*INF-BIOx121 2017*

# RNA-seq
## differential expression analysis

Arvind Sundaram
Sep 18-20, 2017

*RNA-seq analysis*

# Transcriptome

Arvind Sundaram
Sep 18, 2017

# Transcriptome
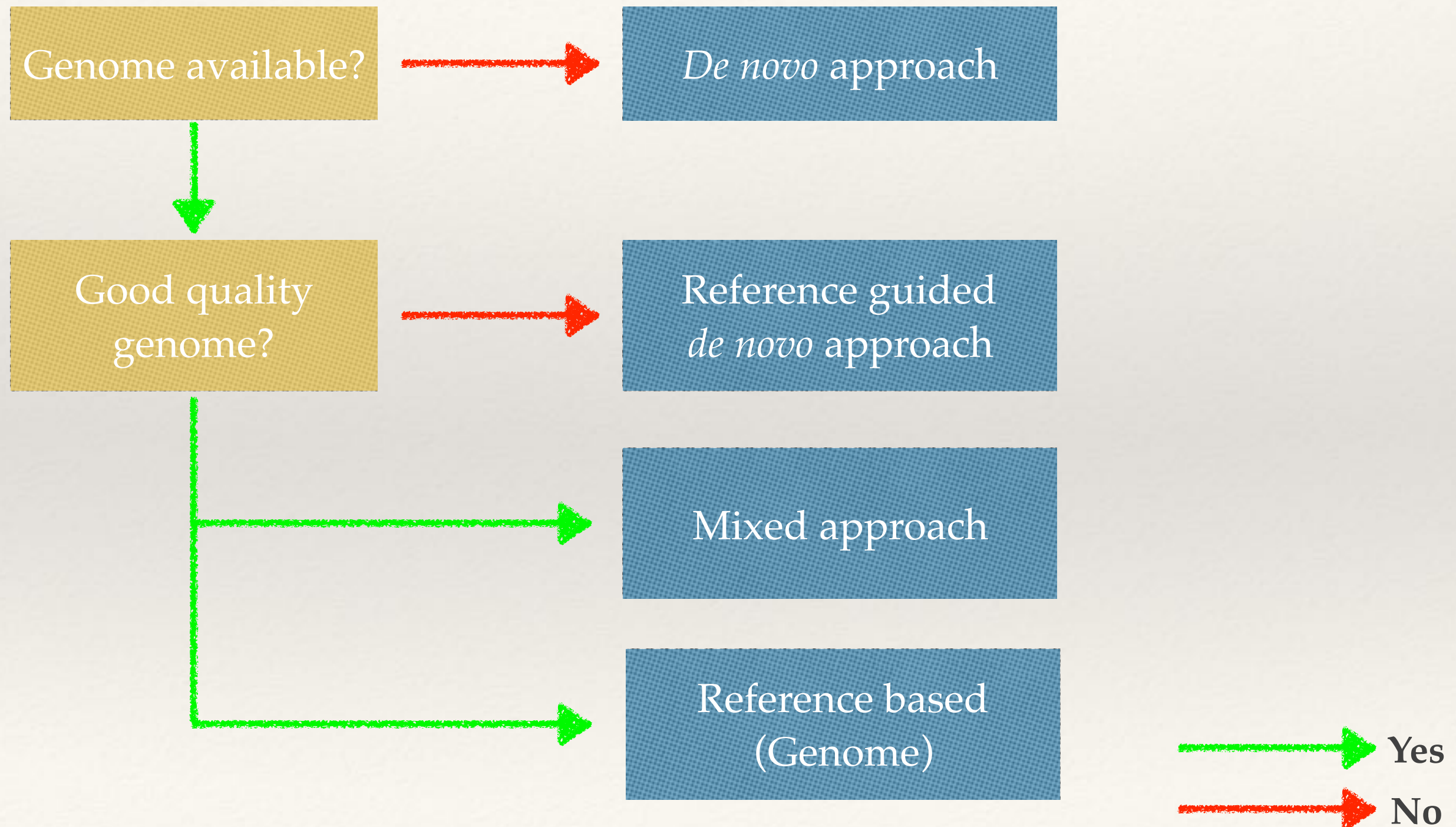
❖ Transcriptome is the total set of transcripts in an organism

❖ Transcriptome changes across cell types and environmental conditions

So….

❖ Transcriptome is a set of (all) RNA molecules in one cell or a population of cells in a given moment

❖ 'Constructing' a global transcriptome aims to capture all possible transcripts found across all cell and tissue types.
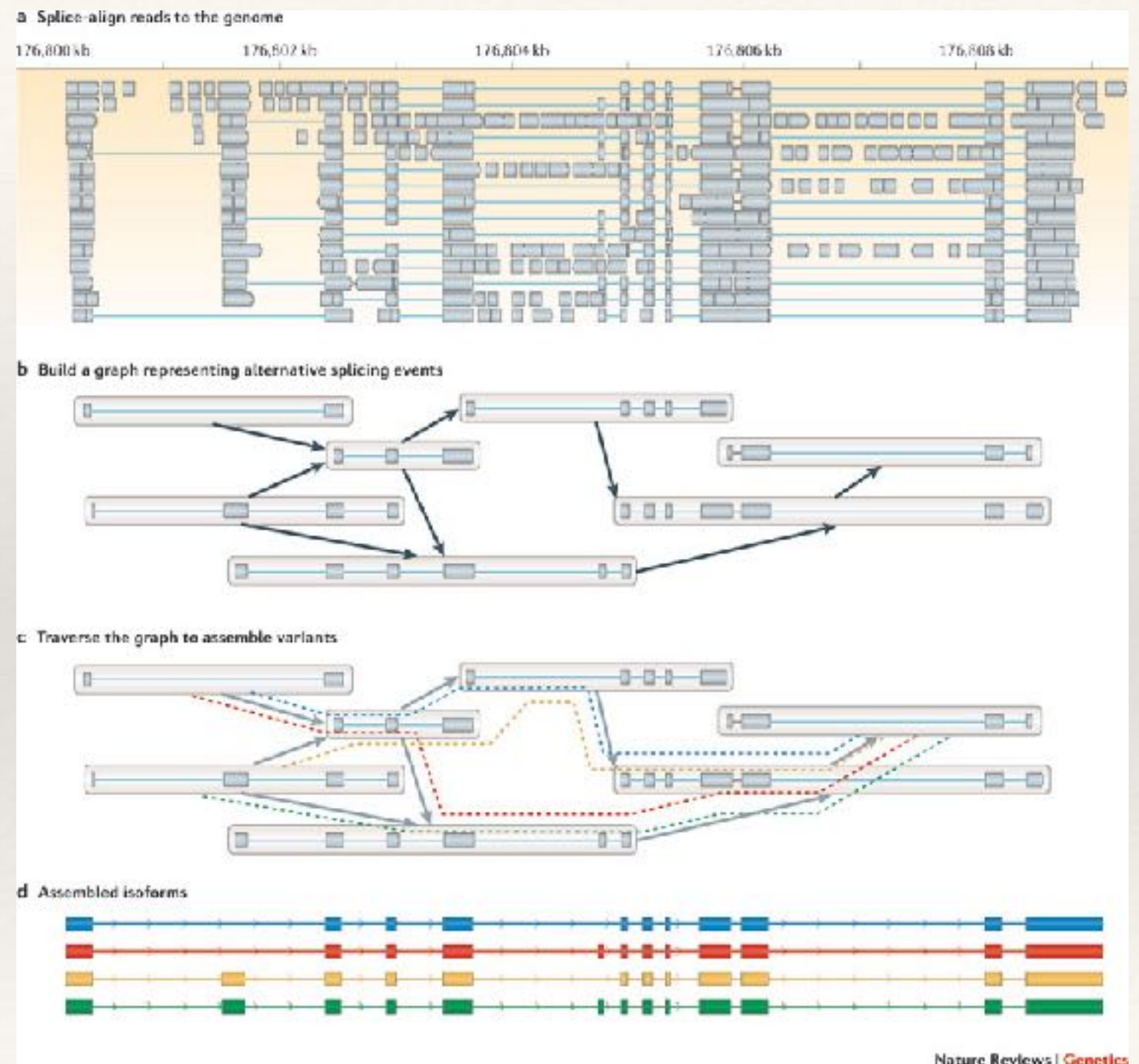
# Assembling a transcriptome

# Reference based transcriptome

- Requires good quality (draft) genome

- Splice-aware aligner

- Improves on existing knowledge

Fast

less CPU

prior knowledge



a  Splice-align reads to the genome

b  Build a graph representing alternative splicing events

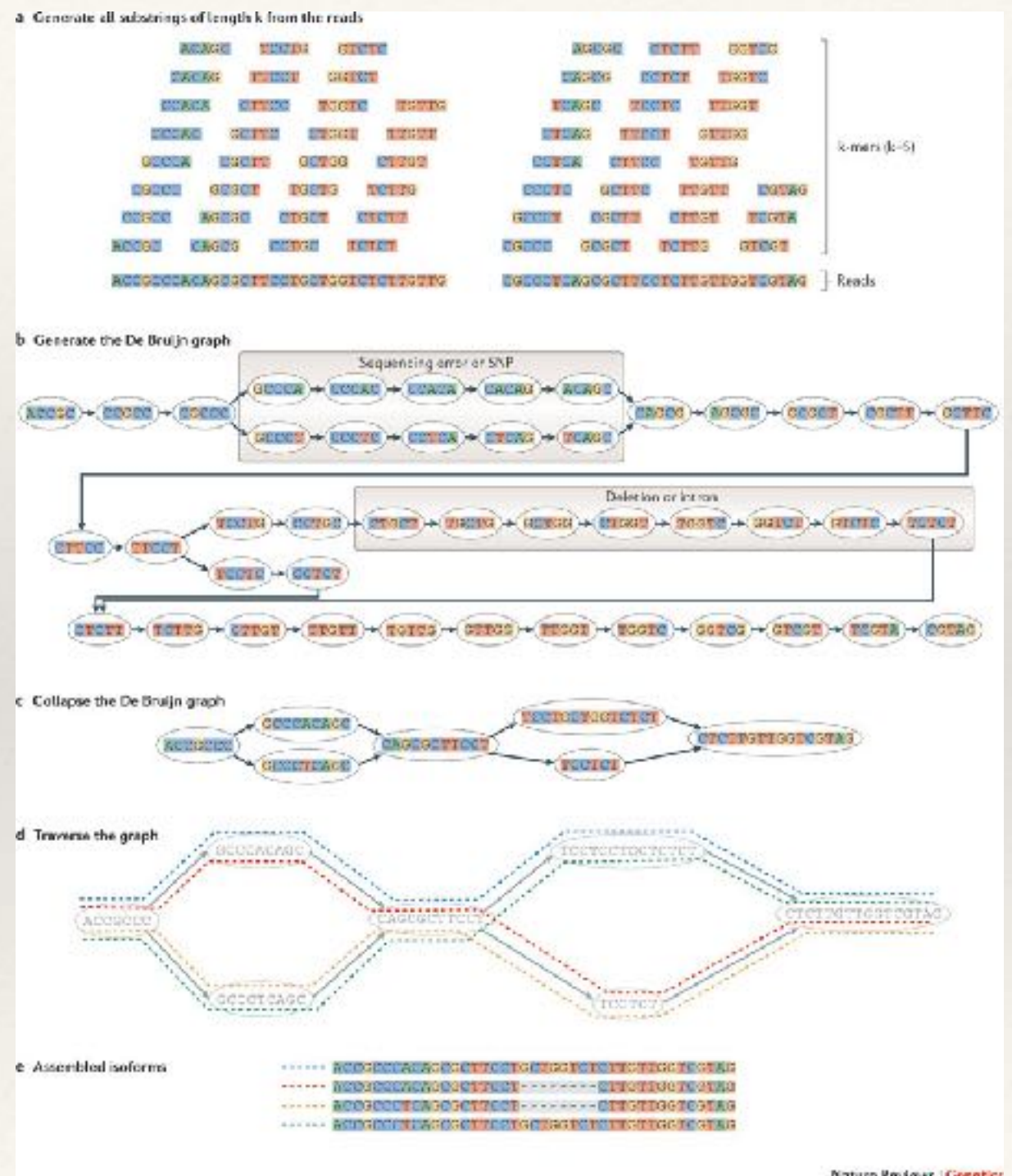c  Traverse the graph to assemble variants

d  Assembled isoforms

# *de novo* transcriptome

- Requires high coverage sequence data

- Transcript detection based on coverage

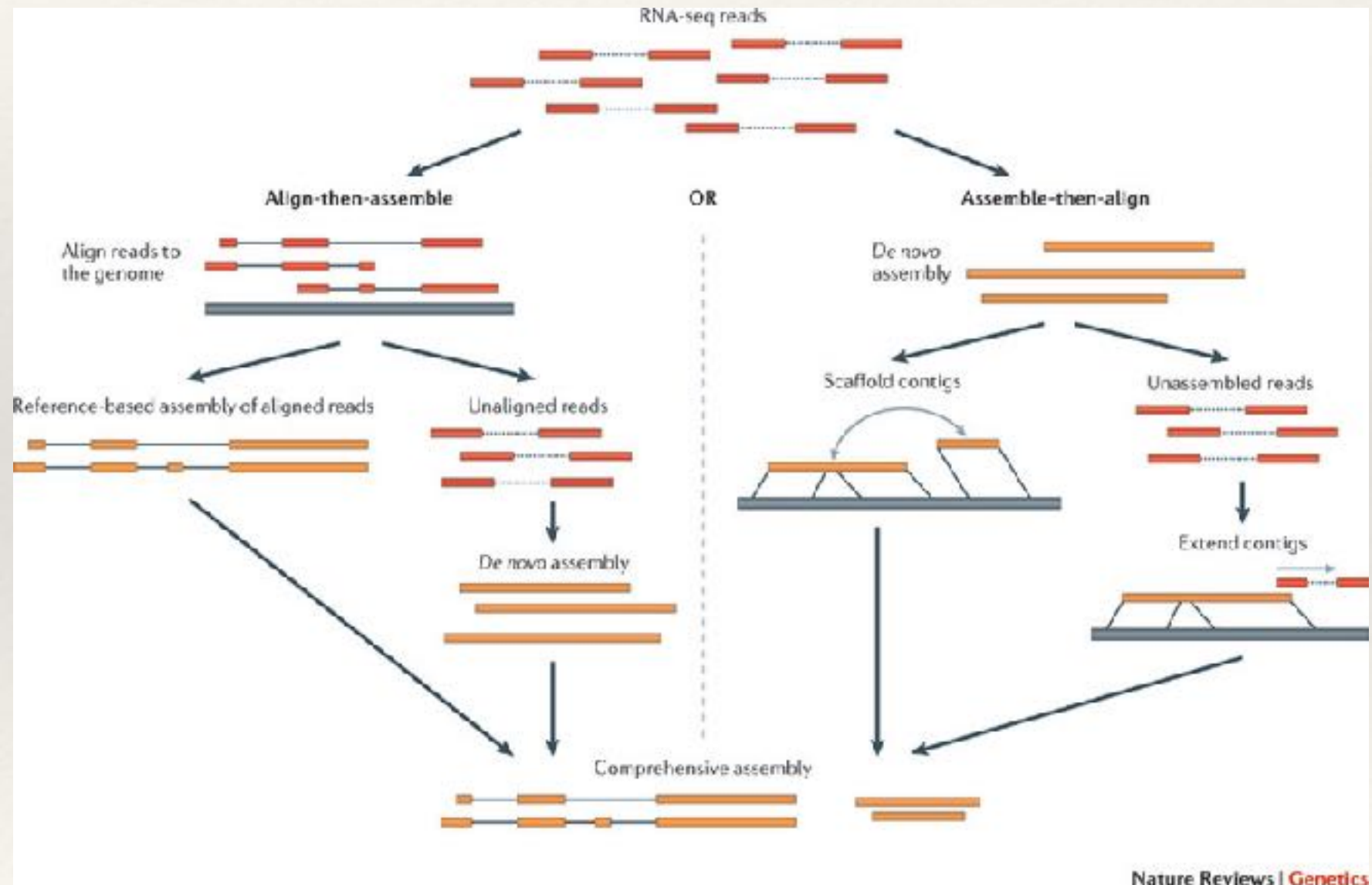- Multiple gene-copies is difficult to resolve

Slow

more CPU

No prior knowledge

# Mixed approach

- ❖ Identifies novel transcripts

- ❖ Polyploidy species



RNA-seq reads

Align-then-assemble   OR   Assemble-then-align

Align reads to the genome

De novo assembly

Reference-based assembly of aligned reads   Unaligned reads

Scaffold contigs   Unassembled reads

De novo assembly

Extend contigs

Comprehensive assembly

Nature Reviews | Genetics

Varies

mixed CPU

prior knowledge

# *de novo* transcriptome

- One will choose *de novo* / mixed approach due to many reasons
  - Non-model species with less genomic resources
  - Improve gene annotation
  - Genes of interest are not well annotated

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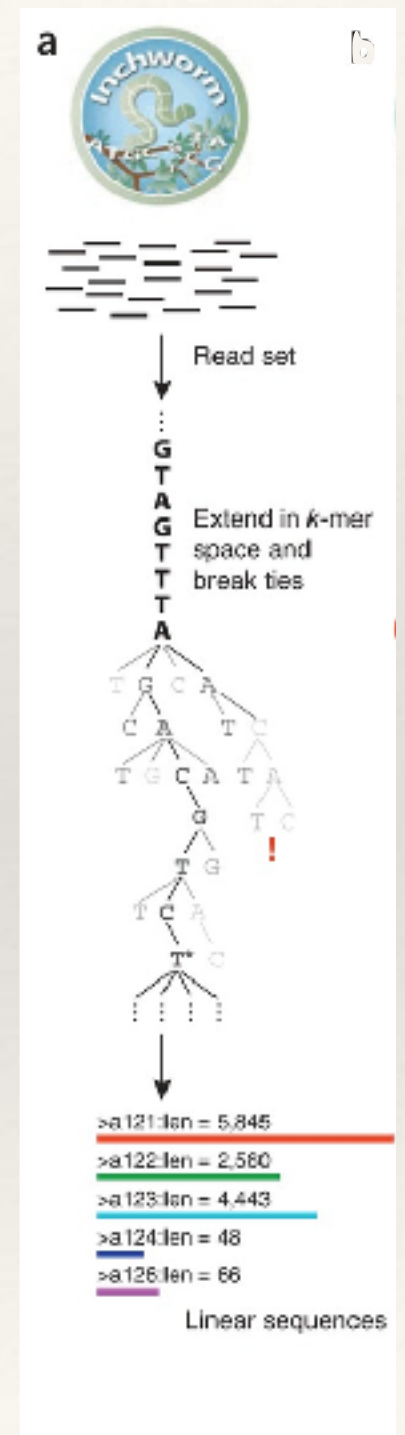
- ❖ Trinity is the best single parameter *de novo* RNA assembly pipeline available

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

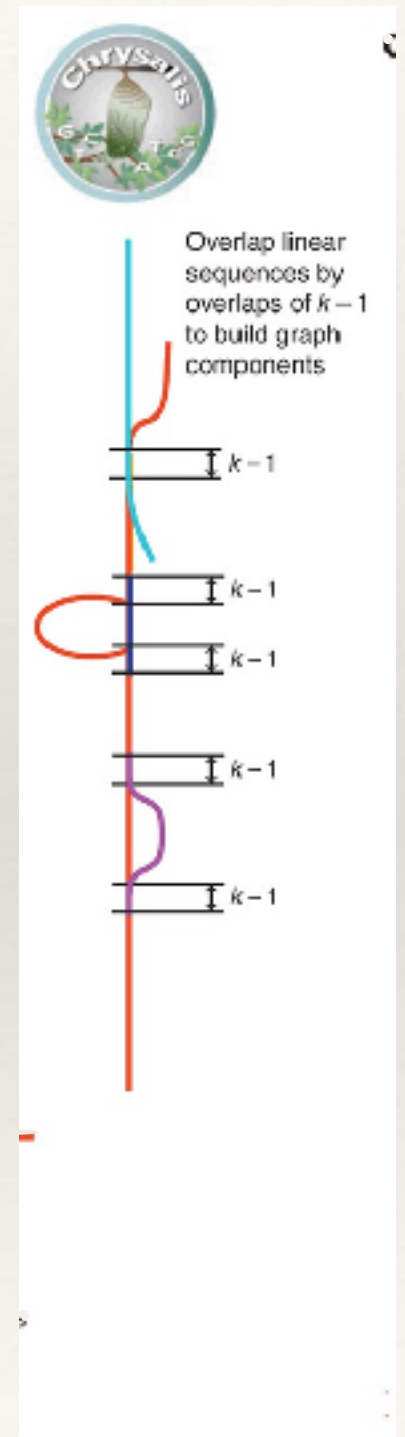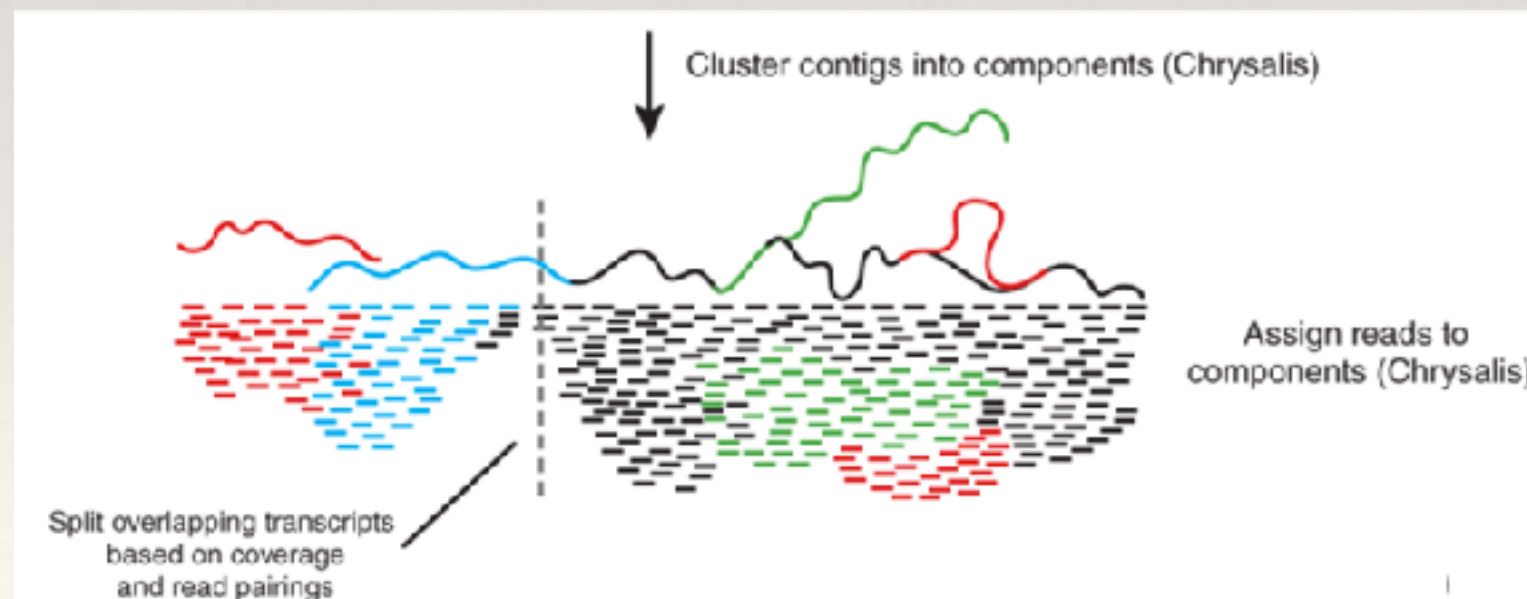- ❖ Contains tools to help with visualisations

# Trinity pipeline -Inchworm



❖ Inchworm assembles the data by greedily searching for paths in a k-mer graph, resulting in a collection of linear contigs, with each k-mer present only once in the contigs

❖ Inchworm does not capture the full complexity of the transcriptome; for example, only one alternatively spliced variant can be reported at full length per locus, with partial sequences reported for unique regions of any alternatively spliced transcripts.
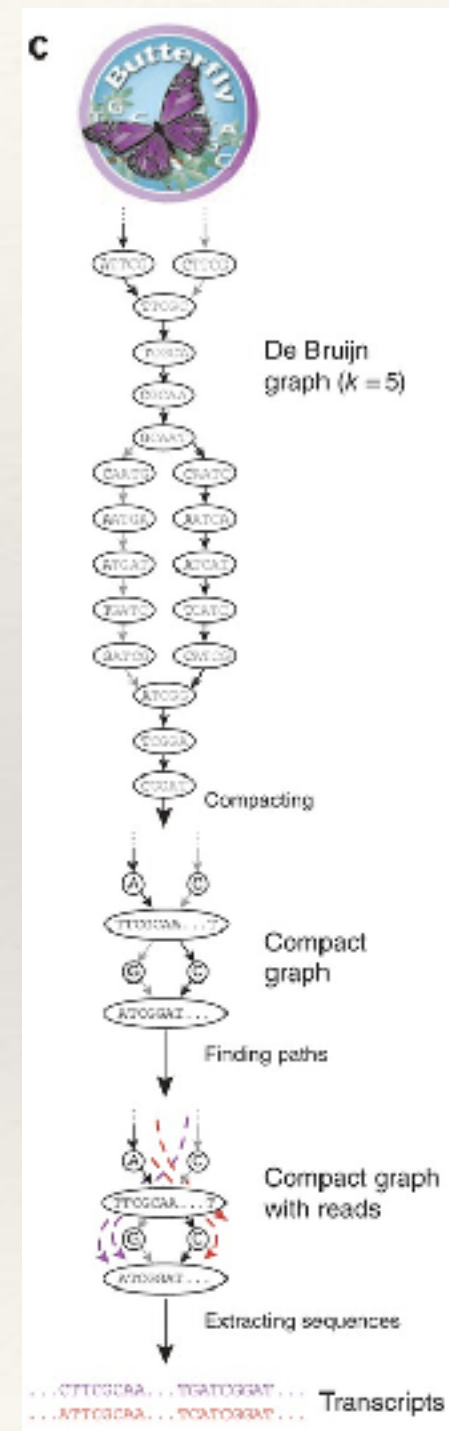
# Trinity pipeline –Chrysalis

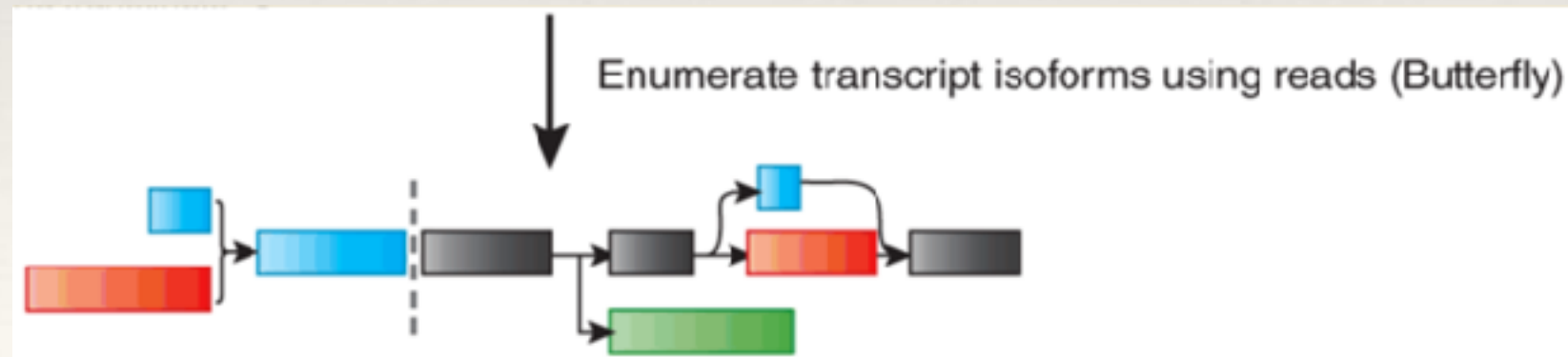❖ Chrysalis clusters minimally-overlapping Inchworm contigs into sets of connected components, and constructs complete de Bruijn graphs for each component

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

$k-1$

$k-1$

$k-1$

$k-1$

$k-1$

Cluster contigs into components (Chrysalis)

Assign reads to components (Chrysalis)

Split overlapping transcripts based on coverage and read pairings

# Trinity pipeline – Butterfly



❖ Butterfly reconstructs plausible full-length, linear transcripts by reconciling the individual de Bruijn graphs generated by Chrysalis with the original reads and paired-ends.

# Trinity assembler

❖ 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

- Trinity Wiki Home
- Installing Trinity
  - Trinity Computing Requirements
  - Accessing Trinity on Publicly Available Compute Resources
  - Run Trinity using Docker
- Running Trinity
  - Genome Guided Trinity Transcriptome Assembly
  - Gene Structure Annotation of Genomes
- Trinity process and resource monitoring
  - Monitoring Progress During a Trinity Run
  - Examining Resource Usage at the End of a Trinity Run
- Output of Trinity Assembly

- Assembly Quality Assessment
  - Counting Full-length Transcripts
  - RNA-Seq Read Representation
  - Contig Nx and ExN50 stats
  - Examine strand-specificity of reads
- Downstream Analyses
  - Transcript Quantification
  - QC Samples and Bio Replicates
  - Differential Expression
  - Coding Region Identification
  - Functional Annotation of Transcripts
  - Gene Ontology term functional category enrichments
- Trinity Tidbits
- Frequently Asked Questions (FAQ)

# Trinity assembler

❖ Assembly algorithms require large amounts of memory

❖ 2/3rds of Trinity is parallelised 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

  ❖ – 30 mins - 1 hour / million reads

# Transcriptome Assembly Quality Assessment

❖ Guide from Trinity

  ❖ https://github.com/trinityrnaseq/trinityrnaseq/wiki/Transcriptome-Assembly-Quality-Assessment

❖ BUSCO

❖ CEGMA (discontinued)

# Genome annotation pipeline

❖ MAKER2

  ❖ Step-by-step annotation guide

  ❖ https://github.com/sujaikumar/assemblage/blob/master/README-annotation.md

❖ Gene prediction based on ESTs and/or transcriptomes

  ❖ AUGUSTUS, GeneMark-ES, SNAP and more

# Genome annotation pipeline

- ❖ BLAST against appropriate databases

  - ❖ NCBI blast+ toolkit

- ❖ Gene ontology

- ❖ KEGG pathways

- ❖ blast2go

# blast2go