# Galaxy / CEGR: Pugh/Mahony Labs, Center for Eukaryotic Gene Regulation

Greg consulted with Dr. Frank Pugh and Dr. Shaun Mahony in the Center for Eukaryotic Gene Regulation to build the Galaxy / CEGR instance (see Figure 1). This environment is hosted within a virtual machine (VM) on ICDS-ACI and has access to cluster nodes for job execution. The environment enables researchers in the Center for Eukaryotic Gene Regulation (CEGR) to perform ChIP-exo analyses on very large datasets produced by the Center's sequencer.
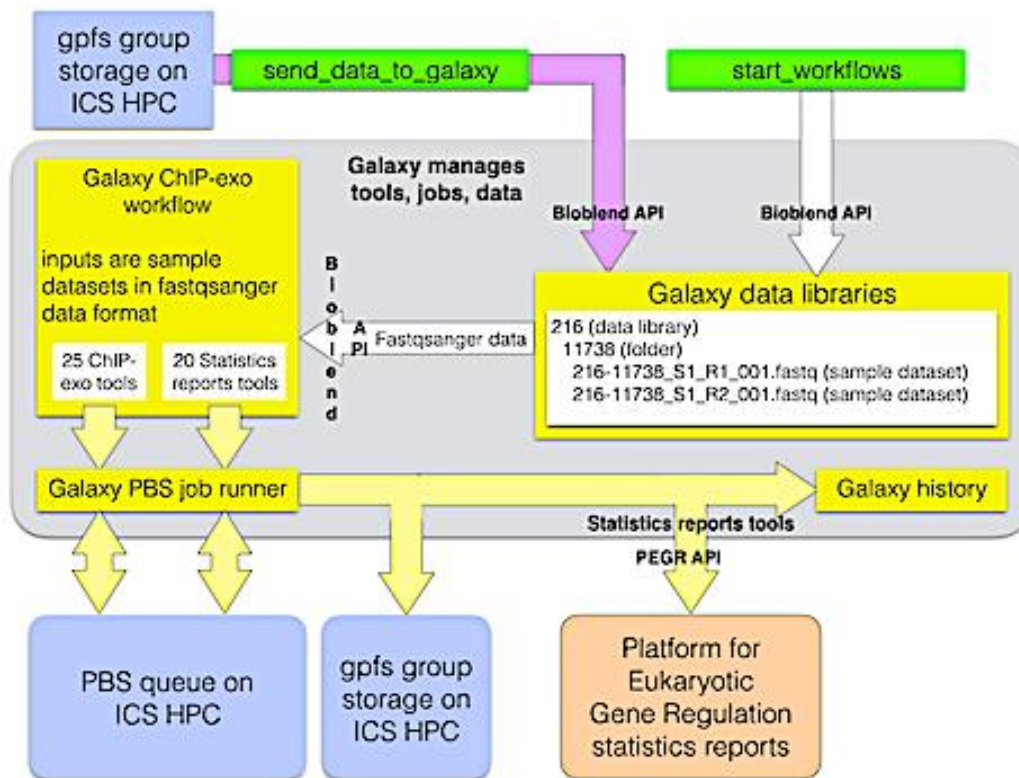


**Figure 1: Galaxy / CEGR Process Flow**

At the time that it was built, the environment consisted of four primary components, an Illumina NextSeq 500 sequencer owned by the Pugh/Mahony labs, the CEGR pre-processing pipeline, the Galaxy ChIP-exo environment and the Platform for Eukaryotic Gene Regulation (PEGR), a web application developed by the labs.

The sequencer produces raw datasets that are converted into the fastqsanger data format and imported into Galaxy data libraries (Galaxy's hierarchical container for datasets). The Galaxy ChIP-exo workflow is automatically executed for each set of samples, with each tool in the workflow submitting jobs to the HPC cluster nodes, producing datasets that are then used as inputs to the next tool in the workflow chain. The Galaxy ChIP-exo workflow includes tools that generate metadata (fine-grained statistics) about datasets produced by specified tools within the workflow. These statistics are sent to PEGR for review.

The CEGR pre-processing pipeline consists of four custom programs, developed to automate the processes used by the labs, ultimately preparing the data for analysis within the Galaxy ChIP-exo instance (the green boxes in Figure 2 below depict these custom programs). Each of these programs includes quality assurance components that automatically halt processing if errors occur, logging the details for review and correction. Each program can be executed independently (assuming that the previous program in the pipeline has completed successfully) allowing for a certain step to be re-executed after corrections are made.
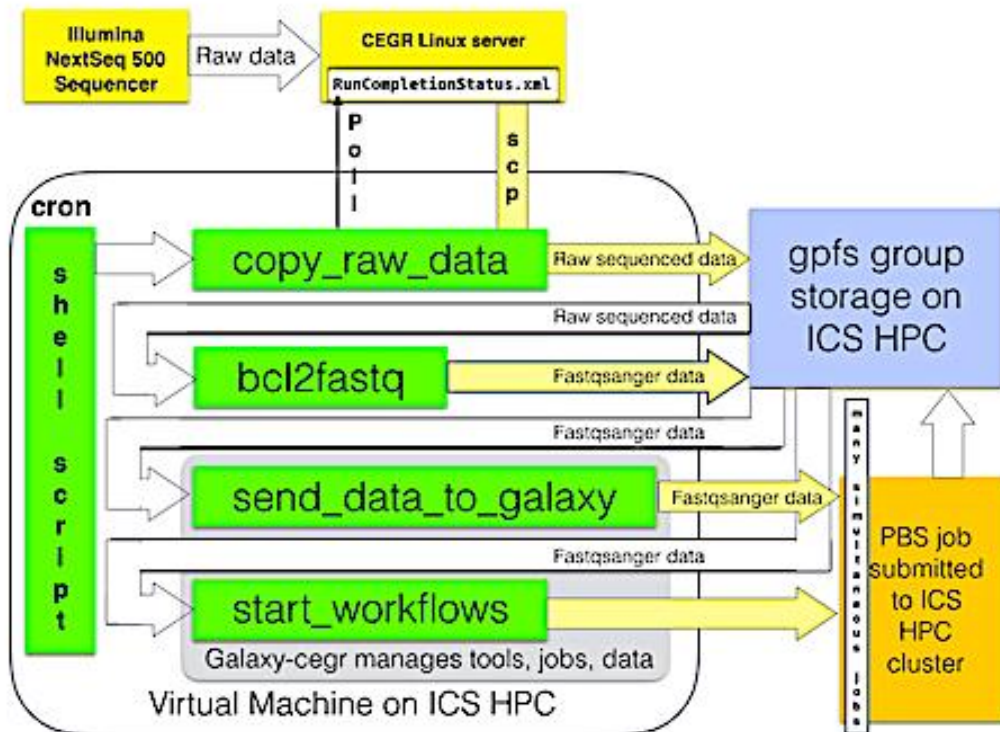


**Figure 2: Galaxy / CEGR Pre-processing Pipeline**

- **copy_raw_data** - Polls the lab's server that contains the raw datasets produced by the sequencer to determine when the sequencing run is complete.  The raw datasets are then copied from the server to the Science Gateway's file store.
- **bcl2fastq** – Converts the raw datasets into the fastqsanger data format required by the initial tools within the Galaxy ChIP-exo workflow.
- **send_data_to_galaxy** – Creates a Galaxy data library for the run's sample datasets.  The sample datasets are imported into appropriate folders from which they can be retrieved for analysis.  The program submits a PBS job to the ICS HPC cluster to import each sample dataset, storing them on the Science Gateway's file store.
- **start_workflows** – Retrieves sample datasets from the run's Galaxy data library folders and provides them as input to the Galaxy ChIP-exo workflow.  Each tool in the workflow performs its function by submitting a PBS job to the HPC cluster, storing the resulting datasets on the Science Gateway's file store.

Galaxy includes a feature rich REST API which is used by the pre-processing pipeline for all direct interaction with Galaxy.

The Galaxy environment is configured with 8 front-end web server processes.  It contains the ChIP-exo workflows for both single and paired reads.  These workflows consist of the chain of tools that perform the ChIP-exo analysis and send the statistics to PEGR (see Figure 3).
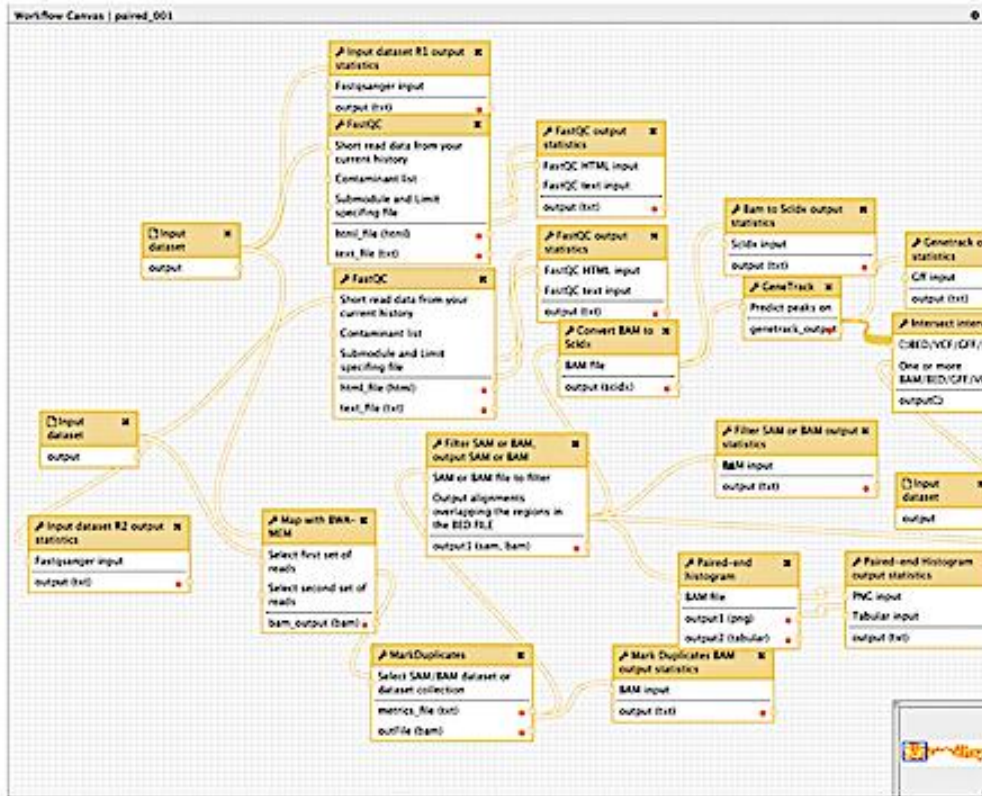
**Figure 3: Galaxy / CEGR ChIP-exo Workflow**

The Galaxy environment is also configured with 8 job handler processes that submit PBS jobs to the ICS HPC cluster nodes, producing datasets that are stored on the Science Gateway's file store.  This configuration allows for load balancing on the web front-end and the many simultaneous PBS jobs submitted to the HPC cluster.

The Platform for Eukaryotic Gene Regulation (PEGR) is a web-based sample tracking application developed by the Pugh/Mahony labs.  Lab technicians enter sample information into PEGR that is used by the CEGR pre-processing pipeline to prepare the data for analysis within Galaxy.  PEGR interacts with Galaxy via the Galaxy REST API.

**Galaxy / CoralSNP: Baums lab, Department of Biology**

Greg consulted with Dr. Iliana Baums in the Department of Biology to develop the Galaxy / CoralSNP analysis environment.  Reliable and standardized identification of genotypes is a critical need for basic research and restoration planning, especially in plants and animals that reproduce asexually and form large clonal populations

like aspens, seagrasses, waterfleas and corals.  Figure 4 provides the general process.

- The scientist collects the coral, extracts the DNA and submits it to the processing facility
- The raw sequenced samples and sample metadata are uploaded to the Galaxy / CoralSNP environment for analysis
- The analysis produces new sample MLGs and genotype information that can be accessed via the Galaxy / CoralSNP interface
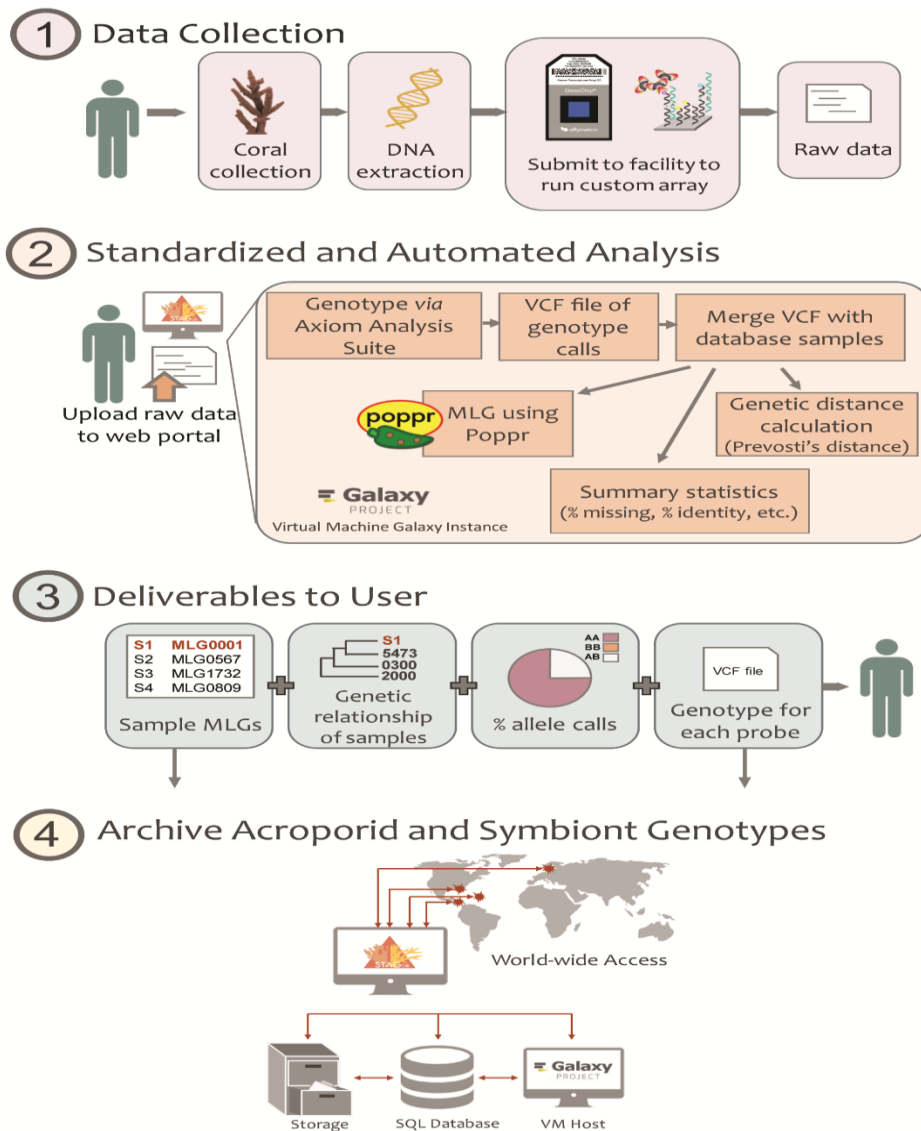


**Figure 4: General process for standardized identification of genotypes**

Galaxy / CoralSNP provides a high-resolution hybridization-based genotype array coupled with a standardized analysis workflow and database for the most speciose genus of coral, Acropora, and their symbionts (see figure 5).
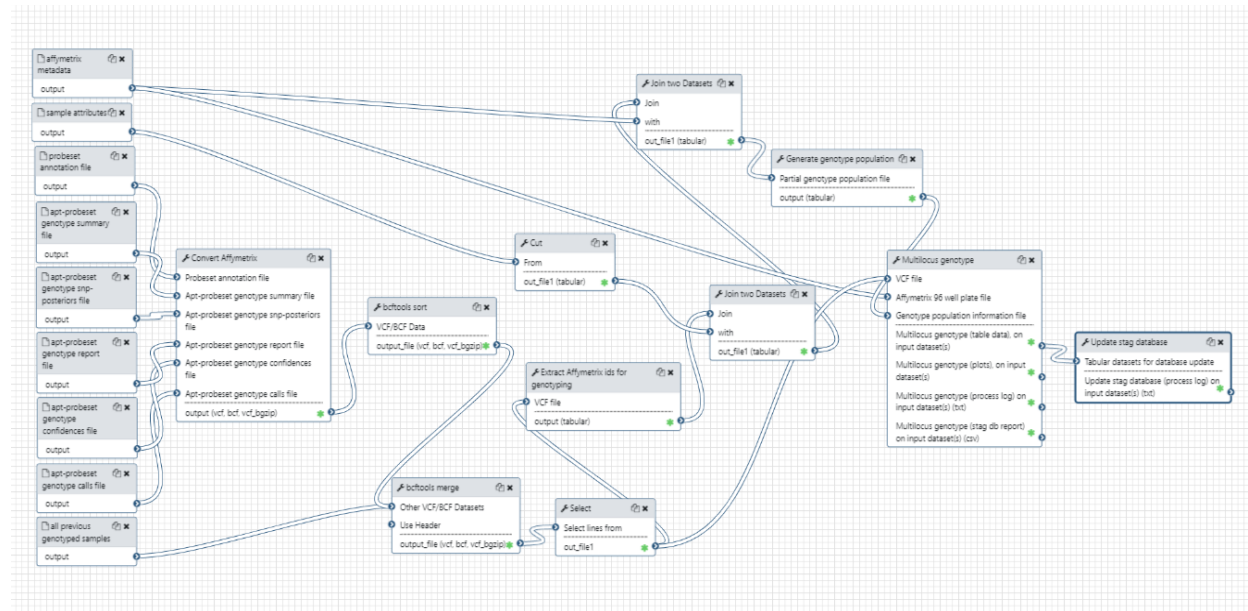


**Figure 5: Galaxy / CoralSNP Workflow**

The process is simple.  A sample metadata file is created by the user from a template and uploaded along with their raw Affymetrix data files into the Galaxy CoralSNP environment.  The appropriate files are selected as inputs to the Queue Genotype Workflow tool (Figure 6) which validates the metadata, executes the CoralSNP workflow (Figure 5) and updates a dataset that contains all previously genotyped samples as well as the STAG database with the samples in the current run.  From the user's perspective, the entire analysis is as simple as uploading data and specifying it as the input to execute a tool.
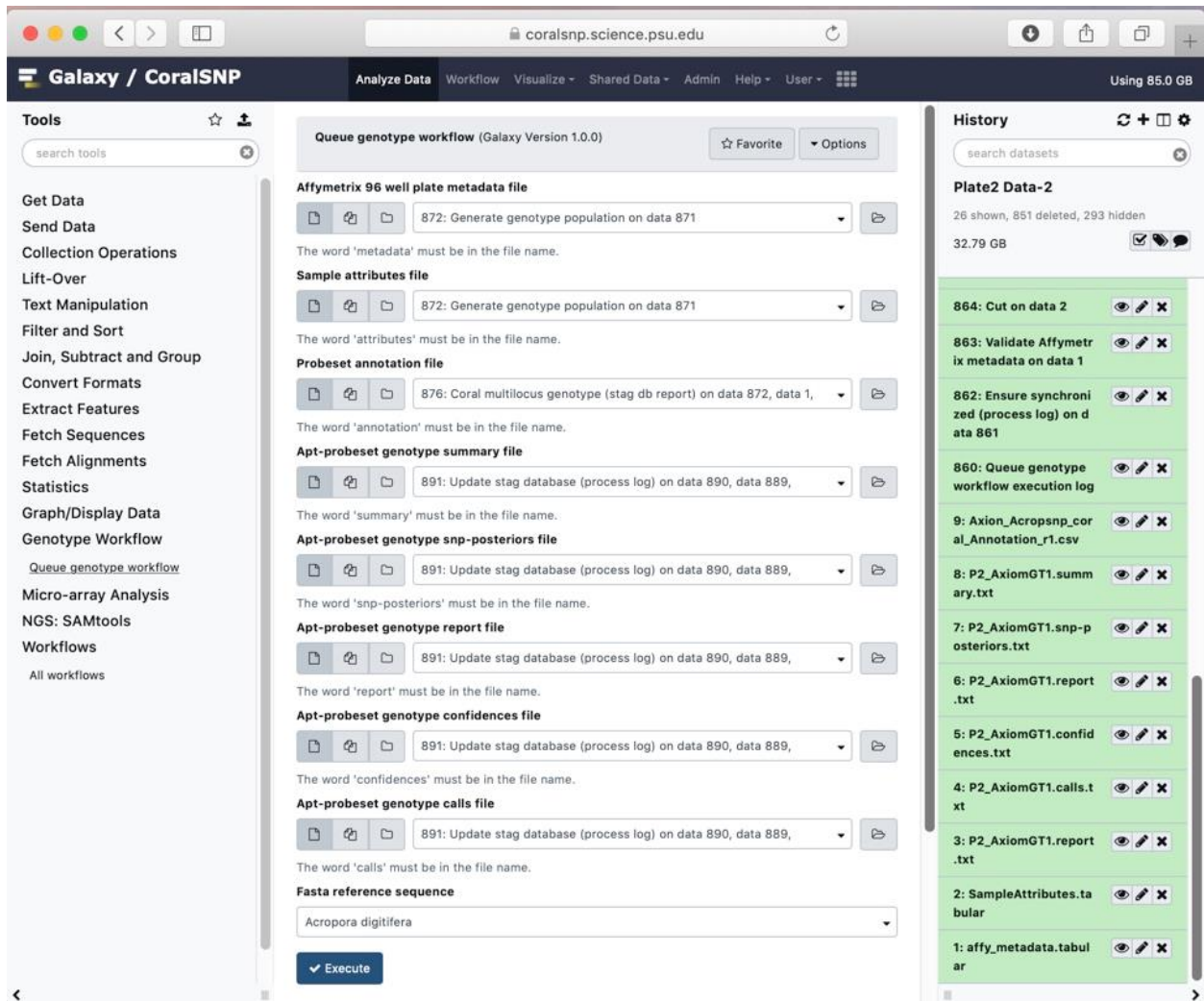
**Figure 6: The Galaxy / CoralSNP Queue Genotype Workflow Tool**

The Queue Genotype Workflow tool shields the complexity of the analysis from the user, and performs its functions via the Galaxy REST API to produce multi-locus genotypes for Caribbean or Pacific acroporids. The SNP array combined with the analysis allows for reliable, standardized identification of host genet and symbiont strains and serves as a template for the development of arrays for additional coral genera. This data can be used for downstream genomic analyses as well as restoration planning.