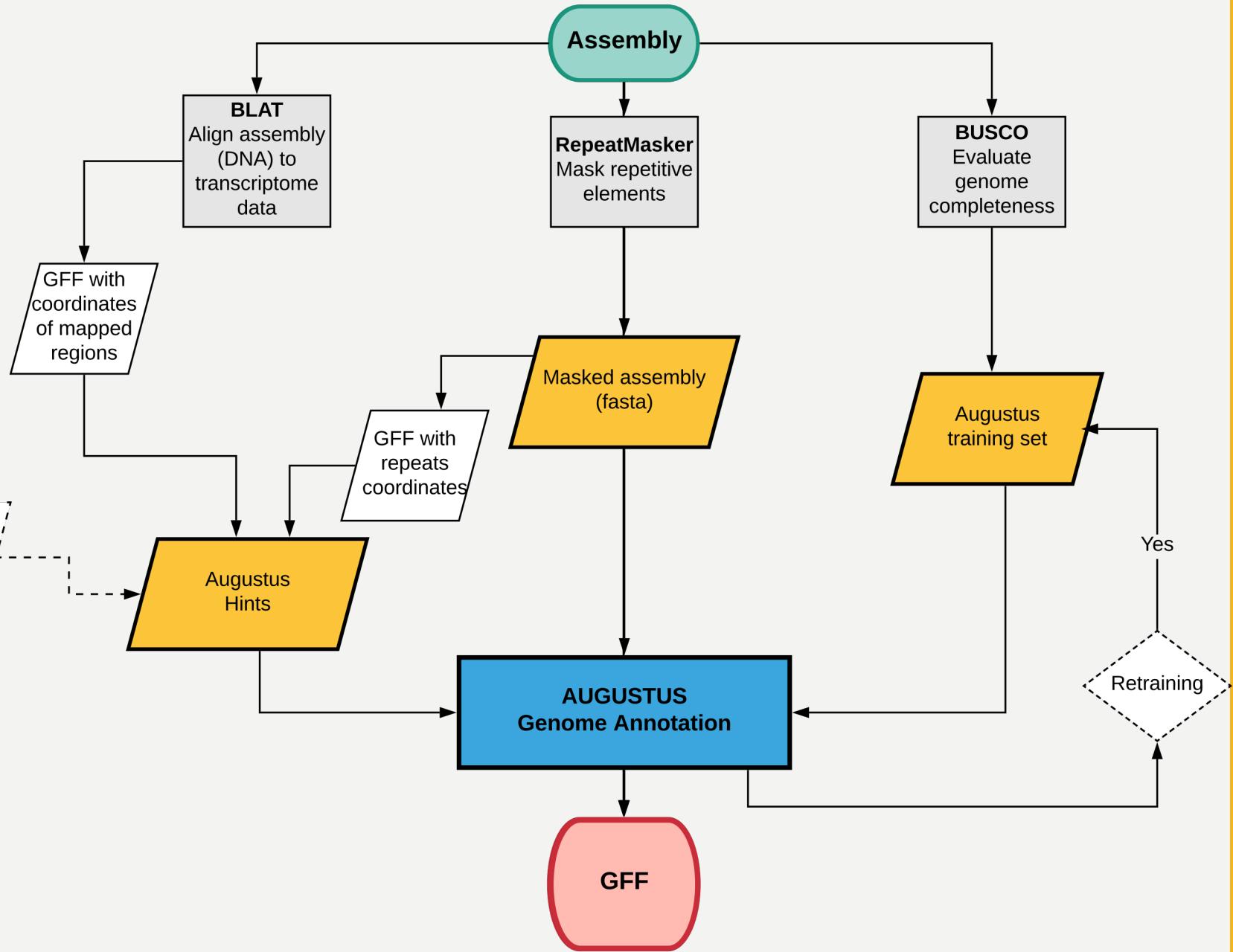


# **GENOME ANNOTATION WORKSHOP**

MIRIAN T. N. TSUCHIYA  
DATA SCIENCE POSTDOCTORAL FELLOW  
DATA SCIENCE LAB - OCIO



**WHERE ARE WE?**



# WHAT DID WE DO

## I. BUSCO

- BUSCO summary
- Retraining parameters

## 2. BLAT

- Hints

## 3. RepeatMasker

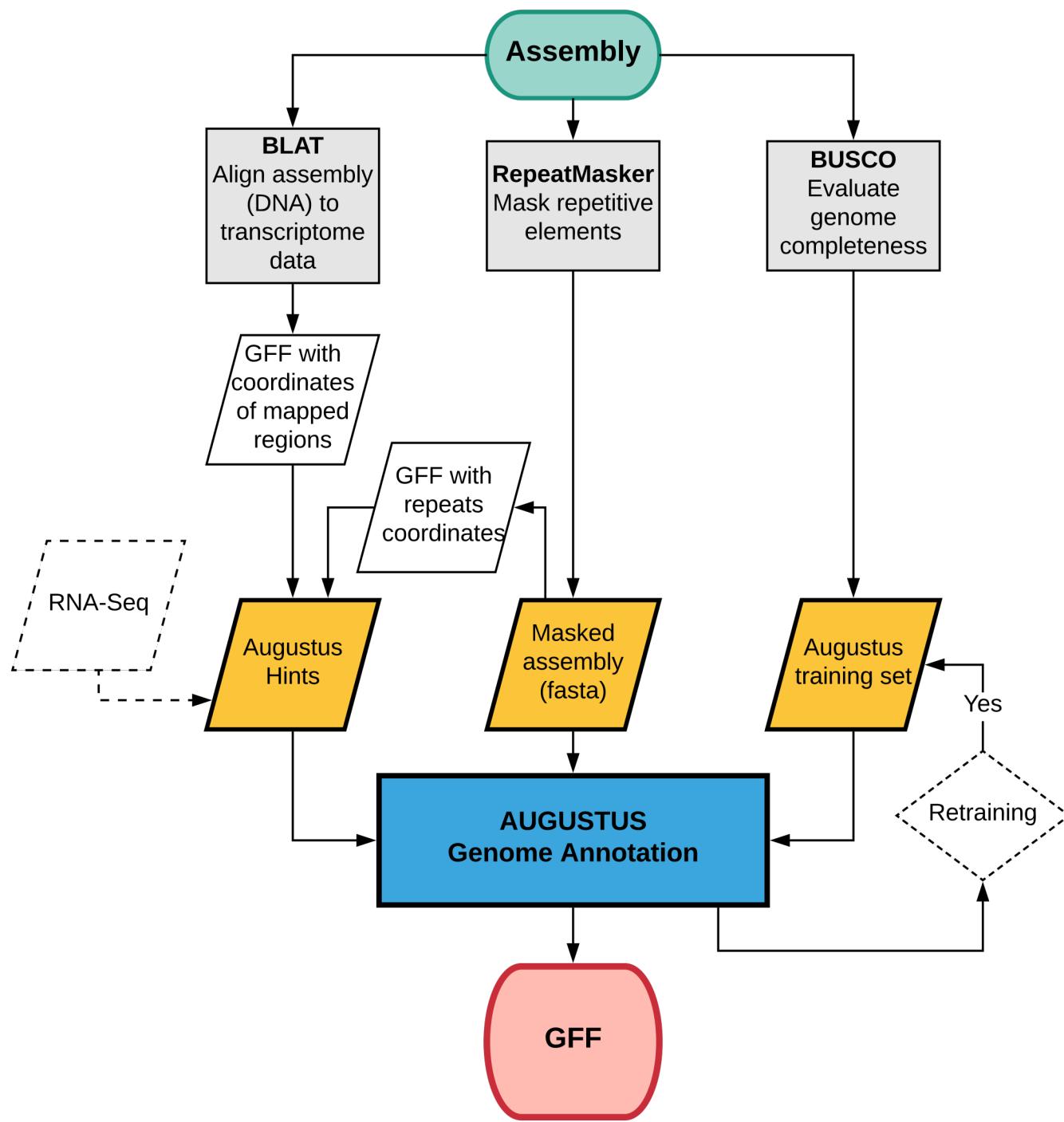
- Masked assembly

# WHAT ARE WE DOING TODAY?

1. Setting up an embarrassingly parallel AUGUSTUS run
2. Combine the results
3. Visualize it all using JBrowse

# AUGUSTUS

**FINALLY!**



# AUGUSTUS

- ab initio (internal) + evidence-driven(external)

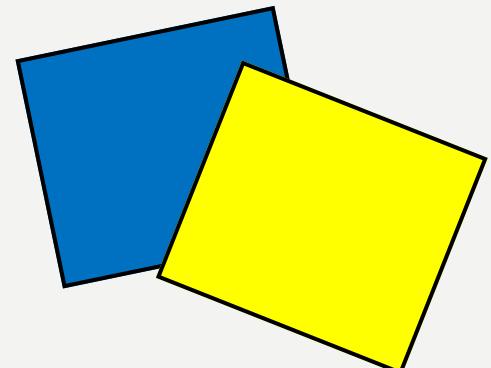
AUGUSTUS is based on a generalized hidden Markov model (GHMM), which defines probability distributions for the various sections of genomic sequences. Introns, exons, intergenic regions, etc. correspond to states in the model and each state is thought to create DNA sequences with certain pre-defined emission probabilities.

# REQUIRED FILES

1. Masked assembly ✓
2. Hints file ✓
3. Retraining parameters ✓
4. Extrinsic file ✓

# TASKS

1. Copy the masked assembly to your augustus folder
2. Download and extract EVidenceModeler to your augustus folder



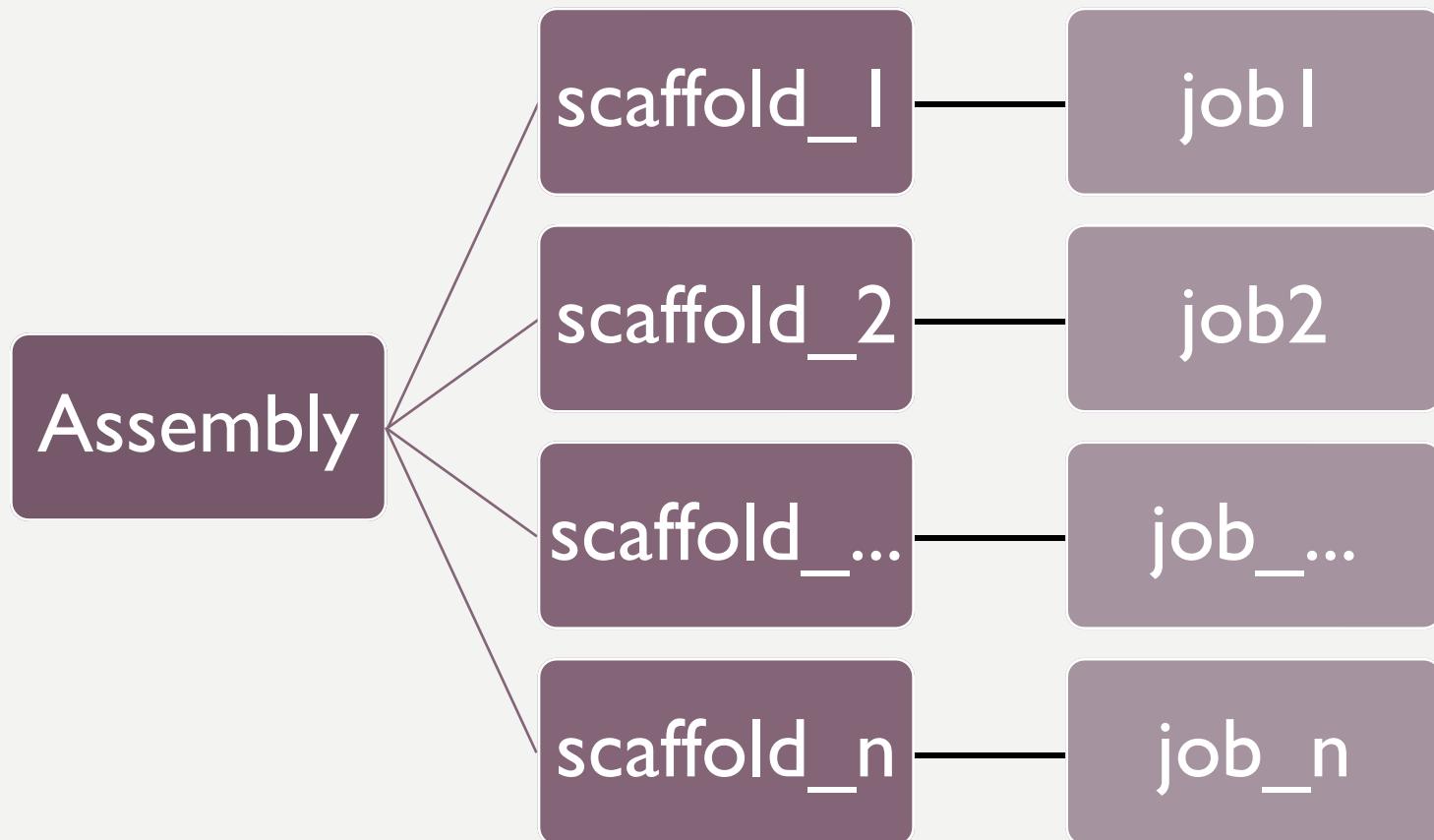
# EMBARRASSINGLY PARALLEL

- AUGUSTUS runs serially (aka one scaffold at a time)



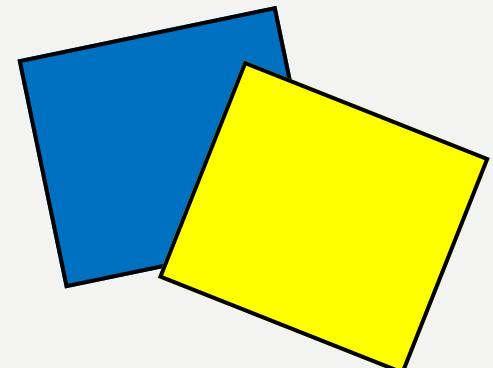
# EMBARRASSINGLY PARALLEL

- But we can “force” it to run in parallel



# TASKS

- Login to the interactive queue
- Run the EVM script `partition_EVM_inputs.pl` from your scaffolds folder
  - Don't forget that you copied the masked assembly to your augustus folder. Adjust the paths accordingly.



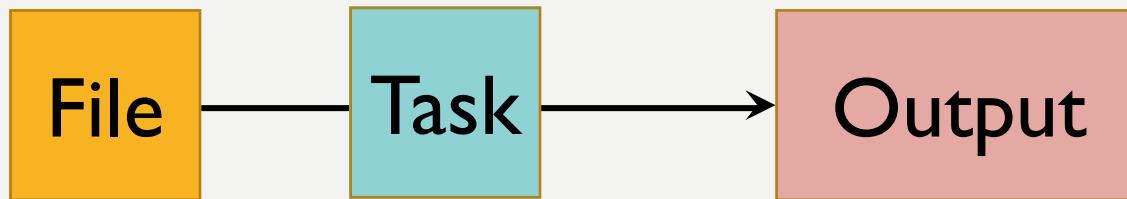
# SCAFFOLDS FOLDER

- How many folders do we have?
- Take a look in the folder `scaffold_110`.
  - What's inside?
  - Use `less` to look at the files.

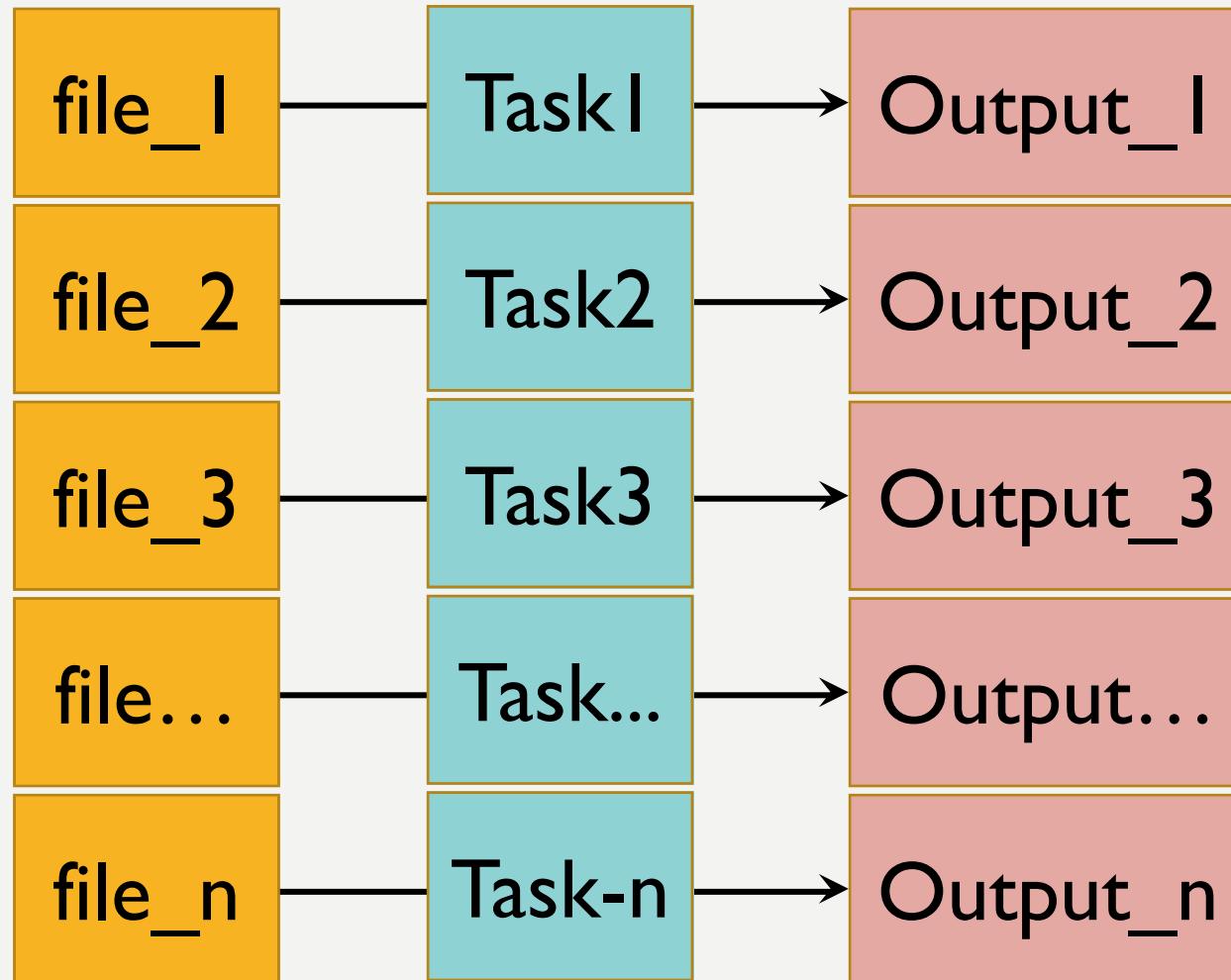
# AUGUSTUS JOB FILE

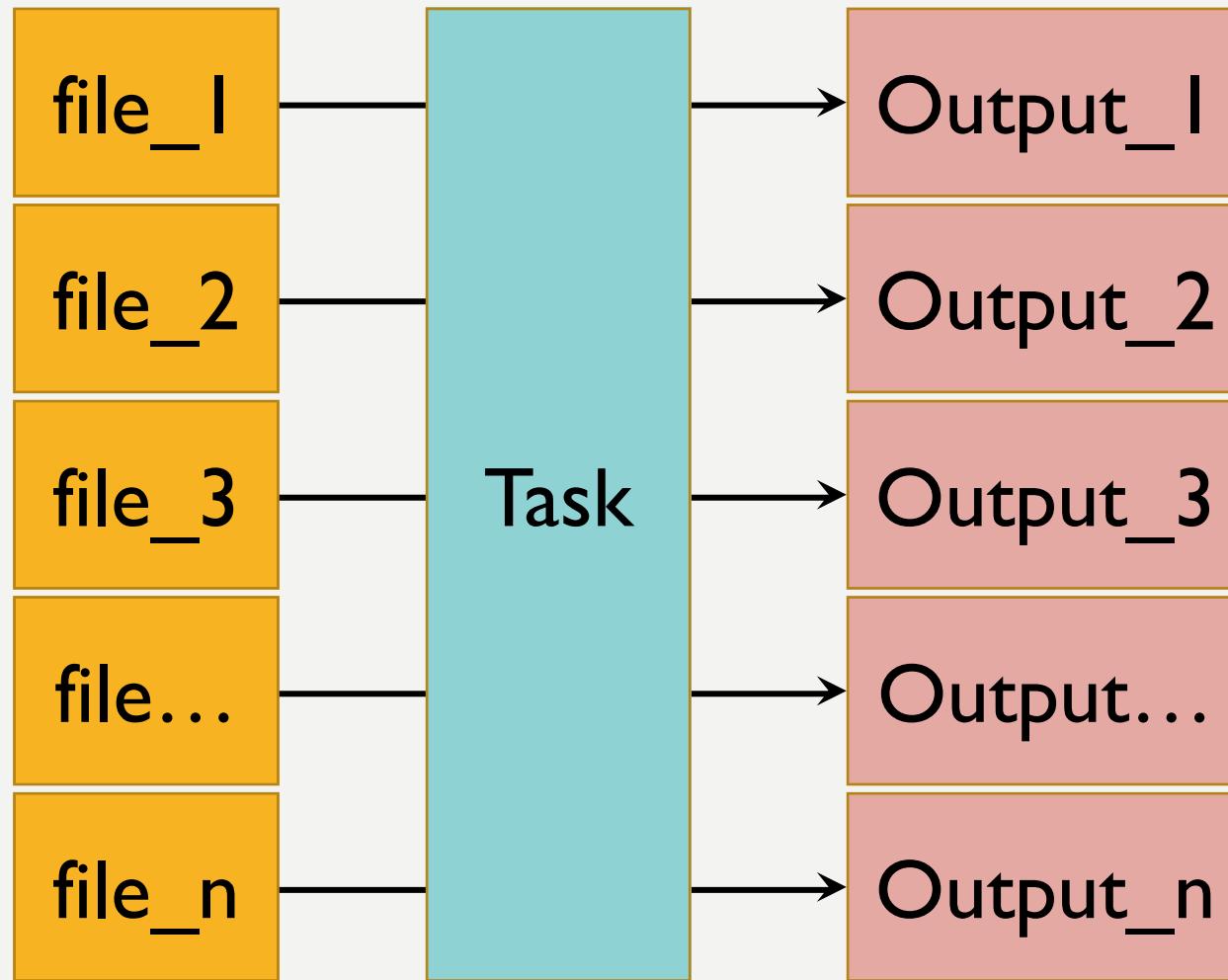
- How are we running augustus on each scaffold?
  - Option 1: Create one job file for each scaffold... manually
  - Option 2: Create one job file and use a loop to submit it.

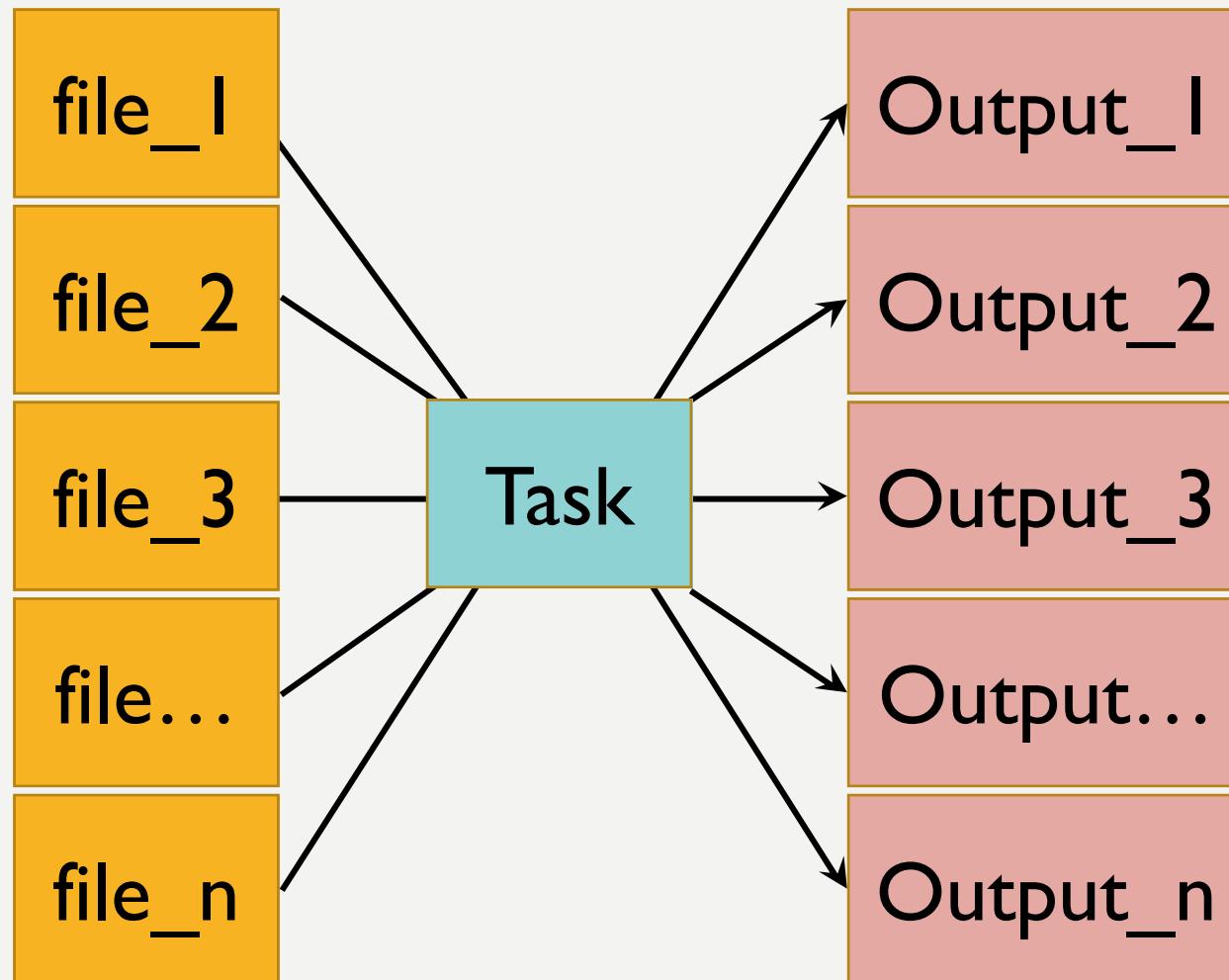
# BRIEF INTRO ABOUT LOOPS



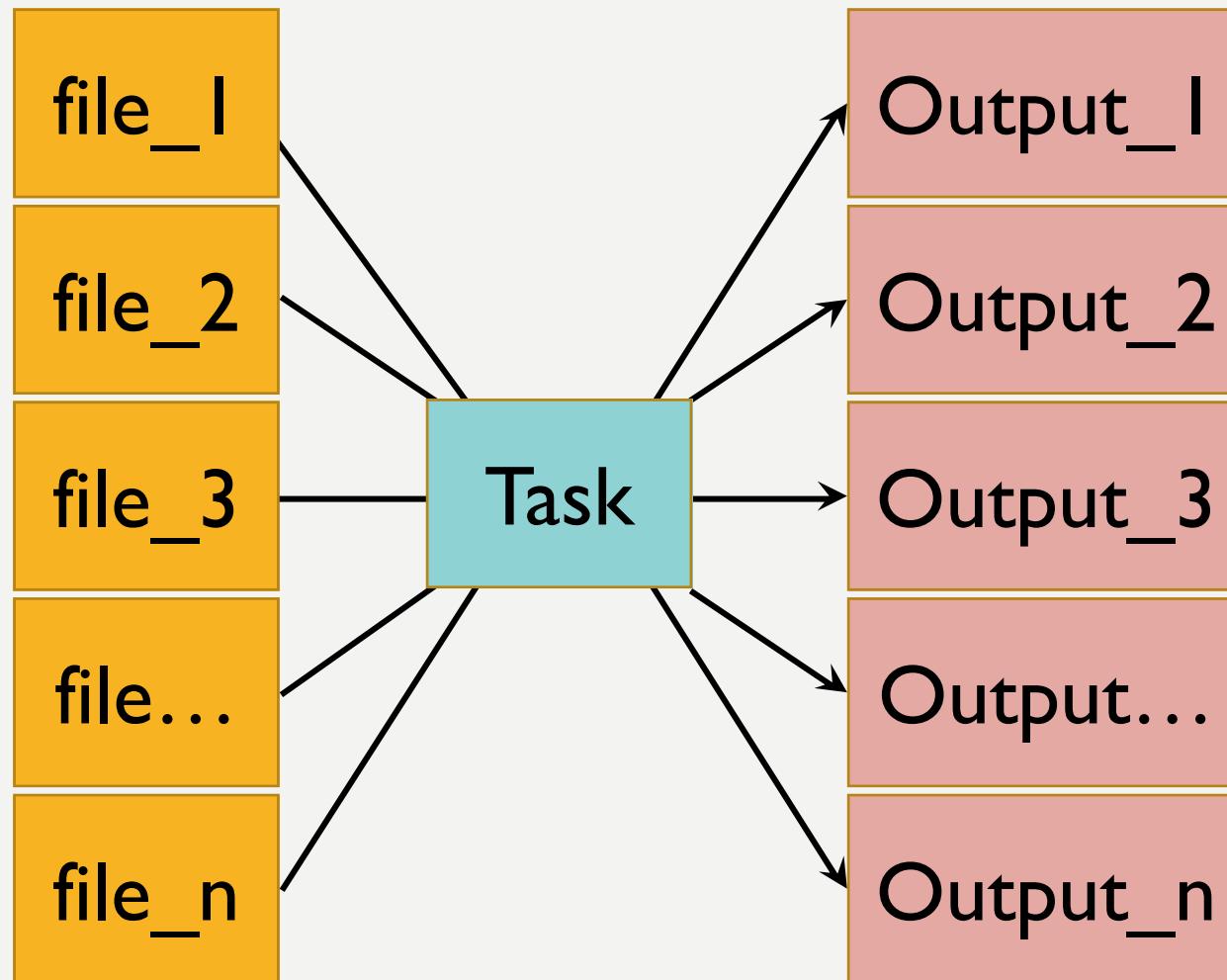
One file = One job







```
for f in file_*; do task > ${f}.output; done
```





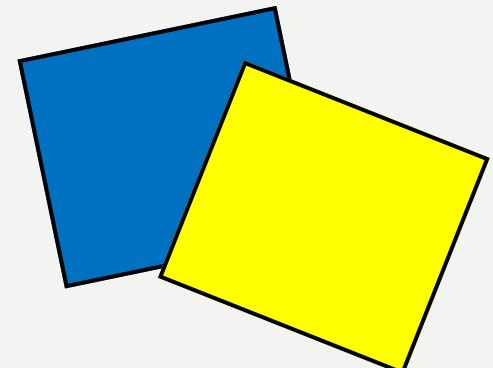
**FOR DOG IN PUPPY\_\*; DO  
PLAY \${DOG} > \${DOG}.HAPPY;  
DONE**

# BACK TO AUGUSTUS

- We will create one job (`augustus.job`) and will submit it using a `for` loop that will iterate over each scaffold - all at the same time.

# TASKS

- Create the augustus job and save it in the jobs folder.
- You will submit the job from the folder scaffolds.

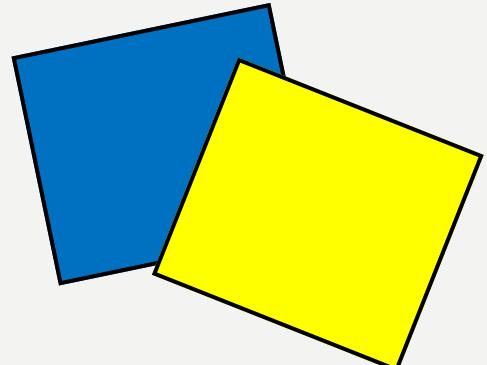


# AUGUSTUS RESULTS

- List the files in the folder outputs. You should have one gff file per scaffold.
- Use cat or less to visualize file contents.

# COMBINING THE RESULTS

- Now we want to combine all the gffs into a single gff file. `Dhydei_augustus_all.gff3`
- We will use the script `'join_aug_pred.pl'` from AUGUSTUS to combine the results.



# FROM THE FOLDER AUGUSTUS

```
cp -r  
/data/genomics/workshops/Gaworkshop/  
augustus/outputs outputs2
```



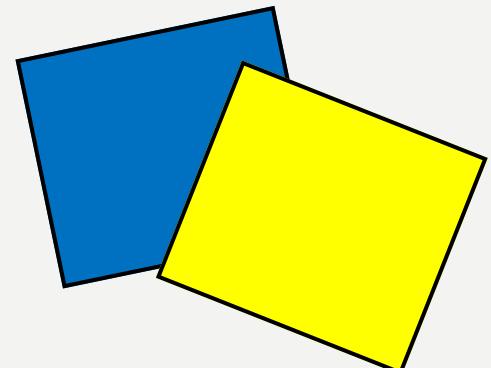
# JBROWSE

# JBROWSE



# TASKS

- Run the JBrowse scripts on your files. You need:
  - Assembly
  - GFF (final)
- Compress the file and copy it to your Desktop.



# JBROWSE

- You can deploy an instance locally (from your computer) or you can use a cloud service (AWS, Azure, etc)

# **SMITHSONIAN GENOME HUB**



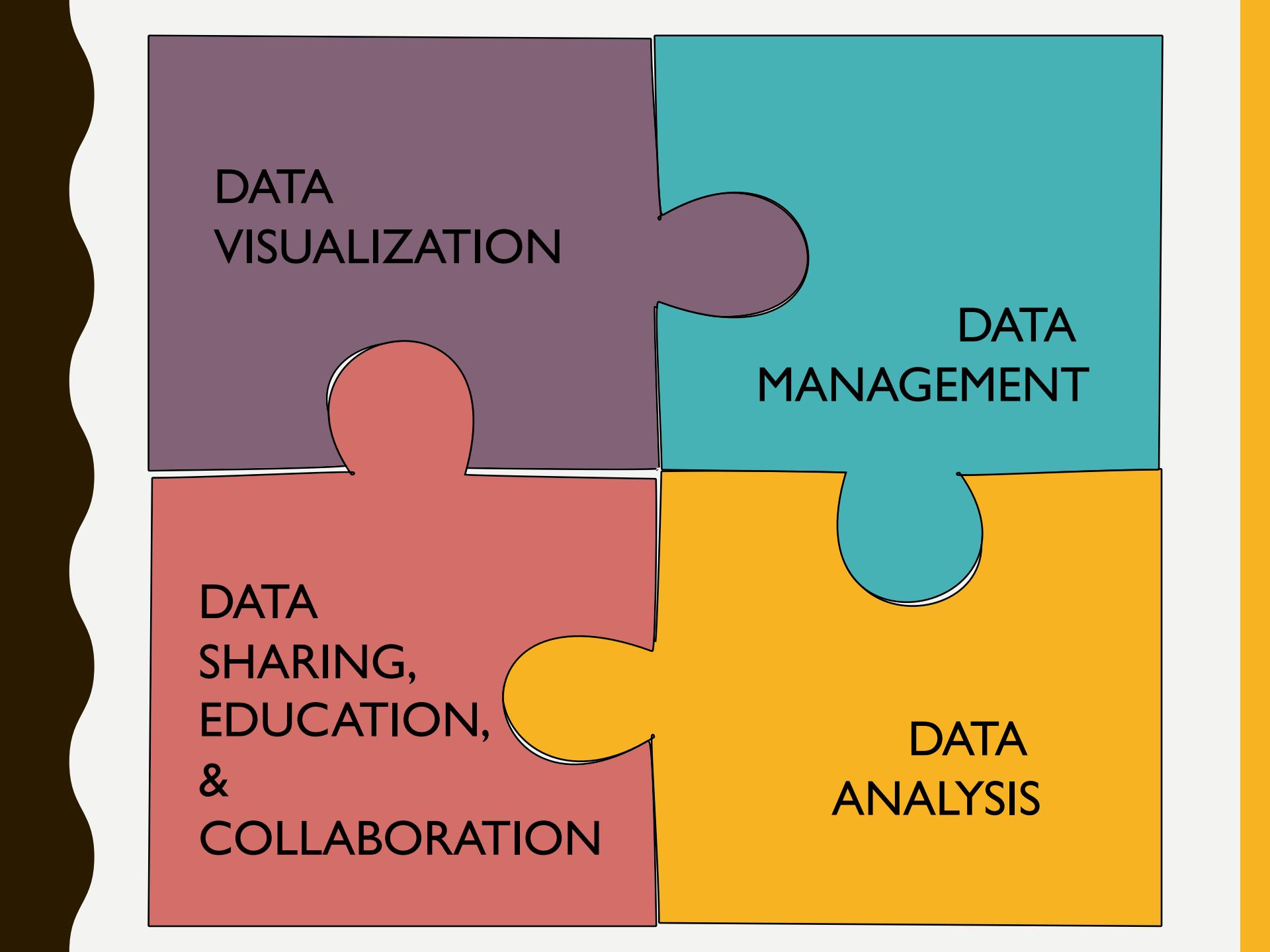
**A CLOUD BASED COLLABORATIVE  
TOOL FOR BIODIVERSITY SCIENTISTS**

species  
flies evolution fungi  
birds endangered lichens  
conservation

# biodiversity

invertebrates plants butterflies fishes  
amphibians non-model  
lichens mammals





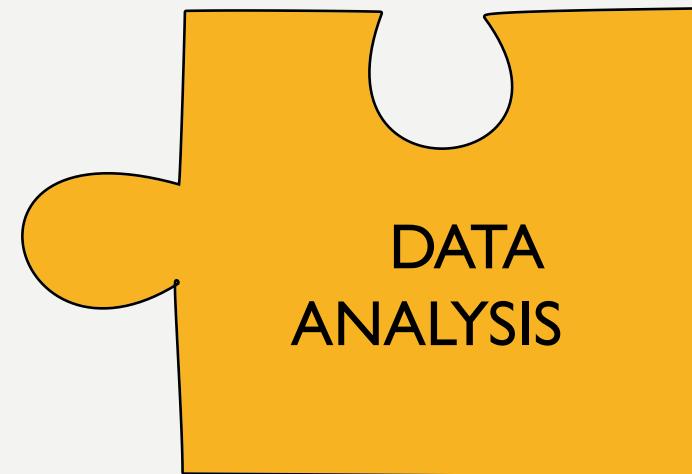
DATA  
VISUALIZATION

DATA  
MANAGEMENT

DATA  
SHARING,  
EDUCATION,  
&  
COLLABORATION

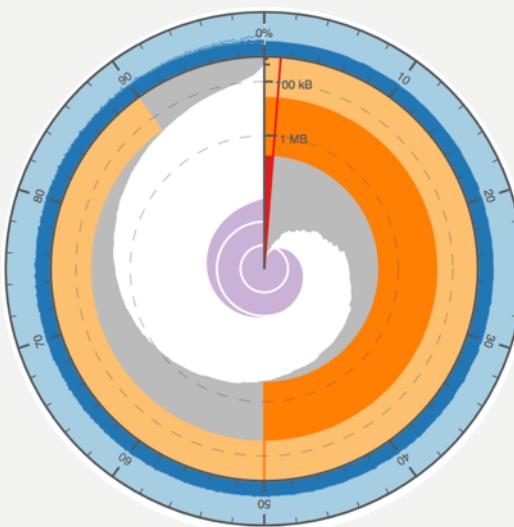
DATA  
ANALYSIS

- User-friendly web interface for data analysis
- Pipelines/workflows are easy to create and share
- Widely used software = Galaxy tools
- WebApollo for manual annotation

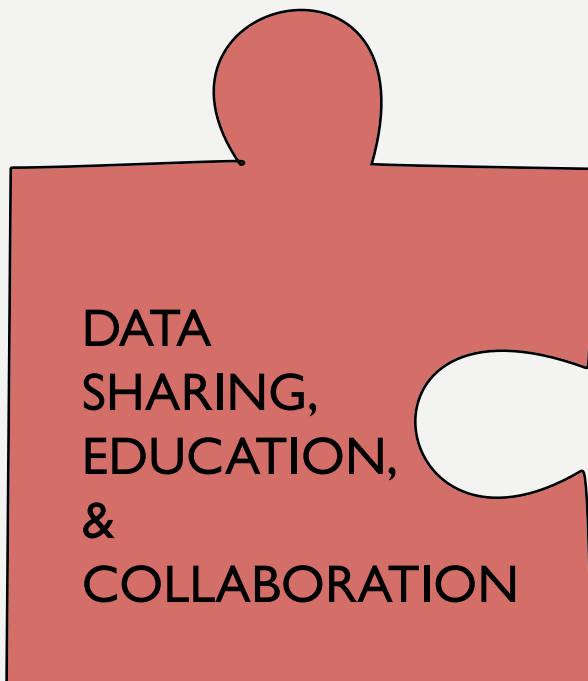


## DATA VISUALIZATION

- Individual and institutional project dashboards
- Genome browsers
- Assembly stats and other graphs useful for publication



- Facilitates collaboration and data sharing among researchers
- Creates learning opportunities for students and researchers



- Promotes diffusion of knowledge to the general public

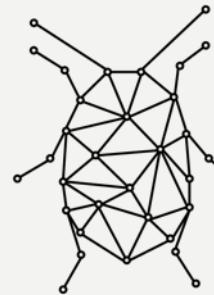
- Repository of all *de novo* genome assemblies and projects

- Reporting
- Managing
- Organization



*Smithsonian's*  
~~MISS PEREGRINE'S~~  
HOME FOR  
~~PECULIAR CHILDREN~~  
*genomes*





OCIO  
DATA  
SCIENCE  
LAB



@SIDatascience



datascience.si.edu



[tsuchiyam@si.edu](mailto:tsuchiyam@si.edu)



@MirianTsuchiya



Hydra help: SI-HPC@si.edu

Bug #415

Bug #416

Bug #417

Bug #418

Bug #419

Bug #420

