

High-Throughput Sequencing and its applications

IN-BIOS5000/9000

Genome Sequencing Technologies, Assembly, Variant Calling and Statistical Genomics
19 October 2020

Torbjørn Rognes

Dept. of Informatics, UiO & Dept. of Microbiology, OUS
torognes@ifi.uio.no



UiO University of Oslo



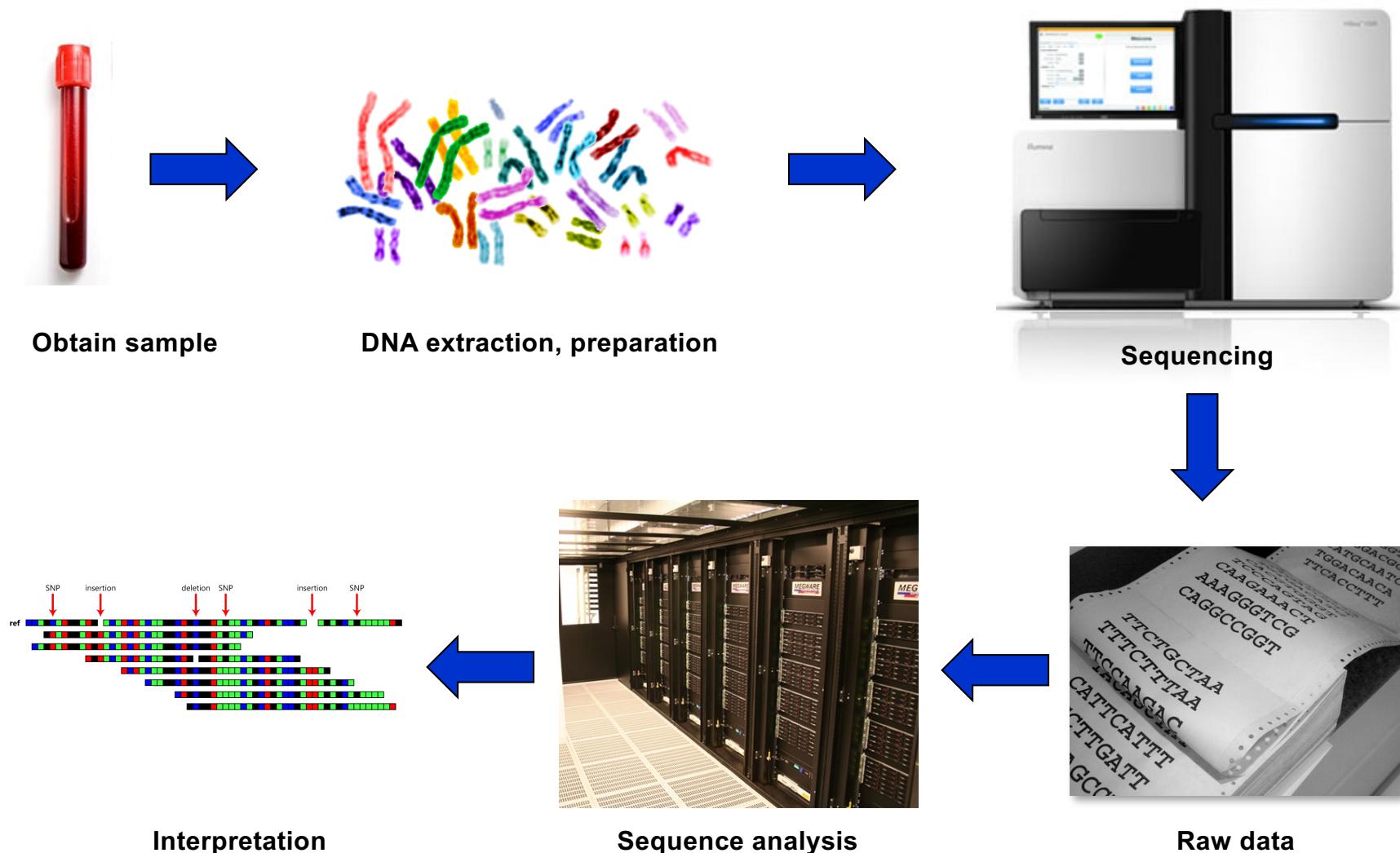
Oslo
University Hospital

Overview

- Sequencing technologies & their general principles
- Important properties & sequencing development
- Sequence quality & file format
- Paired-end reads & mate pair sequencing
- Applications & basic bioinformatics
- Whole genome *de novo* sequencing and assembly
- Resequencing & variant calling
- Other applications
- Challenges

DNA sequencing

High-Throughput Sequencing (HTS), Deep sequencing, Next Generation Sequencing (NGS)



Illumina

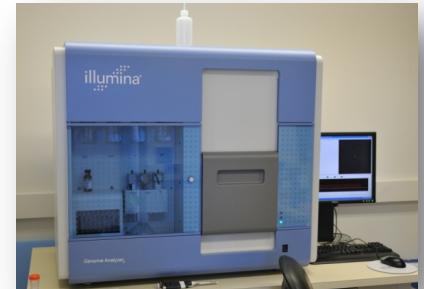
- Sequencing by synthesis using fluorescence
- One fragment = one cluster = one read
- Read lengths up to 250bp, paired-end reads
- Dominant technology today
- Formerly known as Solexa
- HiSeq 2500 specifications:
 - Can sequence entire human genome in 27 hours at 30X coverage (2x100bp)
 - Up to 2x150 bp
 - Run time 7 hours to 11 days
 - Up to 6 billion 100bp reads in 11 days



MiSeq



NovaSeq 6000



GA IIx



HiSeq 2500



Sanger sequencing center

Other sequencing technologies



Roche (454)



ABI (SOLiD)



Ion Torrent

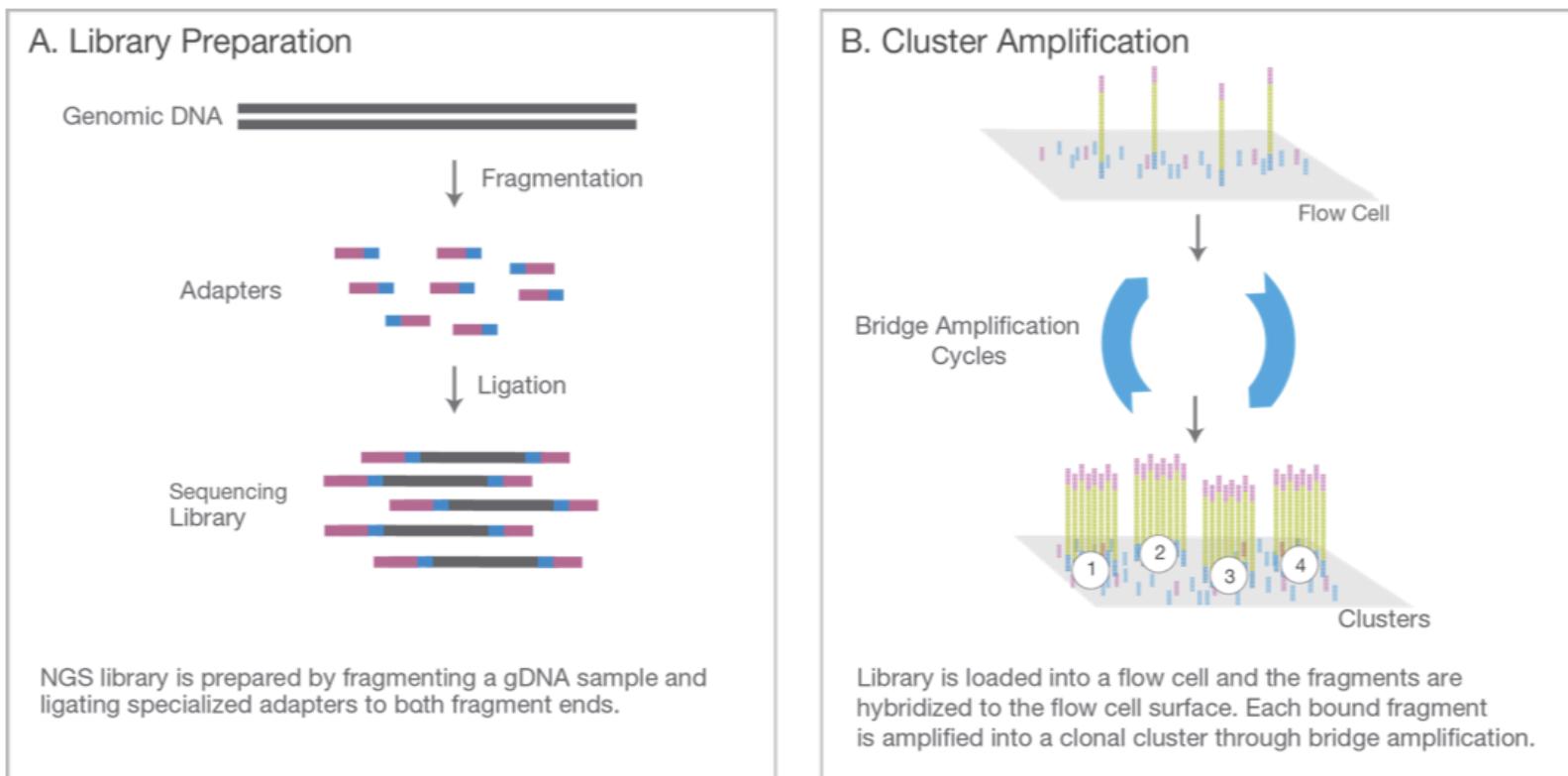


Pacific Biosciences SMRT and Sequel systems



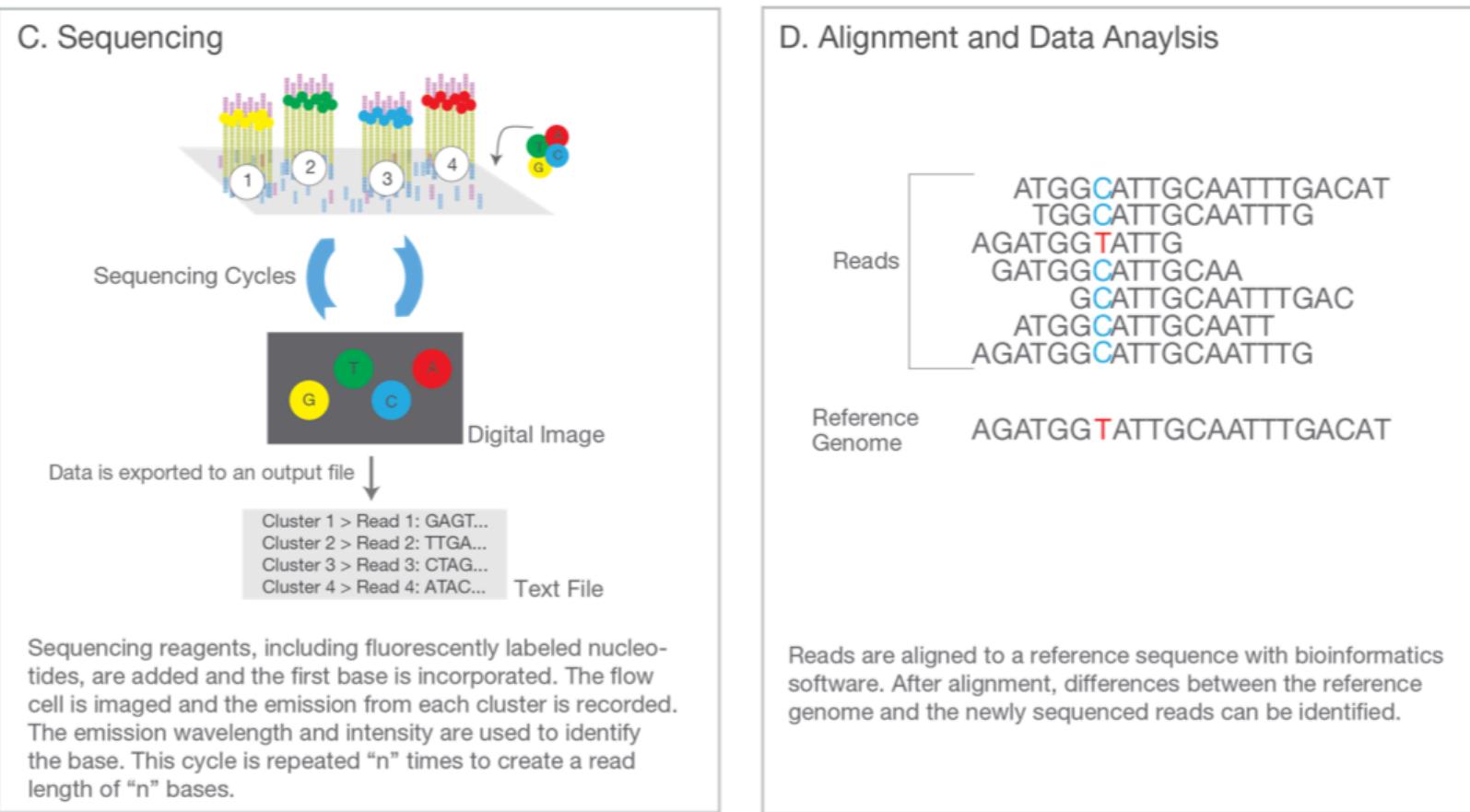
Oxford Nanopore

HTS sequencing principles



Source: Illumina

HTS sequencing principles

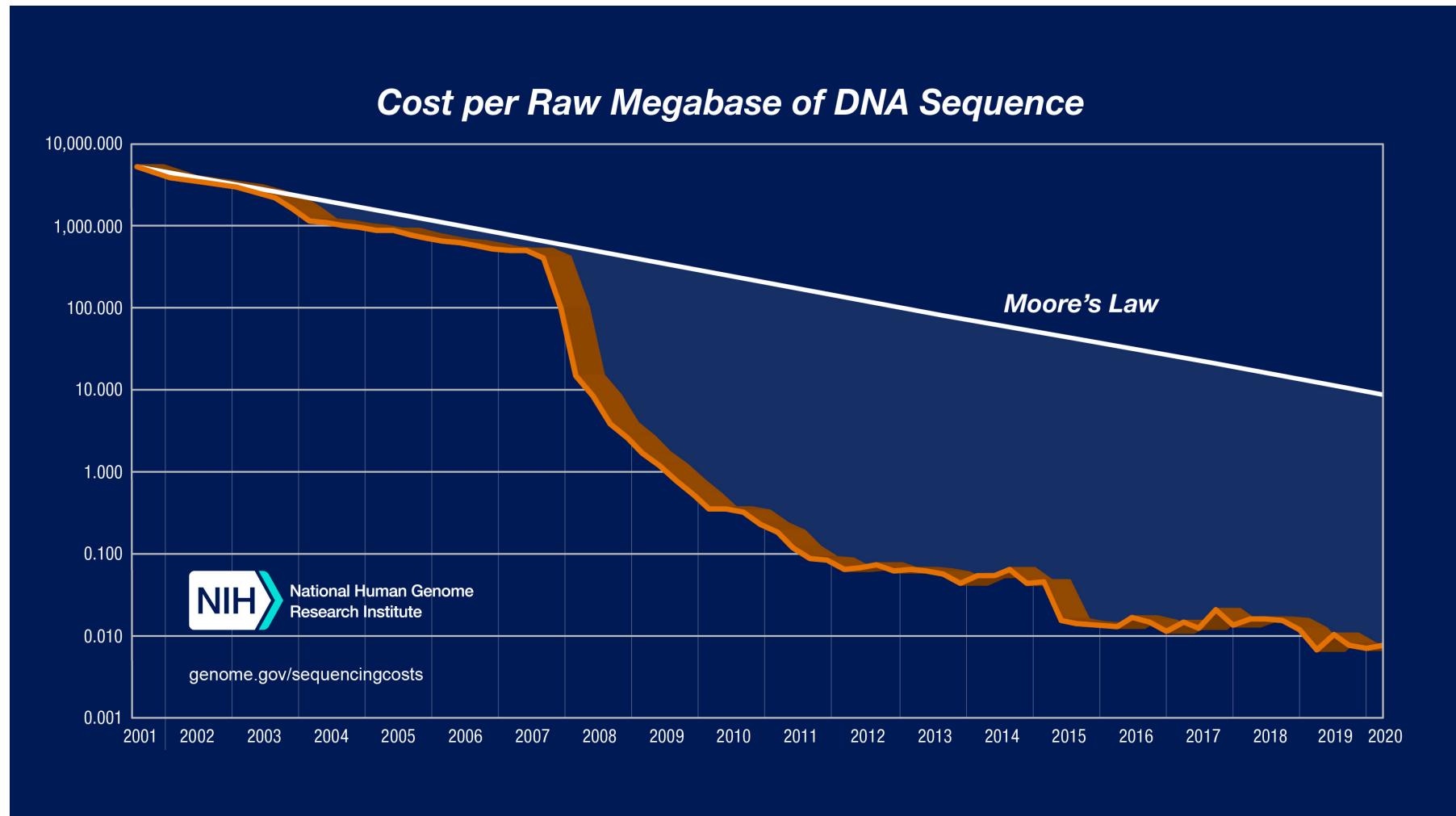


Source: Illumina

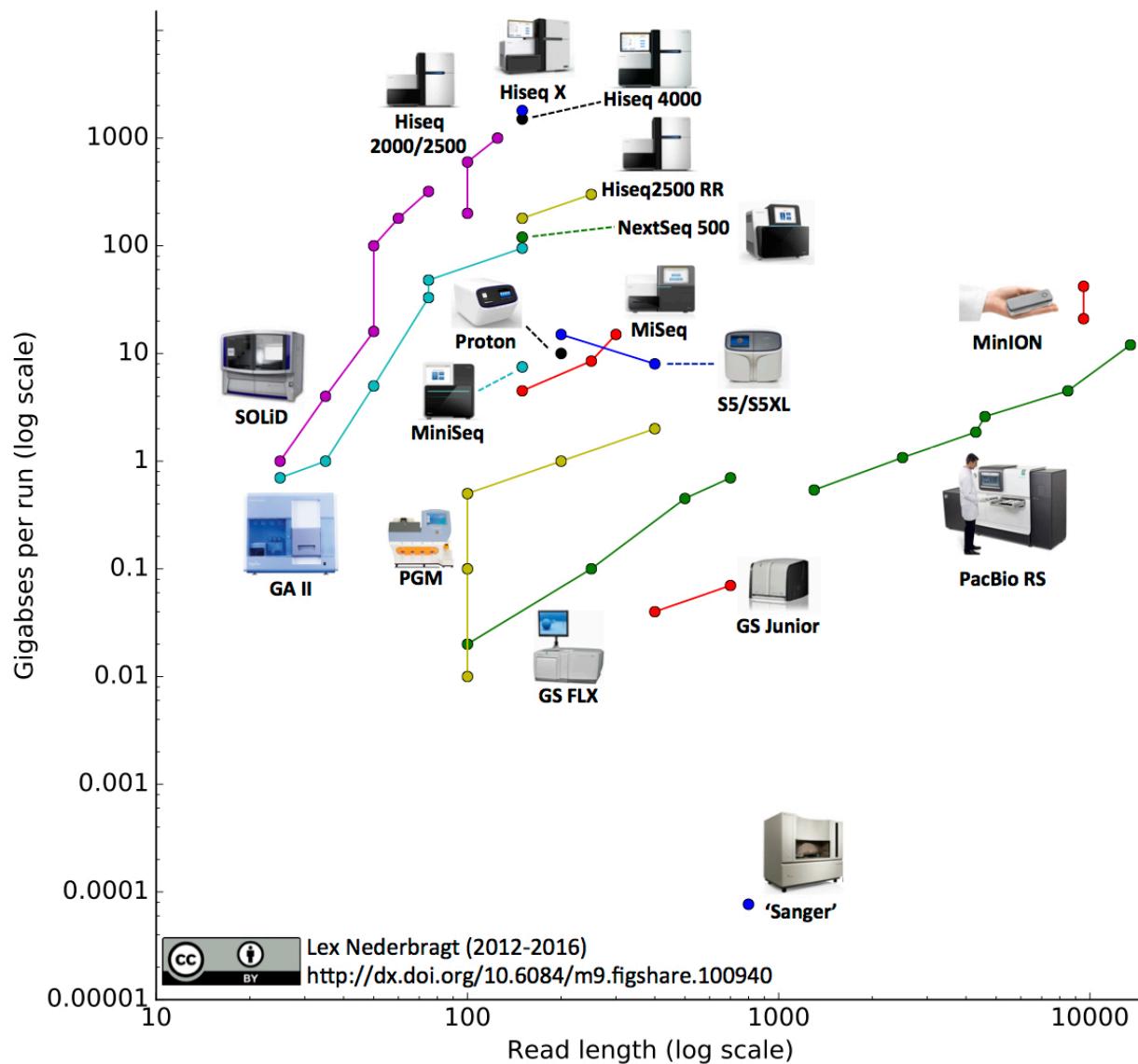
Important technology properties

- Cost
 - Per base
 - Investment
- Read length
- Speed / capacity (bases per day)
- Errors
 - Frequency
 - Profile (indels, substitutions)
 - Random or systematic?
- Paired-end support
- PCR-based?
 - Single molecule
 - PCR amplification step
- Amount of lab work necessary

The cost of sequencing



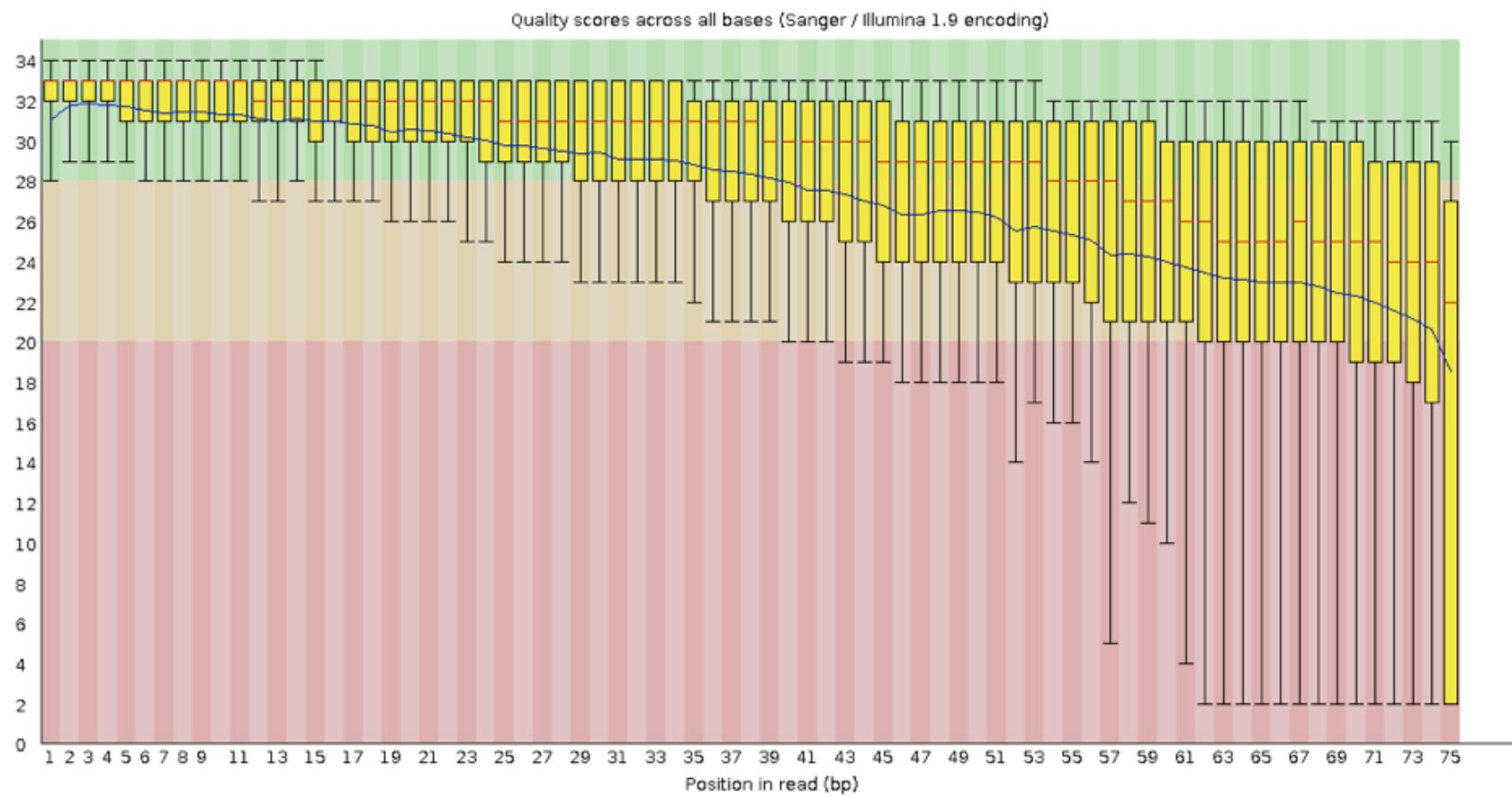
Sequencing technology development



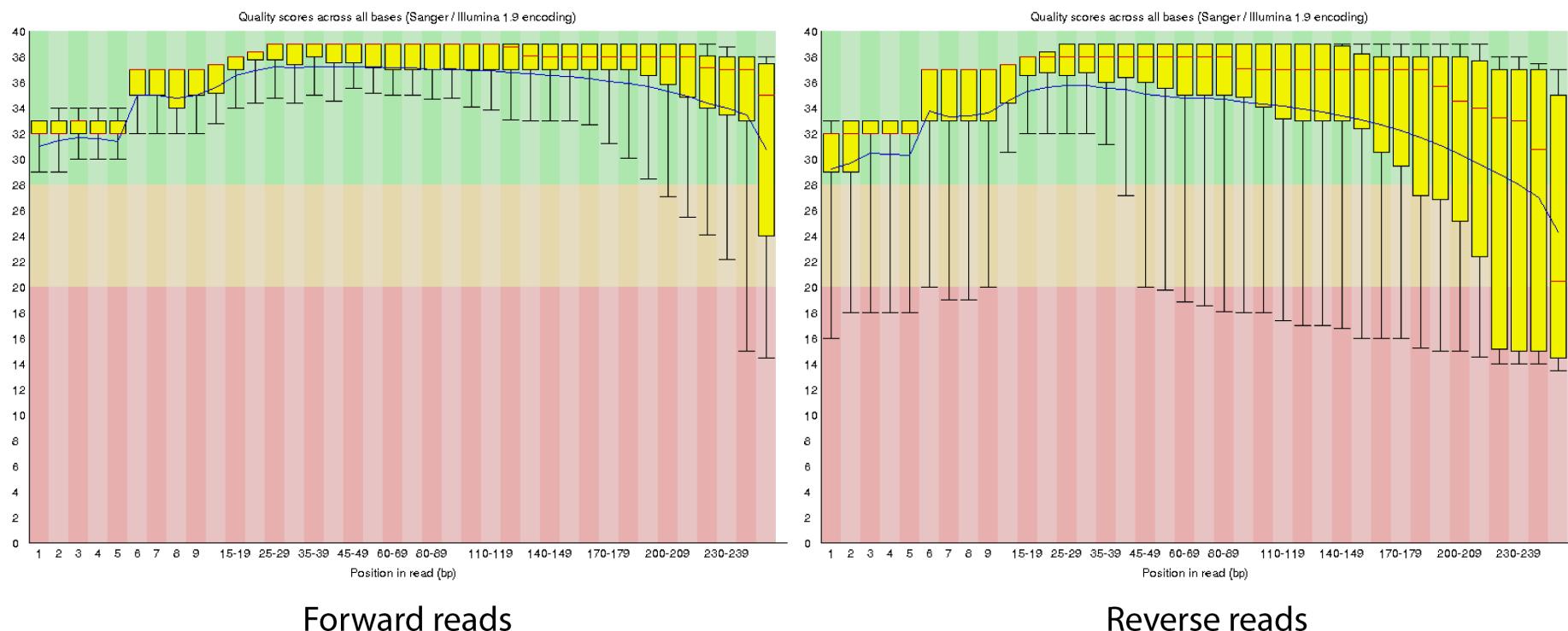
Source: Lex Nederbragt (2012-2016) <https://doi.org/10.6084/m9.figshare.100940>

FASTQC - Quality plot of Illumina reads

✖ Per base sequence quality



Quality plots of Illumina MiSeq reads



The FASTQ format

- A sequence file format in plain text that includes quality scores for each nucleotide in the sequence
- Example:

```
@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGCCGCTGCCGATGGCGTCAAATCCCACC
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```

- The first line starts with a '@' symbol followed by an identifier before the first space.
- The second line contains the actual sequence.
- The third line starts with a '+' symbol, optionally followed by the same identifier as the first line. Identifier rarely used.
- The fourth line contains characters that represent the quality scores for each nucleotide
- In principle, sequence and quality may span multiple lines, but rarely do.

FASTQ quality scores

- Quality value Q (or Phred quality score):
$$Q = -10 \log_{10} p$$
where p = error probability
- The Q values are encoded as ascii characters by adding 33 to the value after rounding to nearest integer:
$$c = 33 + \text{round}(Q)$$
- Only Q values 0-93 used (often just 0-41), corresponding to characters 33-126, and to p values from 1 to $5 \cdot 10^{-10}$
- Example:
 $p=0.0001$ (high quality)
$$Q = -10 \log_{10} 0.0001 = -10 * -4 = 40$$
$$c = 33 + Q = 33 + 40 = 73 = 'I'$$
- Older versions of the format (before 2011) differed slightly

From quality characters to p-values

$$p = 10^{-(c-33)/10}$$

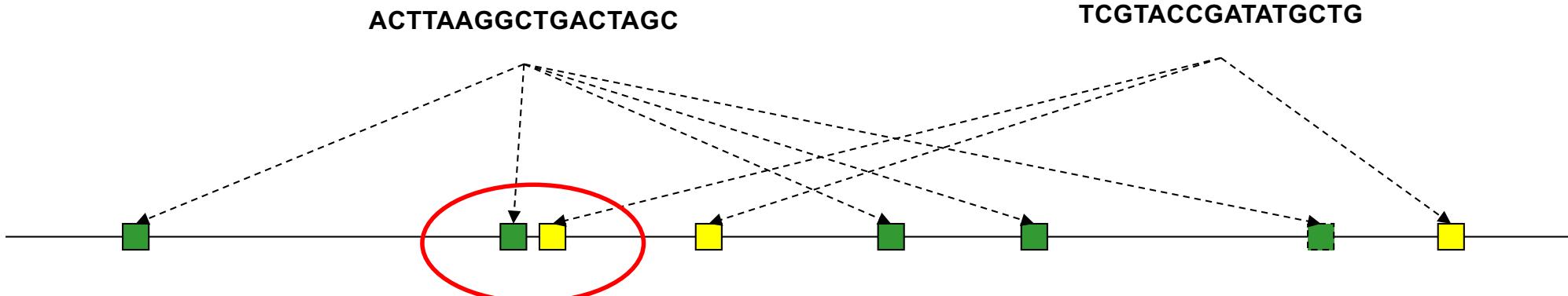
Example:

II9IG9IC

Char	ASCII	Q	p
I	73	40	0.0001
9	57	24	0.0040
G	71	38	0.00016
C	67	34	0.00040

0	<NUL>	32	<SPC>	64	@	96	'
1	<SOH>	33	!	65	A	97	a
2	<STX>	34	"	66	B	98	b
3	<ETX>	35	#	67	C	99	c
4	<EOT>	36	\$	68	D	100	d
5	<ENQ>	37	%	69	E	101	e
6	<ACK>	38	&	70	F	102	f
7	<BEL>	39	'	71	G	103	g
8	<BS>	40	(72	H	104	h
9	<TAB>	41)	73	I	105	i
10	<LF>	42	*	74	J	106	j
11	<VT>	43	+	75	K	107	k
12	<FF>	44	,	76	L	108	l
13	<CR>	45	-	77	M	109	m
14	<SO>	46	.	78	N	110	n
15	<SI>	47	/	79	O	111	o
16	<DLE>	48	0	80	P	112	p
17	<DC1>	49	1	81	Q	113	q
18	<DC2>	50	2	82	R	114	r
19	<DC3>	51	3	83	S	115	s
20	<DC4>	52	4	84	T	116	t
21	<NAK>	53	5	85	U	117	u
22	<SYN>	54	6	86	V	118	v
23	<ETB>	55	7	87	W	119	w
24	<CAN>	56	8	88	X	120	x
25		57	9	89	Y	121	y
26	<SUB>	58	:	90	Z	122	z
27	<ESC>	59	;	91	[123	{
28	<FS>	60	<	92	\	124	
29	<GS>	61	=	93]	125	}
30	<RS>	62	>	94	^	126	~
31	<US>	63	?	95	_	127	

Multiple mapping



- Problem:
 - Short reads (e.g. 100bp) may not map uniquely due to long repeats (>100bp) in the genome
- Solutions:
 - Get longer reads
 - Get paired reads (pairs of reads with fixed distance)

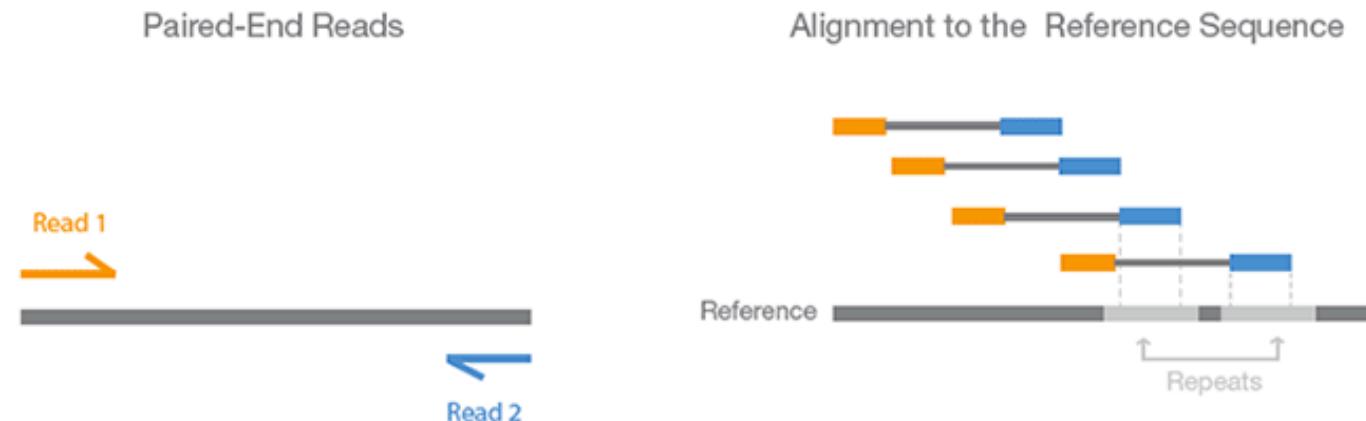
Paired-end / mate pair sequencing

- Paired-end reads or mate pair reads are pairs of reads known to come from two close regions in the genome.
- They are located with an approximate fixed distance from each other.
- Typically paired ends are a ~100-500bp apart, while mate pairs are ~2-10kb apart
- Allows short reads to have a larger "effective" size
- Performed by sequencing fragments from both ends
- Often used with Illumina reads
 - Typically 2 x 100 bp separated by 300bp
 - Allpaths-LG requires ~100 bp from fragments ~ 180bp
- May also overlap (e.g. 2x250bp from 400bp fragments)

Paired-end / mate pair sequencing

- Both ends of fragments of fixed length (a few kbp) may be sequenced
- Gives information on genomic distance between pairs of reads
- May be used to overcome some problems with short reads

Figure 4. Paired-End Sequencing and Alignment



Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

Common HTS applications

De novo genome sequencing	Determining the complete genome sequence of an organism for the first time
Whole genome re-sequencing	Finding polymorphisms (SNPs) and discover mutations in an individual
Exome sequencing	Sequencing only protein-coding regions of a genome from an individual to identify mutations or polymorphisms (SNPs)
Transcriptomics (RNA-seq)	Sequencing of expressed RNA (after reverse transcription to cDNA), (small RNA, mRNA or total RNA) to determine level
Chromatin immunoprecipitation-sequencing (ChIP-Seq) (ChIP-exo)	Mapping of genome-wide protein-DNA interactions
Methylation sequencing (Methyl-Seq)	Determining methylation patterns in the genome (epigenomics) (often on bisulfite-treated DNA)
Metagenomics	Sequencing genomic DNA of multiple species (microorganisms) simultaneously from a certain environment
Metatranscriptomics	Sequencing RNA from multiple species (microorganisms) simultaneously
Amplicon sequencing	Sequencing of genomic regions selected and amplified by PCR, often from multiple species simultaneously

Basic bioinformatics tools

Mapping: Mapping sequence reads to a known genome sequence, used initially in resequencing procedures.

E.g.: BWA, Bowtie, SOAP2, Maq, BFAST, RMAP, ...

Assembly: Assembling together reads into a complete genome sequence, usually divided into a number of contigs and scaffolds. For sequencing entirely new genomes.

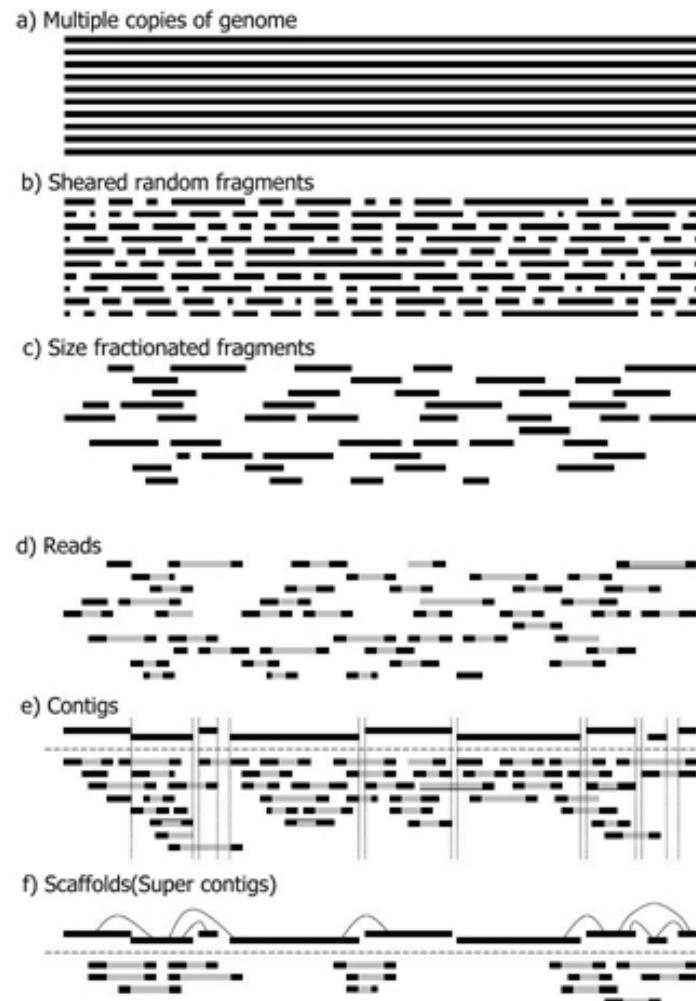
E.g.: Celera/CABOG, Newbler, Phrap, TIGR, Arachne, Velvet, MaSuRCA, SPAdes, ALLPATHS-LG, Abyss, ...

Other:

SNP discovery, Chip-Seq, RNA-Seq, Methyl-Seq, Metagenomics, other mutation/variant discovery

Whole genome *de novo* sequencing

- Whole genome sequencing results in millions of small pieces of the full genome
- The challenge is to puzzle these together in the right order
- Genome sizes ranging from 2Mbp (bacteria) to 3Gbp (human) to 150Gbp (plant)
- Read size from 30 bp to 10000 bp

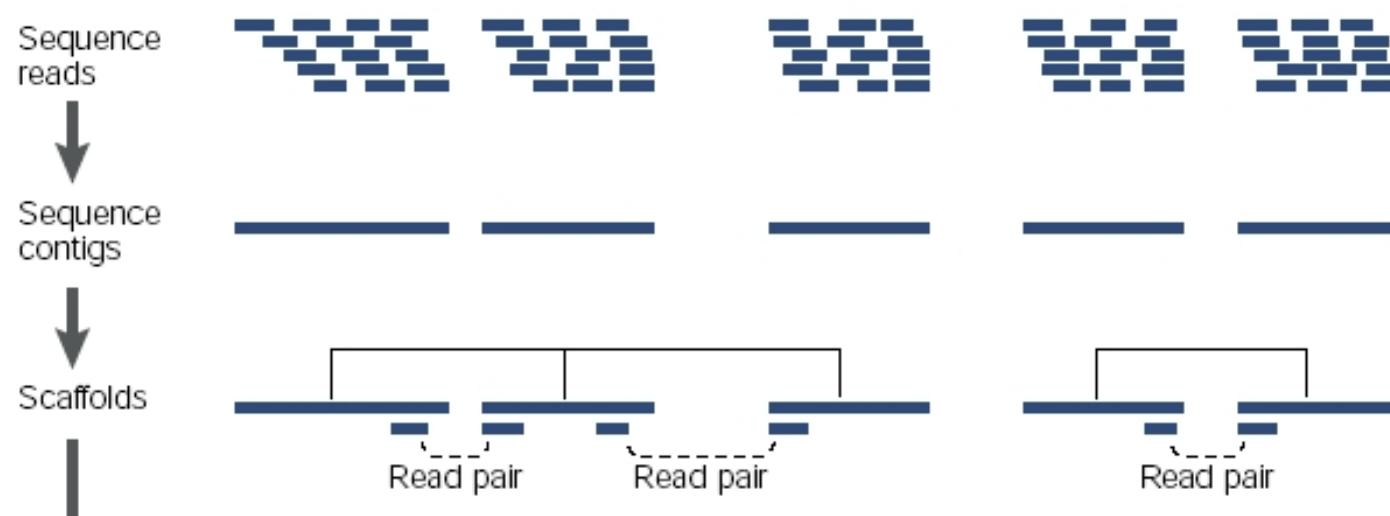


Definitions

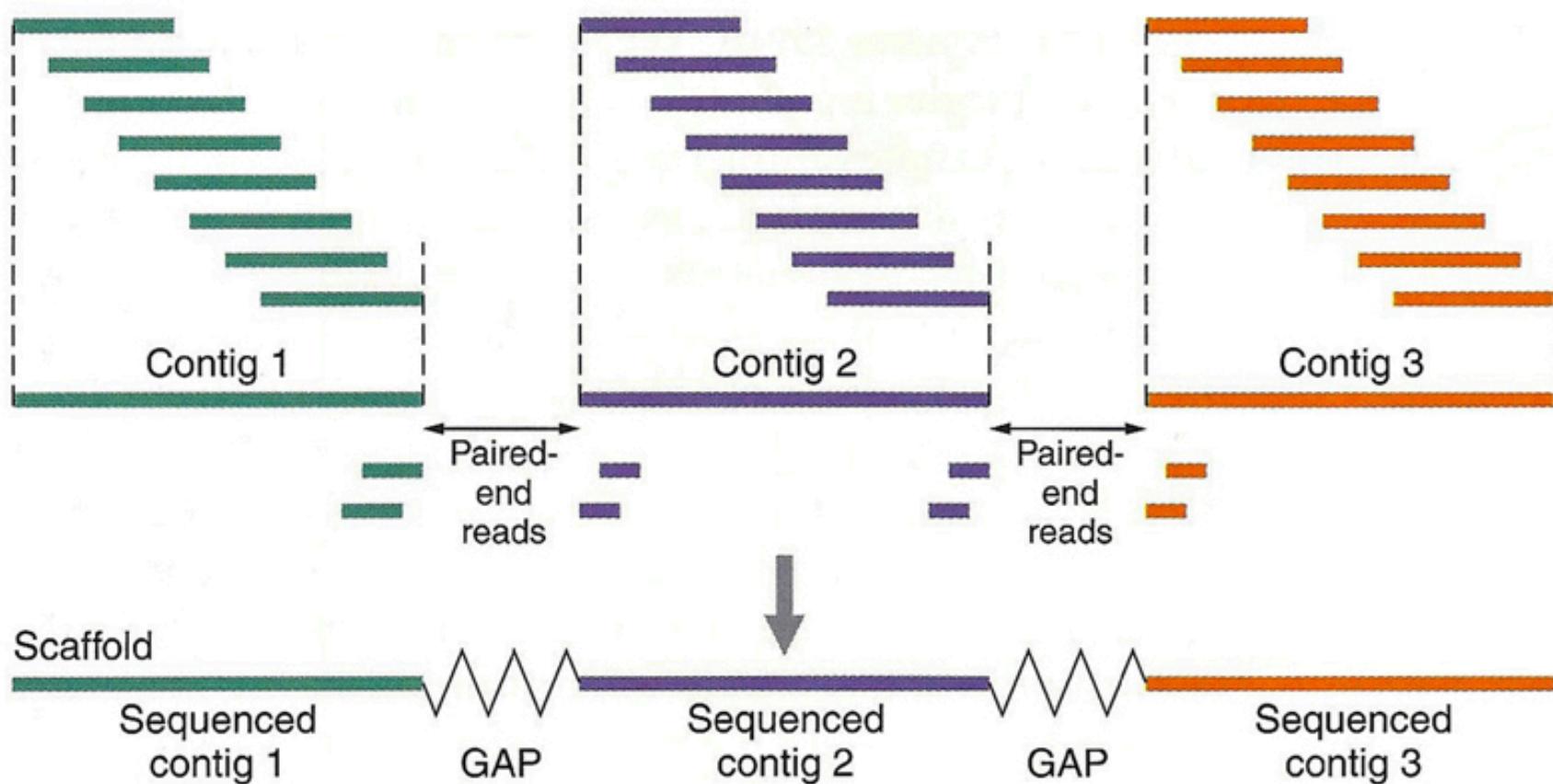
- **Reads** are raw sequences from the sequencers: short continuous sequences
- **Contigs** are longer continuous sequences formed from reads that are partially overlapping. A consensus sequence is based on the reads.
- **Scaffolds** are even longer dis-continuous sequences formed from contigs using information about the distance between contigs and their orientation. Depends upon data from paired-end, mate-pairs or related info
- **An assembly** is the collection of all scaffolds, ideally as few and long and correct as possible

Genome assembly

- Typically whole genome sequencing of novel bacterial species
- Short reads makes eukaryotes hard due to repeats, but not impossible, needs “paired ends”



Paired-end reads span gaps

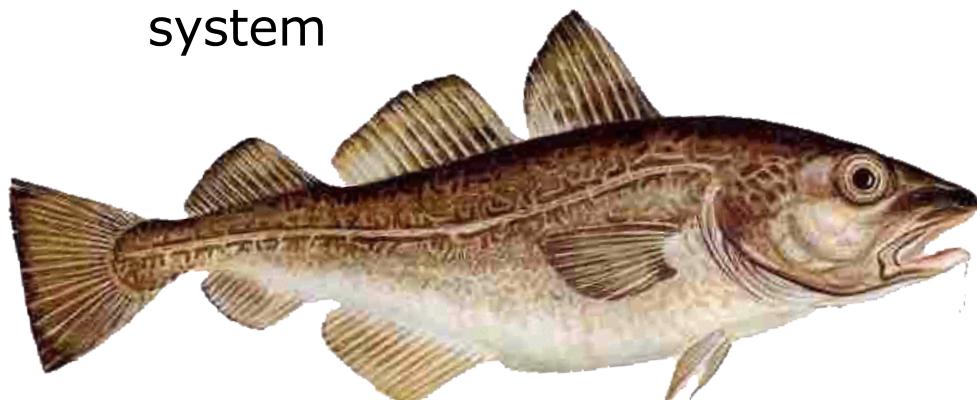


Problematic issues

- Sequencing errors
 - Introduces false sequences into the assembly
 - May be alleviated by higher coverage / larger sequencing depth, or by error detection and correction
- Repeats
 - Genomes often contain many almost identical repeated sequences
 - Repeats longer than the read length makes it impossible to determine the exact location of the read
 - May cause compression or misassemblies
 - May be alleviated by longer reads or paired-end/mate pair reads
- Heterozygosity
 - Diploid organisms (e.g Humans) actually have two “genomes”, not one. Chromosome pairs 1-22 for all, plus XX or XY. One set of chromosomes from our mother and one from our father.
 - The two are mostly identical, but there are some differences
 - Causes “bubbles” in the assembly

The cod genome

- Atlantic cod
- Estimated genome size: 830Mbp
- Sequenced at Dept of Biology, UiO with collaborators
- Started with Roche 454 sequencing
- 40X coverage
- *de novo* assembly using Newbler and Celera Assembler
- 22154 genes identified
- Lacks part of common immune system



LETTER

doi:10.1038/nature10342

The genome sequence of Atlantic cod reveals a unique immune system

Bastiaan Star¹, Alexander J. Nederbragt¹, Sissel Jentoft¹, Unni Grimholt¹, Martin Malmstrom¹, Tone F. Gregers², Trine B. Rourne¹, Jonas Paulsen^{1,3}, Monica H. Solbakken¹, Animesh Sharma³, Ola F. Wetteng^{3,4}, Anders Lanzen^{5,6}, Roger Winer¹, James Knight¹, Jan-Hinnerk Vogel¹⁰, Birtheven Aken¹¹, Øivind Andersen¹¹, Karin Lagesen¹, Ave Tooming-Klerterud¹, Rolf B. Edvardsen¹², Kirubakaran G. Tina^{1,13}, Mari Espelund¹, Chirag Nepal^{1,8}, Christopher Previti¹, Bård Ove Karslien¹⁴, Truls Mounø¹⁵, Morten Skage¹, Paul R. Berg¹, Tor Gjøen¹⁶, Helmer Kuh¹⁶, Jim Thorsen¹⁷, Ketil Malde¹⁷, Richard Reinhard¹⁶, Leif Du¹⁸, Steinar D. Johansen^{4,18}, Steve Seal¹⁹, Sigbjørn Lien¹³, Frank Nilsen¹⁹, Inge Jonassen^{4,18}, Stig W. Omholt^{1,23}, Nils Chr. Stenseth¹ & Kjetill S. Jakobsen¹

Atlantic cod (*Gadus morhua*) is a large, cold-adapted teleost that sustains long-standing commercial fisheries and incipient aquaculture^{1,2}. Here we present the genome sequence of Atlantic cod, showing evidence for complex thermal adaptations in its haemoglobin gene cluster and an unusual genome architecture compared to other teleosts and vertebrates. The genome was initially assembled exclusively by 454 sequencing of shotgun and paired-end libraries, and automated annotation identified 22,154 genes. The major histocompatibility complex (MHC) II is a conserved feature of the adaptive immune system of jawed vertebrates^{3,4}, but we show that Atlantic cod has lost the genes for MHC-II, CD4 and invariant chain (Ii) that are essential for the function of this pathway. Nevertheless, Atlantic cod is not exceptionally susceptible to disease under natural conditions⁵. We find a highly expanded number of MHC-I genes and a unique composition of its Toll-like receptor (TLR) families. This indicates that the Atlantic cod immune system has evolved compensatory mechanisms in both adaptive and innate immunity in the absence of MHC-II. These observations affect fundamental assumptions about the evolution of the adaptive immune system and its components in vertebrates.

We sequenced the genome of a heterozygous male Atlantic cod (NEAC_001, Supplementary Note 1 and 2), applying a whole-genome shotgun approach to 40X coverage (estimated genome size 830 megabases (Mb), Supplementary Note 4 and Supplementary Fig. 2) using 454 technology (Supplementary Note 3). Two programs (Newbler⁶ and Celera⁷, Supplementary Notes 5 and 6) produced assemblies with short contigs yet with scaffold sizes comparable to those of Sanger-sequenced teleost genomes (Supplementary Note 10 and Supplementary Fig. 8). Although fragmentation of the short reads repeats the same address (Supplementary Note 7), the resolved numerous gaps attributable to heterozygosity (Supplementary Note 8). The assemblies differ in scaffold and contig length (Table 1), although their scaffolds align to a large extent (Supplementary Note 9 and Supplementary Fig. 7). We obtained about one million single nucleotide polymorphisms (SNPs) by mapping 454 and Illumina reads from the sequenced individual to the Newbler assembly (Supplementary Note 11). Both assemblies cover more than 98% of the reads from an extensive transcriptome data set, indicating that the proteome is well represented (Supplementary Note 13). The assemblies are consistent with four

independently assembled bacterial artificial chromosome (BAC) insert clones (Supplementary Note 14 and Supplementary Fig. 9), and with the expected insert size of paired BAC-end reads (Supplementary Note 15 and Supplementary Fig. 10).

A standard annotation approach based on protein evidence was complemented by a whole-genome alignment of the Atlantic cod with the stickleback (*Gasterosteus aculeatus*), after repeat-masking 25.4% of the Newbler assembly (Supplementary Note 16 and Supplementary Table 6). In this way, 17,920 out of 20,787 protein-coding stickleback genes were mapped onto reorganized scaffolds (Supplementary Note 17). Additionally, protein-coding genes, pseudogenes and non-coding RNAs were annotated using the standard transcriptome approaches resulting in a final gene set of 22,154 genes (Supplementary Note 7). Comparative analysis of gene ontology classes indicates that the major functional pathways are represented in the annotated gene set (Supplementary Note 18 and Supplementary Fig. 11). We anchored 332 Mb of the Newbler assembly to 23 linkage groups of an existing Atlantic cod linkage map using 924 SNP⁸ (Supplementary Note 19 and Supplementary Table 8). These linkage groups have distinct orthology to chromosomes of other teleosts, on the basis of the number of co-occurring genes, showing that the whole-genome shotgun assembly reflects the expected chromosomal ancestry (Fig. 1, Supplementary Note 20 and Supplementary Table 9).

Table 1 | Assembly statistics

	Number	Bases (Mb)	N50L (bp)*	N50 (bp)†	M _L (bp)‡
Newbler					
Contigs [§]	284,239	636	2,776	50,237	76,504
Scaffolds [§]	6,467	1	687,709	218	4,999,318
Entire assembly	157,887	753	459,495	344	4,999,318
Celera					
Contigs [§]	135,024	555	7,128	19,938	117,463
Scaffolds [§]	3,832	608	488,312	373	2,810,583
Entire assembly	17,039	629	469,840	395	2,810,583

*Number of sequences with lengths in half of the standard total occur.

†Number of sequences with lengths of N50, or longer.

‡Maximum length.

§Sequences with lengths of 500 bp.

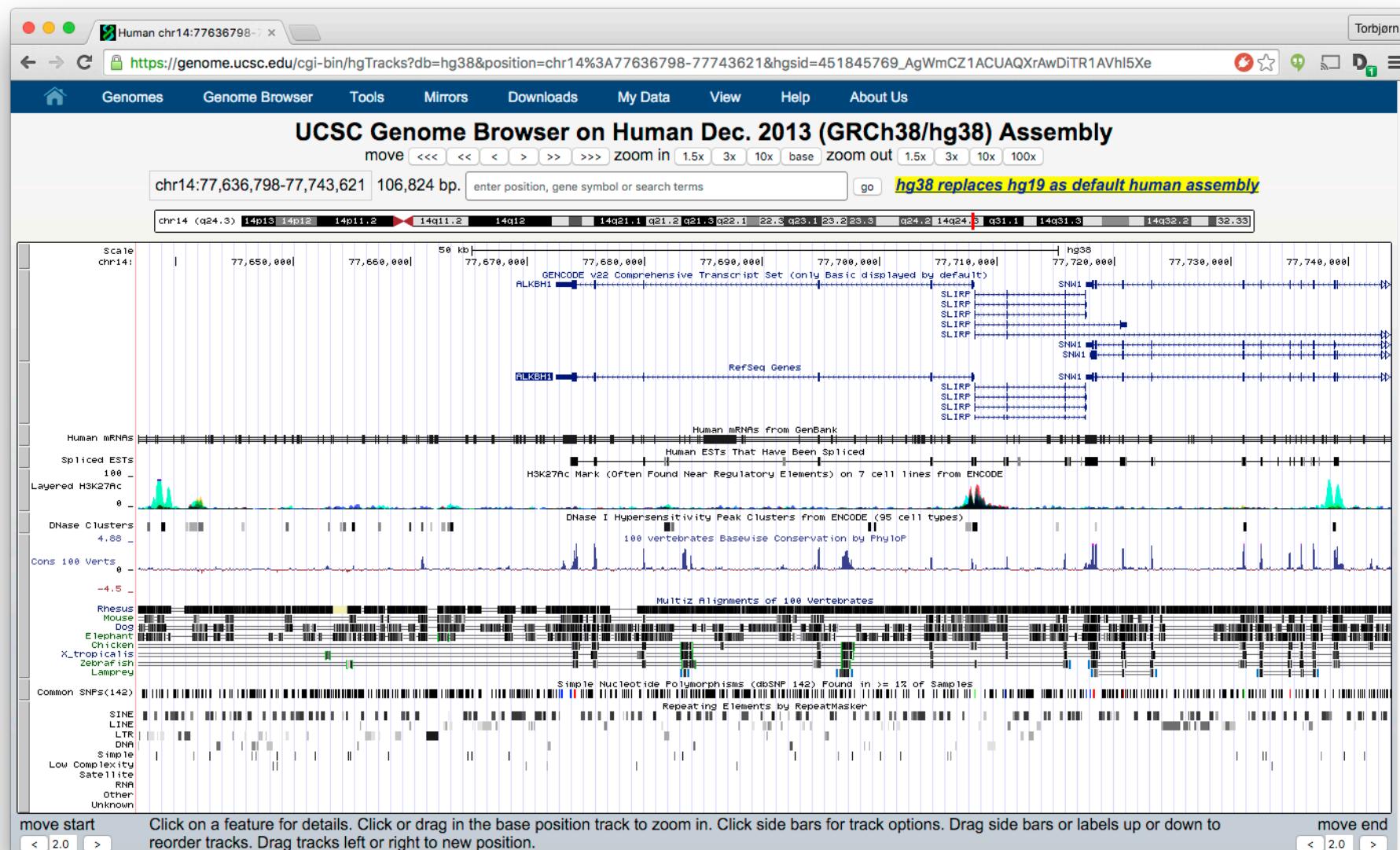
||Sequences with lengths of 500 bp.

¹Centre for Ecological and Evolutionary Synthesis (CEES), Department of Biology, University of Oslo, PO Box 1306, Blindern, N-0316 Oslo, Norway. ²Department of Molecular Biosciences, Centre for Integrative Biogeochemistry, University of Oslo, Blindern, N-0316 Oslo, Norway. ³Bioperformatics Core Facility, Institute for Medical Informatics, Oslo University Hospital, Montebello, N-0310 Oslo, Norway. ⁴Department of Informatics, University of Bergen, N-5020 Bergen, Norway. ⁵Department of Natural Sciences and Technology, Høgskolen i Vestfold University College, P.O. Box 4010, Kongsberg, N-3930 Hænar, Norway. ⁶Computational Biology Unit, Uni Computing Uni Research AS, N-5020 Bergen, Norway. ⁷Life Sciences, 15 Commercial Street, Branford, Connecticut 06405, USA. ⁸Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge CB10 1SD, UK. ⁹Department of Animal and Aquacultural Sciences, University of Oslo, PO Box 1000, N-0316 Oslo, Norway. ¹⁰Max Planck Institute for Molecular Genetics, Ihnestrasse 63-73, D-14195 Berlin-Dahlem, Germany. ¹¹Institute for Basic Sciences and Aquatic Medicine, School of Veterinary Sciences, N-0003 Oslo, Norway. ¹²Department of Animal and Aquacultural Sciences, CEGENE, Centre for Genomic Evolution, Norwegian University of Life Sciences, PO Box 5030, 1432 Ås, Norway. ¹³Faculty of Biosciences and Aquaculture, University of Nordland, N-8049 Bodø, Norway. ¹⁴Department of Pharmaceutical Biosciences, School of Pharmacy, University of Oslo, P.O. Box 1068, Blindern, N-0316 Oslo, Norway. ¹⁵Institute for Basic Sciences and Aquatic Medicine, School of Veterinary Sciences, University of Tromsø, N-9037 Tromsø, Norway. ¹⁶Department of Biomedicine, Faculty of Health Sciences, University of Tromsø, N-9037 Tromsø, Norway. ¹⁷Department of Biology, PO Box 7800, University of Bergen, N-5020 Bergen, Norway.

8 SEPTEMBER 2011 | VOL 477 | NATURE | 207

©2011 Macmillan Publishers Limited. All rights reserved

Genome browsers



Source: genome.ucsc.edu

Mapping reads to a reference genome

Goal: Identify positions in the genome that are most similar to the sequence reads

Input data:

- 10-1000 million reads, each 30-300bp
- Sequencing errors (typ. ~1% error rate)

Reference genome:

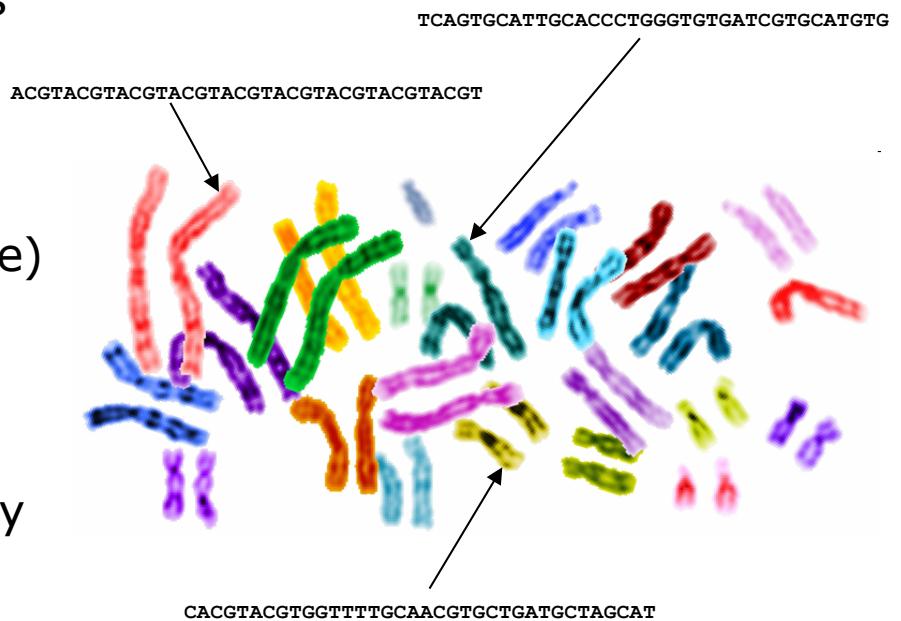
- E.g. human genome, 3 Gbp
- Some genome variation, heterozygosity

Output:

- 0, 1, or more potential genomic locations for each read
- Mapping quality assignment

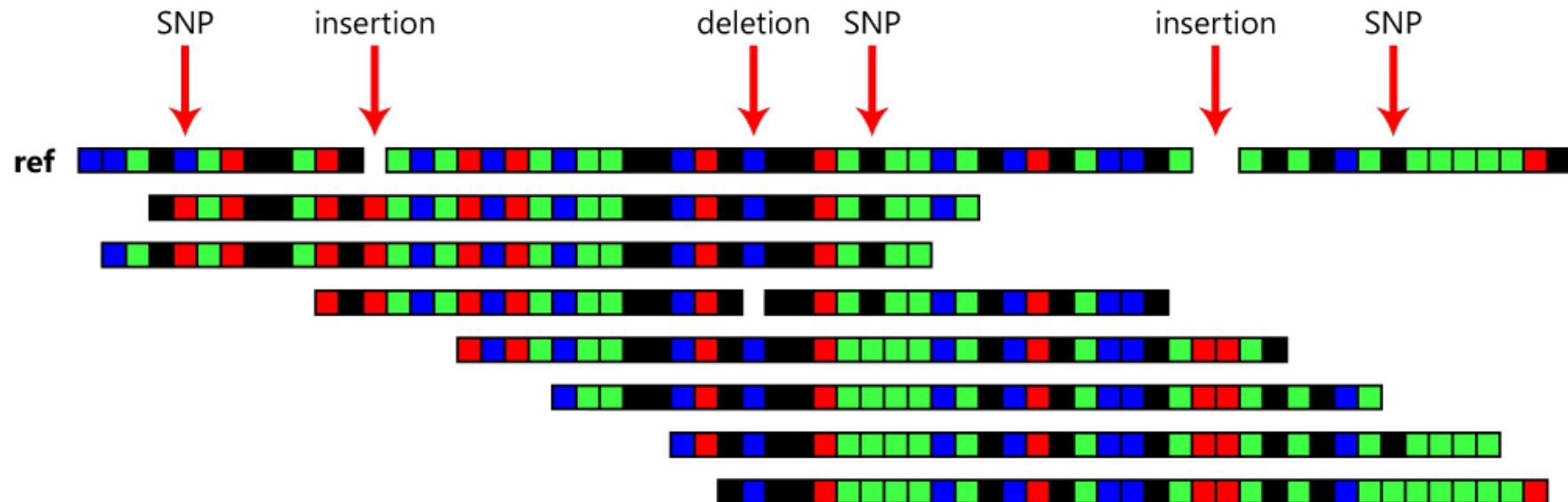
Requirements:

- Sensitivity, specificity, speed, compactness



Resequencing

- Sequencing DNA from a new individual when we already have a reference genome sequence
- Map reads to reference genome instead of assembly



Variation detection by resequencing

- Natural variation discovery
- Mutation detection
- Single Nucleotide Polymorphisms (SNPs)
- Small insertions & deletions (Indels)
- Copy Number Variation (CNV)
- Large inversions, translocations etc
- Requires high coverage, that is, the average number of times each base is sequenced (typically 40X, but may require 100X)

The diagram illustrates a DNA sequence with several variations. The first variation is highlighted with a purple underline and a red box around the last base, labeled 'SNP'. The second variation is highlighted with a red underline and a red arrow pointing to the 'g' base, labeled 'sequencing errors'. The third variation is highlighted with a red underline and a red arrow pointing to the 'a' base, labeled 'sequencing errors'. The fourth variation is highlighted with a red underline and a red arrow pointing to the 'c' base, labeled 'sequencing errors'.

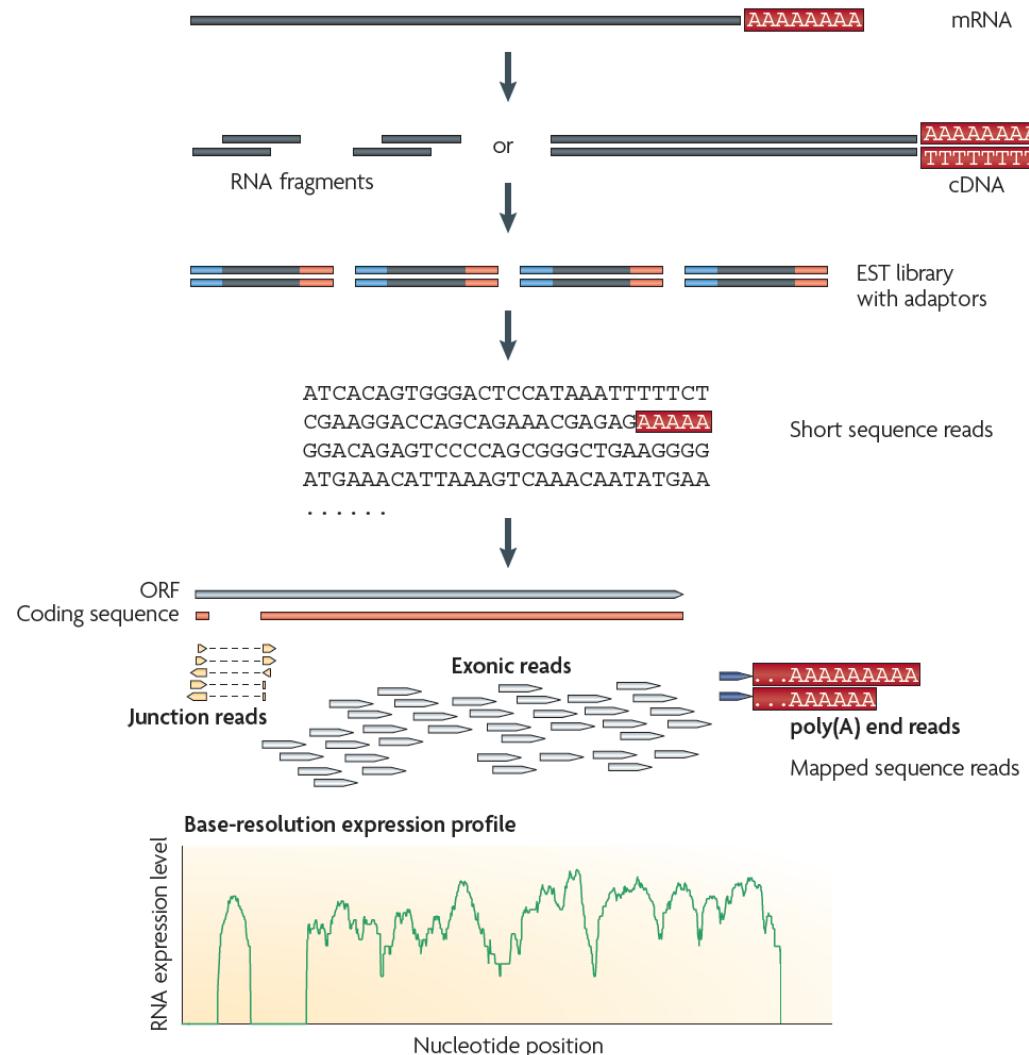
GTTACTGTCGTTGTAATACTCCACGATGTC
GTTACTGTCGTTGTAATACTCCACGATGTC
GTTACTGTCGTTGTAATACTCCACAAATGTC
GTTACTGTCGTTGTAATgCTCCACGATGTC
GTTACTGTCGTTGTAATACTCCACAAATGTC
GTTACTGTCGTTGTAATACTCCACGATGTC
GTTACTGTCGTGGTAATACTCCACaATGTC
GTTACTGTCGTTGTAATACTCCACaATGTC
GTTAaTGTGCGTTGTAATACTCCACGATGTC
GTTACTGTCGTTGTAcTACTCCACGATGTC
GTTACTGTCGTTGTAATACTCCACaATGTC

↑ ↑ ↑ ↑
sequencing errors SNP

Mapping and coverage

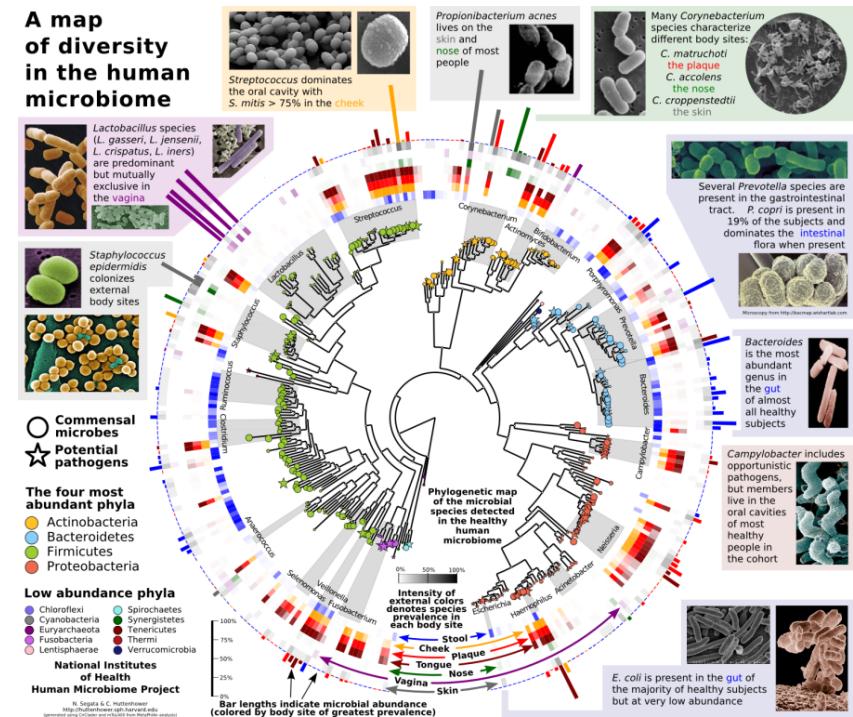
Gene expression (RNA-Seq)

- Gene expression analysis
- “transcriptomics”
- Replaces microarrays
- mRNAs
- Small RNAs (miRNA, piRNA...)
- Splice variants
- Counts the number of reads for each RNA



Metagenomics/metatranscriptomics

- Samples contains collection of DNA/RNA from many microorganisms present in some niche - a microbial community
- Sources: Soil, ocean, mine, human body, the built environment, ...
- Ecological diversity studies
- Clinical studies (e.g. human gut)
- Big data: Many hundred million sequences



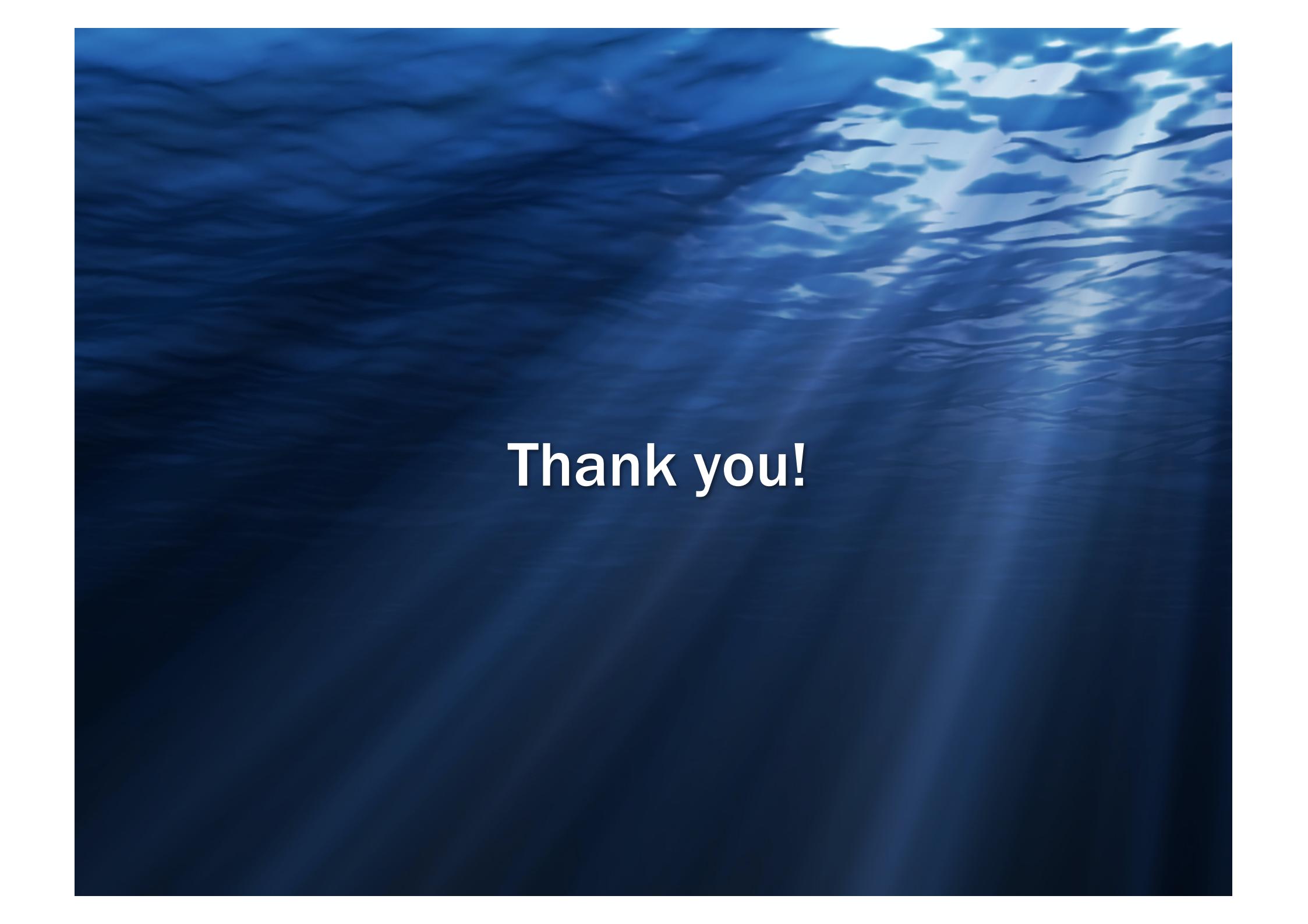
TARA OCEANS



Human Microbiome Project

Challenges

- Cost of actual sequencing is decreasing, but what about the cost of analysis?
- Lack of competent people for bioinformatics analysis
- Large storage needs due to the amounts of data generated. Terabytes of data.
- Compute intensive analysis (read mapping, assembly, etc)
- Security and privacy issues related to sensitive human data



Thank you!