

# Learning Discriminative Representations to Interpret Image Recognition Models

## Supplementary Material

Felipe Torres Figueroa

École Centrale de Marseille

QARMA - Laboratoire d'Informatique et de Systèmes (LIS)

Marseille, September 23<sup>rd</sup> 2024

# Table of Contents

- 1 Opti-CAM: Optimizing saliency maps for interpretability
- 2 CA-Stream: Attention-based pooling for interpretable image recognition
- 3 A learning paradigm for interpretable gradients

# Table of Contents

- 1 **Opti-CAM: Optimizing saliency maps for interpretability**
- 2 CA-Stream: Attention-based pooling for interpretable image recognition
- 3 A learning paradigm for interpretable gradients

# Evaluating Interpretability

## Interpretable Image Recognition

An explanation should demonstrate similar predictive properties to its query:

Input image. ( $I$ )



$$P_i^c = 0.8756$$

Explanation Map. ( $E^c$ )



$$O_i^c = 0.7442$$

$$AD(\%) := \frac{1}{N} \sum_{i=1}^N \frac{[y_i^c - o_i^c]_+}{y_i^c} \cdot 100 \quad (1)$$

$$AI(\%) := \frac{1}{N} \sum_i \mathbb{1}_{y_i^c < o_i^c} \cdot 100 \quad (2)$$

$$AG(\%) := \frac{1}{N} \sum_{i=1}^N \frac{[o_i^c - y_i^c]_+}{1 - y_i^c} \cdot 100 \quad (3)$$

# Interpretability

## Causality Analysis

Saliency guided perturbations reveal the importance of salient regions.

---

### Algorithm 1: Insertion Algorithm

---

**Input:** black-box  $f$ , image  $x$ , saliency map  $s^c$ , number of pixels  $N$  removed per step.

**Output:** insertion score  $ins$ .  $n \leftarrow 0$

$x' \leftarrow \text{Blur}(x)$

$p_n^c \leftarrow f(x)$

**while**  $x \neq x'$  **do**

    According to  $s$ , set the next  $n$  pixels in  $x'$  to corresponding pixels in  $x$

$n \leftarrow n + 1$

$p_n^c \leftarrow f(x')$

$ins \leftarrow \text{AreaUnderCurve}(p_n^c \text{ vs. } i/n, \forall i = 0, \dots, n)$

**return**  $ins$

---

# Interpretability

## Causality Analysis

Saliency guided perturbations reveal the importance of salient regions.

---

### Algorithm 2: Deletion Algorithm

---

**Input:** black-box  $f$ , image  $x$ , saliency map  $s^c$ , number of pixels  $N$  removed per step.

**Output:** deletion score  $del$ .

$n \leftarrow 0$

$p_n^c \leftarrow f(x)$

while  $x$  has non-zero pixels **do**

    According to  $s$ , set the next  $n$  pixels in  $x$  to 0

$n \leftarrow n + 1$

$p_n^c \leftarrow f(x)$

$del \leftarrow \text{AreaUnderCurve}(p_n^c \text{ vs. } i/n, \forall i = 0, \dots, n)$

**return**  $del$

---

# Interpretability

## Weakly Supervised Object Localization

$$\text{OM} := 1 - \left( \max_{B \in \mathbb{B}} \text{IoU}(B, B_p) \right) \mathbb{1}_{c_p=c}, \quad (4)$$

$$\text{LE} := 1 - \max_{B \in \mathbb{B}} \text{IoU}(B, B_p). \quad (5)$$

$$P := \frac{\sum_{\mathbf{p} \in U} S_{\mathbf{p}}^c}{\sum_{\mathbf{p}} S_{\mathbf{p}}^c} \quad (6)$$

$$R := \frac{\sum_{\mathbf{p} \in U} S_{\mathbf{p}}^c}{|U|}. \quad (7)$$

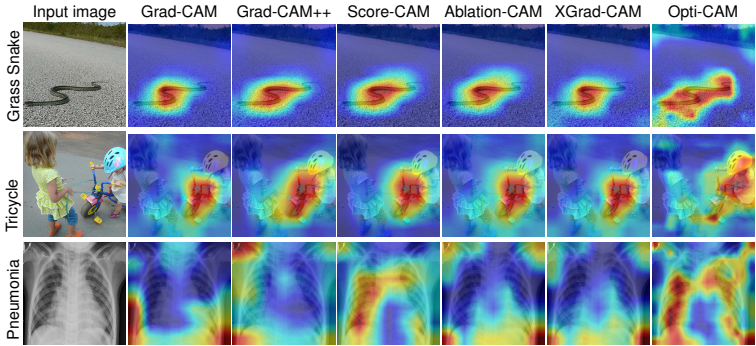
$$\text{BoxAcc}(\eta, \delta) := \max_{B \in \mathbb{B}} \mathbb{1}_{\text{IoU}(B_p^\eta, B) \geq \delta}. \quad (8)$$

$$\text{SP} := \mathbb{1}_{\mathbf{p}^* \in U}. \quad (9)$$

$$\text{SM} := \log \max \left( 0.05, \frac{|B_p|}{hw} \right) - \log p^c, \quad (10)$$

# Results

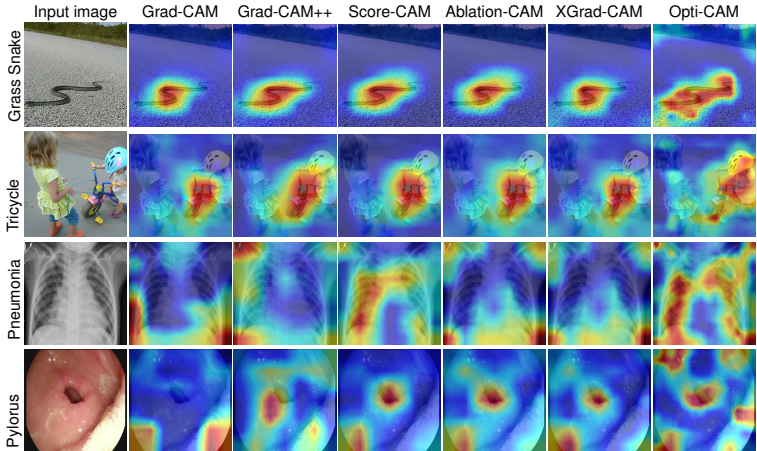
## Qualitative Evaluation





# Results

## Qualitative Evaluation



# Results

## Quantitative Experiments

### Classification Metrics on Transformers:

METHOD	ViT-B				DeiT-B			
	AD↓	AG↑	AI↑	T(s)	AD↓	AG↑	AI↑	T(s)
Fake-CAM	0.3	0.4	48.3	0.00	0.6	0.3	44.6	0.00
Grad-CAM	69.4	2.5	12.4	0.14	33.5	1.7	12.5	0.11
Grad-CAM	86.3	1.5	1.0	0.15	50.7	0.9	7.2	0.13
Score-CAM	32.0	6.2	33.0	23.69	53.6	2.2	12.2	22.47
XGrad-CAM	88.1	0.4	4.3	0.13	80.5	0.3	4.1	0.12
Layer-CAM	82.0	0.2	2.9	0.24	88.9	0.4	2.6	0.24
Experturbation	28.8	6.2	24.4	133.52	60.9	2.0	8.5	129.12
RawAtt	92.6	0.2	2.8	0.02	95.3	0.0	1.8	0.02
Rollout	42.1	5.6	20.9	0.02	55.2	0.8	7.9	0.02
TIBAV	81.7	0.8	5.8	0.16	62.3	0.7	7.1	0.16

# Results

## Quantitative Experiments

### Classification Metrics on Transformers:

METHOD	ViT-B				DeiT-B			
	AD↓	AG↑	AI↑	T(s)	AD↓	AG↑	AI↑	T(s)
Fake-CAM	0.3	0.4	48.3	0.00	0.6	0.3	44.6	0.00
Grad-CAM	69.4	2.5	12.4	0.14	33.5	1.7	12.5	0.11
Grad-CAM	86.3	1.5	1.0	0.15	50.7	0.9	7.2	0.13
Score-CAM	32.0	6.2	33.0	23.69	53.6	2.2	12.2	22.47
XGrad-CAM	88.1	0.4	4.3	0.13	80.5	0.3	4.1	0.12
Layer-CAM	82.0	0.2	2.9	0.24	88.9	0.4	2.6	0.24
Experturbation	28.8	6.2	24.4	133.52	60.9	2.0	8.5	129.12
RawAtt	92.6	0.2	2.8	0.02	95.3	0.0	1.8	0.02
Rollout	42.1	5.6	20.9	0.02	55.2	0.8	7.9	0.02
TIBAV	81.7	0.8	5.8	0.16	62.3	0.7	7.1	0.16
Opti-CAM	<b>0.6</b>	<b>18.0</b>	<b>90.1</b>	16.05	<b>0.9</b>	<b>26.0</b>	<b>83.5</b>	15.17

# Results

## Quantitative Evaluation

### Localization Experiments:

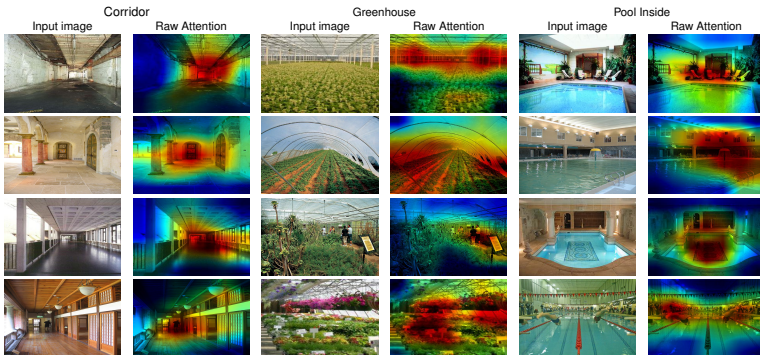
METHOD	ViT-B							DeiT-B						
	OM↓	LE↓	F1↑	BA↑	SP↑	EP↑	SM↓	OM↓	LE↓	F1↑	BA↑	SP↑	EP↑	SM↓
Fake-CAM	62.8	54.0	57.7	47.9	99.8	28.6	0.87	61.4	54.0	57.7	47.9	99.8	28.7	0.83
Grad-CAM	79.6	74.3	29.4	45.0	58.1	31.0	3.27	65.5	60.3	44.3	47.2	62.8	30.2	1.20
Grad-CAM++	84.2	80.6	14.8	23.8	51.4	27.3	4.15	70.6	67.2	34.3	43.6	57.7	30.3	2.14
Score-CAM	77.6	71.6	46.0	54.3	<b>66.1</b>	33.1	3.14	79.9	76.2	31.9	43.8	<b>63.4</b>	32.2	3.14
XGrad-CAM	82.0	76.9	19.6	41.3	52.8	28.5	3.31	82.0	78.4	19.5	44.1	53.4	28.8	3.03
Layer-CAM	70.7	63.9	20.6	50.5	60.7	32.6	1.44	80.2	77.3	17.6	50.8	62.7	35.1	3.15
ExPerturbation	71.5	64.9	35.9	44.6	62.3	<b>35.3</b>	1.34	69.9	64.3	36.2	44.2	63.1	<b>35.5</b>	1.16
RawAtt	72.4	64.8	18.5	50.4	55.4	31.6	1.68	73.5	68.2	5.9	<b>48.1</b>	46.5	27.3	1.91
Rollout	67.6	58.8	36.9	<b>50.7</b>	57.8	30.0	1.16	63.9	57.0	27.8	47.9	36.5	27.2	0.94
TIBAV	70.1	63.1	26.6	<b>58.8</b>	<b>66.1</b>	35.0	1.23	68.2	62.2	28.1	59.6	64.1	33.5	1.08
Opti-CAM (ours)	<b>64.4</b>	<b>54.6</b>	<b>54.5</b>	48.0	58.2	28.7	<b>0.98</b>	<b>62.3</b>	<b>55.1</b>	<b>53.9</b>	48.0	55.1	28.8	<b>0.84</b>

# Table of Contents

- 1 Opti-CAM: Optimizing saliency maps for interpretability
- 2 CA-Stream: Attention-based pooling for interpretable image recognition**
- 3 A learning paradigm for interpretable gradients

# Results

## Qualitative Experiments



# Results

## Quantitative Experiments

NETWORK	METHOD	POOL	AD↓	AG↑	AI↑	I↑	D↓
RESNET-18	Grad-CAM	GAP	17.64	12.73	41.21	63.13	<b>10.66</b>
		CA	<b>16.99</b>	<b>17.22</b>	<b>44.95</b>	<b>65.94</b>	10.68
	Grad-CAM++	GAP	19.05	11.16	37.99	62.80	<b>10.75</b>
		CA	<b>19.02</b>	<b>14.76</b>	<b>40.82</b>	<b>65.53</b>	10.82
	Score-CAM	GAP	13.64	12.98	44.53	62.56	<b>11.37</b>
		CA	<b>11.53</b>	<b>18.12</b>	<b>50.32</b>	<b>65.33</b>	11.51
CONVNEXT-S	Grad-CAM	GAP	42.99	1.69	12.60	48.42	<b>30.12</b>
		CA	<b>22.09</b>	<b>14.91</b>	<b>32.65</b>	<b>84.82</b>	43.02
	Grad-CAM++	GAP	56.42	1.32	10.35	48.28	<b>33.41</b>
		CA	<b>51.87</b>	<b>9.40</b>	<b>20.55</b>	<b>84.28</b>	52.58
	Score-CAM	GAP	74.79	1.29	10.10	47.40	<b>38.21</b>
		CA	<b>64.21</b>	<b>8.81</b>	<b>18.96</b>	<b>82.92</b>	57.46

# Results

## Quantitative Experiments

### PASCAL VOC 2012 - RESNET-50

POOLING

MAP $\uparrow$ 

GAP

78.32

CA

78.35

### INTERPRETABILITY METRICS

METHOD	POOLING	AD $\downarrow$	AG $\uparrow$	AI $\uparrow$	I $\uparrow$	D $\downarrow$
Grad-CAM	GAP	<b>12.61</b>	9.68	27.88	<b>89.10</b>	59.39
	CA	12.77	<b>15.46</b>	<b>34.53</b>	88.53	<b>59.16</b>
Grad-CAM++	GAP	<b>12.25</b>	9.68	27.62	<b>89.34</b>	54.23
	CA	12.28	<b>16.76</b>	<b>34.87</b>	89.02	<b>53.34</b>
Score-CAM	GAP	14.8	6.76	36.41	71.10	<b>39.95</b>
	CA	<b>10.96</b>	<b>21.35</b>	<b>43.82</b>	<b>89.21</b>	51.44



# Results

## Quantitative Experiments

### CUB-200-2011 - RESNET-50

POOLING

Acc↑

GAP

76.96

CA

75.90

### INTERPRETABILITY METRICS

METHOD	POOLING	AD↓	AG↑	AI↑	I↑	D↓
Grad-CAM	GAP	10.87	10.29	45.81	65.71	<b>6.17</b>
	CA	<b>10.44</b>	<b>17.61</b>	<b>53.54</b>	<b>74.60</b>	6.56
Grad-CAM++	GAP	11.35	9.68	44.32	65.64	<b>5.92</b>
	CA	<b>11.01</b>	<b>16.50</b>	<b>51.63</b>	<b>74.64</b>	6.21
Score-CAM	GAP	9.05	10.62	48.90	65.58	5.94
	CA	<b>6.37</b>	<b>19.50</b>	<b>60.41</b>	<b>74.22</b>	<b>2.14</b>

# Results

## Quantitative Experiments

ACCURACY AND PARAMETERS			
PLACEMENT	CLS DIM	#PARAM	ACC↑
$S_0 - S_4$	64	6.96M	<b>74.70</b>
$S_1 - S_4$	256	6.95M	74.67
$S_2 - S_4$	512	6.82M	74.67
$S_3 - S_4$	1024	6.29M	74.67
$S_4 - S_4$	2048	4.20M	74.63

# Results

## Quantitative Experiments

INTERPRETABILITY METRICS						
METHOD	PLACEMENT	AD↓	AG↑	AI↑	I↑	D↓
GRAD-CAM	$S_0 - S_4$	<b>12.54</b>	<b>22.67</b>	48.56	75.53	13.50
	$S_1 - S_4$	12.69	22.65	48.31	75.53	13.41
	$S_2 - S_4$	<b>12.54</b>	21.67	<b>48.58</b>	75.54	13.50
	$S_3 - S_4$	12.69	22.28	47.89	<b>75.55</b>	13.40
	$S_4 - S_4$	12.77	20.65	47.14	74.32	<b>13.37</b>
GRAD-CAM++	$S_0 - S_4$	13.99	19.29	44.60	75.21	13.78
	$S_1 - S_4$	13.99	19.29	44.62	75.21	13.78
	$S_2 - S_4$	13.71	<b>19.90</b>	<b>45.43</b>	75.34	13.50
	$S_3 - S_4$	13.69	19.61	45.04	<b>75.36</b>	13.50
	$S_4 - S_4$	<b>13.67</b>	18.36	44.40	74.19	<b>13.30</b>
SCORE-CAM	$S_0 - S_4$	<b>7.09</b>	23.65	54.20	74.91	14.68
	$S_1 - S_4$	<b>7.09</b>	23.65	54.20	74.92	14.68
	$S_2 - S_4$	<b>7.09</b>	<b>23.66</b>	<b>54.21</b>	74.91	14.68
	$S_3 - S_4$	7.74	23.03	52.92	<b>74.97</b>	14.65
	$S_4 - S_4$	7.52	19.45	50.45	74.19	<b>14.46</b>

# Results

## Quantitative Experiments

ACCURACY AND PARAMETERS						
	REPRESENTATION		#PARAM		ACC↑	
	Class agnostic		32.53M		74.70	
	Class specific		32.59M		74.68	
INTERPRETABILITY METRICS						
METHOD	REPRESENTATION	AD↓	AG↑	AI↑	I↑	D↓
Grad-CAM	Class agnostic	12.54	22.67	48.56	75.53	13.50
	Class specific	12.53	22.66	48.58	75.54	13.50
Grad-CAM++	Class agnostic	13.99	19.29	44.60	75.21	13.78
	Class specific	13.99	19.28	44.62	75.20	13.78
Score-CAM	Class agnostic	7.09	23.65	54.20	74.91	14.68
	Class specific	7.08	23.64	54.15	74.99	14.53

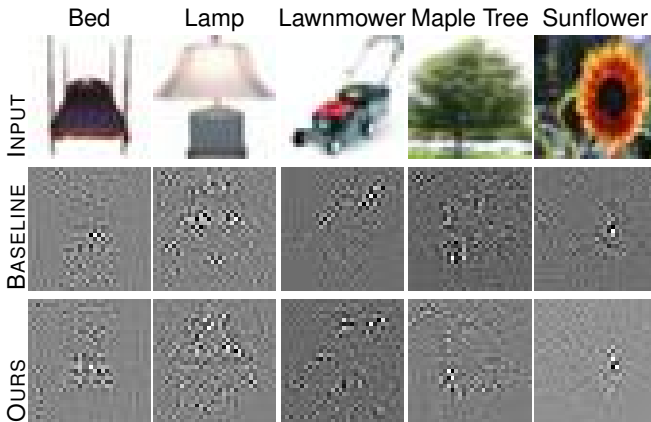
# Table of Contents

- 1 Opti-CAM: Optimizing saliency maps for interpretability
- 2 CA-Stream: Attention-based pooling for interpretable image recognition
- 3 A learning paradigm for interpretable gradients**

# Results

## Qualitative Experiments

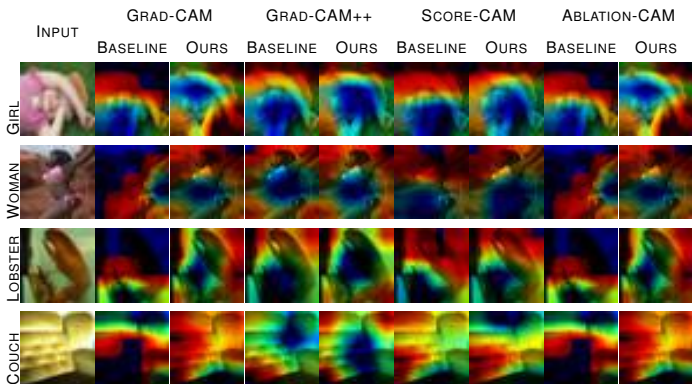
### Denoising effect



# Results

## Qualitative Experiments

### CAM Visualizations



# Results

## Quantitative Experiments

### Recognition Metrics

RECOGNITION METRICS			
MODEL	ERROR	$\lambda$	ACC $\uparrow$
RESNET-18	-	-	<b>73.42</b>
	COSINE	$7.5 \times 10^{-3}$	72.86
MOBILENET-V2	-	-	59.43
	COSINE	$1 \times 10^{-3}$	<b>62.36</b>



# Results

## Quantitative Experiments

### Recognition Metrics

RECOGNITION METRICS			
MODEL	ERROR	$\lambda$	ACC $\uparrow$
RESNET-18	-	-	<b>73.42</b>
	COSINE	$7.5 \times 10^{-3}$	72.86
MOBILENET-V2	-	-	59.43
	COSINE	$1 \times 10^{-3}$	<b>62.36</b>

Recognition properties are maintained.

# Results

## Quantitative Experiments

### Interpretability metrics

MOBILENET-V2						
METHOD	ERROR	AD↓	AG↑	AI↑	INS↑	DEL↓
GRAD-CAM	-	44.64	6.57	25.62	44.64	<b>14.34</b>
	COSINE	<b>40.89</b>	<b>7.31</b>	<b>27.08</b>	<b>45.57</b>	15.20
GRAD-CAM++	-	45.98	6.12	24.10	44.72	<b>14.76</b>
	COSINE	<b>40.76</b>	<b>6.85</b>	<b>26.46</b>	<b>45.51</b>	14.92
SCORE-CAM	-	40.55	7.85	28.57	45.62	<b>14.52</b>
	COSINE	<b>36.34</b>	<b>9.09</b>	<b>30.50</b>	<b>46.35</b>	14.72
ABLATION-CAM	-	45.15	6.38	25.32	44.62	<b>15.03</b>
	COSINE	<b>41.13</b>	<b>7.03</b>	<b>26.10</b>	<b>45.38</b>	15.12
AXIOM-CAM	-	44.65	6.57	25.62	44.64	15.27
	COSINE	<b>40.89</b>	<b>7.31</b>	<b>27.08</b>	<b>45.57</b>	<b>15.20</b>

# Results

## Quantitative Experiments

### Interpretability metrics

MOBILENET-V2						
METHOD	ERROR	AD↓	AG↑	AI↑	INS↑	DEL↓
GRAD-CAM	-	44.64	6.57	25.62	44.64	<b>14.34</b>
	COSINE	<b>40.89</b>	<b>7.31</b>	<b>27.08</b>	<b>45.57</b>	15.20
GRAD-CAM++	-	45.98	6.12	24.10	44.72	<b>14.76</b>
	COSINE	<b>40.76</b>	<b>6.85</b>	<b>26.46</b>	<b>45.51</b>	14.92
SCORE-CAM	-	40.55	7.85	28.57	45.62	<b>14.52</b>
	COSINE	<b>36.34</b>	<b>9.09</b>	<b>30.50</b>	<b>46.35</b>	14.72
ABLATION-CAM	-	45.15	6.38	25.32	44.62	<b>15.03</b>
	COSINE	<b>41.13</b>	<b>7.03</b>	<b>26.10</b>	<b>45.38</b>	15.12
AXIOM-CAM	-	44.65	6.57	25.62	44.64	15.27
	COSINE	<b>40.89</b>	<b>7.31</b>	<b>27.08</b>	<b>45.57</b>	<b>15.20</b>

Interpretable properties are enhanced. Deletion still poses an issue.

# Table of Contents

- 1 Opti-CAM: Optimizing saliency maps for interpretability
- 2 CA-Stream: Attention-based pooling for interpretable image recognition
- 3 A learning paradigm for interpretable gradients

