

Near-optimal probabilistic RNA-seq quantification

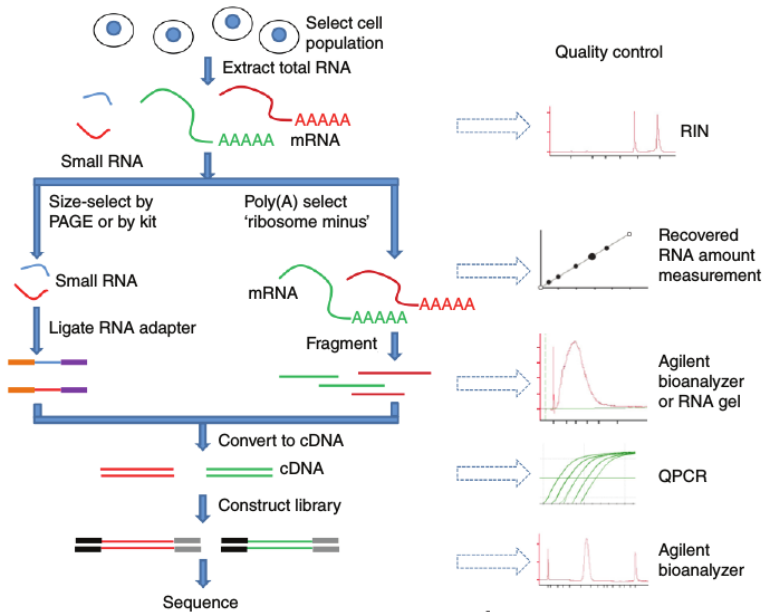
Bray and Pachter *et al.* Nature biotechnology(2016)
doi:10.1038/nbt.3519

Saket Choudhary

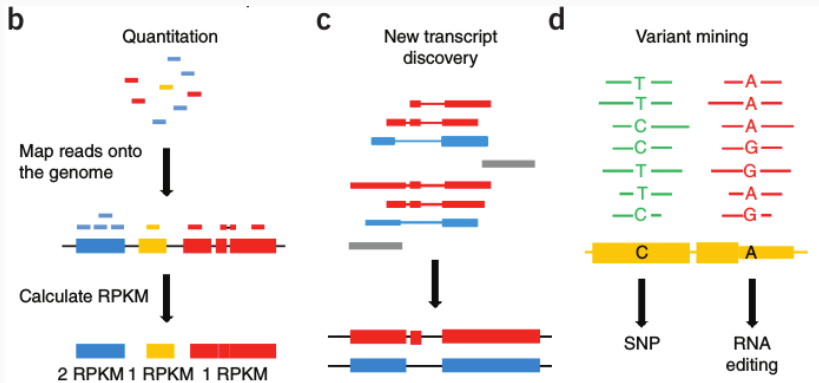
September 25, 2016

RNA-Seq Workflow

a



RNA-Seq Workflow



Zheng and Mortazavi(2012)

- First two steps in typical RNA-Seq processing pipeline:

- First two steps in typical RNA-Seq processing pipeline:
 - Alignment

- First two steps in typical RNA-Seq processing pipeline:
 - Alignment
 - Quantification

- First two steps in typical RNA-Seq processing pipeline:
 - Alignment
 - Quantification
- Alignments are slow and probably not so important

It's all about compatible transcripts

- Circumvent alignment step – Use information from k – mers
- Pseudoalignment: Find *compatible* transcripts for a read, without pinpointing where exactly it aligns



Figure 1: Reads and overlapping transcripts

Method II

sequence

ATGGAAGTCGCGGAATC

7mers

ATGGAAG
TGAAGT
GAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

de Bruijn graph

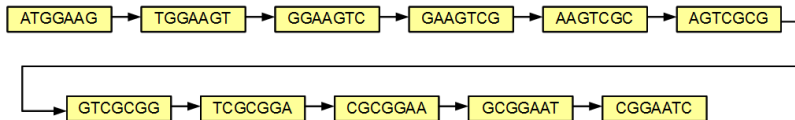


Figure 2: de Bruijn Graph

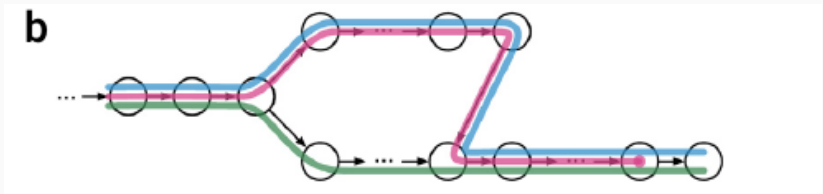


Figure 3: Transcriptome - de Bruijn Graph. Node = k — mers, Path = Transcript

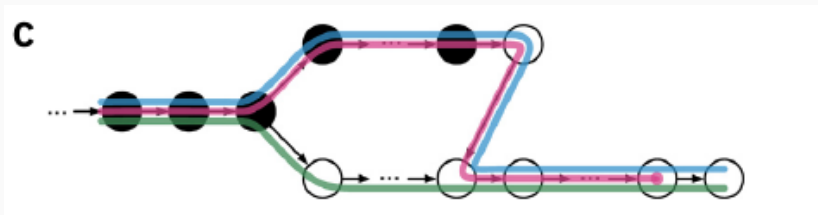


Figure 4: k – mers in read = black nodes

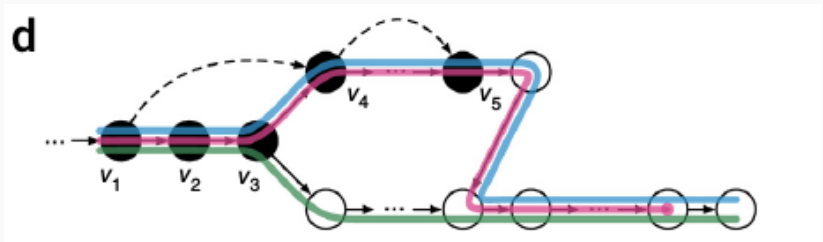


Figure 5: Nodes can be skipped if $k - mers$ did arise from blue transcript

e

$$\begin{array}{c} \text{blue} \\ \text{pink} \\ \text{green} \end{array} v_1 \cap \begin{array}{c} \text{blue} \\ \text{pink} \end{array} v_4 \cap \begin{array}{c} \text{blue} \\ \text{pink} \end{array} v_5 = \begin{array}{c} \text{blue} \\ \text{pink} \end{array}$$

Figure 6: Intersection of k-compatibility class

Method VII

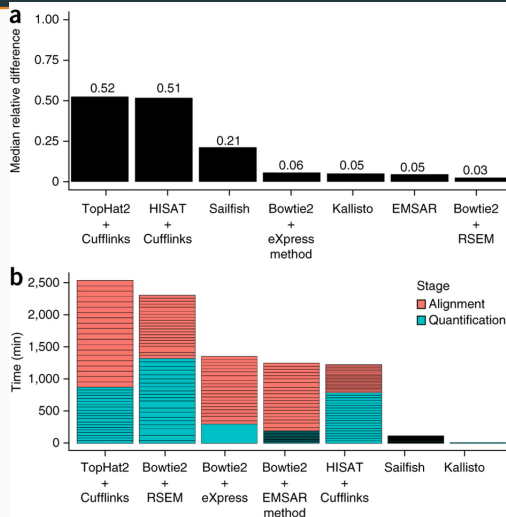


Figure 7: Quantification over a cup of coffee

Questions?

- Better than *Sailfish* that looks up k – mers in reads into k – mers of transcriptome

notes

- Better than *Sailfish* that looks up $k - mers$ in reads into $k - mers$ of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns

notes

- Better than *Sailfish* that looks up $k - mers$ in reads into $k - mers$ of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns
- Key Idea: Find compatible transcript for a read, not where it exactly aligns

notes

- Better than *Sailfish* that looks up $k - mers$ in reads into $k - mers$ of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns
- Key Idea: Find compatible transcript for a read, not where it exactly aligns
- Pseudoalignment: Subset $S \subset T$ such that read r is compatible.

notes

- Better than *Sailfish* that looks up k – mers in reads into k – mers of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns
- Key Idea: Find compatible transcript for a read, not where it exactly aligns
- Pseudoalignment: Subset $S \subset T$ such that read r is compatible.
- Hash k -mers of reads and have a de-bruijn graph of transcriptome assembly handy

notes

- Better than *Sailfish* that looks up $k - mers$ in reads into $k - mers$ of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns
- Key Idea: Find compatible transcript for a read, not where it exactly aligns
- Pseudoalignment: Subset $S \subset T$ such that read r is compatible.
- Hash k -mers of reads and have a de-bruijn graph of transcriptome assembly handy
- T-DBG: nodes are k -mers , each transcript corresponds to a path and path cover induces a k -compatibility class for each k -mer

notes

- Better than *Sailfish* that looks up $k - mers$ in reads into $k - mers$ of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns
- Key Idea: Find compatible transcript for a read, not where it exactly aligns
- Pseudoalignment: Subset $S \subset T$ such that read r is compatible.
- Hash k -mers of reads and have a de-bruijn graph of transcriptome assembly handy
- T-DBG: nodes are k -mers, each transcript corresponds to a path and path cover induces a k -compatibility class for each k -mer
- T-DBG: Colors correspond to transcripts, node corresponds to k -mers, every k -mer receives a color for each transcript it occurs in

notes

- Better than *Sailfish* that looks up $k - mers$ in reads into $k - mers$ of transcriptome
- Pseudoalignment: Find compatible transcript for a read, not where it exactly aligns
- Key Idea: Find compatible transcript for a read, not where it exactly aligns
- Pseudoalignment: Subset $S \subset T$ such that read r is compatible.
- Hash k -mers of reads and have a de-bruijn graph of transcriptome assembly handy
- T-DBG: nodes are k -mers, each transcript corresponds to a path and path cover induces a k -compatibility class for each k -mer
- T-DBG: Colors correspond to transcripts, node corresponds to k -mers, every k -mer receives a color for each transcript it occurs in
- Hash table stores mapping of each k -mer to the contig it is contained in