

Spatial Transcriptomics

Project Proposal

November 8, 2022

Project Sponsor

ENPH 479 - team 2255



Dr. Jia Rui Ding
PhD



Sayem Zaman



Daniel Chen

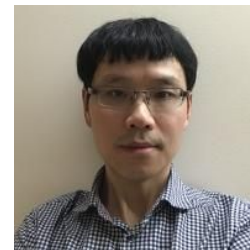


David Wang

Project Sponsor

Dr. Jia Rui Ding

- Assistant Professor in UBC Computer Science Department
- Postdoc Associate, Broad Institute of MIT & Harvard (2017 - 2021)



Research Interest

- Bioinformatics
- Computational Biology
- Machine Learning
- Probabilistic Deep Learning
- Single-Cell Genomics

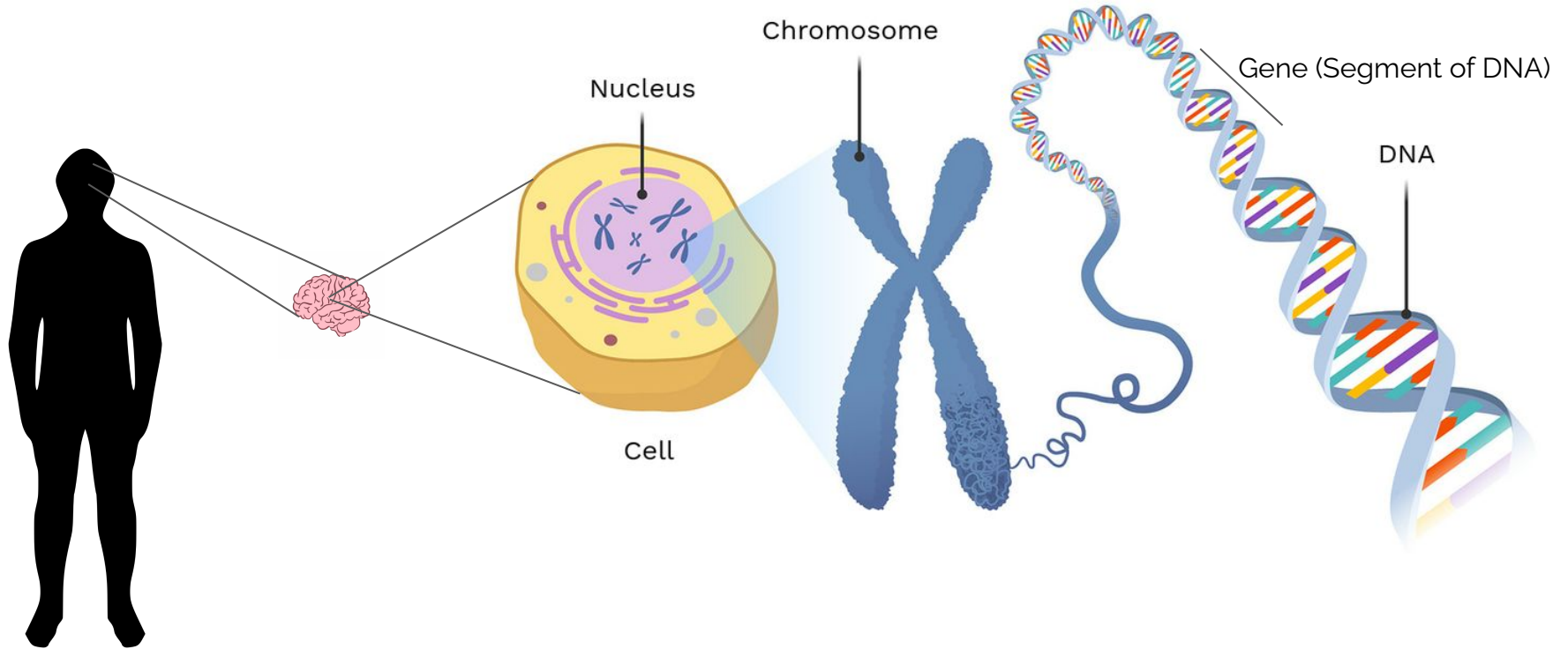
Project Proposal Outline

- Introduction
 - Impacts
 - RNA Sequencing & Spatial Transcriptomics
- Project Detail
 - Project Goal
 - Project Workflow
- Technical Implementation Details
 - VAE & VGAE
 - Map Function between Latent Spaces
 - Testing and Validation
- System Level Diagram
- Project Timeline
- Deliverables

Roughly **9.6 million** people die from cancer every year

1 in every 6 deaths is due to cancer

Cells and DNA



Immunotherapy for treating cancer

Our immune system has the ability to find and destroy cancer cells.

But cancer cells can sometimes **hide from the immune system** and avoid being destroyed. Cancer cells may also stop the immune system from working properly.

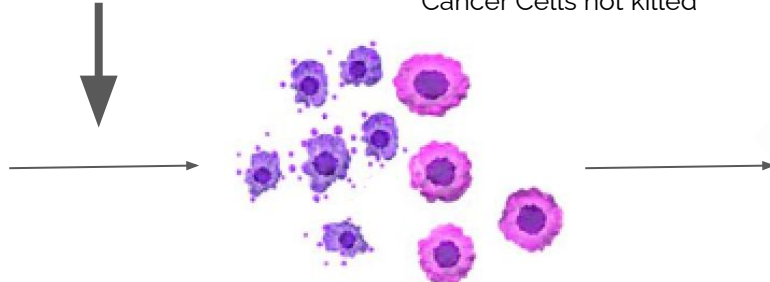
Immunotherapy helps to **strengthen or restore the immune system's natural ability to fight cancer** (with very little damage to patient's body)

Cancer resistance to Immunotherapy

Cancer Cells in a tumor



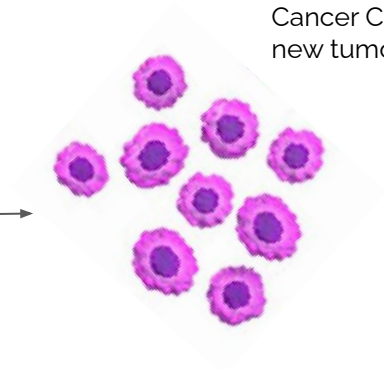
Apply Immunotherapy



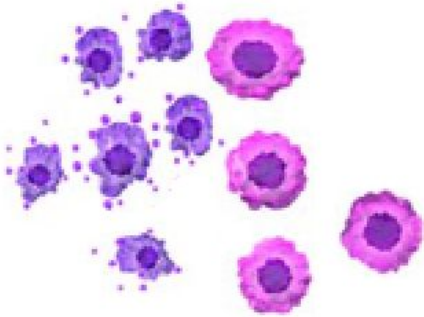
Cancer Cells not killed

Cancer Cells killed by immunotherapy

Cancer Cells divide to form new tumor



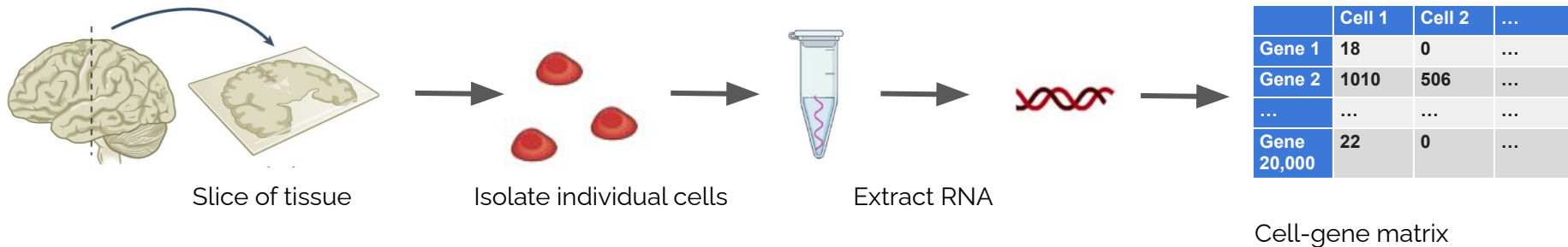
A promising approach...



If we know the **gene expression information** (which genes) and **spatial information** (where they are) of cancer cells not affected by immunotherapy, we can learn to make better treatments

How to get **gene expression** information from cells

Single Cell RNA Sequencing (scRNA-seq)

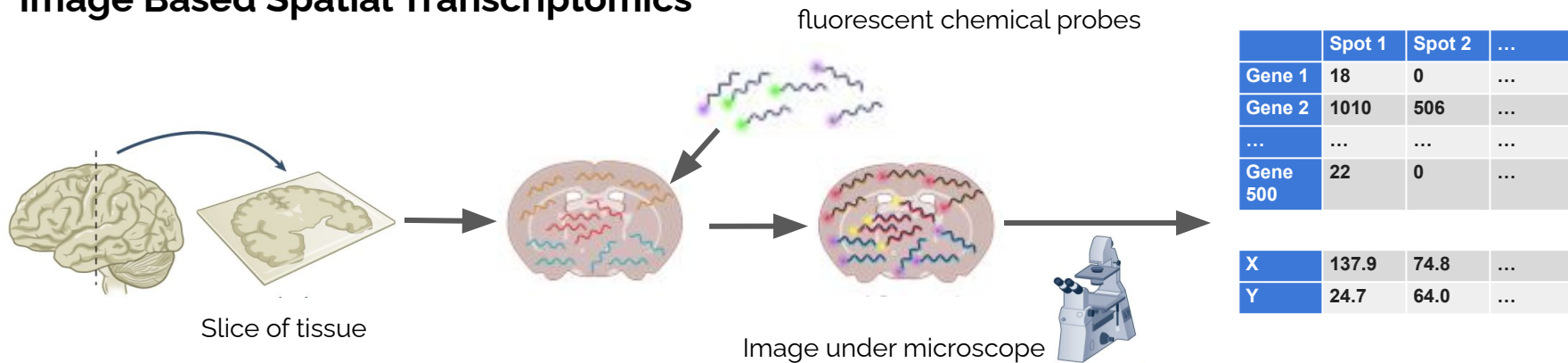


Extracts all genetic data in the cell (i.e. 20,000 genes)

Problem: we don't know what part of the original tissue the genes came from

How to get **spatial** information from cells

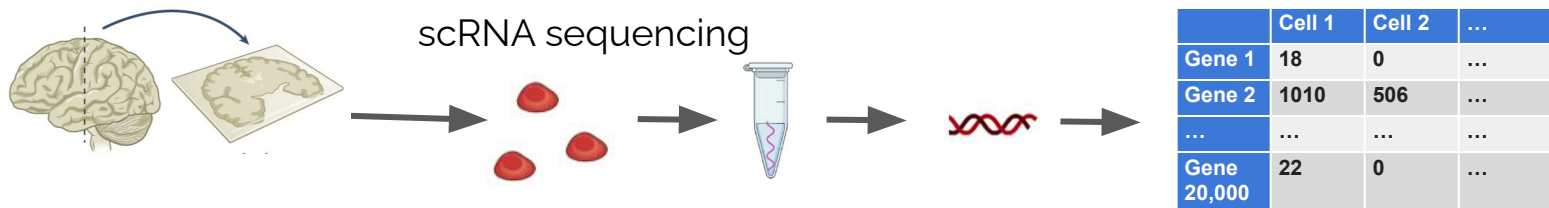
Image Based Spatial Transcriptomics



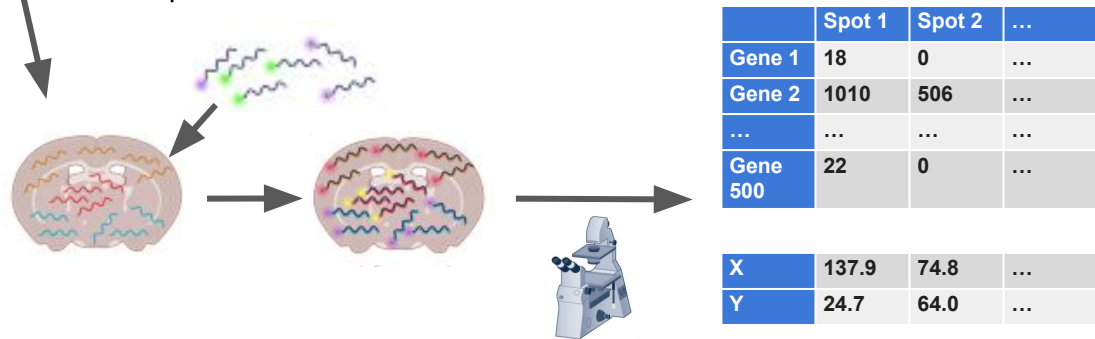
Identifies spatial location of genes, but due to chemical limits only a smaller number of genes can be profiled

Problem: we can only recover a small number of genes and their locations (~500)

Goal: Get both **spatial** and **genetic** information from cells



Spatial
Transcriptomics



Can the relationship be learned?

Project Goal

Our project aims to develop a **probabilistic machine learning model** to map the image-based spatial transcriptomics data to scRNA sequencing data.

- This will allow us to **infer the spatial information of cells based on the genetic information** of the RNA molecules within each cell.

Data Collection

- **Image-Based Spatial Transcriptomics** and **Single-Cell RNA Sequencing**

- High dimensional data (tens of thousands of features)

Generate ML Models

Relationship Between Models

How can we analyze the data using machine learning?

scRNA seq data

	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...

$$\begin{matrix}
 1 & 2 & \dots & n \\
 \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ m \end{matrix} & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}
 \end{matrix}$$

Sparse Matrix (Lots of zeros in data matrix)

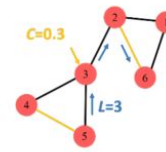
Spatial transcriptomics data

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

$$\begin{matrix}
 1 & 2 & \dots & n \\
 \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ m \end{matrix} & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}
 \end{matrix}$$

X	137.9	74.8	...
Y	24.7	64.0	...

0	1	0	0	0	1
1	0	1	0	0	1
0	1	0	1	1	0
0	0	1	0	1	0
0	0	1	1	0	0
1	1	0	0	0	0



Adjacency Matrix / Graph (Shows the spatial relationship to other cells)

High dimensional dataset

Another high dimensional dataset

Data is very complex and very different, hard to find direct mapping

Can we first transform the data into a more similar space, and still preserve the important features?

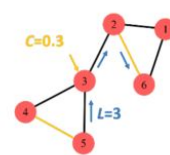
scRNA seq data

	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...

$$\begin{matrix}
 & \begin{matrix} 1 & 2 & \dots & n \end{matrix} \\
 \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ m \end{matrix} & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}
 \end{matrix}$$

Spatial transcriptomics data

0	1	0	0	0	1
1	0	1	0	0	1
0	1	0	1	1	0
0	0	1	0	1	0
0	0	1	1	0	0
1	1	0	0	0	0

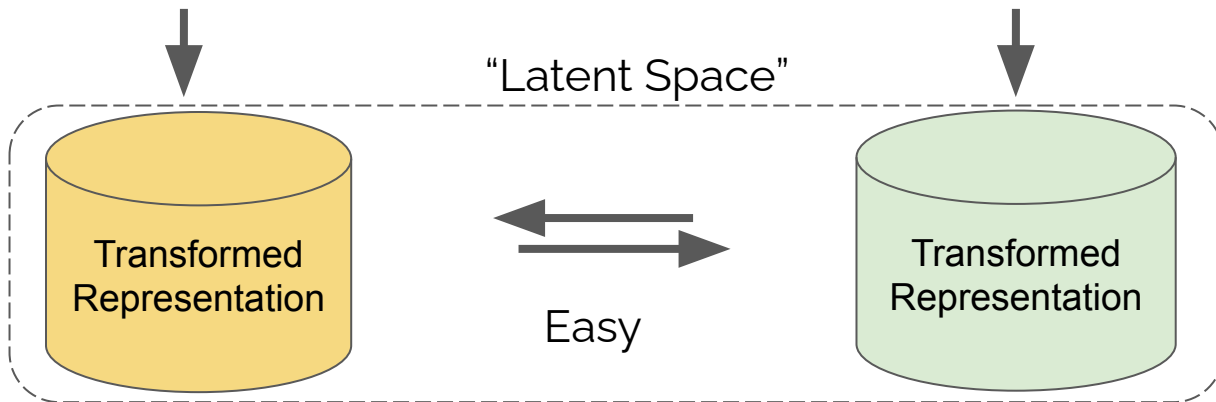


	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...
X	137.9	74.8	...
Y	24.7	64.0	...

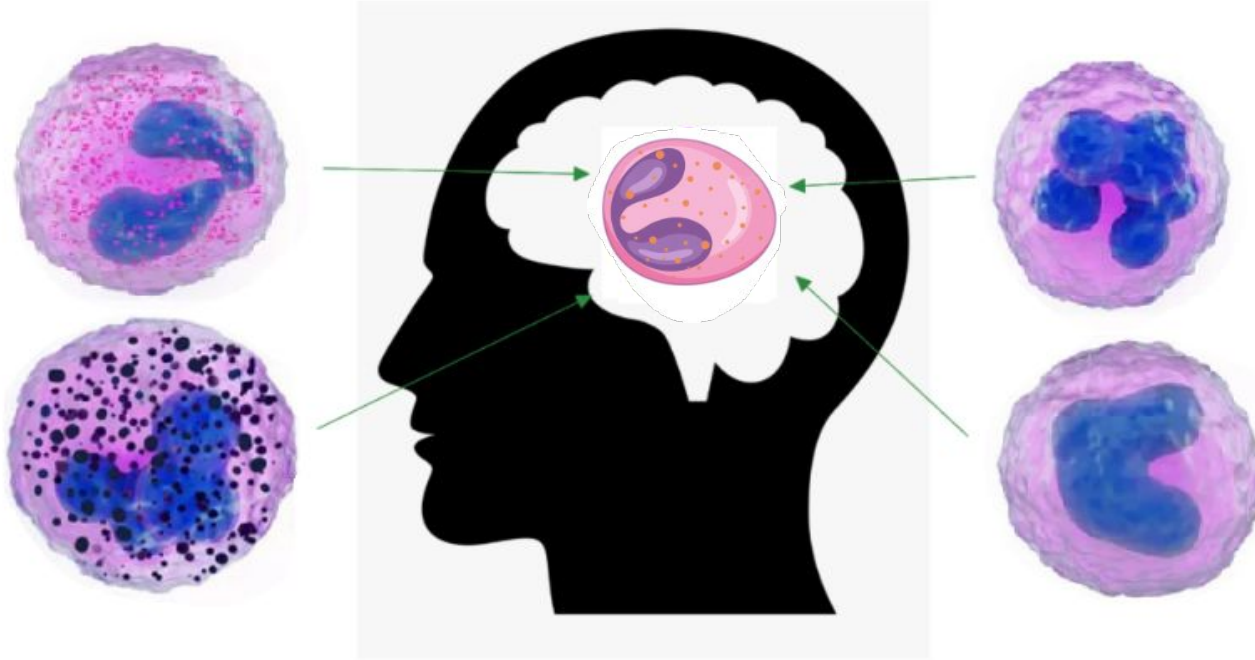
Very Hard



"Latent Space"



Latent space representation



Shared features:

1. Light or pinkish color
2. Relatively round shape
3. Well defined nucleus

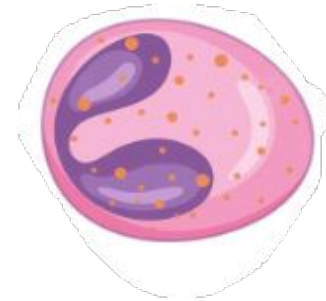
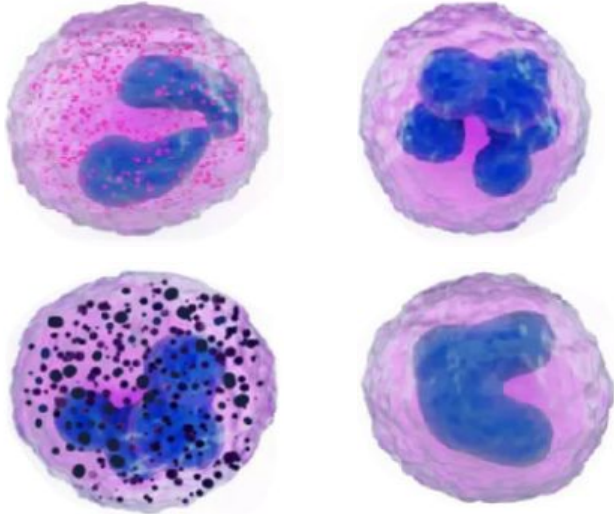
Understanding Latent spaces

A Latent Space is an abstract space that **encodes** a meaningful **compressed** internal representation of externally observed **high dimensional** events

High Multidimensional Space



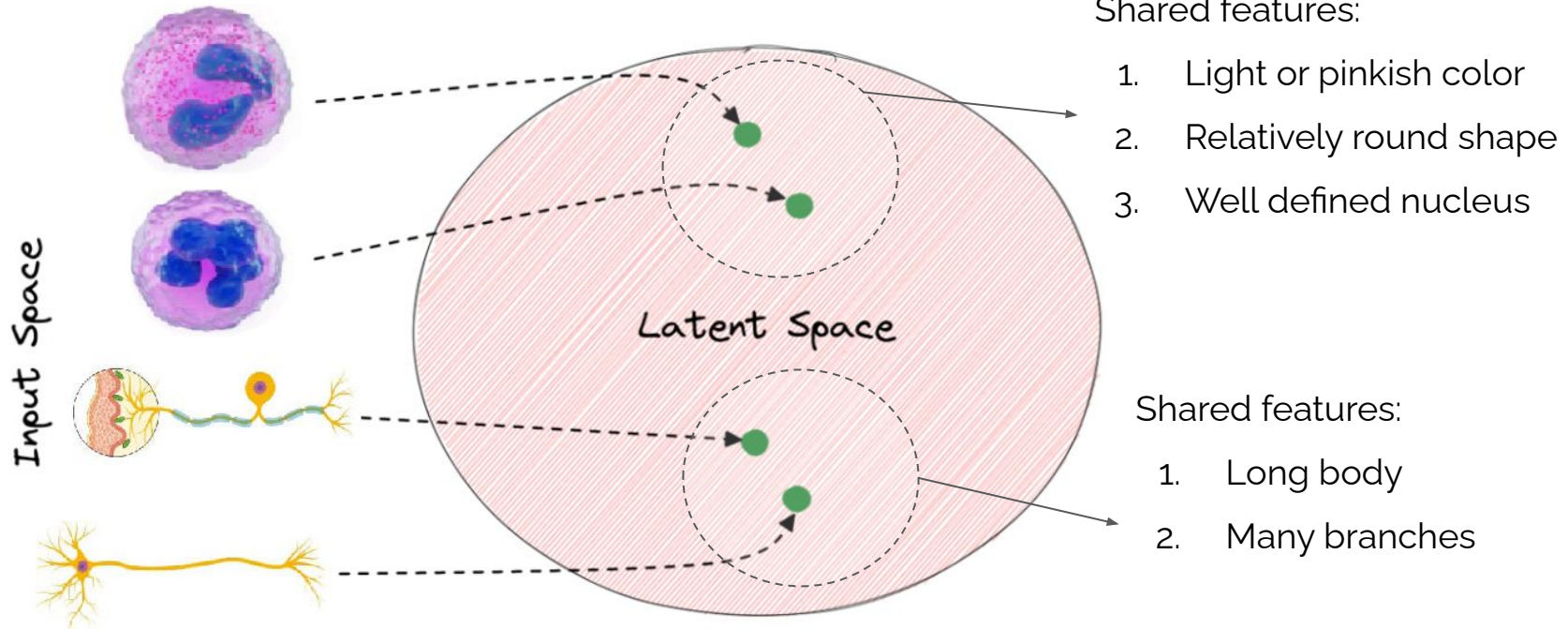
Lower Dimensional Representation



Lots of specific features which are useless for representation

Key features which encapsulate fundamental observations

Understanding Latent spaces



How do we obtain a good latent space for the dataset?

Variational Autoencoders

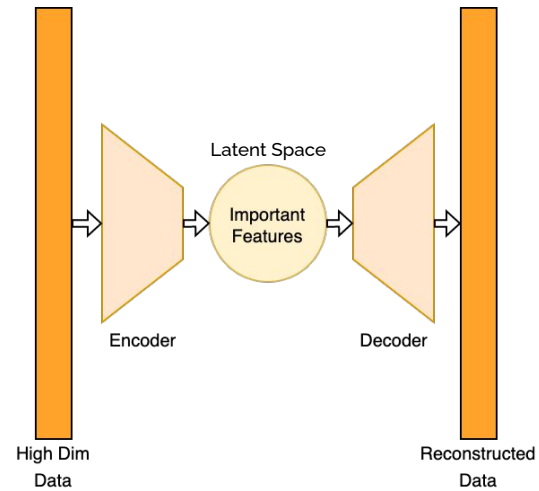
Probabilistic machine learning method that **encodes** and **decodes** a latent space

Encoder

- learns the “**important features**” that are worth preserving

Decoder

- learns how to “**interpret**” coded information to reconstruct the original information



Variational Autoencoders (VAE)

2 parts: encoder and decoder (neural networks)

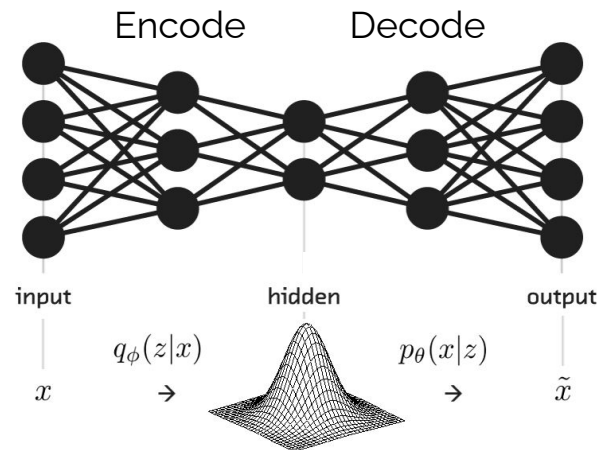
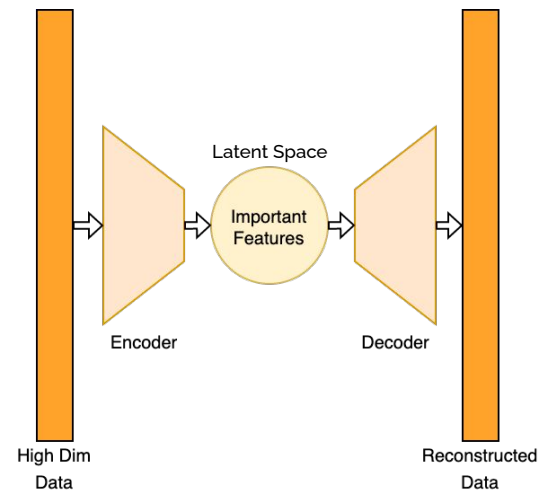
Encoder Network

- Transforms high-dim data into low-dim data

Decoder Network

- Reconstruct high-dim data from low dim data

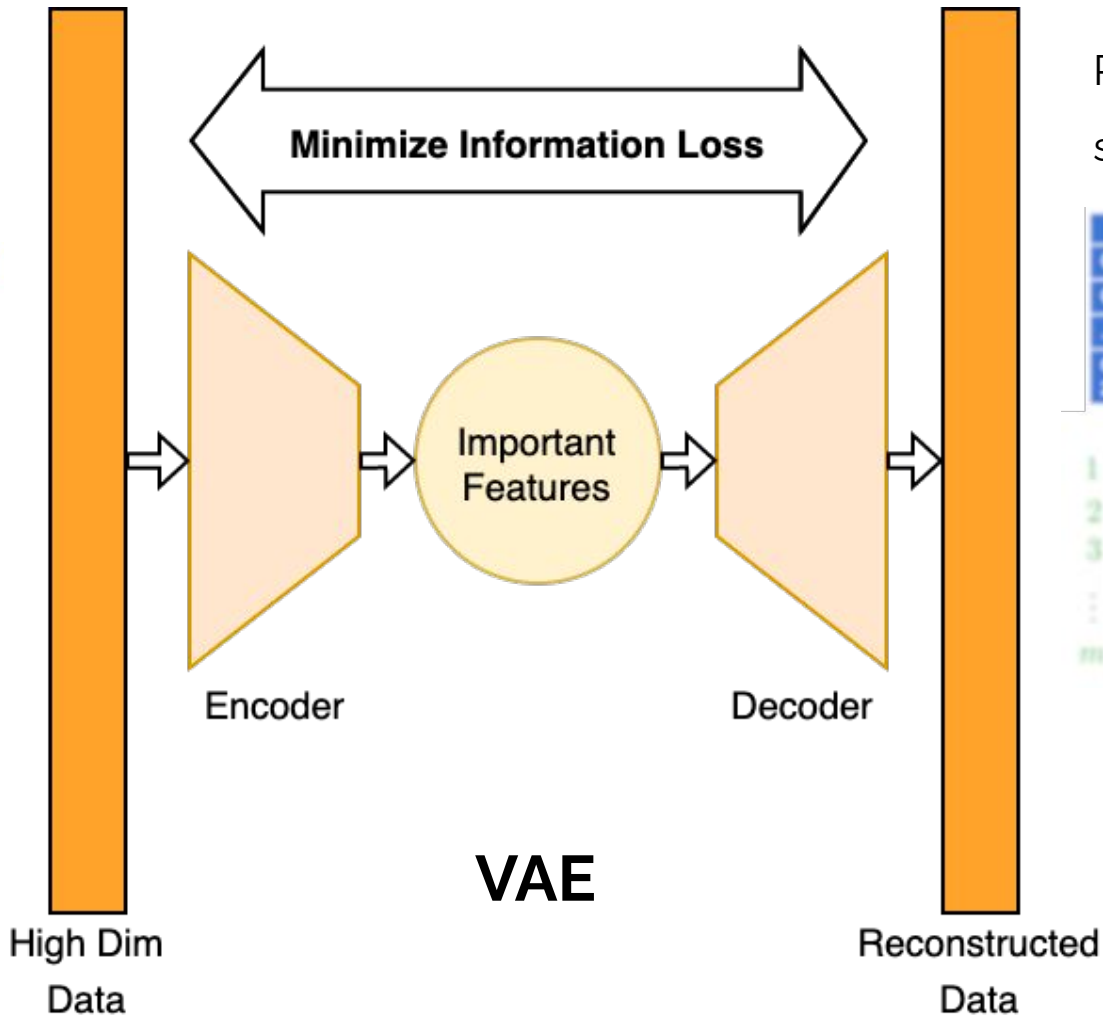
Minimize information loss of reconstructed data



scRNA seq data

	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...

$$\begin{matrix} 1 & 2 & \dots & n \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ m \end{matrix} & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \end{matrix}$$



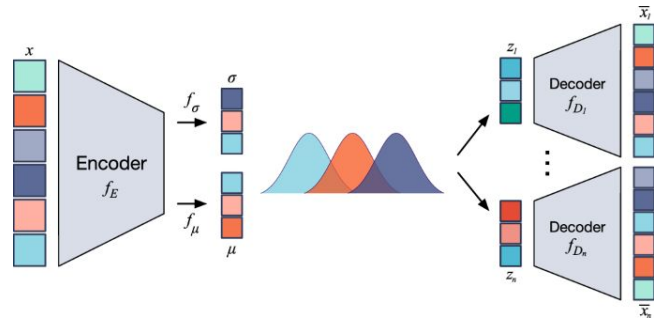
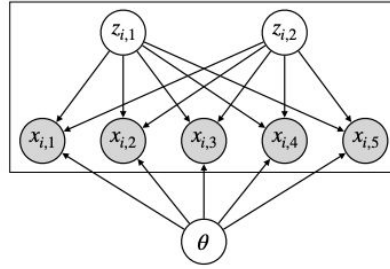
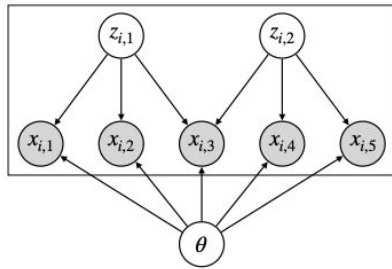
Reconstructed
scRNA seq data

	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...

$$\begin{matrix} 1 & 2 & \dots & n \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ m \end{matrix} & \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \end{matrix}$$

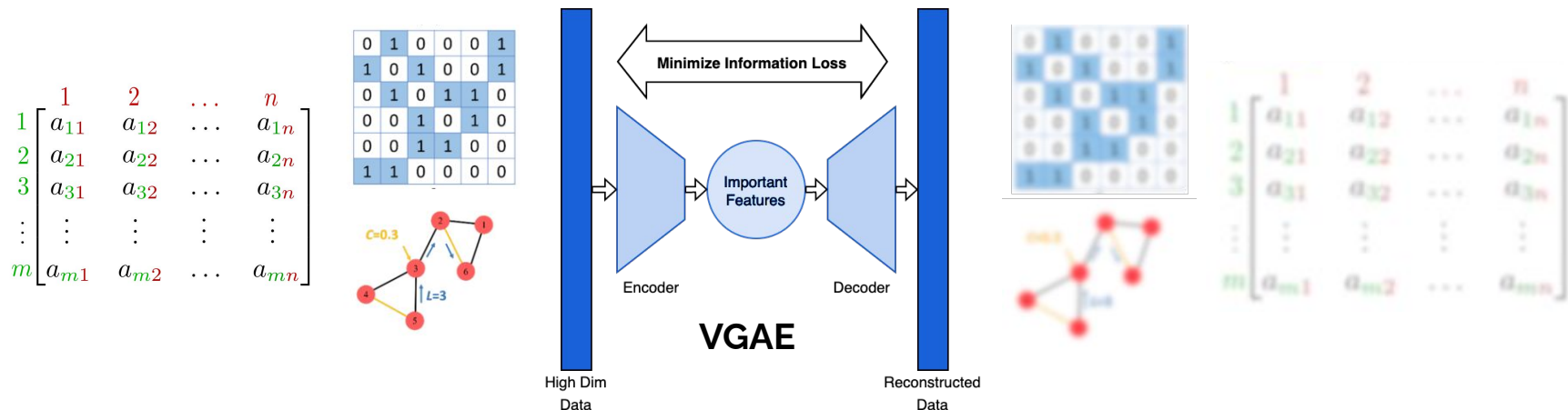
Benchmarking VAE Methods

- Many VAE algorithms in the field
- We will analyze them with the same dataset and pick the one with **best performance (lowest loss)**
- **Iterative Process:** continuously tune model structure and parameters



Variational Graph Autoencoder (VGAE)

- Similar to Variational Autoencoder (VAE)
- Additionally allows **spatial information (x, y)** as input
- Encodes **sequencing data** and **adjacency matrix (spatial graph)** together to a latent space



Data Collection

Image-Based Spatial Transcriptomics and Single-Cell RNA Sequencing

- High dimensional data

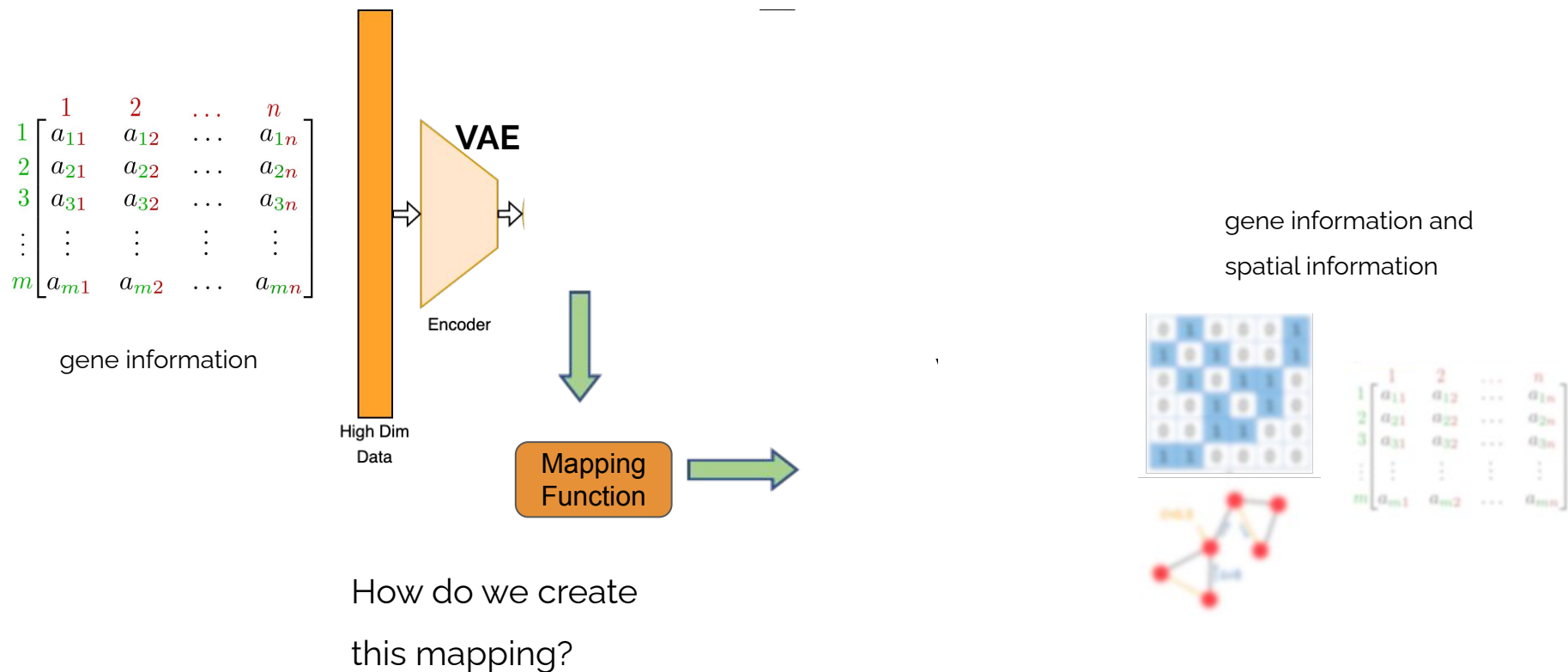
Generate ML Models

Variational Autoencoder and Variational Graph Autoencoder

- Analyze high-dim data
- Learn low-dim features

Relationship Between Models

Connecting between the Latent Spaces



Theories for Bridging Two Latent Spaces

- Many theories developed for non-bioinformatics fields
- Potential Methods:
 - Domain Transfer
 - Latent Mapping
 - Latent Translation
- **Evaluate and apply those implementations on our latent spaces**

Data Collection

Image-Based Spatial Transcriptomics and Single-Cell RNA Sequencing

- High dimensional data

Generate ML Models

Variational Autoencoder and Variational Graph Autoencoder

- Analyze high-dim data
- Learn low-dim features

Relationship Between Models

Mapping Function between two Latent Spaces

- Understand what the features represent
- Minimize reconstruction loss

Validation

Important to verify our model on **new, unseen data**.

- **VAE Model**

- Accuracy of encoding and decoding of scRNA-seq data (reduce information loss)

NEW
scRNA-seq
data

	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...



Encoder

VAE



Decoder



Minimize information loss

Validation

Important to verify our model on **new, unseen data**.

- VAE Model
 - Accuracy of encoding and decoding of scRNA-seq data (reduce information loss)
- **VGAE Model**
 - Accuracy of encoding and decoding of spatial transcriptomics data (reduce information loss)

Minimize information loss



NEW
spatial
transcriptomics
data

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...



Encoder

VGAE



Decoder



Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...

Validation

Important to verify our model on **new, unseen data**.

- VAE Model
 - Accuracy of encoding and decoding of scRNA-seq data (reduce information loss)
- GVAE Model
 - Accuracy of encoding and decoding of spatial transcriptomics data (reduce information loss)
- **Joint Model**
 - Accuracy of latent space mapping between scRNA-seq and spatial transcriptomics data (learning associations)

NEW
spatial
transcriptomics
data

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...



Encoder

VAE

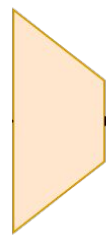


Decoder

VGAE

NEW spatial transcriptomics data

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...



Encoder

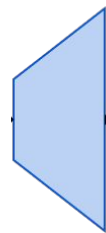
VAE



Mapping Function



VGAE



Decoder



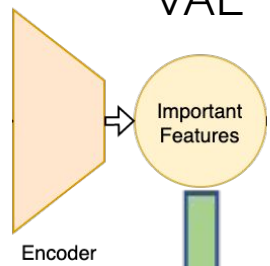
X	137.9	74.8	...
Y	24.7	64.0	...

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

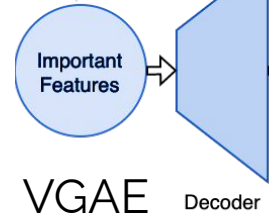
X	137.9	74.8	...
Y	24.7	64.0	...

NEW spatial transcriptomics data

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...



Mapping Function



	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...



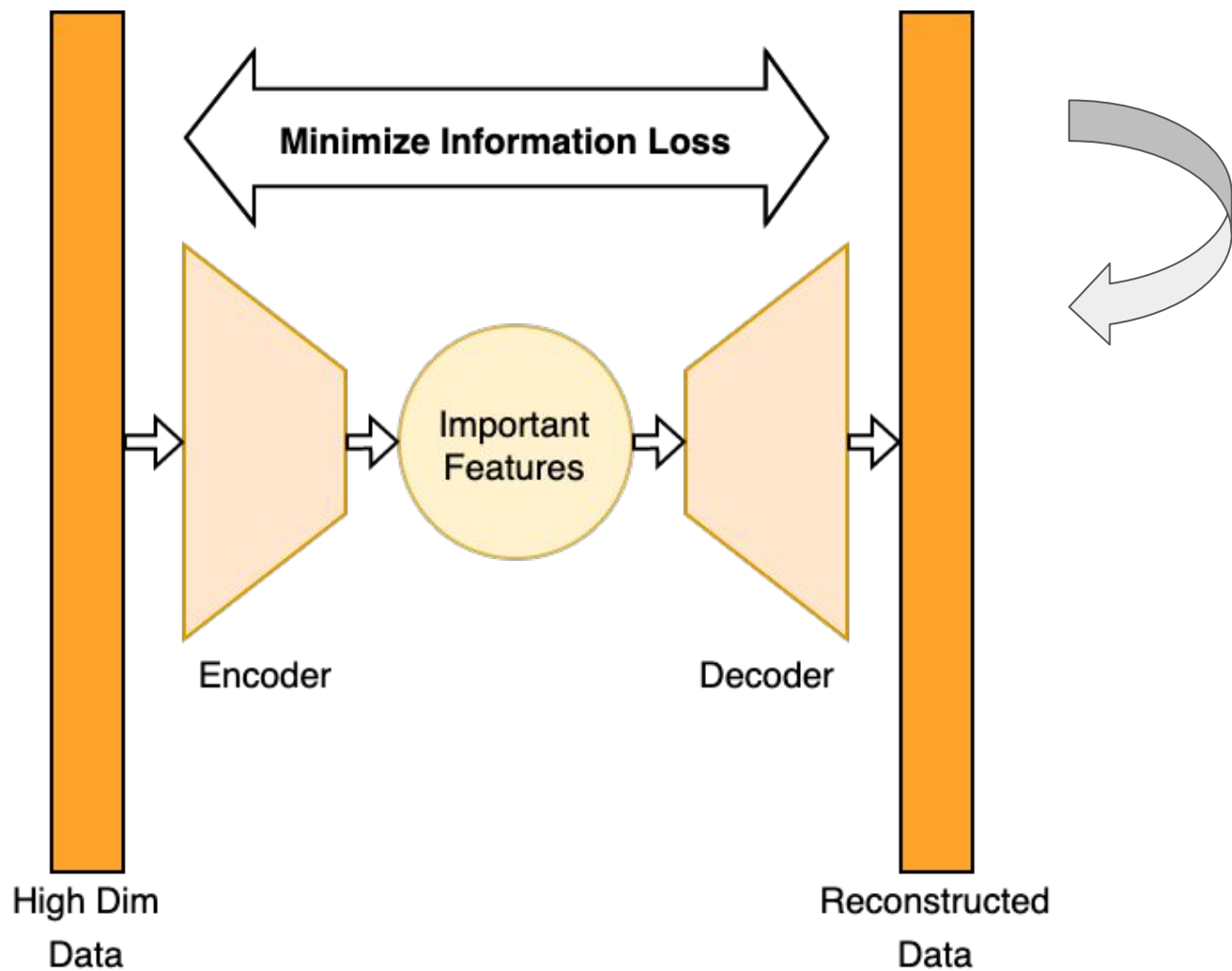
Minimize information loss

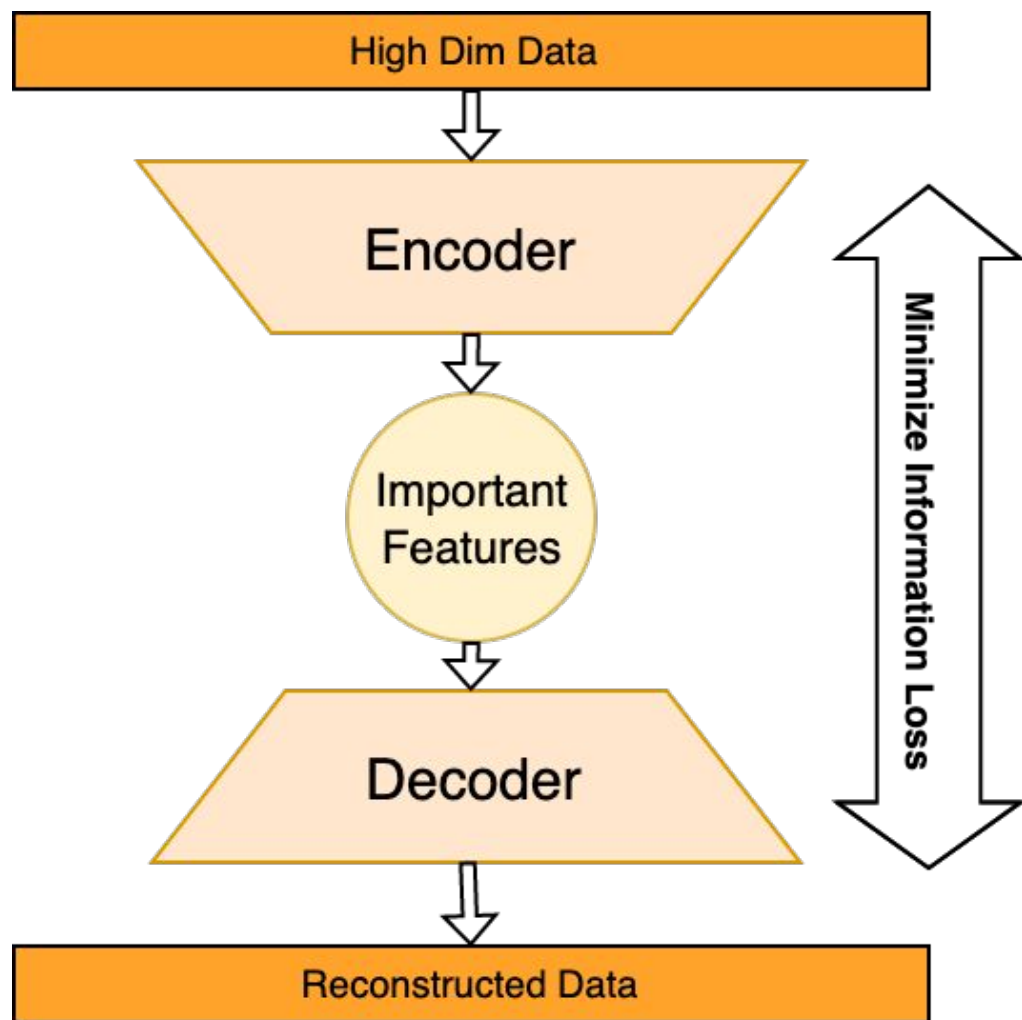
X	137.9	74.8	...
Y	24.7	64.0	...

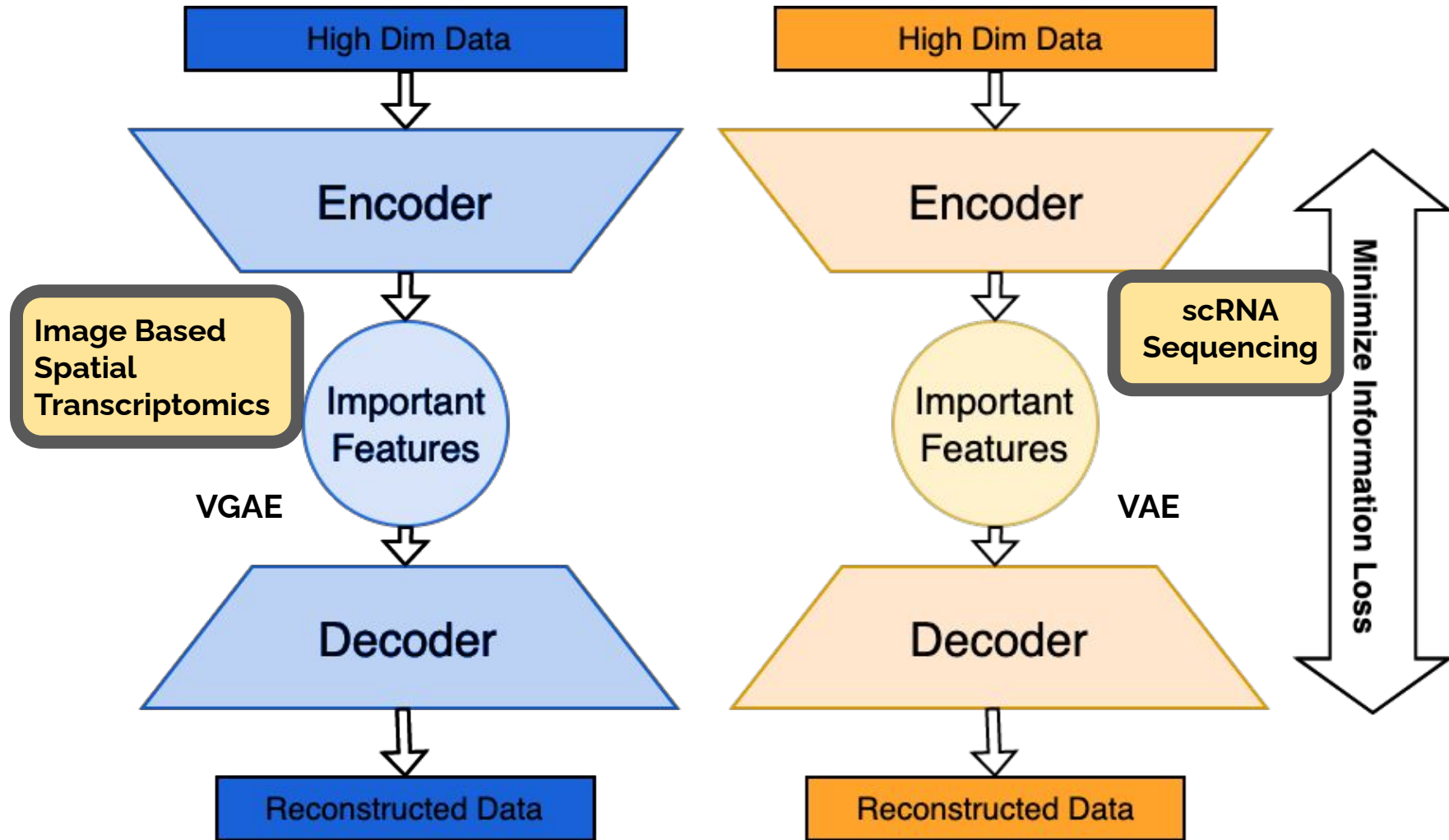
Validation

Important to verify our model on **new, unseen data**.

- VAE Model
 - Accuracy of encoding and decoding of scRNA-seq data (reduce information loss)
- GVAE Model
 - Accuracy of encoding and decoding of spatial transcriptomics data (reduce information loss)
- Joint Model
 - Accuracy of latent space mapping between scRNA-seq and spatial transcriptomics data (learning associations)







	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...

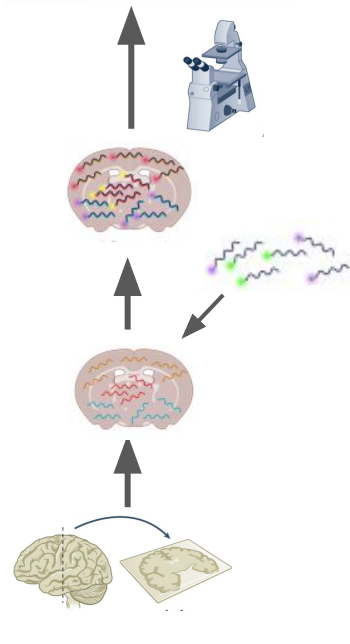
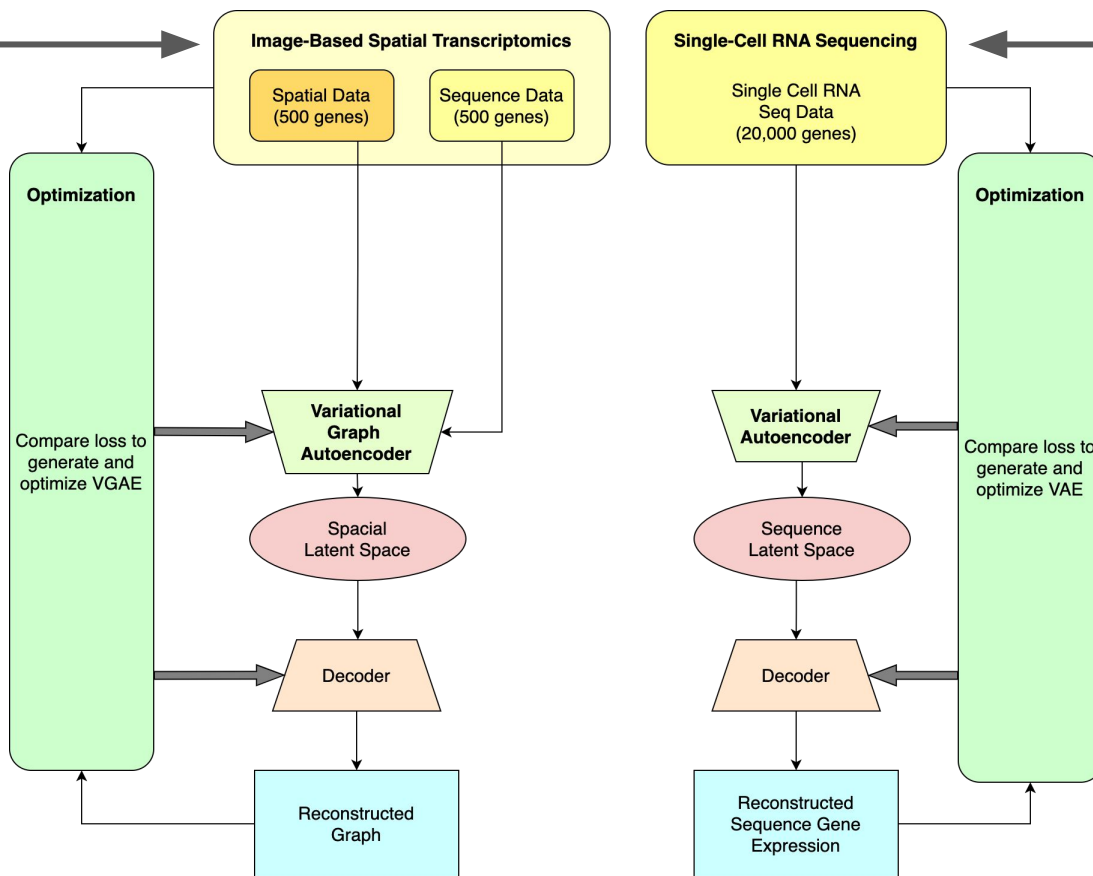
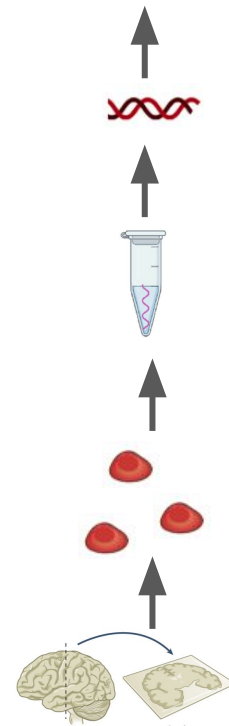


Image Based
Spatial Transcriptomics



	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...



Single Cell
RNA Sequencing

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

	X	Y	...
X	137.9	74.8	...
Y	24.7	64.0	...

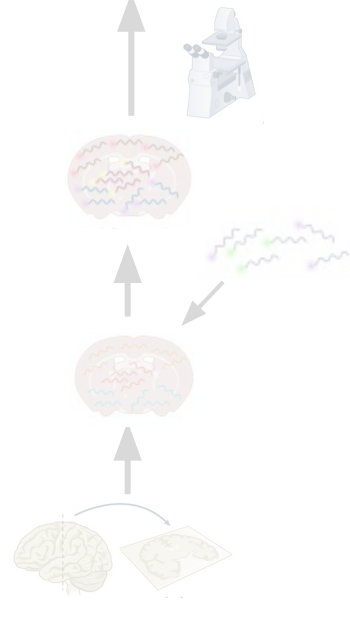
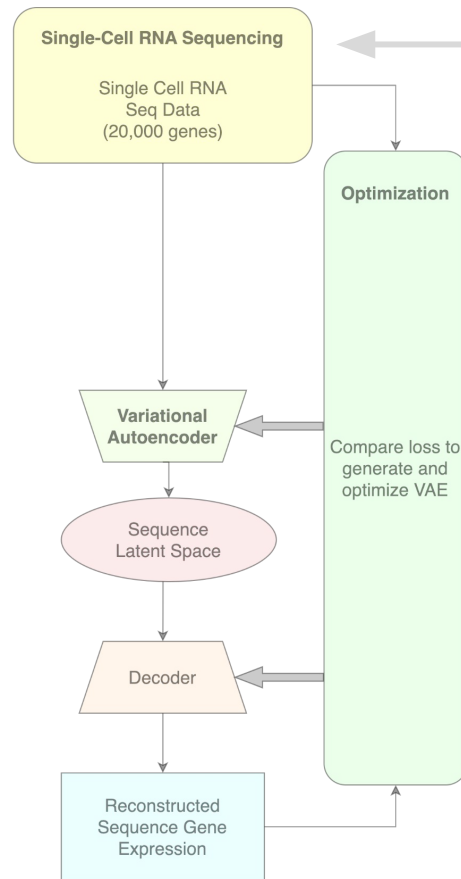
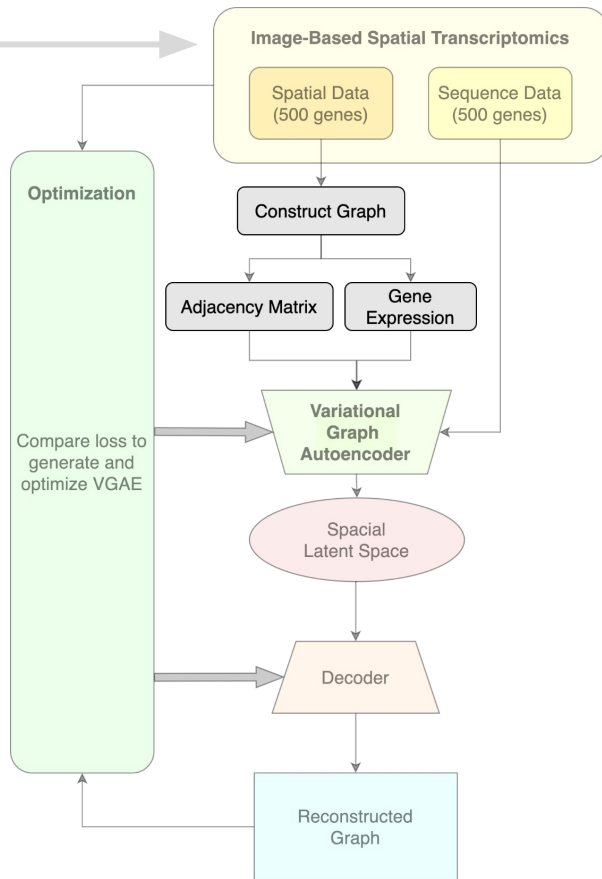
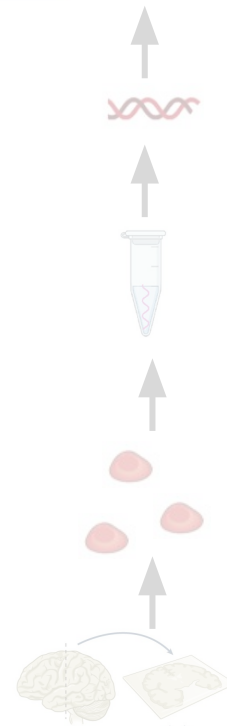


Image Based
Spatial Transcriptomics



	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...



Single Cell
RNA Sequencing

	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...

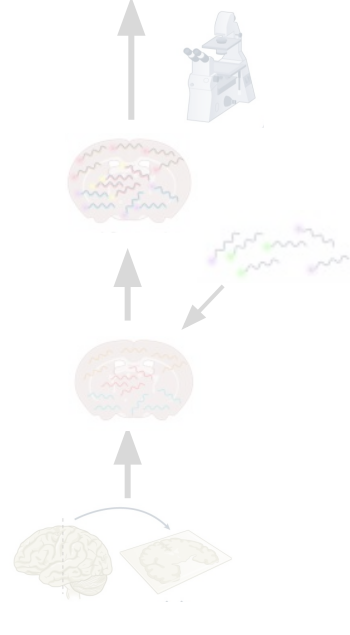
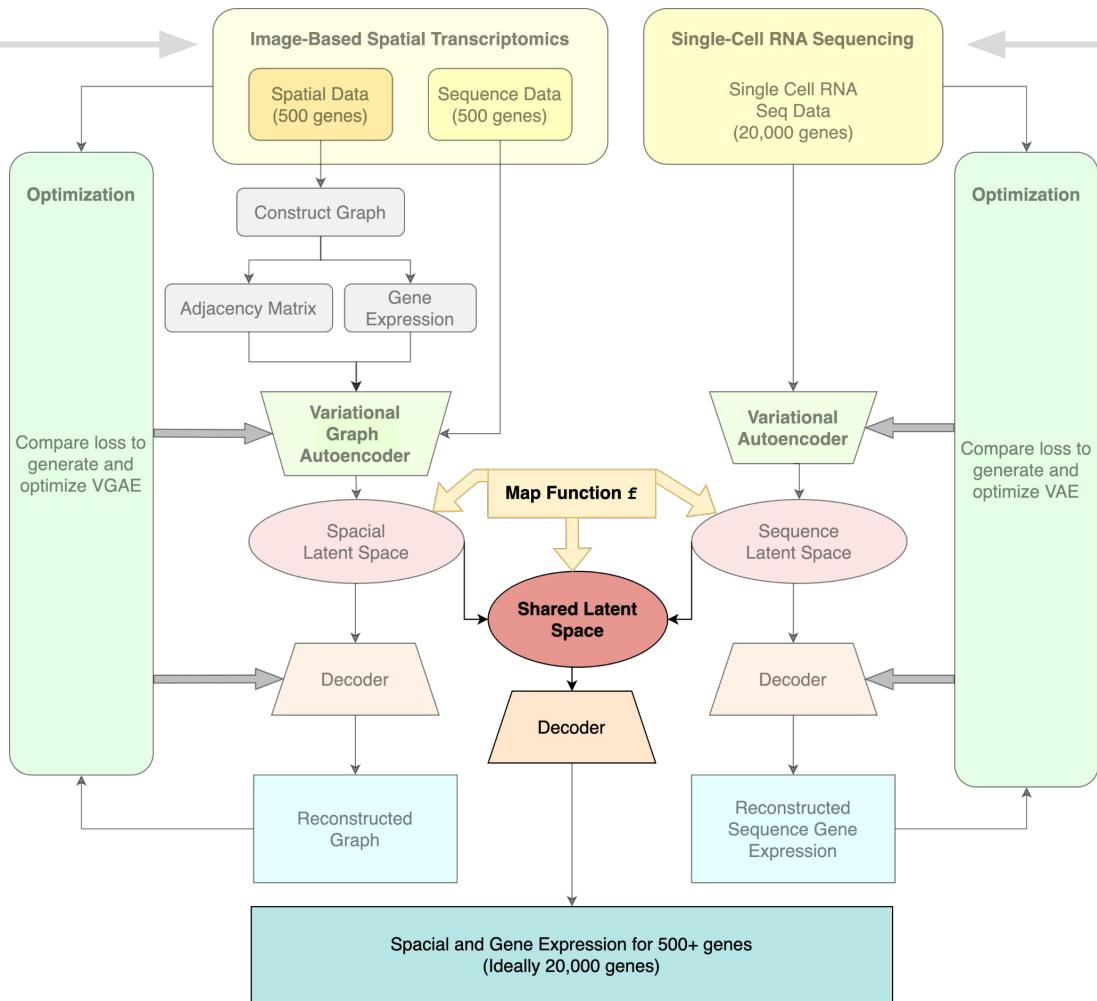
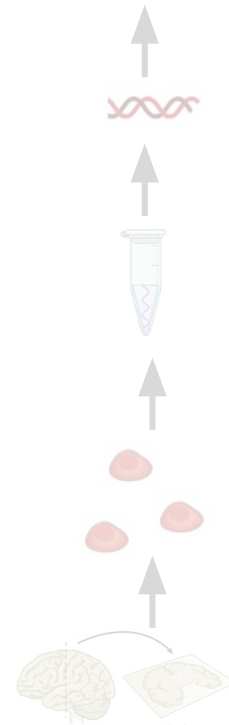


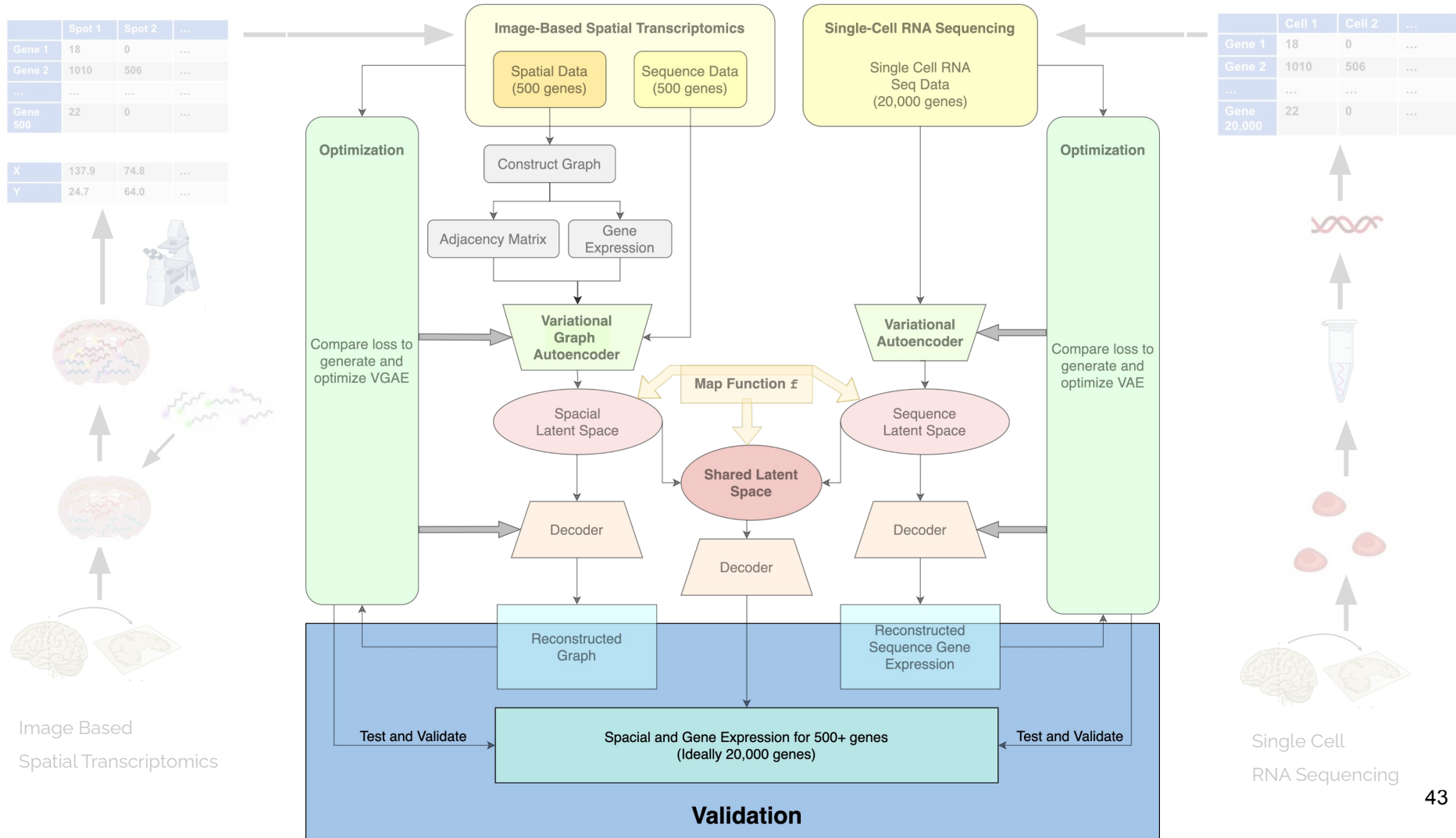
Image Based
Spatial Transcriptomics



	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...



Single Cell
RNA Sequencing



	Spot 1	Spot 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 500	22	0	...

X	137.9	74.8	...
Y	24.7	64.0	...

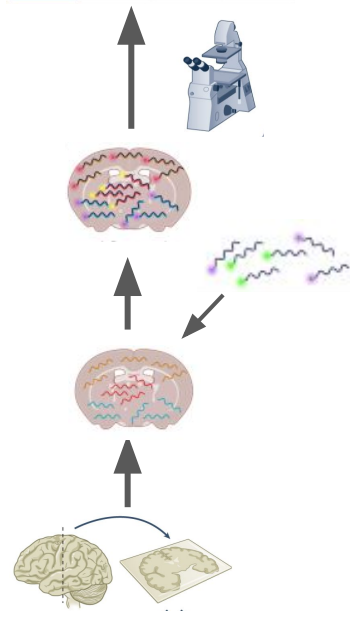
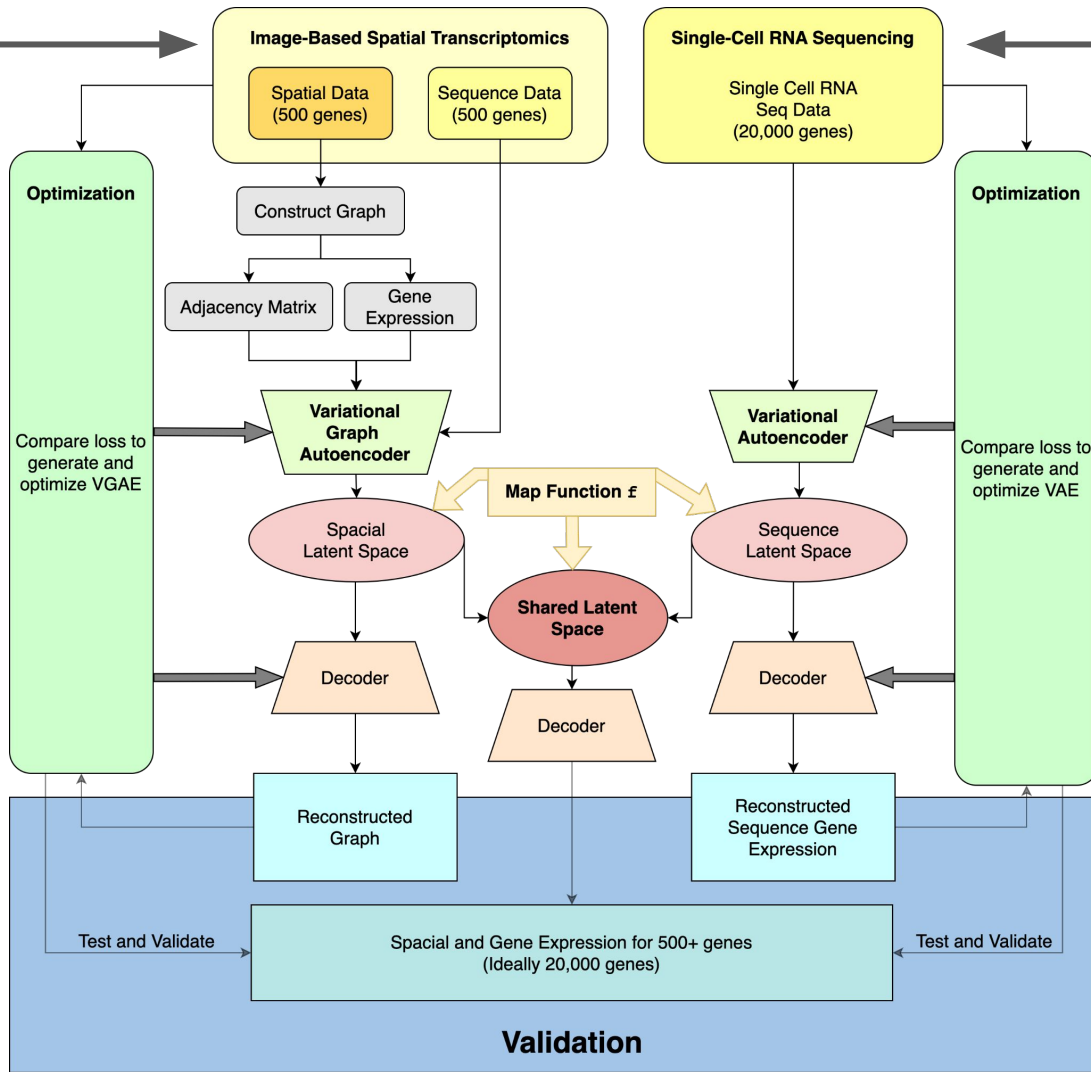
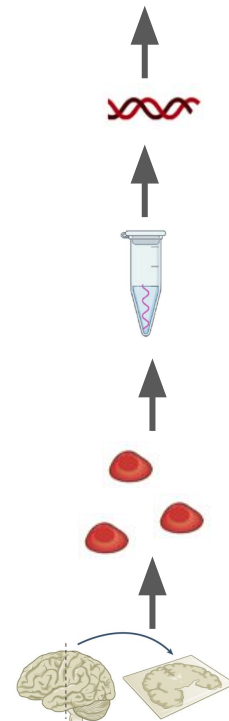


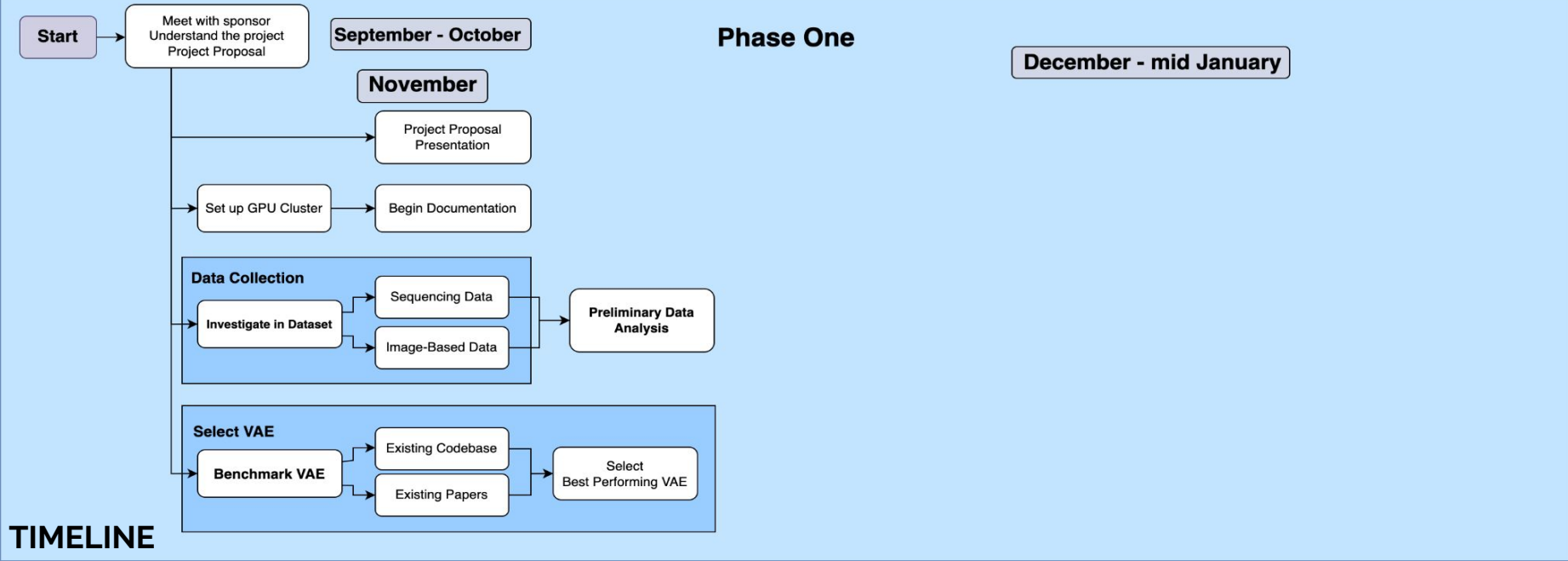
Image Based
Spatial Transcriptomics

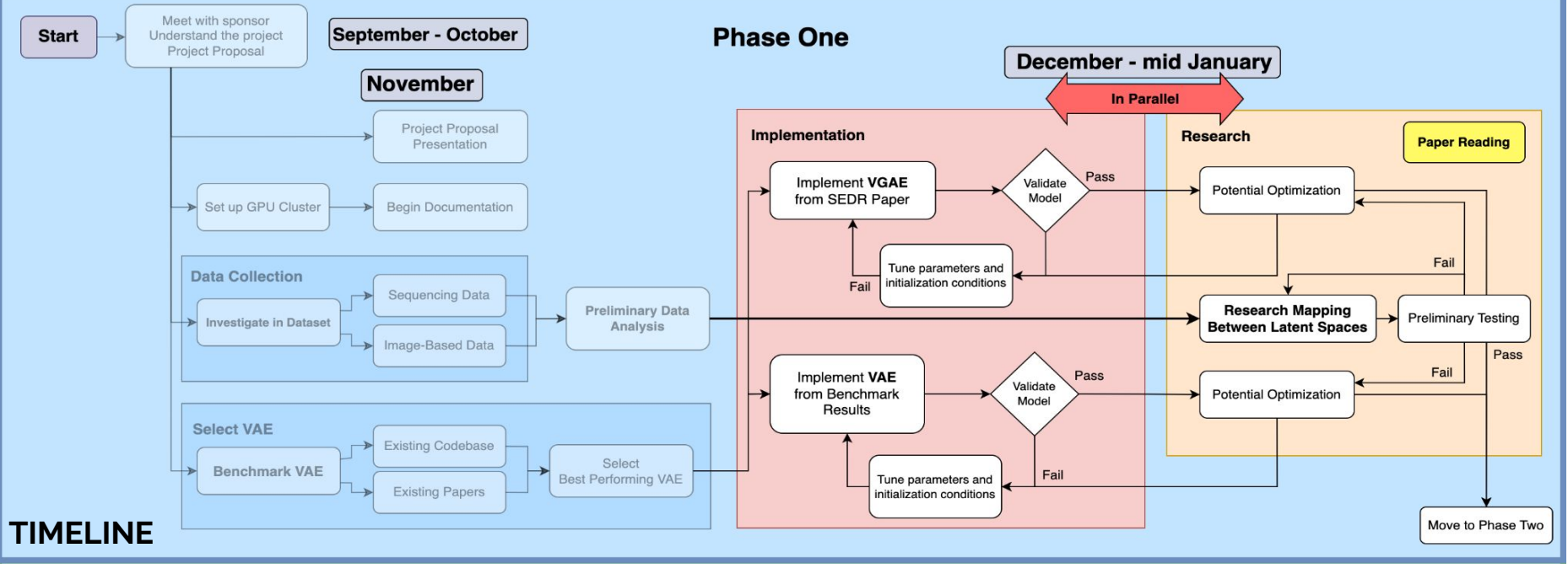


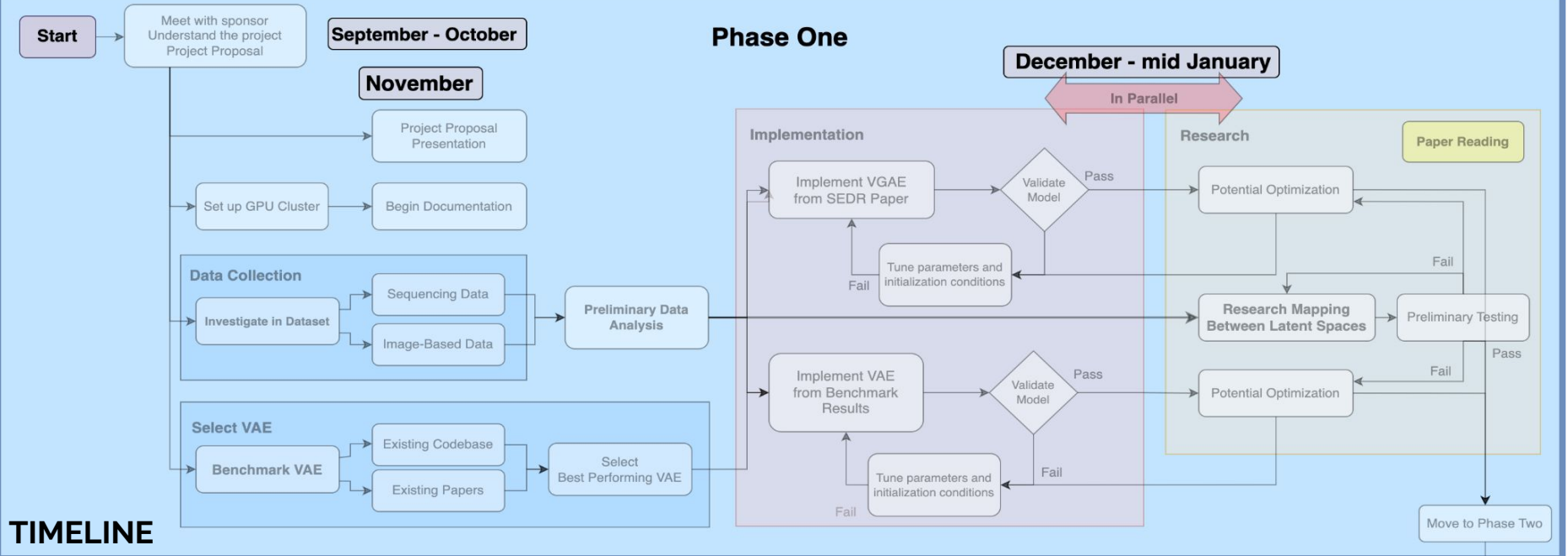
	Cell 1	Cell 2	...
Gene 1	18	0	...
Gene 2	1010	506	...
...
Gene 20,000	22	0	...

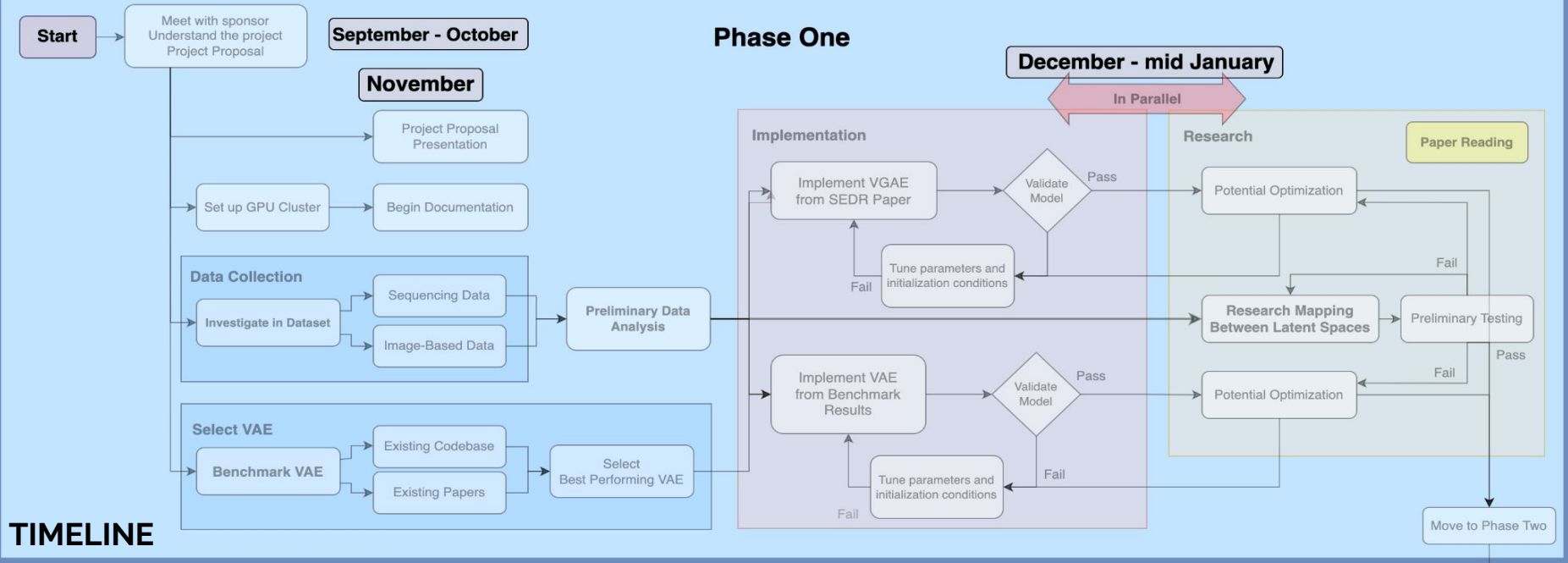


Single Cell
RNA Sequencing

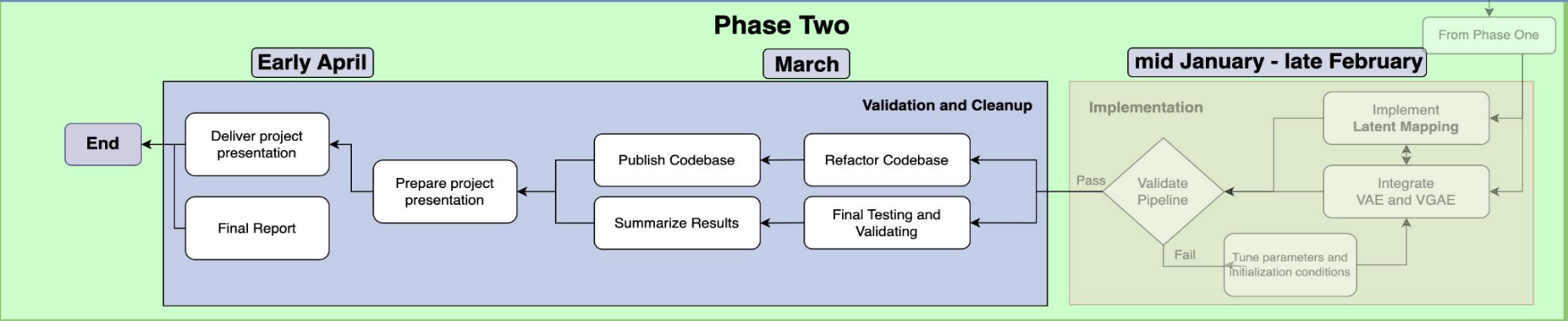


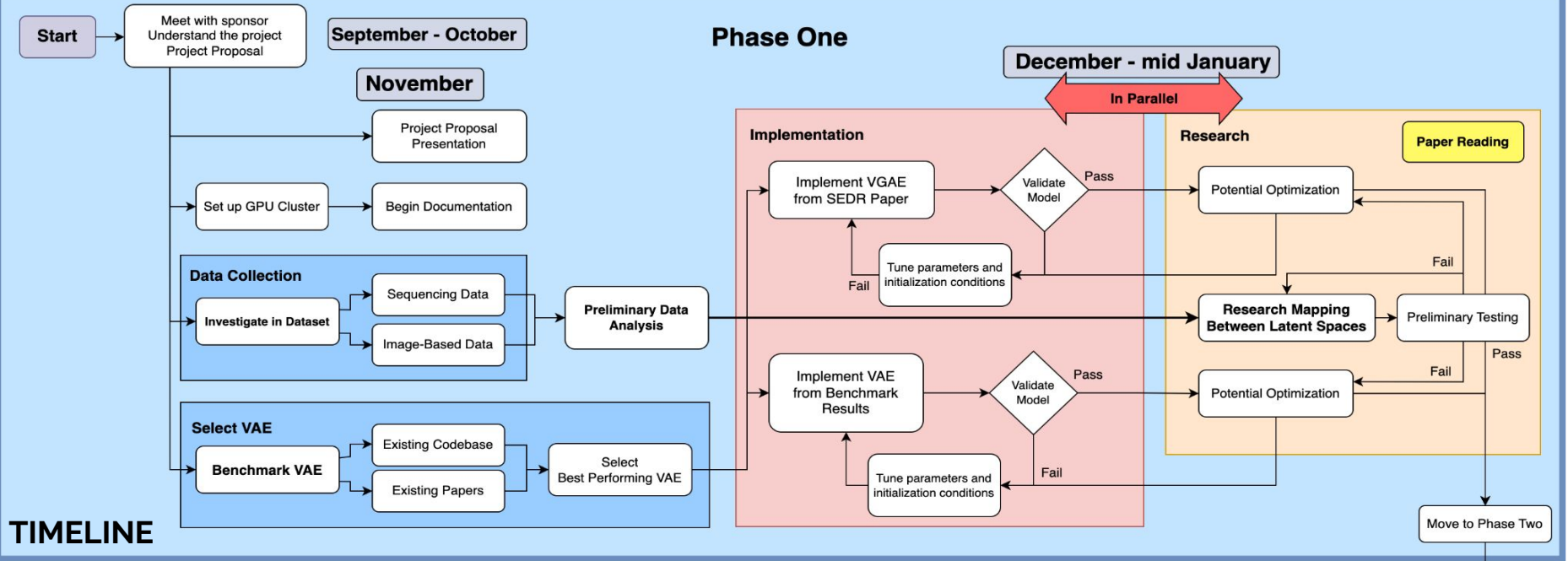




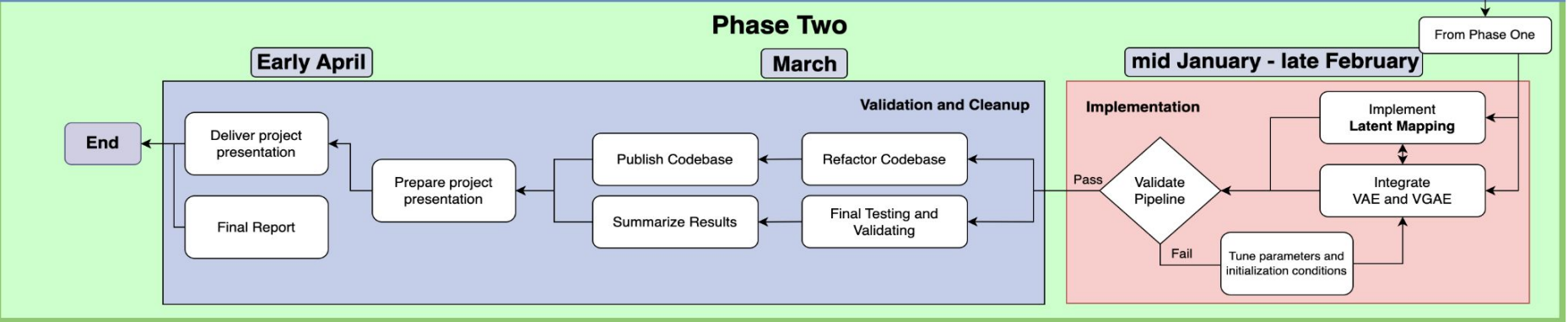


TIMELINE





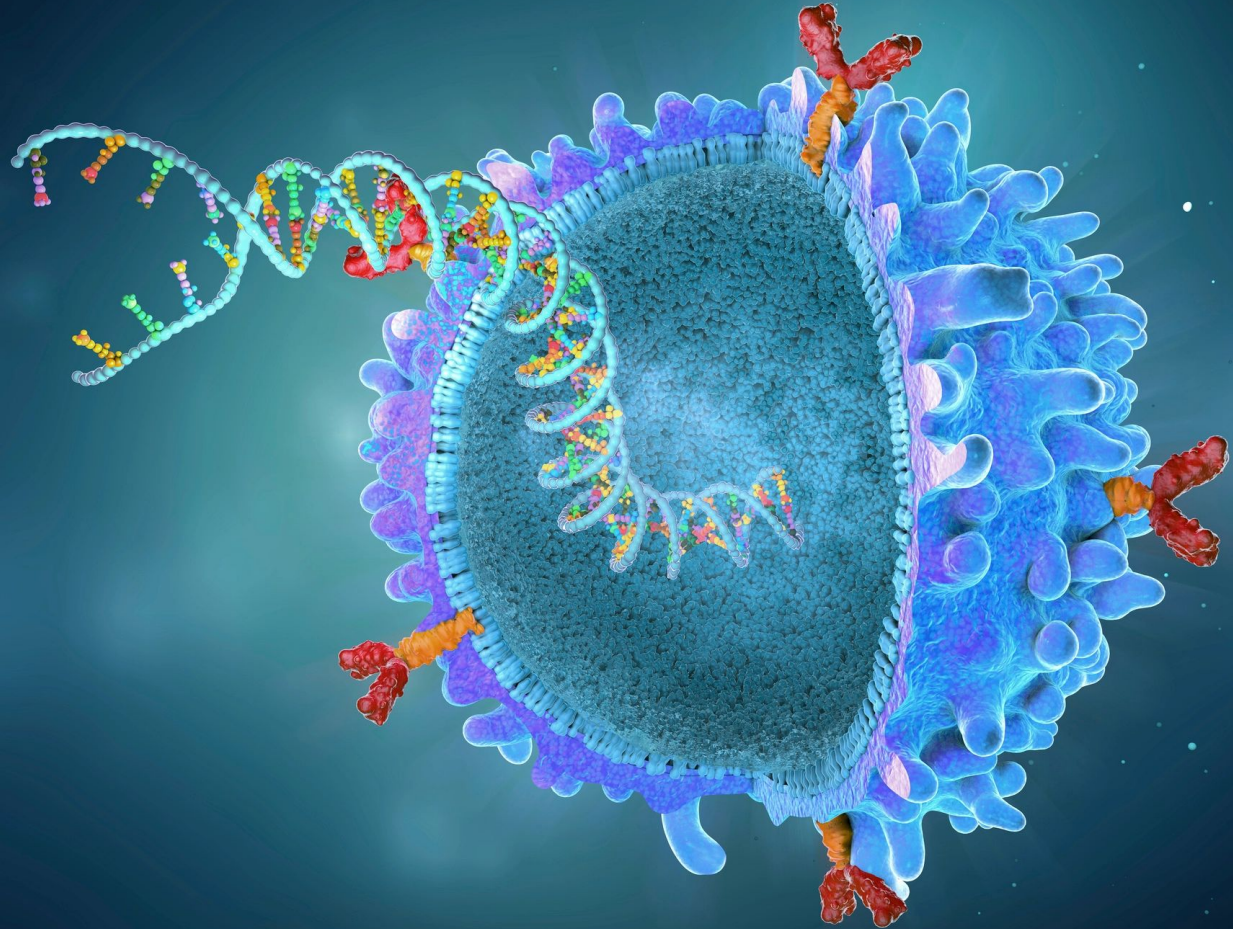
TIMELINE



Deliverables

1. **Fully functional VAE and VGAE pipeline with a map function** between two latent spaces that passes our validation methods, including
 - Successfully reconstructing missing genes for image-based spatial transcriptomics
 - Achieve acceptable accuracy for predicting on validation dataset
2. **Publishing codebase** up to open source standard

**Understanding
cellular composition
provide powerful
insights into
treatment of disease
and cancer**

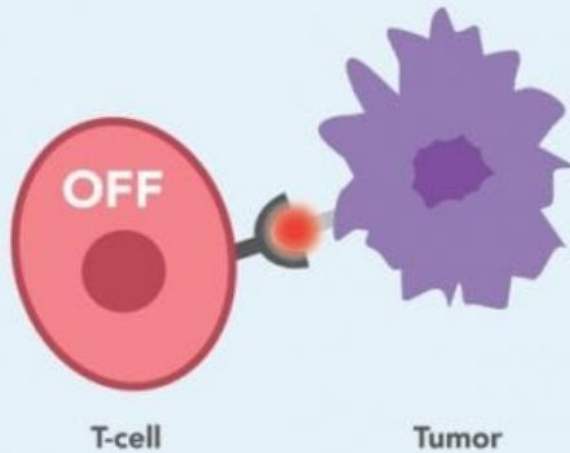


Thank you!

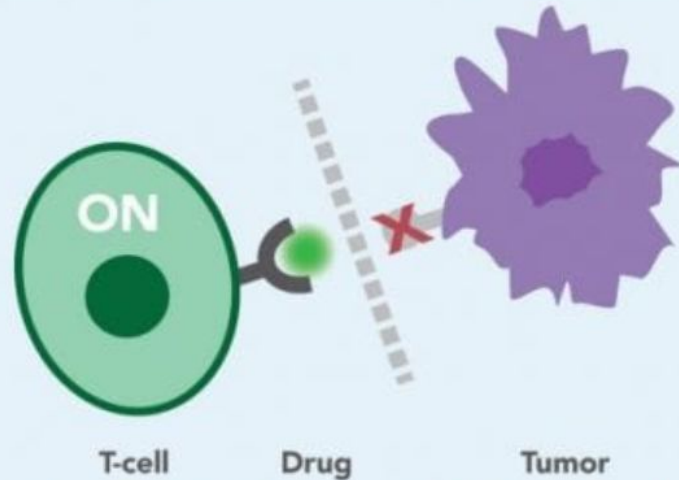
Appendix

How Does Immunotherapy Work?

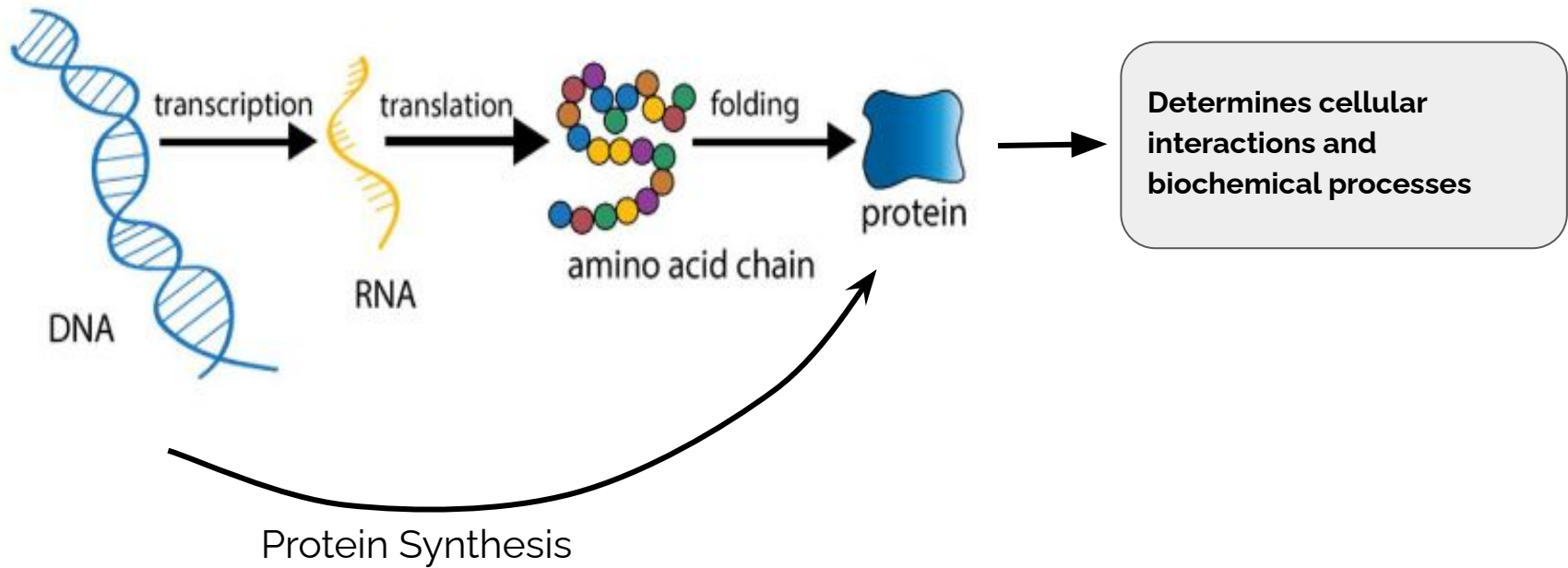
Tumor cells bind to T-cells
to deactivate them



Immunotherapy drugs can block
tumor cells from deactivating T-cells



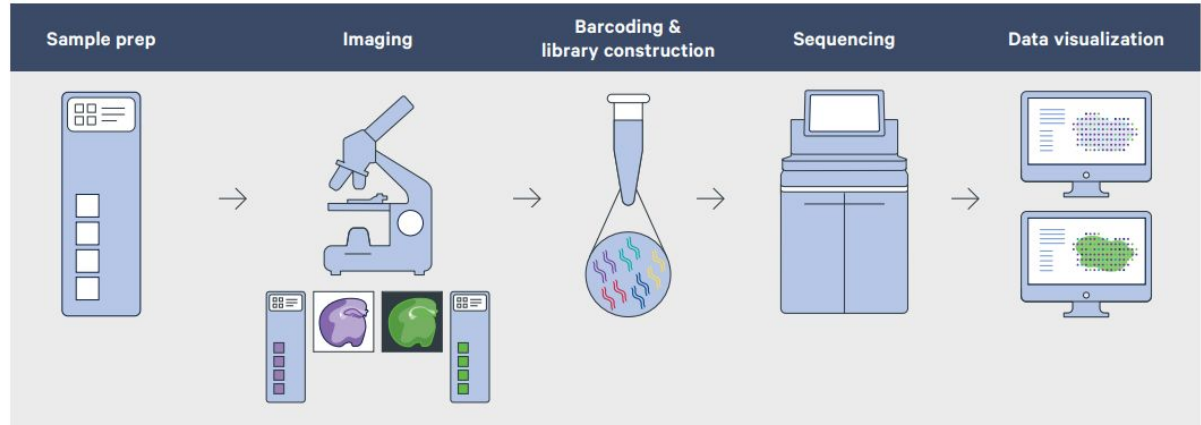
DNA and proteins



Producing spatial transcriptomics data

Another way we can study the RNA molecules inside the cells is by spatial transcriptomics.

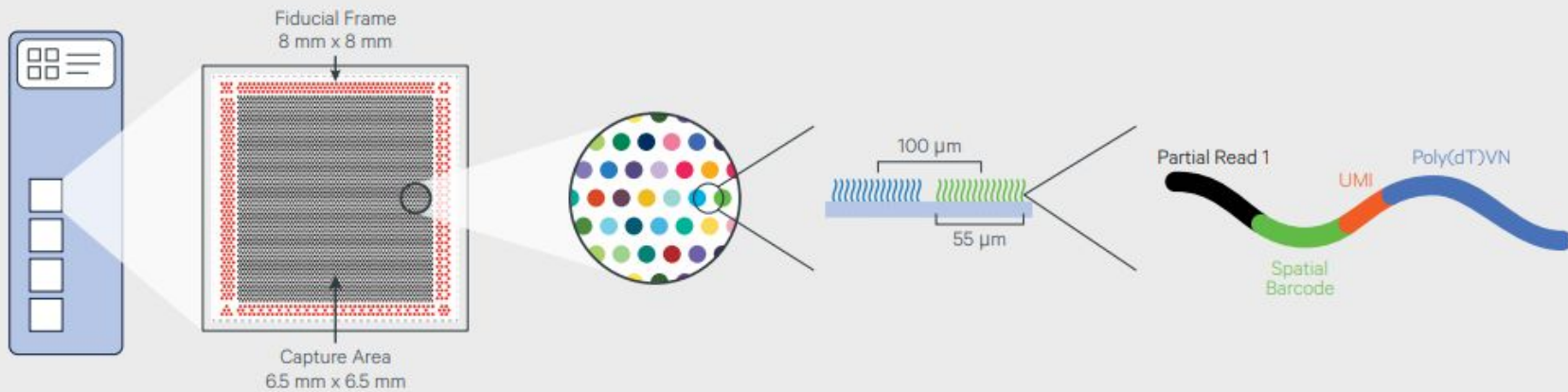
In addition to producing genetic data, just like sequential RNA data generation method, spatial transcriptomics also provides insight into the spatial information of cells (x and y coordinates of cells).



Visium Spatial Gene Expression Slide

Capture Area with ~5,000 barcoded spots

Visium Gene Expression barcoded spots



Variational Graph Auto-Encoders

SEMI-SUPERVISED CLASSIFICATION WITH GRAPH CONVOLUTIONAL NETWORKS

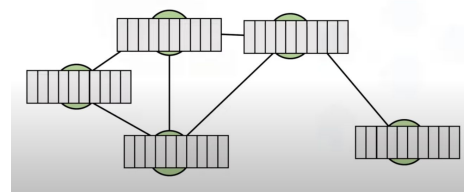
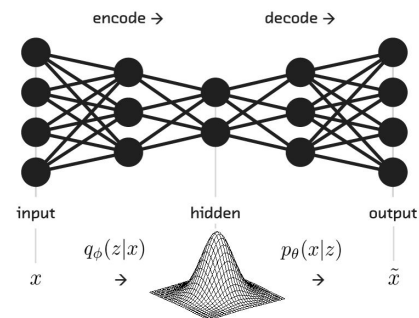
Thomas N. Kipf
University of Amsterdam
T.N.Kipf@uva.nl

Max Welling
University of Amsterdam
Canadian Institute for Advanced Research (CIFAR)
M.Welling@uva.nl

Thomas N. Kipf
University of Amsterdam
T.N.Kipf@uva.nl

Max Welling
University of Amsterdam
Canadian Institute for Advanced Research (CIFAR)
M.Welling@uva.nl

- Variational autoencoders
 - Encodes high dimensional data into lower-dimensional latent space.
 - Reconstruct original data by decoding latent space back into full-dimensional data
- Graph CNNs
 - Nodes contain vector data
 - Graph to contain relationship of each node to its neighbors
- Graph variational autoencoder can encode the data of an adjacency matrix of cellular data

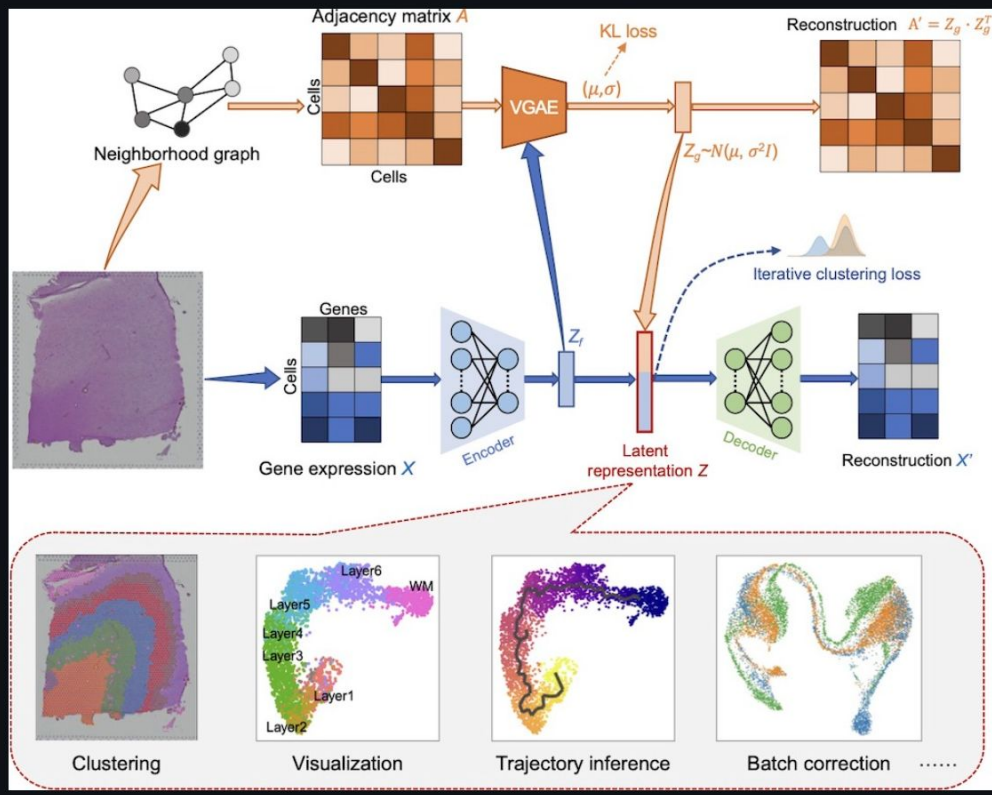


SEDR: paper + code

First step of our project

- Variational graph autoencoder to reconstruct adjacency matrix for spatial transcriptomics data
- Open source github repo with results on various spatial transcriptomics datasets

SEDR (spatial embedded deep representation) learns a low-dimensional latent representation of gene expression embedded with spatial information for spatial transcriptomics analysis. SEDR method consists of two main components, a deep autoencoder network for learning a gene representation, and a variational graph autoencoder network for embedding the spatial information. SEDR has been applied on the 10x Genomics Visium spatial transcriptomics dataset as well as Stereo-seq dataset, and demonstrated its ability to achieve better representation for various follow-up analysis tasks including clustering, visualization, trajectory inference and batch effect correction.



Ann-data

- Both sequence and iBST data are stored in this format.
- Helps us understand what is stored with the data and how to parse it

