# I5k Workspace webinar Changes to gene and protein data access

Monica Poelchau

7/19/2022

# Agenda

- Why did we remove the i5k Workspace gene and protein pages?
  - What exactly did we remove?
  - Why?
  - A brief digression on how we keep i5k Workspace data in sync with NCBI
- The data are still accessible – how can you find what you need?
  - Gene and protein information
  - Dataset-level information
- Q&A to address any questions or concerns.
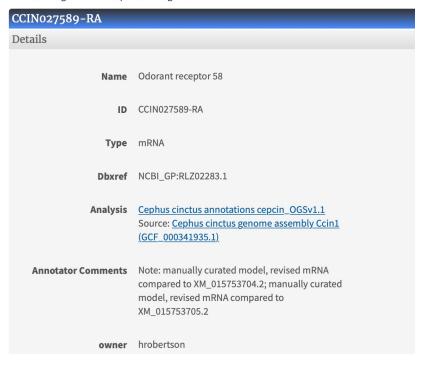
# Gene and protein page removal

- ## What was removed?

  - Pages that showed information at the individual gene and protein level

  - Ability to search for gene and protein information (e.g. "Cimex lectularius heat shock protein")

**CCIN027589, CCIN027589 (gene) Cephus cinctus**

Overview
Sequences
Transcripts

**Transcripts**

The following features are part of this gene:

**CCIN027589-RA**

Details

| | |
|---|---|
| **Name** | Odorant receptor 58 |
| **ID** | CCIN027589-RA |
| **Type** | mRNA |
| **Dbxref** | NCBI_GP:RLZ02283.1 |
| **Analysis** | Cephus cinctus annotations cepcin_OGSv1.1 Source: Cephus cinctus genome assembly Ccin1 (GCF_000341935.1) |
| **Annotator Comments** | Note: manually curated model, revised mRNA compared to XM_015753704.2; manually curated model, revised mRNA compared to XM_015753705.2 |
| **owner** | hrobertson |

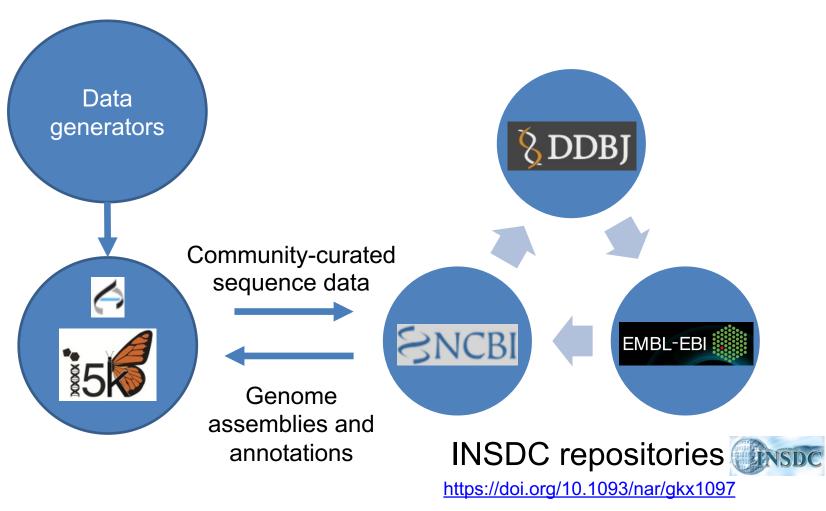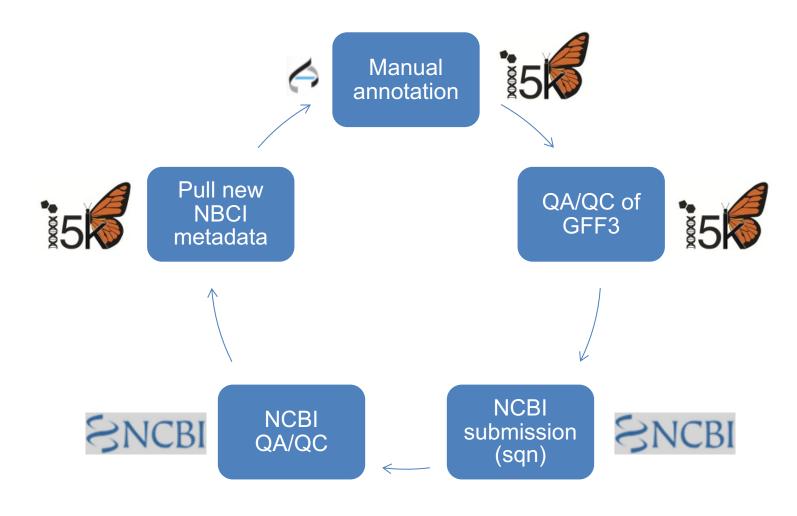# Why did we remove the pages?

1. The data are already available elsewhere at the gene/protein level: :
   1. NCBI
   2. Other insect databases, e.g. HymenopteraMine, VEuPathDB
2. The data are still available at the dataset level:
   1. I5k Workspace@NAL
   2. Ag Data Commons
   3. NCBI
   4. Other insect databases, e.g. HymenopteraMine, VEuPathDB
3. Not a high priority for users
   1. Structured interview results
   2. Fewer page views
4. Challenging to maintain the code
5. Challenging to keep the content up-to-date and comprehensive

# Why are the data at NCBI?

- NCBI, as part of the INSDC, is the primary archive for all sequence data
- Our policy is to submit all sequence data and metadata that can go to NCBI, to NCBI

Data generators

Community-curated sequence data

Genome assemblies and annotations

DDBJ

NCBI

EMBL-EBI

INSDC repositories

https://doi.org/10.1093/nar/gkx1097

# Manual annotation QA/QC and submission

# Result of a successful GenBank submission

# How can I find a specific gene?

**NCBI Protein search**

**NCBI Gene search**

https://www.ncbi.nlm.nih.gov/protein/

https://www.ncbi.nlm.nih.gov/gene/

# How can I find a specific gene?

## Vectors

## Hymenoptera



https://vectorbase.org/vectorbase/app

https://hymenopteramine.rnet.missouri.edu/hymenopteramine/begin.do
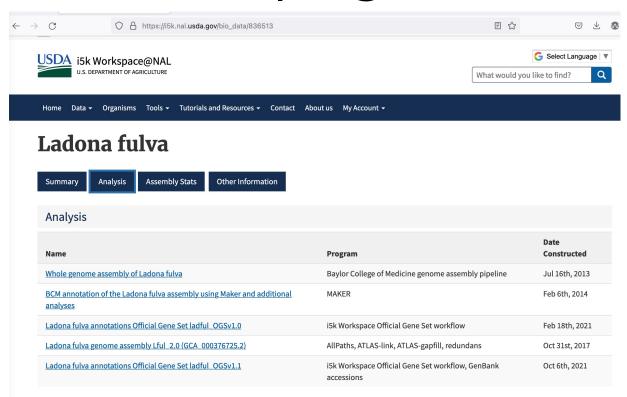
# How can I find a gene set?

**I5k Workspace@NAL search**

**I5k Workspace@NAL browse**



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

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

# How can I find a gene set?

**I5k Workspace@NAL browse**          **I5k Workspace@NAL browse**

# How can I find a gene set?

## NCBI search



https://www.ncbi.nlm.nih.gov/genome/

## NCBI datasets



https://www.ncbi.nlm.nih.gov/datasets/
FYI, don't use Safari

# How can I find a gene set?

## Ag Data Commons search

## Ag Data Commons browse

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

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

# Questions?

# Thank you!

- Contact us:
  - i5k@usda.gov
  - https://i5k.nal.usda.gov/contact

# List of datasets with removed gene pages

- [https://i5k.nal.usda.gov/news/i5k-workspace-removing-all-gene-and-mrna-pages](https://i5k.nal.usda.gov/news/i5k-workspace-removing-all-gene-and-mrna-pages)