### RNA-Seq Alignment/Assembly using STAR CS 4390/5390 and TopHat if we have time.



# The Central Dogma

#### **RNA**

- (non-coding RNA)





**RNA** sequencing







**RNA** sequencing







**RNA** sequencing









#### **RNA** sequencing









#### **RNA** sequencing





Split the read into chunks, and align those seperately

- the chunks end up being small
- can be done at random positions or predicted intron boundaries

- Split the read into chunks, and align those seperately can be done at random positions or predicted intron boundaries
  - the chunks end up being small

Align to the known transcriptome rather than the genome can miss novel transcripts (ones we have not seen before)

Split the read into chunks, and align those seperately can be done at random positions or predicted intron boundaries the chunks end up being small

- Align to the known transcriptome rather than the genome can miss novel transcripts (ones we have not seen before)
- Design a new aligner that takes these complications into account

Split the read into chunks, and align those seperately can be done at random positions or predicted intron boundaries the chunks end up being small

Align to the known transcriptome rather than the genome can miss novel transcripts (ones we have not seen before)

Design a new aligner that takes these complications into account

Some combination of the above techniques

#### **Spliced Transcripts Alignment to a Reference** (STAR)

Designed to align non-contiguous sections of a read, directly to the reference genome

Consists of two major steps:

- Seed searching
- clustering, stitching, and scoring

#### OINFORMATICS ORIGINAL PAPER

Sequence analysis

#### STAR: ultrafast universal RNA-seq aligner

Alexander Dobin<sup>1,\*</sup>, Carrie A. Davis<sup>1</sup>, Felix Schlesinger<sup>1</sup>, Jorg Drenkow<sup>1</sup>, Chris Zaleski<sup>1</sup>, Sonali Jha<sup>1</sup>, Philippe Batut<sup>1</sup>, Mark Chaisson<sup>2</sup> and Thomas R. Gingeras<sup>1</sup> <sup>1</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and <sup>2</sup>Pacific Biosciences, Menlo Park, CA, USA Associate Editor: Inanc Birol

Vol. 29 no. 1 2013, pages 15-21 doi:10.1093/bioinformatics/bts635

Advance Access publication October 25, 2012

Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring  $R[i \dots (i + MML - 1)]$ • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$  $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$ 

• where MML is the Maximal Mapping Length



Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring R[i ... (i + MML - 1)]

- such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$

• where MML is the Maximal Mapping Length

This is similar to a Maximal Exact Match (MEM) or Maximal Unique Map (MUM), the latter requires that only one location in the genome matches that location. In both cases, the value if *i* is not given.



 $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$ 

- $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$
- Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring R[i ... (i + MML - 1)] • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$ • where MML is the Maximal Mapping Length

- map from the start of the read as far as possible
- restart searching from the next position to the right



- $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$
- Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring R[i ... (i + MML - 1)] • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$ • where MML is the Maximal Mapping Length

- map from the start of the read as far as possible restart searching from the next position to the right



- $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$
- Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring R[i ... (i + MML - 1)] • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$ • where MML is the Maximal Mapping Length

- map from the start of the read as far as possible
- restart searching from the next position to the right



- $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$
- Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring R[i ... (i + MML - 1)] • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$ • where MML is the Maximal Mapping Length

- map from the start of the read as far as possible restart searching from the next position to the right



Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring  $R[i \dots (i + MML - 1)]$ • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$  $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$ • where MML is the Maximal Mapping Length

- map from the start of the read as far as possible
- restart searching from the next position to the right





- $R[i \dots (i+MML-1)] = G[j \dots (j_k+MML-1)]$
- Maximal Mappable Prefix (MMP) for read R, read start location i, and genome G: • the longest substring  $R[i \dots (i + MML - 1)]$ • such that there exists some set  $J = \{j_1, j_2, \dots, j_n\}$  where for all  $j_k \in J$ • where MML is the Maximal Mapping Length

The basic algorithm is

- map from the start of the read as far as possible
- restart searching from the next position to the right

The search is implemented using an uncompressed suffix array(s)

one for each chromosome

- one for each chromosome
- MMP mapping comes "free" with binary search

### Seed Search





- one for each chromosome
- MMP mapping comes "free" with binary search







- one for each chromosome
- MMP mapping comes "free" with binary search







- one for each chromosome
- MMP mapping comes "free" with binary search







- one for each chromosome
- MMP mapping comes "free" with binary search







- one for each chromosome
- MMP mapping comes "free" with binary search







- one for each chromosome
- MMP mapping comes "free" with binary search







- one for each chromosome
- MMP mapping comes "free" with binary search

needs to be performed using both the read and it's reverse complement











#### • Easily identifies "multi-mapped" reads since they are together in the SA

- Find mismatches internal to the read





Easily identifies "multi-mapped" reads since they are together in the SA

- Easily identifies "multi-mapped" reads since they are together in the SA Find mismatches internal to the read
- Overcomes issues with poly-A tails, library adaptors, and low quality



### Seed Search

- Easily identifies "multi-mapped" reads since they are together in the SA Find mismatches internal to the read
- Overcomes issues with poly-A tails, library adaptors, and low quality
- Has high speed, but large memory footprint compared to compressed SAs



### Seed Search





- A set of the MMPs are selected as "anchors"
  - "In the current implementation, all the alignments that map less than a user defined value (typically 20-50) are selected as anchors."

A set of the MMPs are selected as "anchors"

- "In the current implementation, all the alignments that map less than a user defined value (typically 20-50) are selected as anchors."
- Anchors are those that map to less than 50 locations in the genome

A set of the MMPs are selected as "anchors"

- "In the current implementation, all the alignments that map less than a user defined value (typically 20-50) are selected as anchors."
- Anchors are those that map to less than 50 locations in the genome
- Alignment "windows" are then defined as regions around anchors • all of the MMPs in those windows will be stitched together linearly the size of these windows determines the maximum intron size

They stitch two MMPs together allowing • one gap (of multiple bases) in the genome, and • minimal mismatches.

 $\Delta$  is the difference in the size of the space between the MMPs in the genome and the read this is how long the gap is going to be

Match counts the number of bases that are better aligned before the gap minus the count of those  $\Delta =: (g_2 - g_1) - (r_2 - r_1)$ that are better placed after



 $Match(r', \Delta) =: \begin{cases} 1 & \text{if } R(r_1 + r') = G(g_1 + r') \& R(r_1 + r' + \Delta) \neq G(g_1 + r' + \Delta) \\ -1 & \text{if } R(r_1 + r') = G(g_1 + r') \& R(r_1 + r' + \Delta) \neq G(g_1 + r' + \Delta) \\ 0 & otherwise \end{cases}$ 



Mate pairs are initially mapped independently as long as they are on the same strand

If the whole read is not covered in 1 window

- find two (or more) non-overlapping windows that map to the read
- create a chimeric mapping



Genomic window 2

window ng



Mate pairs are initially mapped independently as long as they are on the same strand

If the whole read is not covered in 1 window

- find two (or more) non-overlapping windows that map to the read
- create a chimeric mapping



The score of each mapped read is:

#### $S(\mathbb{A}) = mt_{\mathbb{A}} - ms_{\mathbb{A}} - in_{\mathbb{A}} - dl_{\mathbb{A}} - gp_{\mathbb{A}}$

The score of each mapped read is:

 $S(\mathbb{A}) = mt_{\mathbb{A}} - ms_{\mathbb{A}} - in_{\mathbb{A}} - dl_{\mathbb{A}} - gp_{\mathbb{A}}$ match and mismatch scores as we saw before

The score of each mapped read is:

 $S(\mathbb{A}) = mt_{\mathbb{A}} - ms_{\mathbb{A}} - in_{\mathbb{A}} - dl_{\mathbb{A}} - gp_{\mathbb{A}}$ match and mismatch scores affine gaps with different scores for the read and the as we saw before

genome

The score of each mapped read is:

 $S(\mathbb{A}) = mt_{\mathbb{A}} - ms_{\mathbb{A}} - in_{\mathbb{A}} - dl_{\mathbb{A}} - gp_{\mathbb{A}}$ match and mismatch scores affine gaps with different scores for the read and the as we saw before

genome

splice junction score based on the sequences at the end of the junctions: GT/AG, GC/AG, & AT/AC are normal, all others have higher penalties (this is the *P*<sub>gap</sub> score from earlier)



genome

The score of each mapped read is:

 $S(\mathbb{A}) = mt_{\mathbb{A}} - ms_{\mathbb{A}} - in_{\mathbb{A}} - dl_{\mathbb{A}} - gp_{\mathbb{A}}$ match and mismatch scores affine gaps with different scores for the read and the as we saw before

For each read, the mapping with the highest score is retained

splice junction score based on the sequences at the end of the junctions: GT/AG, GC/AG, & AT/AC are normal, all others have higher penalties (this is the *P*<sub>gap</sub> score from earlier)



### Speed vs. Space tradeoff of the SA

Table 1. Mapping speed and RAM benchmarks on the experimental RNA-seq dataset

Aligner	Mapping speed: million read pairs/hour		Peak physical RAM, GB	
	6 threads	12 threads	6 threads	12 threads
STAR	309.2	549.9	27.0	28.4
STAR sparse	227.6	423.1	15.6	16.0
TopHat2	8.0	10.1	4.1	11.3
RUM	5.1	7.6	26.9	53.8
MapSplice	3.0	3.1	3.3	3.3
GSNAP	1.8	2.8	25.9	27.0

## Alignment Quality



MapSplice

Fig. 2. True-positive rate versus false-positive rate (ROC-curve) for simulated RNA-seq data for STAR, TopHat2, GSNAP, RUM and

### Take Aways for STAR

Non-contiguous alignment for RNA-Seq is not a totally solved problem

Algorithm is extendable to longer read lengths since it can ignore poor quality regions and chimeric reads

- STAR is specifically designed to take introns into account during alignment

Large memory consumption, but fast due to the use of uncompressed SAs

#### Creates an alignment in stages:

- find all high quality whole read alignments to the genome
- identifying possible splice locations
- aligning initially unmapped reads to induced sequences

#### DINFORMATICS

#### Sequence analysis

#### **TopHat: discovering splice junctions with RNA-Seq**

Cole Trapnell<sup>1,\*</sup>, Lior Pachter<sup>2</sup> and Steven L. Salzberg<sup>1</sup>

<sup>1</sup>Center for Bioinformatics and Computational Biology, University of Maryland, College Park, MD 20742 and <sup>2</sup>Department of Mathematics, University of California, Berkeley, CA 94720, USA

Received on October 23, 2008; revised on February 24, 2009; accepted on February 26, 2009 Advance Access publication March 16, 2009

Associate Editor: Ivo Hofacker

### TopHat

**ORIGINAL PAPER** 

Vol. 25 no. 9 2009, pages 1105–1111 doi:10.1093/bioinformatics/btp120

# Aligning Reads using Bowtie

Using strict alignment critera, TopHat uses Bowtie to align reads to the whole genome

they include multi-mapped reads

The reads which don't map are called "Initially Unmapped" or IUM reads

- likely don't map because they cross introns
- will be used later to confirm splice junctions

"Low complexity" reads are discarded these are likely highly repetitive



Assemble consensus of covered regions

Ξ

Generate possible splices between neighboring exons

Map reads to possible splices via seed-andextend

qt ag ag



### **Construct potential exons**

Construct the set of mapped sequences

- the "islands" of sequence that map to the genome
- using the assemble functionality of MAQ

In low coverage regions replace any mismatches with the original genomic sequence

Append a small amount of the flanking sequence from the genome to the end of each cluster

 there may be some parts of the end of an exon that are only covered by spliced reads



Ì



- Assemble consensus of covered regions
- Generate possible splices between neighboring exons

Map reads to possible splices via seed-andextend



## Create potential splice locations

Splice junctions usually happen with predictable bases

- consider all possible pairs as potential splice locations
- create a set of new sequences
- store the k-mer surrounding such locations as a seed for mapping



qt ag ag

gt

ag ag

right exon

Map reads to possible splices via seed-andextend

Assemble

Generate possible

splices between

neighboring

exons

consensus of

Ξ



## Create potential splice locations

Creating all possible combinations may be too many • only find regions that are high coverage or at the ends of the potential exons

probability of including potential junction *i* to *j* 

brain RNA

 $\sum_{ij}^{n} = \frac{\sum_{m=1}^{j} d_m}{j-i} \cdot \frac{1}{\sum_{m=0}^{n}}$ coverage at location *m* 



## Index IUM reads

For each unmapped read

- extract all unique k-mers from the "high quality" region
- •here *k*~10



extend

# Align remaining reads

Find all matches between k-mers

- from the IUM set and
- the newly created potential junction locations

Extend left and right from the k-mer to find an alignment





extend



# Align remaining reads

Find all matches between k-mers

- from the IUM set and
- the newly created potential junction locations

Extend left and right from the k-mer to find an alignment

Will miss alignments with mismatches near splice location

Discard junctions that are predicted to be anomalies:

 occurrence rate is lower than some percentage of the coverage of the flanking regions



## Take Aways from TopHat

Uses existing software to do some of the heavy lifting

Strict parameters on the splice junctions make the algorithm fast

Limited in the splice junction sequence

### TopHat2

- Aligns reads in stages of increasing time requirements: • first to the transctiptome (set of all known transcripts),
  - then whole reads to the genome,
  - then in chunks

Kim et al. Genome Biology 2013, 14:R36 http://genomebiology.com/2013/14/4/R36

#### **METHOD**

#### TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions

Daehwan Kim<sup>1,2,3\*</sup>, Geo Pertea<sup>3</sup>, Cole Trapnell<sup>5,6</sup>, Harold Pimentel<sup>7</sup>, Ryan Kelley<sup>8</sup> and Steven L Salzberg<sup>3,4</sup>



**Open Access** 

#### Align to transcriptome

Aligned using Bowtie

The transcriptome is the set of all known transcripts

missing some novel transcripts

This step identifies some of the spliced reads before we go into the TopHat(1) pipeline



iptome index
ne index
nents nome.
ne index
tial splice sites apped positions ager than the ead.
ncatenated em.
flanking index
ne or flanking read alignments.
xtending a few
maniatedu.

#### Align unmapped reads to genome

Same as the first step of TopHat(1)



ptome index
ne.
e index
ents ome.
e index
tial splice sites pped positions ger than the ad
catenated
flanking index
e or flanking ead alignments.

#### Align chunks of unmapped reads

Additional step from TopHat(1)

Split read into subsequences

 align these subsequences using **Bowtie** 



#### Identify potential splice locations

Use these mapped reads to limit the number of potential splice locations

consider more pairs of ending
2-mers



	]
tome	Index
index	
nts	
me.	
index	
al splice	sites
al splice ped po: er than	sites sitions the
al splice ped po: er than d.	sites sitions the
al splice ped po: er than d.	sites sitions the
al splice ped po: er than d. atenate	sites sitions the
al splice ped po: er than d. atenate n. lanking	sites sitions the ed
al splice ped pos er than d. atenate n. anking	sites sitions the ed g index
al splice ped pos er than d. atenate n. anking or flanl ad align	sites sitions the ed g index king ments.
al splice ped pos er than d. atenate n. lanking or flanl ad align	sites sitions the ed g index king ments.
al splice ped pos er than d. atenate n. anking or flan ad align	sites sitions the ed g index king ments.
al splice ped pos er than d. atenate n. anking or flan ad align	sites sitions the ad g index king ments.
al splice ped pos er than d. atenate n. anking or flan ad align	sites sitions the ad g index king ments.

#### Align unaligned chunks

Creating a new index of these flanking regions

 align the unaligned chunks of all of the reads to find the location of the splices



#### **Recreate whole read** alignments

Using the flanking mapping or the original mapping

Stitch together to get the whole read mapping



tomei	ndex
index	
nts	

### Re-map reads

This step will help align small overlaps between the two adjoining exons

similar to the last step in TopHat(1)



### Take Aways for TopHat2

#### Built on Bowtie(2), so the actual mapping is very accurate

Annotation allows TopHat2 to better align reads in known transcripts