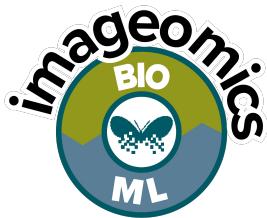




THE OHIO STATE  
UNIVERSITY



ICLR

# A Simple Interpretable Transformer for Fine-grained Image Classification and Analysis



Dipanjyoti Paul



Arpita Chowdhury



Xinqi Xiong



Feng-Ju Chang

...



Charles Stewart



Tanya Berger-Wolf



Yu Su



Wei-Lun Chao



PRINCETON  
UNIVERSITY

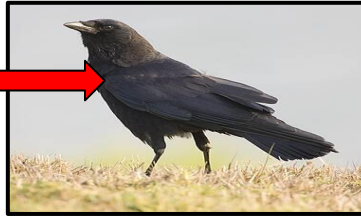


# What Kind of Interpretation we are looking for?

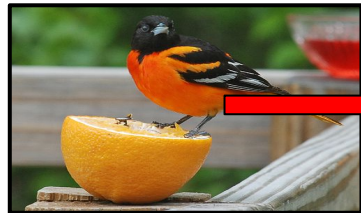
---



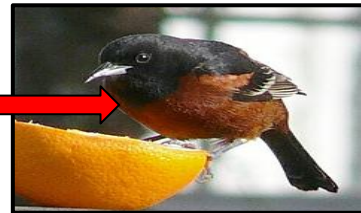
Red-winged  
Blackbird



American Crow



Baltimore Oriole



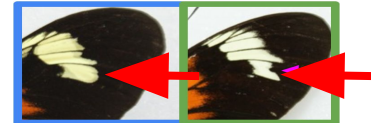
Orchard Oriole



Heliconius  
melpomene



Heliconius  
elevatus



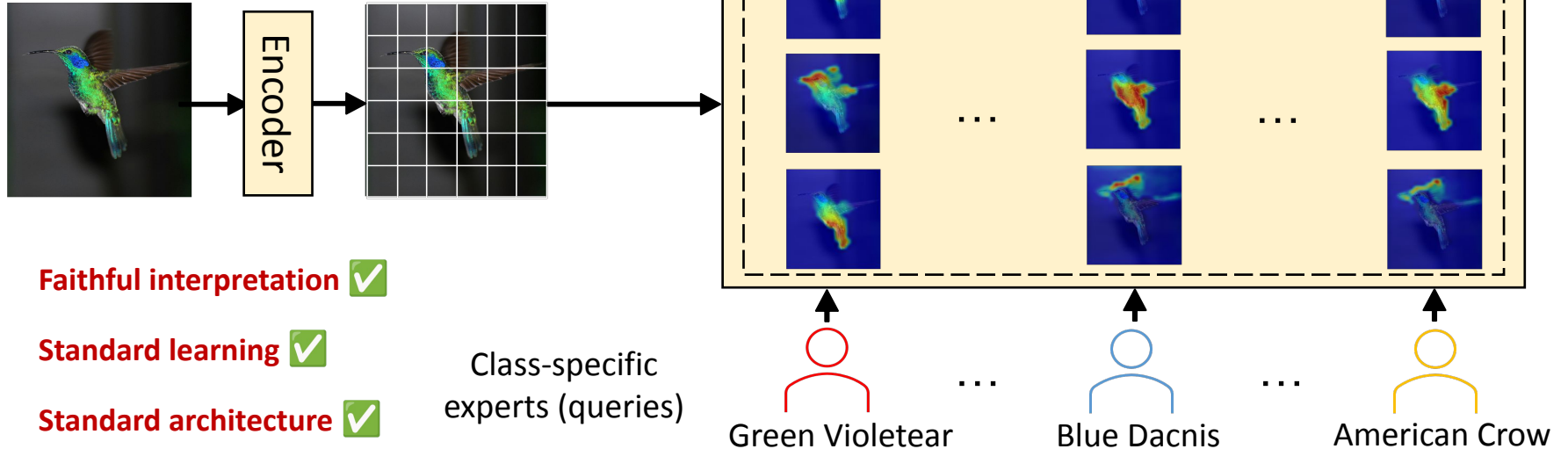
# Motivation (INterpretable TRansformer)

---

- Can we obtain interpretability via **standard neural network architectures** and **standard learning algorithms**?
- Can we have a **faithful** interpretation?

# INTR – Interpretable Transformer for Trait Discovery

**INTR:** inspired by how biologists looking for traits to differentiate species



- **Faithful interpretation** ✓
- **Standard learning** ✓
- **Standard architecture** ✓

# INTR can consistently localize traits

Images of  
Painted Bunting

Trait 1

Trait 2

Trait 3

Trait 4

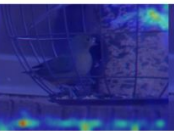
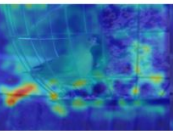
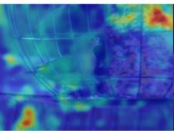
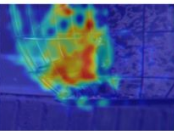
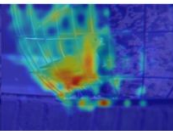
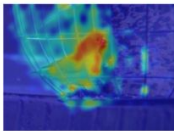
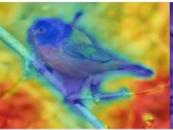
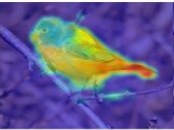
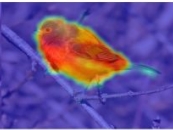
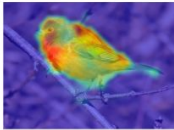
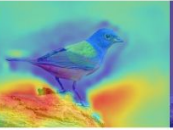
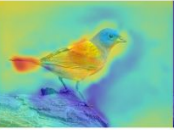
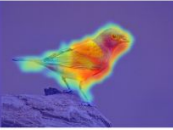
Trait 5

Trait 6

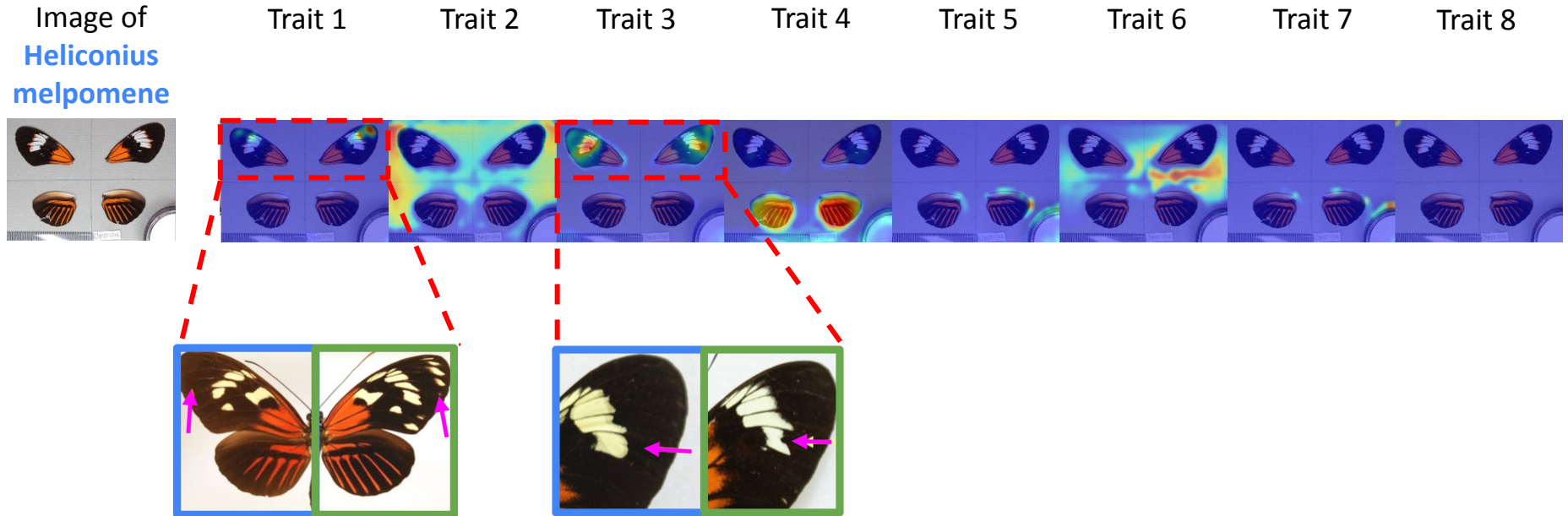
Trait 7

Trait 8

Recognition as  
Painted  
Bunting



# INTR can localize traits for fine-grained species



Source: <https://www.cliniquevetodax.com/Heliconius/index.html>

# INTR can localize traits for fine-grained species

Is this Red-winged Blackbird?

Image

Trait 1

Trait 2

Trait 3

Trait 4

Trait 5

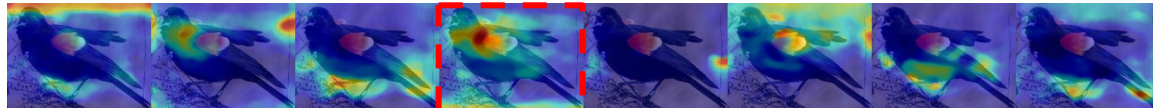
Trait 6

Trait 7

Trait 8



Red-winged Blackbird



American Crow



Is this Orchard Oriole?

Image

Trait 1

Trait 2

Trait 3

Trait 4

Trait 5

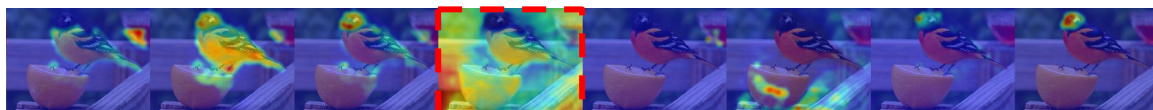
Trait 6

Trait 7

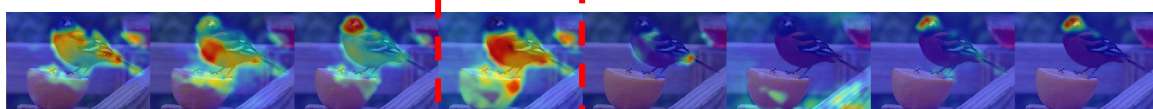
Trait 8



Baltimore Oriole



Orchard Oriole



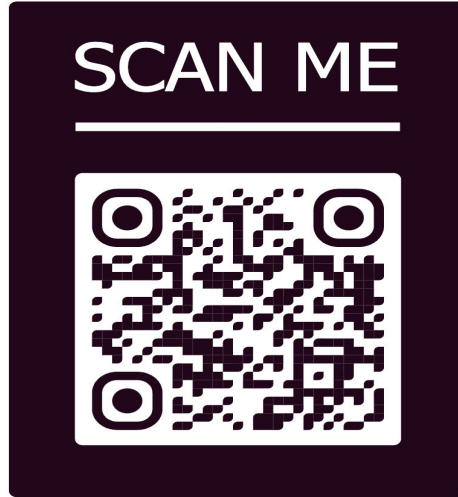
# Summary

---

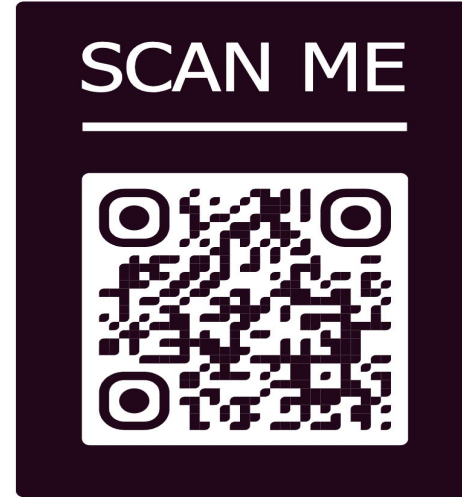
- INTR can able to detect ML traits automatically and provides supporting evidence for its decision
- INTR is build upon standard NN architecture and learning
- Applicability of INTR has been shown on eight datasets
- INTR can offer a new way of thinking about interpretable machine learning



# Thank You!



INTR Github



INTR Paper

