# High-throughput sequencing

Robert Lyle

Department of Medical Genetics

Oslo University Hospital Ullevål

Robert.Lyle@medisin.uio.no

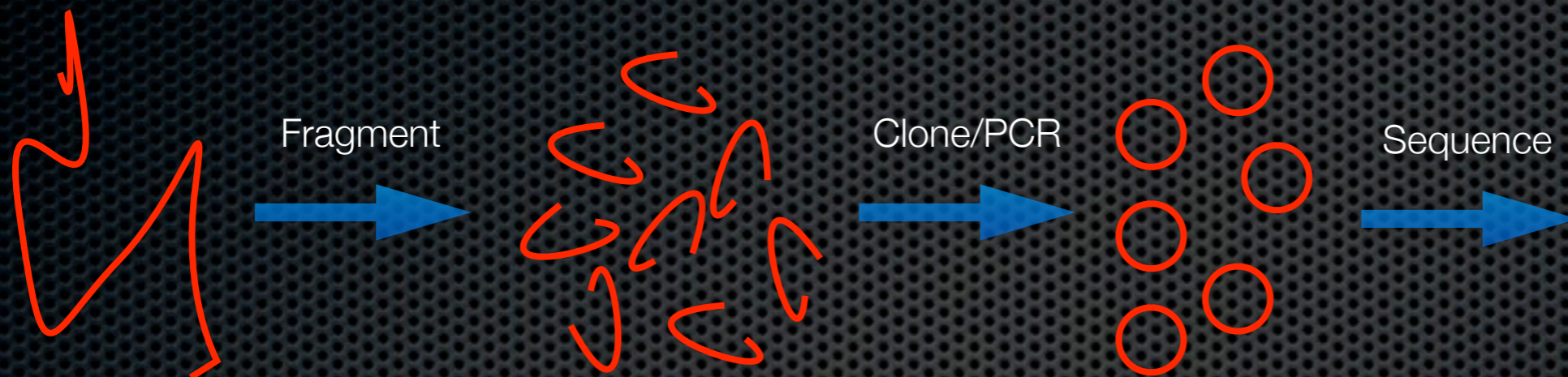# Overview

- Technology
- Data and analysis
- Applications

# Technology

Sequencing past, present and future

# Sequencing: old and next

LTS

Fragment → Clone/PCR → Sequence

**Molecules sequenced**

1, 48, 96...

...unless you have a lot of machines

HTS

Fragment → Array → Sequence

$4 \times 10^5$ - $1 \times 10^9$

...on one machine

**Massively parallel**

# HTS systems available

454      Solexa      SOLiD      HeliScope
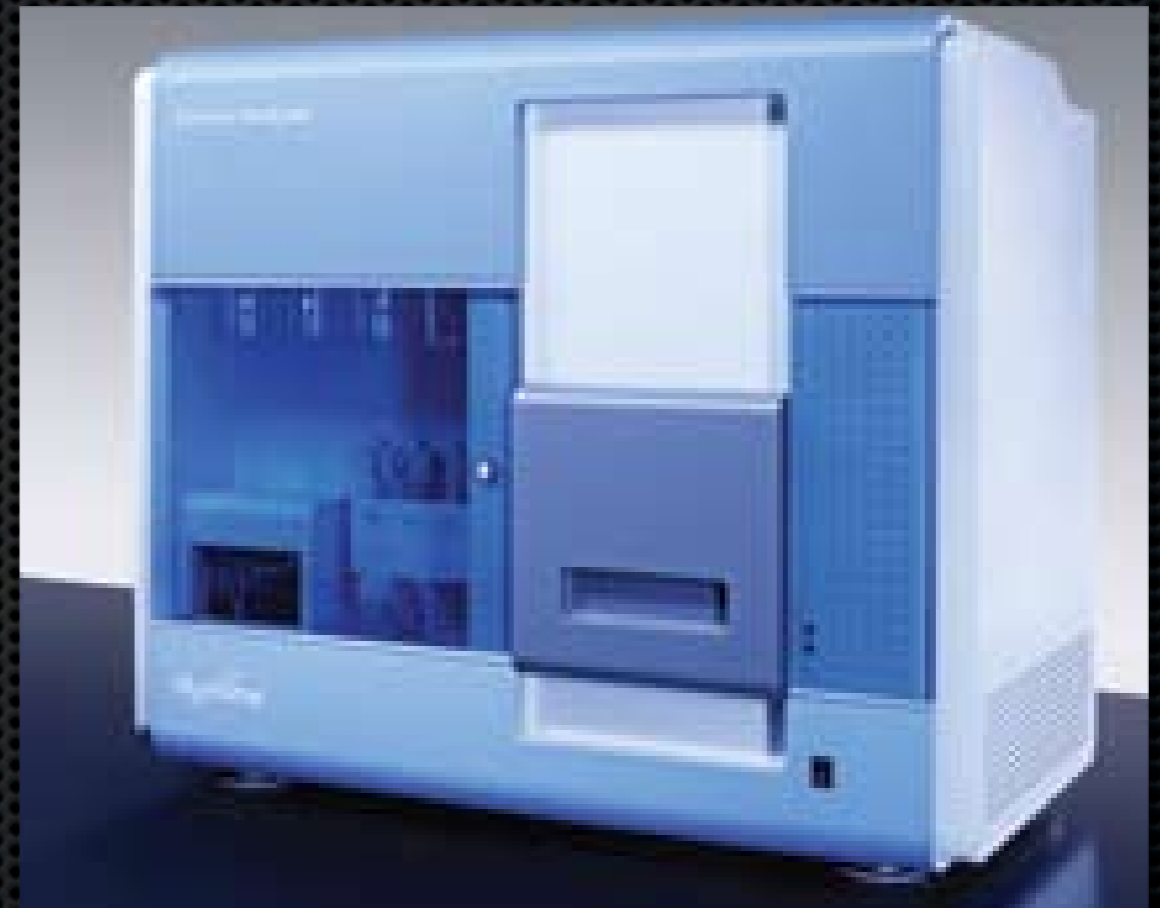


Roche      Illumina      ABI      Helicos

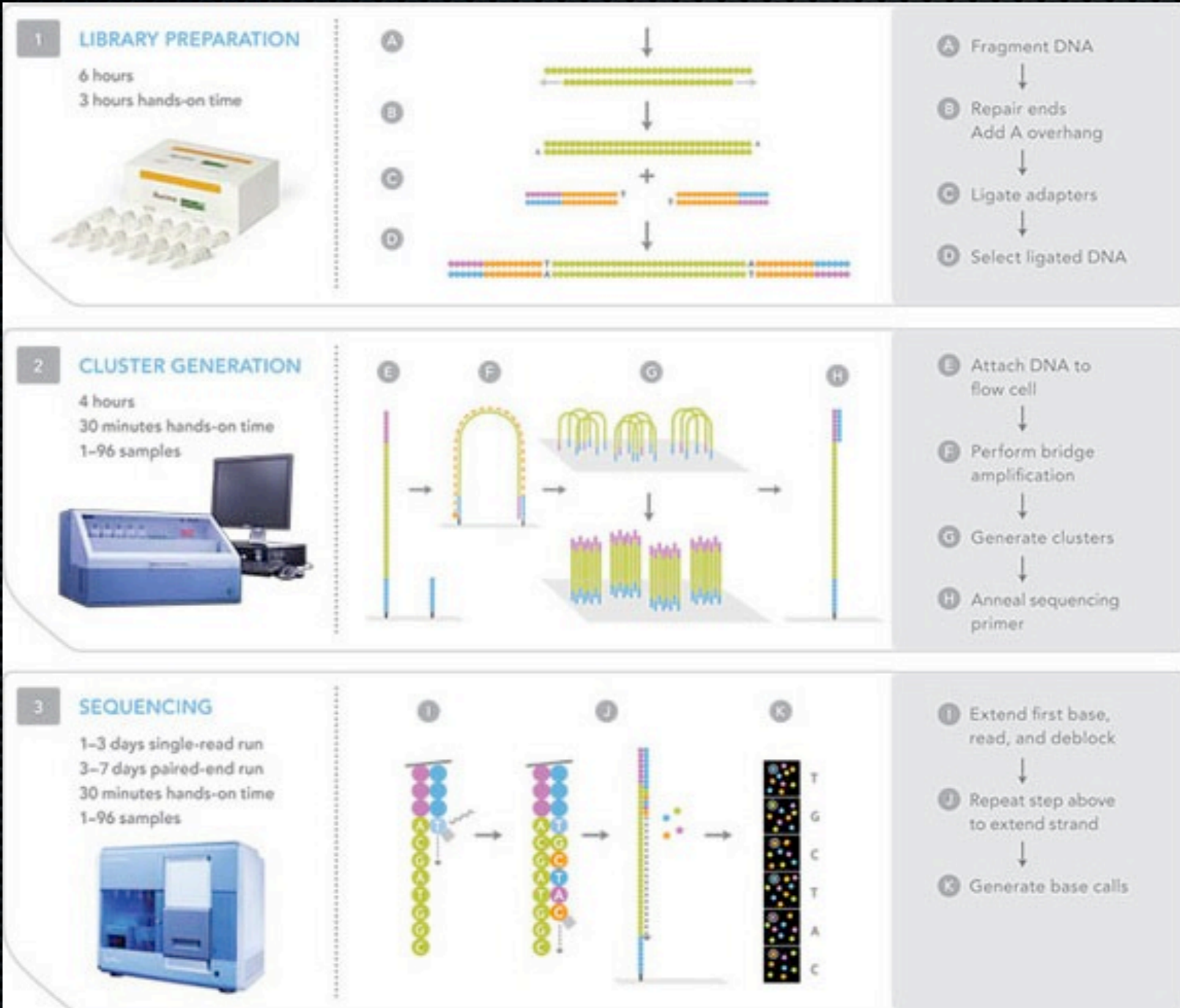## Others in 2010

# High-throughput sequencing core facility
## Department of Medical Genetics, Oslo University Hospital Ullevål

Illumina Genome Analyzer II
- Cluster station
- Paired-end module
- iPAR
- Genome Analyzer Pipeline

- Up to 160 M reads
- Single-end or paired-end
- Read lengths of 18, 35, 50, 75
- Total output up to 15 GB per run



UNIVERSITETET I OSLO

ULLEVÅL universitetssykehus

HELSE SØR-ØST

# www.med.uio.no/ulleval/medgen/sequencing



# sequencing@medisin.uio.no

# Illumina sequencing technology

# Illumina Genome Analyzer IIx and beyond

- Hardware upgrade
  - Increased reagent volume
  - Improved scanning
- Software
  - Improved cluster detection algorithms

Installation
August 2009

- Flow cell
  - Ordered arrays
  - Submicron features

End 2009

# Illumina throughput

# Sequencing throughput in practice

| Genome sequenced (publication year) | HGP (2003) | Venter (2007) | Watson (2008) |
|---|---|---|---|
| Time taken (start to finish) | 13 years | 4 years | 4.5 months |
| Number of scientists listed as authors | > 2,800 | 31 | 27 |
| Cost of sequencing (start to fi nish) | $2.7 billion | $100 million | < $1.5 million |
| Coverage | 8–10 × | 7.5 × | 7.4 × |
| Number of institutes involved | 16 | 5 | 2 |
| Number of countries involved | 6 | 3 | 1 |

ABI/SOLiD - Yoruban (12x) - $60 000

Illumina (2009) - Hsa 30x - $30 000 (?)

# 3rd generation technologies

* Intelligent Biosystems

* Visigen

* Oxford Nanopore

* Reveo

* ZS Genetics

* Complete Genomics (sell whole human genomes in 2009 for $5,000?)

* Pacific Biosciences



Single molecule sequencing
(no amplification bias)

# Data and analysis

# Illumina sequence data

* Random DNA library of short fragments    ~300 bp

* ~100-200 million DNA sequences

* 18, 36, 50, 75, 125 bp long

* Single-end reads

* Paired-end reads

* Run time: 1-10 days

* Data volume: 300 GB.....8 TB

# Data issues

* Up to 4 TB/week

* Data storage and backup

* Network speed

* Security (human data)

# Analysis hardware

| System | Specifications |
|---|---|
| **Pipeline server** | |
| Processor | HP Proliant dl580 g5 rack server (4 quad-core 2.93GHz 64-bit Intel Xeon) |
| Memory | 32 GB |
| Storage | 21 TB (HP 60 MSA) |
| Operating system | Linux |
| **iPAR** | |
| Processor | HP DL 380 (2 × 5460 3.16 GHz) |
| Memory | 16 GB |
| Storage | 3.2 TB  (HP SmartArray P800) |
| Operating system | Linux/XP |
| **Mac Pro (x2)** | |
| Processor | 2 quad-core 2.66 GHz 64-bit Intel Nehalem |
| Memory | 16 GB |
| Storage | 4 TB |
| Operating system | OS X |

NorStore,Titan.....

# Analysis pipeline

Illumina Pipeline 1.4

| SCS | Firecrest | Bustard | GERALD |
|-----|-----------|---------|--------|
| Images | Image Analysis | Base Calling | Aligned Reads |

FASTQ format

```
@EAS54_6_R1_2_1_413_324
CCCTTCTTGTCTTCAGCGTTTCTCC
+
;;3;;;;;;;;;;;;7;;;;;;;88
@EAS54_6_R1_2_1_540_792
TTGGCAGGCCAAGGCCGATGGATCA
+
;;;;;;;;;;;7;;;;;-;;;3;83
@EAS54_6_R1_2_1_443_348
GTTGCTTCTGGCGTGGGTGGGGGGG
+EAS54_6_R1_2_1_443_348
;;;;;;;;;;;;9;7;;.7;39333
```

Other software/analyses

**Integrated solutions**
* CLCbio Genomics Workbench - *de novo* and reference assembly of Sanger, Roche FLX, Illumina, Helicos, and SOLiD data. Commercial next-gen-seq software that extends the CLCbio Main Workbench software. Includes SNP detection, CHiP-seq, browser and other features. Commercial. Windows, Mac OS X and Linux.
* Galaxy - Galaxy = interactive and reproducible genomics. A job webportal.
* Genomatix - Integrated Solutions for Next Generation Sequencing data analysis.
* JMP Genomics - Next gen visualization and statistics tool from SAS. They are working with NCGR to refine this tool and produce others.
* NextGENe - *de novo* and reference assembly of Illumina, SOLiD and Roche FLX data. Uses a novel Condensation Assembly Tool approach where reads are joined via "anchors" into mini-contigs before assembly. Includes SNP detection, CHiP-seq, browser and other features. Commercial. Win or MacOS.
* SeqMan Genome Analyser - Software for Next Generation sequence assembly of Illumina, Roche FLX and Sanger data integrating with Lasergene Sequence Analysis software for additional analysis and visualization capabilities. Can use a hybrid templated/de novo approach. Commercial. Win or Mac OS X.
* SHORE - SHORE, for Short Read, is a mapping and analysis pipeline for short DNA sequences produced on a Illumina Genome Analyzer. A suite created by the 1001 Genomes project. Source for POSIX.
* SlimSearch - Fledgling commercial product.

**Align/Assemble to a reference**
* BFAST - Blat-like Fast Accurate Search Tool. Written by Nils Homer, Stanley F. Nelson and Barry Merriman at UCLA.
* Bowtie - Ultrafast, memory-efficient short read aligner. It aligns short DNA sequences (reads) to the human genome at a rate of 25 million reads per hour on a typical workstation with 2 gigabytes of memory. Uses a Burrows-Wheeler-Transformed (BWT) index. Link to discussion thread here. Written by Ben Langmead and Cole Trapnell. Linux, Windows, and Mac OS X.
* BWA - Heng Lee's BWT Alignment program - a progression from Maq. BWA is a fast light-weighted tool that aligns short sequences to a sequence database, such as the human reference genome. By default, BWA finds an alignment within edit distance 2 to the query sequence. C++ source.
* ELAND - Efficient Large-Scale Alignment of Nucleotide Databases. Whole genome alignments to a reference genome. Written by Illumina author Anthony J. Cox for the Solexa 1G machine.
* Exonerate - Various forms of pairwise alignment (including Smith-Waterman-Gotoh) of DNA/protein against a reference. Authors are Guy St C Slater and Ewan Birney from EMBL. C for POSIX.
* GenomeMapper - GenomeMapper is a short read mapping tool designed for accurate read alignments. It quickly aligns millions of reads either with ungapped or gapped alignments. A tool created by the 1001 Genomes project. Source for POSIX.
* GMAP - GMAP (Genomic Mapping and Alignment Program) for mRNA and EST Sequences. Developed by Thomas Wu and Colin Watanabe at Genentec. C/Perl for Unix.
* gnumap - The Genomic Next-generation Universal MAPper (gnumap) is a program designed to accurately map sequence data obtained from next-generation sequencing machines (specifically that of Solexa/Illumina) back to a genome of any size. It seeks to align reads from nonunique repeats using statistics. From authors at Brigham Young University. C source/Unix.
* MAQ - Mapping and Assembly with Qualities (renamed from MAPASS2). Particularly designed for Illumina with preliminary functions to handle ABI SOLiD data. Written by Heng Li from the Sanger Centre. Features extensive supporting tools for DIP/SNP detection, etc. C++ source
* MOSAIK - MOSAIK produces gapped alignments using the Smith-Waterman algorithm. Features a number of support tools. Support for Roche FLX, Illumina, SOLiD, and Helicos. Written by Michael Strömberg at Boston College. Win/Linux/MacOSX
* MrFAST and MrsFAST - mrFAST & mrsFAST are designed to map short reads generated with the Illumina platform to reference genome assemblies; in a fast and memory-efficient manner. Robust to INDELs and MrsFAST has a bisulphite mode. Authors are from the University of Washington. C as source.
* MUMmer - MUMmer is a modular system for the rapid whole genome alignment of finished or draft sequence. Released as a package providing an efficient suffix tree library, seed-and-extend alignment, SNP detection, repeat detection, and visualization tools. Version 3.0 was developed by Stefan Kurtz, Adam Phillippy, Arthur L Delcher, Michael Smoot, Martin Shumway, Corina Antonescu and Steven L Salzberg - most of whom are at The Institute for Genomic Research in Maryland, USA. POSIX OS required.
* Novocraft - Tools for reference alignment of paired-end and single-end Illumina reads. Uses a Needleman-Wunsch algorithm. Can support Bis-Seq. Commercial. Available free for evaluation, educational use and for use on open not-for-profit projects. Requires Linux or MacOS.
* PASS - It supports Illumina, SOLiD and Roche-FLX data formats and allows the user to modulate very finely the sensitivity of the alignments. Spaced seed intial filter, then NW dynamic algorithm to a SW(like) local alignment. Authors are from CRIBI in Italy. Win/Linux.
* RMAP - Assembles 20 - 64 bp Illumina reads to a FASTA reference genome. By Andrew D. Smith and Zhenyu Xuan at CSHL. (published in BMC Bioinformatics). POSIX OS required.
* SeqMap - Supports up to 5 or more bp mismatches/INDELs. Highly tunable. Written by Hui Jiang from the Wong lab at Stanford. Builds available for most OS's.
* SHRiMP - Assembles to a reference sequence. Developed with Applied Biosystem's colourspace genomic representation in mind. Authors are Michael Brudno and Stephen Rumble at the University of Toronto. POSIX.
* Slider- An application for the Illumina Sequence Analyzer output that uses the probability files instead of the sequence files as an input for alignment to a reference sequence or a set of reference sequences. Authors are from BCGSC. Paper is here.
* SOAP - SOAP (Short Oligonucleotide Alignment Program). A program for efficient gapped and ungapped alignment of short oligonucleotides onto reference sequences. The updated version uses a BWT. Can call SNPs and INDELs. Author is Ruiqiang Li at the Beijing Genomics Institute. C++, POSIX.
* SSAHA - SSAHA (Sequence Search and Alignment by Hashing Algorithm) is a tool for rapidly finding near exact matches in DNA or protein databases using a hash table. Developed at the Sanger Centre by Zemin Ning, Anthony Cox and James Mullikin. C++ for Linux/Alpha.
* SOCS - Aligns SOLiD data. SOCS is built on an iterative variation of the Rabin-Karp string search algorithm, which uses hashing to reduce the set of possible matches, drastically increasing search speed. Authors are Ondov B, Varadarajan A, Passalacqua KD and Bergman NH.
* SWIFT - The SWIFT suit is a software collection for fast index-based sequence comparison. It contains: SWIFT — fast local alignment search, guaranteeing to find epsilon-matches between two sequences. SWIFT BALSAM — a very fast program to find semiglobal non-gapped alignments based on k-mer seeds. Authors are Kim Rasmussen (SWIFT) and Wolfgang Gerlach (SWIFT BALSAM)
* SXOligoSearch - SXOligoSearch is a commercial platform offered by the Malaysian based Synamatix. Will align Illumina reads against a range of Refseq RNA or NCBI genome builds for a number of organisms. Web Portal. OS independent.
* Vmatch - A versatile software tool for efficiently solving large scale sequence matching tasks. Vmatch subsumes the software tool REPuter, but is much more general, with a very flexible user interface, and improved space and time requirements. Essentially a large string matching toolbox. POSIX.
* Zoom - ZOOM (Zillions Of Oligos Mapped) is designed to map millions of short reads, emerged by next-generation sequencing technology, back to the reference genomes, and carry out post-analysis. ZOOM is developed to be highly accurate, flexible, and user-friendly with speed being a critical priority. Commercial. Supports Illumina and SOLiD data.

*De novo* **Align/Assemble**
* ABySS - Assembly By Short Sequences. ABySS is a de novo sequence assembler that is designed for very short reads. The single-processor version is useful for assembling genomes up to 40-50 Mbases in size. The parallel version is implemented using MPI and is capable of assembling larger genomes. By Simpson JT and others at the Canada's Michael Smith Genome Sciences Centre. C++ as source.
* ALLPATHS - ALLPATHS: De novo assembly of whole-genome shotgun microreads. ALLPATHS is a whole genome shotgun assembler that can generate high quality assemblies from short reads. Assemblies are presented in a graph form that retains ambiguities, such as those arising from polymorphism, thereby providing information that has been absent from previous genome assemblies. Broad Institute.
* Edena - Edena (Exact DE Novo Assembler) is an assembler dedicated to process the millions of very short reads produced by the Illumina Genome Analyzer. Edena is based on the traditional overlap layout paradigm. By D. Hernandez, P. François, L. Farinelli, M. Osteras, and J. Schrenzel. Linux/Win.
* EULER-SR - Short read *de novo* assembly. By Mark J. Chaisson and Pavel A. Pevzner from UCSD (published in Genome Research). Uses a de Bruijn graph approach.
* MIRA2 - MIRA (Mimicking Intelligent Read Assembly) is able to perform true hybrid de-novo assemblies using reads gathered through 454 sequencing technology (GS20 or GS FLX). Compatible with 454, Solexa and Sanger data. Linux OS required.
* SEQAN - A Consistency-based Consensus Algorithm for De Novo and Reference-guided Sequence Assembly of Short Reads. By Tobias Rausch and others. C++. Linux/Win.
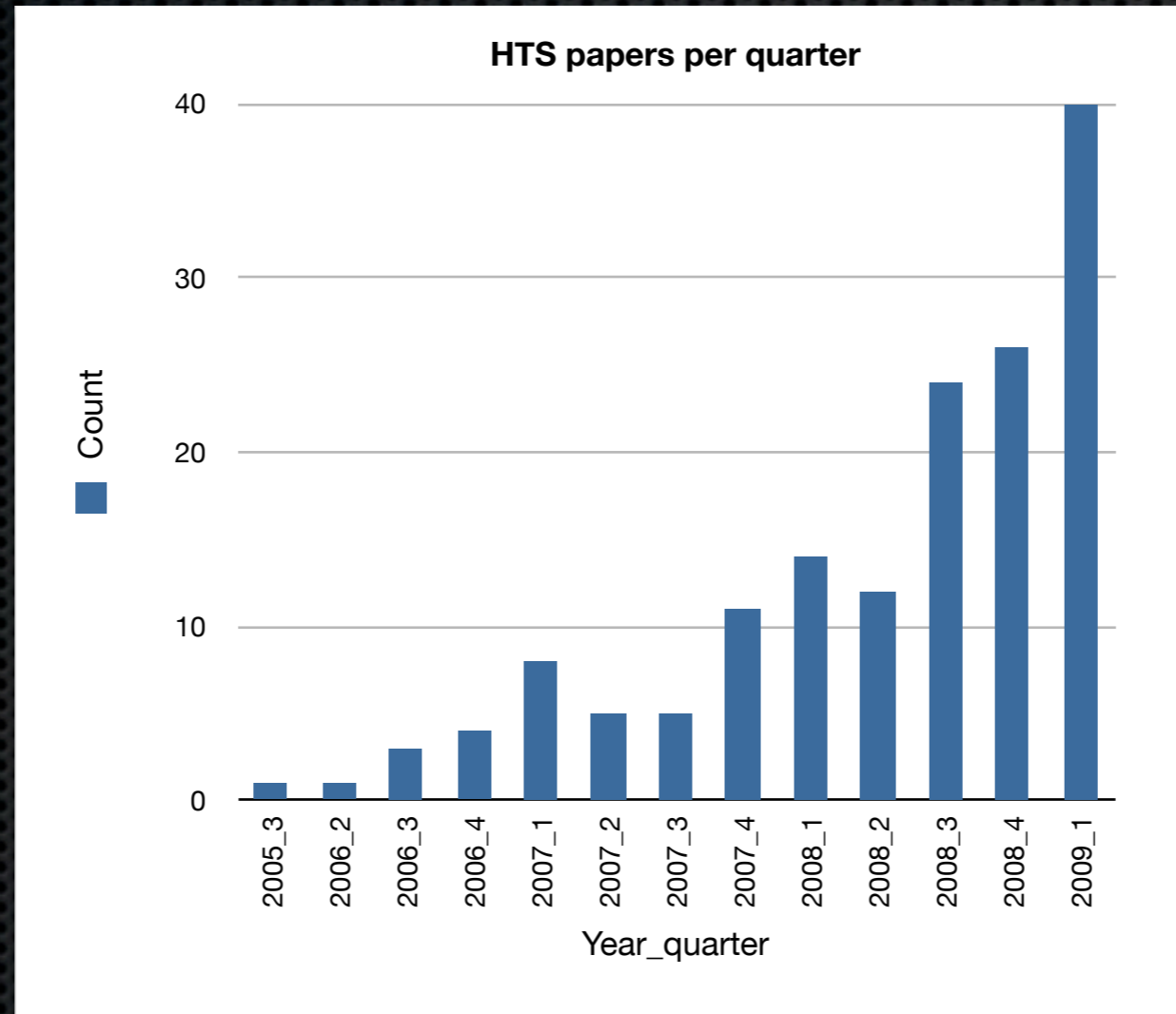
Majority Unix

# CLC genomics workbench

# Applications

# Research publications

| Year_quarter | Count |
|---|---|
| 2006_3 | 3 |
| 2006_4 | 4 |
| 2007_1 | 8 |
| 2007_2 | 5 |
| 2007_3 | 5 |
| 2007_4 | 11 |
| 2008_1 | 14 |
| 2008_2 | 12 |
| 2008_3 | 24 |
| 2008_4 | 26 |
| 2009_1 | 40 |

## HTS papers per quarter

# Applications

| Application | Project |
| --- | --- |
| Resequencing | whole genome<br>linkage/association<br>mutation detection |
| *de novo* sequencing | metagenomics<br>new species |
| Expression | transcriptome<br>SAGE<br>miRNA |
| Epigenetics | DNA methylation<br>ChIP |
| Variation | SNPs<br>CNVs |

# Users I

| User | Institute | Experiment | Species |
|---|---|---|---|
| Gregor/Kristina | IMG | Reseq/Epigen | Human |
| Beate Skinningsrud | IMG | RNAseq | Human |
| Eystein Husebye | UiB | RNAseq | Human |
| Randi Aamodt | CIGENE | DGE | Bee |
| Susanne Lorenz | CIGENE | RNAseq | Salmon |
| Elin Kure | OUSU | miRNA | Human |
| Matthew Kent/Sigbjørn Lien | CIGENE | Genome seq | Cattle/Cod/Pig |
| Arne Klungland | UiO | miRNA | Mouse |
| Gaute Brede | UiTrondheim | miRNA | Human |
| Gregor/SvenOlaf | IMG | reseq | Human |
| Ingar Olsen | OUS Riks | metagenomics | Bacteria |
| Hedda Hovik | UiO | Bact RNA seq | Bacteria |
| Gregor/Robert | IMG/Ullevål | ChIP/Bisulphite | Human |
| Kristina/Robert | IMG/Ullevål | Bisulphite | Human |

# Users II

- Many users

- Many institutes

- Many applications

→ Bioinformatic challenge

# Three research areas...

- 1000 genomes

- Resequencing - finding variants (SNPs, CNVs)

- Epigenetics

# 1000 genomes project

- International consortium

- Sequence 1200 genomes

- Produce a nearly complete catalog of common human genetic variants (defined as frequency 1% or higher; SNPs, CNVs)

  - mutation detection in Mendelian disease

  - accelerate fine-mapping efforts association studies

  - enabling design of next-generation genotyping arrays - improve the power of future genetic association studies

  - improve our ability to "impute" or "predict" untyped genetic variants

- Frequent public data releases

1000genomes.org

# Resequencing

- Compare test sequence to a reference sequence

- Identify genetic variation

  - Single-nucleotide polymorphisms (SNPs)

  - Insertions/deletions

  - Copy-number variation (CNVs)

# Resequencing: mutation detection
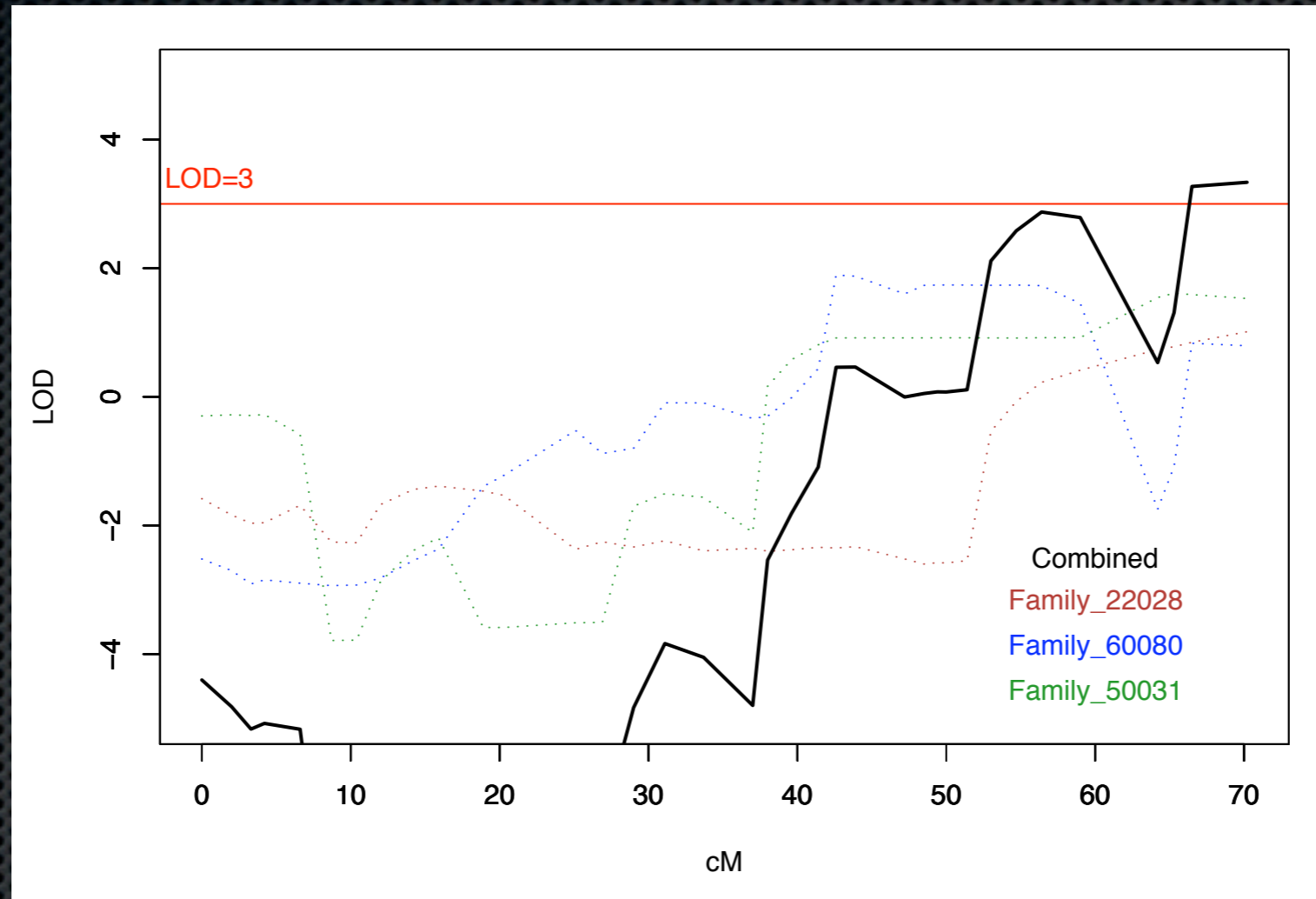
## Genomic region known

- Linkage peak

- Sequence capture - region of interest

## Genomic region unknown

- Rare Mendelian disorders

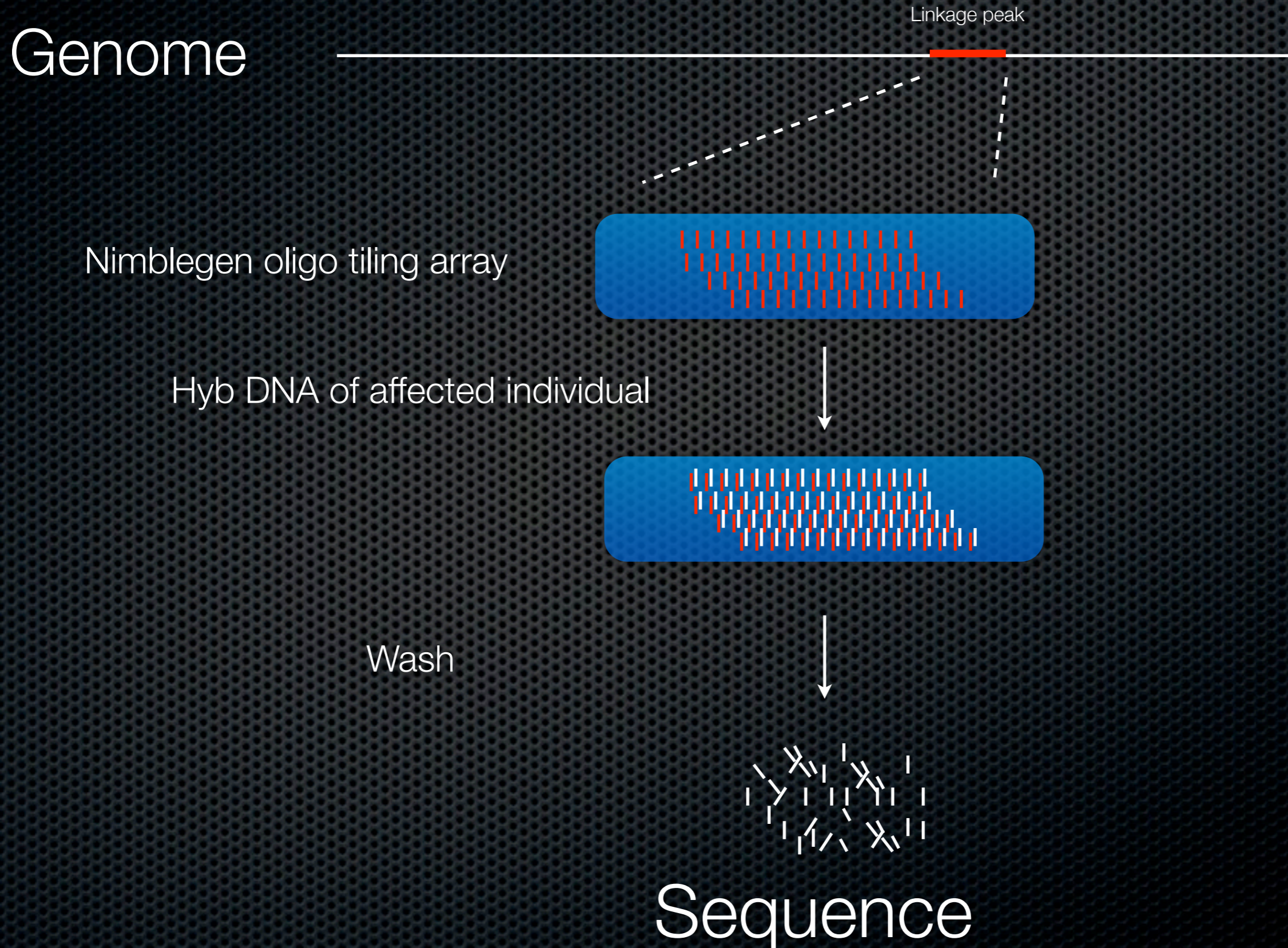- Sequence capture - exome

- RNAseq

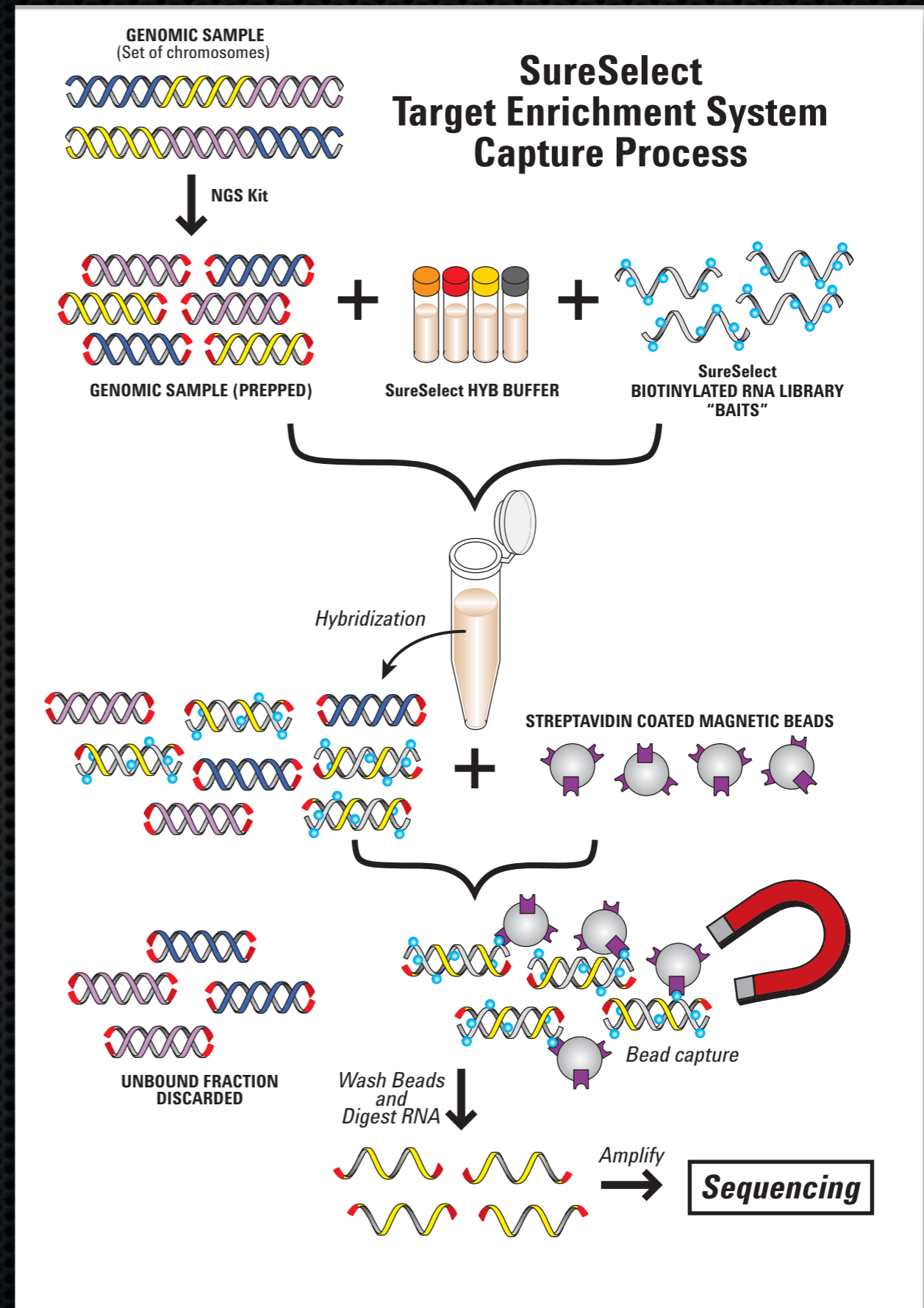# Region known

## Linkage



1-10 Mb?

How can we capture this region to sequence?

# Sequence capture

Genome

Linkage peak

Nimblegen oligo tiling array

Hyb DNA of affected individual

Wash

Sequence

# Agilent SureSelect

- RNA oligonucleotides

- >100 bp
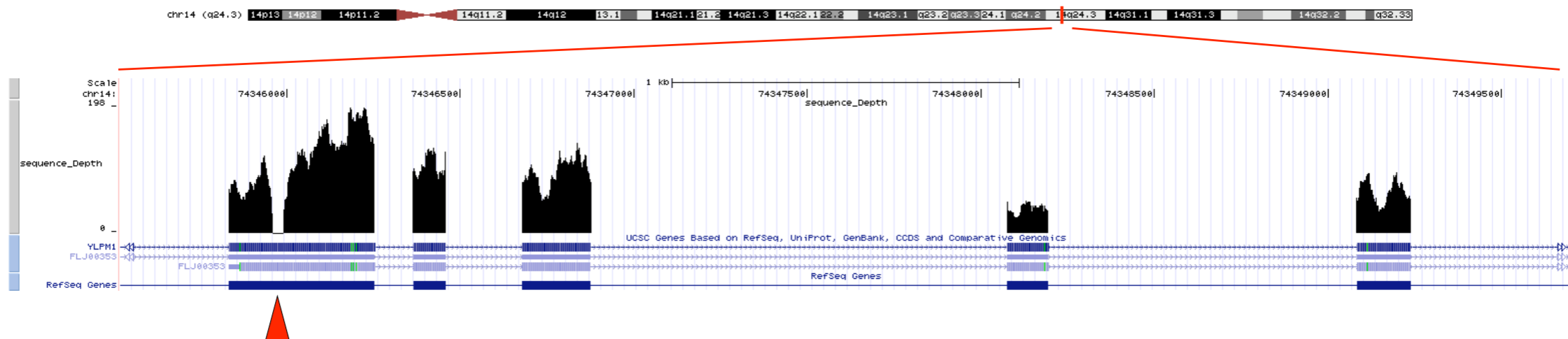
- custom design
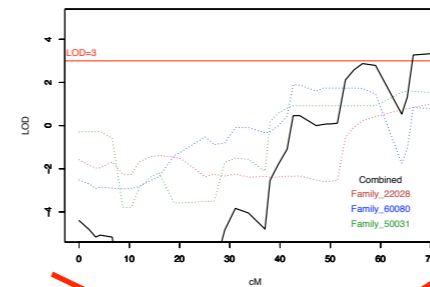
www.agilent.com



SureSelect
Target Enrichment System
Capture Process

# Analyzing resequencing data

- Capture DNA and sequence

- Prepare sequence files (Perl...)

- Align to reference (MAQ etc.)

- Format/filter output files (Perl...)

  - .bed, .gtf

  - View on genome browser

  - identify variants
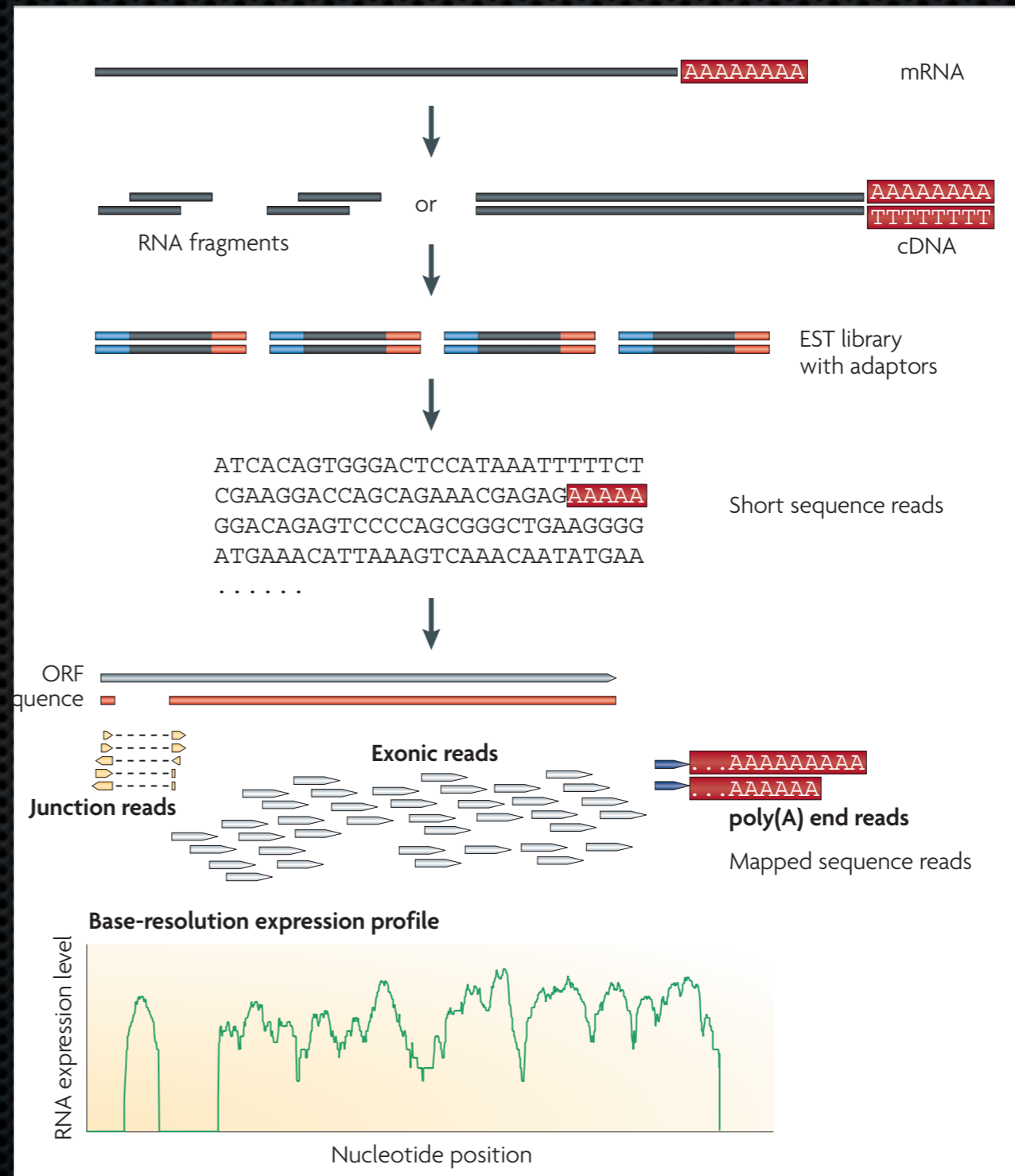
# Analyzing resequencing data



Hsa chr14

13 bp insertion

# Identifying relevant variants is the hard part

# Region unknown

- Sequence capture - exome: sequence all exons

- RNAseq

- Sequence total polyA RNA

- Map reads to reference
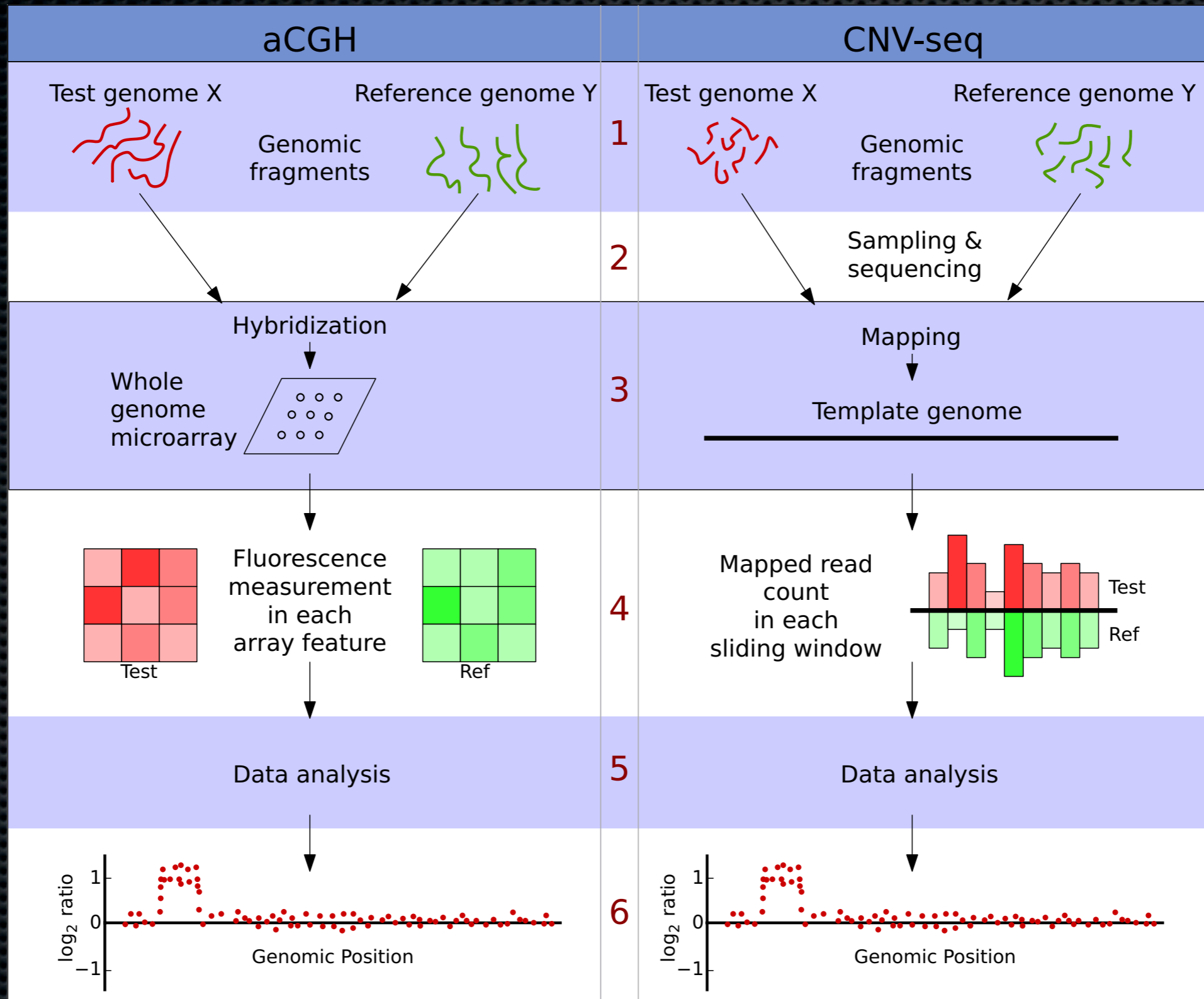
- Identify mutations/variants

# RNAseq data



**Sequence depth**

Gene - APP

Position along Hsa chr21

# Identifying relevant variants is the hard part

# HTS and CNVs
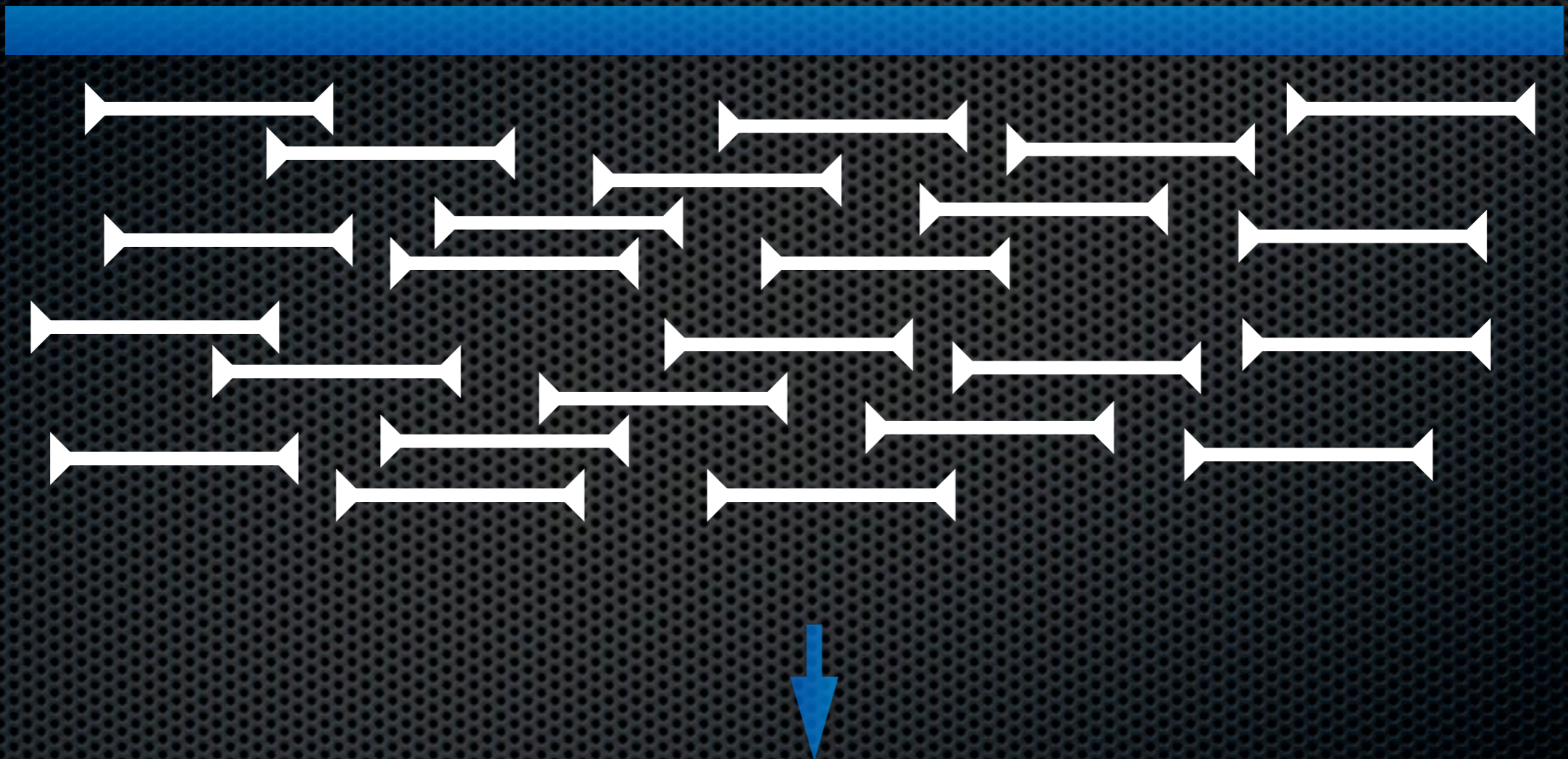
Two strategies to detect CNVs with HTS data

- Read map counting

- Mapping paired-end reads

  - Read map location

  - Read map distance

  - Read map orientation

# Read map counting

# Map paired-end reads

## Reference genome



Consensus sequence
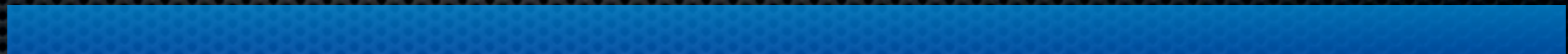identify variants, mutations

# Deletion

Test genome

~300 bp

Map to reference

Reference genome

>>> ~300 bp
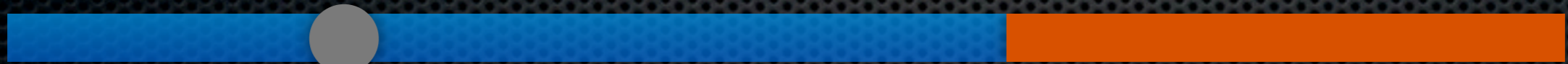
# Duplication

Test genome

~300 bp
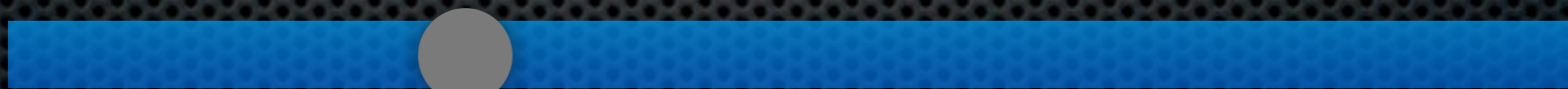
Map to reference

Reference genome

# Balanced translocation

Translocation chromosome
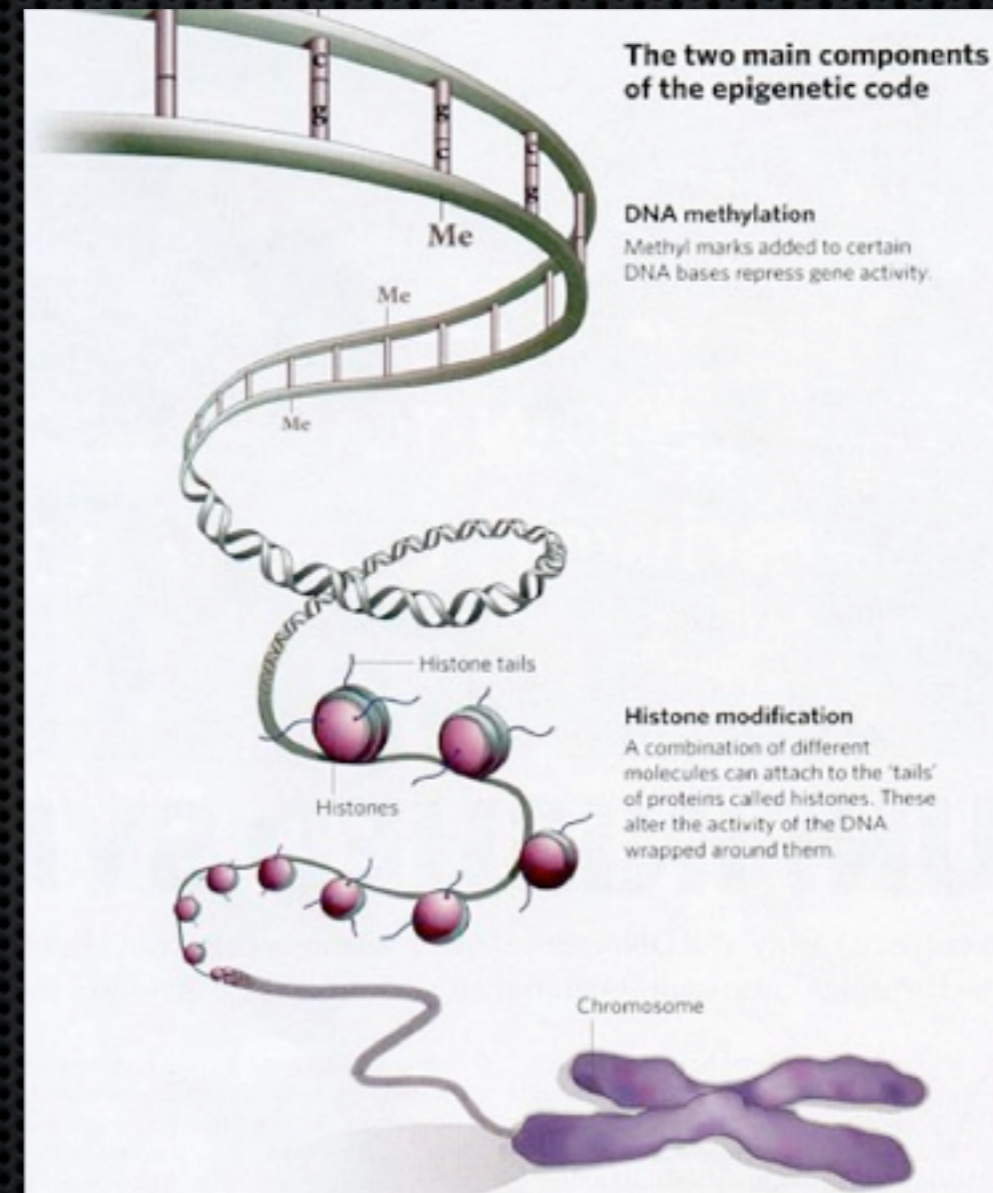
Map to reference

Reference genome

# Detecting all variants

| VARIANT | SINGLE READ | SHORT INSERT PAIRED-ENDS (200–500 bp) | LONG INSERT MATE PAIRS (2–5 kb) | PAIRED-END AND MATE PAIR COMBINED |
|---|---|---|---|---|
| SNP | ++ | ++++ | ++ | ++++ |
| Small indels | ++ | ++++ | ++ | ++++ |
| Insertion | + | +++ | +++ | ++++ |
| Amplification | ++ | +++ | +++ | ++++ |
| Deletion | + | +++ | ++ | ++++ |
| Inversion | + | +++ | ++ | ++++ |
| Complex rearrangement | + | +++ | ++ | ++++ |
| Large rearrangement | + | ++ | +++ | ++++ |

# Epigenetics

* **DNA methylation**

  * CpG dinucleotides

* **Histone modifications**

  * acetylation

  * phosphorylation

  * methylation

  * ubiquitination



**➡ Control of gene expression**

# Epigenetics II
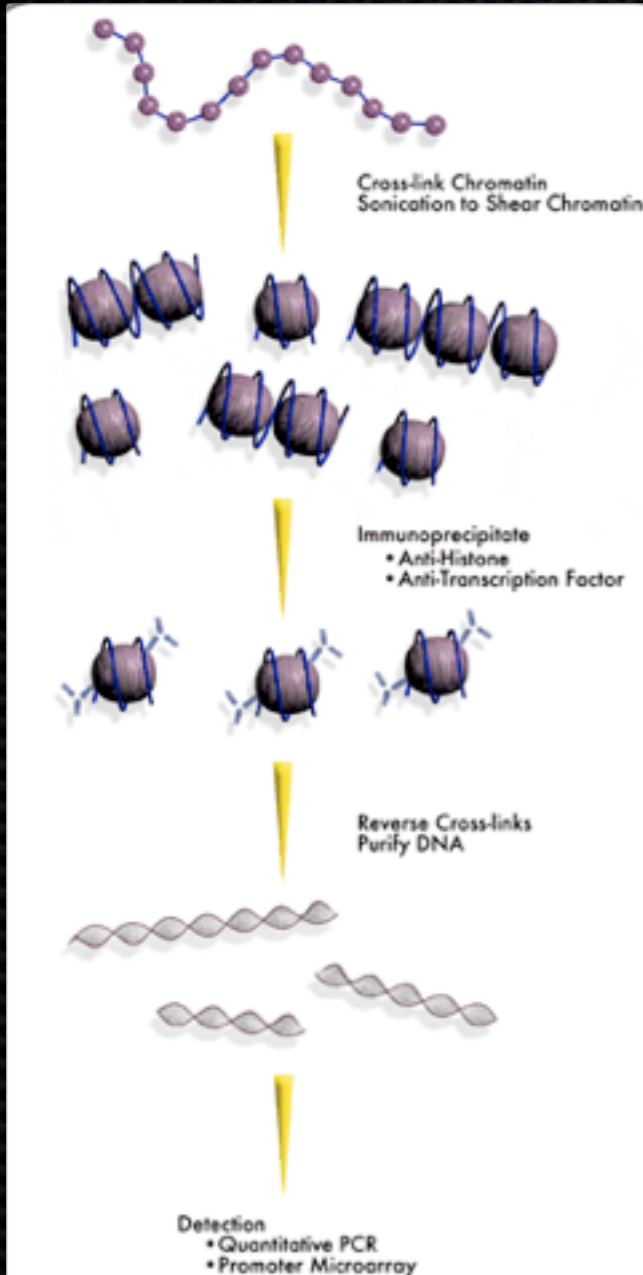
- DNA methylation

  - Long-term epigenetic silencing of specific sequences

  - transposons, imprinted genes, pluripotency genes

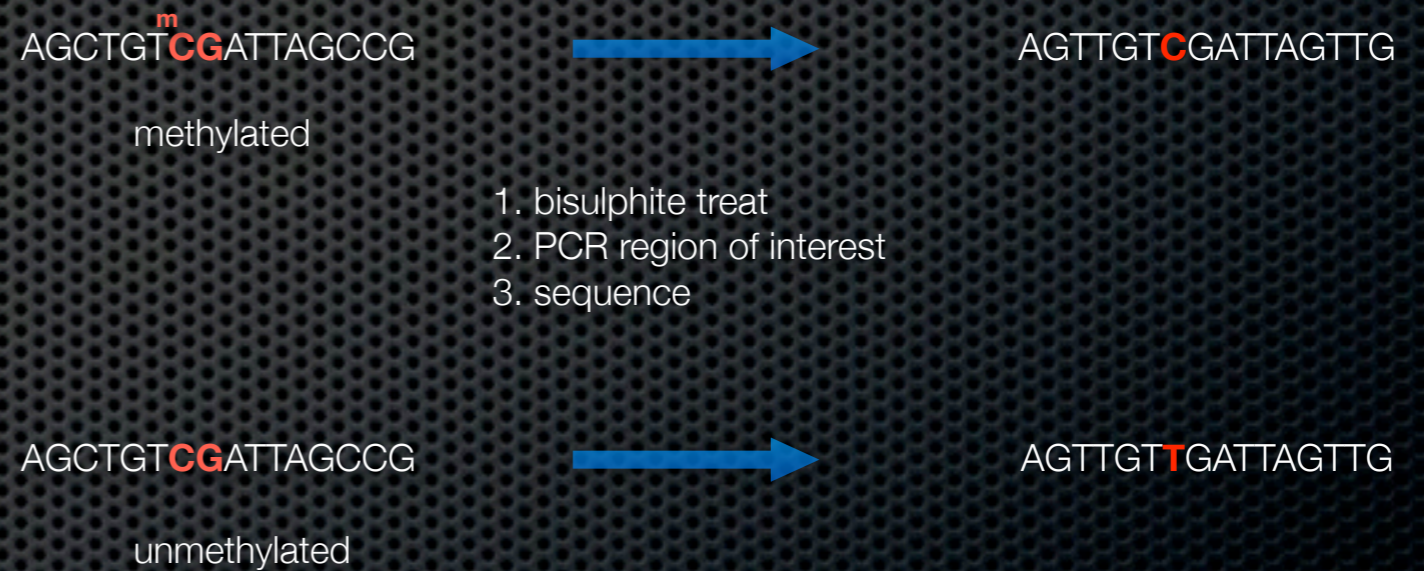- Histone modifications

  - Short term, flexible epigenetic control

→ Control of gene expression

# HTS and epigenetics
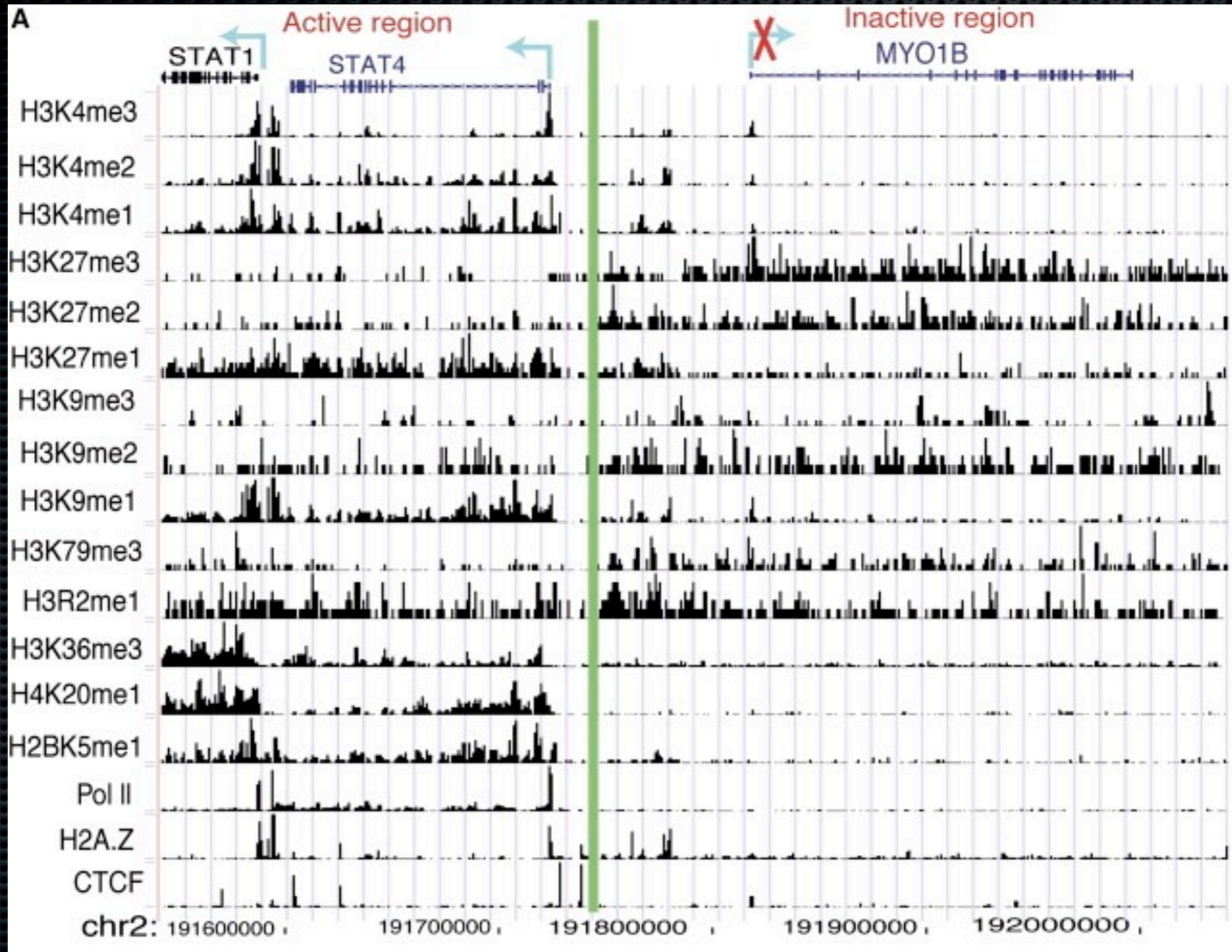
ChIP
chromatin immunoprecipitation

Quantifying DNA methylation

Bisulphite sequencing (BiS)



AGCTGT**CG**ATTAGCCG  →  AGTTGT**C**GATTAGTTG

methylated

1. bisulphite treat
2. PCR region of interest
3. sequence

AGCTGT**CG**ATTAGCCG  →  AGTTGT**T**GATTAGTTG

unmethylated

HTS to identify genome-wide status/variation

# ChIP-seq example

# Summary

- High-throughput sequencing

  - Dramatic increase in sequence production

  - Many applications on one platform

  - Field new and moving very quickly

- Bioinformatics challenges/opportunities

  - Data storage

  - Data analysis

# Visit?

Robert.Lyle@medisin.uio.no