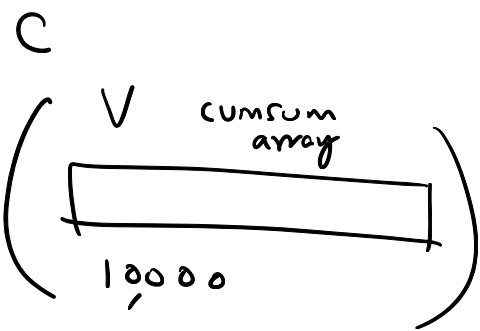
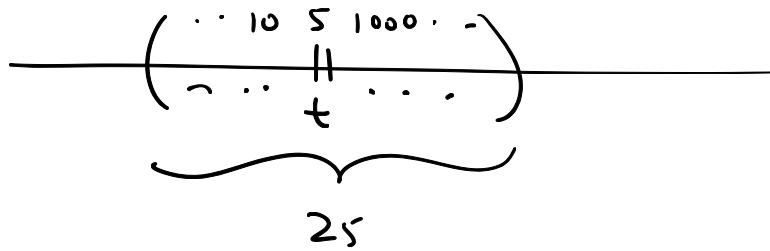


Dataset:

→ for all sequence
for all offsets (0, 1, 2)

S = seq of tokenized (int) kmers

[5, 15, 0, 1000, ...]

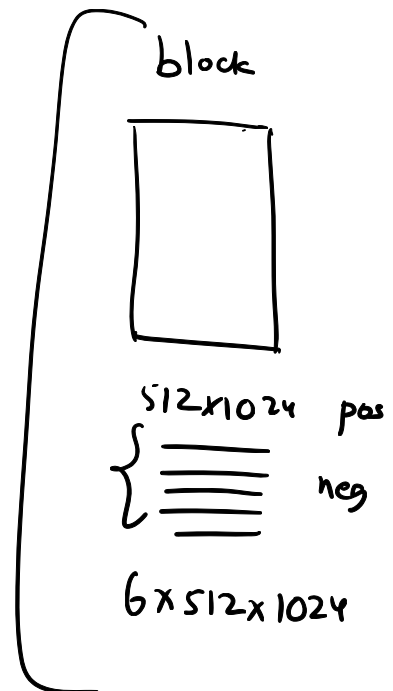


np. searchsorted (C, $\frac{5 \times 512 \times 1024}{}$)

Dataloader

→ num workers = 4

→ batch = 1



worker-id = 0, 1, 2, 3

num workers = 4

array of NegSampler [0] = 5x512x1024

NegSampler [worker-id] = np. searchsorted ()

shuffle the seq idx
↓

for(k, seq) in enumerate(sequences)

if k % num_workers == worker_id:

for offset 0, 1, 2:

tokenize

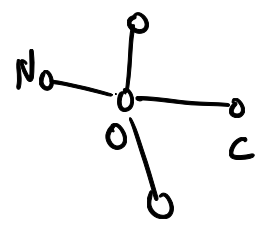
// generate blocks of 6x512x1024

create 512x1024 pos pairs

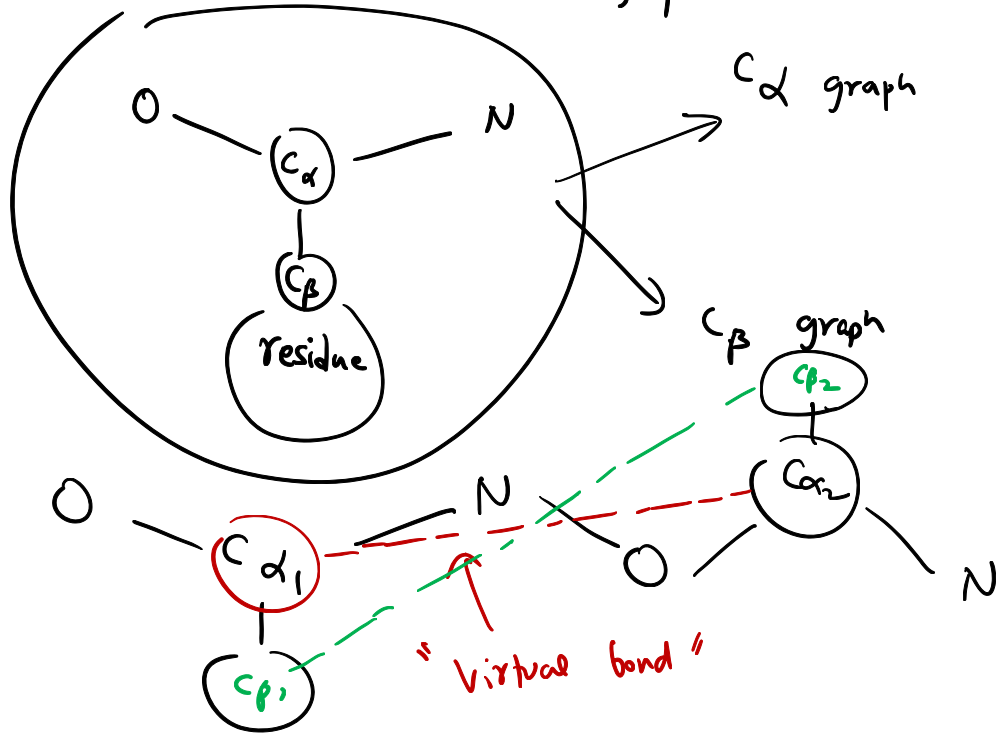
neg sampls [worker_id] = search sorted

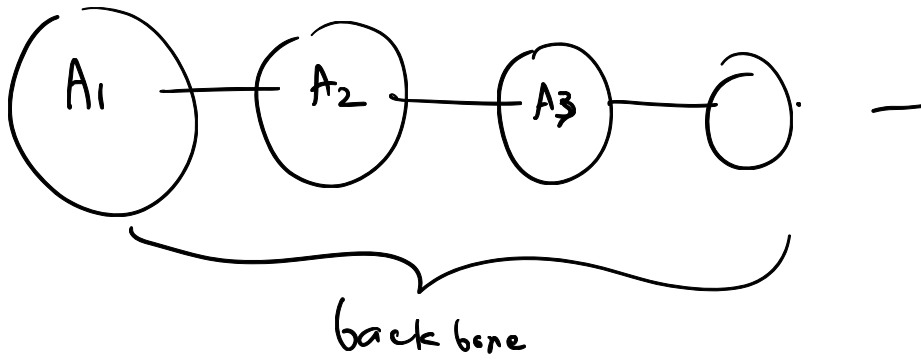
graph neural networks

protein molecule → 3D coords of each atom

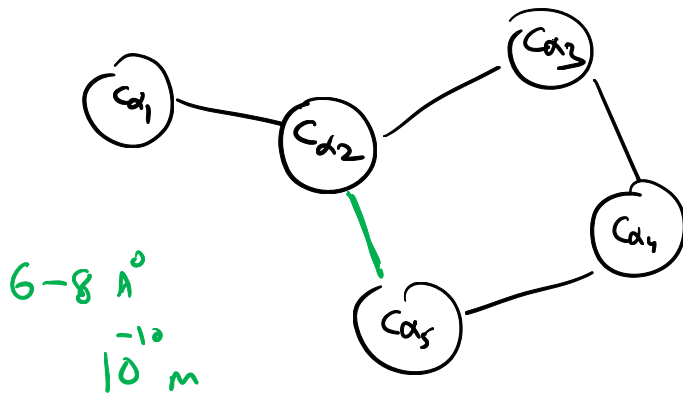


molecular/atom-level graph

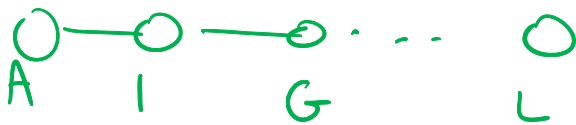
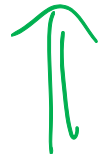




Sequence

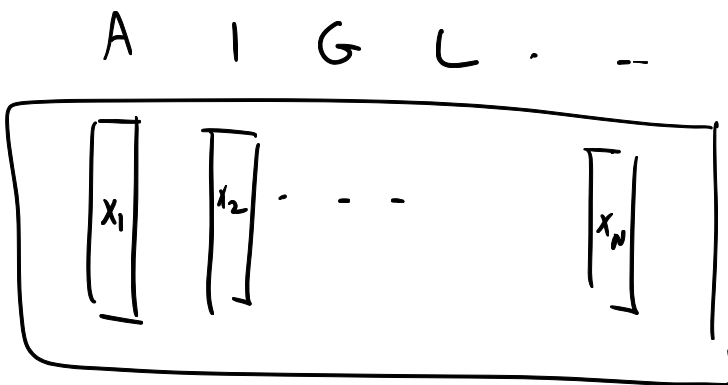


$C\alpha$ - graph



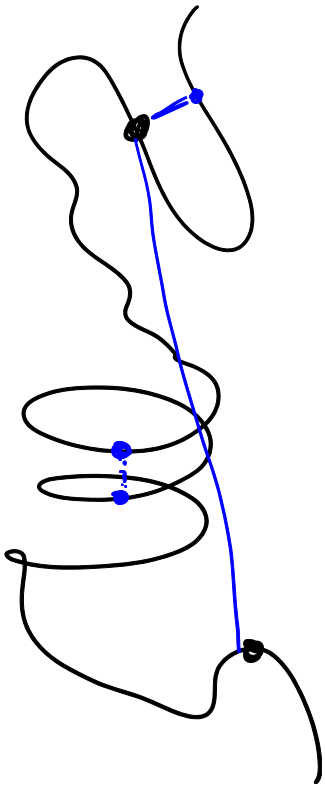
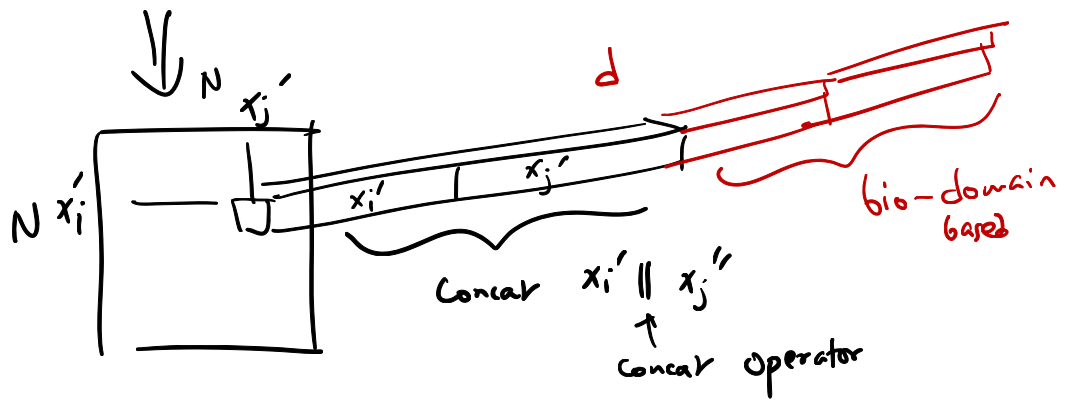
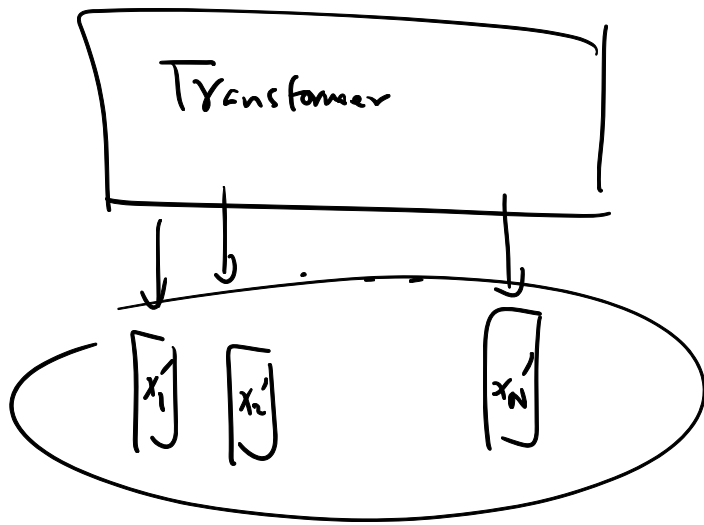
Structure prediction
(machine translation)

Alphafold:



← sequence





a) 0/1 prediction | Binary logistic regression
edge exists or not

b) $d(A_i, A_j) = \text{Real value}$ ← Regression problem

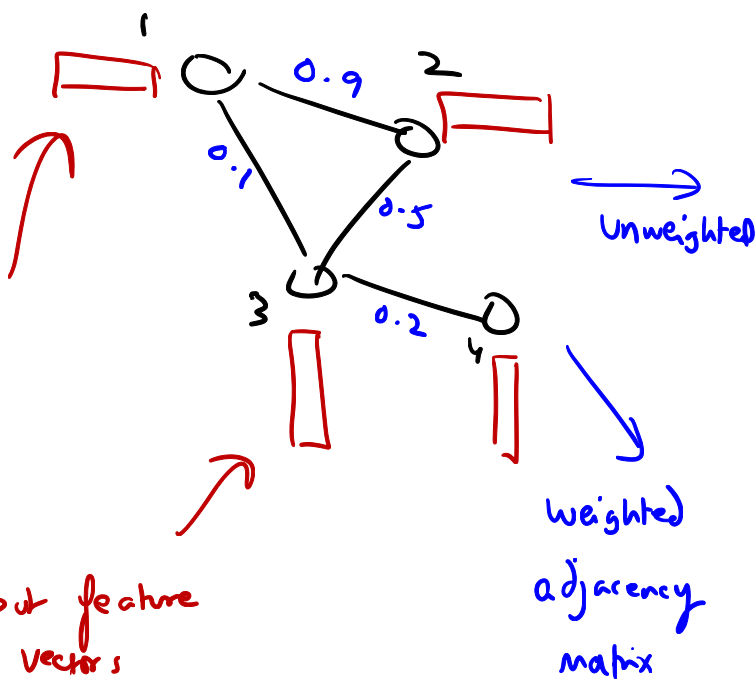
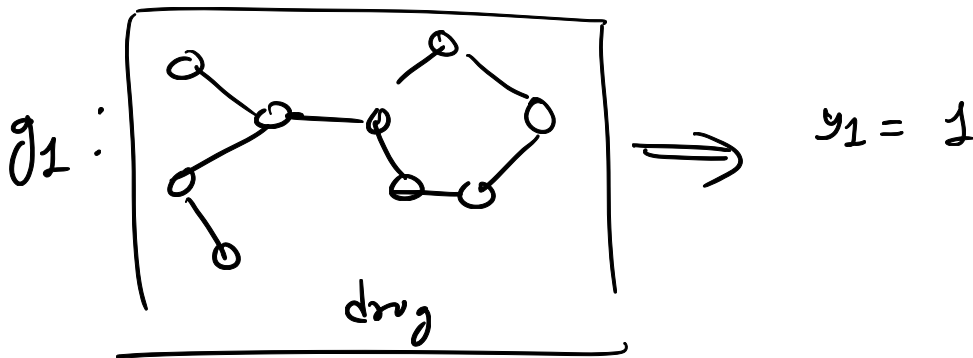
c) Create k -bins

$d(A_i, A_j)$ → 0-2 A
→ 2-4 A
→ 4-6 A
→ ...
→ > 22 A

} k -class problem

Molecular graphs: small chemical compounds

atom-level graphs \rightarrow labels/value per graph



$$A + I$$

	1	2	3	4
1	0	1	1	0
2	1	0	1	0
3	1	1	0	1
4	0	0	1	0

	1	2	3	4
1	1	0	0	0
2	0	1	0	0
3	0	0	1	0
4	0	0	0	1

$a_{ij} = a_{ji}$

Identity

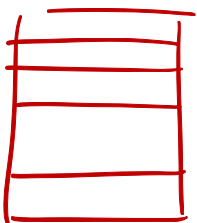
Adjacency value between v_i and v_j

0	0.9	0.1	0
0.9	0	0.5	0
0.1	0.5	0	0.2
0	0	0.2	0

A

Input feature vectors

d



N

d-dim input features per node

d_1			
	d_2		
		d_3	
			d_n

D: degree matrix

$$d(v_i) = \sum_{j=1}^n a_{ij}$$

diagonal matrix

Unweighted:

$D =$

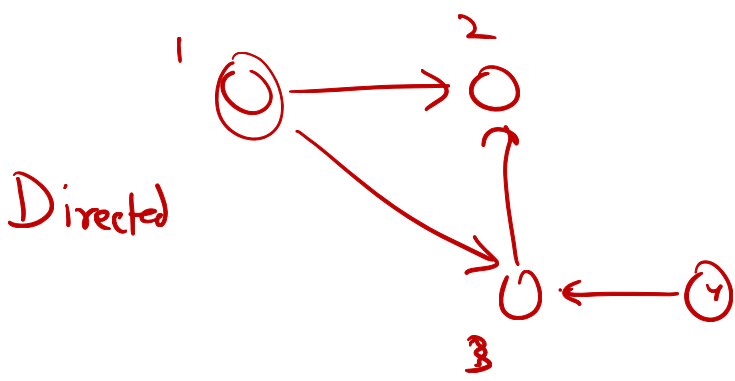
2	0	0	0
0	2	0	0
0	0	3	0
0	0	0	1

weighted

1.0			
	1.4		
		0.8	
			0.2

Undirected graphs: A is symmetric

$a_{ij} = a_{ji}$



A

	1	2	3	4	outdegree
1	0	1	1	0	2
2	0	0	0	0	0
3	0	1	0	0	1
4	0	0	1	0	1
in-degree	0	2	2	0	

D_{in}

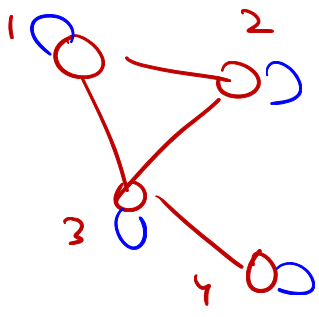
0
2
2
0

D_{out}

2
0
1
1

Transition Matrix

Undirected



$A+I$

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

D

3
3
0
0
4
2

prob
transition
matrix

$$M = m_{ij} = \frac{a_{ij}}{d_i} =$$

$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
0	0	$\frac{1}{2}$	$\frac{1}{2}$

a symmetric

$$= D^{-1} A \quad \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}$$

$$D = \begin{pmatrix} 3 & & & \\ & 3 & & \\ & & 4 & \\ & & & 2 \end{pmatrix}$$

$$D^{-1} = \begin{pmatrix} \frac{1}{3} & & & \\ & \frac{1}{3} & & \\ & & \frac{1}{4} & \\ & & & \frac{1}{2} \end{pmatrix}$$

$$\begin{pmatrix} DD^{-1} = I \\ D^{-1}D = I \end{pmatrix}$$

$\frac{1}{3}$	0
$\frac{1}{3}$	
$\frac{1}{4}$	
	$\frac{1}{2}$

D^{-1}

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

A

Graph Laplacian

$$L = D - A$$

$$\begin{pmatrix} 3 & & & \\ & 3 & & \\ & & 4 & \\ & & & 2 \end{pmatrix} - \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

Normalized Laplacian

$$\begin{aligned}L_a &= \bar{D}^{-1} L &= \bar{D}^{-1} (D - A) \\ & &= \bar{D}^{-1} D - \bar{D}^{-1} A \\ & &= I - M\end{aligned}$$