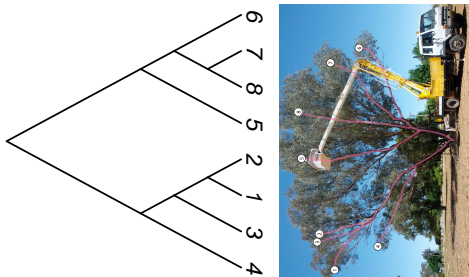


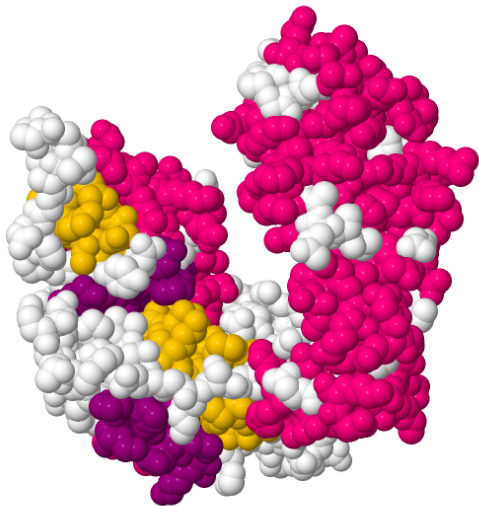
# Calculating a non-model *Eucalyptus* Individual's Somatic Mutation Rate

Adam Orr     @AdamJOrr

9/25/18



# Somatic Mutations Occur During Replication Even Without Exposure to Mutagens

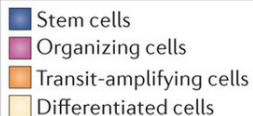


- DNA Polymerase  $\beta$
- Mutation rate  
 $\sim 10^{-9}$

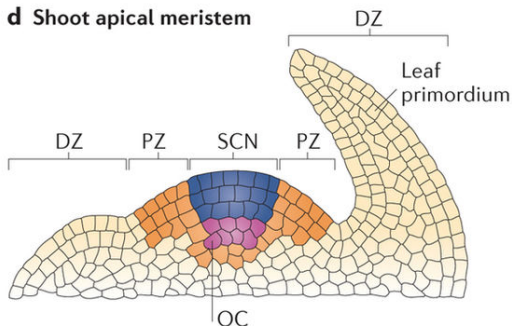
---

PDB: 7ICG

# Plants Grow Directionally



## d Shoot apical meristem



- The genetic structure of the plant *should* mirror its physical structure.

<sup>0</sup>Heidstra & Sabatini (2014) Plant and animal stem cells: similar yet different.

# Do Plants Evolve Differently?

NATURE VOL. 320 27 MARCH 1986

NEWS AND VIEWS

305

## Somatic mutation

# Do plants evolve differently?

from William J. Sutherland and Andrew R. Watkinson

It has recently been stressed that single ornamental plant with variegated leaves plants may consist of a mosaic of genetic clones. It is essential to remove branches

trees. This model also shows that somatic mutation may be an important source of variation within trees and tree populations. It seems likely that this mechanism is unimportant for species with relatively few meristems (such as peas and maize) but more important for long-lived and extensive species. Thus, individual clones of the aspen (*Populus tremuloides*) can cover

The relative importance of somatic and gametic mutations in plants cannot be assessed until the necessary measurements are made. But it is clear that somatic mutation could be important in many plant species.

Plants are weird.

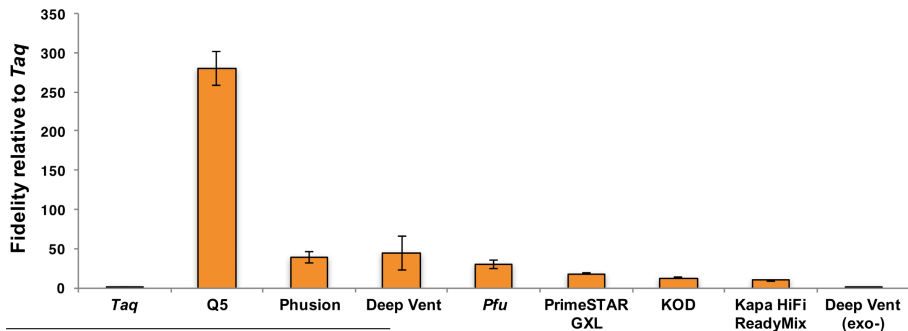
- Segregated germline?
- Dedifferentiation
- Selfing
- Mutational stress response?

# Why are somatic mutations difficult to detect?

Mutations are very rare, but sequencing errors are very common.

**Sequencing error** alone is  $\sim 10^{-2}$  while mutation rate after error-checking is  $\sim 10^{-9}$

- Errors accumulate during PCR prior to sequencing - then propagate.
- *Taq*  $\sim 10^{-4}$
- Technical error from sequencer



<sup>0</sup> Potapov V, Ong JL (2017) Examining Sources of Error in PCR by Single-Molecule Sequencing

# Why Care About Somatic Mutations?

## Disease

- Cancer

## Development

- Understanding the relationship between tissues

## Agriculture

- Looking for interesting phenotypes in clonally reproducing species

## Evolution

- Determining the relationship between somatic and germline mutation rate



---

<sup>0</sup> [https://commons.wikimedia.org/wiki/File:White\\_nectarine\\_and\\_cross\\_section02\\_edit.jpg](https://commons.wikimedia.org/wiki/File:White_nectarine_and_cross_section02_edit.jpg)

# Project Goals

## Is it possible?

Can we detect mutations with sufficient accuracy to reconstruct the physical structure of a tree?

## The Relevant Measurements

What's the mutation rate like? Is there evidence of hypermutation?

# A Genetic Mosaic



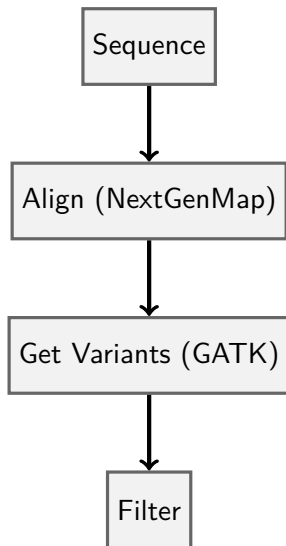
- Edwards identified as mosaic in 1993<sup>1</sup>
- Sheep pen in Yeoval, New South Wales
- Differential oil production gives protection from Christmas beetles

<sup>1</sup>Edwards PB, Wanjura WJ, Brown WV. *Oecologia* 1993, 95:551–557.



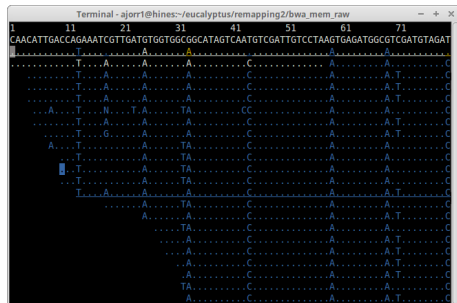
# Study Methodology

- Sequence 8 samples in triplicate
- ~10X coverage for each replicate
- Align sequence to genome of *Eucalyptus grandis*
- Use replicates to remove false positives



# Study Methodology

- Sequence 8 samples in triplicate
- ~10X coverage for each replicate
- Align sequence to genome of *Eucalyptus grandis*
- Use replicates to remove false positives

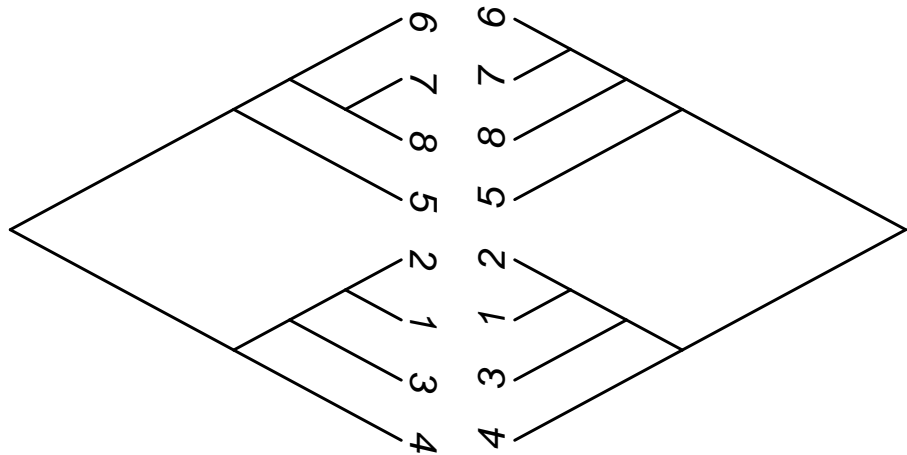


A terminal window titled "Terminal - ajorrt@hines:~/eucalyptus/remapping2/bwa\_mem\_raw" displays sequence alignment results. The top line shows a reference sequence with positions 1, 11, 21, 31, 41, 51, 61, and 71 marked. Below this, multiple lines of alignment data are shown, with columns corresponding to the reference positions. The alignment data consists of letters (A, C, G, T) and dots, representing matches and mismatches between the sequenced samples and the reference genome. The alignment is presented in a grid-like format, with each row representing a different sample or replicate.

# Mutation Pattern Approximately Matches Tree Structure

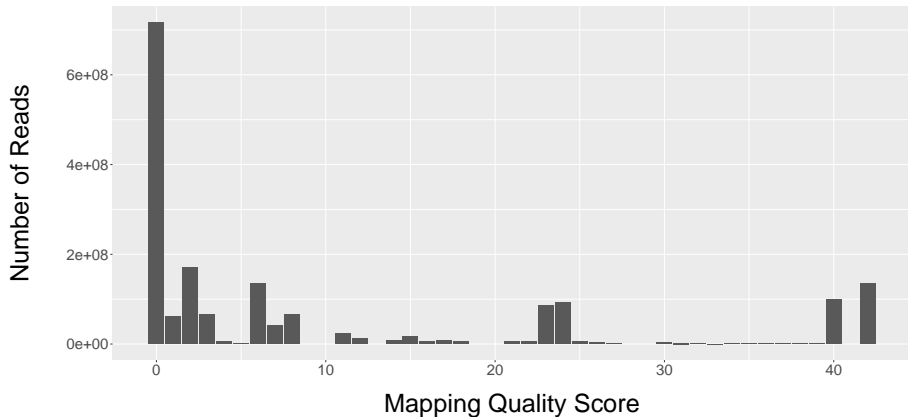
GATK Best Practices Tree

True Tree



# Most Reads Are Not Mapped to the *E. grandis* Reference

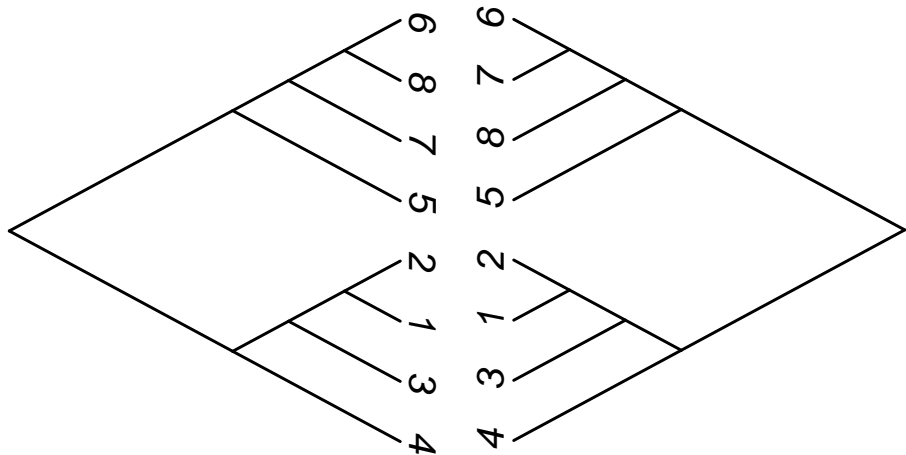
Quality of Reads Mapped



# A Reference-Free Method Performs Similarly

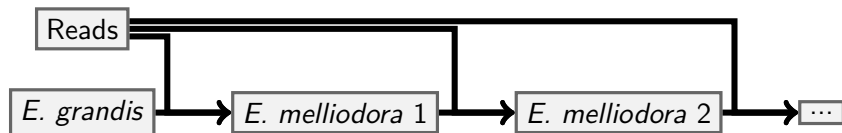
DiscoSNP++ Tree

True Tree

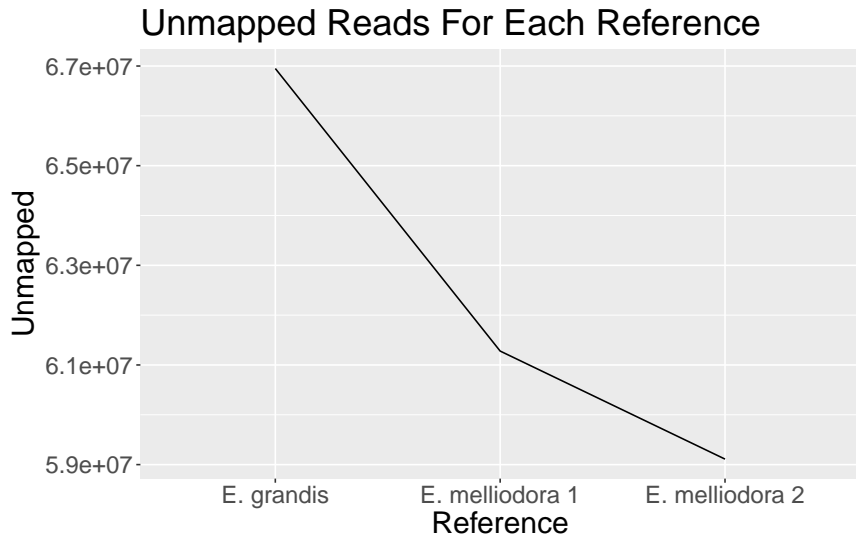


# Approximating a Genome

Use *E. melliodora* genome as a starting place, then generate a new reference and map to that reference.



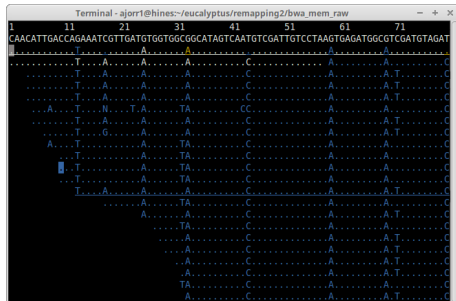
# Our New Reference Has Fewer Unmapped Reads



# Filtering Variants

Remove variants likely from alignment errors:

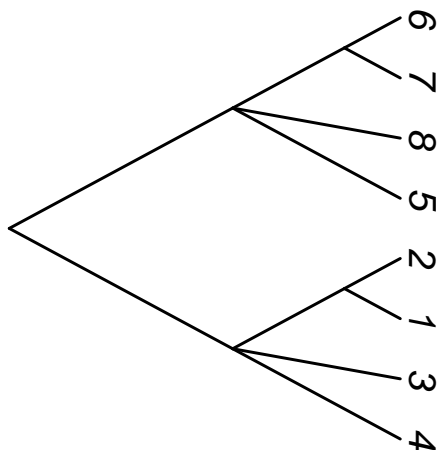
- at sites with excessive depth (>500).
- with excessive levels of heterozygosity.
- within 50 bases of an indel.
- in repeat regions



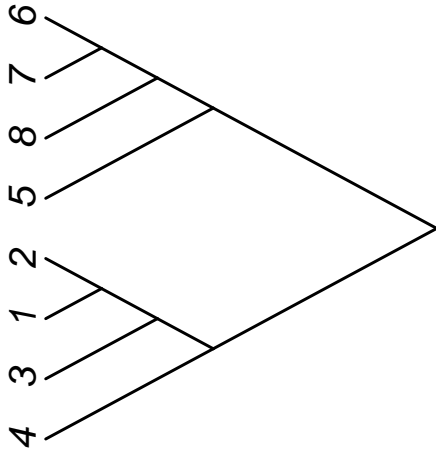


# Removing Variants in Repeat Regions Improves Tree Topology

Predicted Variants



True Tree



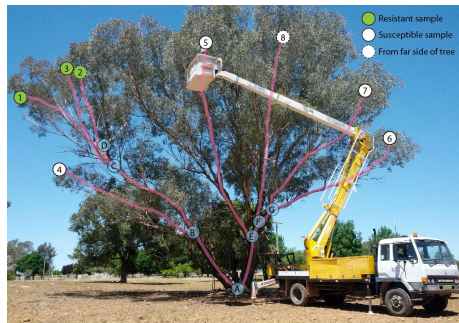
# Using Tree Topology Gives Higher Recall Rate

- *DeNovoGear* is a variant-calling method that uses information in the tree topology to call variants.
- By simulation, we introduced 14000 mutations on the tree

<i>GATK</i>	<i>DeNovoGear</i>
3859 mutations	4193 mutations
27%	30%


# Mutation Rates


- Detected 91 mutations.
- 20 mutations in genes.
- Estimated recall of  $\sim 30\%$ .
- $91 \times \frac{1}{3} = 303$  mutations.
- $2.5 \times 10^{-7}$  mutations per site per genome
- $\sim 1.4$  mutations per meter of length
- Somatic mutations account for  $\sim 55$  mutations per leaf tip.



# Acknowledgements

- Advisor: Reed Cartwright  @MinionLab
- Robert Lanfear, Australian National University  @RobLanfear

Pipeline:  <https://github.com/adamjorr/somatic-variation>

Talk:  <https://github.com/adamjorr/talks>



This work is supported by grants NIH R01-HG007178 and NSF DBI-1356548.