# Lecture 14:

## Language Models and Transformer

Markus Hohle

University California, Berkeley

**Bayesian Data Analysis and Machine Learning for Physical Sciences**
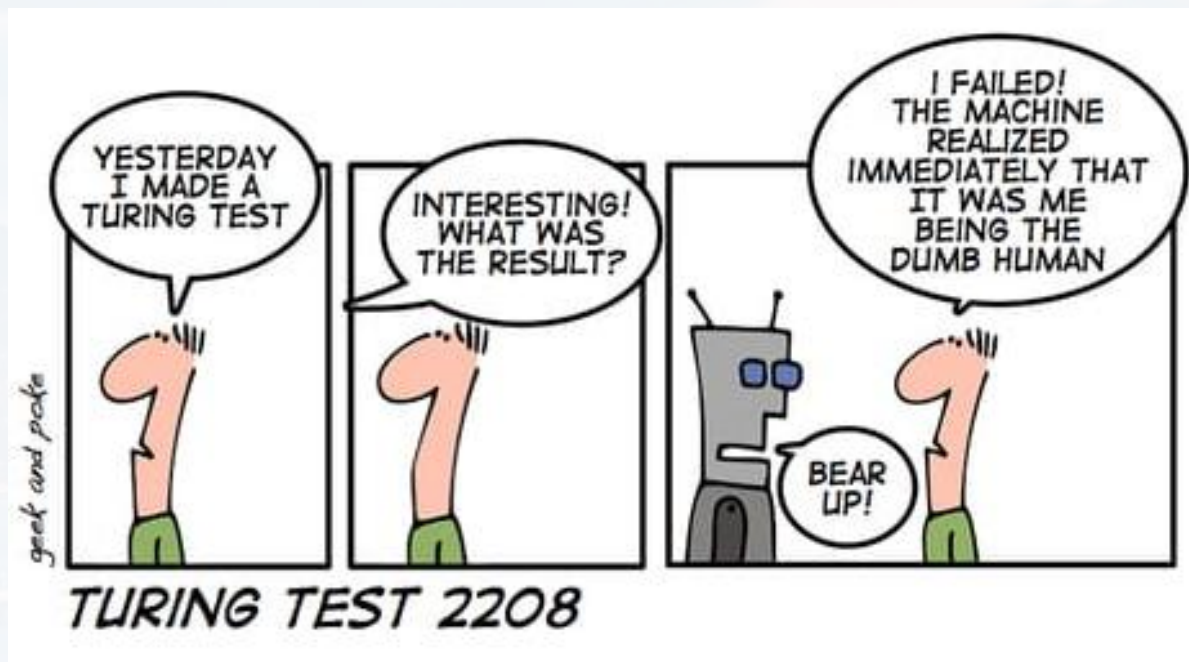
# Bayesian Data Analysis and Machine Learning for Physical Sciences

Berkeley LLM & Transformer

## Outline

- **Introduction**

- **Bigram and MAP**

- **Positional Encoding**

- **Word Embedding**

- **Attention**

- **Transformer Architecture**

# Outline

**corpus:**               (large, representative) data set containing sequences of a language

**token:**                 individual, independent entity of a language

**alphabet/vocabulary:**      set of tokens

| token | size of *alphabet* |
|---|---|
| - letters in a word | - $10^2$ |
| - words in a sentence | - $10^4 \dots 10^6$ |
| **(upper/lower case, cases, gender, tenses, conjugations)** | |
| - amino acids in a protein sequence | - 21 |
| - nucleotides in a DNA/RNA sequence | - 4 |
| - motifs in a DNA/RNA sequence | - $10^4$ |

**corpus:**                     (large, representative) data set containing sequences of a language

**token:**                     individual, independent entity of a language

**alphabet/vocabulary:**       set of tokens

**tokenization**



**token:**    - single atom vs…
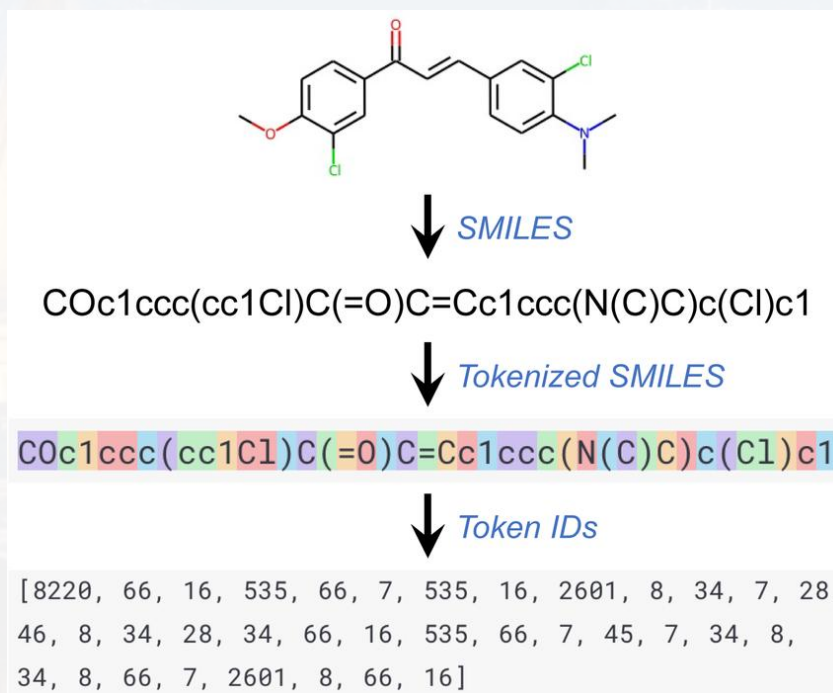
            - …functional group

**corpus:** (large, representative) data set containing sequences of a language

**token:** individual, independent entity of a language

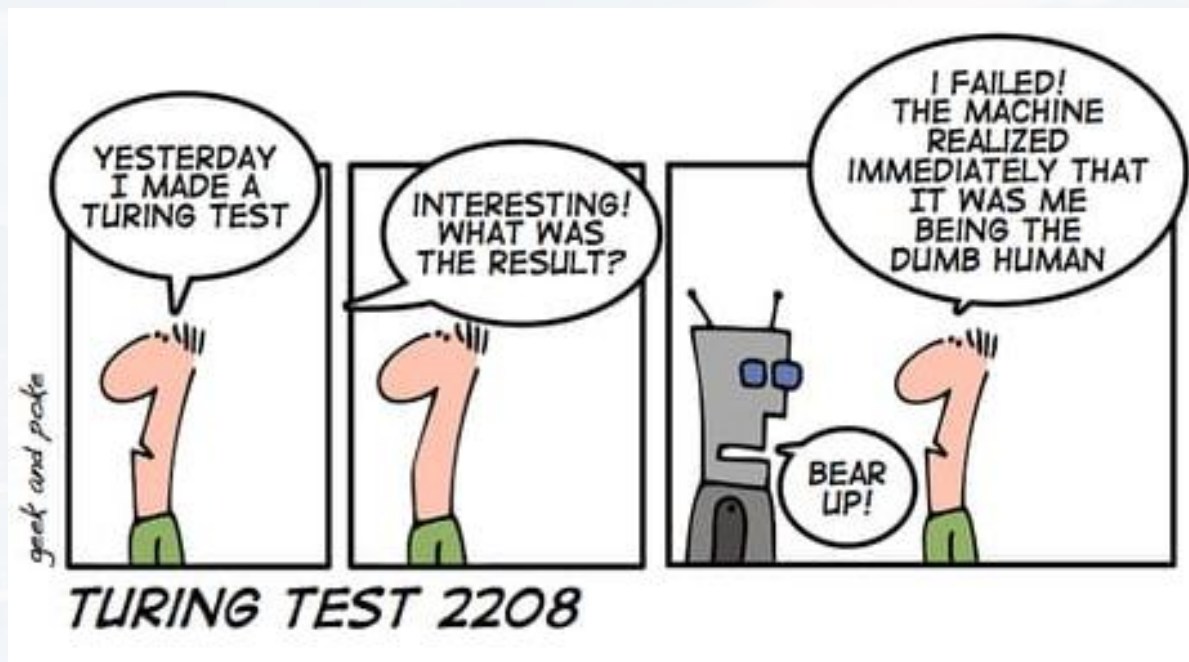**alphabet/vocabulary:** set of tokens

**note:** language models don't know grammar as we do, but they don't need to anyway...

**three things make context (**details: see later**):**

- **word embedding** (relation between similar/different token)

- **positional encoding** (location of token in a sequence)

- **attention** (relation between token within a sequence)

Outline

$$X_1\ X_2\ X_3\ X_4\ X_5\ ....\ X_n \qquad\qquad \text{sequence of } n \textbf{ token } X$$

actually:

$$P(X_1\ X_2\ X_3\ X_4\ X_5\ ....\ X_n) = P(X_n|X_{n-1}\ ...\ X_1)P(X_{n-1}|X_{n-2}\ ...\ X_1)\ ...\ P(X_1)$$
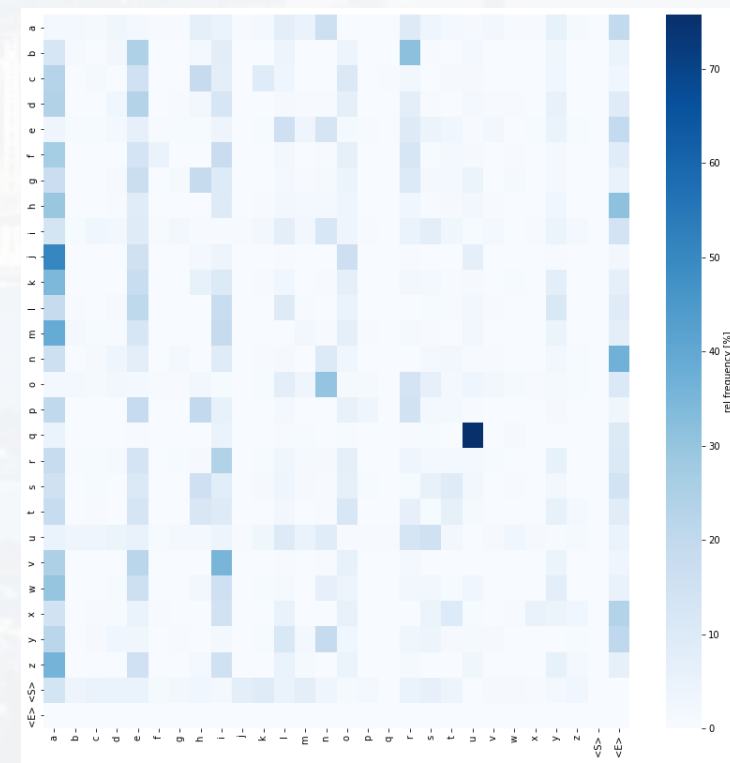
**bigram** (1st order Markov Chain, see e.g. first WhatsApp versions):

$$P(X_1\ X_2\ X_3\ X_4\ X_5\ ....\ X_n) = P(X_n|X_{n-1})P(X_{n-1}|X_{n-2})\ ...\ P(X_1)$$

**P(i|j)**: that token **i** is generated after token **j**

    → N x N transition matrix from frequencies

    → **"bigram" = "two words"**



**frequency matrix of letters in common names**

**bigram** (1st order Markov Chain):

let's build our own bigram model: **generate new names** based on a corpus of names

```
In [15]: words[0:12]
Out[15]:
['emma',
 'olivia',
 'ava',
 'isabella',
 'sophia',
 'charlotte',
 'mia',
 'amelia',
 'harper',
 'evelyn',
 'abigail',
 'emily']
```

see **Andrej Karpathy's GitHub** repository

$$P(X_1\ X_2\ X_3\ X_4\ X_5\ ....X_n) = \boxed{P(X_n|X_{n-1})P(X_{n-1}|X_{n-2})} ...\boxed{P(X_1)}$$

we only need to count **how often** a letter is followed by another

we also need to indicate when a name has **started** and **ended**

$$['<S>'] + ['olivia'] + ['<E>']$$

→ **alphabet:** 26 letters + the two special "letters"

let's create a dictionary first (will help for counting):

**bigram** (1st order Markov Chain):

let's build our own bigram model: **generate new names** based on a corpus of names
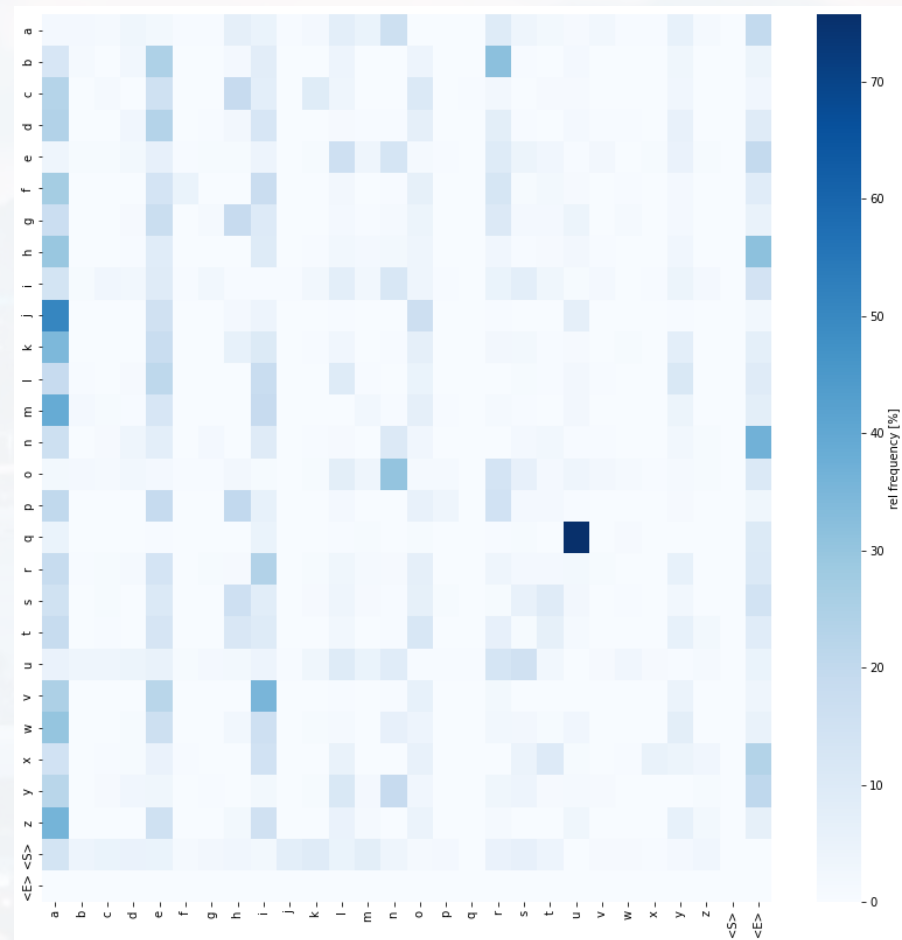
1) dictionary first (will help for counting)

2) count how often letter *i* is followed by letter *j*

        → bigram matrix **N**

3) normalize **N** accordingly

4) begin with a start token

5) draw a letter randomly based on **N**, using

        `np.argmax(np.random.multinomial(1,p))`

6) if next token is stop token → stop

**bigram** (1st order Markov Chain):

let's build our own bigram model: **generate new names** based on a corpus of names

check out **Bigram.ipynb**

**B.SampleNames(15)**     vs totally random     **B.SampleNames(15, False)**

some names are gibberish

some names sound real

some names **are** real

```
In [295]: B.SampleNames(15)
keesa
ann
ja
jon
nma
malynojana
sall
daha
drvah
lzaxi
tyunusthun
jorrwro
ja
asoow
s
```

```
In [296]: B.SampleNames(15,False)

mtkgy
yufexhviovmorhqvikbbbjxebpxwurejaqlzzuwuanxmmomhr<S>uhb
xlmusadjfdzxadaotd
ik<S>vdtydvxevtaselkykcfbamceprtvl
zyr<S>inzoerobzwuovx
eg<S>pbdvikf<S>tomcnkfsjay<S>rikatnaykizszcivpds<S>zj
kh<S>y<S>ualzugqgakakeubjbasc
bblupnibtqmyl<S>vyobf
kybs
rznjgpmlo
tnhoxuckkjjzbwmj<S>vshkycicf<S>kowskphy
rxodh
jvswmzw
jzpcfnpcg
```

Note, there is no conceptual difference between applying our model to *letters in a word* vs *words in a sentence*

caveats: - the bigram model derives $P(X_n)$ from **observed** frequencies

→ essentially **MLE** (problematic if a letter hasn't appeared in the sequence yet

→ $P(X_n)$ assumed to be zero!)

```
Nsam = N/np.sum(N+0.0001, axis = 1, keepdims = True)
S_bi += np.sum(-N[:,i]*np.log(N[:,i]+1e-16))
```

- can we implement something that is closer to:

$$P(X_1\ X_2\ X_3\ X_4\ X_5\ ....X_n) = P(X_n|X_{n-1}...X_1)P(X_{n-1}|X_{n-2}...X_1)...P(X_1)\ \ ?$$
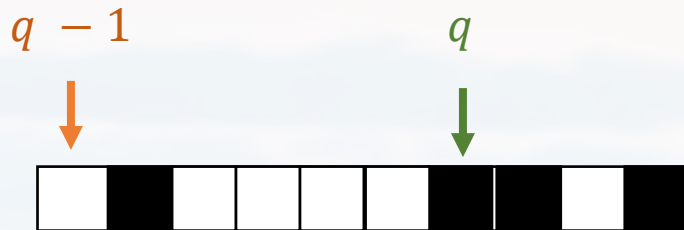
binomial process

$$P(k|n,q) = \binom{n}{k} q^k (1-q)^{n-k}$$

$q - 1$      $q$

$q = ?$

$$P(X_1 X_2 X_3 X_4 X_5 \ldots X_n) = P(X_n|X_{n-1} \ldots X_1)P(X_{n-1}|X_{n-2} \ldots X_1) \ldots P(X_1)$$

$q - 1$      $q$

$$P(k|n, q) = \binom{n}{k} q^k (1-q)^{n-k} \qquad q = ?$$

likelihood function **(here: binomial)**

$$P(q|data\ set) = \frac{\boxed{P(data\ set|q)}\ \boxed{P(q)}\ \text{prior}\ (\ \sim P(X_n|X_{n-1} \ldots X_1)\ )}{\boxed{P(data\ set)}\ \text{evidence}\ \textbf{(const wrt q)}}$$

$$= \frac{1}{\int_0^1 P(q|data\ set)dq}(1-q)^{n-k}q^k$$

$q = const$
before 1$^{st}$ data point
(max entropy!)

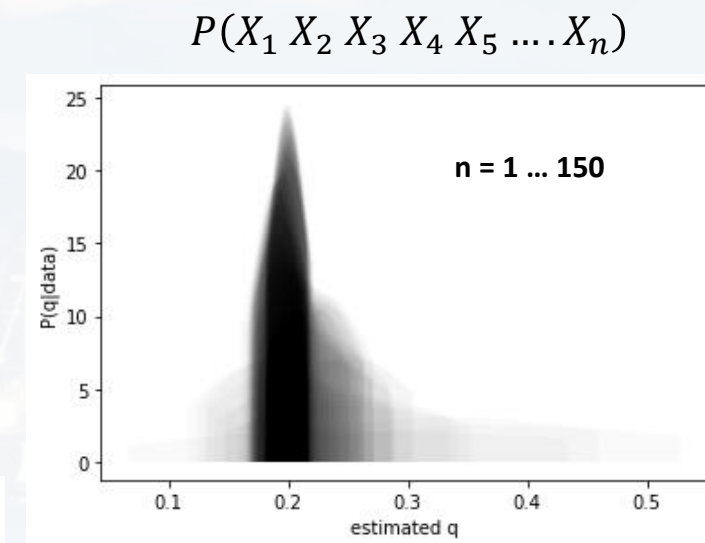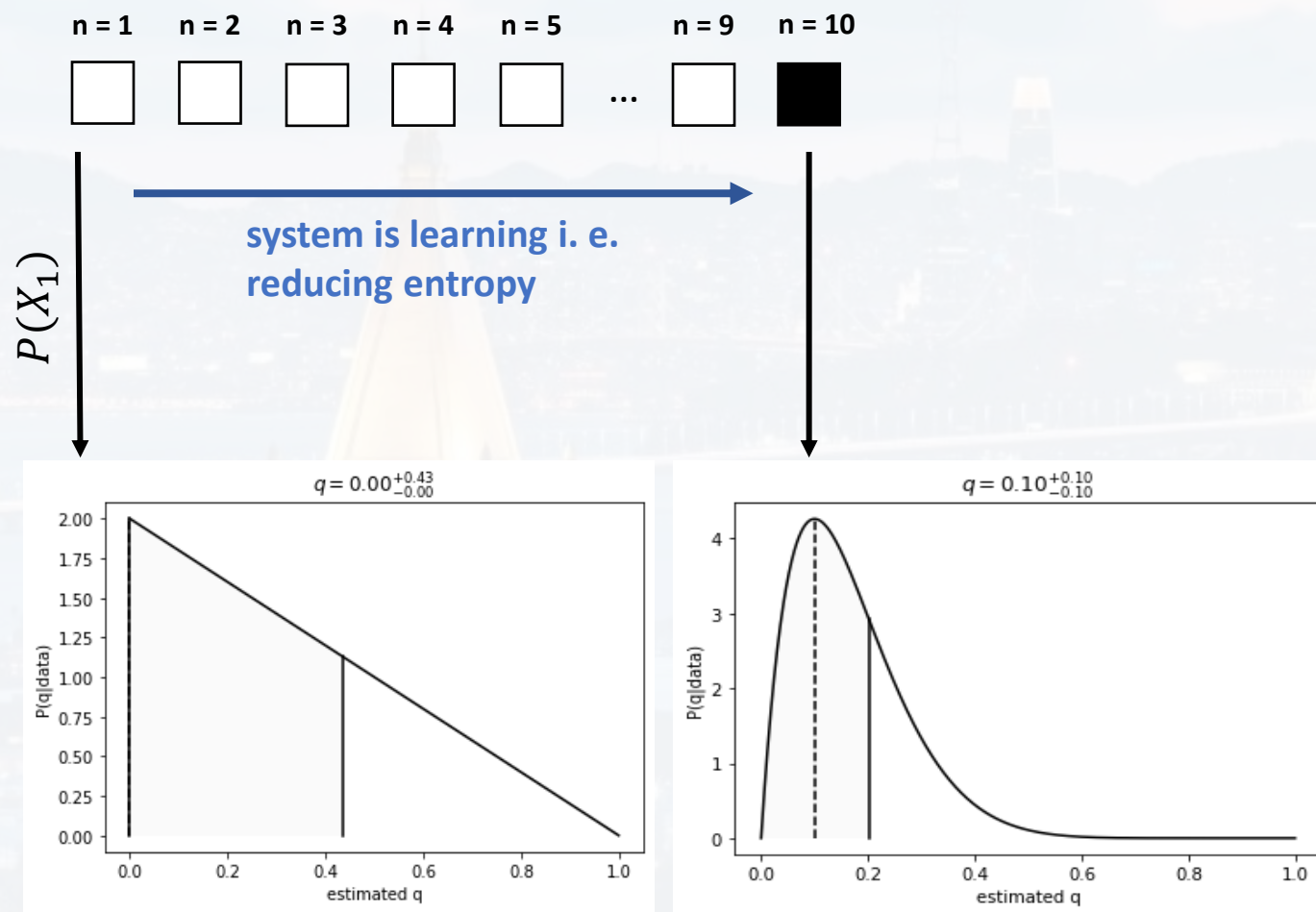$$= \frac{q^{k+\alpha-1}(1-q)^{n-k+\beta-1}}{\int_0^1 q^{k+\alpha-1}(1-q)^{n-k+\beta-1}\ dq}$$

$q = $ conjugate prior
after n$^{th}$ data point

**\*) see lecture 2**

$$P(X_1\, X_2\, X_3\, X_4\, X_5\, \ldots\, X_n) = P(X_n|X_{n-1}\ldots X_1)P(X_{n-1}|X_{n-2}\ldots X_1)\ldots P(X_1)$$

$$P(k|n,q) = \binom{n}{k} q^k (1-q)^{n-k} \qquad\qquad q = ?$$

**B**ayesian **P**arameter **E**stimation*)



n = 1    n = 2    n = 3    n = 4    n = 5      n = 9    n = 10

$P(X_1)$

**system is learning i. e. reducing entropy**

$P(X_1\, X_2\, X_3\, X_4\, X_5\, \ldots\, X_n)$

**n = 1 … 150**

$q = 0.00^{+0.43}_{-0.00}$

$q = 0.10^{+0.10}_{-0.10}$

$$P(X_1\,X_2\,X_3\,X_4\,X_5\,....\,X_n) = P(X_n|X_{n-1}\,...\,X_1)P(X_{n-1}|X_{n-2}\,...\,X_1)\,...\,P(X_1)$$

$$P(k|n,q) = \binom{n}{k} q^k (1-q)^{n-k}$$

$$P(q|data\ set) = \frac{q^{k+\alpha-1}(1-q)^{n-k+\beta-1}}{\boxed{\int_0^1 q^{k+\alpha-1}(1-q)^{n-k+\beta-1}\,dq}}$$

Beta function



more general, we want to learn the probability $P_j(a)$ of letter $a$
at position $j$         $q \rightarrow P_j(a)$

→ multinomial problem

→ conjugate prior is the **Dirichlet distribution**

$$P(sequence) \sim \prod_j \boxed{\prod_a P(a)_j^{\alpha(a)-1}}$$

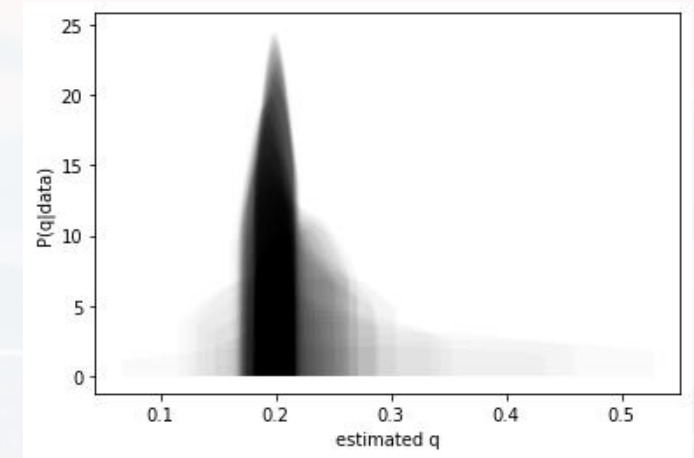equivalent to what was $P(q|data\ set)$
earlier

**\*) see lecture 2**

$$P(X_1\,X_2\,X_3\,X_4\,X_5\,....X_n) = P(X_n|X_{n-1}\,...X_1)P(X_{n-1}|X_{n-2}\,...X_1)\,...P(X_1)$$

**Dirichlet distribution**



$$P(sequence) \sim \prod_j \prod_a P(a)_j^{\alpha(a)-1}$$

**note:**        $\displaystyle\sum_{over\ all\ a} P(a)_j = 1$          → N − dim simplex

**note:**     - we don't need to extract $P(a)$ from the maximum of the pdf given by the BPE posterior

             - we can directly derive the maximum of $P(a)$ from $P(sequence)$
               given the constrain $\sum_{over\ all\ a} P(a)_j = 1$ **(Lagrangian multipliers)**

             - **Maximum a-posteriori (MAP)** approach → see XXmotif (Siebert & Soeding, 2016)
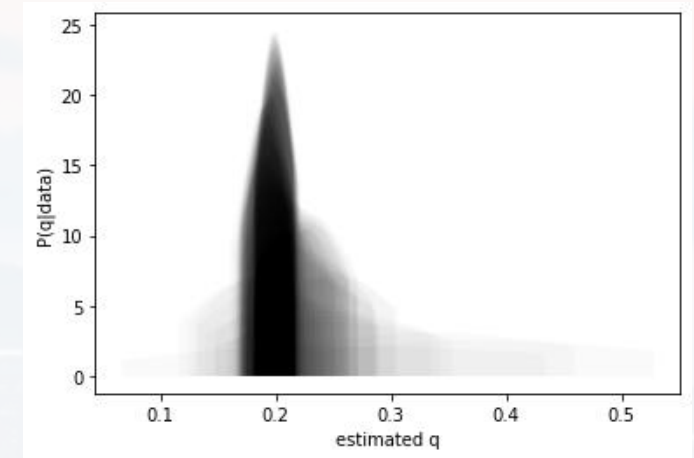
$$P(X_1\,X_2\,X_3\,X_4\,X_5\,....\,X_n) = P(X_n|X_{n-1}\,...X_1)P(X_{n-1}|X_{n-2}\,...X_1)\,...P(X_1)$$

**Dirichlet distribution**

$$P(sequence) \sim \prod_j \prod_a P(a)_j^{\alpha(a)-1}$$

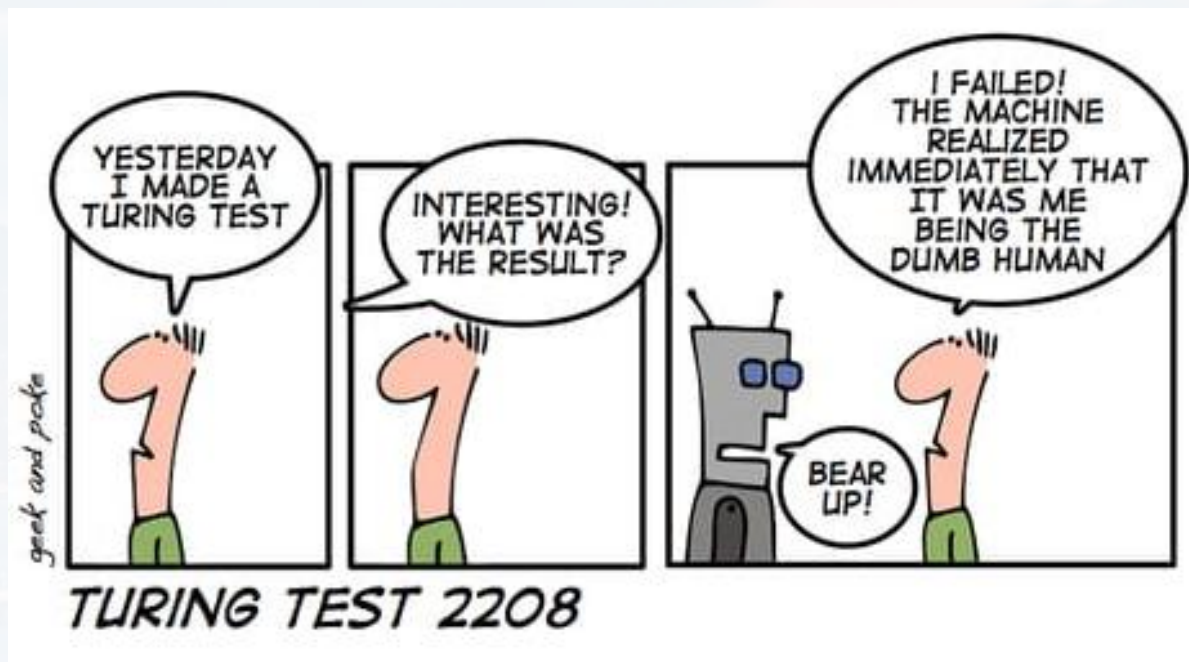<u>**note:**</u>    $\displaystyle\sum_{over\ all\ a} P(a)_j = 1$     → N – dim simplex



**Maximum a-posteriori (MAP)**

- XXmotif (Siebert & Soeding, 2016) significantly outperformed PWMs

- it struggled however with related motifs which where **physically located far apart** from each other

→ solution see later: attention
→ older solutions: LSTMs

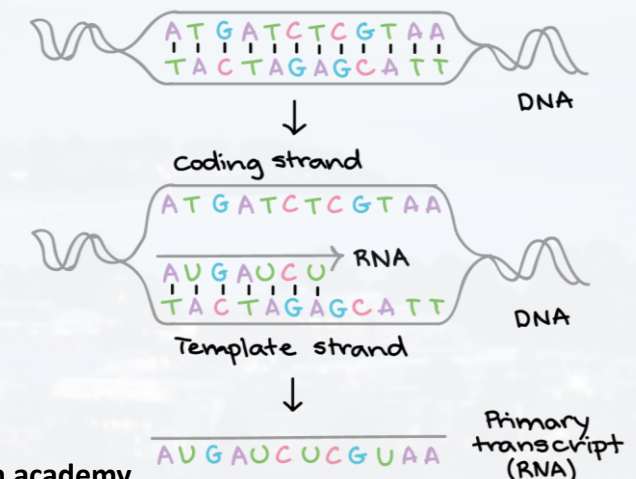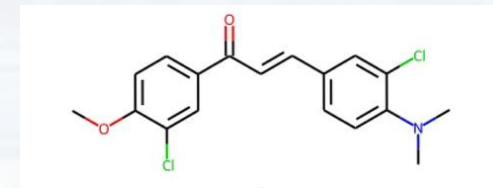## Outline

**three things make context:**

- **positional encoding** (location of token in a sequence)

- **word embedding** (relation between similar/different token)

- **attention** (relation between token within a sequence)

*"The cat jumped on the roof."*

order matters!:

1st:          article
2nd:        noun/subject
3rd:         verb
4th :       noun/object (in English)

→ **positional encoding**



AT GATCTCGT AA
TAC TAGAGCA TT
DNA
↓
Coding strand
AT GATCTCGT AA
AU GAUCU
TAC TAGAGCA TT
RNA
DNA
Template strand
↓
AU GAUCUCGUAA
Primary transcript (RNA)

**Khan academy**

**goal:** find a positional encoding that is

- reasonably simple

- independent from the length of the sequence

- somehow normalized

one idea: n-bit binary encoding

| position | code |
|---|---|
| | 76543210  ← 8bit i.e. eight dimensions |
| 1 | 00000001 |
| 2 | 00000010 |
| 3 | 00000011 |
| 4 | 00000100 |
| 5 | 00000101 |
| 6 | 00000110 |
| 7 | 00000111 |
| 8 | 00001000 |
| 9 | 00001001 |
| 10 | 00001010 |
| 11 | 00001011 |
| 12 | 00001100 |
| 13 | 00001101 |
| 14 | 00001110 |
| 15 | 00001111 |
| 16 | 00010000 |

depending on dimensions (bit)

→ different frequencies

| bit | frequency |
|---|---|
| 0 | 1/2 |
| 1 | 1/4 |
| 2 | 1/8 |
| ... | |

Does that look familiar?
→ *like* Fourier Series

even dimensions: $\quad E(p, 2k) \quad\quad = \sin\left(\dfrac{p}{10,000^{\frac{2k}{d}}}\right)$

odd dimensions: $\quad E(p, 2k+1) = \cos\left(\dfrac{p}{10,000^{2k/d}}\right)$

$p$:        position in sequence
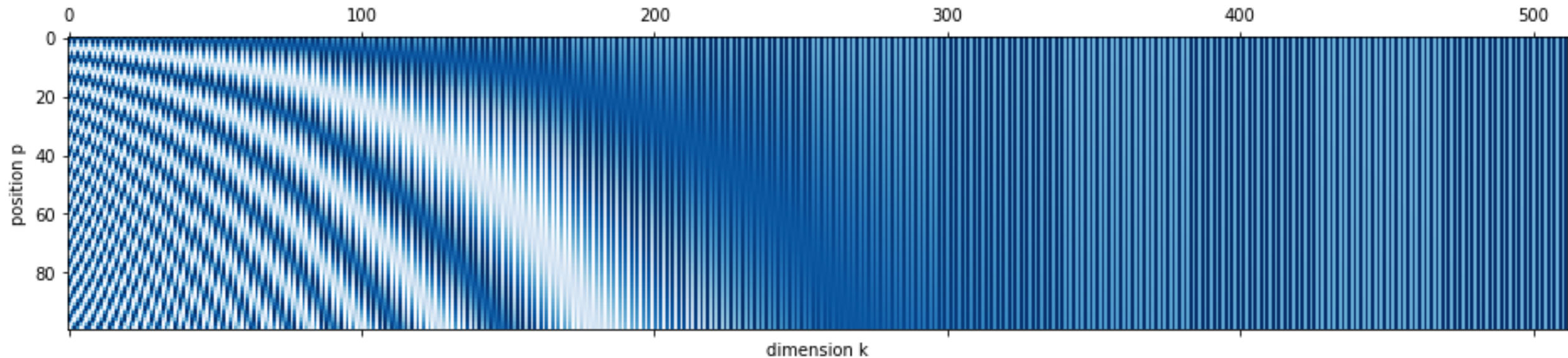$k$:        dimension index
$d$:        number of dimensions
10,000:    an arbitrary number     **Vaswani et al., 2017**

run **PlotPositionEncoding.py**

more info here

even dimensions: $E(p, 2k) = \sin\left(\dfrac{p}{10{,}000^{\frac{2k}{d}}}\right)$

odd dimensions: $E(p, 2k+1) = \cos\left(\dfrac{p}{10{,}000^{2k/d}}\right)$

$p$: position in sequence
$k$: dimension index
$d$: number of dimensions
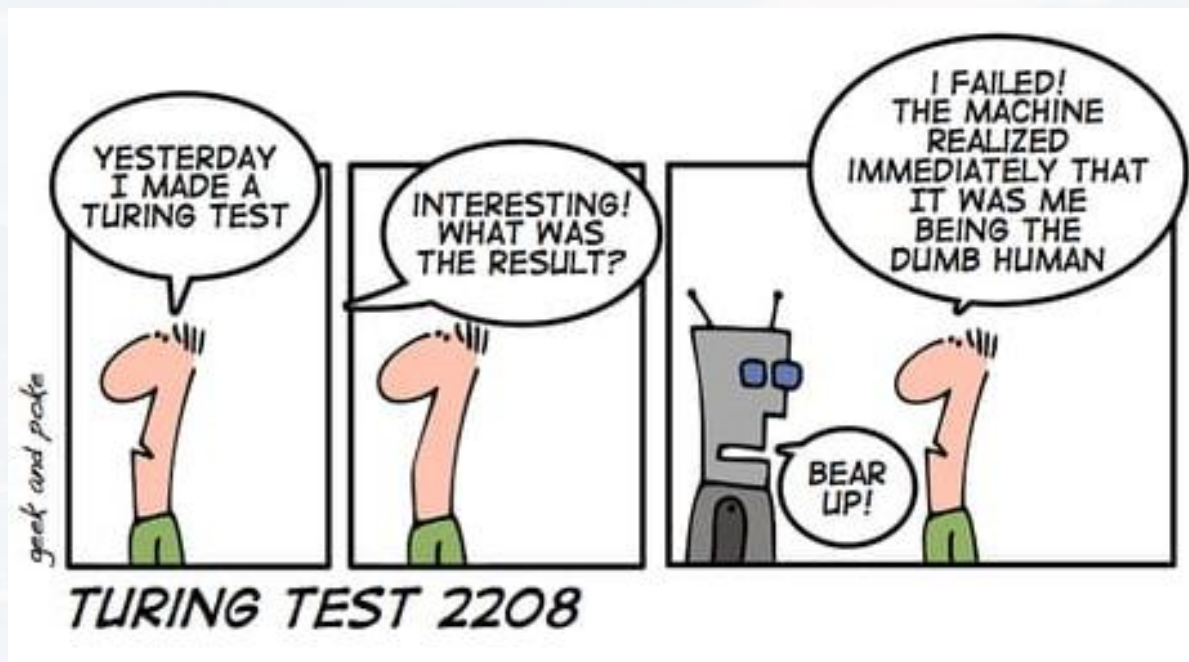10,000: an arbitrary number  **Vaswani et al., 2017**

run **PlotPositionEncoding.py**

**note:**

- easier to handle numerically vs discrete encoding

- $\cos(x + 2\pi k) = \cos(x)$ and $\sin(x + 2\pi k) = \sin(x)$
  → absolute position not relevant but *relative* position

# Outline

**three things make context:**

- **positional encoding** (location of token in a sequence)

- **word embedding** (relation between similar/different token)

- **attention** (relation between token within a sequence)

**problem: turning token (words/letters) into numbers**

single letters:

ACGT                - one – hot works perfectly (four different token)

abcd…            - lower/upper case, special characters (50 different token), one – hot is fine too

words:

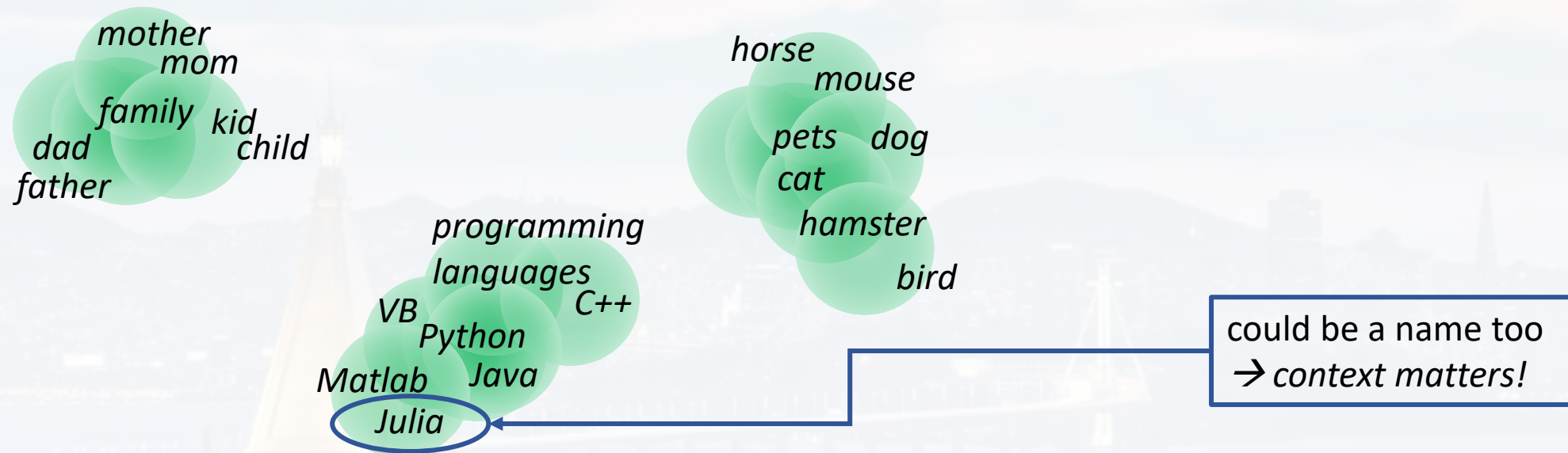actual words      - $10^4 \dots 10^6$ (upper/lower case, cases, gender, tenses, conjugations)

                                  → one – hot doesn't work (matrices would be too large)

                                  → some words have a **similar meaning,** should be **close** in parameter space (**cluster**)
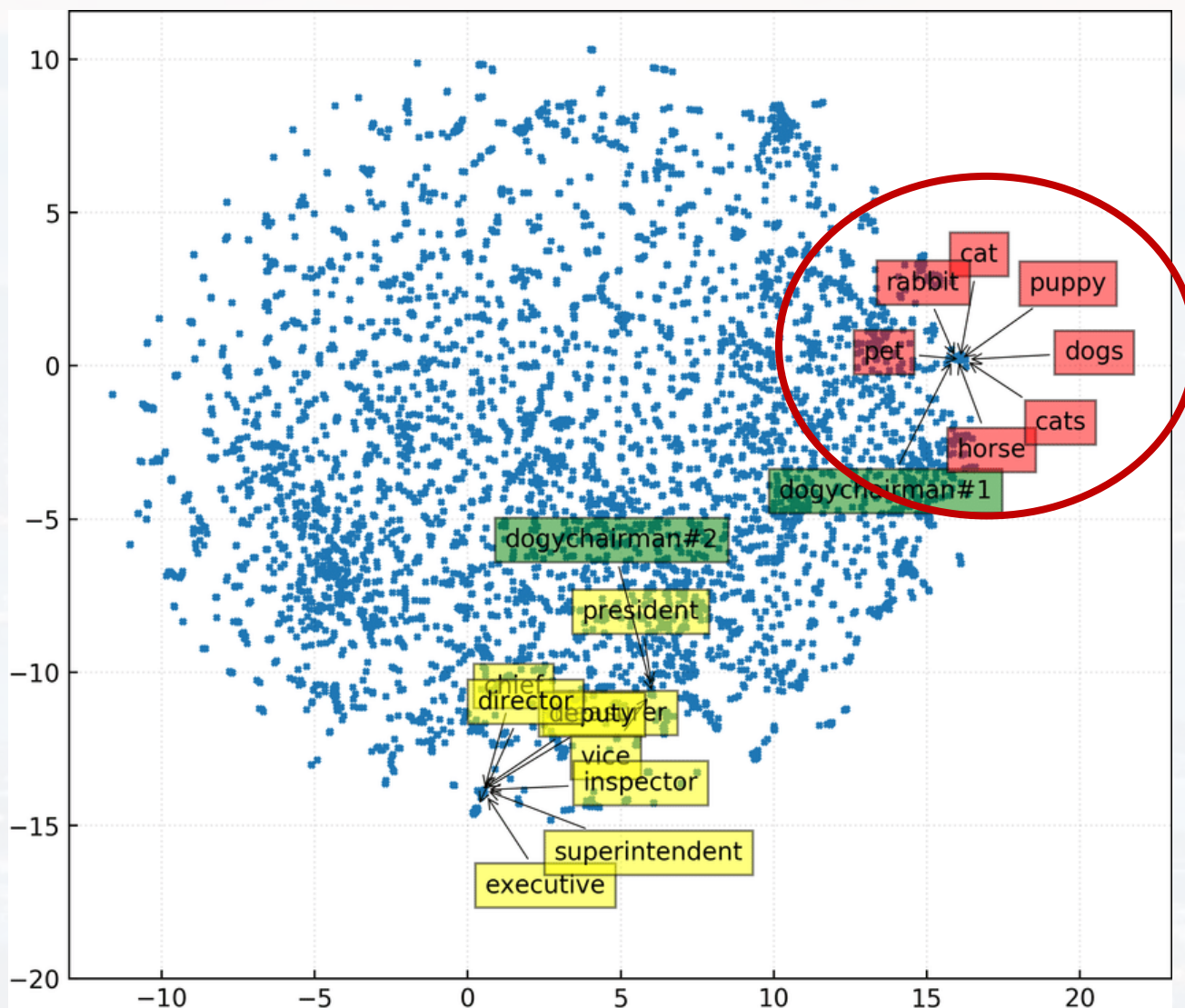
words with a **similar meaning** should form **cluster**

mother
mom
family  kid
dad      child
father

horse
    mouse
pets   dog
cat
hamster
    bird

programming
languages
VB        C++
Python
Matlab  Java
Julia

could be a name too
→ *context matters!*

- embedding, instead of one – hot encoding
- from experience**: N = 30 – 300** dim vector for each token **(which is a lot less than $10^4$ … $10^6$)** is sufficient
- as a result: token with **similar meaning are close** in the vector space!

words with a **similar meaning** should form **cluster**



*"Joint Learning of Sense and Word Embeddings"*
M Alsuhaibani & D Bollegala

**common training set: recorded speeches from the European Parliament:**

*...It seems absolutely disgraceful that we pass legislation and do not adhere to it ourselves. Mrs Lynne, you are quite right and I shall check whether this has actually not been done. I shall also refer the matter to the College of Quaestors, and I am certain that they will be keen to ensure that we comply with the regulations we ourselves vote on.*
*Madam President, Mrs Díez González and I had tabled questions on certain opinions of the Vice-President, Mrs de Palacio, which appeared in a Spanish newspaper.*
*The competent services have not included them in the agenda on the grounds that they had been answered in a previous part-session.*
*I would ask that they reconsider, since this is not the case....*

words of similar meaning should appear in similar environment

→ *target token within a window*

Two common algorithms are **C**ontinuous **B**ag **O**f **W**ords and ***skip gram***

*__C__ontinuous __B__ag __O__f __W__ords*

| | |
|---|---|
| *n*: | **number of unique token from corpus** |
| *m*: | **desired number of dimensions for embedding** |

*The competent services have not included them in the agenda on the grounds that they...*

**target** | *have*

**context window** | *The competent services     not included them*

one – hot encoding

one – hot encoding

**shallow ANN**

input layer of shape *n x m*

hidden layer of shape *m x n*

softmax

output vector of length *n*

target vector of length *n*

learning weights such that both vectors **are identical**
→ **weights[:,i]** of the **hidden layer** make up the embedding vector of ***token i***

*C*ontinuous *B*ag *O*f *W*ords

| | |
|---|---|
| $n$: | **number of unique token from corpus** |
| $m$: | **desired number of dimensions for embedding** |

*The* competent services have not included them in *the agenda on the grounds that they...*

target      *not*

context window      *competent services have     included them in*

one – hot encoding                   one – hot encoding

shallow ANN

input layer of shape *n x m*

hidden layer of shape *m x n*

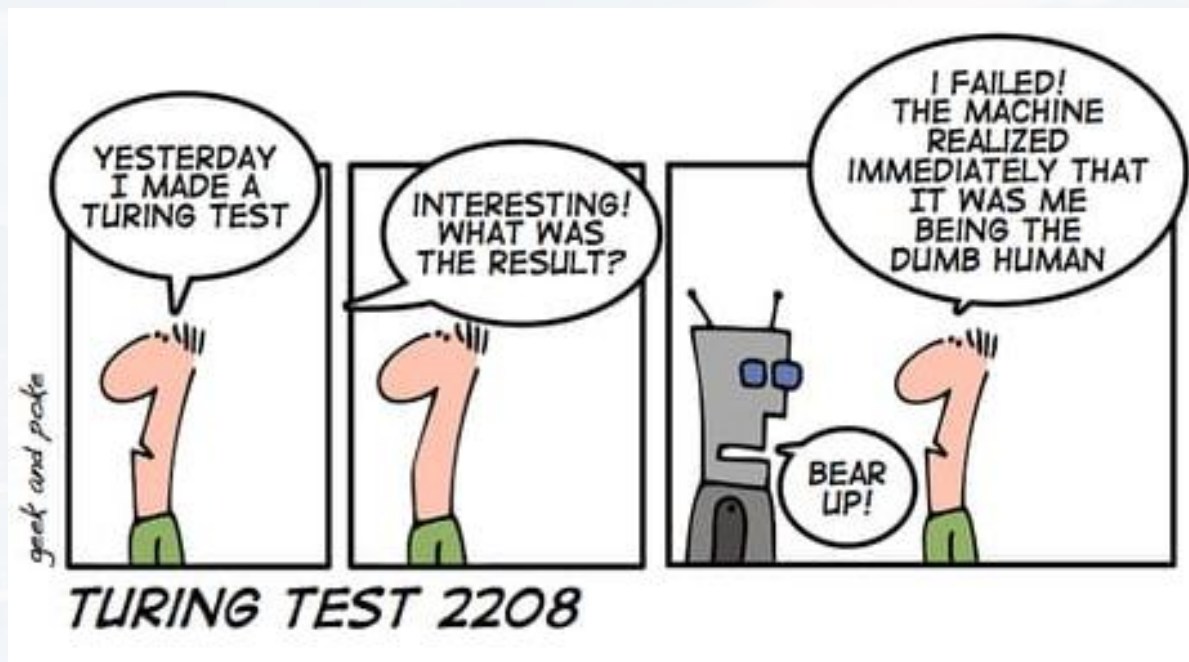softmax

output vector of length *n*

target vector of length *n*

learning weights such that both vectors **are identical**
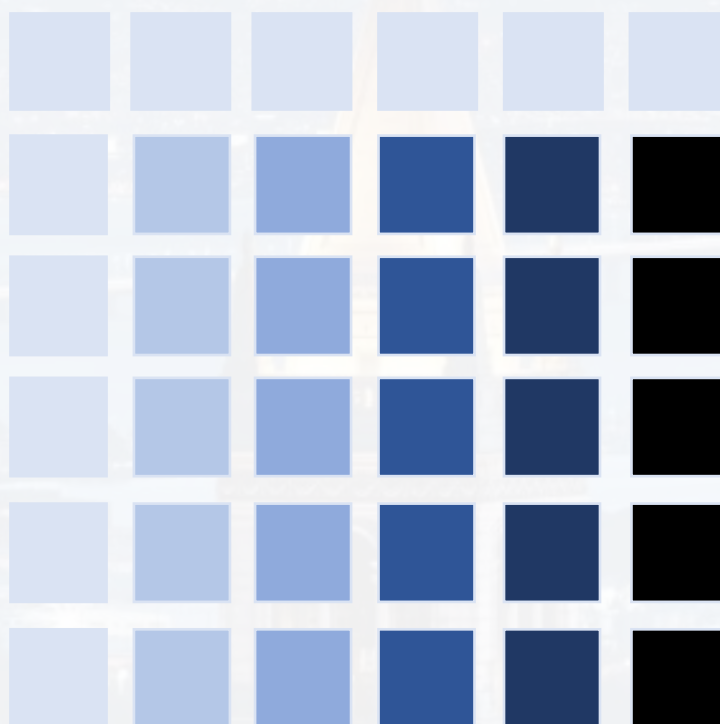→ **weights[:,i]** of the **hidden layer** make up the embedding vector of **token i**

*Continuous Bag Of Words*

| | |
|---|---|
| $n$: | **number of unique token from corpus** |
| $m$: | **desired number of dimensions for embedding** |

*The competent services have not included them in the agenda on the grounds that they...*

target — *included*

context window — *services have not   them in the*

one – hot encoding

... and so on....

one – hot encoding

shallow ANN

input layer of shape $n \times m$

hidden layer of shape $m \times n$

softmax

output vector of length $n$

target vector of length $n$

learning weights such that both vectors **are identical**
→ **weights[:,i]** of the **hidden layer** make up the embedding vector of **token i**

## Outline

**three things make context:**

- **positional encoding** (location of token in a sequence)

- **word embedding** (relation between similar/different token)

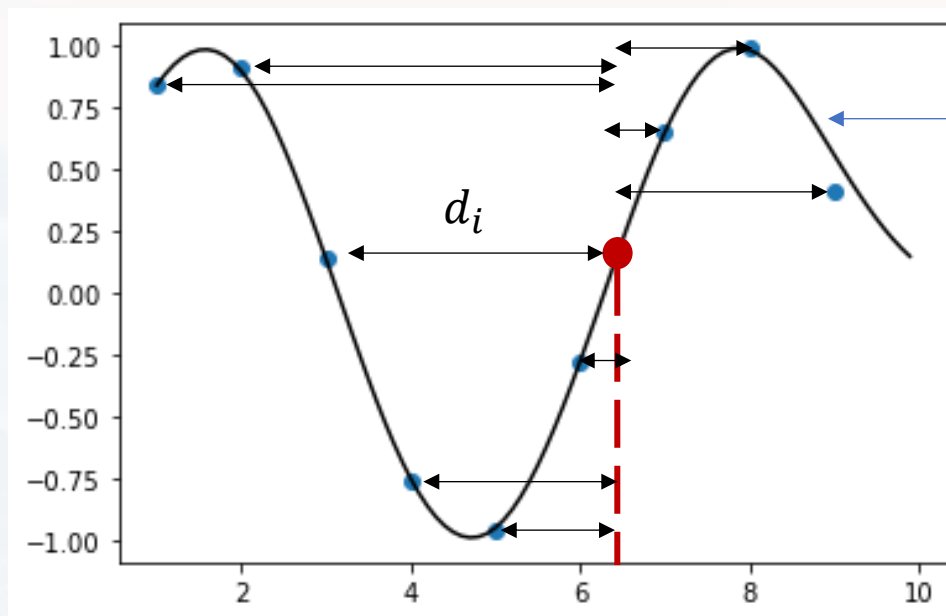- **attention** (relation between token within a sequence)

*"The cat jumped on the roof."*

how the first token influences all other token

how the second token influences all other token

…. and so on

We want to interpolate between the blue dots
→ generating the black line
→ **no curve fitting!**

idea:  - select a point for which we want the interpolation for

- calculate distance $d_i$ to every other point

- each data point should influence the value of the interpolated point

- the closer, the stronger the influence → weighted mean

$$y_{int} \sim \sum_{i=1}^{I} w_i \, y_i \qquad w_i \sim \frac{1}{d_i}$$

calculating distance



```
D = np.tile(V, (1, len(V))) - np.tile(L.transpose(), (len(V), 1))
```

- each data point should influence the value of the interpolated point
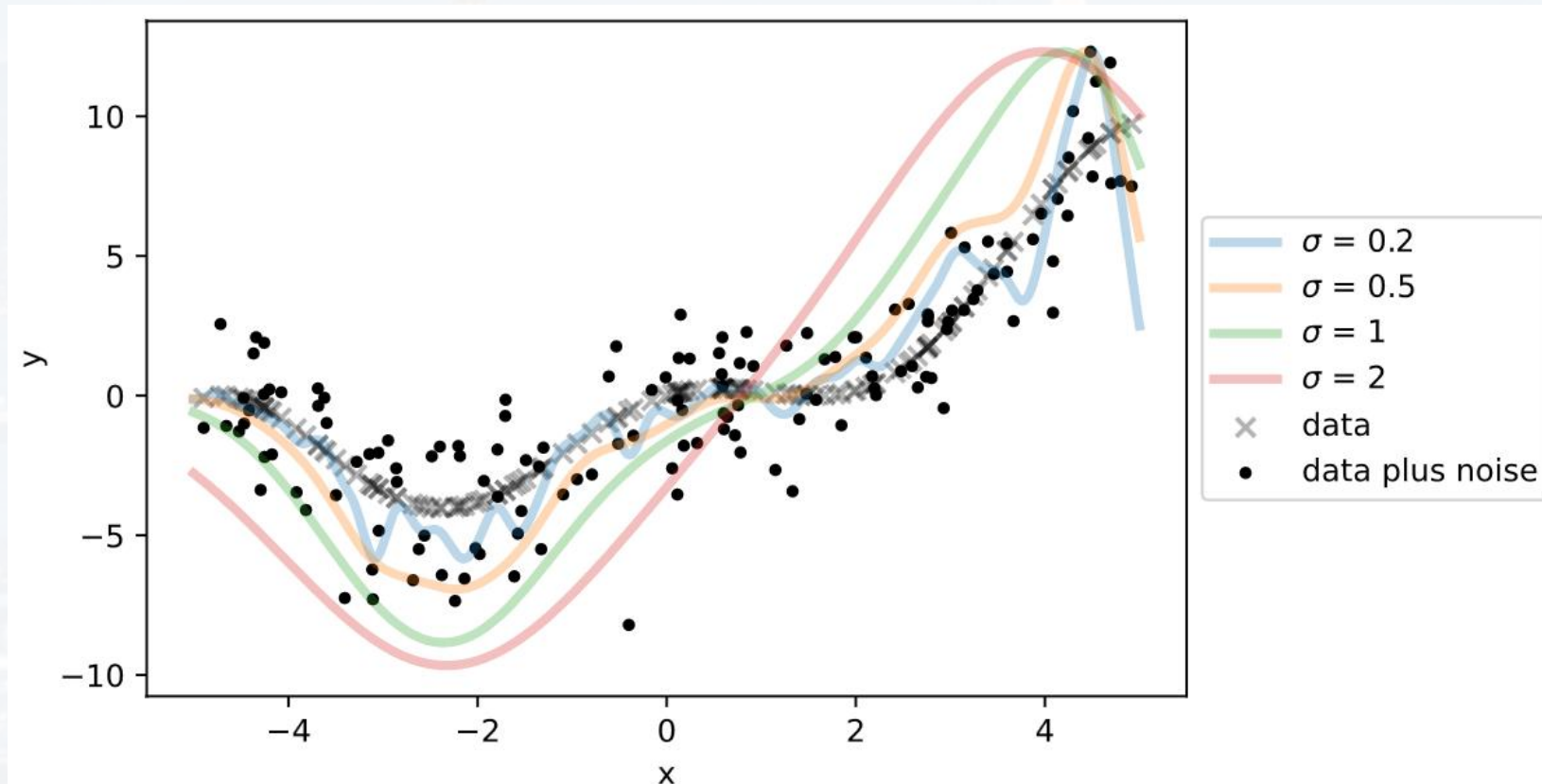- the closer, the stronger the influence → weighted mean

$$y_{int} \sim \sum_{i=1}^{I} w_i \, y_i$$

Gaussian kernel

```
W       = np.exp(-(D**2)/(sigma))
W       = W/np.sum(W + 1e-16, axis = 0)
yint    = np.dot(W.transpose(), y)
```

```
D = np.tile(V, (1, len(V))) -  np.tile(L.transpose(), (len(V), 1))
```

Gaussian kernel

```
W        = np.exp(-(D**2)/(sigma))
W        = W/np.sum(W + 1e-16, axis = 0)
yint     = np.dot(W.transpose(), y)
```
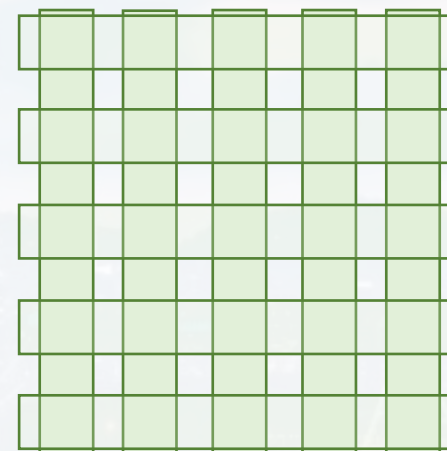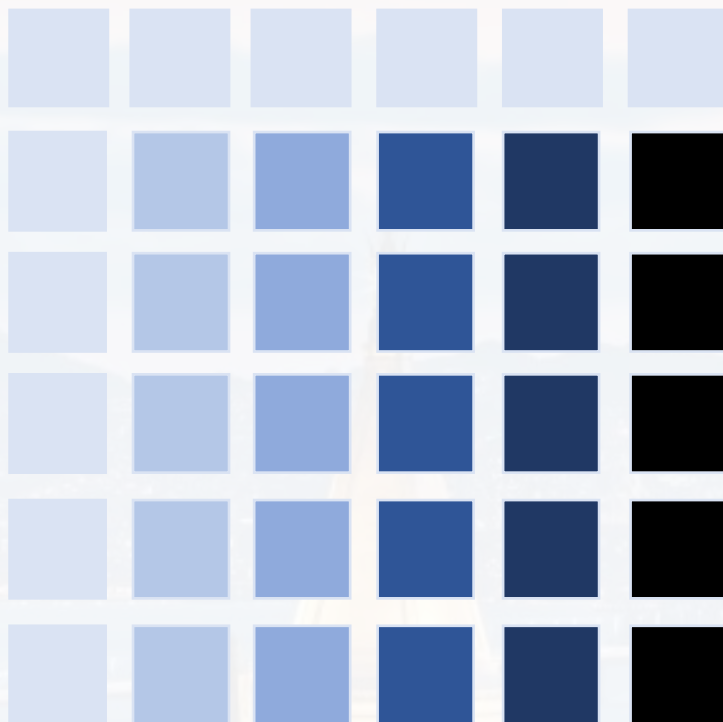
**check out:**

SmoothGaussKernel.py
SmoothExamples.py

*"The cat jumped on the roof."*



Gaussian kernel

```
W         = np.exp(-(D**2)/(sigma))
W         = W/np.sum(W + 1e-16, axis = 0)
yint      = np.dot(W.transpose(), y)
```

**actual attention:**
**these weights are learnable,**
no kernel assumed!

**self attention**

imagine you want to built & train a movie GenAI **that creates movies based on queries**.

**training data**

*"Thrilling horror science fiction movie, plays in space in a distance future"*

$$X = X_1, X_2, \ldots . X_N$$

each token is a vector $X_n$ of length $E$

Alien

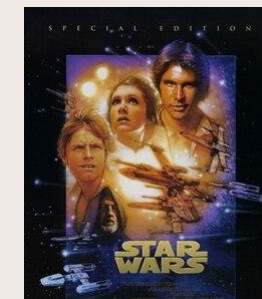*"Entertaining science fiction movie, plays in space in a distance past"*

Star Wars

*"Boring love story, plays on a ship in the past"*

Titanic

**key**

**value**

**self attention**

imagine you want to built & train a movie GenAI **that creates movies based on queries**.

**training data**

each token is a vector $X_n$ of length $E$          $X = X_1, X_2, \ldots . X_N$

Alien

**key**

**value**

*"Boring love story, plays in space."*

**query**

**self attention**        <span style="color:green">**key**</span>    <span style="color:orange">**value**</span>    <span style="color:blue">**query**</span>

each token is a vector $X_n$ of length $E$          $X = X_1, X_2, \ldots . X_N$

*"Boring love story, plays on a ship in the past"*

output:

$$Y_n = \sum_{m=1}^{N} w_{nm} X_m$$

← weights $w_{nm}$

- weights should be learnable

- weight = 0, $X_m$ has no influence on output

- weights should be positive so that neg weights

  don't counteract positive weights

- normalization: $\sum_{m=1}^{N} w_{nm} = 1$

for returning the best suggestion:

comparing key vector $\color{green}X_n$ to query vector $\color{blue}X_m$ via dot product

$$w_{nm} = \frac{\exp(X_n \circ X_m)}{\sum_{\mu=1}^{N} \exp(X_n \circ X_\mu)}$$          softmax

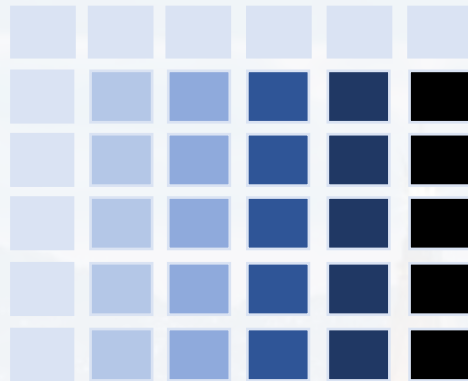**self attention**     **key**   **value**   **query**

each token is a vector $X_n$ of length $E$     $X = X_1, X_2, \ldots X_N$

*"Boring love story, plays on a ship in the past"*



weights $w_{nm}$

output:

$$Y_n = \sum_{m=1}^{N} w_{nm} X_m$$

-   weights should be learnable

-   weight = 0, $X_m$ has no influence on output

-   weights should be positive so that neg weights

   don't counteract positive weights

-   normalization:   $\sum_{m=1}^{N} w_{nm} = 1$

for returning the best suggestion:

comparing key vector $X_n$ to query vector $X_m$ via dot product

$$w_{nm} = \frac{\exp(X_n W(K) \circ X_m W(Q))}{\sum_{\mu=1}^{N} \exp(X_n W(K) \circ X_\mu W(Q))}$$

softmax
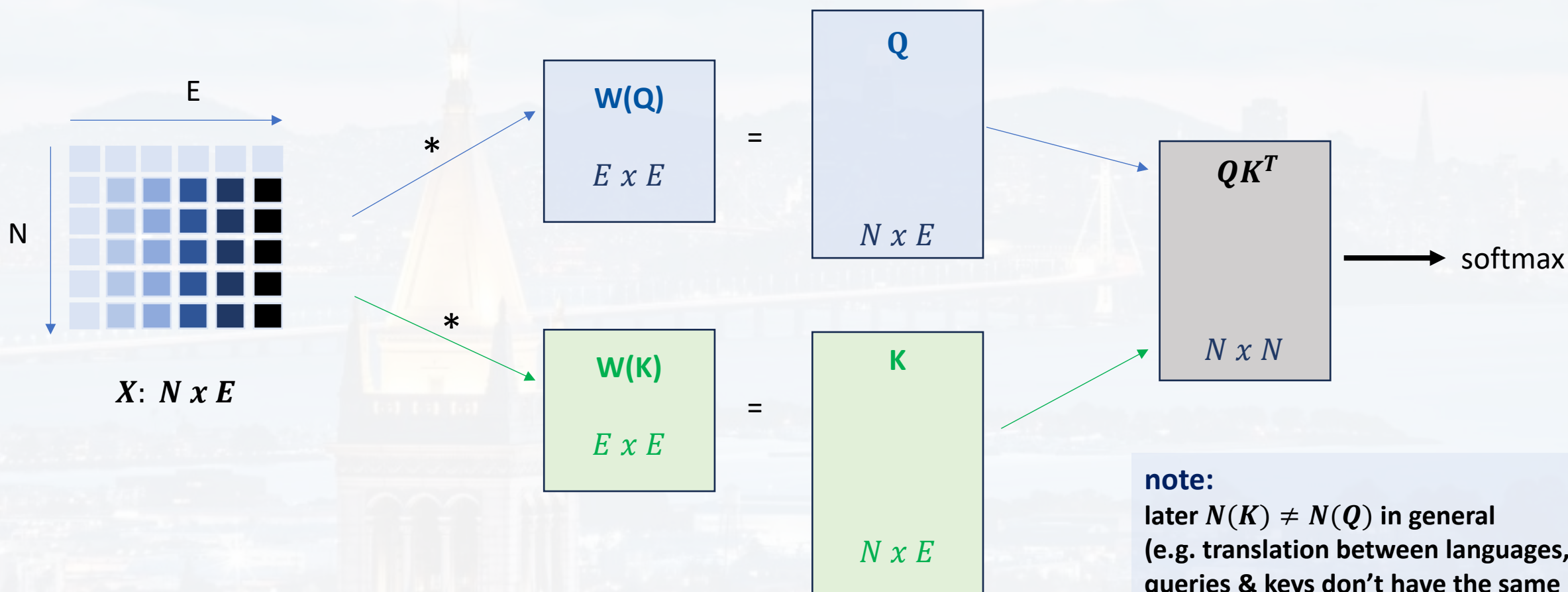
**self attention**   **key**   **value**   **query**

$$w_{nm} = \frac{\exp(X_n W(K) \circ X_m W(Q))}{\sum_{\mu=1}^{N} \exp(X_n W(K) \circ X_\mu W(Q))}$$

N: number of token
E: number of embedding dimensions

output:   $$Y_n = \sum_{m=1}^{N} w_{nm} X_m$$



E

N

*"Boring love story, plays on a ship in the past"*

$X$: $N \times E$

\* **W(Q)** $E \times E$ = **Q** $N \times E$

\* **W(K)** $E \times E$ = **K** $N \times E$

$QK^T$ $N \times N$

softmax

**note:**
later $N(K) \neq N(Q)$ in general
(e.g. translation between languages,
queries & keys don't have the same
number of tokens in general etc)

**self attention**     **key**    **value**    **query**

N: number of token
E: number of embedding dimensions

$$w_{nm} = \frac{\exp(X_n \mathsf{W(K)} \circ X_m \mathsf{W(Q)})}{\sum_{\mu=1}^{N} \exp(X_n \mathsf{W(K)} \circ X_\mu \mathsf{W(Q)})}$$

output: 
$$Y_n = \sum_{m=1}^{N} w_{nm} X_m = \sum_{m=1}^{N} \frac{\exp(X_n \mathsf{W(K)} \circ X_m \mathsf{W(Q)})}{\sum_{\mu=1}^{N} \exp(X_n \mathsf{W(K)} \circ X_\mu \mathsf{W(Q)})} X_m$$

$$= \sum_{m=1}^{N} softmax(\boldsymbol{Q}\boldsymbol{K}^T) X_m$$

$$\rightarrow \sum_{m=1}^{N} softmax(\boldsymbol{Q}\boldsymbol{K}^T) X_m \mathsf{W(V)} = \sum_{m=1}^{N} softmax(\boldsymbol{Q}\boldsymbol{K}^T) \mathsf{V} \qquad \textbf{value}$$

summarizing the characteristics of the movie by using the movie itself

The output would be a movie, generated by weighted contributions of those movies (values), where the keys match well with the query

**self attention**   **key**   **value**   **query**

N: number of token
E: number of embedding dimensions

$$w_{nm} = \frac{\exp(X_n W(K) \circ X_m W(Q))}{\sum_{\mu=1}^{N} \exp(X_n W(K) \circ X_\mu W(Q))}$$

output: $\quad Y_n = \sum_{m=1}^{N} w_{nm} X_m \quad \rightarrow \sum_{m=1}^{N} softmax(QK^T) X_m W(V) = \sum_{m=1}^{N} softmax(QK^T) V$

**note:**

- dot product scales with $E$, therefore normalization: $softmax\left(\frac{QK^T}{\sqrt{E}}\right)$

- cross attention:   - eg. key a phrase in language A, query in language B
  - **encoder/decoder** structure, see next slides

- we want to recognize many underlying pattern → **multiple attention** layers in parallel → transformer

## Outline

N:              number of token
E:              number of embedding dimensions

*"The cat jumped on the roof."*

N

positional encoding                    word embedding

E                                      E                                      E

N        +        N        =        N

transformer

within the transformer:

**attention:**

W(Q)                    Q

query

W(K)                    K

key

W(V)                    V

value

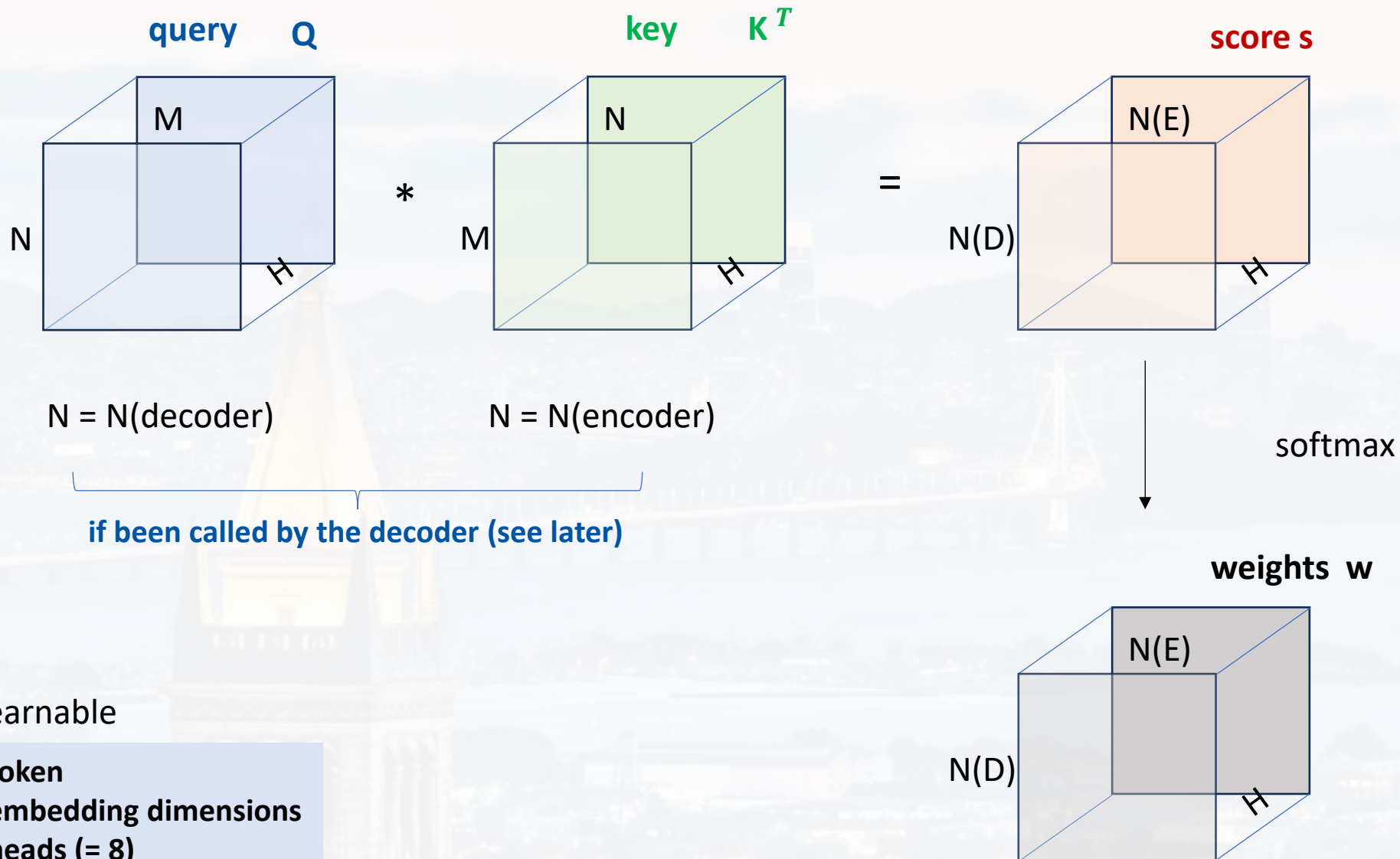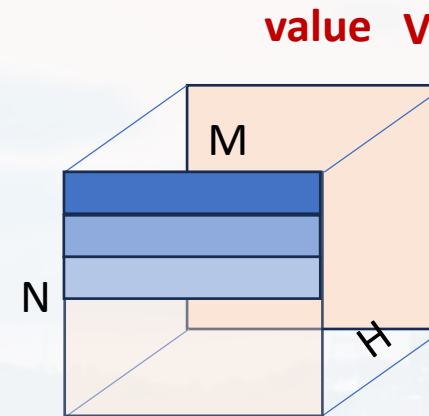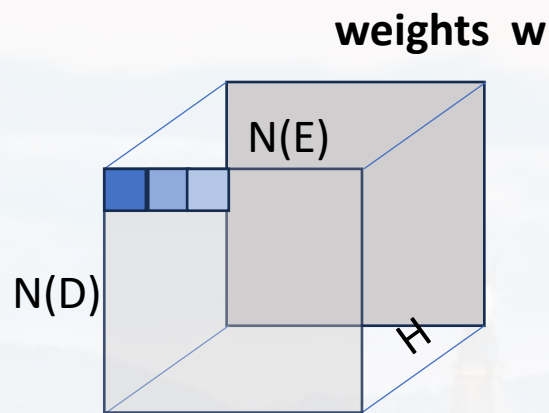W(Q), W(K), W(V): learnable

N:          number of token
E:          number of embedding dimensions
H:          number of heads (= 8)
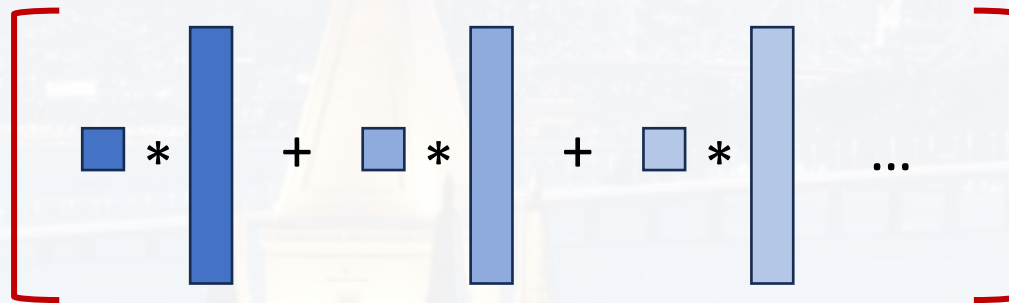M:          head size (= 64)

**within the transformer:**

**attention:**

query   **Q**          **key**    $K^T$          **score s**

M

N        H

\*     N     M      H     =    N(E)    N(D)    H

N = N(decoder)          N = N(encoder)

**if been called by the decoder (see later)**

softmax

**weights  w**

N(E)

N(D)    H

**W(Q), W(K), W(V)**: learnable

| | |
|---|---|
| **N:** | number of token |
| **E:** | number of embedding dimensions |
| **H:** | number of heads (= 8) |
| **M:** | head size (= 64) |

**within the transformer:**

| N: | number of token |
|---|---|
| E: | number of embedding dimensions |
| H: | number of heads (= 8) |
| M: | head size (= 64) |

**weights  w**

**value  V**

**attention:**

N(E)

M

N(D)

N

H

H

N(E)      = N(encoder)

N(D)      = N(decoder) , **see later**

$$\left[ \quad \blacksquare * \blacksquare \quad + \quad \blacksquare * \blacksquare \quad + \quad \blacksquare * \blacksquare \quad \ldots \quad \right]$$

σ

M

$$\sigma[i, :, k] = \sum_{j} w[i, j, k] * v[j, :, k]$$

N

H

**within the transformer:**

σ

M

N

H

| N: | number of token |
|----|-----------------|
| E: | number of embedding dimensions |
| H: | number of heads (= 8) |
| M: | head size (= 64) |

concatenating σ

$$\sigma_{con} = \boxed{\sigma[:,:,0]} \boxed{\sigma[:,:,1]} \ldots \boxed{\sigma[:,:,H]} \quad \downarrow N \quad *$$

$H \cdot M$

E

$H \cdot M$

learnable $\Omega$

**output:**

E

=

N

N = N(dencoder) ,
**if been called by the
decoder (see later)**

**that was attention → now: encoder**

https://jalammar.github.io/illustrated-transformer/



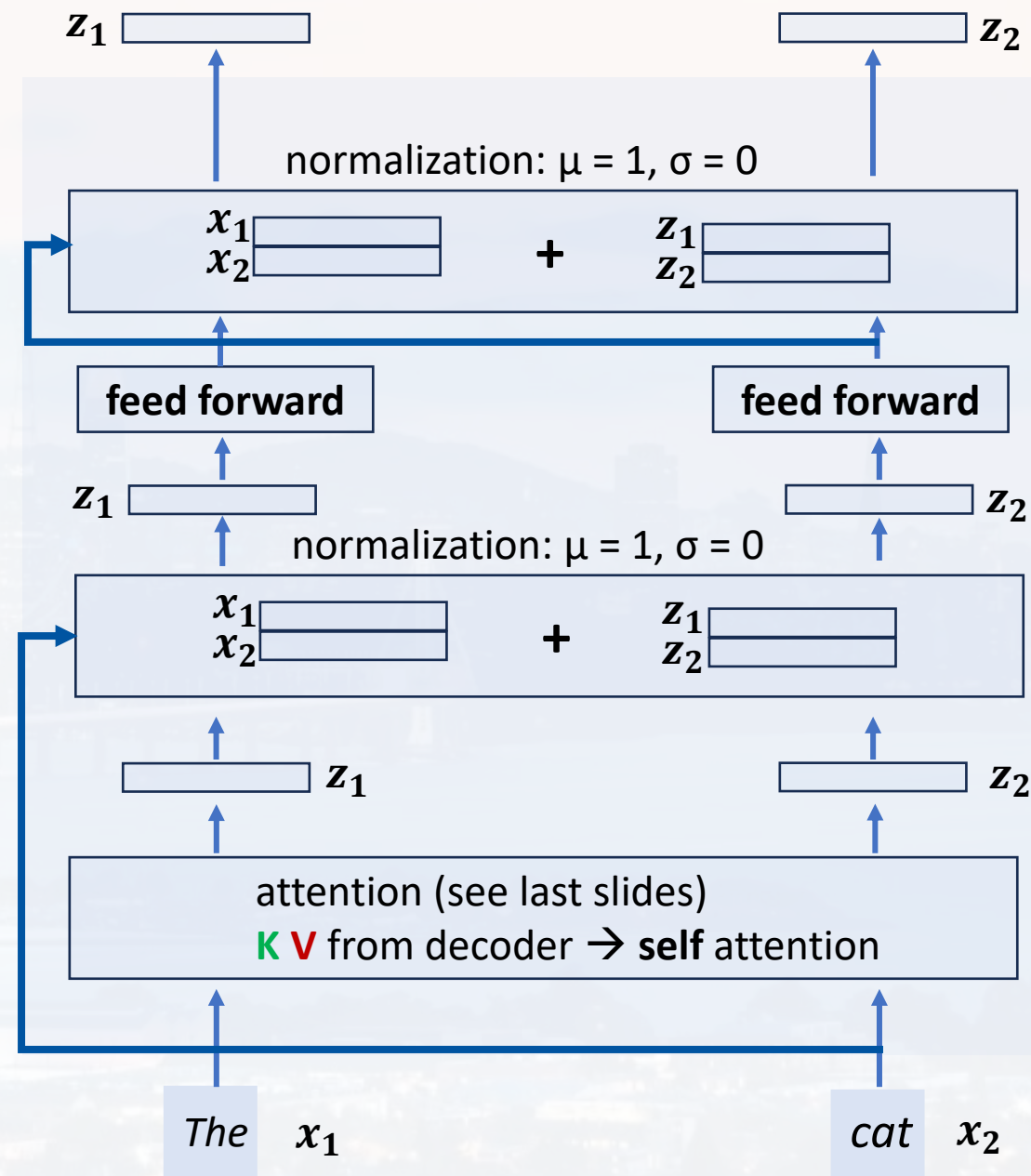ENCODER #1

Add & Normalize

Feed Forward      Feed Forward

Add & Normalize

Self-Attention

positional encoding     +     word embedding

$z_1$                      $z_2$

normalization: $\mu = 1$, $\sigma = 0$

$x_1$
$x_2$    +    $z_1$ $z_2$

**feed forward**               **feed forward**

$z_1$                  $z_2$

normalization: $\mu = 1$, $\sigma = 0$

$x_1$
$x_2$    +    $z_1$ $z_2$

$z_1$                  $z_2$

attention (see last slides)
**K** **V** from decoder → **self** attention

$The$   $x_1$                 $cat$   $x_2$

K V from encoder to decoder → cross attention

generated token: **argmax** from softmax (length = vocab size)

ENCODER #2

ENCODER #1

DECODER #2

DECODER #1

Add & Normalize

Feed Forward

Self-Attention

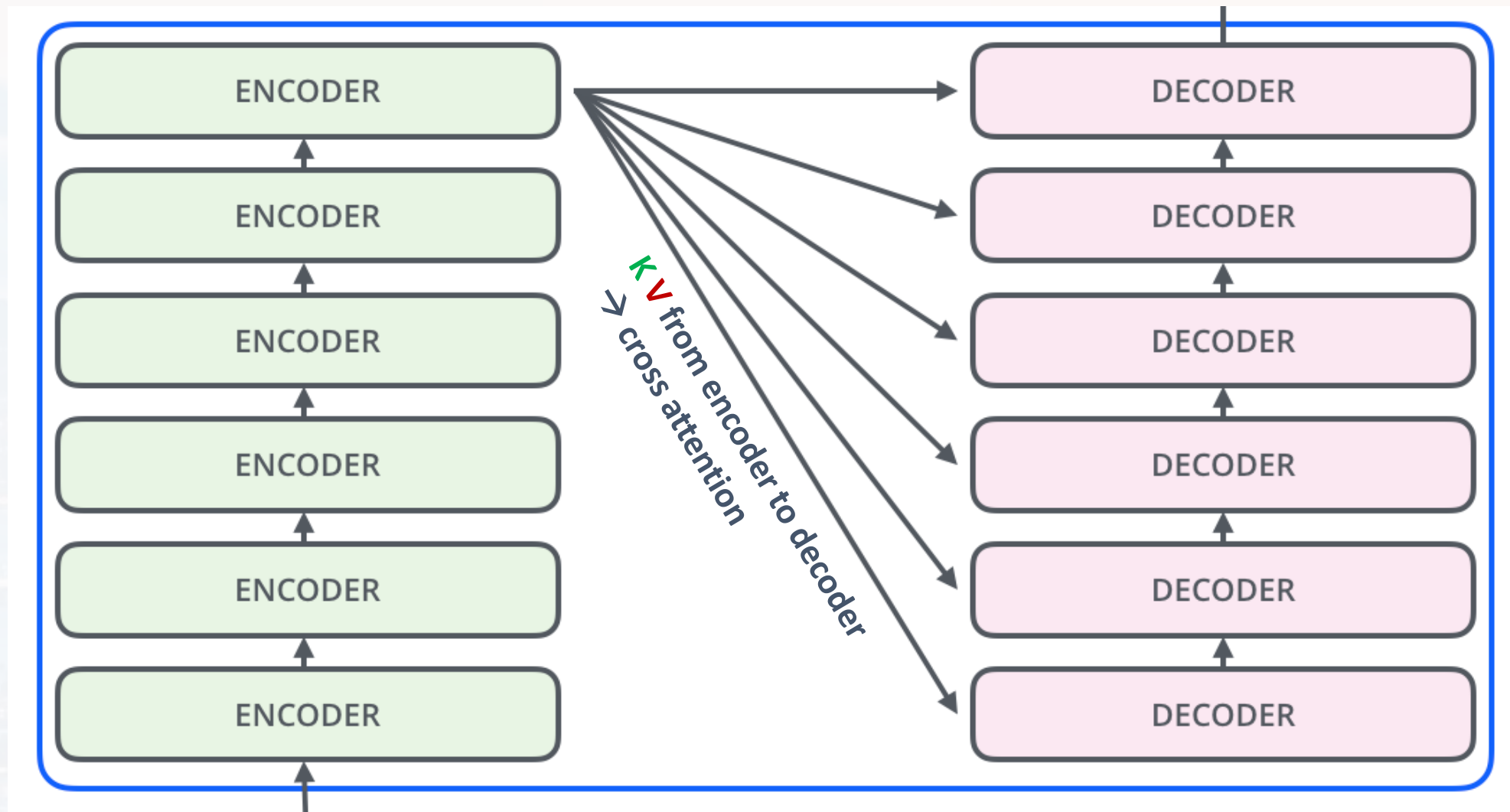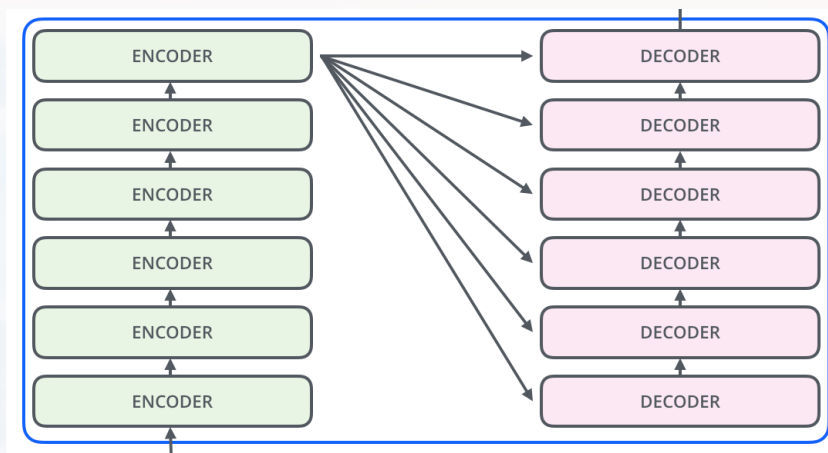Encoder-Decoder Attention

Softmax

Linear

https://jalammar.github.io/illustrated-transformer/

K V from encoder to decoder
→ cross attention



general: $N(E) \neq N(D)$

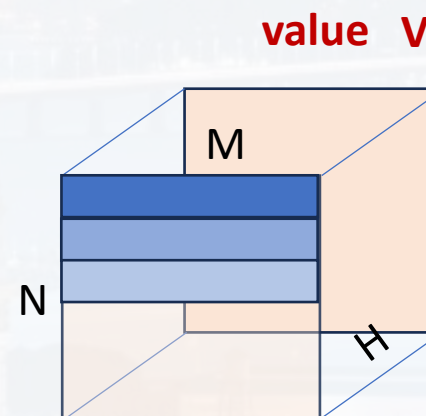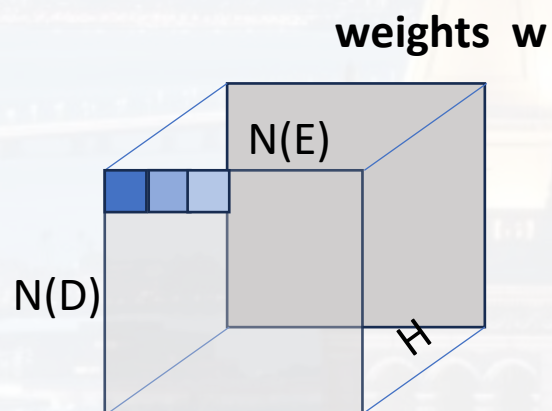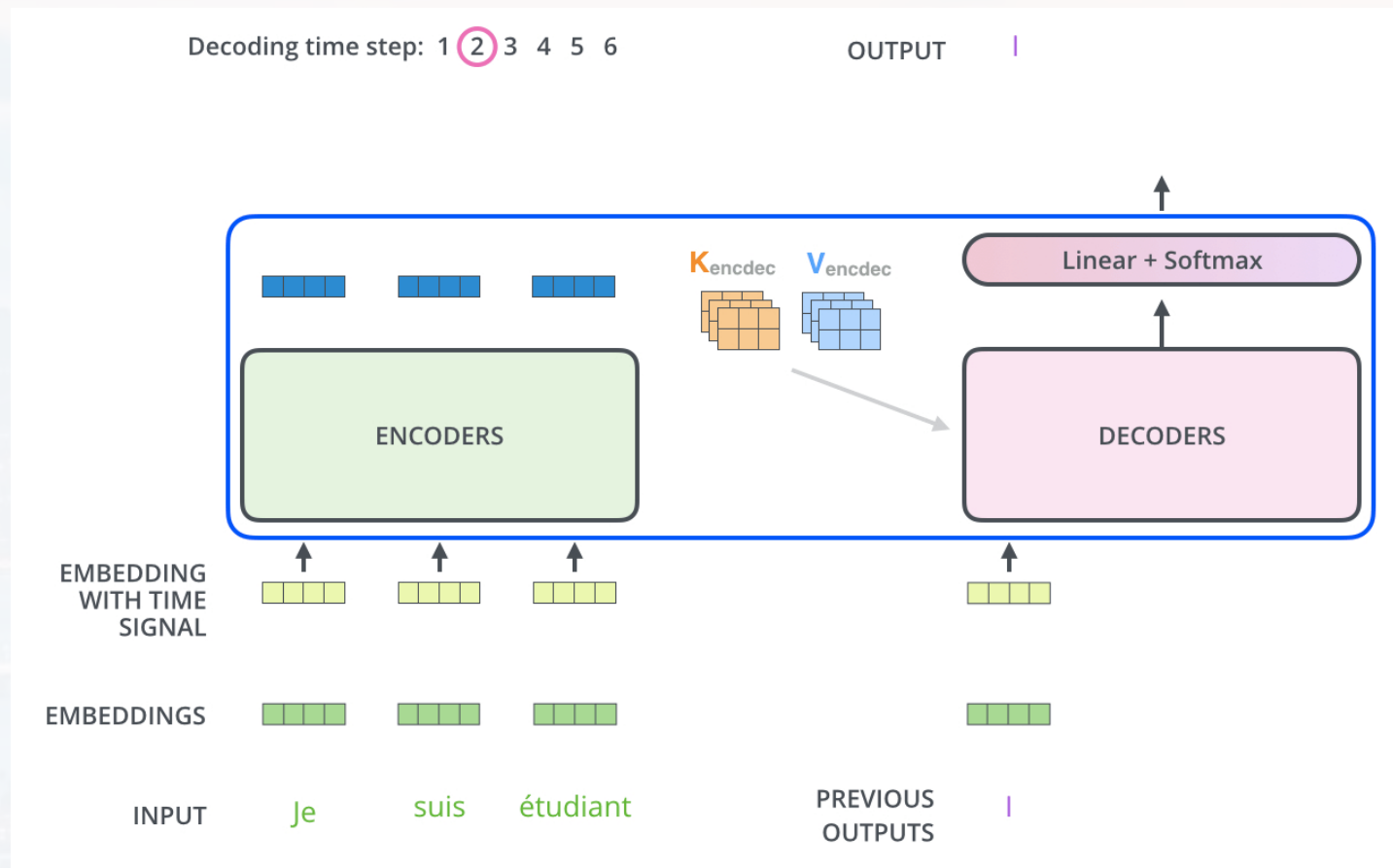| English: | Let there be light. | $N(E) = 4$ |
| German: | Es werde Licht. | $N(D) = 3$ |
| Latin: | Fiat lux. | $N(D) = 2$ |
| Hebrew: | yehi, or | $N(D) = 2$ |

**weights w**

**value V**

**attention:**



$N = N(encoder)$ , **if been called by the decoder**

https://jalammar.github.io/illustrated-transformer/

Attention Is All You Need (Vaswani et al, 2017)

**more about transformers:**

Jay Alammar

Interactive Visualization

transformers intro

building GPT from scratch!

Thank you very much for your



Head 1 of 12