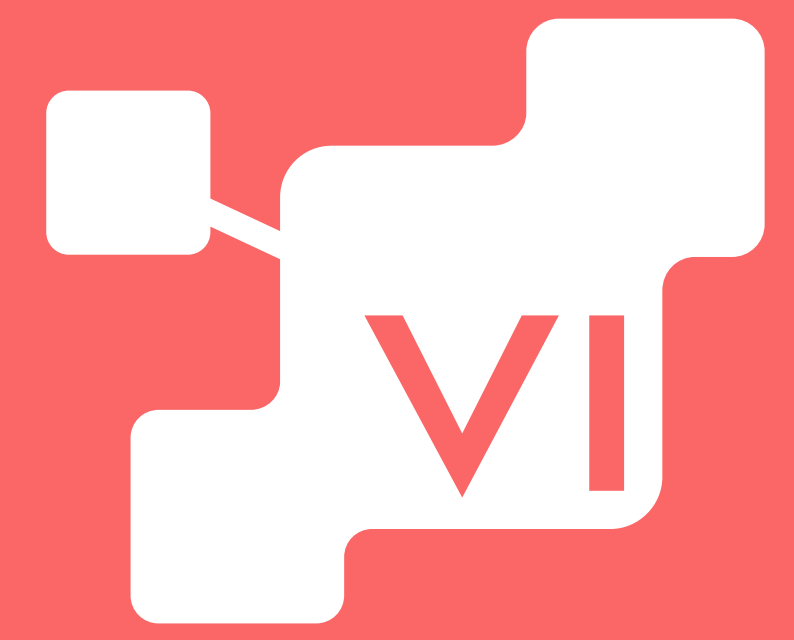


# A Spitting Image: Modular Superpixel Tokenization in Vision Transformers

M. Aasan, O. Kolbjørnsen, A. Schistad Solberg, and A. Ramírez Rivera

University of Oslo, SFI Visual Intelligence, AkerBP



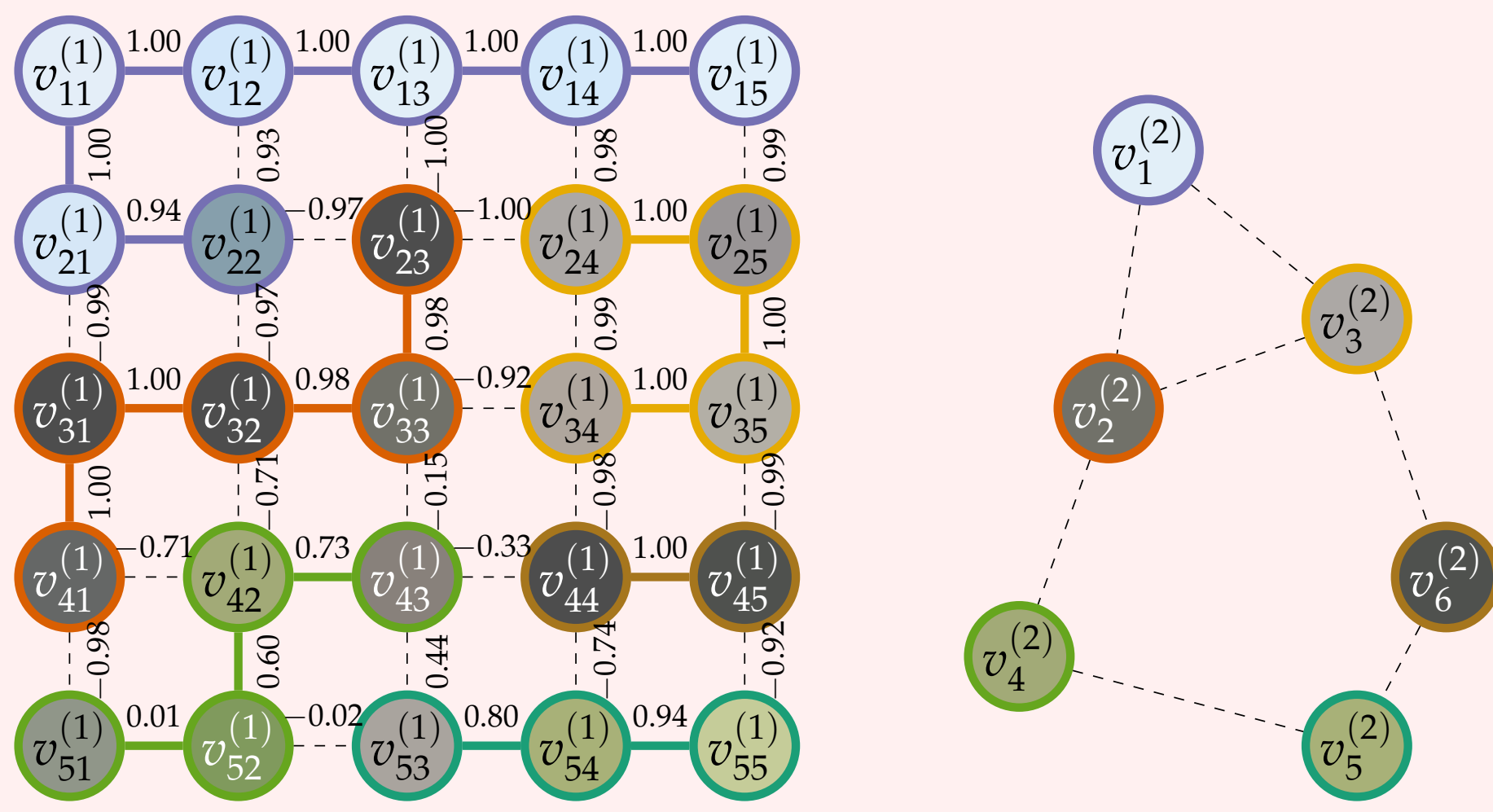
## Motivation

**Goal:** Adaptable tokenization for ViTs

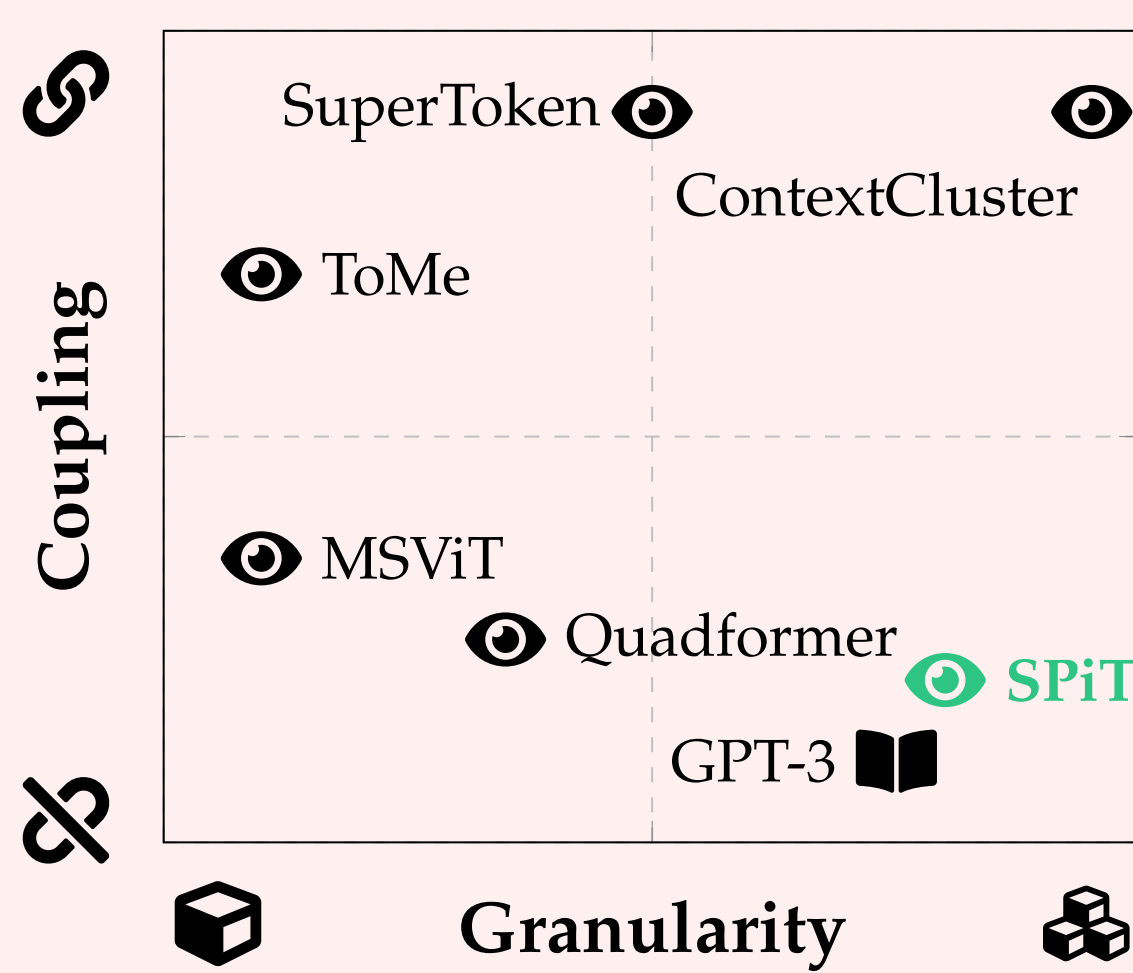
Flexible tokenization can more readily adapt to image content and context

- Superpixels generalize patch-based tokenization in ViTs
- Efficient tokenization via hierarchical graph-based superpixels
- Redined spatial resolution with pixel-level granularity for dense predictions

## Superpixel Aggregation



## Taxonomy



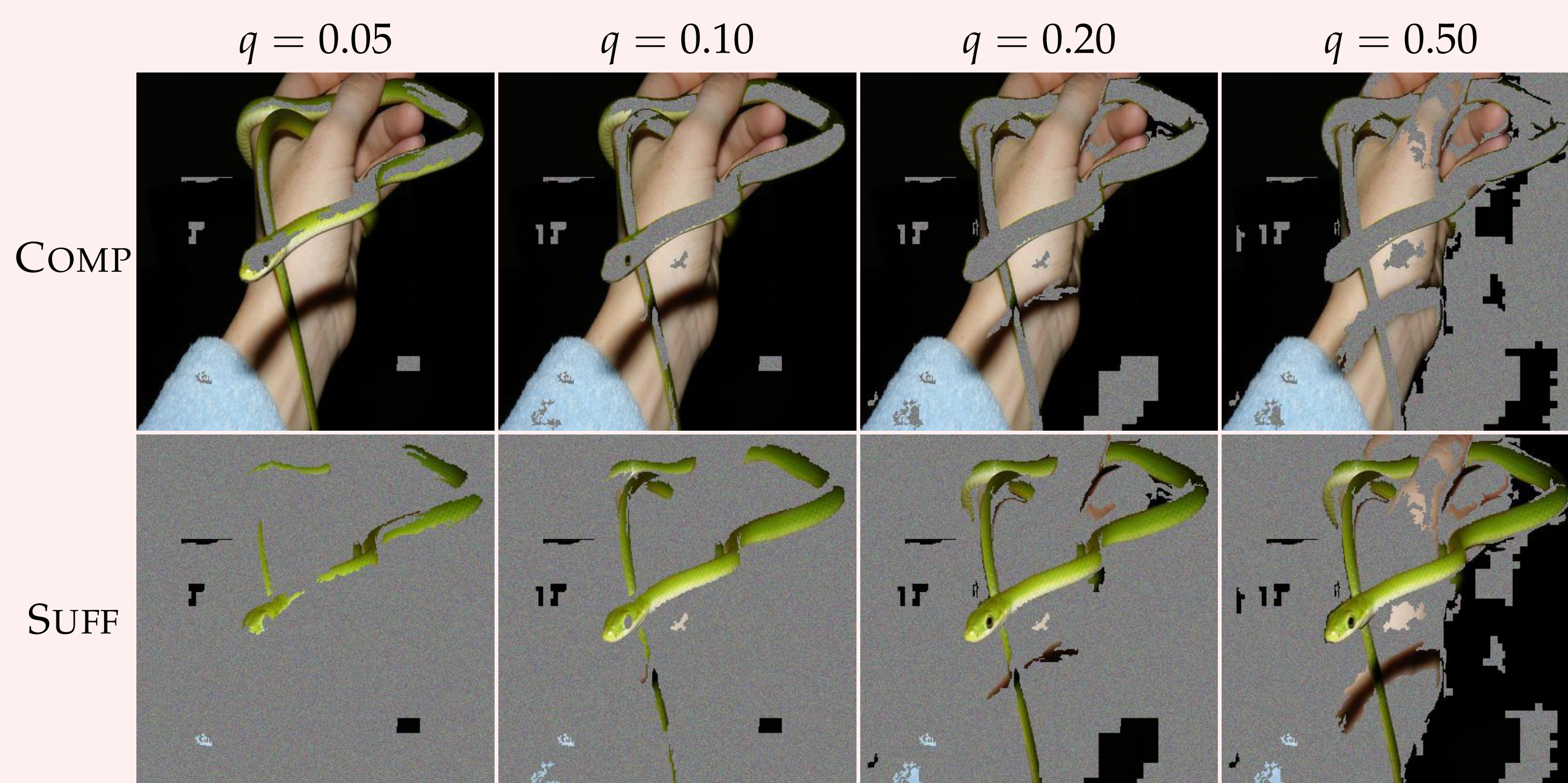
## Classification Accuracy

Model	Grad. Im./s.	INREAL		IN1K		CALTECH256		CIFAR100	
		Lin.	kNN	Lin.	kNN	Lin.	kNN	Lin.	kNN
ViT-B16	✗ 793.04	0.85	0.85	0.80	0.74	0.88	0.88	0.89	0.90
ViT-B16	✓ 721.12	0.85	0.84	<b>0.81</b>	0.75	<b>0.89</b>	0.89	<b>0.90</b>	<b>0.90</b>
RViT-B16	✗ 619.86	0.84	0.83	0.79	0.72	0.87	0.88	0.89	0.84
RViT-B16	✓ 585.64	0.84	0.84	0.79	0.73	0.86	0.86	0.89	0.76
SPiT-B16	✗ 690.72	0.79	0.82	0.76	0.57	0.83	0.83	0.81	0.63
SPiT-B16	✓ 640.59	<b>0.86</b>	<b>0.85</b>	0.80	<b>0.75</b>	0.89	<b>0.89</b>	0.88	0.85

## Occlusion and Faithfulness

$$\text{COMP}_{Q|x,\Phi} = \frac{1}{|Q|} \sum_{q \in Q} (\Phi(x) - \Phi(x \setminus x_{>q})),$$

$$\text{SUFF}_{Q|x,\Phi} = \frac{1}{|Q|} \sum_{q \in Q} (\Phi(x) - \Phi(x \setminus x_{\leq q})).$$



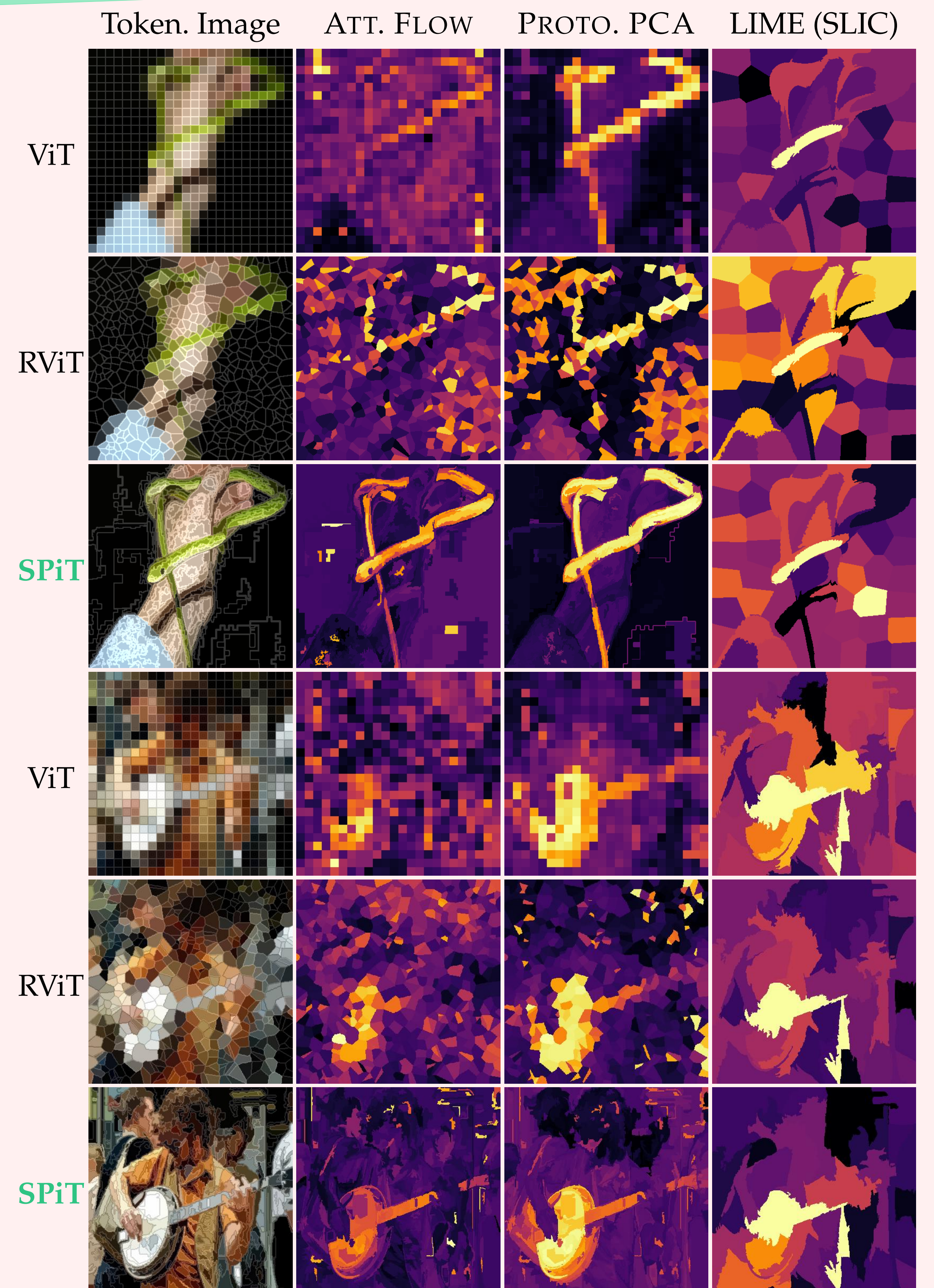
## Faithfulness of Attributions

	ViT-B16		RViT-B16		SPiT-B16	
	COMP ↑	SUFF ↓	COMP ↑	SUFF ↓	COMP ↑	SUFF ↓
LIME/SLIC	0.244	0.543	0.236	0.591	0.244	0.520
ATT.FLOW	0.160	0.664	0.223	0.685	<b>0.259</b>	0.558
PROT.PCA	0.206	0.710	0.209	0.691	0.256	0.592

## Project Page

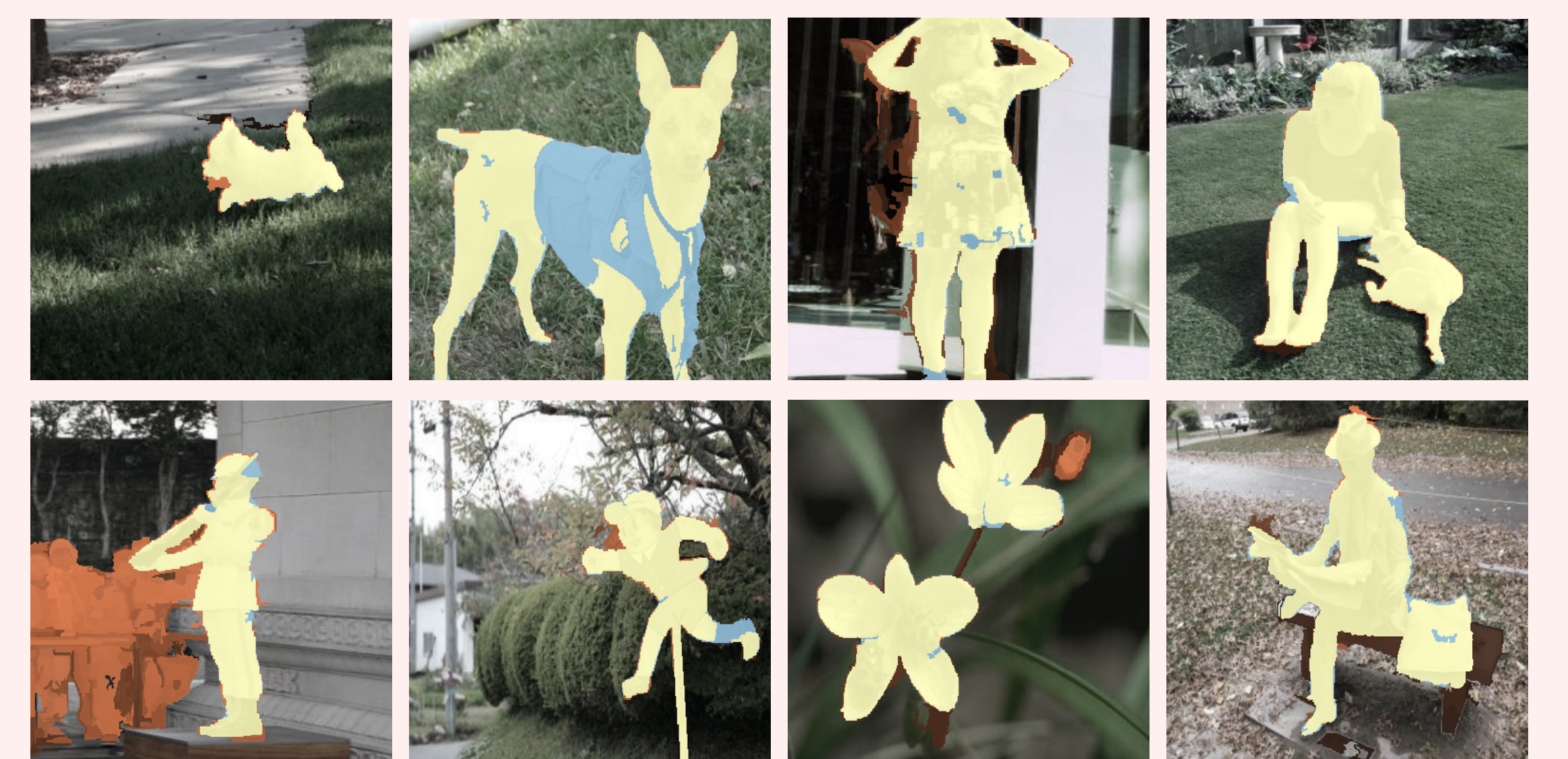


## Attributions



## Out-of-the-box

### Unsupervised Salient Segmentation



True positives False positives False negatives

### Source

### Feature Correspondences

