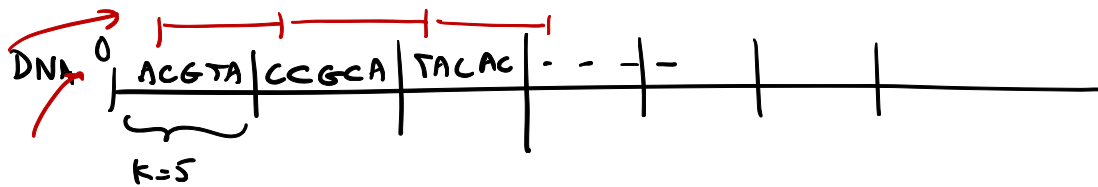


DNA/genome : 3 billion base pairs

Protein sequence : Uniprot \sim 100 million known sequence



k-mers
n-grams

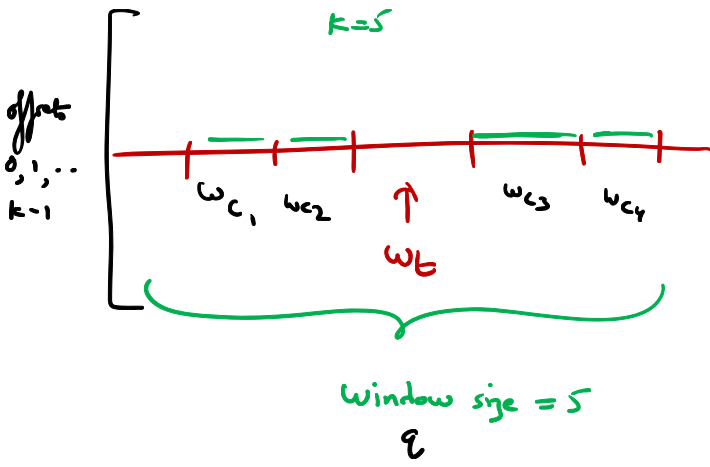


k=5

offset

I) Need 3 ORF $\begin{matrix} \nearrow 0 \\ \rightarrow 1 \\ \searrow 2 \end{matrix}$
↑
Open reading Frame \leftarrow transcription/translation

II) for k-mer
use all k offsets
 $k=0, 1, 2, \dots, k-1$



✗ offsets
✗ possible windows of size q

a) create positive example
 $P(w_{ci} | w_t)$
 $= \sigma(\vec{c}_i^T \vec{t})$

b) create γ negative samples

$$1 - P(w_{uj} | w_t)$$

DNA Vocab: $k=5$
 $4^5 = 2^{10} = 1024$

Practical: $\{A, C, G, T\} \leftarrow$ Base Alphabet

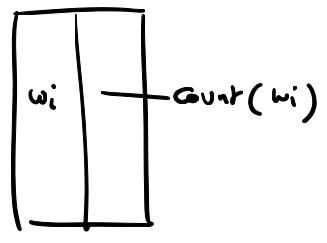
$\left. \begin{matrix} N: \text{unknown} \\ R, S, \dots \end{matrix} \right\} \leftarrow$ expanded Alphabet

Proteins
 $20 \leftarrow$ base alphabet
 $23-25 \leftarrow$ expanded protein alphabet
 $k=5 \quad 20^5$

Negative Sampling

any random word from vocab.

a) create the vocab
 \rightarrow count while tokenizing



$$P(w_i) = \frac{\text{count}(w_i)^\alpha}{\sum_{j=1}^{|V|} \text{count}(w_j)^\alpha}$$

$\alpha = 0.75$
 \leftarrow temp down high freq words
 bump up low freq words

$$CE := \underbrace{\left(\log P(w_{c_i} | w_t) \right)}_{1: +ve} + \underbrace{\sum_{j=1}^r \log (1 - P(w_{n_j} | w_t))}_{r: -ve}$$

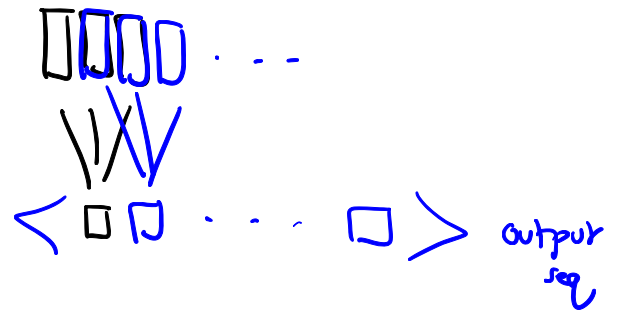
Sum over all output words c_i in the window
Sum over all windows

Self-Attention : Transformer

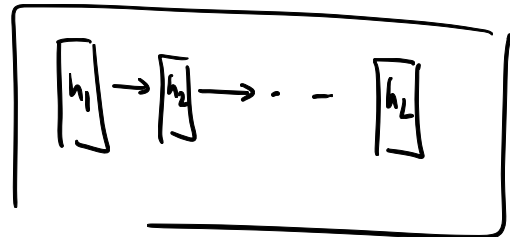
CNN: kernel-size = 3



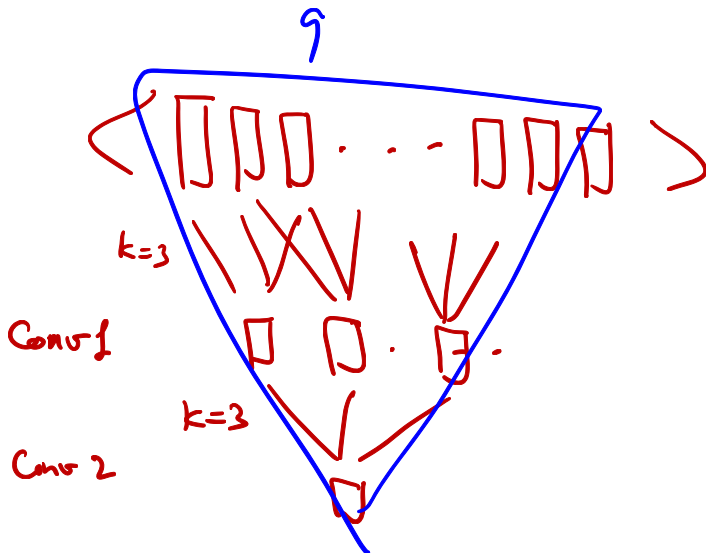
L - len seq
 $x_i \in \mathbb{R}^d$

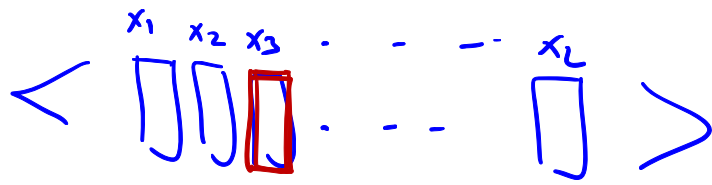


RNN:



CNN: hierarchical dependence





key/context

attention ($\vec{x}_j | \vec{x}_3$) $j=1, \dots, L$ (including $j=3$)

target/query
pair-wise

Compute for all elements as the target/query

$O(L^2)$ cost / Attention matrix

word2Vec Implementation

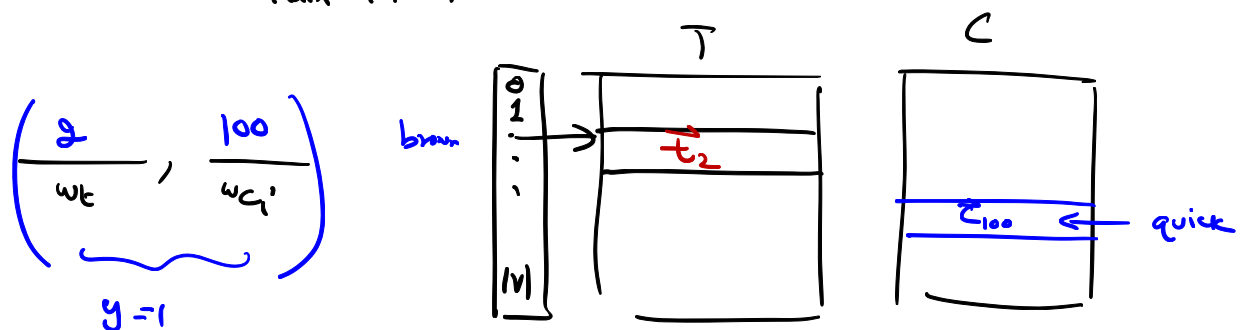
i) Create the vocab: V_k

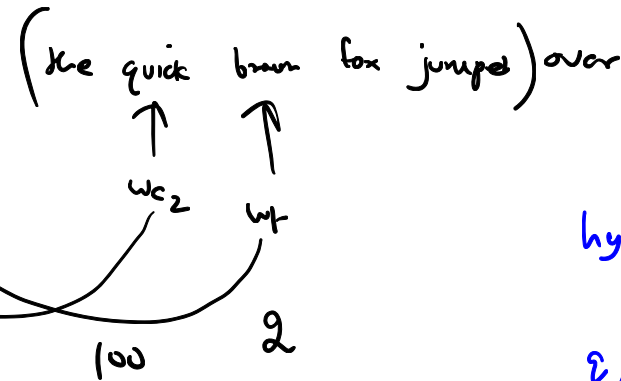
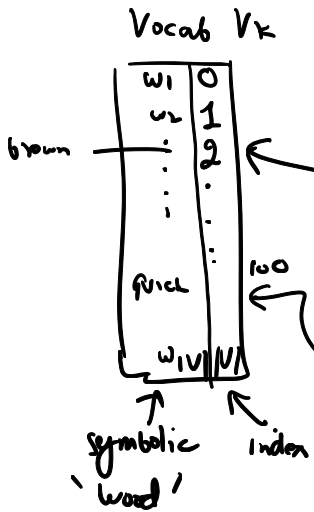
k - mer
word size

Model ():

init $\left[\begin{array}{l} T: V_k \times d \\ C: V_k \times d \\ d: \text{embed dim} \end{array} \right]$ Parameter
300, 512, 768, 1024 ?

forward ($y, x: B \times \underline{\text{index}}$ into the vocab)





- hyperparameters
- k : k-mer
 - ℓ/w : window size
 - r : negative sample
 - d : embed dim

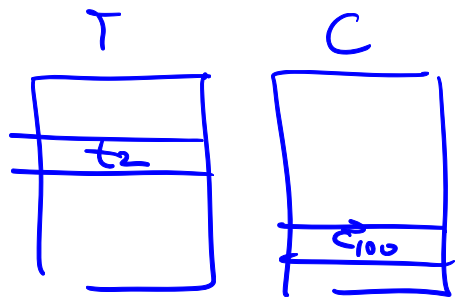
w_t	w_{c_j}	y
2	100	1
2	5	0
2	10000	0
2	13	0
\vdots		

forward(X, y):

$X: B \times 2$
 $y: B \times 1$

$$\log \left(\begin{pmatrix} T \\ T_2 \end{pmatrix}^T \begin{pmatrix} C \\ C_{100} \end{pmatrix} \right)$$

- a) with logits: just do the dot product!
- b) σ / \log



t-SNE to visualize in 2D.

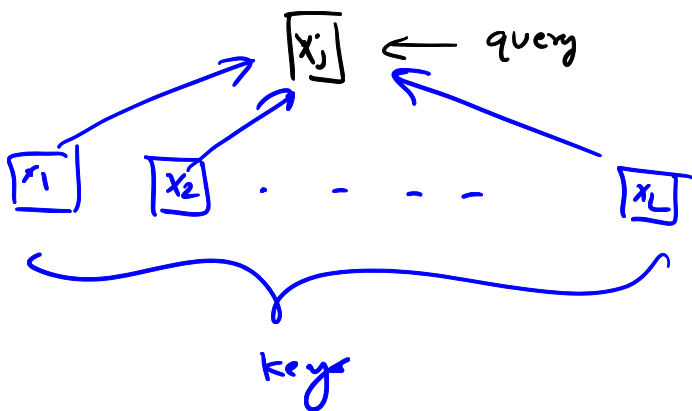
Self-Attention

block (seq of length L)



each element is a token

$$\alpha_{ji} \equiv \text{attention}(\text{key}_i \mid \text{query}_j)$$

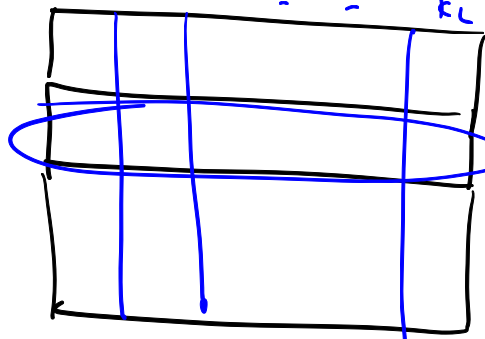


key "attend to"
query x_j

$$\vec{\alpha}_j = \text{Softmax}_{i=1 \dots L} \left(\vec{q}_j^T \vec{k}_i \right)$$

how
for
query j

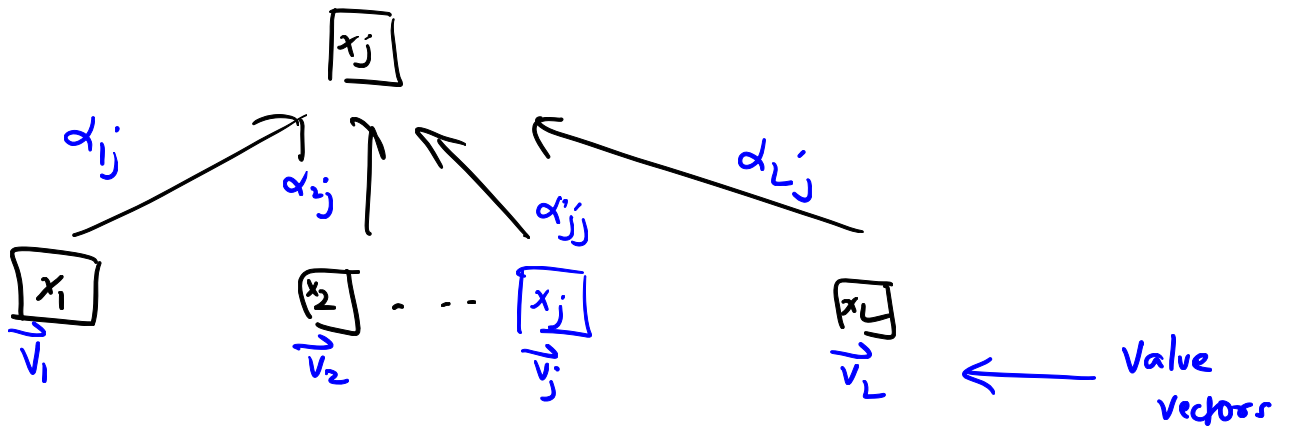
$\vec{\alpha}_j$



probability vector

A

Attention matrix



Output:

$$\vec{v}_j = \sum_{i=1}^L \alpha_{ji} \vec{v}_i$$

Weighted sum of value vectors
proportional to attention

