

Lecture 5:

Transcriptomics and Single-Cell Analysis

From bulk to single-cell • Cell atlas projects • Resolution revolution

Introduction to Biomedical Data Science

Lecture Contents

Part 1: Bulk RNA-seq Analysis

Part 2: Single-Cell Technologies

Part 3: Advanced Methods and Integration

Part 1/3:

Bulk RNA-seq

- Expression profiling
- Differential analysis
- Pathway enrichment
- Time series

RNA-seq Workflow

Experimental Design

Plan your study carefully with appropriate controls and biological questions

Replication Strategies

Biological replicates (≥ 3) are essential for statistical power

Batch Effect Prevention

Randomize sample processing to avoid confounding variables

Power Analysis

Determine sample size needed to detect biological effects

Cost Optimization

Balance sequencing depth and sample number for your budget

 **Key Recommendation**

More biological replicates with moderate depth > few samples with deep sequencing



Library Preparation Methods

PolyA Selection

Enriches mRNA by capturing poly-adenylated transcripts

Ribosomal Depletion

Removes rRNA to capture all RNA types including non-coding

Strand Specificity

Preserves information about which DNA strand was transcribed

UMI Incorporation

Unique Molecular Identifiers enable accurate quantification

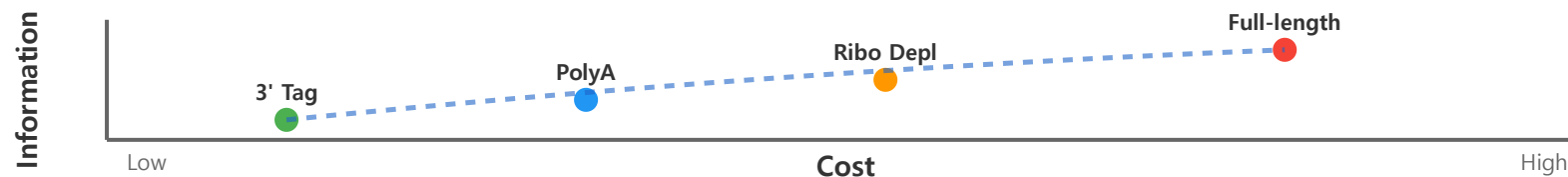
3' Tag-seq

Sequences only 3' ends - cost effective for counting

Full-length Coverage

Complete transcript coverage for isoform analysis

Trade-off: **Cost vs. Information Content**



Normalization Methods

RPKM/FPKM Issues

Reads/Fragments Per Kilobase Million - biased by composition

TPM Calculation

Transcripts Per Million - better for comparison

DESeq2 Normalization

Median-of-ratios method for differential expression

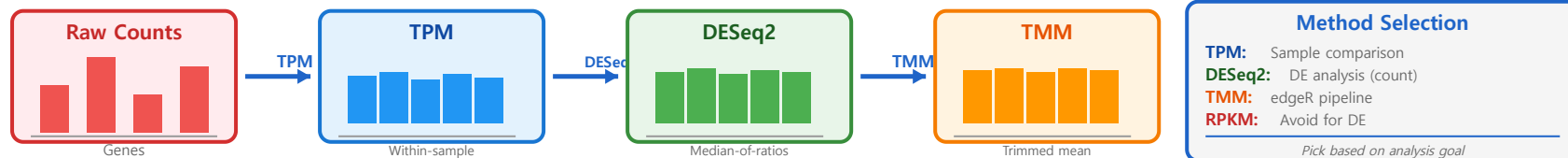
TMM Method

Trimmed Mean of M-values - robust to outliers

Batch Correction

ComBat, limma removeBatchEffect for technical variation

💡 Choose normalization method based on your downstream analysis goals



Differential Expression

Statistical Models

Account for biological and technical variability

Negative Binomial

Models count data with overdispersion

Fold Change Thresholds

Typically $|\log_2\text{FC}| > 1$ for biological significance

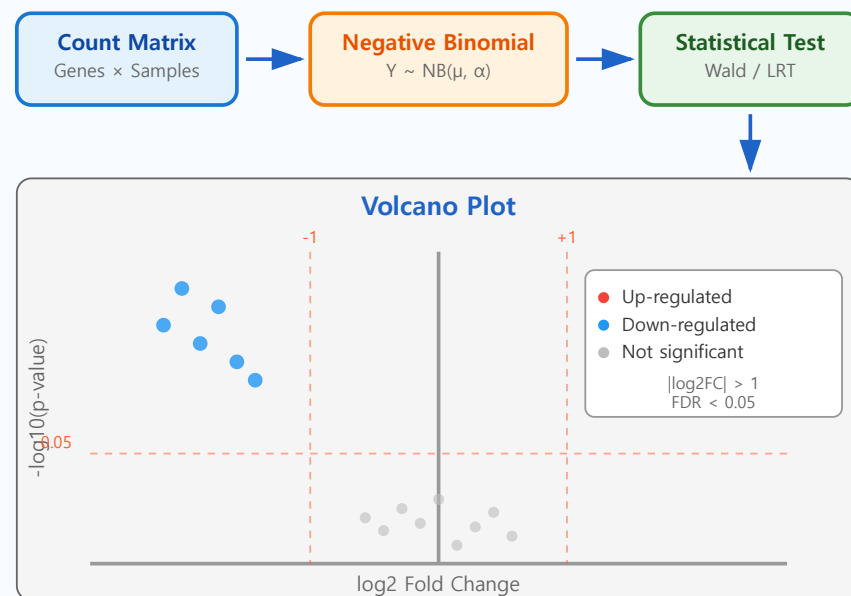
FDR Control

False Discovery Rate < 0.05 for multiple testing

Volcano Plots

Visualize FC vs statistical significance

DE Analysis Pipeline



💡 Balance statistical significance with biological relevance

Statistical Testing

RNA-seq Statistical Methods Comparison

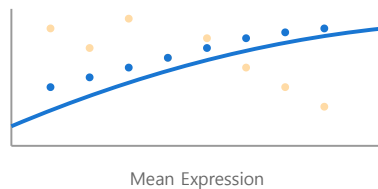
DESeq2

$$Y \sim \text{NB}(\mu, \alpha)$$

Key Features:

- Shrinkage estimation of dispersion
- Size factor normalization
- Wald test / LRT

Dispersion Shrinkage



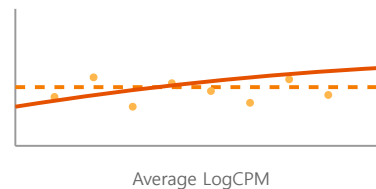
edgeR

$$Y \sim \text{NB}(\mu, \phi)$$

Key Features:

- Empirical Bayes methods
- TMM normalization
- Quasi-likelihood F-test

BCV Plot



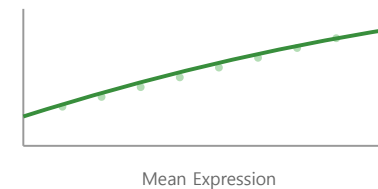
limma-voom

$$\log(Y) \sim N(\mu, \sigma^2)$$

Key Features:

- Transform to log-space
- Precision weights (voom)
- Linear modeling

Mean-Variance Trend

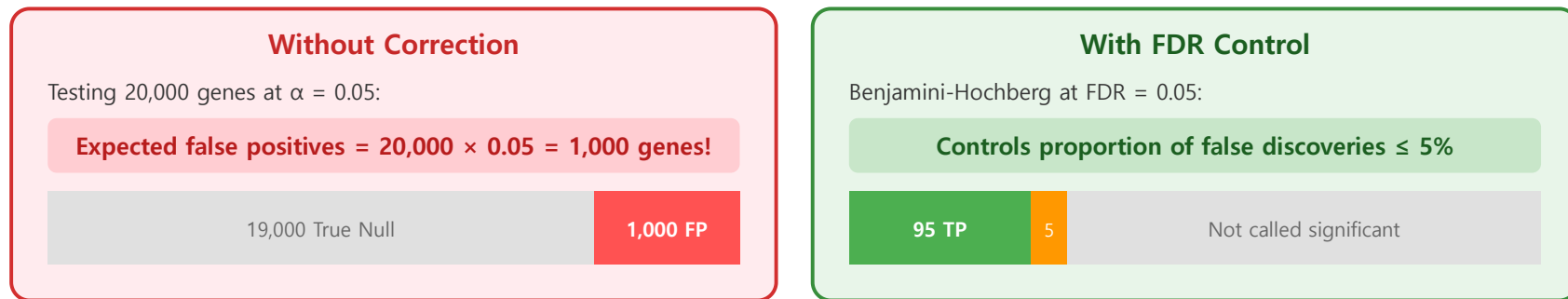


Performance: All methods perform similarly with proper use • Choose based on experimental design and analysis goals

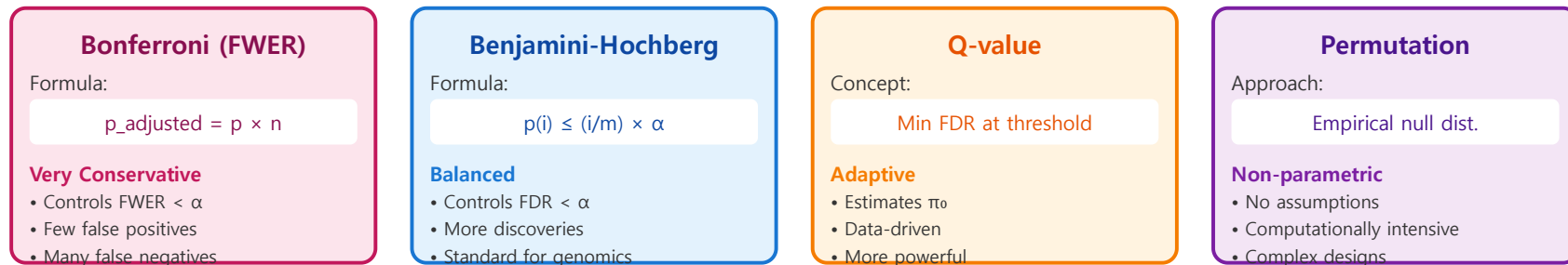
💡 DESeq2 and edgeR are most widely used and well-validated

Multiple Testing Correction

The Multiple Testing Problem



Correction Methods

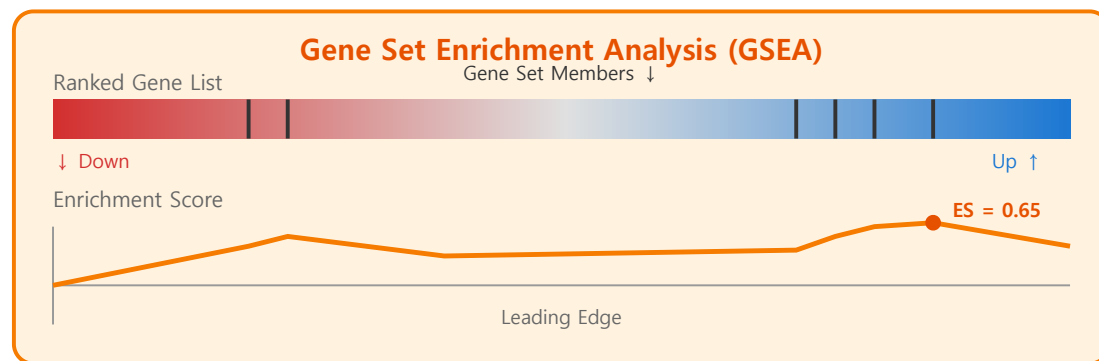
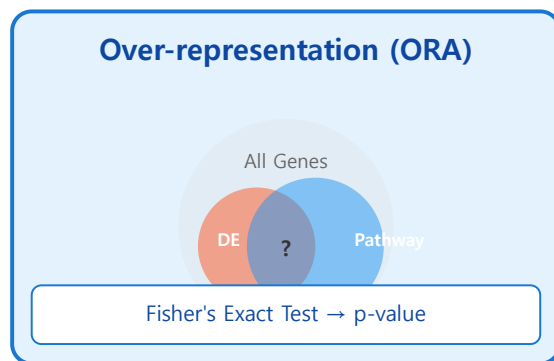


Trade-off: Strict control (Bonferroni) → Few discoveries | Relaxed control (FDR) → More discoveries with controlled error

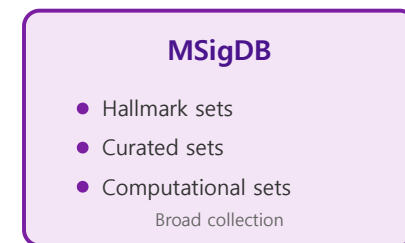
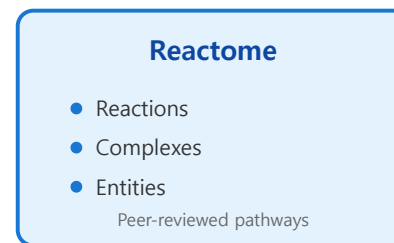
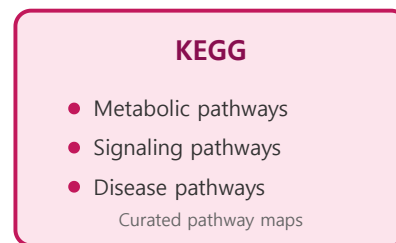
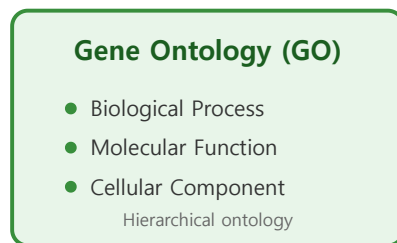
💡 Testing 20,000 genes requires careful multiple testing correction

Pathway Analysis

From Genes to Biological Pathways



Pathway Databases



ORA: Simple but threshold-dependent | GSEA: Uses all genes, more powerful | Choose database based on biological question

💡 Pathways provide biological context for gene expression changes

Part 2/3:

Single-Cell Technologies

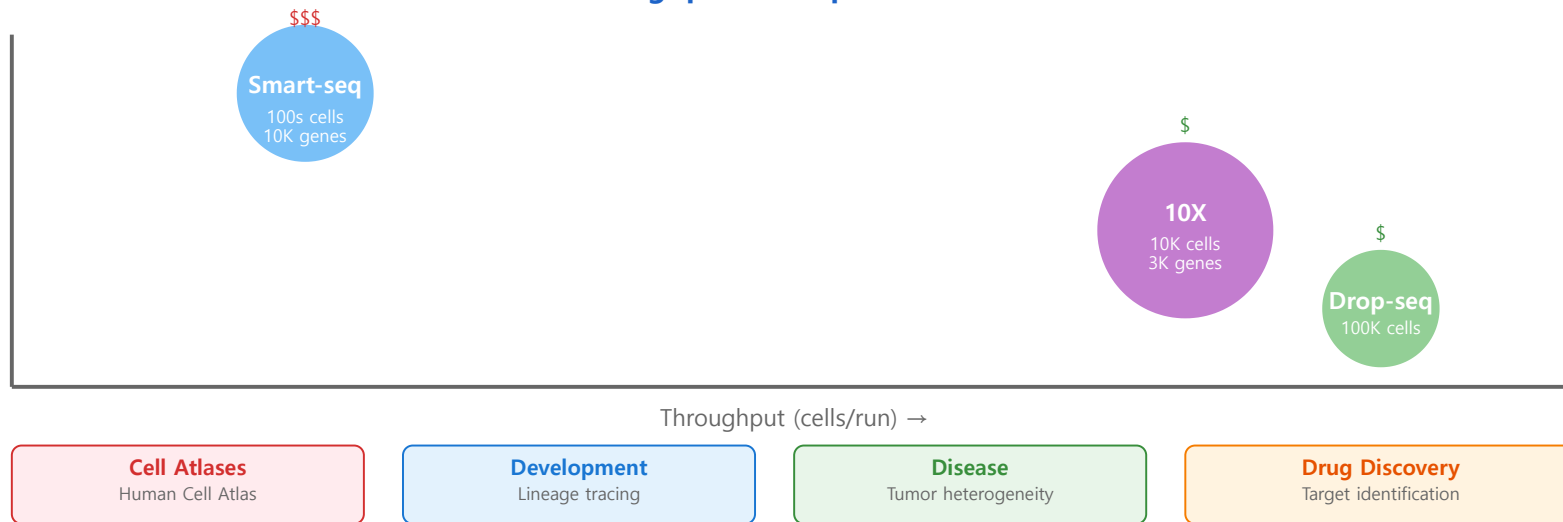
- Technology overview
- Cell isolation
- Quality control
- Analysis challenges

scRNA-seq Overview

Evolution & Comparison of Single-Cell Technologies



Throughput vs Depth Trade-off



💡 Revolution in understanding cellular heterogeneity

Droplet-based Methods

10X Genomics Platform

Most widely used - Chromium platform with GEMs

Drop-seq Principles

Co-encapsulation of cells and barcoded beads

InDrop Technology

Hydrogel beads with photocleavable barcodes

Barcode Design

Cell barcode + UMI for molecular counting

Doublet Detection

Computational and experimental QC for multiplets

💡 High throughput but lower sensitivity per cell

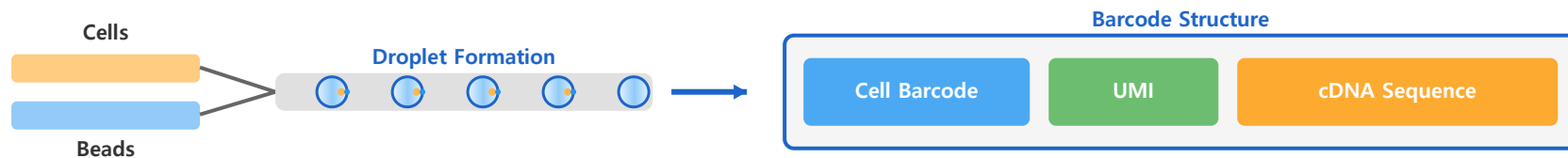


Plate-based Methods

Plate-based scRNA-seq Workflow

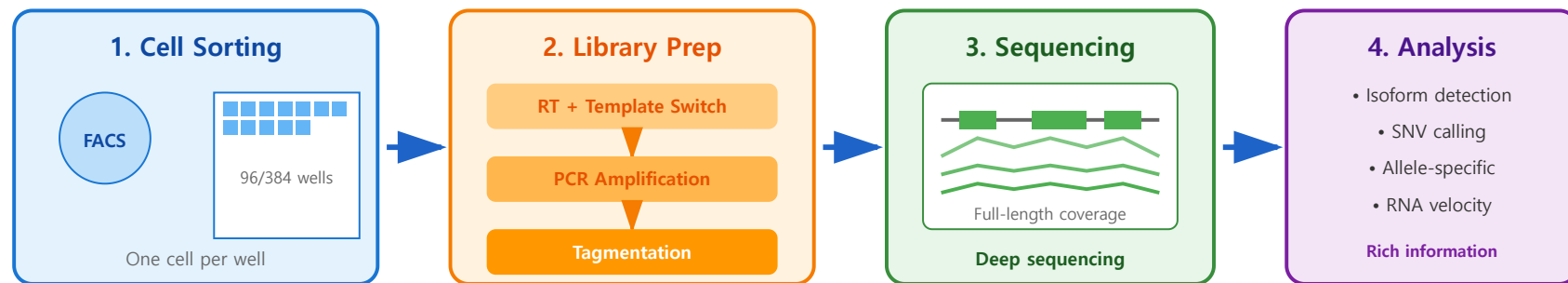


Plate-based Methods Comparison

Smart-seq2/3

- ✓ Full-length transcripts
- ✓ Highest sensitivity
- ✓ Isoform analysis
- ✗ No UMIs
- ✗ Higher cost/cell

MARS-seq

- ✓ UMI incorporation
- ✓ Automated
- ✓ 3' counting
- ✓ Cost-effective
- ✗ 3' bias

CEL-seq2

- ✓ Linear amplification
- ✓ UMIs
- ✓ Low bias
- ✓ Multiplexing
- ✗ Complex protocol

💡 Lower throughput but deeper sequencing per cell

Data Preprocessing

Cell Filtering

Remove low-quality cells and empty droplets

Gene Filtering

Exclude genes detected in too few cells

Normalization Methods

Account for sequencing depth and composition

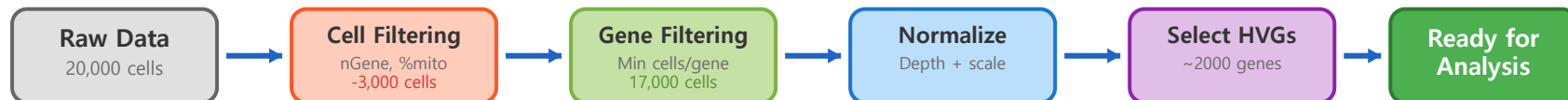
Imputation Strategies

Handle dropout events (use with caution)

Batch Effects

Technical variation from sample processing

💡 Quality control is critical for downstream analysis



QC Metrics: nGene: 200-6000 nUMI: 500-50000 %mito < 10% %ribo: varies Doublets: <5%

Dimensionality Reduction

PCA for scRNA-seq

First step to reduce noise and computational burden

t-SNE Principles

Preserves local structure, stochastic

UMAP Advantages

Faster, preserves global + local structure

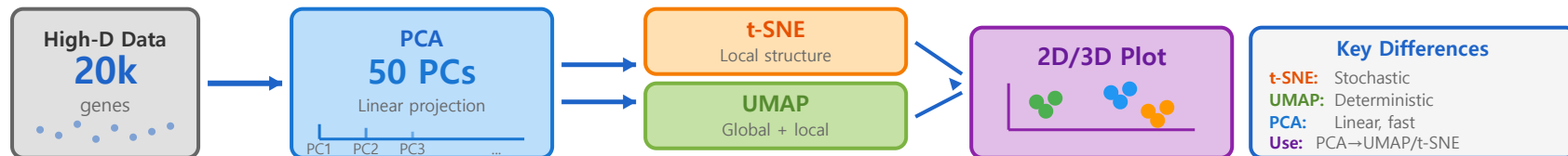
Diffusion Maps

Captures continuous trajectories

Parameter Selection

Perplexity, n_neighbors affect results

💡 Visualization != clustering - use both appropriately



Clustering Methods

Graph-based Clustering

Build kNN graph then find communities

Leiden Algorithm

Improved Louvain with better guarantees

K-means Adaptations

SC3 uses consensus clustering

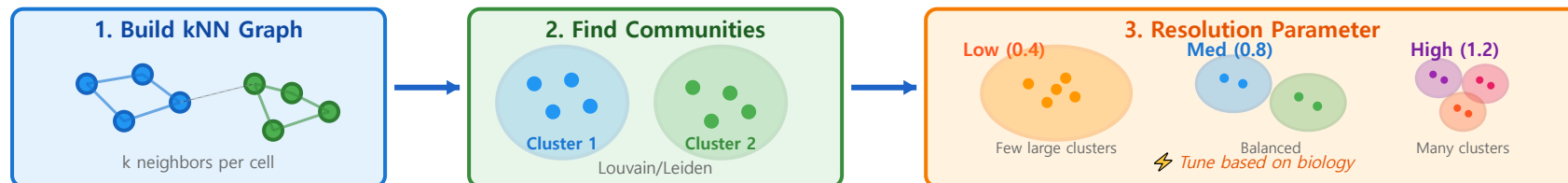
Resolution Selection

Higher resolution = more clusters

Stability Analysis

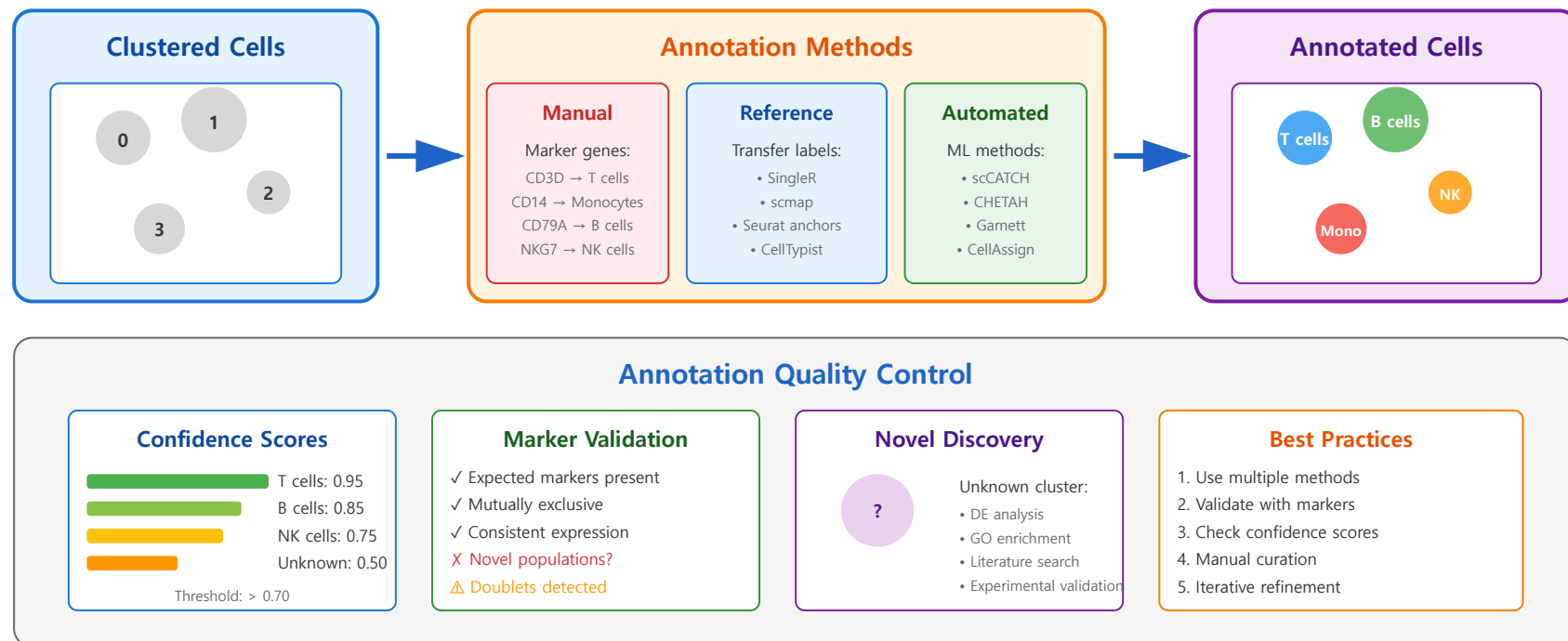
Bootstrap to assess cluster robustness

💡 No single correct clustering - depends on biological question



Cell Type Annotation

Cell Type Annotation Pipeline



💡 Combine automated tools with manual curation

Trajectory Analysis

Pseudotime Inference

Order cells along developmental paths

Branching Processes

Identify cell fate decisions

Monocle Algorithm

Reverse graph embedding

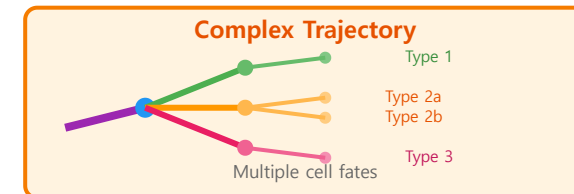
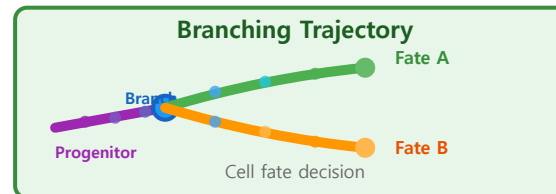
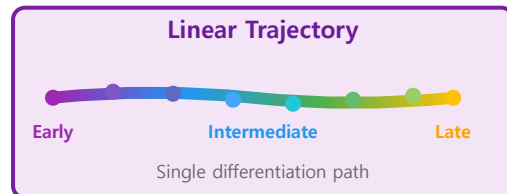
Slingshot Method

Cluster-based trajectory inference

Validation Approaches

Known genes, time-series data

💡 Assumes continuous progression - verify biological relevance



Part 3/3:

Advanced Methods

- Spatial context
- Multi-modal data
- Velocity analysis
- Communication inference

Spatial Transcriptomics

Visium Technology

10X spatial - 55 μ m spots, whole transcriptome

MERFISH Principles

Multiplexed error-robust FISH, subcellular resolution

seqFISH Evolution

Sequential FISH with barcoding, 10,000+ genes

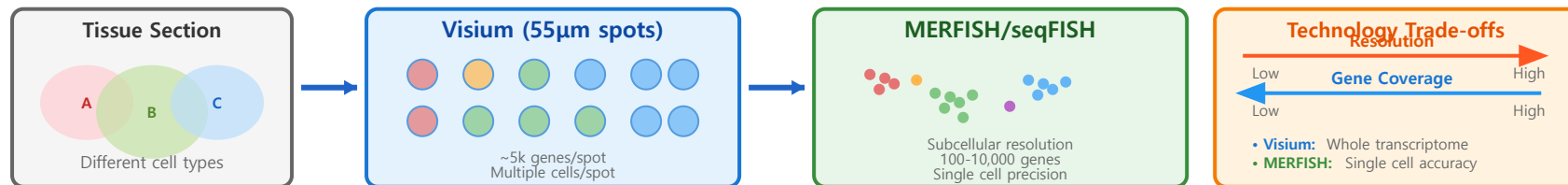
Slide-seq Methods

Bead-based spatial barcoding, 10 μ m resolution

Resolution Trade-offs

Gene coverage vs spatial resolution vs throughput

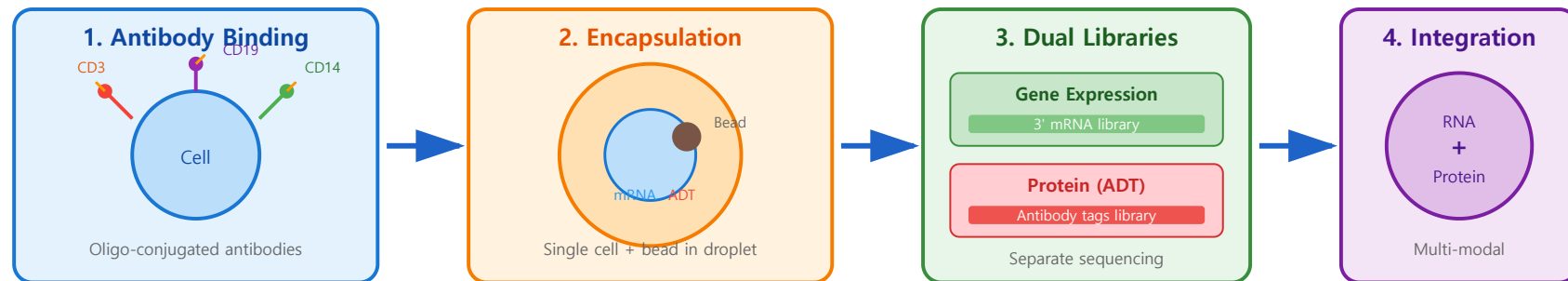
💡 Spatial context reveals tissue architecture and cell-cell interactions



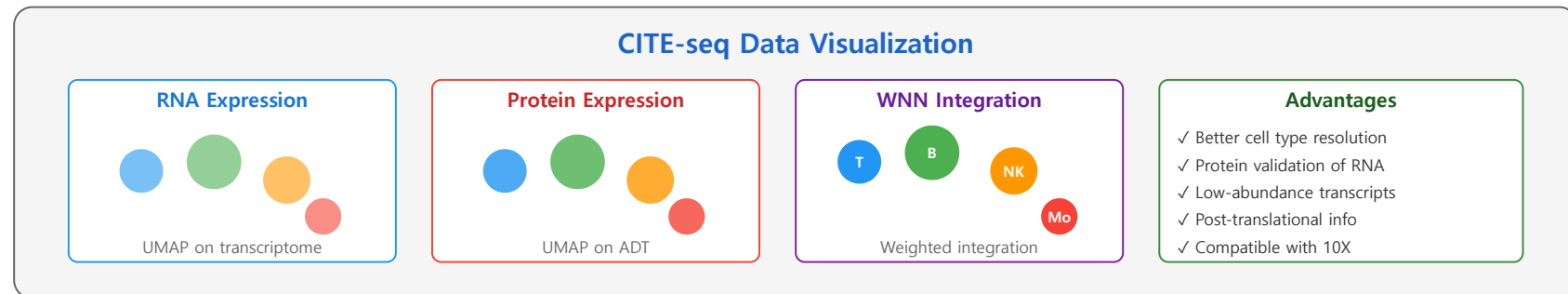
CITE-seq

Cellular Indexing of Transcriptomes and Epitopes by Sequencing

CITE-seq Mechanism



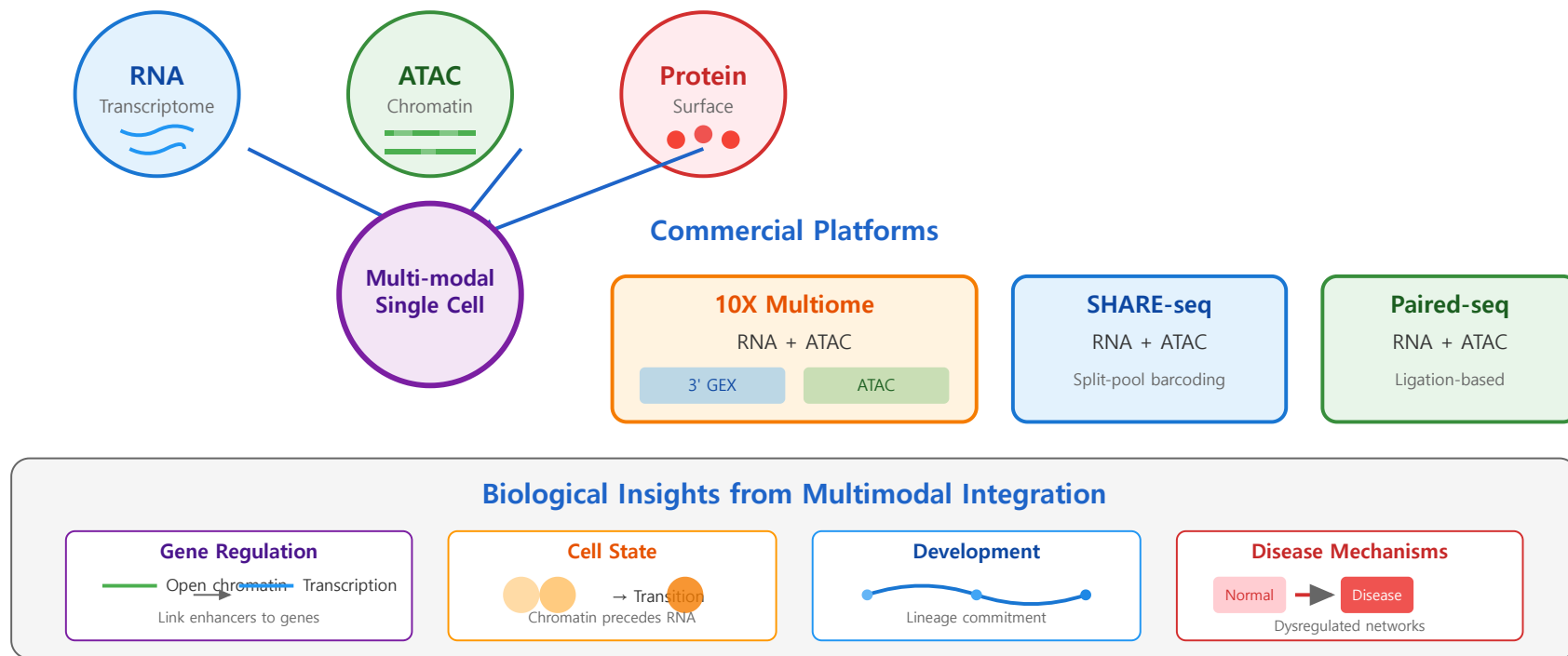
CITE-seq Data Visualization



💡 Bridges transcriptomics and proteomics at single-cell level

Multimodal Omics

Single-Cell Multimodal Technologies



💡 Multi-omics reveals regulatory mechanisms

RNA Velocity

Spliced/Unspliced Ratio

Infer transcriptional dynamics from steady-state

Velocity Estimation

Predict future cell states

Dynamic Models

Account for transcription, splicing, degradation

scVelo Improvements

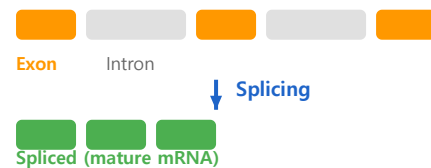
Dynamical model, latent time

Interpretation

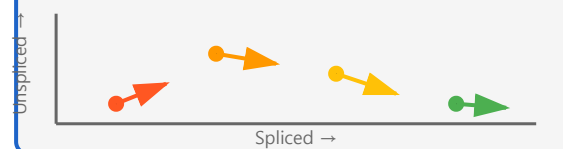
Direction and magnitude of cell state changes

💡 RNA velocity adds temporal dimension to snapshots

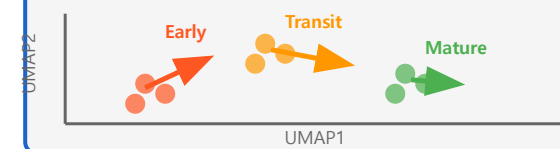
Gene Transcription & Splicing



Unspliced vs Spliced



Velocity Field on UMAP



Cell-Cell Communication

Inferring Cell-Cell Interactions from scRNA-seq

Ligand-Receptor Interaction



Communication Inference Tools

CellPhoneDB

- Statistical framework
- Permutation test
- Multi-subunit
- Spatial optional

NicheNet

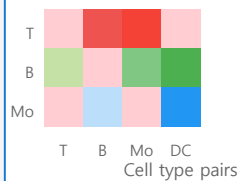
- Gene regulatory
- Prior knowledge
- Target prediction
- Prioritization

CellChat

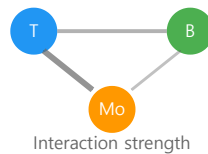
- Network analysis
- Pattern recognition
- Visualization
- Comparison

Typical Analysis Outputs

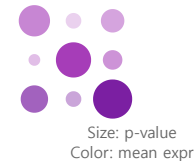
Interaction Heatmap



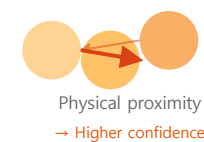
Network Graph



Dot Plot



Spatial Context



💡 Infer cellular communication from expression patterns

Batch Effect Correction

MNN Correction

Mutual nearest neighbors for batch alignment

Harmony Algorithm

Iterative clustering and correction

LIGER Integration

Integrative non-negative matrix factorization

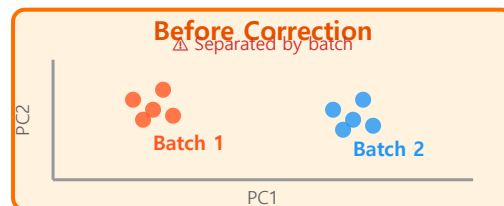
Seurat Integration

Canonical correlation analysis + anchors

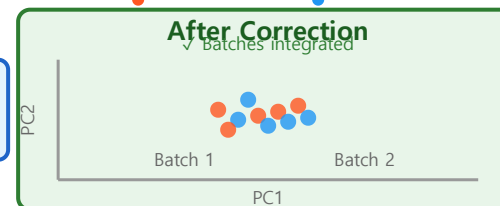
Benchmark Studies

Compare methods on simulated and real data

💡 Critical for multi-sample and multi-technology integration



Correction
MNN/Harmony
LIGER/Seurat



Integration Goals

- ✓ Mix batches
- ✓ Preserve biology
- ✓ Keep cell types
- ✓ Remove technical

Balance is key!

Integration Methods

Anchor-based Methods

Seurat, LIGER find correspondence between datasets

Deep Learning Approaches

scVI, scGAN learn shared latent space

Reference Building

Create comprehensive cell atlases

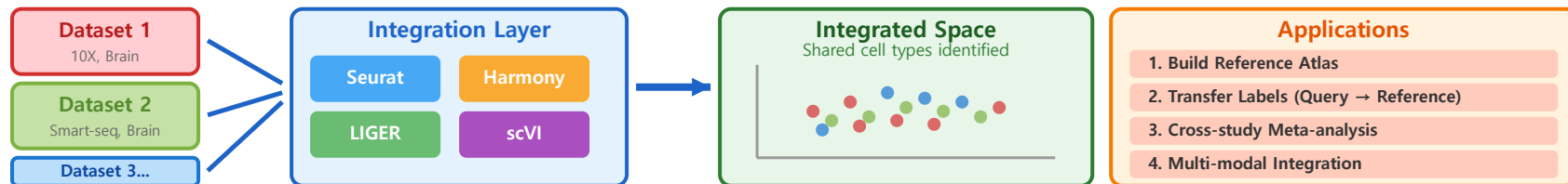
Query Mapping

Project new data onto reference

Performance Metrics

Biological conservation vs batch mixing

💡 Integration enables meta-analysis and transfer learning



Hands-on: Seurat Tutorial

Seurat v5 Standard Workflow

1. Data Loading & QC

```
# Read 10X data
data <- Read10X("filtered_feature_bc_matrix/")
seurat <- CreateSeuratObject(data, min.cells=3)
```

2. QC Filtering

```
# Calculate mitochondrial %
seurat[["percent.mt"]] <- PercentageFeatureSet(seurat, "^MT-")
seurat <- subset(seurat, nFeature_RNA > 200 & percent.mt < 10)
```

3. Normalization & Scaling

```
seurat <- NormalizeData(seurat)
seurat <- FindVariableFeatures(seurat, nfeatures=2000)
seurat <- ScaleData(seurat)
```

4. Dimension Reduction & Clustering

```
seurat <- RunPCA(seurat) %>% RunUMAP(dims=1:30)
seurat <- FindNeighbors(seurat) %>% FindClusters(res=0.5)
DimPlot(seurat, label=TRUE) + NoLegend()
```

Key Visualizations

QC Metrics



UMAP Clustering



Feature Plot



Integration with Harmony/Seurat

```
# Integration of multiple samples
seurat <- IntegrateLayers(seurat, method=HarmonyIntegration)
```



Find Markers & Annotate

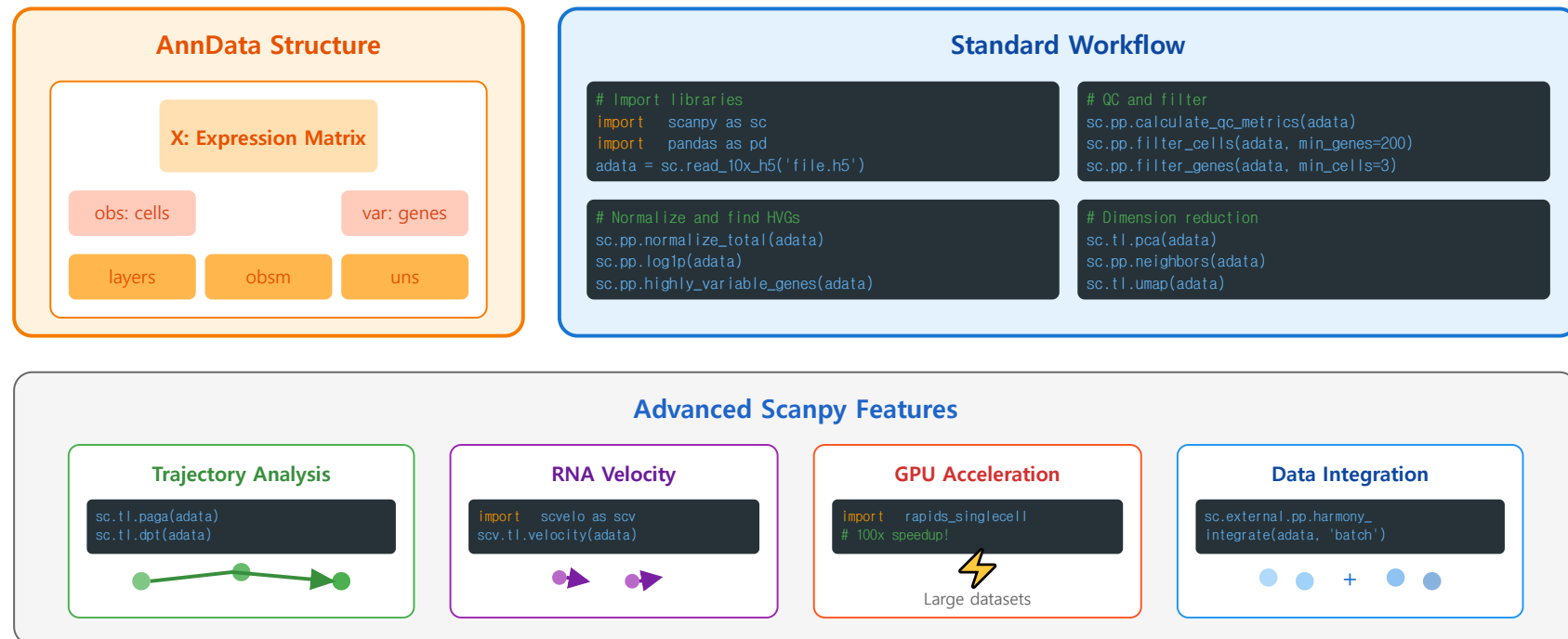
```
markers <- FindAllMarkers(seurat, only.pos=TRUE)
new.ids <- c("T cells", "B cells", "NK", "Monocytes")
```



💡 Most widely used R package for scRNA-seq analysis

Hands-on: Scanpy Analysis

Scanpy: Python-based Single Cell Analysis



💡 Python ecosystem with extensive documentation

Thank you!

Key Applications

- Disease studies - Cell type changes in pathology
- Development biology - Cell fate trajectories
- Drug discovery - Target identification and validation
- Clinical futures - Diagnostic and therapeutic applications

Introduction to Biomedical Data Science