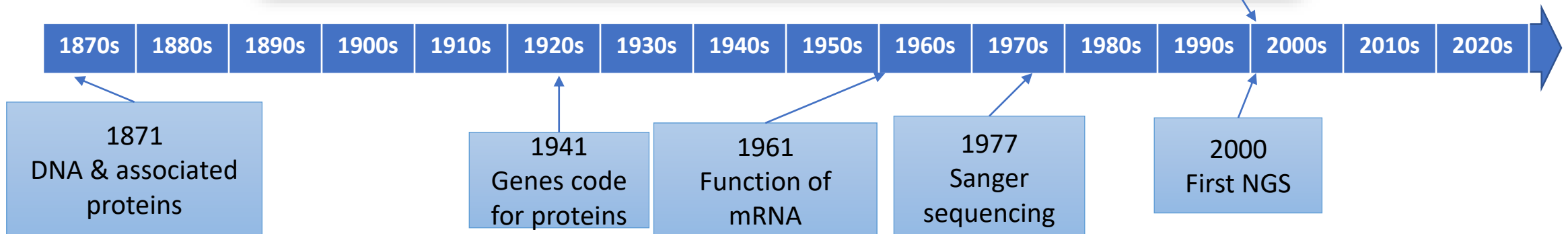# A framework for building a single-cell transcriptome treasure chest

Yiwen Wang
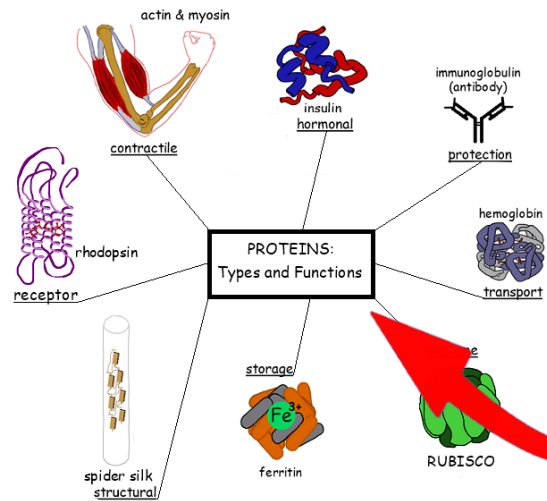
2024/05/17
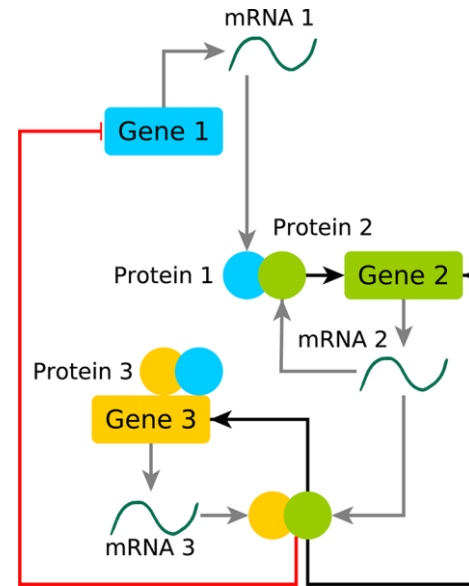
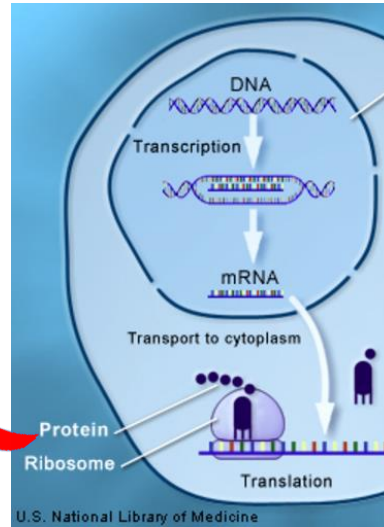# More than 20 years ago, the code was "cracked"



Genetic Code of Human Life Is Cracked by Scientists

| 1870s | 1880s | 1890s | 1900s | 1910s | 1920s | 1930s | 1940s | 1950s | 1960s | 1970s | 1980s | 1990s | 2000s | 2010s | 2020s |

1871
DNA & associated proteins

1941
Genes code for proteins

1961
Function of mRNA

1977
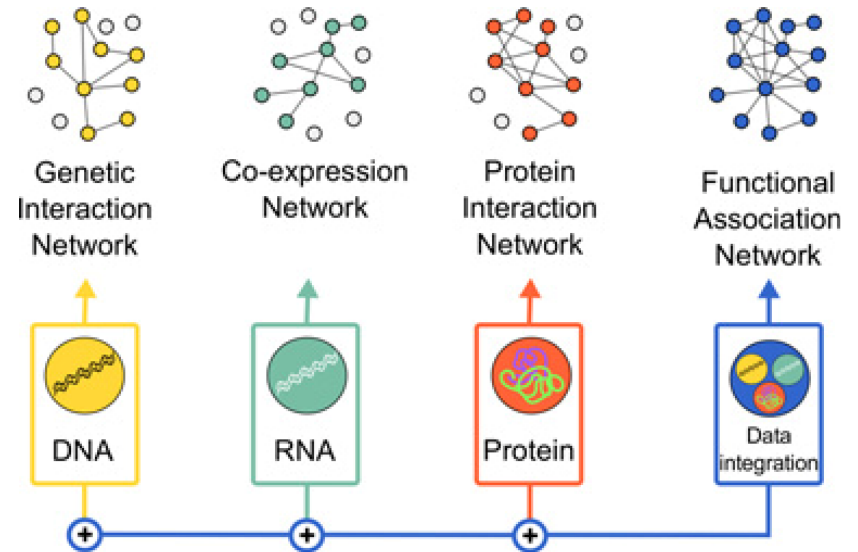Sanger sequencing

2000
First NGS
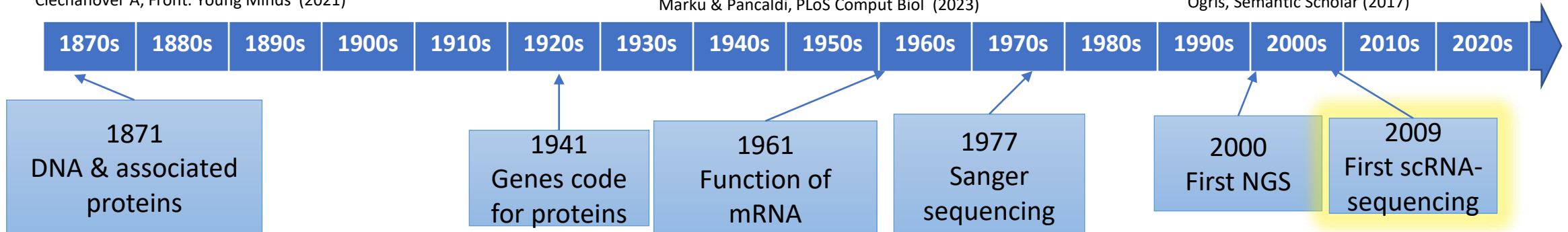
# A key for deciphering how the code is executed: Single-cell transcriptome (SCT)



Ciechanover A, Front. Young Minds (2021)

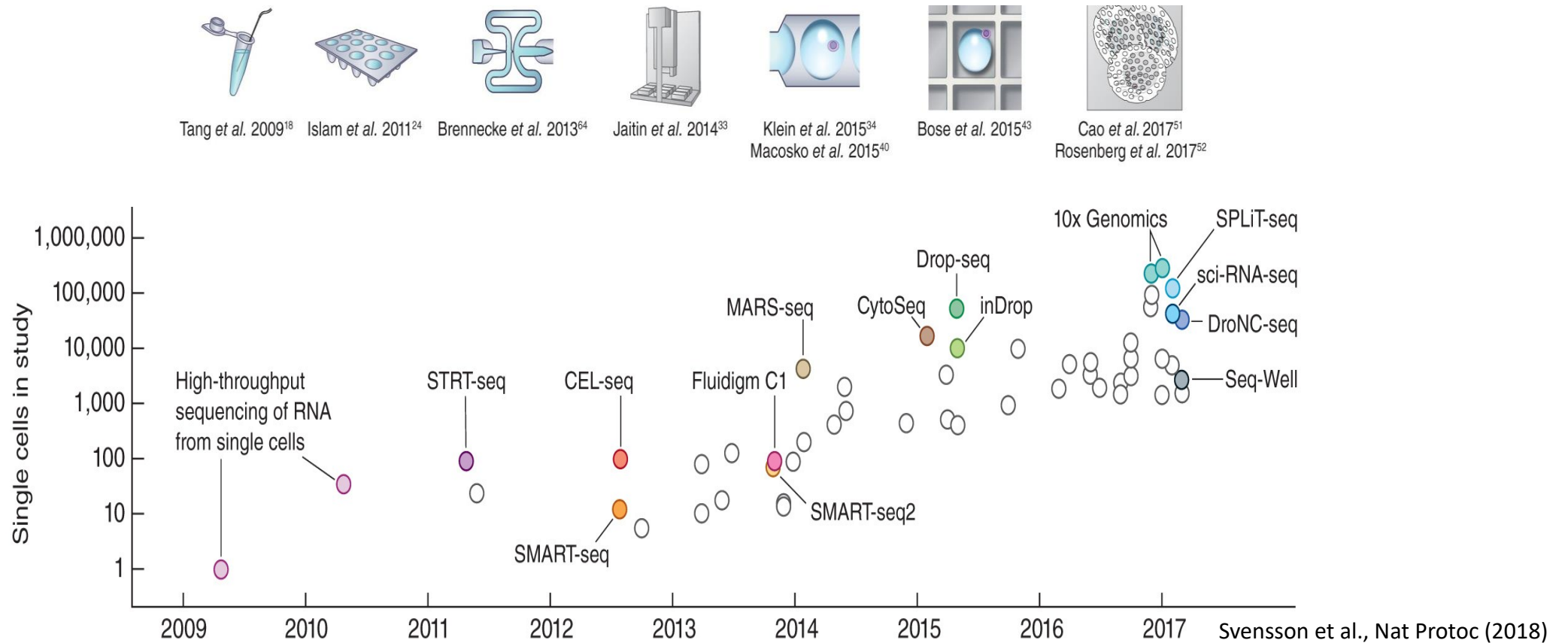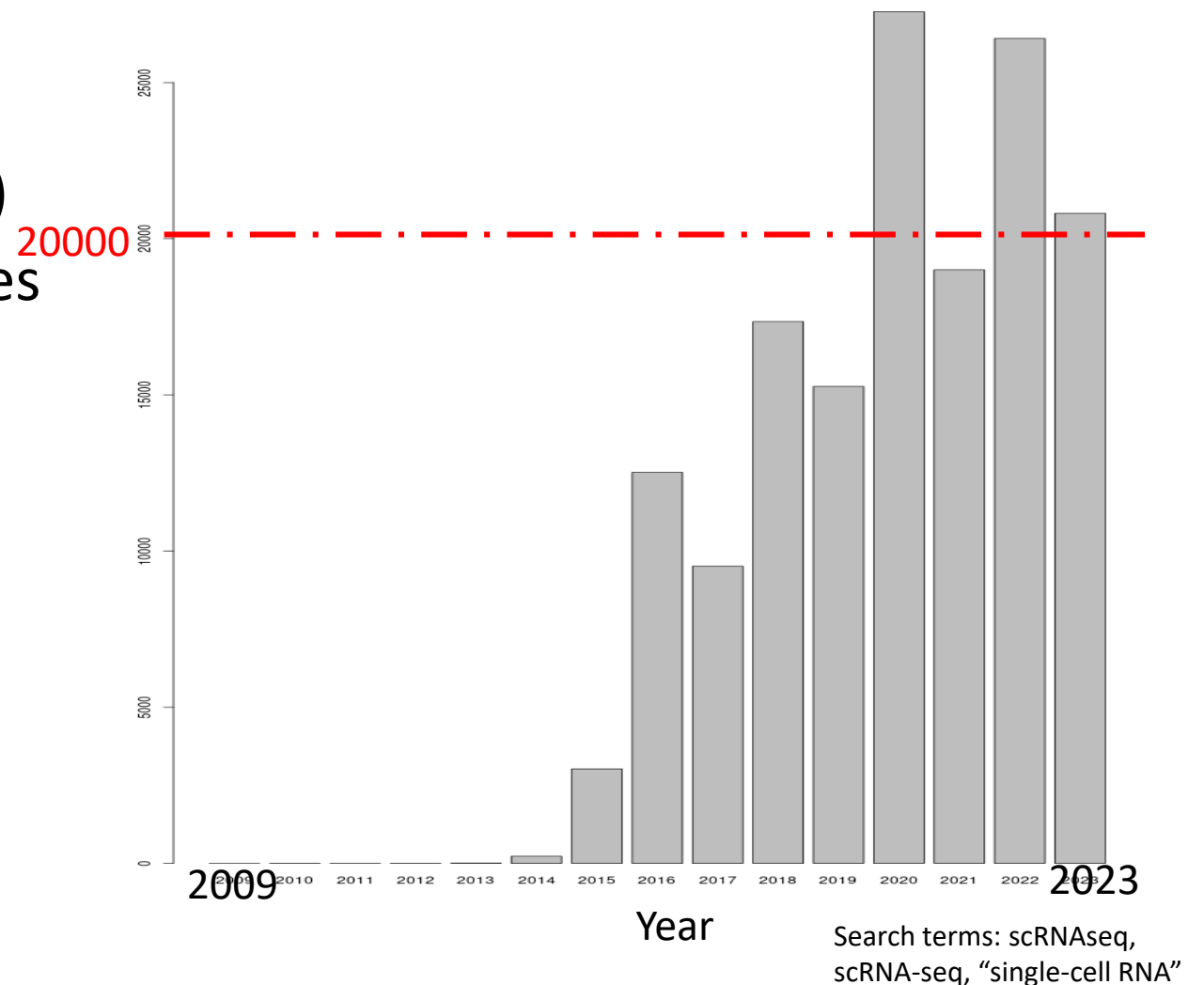Marku & Pancaldi, PLoS Comput Biol (2023)

Ogris, Semantic Scholar (2017)

| 1870s | 1880s | 1890s | 1900s | 1910s | 1920s | 1930s | 1940s | 1950s | 1960s | 1970s | 1980s | 1990s | 2000s | 2010s | 2020s |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|

1871
DNA & associated proteins

1941
Genes code for proteins

1961
Function of mRNA

1977
Sanger sequencing

2000
First NGS

2009
First scRNA-sequencing

# Rapid advances of single-cell RNA sequencing technologies since 2009



Svensson et al., Nat Protoc (2018)

1870s | 1880s | 1890s | 1900s | 1910s | 1920s | 1930s | 1940s | 1950s | 1960s | 1970s | 1980s | 1990s | 2000s | **2010s** | **2020s**

# SCT data in GDS: where the treasure's buried

- More than 20,000 scRNA-seq studies per year since 2022 uploaded to GEO DataSet (GDS)

- SCT treasure hunt: dig out values buried in GDS
  - ✓ Build curated database (treasure chest)
  - ✓ Build bioinfo tools
  - ✓ Build cell atlas
  - ✓ Benchmark
  - ✓ Generate hypothesis
  - ✓ Find supporting evidence



Search terms: scRNAseq, scRNA-seq, "single-cell RNA"

# SCT treasure hunt representative examples

| Curated database (treasure chest) | Build bioinfo tools | Cell atlas |
|---|---|---|
| IO CZ CELL×GENE DISCOVER | SEURAT | HuBMAP / Azimuth |
| Gepliver | CellTypist | THE UNIVERSITY OF TEXAS MD Anderson Cancer Center T Cell Map |
| GREIN | SCENIC | HUSCH |
| SIB IMMUCAN | | Cincinnati Children's ToppCell |
| scLiver DB | | HCCDB |
| TISCH2 | CellChat | |

**METHOD DETAILS**

**Single-cell RNA-seq data source**

To have a comprehensive understanding of immune cells in different repertoire COVID-19 single-cell RNA-seq datasets of multiple compartments, including pe clear cells, bronchoalveolar lavage and lung biopsy, which in total covered ove mild/moderate, 42 severe and 2 convalescent COVID-19 patients. More details c:

**Integration of PBMC datasets and BAL datasets using reciproc**

We input raw count files of 5 preprocessed PBMC datasets into Seurat a

Jin et al., iScience, 2021

**SCENIC runs on the different data sets.** SCENIC was run on all the data sets using the expression matrices provided by the authors (downloaded from GEO or the authors' website), includ-ing only the cells that passed their quality control, and the default gene filtering for GENIE3 (which in all these data sets resulted in
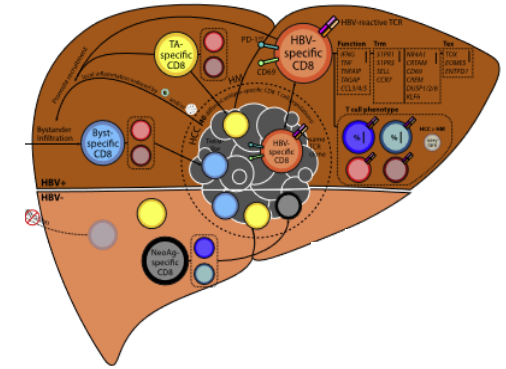
Aibar et al., Nature Methods, 2017

**Lung scRNA-seq dataset atlas.** Nineteen datasets profiling human lung samples using scRNA-seq were downloaded from publicly avail-able sources (links for each source dataset are provided in Supple-mentary Table 2). Low-quality cells were filtered using uniform quality control thresholds; cells with RNA counts between 300 and 100,000 and with mitochondrial read percentages below 20% were retained.

Hao et al., Nature Biotechnology, 2023

# Start SCT treasure hunt for our research

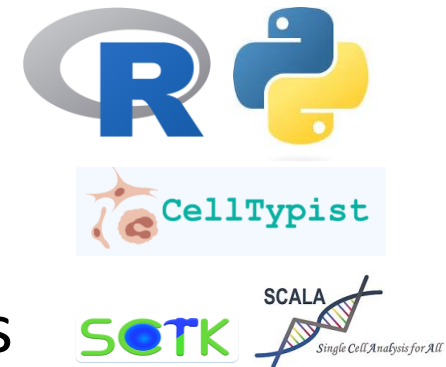| Define | |
|---|---|
| Research topic | Goal |
| scRNA-seq Hepatocellular carcinoma human liver tissue | T-cell responses in HBV vs. non-HBV |



## Data from GDS

- Query on website
- Query by scripts (R, Python, bash…)
- Download

## Pipeline for

- ✓ Re-analysis
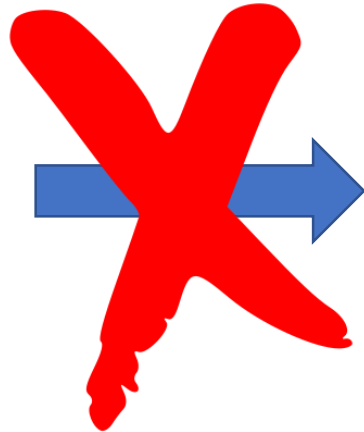- ✓ Benchmark
- ✓ Meta-analysis
- ✓ Integrate and investigate

# Obstacles to SCT treasure hunt

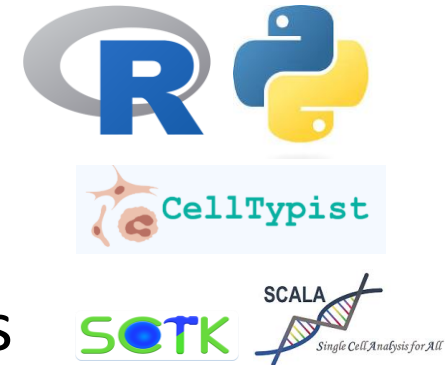| Define | |
|---|---|
| Research topic | Goal |
| scRNA-seq<br>Hepatocellular carcinoma<br>human liver tissue | T-cell responses in HBV vs. non-HBV |

## Data from GDS

- Diverse file formats
- Messy annotations

## Pipeline for

- ✓ Re-analysis
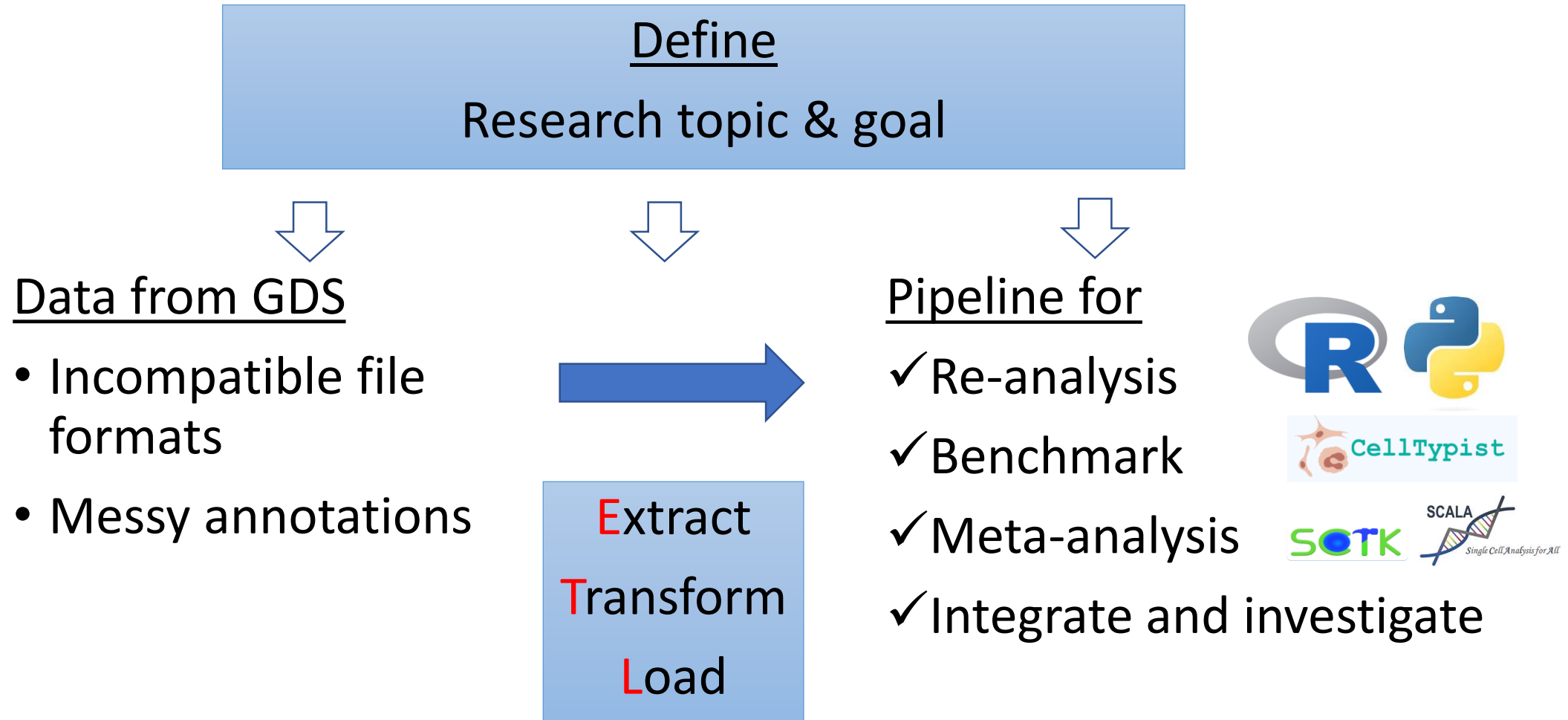- ✓ Benchmark
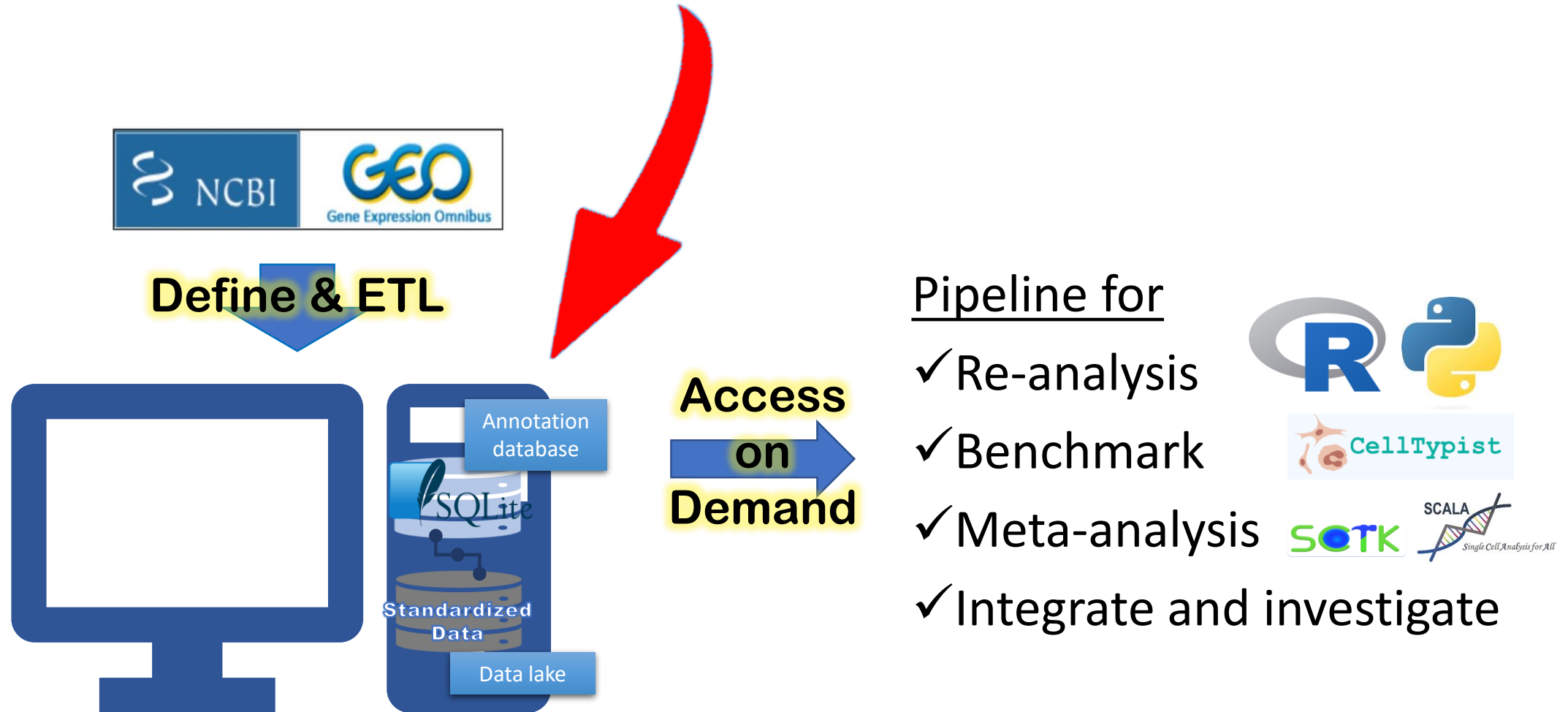- ✓ Meta-analysis
- ✓ Integrate and investigate

# Peek into SCT treasure chests built by others

| Define | |
|---|---|
| **Research topic** | **Goal** |
| scRNA-seq Hepatocellular carcinoma human liver tissue | T-cell responses in HBV vs. non-HBV |

| **Curated database** | **HCC & human liver tissue** | **HBV vs. non-HBV** |
|---|---|---|
| CZ CELL×GENE DISCOVER | 0 | N/A |
| SIB IMMUCAN | 4 | N/A |
| GREIN | 3 | N/A |
| GepLiver | 15 | N/A |
| sc Liver DB | 6 | N/A |

# ETL is critical for SCT treasure hunt

# A framework for building SCT treasure chest



Define  Extract  Transform  Load

✓No server needed
✓Only R programming required
✓Set up in a few days

# Pipeline for building SCT treasure chest

# On-demand access for downstream pipelines



HBV, HCC, liver, healthy margin, T cell

SQLiteStudio

https://inloop.github.io/sqlite-viewer/

Viewer in your browser
(in progress)

SQLite

Standardized Data

Annotation database

Data lake

## Pipeline for

✓Re-analysis

✓Benchmark

✓Meta-analysis

✓Integrate and investigate

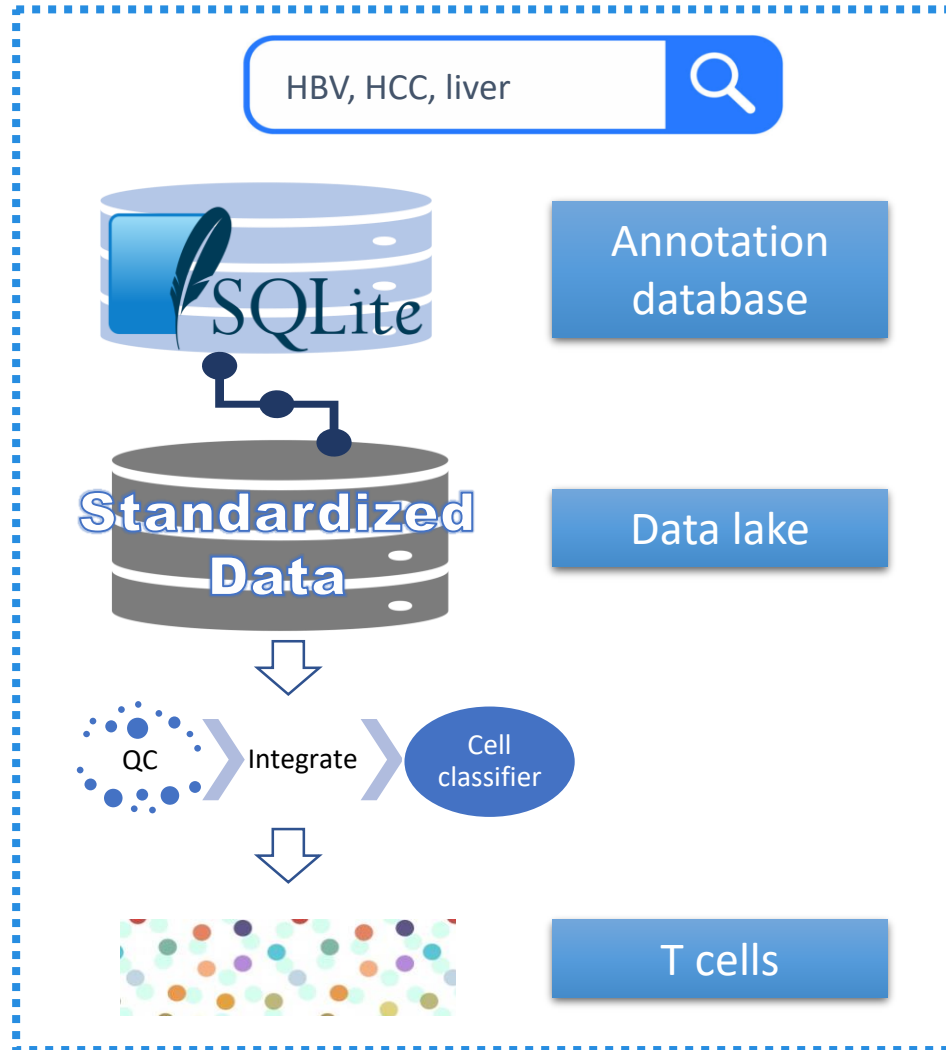# Use case: a SCT treasure chest for studying HBV-related HCC



**61 datasets found on GDS**
queried using keywords:
scRNA-seq & (HCC | HBV)

**167 human subjects**
**in 23 auto-processed datasets**
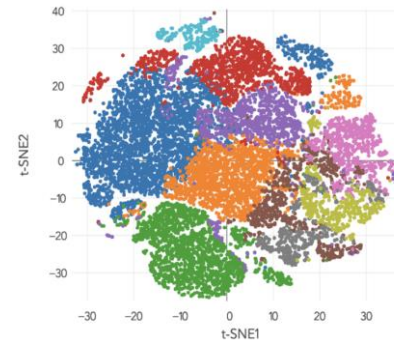grouped by sources and diseases
(4,088,738 cells)

**3,590,144 cells**
**from liver samples**
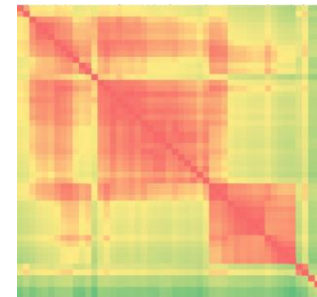grouped by available disease labels

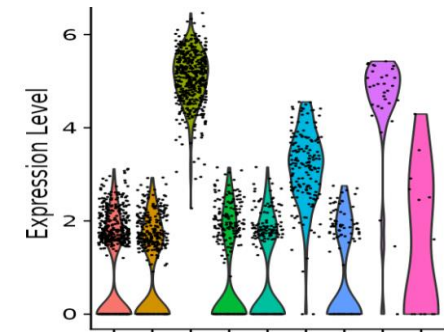# Tap in SCT treasure chest for T-cell responses in HBV-related HCC
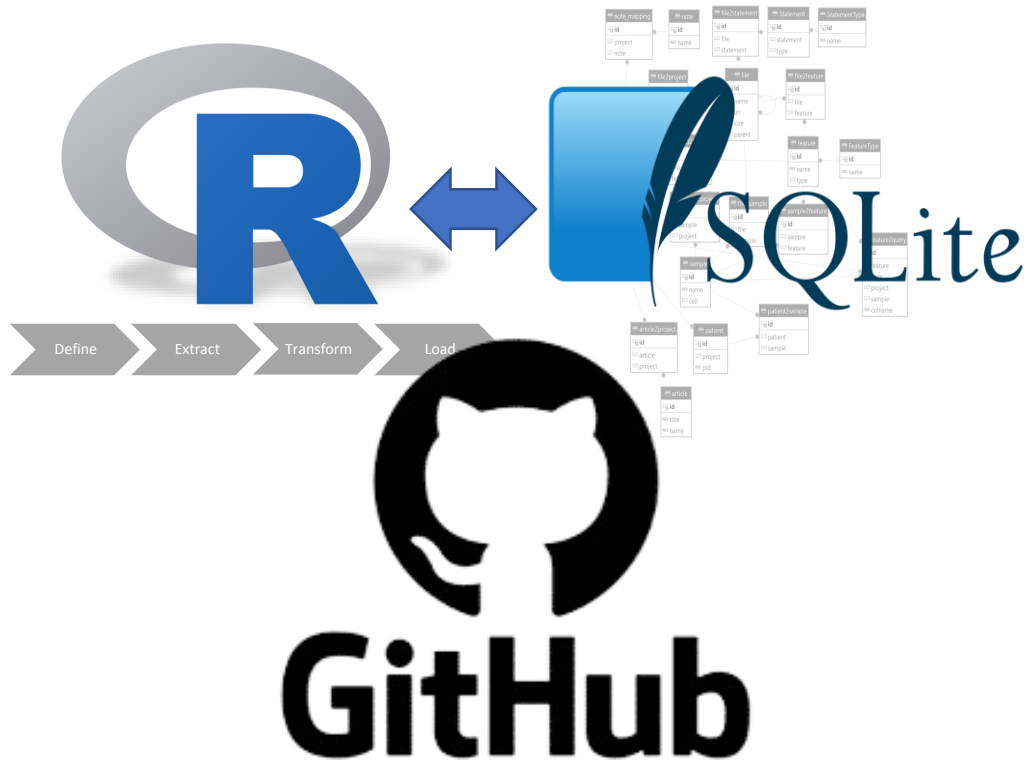


Note: Figures are for illustration purposes only.

# The framework and our curated database will be available online soon

# Contact us if you are interested

- Have a taste with our HBV-related HCC treasure chest
- Build a customized SCT treasure chest for your own research
- Suggest features that might be useful for you
- Join the project

# Thanks for your attention!