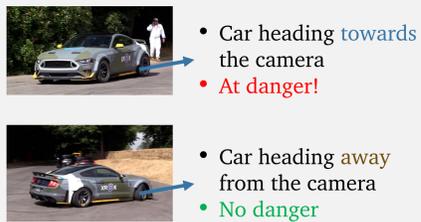




## Motivation



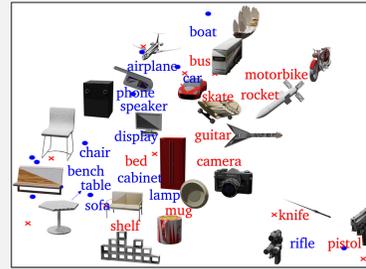
Recognition needs to understand two aspects:

- *What* is the object
- *How* is it presented

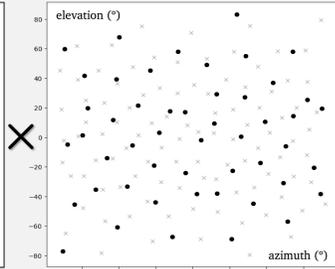
## Benchmark

### Data

13 in-domain & 20 out-of-domain semantic categories

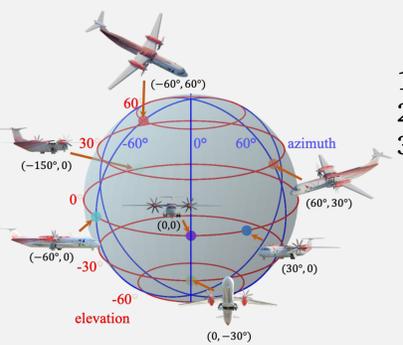


in-domain & out-of-domain pose



Rendered images of objects of different categories and poses. In-domain and out-of-domain splits for evaluating generalization.

### Evaluation metrics



- 1) Semantic classification
- 2) Absolute pose
- 3) Relative pose
  - Canonical pose free
  - Category agnostic
  - Generalize to open categories

## Unsupervised Learning for Both Object Semantics and Pose

Goal: learn **what** is the object (semantics) and **how** it presented (pose)

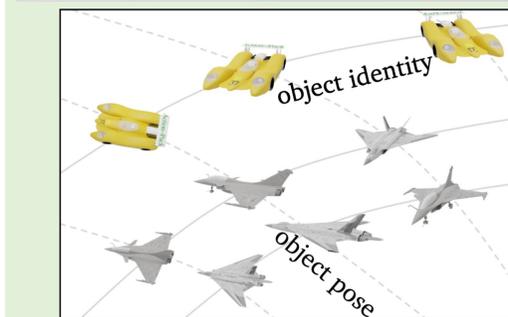
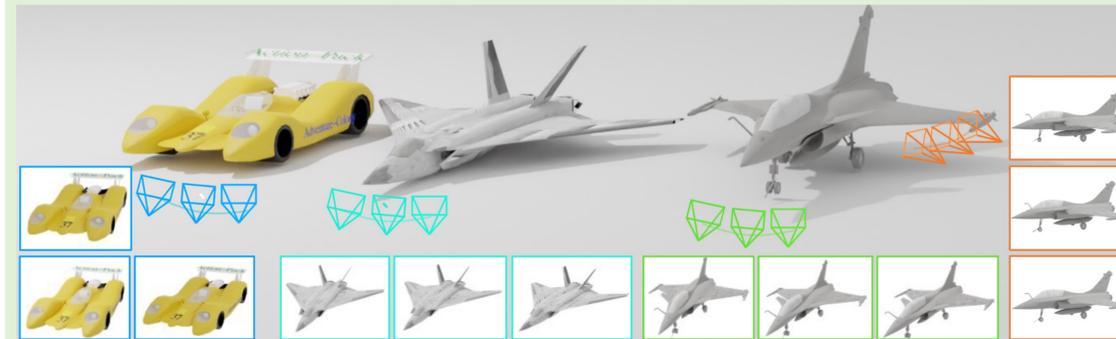


Scenario: a robot moves around in the environment

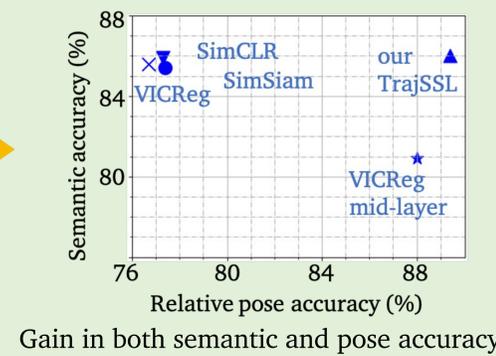
A natural data acquisition scheme:

- No labels
- Adjacent images of the same object from a smooth viewpoint trajectory

Training data: **Image triplets** with small pose changes; **No** semantic or pose labels.



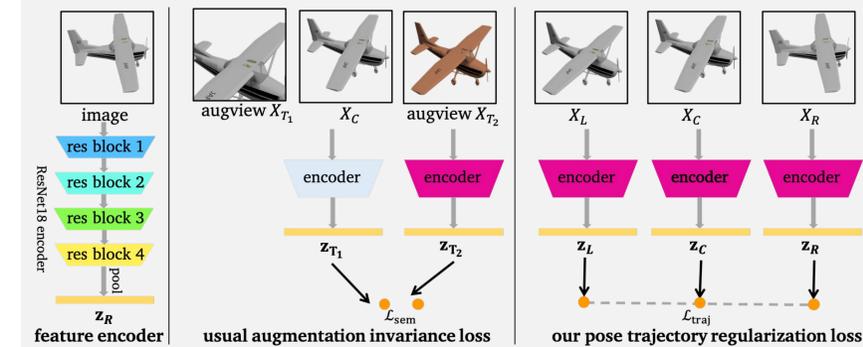
Emergent representation with **disentangled** semantics and pose



Gain in both semantic and pose accuracy.

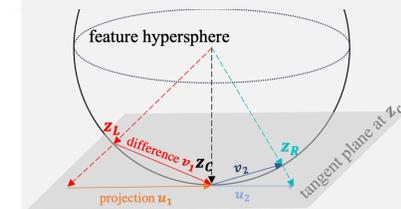
## Methods

### Stage 1: Self-supervised representation learning



Our encoder produce embeddings for a triplet of images  $\{X_L, X_C, X_R\}$  from a sequence with respective poses  $\{p_L, p_C, p_R\}$  forming a trajectory, where pose changes are subtle.

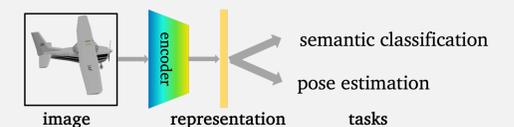
Two unsupervised losses imposed on the embeddings,  $\mathcal{L}_{sem}$  &  $\mathcal{L}_{traj}$ .  $\mathcal{L}_{sem}$  is an invariant loss (e.g. VICReg).



**Trajectory Loss**  
Viewpoint trajectory regularization makes 3 embeddings form a line:  
$$\mathcal{L}_{traj}(z_L, z_C, z_R) = -\frac{\mathbf{u}_1 \cdot \mathbf{u}_2}{\|\mathbf{u}_1\| \|\mathbf{u}_2\|}$$

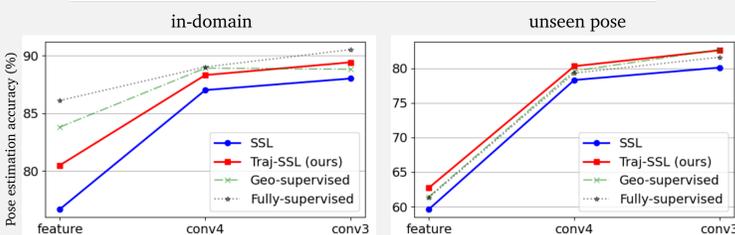
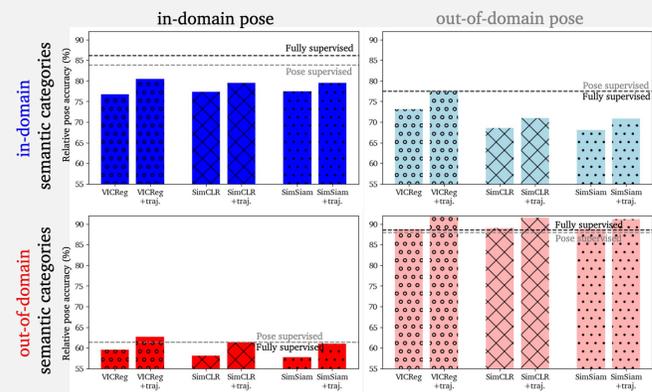
Final loss is a combination:  $\mathcal{L} = \mathcal{L}_{sem} + \mathcal{L}_{traj}(z_L, z_C, z_R)$

### Stage 2: Probing representation to downstream tasks



## Results

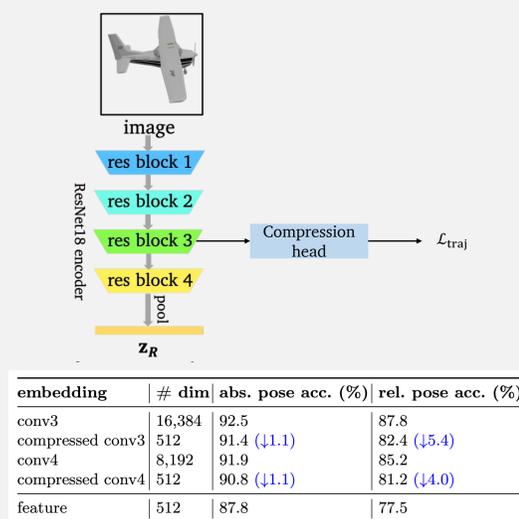
**Generalizability:** Trajectory regularization improves relative pose accuracy for in and out-of-domain data.



**Real Data:** Trajectory regularization improves retrievals with similar pose and appearance.



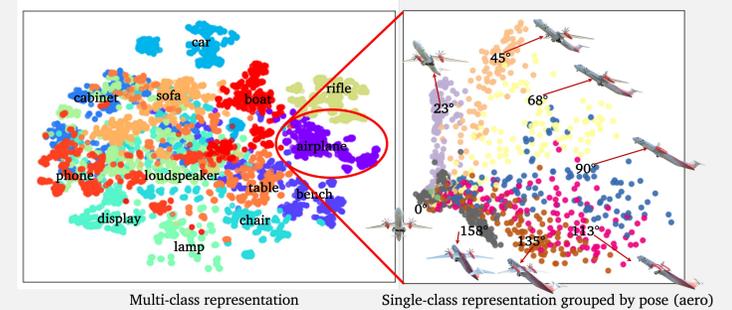
**Compressing mid-layer representation** up to 32x gives small accuracy loss.



**Mid-Layer Representation** improves pose accuracy for in and out-of-domain data.

## Representation Visualization

The **joint semantic-pose embedding**: Images are clustered by semantics; within each semantic cluster, images form mini-cluster by pose.



Comparison: No trajectory loss leads to representation collapse.

