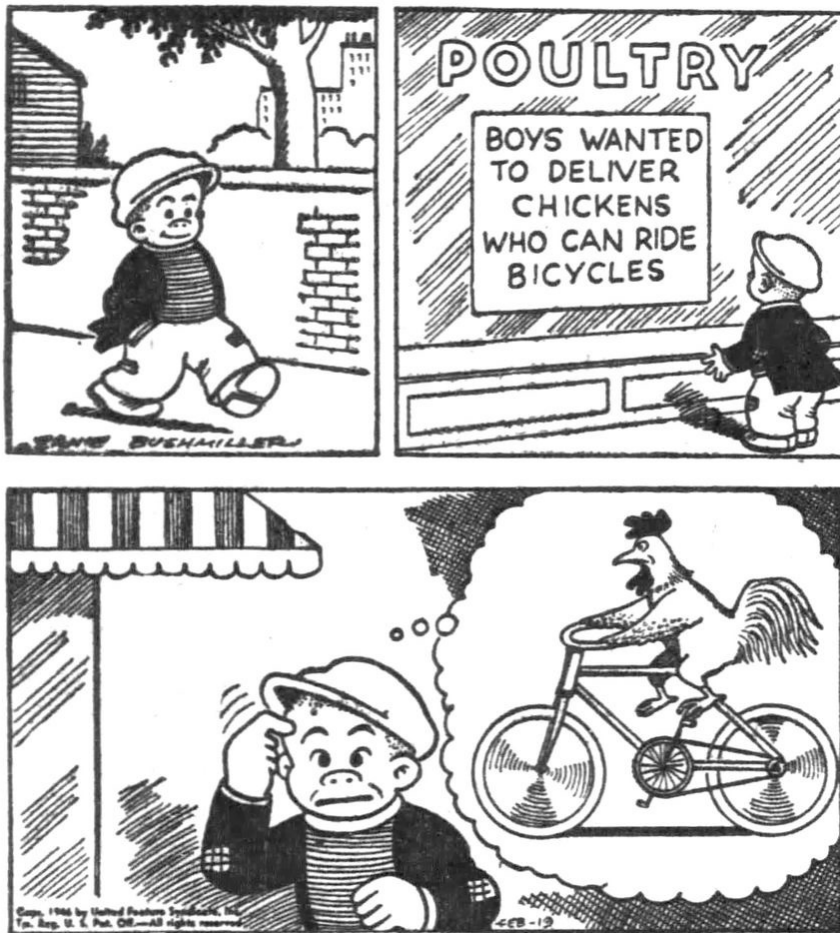# ANNOUNCEMENTS

1. All assignments, long and short, optional and non-optional are released

2. Kaggle competition + dataset released ; "BASELINES" out soon!

---

# Last few lectures ...
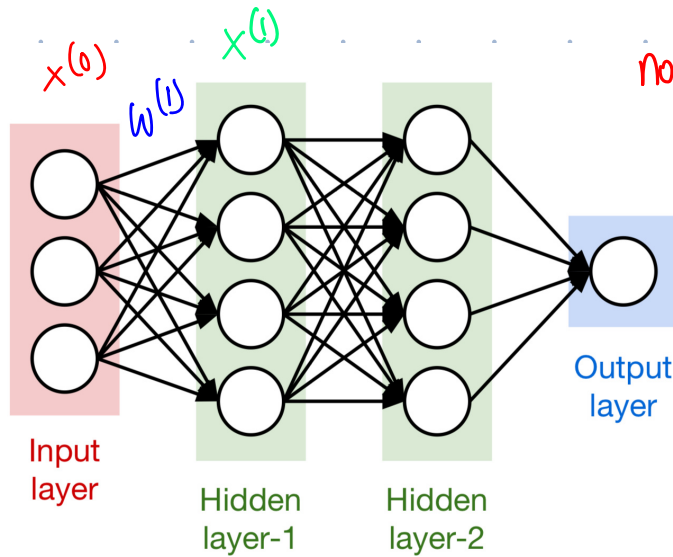
## I. FULLY-CONNECTED NETWORK
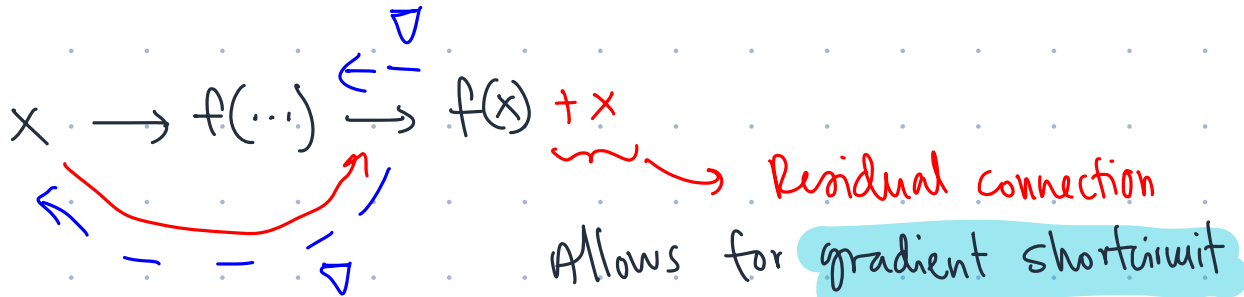
$$x^{(l)} = g\left(W^{(l)} x^{(l-1)} + b^{(l)}\right)$$

↓ non-linearity

$x^{(0)}$   $x^{(l)}$

$W^{(l)}$

Input layer

Hidden layer-1

Hidden layer-2

Output layer

## II. RESIDUAL CONNECTIONS

$$x \longrightarrow f(\cdots) \longrightarrow f(x) + x$$

→ Residual connection

Allows for gradient shortcut

## III: CONVOLUTIONAL NETWORKS —

— Conv filter (5×5×3)

image ←
(32×32×3)

32

32

3

*
*
*
*

4

sense of context becomes extremely relevant!

Q. what are the atomic units, the "pixels" of language?

" berry, the llama was spotted in Baker Lab"

1. words = Berry, the, llama, ... — the underlying dictionary / vocabulary needs to be constantly updated!

2. Characters = a, b, c, ... — scalable from a language perspective

10 words = 10T units of compute

10 byte/chars = 40T units of compute!

3. Subwords = "walk", "ing", "chat", "g", "p", "t"

The atomic units are called TOKENS!

# Vector representations (+ notion of similarity)

Tokenization:  where | are | you | going | ?

Iradhi says — use some form of an encoding

Need to represent in a numerical format, more precisely, we need a vector

$$V = \{ \text{where, when, are, you, going, from, moving, ?} \}$$

$$|V| = 8$$

|  | where | when | are | you | going | from | moving | ? |
|---|---|---|---|---|---|---|---|---|
| "where are you going?" | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 |
| "when are you moving?" | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| "where are you from?" | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |

DOT-PRODUCT SIMILARITY —

$$dot\left(s^{(1)}, s^{(3)}\right) = 4 \quad , \quad dot\left(s^{(1)}, s^{(2)}\right) = 3 \quad \text{"matt"}$$

FLAW — SENTENCE STRUCTURE IS SOMEWHAT CAPTURE, BUT SEMANTICS ARE MISSED!

If $|V| = 1m$, each vector is $1M$, and sparse (= mostly zeros!)

# Capturing "semantics" — a stupid example!

$s^{(1)}$ : this building is made of stone

$s^{(2)}$ : Ezra Cornell sat in Gates Hall

$s^{(3)}$ : this illu is made of stone

$s^{(4)}$ : Tushaar went to his illu

"illyun" — house

Matt said, " ... "

**SIMILAR WORDS APPEAR IN SIMILAR CONTEXTS!**

I                                                    II

" this illu is beautiful "

word

Context

target token

$$S = \left( t_1, \ldots, t_{i-w}, \ldots, t_{i-1}, t_i, t_{i+1}, \ldots, t_{i+w}, \ldots \right)$$

context                    context

# Binary classification that we DON'T care about!

Given a target token, $t_i$
context tokens, $\{t_{i-w}, \ldots, t_{i-1}, t_{i+1}, \ldots, t_{i+w}\} = c_i$

we want to build a binary cls,

$$P(+1 \mid t_i, c_i) \leftarrow \text{ probability that } c_i \text{ is a true context of } t_i$$

"Similar" tokens appear in similar ctxs

in our e.g., — "illu" = target,
$c_i = \{$ this, is, beautiful$\}$ is the true context.

We can use the notion of some similarity b/w $t_i, c_i$

$$P(+1 \mid t_i, c_i) = \prod_{-w \leq k \leq w} P(+1 \mid t_i, t_{i+k})$$

What is the proba that $t_{i+k}$ is a true context token of $t_i = $ "illu"?

ensures normalization b/w 0, 1

$$= \frac{1}{1 + \exp(-\text{dot}(t_i, t_{i+k}))}$$

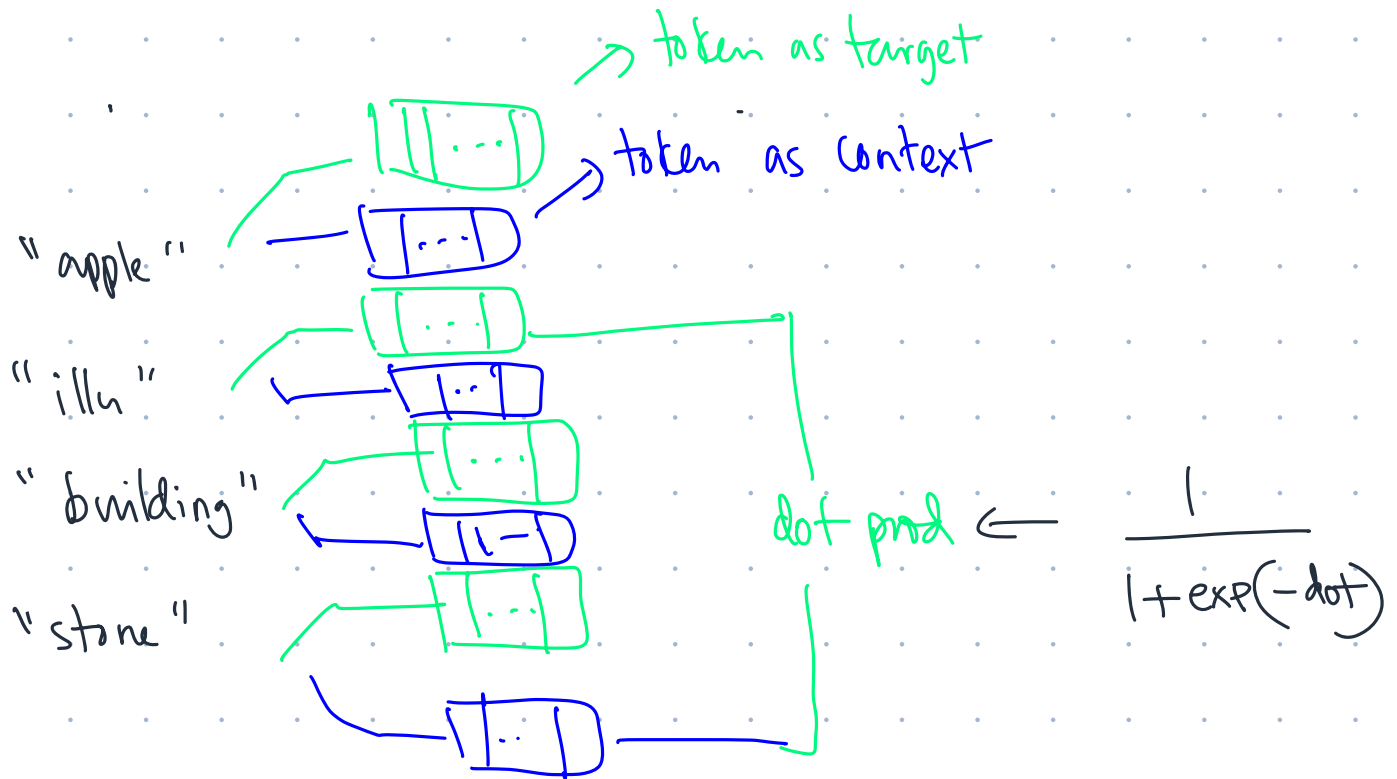Computes similarity b/w target, ctx token

(Similar to logistic regression!)

Q. What should $t_i, t_{i+k}$ be?

IDK what $t_i, t_{i+k}$ are, but we can learn them using GD!

token as target

token as context

"apple"

"illu"

"building"

"stone"

dot prod ← $\dfrac{1}{1 + \exp(-dot)}$

OBJ: maximize $P\left(y \mid \text{true target, true context}\right)$

— Started the idea of learnable embeddings!

<u>Reformulation</u> — Given a target token, can you predict the context tokens?

"illy"
target

$\longrightarrow$ "this"  "building"  "is"

Context

dot prod

$\left(\quad\boxed{|\ |\ |\cdots\ |}\quad\boxed{|\cdots\ |}\ \text{"apple"}\ \right) = \left.\begin{array}{c}\\ \\ \\ \end{array}\right\}$

$\left(\quad\quad\quad\quad\boxed{|\cdots|}\ \text{"building"}\ \right) = \left\}\begin{array}{c}\text{vector of}\\ |V|\end{array}\right.$

$\left(\quad\quad\quad\quad\boxed{|\cdots|}\ \text{"stone"}\ \right) =$

Normalization is done using <u>softmax</u> —  $\dfrac{\exp(z_j)}{\sum\limits_{k=1}^{n} \exp(z_k)}$

Aside: PERPLEXITY

Min. NLL loss — straightforward max log-likelihood
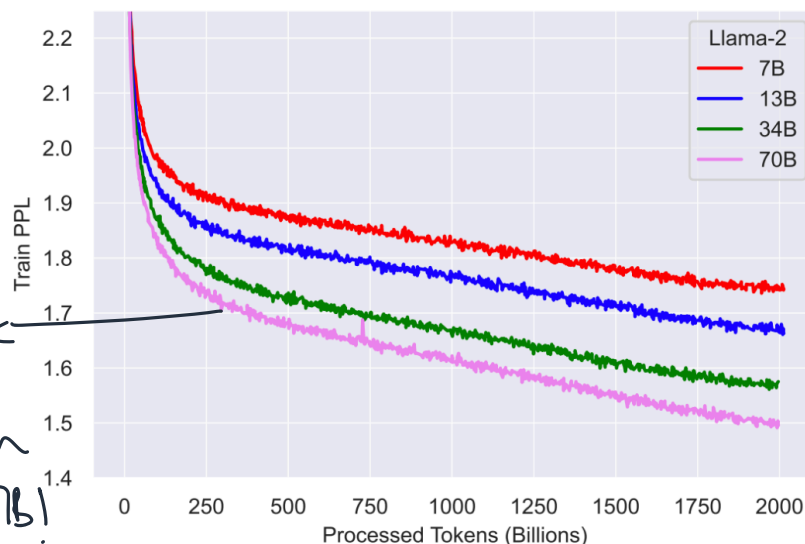take negative → NLL

↓                                    → model predicts
                                         "sentence"
PERPLEXITY —        "Complete this ___ "    true is "confirmation"
                                         RESULT — model is shocked!

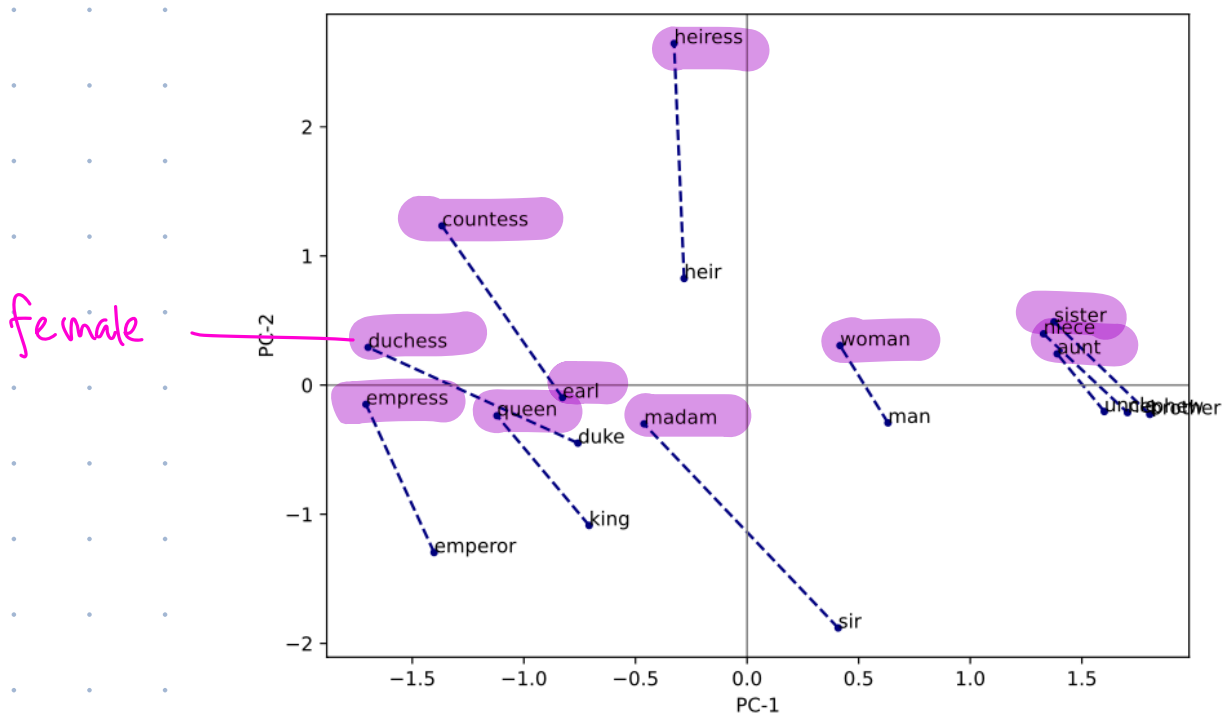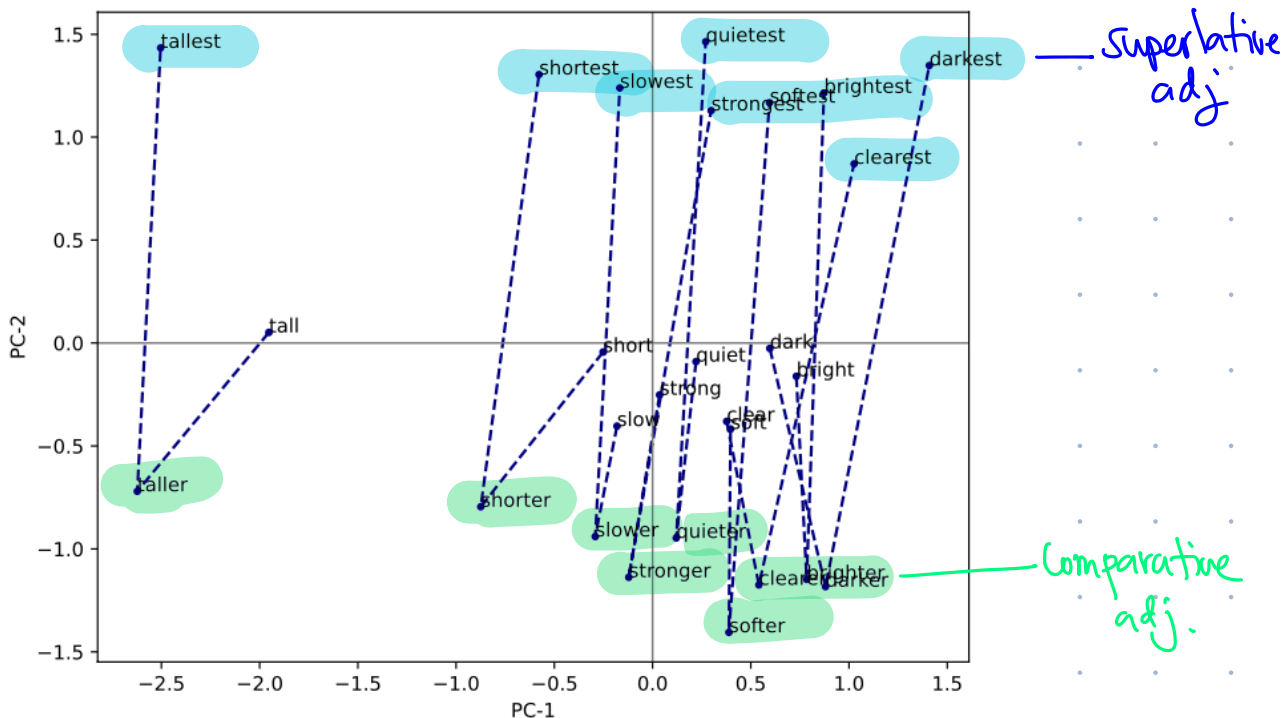| Hyperparams | | | | Dev Set Accuracy | | |
| --- | --- | --- | --- | --- | --- | --- |
| #L | #H | #A | LM (ppl) | MNLI-m | MRPC | SST-2 |
| 3 | 768 | 12 | 5.84 | 77.9 | 79.8 | 88.4 |
| 6 | 768 | 3 | 5.24 | 80.6 | 82.2 | 90.7 |
| 6 | 768 | 12 | 4.68 | 81.9 | 84.8 | 91.3 |
| 12 | 768 | 12 | 3.99 | 84.4 | 86.7 | 92.9 |
| 12 | 1024 | 16 | 3.54 | 85.7 | 86.9 | 93.3 |
| 24 | 1024 | 16 | 3.23 | 86.6 | 87.8 | 93.7 |

perplexity →

(lower is better)

Random NLP Benchmarks

(higher is better)



Llama-70B is known to be much better than Llama-7B!

Llama-2
— 7B
— 13B
— 34B
— 70B

Train PPL
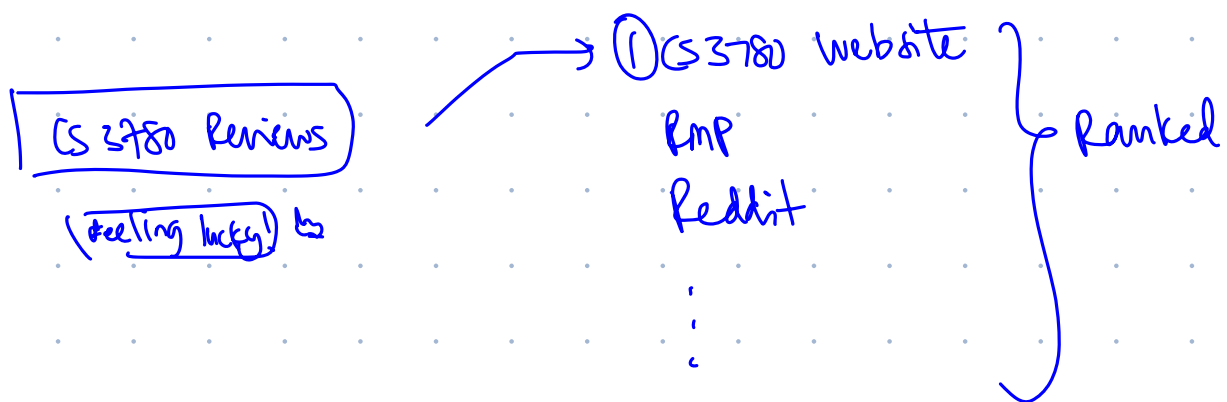
Processed Tokens (Billions)

# Visualizing learned embeddings using PCA



Top plot axes: PC-1 (x-axis), PC-2 (y-axis)

Labels (superlative adj., highlighted blue): tallest, shortest, slowest, quietest, softest, strongest, brightest, darkest, clearest

Labels (base): tall, short, quiet, dark, bright, strong, slow, clear, soft

Labels (comparative adj., highlighted green): taller, shorter, slower, quieter, stronger, clearer, darker, brighter, softer

Annotations: superlative adj., Comparative adj.

Bottom plot axes: PC-1 (x-axis), PC-2 (y-axis)

Labels: heiress, heir, countess, duchess, empress, queen, earl, duke, madam, king, emperor, woman, man, sister, niece, aunt, uncle, brother, nephew, sir

Annotation: female

$$\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{woman} = \overrightarrow{queen}$$

## Static → Dynamic embeddings

Once trained, these learnable embedding are `static`
↪ One embedding for each token

"The **bank** teller"      vs.      "the river **bank**"

↓                               ↙
same embedding representation

NEED   DYNAMIC   EMBEDDINGS!

→ ① CS3780 website
CS 3780 Reviews        RMP            } Ranked
(Feeling lucky!) ↪     Reddit
                       ⋮

Final document comes
    80% from ① + 20% from ② + 30% from ③

# (softmax) self-attention —

" the "    " bank "    " teller "

interested in emb for " bank "

$$\text{dot}(\text{vec}("the"), \text{vec}("bank"))$$ ← dot-prod. similarity b/w "the", "bank"
$$\text{dot}(\text{vec}("bank"), \text{vec}("bank"))$$
$$\text{dot}(\text{vec}("teller"), \text{vec}("bank"))$$

Returned links/ "keys"     Search "query"     use them

to form a weighted average of all embeddings

$$\langle "the", "bank" \rangle \cdot \text{vec}("the") +$$
$$\langle "bank", "bank" \rangle \cdot \text{vec}("bank") +$$
$$\langle "teller", "bank" \rangle \cdot \text{vec}("teller")$$

# POSITIONAL INFORMATION

bank : "The bank teller"    different from   bank: "The river bank"

"The bank teller sat at the river bank"
  (1)   (2)   (3)   (4) (5) (6) (7)  (8)

# Attention computations + Efficiency

$$O = \text{softmax}\left(QK^T\right) V$$

Time — $O(n^2 d)$

Space — $O(n^2)$ to store $QK^T$

X is a "n" long sequence, with d-dimensional tokens

$$X \in \mathbb{R}^{n \times d}$$

$$Q \quad K \quad V \quad \in \mathbb{R}^{n \times d}$$

$$\text{linear}(x) \quad \text{linear}(x) \quad \text{linear}(X)$$