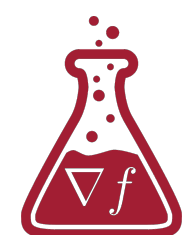# Decomposing and Editing Predictions by Modeling Model Computation

Harshay Shah, Andrew Ilyas, Aleksander Mądry

**ICML 2024**

https://arxiv.org/abs/2404.11534

gradientscience.org/modelcomponents

# Why study model predictions?

**Tinker your ML pipeline** 💻

**Try to get SOTA results** 🔢



**Repeat**

# Why study model predictions?



Core issue: We don't understand **how** models internally turn examples into predictions

# Models as computation graphs



Example $z$          Model output $f(z)$

Any metric that quantifies "correctness"
e.g., cross-entropy loss, correct-class confidence., etc

# Models as computation graphs



**A set of model components $C$**

**Example $z$**

**Model output $f(z)$**

Model $f$ as a computation graph over model components

# Models as computation graphs



Model $f$ as a computation graph over model components

**Examples of model components in common model architectures**

Convolution filters in ResNet models

Weight vectors in MLPs

Attention heads & MLPs in Transformers

Coefficients in linear models

# Models as computation graphs

Model components

## High-level question

Can we somehow understand how model components collectively turn examples into predictions?

Examples

Convolution filters in ResNet models

Weight vectors in MLPs

Attention heads & MLPs in Transformers

Parameters in linear models

# Background: interpreting model components

**Vision models**



Convolution filters learn to detect curves and frequency [Cammarata et al. 2020]



unit 149 "mountain top" (acc lost: train 1.2% val 3.5%)

unit 242 "house" (acc lost: train 1.5% val 2.5%)

Convolution filters in deeper layers detect high-level concepts [Bau et al. 2020]

# Background: Interpreting model components

## Vision models



Convolution filters learn to detect curves and frequency [Cammarata et al. 2020]



unit 149 "mountain top" (acc lost: train 1.2% val 3.5%)

unit 242 "house" (acc lost: train 1.5% val 2.5%)

Convolution filters in deeper layers detect high-level visual concepts [Bau et al. 2020]

## Language models



Induction heads in transformers [Olsson et al. 2022]



**Factual Knowledge**

Q27 Ireland — P36 capital → Q1761 Dublin

Knowledge neurons encode factual knowledge  [Dai et al. 2021]

"Duplicate token head", "Name-mover head", "Backup head", "ML Tea head"... 🤔
[Wang et al. 2022]

Vision models

Language models

**Our goal**

Analyze how *every* model component $c \in C$

contributes to individual predictions $f( \cdot )$

Knowledge neurons encode factual
knowledge [Dai et al. 2021]

"Duplicate token head", "Name-mover head",
"Backup head", "ML Tea head"...
[Wang et al. 2022]

Convolution filters in deeper layers detect
high-level visual concepts [Bau et al. 2020]

# Our work

**Component attribution framework**
Decompose any prediction into "contributions" from every model component

**COAR**: Component Attribution via Regression
A general method for efficient and accurate component attribution

**COAR-Edit**: Model editing using component attributions
Edit model behavior by ablating a targeted subset of components

# The component attribution framework

**Main idea**

> If we can "understand" how all model components shape a prediction
>
> ↓
>
> we should be able to estimate how predictions <u>change</u> in response
>
> to interventions to one or more model components

# The component attribution framework

**Main idea**

If we can "understand" how all model components shape a prediction,

⬇

we should be able to estimate how interventions to model components
<u>change</u> model predictions

**Component ablations as interventions**

A **component ablation** intervenes on the
*parameters* corresponding to one or
more model components.

For instance, zeroing out or adding noise



Ablate([1 **0** 1 **0** ... 1])

↑

Ablation vector **v**

"Cat" (35%)

"Cat" (32%)

# The component attribution framework

**Component ablations as interventions**

A **component ablation** intervenes on the *parameters* corresponding to one or more model components.

Ablate([1 **0** 1 **0** ... 1])
↑
Ablation vector **v**

"Cat" (35%)

"Cat" (32%)

**Component attribution**

A **component attribution** $g$ takes as input an ablation vector $v$ and estimates the effect of the component ablation on a given model prediction.

Ablate($\mathbf{v_1}$)

Ablate($\mathbf{v_2}$)

Ablate($\mathbf{v_k}$)

$\approx g(\mathbf{v}_1)$

$\approx g(\mathbf{v}_2)$

$\approx g(\mathbf{v}_k)$

# Formalizing component attribution

**Fix:**

| Example $z$ | Trained model $f$ | Set of components $C$ |
|---|---|---|
| e.g., from ImageNet | e.g., a ResNet50 | e.g., conv filters in all layers |
| | | `layer7.block3.conv[42]` |

For **any** component ablation $v \in \{0,1\}^{|C|}$

   1. Using $v$, apply component ablation to the model $f$

   2. Evaluate output of ablated model on example $z$ to get $f(z, v)$

**Goal**: Given (any) component ablation $v$, estimate $f(z, v)$ (i.e., without intervening)

# Formalizing component attribution



Component ablation
$v = [\mathbf{0} \; 1 \ldots \mathbf{0} \; 1]$

Ground-truth output of ablated model

Component attribution directly predicts model output $f(\text{🐱})$

# Formalizing component attribution



$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$ $w_7$ $w_8$ $w_9$ $w_{10}$ $w_{11}$ $w_{12}$

Attribution scores $w$ estimate component-wise contributions

$$f(🐱) \approx g_{🐱}(\textcolor{red}{v}) = w^\top \textcolor{red}{v} + b$$

Component attribution directly predicts model output $f(🐱)$

Component ablation
$\textcolor{red}{v} = [\textcolor{red}{0}\ 1 ... \textcolor{red}{0}\ 1]$

Ground-truth output of ablated model

# Formalizing component attribution



**Next:** We want to estimate component attributions that accurately predict how component ablations change model predictions

Ablate($\mathbf{v_i}$)

Each point is a random ablation $v_i$

$$g(\mathbf{v_i}) = b + \mathbf{w}^\top \mathbf{v_i}$$

Component attribution

Attribution vector $w$ quantifies component-wise contributions

# Our work

# **COAR**: Component Attribution via Regression

Cast component attribution into a **supervised learning** problem in two steps

**Step 1/2**

**Construct a dataset of component ablations** by ablating <u>random</u> subsets of components and recording both the ablations and the ablated model's outputs for each example of interest.



Example $z$

Ablate($\mathbf{v_1}$)     Ablate($\mathbf{v_2}$)     Ablate($\mathbf{v_k}$)

$f(z, v_1)$     $f(z, v_2)$     $f(z, v_k)$

Dataset of component ablations $D^{(z)} = \{(v_i, f(z, v_i))\}_{i=1}^{n}$

# **COAR**: <u>Co</u>mponent <u>A</u>ttribution via <u>R</u>egression

Cast component attribution into a **supervised learning** problem in two steps

## Step 2/2

**Fit a linear regression model** that maps an ablation vector $v_i$ to the ablated models' output $f(z, v_i)$. The weights $(w, b)$ of this linear model serve as our component attribution $g^{(z)}(v) = w^\top v + b$

Ground-truth output of ablated model

$$(w^{(z)}, b^{(z)}) = \arg \min_{w,b} \sum_{D^{(z)}} \left( f(z, v_i) - v_i^\top w - b \right)^2$$

Dataset of component ablations

Attribution-based estimate

# **COAR**: Component Attribution via Regression

Does COAR learn **accurate** component attributions?

## Setup

| Example $z$ | Model $f$ | Components $C$ |
|---|---|---|
| from ImageNet | ImageNet-trained ResNet50 | 22,720 conv filters |



layer∗.block∗.conv∗

Evaluating component attributions $g^{(z)}$

1. Sample an (unseen) <u>random</u> ablation vector $v$

2. Check if the attribution-based estimate $g^{(z)}(z)$ predicts ground-truth output $f(z, v)$

# **COAR**: Component Attribution via Regression

Does COAR learn **accurate** component attributions?



ResNet-50 trained on ImageNet

Corr. $\rho(z) = 0.70$

Line $x = y$

COAR estimate $g^{(z)}(0_C)$

True counterfactuals $f_M(z, 0_C)$ on example $z$

Example z

# **COAR**: Component Attribution via Regression



ResNet-50 trained on ImageNet

ResNet-50 trained on ImageNet

**Results consistent across**

Architectures: MLPs, CNNs, Transformers

Language models: GPT-2, Phi-2, Llama-7b

Datasets: ImageNet, CIFAR-10, TinyStories, BoolQ

Internal Infl.
[II]

Neuron Cond.
[NC]

Grad ⊙ Param
[GP]

Leave-one-out
[LOO]

**COAR**
[Ours]

# Our work

**Component attribution framework**
Decompose any prediction into "contributions" from every model component

$\downarrow$

**COAR**: <u>Co</u>mponent <u>A</u>ttribution via <u>R</u>egression
A general method for efficient and accurate component attribution

$\downarrow$

**COAR-Edit**: Model editing using component attributions
Edit model behavior by ablating a targeted subset of components

# **COAR-Edit**: Model editing using COAR attributions

**Component attribution asks**

> *How would model outputs change if we were to ablate a subset of components?*

↕

**Model editing inverts this to**

> *Which components, when ablated, would change model outputs in a specific way?*

# **COAR-Edit**: Model editing using COAR attributions

**Goal**: perform a model edit that improves performance on **target examples** without degrading performance on **reference examples**



$$\texttt{COAR-edit}(z_1, \ldots, z_n) = \textbf{v*}$$

Targeted ablation

A few target and reference examples

$z$

$f(z)$

$\hat{f}(z)$

**Before edit:**

Target accuracy: **60%**
Reference accuracy: **80%**

Effect of model edit

**After edit:**

Target accuracy: **85%**
Reference accuracy: **80%**

# **COAR-Edit**: Model editing using COAR attributions

**Main idea**

Use COAR attributions to identify model components that,
when ablated,  change model behavior in a targeted manner

**No additional training needed** ✅          **Sample-efficient** ✅

# **COAR-Edit**: Model editing using COAR attributions

**Step 1/3**

> Compute COAR attributions for target and reference examples

**Step 2/3**

> For every component, quantify its "importance" to target examples
> *relative* to reference examples with a simple t-test (null: target ~ reference)

**Step 3/3**

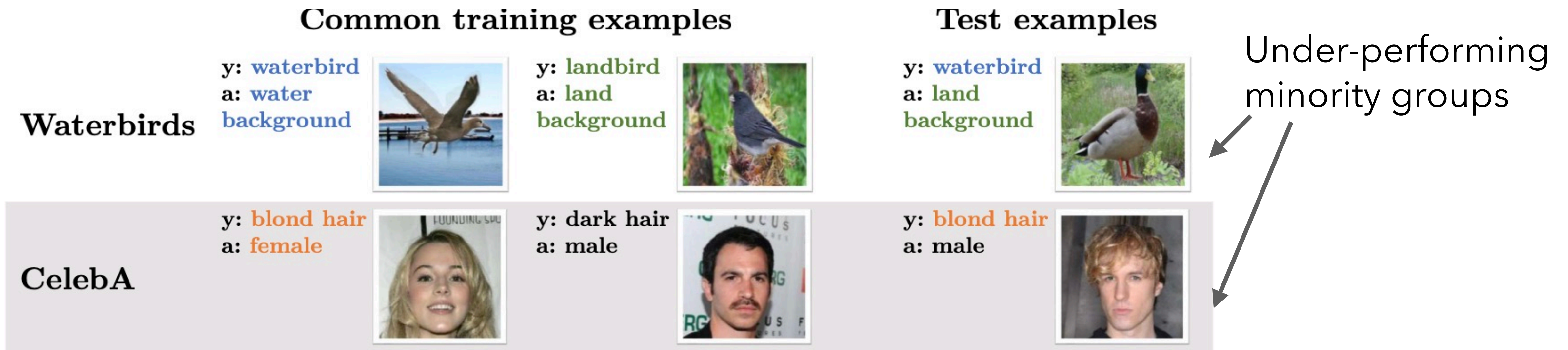> Ablate the bottom-$k$ components with the lowest test statistics to improve
> model performance on the target examples.

# **Case study #1**: Improving group robustness

**Problem**

1. Models latch on to spurious correlations in the training dataset

2. At test time, models performance sucks when spurious correlation is absent



Common training examples     Test examples

Under-performing minority groups

Waterbirds
- y: waterbird / a: water background
- y: landbird / a: land background
- y: waterbird / a: land background

CelebA
- y: blond hair / a: female
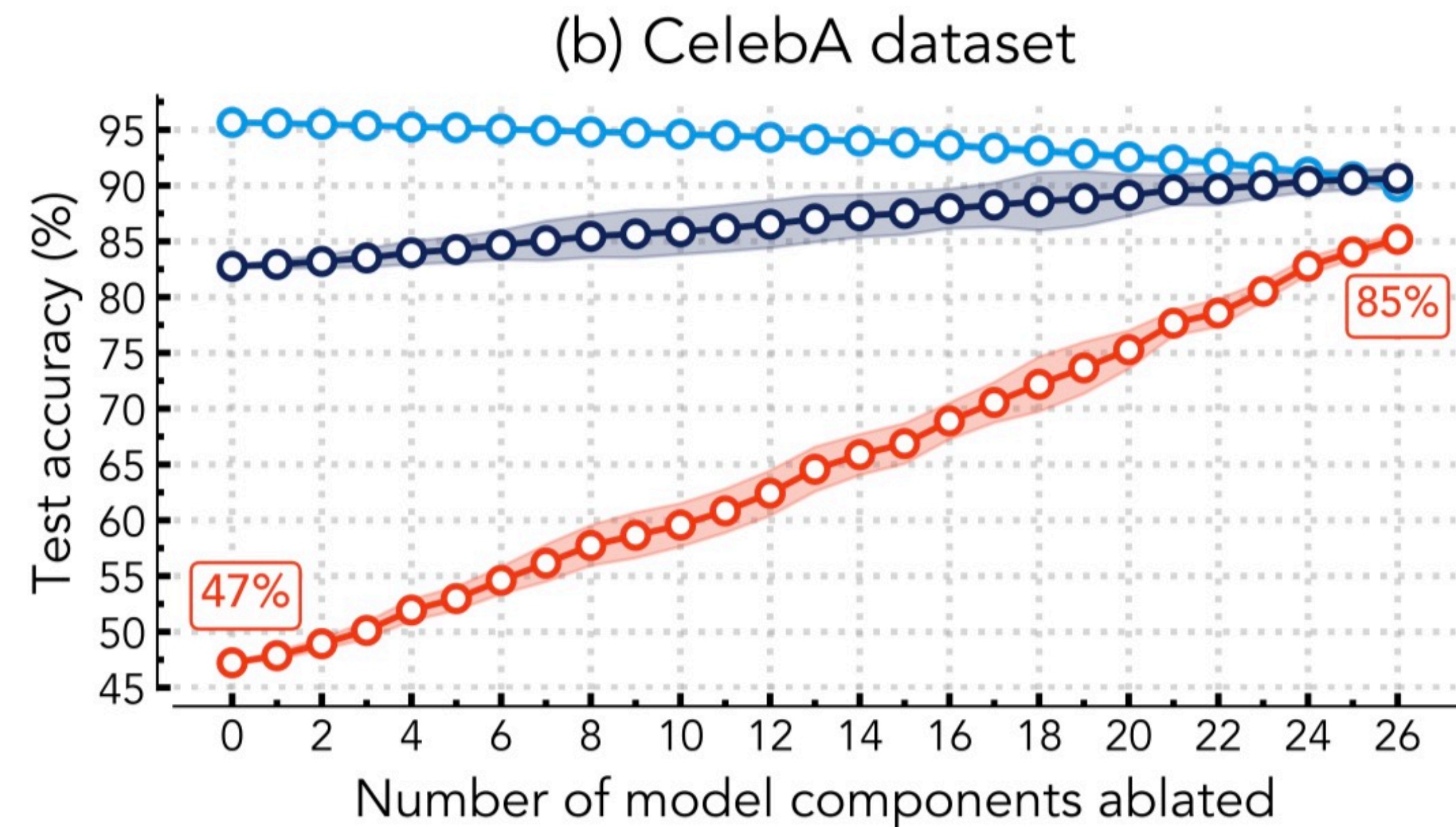- y: dark hair / a: male
- y: blond hair / a: male

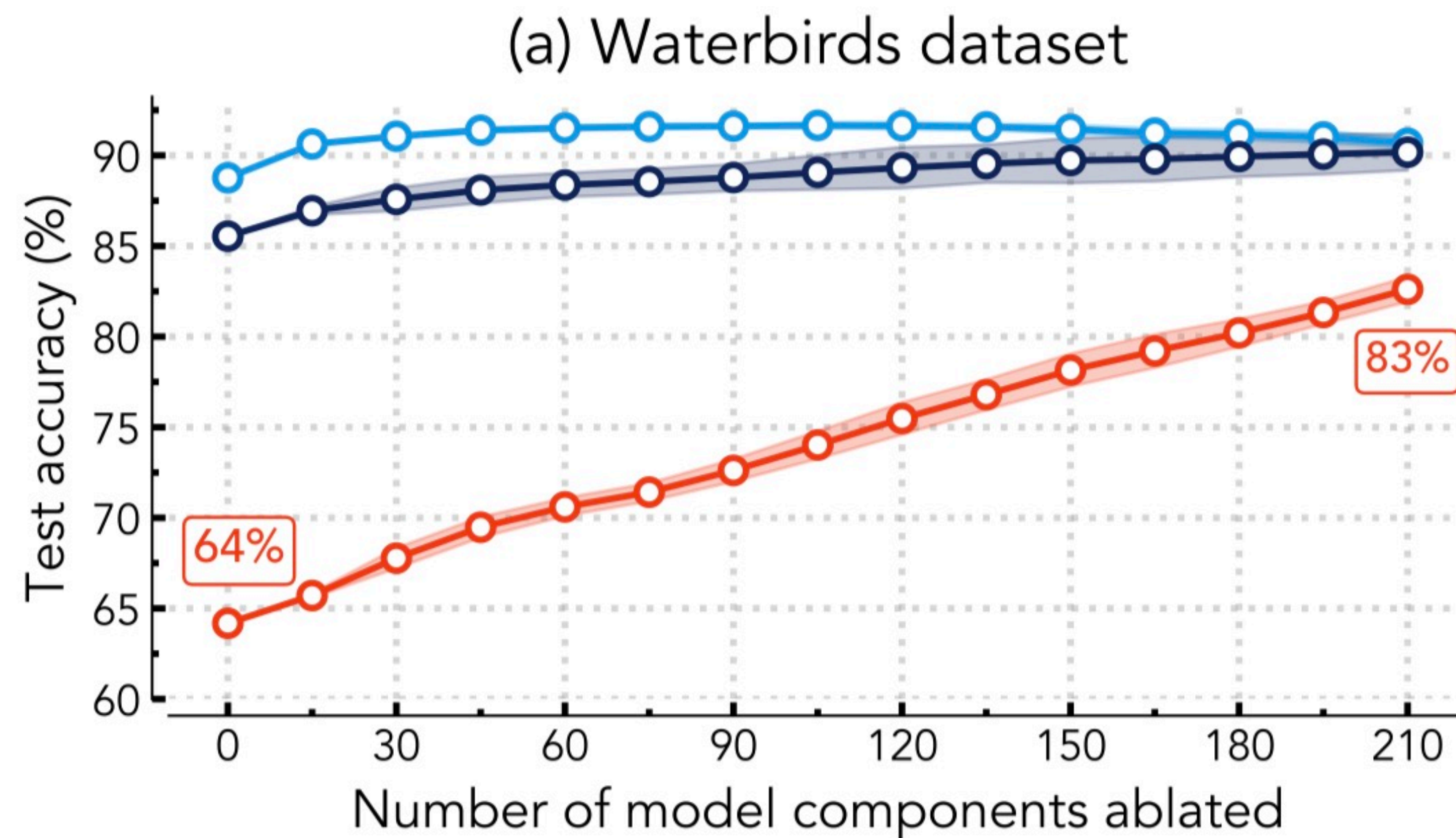[Group robustness benchmarks from Sagawa et al. 2020]

# Case study #1: Improving group robustness

**Applying COAR-Edit**

> **1.** **Target examples**: a few examples (~10) from the majority group(s)
>
> **2.** **Reference examples**: a few examples (~10) from the minority group(s)



(a) Waterbirds dataset

(b) CelebA dataset

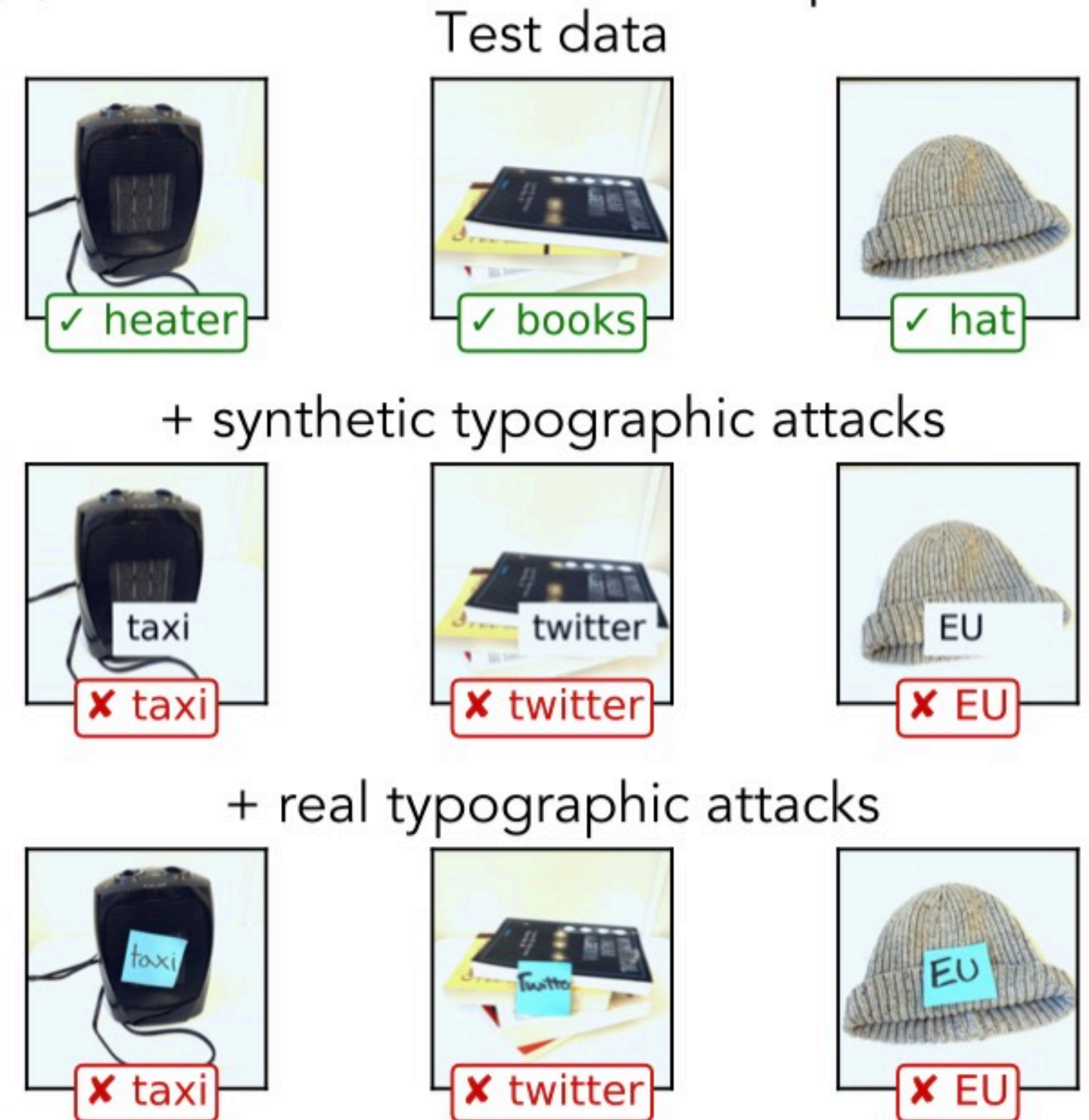Averaged over examples — Averaged over subpopulations — On worst-performing subpopulation.

# **Case study #2**: Robustness to typographic attacks

**Problem**

Zero-shot CLIP classifiers are
sensitive to typographic attacks

[Goh et al. 2021]

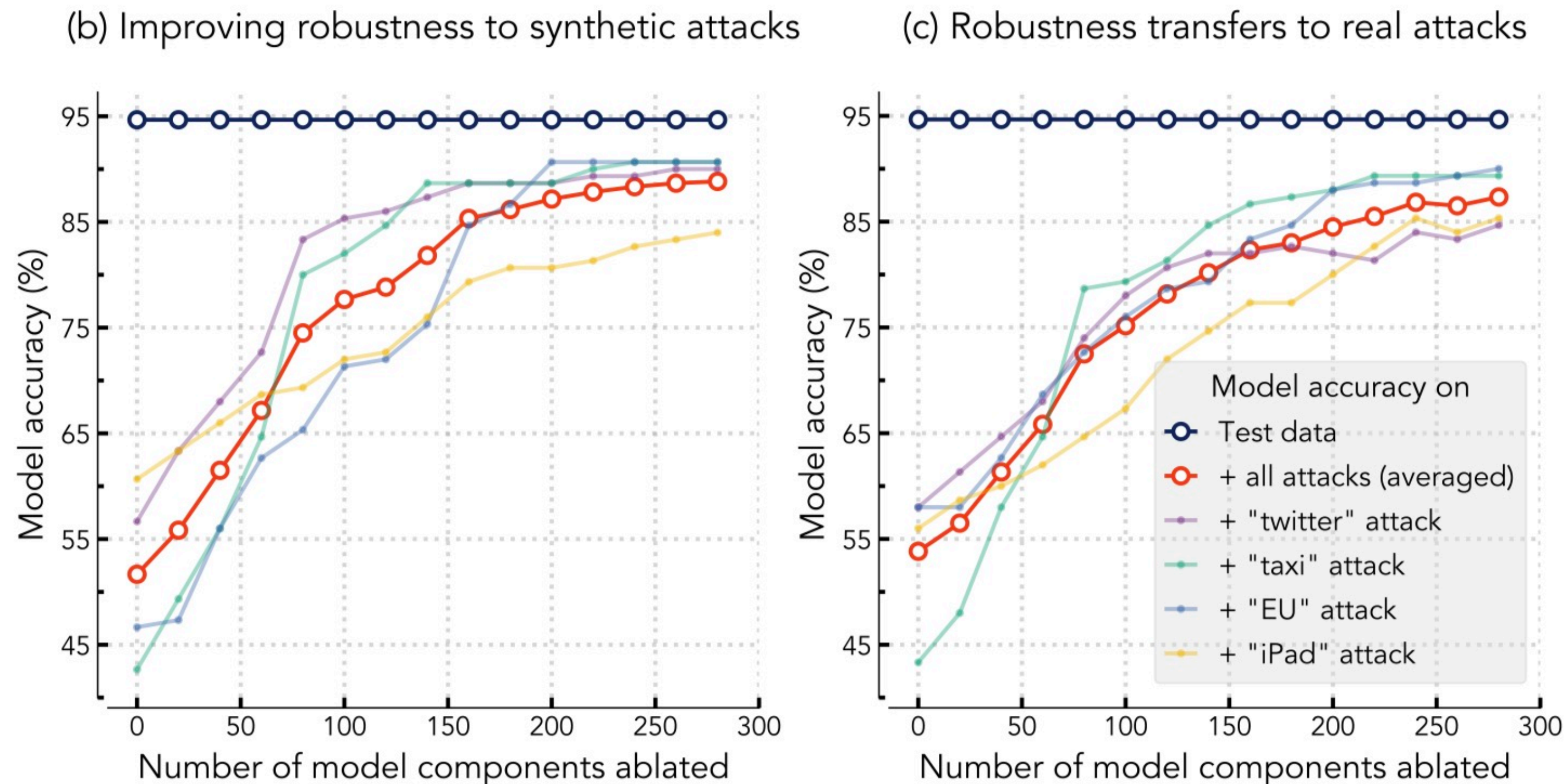Evaluating a CLIP ViT-B/16 model
on images w/ and w/o attacks



(a) Effect of attacks on model predictions

# Case study #2: Robustness to typographic attacks

**Applying COAR-Edit**

> **1.** Target examples: a few examples (~10) with synthetic typographic attacks
>
> **2.** Reference examples: a few examples (~10) without typographic attacks



(b) Improving robustness to synthetic attacks

(c) Robustness transfers to real attacks

Model accuracy on
- Test data
- + all attacks (averaged)
- + "twitter" attack
- + "taxi" attack
- + "EU" attack
- + "iPad" attack

Model accuracy (%)

Number of model components ablated

# Summary

→ Decompose predictions into contributions from every model component

   → How? Use **COAR** to learn component attributions

→ Edit model behavior at the level of examples, subpopulations, and concepts

   → How? Use **COAR-Edit** to identify and ablate a targeted set of model components

Check out our paper for more findings!
`https://arxiv.org/abs/2404.11534`

**@harshays_**