

# GENOME 569

Class 3: NGS read alignment

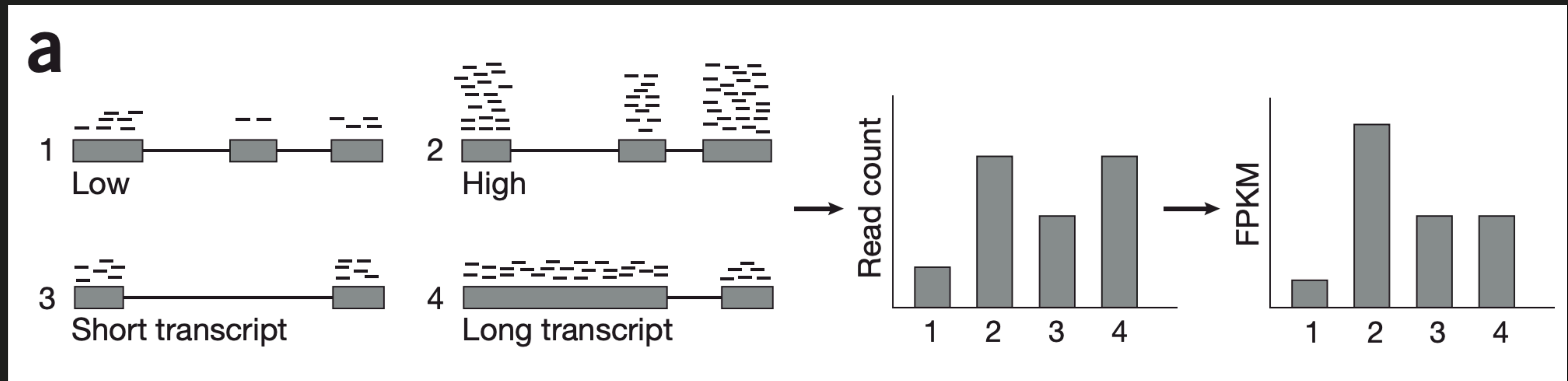
# Discussion about Phase I

Any problems?

What packages did CoPilot use?

What did you have to do “the old fashioned way?”

# Measuring gene expression with NGS



The number of reads from a transcript is proportional to its abundance.

With random RT primers, you also need to correct for transcript length.

# How to map billions of short reads onto genomes

Cole Trapnell & Steven L Salzberg

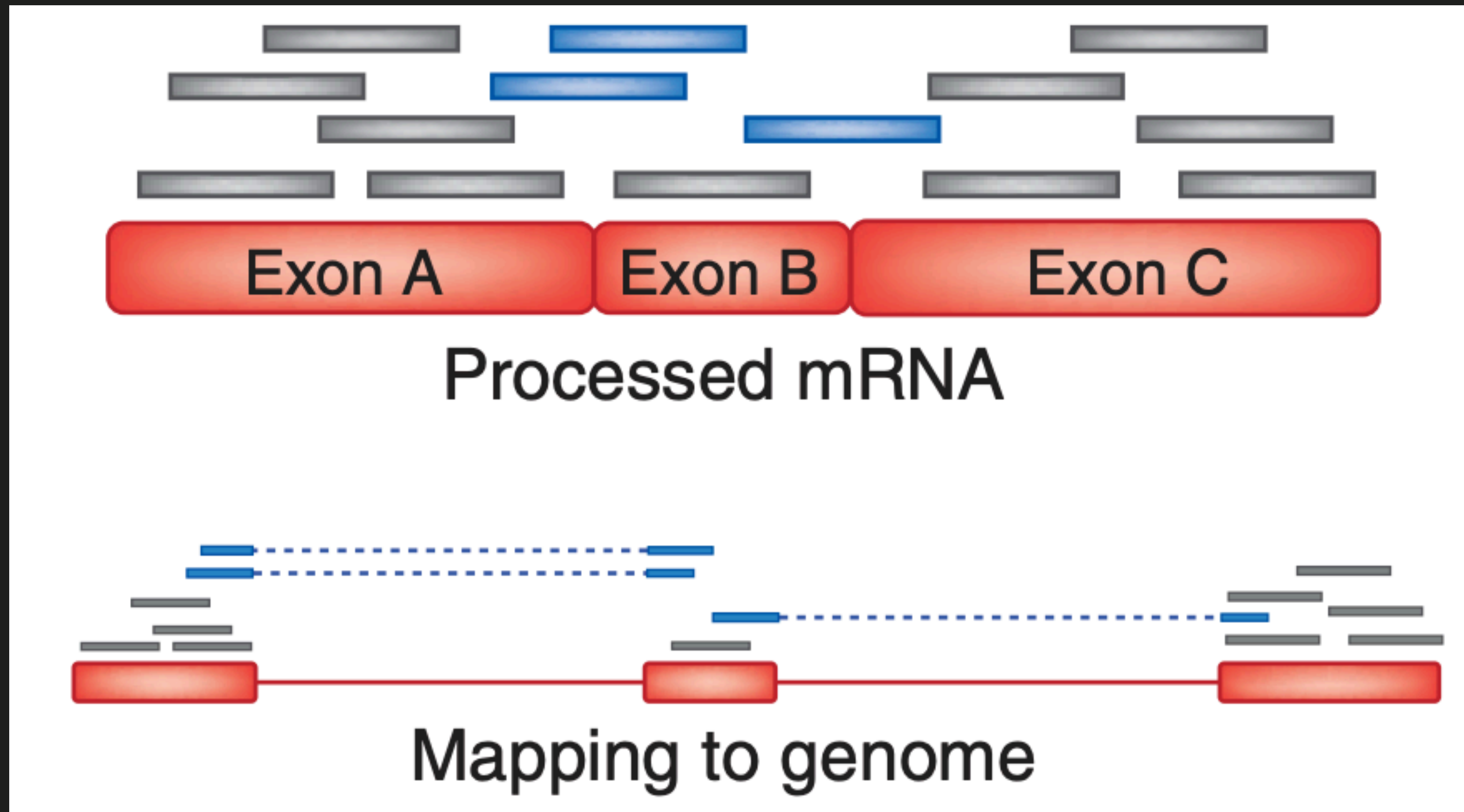
Mapping the vast quantities of short sequence fragments produced by next-generation sequencing platforms is a challenge. What programs are available and how do they work?

A new generation of DNA sequencers that can rapidly and inexpensively sequence billions of bases is transforming genomic science. These new machines are quickly becoming the technology of choice for whole-genome sequencing and for a variety of sequencing-based assays, including gene expression, DNA-protein interaction, human resequencing and RNA splicing studies<sup>1-3</sup>. For example, the RNA-Seq protocol, in which processed mRNA is converted to cDNA and sequenced, is enabling the identification of previously unknown genes and alter-

**Table 1 A selection of short-read analysis software**

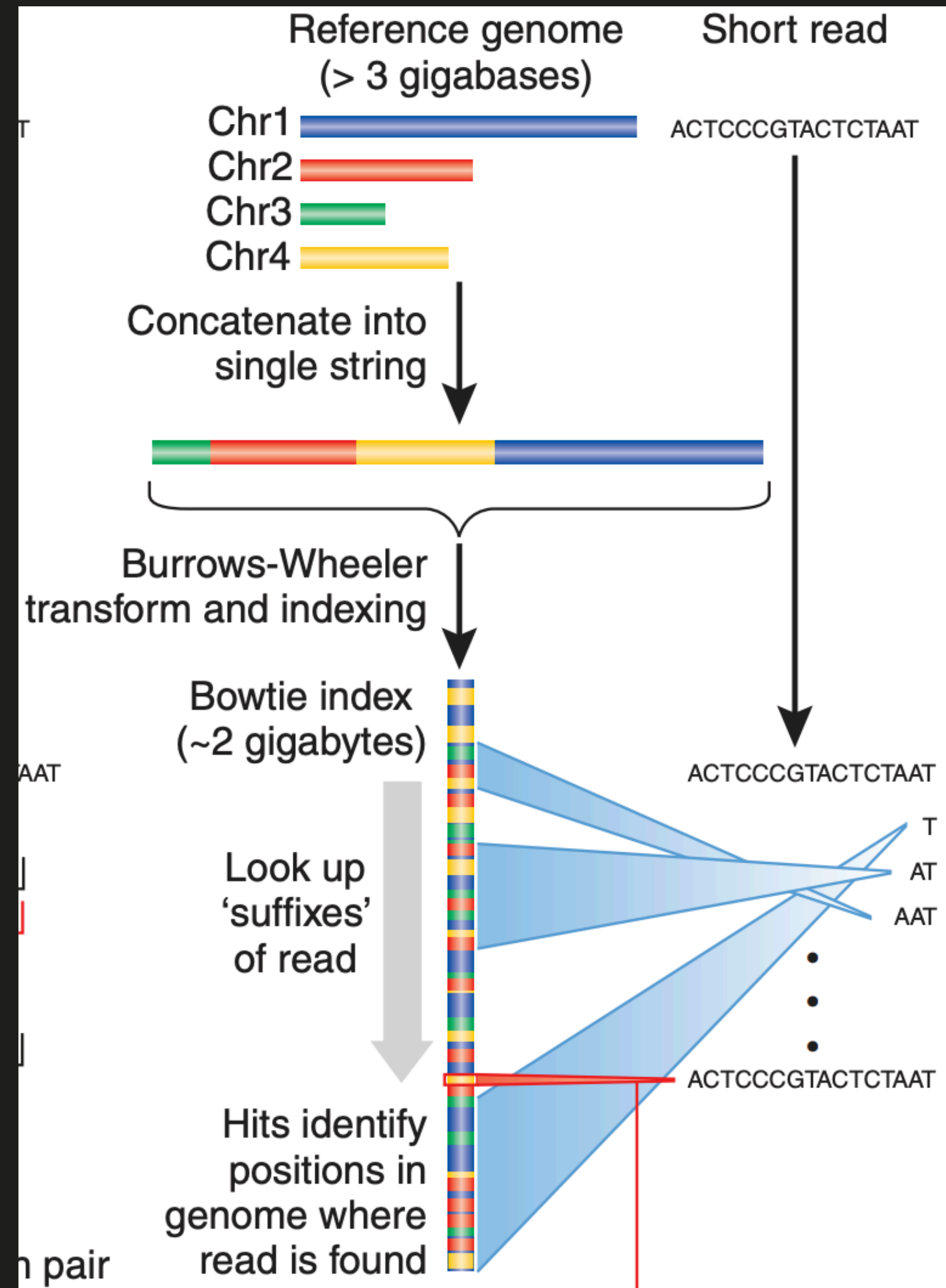
Program	Website	Open source?	Handles ABI color space?	Maximum read length
Bowtie	<a href="http://bowtie.cbcb.umd.edu">http://bowtie.cbcb.umd.edu</a>	Yes	No	None
BWA	<a href="http://maq.sourceforge.net/bwa-man.shtml">http://maq.sourceforge.net/bwa-man.shtml</a>	Yes	Yes	None
Maq	<a href="http://maq.sourceforge.net">http://maq.sourceforge.net</a>	Yes	Yes	127
Mosaik	<a href="http://bioinformatics.bc.edu/marthlab/Mosaik">http://bioinformatics.bc.edu/marthlab/Mosaik</a>	No	Yes	None
Novoalign	<a href="http://www.novocraft.com">http://www.novocraft.com</a>	No	No	None
SOAP2	<a href="http://soap.genomics.org.cn">http://soap.genomics.org.cn</a>	No	No	60
ZOOM	<a href="http://www.bioinfor.com">http://www.bioinfor.com</a>	No	Yes	240

# Aligning RNA-seq reads





# Aligning RNA-seq reads



# Software for mapping NGS reads

Tool	DNA/RNA	Mapping strategy
bwa	DNA	Burroughs Wheeler
Bowtie	DNA	Burroughs Wheeler
TopHat	RNA	Burroughs Wheeler + linkage
HISAT	RNA	Burroughs Wheeler
STAR	RNA	Suffix arrays

Many others (some highly specialized), but these are the most popular

**STAR**



# Problem: map NGS reads to a genome

## Solution:

```
STAR \
  --runThreadN 8 \
  --genomeDir $INDEX \
  --genomeLoad LoadAndKeep \
  --readFilesIn $INPUT/$FILE \
  --readFilesCommand zcat \
  --outFileNamePrefix $OUTPUT/$SAMPLE. \
  --outSAMtype BAM Unsorted \
  --outSAMstrandField intronMotif
```

# Sample problem 1

```
$ module load STAR/latest
```

Download “small\_reads.fastq”

Use STAR to map these reads to the worm genome:

```
wget https://ctrappnell.github.io/genome569/example\_files/small\_reads.fastq
```

```
STAR --genomeDir ~coletrap/teaching/genome569/reference \  
--readFilesIn small_reads.fastq \  
--outFileNamePrefix example/ \  
--outSAMtype BAM Unsorted \  
--outSAMstrandField intronMotif
```

**SAM**

# Problem: store alignments in a standard format

## Solution:

### Sequence Alignment/Map Format Specification

The SAM/BAM Format Specification Working Group

5 Feb 2020

The master version of this document can be found at <https://github.com/samtools/hts-specs>.  
This printing is version dfc3e48 from that repository, last modified on the date shown above.

## 1 The SAM Format Specification

SAM stands for Sequence Alignment/Map format. It is a TAB-delimited text format consisting of a header section, which is optional, and an alignment section. If present, the header must be prior to the alignments. Header lines start with '@', while alignment lines do not. Each alignment line has 11 mandatory fields for essential alignment information such as mapping position, and variable number of optional fields for flexible or aligner specific information.

# Example alignment

```
Coord      12345678901234  5678901234567890123456789012345
ref        AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1          TTAGATAAAGGATA*CTG
+r002           aaaAGATAA*GGATA
+r003           gcctaAGCTAA
+r004                        ATAGCT.....TCAGC
-r003                        ttagctTAGGC
-r001/2                                CAGCGGCAT
```

@HD VN:1.6 SO:coordinate

@SQ SN:ref LN:45

```
r001    99 ref   7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002     0 ref   9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003     0 ref   9 30 5S6M          * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004     0 ref  16 30 6M14N5M       * 0 0 ATAGCTTCAGC *
r003  2064 ref  29 17 6H5M          * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001   147 ref  37 30 9M            = 7 -39 CAGCGGCAT * NM:i:1
```



# Example alignment

```
Coord      12345678901234  5678901234567890123456789012345
ref        AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1          TTAGATAAAGGATA*CTG
+r002            aaaAGATAA*GGATA
+r003            gcctaAGCTAA
+r004                        ATAGCT.....TCAGC
-r003                        ttagctTAGGC
-r001/2                                CAGCGGCAT
```

Read name	Where the Read maps			What's different	The original Read sequence			Additional metadata
@HD VN:1.6 SO:coordinate								
@SQ SN:ref LN:45								
r001	99	ref 7 30	8M2I4M1D3M	= 37 39	TTAGATAAAGGATACTG	*	SA:Z:ref,29,-,6H5M,17,0;	
r002	0	ref 9 30	3S6M1P1I4M	* 0 0	AAAAGATAAGGATA	*		
r003	0	ref 9 30	5S6M	* 0 0	GCCTAAGCTAA	*		
r004	0	ref 16 30	6M14N5M	* 0 0	ATAGCTTCAGC	*		
r003	2064	ref 29 17	6H5M	* 0 0	TAGGC	*	SA:Z:ref,9,+,5S6M,30,1;	
r001	147	ref 37 30	9M	= 7 -39	CAGCGGCAT	*	NM:i:1	



# Key features of SAM

Widely adopted. Nearly every read aligner uses it, many analysis tools accept it as input

“Lossless” - SAM files include all the information in the raw reads (even those that don’t map to the genome)

Can be stored in a binary format and heavily compressed

Can be indexed for fast lookup. You can easily extract all the alignments for a specific locus.

**samtools**

**Problem:** extract alignments in a given locus

**Solution:** `samtools view input.bam <region>`

Regions look like this:

`chr1:19200776-19220776`







```
samtools tview input.bam
```

[illegible]

**BED**



# BED

Browser Extensible Data format introduced by UCSC genome browser

Widely used to annotate genomes with intervals of interest

Simple, tab-delimited text file

Not very extensible, so used for very simple features (e.g. enhancers)

FYI gene models typically stored in GFF or GTF format, which is much more complex.

Chromosome	Start		Stop	Feature name		“Score”	Strand
I	15062083		15063836	WBGene00004512		255	+
I	15064301		15064453	WBGene00004567		255	+
I	15064838		15068346	WBGene00004622		255	+
I	15069280		15071033	WBGene00004513		255	+
MtDNA	898	1593	WBGene00014454	255	+		
MtDNA	10403	11354	WBGene00014472	255	+		
V	17115903		17116021	WBGene00077465		255	+
V	17117863		17117981	WBGene00077466		255	+
V	17118837		17118935	WBGene00077467		255	+

**bedtools**

**Problem:** compute overlap between BED files

**Solution:** `bedtools intersect -a reads.bed -b genes.bed`

This command computes the number of bases in file “B” that are covered by intervals in file “A”.

# Bedtools has many utilities

Command	Function
bamtobed	Convert a BAM file to a BED file
closest	For each interval in one file, find the closest in another
overlap	Compute the overlap between intervals
merge	Merge intervals that overlap
subtract	Remove the overlapping regions from intervals

And many, many, many more functions.  
Many with multiple modes of operation.

# Sample problem 2

Download “genes.bed”

Use bedtools to count the number of reads you just mapped that hit each gene.

```
samtools sort example/Aligned.out.bam > example/sorted.bam
```

```
bedtools intersect -a genes.bed -b example/sorted.bam -wa -c
```



# Sample problem 2

I	11947001	11953126	WBGene00011060	255	+	1
I	11953512	11961984	WBGene00002004	255	−	15
I	11971179	11971797	WBGene00044805	255	−	0
I	11979401	11985612	WBGene00013135	255	+	1

Note that the output is a BED file! The “score” field I told you to ignore earlier has the results.

# Droplet-based Single-cell RNA-seq (10X)

## ARTICLE

Received 20 Sep 2016 | Accepted 23 Nov 2016 | Published 16 Jan 2017

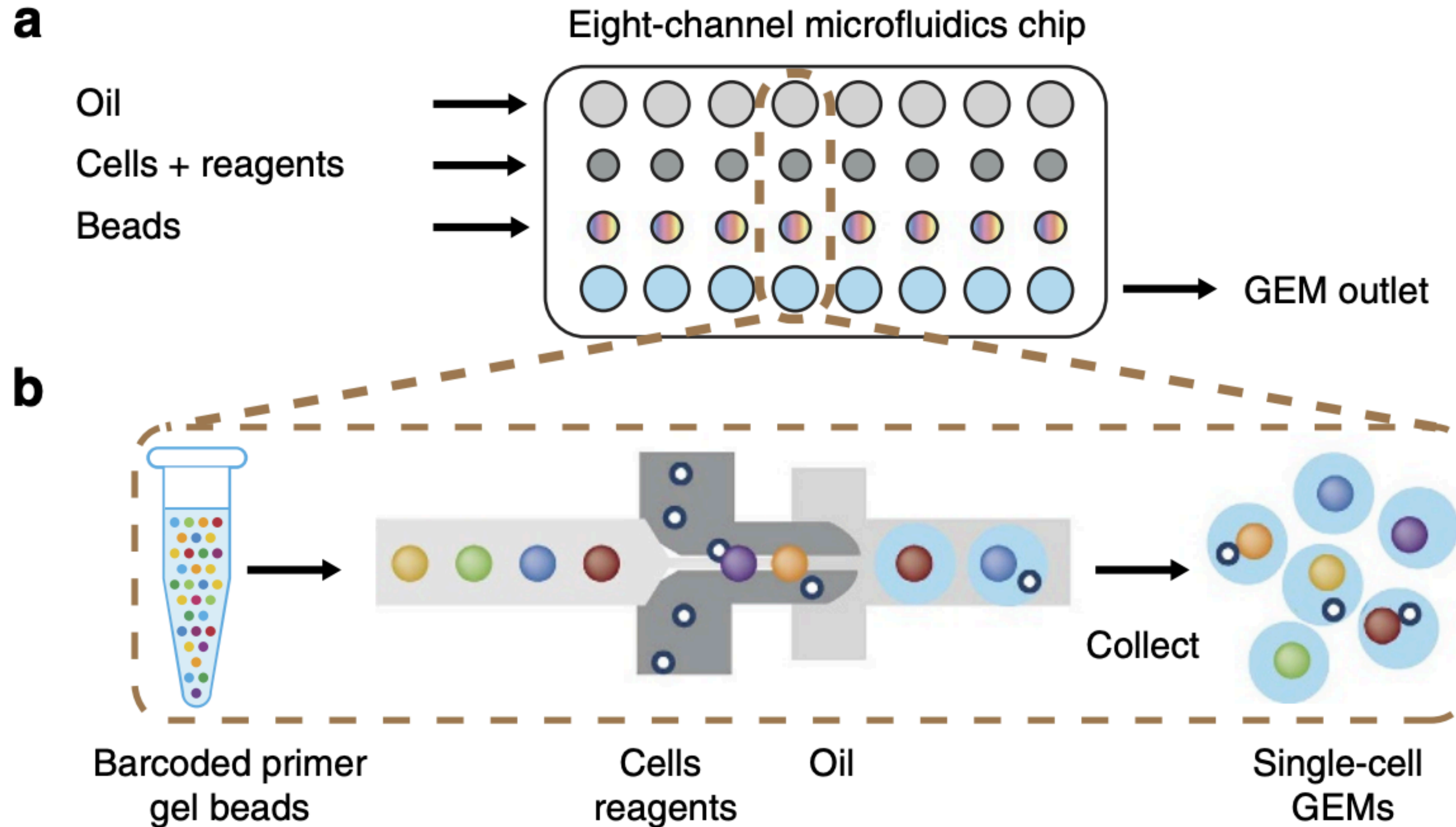
DOI: [10.1038/ncomms14049](https://doi.org/10.1038/ncomms14049)

OPEN

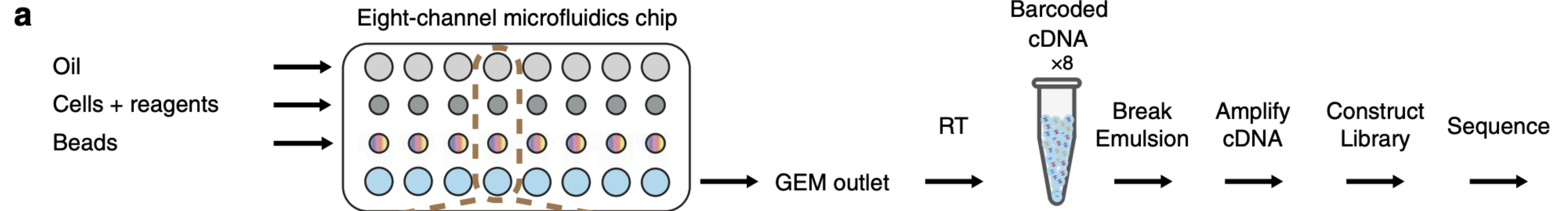
## Massively parallel digital transcriptional profiling of single cells

Grace X.Y. Zheng<sup>1</sup>, Jessica M. Terry<sup>1</sup>, Phillip Belgrader<sup>1</sup>, Paul Ryvkin<sup>1</sup>, Zachary W. Bent<sup>1</sup>, Ryan Wilson<sup>1</sup>, Solongo B. Ziraldo<sup>1</sup>, Tobias D. Wheeler<sup>1</sup>, Geoff P. McDermott<sup>1</sup>, Junjie Zhu<sup>1</sup>, Mark T. Gregory<sup>2</sup>, Joe Shugart<sup>1</sup>, Luz Montesclaros<sup>1</sup>, Jason G. Underwood<sup>1,3</sup>, Donald A. Masquelier<sup>1</sup>, Stefanie Y. Nishimura<sup>1</sup>, Michael Schnall-Levin<sup>1</sup>, Paul W. Wyatt<sup>1</sup>, Christopher M. Hindson<sup>1</sup>, Rajiv Bharadwaj<sup>1</sup>, Alexander Wong<sup>1</sup>, Kevin D. Ness<sup>1</sup>, Lan W. Beppu<sup>4</sup>, H. Joachim Deeg<sup>4</sup>, Christopher McFarland<sup>5</sup>, Keith R. Loeb<sup>4,6</sup>, William J. Valente<sup>2,7,8</sup>, Nolan G. Ericson<sup>2</sup>, Emily A. Stevens<sup>4</sup>, Jerald P. Radich<sup>4</sup>, Tarjei S. Mikkelsen<sup>1</sup>, Benjamin J. Hindson<sup>1</sup> & Jason H. Bielas<sup>2,6,8,9</sup>

# Droplet-based Single-cell RNA-seq (10X)



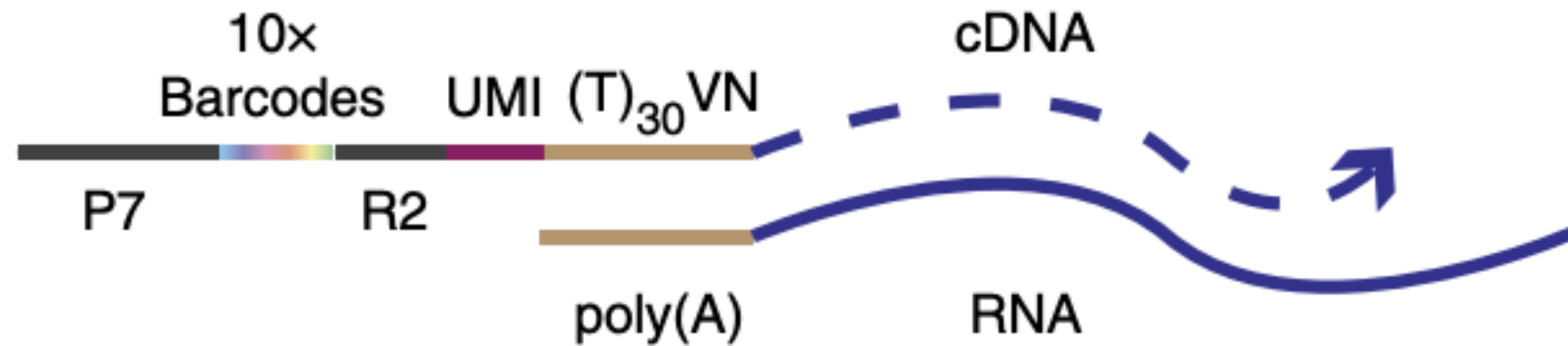
# Droplet-based Single-cell RNA-seq (10X)



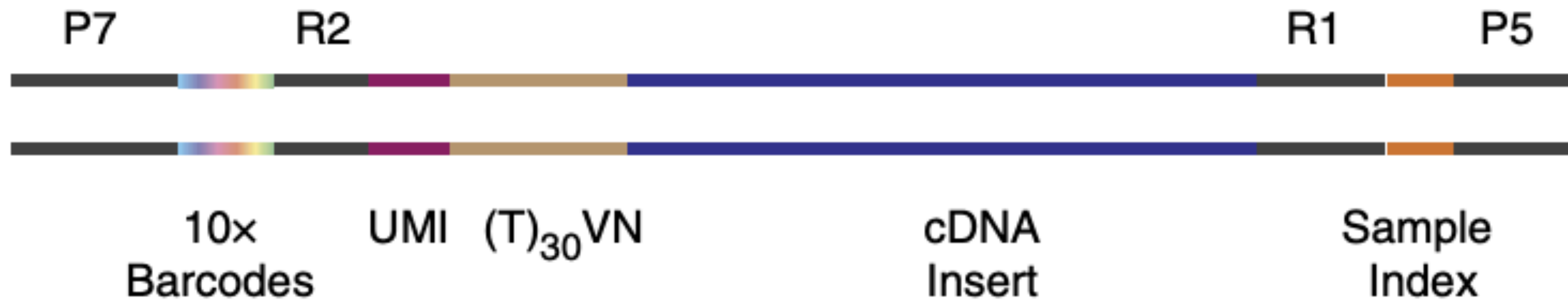


# 10X RNA-seq read structure (v1)

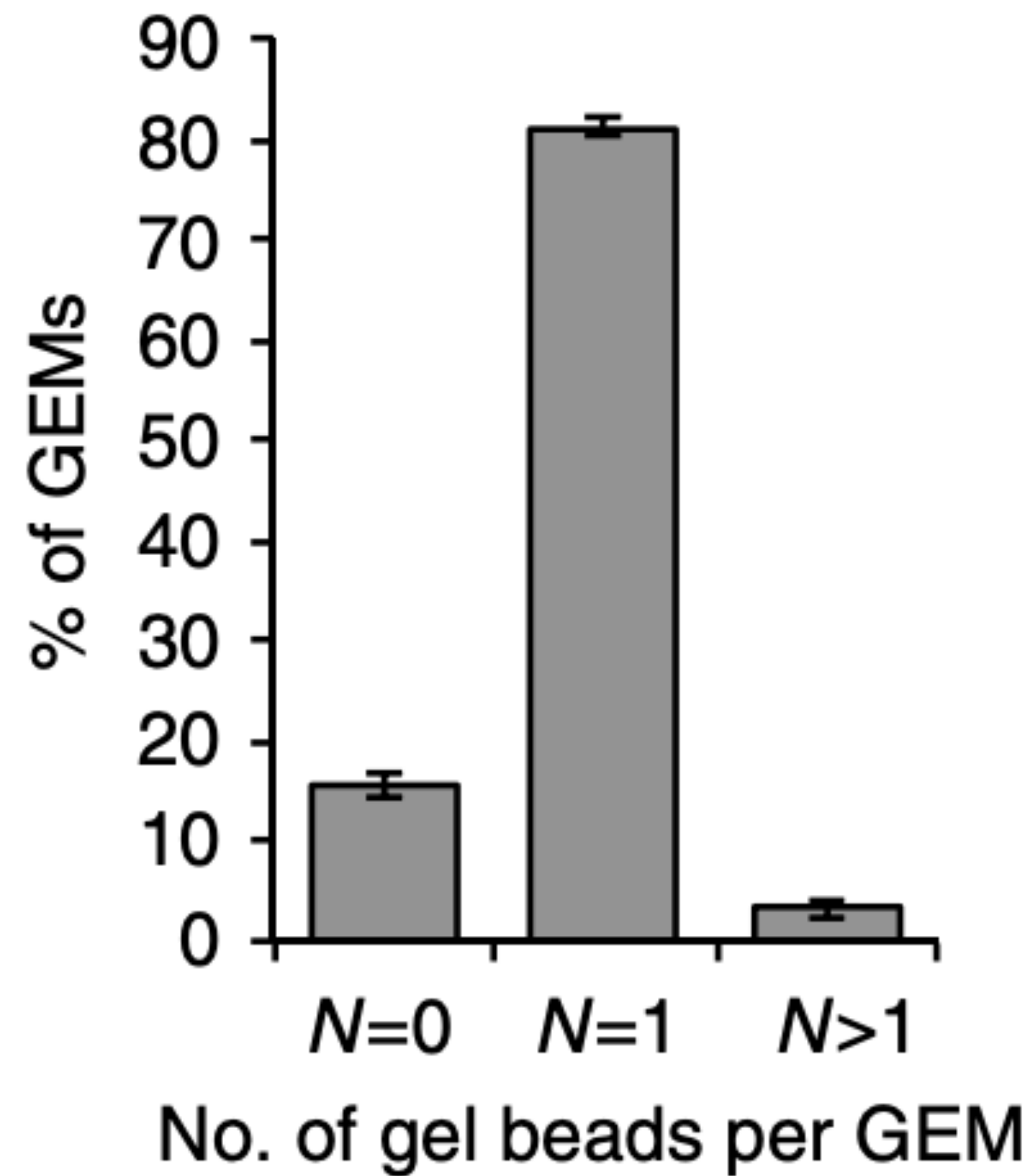
**d**



**e**

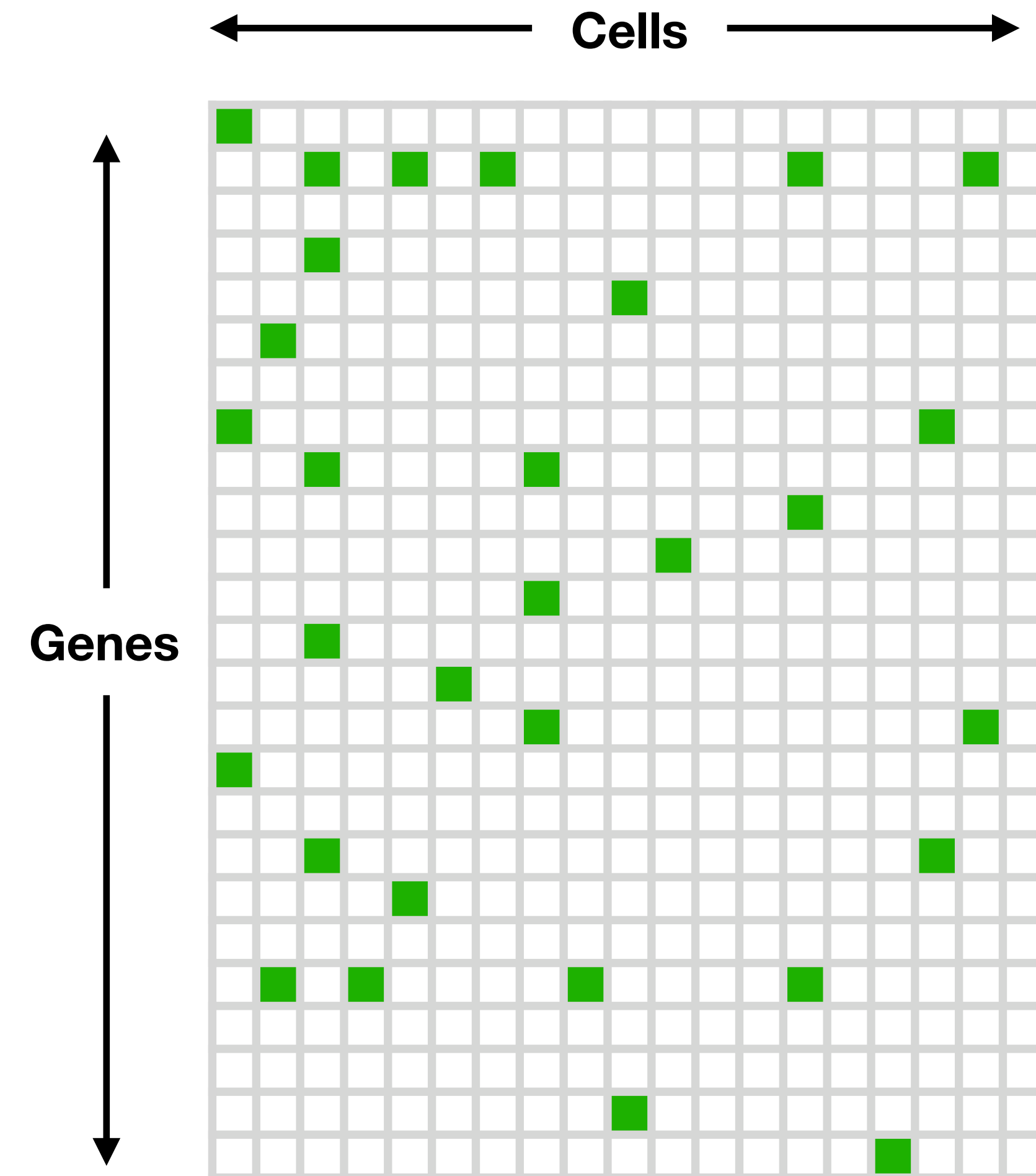
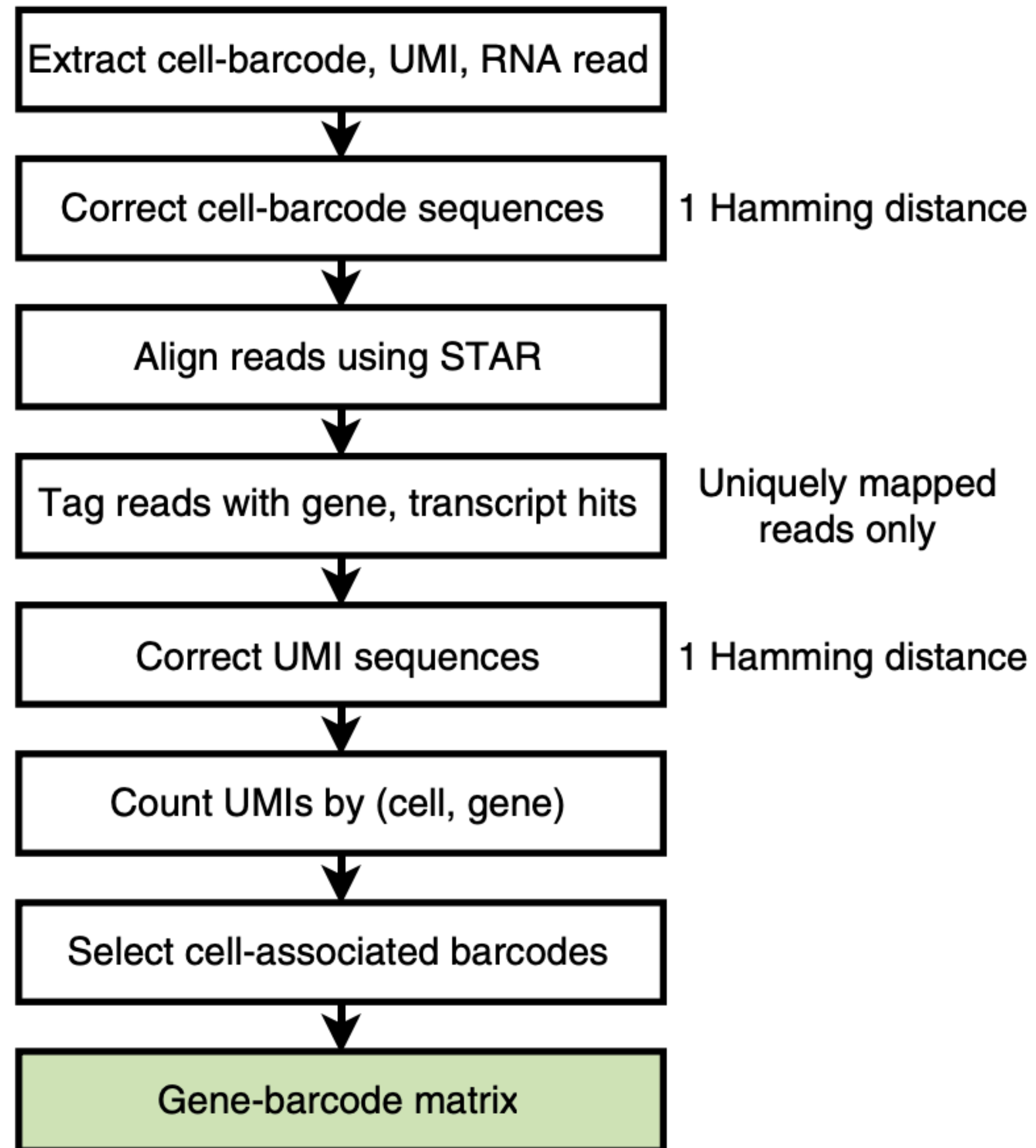


# Droplets contain ~1 cell (“Poisson loading”)





# The 10X bioinformatics workflow



**Most genes detected in few cells - matrix is very sparse!**



bioRxiv

THE PREPRINT SERVER FOR BIOLOGY

HOME | SUBMIT

New Results

 [Follow this preprint](#)

## **STARsolo: accurate, fast and versatile mapping/quantification of single-cell and single-nucleus RNA-seq data**

 Benjamin Kaminow,  Dinar Yunusov,  Alexander Dobin

**doi:** <https://doi.org/10.1101/2021.05.05.442755>

This article is a preprint and has not been certified by peer review [what does this mean?].

# Sample problem 3

Run STARsolo (with the proper arguments for 10X v1 chemistry) on one of the Packer et al 2019 samples.

Explore the output! How many of your reads mapped to the worm genome? How many cells did STARsolo report output for?

# THE END

For now