# Counterfactual Visual Explanations



**Yash Goyal**
(Georgia Tech)

Ziyan Wu
(Siemens)

Jan Ernst
(Siemens)

Dhruv Batra
(Georgia Tech)

Devi Parikh
(Georgia Tech)

Stefan Lee
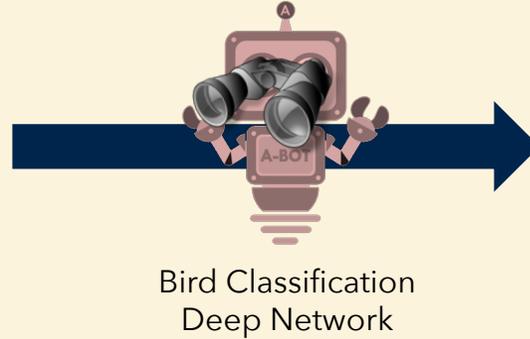(Georgia Tech)

# Counterfactual Visual Explanations



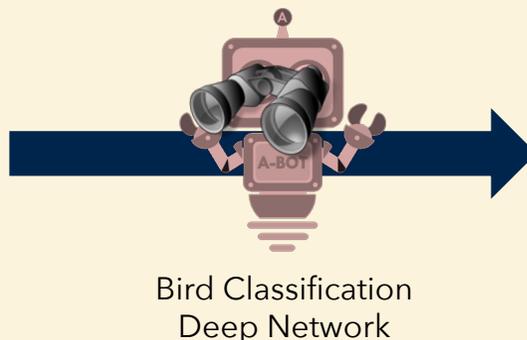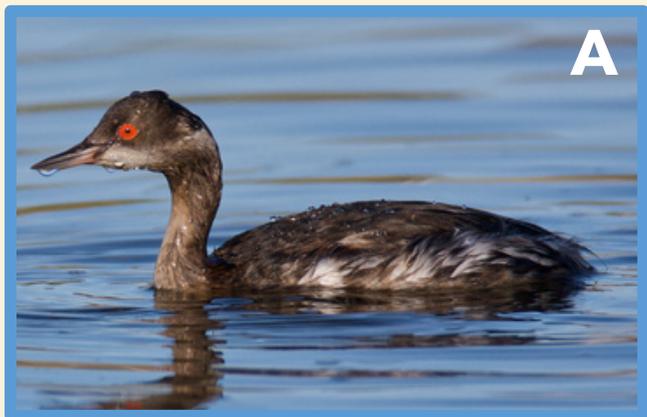**An XAI Question:** Why did the model predict **Eared Grebe** instead of **Horned Grebe**?
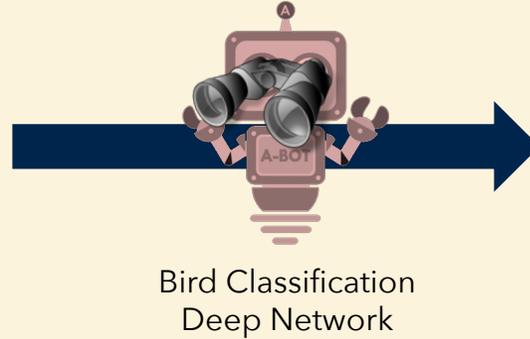
# Counterfactual Visual Explanations



**An XAI Question:** Why did the model predict **Eared Grebe** instead of **Horned Grebe**?

For input X, why did the model predict Y instead of Z?

# Counterfactual Visual Explanations



**An XAI Question:** Why did the model predict **Eared Grebe** instead of **Horned Grebe**?

**For input X, why did the model predict Y instead of Z?**

**Explanation through Counterfactual:**

If **X** was **X***, then the outcome would have been **Z** rather than **Y**.

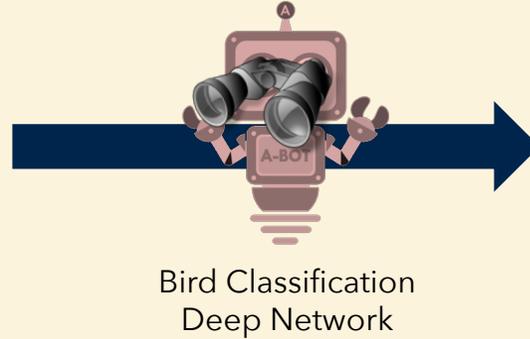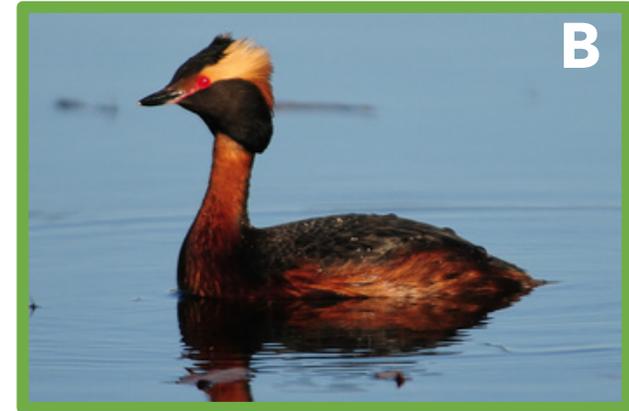# Counterfactual Visual Explanations



**An XAI Question:** Why did the model predict **Eared Grebe** instead of **Horned Grebe**?
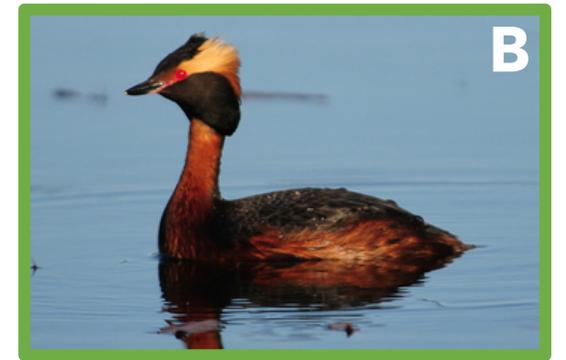
**Explanation through Counterfactual:**

What would have to change in **image A** to make the model predict **Horned Grebe**?

# Counterfactual Visual Explanations



An XAI Question: Why did the model predict **Eared Grebe** instead of **Horned Grebe**?

## Explanation through Counterfactual:

What would have to change in **image A** to make the model predict **Horned Grebe**?

An image where the network predicts Horned Grebe.
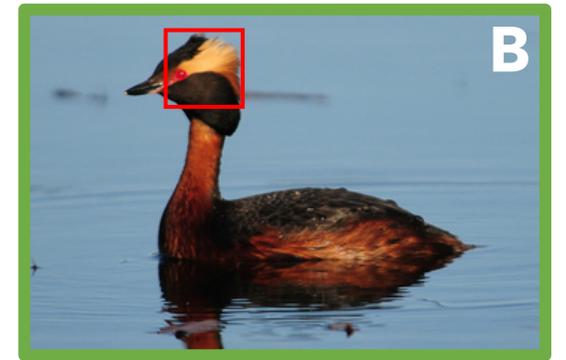
# Counterfactual Visual Explanations

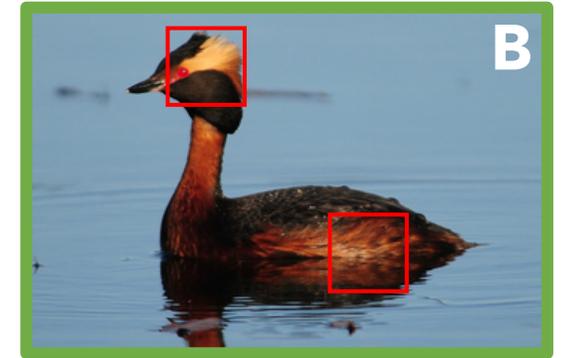What would have to change in **image A** to make the model predict **Horned Grebe**?





An image where the network predicts Horned Grebe.

# Counterfactual Visual Explanations

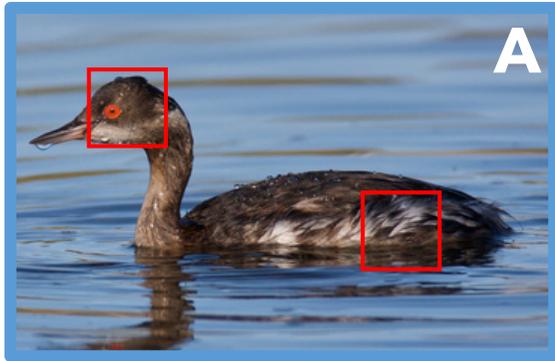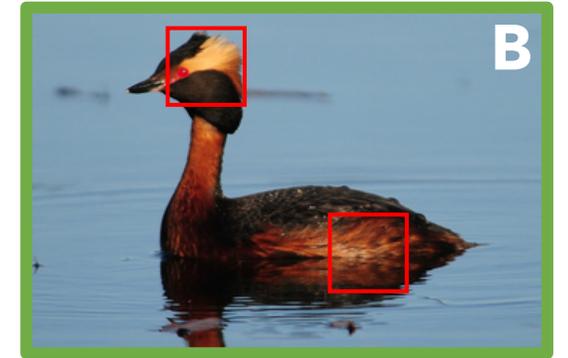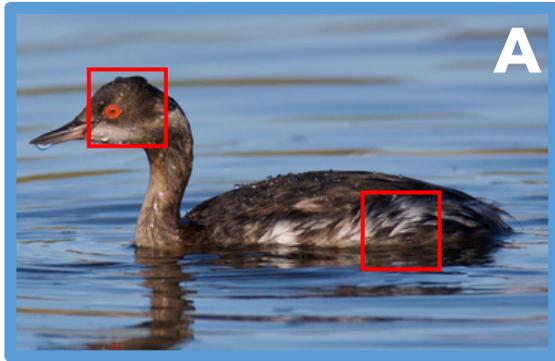What would have to change in **image A** to make the model predict **Horned Grebe**?



If  looked more like 



An image where the network predicts Horned Grebe.

# Counterfactual Visual Explanations

What would have to change in **image A** to make the model predict **Horned Grebe**?



If  looked more like 

If  looked more like 



An image where the network predicts Horned Grebe.

# Counterfactual Visual Explanations

What would have to change in **image A** to make the model predict **Horned Grebe**?



If  looked more like 

If  looked more like 



An image where the network predicts Horned Grebe.

## How can we identify these region pairs important to the model?

# Counterfactual Visual Explanations

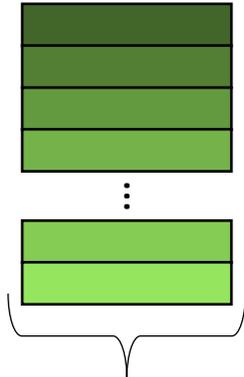**How can we identify these region pairs important to the model?**
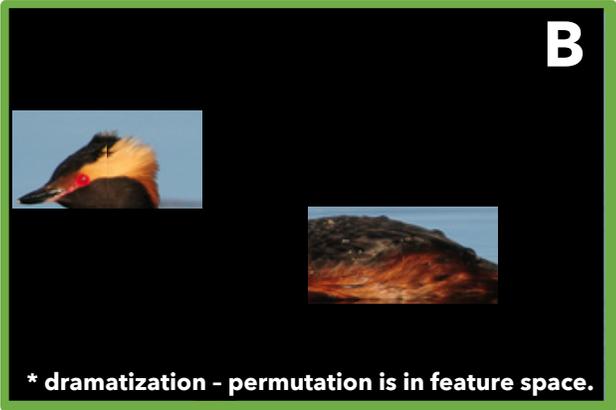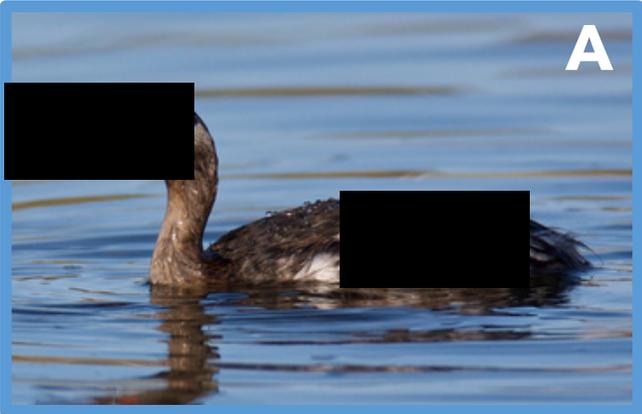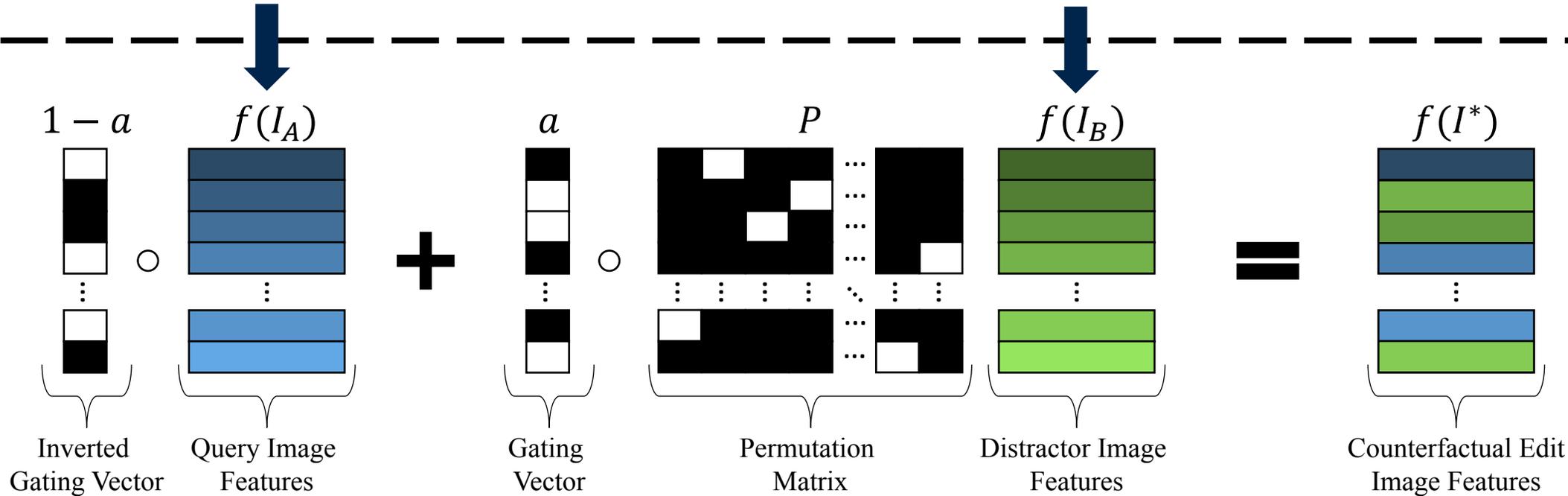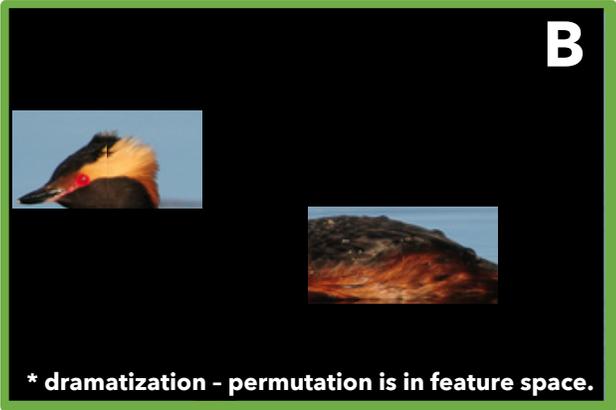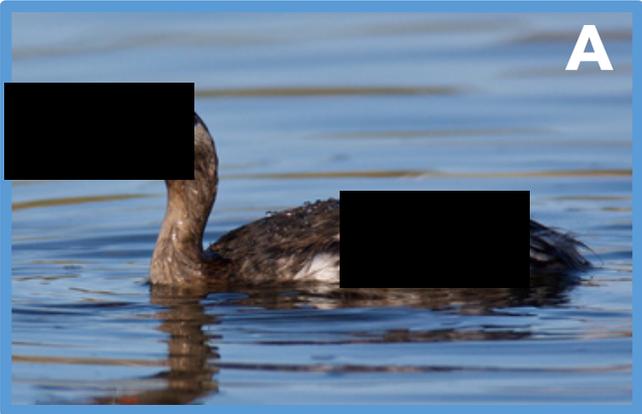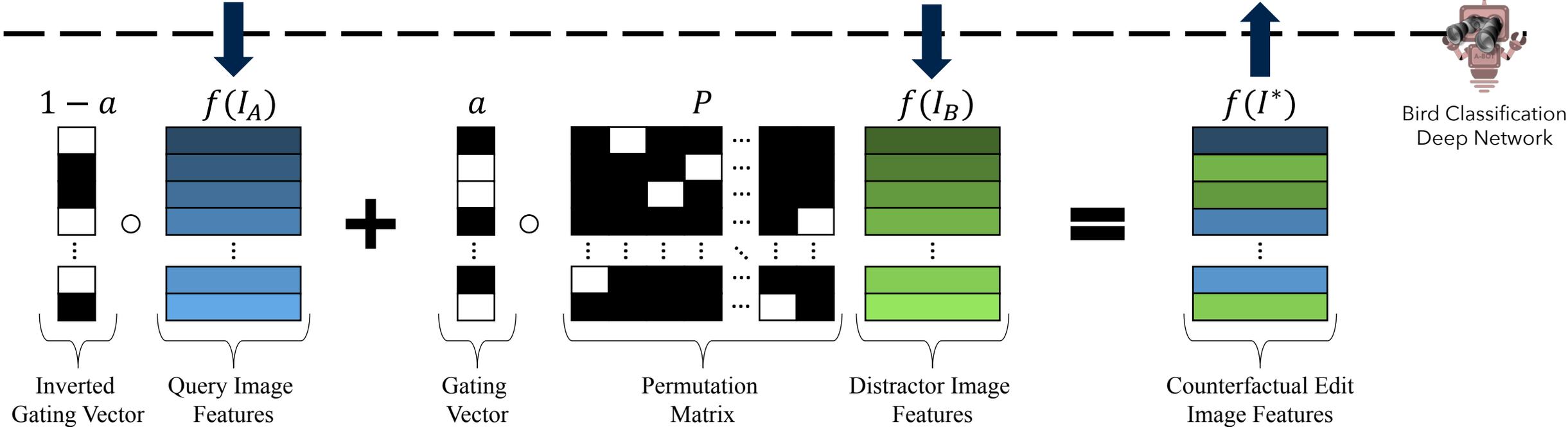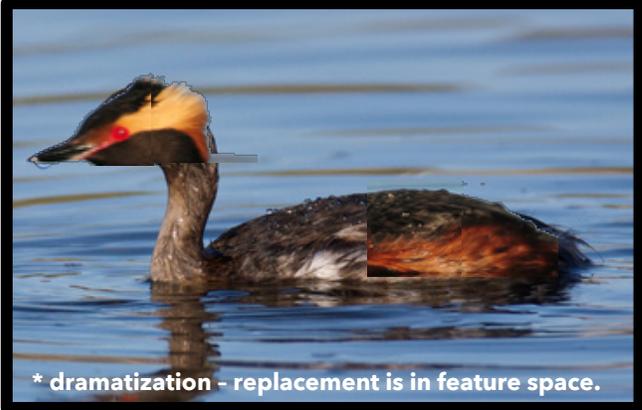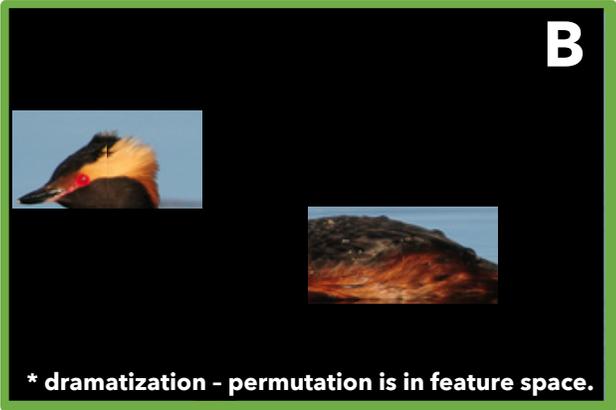


Bird Classification
Deep Network

Image I

# Counterfactual Visual Explanations

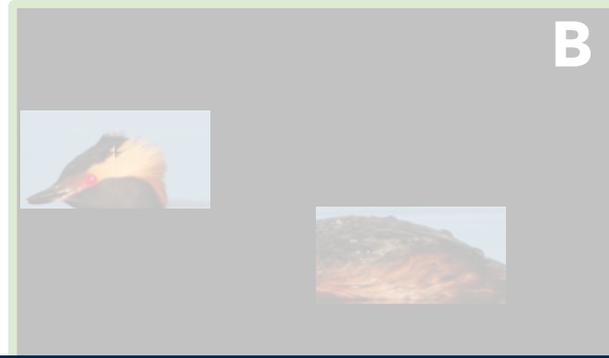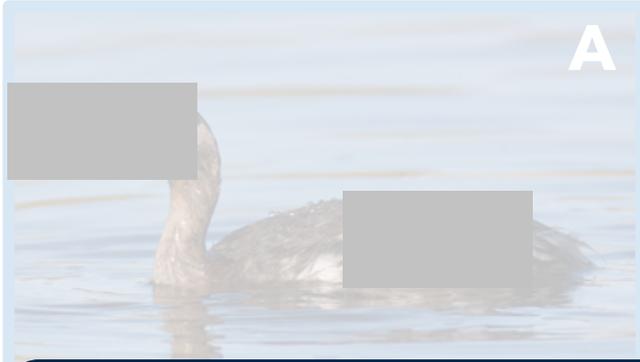## How can we identify these region pairs important to the model?



Bird Classification
Deep Network

Spatial Features

$d$

$h$

$w$

$d$

$hw$

$f(I)$

Spatial Feature Extractor

Image I

# Counterfactual Visual Explanations

## How can we identify these region pairs important to the model?



Bird Classification
Deep Network

Image $I$

Spatial
Features

$h$

$w$

$d$

$d$

$hw$

$\log P(y_1|I)$
$\log P(y_2|I)$
$\vdots$
$\log P(y_{|Y|}|I)$

$f(I)$
Spatial Feature Extractor

$g(f(I))$
Decision Network

# Counterfactual Visual Explanations



A

B

$f(I_A)$

$f(I_B)$

Bird Classification
Deep Network

Query Image
Features

Distractor Image
Features

# Counterfactual Visual Explanations



$f(I_A)$

$P$

$f(I_B)$
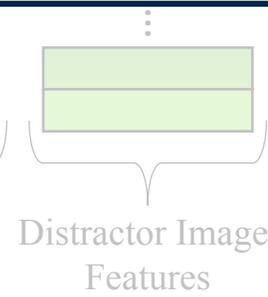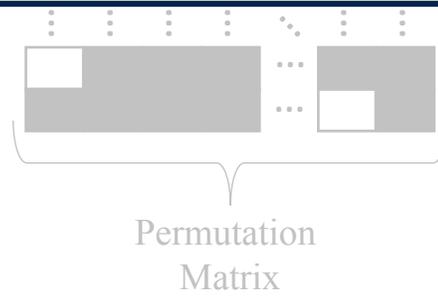
* dramatization – permutation is in feature space.

Bird Classification
Deep Network

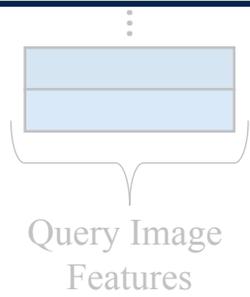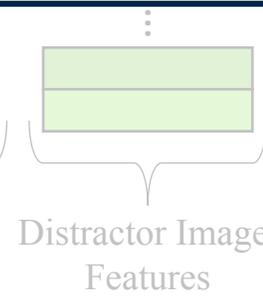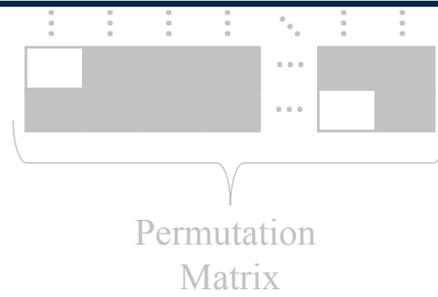Query Image
Features

Permutation
Matrix

Distractor Image
Features

# Counterfactual Visual Explanations

# Counterfactual Visual Explanations



* dramatization - permutation is in feature space.

$1 - a$     $f(I_A)$     $a$     $P$     $f(I_B)$     $f(I^*)$

Bird Classification Deep Network

Inverted Gating Vector    Query Image Features    Gating Vector    Permutation Matrix    Distractor Image Features    Counterfactual Edit Image Features

# Counterfactual Visual Explanations



$$1 - a \quad f(I_A) \qquad a \qquad P \qquad f(I_B) \qquad f(I^*)$$

Bird Classification Deep Network

Inverted Gating Vector | Query Image Features | Gating Vector | Permutation Matrix | Distractor Image Features | Counterfactual Edit Image Features

* dramatization - permutation is in feature space.

* dramatization - replacement is in feature space.
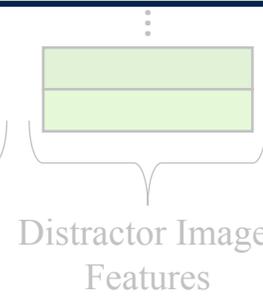
# Counterfactual Visual Explanations
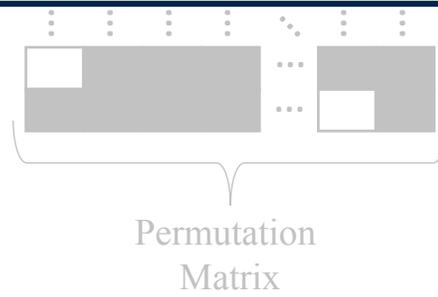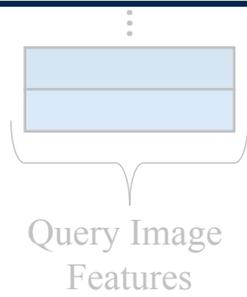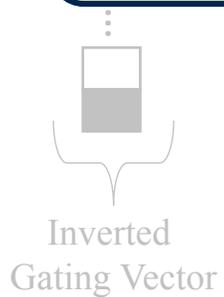


**Counterfactual Visual Explanation Generation:** Find
1) binary gating vector $a$, and
2) a permutation matrix $P$

Inverted Gating Vector | Query Image Features | Gating Vector | Permutation Matrix | Distractor Image Features | Counterfactual Edit Image Features

# Counterfactual Visual Explanations



**Counterfactual Visual Explanation Generation:** Find
1) binary gating vector $a$, and
2) a permutation matrix $P$
such that the model changes its decision to the distractor class

Classification Network

Inverted Gating Vector

Query Image Features

Gating Vector

Permutation Matrix

Distractor Image Features

Counterfactual Edit Image Features

# Counterfactual Visual Explanations



**Counterfactual Visual Explanation Generation:** Find
1) binary gating vector $a$, and
2) a permutation matrix $P$
such that the model changes its decision to the distractor class
***with the fewest edits (i.e. $min \ \|a\|_1$)***

Inverted Gating Vector    Query Image Features    Gating Vector    Permutation Matrix    Distractor Image Features    Counterfactual Edit Image Features

# Results – Single Edit

If the highlighted region in **image A**
looked like the highlighted region in **image B**,
*then image A is more likely to be classified as **class B***.

Query Image A

Distractor Image B



Eared Grebe

Horned Grebe

# Results – Single Edit

If the highlighted region in **image A**
looked like the highlighted region in **image B**,
*then image A is more likely to be classified as **class B**.*



Query Image A

Distractor Image B

Composite Image
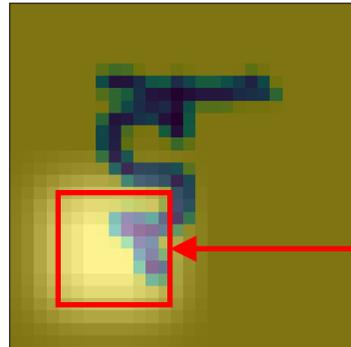(for visualization only)
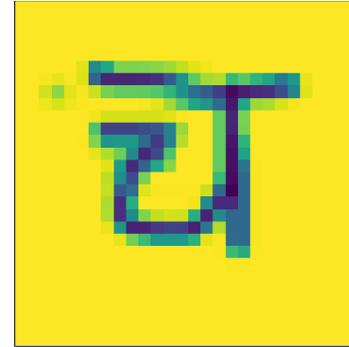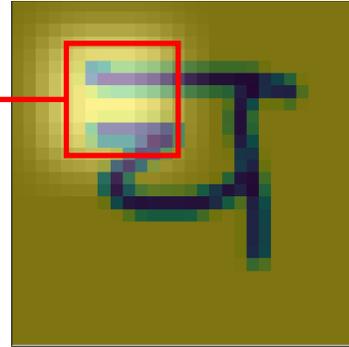
Eared Grebe

Horned Grebe
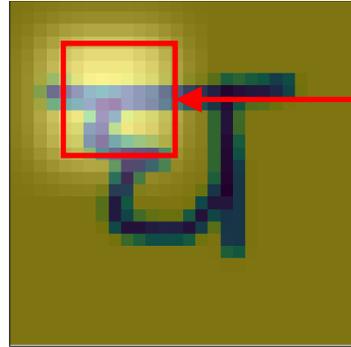
# Results – Single Edit



Query
Image A

Distractor
Image B

Composite Image
(for visualization only)
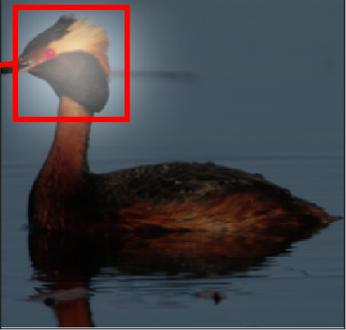
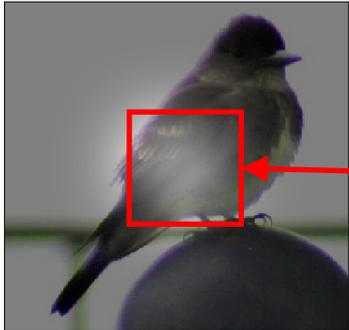# Results – Single Edit

# Results - Single Edit

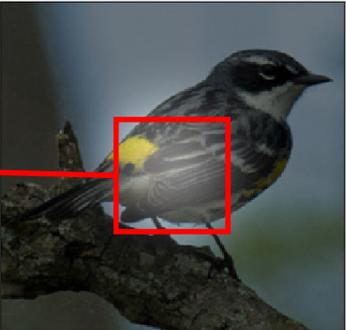Query Image A — Eared Grebe
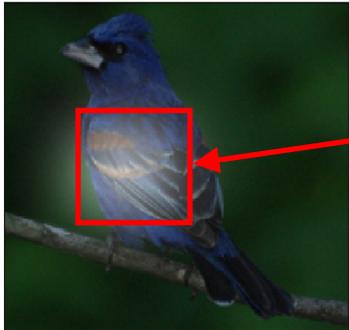
Distractor Image B — Horned Grebe
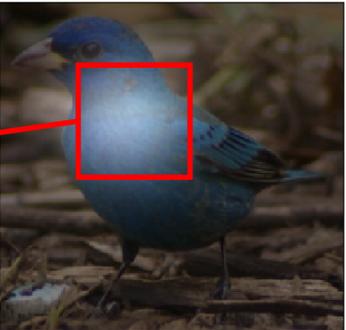
Composite Image (for visualization only)

Olive sided Flycatcher

Myrtle Warbler

Blue Grosbeak

Indigo Bunting

# Machine Teaching – Bird Classification

Do our **counterfactual explanations** help untrained participants learn to identify fine-grained classes?

# Machine Teaching – Bird Classification

Do our **counterfactual explanations** help untrained participants learn to identify fine-grained classes?

# Machine Teaching – Bird Classification

Do our **counterfactual explanations** help untrained participants learn to identify fine-grained classes?
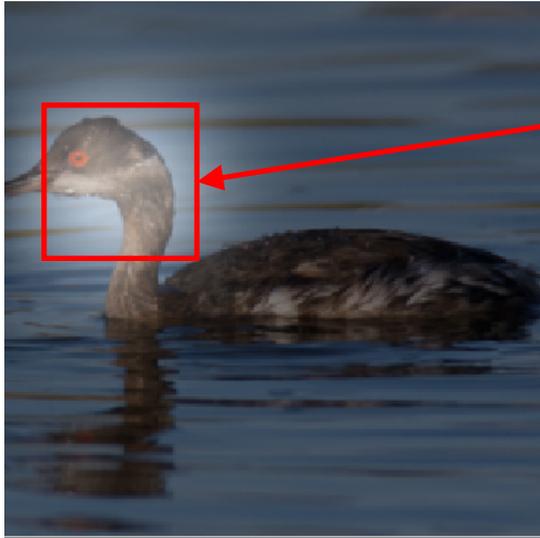
# Machine Teaching – Bird Classification

Do our **counterfactual explanations** help untrained participants learn to identify fine-grained classes?
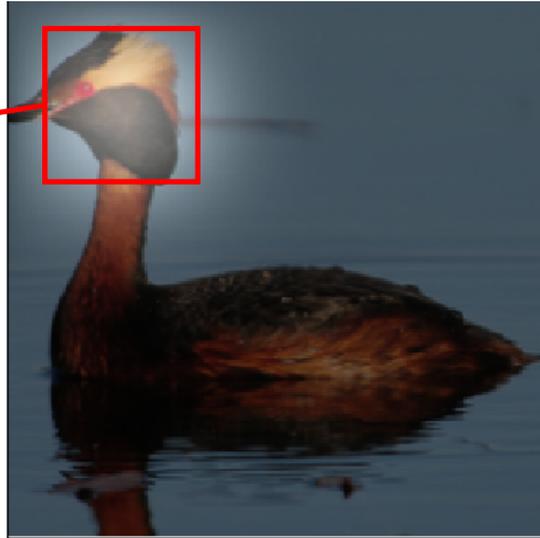
# Counterfactual Visual Explanations
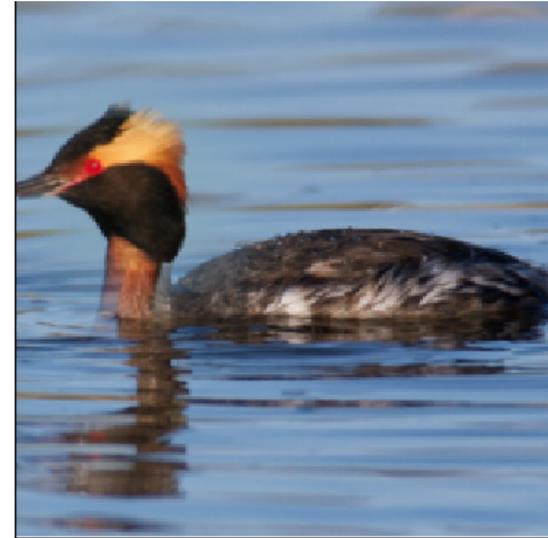


Query Image A

Eared Grebe

Distractor Image B

Horned Grebe

Composite Image
(for visualization only)

# Questions?

Stop by our poster at #149 in Pacific Ballroom!