

# Discovery and annotation of structural variants

**Victor Guryev**

**European Research Institute for the Biology of Ageing (ERIBA)**

**UMC Groningen, Rijksuniversiteit Groningen**

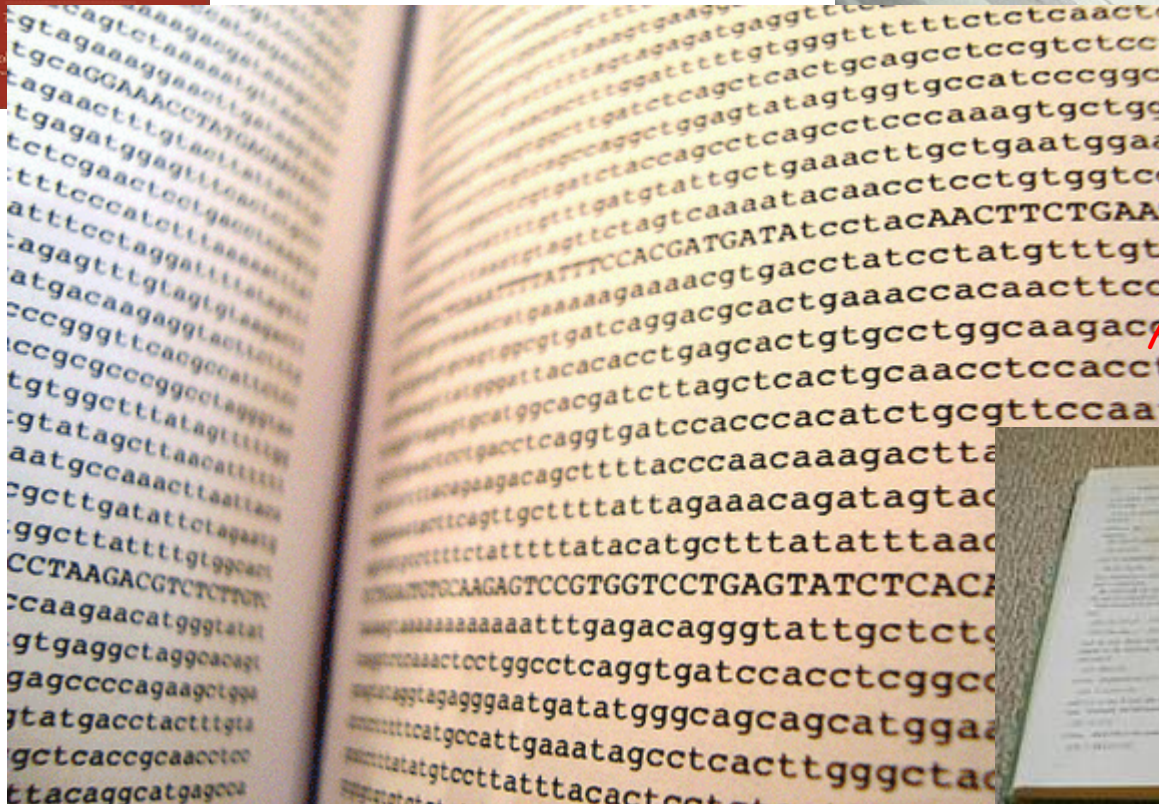
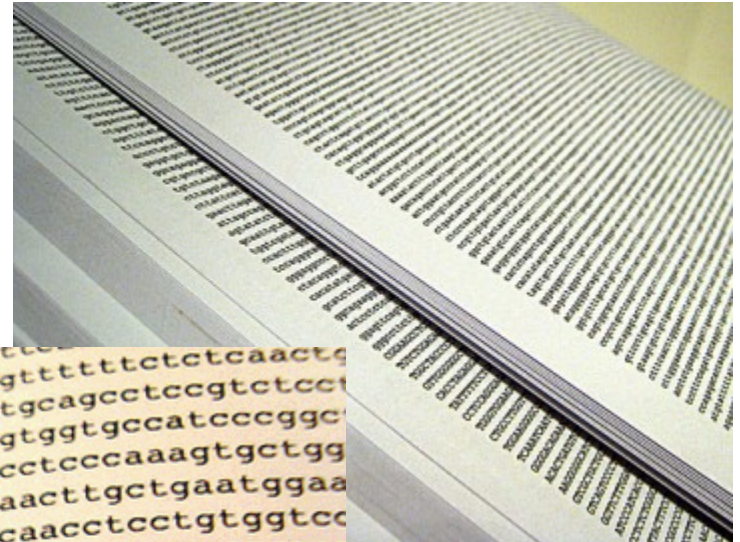
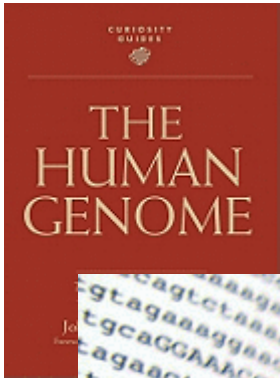
**NBIC**

**High-Throughput Next Generation Biology Course**

**Groningen**

**May 7, 2013**

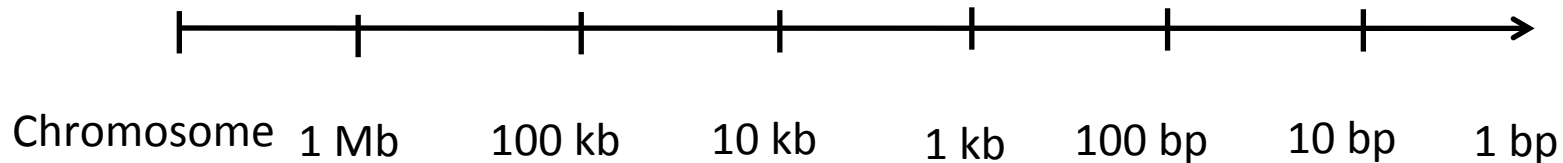
# The human genome 'book'



A

# Scale of genetic differences

Scale



Variants

Short indels

SNPs

**Structural variants**

Tools  
Solutions



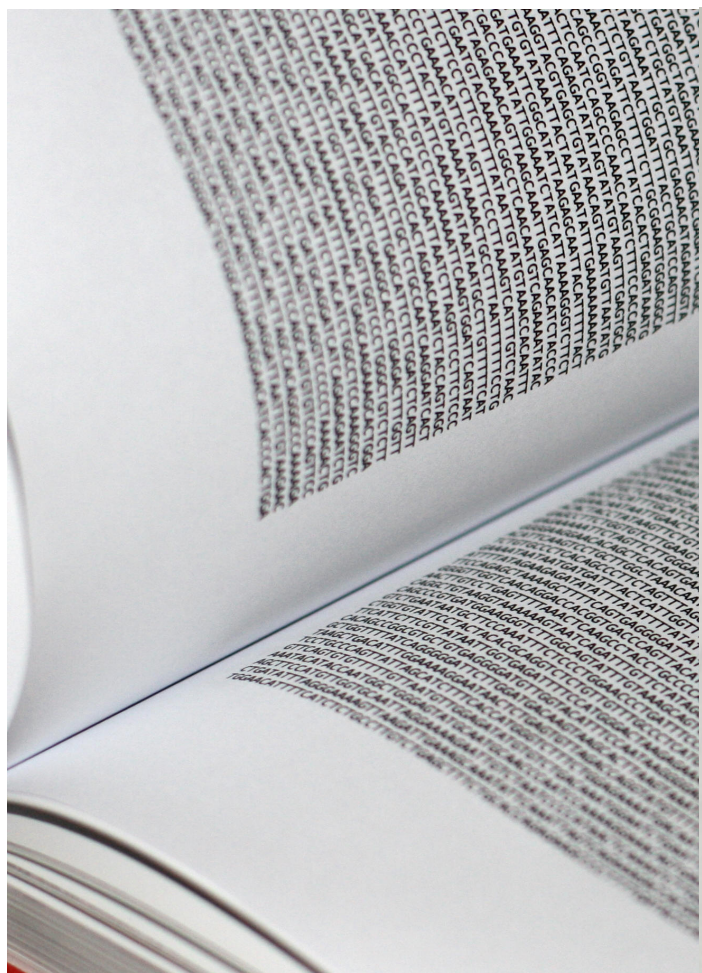
# \$1000 genome and beyond



*We can be confident in predicting that the \$1,000 human genome will be achieved in 2013*

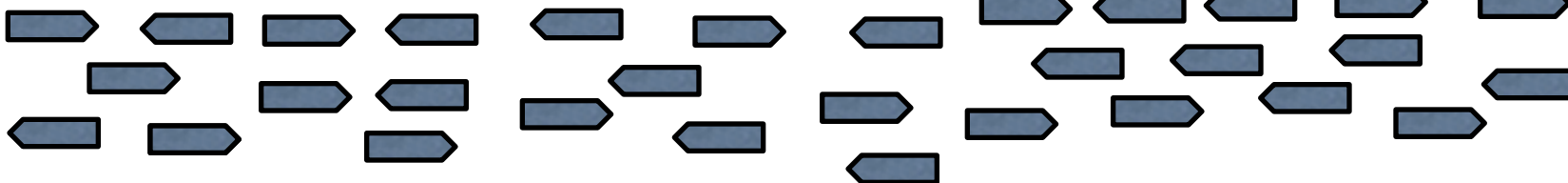
Life sciences are ready for a revolution, but it will require collaboration on many fronts, says **Yang Huanming**, president of BGI (the Beijing Genomics Institute)

# How do we get our NGS genomes?

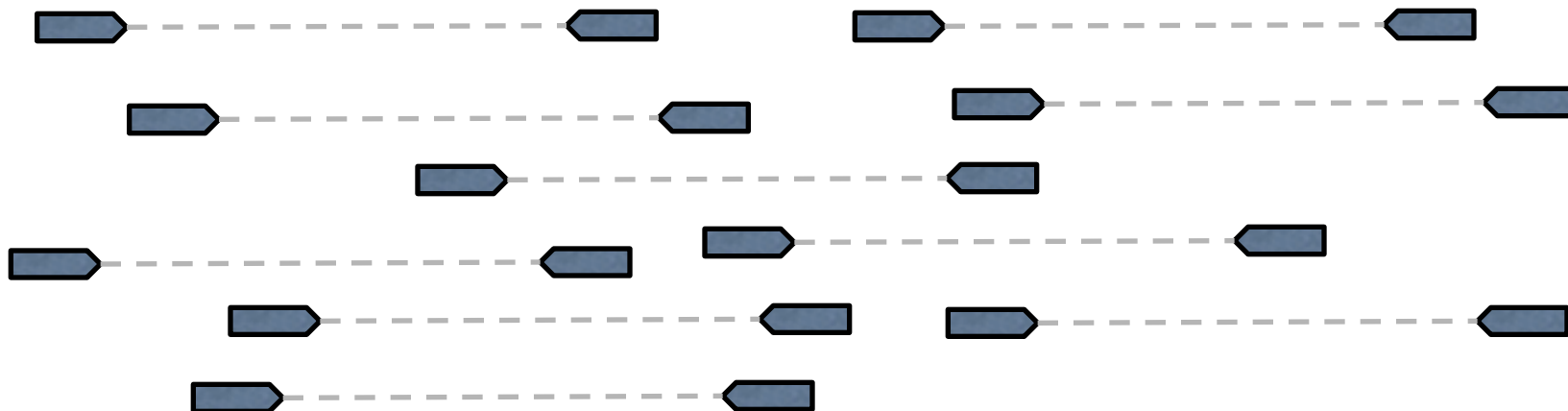


# Fragment and paired-sequencing

chromosome



chromosome



90 bp

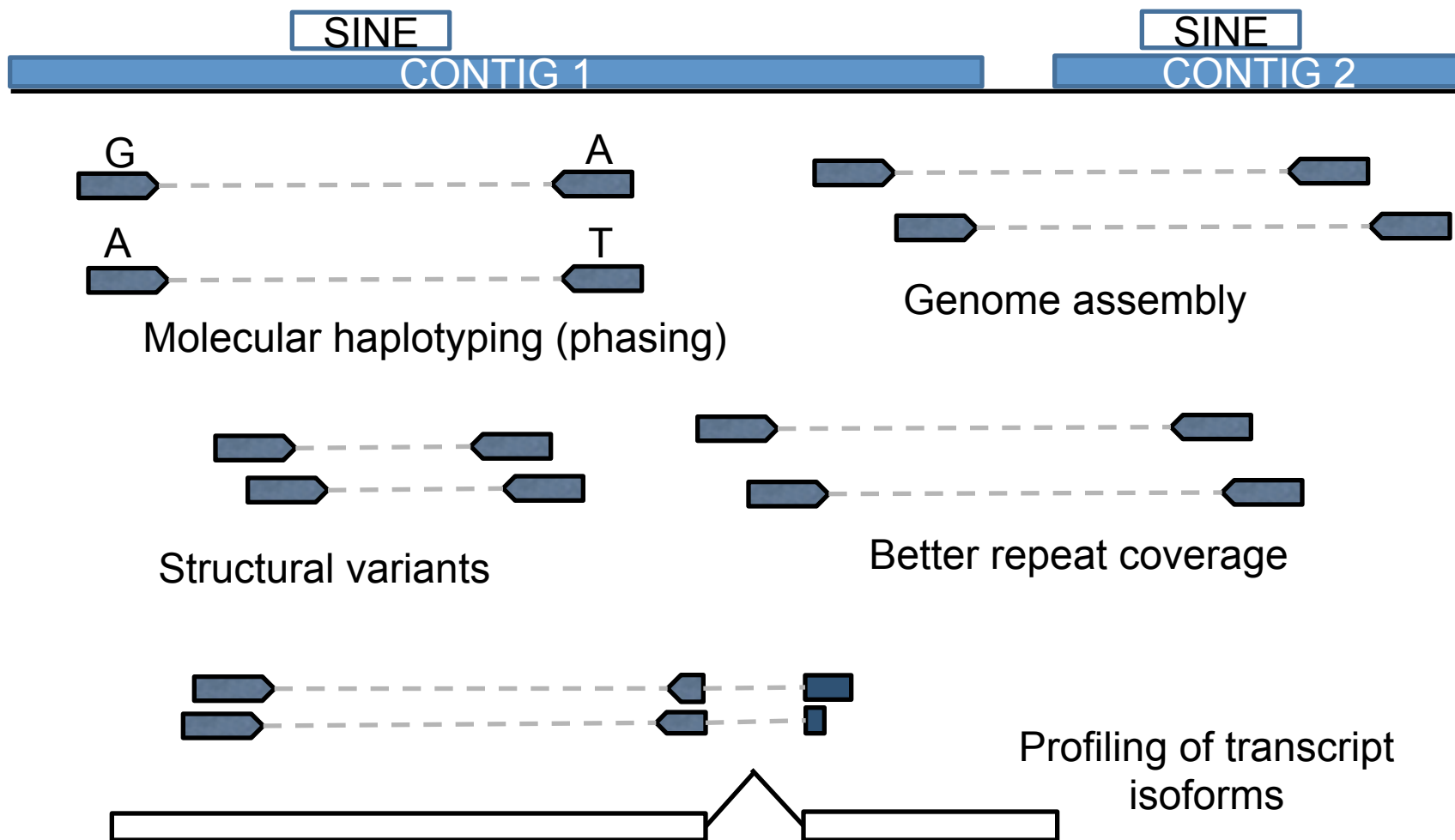


~ 500 bp

# Advantages of paired-sequencing

1) Twice as many bases per slide !

2) Structural information !!!



# SV types and their detection

## Structural Genome Variations (SVs)

ABCD

### Copy-number variants

### Copy-balanced variants

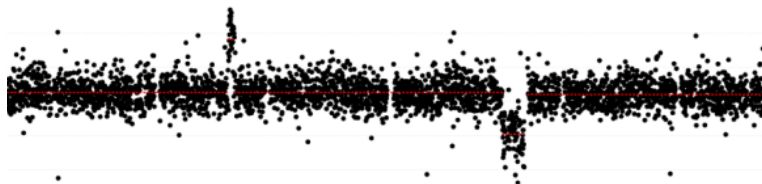
Deletion  
ABD

Duplication  
ABCCCD

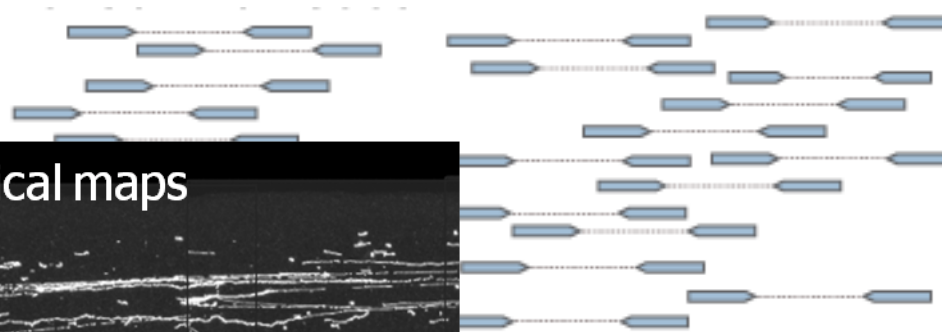
Inversion  
ADCDB

Translocation  
AB CD

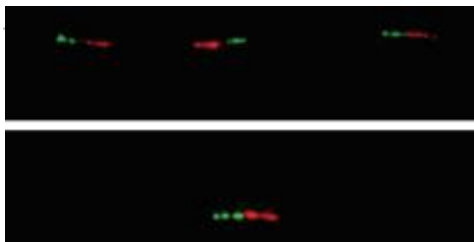
aCGH



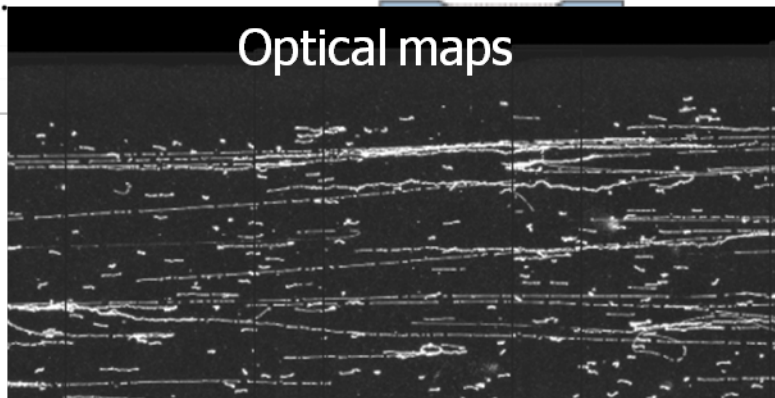
Di-tag fosmid and NGS sequencing



Fibre-FISH



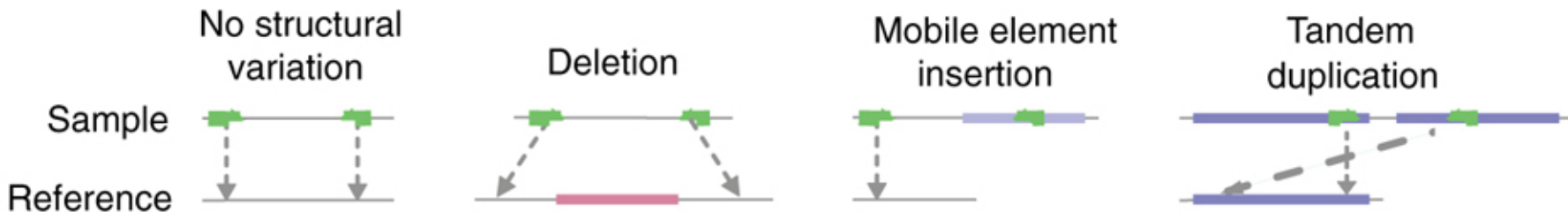
Optical maps





# Approaches for SV detection using NGS data

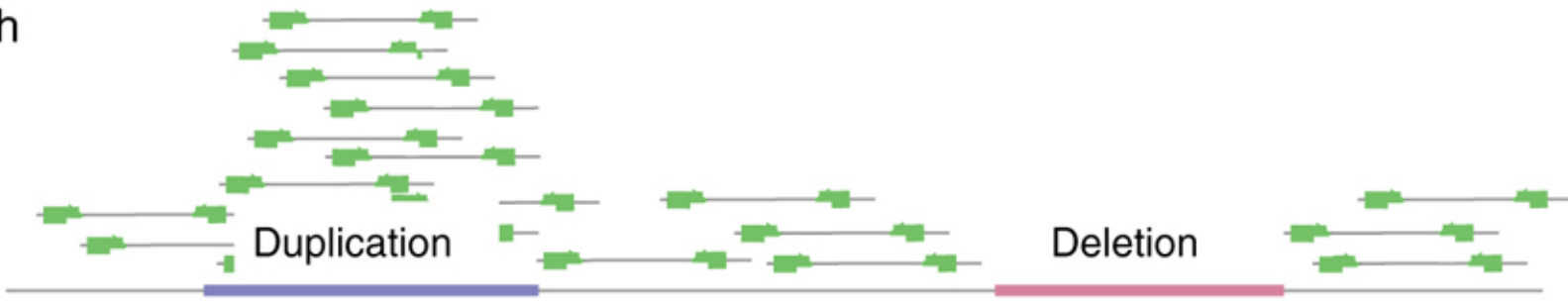
## Read pairs



## Read depth

Sample reads

Reference



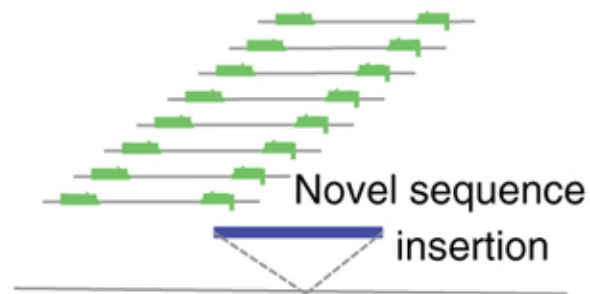
## Split reads

Reference



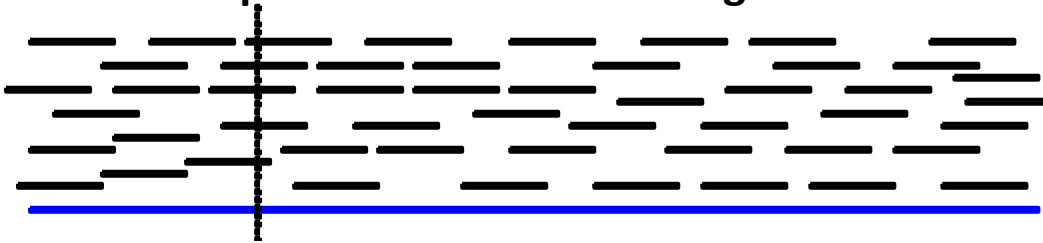
## Assembly

Reference

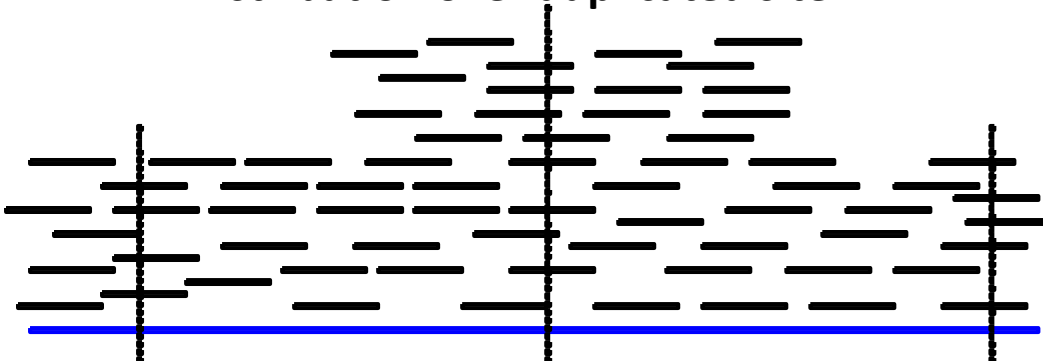


# RD: Read density analysis

Expected distribution of tags



Distribution over duplicated site



Scope:

Copy-number changes

Tools:

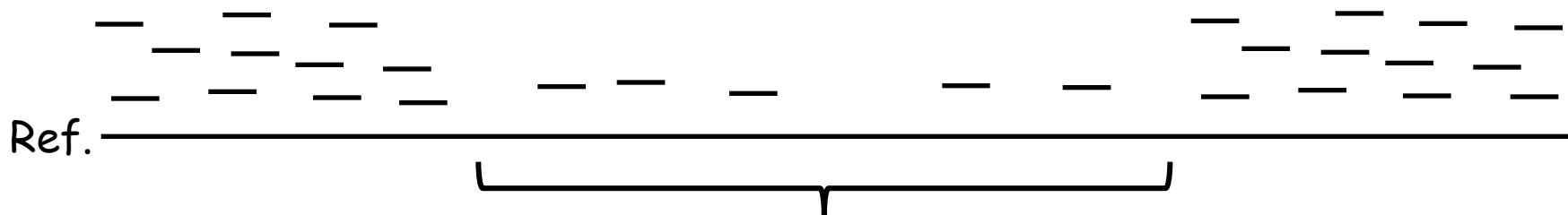
CNV-Seq (Xie & Tammi  
2009)

SegSeq (Chiang et al, 2009)

DWAC-Seq (our tool)

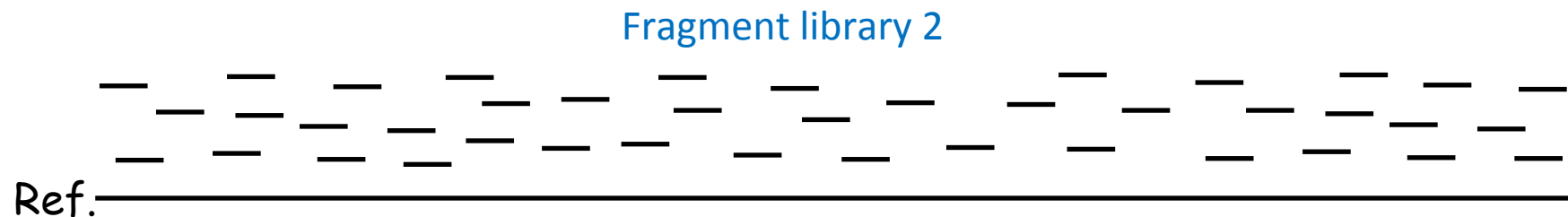
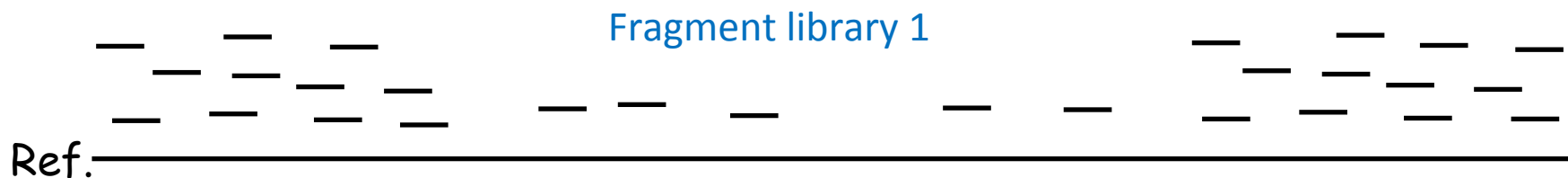
# Real data: non-uniformity of genome coverage

## Fragment library, sample vs genome reference



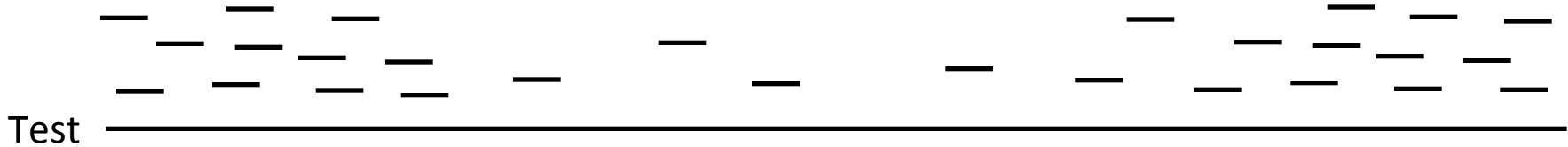
Heterozygous deletion ?  
'Sequenceability' issue ?  
'Mappability' issue ?

## Sample vs sample comparison

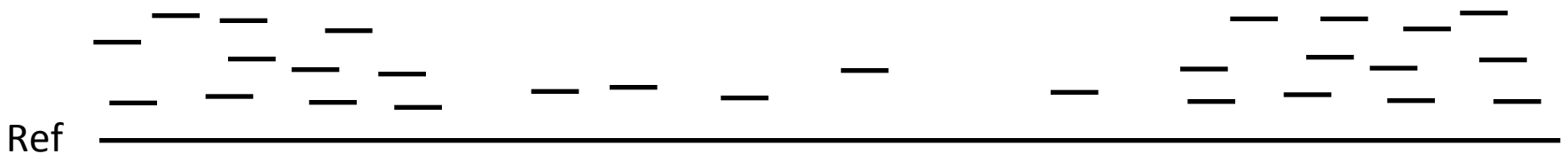


# Read depth analysis with dynamic windows

Fragment library 1



Fragment library 2



Static windows, e.g. 1kb

Window 4



Dynamic windows, each having fixed number reads in the reference set

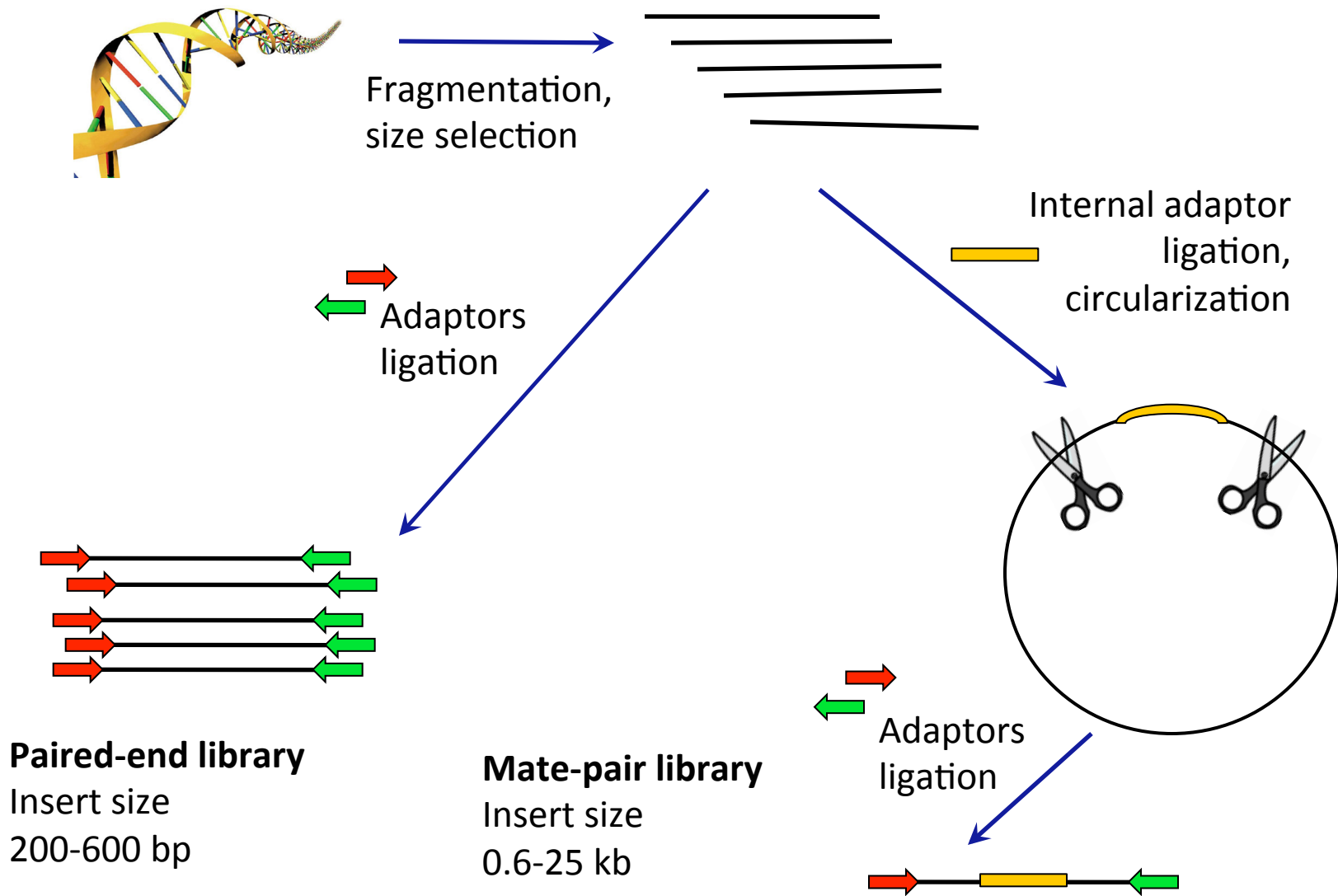
**Step1:** Segmenting the genome, CNV calling

**Step2:** Fine-mapping to determine their exact breakpoints

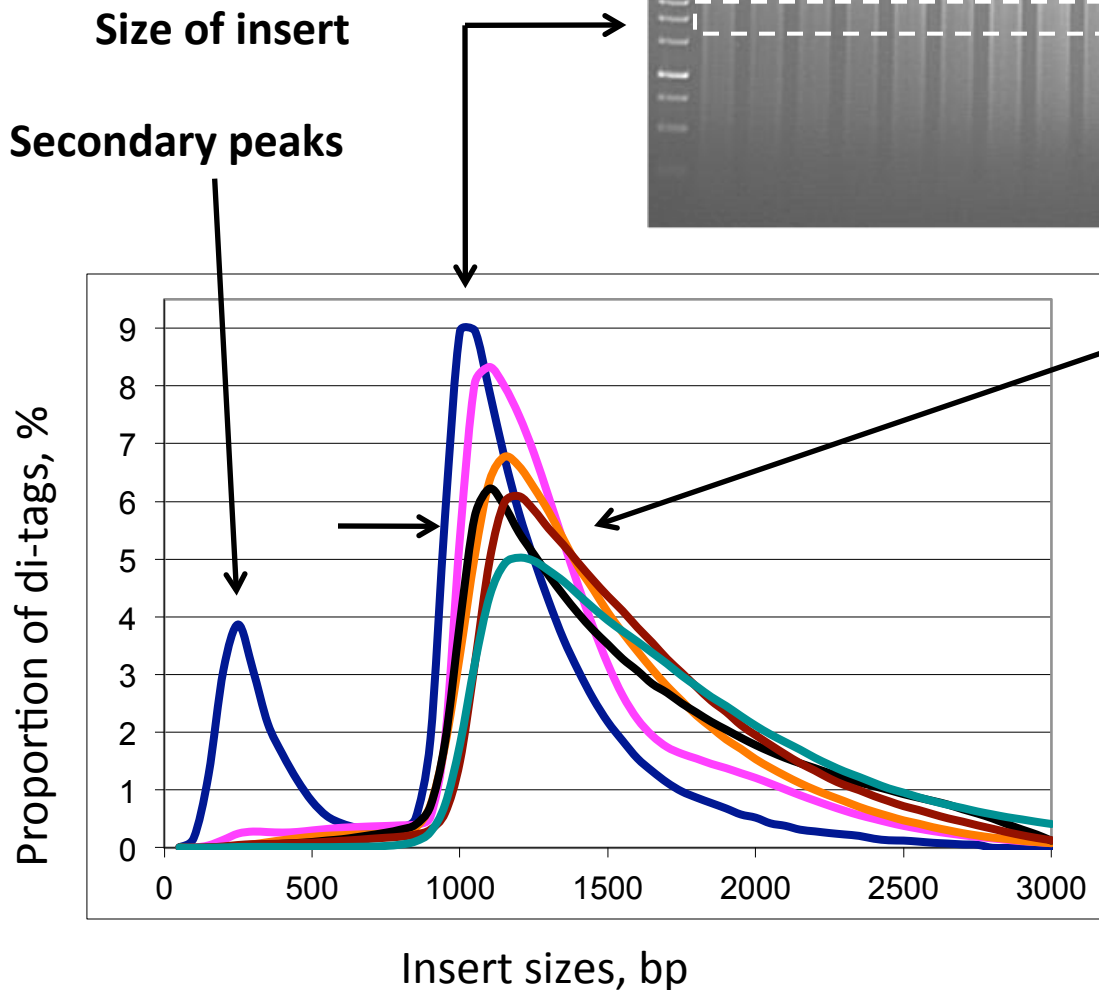
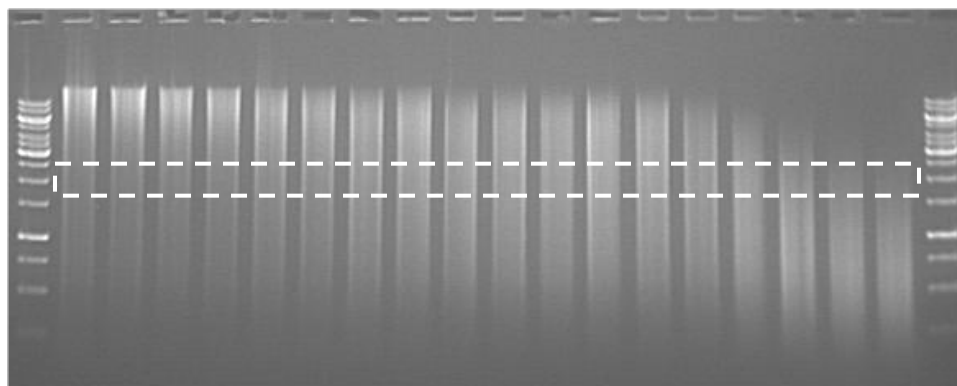
Dynamic-Window Approach for CNV calling using nextgen Sequencing

<http://tools.genomes.nl>

# Paired-end and mate-pair libraries



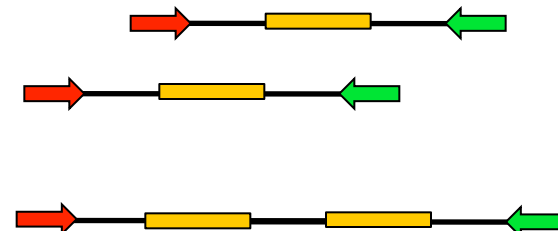
# QC of paired libraries



**Sharpness of size distribution**

**Clonality**  
• Low DNA input/  
overamplification

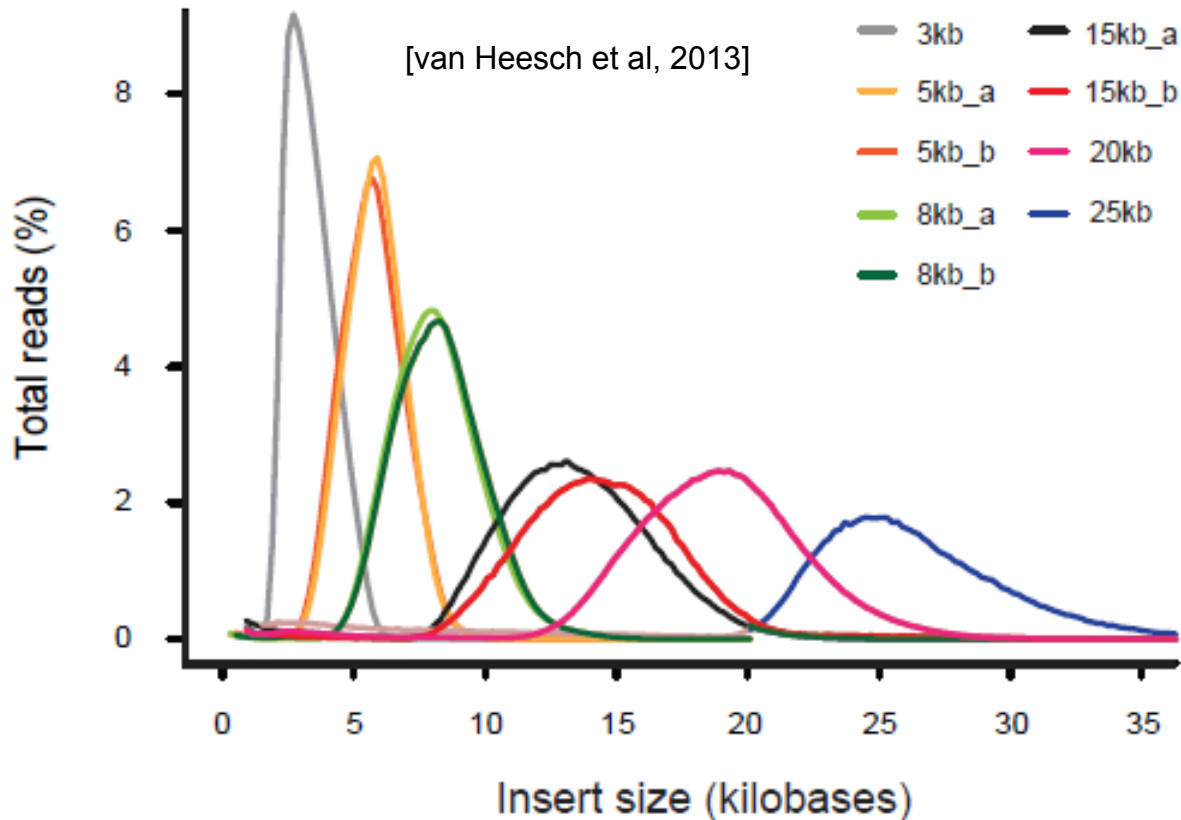
**Chimerism**



# Di-tag libraries size: S, M, L, XL, ...

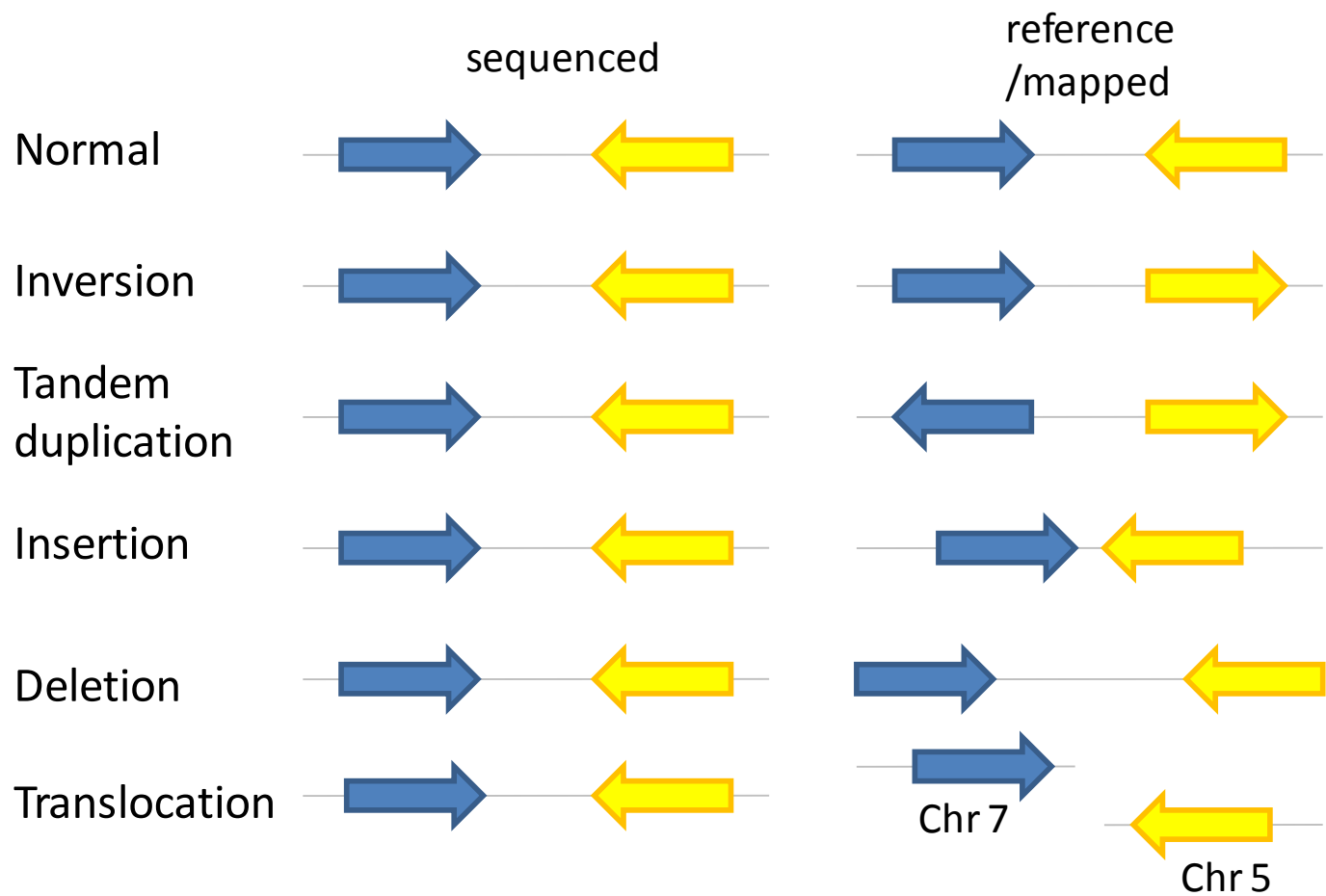
Different insert sizes to cover small and large variants with optimal precision

Di-tag libraries generated: PE (200bp), MP (2,3,5,8,15,20,25 kb)



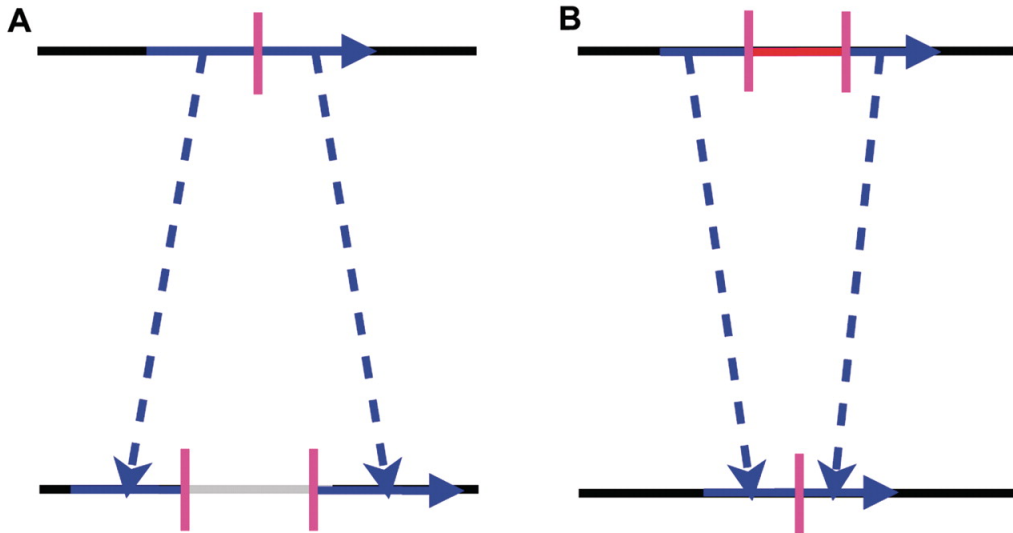
+ Fosmid libraries:  
40kb tags  
Protocol for Illumina  
[Williams et al 2012]

# Signatures of structural variation





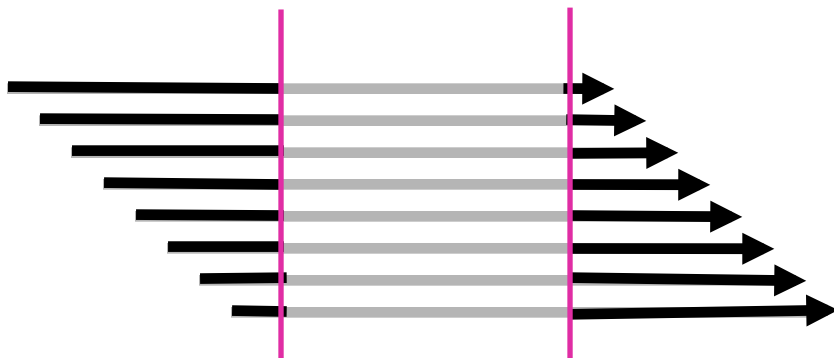
# Split-reads mapping



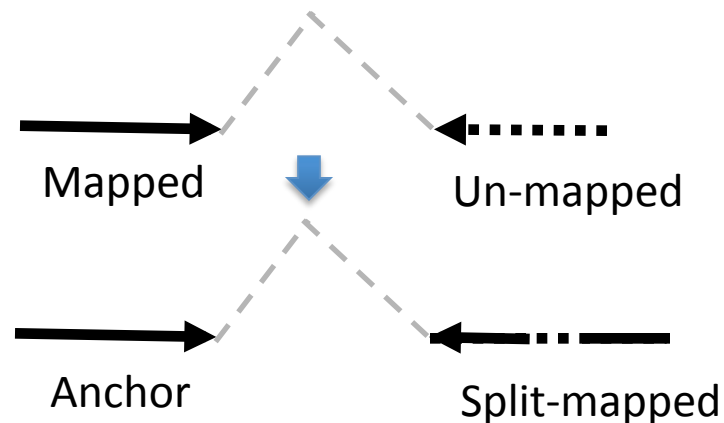
**Scope:** prediction of copy-number and copy-neutral SV at nucleotide resolution

**Tools:**  
Pindel (Ye et al, 2009)  
SRiC (Zhang et al, 2011)

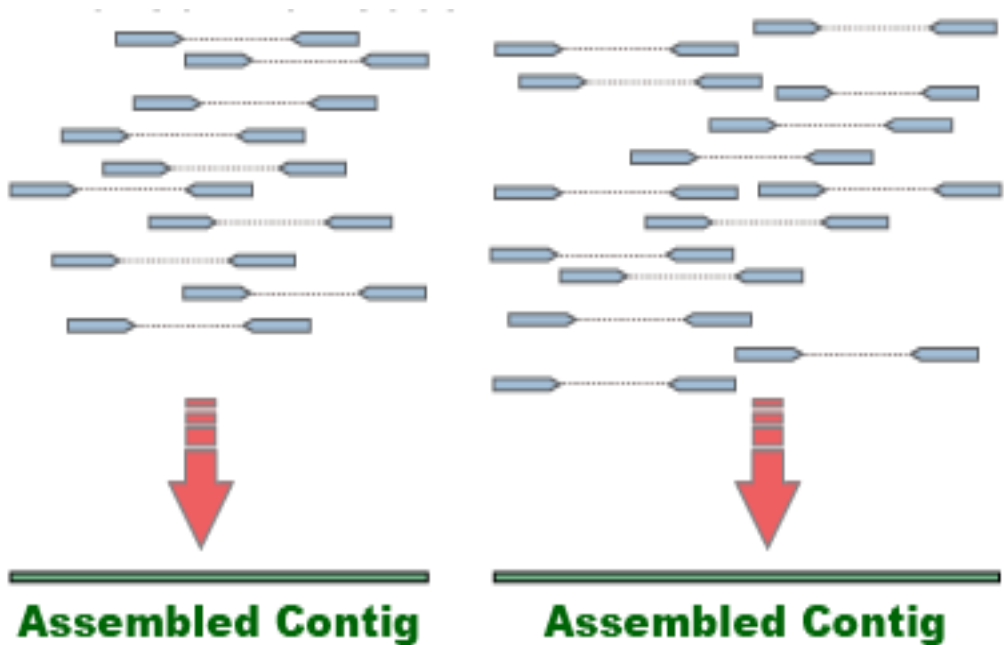
Evidence from multiple reads



Advantage of **paired reads**



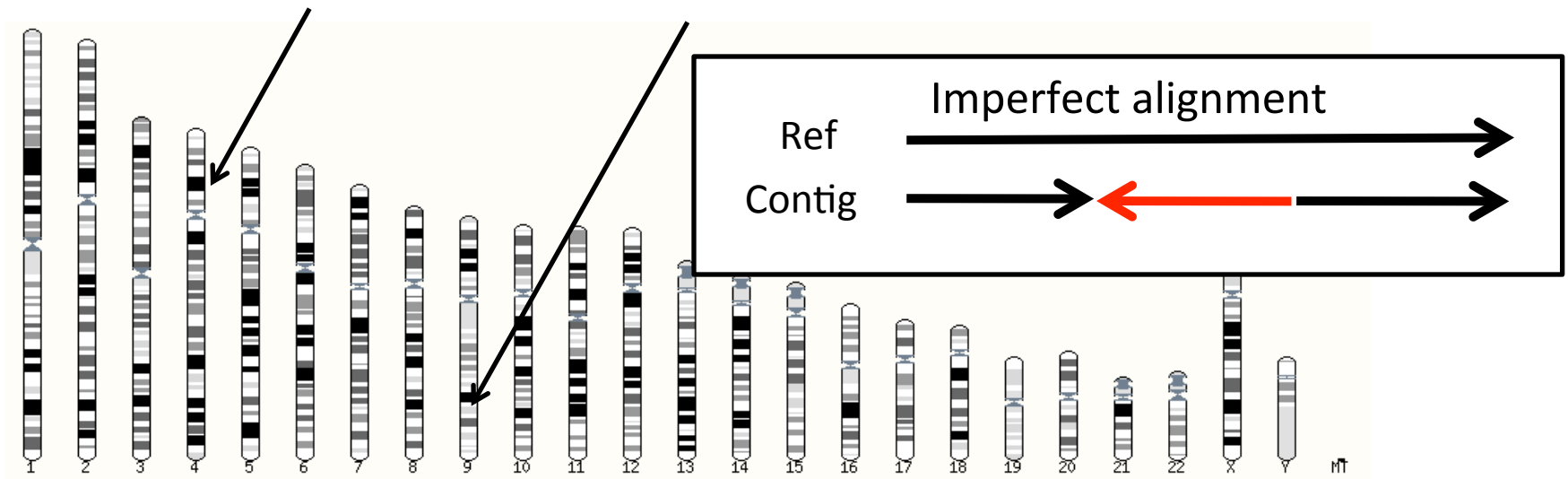
# De novo assembly for SV detection



**Scope:** various types of SVs including large inserts

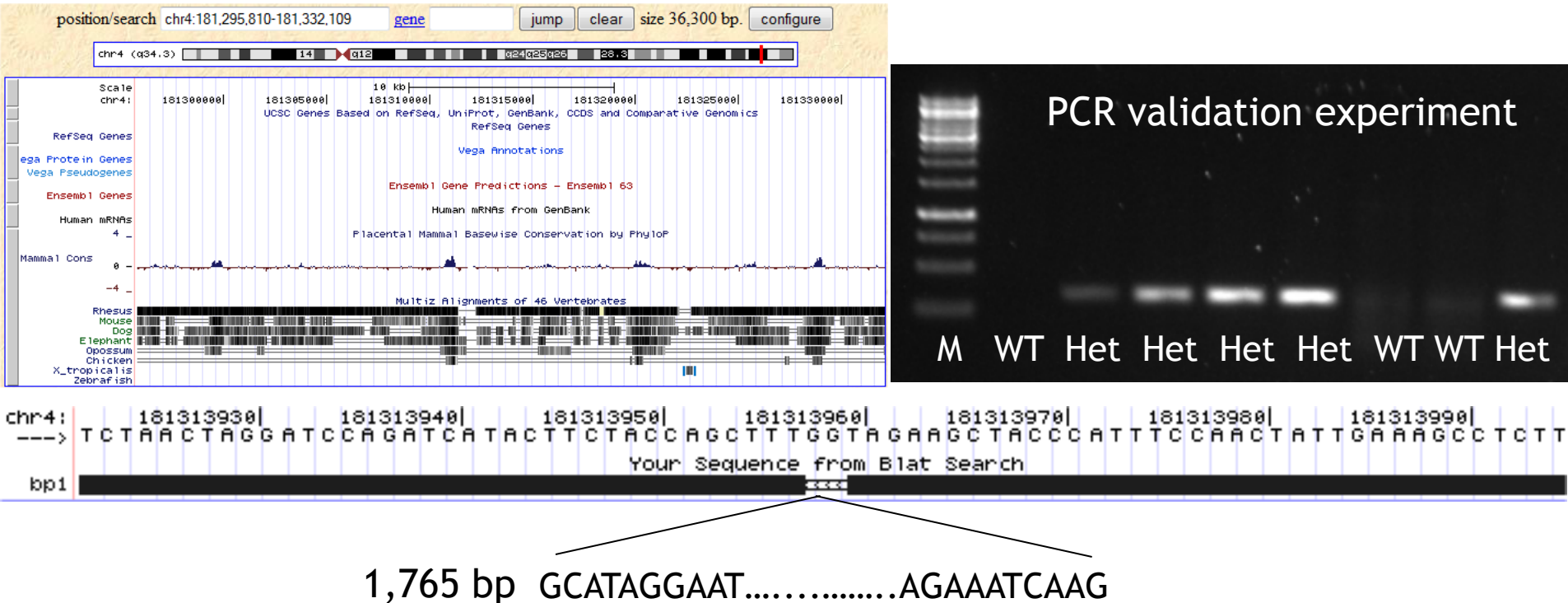
**Tools:** de novo assemblers  
SOAPdenovo, ABYSS, Allpaths-LG, Velvet

BLAST/BLAT search for comparison of contigs and genome reference



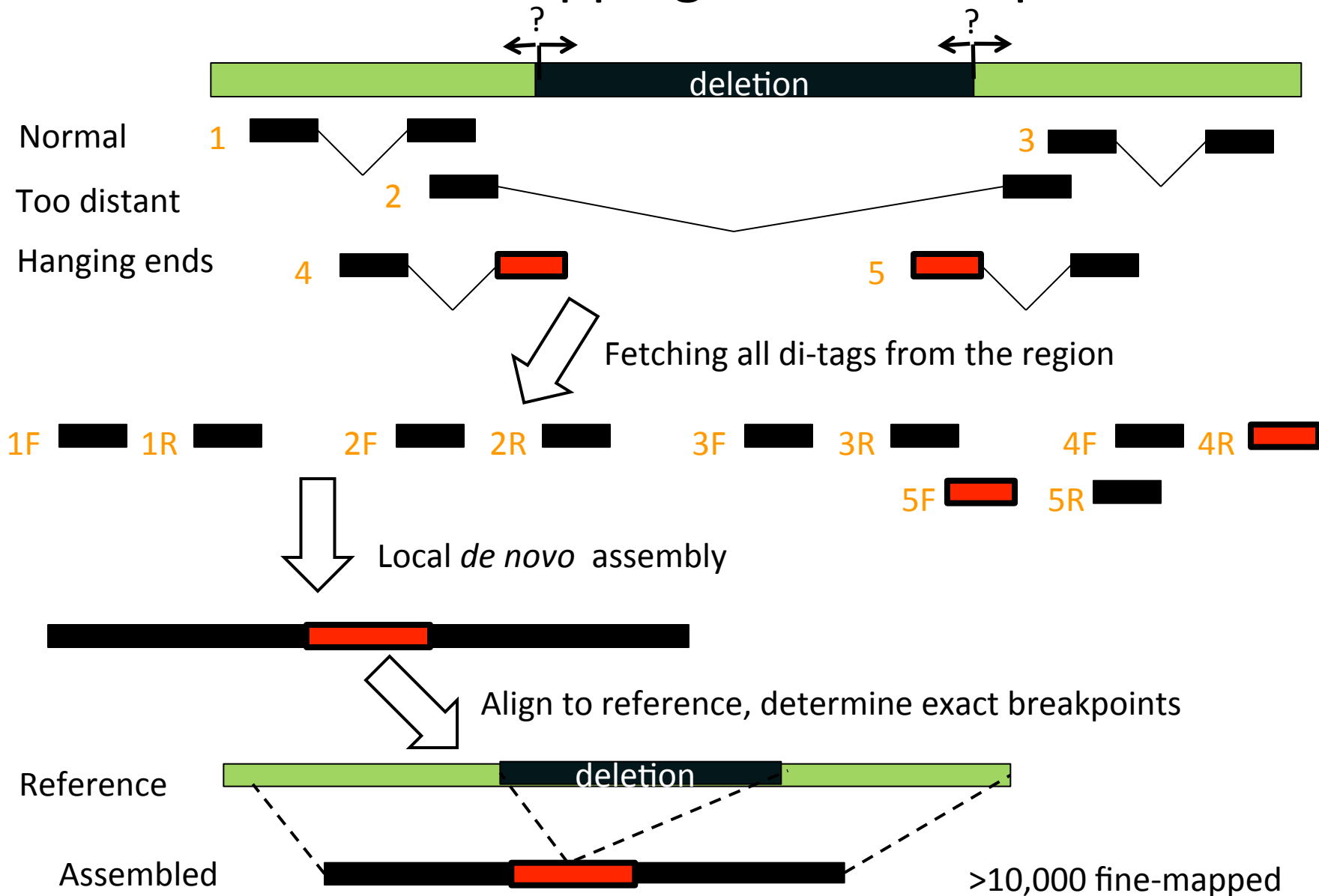
# *De novo* assembly: sample-specific segments

Comparison of individual *de novo* assembly of GoNL data to GRCh37  
 ~ 70 new regions (> 1kb), totaling 235 kb of NEW sequences



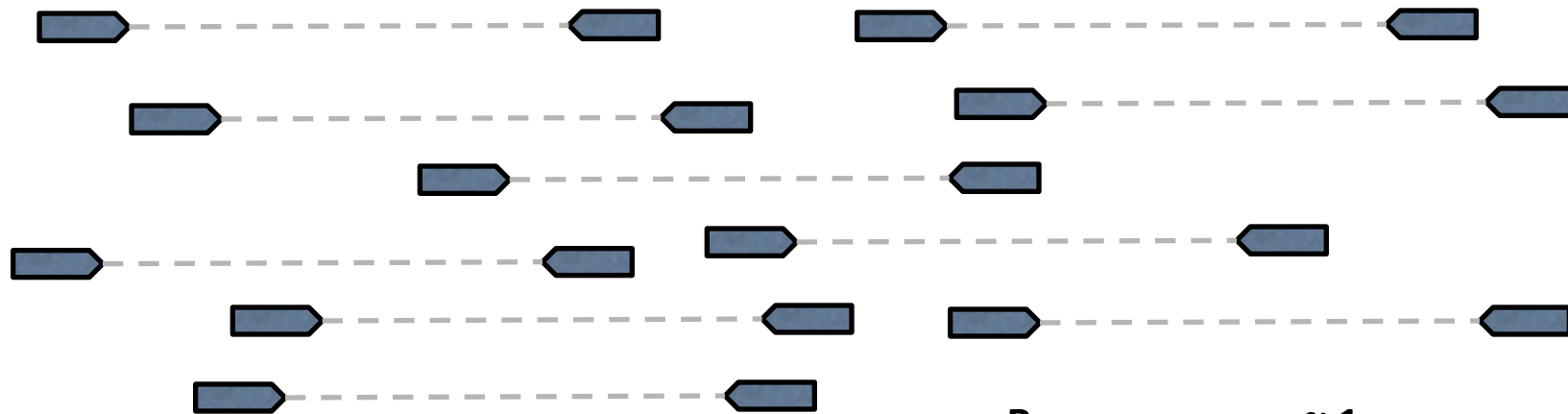
- No significant homology to known sequences on nucleotide and amino-acid levels
- Estimated frequency: ~42% in Dutch population, ~5% in 1000 Genomes

# Local *de novo* assembly: fine-mapping of SV breakpoints



# Base- and physical coverage

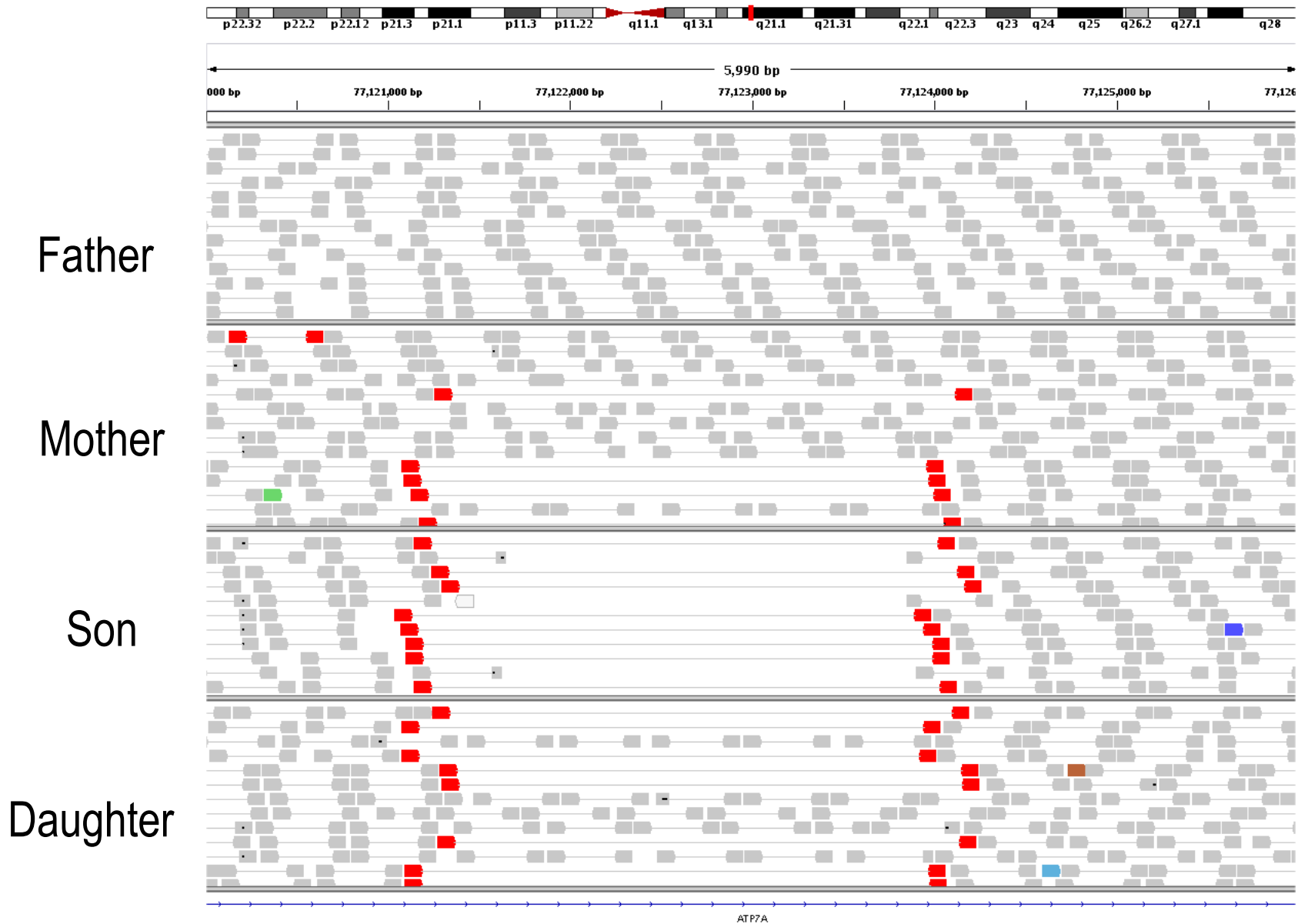
chromosome



Base coverage:  $\sim 1x$ ;  
Physical coverage  $\sim 4x$

Approach	Base coverage	Physical coverage
Depth of coverage	✓	
Discordant pairs		✓
Split-mapping	✓	
<i>De novo</i> assembly	✓	✓

# Composite patterns of SV



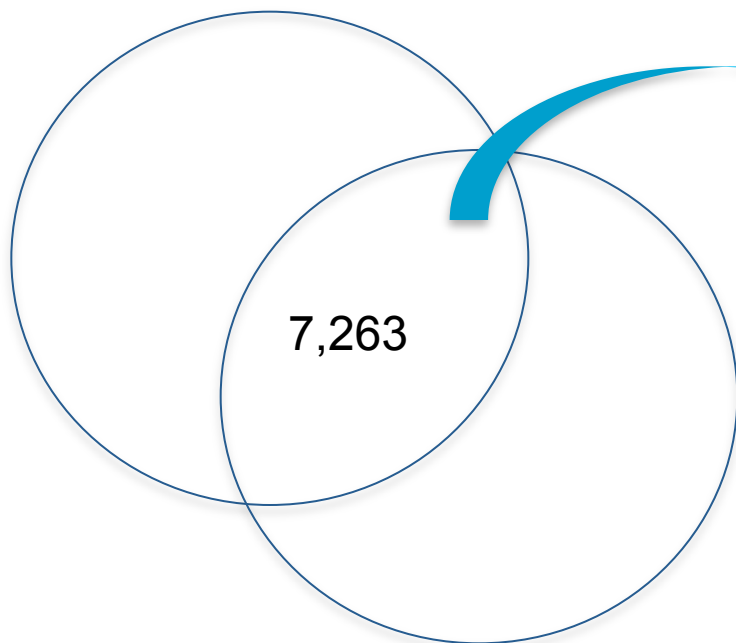
# Combining detection methods: a good idea!

## 1000 Genomes

A Deep Catalog of Human Genetic Variation

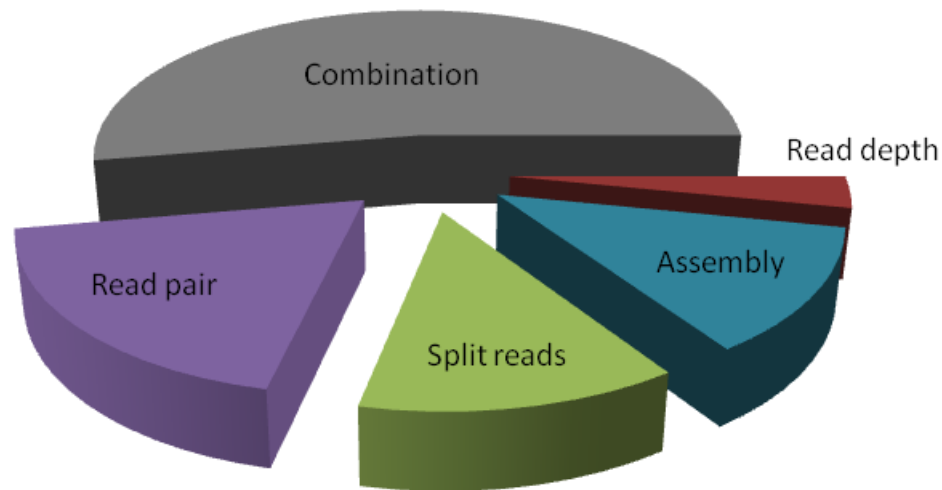
1000 genomes project ( phase 1 )

21,541 deletions ( >30bp )



**Go•NL**  
GENOMEoftheNETHERLANDS

## Detection methods



# NextGen Sequencing: what do we get?

	Genome of NL	1000 genomes
Individuals	769 (250 families)	1092
SNPs	19.8 M	36.7 M
Small indels	1.4 M	1.4 M

## Per individual genome (as compared to reference genome):

3.7M SNPs

360k short indels (1-20bp)

5.2k medium deletions ( 20 – 100 bp)

3.3k large deletions ( 100+ bp)

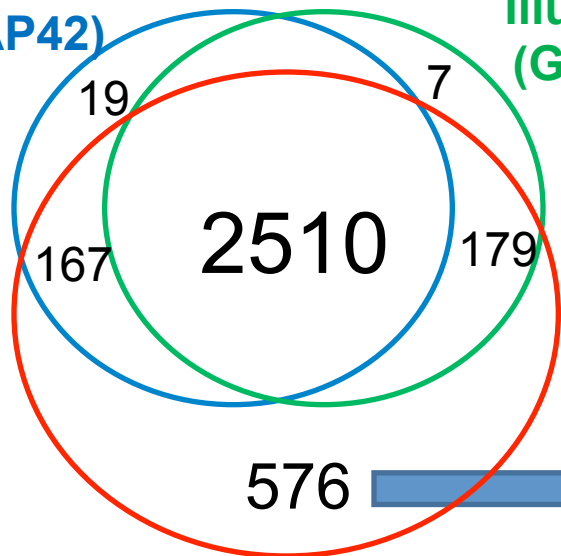


# NextGen Sequencing: what are we missing?

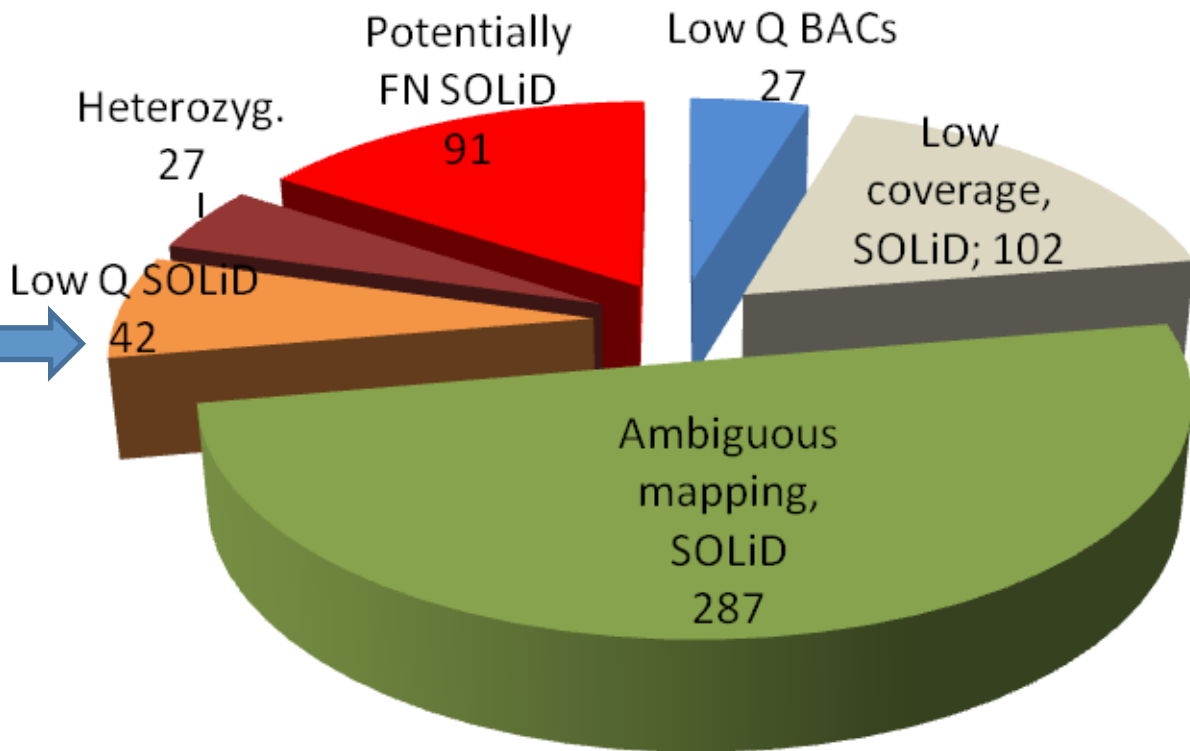
13 BACs (2.17 Mb)  
Sanger Sequencing, Assembly  
NGS: SOLiD and Solexa @ 20-25x

SOLiD  
(SAP42)

Illumina  
(GATK)



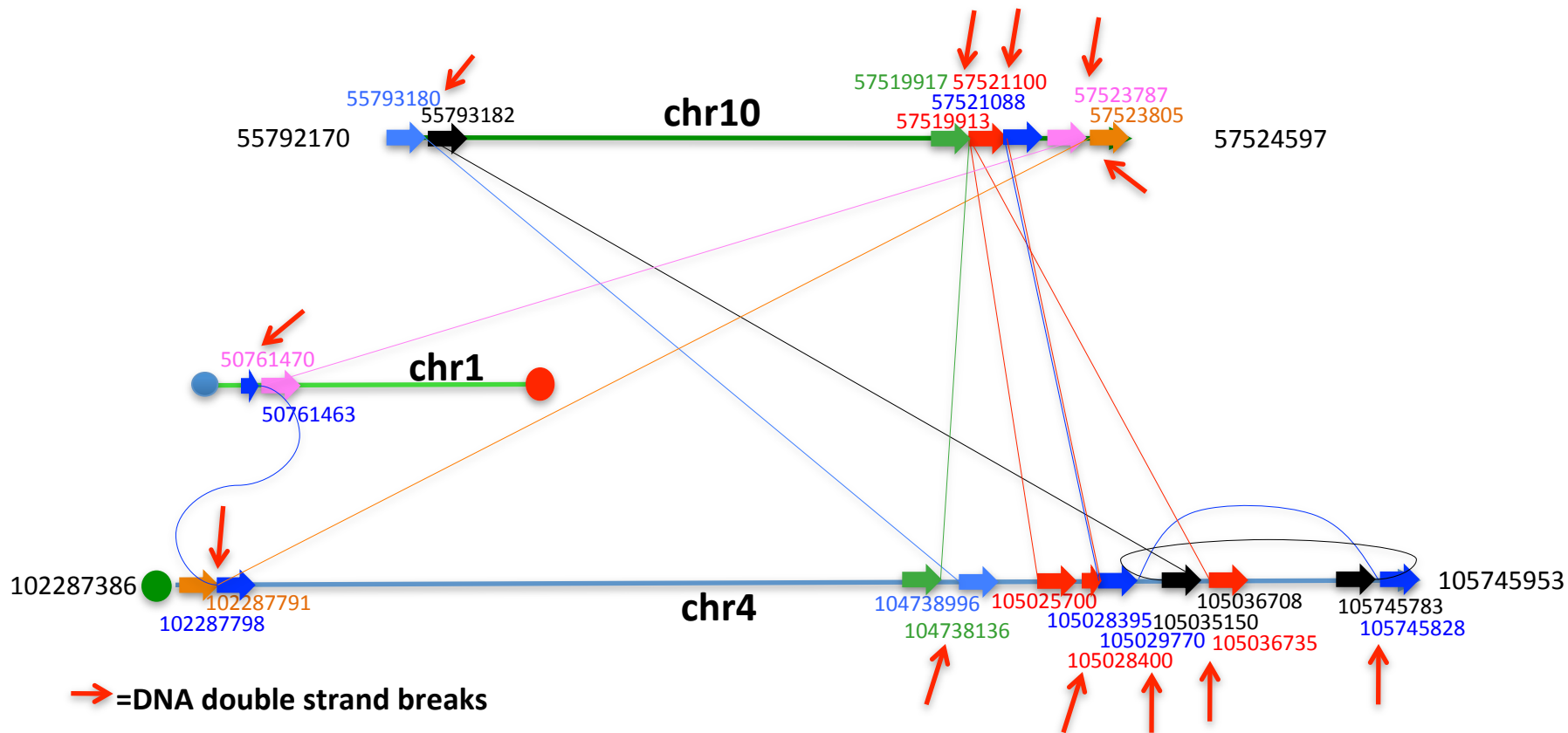
**Finished BACs**



[HS rats sequencing consortium, 2013]

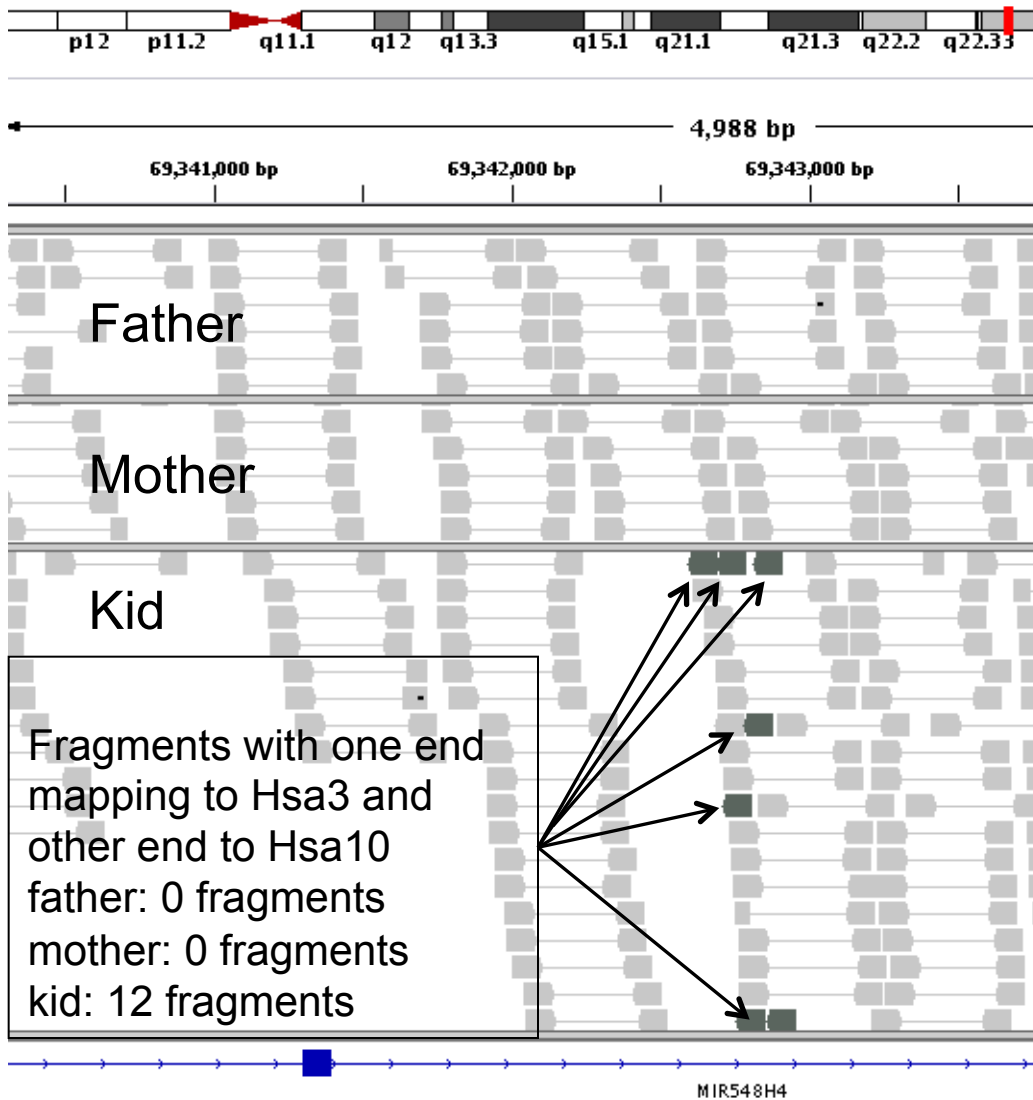
# How SVs arise?

# Complex structural variations, chromotrypsis

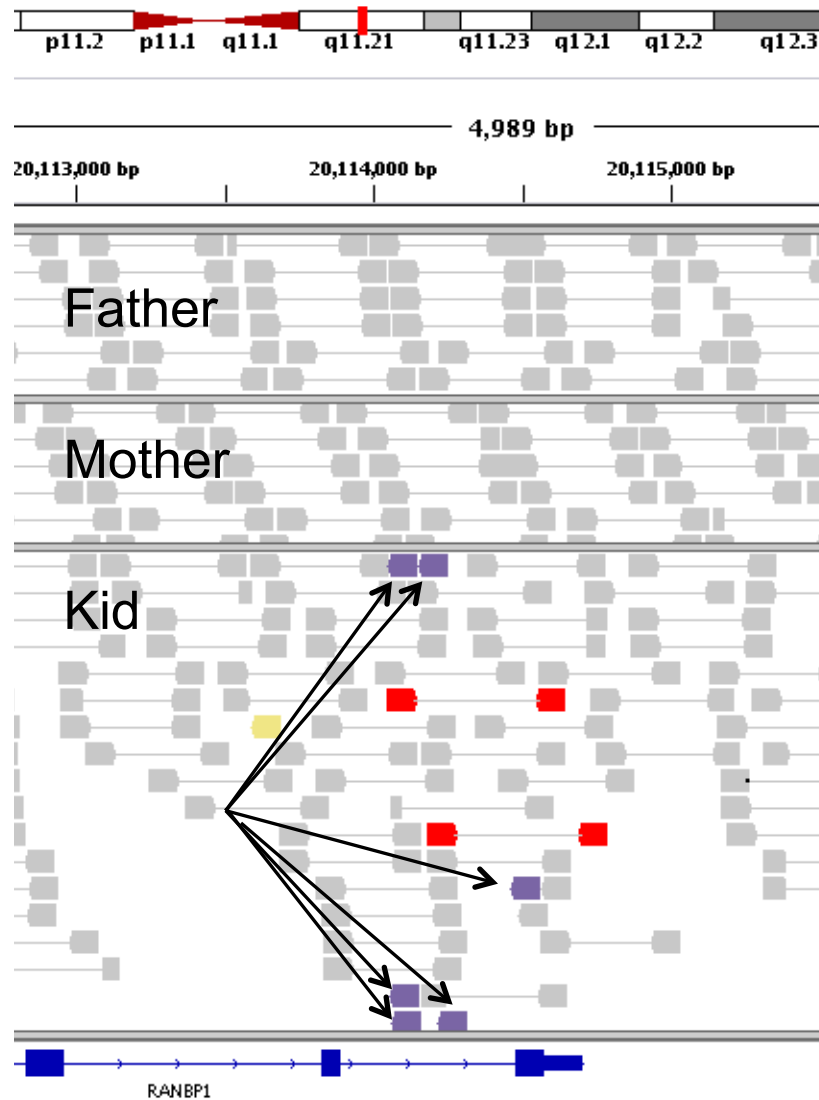


# De novo SVs in healthy individuals

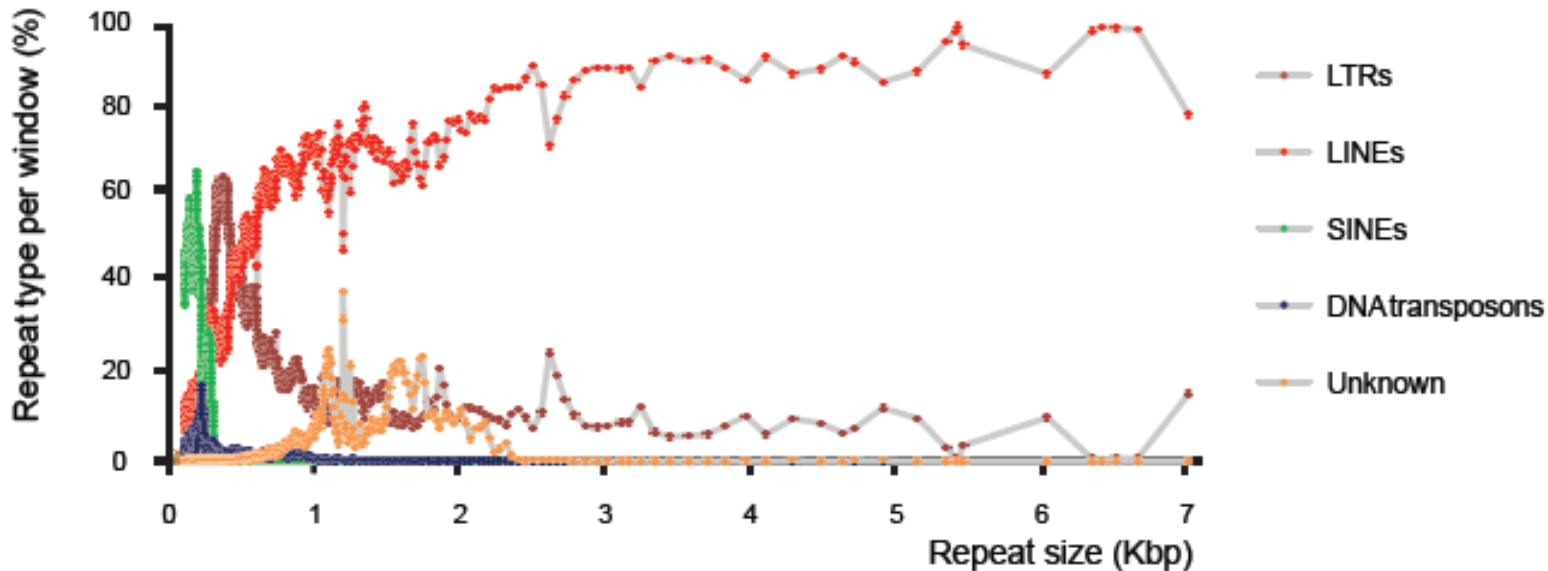
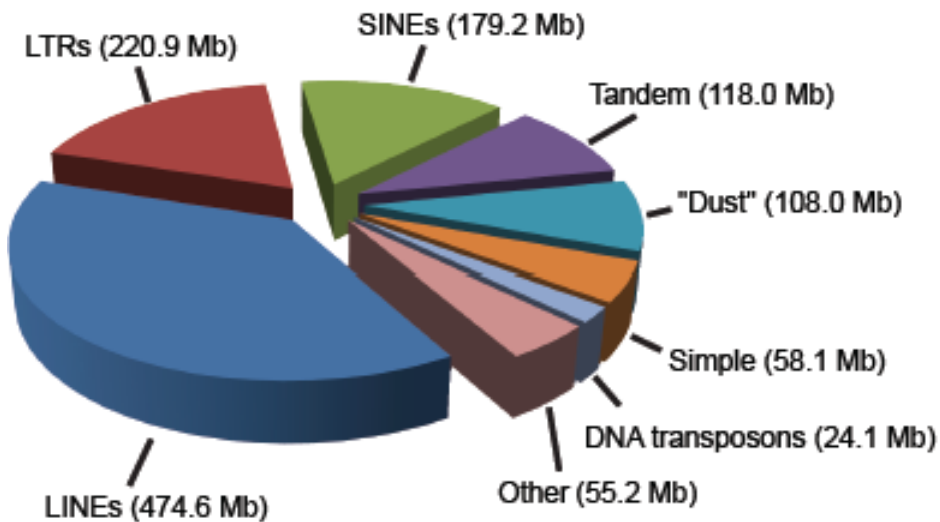
Hsa15



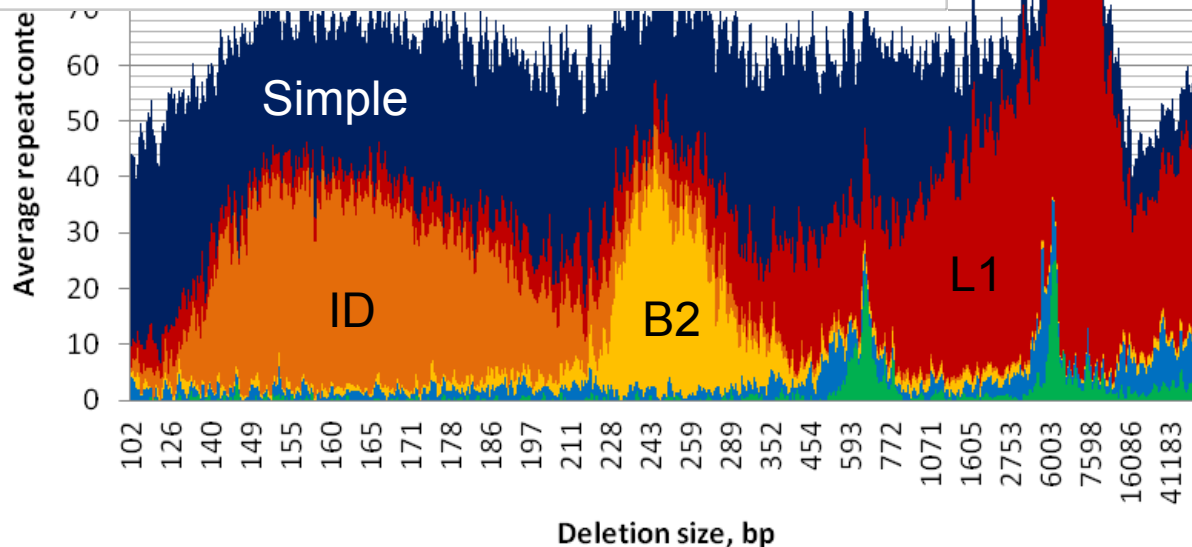
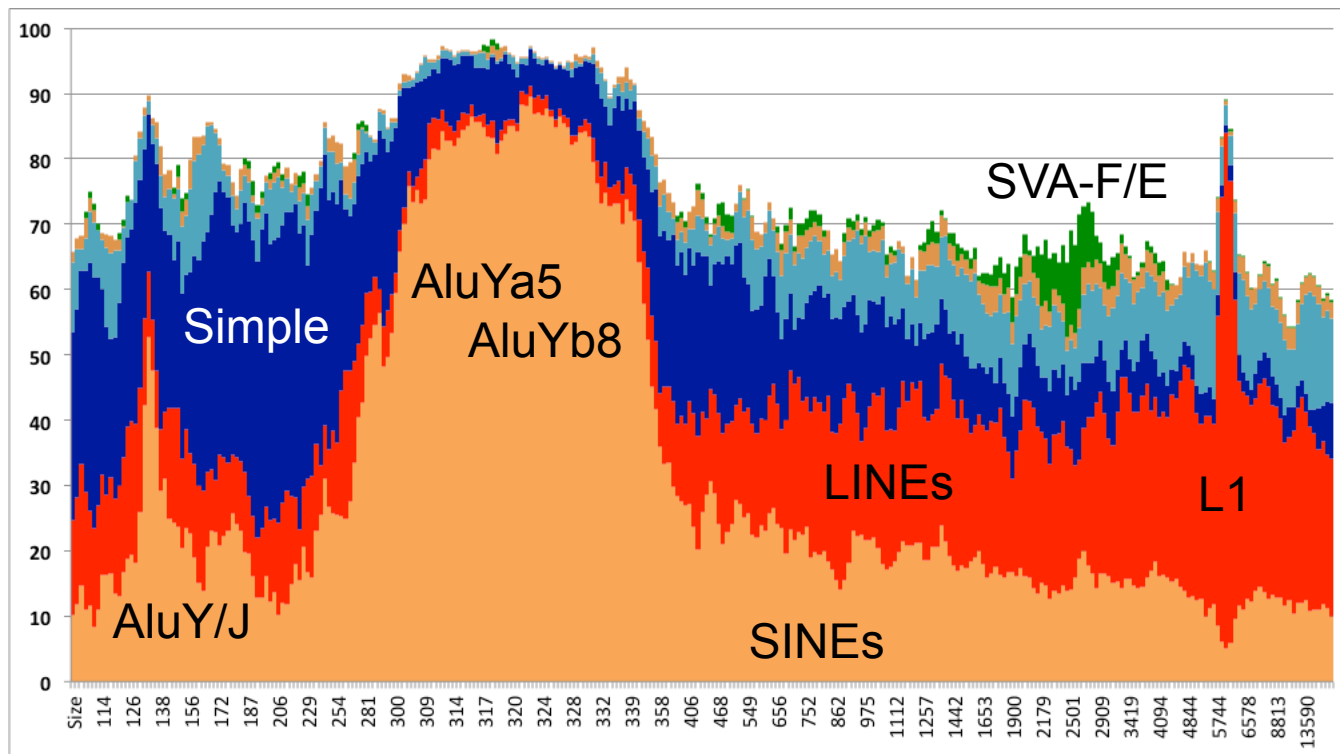
Hsa22



# Repeats in mammalian genomes



# Repeat instability: primary cause of SVs



- Simple repeats
- LINE/L1
- SINE/ID
- SINE/B2
- LTR/ERVK
- LTR/ERV1

# Representing SVs in VCF format

```
##ALT=<ID=DEL,Description="Deletion">
##ALT=<ID=DEL:ME:ALU,Description="Deletion of ALU element">
##ALT=<ID=DEL:ME:L1,Description="Deletion of L1 element">
##ALT=<ID=DUP,Description="Duplication">
##ALT=<ID=DUP:TANDEM,Description="Tandem Duplication">
##ALT=<ID=INS,Description="Insertion of novel sequence">
##ALT=<ID=INS:ME:ALU,Description="Insertion of ALU element">
##ALT=<ID=INS:ME:L1,Description="Insertion of L1 element">
##ALT=<ID=INV,Description="Inversion">
##ALT=<ID=CNV,Description="Copy number variable region">
##FORMAT=<ID=GT,Number=1,Type=Integer,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Float,Description="Genotype quality">
##FORMAT=<ID=CN,Number=1,Type=Integer,Description="Copy number genotype for imprecise events">
##FORMAT=<ID=CNQ,Number=1,Type=Float,Description="Copy number genotype quality for imprecise">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001
1 2827693 . CCGTGGATGCGGGGACCCGCATCCCCTCTCCCTTCACAGCTGAGTGACCCACATCCCCTCTCCCCTCGCA C . PASS \
SVTYPE=DEL;END=2827680;BKPTID=Pindel_LCS_D1099159;HOMLEN=1;HOMSEQ=C;SVLEN=-66 GT:GQ 1/1:13.9
2 321682 . T <DEL> 6 PASS IMPRECISE;SVTYPE=DEL;END=321887;SVLEN=-105;CIPOS=-56,20; \
CIEND=-10,62 GT:GQ 0/1:12
2 14477084 . C <DEL:ME:ALU> 12 PASS IMPRECISE;SVTYPE=DEL;END=14477381;SVLEN=-297;
MEINFO=AluYa5,5,307,+;CIPOS=-22,18;CIEND=-12,32 GT:GQ 0/1:12
3 9425916 . C <INS:ME:L1> 23 PASS IMPRECISE;SVTYPE=INS;END=9425916;SVLEN=6027;CIPOS=-16,22;\
MIINFO=L1HS,1,6025,- GT:GQ 1/1:15
3 12665100 . A <DUP> 14 PASS IMPRECISE;SVTYPE=DUP;END=12686200;SVLEN=21100;CIPOS=-500,500;\
CIEND=-500,500 GT:GQ:CN:CNQ ./.:0:3:16.2
4 18665128 . T <DUP:TANDEM> 11 PASS IMPRECISE;SVTYPE=DUP;END=18665204;SVLEN=76;CIPOS=-10,10;\
CIEND=-10,10 GT:GQ:CN:CNQ ./.:0:5:8.3
```

# Variant annotation: Variant Effect Predictor (VEP)

Ensembl: Data Tools - Assembly converter, ID History converter, Variant Effect Predictor

www.ensembl.org/tools.html

Apple iCloud Facebook Twitter Wikipedia Yahoo! News Popular

Guryev, V (eriba) - Outlook Web App | Inbox (3,826) - victor.guryev@gmail.com - Gmail | Ensembl: Data Tools - Assembly converter, ID Hi... | Capture a Screen Shot with Mac OS X

Ensembl BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors

Search all species...

**Custom Data**

- Add your data
- Attach DAS
- Manage Data
- Features on Karyotype
- Data Converters
  - Assembly Converter
  - ID History Converter
  - Variant Effect Predictor**
  - Region Report

**Variant Effect Predictor:**

This tool takes a list of variant positions and alleles, and predicts the effects of each of these on overlapping transcripts and regulatory regions annotated in Ensembl. The tool accepts substitutions, insertions and deletions as input, see [data formats](#).

Upload is limited to 750 variants; lines after the limit will be ignored. Users with more than 750 variations can split files into smaller chunks, use the standalone [perl script](#) or the [variation API](#). See also [full documentation](#)

**NB:** Ensembl now by default uses Sequence Ontology terms to describe variation consequences. See [this page](#) for details

**Input file**

Species: Human (Homo sapiens): GRCh37

Name for this data (optional):

Paste data:

```
1 881907 881906 -/C +
5 140532 140532 T/C +
```

Upload file: Choose File no file selected

or provide file URL:

Input file format: Ensembl default

**Options**

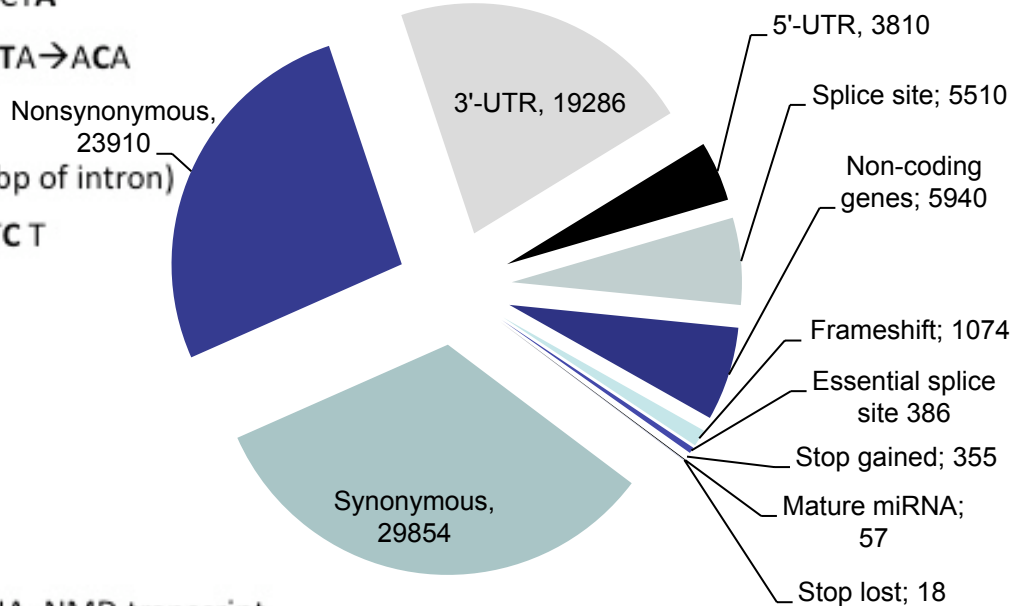
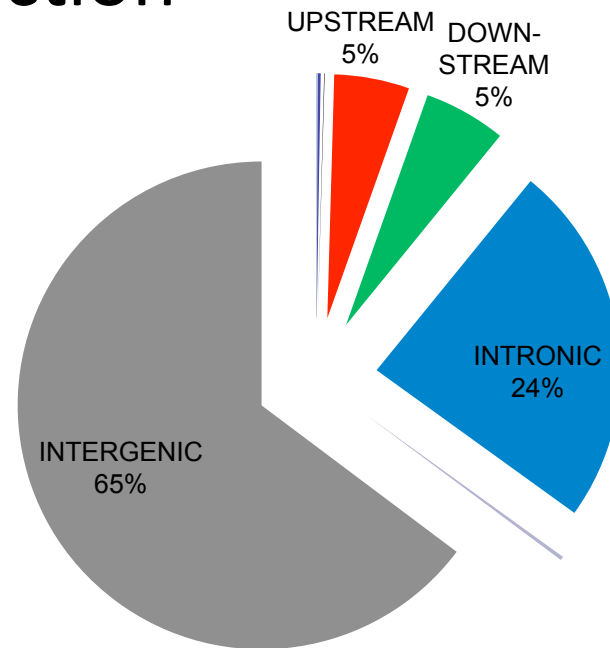
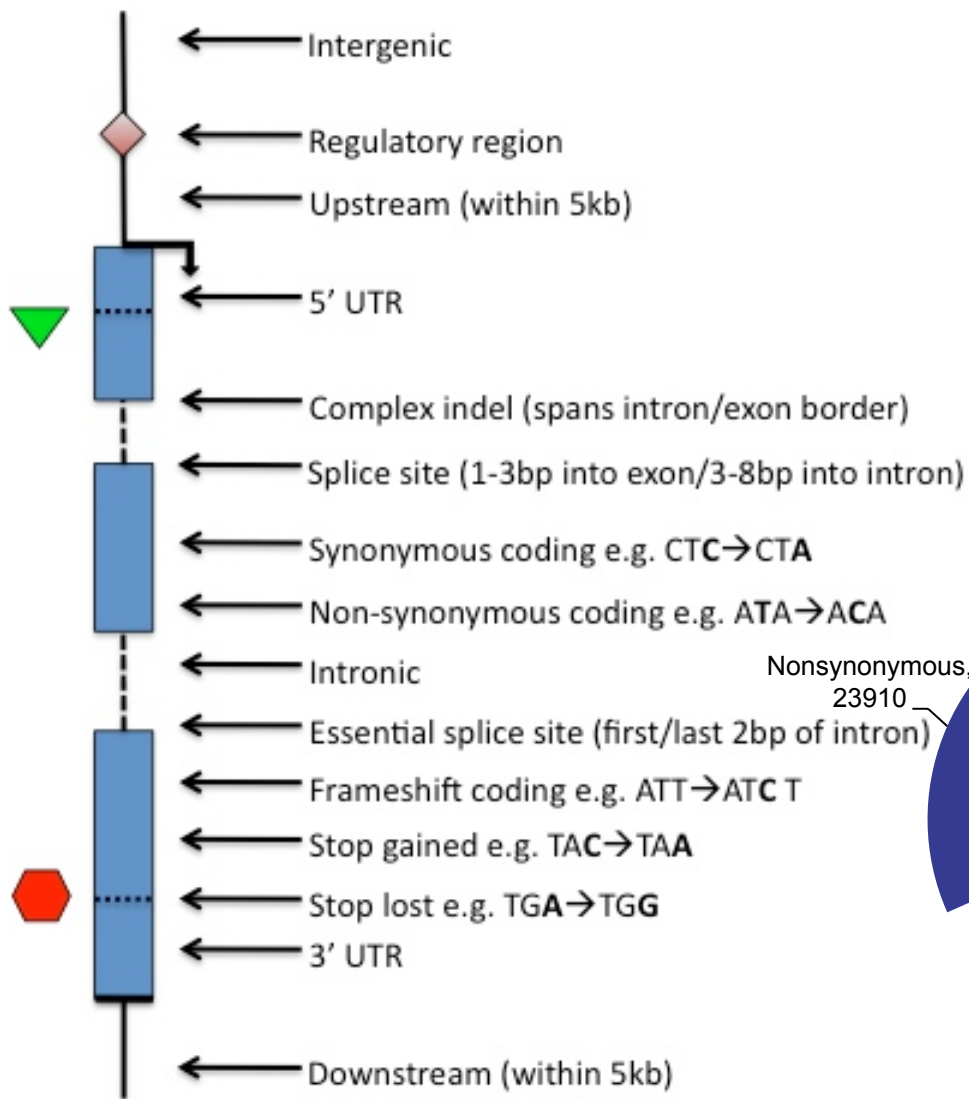
Transcript database to use:

- Ensembl transcripts
- RefSeq and other transcripts

Get regulatory region consequences (human and mouse only):

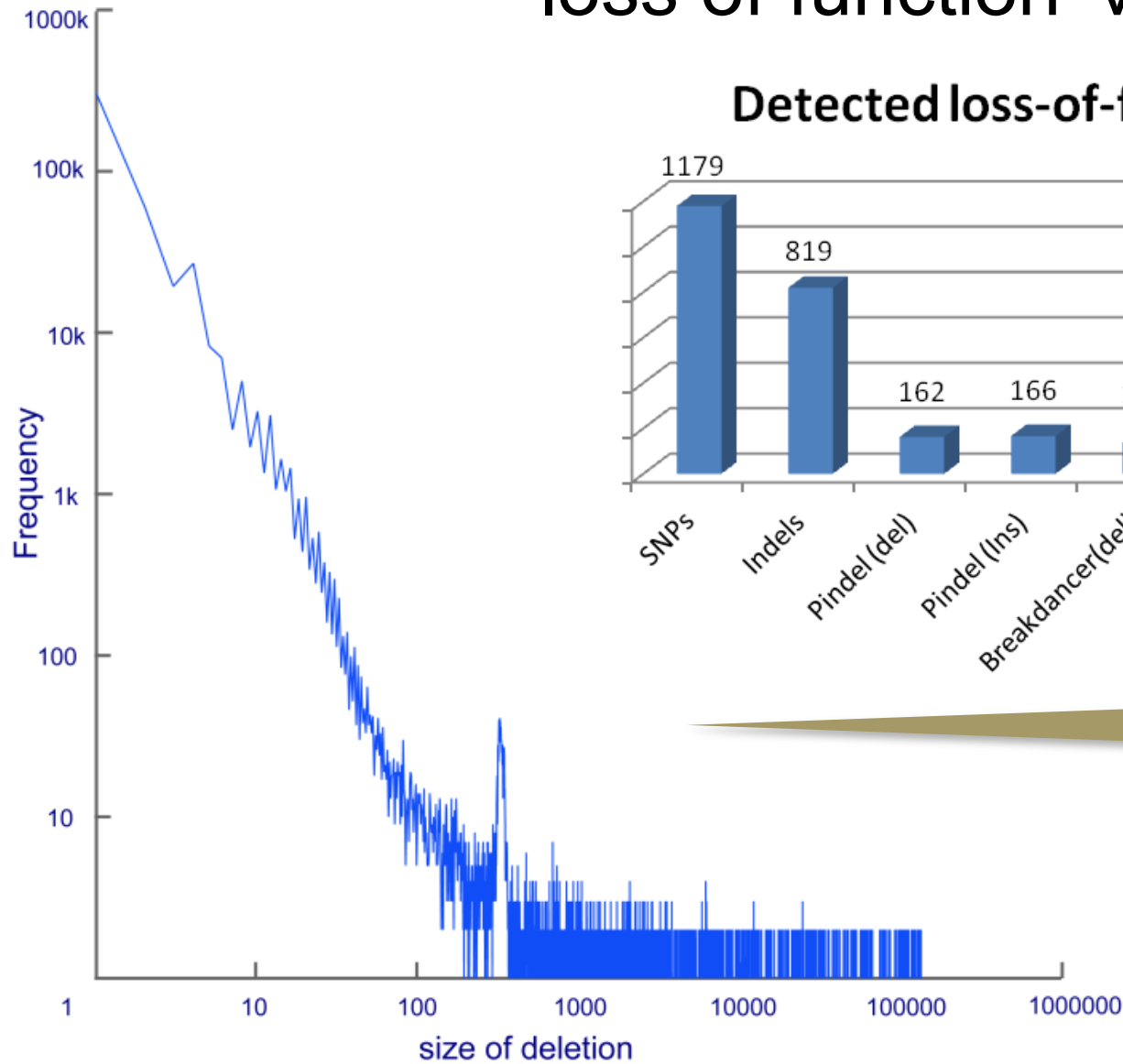


# Genome variation, function

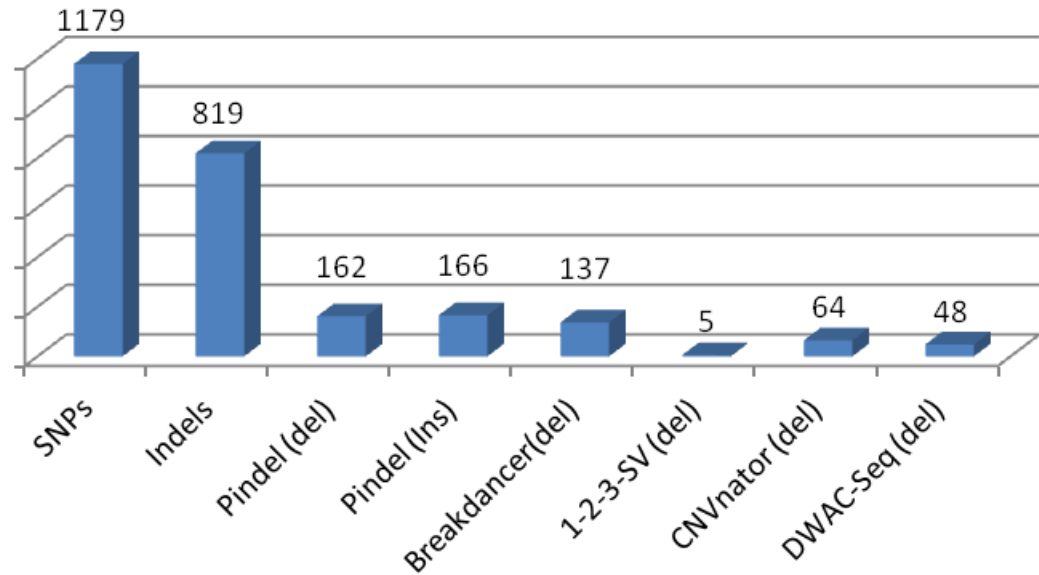


Others: Within non-coding gene, Within mature miRNA, NMD transcript

# Contribution of different variation types to 'loss of function' variants



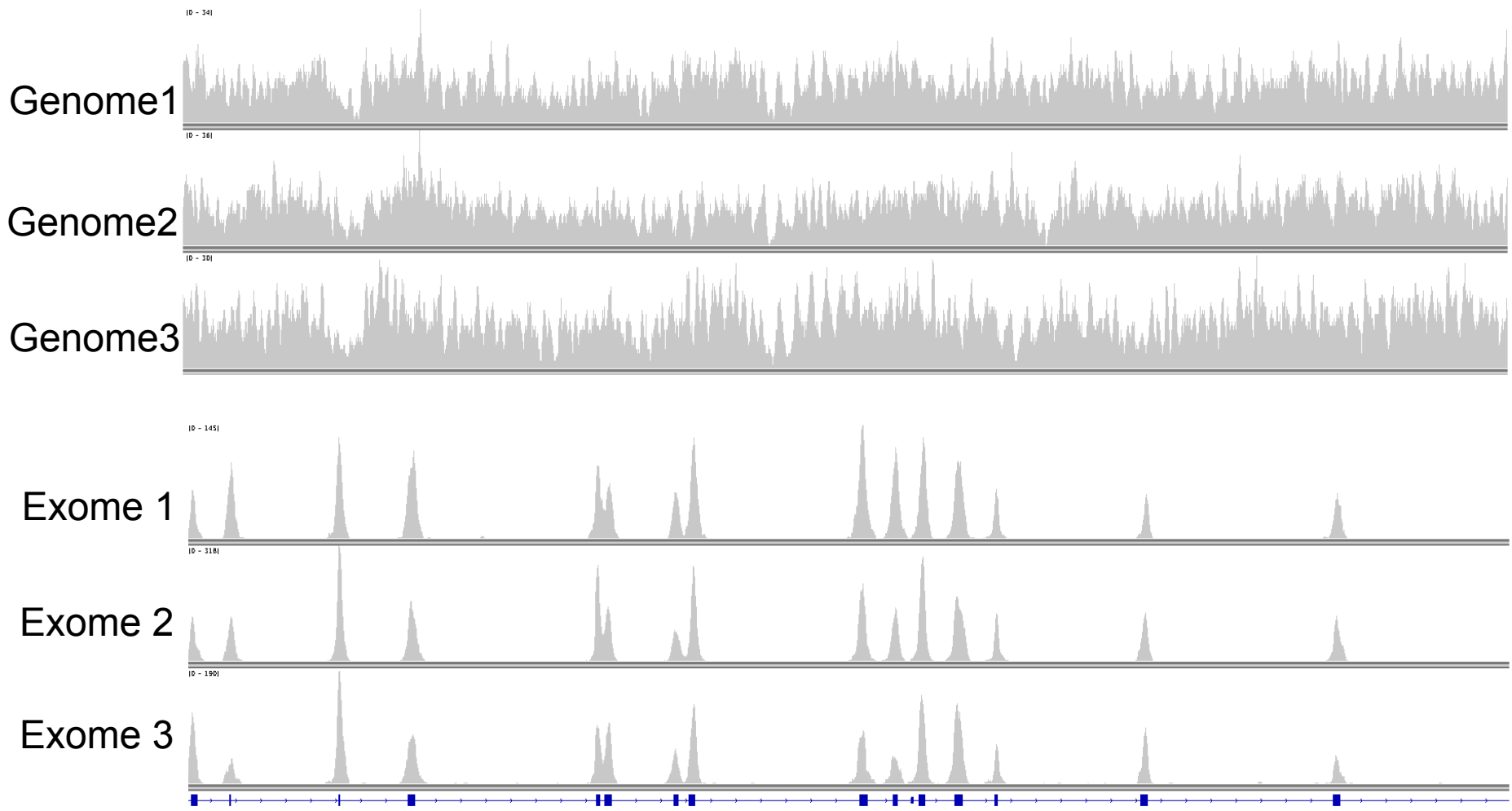
## Detected loss-of-function alleles



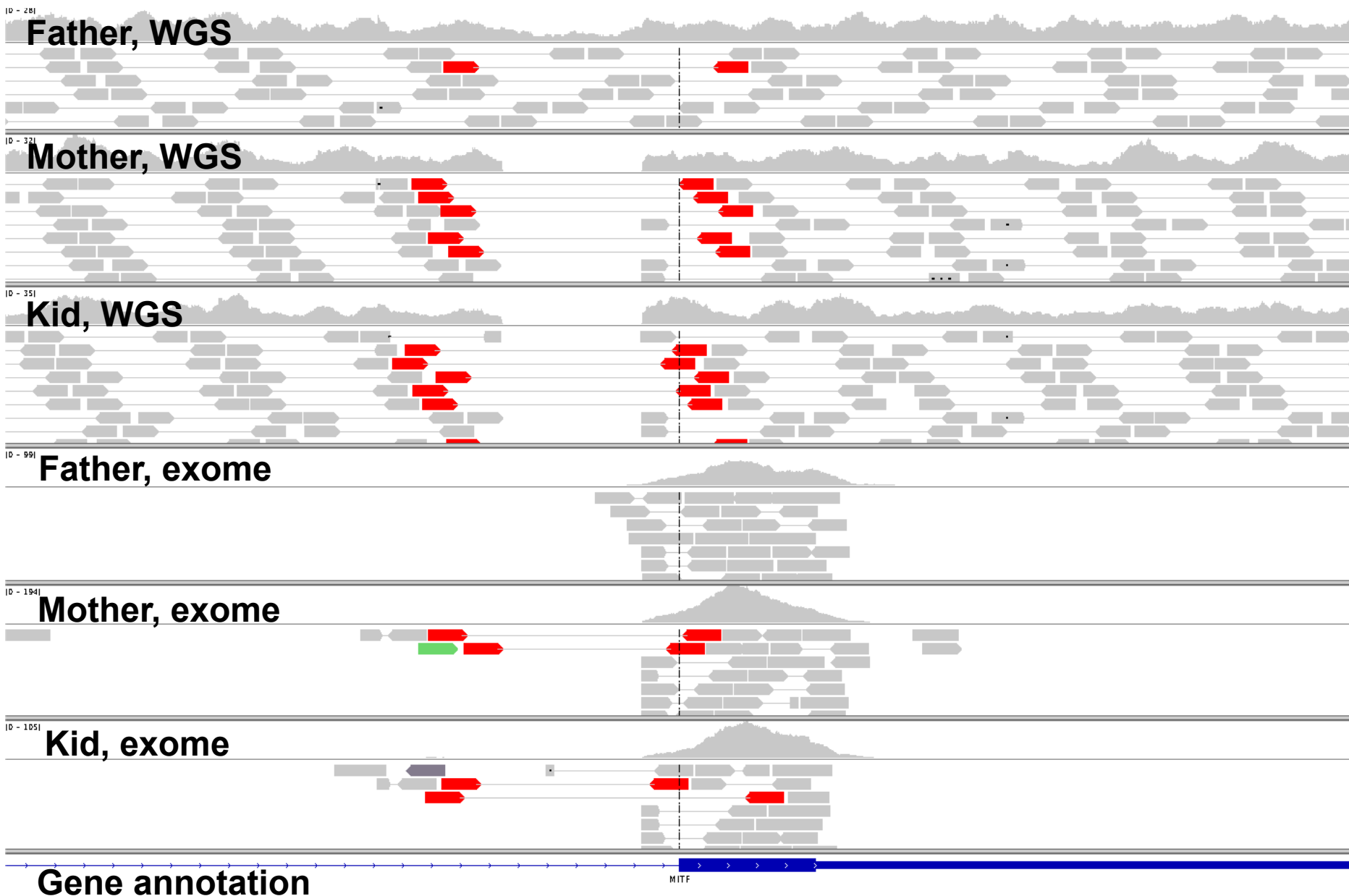
Variant size

Why read whole genome if we only know function of a small proportion of our genomes?

# Whole-genome sequencing vs exome-seq



# Catching SVs with exome data



# Validation strategies

## **Sanger Sequencing** (golden standard)

*de novo* variants

false discovery rate

false negative rate

## **Targeted resequencing by enrichment** ( e.g. Ion Torrent, Fluidigm )

loss-of-function or non-synonymous alleles

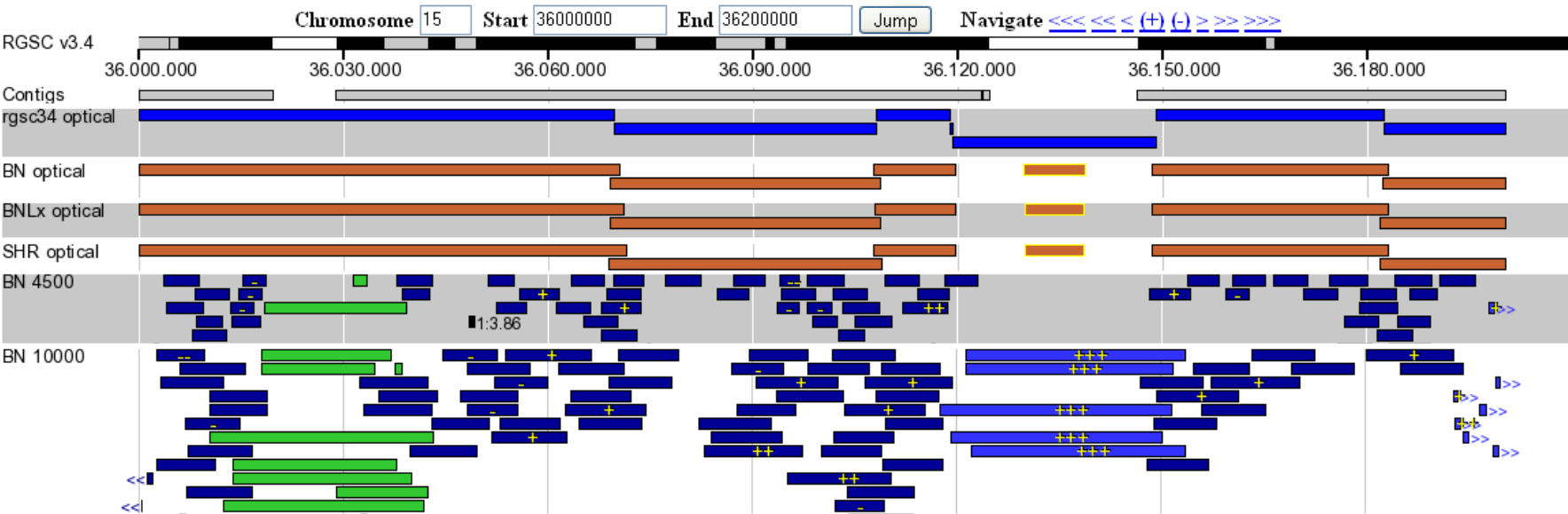
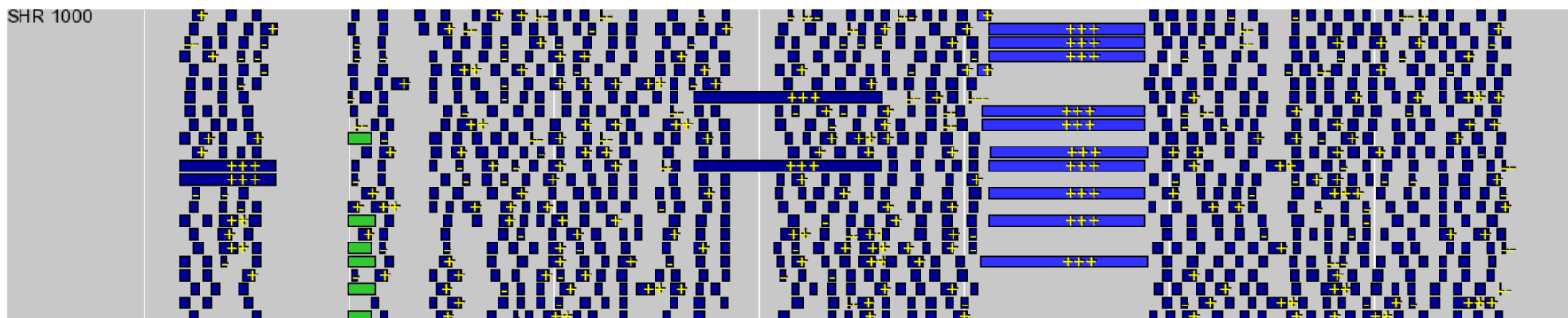
## **SV breakpoint array**

large (>30 bp) SVs with known breakpoints

## **FISH, aCGH**

large SVs and translocations

# Assembly quality is critical



# Take home messages

1. **Prerequisites:** Genome assembly quality
2. **Study design:** # of libraries, type(PE, MP), insert sizes
3. **Quality control:** Insert size distribution, chimerism
4. **Combine methods** for SV discovery: read depth, read-pair, split-reads, *de novo* assembly
5. **Do verifications** (high FP rates)



# Future directions

## **Longer reads:**

PacBio ( long reads, relatively low throughput )

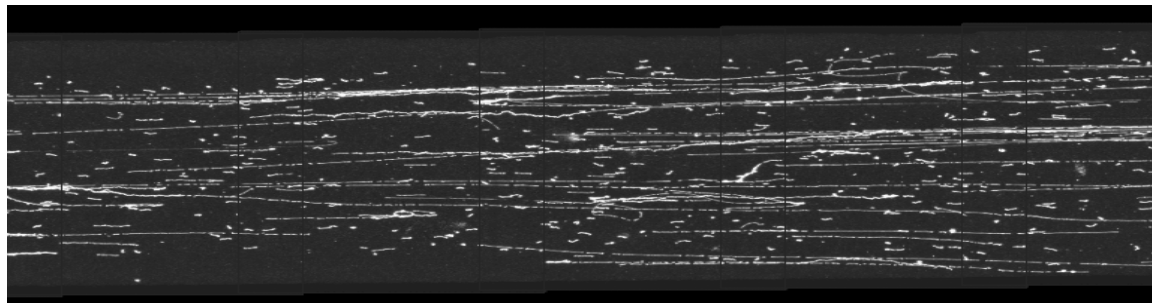
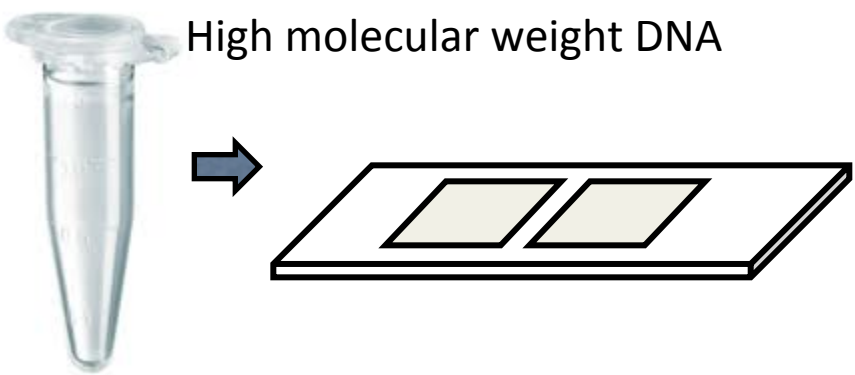
Oxford nanopore ( very long reads, easy preps )

## **Lower prices:**

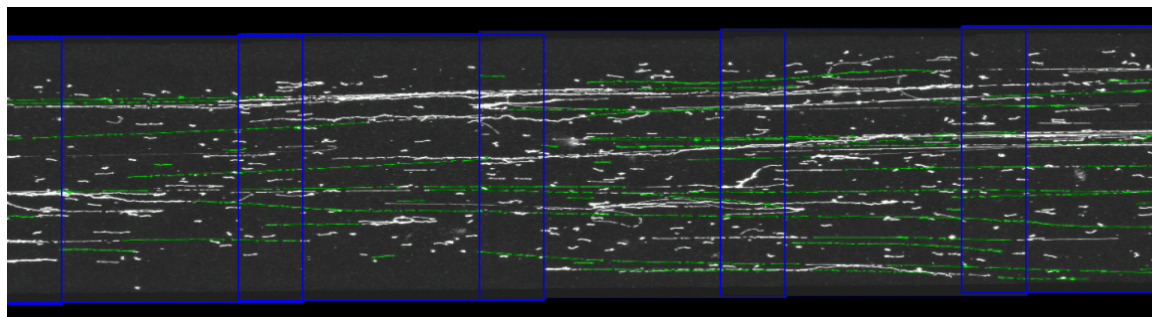
Higher coverage = better calls, more samples, whole genome sequencing

## **Even better and faster data analysis**

# Optical mapping



Swa I restriction of genomic DNA



~ 20% missing cuts

1-5 extra "cuts" / Mbp

# Acknowledgements



UMC Utrecht

Cuppen Group



Hubrecht  
Institute

Human genomes

Genome of the Netherlands



**Vacancies:** PhD student,  
bioinformatics technician

**Topic:** Genome structure and ageing

**Contact:** [v.guryev@umcg.nl](mailto:v.guryev@umcg.nl)



netherlands  
bioinformatics  
centre