



# Visual Attention Emerges from Recurrent Sparse Reconstruction

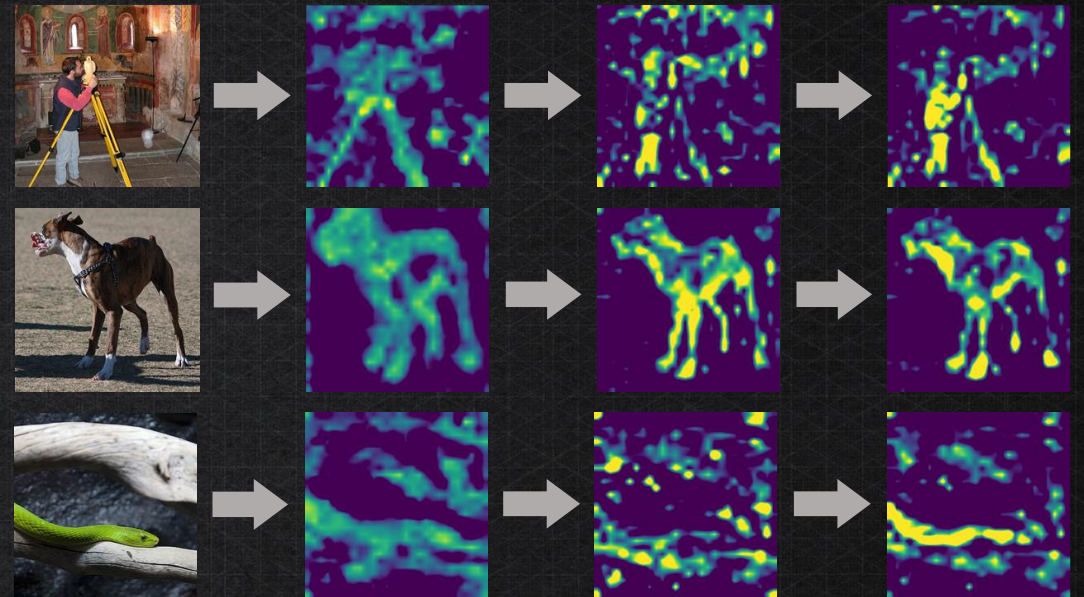
Baifeng Shi

University of California, Berkeley

*Joint work with Yale Song, Neel Joshi, Trevor Darrell, Xin Wang*

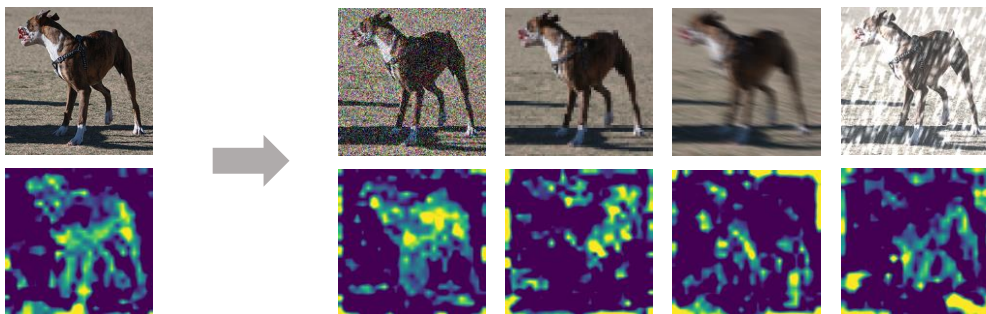
*Appearing at T4V workshop, CVPR 2022, New Orleans*

*and ICML 2022, Baltimore*





# Self-Attention is still Imperfect



- Self-attention based transformers are still not robust enough.
- Here the model failed under image corruption



Attention maps are not matching the intuitive human attention represented by eye fixation maps.

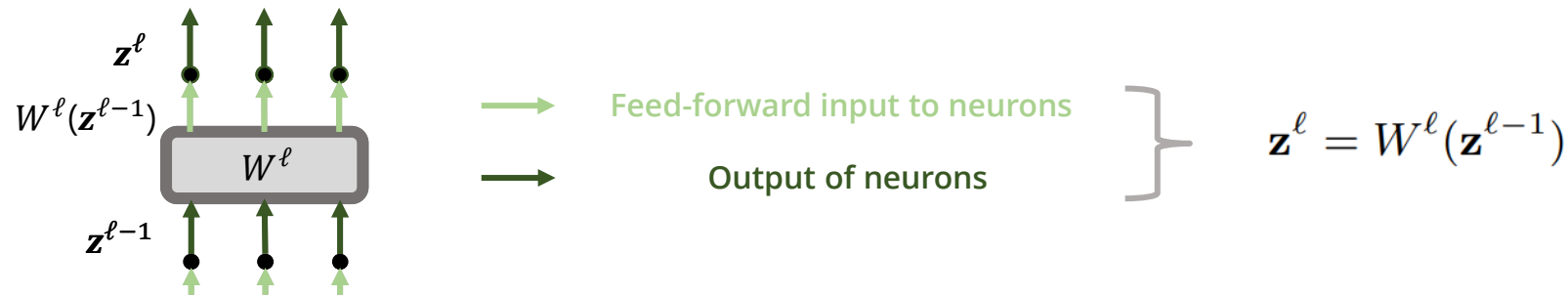
# VARs: Visual Attention emerges from Recurrent Sparse reconstruction

Inspired by two key features in human visual attention: **recurrency** and **sparsity**,

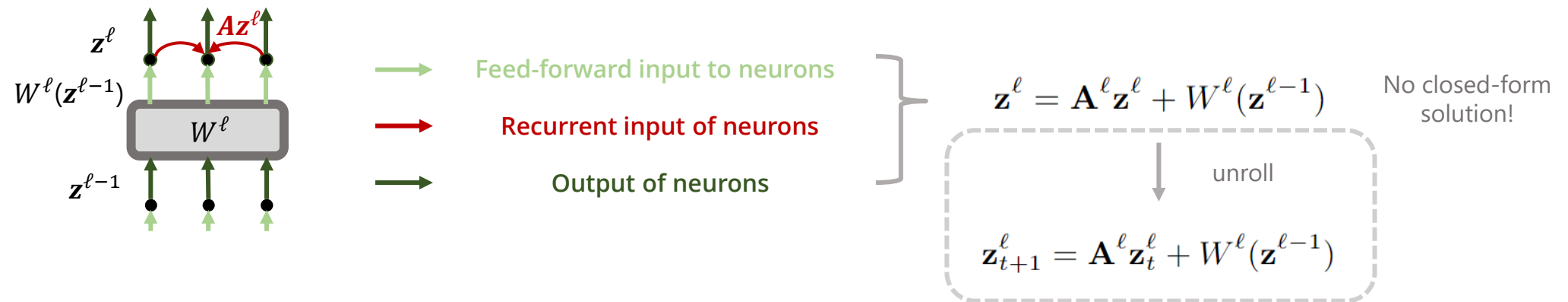
- We propose VARs, Visual Attention from Recurrent Sparse reconstruction, a new attention algorithm
  - which can be plugged into various neural networks to improve model robustness.
- We also show connections between VARs and existing attention algorithms and human attention
  - Self-attention is a special case of VARs.
  - VARs is closer to human attention compared to existing methods.

# Adding Recurrency to Feedforward Networks

## Feedforward Networks

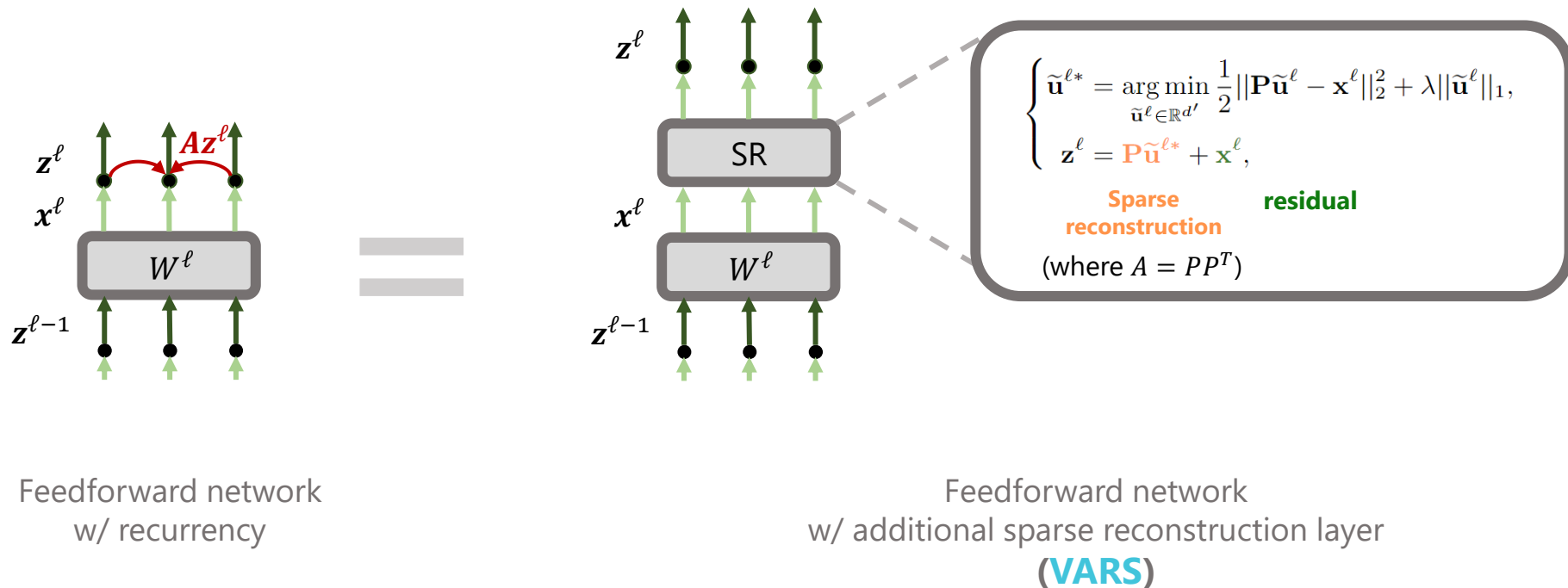


## Feedforward Networks w/ recurrency



# Recurrency Entails Sparse Reconstruction

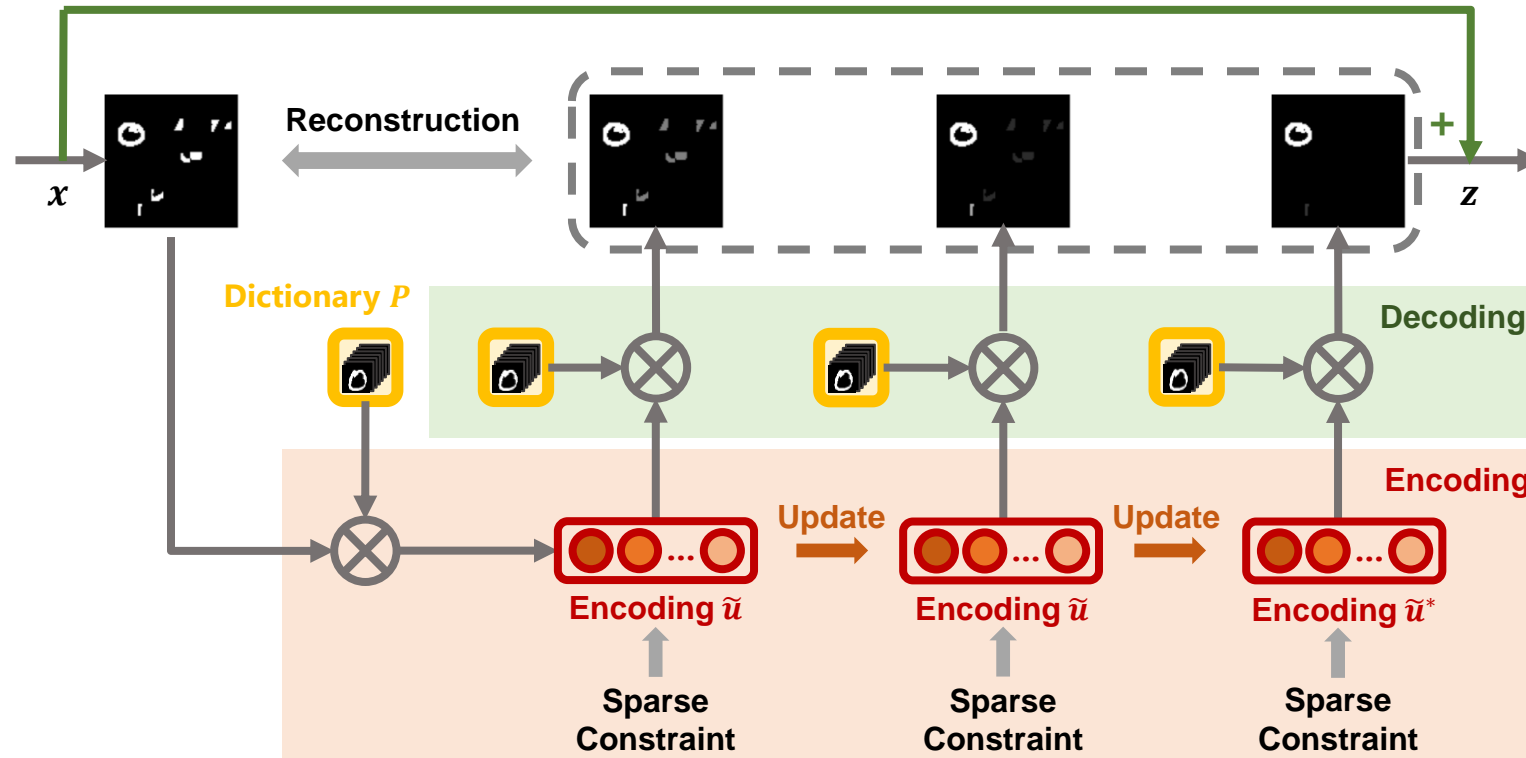
Recurrently connected network = Feedforward network w/ additional sparse reconstruction layer\*.



\* Under certain conditions.

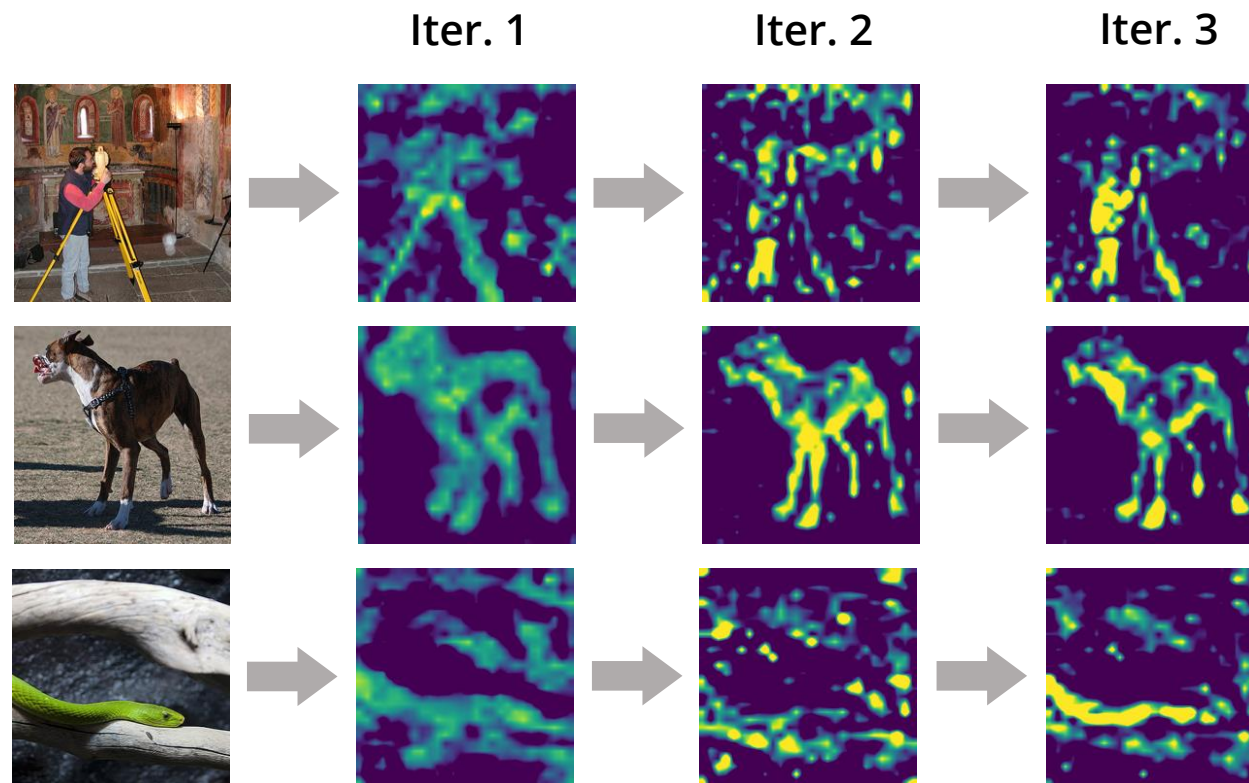


# Attention Emerges from Sparse Reconstruction



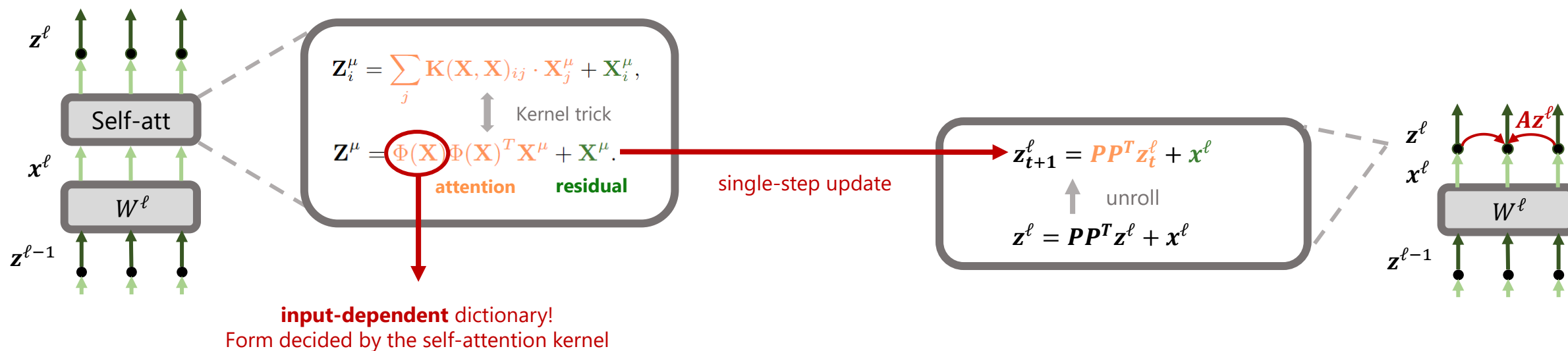
In each VARS block, we iteratively optimize the sparse reconstruction using ISTA (LISTA)

# Iterative Improvement of Attention Maps of VARS



The feature map is slowly focused on the salient objects.

# Self-Attention as a Special Case of VARS





# Performance on Robustness Benchmarks

## Our models:

- VARS-S: **static** dictionary
- VARS-D: **dynamic** dictionary (similar to self-attention)
- VARS-SD: **static + dynamic**

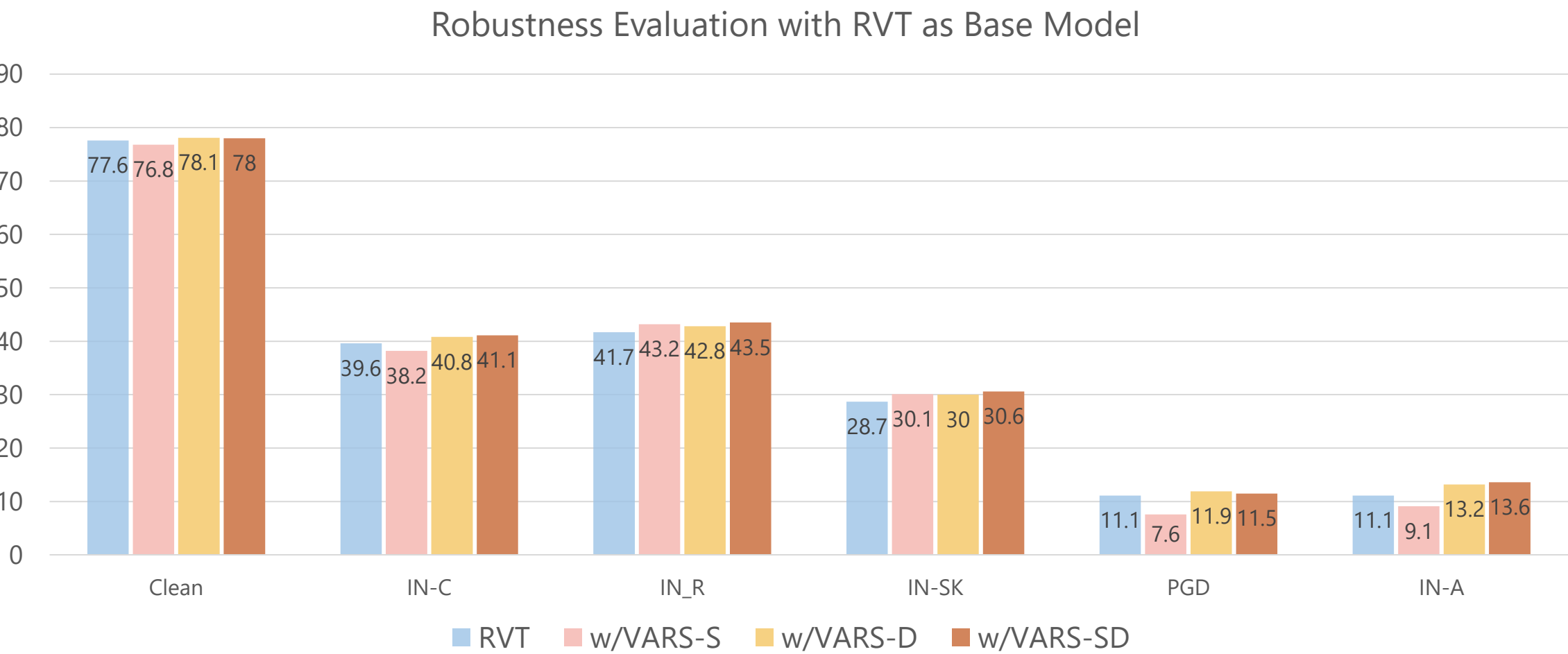
## Baseline Architectures:

- DeiT
- RVT (Robust Vision Transformer)
- GFNet (Global Filter Network)

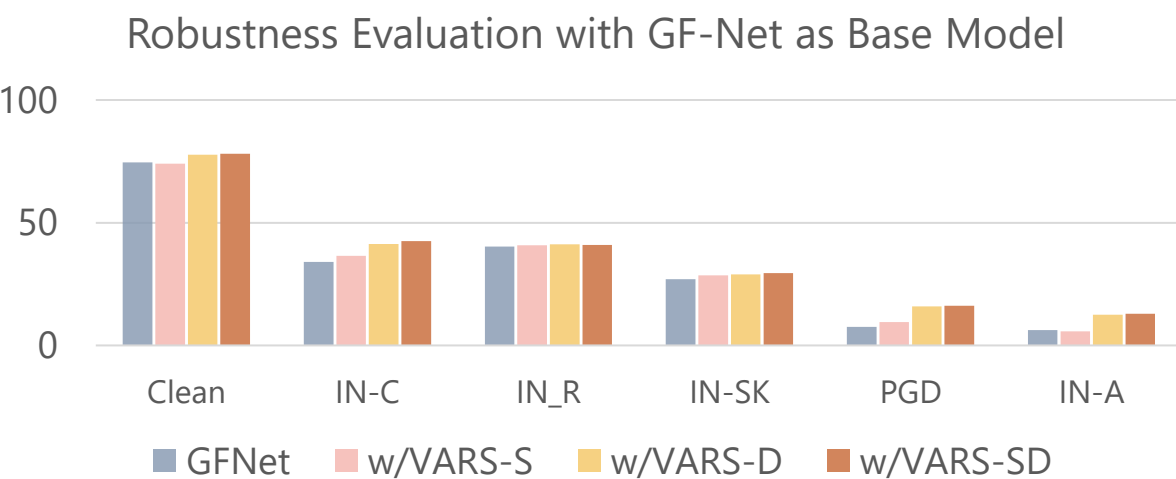
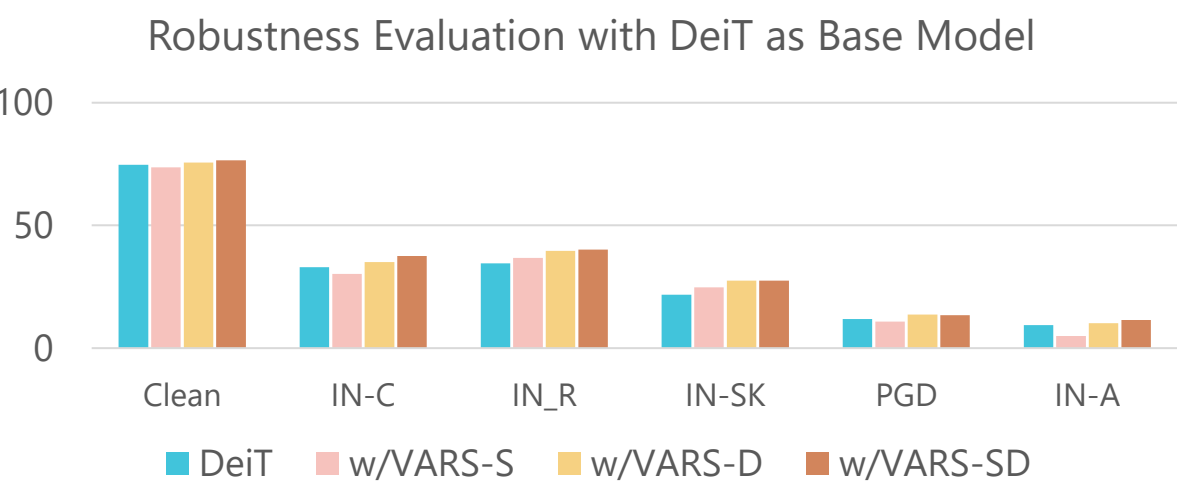
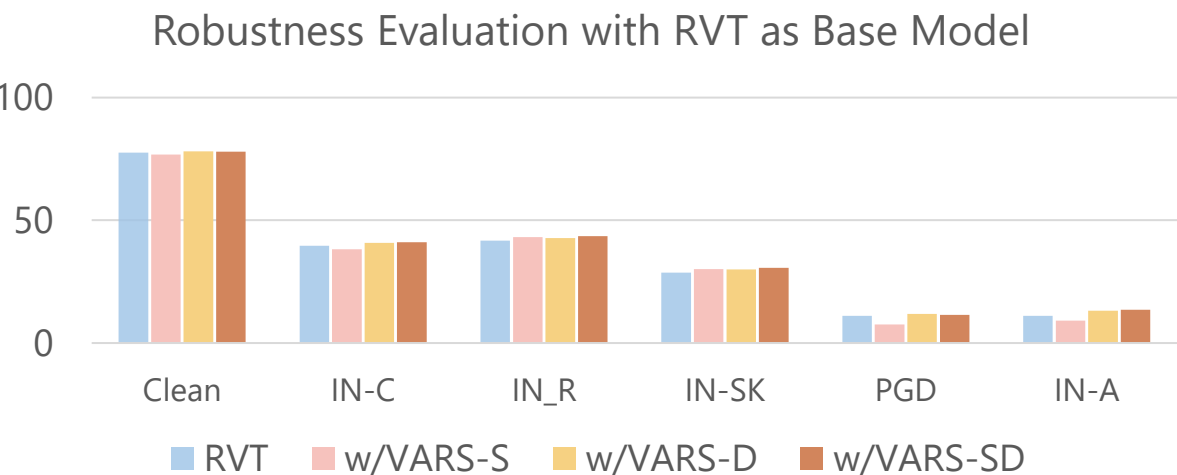
## Datasets:

	Dataset Name	Type
Robustness	ImageNet-C	Natural corruption
	ImageNet-R	Out of distribution
	ImageNet-SK	Out of distribution
	PGD	Adversarial attacks
	ImageNet-A	Natural adv. examples
Others	PACS	Domain generalization
	PASCAL VOC	Semantic segmentation
	MIT1003	Human eye fixation

# Performance on Robustness Benchmarks

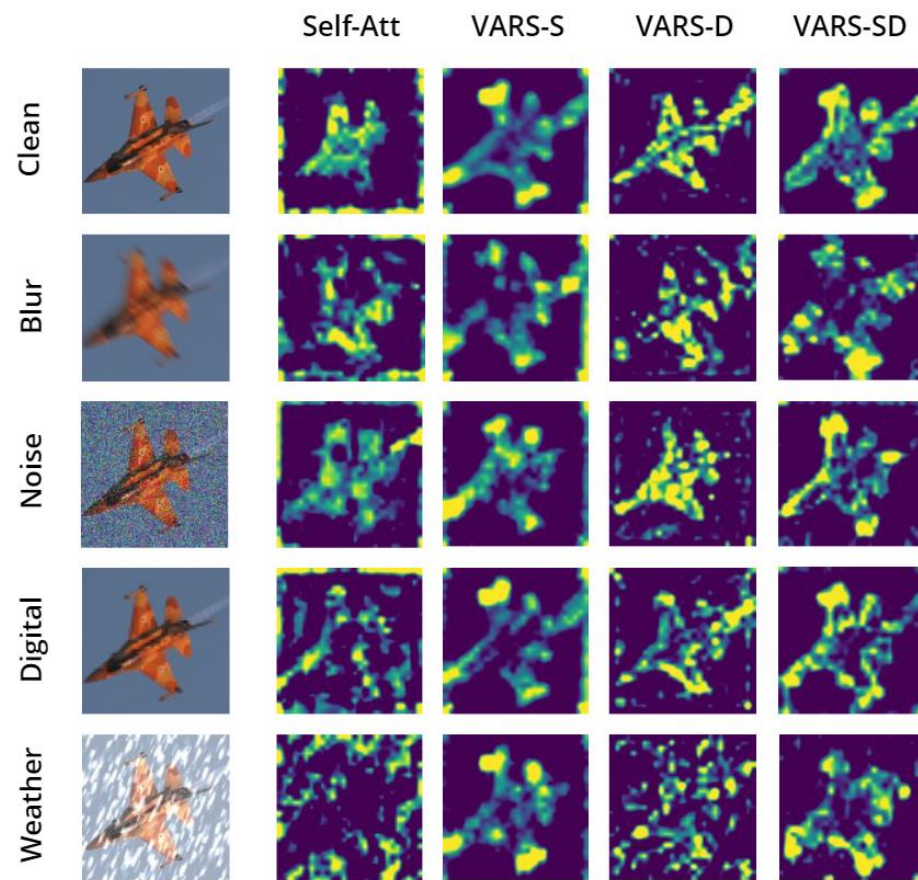


# Performance on Robustness Benchmarks



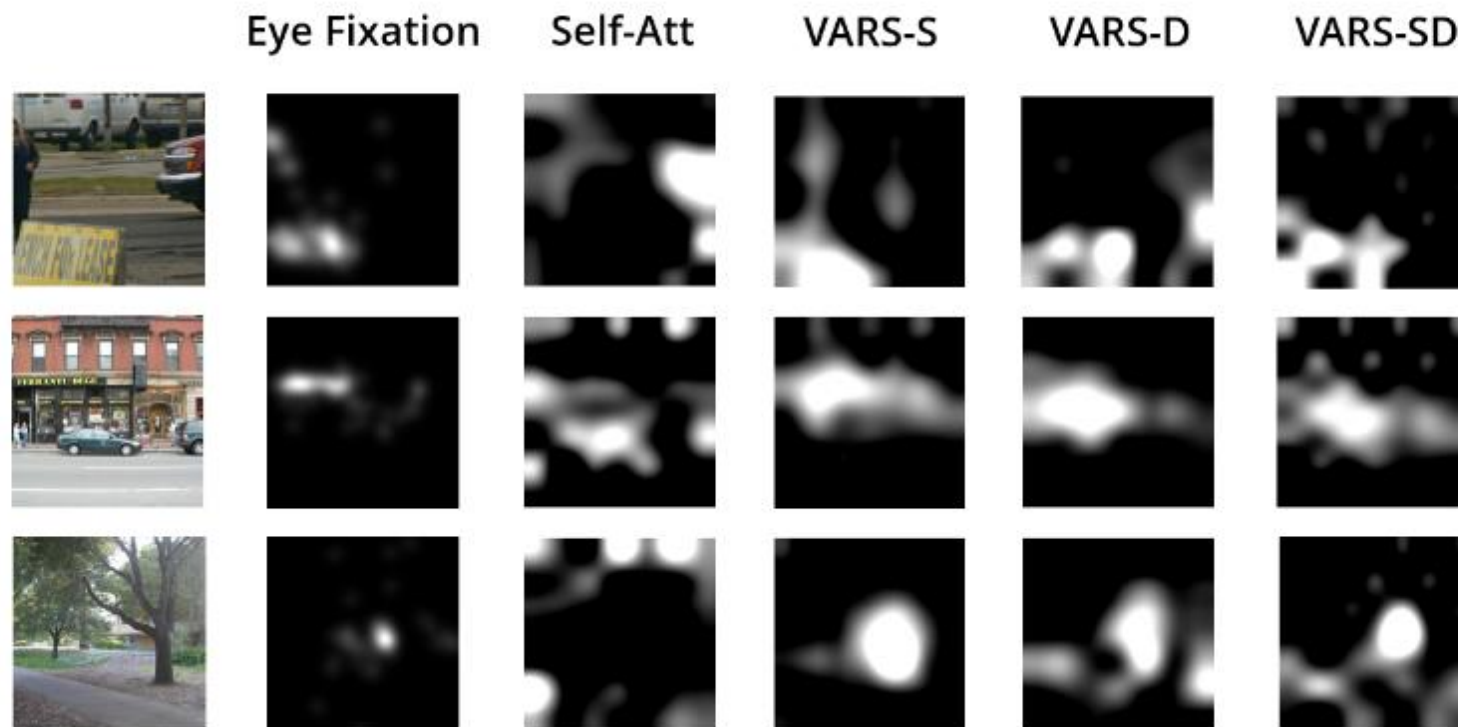


# Performance on Robustness Benchmarks



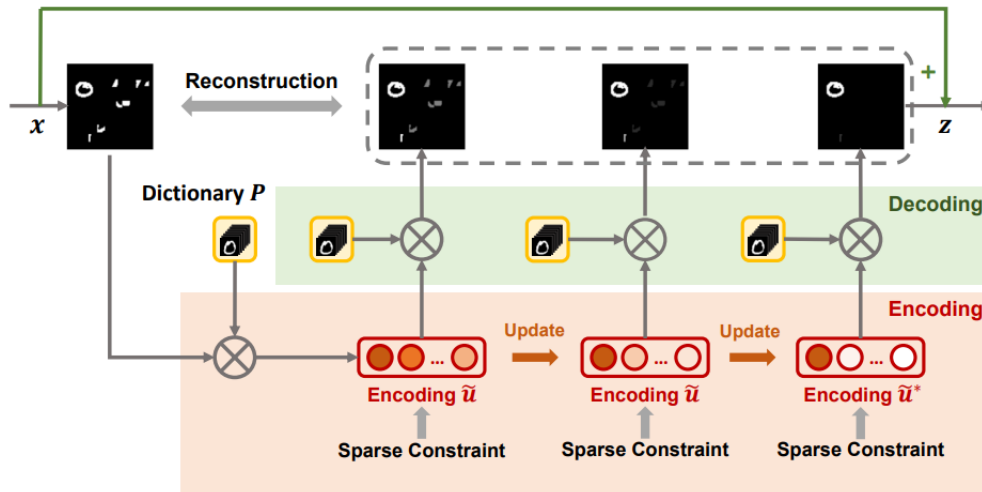
VAR-S with dynamic dictionaries identify better salient regions than the others.

# VARS is More Consistent with Human Eye Fixation



VARS's attention maps are more consistent with human eye fixation than self-attention's.

# Summary



- We propose **VARs**, a new attention design,
  - which can be plugged into various neural networks to improve model robustness.
- We also show that
  - VARs is closer to human attention compared to existing methods.
  - Self-attention is also a special case of VARs.



Thank you!

Baifeng Shi  
bfshi.github.io  
baifeng\_shi@berkeley.edu

