

# Advanced topics in single-cell transcriptomics

**Avi Srivastava**

New York Genome Center

May 27, 2020

# The Molecular Microscope of the (early) 21<sup>st</sup> Century

Sample collection & prep



## Sequencing Output Per Flow Cell

Flow Cell Type	S4
2 × 100 bp	1600–2000 Gb <sup>†</sup>
2 × 150 bp	2400–3000 Gb

4.8-6 Tbases / run

## "Analysis"

```
@SRR1215997.1 HWI-ST1311:58:C132FACXX:3:1101:3092:2249/1  
GAAGATGAGTGCATTGGAGGCTCTGGTATGGAATGATAAAAGTGAAGAATCAGCTCCGCTTGCAGAACTGGCCTATGATCTGGATGTGGATGATG  
+  
7<@DD;DAHDFB?G4A?BAFHFGFEH<C<E<?CFEFBGCH<DDBE<DGCF>?F<D*>FFFHIGH@GGIGII9)=;=?DED>@C>C3@CC>@C:3:3:5@  
@SRR1215997.2 HWI-ST1311:58:C132FACXX:3:1101:3435:2101/1  
CTTNTGACGCACTCCTCTAATTCGCCATATCTGTCTCATCATCCCAAGGTTCACATCTAGTAAGATGGAAGACTGGCAACAAGTGCAGGTTTTGG  
+  
?@@#4ADACFHFBGHIEHIIHHGIIJIIGIGGBGIGE?BBBFGGII./8C8)/@CHGIHHGEHIDEHGCHEDDEFEECCDDCCACCCD>A??BB<  
@SRR1215997.3 HWI-ST1311:58:C132FACXX:3:1101:3410:2170/1  
TGCTCAACGGCTCTCAGCTGGTGCCTGGACTGGCCACAGTGGCCTACGCCTGCTTCCACCGTAGCTGGCACCTGCGCACCAAGT
```

Unparalleled resolution & throughput, but ...

Sequencing is the medium for many different types of assays

Different measurements often require new methods

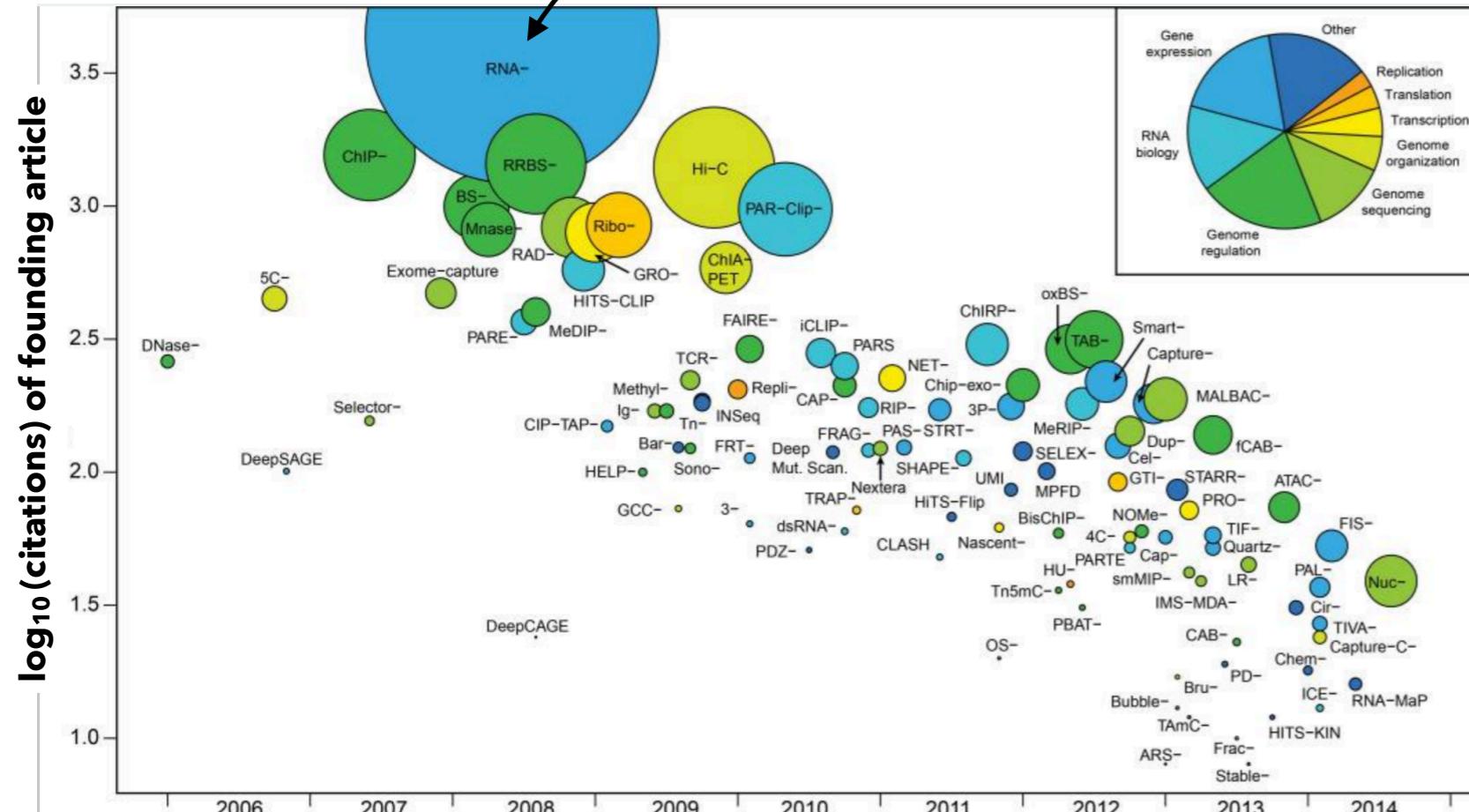
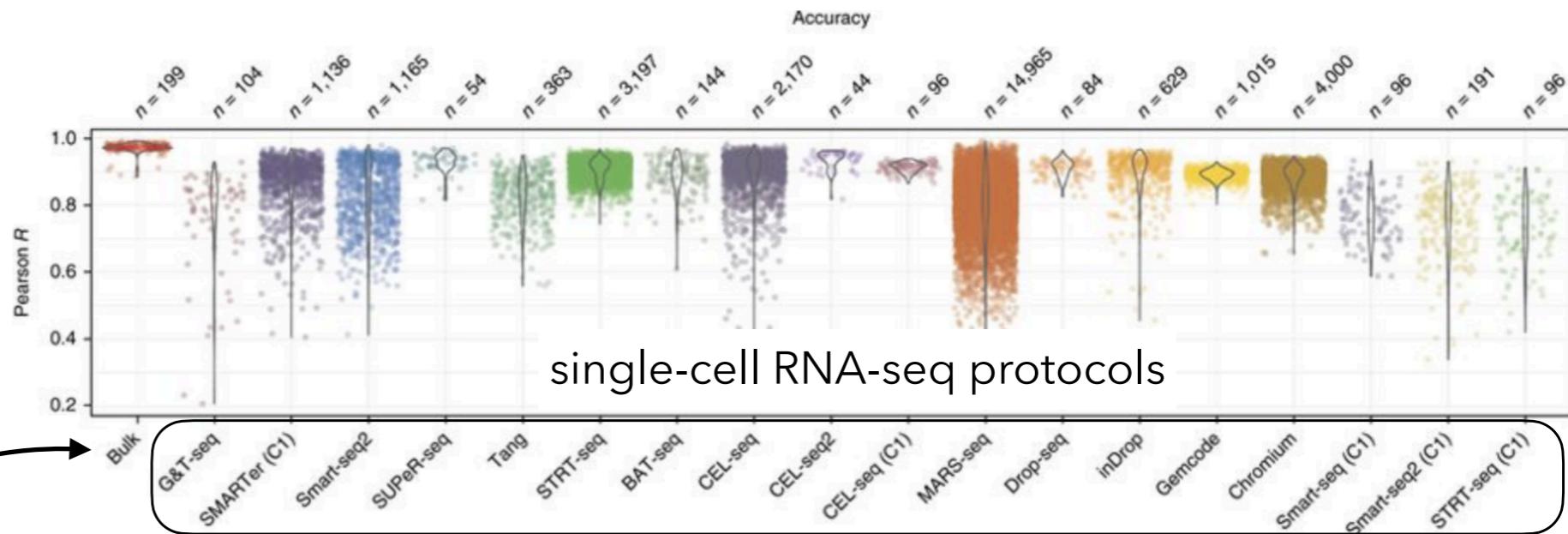
Produces so much data that analysis becomes the bottleneck

Magnitude of data requires fundamentally new approaches

# Challenges of studying RNA

## Understanding new types of data

Even just a single “seq” type  
(e.g. RNA-seq) has  
many different variations that  
cannot all be processed in  
the same way



Proliferation of new technologies

Measuring *different things*

Measuring in *different ways*

High-Throughput sequencing is the common factor (the medium)

# The computational challenge(s)

## Scalability

Muir et al. *Genome Biology* (2016) 17:53  
DOI 10.1186/s13059-016-0917-0

Genome Biology

OPINION

Open Access



### The real cost of sequencing: scaling computation to keep pace with data generation

Paul Muir<sup>1,2,3</sup>, Shantao Li<sup>4</sup>, Shaoke Lou<sup>4,5</sup>, Daifeng Wang<sup>4,5</sup>, Daniel J Spakowicz<sup>4,5</sup>, Leonidas Salichos<sup>4,5</sup>, Jing Zhang<sup>4,5</sup>, George M. Weinstock<sup>6</sup>, Farren Isaacs<sup>1,2</sup>, Joel Rozowsky<sup>4,5</sup> and Mark Gerstein<sup>4,5,7\*</sup>

"This new regime, in which costs scale with the amount of computational processing time, places a premium on driving down the average cost by developing efficient algorithms for data processing."

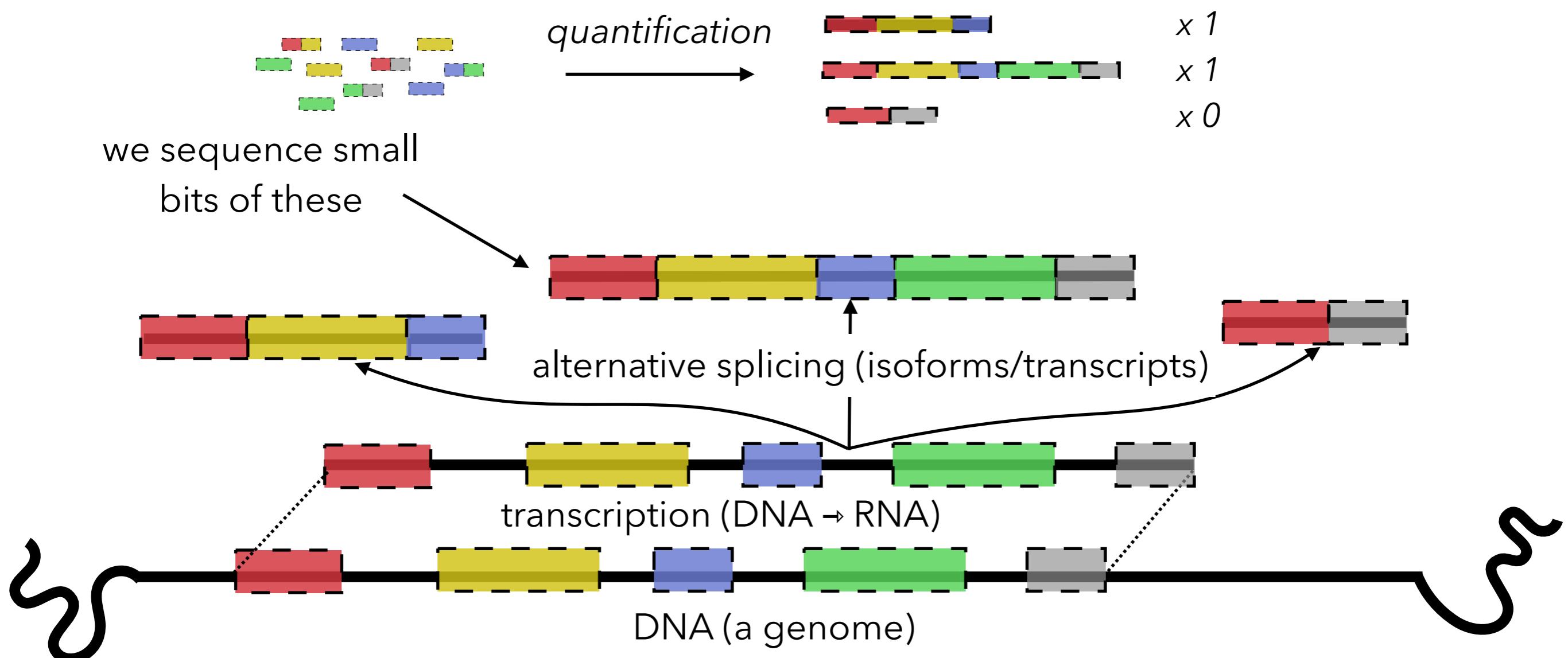
**Also, it's not just "new" data that is the problem:**

**In addition to new data, re-analysis of existing experiments often desired: In light of new annotations, discoveries, and methodological advancements.**

# Expression Study basics in one slide

Lets us (among other things) quantify expression in a tissue:

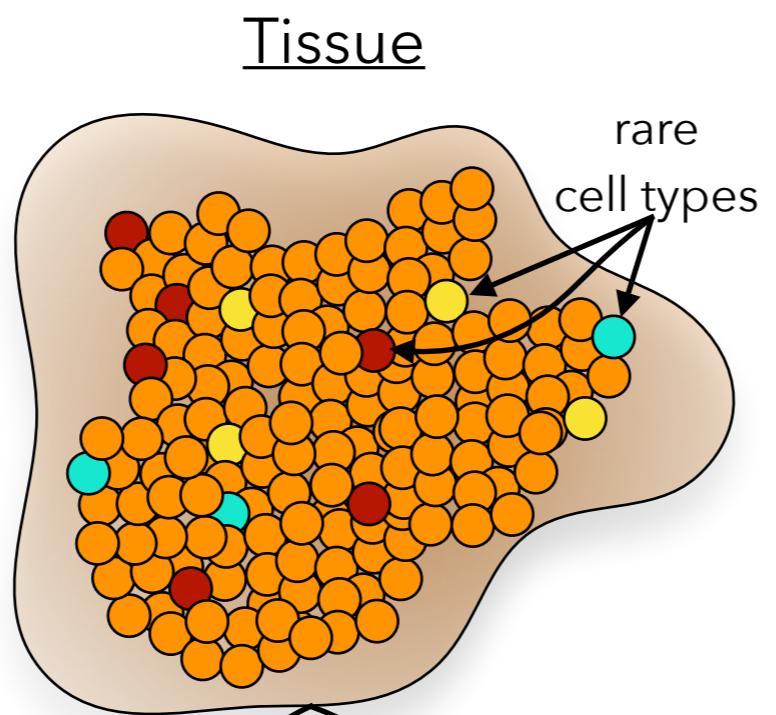
- What genes / transcripts are turned on?
- At what levels are they expressed?
- How do they respond to environment / stimuli?



# RNA-seq → single-cell RNA-seq

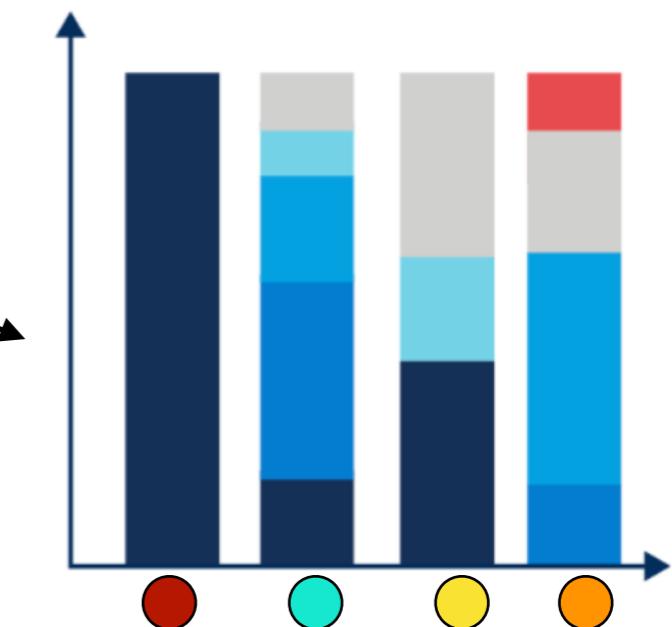
## Bulk RNA-seq

- Typically millions or 10s of millions of cells
- High-fidelity & high-sensitivity
- Measure transcript abundance at the population-level



## single-cell RNA-seq

- Typically tens of thousands of cells
- Low-coverage (reads / cell)
- Measure gene abundance at the single-cell level, specially useful for rare cell types.



involves cell-type identification  
(supervised or unsupervised)

# Single Cell Protocols

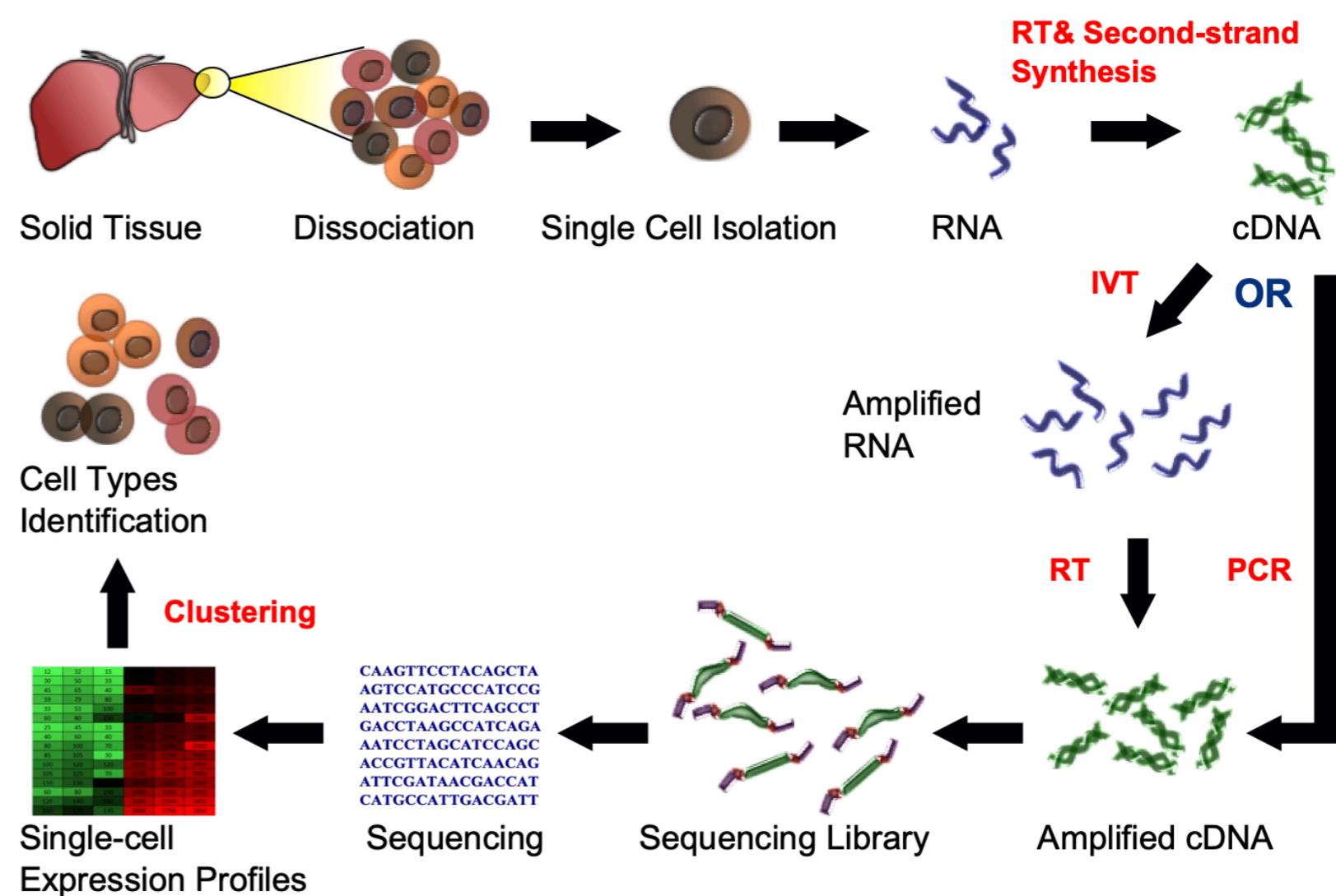
## Single Cell Isolation:

\* **Micro-well:** Requires human pipetting but are surface-markable however has very low-throughput.

\* **Micro-fluid (Fluidigm's C1):** It's based on an Integrated system which has higher capture-rate than micro-well, however throughput is still low.

\* **droplets:** This method has very high-throughput, add cellular-barcoding but has the caveat of higher sequencing cost because of low coverage of transcripts.

## Single Cell RNA-sequencing workflow:

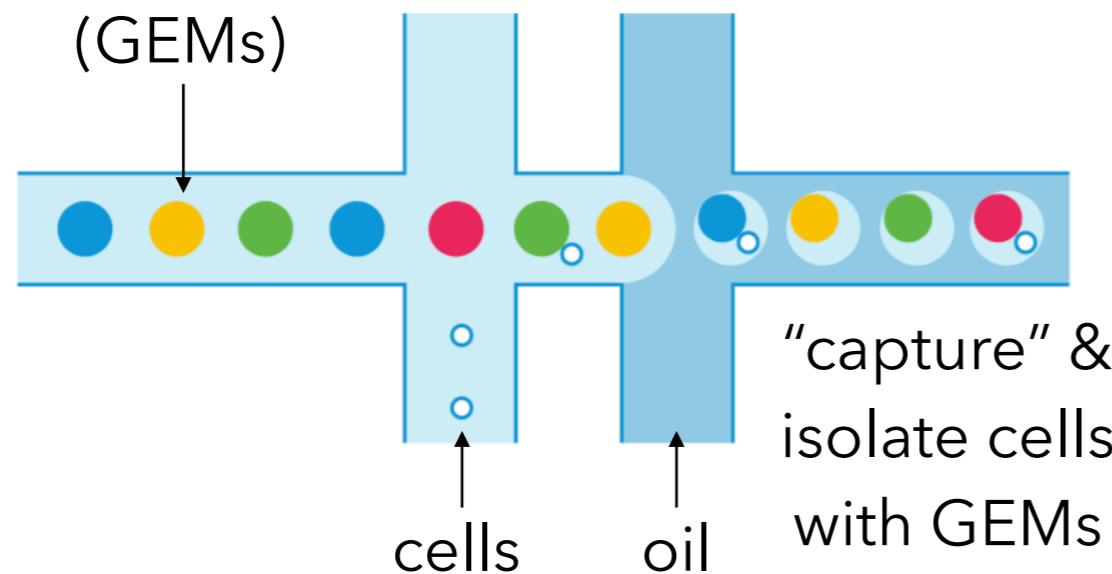


# How does it work?

Many different protocols – here is a gist of droplet-based (microfluidic) techniques.

Gel beads in *EM*ulsion

**isolate**



**collect**

"capture" &  
isolate cells  
with GEMs

**tag**

attach CB & UMI

Cell Barcode (CB)

"What cell (GEM bead)  
did I come from?"

**lib. prep,  
amplification &  
sequencing**

Unique Molecular Identifier (UMI)

"What pre-amplification molecule  
did I come from?"

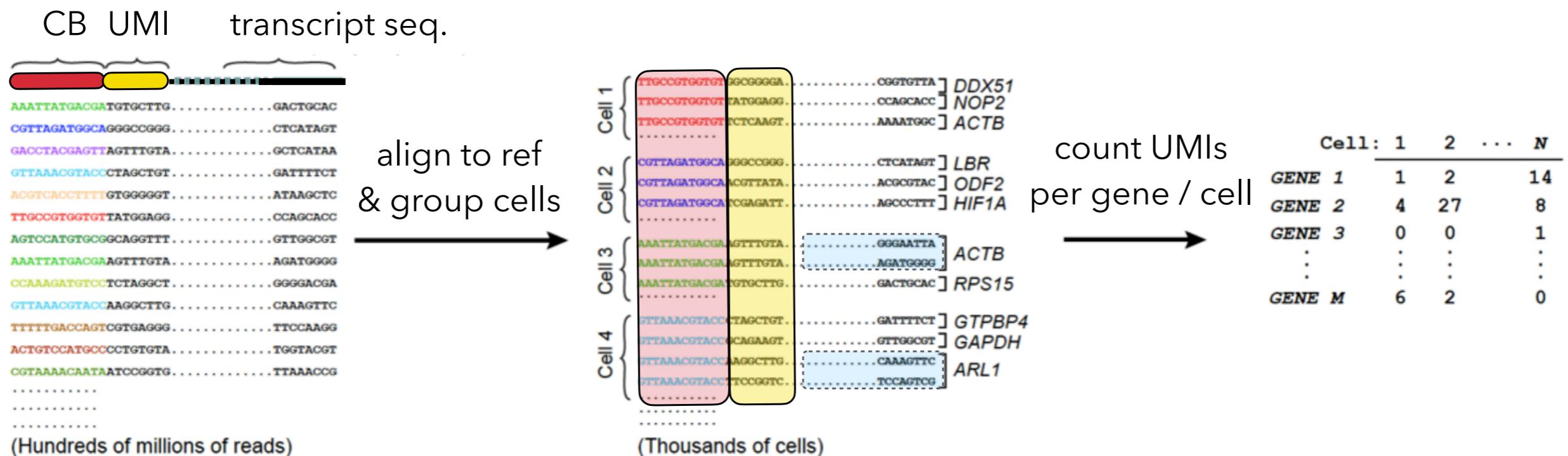


# How does it work?

**In theory** using the CB & UMI to estimate gene expression is easy:

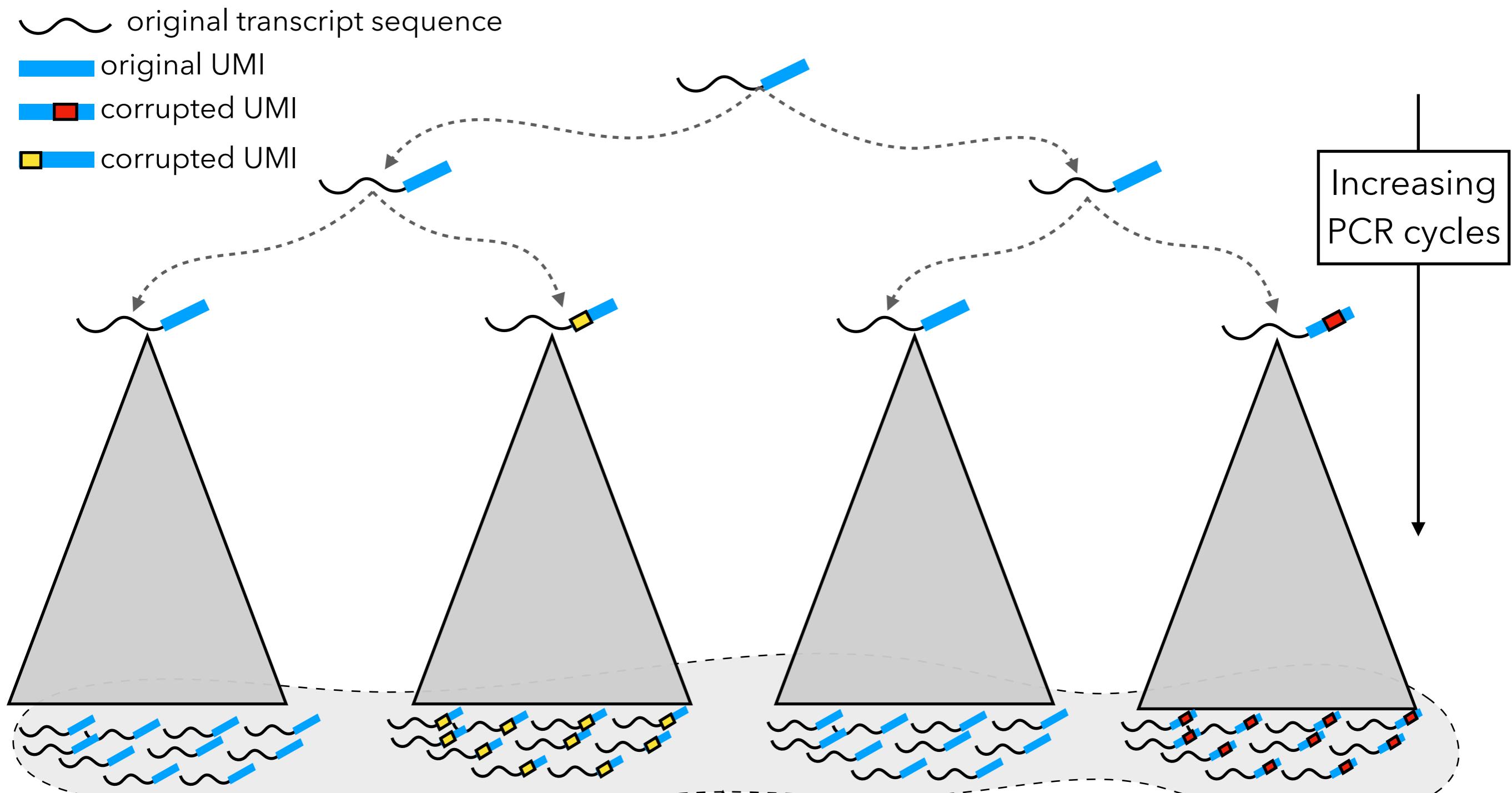
1. Group sequenced reads by their CB (read → cell)
2. Collapse all reads with the same UMI (read in cell → molecule)
3. Count!

... how many lines is that in Pandas?



**In practice** the process is **much** more complicated

# The cost of small input material

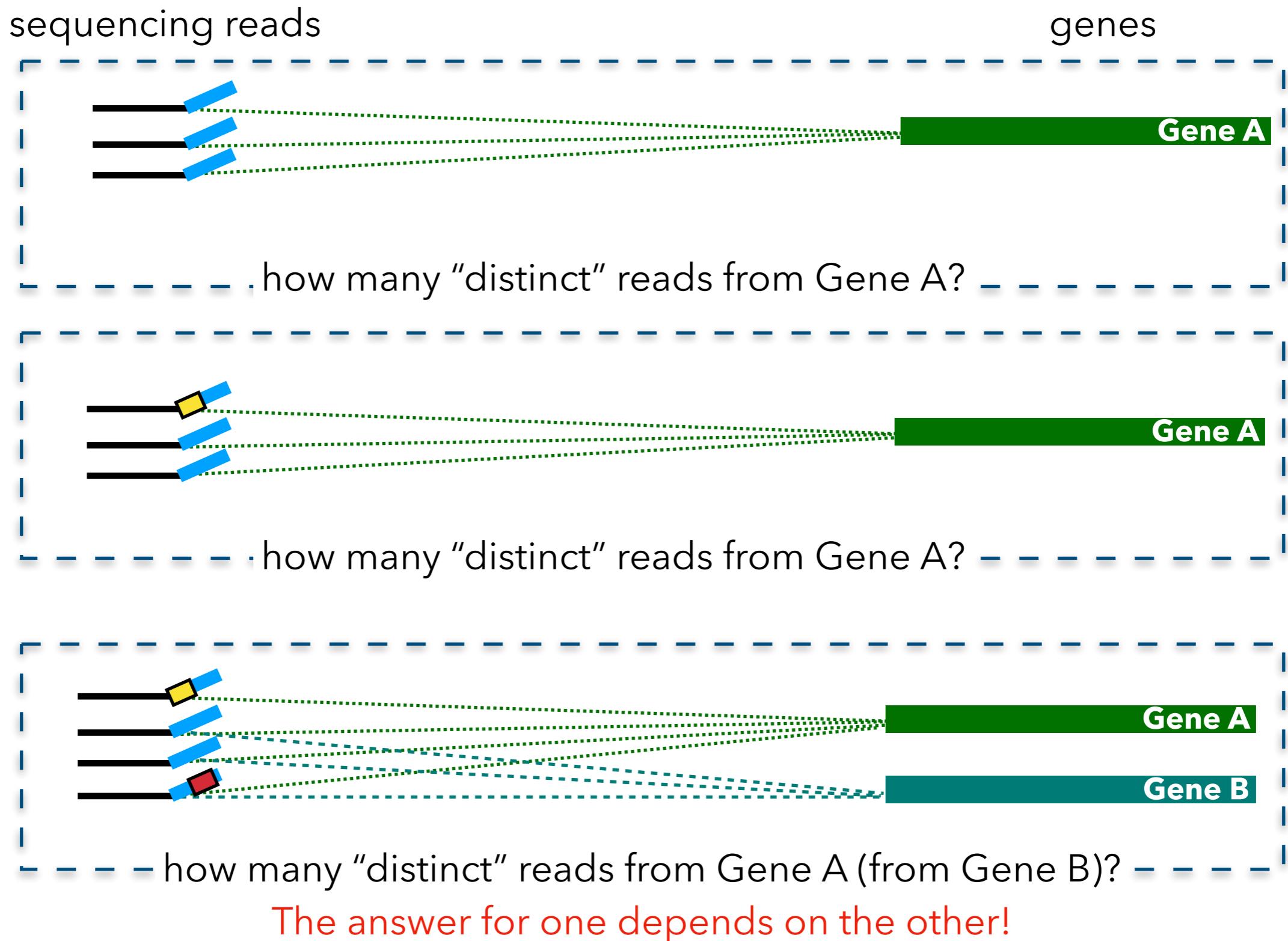


Sample reads from this exponentially amplified pool of molecules

PCR “Duplicates” of the same molecule can now have distinct UMIs – but probably *similar*

Must “collapse” similar UMIs, or we will vastly mis-estimate # of original molecules

# Multi-mapping is the real culprit



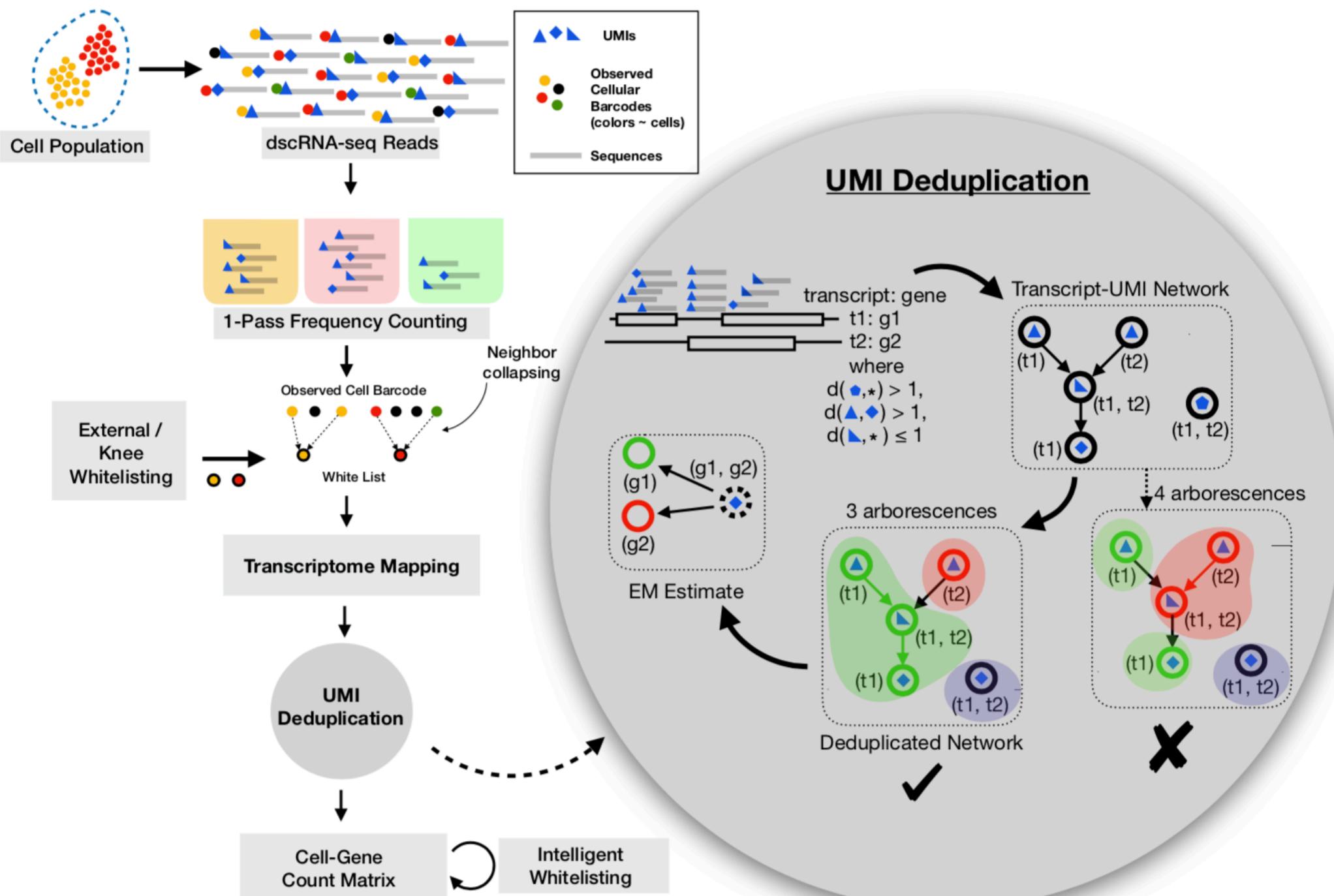
# single-cell quantification is challenging

There are *many* reasons single-cell quantification is challenging!

most of these stem from the difficulty of cell isolation & small quantities of genetic material:

- 1 UMI  $\neq$  1 pre-PCR molecule (often); because of small material per-cell, typical protocols involve many rounds of PCR; UMI tags are subject to PCR and sequencing error.
- 1 Cell Barcode  $\neq$  1 Cell (sometimes); capture-related issues (e.g. doublets & empty droplets)
- “dropout” & bias due to sampling of relatively small number of reads from highly-amplified pool of original molecules
- UMI “collisions” are possible due to limited UMI pool

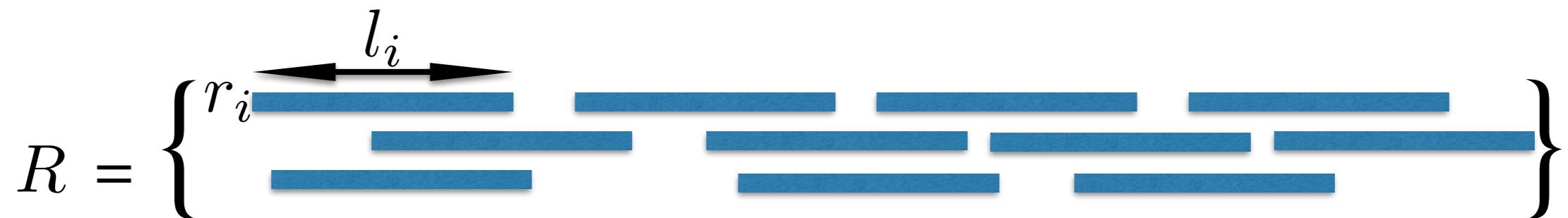
# alevin: dscRNA-seq quantification



# The read-alignment problem

**Given:** A collection of sequencing reads, and some target sequence (e.g. a genome)

**Find:** For each read, all locations where the read is within edit distance  $\epsilon$  of the reference, and the edits that achieve this distance.

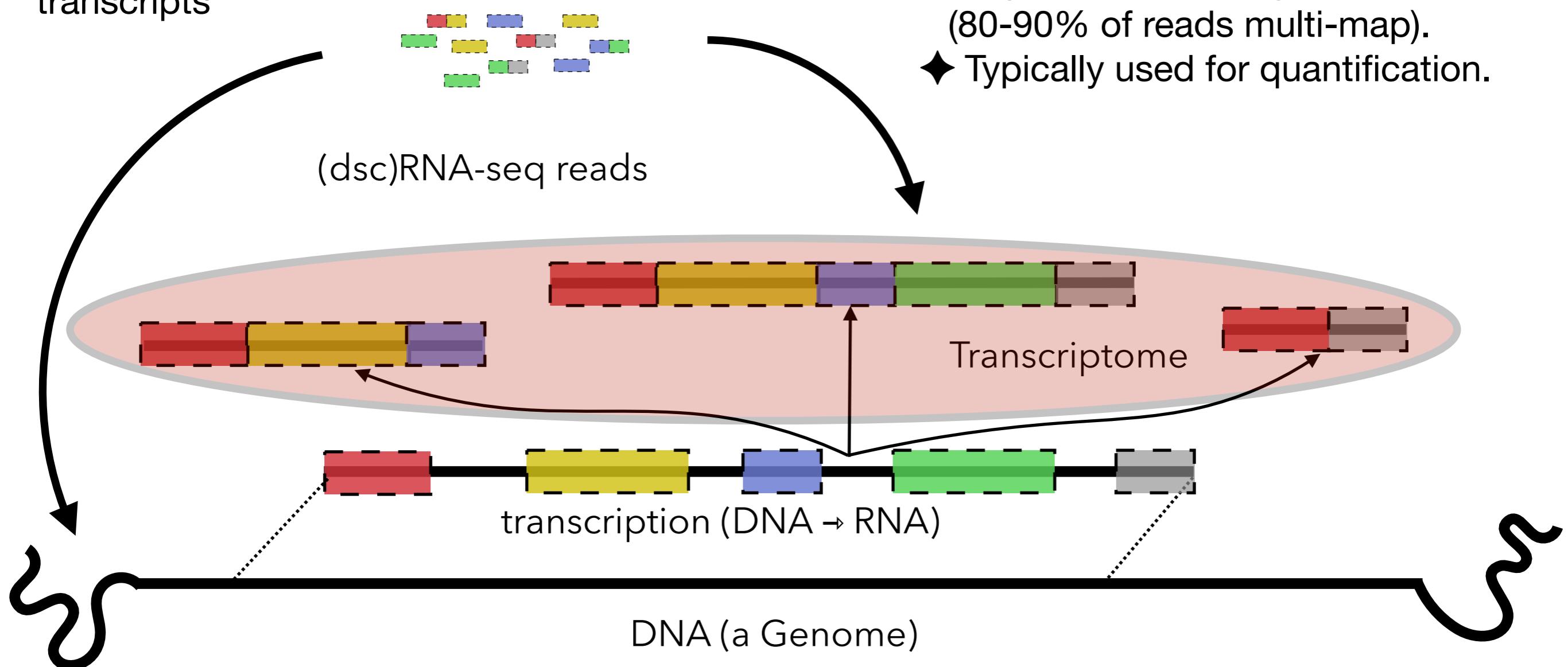


$$d : \Sigma^{|u|} \times \Sigma^{|v|} \rightarrow \mathbb{Z} \text{ and } \epsilon \text{ (maximum edit distance)}$$

# Read-Alignment Strategies

## Genome Mapping

- ◆ Tools like: STAR, HISAT.
- ◆ Maps to full Genome (~3G for Human).
- ◆ Less per read multi-mapping.
- ◆ Typical use case includes finding new transcripts



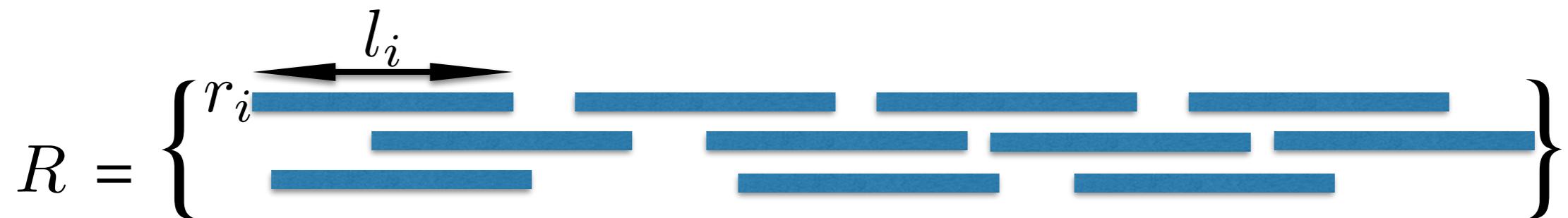
## Transcriptome Mapping

- ◆ Tools like: Bowtie, RapMap, Selective-alignment.
- ◆ Maps to Transcriptome (~300M for Human).
- ◆ High multi-mapping rate (80-90% of reads multi-map).
- ◆ Typically used for quantification.

# The read-MAPPING problem

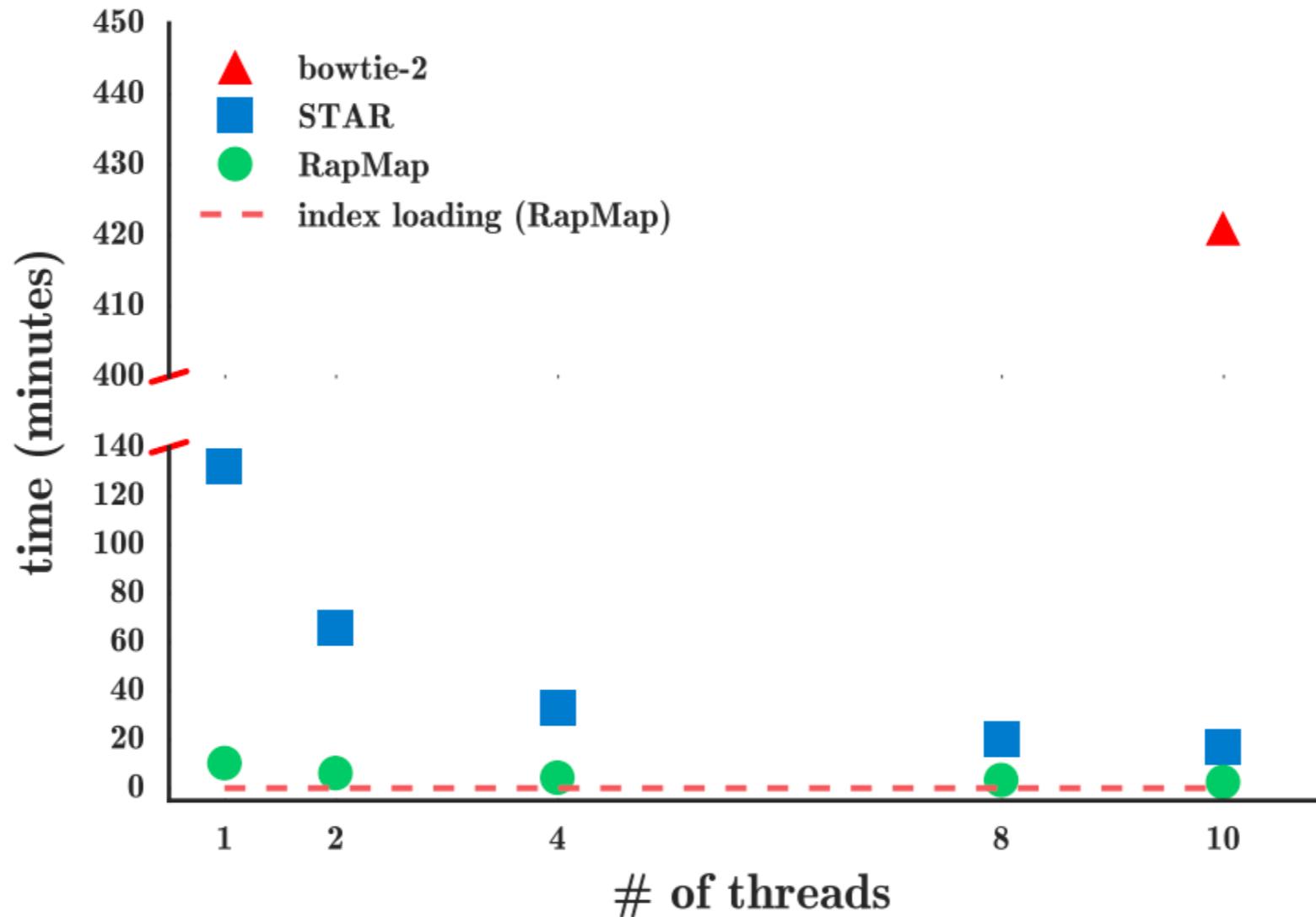
**Given:** A collection of sequencing reads, and some target sequence (e.g. a genome)

**Find:** For each read, all locations where the read is within edit distance  $\epsilon$  of the reference, ~~and the edits that achieve this distance.~~



$$d : \Sigma^{|u|} \times \Sigma^{|v|} \rightarrow \mathbb{Z} \text{ and } \epsilon \text{ (maximum edit distance)}$$

# Why relax the problem ?



- Analysis on 75Million (76bp) PE reads to human transcriptome; RapMap does in just **minutes**.

# Phylogeny of read-alignment

## Aligning (Mapping) NGS Reads

DNA-sequencing

RNA-sequencing

Genome (Spliced)

Transcriptome

- Aligns RNA-seq reads to genome
- Challenge of **Spliced Alignment**
- Example: topHat, STAR, HISAT(1/2)

- Aligns RNA-seq reads to transcriptome
- Challenge of **high multi-mapping rate**
- Example: Bowtie(1/2), BWA(SW/MEM)

Aligner

- Base-to-Base Alignment (CIGAR string)
- Selective-Alignment

Mapper

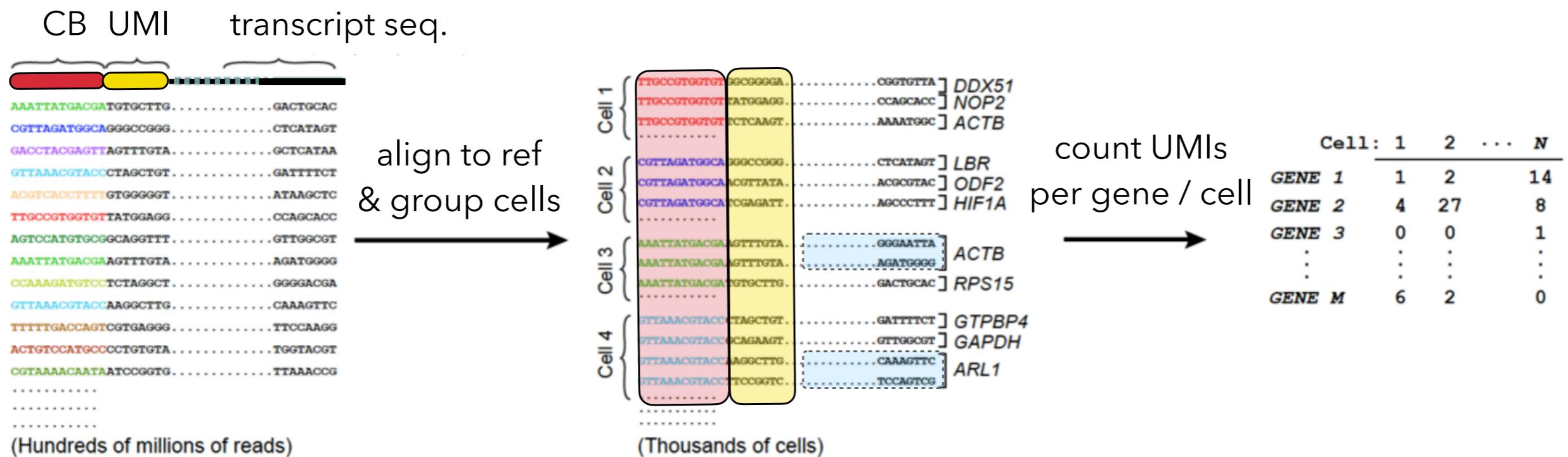
- **NO** CIGAR string
- RapMap

# Recap

**In theory** using the CB & UMI to estimate gene expression is easy:

1. Group sequenced reads by their CB (read → cell)
  2. Collapse all reads with the same UMI (read in cell → molecule)
  3. Count!

... how many lines is that in Pandas?



**In practice** the process is **much** more complicated

# dscRNA-seq quantification methods

## CellRanger, (2017) \*\* – STAR aligner

Zheng, Grace XY, et al. "Massively parallel digital transcriptional profiling of single cells." *Nature communications* 8.1 (2017): 1-12.

## UMI-tools, (2017) – STAR aligner

Smith, Tom, Andreas Heger, and Ian Sudbery. "UMI-tools: modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy." *Genome research* 27.3 (2017): 491-499.

## Alevin, (2018) – Salmon's quasi-mapping / selective -alignment

Srivastava, Avi, et al. "Alevin efficiently estimates accurate gene abundances from dscRNA-seq data." *Genome biology* 20.1 (2019): 65.

## STAR-solo, (2019) \* – STAR aligner

Dobin, Alexander, et al. "STAR: ultrafast universal RNA-seq aligner." *Bioinformatics* 29.1 (2013): 15-21. [ No preprint for single-cell version ]

## Hera-T, (2019) \* – Hera's aligner

Tran, Thang, et al. "Hera-T: an efficient and accurate approach for quantifying gene abundances from 10X-Chromium data with high rates of non-exonic reads." *bioRxiv* (2019): 530501.

## bustools, (2019) \* – Kallisto's pseudo-alignment

Melsted, Páll, et al. "Modular and efficient pre-processing of single-cell RNA-seq." *BioRxiv* (2019): 673285.

This is not an exhaustive list, methods were selected based  
on read-alignment strategy

\* Not Peer-reviewed

\*\* No methods paper

# One shortcoming with current dscRNA-seq processing techniques

Unfortunately, current approaches have no principled way to deal with reads that map between multiple genes ... simply **discard them.**

This may seem like a minor issue, but, in a typical dataset, this is **13-23% of the reads!**

Sample	Percentage
Human PBMC 4k	13.8
Human PBMC 8k	13.8
Mouse Neurons 900	21.6
Mouse Neurons 2k	22.5
Mouse Neurons 9k	17.1

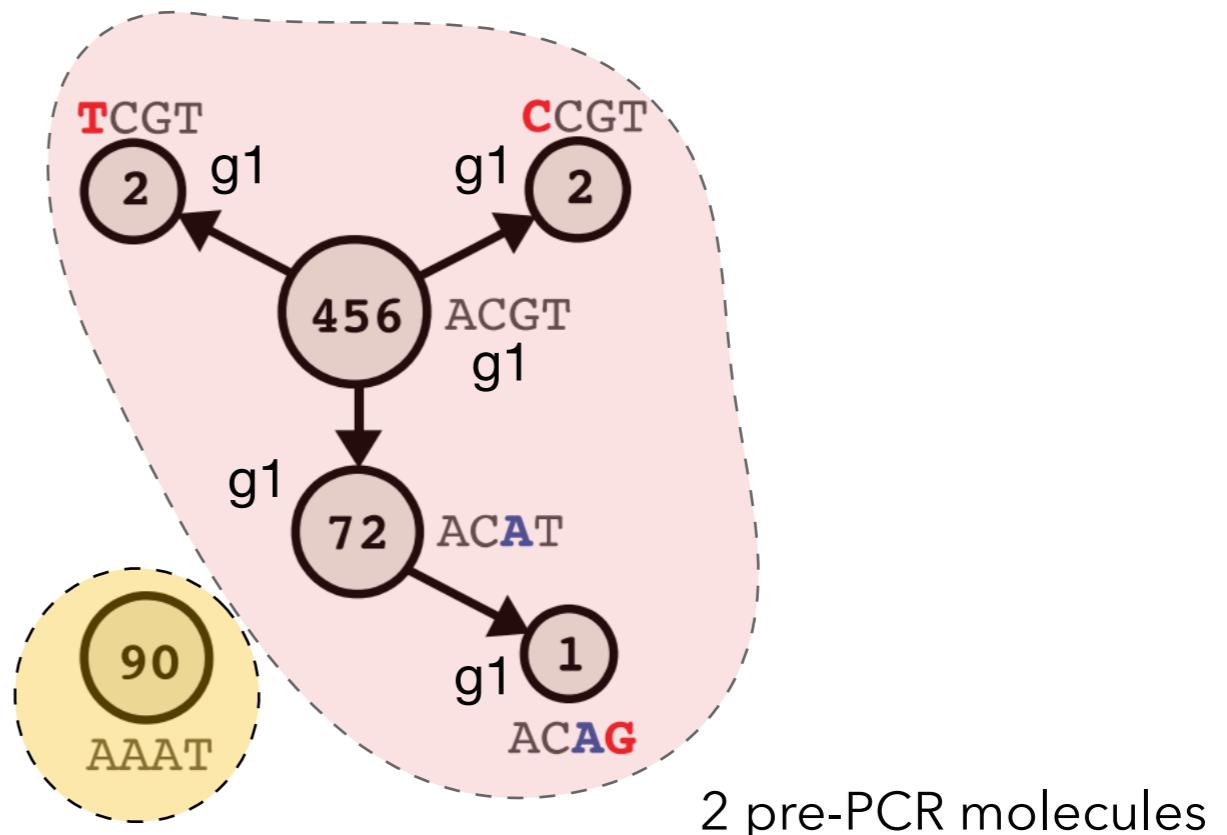
# UMI-Deduplication

**Given:** For each cell, a set of UMI's with their frequencies based on gene mapping.

**Find:** Number of pre-PCR molecules by deduplicating the set of UMIs.

## UMI-tools directional approach

- ✿ Start with highest frequency UMI & recursively add an edge to all the UMIs within 1-edit distance.
- ✿ The direction of the edge is from high frequency node to low frequency node.
- ✿ Collapse the connected component into 1 pre-PCR molecule.



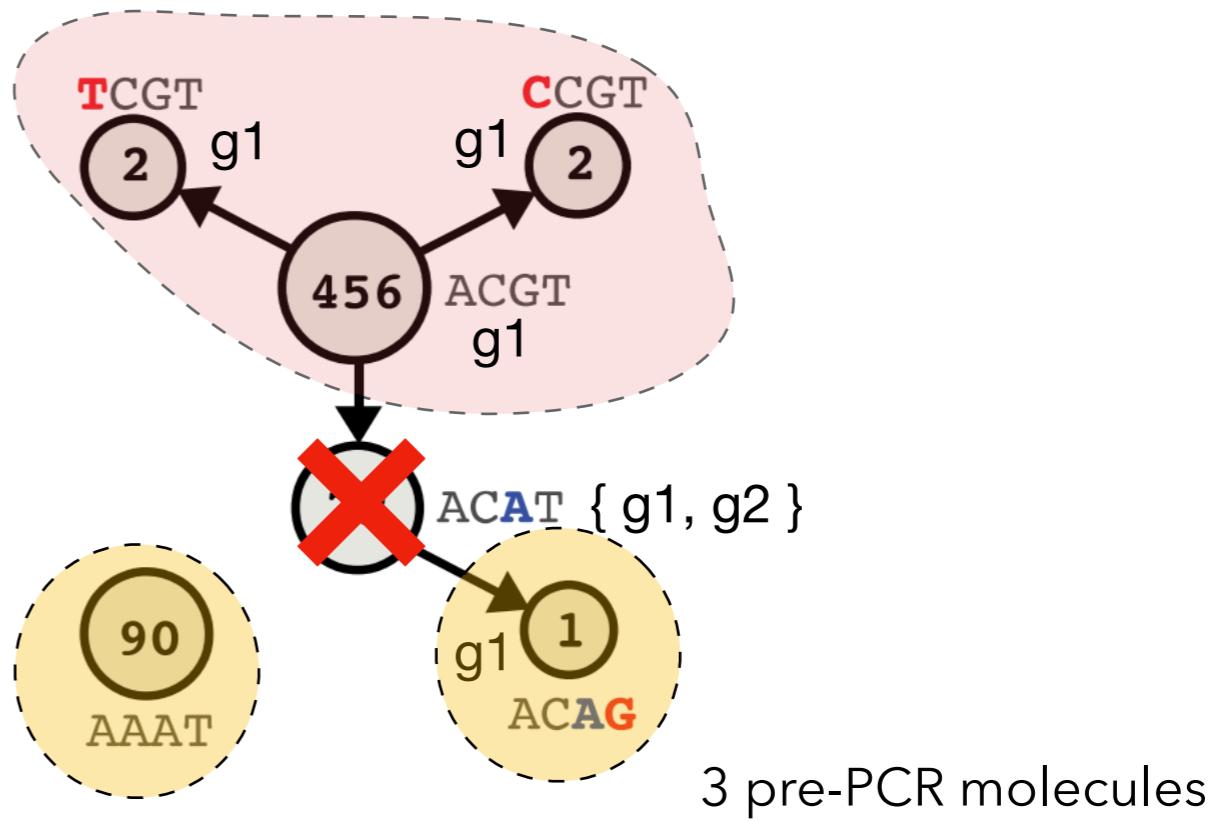
# UMI-Deduplication

**Given:** For each cell, a set of UMI's with their frequencies based on gene mapping.

**Find:** Number of pre-PCR molecules by deduplicating the set of UMIs.

## UMI-tools directional approach

- Start with highest frequency UMI & recursively add an edge to all the UMIs within 1-edit distance.
- The direction of the edge is from high frequency node to low frequency node.
- Collapse the connected component into 1 pre-PCR molecule.
- Dropping multi-mapping reads breaks the network.



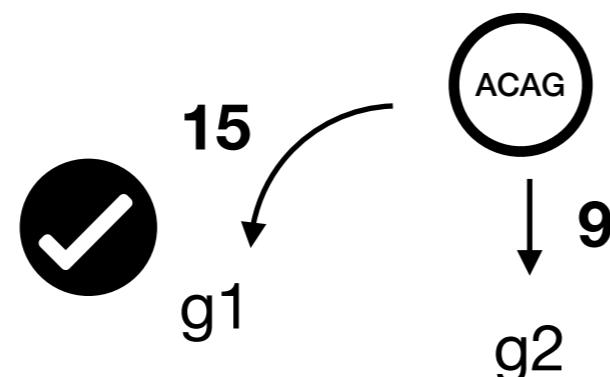
# UMI-Deduplication

**Given:** For each cell, a set of UMI's with their frequencies based on gene mapping.

**Find:** Number of pre-PCR molecules by deduplicating the set of UMIs.

## CellRanger (STAR-solo) approach

- CellRanger follows bottom-up instead of top-down approach i.e. unlike UMI-tools way of assigning UMI to gene, cellranger assign genes to UMIs.
- In the first round, within each cell, cellranger assigns only the “**gene-unique**” mapped read to a UMI, generating a frequency distribution of genes for each UMI.
- In the second round of processing, each UMI is assigned to the **highest frequency** gene from its distribution.



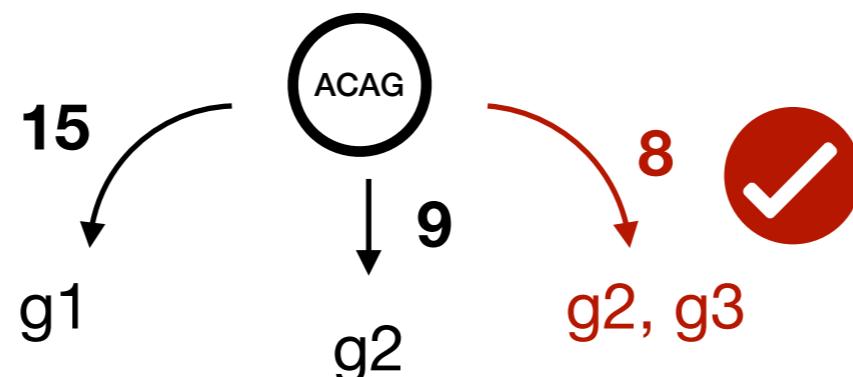
# UMI-Deduplication

**Given:** For each cell, a set of UMI's with their frequencies based on gene mapping.

**Find:** Number of pre-PCR molecules by deduplicating the set of UMIs.

## CellRanger (STAR-solo) approach

- CellRanger follows bottom-up instead of top-down approach i.e. unlike UMI-tools way of assigning UMI to gene, cellranger assign genes to UMIs.
- In the first round, within each cell, cellranger assigns only the “**gene-unique**” mapped read to a UMI, generating a frequency distribution of genes for each UMI.
- In the second round of processing, each UMI is assigned to the **highest frequency** gene from its distribution.



# Parsimonious UMI Graph (PUG) resolution

Represent the UMIs / transcripts relationship as a graph  $G = (V, E)$ .

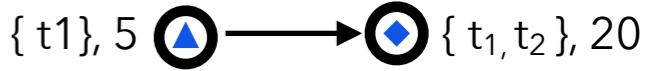
Each **vertex** ( $v \in V$ ) is a tuple:

- $\mathbf{eq}_v$ : the set of transcripts to which the read maps
- $\mathbf{s}_v$  : the UMI sequence tagging the read

Each  $v$  has count,  $c(v)$  – number of “**equivalent**” reads.

“**equivalent**” means aligns to same transcripts and has same UMI.

There is an **edge** ( $e = \{u, v\} \in E$ ) if, for some chosen edit distance  $\tau$ :

- $s_u = s_v$  and  $|eq_u \cap eq_v| > 0$  (*bidirected edge*) 
- $d(s_u, s_v) < \tau$ ,  $c(u) \sim c(v)$  and  $|eq_u \cap eq_v| > 0$  (*bidirected edge*)
- $d(s_u, s_v) < \tau$ ,  $c(u) > 2c(v)+1$  and  $|eq_u \cap eq_v| > 0$  (*directed*  $v \rightarrow u$ ) 

# A parsimony-guided approach

We will attempt to *explain* the PUG using the *minimum number* of pre-PCR molecules. More formally:

**Given:** UMI-resolution graph  $G = (V, E)$

**Find:** a *minimum cardinality* cover by *monochromatic arborescences*.

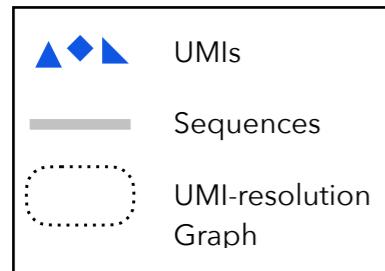
↑  
We seek parsimony

↑  
Each component can be ascribed  
to the reads sequenced from a *single  
initial molecule* (before amplification).

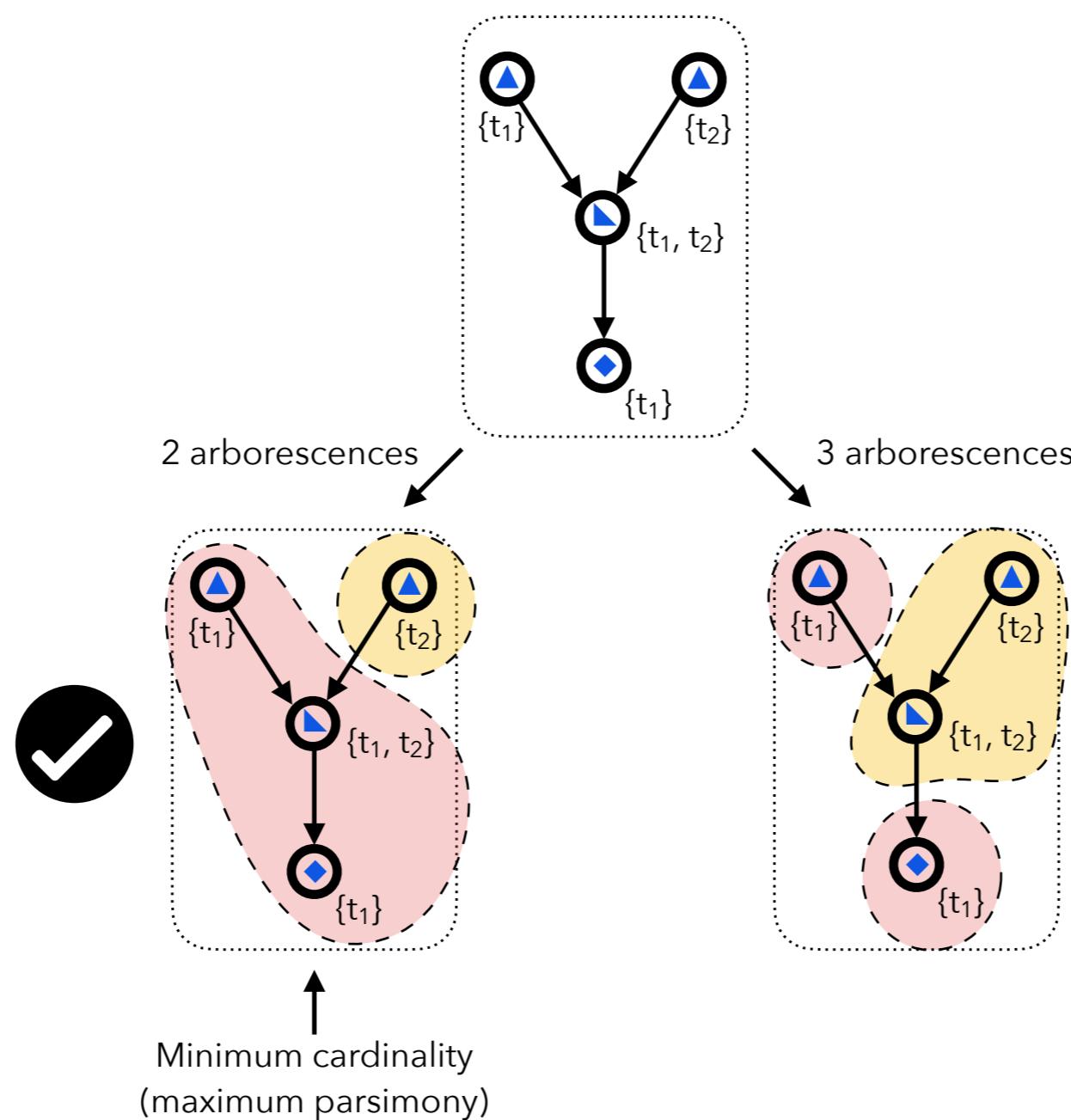
Each *monochromatic arborescence* is a set of vertices that can all  
be described as “reads” coming from the same, original transcript.

The decision problem is *NP-complete* (reduction from DOMINATING  
SET); but experimental instances are usually *simple* and we adopt a  
*greedy heuristic*. (solutions are known to be optimal in > 80% of cases)

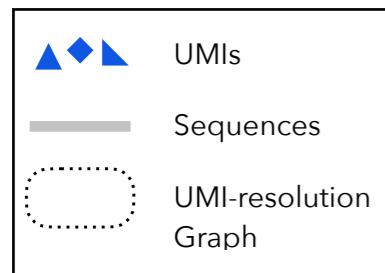
# UMI Resolution



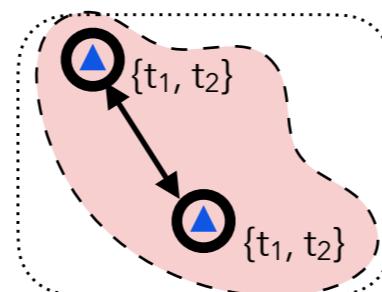
$eq_1 \xrightarrow{\quad} eq_2$  Unidirectional Edge  
 $eq_1 \leftrightarrow eq_2$  Bidirectional Edge  
where  $d(\triangle, \diamond) \leq \tau$  and  $|eq_1 \cap eq_2| > 0$



# UMI Resolution



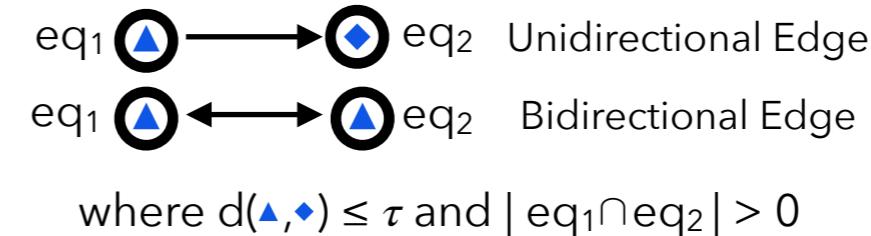
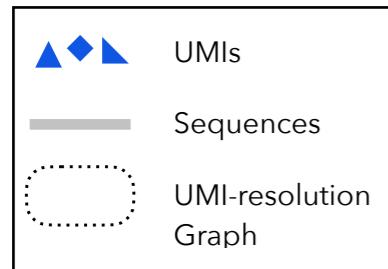
$eq_1 \xrightarrow{\quad} eq_2$  Unidirectional Edge  
 $eq_1 \leftrightarrow eq_2$  Bidirectional Edge  
where  $d(\triangle, \diamond) \leq \tau$  and  $|eq_1 \cap eq_2| > 0$



Equally-parsimonious under  $t_1$  **or**  $t_2$

- Parsimony cannot resolve the gene of origin here
- We treat UMI as gene-ambiguous & resolve via EM-algorithm

# EM & Tiers Characterization



tier 1



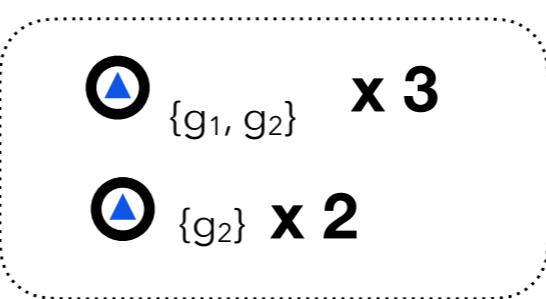
Gene: Count : Tier

$g_1. : 4 : t1$

$g_2. : 2 : t1$

Detailed description: Below the tier 1 cluster, a downward arrow points to a "Gene: Count : Tier" table. The first row shows gene  $g_1$  with a count of 4 and tier  $t1$ . The second row shows gene  $g_2$  with a count of 2 and tier  $t1$ .

tier 2



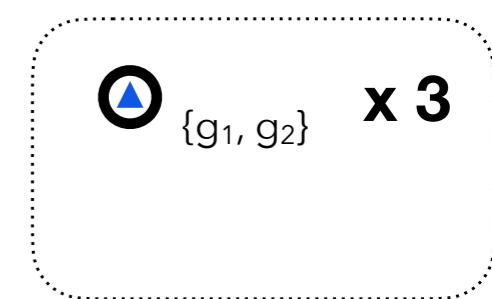
Gene: Count : Tier

$g_1. : 1.5 : t2$

$g_2. : 3.5 : t1$

Detailed description: Below the tier 2 cluster, a downward arrow points to a "Gene: Count : Tier" table. The first row shows gene  $g_1$  with a count of 1.5 and tier  $t2$ . The second row shows gene  $g_2$  with a count of 3.5 and tier  $t1$ .

tier 3



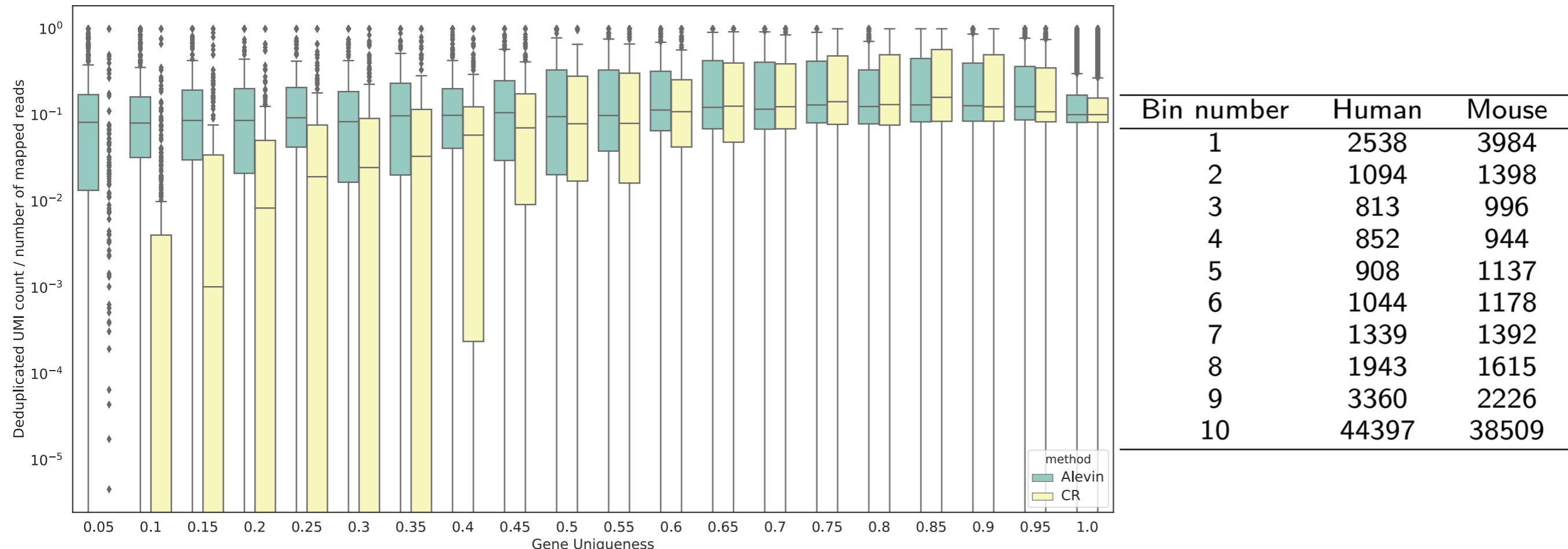
Gene: Count : Tier

$g_1. : 1.5 : t3$

$g_2. : 1.5 : t3$

Detailed description: Below the tier 3 cluster, a downward arrow points to a "Gene: Count : Tier" table. Both gene  $g_1$  and gene  $g_2$  have a count of 1.5 and tier  $t3$ . The entire row for gene  $g_2$  is highlighted with a green oval.

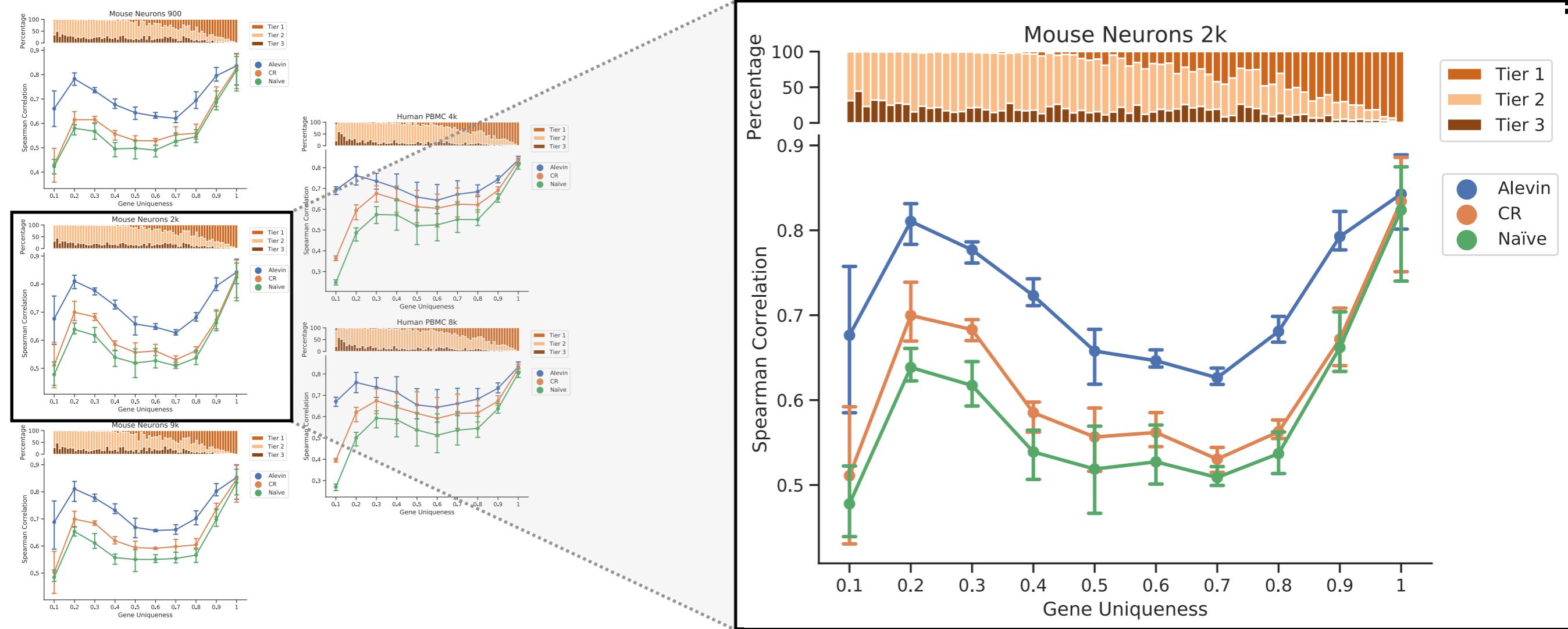
# Discarding gene-ambiguous reads does not affect all genes equally



Stratify genes by sequence uniqueness and look at distribution of # of de-duplicated UMIs per input read (PBMC 4k). Large effect for genes with lots of sequence homology (including important paralog families). Thus, the discarding of gene multi-mapping reads leads to systematic bias.

# Accounting for gene-ambiguous improves correlation with bulk RNA-seq

Trend across 5 public 10x datasets<sup>+</sup>



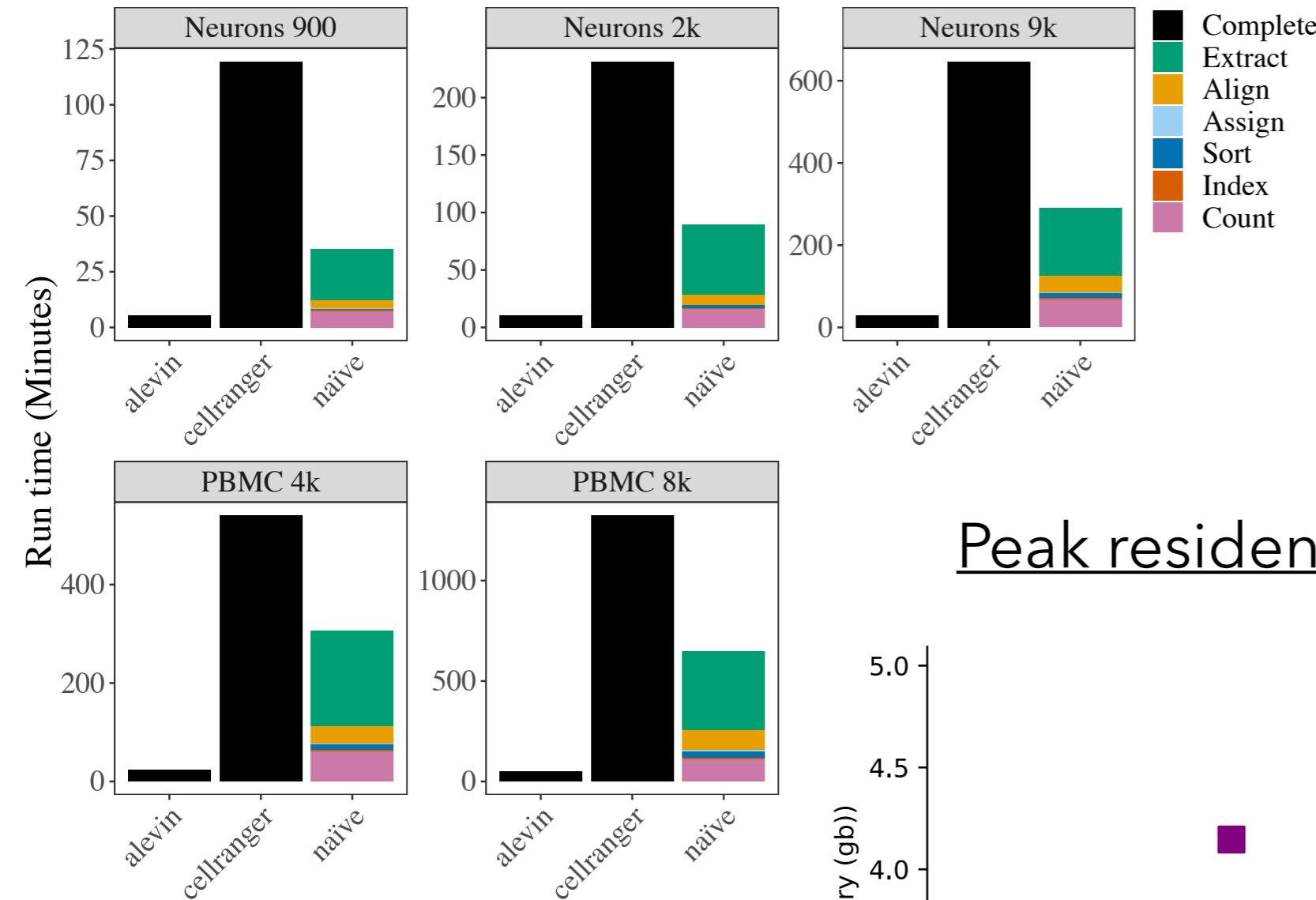
For each scRNA-seq dataset, we obtained multiple bulk samples from the same tissue. We compare gene-level quants in bulk (quantified with RSEM\*) to average expressions across single cells.

\*Li, B., & Dewey, C. N. (2011). RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC bioinformatics*, 12(1), 323.

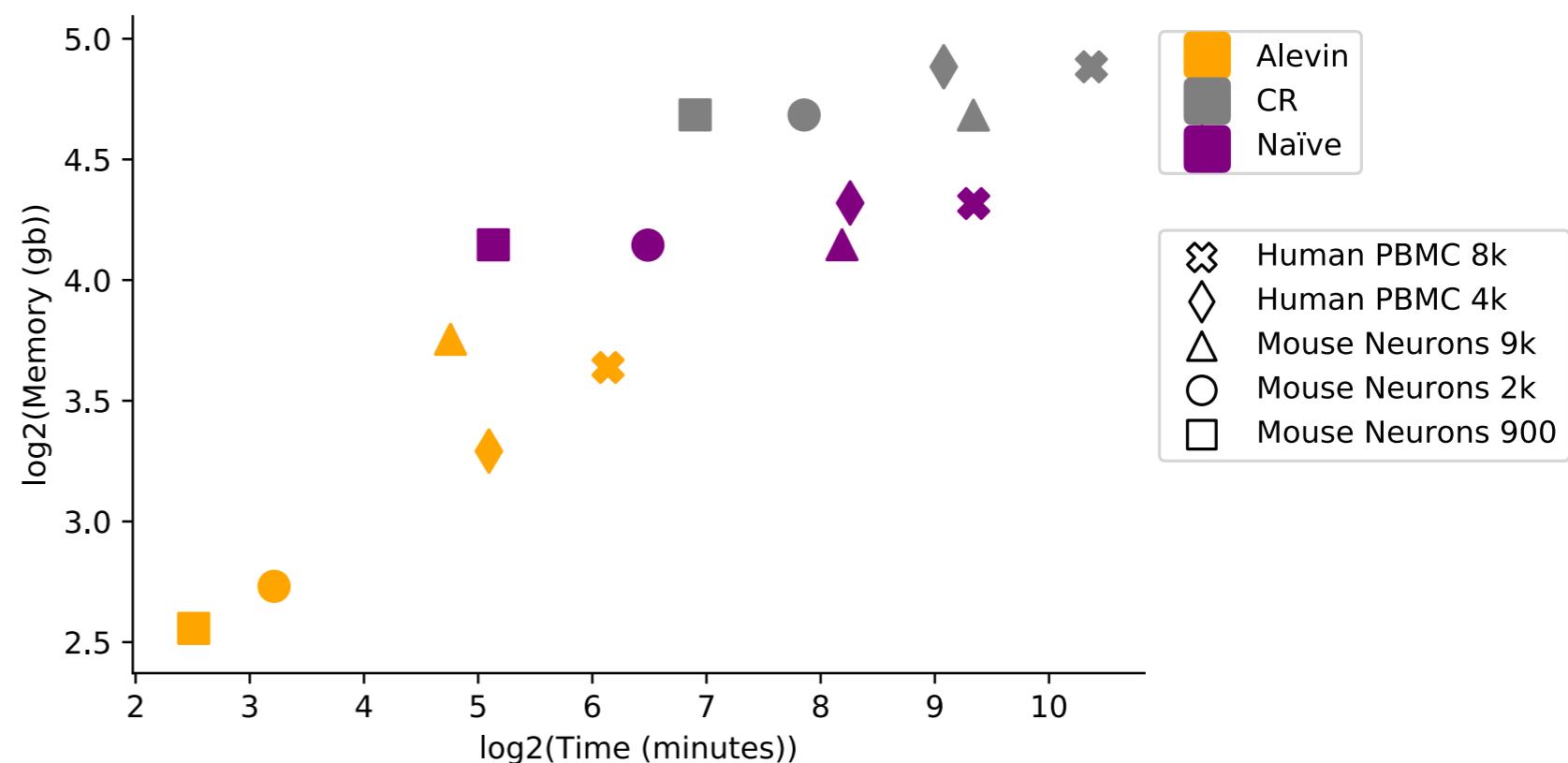
+Zheng, Grace XY, et al. "Massively parallel digital transcriptional profiling of single cells." *Nature communications* 8 (2017): 14049.

# alevin is fast & efficient

## Wall clock time (16 threads)



## Peak resident memory (16 threads)



# alevin-tutorial

<https://combine-lab.github.io/alevin-tutorial/>



## Alevin w/ Feature Barcodes

Feature Barcoding based Single-Cell Quantification with alevin ...

20 Apr 2020 • on [alevin](#)

## Spatial Alevin

Spatial Single-Cell Quantification with alevin ...

6 Apr 2020 • on [alevin](#)

## How to use alevin with iSEE

Preparing alevin output for exploration with iSEE ...

25 Mar 2020 • on [alevin](#)

## Alevin Velocity

RNA Velocity with alevin ...

22 Mar 2020 • on [alevin](#)

## Selective Alignment

Fast is Good but Fast and accurate is better ! ...

30 Oct 2019 • on [selective-alignment](#)