

Learning Structured Video Representations without Supervision

Andrii Zadaianchuk, 2025

Why do we need structured representations?

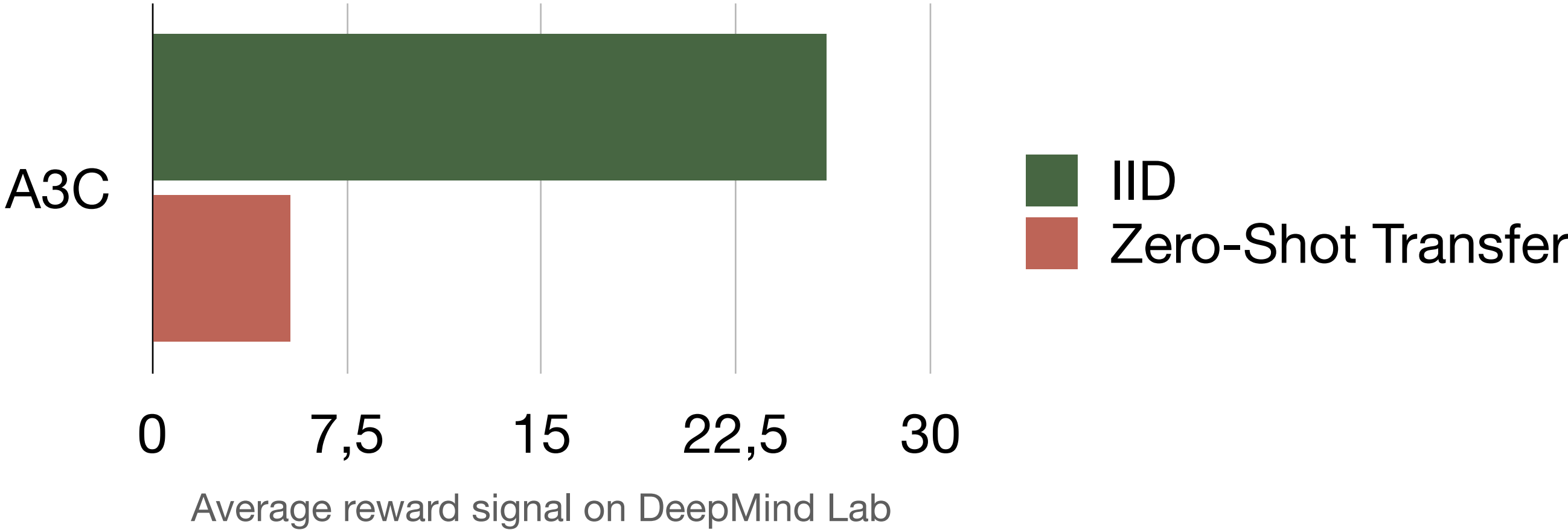
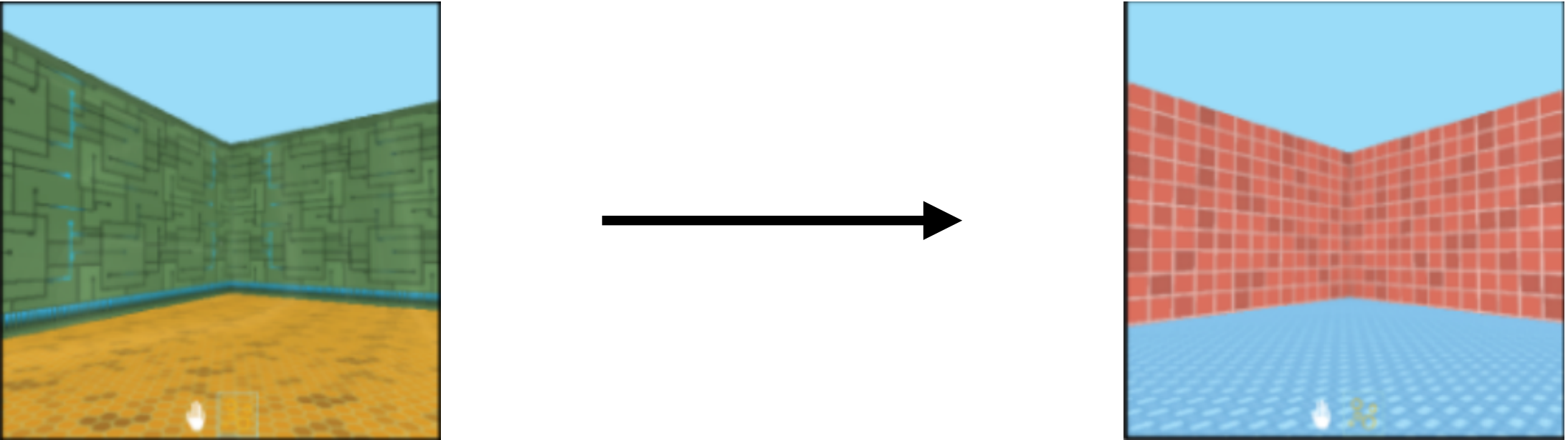


DeepMind Lab
Nav Maze Level 1

Here the player has to navigate a maze with multiple rooms in order to find the goal.

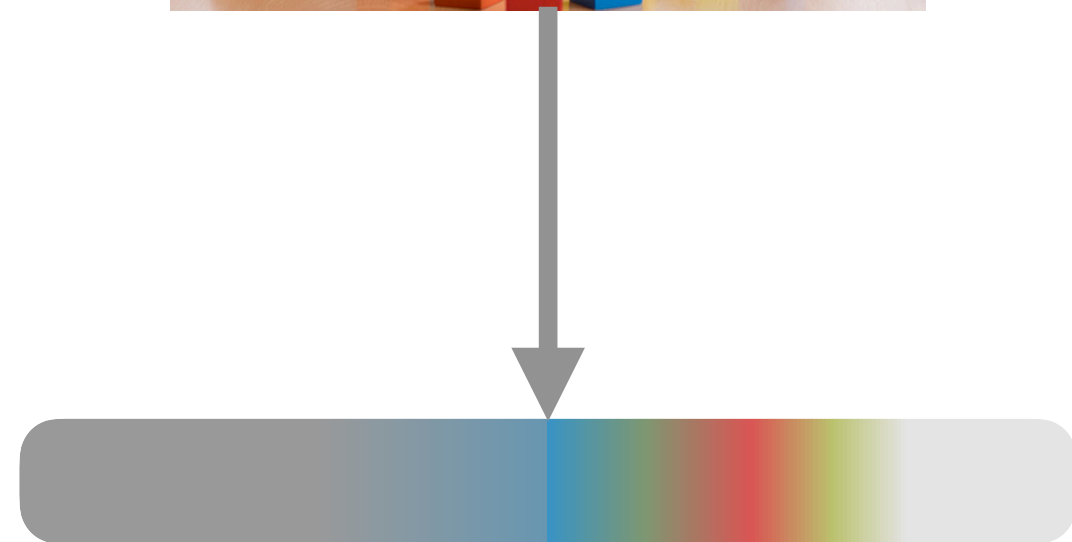
Player has to navigate a maze with multiple rooms in order to find the goal.

What about Zero-Shot Transfer?



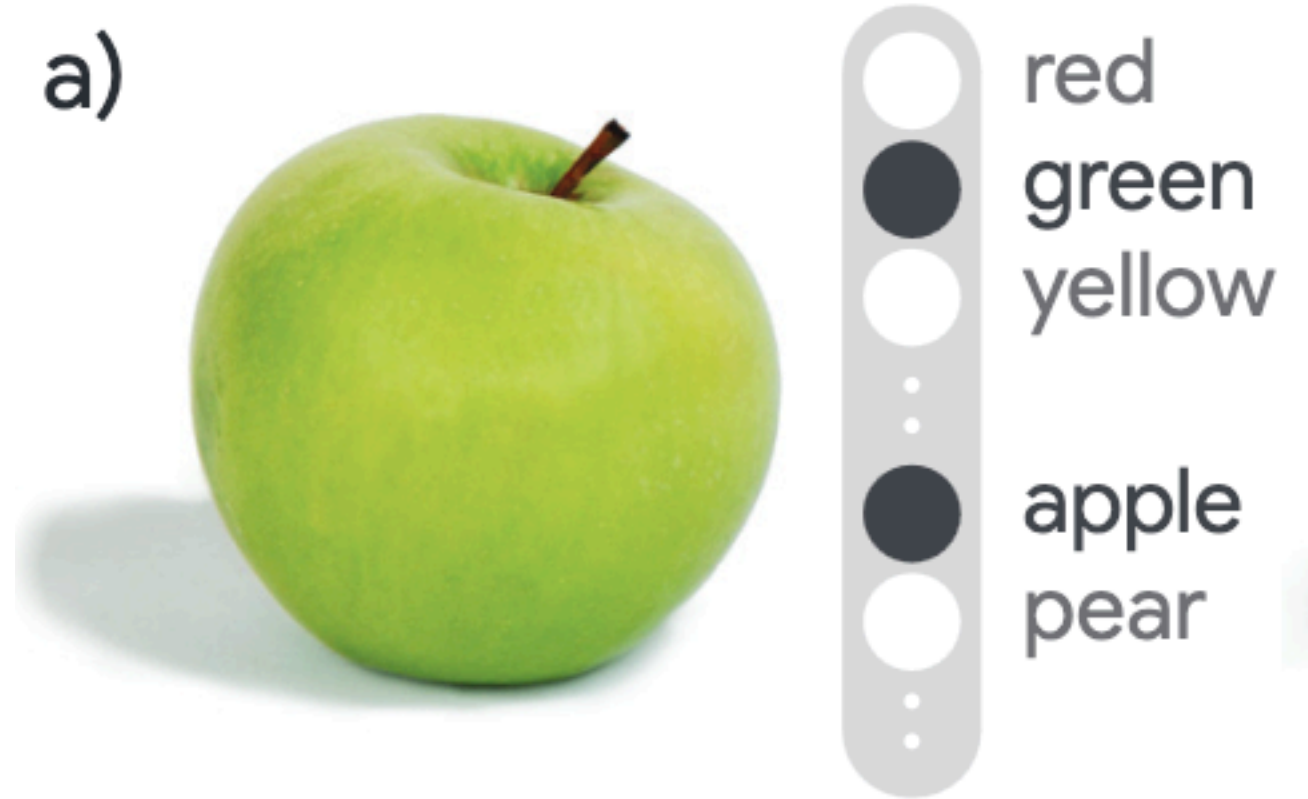
How to represent the world around us robustly?

Representation Learning from Pixels

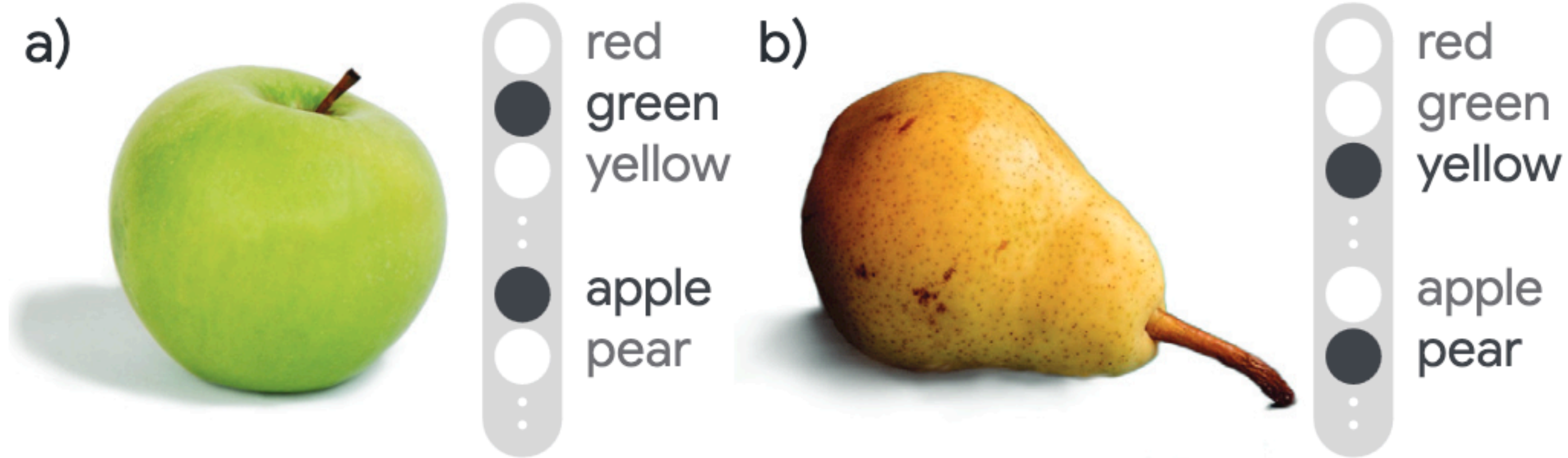


Single vector

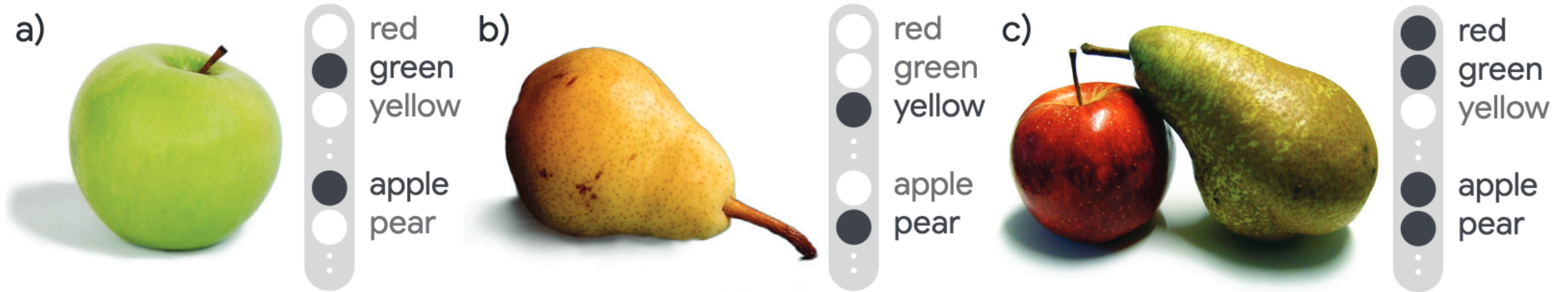
Binding Problem in Distributed Representation



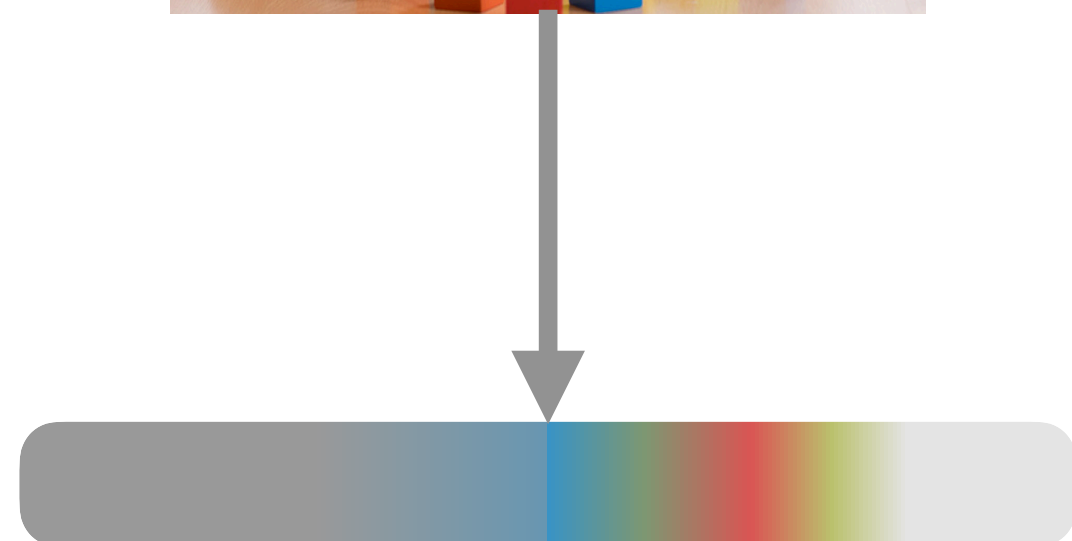
Binding Problem in Distributed Representation



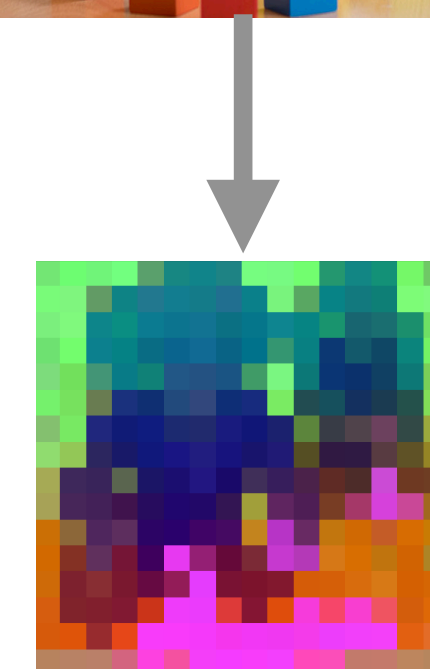
Binding Problem in Distributed Representation



Representation Learning from Pixels



Single vector



Dense grid of features

?

What if we can learn representations that are **structured** similarly to the original scene?

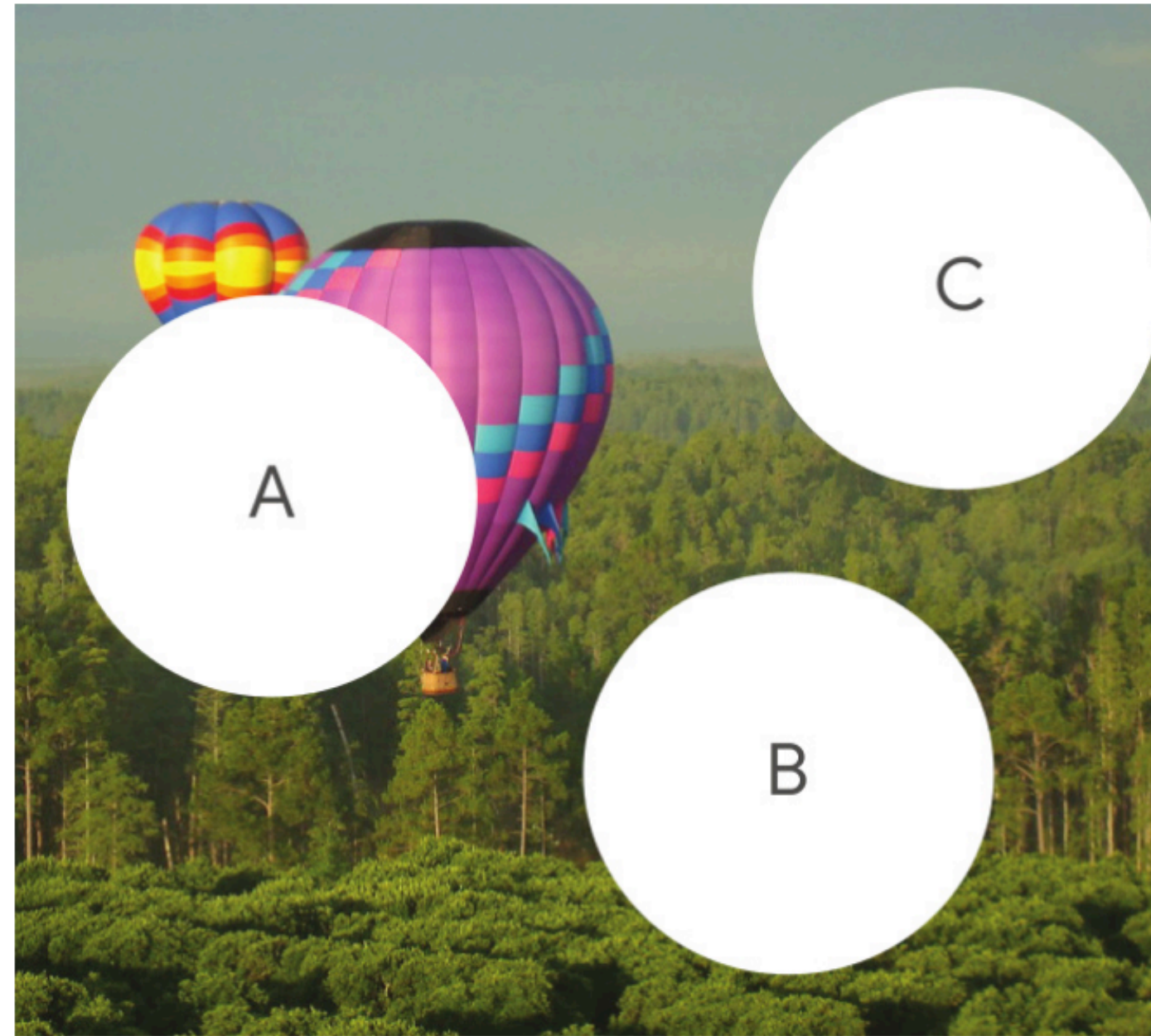
Do we need structured representations if we have scale?

An image showing 12 plain tea cups from the same set, all identical in style and design.



Here is the updated image showing exactly 12 plain tea cups, all identical in design and style.
Let me know if there's anything else you'd like adjusted!

How humans structure information about scene?



We group:

- regions that are largely independent of their context
- regions that exhibit strong internal predictive structure

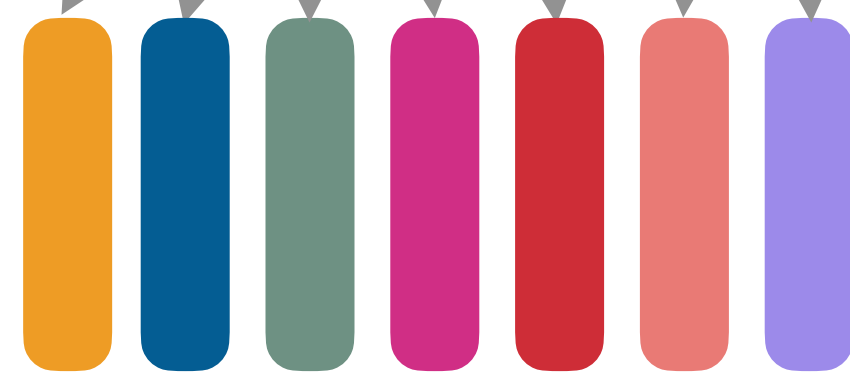
Objects are good candidates for both!

Objects are building blocks of the visual scene?



- Instance segmentation and tracking
- Visual reasoning and planning
- Combinatorial generalization

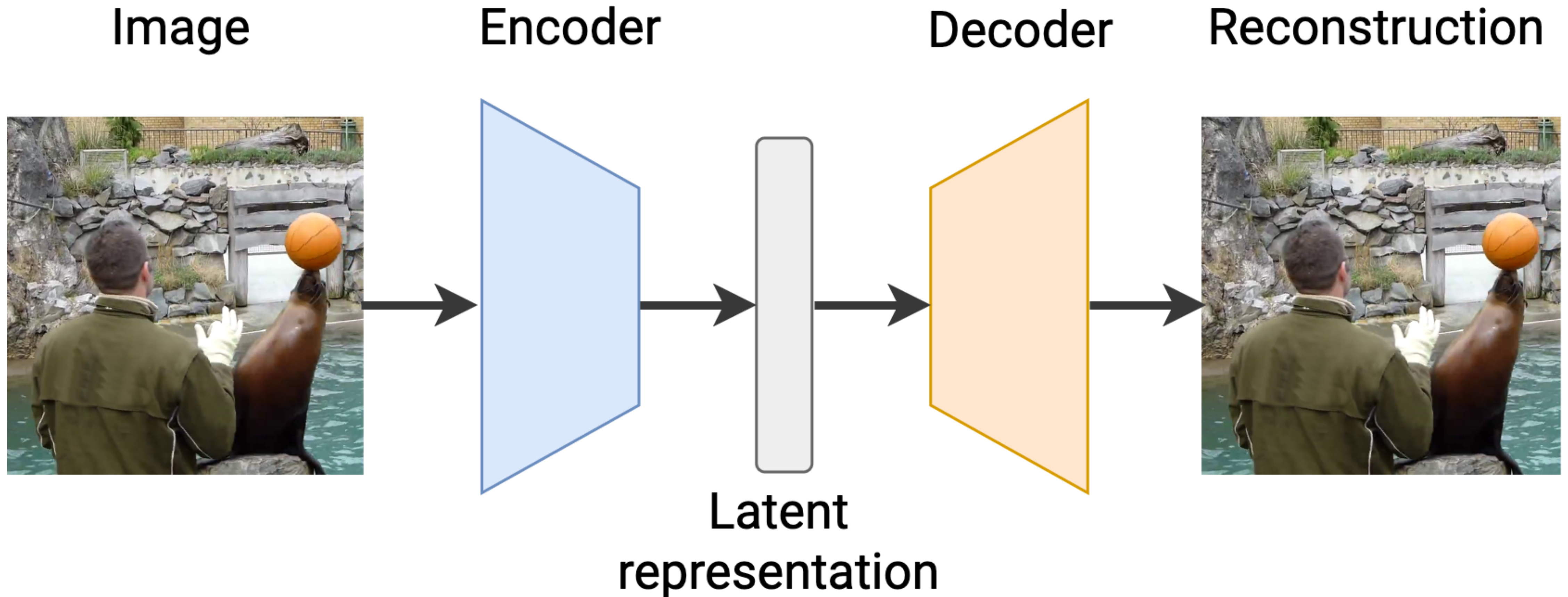
Unsurevised Object-Centric Representations



Object-centric
representations

- Different objects are represented by different vectors
- Those vectors are grounded on particular image segments
- Trained end-to-end with **architectural inductive biases** and **self-supervision objectives**

Unsupervised Representation Learning

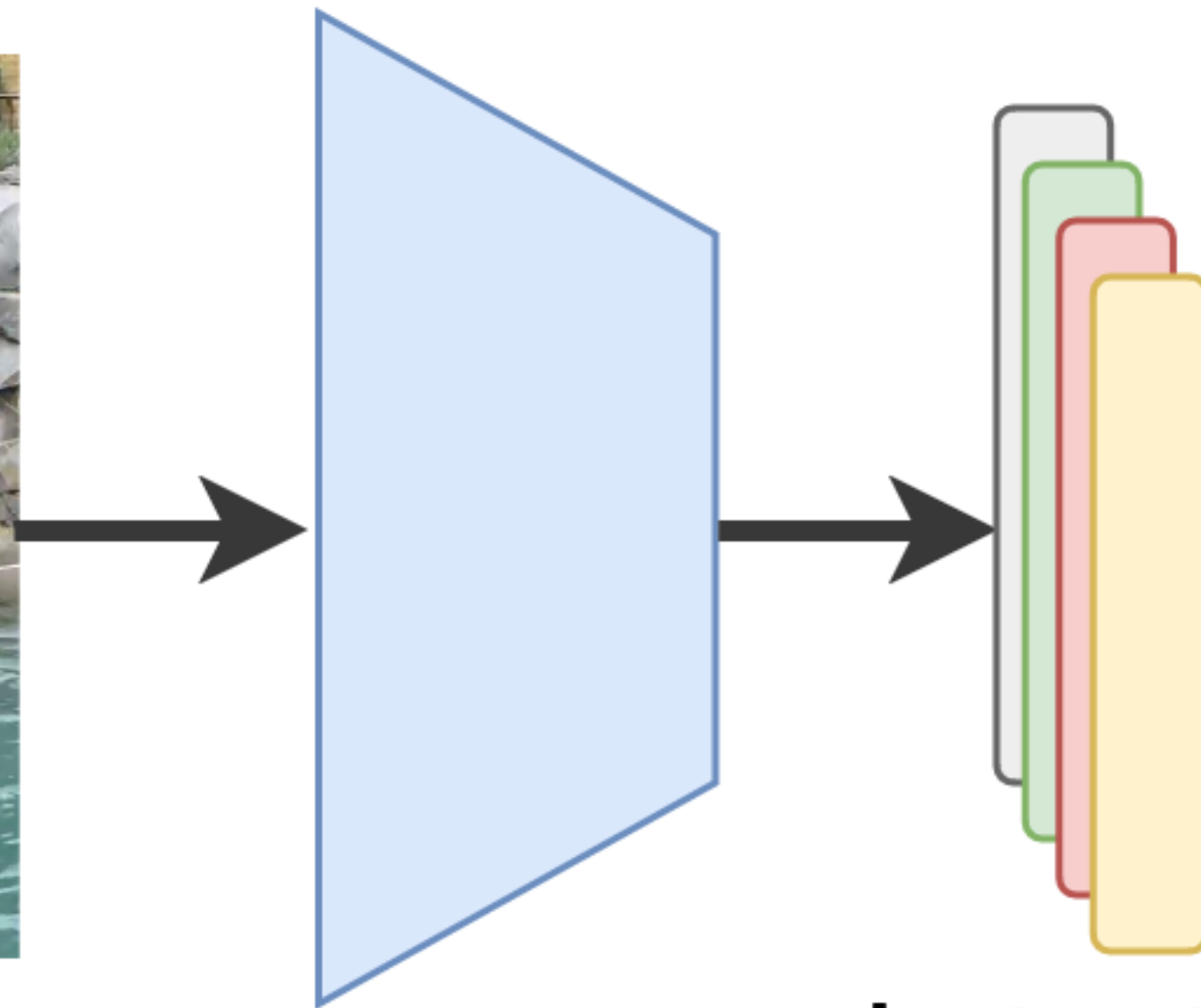


Object-Centric Representation Learning

Image

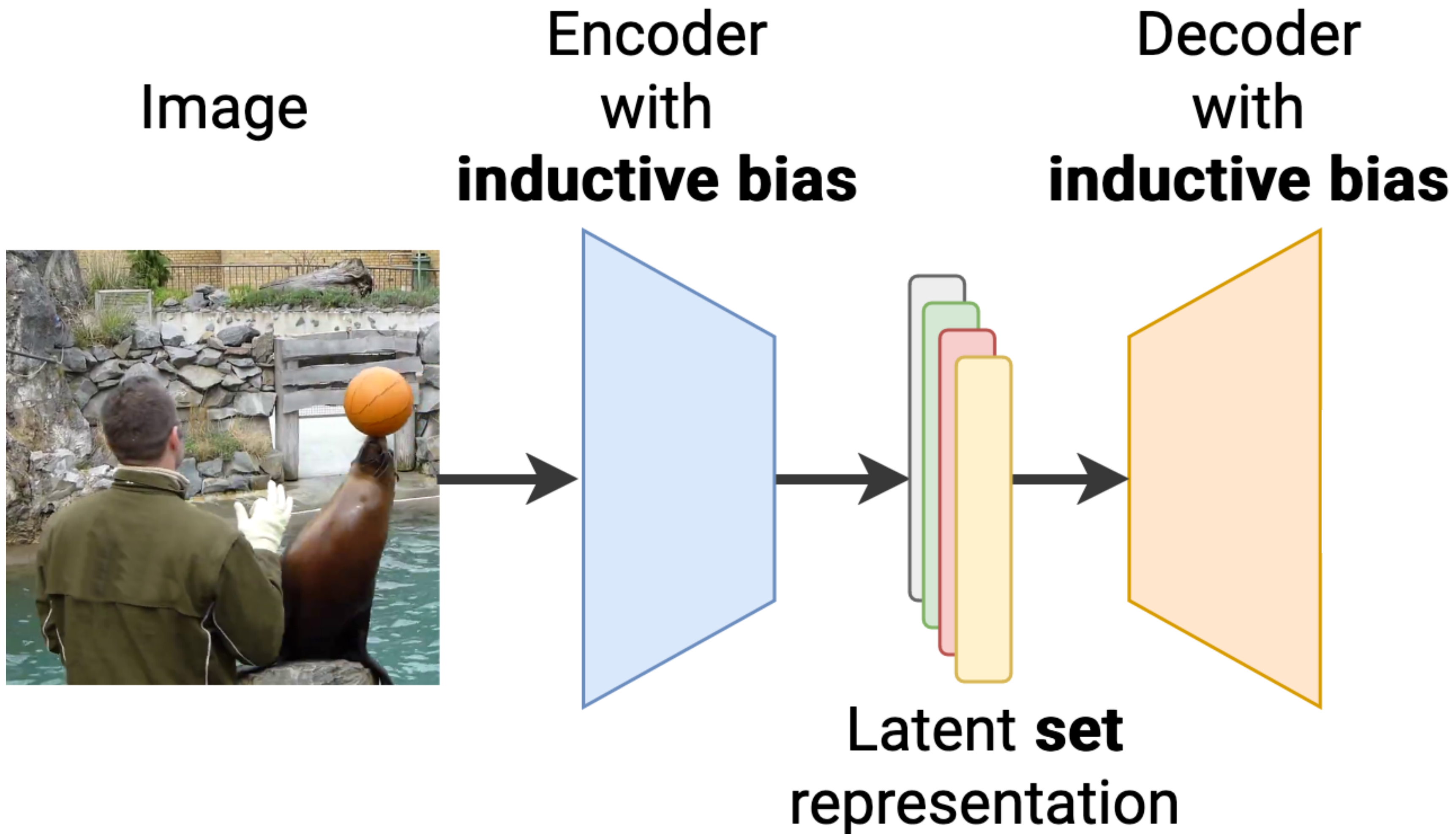


Encoder
with
inductive bias

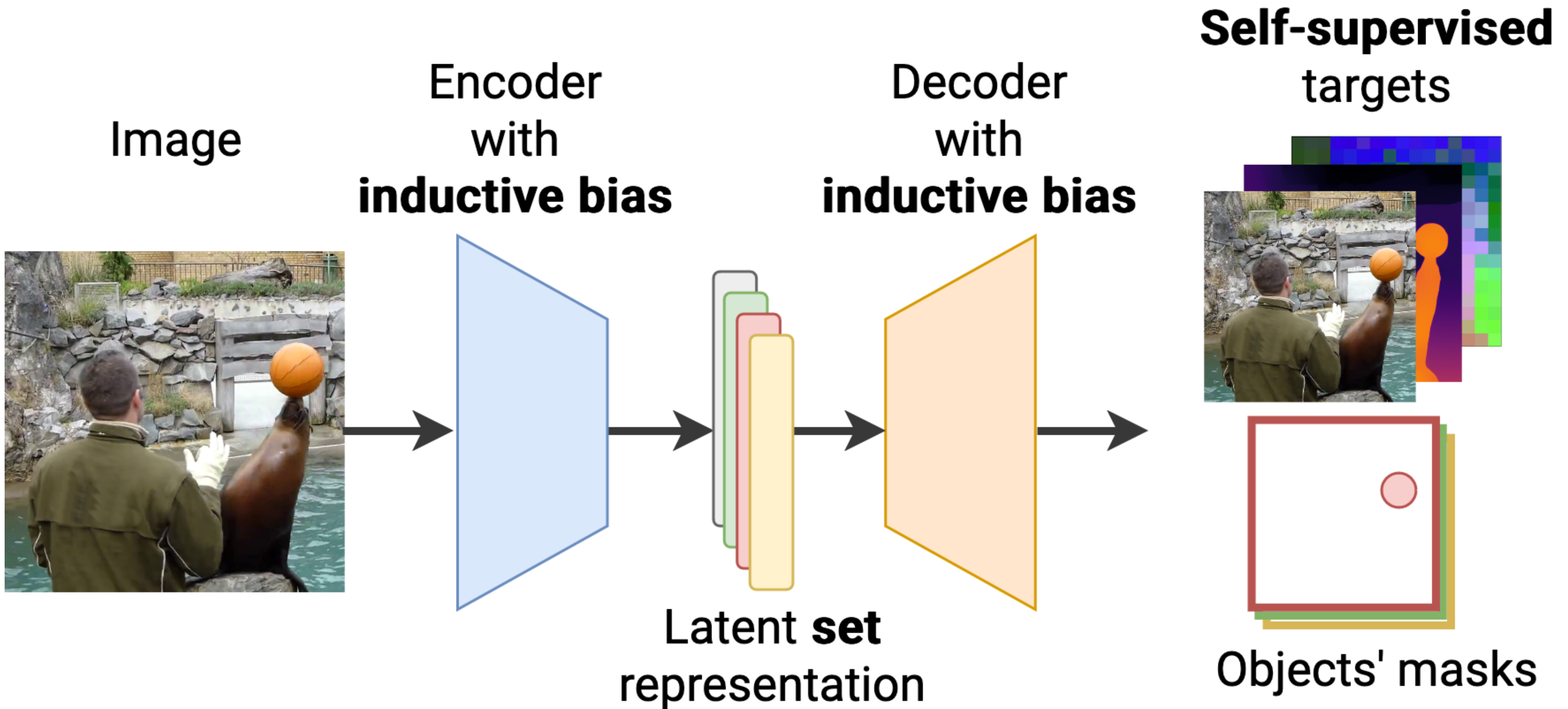


Latent **set**
representation

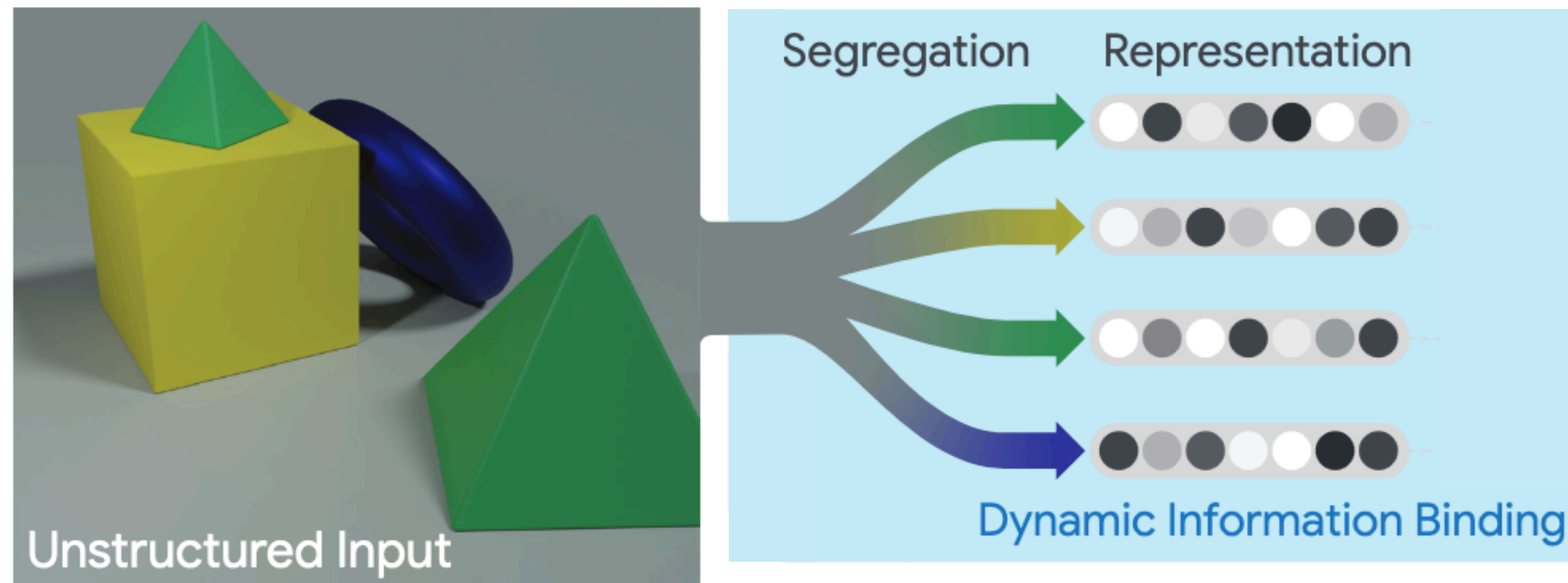
Object-Centric Representation Learning



Object-Centric Representation Learning

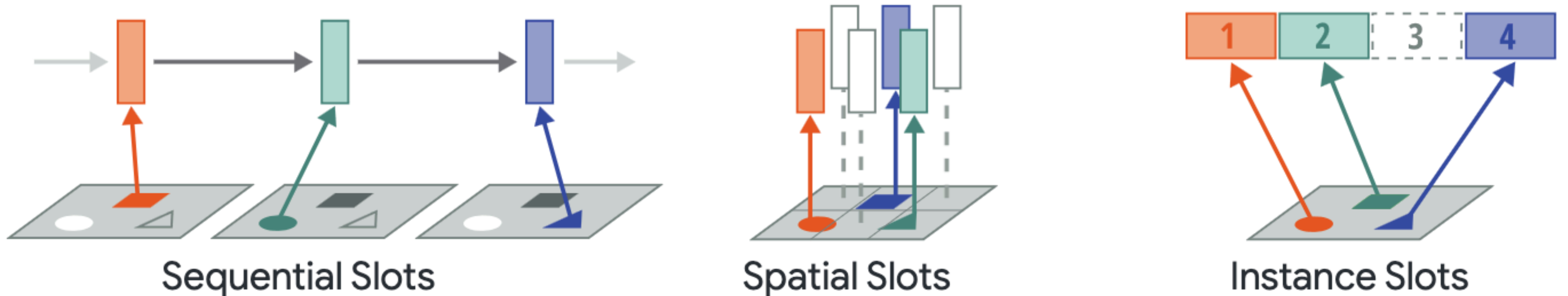


Scene Decomposition Into Objects



- Dense pixels / features should be separated into discrete set of vectors or *slots*
- *Routing problem*: which vector is responsible for which object?

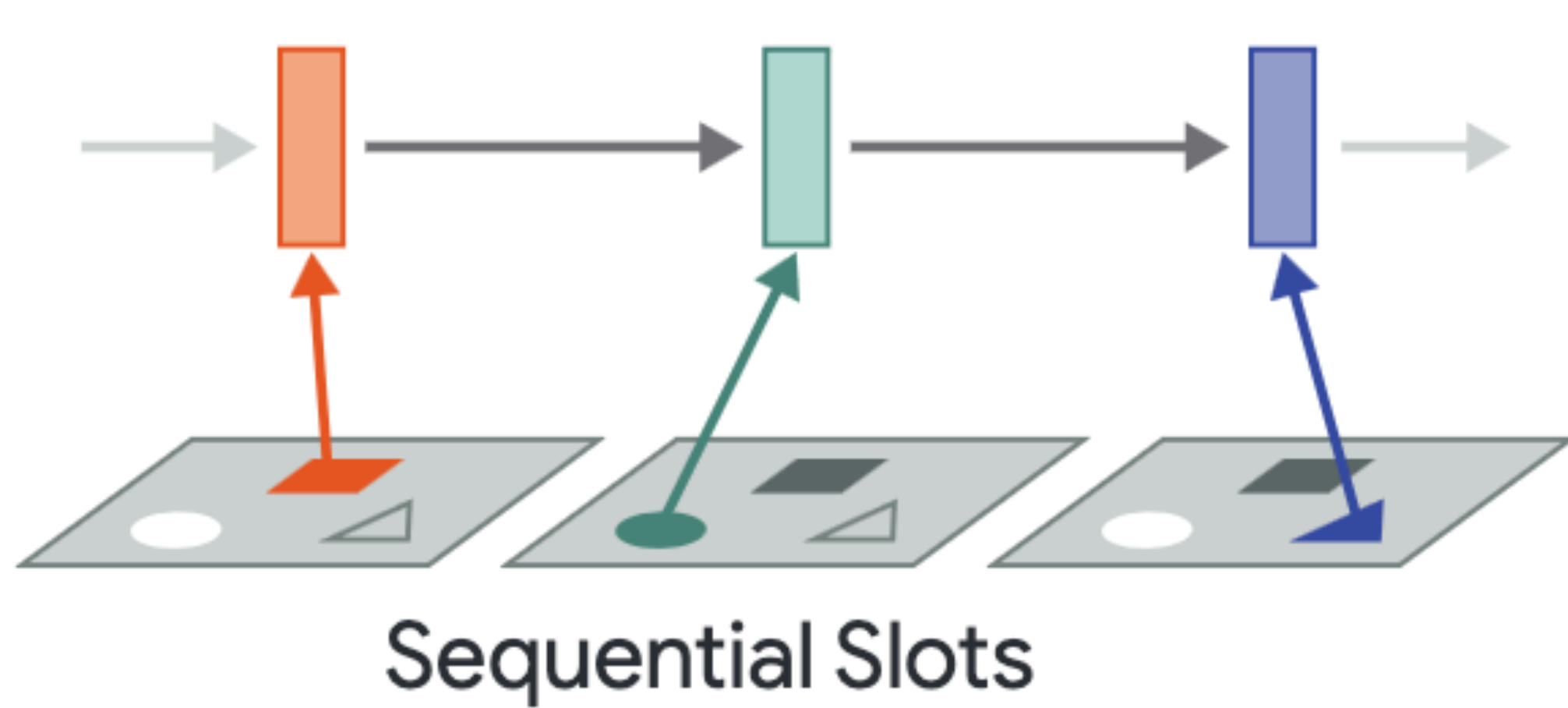
Different Ways to Decompose



Encoder inductive biases could be categorised in terms of encoder outputs named *slots*:

- Sequential slots → ordered **sequence** of vectors
- Spatial slots → sparse **grid** of vectors
- Instance slots → permutation-invariant **set** of vectors

Different Ways to Decompose



AIR [Eslami et al., 2016]

SQAIR [Kosiorek et al., 2018]

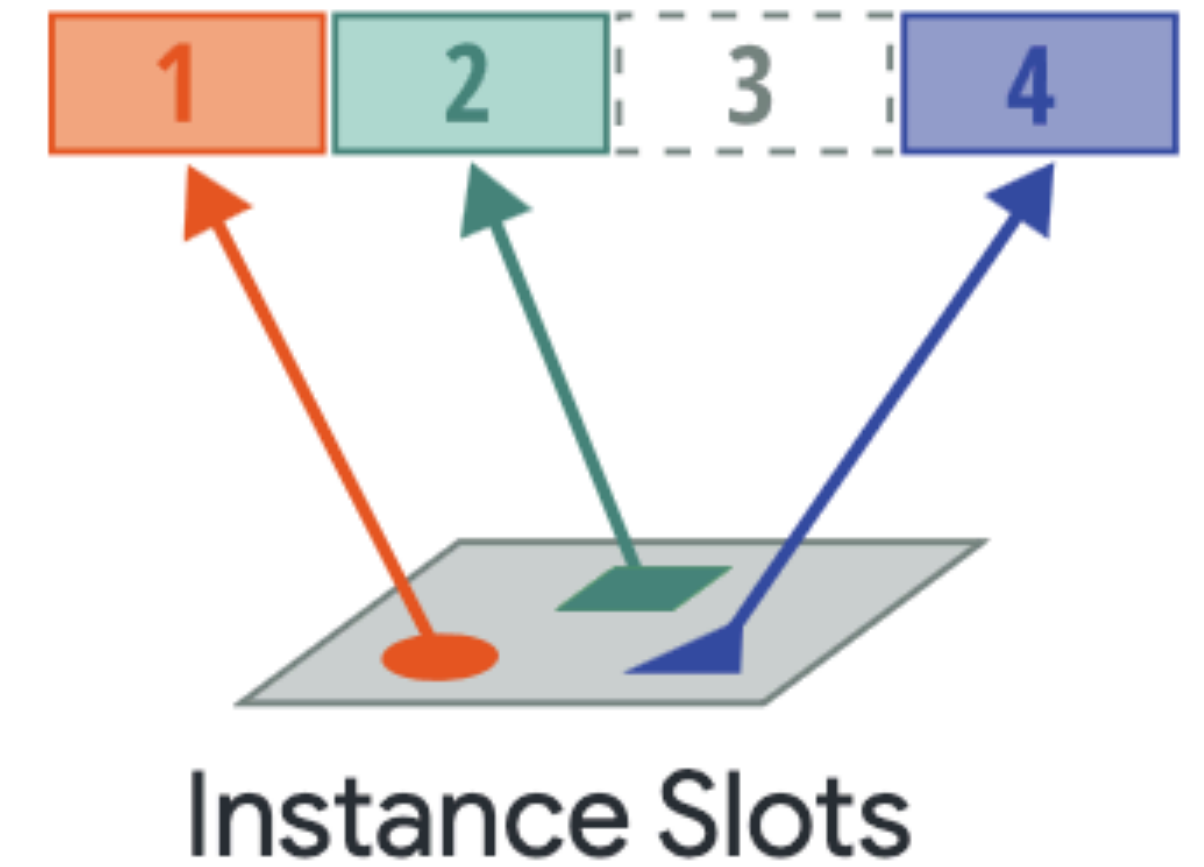
MONet [Burgess et al., 2019]



SPAIR [Crawford & Pineau, 2019]

SPACE [Lin et al., 2020]

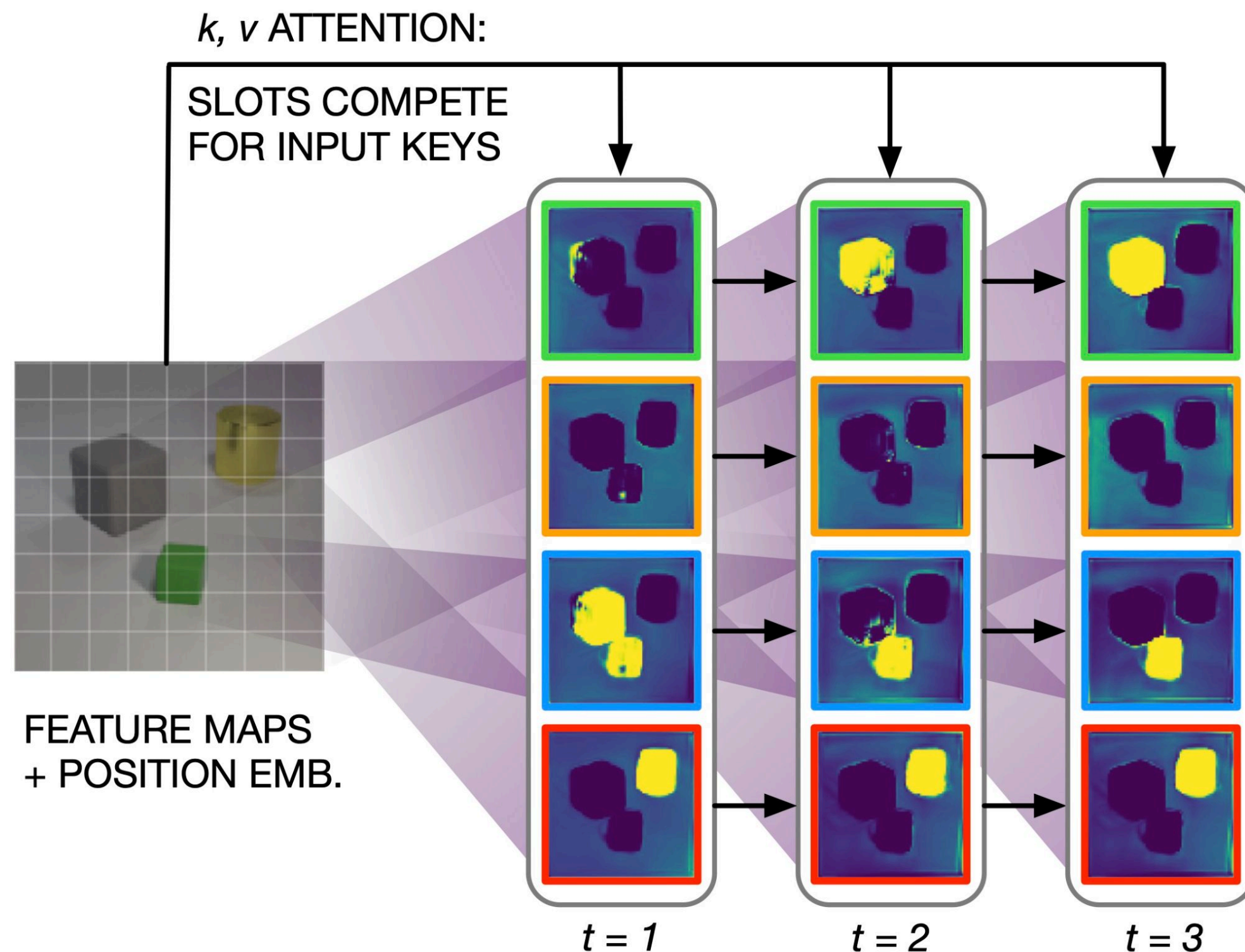
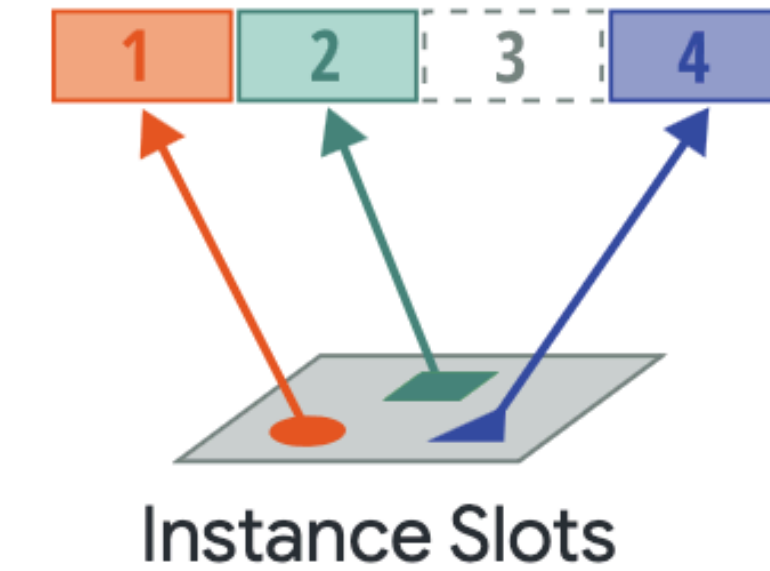
SCALOR [Jiang et al., 2020]



SA [Locatello et al., 2020]

DINOSAUR [Seitzer et al., 2023]

Instance Slots: Slot Attention Encoder



Slot Attention Pseudocode

inputs: feature maps + position embedding

slots \sim normal(mean, std)

for $t = 0 \dots T$:

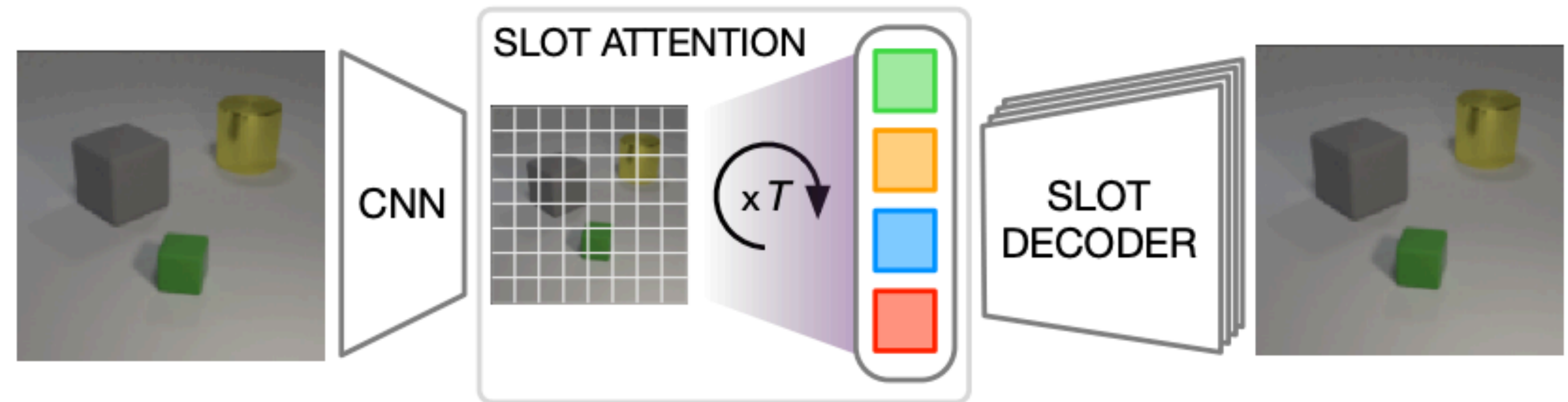
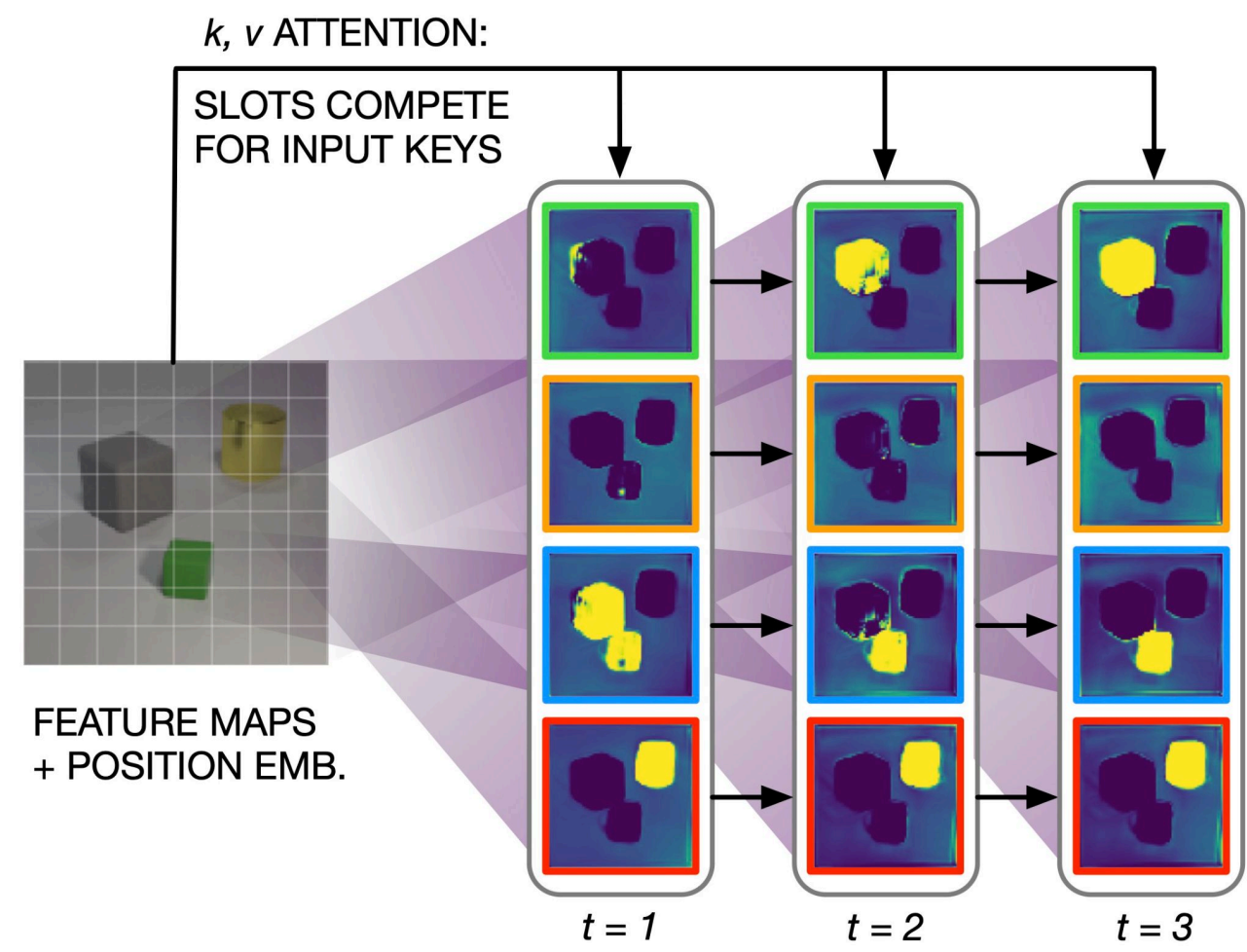
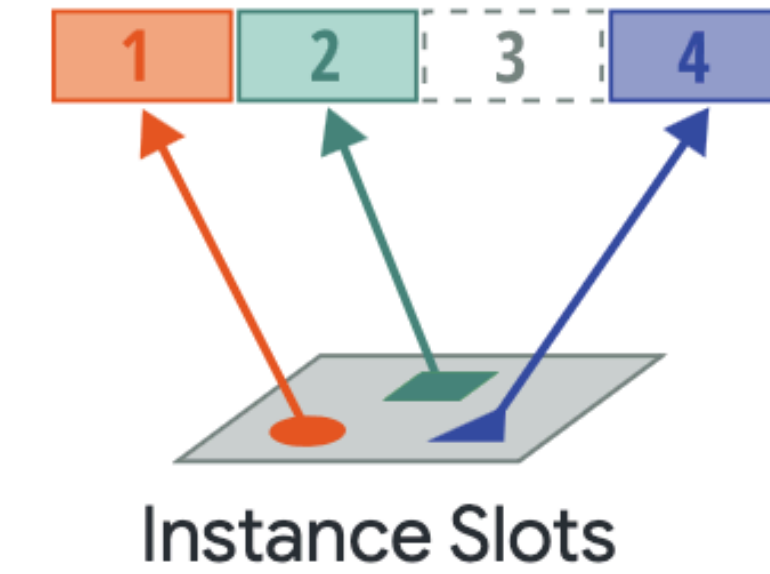
scores = dot(k (inputs), q (slots))

weights = softmax(scores / t, axis='slots')

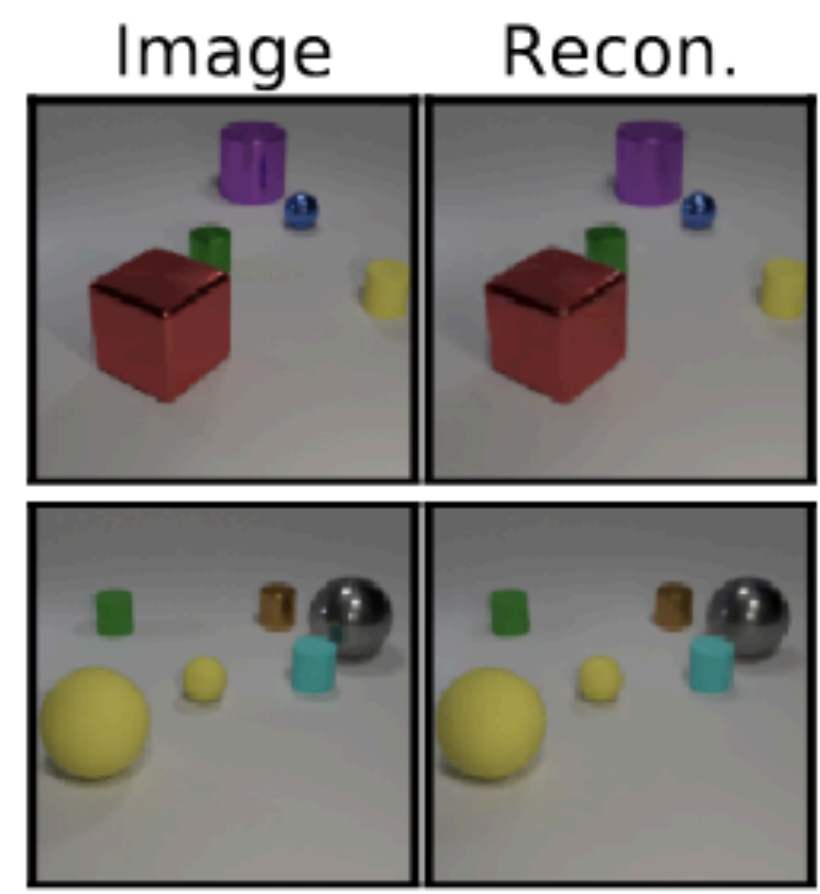
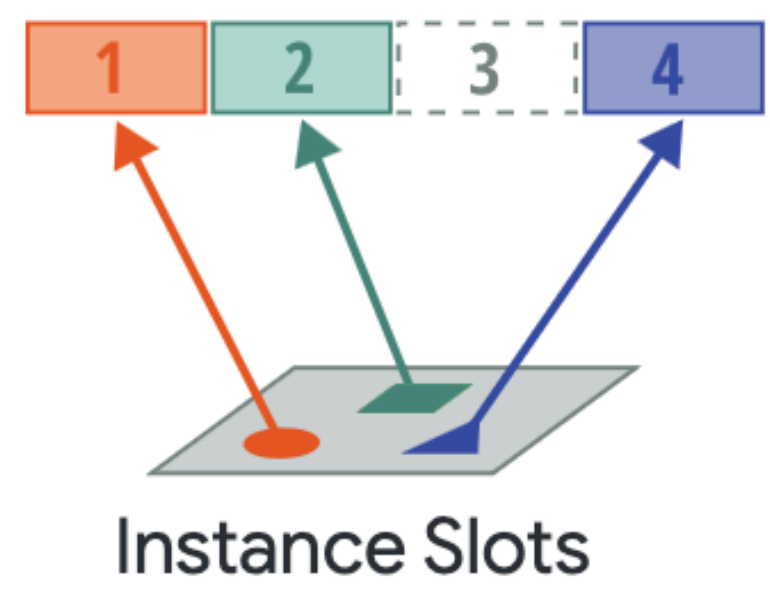
updates = weighted_mean(weights, v (inputs))

slots = gru(slots, updates) # GRU update

Instance Slots: Slot Attention Training



Instance Slots: Slot Attention Results



Discovering Object-Centric Structure from the Real-World Video Data

Object-Centric Learning for Real-World Data

Ground Truth



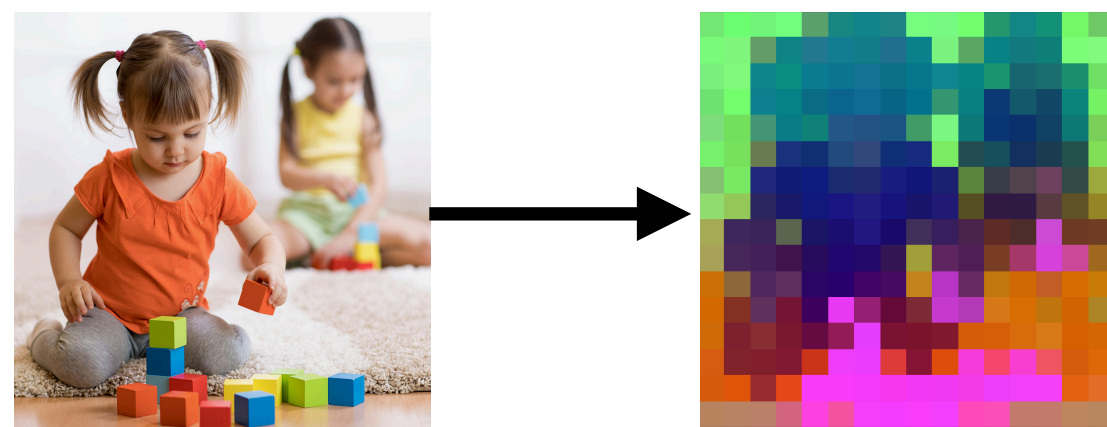
Slot Attention



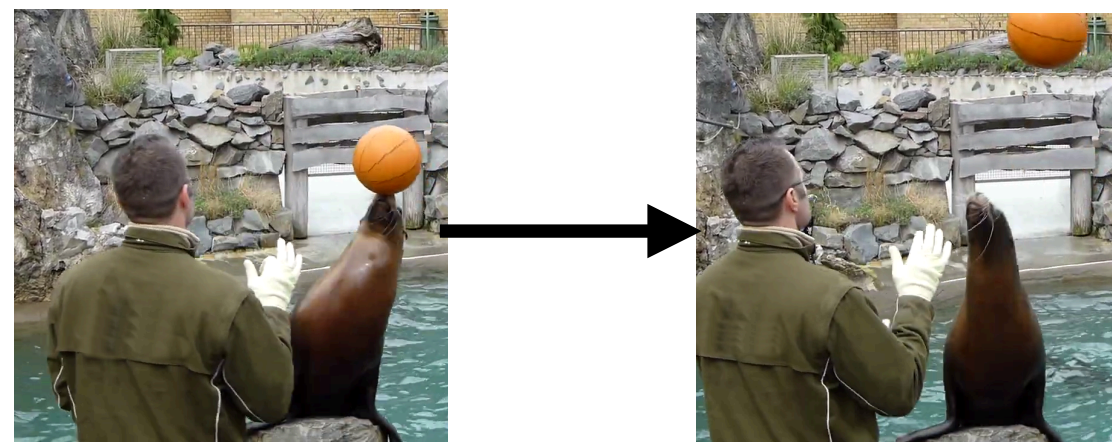
Image reconstruction as the target **is not enough** for grouping real-world scenes

Self-supervised Object-Centric Objectives

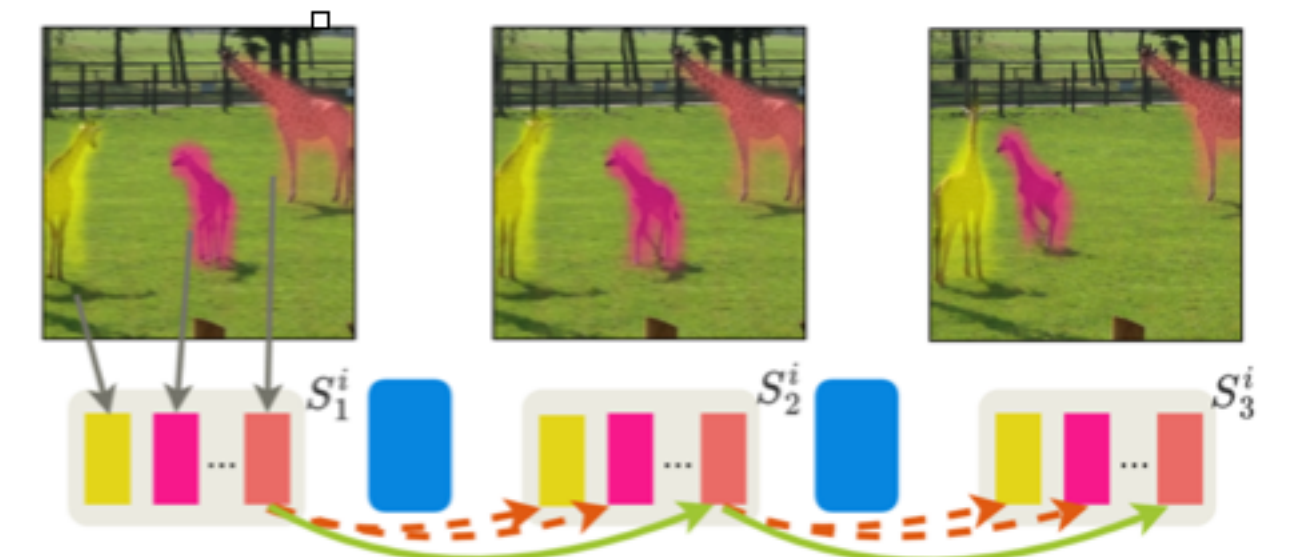
Semantics
Reconstruction



Motion
Prediction

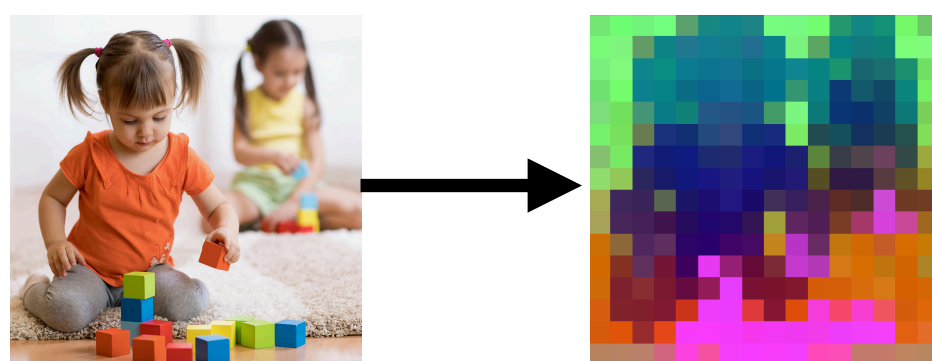


Identity
Preservation



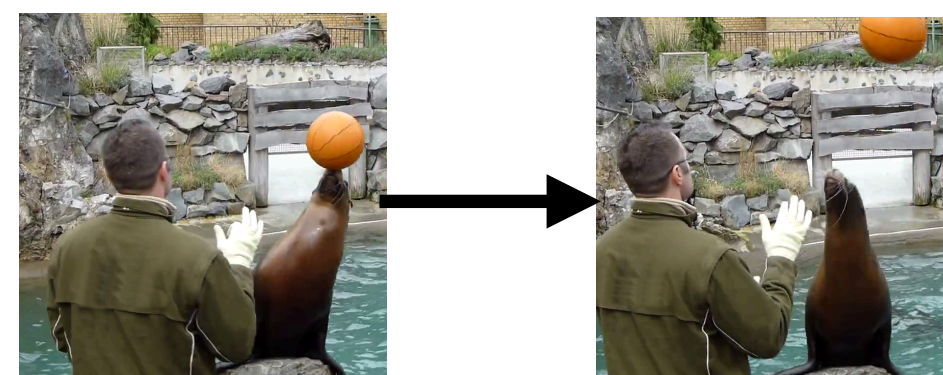
Self-supervised Object-Centric Objectives

Semantics
Reconstruction



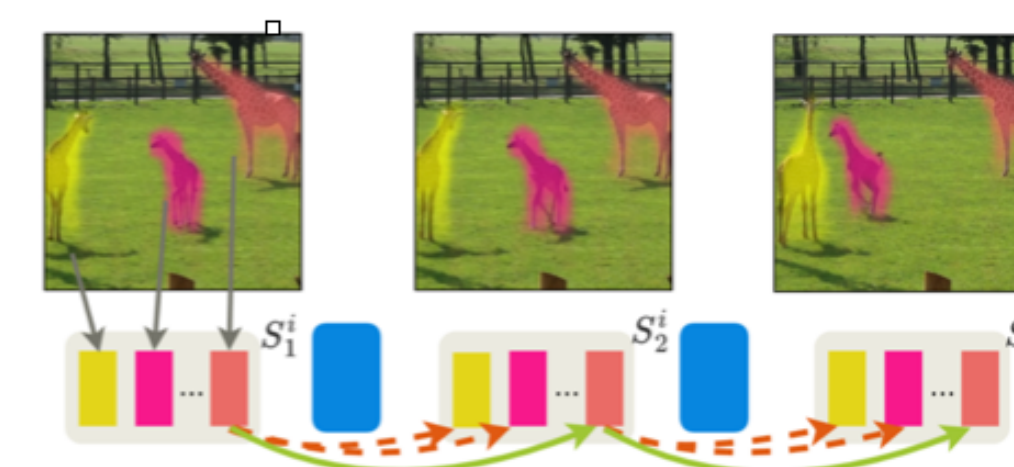
DINOSAUR
[Seitzer et al., 2023]

Motion
Prediction



VideoSAUR
[Zadaianchuk et al., 2023]

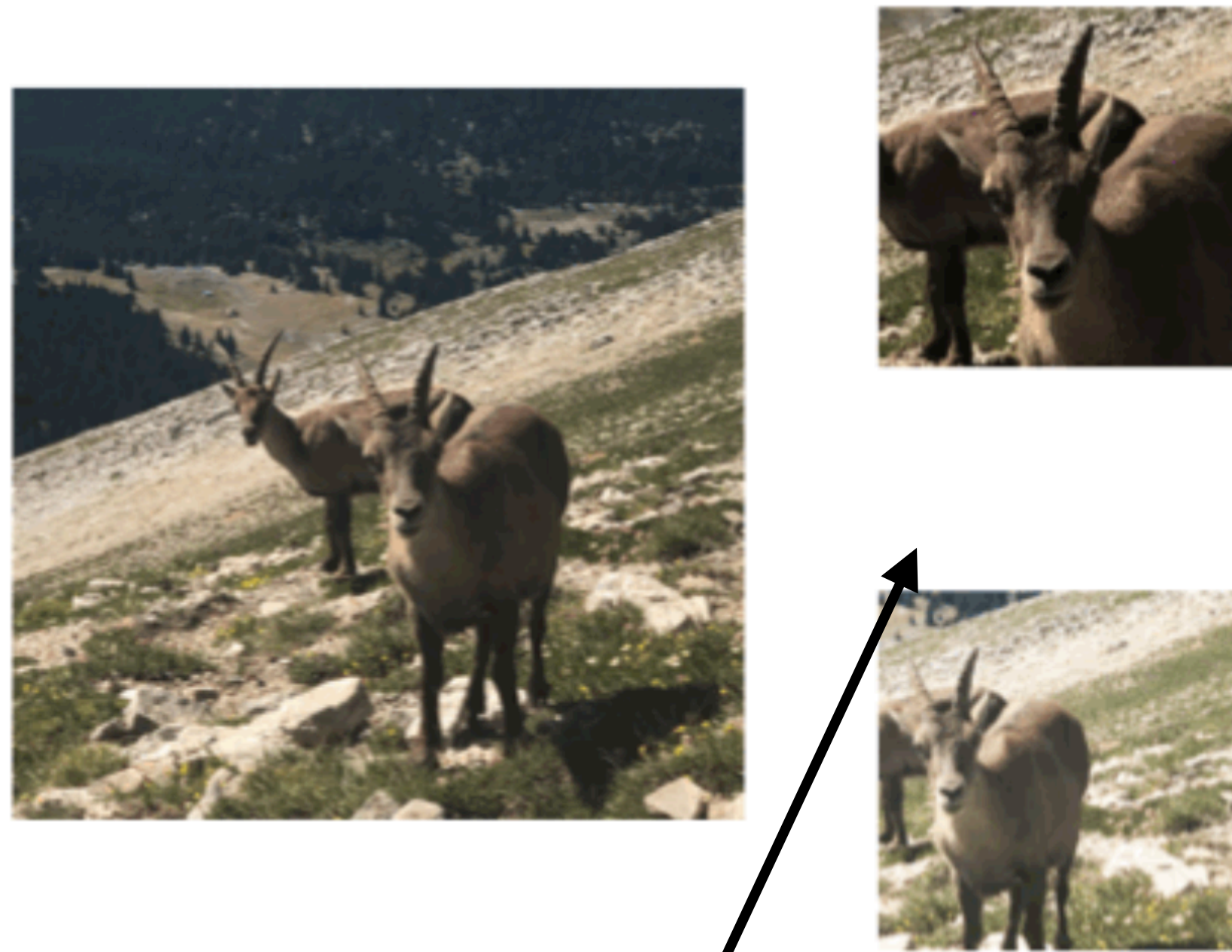
Identity
Preservation



SlotContrast
[Manasyan et al., 2024]

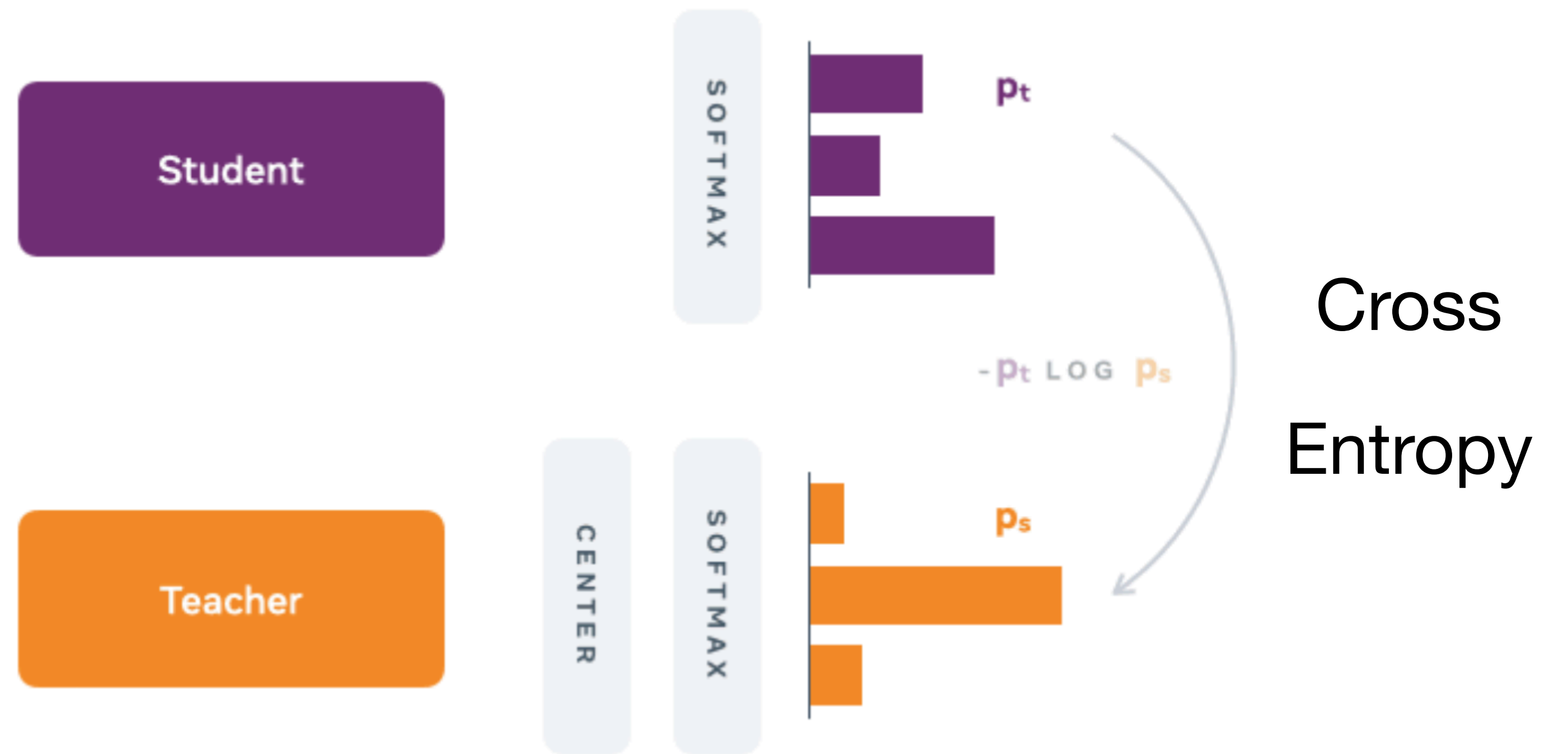
Self-supervised Semantic Features

Self-supervised Semantic Features

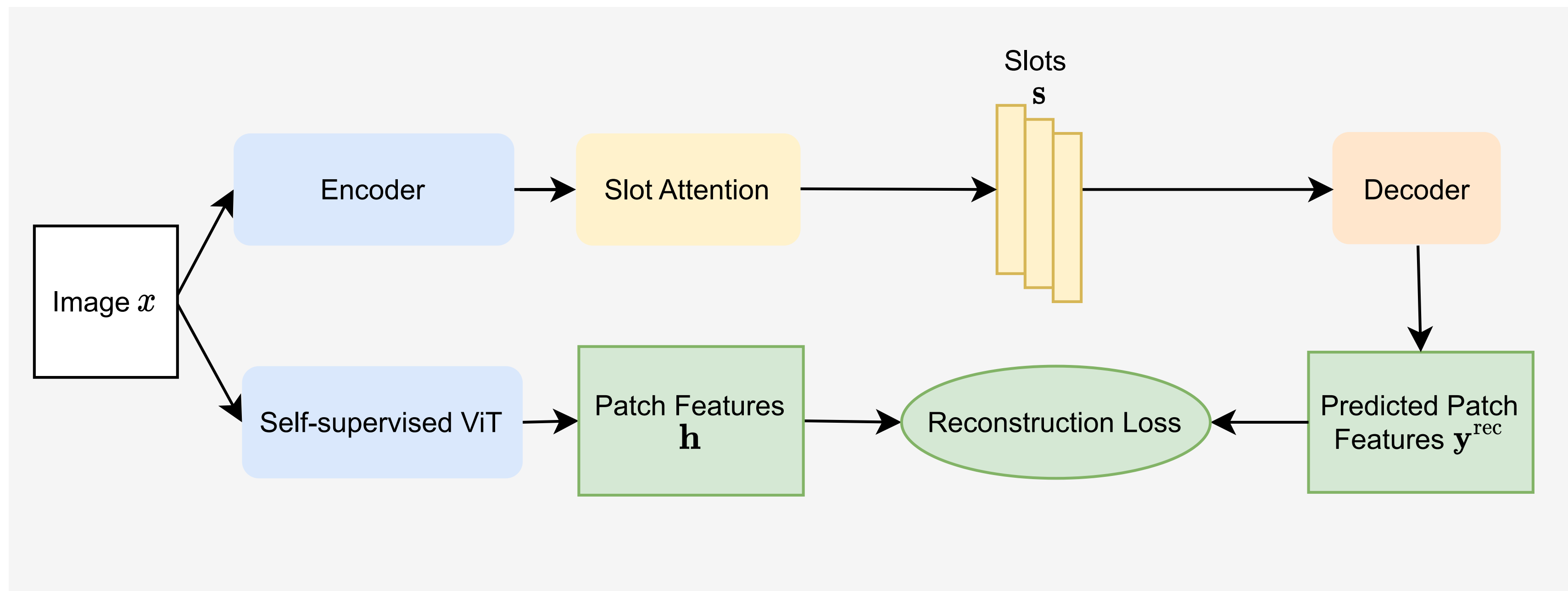


Multi-crop augmentations strategy:

many small crops for student & larger crop for teacher



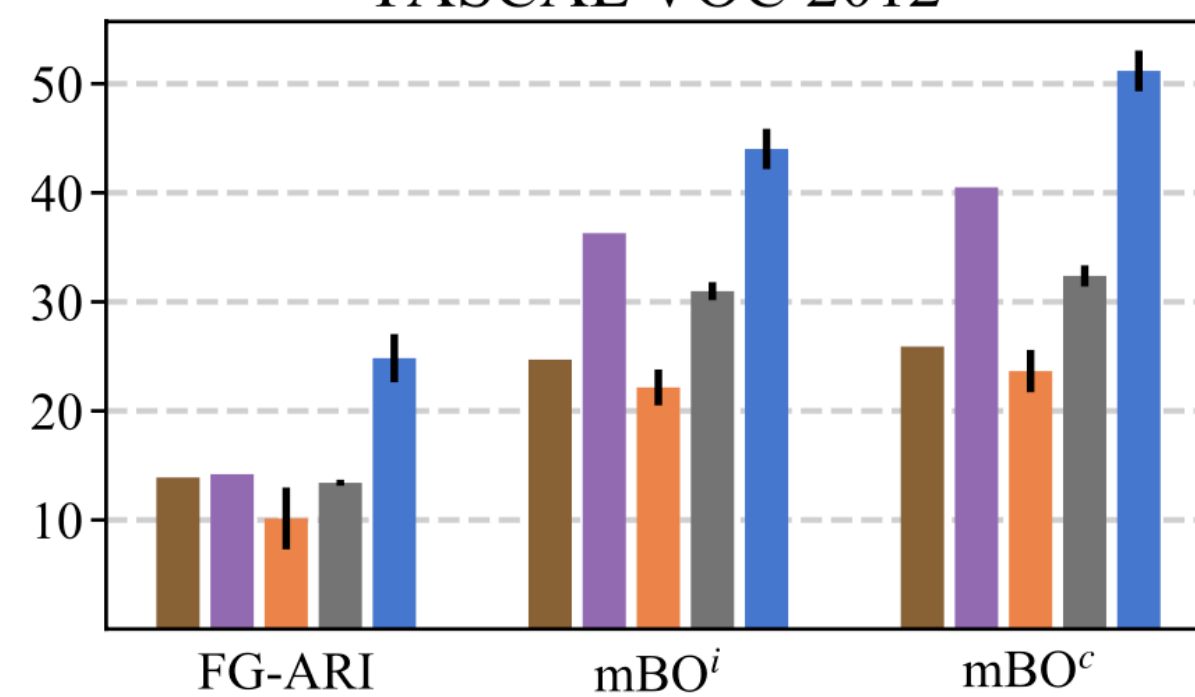
DINOSAUR: Self-supervised Features as Targets



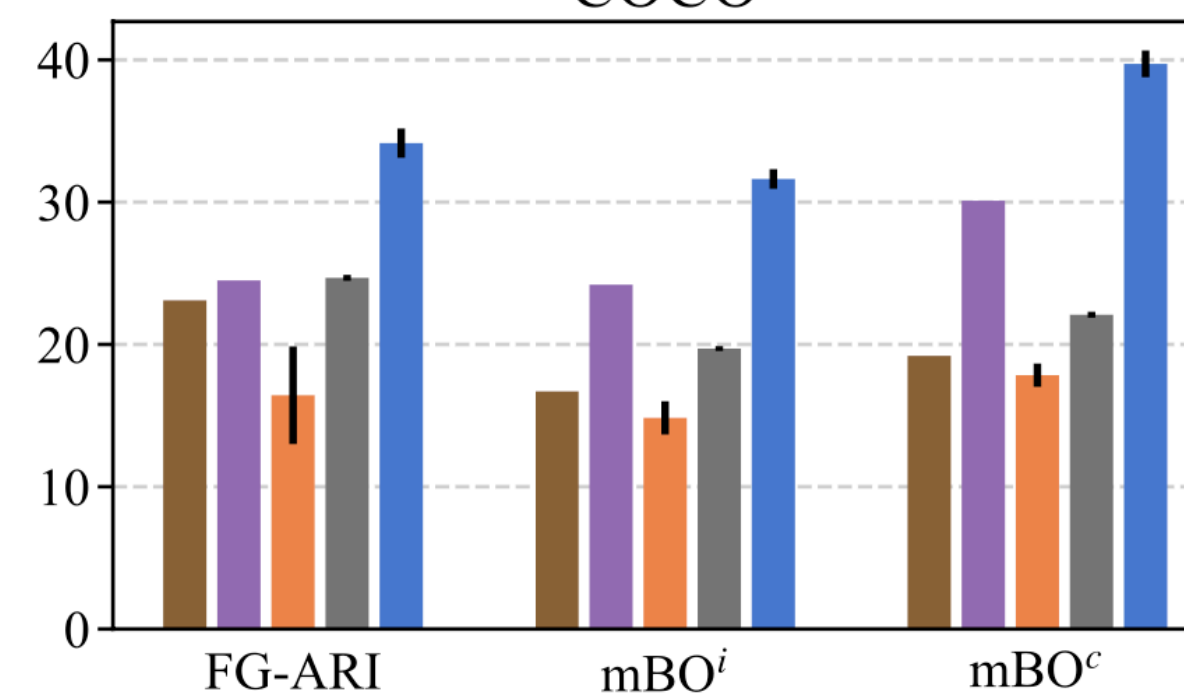
DINOSAUR Results



PASCAL VOC 2012

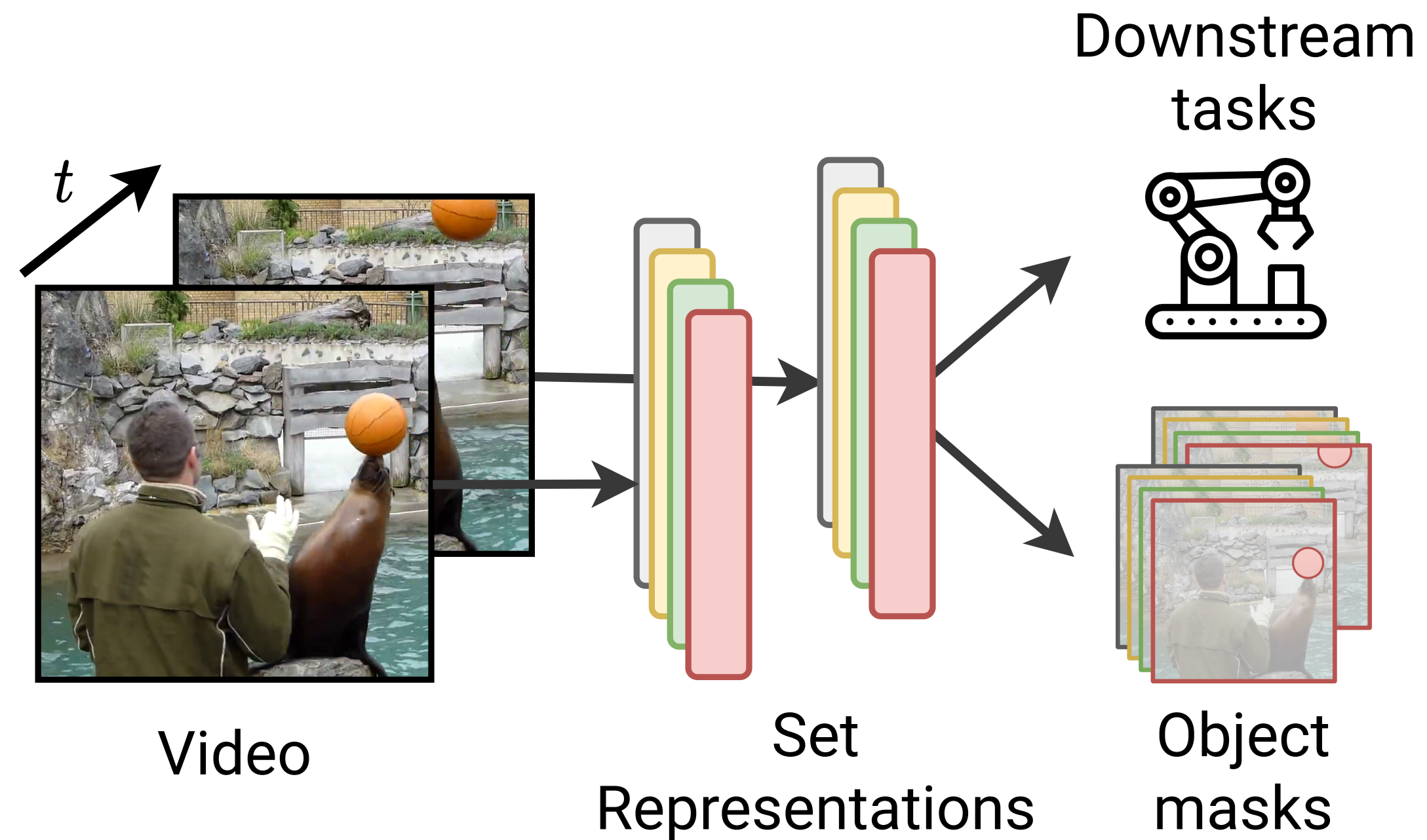


COCO

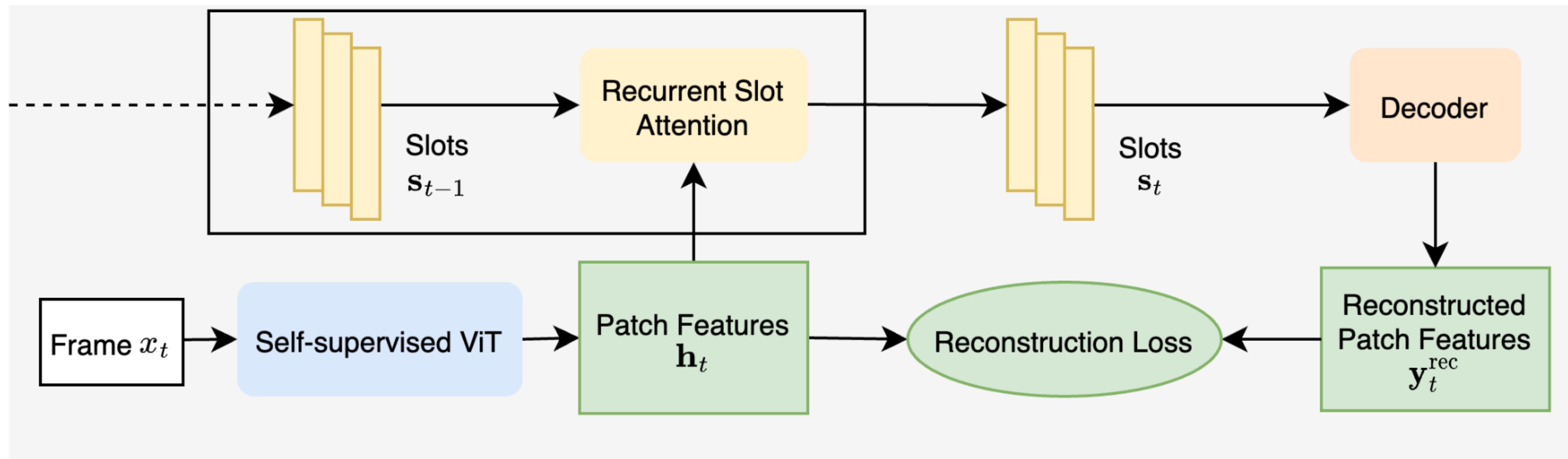


Block Masks
 DINO K-Means
 Slot Attention
 SLATE
 DINOSAUR

Can we extract even better targets from videos?



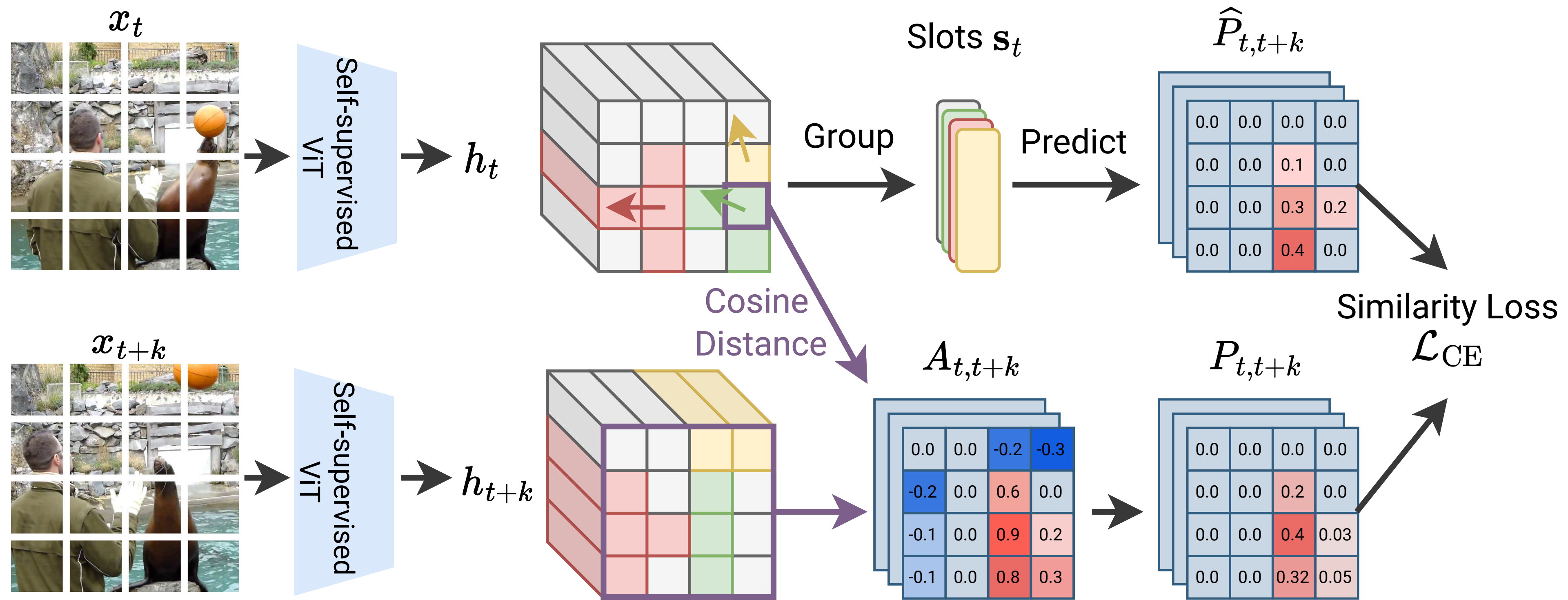
Recurrent Slot Attention for Slot Consistency



- DINO SAUR training objective: reconstruction of the current frame features
- Recurrent SA from SAVi⁹: connects slots from different frames via initialisation of SA iteration

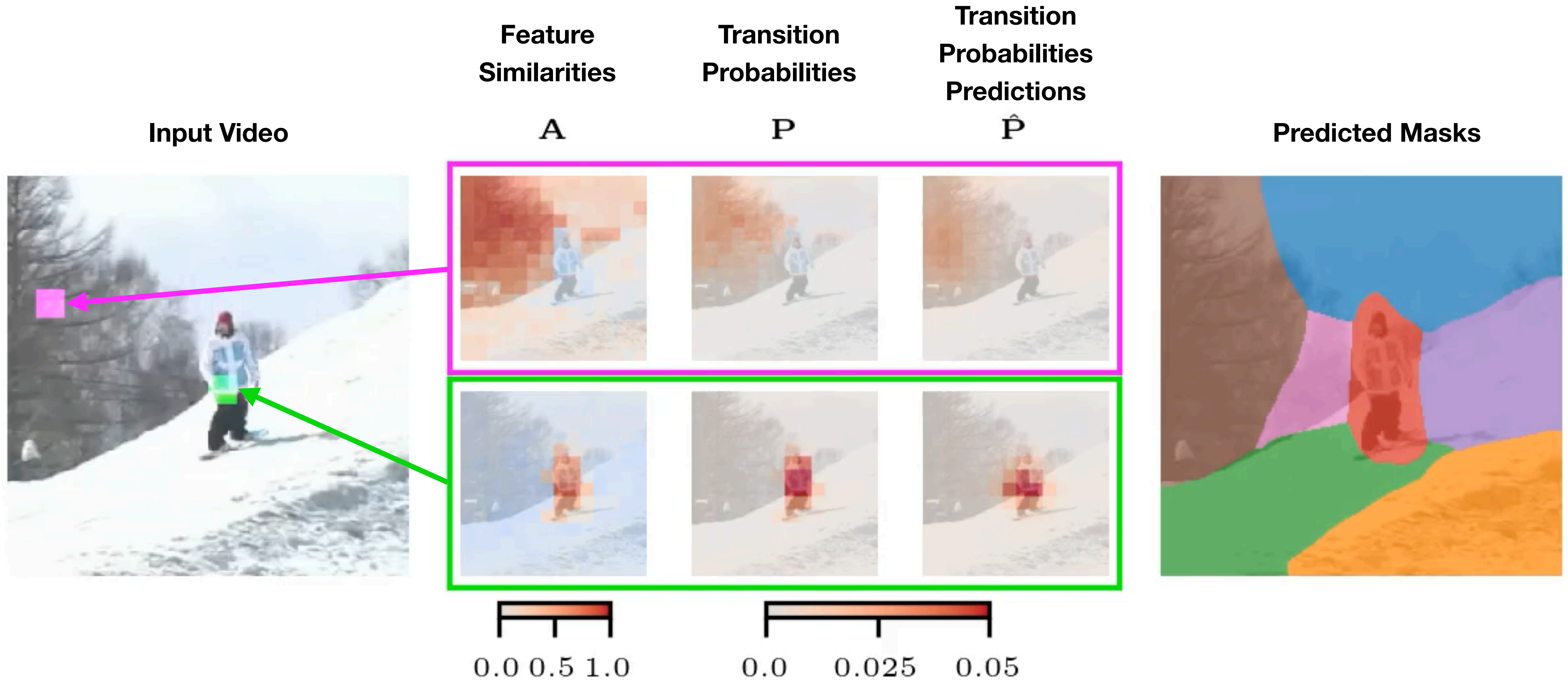
How can we facilitate video object discovery from the **temporal structure of the video**?

Temporal Features Similarity Prediction

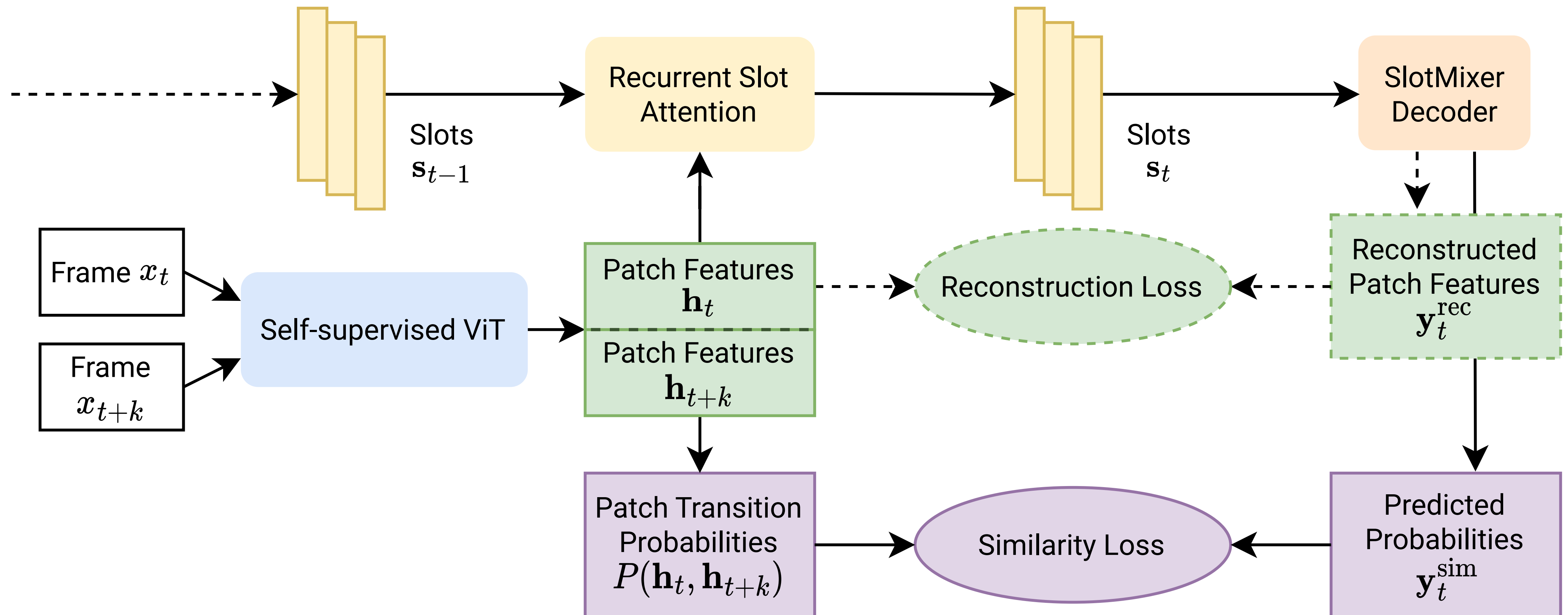


- Successful prediction of the temporal similarity requires **combining semantics and motion information**

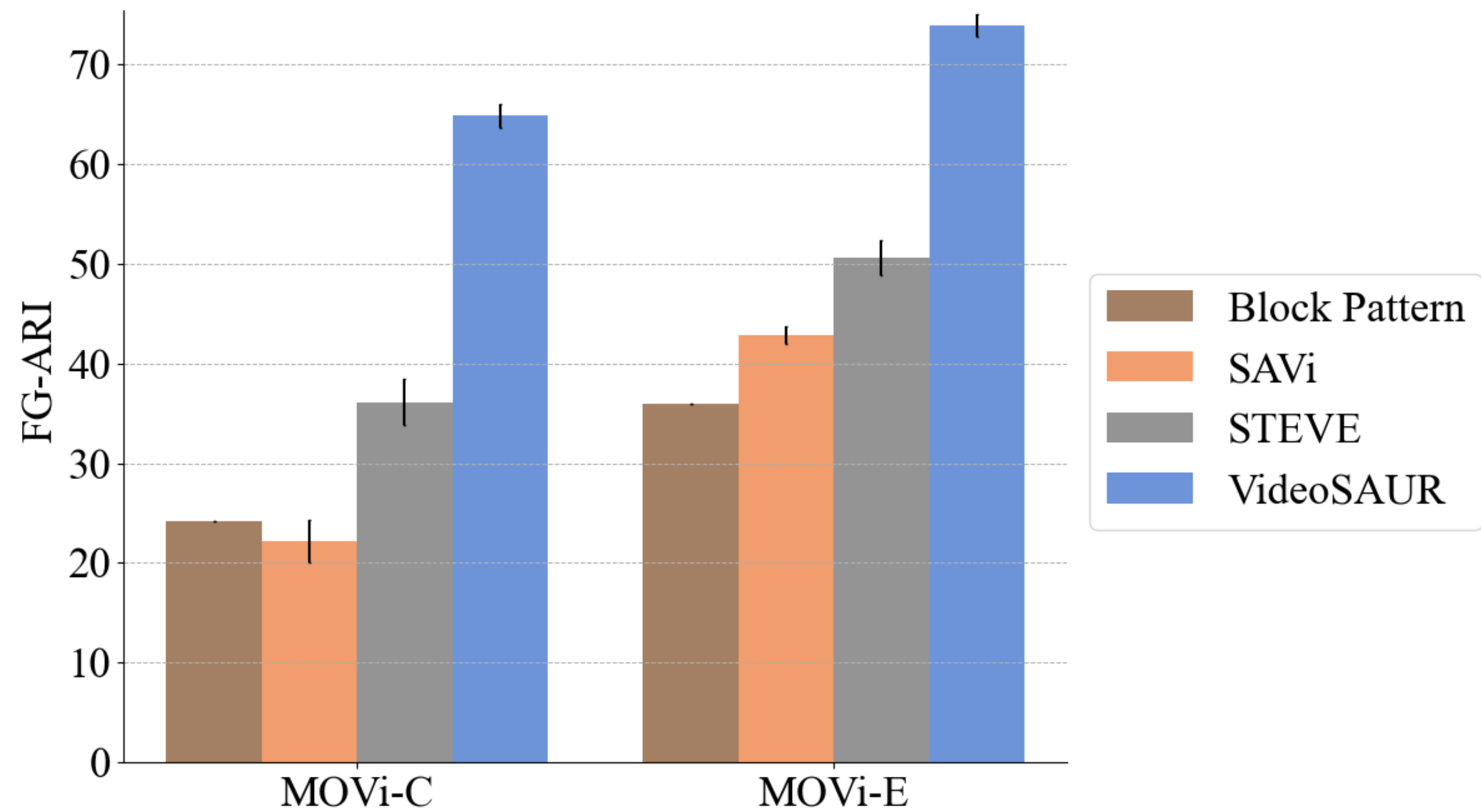
Temporal Features Similarity Prediction



VideoSAUR

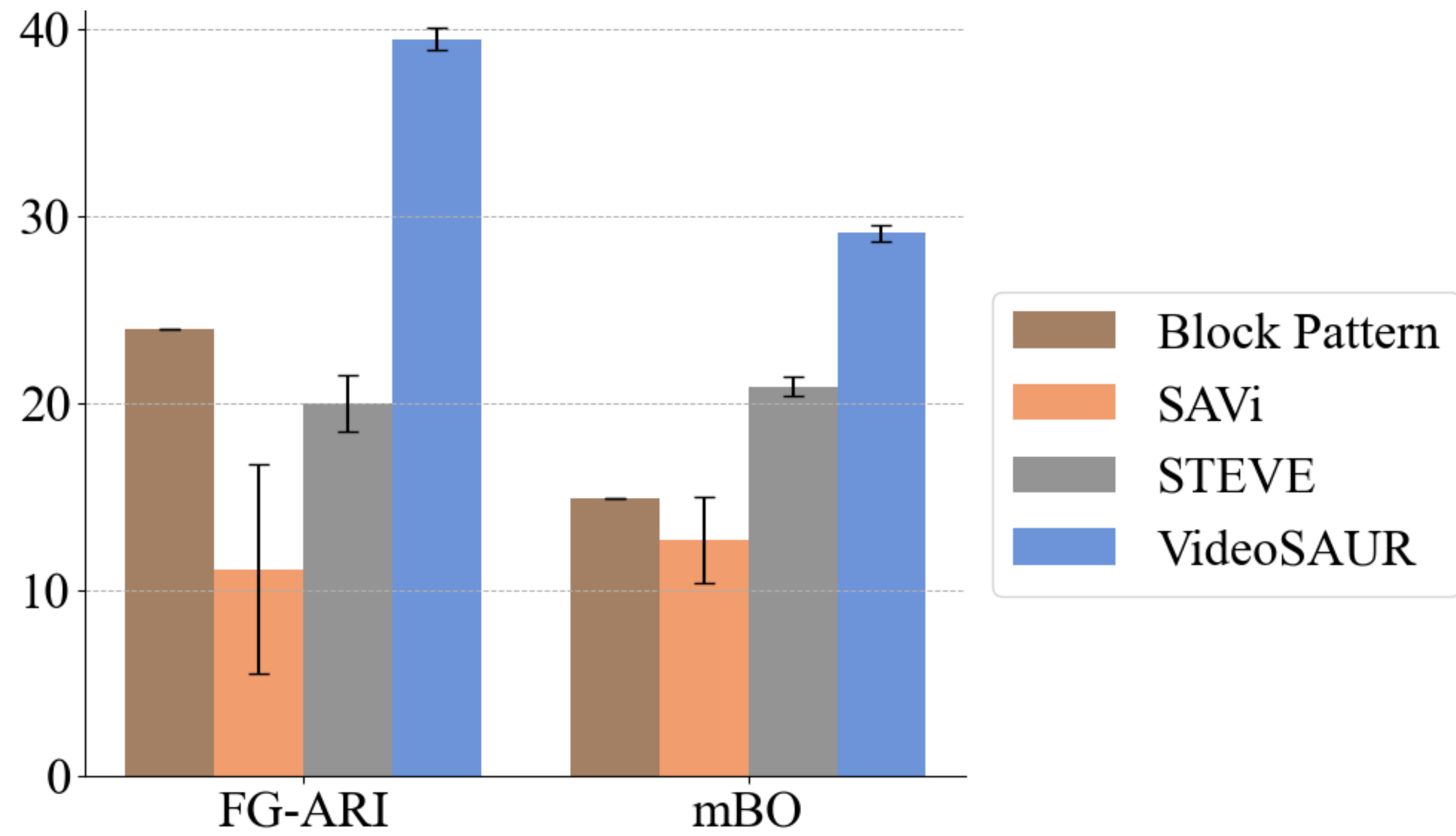
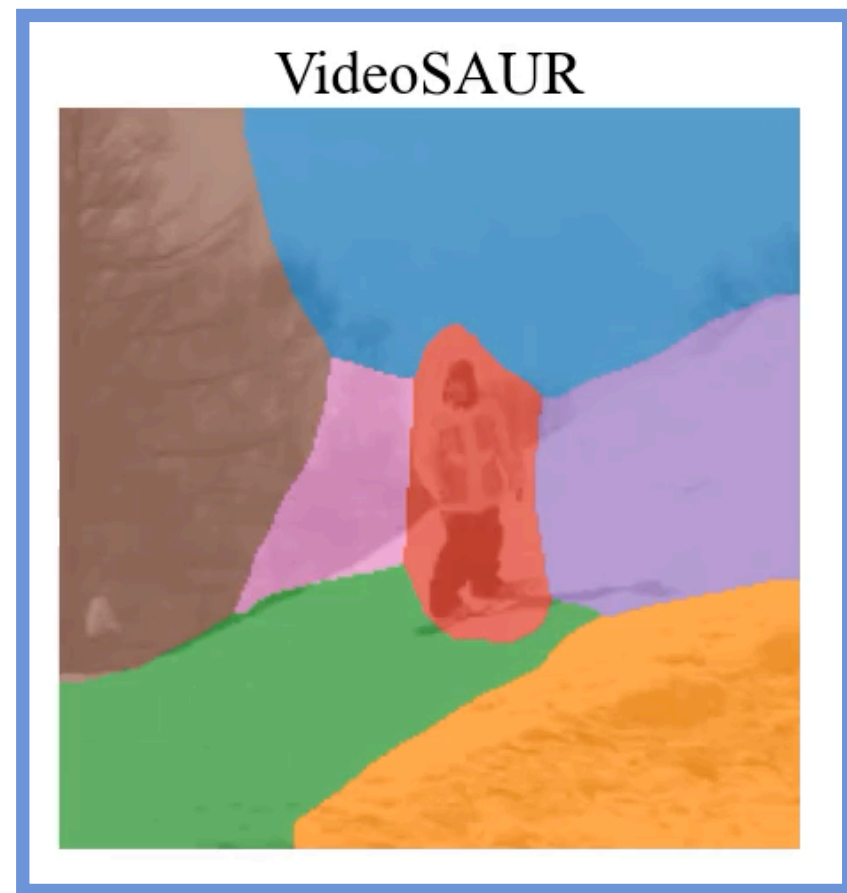
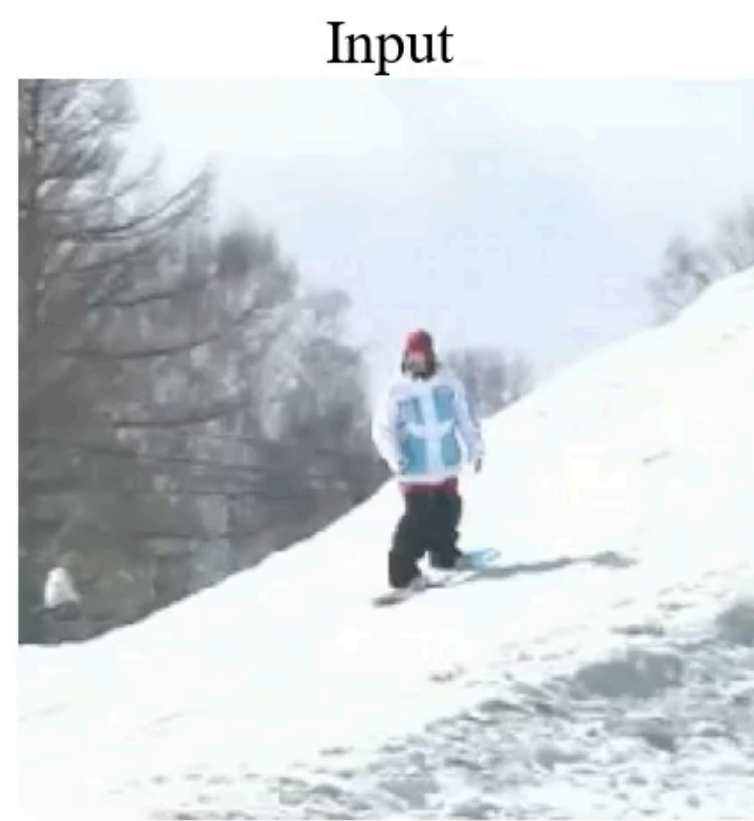


VideoSAUR Results on Synthetic Videos



VideoSAUR Results on Real-World Videos

YouTube-VIS



What about long-term consistency?

We need to maintain a consistent slot for an object throughout a video sequence.

No slot (ID) switches



Reassign original slot: Reappeared objects should get their original slot (**object permanence**).



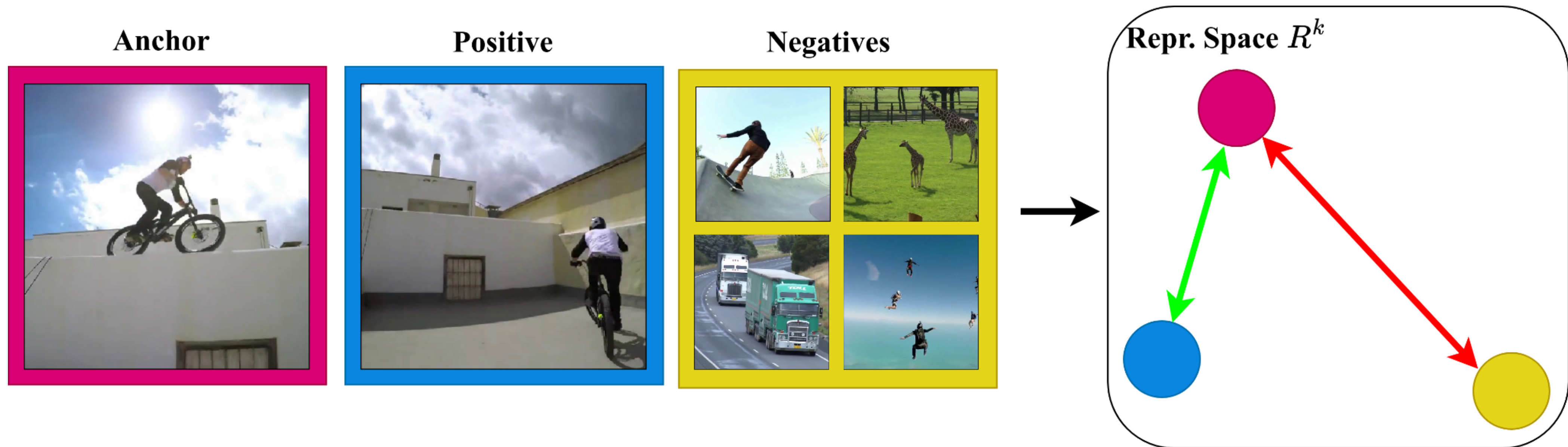
Assign new slots: Newly appearing objects should use unused slots.



Preserve slot assignments: Do not reuse a slot of a disappeared object.



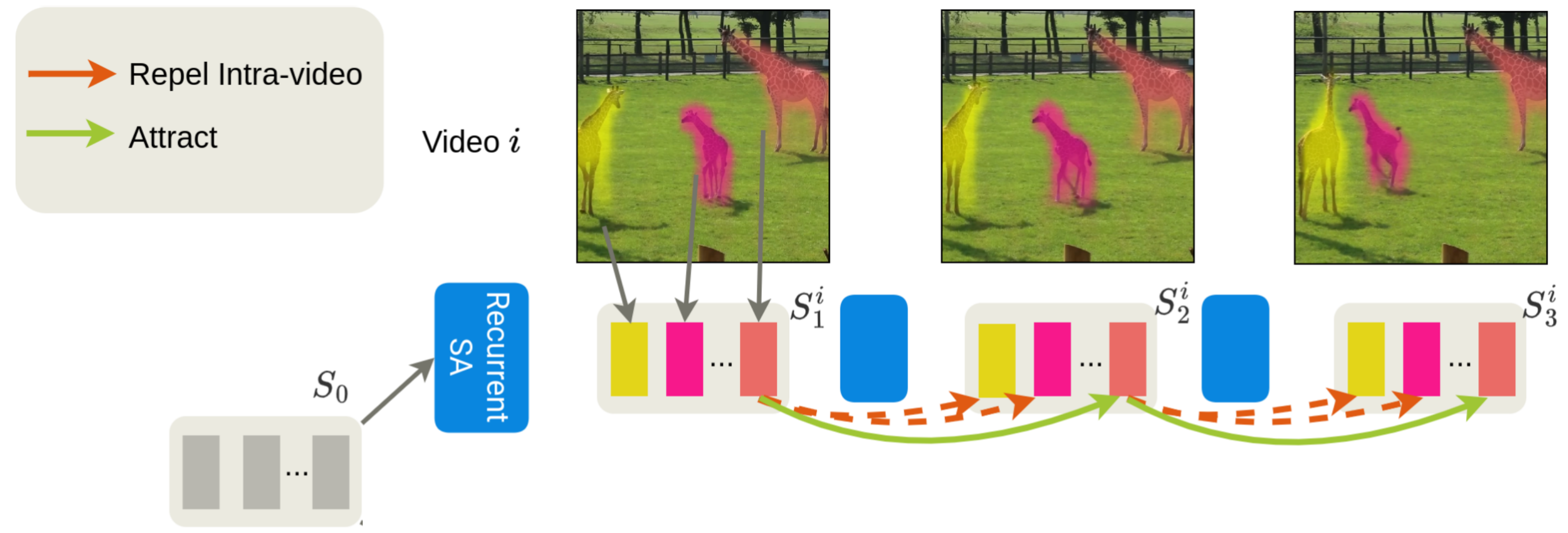
Video-level Contrastive Learning



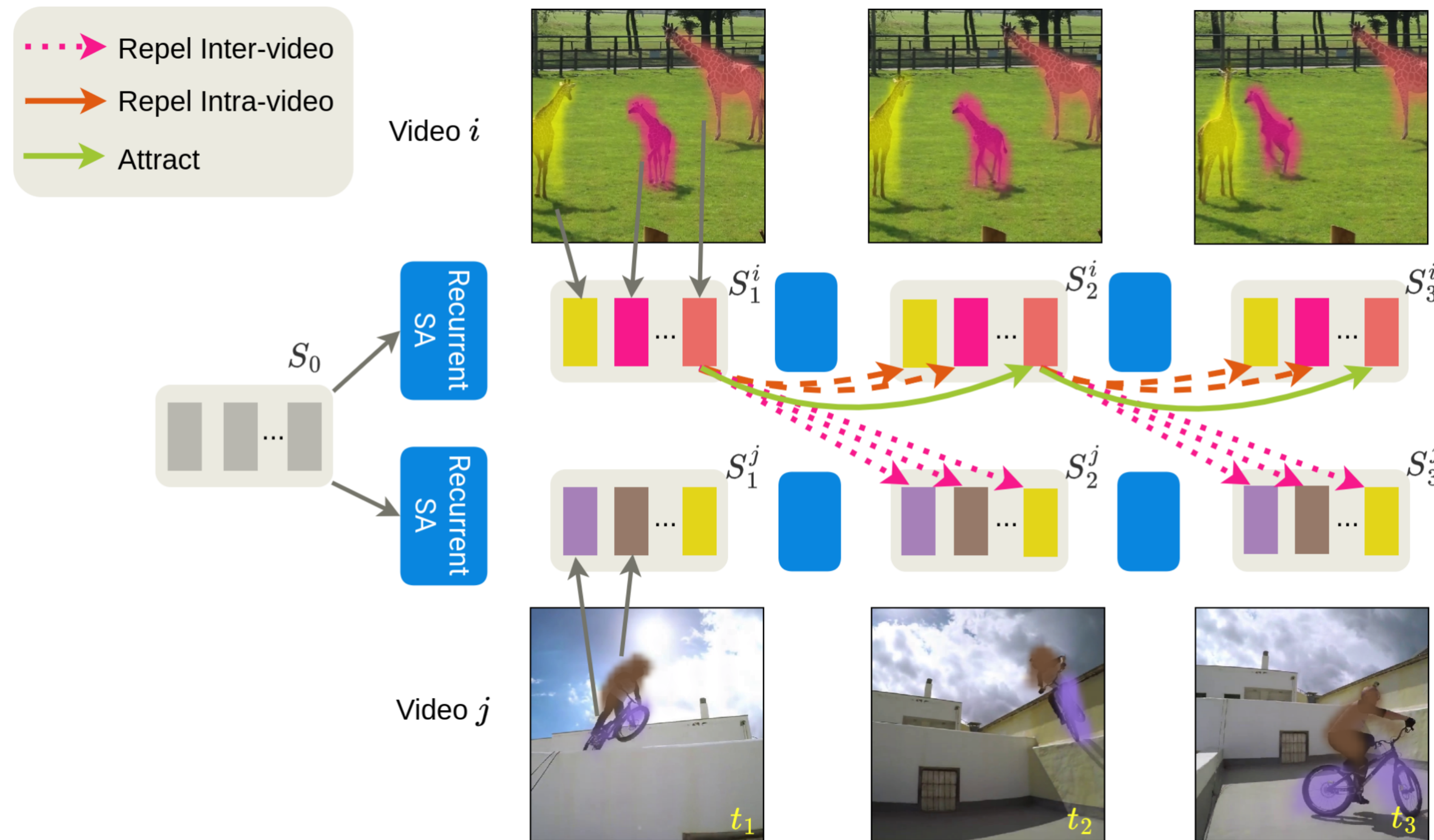
- Attract video frames of the same video
- Rebel frames from different videos in the dataset

How can we incorporate similar **contrastive objective** on more granular **slot representation level**?

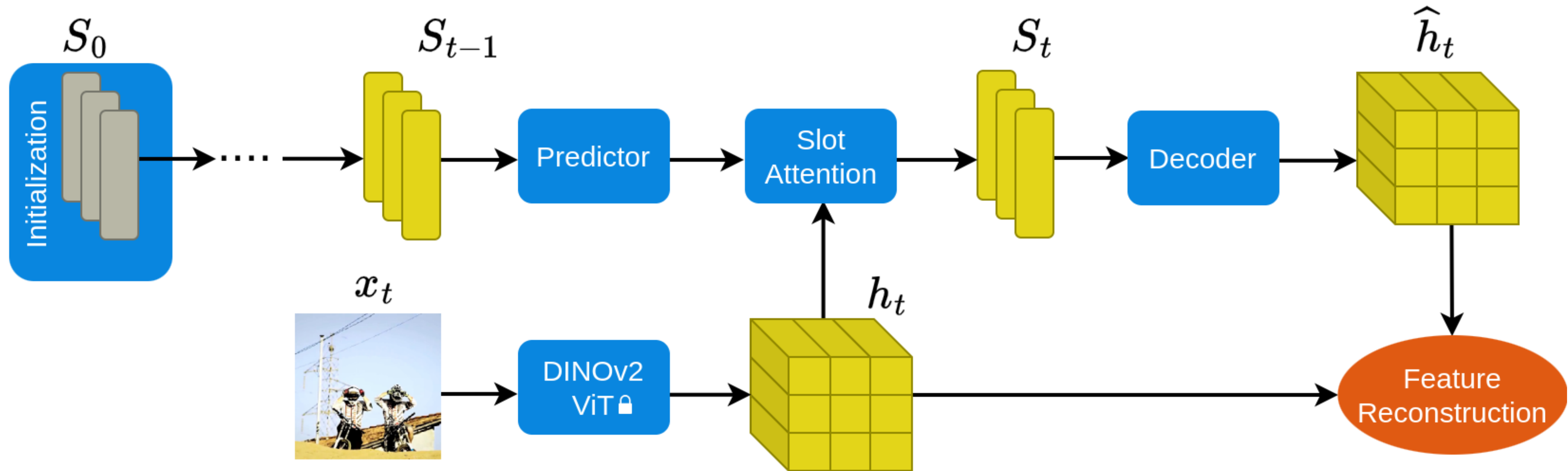
Slot-Slot Contrastive Loss (Intra-Video)



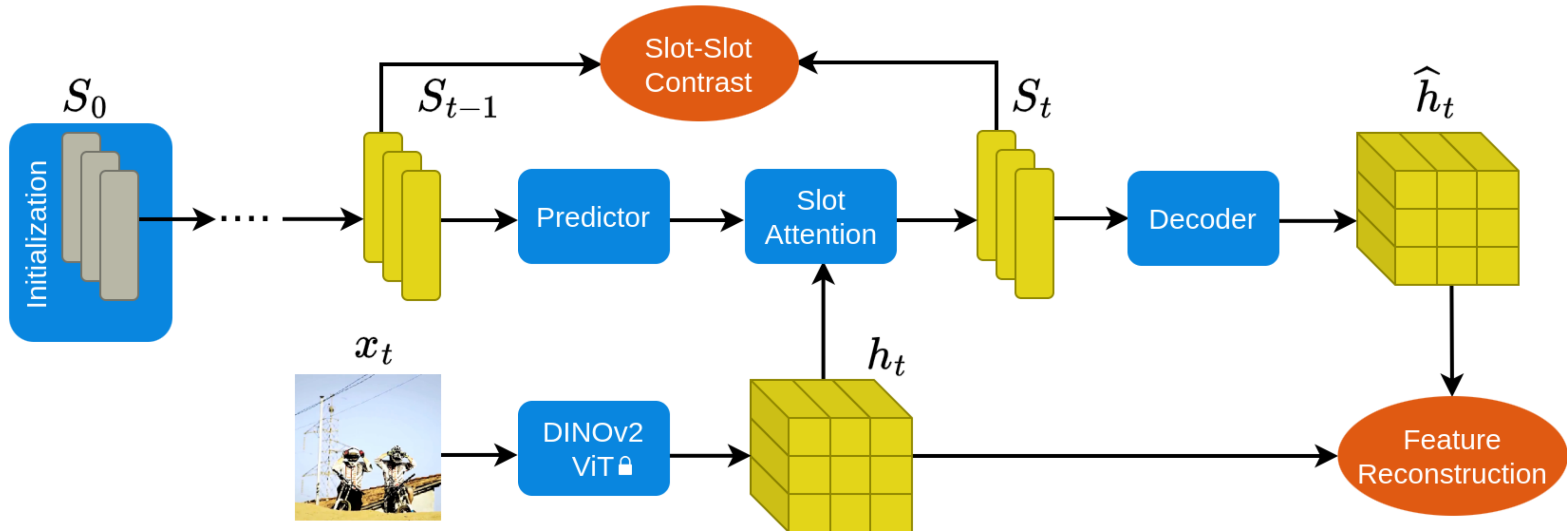
Batch-level Slot-Slot Contrast



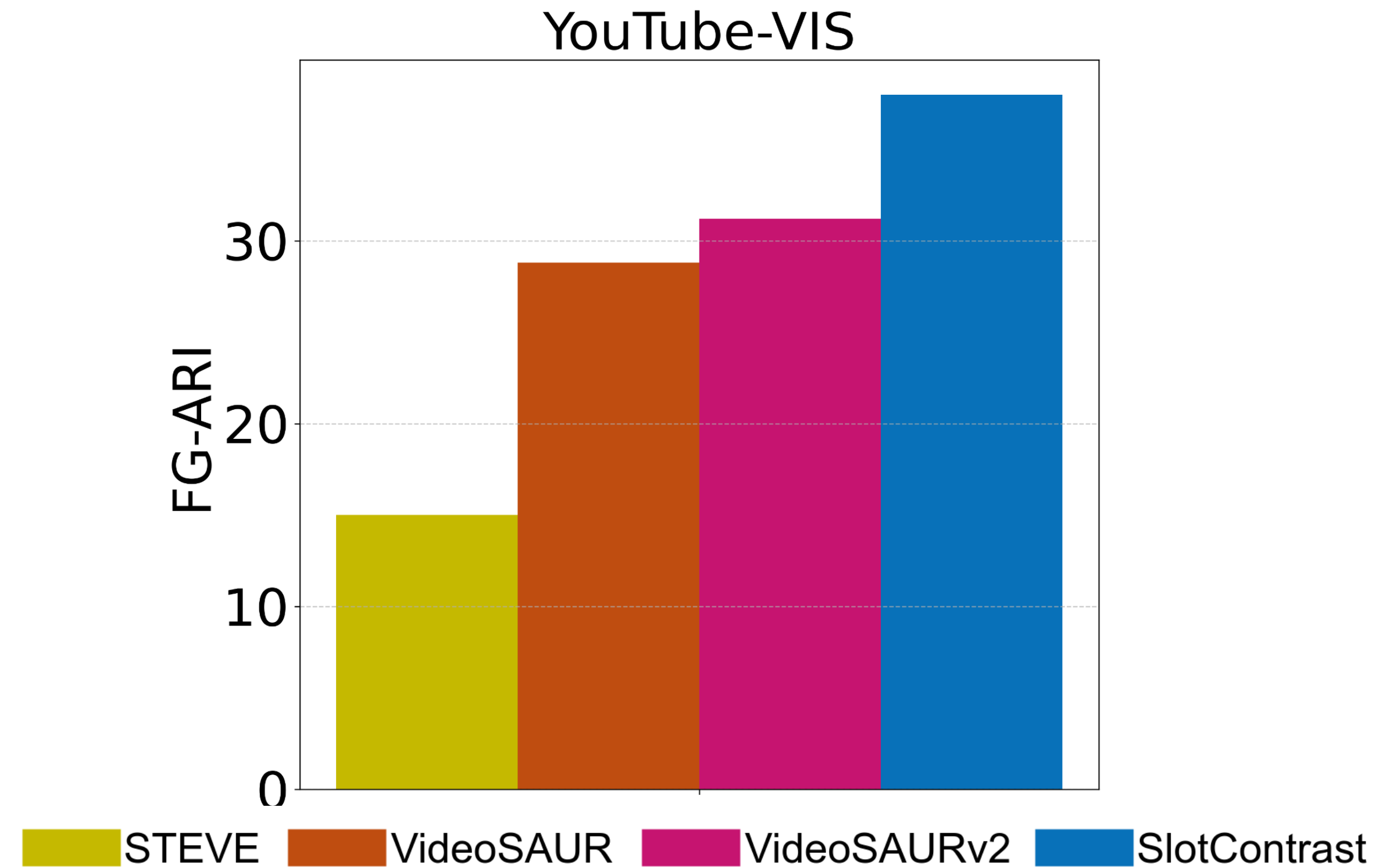
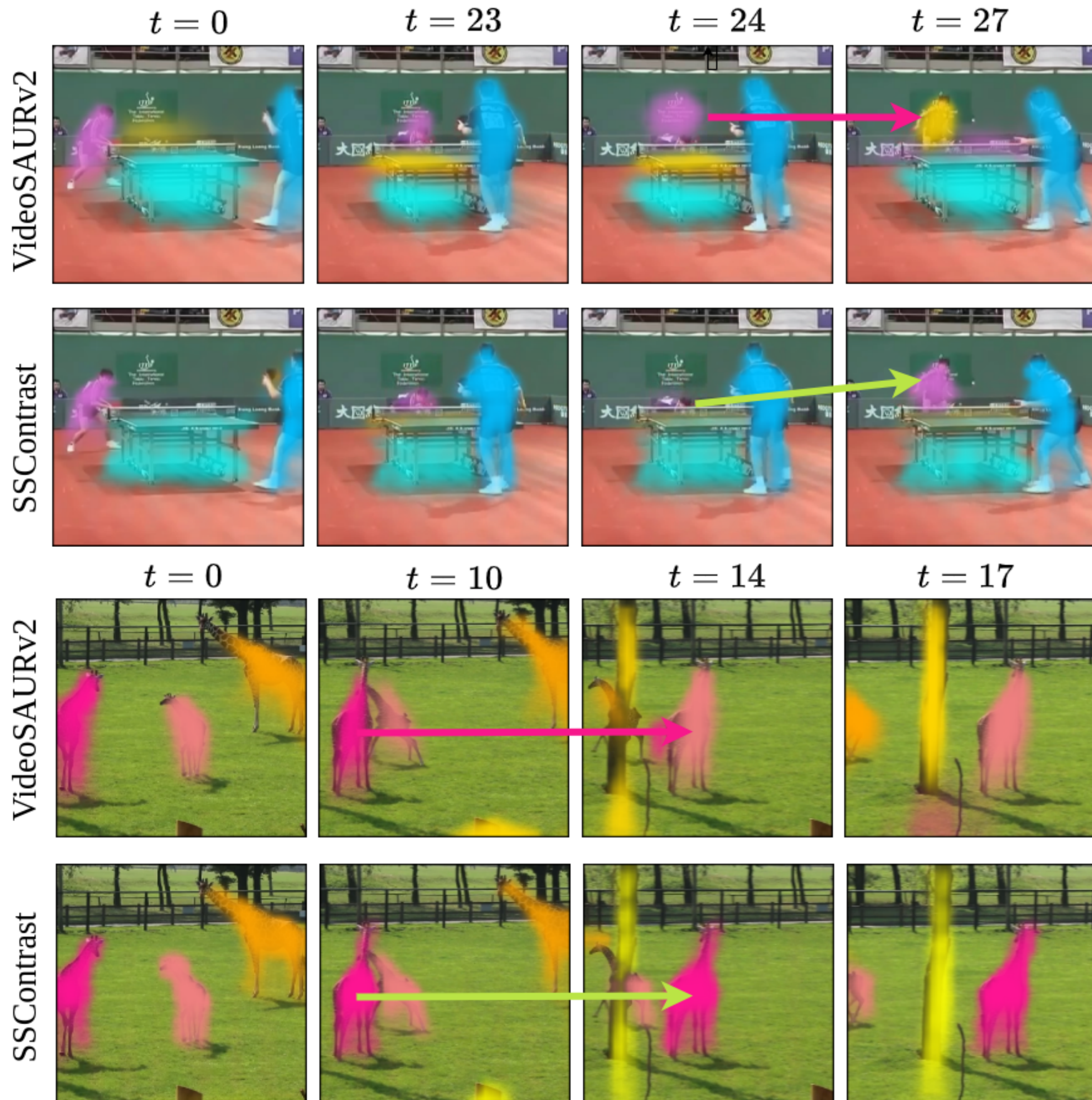
Video Object-Centric Learning Architecture



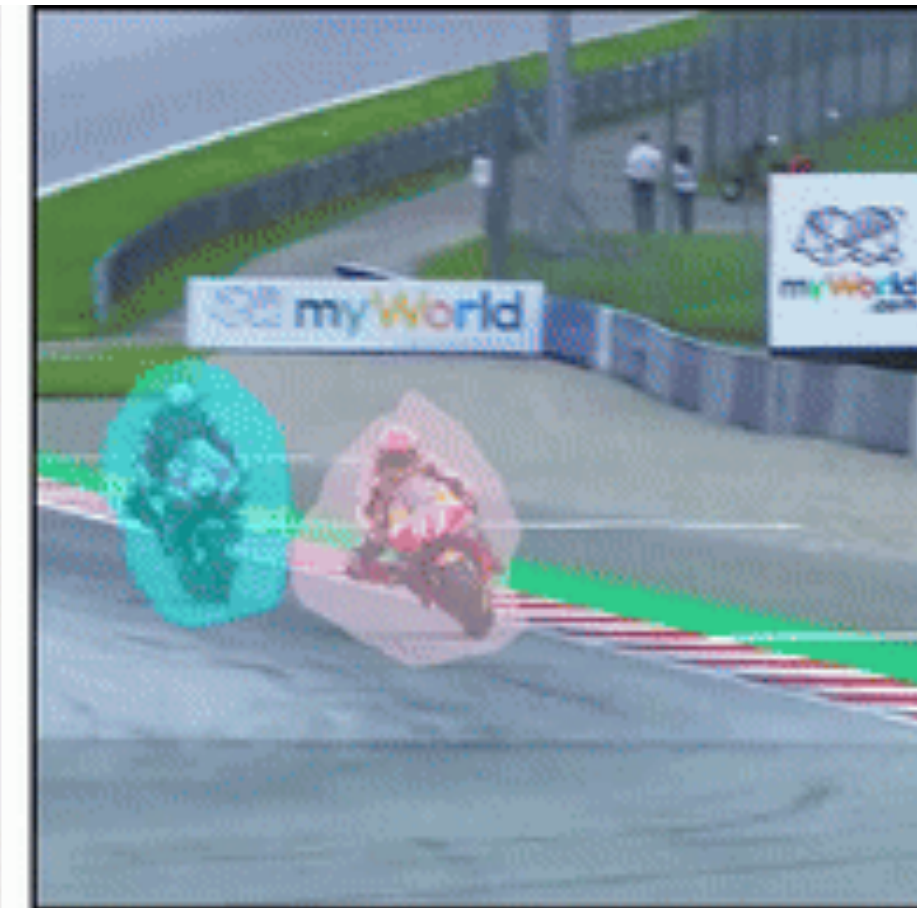
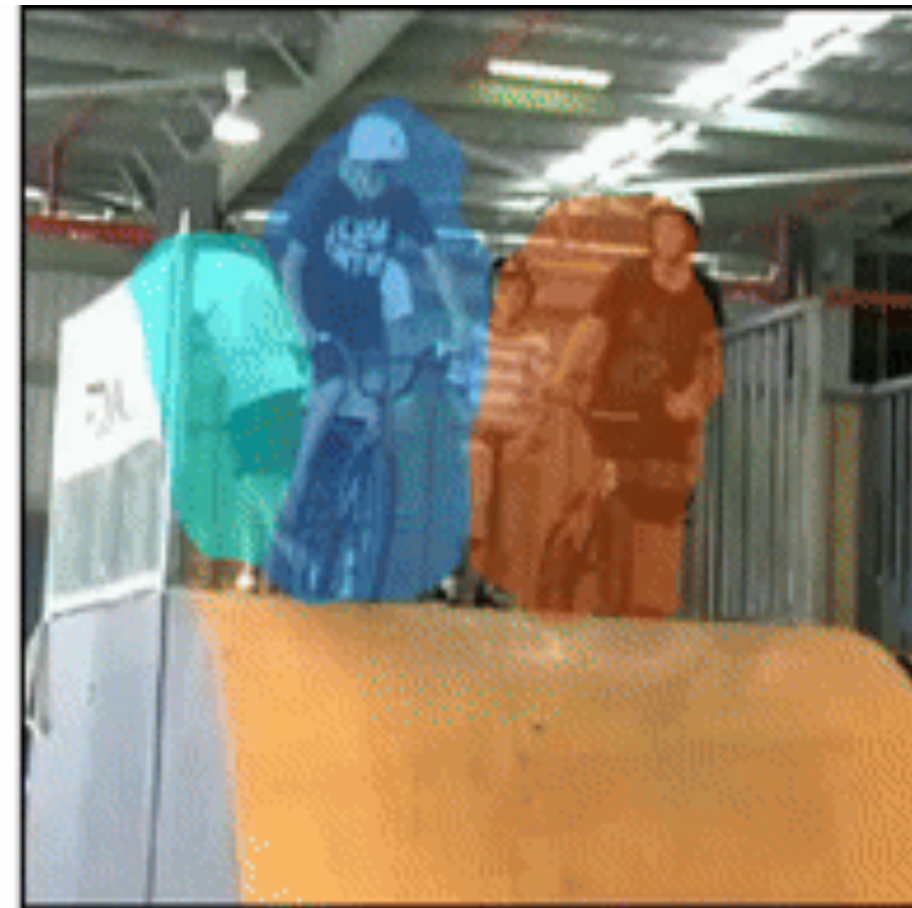
Slot Contrast Architecture



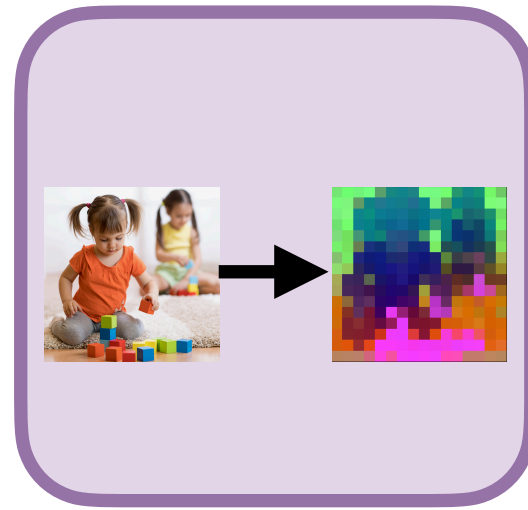
Object Discovery on Real-World Videos



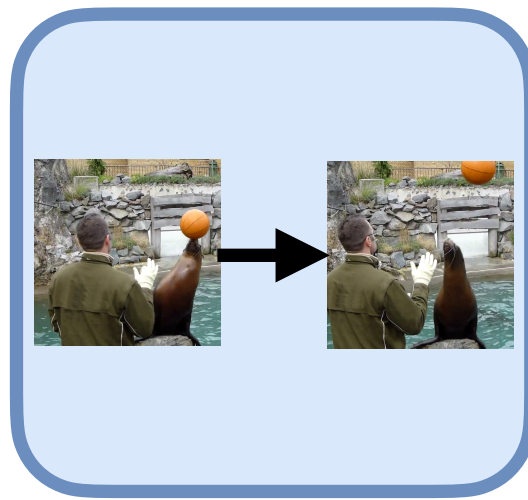
Object Discovery on Real-World Videos



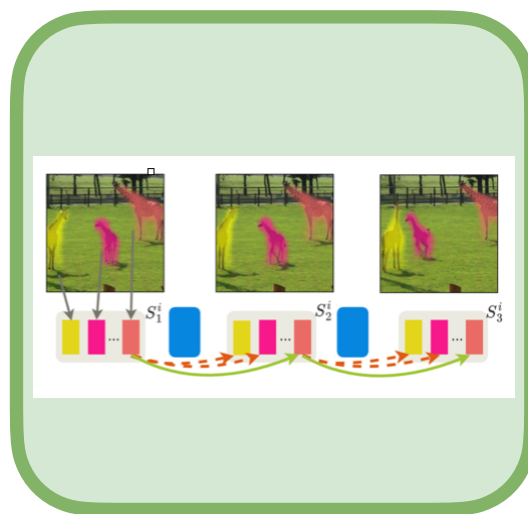
Summary



- **Semantics reconstruction objective** allows to scale object-centric representations to **real-world images**



- **Temporal similarity prediction** further scales object-centric representation to **real-world videos**



- **Slot Contrast loss** further improves **long-term consistency** of learned representations