# Genome project management resources at the National Agricultural Library

Chris Childers and Monica Poelchau

USDA-ARS, National Agricultural Library

USDA

# So you have a genome project. Where will you store your data?

- Make your data available through NCBI when applicable (or other INSDC organizations).

- To make your data even more useful for your community, consider also making it available in a taxon-specific repository.

- Advantages for you:
  - Greater visibility for your dataset
  - Value-added tools for searching and browsing, analysis
  - Curation tools to improve annotation quality
  - Help with data management
  - Increasing mandate from journals and funding bodies to make research data fully accessible post-publication[1, 2]

    [1]http://www.nature.com/authors/policies/data/data-availability-statements-data-citations.pdf
    [2]https://obamawhitehouse.archives.gov/the-press-office/2013/05/09/executive-order-making-open-and-machine-readable-new-default-government-

USDA

# So you have a genome project. Where will you store your data?

- Advantages for the scientific community:
  - Helps facilitate knowledge discovery for humans (and sometimes machines);
  - Easier to find data for comparative analyses;
  - Promotes reproducible research;
  - General repositories (e.g. GenBank) may not meet the needs for storing all data types, in particular for non-standard data types (e.g. phenotypic data).

# Genome data management resources for arthropods – how to choose

- What species is the data from?
  - Many taxon-specific genome databases are here at this workshop
- What kind data do you have?
  - Raw data, genome assemblies, transcriptome assemblies, gene annotations, can and should all be stored at NCBI (or other INSDC organization)
  - Some or all of these data types can also be made accessible at genome databases (just ask)
  - Generic repositories (e.g. Dryad, Ag Data Commons) can be used for data types that don't fit the mold

USDA

# The i5k Workspace@NAL

- We support any 'orphaned' arthropod genome project.
  - Connect researchers to the data
  - Create standardized tools for accessing the data in useful ways
  - Provide resources to facilitate manual curation projects
- Supported data types:
  - Genome assembly
  - Anything that you can map to or predict from the genome assembly
- Main requirements:
  - Genome assembly needs to be in GenBank/ENA/DDBJ
  - Data should be public (no private repositories)
  - Manual annotation only occurs at one genome database at a time

- Research plan
- Genome sequencing
- Genome assembly
- Automated annotation of genome assembly

- Manual Curation
- Official gene set (OGS) generation

- Biological insights/Publication

- Data access for the broader community
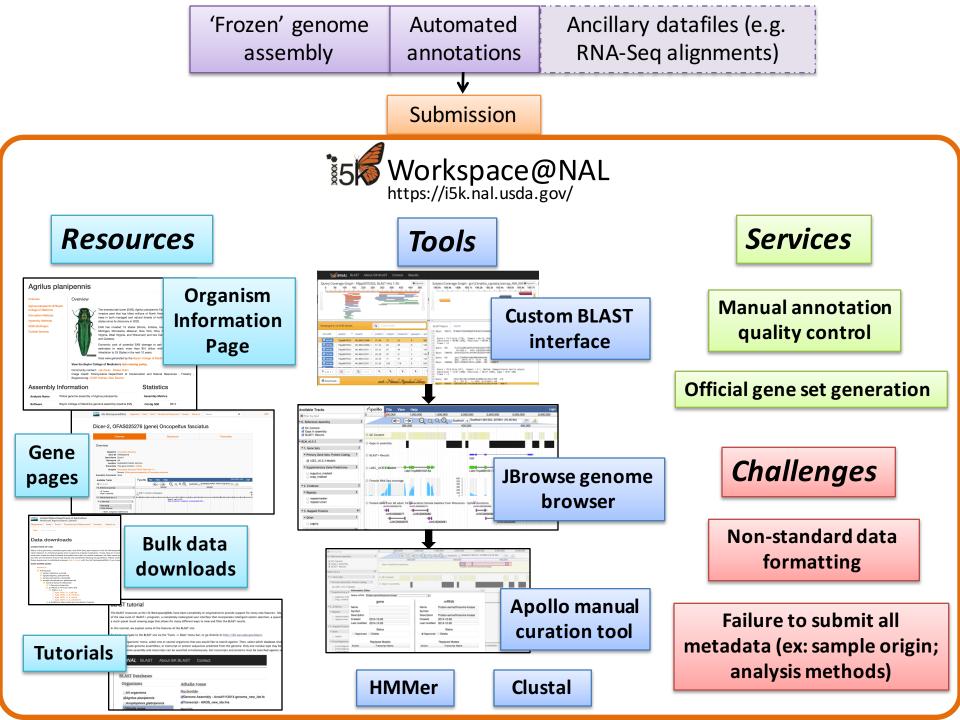- Genome project maintenance

Genome Project Trajectory

USDA

# The i5k Workspace@NAL

Our background:

- Originally set up to support genomes sequenced as part of the i5k initiative

- I5k: International effort to prioritize insect genomes for sequencing; provide guidelines for genome sequencing and curation; and seek funding

- I5k Goal: coordinate the sequencing and assembly of 5000 insect or related arthropod genomes

- Brief introduction to i5k at the beginning of the i5k session on Thursday

USDA

'Frozen' genome assembly | Automated annotations | Ancillary datafiles (e.g. RNA-Seq alignments)

Submission

i5k Workspace@NAL
https://i5k.nal.usda.gov/

**Resources**

Organism Information Page

Gene pages

Bulk data downloads

Tutorials

**Tools**

Custom BLAST interface

JBrowse genome browser

Apollo manual curation tool

HMMer | Clustal

**Services**

Manual annotation quality control

Official gene set generation

**Challenges**

Non-standard data formatting

Failure to submit all metadata (ex: sample origin; analysis methods)

# i5k Workspace content –
# 57 species and counting

| Order | Quantity | Order | Quantity |
|---|---|---|---|
| Amphipoda | 1 | Hemiptera | 7 |
| Araneae | 3 | Hymenoptera | 14 |
| Blattodea | 1 | Lepidoptera | 2 |
| Calanoida | 1 | Odonata | 1 |
| Coleoptera | 7 | Orthoptera | 1 |
| Diplura | 1 | Scorpiones | 1 |
| Diptera | 13 | Thysanoptera | 1 |
| Ephemeroptera | 1 | Trichoptera | 1 |
| Harpacticoida | 1 | | |

- Many other datasets mapped to, or predicted from each genome assembly (gene predictions, transcriptomes, RNA-Seq, etc.)

USDA

# Community annotation at the i5k Workspace

- What is community annotation?
  - Scientists collectively examine and improve gene models (usually computationally predicted)
- Why annotate?
  - Verify quality of automated gene predictions
  - Improve gene models for specific analyses
  - Link gene models to existing literature and ontologies
- Our community: Over 400 registered annotators have curated over 10,000 gene models using the Apollo software
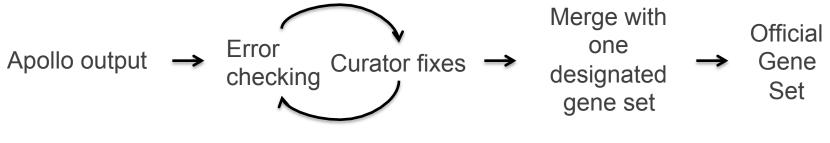
USDA

# Community annotation at the i5k Workspace

Our support for community annotation includes:

- Access to a large community of curators

- Tutorials, guidelines, webinars

- Registration mechanism for new annotators

- One-on-one support

- Software to evaluate changes between curated and original annotations (Chien-Yueh Lee, https://github.com/chienyuehlee/gff-cmp-cat)
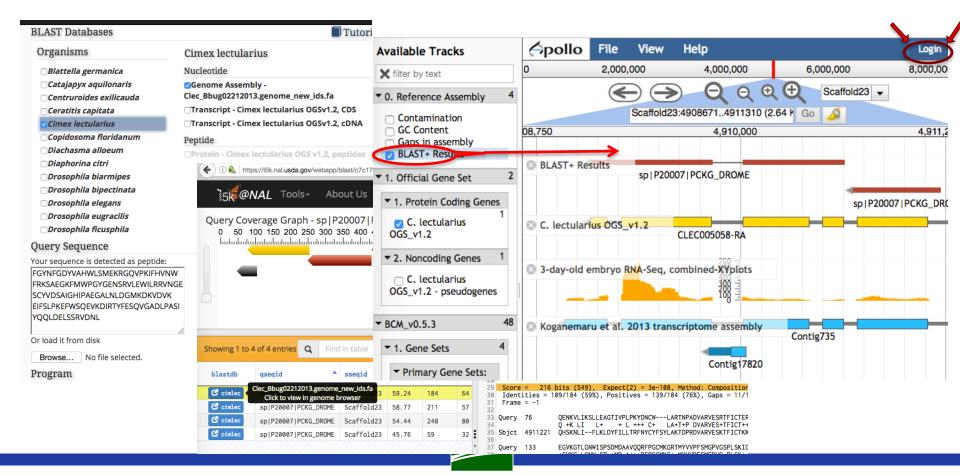
USDA

# QC and OGS pipeline

- QC program corrects common formatting errors from the curation process

- OGS generation program merges curated models with one designated gene set using curator-supplied information

- Still in development, already 6 OGS's produced (Mei-Ju Chen)

Apollo output → Error checking → Curator fixes → Merge with one designated gene set → Official Gene Set

USDA

# Genome already hosted elsewhere?

- You can also use our tools to query the datasets that we host.

# Other resources at the NAL:
# The Ag Data Commons

- Hosts any dataset funded by the USDA

- Landing page

- Citable DOI

- https://data.nal.usda.gov/

- Nine i5k datasets already available

# What we'll talk about tomorrow

1. Background: What is the i5k Workspace?
2. Submitting data
3. Finding data at the i5k Workspace
    1. General search/Content types
    2. Data downloads
    3. BLAST
    4. Clustal(s)
    5. HMMER
    6. Jbrowse/Apollo
4. Improving data at the i5k Workspace via community annotation
    1. See Monica Munoz-Torres' workshop for full use of Apollo

USDA

# Need more information?

i5k Workspace@NAL:

- https://i5k.nal.usda.gov/

- https://github.com/NAL-i5K/

- Poster during the Friday session

The i5k initiative:

- New website: http://i5k.github.io/
  Ag Data Commons:

- https://data.nal.usda.gov/

USDA

# Acknowledgements

The NAL Team

- Gary Moore
- Susan McCarthy
- Yu-yu Lin
- Mei-Ju Chen

Workspace alumni

- Chien-Yueh Lee
- Han Lin
- Jun-Wei Lin
- Vijaya Tsavatapalli

i5k Workspace@NAL advisory committee

- Jay Evans
- Kevin Hackett
- Simon Liu
- Ursula Pieper

- i5k Coordinating Committee
- i5k Pilot Project
- Apollo & JBrowse Development Teams
- GMOD/Tripal community
- The AgBioData consortium
- All of our users and contributors!