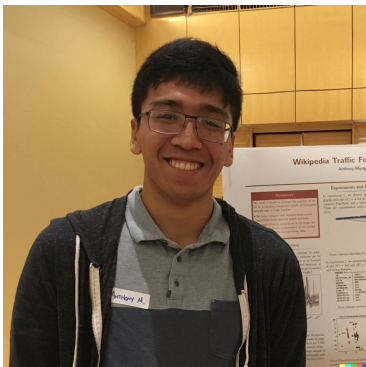# Running a Kaggle Competition Team

## Insights and Lessons from BirdCLEF with Data Science @ Georgia Tech

Anthony Miyaguchi
acmiyaguchi@gatech.edu

Georgia Institute of Technology

2023-05-03

# Who am I?



- A Software Engineer
  - 5 years as a Data Engineer at Mozilla
  - 1.5 years as a Software Engineer at Planet Labs
- OMSCS, matriculated Spring 2022
  - B.S. Computer Science and Engineering from UCLA 2016
  - Graduate Certificate from Stanford Center for Professional Development (SCPD) 2018
- Career focus on scalable data systems and machine learning

# Overview

- What is Kaggle?
- Why Kaggle?
- BirdCLEF 2022
  - Recruitment and Technical Approach
- BirdCLEF 2023
  - Updated Recruitment and Technical Approach
- Thoughts and Advice
- Q&A

# What is Kaggle?



Figure 1: Kaggle: a platform for data science competitions

Figure 2: Wikipedia page view data in the style of the "Web Traffic Time Series Forecasting" Kaggle Competition

# Temporal Regularized Matrix Factorization for High-dimensional Time Series Prediction

**Hsiang-Fu Yu**
University of Texas at Austin
rofuyu@cs.utexas.edu

**Nikhil Rao**
Technicolor Research
nikhilrao86@gmail.com

**Inderjit S. Dhillon**
University of Texas at Austin
inderjit@cs.utexas.edu

## Abstract

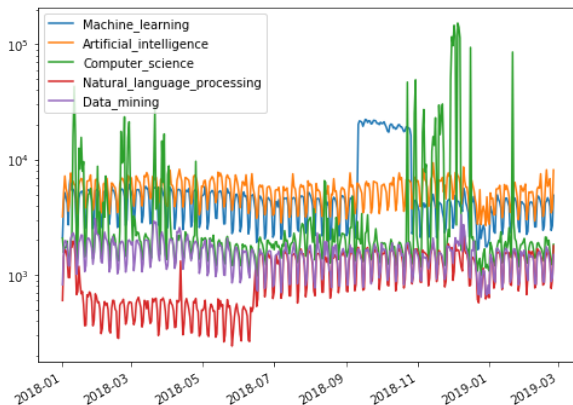Time series prediction problems are becoming increasingly high-dimensional in modern applications, such as climatology and demand forecasting. For example, in the latter problem, the number of items for which demand needs to be forecast might be as large as 50,000. In addition, the data is generally noisy and full of

Figure 3: The first piece of literature I implemented outside of school work as part of the "Web Traffic Time Series Forecasting" Kaggle Competition

Figure 4: "Wikipedia Traffic Forecasting with Graph Embeddings" is a machine learning project during my time in SCPD inspired by the Kaggle competition.

# Why Kaggle?



Figure 5: The Titanic dataset is an iconic Kaggle challenge.

## Reasons to consider Kaggle

- Good start to a **research project**
  - Access to impactful real-world problems
    - bird conservation, cancer detection, etc.
- **Structure and deadlines** to stay on track
- Anyone can participate, with a focus on **practical applications and the potential to win cash prizes**

Figure 6: BirdCLEF is a competition to help with bird conservation.

# DS@GT Competition Team

## Recruitment

- Built a team of 5 people from DS@GT in Spring 2022
- 3 masters students and 2 undergraduates

## Technical Approach

- **Motif Mining** with SiMPLe-Fast
  - Find salient sections of audio for classification
- **Unsupervised Representation Learning** via Triplet Loss
  - Train an embedding for downstream tasks



Figure 7: DS@GT: a student-run data science organization

# Why is audio classification challenging?



Figure 8: xeno-canto is a crowd sourced database of bird sounds.

# Reading the literature

## Motif Mining - SiMPLe-Fast

Silva, D. F., Yeh, C. C. M., Zhu, Y., Batista, G. E., & Keogh, E. (2018).
Fast similarity matrix profile for music analysis and exploration. IEEE
Transactions on Multimedia, 21(1), 29-38.

## Spatial Embeddings - Tile2Vec

Jean, N., Wang, S., Samar, A., Azzari, G., Lobell, D., & Ermon, S. (2018).
Tile2Vec: unsupervised representation learning for spatially distributed data.
arXiv.

# Motif Mining with SiMPLe-Fast



Figure 9: SiMPLe-Fast is a motif mining algorithm used to find all pairs similarity in a time-series.

Figure 10: We train an embedding using mined motifs via a triplet loss

# Motif Mining and Unsupervised Representation Learning for BirdCLEF 2022

Anthony Miyaguchi[1], Jiangyue Yu[1], Bryan Cheungvivatpant[1], Dakota Dudley[1] and Aniketh Swain[1]

[1]*Georgia Institute of Technology, North Ave NW, Atlanta, GA 30332*

**Abstract**
We build a classification model for the BirdCLEF 2022 challenge using unsupervised methods. We implement an unsupervised representation of the training dataset using a triplet loss on spectrogram representation of audio motifs. Our best model performs with a score of 0.48 on the public leaderboard.

Figure 11: Working notes submitted to the CLEF 2022 conference

# Results of BirdCLEF 2022

## Best Working Notes Award

We performed poorly, but we kept good notes and had a unique approach.
We won $2,500 in Google Cloud Platform credits.

## Personal Takeaways

- Building a team is worthwhile and a forcing function for progress
- Managing a team of 5 people is challenging
- 6 weeks is not enough time for meaningful contributions to a large
  codebase

## New year, new team

- Building the team from the get-go
- Proposals, assessments, and interviews
- Reaching out to Slack, EdStem, and DS@GT

## New technical approach

- Retrain using embeddings from older models
- Build a process for machine-assisted dataset annotation
- Toy with sequence models (RNNs, Transformers, etc.)

## BirdCLEF 2023

Kaggle Competition Team, Data Science @ Georgia Tech

Contact:
    Anthony Miyaguchi <acmiyaguchi@gatech.edu>, Project Lead
    Krishi Manek <kmanek3@gatech.edu>, Director of Projects
Last Updated: 2023-01-18

Figure 12: A proposal document is a useful way to document intent and expectations.

# Recruitment: Assessment



distribution of birdcall recordings on the globe

Figure 13: Assessment in a Fall 2022 project group, using BirdCLEF metadata

**DS@GT BirdCLEF 2023 Asssessment Notebook**

embeddings

(a) Plot the BirdNET embeddings in a 2D or 3D scatterplot colored by primary_label.

(b) Describe the geometry of 1D embeddings and how you might use it to visualize data.

k-nn classification

(a) How would you compute the nearest neighbors of all points in the BirdNET embedding dataset?

(b) What percentage of primary labels are correctly predicted by BirdNET for each species?

(c) Compute labels for each point in the dataset using k-nn classification.

transfer learning

(a) Fit BirdNET embeddings to a logistic regression model

(b) Fit BirdNET embeddings to a neural network

(c) Analyze the input data in the embedding space of the second-to-last layer of the model defined in part (b)

Figure 14: Table of Contents from the BirdCLEF 2023 assessment.

## [Closed] Recruiting for DS@GT BirdCLEF 2023 Competition Team #34

**Anthony Miyaguchi**
3 months ago in **Seeking Teammates**

♡ 8

I'm Anthony Miyaguchi, an OMSCS student in my 3rd semester and a professional software engineer. I ran a projects group last year for BirdCLEF 2022 as part of the Data Science @ Georgia Tech (DS@GT) club, where we won best working notes in the Kaggle competition and $2,500 in GCP credits. This year, I am recruiting 2-3 team members for the BirdCLEF 2023 competition, which will open sometime in February. The goal is to win the working notes competition this year and to present our work at CLEF 2023 in Thessaloniki, Greece.

Figure 15: A post on the OMSCS Research EdStem board.

### Reach out to other students!

- Slack: OMSCS Study Group and OMSA Study Group
- EdStem: OMSCS Research
- Clubs: Data Science @ Georgia Tech

# Organization

## Team Structure

- Weekly meetings over Zoom once a week for 45 minutes
- Slack channel for communication
- Shared GitHub and Google Cloud Project

## Being a Project Lead

Leading a team is challenging. Every team is different, and every leader has their own style.

# Technical approach

## Outline

- Building data pipelines with Luigi
- BirdNET embeddings
- Sound Separation with MixIT
- Automated dataset annotation
- Sequence models with embeddings

# Reading the literature, yet again

## Domain specific deep learning model - BirdNET

Kahl, S., Wood, C. M., Eibl, M., & Klinck, H. (2021). BirdNET: A deep learning solution for avian diversity monitoring. Ecological Informatics, 61, 101236.

## Sound separation - MixIT

Denton, T., Wisdom, S., & Hershey, J. R. (2022, May). Improving bird classification with unsupervised sound separation. In ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 636-640). IEEE.
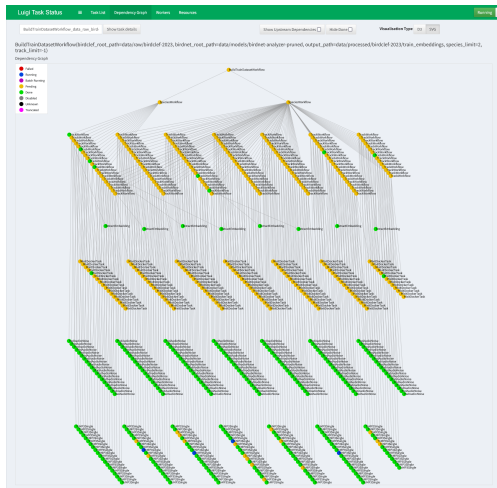
# Building data pipelines with Luigi



Figure 16: Luigi is a Python library for building data pipelines.
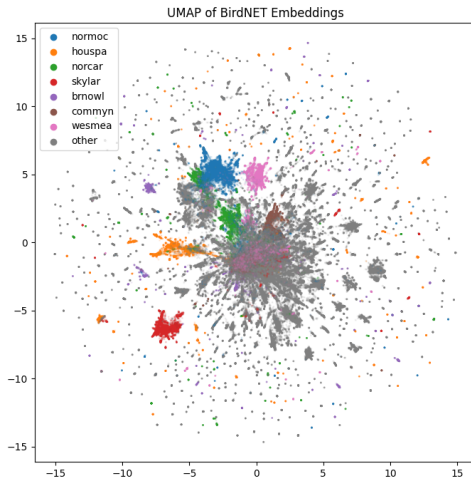
# BirdNET embeddings



Figure 17: We can use the BirdNET embedding space for search and nearest neighbor queries.
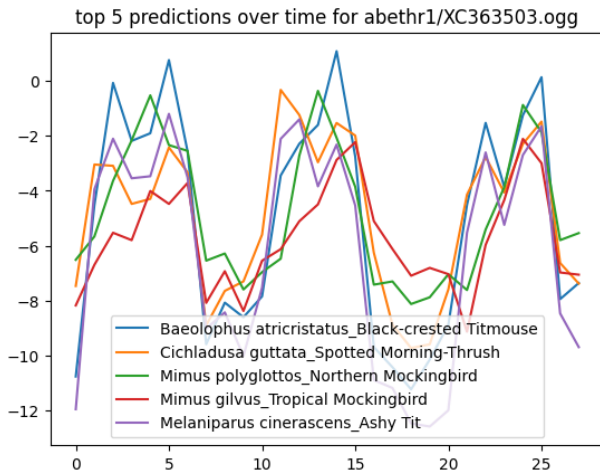
# BirdNET predictions for annotation



Figure 18: The BirdNET predictions can help with data annotation.

# Sound Separation with MixIT
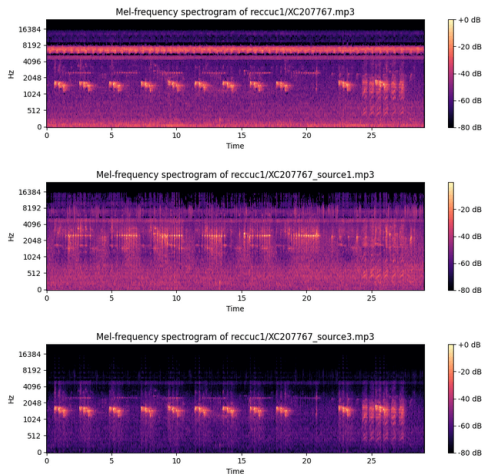


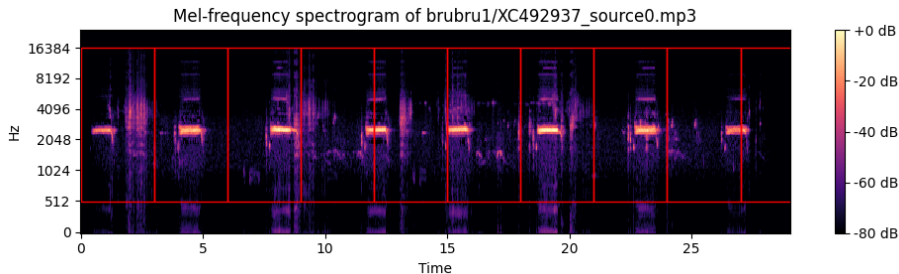Figure 19: MixIT is a sound separation algorithm.

Figure 20: Chunked spectrogram of a bird call.

# Sequence models with embeddings



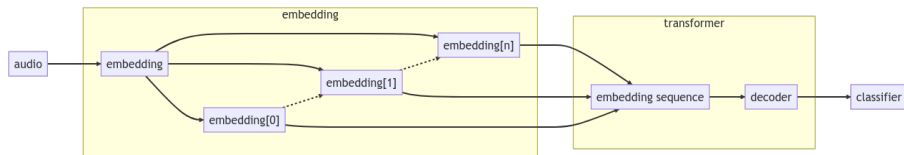Figure 21: We experiment with embeddings in a sequence model (e.g. Transformers) to imbue temporal context.

Figure 22: Bird conservation is a worthy cause and a great opportunity to learn.

## Interested in joining?

Talk to me if you're interesting in jumping in last minute.

# Advice for myself one year ago

## Building a team is worthwhile

- A strong team can help you achieve more than you could on your own. It's also an opportunity to connect with other students.

## Be prepared to learn how to lead a team

- Effective communication and clear timelines are key to keeping the team on track
- Remember that everyone on the team is capable and valuable, and make an effort to recognize and appreciate their contributions

## Reach out to OMSCS and OMSA early

- Working professionals have *a lot* to bring to the table.

Figure 23: Be on the lookout!

There's an abundance of opportunities for OMSCS students to collaborate with other students.

# Thank you to everyone involved

## DS@GT Leadership

- Pulak Agarwal
- Krishi Manek

## BirdCLEF 2022

- Jiangyue Yu
- Bryan Cheungvivatpant
- Dakota Dudley
- Aniketh Swain

## BirdCLEF F22 EDA

- Jinsong Zhen
- Kien Tran
- Siying Liu
- Muskaan Gupta
- Xinjin Li

## BirdCLEF 2023

- Chris Hayduk
- Erin Middlemas
- Grant Williams
- Nathan Zhong

# Links and Resources

- Working Notes, "Motif Mining and Unsupervised Representation Learning for BirdCLEF 2022"
- DS@GT, Kaggle Competition Team Proposal, BirdCLEF 2022
- DS@GT, Project Group Proposal, BirdCLEF EDA Fall 2022
- DS@GT, Kaggle Competition Team Proposal, BirdCLEF 2023
- DS@GT, Assessment, BirdCLEF EDA Fall 2022
- DS@GT, Assessment, BirdCLEF 2023
- BirdCLEF Motif Viewer, Barn Owl, XC138041
- BirdCLEF 2023 MixIT Exploration, Red-chested Cuckoo, 2FXC207767

# Thank you!

## Time for Questions and Answers



Figure 24: Q&A