

# STELLAR: Learning Sparse Visual Representations via Spatial-Semantic Factorization

**ICML 2026**

*Theodore Zhengde Zhao, Sid Kiblawi, Jianwei Yang, Naoto Usuyama,  
Reuben Tan, Noel C Codella, Tristan Naumann, Hoifung Poon, Mu Wei*



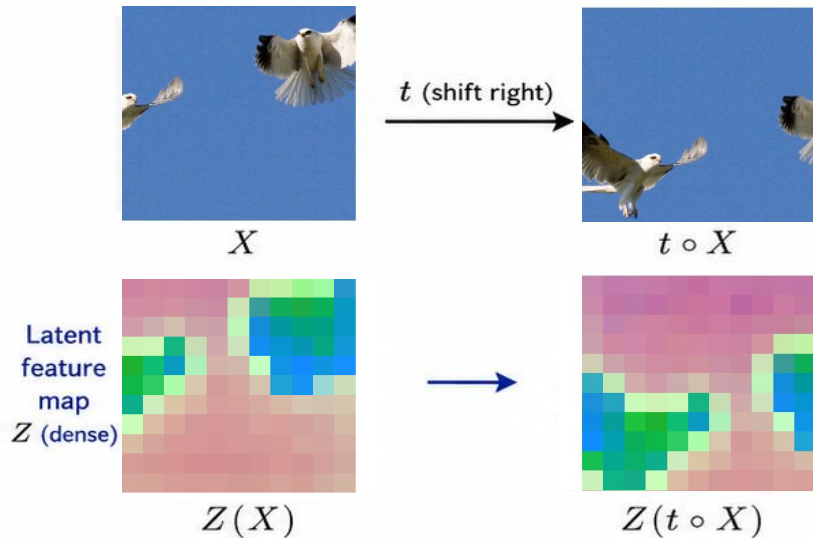
# Research Question

- How to learn “good” visual representation from image data
  - Rich abstract semantics for understanding tasks
  - Contains sufficient information to recover the reality (reconstruction)
- AE, VAE, MAE, DINO, Diffusion, SimCLR, BYOL...
  - Reconstruction-centric SSL: preserve details, but often weaker semantics
  - Semantic-centric SSL: strong semantics from invariance (joint-embedding), abandoning pixel reconstruction
- Take a step back:
  - 1. What’s a “good” visual representation...
  - 2. ... and how to learn it from image data

# A Thought Experiment...

How should the latent representation change under transformation?

## Reconstruction wants equivariance



For faithful reconstruction, the latent should track spatial changes.

$$D(Z(t \circ X)) \approx t \circ X$$

latent changes with transformation

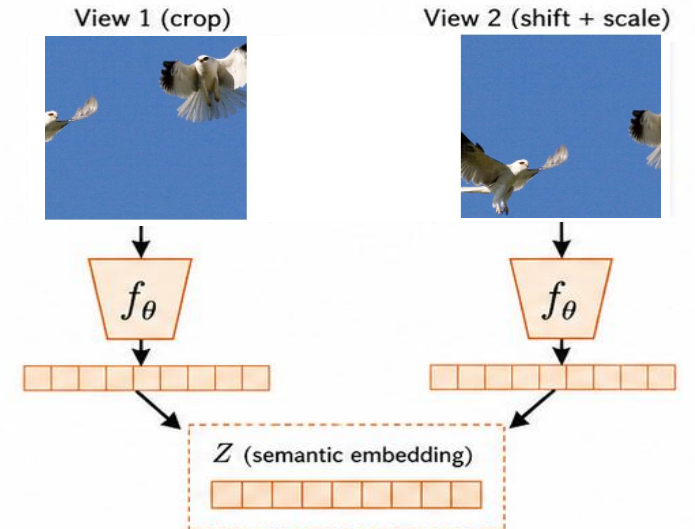
## Invariance–Equivariance Paradox



A single dense latent cannot be both fully invariant and spatially equivariant.

Dense representation entangles semantics and geometry.

## Semantic SSL wants invariance

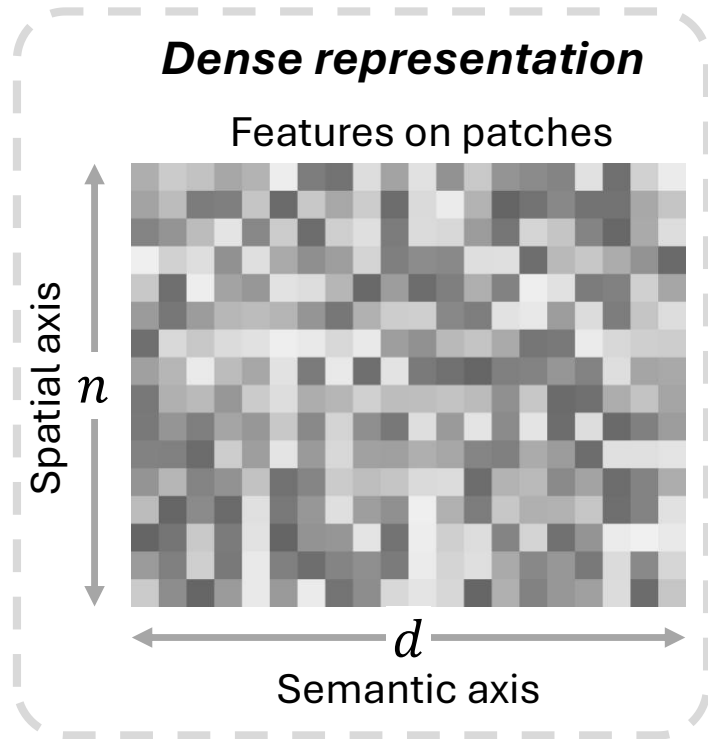


To learn high-level semantics, alignment pushes features to ignore spatial variation.

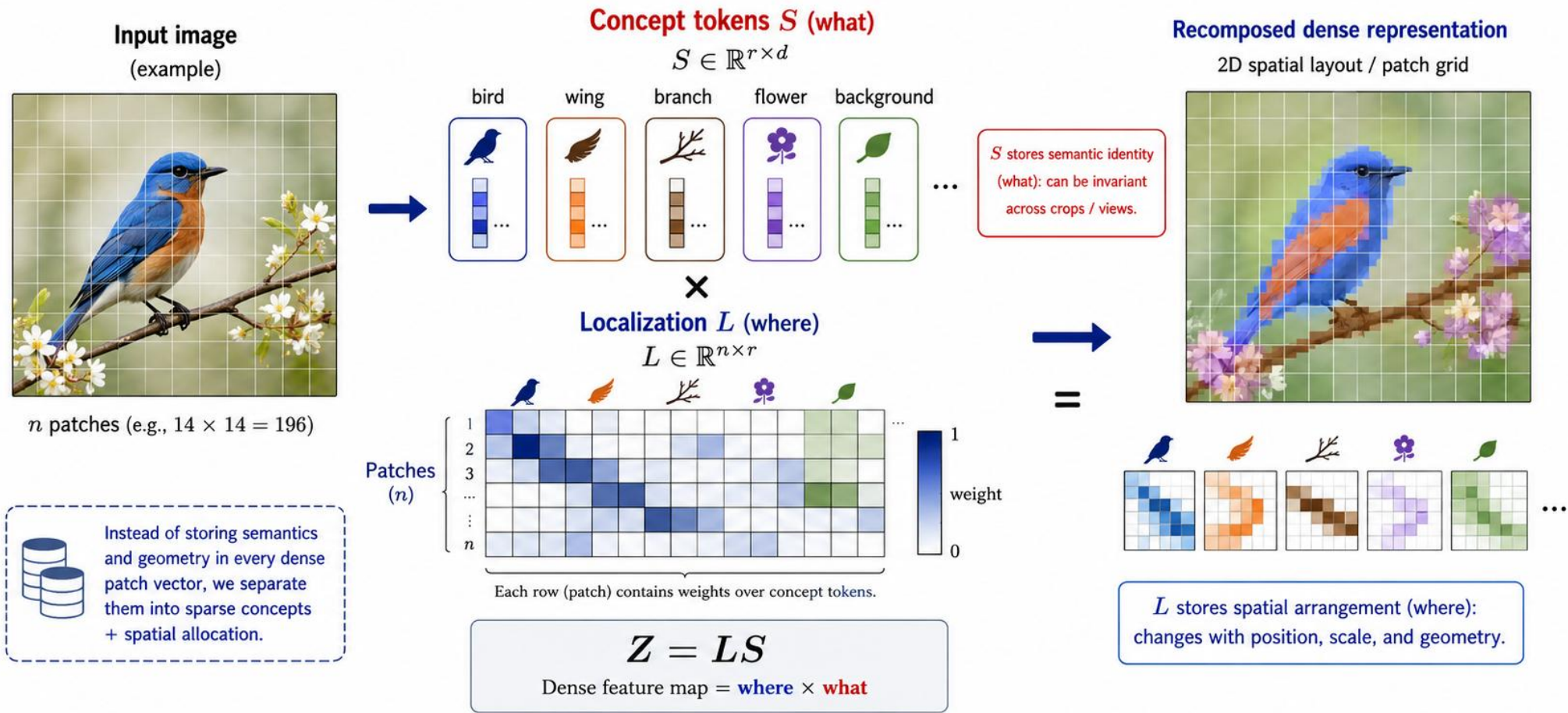
$$Z(t \circ X) \approx Z(X) \quad \text{or} \quad \left\| \frac{\partial Z(t \circ X)}{\partial t} \right\| \approx 0$$

latent stays the same across transformations

# Spatial-Semantics Factorization



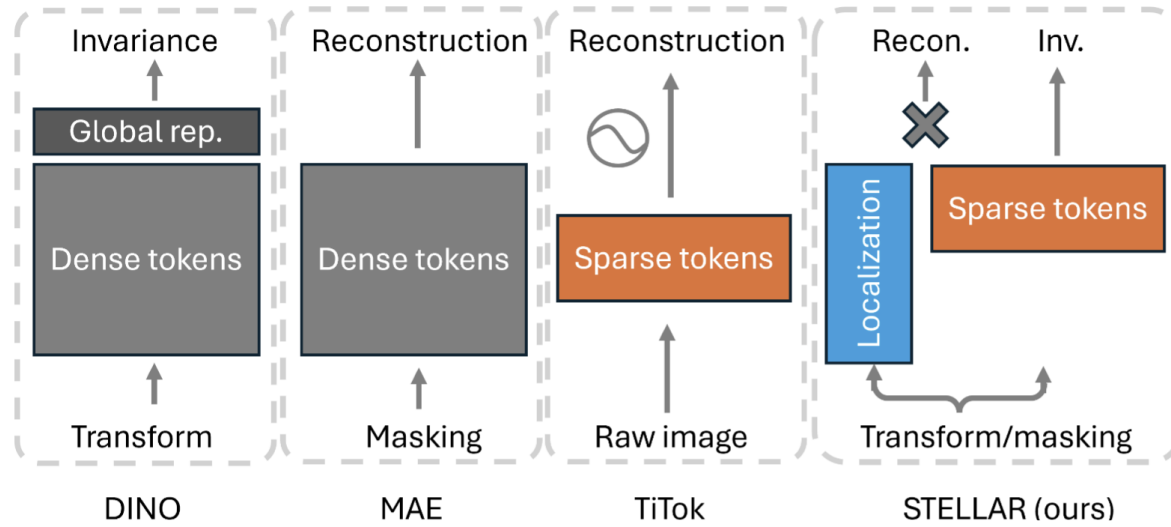
# Modeling Image with Sparse Concepts



## Takeaway:

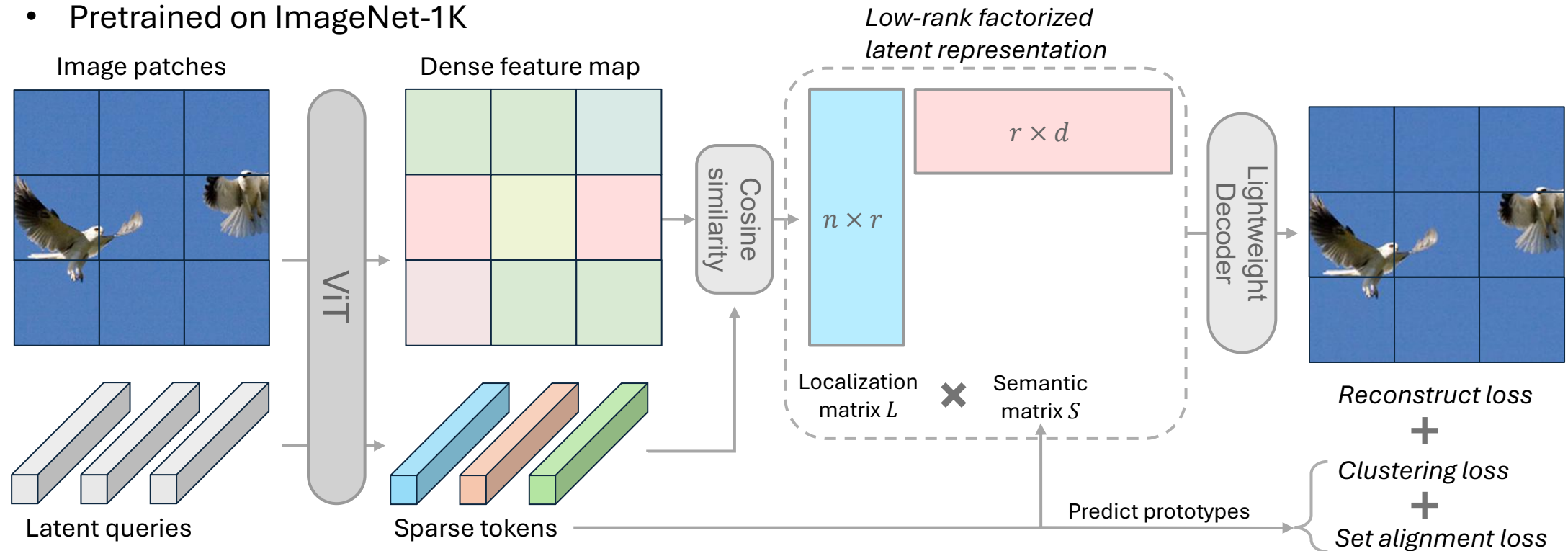
An image can be described by a few concept tokens, then placed back onto the 2D grid.

# Jointly Self-supervised Learning



# Simple Architecture

- Encoder: standard ViT + learnable queries
- Spatial localization from cosine similarity
- Decoder: 6-layer ViT predicting VQ-GAN code
- Pretrained on ImageNet-1K

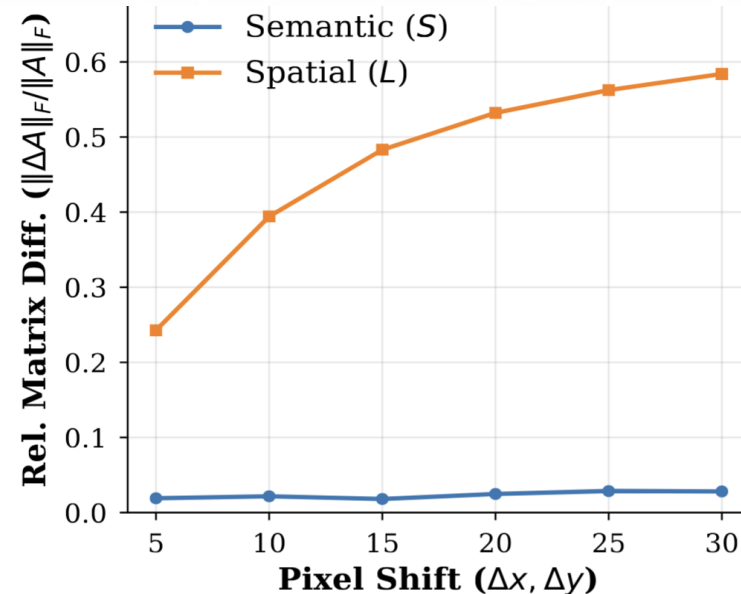
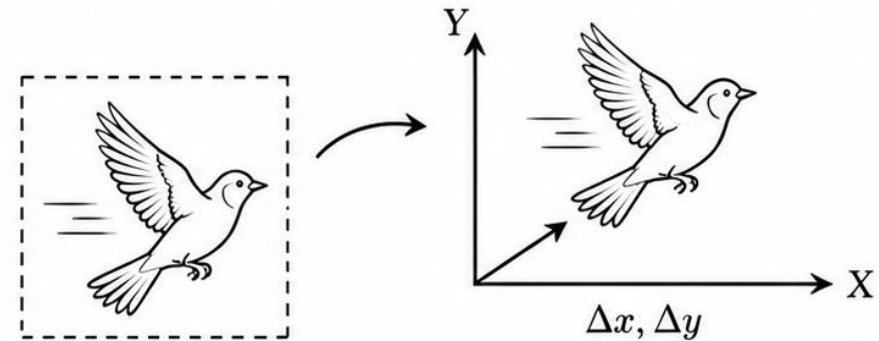


# Exp: Equivariance Partition

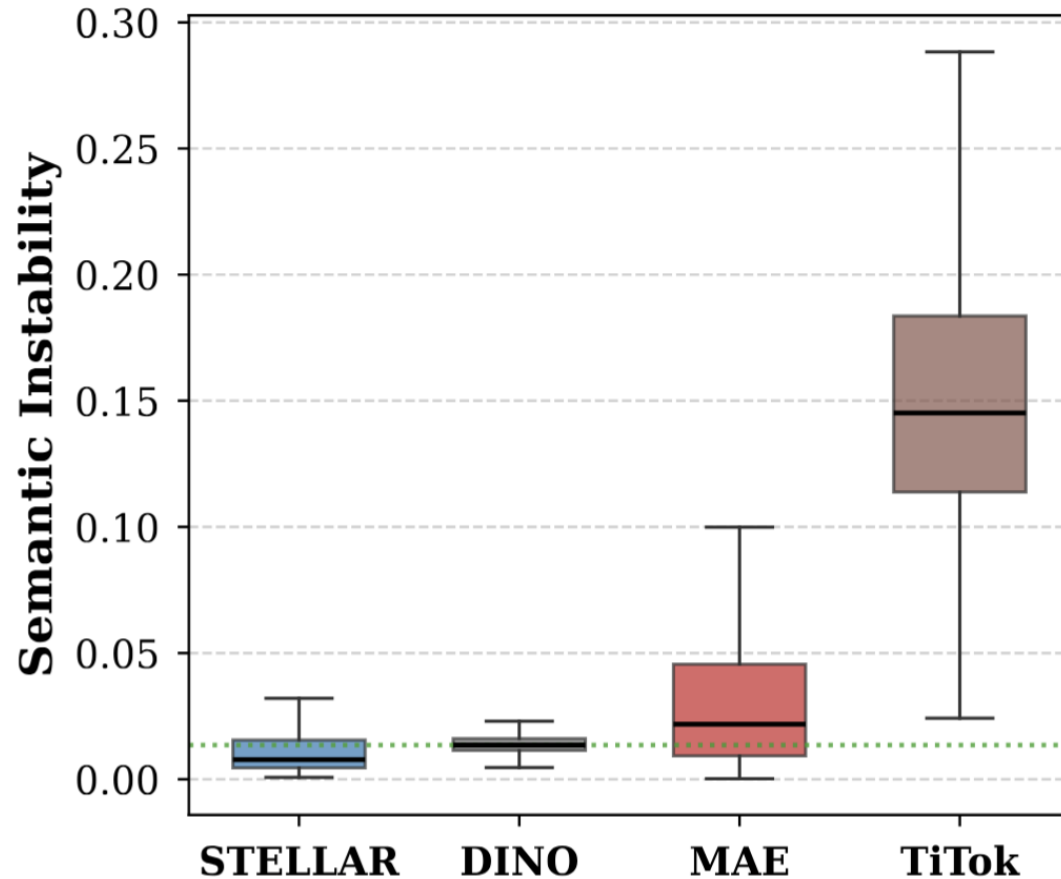
- Measure representation change under spatial transformation

$$\frac{\partial Z}{\partial \theta} = \underbrace{\left( \frac{\partial L}{\partial \theta} \right)}_{\text{Spatial Equivariance}} S + L \underbrace{\left( \frac{\partial S}{\partial \theta} \right)}_{\text{Semantic Variance} \approx 0}$$

By disentangling spatial and semantic information, we offload spatial equivariance entirely to the localization matrix (**L**), while the semantic tokens (**S**) remain invariant.



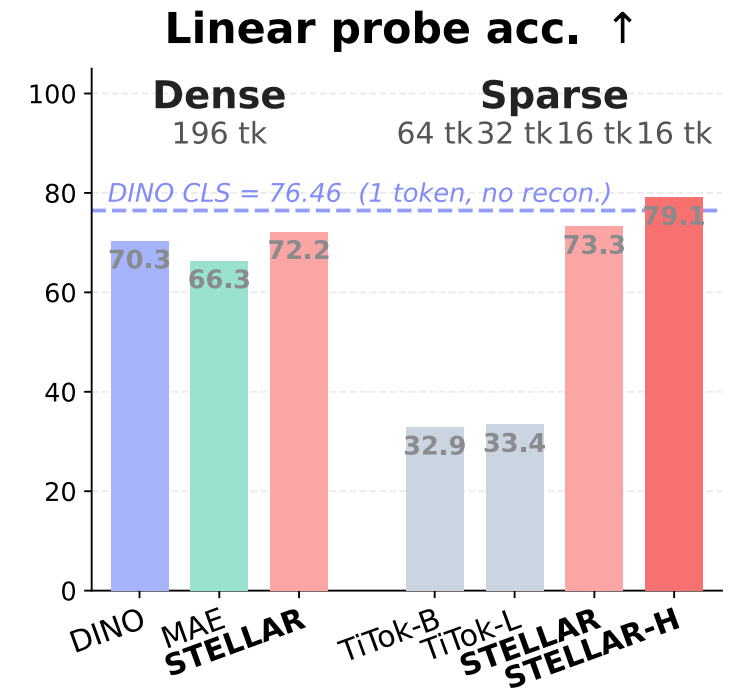
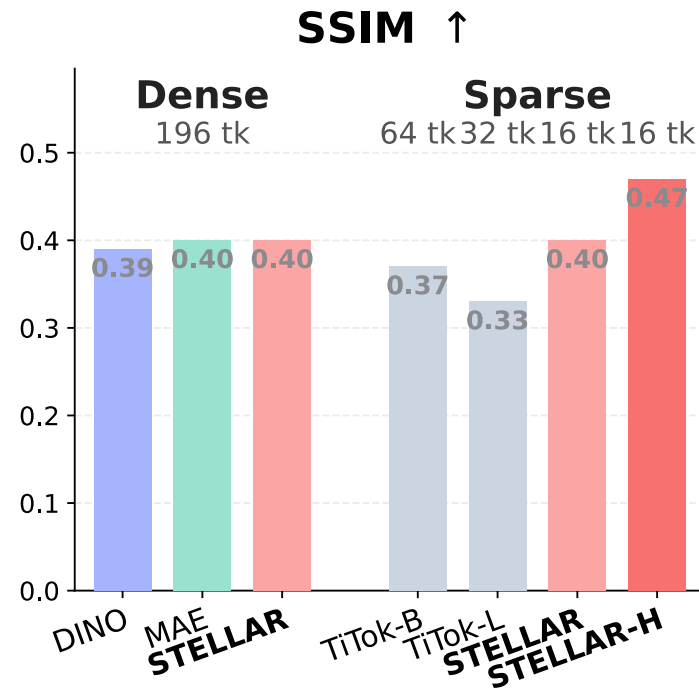
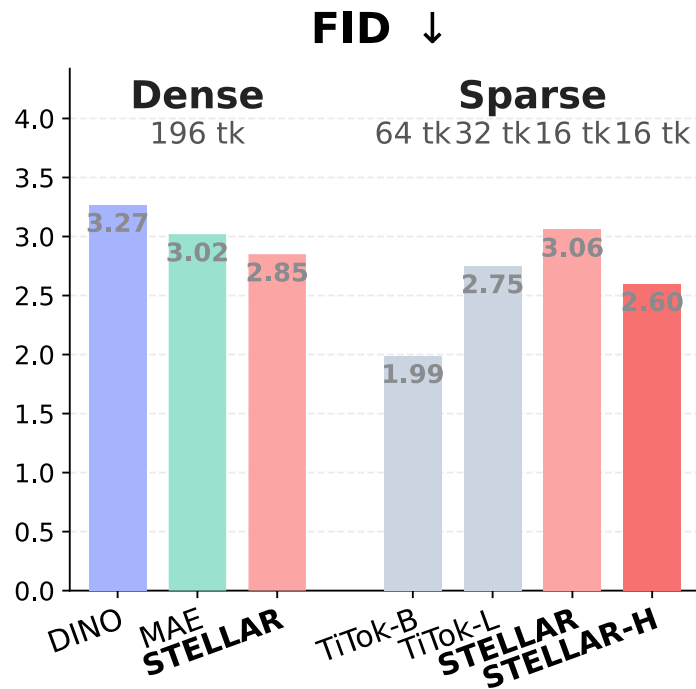
# Exp: Semantic Invariance



- Perform random crop and measure change in the latent representation in terms of cosine distance
- STELLAR has DINO-level semantic invariance, while other reconstruction methods have higher variance

# Exp: Unified Representation

- STELLAR achieves better reconstruction and spatial consistency than other dense and sparse representation
- Linear probing Acc outperformed other reconstruction feasible representations



# Exp: Downstream Transfer

- Top dense prediction performance with linear probing
- Global semantics outperformed reconstruction SSL methods

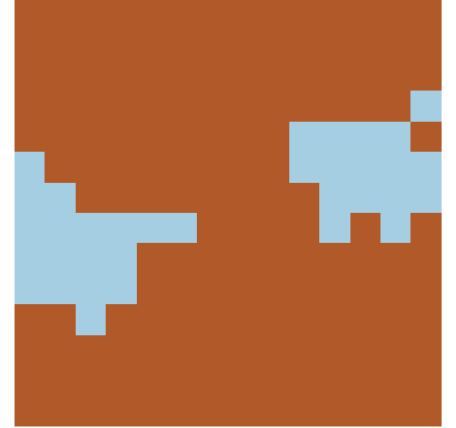
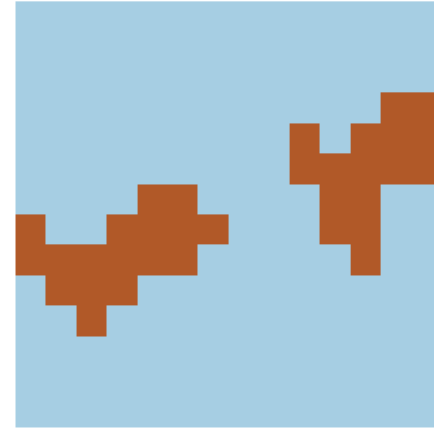
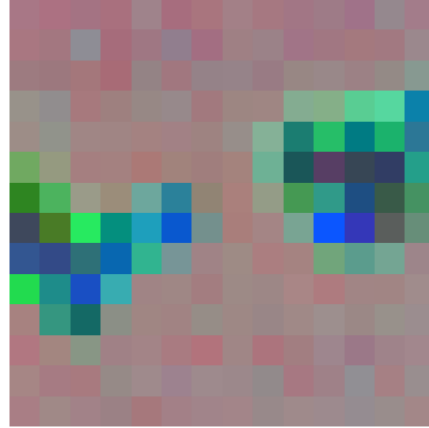
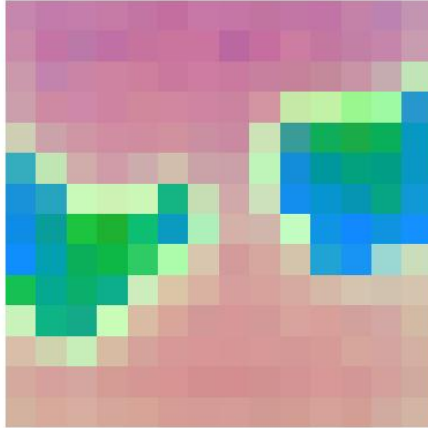
Model	Arch.	SSL Type		Segmentation (mIoU)			Classification (Acc)			
		Target	Method	ADE20K	CitySc	VOC	IN1K	Pets	Food	GlaS
<i>Semantic-Centric (Joint Embedding / Invariance)</i>										
BYOL	RN-50	GLOBAL	DISTILL	18.43	<u>18.66</u>	63.89	<u>70.39</u>	<u>82.77</u>	<u>64.57</u>	<u>95.00</u>
MoCo v3	ViT-B	GLOBAL	CONTR.	29.45	25.13	74.08	74.31	91.14	77.47	<b>97.50</b>
DINO	ViT-B	GLOBAL	DISTILL	26.87	26.82	79.29	<u>76.46</u>	<b>93.84</b>	<u>79.28</u>	95.00
MSN	ViT-B	GLOBAL	MASKING	26.66	25.39	68.59	73.65	75.91	68.93	92.50
DENSECL	RN-50	DENSE	CONTR.	<u>23.08</u>	18.63	<u>70.95</u>	61.10	72.99	59.16	85.00
DATA2VEC	ViT-B	DENSE	LAT-MIM	22.03	23.49	61.33	54.90	26.47	34.40	73.75
SIAMESEIM	ViT-B	DENSE	LAT-MIM	29.24	26.52	81.38	74.97	91.61	71.01	91.25
I-JEPA	ViT-H	DENSE	LAT-MIM	21.57	18.59	74.13	71.72	84.68	70.34	87.50
iBOT	ViT-B	GL+DE	DIST+MIM	<u>31.78</u>	25.69	77.06	76.40	92.40	78.08	96.25
iBOT	ViT-L	GL+DE	DIST+MIM	<u>33.26</u>	26.37	77.57	<u>78.53</u>	92.12	<b>81.07</b>	96.25
<i>Image-Centric (Reconstruction)</i>										
BEIT	ViT-B	DENSE	TOK MIM	11.58	18.90	27.44	32.94	36.20	54.49	90.00
BEIT	ViT-L	DENSE	TOK MIM	12.64	20.37	25.48	36.77	36.71	56.03	90.00
SIMMIM	SWIN-B	DENSE	PIX MIM	12.46	17.23	35.14	24.77	27.39	40.94	77.50
MAE	ViT-B	DENSE	PIX MIM	30.91	<u>29.44</u>	76.43	66.32	81.58	70.40	93.75
MAE	ViT-L	DENSE	PIX MIM	<u>34.36</u>	<u>32.53</u>	77.79	73.09	84.30	76.22	95.00
MAE	ViT-H	DENSE	PIX MIM	36.16	<b>35.21</b>	78.07	75.22	84.96	<u>78.36</u>	<u>95.00</u>
SEMMAE	ViT-B	DENSE	PIX MIM	3.52	25.48	48.33	43.84	56.99	58.90	92.50
TiTOK-64	ViT-B	SPARSE	SPRS REC	–	–	–	32.87	42.06	43.68	<b>97.50</b>
TiTOK-32	ViT-L	SPARSE	SPRS REC	–	–	–	33.42	27.83	38.83	78.75
<i>Our Method (Sparse Factorized Modeling)</i>										
<b>STELLAR</b>	ViT-B	SPARSE	INV+REC	31.33	27.74	<u>81.83</u>	73.26	89.70	74.09	95.00
<b>STELLAR</b>	ViT-L	SPARSE	INV+REC	34.02	31.32	<b>85.90</b>	76.94	<u>92.53</u>	74.78	<b>97.50</b>
<b>STELLAR</b>	ViT-H	SPARSE	INV+REC	<b>36.66</b>	33.30	<u>85.66</u>	<b>79.10</b>	<u>92.53</u>	77.43	92.50

# Visualization

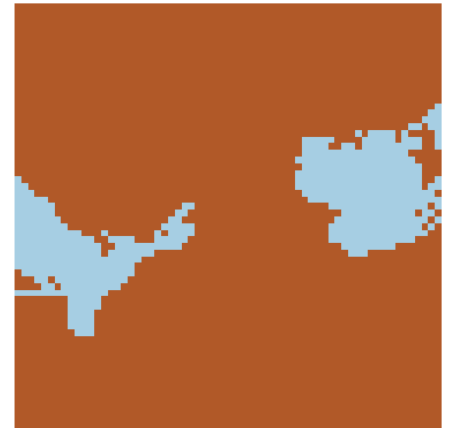
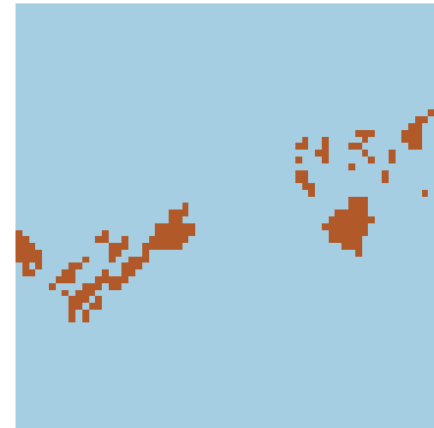
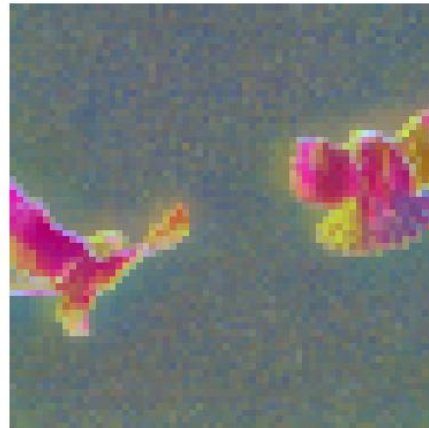
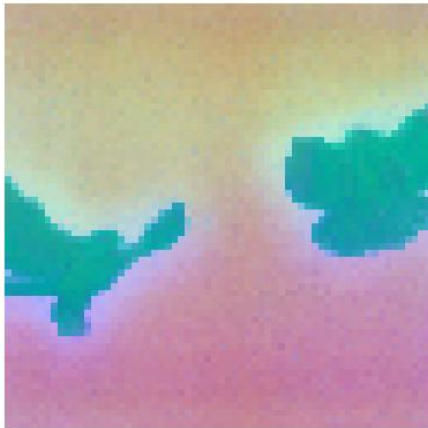
Feature map

Low-rank approx.

Example concept localization (threshold  $1/r$ )

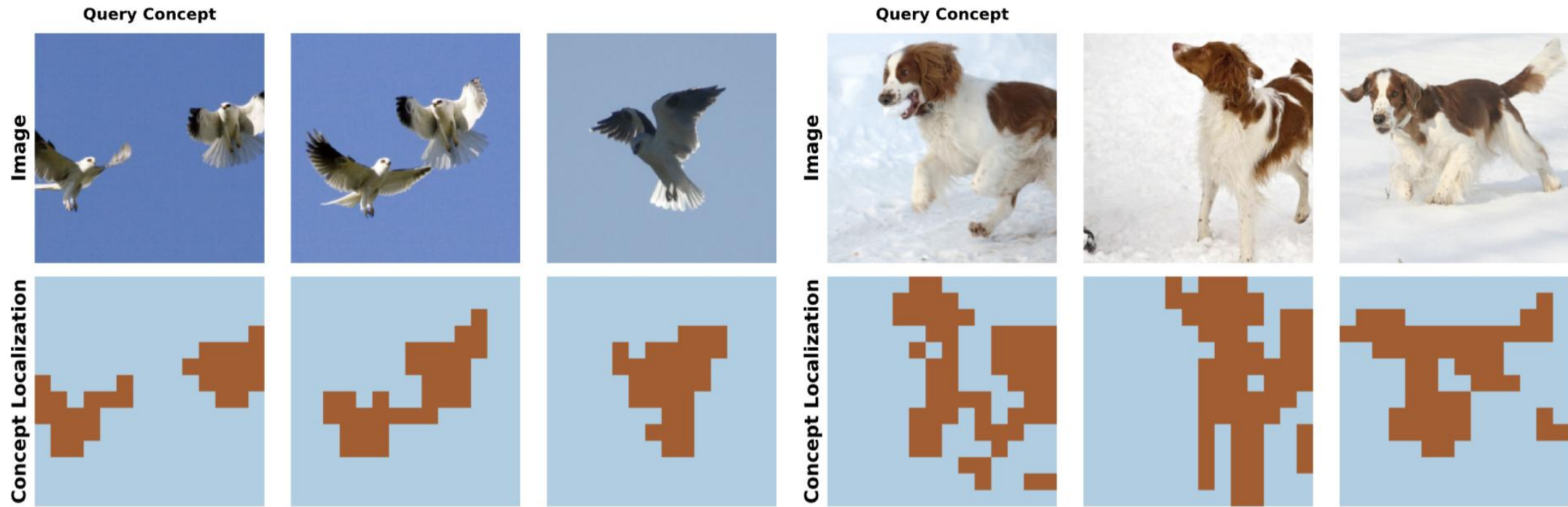


Directly scaling up res.  
to 1024 at inference



# Concept Retrieval

- Retrieve the closest tokens from the query concept, and visualize the image and the localization of the token in the image



# Sparse Reconstruction

Original



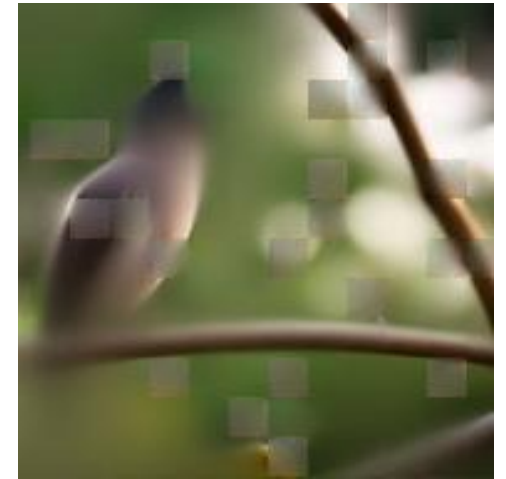
STELLAR 16 tokens



TiTok 32 tokens



MAE 32 tokens

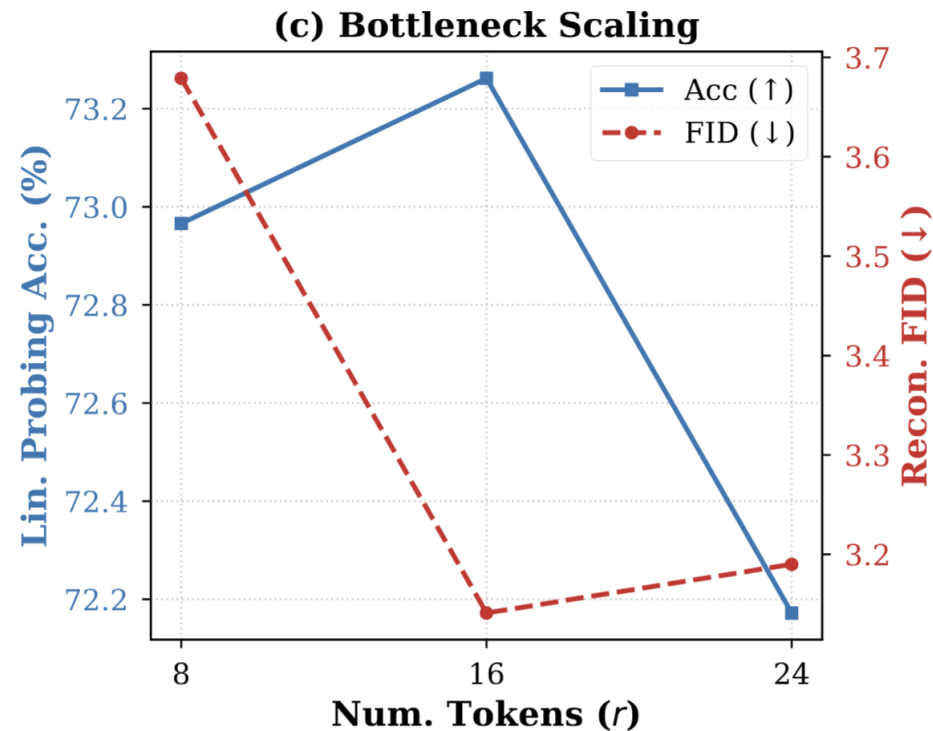


# Ablation Studies

	Recon.	Cluster	Set Align	CLS Align	KoLeo	rFID ↓	IN1K ↑	ADE ↑
DEFAULT	✓	✓	✓	✓	✓	<b>3.14</b>	<b>73.26</b>	<b>31.33</b>
<i>Impact of Individual Components</i>								
(A)	✗	✓	✓	✓	✓	—	72.44 (-0.82)	29.94 (-1.39)
(B)	✓	✗	✗	✗	✓	3.21 (+0.07)	52.07 (-21.19)	20.46 (-10.87)
(C)	✓	✓	✗	✗	✓	8.95 (+5.81)	2.73 (-70.53)	1.93 (-29.39)
(D)	✓	✗	✓	✓	✓	3.62 (+0.48)	42.14 (-31.12)	18.90 (-12.43)
(E)	✓	✓	✓	✗	✓	3.26 (+0.12)	70.79 (-2.47)	30.20 (-1.12)
(F)	✓	✓	✓	✓	✗	3.25 (+0.11)	72.05 (-1.21)	30.10 (-1.23)

# Effect of Number of Tokens/Rank

- Semantic favors less tokens (highly concentrate)
- Reconstruction favors more tokens(expressiveness)



# Takeaways

- Factorized latent supports joint learning of semantics and reconstruction
- Represent one image using **16 tokens** with **2.60 FID** and **79.1 lin. Acc.**
- Efficient latent representation (90% token reduction)

# Thanks!

Q&A