# Canadian Bioinformatics Workshops

www.bioinformatics.ca

Afrikaans български Català Dansk Deutsch Ελληνικά English English (CA) English (GB) English (US) Esperanto Castellano Castellano (AR) Español (CL) Castellano (CO) Español (Ecuador) Castellano (MX) Castellano (PE) Euskara Suomeksi français français (CA) Galego עברית hrvatski Magyar Italiano 日本語 한국어 Macedonian Melayu Nederlands Norsk Sesotho sa Leboa polski Português română slovenski jezik српски srpski (latinica) Sotho svenska 中文 華語 (台灣) isiZulu

# creative commons

## Attribution-Share Alike 2.5 Canada

### You are free:

**to Share** — to copy, distribute and transmit the work

**to Remix** — to adapt the work

*Free Cultural Works* **APPROVED FOR**

### Under the following conditions:

**Attribution**. You must attribute the work in the manner specified by the author or licensor (but not in any way that suggests that they endorse you or your use of the work).

**Share Alike**. If you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar licence to this one.

- For any reuse or distribution, you must make clear to others the licence terms of this work.
- Any of the above conditions can be waived if you get permission from the copyright holder.
- The author's moral rights are retained in this licence.

Disclaimer

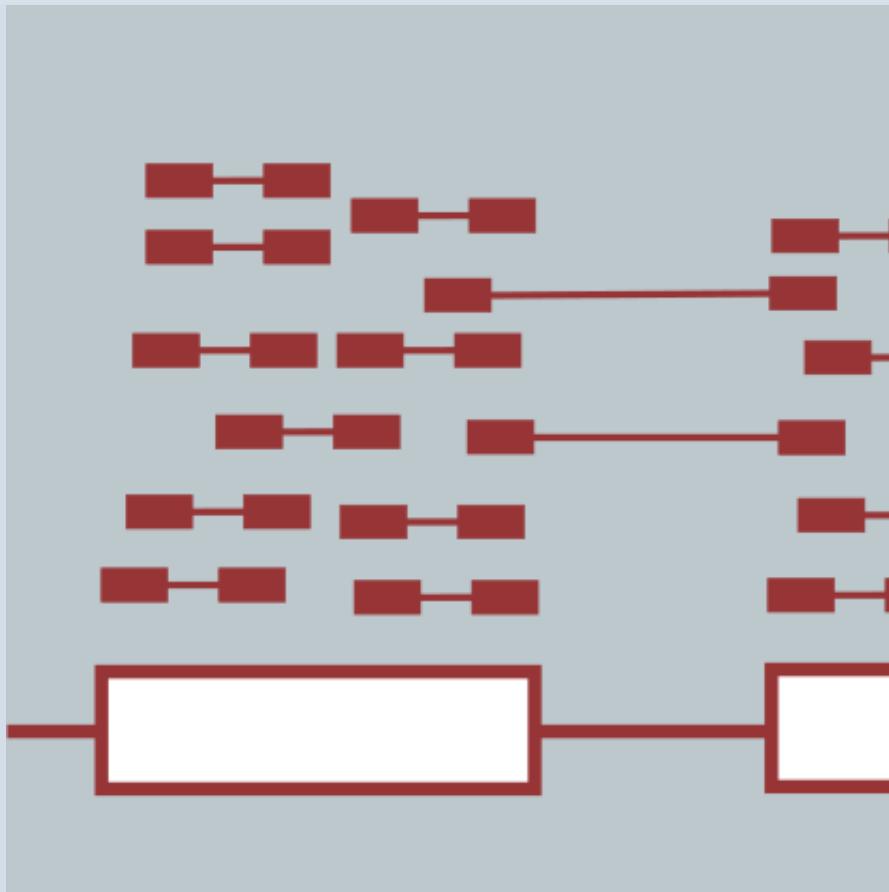**Your fair dealing and other rights are in no way affected by the above.**

This is a human-readable summary of the Legal Code (the full licence) available in the following languages:
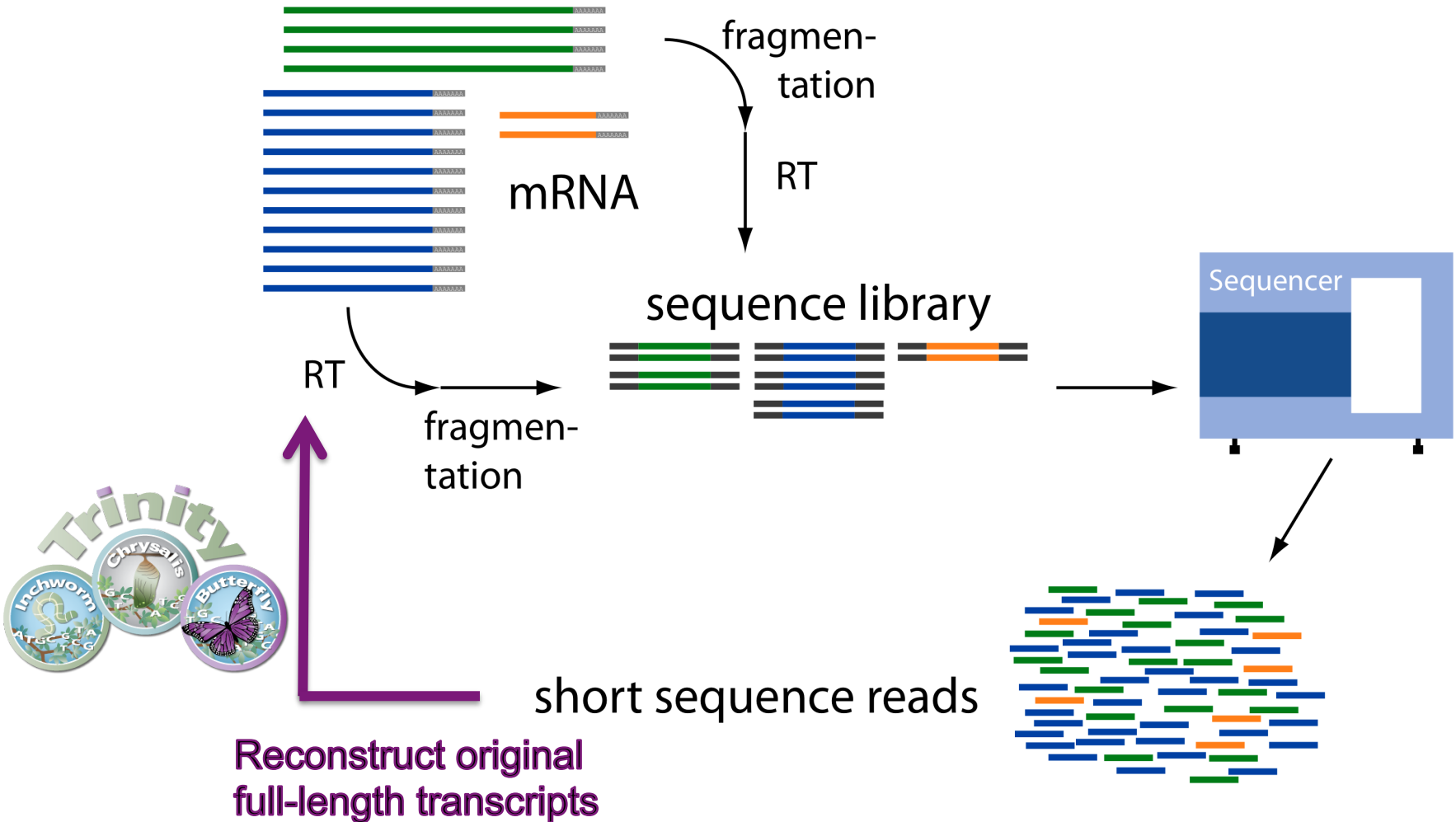English French

Learn how to distribute your work using this licence

# 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.

- Learn various ways to assess the quality of an assembled transcriptome.

# Assembly Required



mRNA

fragmen-
tation

RT

RT

fragmen-
tation

sequence library

Sequencer

short sequence reads

**Reconstruct original
full-length transcripts**

Adapted from G. Raetsch

**bio**informatics.ca

# Transcript Reconstruction from RNA-Seq Reads

RNA-Seq reads

# Transcript Reconstruction from RNA-Seq Reads



RNA-Seq reads

Align reads to genome

Genome

# Transcript Reconstruction from RNA-Seq Reads



RNA-Seq reads

Align reads to genome

Genome

Assemble transcripts from spliced alignments

# 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

# Transcript Reconstruction from RNA-Seq Reads

RNA-Seq reads

**Many tools to choose among:**

Align reads to genome

TopHat
STAR
HISAT
GSNAP
…

Assemble transcripts *de novo*

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

Genome

Assemble transcripts from spliced alignments

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

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

GTC

AGTCA

GATC

CG

GATTACA

CGATCA

AGC

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

# Graph Data Structures Commonly Used For Assembly



RNA-Seq reads

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

Reads to Graph

GTC

AGTCA

GATC

CG

GATTACA

CGATCA

AGC

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

**GATCGTCCGAGCGATTACA**

# 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
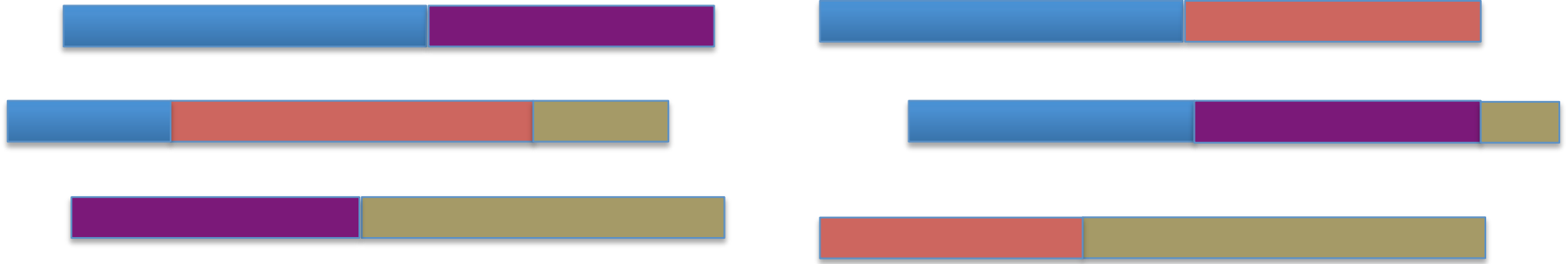
Node = read
Edge = overlap

Transcript B

# Finding pairwise overlaps between *n* reads involves ~ $n^2$ comparisons.



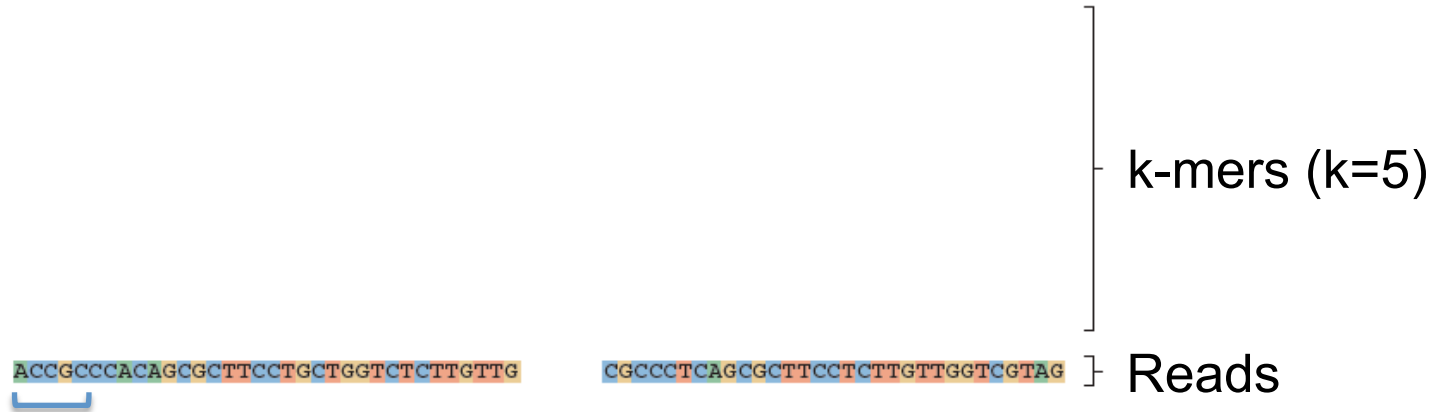*Impractical for typical RNA-Seq data (50M reads)*

# No genome to align to… De novo assembly required



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

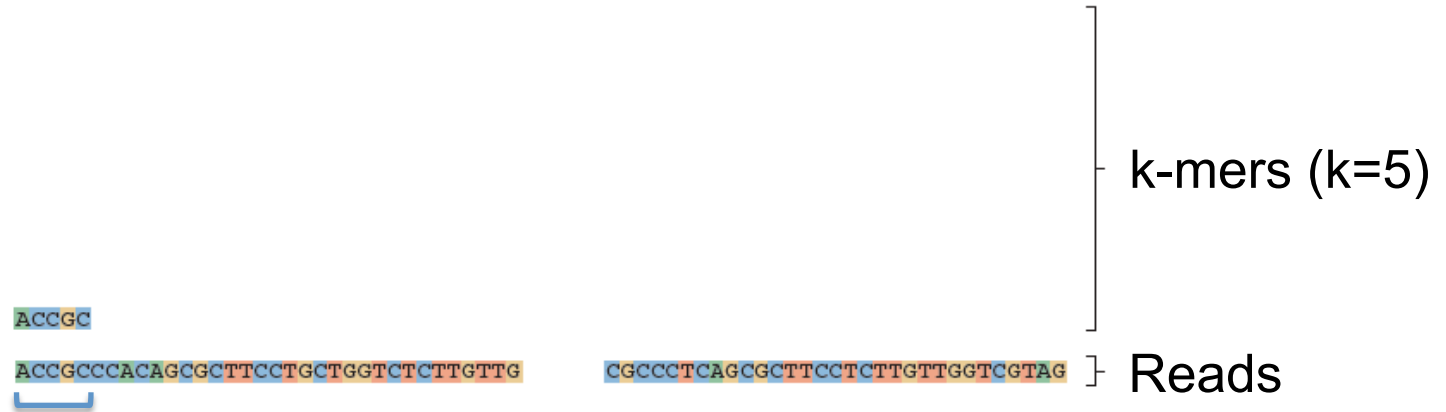## Use a de Bruijn graph

# 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

# Sequence Assembly via De Bruijn Graphs

**Generate all substrings of length k from the reads**
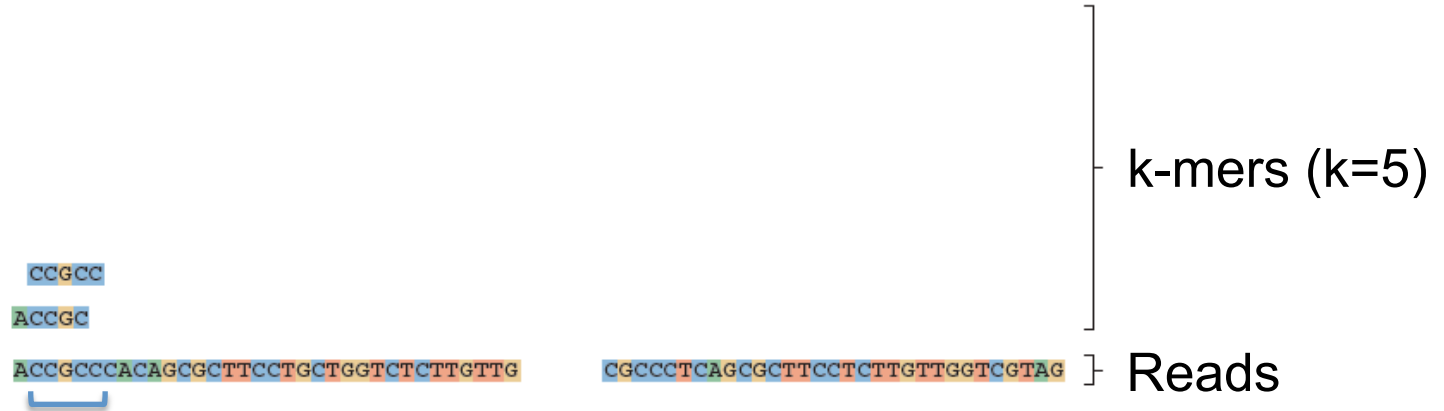
k-mers (k=5)

ACCGC

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG          CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG          Reads

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

# 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

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)

CCGCC

ACCGC

ACCGCCCACAGCGCTTCCTGCTGGTCTCTTGTTG    CGCCCTCAGCGCTTCCTCTTGTTGGTCGTAG    Reads

## Construct the de Bruijn graph

ACCGC

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)

**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)

**Module**

**bio**informatics.ca

# 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)

**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

**bio**informatics.ca

# 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

Thousands of disjoint graphs

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

# Inchworm Algorithm

- Decompose all reads into overlapping Kmers => hashtable(kmer, count)

Read:   **AATGTGAAAACTGGATTACATGCTGGTATGTC**…

**AATGTGA**

**ATGTGAA**   Overlapping kmers of length (k)

**TGTGAAA**

…

**Kmer Catalog (hashtable)**

| Kmer | Count among all reads |
|---------|-----------------------|
| AATGTGA | 4 |
| ATGTGAA | 2 |
| TGTGAAA | 1 |
| GATTACA | 9 |

**bio**informatics.ca

# Inchworm Algorithm

- Decompose all reads into overlapping Kmers => hashtable(kmer, count)

- Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.

**GATTACA**
9

**Kmer Catalog (hashtable)**

| Kmer | Count among all reads |
|---------|-----------------------|
| AATGTGA | 4 |
| ATGTGAA | 2 |
| TGTGAAA | 1 |
| GATTACA | 9 |

**bio**informatics.ca

# Inchworm Algorithm

- Decompose all reads into overlapping Kmers => hashtable(kmer, count)

- Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.

- Extend kmer at 3' end, guided by coverage.

# Inchworm Algorithm

**bio**informatics.ca
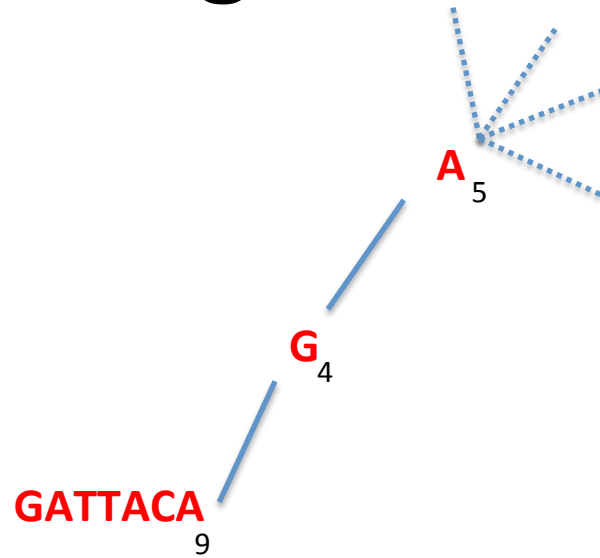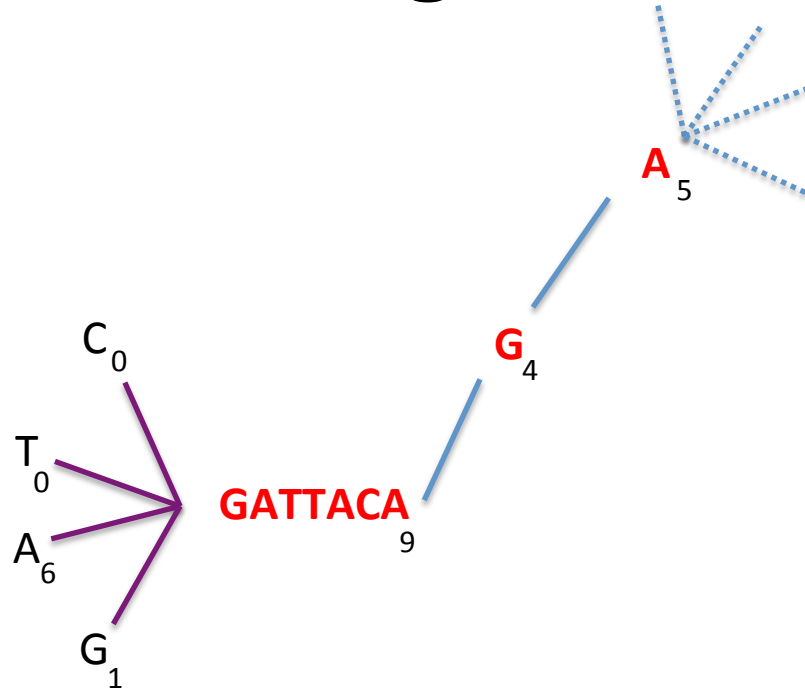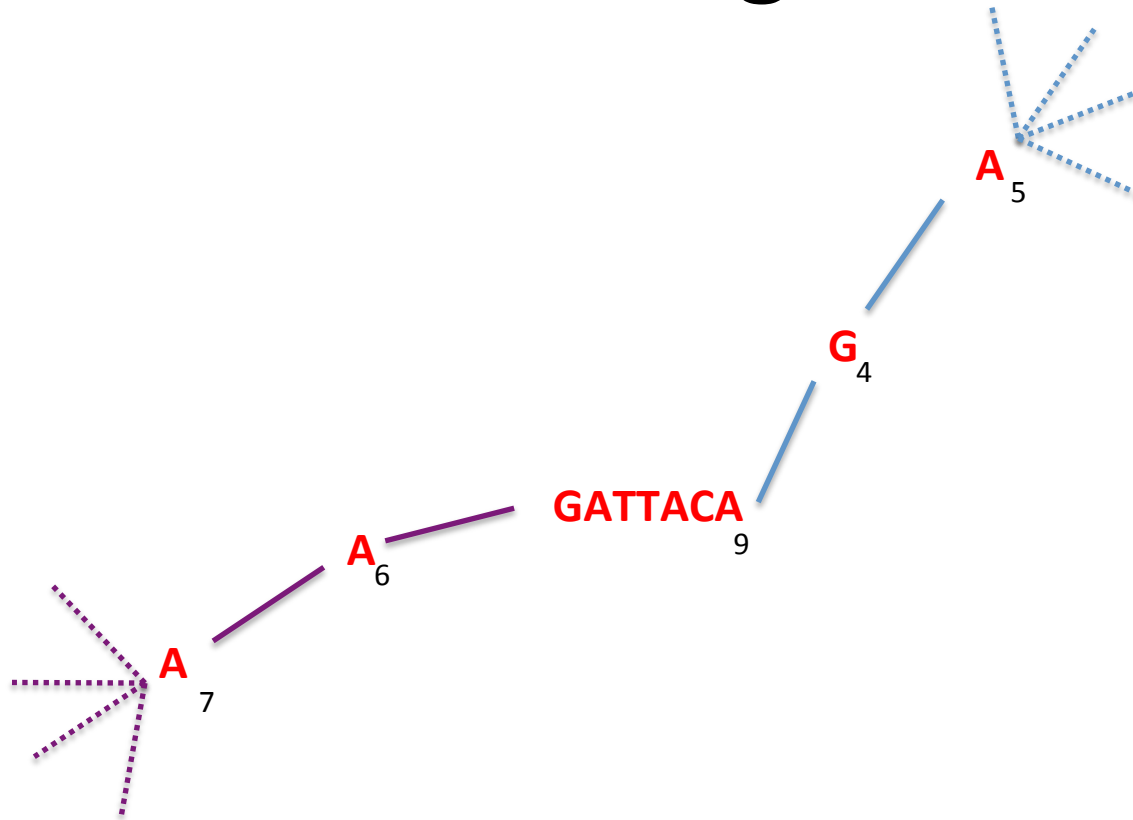
# Inchworm Algorithm

# Inchworm Algorithm

# Inchworm Algorithm

GATTACA$_9$

G$_4$

A$_1$

T$_0$

C$_4$

# Inchworm Algorithm

$G_4$

$A_1$

**GATTACA**$_9$

$T_0$

$C_4$

# Inchworm Algorithm

# Inchworm Algorithm



$G_0$

$A_5$

$T_1$

$G_4$

$C_0$

$A_1$

GATTACA$_9$

$T_0$

$G_1$

$C_4$

$A_1$

$C_1$

$T_1$

# Inchworm Algorithm

**A**$_5$

**G**$_4$

**GATTACA**$_9$

# Inchworm Algorithm

# Inchworm Algorithm



$A_5$

$G_4$

GATTACA$_9$

$A_6$

$A_7$

Report contig:     ….AAGATTACAGA….

Remove assembled kmers from catalog, then repeat the entire process.

Expressed isoforms

Isoform A

Isoform B

# Inchworm Contigs from Alt-Spliced Transcripts

Expressed isoforms

Expression

Isoform A

(low)

Isoform B

(high)

Graphical
representation

**bio**informatics.ca

# Inchworm Contigs from Alt-Spliced Transcripts

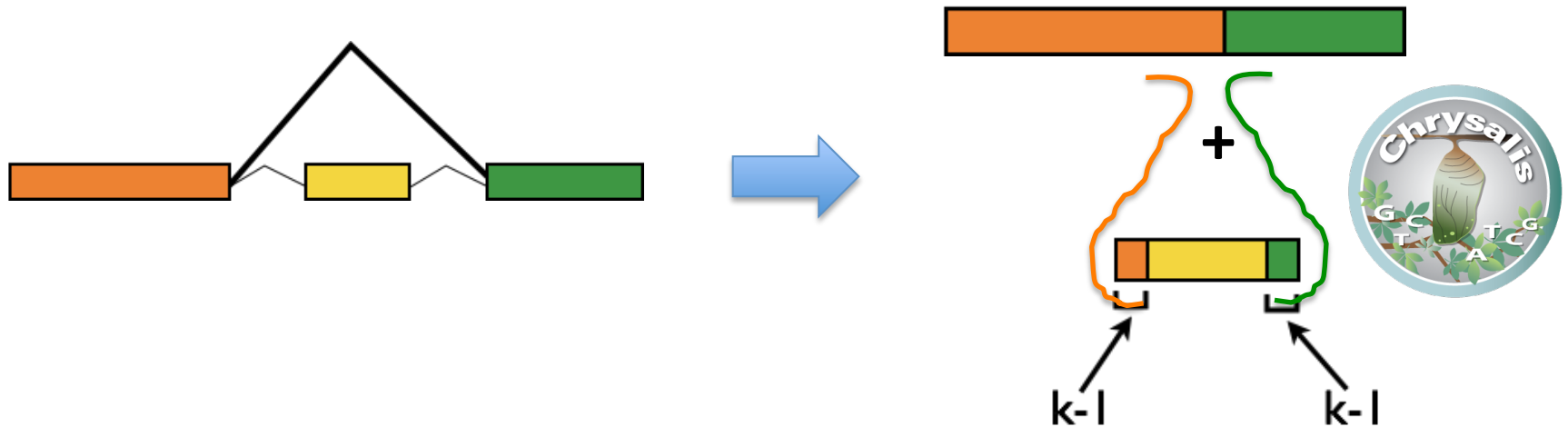# Inchworm Contigs from Alt-Spliced Transcripts



No k-mers in common

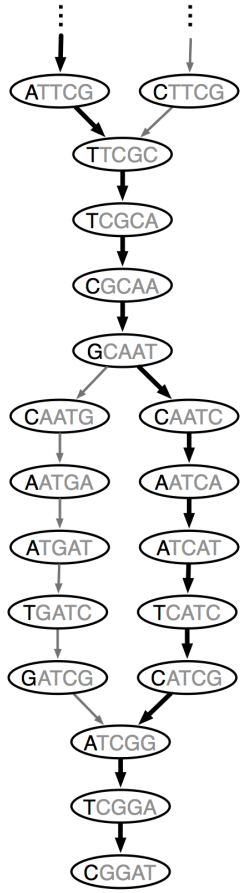# Inchworm Contigs from Alt-Spliced Transcripts

# Chrysalis Re-groups Related Inchworm Contigs



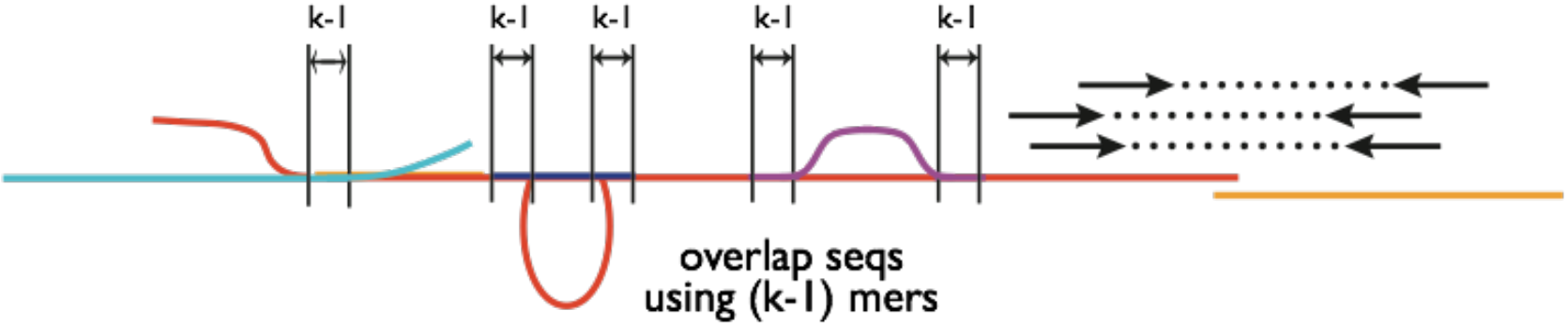Chrysalis uses (k-1) overlaps and read support to link related Inchworm contigs

# 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)
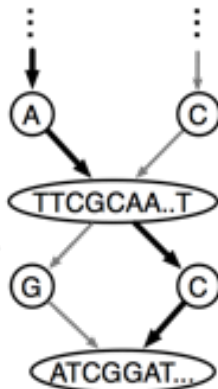
overlap seqs
using (k-I) mers

Thousands of Chrysalis Clusters

# Butterfly

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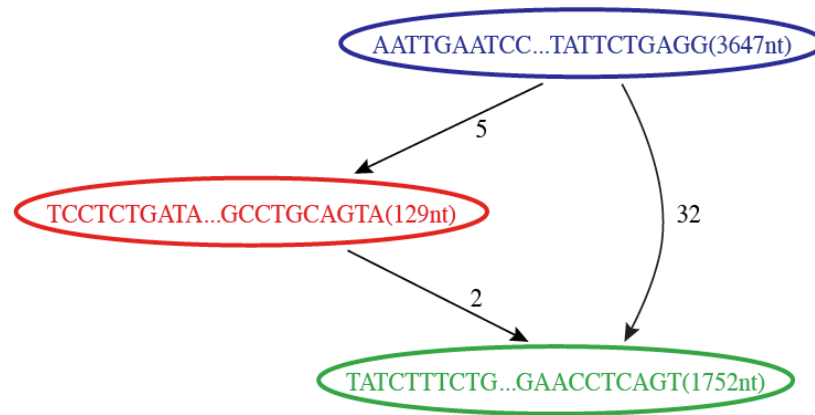
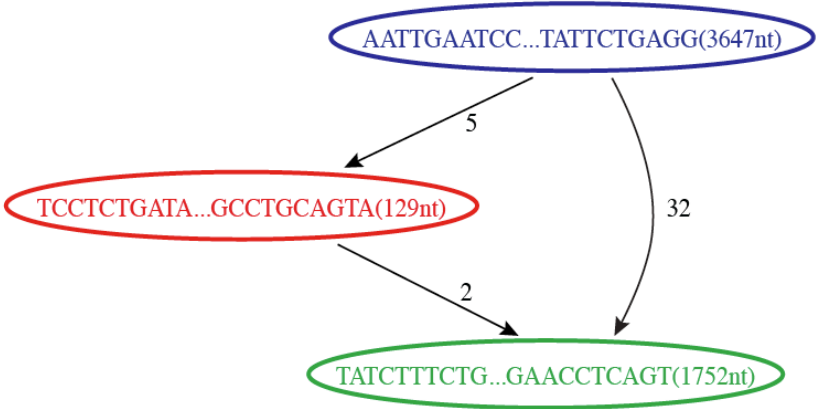..CTTCGCAA..TGATCGGAT...

..ATTCGCAA..TCATCGGAT...

# Butterfly Example 1:
## Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph

# Reconstruction of Alternatively Spliced Transcripts
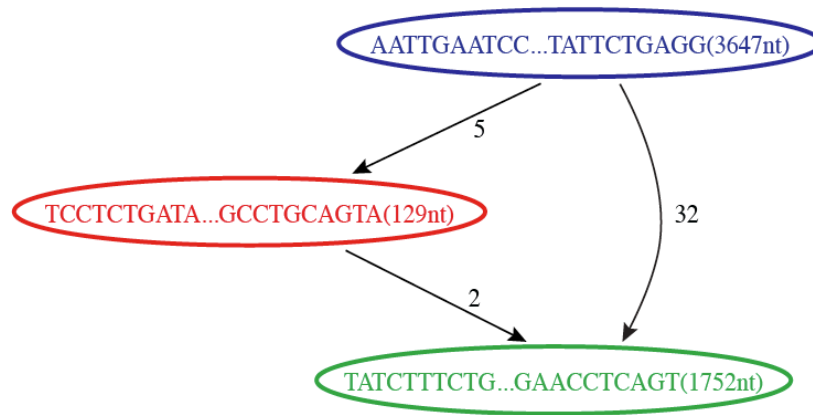
Butterfly's Compacted
Sequence Graph



Reconstructed Transcripts

# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph



Reconstructed Transcripts

**bio**informatics.ca
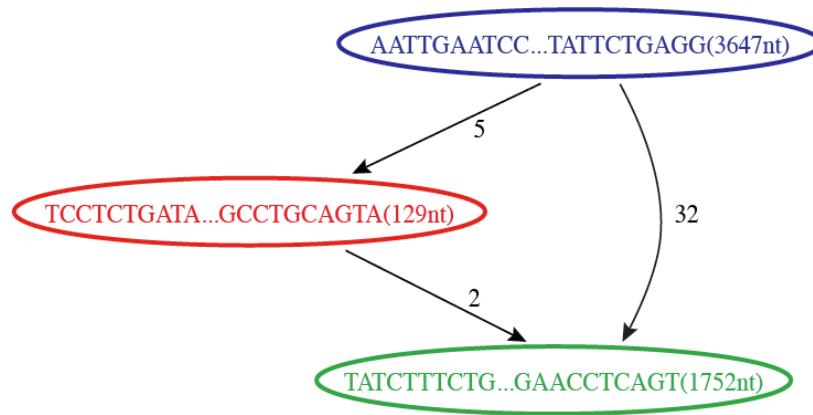
# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted
Sequence Graph



Reconstructed Transcripts



Aligned to Mouse Genome



Naa25 Nalpha acteyltransferase 25 (Reference structure)
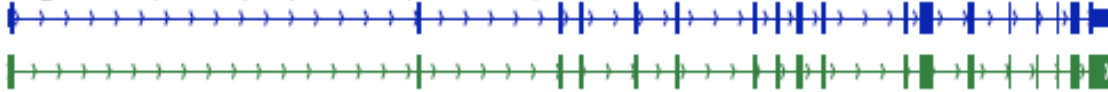
**bio**informatics.ca

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

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



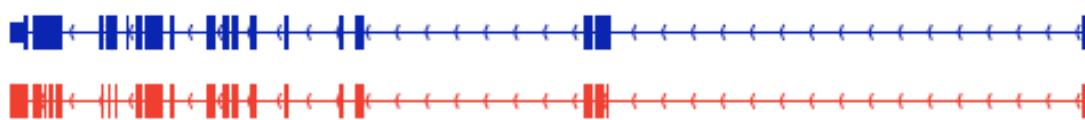chr7:148,744,197–148,821,437

NM_007459; Ap2a2 adaptor protein complex AP-2, alpha 2 subunit

chr7:52,150,889–52,189,508

NM_001077264; Ap2a1 adaptor protein complex AP-2, alpha 1 subunit

# 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

**BROAD**
I N S T I T U T E

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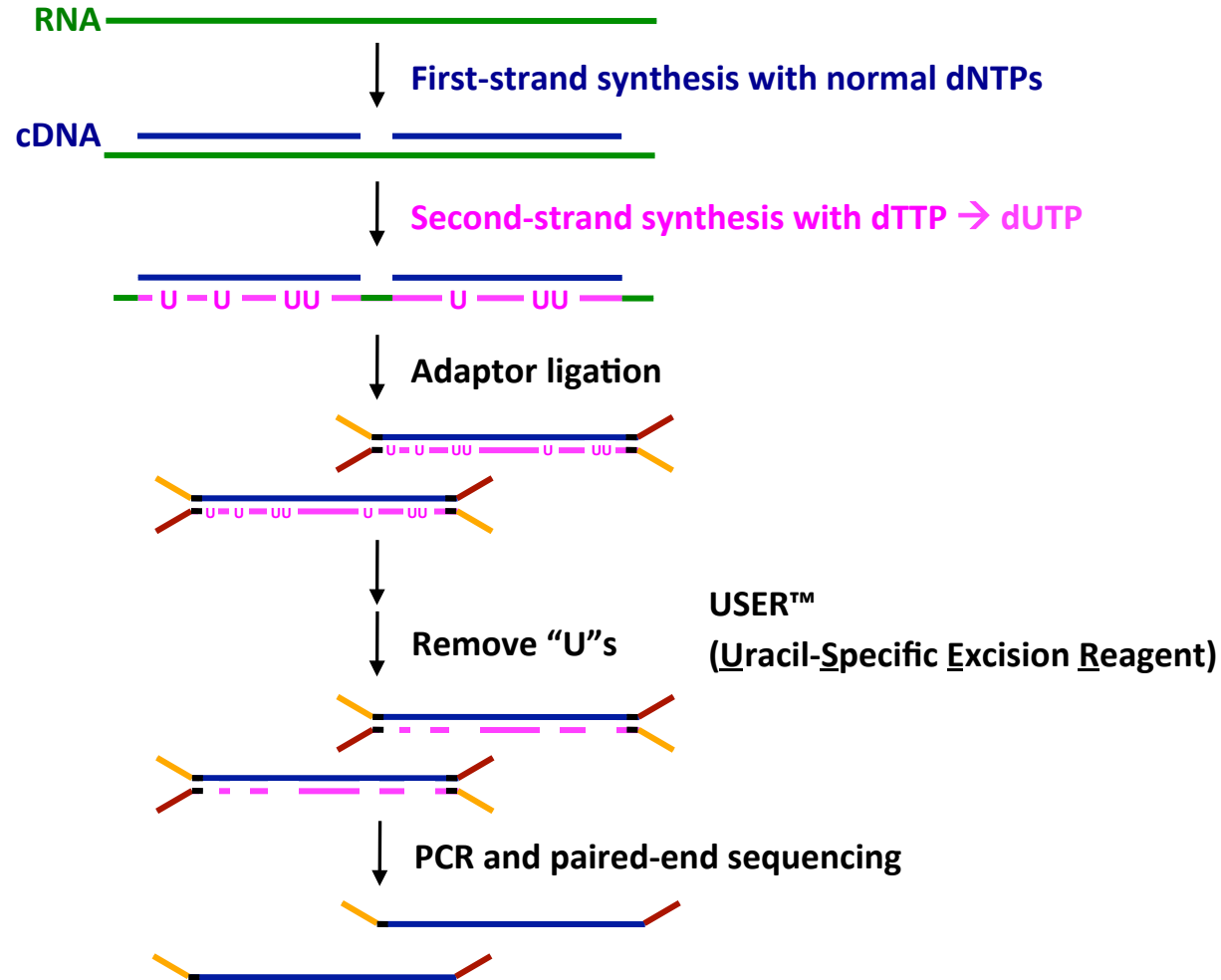
Nevertheless, direct information on the originating strand can substantially enhance the value of an RNA-seq experiment. For example, such information would help to accurately identify anti-

**'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 overlapping transcripts. These tasks are particularly challenging in small microbial genomes, prokaryotic and eukaryotic, in which
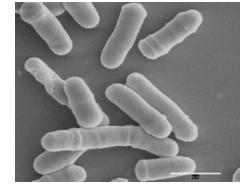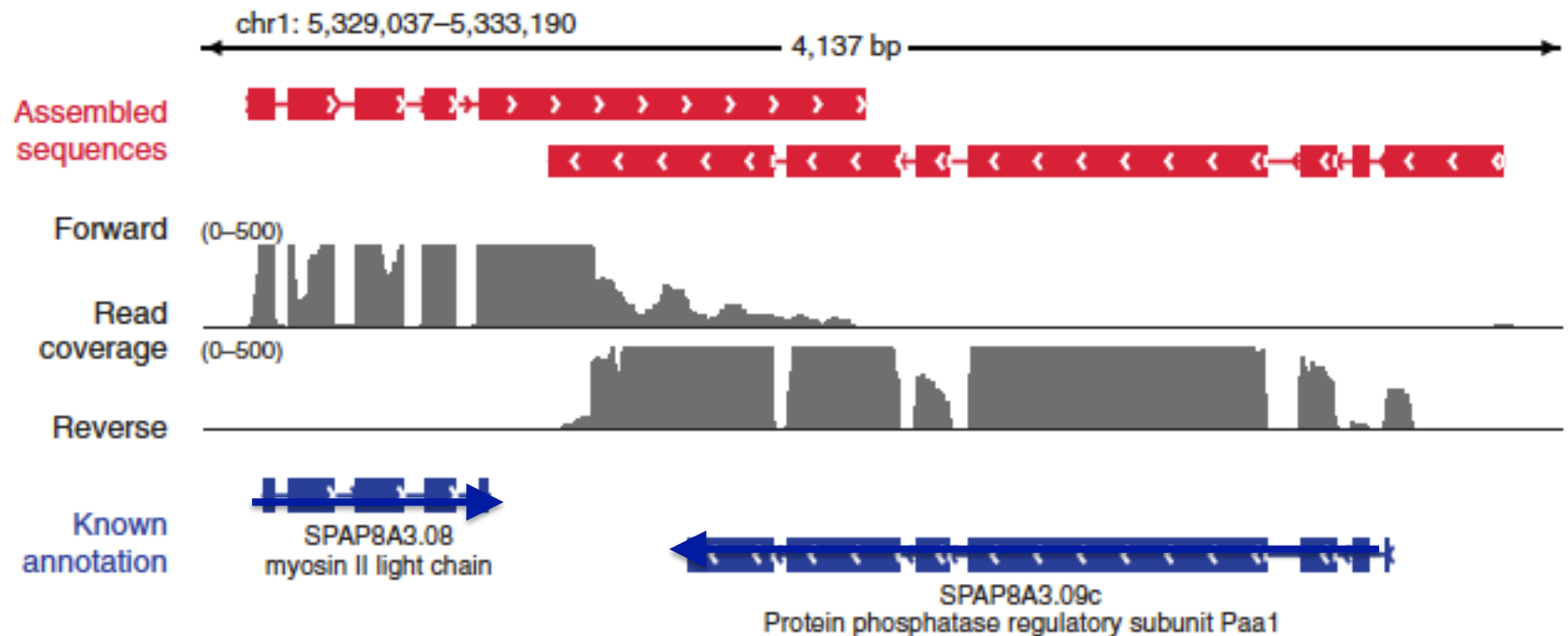
# dUTP 2^nd Strand Method:  Our Favorite



**First-strand synthesis with normal dNTPs**

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

**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

**bio**informatics.ca

# Overlapping UTRs from Opposite Strands

*Schizosacharomyces pombe*
(fission yeast)

**bio**informatics.ca

# Antisense-dominated Transcription

**bio**informatics.ca

# Trinity output: a multi-fasta file

**bio**informatics.ca

# We are on a Coffee Break & Networking Session