COURSE ON

**2ND EDITION**

# COMPUTATIONAL BIOSCIENCES USING HPC SYSTEMS

6, 7, 8 FEBRUARY 2024

@ NOVA SCHOOL OF
SCIENCE AND TECHNOLOGY

UCIBIO · i4HB Institute for Health and Bioeconomy · INCD Infraestrutura Nacional de Computação Distribuída · LAQV requimte · LIP · EURO PORTUGAL · PRR Plano de Recuperação e Resiliência

# MODULE 1 - OMICS
## Transcriptome assembly

Pedro M. Costa
(pmcosta@fct.unl.pt)

NOVA SCHOOL OF
SCIENCE & TECHNOLOGY
DEPARTMENT OF LIFE SCIENCES

- RNA-Seq is quantitative.

- Can yield more than 100K validated transcripts.

- It is usually non-targeted (i.e. "transcriptome-wide")

- Depending on sequencing depth and length, may not yield full-length mRNAs.

For instance: 10-20 M reads, 150 bp single-end reads is the basic for expression analysis IF the transcriptome is reasonably annotated (unlikely in marine organisms). 100 M reads, 150-300 bp paired end is great for quantification AND characterisation of mRNAs (but it is also very expensive and challenging computationally).



**MiniSeq System**
Power and simplicity for targeted sequencing.

**MiSeq Series**
Small genome and targeted sequencing.

**NextSeq Series**
Everyday genome, exome transcriptome sequencing, and more.

**HiSeq Series**
Production-scale genome, exome, transcriptome sequencing, and more.

**HiSeq X Series**
Population- and production-scale human whole-genome sequencing.

**NovaSeq Series**
Population- and production-scale genome, exome, transcriptome sequencing, and more.

**Figure 6: Sequencing Systems for Virtually Every Scale** — Illumina offers innovative NGS platforms that deliver exceptional data quality and accuracy over a wide scale, from small benchtop sequencers to production-scale sequencing systems.
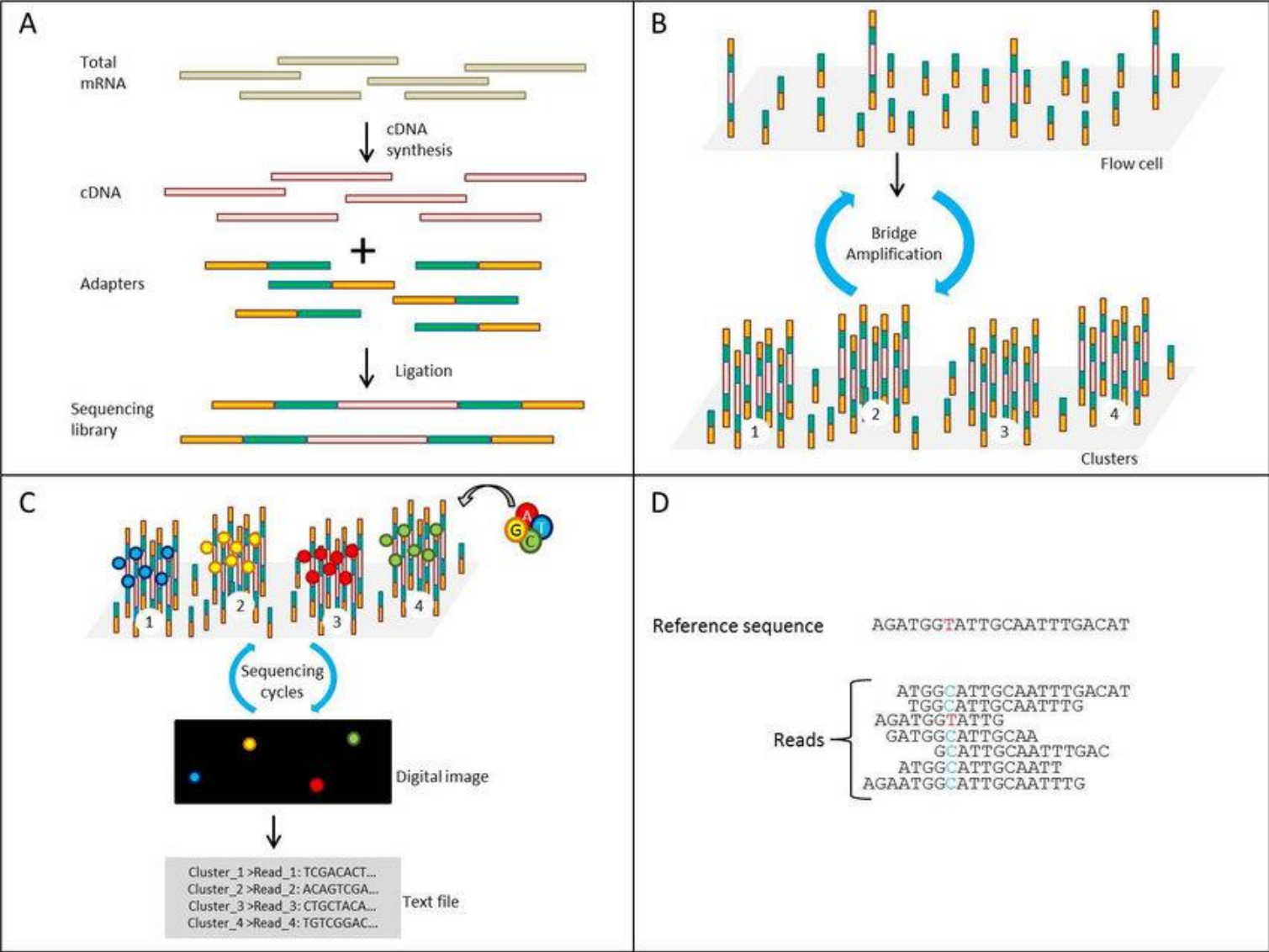
Martins et al. (2019). Int. J. Human Environ. Health 16, 4718. (doi: 10.3390/ijerph16234718)

This is called Next-Generation Sequencing (NGS). A similar process applies to genomes.
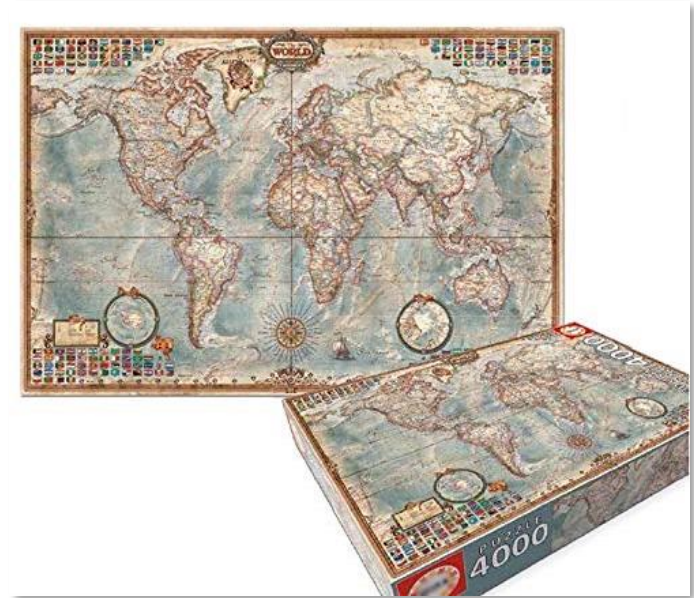
P. Juarez©

Model organisms such as humans, rats, mice, zebrafish and a few other benefit from a high degree of genomic resources, including available transcriptomes/genomes against which RNA-Seq raw data can be **mapped**.



Non-model/novel organisms have limited or null information on gene, peptide or mRNA sequences. In these cases, the transcriptome needs to be *de novo* **assembled**. Pretty much like a 10K+ pieces without a reference photo…

Genome: ACGATGCTCAGACCCCCCCCC

Short reads: ACGATGCTCAGA    CTCAGACCC    AGACCCC    CCCCCCC

k-mers:
ACGAT          CTCAG          AGACC          CCCCC
CGATG          TCAGA          GACCC          CCCCC
GATGC          CAGAC          ACCCC          CCCCC
ATGCT          AGACC
TGCTC          GACCC
GCTCA
CTCAG
TCAGA

Commonly, k = 25 is the target

De Bruijn graph:



Assembled Contigs: ACGATGCTCAGACCCC

## Two Paradigms for Assembly



**Assembly of Large Genomes using Second Generation Sequencing**
Schatz MC, Delcher AL, Salzberg SL (2010) *Genome Research* 20, 1165-73.

6

**a**

Read set

Extend in *k*-mer space and break ties

>a121:len = 5,845
>a122:len = 2,560
>a123:len = 4,443
>a124:len = 48
>a126:len = 66

Linear sequences

**b**

Overlap linear sequences by overlaps of *k* − 1 to build graph components

*k* − 1
*k* − 1
*k* − 1
*k* − 1
*k* − 1

**c**

De Bruijn graph (*k* = 5)

Compacting

Compact graph

Finding paths

Compact graph with reads

Extracting sequences

...CTTCGCAA...TGATCGGAT...
...ATTCGCAA...TCATCGGAT...   Transcripts

**Figure 1** Overview of Trinity. (**a**) Inchworm assembles the read data set (short black lines, top) by greedily searching for paths in a *k*-mer graph (middle), resulting in a collection of linear contigs (color lines, bottom), with each *k*-mer present only once in the contigs. (**b**) Chrysalis pools contigs (colored lines) if they share at least one *k* − 1-mer and if reads span the junction between contigs, and then it builds individual de Bruijn graphs from each pool. (**c**) Butterfly takes each de Bruijn graph from Chrysalis (top), and trims spurious edges and compacts linear paths (middle). It then reconciles the graph with reads (dashed colored arrows, bottom) and pairs (not shown), and outputs one linear sequence for each splice form and/or paralogous transcript represented in the graph (bottom, colored sequences).

**Training set**

Raw Fastq data set

?M reads (<10)

36 bp paired-end sequencing
Illumina platform
(but we will work only with
L)

GEO GSM767958

Today, we will put aside
*quantification* and focus on
*identification*



Cebrià et al. (2016). Regeneration and Growth as Modes of Adult Development: the Platyhelminthes as a Case Study. Doi: 10.1007/978-3-7091-1871-9_4

- De novo transcriptome assembly from *SmedIllumina_R1.fastq.gz* using Trinity

- Basic quality assessment using *TrinityStats* (*Nx* analyses)

- Predict ORFs using *Transdecoder*

- Annotate the resulting ORFs using *Pblast* (homology-matching):
  - -Against UniProt
  - -Against Uniprot (Human proteome only)

Inter-specific RNA contamination

Fortunately there are tools to clean and normalise data before assembly and to check its quality afterwards…

Adapters and low-quality reads must be removed

Sequences too short

Over-represented genes (i.e. high vs low expression genes)

Heterogenous representation of the transcriptome

Gene-dense genomes resulting in a large number of transcriptional variants (for instance, resulting from chromosome duplication, cryptic genes…)

Transcripts

E.g., Transdecoder → Predicted Open Reading Frames (ORFs)

E.g., Blast; HMMER → Annotation

Available databases include UniProt, RefSeq, Pfam (conserved domains)...

Topology-based pathway enrichment

E.g., Blast2GO; DAVID; several R packages → Gene ontology / Gene enrichment

Available databases include KEGG, Reactome...

Q: Can we quantify expression when dealing with *de novo* assembly? Yes!

First we use a programme like Trinity to assembly the transcriptome of our tart organism/tissue/organ (this can/should be done using several samples). We will then use this transcriptome as reference for mapping using, e.g. Kallisto.

Sample1.fastqc
Sample2.fastqc
Sample3.fastqc

Trinity → Reference transcriptome

Kallisto + R

DEGs → BLAST → Annotated DEGs

*Transdecoder* workflow

Homology-matching can be done with cDNA or AA sequences. However, the later can filter sequences by isolating coding from non-coding and reducing variability.

BLAST (Basic Local Alignment Search Tool) has command line versions that enable batch searches for a large number of queries.

https://www.uniprot.org/

You can easily customise and download a database from UniProt's website

https://www.genome.jp/kegg/