# Enabling Biodiversity Research with Open Source Workflow, GIS and Metadata Tools

CJ Grady,
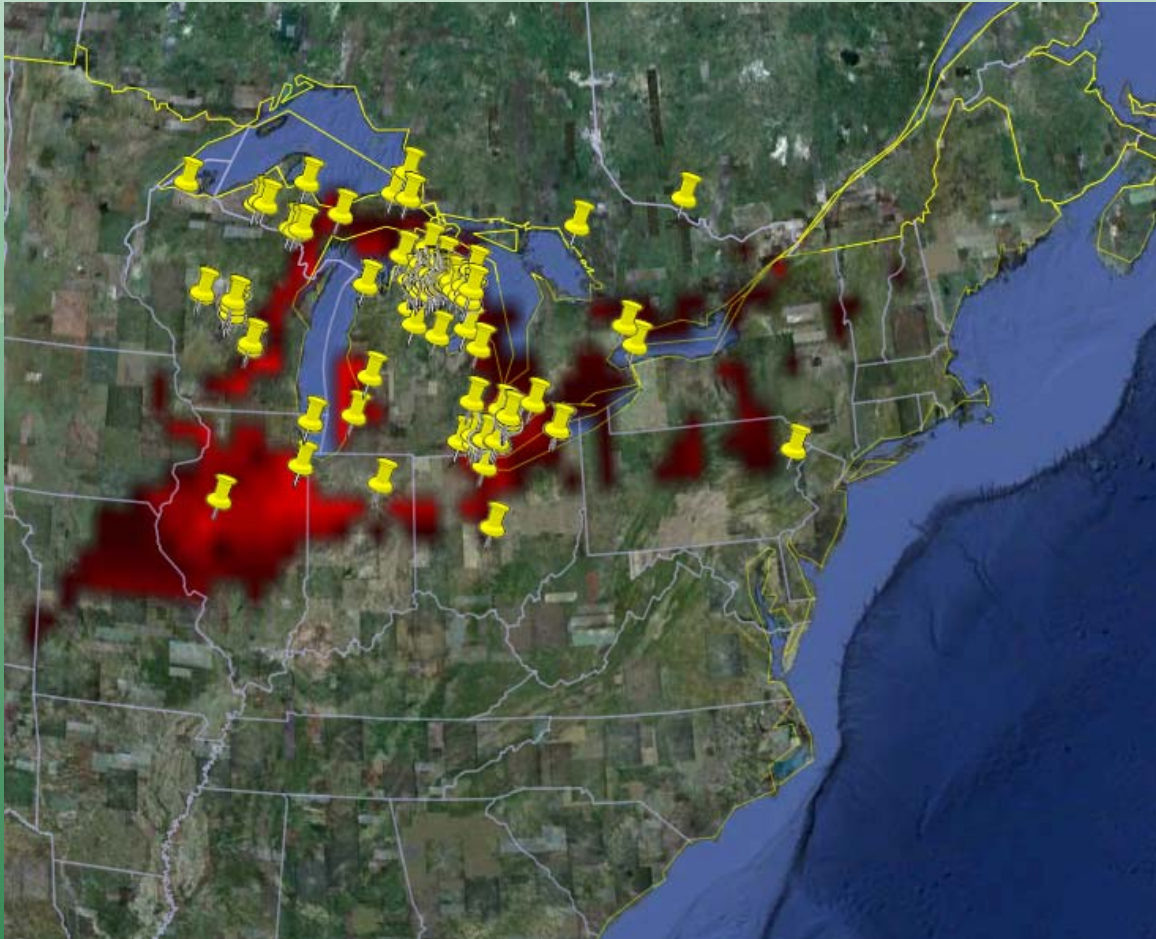
Jim Beach, Aimee Stewart, Jeff Cavner

University of Kansas Biodiversity Institute

http://lifemapper.org

# What is Lifemapper?



Current museum (GBIF) vouchered occurrences for Dendroica kirtlandii

# What is Lifemapper?

**Integration**



**Work Flow Tools**



**Archiving**



EEML

**Repeatable Transparent Science**

Ellison, A. 2010. J. Ecology

# Background

- Workflows
  - A series of connected steps describing a process
- Metadata
  - Information that describes a data set or process.
  - Data about data
- Open Source Software
  - Software that includes access to the source code used to create it. This is provided to encourage study, contribution and sharing.
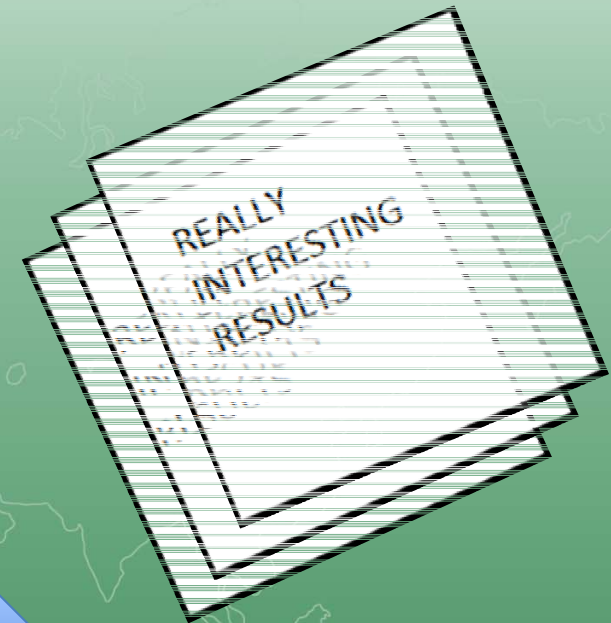
# So What's the Problem?

- Difficult and time consuming to assemble biodiversity experiments by hand

- Scientists often don't have adequate computing resources

- Experiments can be difficult or impossible to reproduce

# Study of Experiment Reproducibility



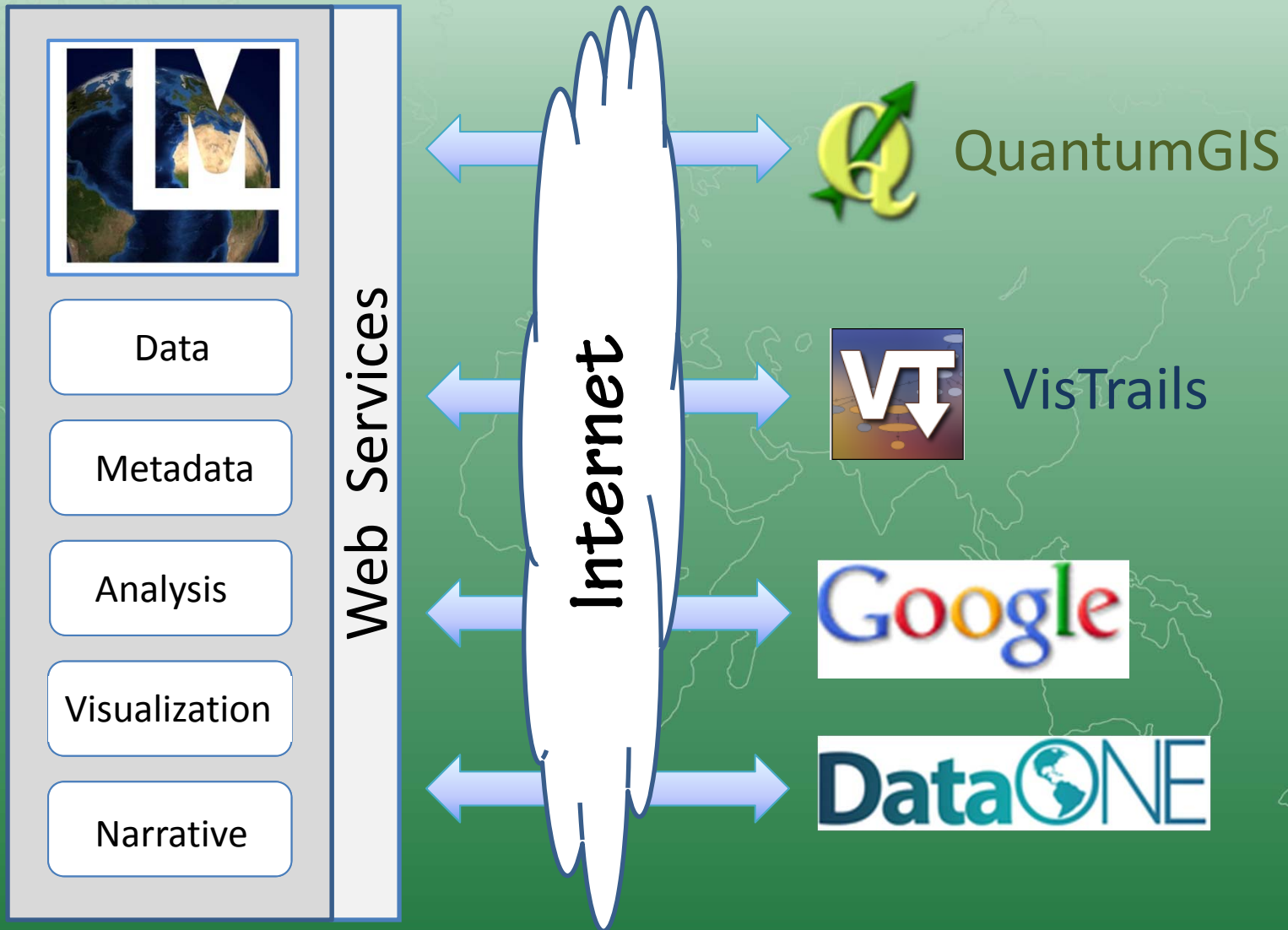Ellison, Aaron. 2010. Repeatability and transparency in ecological research. Ecology 90.

# What we have done

- Metadata for all of our Species Distribution Modeling services

- Simple process metadata
  - Documents how an experiment is ran through our cluster including what versions of software
  - Also describes what web services would be called to execute the experiment again
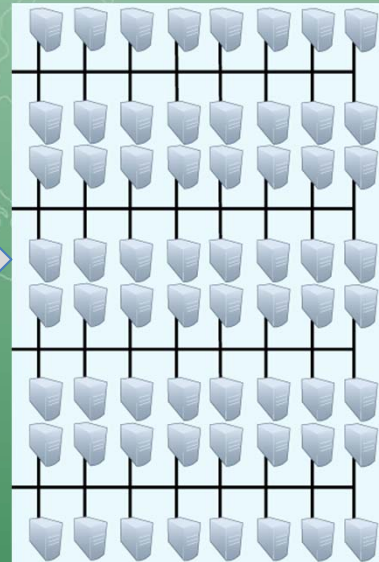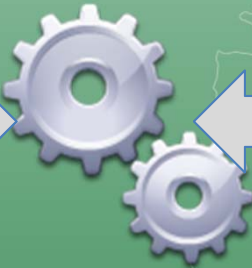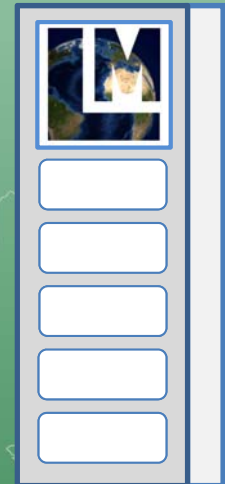
# Web Services

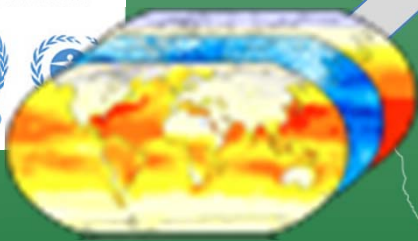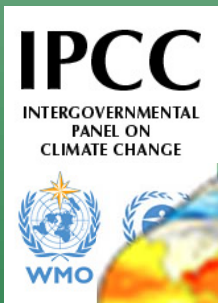# Lifemapper Backend



Pipeline

Data Archive

Compute Resources

# Species Distribution Modeling



Species Occurrence Data

Environmental Data

SDM Modeling Algorithm

Predicted Habitat
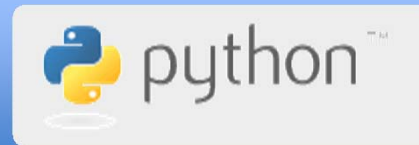
# What we are doing

- Publishing metadata to a public repository
- Client extensions
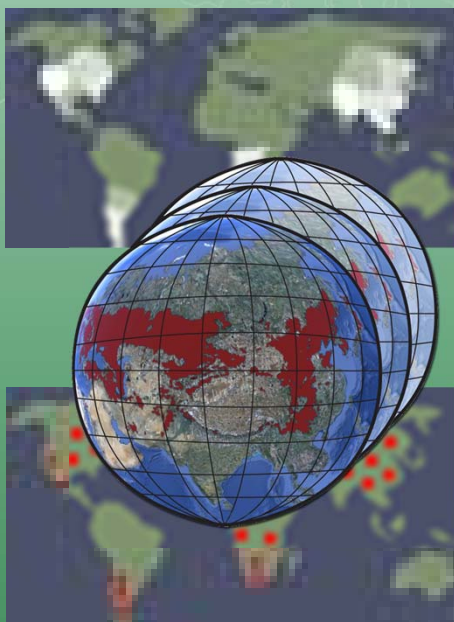- Lifemapper Range and Diversity

# Clients



Lifemapper Web Services
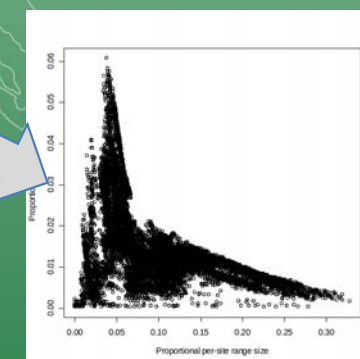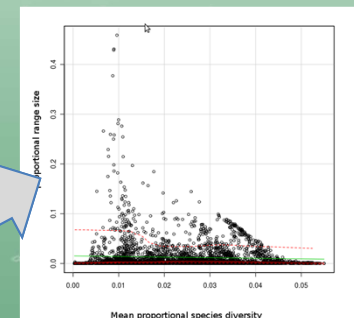
# LmRAD:
# Range and Diversity



Species Habitat Data

Presence Absence Matrix (PAM)

Range and Diversity
Quantifications

# Species Distribution and Diversity

- Biodiversity patterns
  - Species abundance, distribution and diversity
  - Multiple scales and extents
  - Used for conservation and management decisions
- Challenges
  - Large extents (> 10,000 km2)
  - Fine resolution (< 1000 m2 ≈ 30m x 30m)
  - Many species (10,000 +)

# QGIS with
# Lifemapper MacroEcology plug-in

# Hexagonal Grid

# Mouse 'Presence'

# Sites / Species Plot

# VisTrails

# Reproducibility

- Simple process metadata
- Process metadata extensions
- Lifemapper client metadata reader

# Collaborations

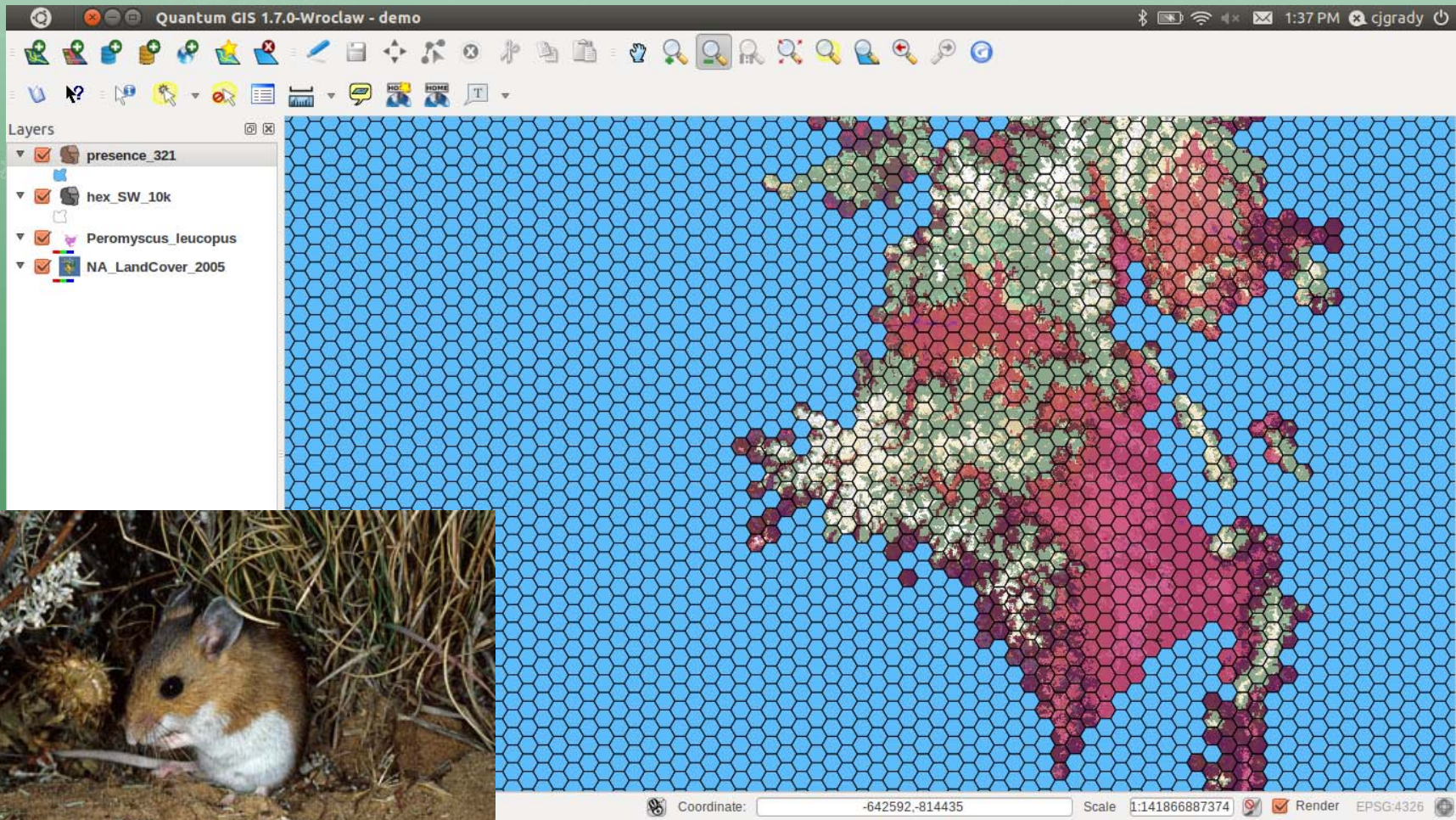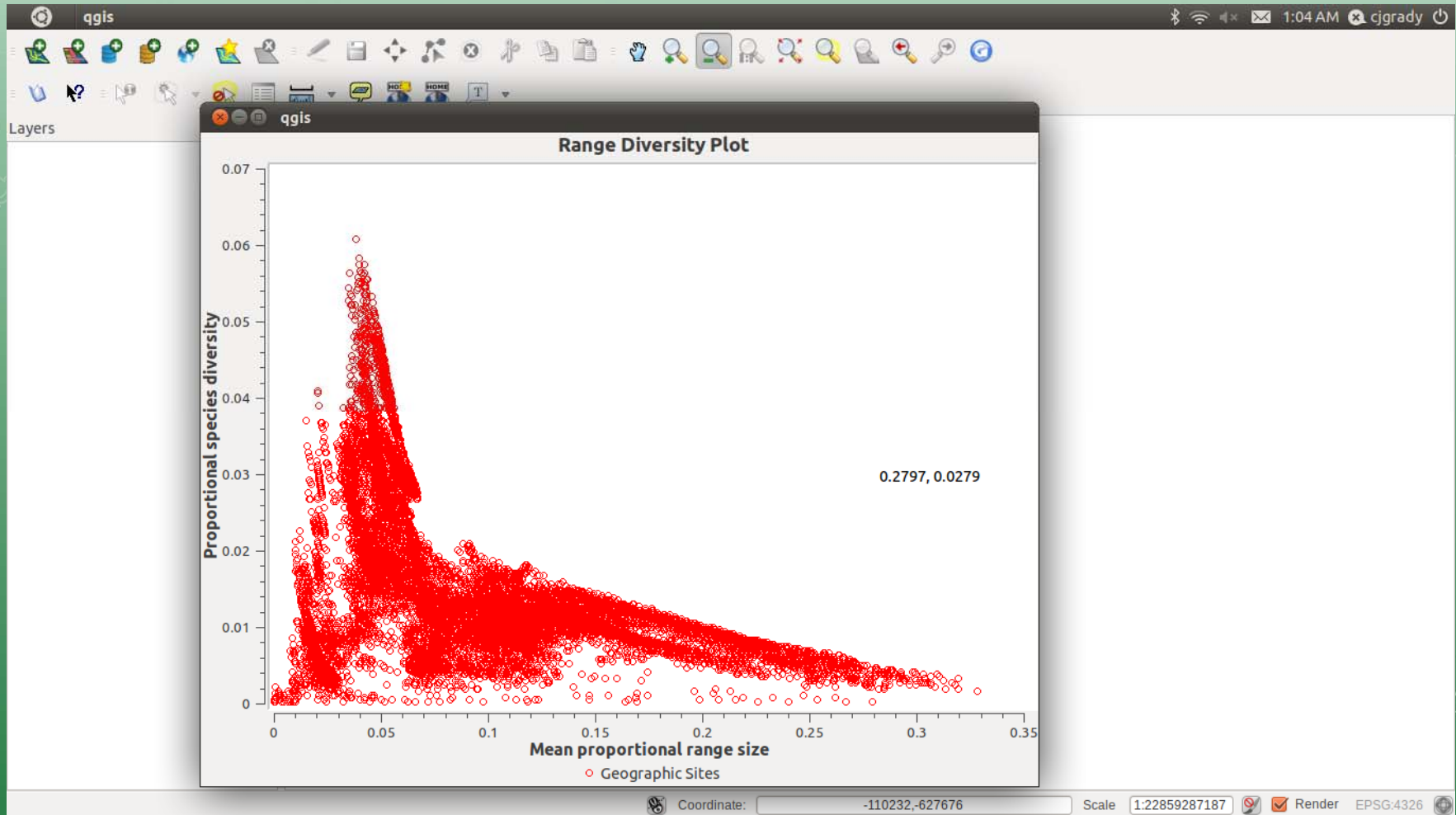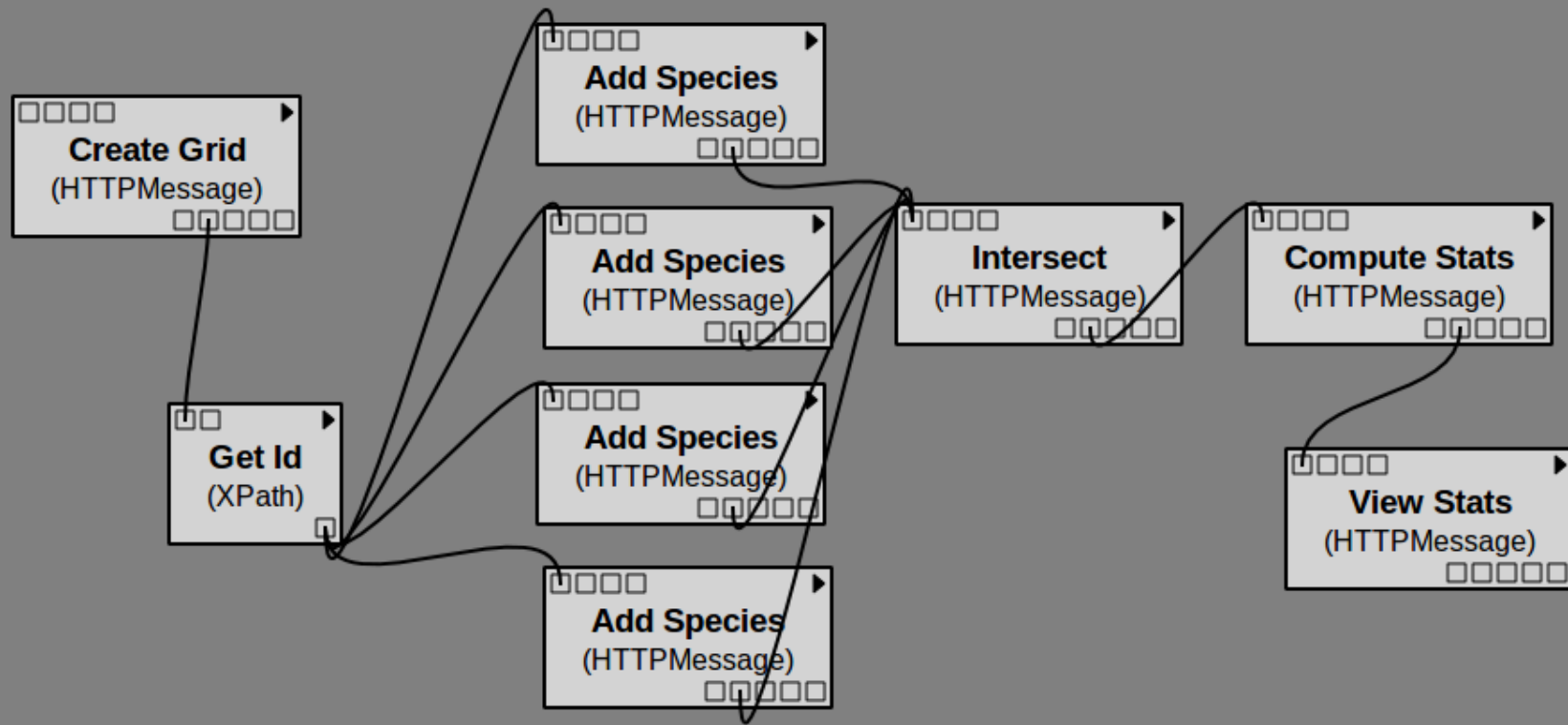- KU Biodiversity Institute

- NSF
  Cyber-Commons

- Change Thinking

- CI Team

http://www.youtube.com/watch?v=VCFixtqlimg

# Future Directions

- Publish metadata through standard APIs
- Contribute process metadata extensions to community
- Gesture based interface
- Explore extensions into cloud and other grid computing environments

# Parallel Processing

# Summary

- Provide end-users with clients to assemble and manage biodiversity modeling experiments

- Allow users to harness the computing power available through our cluster to perform computationally intensive tasks

- Include process metadata to document how an experiment was performed
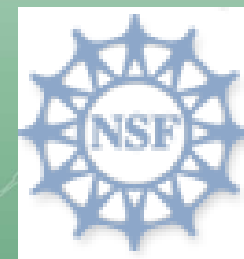
# Funding

U.S. National Science Foundation

NSF EPSCoR 0553722

NSF EPSCoR 0919443

EHR/DRL 0918590

BIO/DBI 0851290

OCI/CI-TEAM 0753336

# Questions

- [cjgrady@ku.edu](mailto:cjgrady@ku.edu)
- http://lifemapper.org