

Genome & Transcriptome studies using Next Generation Sequencing (NGS) Technologies : Past and Present



Fatma Guerfali, PhD
Institut Pasteur de Tunis

fatma.guerfali@gmail.com

Institut Pasteur de Tunis

1893 – 3rd IP founded

Over a century at the service
of Public Health



Vaccination centre

Major discoveries endowed
IPT with an uncontested
international renown



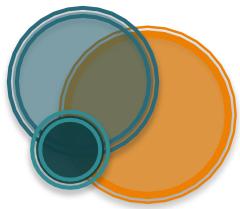
Major public health
impact: Eradication of
trachoma, malaria and
schistosomiasis

Charles Nicolle
(Nobel prize in Medicine, 1928)

1956 (independence)
Preparation of vaccines
Epidemiology
Diagnostic centre
Education



- Next-generation sequencing (**NGS**), or high-throughput (**HT**) sequencing = **catch-all term** describing different modern sequencing technologies used by different platforms.
- Many variations
 - 'DNA-Seq' and other related 'seq' technologies allow to cover genome complexity : genomic DNA-Seq, Methyl-Seq, ChIP-Seq, exome sequencing...



INTRODUCTION

BASIC CONCEPTS OF NGS



Part 1 Genomes and Transcriptomes : From central dogmas to ongoing discoveries

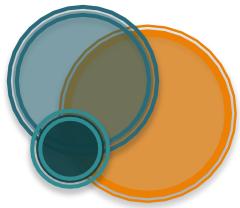
(Central dogma of molecular Biology, Human Genome Project, GENCODE, ENCODE...)

Part 2 Advances in Sequencing Technologies

(Examples of technologies developed for Closed/Open systems for gene expression analysis, Next-Generation Sequencing platforms and technologies ...)

Part 3 Overview of NGS (DNA / RNA Seq) Protocols and related file formats

(Overview of protocols for Genomic and Transcriptomic analysis of standard samples using Next-Generation Sequencing and overview of the different formats generated at each step...)



Part 1 Genomes and Transcriptomes : From central dogmas to ongoing discoveries

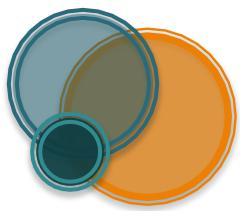
(Central dogma of molecular Biology, Human Genome Project, GENCODE, ENCODE...)

Part 2 Advances in Sequencing Technologies

(Examples of technologies developed for Closed/Open systems for gene expression analysis, Next-Generation Sequencing platforms and technologies ...)

Part 3 Overview of NGS (DNA / RNA Seq) Protocols and related file formats

(Overview of protocols for Genomic and Transcriptomic analysis of standard samples using Next-Generation Sequencing and overview of the different formats generated at each step...)





GENOMES & TRANSCRIPTOMES

CENTRAL DOGMAS

- Proposed by Francis Crick 1958
- DNA holds the coded hereditary information in the nucleus
- This code is expressed at the ribosome during protein synthesis in the cytoplasm
- The protein produced by the genetic information is what is influenced by natural selection
- If a protein is modified it cannot influence the gene that codes for it. Therefore only one way flow of information was thought to occur:

DNA→RNA→Protein

GENOMES & TRANSCRIPTOMES

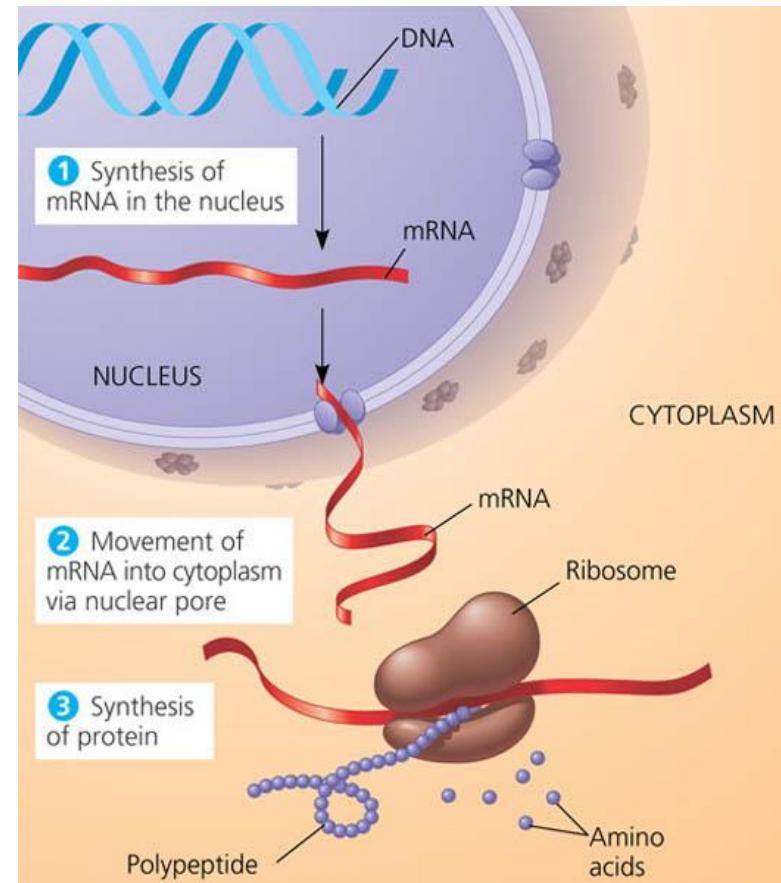
CENTRAL DOGMAS

What is a gene? (definition until less than 10 years ago)

Gene

Is defined as a DNA portion whose corresponding mRNA encodes a protein.

= RNA only considered to be a "bridge" in the transfer of biological information between DNA and proteins



PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Costa et al., 2010)

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

BASIC CONCEPTS OF GENE AND GENE EXPRESSION

What is a gene ?

Prokaryotes

"cistron" = gene

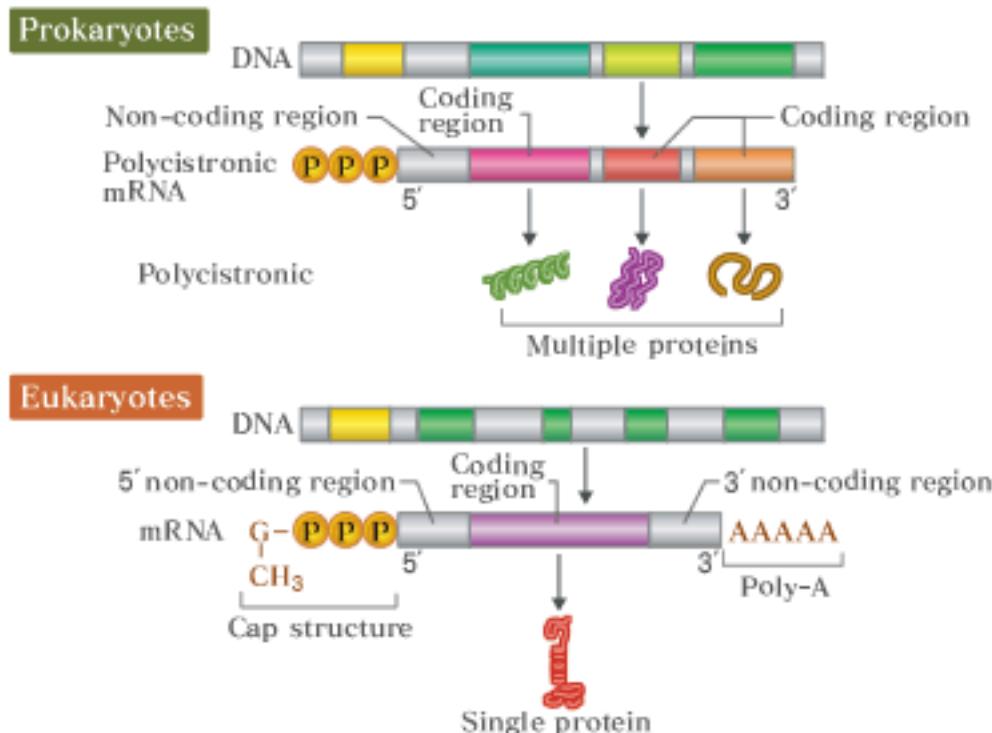
"operon" = a gene unit controlled for gene expression by the same regulating region (e.g. Lactose, Histidine : if present in the medium, the prokaryotic cell suppress all genes related to its synthesis by the genes in operon).

Eukaryotes

In general: Eukaryotes do not have operons and hence no polycistronic RNA is generated

→ Exceptions !!!! (Leishmania)

PART
1



GENOMES & TRANSCRIPTOMES

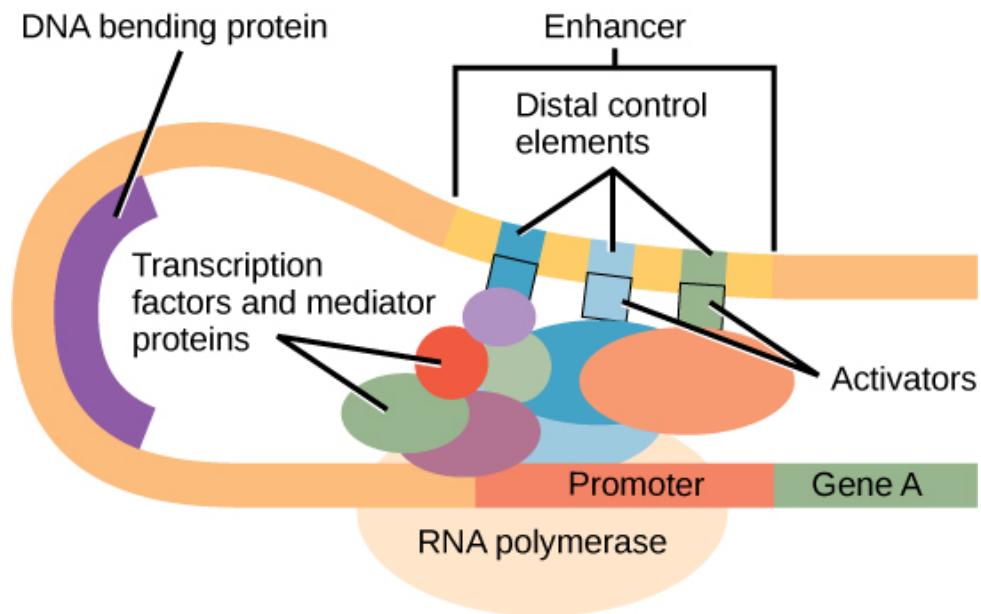
BASIC CONCEPTS OF GENE AND GENE EXPRESSION

► What controls gene expression?

Transcription factors (TFs) = proteins that bind to the DNA and help to control gene expression. We call the sequences to which they bind **transcription factor binding sites (TFBSs)**, which are a type of *cis*-regulatory sequence.

There are two main type of *cis*-regulatory elements:

- promoters
- *cis*-regulatory modules (sometimes called “enhancers”).



<https://www.boundless.com/biology/textbooks/boundless-biology-textbook/>

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

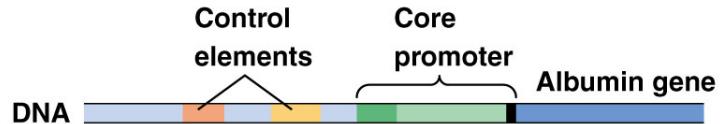
NGS

FATMA GUERFALI

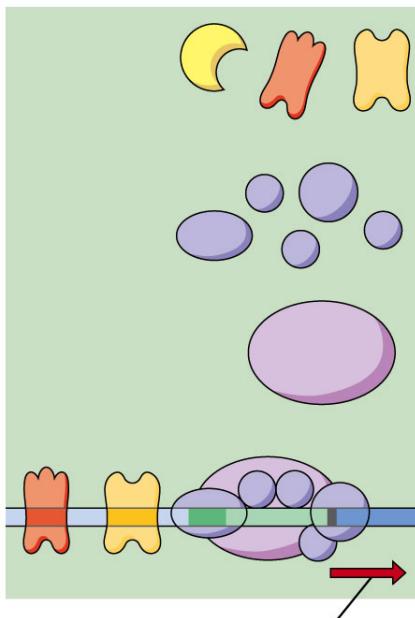
GENOMES & TRANSCRIPTOMES

BASIC CONCEPTS OF GENE AND GENE EXPRESSION

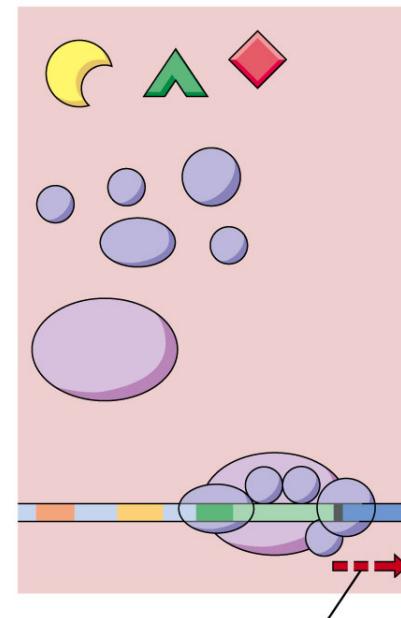
► What controls gene expression?



(a) Liver cell nucleus



(b) Brain cell nucleus



It is the set of regulatory Transcription factors (TFs) that can be different in different tissues

PART
1

Albumin gene transcribed
at high level

© 2012 Pearson Education, Inc.

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

http://www.mun.ca/biology/desmid/brian/BIOL2060/BIOL2060-23/23_23.jpg

NGS

FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

BASIC CONCEPTS OF GENE AND GENE EXPRESSION



What controls gene expression?

Prokaryotic Gene Regulation

- Regulation of the lac operon (dual control: repression and promotion)
- Regulation of the trp operon (a "riboswitch")

Eukaryotic Gene Regulation

- Genomic Control (DNA rearrangements...)
- Transcriptional Control (TFs, Enhancers...)
- Translational Control (isoforms, dose-dependant synthesis (iron...), siRNAs/miRNAs...)
- Post-translational Control (active/inactive forms (kinases...), ubiquitin degradation...)

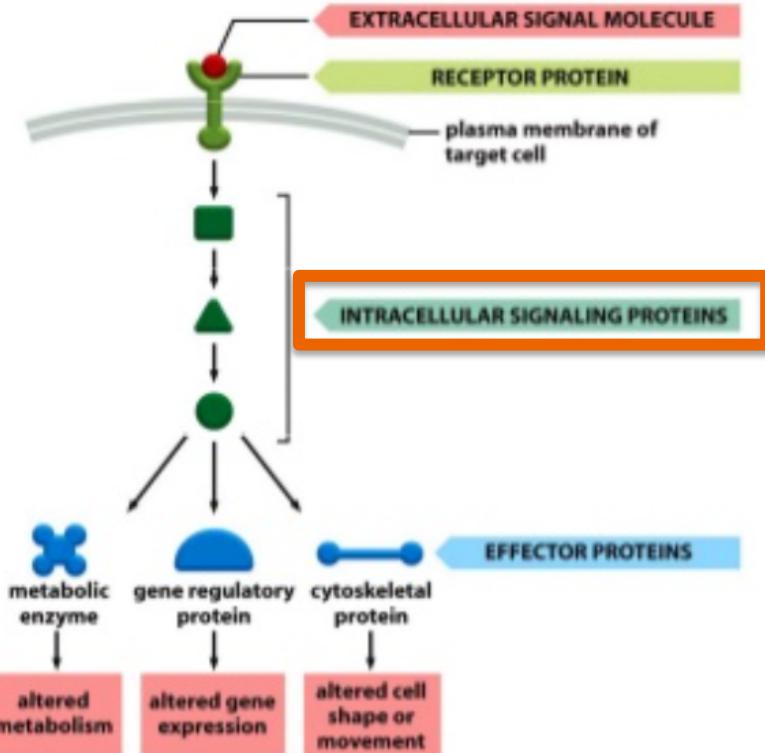
GENOMES & TRANSCRIPTOMES

BASIC CONCEPTS OF GENE AND GENE EXPRESSION

► What controls gene expression?

Overview of Cell Signaling

- Sources of extracellular signal
 - Non-cellular environment
 - Cellular environment (cell-cell communication)
 - Hundreds of types of signals
- Cells signaling
 - Stimulus sensing; communication
 - Information processing; decision making
- ↓Receptors
↓Signal transducers
↓Effector proteins
- Signaling pathways regulate nearly all cellular functions.



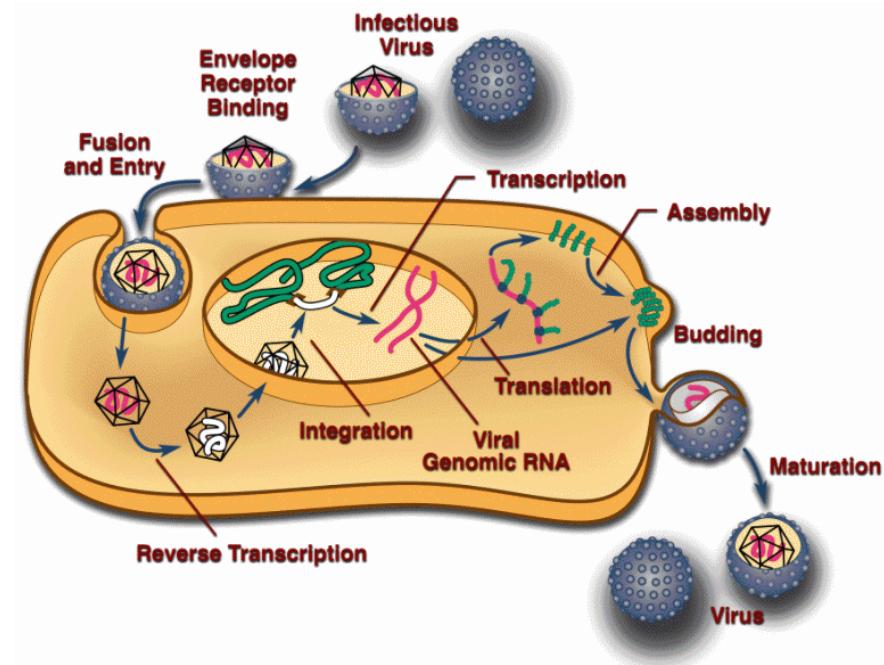
Alberts MBoC 5e

PART
1

GENOMES & TRANSCRIPTOMES

CENTRAL DOGMAS

- 1960s-70s: work on retroviruses changed the “central dogma”
 - Temin (Rous sarcoma virus (RSV))
 - Baltimore (Rauscher murine leukemia virus (R-MLV))
- Retroviruses (e.g. HIV) carry RNA as their genetic information
- When they invade their host cell they convert their RNA into a DNA copy using reverse transcriptase



<https://www.extremetech.com/extreme/247618-ancient-dna-showcases-war-viruses-hominid-ancestors>

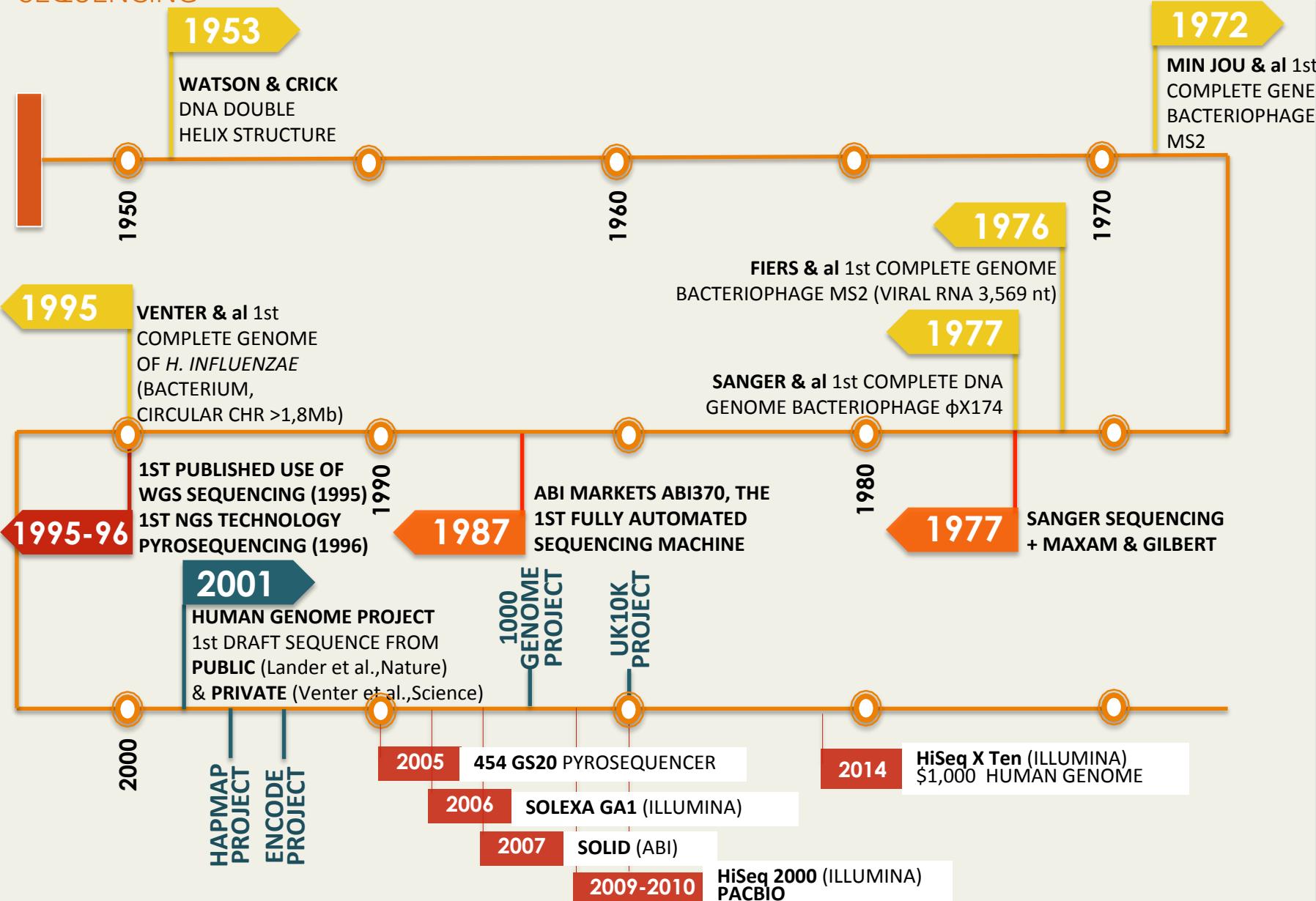
- Thus the central dogma has been modified:

$$\text{DNA} \leftrightarrow \text{RNA} \rightarrow \text{Protein}$$

- This has helped to explain an important paradox in the evolution of life.

GENOMES & TRANSCRIPTOMES

SEQUENCING



GENOMES & TRANSCRIPTOMES

HGP: THE METHODS



- The Human Genome Project (HGP)= a 13-years (1990-April 14, 2003) international effort to sequence the 3 billion "letters" of human DNA.
- \$ 3 billion project, led by the U.S. DoE and the NIH.
- International Human Genome Sequencing Consortium (IHGSC)= group of publicly funded researchers
- 20 main sequencing centers in US, UK, Japan, France, China, Germany

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

<http://www.genome.gov/sequencingcosts/>
www.sanger.ac.uk

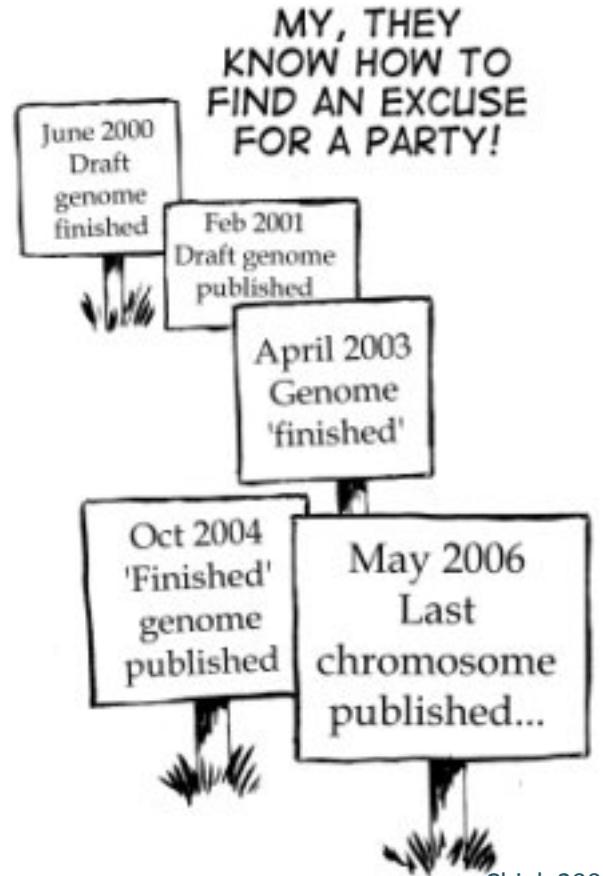
NGS
FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

HGP: THE METHODS

- 2 groups competing for sequencing:
 - Public
 - Private (Celera Genomics)
- Opposing philosophies :
 - Public : HGP Bermuda Agreement (1996)
→ all information from the project would be made freely available to all within 24h.
 - Private
→ access restricted to paying customers !
- In February 2001, drafts of the human genome sequence were published simultaneously by both public-private groups in separate articles (Lander et al (IHGSC),, Feb 2001 Nature) (Venter et al., Feb 2001 Science)



Chial, 2008

<http://www.genome.gov/sequencingcosts/>
<http://www.yourgenome.org/>
www.sanger.ac.uk

NGS

FATMA GUERFALI

PART
1

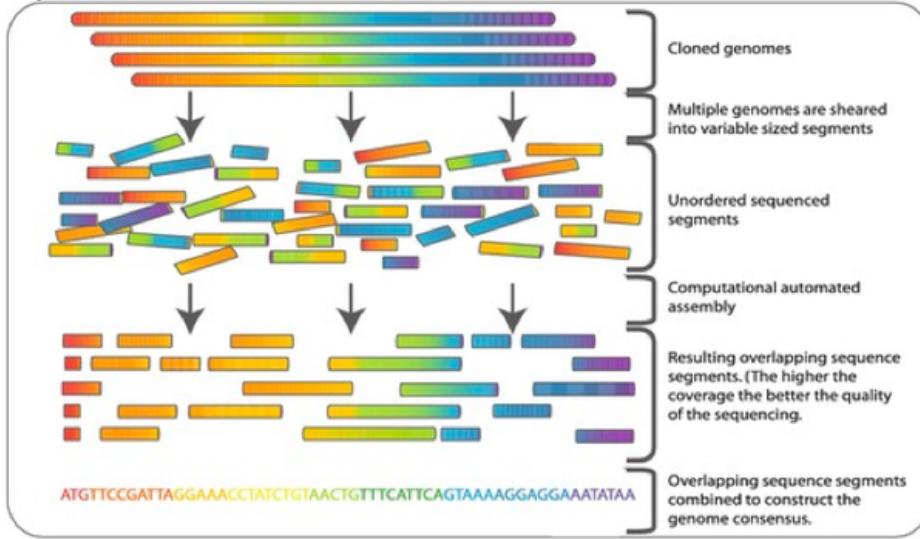
OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA



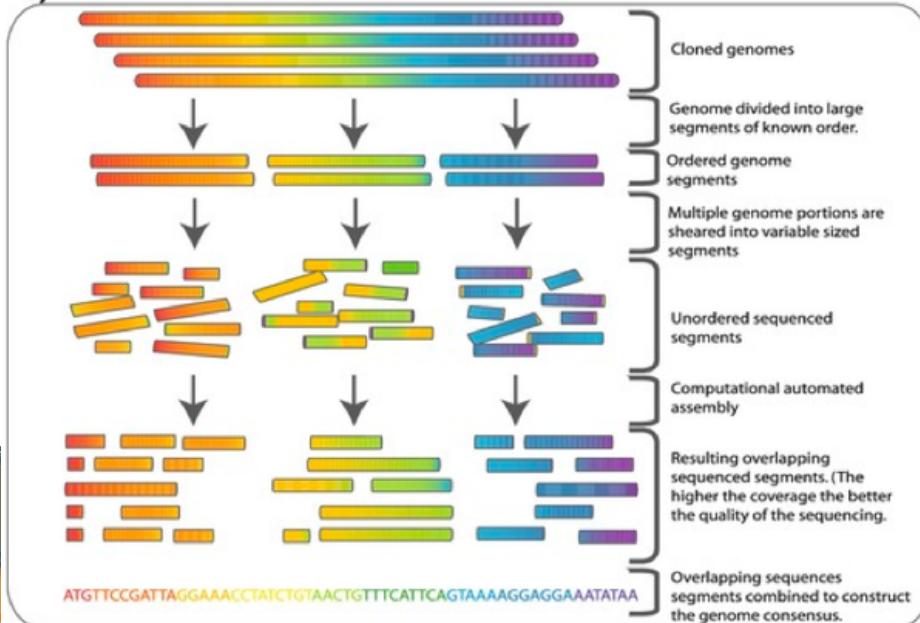
GENOMES & TRANSCRIPTOMES

HGP: THE METHODS

a)



b)



WHOLE GENOME SHOTGUN

(Celera genomics, Private)

- Genome sheared randomly into small fragments (appropriately sized for sequencing)
- Reassembly

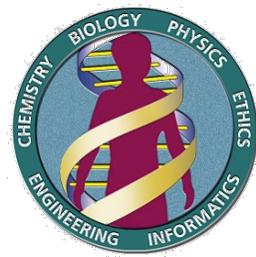
HIERARCHICAL GENOME SHOTGUN (TOP-DOWN SEQUENCING)

- low-resolution physical map of the genome is made prior to sequencing.
- Genome fragmented into large segments.
- Order of these segments is deduced
- They are further sheared into fragments appropriately sized for sequencing.
- Assembly relying on the physical map
- Finishing phase : filling in gaps and resolving DNA sequences in ambiguous areas



GENOMES & TRANSCRIPTOMES

HGP: THE METHODS



Variable degrees of completion of published genomes

- **Draft Sequencing**

- high-throughput or shotgun phase (whole genome or clone-based approach)
- Assembly using specific algorithms (whole-genome or single-clone assembly)
→ lower accuracy than finished sequence; some segments are missing or in the wrong order or orientation.

- **Finishing**

- Accuracy in bases identification + Quality Check + few if any gaps.
- Contiguous segments of sequence are ordered and linked to one another
- No ambiguities or discrepancies about segments order and orientation

- **Complete Genome**

A Genome represented by a single contiguous sequence with no ambiguities

→ The sequences available are *finished to a certain high quality*.

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Mardis et al., 2002)
<http://www.informatics.jax.org/>
NGS
FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE



- The HGP project required that all human genome sequence information be **freely and publicly available**. The existing **DNA sequences** have been stored in databases available to anyone willing to exploit and analyze them.
- Dedicated databases house various data for model organisms such as sequences of known and hypothetical genes and proteins (**GenBank**, **NCBI**). Other databases (**Ensembl** <http://www.ensembl.org>) present additional data and annotation as well as powerful tools for visualizing and searching it.
- Community efforts for non-model organisms like Eukaryotic Pathogens : **EuPathDB** (<http://eupathdb.org/eupathdb/>).

Computer programs have been developed to analyze and interpret the data.

PART
1

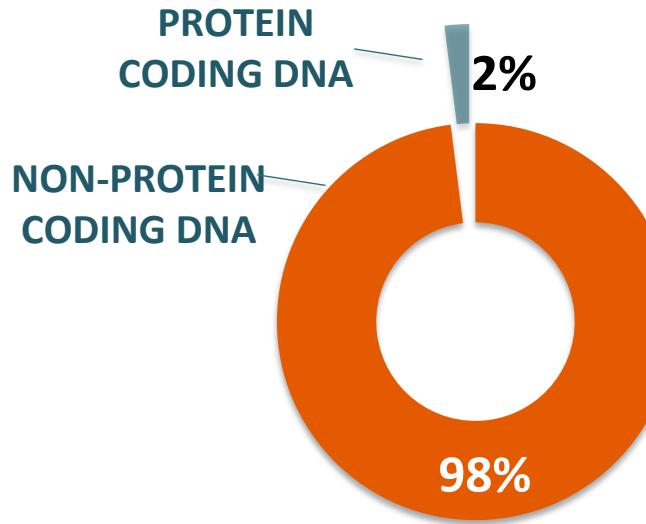
OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

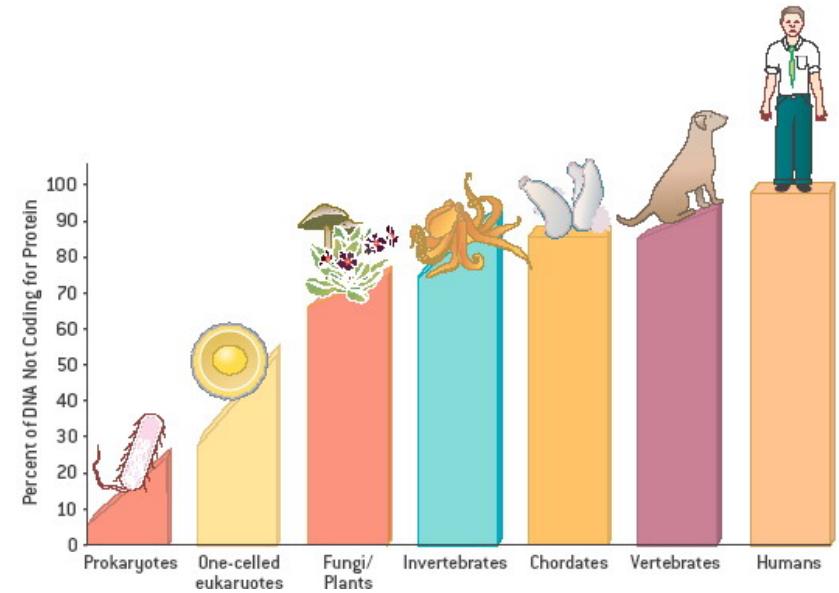
HGP: THE HERITAGE



$\approx 20,000$

The human genome contains only about 20,000 protein-coding genes : sequence alone is not enough to explain the whole complexity !

The proportion of non-coding DNA increases with organism complexity



(Mattick, 2011)

<http://www.yourgenome.org/stories/how-is-the-completed-human-genome-sequence-being-used>

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

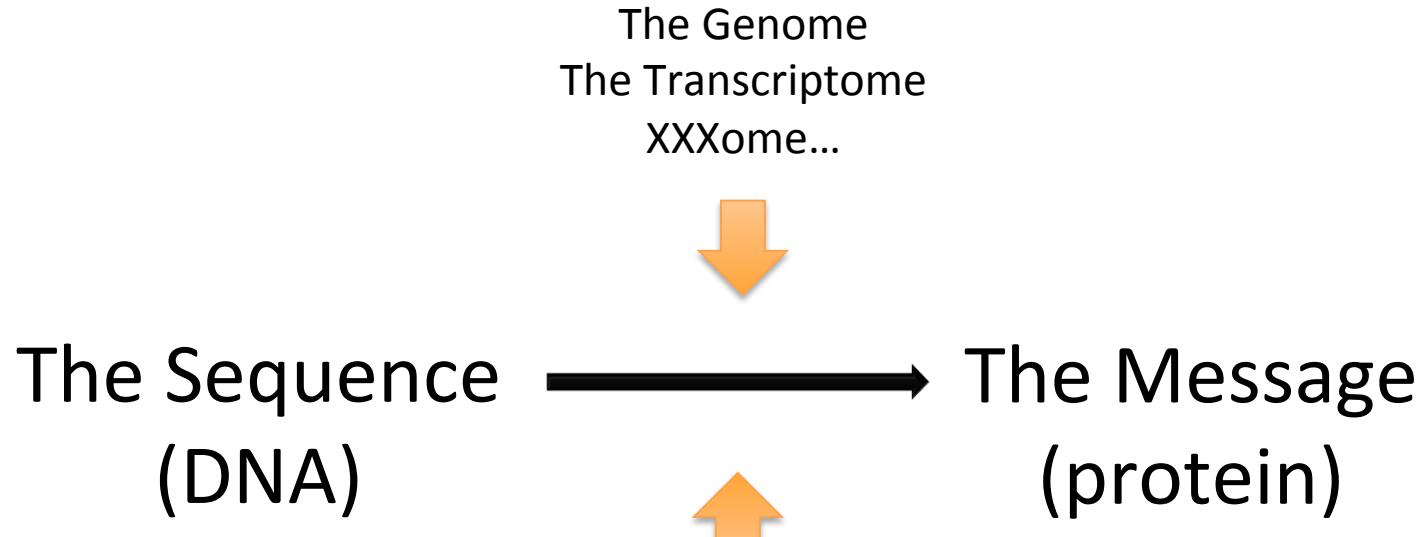
NGS

FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE



PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

...

NGS

FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

- = It appears that the *genetic programming* of humans and other complex organisms has been *misunderstood for the past 50 years*, due to the assumption that most genetic information is transacted by proteins.
- The majority of these sequences are dynamically transcribed, mainly into non-protein-coding RNAs, with tens if not hundreds of thousands that show specific expression patterns and subcellular locations, as well as many classes of small regulatory RNAs.

Only infrastructural RNAs where known ! (rRNAs, tRNAs,...)

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Mattick, 2011)

NGS

FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

- The emerging evidence indicates that these RNAs control the epigenetic states that underpin development, and that many are dysregulated in cancer and other complex diseases.
- Moreover it appears that animals, particularly primates, have evolved plasticity in these RNA regulatory systems, especially in the brain.
- Thus, it appears that what was dismissed as 'junk' because it was not understood holds the key to understanding human evolution, development, and cognition.

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Mattick, 2011)

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

- “The Encyclopedia of DNA Elements (ENCODE) Project aims to provide a more biologically informative representation of the human genome by using high-throughput methods to identify and catalogue the functional elements encoded.”

● ENCODE = Encyclopedia of DNA Elements

- ENCODE is a public research consortium launched by the US National Human Genome Research Institute (NHGRI) in September 2003.
- The aim was to find the functional elements of the human genome.
- All data that have been obtained in this project have been made available in public databases / repositories.

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Birney et al., 2007)

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

Non-Coding: the ENCODE Project

- Among the major key findings:

- ≈ 20% of ncDNA in the human genome is functional (regulation of gene expression for coding genes).
- ≈ 60% is transcribed without any known function yet.
- The expression of each coding gene is under the control of several regulation sites that can be located near or far from the gene = unexpected complexity !

→ Much more complex transcriptional “landscape” than what was previously expected !!!

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

NGS
FATMA GUERFALI



GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

Coding: the GENCODE Project

- The human genome has been the focus of intensive manual annotation: GENCODE

Harrow, J. et al. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res.* 22, 1760–1774 (2012)

The GENCODE Consortium aims to identify all gene features in the human genome using a combination of:

- computational analysis
- manual annotation
- experimental validation.

PART
1

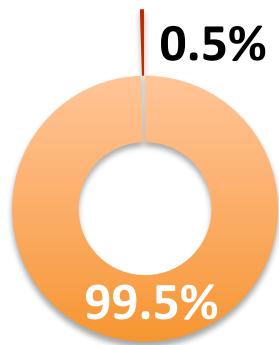
OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

NGS
FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

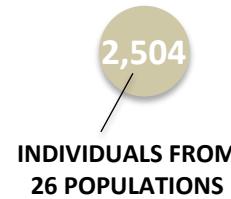
HGP: THE HERITAGE

HAPMAP



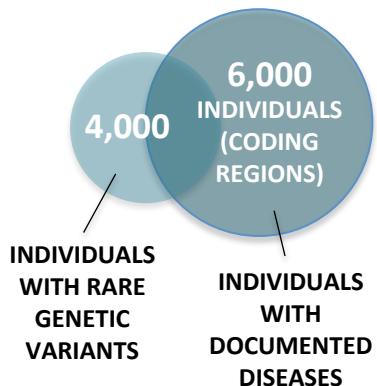
CATALOG OF COMMON GENETIC VARIATION AMONG HUMANS TO IDENTIFY DISEASE-RELATED GENES.
→ THE DNA SEQUENCE OF ANY TWO PEOPLE IS 99.5% IDENTICAL.

1,000 GENOMES



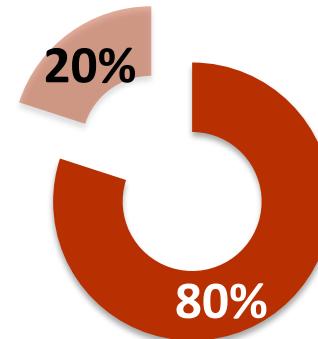
RESOURCE OF HUMAN GENETIC VARIATION

UK10K



LINK BETWEEN GENETIC VARIANTS AND RARE DISEASES

ENCODE



80 % OF THE HUMAN GENOME SEQUENCE IS CLOSE TO REGIONS THAT CONTROL BIOLOGICAL FUNCTION.
→ “JUNK DNA” IN THE HUMAN GENOME, IS ACTUALLY FUNCTIONAL.

(Mattick, 2011)

<http://www.yourgenome.org/stories/how-is-the-completed-human-genome-sequence-being-used>

NGS

FATMA GUERFALI

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA



GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE



dbSNP: Database for Short Genetic Variations

Catalog of nucleotide changes for human and other model organisms

<https://www.ncbi.nlm.nih.gov/snp/>

National Center for Biotechnology Information • National Library of Medicine • National Institutes of Health • Department of Health and Human Services

Scope and Access

The NCBI Short Genetic Variations database, commonly known as dbSNP, catalogs short variations in nucleotide sequences from a wide range of organisms. These variations include single nucleotide variations, short nucleotide insertions and deletions, short tandem repeats and microsatellites. Short Genetic Variations may be common, thus representing true polymorphisms, or they may be rare. Some rare human entries have additional information associated with them, including disease associations, genotype information and allele origin, as some variations are somatic rather than germline events.



PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

ftp://ftp.ncbi.nih.gov/pub/factsheets/Factsheet_SNP.pdf

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

D dbSNP

A SNP **B** human[orgn] AND HFE[gene]

C Display Settings: Summary, 20 per page, Sorted by SNP_ID

E Help

F **G** graphical presentation under the context of genome or mRNA sequences

H

I Find related data

J Filters: Manage Filters

K Send to:

L Database: Select

M Find items

N Search details
"Homo sapiens" [Organism]
AND HFE [gene]

O Search See more...

P Recent activity Turn Off Clear

Q human[orgn] AND HFE[gene] (841)

R Format Items per page Sort by

<input type="radio"/> Summary	<input type="radio"/> 5	<input type="radio"/> Default order
<input type="radio"/> Graphic Summary	<input type="radio"/> 10	<input type="radio"/> Organism
<input type="radio"/> FASTA	<input type="radio"/> 20	<input checked="" type="radio"/> SNP_ID
<input checked="" type="radio"/> FlatFile	<input type="radio"/> 50	<input type="radio"/> Success Rate
<input type="radio"/> Chromosome Report	<input type="radio"/> 100	<input type="radio"/> Heterozygosity
<input type="radio"/> Old Summary	<input type="radio"/> 200	<input type="radio"/> Chromosome Base Position
<input type="radio"/> dbSNP Batch Report		

S Apply

T Page 1 of 43

U Next >

V Last >

W << First

X < Prev

Y > Next

Z >> Last

aa rs1799945 [Homo sapiens]

bb 1. TGACCAGCTTCTGTGTTCTATGAT [C/G] ATGAGAGTCGCCGTGGAGCCCCG

cc Chromosome: 6:26090951

dd Gene: HFE (GeneView)

ee Functional Consequence: intron variant, missense, nc transcript variant

ff Allele Origin: G(germline)/C(germline)

gg Clinical significance: other

hh Validated: by 1000G, by cluster, by frequency, by hapmap

ii Global MAF: G=0.0731/366

jj HGVS: NC_000006.11:g.26091179C>G, NC_000006.12:g.26090951C>G, NG_008720.2:g.8671C>G, NM_00410.3:c.187C>G, NM_001300749.1:c.187C>G, NM_139003.2:c.187C>G, NM_139004.2:c.187C>G, NM_139006.2:c.187C>G, NM_139007.2:c.77-363C>G, NM_139008.2:c.77-363C>G, NM_139009.2:c.118C>G, NM_139010.2:c.77-1734C>G, NM_139011.2:c.77-2168C>G, NP_000401.1:p.His63Asp, NP_001287678.1:p.His63Asp, NP_620572.1:p.His63Asp, NP_620573.1:p.His63Asp, NP_620575.1:p.His63Asp, NP_620578.1:p.His40Asp, XM_005249040.1:c.187C>G, XM_011514543.1:c.187C>G, XM_011514544.1:c.187C>G, XP_005249097.1:p.His63Asp, XP_011512845.1:p.His63Asp, XP_011512846.1:p.His63Asp, XR_241893.1:n.309C>G, XR_241893.2:n.309C>G, XR_241894.1:n.434C>G

kk PubMed Varview Protein3D OMIM

ll rs1800562 [Homo sapiens]

mm 2. CCTGGGA[GAGATATACTG[A/G]CCAGGTGGAGCACCCAGGCCTGGAT

nn Chromosome: 6:26092913

oo Gene: HFE (GeneView)

pp Functional Consequence: intron variant, missense, nc transcript variant

qq Allele Origin: G(germline)/A(germline)

rr Clinical significance: Pathogenic

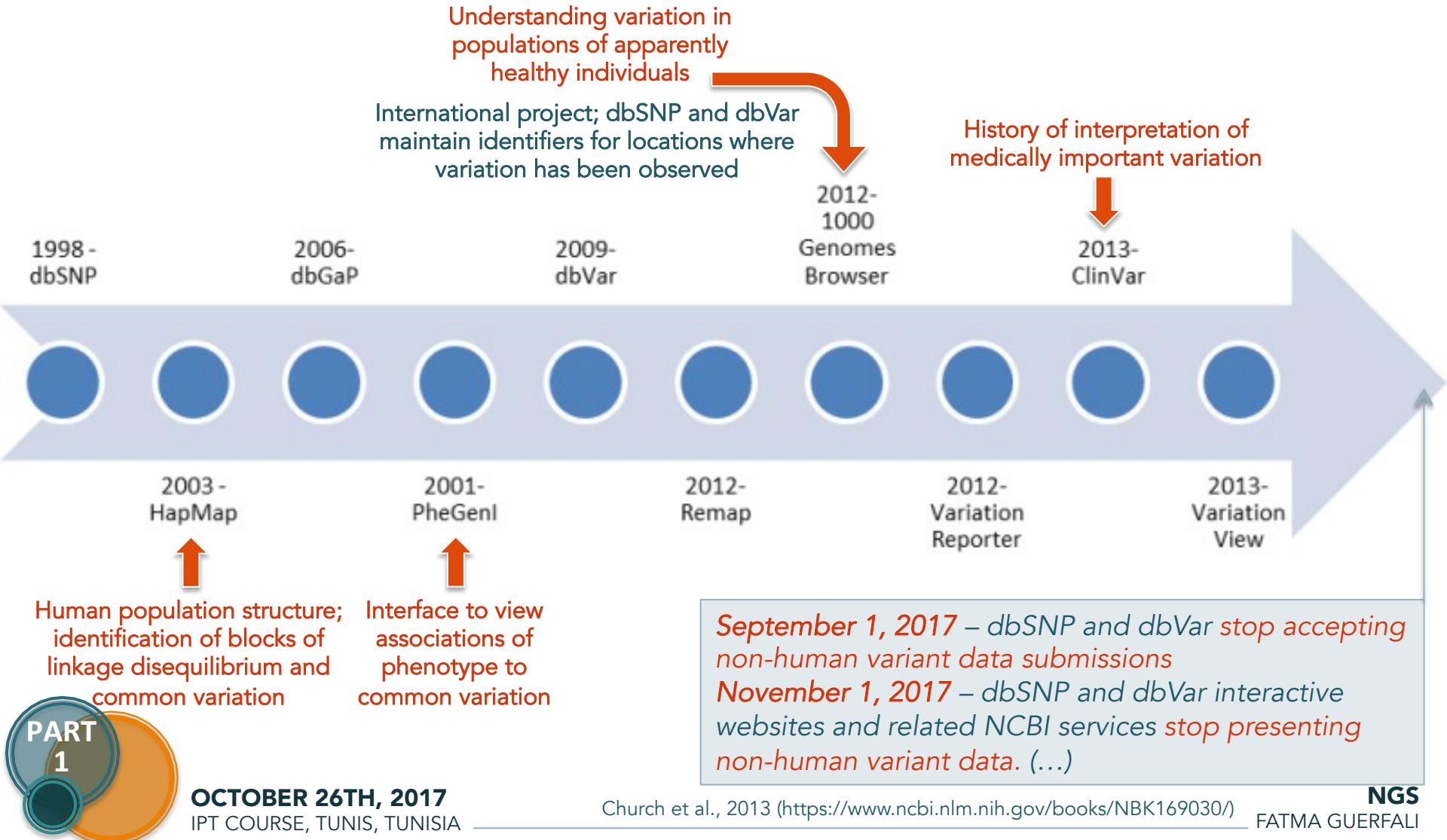
ss Validated: by 1000G, by cluster, by frequency, by hapmap

tt Global MAF: A=0.0126/63

uu HGVS: NC_000006.11:g.26093141G>A, NC_000006.11:g.26093141G>A, NG_008720.2:g.10633G>A, NM_00410.3:c.845G>A, NM_001300749.1:c.845G>A,

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

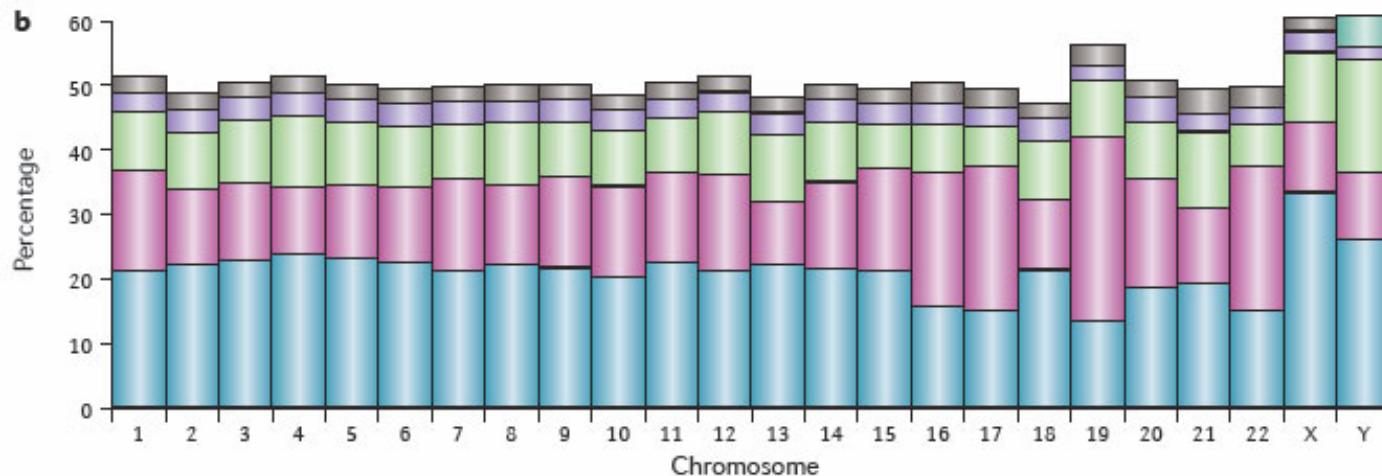


GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

- About 50% of the Human genome is composed of repeats

a Repeat class	Repeat type	Number (hg19)	Cvg	Length (bp)
Minisatellite, microsatellite or satellite	Tandem	426,918	3%	2–100
SINE	Interspersed	1,797,575	15%	100–300
DNA transposon	Interspersed	463,776	3%	200–2,000
LTR retrotransposon	Interspersed	718,125	9%	200–5,000
LINE	Interspersed	1,506,845	21%	500–8,000
rDNA (16S, 18S, 5.8S and 28S)	Tandem	698	0.01%	2,000–43,000
Segmental duplications and other classes	Tandem or interspersed	2,270	0.20%	1,000–100,000



PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

Treangen & Salzberg, 2012

Nature Reviews | Genetics

NGS

FATMA GUERFALI

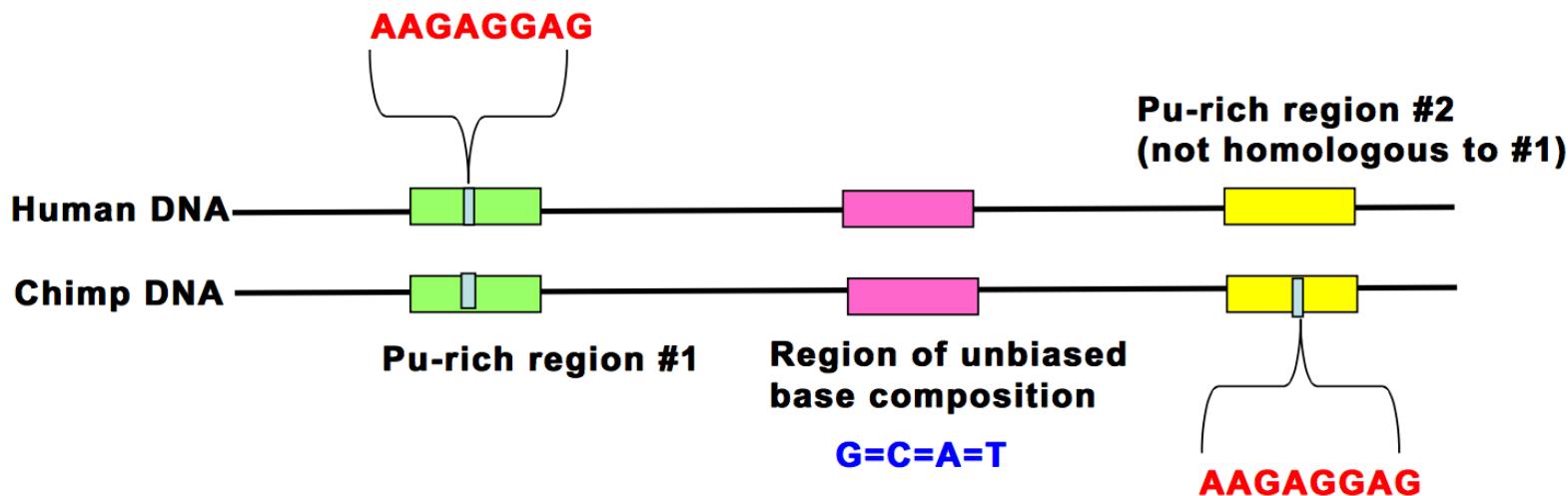
GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

- Low complexity regions

Molecular structural region of biased composition—e.g.:

- homopolymeric runs
- short-period repeats (simple tandem repeats, polypurine, AT-rich...)
- subtle overrepresentation of one or a few residues
- ...



When a sequence is of low complexity (and/or short length), a high % identity with another sequence does not reflect necessarily a shared evolutionary origin.

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

Lewis, 2015 (<http://slideplayer.com/slide/4577926/>)

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

General Overview

- Only 2% of the Human genome is coding
- Gaps in sequences remain
- Annotated and 1000s of previously unannotated RNAs, including Families (tens to hundreds) of Non-coding RNA : lncRNAs, snRNAs (microRNAs, snoRNAs, snRNAs), pseudogenes
- ¾ of the human genome is capable of being transcribed
- Pseudogenes can be transcribed
- Over 100,000 splice junctions discovered that were not incorporated into transcript models (strand-specific RNA-seq).

"These observations, taken together, prompt a redefinition of the concept of a gene."

(Adapted from Djebali, S. et al. 2012)

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

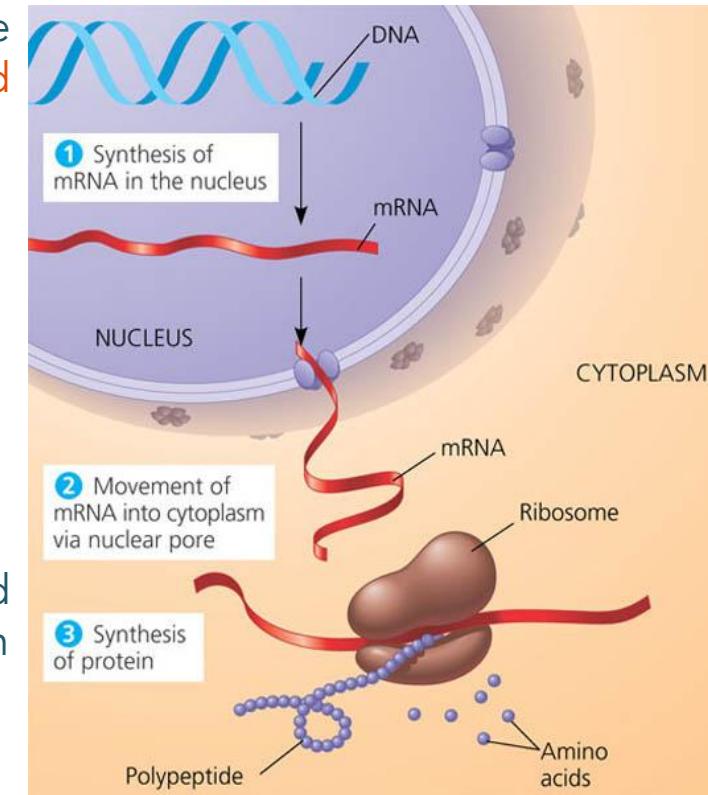
RNA is no longer considered to be a “bridge” in the transfer of biological information between DNA and proteins

Transcriptome consist of

- Coding RNA (mRNA (2–4%)
- Previously known non-coding RNAs
 - rRNA (80–90%)
 - tRNA (5–15%)
- Small fraction of intragenic (i.e., intronic) and intergenic non-coding RNA (1%, ncRNA) with **undefined regulatory functions**.

NB: Intragenic & Intergenic sequences

- Are enriched in repetitive elements
- have long been considered genetically inert



(Costa et al., 2010)

PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

Transcriptome = The complete collection of transcribed elements of the genome (Affymetrix, 2004)

NGS

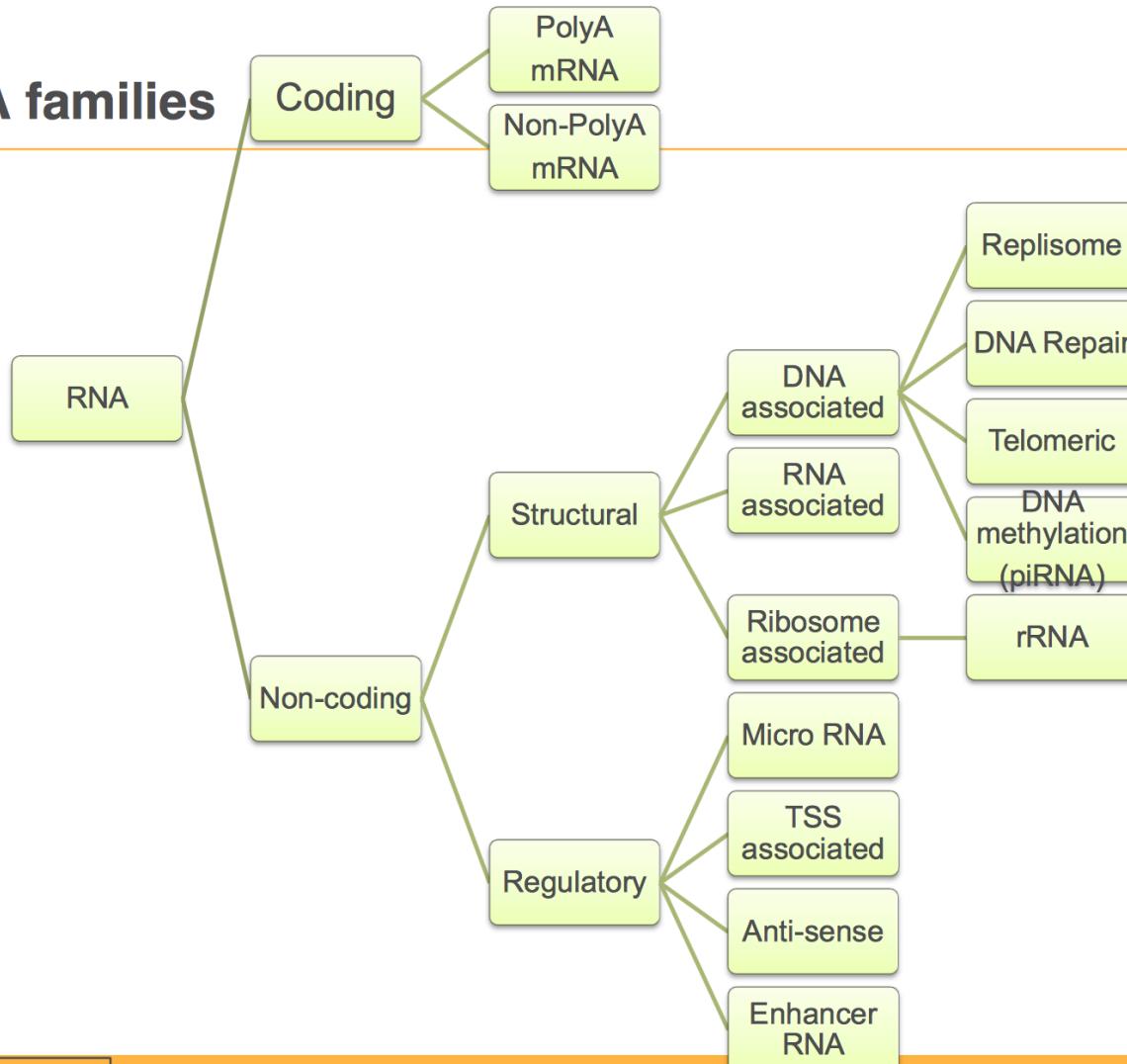
FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

An overview of the eukaryotic transcriptome through examples of its products

RNA families



PART
1

IPT COURSE, TUNIS, TUNISIA

Abizar Lakdawalla, PhD (illumina.com)

illumina®

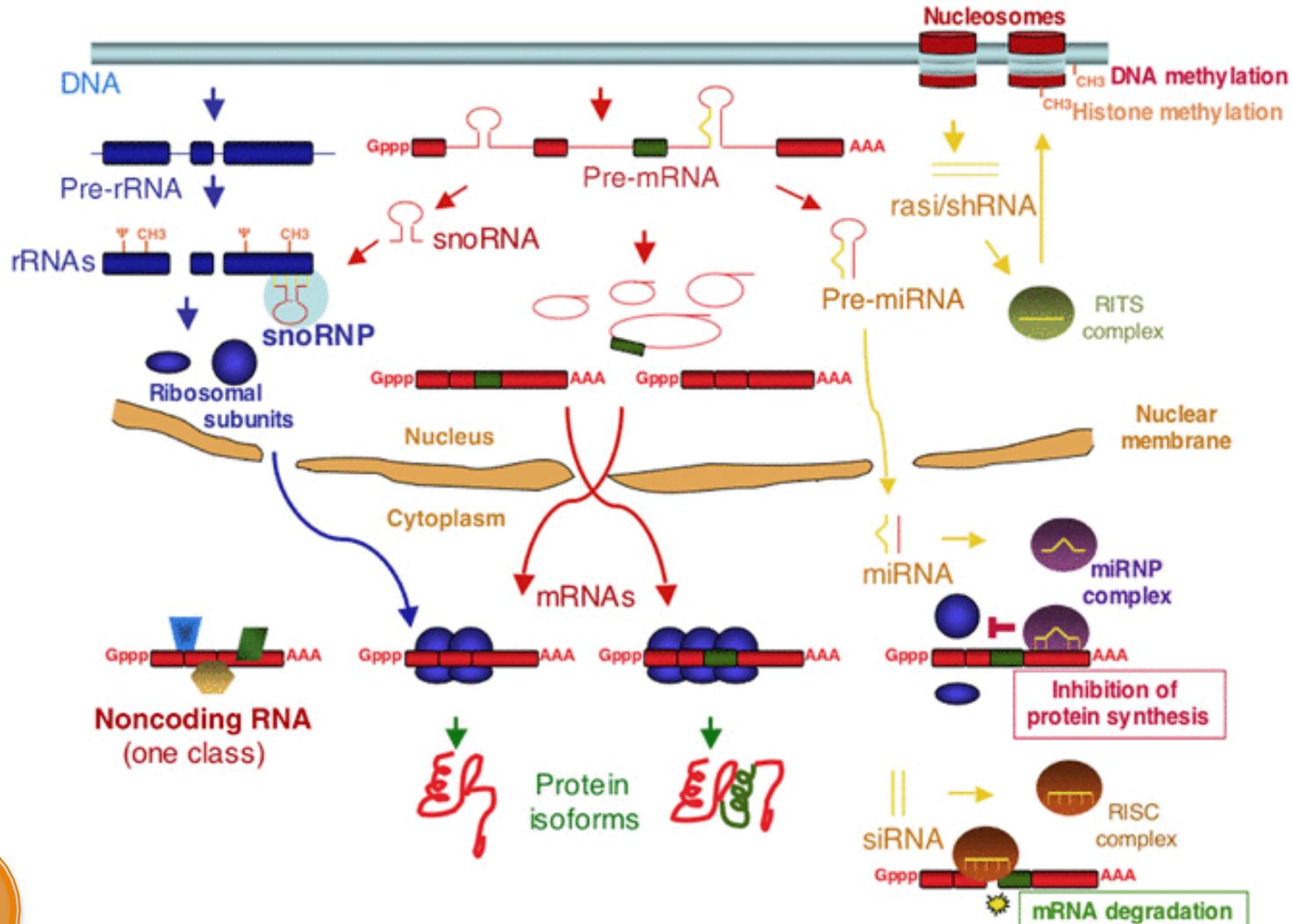
NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

An overview of the eukaryotic transcriptome through examples of its products



PART
1

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

An overview of the eukaryotic transcriptome through examples of its products

REVIEW

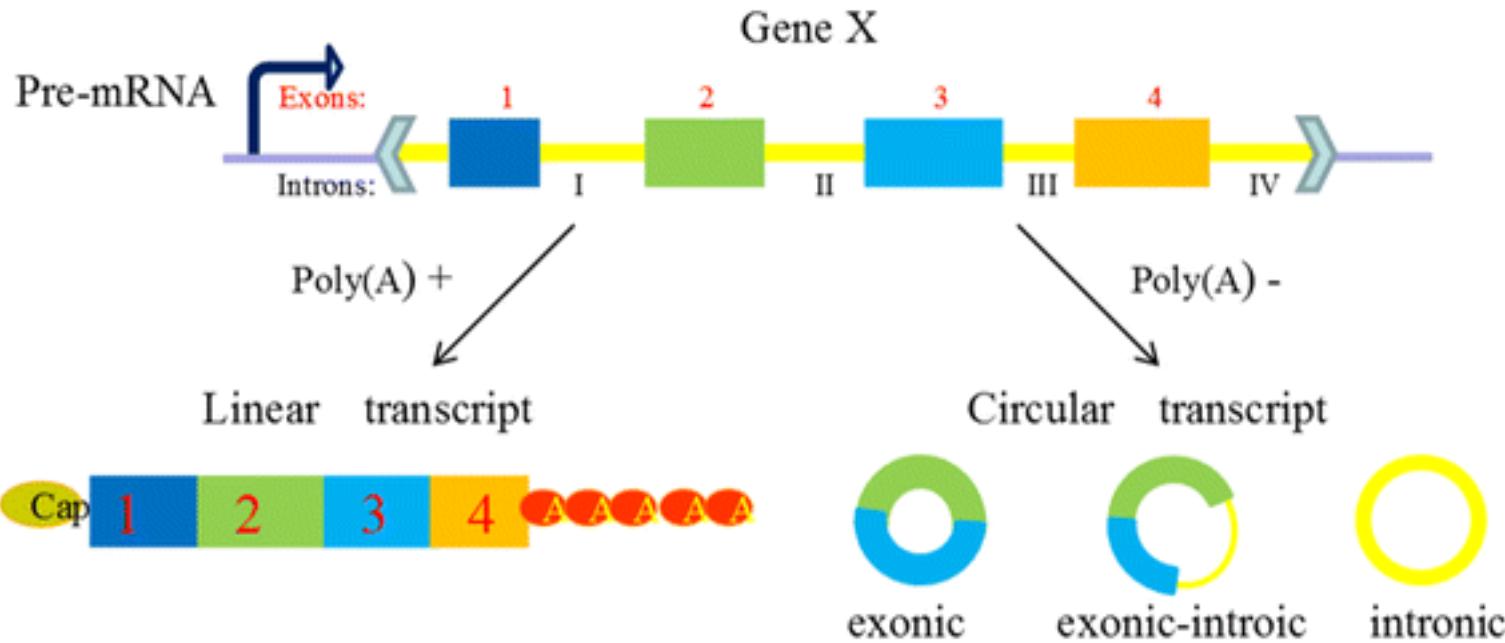
OPEN ACCESS

Circles reshaping the RNA world: from waste to treasure

Jing Liu, Tian Liu, Xiaman Wang and Aili He 

Molecular Cancer 2017 16:58 | DOI: 10.1186/s12943-017-0630-y | © The Author(s). 2017

Received: 25 October 2016 | Accepted: 2 March 2017 | Published: 9 March 2017



PART
1

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Liu et al., 2017)

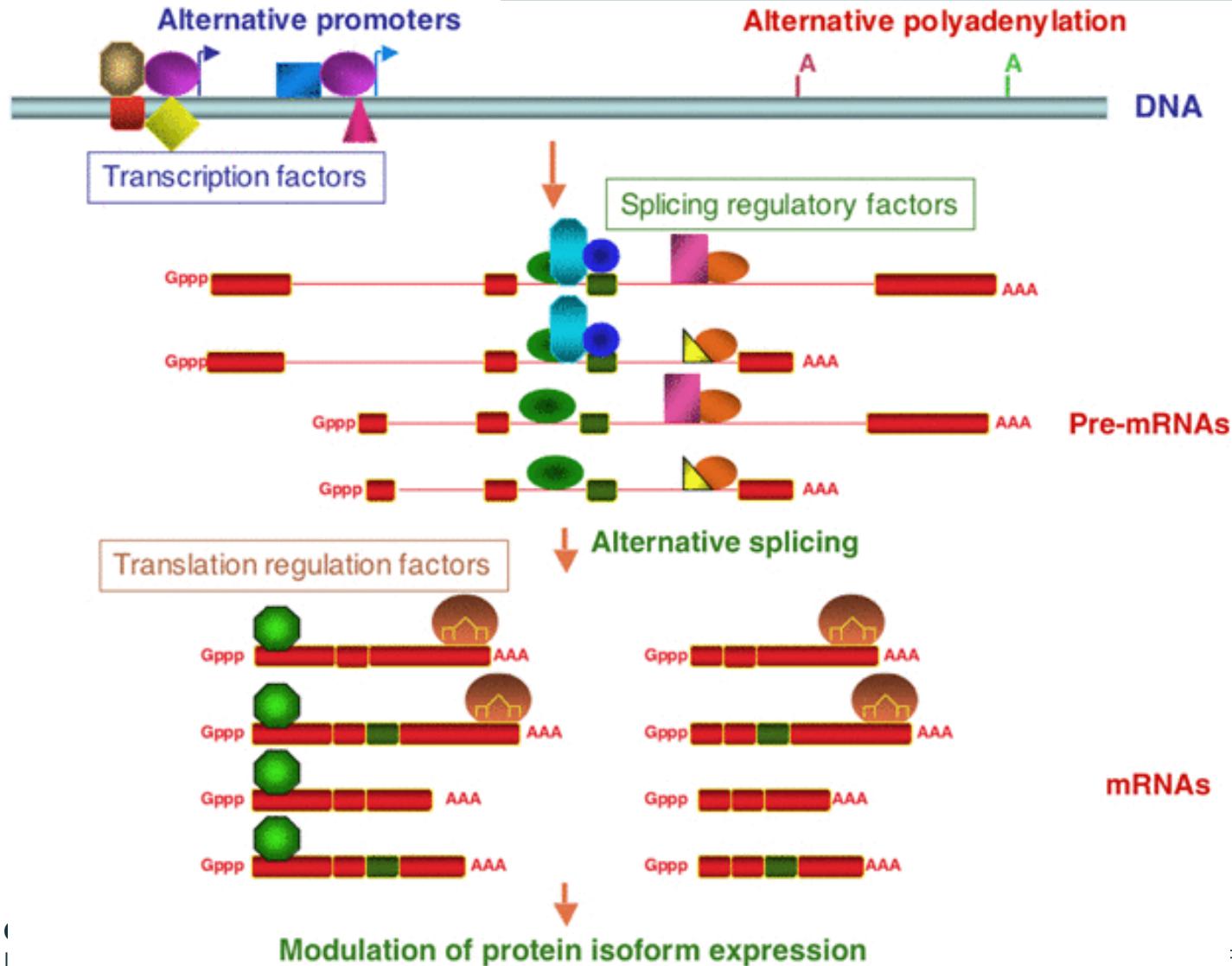
NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

Transcript diversity and combinatorial regulation of gene expression



PART
1

NGS

FATMA GUERFALI

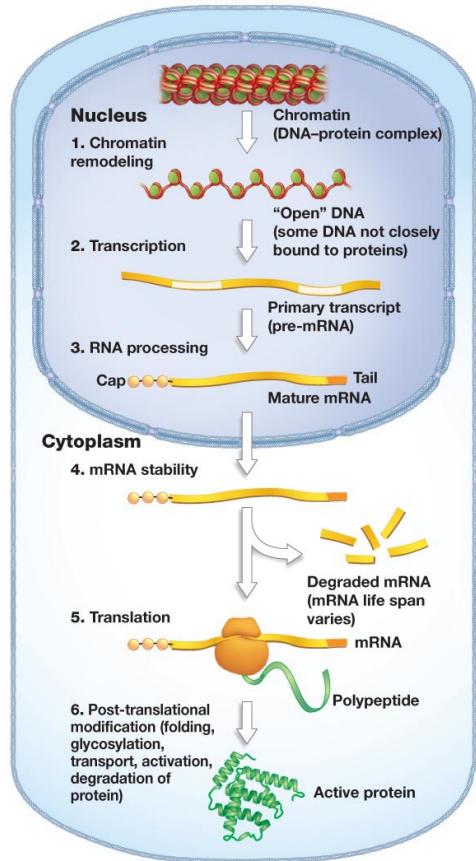
(Soares & Valcárcel, 2006)

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

The regulation processes are complex, with several coding levels

Some methods used by eukaryotes for gene regulation:

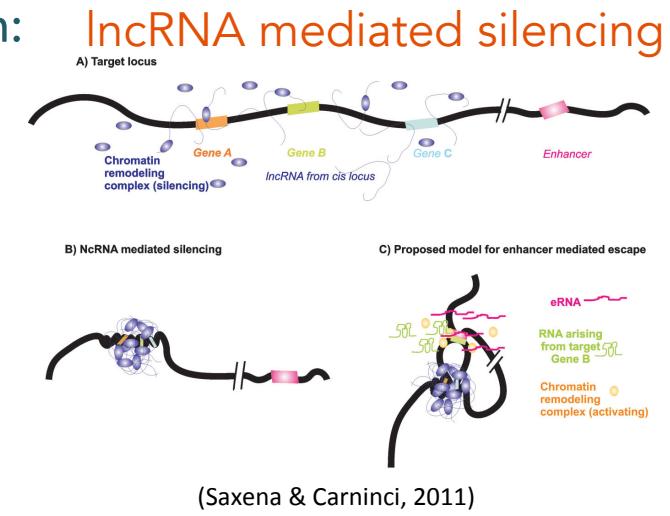


<http://www.uic.edu>
OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

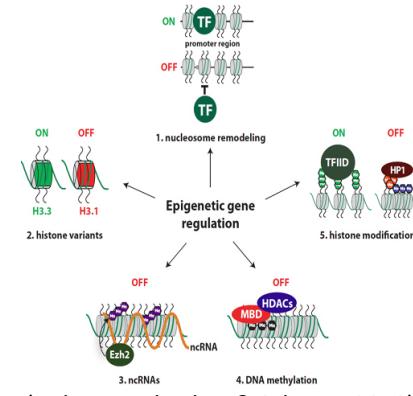
Acetylation
Ubiquitination
Methylation
Phosphorylation
Degradation
(...)

+

Response to stimuli
Transcription initiation



Epigenetic regulation



(Hahn, Dambacher & Schotta, 2010)

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

TAKE-HOME MESSAGE (OMICS : EXAMPLES)

With the increasing amount of data generated, people have been talking more and more about « omics »

Genomics	Studies the genomes of organisms.
Transcriptomics	Studies the transcriptome : it is the set of all RNA molecules, including mRNA, rRNA, tRNA, and other non-coding RNA produced in one or a population of cells.
Proteomics	Studies the proteome : it is the entire set of proteins, including the modifications made to a particular set of proteins, produced by an organism or system. Large-scale study of proteins, particularly their structures and functions. Mass spectrometry techniques are used
Metabolomics	Scientific study of chemical processes involving metabolites. It is a "systematic study of the unique chemical fingerprints that specific cellular processes leave behind", the study of their small-molecule metabolite profiles
Pharmacogenomics	investigates the effect of the sum of variations within the genome on drugs
Proteogenomics	An emerging field of biological research at the intersection of proteomics and genomics. Proteomics data used for gene annotations.
Immunoproteomics	study of large sets of proteins (proteomics) involved in the immune response
Functional genomics	Gene and protein functions and interactions (uses microarray kind of techniques)

Part 1 Genomes and Transcriptomes : From central dogmas to ongoing discoveries

(Central dogma of molecular Biology, Human Genome Project, GENCODE, ENCODE...)

Part 2 Advances in Sequencing Technologies

(Examples of technologies developed for Closed/Open systems for gene expression analysis, Next-Generation Sequencing platforms and technologies ...)

Part 3 Overview of NGS (DNA / RNA Seq) Protocols and related file formats

(Overview of protocols for Genomic and Transcriptomic analysis of standard samples using Next-Generation Sequencing and overview of the different formats generated at each step...)

- ▶ In order to cover all these complex events, researchers seek to study entire genomes and transcriptomes.

How is it possible nowadays to do so?

→ Sequencing using NGS



GENOMES & TRANSCRIPTOMES

SEQUENCING

What is sequencing about ?

Determining the precise order of the smallest units of a query. Examples:

DNA



RNA



Protein



ChIP



determine the
precise order
of nucleotides



determine the
precise order
of nucleotides



determine the
precise order of
amino acids



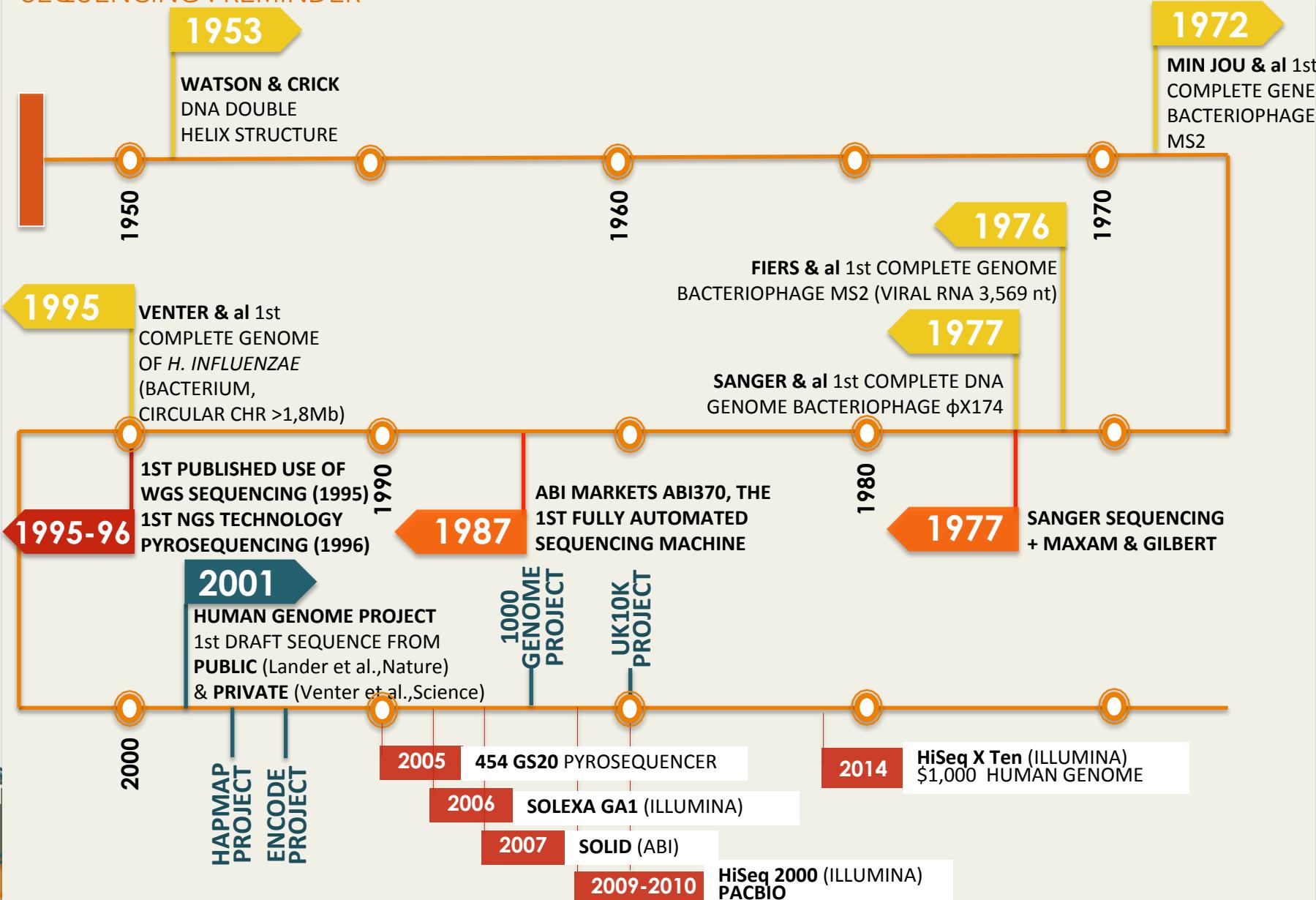
determine
genome-wide
the DNA binding
sites for DNA-
associated proteins
(TFs, ...)

PART
2

- ▶ Advances in Sequencing Technologies
- ▶ Next Generation Sequencing Platforms & Technologies

GENOMES & TRANSCRIPTOMES

SEQUENCING : REMINDER



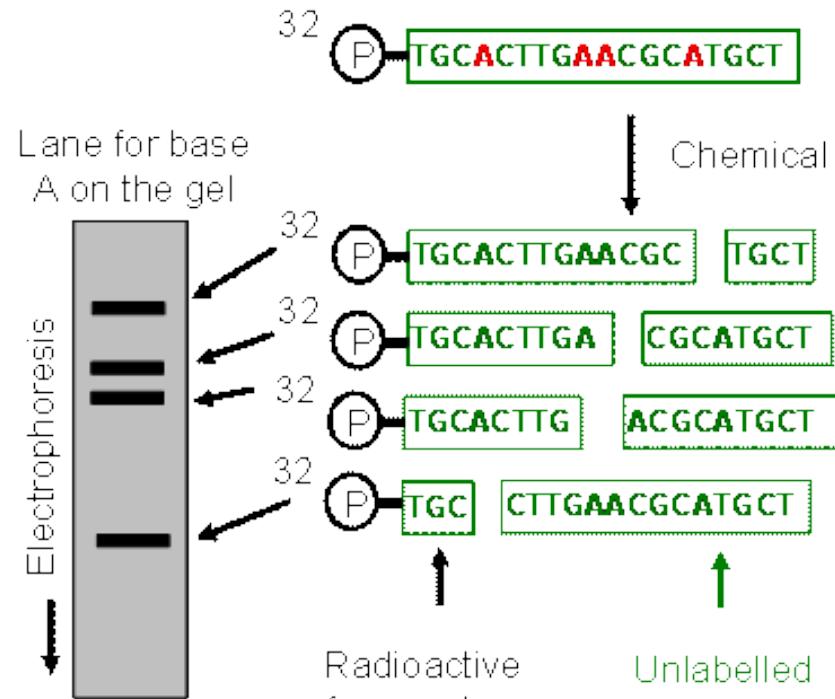
Maxam & Gilbert Sequencing (pre-NGS)

- The DNA fragment is labelled with ^{32}P at its 5' end.

- Then chemicals are used that break the DNA preferentially at each (or two) of the four nucleotide bases under conditions in which only one break per chain is made.

- After gel electrophoresis and autoradiography only the fragments possessing 5'-terminal ^{32}P -phosphate group show up on the gel.

- Time-consuming and requires handling of toxic chemicals.



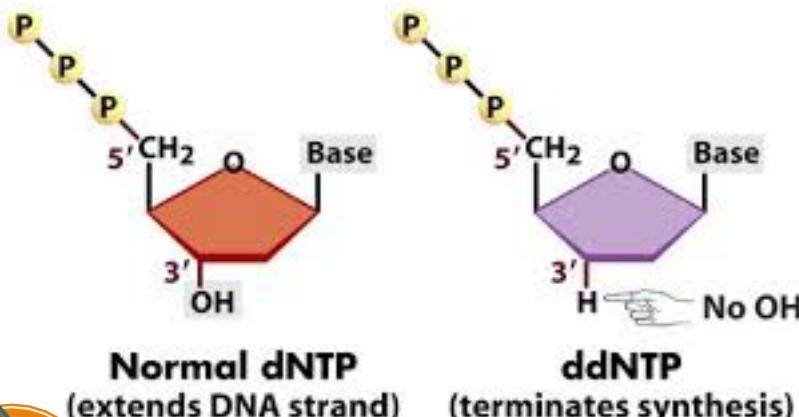
wiki.biomine.skelleftea.se

SEQ TECHNOLOGIES

ADVANCES MADE IN SEQUENCING APPROACHES AND TECHNOLOGIES

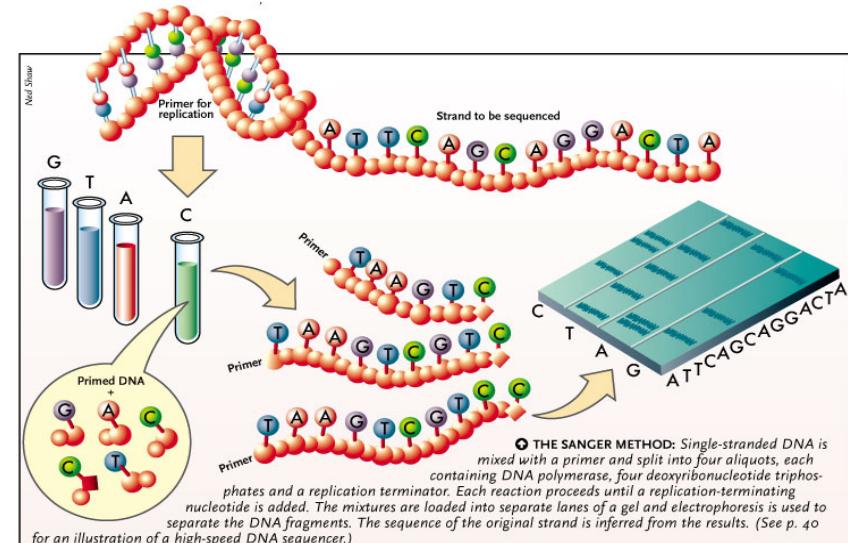
Sanger Sequencing (pre-NGS)

- Sanger sequencing is the method developed by Frederick Sanger in 1977. This method involves copying single-stranded DNA with chemically altered bases called dideoxynucleotides (ddNTPs).
- ddNTPs when incorporated at the 3' end of the growing chain, terminate the chain selectively at A, C, G, or T. The terminated chains are then resolved by capillary electrophoresis.
- <https://www.youtube.com/watch?v=FvHRio1yyhQ>



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA



NGS

FATMA GUERFALI
www.the-scientist.com

SEQ TECHNOLOGIES

EXAMPLE OF TECHNIQUES DEVELOPED TO STUDY GENE EXPRESSION

- The principle of Sanger sequencing revolutionized the approaches developed to study genes and gene expression
- High accuracy
- Remained the most widely used sequencing method for about 40 years, until recently supplanted by "Next Generation Sequencing" (NGS) methods
- NGS are nowadays widely used, especially for large-scale automated genome sequencing
- However, Sanger sequencing remains the method of choice for small scale sequencing (few genes approach) and validation of microarray or NGS results, or for obtaining long reads (> 500 nucleotides)

PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

https://en.wikipedia.org/wiki/Sanger_sequencing

NGS

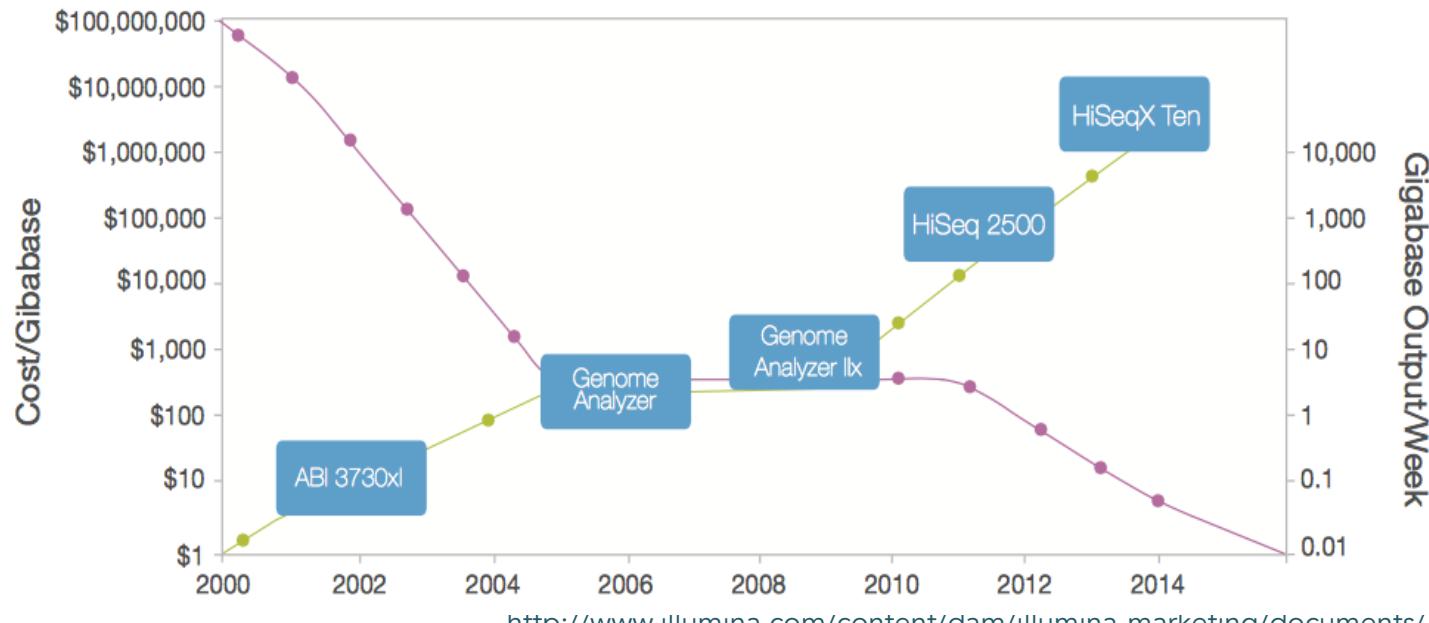
FATMA GUERFALI

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

- 1977: Sanger Sequencing
- 1987: First automated sequencing instrument based on capillary electrophoresis (AB370, 1987 ; AB3730xl, 1998) → primary instrument for HGP (public and private).
→ "First generation instruments" : 84Kb / run
- 2005: Genome Analyzer
- "Next generation instruments" start : 1Gb / run ----- 2014: 1.8Tb

Sequencing
Cost and Data
Output Since
2000



http://www.illumina.com/content/dam/illumina-marketing/documents/products/illumina_sequencing_introduction.pdf

PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

NGS

FATMA GUERFALI

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

- NGS

Sequence A

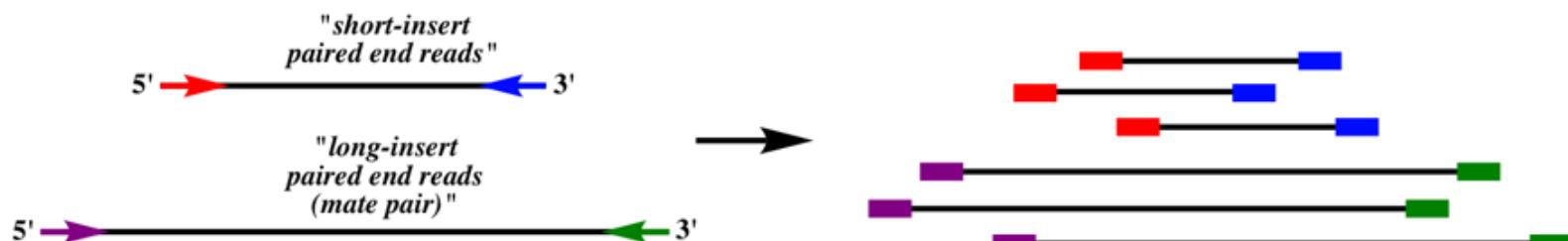
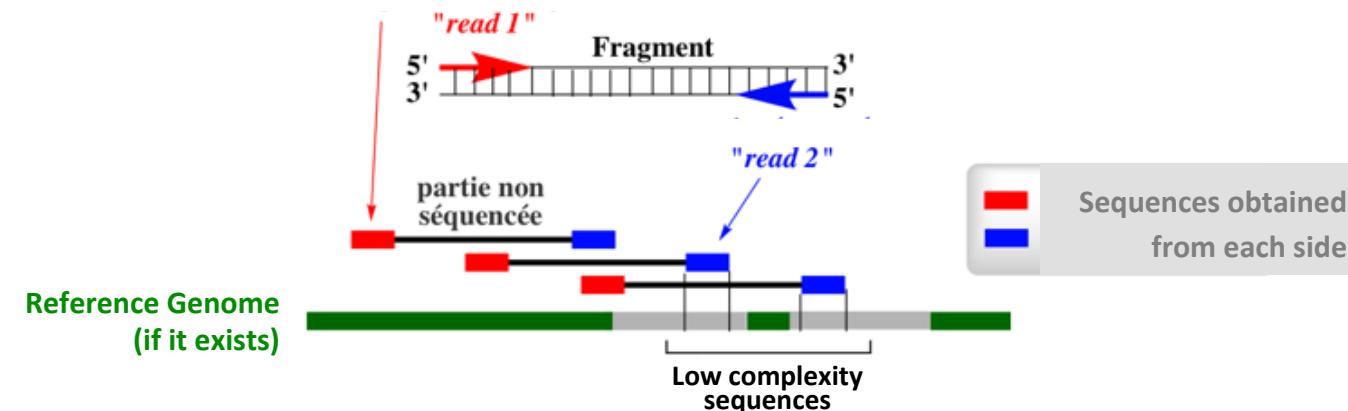
TTAGCGATGATTCTCGATGC|GGTTCCAATTGC

Sequence B

ATTCGGAATGCATC TTAGCGATGATTCTCGATGC

Contig

ATTCGGAATGCATCTTAGCGATGATTCTCGATGCGGTCCAATTGC



E. Jaspard (2014)

PART
2

Next Generation of Sequencing Technologies

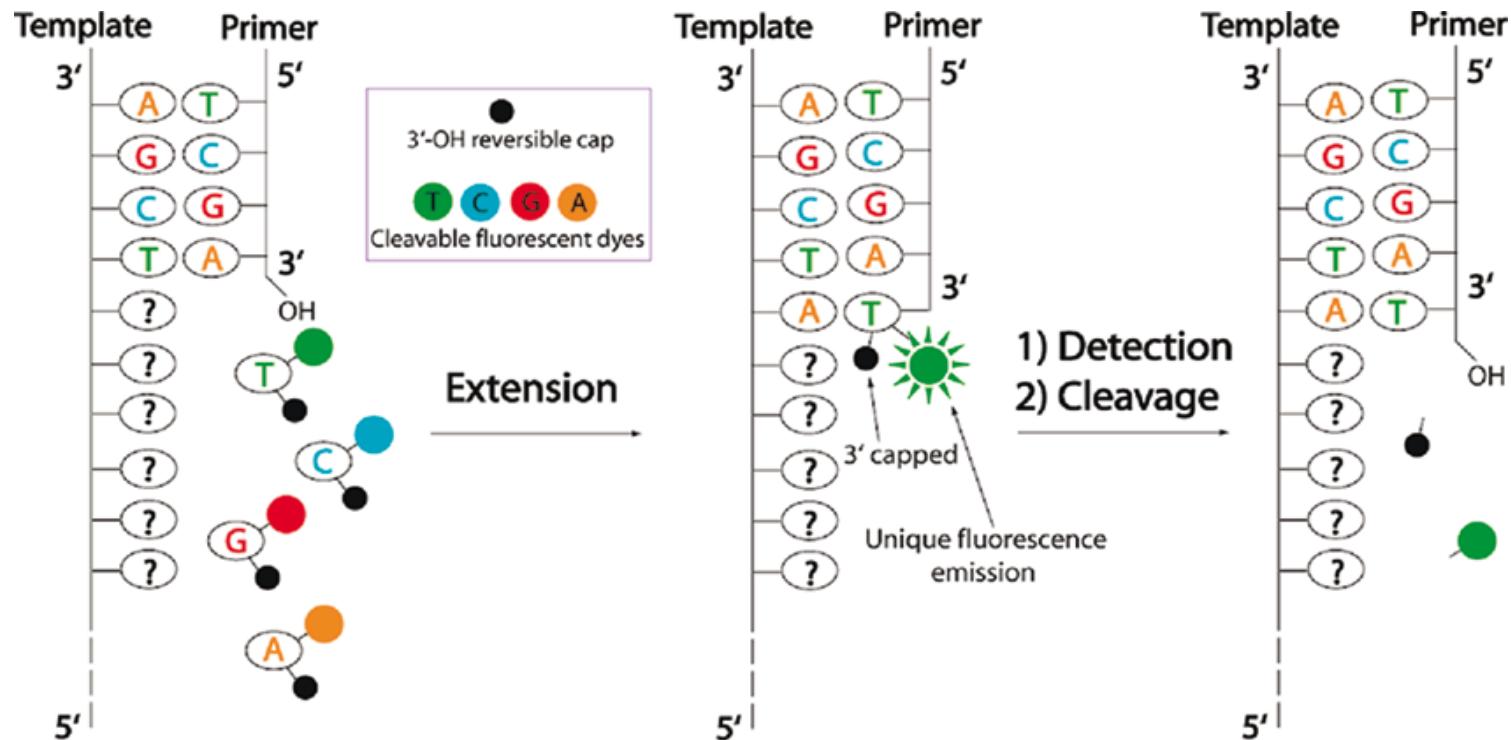
Technology	Company	Support	Chemistry
Massively Parallel Sequencing			
Solexa	Illumina	Bridge PCR on flowcell	Seq-By-Synthesis
454	Roche Applied Science	emPCR on beads	Pyrosequencing
SOLiD	AB / Life Technologies	emPCR on beads	Seq-By-Ligation
Ion Torrent	Life Technologies	emPCR on beads	Proton detection
Single Molecule Sequencing			
PacBio SMRT	Pacific Biosciences	Pol performance	Real-time-Seq
Nanopore	Oxford Nanopore Tech/McNally	Translocation	NA

- The principle of "**Sequencing By Synthesis**" (more commonly **SBS**)
= tracking the addition of labeled nucleotides as the DNA chain is copied.
- This method involves to:
 - Take a single-stranded DNA to sequence
 - Synthesize its complementary strand enzymatically
 - Detect light emitted sequentially
- **The DNA sequence is determined by the light emitted during the incorporation of nucleotides**, knowing that only one of the possible four bases A/T/C/G is incorporated.

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

- The DNA template is immobilized.
- Solutions of A, nucleotides C , G and T sequentially added and removed.
- Light is generated when a nucleotide complements the unpaired base.
- Chemiluminescent signal detected to determine the sequence.



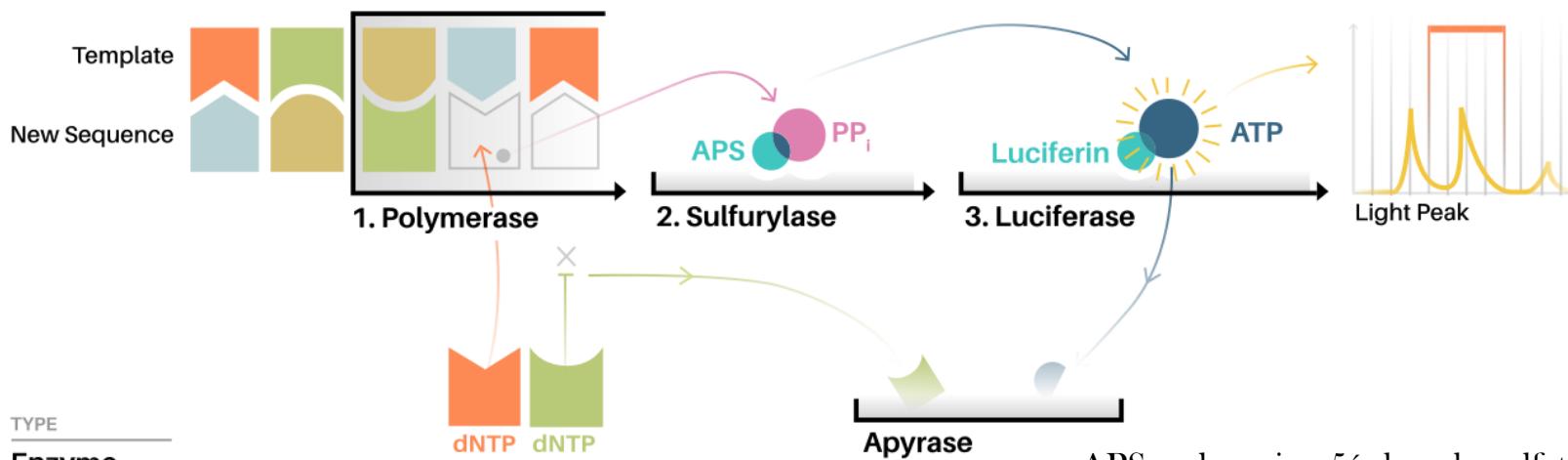
PART
2

- **Pyrosequencing** : Basically, the process allows for sequencing a single stranded DNA by synthesizing the complementary strand over its entire length, a pair of bases at a time, and detection at each step of the base actually added.
- Pyrosequencing is based on the **detection of the activity of a DNA polymerase** (the enzyme for DNA synthesis).
- The DNA sequence is determined by the light (**luciferase activity**) emitted during the incorporation of nucleotides, knowing that only one of the possible four bases A/T/C/G is incorporated.

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

- Incorporation of the complementary dNTPs (added sequentially) by DNAPol releases pyrophosphate (PPi).
- ATP sulfurylase converts PPi to ATP in the presence of APS.
- ATP = substrate for the luciferase-mediated conversion of luciferin to oxyluciferin This conversion generates visible light in amounts proportional to the amount of ATP detected by a camera and analyzed in a pyrogram.
- Unincorporated nucleotides and ATP are degraded by the apyrase, and the reaction can restart with another nucleotide



TYPE

Enzyme

Catalyst

Label

PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(<https://en.wikipedia.org>)

APS : adenosine 5' phosphosulfate

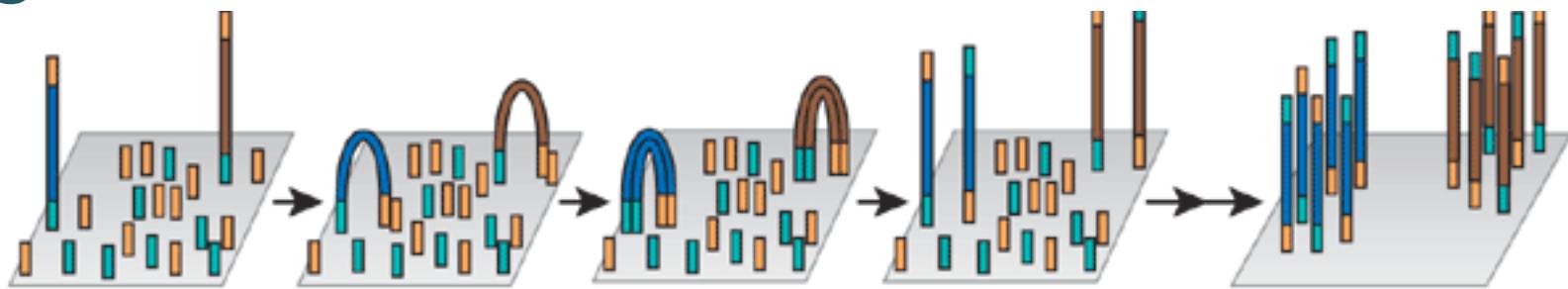
PPi : Pyrophosphate

DNApol: DNA Polymerase

NGS

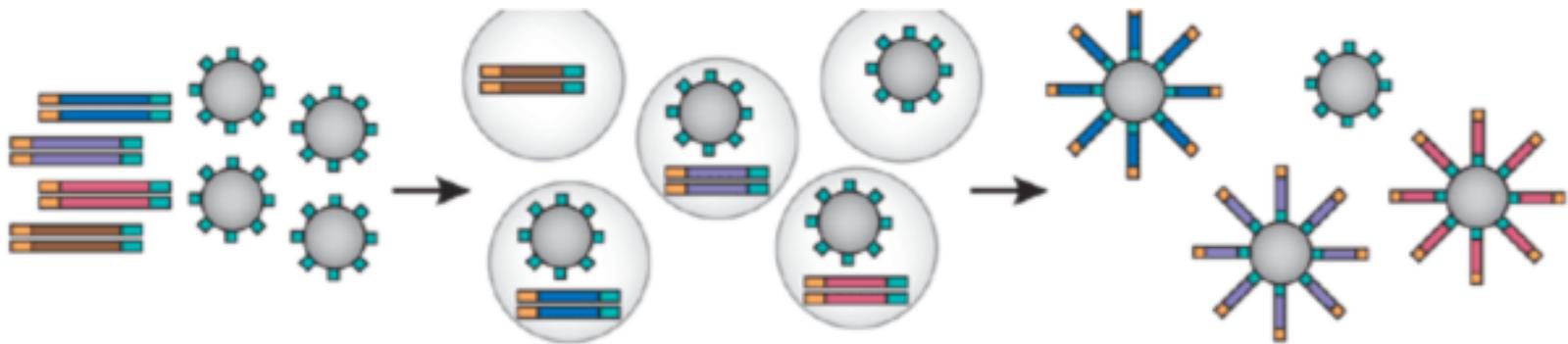
FATMA GUERFALI

Bridge PCR



- The **adaptor-flanked shotgun library** is **PCR amplified** on a flow cell
- both **primers coat** the surface of a solid substrate
- Amplification products from any given member of the library remain locally fixed near the point of origin = **cluster**
- The PCR produces **clonal clusters** that **contains copies** of a single DNA.

emPCR



- The **adaptor-flanked shotgun library** is **PCR amplified** in the context of a **water-in-oil emulsion**.
- **PCR primer** is **5'-attached** on micron-scale **beads**.
- 1 bead-containing compartments = 0 or 1 template DNA.
- **PCR amplicons** are **captured** to the surface of the bead.
- 1 clonally amplified bead = PCR products corresponding to amplification of a single molecule from the library.

Next Generation of Sequencing Technologies

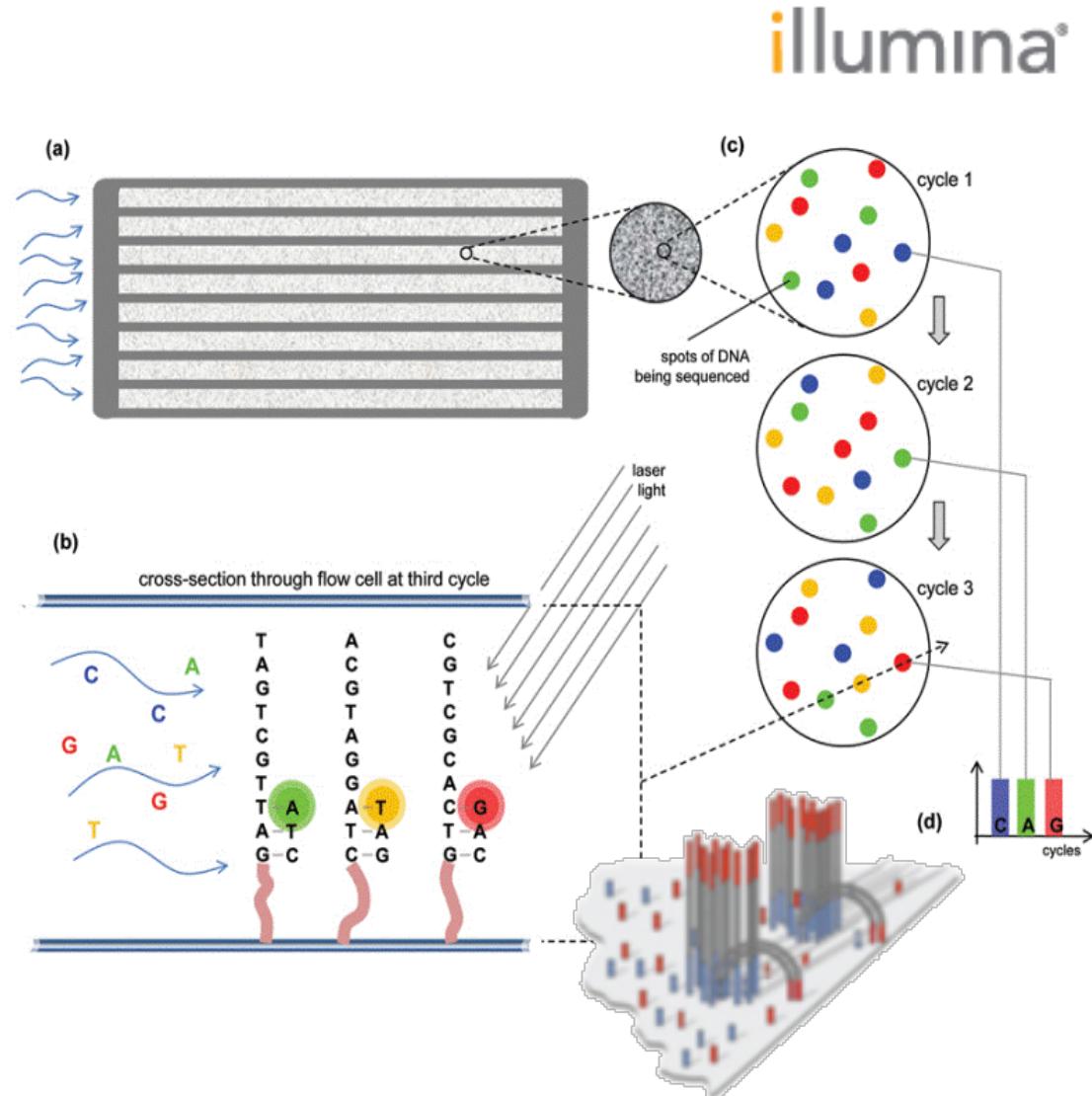
Technology	Company	Support	Chemistry
Massively Parallel Sequencing			
Solexa	Illumina	Bridge PCR on flowcell	Seq-By-Synthesis
454	Roche Applied Science	emPCR on beads	Pyrosequencing
SOLiD	AB / Life Technologies	emPCR on beads	Seq-By-Ligation
Ion Torrent	Life Technologies	emPCR on beads	Proton detection
Single Molecule Sequencing			
PacBio SMRT	Pacific Biosciences	Pol performance	Real-time-Seq
Nanopore	Oxford Nanopore Tech/McNally	Translocation	NA

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

Solexa (Illumina)

- The input sample must be cleaved into short sections.
- Fragments are ligated to adaptors and annealed to the slide using the adaptors.
- Fragments are separated into single strands to be sequenced.
- Nucleotides are modified so that each emits a different coloured light when excited by a laser.
- they have a terminator, so that only 1 base is added at a time.
- PCR, process repeated in cycles, images analyzed (SBS).



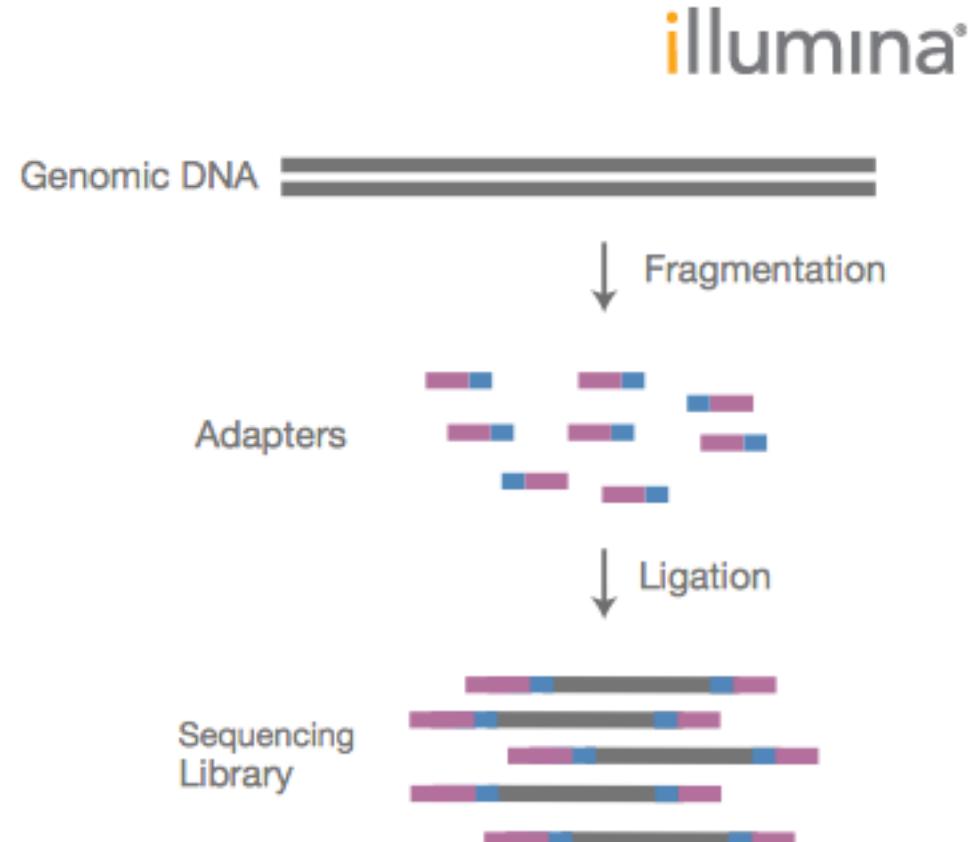
PART
2

1. Library Preparation

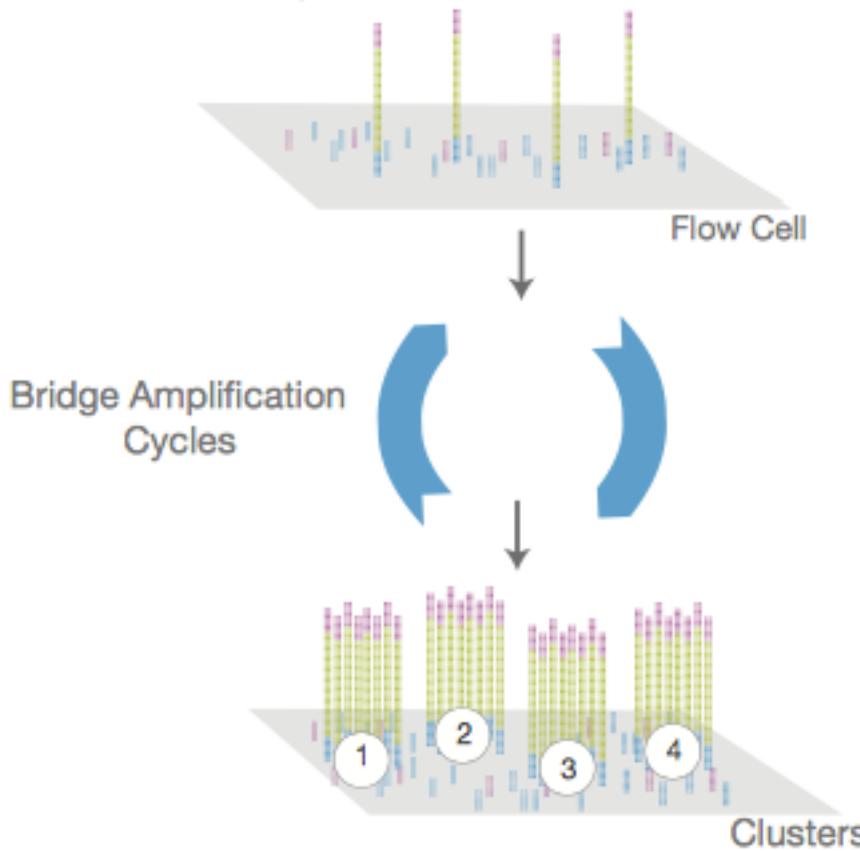
- Random fragmentation of the DNA or cDNA sample
- 5' and 3' adapter ligation.

N B : Alternatively, "tagmentation" combines the fragmentation and ligation reactions into a single step that greatly increases the efficiency of the library preparation process.

- Adapter-ligated fragments are then PCR amplified and gel purified.



NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.



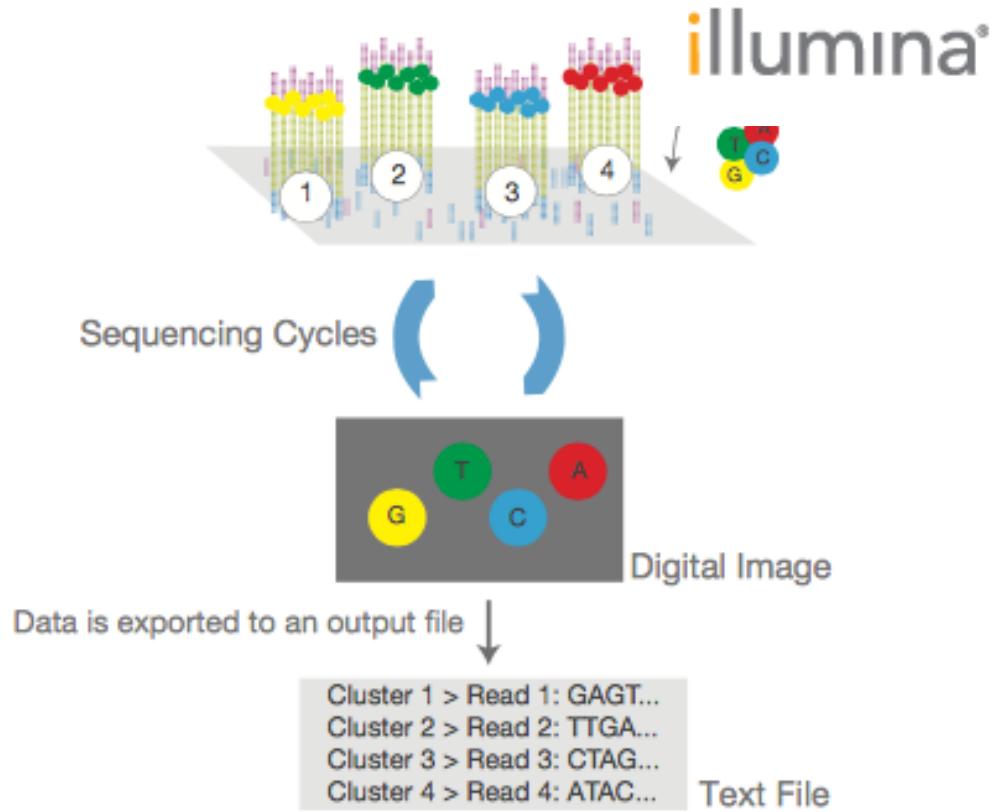
Library is loaded into a flow cell and the fragments hybridize to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

2. Cluster Generation

- library is loaded into a flow cell
- fragments are captured on a surface-bound oligos complementary to the library adapters.
- Each fragment is then amplified into distinct, clonal clusters through bridge amplification.
- When cluster generation is complete, the templates are ready for sequencing.

3. Sequencing

- Illumina SBS technology :
= reversible terminator-based method detecting single bases as they are incorporated into DNA template strands.
- All 4 reversible terminator-bound dNTPs are present during each sequencing cycle
= natural competition minimizes incorporation bias and greatly reduces raw error rates.



Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.



Reads	ATGGCATTGCAATTGACAT TGGCATTGCAATTG AGATGGTATTG GATGGCATTGCAA GCATTGCAATTGAC ATGGCATTGCAATT AGATGGCATTGCAATTG
Reference Genome	AGATGGTATTGCAATTGACAT

Reads are aligned to a reference sequence with bioinformatics software. After alignment, differences between the reference genome and the newly sequenced reads can be identified.

4. Alignment & Data Analysis

- The newly identified sequence reads are aligned to a reference genome.
- Following alignment, many **variations of analysis** are possible such as single nucleotide polymorphism (SNP) or insertion-deletion (indel) identification, read counting for RNA methods, phylogenetic or metagenomic analysis,...
- SBS sequencing: www.illumina.com/SBSvideo.

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

Solexa (Illumina)



MiniSeq System



MiSeq Series



NextSeq Series



HiSeq Series



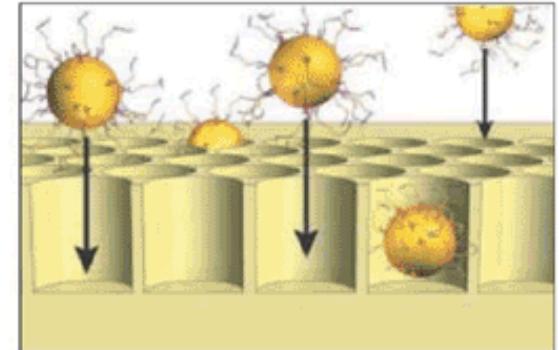
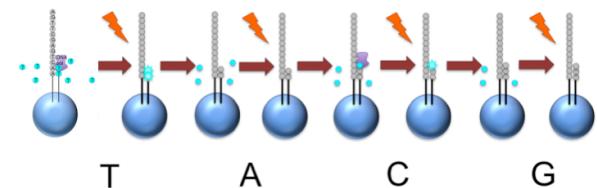
HiSeq X Series*

	MiniSeq System	MiSeq Series	NextSeq Series	HiSeq Series	HiSeq X Series*
Key Methods	Amplicon, targeted RNA, small RNA, and targeted gene panel sequencing.	Small genome, amplicon, and targeted gene panel sequencing.	Everyday exome, transcriptome, and targeted resequencing.	Production-scale genome, exome, transcriptome sequencing, and more.	Population- and production-scale whole-genome sequencing.
Maximum Output	7.5 Gb	15 Gb	120 Gb	1500 Gb	1800 Gb
Maximum Reads per Run	25 million	25 million [†]	400 million	5 billion	6 billion
Maximum Read Length	2 × 150 bp	2 × 300 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp
Run Time	4–24 hours	4–55 hours	12–30 hours	<1–3.5 days (HiSeq 3000/HiSeq 4000) 7 hours–6 days (HiSeq 2500)	<3 days
Benchtop Sequencer	Yes	Yes	Yes	No	No
System Versions	<ul style="list-style-type: none"> • MiniSeq System for low-throughput targeted DNA and RNA sequencing 	<ul style="list-style-type: none"> • MiSeq System for targeted and small genome sequencing • MiSeq FGx System for forensic genomics • MiSeqDx System for molecular diagnostics 	<ul style="list-style-type: none"> • NextSeq 500 System for everyday genomics • NextSeq 550 System for both sequencing and cytogenomic arrays 	<ul style="list-style-type: none"> • HiSeq 3000/HiSeq 4000 Systems for production-scale genomics • HiSeq 2500 Systems for large-scale genomics 	<ul style="list-style-type: none"> • HiSeq X Five System for production-scale whole-genome sequencing • HiSeq X Ten System for population-scale whole-genome sequencing

PART
2

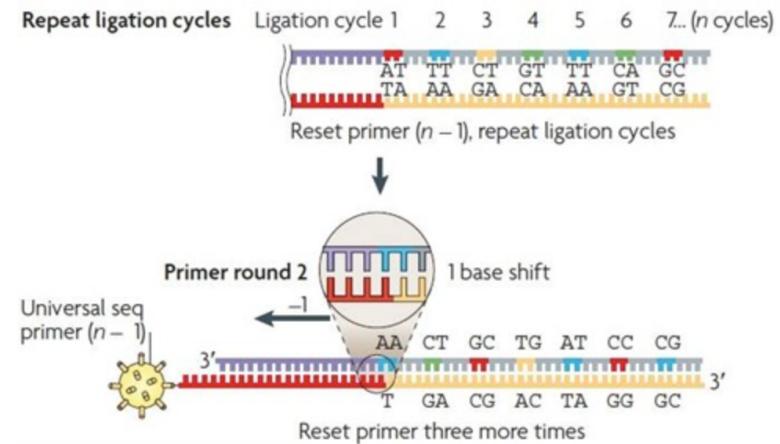
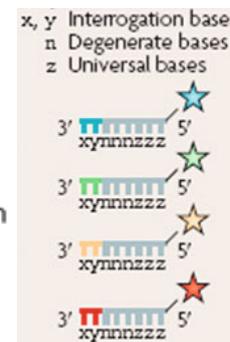
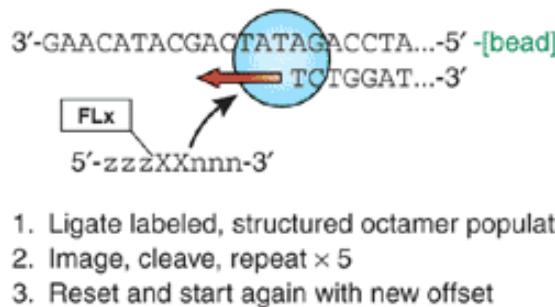


- As in Illumina, the DNA is fragmented.
- Adaptors added, end annealed to beads.
1 DNA fragment = 1 bead (theory).
- Fragments amplified by PCR using adaptor-specific primers.
- The sequence can then be determined computationally.
- Longer reads than Illumina, different lengths.



SOLID = Sequencing by Oligonucleotide Ligation and Detection

- Sequencing uses a **ligase**, rather than a polymerase → **Sequencing-By-Ligation**
- Each sequencing cycle introduces a **partially degenerated population** of fluorescently **labeled octamers**. The population is structured such that the **label correlates with the identity of the central 2 bp in the octamer**.
- After ligation and imaging in four channels, the **labeled portion of the octamer** (that is, 'zzz') is cleaved leaving a free end for another cycle of ligation.
→ Reading errors reduced compared to pyrosequencing.

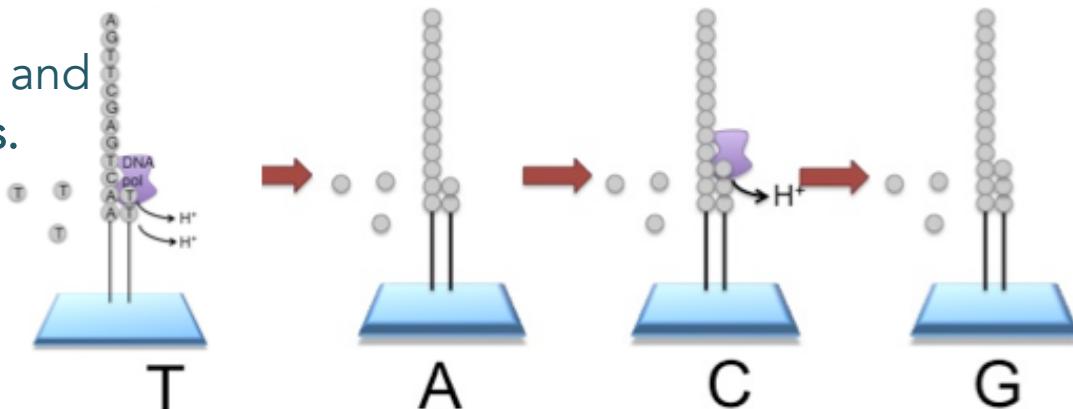


Adapted from (Shendure & Ji, 2008)

Adapted from (<http://www.springer.com/978-1-4614-7725-9>)

Adapted from (<http://www.oezratty.net/>)

- As in other kinds of NGS, the input DNA is fragmented.
- Unlike other methods, Ion Torrent do not use optical signalling.
- Adaptors are added and one molecule is placed onto a bead.
- Amplification on the bead by emulsion PCR. Each bead is placed into 1 well of a slide.
- The pH is detected in each of the wells, as each H⁺ ion released will decrease the pH. The changes in pH allow us to determine if that base, and how many thereof, was added to the sequence read.
- The dNTPs are washed away, and the process is repeated in cycles.





Targeted DNA
Sequencing



Targeted RNA
Sequencing



Microbial
Sequencing



Ion PGM System performance specifications

	Ion 314™ Chip v2 or Ion 314 Chip v2 BC	Ion 316™ Chip v2 or Ion 316 Chip v2 BC	Ion 318™ Chip v2 or Ion 318 Chip v2 BC	
Output*	200 base 400 base [†]	30–50 Mb 60–100 Mb	300–600 Mb 600 Mb–1Gb	600 Mb–1Gb 1.2–2 Gb
Reads		400–550 thousand	2–3 million	4–5.5 million
Run time	200 base 400 base	2.3 hr 3.7 hr	3.0 hr 4.9 hr	4.4 hr 7.3 hr
Research areas	Cancer research, inherited disease research, microbial genomics, stem cell research, agriculture, epigenomics, metagenomics, forensic science, and ancient DNA genomics			
Key applications	Targeted DNA sequencing, copy number analysis, targeted RNA sequencing, small RNA sequencing, <i>de novo</i> microbial sequencing, bacterial typing research, viral typing research, ChIP sequencing, methylation analysis, SNP verification, and genotyping by sequencing			



The Ion Proton™ System

Rapid, high-throughput
benchtop sequencing

The Ion Proton™ System minimizes the high cost and complexity of high-throughput sequencing, bringing it to your research lab, on your budget, and on your schedule. This system enables high-quality exome and transcriptome sequencing in just a few days.

Ion Proton™ System performance specifications with Ion PI™ Chip

Output	Up to 10 Gb		
Reads	60–80 million reads		
Read length	Up to 200-base reads		
Run time	2–4 hours		
Research areas	Cancer Genetic disorders	Agriculture Stem cells	Epigenetics
Key applications	Human-scale genome sequencing Exome sequencing	Transcriptome sequencing Copy number analysis	ChIP sequencing Targeted RNA sequencing

Ion PGM™ Chip	Run time		Output	
	200 bp read	400 bp read	200 bp read	400 bp read
Ion 314™ Chip v2	2.3 hr	3.7 hr	30–50 Mb	60–100 Mb
Ion 316™ Chip v2	3.0 hr	4.9 hr	300–500 Mb	600 Mb–1 Gb
Ion 318™ Chip v2	4.4 hr	7.3 hr	600 Mb–1 Gb	1.2–2 Gb
Ion Proton™ Chip	Run time		Output	
	200 bp read		200 bp read	
Ion PI™ Chip	2–4 hr		Up to 10 Gb	

Next Generation of Sequencing Technologies

Technology	Company	Support	Chemistry
Massively Parallel Sequencing			
Solexa	Illumina	Bridge PCR on flowcell	Seq-By-Synthesis
454	Roche Applied Science	emPCR on beads	Pyrosequencing
SOLiD	AB / Life Technologies	emPCR on beads	Seq-By-Ligation
Ion Torrent	Life Technologies	emPCR on beads	Proton detection
Single Molecule Sequencing			
PacBio SMRT	Pacific Biosciences	Pol performance	Real-time-Seq
Nanopore	Oxford Nanopore Tech/McNally	Translocation	NA

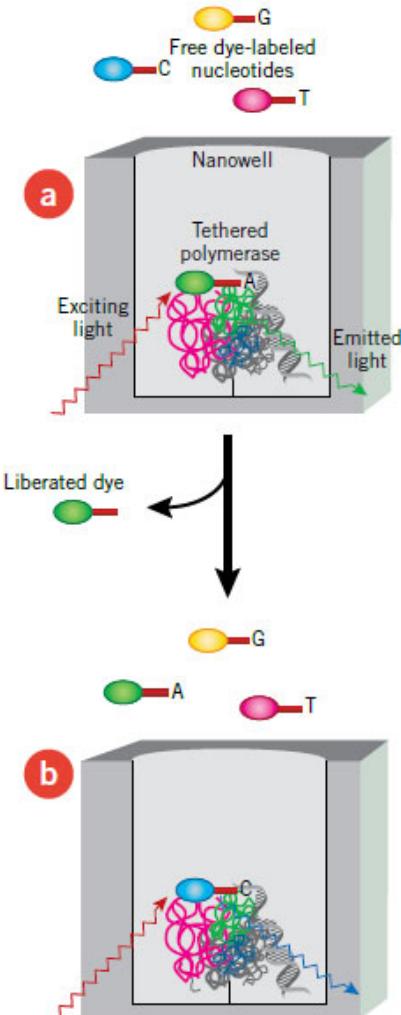
Single Molecule Real Time (SMRT) DNA Seq technology

SMRT Sequencing enables real-time observation of DNA synthesis. SMRT Sequencing is built upon 2 key innovations:

- zero-mode waveguides (ZMWs)
- phospholinked nucleotides.

a → DNA polymerase molecule is tethered to the bottom of a nanowell → ZMW design ensures only one nucleotide-linked dye can be directly excited at a time.

b → Each incorporated phospholinked nucleotide will reside on the enzyme's active site for a few milliseconds, which is enough time for a fluorescent signal to be recorded.

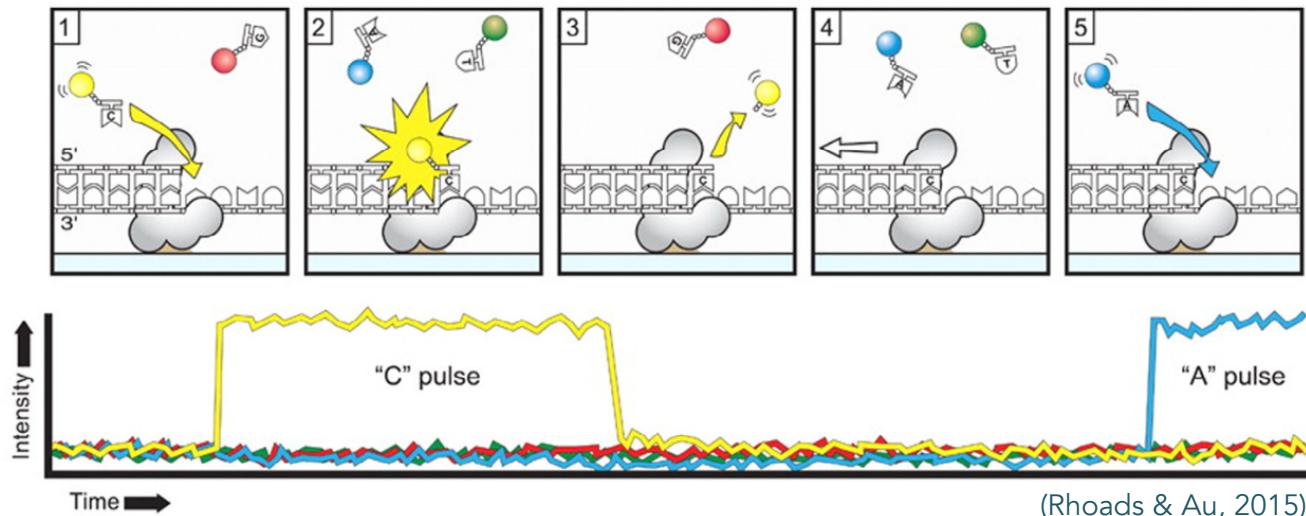
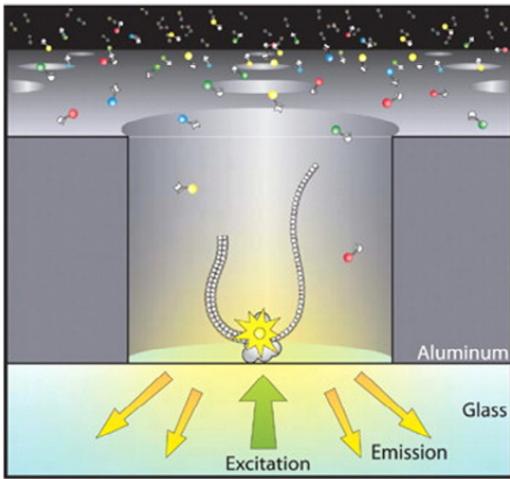


(Osherovich et al., 2010)

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

PacBio

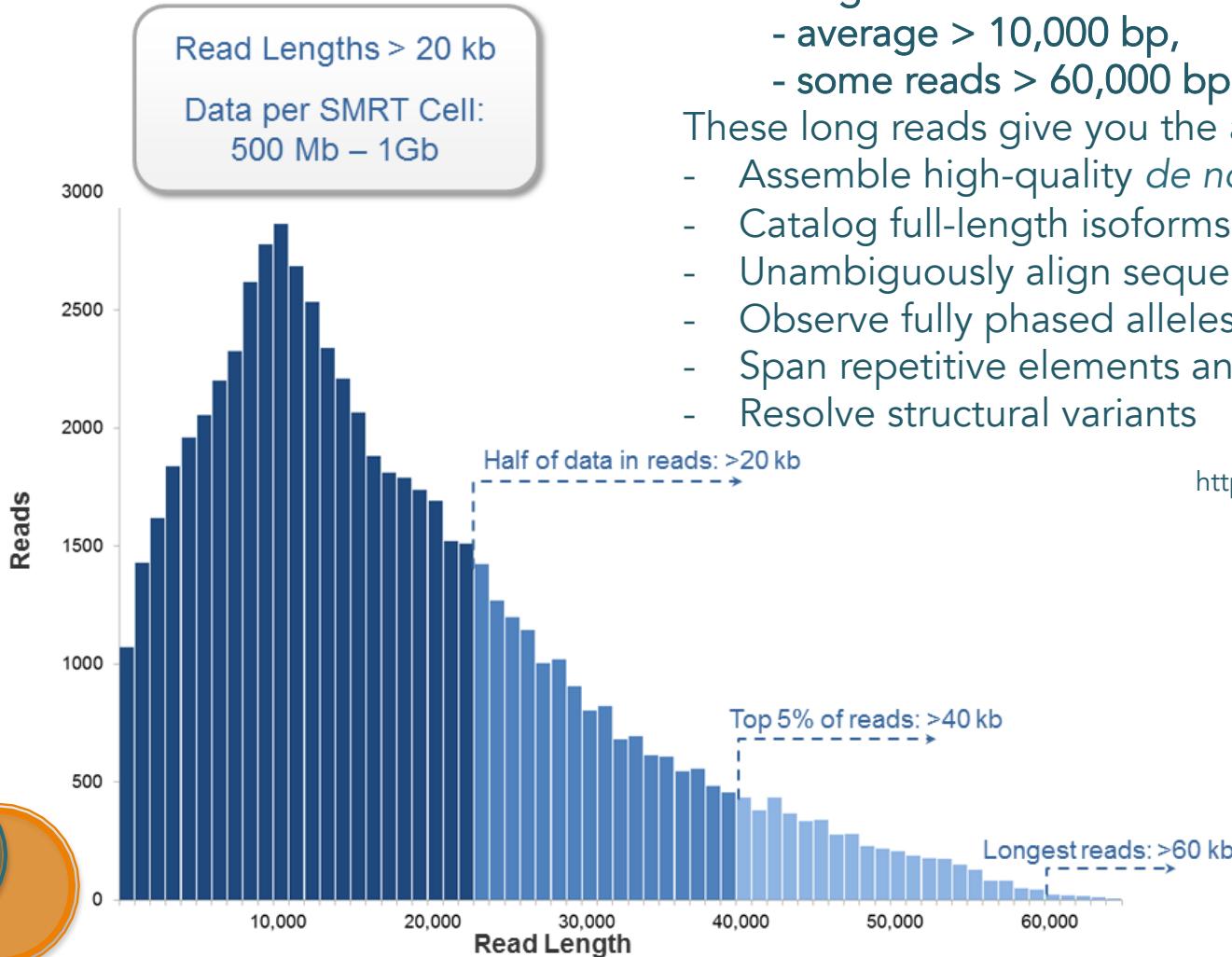


PCR	# PCR	% closed by Sanger	% closed by PacBio
<2.5 kb	246	64	73
>2.5 kb	113	0	88
hard stop	3	0	100

362 PCR products (each covering a different gap) were sequenced with both Sanger and PacBio technologies. While the majority of gaps less than 2.5 kb were closed with both Sanger (64%) and PacBio (73%) technologies, none of the gaps larger than 2.5 kb were closed with a single round of Sanger technology. Three hard stop gaps were all closed by PacBio.

2

Single Molecule Real Time (SMRT) DNA Seq technology



Read length:

- average > 10,000 bp,
- some reads > 60,000 bp

These long reads give you the ability to:

- Assemble high-quality *de novo* genomes
- Catalog full-length isoforms
- Unambiguously align sequences
- Observe fully phased alleles
- Span repetitive elements and complex regions
- Resolve structural variants

[http://www.pacb.com/smrt-science/
smrt-sequencing/read-lengths/](http://www.pacb.com/smrt-science/smrt-sequencing/read-lengths/)

Single Molecule Real Time (SMRT) DNA Seq technology

Sequel System: high-throughput, cost-effective access to SMRT Sequencing



PacBio RS II: the original long-read sequencer

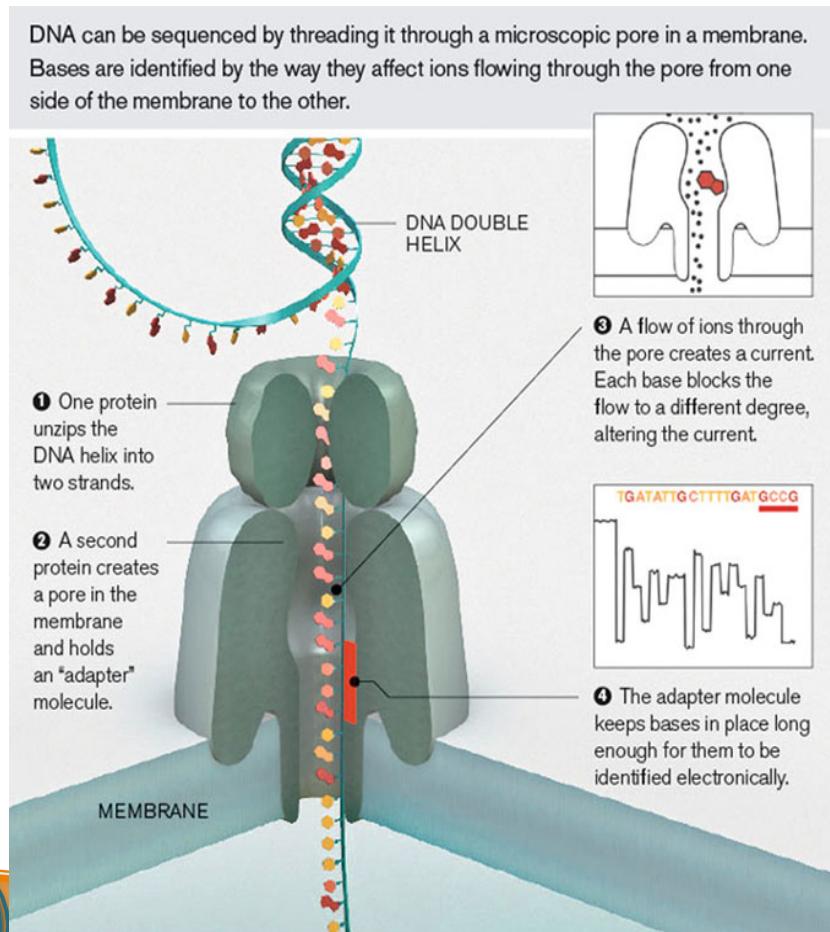


The Sequel System delivers $\geq 7\times$ more reads with 1 million zero-mode waveguides (ZMWs) per SMRT Cell. The Sequel System is ideal for projects such as rapidly and cost-effectively generating high-quality whole genome de novo assemblies.

The PacBio RS II is ideal for whole genome sequencing of small genomes, targeted sequencing, complex population analysis, RNA sequencing of targeted transcripts, and microbial epigenetics.

<http://www.pacb.com/>

Single Molecule Real Time (SMRT) DNA Seq technology



Schematic representation of nanopore sequencing system.

- Upper protein → ssDNA.
 - 2nd protein
 - forms a nanopore in a membrane.
 - contains an adaptor molecule that reduces the speed of passing DNA through the pore.
 - Each base obstructs the flow to a different degree.
- uses an electronic rather than an optical signal to identify DNA bases

Single Molecule Real Time (SMRT) DNA



Weight

87g (103g with a flow cell)

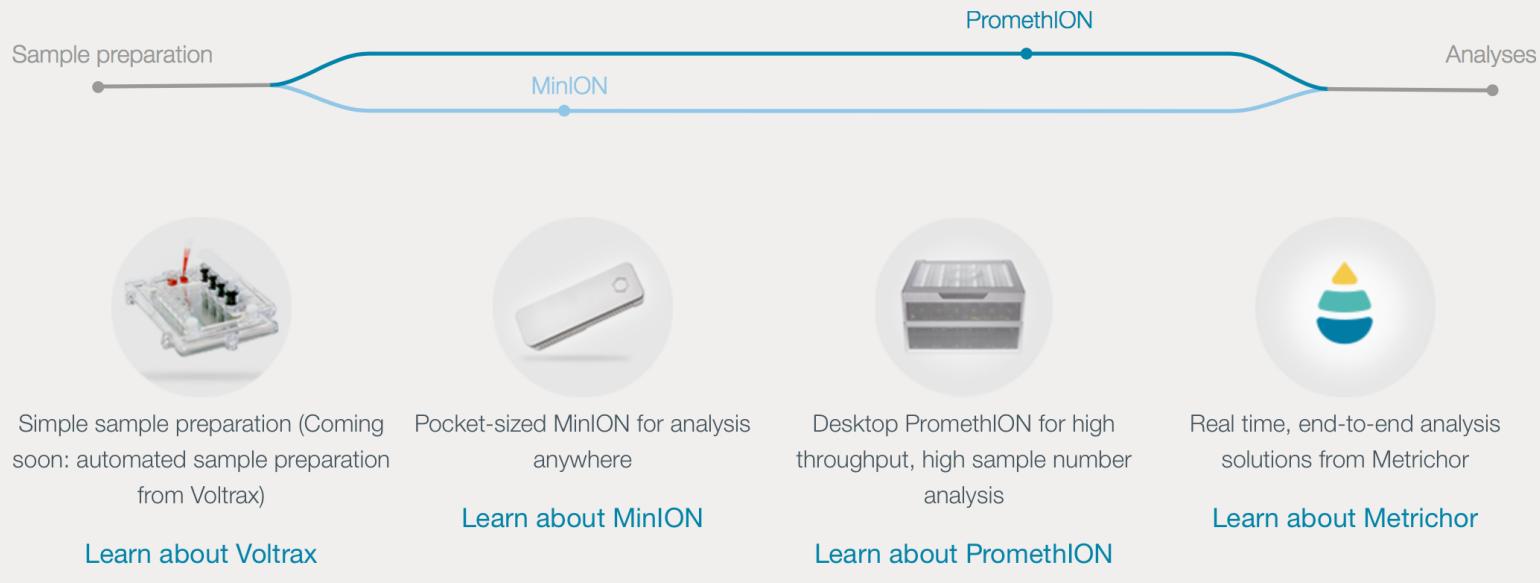


Size

W 105, H 23, D 33mm



Simple sequencing workflows

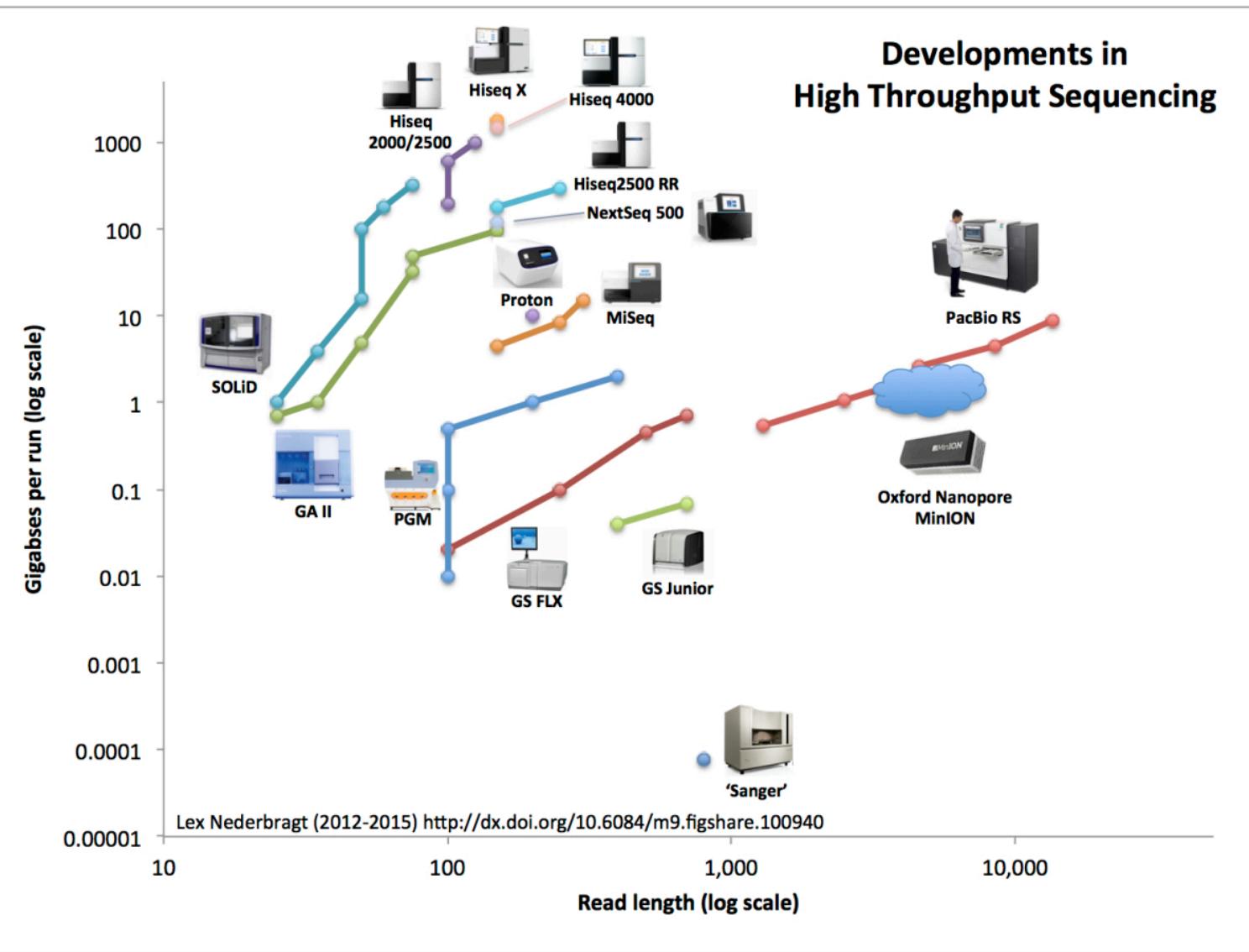


Read length : The longest read reported by a MinION user to date is more than 200Kb

Portability : The MinION can be used outside the traditional lab environment to perform field-based work

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING



PART
2

The four main advantages of NGS over classical Sanger sequencing

- **Speed**

NGS is quicker than Sanger sequencing in two ways.

- Chemical reaction may be combined with the signal detection, whereas in Sanger sequencing these are two separate processes.
- NGS is massively parallel.

- **Cost**

The human genome sequence cost 3 billion \$.

Sequencing a human genome with Illumina allows to reach the \$1,000 expected.

- **Sample size**

needs significantly less starting amount of DNA/RNA

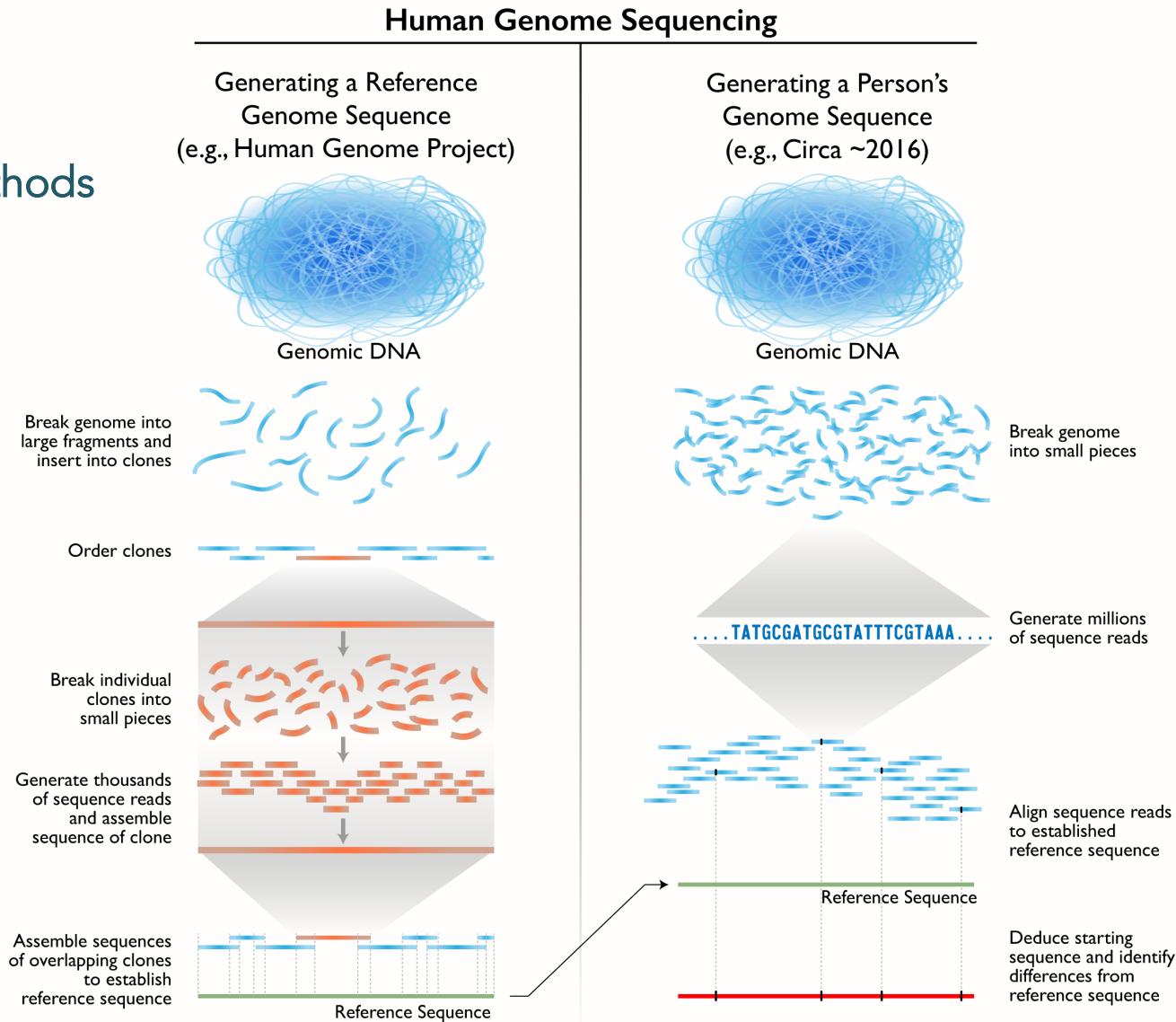
- **Accuracy**

More repeats than with Sanger sequencing → greater coverage, higher accuracy and sequence reliability (individual reads less accurate for NGS).

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING

Comparison of human genome sequencing methods HGP vs. ~ 2016



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

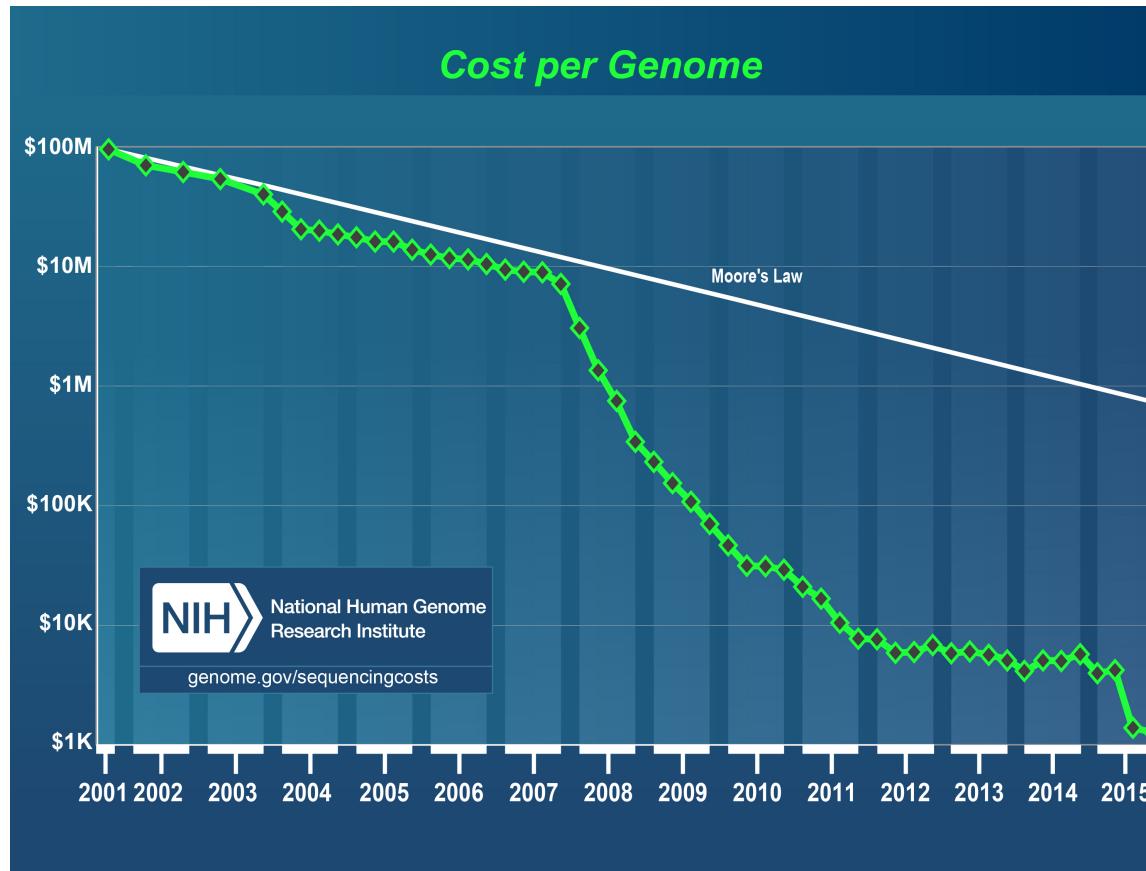
<http://www.genome.gov/sequencingcosts/>

NGS

FATMA GUERFALI

Sequencing costs vs. technological advances

Data from the NHGRI (National Human Genome Research Institute)
→ Genome Sequencing Program (GSP)



Sequencing costs vs. technological advances

HiSeq X Ten System

Highlights

- \$1000 Genome Is a Reality**
HiSeq X Ten System is the first and only platform to break the \$1000 barrier human whole-genome sequencing
- Population- and Production-Scale Whole-Genome Sequencing**
HiSeq X Ten System delivers > 18,000 human genomes per year; HiSeq X Five System delivers > 9000 human genomes per year
- Proven Performance**
Take advantage of industry-leading data quality with the highly accurate Illumina sequencing by synthesis technology
- Species Expansion**
Now enabling cost-effective whole-genome sequencing of nonhuman species with unrivaled throughput

Parameter	Specification
Output per Run	Dual flow cell: 1.6-1.8 Tb
Single Reads Passing Filter	Dual flow cell: 5.3-6 billion
Supported Read Length	2 × 150 bp
Run Time	< 3 days
Quality	≥ 75% of bases above Q30 at 2 × 150 bp
Supported Library Preparation	TruSeq DNA PCR-Free Library Prep Kit TruSeq Nano DNA Library Prep Kit

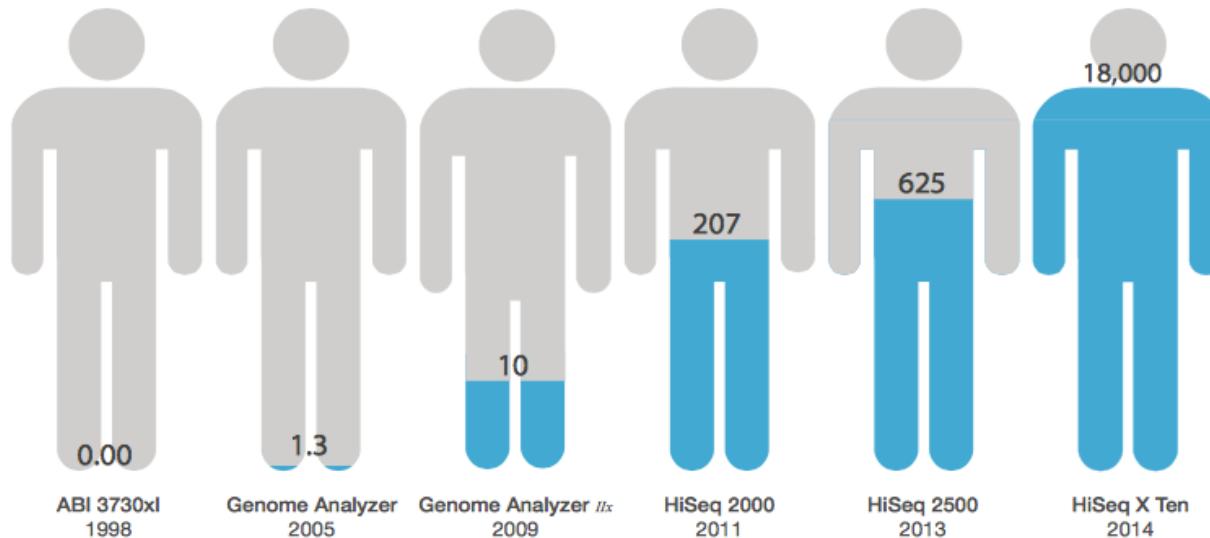


Sequencing costs vs. technological advances

1993-2001: First human genome: required 13-15 years. Cost \approx 3 billion \$.

2014: HiSeq X Ten can sequence over 45 human genomes in a single day for approximately \$1000 each $\rightarrow \approx 18,000$ genomes / year.

Human Genomes Sequenced Annually



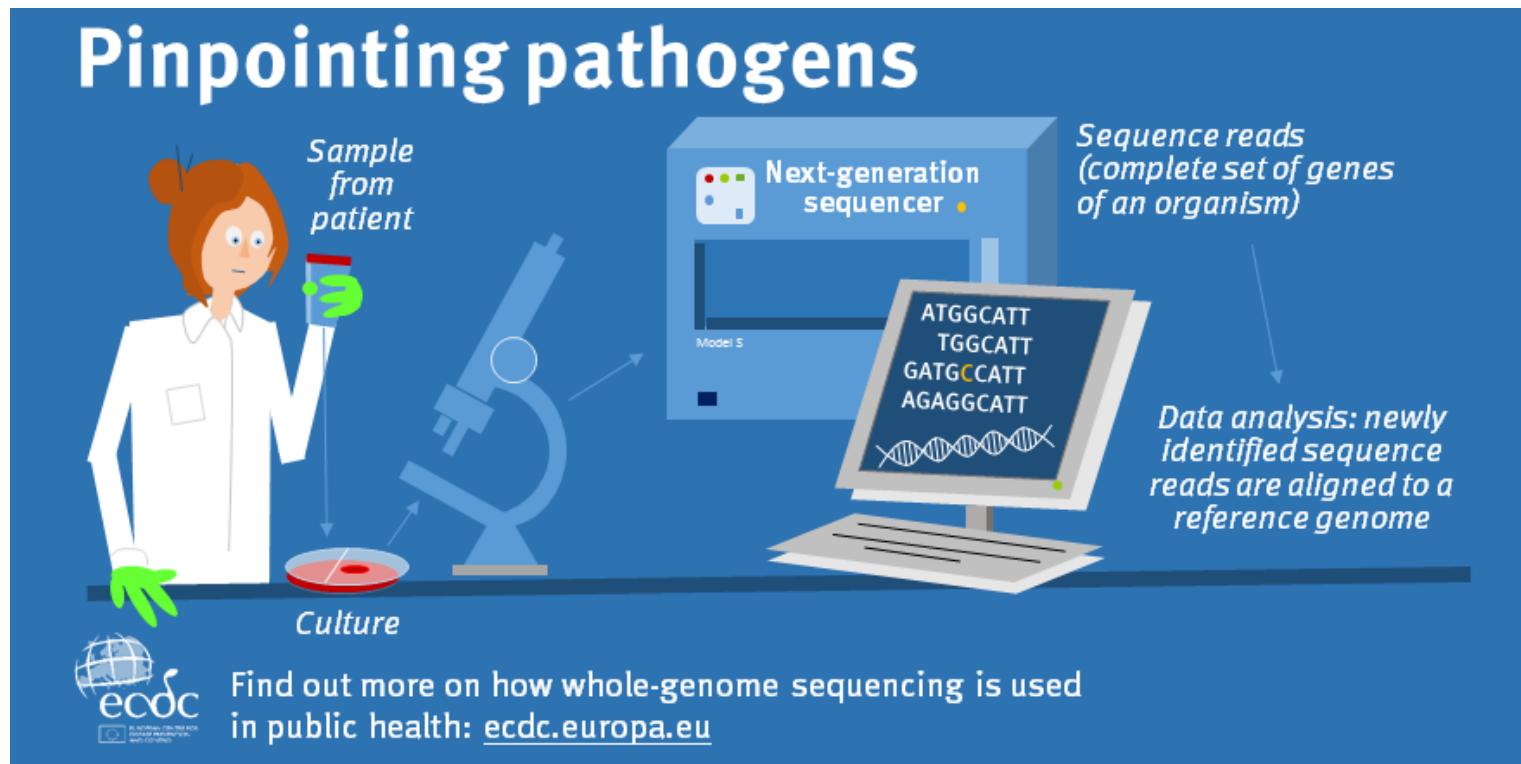
"Beyond the massive increase in data output, the introduction of NGS technology has transformed the way scientists think about genetic information."

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL_LIFE APPLICATIONS

NGS allows you to perform high-definition sequence-driven science
→ Epidemiological outbreaks...

ECDC (European Center for Disease Prevention and Control) : roadmap for integration of molecular typing and genomic typing into European-level surveillance and epidemic preparedness



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

NGS
<http://ecdc.europa.eu/en/healthtopics/microbiology/microbiology-activities/Pages/molecular-typing.aspx> FATMA GUERFALI

NGS allows you to perform high-definition sequence-driven science

→ Epidemiological outbreaks...

- *de novo* genome sequencing
- Genome assembly
- *In Vitro* Diagnostics
- Zoom in to deeply sequence target regions
- Simultaneously sequence enormous numbers of samples using multiplex sequencing with DNA barcode tags
-

(Fournier et al., 2013)

Table 1 | Examples of infectious disease outbreaks that were investigated using next-generation sequencing

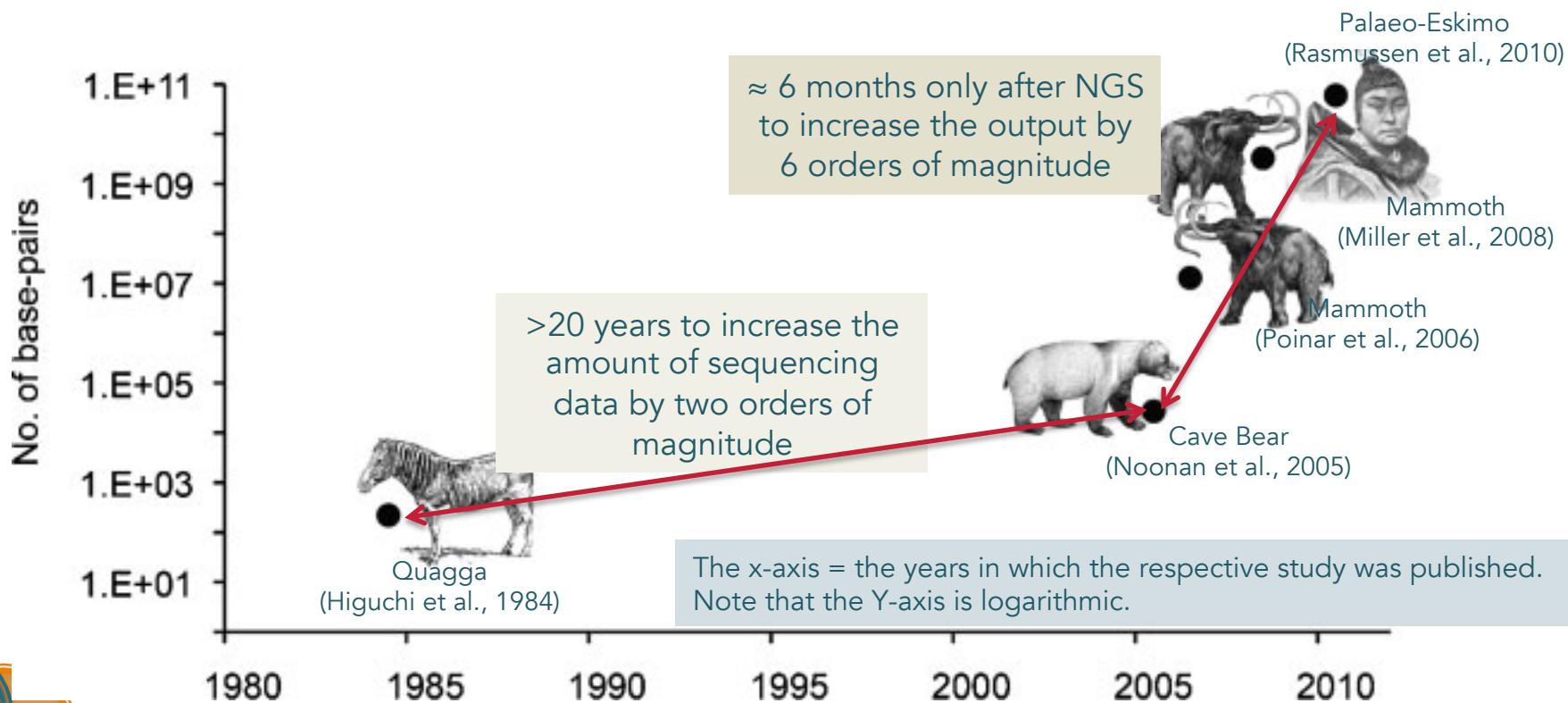
Microorganism	Location	Year	Reference
Carbapenem-resistant <i>Klebsiella pneumoniae</i>	USA	2011	112
<i>Clostridium difficile</i>	Worldwide	2013	113
<i>Escherichia coli</i> O104:H4	Germany	2011	114,115
<i>Legionella pneumophila</i> serogroup 1	United Kingdom	2013	116
Methicillin-resistant <i>Staphylococcus aureus</i> (MRSA)	United Kingdom	2009	117
<i>Mycobacterium tuberculosis</i>	Canada	2006–2008	118
<i>Vibrio cholerae</i> O1 biovar El Tor	Haiti	2010–2011	119
Arenavirus	Australia	2008	120
Bas-Congo virus	Democratic Republic of the Congo	2009	121
Influenza A virus H1N1	Worldwide	2009	122

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL_LIFE APPLICATIONS

NGS is used in many specific fields such as Paleogenomics

Ancient DNA sequencing (extinct organisms, limited sample)
→ increased output from ancient DNA after NGS



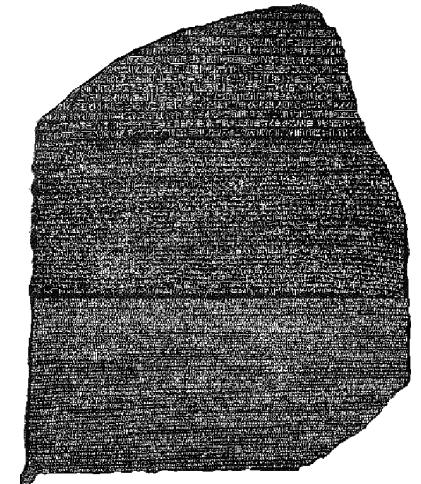
PART
2

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL_LIFE APPLICATIONS

NGS is used in many specific fields such as
Comparative Genomics

Comparative genomics is based on homology and
evolutionary dynamics between organisms



COMPARING CLOSELY RELATED SPECIES



COMPARING EVOLUTIONARILY DISTANT SPECIES



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

<http://www.beller.no>
<http://www.notcot.org/>

NGS
FATMA GUERFALI

Moreover: thinking at a population scale is possible !

- NGS allows you to think at the individual and population scales.
- Explain what determines the complexity of living organisms, etc...

Humans as an example:



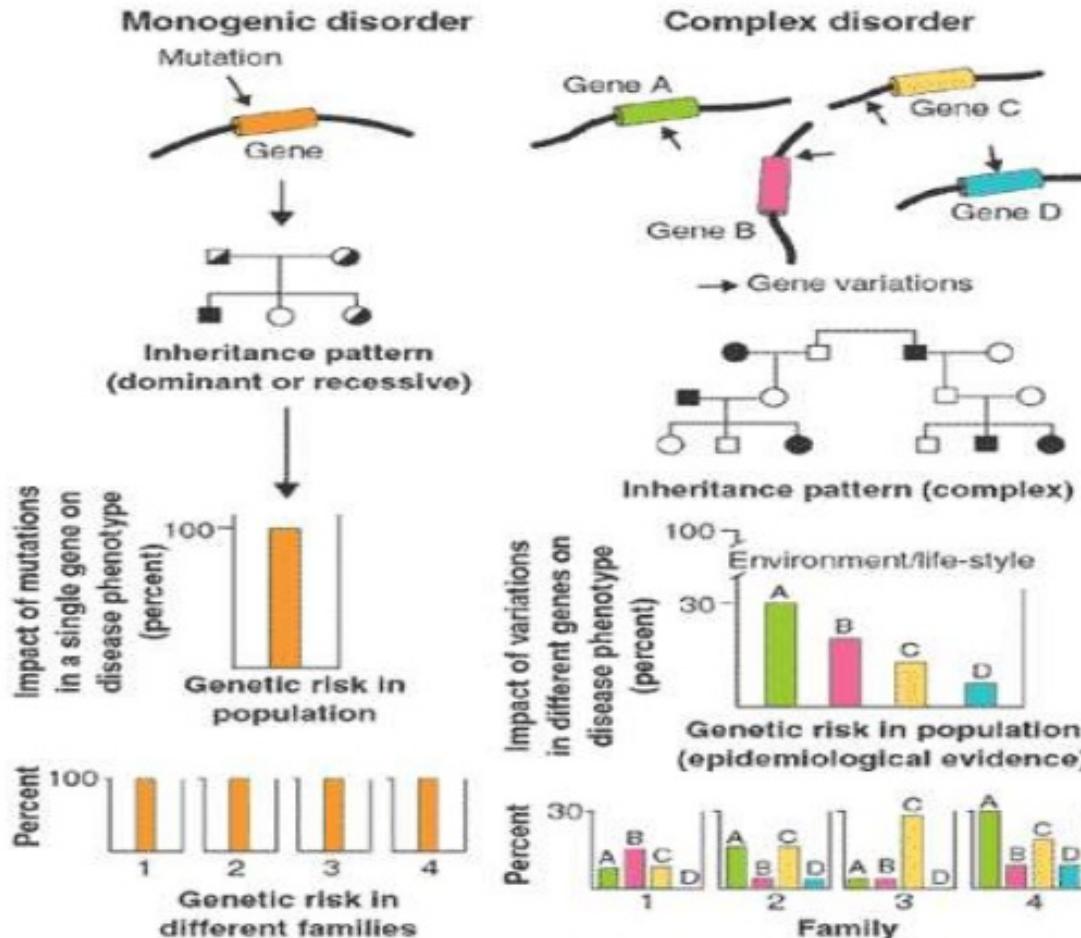
No two individuals are **similar**:

- What determines their different characteristics ?
- What determines their different susceptibility to diseases?

SEQ TECHNOLOGIES

ADVANCES MADE TO STUDY GENES AND GENE EXPRESSION

Monogenic vs. Complex Disorders

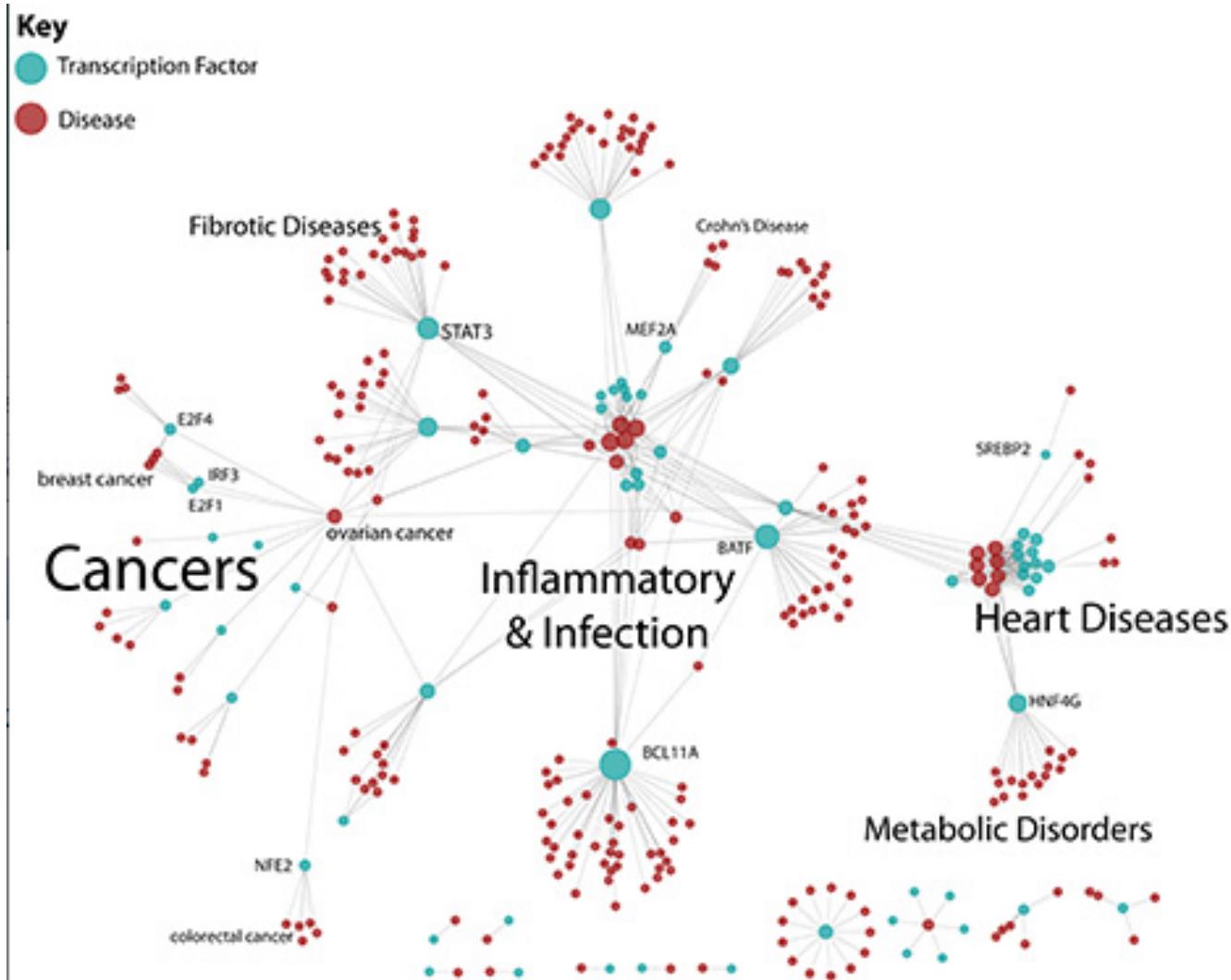


→ Link between thousands of genetic variants to complex diseases possible !

PART
2

SEQ TECHNOLOGIES

ADVANCES MADE TO STUDY GENES AND GENE EXPRESSION

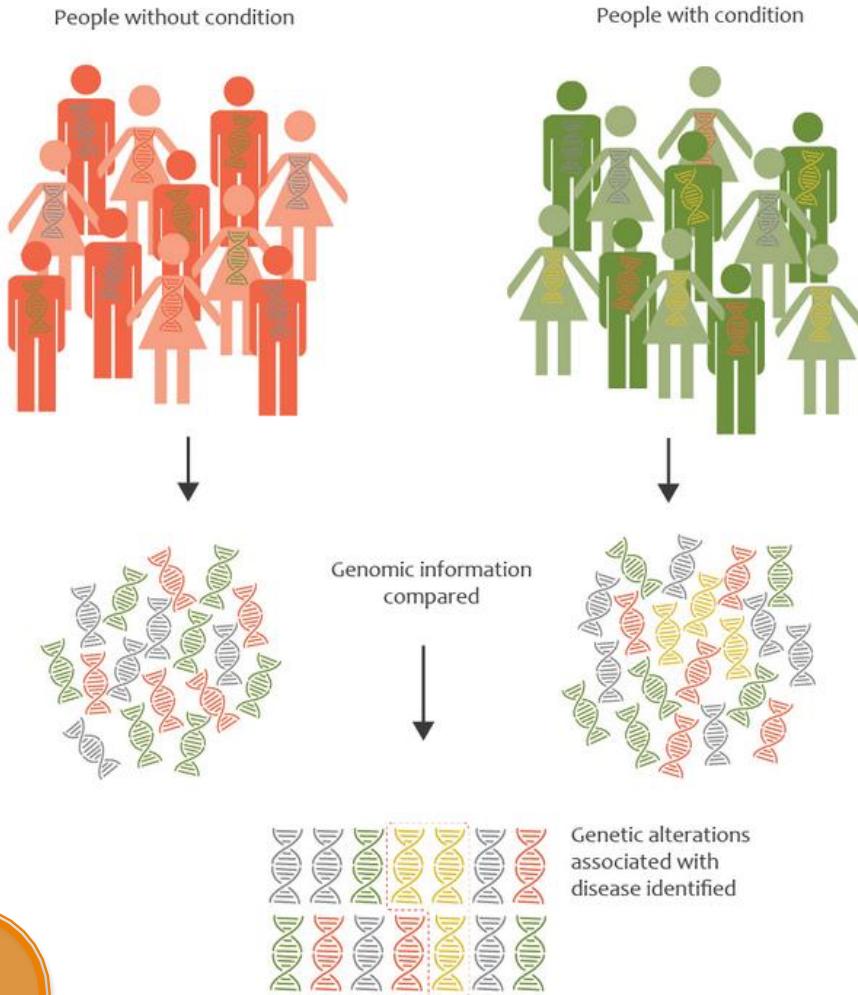


PART
2

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

How researchers compare genomic information to identify genetic alterations



Genomics

- **Genomics** is the study of genomes, including large chromosomal segments containing many genes.
- The *initial phase of genomics* aims to map and sequence an initial set of entire genomes.
- *Functional genomics* aims to deduce information about the function of DNA sequences.
 - Should continue long after the initial genome sequences have been completed.

PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

<https://twitter.com/genomicsedu/status/525219646476410880>

NGS

FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

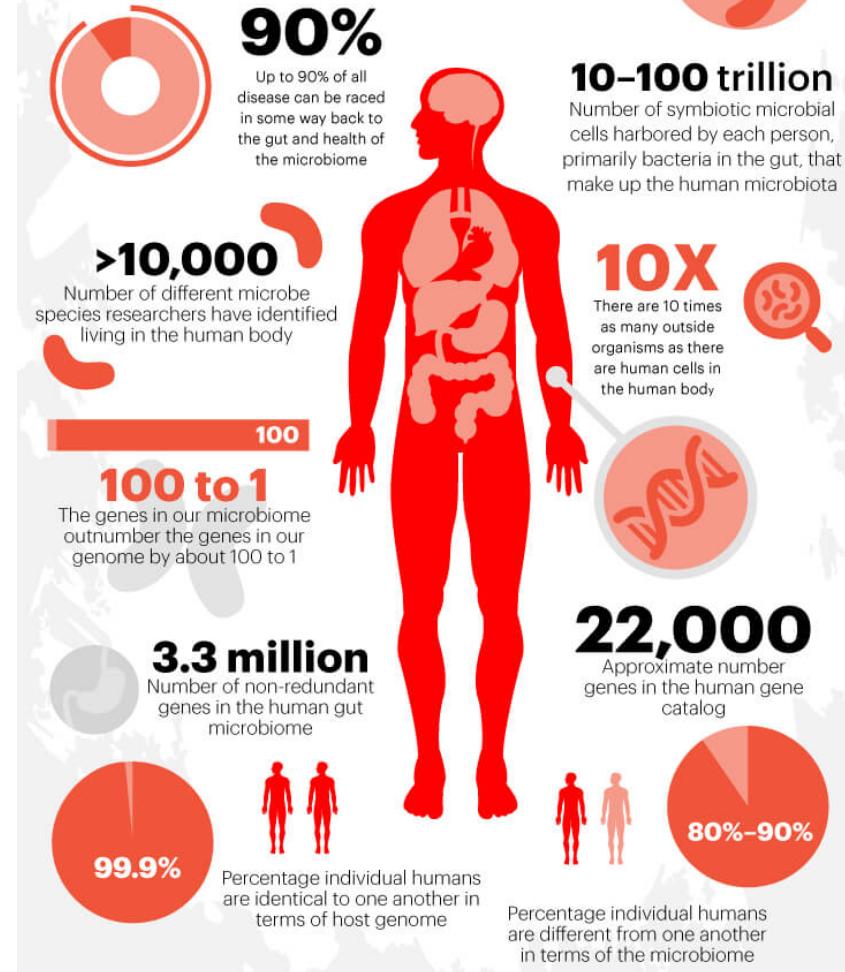
- “microbiome” : “micro” means small and “biome” means a habitat of living things.
- Popular among researchers : up to 90% of all diseases can be traced back directly or indirectly to health of the microbiome.
- Our microbiome is home to **trillions of microbes**, diverse organisms that help govern nearly every function of the human body in some way.

PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

<https://draxe.com/microbiome/>

The Importance of the MICROBIOME by the Numbers



NGS

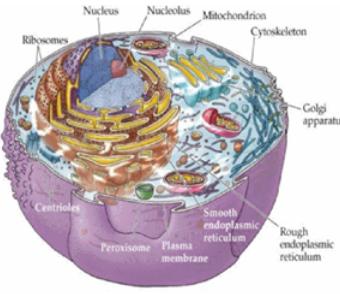
FATMA GUERFALI

GENOMES & TRANSCRIPTOMES

HGP: THE HERITAGE

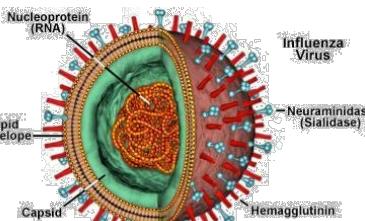


EUKARYOTES 3694

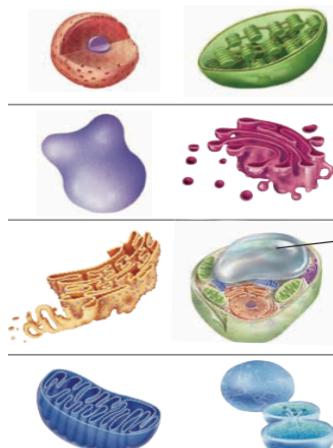


Eukaryote

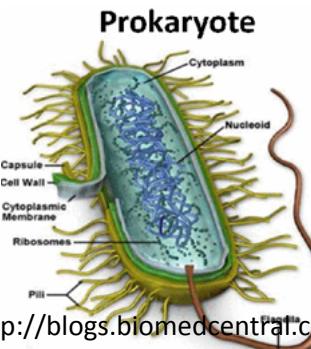
VIRUSES 5884



ORGANELLES 8634



PROKARYOTES 74156



Perman, 2014: <http://blogs.biomedcentral.com/>
<https://www.ncbi.nlm.nih.gov/genome/browse/>
<https://www.quora.com/>

FATMA GUERFALI

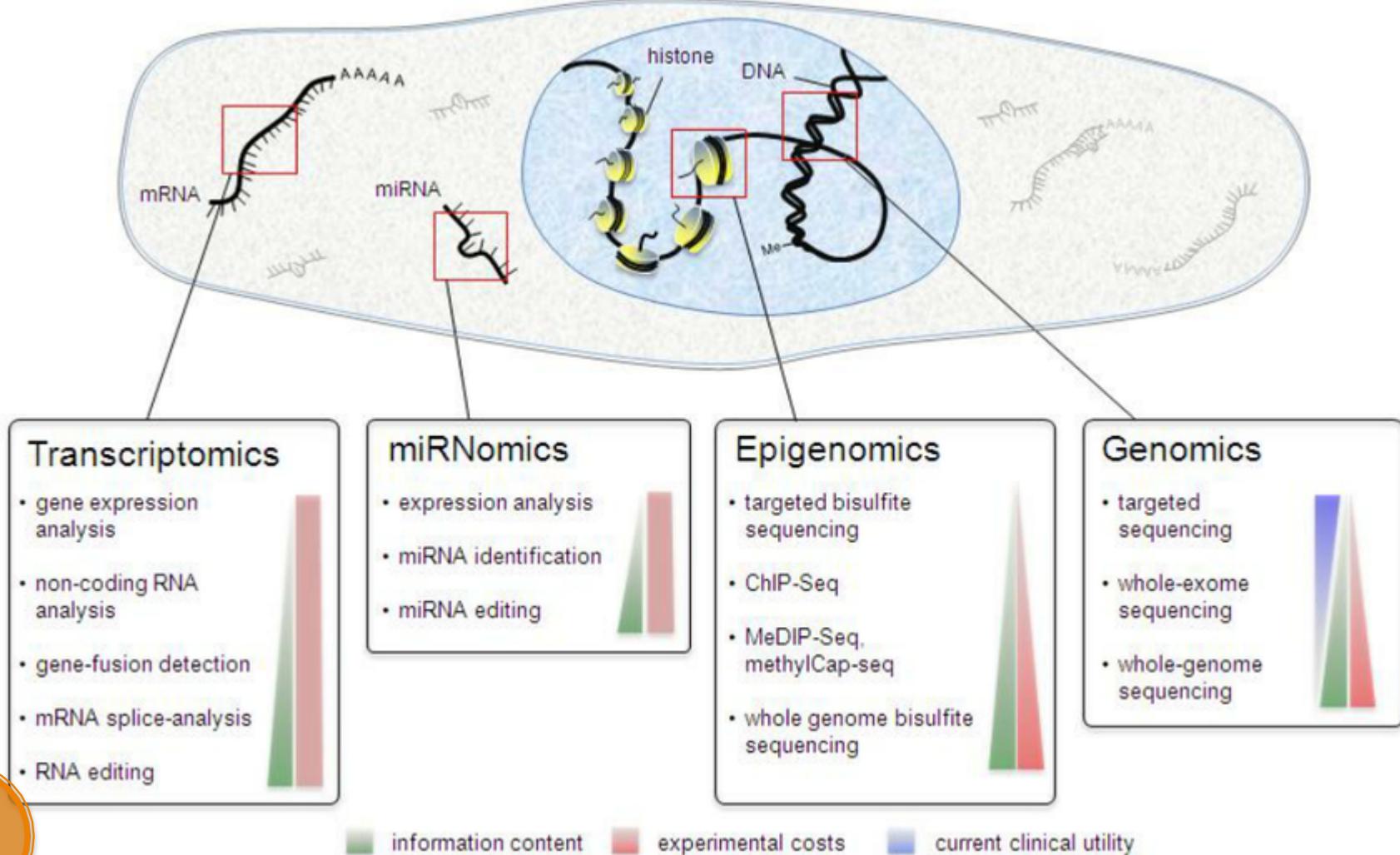
PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL_LIFE APPLICATIONS

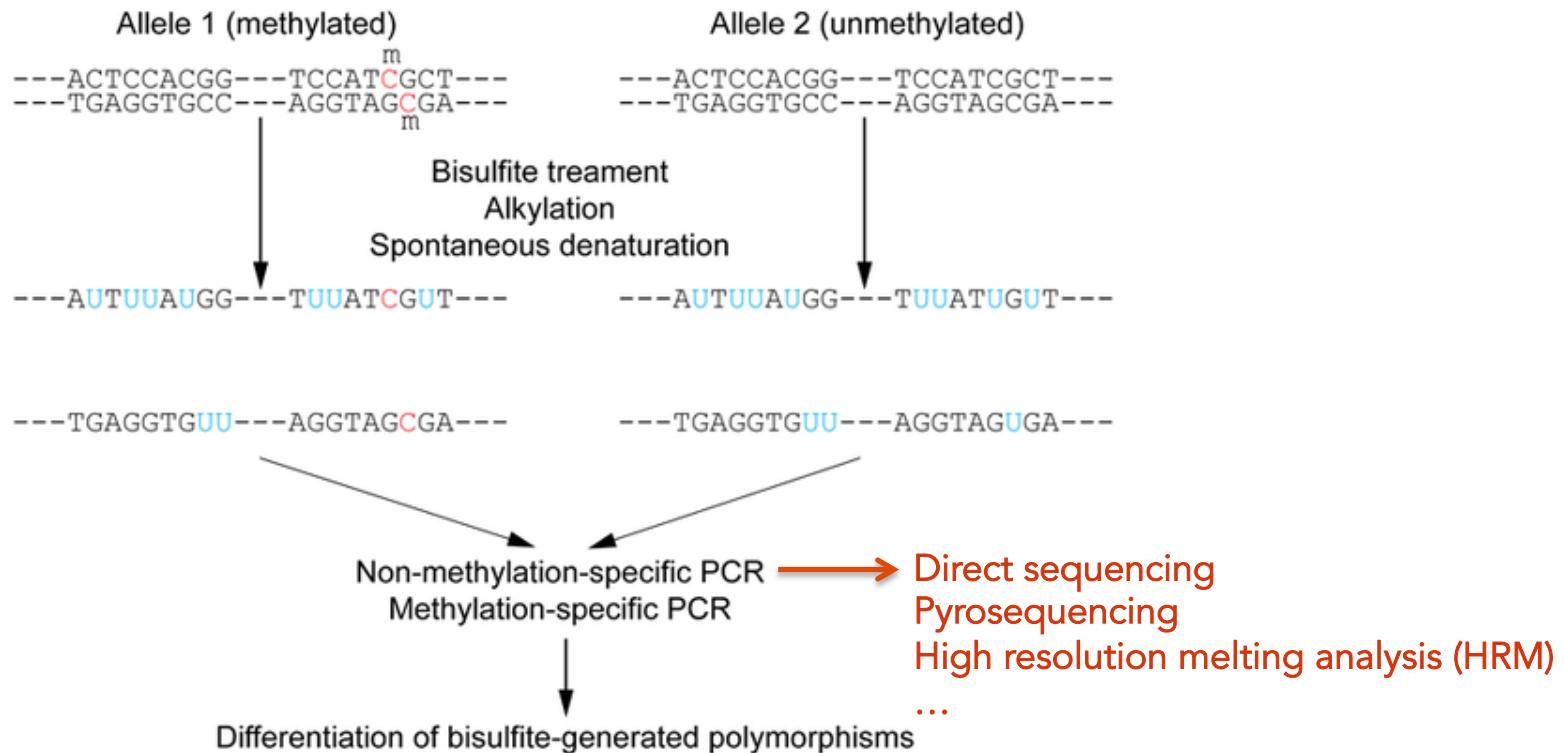
Sequencing today : a broad range of applications thanks to NGS



SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL_LIFE APPLICATIONS

Sequencing today : a broad range of applications thanks to NGS



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

https://en.wikipedia.org/wiki/Bisulfite_sequencing

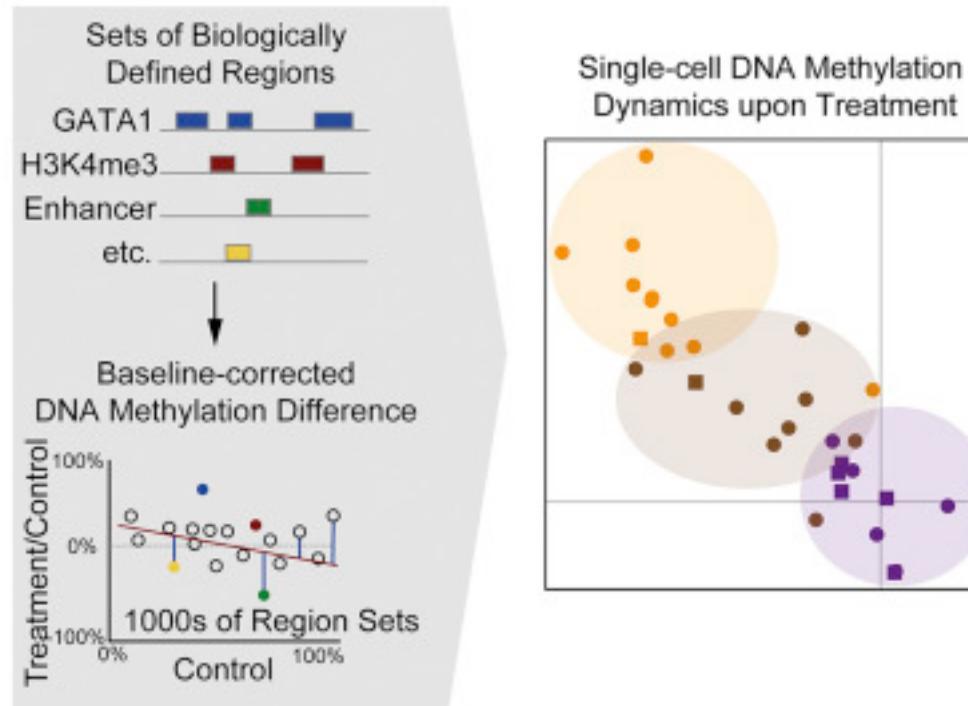
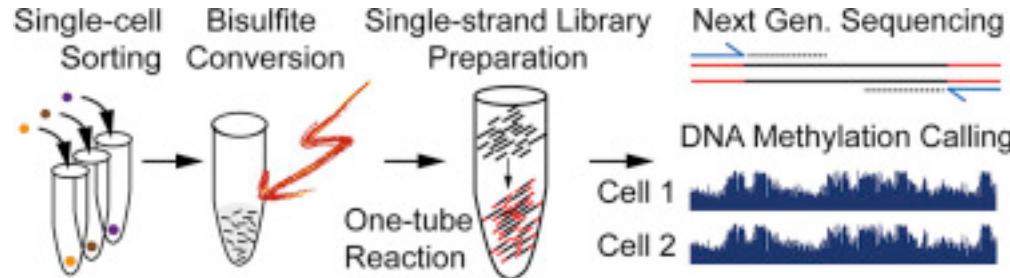
NGS

FATMA GUERFALI

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL_LIFE APPLICATIONS

Sequencing today : a broad range of applications thanks to NGS



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

https://en.wikipedia.org/wiki/Single_cell_sequencing

NGS

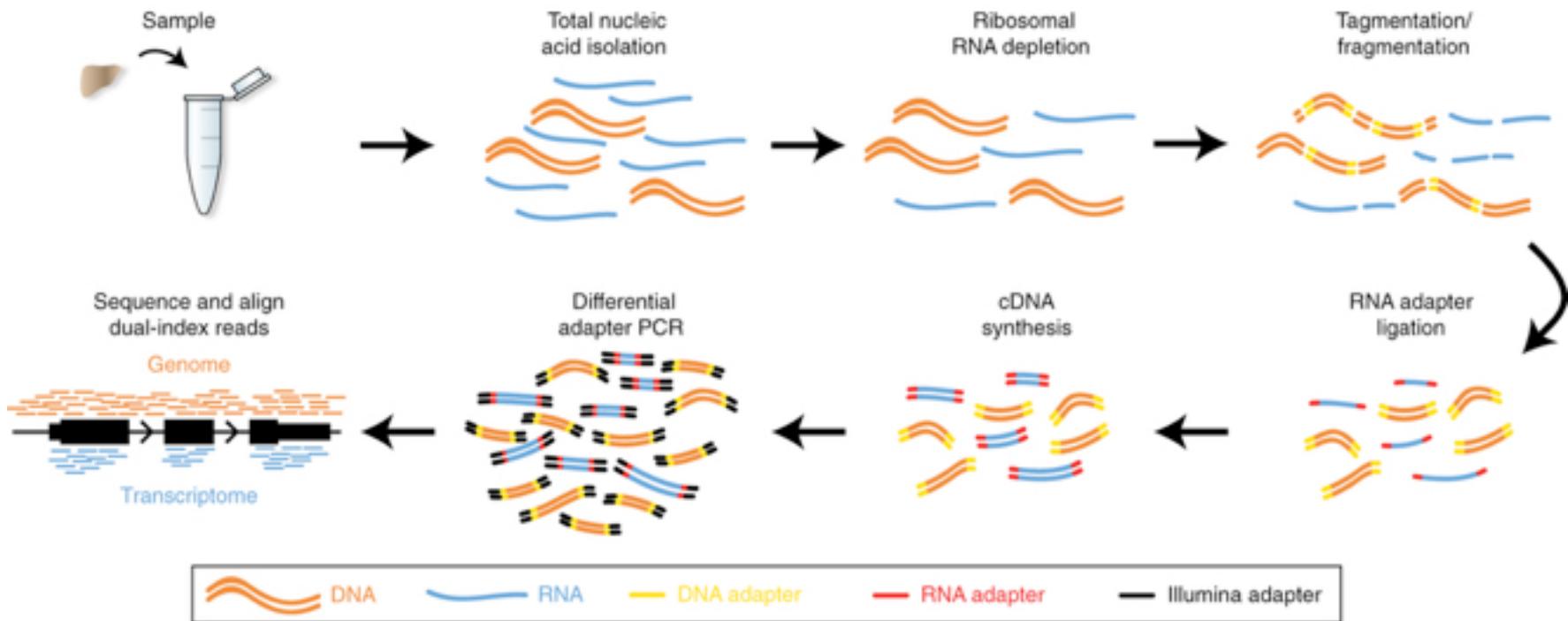
FATMA GUERFALI

SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING : REAL-LIFE APPLICATIONS

- NGS: New improvements

Combined DNA-seq and RNA-seq (« Simul-seq »)



PART
2

OCTOBER 26TH, 2017
IPT COURSE, TUNIS, TUNISIA

(Reuter et al., 2016)

NGS

FATMA GUERFALI

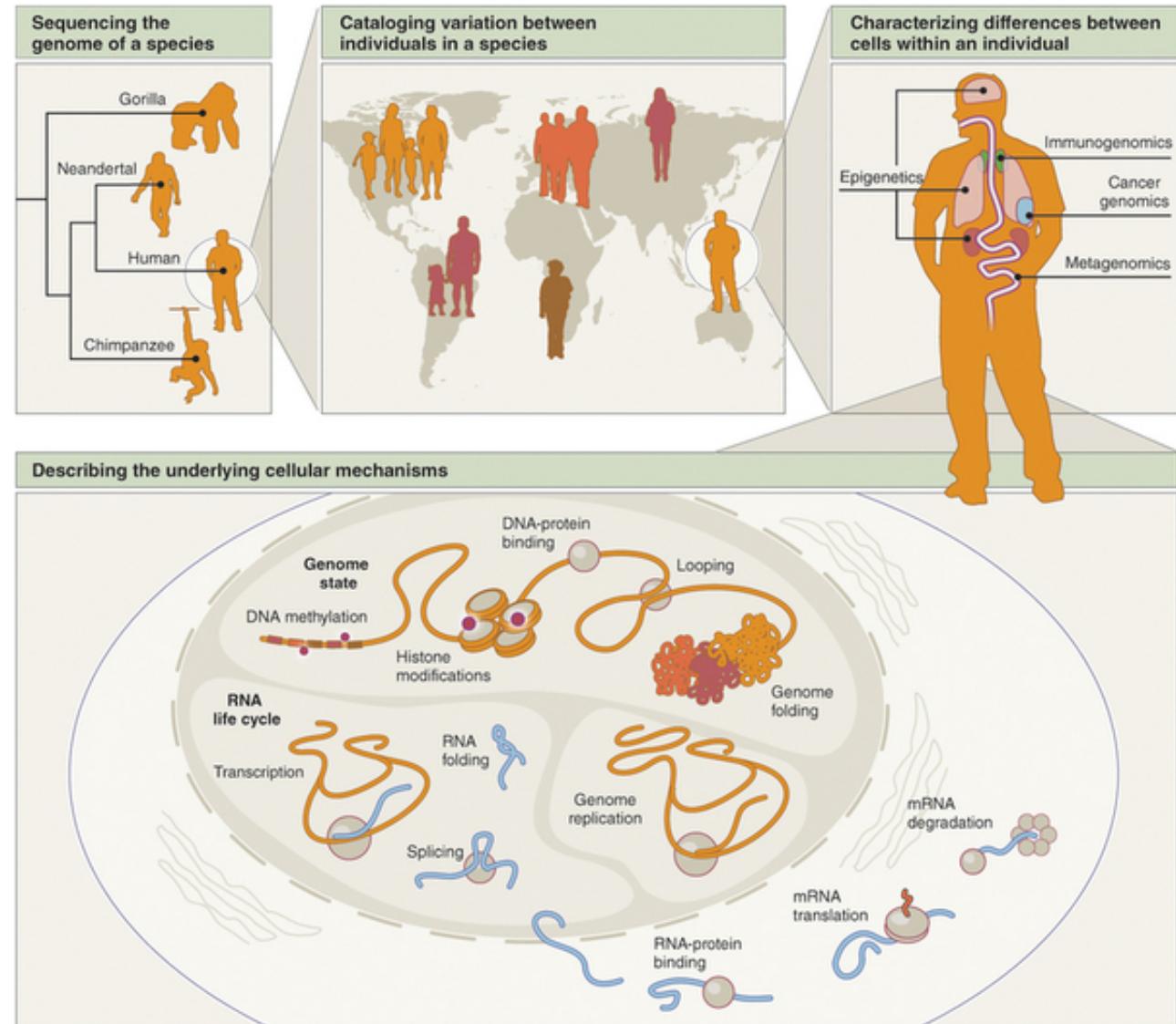
SEQ TECHNOLOGIES

NEXT-GENERATION SEQUENCING: REAL-LIFE APPLICATIONS

- NGS

Where NGS
can help?

at all levels of
information
required !

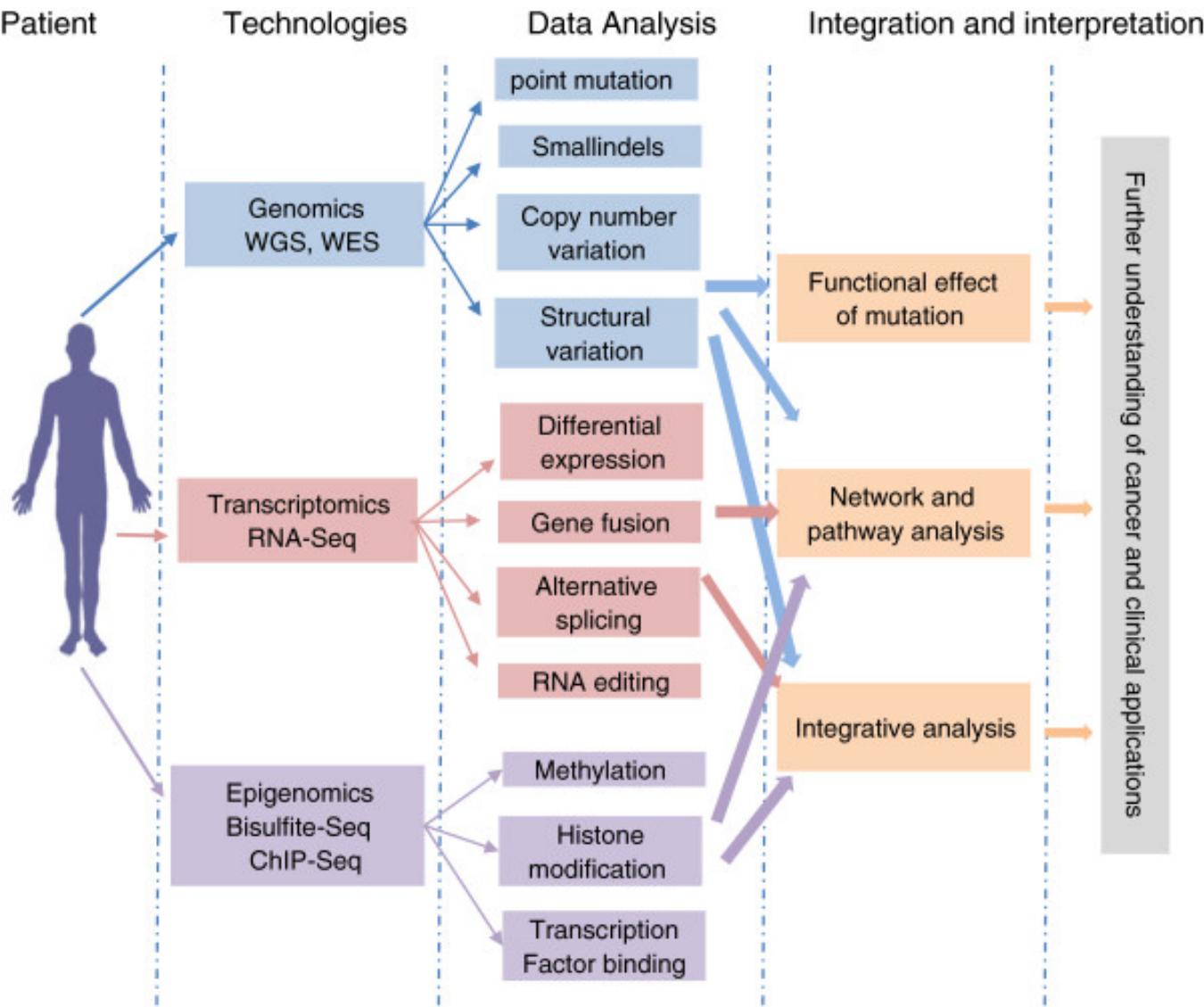


PART
2

SEQ TECHNOLOGIES

TAKE-HOME MESSAGES : NEXT-GENERATION SEQUENCING

Sequencing today : Patient
a broad range of
applications thanks
to NGS



PART
2

OCTOBER 26TH, 2011
IPT COURSE, TUNIS, TUNISIA

(Shyr and Liu, 2013)

FATMA GUERFALI



High-throughput Approach

→ The results of HT experiments provide **starting points** of knowledge, and **are not the end points !!!**

Very useful to:

- Identify genetic variants
- Understand the effect of a particular treatment
- Understand the real interactions between genes expressed in a particular condition
- Drug design

...

But...

NGS or previous technologies?

There is no single method better than any other one, all depends on your question of interest !

One can envision a powerful symbiosis between previous technologies (microarrays...) and NGS technologies.

- Arrays may be best suited in classifying cohorts of samples, such as tumor tissues.
- Once samples of interest are defined, NGS could be used to provide comprehensive deep-sequence analysis of either genomic DNA to identify mutations, or RNA to report differences at the transcriptome level.