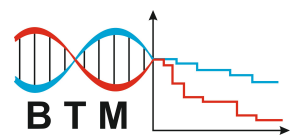


Improving detection of gene duplications in whole-genome sequencing data using allelic depth imbalance

Paweł Sztromwasser
BTM, Medical University of Łódź



Copy number variants

normal



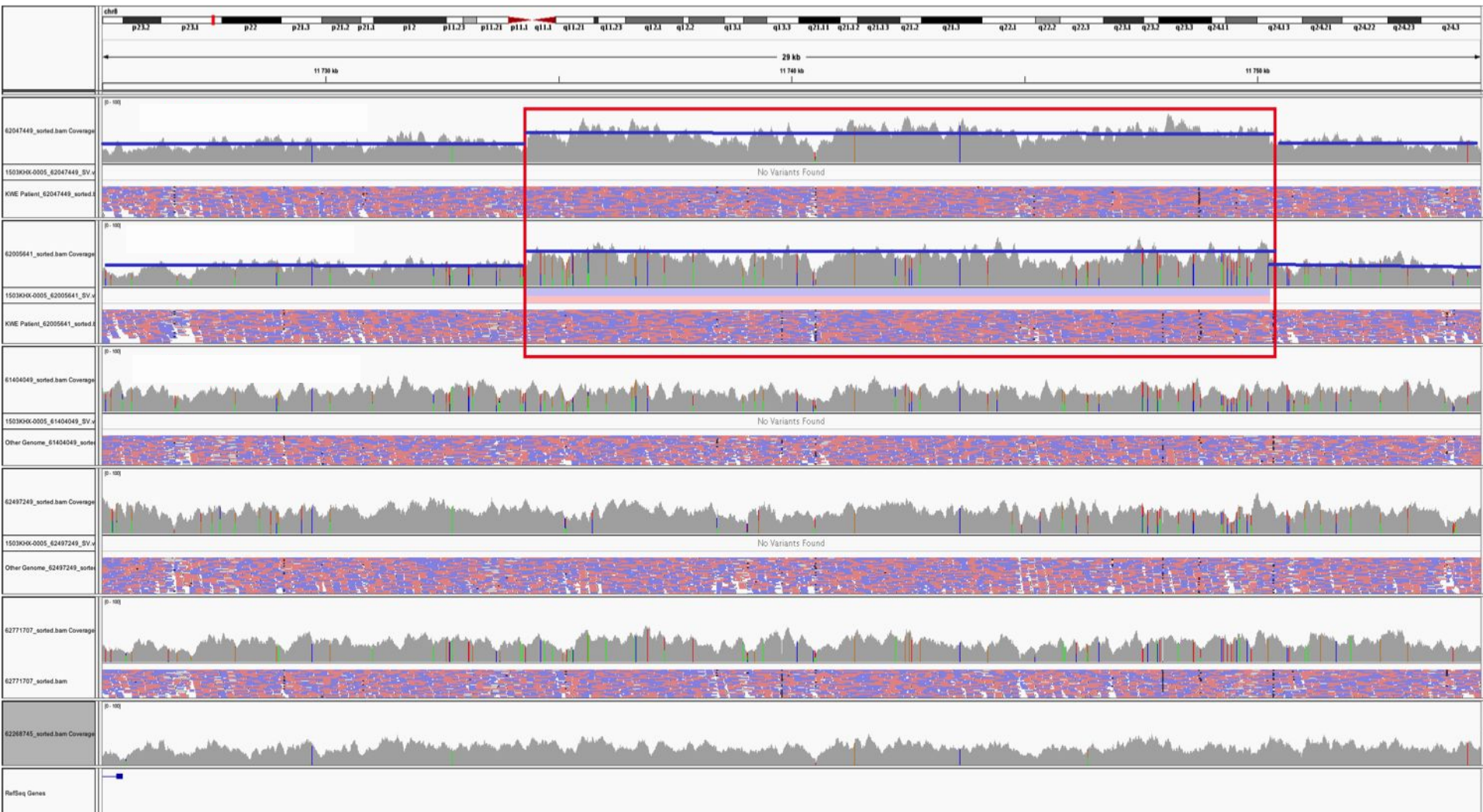
deletion



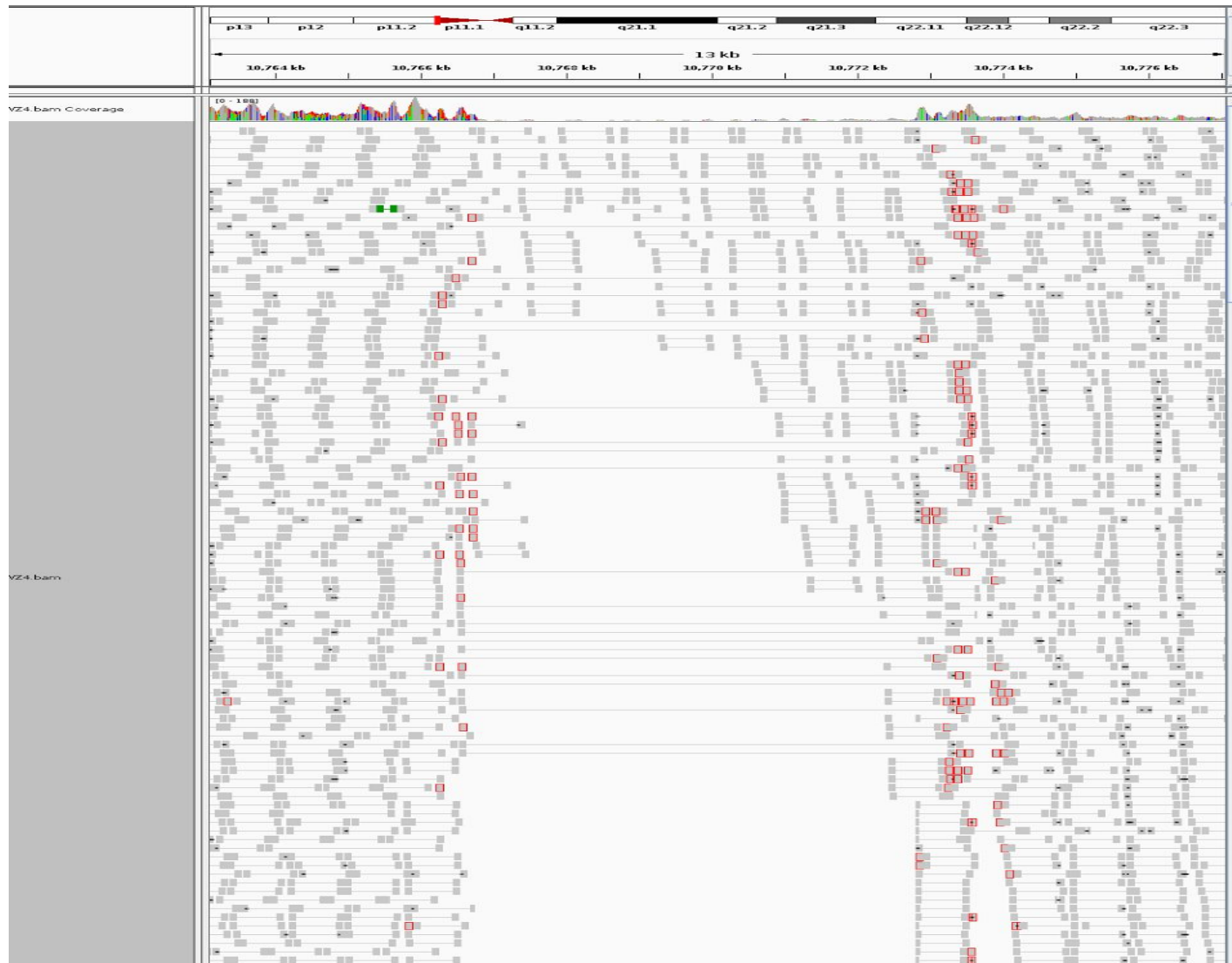
duplication

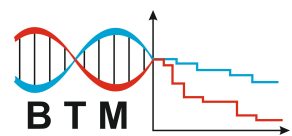


How to detect them in NGS data?



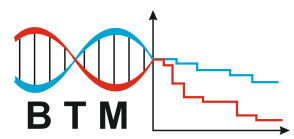
How to detect them in NGS data?





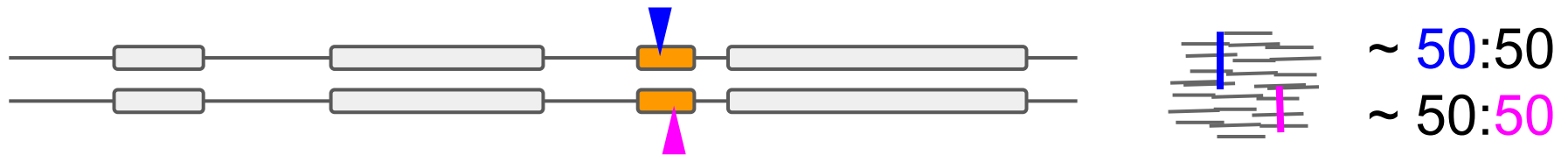
Benchmark results

	LUMPY, TrioCNV, ERDS	
	Precision (PPV)	sensitivity
1000 Genomes deletions (1947)	2.0 - 50.6 %	28.6 - 83.0%
1000 Genomes duplications (90)	1.4 - 6.5 %	7.8 - 43.3 %

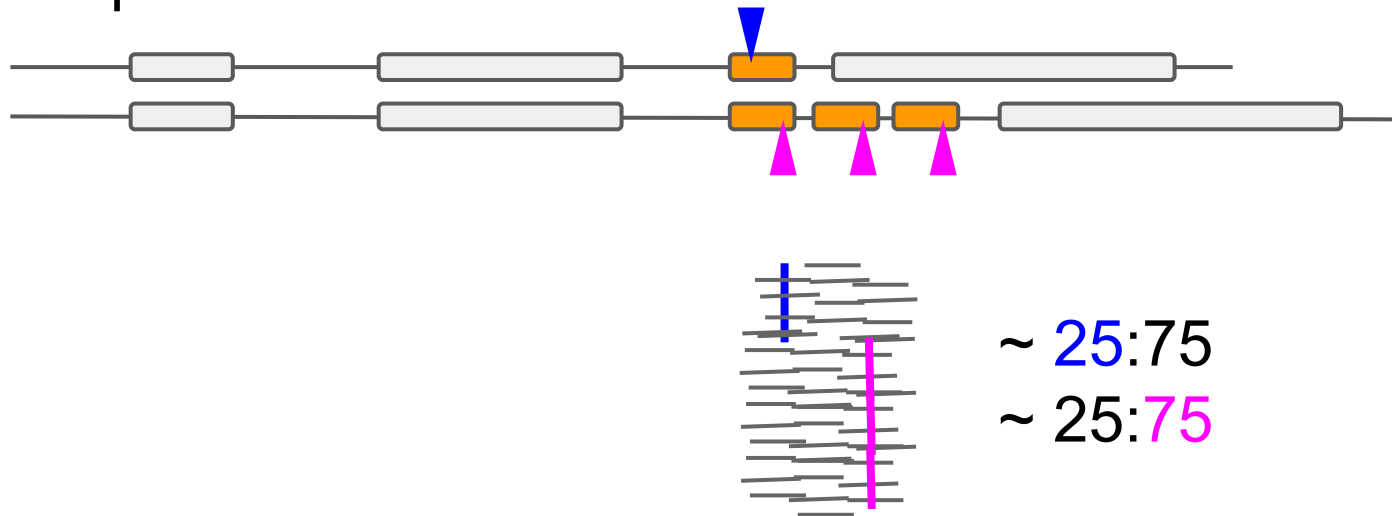


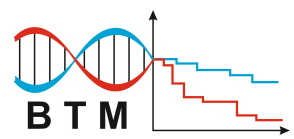
Allelic depth imbalance (ADI)

normal



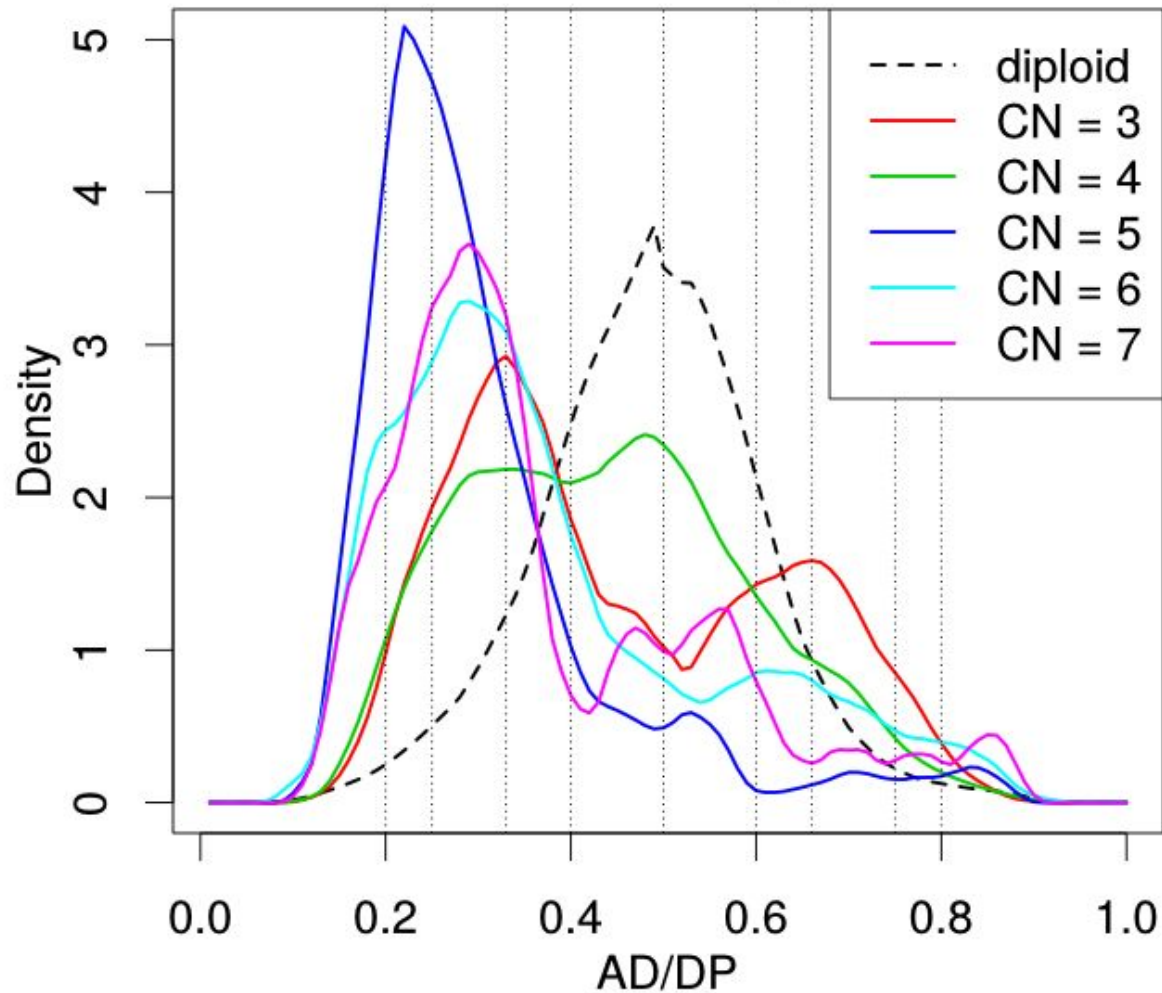
duplication

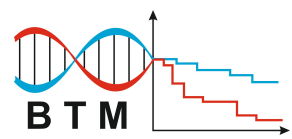




Proof of concept

1000 Genomes duplications





Allelic Depth Imbalance (ADI) score

χ - genomic interval, potential duplication

n - number of het variants overlapped by χ

$$(1) \quad ADI_x = \sum_{i=1}^n \left| 0.5 - \frac{AD_i}{DP_i} \right|$$

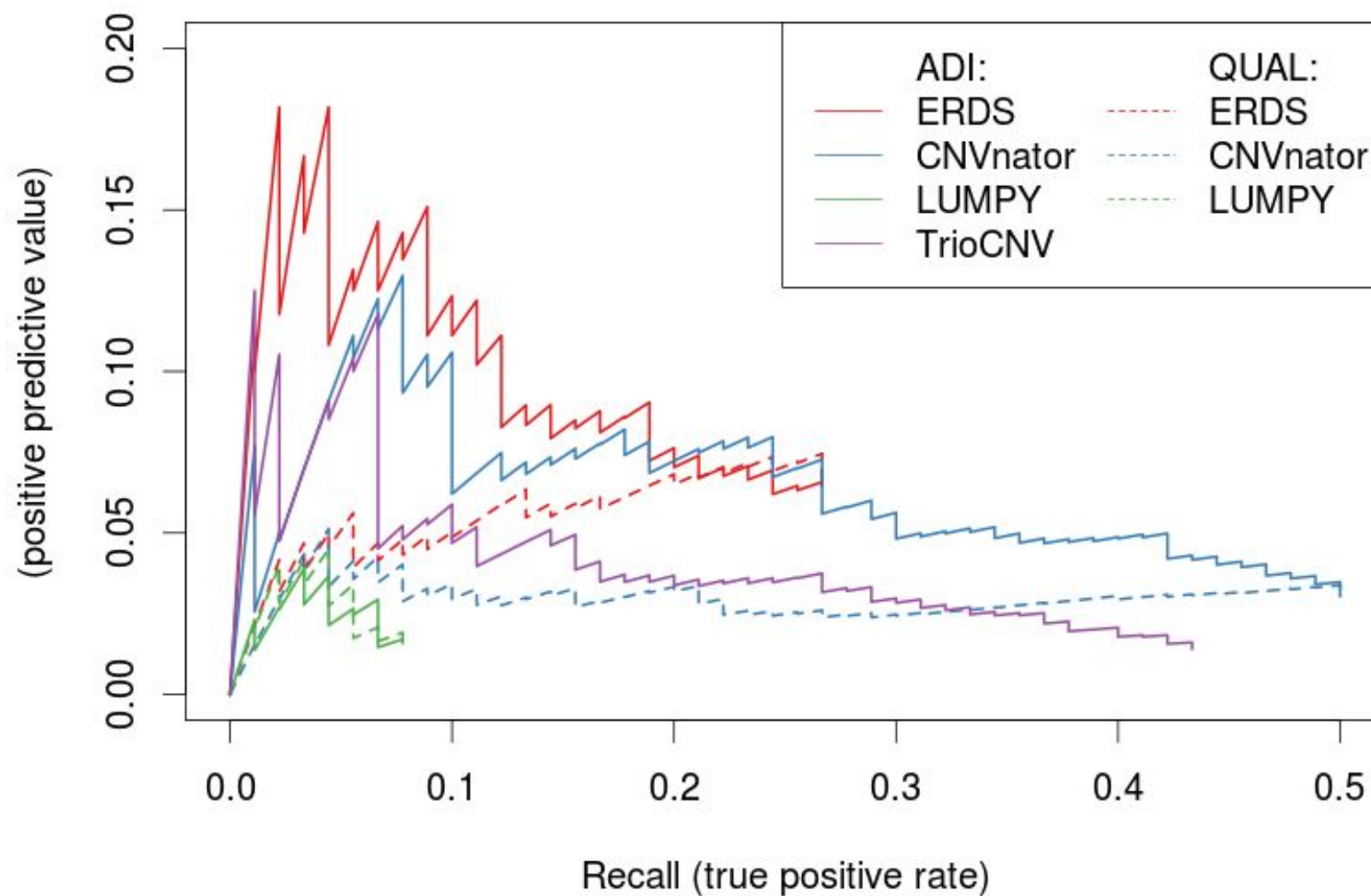
alternative allele depth

total depth

$$(2) \quad B_n = \{b_i : i = 1..1000 \text{ and } b_i \text{ overlaps } n \text{ heterozygous SNVs}\}$$
$$ADI_{B_n} = \{ADI_{b_i} : b_i \in B_n\}$$
$$ADIscore(x, B_n) = rank(ADI_x, ADI_{B_n})$$

Results

Prioritizing calls



Acknowledgements



Inge Jonassen

Kjell Petersen

Vidar Steen

Stefan Johansson

Tomasz Stokowy

Wojciech Fendler



Project funded by Bergen Research Foundation and National Science Center in Poland
(POLONEZ 2016/23/P/NZ2/04251)

This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 665778.

Materials

1. NA12878 genome

sequenced on Illumina HiSeq XTen, 150bp pair-end, 124GBps, 29x
(NA12878D library provided by DNAnexus)

<https://kccg.garvan.org.au/confluence/pages/viewpage.action?pageId=31592745>

2. CEU trio (NA12878, NA12891, and NA12892)

sequenced on Illumina HiSeq 2000, 101bp pair-end, 245-290GBps, 37-64x
(by 1000Genomes Consortium)

ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/technical/working/20120117_ceu_trio_b37_decoy/

Reference CNVs

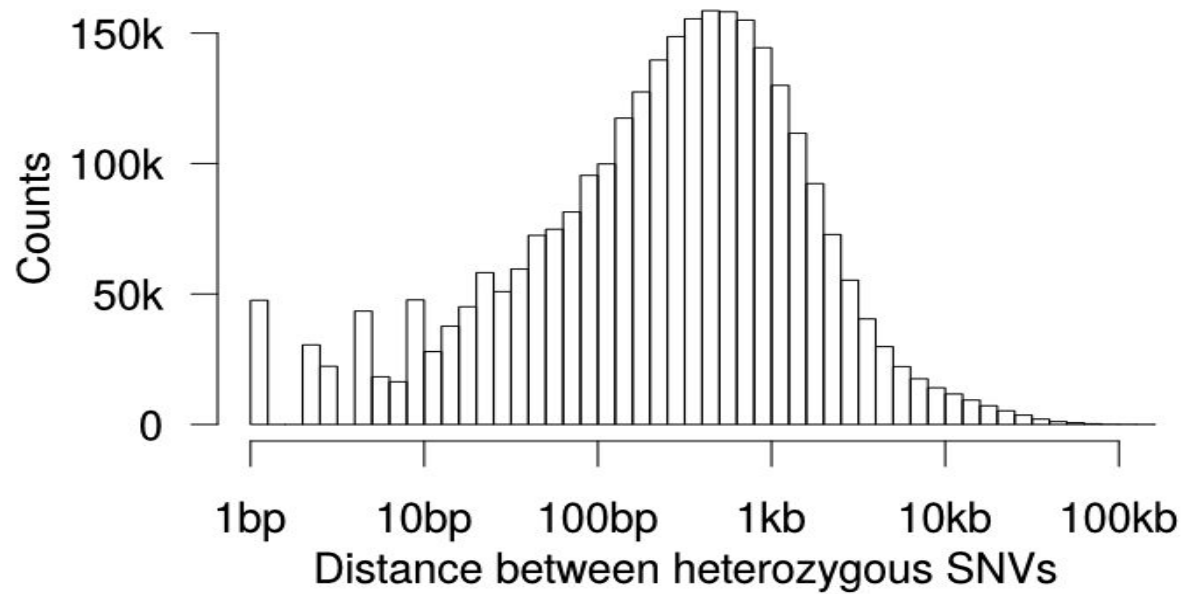
Copy number → ↓ CNV set	0	1	2	3	4	5	6	7
1000 Genomes	624	1411	135	38	57	1	3	1
Conrad (2010)	137	352	(3975)	89	134	13	1	1
Mills (2011) GS	617		-	271*				

*) Conrad + 24 from McCarroll et al, 2008

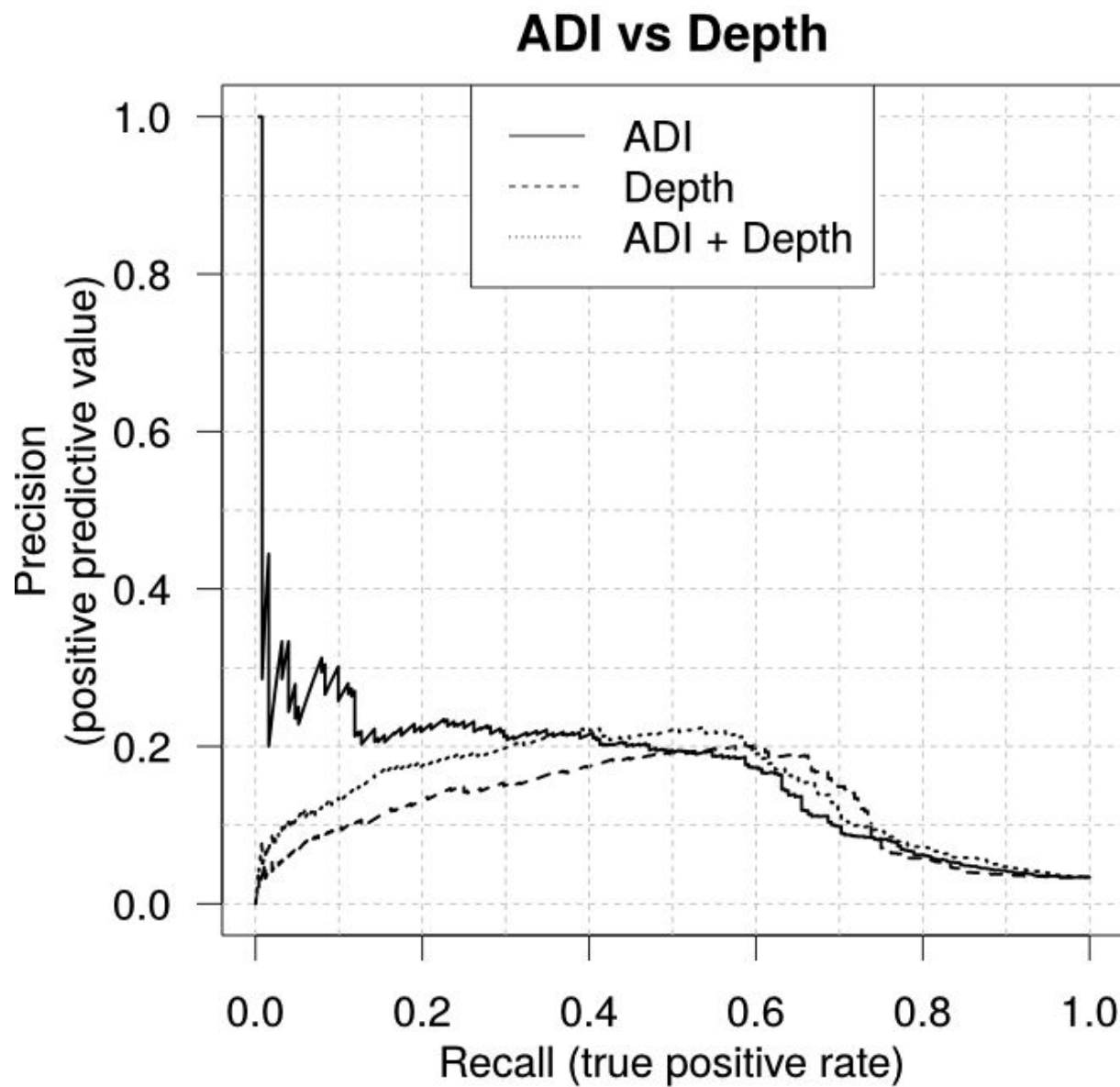
10 out of 238 Conrad variants are overlapping a variant in the 1000 Genomes

	1000 Genomes (100)	Conrad (238)	Mills (271)
1000 Genomes	-	10 (4%)	18 (7%)
Conrad	9 (9%)	-	243 (90%)
Mills	14 (14%)	238 (100%)	-

Can it work?



Performance



Performance

	1000 Genomes	
	precision	sensitivity
LUMPY	1.6%	7.8%
TrioCNV	1.4%	43.3%
ERDS	6.5%	26.7%

