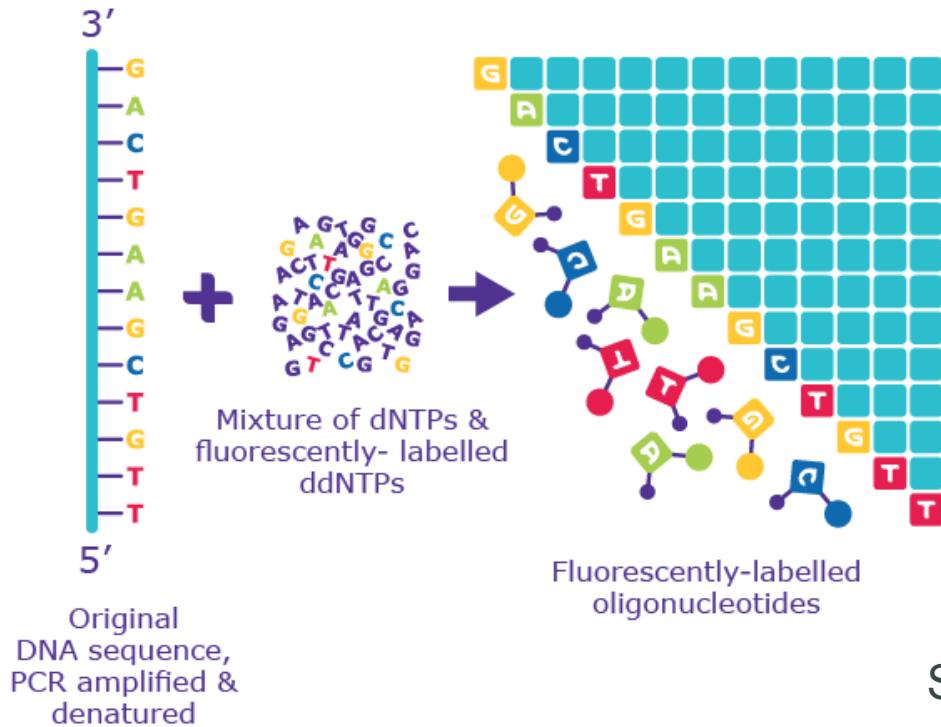**Introduction to Bioinformatics (LF:DSIB01)**

# Week 4 : Next Generation Sequencing: techniques and data

# Nucleic Acid Sequencing History



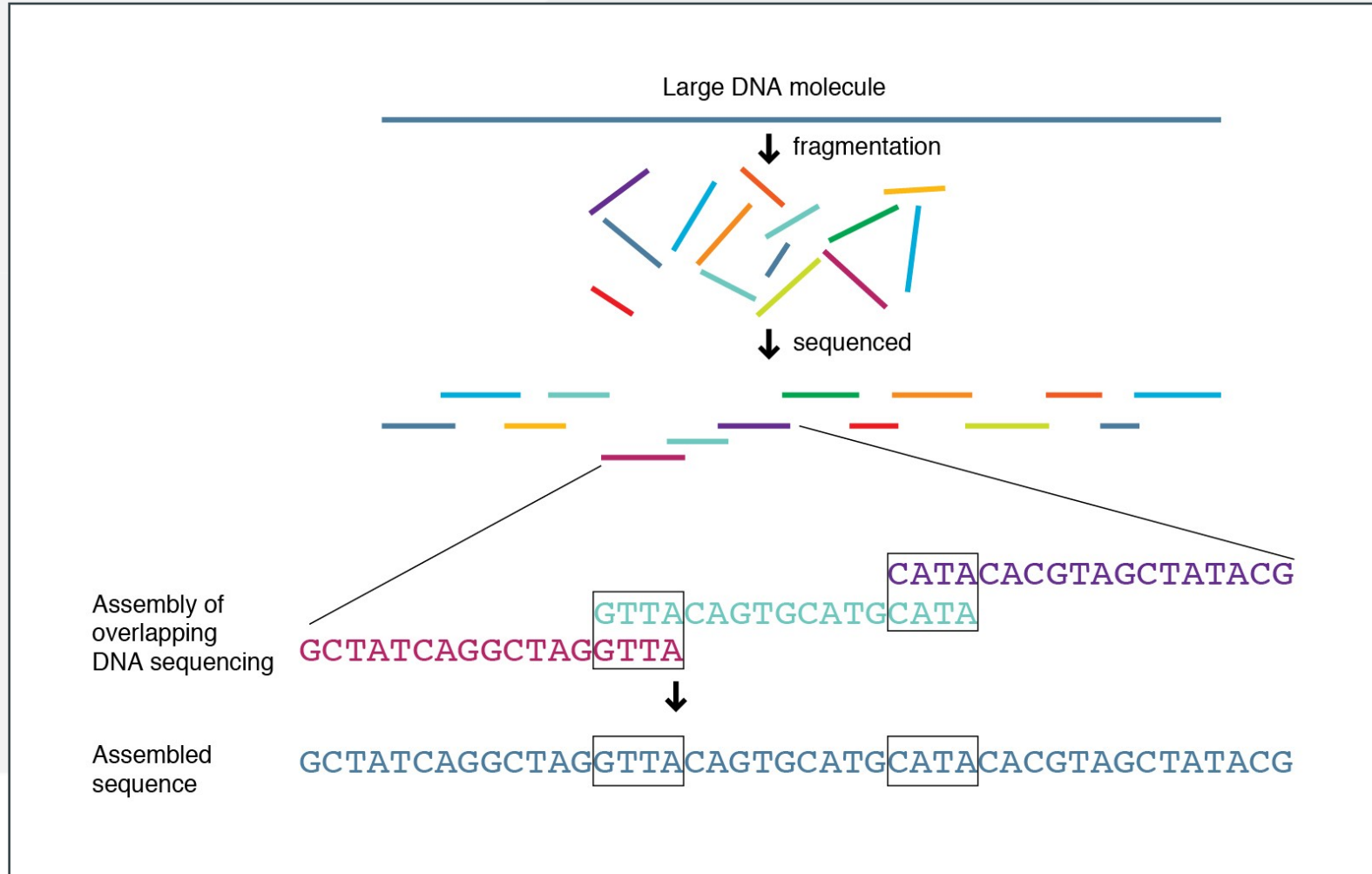**1** PCR with fluorescent, chain-terminating ddNTPs

3'
G
A
C
T
G
A
A
G
C
T
G
T
T
5'

Original DNA sequence, PCR amplified & denatured

Mixture of dNTPs & fluorescently- labelled ddNTPs

Fluorescently-labelled oligonucleotides

**2** Size separation by capillary gel electrophoresis

Large fragments

Small fragments

Laser beam

Photomultiplier

**3** Laser excitation & detection by sequencing machine

Output chromatogram

Sanger Sequencing – Bacteriophage genome sequenced: 1977

CEITEC

# Nucleic Acid Sequencing History



Large DNA molecule

↓ fragmentation

↓ sequenced

Assembly of overlapping DNA sequencing

CATACACGTAGCTATACG

GTTACAGTGCATGCATA

GCTATCAGGCTAGGTTA

↓

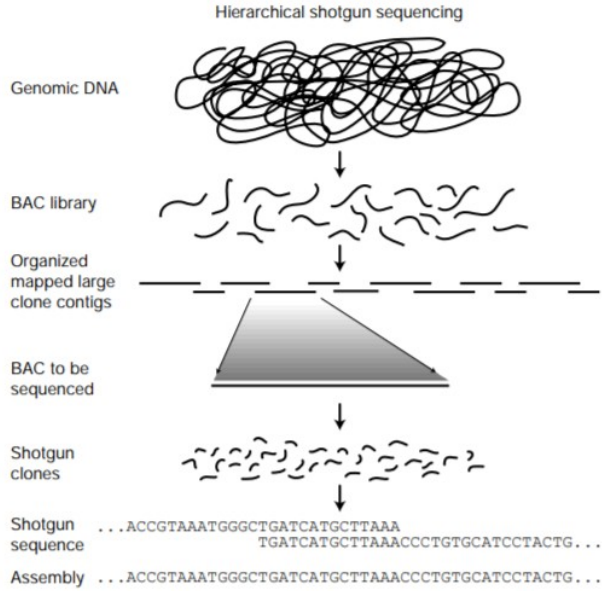Assembled sequence

GCTATCAGGCTAGGTTACAGTGCATGCATACACGTAGCTATACG

1995 – Shotgun Sequencing

## Whole-Genome Random Sequencing and Assembly of *Haemophilus influenzae* Rd

Robert D. Fleischmann, Mark D. Adams, Owen White, Rebecca A. Clayton, Ewen F. Kirkness, Anthony R. Kerlavage, Carol J. Bult, Jean-Francois Tomb, Brian A. Dougherty, Joseph M. Merrick, Keith McKenney, Granger Sutton, Will FitzHugh, Chris Fields,* Jeannine D. Gocayne, John Scott, Robert Shirley, Li-Ing Liu, Anna Glodek, Jenny M. Kelley, Janice F. Weidman, Cheryl A. Phillips, Tracy Spriggs, Eva Hedblom, Matthew D. Cotton, Teresa R. Utterback, Michael C. Hanna, David T. Nguyen, Deborah M. Saudek, Rhonda C. Brandon, Leah D. Fine, Janice L. Fritchman, Joyce L. Fuhrmann, N. S. M. Geoghagen, Cheryl L. Gnehm, Lisa A. McDonald, Keith V. Small, Claire M. Fraser, Hamilton O. Smith, J. Craig Venter†

An approach for genome analysis based on sequencing and assembly of unselected pieces of DNA from the whole chromosome has been applied to obtain the complete nucleotide sequence (1,830,137 base pairs) of the genome from the bacterium *Haemophilus influenzae* Rd. This approach eliminates the need for initial mapping efforts and is therefore applicable to the vast array of microbial species for which genome maps are unavailable. The *H. influenzae* Rd genome sequence (Genome Sequence DataBase accession number L42023) represents the only complete genome sequence from a free-living organism.
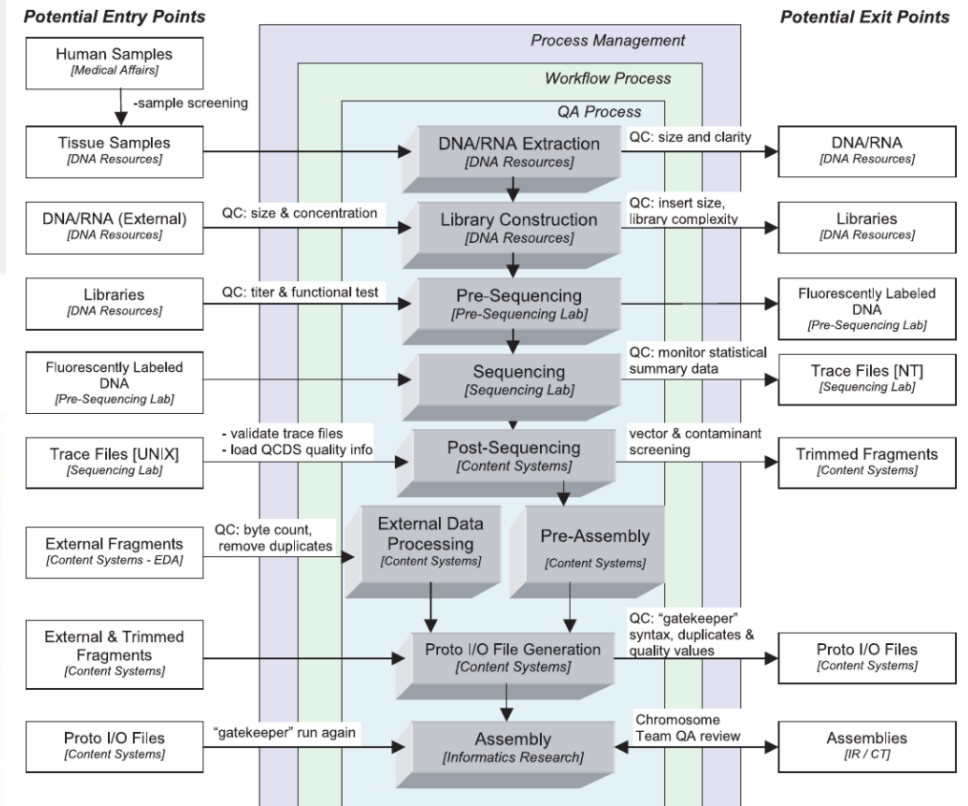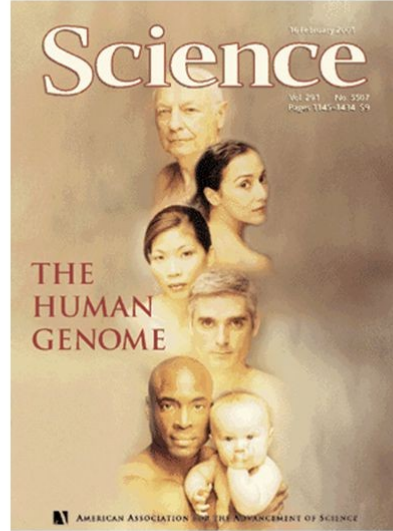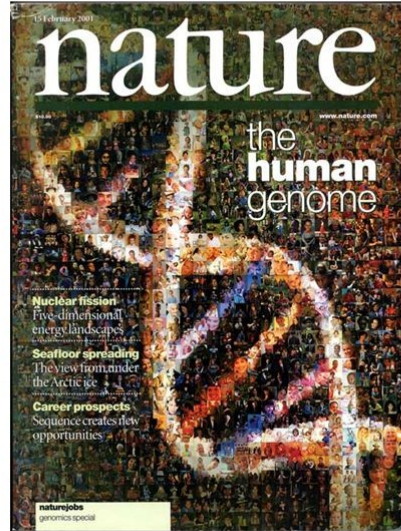
CEITEC

# Nucleic Acid Sequencing History



Figure 2 Idealized representation of the hierarchical shotgun sequencing strategy. A library is constructed by fragmenting the target genome and cloning it into a large-fragment cloning vector; here, BAC vectors are shown. The genomic DNA fragments represented in the library are then organized into a physical map and individual BAC clones are selected and sequenced by the random shotgun strategy. Finally, the clone sequences are assembled to reconstruct the sequence of the genome.

**Genome sequencing**

February 2001 – Publication of the first draft of the human genome





# Initial sequencing and analysis of the human genome

**International Human Genome Sequencing Consortium***

* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

6

2001
Human Genome Sequenced
Cost: ~300M USD

# Nucleic Acid (not-)Sequencing History



Prepare RNA from

"Normal"   Tumor

Label with Biotin

Tumor

Normal

2000s - Microarrays

# Nucleic Acid Sequencing History



**Cost per Genome**



NIH National Human Genome Research Institute

genome.gov/sequencingcosts

Moore's Law

2006 – Solexa Genome Analyser
2007 – Solexa bought by Illumina



Next Generation Sequencing
New Generation Sequencing
**NGS**

**Realistic goal in three-five years**

Sequence the entire human genome in a few days for $1000 (Era of Personal Genomics)

HOWEVER, speed of sequencing does not necessarily mean an **understanding** of the genetic information or DNA structure!

## We are building a research program of 1,000,000+ people.

The *All of Us* Research Program is an ambitious effort to gather health data from one million or more people living in the United States to accelerate research that may improve health.

**OPPORTUNITIES FOR RESEARCHERS**

environment
lifestyle
biology

Research focuses on the intersection of three factors

2015

CEITEC

# NGS Data analysis workflow

Sequencing → Adaptor Trimming → Quality Control → Alignment to Reference

Today's Practical

Next Lecture + Practical

# Long Read Sequencing



Long Reads – Low per read accuracy

# RNA-Seq

# RNA-Seq Analysis

- Alignment to transcriptome or genome (gapped)

- Poly-A selection or Ribosomal RNA depletion

- Can be used to quantify RNA, or to identify structural differences (e.g. splicing)

- Usual downstream analysis: Fold Change between conditions

# Immunoprecipitation based techniques



ChIP-Seq : DNA Binding Proteins

CLIP-Seq : RNA Binding Proteins

# CLIP and complementary methods

Markus Hafner [1], Maria Katsantoni [2,3], Tino Köster [4], James Marks [1],
Joyita Mukherjee [5,6], Dorothee Staiger [4], Jernej Ule [5,6,7 ✉] and Mihaela Zavolan [2,3]

## a

| | No pretreatment | | | | 4SU treatment | |
|---|---|---|---|---|---|---|
| Pretreatment | | | | | | |
| Cross-linking | | UVC | | | UVA/B | |
| Mild RNase digestion | HITS-CLIP/ CRAC | iCLIP | irCLIP | eCLIP | PAR-CLIP | Proximity-CLIP [a] |
| Immuno-precipitation | | | | | | |
| 3' Adapter ligation | | | | | | Mass spectrometry / RNA-seq |
| Gel purification | | | | | | |
| Protein digestion | | | | | | |
| 5' Adapter ligation | | x | x | x | | |
| Reverse transcription | | | | | | |
| Adapter ligation/ circularization | x | | | | x | x |
| Linearization | x | x | | x | x | x |
| PCR amplification | | | | | | |
| High-throughput sequencing | | | | | | |
| Peak calling | Coverage/ mutation | Reverse transcription stops | Reverse transcription stops | Reverse transcription stops | T to C mutations | T to C mutations |

## b

RBP–ADARcd → Standard RNA isolation → TruSeq library preparation → Standard RNA-seq → A to G mutations

A → I

## c

Cross-link site

RNA

UMI Adapter

Possible cDNAs:
- Truncation
- Read-through
- Substitution/insertion
- Deletion

## d

**PCR amplification**

**Duplicate removal**

**Adapter removal**

# Rib



Ribo-seq

RNA-seq

Nuclease digestion
Monosome isolation

RNA purification

Ribosome footprints

Library construction

Deep sequencing

Actively translated regions show strong 3-nt periodicity

Novel peptide identification

Translation efficiency

RNA purification
Fragmentation

Library construction

Deep sequencing

RNA-seq does not show 3-nt periodicity

$$\text{Translation efficiency} = \frac{\text{Ribo-seq levels}}{\text{RNA-seq levels}}$$

Adapted from Hsu *et al.* 2016

Identification of actively translated RNA

Exact position of Ribosome

Rich conditions

Starvation

ABP140   MET7   SSP2   PUS7

# ATAC-Seq



Identification of accessible chromatin areas

# Shape-Seq

# … and others

**RNA Transcription**
- Chromatin Isolation by RNA Purification (ChIRP-Seq)
- Global Run-on Sequencing (GRO-Seq)
- Ribosome Profiling Sequencing (Ribo-Seq)/ARTseq™
- RNA Immunoprecipitation Sequencing (RIP-Seq)
- High-Throughput Sequencing of CLIP cDNA library (HITS-CLIP) or
- Crosslinking and Immunoprecipitation Sequencing (CLIP-Seq)
- Photoactivatable Ribonucleoside–Enhanced Crosslinking and Immunoprecipitation (PAR-CLIP)
- Individual Nucleotide Resolution CLIP (iCLIP)
- Native Elongating Transcript Sequencing (NET-Seq)
- Targeted Purification of Polysomal mRNA (TRAP-Seq)
- Crosslinking, Ligation, and Sequencing of Hybrids (CLASH-Seq)
- Parallel Analysis of RNA Ends Sequencing (PARE-Seq) or
- Genome-Wide Mapping of Uncapped Transcripts (GMUCT)
- Transcript Isoform Sequencing (TIF-Seq) or
- Paired-End Analysis of TSSs (PEAT)

**RNA Structure**
- Selective 2'-Hydroxyl Acylation Analyzed by Primer Extension Sequencing (SHAPE-Seq)
- Parallel Analysis of RNA Structure (PARS-Seq)
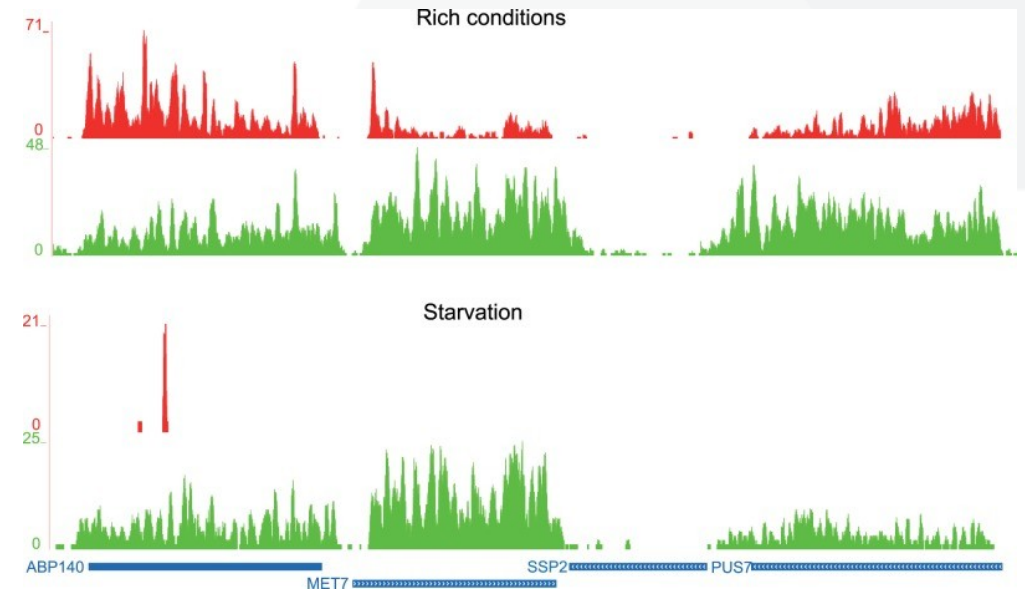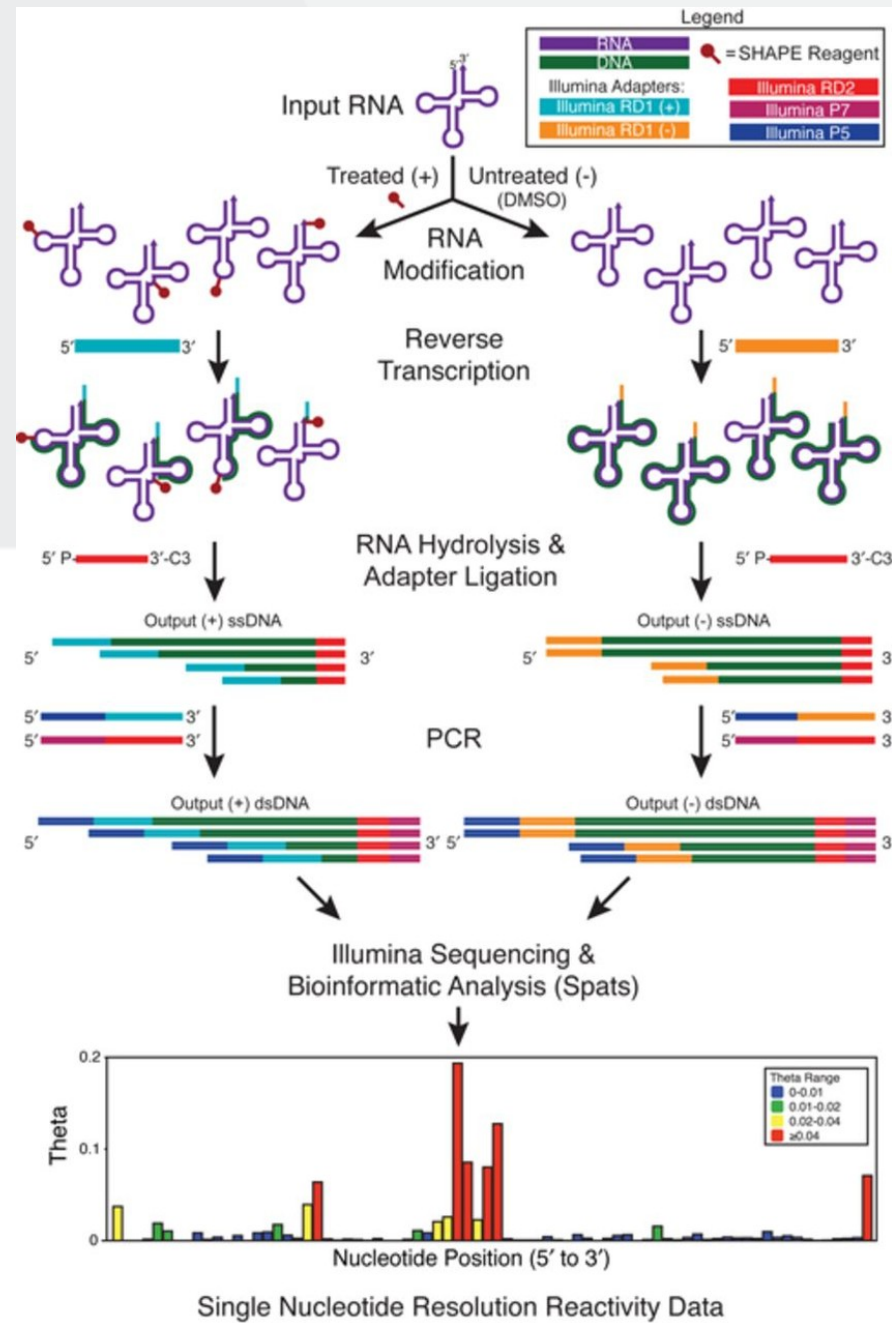- Fragmentation Sequencing (FRAG-Seq)
- CXXC Affinity Purification Sequencing (CAP-Seq)
- Alkaline Phosphatase, Calf Intestine-Tobacco Acid Pyrophosphatase Sequencing (CIP-TAP)
- Inosine Chemical Erasing Sequencing (ICE)
- m6A-Specific Methylated RNA Immunoprecipitation Sequencing (MeRIP-Seq)

**Low-Level RNA Detection**
- Digital RNA Sequencing
- Whole-Transcript Amplification for Single Cells (Quartz-Seq)
- Designed Primer–Based RNA Sequencing (DP-Seq)
- Switch Mechanism at the 5' End of RNA Templates (Smart-Seq)
- Switch Mechanism at the 5' End of RNA Templates Version 2 (Smart-Seq2)
- Unique Molecular Identifiers (UMI)
- Cell Expression by Linear Amplification Sequencing (CEL-Seq)
- Single-Cell Tagged Reverse Transcription Sequencing (STRT-Seq)

**Low-Level DNA Detection**
- Single-Molecule Molecular Inversion Probes (smMIP)
- Multiple Displacement Amplification (MDA)
- Multiple Annealing and Looping–Based Amplification Cycles (MALBAC)
- Oligonucleotide-Selective Sequencing (OS-Seq)
- Duplex Sequencing (Duplex-Seq)

**DNA Methylation**
- Bisulfite Sequencing (BS-Seq)
- Post-Bisulfite Adapter Tagging (PBAT)
- Tagmentation-Based Whole Genome Bisulfite Sequencing (T-WGBS)
- Oxidative Bisulfite Sequencing (oxBS-Seq)
- Tet-Assisted Bisulfite Sequencing (TAB-Seq)
- Methylated DNA Immunoprecipitation Sequencing (MeDIP-Seq)
- Methylation-Capture (MethylCap) Sequencing or
- Methyl-Binding-Domain–Capture (MBDCap) Sequencing
- Reduced-Representation Bisulfite Sequencing (RRBS-Seq)

**DNA-Protein Interactions**
- DNase I Hypersensitive Sites Sequencing (DNase-Seq)
- MNase-Assisted Isolation of Nucleosomes Sequencing (MAINE-Seq)
- Chromatin Immunoprecipitation Sequencing (ChIP-Seq)
- Formaldehyde-Assisted Isolation of Regulatory Elements (FAIRE-Seq)
- Assay for Transposase-Accessible Chromatin Sequencing (ATAC-Seq)
- Chromatin Interaction Analysis by Paired-End Tag Sequencing (ChIA-PET)
- Chromatin Conformation Capture (Hi-C/3C-Seq)
- Circular Chromatin Conformation Capture (4-C or 4C-Seq)
- Chromatin Conformation Capture Carbon Copy (5-C)

**Sequence Rearrangements**
- Retrotransposon Capture Sequencing (RC-Seq)
- Transposon Sequencing (Tn-Seq) or Insertion Sequencing (INSeq)
- Translocation-Capture Sequencing (TC-Seq)

Link

CEITEC

@CEITEC_Brno

www.ceitec.eu

CEITEC