# Manifold Learning and Artificial Intelligence Lecture 8

## A Path to Protein Design and Disease Mechanism (1)
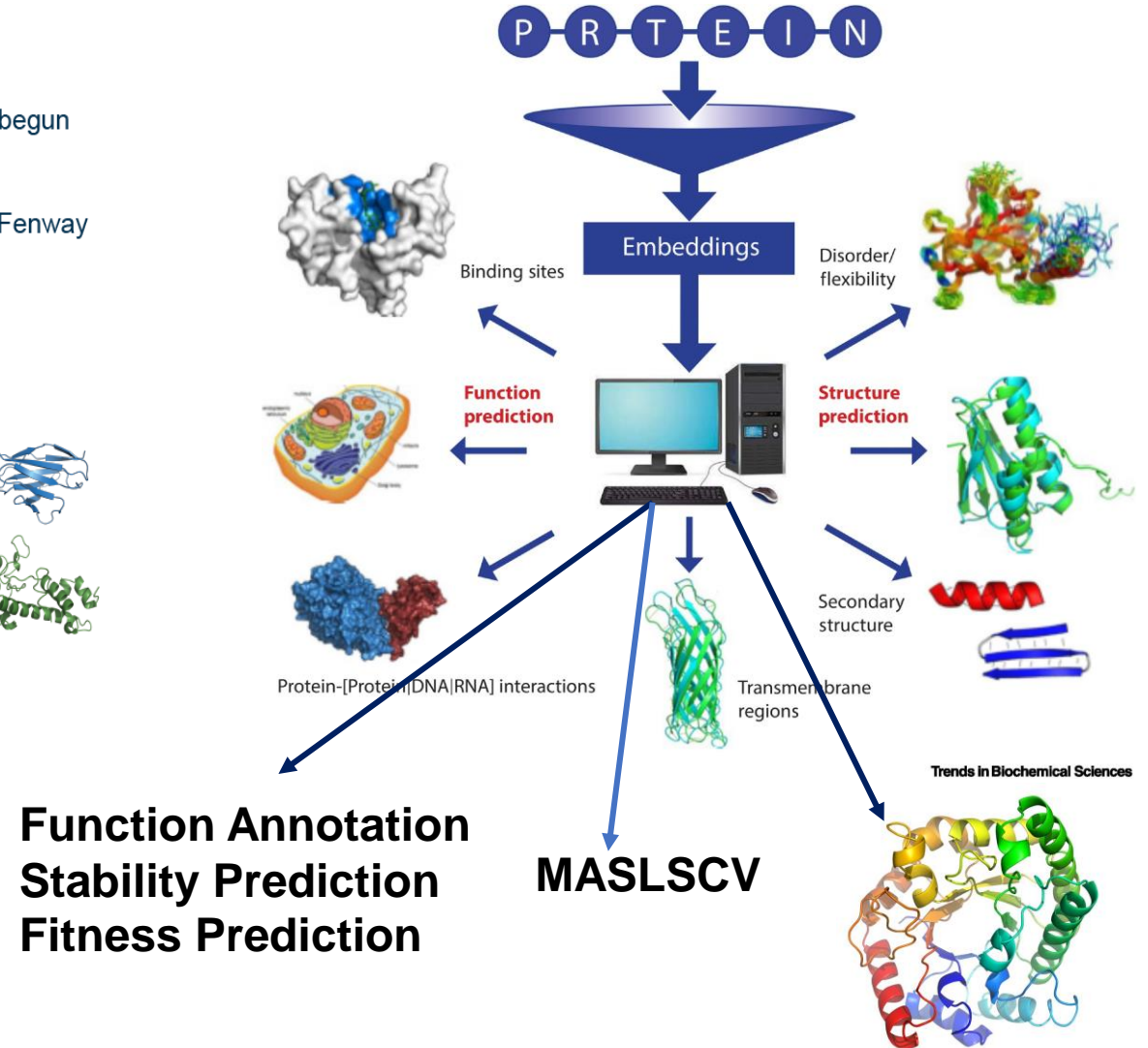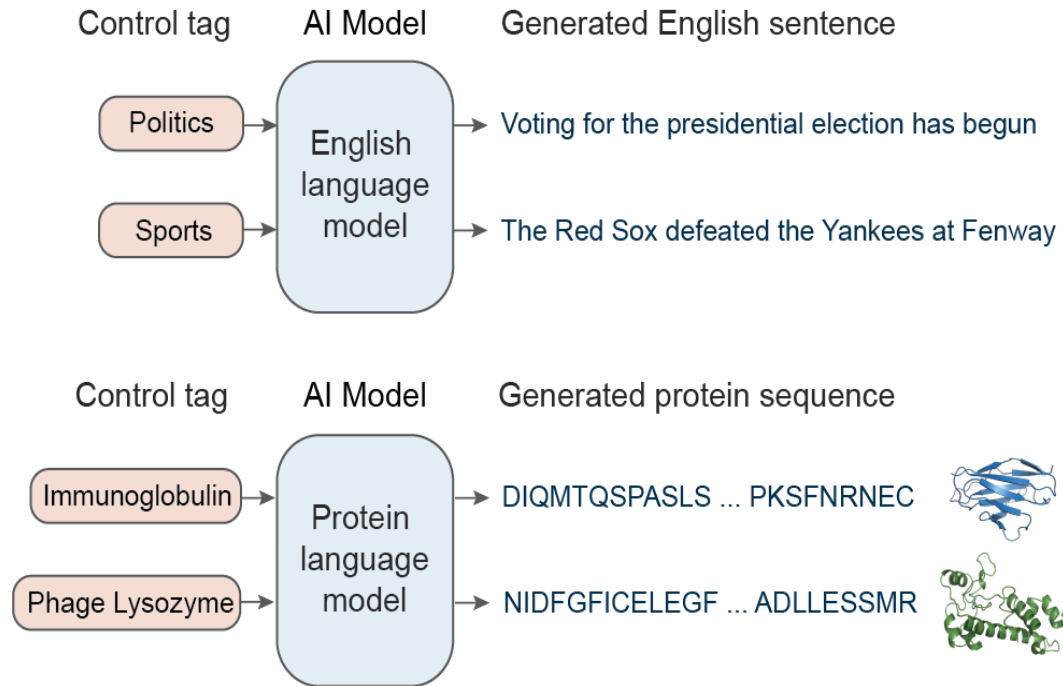## Protein Language Model

Momiao Xiong, University of Texas School of Public Health

- Time: 9:00 pm, US East Time, 01/14/2023
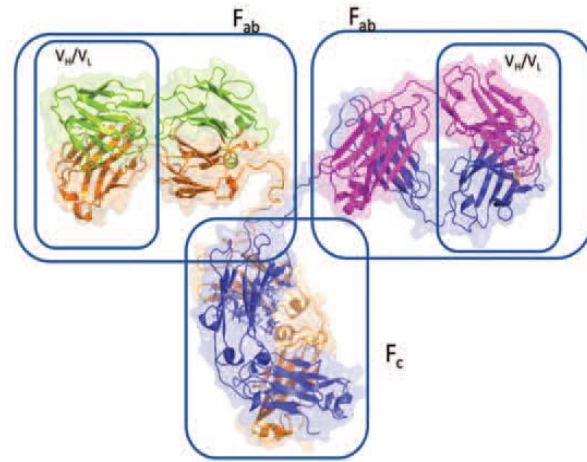- 10:00 am, Beijing Time. 01/15/2023
- Zoom

https://uwmadison.zoom.us/w/93316139423?tk=wfbmsTfN2fgERto_HI1WKtBzh94d3HO02XV
Cexqd8.DQMAAAAVuhNJnxZ2dGVCbllIYlR3ZWJBNHl5LXlYMjBnAAAAAAAAAAAAAAAAAAA
AAAAAA&pwd=Q0NVWFYvRFg5RmxCNkwxMmYrbW41dz09
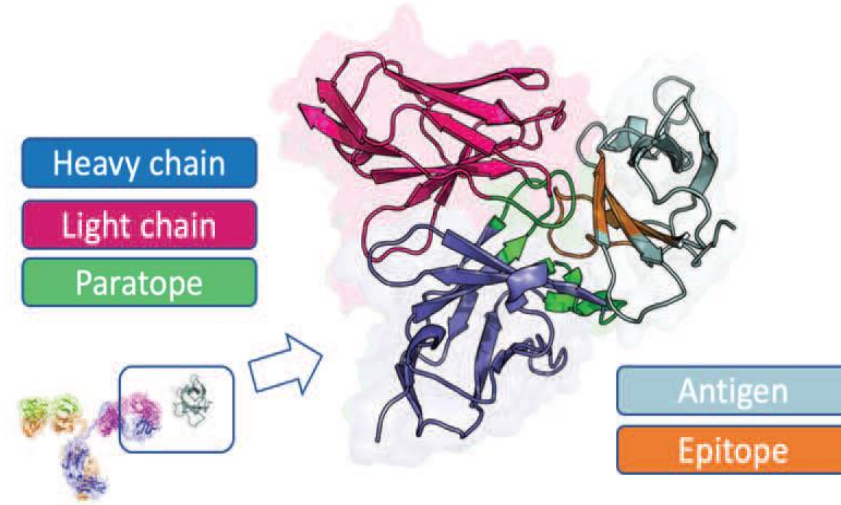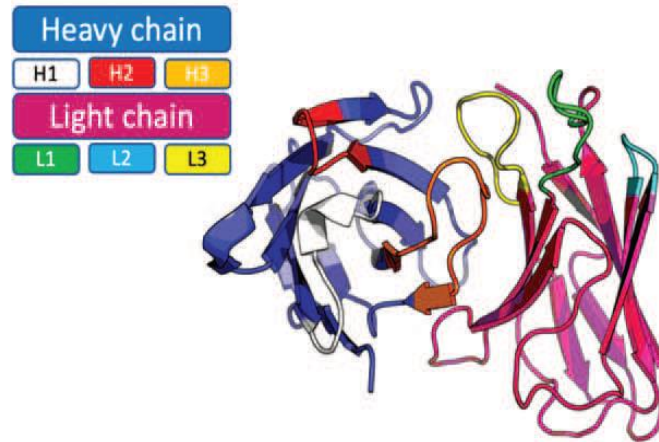
Github Address: https://ai2healthcare.github.io/
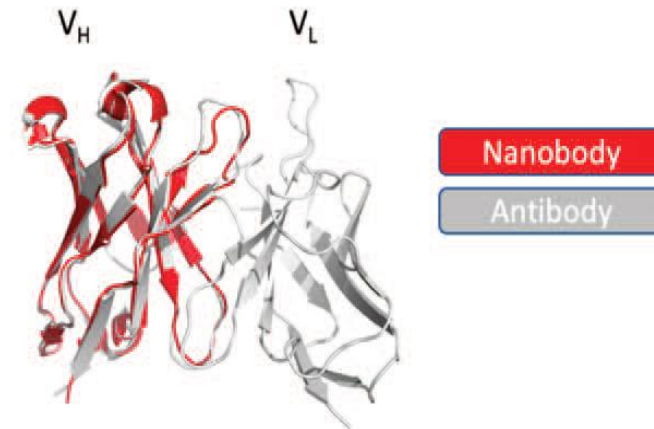
# Introduction

# Antibody



A) IgG Molecule
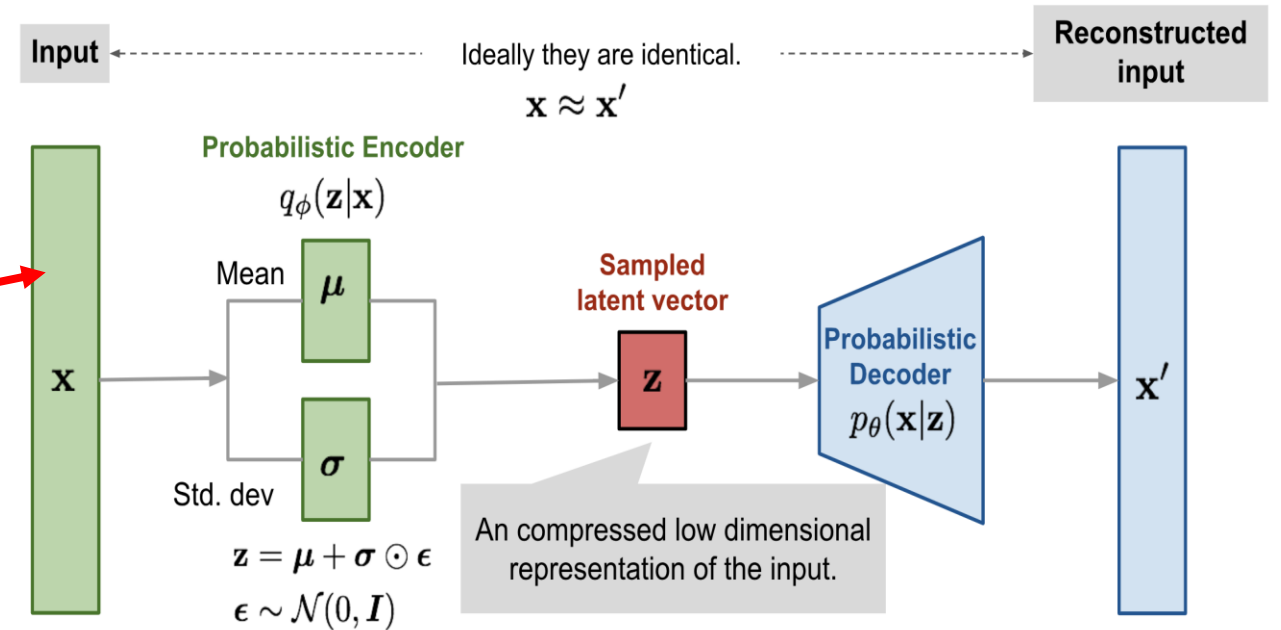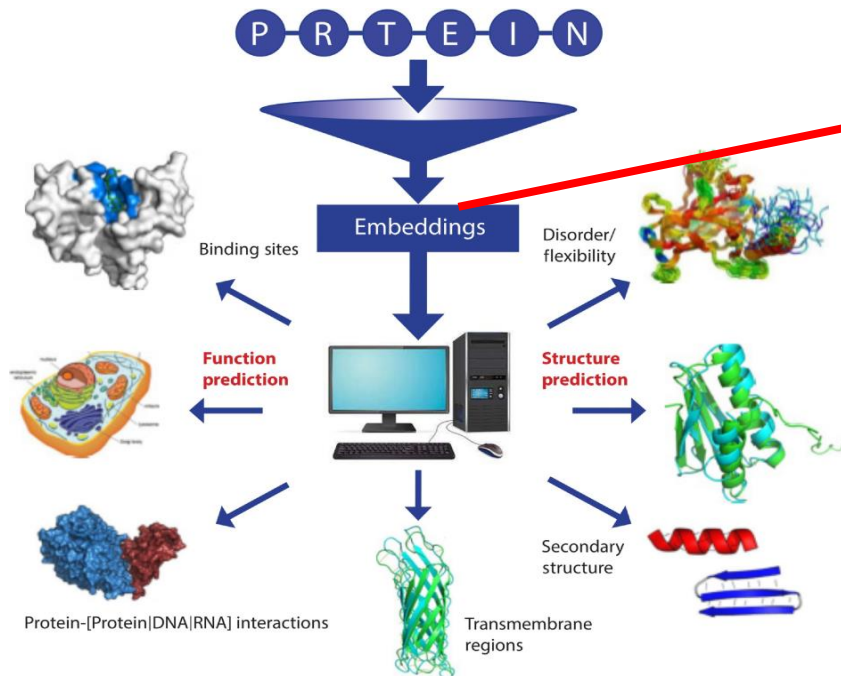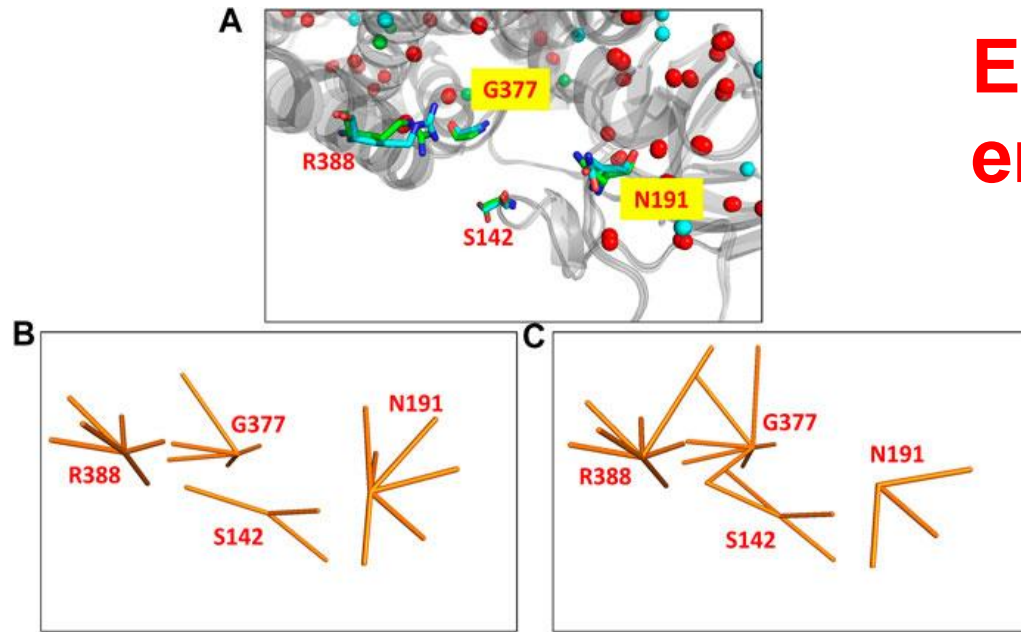
B) Antibody-Antigen Binding

C) Variable Region $V_H/V_L$

D) Antibody vs Nanobody (VHH)

# Extract Protein Variation into embedding and reduce it via VAE



A

B                    C

P R T E I N

Embeddings

Binding sites

Disorder/
flexibility

Function
prediction

Structure
prediction

Protein-[Protein|DNA|RNA] interactions

Transmembrane
regions

Secondary
structure

Input · · · · · · · · · · Ideally they are identical. · · · · · · · · · · Reconstructed input

$$\mathbf{x} \approx \mathbf{x}'$$

**Probabilistic Encoder**

$$q_\phi(\mathbf{z}|\mathbf{x})$$

Mean    $\boldsymbol{\mu}$

$\mathbf{x}$

Std. dev    $\boldsymbol{\sigma}$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{I})$$

Sampled
latent vector

$\mathbf{z}$

An compressed low dimensional representation of the input.

**Probabilistic Decoder**

$$p_\theta(\mathbf{x}|\mathbf{z})$$

$\mathbf{x}'$

# Joint Amino Acid Variation and Expression Embedding



Trends in Biochemical Sciences

# Normal



# Disease

$$X = A^T X + Z$$

**Infer Protein Causal Networks**

**Drug Target Causal Networks**

# Graphic Neural Networks (GNN)

$\begin{bmatrix} y_{11} \\ y_{12} \end{bmatrix}$ **1**

**2** $\begin{bmatrix} y_{21} \\ y_{22} \end{bmatrix}$

$\begin{bmatrix} y_{31} \\ y_{32} \end{bmatrix}$ **3**

**4** $\begin{bmatrix} y_{41} \\ y_{42} \end{bmatrix}$

**5**

$\begin{bmatrix} y_{51} \\ y_{52} \end{bmatrix}$

Probably the most common application of representing data with graphs is using molecular graphs to represent chemical structures

The caffeine molecule

chemical name: 1, 3, 7-trimethylxanthine
chemical formula: $C_8H_{10}N_4O_2$

C — carbon atom
H — hydrogen atom
N — nitrogen atom
O — oxygen atom
CH3 — methyl radical

© 2010 Encyclopædia Britannica, Inc.

# Graphic Neural Network (GNN)

- **Tasks of GNN: Graph level, node level and edge level**

  - **Aggregation**

$$M_{jv}^{l-1} = MSG(h_j^{l-1}, h_v^{l-1}, e_{jv}^{l-1})$$

$$AGG_v^l = AGG^l(\{ M_{jv}^{l-1} | j \in \mathcal{N}(v) \})$$
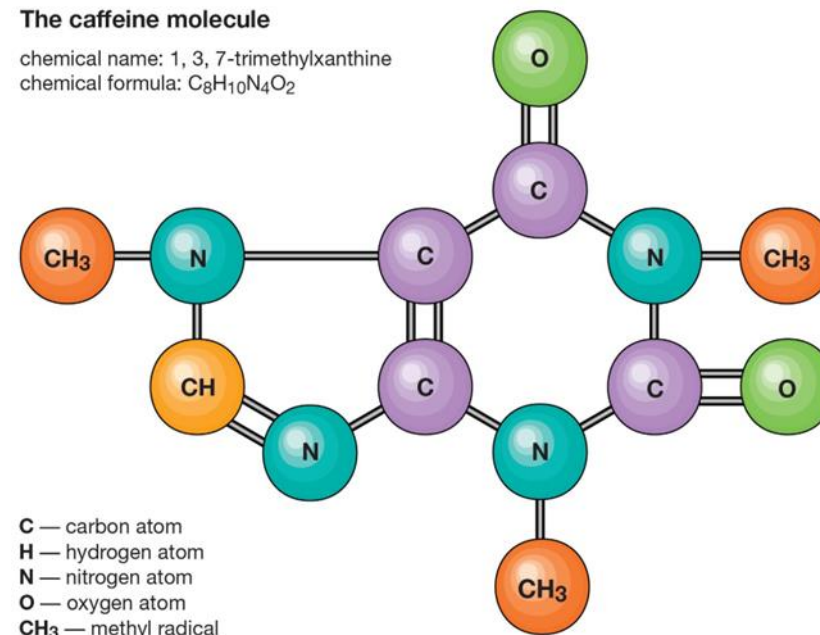
$$h_v^l = Combine^l(h_v^{l-1}, AGG_v^l)$$



1) average messages from neighbors

TARGET NODE

INPUT GRAPH

$h_v^{l-1}$

$AGG_v^l$

$Combine^l$

2) apply neural network

Sharma 2020, Introduction to Graph Neural Networks

# Pipelines of Graphic Neural Networks



$$h_G = READOUT(\{h_v^L, v \in \mathcal{V}\})$$

Distill, 2021; A Gentle Introduction to Graph Neural Networks. https://distill.pub/2021/gnn-intro

# Directed Acyclic Graph Neural Networks

- **A DAG is a directed graph without cycles**

- **updating node representations based on those of all their predecessors sequentially, such that nodes without successors digest the information of the entire graph.**

$$AGG_v^l = \sum_{u \in \mathcal{P}(v)} \alpha_{vu}^l (h_v^{l-1}, h_u^l) h_u^l$$

$$\alpha_{vu}^l (h_v^{l-1}, h_u^l) = \underset{u \in \mathcal{P}(v)}{\text{softmax}} ((w_1^l)^T h_v^{l-1}$$

$$+ (w_2^l)^T h_u^l + (w_3^l)^T y(u, v))$$

$$h_v^l = F^l (h_v^{l-1}, AGG_v^l) = GRU^l (h_v^{l-1}, AGG_v^l)$$

$$h_G = FC(\max - \text{pool}(\|_0^L h_v^l, )$$
$$\underset{v \in \mathcal{T}}{}$$

$$\| \max - \text{pool}(\|_0^L \tilde{h}_u^l))$$
$$\underset{u \in \mathcal{S}}{}$$



FIGURE 1. An example phenotype for a Directed Acyclic Graph Neural Network (DAG-NN).

$$\mathcal{P}(v) = set\ of\ preceeding\ nodes$$

This also allows producing a single output for the whole graph

# Embedding Test Statistics

$$T = \left(\overline{Z}_m - \overline{Z}_w\right)^T \Lambda^{-1}\left(\overline{Z}_m - \overline{Z}_w\right), \mathbf{T} \sim \chi^2_{(k)}$$

**K: dimension of Embedding**

# Causal Test for Embedding

# Test Causal Relationship between Drug and Disease GNNs

**Causal Model:**

$$Y = f_y(X, N_y), X \perp\!\!\!\perp N_y$$

**Fake Sample**

$Y'=G(Z,X)$

$G(Z,X)$ Neural Networks

Z

**Random Noise**

X

Discriminator Neural Networks

**Real or Fake**

X

Y (Effect)

**Drug GNN**

**Disease GNN**

**Real Data**

Y=Disease GNN

X=Drug GNN

# Test for Causation between Drug and Disease Networks

- **Define**

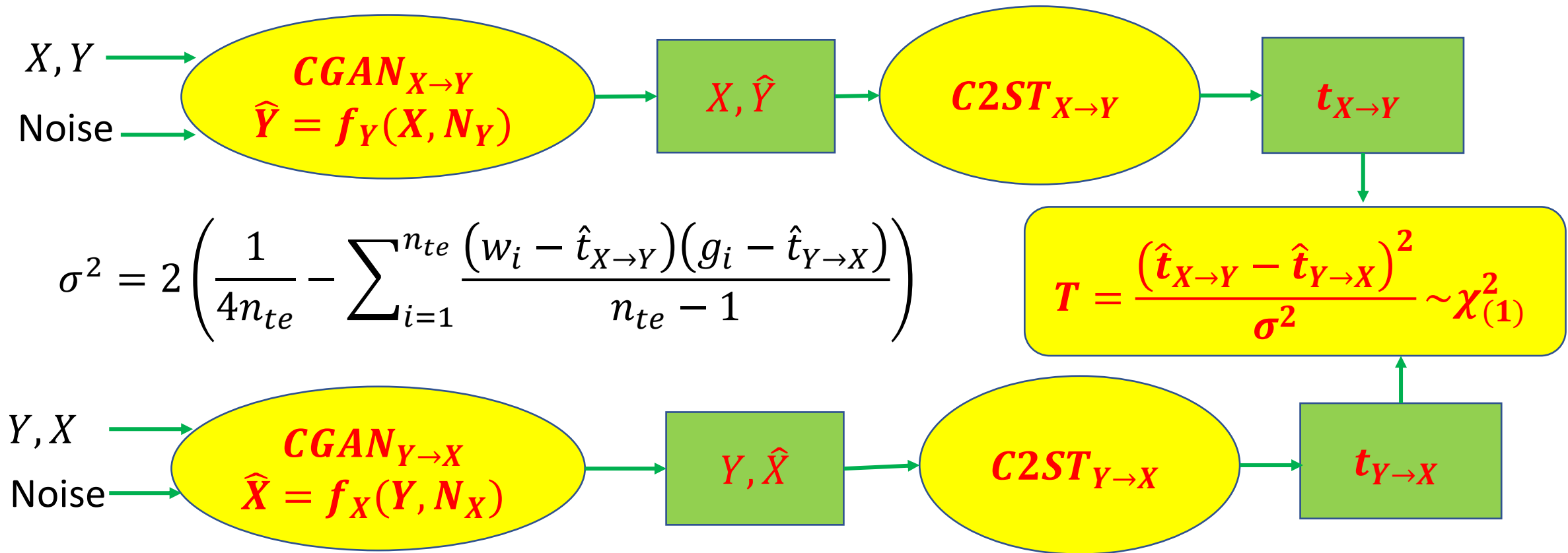$Y$ =output of Disease GNN,

$X$ =output of drug GNN



$$\sigma^2 = 2\left(\frac{1}{4n_{te}} - \sum_{i=1}^{n_{te}} \frac{(w_i - \hat{t}_{X\to Y})(g_i - \hat{t}_{Y\to X})}{n_{te}-1}\right)$$

$$T = \frac{(\hat{t}_{X\to Y} - \hat{t}_{Y\to X})^2}{\sigma^2} \sim \chi^2_{(1)}$$

Xiong MM (2022) Artificial Intelligence and Causal Inference. CRC Press

Test Causation of Protein on Disease (or Trait)

# Test Causation of Gene with Trait

GPN (Genomic Pre-trained Network)



$Y = $ Trait,
$X = $ Embedding (DNAs)

$$T = \frac{\left(\hat{t}_{X \to Y} - \hat{t}_{Y \to X}\right)^2}{\sigma^2} \sim \chi^2_{(1)}$$

# A. Training



**Log-likelihood over masked positions** (to be maximized)

$$L = \log p_1(R) + \dots + \log p_{l-1}(C)$$

$\nabla L$  **Gradient computation**

**Probabilities over amino-acids at** each residue position

| | | | | | |
|---|---|---|---|---|---|
| **$p_1$** | **$p_2$** | **$p_3$** | | **$p_{l-1}$** | **$p_l$** |
| A 0.1 | A 0.5 | A 0.3 | | A 0.1 | A 0.0 |
| Q 0.2 | Q 0.1 | Q 0.2 | | Q 0.0 | Q 0.9 |
| N 0.1 | N 0.0 | N 0.5 | ... | N 0.0 | N 0.1 |
| R 0.0 | **R 0.1** | R 0.0 | | R 0.0 | R 0.0 |
| C 0.0 | C 0.1 | C 0.0 | | **C 0.9** | C 0.0 |

**Test Causal of Protein (Gene) with Trait (Disease)**

**Semantic Change**

Feed Forward & Softmax

$$w_1|Z_{norm} - Z_{disease}| + w_2 Z$$

**L Transformer Layers**

Transformer Layer

**Parameters Update**

$Z$

Transformer Layer

$Y = \text{Trait},$
$X = w_1|Z_{norm} - Z_{disease}| + w_2 Z$

Transformer Layer

**Randomly Masked** Spike Protein Sequence

A / N ... / Q

Noise

$$CGAN_{X \to Y}$$
$$\widehat{Y} = f_Y(X, N_Y)$$

$X, \widehat{Y}$

$$C2ST_{X \to Y}$$

$t_{X \to Y}$

$$T = \frac{(\hat{t}_{X \to Y} - \hat{t}_{Y \to X})^2}{\sigma^2} \sim \chi^2_{(1)}$$

$Y, X$
Noise

$$CGAN_{Y \to X}$$
$$\widehat{X} = f_X(Y, N_X)$$

$Y, \widehat{X}$

$$C2ST_{Y \to X}$$

$t_{Y \to X}$

# A Path to Uncovering Mechanism of Complex Trait

**A** Probe-based spatial transcriptome (merFISH, STARmap, and seqFISH)

Cells (n)

Target genes (t): ~hundreds

t × n

Spatial location: hundreds

Single-cell RNA-seq

Cells (n')

All genes (m): ~23,000

m × n'

Without spatial location

Inferred all-gene spatial transcriptome

Cells (n)

Target genes:

Nontarget genes:

m × n

Spatial location: all genes

**Manifold alignment**

**B**

Dimension reduction

Expression stabilization

Normalization Two rounds integration 1st: Reciprocal PCA

Infer expression by weighted KNN

Cell location from ST

Accurate inference of genome-wide spatial expression with iSpatial

# 18. Protein Language Model

**ESM2:Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences**

- **the UniParc database**

- **Number of amino acid sequences**: 250 millions

**ProGen: Language Modeling for Protein Generation**

Ali Madani et al. 2020

- **1.2 billion parameter conditional language**

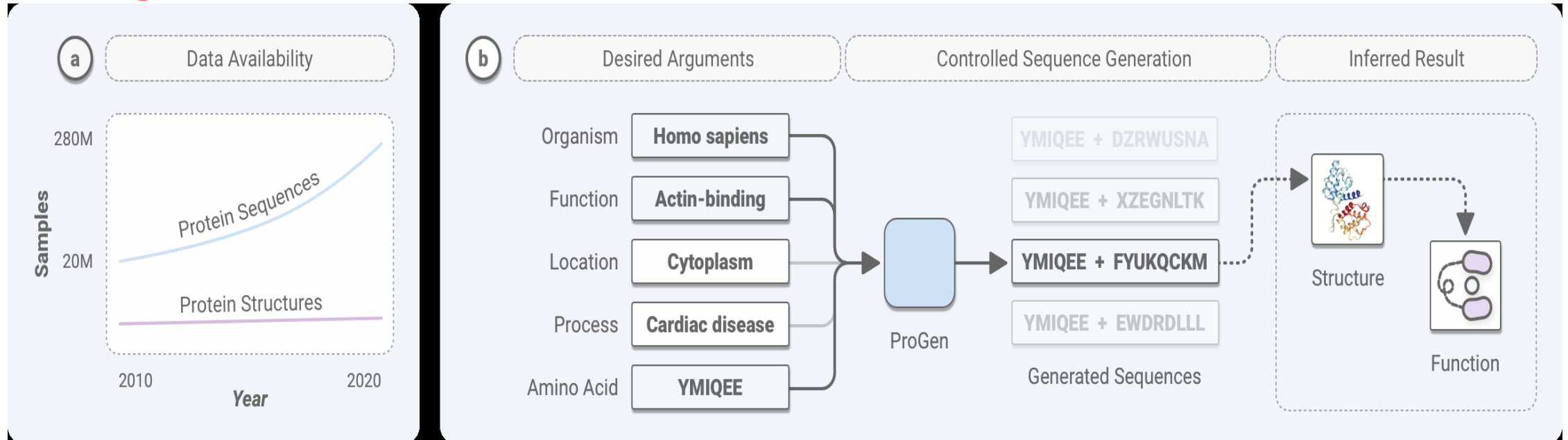- **Number of amino acid sequence:** 280 millions

Generative modeling for protein engineering is key to solving fundamental problems in synthetic biology, medicine, and material science

# 18.1. ProGene

## 18.1.1. Generating proteins with desired properties

- the development of new enzymes, antibody, therapies, and sensors

- However, leading experimental techniques for protein engineering such as directed evolution (Arnold, 1998) still rely on heuristics and random mutations to select initial sequences for rounds of evolution.

- The raw amino acid sequence encodes a protein. This chain of amino acids folds in ways that exhibit local (secondary) and global (tertiary) structure, which in turn determines unique functions.

- Unfortunately, obtaining three-dimensional structural information for proteins is expensive and time consuming. Consequently, there are three orders of magnitude more raw sequences than there are sequences with structural annotations, and protein sequence data grow exponentially.

- **By conditioning on these tags, ProGen provides a new method for protein generation that can be tailored for desired properties**



a) Protein sequence data is growing exponentially as compared to structural data. b) We utilize protein sequence data along with and keyword tags to develop a conditional language model: ProGen.

Ali Madani et al. 2020    ProGen: Language Modeling for Protein Generation

# 18.1.2. Methods

- **Notations**

    **Amino Acid Sequence:** $a = \left(a_1, a_2, \quad , a_{n_a}\right)$

    **Conditional Tag:** $c = \left(c_1, c_2, \ldots, c_{n_c}\right), \ n = n_a + n_c$

    **Sequence:** $x = (c; a)$

    **Distribution:** $P(x) = \prod_{i=1}^{n} P(x_i | x_{<i})$

    **Dataset:** $D = \left[x^1, \ldots, x^{|D|}\right]$

- **Loss Function**

$$L(D) = -\sum_{k=1}^{|D|} \frac{1}{n_i} \sum_{i=1}^{n_i} \log P_\theta\left(x_i^k \mid x_{<i}^k\right)$$

$$p(x) = \prod_{i=1}^{L} p(x_i | x_1 \dots x_{i-1})$$

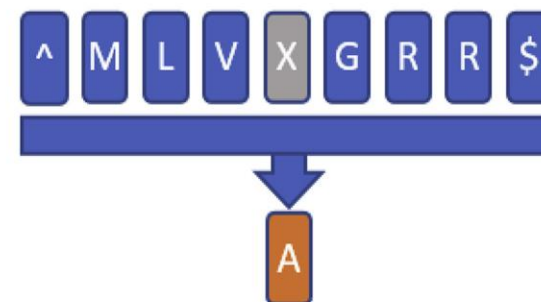$$p(x) = \prod_{i=1}^{L} p(x_i | x_1 \dots x_{i-1}) p(x_i | x_{i+1} \dots x_L)$$

$$p(x) = \prod_{i=1}^{L} p(x_i | x_1 \dots x_{i-1}, x_{i+1} \dots x_L)$$



Processes sequence in one direction

Processes sequence in each direction independently

Processes whole sequence

$$p(x_i = A | x_1 \dots x_{i-1})$$

$$p(x_i = A | x_1 \dots x_{i-1}) p(x_i = A | x_{i+1} \dots x_L)$$

$$p(x_i = A | x_1 \dots x_{i-1}, x_{i+1} \dots x_L)$$

$$x_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \vdots & \vdots \\ 0 & \cdots & 1 \end{bmatrix}$$

$$= (.)_{n \times d}$$

Output Probabilities

Softmax
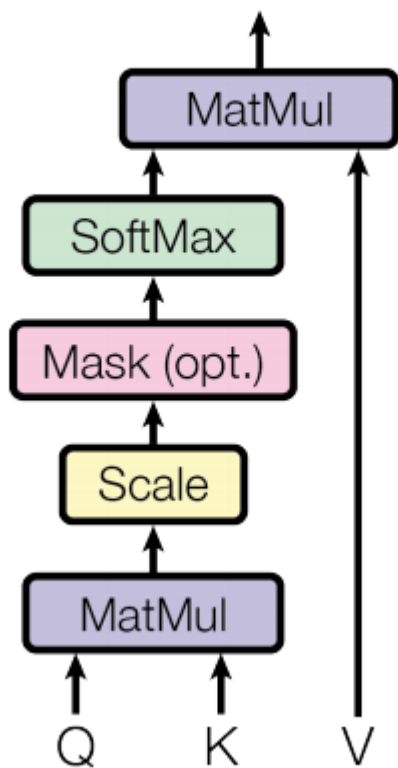
Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

Ashish Vaswani et al. 2017

**Attention Is All You Need**

arXiv:1706.03762
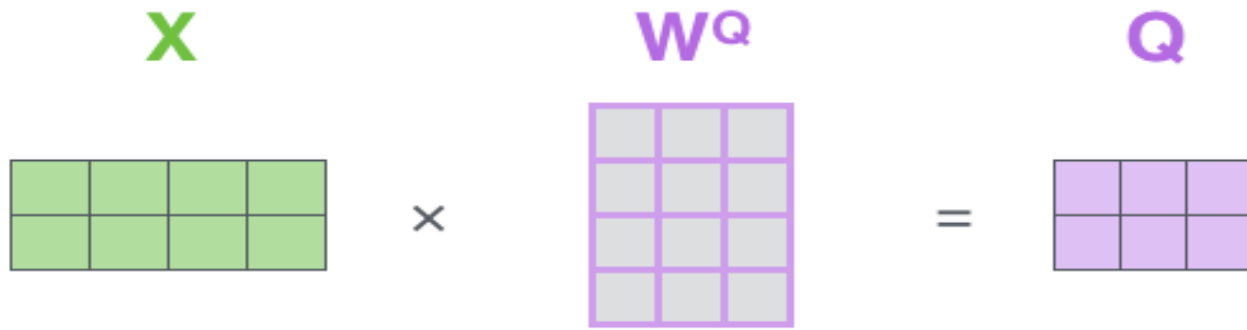
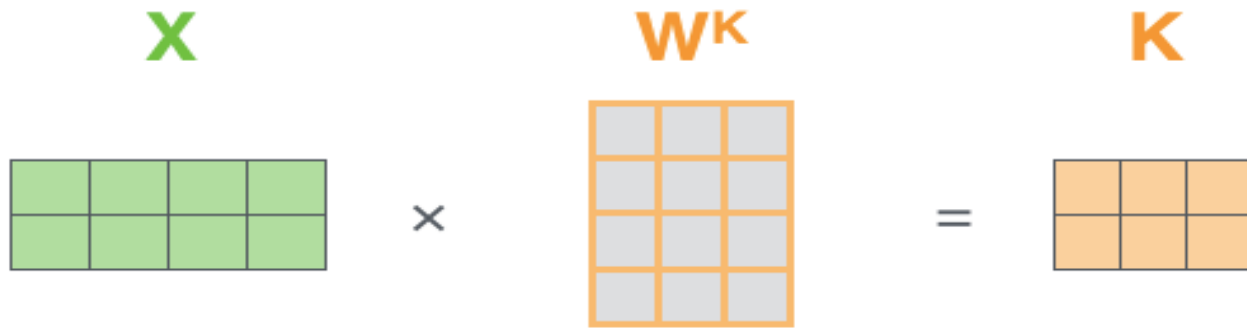## Attention Mechanism

Scaled Dot-Product Attention



**Matrices:**

queries $Q \in R^{N \times D_k}, K \in R^{M \times D_{k'}}V \in R^{M \times D_v}$

Attention $(Q, K, V) = softmax\left(\dfrac{QK^T}{\sqrt{D_k}}\right)V = AV$

$$A = \left(\alpha_{ij}\right)_{N \times M}$$

$$h_i = \sum_j \alpha_{ij}V_j$$

$$\alpha_{ij} = \frac{\exp\{\sum_{l=1}^{D_k} q_{il}k_{jl}\}}{\sum_{j'=1}^{M} \exp\{\sum_{l=1}^{D_k} q_{il}k_{j'l}\}}$$
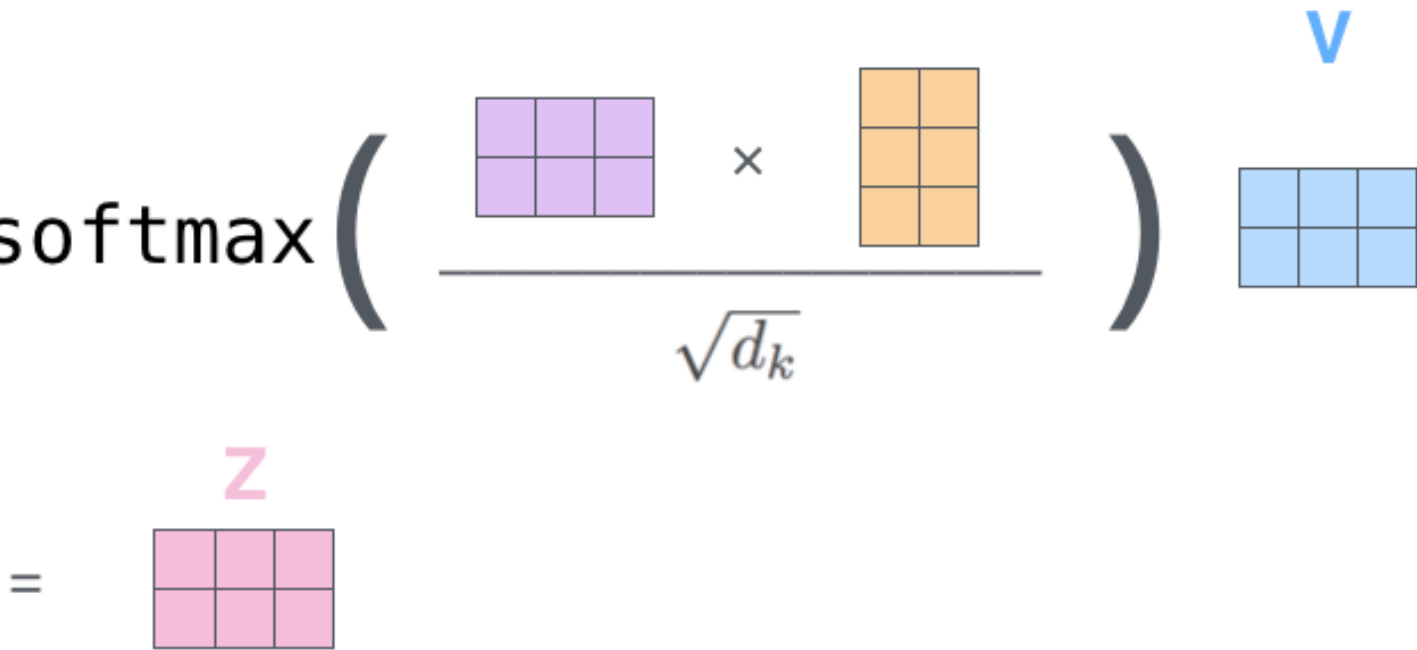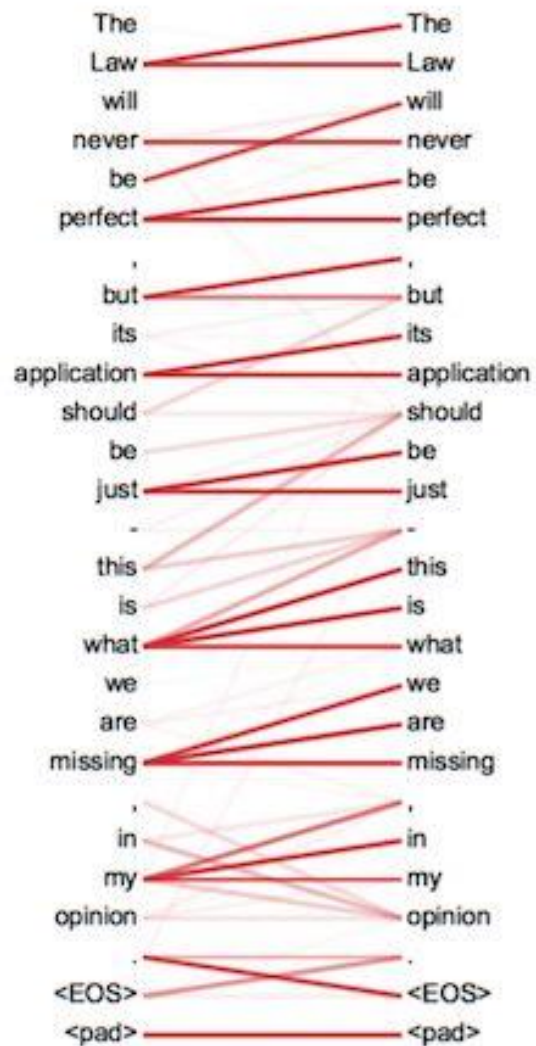
$$W^Q = W_j^1, Q = XW_j^1$$

$$W^k = W_j^2, K = XW_j^2$$

$$W^V = W_j^3, V = XW_j^3$$

$$\text{Attention}\,(Q, K, V) = \mathtt{softmax}\left(\frac{\boxed{Q} \times \boxed{K^T}}{\sqrt{d_k}}\right)\boxed{V}$$

$$= \boxed{Z}$$

$$h_j = Z = Attention\,(XW_j^1, XW_j^2, XW_j^3)$$

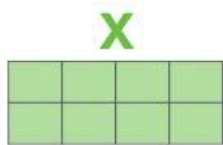# Multihead

1) This is our input sentence*

2) We embed each word*

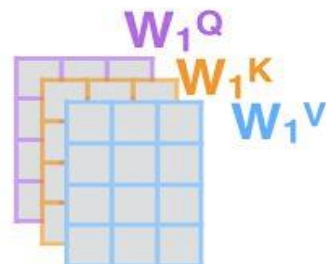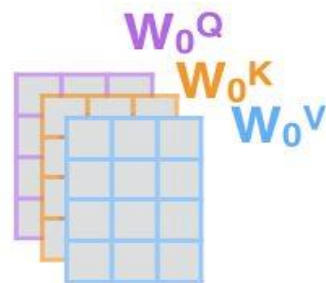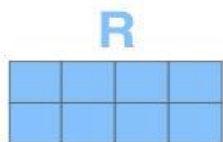3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

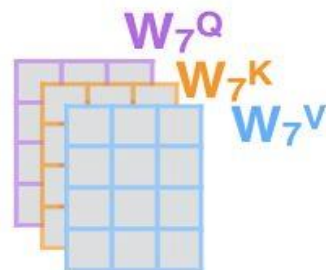5) Concatenate the resulting Z matrices, then multiply with weight matrix W° to produce the output of the layer

**X**

Thinking Machines

**$W_0^Q$**
**$W_0^K$**
**$W_0^V$**

$$h_j = z_j = Attention\ (XW_j^1, XW_j^2, XW_j^3)$$
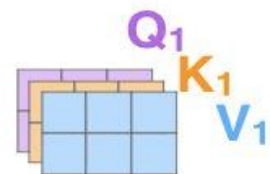
**$Q_0$**
**$K_0$**
**$V_0$**

**$Z_0$**

$= h_0$

**W°**

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one
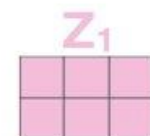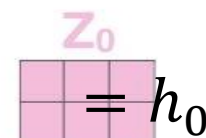
**$W_1^Q$**
**$W_1^K$**
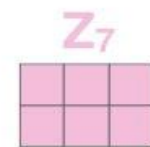**$W_1^V$**

**$Q_1$**
**$K_1$**
**$V_1$**

**$Z_1$**

**Z**

**R**

...

**$W_7^Q$**
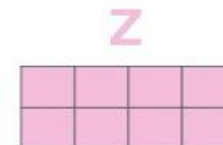**$W_7^K$**
**$W_7^V$**

...

**$Q_7$**
**$K_7$**
**$V_7$**

...

**$Z_7$**

$$Z = MultiHead(X, m)$$
$$= [h_0, \ldots, h_{m-1}]W_0$$

# Residual Connection and Layer Normalization

$$H = Multihead\ (X, m) + X$$

$$\bar{H} = LayerNorm\ (H)$$

## Position-wise FFN.

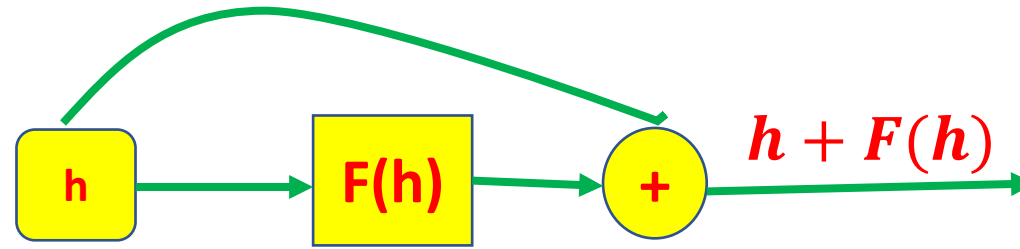**The position-wise FFN is a fully connected feed-forward module that operates separately and identically on each position**

$$FFN(h) = ReLU(\bar{H}W^1 + b^1)W^2 + b^2$$

**the outputs of previous layer:**

$$\bar{H}_i = [\bar{h}_1, \dots, \bar{h}_m], W^1 \in R^{m \times D_f}, W^2 \in R^{D_f \times m}, b^1 \in R^{D_f}, b^2 \in R^m$$

# Residual Connection and Normalization



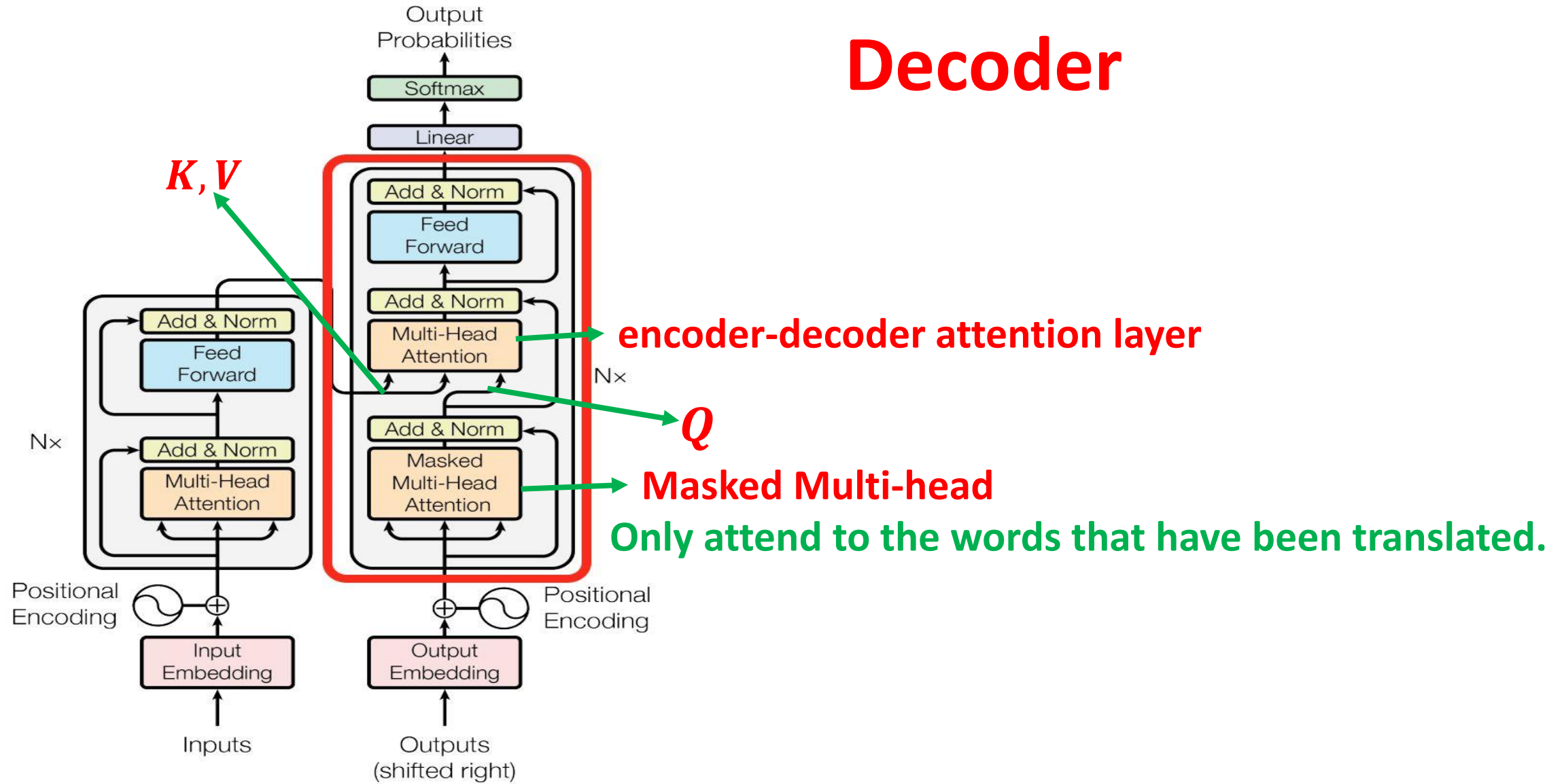$$h' = LayerNorm(SelfAttention(h^{(l)}) + h^{(l)})$$

$$h^{(l+1)} = LayerNorm(FFN(h') + h')$$

$$X_{i+1} = LayerNorm(FFN(\bar{H}_i) + \bar{H}_i)$$

**Decoder**

$K, V$

encoder-decoder attention layer

$Q$

Masked Multi-head

Only attend to the words that have been translated.

# Scores

$$Scores\ (X_0) = LayerNorm(X_L)W_{Vocab}$$

During training, these scores are the inputs of a cross entropy loss function. During generation, the scores corresponding to the final token are normalized with a softmax, yielding a distribution for sampling a new token.

# TRANSFORMER MODELS: AN INTRODUCTION AND CATALOG

Xavier Amatriain

Los Gatos, CA 95032

xavier@amatriain.net

# Cover (One Topic at One Lecture)

- **ProGen: Language Modeling for Protein Generation**

- **ProGen Generate functional protein sequences**

- **ESM2: Language  Model  generalize beyond natural proteins**

- **Language Models and Diffusion Process**

- **Language Model for Antibody Design**

- **MULTI-LEVEL PROTEIN STRUCTURE PRE-TRAINING WITH PROMPT LEARNING**

- **DNA Language Model**

- **Protein Language Models for Protein Docking**

- **Embeddings from  language models predict conservation and variant effects**

- **Table Data, Language and Omics Data Embedding**

- **A new path to uncovering Disease Mechanism using language models, Causal Inference and Omics**

- **A New Paradigm for Drug Discovery and Drug Repurposing**