

Computational Biology

Genetics \neq Genomics

 talking
about individual
genes.

CRISPR - CAS9

LV0 \rightarrow too much abstraction of biological situation.
 \rightarrow very minimalist system.

LV1 \rightarrow Using general published tools to understand particular situation.

LV2 \rightarrow

LV3 \rightarrow Computer comes up with hypothesis itself.

Sequences. → of nucleotides that form DNA and RNA

PDB - datasets - protein sequences.

SRB (NCBI) - Genome raw sequences .

limitations of sequencing.

- short reads
- imperfect reads
- Biased reads .

→ Microarray

Connectivity Map → MIT Project for
research of gene expression under different
environment.

level-1 Utilizing already published datasets using
established analysis methods to point
out defects.

BLAST - basic local alignment search tool

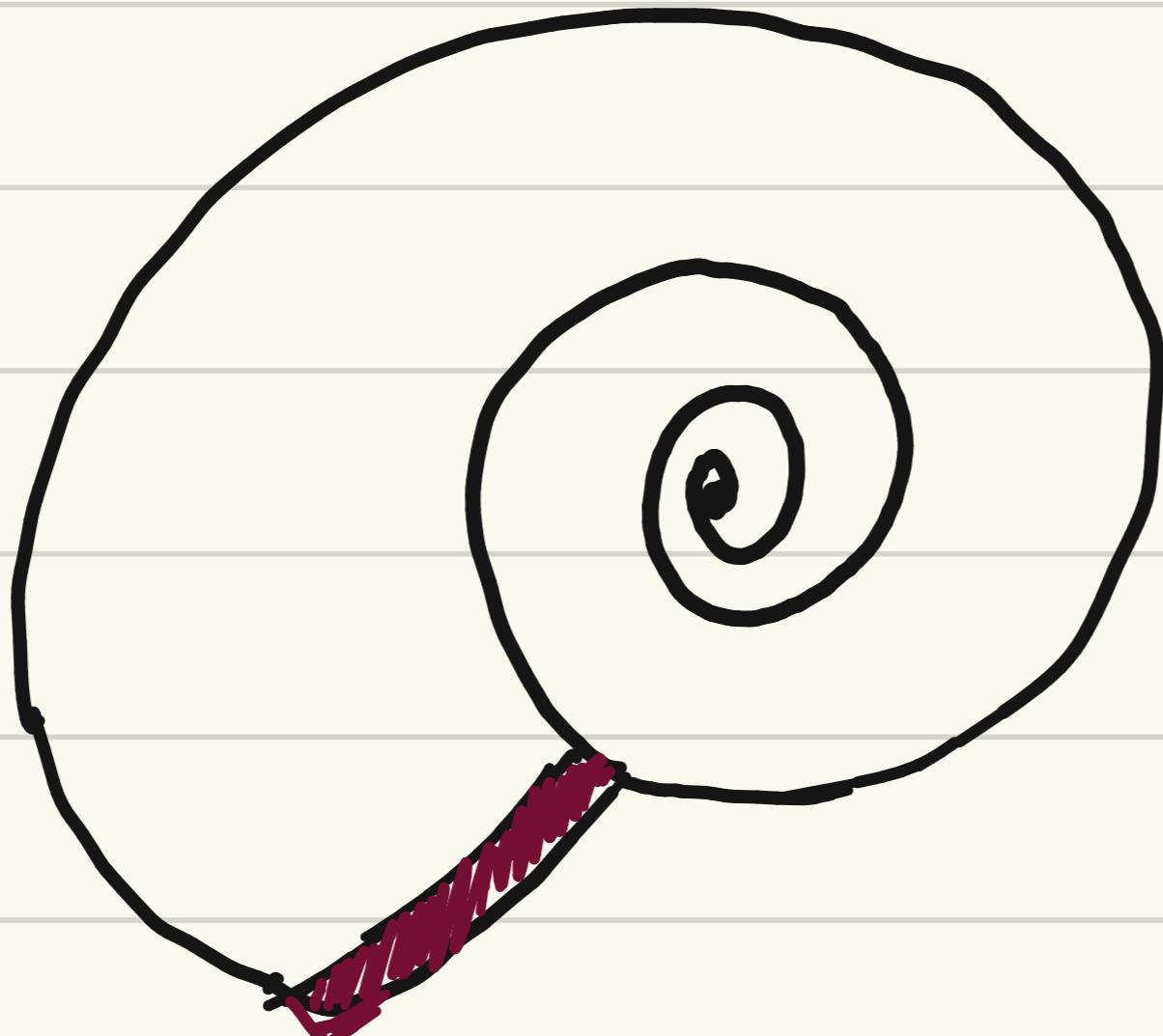
↑
online repository
for finding matching sequences
from the natural world.

level 3 - "Make biological discoveries using Comp."

looking for micro-structure in genomic data
to infer things like, cell stress

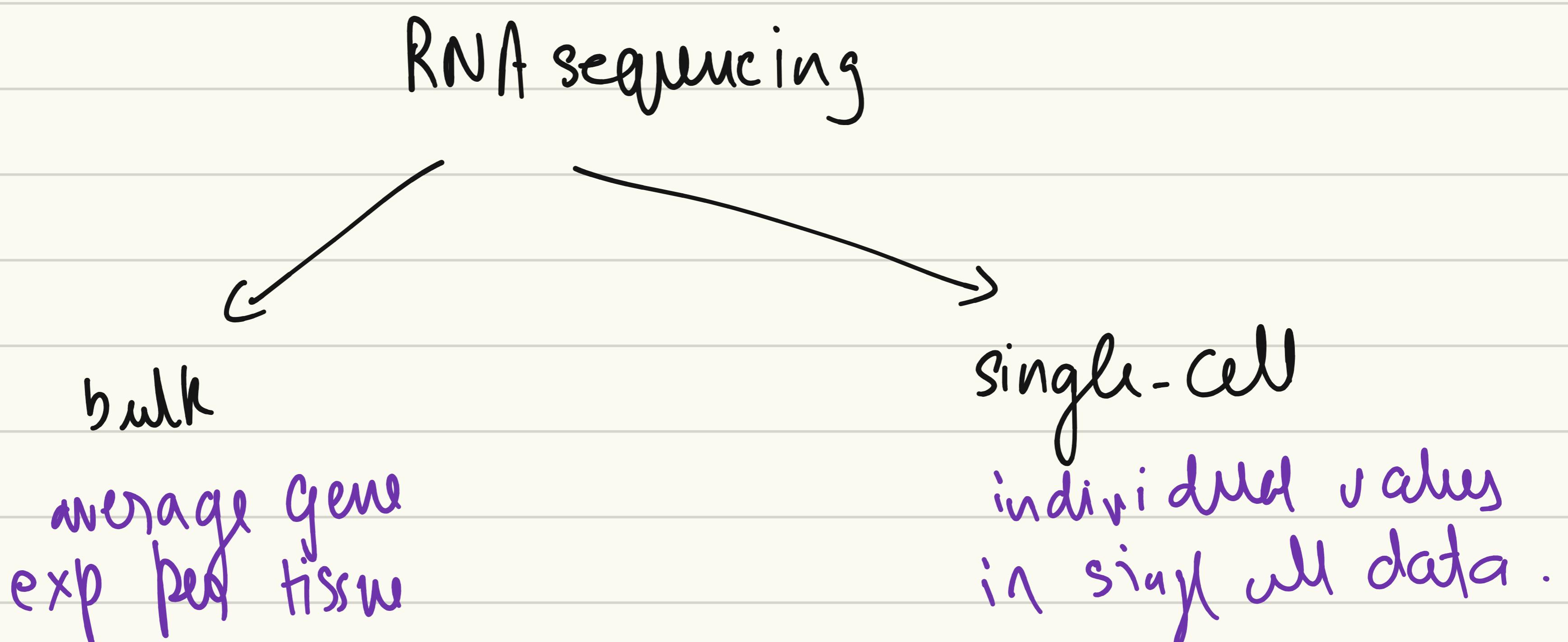
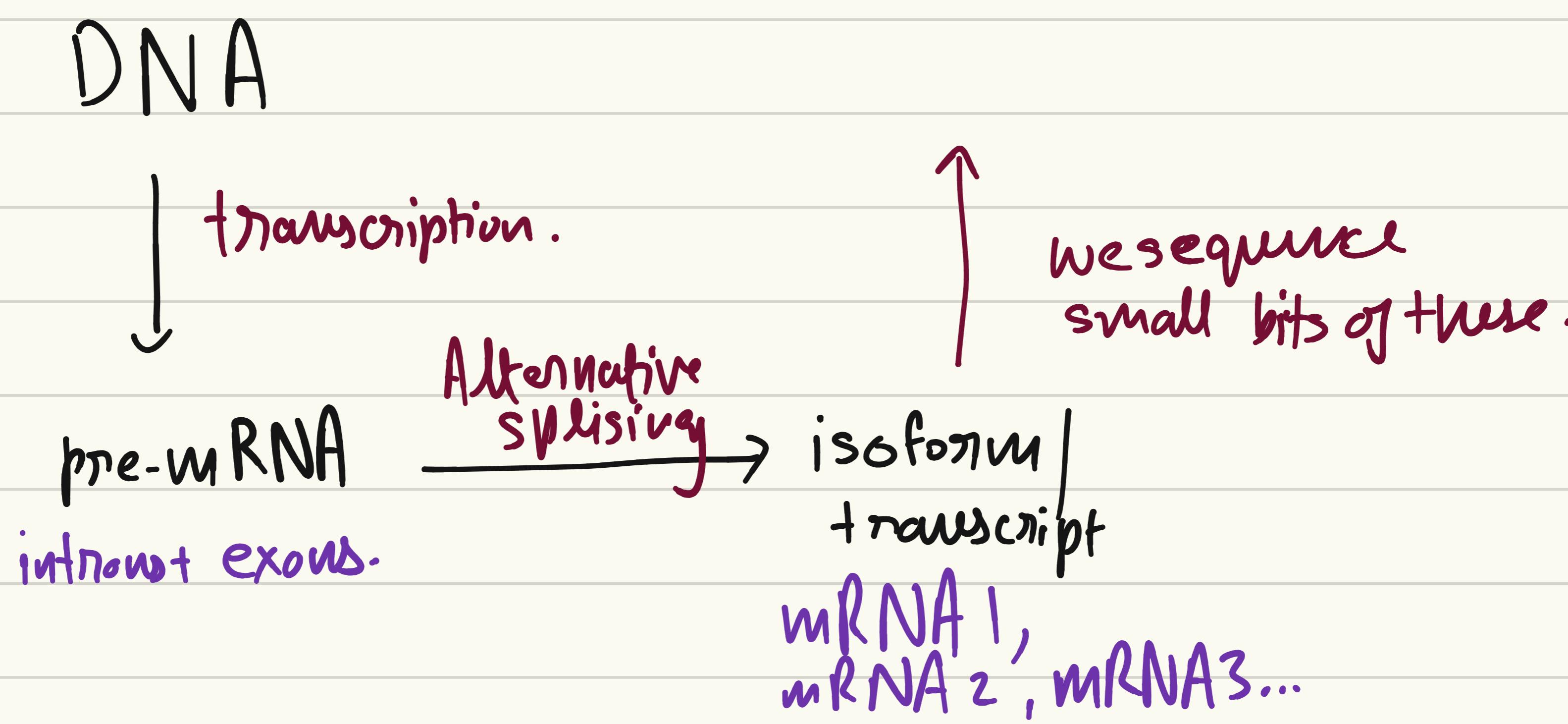
Human reference genome → first generation sequencing
by Human Genome project.

25,000 genes → 1,00,000 Proteins.



Single gene gives
rise to multiple
proteins

→ Not all exons are completely present in mRNA structure.



2016

ten years of next-generation sequencing



Paradigms (new)

Short read sequencing

- lower cost
- higher accuracy
- population level research

long read sequencing

- research (genome assembly)

Research field since 10 years

Latest →

- Database generation
- Routine clinical sequencing
- Pathogen DNA monitoring
- Population - level studies

CONS → - Error rates are still high

- Ethical and ownership issues.

Purpose of research - how genome sequence variation
underline phenotype / disease

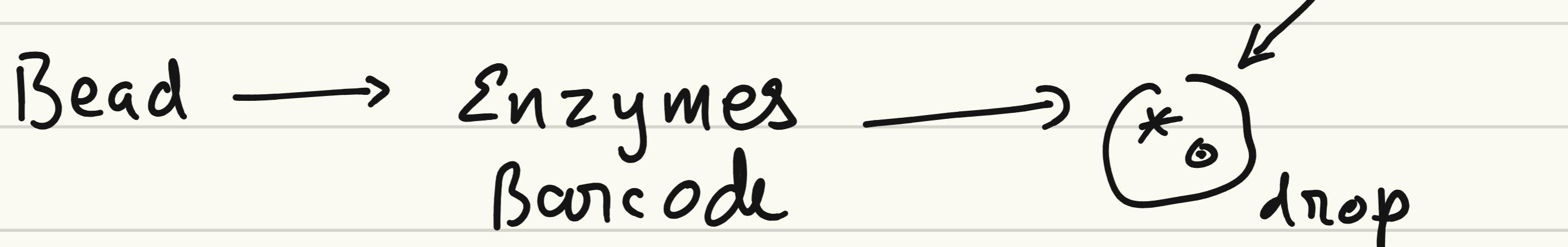
→ Read lengths affect the ability to resolve complex information.

- Sequence depth affects - the ability to detect low-frequency variants.
 - high error rates can give false-calling
 - There are also inherent biasing in sequences.
 - there are often annotation errors in datasets
 - Open data standards are used.
 - high quality genome-databases have comprehensive metadata
-

Lecture - 9/8

mRNA single cell sequencing exp.

- in 2015 Drop-seq was discovered which led to 100x increase in throughput.
- Droplet based scRNA-seq.



barcodes are unique sequences that are added to each RNA molecule of a cell.

AMPlification (PCR)

increases amount of the genetic material to work on.

UMI-count - total no. of unique molecules of a particular mRNA molecule.

greater transcriptome coverage in a cell is more important than no. of cells sequenced.

there are tools to find out how many cells one needs to sequence to find a particular rare gene.

technical replicates

vs.

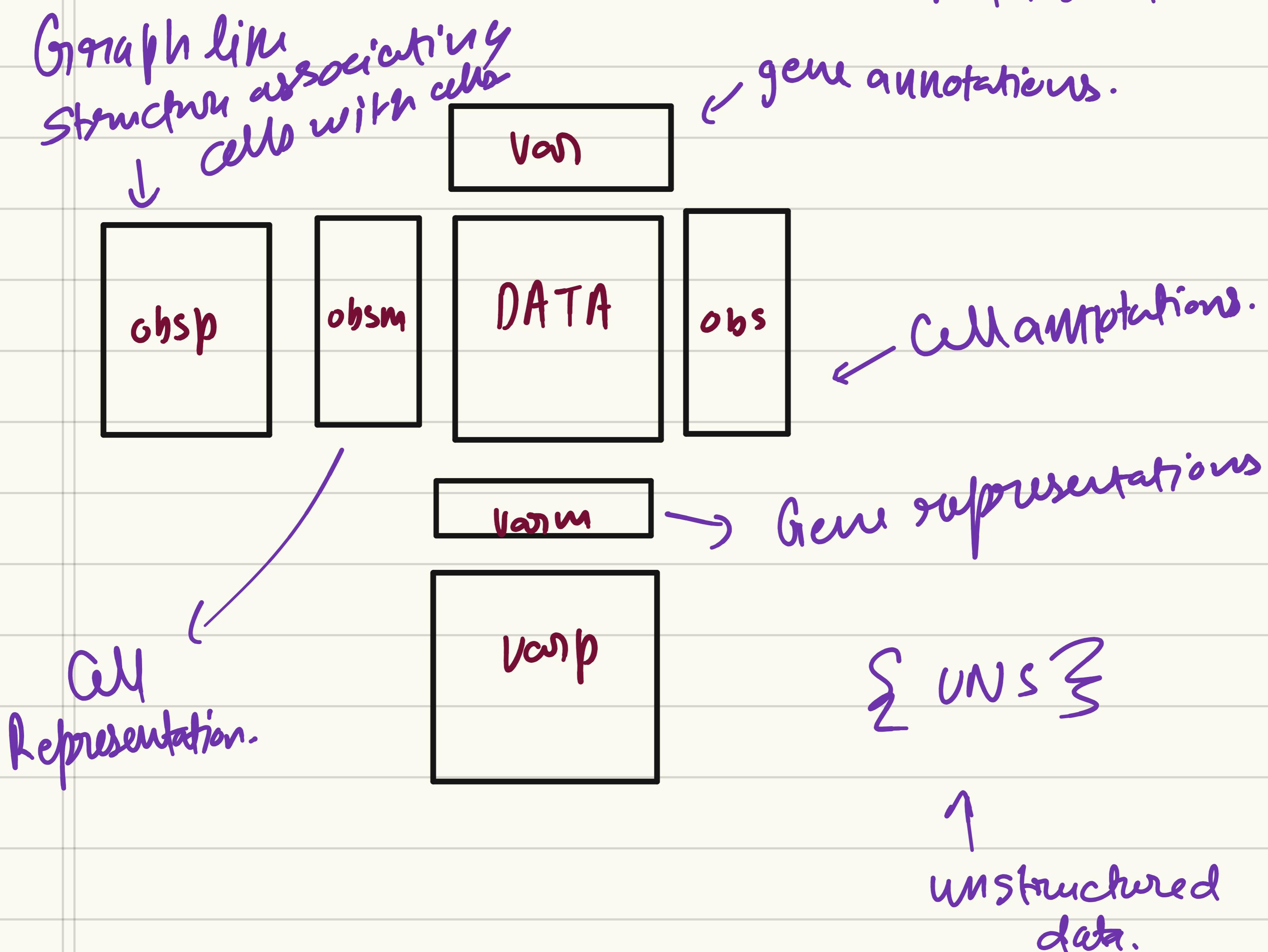
Biological replicates

Sequencing Many times so that accuracy increases

biological variation in the incoming material.

- 1 Dimension Reduction on cell gene matrix
- 2 Clustering
- 3 trajectory infer.
- 4 Annotation.
- 5 Differential Gene Exp.
- 6 Gene regulatory inference.
- 7 Cell-cell Communication infer.
- 8 gene expression program inf
- 9 pathway analysis.

SCAN Py ← Storing unimodal data with Ann Data



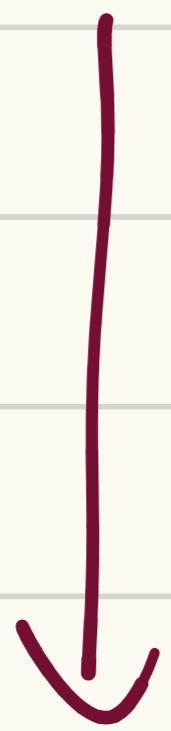
quality control → Normalization



feature selection. (more influenced by Biologics.)

Cell-Cell
distance.
Computation
(cells as nodes
values as edges)

← Dimensionality reduction. (PCA? +)



Unsupervised
Clustering.

Quality control

- lot of entries are zero
- performing QC also interferes with biologics.

Ensure that all cellular barcode data corresponds to viable cells.

lostinfo

Automatic thresholding

- MAD (Mean absolute deviation)

Ambient RNA Problem

some of RNA is expressed so much that they end up in every droplet.

- Ambient RNA is treated as noise.
- We separate this noise on distribution basis.

Doublets (two cells with same barcode)

- We introduce two cell's doublets artificially, we then carry out the principal component analysis.
- doublets mostly end up in one particular group in PCA.
- doublet detection is still an area of research.

NORMALIZATION

- make profile of each cell comparable.

Assumption's :
→ all cells in dataset initially contain equal no. of RNA molecule.
→ Count depth difference is only due to sampling.

Scaling is the solution of this

$$CPM_{ij} = \frac{g_{ij}}{R_j} \times sf$$

(counts per Million)

gene UMI count

total UMI Count.

- logarithmic transformation is also done.

$$\log(CPM_{ij} + 1)$$

- Regressing out biological effects.

effect of different life cycles stages of cell on data.

- We take cell populations at different stages of its life and check difference in PCA.

Feature selection.

Selecting suitable features.

- takes more time to process more features.
- we have to remove bad features
- Redundancy is also removed.
- dimensionality reduces.

Feature selection aims to exclude uninformative genes which do not represent meaningful biological variation.

Genes are binned by their mean expression and the genes with highest

- highly informative genes across cells will have a high deviance value.

Dimensionality Reduction.

→ Visualizing the dataset.

Why? → data - high dim.

- Noisy
- non-trivial dist.
- Non-linear Manifold.

PCA → linear combinations of existing dimensions to create new dimensions (axes) that have highest variations.

16/8

Curse of Dimensionality.

as the Number of genes increase
the variance that can be explained
by any one principal component
reduces.

Batch effect - Cluster Variation due to environmental factors.

Biological challenge. -

Biology messing with cell identity.

for cancer research biological research is important.

how to validate clustering from computation

→ Experimental differentiation. seems to be the best option.

→ care has to be taken to maintain quality and information of samples.

Clustering pitfall : overclustering
Cluster splitting

Clustering approach.

1 k-means : CS771

only K-clusters.
use \sqrt{d} distances from k-means
to classify.

2 hierarchical clustering :

We're going to look for 2 datapoints which are closest to each other.

then we repeat the process.

→ gives a tree structure.

→ cutting the tree at anyone level gives the clusters.

3 Graph based clustering.

→ very high data dimensionality.

→ may not be simple Gaussian distribution.

Using a Graph

$\{V, E\}$

vertex and edges.

Adjacency Matrix

→ sparse adjacency matrix

→ data matrix is also sparse.

dimensionality : $V \times V$

- Data still exists that what is the best measure of distance.
- Learn Manifold from data and then get distance in that manifold

sparse x sparse
data is most useful
for graph based
clustering.

Graph type

- K-nearest Neighbors.
- Shared - Nearest neighbour graph.
 - obtained from K-NN
 - similarity is based on nearest shared neighbours.

Slides have more information.

distance ← → Rank

shared neighbors.

edge \Rightarrow at least one shared neighbor in K-NN

edge strength \Rightarrow Rank

how to define a Community

Clique - a maximum complete subgraph

quasi-Clique - γ -dense
Close enough but not necessarily complete.

$$\gamma = \frac{E_s}{\frac{v_s(v_s-1)/2}{\text{total possible edges.}}} \quad \begin{matrix} \text{actual no.} \\ \text{of edges.} \end{matrix}$$

SNN-Clique
Algorithm.

- extract local maximal quasi cliques.
- merge quasi-cliques to get communities.

Graph Cuts

Cut the graph at edges to define a community

→ Change objective fun. to account Size and Volume of the community.

→ Ratio Cut
number based

→ Normalized Cut.
volume based.

→ finding minimum on normalized ratio is NP-hard.

cut

Modularity optimization

measure of structure.

density of connection

within a module

SEURAT (R package)

shared overlap
in local
neighborhood

k-nearest
neighbour

→ Jaccard distance.

Very
noisy.

modularity optimization

LOUVAIN ALGORITHM

weight of connection (P, Q)

$$Q = \frac{1}{2m} \sum_{P,Q} [A_{PQ} - \frac{K_P K_Q}{2m}] \delta(C_P, C_Q)$$

$$K_P = \sum_Q A_{PQ}$$

probability of
a connection by chance.

$C_p \leftarrow$ cluster to which p is assigned.

$$\delta(C_P, C_Q) = 1 \quad \text{if } C_P = C_Q \\ \text{otherwise } 0$$

$$M = \frac{1}{2} \sum_{PQ} A_{PQ}$$

LEIDEN ALGORITHM

Measures
of Similarity.

{ ARI -
NMI -

CellSIVS ← clustering for rare all
population.

Cell type annotations

→ process of labeling clusters on basis
of phenotype.

Association in each clusters a set of
genes which are highly expressed
in this cluster but not otherwise.

Scanpy is a Python package for single-cell RNA sequencing (scRNA-seq) analysis. It provides a user-friendly interface for preprocessing, analyzing, and visualizing scRNA-seq data. Scanpy offers a wide range of functionalities, including quality control, normalization, dimensionality reduction, clustering, differential expression analysis, trajectory inference, and cell type annotation. It also integrates with other popular scRNA-seq analysis tools, such as Seurat and Cell Ranger. Scanpy is widely used in the single-cell genomics community and has become a standard tool for scRNA-seq analysis.

Scanpy is a software tool that helps scientists analyze data from single-cell RNA sequencing experiments. This type of experiment allows researchers to study the gene expression of individual cells, which can provide valuable insights into how cells function and interact with each other.

Scanpy makes it easier for scientists to process and analyze this data by providing a user-friendly interface. It offers a variety of useful features, such as checking the quality of the data, normalizing it to make comparisons easier, reducing its complexity to identify patterns, grouping similar cells together, identifying genes that are differentially expressed between cell types, predicting cell development paths, and labeling cell types.

In addition, Scanpy can work together with other popular tools used in single-cell RNA sequencing analysis, such as Seurat and Cell Ranger. This makes it even more powerful and flexible for researchers.

Overall, Scanpy is widely used in the scientific community and has become a standard tool for analyzing single-cell RNA sequencing data. It helps researchers make sense of complex data and gain a better understanding of how cells work.

AlphaFold - Deepmind.

Date: |

Page: |

simulating protein folding using AI

→ Physics based modelling tool - Rosetta.

→ DATA - Worldwide Protein database.

(All known protein structure):

Multiple sequence alignment

"3D structure of a protein
is much more stable than
than sequence that code for protein"

AlphaFold - 1

Same protein
(different species)

↓ from similar evolutionary history

different ~~species~~ sequence
but same function and structure.

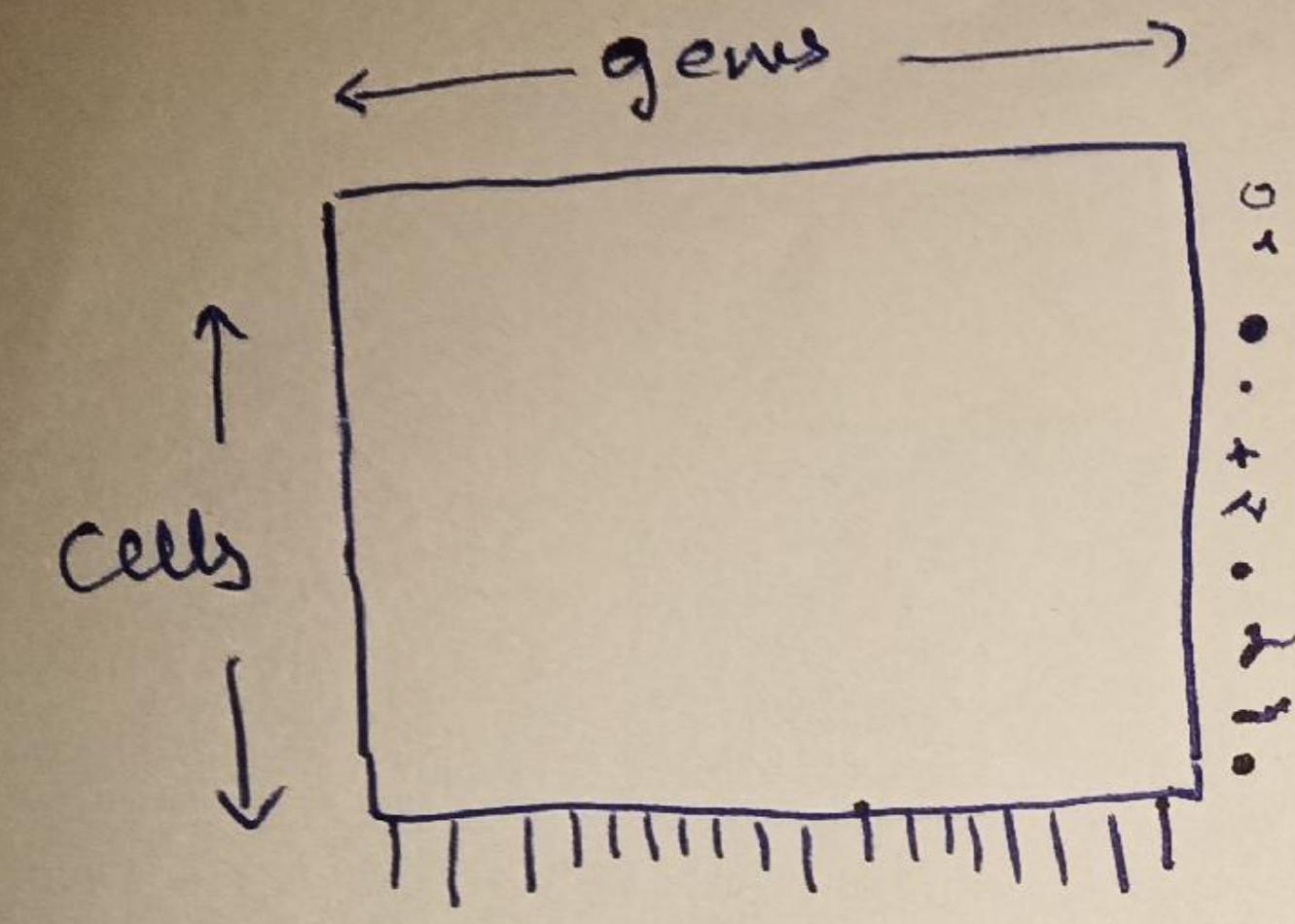
Distogram - 2D array of distance between every pair
pair of residues in protein.

(independent of, rotations / translations)

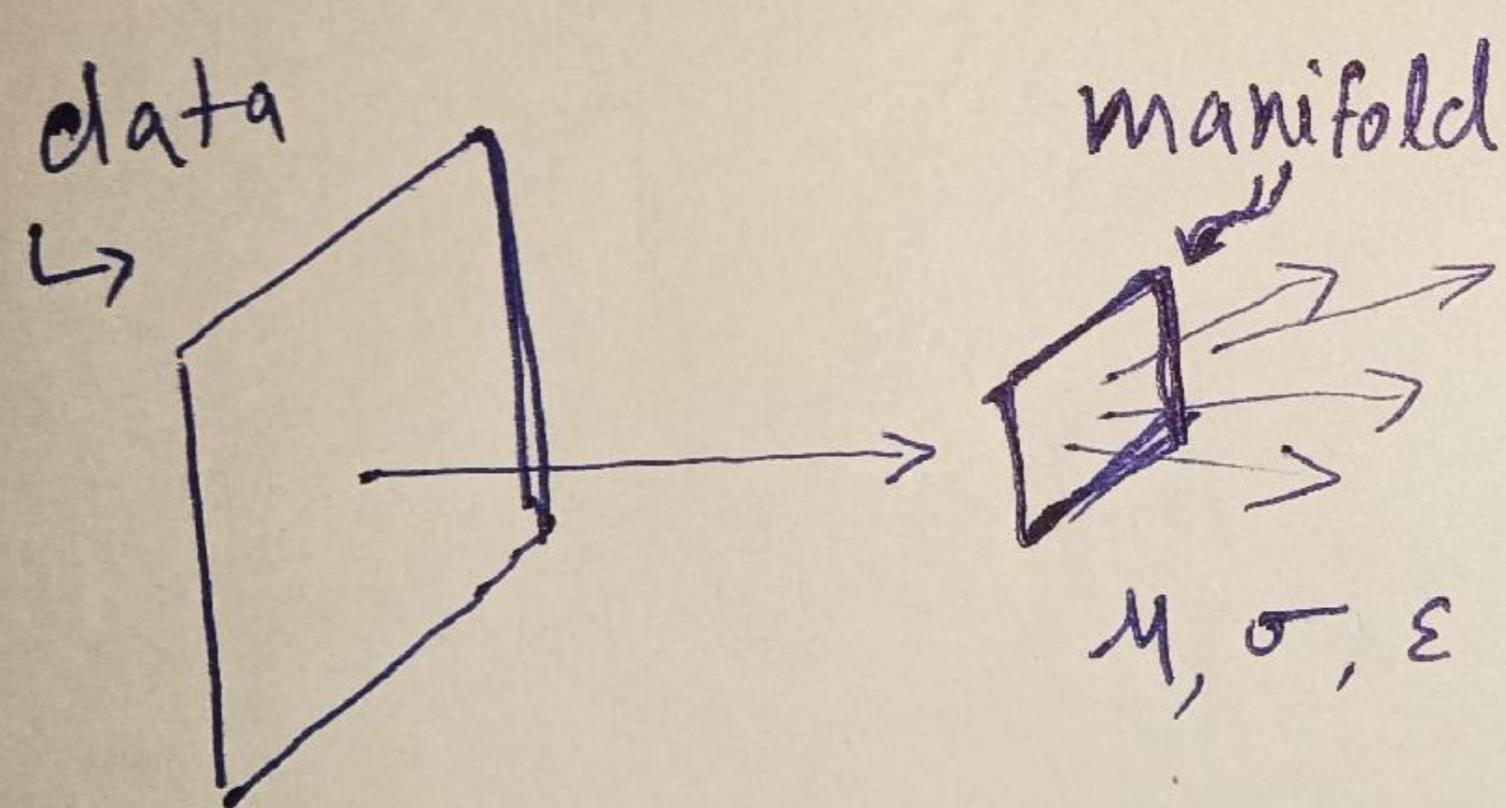
Sequence + MSA features → DNN → Distogram → Gradient descent
on protein specific potential
and prediction scores.

AlphaFold - 2

Sequence → MSA → Evoformer.



1 2 3 4 5
3 4 5 6 7
1 2 3 4 5 6
.....
.....



measurement (hyperparameters, λ)

maximum
possible correlation
v. causation

- learning dimensions that capture maximum variance.
- change in axis which we represent the data.
- implement w_i such that the variance $w_i^T w_i$ is maximised.

- $\arg \max_{w_i} (w_i^T w_i) \leftarrow$ objective function
- Find out: first principal component.

$d=20 \Rightarrow$ take top 20 principal components. (eigenvalues)

etc.
PCA is identical to classical scaling.

\uparrow
Euclidean
distance between an
higher and lower dimension

→ PCA retains only large pairwise distances, but local distances are lost.

ISOMAP → preserves pairwise geodesic distances.

\uparrow
distance measured
over a manifold.

- first compute k-NN graph using euclidean distance.
- then apply preservation of pairwise geodesic distance.

CS690 (1 Sept)

Mutual nearest neighbour

Mutually similar expression profile between different batches.

pair specific batch correction vector = vector diff b/w expression profiles of paired cells

↓

cell-specific batch correction

Segment V3

Canonical Correlation analysis(cca)

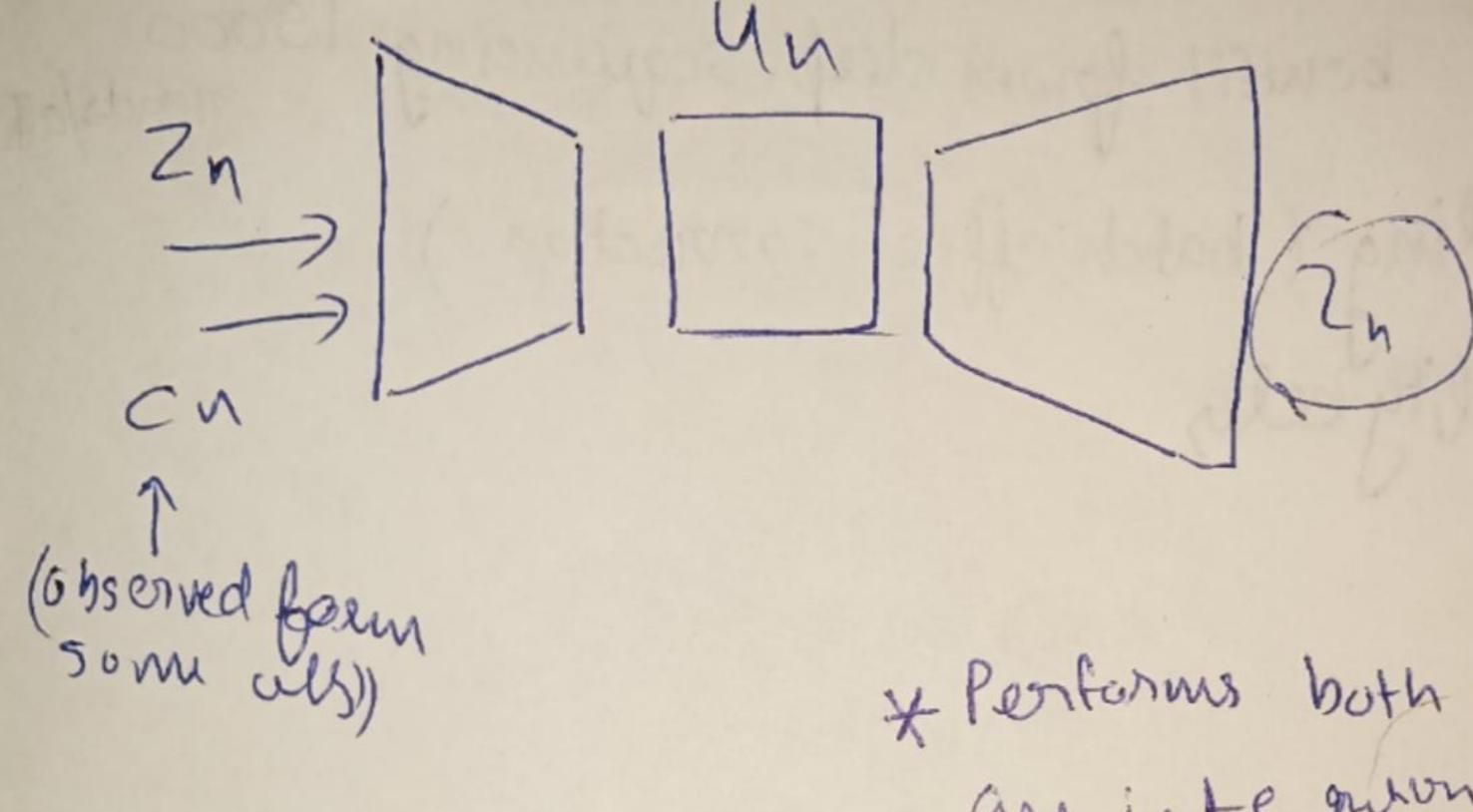
→ Correlated sources of variations between two datasets

L2 Normalization → Correcting in differences.
↳ in global scale
Vector normalization.

If there is a particular type of cell that does not have any correspondence with the other batch MNN should not classify these.

Supervised integration (SCANVI)

$z_n \Rightarrow$ latent space corresponding to cell



→ Any integration method has to ensure that biological trajectory trajectory

KBET - K nearest neighbor batch effect.

- if well mixed data
- then if you choose one neighbor b/w all

Batch effect (additive effect)

CCA

↓

L2

↓

Anchors

(high scoring)

Once the data is properly integrated and clustered we can use the marker genes to annotate cells.

LABEL transfer. (surat accuracy pretty good)
transfer learning,

Reference data → Query

Harmony (integration method)

→ fuzzy clustering

* for simple integration problems harmony well works
Clusters +
Multiple batches → diversity within cluster (max)

UMI - Unique molecular identifier.

Sparse gene expression matrix

→ substantial benefit from deep sequencing 15000 reads/all

→ Surface labeling (batch effect correction.)

→ low quality cells

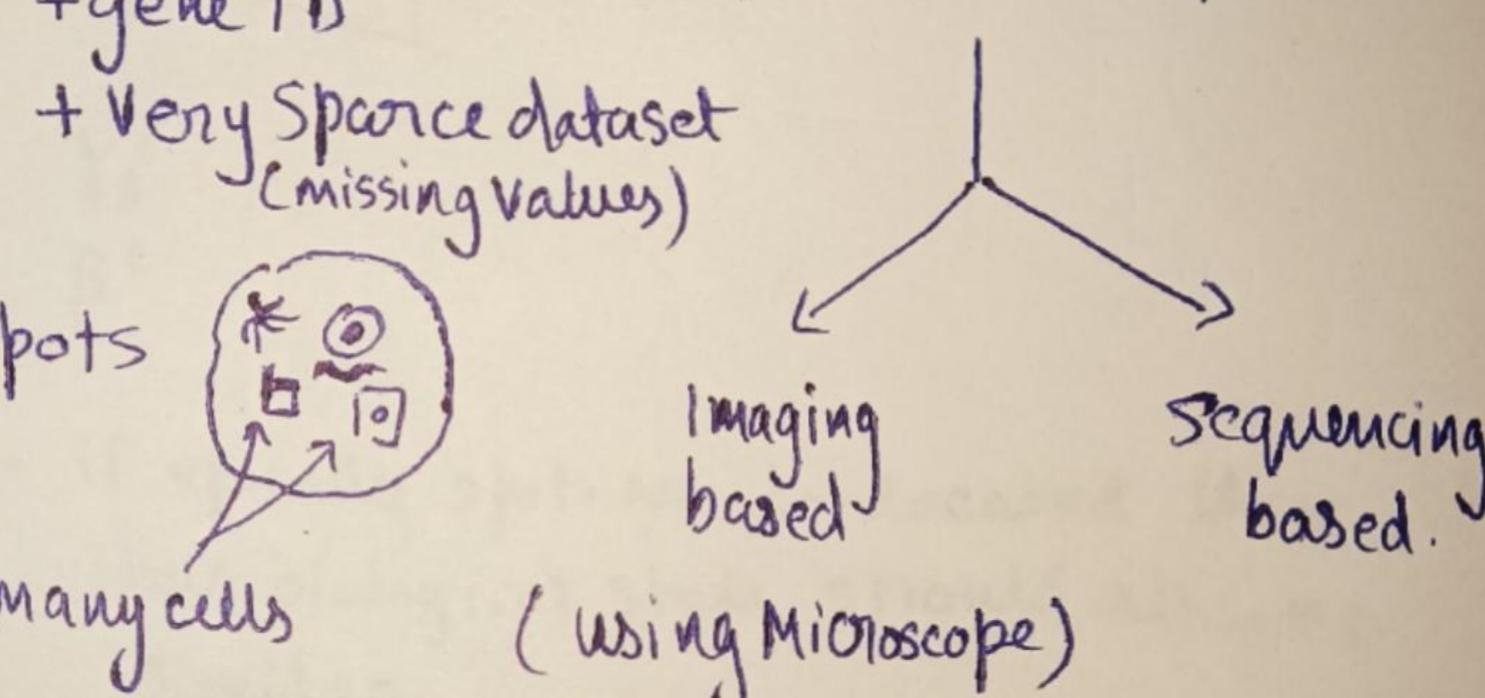
8 CS690 (8 sept)

NGS tech for spatial data

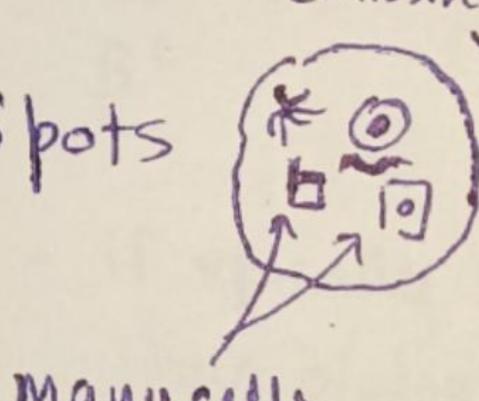
→ 10x Visium (Most used)

→ Stereo-seq (Most advanced)

General problem → image segmentation
+ gene ID
+ Very sparse dataset
(missing values)



Basic Unit → Spots



Many cells

(using Microscope)

Spot size → 10xVisium (2019)

Eg: FISH

= 50 μm (~10 cells) (Less genes)

→ slide-seq

= 10 μm (~1-2 cells)

Spatially Variable Genes ← genes selectively expressed
based on spatial params.

Cell type
Deconvolution

Supervised (sc-RNA-data)
Unsupervised.
(topic modelling / VAE / HMRF)

Principle

①

Tissue → SC (Single cell)

+ ST (spatial)

Spot gene expression signature = linear combination of
(Single cell signature)

RCTD (method)

i - spots

j - genes.

k - cell types.

$$\rightarrow \gamma_{ij} = \text{Poisson}(N_i, \lambda_{ij})$$

$$\log(\lambda_{ij}) = \alpha_i + \log\left(\sum_{k=1}^K \beta_{ik} u_{kj}\right) + \gamma_j$$

↑
Spot variation

↑
comes from
se-seq data.

Cell2Loc (Method) (supervised.)

→ Probabilistic (likelihood, prior)

↓
Posterior
for LV ← cell type
proportions

DestVI (latent variable model (LVM))

→ depending on spatial context reference signature
can vary.

→ Sharing of latent space in sc and ST

Cell type sig(SC) → γ_n | D

Spatial Transcriptome → γ_B^c
 β^c

Splice Mix - if spatially spots are co-located then
their biological state should also be
similar.

HMRF ← Hidden Markov random field.

{ there is some sort of function that relates the
hidden biological state to spatial information.

Markov - Current state is only dependent on neighbouring
state.

Can be arbitrary
and complicated

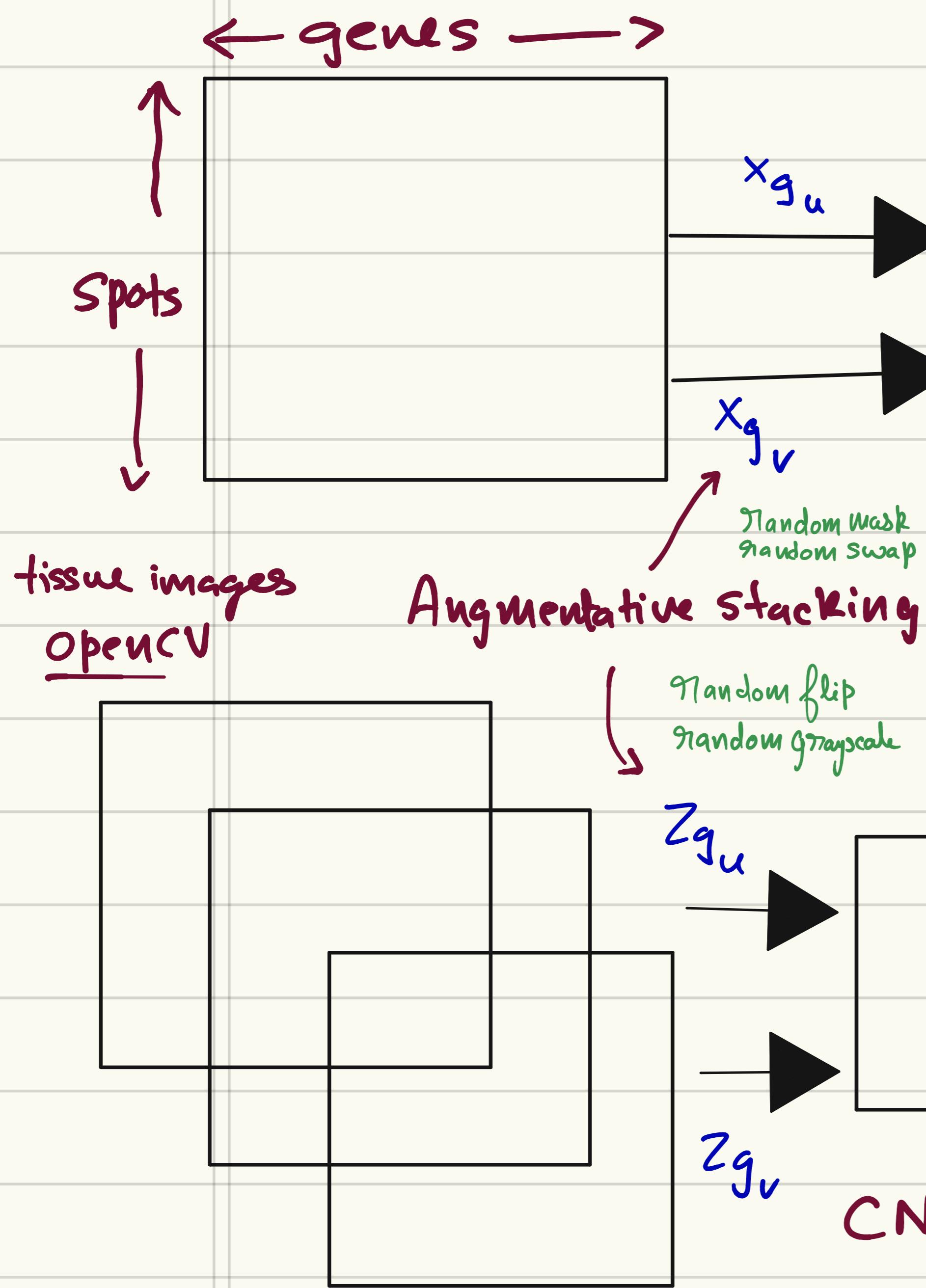
$h_n | h_{n-1}$ → Markov property
On $| h_n$ → Markov property.

Assignment - 2

Clustering and classification using ConGI

Contrastive learning.
in latent space

Spot associated
gene matrix



Non-normalized temp-scaled cross-entropy loss

$$\text{Sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

$$l_{i,j} = -\log \frac{e^{\text{sim}(h_i, h_j)/\tau}}{\sum_{k=1}^{2N} e^{\text{sim}(h_i, h_k)/\tau}}$$

τ temp parameter

Shared Space

SimCLR contrastive loss.

$$L = L_{zi} + \lambda_1 L_{g2g} + \lambda_2 L_{gi}$$

"Return only encoder training, throw away projection head params!"

L = averaged pairwise loss

$$L = \frac{1}{2N} \sum_{K=1}^N [l(z_{i-1}, z_K) + l(z_i, z_{i-1})]$$

final representation $z = z_g + \alpha z_i$ ($\alpha = 0.1$)

minibatch (N) pairs ($2N$)

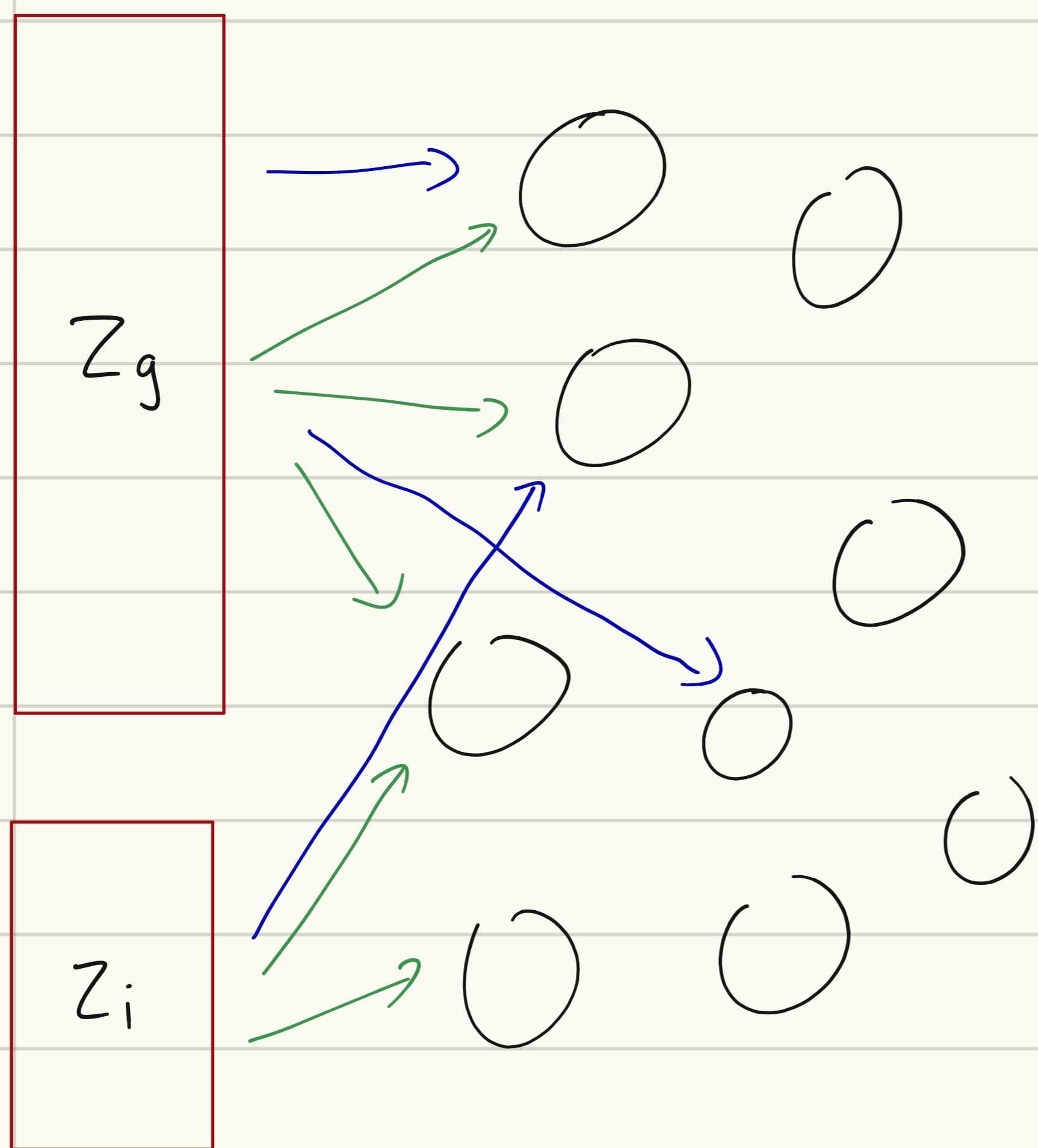
(pretrained)

Shared representation = positive learning
Unshared = negative

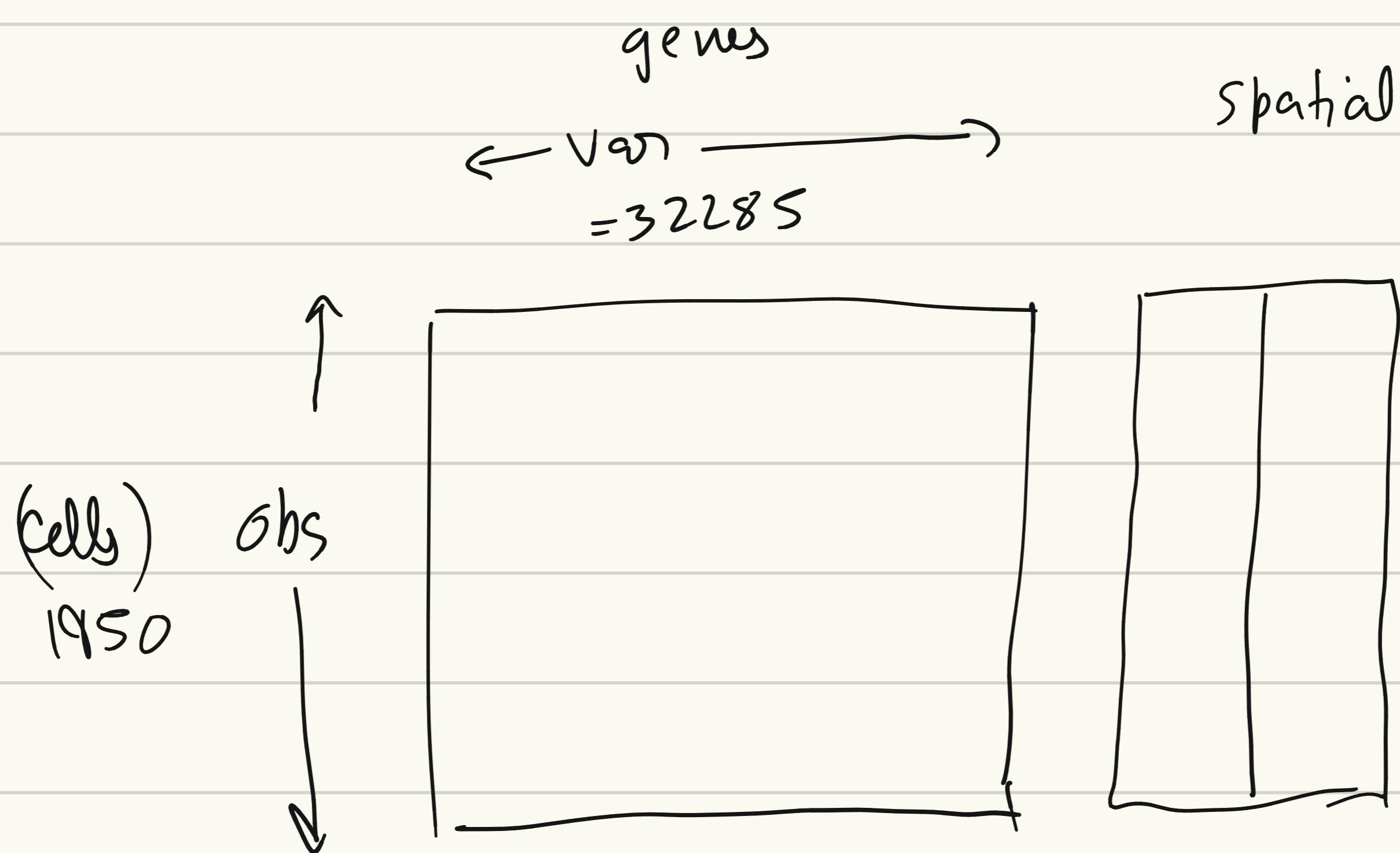
Optimizer = LARS (layerwise adaptive rate scaling)

"Magnitude of update depends on weights norm for better training"
"LARS uses separate learning rate for each layer and not each weight."

Spatial domains



Clustering + refinement
 (mclust, leiden) (against SpaGCN)
 optional



Counts

position

$$x = \text{pos}(x)$$
$$y = \text{pos}(y)$$

labels (dropna)

label_ids(x xy)

$$\text{lbl_x} = \text{lbl_x} + 0.5$$

$$\text{lbl_y} = \text{lbl_y} + 0.5$$

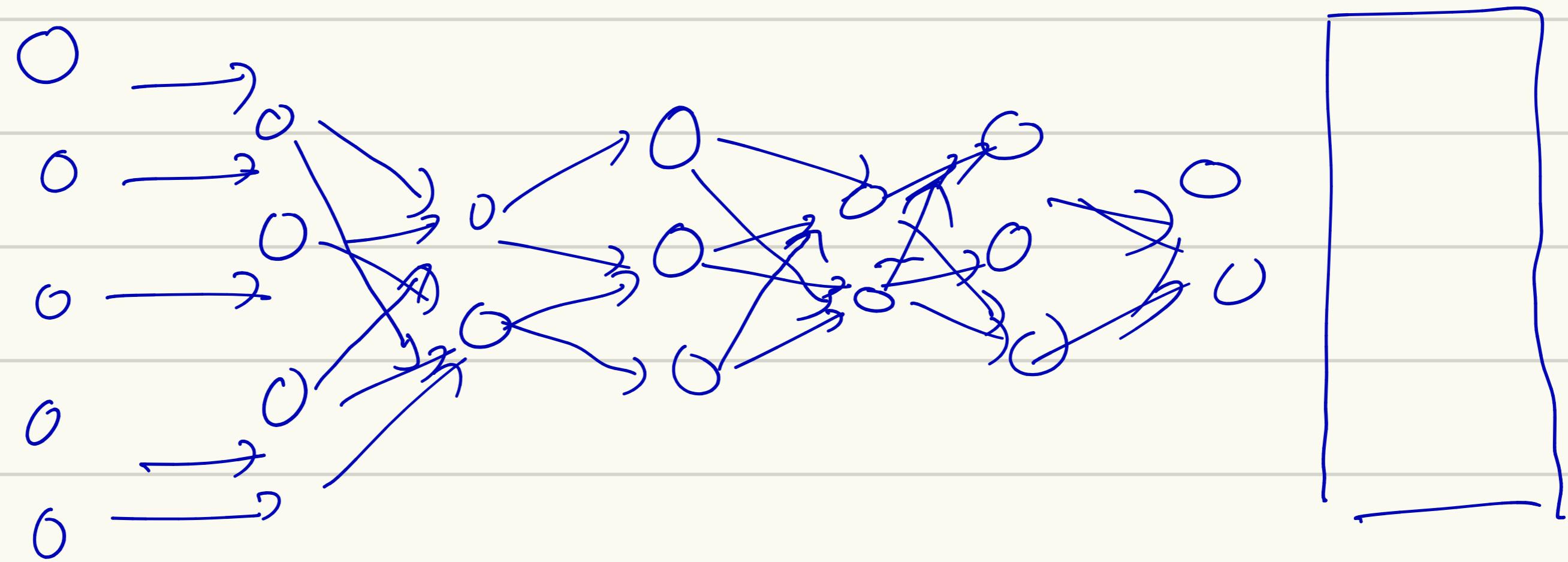
meta_pos \leftarrow joined (Counts + posids)
meta-label

adata.obs[spatial]

adata . obs (label)

GPCR - g-protein coupled receptor.

(Everything is protein)



→ gene expression leads on method.

LIANA +

FASTA - text format for representing nucleotide sequence.

FASTQ -

Sequence + Base quality.

NCBI SRA - database in FASTQ format.

$$Q = -10 \log_{10} p$$

↑
base quality score

probability that base call is incorrect.

Read alignment

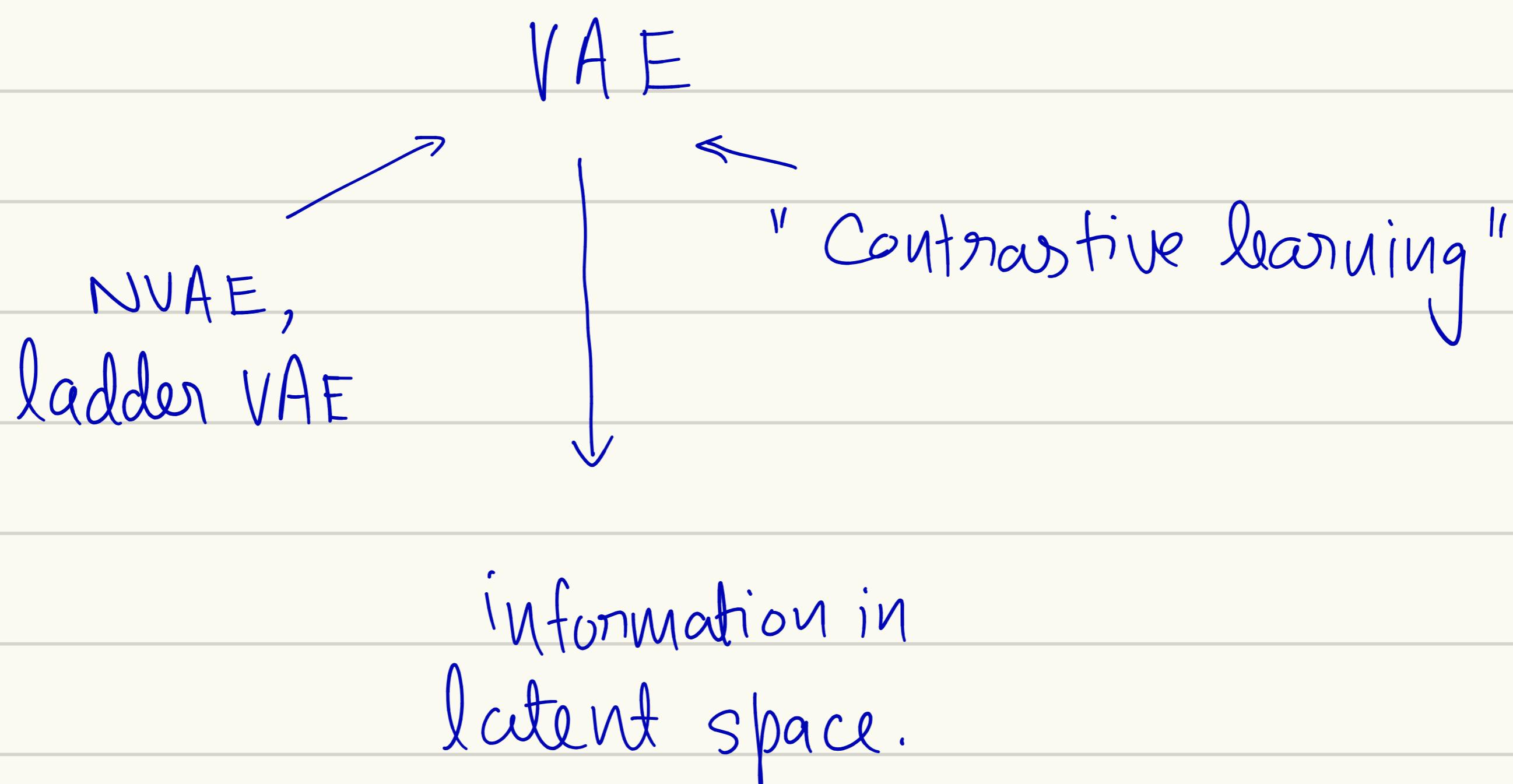
Genome scale index structures.

Genome scale index

i) K-mer

Suffix tree

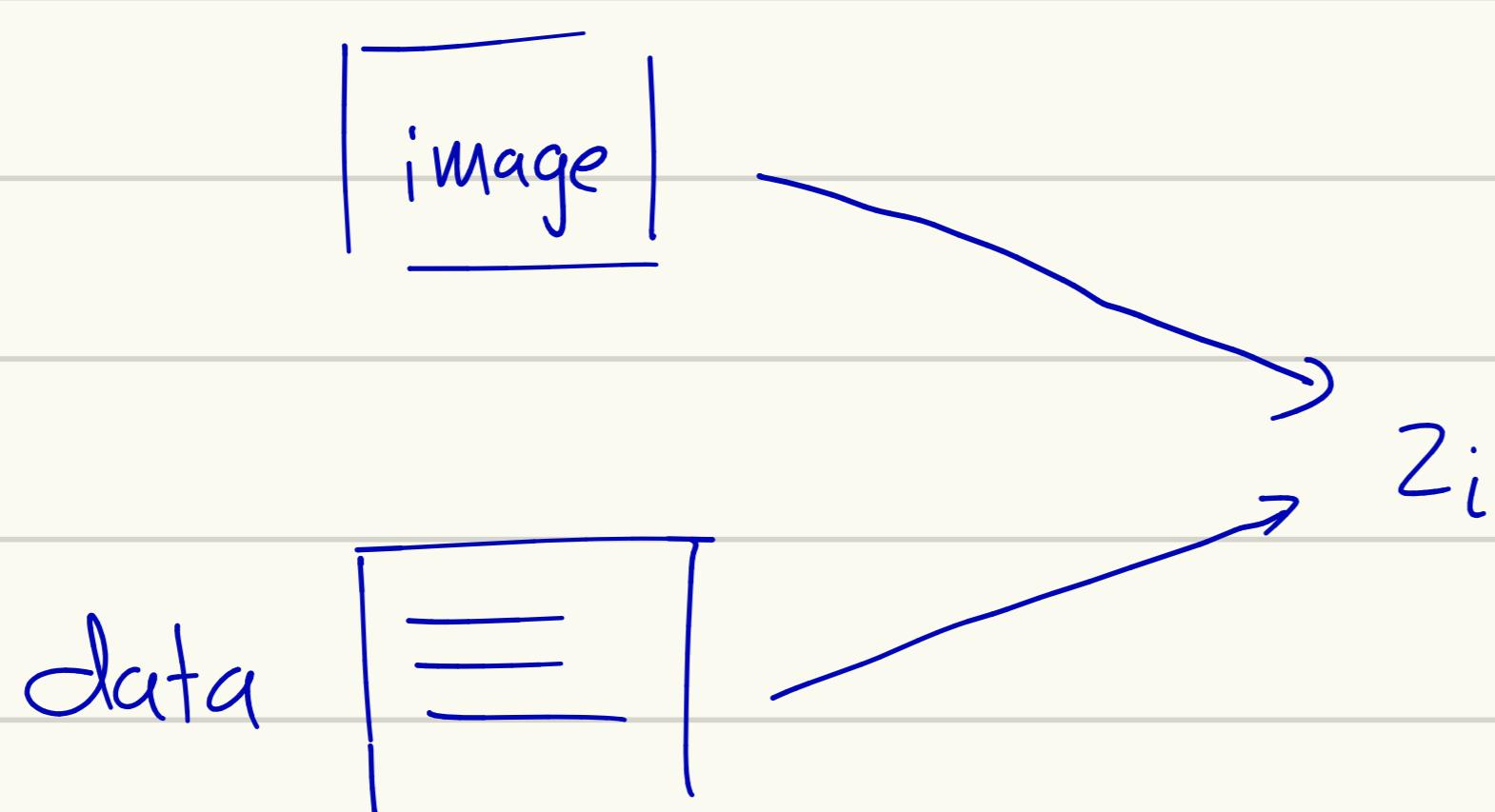
→ Different biological properties learned by different levels of latent variables.



Project 1 :

Mosaic - Using "Uncommon features" to learn latent representation.

Project 3: Multimodal VAE



Project 4: graph neural network.

- do not have very good imputation methods for multiomics data.
- dataset
 - ↑
graph provided ← "missing data" / Smoothness prob.
- MAGIC - "diffusion map approach"

MAGIC
+
GNN → GREP

- Read literature
- Base Model Ready (Comparable to State of Art)
- Better Model

1st week Nov, Review lecture

Use github - Every week discussion
www (friday)

Suffix Array.

"Spatially more efficient and longer search times"

→ Suffix array to store whole genome - 12 Gib

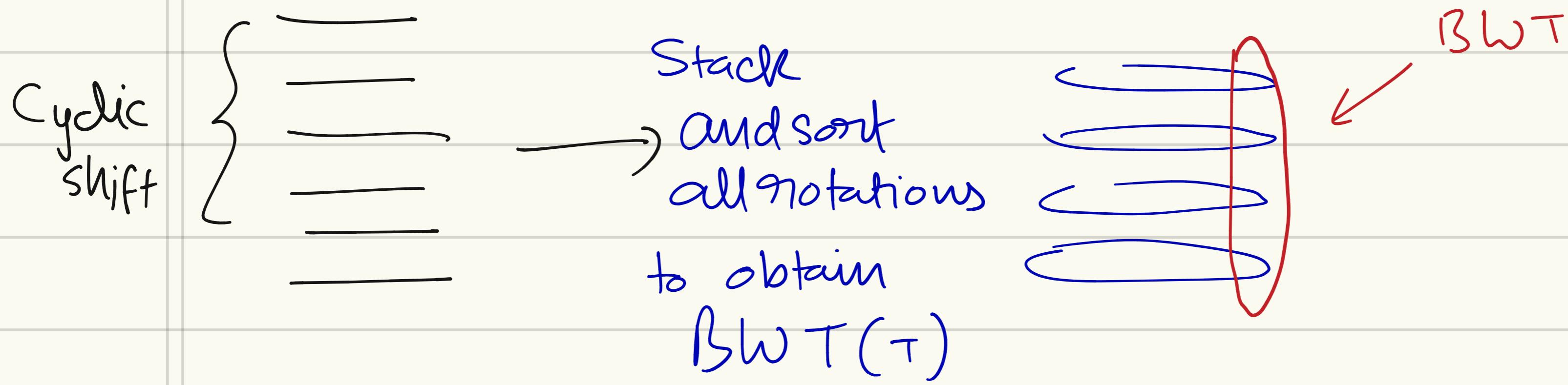
→ (text) $T \rightarrow [0, m]$ $m = \text{len}(T)$

take all suffixes → sort lexicographically.

→ "multi key quick sort"
new paper.

→ space efficient linear time construction of SA

Burrows's Wheeler transform



LF mapping \leftarrow last to first mapping.

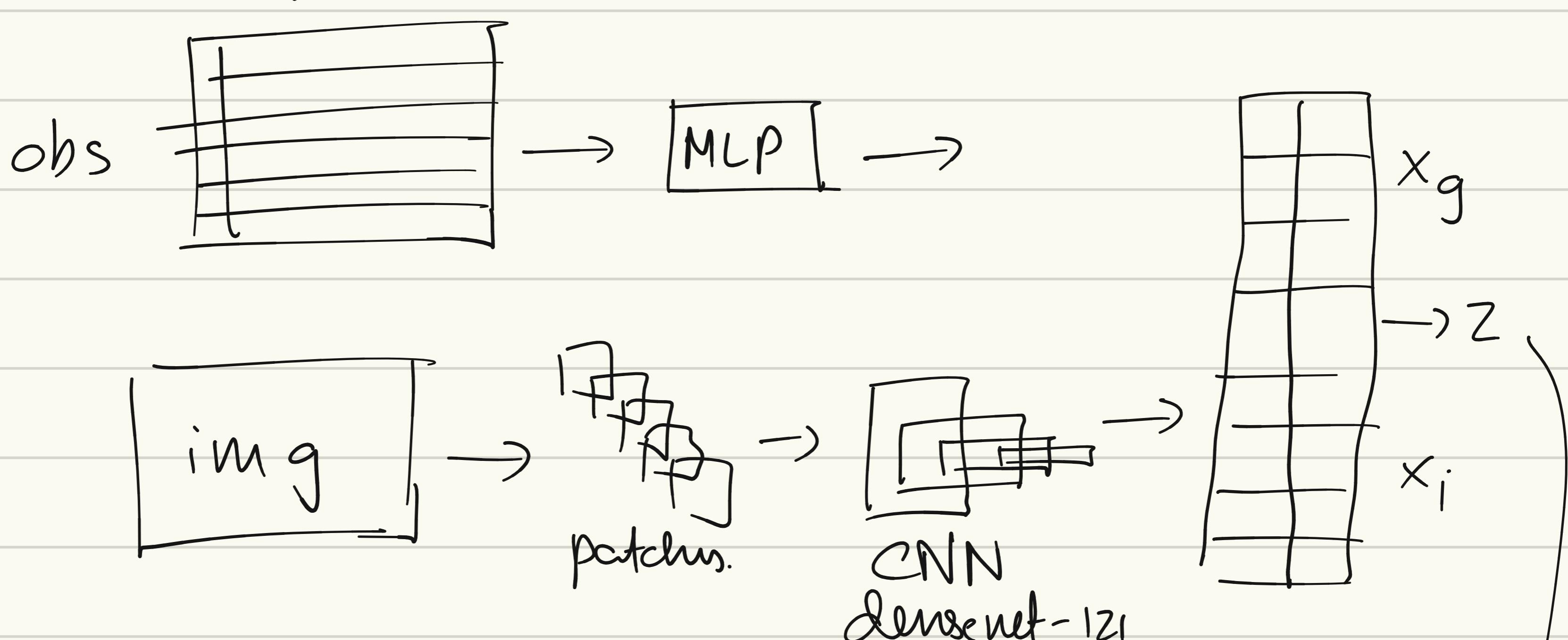
Imputation + Data cleaning.

- BERT based pattern recognition and generation.
- Graph representation learning with IGRAM
- Casually aware imputation
- Vertex and Edges based Friend network reinforcement
 - tracking closely related observations together.
 - optimizing friend network and observations together.

CONGI

- folder containing the data has now been properly passed through library.
- Empty images in patches of spatial ('obs') have been sorted.
Cell positions wrt image more clear.
- Trained on 10 epochs obtained combined embedding structure.

Var.



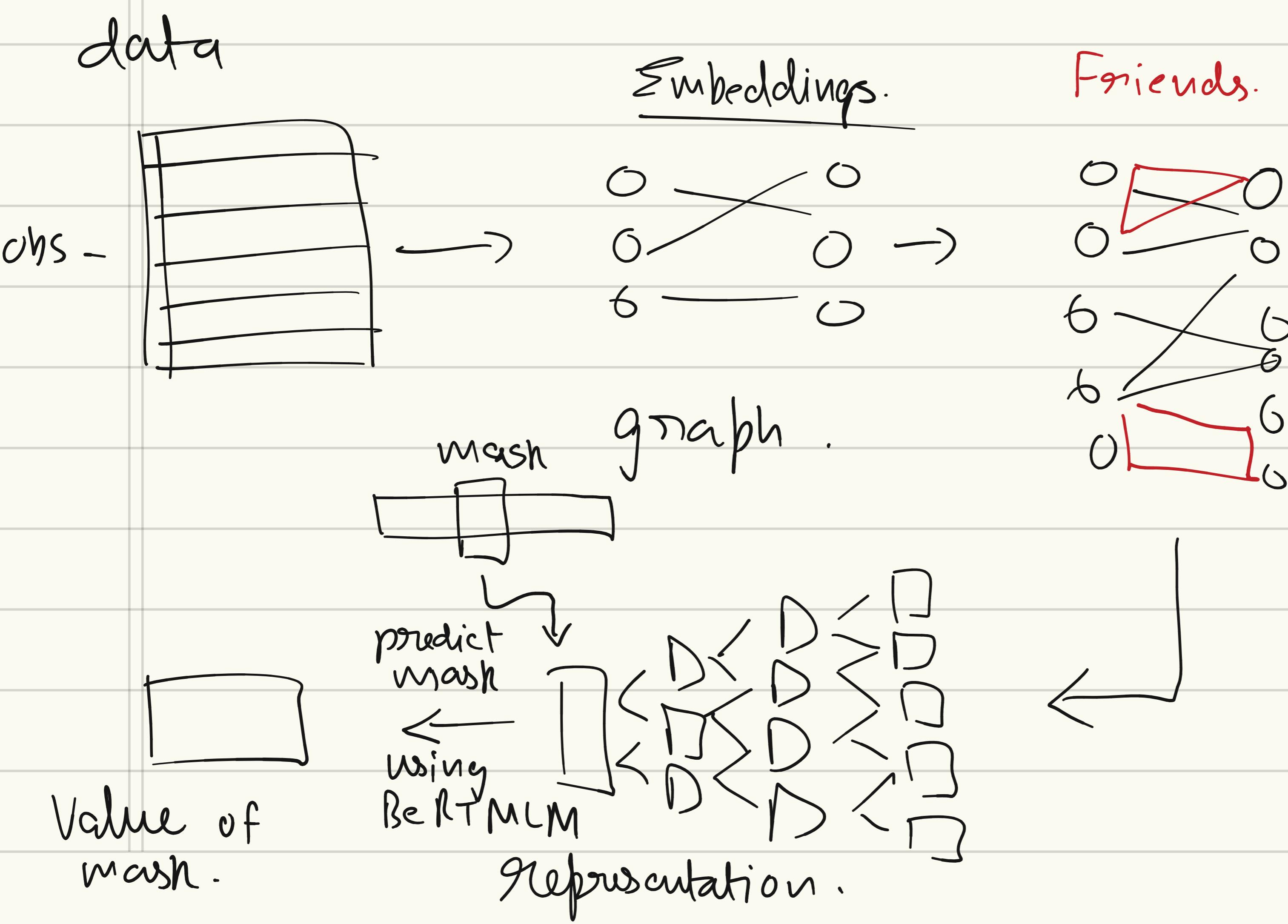
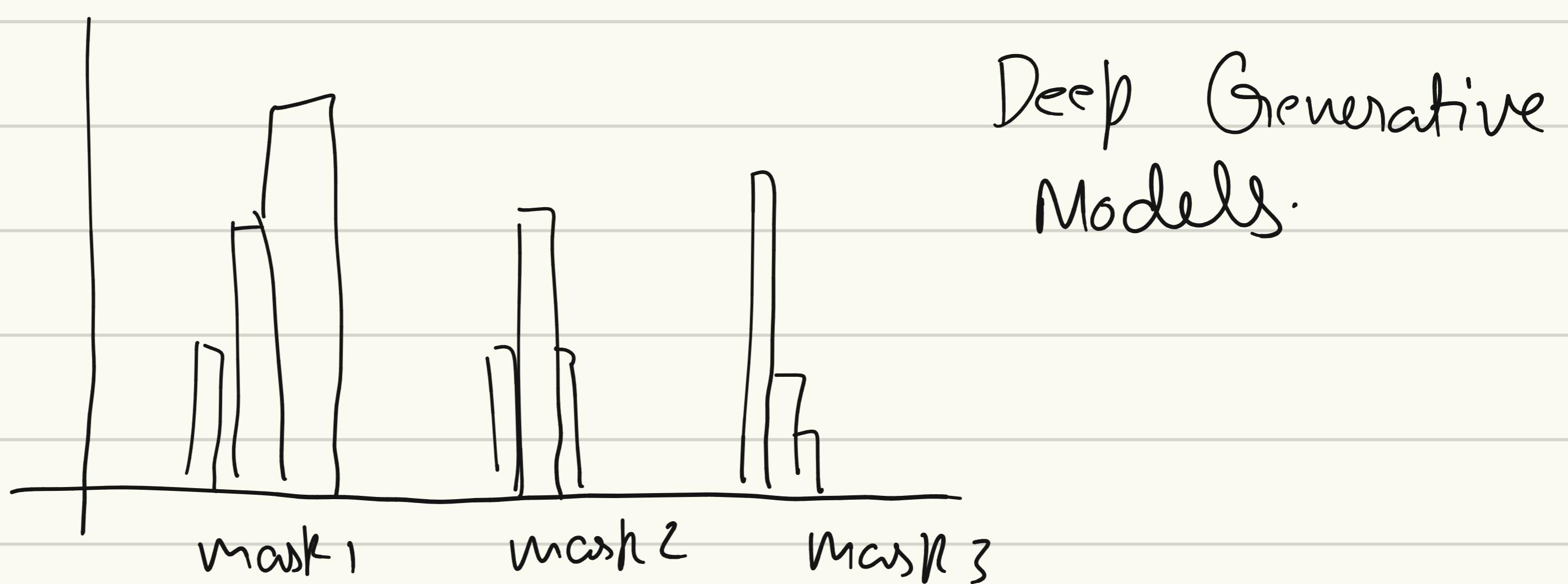
$$z = x_g + \alpha x_i$$

+ training takes some time

10-15 min
for 10 epochs

z - has structure
that can be studied
for plotting.

- Molcst using R-py2 package is unable to load
- thus ani and cluster prediction's are unable to learn, NO PLOTS.



SCBERT

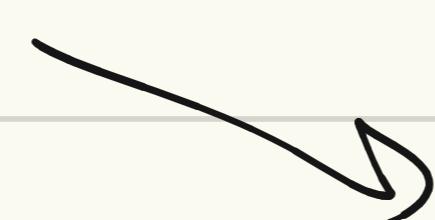
GRAPHE

+

MAGIC

2 Graphs.

Latent

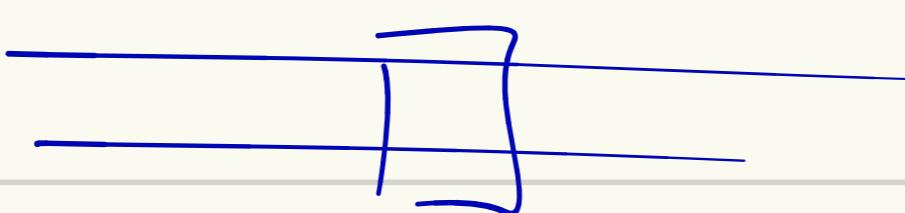


space . - based .

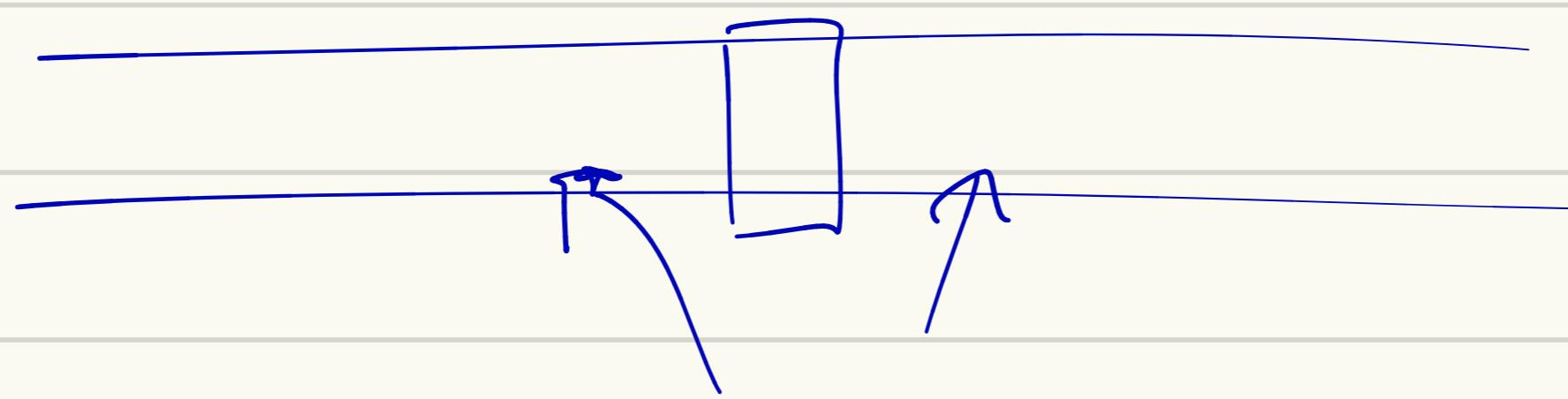
BWT based. Alignment-

Read alignment

Read



↑ match .



Nearly Genome should
be equivalent -

BWT and LF mapping.

- Go To Method Ultrafast human genome alignment
is cite.
- quality aware back tracking
- if there is a multiple candidate subs.
- K-mismatch occurrences of P and T

RNA Alignment

- STAR (suffix map)
- HISAT

RSEM

Gene alignment, Fragment alignment,
Sampling from latent space.

Transcript quantification:-

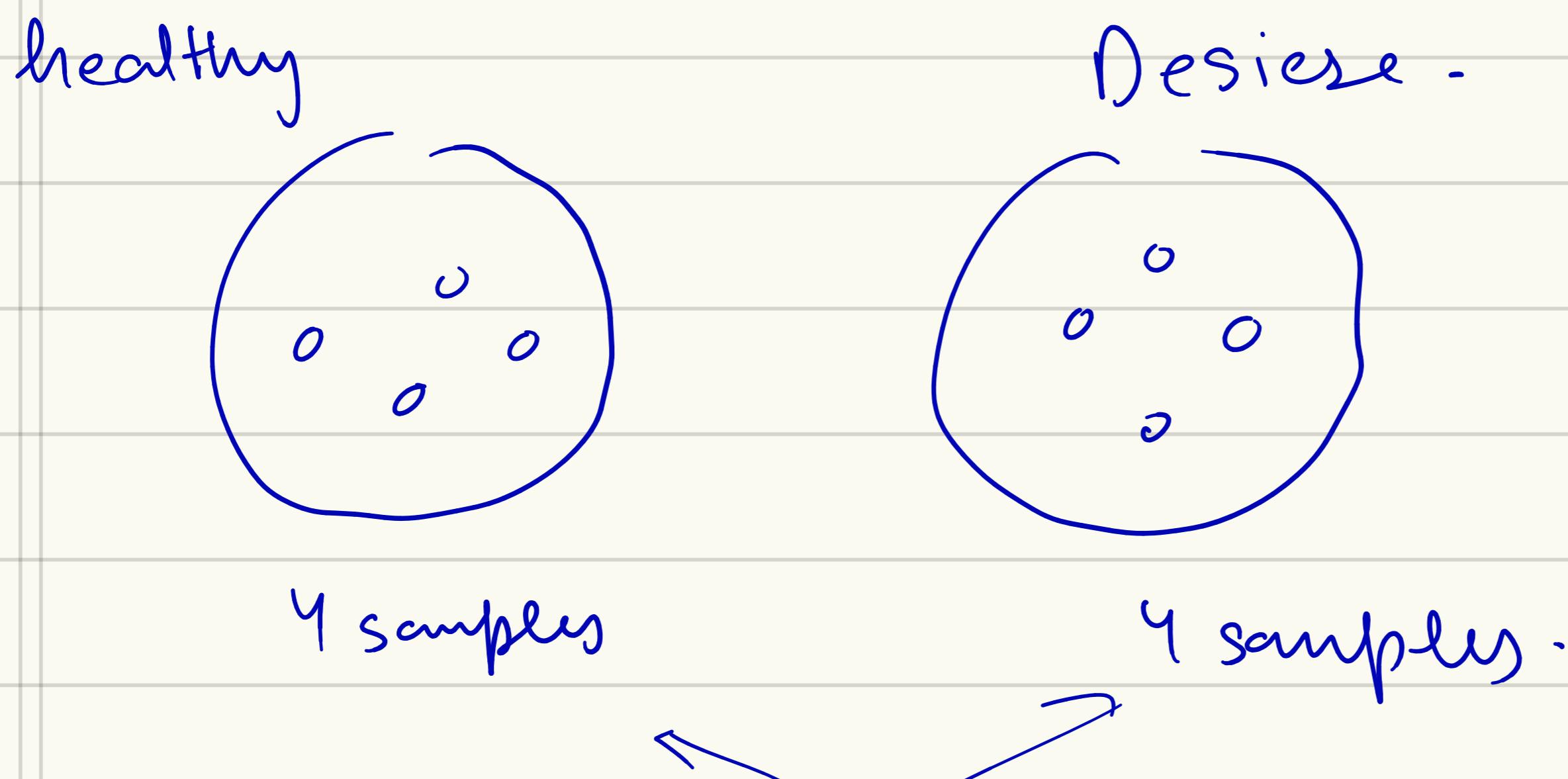


Scribe lecture

3 Nov

Differential Expression.

if we have samples from different conditions.



how do they differ.

Ans: Diff Gene Exp. (DE-Analysis)

Biases:
in RNA

1. Fragment length preference

Some fragments are preferably sampled.

2. Positional bias.

When sampling, selecting. there are preferences in positional.

there's a higher probability of getting good fragments from end of transcript.

3. Sequence based bias.

primers efficiency of binding to molecules differ.

4. Fragment GC bias.

GC - content of organism.

GC content high. - reads less likely.
to getting sequenced.

RSEM - already accounts for fragment-length and positional bias.

- specific variables to account them.

Transcript length:

property vs. prob. of getting sequenced.

→ give length → effective length.

→ soft assignment | probabilities change.

Sum of bias terms
across transcript.

Formulate and test hypothesis.

→ identify genes which behave differently
from one conditions to other.

Two diff. samples.

two diff. genes expression.

ℳ Null hypothesis.

Single cell data - No problem. Each cell gives a lot of samples. (replicates)

Single gene - lookout for variations in other genes also.

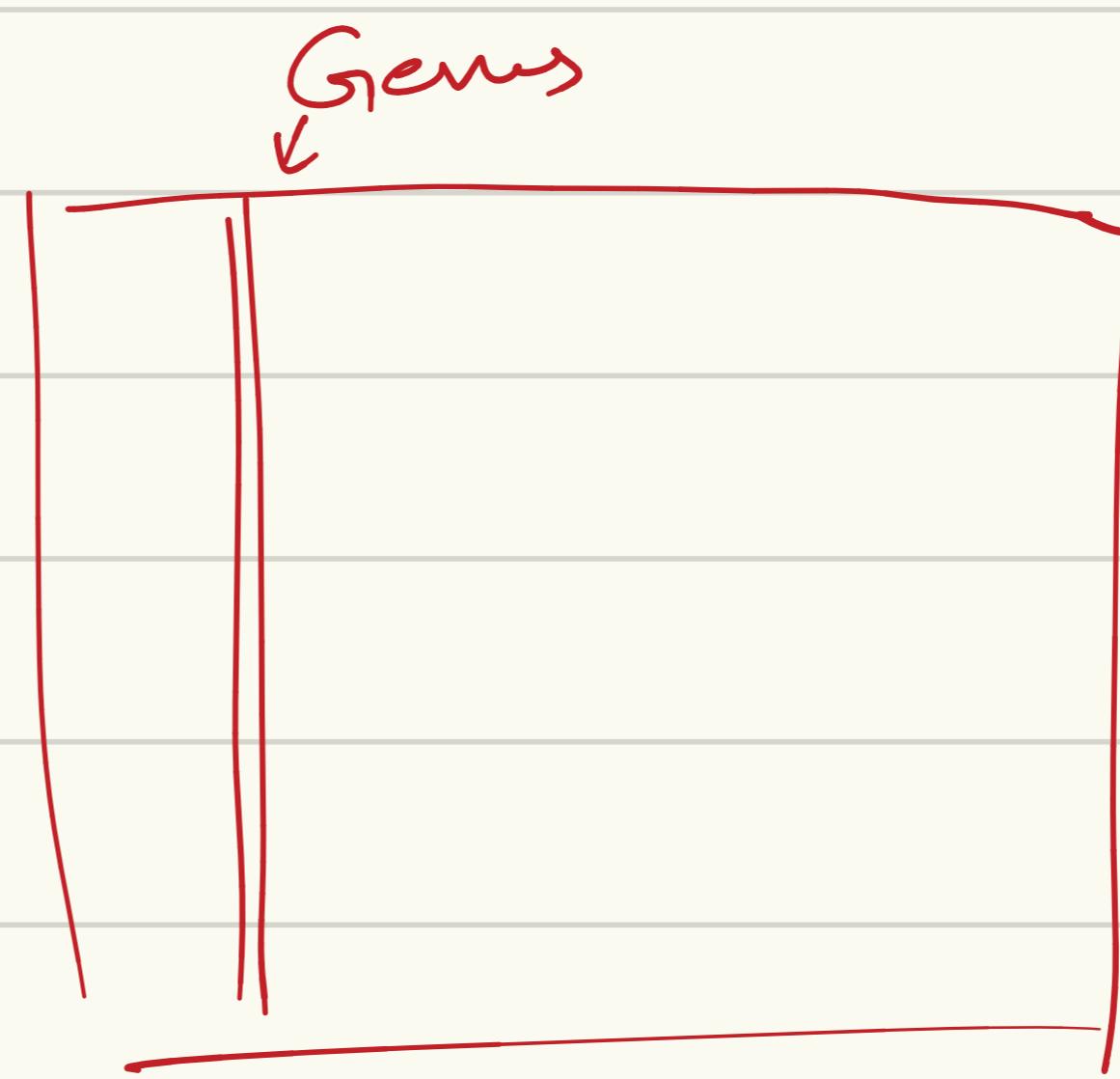
Normalization - (library sized normalization)
for DGE

Bulk seq. data:-

Sequencing depth might be different

Normalization Assumption -

different genes may be assigned.



Rank invariance -

→ there are some genes which will not change across samples.

→ house keeping genes.

Quantile Norm.

no difference across different trials.

Scale Factor Norm.

invariant SET Norm.

Scale Fac Norm.

Same distribution

but different scale factors.



Eg lib size.

Assumption - total minna quant, same
for all samples -

Mean Normalization -

results
in
comparable
datasets.

$$\bar{y}'_j = \bar{y}_j - M_j + M$$

Variance Normalization :-

$$\hat{y}_j^i = \frac{(y_j^i - M_j) \sqrt{V}}{\sqrt{V_j}}$$

$T_i \leftarrow$ eff. no. of trans. per gene (i')

$$\uparrow \\ (x10^6) = TPM \leftarrow$$

Not the only way

RPKM \leftarrow Comparability b/w diff genes.

$$\uparrow \\ RPKM(G_i) = \frac{R}{L} \frac{10^6 \cdot 10^3}{M}$$

15 year old

\downarrow
FPKM (Fragments)

Count No. of fragments that are coming directly from the experimental data.

$$FPKM(G_i) = \frac{C}{L} \frac{10^6 \cdot 10^3}{M}$$

{Set of genes} \leftarrow highly expressed in one experimental conditions.

library size same \hookrightarrow Sample A (500 genes)

\curvearrowright Sample B



build statistical model.

problem.

\rightarrow Exp. conditions.

Gene

A

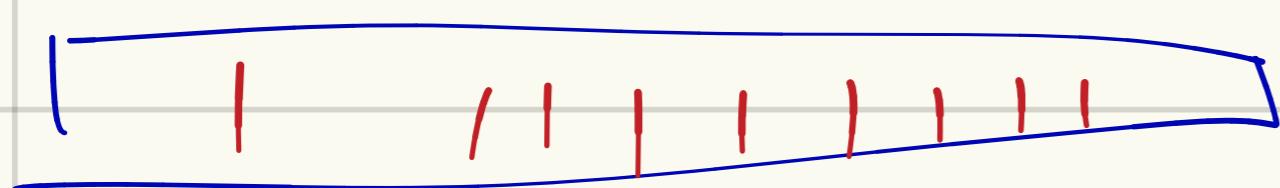
B

20 Mill

20 mill;

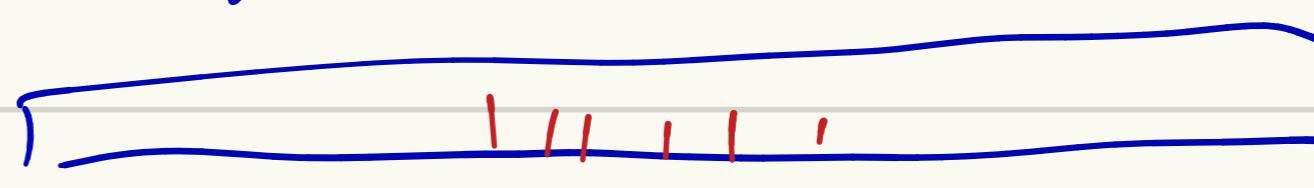
50 genes. \leftarrow same \rightarrow 50 genes.

50 genes \leftarrow diff \rightarrow 0 genes.



100 reads

spread.



spread
across
50 gene

TRIM mean of M (TMM)

$$E[Y_{gK}] = \frac{M_{gK} L_g N_K}{S_K}$$

$$M_g = \log_2 \frac{Y_{gK}/N_K}{J_{gK}/N_{K'}}$$

II II II

Trim of high variability genes (A_g)
and outlier mean value.
then normalise on the
basis of the remaining genes.

Spiking ← molecules whose exact sequence
concentration is known.
(calibration uses)

Differential Expression

for different genes, when the expression levels are statistically significantly different.

→ Statistical tests: performed to determine DE.

Example: If we have ~~the~~ genome data from cells under different conditions (Eg. disease / heat stress)
↳ "How would we, by looking at the expression data, find out what genes are expressing differently?"

→ DE genes could tell us about cells ~~in~~ samples as well. (Eg - classifying diseased cells).

→ DE genes can help us understand the mechanisms behind certain biological processes.

→ DE gene expression could be because of ARTIFACTS
- experimental / technological.

→ Calculate p-values to see if expression data distribution could have been altered by "accident".

↳ Example of A vs B

→ Simple statistical tests are useful when you have a large number of replicates of the dataset.

→ Also, ~~the~~ experiments are expensive, and we have 2-3 replicates.

→ Checking spot size, we want to find out if any variation in read counts is "natural" or "by chance".

Candidates for Read Counts:

→ Poisson: Mean = Variance = ~~1~~.

↳ Underestimates variations

→ Overdispersion (Real data: Variance > Mean)

Note: In large samples the count value might be 100, in some others it might be 1.

Negative binomial:

→ How many tosses before 2 ~~successes~~ heads?

→ Better fit for datasets with greater mean expression levels.

$$P(X=k) = {}^{n-1}C_{k-1} (1-p)^k p^r$$

→ Mean and Variance can be modeled using a dispersion parameter.

$$\text{Let } \mu \text{ } \boxed{\mu + \alpha \mu^2}$$

→ PDE seq: $p(j) \rightarrow$ Conditions
 $j \rightarrow$ gene

+ one concentration of fragments from

gene i is an unknown (latent) variable

$\mu_{ij} \rightarrow$ shot wise

$$\left[s_j^2 r_{ij} p(j) \right] \rightarrow \text{Raw variance -}$$

→ pools information from
genes with similar expression levels.

size factors

Parameters

Expression strength
parameters.

Dependence of raw variance on expected
mean expression,

Simple Assumption:

If gene i is not differentially expressed -
→ the count of reads from i should be proportional
only to the size factors.

⇒ A few DE genes might have an outsized effect -

Common scale: denominator is the estimated
size factor.

⇒ $p \rightarrow$ p_i calculate sample variations,
normalized by size factors.

⇒ for a small number of replicates, we use a
local regression for the smooth dependence
function.

NPBT (Nb DE genes)

→ fb: gene expression is absent.

p-value → Probability of our data being generated under the null hypothesis.

→ Depends ~~strongly~~ on our underlying model / distribution.

In our case, the test statistic is the count in each condition.

Log · fold change?

Multiplet test p-value problem:

↳

DESeq2

EdgeR

DEMMIA

Volcano Plot Explanation:

MA plot

↳ Log fold change vs Mean Expression

→

Pathway Analysis Methods:

↳ GSEA

DAVID

CyberGSE

→ good for many areas

