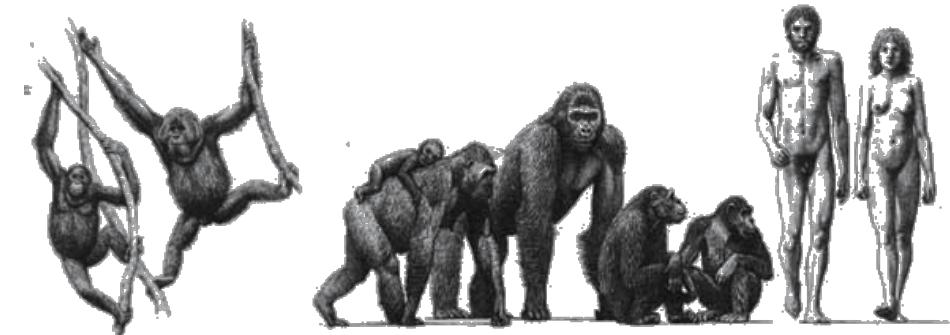


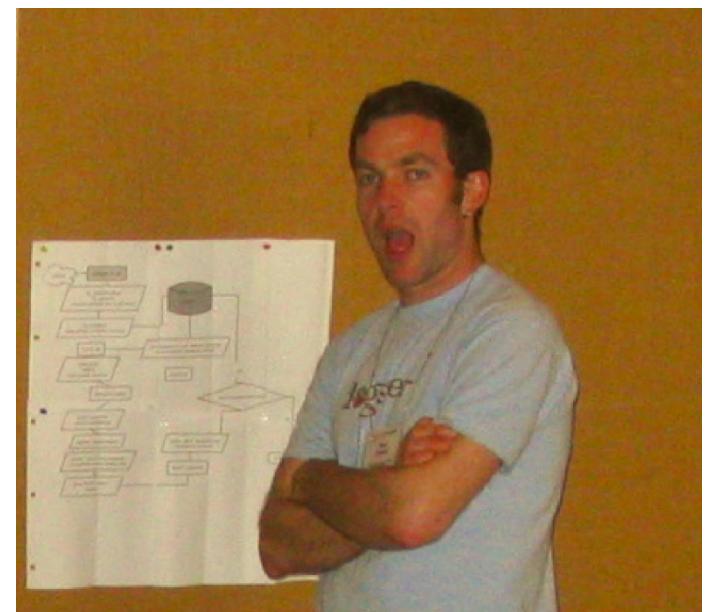
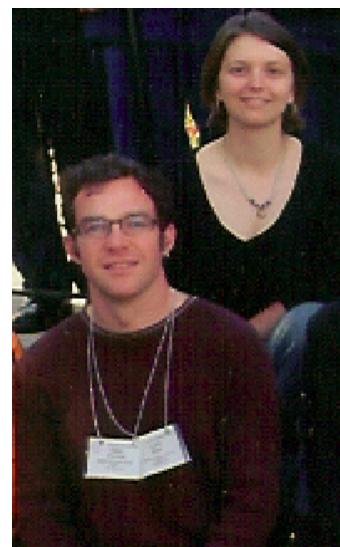
Structural variation

Programming for Biology
CSH, October 2014

Tomas Marques-Bonet
ICREA Research Professor
Institut de Biologia Evolutiva

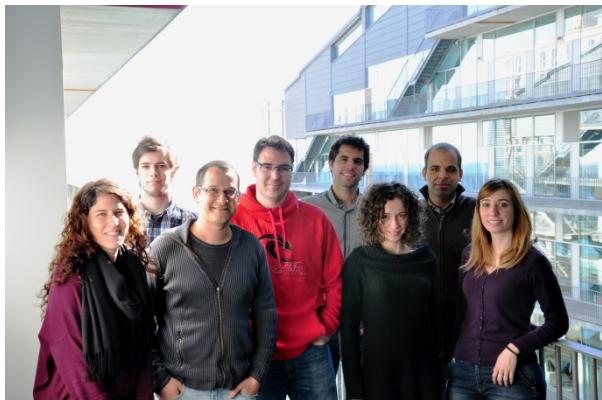


11 years from my 1st CSH!



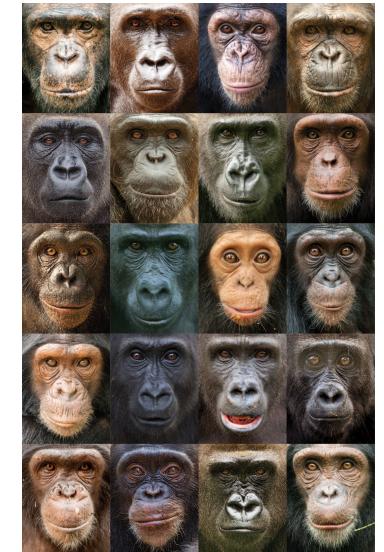
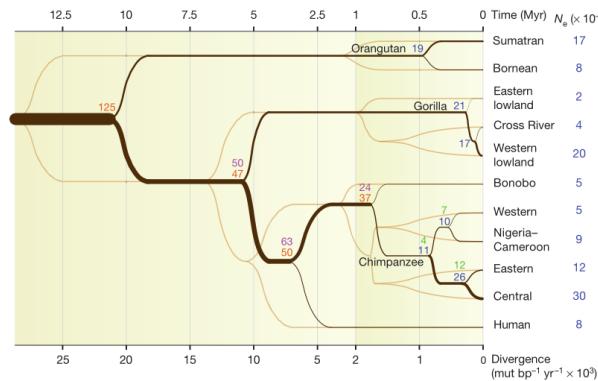
Who we are?

- Evolutionary genomics
 - Barcelona, Biomedical Research Park



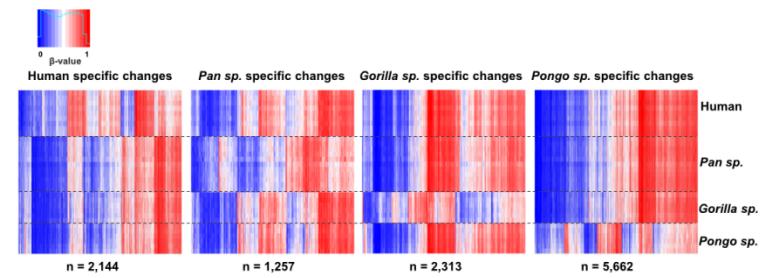
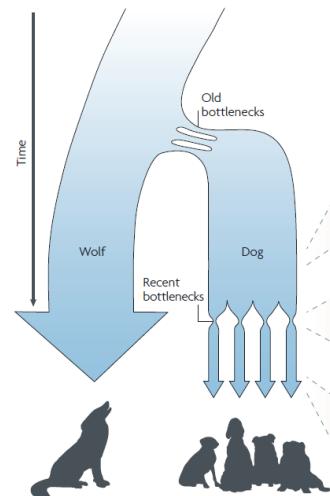
What do we do?

- Natural selection on human evolution



- Transcriptome and Epigenetics in Primates

- Canid evolution



Continuum of Genomic Variation

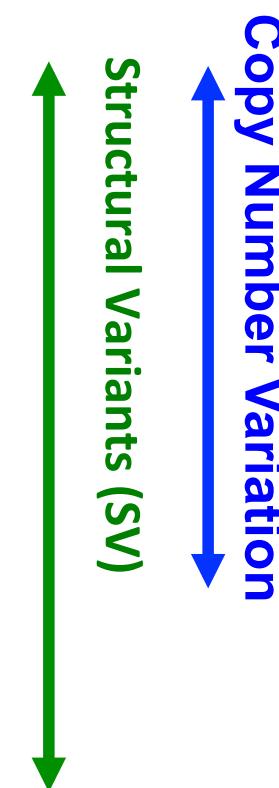
*Single
nucleotide*

- Single base-pair changes

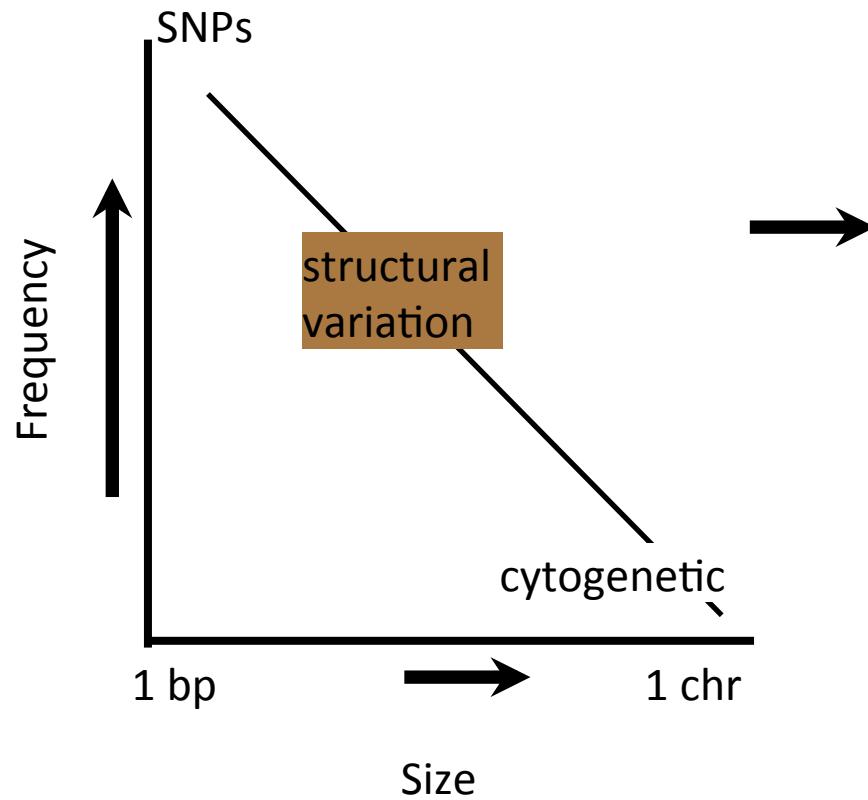


- Cpg Methylation
- Small insertions/deletions
- Mobile elements
- Large-scale genomic copy number variation (>10 kb)
- Local Rearrangements

Chromosome • Chromosomal variation

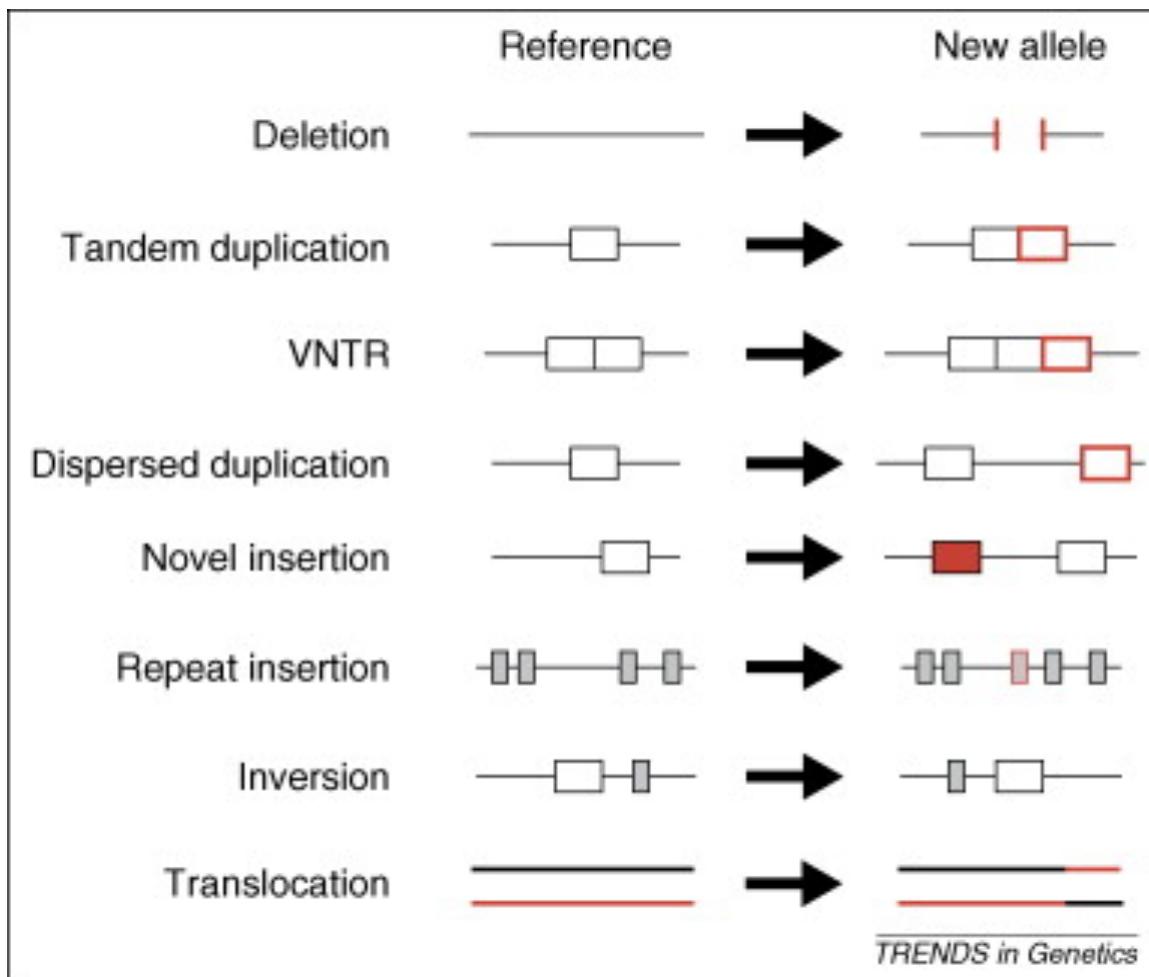


Genomic Structural Variation



- Gene-altering, *e.g.* immune response, drug metabolism
- Abundant: majority of human heterozygosity
- Numerous plausible functional consequences

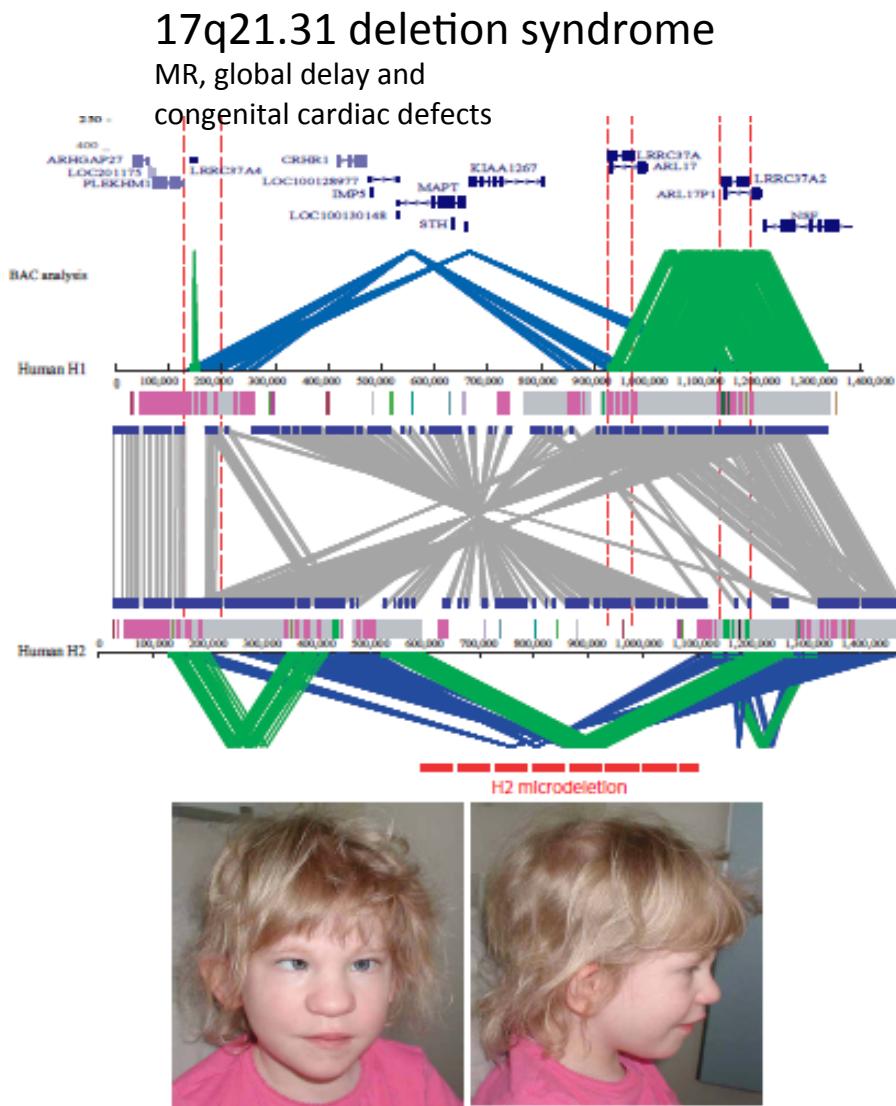
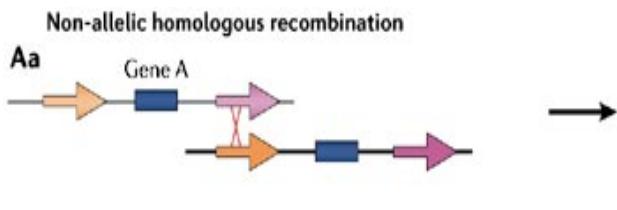
Types of Structural Variation



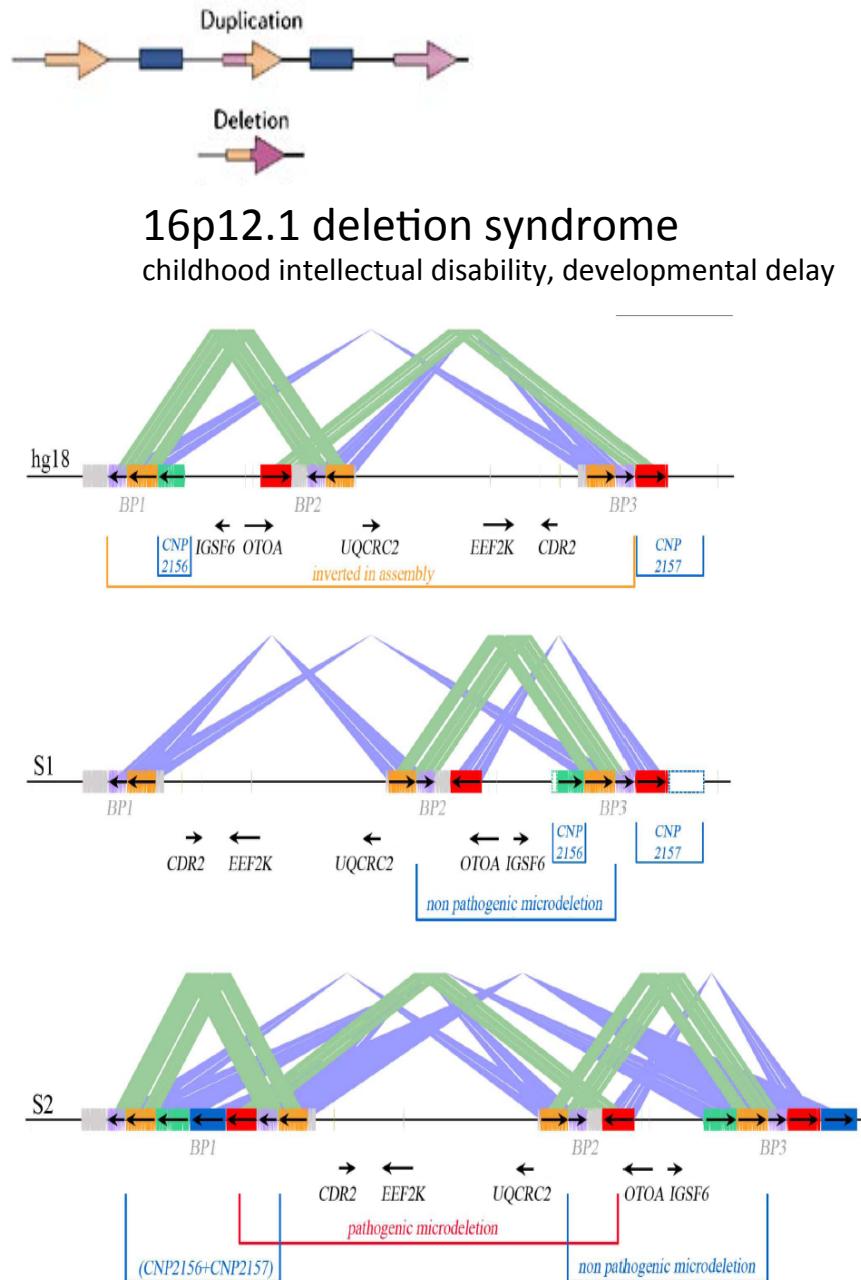
Hurles et al. 2008

Why Study Structural Variation?

- Common in “normal” human genomes-- major cause of phenotypic variation
- Common in certain diseases, particularly cancer
- Now showing up in rare disease; autism, schizophrenia



Zody et al. Nature Genetics (2008)



Antonacci et al. Nature Genetics (2010)

Challenges of CNV studies

- Often involves repeated regions
- Rearrangements are complex
- Can involve highly repetitive elements

Methods to Find SVs

Experimental approach

ArrayCGH (SNP based and genomic)

Sequence based

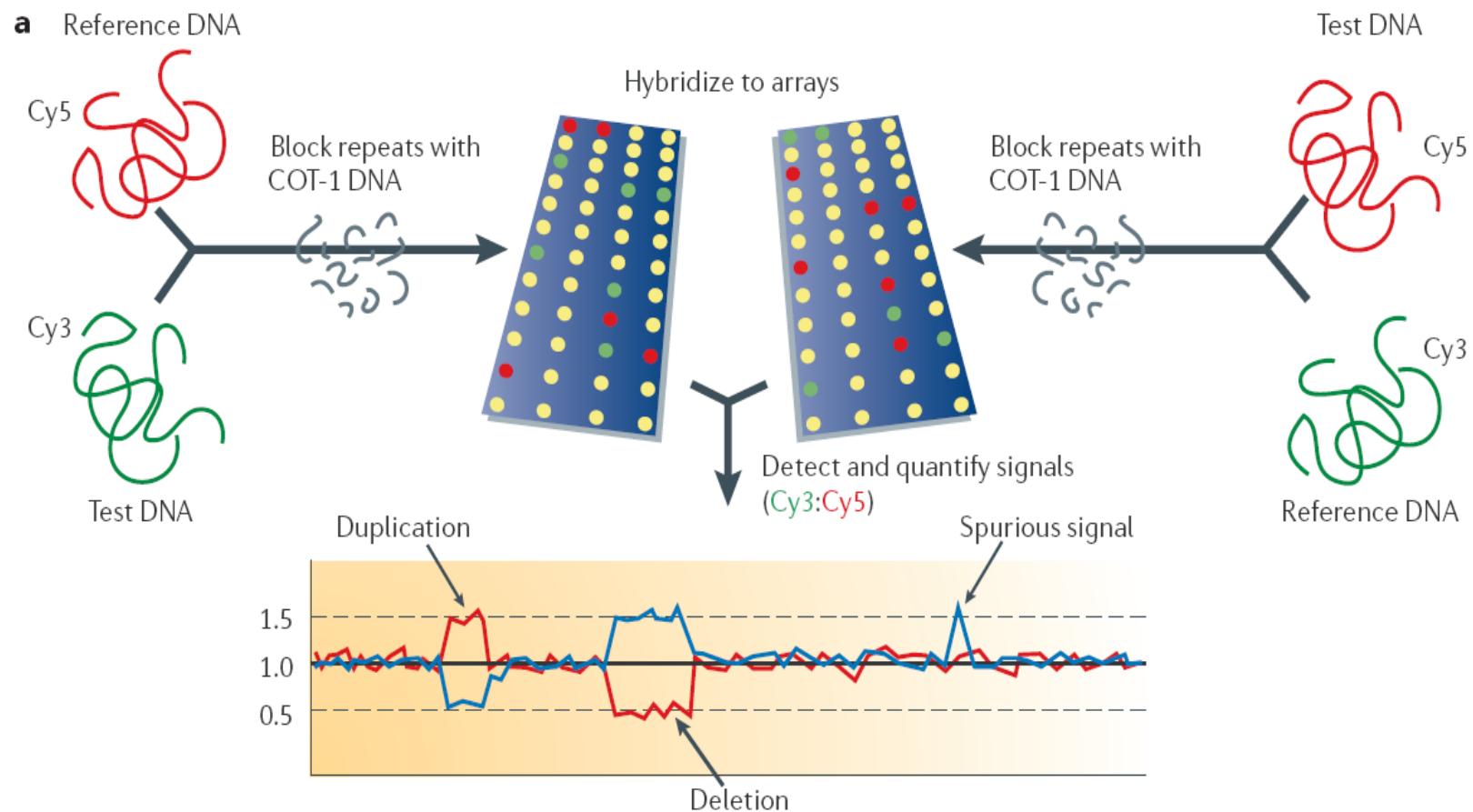
Local and *de novo* assembly

Read pair analysis

Read depth analysis

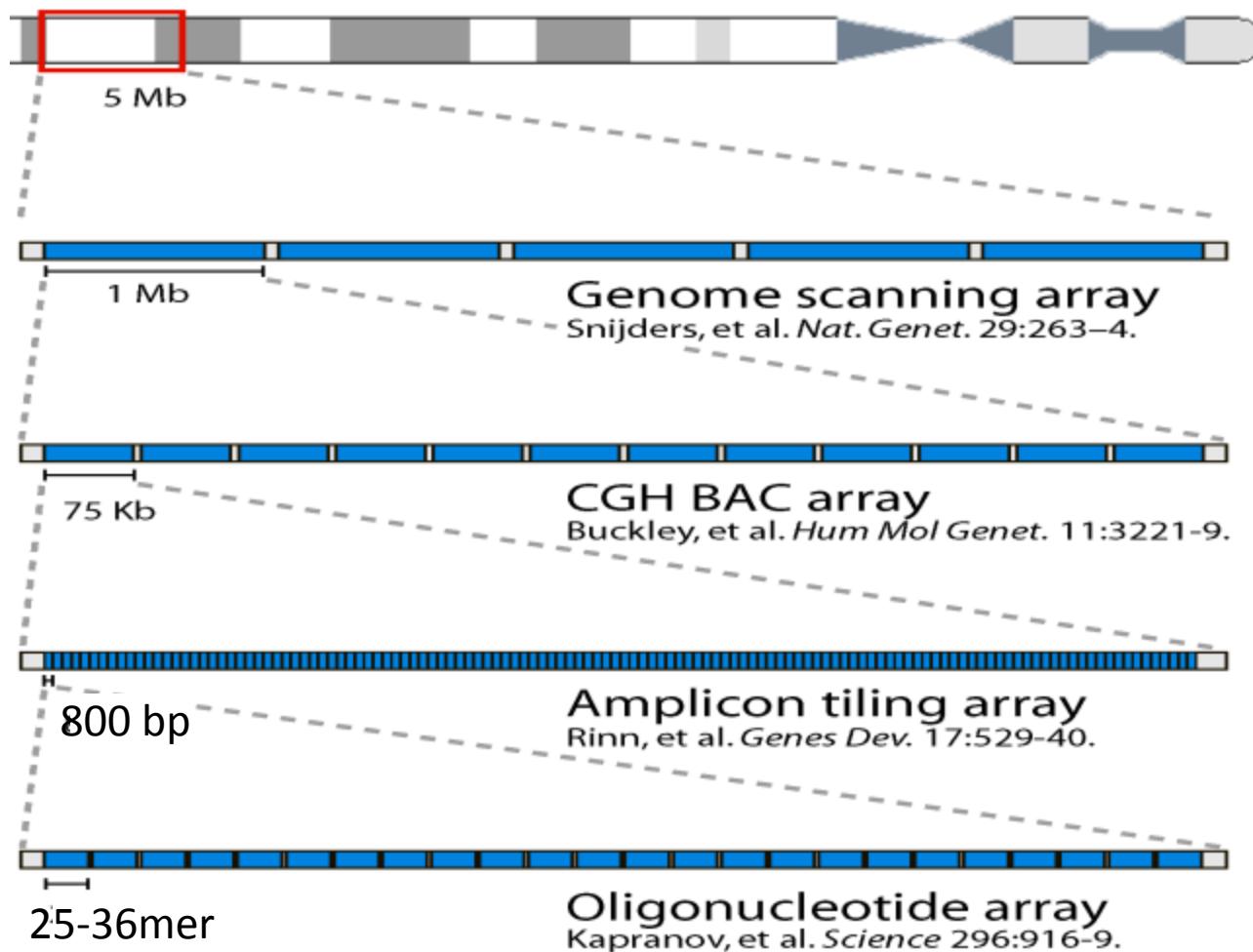
Split read analysis

METHOD 1: Copy Number Variation: Array Comparative Genomic Hybridization



Modified: Feuk et al. *Nat Rev Genet* 2006

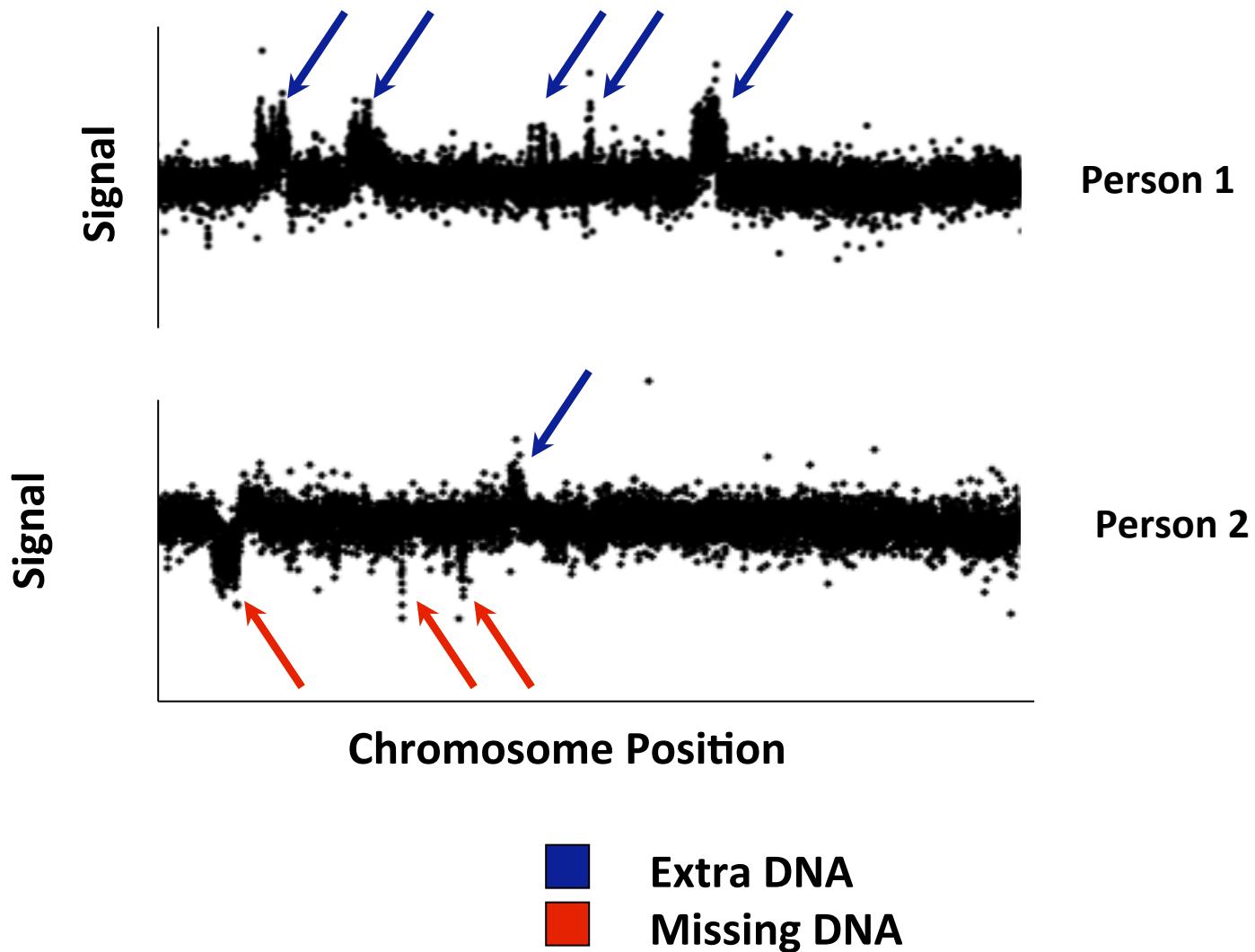
Genome Tiling Arrays



Typical Analysis Procedure

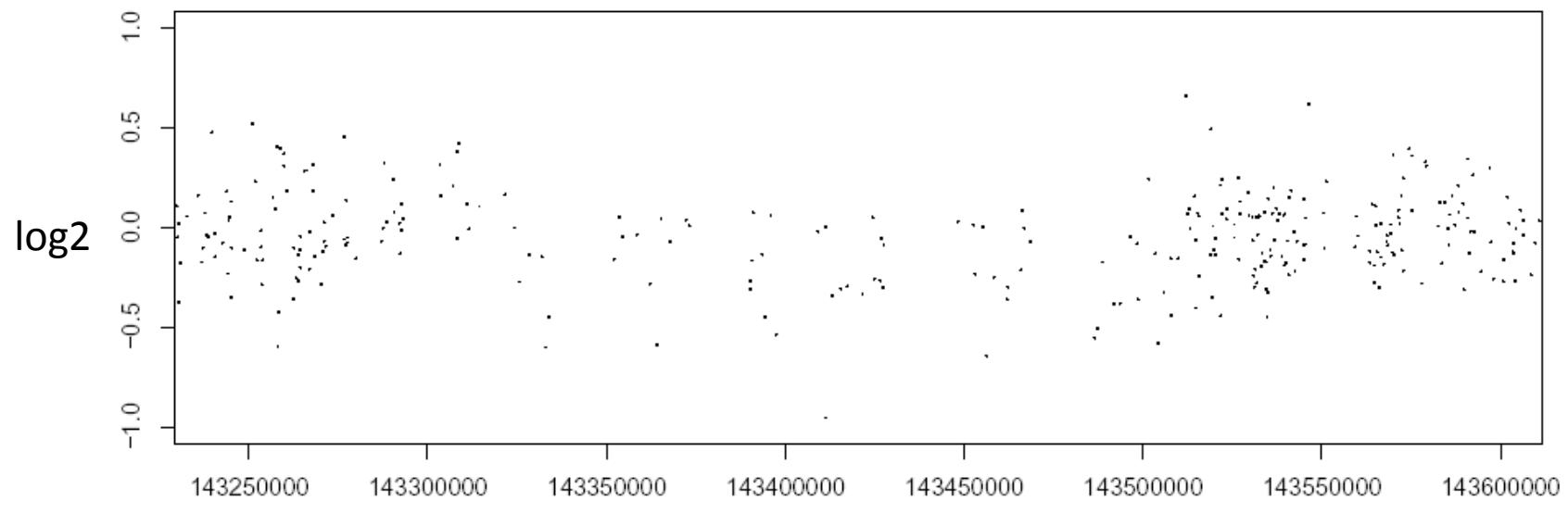
- For each probe, calculate a log2 ratio of test/reference
 - Log2 serves to center values around 0
 - Hemizygous deletion in test: $\log_2(\text{test}/\text{reference}) = \log_2(1/2) = -1$
 - Duplication in test:
 $\log_2(\text{test}/\text{reference}) = \log_2(3/2) = 0.59$
 - Homozygous duplication:
 $\log_2(\text{test}/\text{reference}) = \log_2(4/2) = 1$

Copy Number Variations in the Human Genome

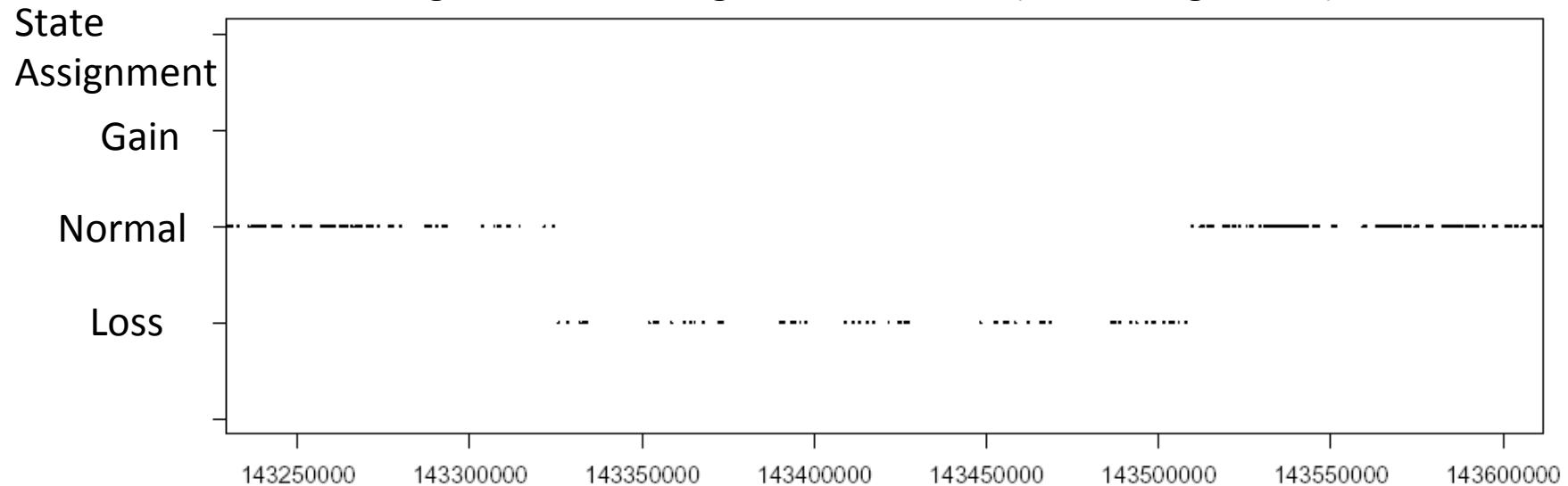


3-State HMM

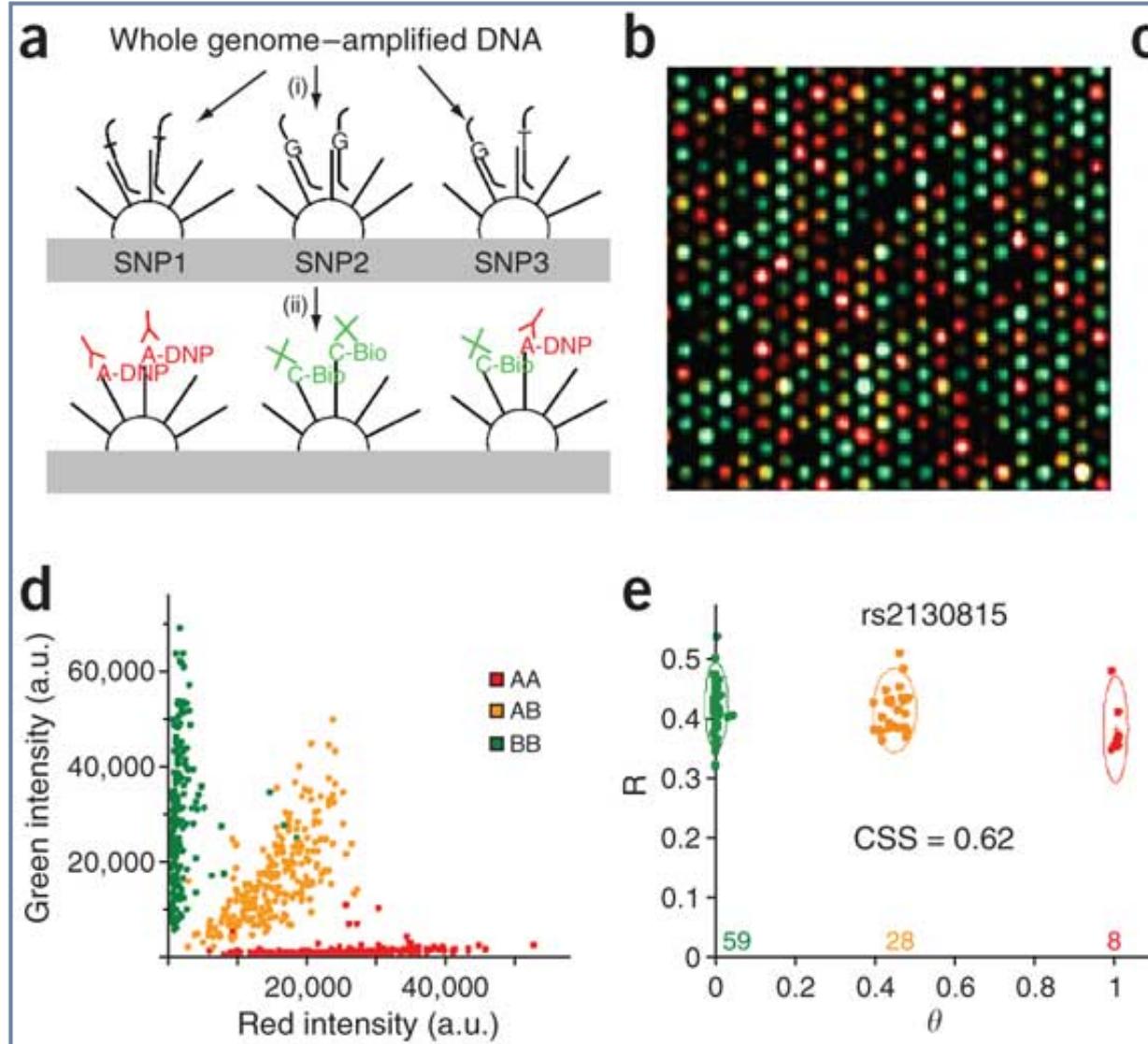
chr07.64797_143243594_143597470



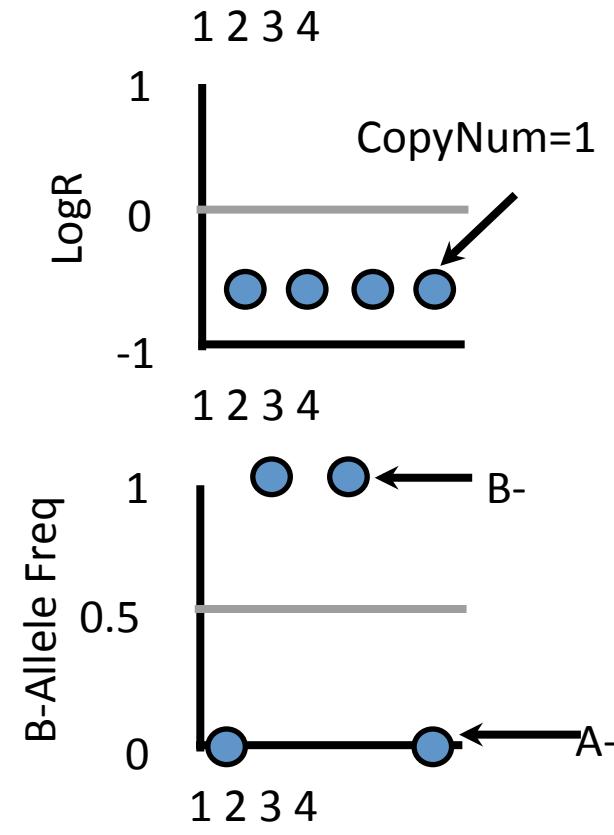
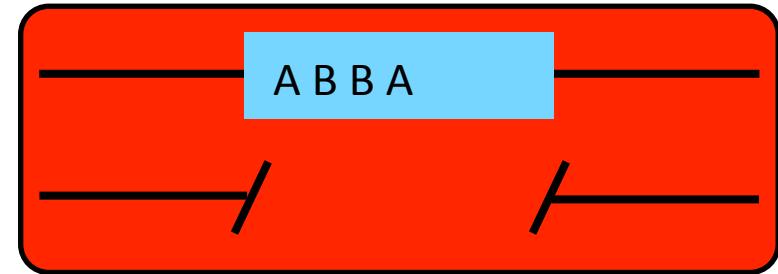
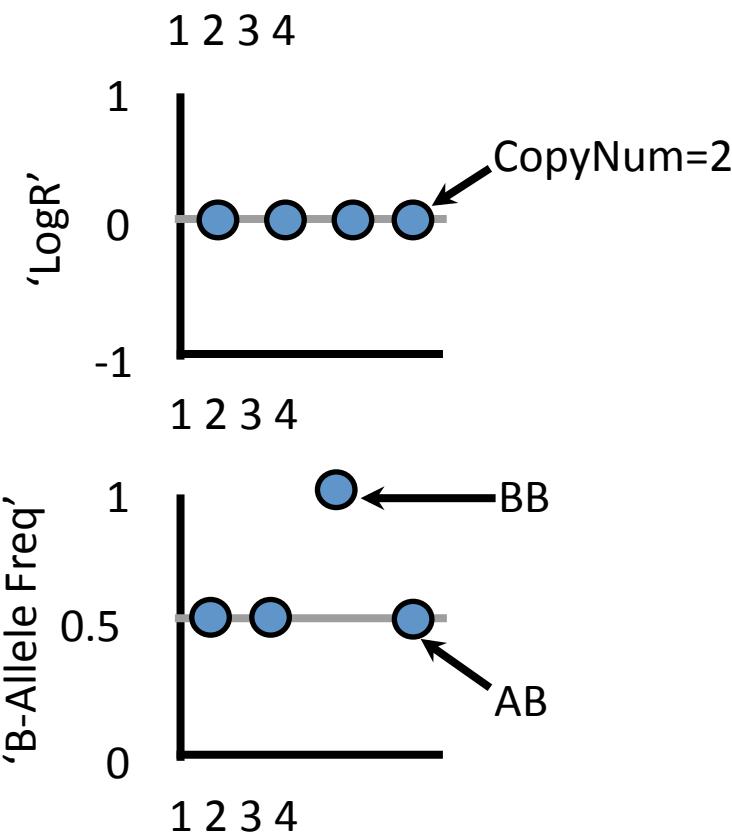
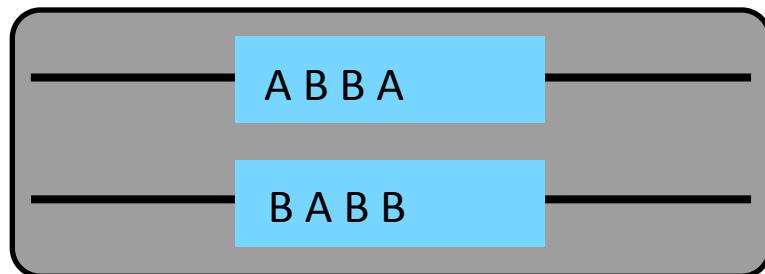
Segmentation using a 3-state HMM (Viterbi Algorithm)

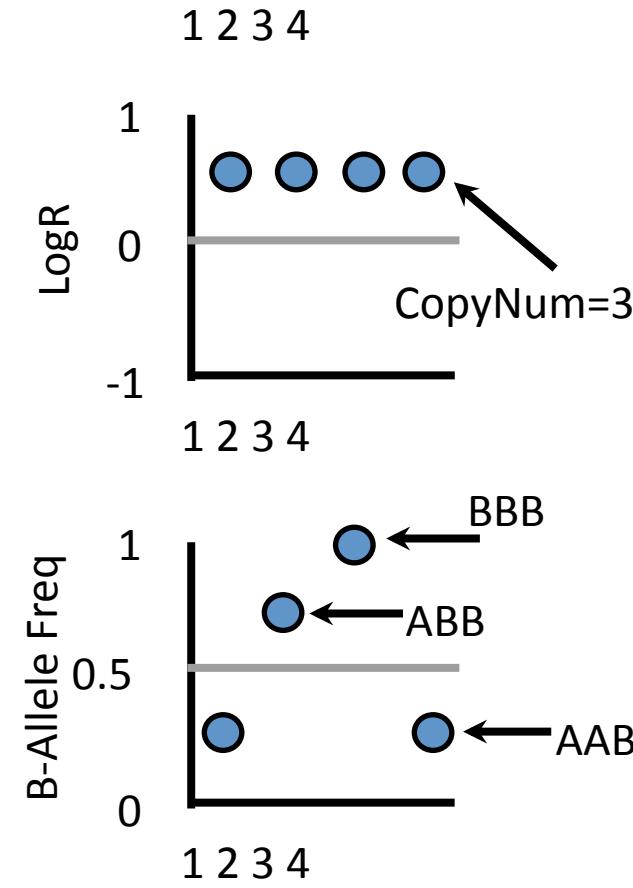
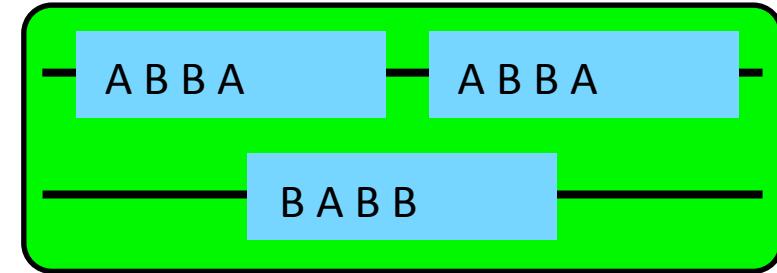
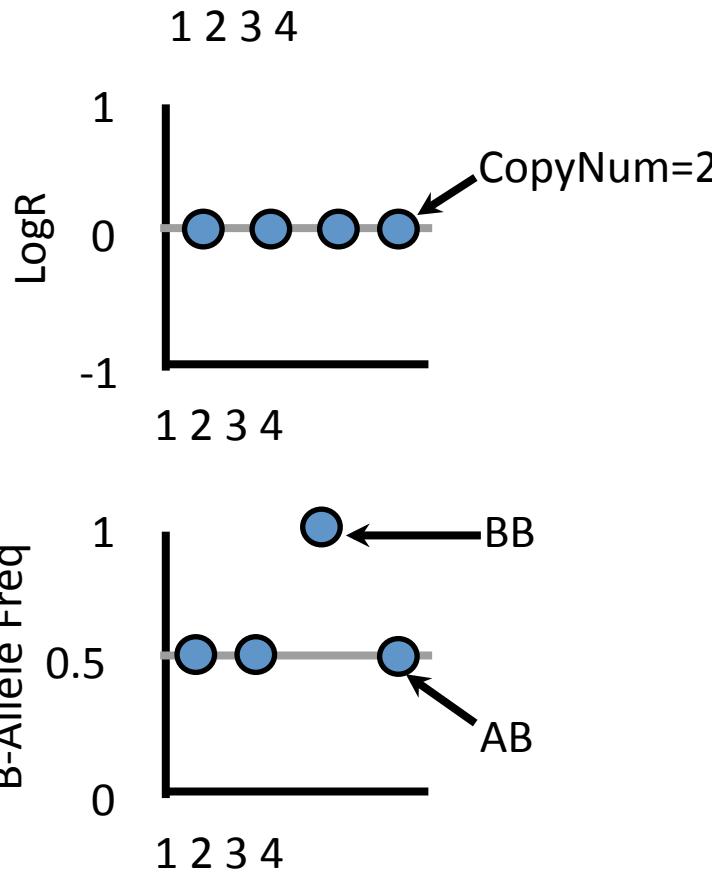
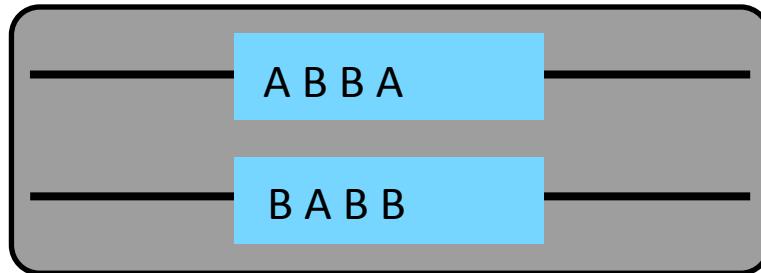


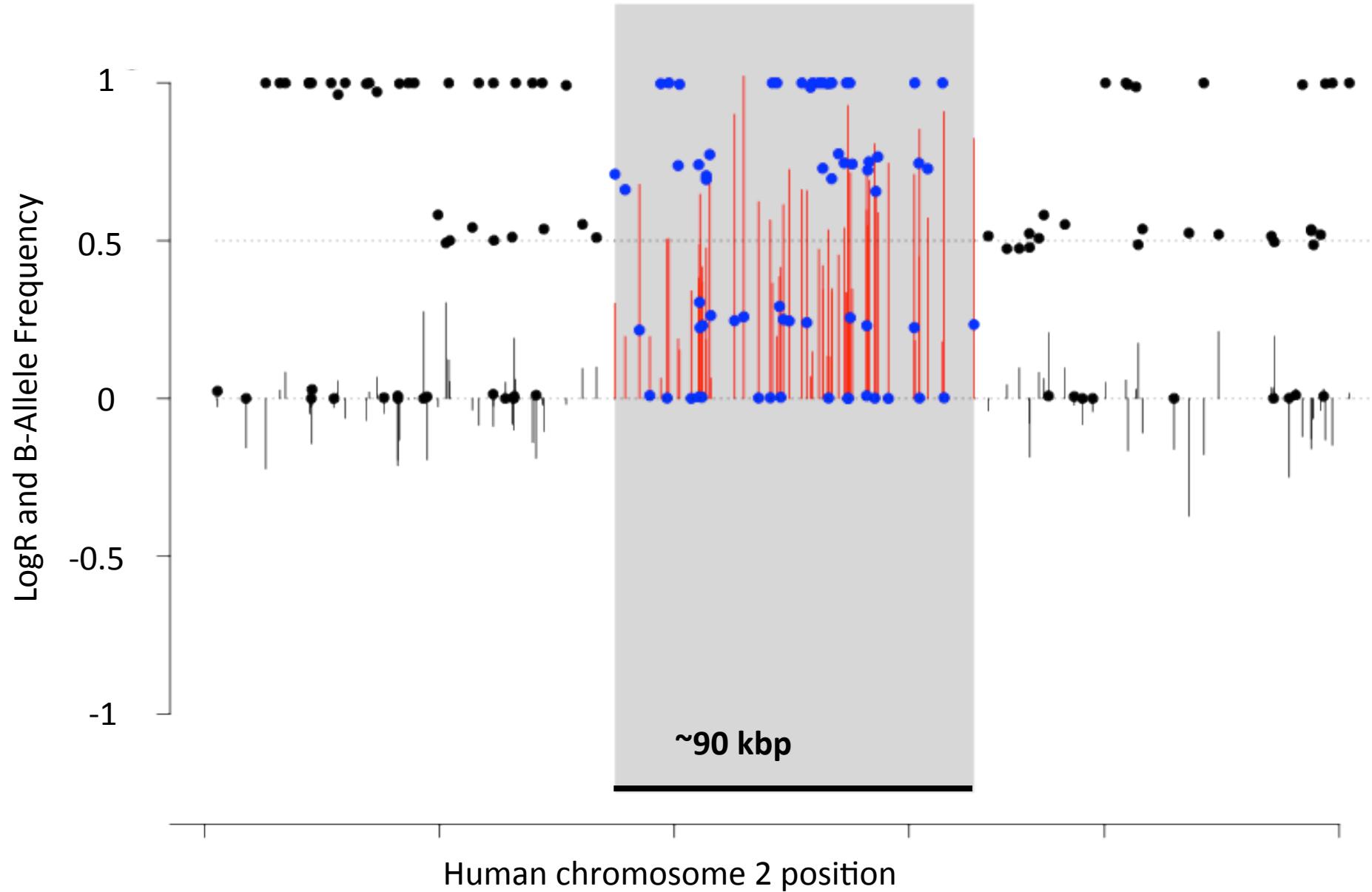
METHOD 2: Copy Number Variation: SNP genotyping Array



SNP Fluorescence-Based Deletion Discovery







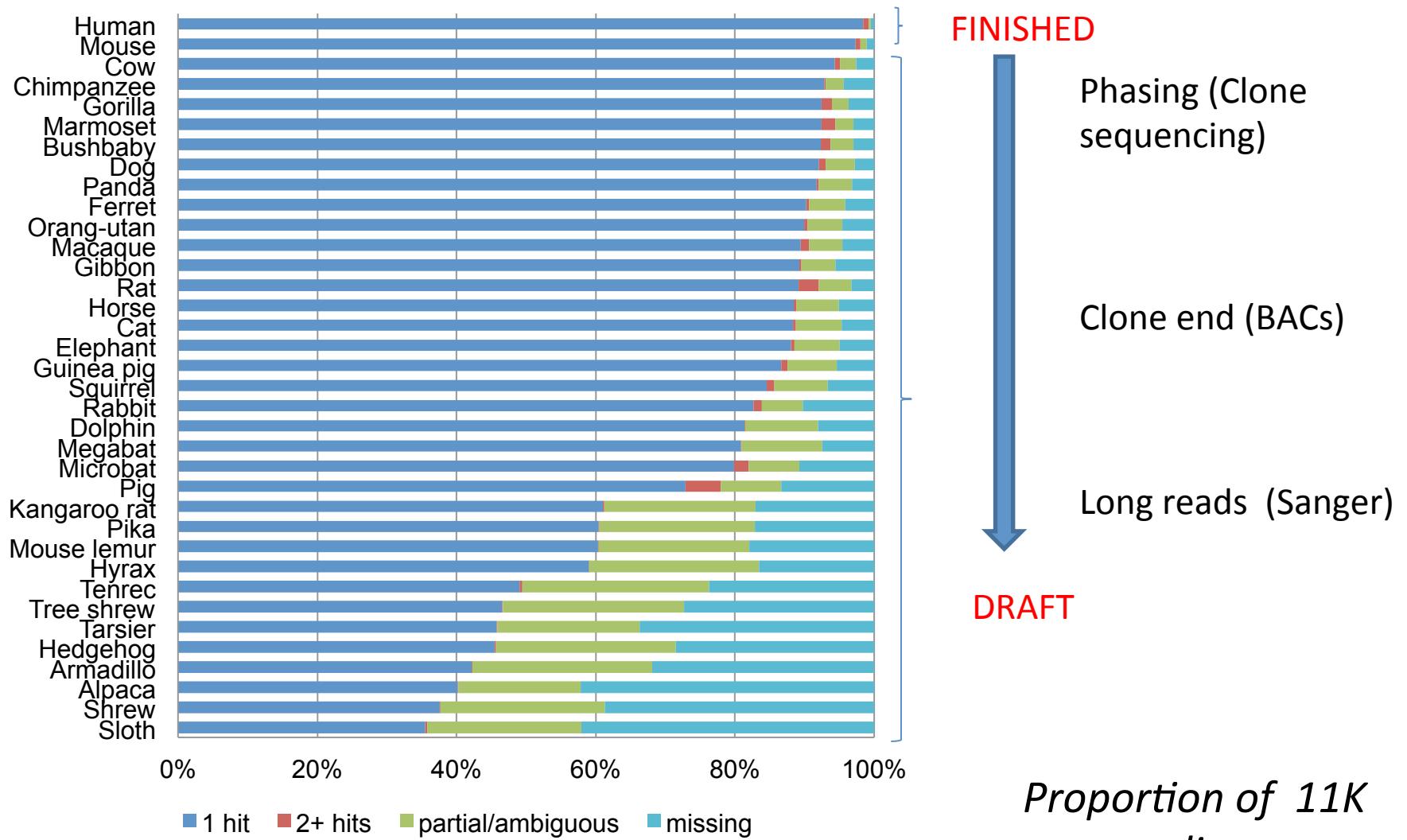
Sequencing Methods

- Going Backwards... Sanger, 454, Illumina.....
- CNV and SV are hotspots of research... but reality is:
 - Limitations of the methods
 - Indirect methods. ALL have problems!!
 - What do we want?
 - Clone sequencing/Phasing (**Moleculo?**)
 - Finish sequence and better assemblies (**PACBIO?**)

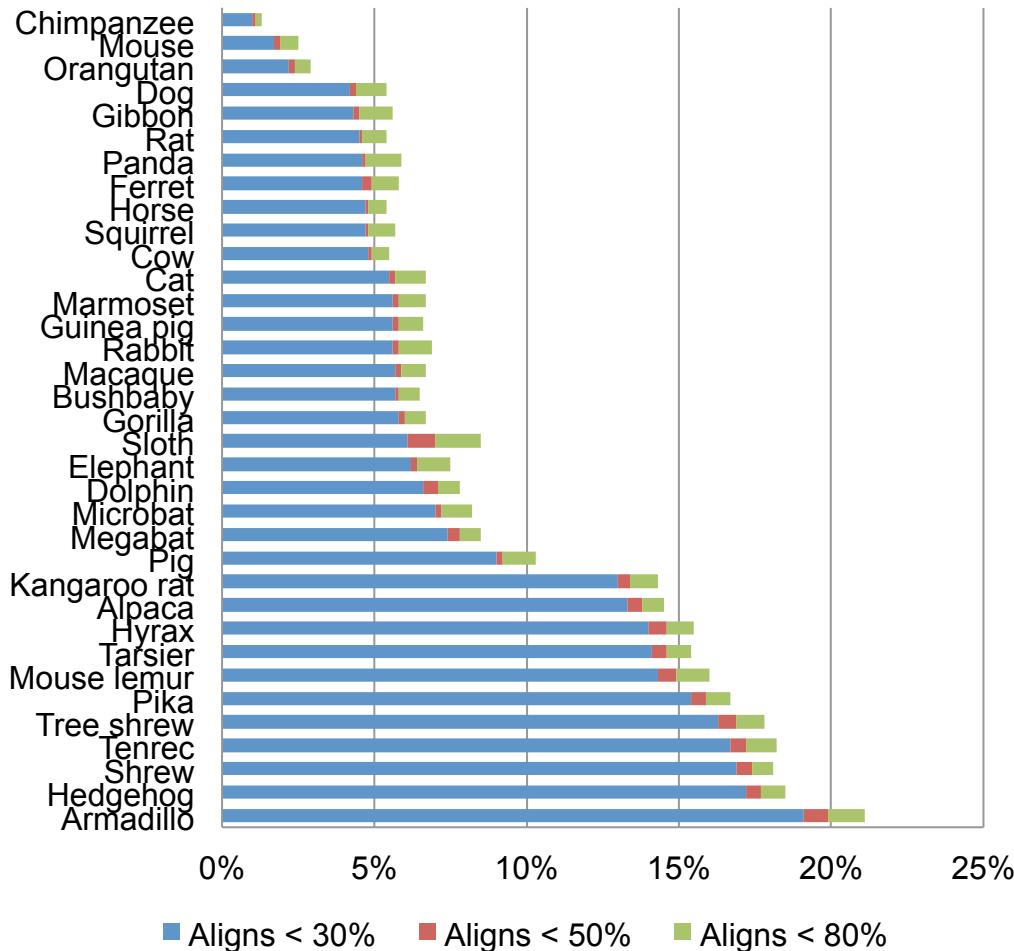
De novo assemblies

- Theory vs. Reality
- Most assemblies (even with Sanger technology!) are collapsed.

Quality of “old days” assemblies

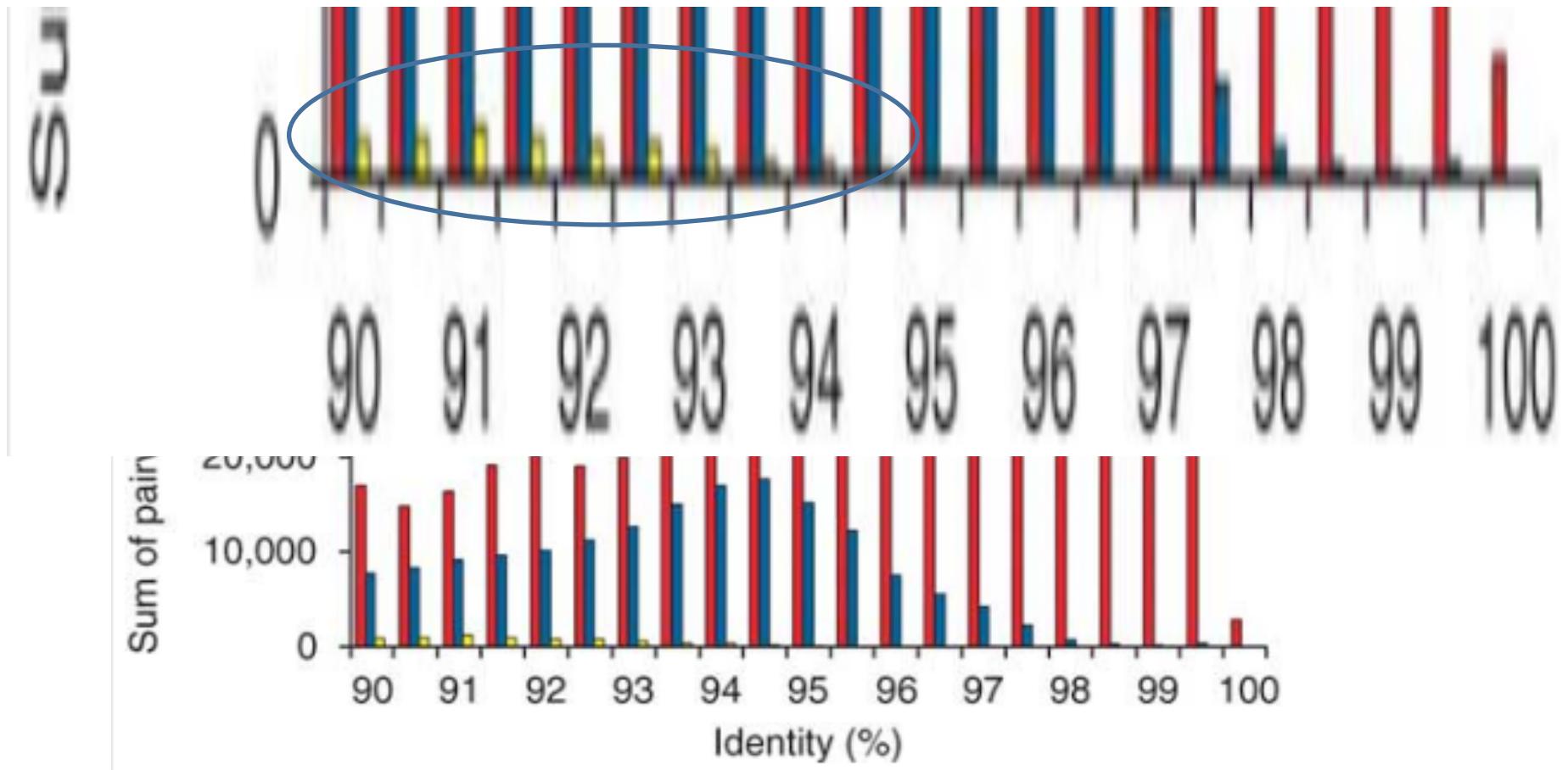


Quality of assemblies (II)



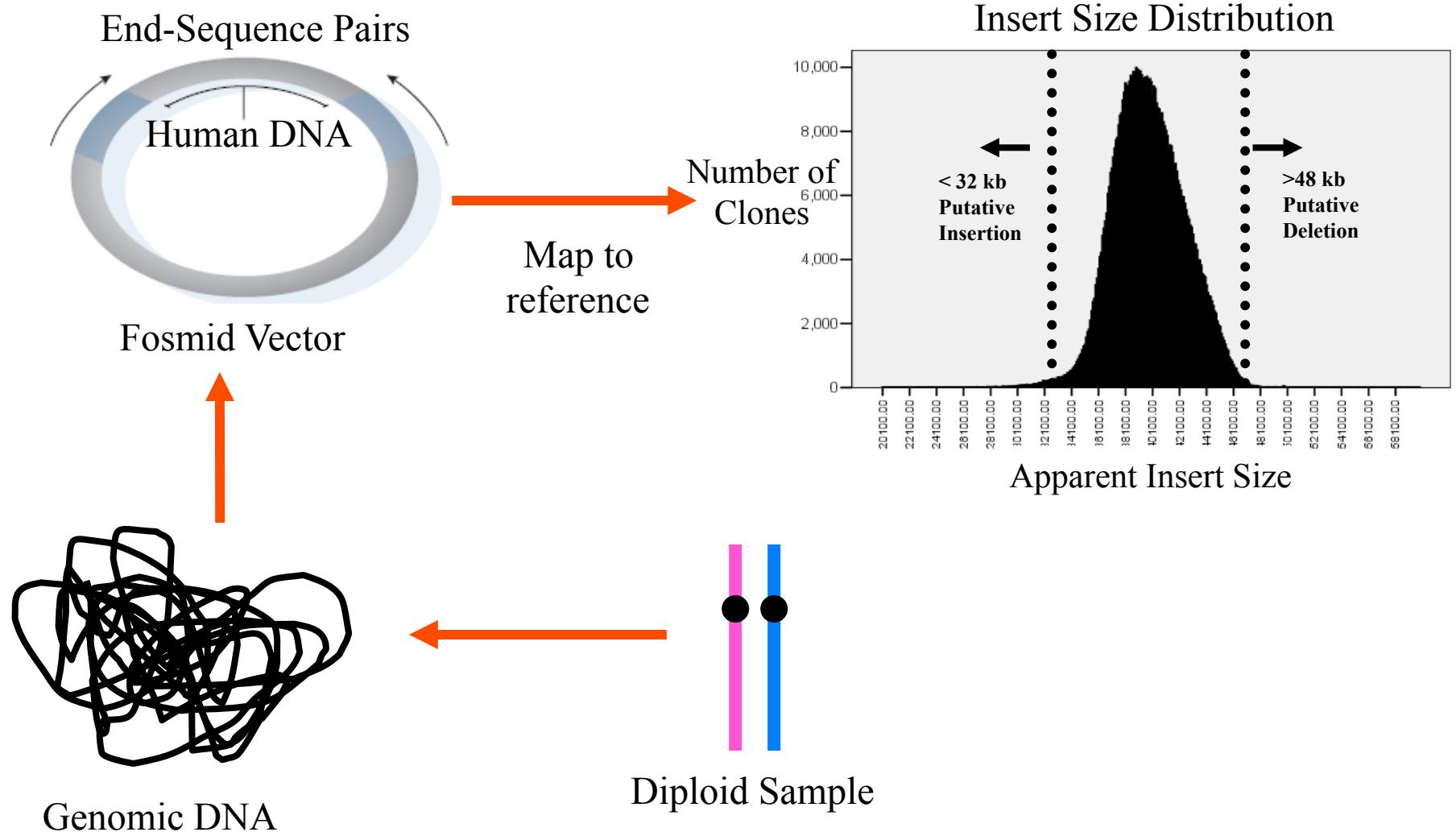
Incomplete
representation of
human genes

Limitations of NGS assemblies



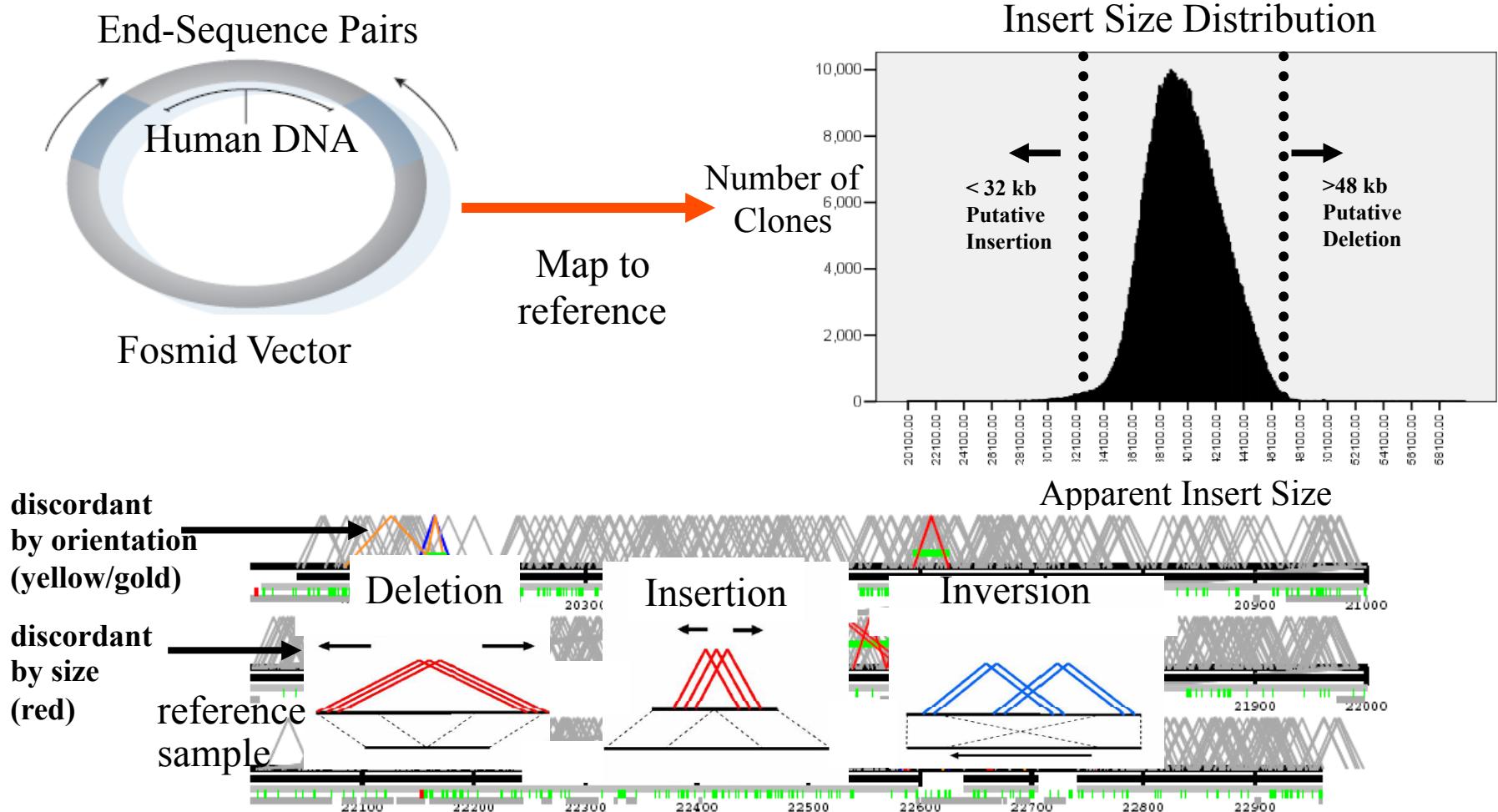
Alkan et al. Nature Methods 2010

Method 2: End-Sequence Pair (ESP) Analysis



Tuzun *et al.* (2005)

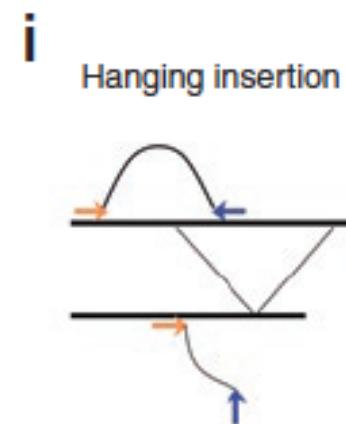
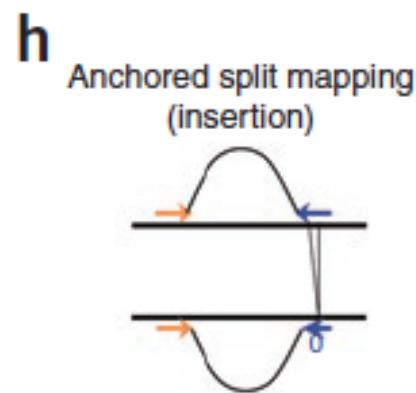
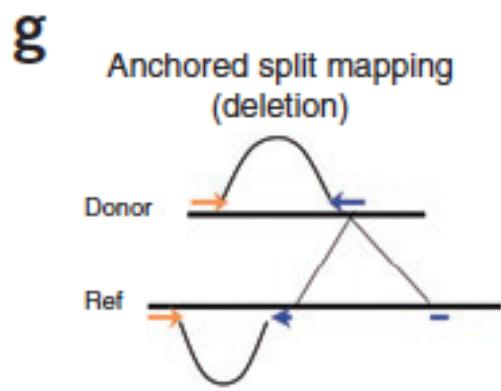
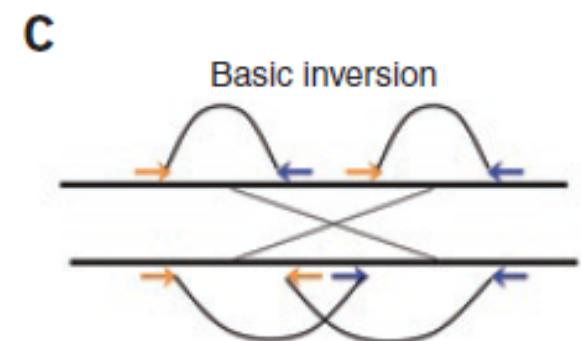
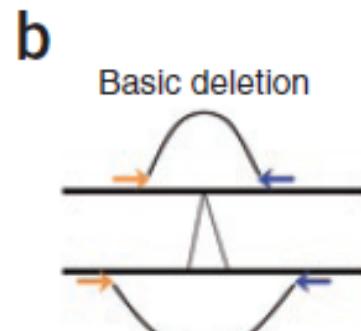
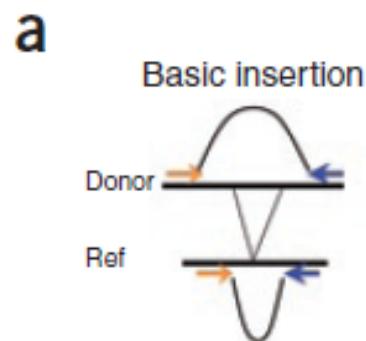
Method 2: End-Sequence Pair (ESP) Analysis



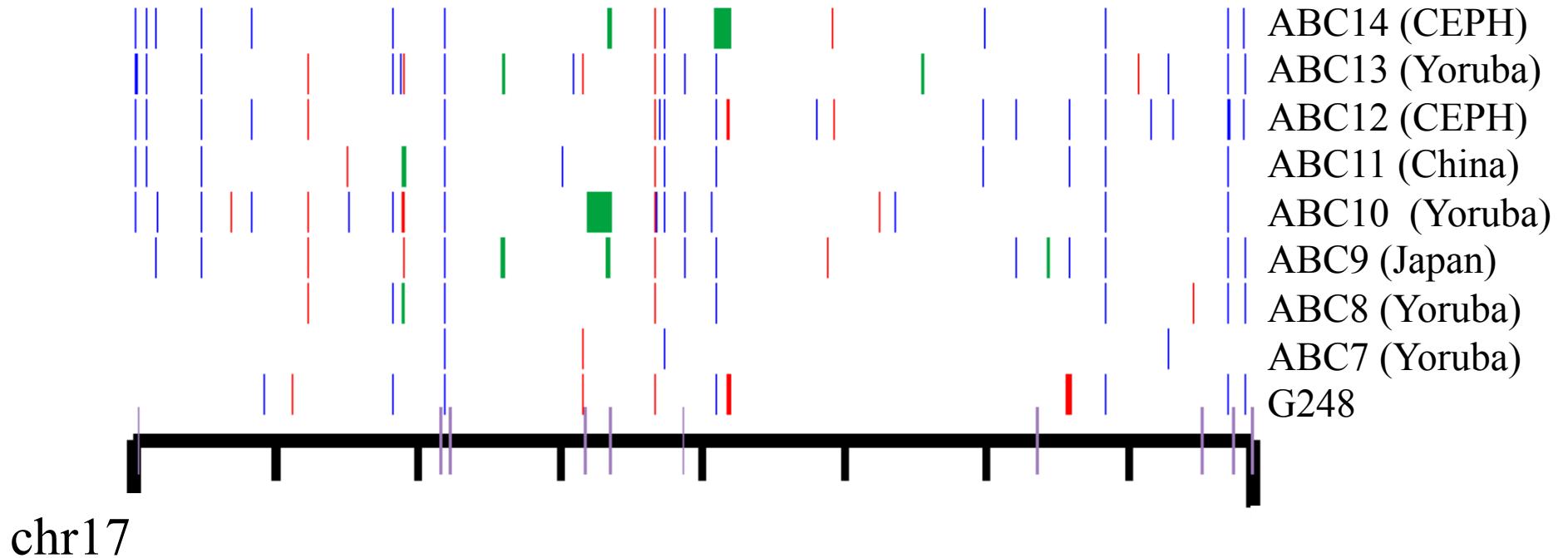
Tuzun *et al.* (2005)

What can we find?

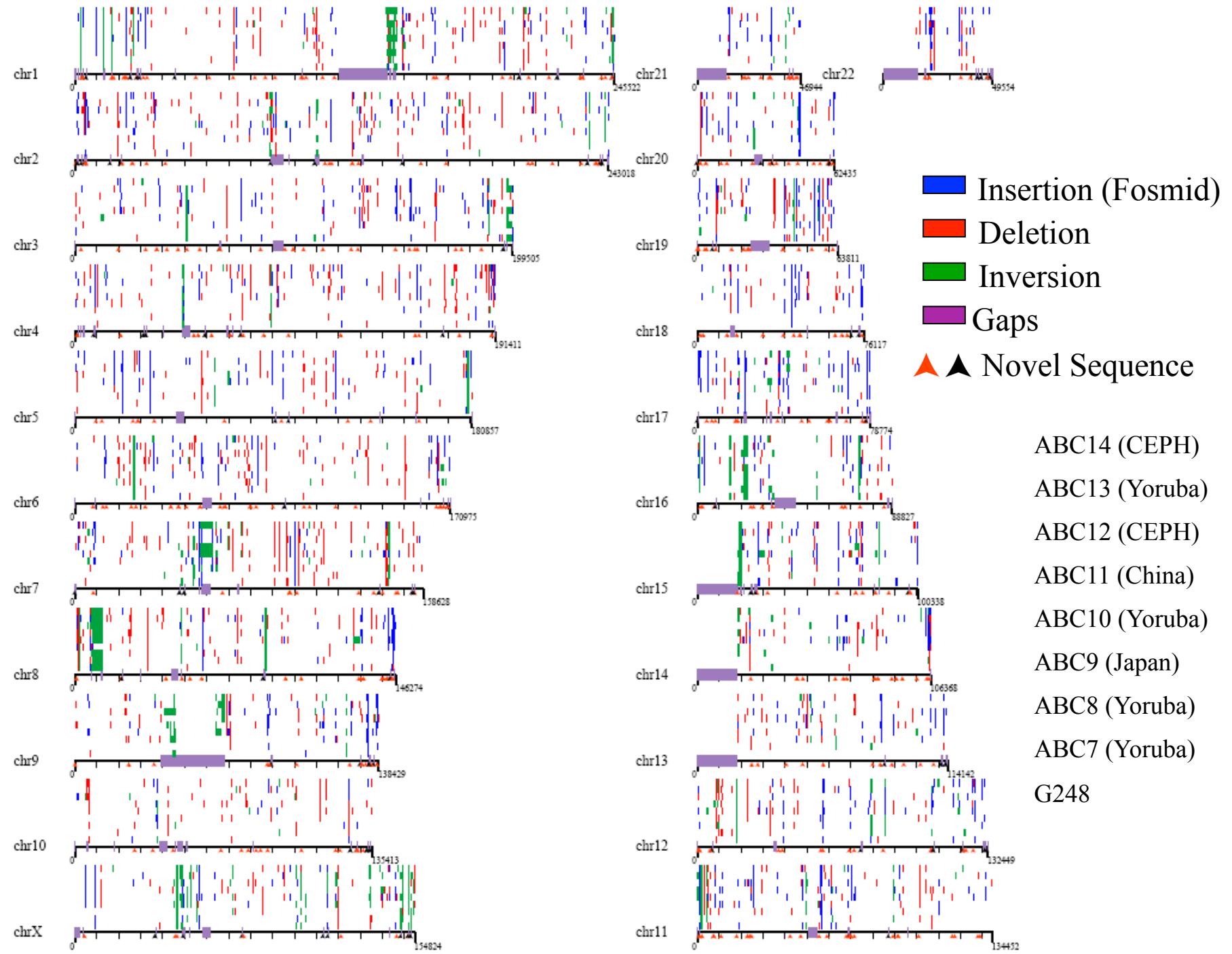
Structural variation detection:



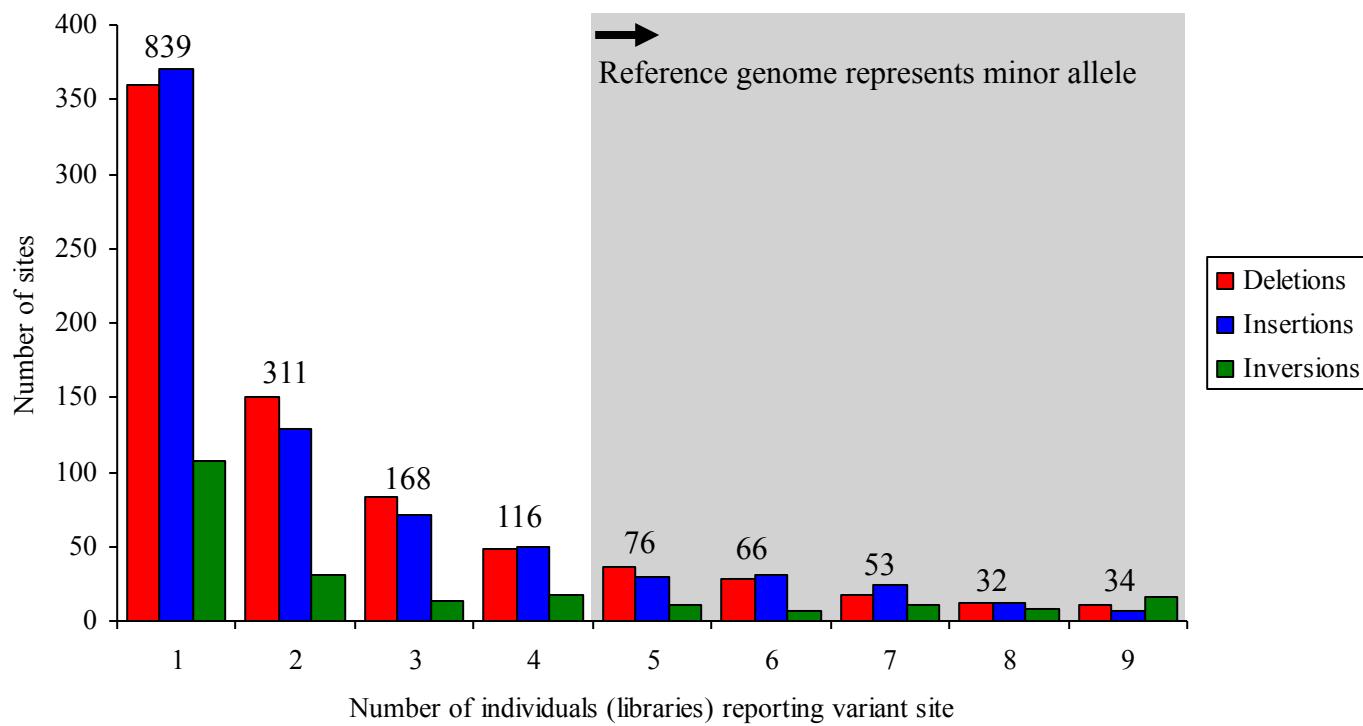
Map of Validated Variants



- Genome wide map of variants
- Ability to resolve structure of individual haplotypes

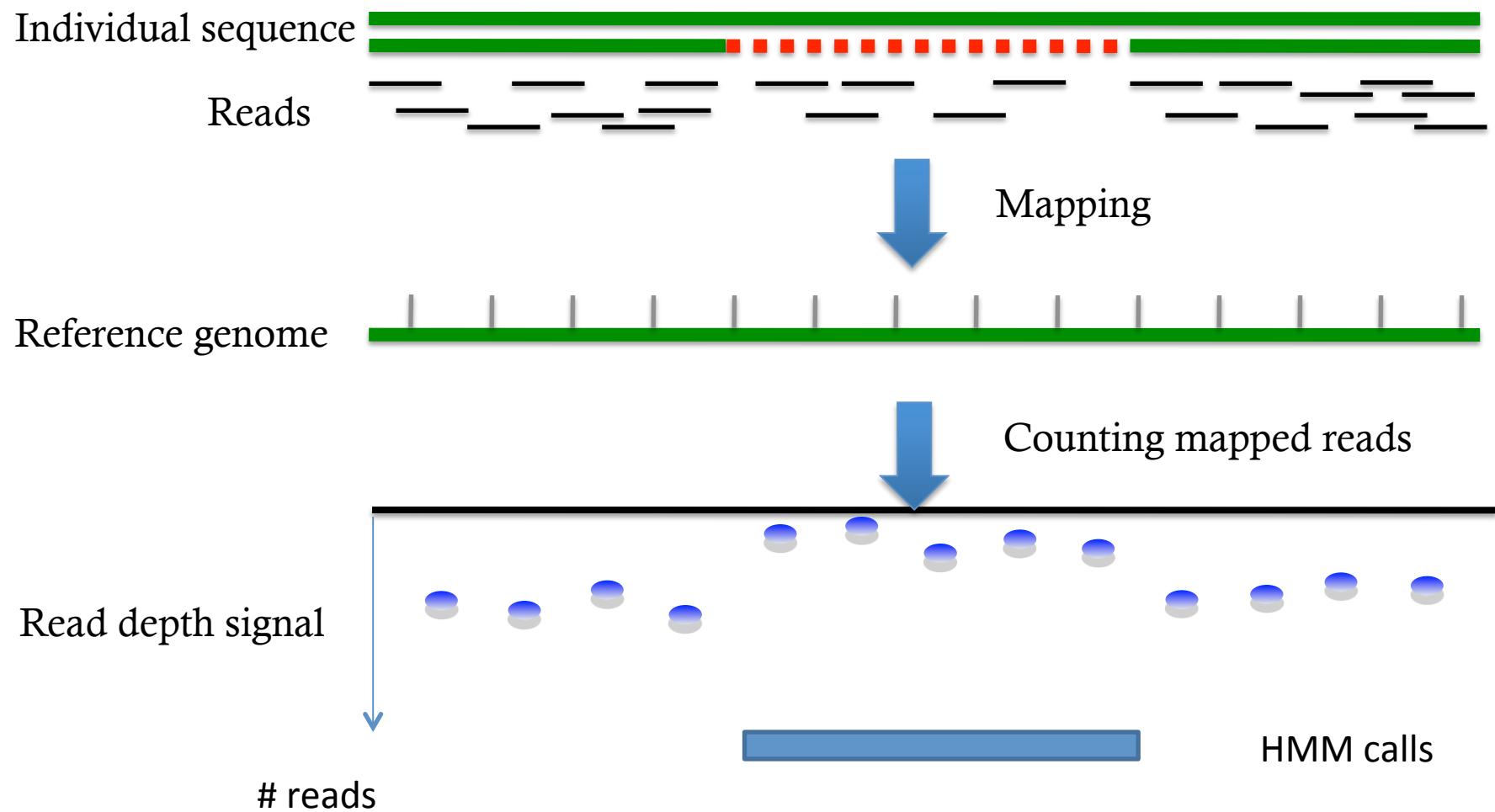


Frequency of Validated Sites

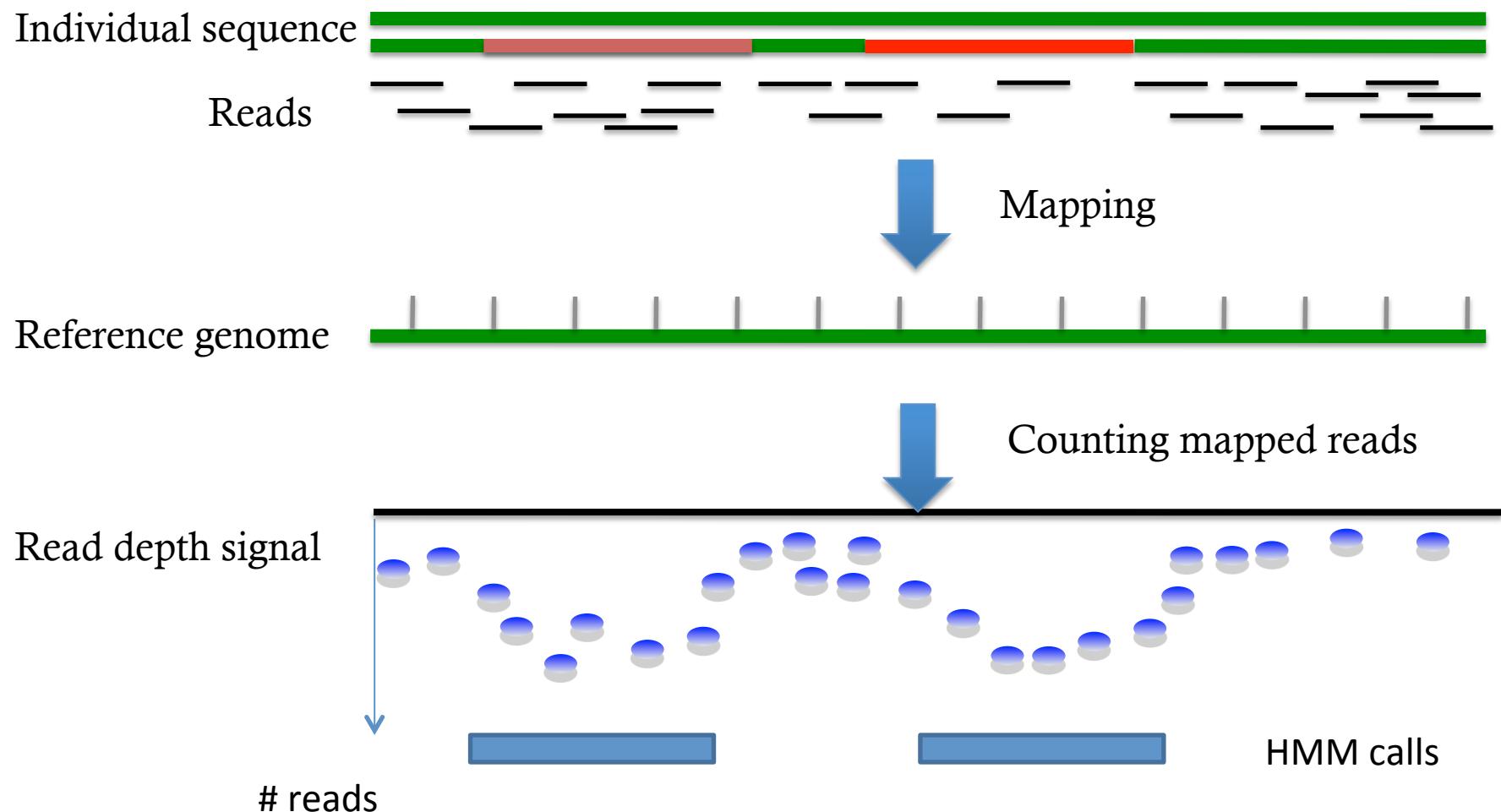


261 (15%) sites where reference genome represents a minor allele

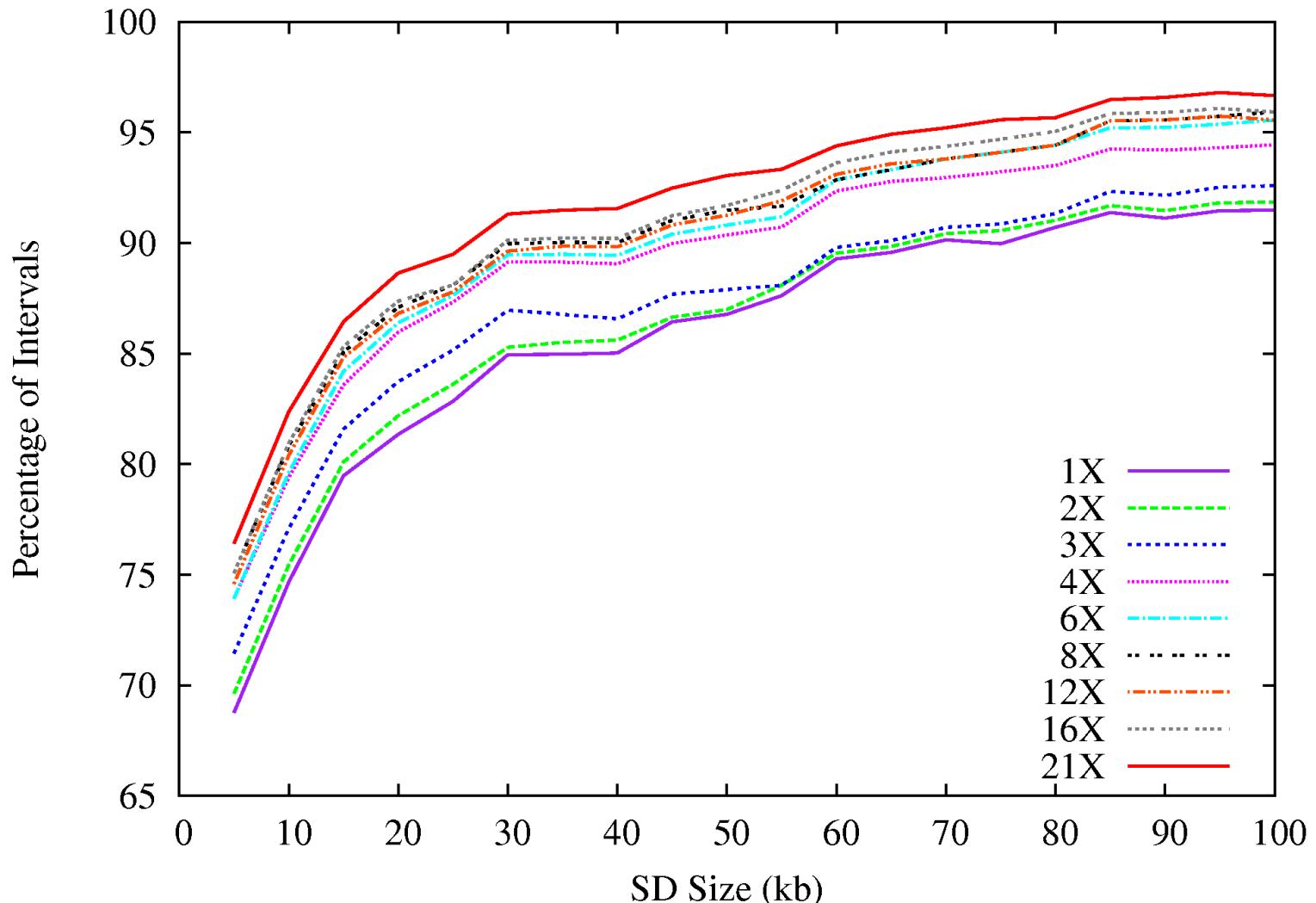
Method 3: Sequence Read Depth Analysis



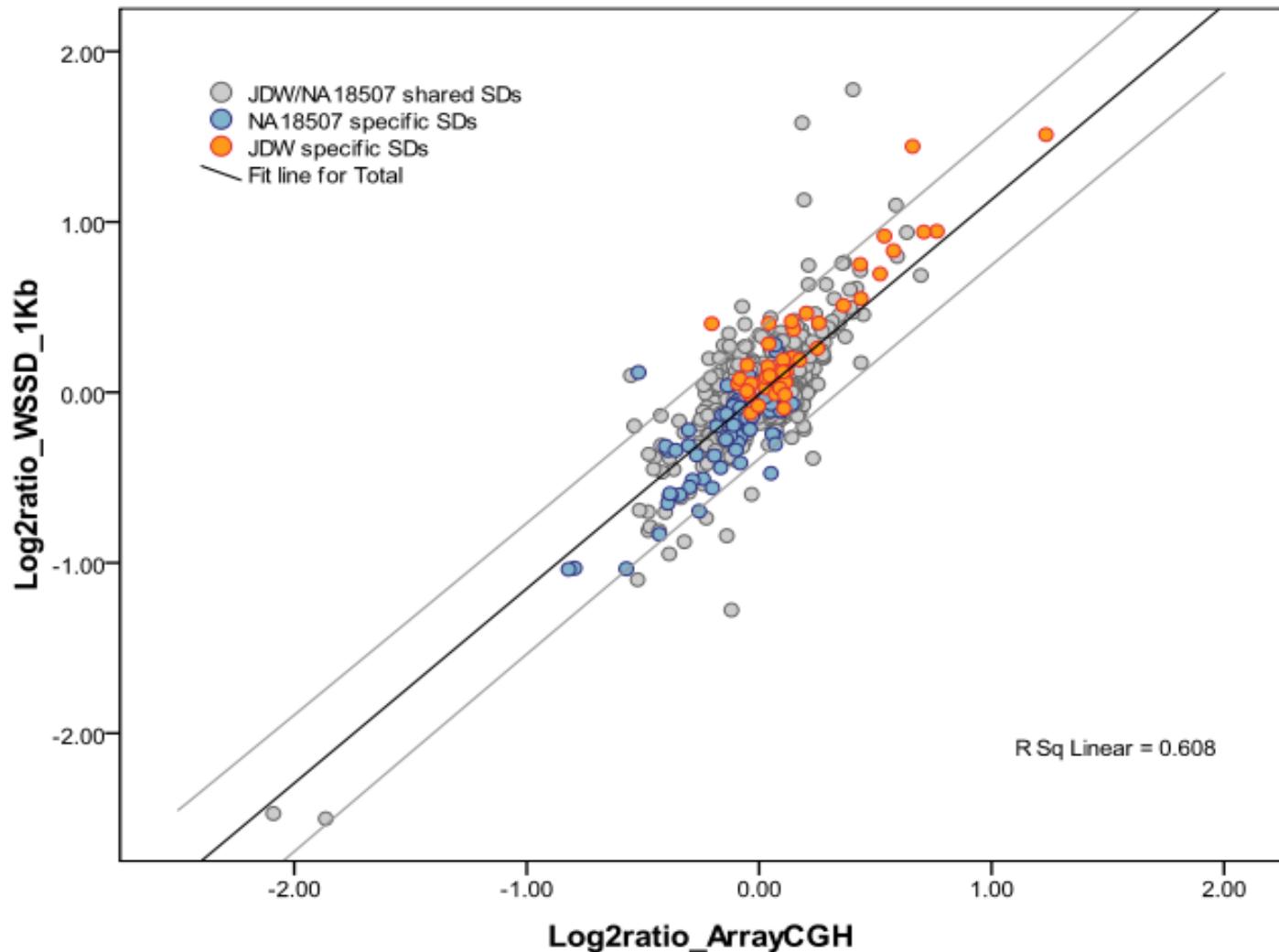
Method 3: Sequence Read Depth Analysis



Sequence coverage and detection power

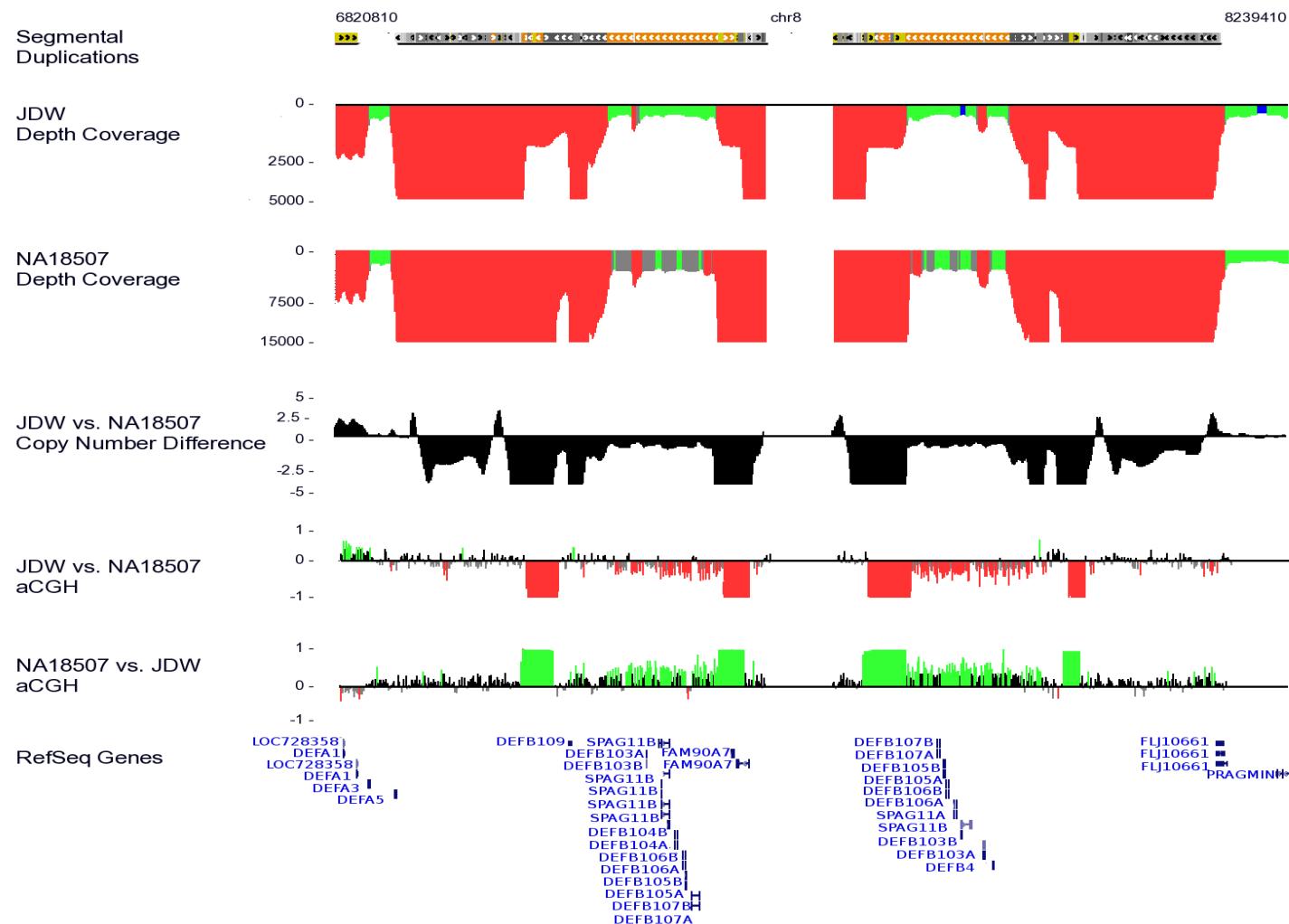


Validation of copy-number estimations



Alkan et al., *Nature Genetics*, 2009

Defensin gene cluster + *FAM90A7*

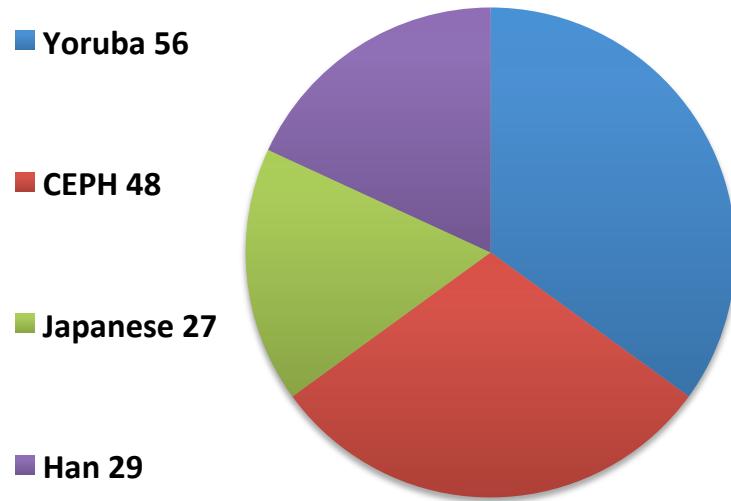


Associated with psoriasis and
Crohn's disease

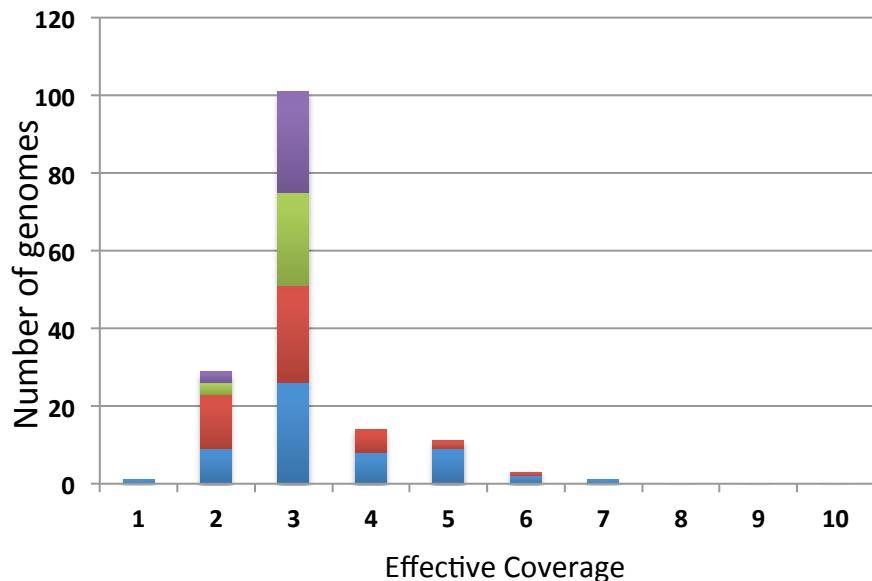
Alkan et al., Nature Genetics, 2009

Scaling up: 1000 Genomes and more

Individuals sequenced in Pilot 1



Histogram of Pilot 1 Illumina effective coverage



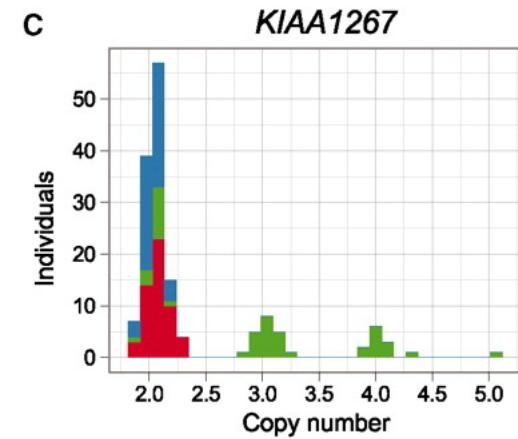
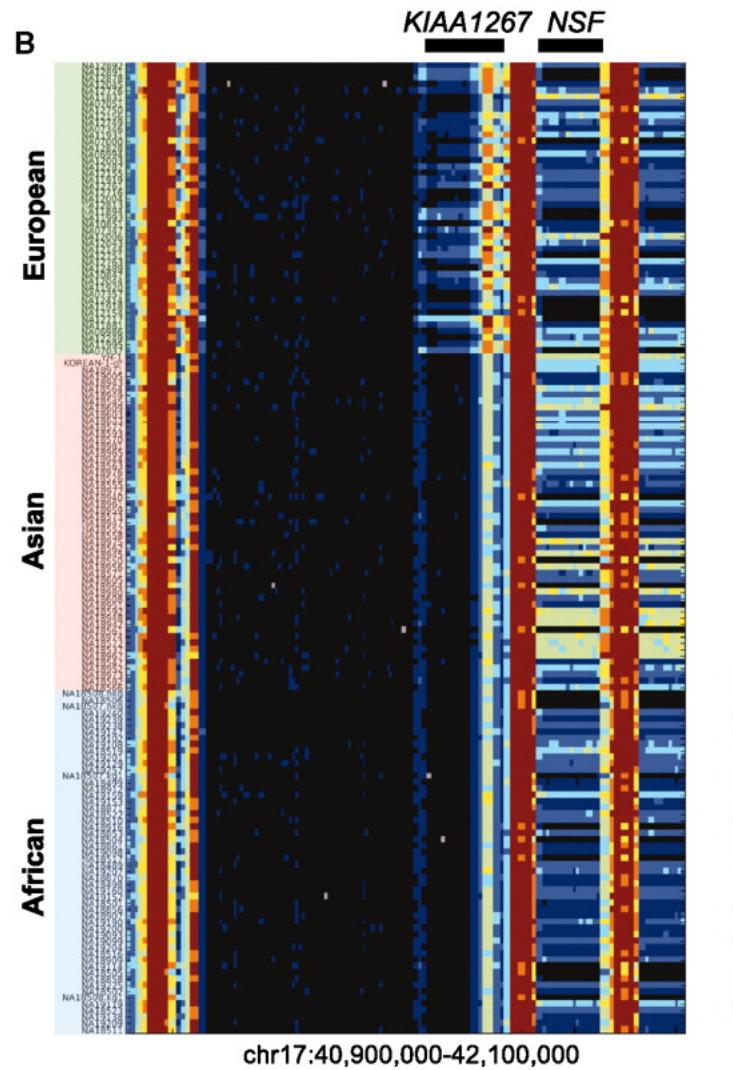
Individuals sequenced in Pilot 2

ID	Effective Coverage	Population
NA19240	24	YORUBA
NA19239	19	YORUBA
NA19238	13	YORUBA
NA12891	21	CEPH
NA12892	18	CEPH
NA12878	22	CEPH

Other Genomes

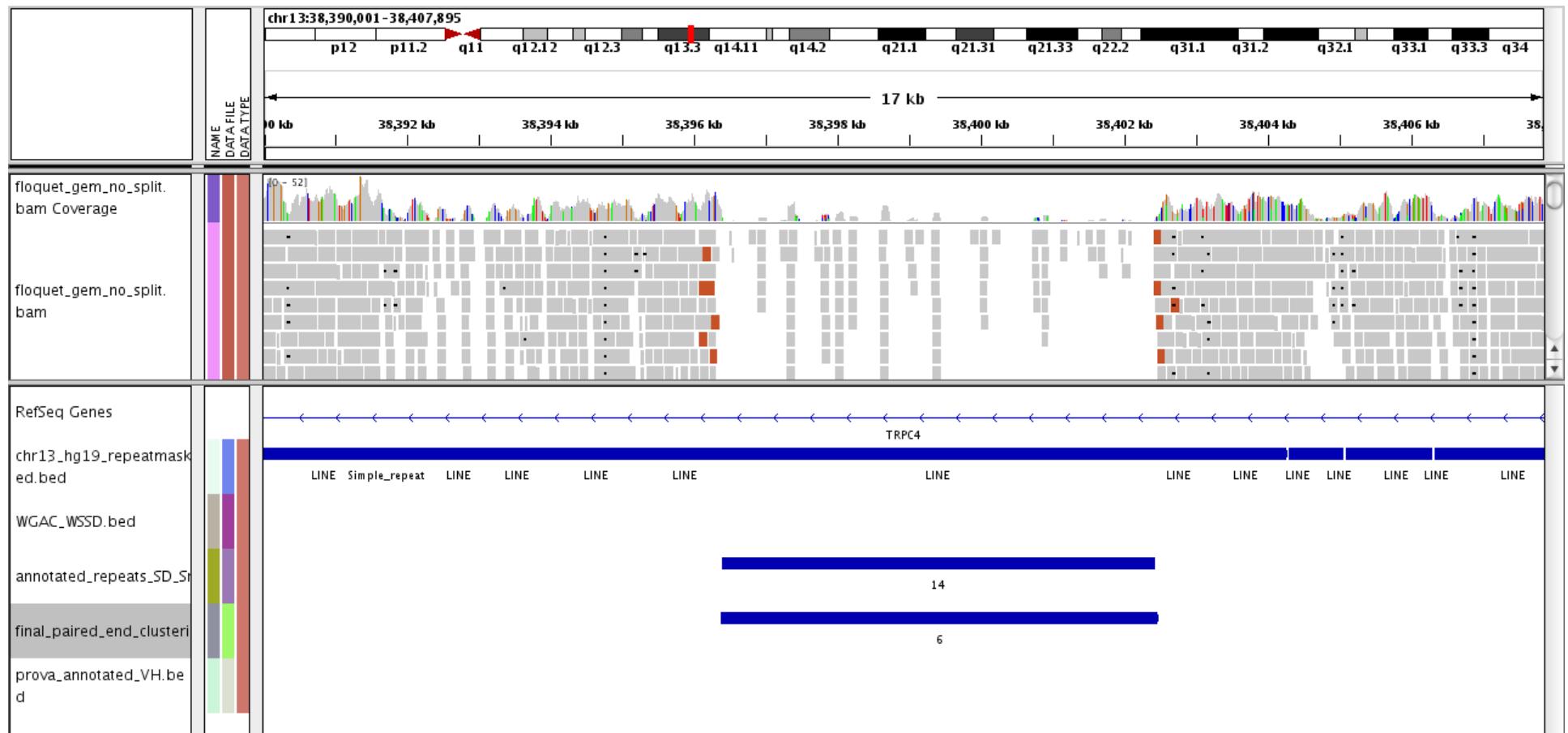
ID	Effective Coverage	Population
YH-1	22	HAN CHINESE°
NA18507	29	YORUBA ‡
NA18506	30	YORUBA *
NA18508	25	YORUBA *
KOREAN	12	KOREAN ♦

Copy number variation in human populations

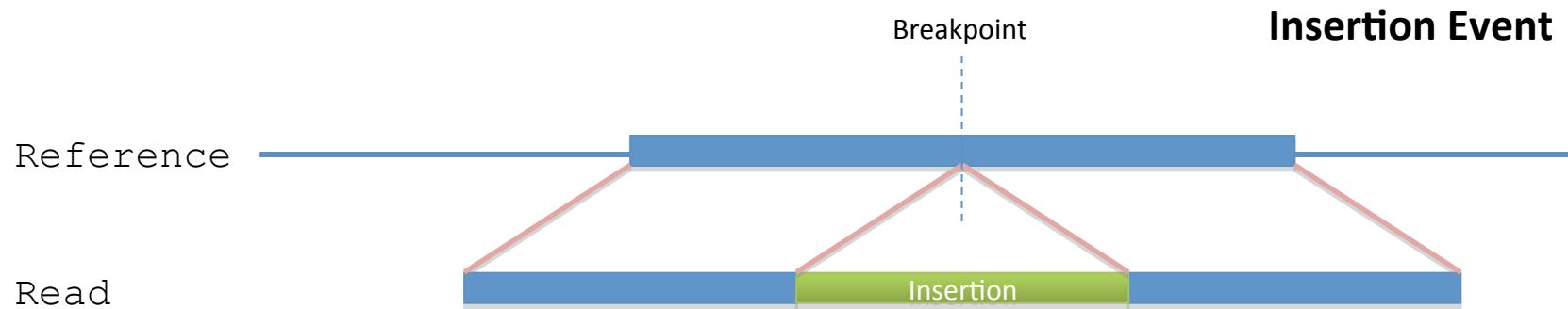
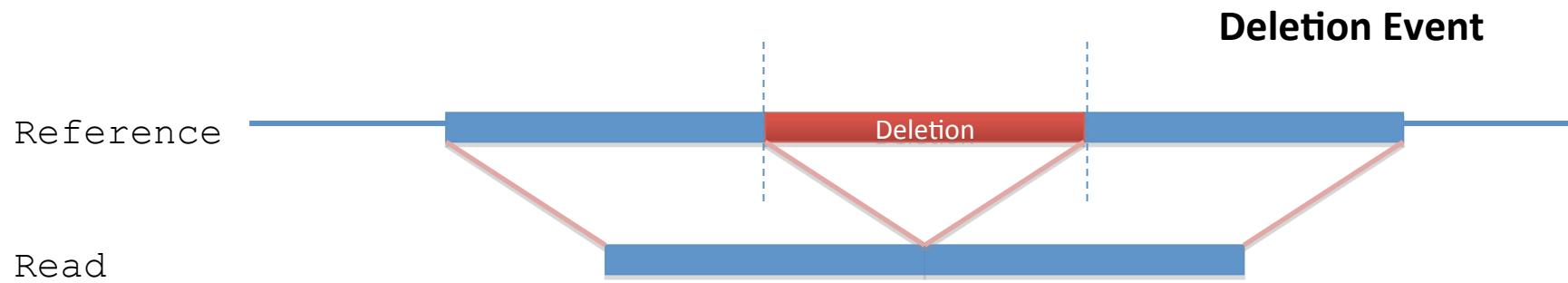


Sudmant et al. Science 2010

Deletions

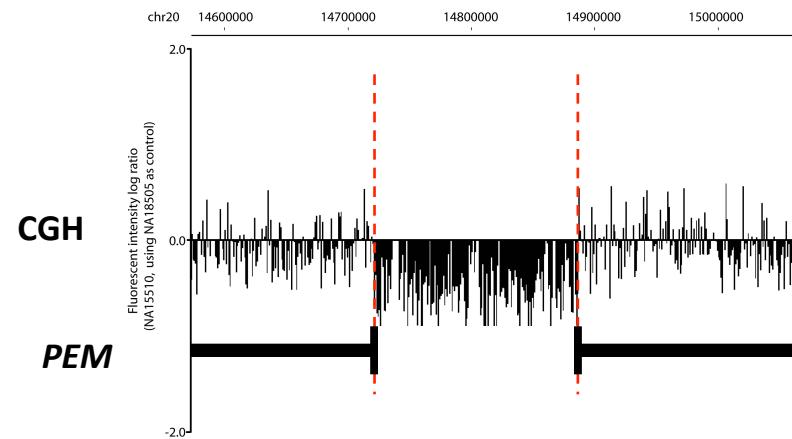


Split-read Analysis



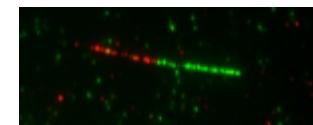
Experimental Validation

A) CGH

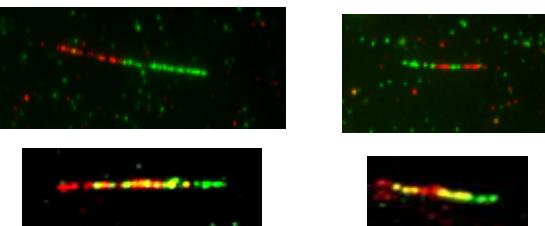


B) Fiber-FISH (For inversions)

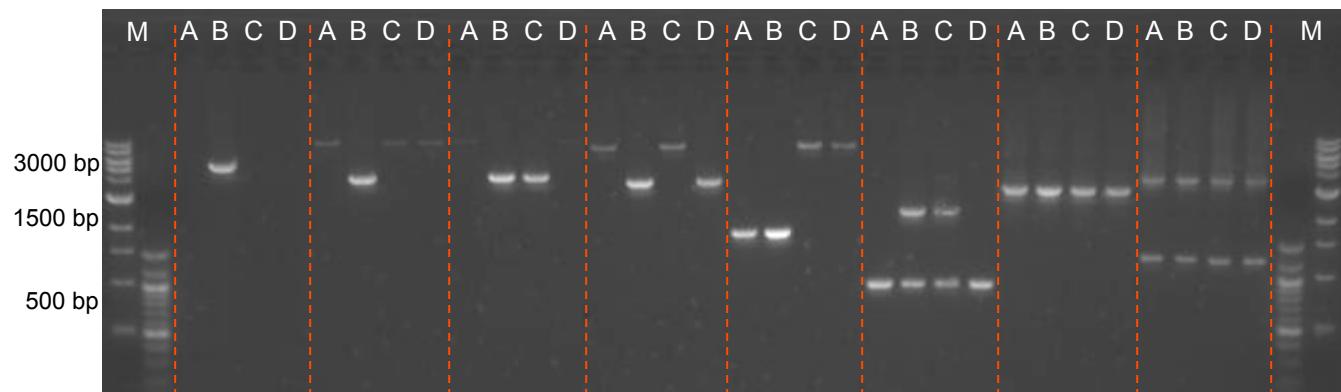
Without inversion



With inversion



C) PCR



Methods to Find SVs

Experimental approach

ArrayCGH (SNP based and genomic)

Based on ratios, Saturate quite fast, poor breakpoint resolution

Sequence based

Read pair analysis

Deletions, small novel insertions, inversions, transposons

Size and breakpoint resolution dependent to insert size

Read depth analysis

Deletions and duplications

Relatively poor breakpoint resolution

Split read analysis

Small novel insertions/deletions, and mobile element insertions

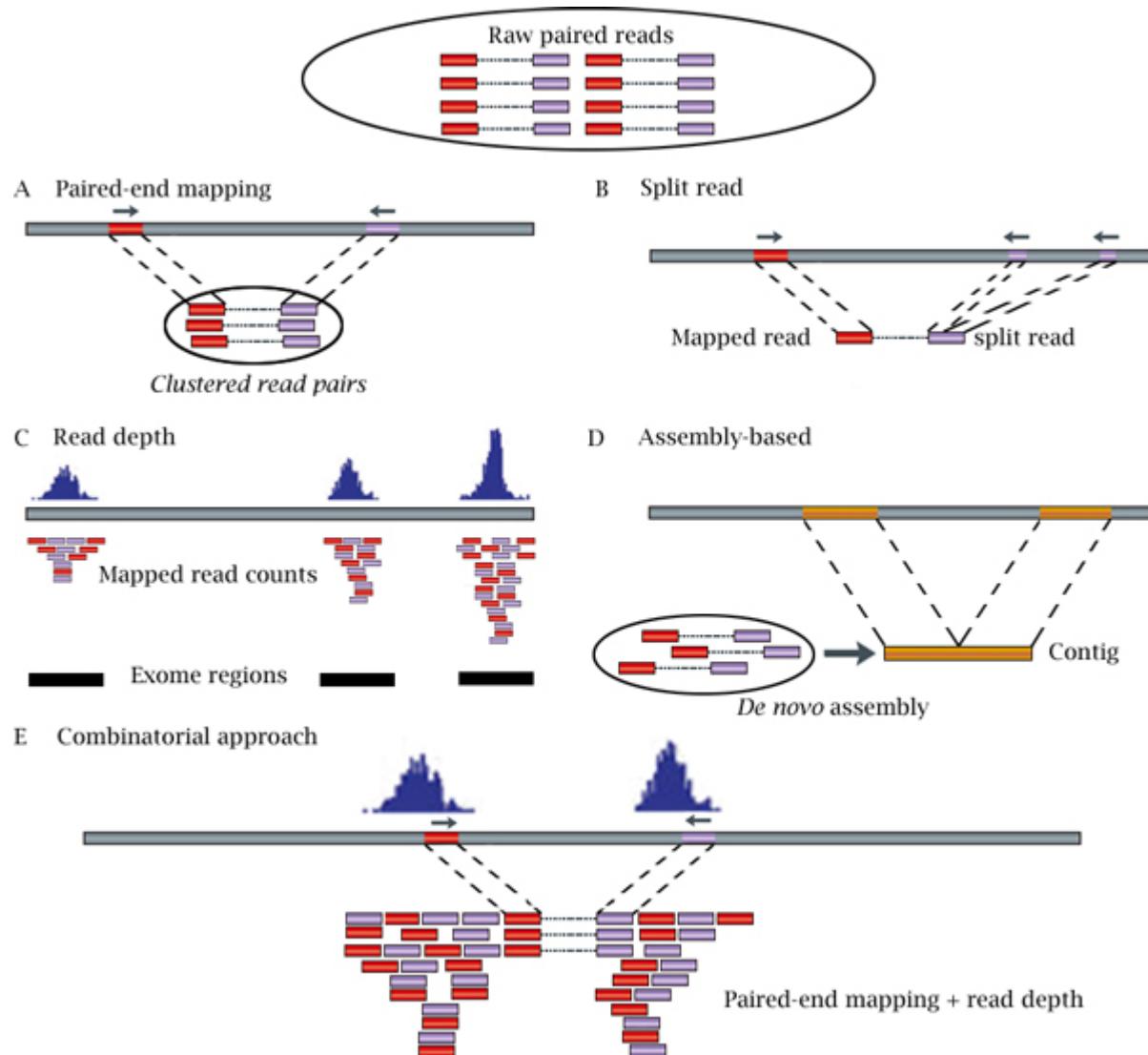
1bp breakpoint resolution

Local and *de novo* assembly

SV in unique segments

1bp breakpoint resolution

Review software



Software I

Method	Reference	Language	Control required?	Input format	GC correction	single-end/pair-end	Methodology characteristics
CNV-seq	[15]	R, perl	Yes	hits	No	single-end	statistical testing
FREEC	[21]	C	Optional	SAM,BAM,bed,etc.	Optional	both	LASSO regression
readDepth	[22]	R	No	bed	Yes	both	CBS, LOESS regression
CNVnator	[23]	C	No	BAM	Yes	both	mean shift algorithm
SegSeq	[14]	Matlab	Yes	bed	No	single-end	statistical testing,CBS
EWT (RDXplorer)	[11]	R, python	No	BAM	Yes	single-end	statistical testing
cnD	[16]	D	No	SAM,BAM	No	both	HMM, Viterbi algorithm
CNVer	[17]	C	No	BAM	Yes	pair-end	maximum-likelihood, graphic flow
CopySeq	[18]	Java	No	BAM	Yes	pair-end	MAP estimator
rSW-seq	[19]	NA	Yes	NA	Yes	single-end	Smith-Waterman algorithm
CNAseg	[20]	R	Yes	BAM	No	pair-end	wavelet transform and HMM
CNAnorm	[24]	R	Yes	SAM,BAM	Yes	both	linear regression or CBS
cn.MOPS	[26]	R, C++	multiple samples	BAM or data matrix	No	both	mixture of Poissons, MAP, EM, CBS
JointSLM	[27]	R, Fortran	multiple samples	data matrix	Yes	both	HMM, ML estimator, Viterbi algorithm

doi:10.1371/journal.pone.0059128.t001

break point position estimation: readDepth = EWT>CNVnator>FREEC>CNV-seq>SegSeq;
copy number estimation: CNVnator>CNV-seq>readDepth>FREEC>EWT>SegSeq;

Zhao et al. BMC Bioinformatics 2013

Duan et al. Plos One 2013