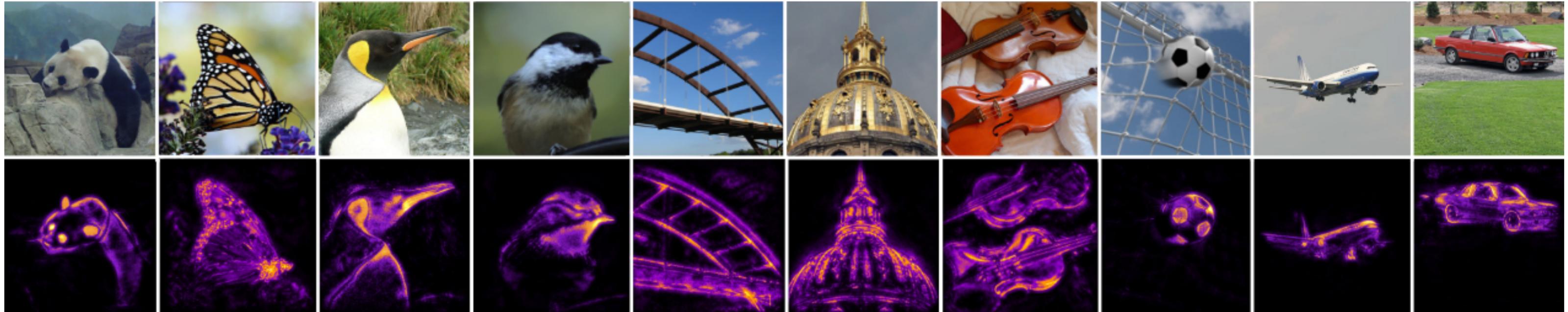




DAVE : Distribution-aware Attribution via ViT Gradient Decomposition

arxiv: <https://arxiv.org/abs/2602.06613>



Adam Wróbel
GMUM, Jagiellonian University



Siddhartha Gairola
MPI for Informatics



Jacek Tabor
GMUM, Jagiellonian University



Bernt Schiele
MPI for Informatics



Bartosz Zieliński
GMUM, Jagiellonian University

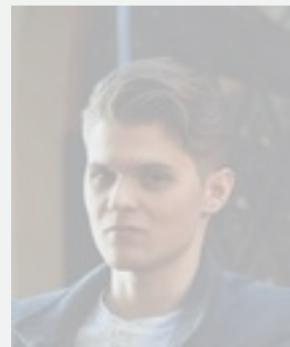


Dawid Rymarczyk
GMUM, Jagiellonian University



DAVE: Distribution-aware Attribution via ViT Gradient Decomposition

arxiv: <https://arxiv.org/abs/2602.06613>



Adam Wróbel
GMUM, Jagiellonian University



Siddhartha Gairola
MPI for Informatics



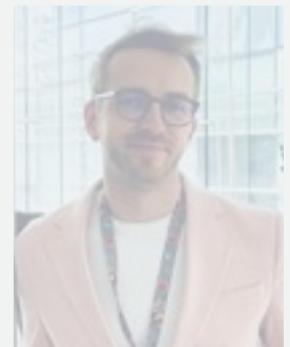
Jacek Tabor
GMUM, Jagiellonian University



Bernt Schiele
MPI for Informatics



Bartosz Zieliński
GMUM, Jagiellonian University



Dawid Rymarczyk
GMUM, Jagiellonian University

Issues with existing Attribution Methods for Vision Transformers

Vision Transformers [1] or ViTs are highly prevalent in computer vision.

However getting good explanations still remains a challenge:

◆ gradient- and attention-based explanations exhibit structured, architecture-induced artifacts

Issues with existing Attribution Methods for Vision Transformers

Vision Transformers [1] or ViTs are highly prevalent in computer vision.

However getting good explanations still remains a challenge:

- ◆ gradient- and attention-based explanations exhibit structured, architecture-induced artifacts
- ◆ leading either to unstable pixel-level attributions or to coarse patch-level explanations
- ◆ thus lack fine-grained visual evidence

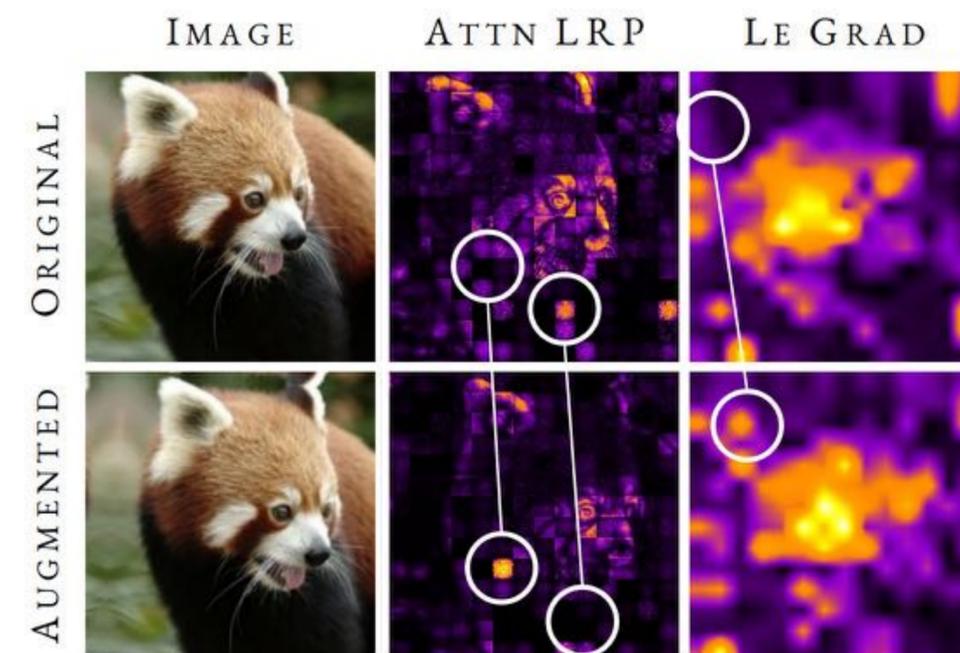
Issues with existing Attribution Methods for Vision Transformers

Vision Transformers [1] or ViTs are highly prevalent in computer vision.

However getting good explanations still remains a challenge:

- ◆ gradient- and attention-based explanations exhibit structured, architecture-induced artifacts
- ◆ leading either to unstable pixel-level attributions or to coarse patch-level explanations
- ◆ thus lack fine-grained visual evidence

Under small augmentations
(5° rotation, 20px horizontal and 8px vertical shift)



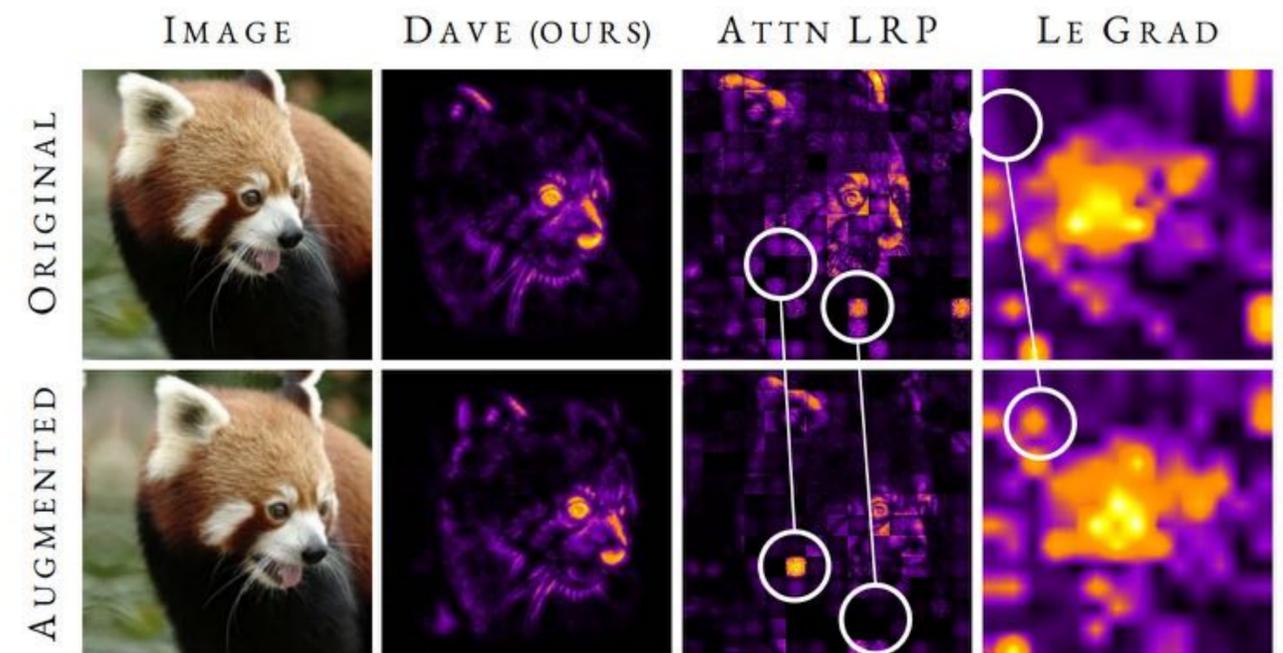
Issues with existing Attribution Methods for Vision Transformers

Vision Transformers [1] or ViTs are highly prevalent in computer vision.

However getting good explanations still remains a challenge:

- ◆ gradient- and attention-based explanations exhibit structured, architecture-induced artifacts
- ◆ leading either to unstable pixel-level attributions or to coarse patch-level explanations
- ◆ thus lack fine-grained visual evidence

Under small augmentations
(5° rotation, 20px horizontal and 8px vertical shift)



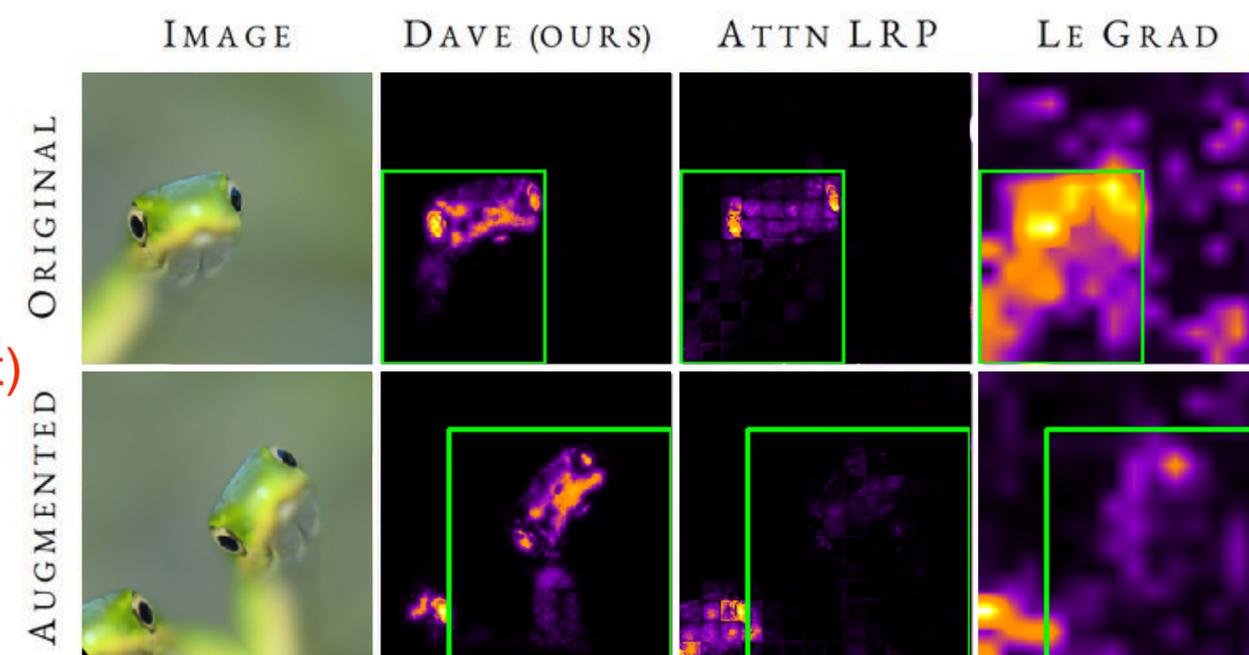
Issues with existing Attribution Methods for Vision Transformers

Vision Transformers [1] or ViTs are highly prevalent in computer vision.

However getting good explanations still remains a challenge:

- ◆ gradient- and attention-based explanations exhibit structured, architecture-induced artifacts
- ◆ leading either to unstable pixel-level attributions or to coarse patch-level explanations
- ◆ thus lack fine-grained visual evidence

Under stronger augmentations
(-40° rotation, 50px horizontal and -30px vertical shift)



What does DAVE solve?

We fix a persistent issue in ViT explainability: **unstable, artifact-heavy pixel attributions.**

DAVE yields:

✓ **fine-grained pixel-level maps**

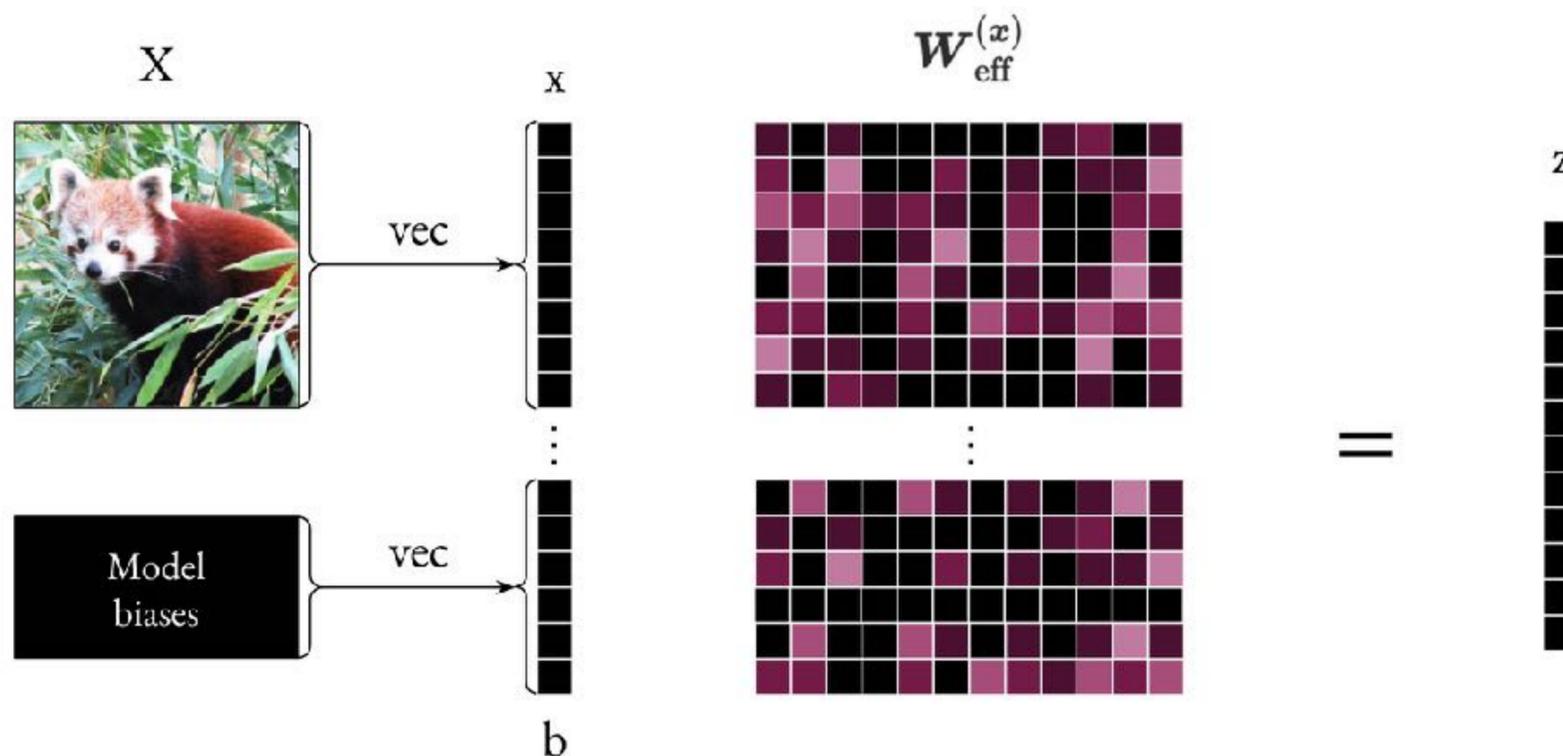
✓ **without patch-grid saliency.**



DAVE - The Method

DAVE interprets an explanation (or attribution) as the **stable** and **locally equivariant** effective transformation that a ViT applies to its input.

ViT is a Dynamic Linear Matrix $\mathbf{z} = \text{ViT}(\mathbf{X}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} [\mathbf{x} | \mathbf{b}]$



- If every base layer is of a form:

$$\text{vec}(\mathbf{Z}) = \mathbf{W}(\mathbf{x})\mathbf{x} + \mathbf{b}$$

- Then the whole ViT is dynamic linear:

$$\text{ViT}(\mathbf{X}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} [\mathbf{x} | \mathbf{b}]$$

DAVE - The Method

Model each ViT layer as an input-dependent linear operator $L(X)$ applied to X .

(self-attention, layer-norm, residual connections, mlp blocks, activation functions)

Then the input-gradient decomposes into:

(1) effective transformation $L(X)$

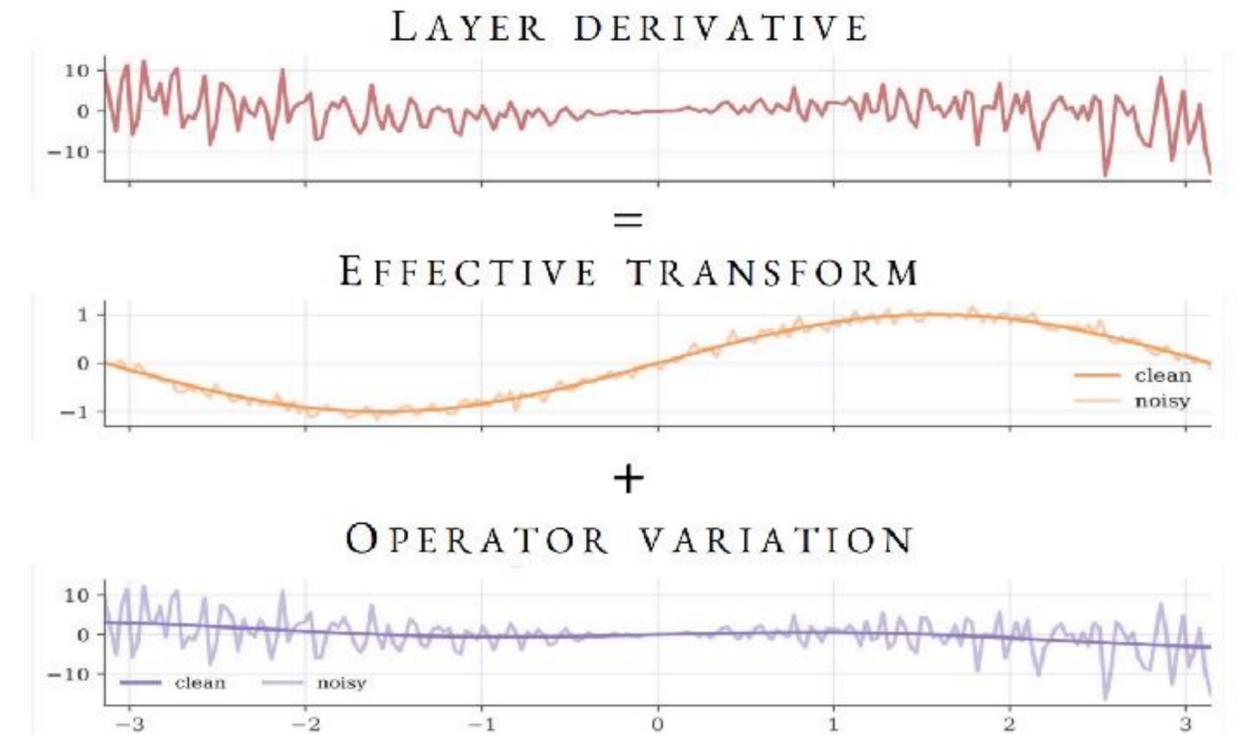
(2) operator variation (how L changes w.r.t. X)

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

DAVE - The Method

! Why gradients break ?

Operator variation can amplify tiny perturbations
 → high-frequency junk in attributions.



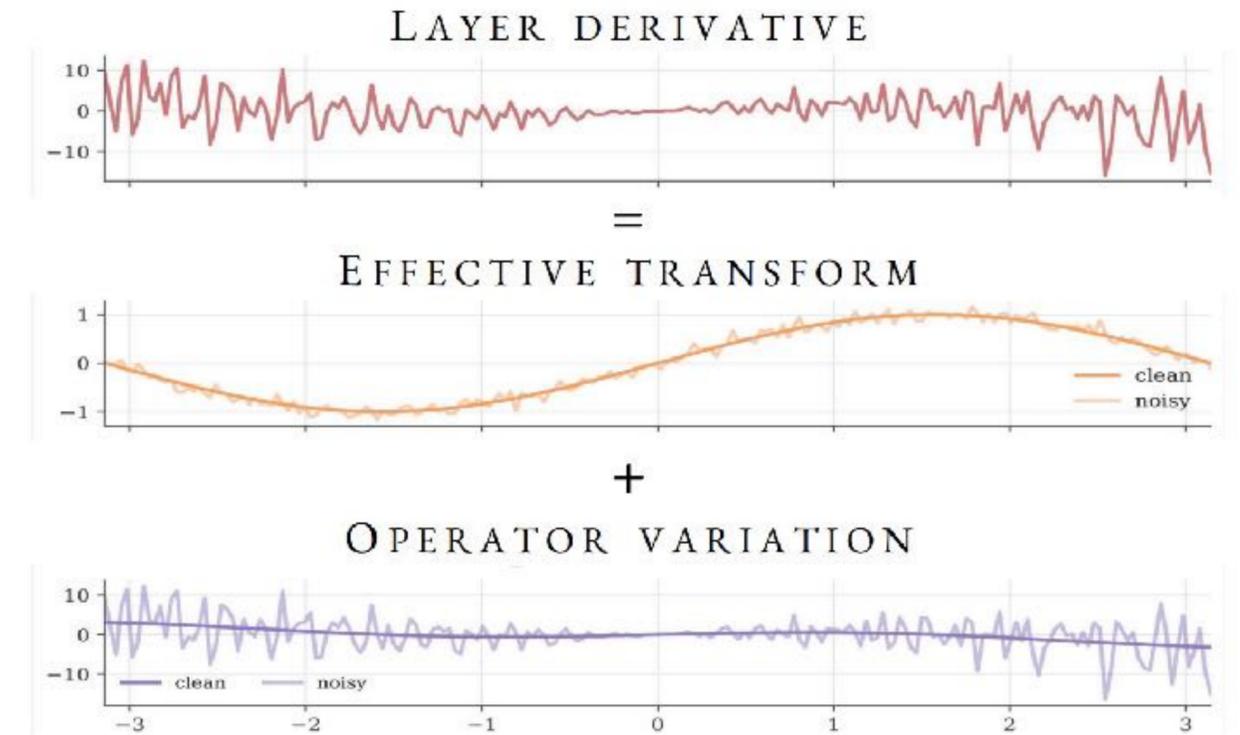
$$\underbrace{D_{\mathbf{X}}F}_{\text{layer derivative}} = \underbrace{L(\mathbf{X})}_{\text{effective transformation}} + \underbrace{\left((D_{\mathbf{X}}L(\mathbf{X}))(\cdot) \right) \mathbf{X}}_{\text{operator variation}}$$

DAVE - The Method

! Why gradients break ?

Operator variation can amplify tiny perturbations
→ high-frequency junk in attributions.

DAVE drops this term and keeps the effective transformation as a cleaner attribution operator.



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

The 'operator variation' term in the equation is crossed out with a large red 'X'.

DAVE - The Method

Effective weights and how to find them efficiently

- **Detach** all dynamic matrices on the forward pass (attention, GELU multiplier, etc.)
- Compute “detached” **backward pass**
- Like: CoDA, B-cos

$$f(\mathbf{x}) = \mathbf{W}(\mathbf{x})\mathbf{x}$$

$$J_f(\mathbf{x}) = \underbrace{\mathbf{W}(\mathbf{x})}_{\text{Our ViT matrix}} + \underbrace{\sum_{i=1}^d x_i \frac{\partial \mathbf{W}}{\partial x_i}(\mathbf{x})}_{\text{JVP}}$$

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

DAVE - The Method

We base our method on a **working hypothesis**: for a well-trained model

$$C(\mathbf{x}) = C^*(\mathbf{x}) + n_\tau(\mathbf{x}) + n_\epsilon(\mathbf{x})$$

The diagram illustrates the decomposition of the cost function $C(\mathbf{x})$ into four components. Arrows point from the labels below to the corresponding terms in the equation above:

- Dynamic-linear representation** points to $C^*(\mathbf{x})$.
- Attribution** points to $n_\tau(\mathbf{x})$.
- Spatial-related global noise** points to $n_\tau(\mathbf{x})$.
- high-frequency local noise** points to $n_\epsilon(\mathbf{x})$.

where:

$$C(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

DAVE - The Method

We base our method on a **working hypothesis**: for a well-trained model

$$C(\mathbf{x}) = C^*(\mathbf{x}) + n_\tau(\mathbf{x}) + n_\epsilon(\mathbf{x})$$

The diagram illustrates the decomposition of the cost function $C(\mathbf{x})$ into three components. The equation $C(\mathbf{x}) = C^*(\mathbf{x}) + n_\tau(\mathbf{x}) + n_\epsilon(\mathbf{x})$ is shown at the top. Below it, four labels are connected to the terms in the equation by arrows: 'Dynamic-linear representation' points to $C(\mathbf{x})$, 'Attribution' points to $C^*(\mathbf{x})$ (which is enclosed in a dashed green box), 'Spatial-related global noise' points to $n_\tau(\mathbf{x})$, and 'high-frequency local noise' points to $n_\epsilon(\mathbf{x})$.

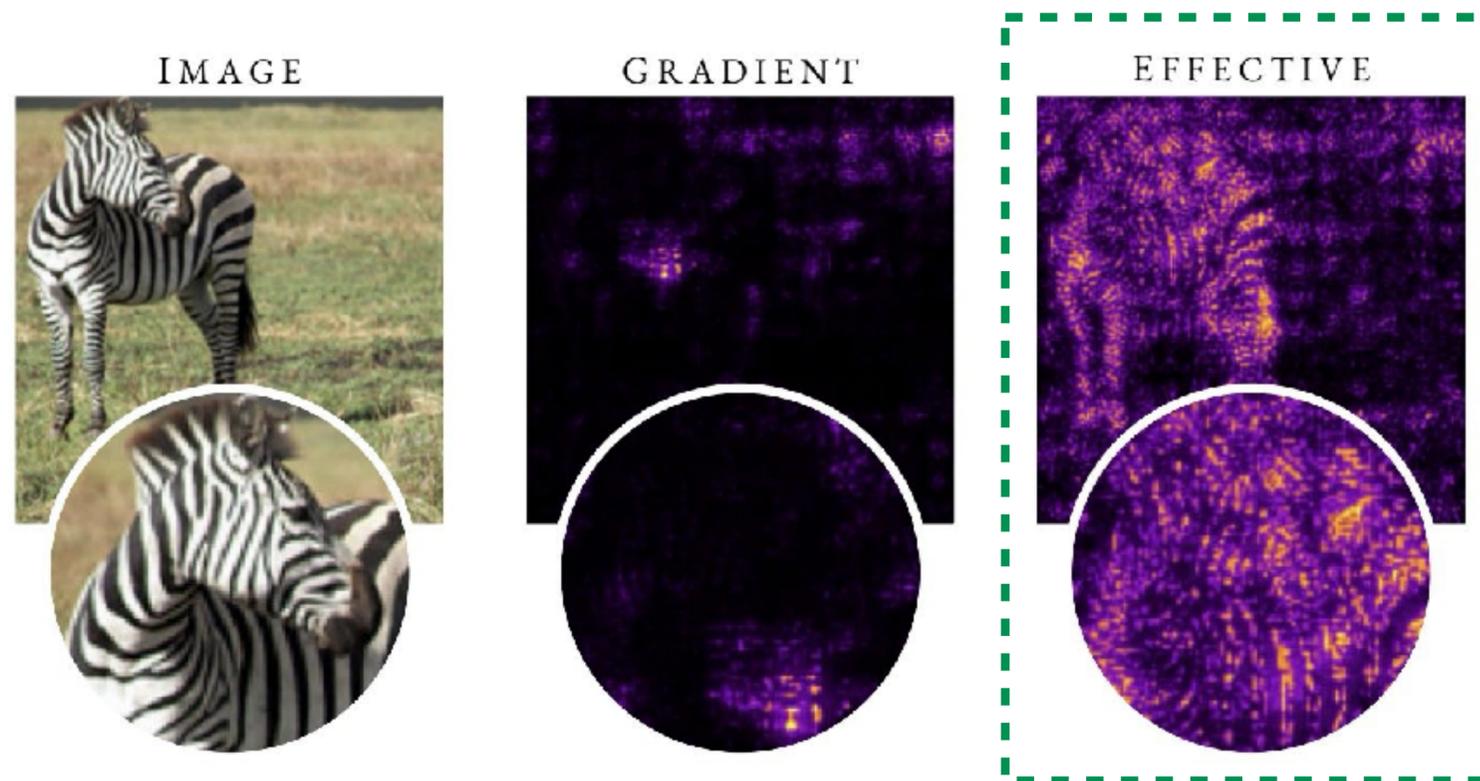
where:

$$C(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

DAVE - The Method

Remove grid artifacts

Even the effective transformation can carry architecture-induced grid patterns.

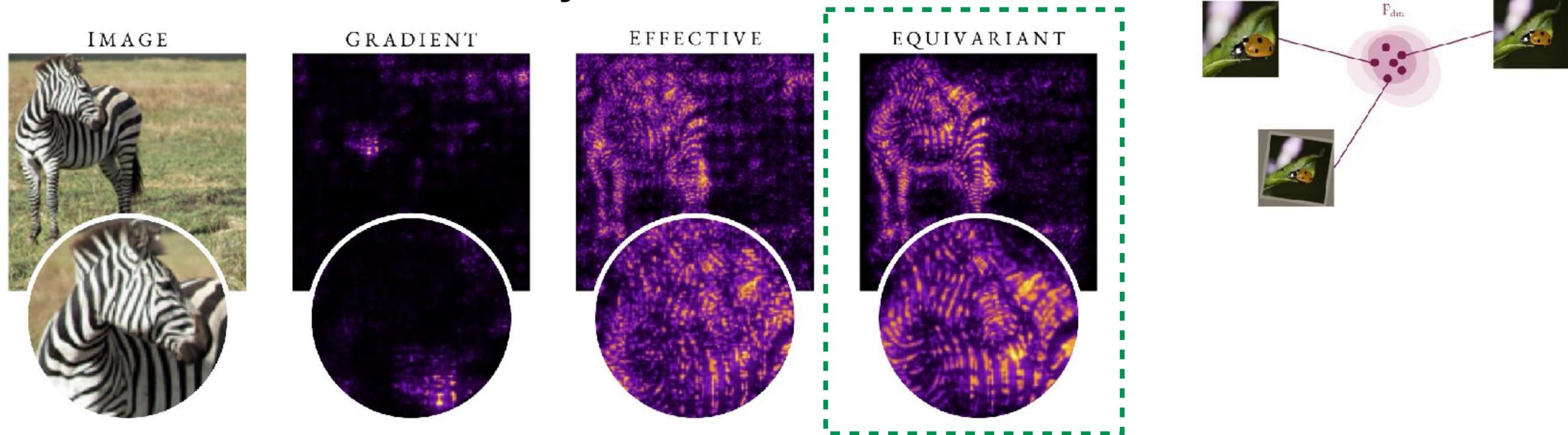


DAVE - The Method

Remove grid artifacts

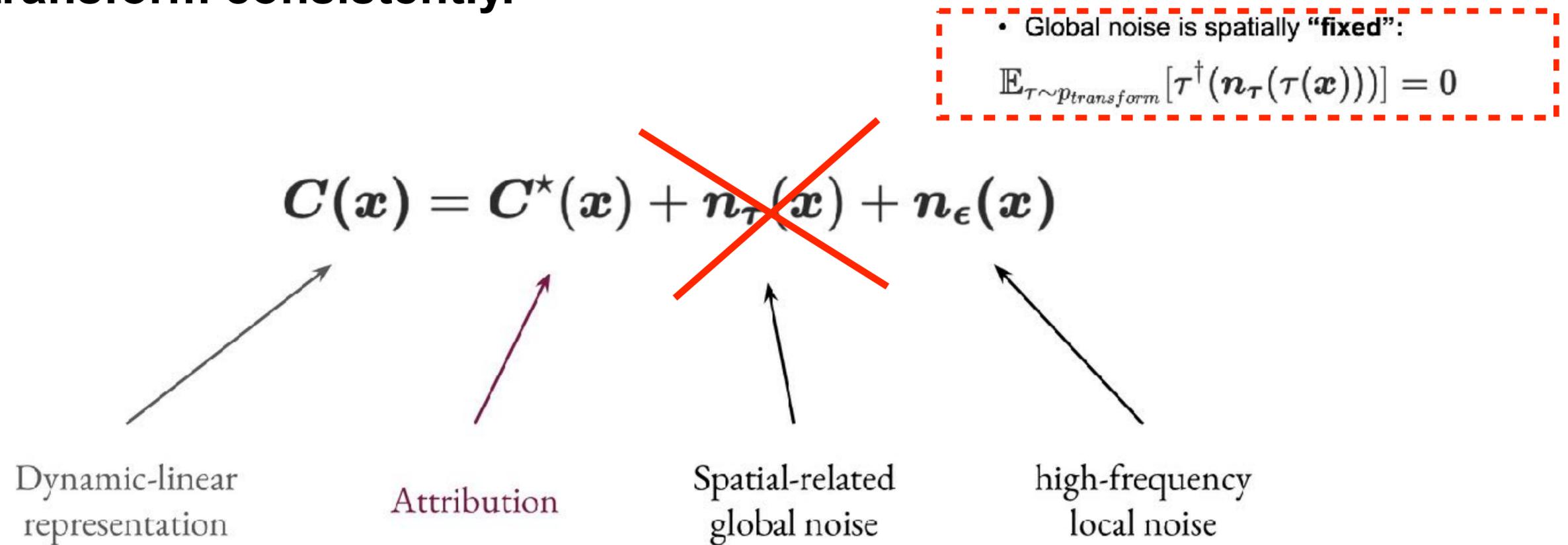
Even the effective transformation can carry architecture-induced grid patterns.

DAVE filters them by enforcing **local equivariance: under small spatial transforms, the attribution must transform consistently.**



DAVE - The Method

DAVE filters them by enforcing **local equivariance**: under small spatial transforms, the **attribution must transform consistently**.



where:

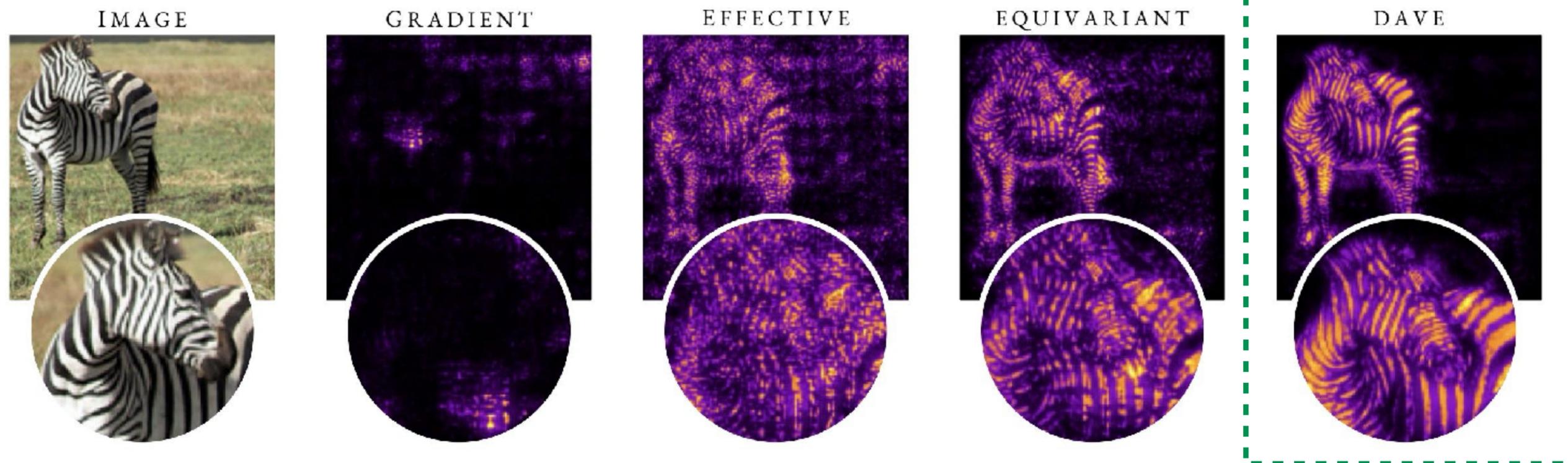
$$C(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

DAVE - The Method

■ Stabilize (low-pass)

DAVE adds low-pass filtering by **averaging** the equivariant effective transformation under **small input perturbations** (Gaussian smoothing in expectation).

This removes components unstable to tiny input changes. (See below fig., last column.)



DAVE - The Method

DAVE adds low-pass filtering by **averaging** the equivariant effective transformation under **small input perturbations** (Gaussian smoothing in expectation).

removed by low-pass filtering

$$C(\mathbf{x}) = C^*(\mathbf{x}) + n_\tau(\mathbf{x}) + n_\epsilon(\mathbf{x})$$

Dynamic-linear representation

Attribution

Spatial-related global noise

high-frequency local noise

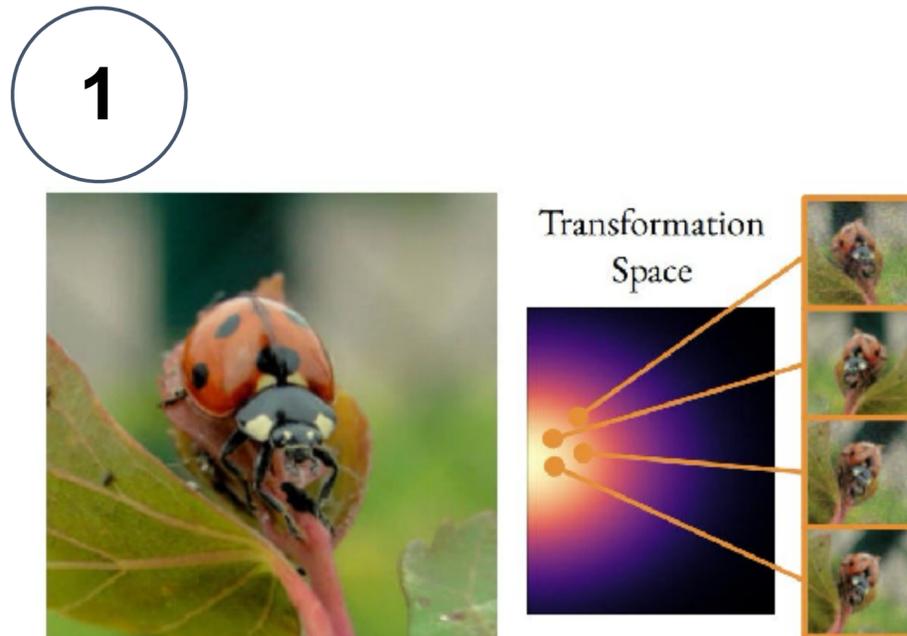
where:

$$C(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

DAVE - The Method

🔧 The pipeline

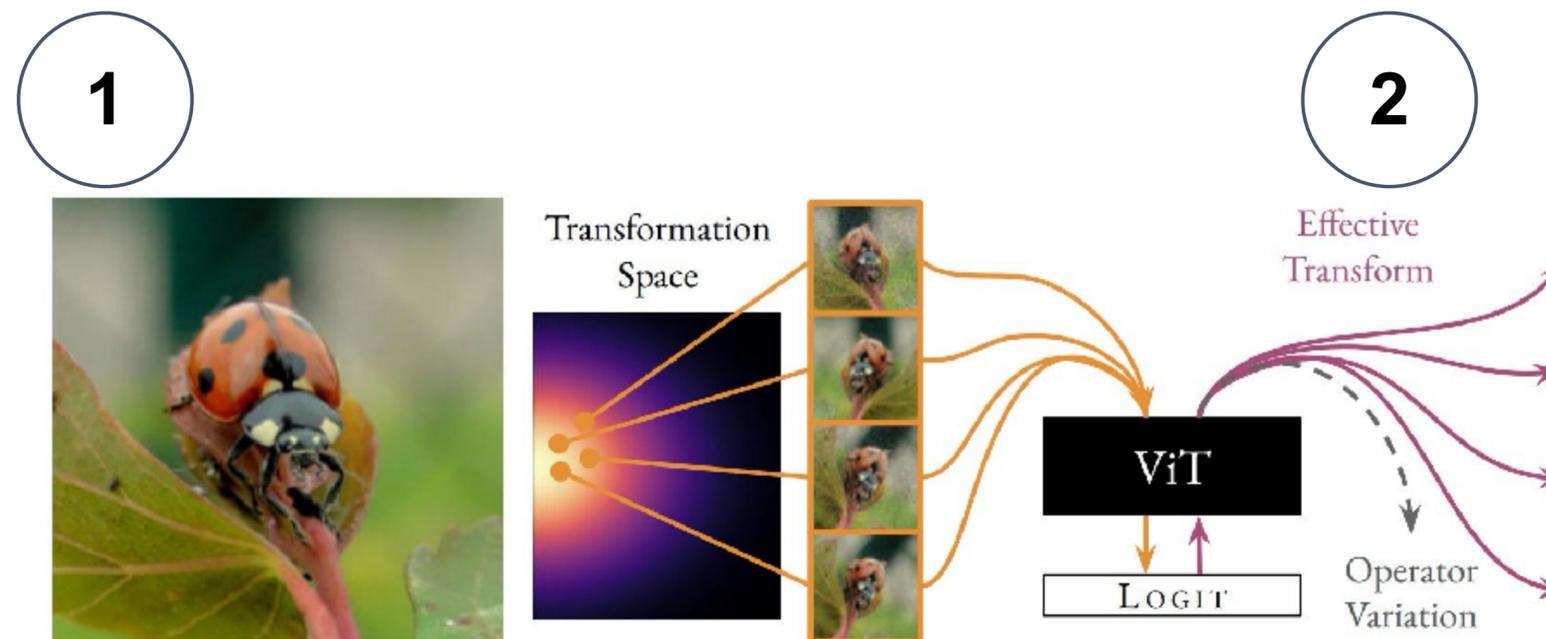
1. Sample small spatial transforms + noise,



DAVE - The Method

🔧 The pipeline

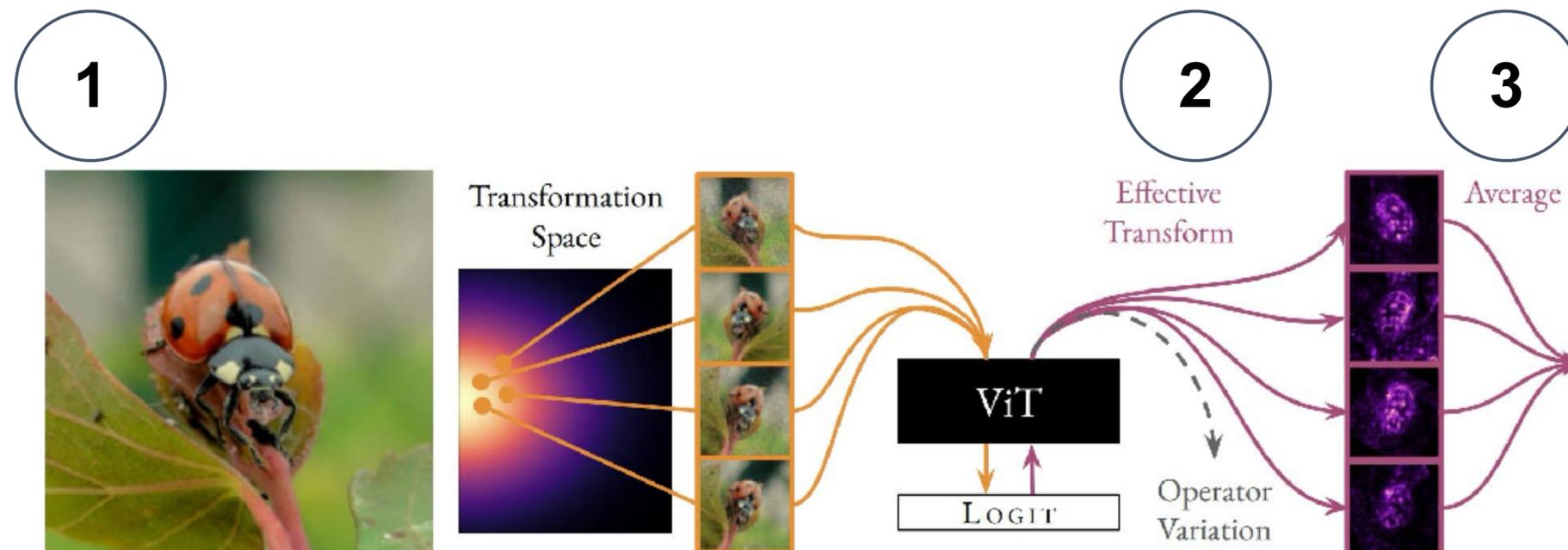
1. Sample small spatial transforms + noise,
2. compute effective transformation (conditioned forward blocks gradients through conditioning),



DAVE - The Method

🔧 The pipeline

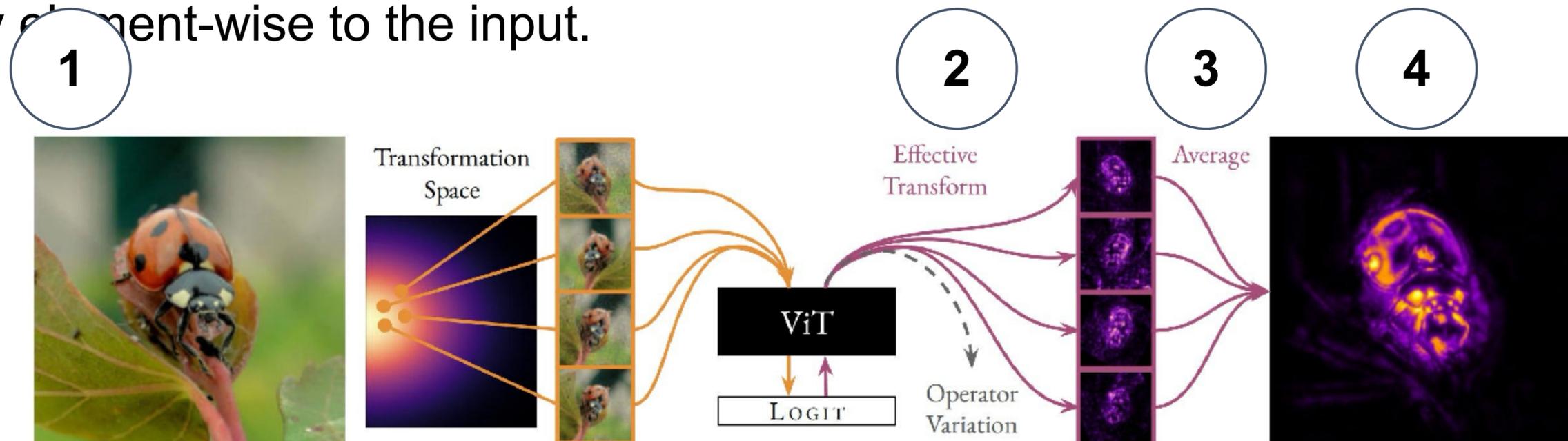
1. Sample small spatial transforms + noise,
2. compute effective transformation (conditioned forward blocks gradients through conditioning),
3. inverse-transform & average,



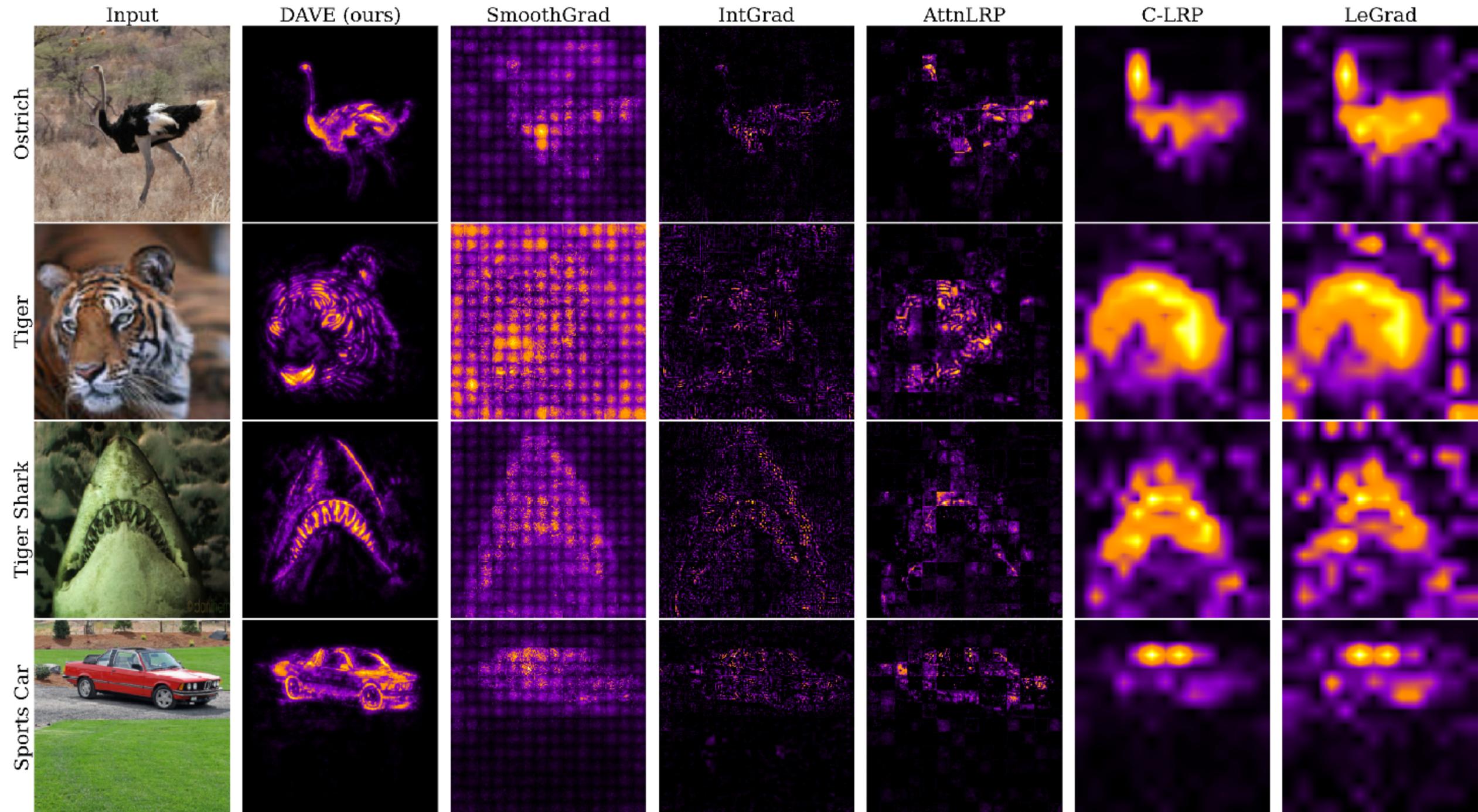
DAVE - The Method

🔧 The pipeline

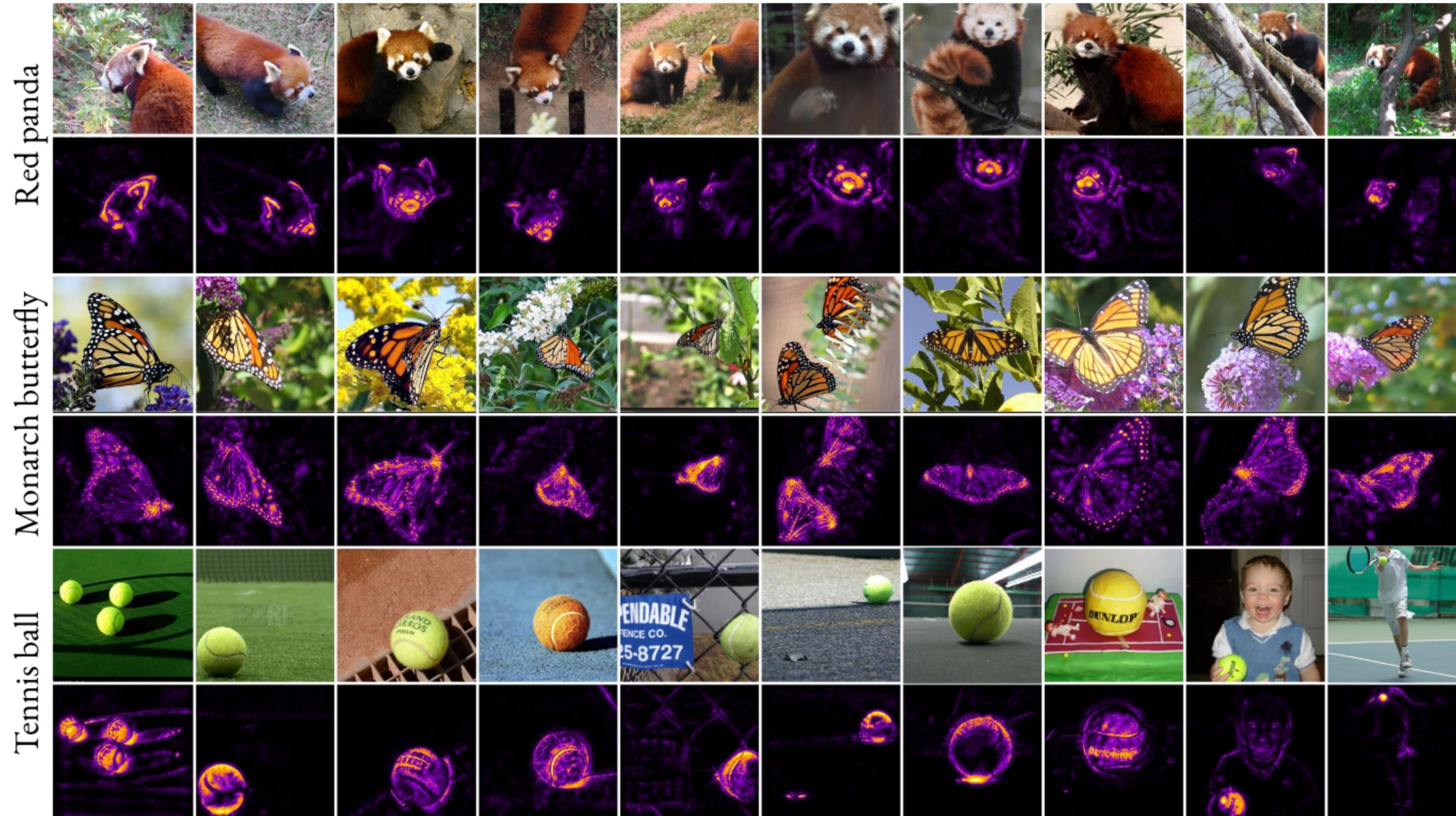
1. Sample small spatial transforms + noise,
2. compute effective transformation (conditioned forward blocks gradients through conditioning),
3. inverse-transform & average,
4. then apply element-wise to the input.



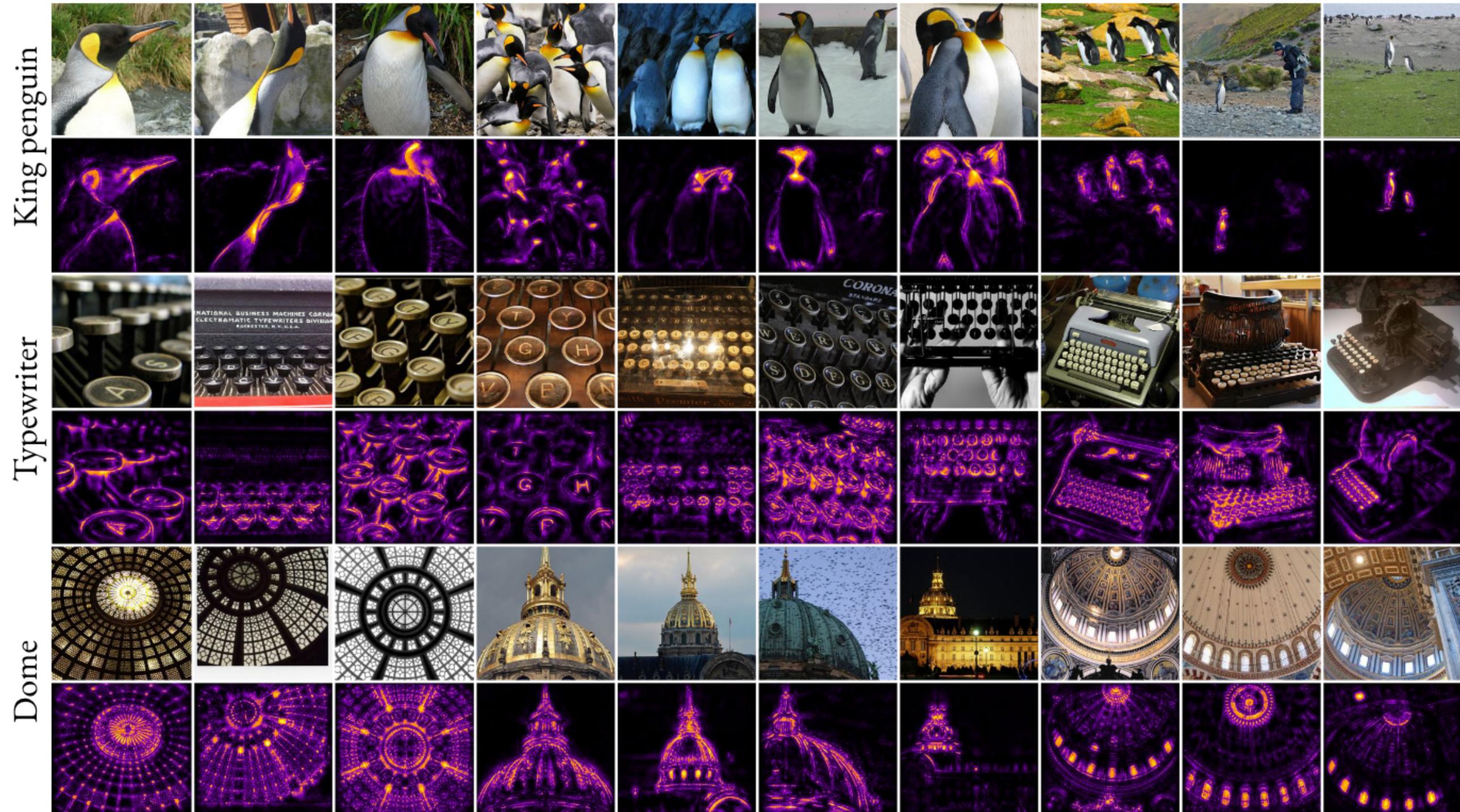
Qualitative Results: DAVE vs existing popular methods



Qualitative Results: class consistency of DAVE attributions



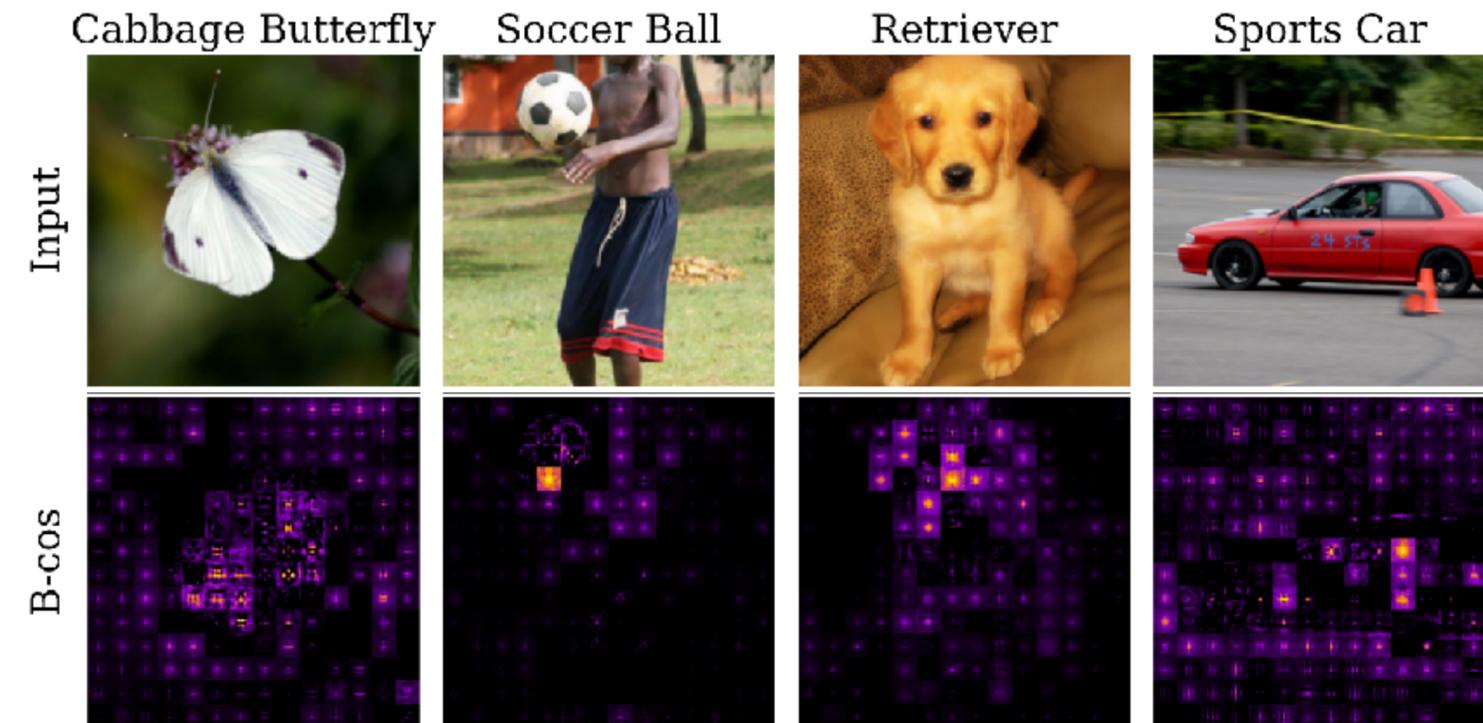
Qualitative Results: class consistency of DAVE attributions



Qualitative Results: DAVE improves inherently-interpretable B-cos Attributions

B-cos ViTs rely on a **conv-stem** for good explanations,

- it is not ideal since most modern ViT architectures do not rely on a conv-stem,
- it can also introduce unwanted inductive biases and move away from a purely transformer based architecture.

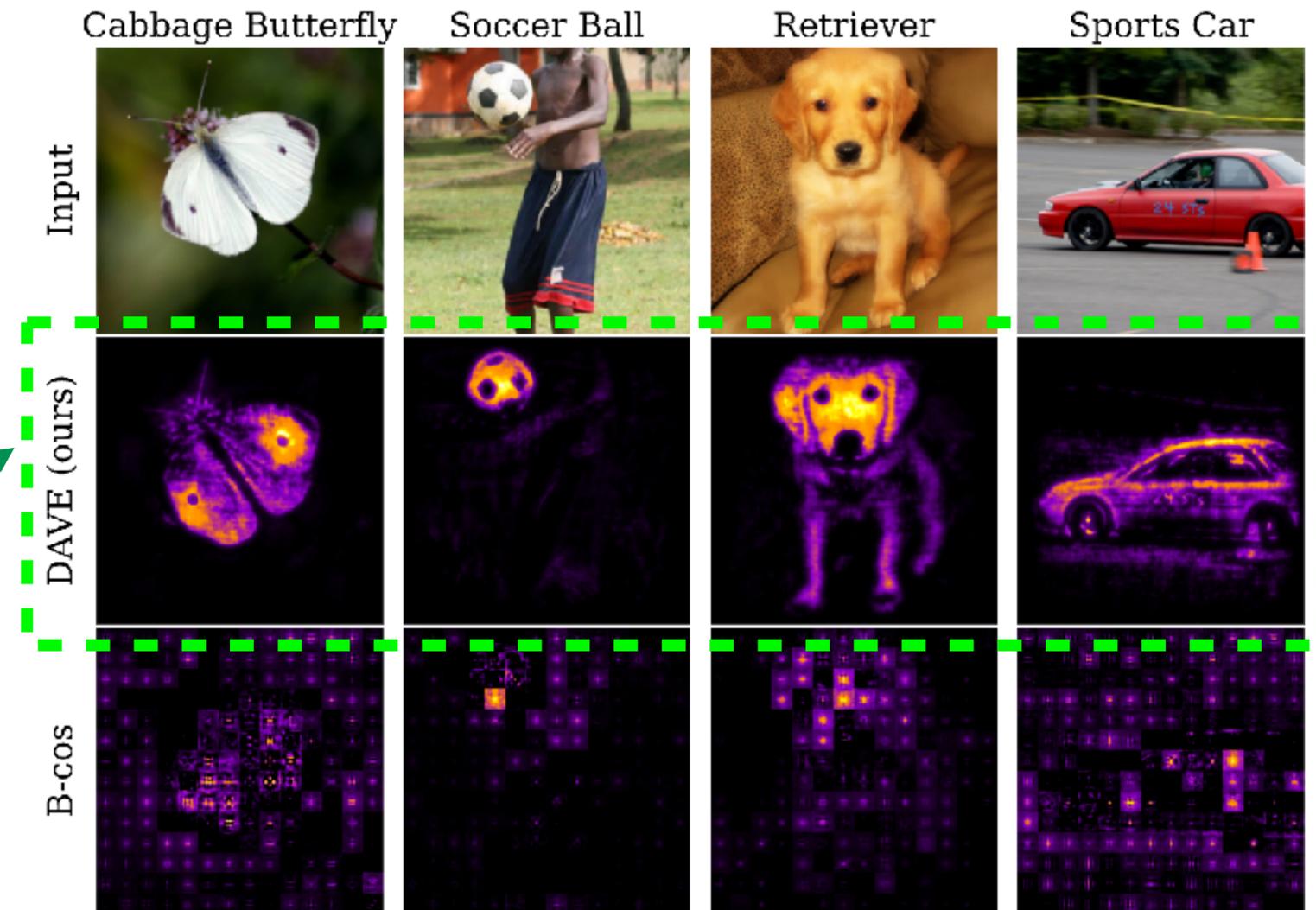


Qualitative Results: DAVE improves inherently-interpretable B-cos Attributions

B-cos ViTs rely on a **conv-stem** for good explanations,

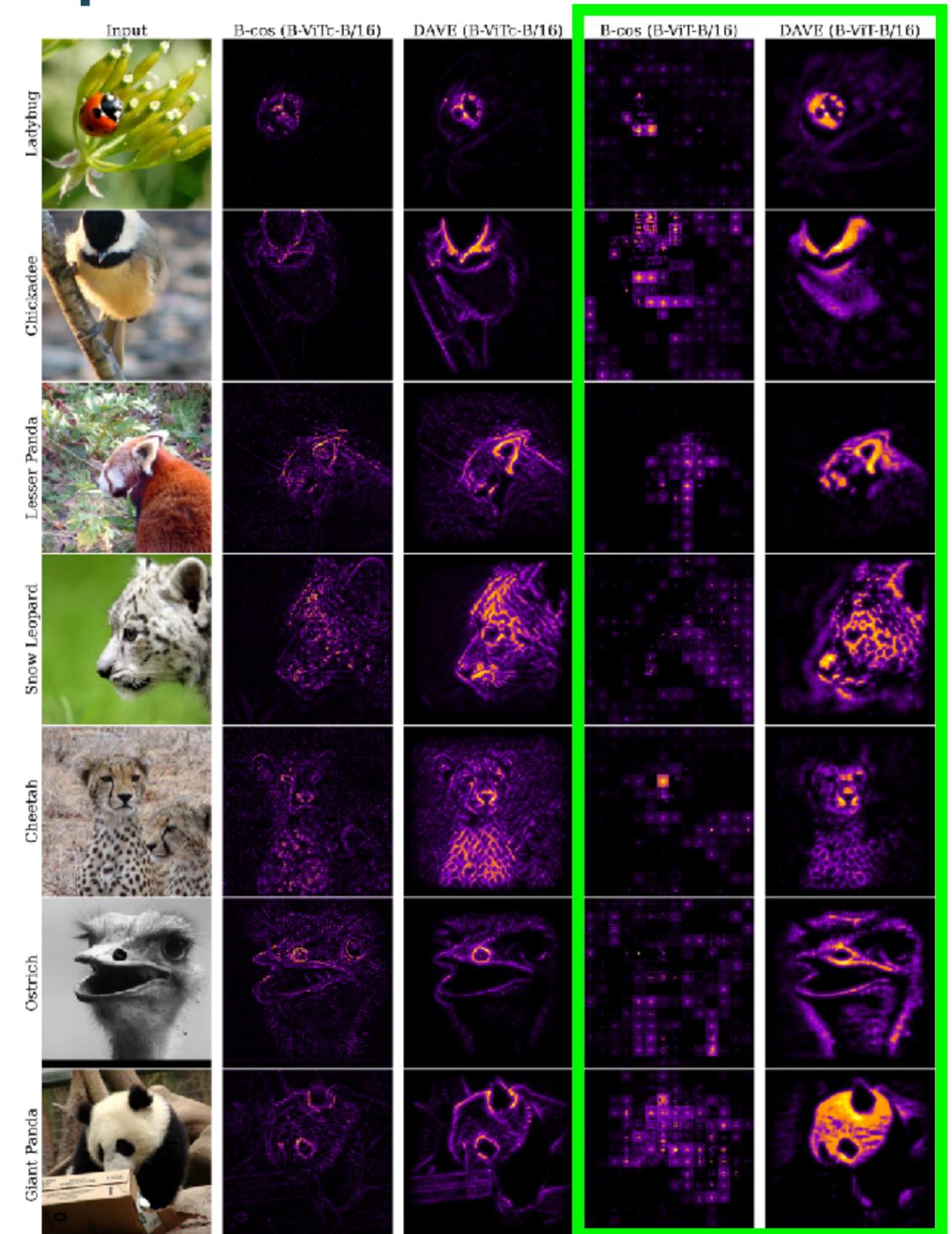
- it is not ideal since most modern ViT architectures do not rely on a conv-stem,
- it can also introduce unwanted inductive biases and move away from a purely transformer based architecture.

✓DAVE, alleviates this reliance on a conv-stem.

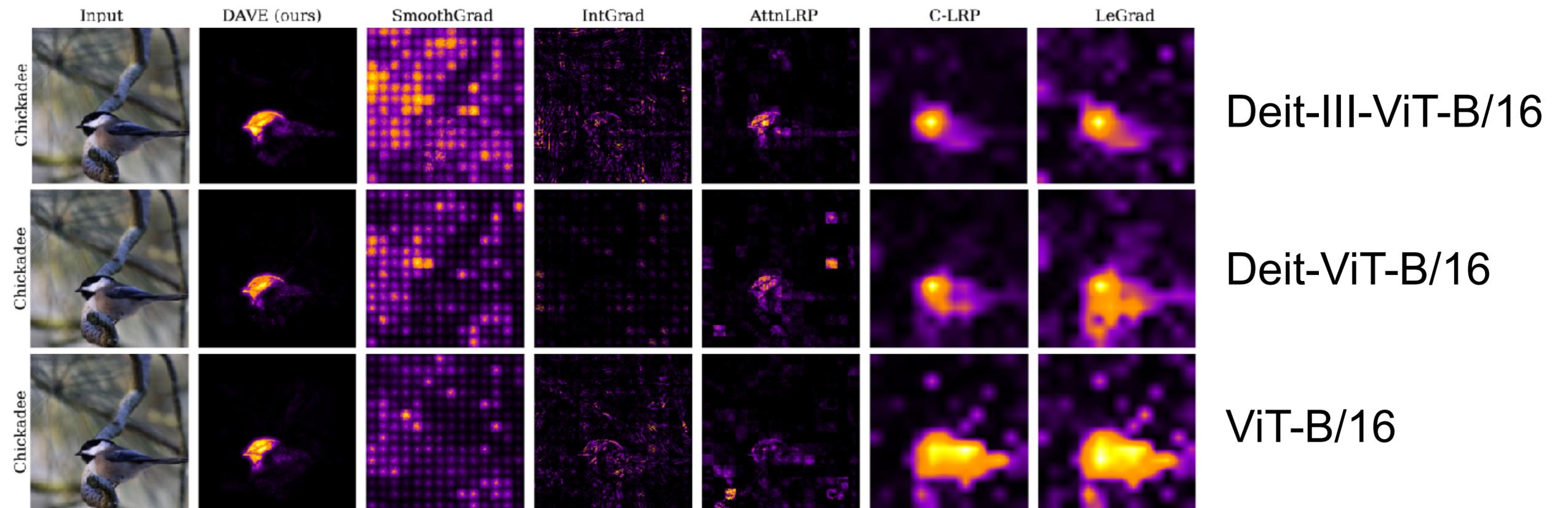


Qualitative Results: DAVE improves inherently-interpretable B-cos Attributions

Can even be argued that the attributions on pure B-cos ViTs (w/o conv-stem) are nicer (at least qualitatively) ?

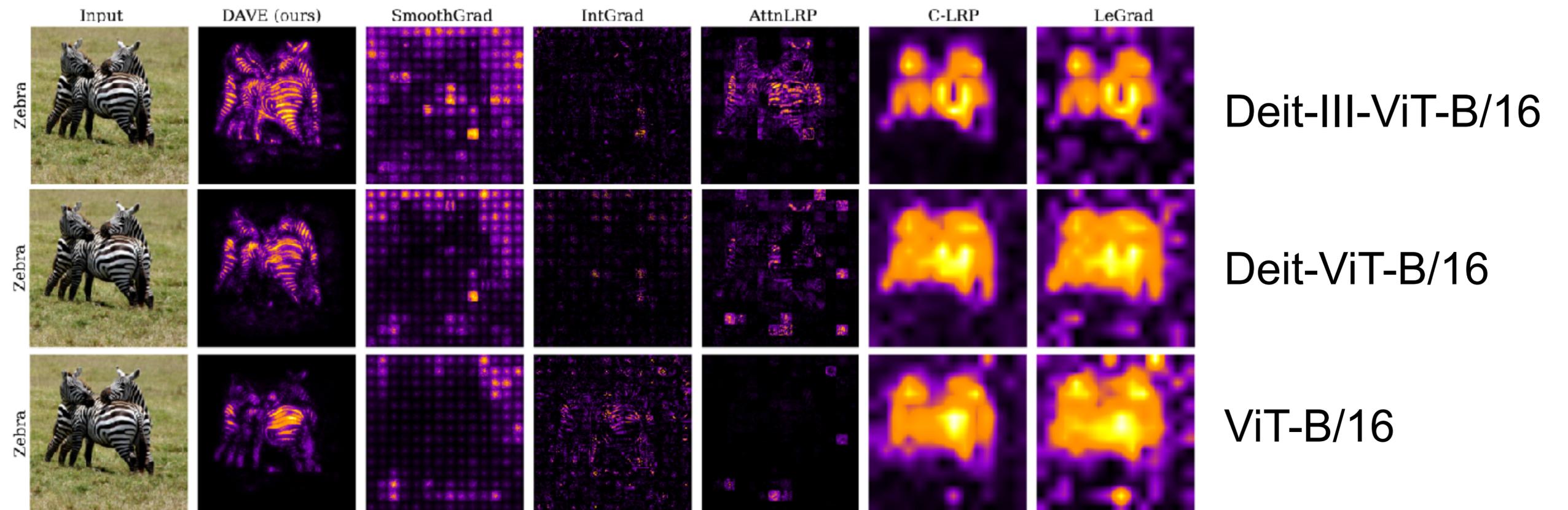


Qualitative Results: consistency across diverse pre-training



✓ DAVE, yields consistent attributions across pre-trained backbones.

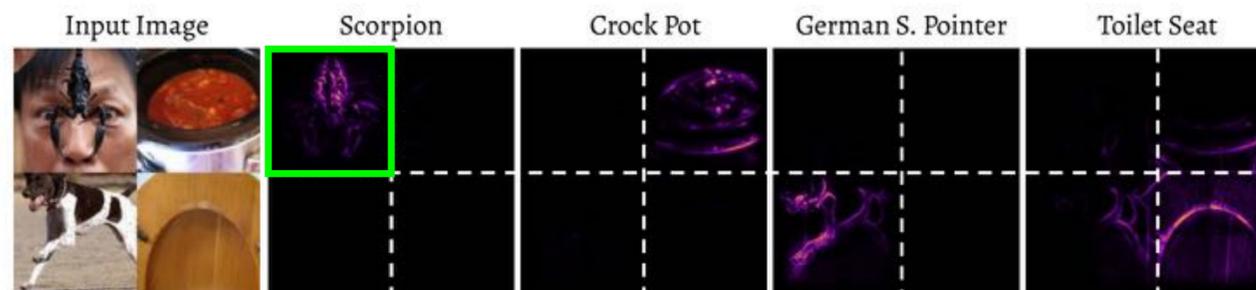
Qualitative Results: consistency across diverse pre-training



✓ DAVE, yields consistent attributions across pre-trained backbones.

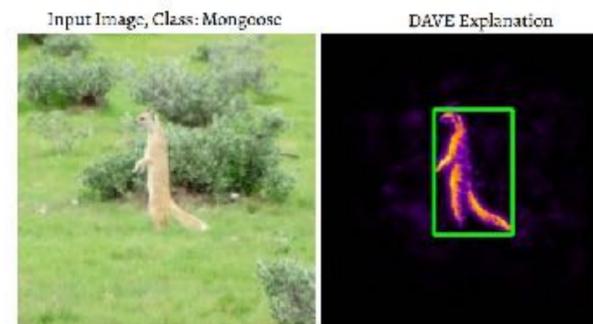
Quantitative Results: localisation via Grid PG and Energy PG

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
DAVE (ours)	60.19	63.52	65.76	51.33
Δ	+1.79	+8.05	+11.92	+1.35

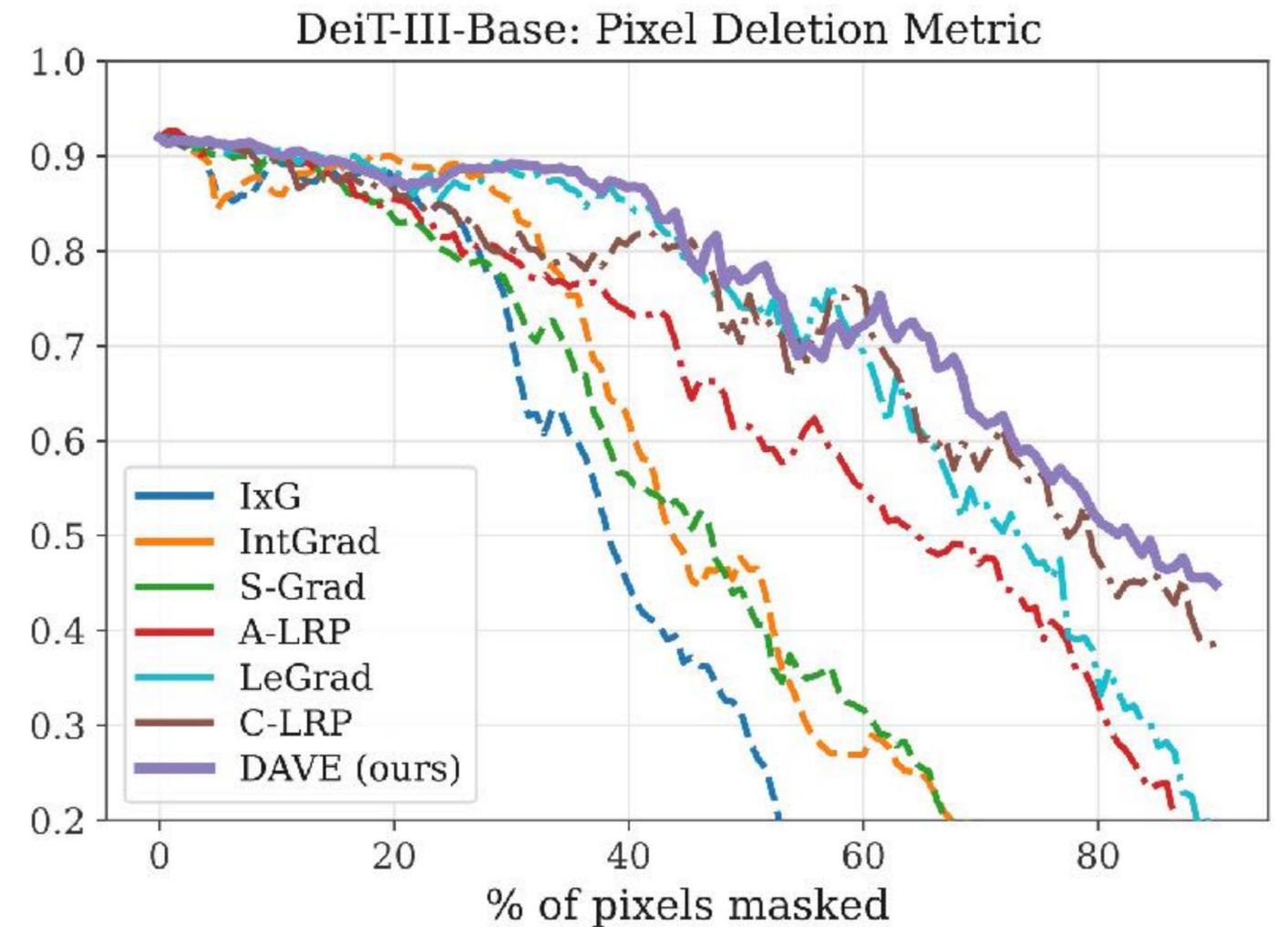
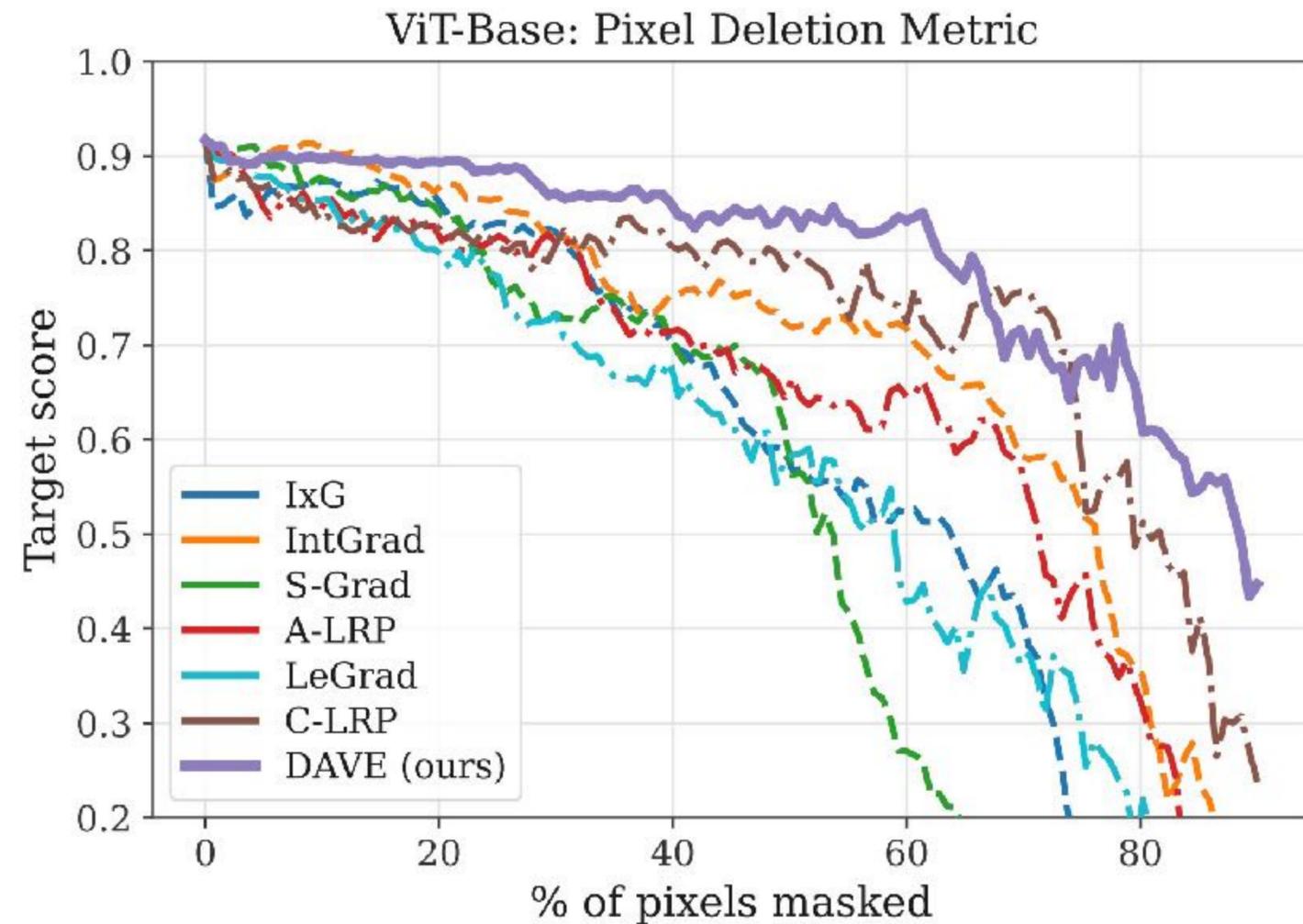


Quantitative Results: localisation via Grid PG and Energy PG

Method	GridPG (%)				EnergyPG (%)			
	ViT-B	DeiT-B	D-III-B	DINO-B	ViT-B	DeiT-B	D-III-B	DINO-B
I×G	32.67	30.25	30.01	33.28	55.33	62.72	67.32	69.98
IntGrad	39.86	36.11	31.68	36.98	58.51	64.22	68.12	72.15
S-Grad	34.27	30.18	31.48	33.13	56.02	60.78	68.10	70.93
LeGrad	47.71	42.58	34.62	28.96	80.06	77.83	77.54	82.26
A-LRP	58.40	54.63	53.84	37.49	60.75	68.16	77.65	75.98
C-LRP	54.98	55.47	52.27	49.99	80.82	79.62	81.94	81.56
DAVE (ours)	60.19	63.52	65.76	51.33	78.60	82.23	82.43	83.38
Δ	+1.79	+8.05	+11.92	+1.35	-2.22	+2.61	+0.49	+1.12



Quantitative Results: pixel perturbation (pixel deletion)



Next steps

Current results only for vision models (classification), also works on self-supervised linear probed models (DINO family):

- ◆ extension to vision language models (CLIP, SigLIP, Llava)
- ◆ extension to convolutional models
- ◆ human study - reviewers ask for that
- ◆ release nicely usable library
- ◆ bias contributions are still not accounted for, so still a problem.

Fun-Fact about DAVE

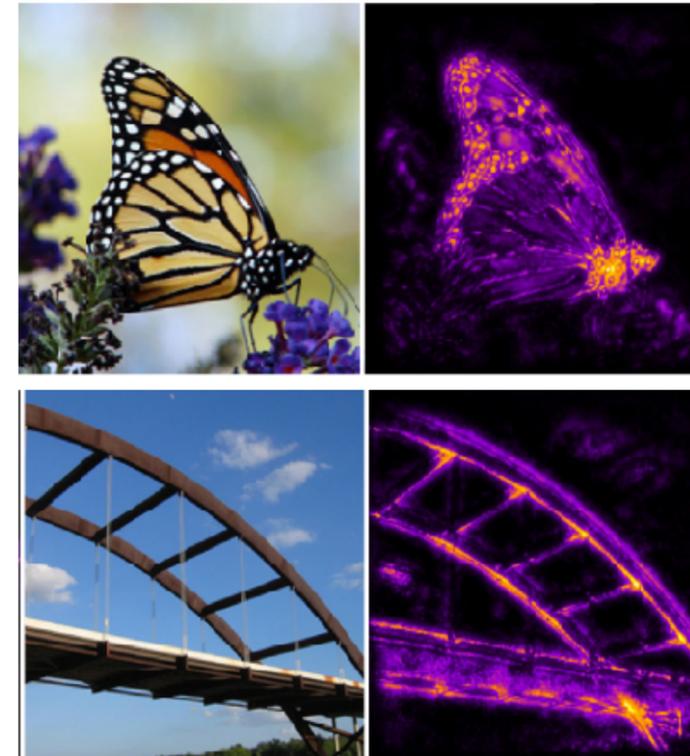
Adam coined the name DAVE - yes it's actually an acronym for the method/title.

But it's also a tribute to **Michelangelo's David**, once hailed as the most precise human sculpture thus fitting, since our DAVE attributions are **highly precise and class-specific**. 🗿 ✨



David by Michelangelo

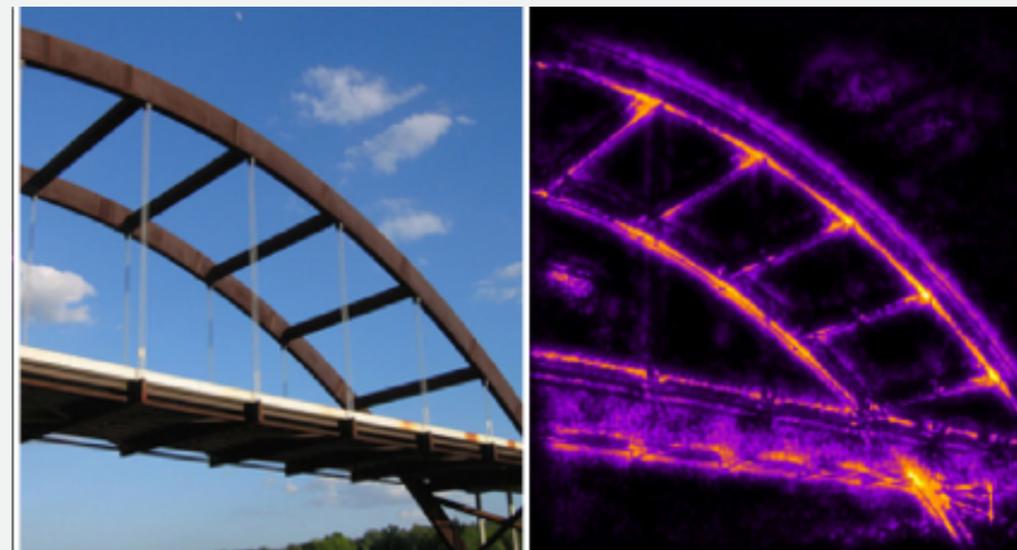
Galleria dell'Accademia, Florence, Italy



David by Michelangelo

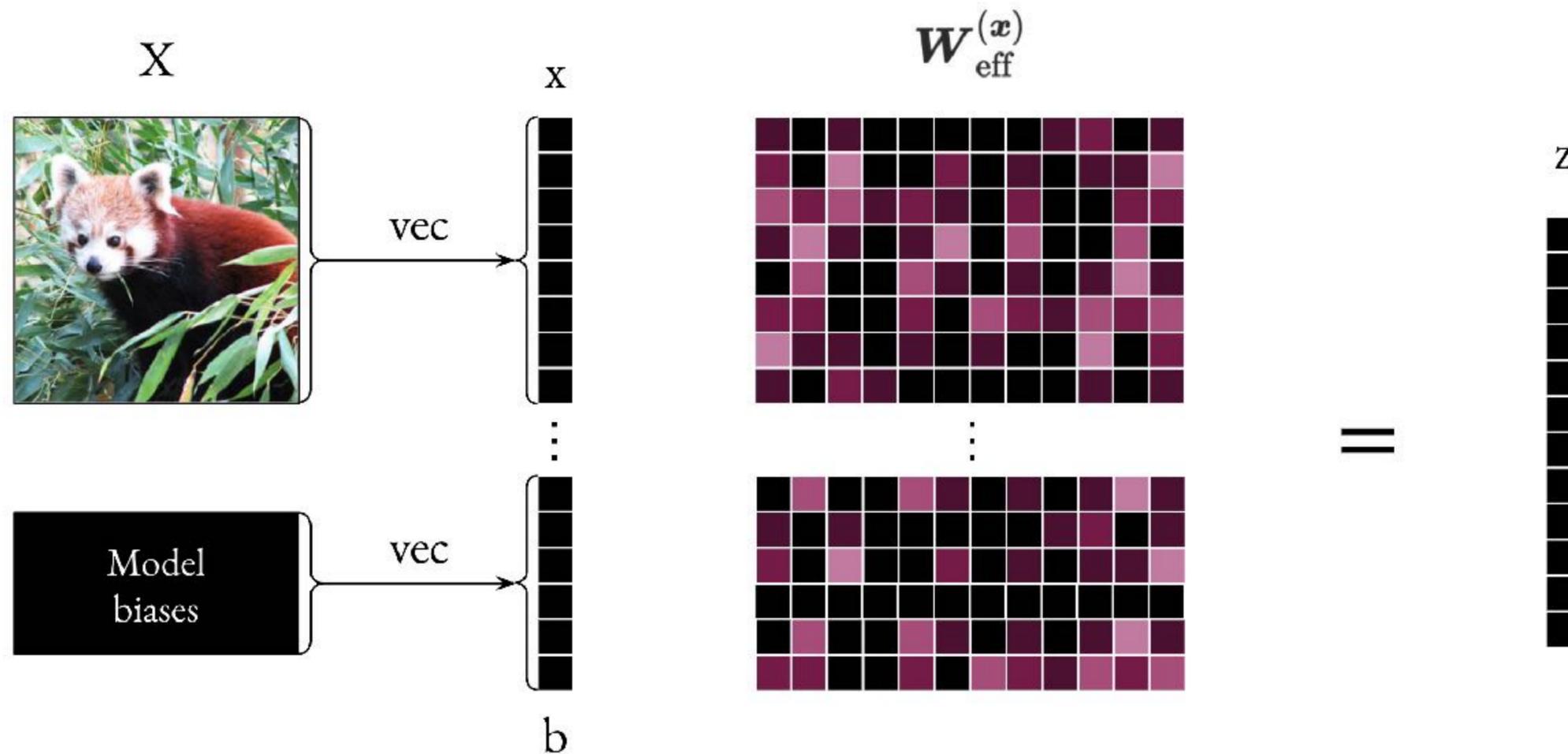
<https://github.com/a-vrobell/DAVE>

Questions, comments, feedback!



DAVE - Better Understanding

ViT is a Dynamic Linear Matrix $z = \text{ViT}(\mathbf{X}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} [\mathbf{x} | \mathbf{b}]$



DAVE - Better Understanding

Layer by layer: general assumption

- If every base layer is of a form:

$$\text{vec}(\mathbf{Z}) = \mathbf{W}(\mathbf{x})\mathbf{x} + \mathbf{b}$$

- Then the whole ViT:

$$\text{ViT}(\mathbf{X}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} [\mathbf{x} | \mathbf{b}]$$

DAVE - Better Understanding

Layer Norm

$$\text{LN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \mu(\mathbf{x})}{\sigma(\mathbf{x})} + \beta$$

$$\text{LN}(\mathbf{x}) = \underbrace{\left(\frac{1}{\sigma(\mathbf{x})} \Gamma H \right)}_{M(\mathbf{x})} \mathbf{x} + \beta$$

Where:

$$H = I - \frac{1}{d} \mathbf{1}\mathbf{1}^\top$$

$$\Gamma = \text{diag}(\gamma)$$

DAVE - Better Understanding

Layer Norm

For a single token:

$$\text{LN}(\mathbf{x}) = \underbrace{\left(\frac{1}{\sigma(\mathbf{x})} \Gamma H \right)}_{M(\mathbf{x})} \mathbf{x} + \beta$$

Where:

$$H = I - \frac{1}{d} \mathbf{1}\mathbf{1}^\top$$

$$\Gamma = \text{diag}(\gamma)$$

DAVE - Better Understanding

Layer Norm

For a single token:

$$\text{LN}(\mathbf{x}) = \underbrace{\left(\frac{1}{\sigma(\mathbf{x})} \Gamma H \right)}_{M(\mathbf{x})} \mathbf{x} + \beta$$

For vectorized tokens:

Blocked-diagonal matrix

Where:

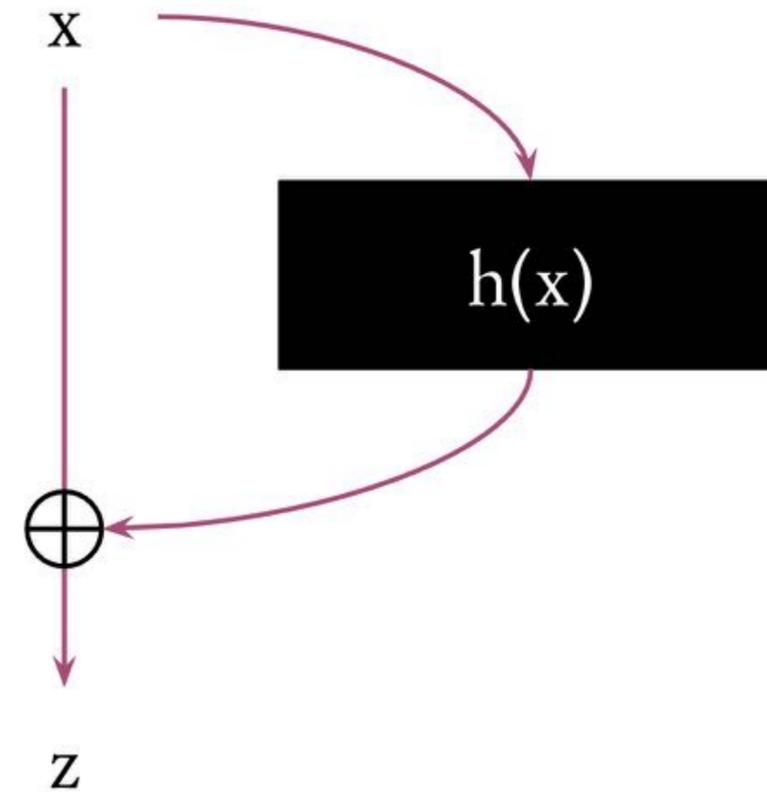
$$H = I - \frac{1}{d} \mathbf{1}\mathbf{1}^\top$$

$$\Gamma = \text{diag}(\gamma)$$

DAVE - Better Understanding

Residual connections

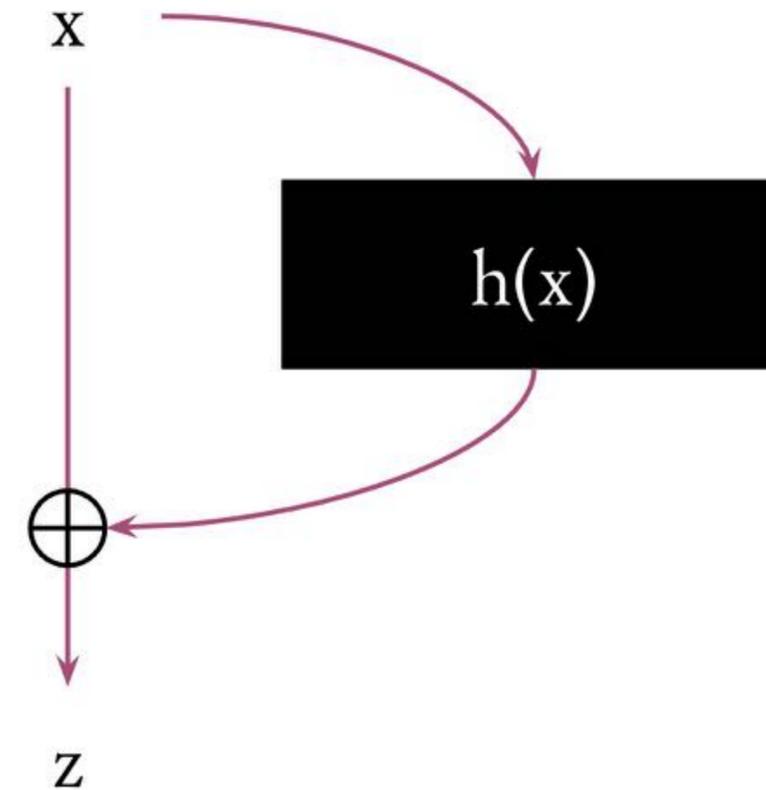
$$z = x + h(x)$$



DAVE - Better Understanding

Residual connections

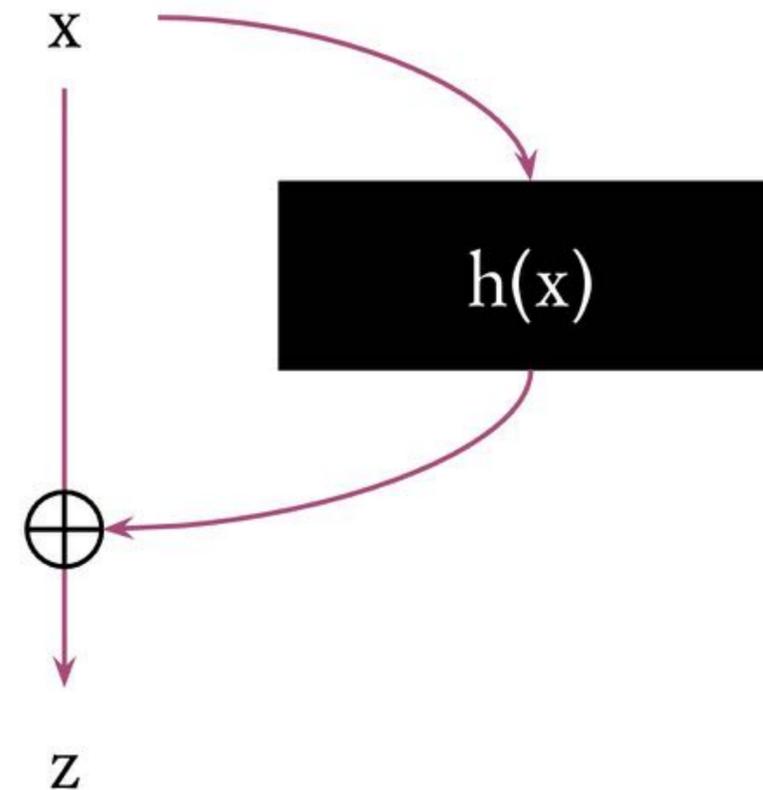
$$\begin{aligned}z &= \mathbf{x} + h(\mathbf{x}) \\ &= \mathbf{x} + [\mathbf{W}(\mathbf{x})\mathbf{x} + \mathbf{b}]\end{aligned}$$



DAVE - Better Understanding

Residual connections

$$\begin{aligned}z &= \mathbf{x} + h(\mathbf{x}) \\ &= \mathbf{x} + [\mathbf{W}(\mathbf{x})\mathbf{x} + \mathbf{b}] \\ &= [\mathbf{I} + \mathbf{W}(\mathbf{x})]\mathbf{x} + \mathbf{b}\end{aligned}$$



DAVE - The Method

$$\mathbf{f}(\mathbf{x}) = \mathbf{W}(\mathbf{x})\mathbf{x}$$

$$\mathbf{J}_f(\mathbf{x}) = \underbrace{\mathbf{W}(\mathbf{x})}_{\text{Our ViT matrix}} + \underbrace{\sum_{i=1}^d x_i \frac{\partial \mathbf{W}}{\partial x_i}(\mathbf{x})}_{\text{JVP}}$$

Our ViT
matrix

JVP

DAVE - The Method

Effective weights and how to find them efficiently

- **Detach** all dynamic matrices on the forward pass (attention, GELU multiplier, etc.)
- Compute “detached” **backward pass**
- Like: CoDA, B-cos

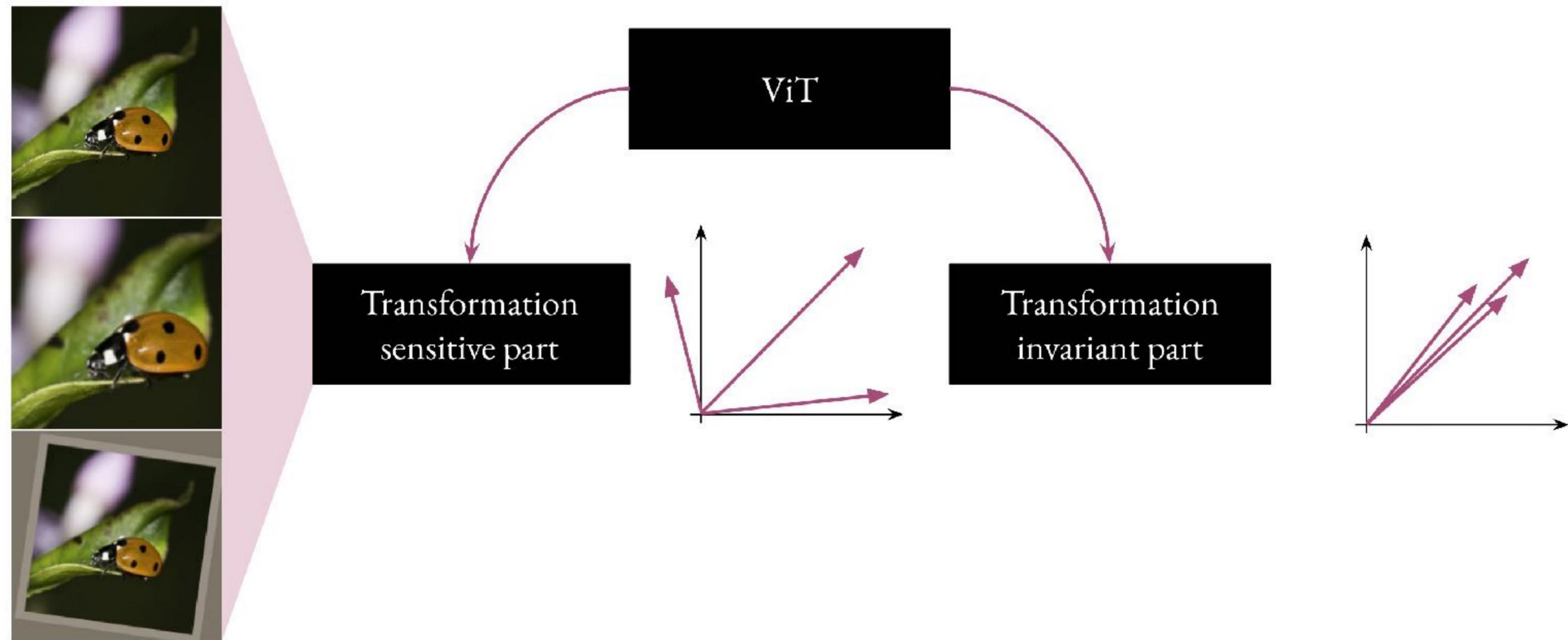
$$f(\mathbf{x}) = \mathbf{W}(\mathbf{x})\mathbf{x}$$

$$J_f(\mathbf{x}) = \underbrace{\mathbf{W}(\mathbf{x})}_{\text{Our ViT matrix}} + \underbrace{\sum_{i=1}^d x_i \frac{\partial \mathbf{W}}{\partial x_i}(\mathbf{x})}_{\text{JVP}}$$

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

DAVE - Better Understanding

Intuition



DAVE - Better Understanding

Assumption: For a well-trained model

$$C(\mathbf{x}) = C^*(\mathbf{x}) + n_{\tau}(\mathbf{x}) + n_{\epsilon}(\mathbf{x})$$

Dynamic-linear
representation



Attribution

Spatial-related
global noise

high-frequency
local noise

DAVE - Better Understanding

Assumption: For a well-trained model

- Global noise is spatially “**fixed**”:

$$\mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger(\mathbf{n}_\tau(\tau(\mathbf{x})))] = \mathbf{0}$$

- Attribution is **transformation-invariant**:

$$\mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger(\mathbf{C}^*(\tau(\mathbf{x})))] \approx \mathbf{C}^*(\mathbf{x})$$

- Local noise has significantly **higher frequency band** than attribution

DAVE - Better Understanding

Separating spatial noise

$$\mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger(\mathbf{C}(\tau(\mathbf{x})))] =$$

DAVE - Better Understanding

Separating spatial noise

$$\mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger(\mathbf{C}(\tau(\mathbf{x})))] =$$

$$\mathbb{E}_{\tau \sim p_{transform}} \tau^\dagger[\mathbf{C}^*(\tau(\mathbf{x})) + \mathbf{n}_\tau(\tau(\mathbf{x})) + \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

DAVE - Better Understanding

Separating spatial noise

$$\mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger (\mathbf{C}(\tau(\mathbf{x})))] =$$

$$\mathbb{E}_{\tau \sim p_{transform}} \tau^\dagger [\mathbf{C}^*(\tau(\mathbf{x})) + \mathbf{n}_\tau(\tau(\mathbf{x})) + \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

$$\mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger \mathbf{C}^*(\tau(\mathbf{x}))] + \mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger \mathbf{n}_\tau(\tau(\mathbf{x}))] + \mathbb{E}_{\tau \sim p_{transform}} [\tau^\dagger \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

DAVE - Better Understanding

Separating spatial noise

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger (\mathbf{C}(\tau(\mathbf{x})))] =$$

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} \tau^\dagger [\mathbf{C}^*(\tau(\mathbf{x})) + \mathbf{n}_\tau(\tau(\mathbf{x})) + \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{C}^*(\tau(\mathbf{x}))] + \mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{n}_\tau(\tau(\mathbf{x}))] + \mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

Plugging in the 1st and 2nd assumption:

- Global noise is spatially **“fixed”**:

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger (\mathbf{n}_\tau(\tau(\mathbf{x})))] = 0$$

- Attribution is **transformation-invariant**:

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger (\mathbf{C}^*(\tau(\mathbf{x})))] \approx \mathbf{C}^*(\mathbf{x})$$

DAVE - Better Understanding

Separating spatial noise

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger (\mathbf{C}(\tau(\mathbf{x})))] =$$

$$\mathbb{E}_{\tau \sim p_{\text{transform}}} \tau^\dagger [\mathbf{C}^*(\tau(\mathbf{x})) + \mathbf{n}_\tau(\tau(\mathbf{x})) + \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

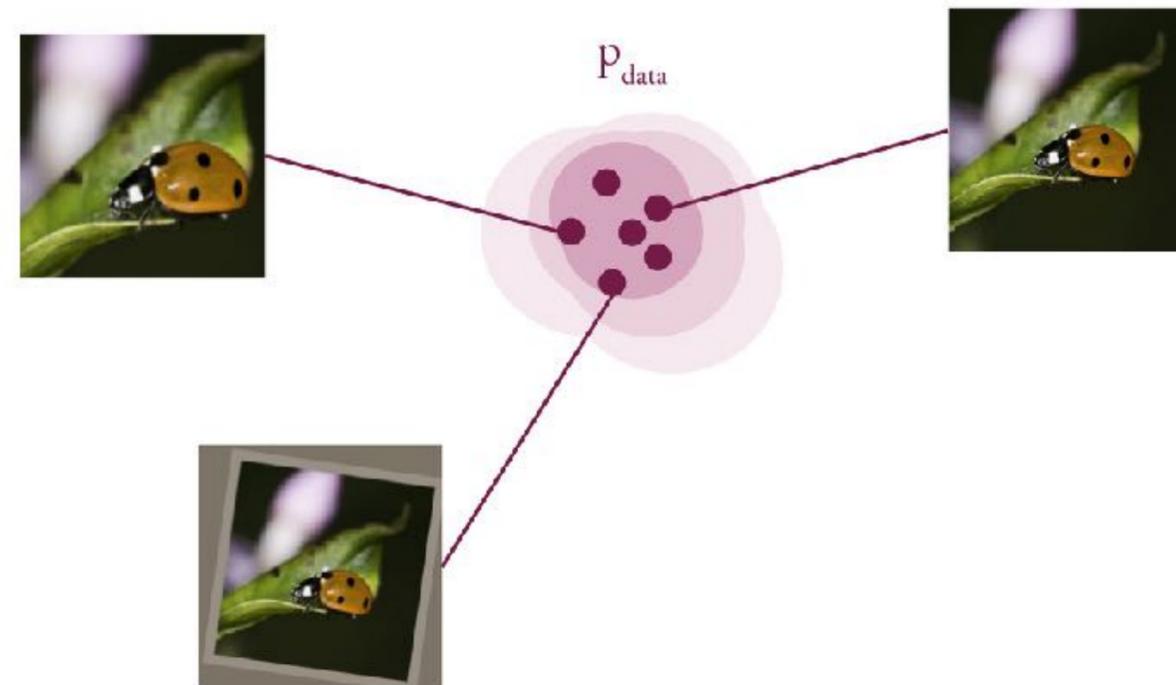
$$\mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{C}^*(\tau(\mathbf{x}))] + \mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{n}_\tau(\tau(\mathbf{x}))] + \mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{n}_\epsilon(\tau(\mathbf{x}))] =$$

$$\mathbf{C}^*(\mathbf{x}) + \mathbb{E}_{\tau \sim p_{\text{transform}}} [\tau^\dagger \mathbf{n}_\epsilon(\tau(\mathbf{x}))]$$

DAVE - Better Understanding

Constructing neighbourhood for reliable and robust features:

sample points as small augmentations applied to the input such that the output logic distribution is preserved



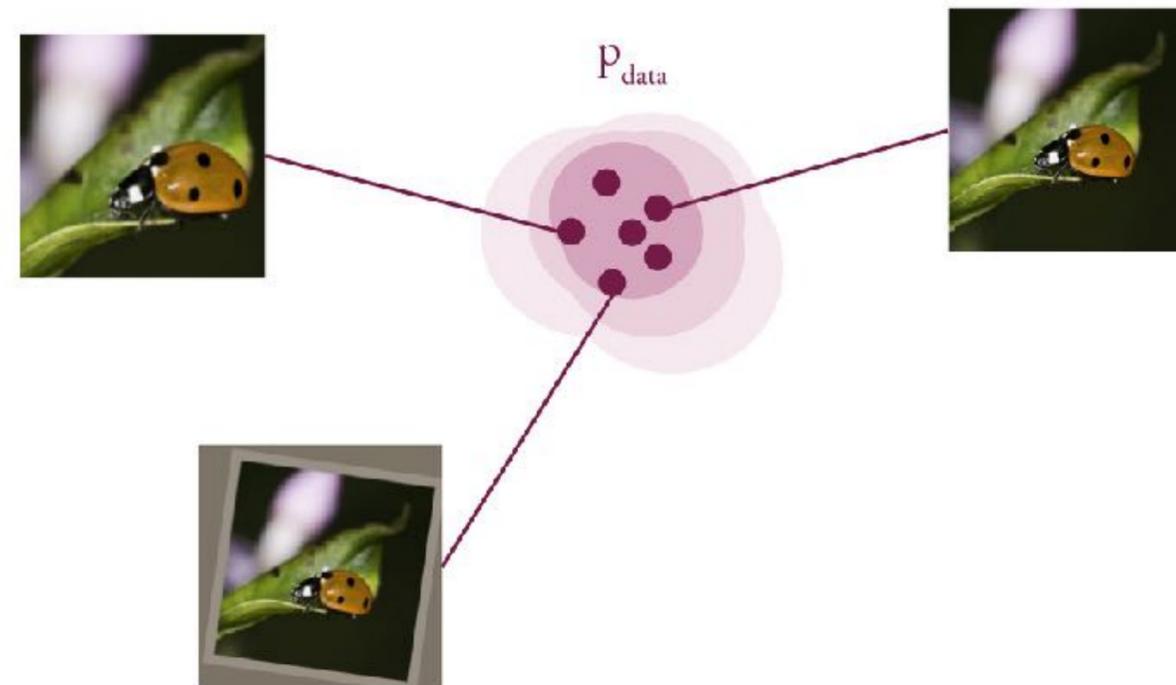
DAVE - Better Understanding

Transformations as layers for attribution pushback:

Each transformation is just **another layer** $\tau(\mathbf{x}) = \mathbf{W}\mathbf{x}$, $\mathbf{x} = \text{vec}(\mathbf{X})$

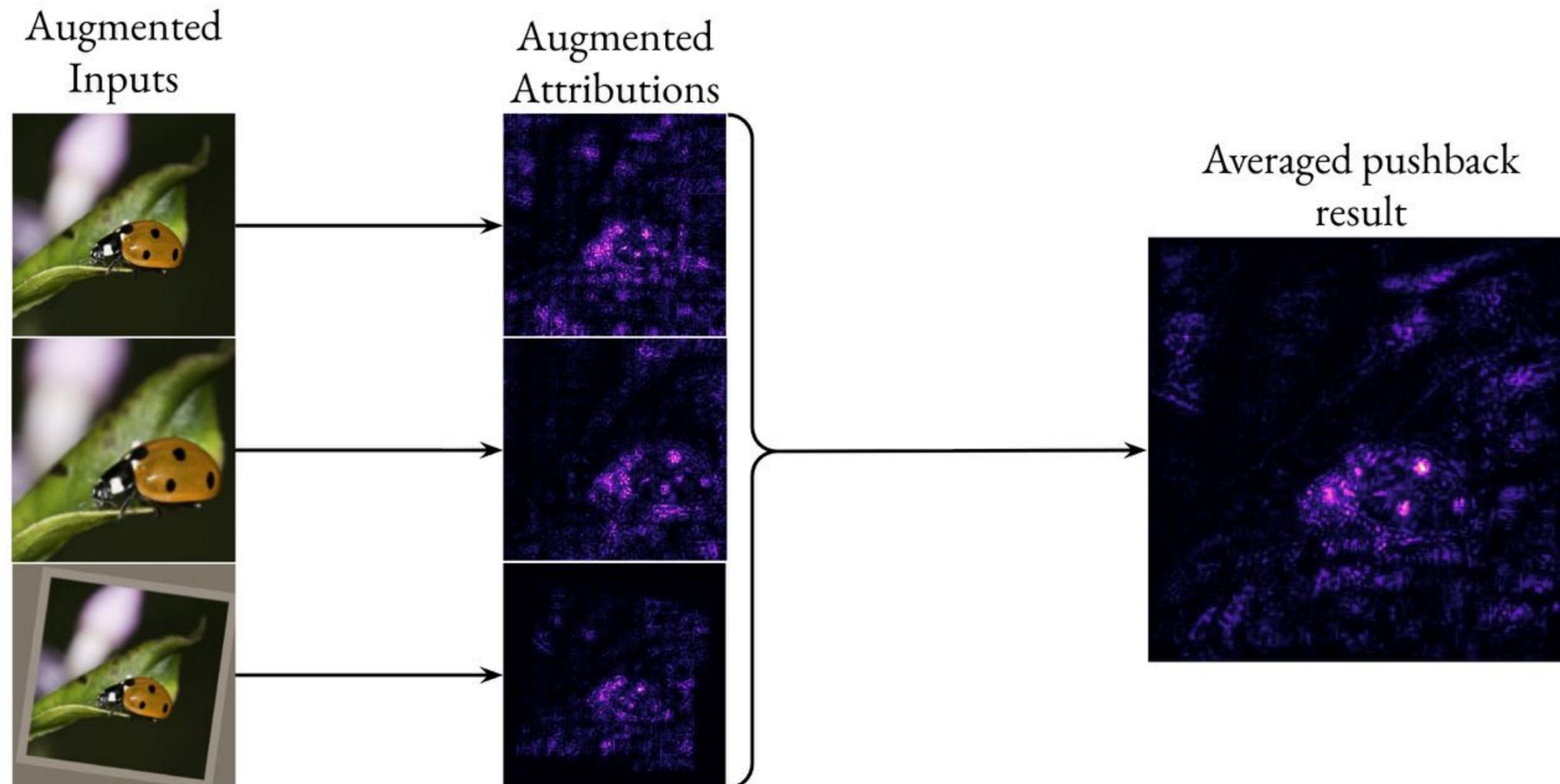
Transformations

- ◆ Rotation
- ◆ Offset
- ◆ Padding
- ◆ Scaling
- ◆ Blurring



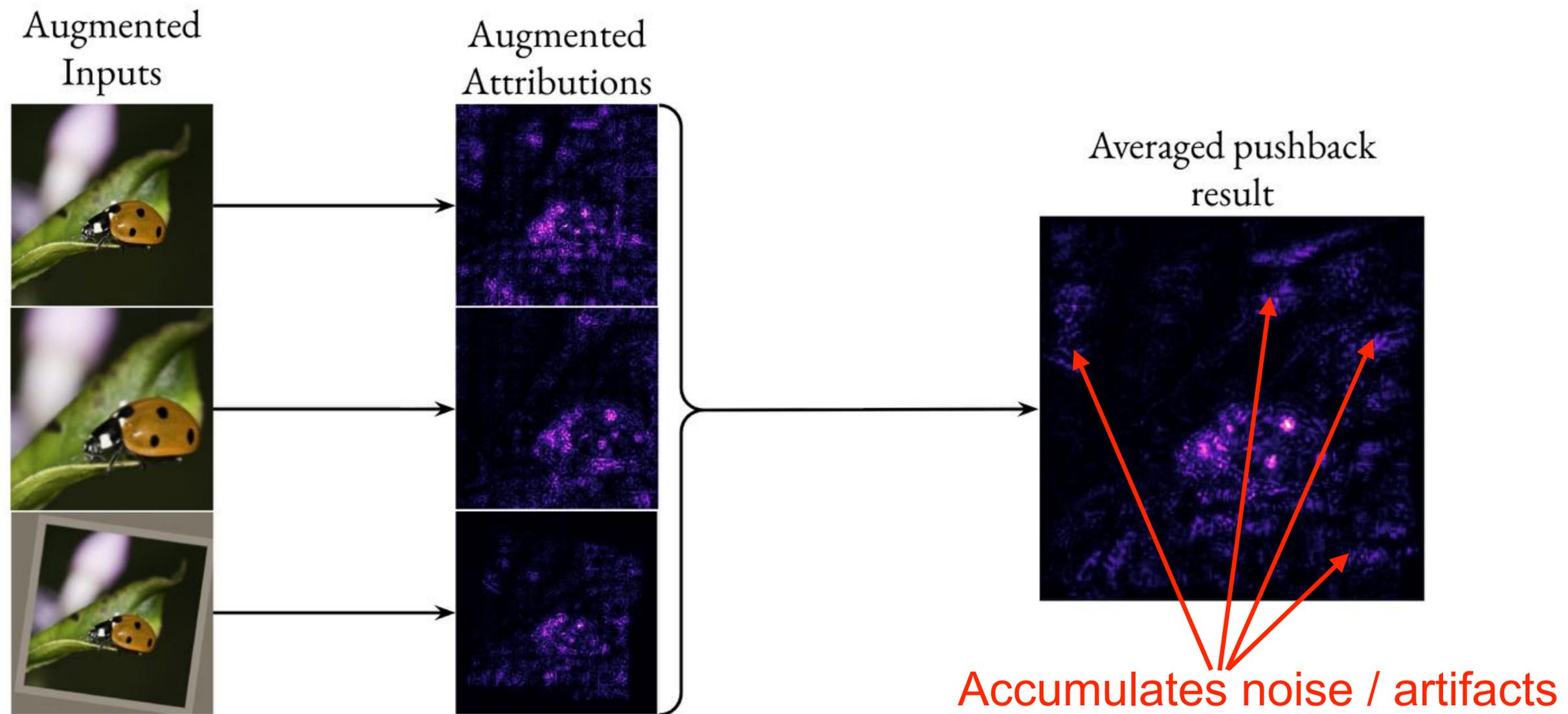
DAVE - Better Understanding

Effective input **equivariant** attributions



DAVE - Better Understanding

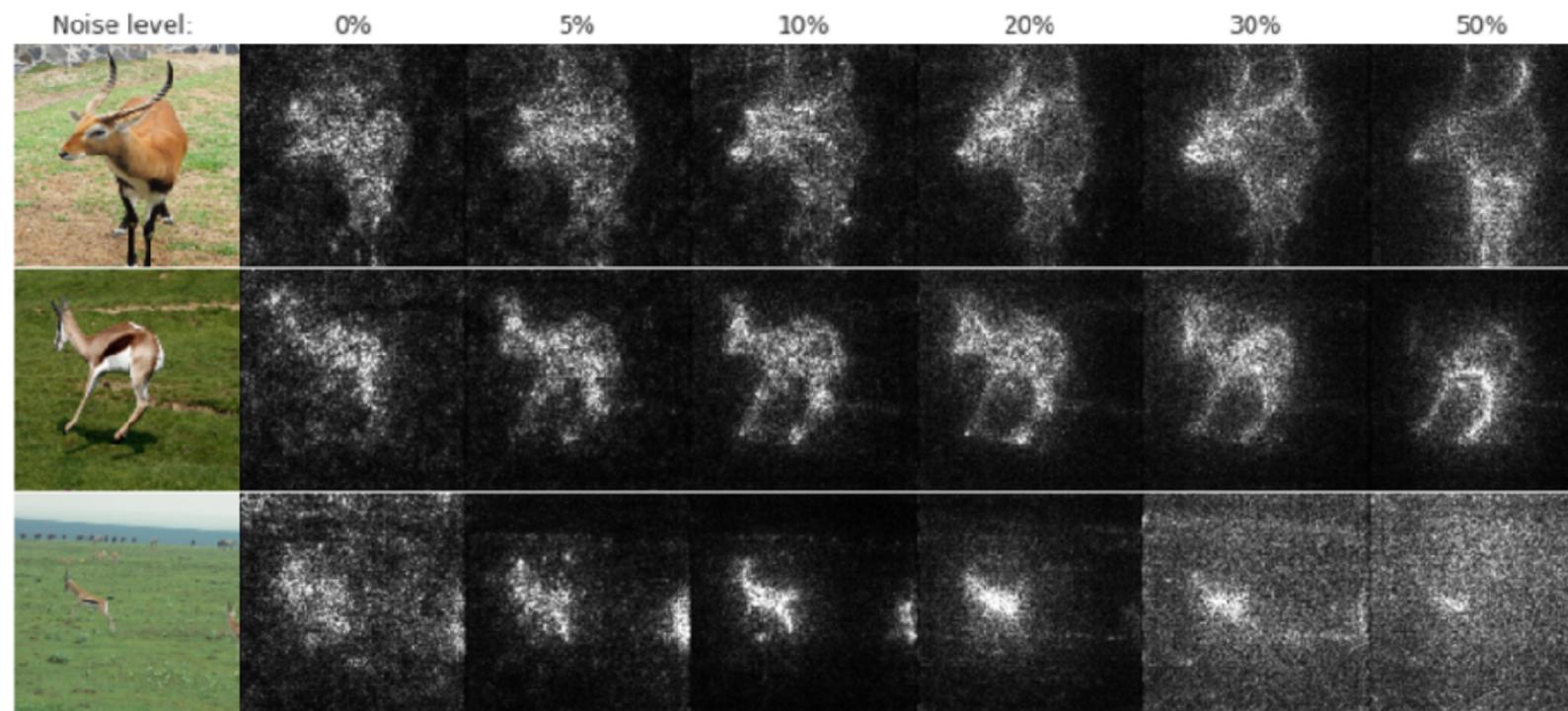
Effective input **equivariant** attributions



DAVE - Better Understanding

Separating local noise: removing artifacts

SmoothGrad: removing noise by adding noise



DAVE - Better Understanding

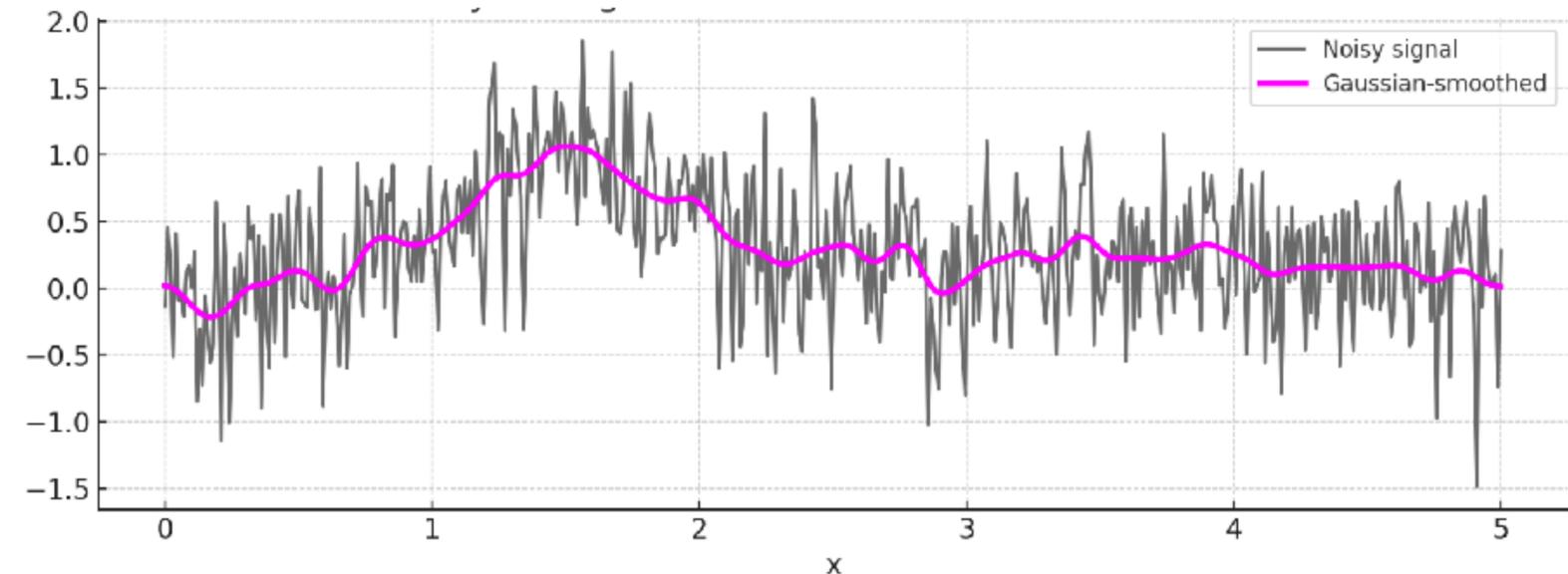
Adding noise = low-pass filtering (in attribution space)

Our attribution:

$$\mathbf{C}(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

Adding noise:

$$\tilde{\mathbf{C}}(\mathbf{x}) = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\mathbf{C}(\underbrace{\mathbf{x} + \epsilon}_{\mathbf{z}})]$$



DAVE - Better Understanding

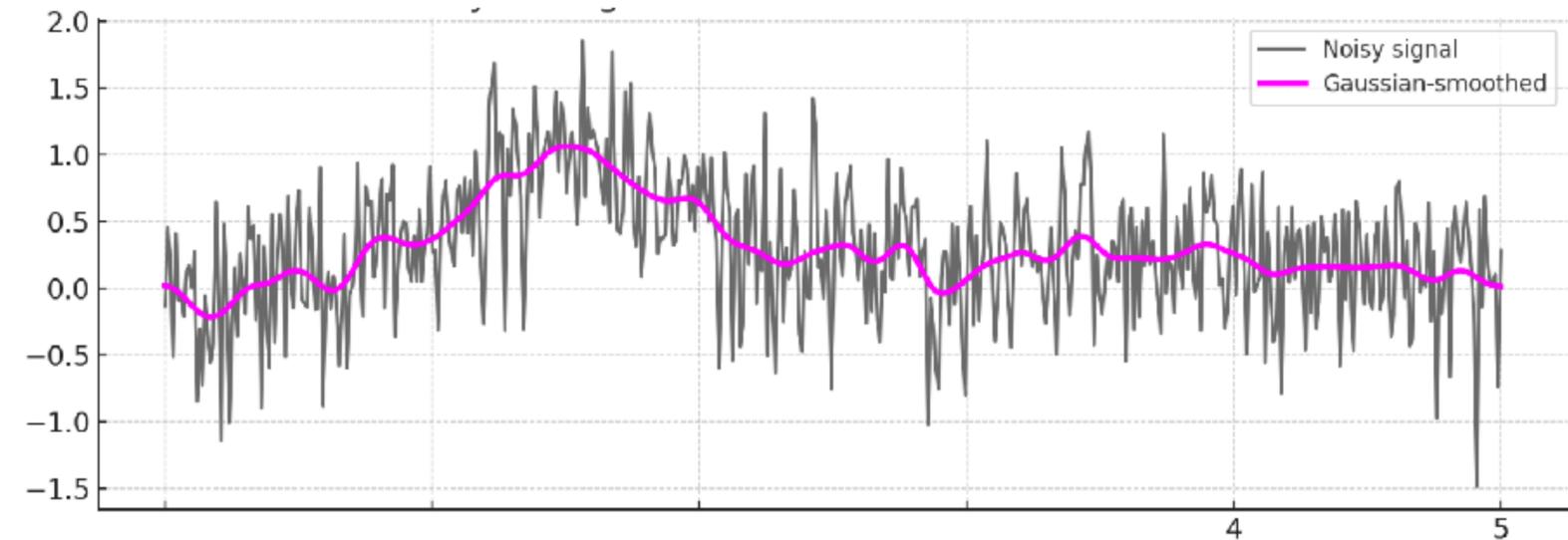
Adding noise = low-pass filtering (in attribution space)

Our attribution:

$$\mathbf{C}(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

Adding noise:

$$\tilde{\mathbf{C}}(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\mathbf{C}(\underbrace{\mathbf{x} + \boldsymbol{\epsilon}}_{\mathbf{z}})] = \int \mathbf{C}(\mathbf{z}) \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}) d\boldsymbol{\epsilon} =$$



DAVE - Better Understanding

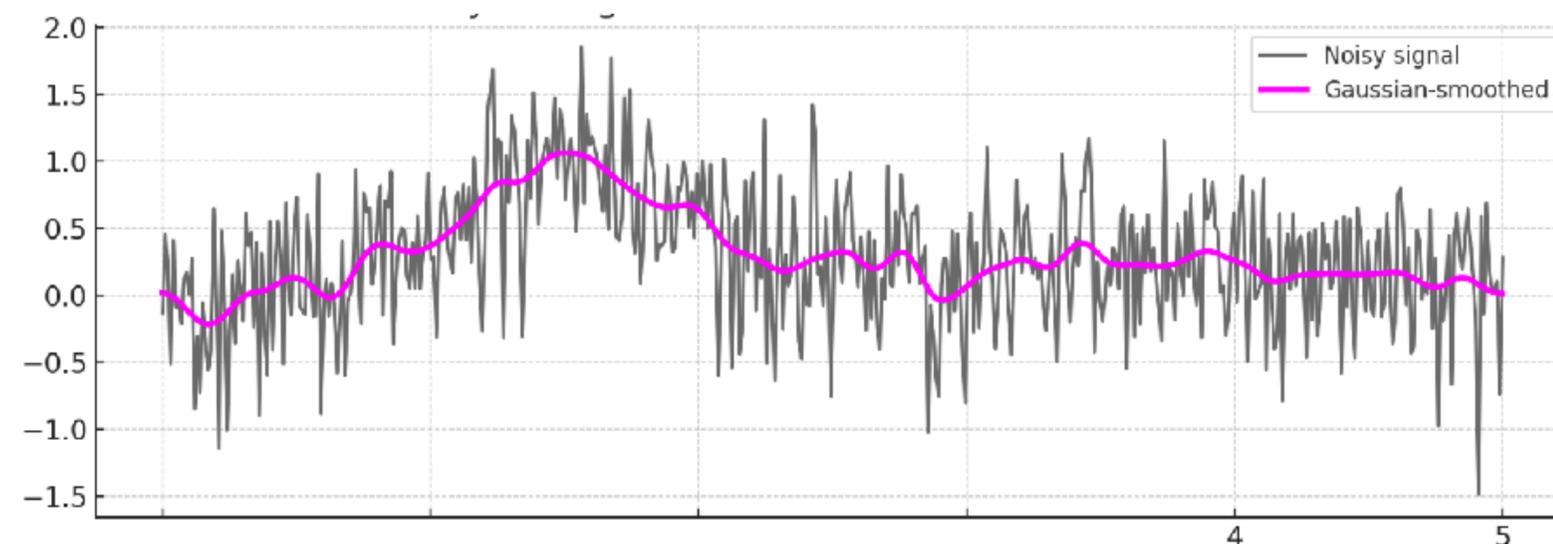
Adding noise = low-pass filtering (in attribution space)

Our attribution:

$$\mathbf{C}(\mathbf{x}) = \mathbf{W}_{\text{eff}}^{(\mathbf{x})} \mathbf{x}$$

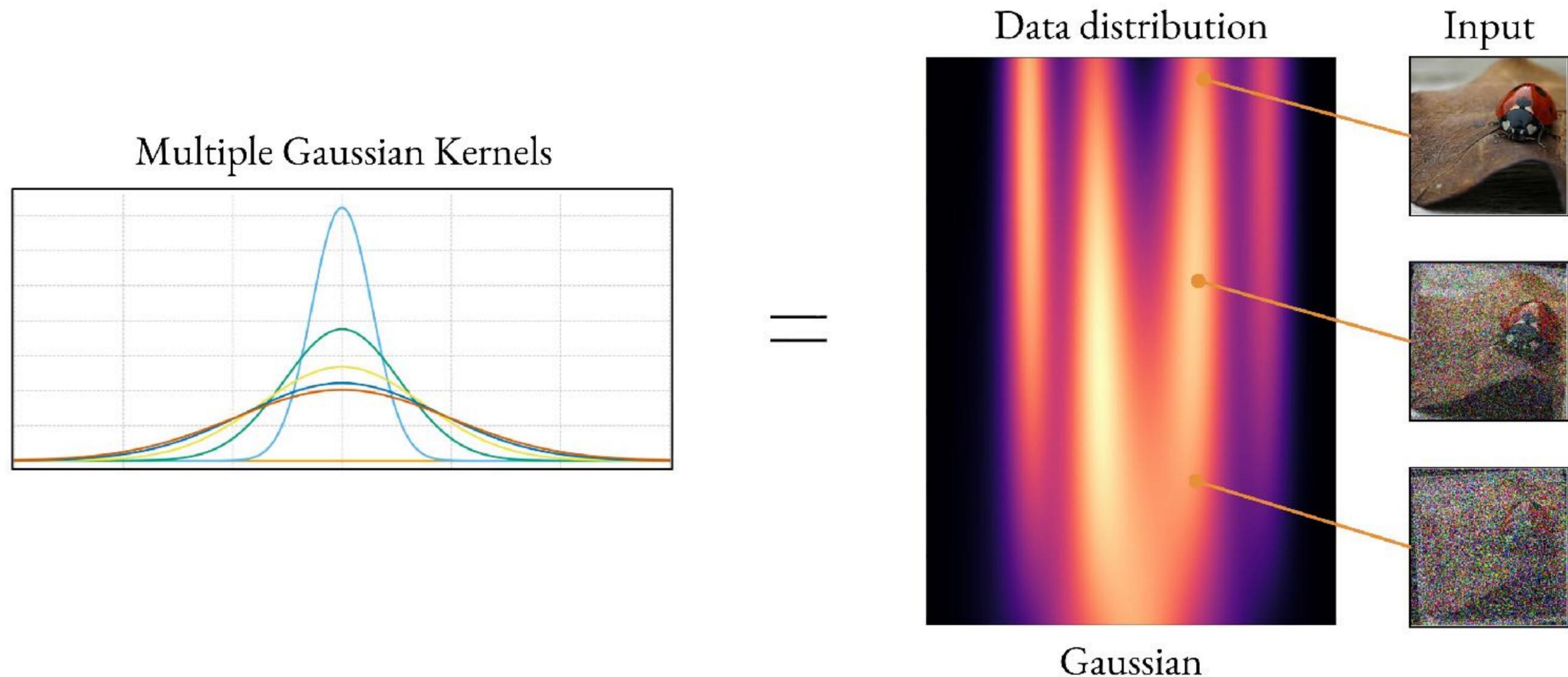
Adding noise:

$$\begin{aligned} \tilde{\mathbf{C}}(\mathbf{x}) &= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\mathbf{C}(\underbrace{\mathbf{x} + \boldsymbol{\epsilon}}_z)] = \int \mathbf{C}(\mathbf{z}) \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}) d\boldsymbol{\epsilon} = \\ &= \int \mathbf{C}(\mathbf{z}) \mathbf{G}(\mathbf{z} - \mathbf{x}) d\mathbf{z} = (\mathbf{C} * \mathbf{G})(\mathbf{x}) \end{aligned}$$



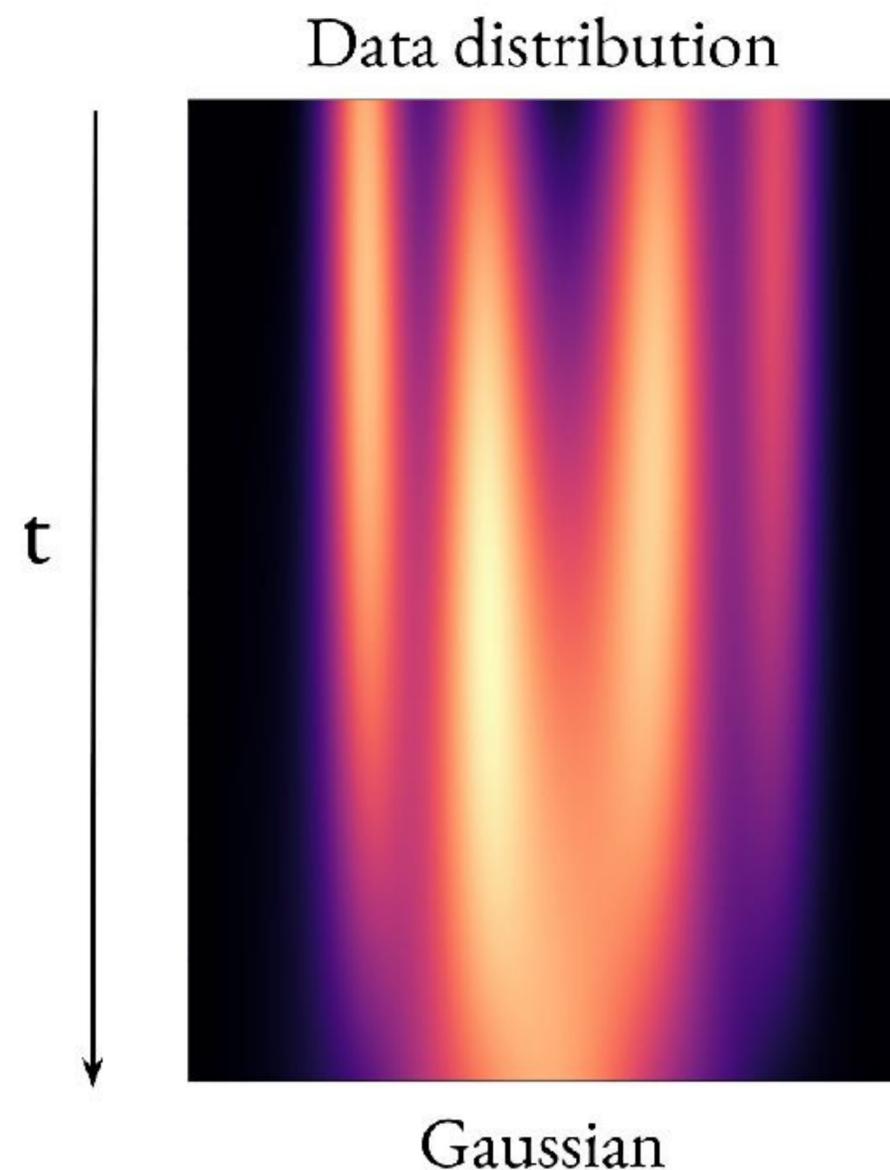
DAVE - Better Understanding

Let's apply multiple low-pass filters: noise interpolation path



DAVE - Better Understanding

Let's apply multiple low-pass filters: noise interpolation path



- Variance-preserving path (chosen empirically):

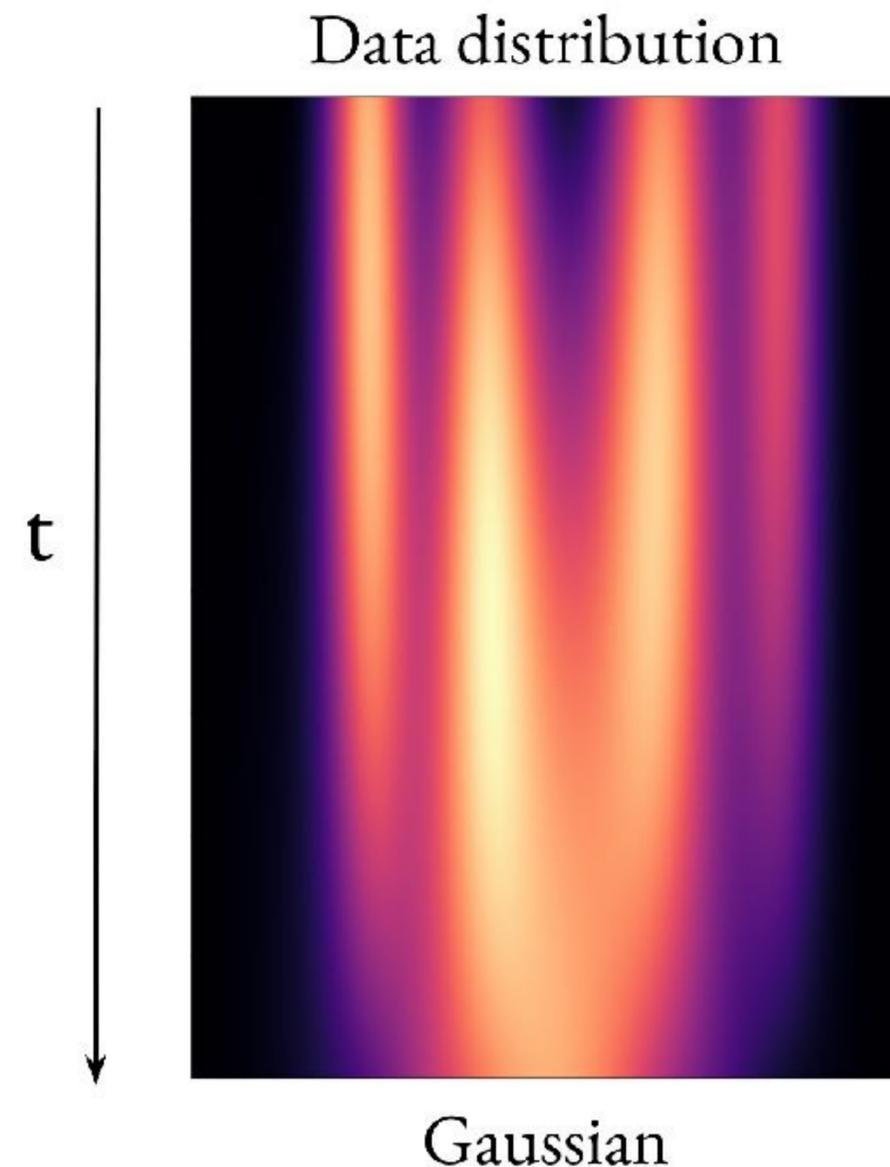
$$\mathbf{x}_t = \cos(t)\mathbf{x} + \sin(t)\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Final aggregation:

$$\begin{aligned} \tilde{\mathbf{C}}(\mathbf{x}) &= \mathbb{E}_{t,\boldsymbol{\epsilon}}[\mathbf{C}(\mathbf{x}_t)] = \\ &= \mathbb{E}_t[(\mathbf{C} * \mathbf{G}_t)(\cos(t)\mathbf{x})] \end{aligned}$$

DAVE - Better Understanding

Let's apply multiple low-pass filters: noise interpolation path



- Variance-preserving path (chosen empirically):

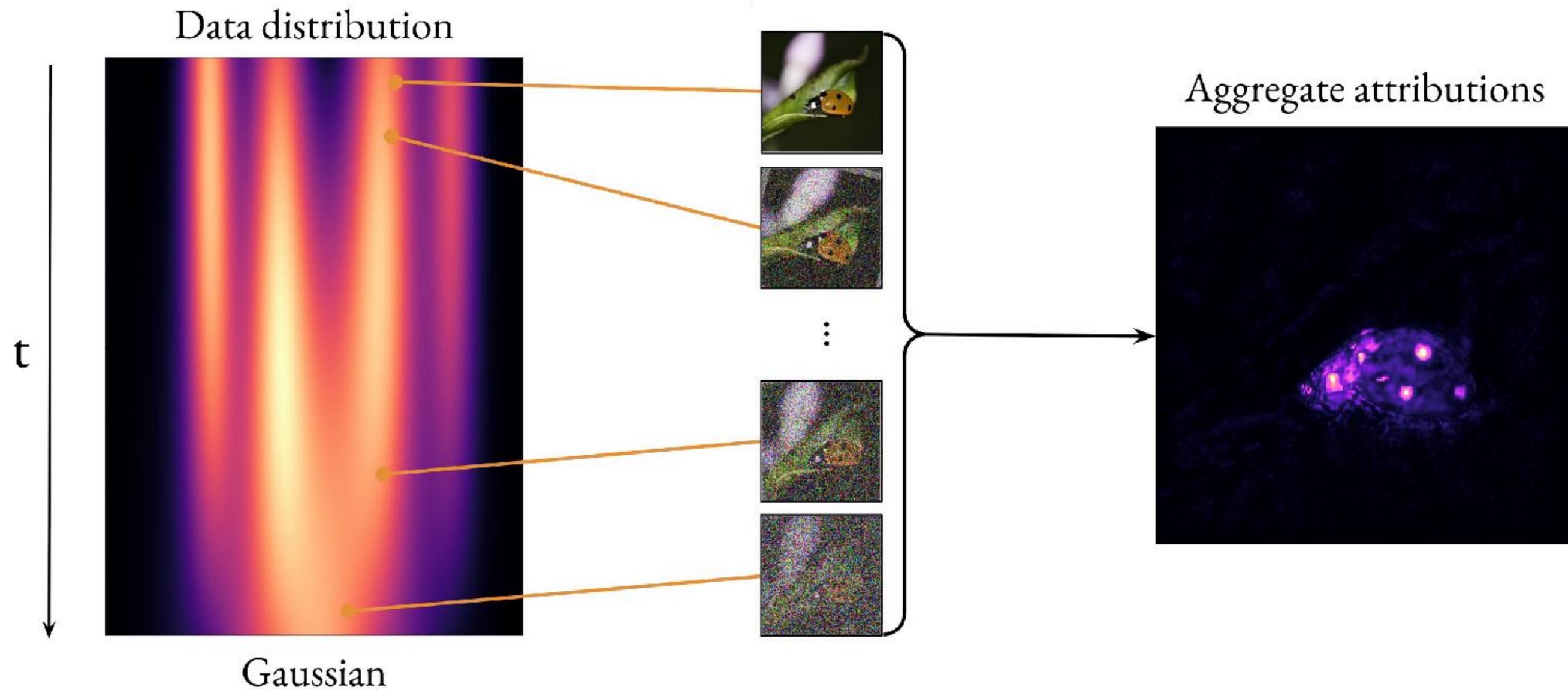
$$\mathbf{x}_t = \cos(t)\mathbf{x} + \sin(t)\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Final aggregation:

$$\begin{aligned} \tilde{\mathbf{C}}(\mathbf{x}) &= \mathbb{E}_{t,\boldsymbol{\epsilon}}[\mathbf{C}(\mathbf{x}_t)] = \\ &= \mathbb{E}_t[(\mathbf{C} * \mathbf{G}_t)(\cos(t)\mathbf{x})] \end{aligned}$$

DAVE - Better Understanding

Final pipeline: augment and interpolate noise over the path



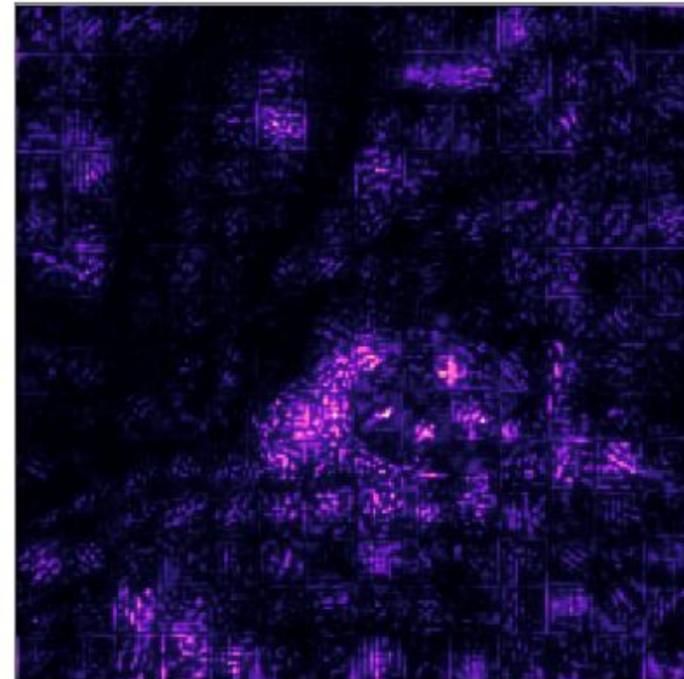
DAVE - Better Understanding

Final pipeline: augment and interpolate noise over the path

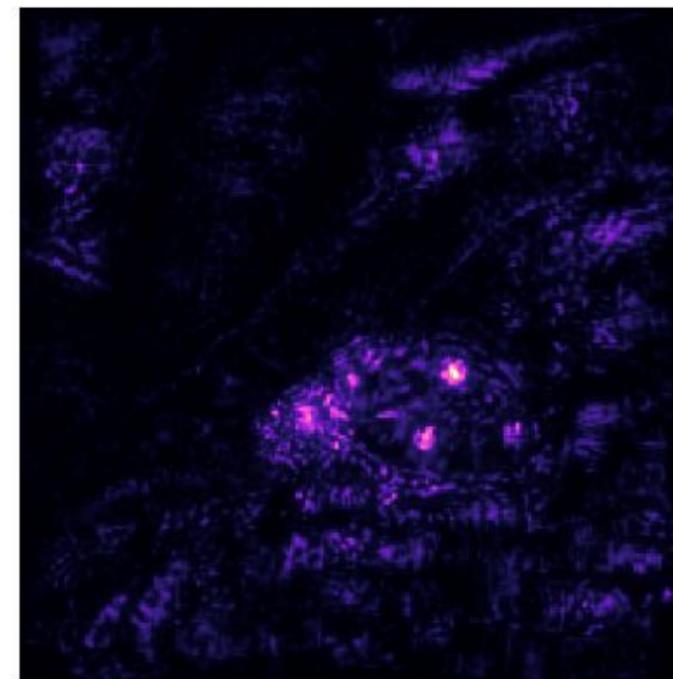
Input Image



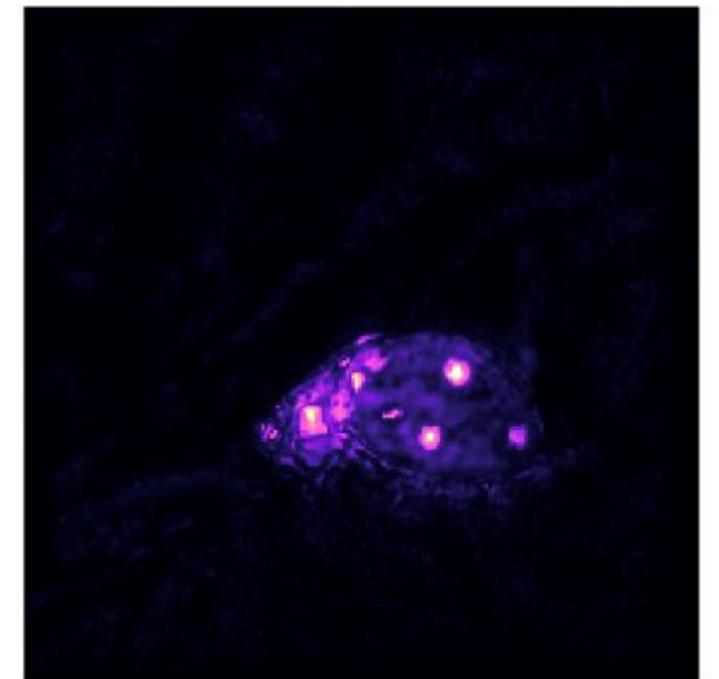
Single point



+ Input
Augmentations



+ Noise
interpolation path



DAVE - The Method

DAVE interprets attribution as the stable and locally equivariant effective transformation that a Vision Transformer applies to its input.

Under architectural assumptions of ViTs, it extracts this transformation from the input gradient by decomposing it into a direct input–output transformation and a local variation term.

The method discards the local variation term, which captures input-dependent sensitivity of internal model mechanisms, and aggregates the remaining transformation over a distribution of inputs.

This distribution is constructed to suppress locally non-equivariant components of the transformation, as well as high-frequency operator fluctuations, while preserving consistent inputdependent attribution structure.