

Workshop UCIBIO

MODULE 2 Transcriptome assembly

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

13 / 9 / 2022

RNA-Seq. Basics

- 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).

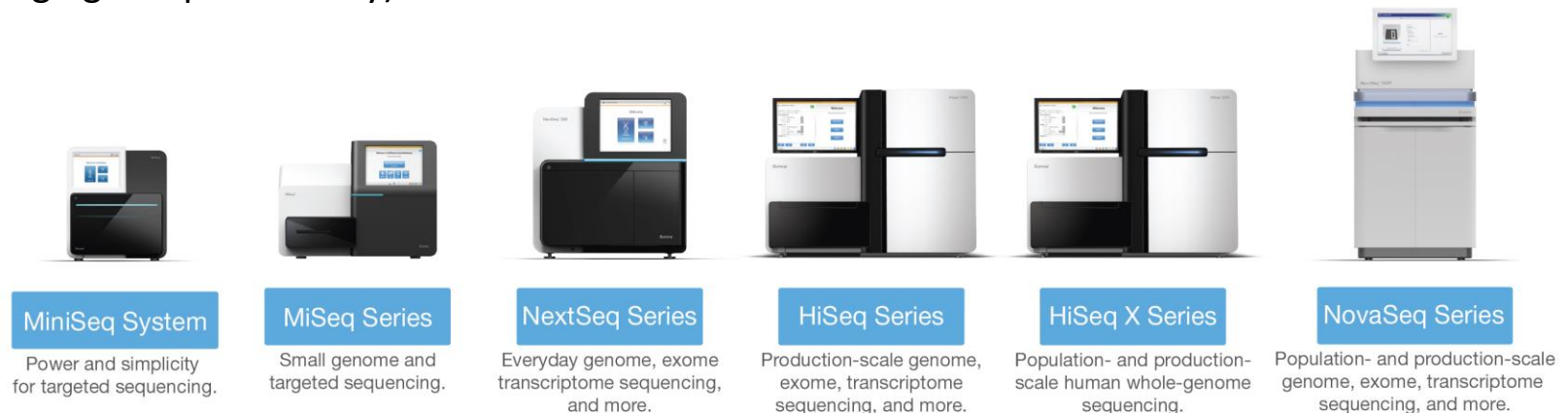
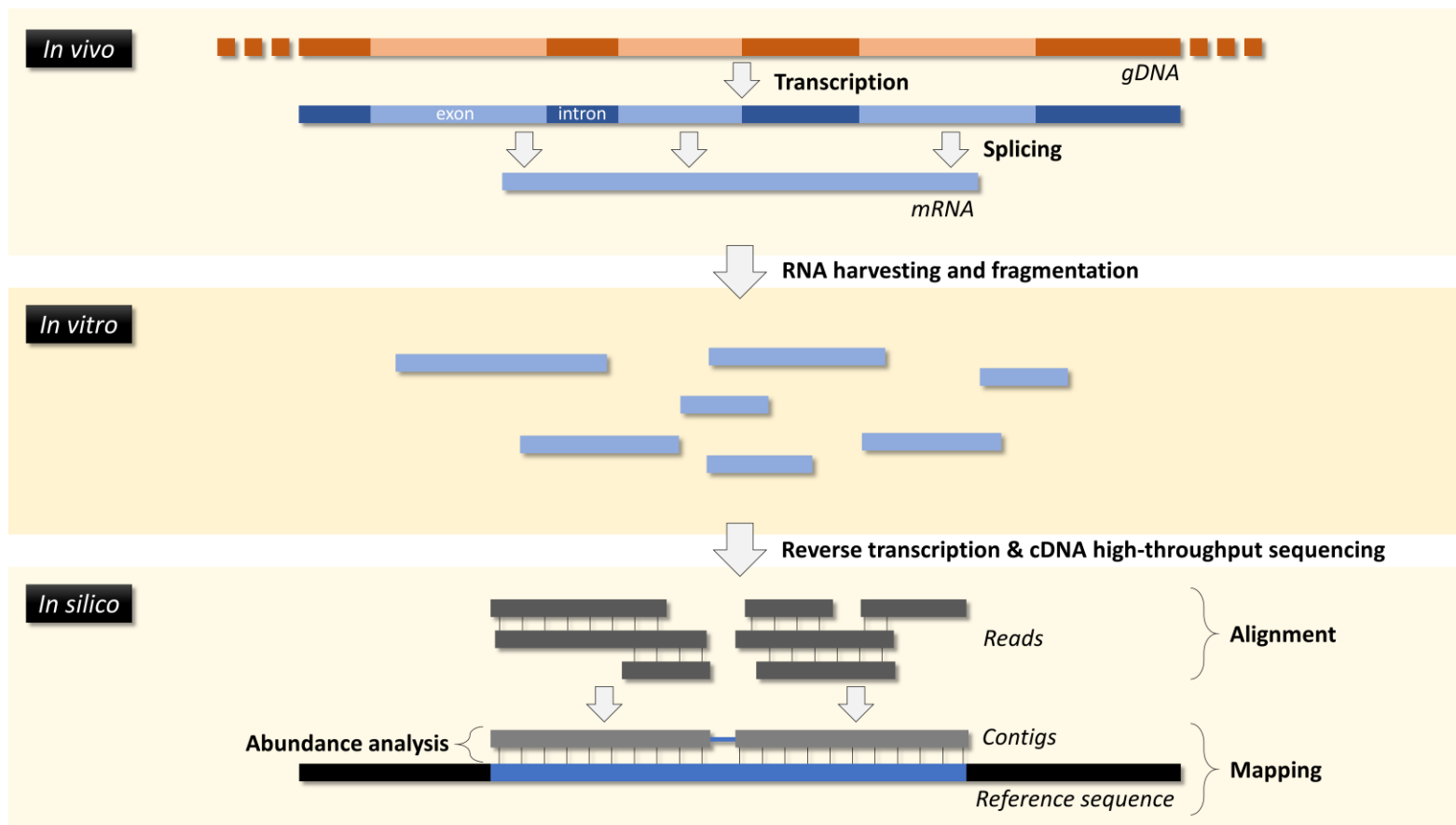


Figure 6: Sequencing Systems for Virtually Every Scale—Illumina offers innovative NGS platforms that deliver exceptional data quality and accuracy over a wide scale, from small benchtop sequencers to production-scale sequencing systems.

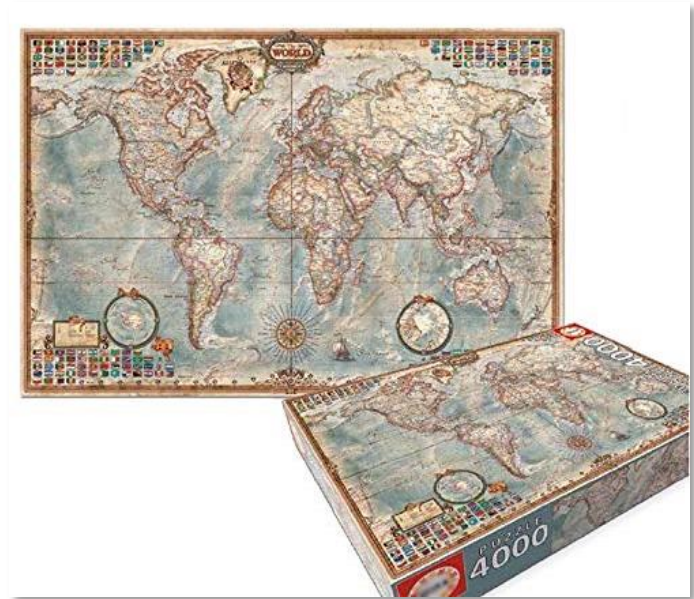


Martins et al. (2019). Int. J. Human Environ. Health 16, 4718. (doi: 10.3390/ijerph16234718)

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

Transcriptome mapping and assembly. *Model vs. non-model organisms*

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...



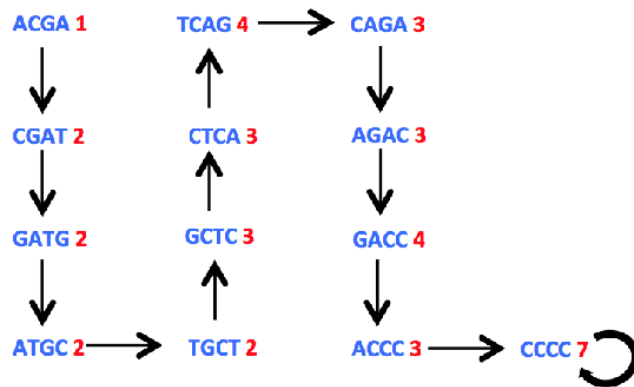
Transcriptome mapping and assembly. *K*-mer

Genome: ACGATGCTCAGACCCCCCCC
Short reads: ACGATGCTCAGA CTCAGACCC AGACCCC CCCCCC

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

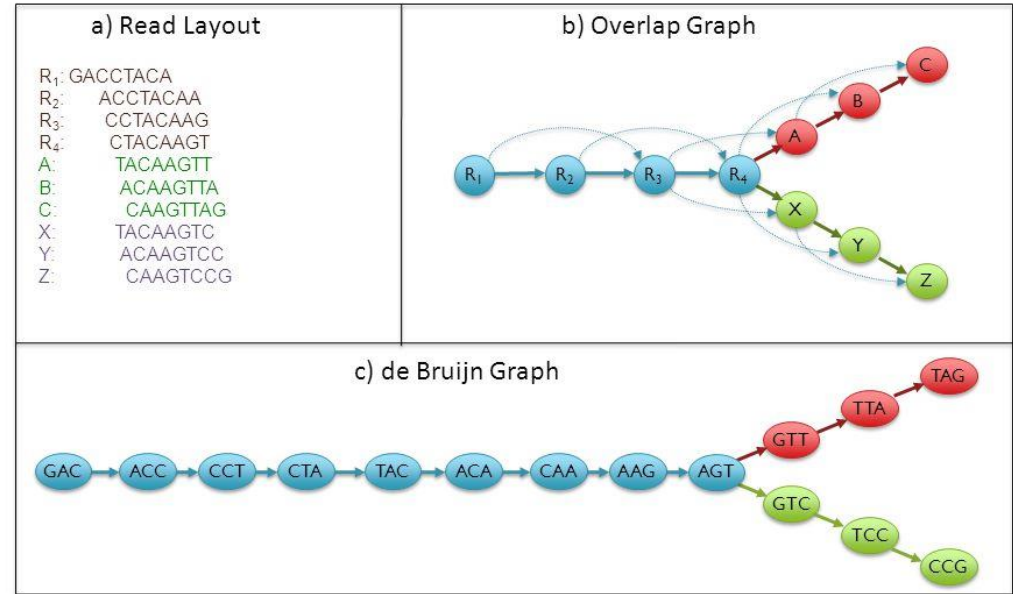
Commonly, $k = 25$ is the target

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.

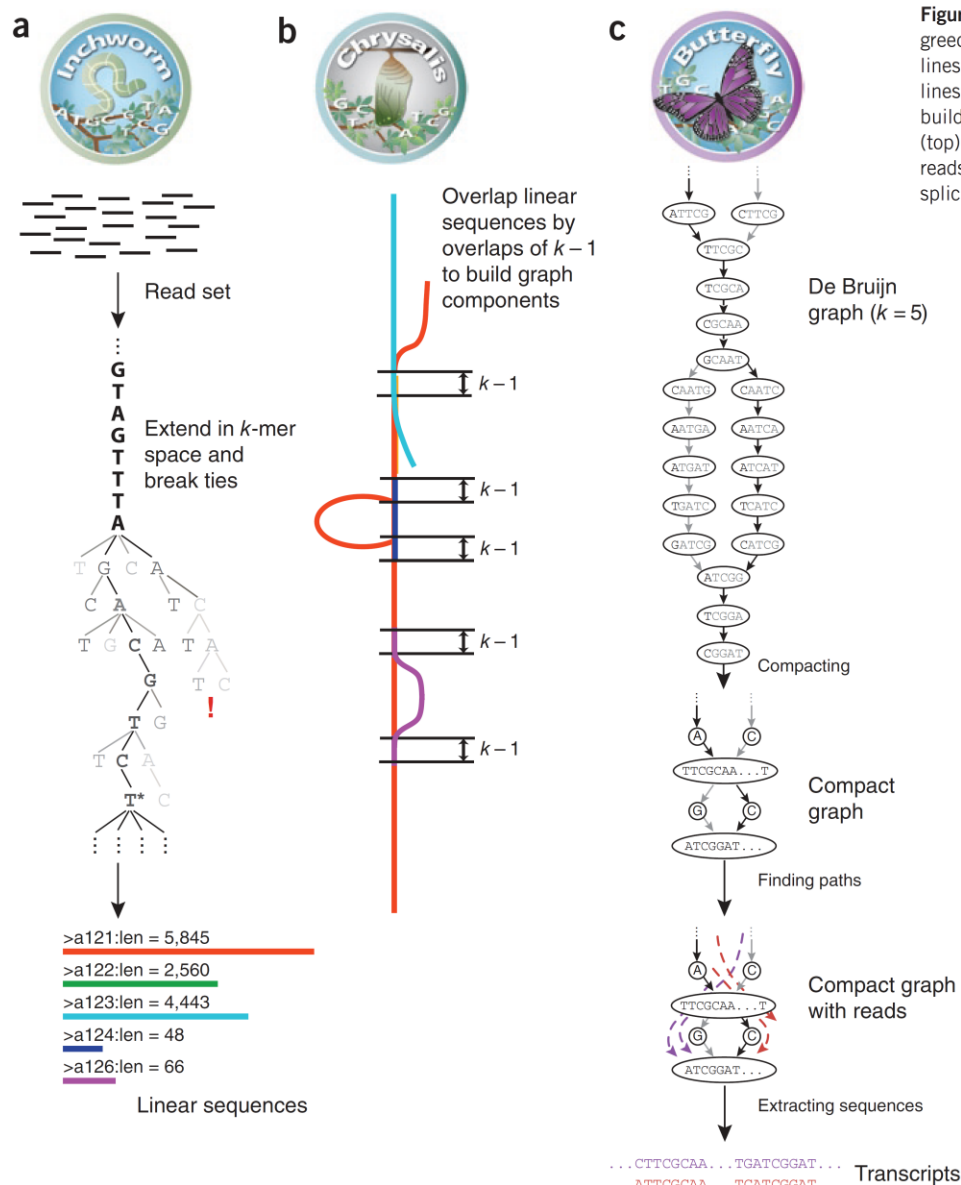


Figure 1 Overview of Trinity. (a) Inchworm assembles the read data set (short black lines, top) by greedily searching for paths in a k -mer graph (middle), resulting in a collection of linear contigs (color lines, bottom), with each k -mer present only once in the contigs. (b) Chrysalis pools contigs (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).

Grabherr et al. (2011) Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol* 29, 644-652. <https://doi.org/10.1038/nbt.1883>

Trinity. Constraints to be considered in novo assembly

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

- Inter-specific RNA contamination

- Adapters and low-quality reads must be removed

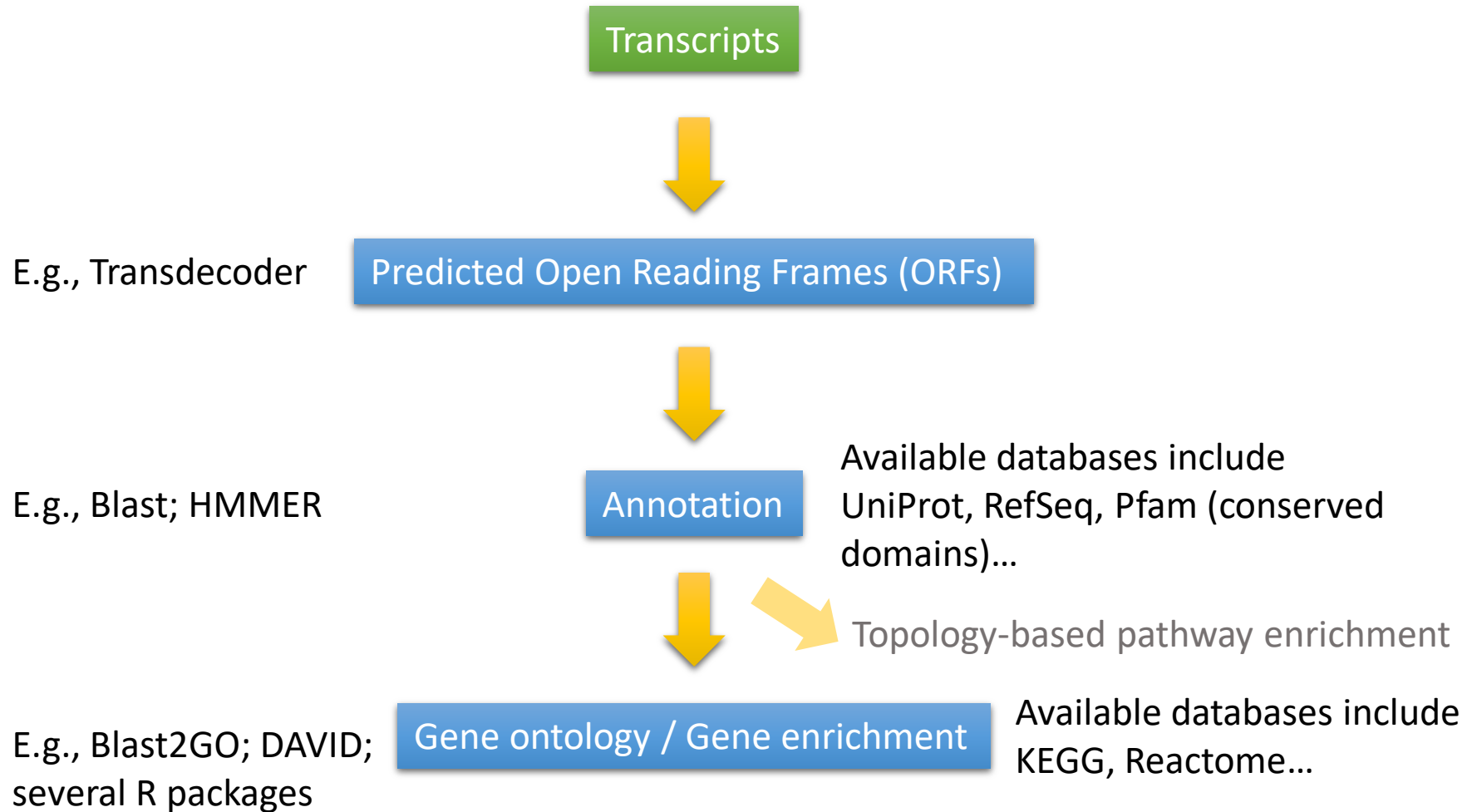
- Sequences too short

- 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...)

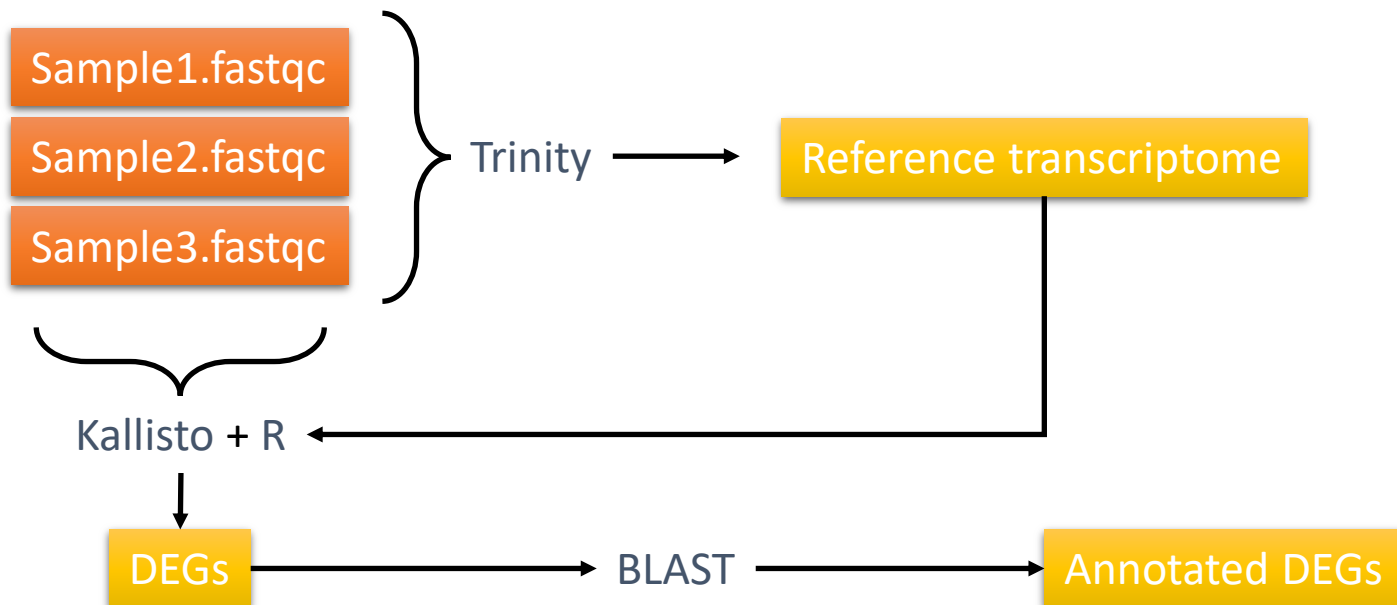
Transcriptome annotation. *Standard procedure*



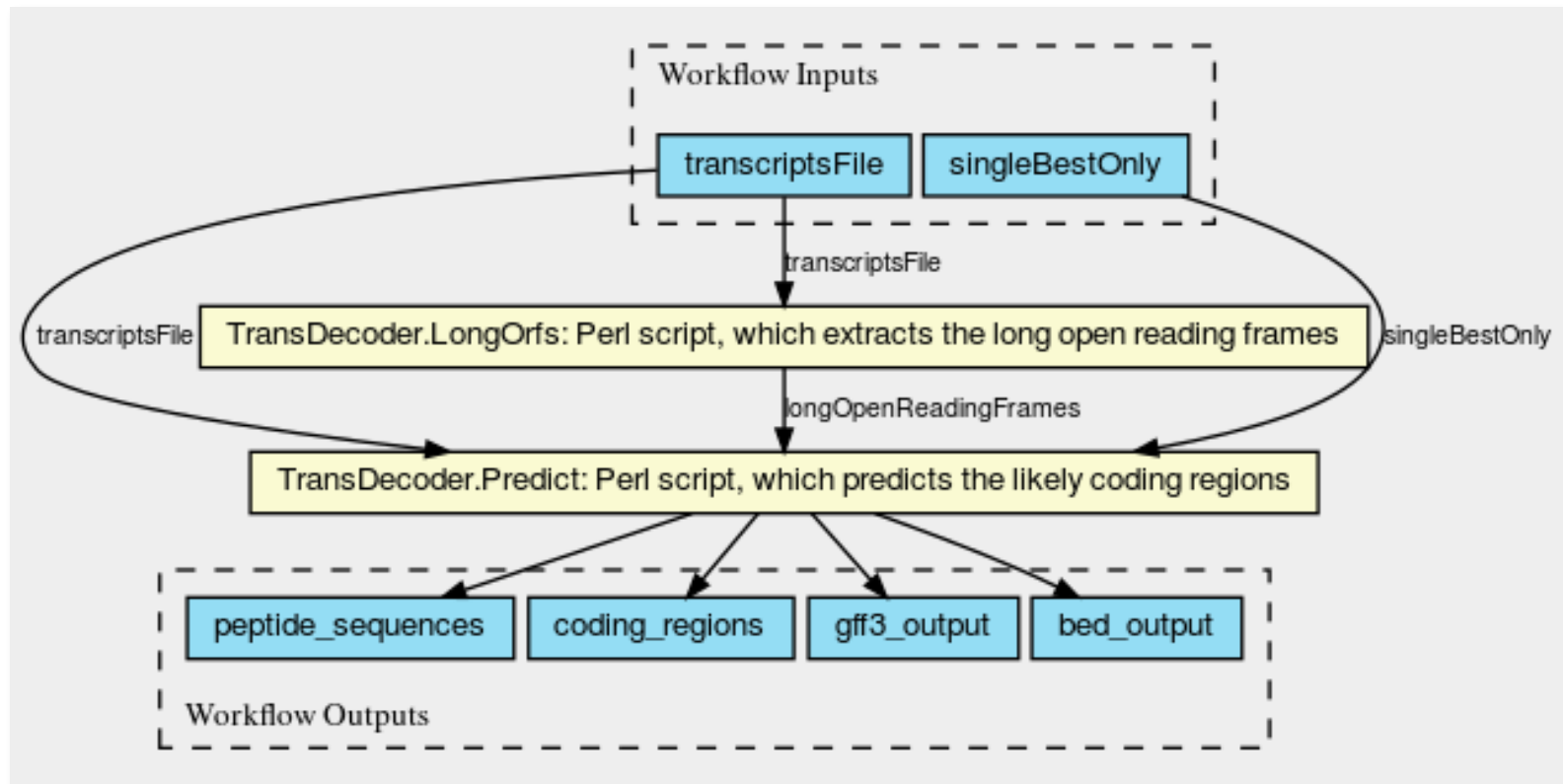
Transcriptome annotation. Standard procedure

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.



Pre-annotation tools. Predicting ORFs



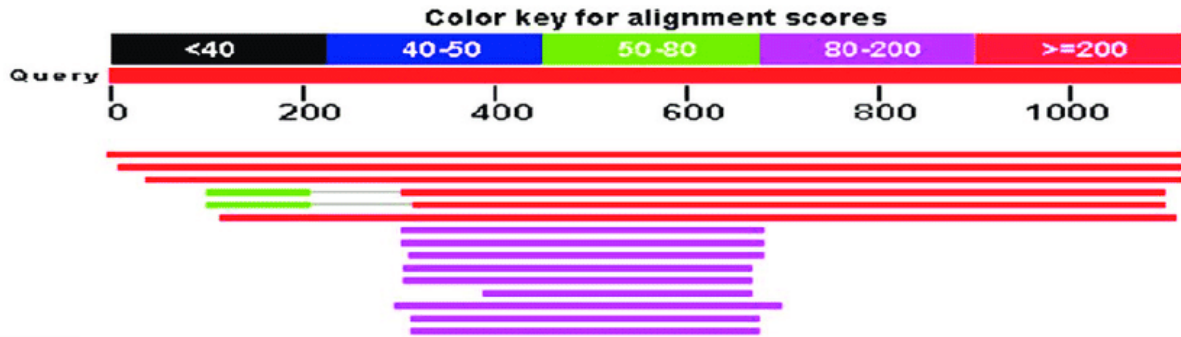
Transdecoder workflow

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.

Annotation tools. Protein BLAST

Distribution of 102 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



and for links to other resources: [U](#) UniGene [E](#) GEO [G](#) Gene [S](#) Structure [M](#) Map Viewer

Sequences producing significant alignments:
(Click headers to sort columns)


Accession	Description
NC_009396.1	Leishmania infantum JPCMS chromosome 12, complete sequence >emb AM502230.1 Leishmania infantum chromosome 12
NC_007253.1	Leishmania major strain Friedlin chromosome 12, complete sequence >emb CT005251.1 Leishmania major strain Friedlin,
NC_009304.1	Leishmania braziliensis MHOM/BR/75/M2904 chromosome 12 >emb AM494949.1 Leishmania braziliensis chromosome 12

- ☐ PREDICTED: lactase-phlorizin hydrolase [Macaca fascicularis]
- ☐ hypothetical protein EGK_05718 [Macaca mulatta]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Macaca mulatta]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Papio anubis]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Macaca nemestrina]
- ☐ hypothetical protein EGM_05165 [Macaca fascicularis]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Chlorocebus sabaeus]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Mandrillus leucophaeus]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Rhinopithecus roxellana]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Cercopithecus atys]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Callithrix jacchus]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Saimiri boliviensis boliviensis]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Aotus nancymae]
- ☐ PREDICTED: lactase-phlorizin hydrolase [Colobus angolensis palliatus]
- ☐ PREDICTED: LOW QUALITY PROTEIN: lactase-phlorizin hydrolase [Pan troglodytes]

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


	Max score	Total score	Query cover	E value	Ident	Accession
	4011	4011	100%	0.0	99%	EAX11622.1
	4011	4011	100%	0.0	100%	NP_002290.2
	4009	4009	100%	0.0	99%	AAA59504.1
	4009	4009	100%	0.0	99%	CAA30801.1
	3969	3969	100%	0.0	99%	XP_003822858.1
	3930	3930	100%	0.0	98%	XP_003267652.1
	3891	3891	100%	0.0	96%	XP_004032645.1
	3886	3886	100%	0.0	97%	XP_002812489.1
	3835	3835	100%	0.0	96%	XP_005573098.1
	3834	3834	100%	0.0	96%	EHH22449.1
	3833	3833	100%	0.0	96%	XP_014965495.1
	3833	3833	100%	0.0	96%	XP_003909221.1
	3832	3832	100%	0.0	96%	XP_011758105.1
	3829	3829	100%	0.0	96%	EHH55875.1
	3828	3828	100%	0.0	96%	XP_007963046.1
	3825	3825	100%	0.0	96%	XP_011825664.1
	3823	3823	100%	0.0	95%	XP_010385578.1
	3821	3821	100%	0.0	96%	XP_011925242.1
	3741	3741	100%	0.0	93%	XP_002749525.1
	3723	3723	100%	0.0	93%	XP_003922057.1
	3682	3682	100%	0.0	92%	XP_012332156.1
	3547	3547	100%	0.0	90%	XP_011793136.1
	3491	3694	95%	0.0	98%	XP_009441718.1

Databases for ORF annotation. *Uniprot*



UniProtKB ▾

BLAST Align Retrieve/ID mapping Peptide search SPARQL


**We will be switching to the new UniProt website in a few weeks. Please explore and share your feedback.**

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.


UniProtKB

UniProt Knowledgebase


Swiss-Prot (566,996)

 Manually annotated and reviewed. Records with information extracted from literature and curator-evaluated computational analysis.

TrEMBL (230,328,648)


 Automatically annotated and not reviewed. Records that await full manual annotation.

UniRef




The UniProt Reference Clusters (UniRef) provide clustered sets of sequences from the UniProt Knowledgebase (including isoforms) and selected UniParc records.

UniParc



UniParc is a comprehensive and non-redundant database that contains most of the publicly available protein sequences in the world.


Proteomes




A proteome is the set of proteins thought to be expressed by an organism. UniProt provides proteomes for species with completely sequenced genomes.

Supporting data

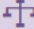
Literature citations




Cross-ref. databases




Taxonomy




Diseases




Subcellular locations




Keywords




Getting started


 **Text search**

Our basic text search allows you to search all the resources available


 **BLAST**

Find regions of similarity between your sequences

 **UniProt data**

 **Download latest release**

Get the UniProt data

 **Statistics**

View Swiss-Prot and TrEMBL statistics

<https://www.uniprot.org/>

12

Databases for ORF annotation. *Uniprot*

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

UniProtKB chondrichthyes

Advanced Search

BLAST Align Retrieve/ID mapping Peptide search SPARQL Help Contact

UniProt BETA We will be switching to the new UniProt website in a few weeks. Please explore and share your feedback. Take me to the new website.

UniProtKB 2022_01 results

UniProtKB consists of two sections:

- Reviewed (Swiss-Prot) - Manually annotated**
Records with information extracted from literature and curator-evaluated computational analysis.
- Unreviewed (TrEMBL) - Computationally analyzed**
Records that await full manual annotation.

The UniProt Knowledgebase (UniProtKB) is the central hub for the collection of functional information on proteins, with accurate, consistent and rich annotation. In addition to capturing the core data mandatory for each UniProtKB entry (mainly, the amino acid sequence, protein name or description, taxonomic data and citation information), as much annotation information as possible is added.

Help UniProtKB help video Other tutorials and videos Downloads

Filter by

- Reviewed (281) Swiss-Prot
- Unreviewed (142,194) TrEMBL

Popular organisms

- TETCF (46)
- TORMA (88)
- SQUAC (231)
- PORAF (7)
- CHIPU (33,574)
- Other organisms

Search terms

Filter "chondrichthyes" as:

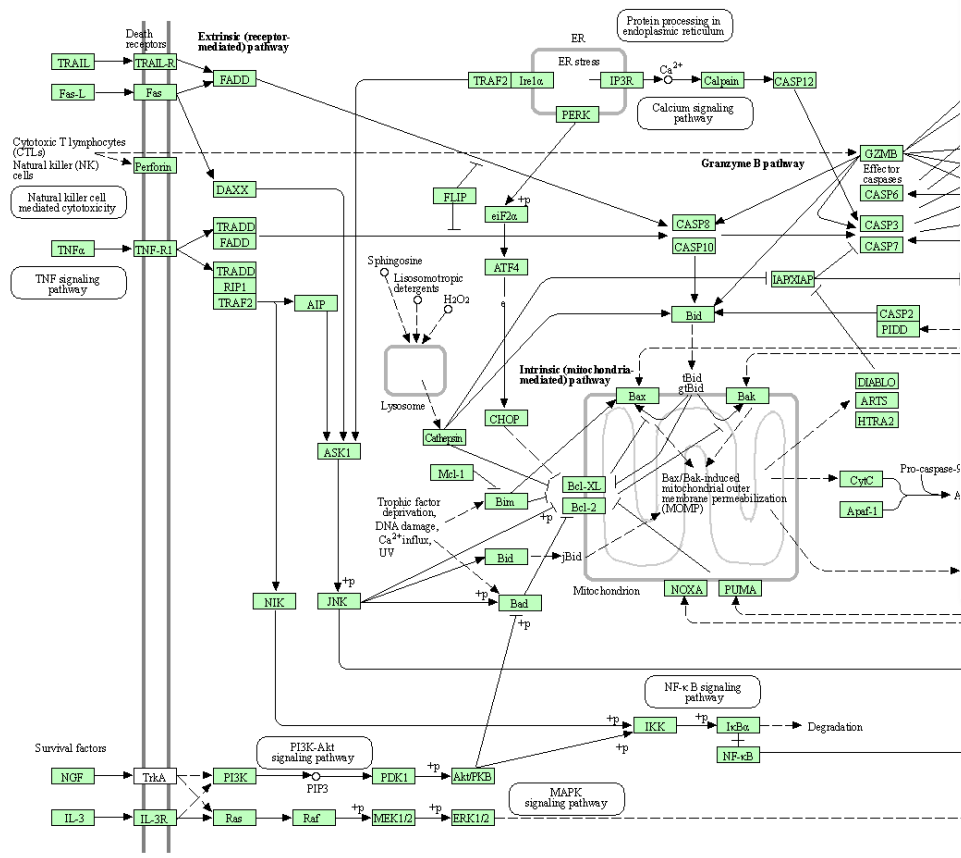
1 to 25 of 142,475 Show 25

Entry	Accession	Gene names	Organism	Length
<input type="checkbox"/> P02712	CHRN1	Acetylcholine receptor subunit beta	Tetronarce californica (Pacific electric ray) (Torpedo californica)	493
<input type="checkbox"/> P26362	CFTR	Cystic fibrosis transmembrane conductance regulator	Squalus acanthias (Spiny dogfish)	1,492
<input type="checkbox"/> Q91437	CAD	Calcitonin receptor-like receptor	Squalus acanthias (Spiny dogfish)	2,242
<input type="checkbox"/> P04058	ACHE	Acetylcholinesterase	Tetronarce californica (Pacific electric ray) (Torpedo californica)	586
<input type="checkbox"/> P07692	ACHE	Acetylcholinesterase	Torpedo marmorata (Marbled electric ray)	590
<input type="checkbox"/> P55013	SLC12A2	Solute carrier family 12 member 2	Squalus acanthias (Spiny dogfish)	1,191
<input type="checkbox"/> P02718	CHRN	Acetylcholine receptor subunit delta	Tetronarce californica (Pacific electric ray) (Torpedo californica)	522
<input type="checkbox"/> P02714	CHRG	Acetylcholine receptor subunit gamma	Tetronarce californica (Pacific electric ray) (Torpedo californica)	506
<input type="checkbox"/> P84232	H3.2	Histone H3.2	Poroderma africanum (Striped catshark) (Squalus africanus)	136
<input type="checkbox"/> Q73925	KCNQ1	Potassium voltage-gated channel subunit	Squalus acanthias (Spiny dogfish)	660

Post-annotation tools. Gene enrichment

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

APOPTOSIS



04210 6/20/18
(c) Kanehisa Laboratories



KEGG [Help](#)
[Japanese](#)

KEGG Home
[Release notes](#)
[Current statistics](#)

KEGG Database
[KEGG overview](#)
[Searching KEGG](#)
[KEGG mapping](#)
[Color codes](#)

KEGG Objects
[Pathway maps](#)
[Brite hierarchies](#)
[KEGG DB links](#)

KEGG Software
[KEGG API](#)
[KGML](#)

KEGG FTP
[Subscription](#)
[Background info](#)

[GenomeNet](#)
[DBGET/LinkDB](#)

[Feedback](#)
[Copyright request](#)
[Kanehisa Labs](#)

KEGG: Kyoto Encyclopedia of Genes and Genomes

KEGG is a database resource for understanding high-level functions and utilities of the biological system, such as the cell, the organism and the ecosystem, from molecular-level information, especially large-scale molecular datasets generated by genome sequencing and other high-throughput experimental technologies. See [Release notes](#) (April 1, 2022) for new and updated features.

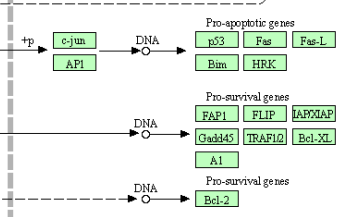
[News](#) [Background info](#) [updated](#)

Main entry point to the KEGG web service
KEGG2 [KEGG Table of Contents](#) [\[Update notes\]](#) [\[Release history\]](#)

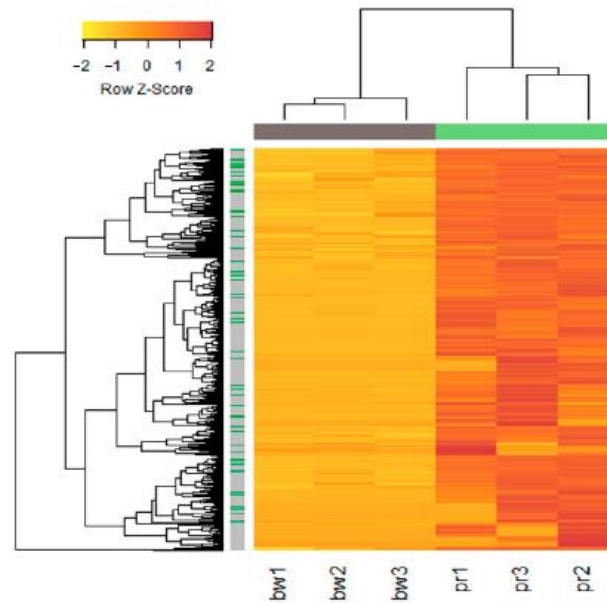
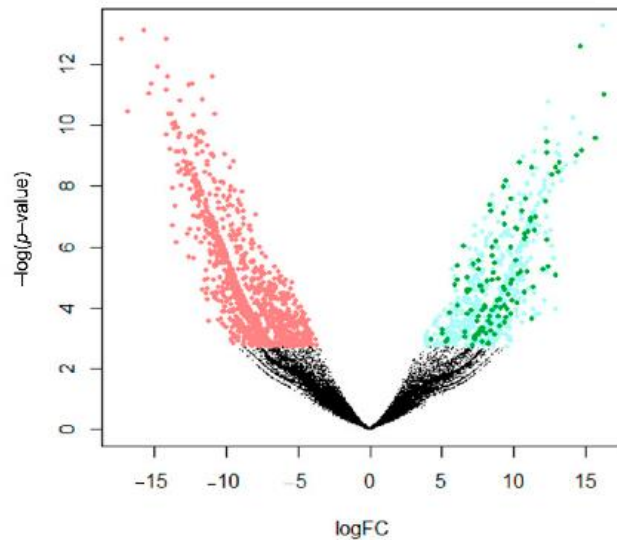
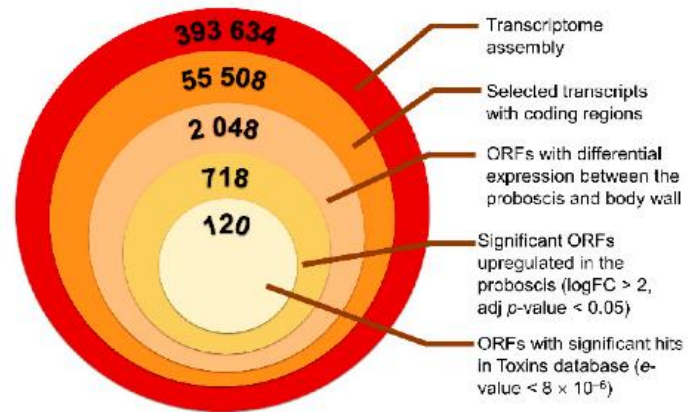
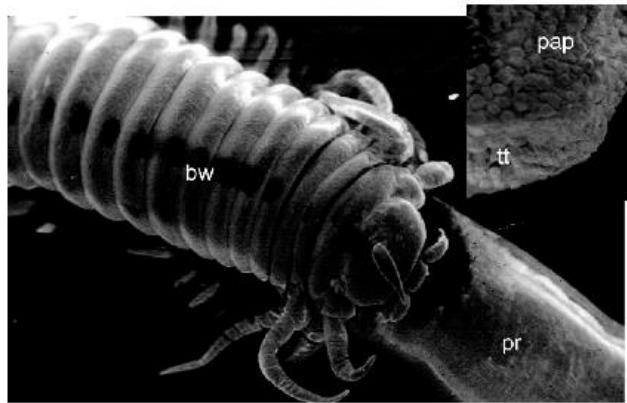
- Data-oriented entry points**
- | | | |
|-----------------------|--|--|
| KEGG PATHWAY | KEGG pathway maps | Pathway
Brite
Brite table
Module
Network
KO (Function)
Organism
Virus
Compound
Disease (ICD)
Drug (ATC)
Drug (Target)
Antimicrobials |
| KEGG BRITE | BRITE hierarchies and tables | |
| KEGG MODULE | KEGG modules | |
| KEGG ORTHOLOGY | KO functional orthologs [Annotation] | |
| KEGG GENES | Genes and proteins [SeqData] | |
| KEGG GENOME | Genomes [KEGG Virus Taxonomy] | |
| KEGG COMPOUND | Small molecules | |
| KEGG GLYCAN | Glycans | |
| KEGG REACTION | Biochemical reactions [RModule] | |
| KEGG ENZYME | Enzyme nomenclature | |
| KEGG NETWORK | Disease-related network variations | |
| KEGG DISEASE | Human diseases | |
| KEGG DRUG | Drugs [New drug approvals] | |
| KEGG MEDICUS | Health information resource [Drug labels search] | |

Organism-specific entry points
KEGG Organisms Enter org code(s) [hsa](#) [hsa eco](#)

- Analysis tools**
- | | |
|--------------------|--|
| KEGG Mapper | KEGG PATHWAY/BRITE/MODULE mapping tools |
| BlastKOALA | BLAST-based KO annotation and KEGG mapping |
| GhostKOALA | GHOSTX-based KO annotation and KEGG mapping |
| KofamKOALA | HMM profile-based KO annotation and KEGG mapping |
| BLAST/FASTA | Sequence similarity search |
| SIMCOMP | Chemical structure similarity search |



De novo assembly and annotation. Example



Training set

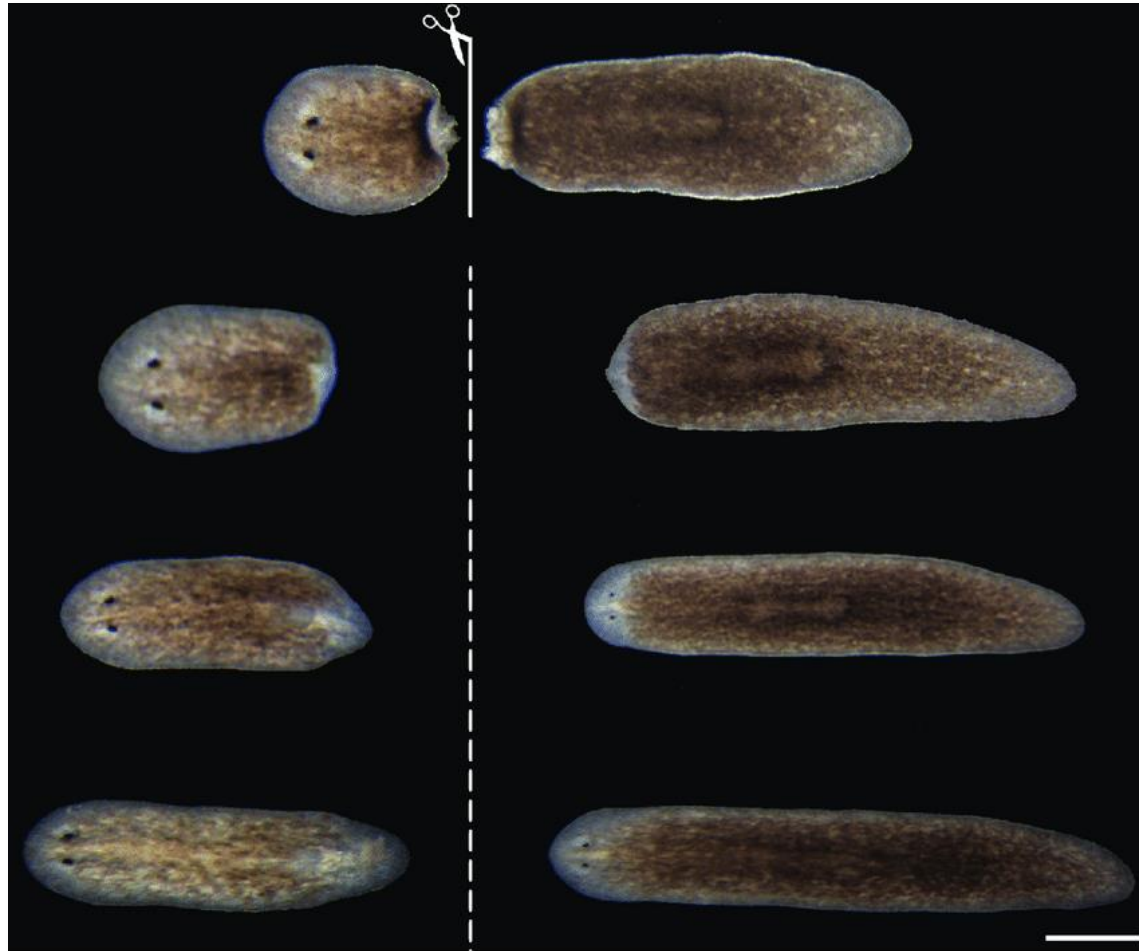
Raw Fastq data set

?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

Transcriptome assembly and annotation. *Objectives*

- 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)