# Short lecture
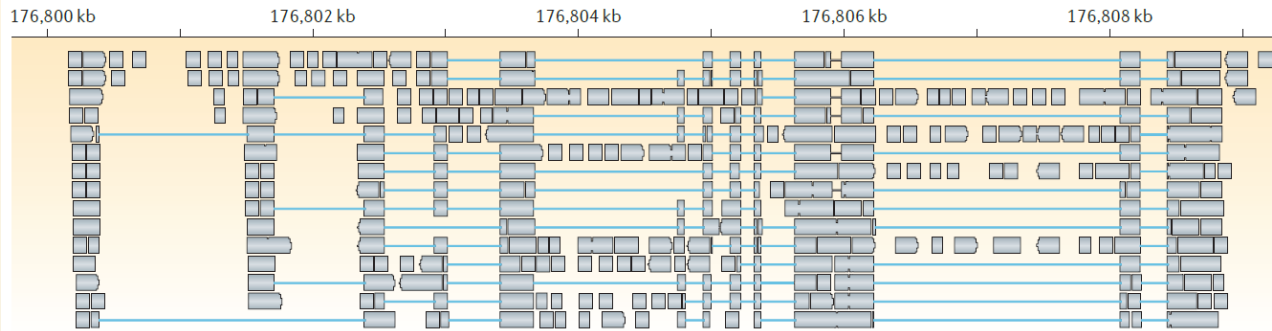## – How to make a transcriptome

# What to consider I

- What do I want?

- What will I use it for?

- Which resources are available for your species (very closely related species)?

- What kind of data do I have?

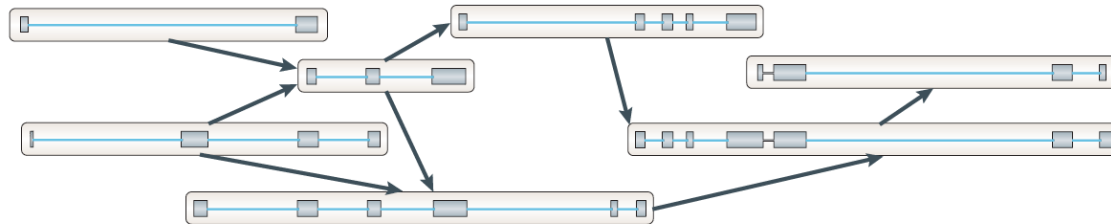- 2n or xn ploidy?

# The strategies I

- Reference based (*ab initio*)
  - Maps RNAseq reads back towards reference genome and builds transcripts
  - Needs a certain amount of splice-junction covering reads
- *De novo* (with/without genome guiding)
  - Assembly of RNAseq reads only
  - Guided: reads are clustered according to chromosome / scaffold prior to assembly
- Mixed approach
  - Merging several assemblies to one
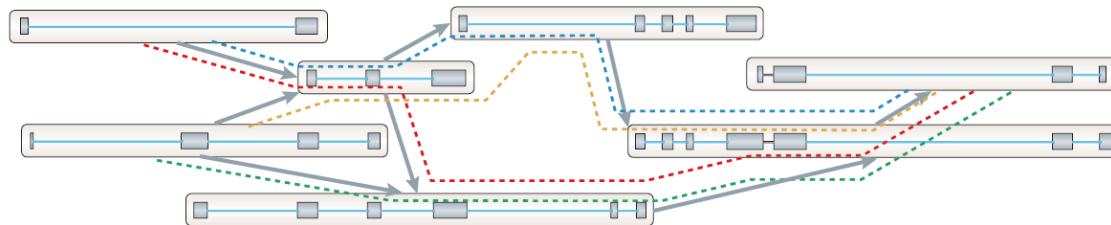
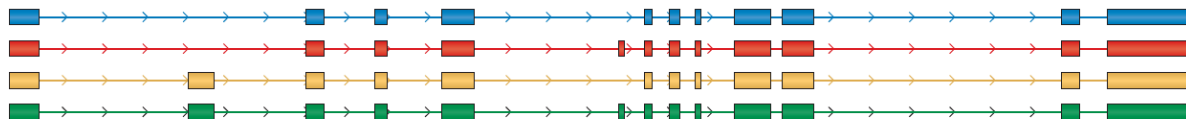# Reference based



**a** Splice-align reads to the genome

**b** Build a graph representing alternative splicing events

**c** Traverse the graph to assemble variants

**d** Assembled isoforms
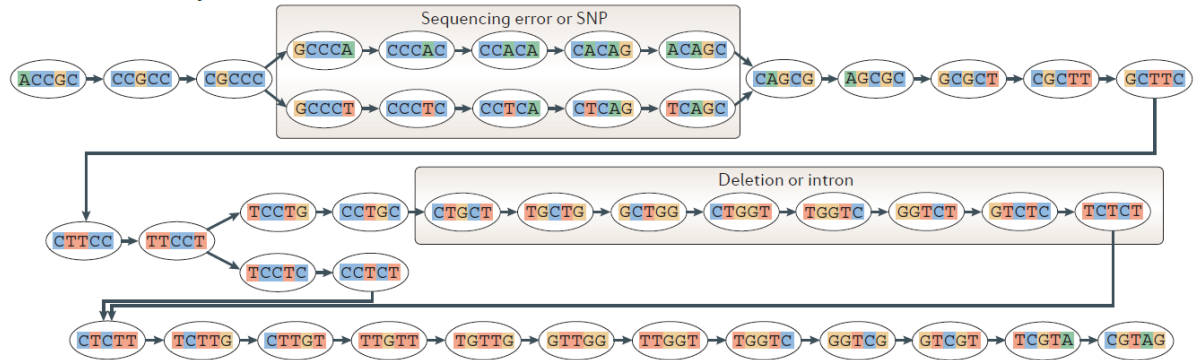
# Reference based II

- Benefits
  - Time efficient / single computer job
  - Requires less coverage of samples
  - Artefacts / contaminations does not align to the reference
  - Low abundance / novel isoforms are resolved
- Complications
  - Depends on quality of reference
  - Gene dense organisms
  - Higher eukaryotes with complex splice variants – especially *trans*-splicing
  - Software settings may discard splice variants / transcripts
  - Different treatment of multi-mapping reads
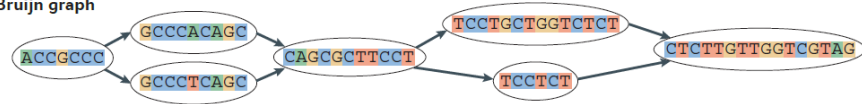
# *De novo*



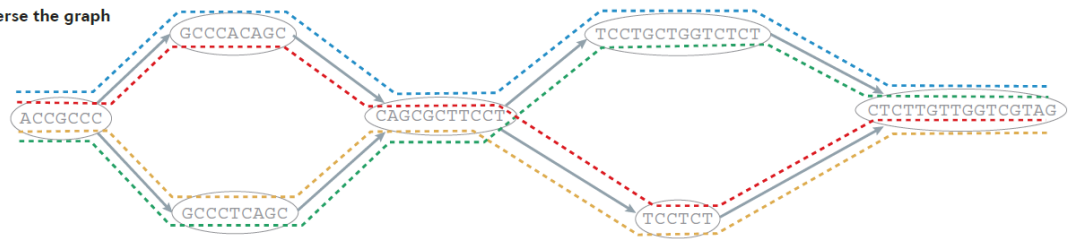**a** Generate all substrings of length k from the reads

**b** Generate the De Bruijn graph

Sequencing error or SNP

Deletion or intron

**c** Collapse the De Bruijn graph

**d** Traverse the graph

**e** Assembled isoforms

# *De novo* II

- Benefits
  - No reference needed
  - Detects all transcripts (coverage dependent)
  - No knowledge/prediction of splice sites needed
  - Complex splice patterns can be resolved
- Complications
  - Requires lots of computing power
  - Requires more coverage to resolve transcripts
  - Sensitive to read errors and artefacts / contaminations
  - Paralog ("gene copies") resolution is an issue

# Mixed approach

- *De novo* and *ab initio* assembly concatenation
- Multiple kmer strategy

- Who benefits from a mixed approach?
  - Gene dense eukaryotes
  - Polyploid species
  - When the aim is to make a really good reference transcriptome

# How to make it I

- Use all available data

- Consider normalization to shorten computation time and increase chances of resolving less abundant transcripts



Mouse RNA-seq Trinity assembly

# How to make it II

- Consider the assembly algorithm
  - Large datasets with short reads benefits from using De bruijn graph based assembly programs (more than a hundred million read pairs)
  - Small datasets with short reads benefits from using Overlap-Layout-Consensus (OLC) based assembly programs



Overlaps identified

Reads connected by overlaps

# How to make III

- Do you have a an organism known to be gene dense with overlapping UTRs?
    - Select a program with options like jaccard clip to improve algorithm
    - The cost is more computation time so do not use it unless necessary

# What to consider – INFBIO case

- What do I want? – transcriptome
- What will I use it for? – differential expression
- Which resources are available for your species (very closely related species)? – genome
- What kind of data do I have? – Illumina PE
- 2n or ploidy? – 2n

# Choosing our strategy

# Choosing our strategy



- Non-model species
- Genome resource
- Genes of interest are <u>poorly annotated</u>
- Genes of interest are <u>in fragmented area</u> of genome

# Trinity assembler

- Trinity is the best single parameter *de novo* RNA assembly pipeline available

- Good on splice variants, full length transcripts and resolution of lowly expressed transcripts

- Contains tools to help with visualizations

# Trinity pipeline - Inchworm

- It employs a greedy kmer based approach to reconstruct the best representative for a transcriptionally active region (often full-length dominant isoform).



sequence     **ATGGAAGTCGCGGAATC**

7mers

```
ATGGAAG
 TGGAAGT
  GGAAGTC
   GAAGTCG
    AAGTCGC
     AGTCGCG
      GTCGCGG
       TCGCGGA
        CGCGGAA
         GCGGAAT
          CGGAATC
```



Piece RNA-Seq reads into contigs (Inchworm)

a

Read set

Extend in *k*-mer space and break ties

>a121:len=5845
>a122:len=2560
>a123:len=4443
>a124:len=48
>a126:len=66

Linear sequences

# Trinity pipeline

- Chrysalis clusters Inchworm related contigs into components (alternatively spliced variants)

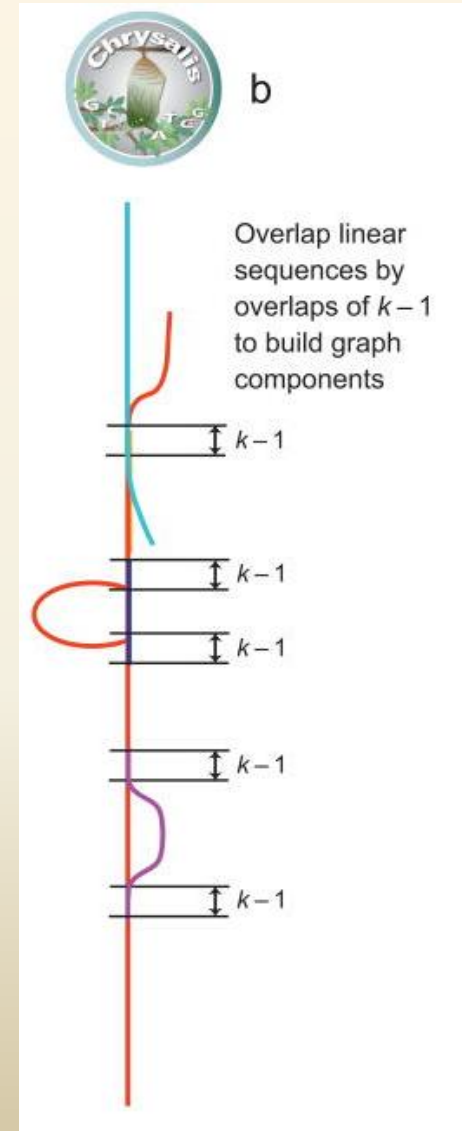- Then a De Bruijn graph is made for each component



Cluster contigs into components (Chrysalis)

Assign reads to components (Chrysalis)

Split overlapping transcripts based on coverage and read pairings



b

Overlap linear sequences by overlaps of $k - 1$ to build graph components

$k - 1$

$k - 1$

$k - 1$

$k - 1$

$k - 1$

# Trinity pipeline

- Butterfly analyzes the paths taken by reads and read pairings in the graphs and reports all plausible transcripts including splice variants and transcripts derived from paralogs (duplicated genes)



De Bruijn graph (k = 5)

Compacting

Compact graph

Finding paths

Compact graph with reads

Extracting sequences

Transcripts



Enumerate transcript isoforms using reads (Butterfly)

# Trinity computation requirements

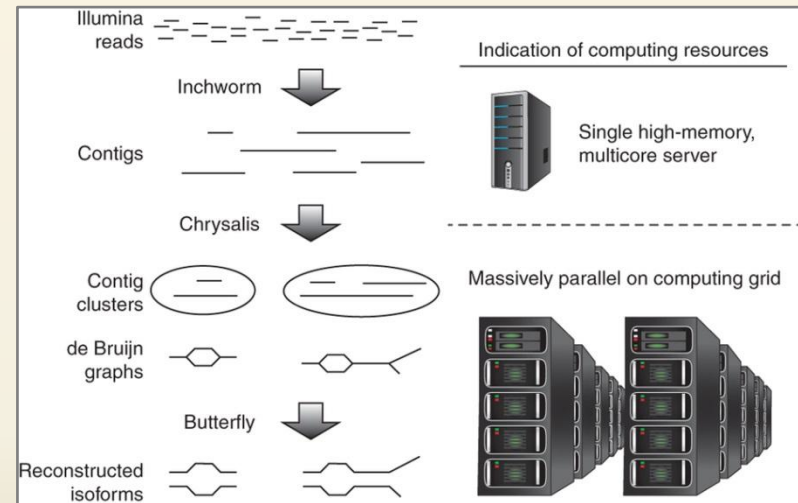- Assembly algorithms require large amounts of memory
- 2/3rds of Trinity is parallelized to save computation time
- Estimate at least 1 week of trial/error/final computation
- Remember to calculate memory/time requirements before starting!
  - 1Gb RAM / million reads
  - ½ - 1 hour / million reads

# What to expect

- Significantly more transcripts than predicted in the same or closely related species!

- Low coverage over splice junctions, sequencing errors and heterozygosity restricts full-length transcript reconstruction

```
###################################
## Counts of transcripts, etc.
###################################
Total trinity 'genes':   320520
Total trinity transcripts:        468626
Percent GC: 47.31


############################################
Stats based on ALL transcript contigs:
############################################

        Contig N10: 3657
        Contig N20: 2645
        Contig N30: 2042
        Contig N40: 1597
        Contig N50: 1235

        Median contig length: 459
        Average contig: 784.28
        Total assembled bases: 367534825


##################################################
## Stats based on ONLY LONGEST ISOFORM per 'GENE':
##################################################

        Contig N10: 3360
        Contig N20: 2278
        Contig N30: 1635
        Contig N40: 1193
        Contig N50: 880

        Median contig length: 382
        Average contig: 634.85
        Total assembled bases: 203483069
```

# How to make it comparable

- Trinity comes with:
  - Full length estimation (BLAST based)
  - Abundance estimation (simple expression analysis)
- Consider mapping transcripts towards a reference genome if available