



## Improving tools for citizen scientists

Wouter Koch Digital transformation projects 02.12.2022



A person who...



A person who...

is there with the species...



A person who...

is there with the species...

observes it...



A person who...

is there with the species...

observes it...

recognizes it...



A person who...

is there with the species...

observes it...

recognizes it...

and reports it



A person who...

is there with the species...

observes it...

recognizes it...

and reports it

All biases, especially in CS



What is lacking in the available citizen science data? What can be done to collect better data?



### Paper II

Identification is fundamental

Al is helping more and more

Data are taxonomically biased

# Maximizing citizen scientists' contribution to automated species recognition

Wouter Koch, Laurens Hogeweg, Erlend B. Nilsen & Anders G. Finstad

Scientific Reports, 2022: doi:10.1038/S41598-022-11257-X

## Paper II

Maximizing citizen scientists' contribution to automated species recognition

Identification is fundamental

Al is helping more and more

Data are taxonomically biased

data? Which images would help Al models the most?

How does that affect image

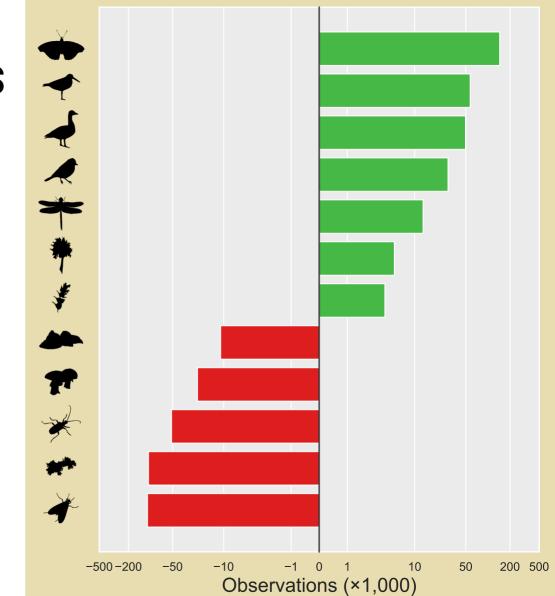
Wouter Koch, Laurens Hogeweg, Erlend B. Nilsen & Anders G. Finstad Scientific Reports, 2022: doi:10.1038/S41598-022-11257-X

### **Taxonomic bias**

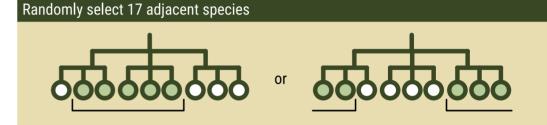
Known issue, e.g. Troudet et al.

We tested and confirm this within Norwegian Citizen Science image data

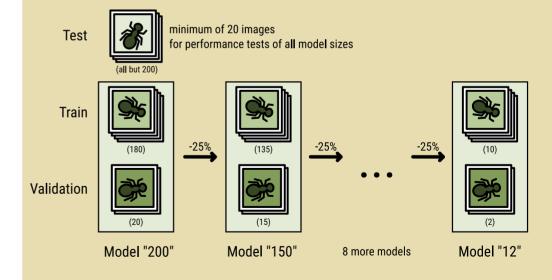
Selection of 12 taxonomic orders



Goal: find the *Value of Information* of adding images per taxon

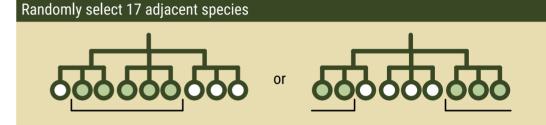


#### For every selected species, divide images for model training

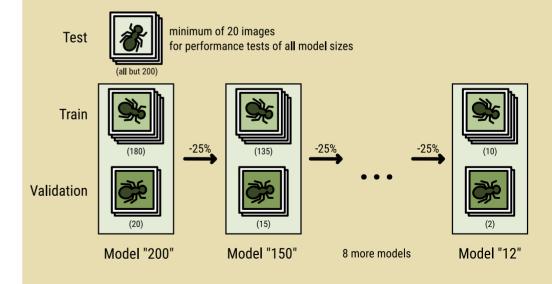


Goal: find the *Value of Information* of adding images per taxon

Take 17 species for one order



#### For every selected species, divide images for model training



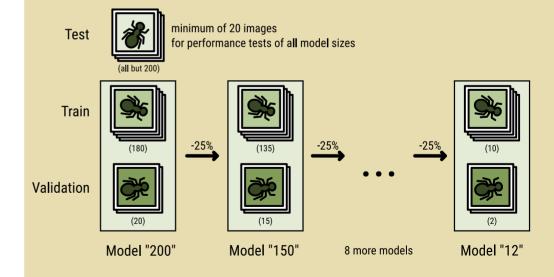
Goal: find the *Value of Information* of adding images per taxon

Take 17 species for one order

Train on 200, 150 ... 12 images per species (25% reduction)

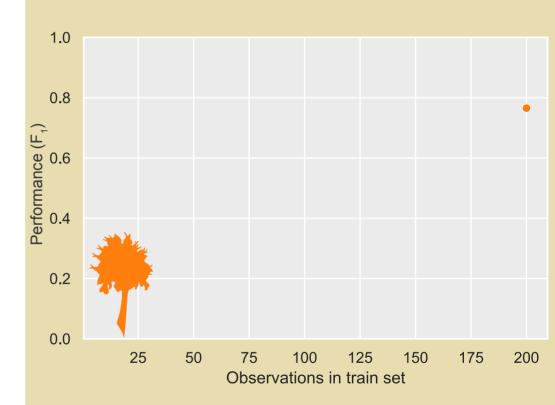
## Randomly select 17 adjacent species or or

#### For every selected species, divide images for model training



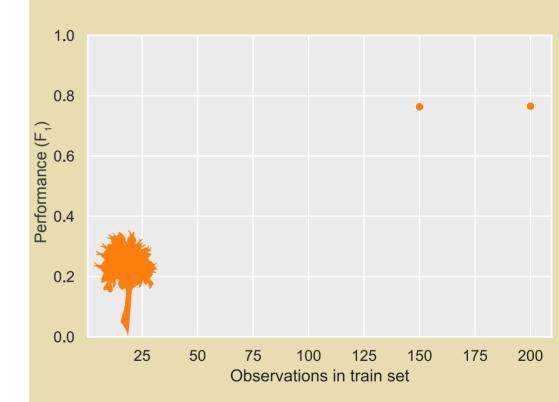
Goal: find the *Value of Information* of adding images per taxon

Take 17 species for one order



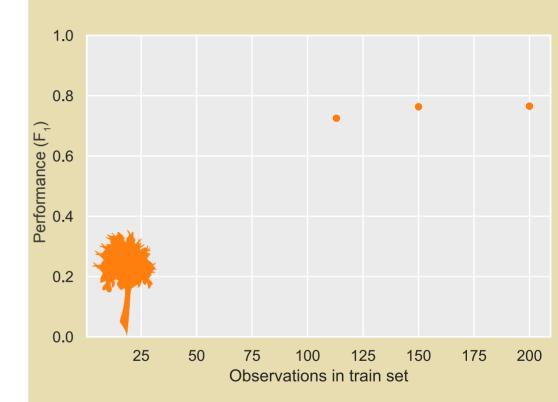
Goal: find the *Value of Information* of adding images per taxon

Take 17 species for one order



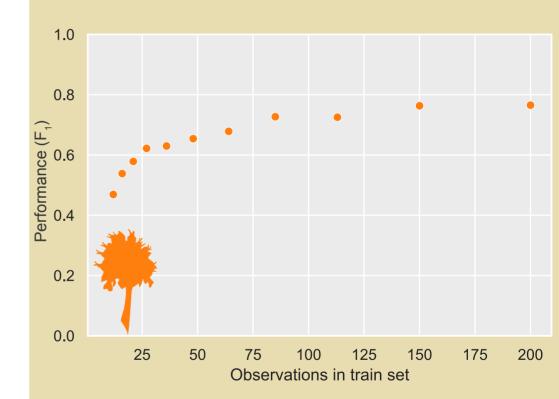
Goal: find the *Value of Information* of adding images per taxon

Take 17 species for one order



Goal: find the *Value of Information* of adding images per taxon

Take 17 species for one order

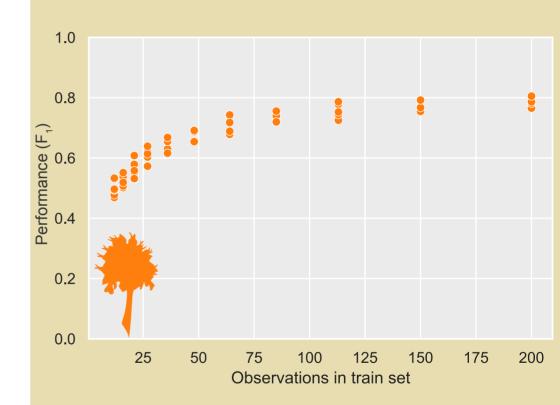


Goal: find the *Value of Information* of adding images per taxon

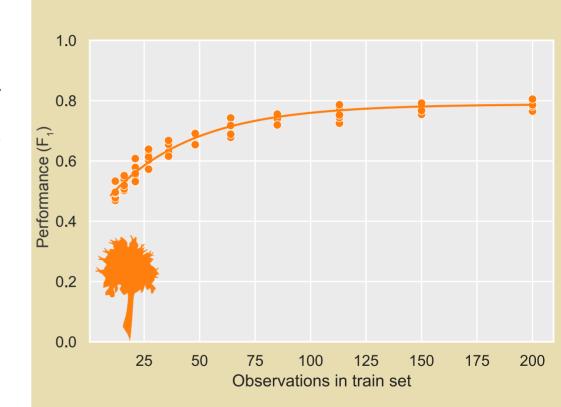
Take 17 species for one order

Train on 200, 150 ... 12 images per species (25% reduction)

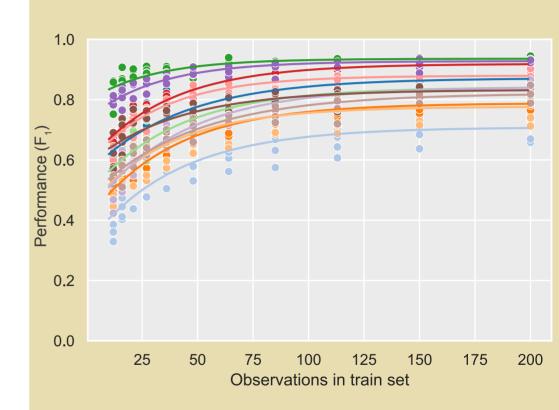
5 times per order



Gives performance curves over images per species

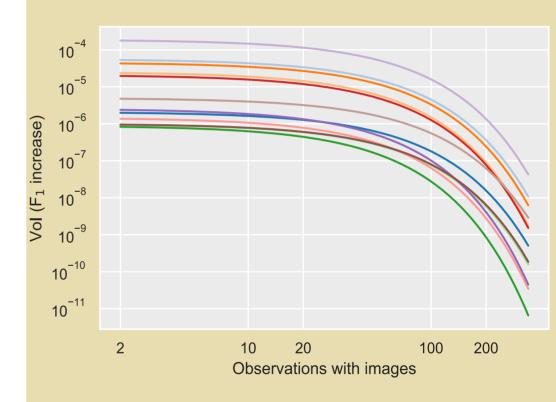


Gives performance curves over images per species



Gives performance curves over images per species

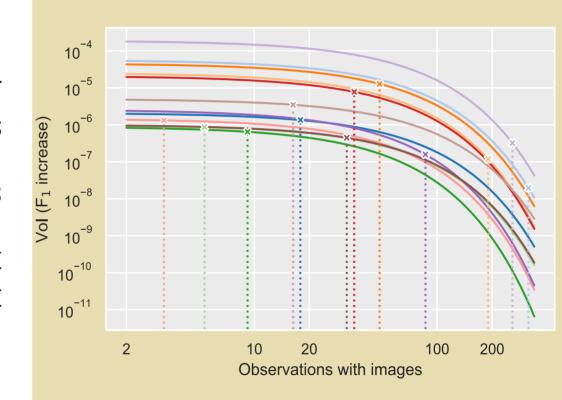
Get the slope of the fitted curves



Gives performance curves over images per species

Get the slope of the fitted curves

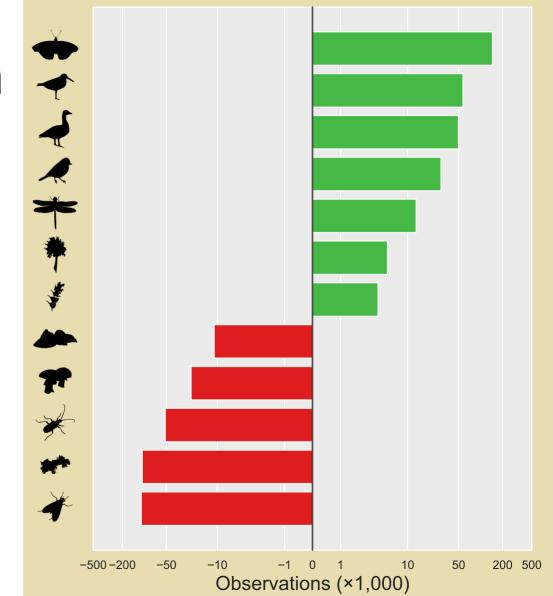
Look up the Value of Information at the current amount



Gives performance curves over images per species

Get the slope of the fitted curves

Look up the Value of Information at the current amount



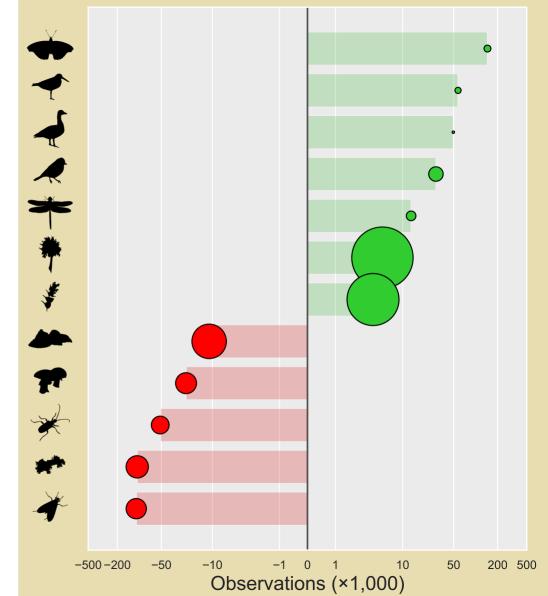
Gives performance curves over images per species

Get the slope of the fitted curves

Look up the Value of Information at the current amount

Value not just in most scarce

Opportunity to focus resources



### Paper III

Taxonomic bias is equated to popularity

## Recognizability bias in citizen science photographs

Wouter Koch, Laurens Hogeweg, Erlend B. Nilsen, Robert B. O'Hara & Anders G. Finstad

Preprint on bioRχiv: doi:10.1101/2022.06.25.497604

## Paper III

Recognizability bias in citizen science photographs

Taxonomic bias is equated to popularity

Wouter Koch, Laurens Hogeweg, Erlend B. Nilsen, Robert B. O'Hara & Anders G. Finstad

What about recognizability?

Preprint on bioRχiv: doi:10.1101/2022.06.25.497604

Is AI biased in the same way as humans?

What does that mean for training AI models?

# Performance vs popularity

Same 12 orders as in paper II

Deep dive into birds: most prevalent and good standardized trait data

Train AI models with 200 images per species: model does not "know" which are reported more



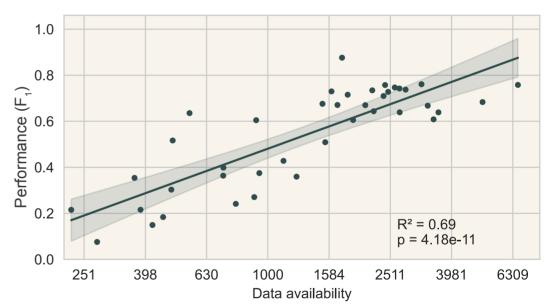
# Performance vs popularity

Same 12 orders as in paper II

Deep dive into birds: most prevalent and good standardized trait data

Train AI models with 200 images per species: model does not "know" which are reported more

Still, model does better on more popular species. Why?



Close by or better zoom?

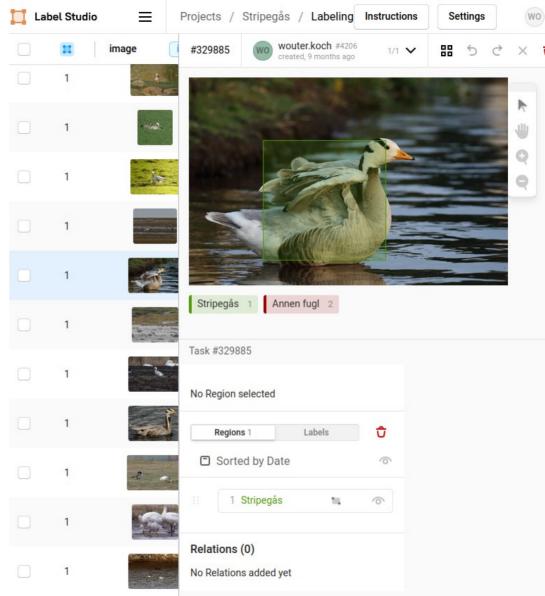
Fewer other species?



Close by or better zoom?

Fewer other species?

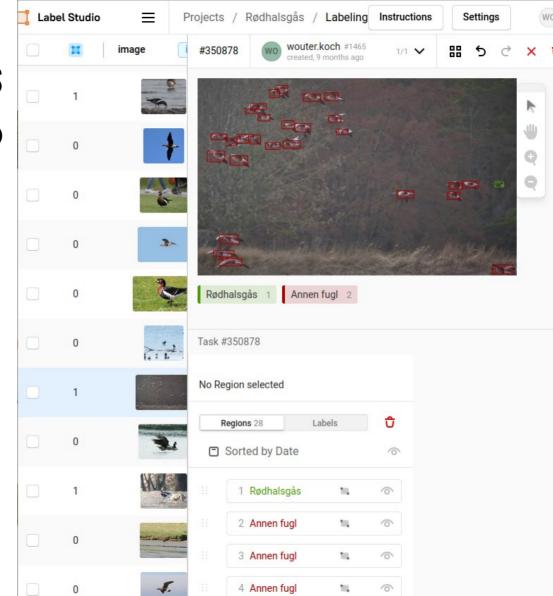
Pictures annotated with help from volunteers



Close by or better zoom?

Fewer other species?

Pictures annotated with help from volunteers

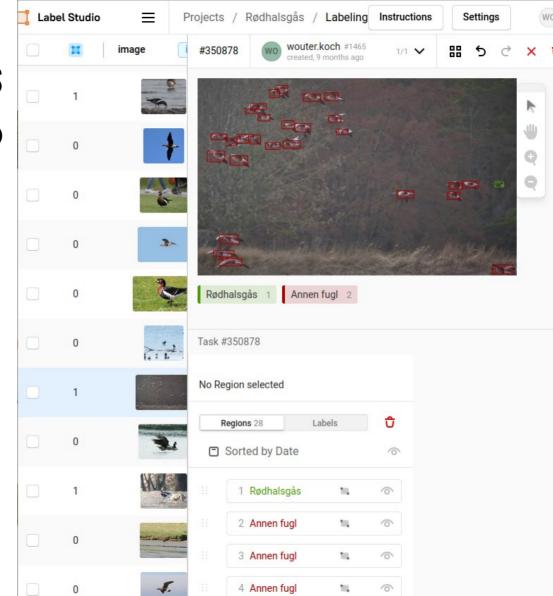


Close by or better zoom?

Fewer other species?

Pictures annotated with help from volunteers

Provides image quality measures



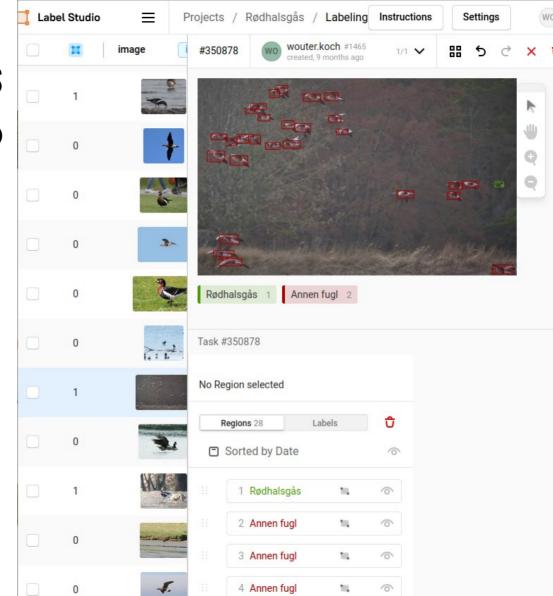
Close by or better zoom?

Fewer other species?

Pictures annotated with help from volunteers

Provides image quality measures

Does not correlate with performance



#### Is it the kind of bird?

Larger birds, different habitat, migratory, other behavior?



#### Is it the kind of bird?

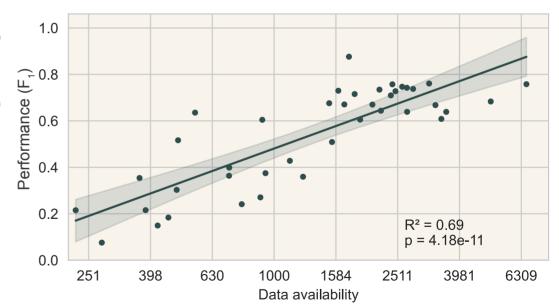
Larger birds, different habitat, migratory, other behavior?

None of these correlate with the Al model performance



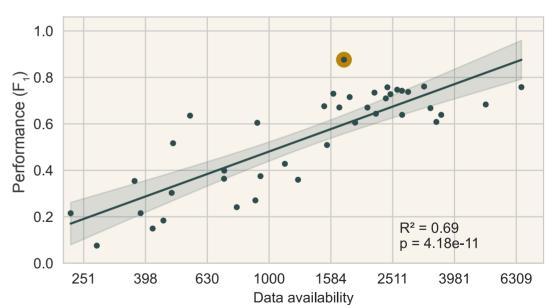
The AI model finds some species easier to recognize

People do too, and report more



The AI model finds some species easier to recognize

People do too, and report more

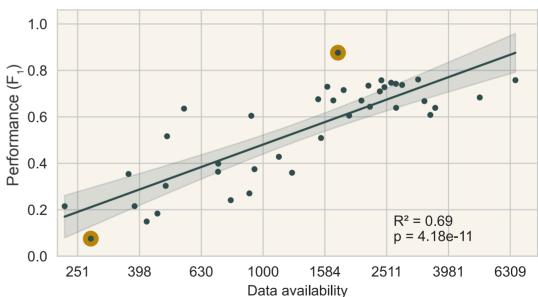


The AI model finds some species easier to recognize

People do too, and report more



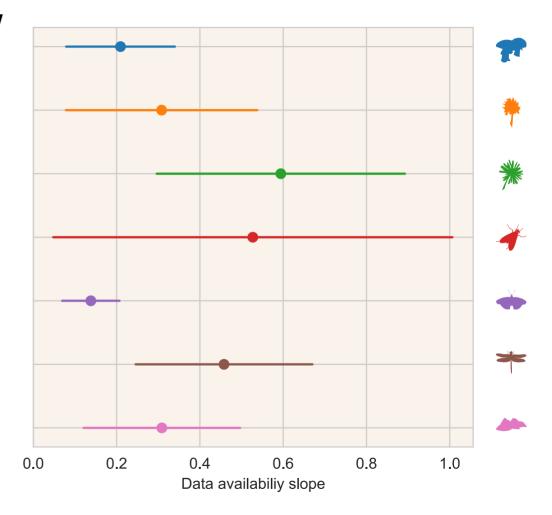




The AI model finds some species easier to recognize

People do too, and report more

Same pattern for other orders

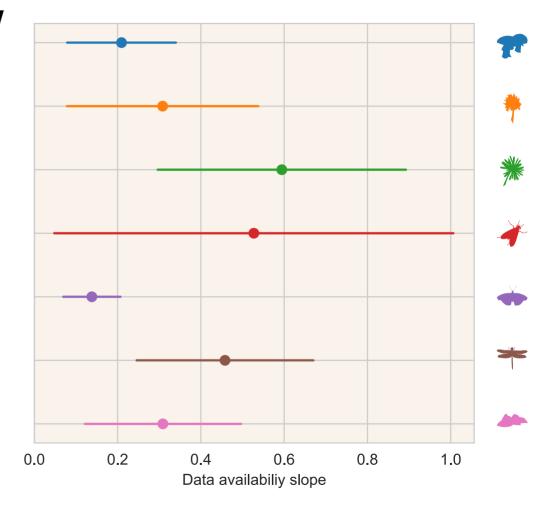


The AI model finds some species easier to recognize

People do too, and report more

Same pattern for other orders

Increased data likely of easier species



## **Final thoughts**

Citizen science = valuable data + engagement

We need to cope with bias and other issues in citizen science

But there is more we can do than simply accepting the data "as is"

