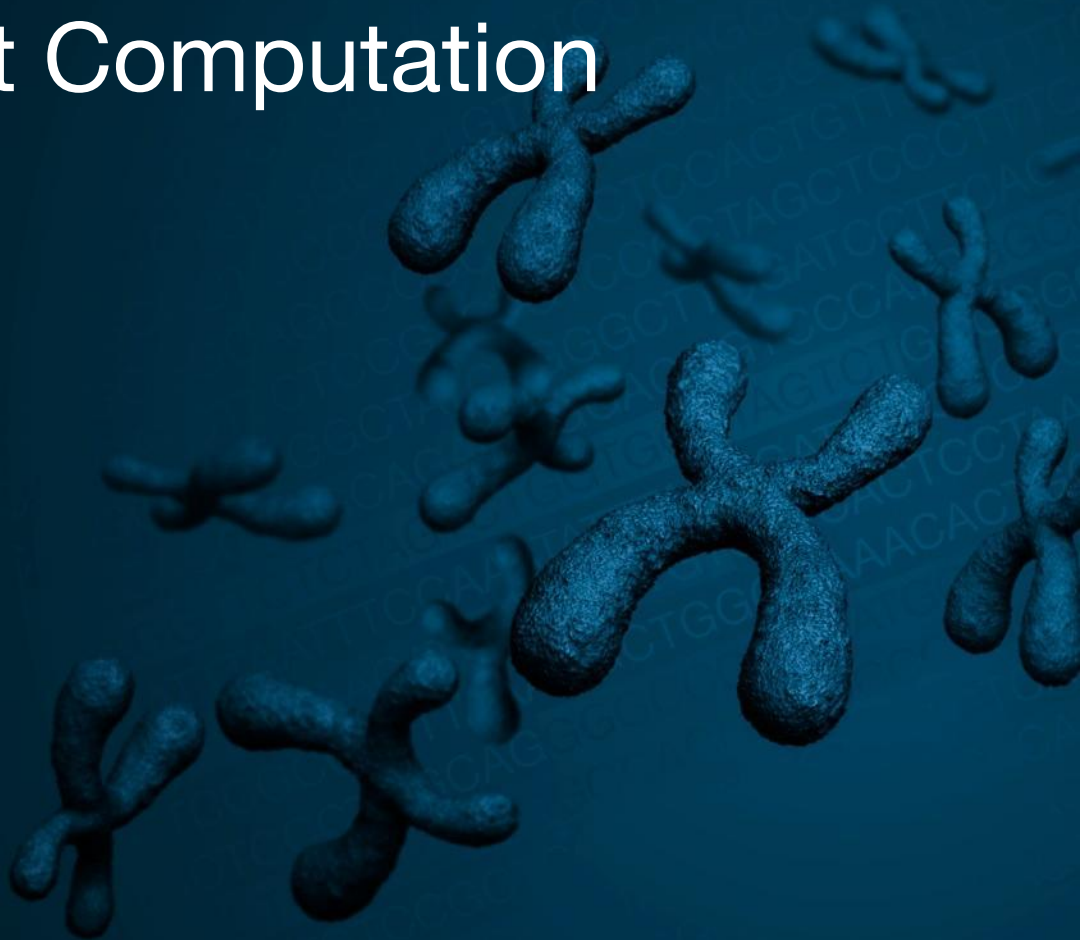


# eHive: a System for Massive High-throughput Computation

Brandon Walts  
Leo Gordon  
Matthieu Muffato  
Andrew Yates  
Paul Flicek



# Motivation

- Problem space
  - Large data sets, growing faster than Moore's law
  - Many analyses are easy to run in parallel
  - Analysis tools tend to follow the UNIX philosophy - do one thing but do it well
  - Lots of data handling between analyses
  - Provenance and reproducibility are important
  - Regular repeats of analysis as part of a production cycle
- Infrastructure
  - Large compute farm, managed by a scheduler (LSF)
  - Data in a mix of RDBMS and flat files

# Thus, eHive

- First release in 2004
- Currently controls 450 cpu-years per year of compute for Ensembl
- Adopted by several institutions outside of EMBL-EBI

# What's with the name?

eHive's approach to computation is based on a swarm of autonomous agents - naturally leading to a beehive metaphor:

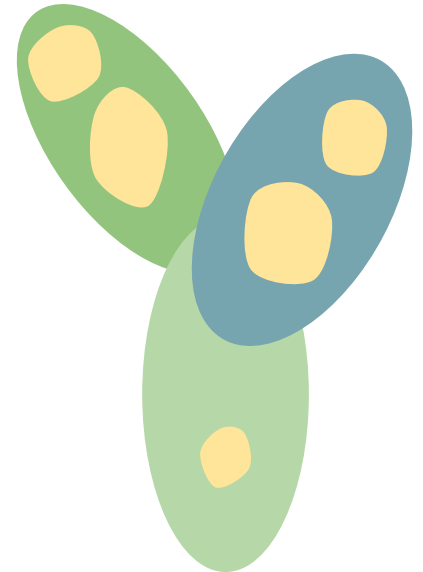
# What's with the name?

- Independent agents ("workers") perform computation.



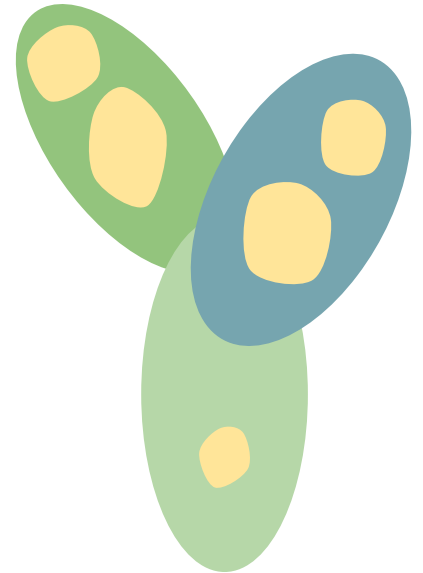
# What's with the name?

- Agents have access to resources of different types ("meadows").



# What's with the name?

- There is an overseeing process ("beekeeper"), but it is lightweight -- concerned with managing worker population and identifying problems.

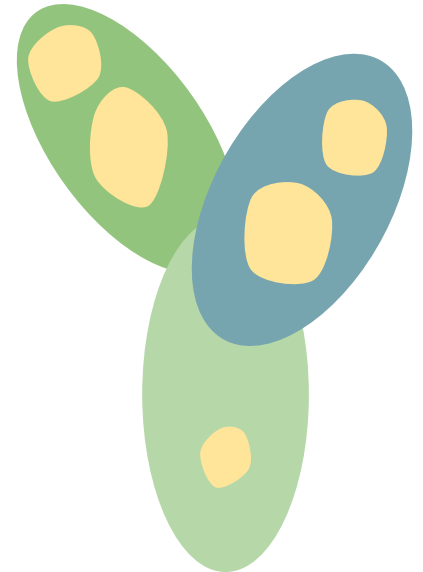


# What's with the name?

- There is a central database ("blackboard") that workers update to coordinate their activity



analysis	job	state
download	1	done
align	2	ready



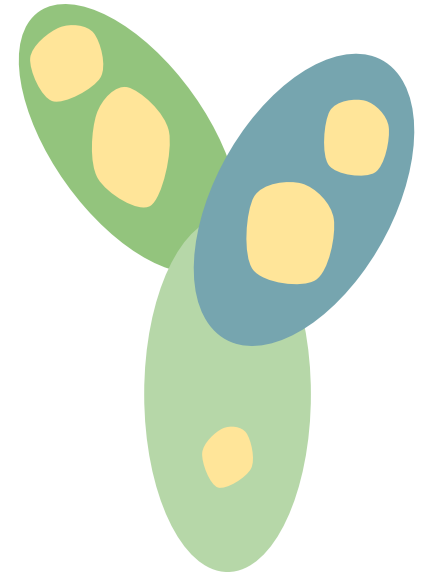


# eHive lifecycle

- The beekeeper checks the current job list and worker population, creating new workers if necessary

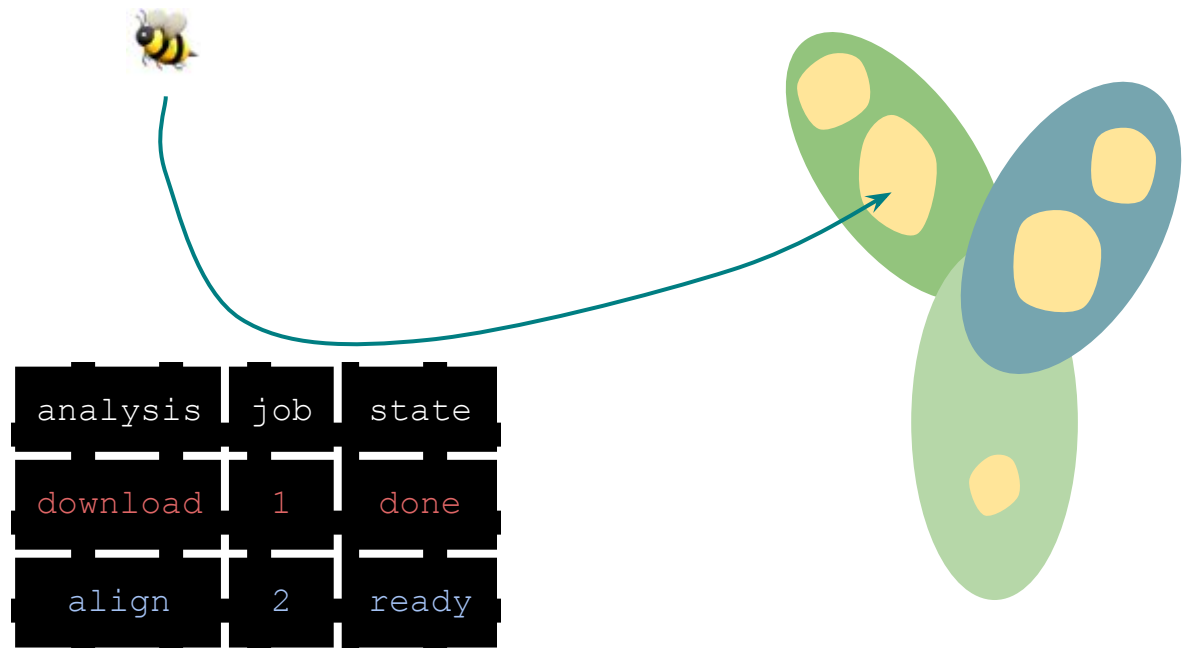


analysis	job	state
download	1	done
align	2	ready



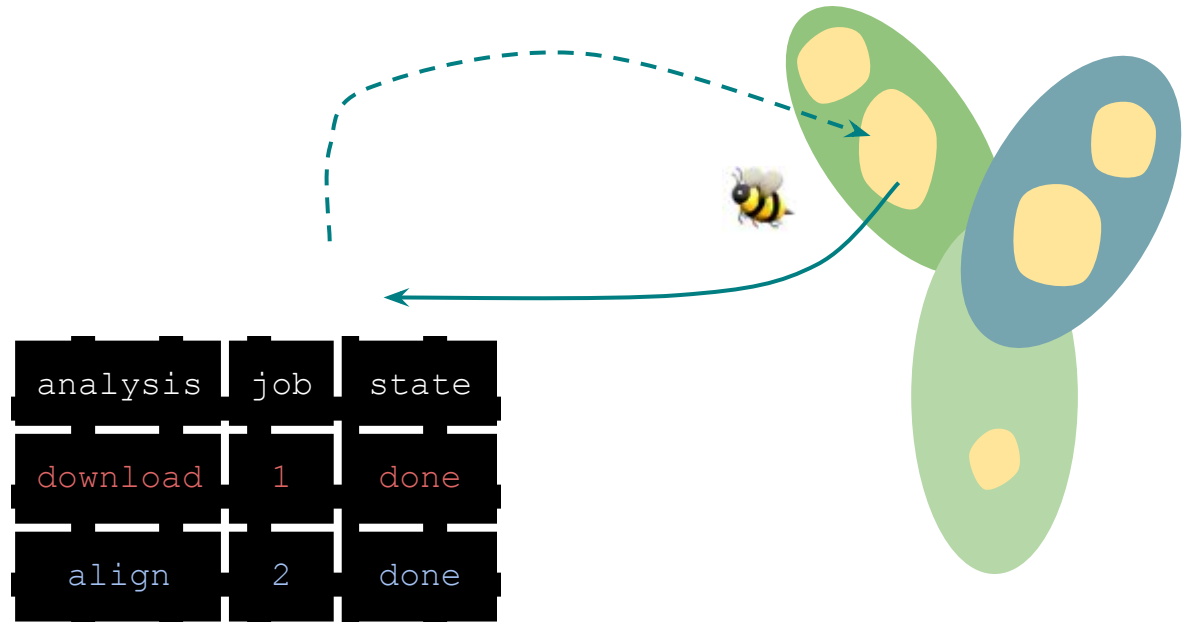
# eHive lifecycle

- Worker checks for a job it is able to execute, claims it, specializes if necessary, and begins execution



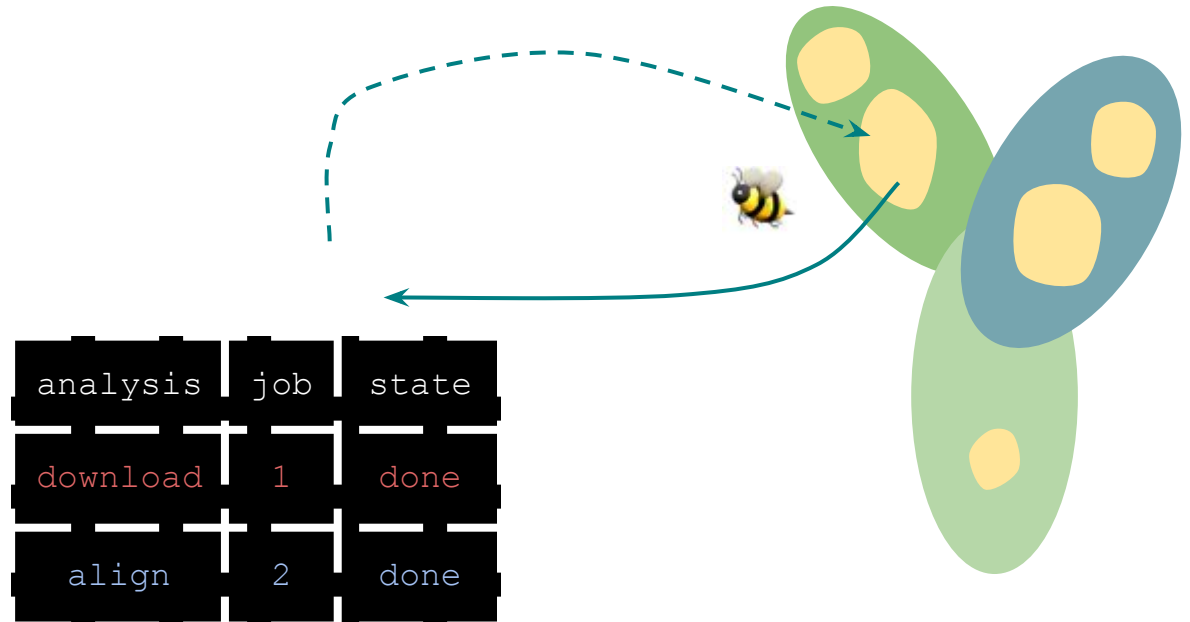
# eHive lifecycle

- Worker completes running job. Updates the job list, then checks to see if there is more work to do.



# eHive lifecycle

- For short jobs, workers can claim a batch to do in one cycle, reducing dispatch and startup overhead



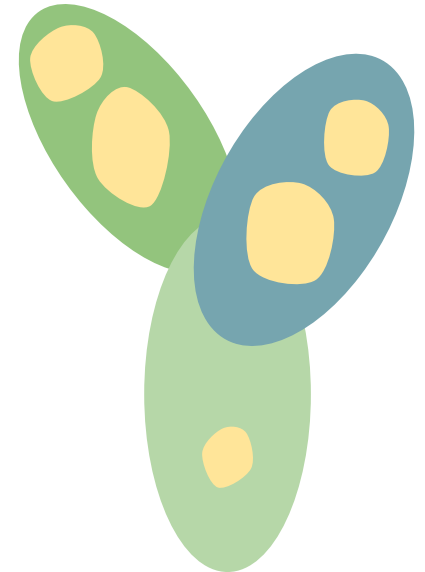
# eHive lifecycle

- If errors occur, the beekeeper notes this and can partially or completely stop workflow execution

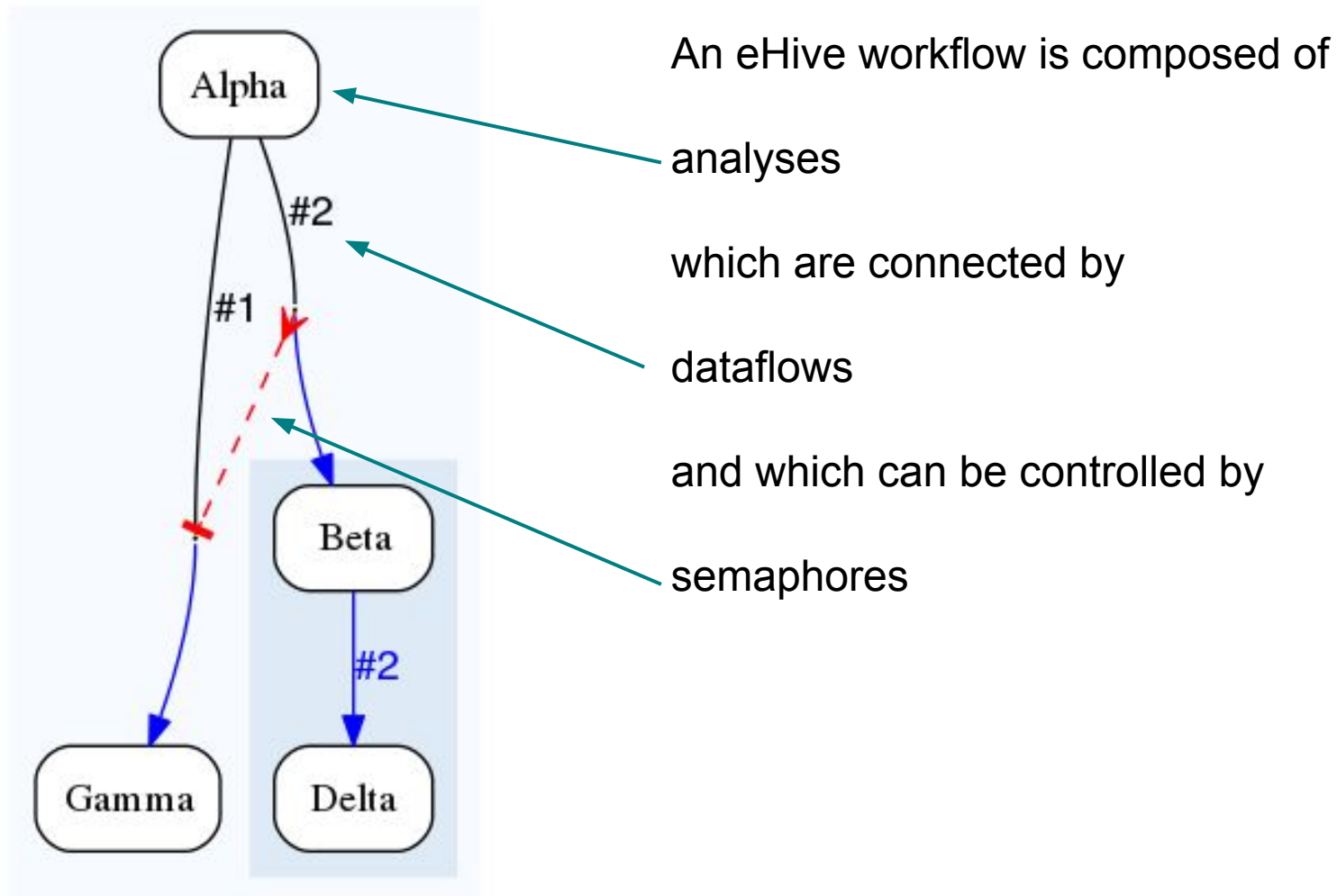


STOP

analysis	job	state
download	1	done
align	2	FAIL

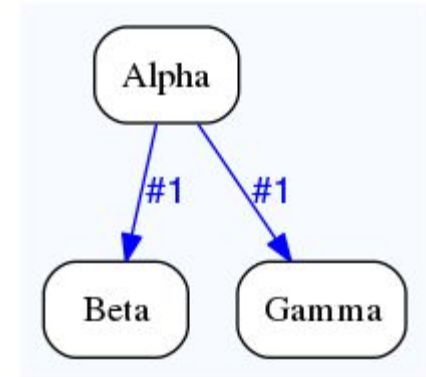


# Workflow structure

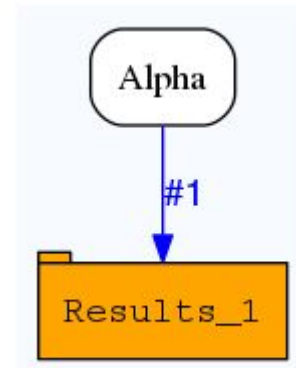


# Dataflows and events

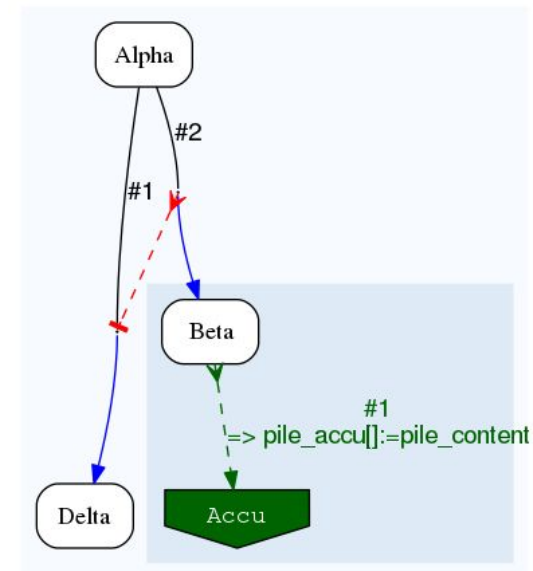
- When each job runs, it can generate zero or more dataflow events.
- Different events can be transmitted on different "branches."
- The consequences of these events are determined by how they are wired in the workflow



Create ("seed") new jobs



Store data in a table



Store data in an "accu" data structure

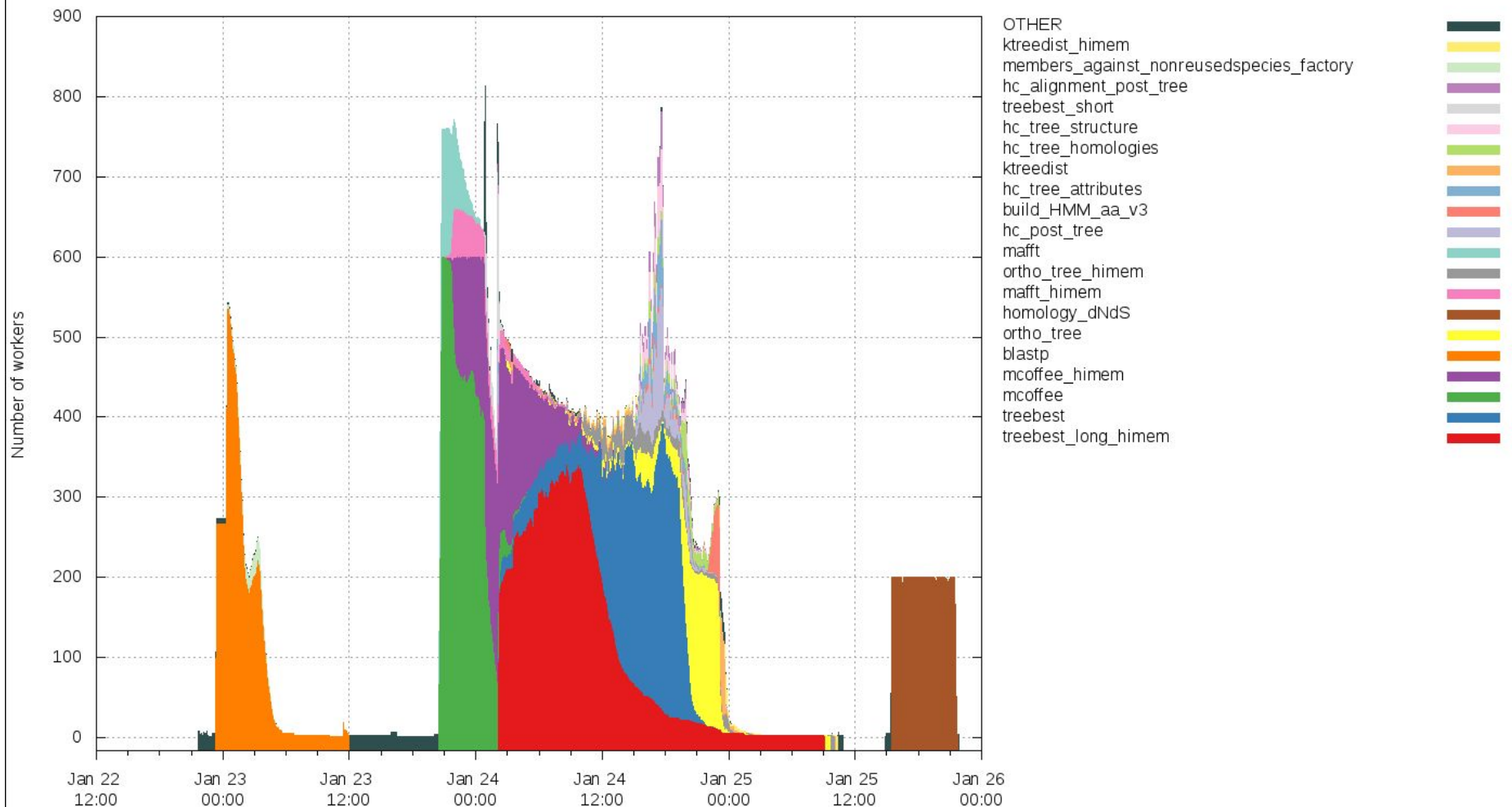
# What's in the box?

- eHive code as a collection of Perl modules
- Scripts to instantiate and execute workflows
- Visualization and debugging utilities
  - Workflow structure
  - Resource usage
- Interfaces for different schedulers
  - LSF is officially supported
  - SGE, PBS Pro, and HTCondor reference implementations available
- guiHive - web based workflow management tool

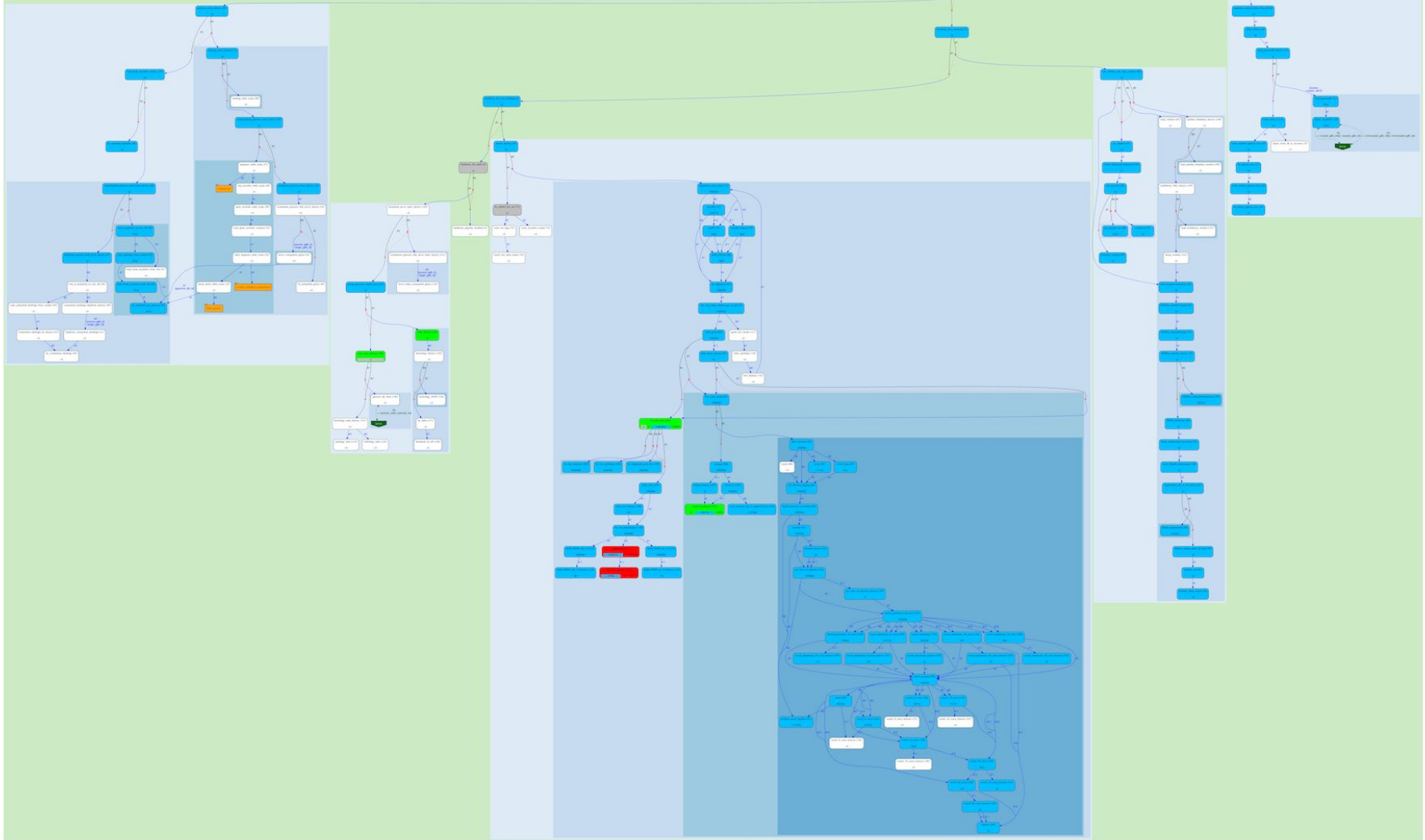


# Workflow analytics

Profile of the 20 top-analysis of mysql://ensro@compara1/mm14\_protein\_trees\_79 to 2015-01-26



# Use cases...



# Obtaining eHive and guiHive



<https://ensembl-hive.readthedocs.io/>



<https://github.com/Ensembl/ensembl-hive>  
<https://github.com/Ensembl/guiHive>



<https://hub.docker.com/r/ensemblorg/guihive/>

# Ensembl Acknowledgements

## The Entire Ensembl Team

Bronwen L. Aken<sup>1</sup>, Premanand Achuthan<sup>1</sup>, Wasiu Akanni<sup>1</sup>, M. Ridwan Amode<sup>1</sup>,  
Friederike Bernsdorff<sup>1</sup>, Jyothish Bhai<sup>1</sup>, Konstantinos Billis<sup>1</sup>, Denise Carvalho-Silva<sup>1</sup>,  
Carla Cummins<sup>1</sup>, Peter Clapham<sup>2</sup>, Laurent Gil<sup>1</sup>, Carlos García Girón<sup>1</sup>, Leo Gordon<sup>1</sup>,  
Thibaut Hourlier<sup>1</sup>, Sarah E. Hunt<sup>1</sup>, Sophie H. Janacek<sup>1</sup>, Thomas Juettemann<sup>1</sup>,  
Stephen Keenan<sup>1</sup>, Matthew R. Laird<sup>1</sup>, Ilias Lavidas<sup>1</sup>, Thomas Maurel<sup>1</sup>, William McLaren<sup>1</sup>,  
Benjamin Moore<sup>1</sup>, Daniel N. Murphy<sup>1</sup>, Rishi Nag<sup>1</sup>, Victoria Newman<sup>1</sup>, Michael Nuhn<sup>1</sup>,  
Chuang Kee Ong<sup>1</sup>, Anne Parker<sup>1</sup>, Mateus Patricio<sup>1</sup>, Harpreet Singh Riat<sup>1</sup>, Daniel Sheppard<sup>1</sup>,  
Helen Sparrow<sup>1</sup>, Kieron Taylor<sup>1</sup>, Anja Thormann<sup>1</sup>, Alessandro Vullo<sup>1</sup>, Brandon Walts<sup>1</sup>,  
Steven P. Wilder<sup>1</sup>, Amonida Zadissa<sup>1</sup>, Myrto Kostadima<sup>1</sup>, Fergal J. Martin<sup>1</sup>,  
Matthieu Muffato<sup>1</sup>, Emily Perry<sup>1</sup>, Magali Ruffier<sup>1</sup>, Daniel M. Staines<sup>1</sup>, Stephen J. Trevanion<sup>1</sup>,  
Fiona Cunningham<sup>1</sup>, Andrew Yates<sup>1</sup>, Daniel R. Zerbino<sup>1</sup> and Paul Flicek<sup>1,2,\*</sup>

<sup>1</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK and <sup>2</sup>Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

## Funding



Co-funded by the  
European Union