# BIOL 343
# Applied Bioinformatics I

## Alignment/Mapping

Dr. Nic Wheeler

# Learning Objectives

You will be able to:

1.

Bharti 2021, Next Generation Sequencing and Analysis

# Alignment is the most important step in RNA-seq analysis
## Counting (also important) and DEG ID relies on high-confidence mapping

- Recall the goal of our RNA-seq experiments…

  - Treatment vs Control

  - Mutant vs Wild type

  - ***Identify differentially expressed genes (DEGs)***

- DEGs will be identified using statistical tests comparing ***expression values*** of transcripts/genes

- Expression values will be calculated based on the number of reads that ***align/map*** to a given genomic locus
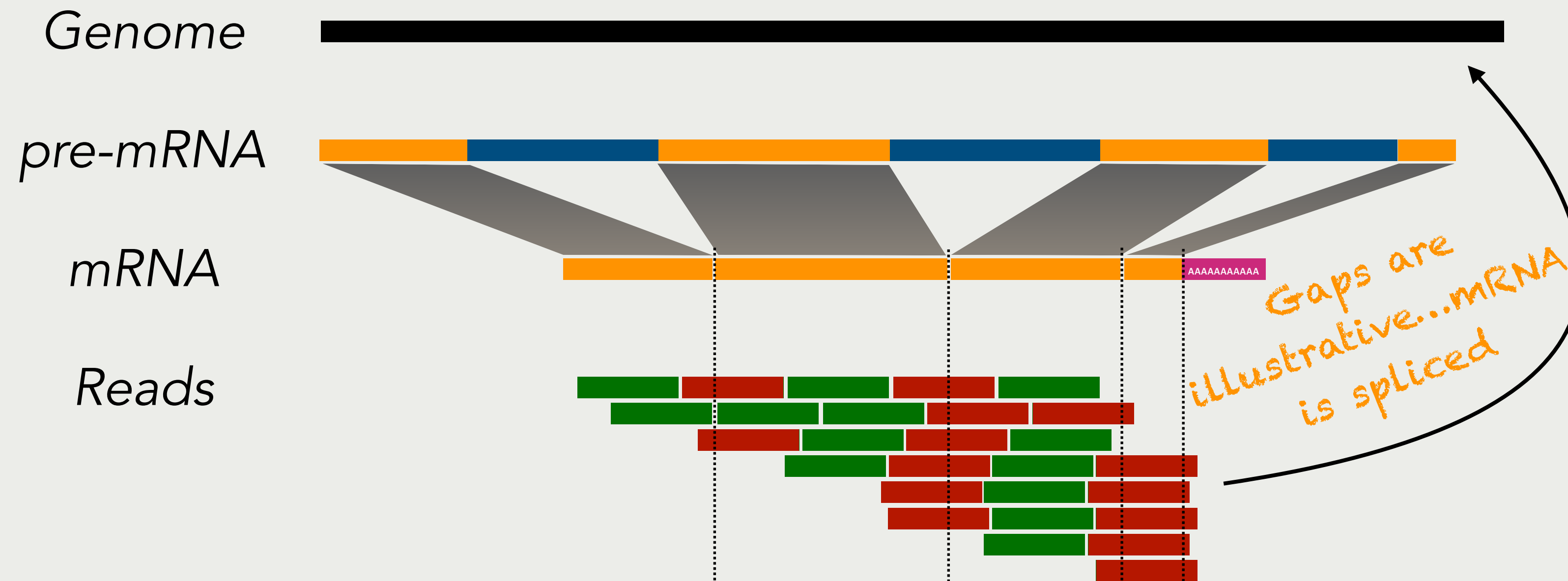
# Types of alignment algorithms
## Needleman-Wunsch and…

- Needleman-Wunsch (global alignment)

  - Dynamic programming

  - Mismatch penalty (transitions or transversions)

  - Gap penalty

- Problems - not global alignment, reference genomes are *huge* strings with lots of repetition, reads are likely to align many locations, and reads will align with massive gaps if spanning an intron

- Solution - Suffix array (STAR) *or* Burrow-Wheelers transform and FM-index (HISAT)
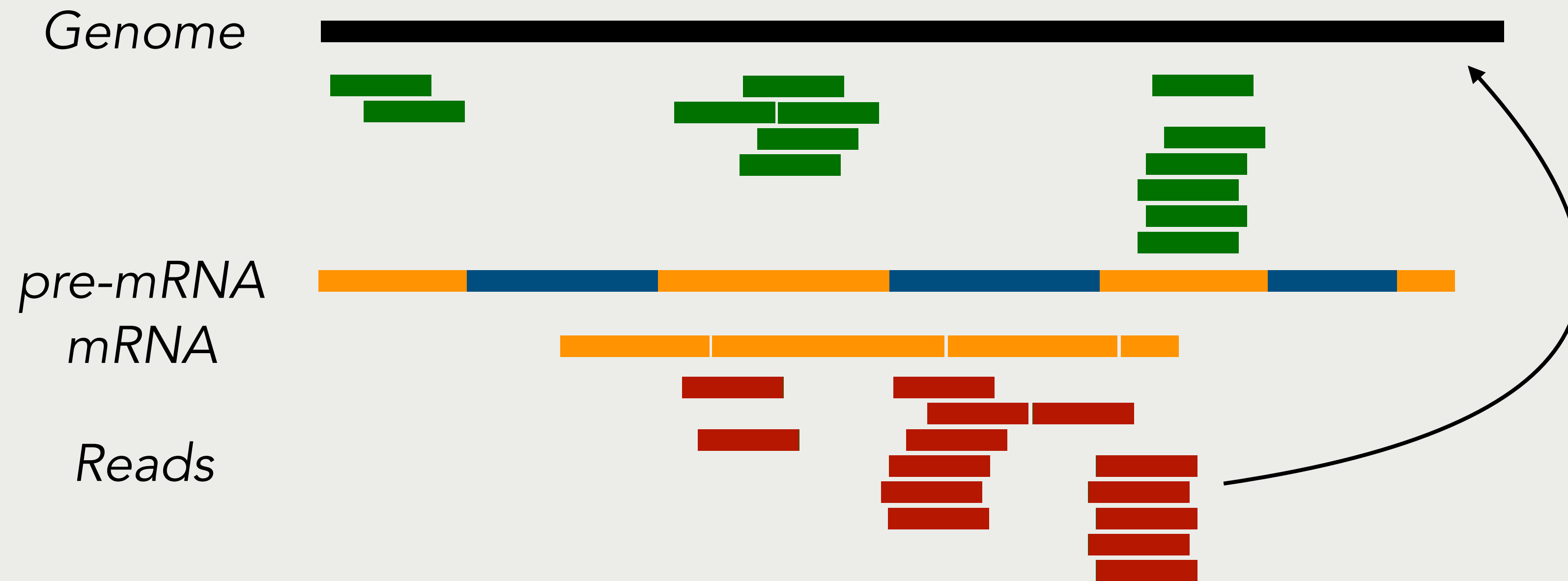
# Splice-aware alignment
## Gaps are large and encouraged

Genome

pre-mRNA

mRNA

AAAAAAAAAA

Gaps are illustrative...mRNA is spliced
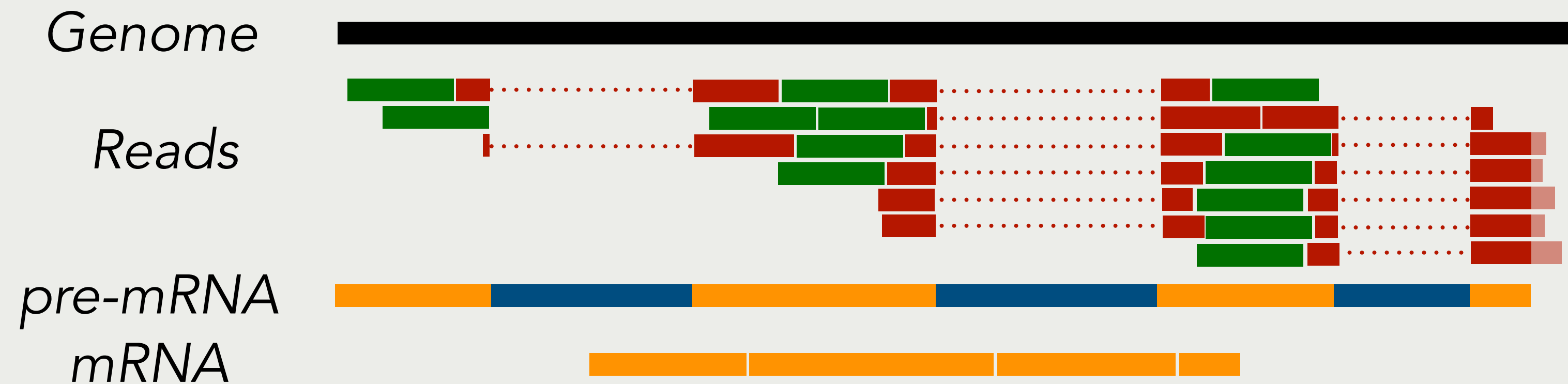
Reads

*But, reads aren't aligned to a transcriptome (mRNAs), but a genome*

# Splice-aware alignment
## Gaps are large and encouraged

Genome

pre-mRNA
mRNA

Reads

*But, reads aren't aligned to a transcriptome (mRNAs), but a genome*

# Two (main) approaches to splice-aware alignment
## STAR and HISAT

**STAR**

Spliced Transcripts Alignment to a Reference

Published in 2013

40574 citations

Requires a lot of RAM; ultra fast

**HISAT**

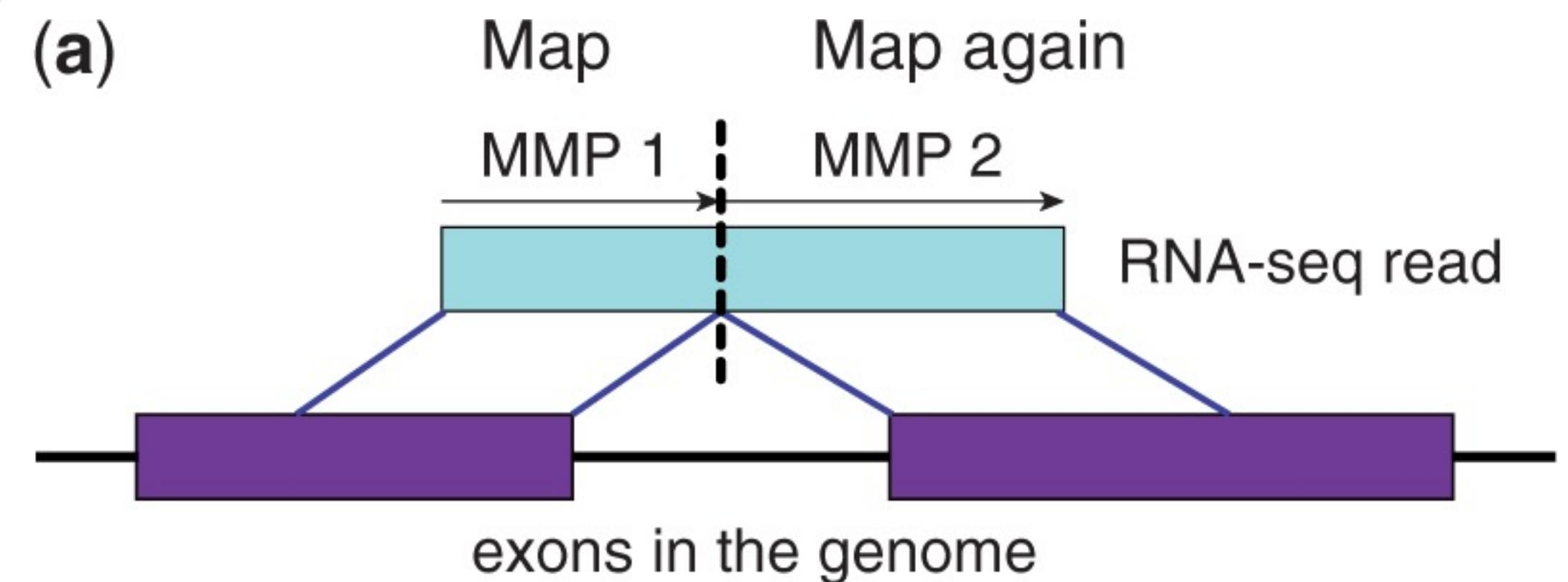Hierarchical Indexing for Spliced Alignment of Transcripts

Published in 2015

17667 citations

Less RAM needed; still fast

# Two (main) approaches to splice-aware alignment
## STAR

1. Find the Maximal Mappable Prefix of the read

   - MMP 1 will map to a splice donor

2. Find the MMP of the remainder of the read

   - MMP 2 will map to a splice acceptor

- Uses a suffix array of the reference genome



*Dobin 2013, Bioinformatics*

# Two (main) approaches to splice-aware alignment
## STAR

- Uses a suffix array of the reference genome

  - Every substring of the genome sorted lexicographically

  - Given a search string *P*, two binary searches to find the boundaries

    - *gtg* - binary search to find boundary 1 at index 5, binary search to find boundary 2 at index 9

- Many developments (ongoing) in 1) generating the SA and 2) searching the SA

# Two (main) approaches to splice-aware alignment
## STAR

- Only two binary searches - ultrafast!

- Suffix array of a large genome - very big RAM!

- Generating suffix array - kinda slow!



*Shrestha 2014, Briefings in Bioinformatics*

# Two (main) approaches to splice-aware alignment
## STAR

- Other advantages to this approach:

  - Robust to mismatches - MMPs can be extended

  - Can trim (**soft clip**) if extension of MMP results in many mismatches



(b) Map MMP 1 → Extend
mismatches

(c) Map MMP 1 → Trim
A-tail, or adapter, or poor quality tail