# Canadian Bioinformatics Workshops

www.bioinformatics.ca

bioinformaticsdotca.github.io

Supported by

OICR
Ontario Institute
for Cancer Research

**bio**informatics.ca

# Learning Objectives of Module

- Understand the challenges involved in reconstructing transcripts from RNA-Seq data

- Become familiar with computational algorithms and data structures leveraged for transcript assembly

- Appreciate the importance of strand-specific RNA-Seq data for transcript reconstruction

- Differentiate between differential gene expression and differential transcript usage.

**bio**informatics.ca

# Assembly Required



fragmen-
tation

RT

mRNA

sequence library

Sequencer

RT

fragmen-
tation

short sequence reads

Reconstruct original
full-length transcripts

Adapted from G. Raetsch

**bio**informatics.ca

# Transcript Reconstruction from RNA-Seq Reads

RNA-Seq reads



## Advancing RNA-Seq analysis

Brian J Haas & Michael C Zody                                    Nature Biotech, 2010

**New methods for analyzing RNA-Seq data enable *de novo* reconstruction of the transcriptome.**

bioinformatics.ca

# Transcript Reconstruction from RNA-Seq Reads

# Transcript Reconstruction from RNA-Seq Reads

# Transcript Reconstruction from RNA-Seq Reads



**Non-model organisms: "I don't have a reference genome!"**

# Transcript Reconstruction from RNA-Seq Reads

RNA-Seq reads

Assemble transcripts
*de novo*

# Transcript Reconstruction from RNA-Seq Reads



RNA-Seq reads

Assemble transcripts *de novo*

Align transcripts to genome

**bio**informatics.ca

# Transcript Reconstruction from RNA-Seq Reads
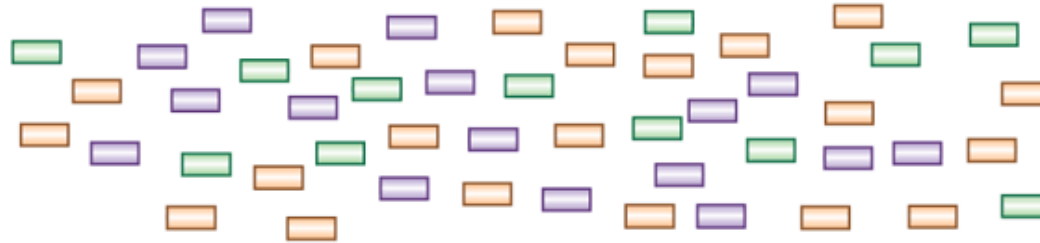
RNA-Seq reads

**Many tools to choose among:**

Align reads to genome

TopHat
STAR
HISAT
GSNAP
…

Genome

Assemble transcripts from spliced alignments

**Cufflinks**
**Stringtie**
**IsoLasso**
**Bayesembler**
**Trip**
**Traph**
**CEM**
**TransComb**
**…**

Assemble transcripts *de novo*

**Trinity**
**Oases**
**SoapDenovoTrans**
**AbyssTrans**
**IDBA-Tran**
**Shannon**
**BinPack**
**Bridger**
**…**

GMAP
BLAT
AAT
Spidey
Sim4
…

Align transcripts to genome

# Graph Data Structures Commonly Used For Assembly

RNA-Seq reads



- Sequence
- Order
- Orientation (+, -)
- Overlap

Reads to Graph

Nodes = sequence (+/-)
Edges = order, overlap

**bio**informatics.ca

# Graph Data Structures Commonly Used For Assembly



RNA-Seq reads

- Sequence
- Order
- Orientation (+, -)
- Overlap

Reads to Graph

Nodes = sequence (+/-)
Edges = order, overlap

**GATCGTCCGAGCGATTACA**

# Genome-Guided Transcript Reconstruction

## Splice-align reads to the genome

From Martin & Wang. Nature Reviews in Genetics. 2011

**bio**informatics.ca

# Genome-Guided Transcript Reconstruction

**Splice-align reads to the genome**



Alignment segment piles  =>   exon regions

From Martin & Wang. Nature Reviews in Genetics. 2011    **bio**informatics.ca

# Genome-Guided Transcript Reconstruction

**Splice-align reads to the genome**



Large alignment gaps   =>   introns

# Genome-Guided Transcript Reconstruction

**Splice-align reads to the genome**



Overlapping but different introns =  evidence of alternative splicing
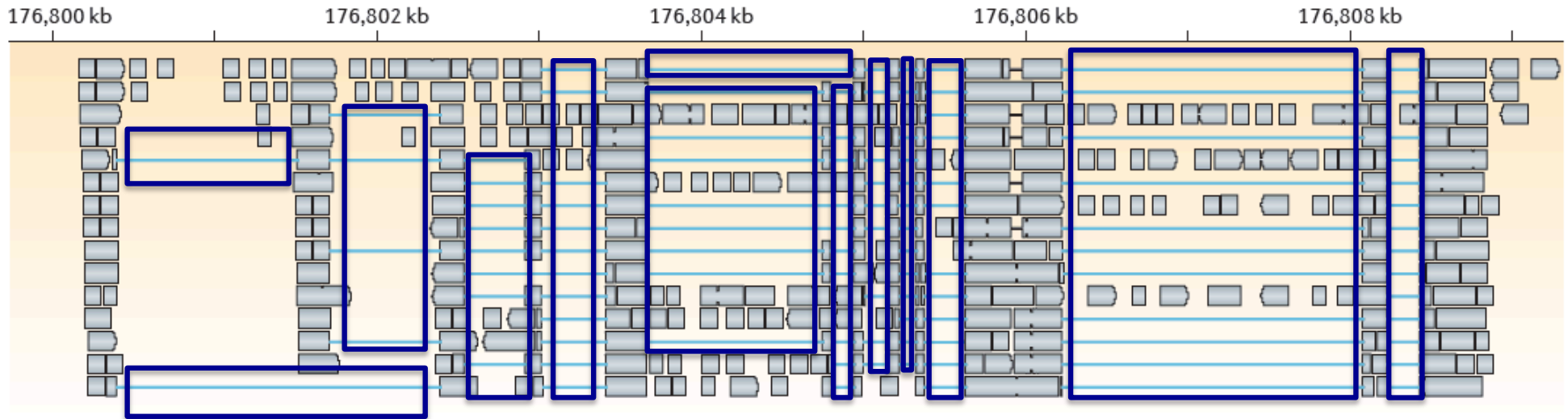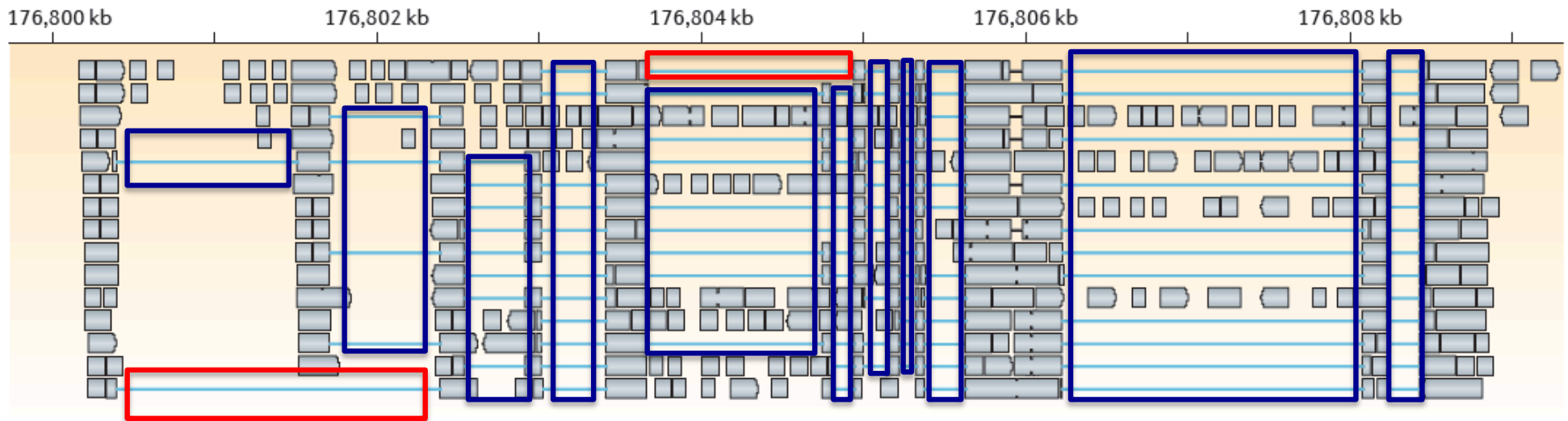
# Genome-Guided Transcript Reconstruction

## Splice-align reads to the genome



From Martin & Wang. Nature Reviews in Genetics. 2011    **bio**informatics.ca

# Genome-Guided Transcript Reconstruction

## Splice-align reads to the genome



Individual reads can yield multiple exon and intron segments (splice patterns)

From Martin & Wang. Nature Reviews in Genetics. 2011     **bio**informatics.ca

# Genome-Guided Transcript Reconstruction

**Splice-align reads to the genome**



Nodes = unique splice patterns
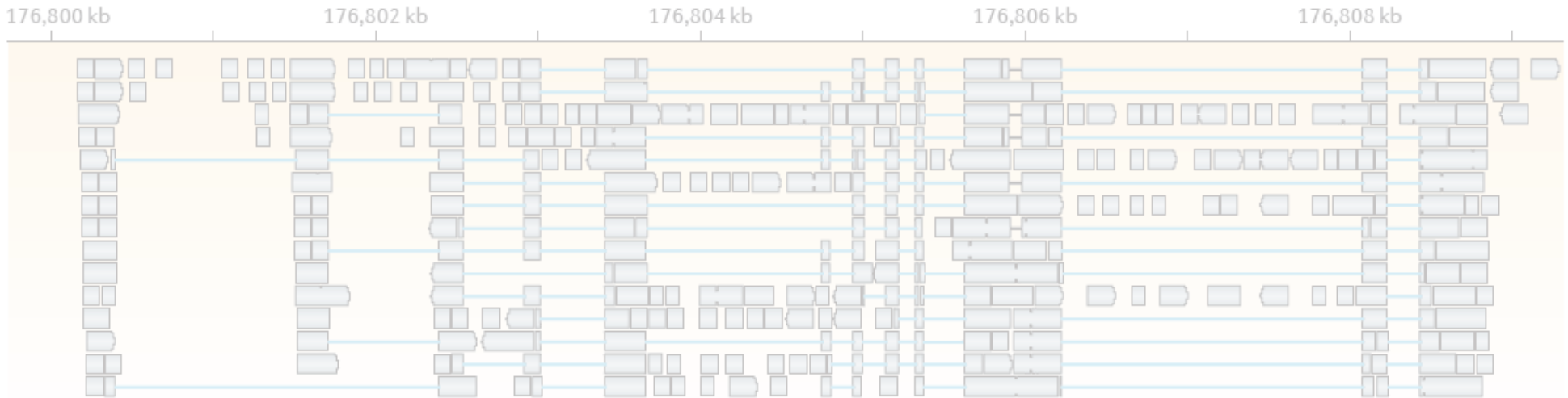
# Genome-Guided Transcript Reconstruction

**Splice-align reads to the genome**



**Construct graph from unique splice patterns of aligned reads.**



Nodes = unique splice patterns
Edges = compatible patterns

From Martin & Wang. Nature Reviews in Genetics. 2011    **bio**informatics.ca

# Genome-Guided Transcript Reconstruction

## Splice-align reads to the genome



176,800 kb          176,802 kb          176,804 kb          176,806 kb          176,808 kb

## Construct graph from unique splice patterns of aligned reads.
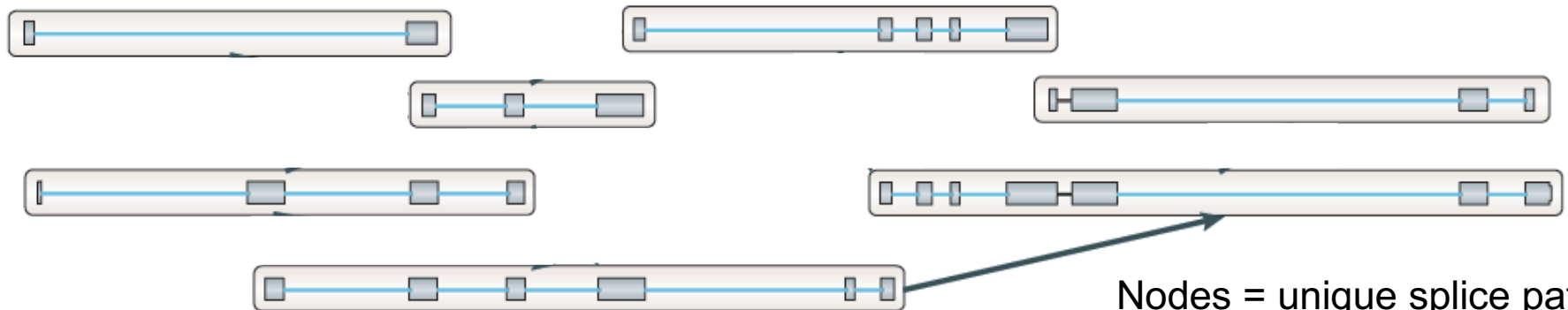


Nodes = unique splice patterns
Edges = compatible patterns

# Genome-Guided Transcript Reconstruction

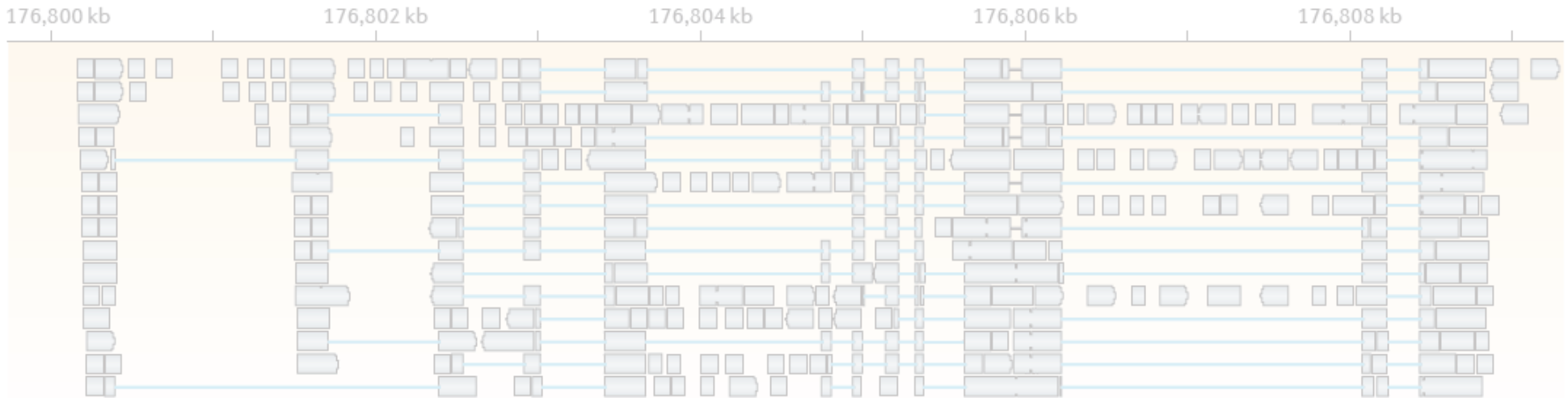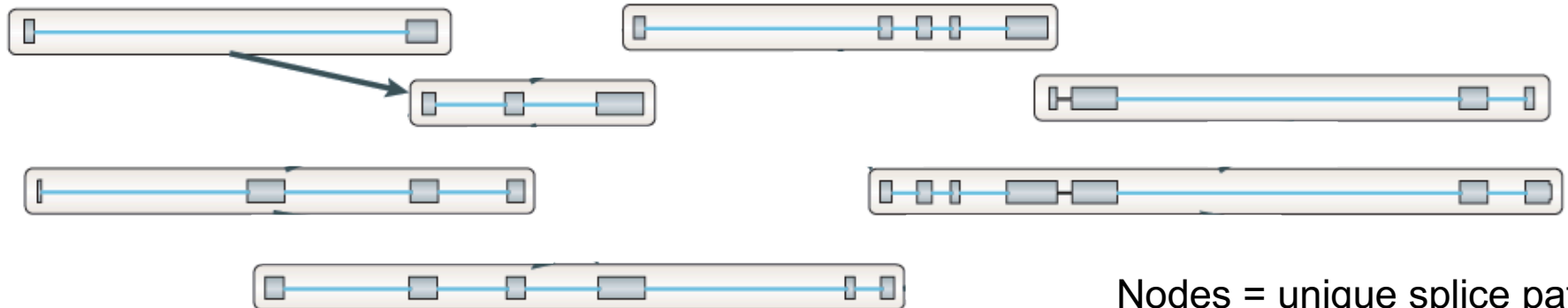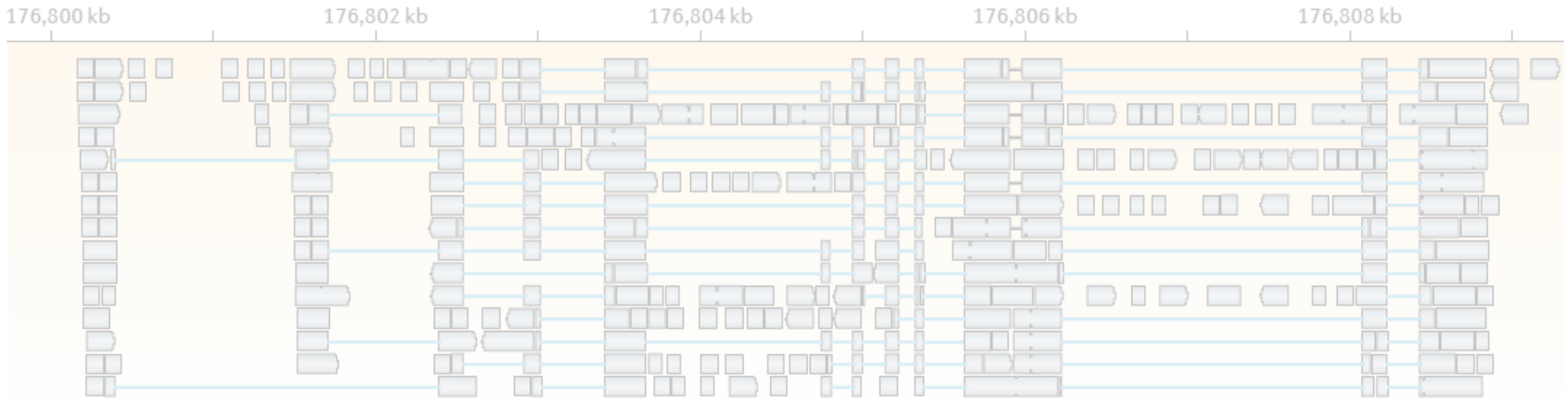## Splice-align reads to the genome



## Construct graph from unique splice patterns of aligned reads.



Nodes = unique splice patterns
Edges = compatible patterns

# Genome-Guided Transcript Reconstruction

**Traverse paths through the graph to assemble transcript isoforms**

# Genome-Guided Transcript Reconstruction

**Traverse paths through the graph to assemble transcript isoforms**



**Reconstructed isoforms**



From Martin & Wang. Nature Reviews in Genetics. 2011 **bio**informatics.ca

# What if you don't have a high quality reference genome sequence?

**Genome-free de novo transcript reconstruction to the rescue.**

**bio**informatics.ca

# Read Overlap Graph:   Reads as nodes, overlaps as edges

# Read Overlap Graph: Reads as nodes, overlaps as edges

Node = read
Edge = overlap

**bio**informatics.ca

# Read Overlap Graph:    Reads as nodes, overlaps as edges



Transcript A

Generate consensus sequence where reads overlap

Node = read
Edge = overlap

Transcript B

bioinformatics.ca

# Finding pairwise overlaps between *n* reads involves ~ *n²* comparisons.



*Impractical for typical RNA-Seq data (50M reads)*

# No genome to align to… De novo assembly required



Want to avoid $n^2$ read alignments to define overlaps

# Use a de Bruijn graph

**bio**informatics.ca

# Sequence Assembly via De Bruijn Graphs

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

k-mers (k=5)

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG          CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG  ⊦ Reads

From Martin & Wang, Nat. Rev. Genet. 2011

**bio**informatics.ca

# Sequence Assembly via De Bruijn Graphs

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

k-mers (k=5)

ACCGC

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG    CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG   Reads

From Martin & Wang, Nat. Rev. Genet. 2011

**bio**informatics.ca

# Sequence Assembly via De Bruijn Graphs

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

k-mers (k=5)

ACCGC

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG        CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG    Reads

**Construct the de Bruijn graph**

ACCGC

From Martin & Wang, Nat. Rev. Genet. 2011          Nodes = unique k-mers, Edges = overlap by (k-1)

# Sequence Assembly via De Bruijn Graphs

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

k-mers (k=5)

CCGCC

ACCGC

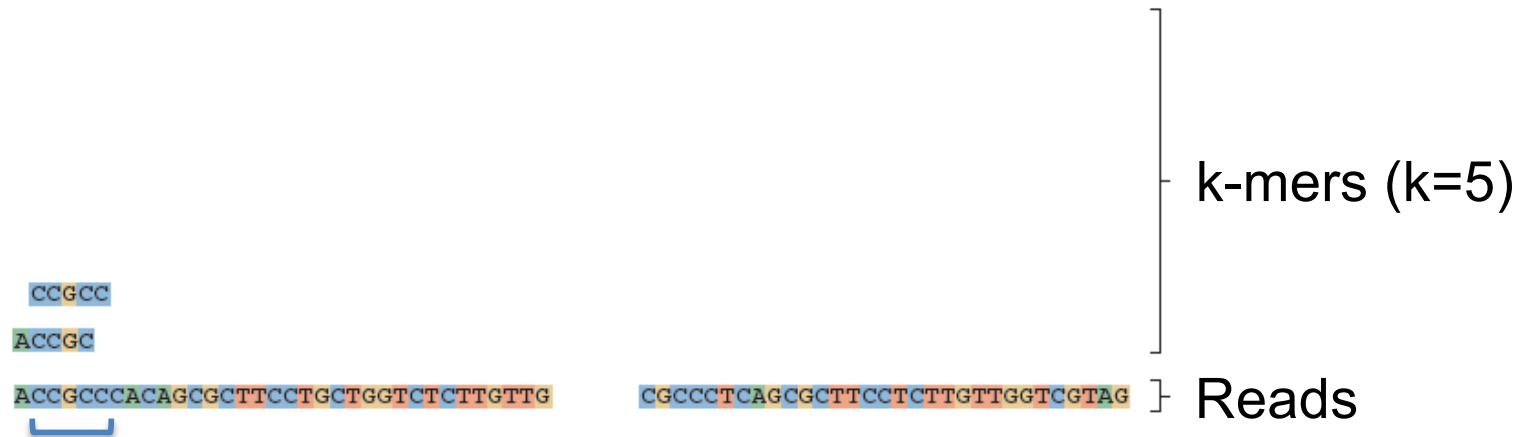ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG          CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG        Reads

**Construct the de Bruijn graph**

ACCGC

Nodes = unique k-mers, Edges = overlap by (k-1)

# Sequence Assembly via De Bruijn Graphs

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

k-mers (k=5)

(k-1) overlap

CCGCC

ACCGC

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG

CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG ⊦ Reads

**Construct the de Bruijn graph**

(ACCGC) (CCGCC)

Nodes = unique k-mers, Edges = overlap by (k-1)

**bio**informatics.ca

# Sequence Assembly via De Bruijn Graphs

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

k-mers (k=5)

(k-1) overlap

CCGCC
ACCGC

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG          CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG    Reads
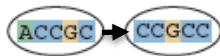
**Construct the de Bruijn graph**

ACCGC → CCGCC



Nodes = unique k-mers, Edges = overlap by (k-1)

# Sequence Assembly via De Bruijn Graphs
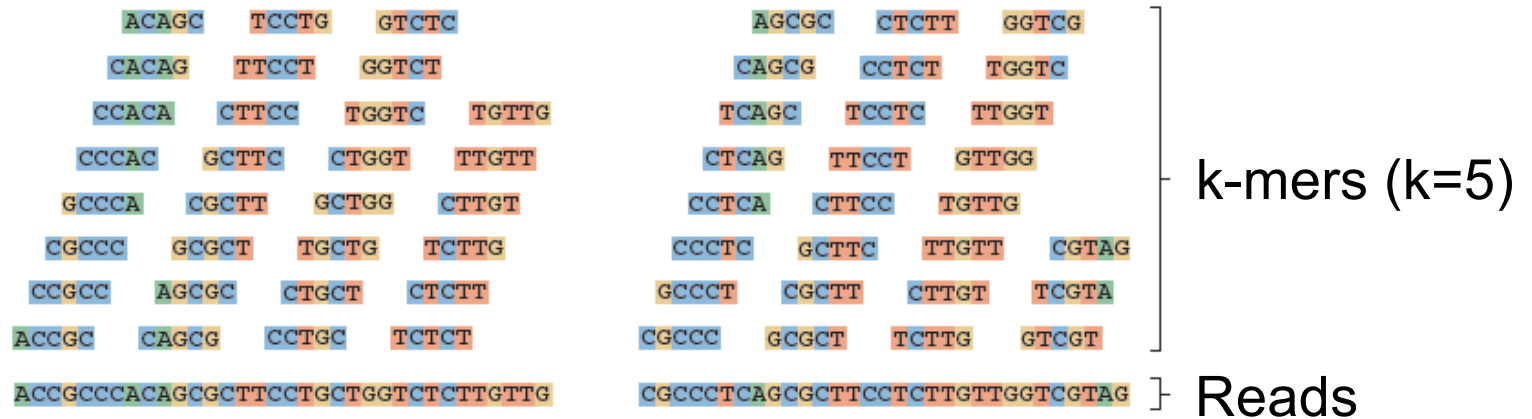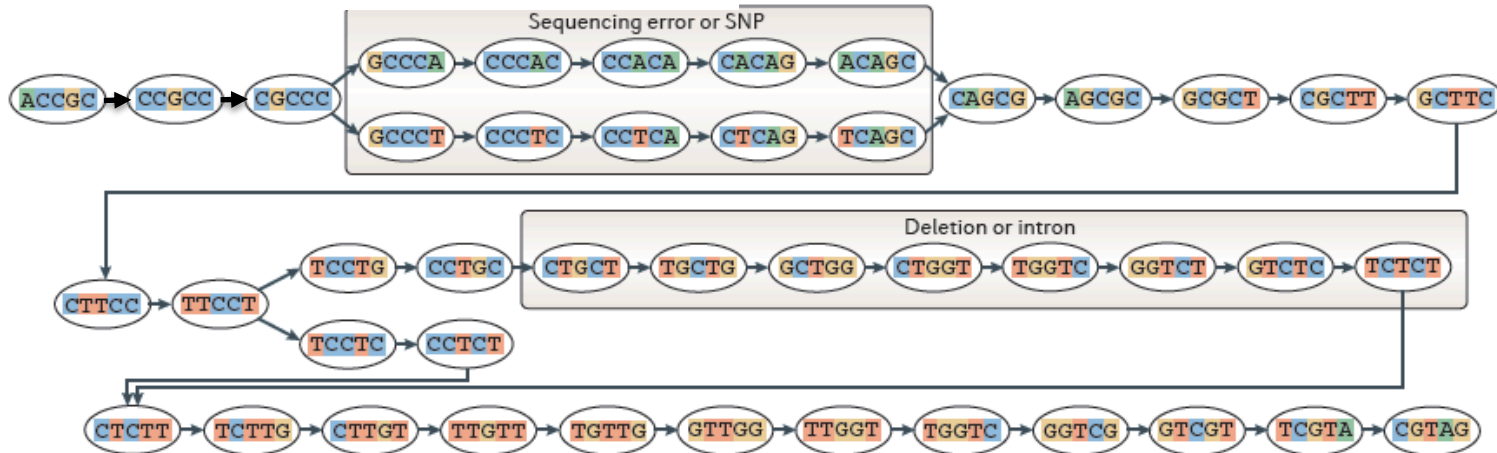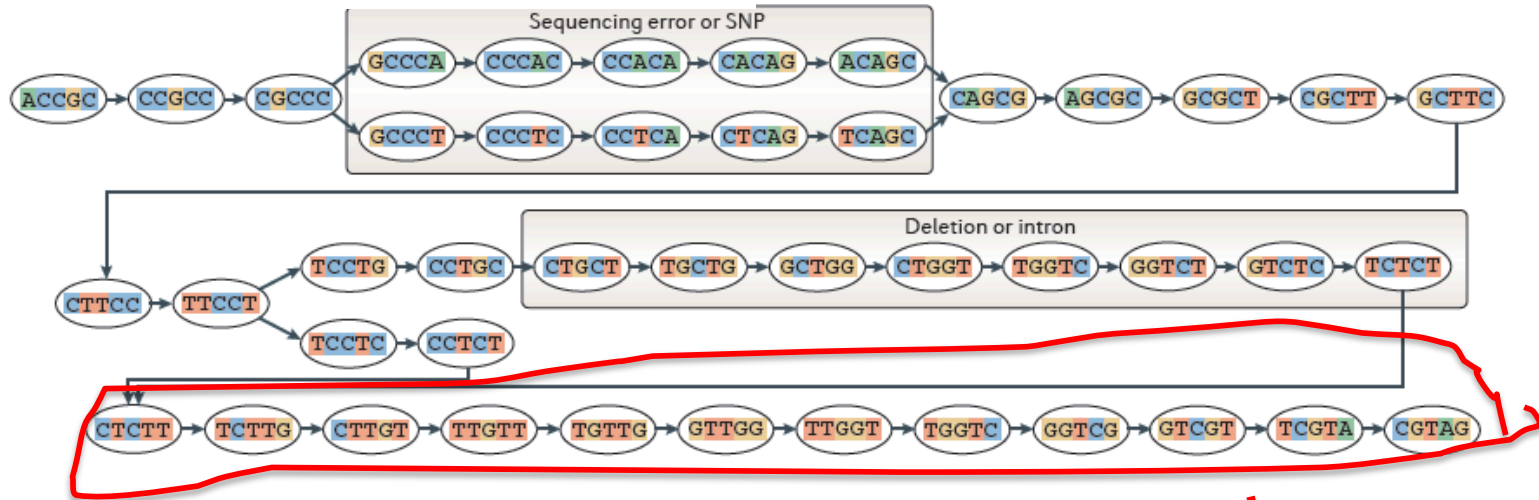
**Generate all substrings of length k from the reads**



k-mers (k=5)

Reads

**Construct the de Bruijn graph**

Nodes = unique k-mers, Edges = overlap by (k-1)

**Module 7**

**bio**informatics.ca

# Construct the de Bruijn graph



# Collapse the de Bruijn graph



From Martin & Wang, Nat. Rev. Genet. 2011

**bio**informatics.ca

# Collapse the de Bruijn graph



# Traverse the graph



# Assemble Transcript Isoforms



From Martin & Wang, Nat. Rev. Genet. 2011

**bio**informatics.ca

# Contrasting Genome and Transcriptome *De novo* Assembly

## Genome Assembly

- Uniform coverage

- Single contig per locus

- Assemble small numbers of large Mb-length chromosomes
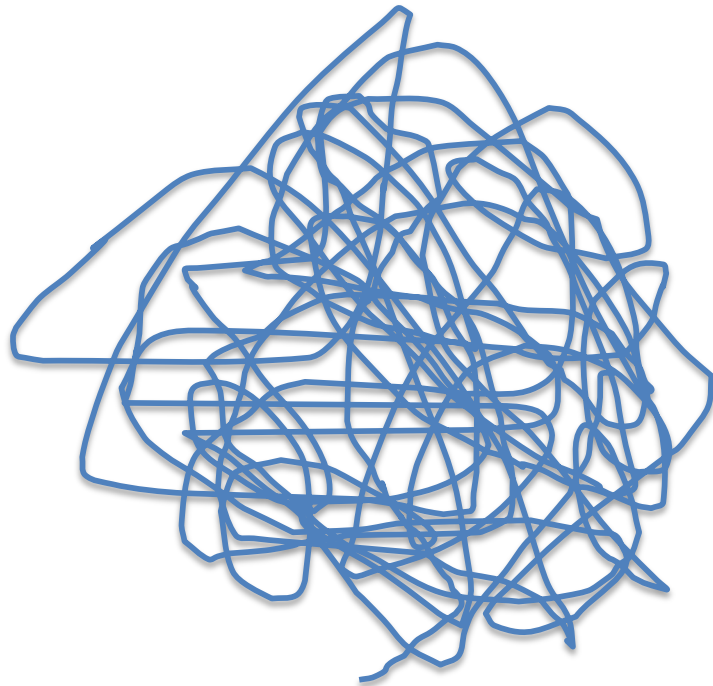
- Double-stranded data

## Transcriptome Assembly

- Exponentially distributed coverage levels

- Multiple contigs per locus (alt splicing)

- Assemble many thousands of Kb-length transcripts

- Strand-specific data available

**Trinity**

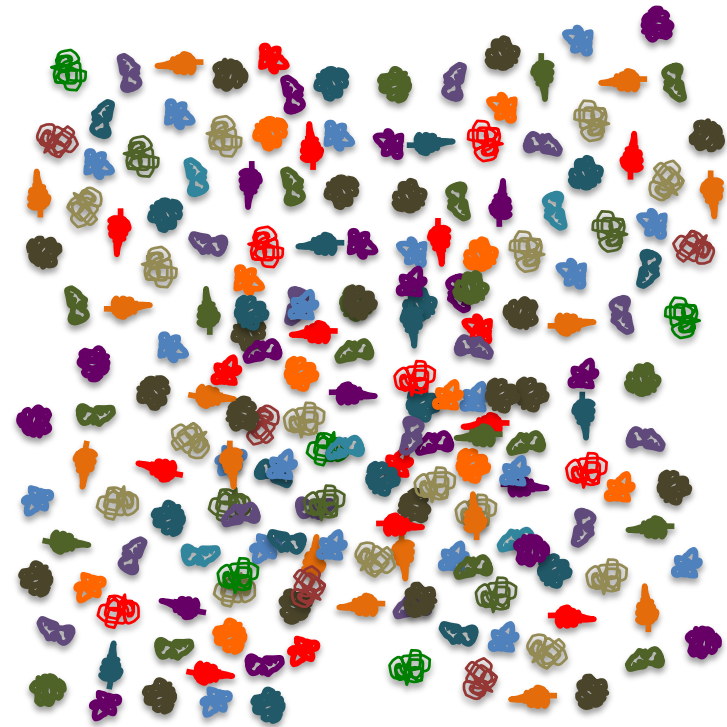# Trinity Aggregates Isolated Transcript Graphs
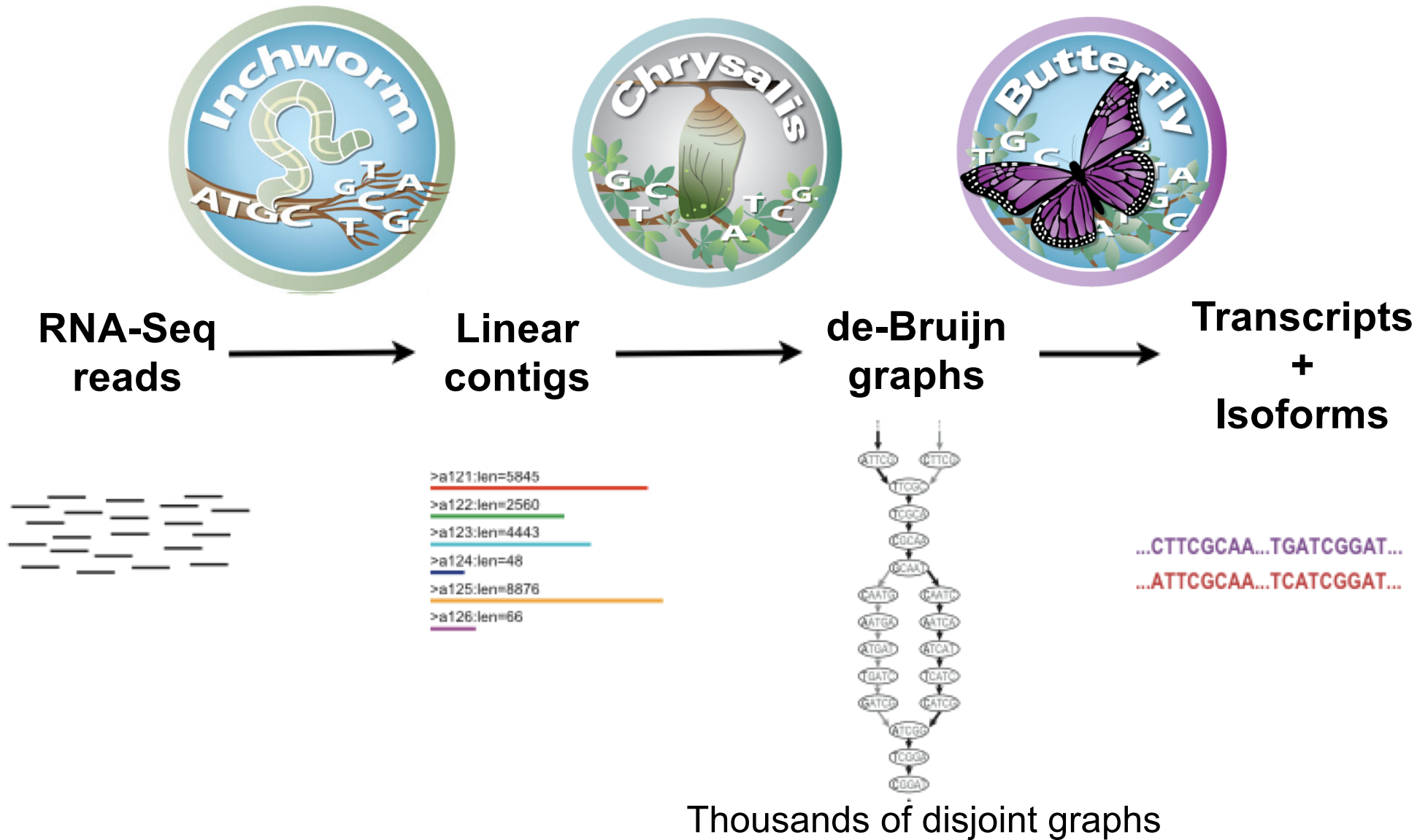
**Genome Assembly**
Single Massive Graph

**Trinity Transcriptome Assembly**
Many Thousands of Small Graphs



Entire chromosomes represented.

Ideally, one graph per expressed gene.

**bio**informatics.ca

# Trinity – How it works:



RNA-Seq reads → Linear contigs → de-Bruijn graphs → Transcripts + Isoforms

>a121:len=5845
>a122:len=2560
>a123:len=4443
>a124:len=48
>a125:len=8876
>a126:len=66

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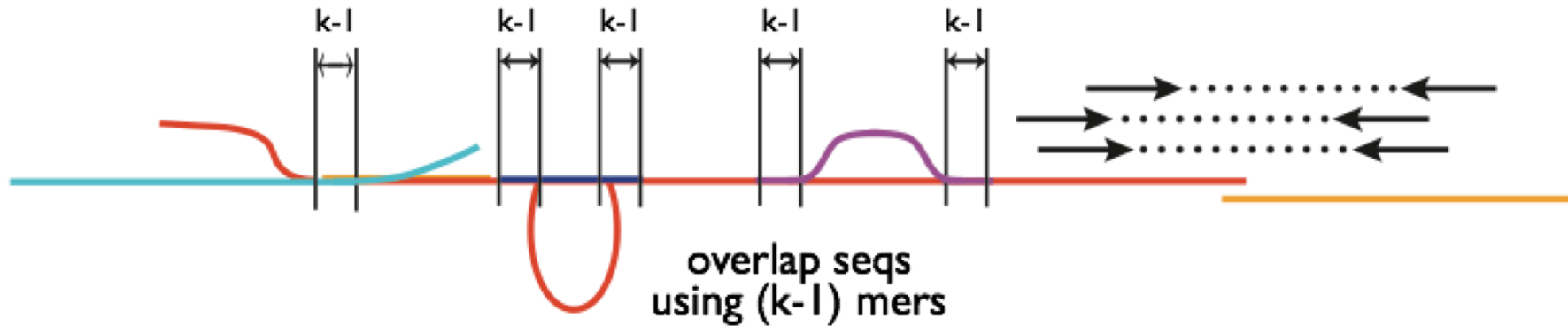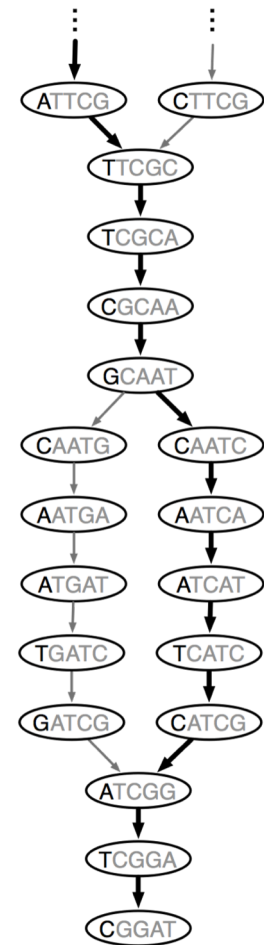...ATTCGCAA...TCATCGGAT...

Thousands of disjoint graphs

# Chrysalis

>a121:len=5845

>a122:len=2560

>a123:len=4443

>a124:len=48

>a125:len=8876

>a126:len=66

Integrate isoforms
via k-1 overlaps

Build de Bruijn Graphs
(ideally, one per gene)

overlap seqs
using (k-1) mers
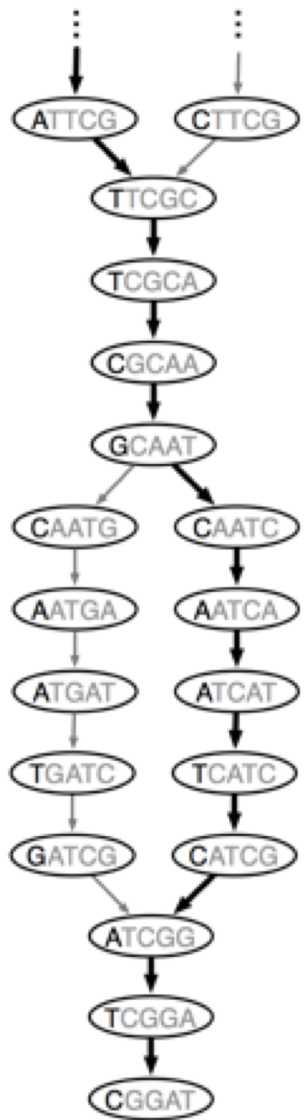
Thousands of Chrysalis Clusters

Butterfly

de Bruijn graph → compacting → compact graph → finding paths → compact graph with reads → extracting sequences → sequences (isoforms and paralogs)

..CTTCGCAA..TGATCGGAT...

..ATTCGCAA..TCATCGGAT...

# Butterfly Example 1:
# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph

# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph



Reconstructed Transcripts

**bio**informatics.ca

# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph



Reconstructed Transcripts

**bio**informatics.ca

# Reconstruction of Alternatively Spliced Transcripts



Butterfly's Compacted Sequence Graph

Reconstructed Transcripts

Aligned to Mouse Genome

Naa25 Nalpha acteyltransferase 25 (Reference structure)

**bio**informatics.ca

# Butterfly Example 2:
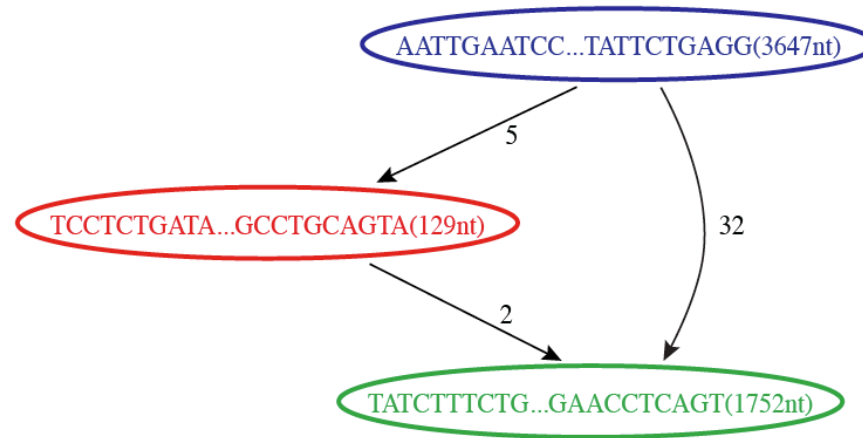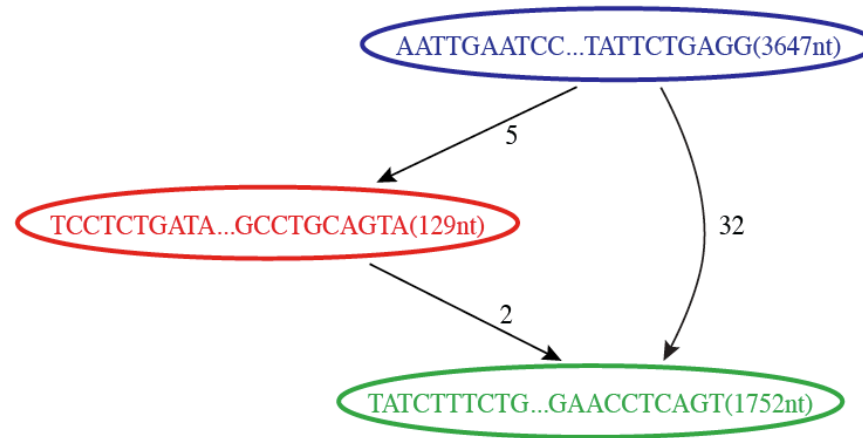# Teasing Apart Transcripts of Paralogous Genes

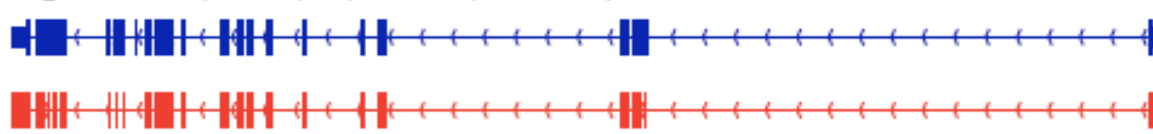# Butterfly Example 2:
# Teasing Apart Transcripts of Paralogous Genes

# Strand-specific RNA-Seq is Preferred

Computationally: fewer confounding graph structures in de novo assembly:

ex.  Forward != reverse complement

(GGAA != TTCC)

Biologically: separate sense vs. antisense transcription

## Comprehensive comparative analysis of strand-specific RNA sequencing methods

Joshua Z Levin[1,6], Moran Yassour[1–3,6], Xian Adiconis[1], Chad Nusbaum[1], Dawn Anne Thompson[1], Nir Friedman[3,4], Andreas Gnirke[1] & Aviv Regev[1,2,5]
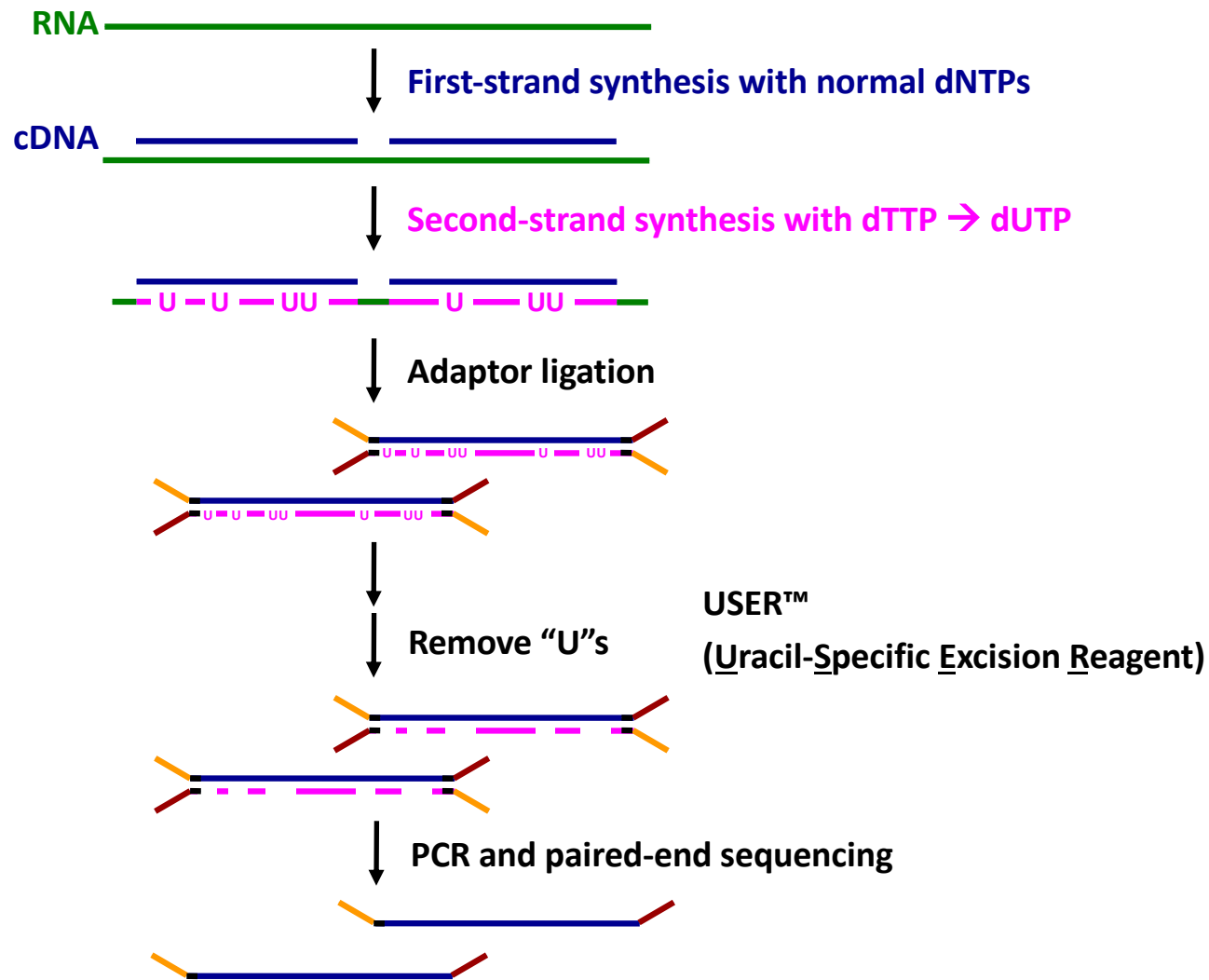
Strand-specific, massively parallel cDNA sequencing (RNA-seq) is a powerful tool for transcript discovery, genome annotation and expression profiling. There are multiple published methods

Nevertheless, direct information on the originating strand can substantially enhance the value of an RNA-seq experiment. For example, such information would help to accurately identify anti-

'dUTP second strand marking' identified as the leading protocol

computational pipeline to compare library quality metrics from any RNA-seq method. Using the well-annotated *Saccharomyces cerevisiae* transcriptome as a benchmark, we compared seven library-construction protocols, including both published and

boundaries of adjacent genes transcribed on opposite strands and resolve the correct expression levels of coding or noncoding over-lapping transcripts. These tasks are particularly challenging in small microbial genomes, prokaryotic and eukaryotic, in which

**Module**

**bio**informatics.ca

# dUTP 2^nd Strand Method:  Our Favorite

RNA ─────────────────

↓ **First-strand synthesis with normal dNTPs**

cDNA ─────────────────

↓ **Second-strand synthesis with dTTP → dUTP**

U ─ U ── UU ──── U ── UU

↓ **Adaptor ligation**

↓ **Remove "U"s**   **USER™**
**(Uracil-Specific Excision Reagent)**

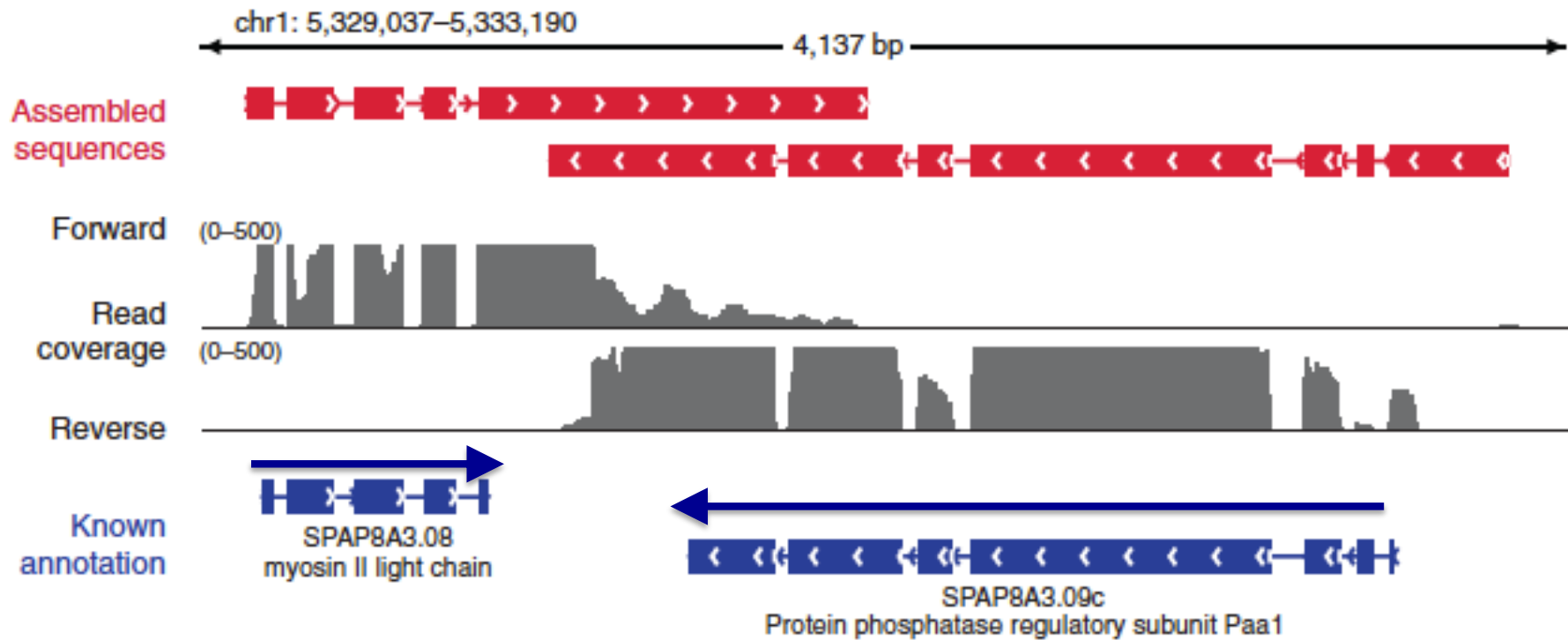↓ **PCR and paired-end sequencing**

**Modified from Parkhomchuk *et al.* (2009) *Nucleic Acids Res.* 37:e123**
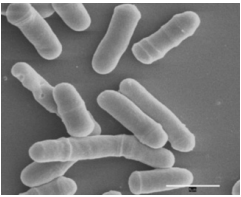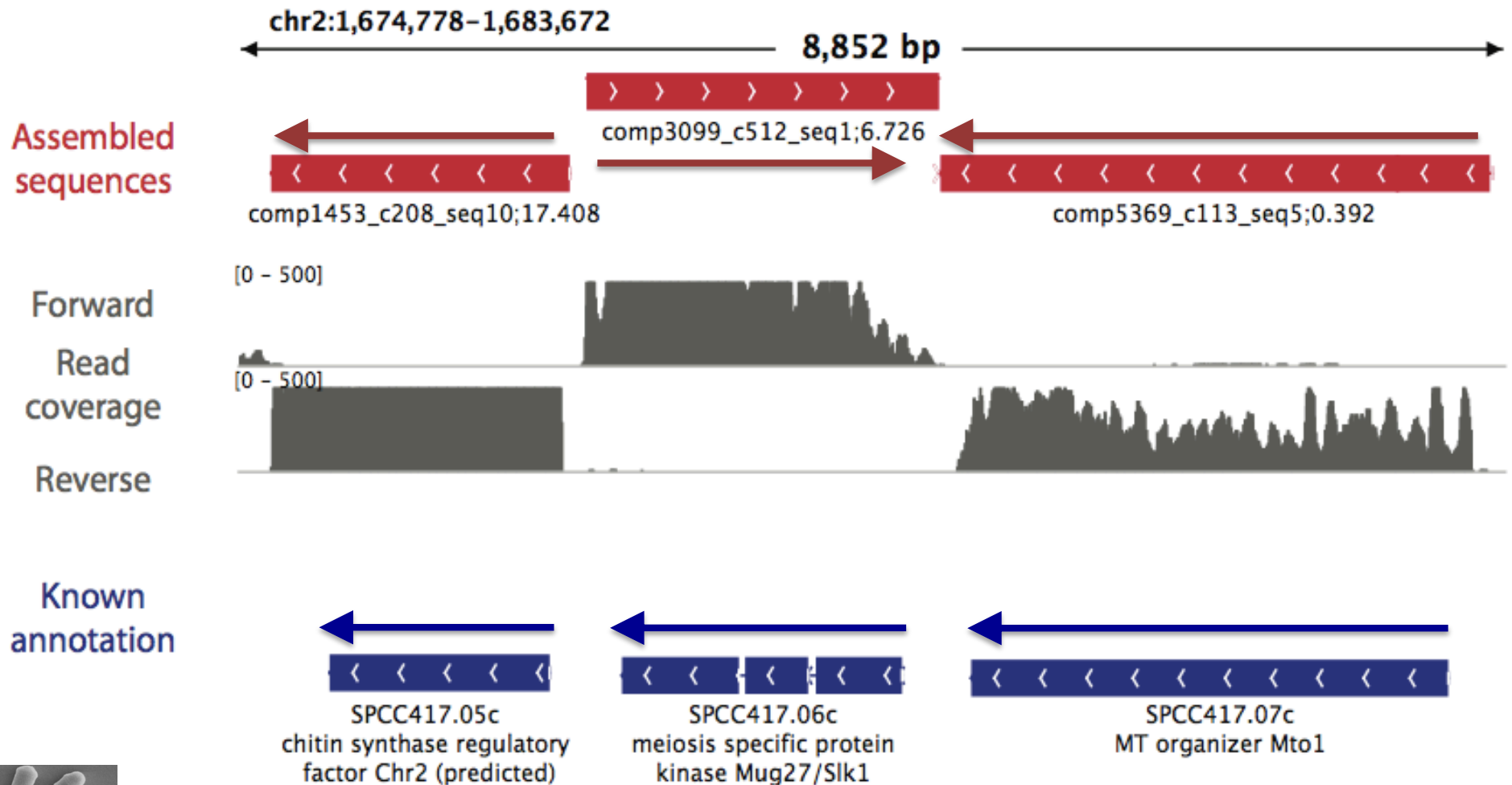
Slide from J. Levin

# Overlapping UTRs from Opposite Strands



*Schizosacharomyces pombe*
(fission yeast)

# Antisense-dominated Transcription

# Trinity output: a multi-fasta file
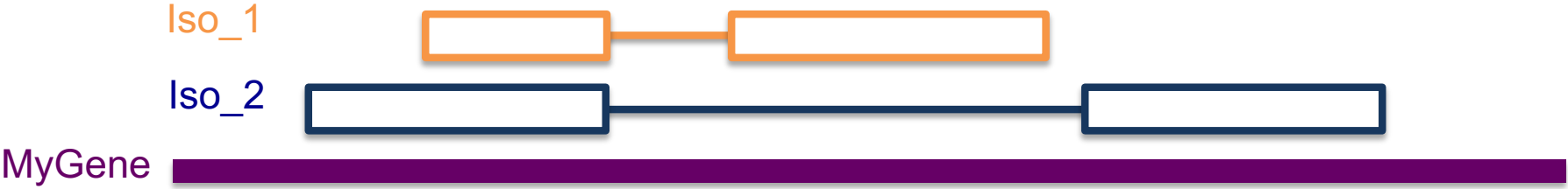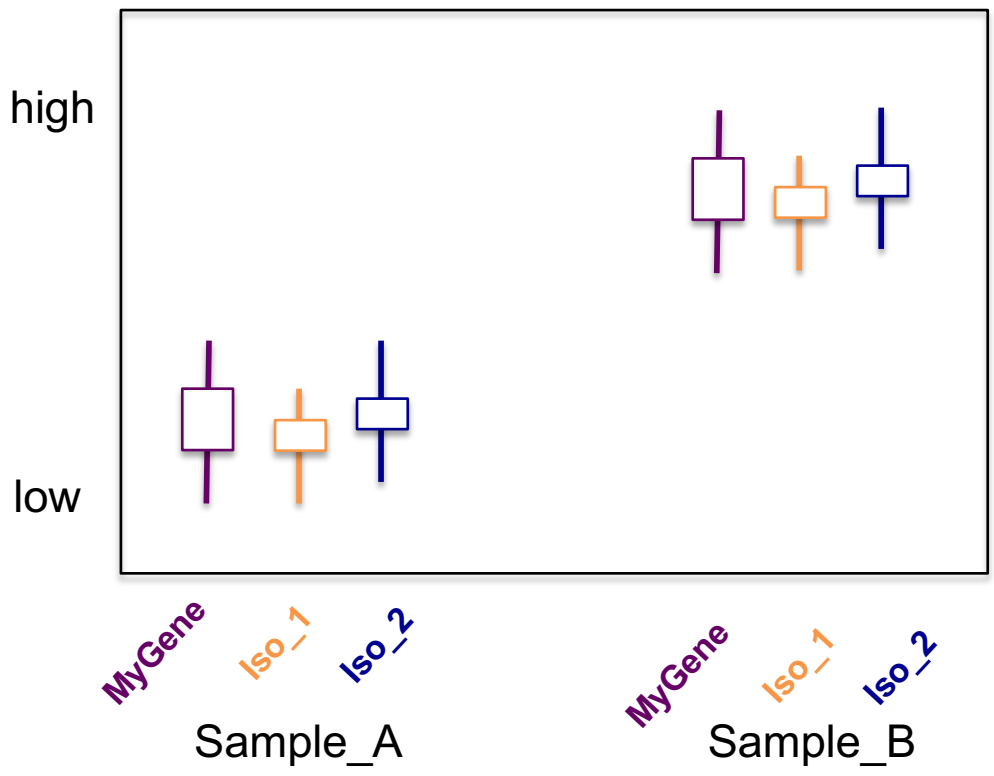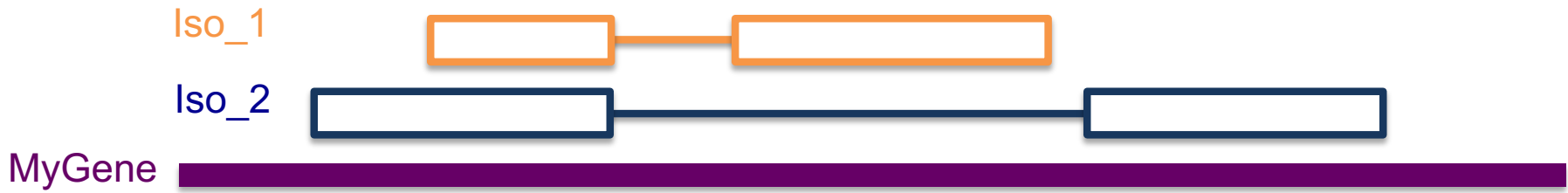
# Flavors of Differential Expression Analyses

- Transcripts:
  - Differential Transcript Expression (DTE)
  - Differential Transcript Usage (DTU)
  - Differential Exon Usage (DEU)
- Gene:
  - Differential Gene Expression (DGE)
  - Gene Differential Expression (GDE) ?

**bio**informatics.ca

# Differential Gene Expression (DGE) and Differential Transcript Expression (DTE)
## (Example 1)

# Differential Gene Expression (DGE) and Differential Transcript Expression (DTE) (Example 1)



| Feature | Diff Expressed? |
|---------|-----------------|
| MyGene | Yes |
| Iso_1 | Yes |
| Iso_2 | Yes |
| Diff. Transcript Usage ? (eg. Isoform switching) | No |

Differential Gene Expression (DGE) and Differential Transcript Expression (DTE)
(Example 2)

Differential Gene Expression (DGE) and Differential Transcript Expression (DTE)
(Example 3)

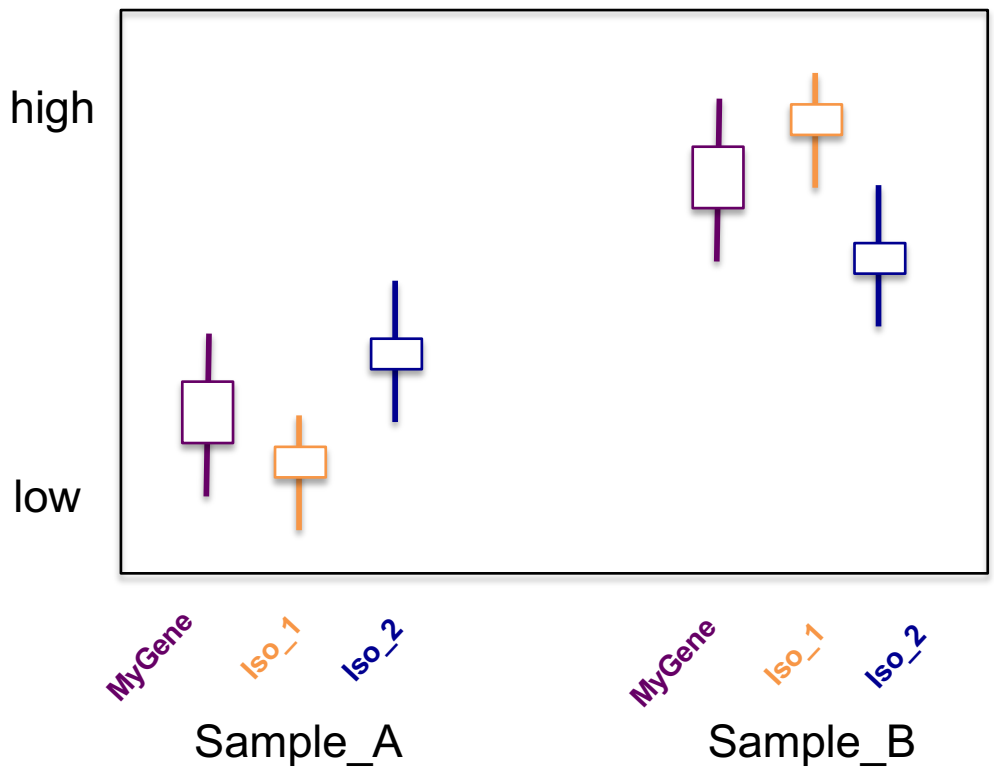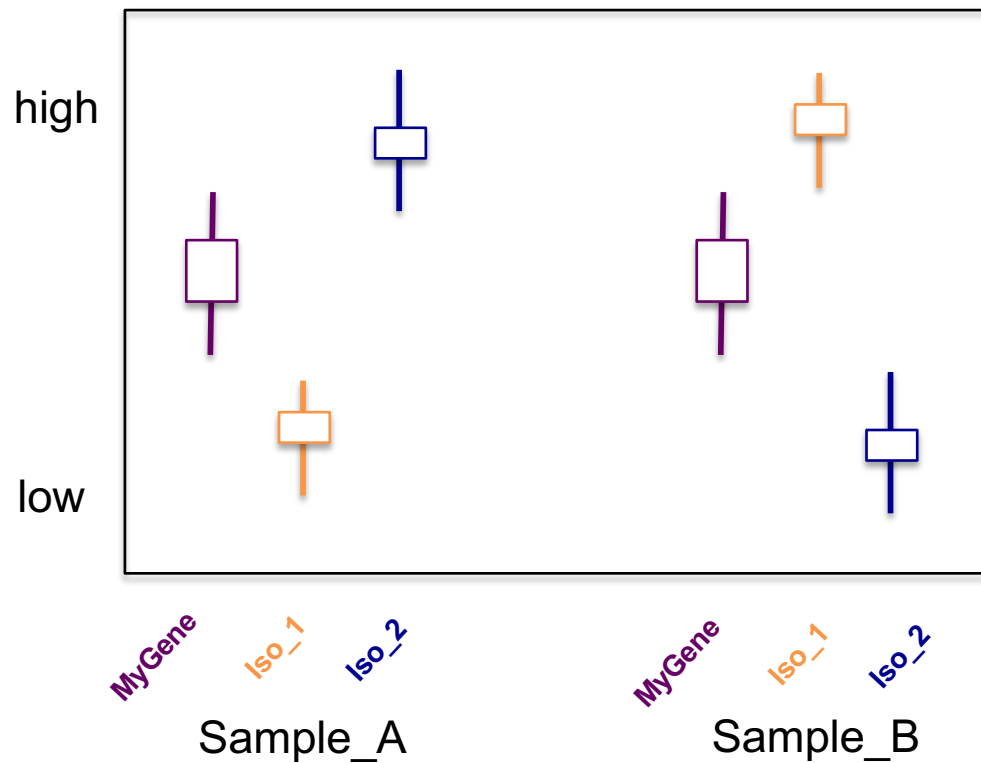| Feature | Diff Expressed? |
|---|---|
| MyGene | No |
| Iso_1 | Yes |
| Iso_2 | Yes |
| Diff. Transcript Usage ? (eg. Isoform switching) | Yes |

# Differential Gene Expression (DGE) and Differential Transcript Expression (DTE)
## (Example 3)

From Gene-level view (DGE):  not apparent
From Transcript-level view (GDE):  Yes, gene should be acknowledged as having changed.

Prevailing viewpoint:
DTE or DTU -> Gene Diff Expressed (GDE)

| Feature | Diff Expressed? |
|---|---|
| MyGene | No? |
| Iso_1 | Yes |
| Iso_2 | Yes |
| Diff. Transcript Usage ? (eg. Isoform switching) | Yes |

high

low

MyGene  Iso_1  Iso_2        MyGene  Iso_1  Iso_2

Sample_A              Sample_B

# Clarifying view:   (DTE or DTU or DGE) as special cases of Ge



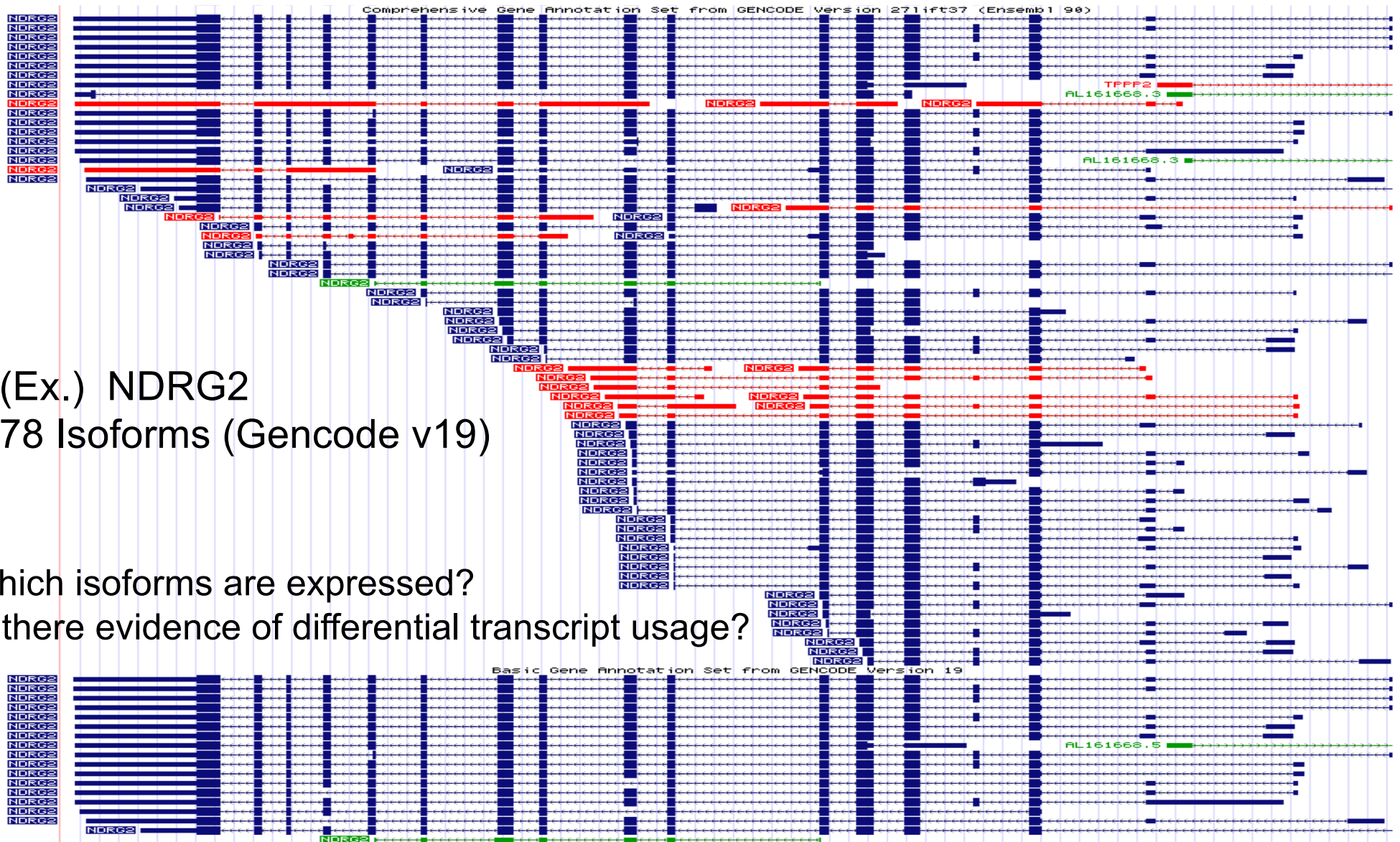DTE:  differential transcript expression
DTU:  differential transcript usage
DGE: differential gene expression (gene-level analysis)
GDE: gene differential expression (transcript-level analysis)

Ntranos, Yi, et al., 2018 – see supp.

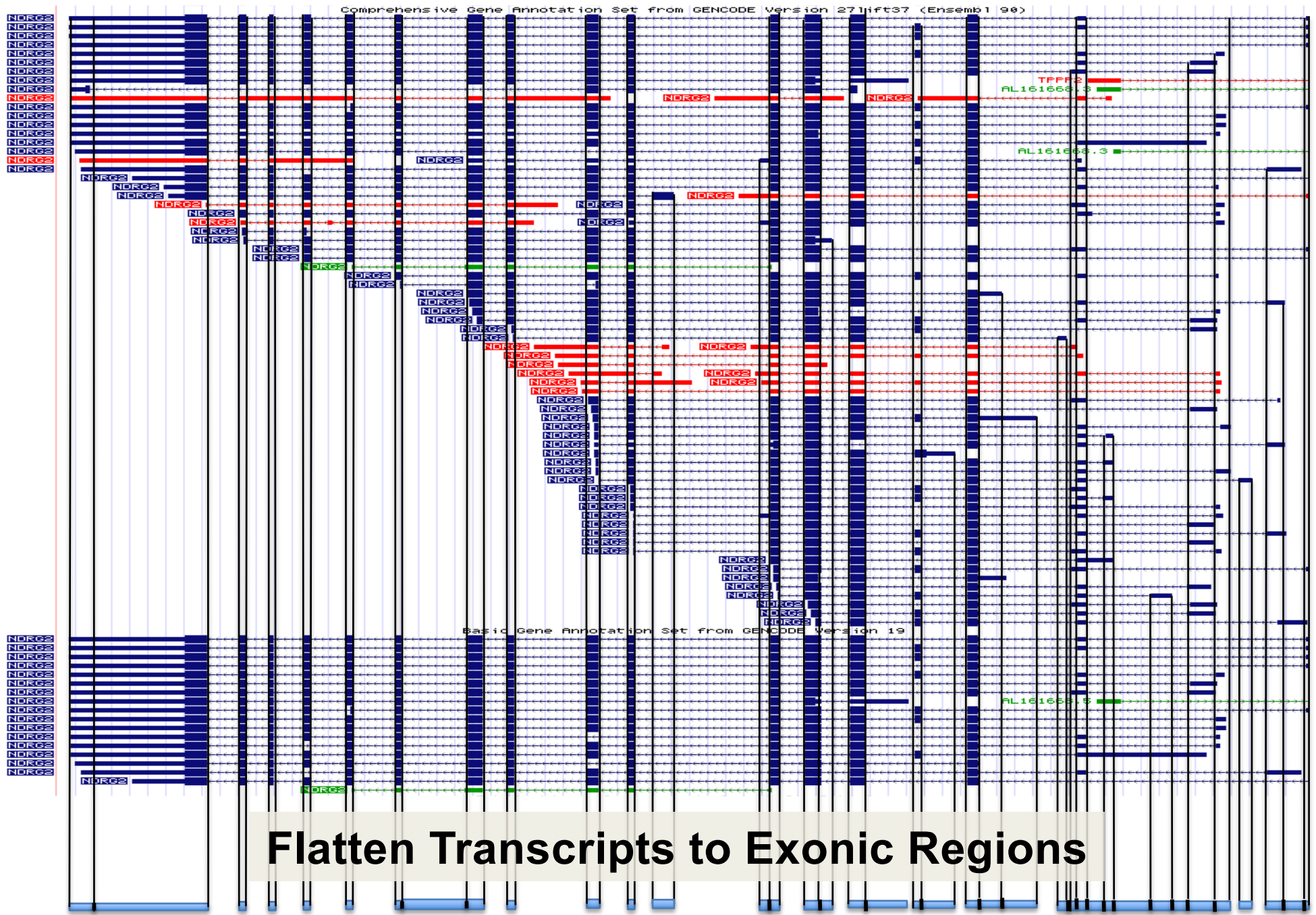See Lior Pachter's blog post:   https://liorpachter.wordpress.com/2019/01/07/fast-and-accurate-gene-differential-expression-by-testing-transcript-compatibility-counts/

# High Confidence Differential Transcript Expression is Difficult to Attain With Many Candidate Isoforms



(Ex.)  NDRG2
78 Isoforms (Gencode v19)

Which isoforms are expressed?
Is there evidence of differential transcript usage?

# Measure Differential Transcript Usage (DTU) via Differential Exon Usage (DEU)



Comprehensive Gene Annotation Set from GENCODE Version 27lift37 (Ensembl 90)

Basic Gene Annotation Set from GENCODE Version 19

## Flatten Transcripts to Exonic Regions

# Measure Differential Transcript Usage (DTU) via Differential Exon Usage (DEU)

**bio**informatics.ca

## Detecting differential usage of exons from RNA-seq data

Simon Anders,[1,2] Alejandro Reyes,[1] and Wolfgang Huber

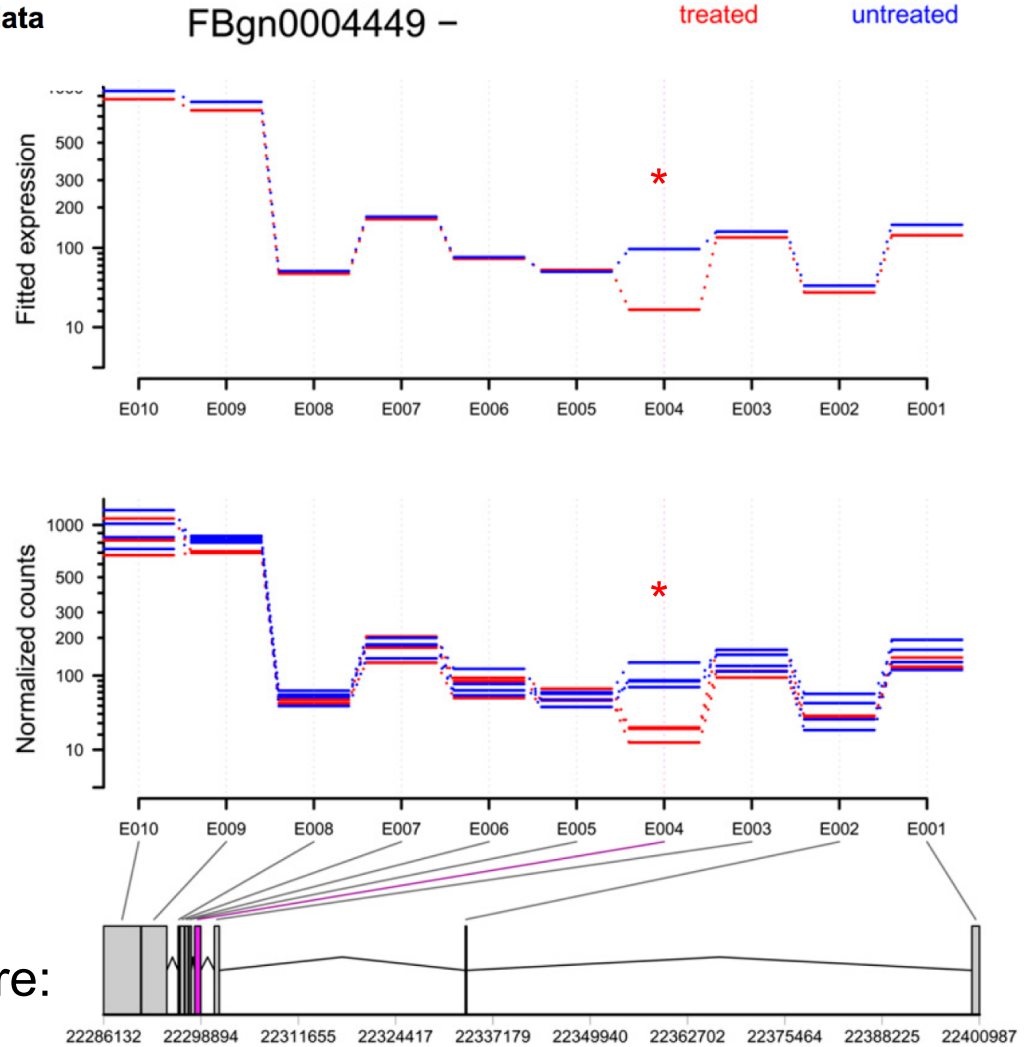Averaged Replicates

Each Replicate

Flattened gene structure:



**Figure 3.** The treatment of knocking down the splicing factor *pasilla* affects the fourth exon (counting bin E004) of the gene *Ten-m* (CG5723). (*Top* panel) Fitted values according to the linear model; (*middle* panel) normalized counts for each sample; (*bottom* panel) flattened gene model. (Red) Data for knockdown samples; (blue) control.

Module                                            bioinformatics.ca

# Enabling Differential Transcript Usage Analysis for De novo Transcriptome Assemblies

**METHOD**

**Open Access**

CrossMark

## SuperTranscripts: a data driven reference for analysis and visualisation of transcriptomes

Nadia M. Davidson[1,2]*, Anthony D. K. Hawkins[1] and Alicia Oshlack[1,2]*

# Enabling Differential Transcript Usage Analysis for De novo Transcriptome Assemblies
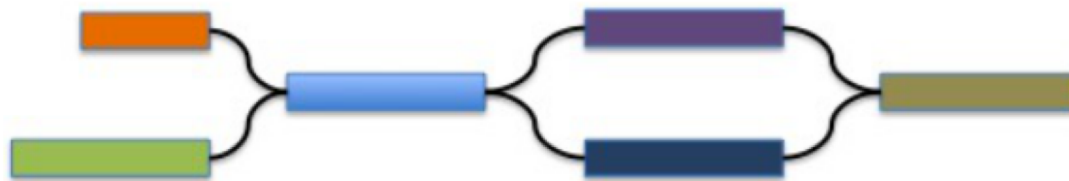
**METHOD**  **Open Access**

CrossMark

## SuperTranscripts: a data driven reference for analysis and visualisation of transcriptomes

Nadia M. Davidson[1,2]*, Anthony D. K. Hawkins[1] and Alicia Oshlack[1,2]*

Transcript splice graph:



Similar method and protocols now integrated into Trinity:
https://github.com/trinityrnaseq/trinityrnaseq/wiki/SuperTranscripts

# Enabling Differential Transcript Usage Analysis for De novo Transcriptome Assemblies
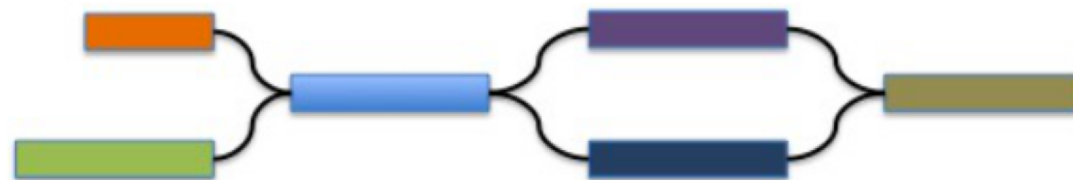
## SuperTranscripts: a data driven reference for analysis and visualisation of transcriptomes

Nadia M. Davidson[1,2]*, Anthony D. K. Hawkins[1] and Alicia Oshlack[1,2]*

Transcript splice graph:



Linearize graph via topological sorting
or graph multiple alignment

SuperTranscript:

DEXseq for DTU,
GATK for Variant Detection

Similar method and protocols now integrated into Trinity:
https://github.com/trinityrnaseq/trinityrnaseq/wiki/SuperTranscripts

**Module**

**bioinformatics**.ca

# Time for Transcript Reconstruction Lab

**bio**informatics.ca

# We are on a Coffee Break & Networking Session

Workshop Sponsors: