# Directions in Interpretability

Ruth Fong

HEIBRiDS lecture

November 14, 2022

Slides and links available at ruthfong.com

PRINCETON
UNIVERSITY

# What is interpretability?

Research focused on explaining **complex AI systems** in a **human-interpretable** way.

# Why interpretability?

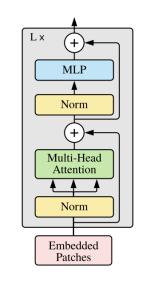- 🔬 Science

- 🤝 Trust

- 🤖 Learning

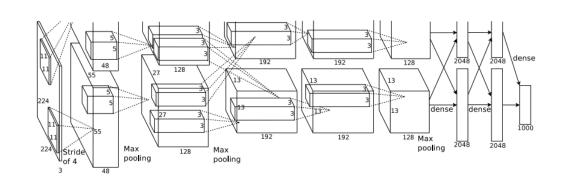# An incomplete retrospective: the first decade of deep learning



Monet → photo

**GANs (2014–2018)**
GAN, ProGAN, CycleGAN



**Transformers (2017–now)**
Transformer, BERT, ViT

IM GENET

2012 ————————————————————————————————→ 2022



**CNNs (2012–2016)**
AlexNet, VGG16,
GoogLeNet, ResNet50



**Self-supervised learning (2016–now)**
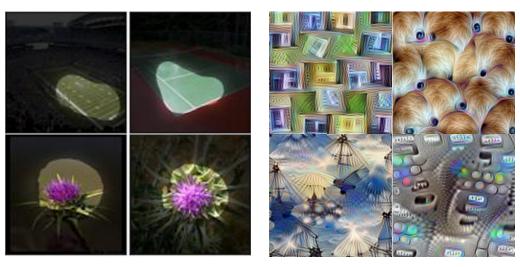Colorization, MOCO, SWaV



**Diffusion models (2020–now)**
DDPM, DALL-E 2, Imagen

4

[Krizhevsky et al., NeurIPS 2012; Zhu* & Park* et al., ICCV 2017; Zhang et al., ECCV 2016;
Dosovitskiy* et al., ICLR 2021; Ramesh et al., arXiv 2022]

# An incomplete retrospective: the first decade of interpretability



**Feature visualization (2013–2018)**
Activation Max., Feature Inversion,
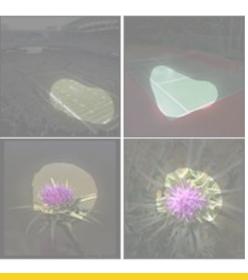Net Dissect, Feature Vis.

**Attribution heatmaps (2013–2019)**
Gradient, Grad-CAM,
Occlusion, Perturbations, RISE

**Interpretable-by-design (2020-now)**
Concept Bottleneck, ProtoPNet,
ProtoTree

2022

[Selvaraju et al., ICCV 2017; Fong* & Patrick* et al., ICCV 2019;
Bau* & Zhou* et al., CVPR 2017; Olah et al., Distill 2017; Koh*, Nguyen*, Tang* et al., ICML 2020]

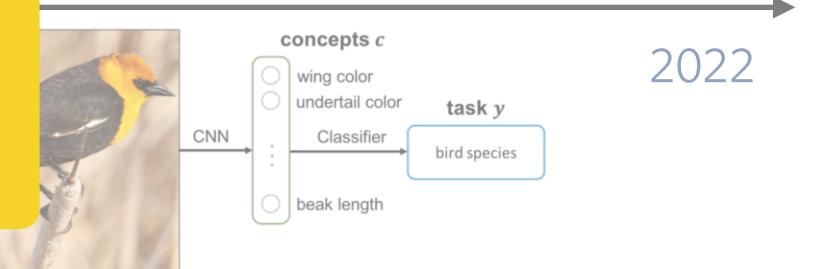# An incomplete retrospective: the first decade of interpretability

Primarily focused on understanding and approximating **CNNs**

*Exceptions:*
*GANPaint [Bau et al., ICLR 2019]*
*Transformer Circuits [Elhage et al., 2021]*

2022

Orig Img    Mask

concepts $c$
wing color
undertail color
CNN    Classifier    task $y$
bird species
beak length

**Attribution heatmaps (2013-2019)**
Gradient, Grad-CAM,
Occlusion, Perturbations, RISE

**Interpretable-by-design (2020-now)**
Concept Bottleneck, ProtoPNet,
ProtoTree

[Selvaraju et al., ICCV 2017; Fong* & Patrick* et al., ICCV 2019;
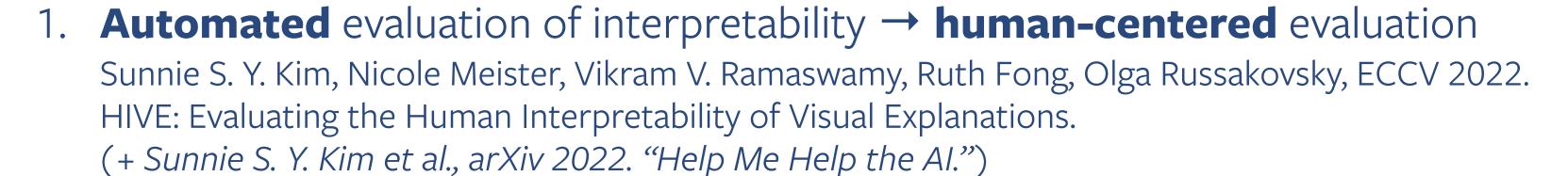Bau* & Zhou* et al., CVPR 2017; Olah et al., Distill 2017; Koh*, Nguyen*, Tang* et al., ICML 2020]

# Directions for the next decade of interpretability

1. Develop interpretability methods for **diverse domains**

   - Beyond CNN classifiers: self-supervised learning, generative models, etc.

2. Center **humans** throughout the development process

   - In design, co-develop methods with real-world stakeholders.

   - In evaluation, measure human interpretability and utility of methods.

   - In deployment, package interpretability tools for the wider community.

# Roadmap

1. **Automated** evaluation of interpretability → **human-centered** evaluation
   Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, Olga Russakovsky, ECCV 2022.
   HIVE: Evaluating the Human Interpretability of Visual Explanations.
   (+ *Sunnie S. Y. Kim et al., arXiv 2022. "Help Me Help the AI."*)

2. Explanations via **labelled attributes** → explanations via **labelled attributes and unlabelled features**
   Vikram V. Ramaswamy, Sunnie S. Y. Kim, Nicole Meister, Ruth Fong, Olga Russakovsky, arXiv 2022.
   ELUDE: Generating Interpretable Explanations via a Decomposition into Labelled and Unlabelled Features.
   (+ *Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors in Concept-based Explanations.*)

3. Interpretability of **supervised** models → interpretability of **self-supervised** models
   Iro Laina, Ruth Fong, Andrea Vedaldi, NeurIPS 2020.
   Quantifying Learnability and Describability of Visual Concepts Emerging in Representation Learning.

4. **Interpretability** in ML + CV → **interdisciplinary** research (interpretability + X)
   (+ *Nicole Meister* and Dora Zhao* et al., arXiv 2022. Gender Artifacts in Visual Datasets.*)
   (+ *Indu Panigrahi et al., arXiv 2022. Improving Fine-Grain Segmentation via Interpretable Modifications.*)

5. **Static** visualizations → **interactive** visualizations
   Ruth Fong, Alexander Mordvintsev, Andrea Vedaldi, Chris Olah, VISxAI 2021.
   Interactive Similarity Overlays.
   (+ *Devon Ulrich and Ruth Fong, in prep. Interactive Visual Feature Search.*)
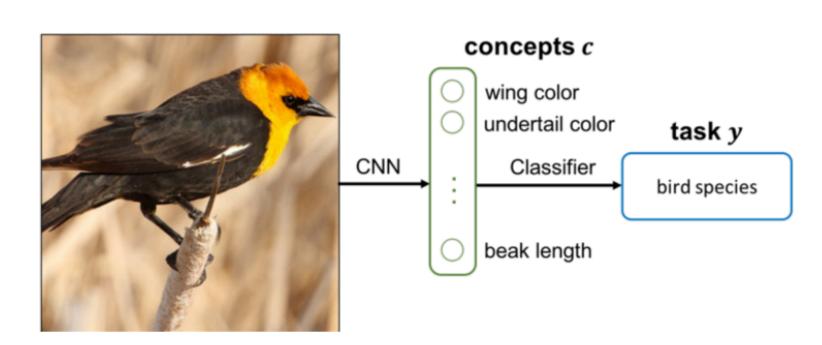
# Roadmap



Sunnie S. Y. Kim

1. **Automated** evaluation of interpretability → **human-centered** evaluation
   Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, Olga Russakovsky, ECCV 2022.
   HIVE: Evaluating the Human Interpretability of Visual Explanations.
   (+ *Sunnie S. Y. Kim et al., arXiv 2022. "Help Me Help the AI."*)

2. Explanations via **labelled attributes** → explanations via **labelled attributes and unlabelled features**
   Vikram V. Ramaswamy, Sunnie S. Y. Kim, Nicole Meister, Ruth Fong, Olga Russakovsky, arXiv 2022.
   ELUDE: Generating Interpretable Explanations via a Decomposition into Labelled and Unlabelled Features.
   (+ *Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors in Concept-based Explanations.*)

3. Interpretability of **supervised** models → interpretability of **self-supervised** models
   Iro Laina, Ruth Fong, Andrea Vedaldi, NeurIPS 2020.
   Quantifying Learnability and Describability of Visual Concepts Emerging in Representation Learning.

4. **Interpretability** in ML + CV → **interdisciplinary** research (interpretability + X)
   (+ *Nicole Meister\* and Dora Zhao\* et al., arXiv 2022. Gender Artifacts in Visual Datasets.*)
   (+ *Indu Panigrahi et al., arXiv 2022. Improving Fine-Grain Segmentation via Interpretable Modifications.*)

5. **Static** visualizations → **interactive** visualizations
   Ruth Fong, Alexander Mordvintsev, Andrea Vedaldi, Chris Olah, VISxAI 2021.
   Interactive Similarity Overlays.
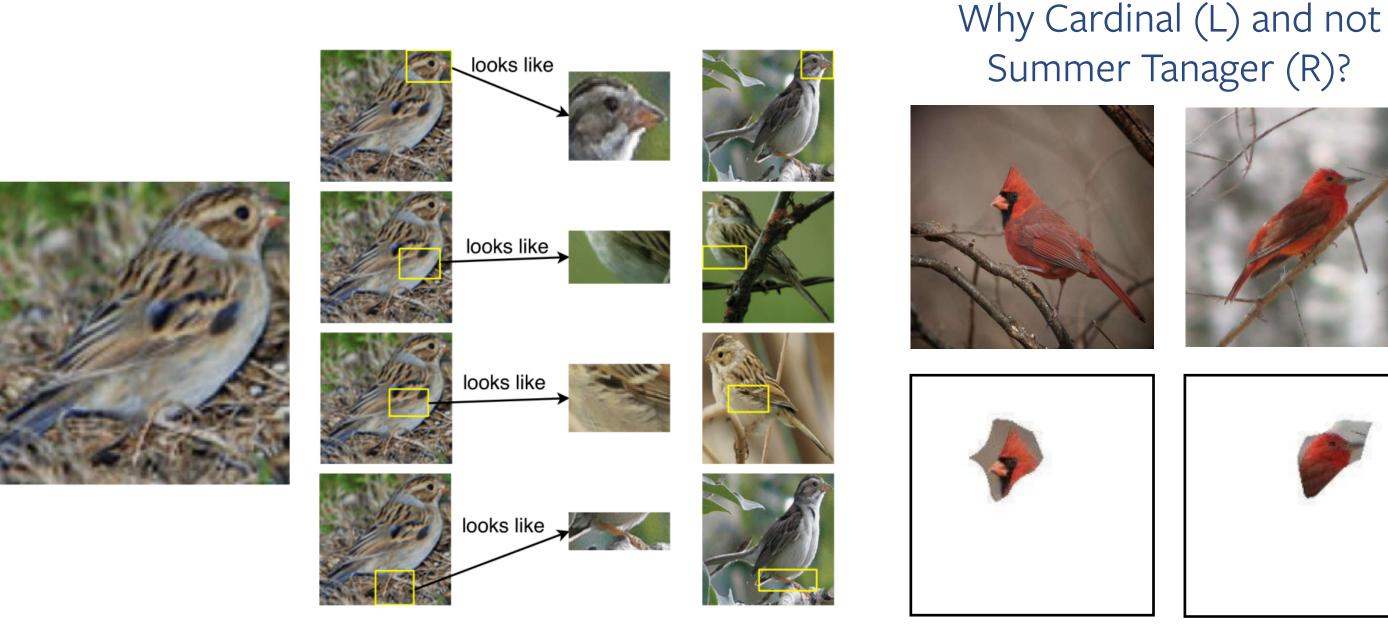   (+ *Devon Ulrich and Ruth Fong, in prep. Interactive Visual Feature Search.*)

# Explanation form factors: Why did the model predict Y?



**Heatmap** explanations
(e.g. Grad-CAM)



**Concept**-based explanations
(e.g. Concept Bottleneck)



**Prototype** explanations
(e.g. ProtoPNet)

Why Cardinal (L) and not
Summer Tanager (R)?



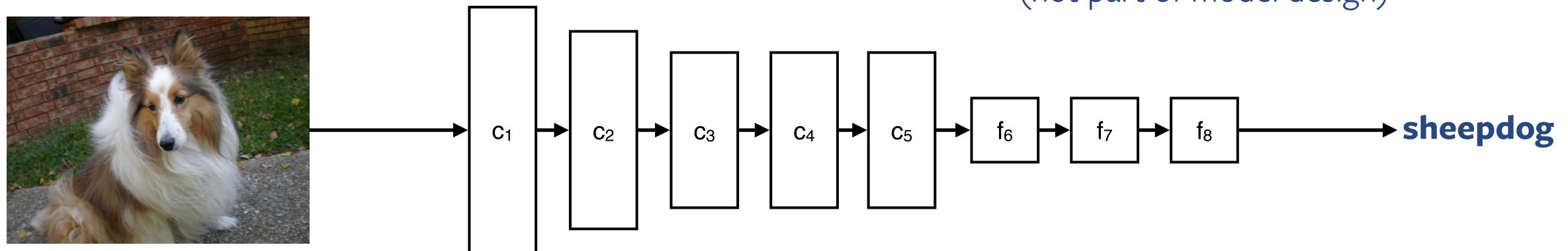**Counterfactual** explanations
(e.g. SCOUT)

[Selvaraju et al., ICCV 2017; Koh*, Nguyen*, Tang* et al., ICML 2020;
Chen* & Li* et al., NeurIPS 2019; Wang & Vasconcelos, CVPR 2020]
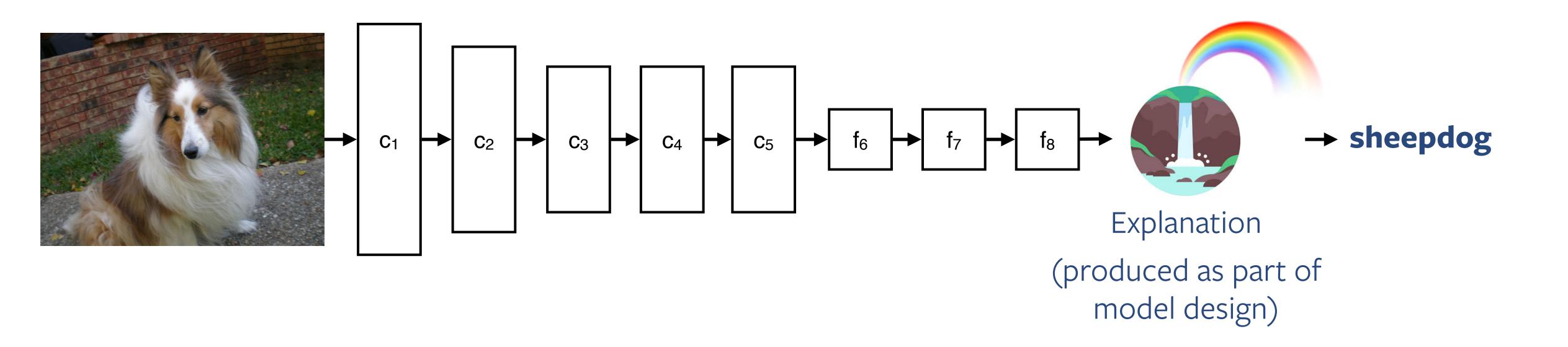
# Post-hoc explanations



Explanation
(not part of model design)

$c_1$ → $c_2$ → $c_3$ → $c_4$ → $c_5$ → $f_6$ → $f_7$ → $f_8$ → **sheepdog**

# Interpretable-by-design models



$c_1$ → $c_2$ → $c_3$ → $c_4$ → $c_5$ → $f_6$ → $f_7$ → $f_8$ →

Explanation
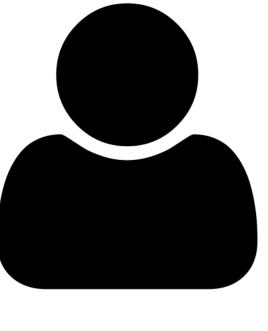(produced as part of
model design)

→ **sheepdog**

# Current metrics focus on heatmap evaluation

- Weak localization performance [Zhang et al., ECCV 2016]
- Perturbation analysis
  - Deletion game [Samek et al., TNNLS 2017]
  - Retrain with removed features [Hooker et al., NeurIPS 2019]
- Sensitivity to...
  - output neuron [Rebuffi*, Fong*, Ji* et al., CVPR 2020]
  - model parameters [Adebayo et al., NeurIPS 2018]
- ...

- Sheng & Huang, HCOMP 2020
  Guess the incorrectly predicted label
- Nguyen et al., NeurIPS 2021
  Is this prediction correct?
- Colin* & Fel* et al., arXiv 2021
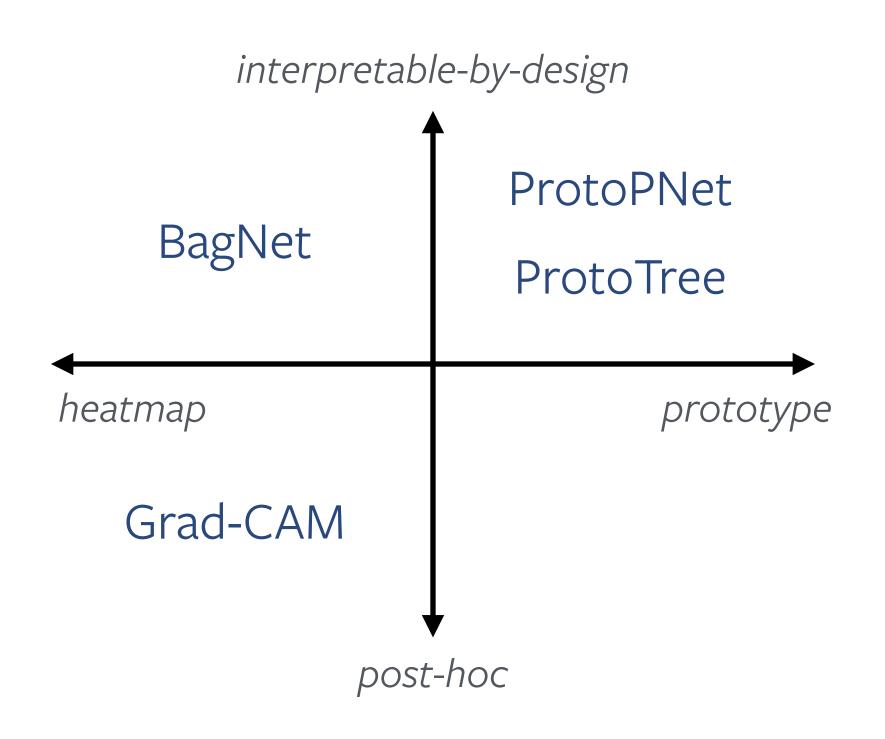  What did the model predict (choose one of two)?

Automatic

Human

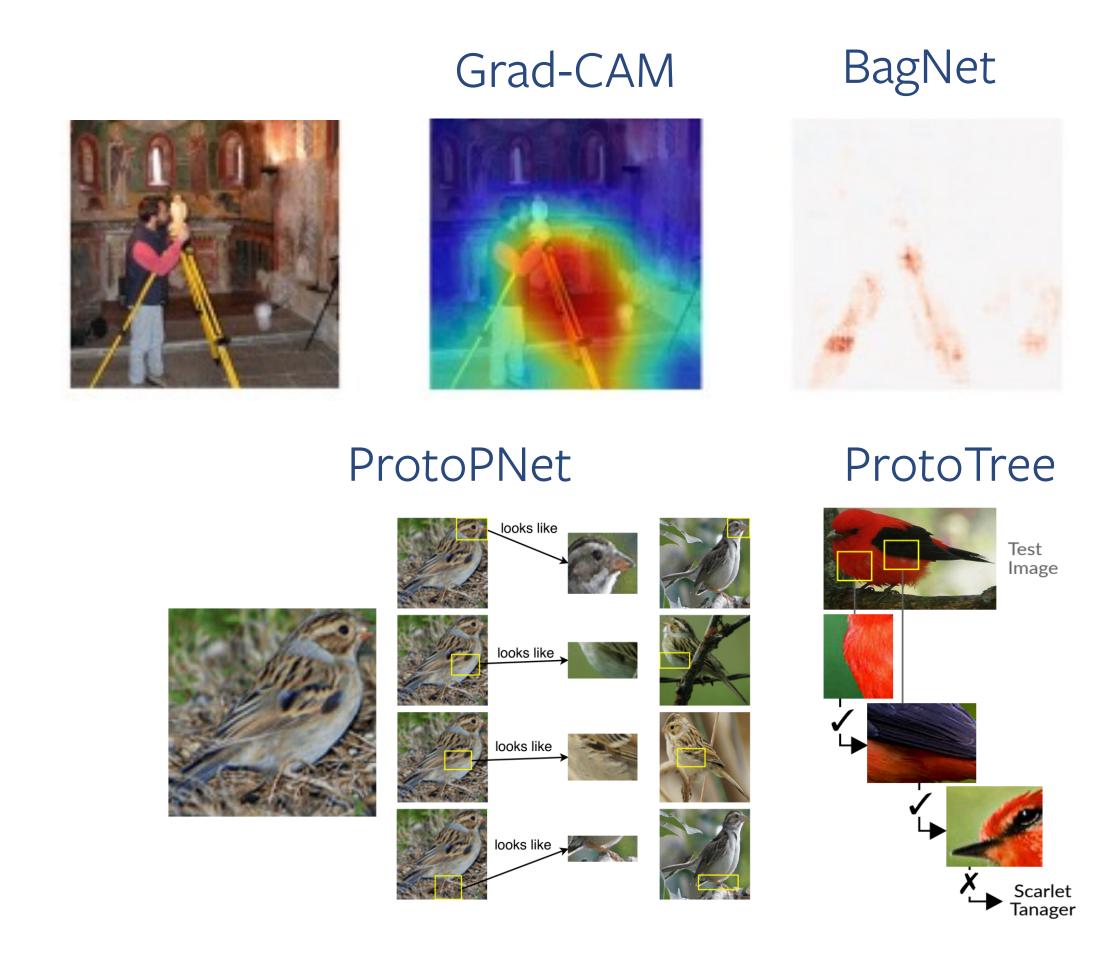# HIVE: Evaluating the Human Interpretability of Visual Explanations

1. Within method → **Cross-method comparison**

2. Automated evaluation → **Human-centered evaluation**

3. Intuition-based reasoning → **Falsifiable hypothesis testing**

# Our contributions

- Novel human study design for evaluating 4 diverse interpretability methods
  - **First human study** for interpretable-by-design and prototype methods
- Quantify the utility of explanations in distinguishing between **correct and incorrect predictions**
- Quantify how users would trade off between **interpretability and accuracy**
- **Open-source** HIVE studies to encourage reproducible research

# 1. Cross-method comparison



interpretable-by-design

ProtoPNet

BagNet

ProtoTree

heatmap                    prototype

Grad-CAM

post-hoc

Grad-CAM          BagNet



ProtoPNet                  ProtoTree



[Selvaraji et al., ICCV 2017; Brendel & Bethge, ICLR 2019; Chen* & Li* et al., NeurIPS 2019, Nauta et al., CVPR 2021]

# 2. Human-centered evaluation
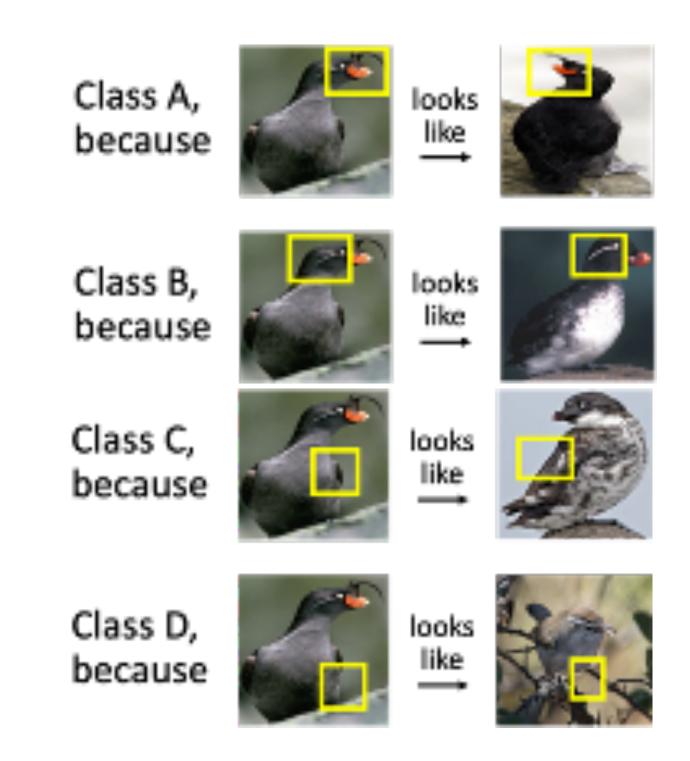
## Agreement task

How confident are you in the model's prediction?



*Experimental set-up: AMT studies with N=50 participants each*

## Distinction task

Which class do you think is correct?



[Sunnie S. Y. Kim et al., ECCV 2022. HIVE.; Chen* & Li* et al., NeurIPS 2019] 17

# 2. Human-centered evaluation

**Agreement task**

How confident are you in the model's prediction?

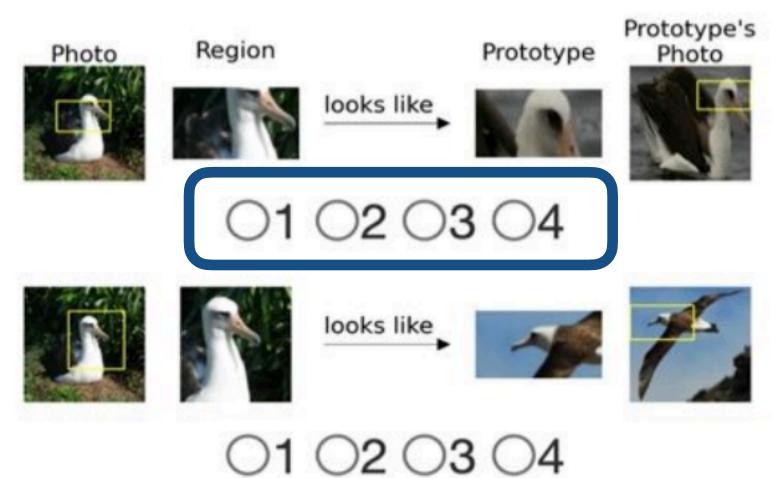**Finding #1:** Prototype similarities often **do not align** with human notions of similarity.

**Task: Rate the similarity of each row's prototype-region pair on a scale of 1-4.**

(1: Not Similar, 2: Somewhat Not Similar, 3: Somewhat Similar, 4: Similar)

Shown below is the model's explanation for its prediction (all prototypes and their source photos are from **Species 2**).

Photo    Region    looks like    Prototype    Prototype's Photo

○1 ○2 ○3 ○4

○1 ○2 ○3 ○4

**Q. What do you think about the model's prediction?**

○ Fairly confident that prediction is *correct*
○ Somewhat confident that prediction is *correct*
○ Somewhat confident that prediction is *incorrect*
○ Fairly confident that prediction is *incorrect*

[Sunnie S. Y. Kim et al., ECCV 2022. HIVE.; Chen* & Li* et al., NeurIPS 2019] 18

# 2. Human-centered evaluation

**Agreement task**

How confident are you in the model's prediction?

**Finding #1:** Prototype similarities often **do not align** with human notions of similarity.

**Finding #2:** Agreement task reveals **confirmation bias**.
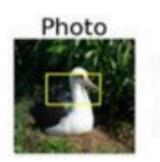
**More than 50%** were fairly or somewhat confident that a prediction is correct (even for incorrect predictions).

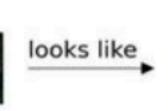**Task: Rate the similarity of each row's prototype-region pair on a scale of 1-4.**

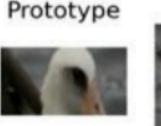(1: Not Similar, 2: Somewhat Not Similar, 3: Somewhat Similar, 4: Similar)

Shown below is the model's explanation for its prediction (all prototypes and their source photos are from **Species 2**).
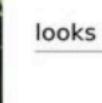
Photo  Region  looks like  Prototype  Prototype's Photo

○1 ○2 ○3 ○4

looks like

○1 ○2 ○3 ○4

**Q. What do you think about the model's prediction?**
☑ Fairly confident that prediction is *correct*
☑ Somewhat confident that prediction is *correct*
○ Somewhat confident that prediction is <u>incorrect</u>
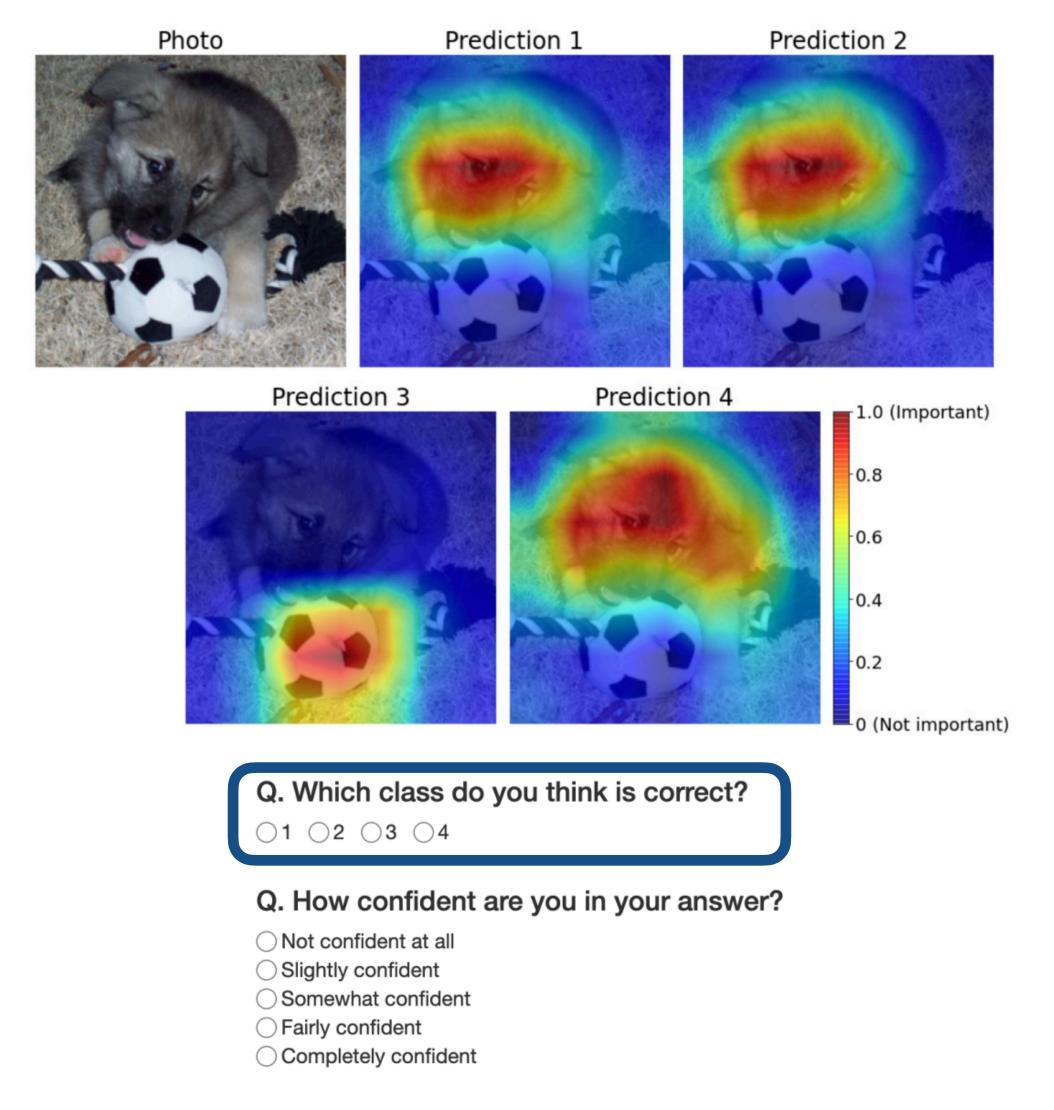○ Fairly confident that prediction is <u>incorrect</u>

[Sunnie S. Y. Kim et al., ECCV 2022. HIVE.; Chen* & Li* et al., NeurIPS 2019] 19

# 2. Human-centered evaluation

**Distinction task**

Which class do you think is correct?

**Finding #3:** Participants struggle to identify the **correct class**, esp. for incorrect predictions.

For incorrect predictions, correctly answered around 25% of the time (**random guessing**).

**Goal:** Interpretability should help humans identify and explain model errors.



[Sunnie S. Y. Kim et al., ECCV 2022. HIVE.; Selvaraju et al., ICCV 2017] 20

# 3. Falsifiable hypothesis testing

**Finding #1:** Prototype similarities often **do not align** with human notions of similarity.

**Finding #2:** Agreement task reveals **confirmation bias**.

**Finding #3:** Participants struggle to identify the **correct class**, esp. for incorrect predictions.

# 3. Falsifiable hypothesis testing

**Finding #1:** Prototype similarities often **do not align** with human notions of similarity.

**Finding #2:** Agreement task reveals **confirmation bias**.

**Finding #3:** Participants struggle to identify the **correct class**, esp. for incorrect predictions.
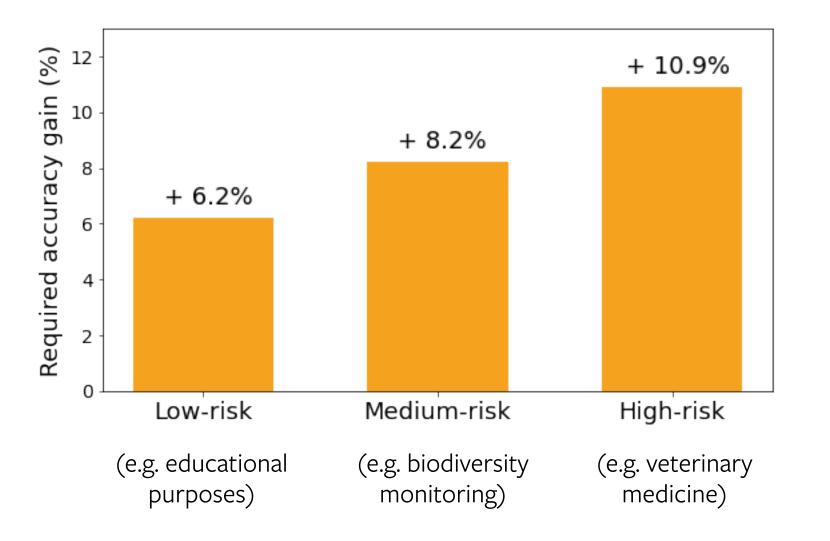
**Finding #4:** Participants prefer interpretability over accuracy, esp. in high-risk settings.

**Follow up: Kim et al., arXiv 2022.**
"Help Me Help the AI": Understanding How Explainability Can Support Human-AI Interaction.
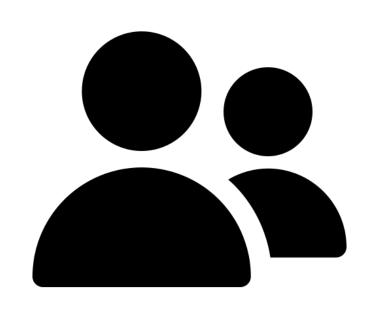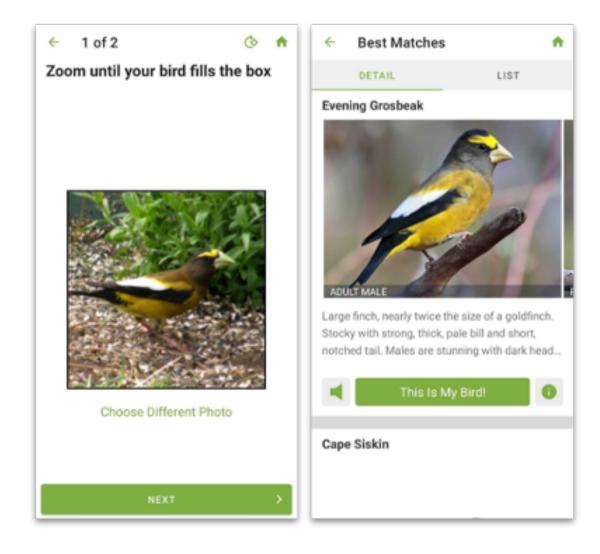
**Interpretability-accuracy tradeoff**

Q: What is the minimum accuracy of a baseline model that would convince you to use it over a model with explanations?



[Sunnie S. Y. Kim et al., ECCV 2022. HIVE.]

# Follow up: "Help Me Help the AI" — interview study with Merlin users



What **kind of explanation** best explains this prediction?

Interview

Merlin app

Prototypes

0.9 similar
0.7 similar
0.6 similar
0.6 similar
0.9 similar
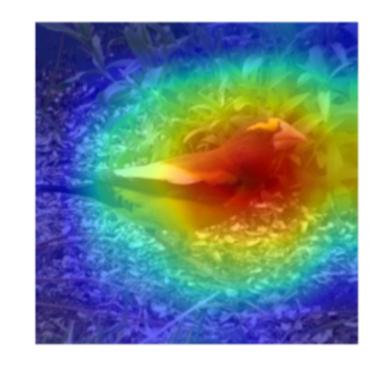
**Score for Evening Grosbeak = 1.7**

= - 1.2 long beak
+ 1.1 yellow beak
+ 0.8 black feathers
- 0.7 white body
+ 0.5 yellow body
+ 0.1 round body
…

Concepts

Heatmaps

Examples

Sunnie S. Y. Kim, Elizabeth Anne Watkins, Olga Russakovsky, Ruth Fong, Andrés Monroy-Hernández, arXiv 2022.
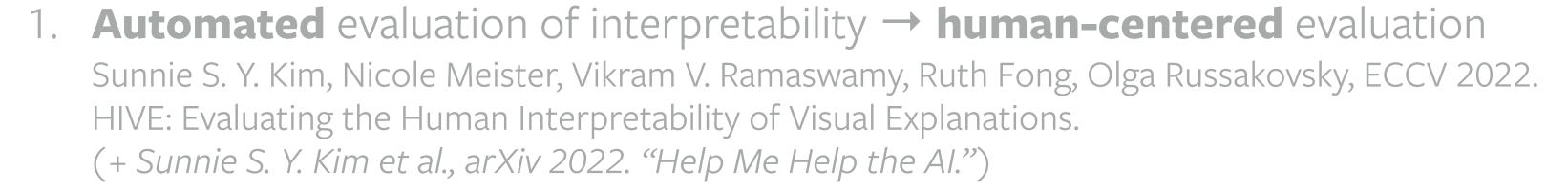"Help Me Help the AI": Understanding How Explainability Can Support Human-AI Interaction.

# Challenges for human evaluation

- Skill cost: web development skills
- Financial cost: budget for AMT experiments
- Time cost: human study design and iteration (e.g. task feasibility, IRB approval, quality control)

**Takeaway:** As a research community, invest in and reward human evaluation studies (like dataset development).
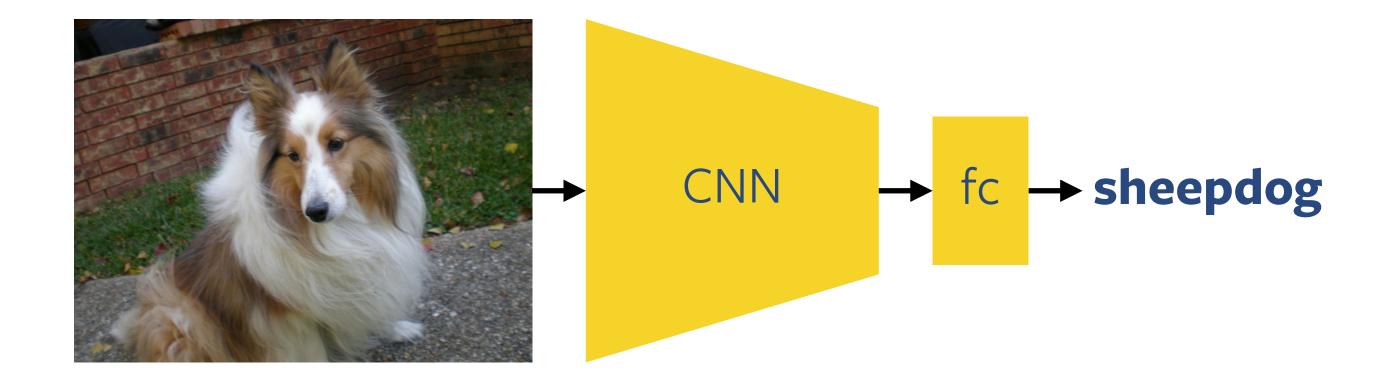
# Roadmap



Vikram V.
Ramaswamy

1. **Automated** evaluation of interpretability → **human-centered** evaluation
   Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, Olga Russakovsky, ECCV 2022.
   HIVE: Evaluating the Human Interpretability of Visual Explanations.
   (+ *Sunnie S. Y. Kim et al., arXiv 2022. "Help Me Help the AI."*)

2. Explanations via **labelled attributes** → explanations via **labelled attributes and unlabelled features**
   Vikram V. Ramaswamy, Sunnie S. Y. Kim, Nicole Meister, Ruth Fong, Olga Russakovsky, arXiv 2022.
   ELUDE: Generating Interpretable Explanations via a Decomposition into Labelled and Unlabelled Features.
   (+ *Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors in Concept-based Explanations.*)

3. Interpretability of **supervised** models → interpretability of **self-supervised** models
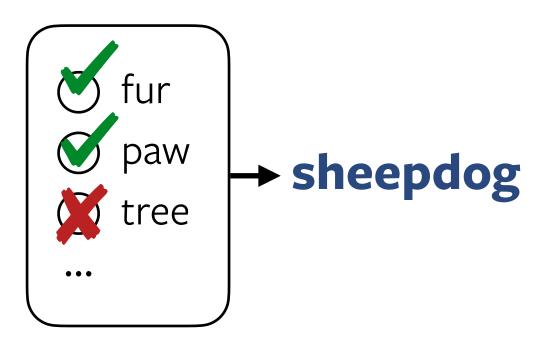   Iro Laina, Ruth Fong, Andrea Vedaldi, NeurIPS 2020.
   Quantifying Learnability and Describability of Visual Concepts Emerging in Representation Learning.

4. **Interpretability** in ML + CV → **interdisciplinary** research (interpretability + X)
   (+ *Nicole Meister\* and Dora Zhao\* et al., arXiv 2022. Gender Artifacts in Visual Datasets.*)
   (+ *Indu Panigrahi et al., arXiv 2022. Improving Fine-Grain Segmentation via Interpretable Modifications.*)

5. **Static** visualizations → **interactive** visualizations
   Ruth Fong, Alexander Mordvintsev, Andrea Vedaldi, Chris Olah, VISxAI 2021.
   Interactive Similarity Overlays.
   (+ *Devon Ulrich and Ruth Fong, in prep. Interactive Visual Feature Search.*)

# Concept-based explanations

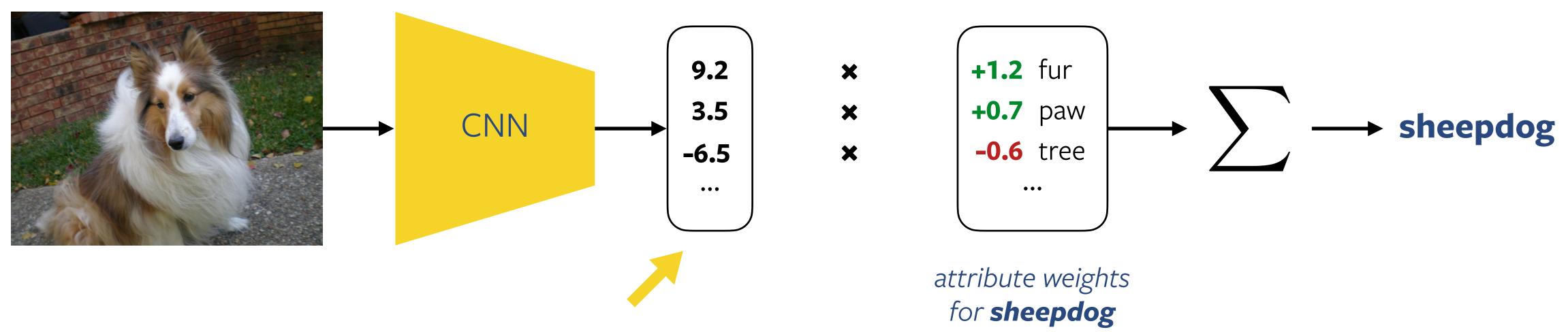Why did the model predict **sheepdog**?



CNN → fc → **sheepdog**

Concept-based explanation

✓ fur
✓ paw
✗ tree
...

→ **sheepdog**

**Pro:** Labelled concepts are interpretable to humans

# Concept Bottleneck: Linear Combination of Labelled Attributes

Predict present or absence of attribute

Linearly combine with attribute weights



CNN

9.2
3.5
-6.5
...

✖
✖
✖

**+1.2** fur
**+0.7** paw
**-0.6** tree
...

$\sum$ → **sheepdog**

*attribute weights for* ***sheepdog***

**Con:** Problems with predicting fractional values
- hard to interpret
- can encode hidden information

[Koh*, Nguyen*, Tang* et al., ICML 2020] 27

# Concept Bottleneck: Linear Combination of Labelled Attributes

Predict present or absence of attribute

Linearly combine with attribute weights



CNN

1
1
0
...

× × ×

+1.2 fur
+0.7 paw
-0.6 tree
...

$\sum$ → **sheepdog**

*attribute weights
for **sheepdog***

**Con:** Problems with predicting fractional values
- hard to interpret
- can encode hidden information

[Koh*, Nguyen*, Tang* et al., ICML 2020] 28

# ELUDE: **E**xplanation via a **L**abelled and **U**nlabelled **DE**composition of features



**Goal:** Approximate behavior of original CNN

[Vikram V. Ramaswamy et al., arXiv 2022. ELUDE.]

# ELUDE: Decomposition of labelled and unlabelled features



**Goal:** Approximate behavior of original CNN

1. Linearly combine **ground-truth, labelled attributes**

2. Learn remaining **unlabelled features as low-rank space**

feature activations

feature weights for **sheepdog**

ground-truth presence/absence of attributes

attribute weights for **sheepdog**

$8.2$
$4.5$
$-7.6$
...

**+1.1** $f_1$
**−0.3** $f_2$
**−0.7** $f_3$
...

**sheepdog**

$1$
$1$
$0$
...

**+1.2** fur
**+0.7** paw
**−0.6** tree
...

[Vikram V. Ramaswamy et al., arXiv 2022. ELUDE.]

30

**Attributes only:** % of model explained via labelled attributes decreases as task complexity increases

| Task | % Explained |
|---|---|
| 2-way scene classification (indoor vs. outdoor) | 95.7 |
| 16-way scene classification (home/hotel, workplace, etc.) | 46.2 |
| 365-way scene classification (airfield, bowling alley, etc.) | 28.8 |

Without fractional values encoding hidden information, attribute-only approaches are limited.

[Vikram V. Ramaswamy et al., arXiv 2022. ELUDE.]

# **Attributes only:** % of model explained via labelled attributes decreases as task complexity increases

| Scene group | TPR |
|---|---|
| home/hotel | 99.0 |
| comm-buildings/towns | 93.5 |
| water/ice/snow | 60.6 |
| forest/field/jungle | 40.2 |
| workplace | 14.2 |
| shopping-dining | 12.4 |
| cultural/historical | 6.5 |
| cabins/gardens/farms | 4.7 |
| outdoor-transport | 3.2 |
| indoor-transport | 0.0 |
| indoor-sports/leisure | 0.0 |
| indoor-cultural | 0.0 |
| mountains/desert/sky | 0.0 |
| outdoor-manmade | 0.0 |
| outdoor-fields/parks | 0.0 |
| industrial-construction | 0.0 |

Without fractional values encoding hidden information, attribute-only approaches are limited.

[Vikram V. Ramaswamy et al., arXiv 2022. ELUDE.]

# **Features + attributes:** Unlabelled features correspond to human-interpretable concepts



bowling alleys?

people eating?

outdoor sports fields?

castle-like buildings?

| Scene group | TPR |
| --- | --- |
| home/hotel | 99.0 |
| comm-buildings/towns | 93.5 |
| water/ice/snow | 60.6 |
| forest/field/jungle | 40.2 |
| workplace | 14.2 |
| shopping-dining | 12.4 |
| cultural/historical | 6.5 |
| cabins/gardens/farms | 4.7 |
| outdoor-transport | 3.2 |
| indoor-transport | 0.0 |
| indoor-sports/leisure | 0.0 |
| indoor-cultural | 0.0 |
| mountains/desert/sky | 0.0 |
| outdoor-manmade | 0.0 |
| outdoor-fields/parks | 0.0 |
| industrial-construction | 0.0 |

attributes only

[Vikram V. Ramaswamy et al., arXiv 2022. ELUDE.] 33
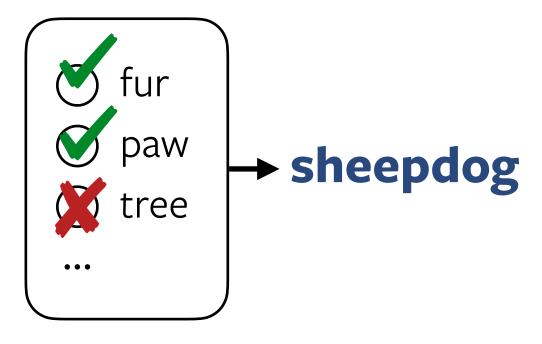
# Follow up: Overlooked factors in concept-based explanations

**Factor #1:** Probe dataset choice matters (i.e. different datasets → different explanations).

**Factor #2:** Some concepts used in explanations are harder to learn than output classes.

**Factor #3:** Humans can reason with a small amount of concepts (i.e. max 32 concepts).



✅ fur
✅ paw
❌ tree
...

→ **sheepdog**

Vikram V. Ramaswamy, Sunnie S. Y. Kim, Ruth Fong, Olga Russakovsky, arXiv 2022.
Overlooked Factors in Concept-based Explanations: Dataset Choice, Concept Salience, and Human Capability.

# Follow up: Overlooked factors in concept-based explanations

**Factor #1:** Probe dataset choice matters (i.e. different datasets → different explanations).

**Factor #2:** Some concepts used in explanations are harder to learn than output classes.

**Factor #3:** Humans can reason with a small amount of concepts (i.e. max 32 concepts).

**Suggestion:** Choose a probe dataset with a similar distribution to that of the training dataset.

**Training dataset: Places365**



hockey arena

**Probe dataset:**

**ADE20k**
{grandstand, goal, ice rink, scoreboard}

**Pascal**
{plaything, road}

Concepts used to explain **hockey arena** differ based on probe dataset.

Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors.

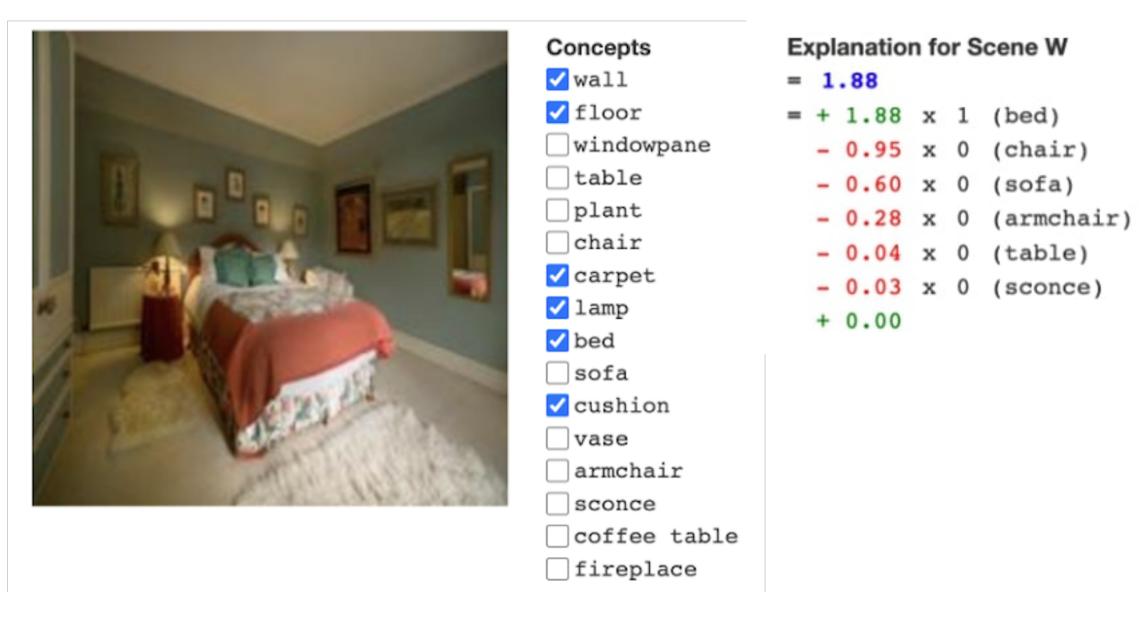# Follow up: Overlooked factors in concept-based explanations

**Factor #1:** Probe dataset choice matters (i.e. different datasets → different explanations).

**Factor #2:** Some concepts used in explanations are harder to learn than output classes.

**Factor #3:** Humans can reason with a small amount of concepts (i.e. max 32 concepts).

**Suggestion:** Only use easily learnable concepts in concept-based explanations.

**Training dataset:
Places365**



bathroom
(norm AP = 43.3)

**Probe dataset:
Broden**

| Concept | norm AP |
|---|---|
| toilet | 39.9 |
| shower | 18.8 |
| countertop | 12.6 |
| bathtub | 11.1 |
| screen door | 9.6 |

The class **bathroom** is easier to learn than the concepts used to explain it.

Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors. 36

# Follow up: Overlooked factors in concept-based explanations

**Factor #1:** Probe dataset choice matters (i.e. different datasets → different explanations).

**Factor #2:** Some concepts used in explanations are harder to learn than output classes.

**Factor #3:** Participants can reason with a small amount of concepts (i.e. max 32 concepts).

1. Which scene do you think the model predicts?
2. How many concepts would you prefer?
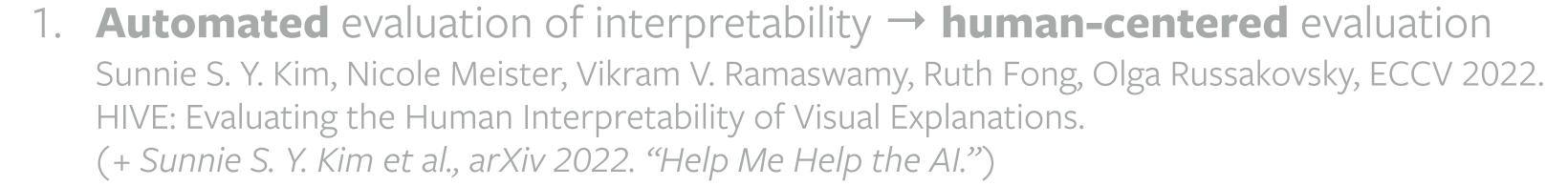
AMT human study
(*N = 125 participants*)



Participants struggle to identify concepts as the number of concepts increases.
(71.7% for 8 concepts; 56.8% for 32 concepts)

Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors.  37

# Challenges for concept-based methods

- Attributes-only approaches are incomplete
- Develop more methods to explain the "remainder"
  - Interpretable Basis Decomposition (IBD) [Zhou et al., ECCV 2018]
  - Automatic Concept-based Explanations (ACE) [Ghorbani et al., NeurIPS 2019]
  - ConceptSHAP [Yeh et al., NeurIPS 2020]
- Ensure that concept-based explanations are truly human-interpretable

**Takeaway:** Be realistic about the benefits and limitations of an interpretability method and work towards addressing the limitations.

# Roadmap


Iro Laina

1. **Automated** evaluation of interpretability → **human-centered** evaluation
   Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, Olga Russakovsky, ECCV 2022.
   HIVE: Evaluating the Human Interpretability of Visual Explanations.
   (+ *Sunnie S. Y. Kim et al., arXiv 2022. "Help Me Help the AI."*)

2. Explanations via **labelled attributes** → explanations via **labelled attributes and unlabelled features**
   Vikram V. Ramaswamy, Sunnie S. Y. Kim, Nicole Meister, Ruth Fong, Olga Russakovsky, arXiv 2022.
   ELUDE: Generating Interpretable Explanations via a Decomposition into Labelled and Unlabelled Features.
   (+ *Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors in Concept-based Explanations.*)

3. Interpretability of **supervised** models → interpretability of **self-supervised** models
   Iro Laina, Ruth Fong, Andrea Vedaldi, NeurIPS 2020.
   Quantifying Learnability and Describability of Visual Concepts Emerging in Representation Learning.

4. **Interpretability** in ML + CV → **interdisciplinary** research (interpretability + X)
   (+ *Nicole Meister* and Dora Zhao* et al., arXiv 2022. Gender Artifacts in Visual Datasets.*)
   (+ *Indu Panigrahi et al., arXiv 2022. Improving Fine-Grain Segmentation via Interpretable Modifications.*)

5. **Static** visualizations → **interactive** visualizations
   Ruth Fong, Alexander Mordvintsev, Andrea Vedaldi, Chris Olah, VISxAI 2021.
   Interactive Similarity Overlays.
   (+ *Devon Ulrich and Ruth Fong, in prep. Interactive Visual Feature Search.*)

# Supervised Learning



$$\left( \quad x \quad , \quad \text{sheepdog} \quad \right)$$

$$\quad x \qquad\qquad y$$

# Self-Supervised Learning



**x**

# Visual Concept

# Visual Concept



## Query

# Self-Supervised Learning



Top-1 Accuracy (%) on ImageNet

2016-2018

2019

2020

Colorization (Zhang et al.)
DeepCluster (Caron et al.)
Jigsaw (Kolesnikov et al.)
CPC (van den Oord et al.)
BigBiGAN (Donahue et al.)
MoCo (He et al.)
SeLa (Asano et al.)
PIRL (Misra et al.)
CMC (Tian et al.)
SeLa-**v2** (Asano et al.)
SimCLR (Chen et al.)
DeepCluster-**v2** (Caron et al.)
MoCo-**v2** (He et al.)
SimCLR-**v2** (Chen et al.)
BYOL (Grill et al.)
SwAV (Caron et al.)

# Self-Supervised Learning

Unlabelled data



Learn clusters

(e.g. DeepCluster, SeLa, SwaV)

Learn features

k-means

(e.g. SimCLR, MoCo, …)

cluster 1

cluster 2

cluster K

# Learnability



(A)

(B)

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]

# Learnability



(A)

white animal in snow

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]

# Describability



" dessert with chocolate sauce

(A)

(B)

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]

# Describability

> " dessert with chocolate sauce

Manual

(A)

(B)

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]

# Describability

> " dessert with chocolate sauce

**Manual**    **OR**    **Automatic**

(A)

(B)

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]

# Evaluation

## Learnability



**ImageNet cluster purity:**
how correlated is a cluster's contents
to a single ImageNet label?

*purity = 1* → *cluster only contains images*
*from a single ImageNet label*

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]
[Asano et al., ICLR 2020; He et al., CVPR 2020]

# Evaluation

Learnability

Describability



[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]
[Asano et al., ICLR 2020; He et al., CVPR 2020]

# Findings

**Follow up: Laina et al., ICLR 2022.**
Measuring the Interpretability of Unsupervised
Representations via Quantized Reverse Probing.

ImageNet cluster purity

**SeLa: cluster 393 (0.668)**
a newborn baby lying on a bed

**SeLa: cluster 332 (0.542)**
a snake on a hand

**MoCo: cluster 2335 (0.459)**
view of the mountains from the lake



98.3%   100.0%

93.3%   95.0%

[Iro Laina, et al., NeurIPS 2020. Quantifying Learnability and Describability.]
[Asano et al., ICLR 2020; He et al., CVPR 2020]

53

# ML fairness cross-talk: Gender artifacts in CV


Nicole Meister


Dora Zhao



**Average pose**

F    M

Area    Distance    Aspect

**Average color**

F    M

Color Channel Value

Red    Green    Blue

*Differences in top 20 female vs. male\* predicted images.*

1. **Resolution & Color**

2. **Person & Background**

Full   Full NoBg   MaskSegm   MaskRect   MaskSegm NoBg   MaskRect NoBg

3. **Contextual Objects**

Horse    Oven    Skateboard    Skateboard

Gender artifacts are **everywhere** in visual datasets.

*(\* binary perceived gender expression; we do not condone gender prediction.)*

Nicole Meister\*, Dora Zhao\*, Angelina Wang, Vikram V. Ramaswamy, Ruth Fong, Olga Russakovsky, arXiv 2022.
Gender Artifacts in Visual Datasets.

# Extending Interpretability to Geosciences



Indu Panigrahi    Elizabeth Barnes



Understand and improve
a coral reef fossil segmentation model
(our work)

Identify important regions in the world that
reliably predict seasonal climate
(Elizabeth Barnes' group at Colorado State)

Indu Panigrahi et al., arXiv 2022. Improving Fine-Grain Segmentation via Interpretable Modifications.
Zachary M. Labe and Elizabeth A. Barnes, JAMES 2021. Detecting Climate Signals Using Explainable AI.

# Challenges for novel frontiers in deep learning

- Need to contextualize interpretability to the novel frontiers
- Lack of access to standardized implementations

**Takeaway:** Collaboration and buy-in from novel research areas is crucial for interpretability in those frontiers.

# Roadmap

1. **Automated** evaluation of interpretability → **human-centered** evaluation
   Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, Olga Russakovsky, ECCV 2022.
   HIVE: Evaluating the Human Interpretability of Visual Explanations.
   (+ *Sunnie S. Y. Kim et al., arXiv 2022. "Help Me Help the AI."*)

2. Explanations via **labelled attributes** → explanations via **labelled attributes and unlabelled features**
   Vikram V. Ramaswamy, Sunnie S. Y. Kim, Nicole Meister, Ruth Fong, Olga Russakovsky, arXiv 2022.
   ELUDE: Generating Interpretable Explanations via a Decomposition into Labelled and Unlabelled Features.
   (+ *Vikram V. Ramaswamy et al., arXiv 2022. Overlooked Factors in Concept-based Explanations.*)

3. Interpretability of **supervised** models → interpretability of **self-supervised** models
   Iro Laina, Ruth Fong, Andrea Vedaldi, NeurIPS 2020.
   Quantifying Learnability and Describability of Visual Concepts Emerging in Representation Learning.

4. **Interpretability** in ML + CV → **interdisciplinary** research (interpretability + X)
   (+ *Nicole Meister\* and Dora Zhao\* et al., arXiv 2022. Gender Artifacts in Visual Datasets.*)
   (+ *Indu Panigrahi et al., arXiv 2022. Improving Fine-Grain Segmentation via Interpretable Modifications.*)

5. **Static** visualizations → **interactive** visualizations
   Ruth Fong, Alexander Mordvintsev, Andrea Vedaldi, Chris Olah, VISxAI 2021.
   Interactive Similarity Overlays.
   (+ *Devon Ulrich and Ruth Fong, in prep. Interactive Visual Feature Search.*)

# Interpretability Tools



Orig Img      Mask      Grad CAM

Net Dissect     Activation Maximization     Feature Vis

Current tools render **static images**.         Future tools should be **interactive**!

[Fong et al., ICCV 2019; Selvaraju et al., ICCV 2017; Bau et al., CVPR 2017;
Mahendran & Vedaldi, IJCV 2016; Olah et al., Distill 2018; Fong et al., VISxAI 2021]

# Interpretability: Interactive, Exploratory, Easy-to-use



$c_1$ $c_2$ $c_3$ $c_4$ $c_5$ $f_6$ $f_7$ $f_8$ **sheepdog**

How can we **easily explore** hypotheses about the model?

# Interactive Similarity Overlays

# Spatial Activations



$f_a$ $f_b$ **golden retriever**

# Spatial Activations



$f_a$

$f_b$ → **golden retriever**

# Interactive Similarity Overlays



$$a_{6,5} = [17.7, 0, 103.4, 6.81, 0, 0, 0, 0, 32.0, 0, 0, 0, ...]$$

# Interactive Similarity Overlays



[Fong et al., VISxAI 2021. Interactive Similarity Overlays.] 64

# Demo: Interactive Similarity Overlays

bit.ly/interactive_overlay



Interactive visualizations empower practitioners to easily explore model behavior.

[Fong et al., VISxAI 2021. Interactive Similarity Overlays.] 65

# Interactive Similarity Overlays

An interactive tool for understanding what neural networks consider similar and different.



**Hover over different parts of the above images**. This interactive visualization shows how similar (or different) a neural network considers different image patches to the current image patch (highlighted in yellow). Try hovering over animal features (e.g., noses, eyes, faces) and background regions.

*This article is best viewed in Google Chrome.*

**Layers with different spatial resolutions.**



The location of the highlighted image patch (in yellow) has been synchronized across images, such that the overlays show similarity scores with respect to each image's highlighted patch (i.e., no similarity scores were computed between images). Consider exploring edges in mixed3b layers and semantic features (e.g., objects and object parts, like noses and eyes) in mixed4e and mixed5b layers.

Interactive Overlays: Basic Examples (TensorFlow) ☆

File   Edit   View   Insert   Runtime   Tools   Help   Cannot save changes

+ Code   + Text   ⬠ Copy to Drive

```python
# Get images
img_urls = ["https://raw.githubusercontent.com/ruthcfong/interactive_overlay/master/images/dog_cat.jpeg",
            "https://raw.githubusercontent.com/ruthcfong/interactive_overlay/master/images/flowers.jpeg",
            "https://raw.githubusercontent.com/ruthcfong/interactive_overlay/master/images/pig.jpeg",
            "https://raw.githubusercontent.com/ruthcfong/interactive_overlay/master/images/bowtie_guy.jpeg",
            "https://raw.githubusercontent.com/ruthcfong/interactive_overlay/master/images/beer.jpeg",
            "https://raw.githubusercontent.com/ruthcfong/interactive_overlay/master/images/chain.jpeg"]
imgs = [load(url) for url in img_urls]

model = models.InceptionV1()
model.load_graphdef()
```

```python
acts = get_acts(model, imgs[0], "mixed4d")
grid = np.hstack(np.hstack(cossim_grid(acts, acts)))
colored_grid = add_color_index(grid, acts.shape[0])
```

```python
lucid_svelte.CossimOverlay({
    "image_url": _image_url(imgs[0]),
    "masks_url": _image_url(colored_grid),
    "size": 224,
    "N": acts.shape[0],
})
```
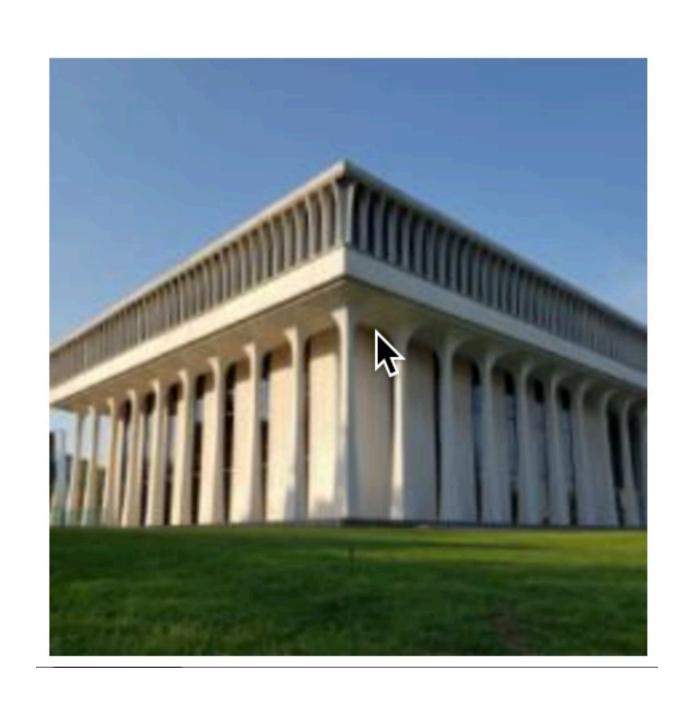


6,4

# Preview: Interactive Visual Feature Search

bit.ly/interactive_search

Devon Ulrich

# Challenges for interactive visualizations

- Skills cost: web development skills

  - 📈 HuggingFace Spaces, Gradio, Streamlit

- Potential misuse: Intuition-based insights should be validated via quantitative experiments

- Poor incentives: software tooling for research is often not rewarded

- Inadequate publishing structures: Sparse publishing venues for interactive articles and/or visualizations

  - 📉 Distill journal hiatus

  - 📈 CVPR demo track

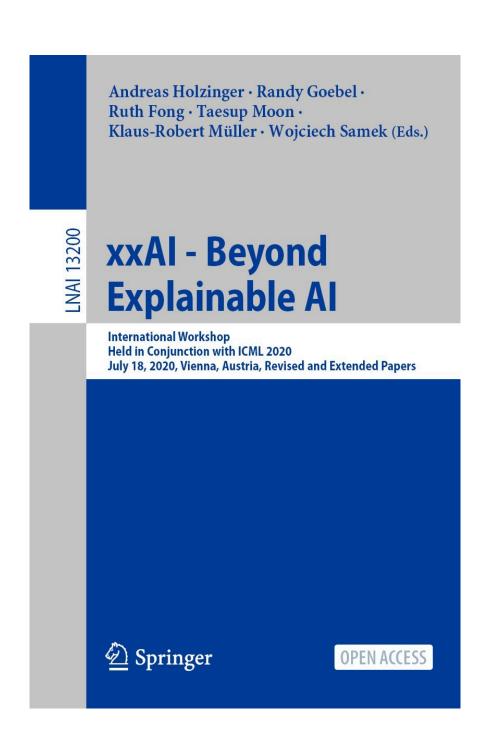- Lack of cross-talk: HCI and AI communities are developing interpretability tools fairly independently

**Takeaway:** Relevant research communities should collectively invest in and reward software tooling for research, particularly interactive tools.

# Takeaways from challenges in interpretability

- **Human studies:** As a research community, invest in and reward human evaluation studies (like dataset development).

- **(Concept-based) interpretability:** Be realistic about the benefits and limitations of an interpretability method and work towards addressing the limitations.

- **New frontiers:** Collaboration and buy-in from novel research areas is crucial for interpretability in those frontiers.

- **Interactive visualizations:** Relevant research communities should collectively invest in and reward software tooling for research, particularly interactive tools.
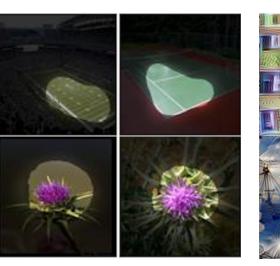
# Directions for the next decade of interpretability

1. Develop interpretability methods for **diverse domains**

   - Beyond CNN classifiers: self-supervised learning, generative models, etc.

2. Center **humans** throughout the development process

   - In design, co-develop methods with real-world stakeholders.

   - In evaluation, measure human interpretability and utility of methods.

   - In deployment, package interpretability tools for the wider community.

ICML 2020 workshop on XXAI

Andreas Holzinger · Randy Goebel ·
Ruth Fong · Taesup Moon ·
Klaus-Robert Müller · Wojciech Samek (Eds.)

LNAI 13200

**xxAI - Beyond
Explainable AI**

International Workshop
Held in Conjunction with ICML 2020
July 18, 2020, Vienna, Austria, Revised and Extended Papers

Springer          OPEN ACCESS
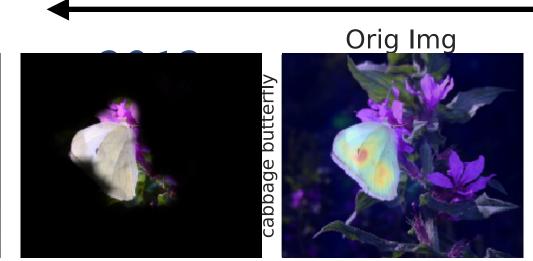
# An incomplete retrospective: the first decade of interpretability



Primarily focused on understanding and approximating **CNNs**

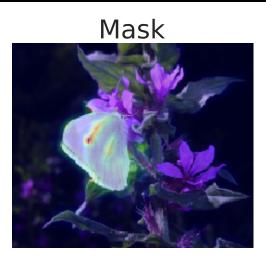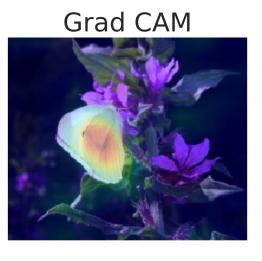**Feature visualization (2013–2018)**
Activation Max., Feature Inversion,
Net Dissect, Feature Vis.

2022

Orig Img      Mask      Grad CAM

concepts c
wing color
undertail color
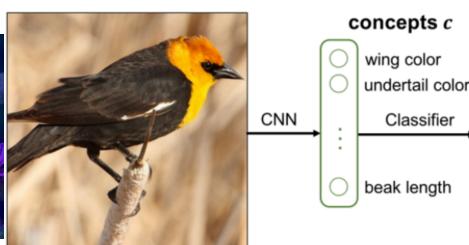CNN        Classifier        task y
bird species
beak length

**Attribution heatmaps (2013–2019)**
Gradient, Grad-CAM,
Occlusion, Perturbations, RISE

**Interpretable-by-design (2020–now)**
Concept Bottleneck, ProtoPNet,
ProtoTree

[Selvaraju et al., ICCV 2017; Fong* & Patrick* et al., ICCV 2019;
Bau* & Zhou* et al., CVPR 2017; Olah et al., Distill 2017; Koh*, Nguyen*, Tang* et al., ICML 2020]

# Into the future: the next decade of interpretability

???

2022                                                        2032

Devon Ulrich  Dora Zhao  Nicole Meister  Sunnie S. Y. Kim  Vikram V. Ramaswamy  Angelina Wang  Ryan A. Manzuk
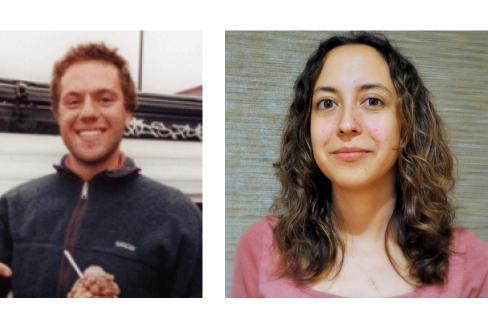
Iro Laina  Andrea Vedaldi  Elizabeth Anne Watkins  Andrés Monroy-Hernández  Chris Olah  Alex Mordvintsev  Adam C. Maloof  Olga Russakovsky

We're hiring postdocs!
bit.ly/vai-lg-postdoc

# Thank You