

Methodologies for Structural Variant detection

Fritz Sedlazeck & Luis Paulin

Dec,13, 2023



@sedlazeck



RICE

Recap from yesterday

Long read sequencing Structural Variants

REVIEWS |

- More comprehensive
 - LR: 20-23k SV / human
 - SR: ~10-12k SV / human
- Access in repetitive regions
 - 193 medical genes (Mandelker 2019)
 - 386 medically relevant genes (Wagner 2022)
 - Centromere, telomeres (e.g. T2T)
- Assembly/ Phasing
 - N50, no gaps, phased, etc.

COMPUTATIONAL TOOLS

Piercing the dark matter: bioinformatics of long-range sequencing and mapping

Fritz J. Sedlazeck¹, Hayan Lee¹, Charlotte A. Darby¹ and Michael C. Schatz^{3,4*}

Abstract | Several new genomics technologies have become available that offer long-read sequencing or long-range mapping with higher throughput and higher resolution analysis

nature methods

ARTICLE

<https://doi.org/10.1038/s41592-022-01750-w>

Accurate detection of complex structural variations using single-molecule sequencing

Fritz J. Sedlazeck^{1,6*}, Philipp Rescheneder^{2,6}, Moritz Smolka¹, Han Fang³, Maria Nattestad^{1,3}, Arndt von Haeseler^{2,4} and Michael C. Schatz^{1,3,5*}

REVIEW

Open Access

Structural variant calling: the long and the short of it



Medhat Mahmoud^{1†}, Nastassia Gobet^{2,3†}, Diana Ivette Cruz-Dávalos^{3,4}, Ninon Mounier^{3,5}, Christophe Dessimoz^{2,3,4,6,7*} and Fritz J. Sedlazeck¹

Sniffles: HG002 v0.6 and CMRG



A robust benchmark for detection of germline large deletions and insertions

Justin M. Zook^{①,2}, Nancy F. Hansen^②, Nathan D. Olson^①, Lesley Chapman¹, James C. Mullikin^②, Chunlin Xiao³, Stephen Sherry³, Sergey Koren^②, Adam M. Phillippy^②, Paul C. Boutros^④, Sayed Mohammad E. Sahraeian⁵, Vincent Huang^⑥, Alexandre Rouette⁷, Noah Alexander⁸, Christopher E. Mason^{⑨,10,11,12}, Iman Hajirasouliha⁹, Camir Ricketts⁹, Joyce Lee^⑬, Rick Tearle¹⁴, Ian T. Fiddes¹⁵, Alvaro Martinez Barrio^⑯, Jeremiah Wala¹⁶, Andrew Carroll¹⁷, Noushin Ghaffari¹⁸, Oscar L. Rodriguez¹⁹, Ali Bashir¹⁹, Shaun Jackman²⁰, John J. Farrell²¹, Aaron M. Wenger²², Can Alkan^{⑩,23}, Arda Soylev^{⑯,24}, Michael C. Schatz²⁵, Shilpa Garg²⁶, George Church^{⑯,26}, Tobias Marschall^{⑯,27}, Ken Chen^{⑯,28}, Xian Fan²⁹, Adam C. English³⁰, Jeffrey A. Rosenfeld^{⑯,31,32}, Weichen Zhou^{⑯,33}, Ryan E. Mills³³, Jay M. Sage³⁴, Jennifer R. Davis³⁴, Michael D. Kaiser³⁴, John S. Oliver³⁴, Anthony P. Catalano³⁴, Mark J. P. Chaisson³⁵, Noah Spies³⁶, Fritz J. Sedlazeck^{⑯,37} and Marc Salit³⁶



Curated variation benchmarks for challenging medically relevant autosomal genes

Justin Wagner¹, Nathan D. Olson^①, Lindsay Harris¹, Jennifer McDaniel¹, Haoyu Cheng^②, Arkarachai Fungtammasan³, Yih-Chii Hwang³, Richa Gupta^③, Aaron M. Wenger⁴, William J. Rowell^④, Ziad M. Khan^⑤, Jesse Farek⁵, Yiming Zhu⁵, Aishwarya Pisupati^⑤, Medhat Mahmoud^⑤, Chunlin Xiao⁶, Byunggil Yoo⁷, Sayed Mohammad Ebrahim Sahraeian⁸, Danny E. Miller^{⑨,10}, David Jáspez^⑪, José M. Lorenzo-Salazar^⑪, Adrián Muñoz-Barrera^⑪, Luis A. Rubio-Rodríguez^⑫, Carlos Flores^{⑬,12,13}, Giuseppe Narzisi^⑭, Uday Shanker Evanji^⑭, Wayne E. Clarke^⑮, Joyce Lee^⑯, Christopher E. Mason^⑯, Stephen E. Lincoln¹⁷, Karen H. Miga^{⑯,18}, Mark T. W. Ebbert^{⑯,20,21}, Alaina Shumate^{22,23}, Heng Li², Chen-Shan Chin^{⑯,3,24}, Justin M. Zook^{⑯,1,24} and Fritz J. Sedlazeck^{⑯,5,24}

Parameters: Default Tuned Coverage ★ 05 □ 10 ○ 20 △ 30 ◇ 50 ■ GT-F1 < 0.5

HG002 ONT GIAB Tier 1



E

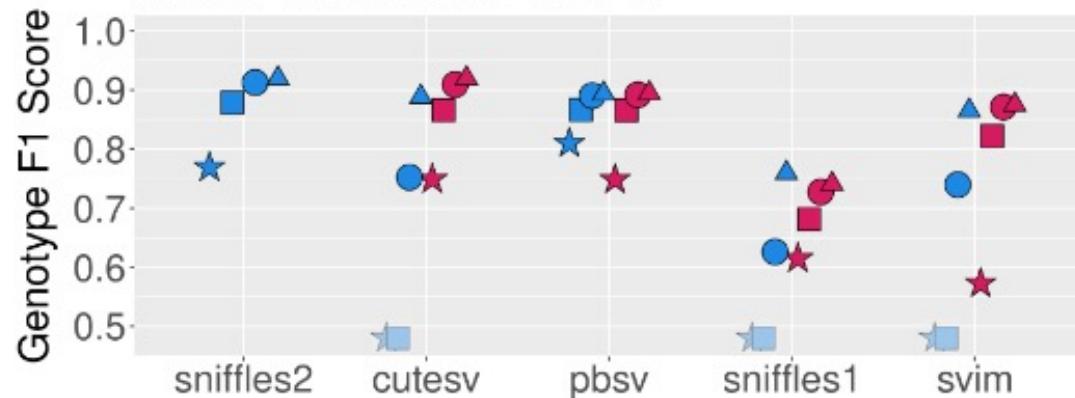
HG002 ONT CMRG



Sniffles: HG002 v0.6 and CMRG

B

HG002 HIFI GIAB Tier 1



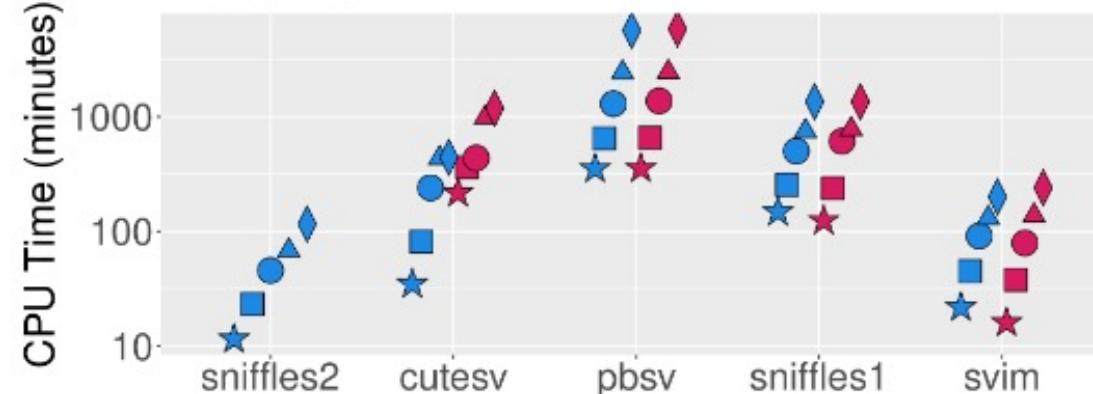
Parameters: Default Tuned Coverage 05 Coverage 10 Coverage 20 Coverage 30 Coverage 50 GT-F1 < 0.5

HG002 ONT GIAB Tier 1



C

HG002 ONT GIAB Tier 1



D

HG002 ONT CMRG



New Applications in SV detection

1. Germline SV
2. Population scale

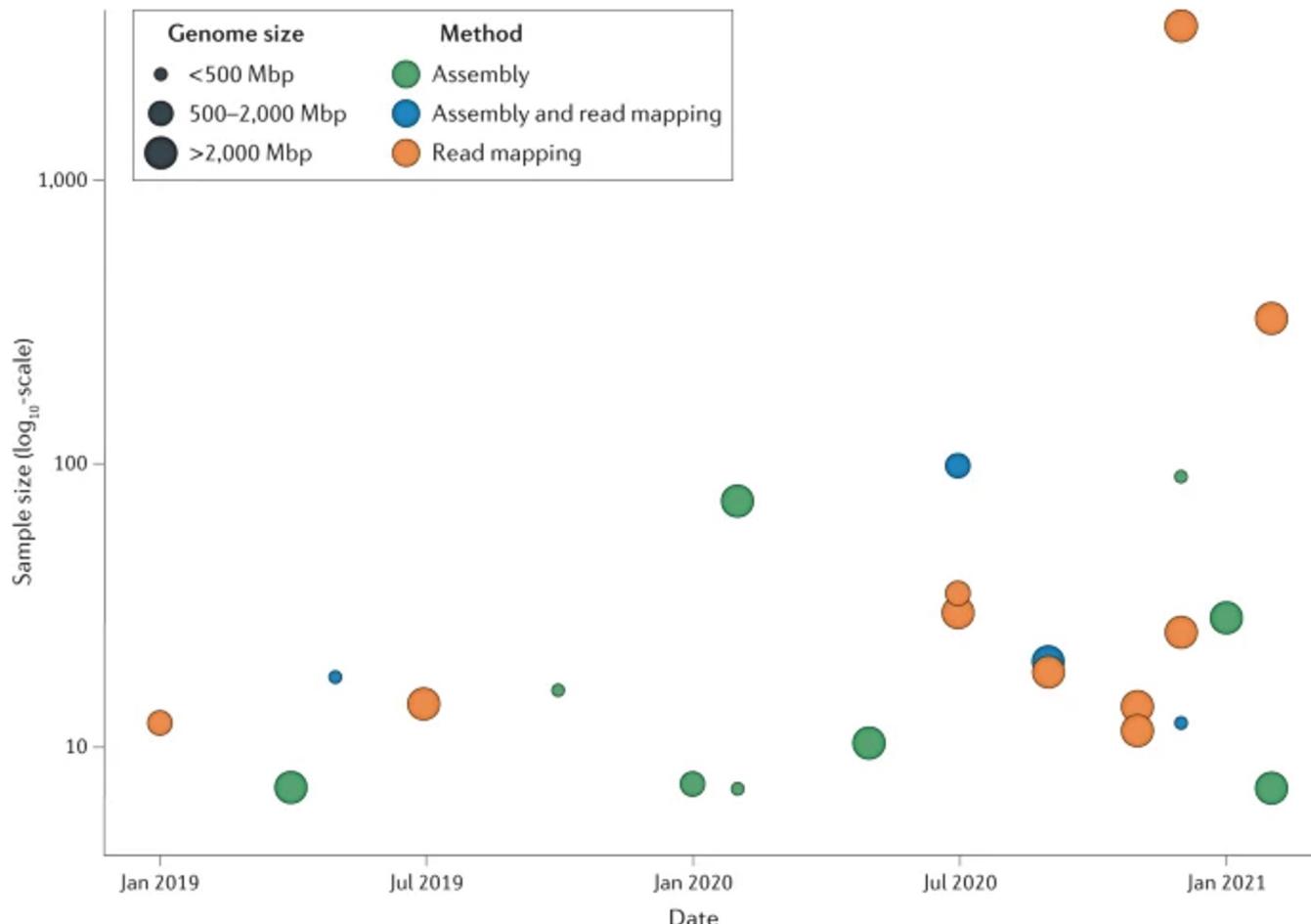
[nature](#) > [nature reviews genetics](#) > [review articles](#) > [article](#)

Review Article | Published: 28 May 2021

Towards population-scale long-read sequencing

Wouter De Coster, Matthias H. Weissensteiner & Fritz J. Sedlazeck 

Fig. 1: Overview of population-scale studies using long-read sequencing.



Full SV Genotyping: From family to population scale

cuteSV: SV calling → merging → re-genotyping → merging → population VCF

~36 CPU hours

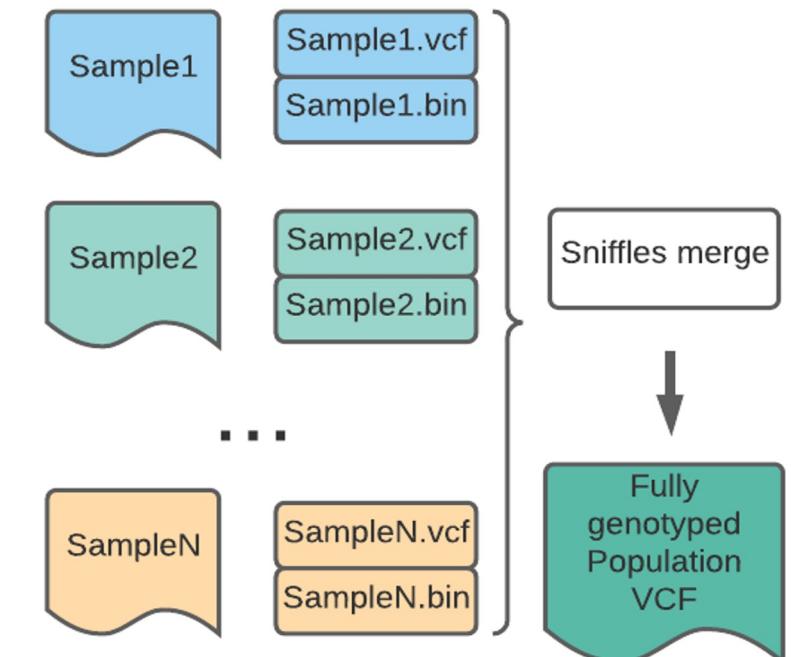
Sniffles2: SV calling → merging → population VCF

65 seconds (**>2000x faster merging**)

Solves n+1 problem

Scaling up to population level

Improves tumor vs. normal



Sniffles2 vs cuteSV: Family Genotyping Accuracy

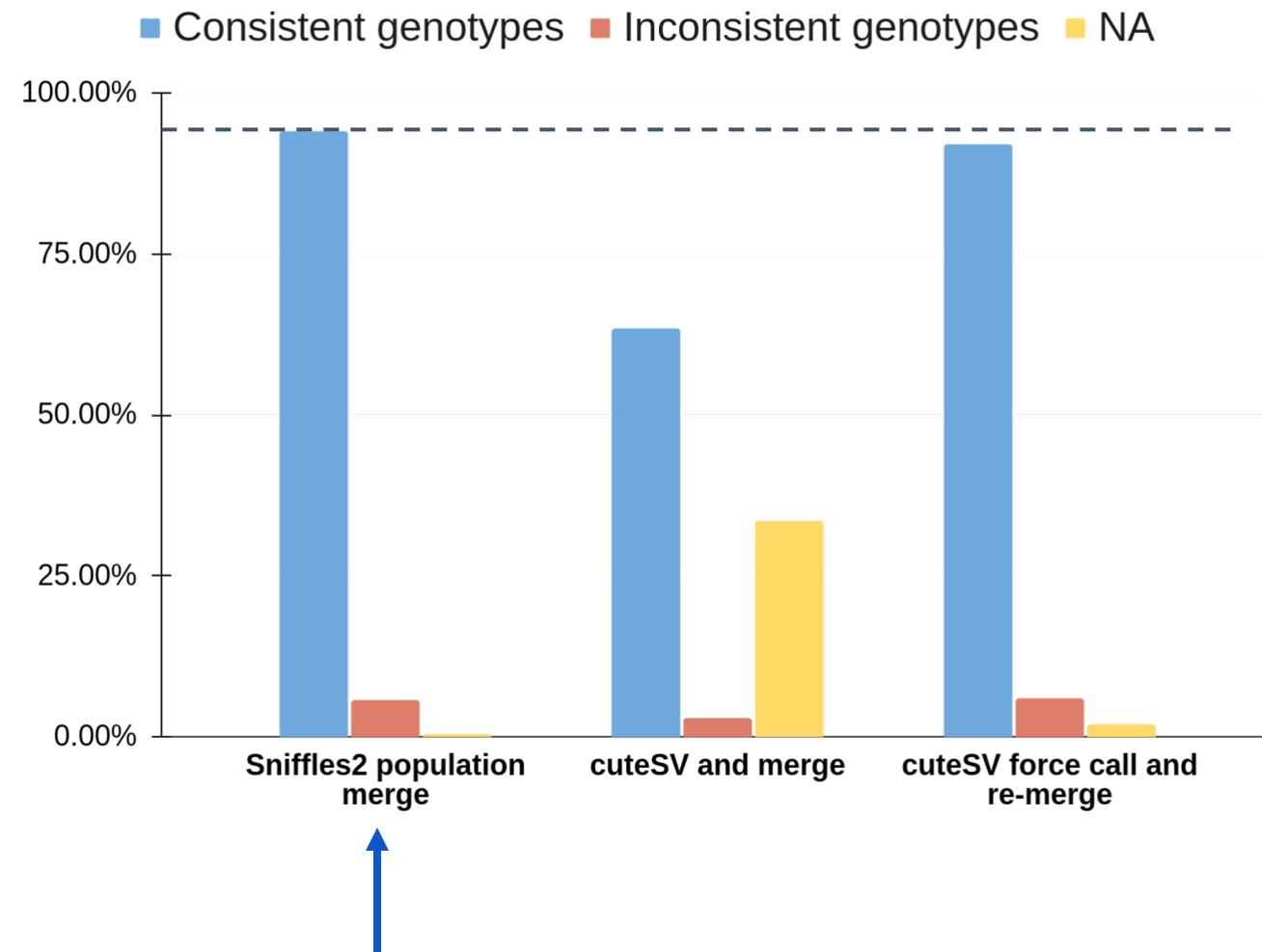
For a family trio*, **Sniffles2**:

- highest fraction of mendelian concordant genotypes (**blue**)
- fewer incomplete genotypes (**yellow**)
- Comparable, yet lower number of non-concordant genotypes (**red**)

Stress test: merging 768 genomes:

15.03 CPU hours

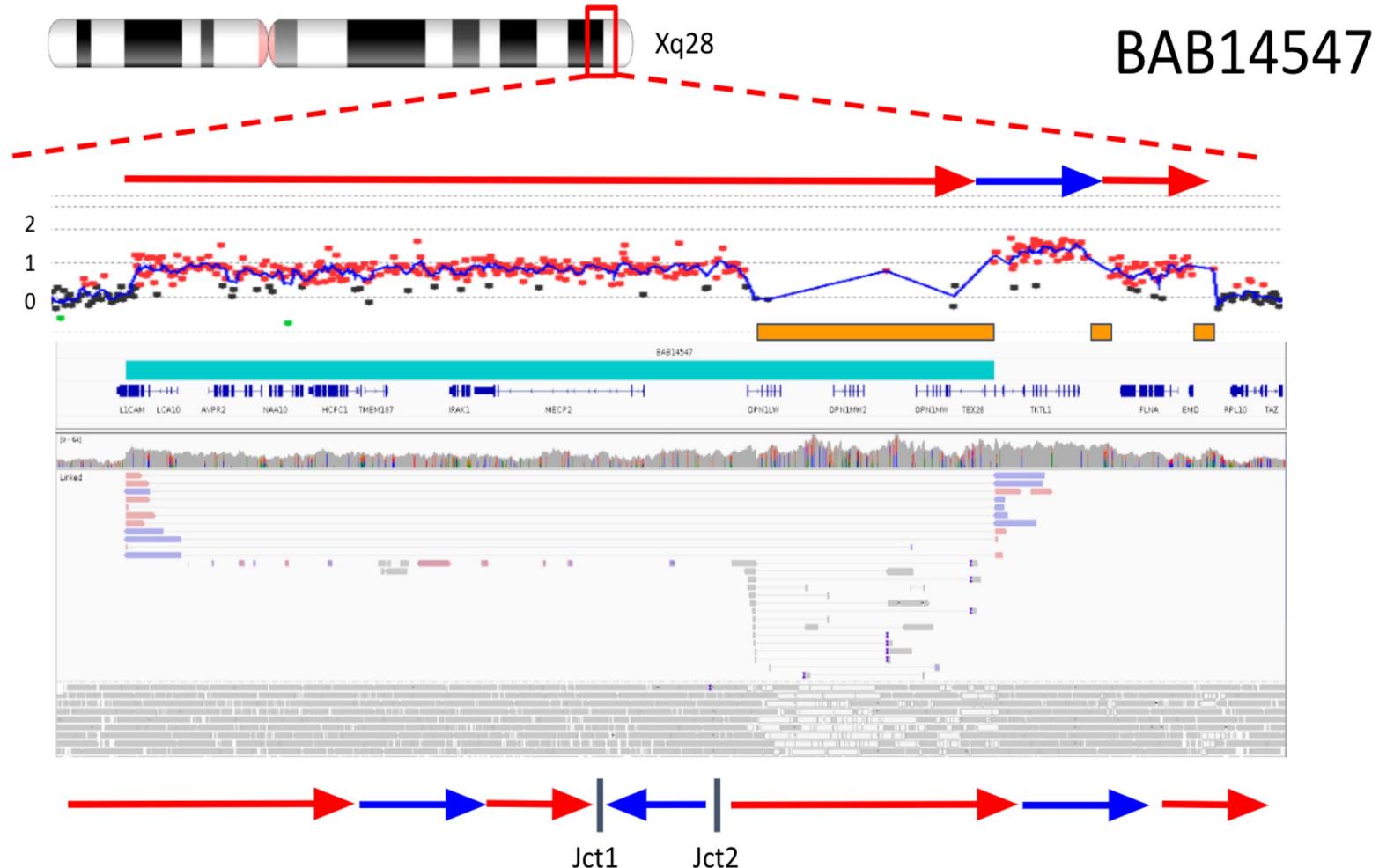
~2x faster than trio with cuteSV



*Benchmark data: HG002/3/4 family trio, ONT.

Sniffles2: Resolving SVs in *MECP2* Duplication Syndrome (MDS)

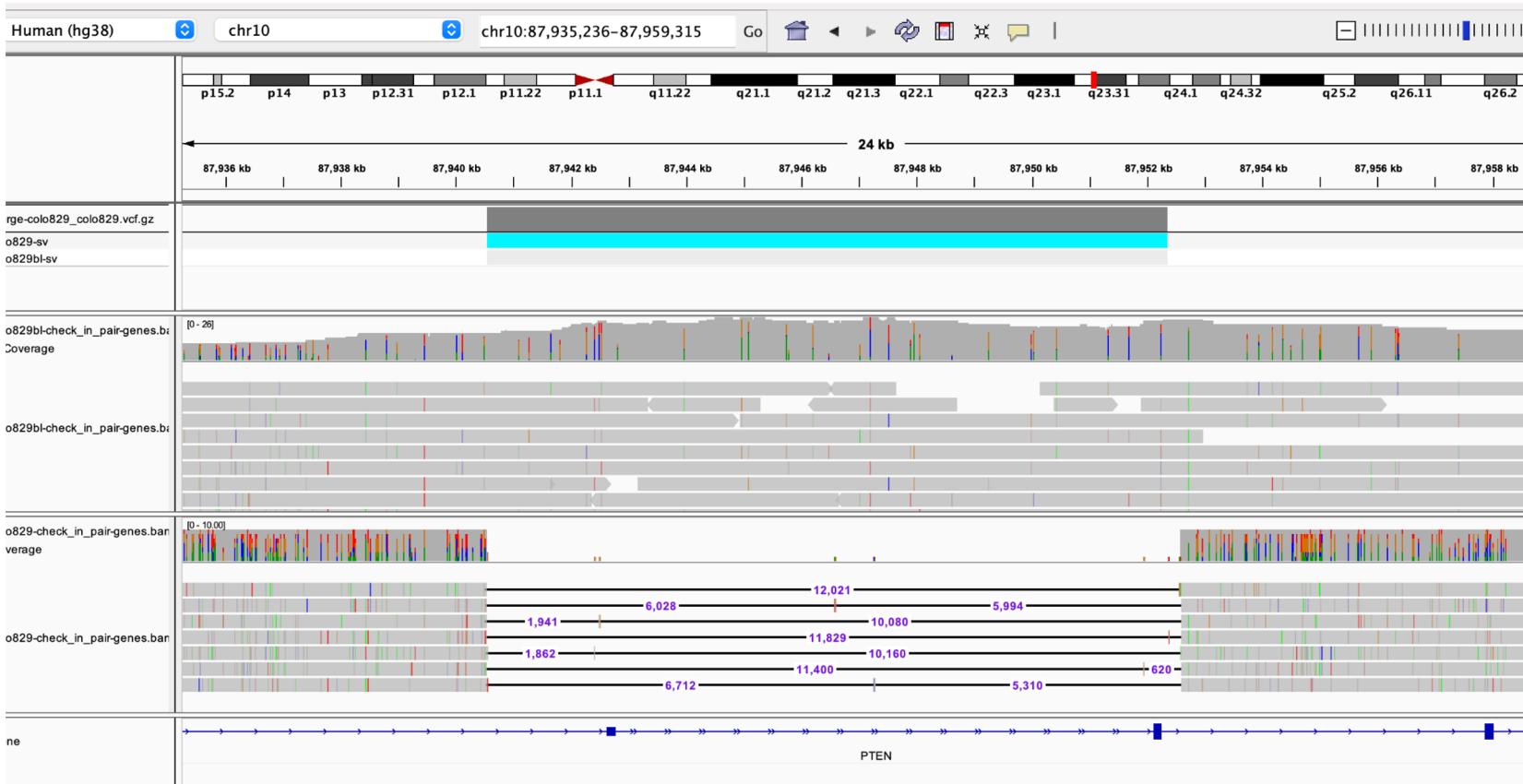
- *MECP2*: profound neurologic and developmental delay in affected males.
- SV resolution can improve M participant outcomes.



in collaboration with Claudia Carvalho (PNRI)

Tumor vs normal: colo829

- Improves SV prioritization
- Mutations of PTEN are a step in the development of many cancers



Applications: research groups

- Gregor
 - Solving unsolved mendelian diseases
 - Complex variants in hard to assess regions
- All of Us
 - 1 million Illumina clinical WGS & 2 million arrays
 - Report findings back to participants
 - ONT will be applied on a subset for research
- Emirates (G42)
 - 85,000 ONT WGS genomes sequenced
 - Annotation resource

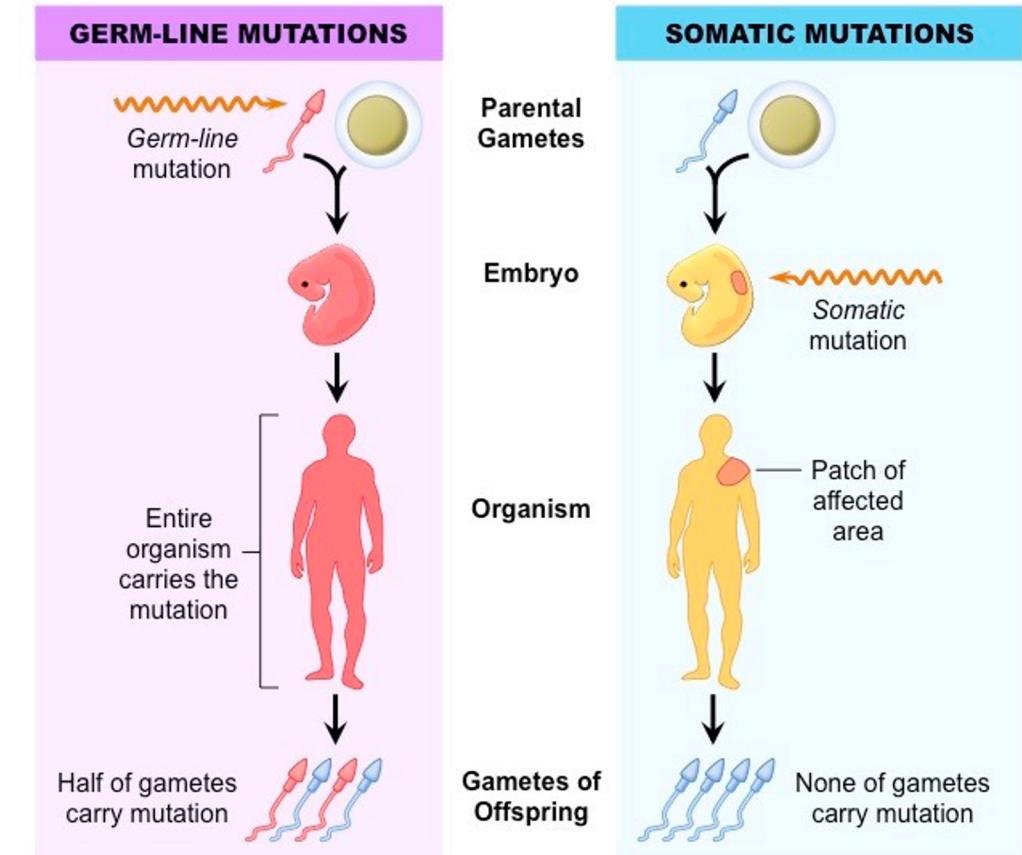


New Applications in SV detection

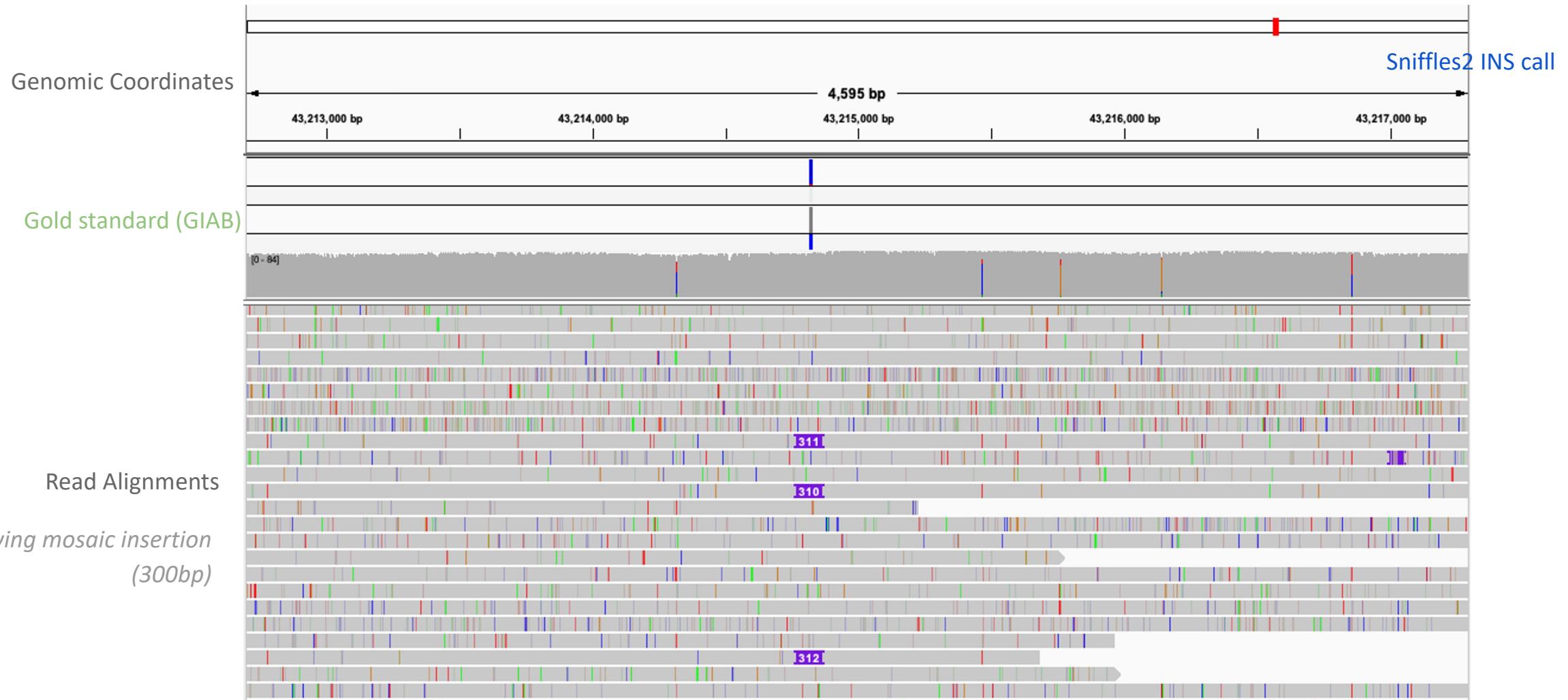
1. Germline SV
2. Population scale
3. **Somatic SVs and human disease:**
 - Neurodegenerative disorders -
accounting for non-heritable disease risk?
 - Cancer drivers (subclonal level)

Review

Somatic mutations in neurodegeneration: An update

Christos Proukakis 

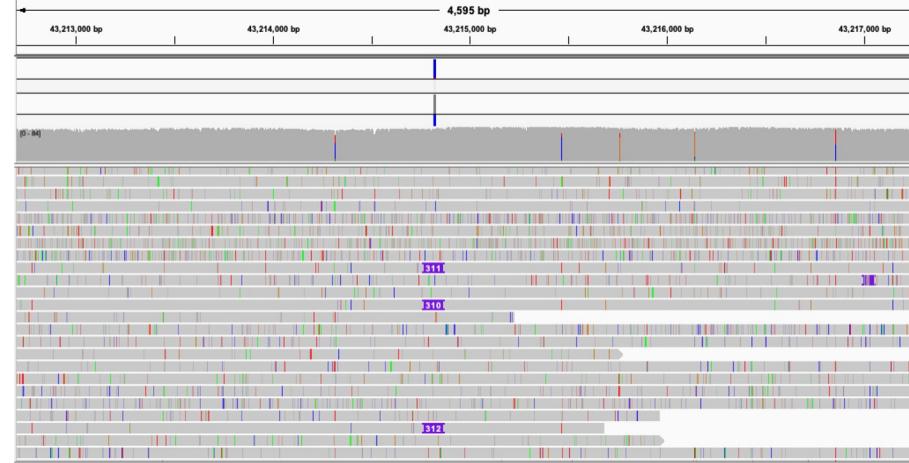
Detecting rare SVs with Sniffles2: Mosaic



Data: Real data spike-in mosaicism of HG002 into HG004

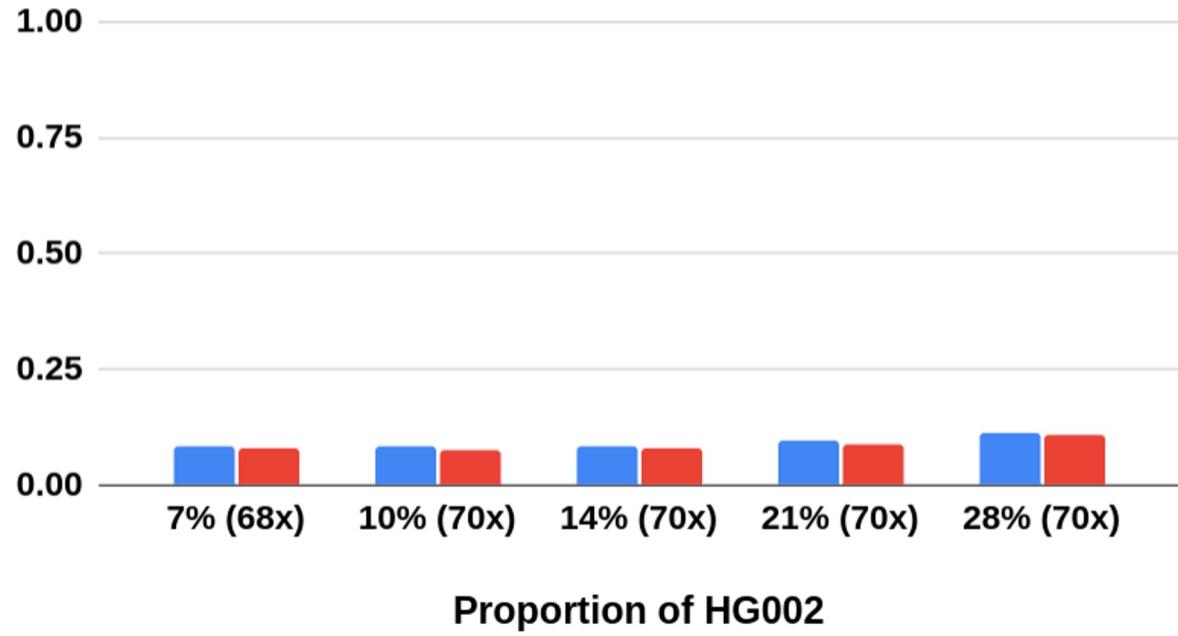
Sniffles 2 mosaic

HG002 (low coverage) HG00733 (high coverage)



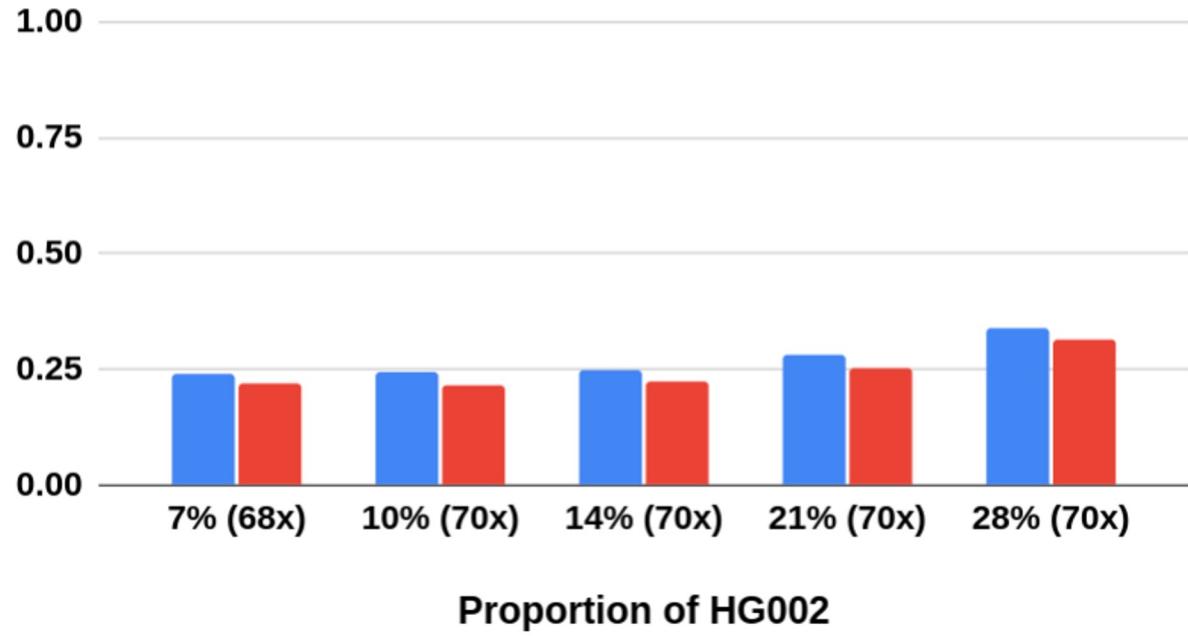
Precision

■ Default ■ CuteSV



Recall

■ Default ■ CuteSV

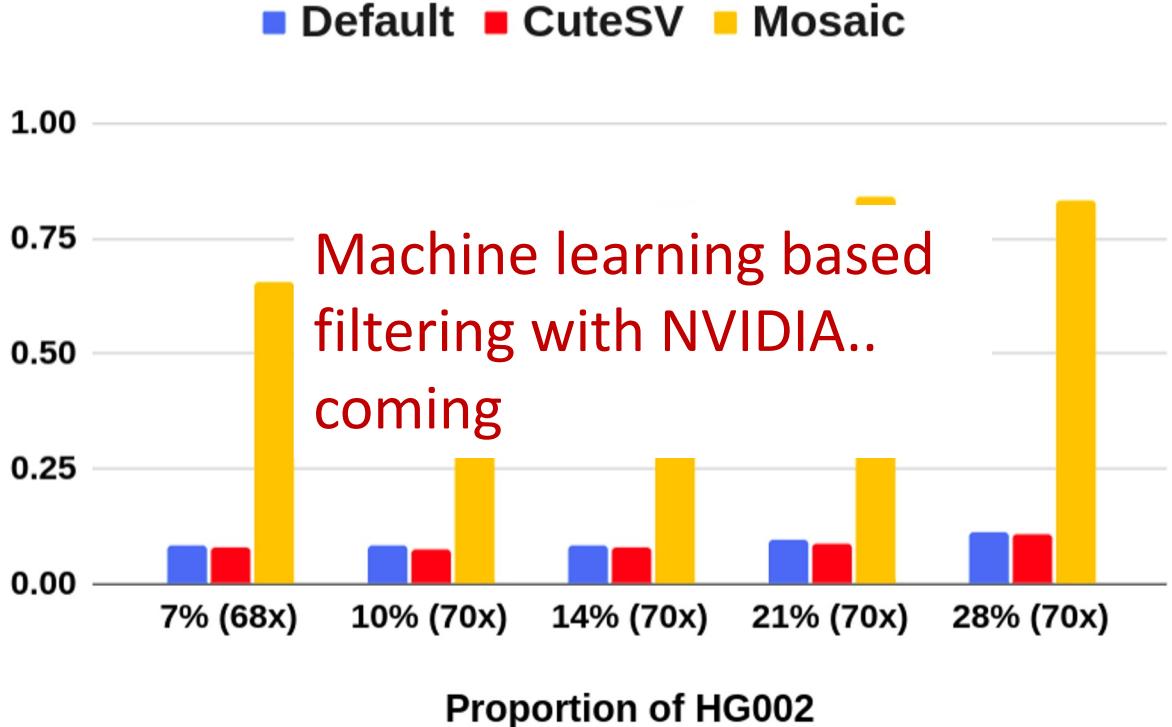


Sniffles 2 mosaic

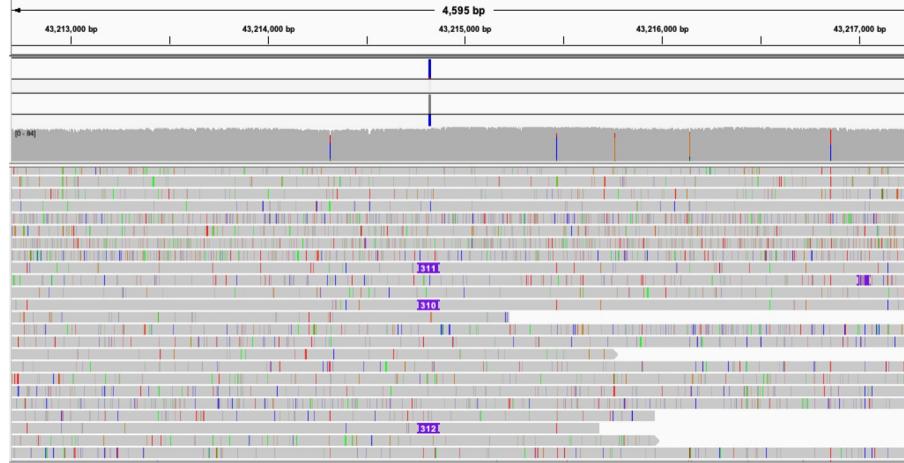
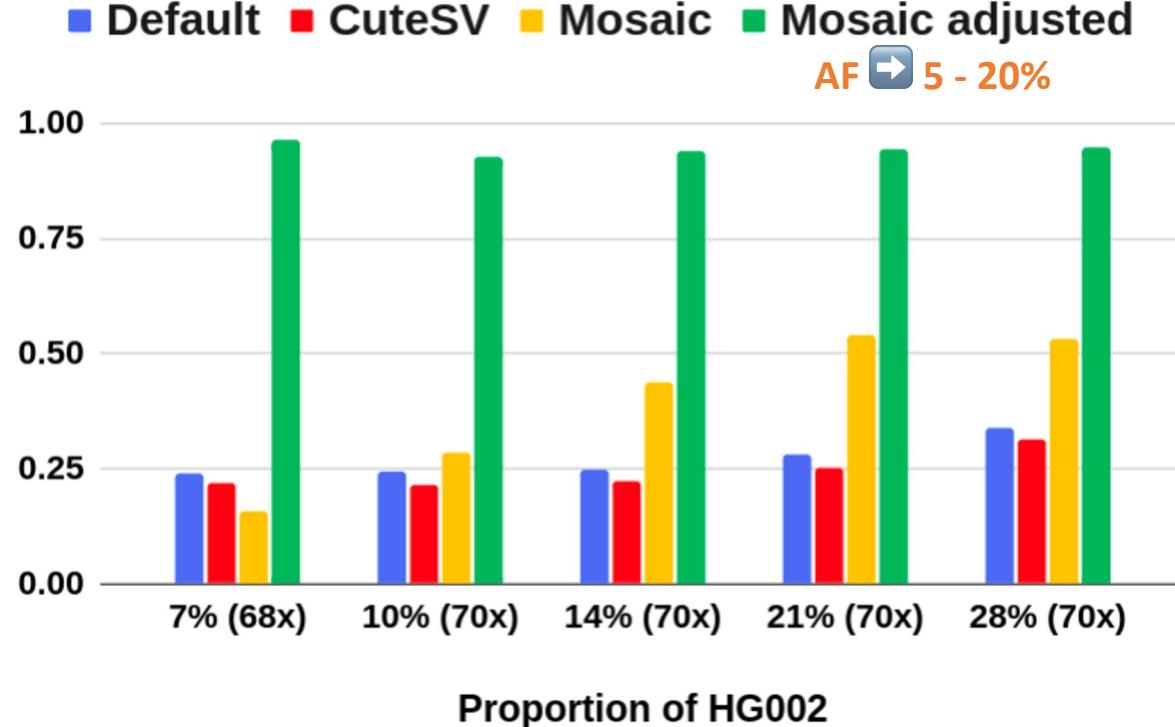
HG002 (low coverage) HG00733 (high coverage)

Mosaic mode identifies rare SVs (down to 2-3 reads only)

Precision



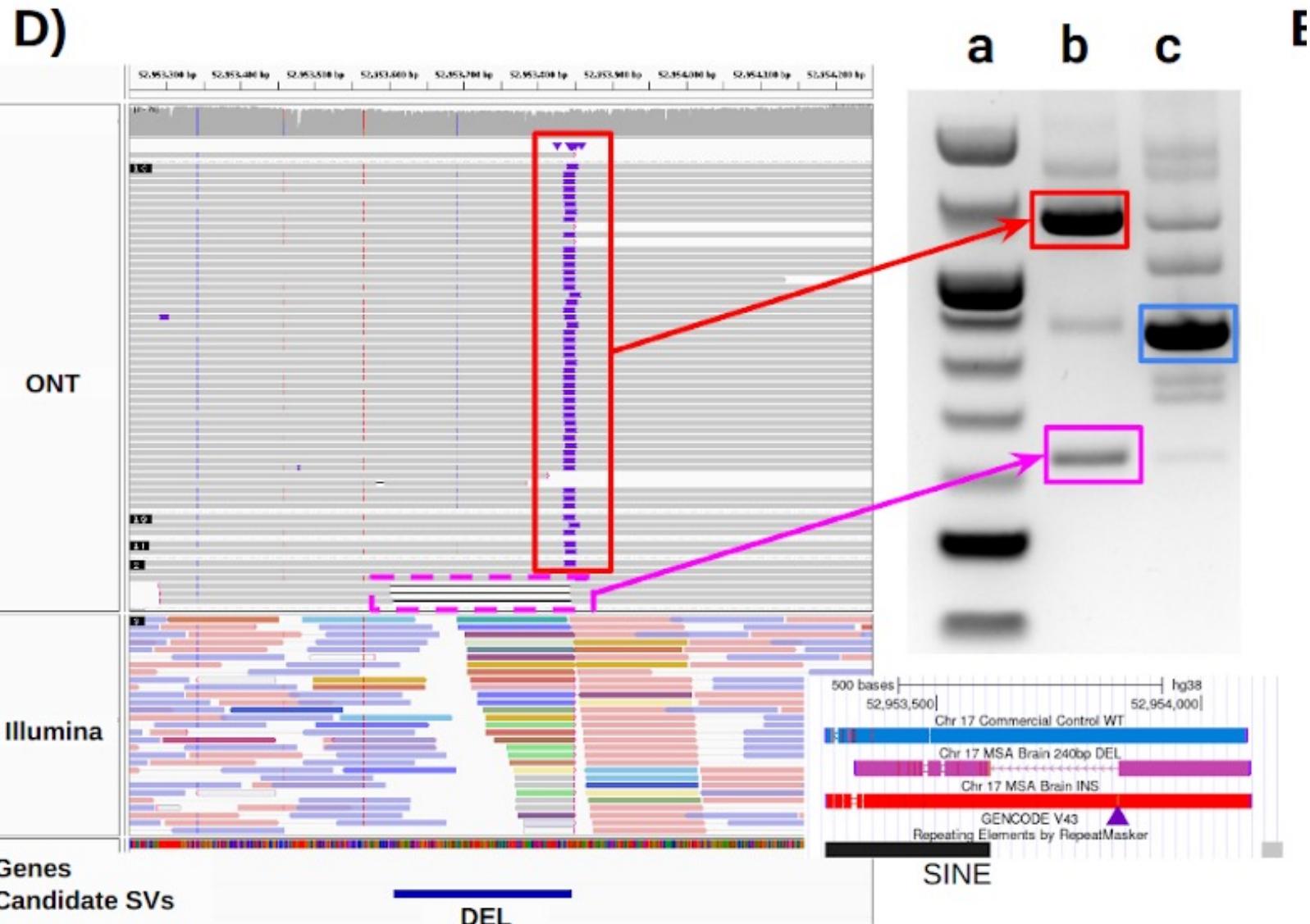
Recall



Sniffles 2 mosaic

55x MSA sample:

- Rare neurodegenerative disorder
- Progressive autonomic dysfunction
- Parkinson-like symptoms
- 26 Alu -Alu → mosaic del
- 125 Ins → mosaic del



Applications: research groups including somatic variants

- Center for Alzheimer's and Related Dementias (CARD)
 - 4,000 ONT genomes
 - 3 different neurological diseases + 1 control group
 - Variant and epigenetic data resource



- Canada's Michael Smith Genome Sciences Centre (Marathon of Hope)
 - Hundreds of ONT cancer + normal samples
 - Illumina RNA seq
- Genomics England
 - Developing cancer pipeline



Acknowledgments



 @sedlazeck



Today hands on

- [https://github.com/fritzsedlazeck/teaching material/blob/main/2023 SV workshop/Day3.md](https://github.com/fritzsedlazeck/teaching_material/blob/main/2023_SV_workshop/Day3.md)
- We will go over the individual sections.