

Genotyping structural variation in variation graphs with the vg toolkit

This manuscript ([permalink](#)) was automatically generated from [jmonlong/manu-vgs@2e0caef](#) on January 22, 2019.

Authors

 Glenn Hickey^{1, },  David Heller^{1, },  Jean Monlong^{1, },  Benedict Paten^{1, †}

 — These authors contributed equally to this work

† — To whom correspondence should be addressed: bpaten@ucsc.edu

1. UC Santa Cruz Genomics Institute, University of California, Santa Cruz, California, USA

Abstract

Introduction

Structural variation (SV) represents genomic mutation involving 50 bp or more and can take several forms, such as for example deletions, insertions, inversions, or translocations. Although whole-genome sequencing (WGS) made it possible to assess virtually any type of structural variation, many challenges remain. In particular, SV-supporting reads are difficult to map to reference genomes. Multi-mapping, caused by widespread repeated sequences in the genome, is another issue because it often resembles SV-supporting signal. As a result, many SV detection algorithms have been developed and multiple methods must usually be combined to minimize false positives. Several large-scale projects used this ensemble approach, cataloging tens of thousands of SV in humans[1,2]. SV detection from short-read sequencing remains laborious and of lower accuracy, explaining why these variants and their impact have been under-studied as compared to single-nucleotide variants (SNVs) and small insertions/deletions (indels).

Over the last few years, exciting developments in sequencing technologies and library preparation made it possible to produce long reads or retrieve long-range information over kilobases of sequence. These approaches are maturing to the point where it is feasible to analyze the human genome. This multi-kbp information is particularly useful for SV detection and de novo assembly. In the last few years, several studies using long-read or linked-read sequencing have produced large catalogs of structural variation, the majority of which were novel and sequence-resolved[3,4,5,6] (*REF_PETER_SOON*). These technologies are also enabling high-quality de novo genome assemblies to be produced[3,7], as well as large blocks of haplotype-resolved sequences[8]. These technological advances promise to expand the amount of known genomic variation in humans in the near future.

In parallel, the reference genome is evolving from a linear reference to a graph-based reference that contains known genomic variation[10,11,9]. By having variants in the graph, mapping rates are increased and variants are more uniformly covered, including indels and variants in complex regions[10]. Both the mapping and variant calling become variant-aware and benefit in terms of accuracy and sensitivity. In addition, different variant types are called simultaneously by a unified framework. Graphs have also been used locally, i.e. to call variants at the region level.

GraphTyper[12] and BayesTyper[13] both construct variation graphs of small regions and use them for variant genotyping. Here again, the graph-approach showed clear advantages over standard approaches that use the linear reference. Other SV genotyping approaches compare read mapping in the reference genome and a sequence modified with the SV. For example SMRT-SV was designed to genotype SVs identified on PacBio reads[4], SVTyper uses paired-end mapping and split-read mapping information[14], and Delly provides a genotyping feature in addition to its discovery mode[15].

Results

Methods

Discussion

References

1. An integrated map of structural variation in 2,504 human genomes

Peter H. Sudmant, Tobias Rausch, Eugene J. Gardner, Robert E. Handsaker, Alexej Abyzov, John Huddleston, Yan Zhang, Kai Ye, Goo Jun, ... Jan O. Korbel

Nature (2015-10) <https://doi.org/73c>

DOI: [10.1038/nature15394](https://doi.org/10.1038/nature15394) · PMID: [26432246](https://pubmed.ncbi.nlm.nih.gov/26432246/) · PMCID: [PMC4617611](https://pubmed.ncbi.nlm.nih.gov/PMC4617611/)

2. Whole-genome sequence variation, population structure and demographic history of the Dutch population

Laurent C Francioli, Androniki Menelaou, Sara L Pulit, Freerk van Dijk, Pier Francesco Palamara, Clara C Elbers, Pieter BT Neerincx, Kai Ye, Victor Guryev, ... Cisca Wijmenga

Nature Genetics (2014-06-29) <https://doi.org/f6bxm8>

DOI: [10.1038/ng.3021](https://doi.org/10.1038/ng.3021) · PMID: [24974849](https://pubmed.ncbi.nlm.nih.gov/24974849/)

3. Resolving the complexity of the human genome using single-molecule sequencing

Mark J. P. Chaisson, John Huddleston, Megan Y. Dennis, Peter H. Sudmant, Maika Malig, Fereydoun Hormozdiari, Francesca Antonacci, Urvashi Surti, Richard Sandstrom, Matthew Boitano, ... Evan E. Eichler

Nature (2014-11-10) <https://doi.org/w69>

DOI: [10.1038/nature13907](https://doi.org/10.1038/nature13907) · PMID: [25383537](https://pubmed.ncbi.nlm.nih.gov/25383537/) · PMCID: [PMC4317254](https://pubmed.ncbi.nlm.nih.gov/PMC4317254/)

4. Discovery and genotyping of structural variation from long-read haploid genome sequence data

John Huddleston, Mark J.P. Chaisson, Karyn Meltz Steinberg, Wes Warren, Kendra Hoekzema, David Gordon, Tina A. Graves-Lindsay, Katherine M. Munson, Zev N. Kronenberg, Laura Vives, ... Evan E. Eichler

Genome Research (2016-11-28) <https://doi.org/f9x79h>

DOI: [10.1101/gr.214007.116](https://doi.org/10.1101/gr.214007.116) · PMID: [27895111](https://pubmed.ncbi.nlm.nih.gov/27895111/) · PMCID: [PMC5411763](https://pubmed.ncbi.nlm.nih.gov/PMC5411763/)

5. Mapping and phasing of structural variation in patient genomes using nanopore sequencing

Mircea Cretu Stancu, Markus J. van Roosmalen, Ivo Renkens, Marleen M. Nieboer, Sjors Middelkamp, Joep de Ligt, Giulia Pregno, Daniela Giachino, Giorgia Mandrile, Jose Espejo Valle-Inclan, ... Wigard P. Kloosterman

Nature Communications (2017-11-06) <https://doi.org/gftpt9>

DOI: [10.1038/s41467-017-01343-4](https://doi.org/10.1038/s41467-017-01343-4) · PMID: [29109544](https://pubmed.ncbi.nlm.nih.gov/29109544/) · PMCID: [PMC5673902](https://pubmed.ncbi.nlm.nih.gov/PMC5673902/)

6. Genome-wide reconstruction of complex structural variants using read clouds

Noah Spies, Ziming Weng, Alex Bishara, Jennifer McDaniel, David Catoe, Justin M Zook, Marc Salit, Robert B West, Serafim Batzoglou, Arend Sidow

Nature Methods (2017-07-17) <https://doi.org/gbnhww>

DOI: [10.1038/nmeth.4366](https://doi.org/10.1038/nmeth.4366) · PMID: [28714986](https://pubmed.ncbi.nlm.nih.gov/28714986/) · PMCID: [PMC5578891](https://pubmed.ncbi.nlm.nih.gov/PMC5578891/)

7. Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain, Sergey Koren, Karen H Miga, Josh Quick, Arthur C Rand, Thomas A Sasani, John R Tyson, Andrew D Beggs, Alexander T Dilthey, Ian T Fiddes, ... Matthew Loose

Nature Biotechnology (2018-01-29) <https://doi.org/gczffw>

DOI: [10.1038/nbt.4060](https://doi.org/10.1038/nbt.4060) · PMID: [29431738](https://pubmed.ncbi.nlm.nih.gov/29431738/) · PMCID: [PMC5889714](https://pubmed.ncbi.nlm.nih.gov/PMC5889714/)

8. Phased diploid genome assembly with single-molecule real-time sequencing

Chen-Shan Chin, Paul Peluso, Fritz J Sedlazeck, Maria Nattestad, Gregory T Concepcion, Alicia Clum, Christopher Dunn, Ronan O'Malley, Rosa Figueroa-Balderas, Abraham Morales-Cruz, ... Michael C Schatz

Nature Methods (2016-10-17) <https://doi.org/f9fv4w>

DOI: [10.1038/nmeth.4035](https://doi.org/10.1038/nmeth.4035) · PMID: [27749838](https://pubmed.ncbi.nlm.nih.gov/27749838/) · PMCID: [PMC5503144](https://pubmed.ncbi.nlm.nih.gov/PMC5503144/)

9. Genome graphs and the evolution of genome inference

Benedict Paten, Adam M. Novak, Jordan M. Eizenga, Erik Garrison

Genome Research (2017-03-30) <https://doi.org/f95nhd>

DOI: [10.1101/gr.214155.116](https://doi.org/10.1101/gr.214155.116) · PMID: [28360232](https://pubmed.ncbi.nlm.nih.gov/28360232/) · PMCID: [PMC5411762](https://pubmed.ncbi.nlm.nih.gov/PMC5411762/)

10. Variation graph toolkit improves read mapping by representing genetic variation in the reference

Erik Garrison, Jouni Sirén, Adam M Novak, Glenn Hickey, Jordan M Eizenga, Eric T Dawson, William Jones, Shilpa Garg, Charles Markello, Michael F Lin, ... Richard Durbin

Nature Biotechnology (2018-08-20) <https://doi.org/gd2zqs>

DOI: [10.1038/nbt.4227](https://doi.org/10.1038/nbt.4227) · PMID: [30125266](https://pubmed.ncbi.nlm.nih.gov/30125266/) · PMCID: [PMC6126949](https://pubmed.ncbi.nlm.nih.gov/PMC6126949/)

11. Fast and accurate genomic analyses using genome graphs

Goran Rakocevic, Vladimir Semenyuk, Wan-Ping Lee, James Spencer, John Browning, Ivan J. Johnson, Vladan Arsenijevic, Jelena Nadj, Kaushik Ghose, Maria C. Suci, ... Deniz Kural

Nature Genetics (2019-01-14) <https://doi.org/gftd46>

DOI: [10.1038/s41588-018-0316-4](https://doi.org/10.1038/s41588-018-0316-4) · PMID: [30643257](https://pubmed.ncbi.nlm.nih.gov/30643257/)

12. GraphTyper enables population-scale genotyping using pangenome graphs

Hannes P Eggertsson, Hakon Jonsson, Snaedis Kristmundsdottir, Eiríkur Hjartarson, Birte Kehr, Gisli Masson, Florian Zink, Kristján E Hjorleifsson, Aslaug Jonasdottir, Adalbjörg Jonasdottir, ... Bjarni V Halldorsson

Nature Genetics (2017-09-25) <https://doi.org/gbx7v6>

DOI: [10.1038/ng.3964](https://doi.org/10.1038/ng.3964) · PMID: [28945251](https://pubmed.ncbi.nlm.nih.gov/28945251/)

13. Accurate genotyping across variant classes and lengths using variant graphs

Jonas Andreas SibbesenLasse Maretty, Anders Krogh

Nature Genetics (2018-06-18) <https://doi.org/gdndnz>

DOI: [10.1038/s41588-018-0145-5](https://doi.org/10.1038/s41588-018-0145-5) · PMID: [29915429](https://pubmed.ncbi.nlm.nih.gov/29915429/)

14. SpeedSeq: ultra-fast personal genome analysis and interpretation

Colby Chiang, Ryan M Layer, Gregory G Faust, Michael R Lindberg, David B Rose, Erik P Garrison, Gabor T Marth, Aaron R Quinlan, Ira M Hall

Nature Methods (2015-08-10) <https://doi.org/gcpgfh>

DOI: [10.1038/nmeth.3505](https://doi.org/10.1038/nmeth.3505) · PMID: [26258291](https://pubmed.ncbi.nlm.nih.gov/26258291/) · PMCID: [PMC4589466](https://pubmed.ncbi.nlm.nih.gov/PMC4589466/)

15. DELLY: structural variant discovery by integrated paired-end and split-read analysis

T. Rausch, T. Zichner, A. Schlattl, A. M. Stutz, V. Benes, J. O. Korb

Bioinformatics (2012-09-07) <https://doi.org/f38r2c>

DOI: [10.1093/bioinformatics/bts378](https://doi.org/10.1093/bioinformatics/bts378) · PMID: [22962449](https://pubmed.ncbi.nlm.nih.gov/22962449/) · PMCID: [PMC3436805](https://pubmed.ncbi.nlm.nih.gov/PMC3436805/)