

block b6

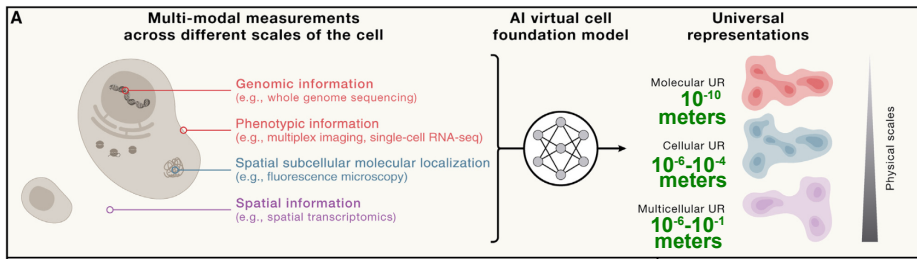
The AI Virtual Cell

block b6  
The AI Virtual Cell

“A conceptual feat at present”

*Lin Tang, Nat Meth Editor, Dec 2025.*

# Multiple scales to unify Multiple data types

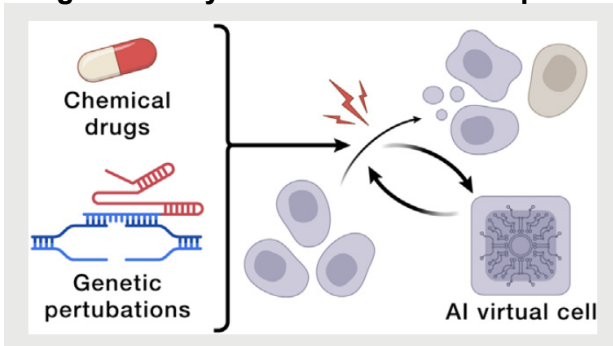


“How to build the virtual cell with AI: priorities and opportunities”, *Bunne et al, Cell Dec2024*

## Universal Representations (URs) Universal Embeddings

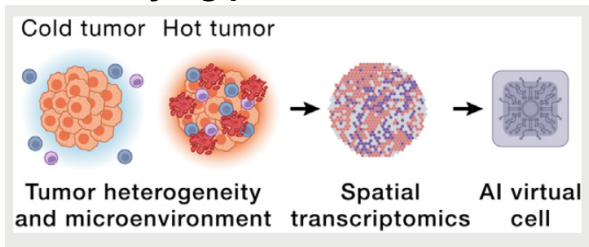
# Desired Capabilities

drug-discovery and cell-based Therapeutics



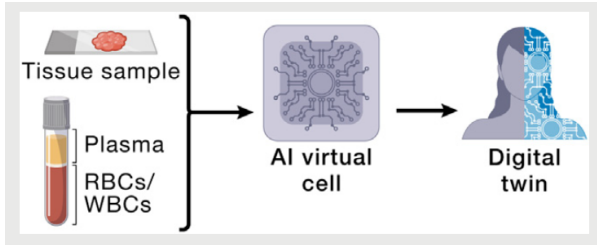
# Desired Capabilities

## identifying pan-cancer markers



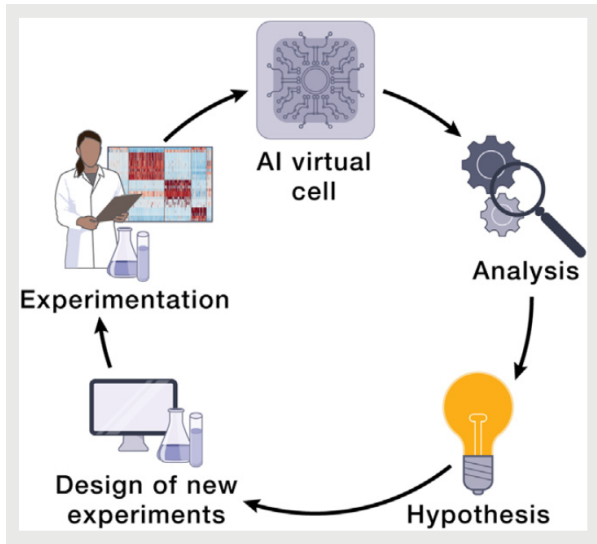
# Desired Capabilities

cell models for individual patients

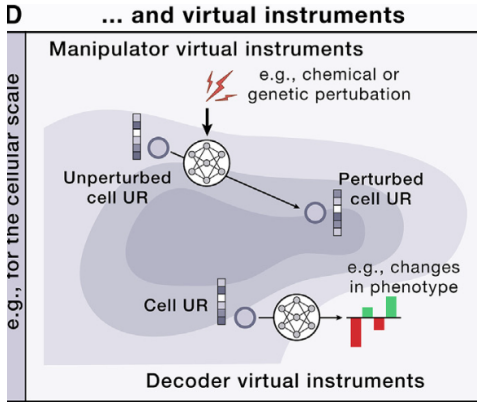


# Desired Capabilities

## HYPOTHESIS GENERATING VA



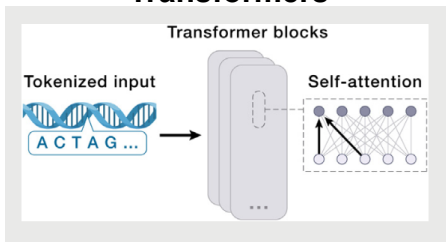
# Desired Capabilities



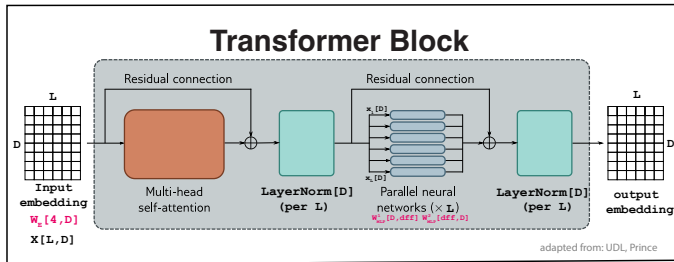


# Methods

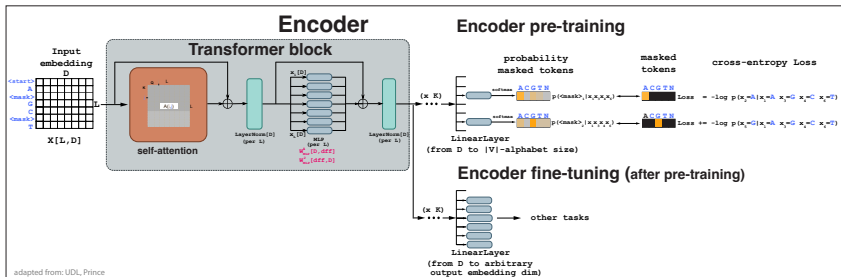
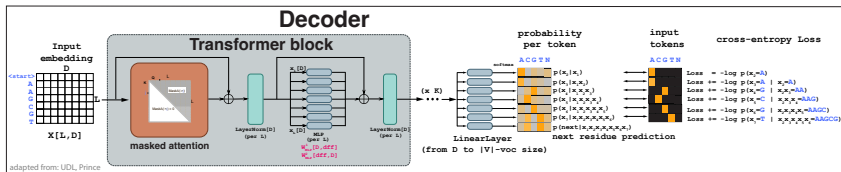
## Transformers



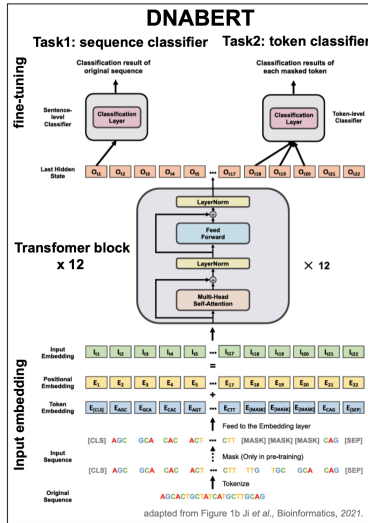
# Methods



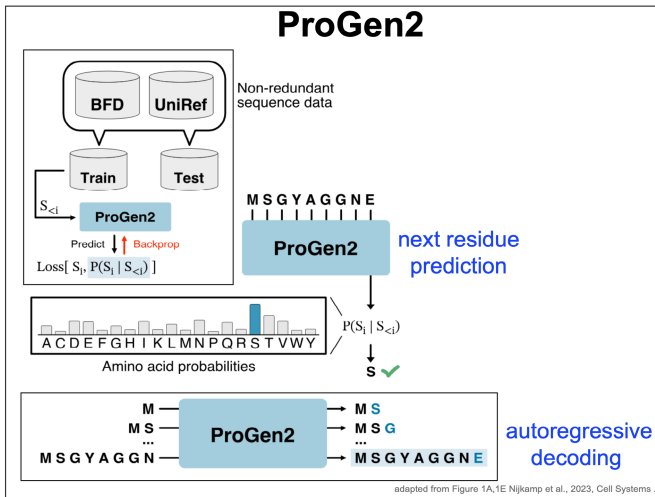
# Methods



# Methods

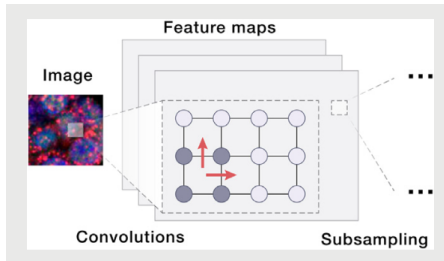


# Methods



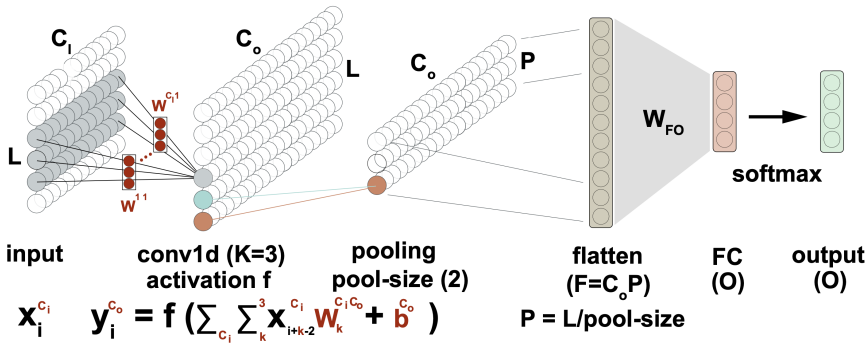
# Methods

## CNNs

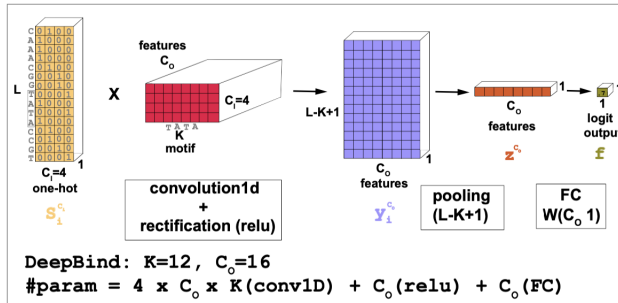
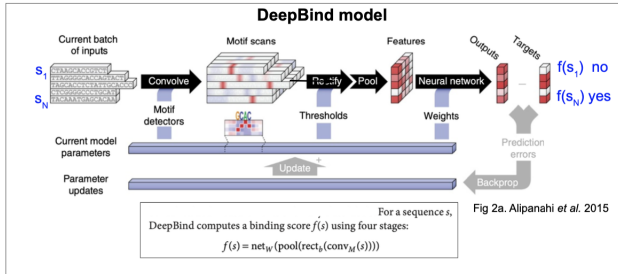


# Methods

## Convolutional Network 1D

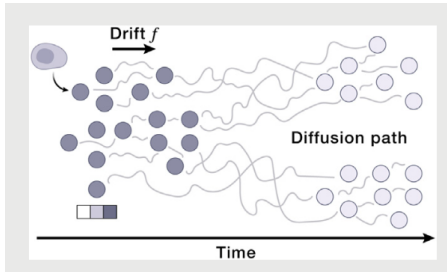


# Methods

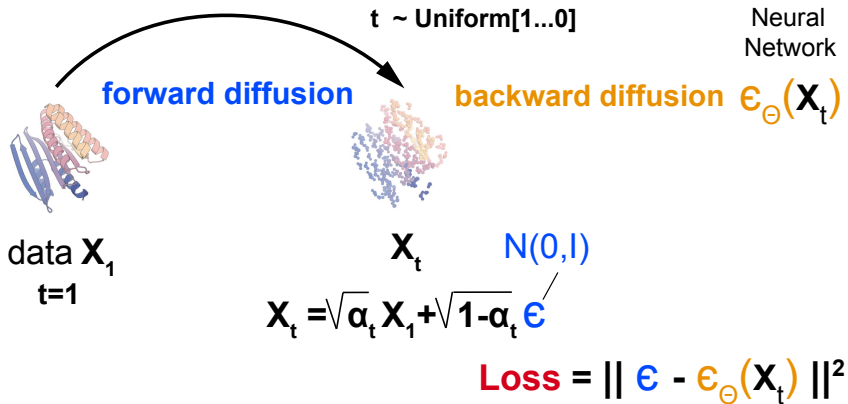


# Methods

## Denoising Diffusion Probabilistic models

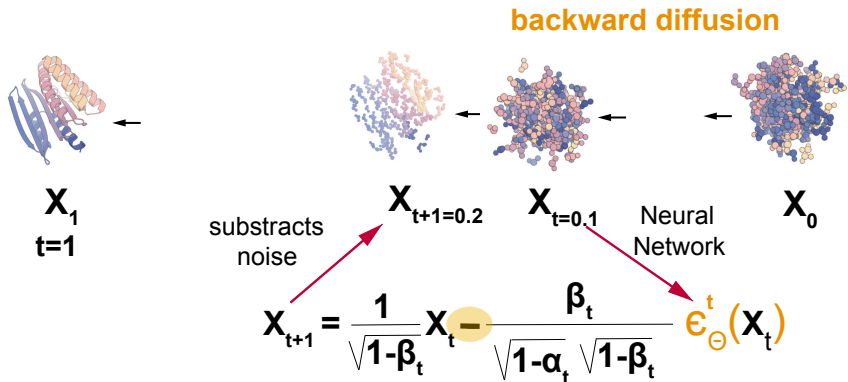


# DIFFUSION Training



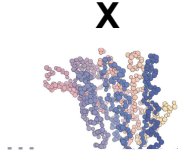
# Methods

## DIFFUSION Sampling



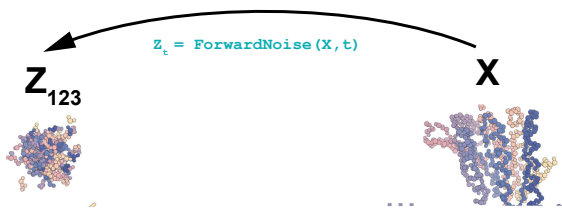
$\hat{X}_0$  (prediction)

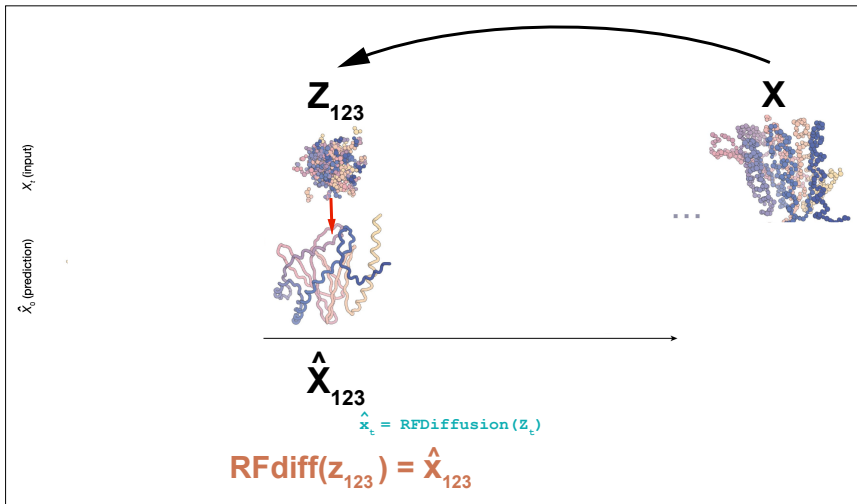
$X_t$  (input)

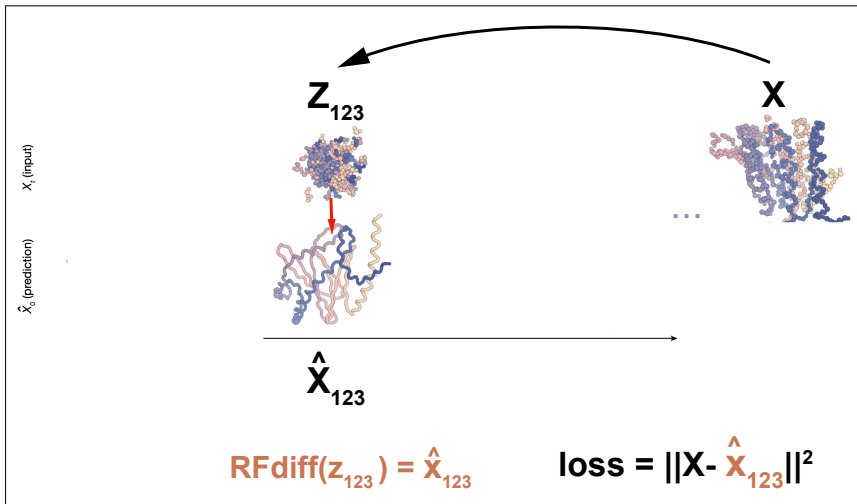


$X_t$  (input)

$\hat{X}_0$  (prediction)







$Z_{200}$

$X_t$  (input)



$\hat{X}_0$  (prediction)

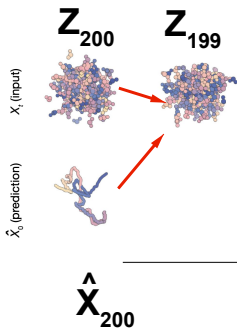


$\hat{X}_{200}$

$\hat{x}_t = \text{RFDiffusion}(Z_t)$

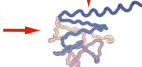
$$\text{RFDiff}(z_{200}) = \hat{X}_{2000}$$

$$z_{t-1} = \text{ReverseStep}(z_t, \hat{x}_t)$$



$$\text{RFdiff}(z_{200}) = \hat{x}_{200}$$

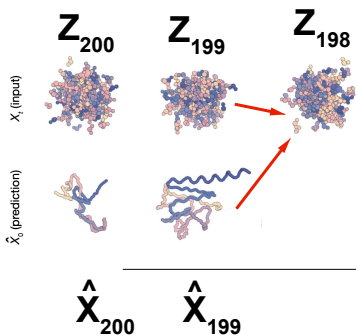
$$z_{199} = \text{interpol}(z_{200}, \hat{x}_{200})$$

$Z_{200}$  $Z_{199}$  $X_t$  (input) $\hat{X}_t$  (prediction) $\hat{X}_{200}$ 

$$\hat{x}_t = \text{RFDiffusion}(Z_t, \hat{x}_{t+1})$$

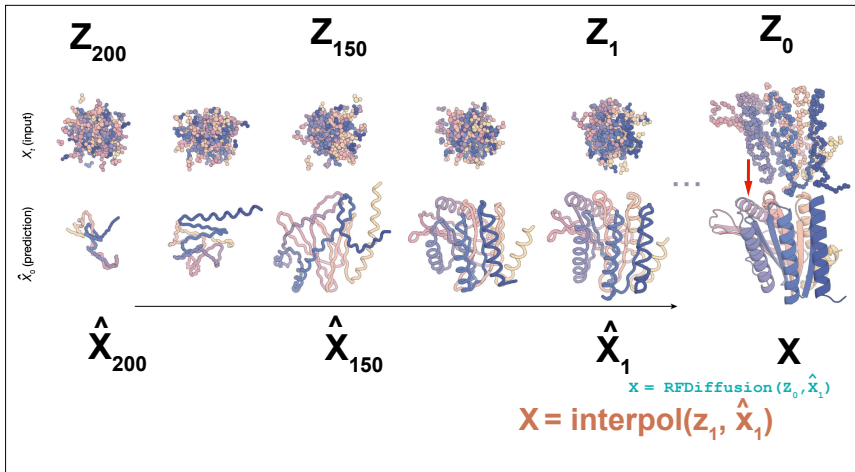
$$\text{RFdiff}(z_{199}, \hat{x}_{200}) = \hat{x}_{199}$$

$$z_{t-1} = \text{ReverseStep}(z_t, \hat{x}_t)$$



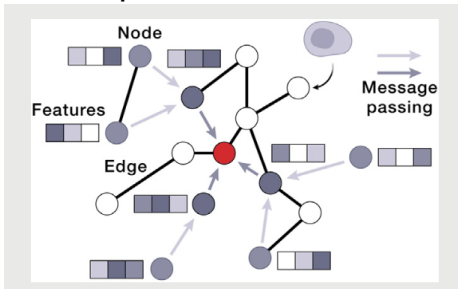
$$\text{RFdiff}(z_{199}, \hat{x}_{200}) = \hat{x}_{199}$$

$$z_{198} = \text{interpol}(z_{199}, \hat{x}_{199})$$



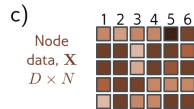
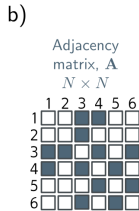
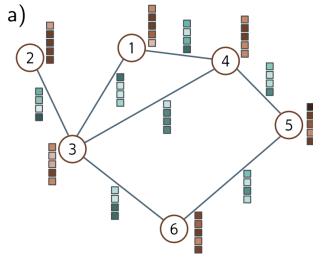
# Methods

## Graphical Neural Networks



# Methods

## Graphical Neural Networks

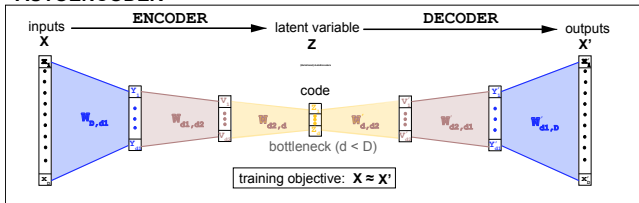


gene/protein networks  
drug/drug interaction networks  
cell/cell interaction networks

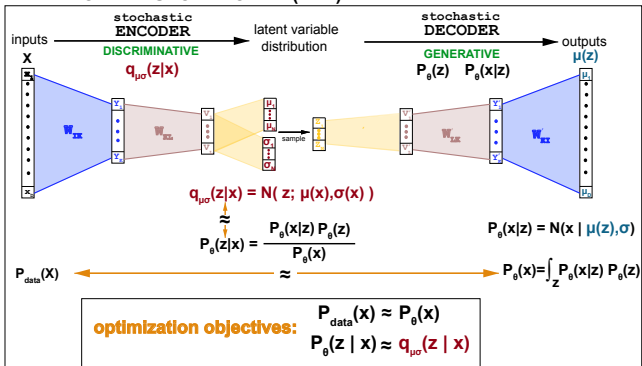
**Graphical Convolutional Networks**

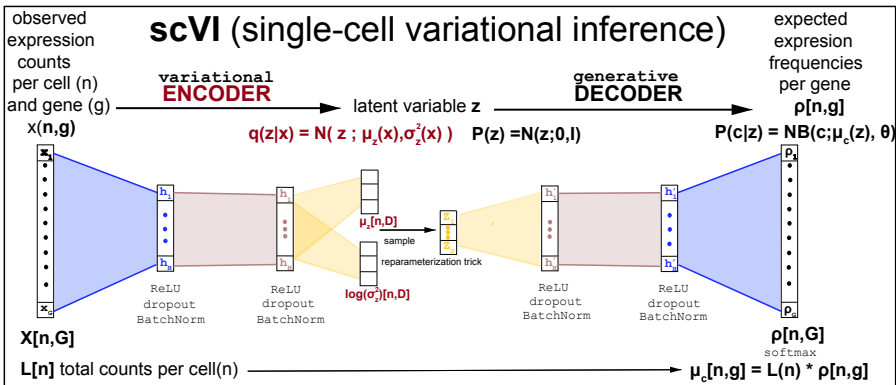
# Methods

## AUTOENCODER



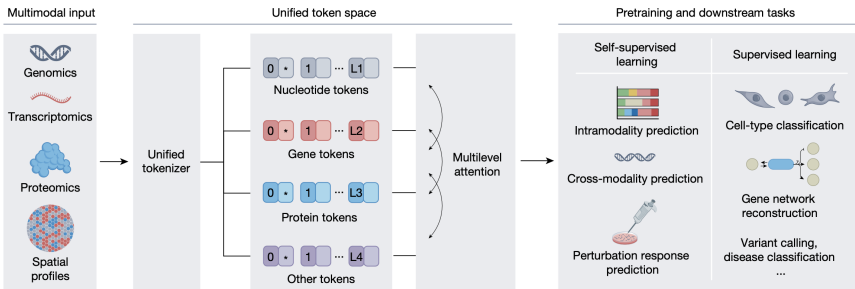
## VARIATIONAL AUTOENCODER (VAE)





# Multimodal Foundational Models (MFMs)

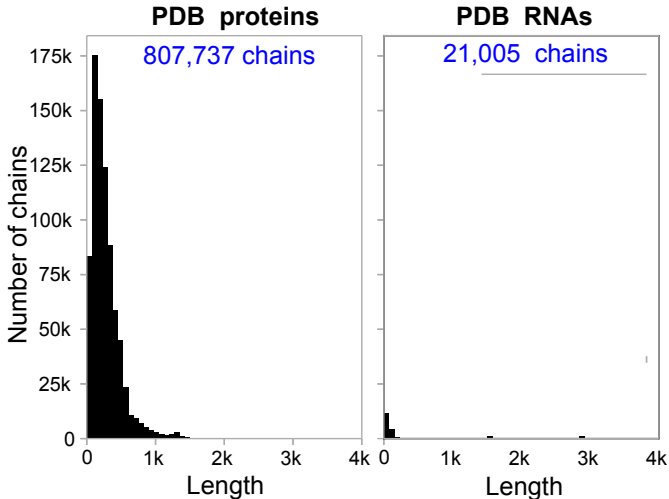
## a Computational components for multi-modal foundation models



“Towards multimodal foundation models in molecular cell biology”, *Cui et al, Nat Perspectives, 2023-25*

# Challenges - Always more data

## Protein Data Bank (PDB) 2024



# Challenges - Always more data

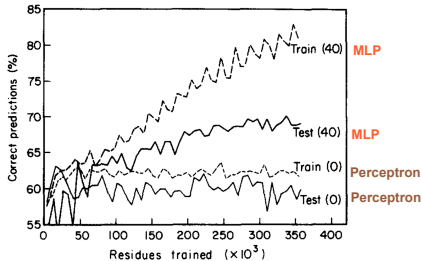
“Calling all data”, Nath Meth Editorial, July 2025.

- ▶ Data generation w/o hypothesis hard to fund
- ▶ Data generation w/o hypothesis hard to publish
- ▶ Data storage costs money
- ▶ New experimental technologies to explore more of the data space needed

**Challenges - Always more  
computational resources**

# Challenges - Always less data leakage:

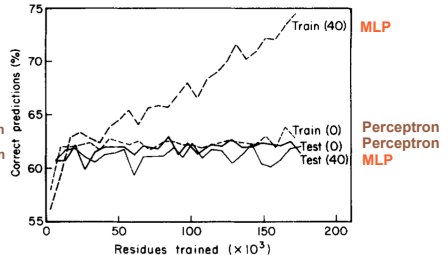
## Rigorous evaluation



Qian-Sejnowski Figure 14

Train set = Globin proteins  
Test set = homologous  
Globin proteins

train/test overlap



Qian-Sejnowski Figure 8

Train set = Globin proteins  
Test set = non-homologous  
Globin proteins

train/test \*less\* overlap

# Challenges - Always less hallucination:

To generate samples based on data

$$P_{\theta} \approx P_{data}$$

then

$$[P_{\theta} \sim] x_{pred} \approx x_{pdata} [\sim P_{data}]$$

# Challenges - Always more ethical

1. Models should be open source
2. Data should be publicly accessible
3. Data privacy

**Challenges - Always more Bayesian**

**Challenges - WHAT TO LEARN  
(teach)?**

Thank You

