
CS681: Advanced Topics in Computational Biology

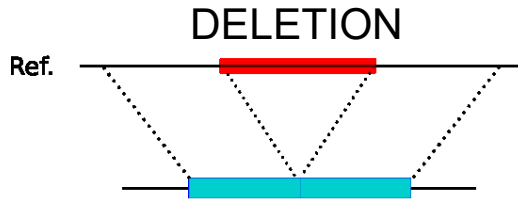
Can Alkan

EA509

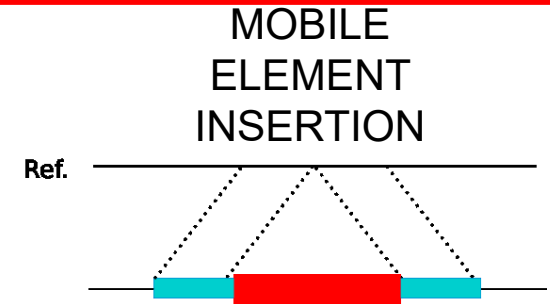
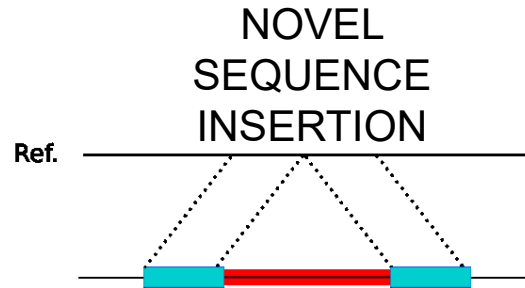
calkan@cs.bilkent.edu.tr

<http://www.cs.bilkent.edu.tr/~calkan/teaching/cs681/>

Structural Variation Classes

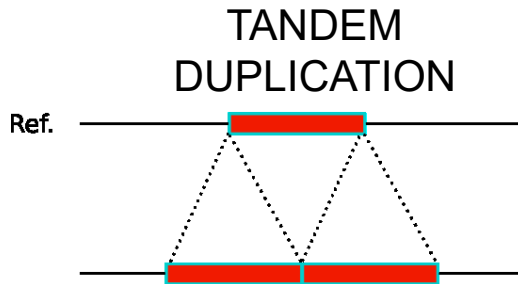


Autism, mental retardation, Crohn's

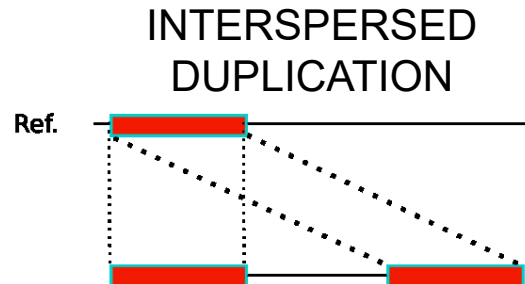


Alu/L1/SVA

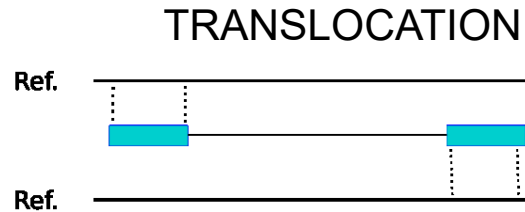
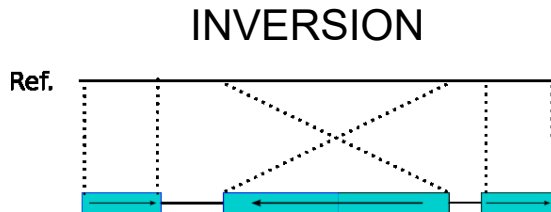
Haemophilia



Schizophrenia, psoriasis



CNV: Copy number variants



Chronic myelogenous leukemia

Balanced rearrangements

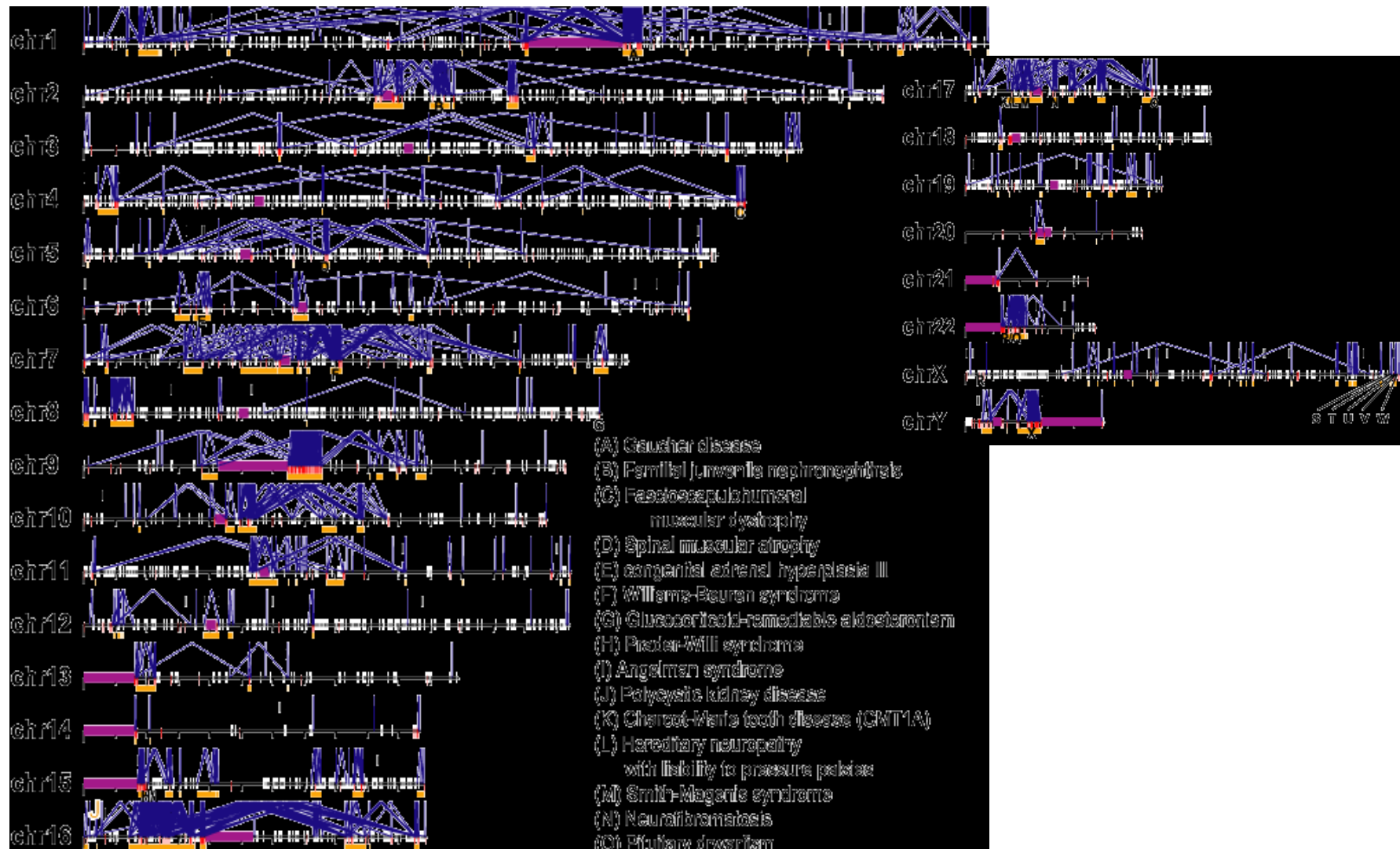
Structural variation discovery with HTS data

- ❑ SVs: genomic alterations > 50 bp.
 - ❑ Databases:
 - dbVar: <http://www.ncbi.nlm.nih.gov/dbvar/>
 - DGV: <http://projects.tcag.ca/variation/>
 - ❑ Input: sequence data and reference genome
 - ❑ Output: set of SVs and their genotypes (homozygous/heterozygous)
 - ❑ Often there are errors, filtering required
 - ❑ SV detection methods can be based on statistical analysis or combinatorial optimization
 - ❑ Tools:
 - ❑ Illumina: TARDIS, LUMPY, DELLY, Manta, TIDDIT, Genome STRiP, etc.
 - ❑ Long reads: Sniffles, cuteSV, etc.
-

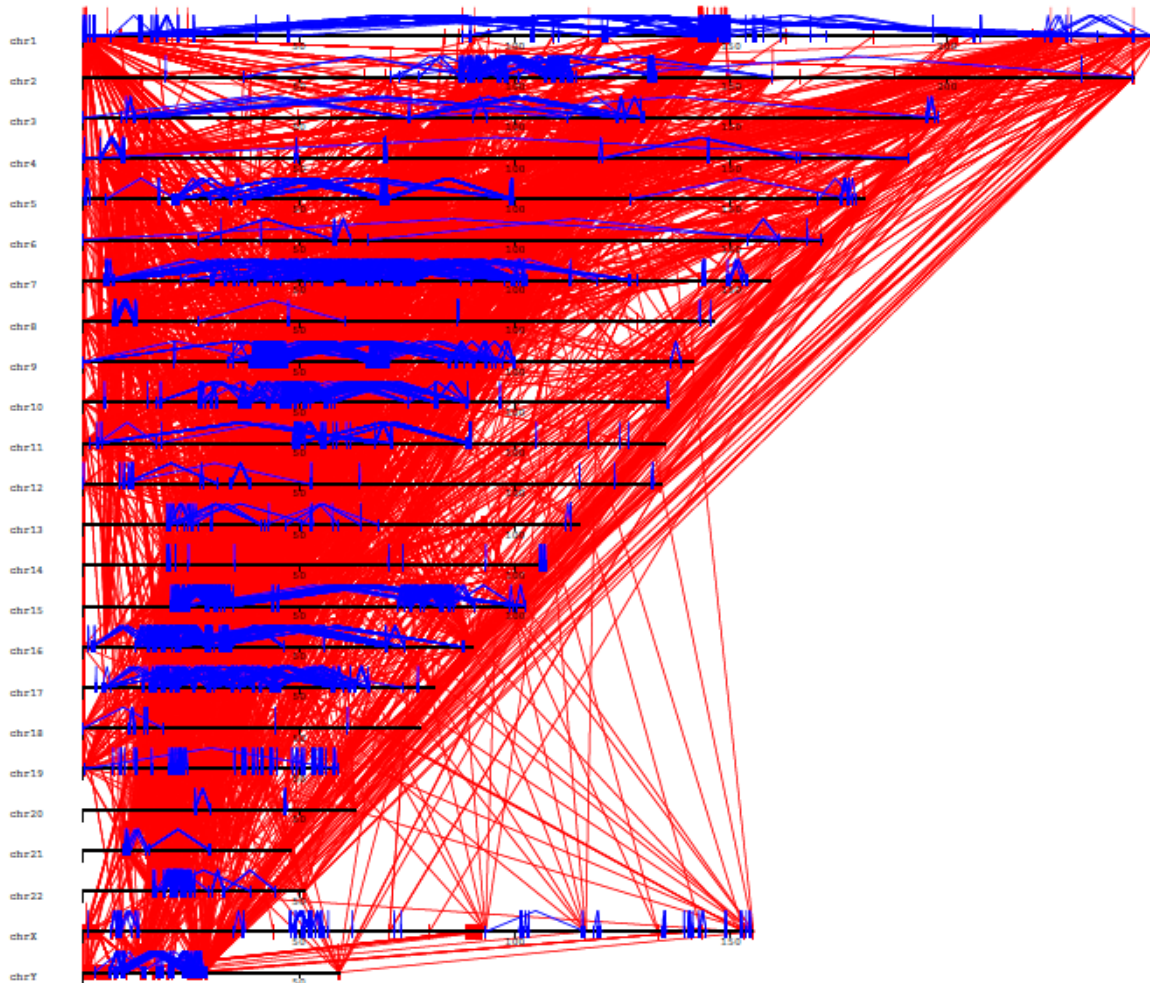
Challenges

- Most SVs are embedded within or around segmental duplications or long repeats
 - If you use unique mapping, you will lose sensitivity
 - Ambiguous mapping of reads will increase false positives
 - Reference genome is incomplete; missing portions are duplications which cause more problems in accurate detection
 - Many SVs are complex; many rearrangements at the same site
 - CNV discovery is heavily studied but still not perfect; detection of balanced rearrangements are still problematic
-

Duplications and CNV hotspots



Duplications: inter & intra



- 51,599 pairs of SDs
 - 18,559 pairs intrachromosomal
 - 32,740 pairs interchromosomal
- Non-redundant corresponds to 166 Mb (~5% of genome)

Genome-wide SV Discovery Approaches

Hybridization-based

- lafrate et al., 2004, Sebat et al., 2004
- SNP microarrays: McCarroll *et al.*, 2008, Cooper *et al.*, 2008, Itsara *et al.*, 2009
- Array CGH: Redon *et al.* 2006, Conrad *et al.*, 2010, Park *et al.*, 2010, WTCCC, 2010

Single molecule analysis

- **Optical mapping:** Teague et al., 2010

Sequencing-based

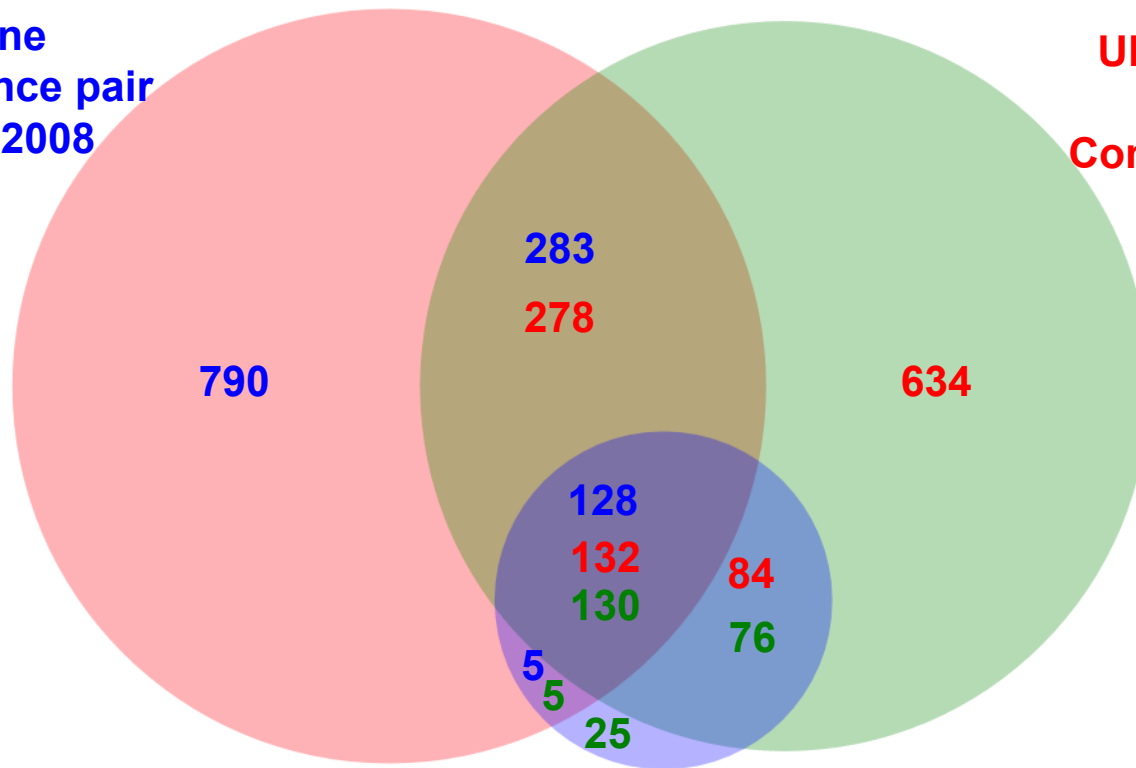
- Read-depth: Bailey et al, 2002
- Fosmid ESP: Tuzun *et al.* 2005, Kidd *et al.* 2008
- Sanger sequencing: Mills *et al.*, 2006
- Next-gen sequencing: Korbelt *et al.* 2007, Yoon *et al.*, 2009, Alkan et al., 2009, Hormozdiari *et al.* 2009, Chen *et al.* 2009,
 - 1000 Genomes Project

Detection diversity

Gains & Losses > 5 Kbp in the same 5 individuals

**Fosmid clone
End-sequence pair
Kidd et al., 2008
(N = 1,206)**

**Ultra-dense tiling
array CGH
Conrad et al., 2010
(N = 1,128)**



**Affymetrix 6.0 SNP microarray
McCarroll et al., 2008 (N = 236)**

Sequencing technologies

Short-Read

Illumina

- 100-200bp
- Paired-end
- Billions of reads
- < 0.1% error



Long Read



PacBio and Oxford Nanopore

- > 10 Kb, up to 1 Mb
- Single-end
- Hundreds of millions of reads
- 12-20% error – indel dominated

Long Range



10X + Illumina

- 100-200bp
- Paired-end
- Billions of reads
- < 0.1% error
- Barcoded: 30-50 Kb molecule range

Sequencing technologies - algorithms

Short-Read

ILLUMINA

TARDIS

DELLY

LUMPY

Manta

Pindel

CNVnator



Long Read



PacBio and Oxford Nanopore

SMRT-SV

Sniffles

PBHoney

Picky

Multiplatform (Long + Short read)

HySa

CORGi

pbsv

NanoSV

SVIM

MultiBreak-SV

Long Range



10X + Illumina

VALOR

GROC-SVs

NAIBR

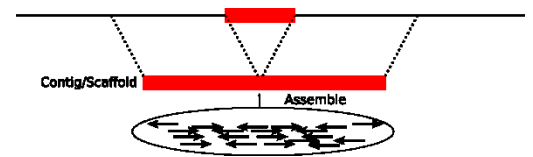
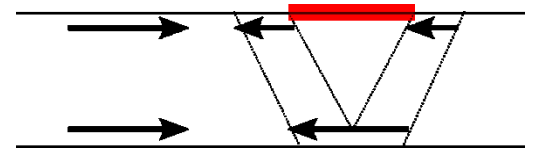
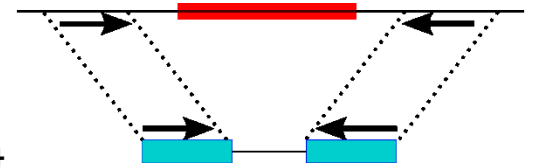
LongRanger

LinkedSV

ZoomX

Sequence signatures of structural variation

- Read pair analysis
 - Deletions, small novel insertions, inversions, transposons
 - Size and breakpoint resolution dependent to insert size
- Read depth analysis
 - Deletions and duplications only
 - Relatively poor breakpoint resolution
- Split read analysis
 - Small novel insertions/deletions, and mobile element insertions
 - 1bp breakpoint resolution
- Local and *de novo* assembly
 - SV in unique segments
 - 1bp breakpoint resolution



SV by sequencing: first algorithms

Recent Segmental Duplications in the Human Genome

Read Depth

1342

Jeffrey A. Bailey,¹ Zhiping Gu,² Royden A. Clark,¹ Knut Reinert,²
Rhea V. Samonte,¹ Stuart Schwartz,¹ Mark D. Adams,²
Eugene W. Myers,² Peter W. Li,² Evan E. Eichler^{1*}

Science, 2002

nature
genetics

↗



Read Pair

1138

Fine-scale structural variation of the human genome

Eray Tuzun^{1,5}, Andrew J Sharp^{1,5}, Jeffrey A Bailey^{2,5}, Rajinder Kaul³, V Anne Morrison¹,
Lisa M Pertz², Eric Haugen³, Hillary Hayden³, Donna Albertson⁴, Daniel Pinkel⁴, Maynard V Olson³ &
Evan E Eichler¹

Nature Genetics, 2005



Split read

592

An initial map of insertion and deletion (INDEL) variation in the human genome

Ryan E. Mills,^{1,2} Christopher T. Luttig,¹ Christine E. Larkins,³ Adam Beauchamp,⁴
Circe Tsui,^{1,2} W. Stephen Pittard,^{2,5} and Scott E. Devine^{1,2,3,4,6}

Genome Research, 2006



All these first algorithms used Sanger sequence, but laid out the basic principles for HTS analysis

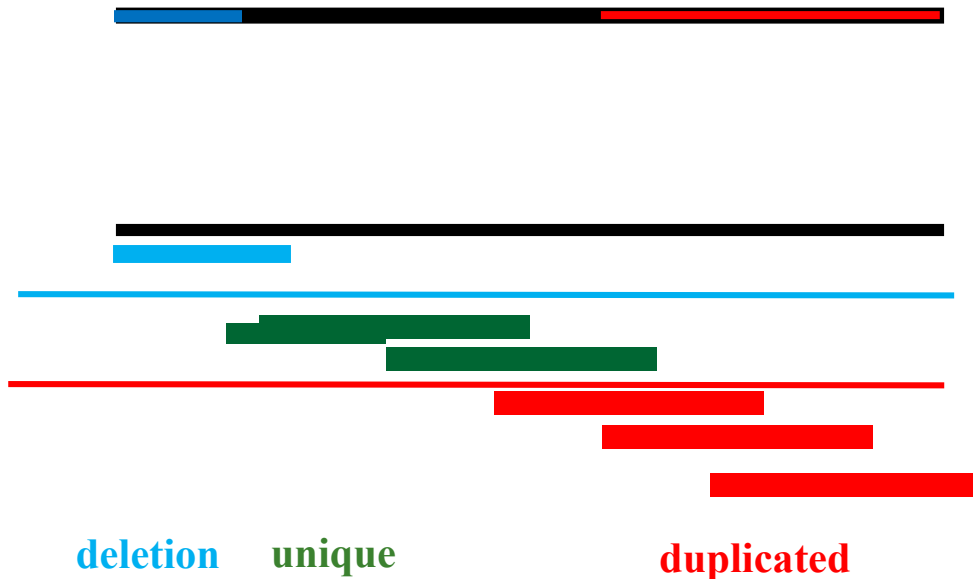
Read depth based algorithms

- Assume random (Poisson) distribution in read depth
 - Multiple mapping:
 - WSSD (whole genome shotgun sequence detection)
 - Unique mapping:
 - Low resolution: Campbell et al. Nat Genet 2008, Chiang et al. Nat Meth, 2009 (SegSeq)
 - High(er) resolution: CNVnator, EWT (RDXplorer)
-

Read depth analysis: WSSD

- Uses database of random reads to confirm duplicated nature of the sequence
 - increased # of copies => increased number of reads
 - decreased # of copies => decreased number of reads
- Compute depth-of-coverage in 5kb windows (sliding by 1kb); select regions with increased depth as **duplications**, regions with reduced depth as **deletions** (WSSD method)

Sequence to Test



Random Genome Sample (Whole-Genome Shotgun Sequence)



Multiple vs. unique mapping

Genome

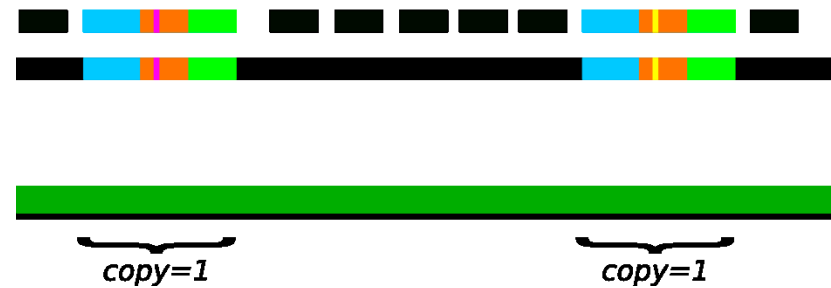
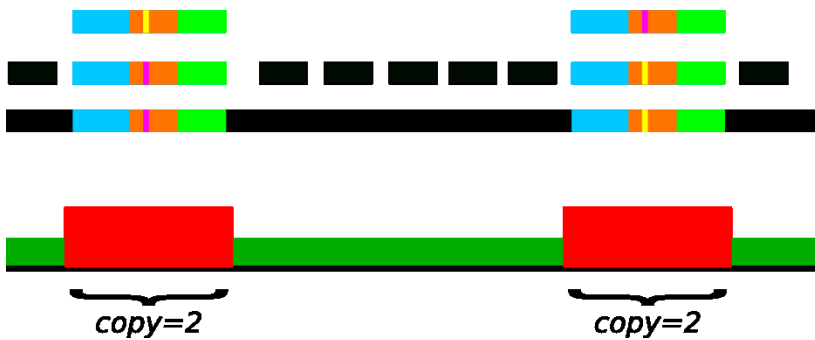
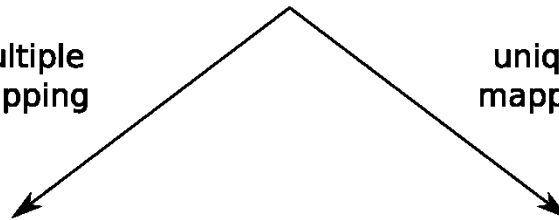


Reads

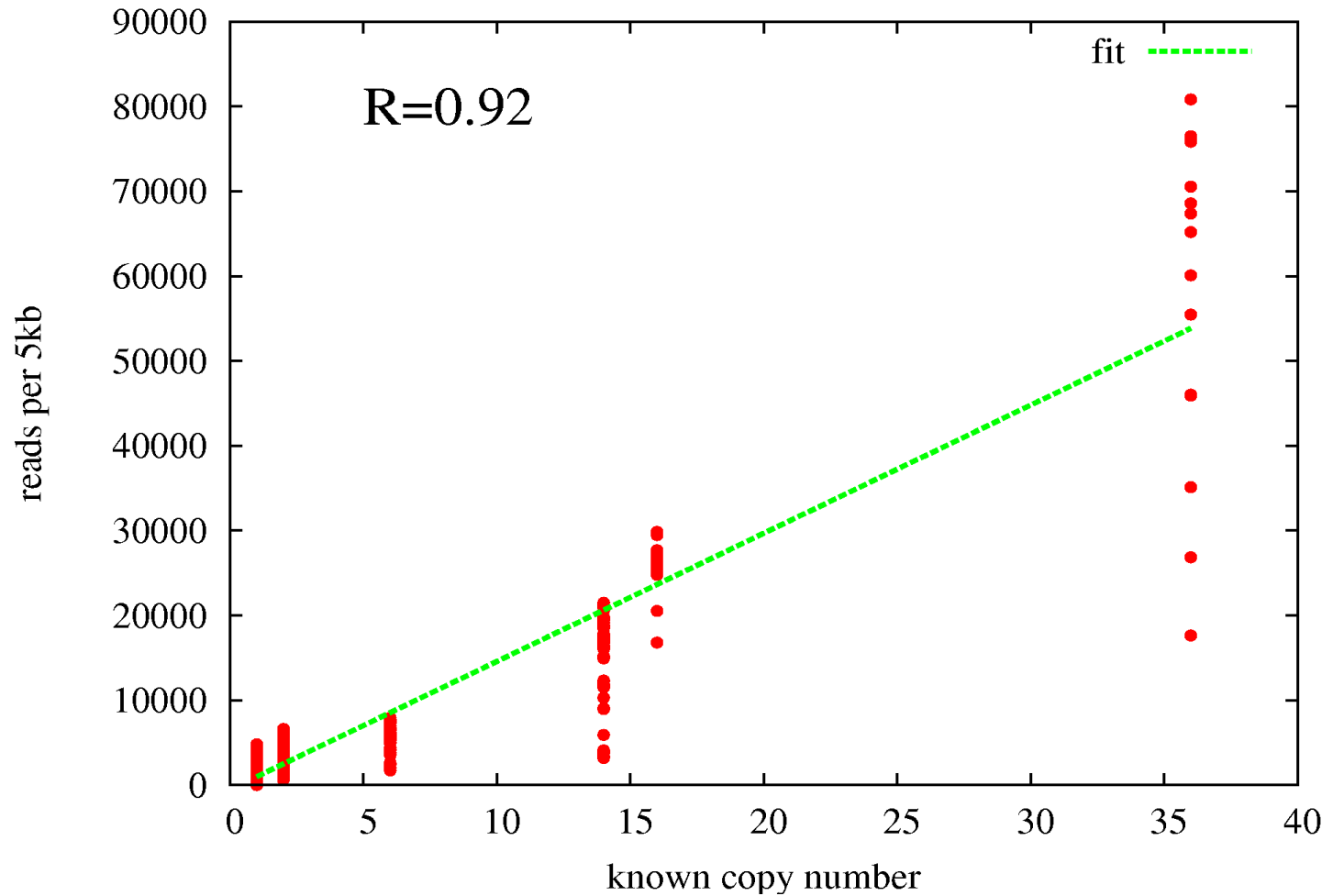


multiple
mapping

unique
mapping



Read depth - Copy number correlation

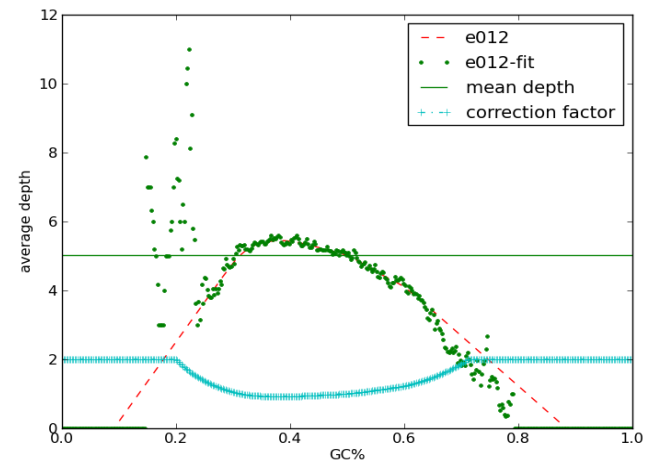
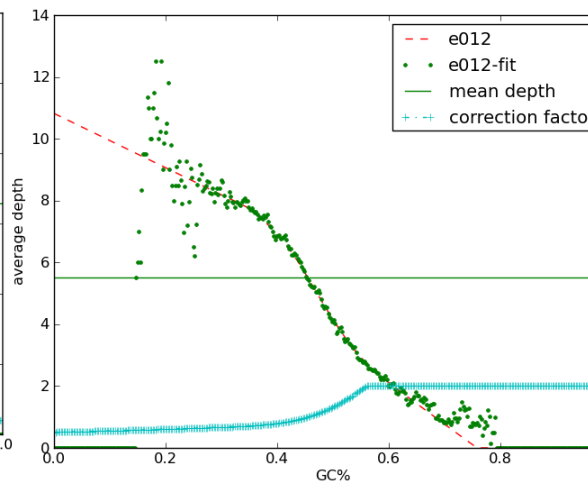
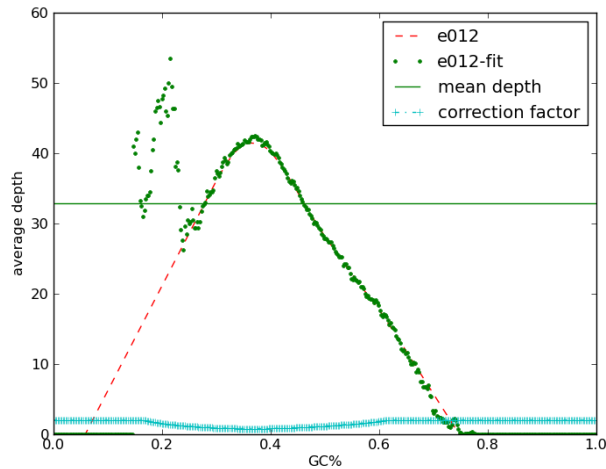


WSSD-HTS: mrCaNaVaR

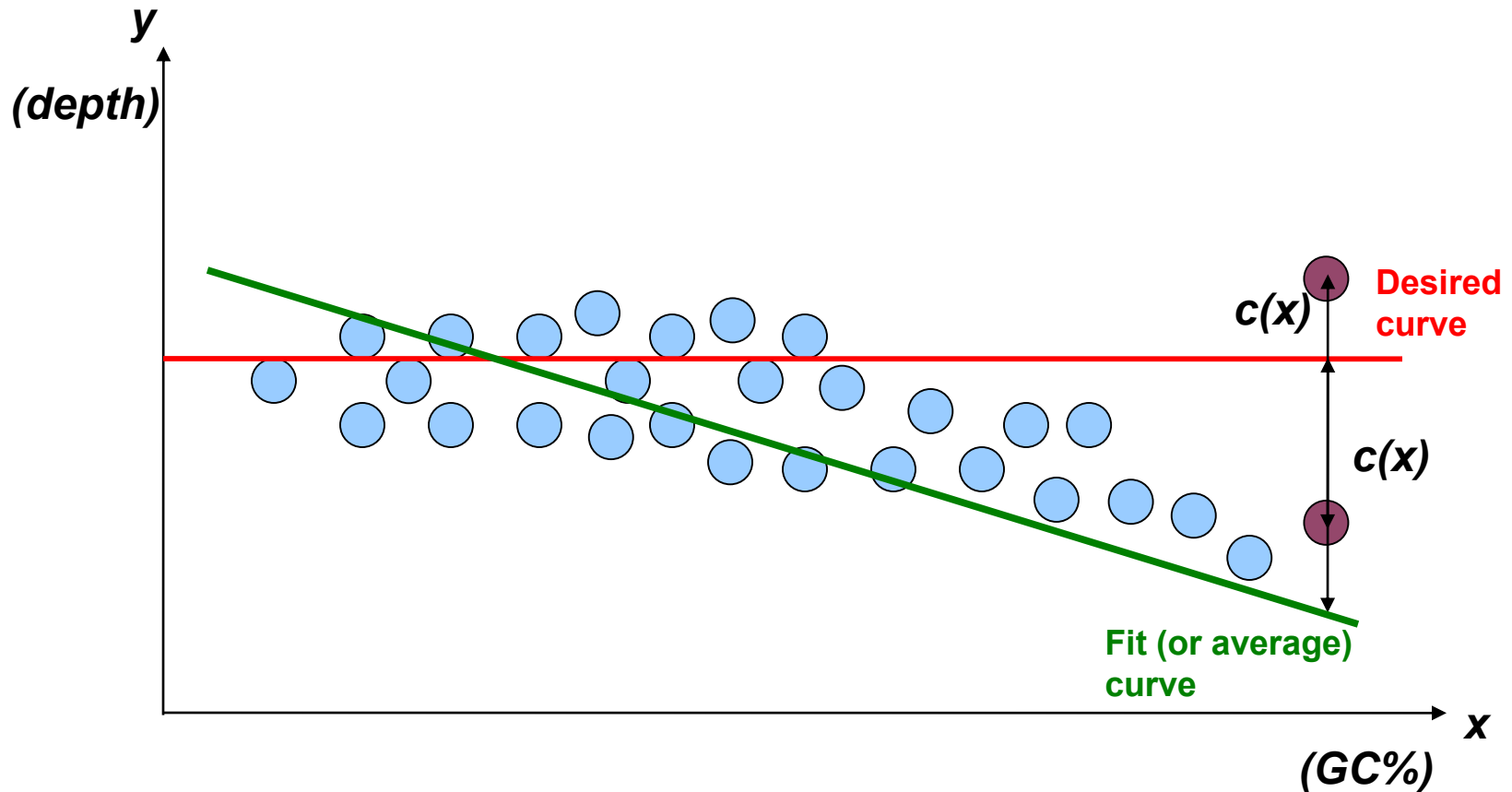
- HTS specific problems
 - Short reads: MegaBLAST is replaced by mrFAST / mrsFAST
 - Common repeats: all repeats need to be masked
 - GC % bias needs to be fixed
- Improvement
 - Absolute copy number detection in 1 kb non-overlapping windows
 - Genotyping highly identical paralogs

Read depth distribution

- Read depth doesn't really follow Poisson distribution
 - Biases against high and low GC %



GC% correction: LOESS



$$y' = y - c(x)$$
$$c(x) = f(x) - e(x)$$

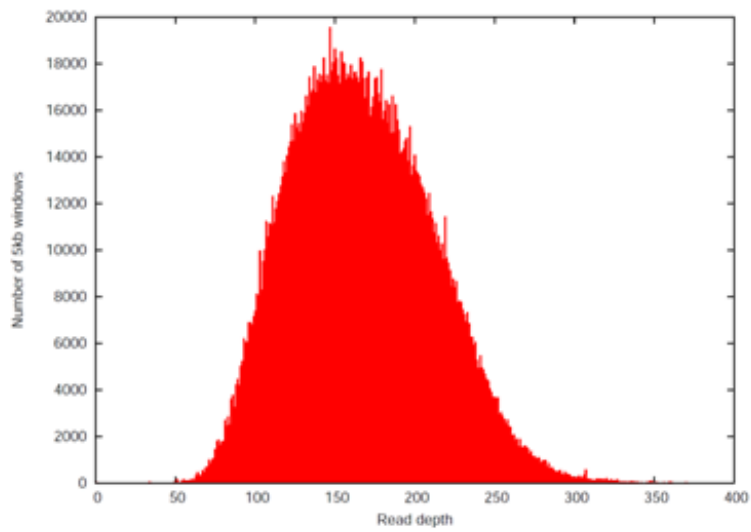
GC% correction (modified LOESS)

$$k_{gc} = \mu_{total} / \mu_{gc}$$

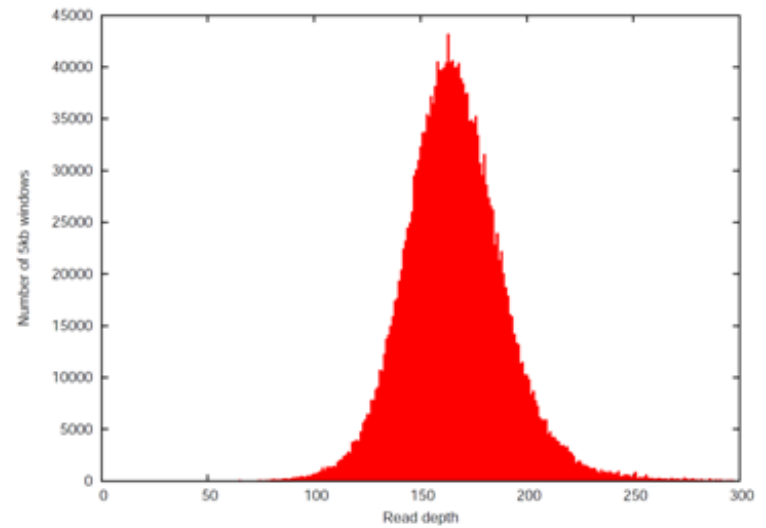
$$d'_{gc} = d_{gc} k_{gc}$$

The version in SegSeq and CNVnator

GC% correction

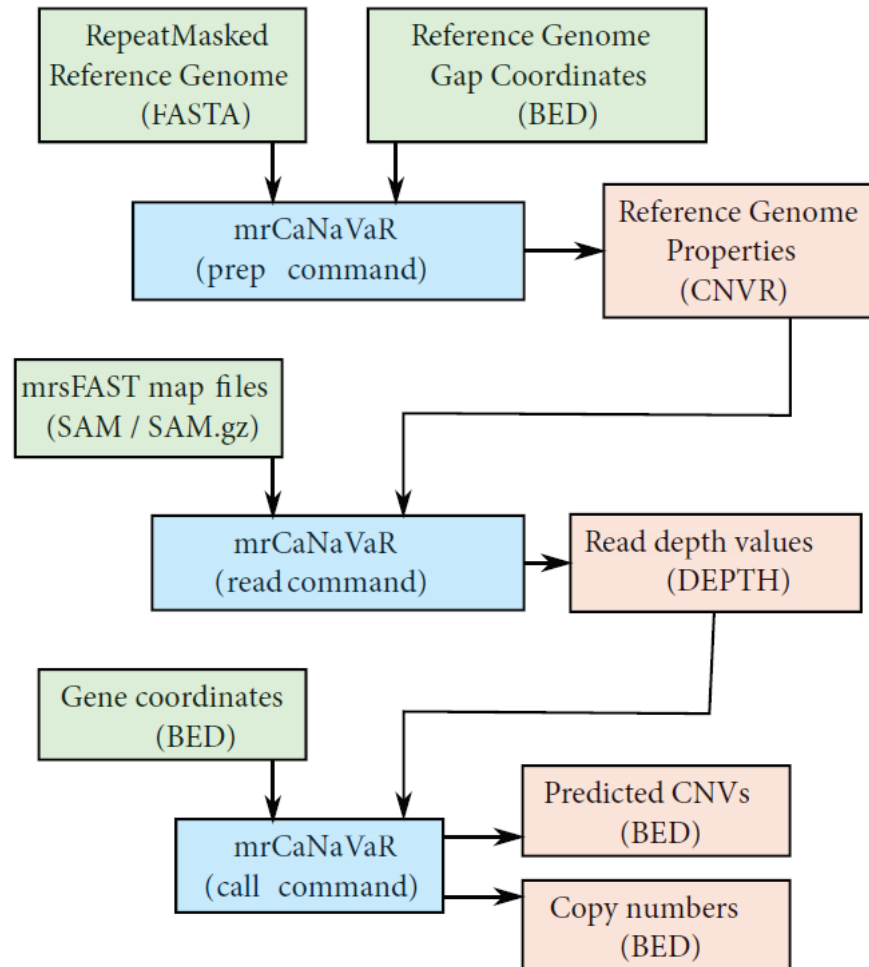


Before GC correction

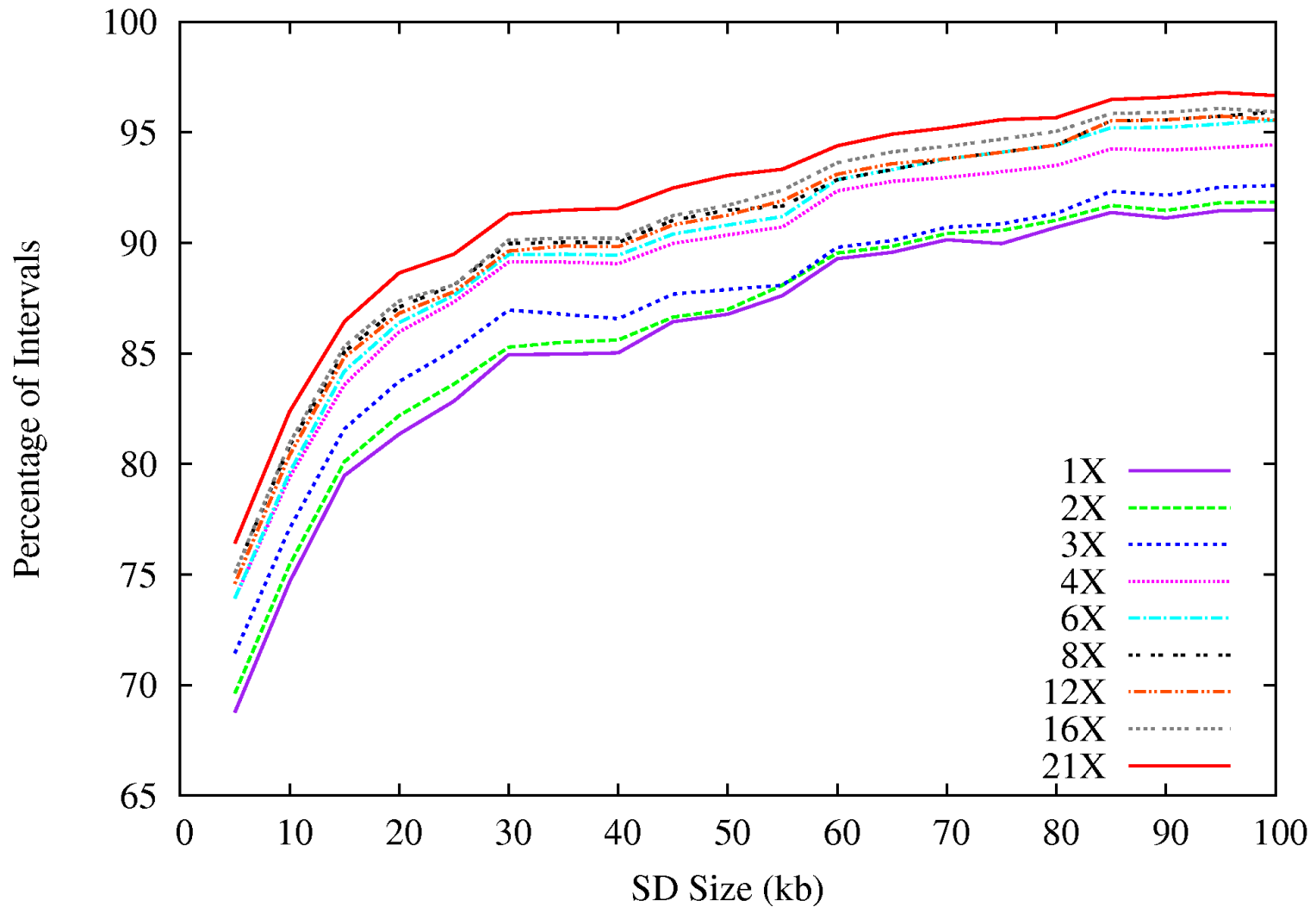


After GC correction

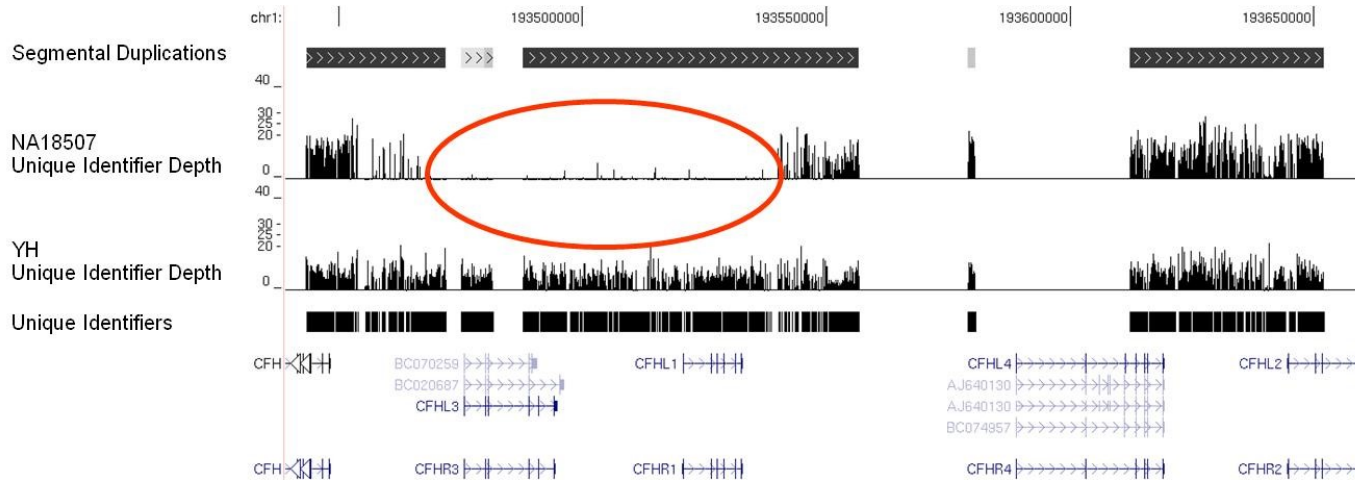
WSSD-HTS: mrCaNaVaR



Sequence coverage and detection power



Differentiating Paralogous Genes

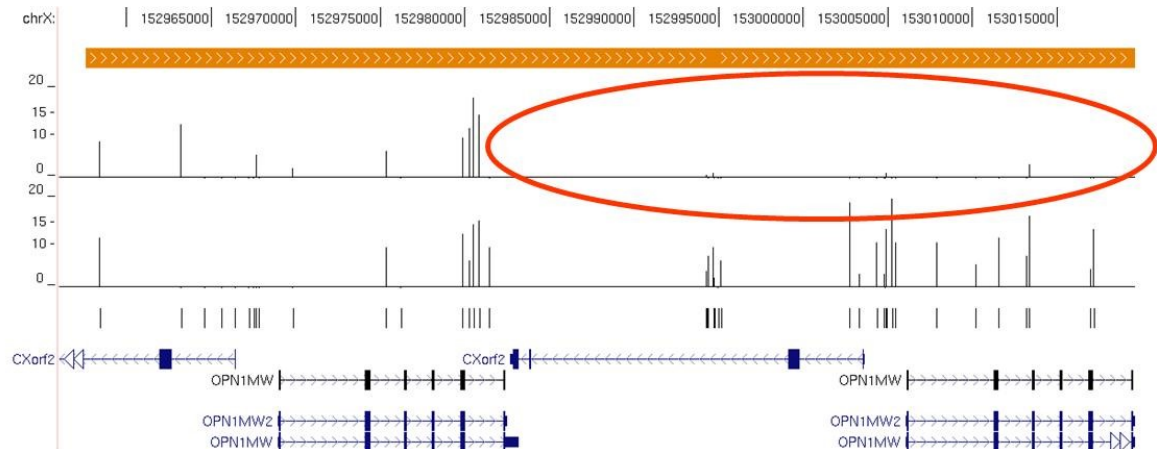


Associated with psoriasis and Crohn's disease

CFHR

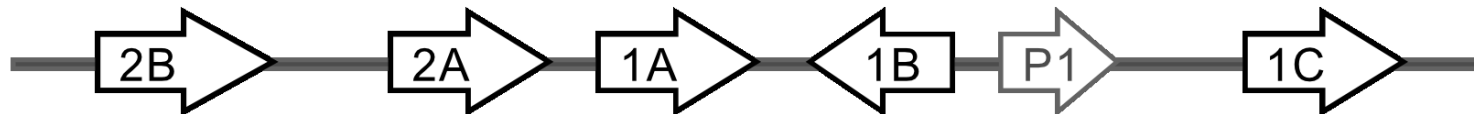
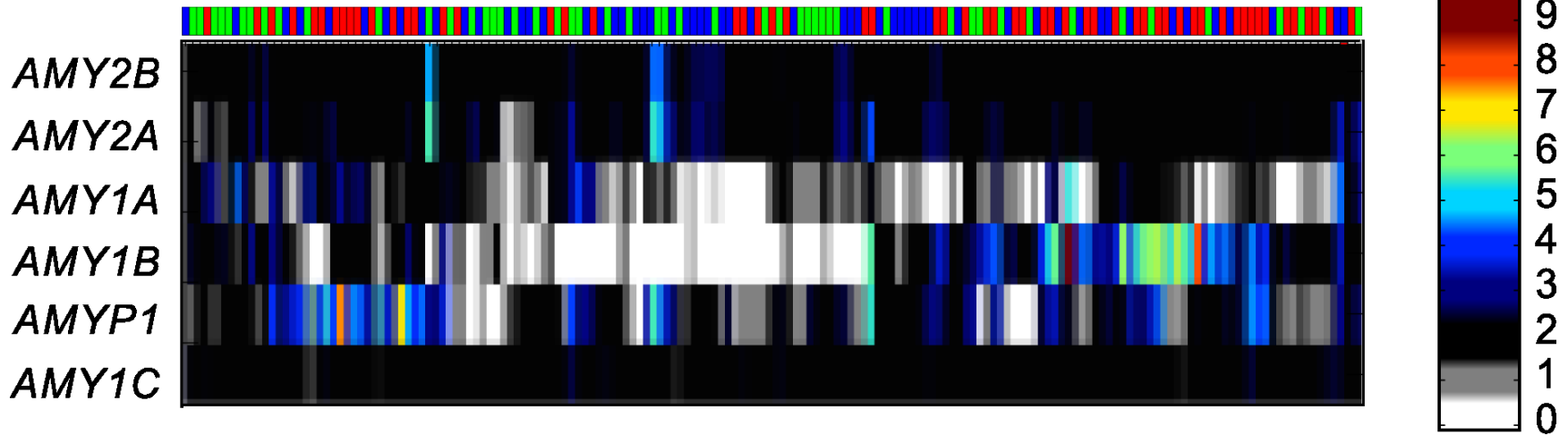
Associated with color blindness

opsin



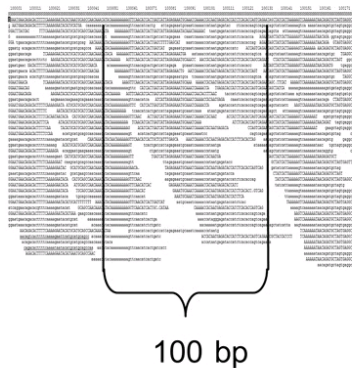
Singly Unique Identifiers (SUNs)

Copy 1 AT**A**CTAGGCATATAATATCCGACGATATACATATA**G**ATGTTAG
Copy 2 ATGCTAGGCAT**G**TAATATCCGACG**A**CATACATATACATGTTAG
Copy 3 AT**A**CTAGGCATATA**A**CATCCGACGATATACATATACATGTTAG
Copy 4 ATGCTA**C**GCATATAATATCC**C**ACGATATACATATACATGTTAG
Copy 5 ATGCTA**C**GCATATAATATCCGACGATATACATATACAT**G**ATAG
Copy 6 AT**A**CTAGGCAT**G**TAATATCCGACGATATAC**- -**ATACATGTTAG

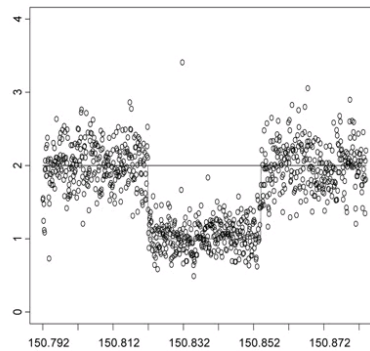


Event-Wise Testing (EWT)

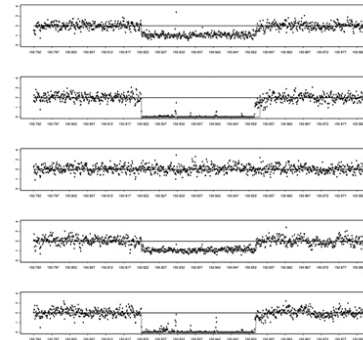
(A) Estimation of Read Depth



(B) Event detection



(C) Comparison of multiple genomes



Polymorphic events
(CNVs)

Monomorphic events
(e.g. segmental
duplications or rare
variants in the reference
genome)

- Unique mappings are used
- No masking
- Window size 100 bp
- Probabilistic analysis

Event-Wise Testing (EWT)

- Read counts are converted to Z score:
 - $z_i = (RC_i - \mu_i) / \sigma_i$
- Upper and lower tail probabilities
 - $p_i^U = P(Z > z_i)$
 - $p_i^L = P(Z < z_i)$
- Unusual events for interval A , $l = |A|$; L number of windows in chromosome; FPR: false positive rate

$$\max \{ p_i^U \mid i \in A \} < \left(\frac{FPR}{L/l} \right)^{\frac{1}{l}}$$

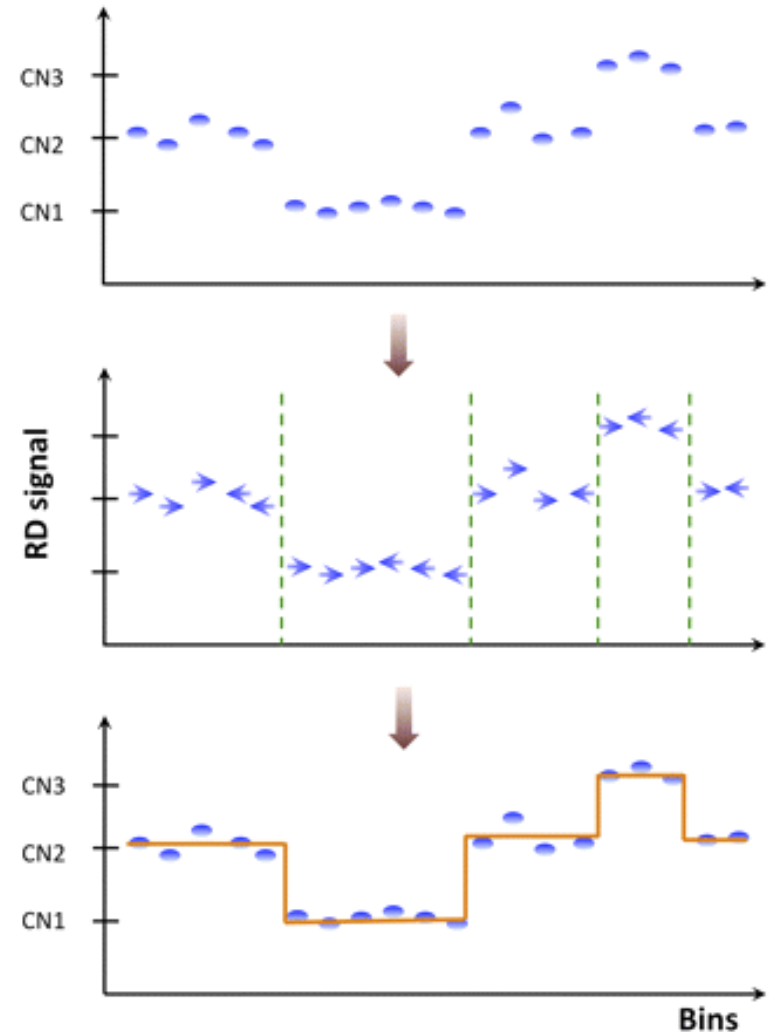
Duplication

$$\max \{ p_i^L \mid i \in A \} < \left(\frac{FPR}{L/l} \right)^{\frac{1}{l}}$$

Deletion

CNVnator

- Unique mappings
- Mappings with low MAPQ are discarded
- Partitioning is based on mean-shift technique developed for image processing



CNVs with exome sequencing

- Exome sequencing: capture only coding exons from DNA and sequence
 - 1.5% of total genome
 - Good for protein coding variants but misses regulatory sequence, introns, etc.
 - Whole genome sequencing generates random data, but exome does not
 - Capture efficiency changes for *every* exon (n~200,000)
 - CNVs from exomes: ExomeCNV, FREEC, CoNIFER
-

READ PAIRS + SPLIT READS
