# Using genome browsers constructed by G-OnRamp to provide students with a Course-based Undergraduate Research Experience in genome annotation

## Wilson Leung
## Washington University in St. Louis

http://g-onramp.org

01/2020

# Creating genome browsers requires *substantial bioinformatics expertise*

**Manage multiple data formats**

```
$ faToTwoBit newGenome.fa newGenome.2bit
```
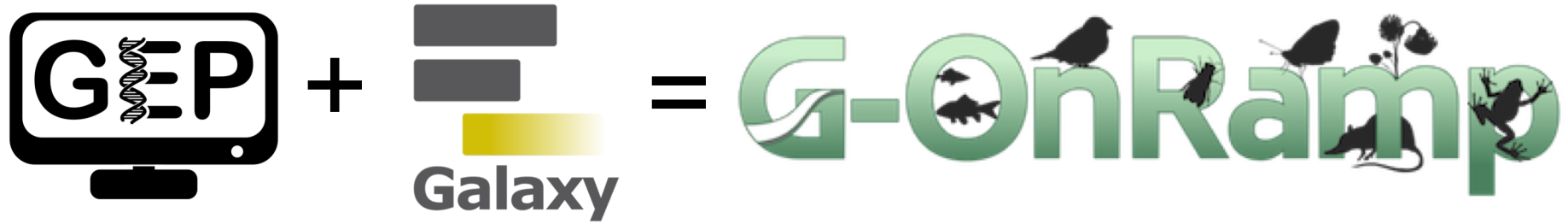
FASTA    twoBit

**Manage input and output of multiple tools**

```
$ twoBitInfo newGenome.2bit stdout | \
    sort -k 2,2nr > newGenome.chrom.sizes
```

**Set up genome database**

```
INSERT INTO defaultDb (genome, name)
VALUES ("Ganaspis species 1", "ganSpe1");
```
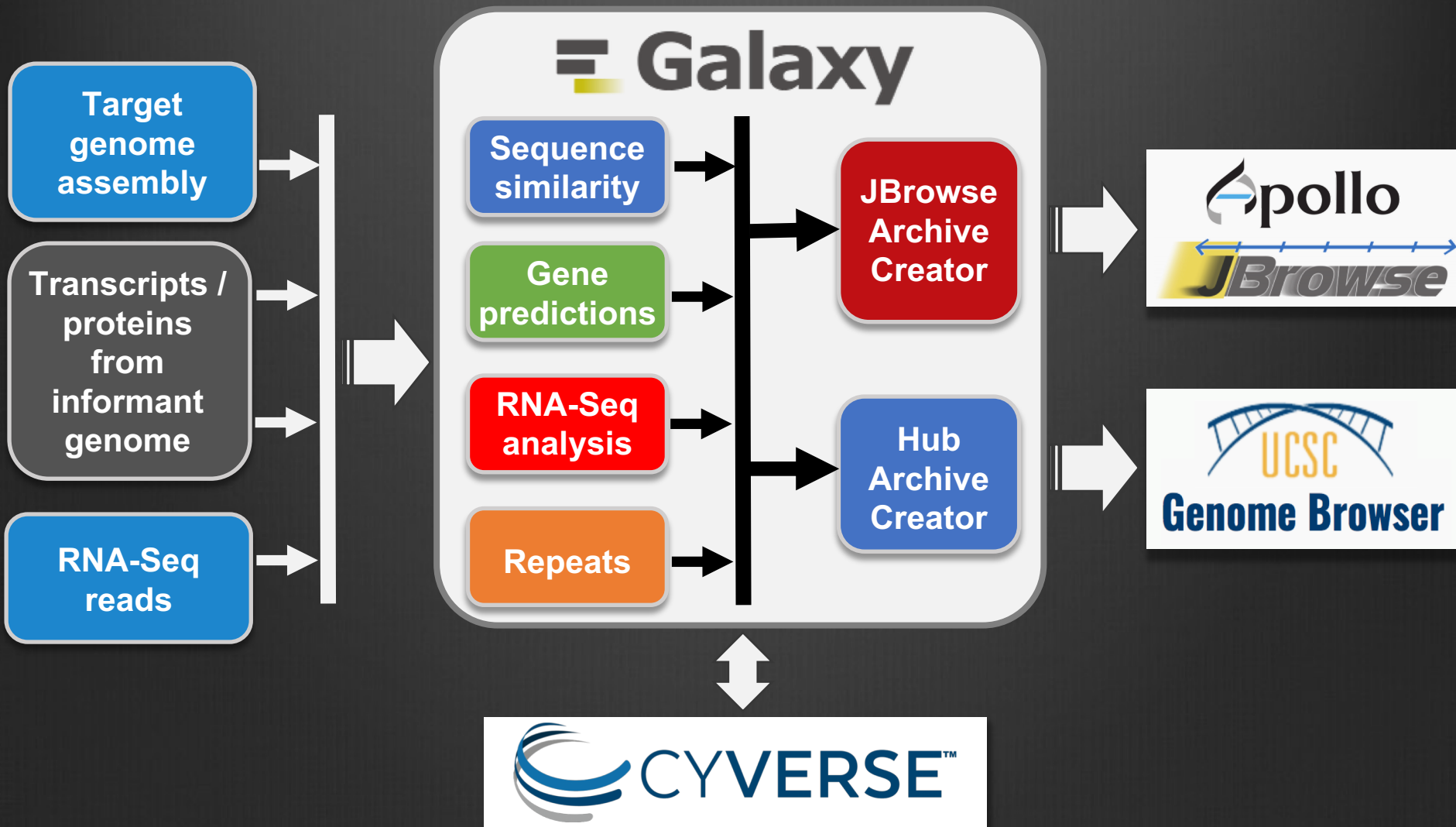
**Step 0: Download and install bioinformatics tools**

## Goals:

- Enable researchers and biology faculty to create **genome browsers** for their favorite **eukaryotic species**

- Enable faculty to use the genome browsers created by G-OnRamp to **engage students in genomics research**

# Use G-OnRamp to create genome browsers for **eukaryotic genomes**

# It is **easy to get started** with G-OnRamp

**Workflow: G-OnRamp production workflow for UCSC**

**Run Workflow**

**History Options**

**Send results to a new history**

Yes | No

📄 1: Target genome

📄  📑   1: Gan_sp1-scaffolds.fa

**Target genome assembly**

📄 2: Informant mRNA GenBank records

📄  📑   5: GCF_000001215.4_Release_6_plus_ISO1_MT_rna.gbff

**Transcripts / proteins from informant genome**

📄 3: Informant protein sequences

📄  📑   4: GCF_000001215.4_Release_6_plus_ISO1_MT_protein.faa

📄 4: RNA-Seq: Forward reads

📄  📑   2: SRR805628_1.fastq

**RNA-Seq reads**

📄 5: RNA-Seq: Reverse reads

📄  📑   3: SRR805628_2.fastq

# UCSC Assembly Hub for *Ganaspis species 1*

# G-OnRamp is composed of **sub-workflows**

## Sub-workflows

**Sequence similarity**
- NCBI BLAST+
- UCSC BLAT

**Gene predictions**
- Augustus
- GlimmerHMM
- SNAP

**RNA-Seq analysis**
- HISAT2
- regtools
- StringTie

**Repeats identification**
- TRF
- WindowMasker

**Create Assembly Hubs**
- Hub Archive Creator
- JBrowse Archive Creator

**Apollo interactions**
- Delete an Apollo Record
- Apollo User Manager
- Create or Update Organism

# Galaxy PROJECT
## an open, web-based platform for bioinformatics analyses

## ⚙ Accessible
- ⚙ Does not require programming experience

## ⚙ Reproducible
- ⚙ Easily repeat analyses that contain multiple steps

## ⚙ Transparent
- ⚙ Share and publish workflows and results

**Publications each year**



**Galaxy Help: New Topics and Responses per month**
● New Topics   ● New Responses



Month

# Galaxy automatically keeps track of **each step of the analysis** (History)



- Keep track of the metadata associated with each Dataset

- Examine details of each step of the analysis

- Easily repeat each step of the analysis

# Bioinformatics tools often have
## different user interfaces

```
% windowmasker –help
USAGE
  windowmasker [–h] [–help] [–xmlhelp] [–ustat unit_counts]
    [–in input_file_name] [–out output_file_name] [–checkdup check_duplicates]
    [–fa_list input_is_a_list] [–mem available_memory] [–meta info_string]
    [–unit unit_length] [–genome_size genome_size] [–window window_size]
    [–t_extend T_extend] [–t_thres T_threshold] [–set_t_high score_value]
    [–set_t_low score_value] [–parse_seqids] [–outfmt output_format]
    [–t_high T_high] [–t_low T_low] [–infmt input_format]
    [–exclude_ids exclude_id_list] [–ids id_list] [–text_match text_match_ids]
    [–sformat unit_counts_format] [–smem available_memory] [–dust use_dust]
    [–dust_level dust_level] [–mk_counts] [–convert] [–version–full]
```

**Specify input file with -in flag**

```
faToTwoBit – Convert DNA from fasta to 2bit format
usage:
    faToTwoBit in.fa [in2.fa in3.fa ...] out.2bit
```

**Specify one or more input files without a flag**

# Galaxy provides a **standardized interface** for specifying inputs, parameters, and outputs

# Create a sub-workflow by adding tools and specifying the **input and output datasets**

Combine sub-workflows to create the large workflow

# G-OnRamp training materials
(http://g-onramp.org/training)


http://g-onramp.org/genome-browsers

- 6 G-OnRamp workshops from 2016–2018:
  - 65 participants
  - 40+ institutions
  - Half are from Primarily Undergraduate Institutions (PUIs)

- Created genome browsers for 18 species
  - Available through the CyVerse Data Store

# Genomics Education Partnership (http://gep.wustl.edu)



- Integration of **genomics** and research thinking into the **undergraduate biology curriculum**

- Creation of **student-scientist partnerships**

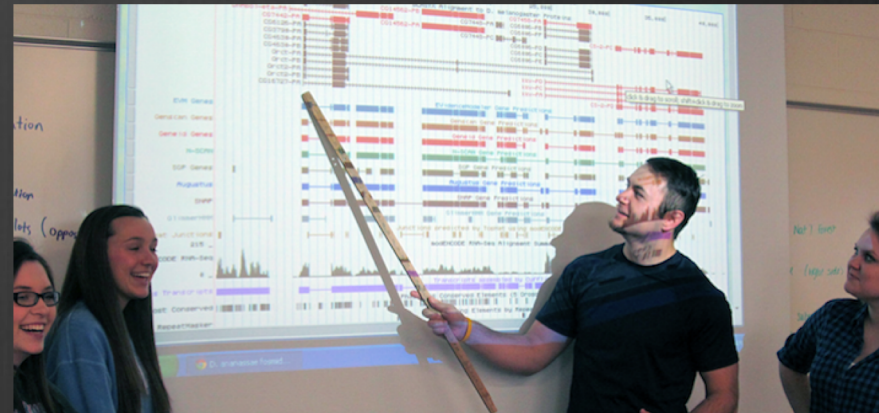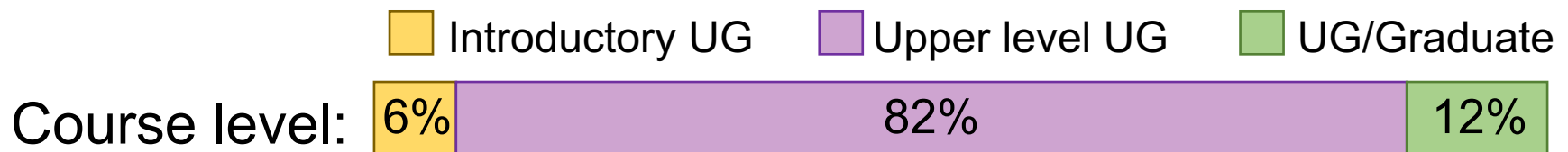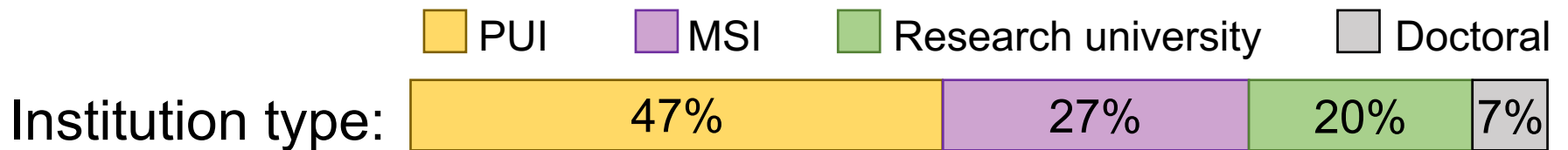- **Publication** of research in genomics and in science education

# Comparative annotation of **four parasitoid wasp species**



- Goal: understand how **venom proteins** from parasitoid wasps manipulate the signal transduction pathways and second messenger system of their hosts
  - Dr. Nathan T. Mortimer (Illinois State University)

- Engaged >200 GEP students from 15 institutions:

| | PUI | MSI | Research university | Doctoral |
|---|---|---|---|---|
| Institution type: | 47% | 27% | 20% | 7% |

| | Introductory UG | Upper level UG | UG/Graduate |
|---|---|---|---|
| Course level: | 6% | 82% | 12% |

# The Genome Browsers produced by G-OnRamp **work well in the classroom**

## Mean annotation post-test scores



## Responses to SURE survey questions



Ability to analyze data

Understanding science

Wasp projects (N = 181–195)

Other GEP projects (N = 1200–1270)

# G-OnRamp deployment options (http://g-onramp.org/deployments)

- G-OnRamp virtual appliance
  - Suitable for local testing and training
  - Freely available

- G-OnRamp on Amazon Web Services (AWS)
  - Production analysis of whole genome assemblies
  - CloudLaunch: https://launch.usegalaxy.org

Get started with G-OnRamp

| G-OnRamp Ubuntu Virtual Machine Image ❯ | CloudLaunch Deployment ❯ | G-OnRamp Training Materials ❯ |

# Summary

- G-OnRamp provides a web-based platform for creating genome browsers for eukaryotic genomes

- Faculty have successfully used the genome browsers created by G-OnRamp to engage students in genomics CUREs

**GEP is recruiting Science Partners**
**http://gep.wustl.edu/contact_us**

**Visit the G-OnRamp poster PO0137**

# Related posters and sessions

**Genomics Education Partnership**

**Poster PE0138**

**Posters PE0141 and PE0142**

**Galaxy Session:**
**1/14 @4:00pm (California)**

# Acknowledgements



Washington University in St. Louis

Sarah C. R. Elgin

Yating Liu

Oregon Health & Science University

Jeremy Goecks

Luke Sargent

Illinois State University
*Illinois' first public university*

Nathan T. Mortimer

Grinnell College

David Lopatto

# Questions?



http://g-onramp.org