# Apache Spark & Citizen Science

Using eBird Data to Predict Bird Abundance at Scale

**Tom Auer (mta45@cornell.edu)**, Daniel Fink, and Steve Kelling
Cornell Lab of Ornithology

SPARK SUMMIT 2017

# Background

**The Cornell Lab of Ornithology**

Photo: www.jonreis.com

Our mission: To interpret and conserve the earth's biological diversity through research, education, and citizen science *focused on birds.*
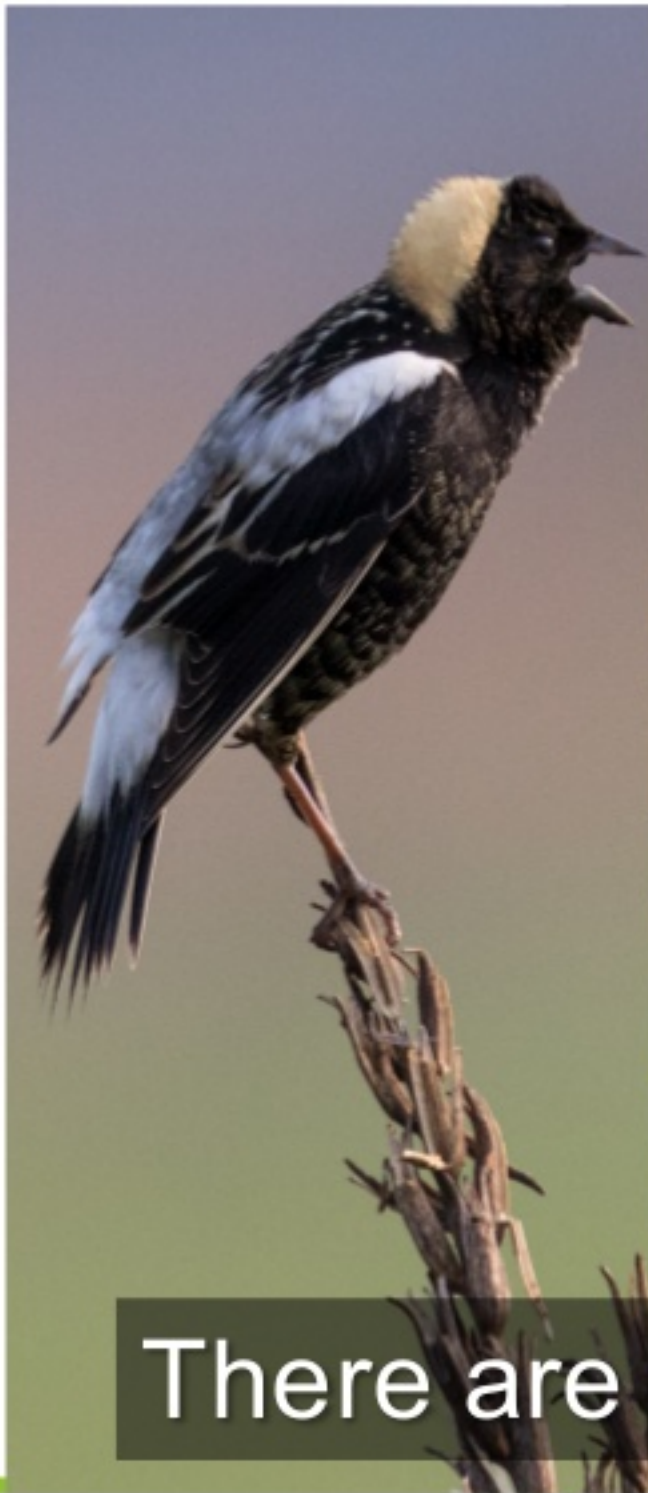
Why Birds?

Why Birds?

There are > 10,000 Species

They are found in all environments

Why Birds?

Why Birds?

They can adapt to novel environments

Indian Vulture

India & Pakistan
30,000,000 in 1990
functionally extinct today

Why Birds?

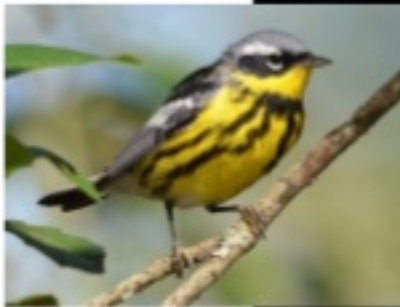They are the most easily observed, counted, and studied of all widespread animal groups

# eBird

eBird

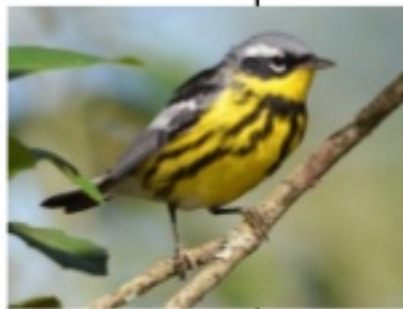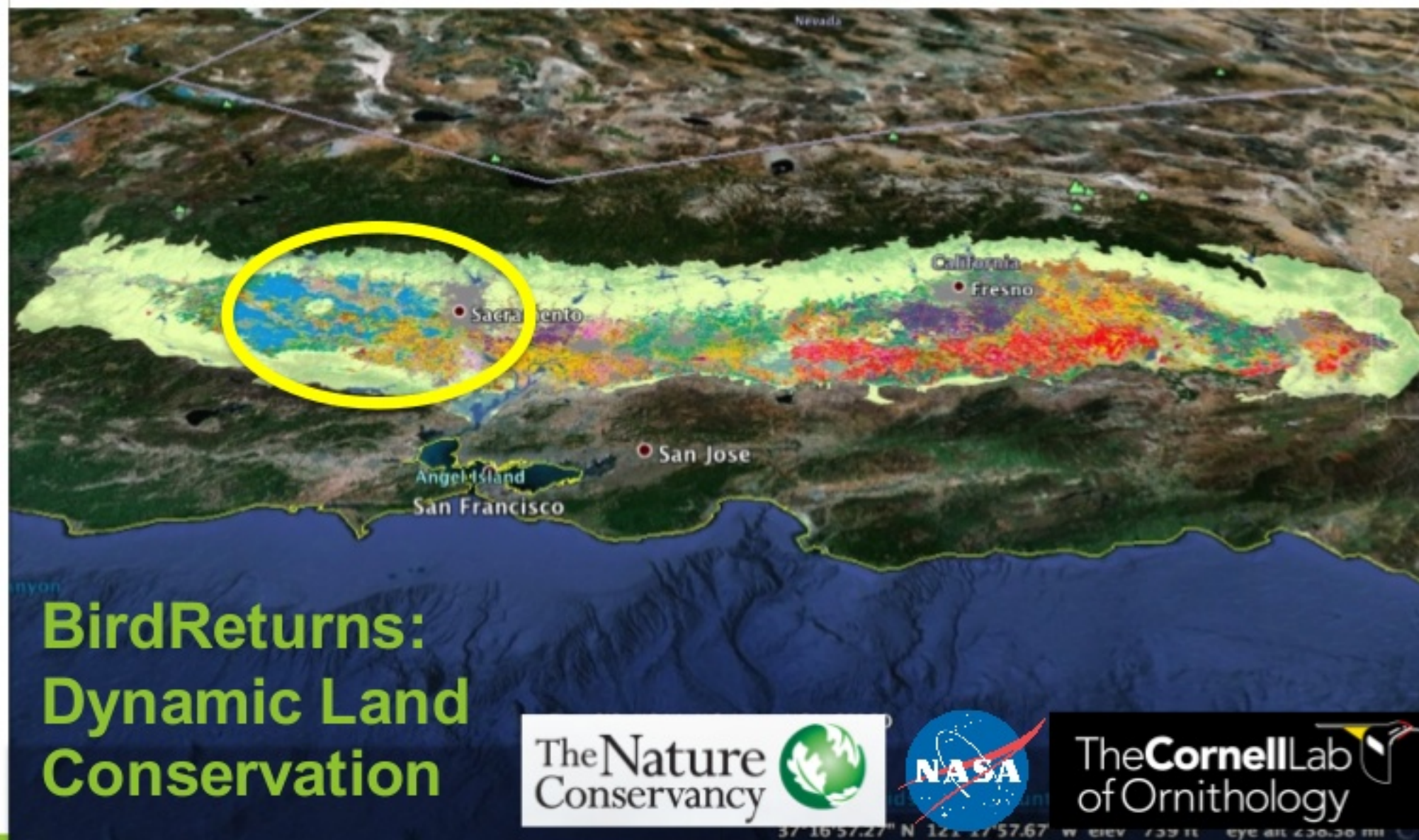**30 million hours** collecting bird observations

Microsoft
Azure

nnot currently be displayed.

BirdReturns:
Dynamic Land
Conservation

The Nature Conservancy · NASA · The Cornell Lab of Ornithology
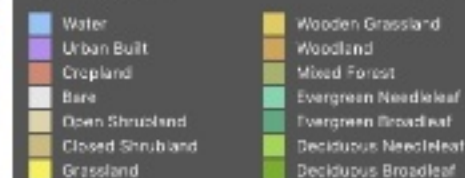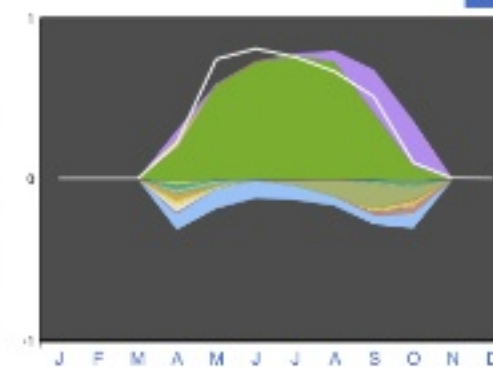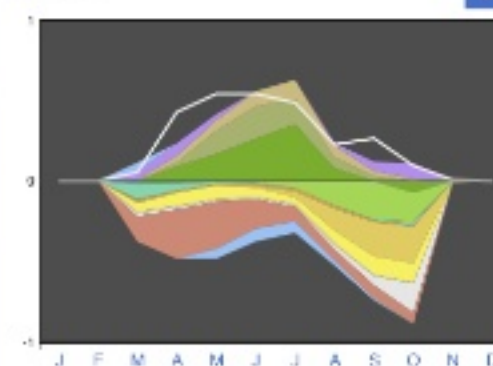
# Technical Experience

June 6
Magnolia Warbler

# Linking Populations and Environment



eBird Observations

Machine Learning

Elevation · Habitat · Search Effort

Scaling: Ecological Challenges

Tree Swallow, March

# SpatioTemporal Exploratory Model (STEM)

## 1. Divide

- Partition extent into regions
- Train and predict models within regions



100 randomized replicates

## 2. Recombine

- Average predictions across all models for each location within regions

Region

Location

**Step 1:**
**Divide into Regions**

Train
Test
Prediction

Randomized
Partitions
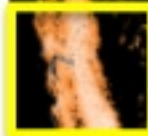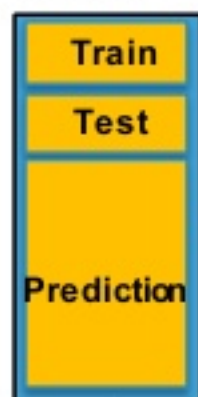
**Step 2:**
**Fit & Predict in each Region**

Region 1

~20,000

**Step 3:**
**Summarize by Location**

**Step 4:**
**Write Out**

.gz

12m rows
2gb

3.1b rows
750gb

1b rows
100gb

52m rows
12gb

# Code

# Current Modeling

- Sampling to address class imbalance
- First stage: binary response GBM
  - Calibrate with GAM
- Second stage: Poisson response GBM

*Models use weights*

# Future Modeling

- "Occupancy" Models
- Semi-parametric learning: GamboostLSS
- Statistical/Machine Learning models: suRFing
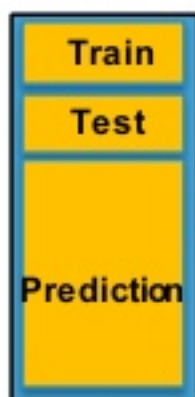
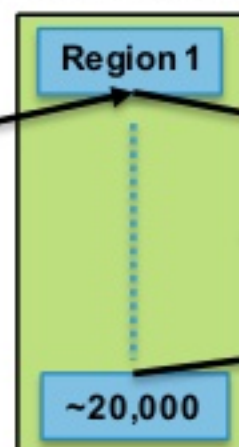# What have we tried?

- HPC Parallelization
- Hadoop MapReduce
- ~~SparkR~~
- Spark 2.x

# RDD pipe()

Stage 0:
Divide into Regions

Stage 1:
Fit & Predict in each Region

Stage 2:
Summarize by Location

Stage 3:
Write Out

Train

Test

Prediction

Randomized
Partitions

Region 1

~20,000

.gz

12m rows
2gb

3.1b rows
750gb

1b rows
100gb

52m rows
12gb

```
.map(lambda p: (p[0].split("\t")[0], p[0].split("\t")[1]+"EOL")) \
.groupByKey().mapValues(list) \
```

# Spark!

- Fast: ~25% faster than MapReduce
- Portable: HPC, Azure
- Scalable: data volume doubled

**Cornell University**
Center for Advanced Computing

**Microsoft Azure**

# What's next?

- RDD pipe()?
- Spark DataFrames
- More Spark!

# Summary

RDD pipe() allows us to keep our code base within our community language and use new R modeling libraries, while leveraging the speed of Spark for parallelizing our modeling workflow to address ecological challenges.

SPARK
SUMMIT
2017