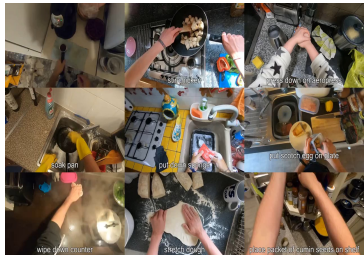


Human data

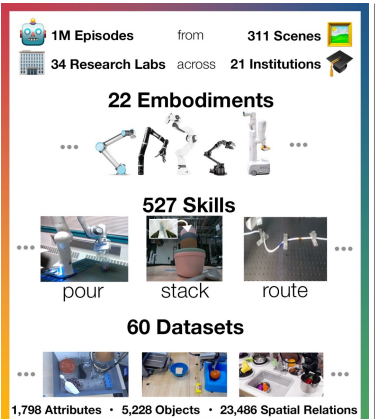
EPICKICHEN



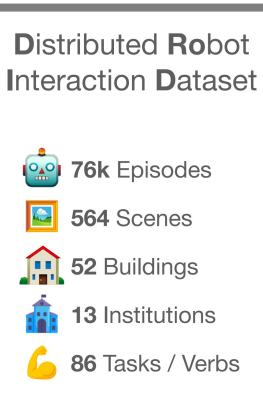
100h activities
20M frames
45 Kitchens
300 objects
97 Tasks / Verbs

Robotics data

Open-X

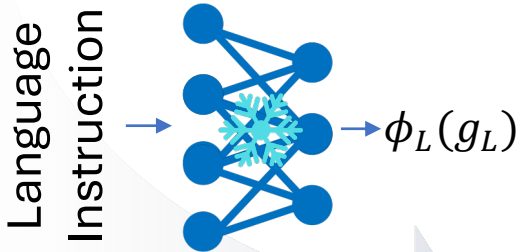


DROID

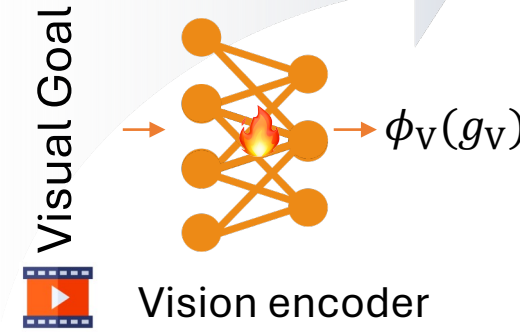


DecisionNCE

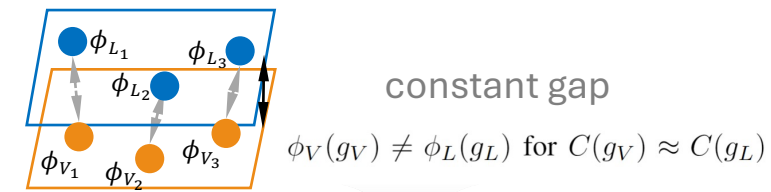
Language Instruction



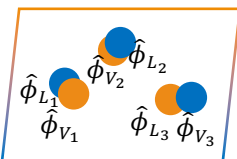
Visual Goal



Bridge the modality gap



Collapse



1. Centralize

$$\hat{\phi}_V \leftarrow \phi_V - \mathbb{E}_{g_V}[\phi_V(g_V)]$$

$$\hat{\phi}_L \leftarrow \phi_L - \mathbb{E}_{g_L}[\phi_L(g_L)]$$

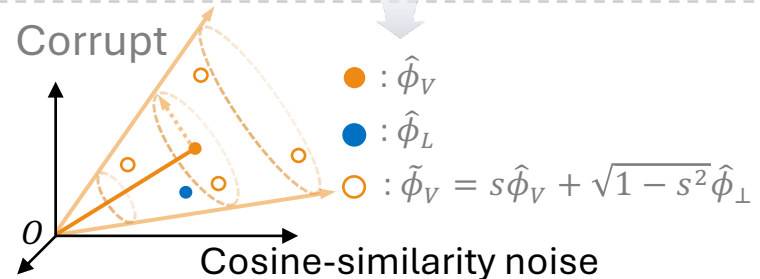
2. Delete

$$\hat{\phi}_V \leftarrow \text{del}(\phi_V, \arg\max_i \|\phi_V^i - \hat{\phi}_V\|)$$

$$\hat{\phi}_L \leftarrow \text{del}(\phi_L, \arg\max_i \|\phi_L^i - \hat{\phi}_L\|)$$

$$\phi_V(g_V) \approx \phi_L(g_L) \text{ for } C(g_L) \approx C(g_V)$$

Corrupt



Unimodal Task Train

Visual Goal



II. Modality gap reduction

Embedded task goal

$\tilde{\phi}_V(g_V)$

Film

Robot observations



Resnet34

MLPs

Diffusion Loss

Robot Policy

Multimodal Task Eval

Visual Goal



II. Modality gap reduction (w/o Corrupt)

Embedded task goal

$\hat{\phi}_V(g_V)$

$\hat{\phi}_L(g_L)$

II. Modality gap reduction (w/o Corrupt)



Language Instruction



- Pick up the duck
- Open the drawer
- Move the pot
- Fold the cloth

I. Robotics multimodal encoder pretrain

II. Modality gap reduction

III. Robot policy train and evaluation