

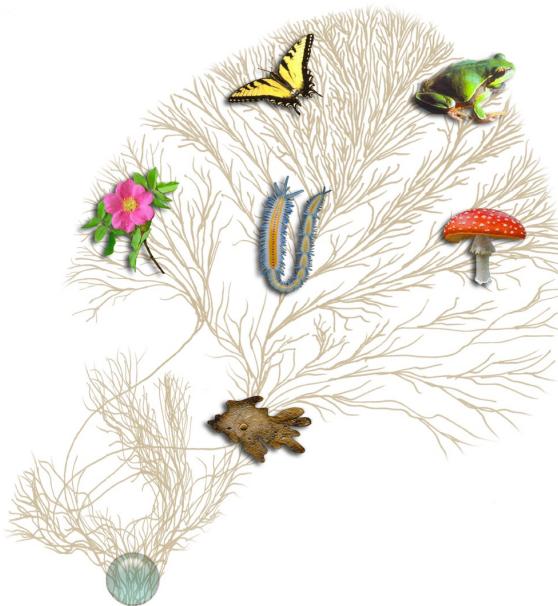


Intro to Transcriptomics and Transcriptome Assembly

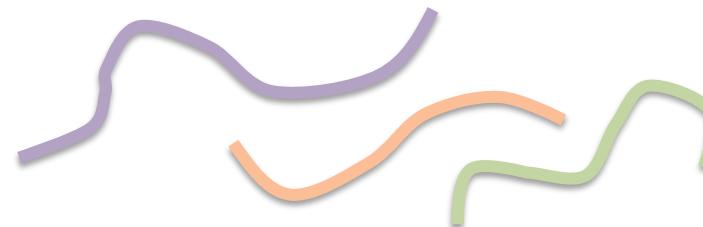
Programming for Biologists
CSHL 2023

Brian Haas
Broad Institute

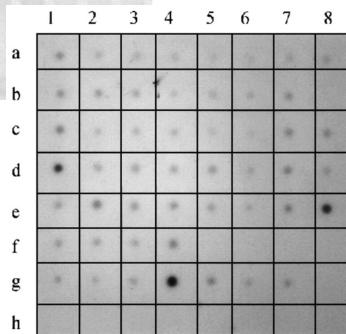
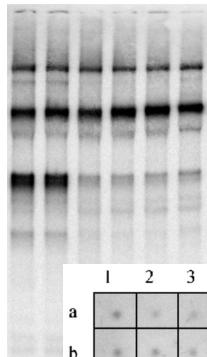
Biological Investigations Empowered by Transcriptomics



Extract RNA,
... some protocol for processing, ...

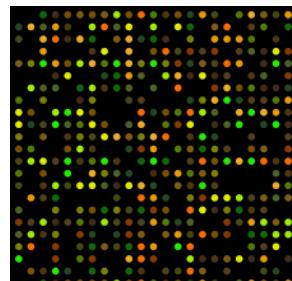


Northern

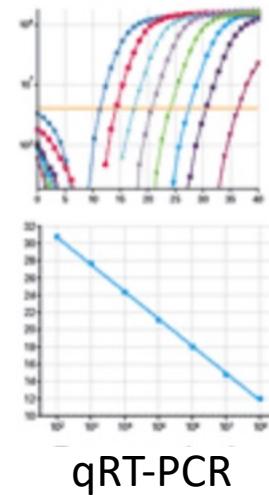


Dot Blot

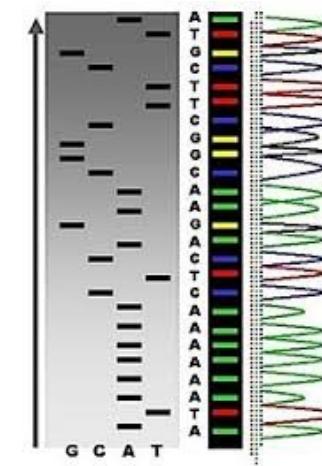
Analysis Method
(pick your favorite)



Microarray



qRT-PCR



Sanger Sequencing



MinION MkI: portable, real time biological analyses



MinION



Other...

Historical Timeline to Modern Transcriptomics (from 1970)

Reverse Transcription (1970)

Northern Blot
Sanger Sequencing
(1977)

Expressed Sequence Tags (1992)

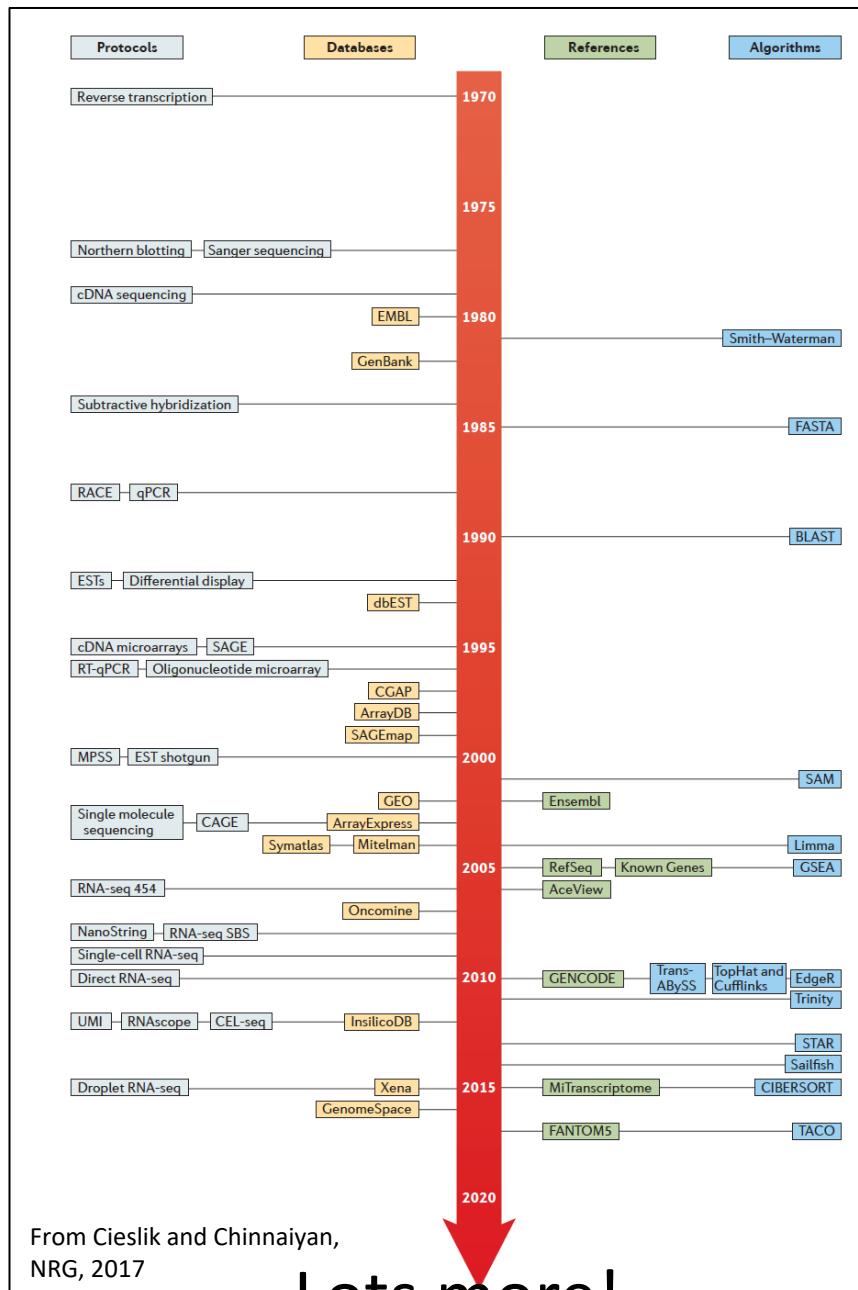
cDNA microarrays (1995)

RNA-Seq (2006-2008)

PacBio IsoSeq (2014)

Droplet single cell RNA-Seq (2015)

Direct RNA Seq Nanopore (2018)



Note: Just a small sampling of what's available.

Smith Waterman (1981)

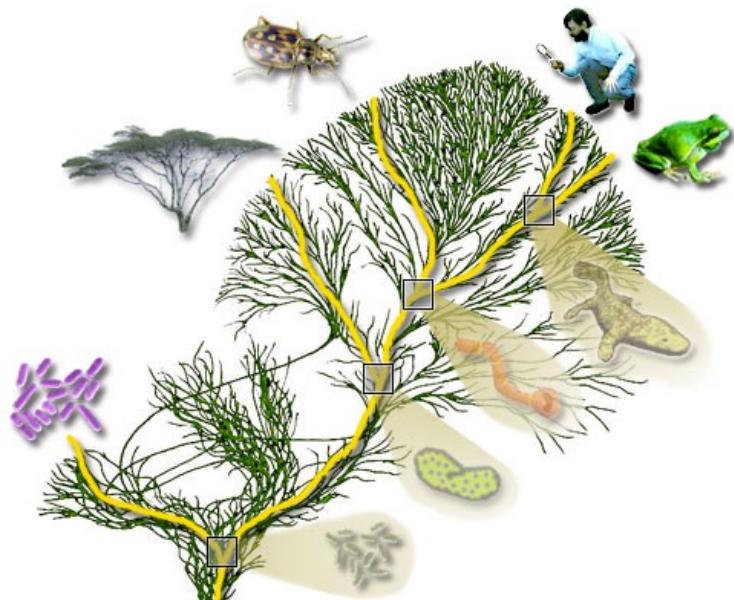
BLAST (1990)

Tophat/Cufflinks (2010)

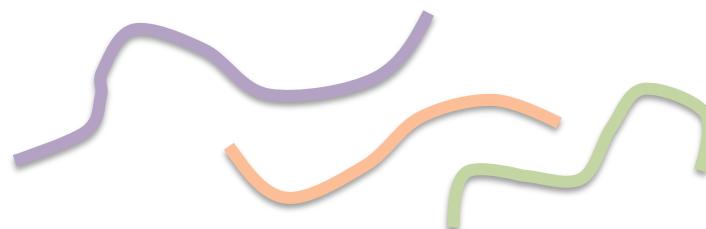


Kallisto (2016)
Salmon (2017)

Modern Transcriptome Studies Empowered by RNA-seq



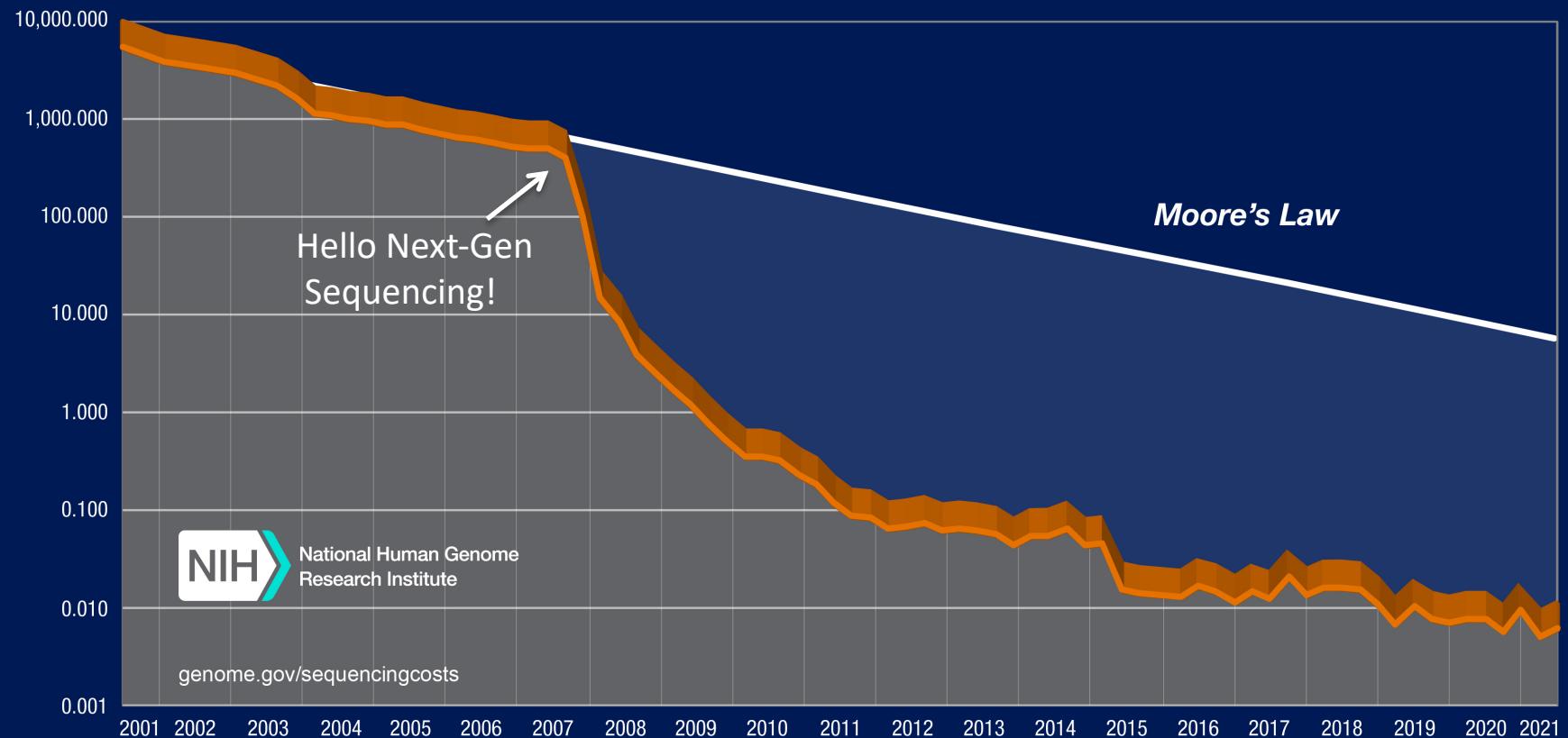
Extract RNA, convert to cDNA



Next-gen Sequencer
(pick your favorite)

Millions to Billions of Reads

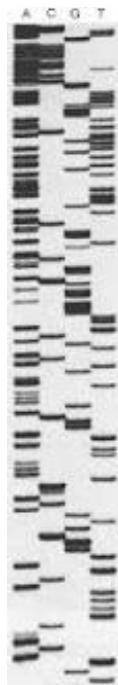
Cost per Raw Megabase of DNA Sequence



From <https://www.genome.gov/sequencingcostsdata/>

Personal Reflections...

Circa 1995



Generating RNA-Seq: How to Choose?

| Platform | iSeq Project Firefly 2018 | MinSeq | MiSeq | Next Seq 550 | HiSeq 2500 RR | Hiseq 2500 V3 | HiSeq 2500 V4 | HiSeq 4000 | HiSeq X | Nova Seq S1 2018 | Nova Seq S2 | Nova Seq S4 | 5500 XL | 318 HiQ 520 | Ion 530 | Ion Proton P1 | PGM HiQ 540 | RS P6-C4 | Sequel | R&D end 2018 | Smidg ION RnD | Mini ION R9.5 | Grid ION X5 | PromethION RnD | PromethION theoretical | QiaGen Gene Reader | BGI SEQ 500 | BGI SEQ 50 | # |
|-----------------------------------|---------------------------|--------|-------|--------------|---------------|---------------|---------------|------------|---------|------------------|-------------|-------------|---------|-------------|------------|---------------|-------------|----------|--------|--------------|---------------|---------------|-------------|----------------|------------------------|--------------------|-------------|------------|----|
| Reads: (M) | 4 | 25 | 25 | 400 | 600 | 3000 | 4000 | 5000 | 6000 | 3300 | 6600 | 20000 | 1400 | 3-5 | 15-20 | 165 | 60-80 | 5.5 | 38.5 | -- | -- | -- | -- | -- | -- | 400 | 1600 | 1600 | -- |
| Read length: (paired-end*) | 150* | 150* | 300* | 150* | 100* | 100* | 125* | 150* | 150* | 150* | 150* | 150* | 60 | 200 400 | 200 400 | 200 | 200 | 15K | 12K | 32K | -- | -- | -- | -- | -- | -- | 100* | 50 | -- |
| Run time: (d) | 0.54 | 1 | 2 | 1.2 | 1.125 | 11 | 6 | 3.5 | 3 | 1.66 | 1.66 | 1.66 | 7 | 0.37 | 0.16 | -- | 0.16 | 4.3 | -- | -- | -- | 2 | 2 | 2 | -- | -- | 1 | 0.4 | -- |
| Yield: (Gb) | 1 | 7.5 | 15 | 120 | 120 | 600 | 1000 | 1500 | 1800 | 1000 | 2000 | 6000 | 180 | 1.5 | 7 | 10 | 12 | 12 | 5 | 150 | 4 | 8 | 40 | 2400 | 11000 | 80 | 200 | 8 | -- |
| Rate: (Gb/d) | 1.85 | 7.5 | 7.5 | 100 | 106.6 | 55 | 166 | 400 | 600 | 600 | 1200 | 3600 | 30 | 5.5 | 50 | -- | 93.75 | 2.8 | -- | -- | -- | 4 | 20 | 1200 | 5500 | -- | 200 | 20 | -- |
| Reagents: (\$K) | 0.1 | 1.75 | 1 | 5 | 6.145 | 23.47 | 29.9 | -- | -- | -- | -- | -- | 10.5 | 0.6 | -- | 1 | 1.2 | 2.4 | -- | 1 | -- | 0.5 | 1.5 | -- | -- | 0.5 | -- | -- | -- |
| per-Gb: (\$) | 100 | 233 | 66 | 50 | 51.2 | 39.1 | 31.7 | 20.5 | 7.08 | 18 | 15 | 5.8 | 58.33 | -- | -- | 100 | -- | 200 | 80 | 6.6 | -- | 62.5 | 37.5 | 20 | 4.3 | -- | -- | -- | -- |
| hg-30x: (\$) | 12000 | 28000 | 8000 | 5000 | 6144 | 4692 | 3804 | 2460 | 849.6 | 1800 | 1564 | 700 | 7000 | -- | -- | 12000 | -- | 24000 | 9600 | 1000 | -- | 7500 | 4500 | 2400 | 500 | -- | 600 | -- | |
| Machine: (\$) | 30K | 49.5K | 99K | 250K | 740K | 690K | 690K | 900K | 1M | 999K | 999K | 999K | 595K | 50K | 65K | 243K | 242K | 695K | 350K | 350K | -- | -- | 125K | 75K | 75K | -- | 200K | -- | |

#Page maintained by http://twitter.com/albertvilella http://tinyurl.com/ngslytics #Editable version: http://tinyurl.com/ngsspecsshared

#curl "https://docs.google.com/spreadsheets/d/1GMMfhLyLK0-q8Xklo3YxlWaZA5vVMuhU1kg41g4xLkXc/export?gid=4&format=csv" | grep -v '^#' | grep -v '\"' | column -t -s\| less -S

Stats circa 2018

For current, see: <https://tinyurl.com/wbgcs65>



*Not all shown at scale

Generating RNA-Seq: How to Choose?

| Platform | Project Firefly 2018 | MiniSeq | MiSeq | Next Seq 550 | HiSeq 2500 RR | Hiseq 2500 V |
|----------------------------|----------------------|---------|-------|--------------|---------------|--------------|
| Reads: (M) | 4 | 25 | 25 | 400 | 600 | 300 |
| Read length: (paired-end*) | 150* | 150* | 300* | 150* | 100* | 100 |
| Run time: (d) | 0.54 | 1 | 2 | 1.2 | 1.125 | 1 |
| Yield: (Gb) | 1 | 7.5 | 15 | 120 | 120 | 60 |
| Rate: (Gb/d) | 1.85 | 7.5 | 7.5 | 100 | 106.6 | 5 |
| Reagents: (\$K) | 0.1 | 1.75 | 1 | 5 | 6.145 | 23.4 |
| per-Gb: (\$) | 100 | 233 | 66 | 50 | 51.2 | 39. |
| hg-30x: (\$) | 12000 | 28000 | 8000 | 5000 | 6144 | 469 |
| Machine: (\$) | 30K | 49.5K | 99K | 250K | 740K | 690K |

#Page maintained by <http://twitter.com/albertvilella> http://
<https://docs.google.com/spreadsheets/d/1GMMfhylKQ-q8>



| Plat | Mini ION R9.5 | Grid ION X5 | Prome thION RnD | Prome thION theor etical | QiaGen Gene Reader | BGI SEQ 500 | BGI SEQ 50 | # |
|------|---------------|-------------|-----------------|--------------------------|--------------------|-------------|------------|----|
| -- | -- | -- | -- | -- | 400 | 1600 | 1600 | -- |
| -- | -- | -- | -- | -- | 100* | 50 | -- | -- |
| -- | 2 | 2 | 2 | -- | -- | 1 | 0.4 | -- |
| 4 | 8 | 40 | 2400 | 11000 | 80 | 200 | 8 | -- |
| -- | 4 | 20 | 1200 | 5500 | -- | 200 | 20 | -- |
| -- | 0.5 | 1.5 | -- | -- | 0.5 | -- | -- | -- |
| -- | 62.5 | 37.5 | 20 | 4.3 | -- | -- | -- | -- |
| -- | 7500 | 4500 | 2400 | 500 | -- | 600 | -- | -- |
| -- | -- | 125K | 75K | 75K | -- | 200K | -- | -- |



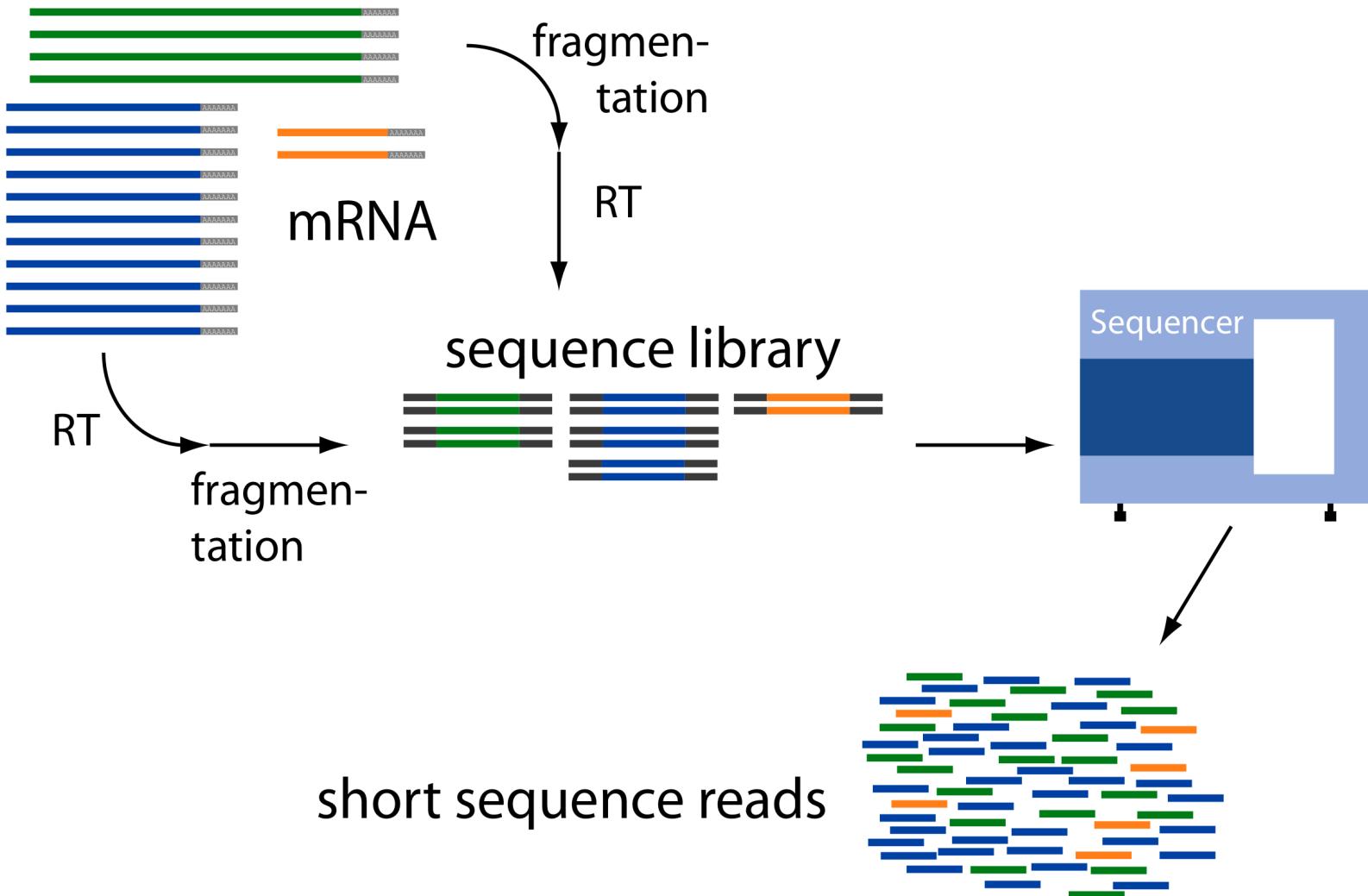
Thx Joshua Levin, for the cartoon. ☺



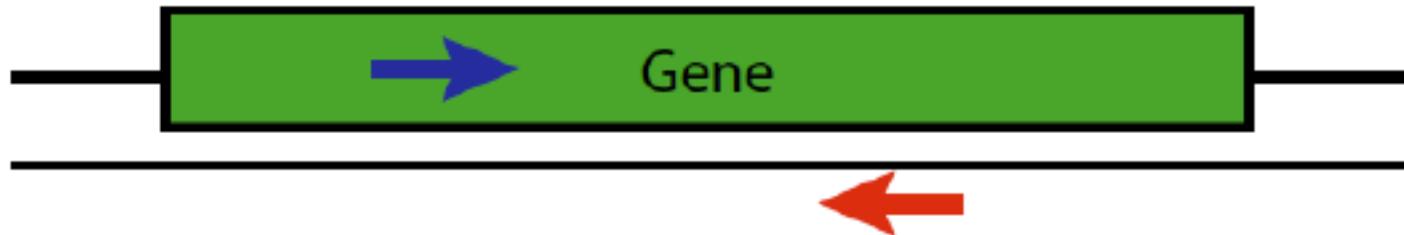
Each has pros/cons



Overview of RNA-Seq



Paired-end Sequences

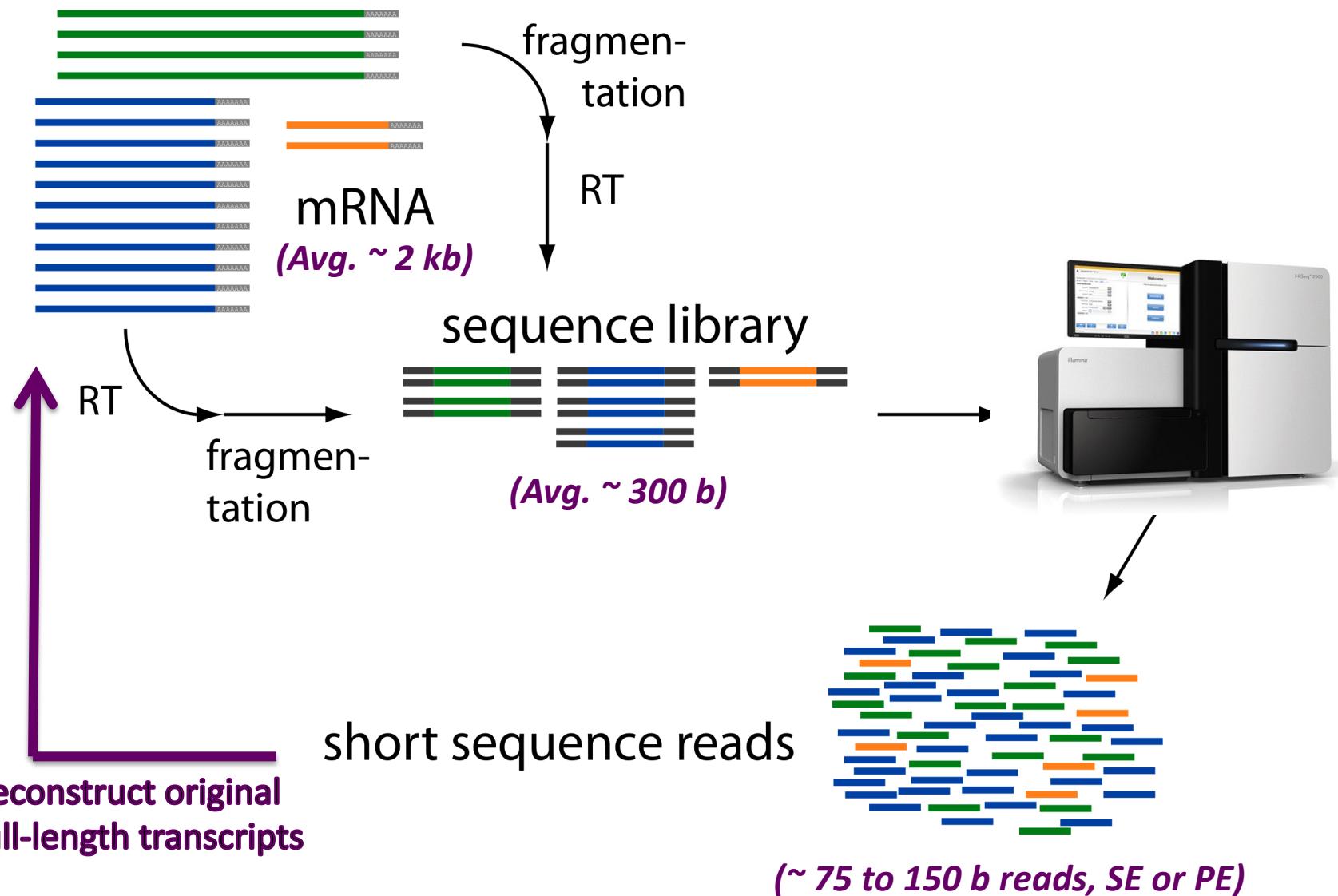


Two FastQ files, read name indicates
left (/1) or right (/2) read of paired-end

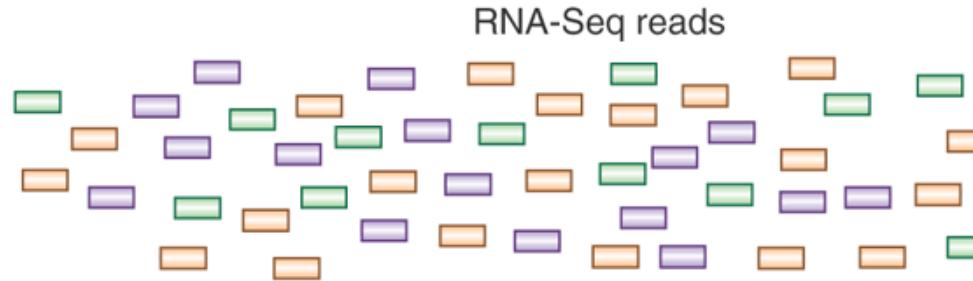
```
@61DFRAAXX100204:1:100:10494:3070/1
AAACAACAGGGCACATTGTCACTCTTGTATTGAAAAACACTTCCGGCCAT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCC@ @CACCCCCA
```

```
@61DFRAAXX100204:1:100:10494:3070/2
CTCAAATGGTTAATTCTCAGGCTGCAAATATTGTTCAAGGATGGAAGAAC
+
C<CCCCCCCCACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBCCCC
```

RNA-Seq Challenge: Transcript Reconstruction



Transcript Reconstruction from RNA-Seq Reads



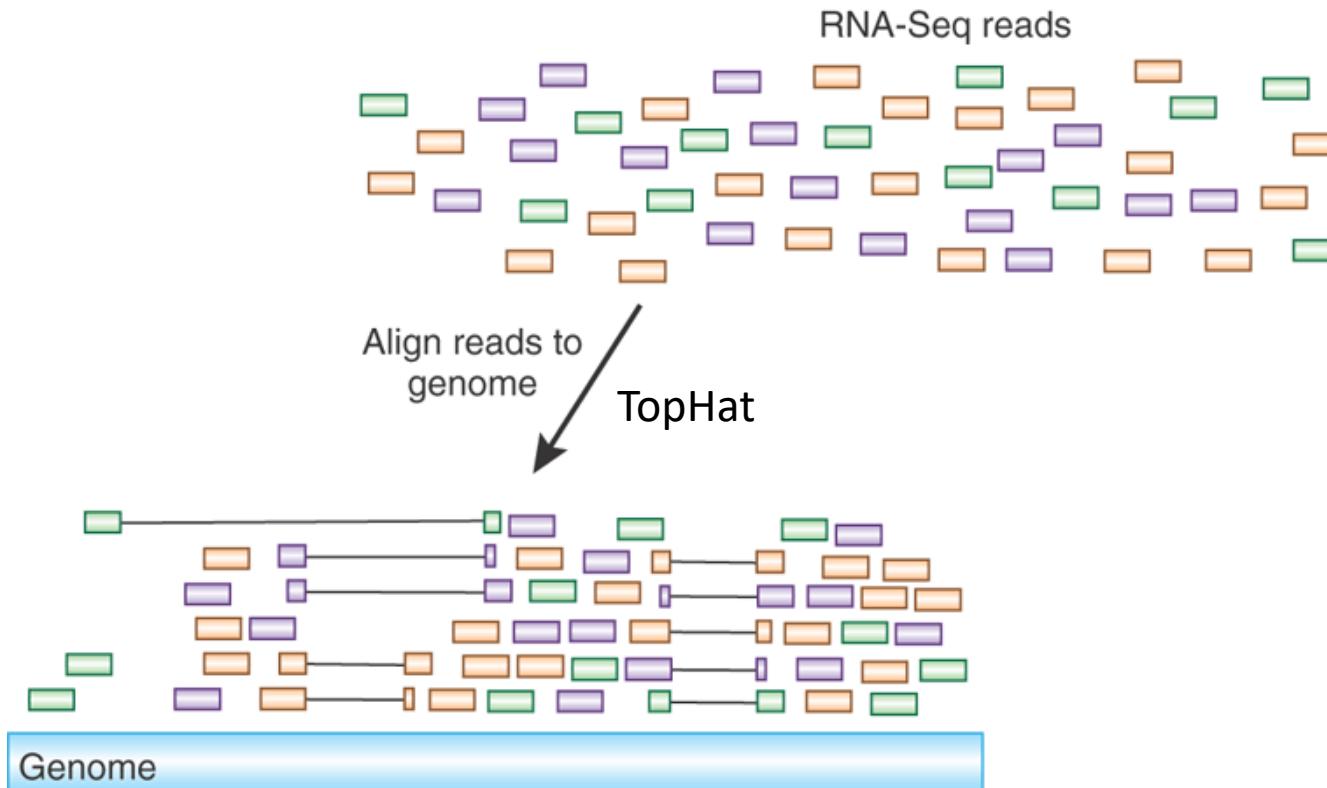
Advancing RNA-Seq analysis

Brian J Haas & Michael C Zody

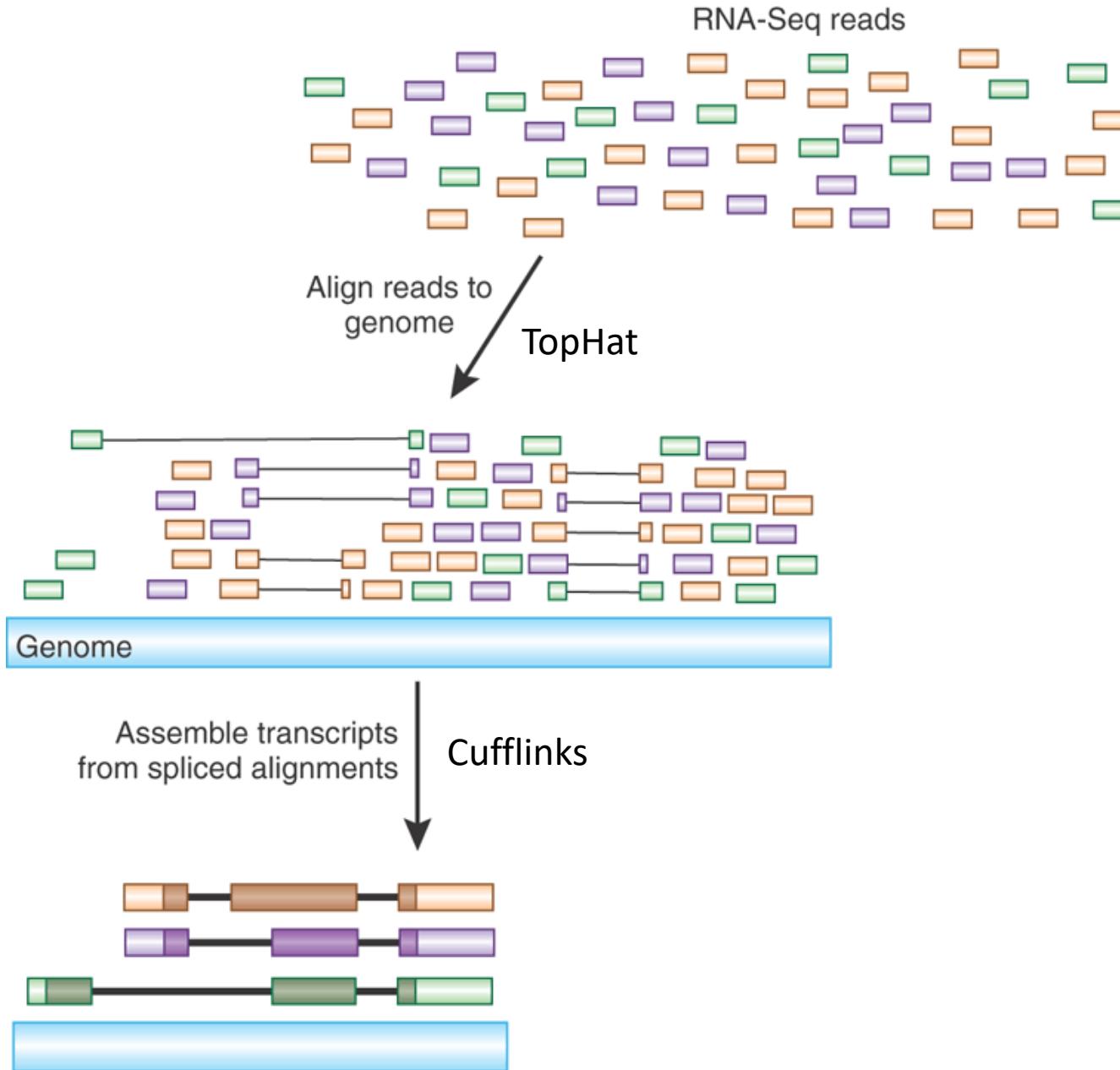
Nature Biotech, 2010

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

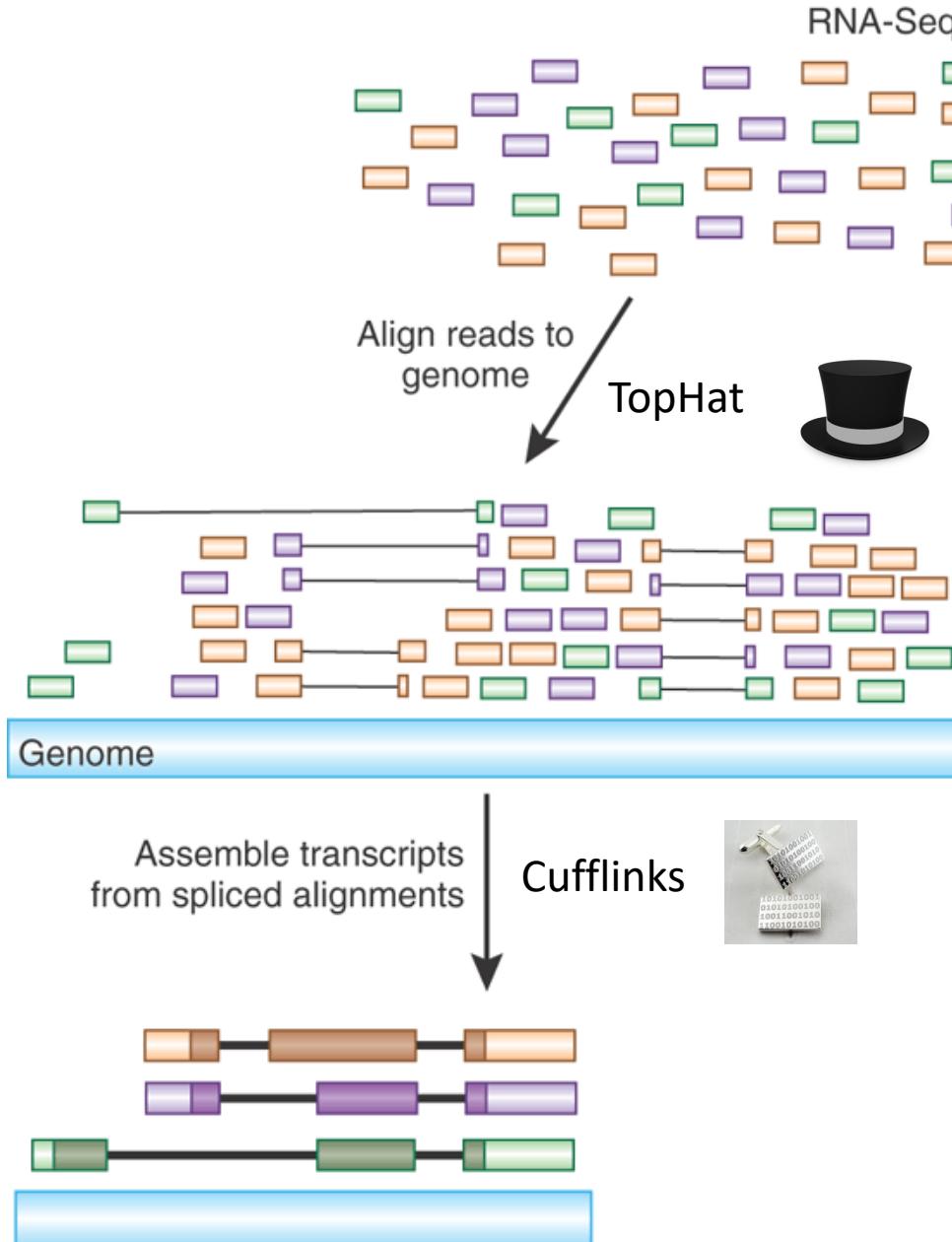
Transcript Reconstruction from RNA-Seq Reads



Transcript Reconstruction from RNA-Seq Reads



Transcript Reconstruction from RNA-Seq Reads



The Tuxedo Suite:
End-to-end **Genome**-based
RNA-Seq Analysis
Software Package

NATURE PROTOCOLS | PROTOCOL

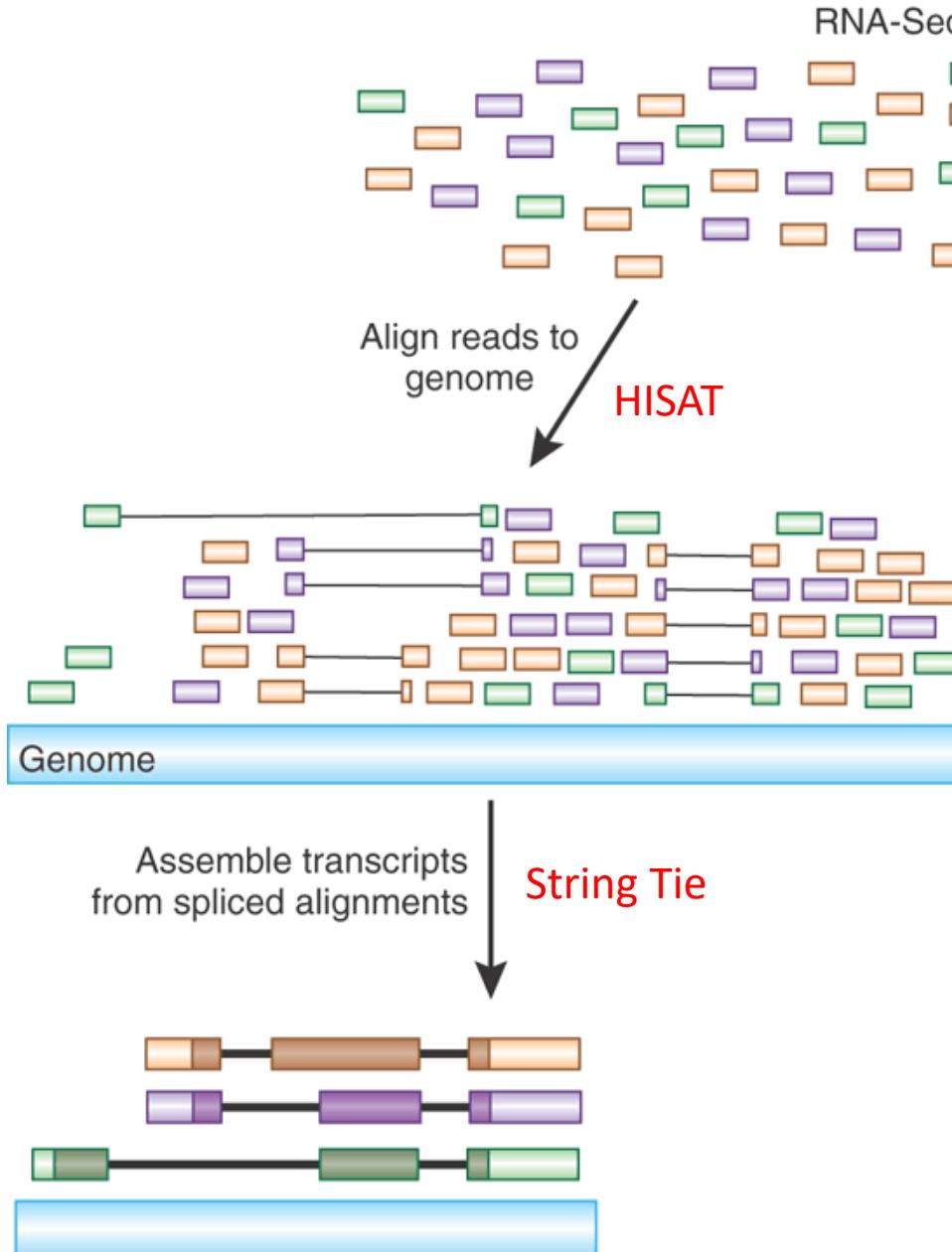
Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks

Cole Trapnell, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R Kelley, Harold Pimentel, Steven L Salzberg, John L Rinn & Lior Pachter

Affiliations | Contributions | Corresponding author

Nature Protocols 7, 562–578 (2012) | doi:10.1038/nprot.2012.016
Published online 01 March 2012

Transcript Reconstruction from RNA-Seq Reads



The “New Tuxedo” Suite:

End-to-end Genome-based
RNA-Seq Analysis
Software Package

NATURE PROTOCOLS | PROTOCOL



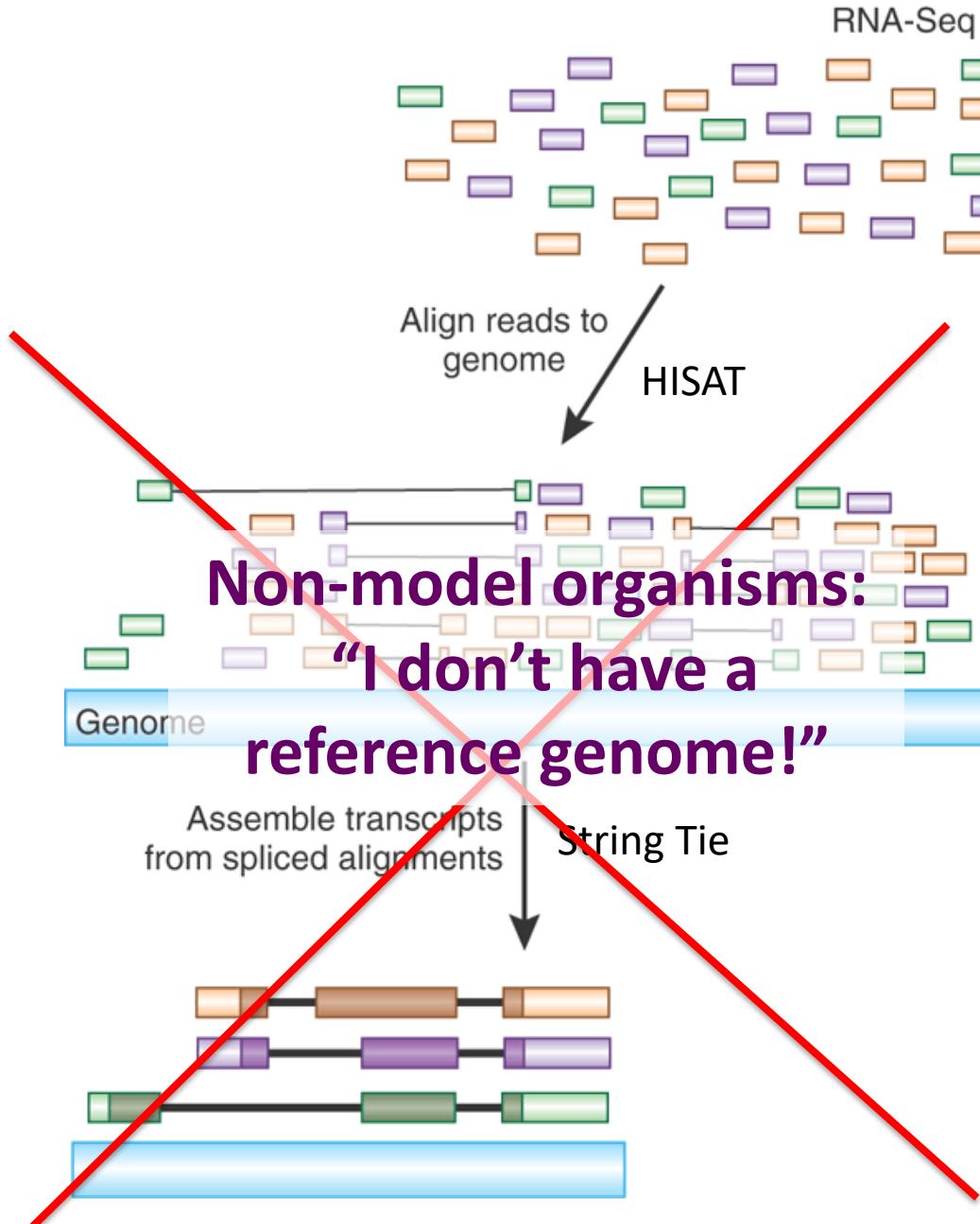
Transcript-level expression analysis of
RNA-seq experiments with HISAT,
StringTie and Ballgown

Mihaela Pertea, Daehwan Kim, Geo M Pertea, Jeffrey T Leek & Steven L Salzberg

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)

Nature Protocols 11, 1650–1667 (2016) | doi:10.1038/nprot.2016.095
Published online 11 August 2016

Transcript Reconstruction from RNA-Seq Reads



The “New Tuxedo” Suite:
End-to-end Genome-based
RNA-Seq Analysis
Software Package

NATURE PROTOCOLS | PROTOCOL

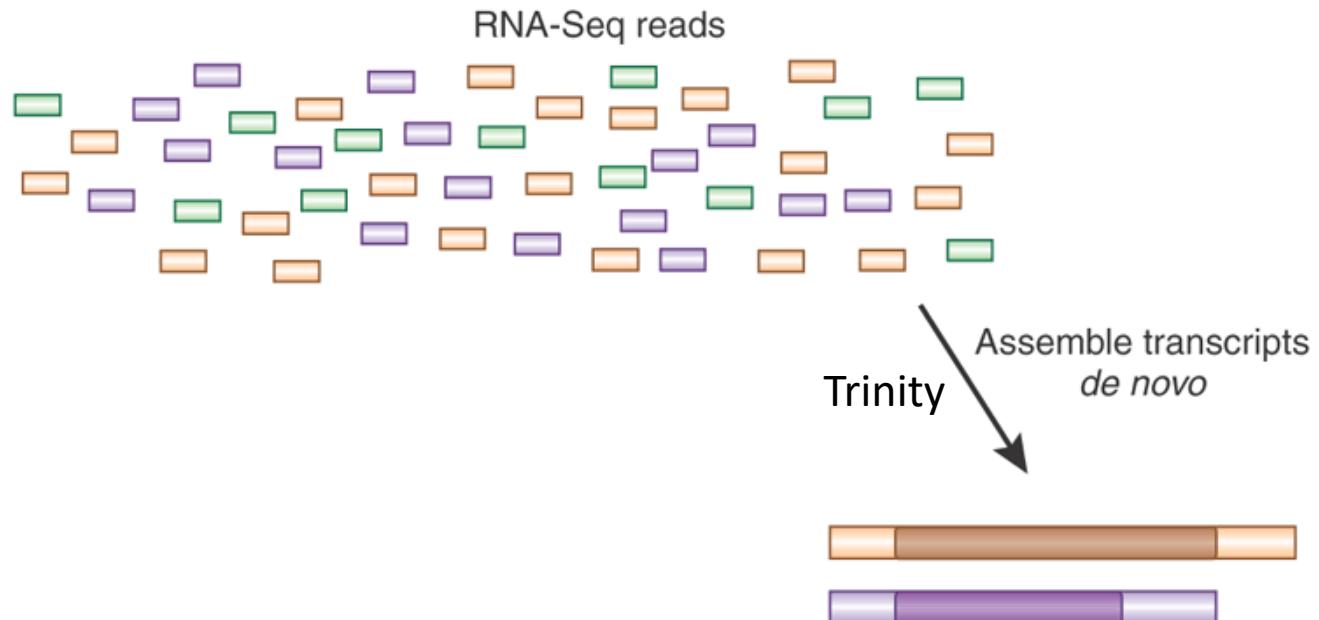
Transcript-level expression analysis of
RNA-seq experiments with HISAT,
StringTie and Ballgown

Mihaela Pertea, Daehwan Kim, Geo M Pertea, Jeffrey T Leek & Steven L Salzberg

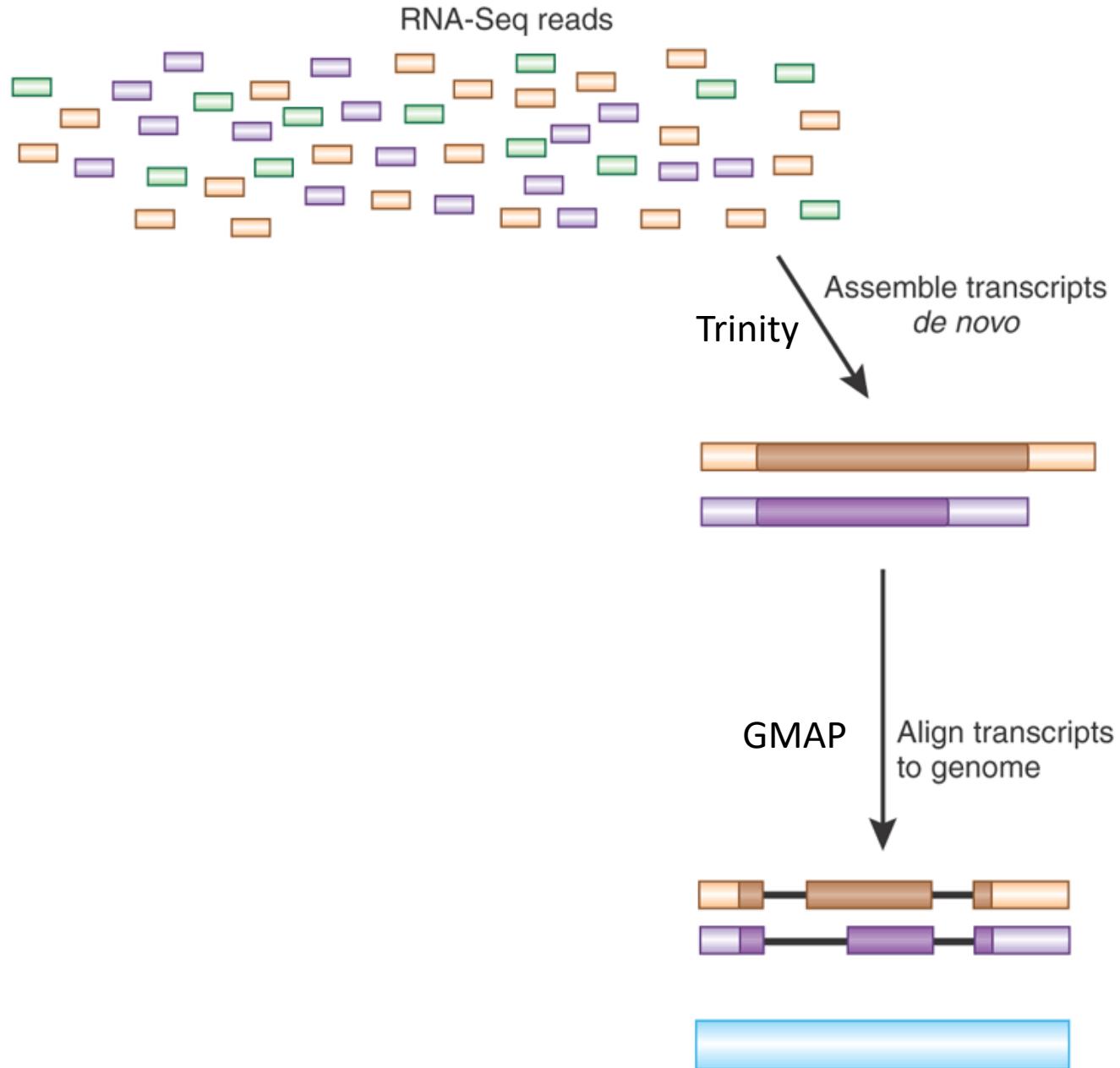
Affiliations | Contributions | Corresponding author

Nature Protocols 11, 1650–1667 (2016) | doi:10.1038/nprot.2016.095
Published online 11 August 2016

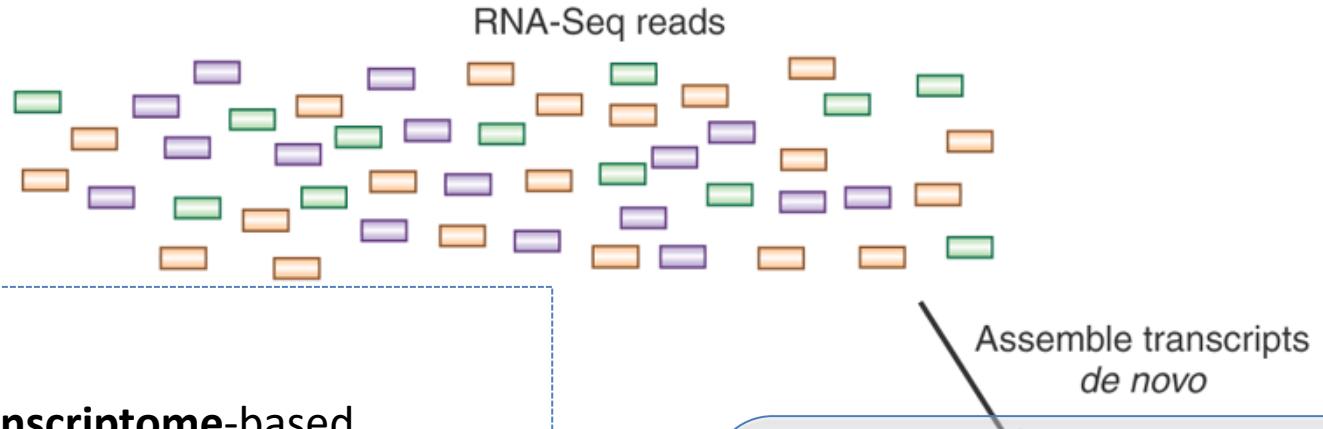
Transcript Reconstruction from RNA-Seq Reads



Transcript Reconstruction from RNA-Seq Reads



Transcript Reconstruction from RNA-Seq Reads



End-to-end Transcriptome-based
RNA-Seq Analysis
Software Package

NATURE PROTOCOLS | PROTOCOL

De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis

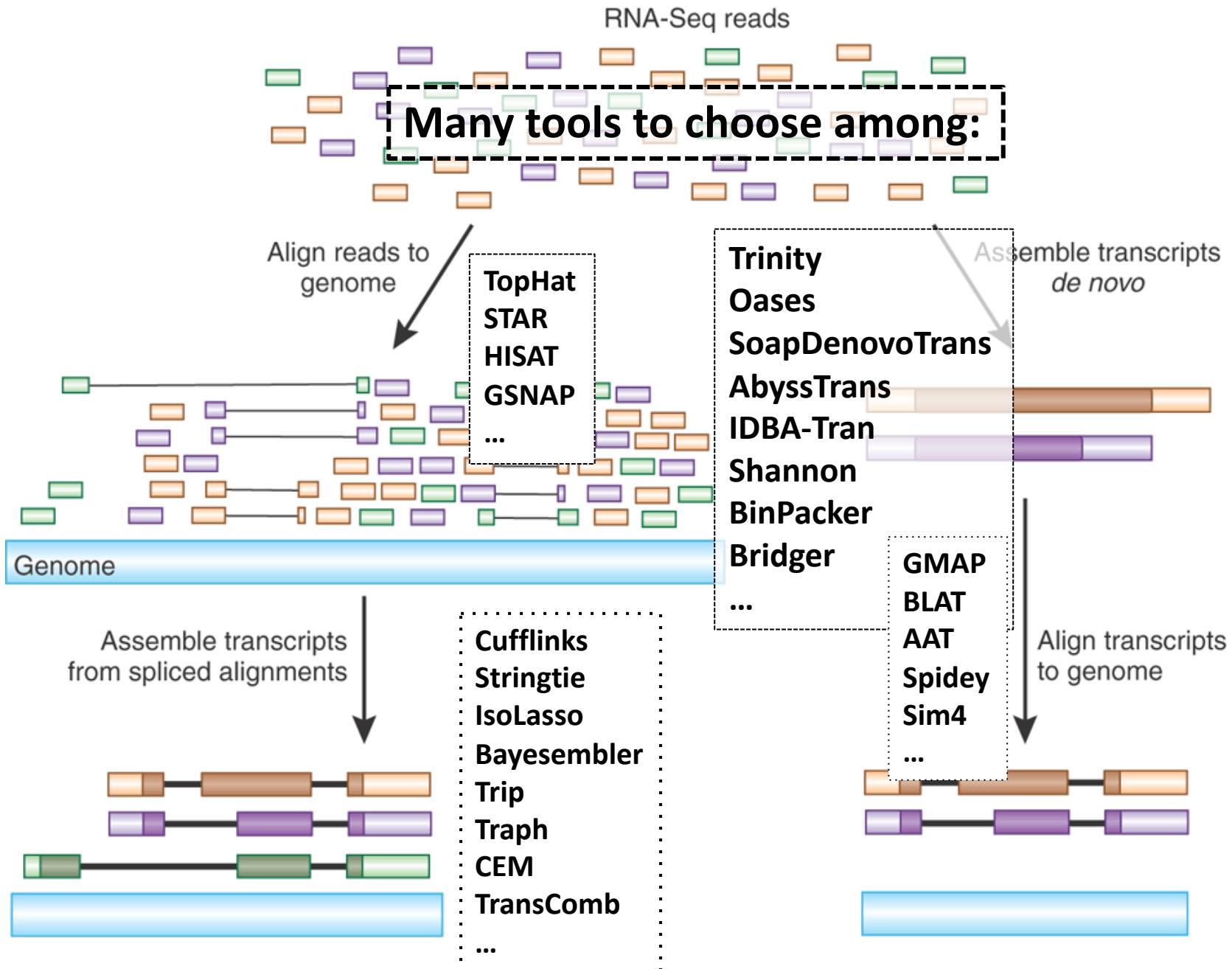
Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes, Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman, Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman & Aviv Regev

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

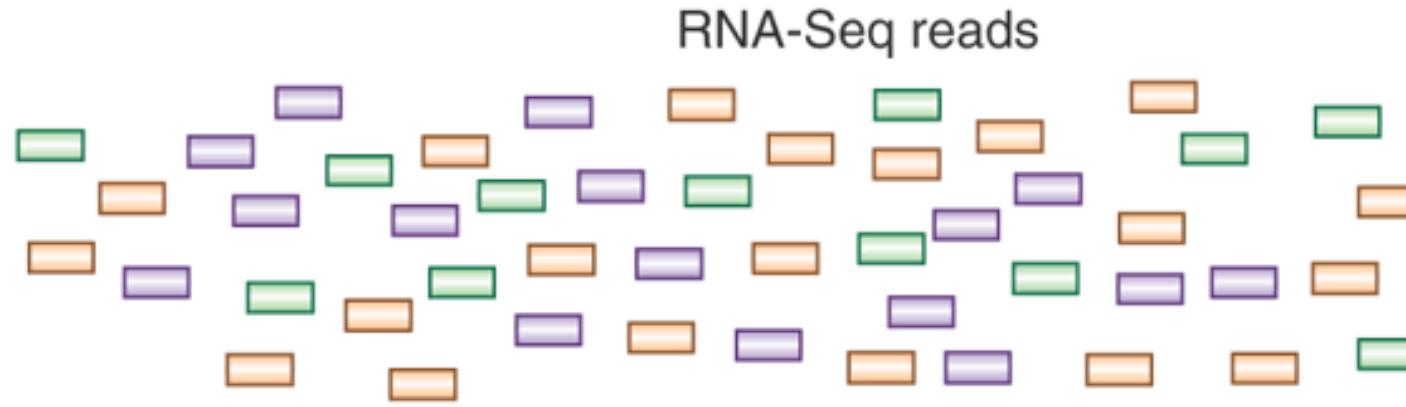
Nature Protocols 8, 1494–1512 (2013) | doi:10.1038/nprot.2013.084

Published online 11 July 2013

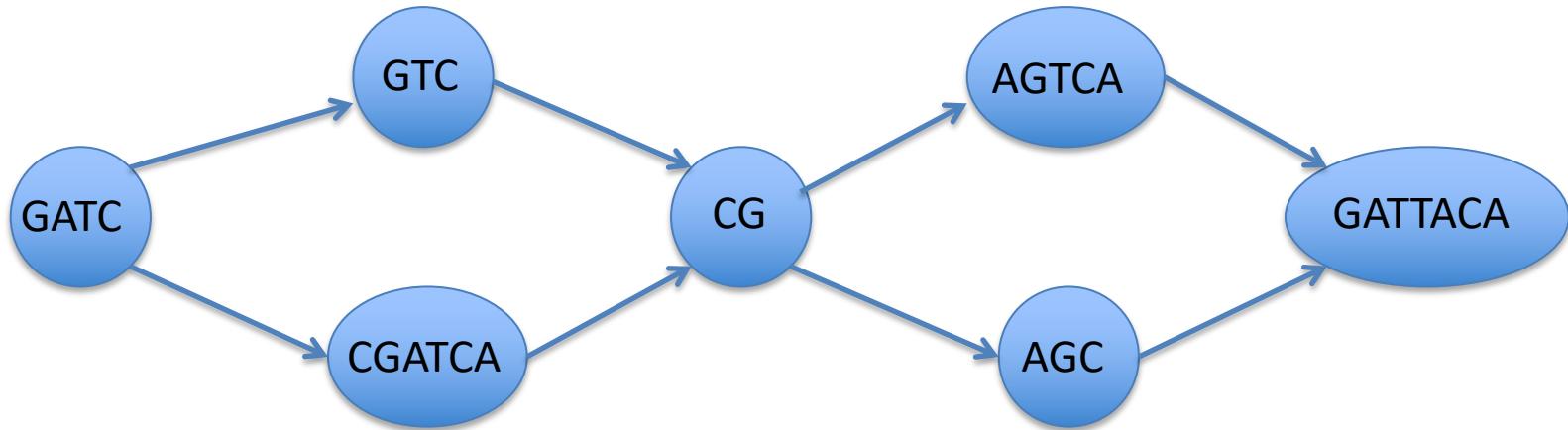
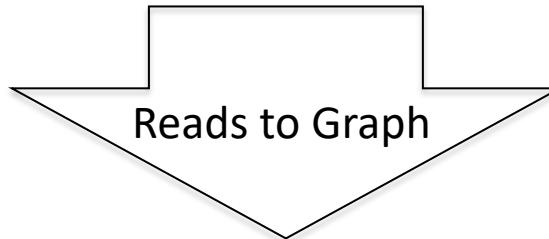
Transcript Reconstruction from RNA-Seq Reads



Graph Data Structures Commonly Used For Assembly

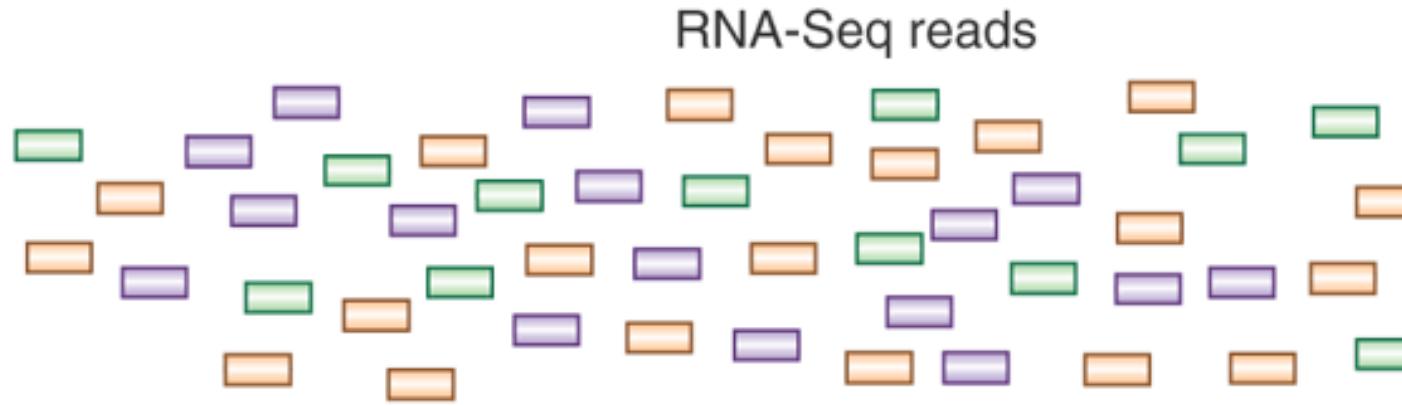


- Sequence
- Order
- Orientation (+, -)
- Overlap

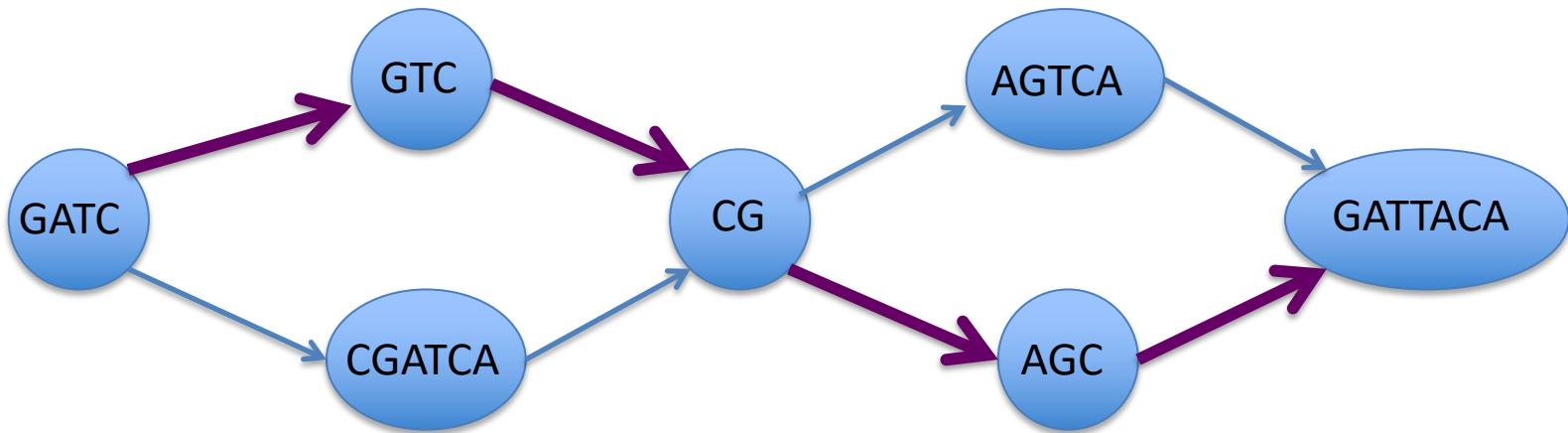
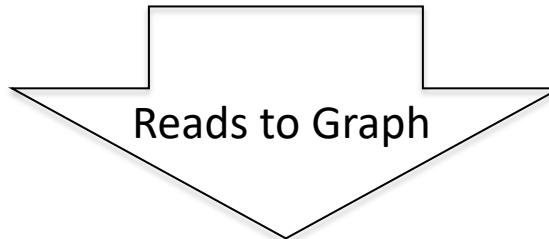


Nodes = sequence (+/-)
Edges = order, overlap

Graph Data Structures Commonly Used For Assembly



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

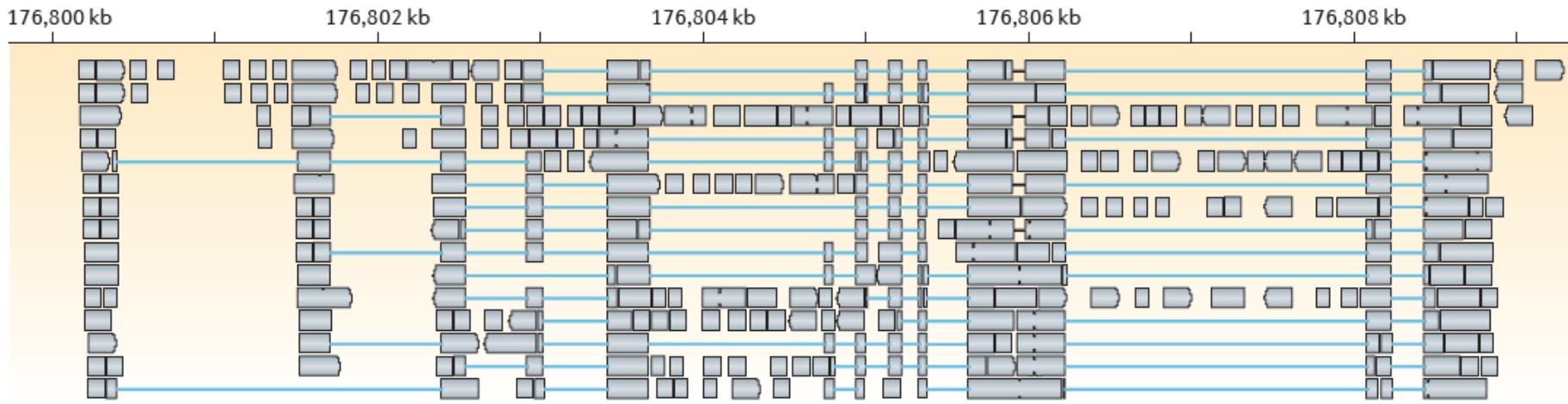


GATCGTCCGAGCGATTACA

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

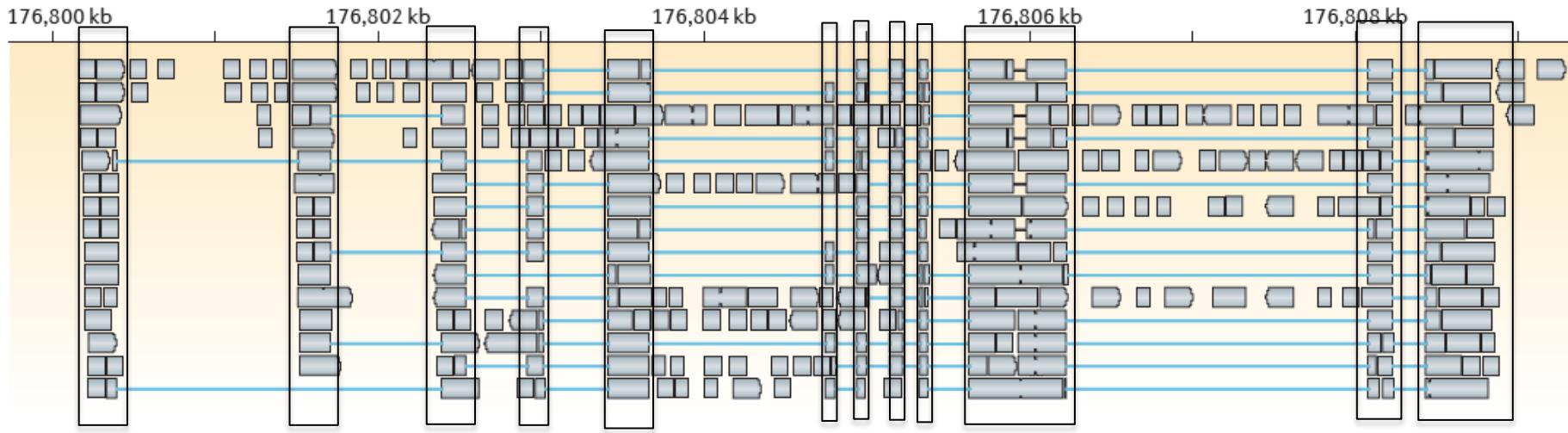
Genome-Guided Transcript Reconstruction

Splice-align reads to the genome



Genome-Guided Transcript Reconstruction

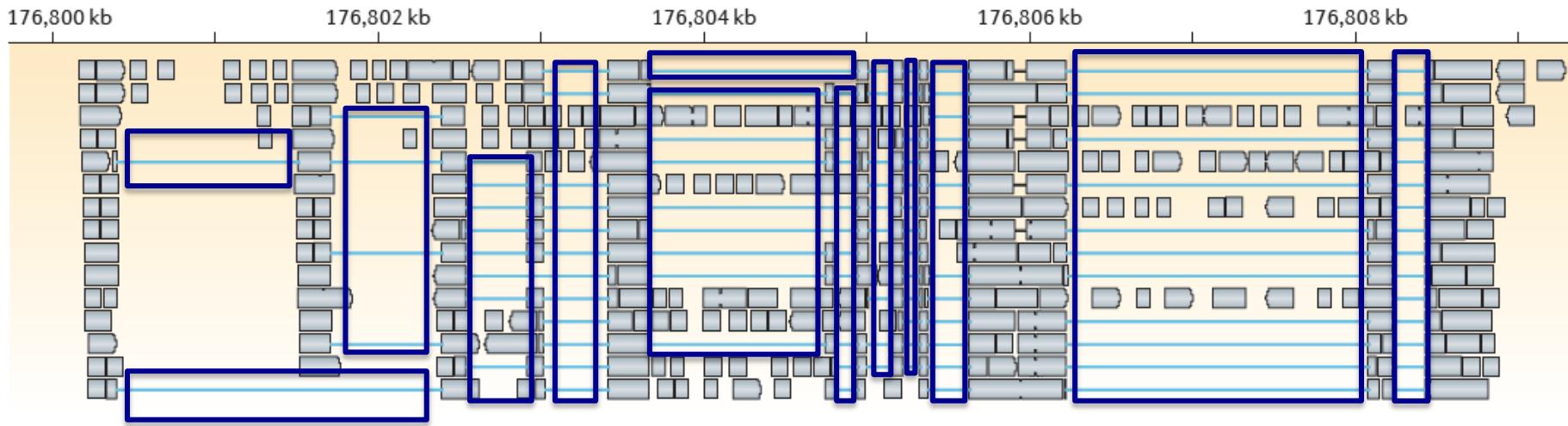
Splice-align reads to the genome



Alignment segment piles => exon regions

Genome-Guided Transcript Reconstruction

Splice-align reads to the genome



Large alignment gaps => introns

Genome-Guided Transcript Reconstruction

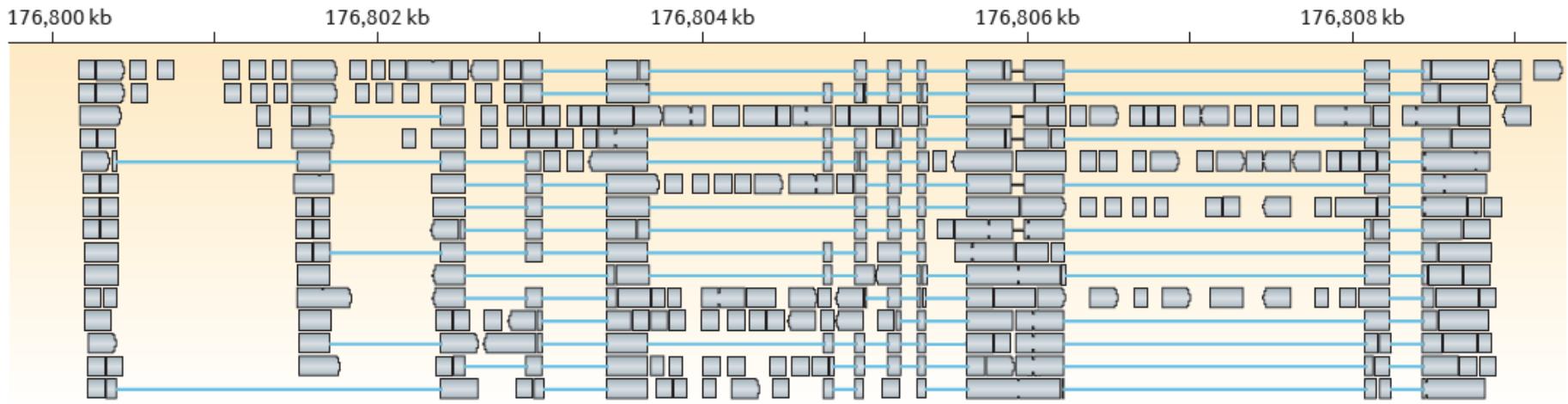
Splice-align reads to the genome



Overlapping but different introns = evidence of alternative splicing

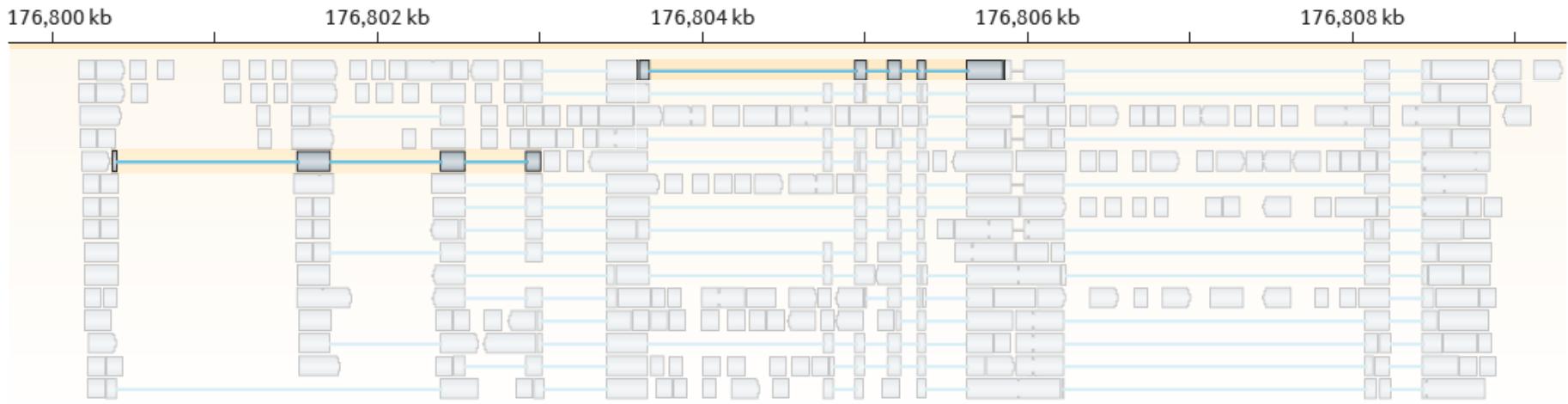
Genome-Guided Transcript Reconstruction

Splice-align reads to the genome



Genome-Guided Transcript Reconstruction

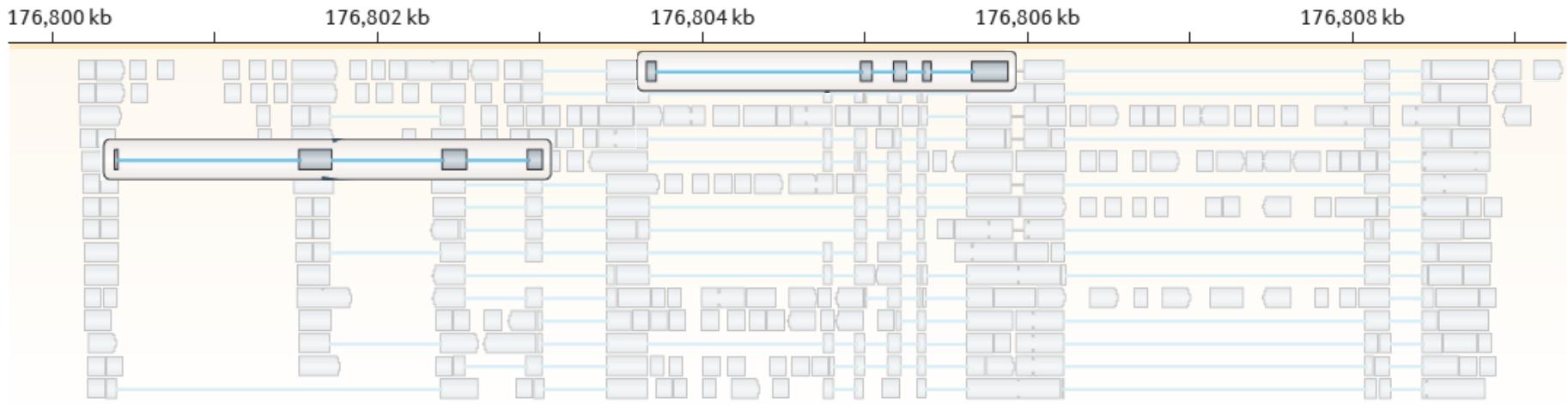
Splice-align reads to the genome



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

Genome-Guided Transcript Reconstruction

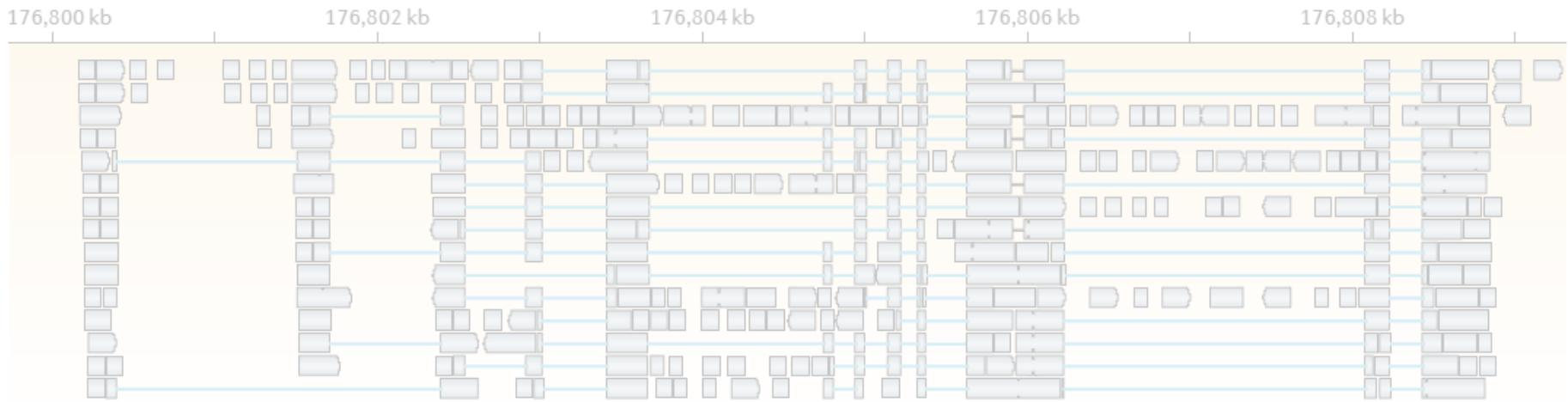
Splice-align reads to the genome



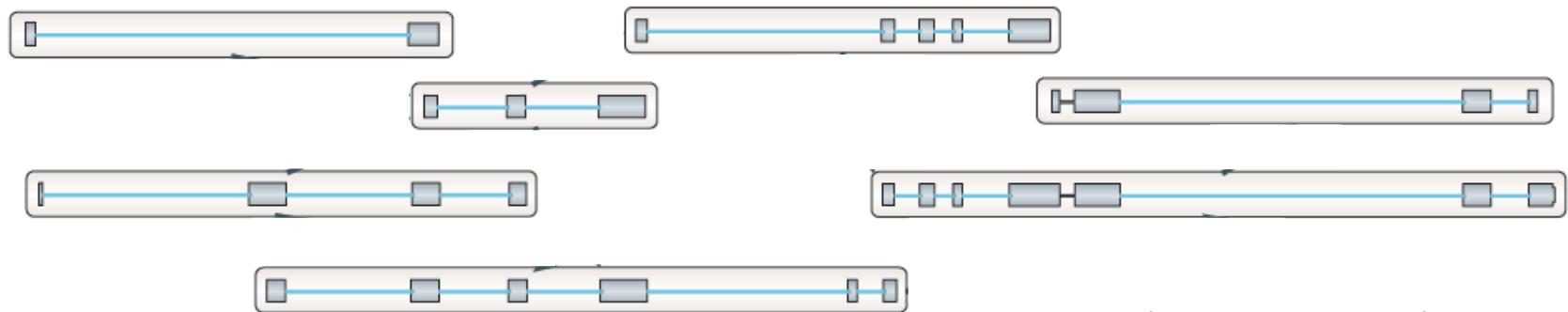
Nodes = unique splice patterns

Genome-Guided Transcript Reconstruction

Splice-align reads to the genome

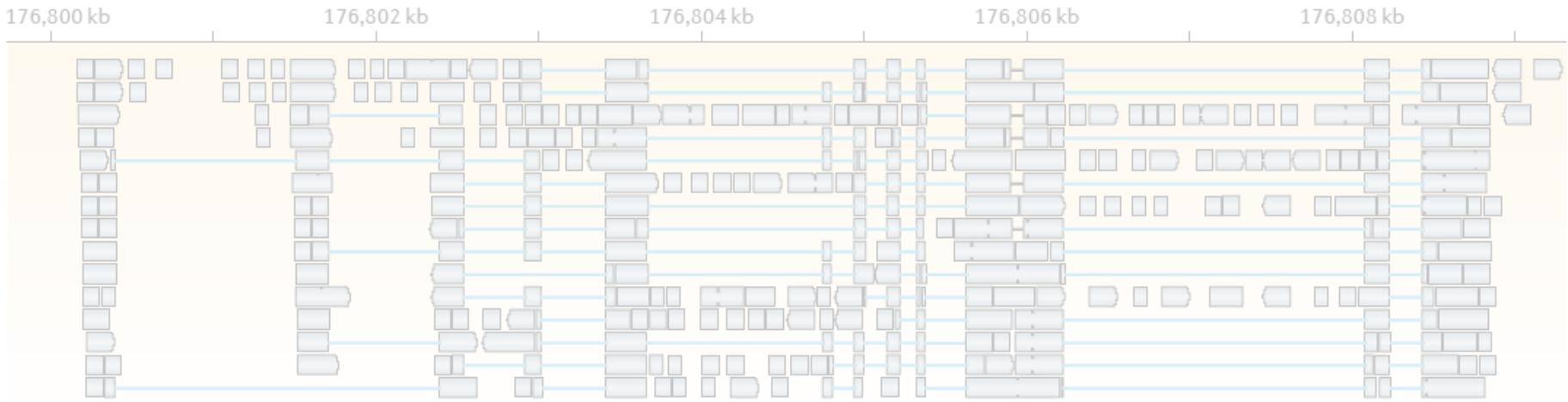


Construct graph from unique splice patterns of aligned reads.

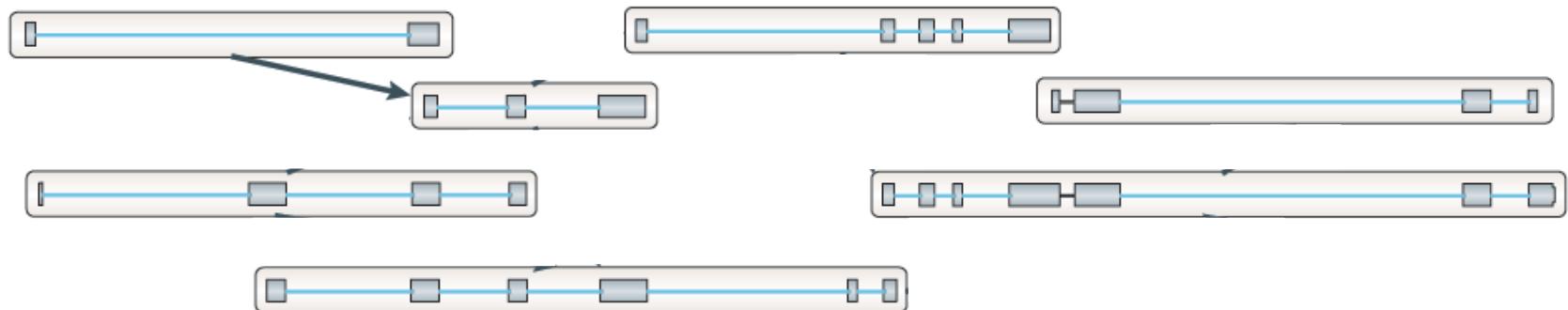


Genome-Guided Transcript Reconstruction

Splice-align reads to the genome



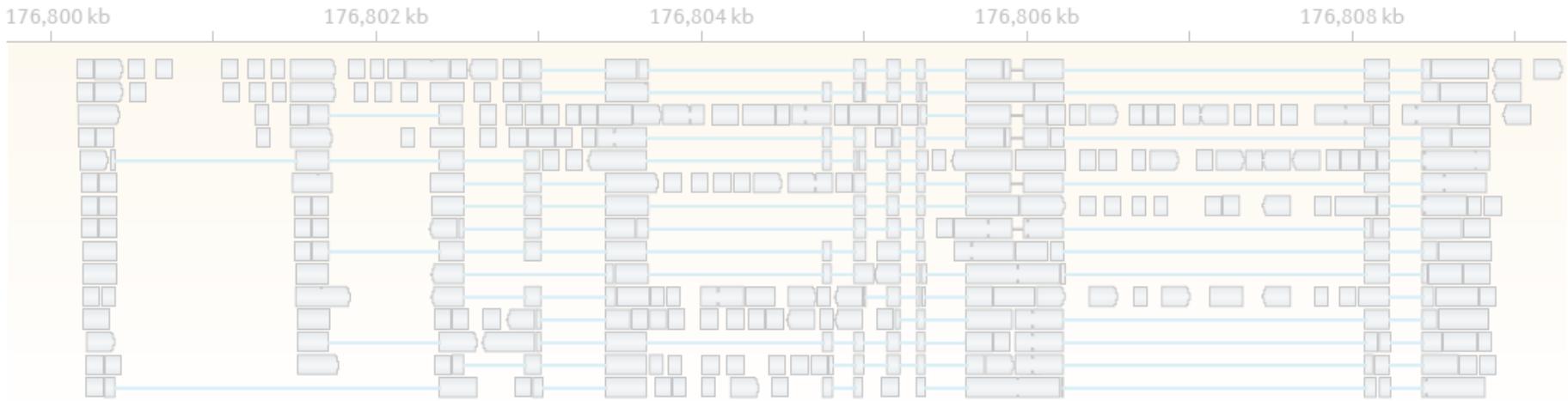
Construct graph from unique splice patterns of aligned reads.



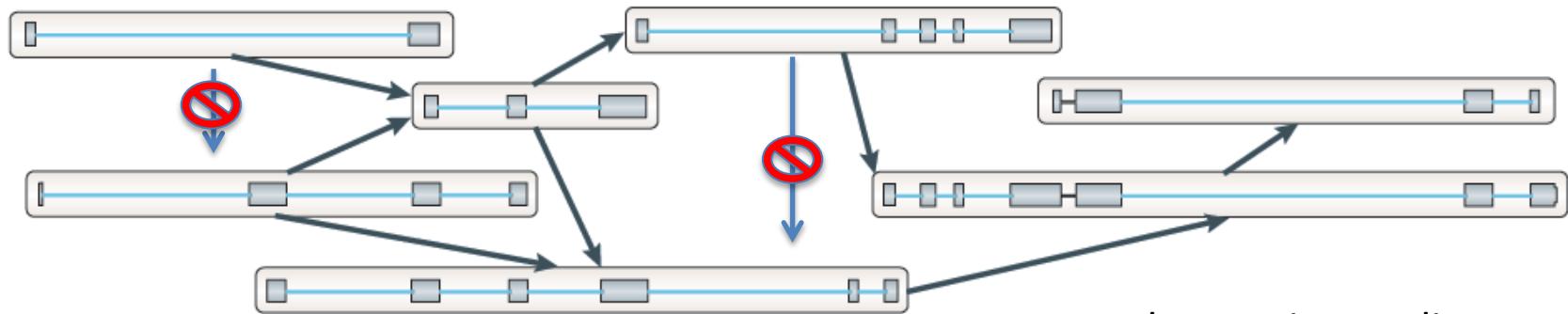
Nodes = unique splice patterns
Edges = compatible patterns

Genome-Guided Transcript Reconstruction

Splice-align reads to the genome

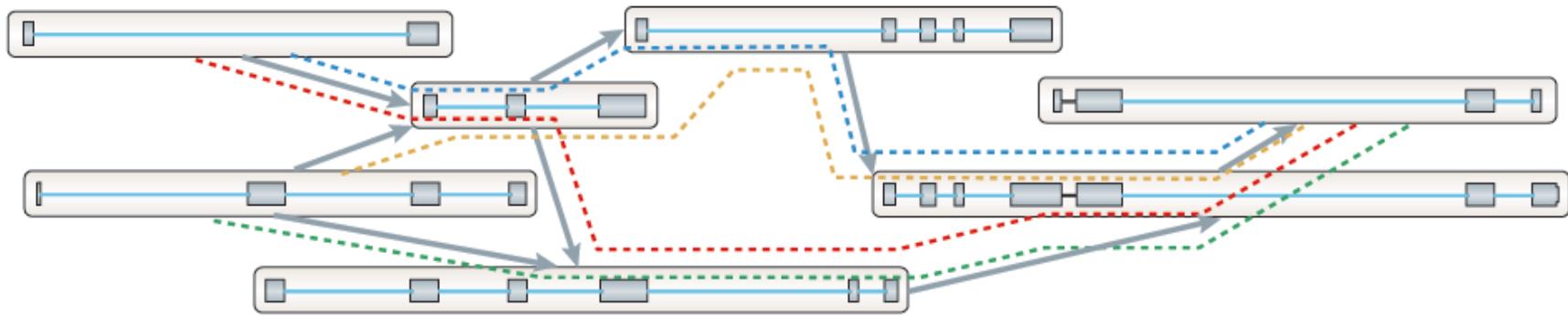


Construct graph from unique splice patterns of aligned reads.



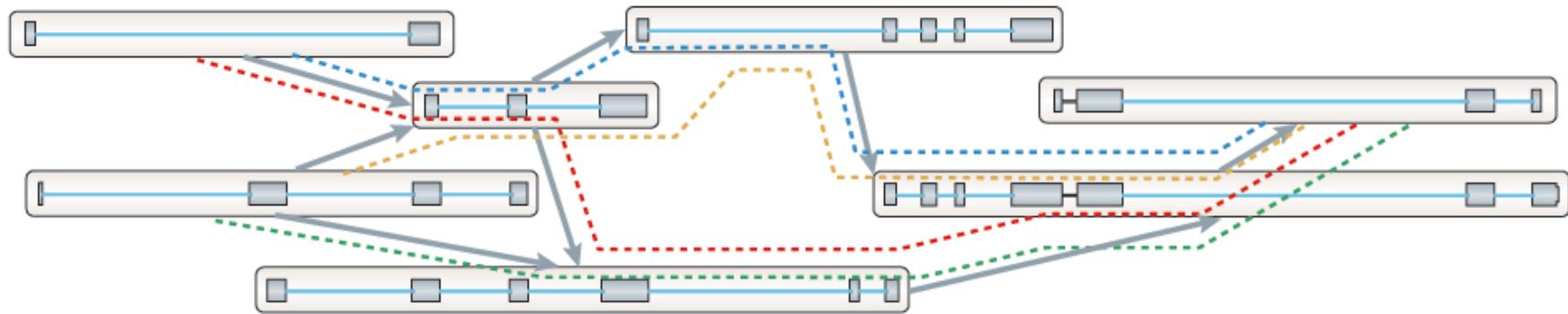
Genome-Guided Transcript Reconstruction

Traverse paths through the graph to assemble transcript isoforms

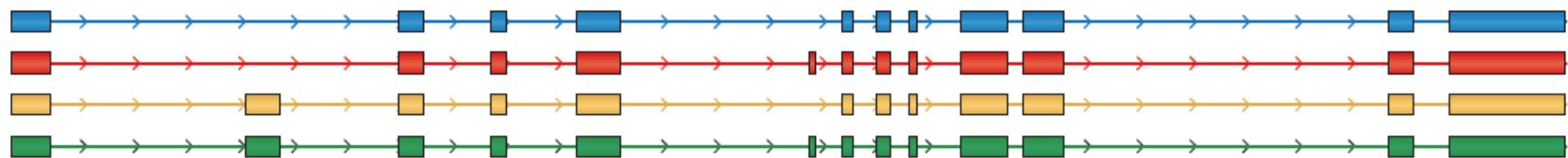


Genome-Guided Transcript Reconstruction

Traverse paths through the graph to assemble transcript isoforms

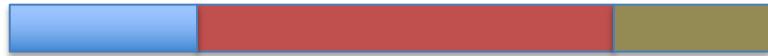
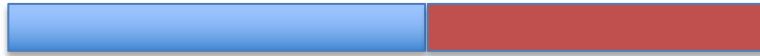


Reconstructed isoforms

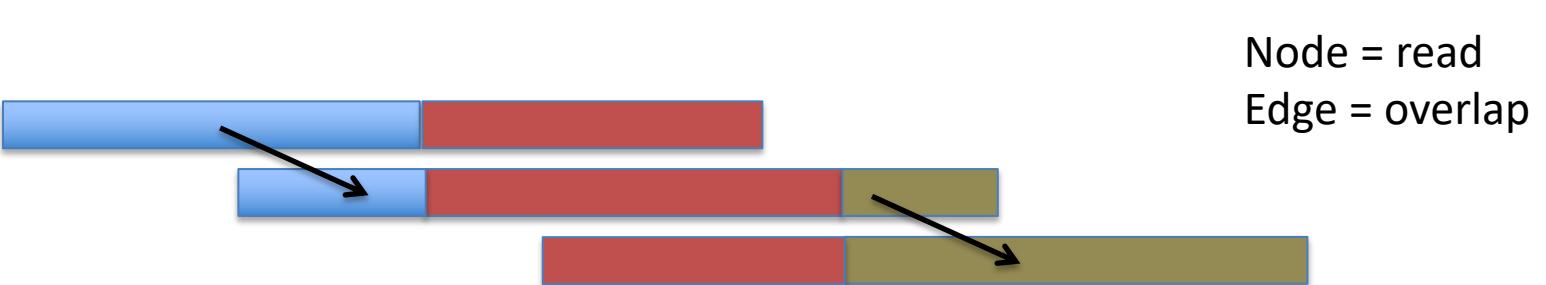


What if you don't have a high quality reference genome sequence?

Read Overlap Graph: Reads as nodes, overlaps as edges



Read Overlap Graph: Reads as nodes, overlaps as edges



Read Overlap Graph: Reads as nodes, overlaps as edges

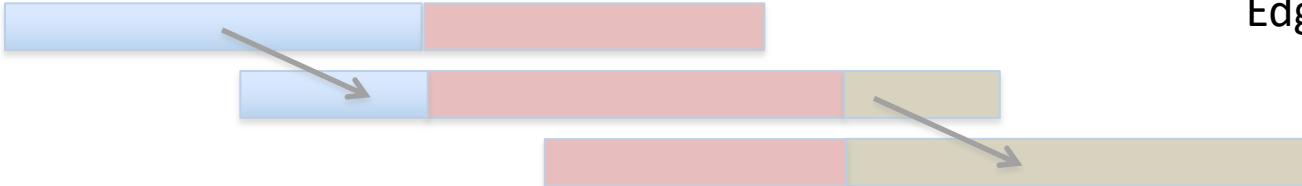


Transcript A



Generate consensus sequence where reads overlap

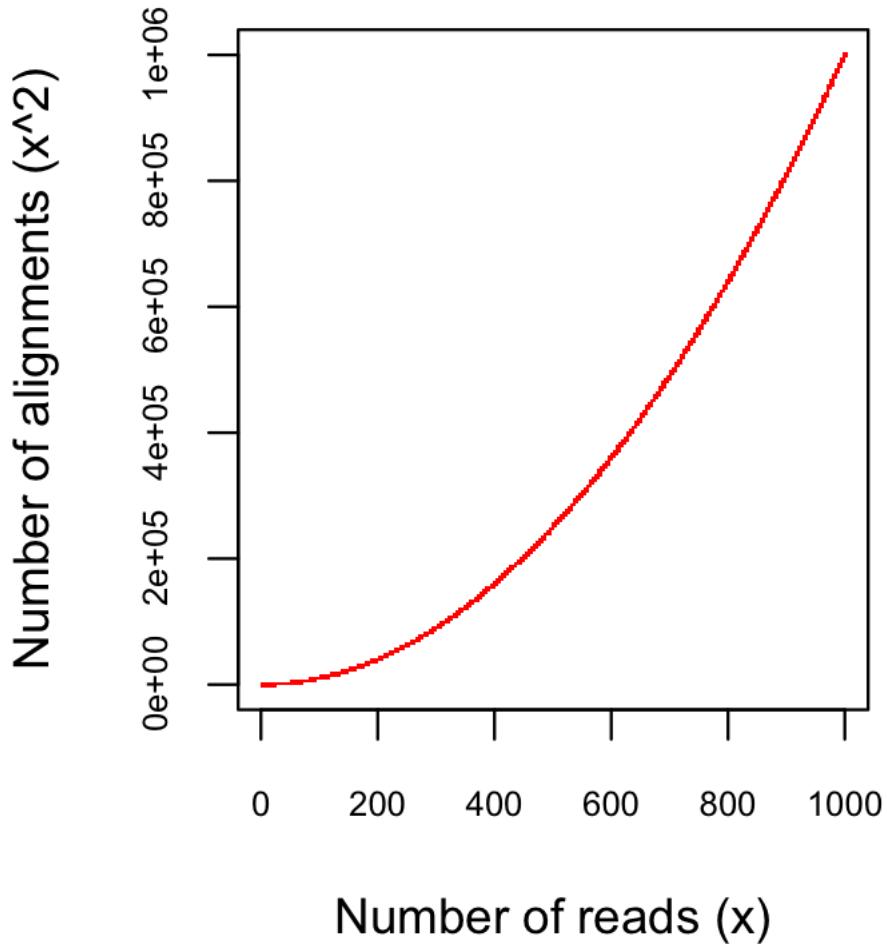
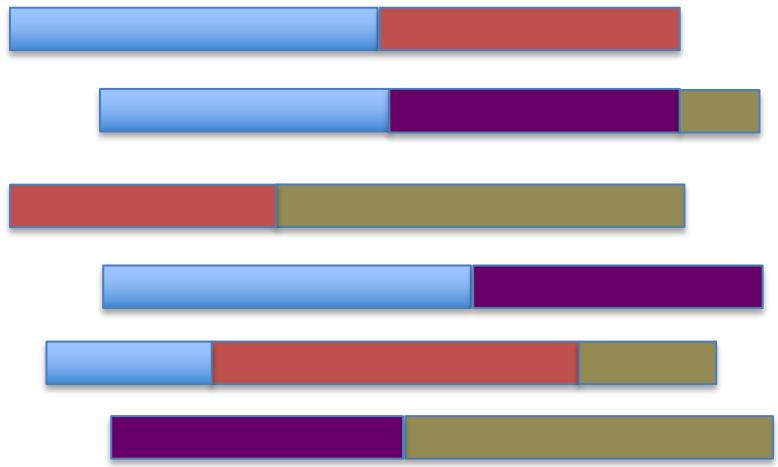
Node = read
Edge = overlap



Transcript B



Finding pairwise overlaps between n reads involves $\sim n^2$ 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



Sequence Assembly via De Bruijn Graphs

Generate all substrings of length k from the reads



Sequence Assembly via De Bruijn Graphs

Generate all substrings of length k from the reads



Construct the de Bruijn graph



Nodes = unique k-mers

Sequence Assembly via De Bruijn Graphs

Generate all substrings of length k from the reads



Construct the de Bruijn graph



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

Sequence Assembly via De Bruijn Graphs

Generate all substrings of length k from the reads



Construct the de Bruijn graph



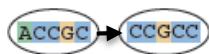
Nodes = unique k-mers
Edges = overlap by (k-1)

Sequence Assembly via De Bruijn Graphs

Generate all substrings of length k from the reads



Construct the de Bruijn graph



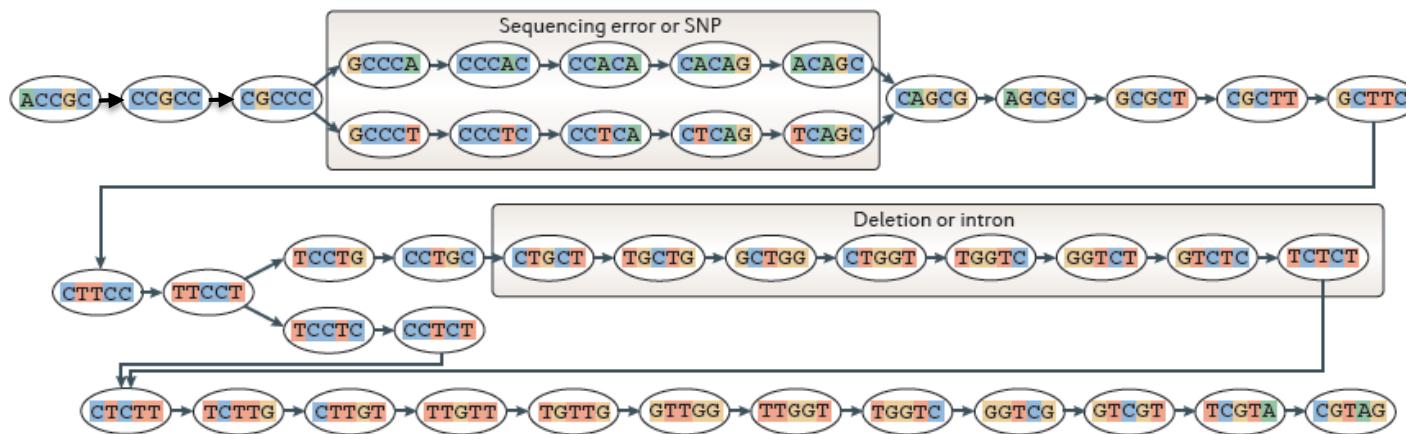
Nodes = unique k-mers
Edges = overlap by (k-1)

Sequence Assembly via De Bruijn Graphs

Generate all substrings of length k from the reads

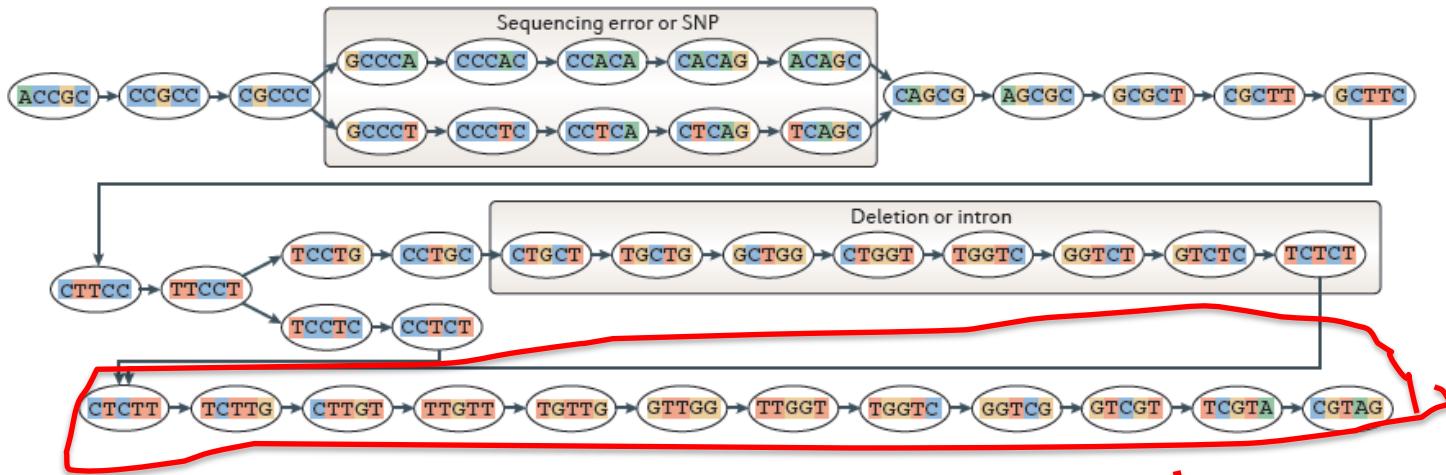
| | | | | | | | |
|------------------------------------|-------|-------|-------|--------------------------------|-------|-------|--------------|
| ACAGC | TCCTG | GTCTC | | AGCGC | CTCTT | GGTCG | k-mers (k=5) |
| CACAG | TTCCT | GGTCT | | CAGCG | CCTCT | TGGTC | |
| CCACA | CTTCC | TGGTC | TGTTG | TCAGC | TCCTC | TTGGT | |
| CCCAC | GCTTC | CTGGT | TTGTT | CTCAG | TTCCT | GTTGG | |
| GCCCA | CGCTT | GCTGG | CTTGT | CCTCA | CTTCC | TGTTG | |
| CGCCC | GCGCT | TGCTG | TCTTG | CCCTC | GCTTC | TTGTT | |
| CCGCC | AGCGC | CTGCT | CTCTT | GCCCT | CGCTT | CTTGT | |
| ACCGC | CAGCG | CCTGC | TCTCT | CGCCC | GCGCT | TCTTG | |
| ACCGCCCCACAGCGCTTCCTGCTGGTCTCTTGTG | | | | CGCCCTCAGCGCTTCCTCTTGTGGTCGTAG | | | |
| | | | | | | | Reads |

Construct the de Bruijn graph

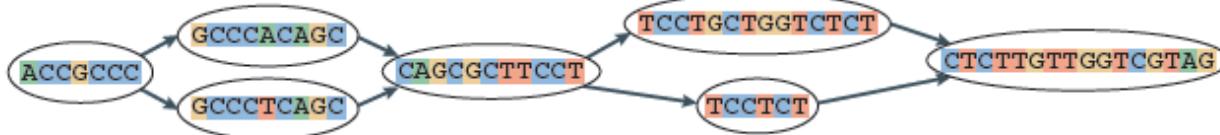


Nodes = unique k-mers
Edges = overlap by (k-1)

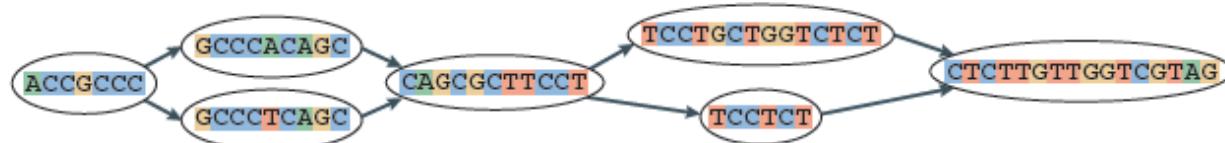
Construct the de Bruijn graph



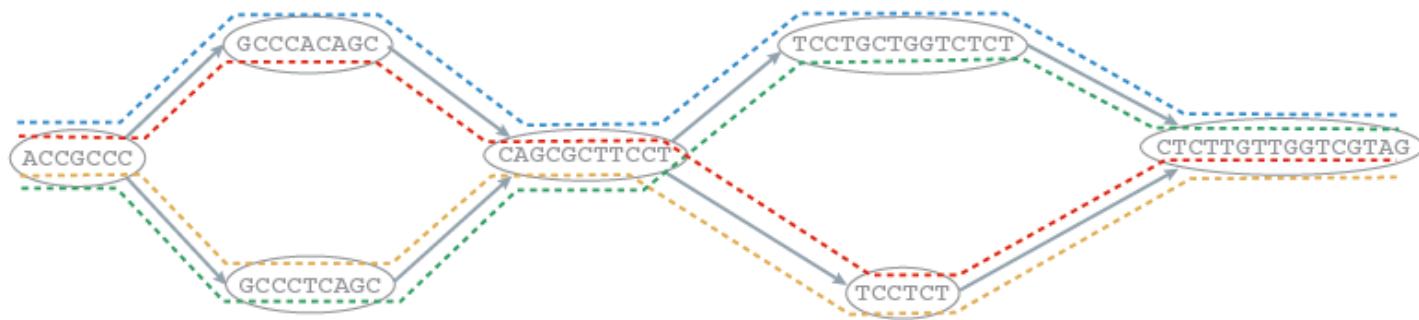
Collapse the de Bruijn graph



Collapse the de Bruijn graph



Traverse the graph



Assemble Transcript Isoforms

— ACCGGCCACAGCGCTTCCTGCTGGTCTCTTGGTGGT CGTAG
- - - ACCGGCCACAGCGCTTCCT - - - CTTGGTGGT CGTAG
--- ACCGGCCCTCAGCGCTTCCT --- - CTTGGTGGT CGTAG
---- ACCGGCCCTCAGCGCTTCCTGCTGGTCTCTTGGTGGT CGTAG

Contrasting Genome and Transcriptome *De novo* Assembly

Genome Assembly

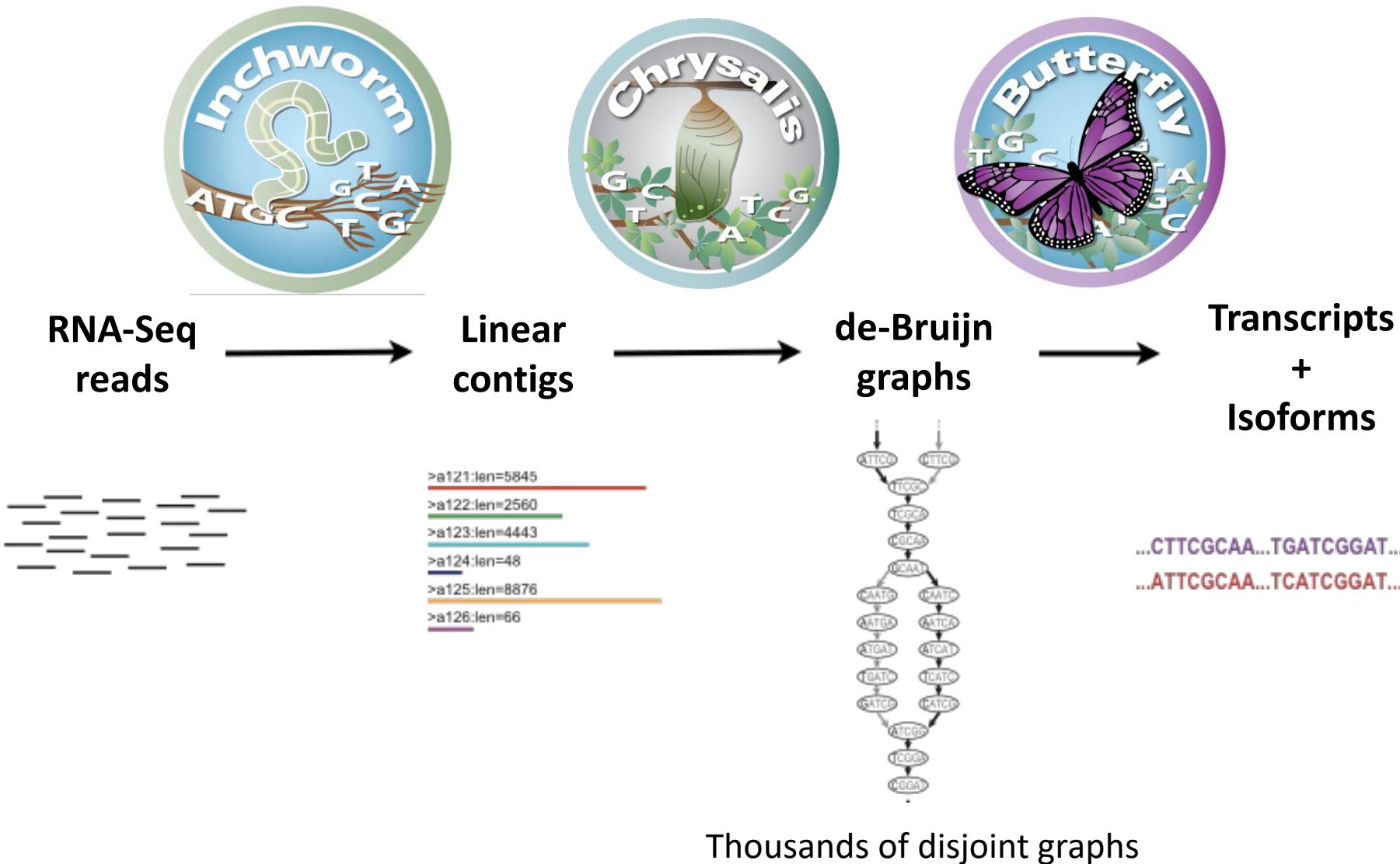
- Uniform coverage
- Single contig per locus
- Assemble small numbers of large Mb-length chromosomes
- Double-stranded data

Transcriptome Assembly

- Exponentially distributed coverage levels
- Multiple contigs per locus (alt splicing)
- Assemble many thousands of Kb-length transcripts
- Strand-specific data available



Trinity – How it works:





Inchworm Algorithm

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

Read: **AATGTGAAACTGGATTACATGCTGGTATGTC...**

AATGTGA

ATGTGAA

Overlapping kmers of length (k)

TGTGAAA

...

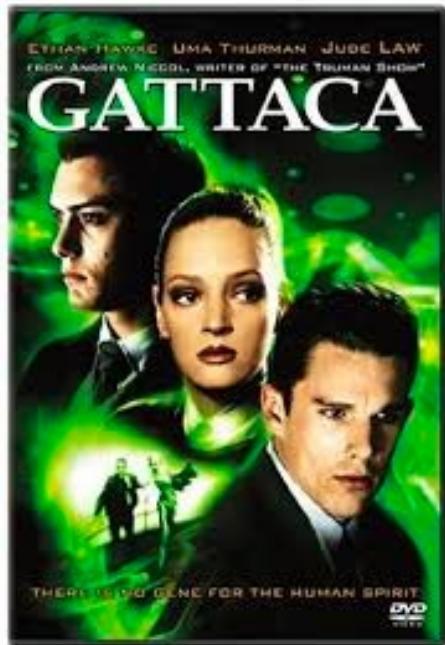
Kmer Catalog (hashtable)

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



Inchworm Algorithm

- Decompose all reads into overlapping Kmers => hashtable(kmer, count)
- Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.



<https://en.wikipedia.org/wiki/Gattaca>

GATTACA
9

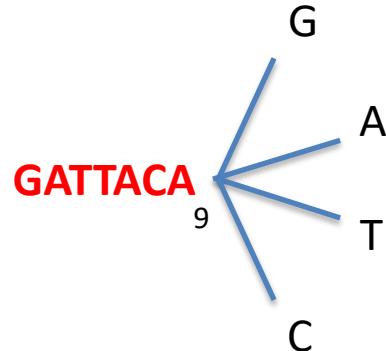
Kmer Catalog (hashtable)

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



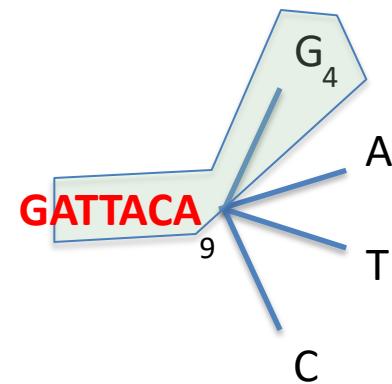
Inchworm Algorithm

- 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.



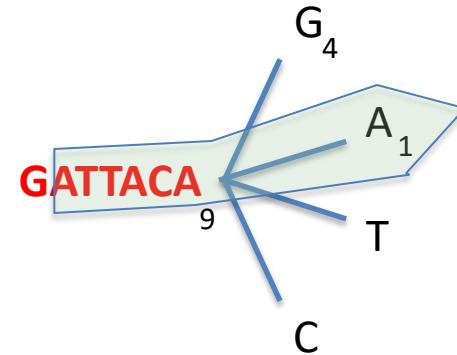


Inchworm Algorithm



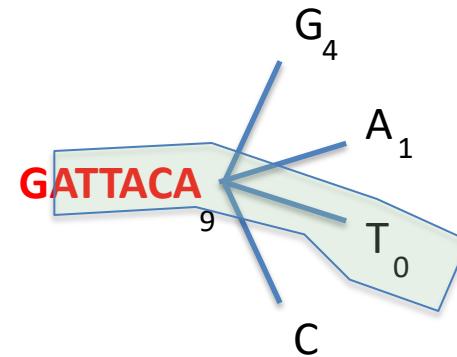


Inchworm Algorithm



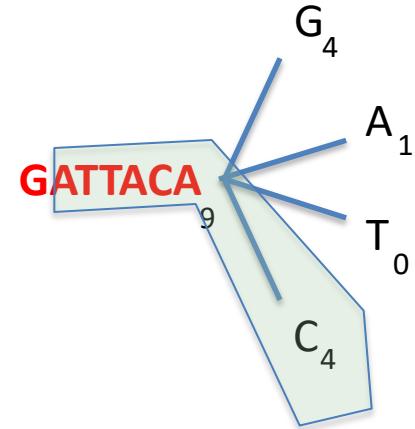


Inchworm Algorithm



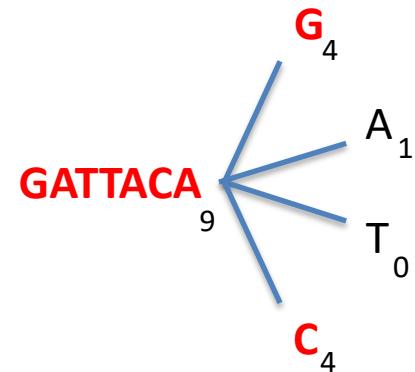


Inchworm Algorithm



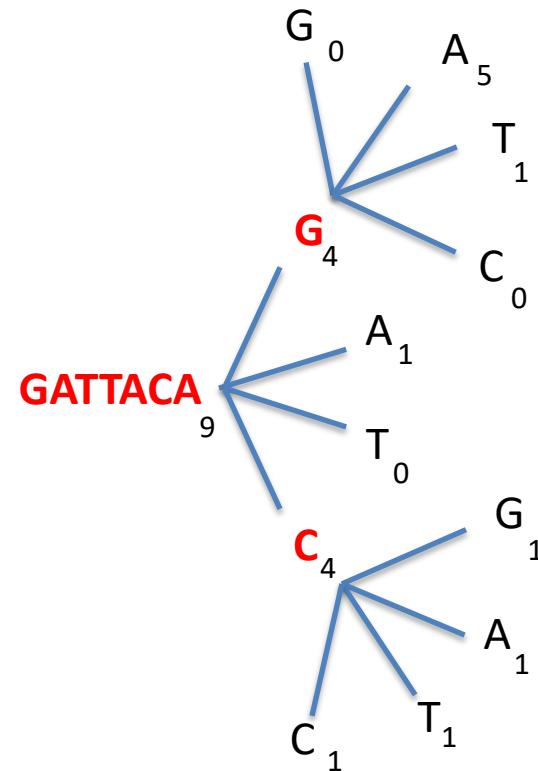


Inchworm Algorithm



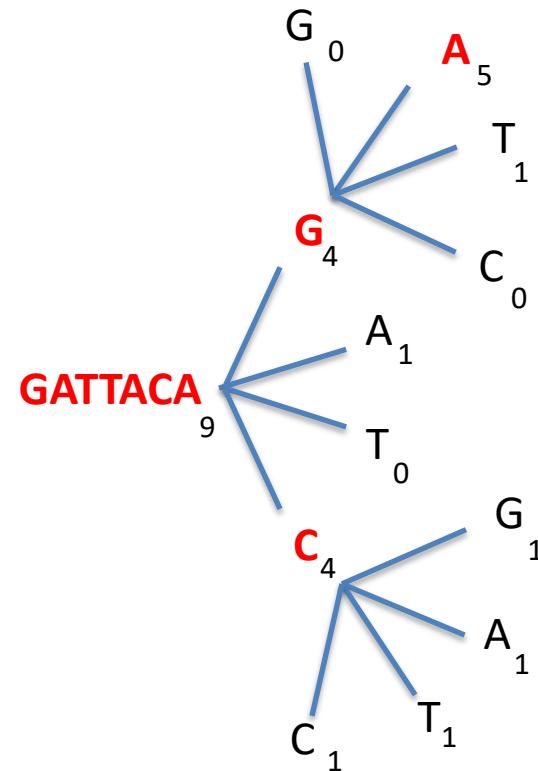


Inchworm Algorithm



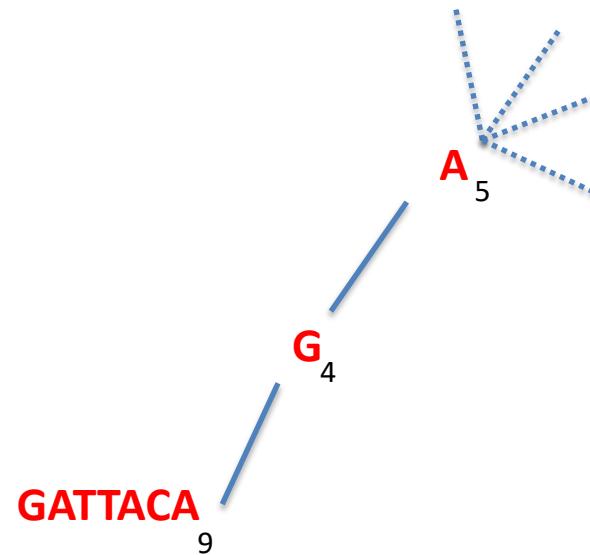


Inchworm Algorithm



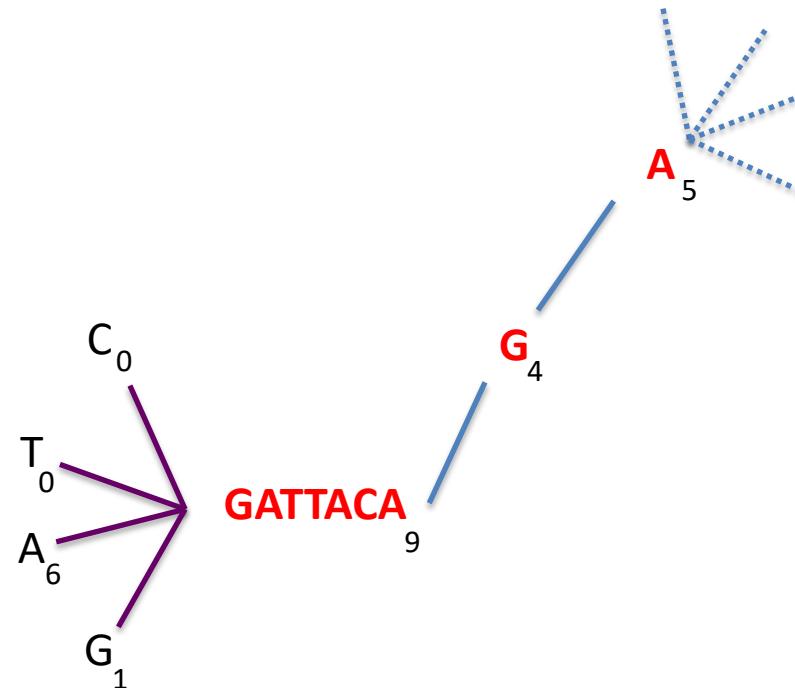


Inchworm Algorithm



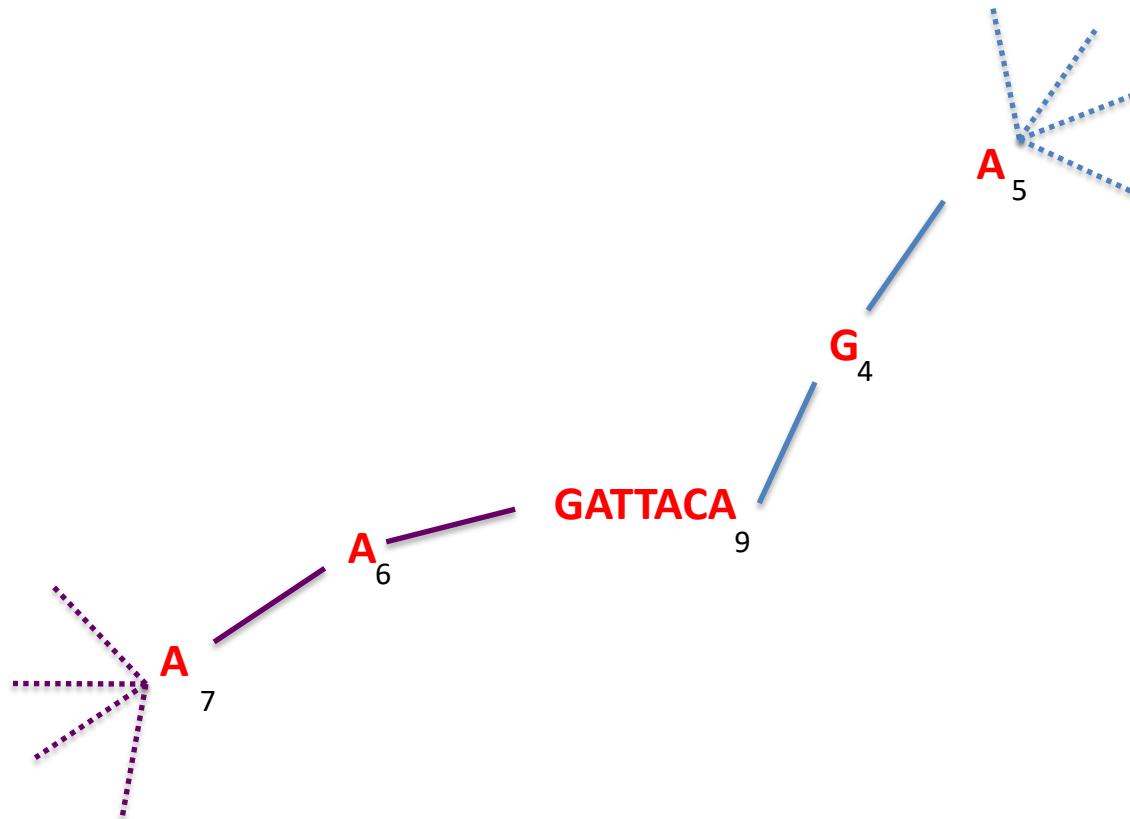


Inchworm Algorithm





Inchworm Algorithm



Report contig:**AAGATTACAGA**....

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

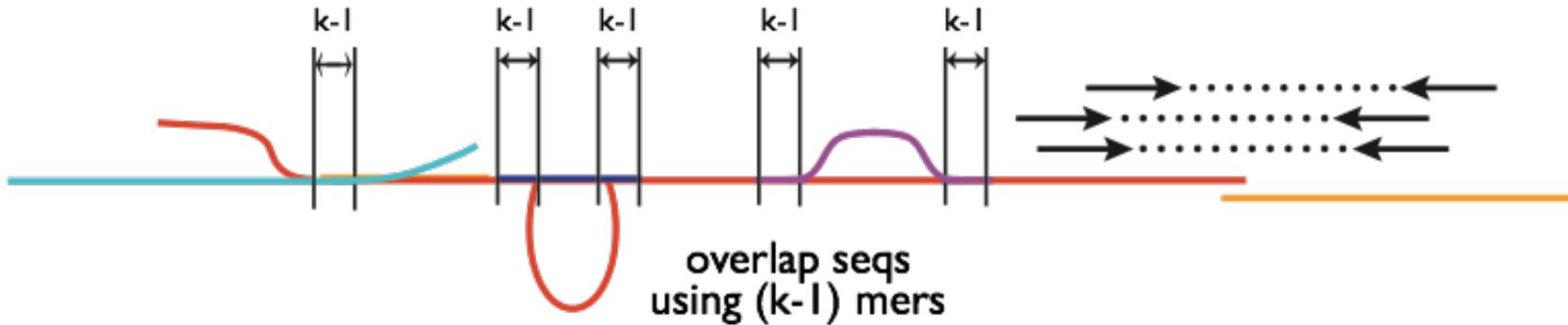
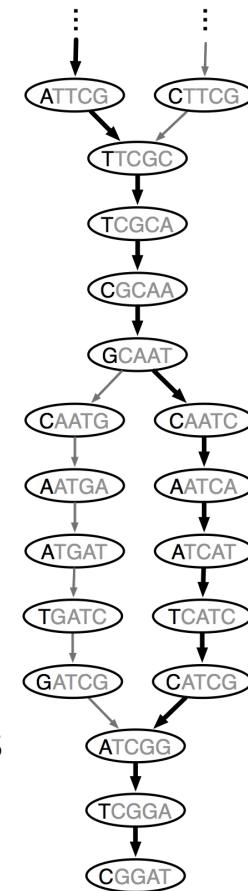
Chrysalis

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

Integrate isoforms via k-1 overlaps

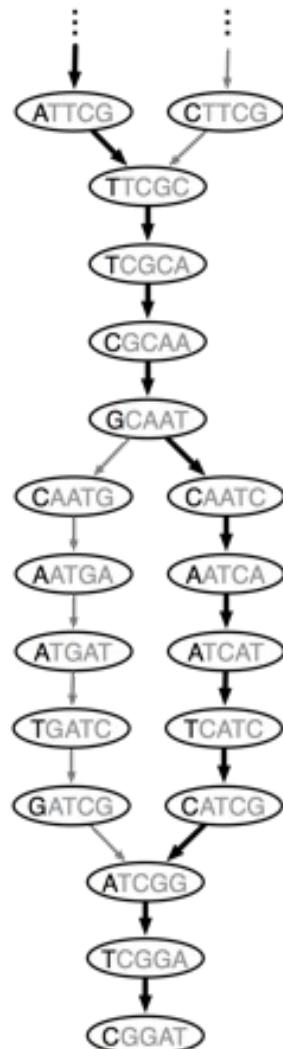


Build de Bruijn Graphs (ideally, one per gene)



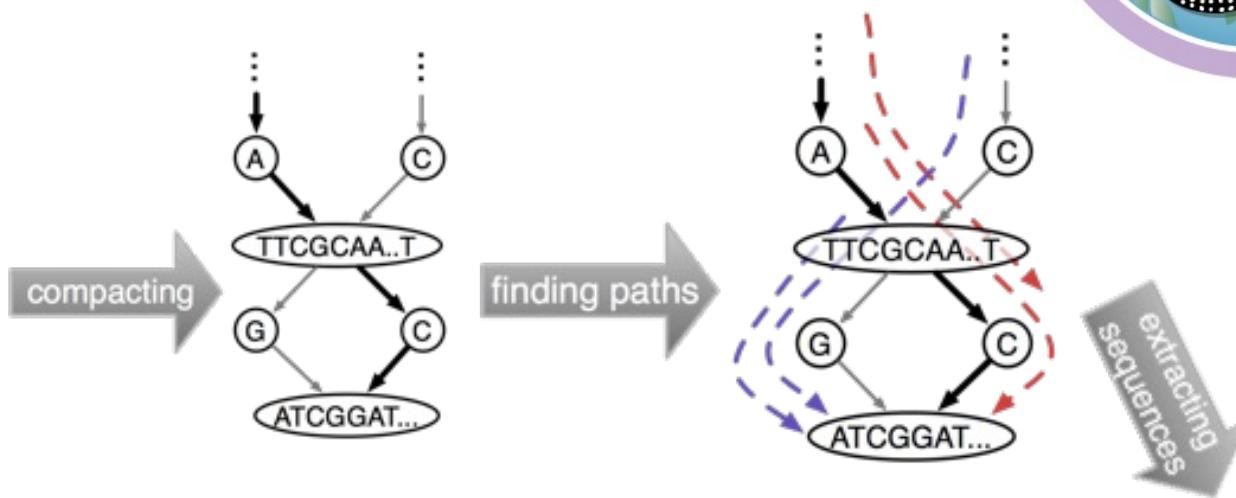


Thousands of Chrysalis Clusters



de Bruijn
graph

Butterfly



compact
graph

compact
graph with
reads

sequences
(isoforms and paralogs)



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

Trinity output: A multi-fasta file

NATURE PROTOCOLS | PROTOCOL

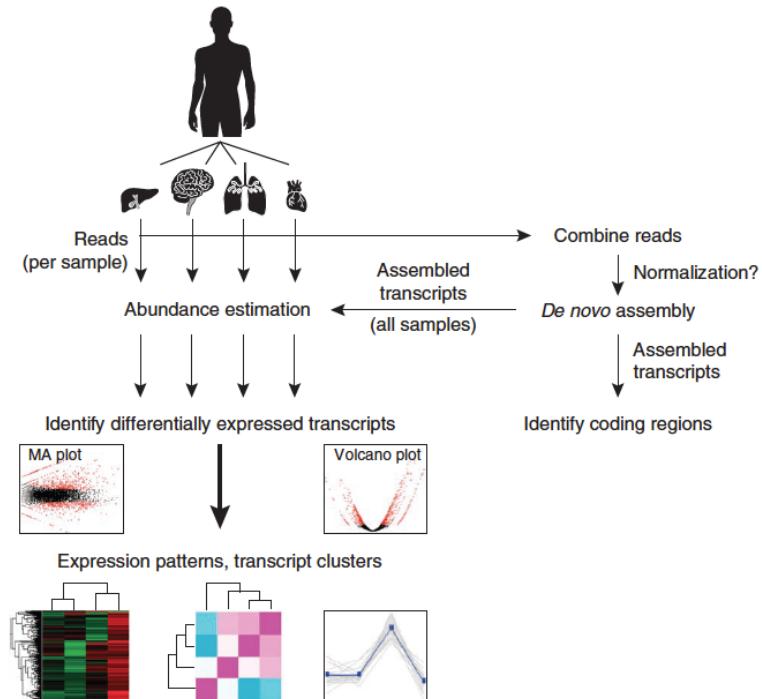
De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis

Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes, Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman, Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman & Aviv Regev

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

Nature Protocols 8, 1494–1512 (2013) | doi:10.1038/nprot.2013.084

Published online 11 July 2013



RNA-Seq De novo Assembly Using Trinity

► Pages 27



Quick Guide for the Impatient

Trinity assembles transcript sequences from Illumina RNA-Seq data.

Download Trinity [here](#).

Build Trinity by typing 'make' in the base installation directory.

Assemble RNA-Seq data like so:

```
Trinity --seqType fq --left reads_1.fq --right reads_2.fq --CPU 6 --max_memory 20G
```

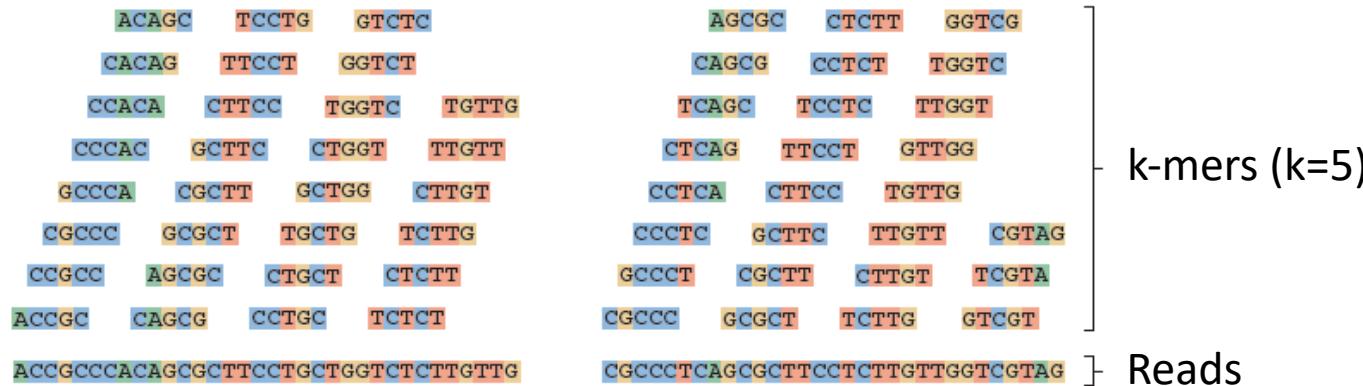
Find assembled transcripts as: 'trinity_out_dir/Trinity.fasta'

Use the documentation links in the right-sidebar to navigate this documentation, and contact our [Google group for technical support](#).

- [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](#)

ProgForBio Exercise : Build a program that counts k-mers

Generate all substrings of length k from the reads



Using a kmer size of 8 and reporting the top 10 kmers and their counts:

```
kmer_counter.py 8 reads.left.fq 10
```

| | |
|----------|------|
| TTTTTTTT | 1743 |
| AAAAAAA | 1204 |
| CCAGCCTT | 666 |
| GAAGCTGG | 627 |
| CAGGCAGG | 549 |
| TTTGCTGT | 536 |
| GCCTGCTG | 533 |
| CTTGGTCT | 525 |
| AAGCTGGA | 518 |
| TTTTATTT | 504 |

https://github.com/trinityrnaseq/CSHLProgForBio/tree/main/Exercise-counting_kmers

Build in 3 phases – demonstrating function reuse

A. retrieving the sequences from a fastq file

`fastq_file_to_sequence_list.py`

B. extracting kmers from an individual sequence

`sequence_to_kmer_list.py`

C. counting the kmers among all sequences

`count_kmers_from_fastq.py`

Functions imported
and reused

Useful structure of a python script for shareable functions

my_module.py (can run directly or use as a function library)

```
#!/usr/bin/env python3

import sys

def decorate_name_val(name_val):
    decorated = " ".join([get_unicorn(), name_val, get_unicorn()])
    return(decorated)

def get_unicorn():
    # returns unicorn symbol
    return("\U0001f984")

if __name__ == '__main__':
    print("my_module.py running as driver: ", decorate_name_val(__name__))
```

executed (base) wm4ca-d15:whatsupwith__name__? bhaas\$./my_module.py
my_module.py running as driver: 🦄 __main__ 🦄

When a script is executed, the special `__name__` variable is automatically set to '`__main__`' within the scope of that file.

Useful structure of a python script for shareable functions

my_module.py (can run directly or use as a function library)

```
#!/usr/bin/env python3

import sys

def decorate_name_val(name_val):
    decorated = " ".join([get_unicorn(), name_val, get_unicorn()])
    return(decorated)

def get_unicorn():
    # returns unicorn symbol
    return("\U0001f984")

if __name__ == '__main__':
    print("my_module.py running as driver: ", decorate_name_val(__name__))
```

my_script.py imported

```
#!/usr/bin/env python

from my_module import *

print("my_script.py: checking __name__ value:", decorate_name_val(__name__))
```

Uses imported `decorate_name_val()`

```
(base) wm4ca-d15:whatsupwith__name__? bhaas$ ./my_module.py
my_module.py running as driver: 🦄 __main__ 🦄
```

```
(base) wm4ca-d15:whatsupwith__name__? bhaas$ ./my_script.py
my_script.py: checking __name__ value: 🦄 __main__ 🦄
```

Build in 3 phases – demonstrating function reuse

A. retrieving the sequences from a fastq file

`fastq_file_to_sequence_list.py`

B. extracting kmers from an individual sequence

`sequence_to_kmer_list.py`

C. counting the kmers among all sequences

`count_kmers_from_fastq.py`

Functions imported
and reused

Part A: Retrieve sequences from a fastq file ↗

Write a python script that retrieves a list of all read sequences from a fastq file.

A script [fastq_file_to_sequence_list.py](#) is provided as a starting point. Fill in the missing code.

A fastq file <reads.fq> is provided as input.

The script usage is:

```
usage: ./fastq_file_to_sequence_list.py filename.fastq num_seqs_show
```

Running it like so:

```
fastq_file_to_sequence_list.py reads.fq 10
```

Should produce the following output:

```
['ACTGCATCCTGGAAAGAATCAATGGTGGCCGGAAAGTGTTTCAAATACAAGAGTGACAATGTGCCCTGTTGTT',  
'GTAATTCGCTACCTGCCACAGTGGCTCACCTGCTTAGAGGACAGGGAAAGGACCCCTAAAGGTAGGCTGATGC',  
'CTGGGCTGCAGCTAACGTTCTGCATCCTCCTTCTGCTTGTGGCTGGGAAGAAGACAATGTTGTCATGGCTGG',  
'CACGTTTCTAACGAGTTTACCAAGATCGTCTAACGCTCATTGTCTTGTACACACCAGTAAAGCTGGCA',  
'TGCTCATTGTCTTGTACACACCAGTAAAGCTGGCAAAATATCATCCAAAGTACATCGCTGAGAACTCCTA',  
'CCCACCTGAAAACATTTCTACATCCACTGTTATGGAATGCTTGATAAGCTTTCATTCTAACCATCAGAGCAC',  
'TCTGAATAAGTCCTGCCACCAATGTTTCATAAGTGTGGCATATGTTTCATTATTCAAACATTACTGTTAAG',  
'CTCCGTTTTGAGAGTGCAACACATAGATACTGCTTGATAGCATTAAACATCTCATTGCTGAAACAGG',  
'GCCTGAGTGTGCAAAATCTTCAGAGTAAGAATACCATAGTTGCTAAATATCTTACCATGAGCAATAATTTTT',  
'TCTGGTGCAGCTAGATGGAATACTGAGAAAATGTTCTCCATCCTGAACGAATATTGCAGCCTGAGAATTAACCA']
```

A. fastq_file_to_sequence_list.py

Reusable function part:

```
6    ## method: seq_list_from_fastq_file(fastq_filename)
7    ##
8    ## Extracts the sequence lines from a fastq file and returns a list
9    ## of the sequence lines
10   ##
11   ##
12   ## input parameters:
13   ##
14   ## fastq_filename : name of the fastq file (type: string)
15   ##
16   ## returns seq_list : list of read sequences.
17   ##                      ie. ["GATCGCATAG", "CGATGCAG", ...]
18
19  def seq_list_from_fastq_file(fastq_filename):
20
21      seq_list = list()
22
23      ## begin your code
24
25
26
27
28
29
30      ## end your code
31
32      return seq_list
33
```

Driver part for testing:

```
35
36  def main():
37
38      programe = sys.argv[0]
39
40      usage = "\n\n\tusage: {} filename.fastq num_seqs_show\n\n".format(programe)
41
42      if len(sys.argv) < 3:
43          sys.stderr.write(usage)
44          sys.exit(1)
45
46      # capture command-line arguments
47      fastq_filename = sys.argv[1]
48      num_seqs_show = int(sys.argv[2])
49
50      seq_list = seq_list_from_fastq_file(fastq_filename)
51
52      print(seq_list[0:num_seqs_show])
53
54      sys.exit(0) # always good practice to indicate worked ok!
55
56
57
58  if __name__ == '__main__':
59      main()
60
```

Build in 3 phases – demonstrating function reuse

A. retrieving the sequences from a fastq file

`fastq_file_to_sequence_list.py`

B. extracting kmers from an individual sequence

`sequence_to_kmer_list.py`

C. counting the kmers among all sequences

`count_kmers_from_fastq.py`

Functions imported
and reused

Part B: Extracting kmers from a sequence ↗

Write a python script to extract all kmers of a specified length from a nucleotide sequence.

A script [fastq_file_to_sequence_list.py](#) is provided as a starting point. Fill in the missing code.

The script usage is:

```
usage: ./sequence_to_kmer_list.py sequence kmer_length
```



Running it like so:

```
sequence_to_kmer_list.py ACTGCATCCTGGAAAGAATCAATGGTGGCCGGAAAGTGTTTCAAATACAAGAGTGACAATGTGCCCTGTTGTT 6
```



Should produce the following output:

```
['ACTGCA', 'CTGCAT', 'TGCATC', 'GCATCC', 'CATCCT', 'ATCCTG', 'TCCTGG', 'CCTGGA', 'CTGGAA', 'TGGAAA', 'GGAAAG', 'GAAAGA', 'AAAGAA', 'AAGAAT', 'AGAACAT', 'GAATCA', 'AATCAA', 'ATCAAT', 'TCAATG', 'CAATGG', 'AATGGT', 'ATGGTG', 'TGGTGG', 'GGTGGC', 'GTGGCC', 'TGGCCG', 'GGCCGG', 'GCCGGG', 'CCGGAA', 'CGGAAA', 'GGAAAG', 'GAAAGT', 'AAAGTG', 'AAGTGT', 'AGTGTT', 'GTGTTT', 'TGTTTT', 'GTTTTT', 'TTTTTC', 'TTTTCA', 'TTTCAA', 'TTCAAA', 'TCAAAT', 'CAAATA', 'AAATAC', 'AATACA', 'ATACAA', 'TACAAG', 'ACAAGA', 'CAAGAG', 'AAGAGT', 'AGAGTG', 'GAGTGA', 'AGTGAC', 'GTGACA', 'TGACAA', 'GACAAT', 'ACAATG', 'CAATGT', 'AATGTG', 'ATGTGC', 'TGTGCC', 'GTGCC', 'TGCCCT', 'GCCCTG', 'CCCTGT', 'CCTGTT', 'CTGTTG', 'TGTTGT', 'GTTGTT', 'TTGTTT']
```

B. sequence_to_kmer_list.py

Reusable function part

```
6
7 ## method: sequence_to_kmer_list(sequence, kmer_length)
8 ##
9 ## Extracts all kmers of a specified length from a sequence
10 ##
11 ## ie. sequence: GATCGATCGATCGA
12 ## and given kmer_length = 4
13 ## would yield
14 ##             GATC
15 ##             ATCG
16 ##             TCGA
17 ##             .... and so forth
18 ##
19 ## input parameters:
20 ##
21 ## sequence : nucleotide sequence (type: string)
22 ##
23 ## returns kmer_list : list of kmer sequences.
24 ##                      ie. ["GATC", "ATCG", ...]
25
26 def sequence_to_kmer_list(sequence, kmer_length):
27
28     kmers_list = list()
29
30     ## begin your code
31
32
33
34
35
36     ## end your code
37
38
39     return kmers_list
40
```

Driver part for testing:

```
41
42 def main():
43
44     programe = sys.argv[0]
45
46     usage = "\n\n\tusage: {} sequence kmer_length\n\n".format(programe)
47
48     if len(sys.argv) < 3:
49         sys.stderr.write(usage)
50         sys.exit(1)
51
52     # capture command-line arguments
53     sequence = sys.argv[1]
54     kmer_length = int(sys.argv[2])
55
56     kmers = sequence_to_kmer_list(sequence, kmer_length)
57
58     print(kmers)
59
60     sys.exit(0) # always good practice to indicate worked ok!
61
62
63
64 if __name__ == '__main__':
65     main()
66
```

Build in 3 phases – demonstrating function reuse

A. retrieving the sequences from a fastq file

`fastq_file_to_sequence_list.py`

B. extracting kmers from an individual sequence

`sequence_to_kmer_list.py`

C. counting the kmers among all sequences

`count_kmers_from_fastq.py`

Functions imported
and reused

Part C: Counting all kmers from all sequences in a fastq file

Now, let's count all kmers in all sequences. We can leverage each of the methods implemented above. Because of the way we wrote the above scripts, we can leverage them as a code library and simply import them for use in a new script.

Use the script [count_kmers_from_fastq.py](#) as the starting point. You'll see at the top of this script:

```
from sequence_to_kmer_list import *
from fastq_file_to_sequence_list import *
```



Those lines import the methods we implemented earlier so that we can just reuse them without having to rewrite or copy/paste any code in this new script.

The usage of our script is:

```
usage: ./count_kmers_from_fastq.py filename.fastq kmer_length num_top_kmers_show
```



And when we run it like so:

```
count_kmers_from_fastq.py reads.fq 6 10
```



It should produce the output:

```
TTTTTT: 3085
CTTCTT: 2550
AAAAAA: 2498
CTGCTG: 2446
AGCTGG: 2400
CAGCAG: 2265
CAGCTG: 2243
TCTTCT: 2208
CTGGAG: 2174
TGCTGT: 2156
```



C. count_kmers_from_fastq.py

Main code block

Import earlier functions

```
1  #!/usr/bin/env python
2
3  import os, sys
4
5  from sequence_to_kmer_list import *
6  from fastq_file_to_sequence_list import *
7
```

New function to implement

```
8
9  ## method: count_kmers(kmer_list)
10 ##
11 ## Counts the frequency of each kmer in the given list of kmers
12 ##
13 ## input parameters:
14 ##
15 ## kmer_list : list of kmers (type: list)
16 ##           ie. ["GATC", "TCGA", "GATC", ...]
17 ##
18 ##
19 ## returns kmer_counts_dict : dict containing ( kmer : count )
20 ##           ie. { "GATC" : 2,
21 ##                  "TCGA" : 1,
22 ##                  ...
23 ##
24
25  def count_kmers(kmer_list):
26
27      kmer_count_dict = dict()
28
29      #####
30      ## Step 2:
31      ## begin your code
32
33
34
35  Step 2
36
37
38
39
40      ## end your code
41      #####
42
43
44      return kmer_count_dict
```

```
45
46  def main():
47
48      programe = sys.argv[0]
49
50      usage = "\n\nusage: {} filename.fastq kmer_length num_top_kmers_show\n\n".format(
51          programe
52      )
53
54      if len(sys.argv) < 4:
55          sys.stderr.write(usage)
56          sys.exit(1)
57
58      # capture command-line arguments
59      fastq_filename = sys.argv[1]
60      kmer_length = int(sys.argv[2])
61      num_top_kmers_show = int(sys.argv[3])
62
63      seq_list = seq_list_from_fastq_file(fastq_filename)
64
65      all_kmers = list()
66
67      #####
68      ## Step 1:
69      ## begin your code, populate 'all_kmers' list with the
70      ## collection of kmers from all sequences
71
72
73      ## end your code
74      #####
75
76      kmer_count_dict = count_kmers(
77          all_kmers
78      ) # see step 2 above. You implement this. :-)
79
80
81      unique_kmers = list(kmer_count_dict.keys())
82
83      #####
84      ## Step 3: sort unique_kmers by abundance descendingly
85      ## (Note, you can run and test without first implementing Step 3)
86      ## begin your code      hint: see the built-in 'sorted' method documentation
87
88      ## end your code
89
90      ## printing the num top kmers to show
91      top_kmers_show = unique_kmers[0:num_top_kmers_show]
92
93
94      for kmer in top_kmers_show:
95          print("{}: {}".format(kmer, kmer_count_dict[kmer]))
96
97
98      sys.exit(0) # always good practice to indicate worked ok!
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113  if __name__ == "__main__":
114      main()
```

Step 1

Step 3

Extra Credit

It should produce the output:

```
TTTTTT: 3085
CTTCTT: 2550
AAAAAA: 2498
CTGCTG: 2446
AGCTGG: 2400
CAGCAG: 2265
CAGCTG: 2243
TCTTCT: 2208
CTGGAG: 2174
TGCTGT: 2156
```



Extra credit section: ↴

If you've accomplished the above, here's another challenge!

Note that the top-most kmer is of low complexity. If we are going to perform downstream operations like assembly and want to start with a seed kmer, we might want to avoid low complexity kmers as they would lack specificity.

Challenge: include another method that computes the complexity of each kmer using Shannon's Entropy (example: see: https://en.wikipedia.org/wiki/Sequence_logo#Logo_creation), and picture the kmer as representing one column of the seqlogo for which you would get one entropy calculation.

Add the entropy value as another column in the above printing.