# High-throughput sequencing

Robert Lyle

Department of Medical Genetics

Oslo University Hospital

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

$4\times10^{5}$ - $1\times10^{9}$

...on one machine

**Massively parallel**

# HTS systems available



454 — Roche
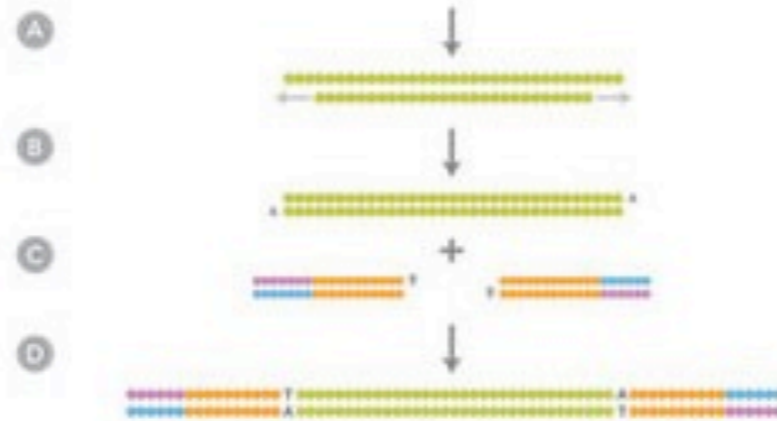
Solexa — Illumina

SOLiD — ABI

HeliScope — Helicos

Others in 2011
(Pacific BioSciences, Ion Torrent)

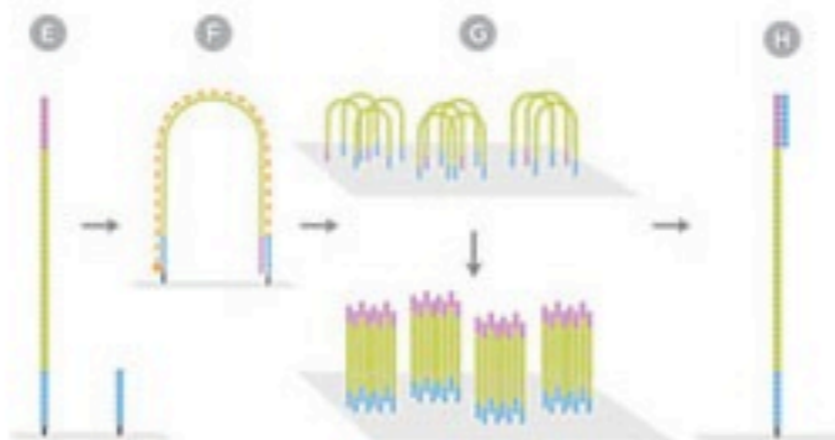# Illumina sequencing technology

**1. Library preparation**

6 hours
3 hours hands-on time

A Fragment DNA
B Repair ends
  Add A overhang
C Ligate adapters
D Select ligated DNA

**2. Cluster generation**

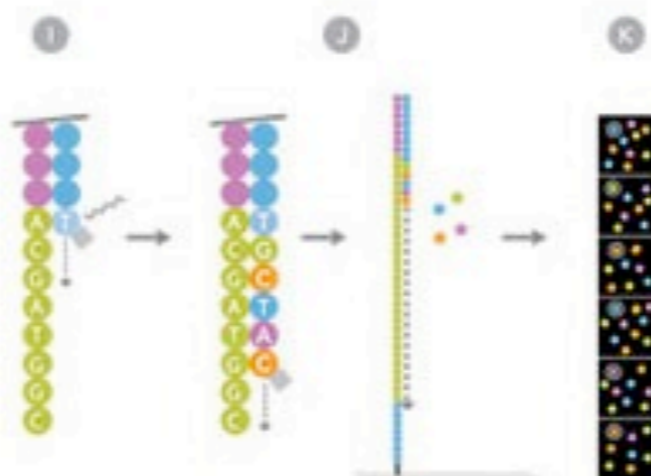4 hours
30 minutes hands-on time
1–96 samples

E Attach DNA to flow cell
F Perform bridge amplification
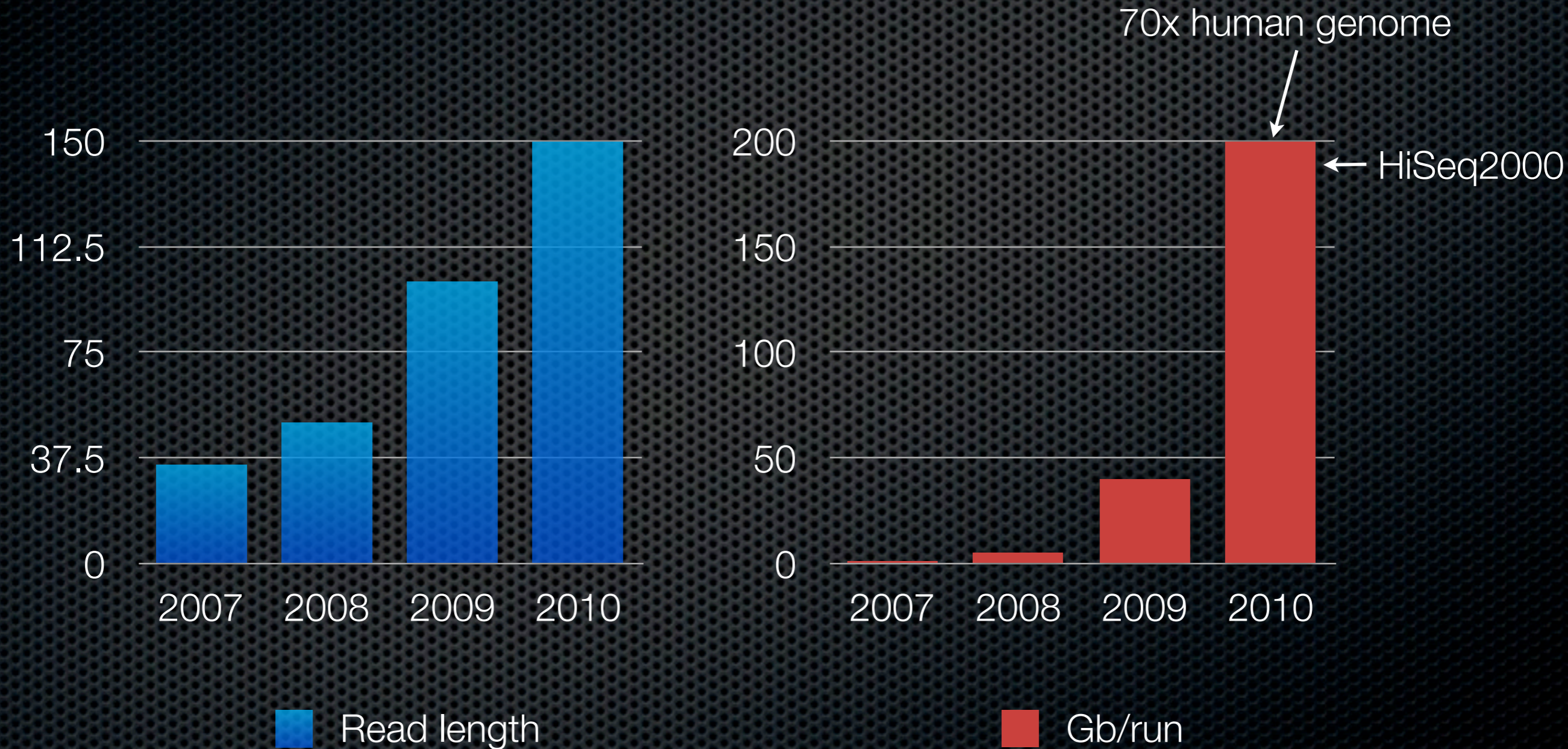G Generate clusters
H Anneal sequencing primer

**3. Sequencing**

1–3 days single-read run
3–7 days paired-end run
30 minutes hands-on time
1–96 samples

I Extend first base, read, and deblock
J Repeat step above to extend strand
K Generate base calls

Illumina throughput

# User contact

http://www.sequencing.uio.no



## Services

454

Illumina

Sample delivery form
FAQ

post@sequencing.uio.no

National conferences
NSB: talks, booth
NSHG, seminars etc.

Contact

# Illumina platform

| Instruments | Illumina GAIIx (2) (HiSeq2000) |
|---|---|
| People | 1 Daily leader 1 PostDoc 2 Technicians 2 Bioinformaticians |
| Data storage | ~60 TB local NorStore Secure storage... |

# Platform services

| User | Sample | DNA, RNA |
|---|---|---|
| **Platform** | Sequencing | QC<br>library preparation<br>sequencing |
| | Costs | Illumina reagents<br>QC reagents<br>20% platform fee<br>(No staff/platform costs) |
| | Bioinformatics | Basic run information, QC<br>Alignment to reference genome<br>? |

# Applications run on Illumina node

| Application | Project | Sample | Protocol |
| --- | --- | --- | --- |
| Resequencing | whole genome linkage/association mutation detection | Genomic DNA | sequence capture, exome sequencing |
| *de novo* sequencing | metagenomics new species | Genomic DNA | SE, PE, mate-pair |
| Expression | transcriptome miRNA | mRNA, miRNA | RNAseq, miRNA |
| Epigenetics | DNA methylation chromatin structure | Genomic DNA | Bisulphite sequencing (RRBS), ChIP, MeDIP |

1x36 bp  ->  2x108 bp

# Runs overview

# Nationwide users

- 🔴 Run
- 🟡 Scheduled
- 🔵 Contacted

University of Tromsø, Tromsø

Institute of Food, Fisheries and Aquaculture Research (NOFIMA), Tromsø

Oslo University Hospital

Norwegian University of Science and Technology (NTNU), Trondheim

University of Oslo

National Institute of Nutrition and Seafood Research (NIFES), Bergen

Norwegian School of Veterinary Science

University of Bergen, Bergen

Institute for Forestry and Landscape

Telemark Hospital, Skien

CIGENE, Norwegian University of Life Sciences

# Data and analysis

# Illumina sequence data

* Random DNA library of short fragments   ~300 bp

* ~100-300 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)


* Data law - 'return of results'

* Bioinformatics

# Users

- Many users

- Many institutes

- Many applications

→ Bioinformatic challenge

# User data storage - Phases

| Phase | Provision | Timeline |
|-------|-----------|----------|
| 0 | Storage/backup of non-sensitive data from NSC (NorStore) | Complete 12.2009 |
| 1 | High capacity secure storage of coded but indirectly identifiable data at local level (OUS, UiO)<br>Establish routines for backup | Complete 6.2010 |
| 2 | Robust secure solution for HTS data at the national level<br>Exchange of data collaborators<br>Infrastructure for analysis through the Bioinformatics platform | Start 11.2010 |
| 3 | Clinical usage, secure handling of person-identifiable data being part of patient journals | ? |

# Illumina analysis pipeline

## Illumina Pipeline 1.4

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



```
@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
```

FASTQ format

Other software/
analyses

| Run size | Images | Intensity files | Sequence files | Analysis output | Total |
|---|---|---|---|---|---|
| **Minimum** | 700 GB | 100 GB | 100 GB | 100 GB | **1 TB** |
| **Average** | 1.4 TB | 200 GB | 200 GB | 200 GB | **2 TB** |
| **Maximum** | 3.5 TB | 500 GB | 500 GB | 500 GB | **5 TB** |

**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 Solution for Next Generation Sequencing data analysis.
* JMP Genomics - Next gen visualization and statistics tool from SAP. They are working with NCBR 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 Mac OS X.
* PASS - It supports Illumina, SOLiD and Roche-FLX data formats and allows the user to modulate very finely the sensitivity of the alignments. Spaced sequential 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. 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 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.

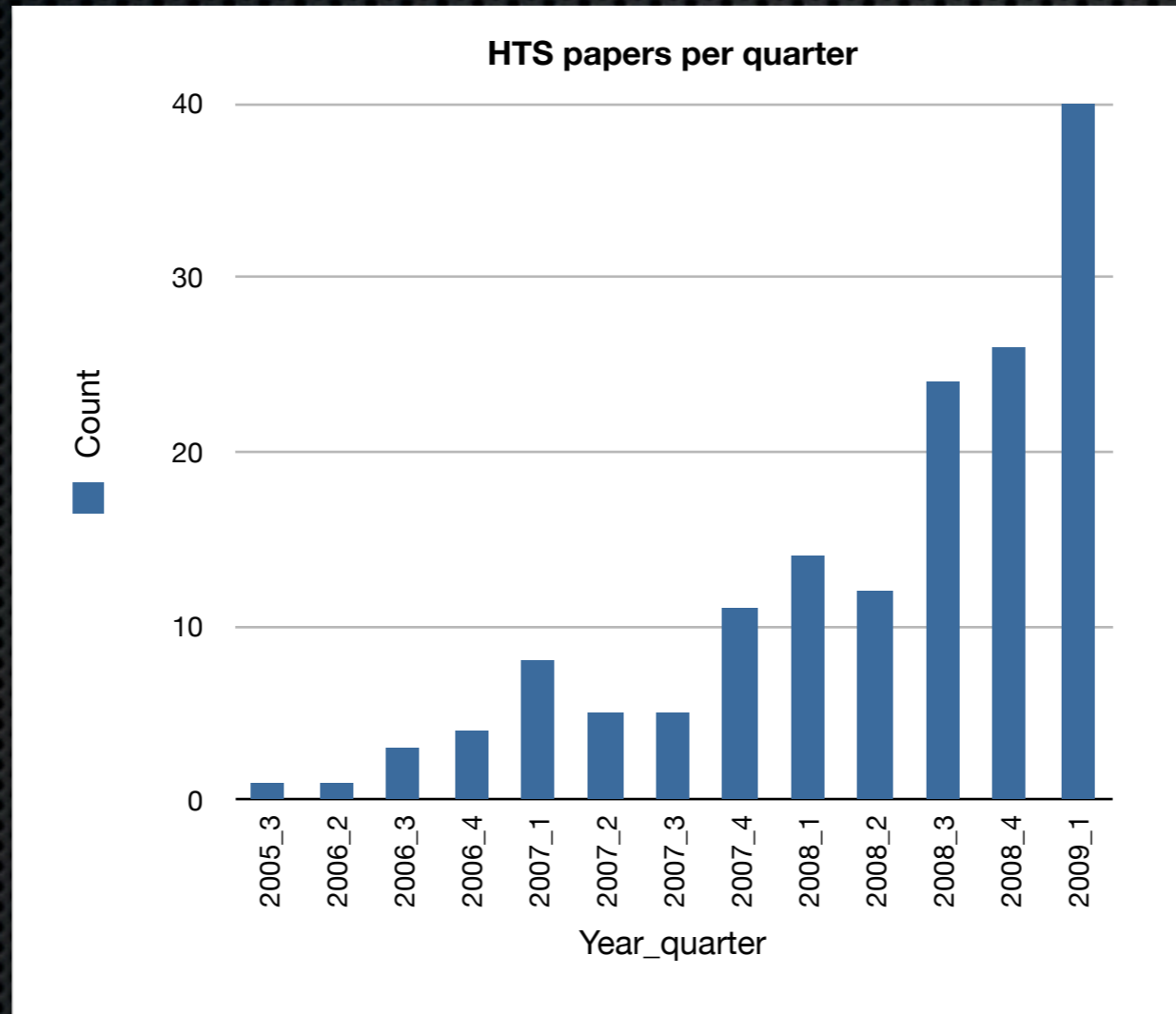# Analysis hardware

| 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.....

# Break?

# Applications

# Research publications

| | |
|---|---|
| 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 |

# Resequencing

* Compare test sequence to a reference sequence

    * Mendelian (linkage)

    * Association studies

    * Exome sequencing

* 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?

# Agilent SureSelect

- RNA oligonucleotides

- >100 bp

- custom design

www.agilent.com

# Analyzing resequencing data

* Capture DNA and sequence

* Prepare sequence files (Perl...)

* Align to reference (MAQ, BWA etc.)

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

  * .bed, .gtf

  * View on genome browser

  * identify variants

# 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

# Aim

```
@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
```

FASTQ format

Compare to
reference

R|G

Sequence

Mutation

# Tools for HTS

- Alignment - MAQ, BWA

- Filtering, sorting etc. - SAMTools, BEDTools

- Viewing - BED, GFF, UCSC browser

- Perl, unix scripts



xkcd.com

# Finding mutations?

- Which variants are deleterious?

- Novel? (dbSNP, 1000genomes, HGMD)

- Synonymous/non-synonymous?

- Conserved?

- Alter protein structure?



PolyPhen2
MutationTaster
ANNOVAR
SeattleSeq Annotation

# 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

# Crude analysis pipeline

| Program/script | Description | Output |
|---|---|---|
| maq fasta2bfa | Prepare ref sequences | bfa |
| maq fastq2bfq | Convert FASTQ reads to BFQ format | bfq |
| maq map | Align | MAQ aln |
| maq assemble | Assemble | MAQ cns |
| maq cns2snp | Call SNPs | MAQ SNP |
| awk '$2>=29621176 && $2<=39095041' | Filter for ROI | MAQ SNP |
| maq.pl SNPfilter | Q filter SNPs | MAQ SNP |
| maq cns2view | MAQ file for ROI | MAQ SNP |
| maqview2bed.pl | bed file for ROI | bed |
| maqsnp2bed.pl | bed file for SNPs | bed |
| maqsnp2snpnexus.pl | Input for SNPnexus | SNPnexus input |
| parseSNPnexus.pl | Parse SNPnexus output | SNPnexus output |
| bases2nexus.pl | Variants file | bases file |
| maqCoverageSummary.pl | Sequence coverage | bed, pdf |
| coverage_v4.pl | Sequence coverage | bed |
| lowCoverageSummary.pl | Sequence coverage | bed |

# Read depth statistics



R

# UCSC Genome browser



http://genome.ucsc.edu/

# Galaxy



http://main.g2.bx.psu.edu/

# Viewing data



SNP locations

Quality score

Read depth

Genes

# Analyzing resequencing data



Hsa chr14

13 bp insertion

# Variants file



Identifying relevant variants is the hard part

# CLC genomics workbench



Commercial software

# Detecting all variants

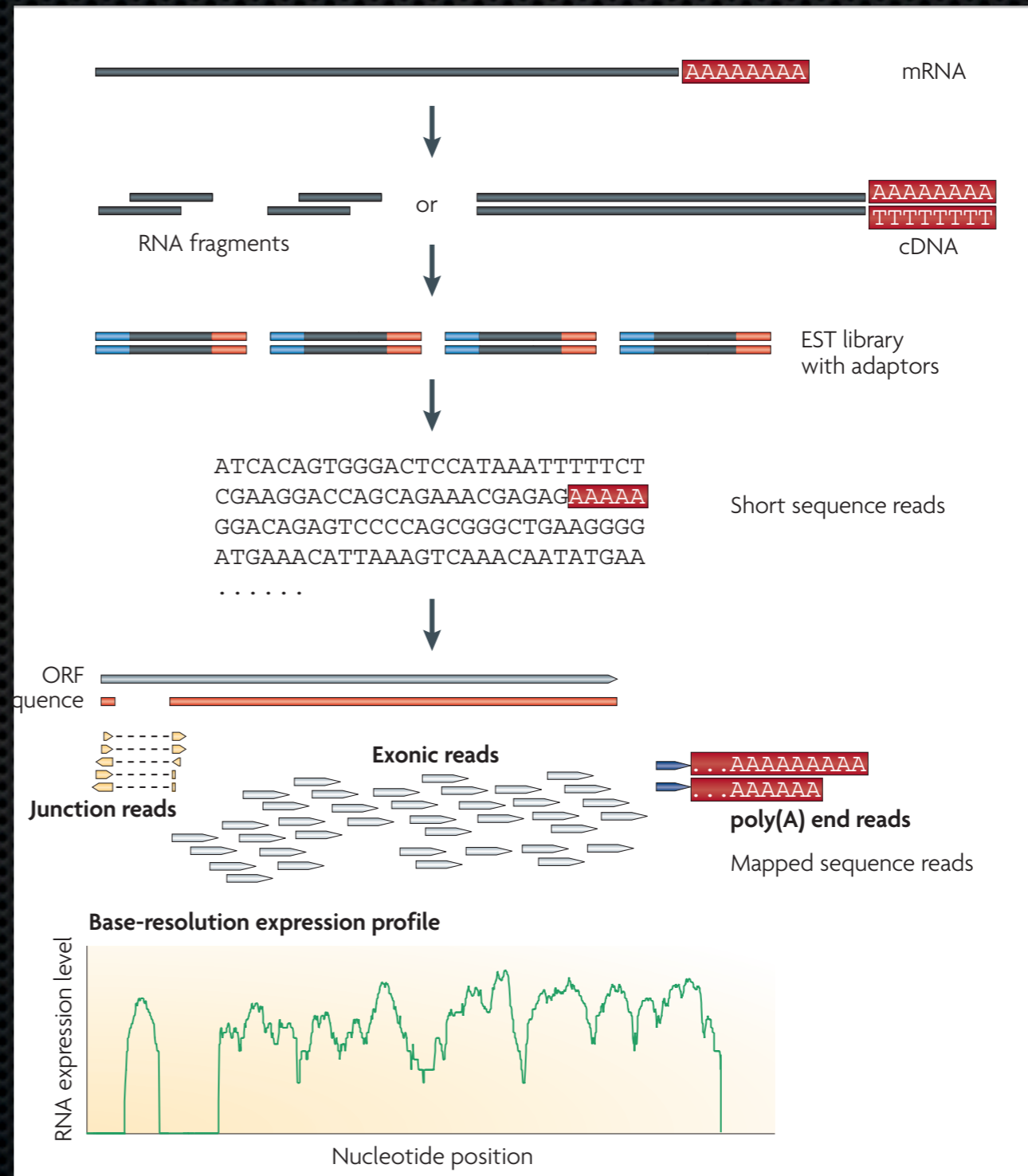| VARIANT | SINGLE READ | SHORT INSERT PAIRED-ENDS (200–500 bp) | LONG INSERT MATE PAIRS (2–5 kb) | PAIRED-END AND MATE PAIR COMBINED |
|---|---|---|---|---|
| SNP | ++ | ++++ | ++ | ++++ |
| Small indels | ++ | ++++ | ++ | ++++ |
| Insertion | + | +++ | +++ | ++++ |
| Amplification | ++ | +++ | +++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

# Region unknown

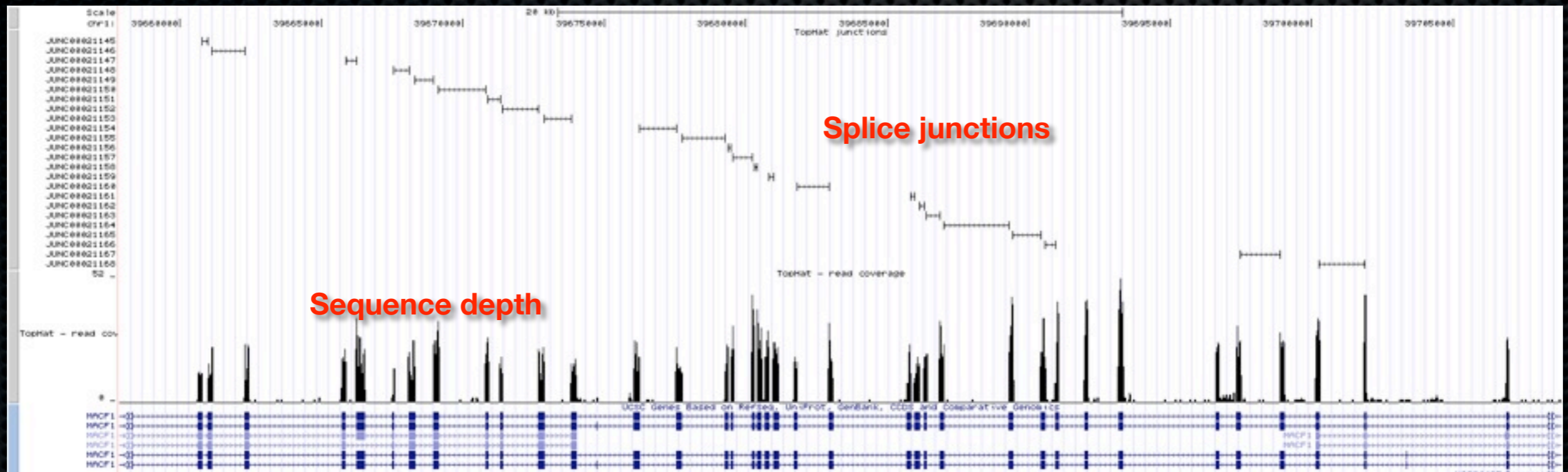- Sequence capture - exome: sequence all exons

- RNAseq

- Sequence total polyA RNA

- Map reads to reference

- Identify mutations/variants

# RNAseq data



Position along Hsa chr1

Identifying relevant variants is the hard part

# Epigenetics

- **DNA methylation**
  - CpG dinucleotides

- **Histone modifications**
  - acetylation
  - phosphorylation
  - methylation
  - ubiquitination



The two main components of the epigenetic code

**DNA methylation**
Methyl marks added to certain DNA bases repress gene activity.

Histone tails

Histones

**Histone modification**
A combination of different molecules can attach to the 'tails' of proteins called histones. These alter the activity of the DNA wrapped around them.

Chromosome

→ Control of gene expression

# Epigenetics II

* DNA methylation

  * Long-term epigenetic silencing of specific sequences

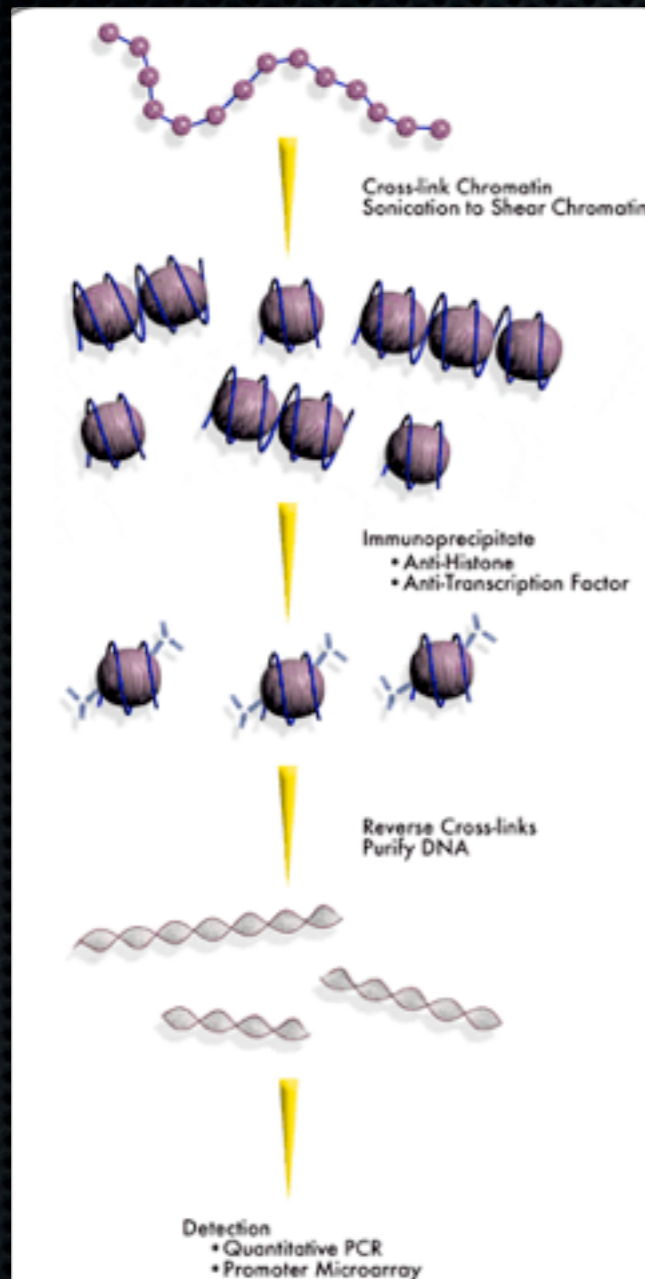  * transposons, imprinted genes, pluripotency genes

* Histone modifications

  * Short term, flexible epigenetic control

→ Control of gene expression

# HTS and epigenetics

ChIP
chromatin immunoprecipitation

Quantifying DNA methylation

Bisulphite sequencing (BiS)



AGCTGGT**CG**ATTAGCCG   →   AGTTGT**C**GATTAGTTG

methylated

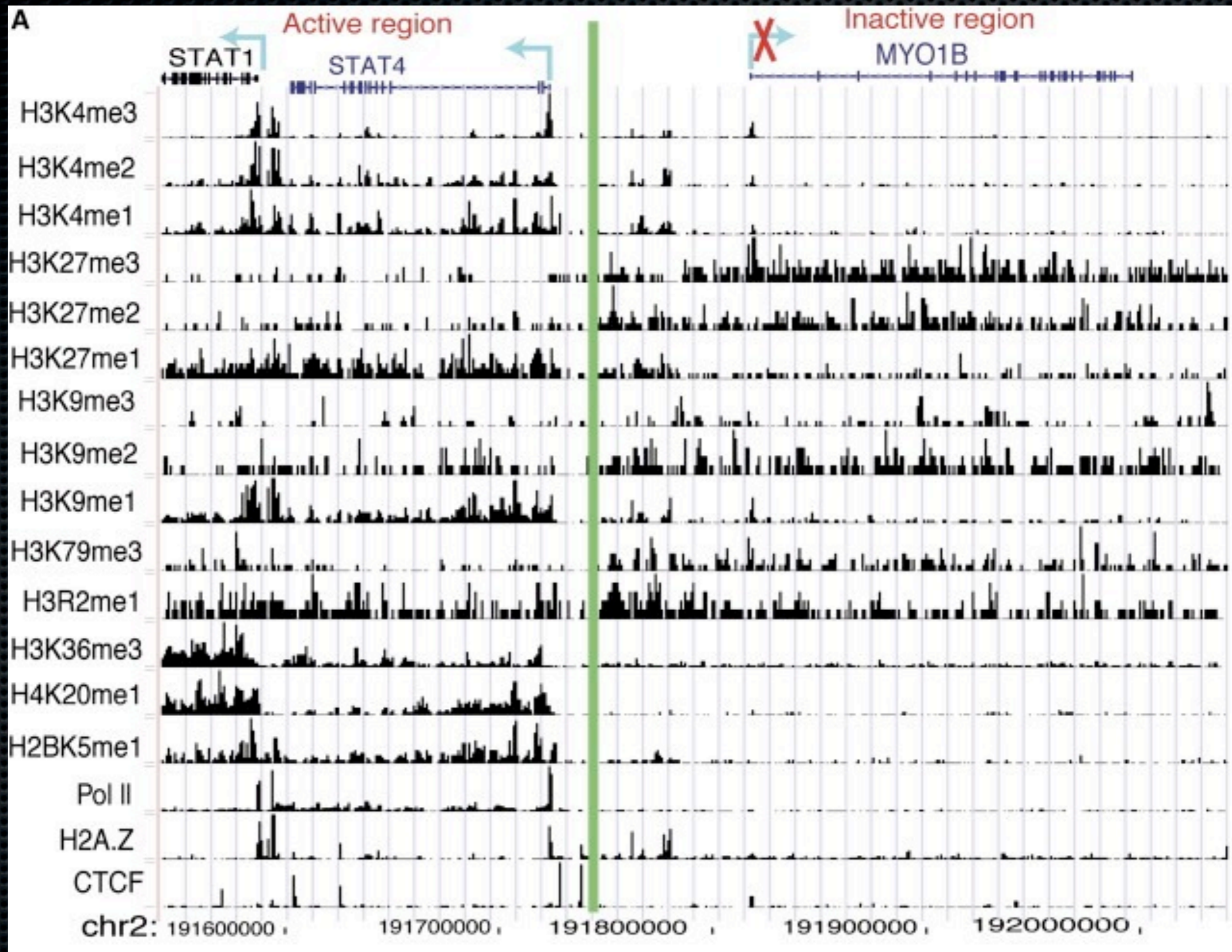1. bisulphite treat
2. PCR region of interest
3. sequence

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

unmethylated

HTS to identify genome-wide status/variation

# ChIP-seq example



Barski et al, Cell 129, 823–837, May 18, 2007

# 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