# PlantTribes: Galaxy tools for comparative gene family analysis in plant genomiCS

Eric Kenneth Wafula[1], Greg Von Kuster[1], Joshua P. Der[2], Loren Honaas[1], Saravanaraj Ayyampalayam[3], Norman Wickett[4], Jim Leebens-Mack[5] and Claude dePamphilis[1]

[1]Penn State University, University Park, PA, [2]California State University, Fullerton, Fullerton, CA, [3]The University of Georgia, Athens, GA,
[4]Chicago Botanic Garden, Glencoe, IL, [5]University of Georgia, Athens, GA
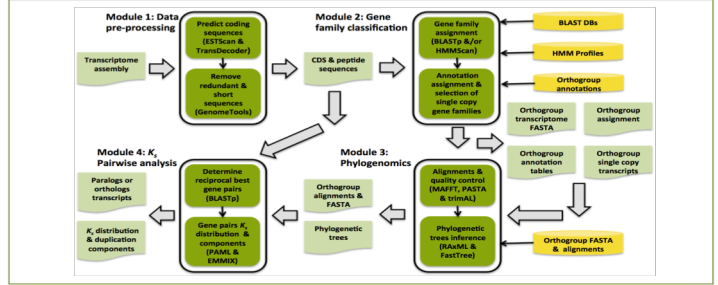
## Abstract

Galaxy PlantTribes is a collection of automated modular analysis pipelines that utilize objective classifications of complete protein sequences from sequenced genomes for comparative and evolutionary analyses of genome-scale gene families and transcriptomes. It post-processes de novo assembly transcripts into putative coding sequences and their corresponding amino acid translations, estimates paralogous/orthologous pairwise synonymous/non-synonymous substitution rates for a set of gene sequences, classifies gene sequences into pre-computed orthologous plant gene family clusters, and builds gene family multiple sequence alignments and their corresponding phylogenies. A user provides de novo assembly transcripts and Galaxy PlantTribes produces: (1) predicted coding sequences and their corresponding translations, (2) a table of pairwise synonymous/non-synonymous substitution rates for either orthologous or paralogous transcript pairs, (3) results of significant duplication components in the distribution of Ks (synonymous substitutions) values, (4) a summary table for transcripts classified into orthologous plant gene family clusters with their corresponding functional annotations, (4) gene family amino acid and nucleotide fasta sequences, (6) multiple sequence alignments, and (5) inferred maximum likelihood phylogenies. Optionally, a user can provide an external gene family scaffold and/or externally predicted coding sequences derived from a transcriptome assembly or gene predictions from a sequenced genome. Galaxy PlantTribes is freely available on the Galaxy main portal (https://usegalaxy.org). In addition, the standalone version of the pipeline is available for download on GitHub (https://github.com/dePamphilis/PlantTribes) as a command-line interface for batch processing of many datasets.

## Workflow



## Post Assembly QC



**AssemblyPostProcessor** post-processes de novo assembled transcripts into putative coding sequences and their corresponding amino acid translations and optionally assigns transcripts to circumscribed gene families ("orthogroups"). After transcripts have been assigned to gene families, overlapping contigs can be identified and merged to reduce fragmentation in the de novo assembly.

## Assembly Classification



**GeneFamilyClassifier** classifies gene coding sequences either produced by the **AssemblyPostProcessor** or from an external source into pre-computed orthologous gene family clusters (orthogroups) of a PlantTribes scaffold. Classified sequences are then assigned with the corresponding orthogroups metadata that includes gene counts of backbone taxa, super clusters ("super orthogroups") at multiple stringencies, and functional annotations.

## Alignment Estimation and QC



**GeneFamilyAligner** estimates protein and codon multiple sequence alignments of integrated orthologous gene family fasta files produced by the **GeneFamilyIntegrator** (not shown). The **GeneFamilyIntegrator** integrates PlantTribes scaffold orthogroup backbone gene models with gene coding sequences classified into the scaffold by the **GeneFamilyClassifier**.

## Alignment Visualization and Editing



Orthogroup fasta multiple sequence alignments produced by the **GeneFamilyAligner** can be visualized in Galaxy using the MSAVierwer plugin and manually edited with Jalview .

## Phylogenetic Inference



**GeneFamilyPhylogenyBuilder** performs gene family phylogenetic inference of multiple sequence alignments produced by the **GeneFamilyAligner.** Orthogroup maximum likelihood (ML) phylogenetic trees are inferred using either RAxML or FastTree algorithms. Trees are rooted using the most distant taxon present in the orthogroup backbone taxa if root rooting order list is not provided.

## Phylogenetic Visualization and Editing



Orthogroup newick phylogentic trees produced by the **GeneFamilyPhylogenyBuilder** can be visualized in Galaxy using either the Phylogenetic Tree Visualization plugin or the PHYLOViZ plugin. The Phylogenetic Tree Visualization plugin provides several options of tree rendering.

## Selection Analysis



**KaKsAnalysis** estimates paralogous and orthologous pairwise synonymous (Ks) and non-synonymous (Ka) substitution rates for a set of gene coding sequences either produced by the **AssemblyPostProcessor** or from an external source.

## Whole Genome Duplication Estimation



The resulting set of estimated paralogous and orthologous pairwise Ks values can be clustered into components using a mixture of multivariate normal distributions to identify significant duplication event(s) in a species or a pair of species. The **KsDistribution** (not shown) can be to plot the distribution of Ks rates and fit the estimated significant normal mixtures component(s) onto the distribution.

## Application

1. Targeted gene family assembly
2. Improving genome annotation quality
3. Transcriptome coverage and functional annotation
4. Species tree inference using single copy genes
5. Gene tree – species tree reconciliation
6. Timing of gene duplications and polyploidy
7. Ancestral gene content reconstruction
8. Gene family expansions/contractions
9. Timing of new gene function evolution
10. Studies of horizontal transfers among species

## Acknowledgement