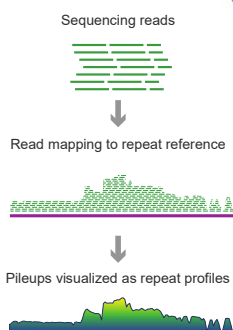


# RepeatProfiler: a pipeline for visualization and comparative analysis of repetitive DNA profiles

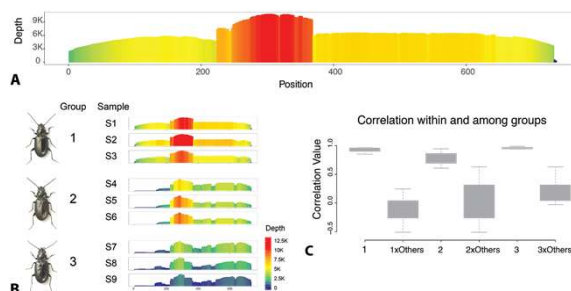
Sherif Negm, Anya Greenberg, Amanda Larracuenta & John Sproul  
University of Rochester, Department of Biology

## Overview

Repetitive DNA can reveal signals of evolutionary history over short time scales that may not be evident in sequences from slower-evolving genomic regions<sup>1</sup>. Many tools for studying repeats are directed toward organisms with existing genomic resources, including genome assemblies and repeat libraries. However, signals in repeat variation may prove especially valuable in disentangling evolutionary histories in diverse non-model groups, for which genomic resources are limited. Here we introduce RepeatProfiler, a tool for generating, visualizing, and comparing repetitive DNA profiles from low-coverage, short-read sequence data. RepeatProfiler facilitates comparative study of repeats across samples and repeats to provide a high-resolution data source for studies on species delimitation, genome evolution, and repeat biology over short evolutionary time scales.

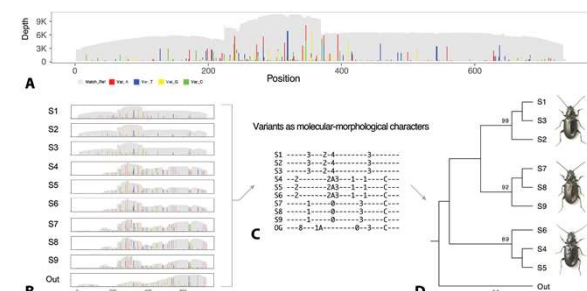


## Standard Output



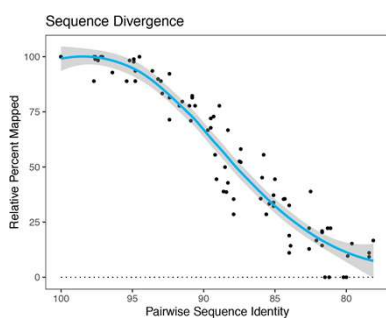
RepeatProfiler produces visually enhanced read-depth profiles (A) to simplify visual comparison of patterns across samples (B). The pipeline also enables statistical comparison of profile shape within and between user-defined sample groups. This correlation analysis measures the degree of similarity in the variable of coverage depth for each position across the reference sequence between two samples. We use Spearman's rank correlation coefficient to calculate correlation values for pairwise comparisons of all samples, and summarize within- and between-group comparisons with boxplots (C).

## Phylogenetic analysis of variant signatures



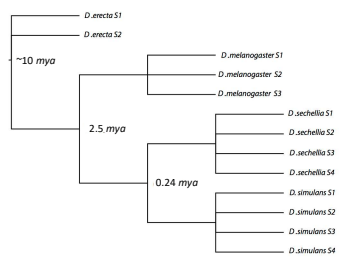
The pipeline also facilitates comparison of profiles across samples in a phylogenetic framework. The pipeline summarizes information contained in variant-enhanced profiles (A) by identifying abundant variants and encoding those variants in a phylip file as molecular-morphological characters (C) which can then be analyzed as morphological data using phylogenetic software. This example shows profiles of nine individual samples belonging to three *Bembidion* ground beetle species and one outgroup (B). The phylogenetic tree (D) correctly groups samples as species based on the variant in profiles (B, C).

## Validating RepeatProfiler

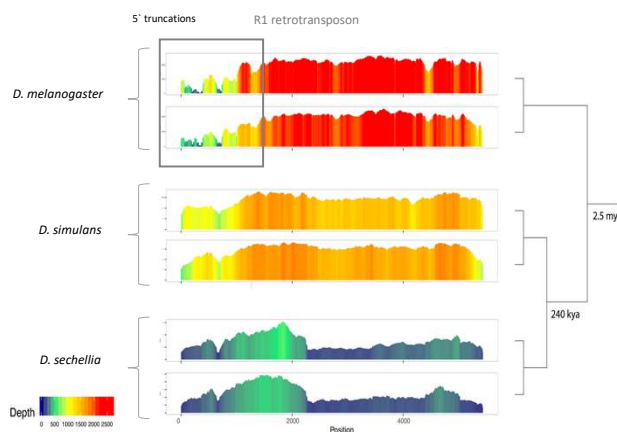


RepeatProfiler's default read mapping<sup>3</sup> settings allow for considerable divergence between the reference sequence and the reads being mapped, enabling comparative study across clades showing shallow to moderate sequence divergence among samples.

We validated our approach to comparing variant signatures through phylogenetic analysis by using RepeatProfiler to analyze 37 abundant transposable elements (TEs) in four closely related *Drosophila* species. We conducted phylogenetic analysis of phylip files summarizing variants for each TE as morphological data in IQ-Tree<sup>4</sup> and found that trees show good phylogenetic signal over short divergence times. This is a consensus tree of all 37 TE trees showing correct sample groupings and branching pattern among samples.<sup>2</sup>

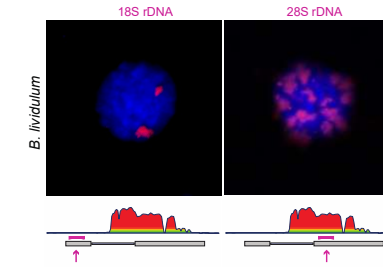


## Species Boundaries and Repeat Biology



Repetitive DNA profiles of two individuals for each of three *Drosophila* species showing strong species-specific signatures despite their recent evolutionary divergence (see tree on right with dates from<sup>2</sup>). In addition to providing evidence of species boundaries, repeat profiles can hold information related to repeat biology. This R1 element shows evidence of 5' truncations (i.e., reduced coverage in gray box) suggesting this element is recently active in this species, as active non-LTR elements are known to accumulate 5' truncations.

## RepeatProfiler and Genome Architecture



Fluorescence in situ hybridization shows that repeat profiles can provide evidence of genome-scale reorganization of repeat architecture over short time-scales. These FISH images show the recent mobilization or rDNA-like sequences that has led to genome-wide reorganization of repeats across in ground beetle species.<sup>1</sup>

## Features & Summary

- Produces visually enhanced repeat profiles across repeats, and samples
- Quantitative comparative analysis of profile features across samples
- High-resolution data source for species delimitation and genome evolution
- Provides insight's into of repeat biology and genome architecture
- Installation on Mac, Linux, and Windows using Homebrew and Docker
- Available from <https://github.com/johnssproul/RepeatProfiler>

Dependencies: R<sup>5</sup>, python, Bowtie2<sup>6</sup>, SAMtools<sup>7</sup>, & R packages (ggplot2<sup>8</sup>, ggpubr, scales, reshape2)