# A Pipeline for Processing Images in iDigBio to Train Convolutional Neural Networks

Gaurav Yeole, Matthew Collins, Alex Thompson, Renato Figueiredo

iDigBio
Integrated Digitized Biocollections

ACIS

# iDigBio data and portal

# 27M+ images stored



Wealth of data; information largely untapped
Metadata available; image processing – human eye+brain

# iDigBio storage back-end
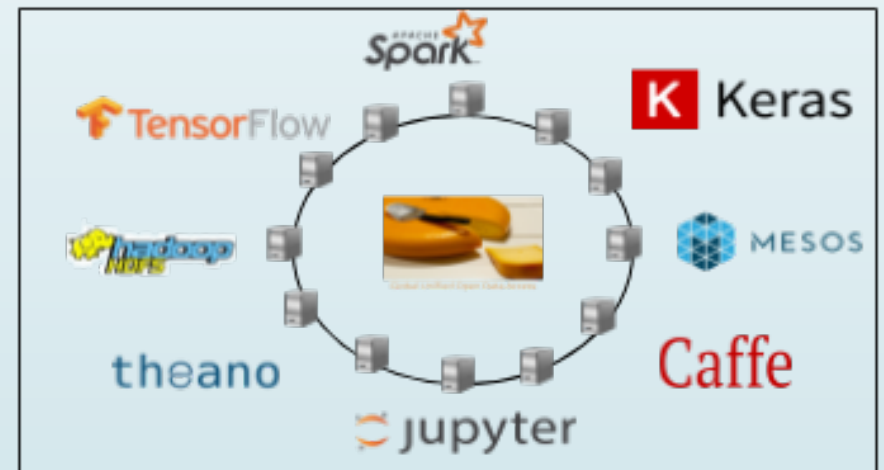
- 115+ million specimen records
- 27+ million media records
- 120 terabytes of media files
- Hosted on Ceph Ditributed storage system
- Opportunity: machine learning processing near data



iDigBio Database of Images

# GUODA – Global Unified Open Data Access

- Apache Mesos – distributed systems kernel
- Apache Spark – fast and general engine for large-scale data processing
- Jupyter Notebook Interface for Python and R
- Installed frameworks for deep learning –
  - TensorFlow
  - Keras
  - Theono
  - Caffe



GUODA cluster

# Applications of Artificial Intelligence (AI)
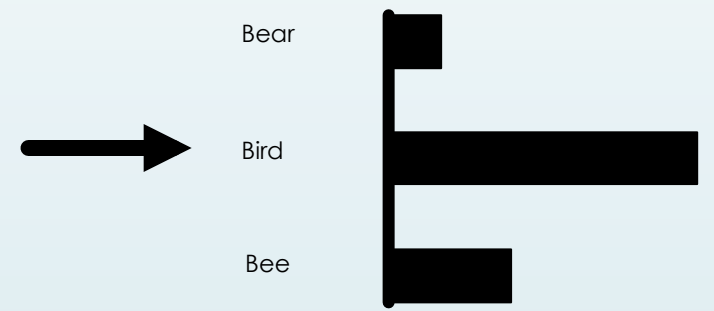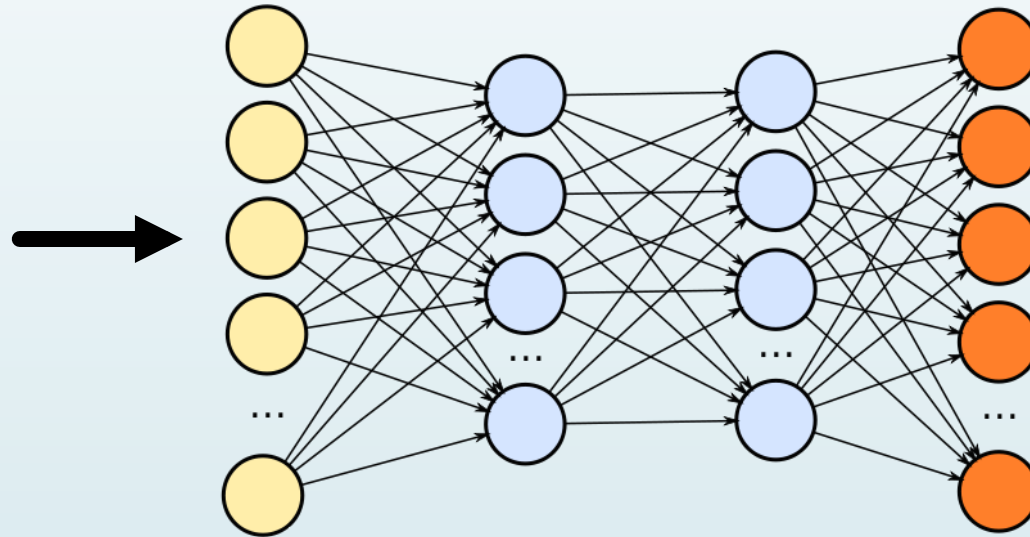
- Decision-making
- System Identification and Control
  - Vehicle control, process control, natural resource management
- Pattern Recognition
  - Face identification, object recognition, etc.
- Classification
  - Image classification, e-mail spam filtering

- **Emerging applications of AI for specimen images?**

# Use Case – Find Mercury Contamination on Herbarium Specimens

- Mercury salt used to specimens as insecticides

- Mercury marking is visible on the specimens

- Mercury vapor and mercury compounds pose threat to human health.

- Mercury contaminated images can be classified using convolutional neural networks.

- E. Schuettpelz, P. Frandsen, R. Dikow, A. Brown, S. Orli, M. Peters, A. Metallo, V. Funk, L. Dorr. (2017). Applications of deep convolutional neural networks to digitized natural history collections. Biodiversity Data Journal. 5. e21139. 10.3897/BDJ.5.e21139.

Convolutional Neural Networks (CNN)

# Accessing iDigBio Images

# Spark dataframe of iDigBio media

```
+-------------------+------------+----------------------+-------------------+
|          accessuri|contaminated|image_nparray(accessuri)|     vector_images|
+-------------------+------------+----------------------+-------------------+
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          1|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
|http://collection...|          0|   [255.0, 216.0, 25...|[255.0,216.0,255....|
+-------------------+------------+----------------------+-------------------+
only showing top 20 rows
```

# Image Processing Pipeline



iDigBio Database of Images

GUODA cluster

Spark Dataframe

Internet

iDigBio Portal

GUODA Jupyter Notebook

# Opportunities for collaboration

**Founded in 2002**

**Focused on researchers and institutions on the Pacific Rim**

**Open Community of Practice**

**Engages "Long Tail" science communities**

**27 Institutional Members**

# AI-focused high-performance clusters

ABCI – top #5 in the world; similar system at NCHC, Taiwan

# Future Work

- Join specimen metadata from iDigBio with images in dataframes which can serve as image and label pairs to train more exciting deep learning models

- Build another audio signal processing pipeline for existing audio data in iDigBio

- Make this pipeline portable and make it compatible to run on infrastructures such as NSF's XSEDE, AIST's ABCI

# Acknowledgements

- Sylvia Orli, Paul Frandsen, Rebecca Dikow, Smithsonian Institution