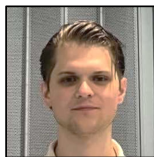


# — DAVE —

Distribution-Aware Attribution via ViT Gradient DEcomposition



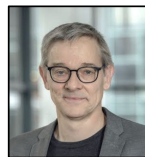
Adam Wróbel \*



Siddhartha Gairola \*



Jacek Tabor



Bernt Schiele



Bartosz Zieliński



Dawid Rymarczyk

# ViT Attribution Gap

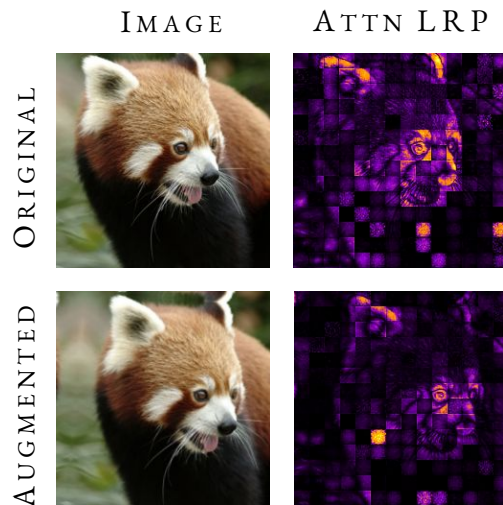
## ARCHITECTURE - INDUCED ARTIFACTS

- **ViTs** dominate modern vision.
- Pixel-level attributions suffer from **artifacts**.
- Patch-level attributions **lack precision**.
- Stable high-resolution explanations are missing.
- **DAVE**: artifact-free, pixel-level attribution.

# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

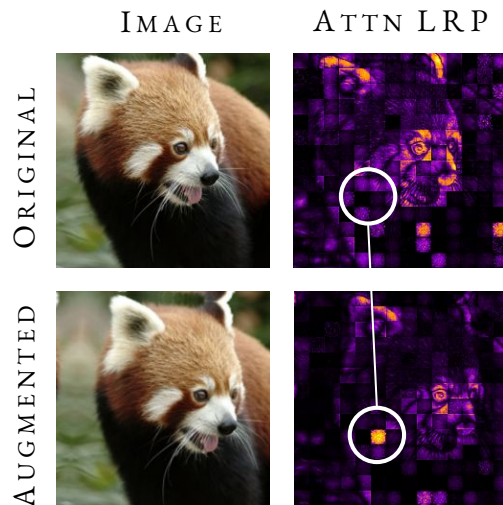
- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- Patch-level attributions lack precision.
- Stable high-resolution explanations are missing.
- **DAVE**: artifact-free, pixel-level attribution.



# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

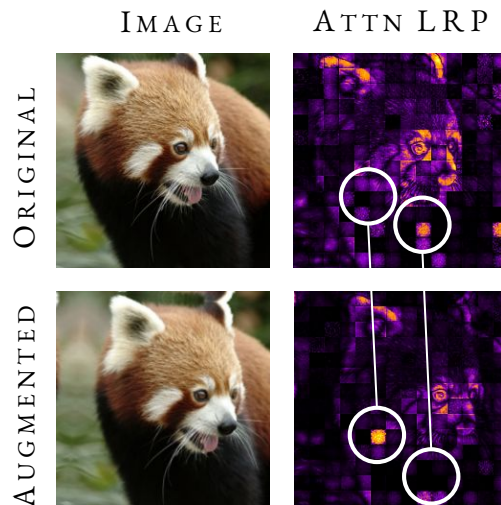
- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- Patch-level attributions lack precision.
- Stable high-resolution explanations are missing.
- **DAVE**: artifact-free, pixel-level attribution.



# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

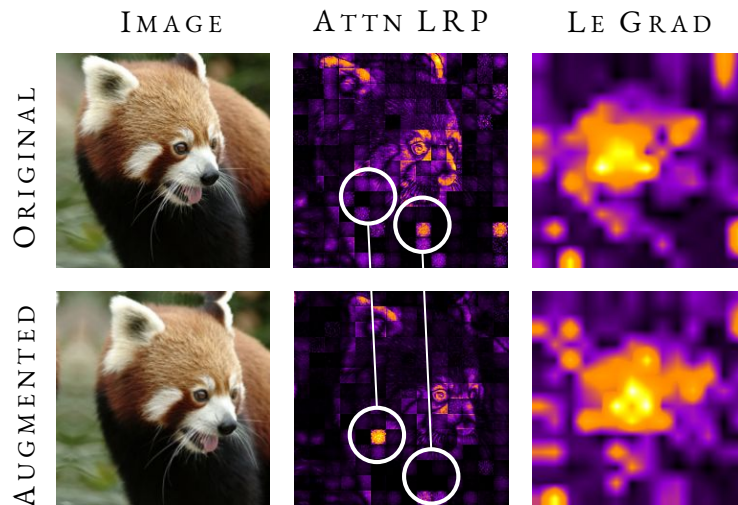
- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- Patch-level attributions lack precision.
- Stable high-resolution explanations are missing.
- **DAVE**: artifact-free, pixel-level attribution.



# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

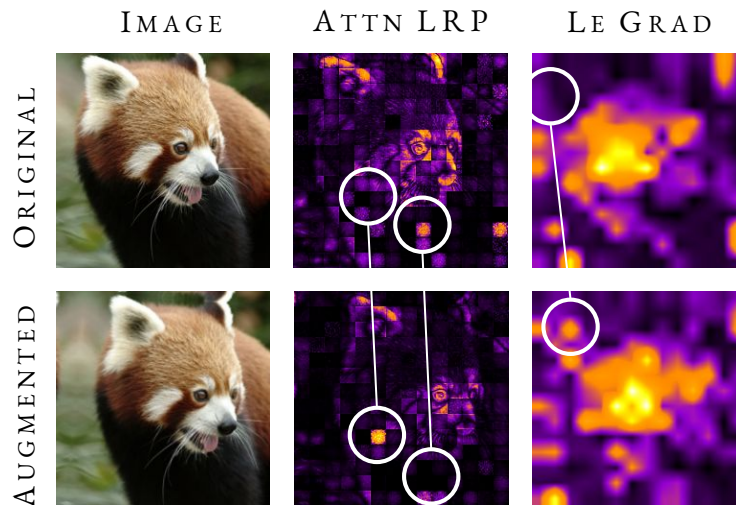
- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- **Patch-level** attributions **lack precision**.
- Stable high-resolution explanations are missing.
- **DAVE**: artifact-free, pixel-level attribution.



# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

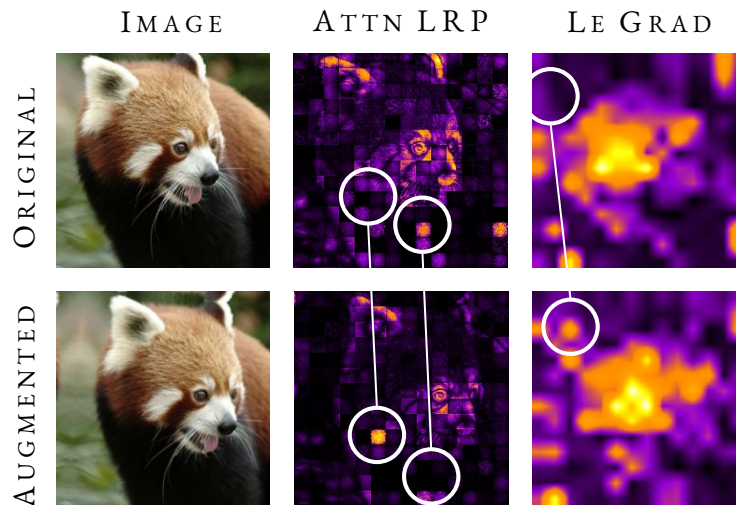
- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- **Patch-level** attributions **lack precision**.
- Stable high-resolution explanations are missing.
- **DAVE**: artifact-free, pixel-level attribution.



# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

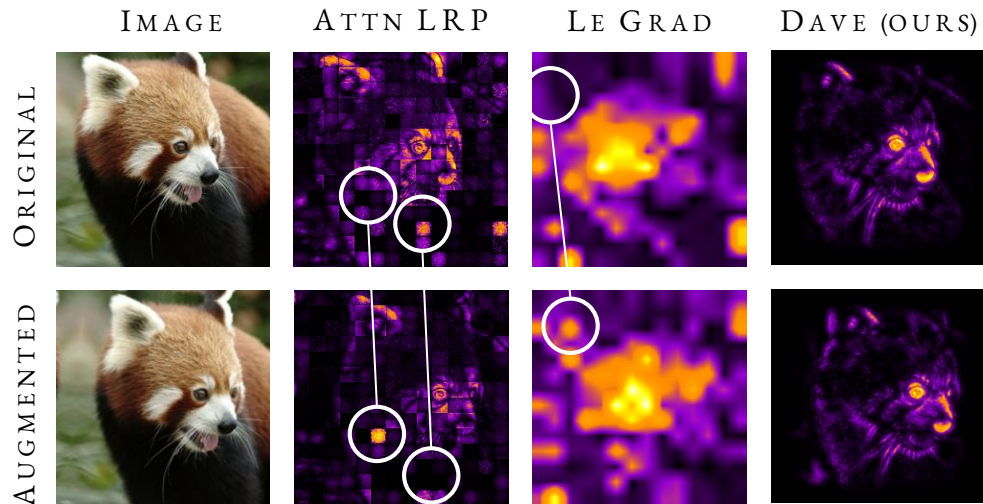
- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- **Patch-level** attributions **lack precision**.
- Stable high-resolution explanations are missing.
- **DAVE: artifact-free, pixel-level attribution.**



# ViT Attributions Gap

## ARCHITECTURE - INDUCED ARTIFACTS

- **ViTs** dominate modern vision.
- **Pixel-level** attributions suffer from **artifacts**.
- **Patch-level** attributions **lack precision**.
- Stable high-resolution explanations are missing.
- **DAVE**: **artifact-free**, **pixel-level** attribution.



# DAVE Construction

## STEP I: EFFECTIVE TRANSFORMATION

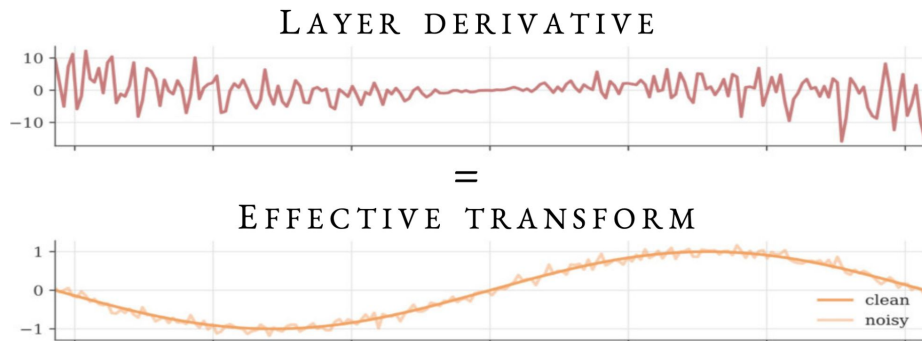
$D_x F$   
layer  
derivative



# DAVE Construction

## STEP I: EFFECTIVE TRANSFORMATION

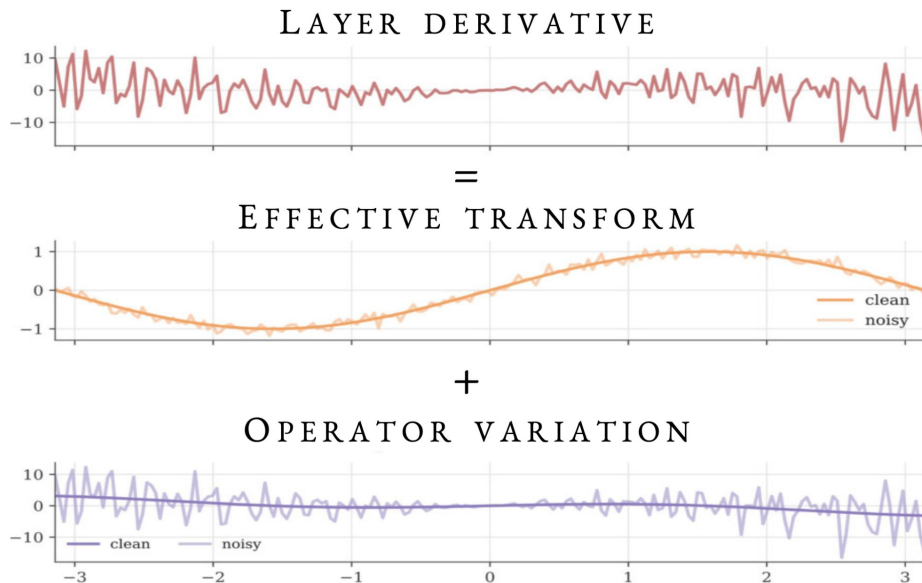
$$\underbrace{D_{\mathbf{X}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}}$$



# DAVE Construction

## STEP I: EFFECTIVE TRANSFORMATION

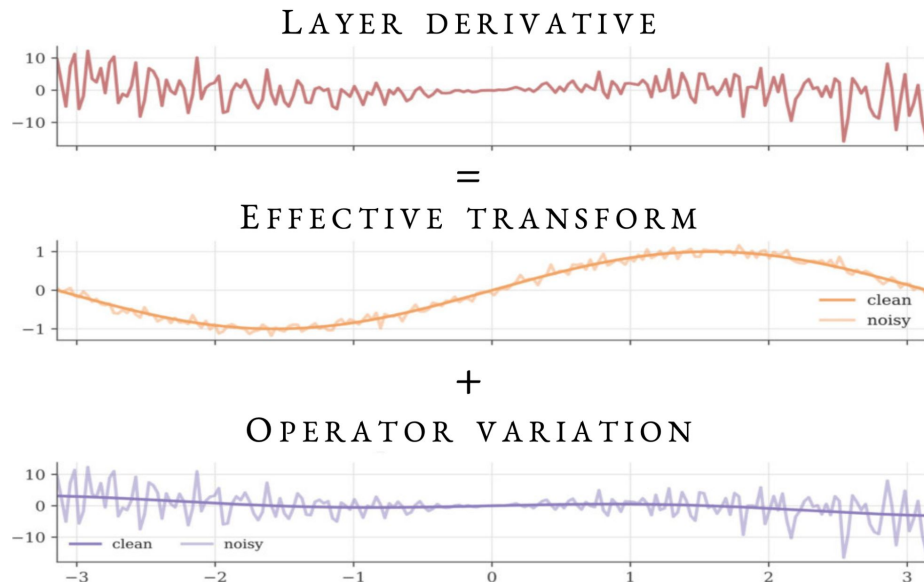
$$\underbrace{D_{\mathbf{X}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}} + \underbrace{\left( (D_{\mathbf{X}}L(\mathbf{X})(\cdot)) \mathbf{X} \right)}_{\text{operator variation}}$$



# DAVE Construction

## STEP I: EFFECTIVE TRANSFORMATION

$$\underbrace{D_{\mathbf{X}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}} + \underbrace{\cancel{((D_{\mathbf{X}}L(\mathbf{X})(\cdot))\mathbf{X})}}_{\text{operator variation}}$$

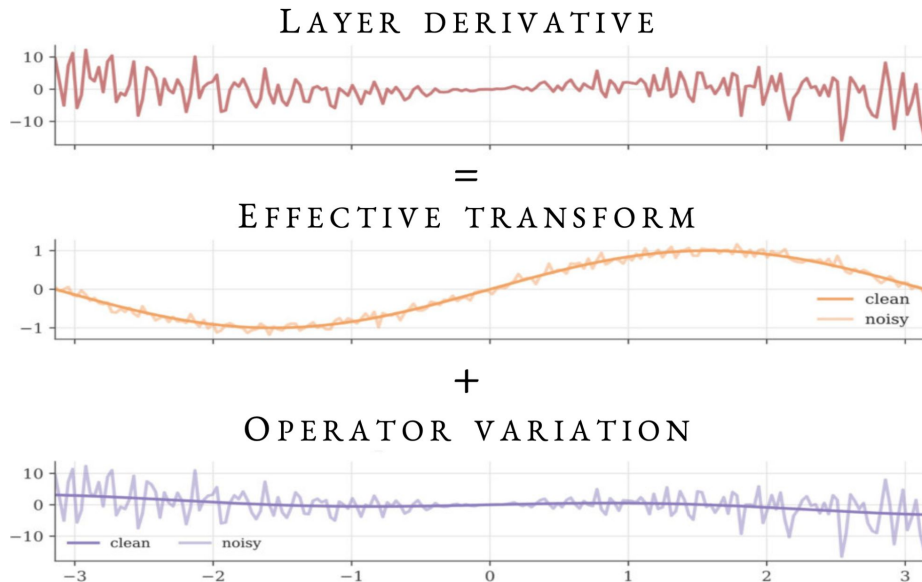


# DAVE Construction

## STEP I: EFFECTIVE TRANSFORMATION

$$\underbrace{D_{\mathbf{X}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}} + \underbrace{\cancel{((D_{\mathbf{X}}L(\mathbf{X})(\cdot))\mathbf{X})}}_{\text{operator variation}}$$

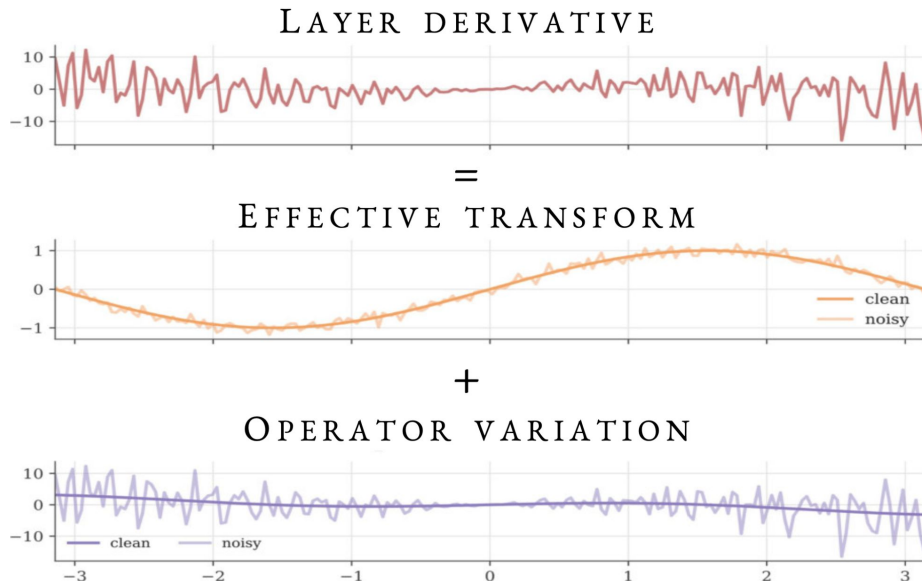
IMAGE



# DAVE Construction

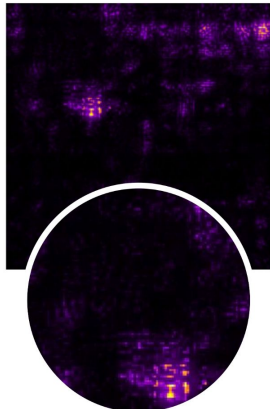
## STEP I: EFFECTIVE TRANSFORMATION

$$\underbrace{D_{\mathbf{x}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}} + \underbrace{\cancel{\left( (D_{\mathbf{x}}L(\mathbf{X})(\cdot)) \mathbf{X} \right)}}_{\text{operator variation}}$$



IMAGE

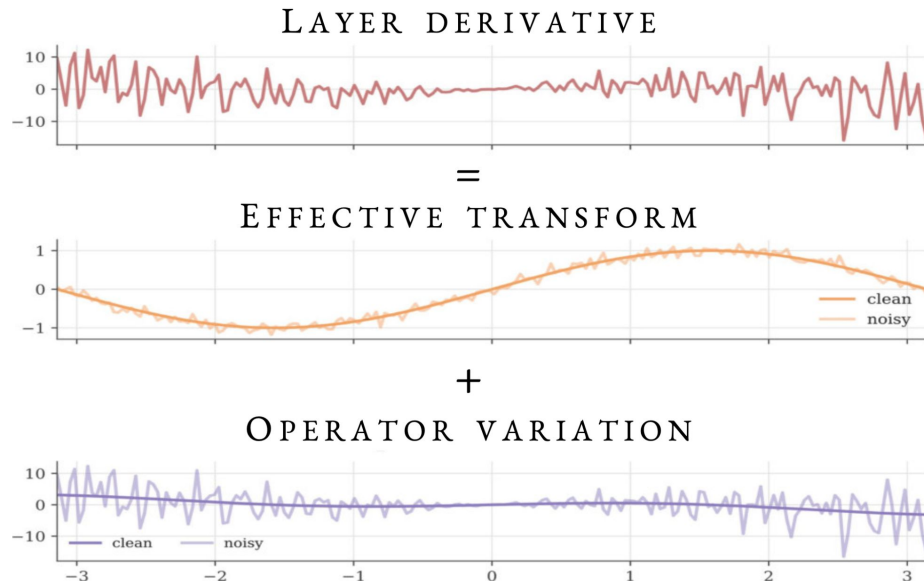
GRADIENT



# DAVE Construction

## STEP I: EFFECTIVE TRANSFORMATION

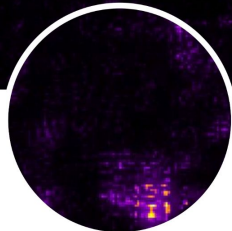
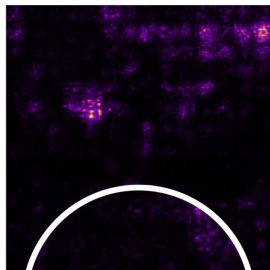
$$\underbrace{D_{\mathbf{X}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}} + \underbrace{\cancel{((D_{\mathbf{X}}L(\mathbf{X})(\cdot))\mathbf{X})}}_{\text{operator variation}}$$



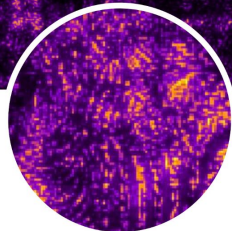
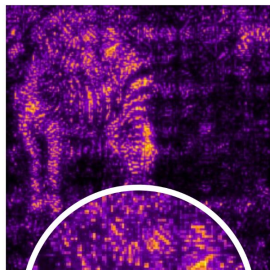
IMAGE



GRADIENT



EFFECTIVE



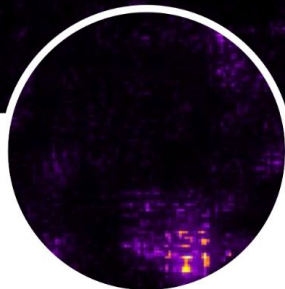
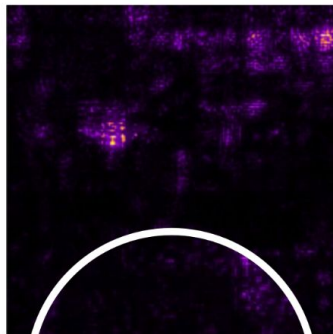
# DAVE Construction

## STEP II: EQUIVARIANT TRANSFORMATION

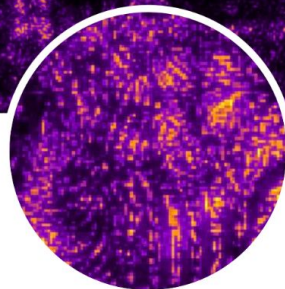
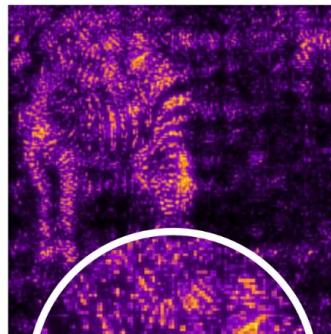
IMAGE



GRADIENT



EFFECTIVE



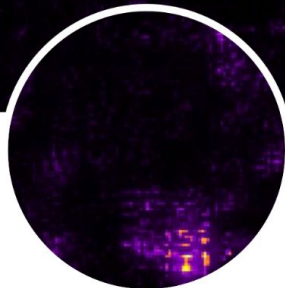
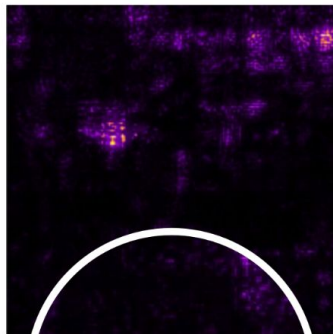
# DAVE Construction

## STEP II: EQUIVARIANT TRANSFORMATION

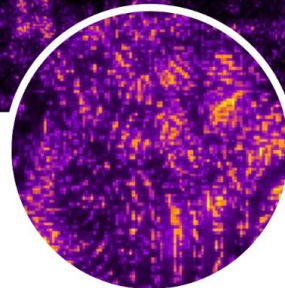
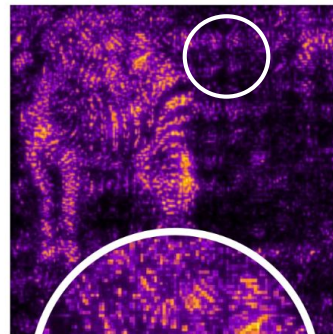
IMAGE



GRADIENT



EFFECTIVE



# DAVE Construction

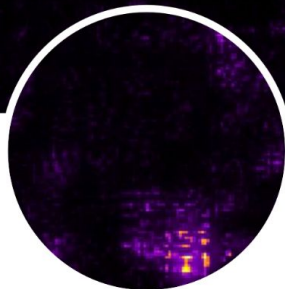
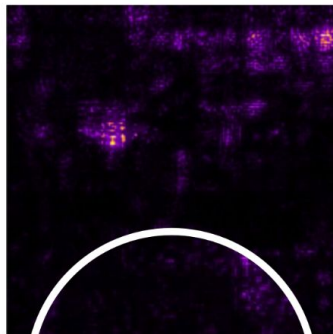
## STEP II: EQUIVARIANT TRANSFORMATION

$$\mathbf{W}_L^{eq}(\mathbf{X}) = \int_G [\tau^{-1} \circ \mathbf{W}_L \circ \tau](\mathbf{X}) d\nu(\tau)$$

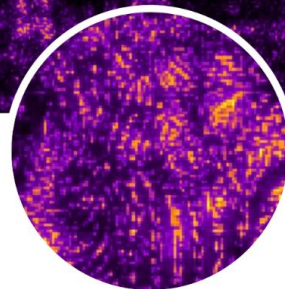
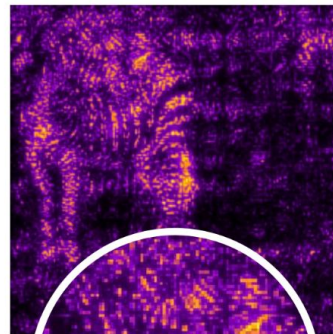
IMAGE



GRADIENT



EFFECTIVE



# DAVE Construction

## STEP II: EQUIVARIANT TRANSFORMATION

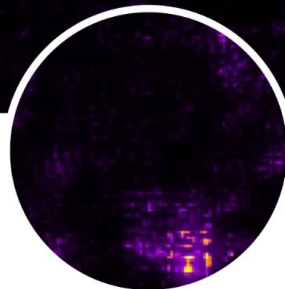
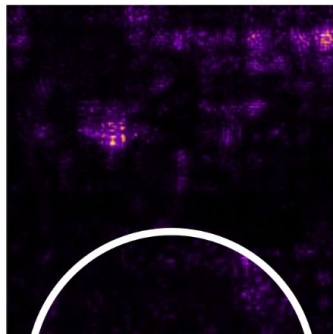
$$W_L^{eq}(\mathbf{X}) = \int_G [\tau]^{-1} \circ W_L \circ [\tau](\mathbf{X}) d\nu([\tau])$$

Rotations / Translations

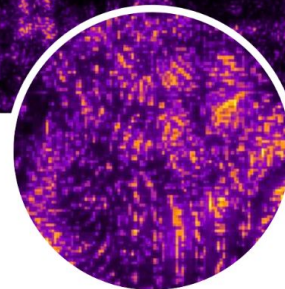
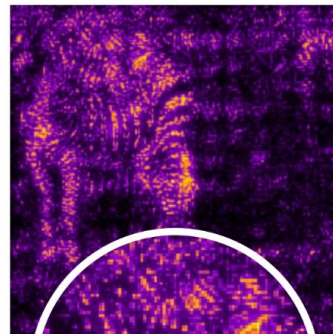
IMAGE



GRADIENT



EFFECTIVE



# DAVE Construction

## STEP II: EQUIVARIANT TRANSFORMATION

$$W_L^{eq}(\mathbf{X}) = \int_G [\tau^{-1} \circ W_L \circ \tau](\mathbf{X}) d\mathcal{U}(\tau)$$

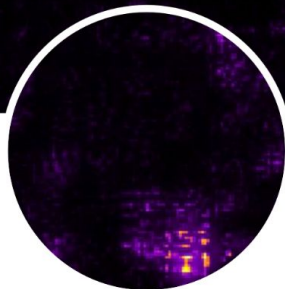
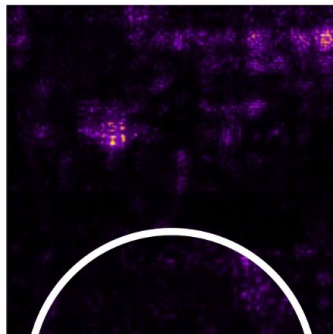
Rotations / Translations

Local neighborhood

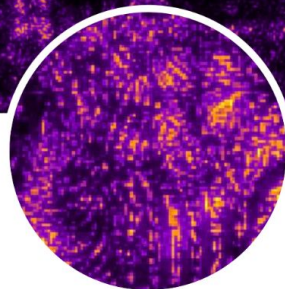
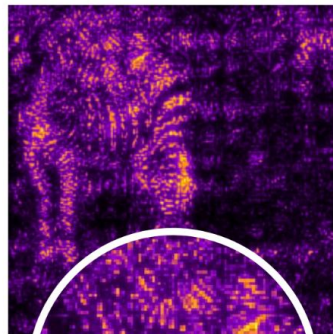
IMAGE



GRADIENT



EFFECTIVE



# DAVE Construction

## STEP II: EQUIVARIANT TRANSFORMATION

$$W_L^{eq}(\mathbf{X}) = \int_G [\tau^{-1} \circ W_L \circ \tau](\mathbf{X}) d\mathcal{U}(\tau)$$

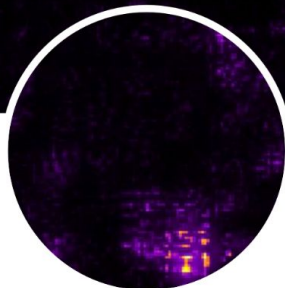
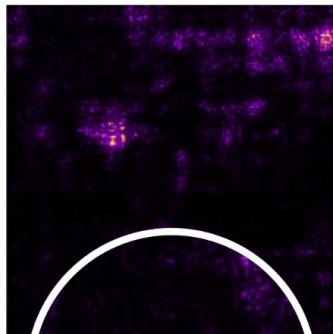
Rotations / Translations

Local neighborhood

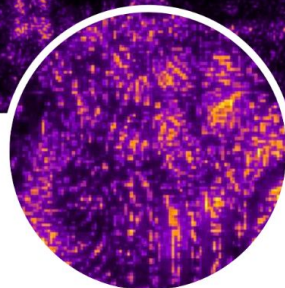
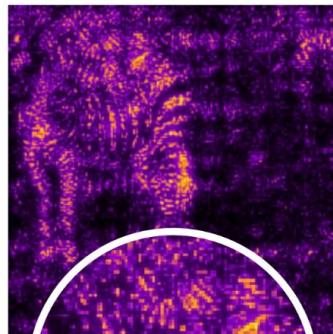
IMAGE



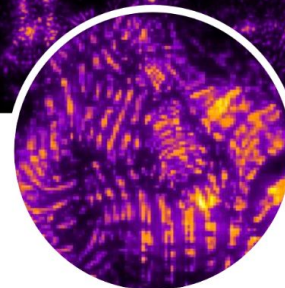
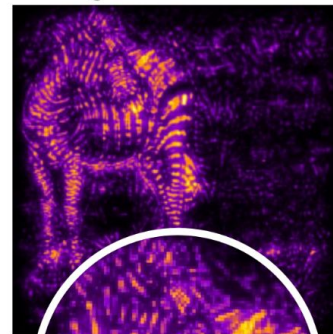
GRADIENT



EFFECTIVE



EQUIVARIANT



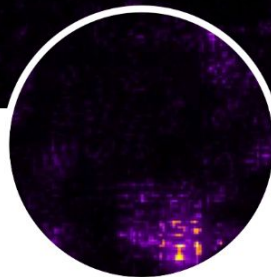
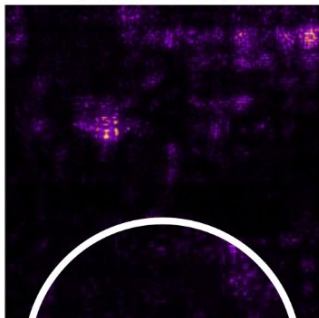
# DAVE Construction

## STEP III: LOW-PASS FILTER

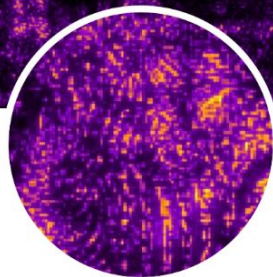
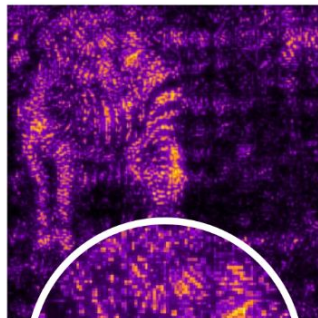
IMAGE



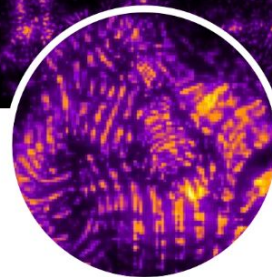
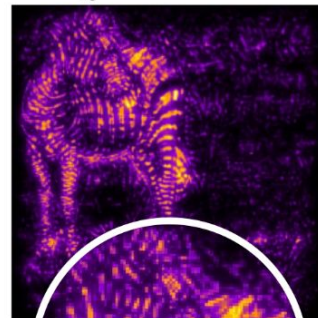
GRADIENT



EFFECTIVE



EQUIVARIANT



# DAVE Construction

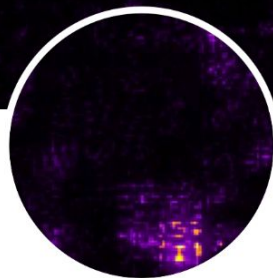
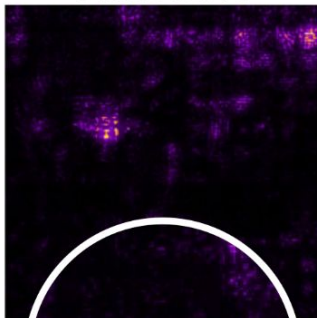
## STEP III: LOW-PASS FILTER

$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, \Sigma)} [\mathbf{W}_L^{\text{eq}}(\mathbf{X} + \epsilon)]$$

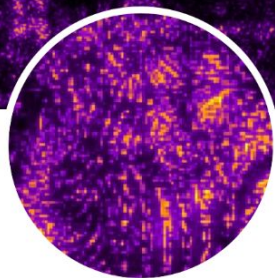
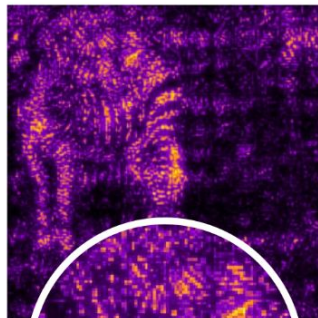
IMAGE



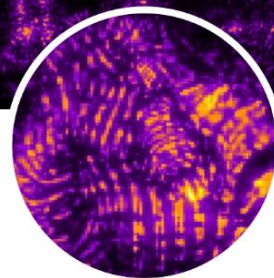
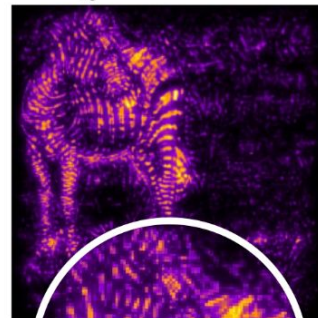
GRADIENT



EFFECTIVE



EQUIVARIANT



# DAVE Construction

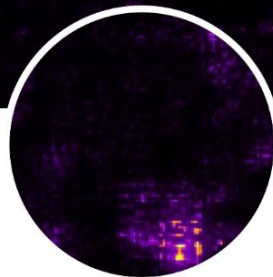
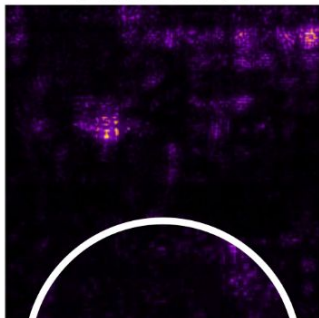
## STEP III: LOW-PASS FILTER

$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, \Sigma)} [\mathbf{W}_L^{\text{eq}}(\mathbf{X} + \epsilon)]$$

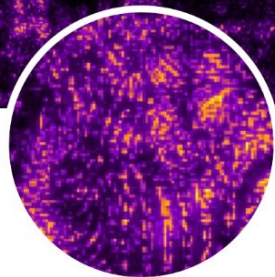
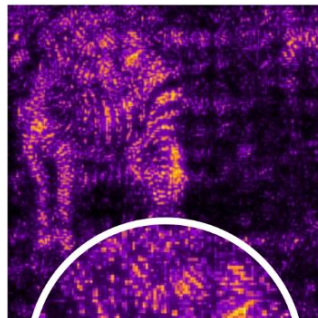
IMAGE



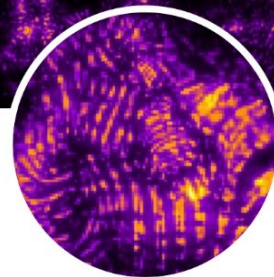
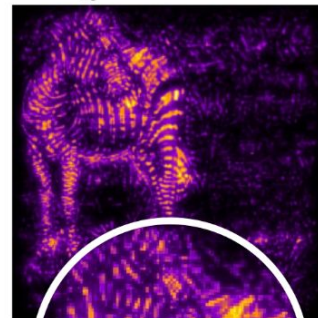
GRADIENT



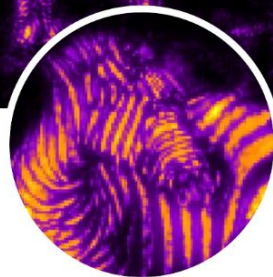
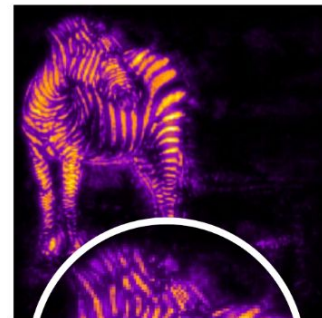
EFFECTIVE



EQUIVARIANT



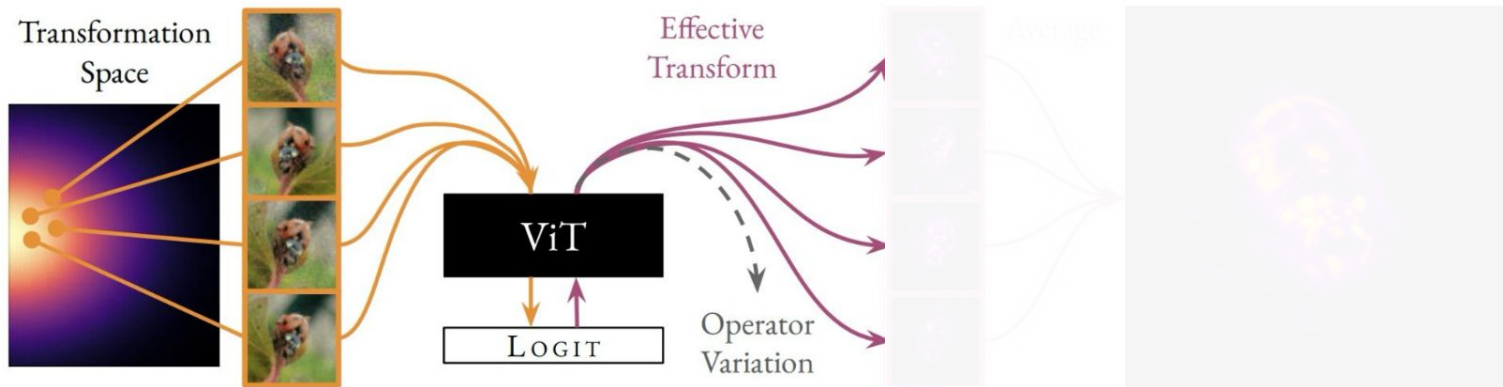
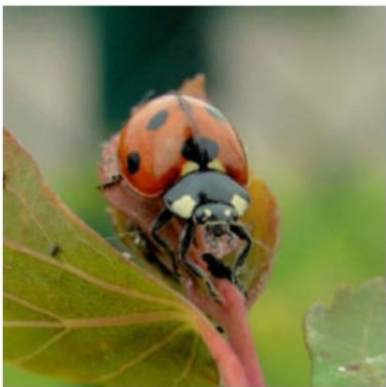
DAVE



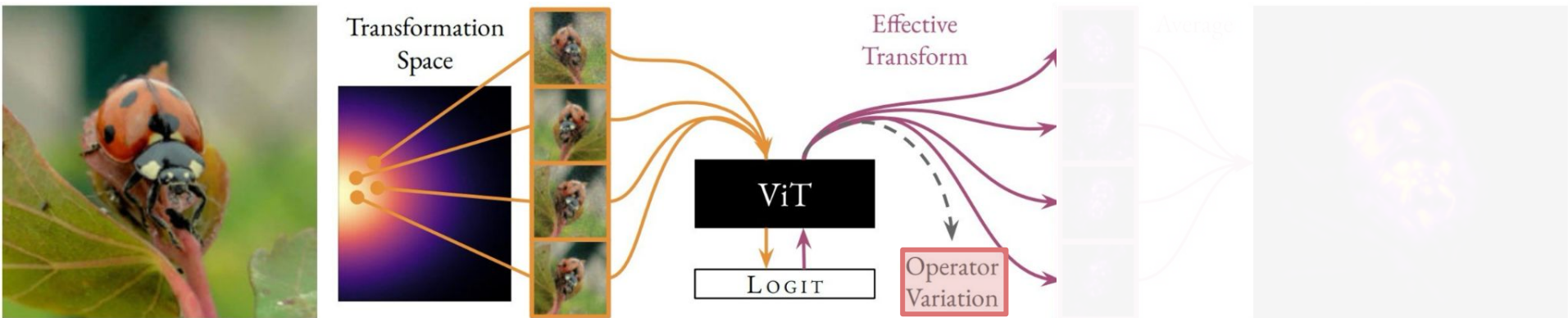
# DAVE Implementation



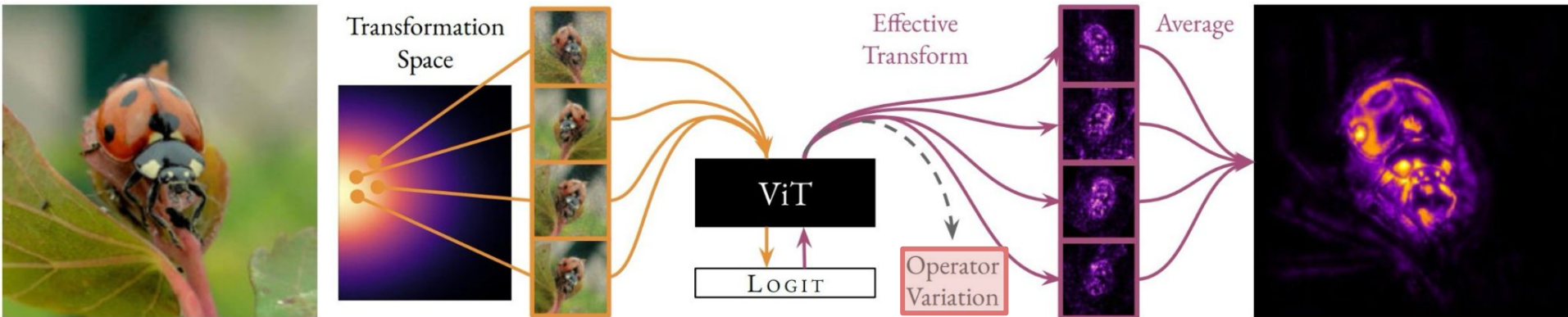
# DAVE Implementation



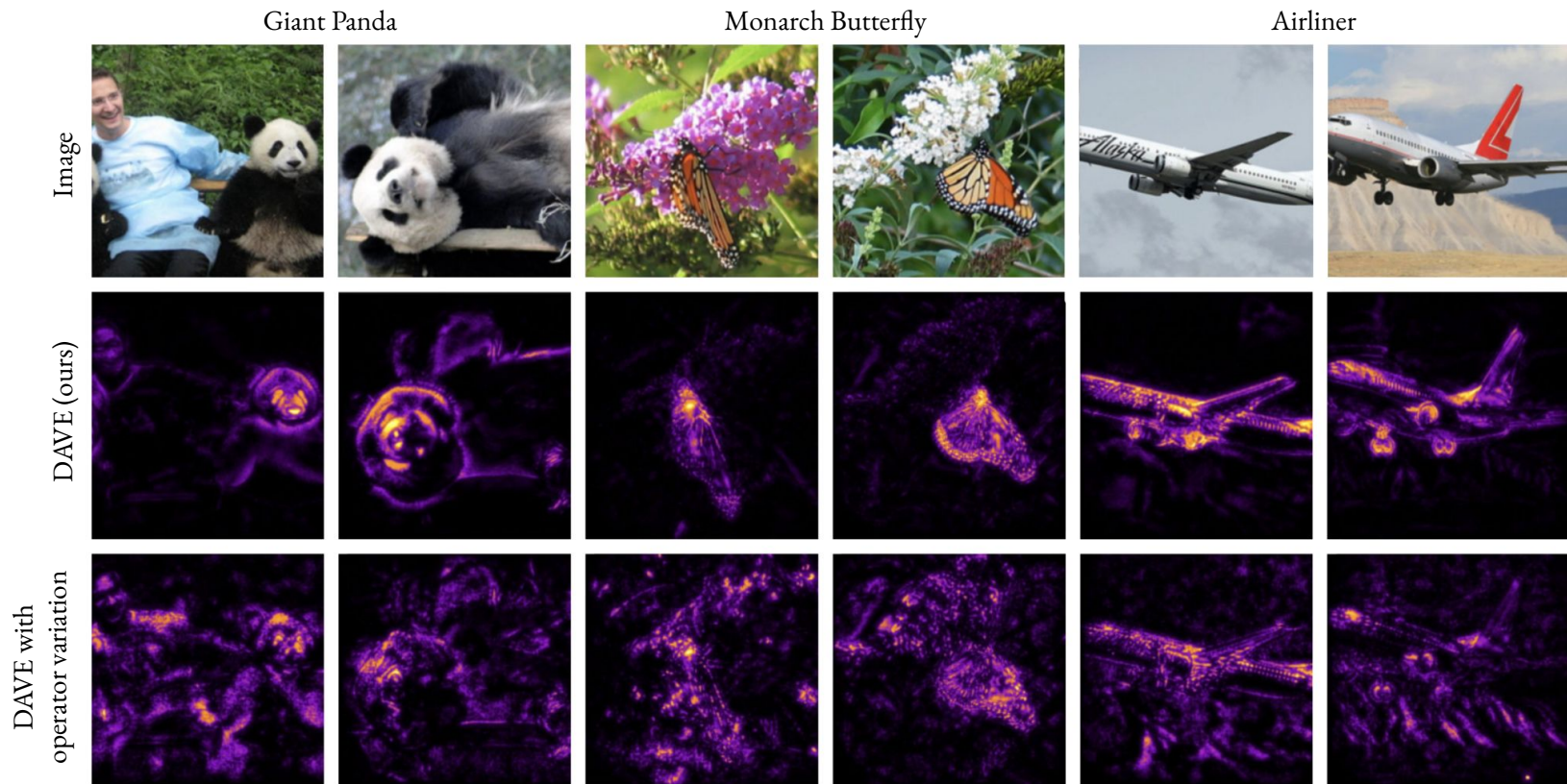
# DAVE Implementation



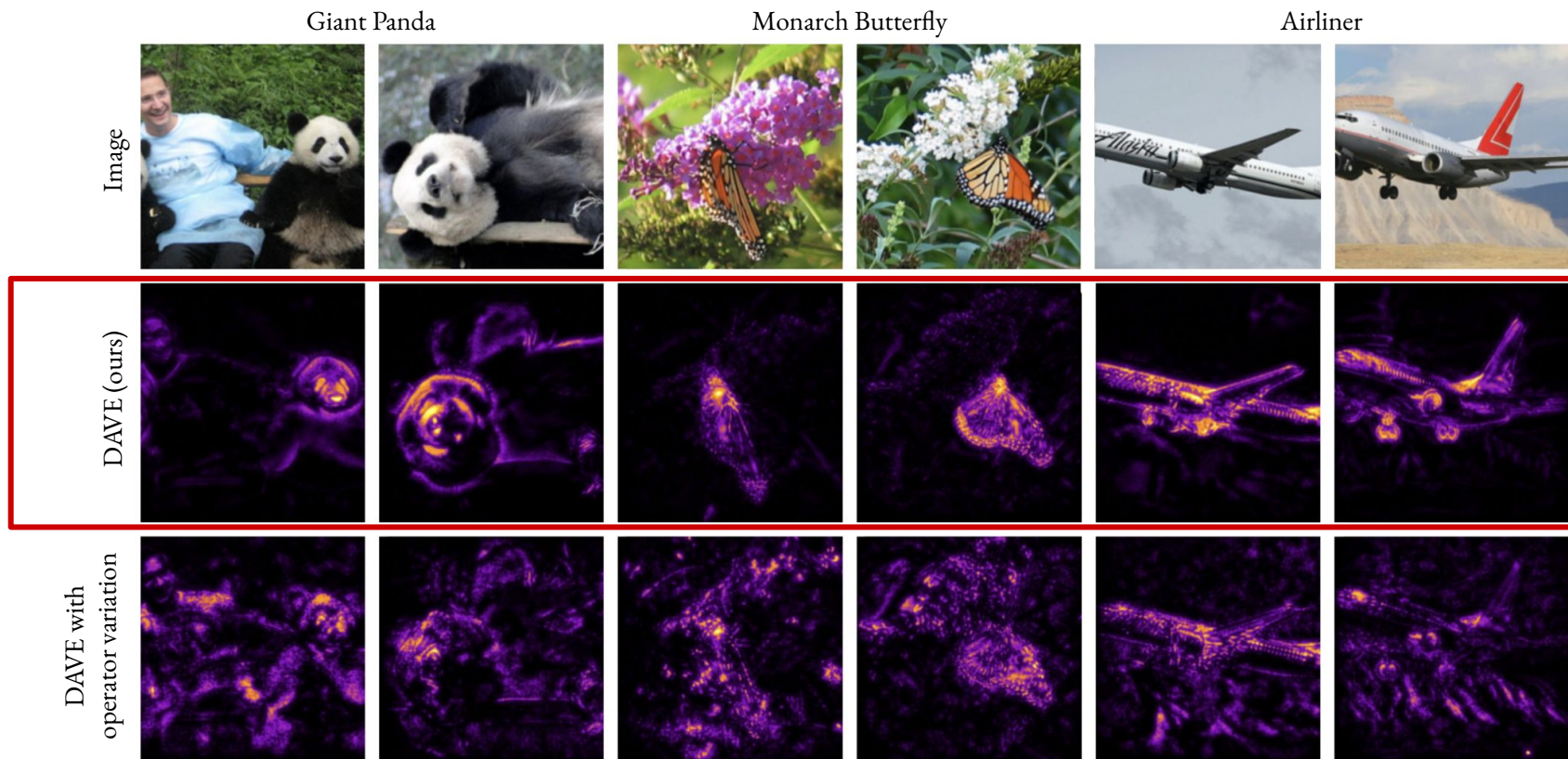
# DAVE Implementation



# Operator Variation Effect



# Operator Variation Effect



# Operator Variation Effect

Giant Panda

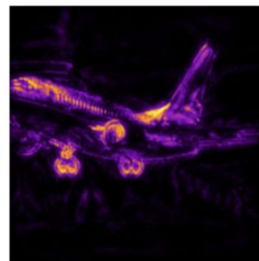
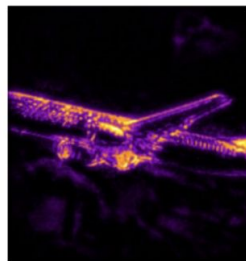
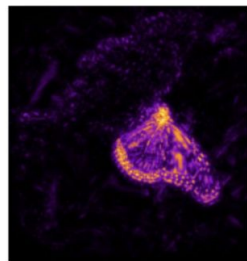
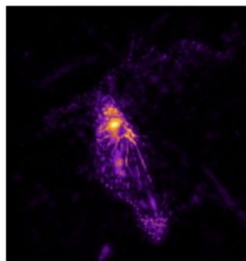
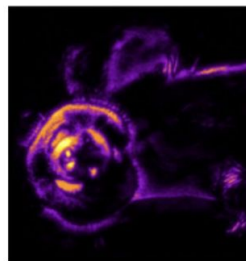
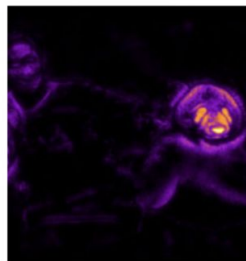
Monarch Butterfly

Airliner

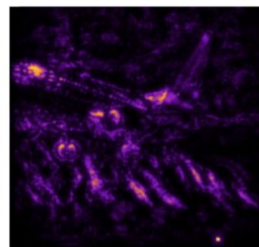
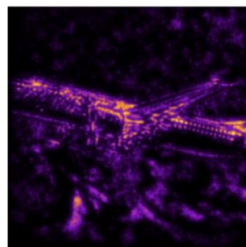
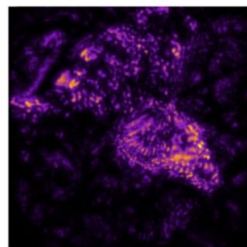
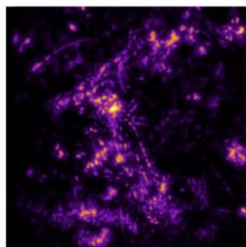
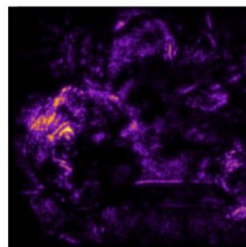
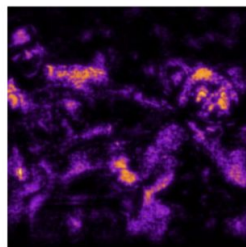
Image



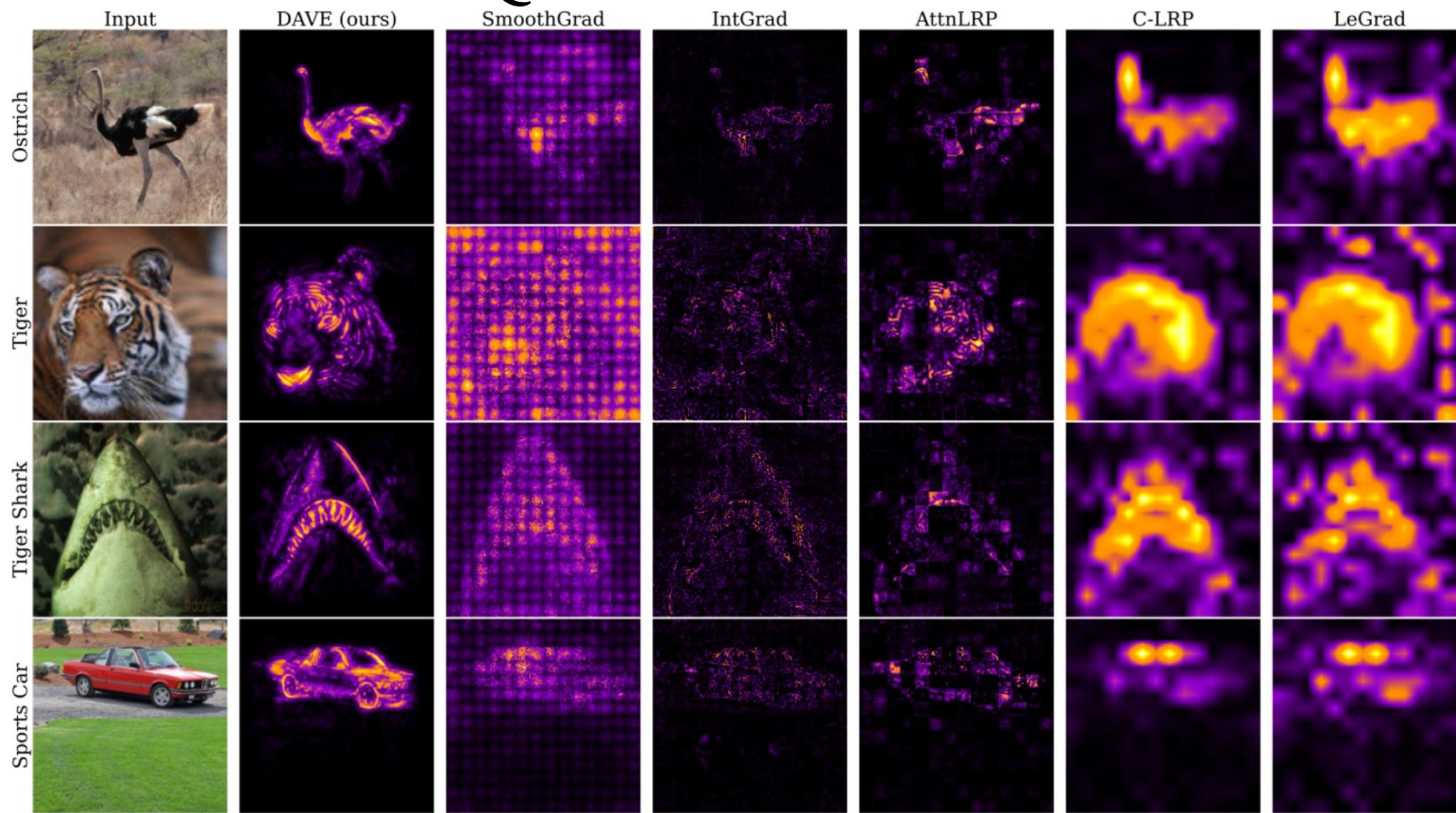
DAVE (ours)



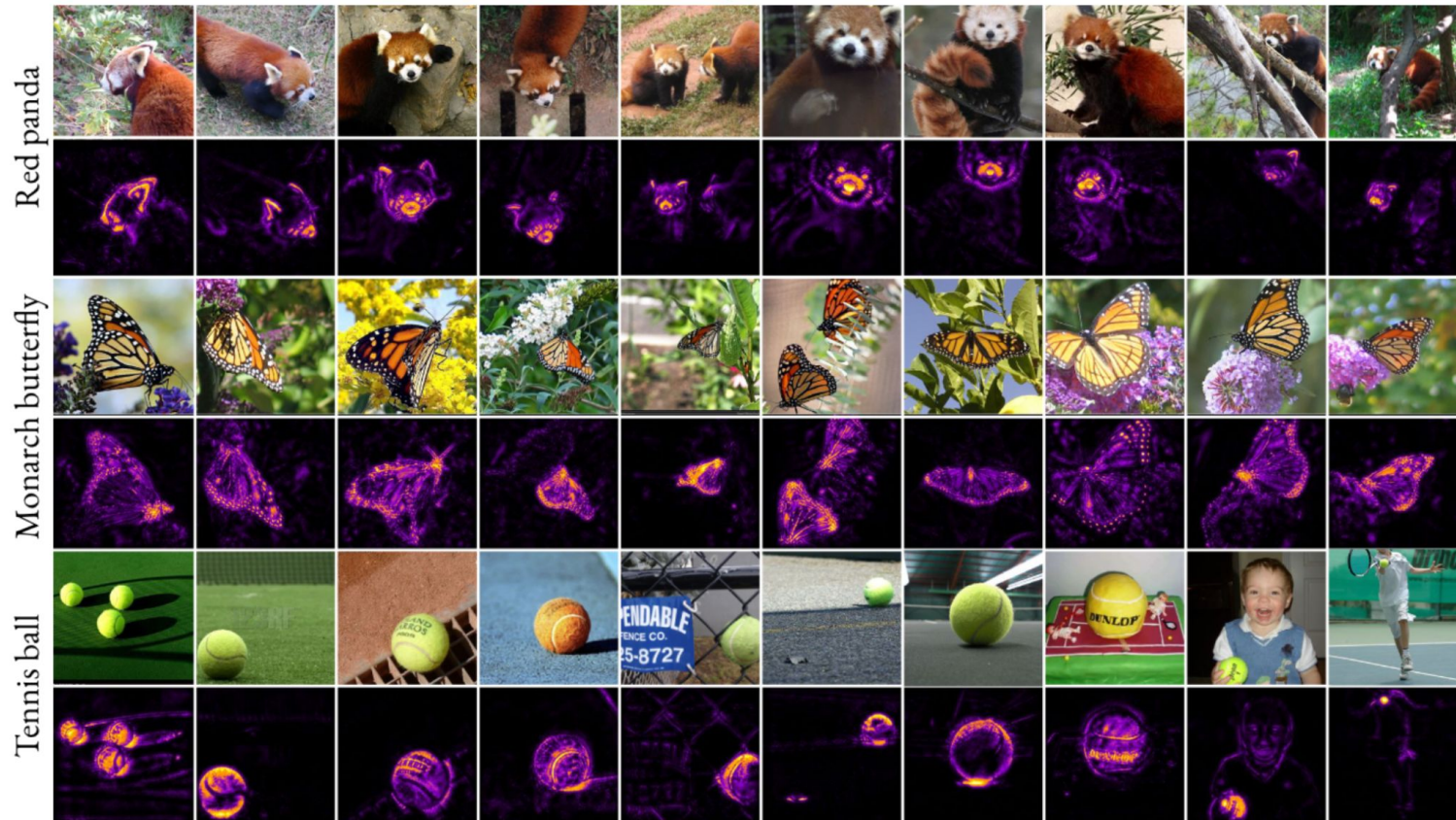
DAVE with operator variation



# Qualitative Results



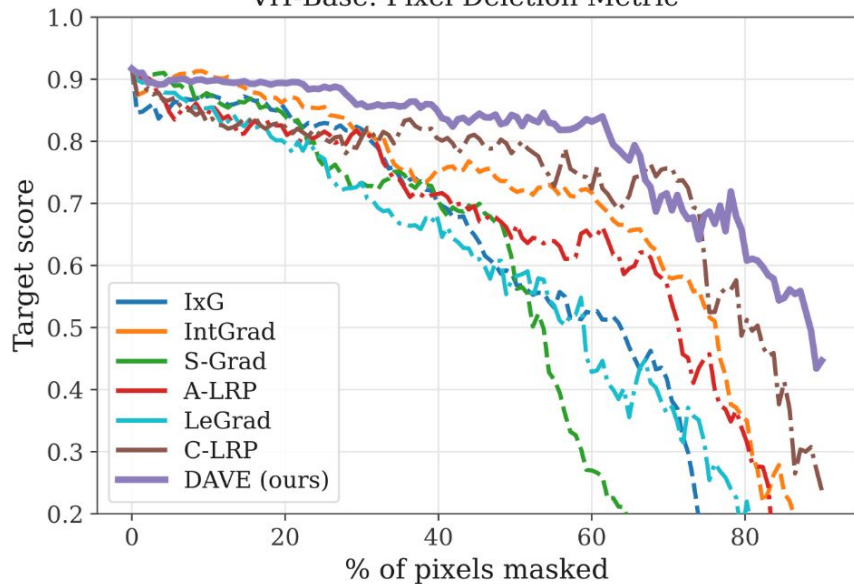
# Class Consistency



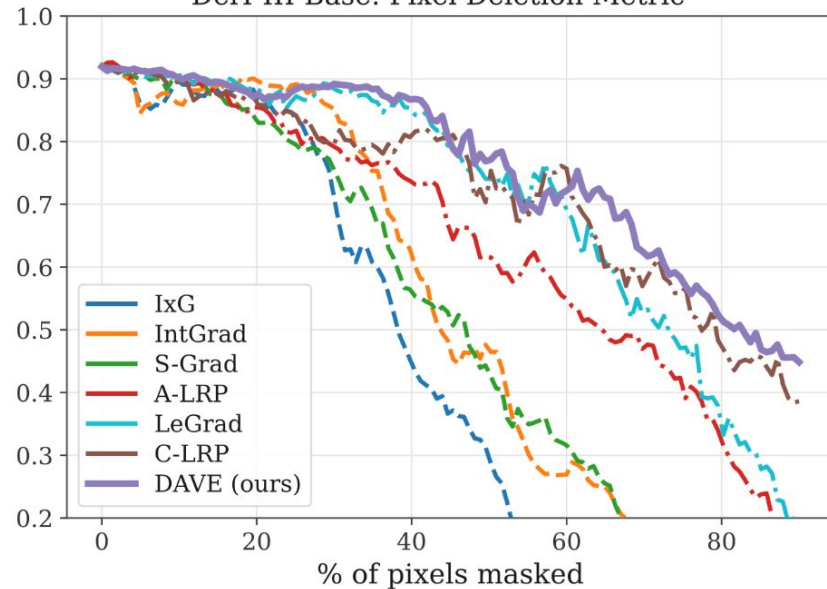
# Quantitative Results

## PIXEL DELETION

ViT-Base: Pixel Deletion Metric



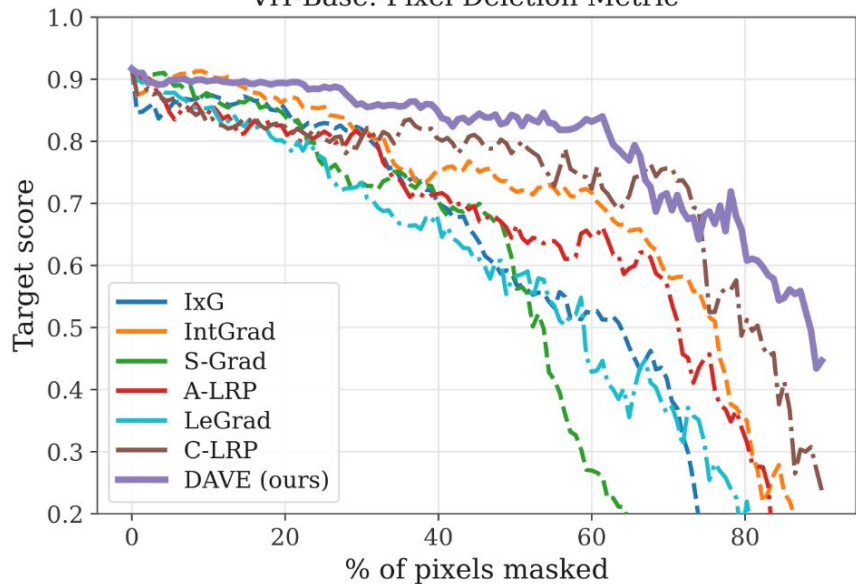
DeiT-III-Base: Pixel Deletion Metric



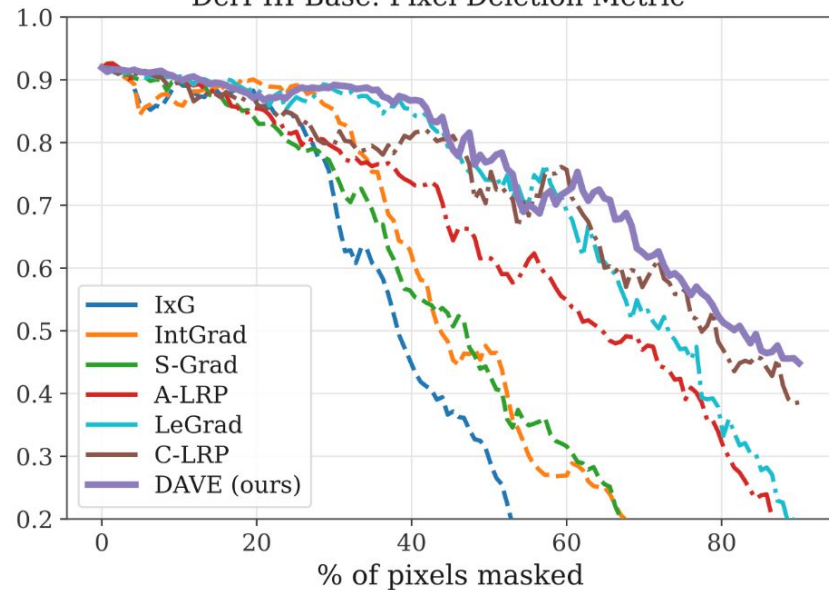
# Quantitative Results

## PIXEL DELETION

ViT-Base: Pixel Deletion Metric



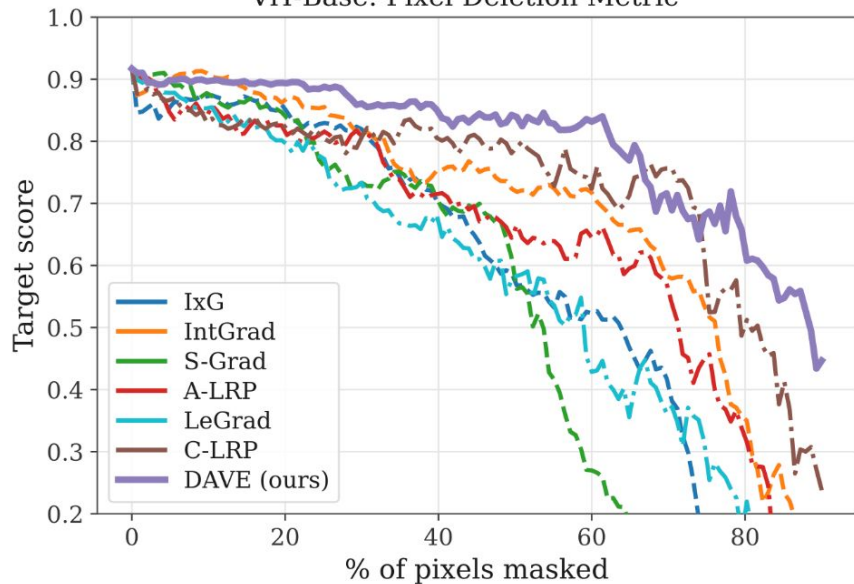
DeiT-III-Base: Pixel Deletion Metric



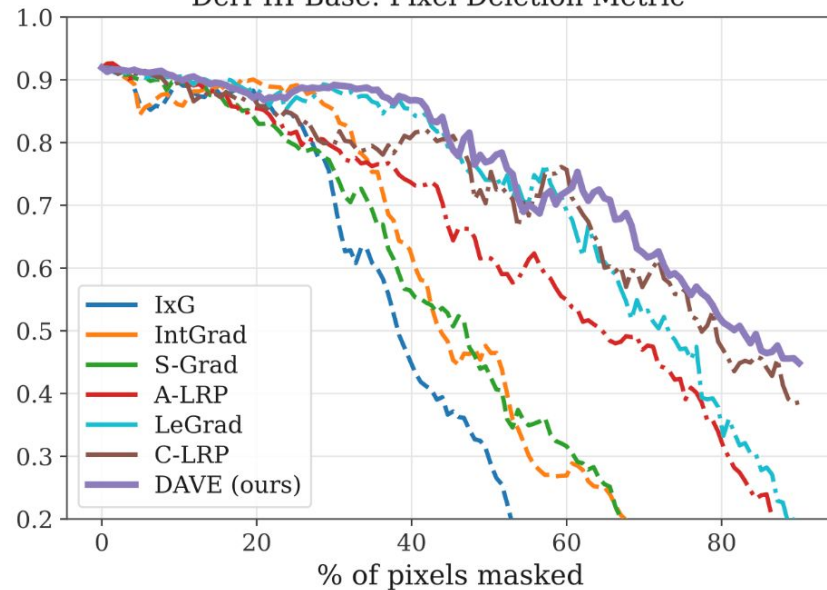
# Quantitative Results

## PIXEL DELETION

ViT-Base: Pixel Deletion Metric



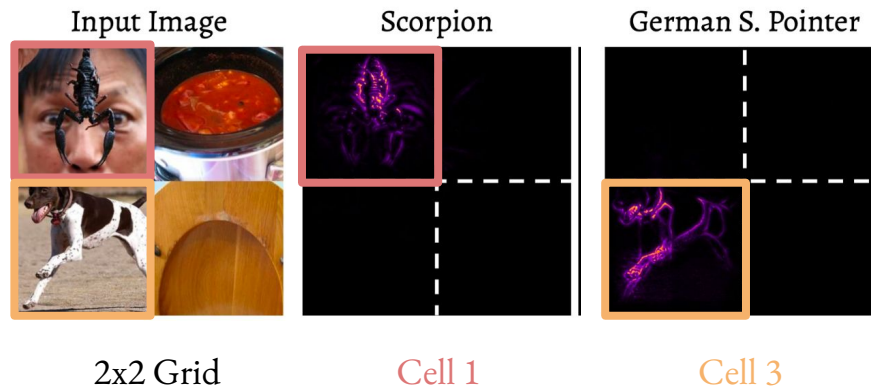
DeiT-III-Base: Pixel Deletion Metric



# Quantitative Results

## LOCALIZATION

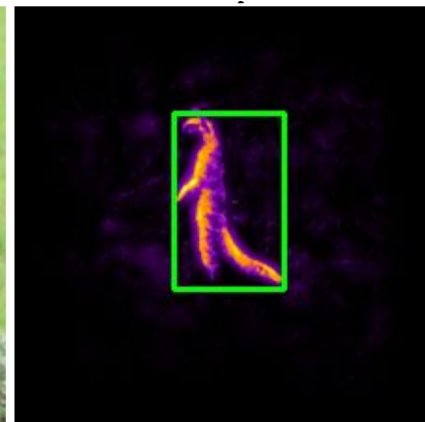
Method	GridPG (%)			
	ViT-B	DeiT-B	D-III-B	DINO-B
I×G	32.67	30.25	30.01	33.28
IntGrad	39.86	36.11	31.68	36.98
S-Grad	34.27	30.18	31.48	33.13
LeGrad	47.71	42.58	34.62	28.96
A-LRP	58.40	54.63	53.84	37.49
C-LRP	54.98	55.47	52.27	49.99
<b>DAVE (ours)</b>	<b>60.19</b>	<b>63.52</b>	<b>65.76</b>	<b>51.33</b>
$\Delta$	<b>+1.79</b>	<b>+8.05</b>	<b>+11.92</b>	<b>+1.35</b>



# Quantitative Results

## LOCALIZATION

Method	EnergyPG (%)			
	ViT-B	DeiT-B	D-III-B	DINO-B
I×G	55.33	62.72	67.32	69.98
IntGrad	58.51	64.22	68.12	72.15
S-Grad	56.02	60.78	68.10	70.93
LeGrad	80.06	77.83	77.54	82.26
A-LRP	60.75	68.16	77.65	75.98
C-LRP	<b>80.82</b>	79.62	81.94	81.56
DAVE (ours)	78.60	<b>82.23</b>	<b>82.43</b>	<b>83.38</b>
$\Delta$	-2.22	+2.61	+0.49	+1.12



# Takeaways

- ViT attributions suffer from **architecture-induced artifacts**.
- ViT gradients mix **effective transform** and **operator variation**.
- DAVE isolates **stable, transformation-consistent** signal from **effective transform**.
- DAVE achieves **state-of-the-art** pixel-level attribution quality.

# Takeaways

- ViT attributions suffer from **architecture-induced artifacts**.
- ViT gradients mix **effective transform** and **operator variation**.
- DAVE isolates **stable, transformation-consistent** signal from effective transform.
- DAVE achieves **state-of-the-art** pixel-level attribution quality.

# Takeaways

- ViT attributions suffer from **architecture-induced artifacts**.
- ViT gradients mix **effective transform** and **operator variation**.
- DAVE isolates **stable, transformation-consistent** signal from **effective transform**.
- DAVE achieves **state-of-the-art** pixel-level attribution quality.

# Takeaways

- ViT attributions suffer from **architecture-induced artifacts**.
- ViT gradients mix **effective transform** and **operator variation**.
- DAVE isolates **stable, transformation-consistent** signal from **effective transform**.
- DAVE achieves **state-of-the-art** pixel-level ViT attribution quality.