

Genome variation graphs with the vg toolkit

Jean Monlong

Updates from the GRC & GIAB
Oct 15, 2019



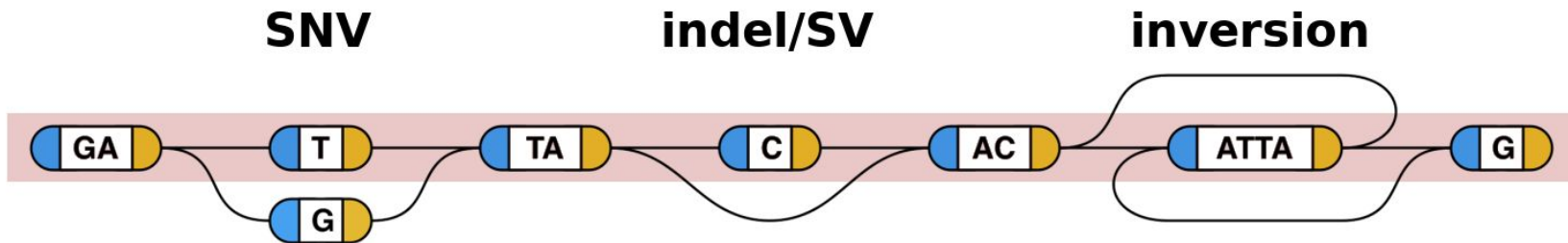
UNIVERSITY OF CALIFORNIA

SANTA CRUZ

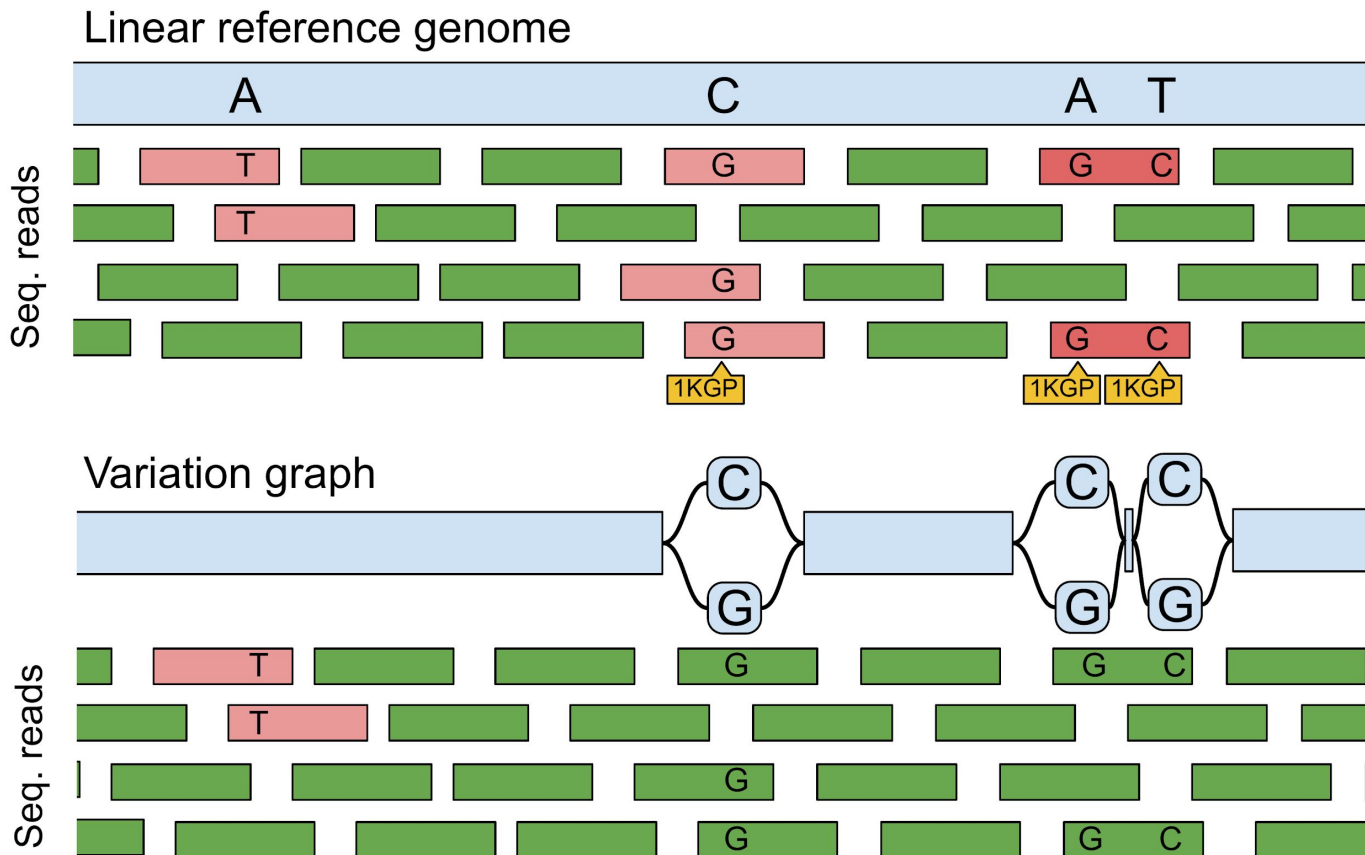
Genomics
Institute

Variation Graphs

An approach to incorporating information on human diversity into the genomic reference.



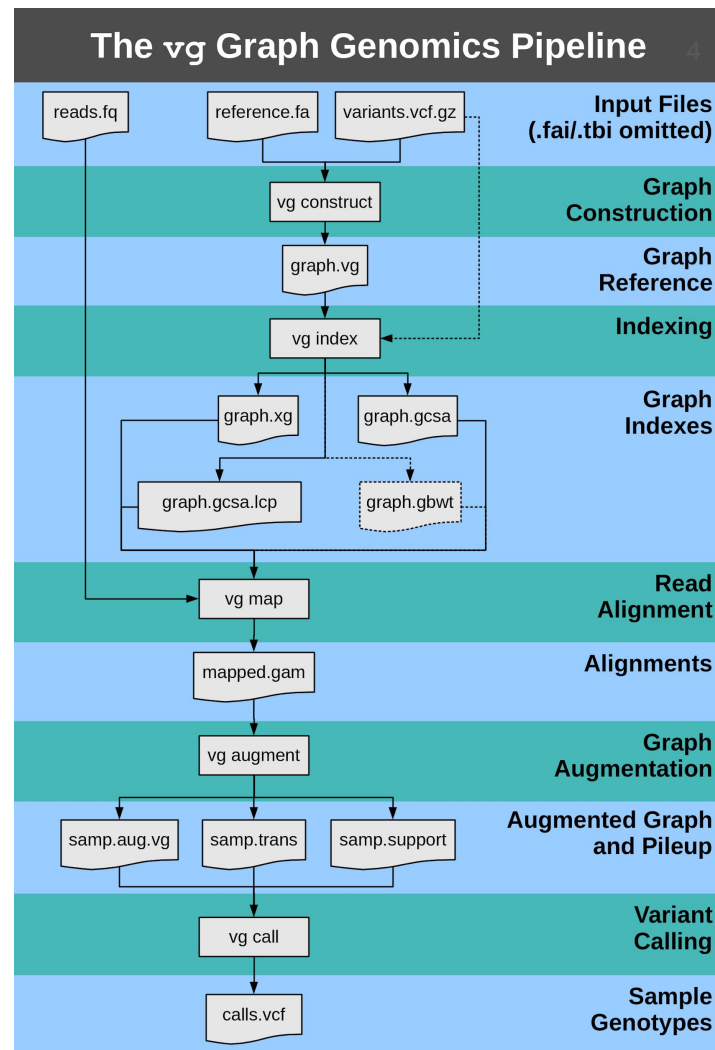
Sequencing reads map better on variation graphs



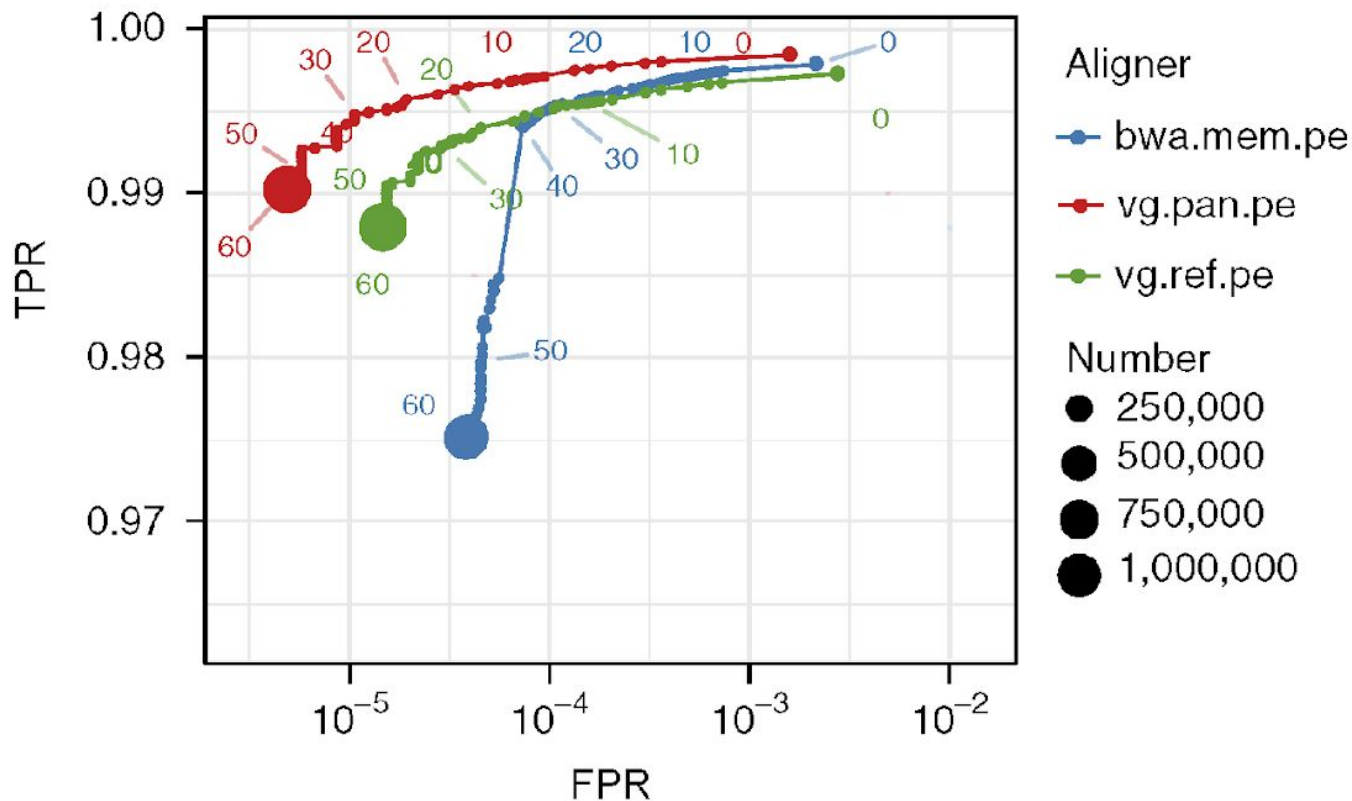


is a complete,
open source solution
for graph construction,
read mapping,
and variant calling.

<https://github.com/vgteam/vg>

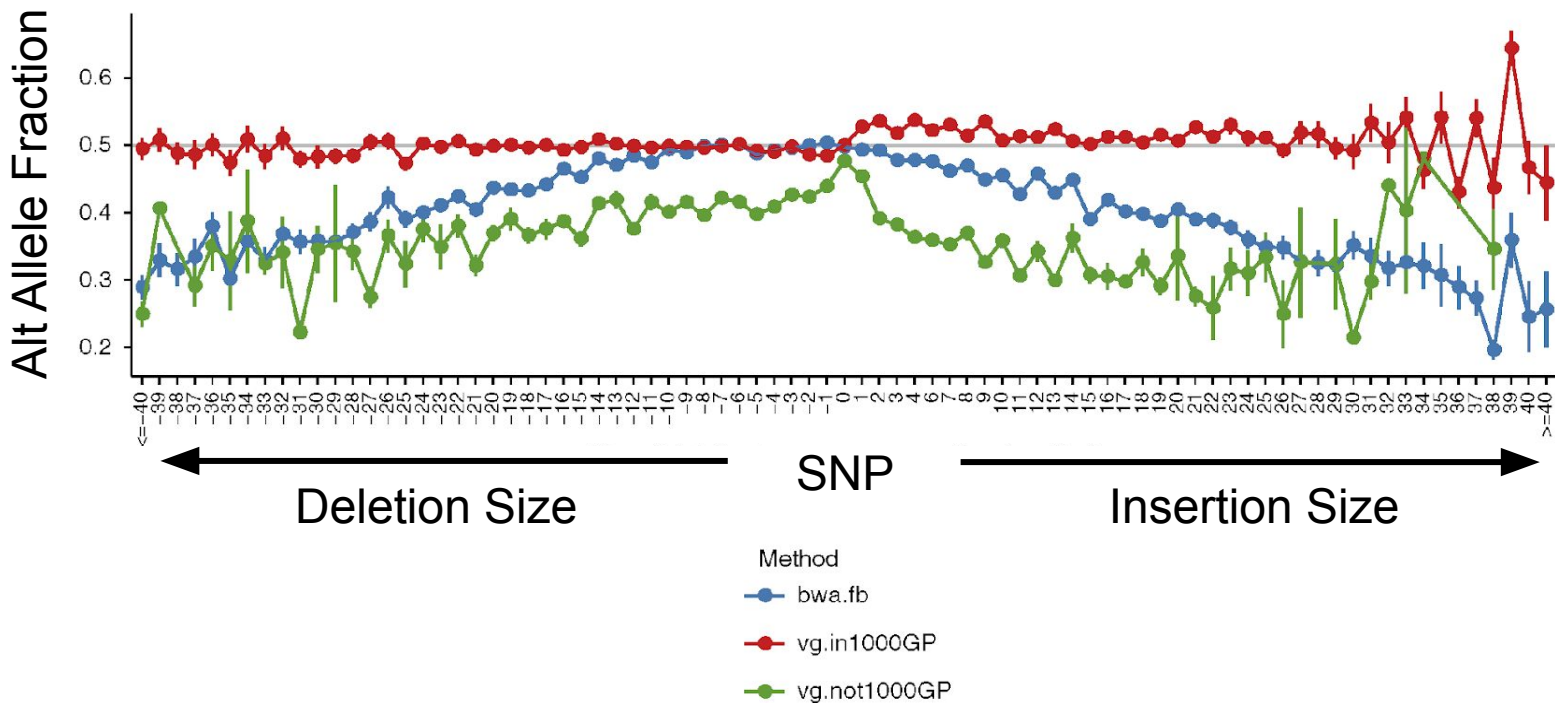


Better read mapping in regions with variants



Garrison et al (2018). Variation graph toolkit improves read mapping by representing genetic variation in the reference. Nature Biotechnology.

Better allele balances at heterozygous sites



Garrison et al (2018). Variation graph toolkit improves read mapping by representing genetic variation in the reference. Nature Biotechnology.

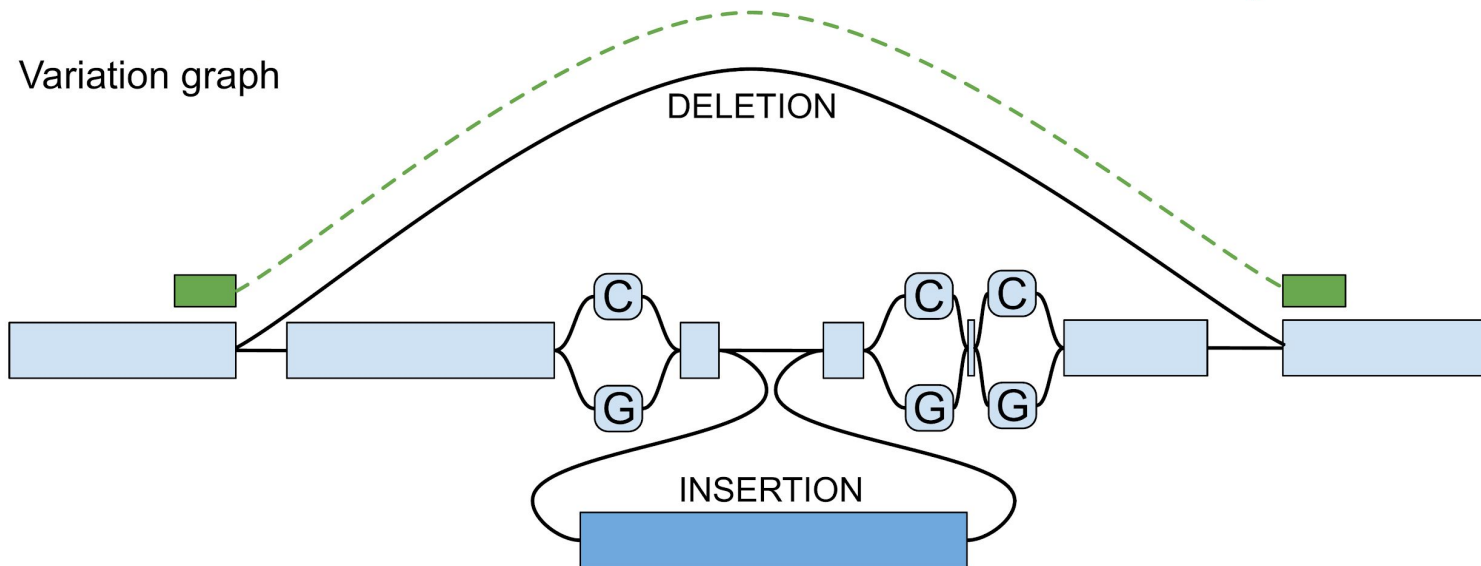
Structural Variants (SVs)

Genomic variants >50bp. E.g. insertions, deletions, inversions.

Linear reference genome



Variation graph

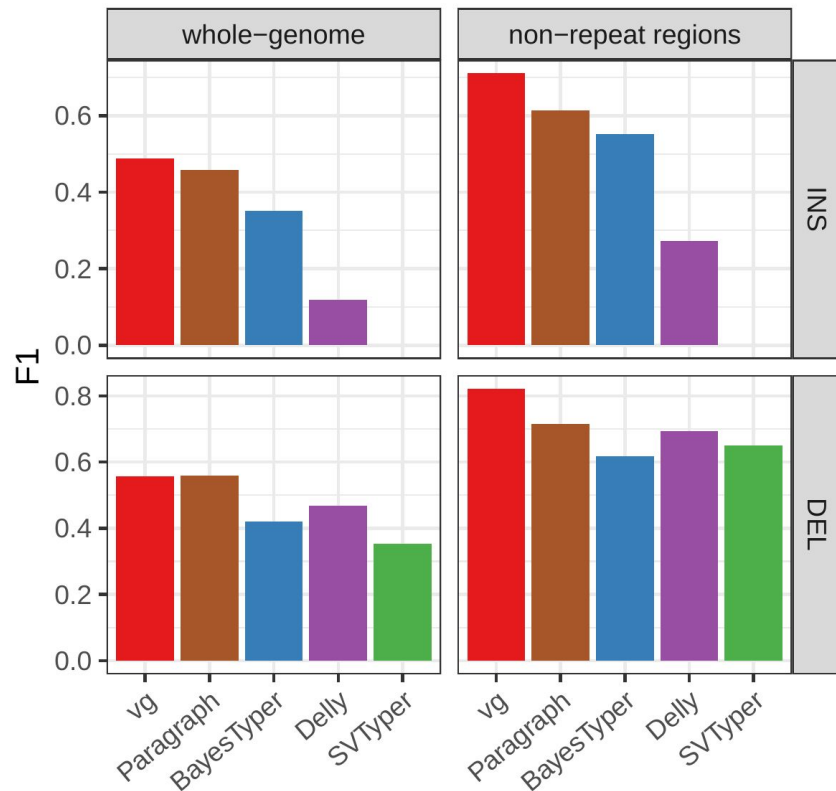


Genotyping SVs from short-read data

High-quality SV catalogs from long-read sequencing studies (HGSVC 2019, GIAB 2019, SVPOP 2019).

Graph-based SV genotypers:
vg, Paragraph, BayesTyper

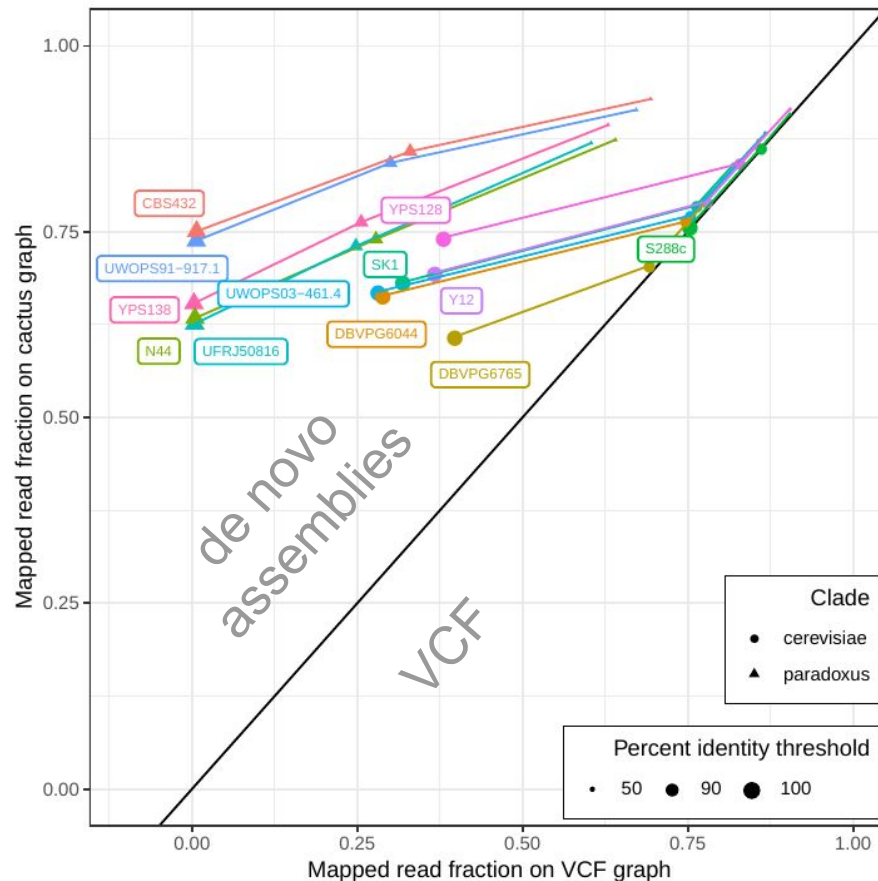
Traditional SV genotypers:
Delly, SVTyper



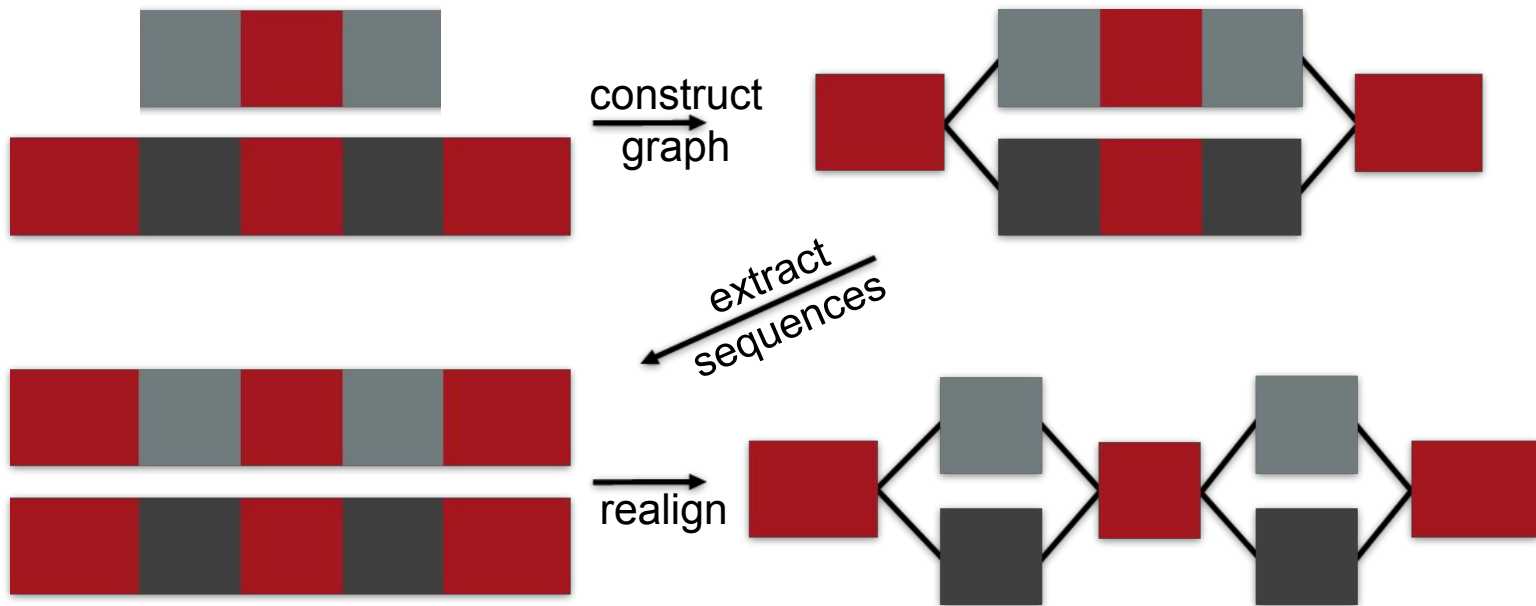
Graph from de novo assemblies

Experiment with 12 yeast strains.

- better read mapping.
- SV better supported by reads.

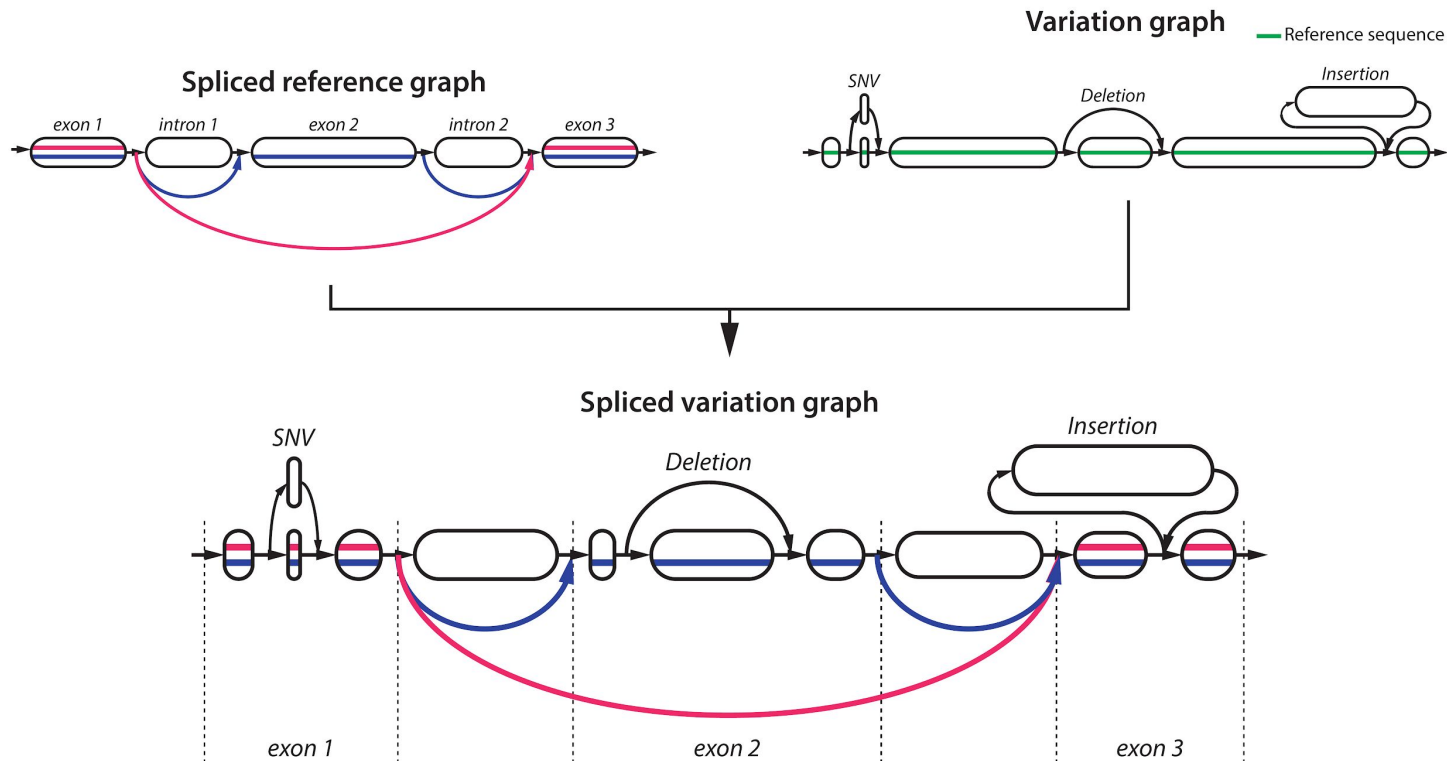


Variant normalization



Transcriptome analysis with variation graphs

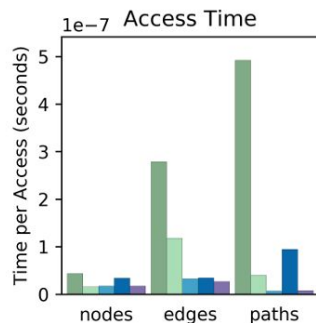
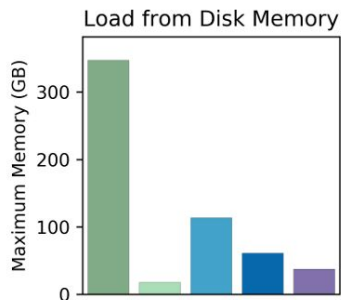
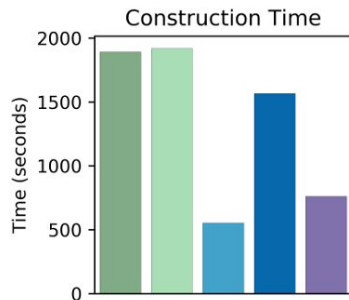
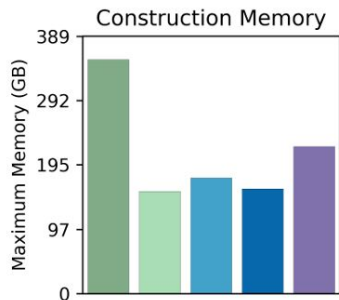
<https://github.com/jonassibbesen/rpvq>



New graph formats optimized for memory and speed

<https://github.com/vgteam/libbdsg>

vg pg hg og xg



VG: Legacy

PG (PackedGraph): Small

HG (HashGraph): Fast

OG (Optimized Dynamic Graph Index):
Balanced

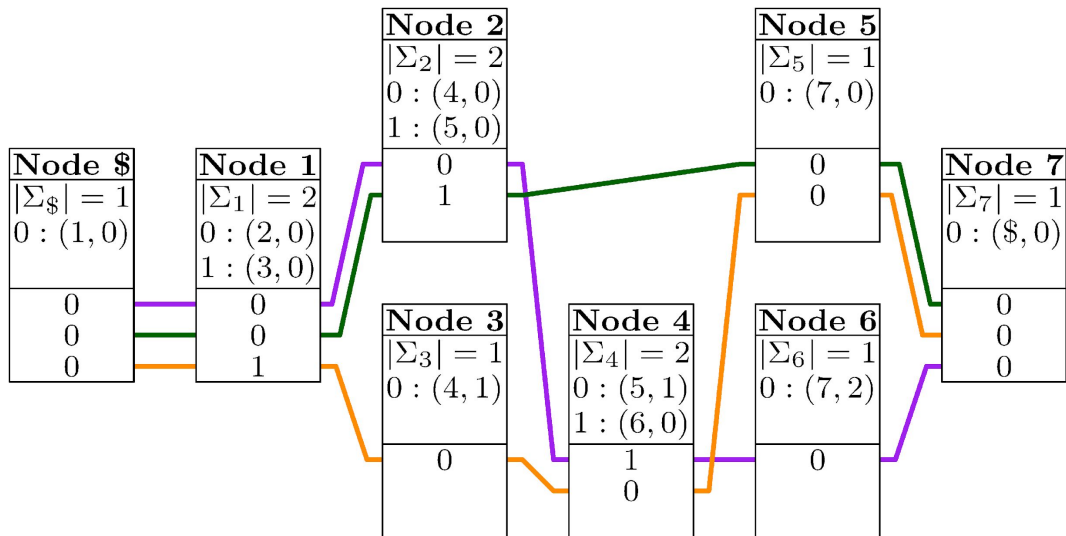
XG: Small and fast, but static

Graph Burrows-Wheeler Transform

<https://github.com/jltsiren/gbwtgraph>

Stores 2,504 haplotypes from 1K Genomes Project in 14.6 GiB

Can scale past tens of thousands of haplotypes



Faster short read mapping with giraffe

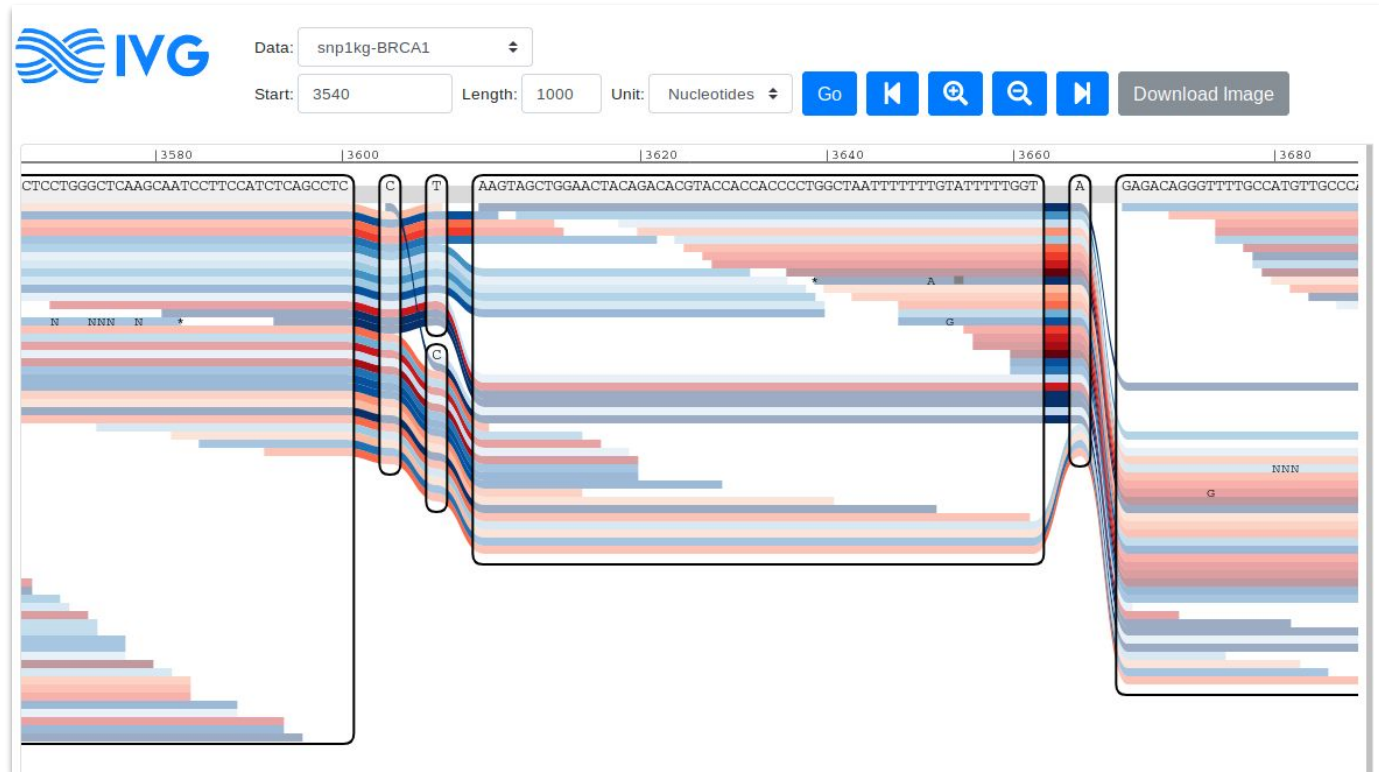
Assuming that most indels are in the graph and a large haplotype database.

- Minimizer-based indexing.
- Restrict to known haplotypes.
- Gapless extension.
- Faster clustering using a snarl tree.



Visualization

<https://github.com/vgteam/sequenceTubeMap>



Beyer et al. (2019). Sequence tube maps: making graph genomes intuitive to commuters. *Bioinformatics* (Oxford, England).

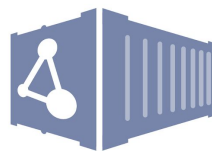
Pipelines in Toil and WDL

<https://github.com/vgteam/toil-vg>

https://github.com/vgteam/vg_wdl



	toil-vg	WDL
Graph construction	x	
Read mapping	x	x
Variant calling(SNV/indels & SVs)	x	x



Dockstore
Create, Share, Use



Acknowledgements

Benedict Paten	Jonas Sibbesen
Adam Novak	Xian Chang
Erik Garrison	Charles Markello
Jordan Eizenga	Yohei Rosen
Glenn Hickey	Robin Rounthwaite
Jouni Siren	Susanna Morin
David Heller	Emily Fujimoto



<https://github.com/vgteam/vg>

Eric Dawson
Mike Lin
Wolfgang Beyer

Richard Durbin
Daniel Zerbino



Check out Charles Markello' Platform Talk (PgmNr 15) on applying vg for rare variant discovery!