# UNIT 3  TRANSCRIPTOMICS

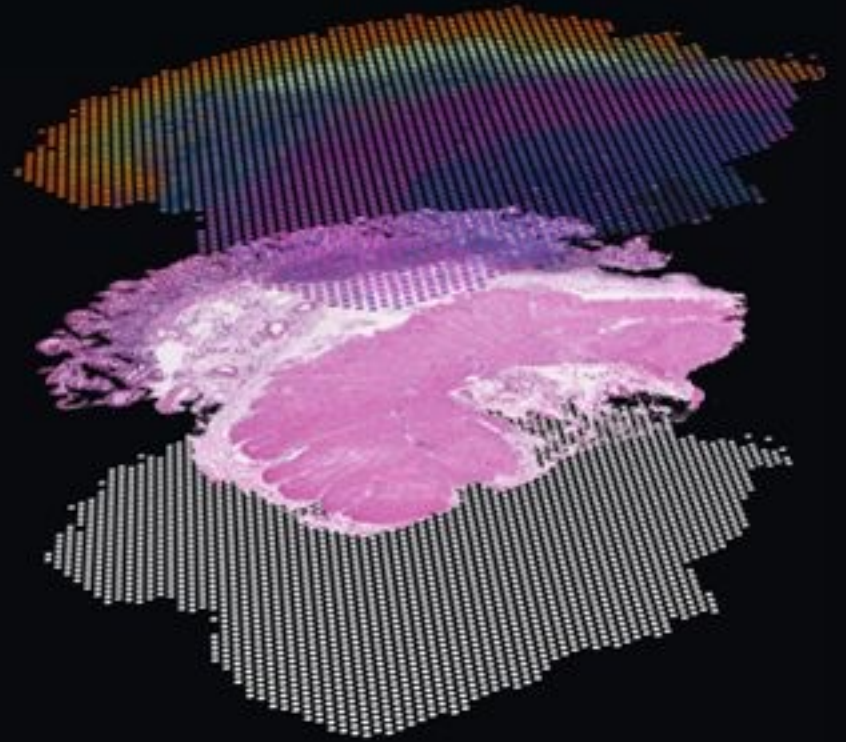# Spatial Transcriptomics

April 28, 2022

# Outline

- Spatial Transcriptomics
  - Sequencing based techniques
    - 10X Visium
  - Imaging based techniques
    - MERFISH
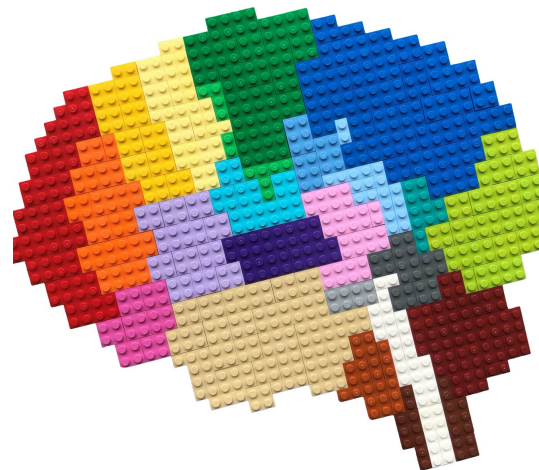- Encoding of sequence data
  - Hemming code
  - One Hot
  - Simplex encoding

# Single-cell and Spatial Transcriptomics

Single-cell transcriptomics

Bulk transcriptomics

Physiological reconstruction

Spatial transcriptomics

Image credit: Bo Xia @BoXia7

# Dimensionalities in transcriptomes

- Samples
- Transcripts / Genes
- Cells / Nucleus
- Spatial Locations
- Time / Differentiation stage

# Spatial transcriptomics technologies

- **Sequencing based**

- Major steps
  - 1. Dissection, capturing
  - 2. Barcoding, sequencing

- Examples
  - 10X Visium
  - Slide-seq
  - Nanostring GeoMx

- **Imaging based**

- Major steps
  - 1. Target and probe design
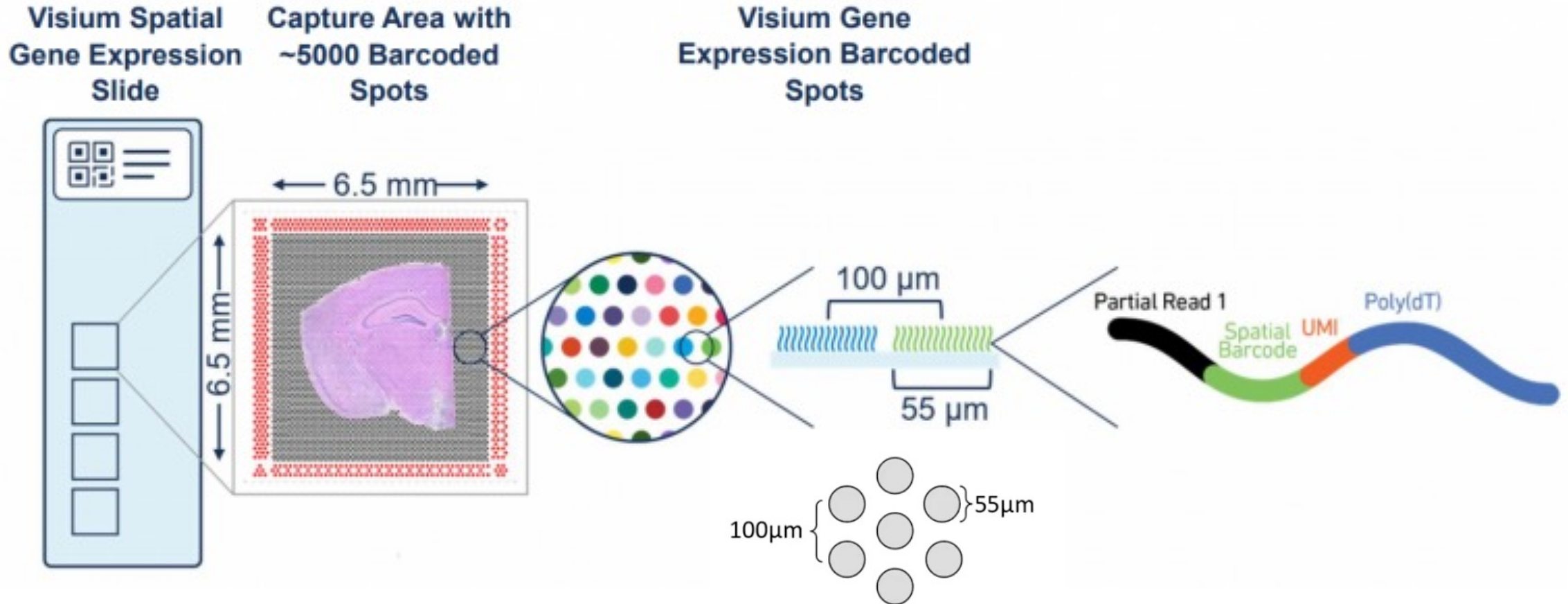  - 2. Fluorescence in situ hybridization (FISH)

- Examples
  - MERFISH
  - seqFISH

# Spatial transcriptomics technologies

| | **Sequencing based** | **Imaging based** |
|---|---|---|
| **Pros** | • Transcriptome-wide coverage<br><br>• Easy scale-up<br><br>• Sequencing data analysis | • Single-cell/single-molecule<br><br>• High spatial resolution (<1$\mu$m)<br><br>• Continuous spatial locations |
| **Cons** | • Fixed spatial dissection<br><br>• Low spatial resolution (~100$\mu$m)<br><br>• Not single-cell | • Coverage restricted to probes<br><br>• More difficult experiments<br><br>• Challenging data analysis |

# 10X Genomics - Visium


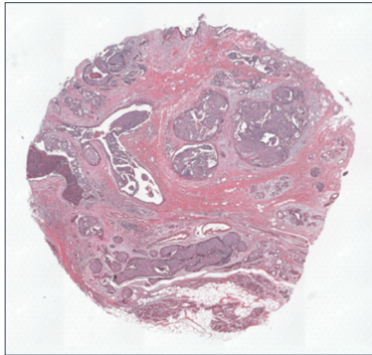
10X Genomics

# 10X Genomics - Visium



10X Genomics

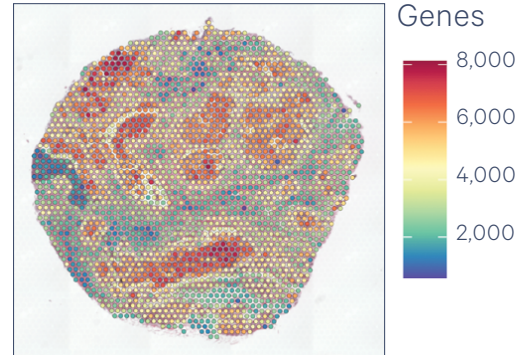# 10X Genomics - Visium

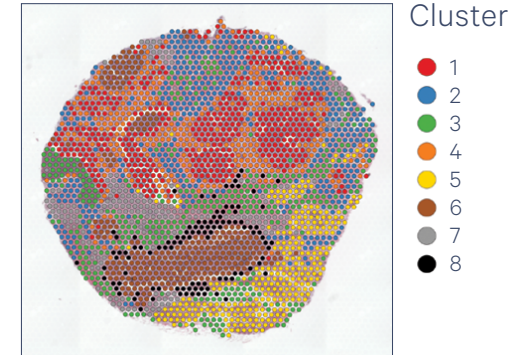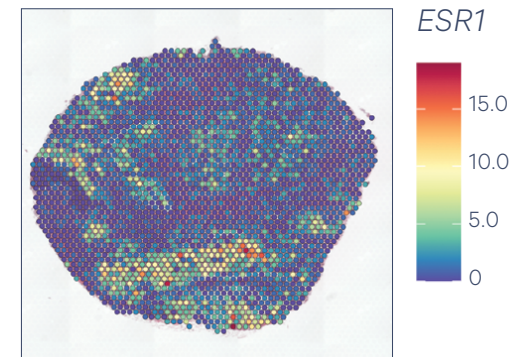**Interrogation of ~18,000 genes in a human breast ductal carcinoma in situ FFPE sample**
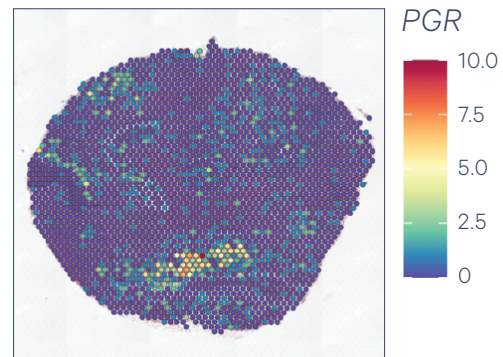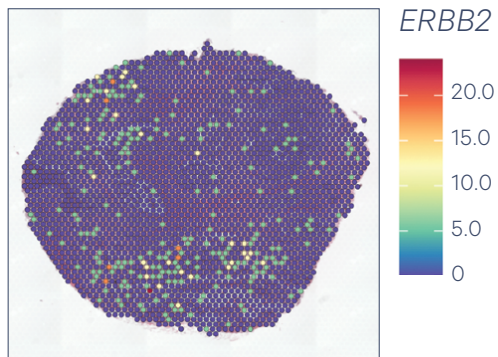


**A.** H&E

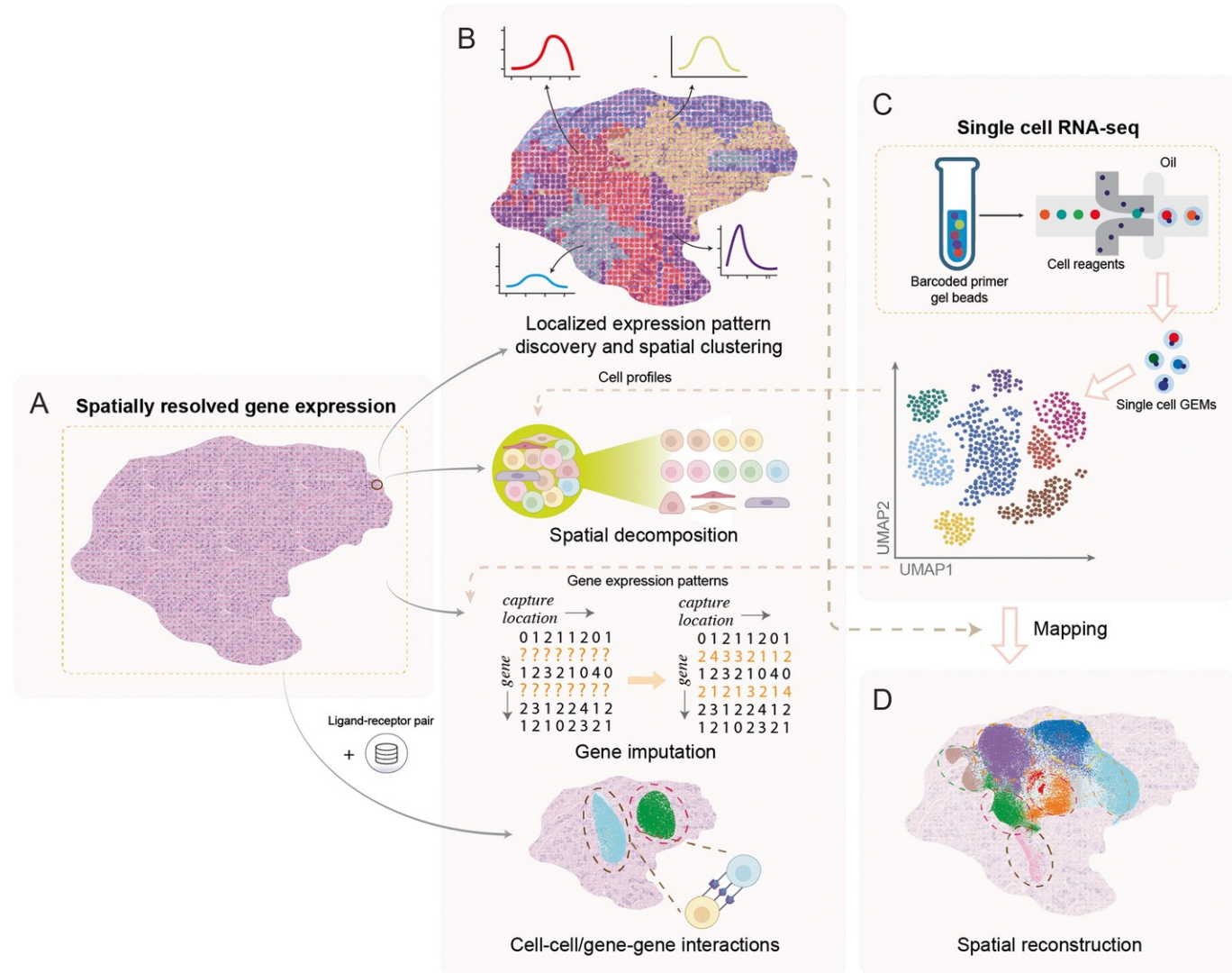**B.** Total genes

**C.** Spot clusters

**D.** Three key breast cancer biomarkers

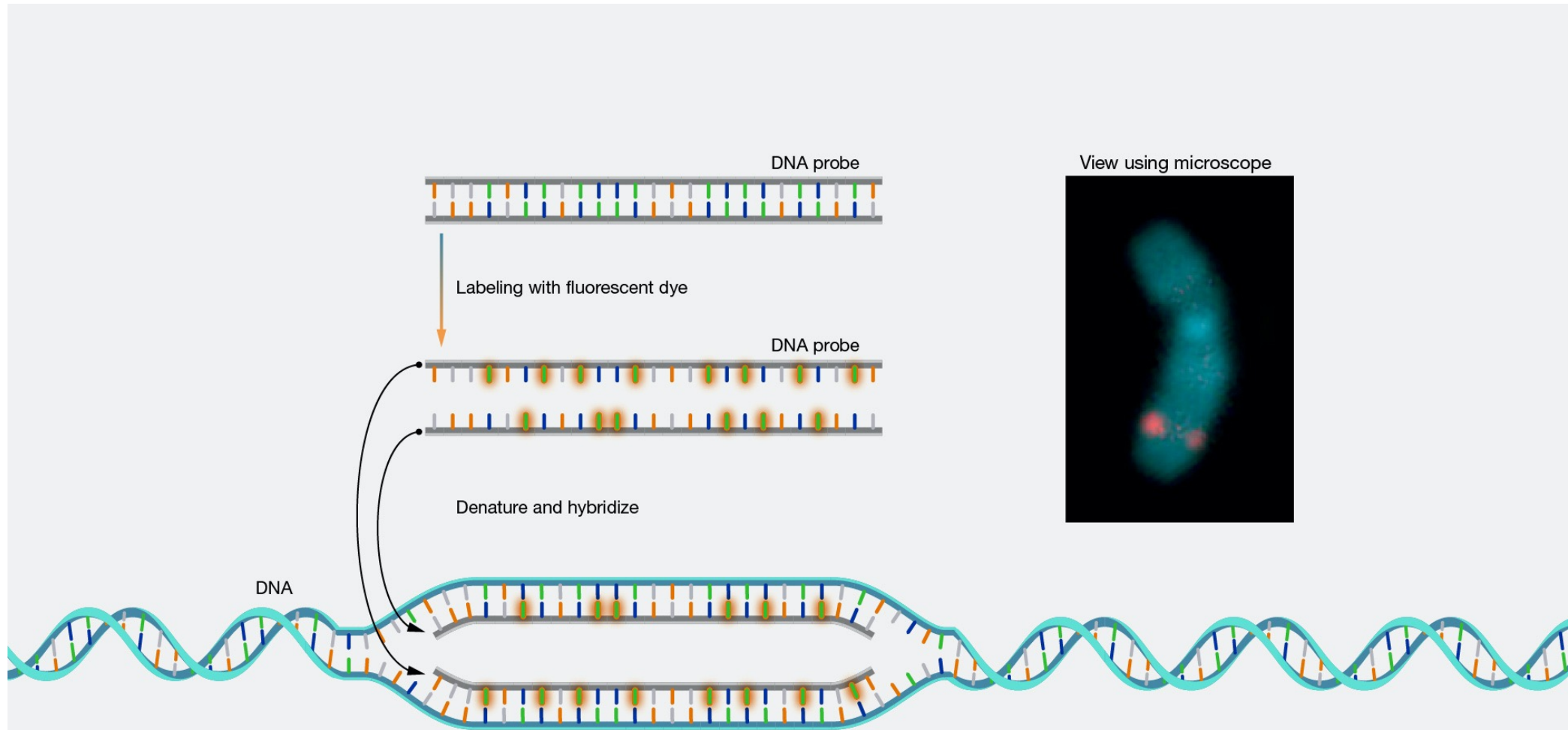10X Genomics

# Computational Problems

- Localized gene expression profiling

- Spatial clustering

- Spatial decomposition and gene imputation

- Spatial location reconstruction for scRNA-seq

- Cellular interaction or gene interaction inference

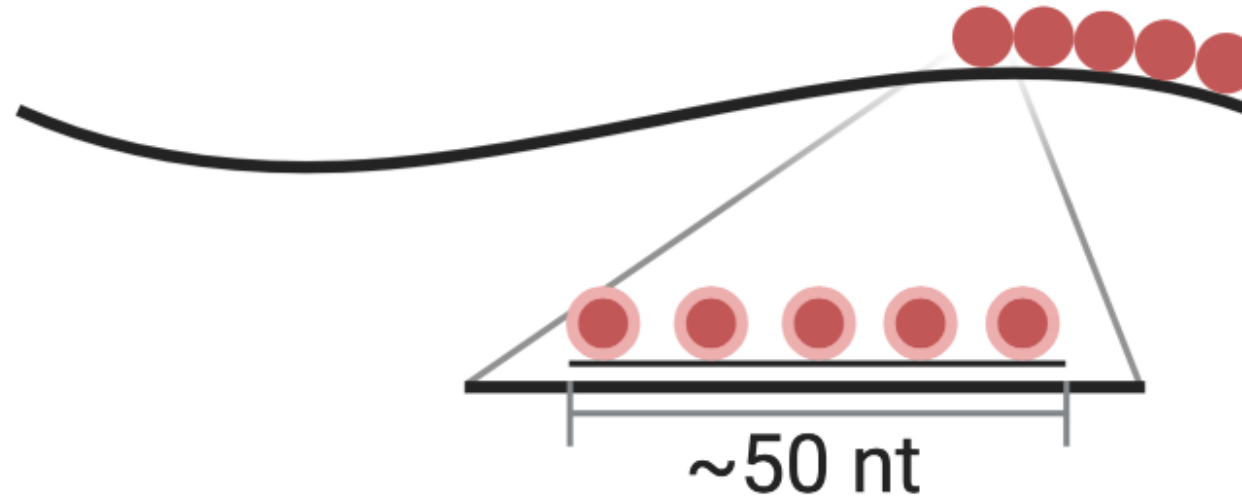Zeng et al. Genome Bio 2022

# **MERFISH**

Multiplexed Error-Robust Fluorescence In Situ Hybridization

# FISH



DNA probe

Labeling with fluorescent dye

DNA probe

Denature and hybridize

DNA

View using microscope

# Single-Molecule FISH (smFISH)

A 1998 smFISH

~50 nt

B 2008 smFISH

17-22 nt

48×

Lior Pachter Lab

# seqFISH



2014 seqFISH

seqFISH error correction

Lior Pachter Lab

# MERFISH



2015 MERFISH
100101001000000, first two rounds

Hyb. round 1 · Photo-bleaching · Hyb. round 2

Hyb. round 4 · Hyb. round 6 · Hyb. round 9

MERFISH error correction

100101001000000
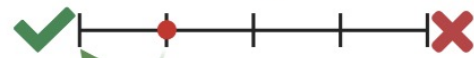Original

100101000000000
✓ ——————— ✗
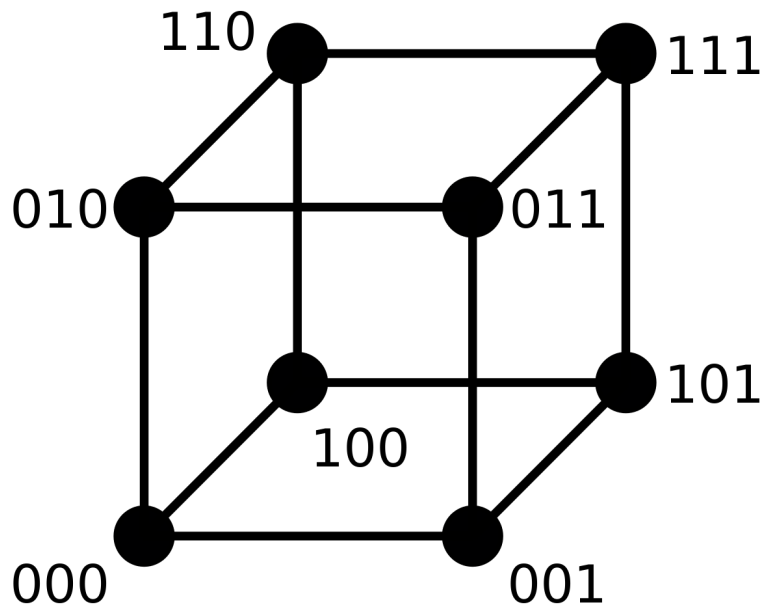Hamming distance
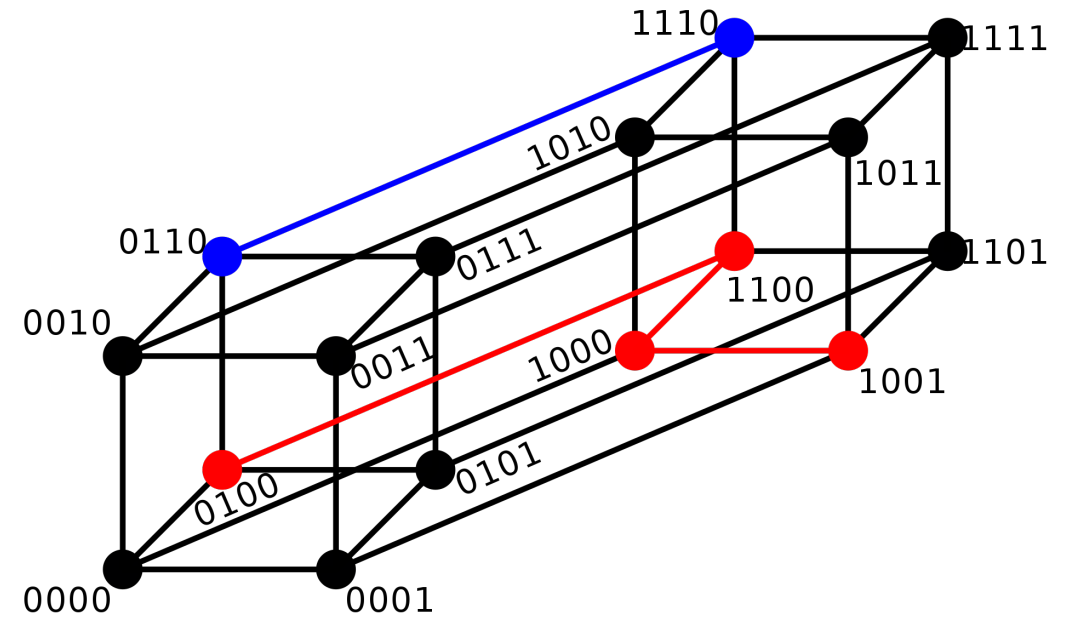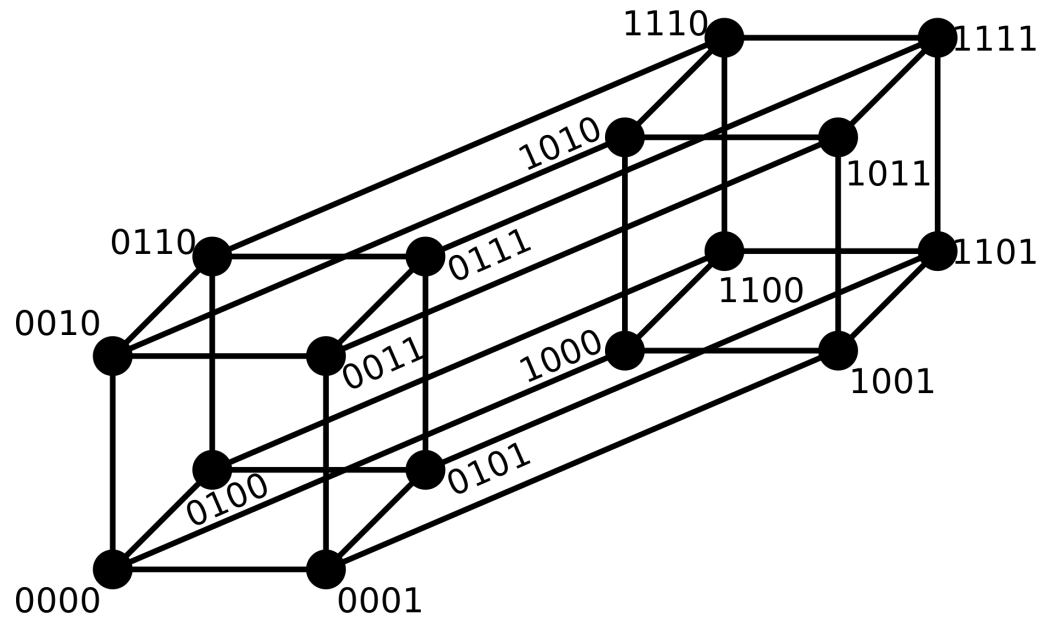
100100000000000
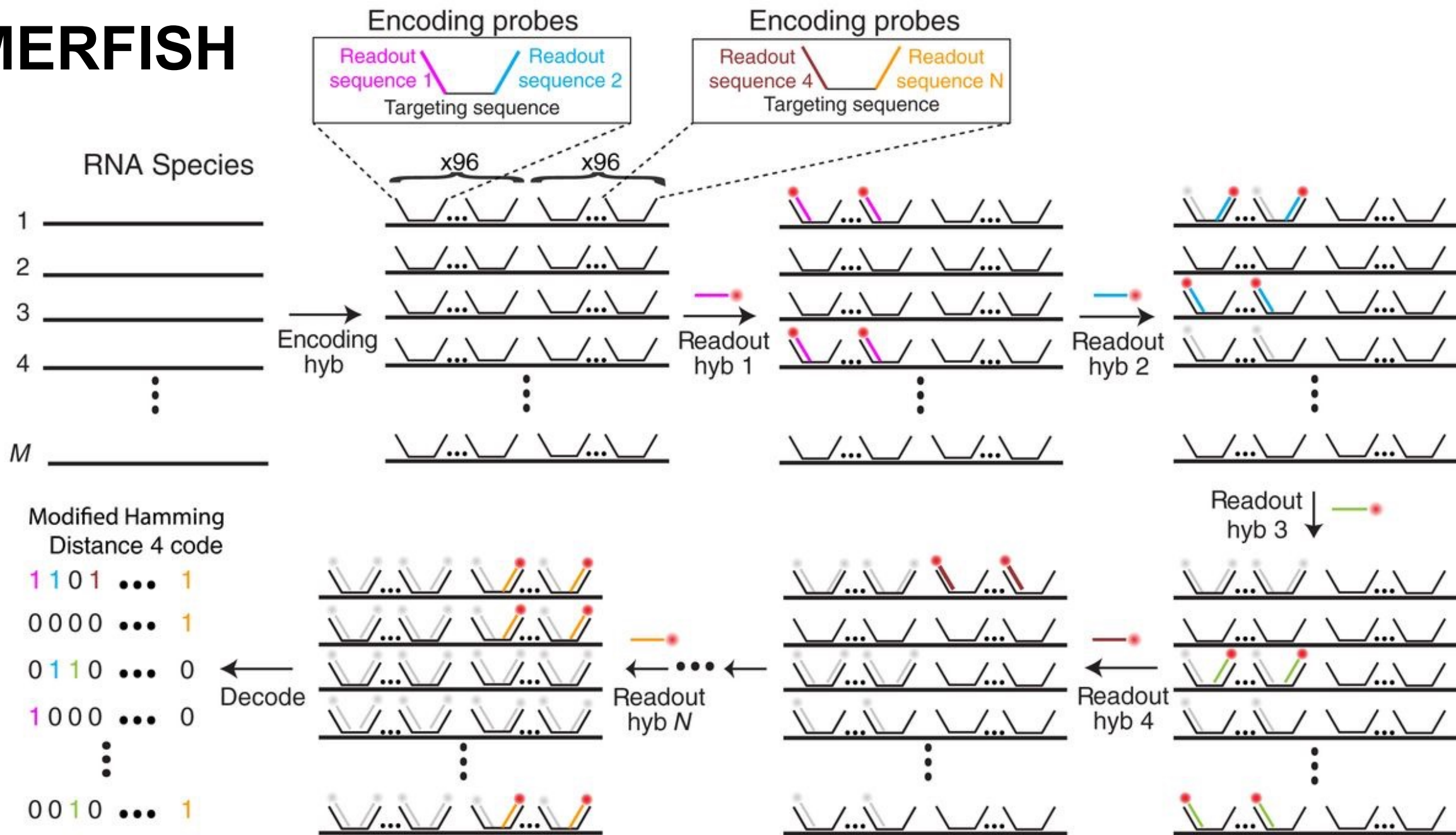✓ ——————— ✗
?
e.g. 100100000101000

# Hamming Distance



100→011 has Hamming distance 3
010→111 has Hamming distance 2

# Hamming Distance



0100→1001 has Hamming distance 3
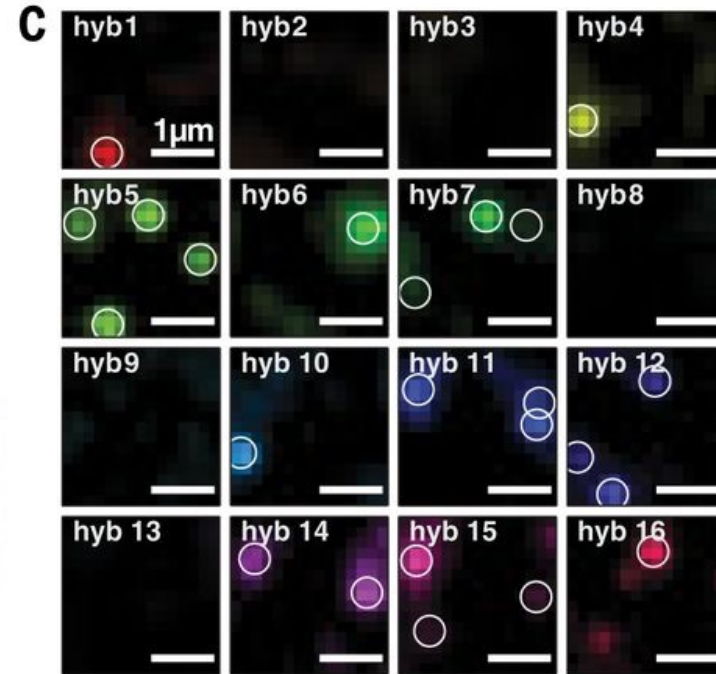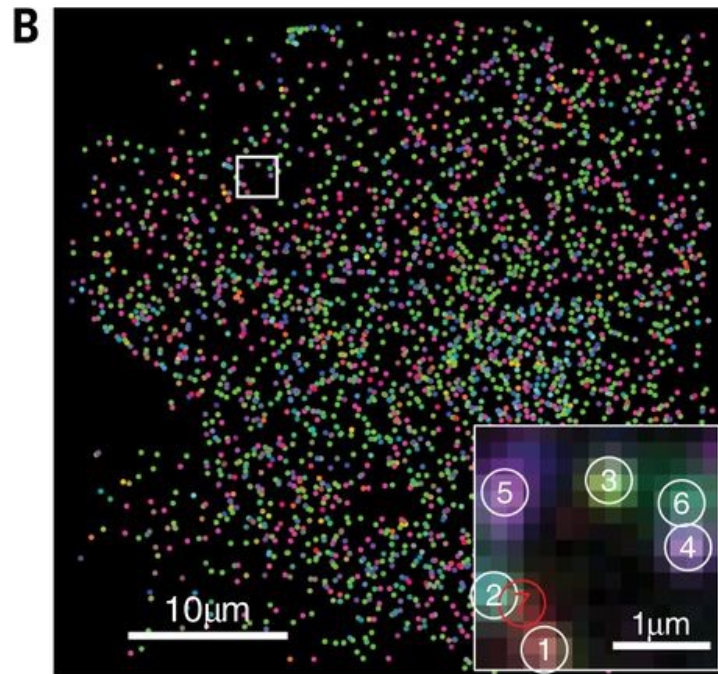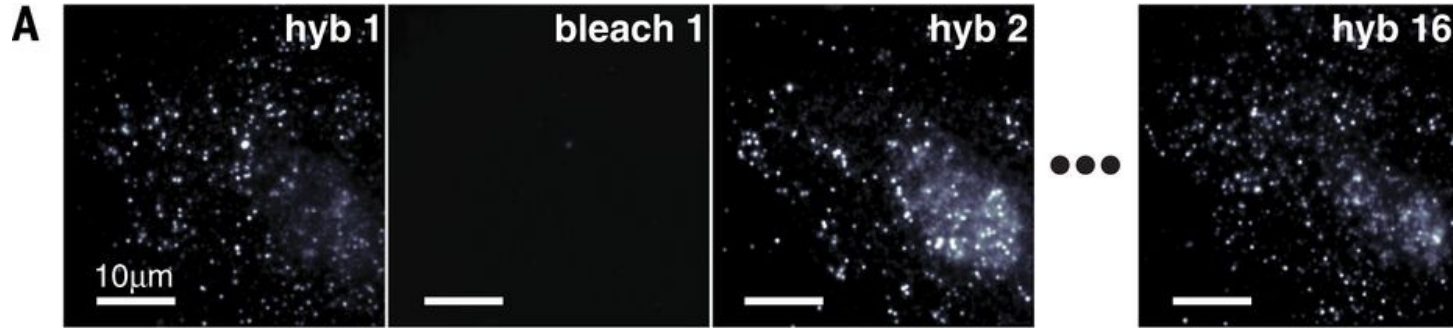0110→1110 has Hamming distance 1

# MERFISH

Chen et al. Science 2015

# MERFISH

Chen et al. Science 2015

# Gray Code and One-Hot Code

| Decimal | Binary | Gray | Decimal of Gray | One-Hot |
|---|---|---|---|---|
| 0 | 0000 | 0000 | 0 | 0000000000000001 |
| 1 | 0001 | 0001 | 1 | 0000000000000010 |
| 2 | 0010 | 0011 | 3 | 0000000000000100 |
| 3 | 0011 | 0010 | 2 | 0000000000001000 |
| 4 | 0100 | 0110 | 6 | 0000000000010000 |
| 5 | 0101 | 0111 | 7 | 0000000000100000 |
| 6 | 0110 | 0101 | 5 | 0000000001000000 |
| 7 | 0111 | 0100 | 4 | 0000000010000000 |
| 8 | 1000 | 1100 | 12 | 0000000100000000 |
| 9 | 1001 | 1101 | 13 | 0000001000000000 |
| 10 | 1010 | 1111 | 15 | 0000010000000000 |
| 11 | 1011 | 1110 | 14 | 0000100000000000 |
| 12 | 1100 | 1010 | 10 | 0001000000000000 |
| 13 | 1101 | 1011 | 11 | 0010000000000000 |
| 14 | 1110 | 1001 | 9 | 0100000000000000 |
| 15 | 1111 | 1000 | 8 | 1000000000000000 |

# One-hot encoding for DNA sequences

| | C | G | A | T | A | A | C | C | G | A | T | A | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| C | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| G | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| T | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |

# Simplex Encoding (Hadamard)

```
z
         y
              T
G
         (0,0,0)    x
                    A =   1 -1 -1
    C                C =  -1  1 -1
              A      G =  -1 -1  1
                     T =   1  1  1
```

```
AA = A⊗A =   1 -1 -1 -1  1  1 -1  1  1
AC = A⊗C =  -1  1 -1  1 -1  1  1 -1  1
AG = A⊗G =  -1 -1  1  1  1 -1  1  1 -1
AT = A⊗T =   1  1  1 -1 -1 -1 -1 -1 -1
CA = C⊗A =  -1  1  1  1 -1 -1 -1  1  1
CC = C⊗C =   1 -1  1 -1  1 -1  1 -1  1
CG = C⊗G =   1  1 -1 -1 -1  1  1  1 -1
CT = C⊗T =  -1 -1 -1  1  1  1 -1 -1 -1
GA = G⊗A =  -1  1  1 -1  1  1  1 -1 -1
GC = G⊗C =   1 -1  1  1 -1  1  1 -1  1
GG = G⊗G =   1  1 -1  1  1 -1 -1 -1  1
GT = G⊗T =  -1 -1 -1 -1 -1 -1  1  1  1
TA = T⊗A =   1 -1 -1  1 -1 -1  1 -1 -1
TC = T⊗C =  -1  1 -1 -1  1 -1 -1  1 -1
TG = T⊗G =  -1 -1  1 -1 -1  1 -1 -1 -1
TT = T⊗T =   1  1  1  1  1  1  1  1  1
```

$$p = 1 + 3k + 9(k - 1) = 12k - 8$$

# Simplex encoding reduces dimensionality

| k | naïve k-mer ($4^k$) | Simplex encoding (12k-8) |
|---|---|---|
| 4 | 256 | 40 |
| 6 | 4096 | 64 |
| 8 | 65536 | 88 |
| 10 | 1048576 | 112 |

# SELMA (Simplex-Encoded Linear Model for Accessible chromatin) improves cleavage bias estimation

Hu *et al.*, *under review. bioRxiv* 2021

Genome Biology

**SHORT REPORT**                                                    **Open Access**

Check for
updates

# Exaggerated false positives by popular differential expression methods when analyzing human population samples

Yumei Li[1†], Xinzhou Ge[2†], Fanglue Peng[3], Wei Li[1*] and Jingyi Jessica Li[2,4,5,6,7*] iD

*Correspondence:
wei.li@uci.edu; lijy03@g.
ucla.edu
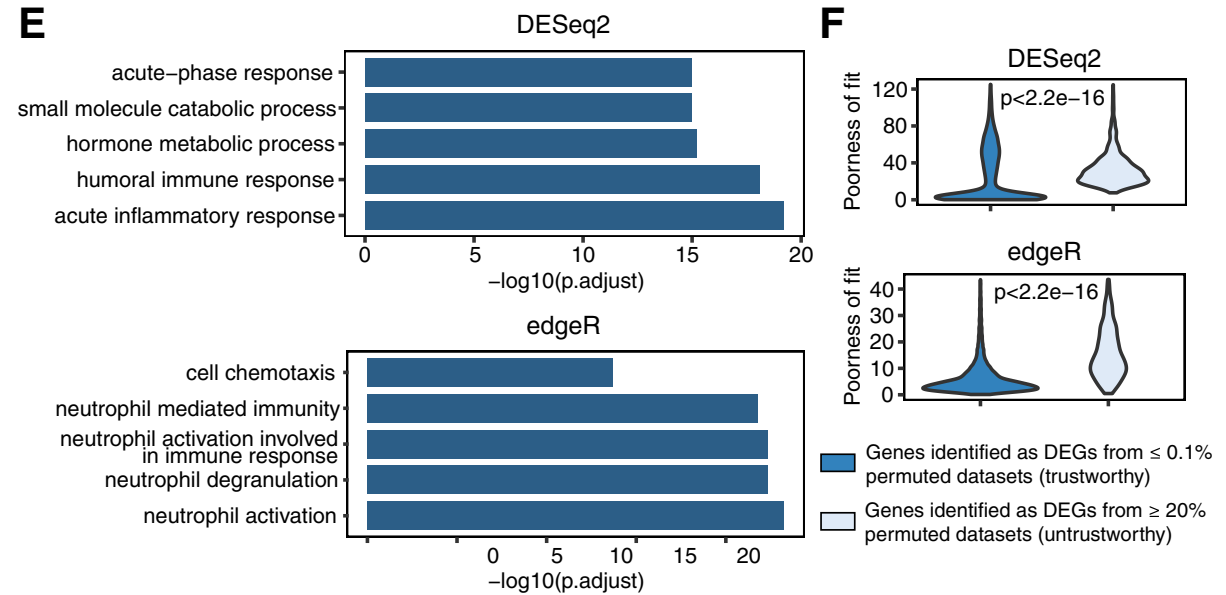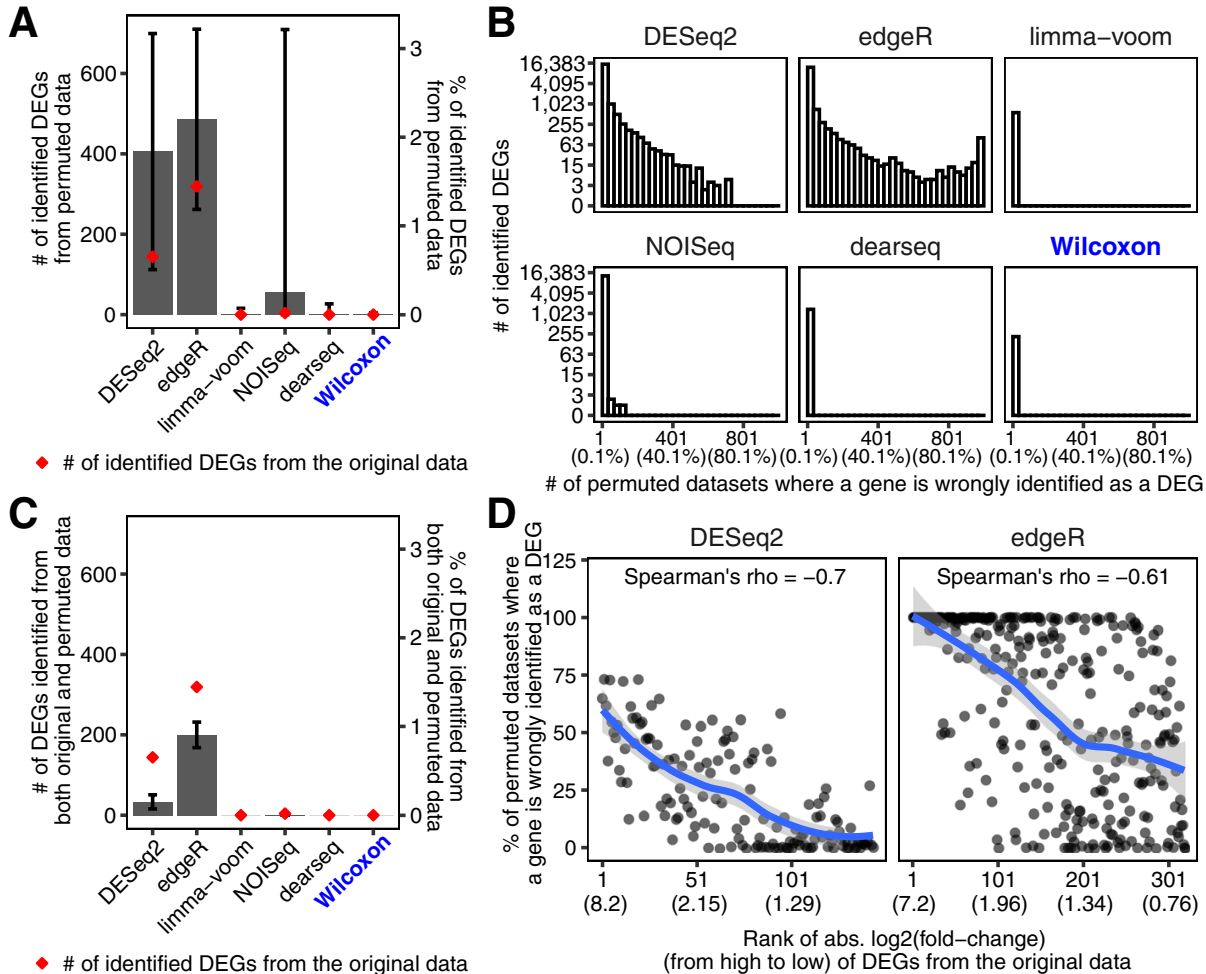†Yumei Li and Xinzhou Ge
contributed equally to this
work.
[1] Division of Computational
Biomedicine, Department
of Biological Chemistry,
School of Medicine,
University of California, Irvine,
Irvine, CA 92697, USA
[2] Department of Statistics,
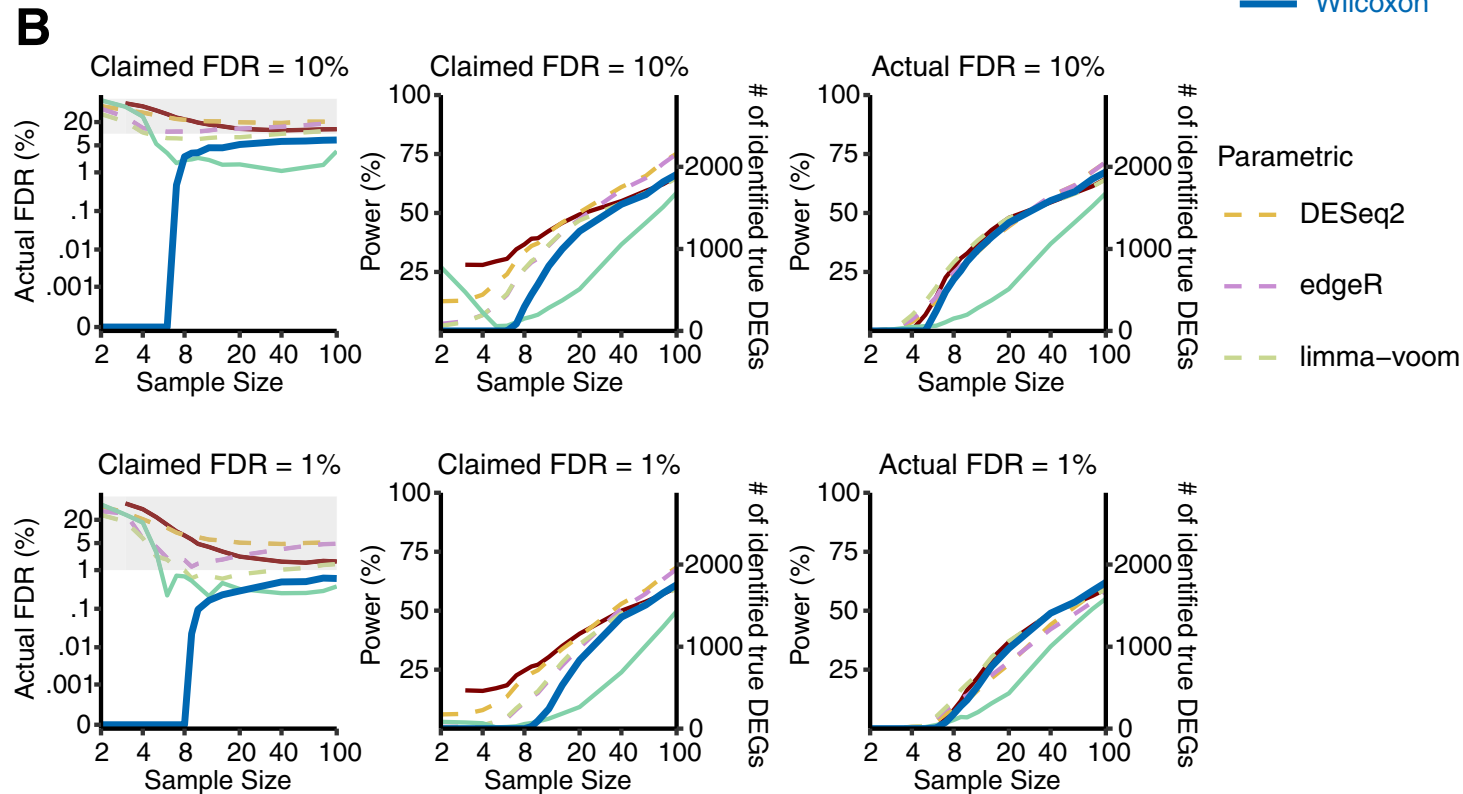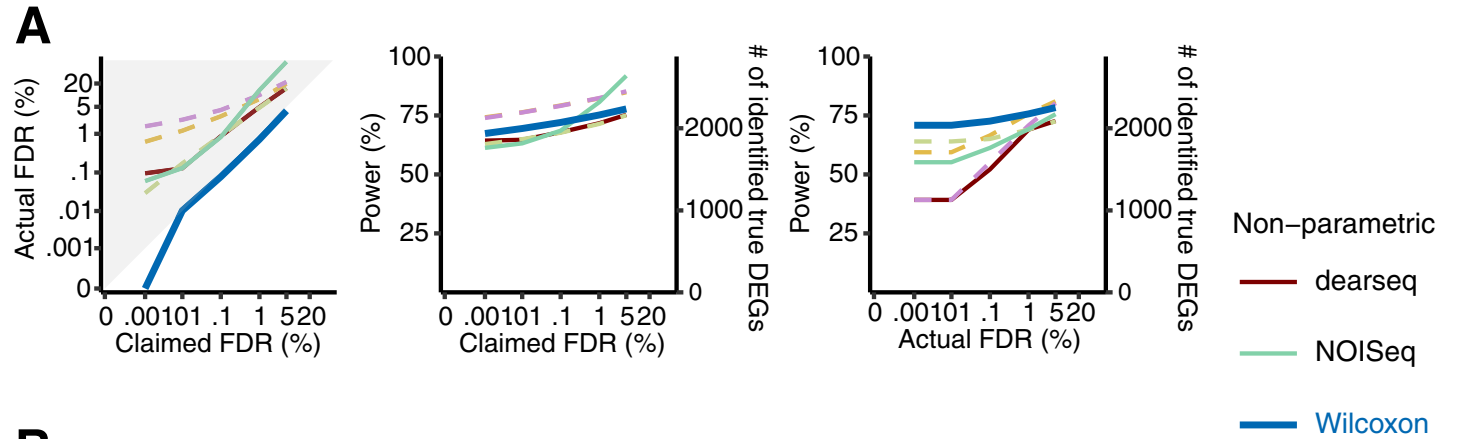University of California, Los

## Abstract

When identifying differentially expressed genes between two conditions using human population RNA-seq samples, we found a phenomenon by permutation analysis: two popular bioinformatics methods, DESeq2 and edgeR, have unexpectedly high false discovery rates. Expanding the analysis to limma-voom, NOISeq, dearseq, and Wilcoxon rank-sum test, we found that FDR control is often failed except for the Wilcoxon rank-sum test. Particularly, the actual FDRs of DESeq2 and edgeR sometimes exceed 20% when the target FDR is 5%. Based on these results, for population-level RNA-seq studies with large sample sizes, we recommend the Wilcoxon rank-sum test.
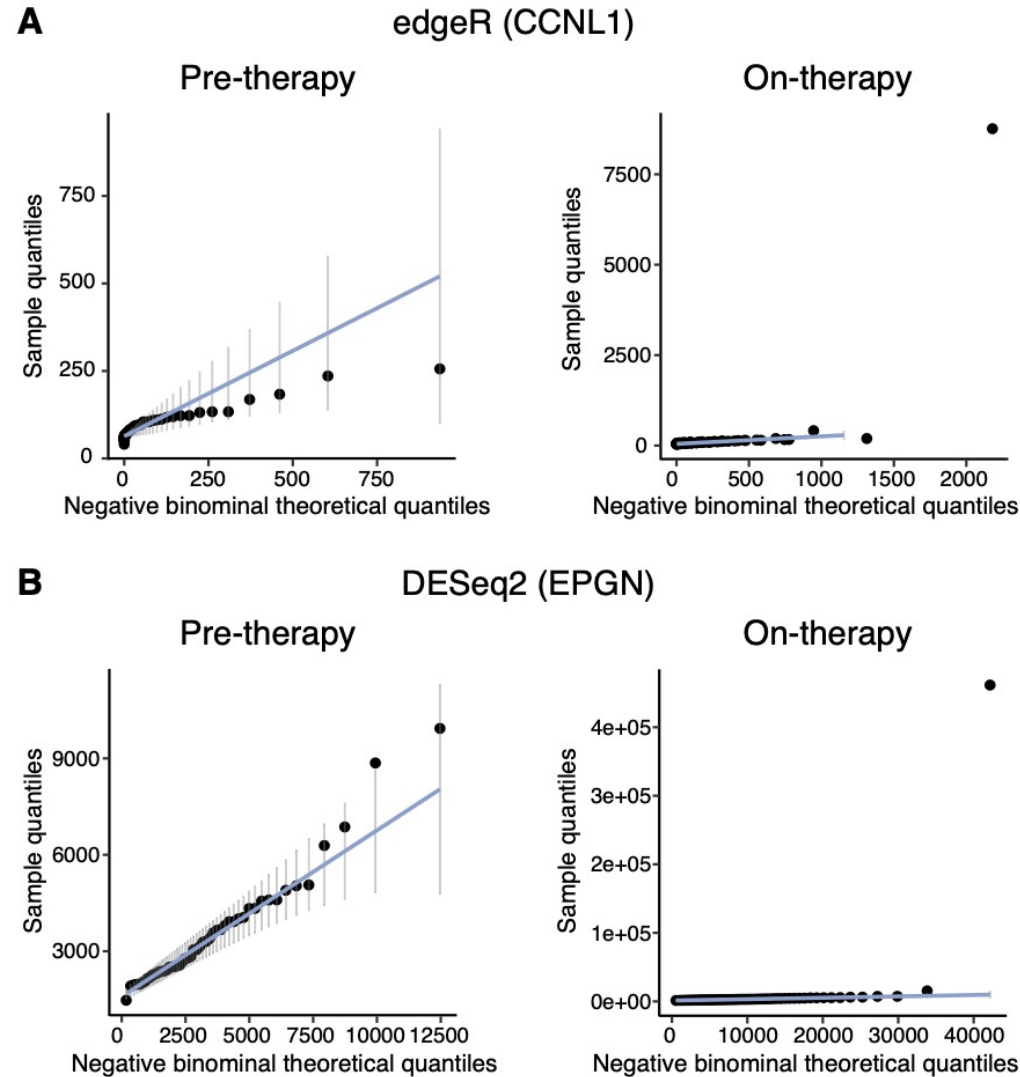
27

**A**

**B**

**C**

**D**

**E**

**F**

# Wilcoxon rank-sum test is better when sample size > 8

# Gene expression can deviate from NB distribution

# Summary

- Spatial transcriptomics techniques
- Encoding strategies
- Differential gene expression