F Fused Heads
- Novel Contribution



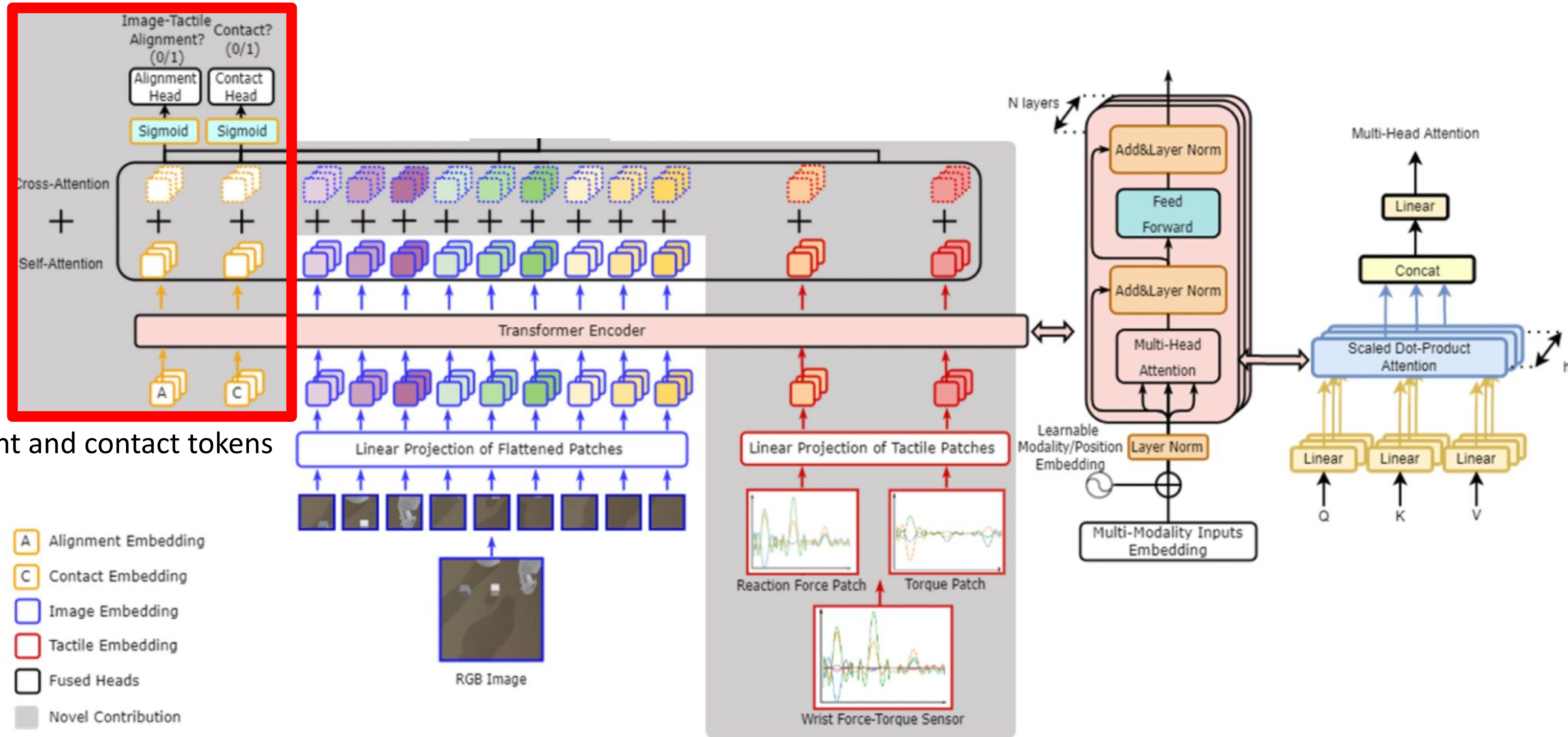


A shared (not separate) visuotactile encoder





A shared (not separate) visuotactile encoder



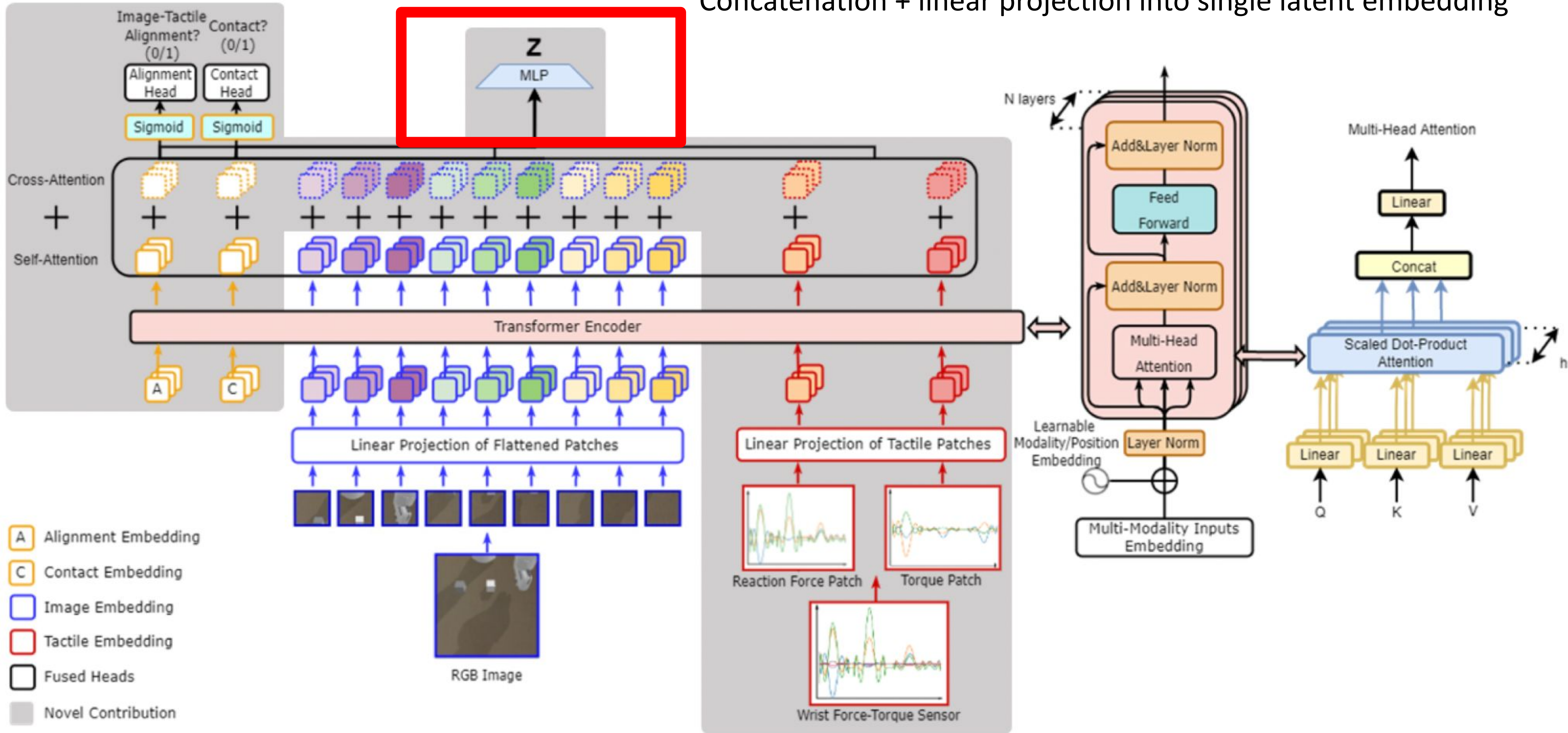
Alignment and contact tokens





A shared (not separate) visuotactile encoder

Concatenation + linear projection into single latent embedding





AugInsert: Learning Robust Visual-Force Policies via Data Augmentation for Object Assembly Tasks

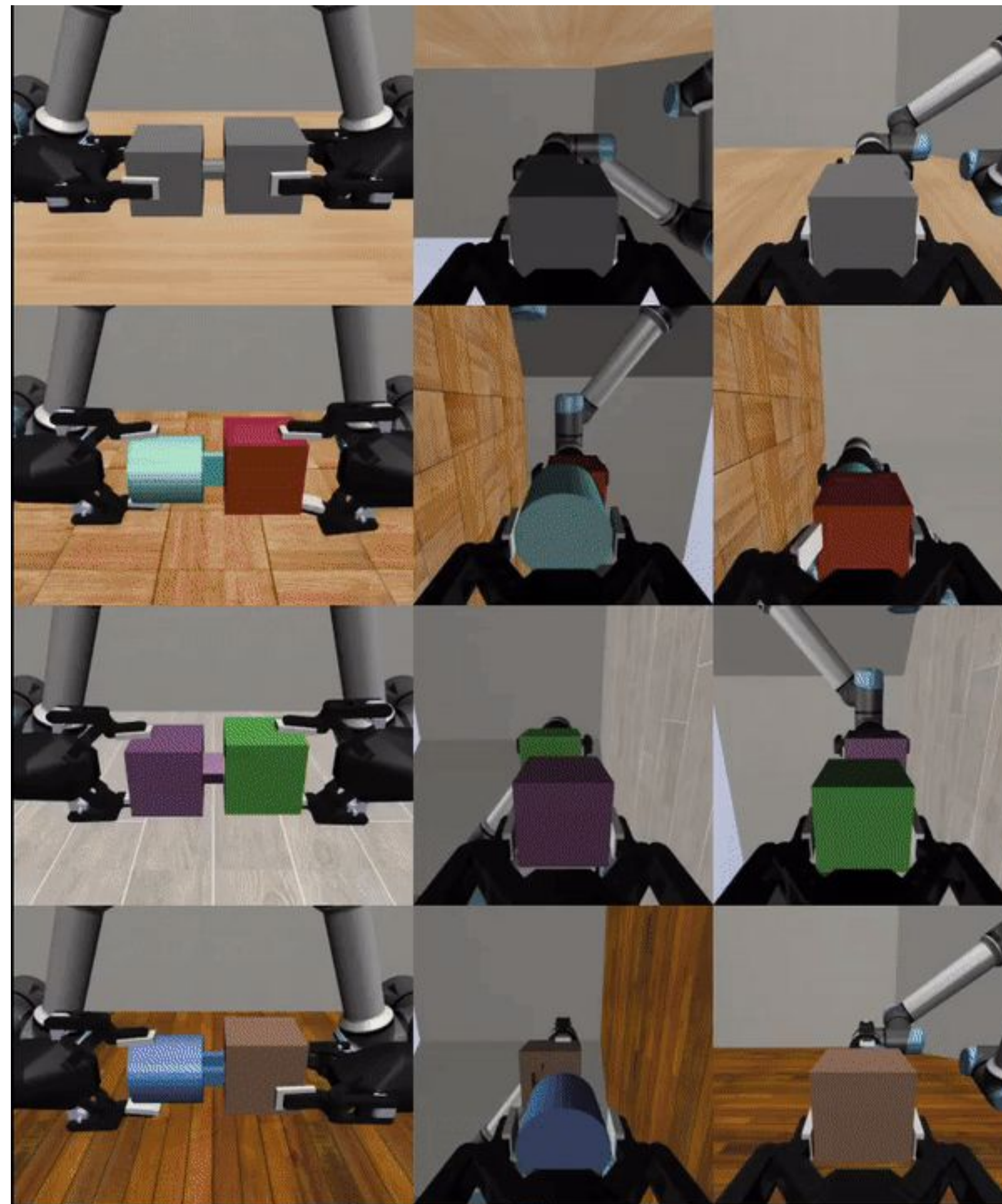
Ryan Diaz, Adam Imdieke, Vivek Veeriah, Karthik Desingh

arXiv Preprint 2024 [Under Review]



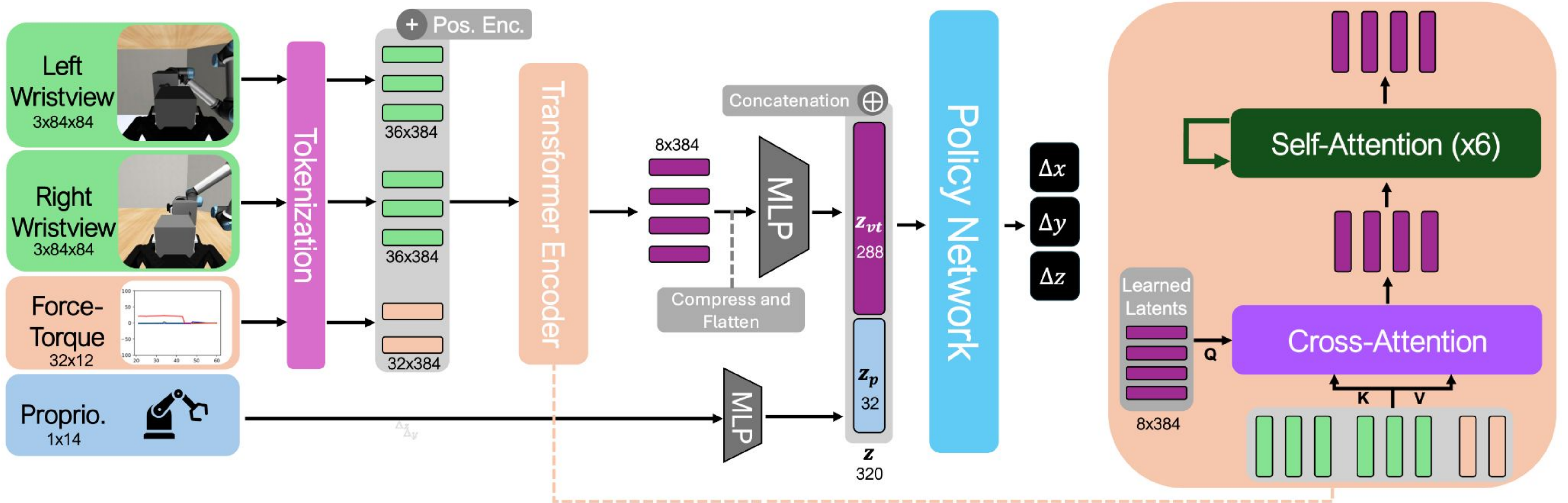


Different flavors of the same task



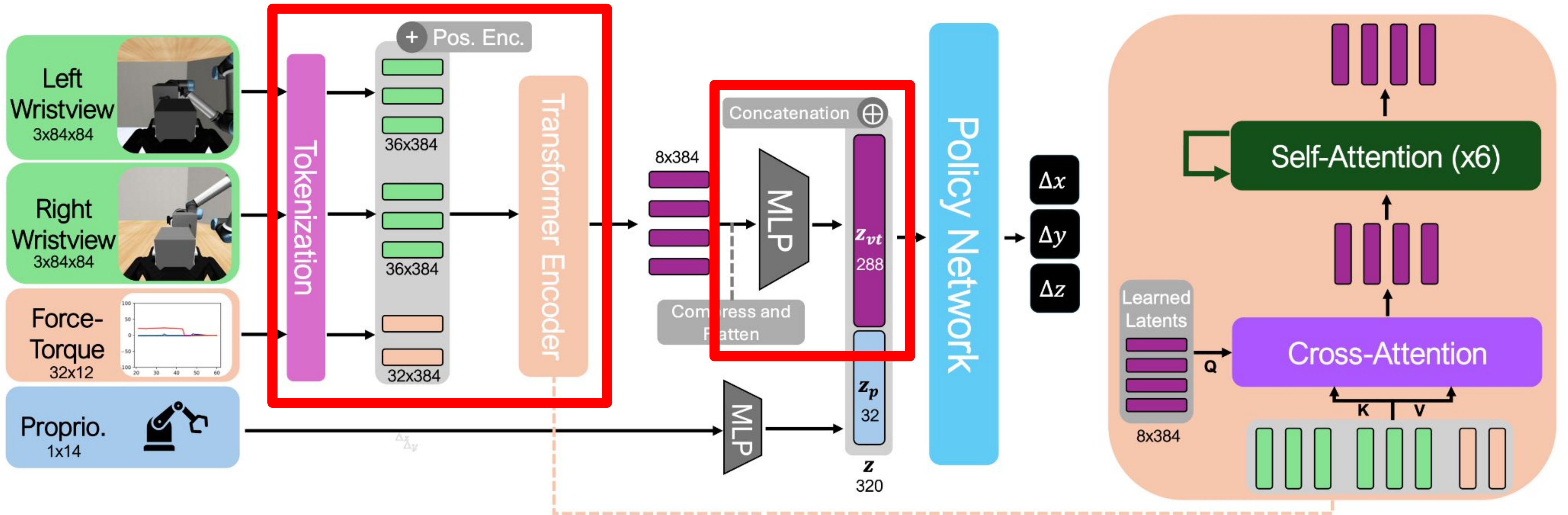


Multisensory Encoding Architecture





Multisensory Encoding Architecture

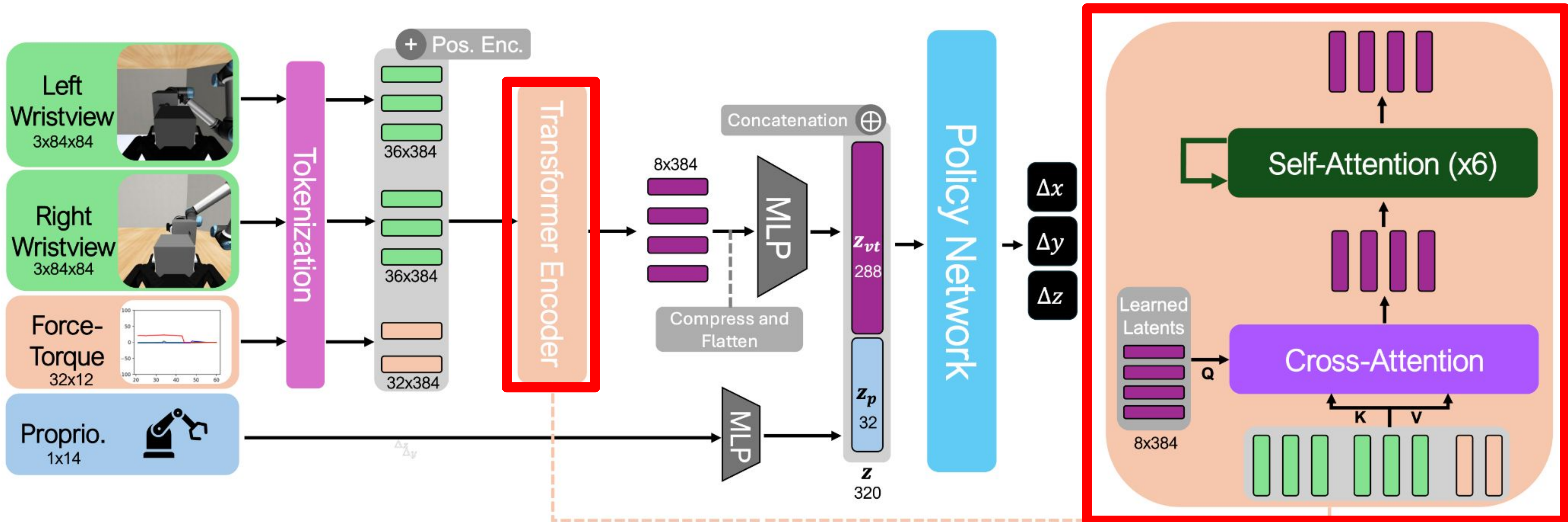


Tokenization and output representation inspired by Visuotactile Transformers





Multisensory Encoding Architecture



Latent cross-attention inspired by Perceiver/PerceiverIO





Hearing Touch: Audio-Visual Pretraining for Contact-Rich Manipulation

Jared Mejia, Victoria Dean, Tess Hellebrekers, Abhinav Gupta

ICRA 2024



Motivation

Two key ingredients for improved manipulation: **pretraining on large datasets** and using **multisensory input** (with tactile data)



Motivation

Two key ingredients for improved manipulation: **pretraining on large datasets** and using **multisensory input** (with tactile data)

How can we combine the two?



Motivation

Two key ingredients for improved manipulation: **pretraining on large datasets** and using **multisensory input** (with tactile data)

How can we combine the two?

Large internet-scale image datasets exist, but not much for tactile



Motivation

Two key ingredients for improved manipulation: **pretraining on large datasets** and using **multisensory input** (with tactile data)

How can we combine the two?

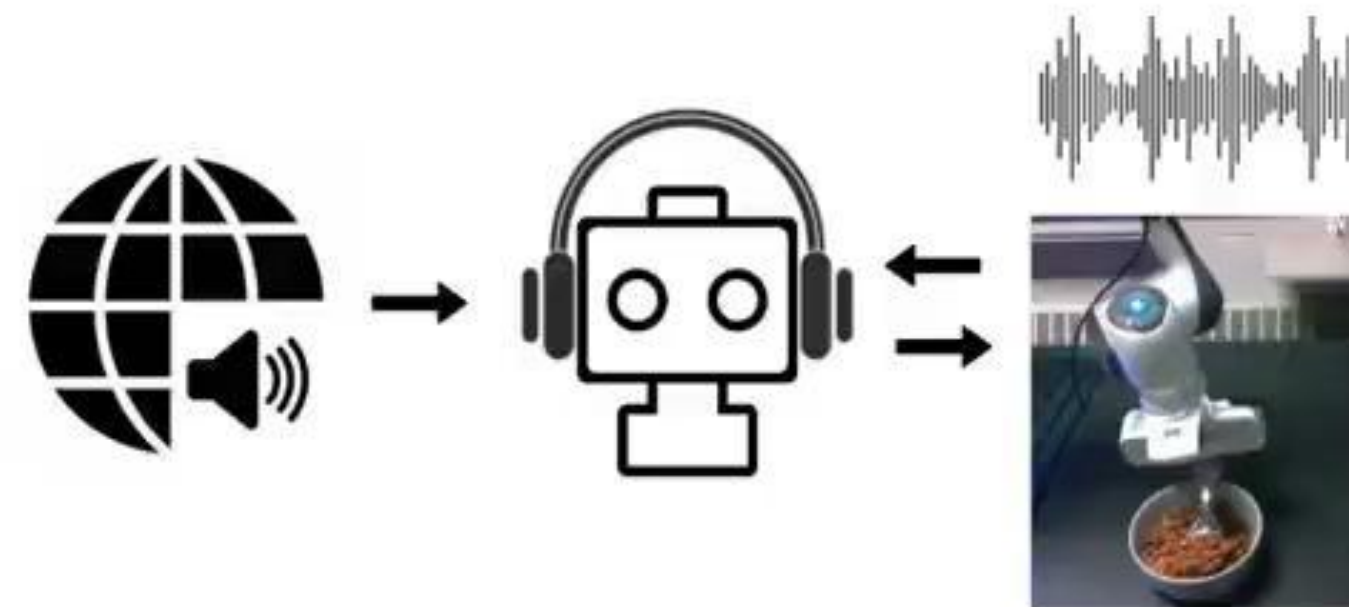
Large internet-scale image datasets exist, but not much for tactile

Idea: Leverage *contact audio* as a tactile sensing mode to enable the use of internet-scale audio datasets for pretraining

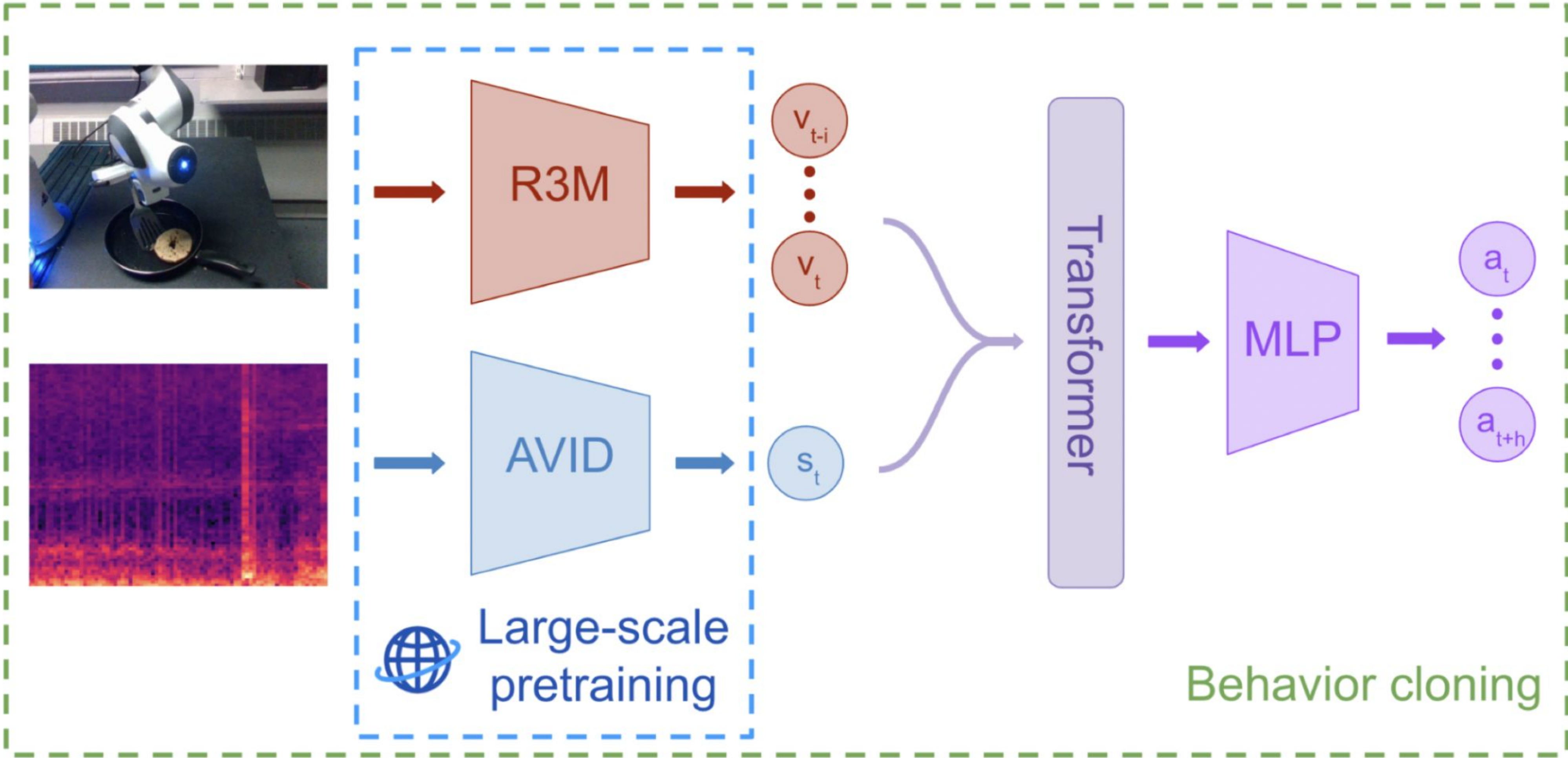


Sound in Robotic Manipulation

Hearing Touch: Audio-Visual Pretraining for Contact-Rich Manipulation



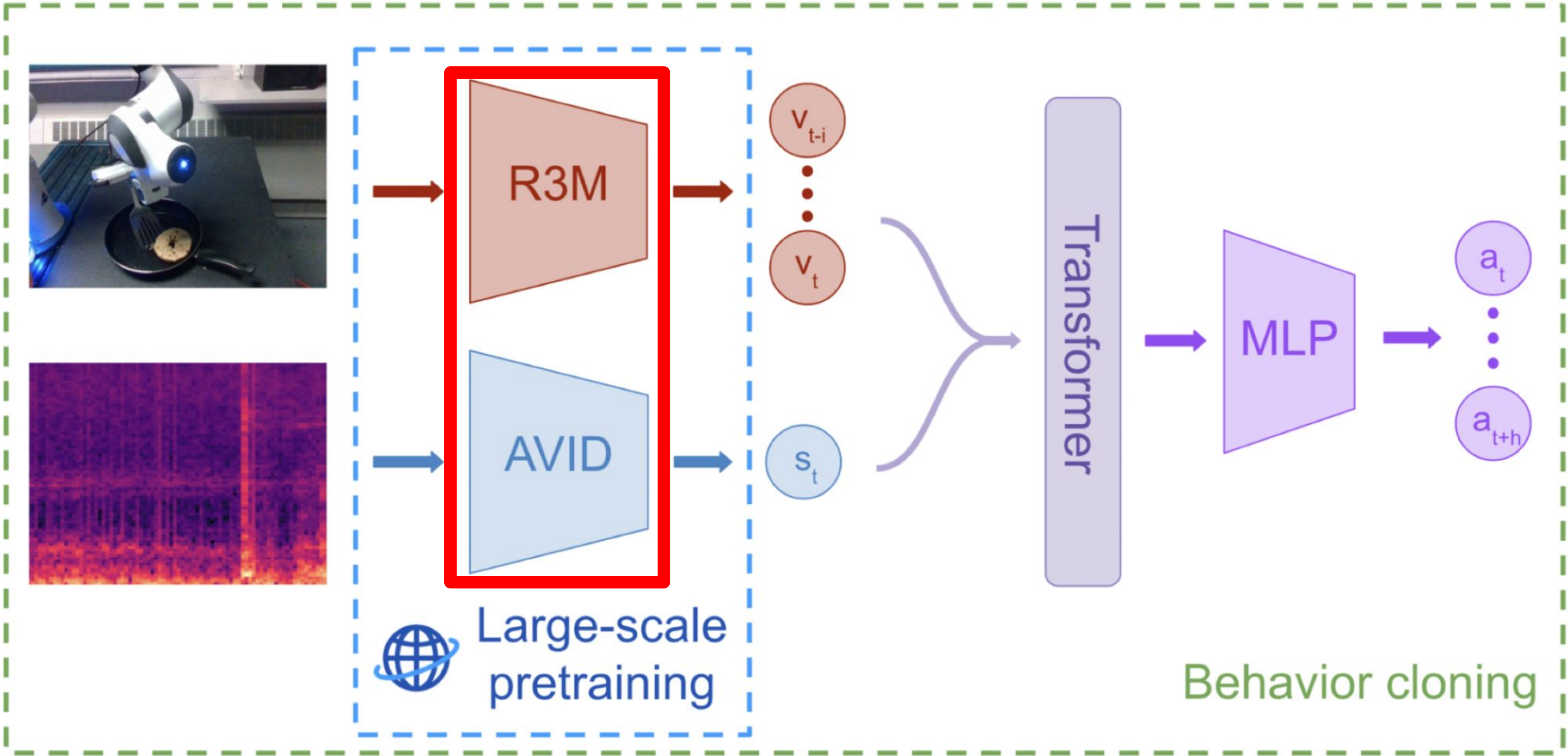
Framework





Framework

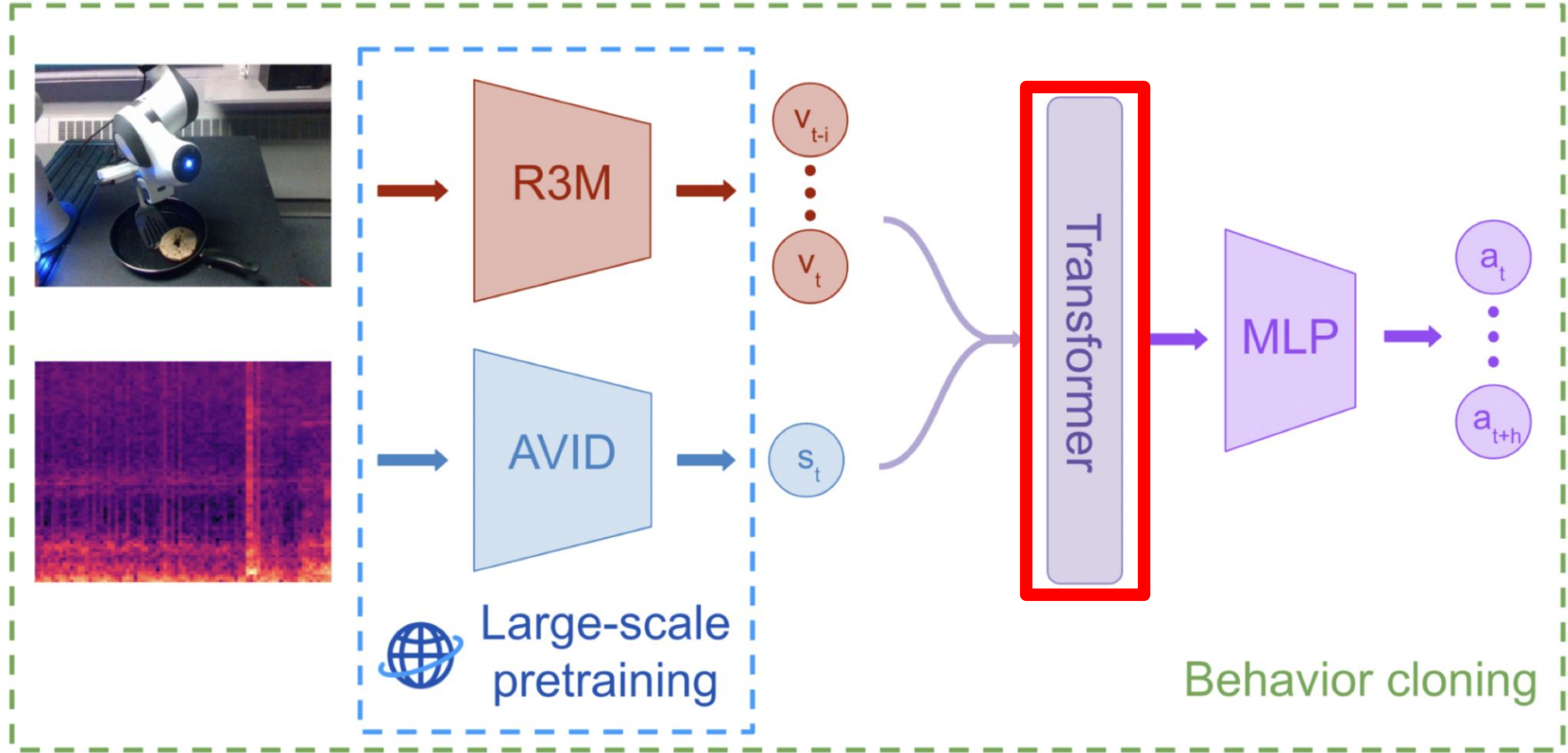
R3M: Trained on **Ego4D** (>3670 hours of egocentric human video)
AVID: Trained on **Audioset** (>2 million 10-sec audio clips from YouTube)





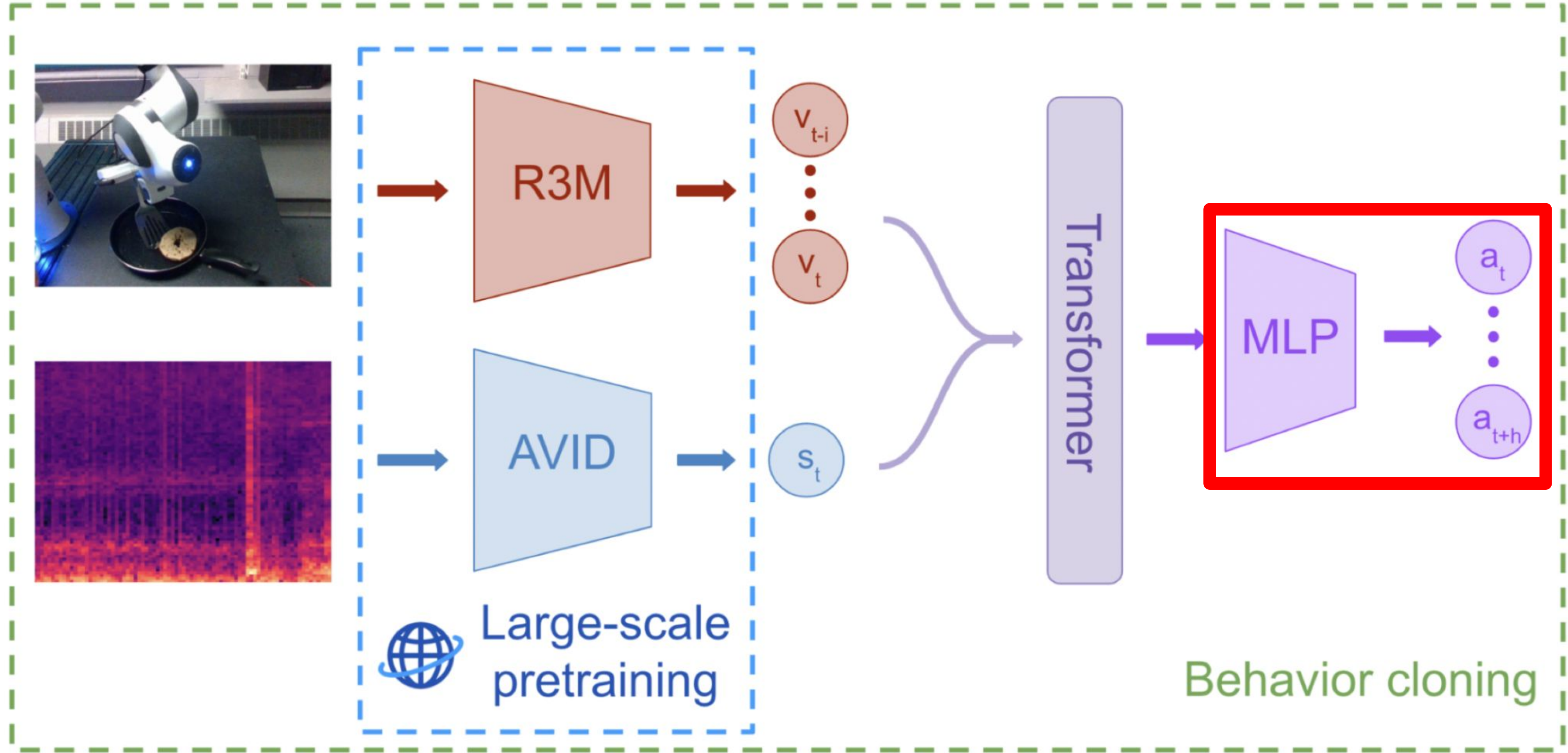
Framework

Multisensory transformer for modality fusion



Framework

MLP policy network trained via behavior cloning





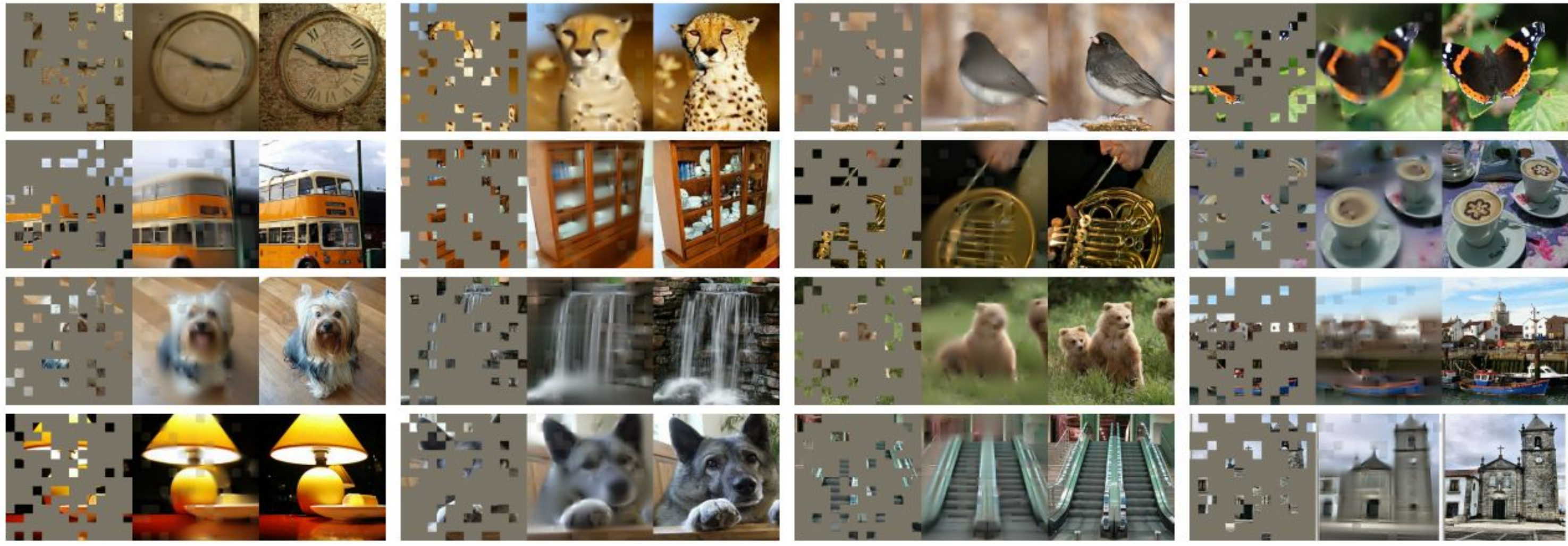
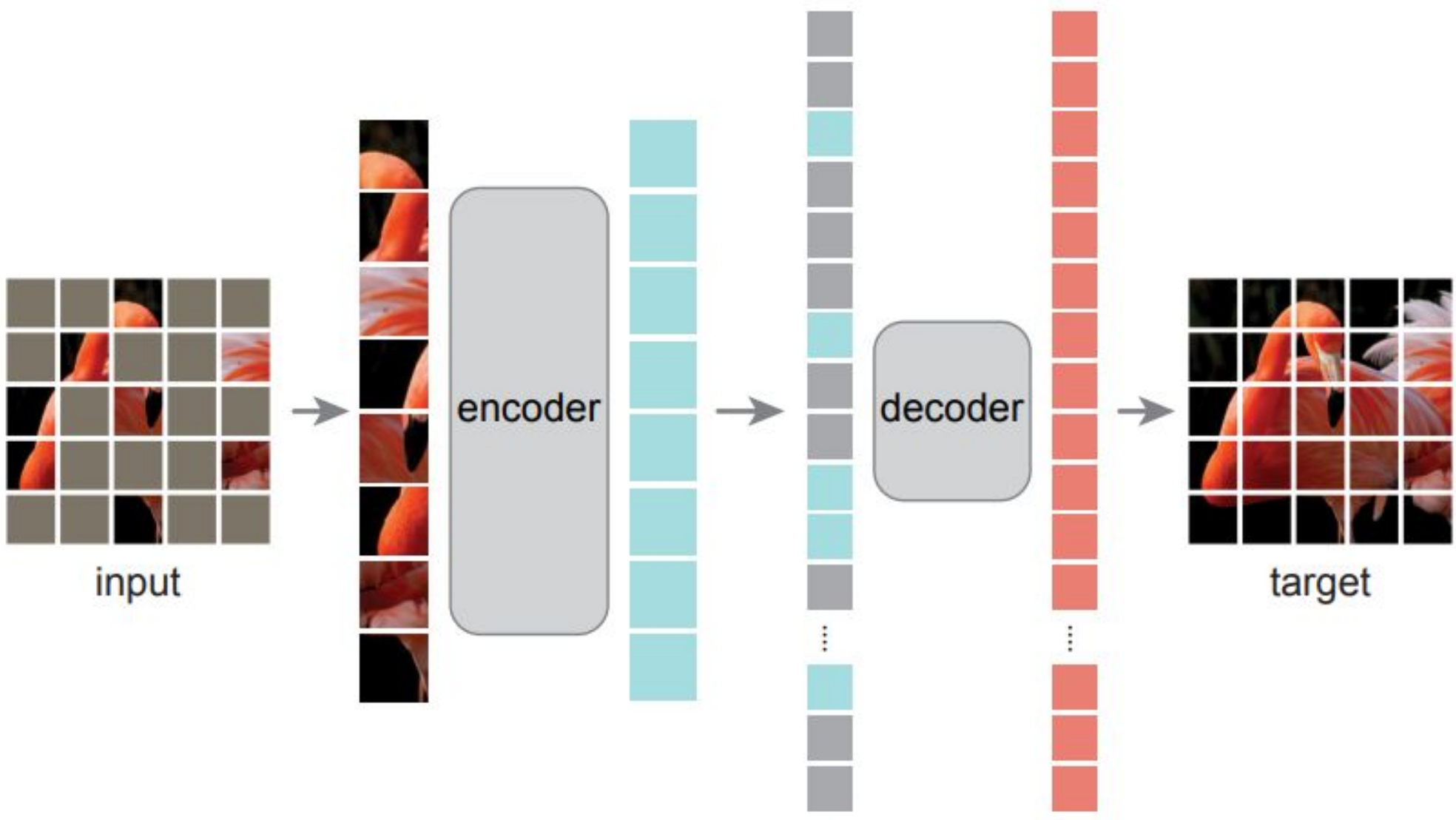
MultiMAE: Multi-modal Multi-task Masked Autoencoders

Roman Bachmann, David Mizrahi, Andrei Atanov, Amir Zamir

ECCV 2022



Masked autoencoder



motivation:

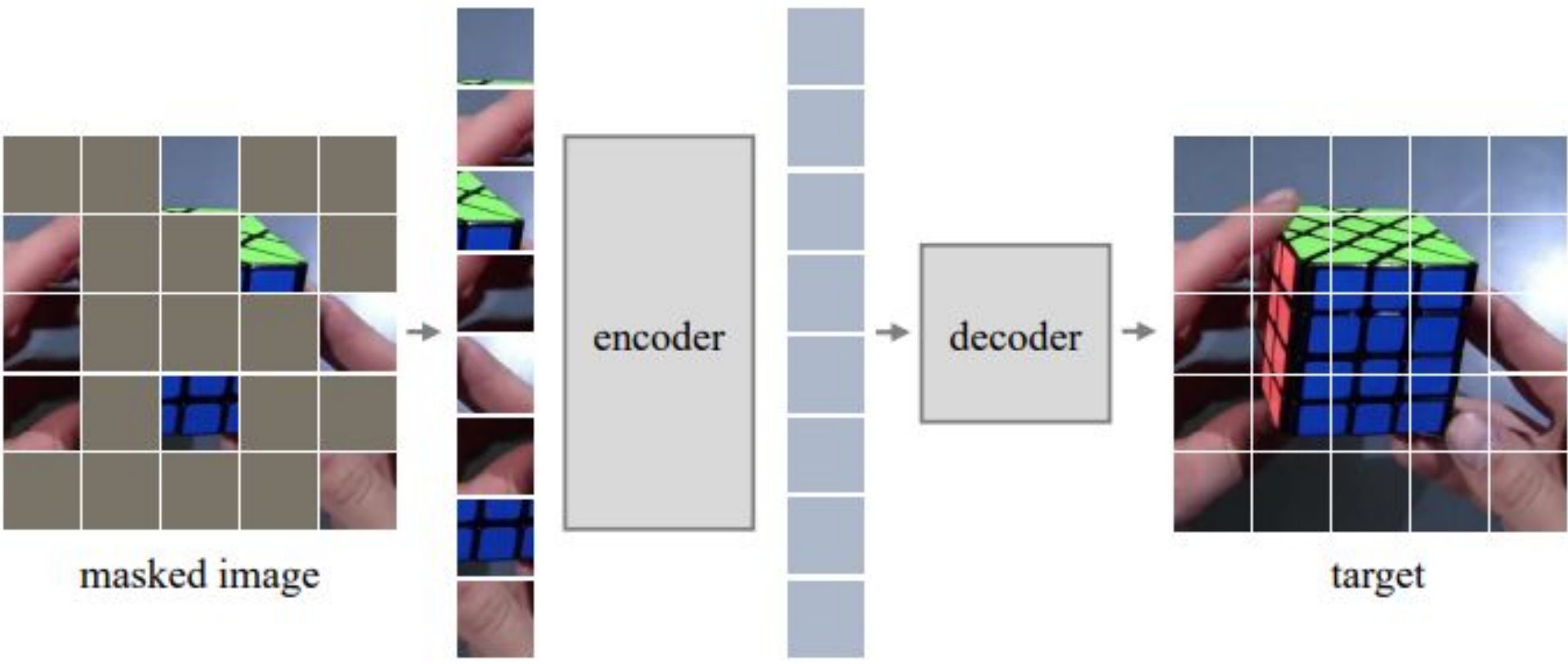
- 1. image is heavily redundant
- 2. computationally efficient

implementation:

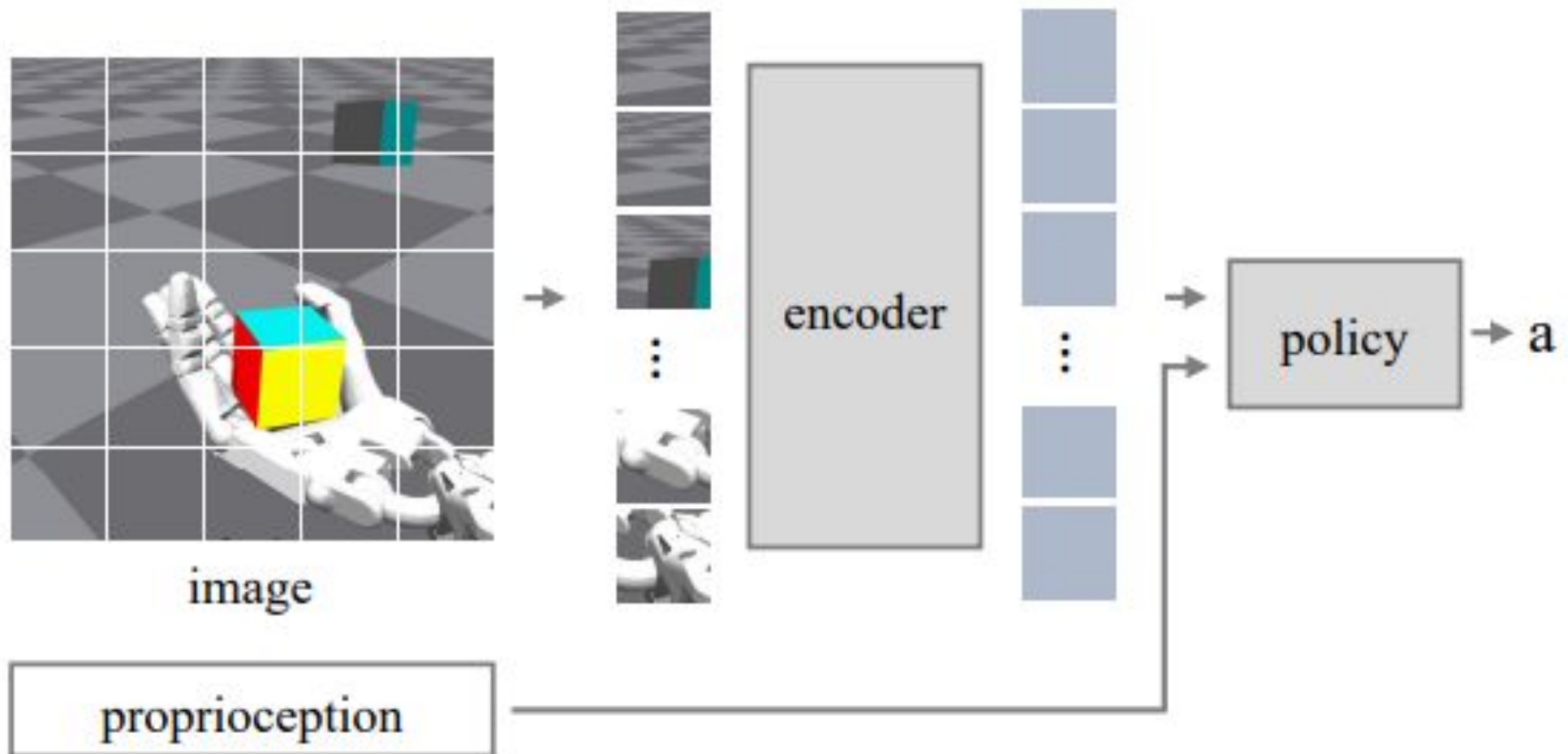
- 1. mask rate: 75%
- 2. a light weight decoder
- 3. loss only calculate on masked patches



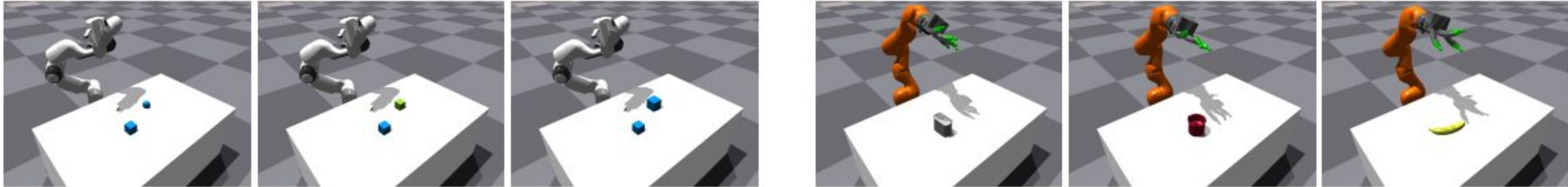
Masked Visual Pre-training for Motor Control



(a) masked visual pretraining

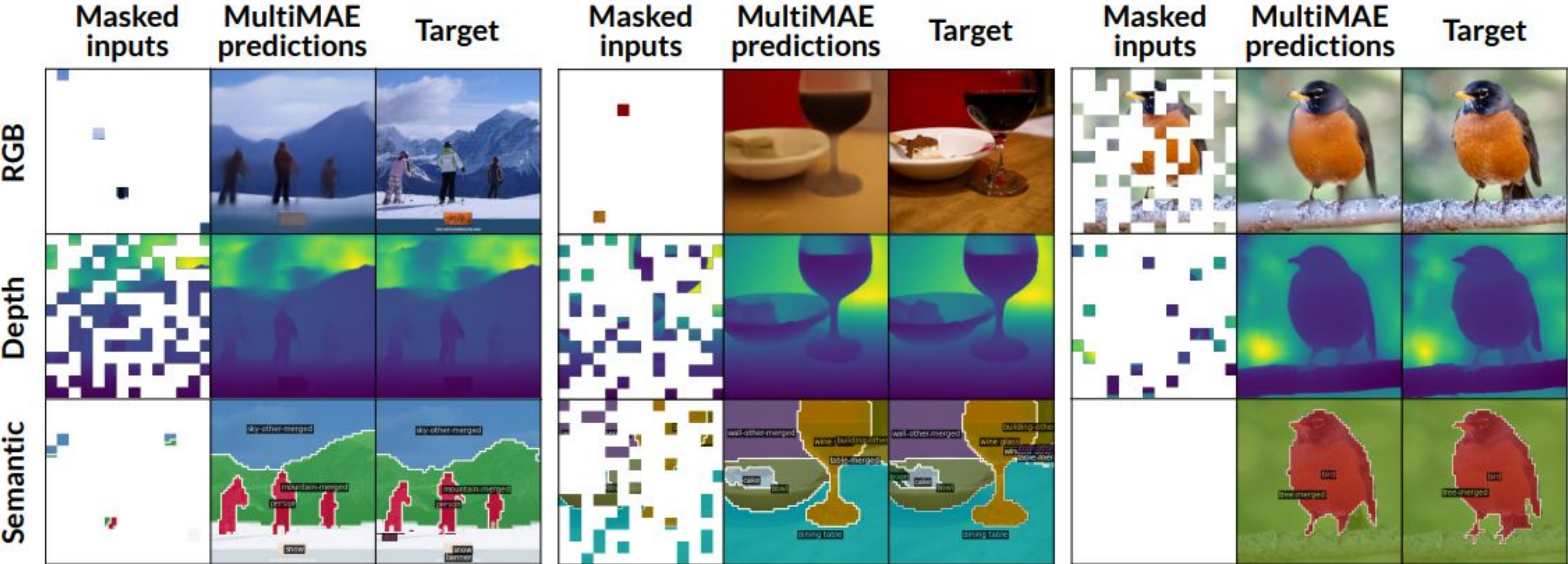


(b) learning motor control

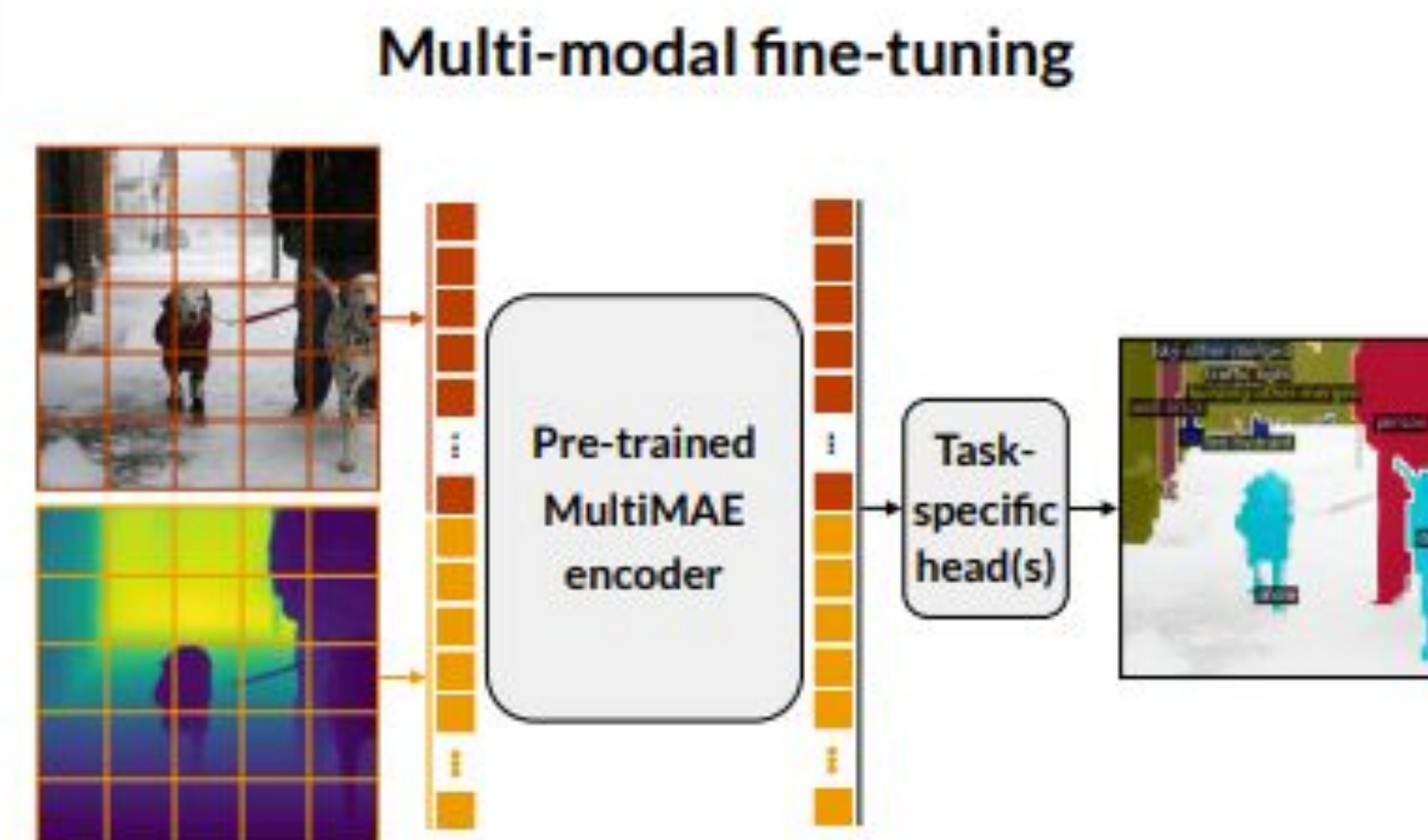
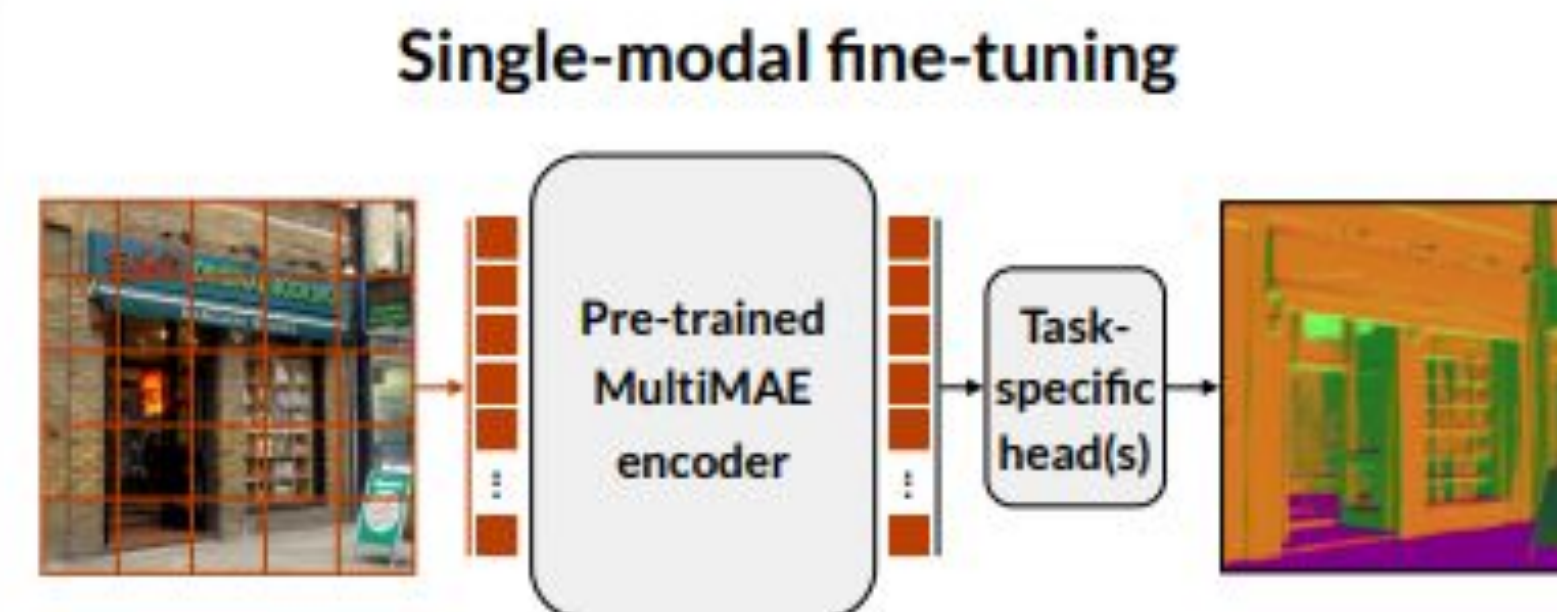
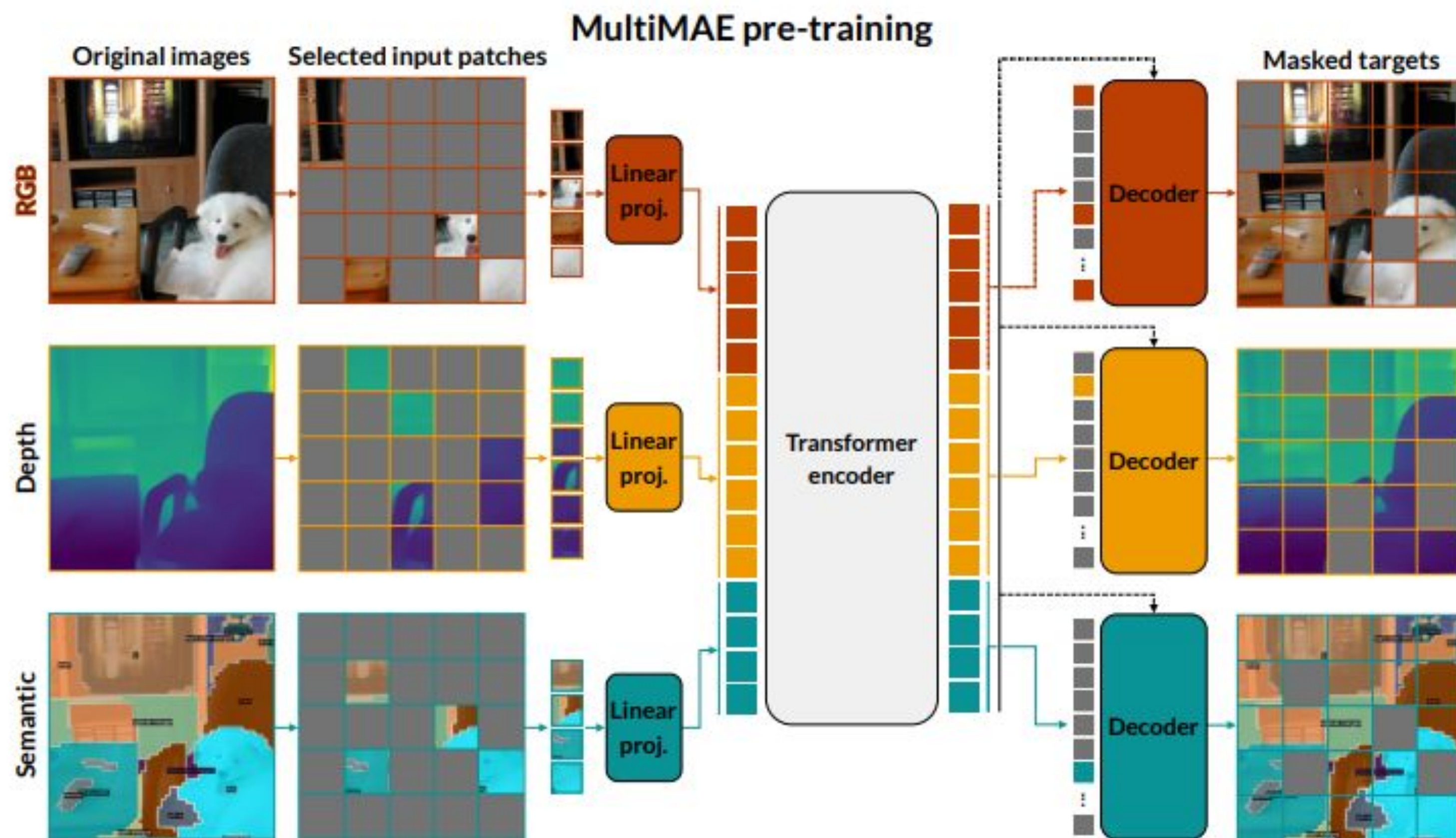




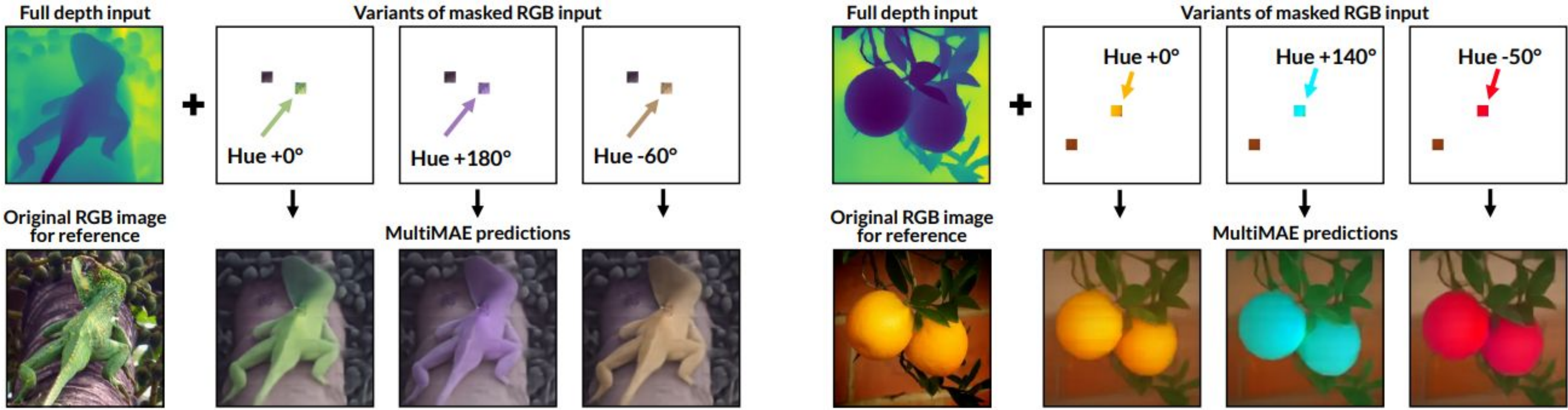
MultiMAE pre-training objective



MultiMAE pre-training



Demonstration of cross-modal interaction





Thank you!



Next Lecture:

Student Lecture 6

Diffusion Models and Policy Learning





DR

DeepRob

[Group 5] Lecture 5
Multisensory Learning + Manipulation
by Mason Hawver, Ryan Diaz, and Hanchen Cui
University of Minnesota




Image Source: D. F. Gomes, P. Paoletti, and S. Luo, "Generation of gelsight tactile images for sim2real learning," IEEE Robotics and Automation Letters, vol. 6, no. 2, pp. 4177–4184, 2021.