

Large Language Models

- Web-Search - Recap
- Statistical Language models
- N-Gram Language Models & Translation Models
- Word Embeddings
- Contextual Word Embeddings

Credits:

Yoav Goldberg, Ido Dagan, Reut Tsarfaty , Moshe Koppel, Wei Song,
David Bamman, Ed Grefenstette, Chris Manning, Tsvi Kuflik,
Hinrich Schütze, Christina Lioma and more

Development:
Moshe Friedman

Web Search - Recap

The Paradigm Shift

Classic IR (Library Science)

Context: A lawyer searching Case Law.

Problem: **Scarcity**. There might be only one relevant precedent.

Goal: Find everything. Missing a document is expensive.

Metric: High Recall.

Web Search

Context: "Chocolate Cake Recipe".

Problem: **Abundance**. There are 5 million relevant recipes.

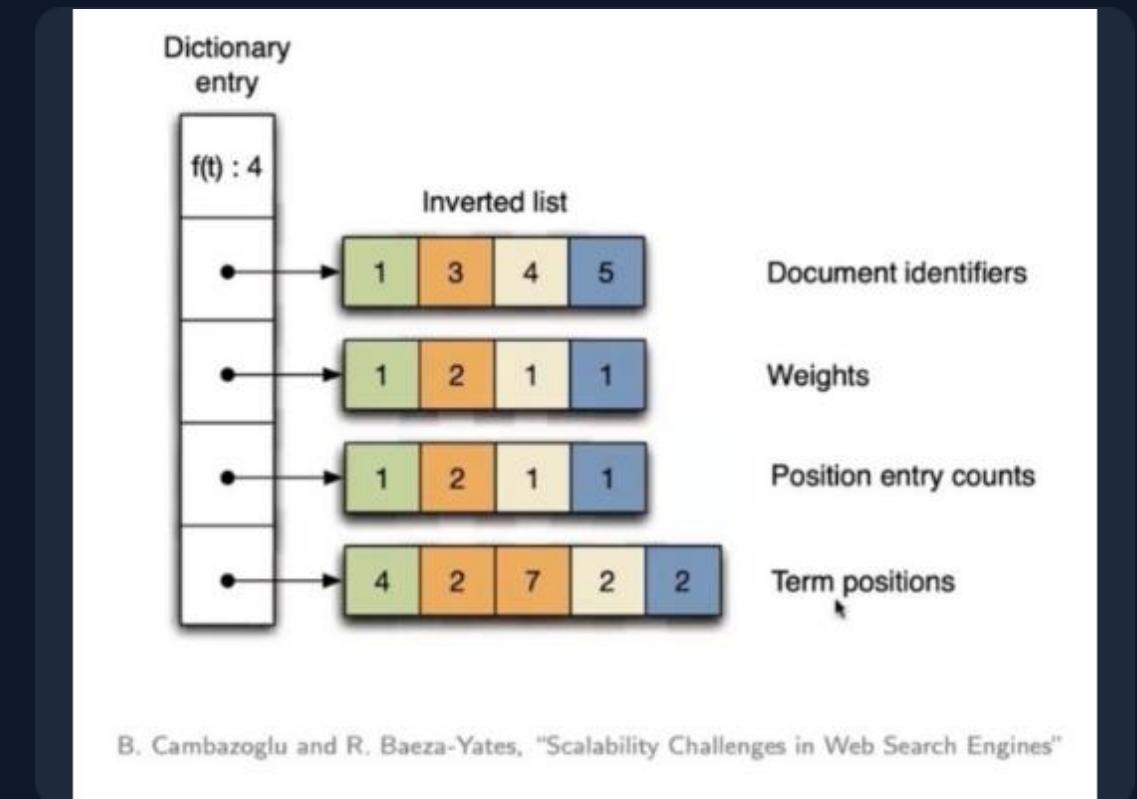
Goal: Find the *best* one immediately.

Metric: High Precision (specifically at the top).

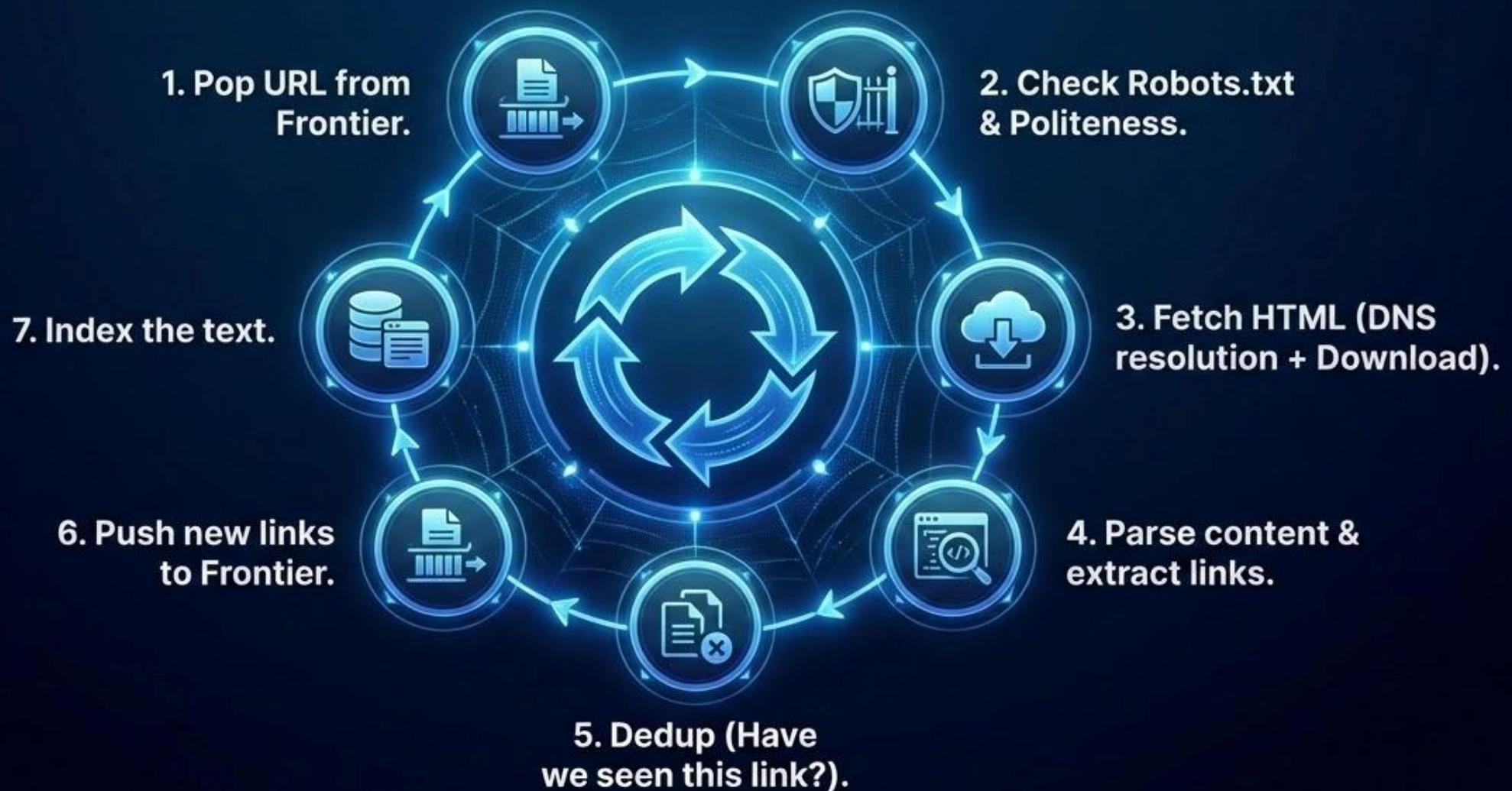
Classic IR Pipeline

The standard process for static document collections
(e.g., legal corpuses).

1. **Tokenization:** Splitting text into words.
2. **Normalization:** Lowercasing, removing punctuation.
3. **Stopword Removal:** Removing "the", "and" & "is".
4. **Stemming/Lemmatization:** "Running" -> "Run".
5. **Indexing:** Creating the Inverted Index.



The Spider's Loop



Anatomy of a URL

Uniform Resource Locators (URLs) are the addresses of the web. Understanding their structure is vital for crawling.

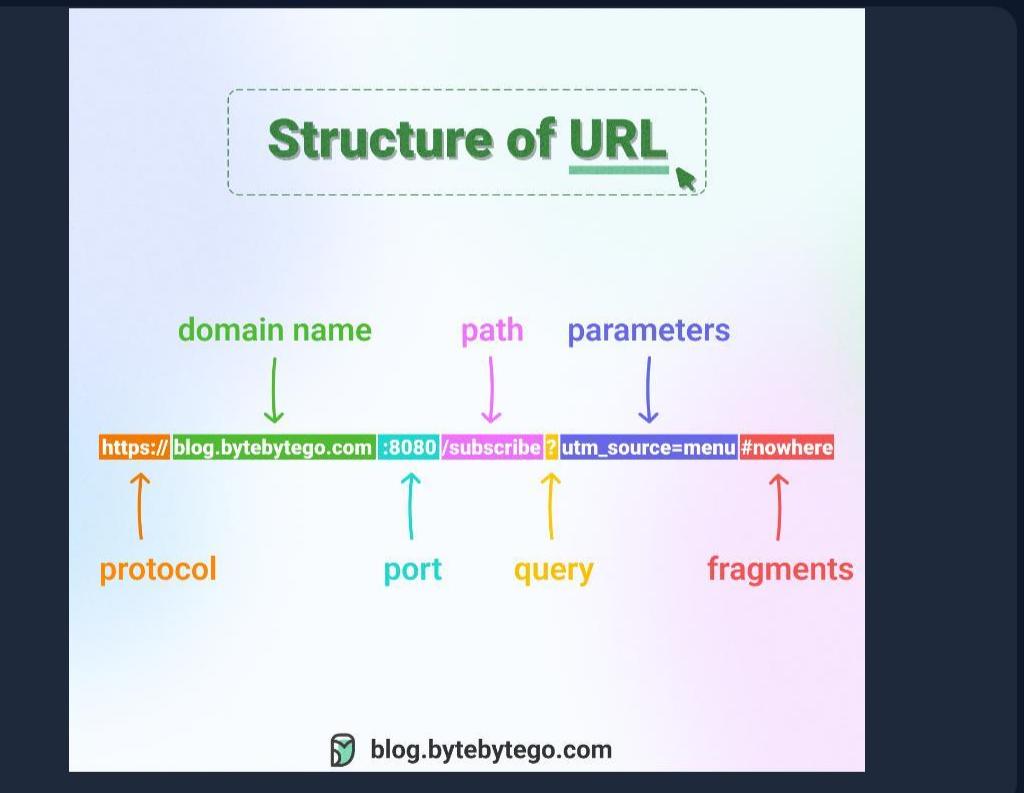
`https://` **Protocol:** How to transfer (HTTP/S).

`www.example.com` **Host/Domain:** The server address.

`/path/to/page` **Path:** Resource location on server.

`?id=123` **Query:** Parameters for dynamic pages.

`#section` **Fragment:** Anchor within the page (not sent to server).



Politeness & Robots.txt

DoS vs DDoS

DoS: Single crawler hammering a server. Crash caused by speed.

DDoS: Distributed crawler (1000 machines) hitting a server at once.

Solution: Crawl-delay and Robots.txt.

```
User-agent: * Disallow: /admin/ Crawl-delay: 5  
User-agent: BadBot Disallow: /
```

Crawling Strategy: BFS

Breadth-First Search

Goal: General Coverage.

Logic: Use a FIFO Queue. Crawl level d before level $d+1$.

Why? Prevents getting stuck in "rabbit holes" (infinite calendars or archives) on a single site.

Example

- Seed: Home (d=0)
- Q: [Sports, News] (d=1)
- Visit Sports: Adds [Article A, Article B]
- Q: [News, Article A, Article B]
- Visit News: (Crucial: We visit News *before* Article A).

The "Seed" Problem

You can't crawl the web if you don't know where to start. The web graph must be traversed from known entry points.

What makes a good seed?

- **High Out-Degree:** Links to many other pages.
- **Trustworthy:** Not spam.
- **Centrality:** Stable and long-standing.
- **Topical Coverage:** Represents the domain well.



Strategies for Finding Seeds

Directories

DMOZ (archived), BOTW, and niche directories are human-curated lists of quality sites.

Institutions

University homepages (.edu) and Government portals (.gov) are high-trust hubs.

Inverse Links

Use existing indices (like Ahrefs/Majestic) to find who links to known authorities.

#Social

Highly shared URLs on Twitter/Reddit often point to fresh, relevant content.

Curated Lists

"Best of" blog posts and Wikipedia "External Links" sections.

HITS

Run HITS on a small set to identify "Hubs", then use those Hubs as seeds for a larger crawl.

The Language of the Web

HTML (Structure)

The Nouns & Verbs.

Title Link

This is what crawlers care about most.

CSS (Presentation)

The Adjectives.

.price { color: green; }

Crawlers use this to detect hidden text (spam).

JS (Behavior)

The Actions.

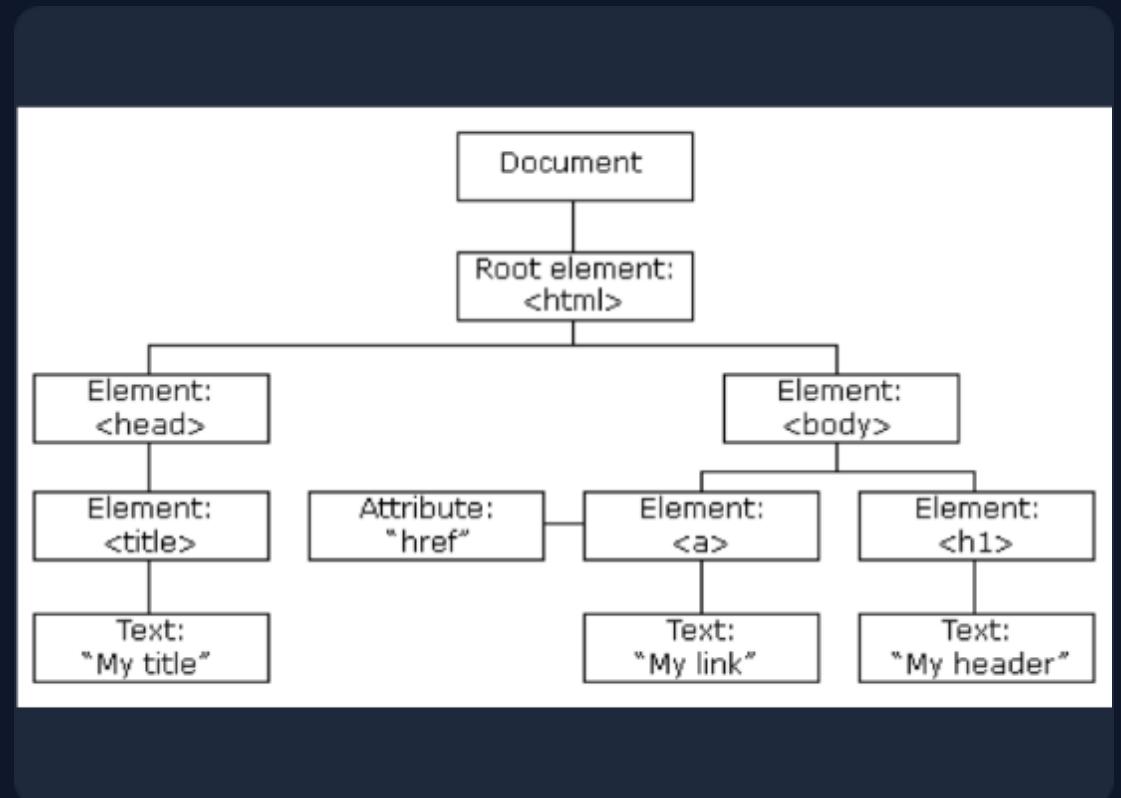
btn.onclick = () => {
 fetch('/data') }

Challenge: Crawlers must execute JS to see content on modern sites.

HTML: The Language of Search

Search engines don't "see" pages like humans. They parse the **DOM** (Document Object Model).

- **Tags:** </div>, provide semantic weight.
- **Links:** are the edges of the web graph.
- **Meta:** provides snippets.



Scraping vs. Crawling

Crawling (Discovery)

Goal: Discovery & Indexing.

Scope: Broad / The whole web.

Output: A list of URLs and cached text.

Example: Googlebot, Bingbot.

Tool: requests, Scrapy.

Scraping (Extraction)

Goal: Extraction.

Scope: Narrow / Specific sites.

Output: Structured data (prices, names).

Example: Price comparison, Lead gen.

Tool: BeautifulSoup (Static), Selenium (Dynamic).

Static Scraping

The Mechanics

Pros: Very fast, low CPU usage.

Cons: Cannot see content loaded by JavaScript.

Libraries: requests, BeautifulSoup.

```
import requests from bs4
import BeautifulSoup

# 1. Fetch
raw HTML response =
requests.get("https://example.com")

# 2. Parse
soup = BeautifulSoup(response.text,
'html.parser')

# 3. Extract specific element (e.g. $19.99
price)
price = soup.find('span', class_='price').text
print(price) # Output: $19.99
```

Dynamic Scraping

The Mechanics

Pros: WYSIWYG. Can handle infinite scroll and SPAs (React/Vue).

Cons: Slow. Launches a full browser instance.

Libraries: Selenium, Puppeteer.

```
from selenium import webdriver
from selenium.webdriver.common.by import
By
import time

# 1. Launch Browser
driver = webdriver.Chrome()
driver.get("https://example.com/dynamic")

# 2. Wait for JS to render
time.sleep(2)

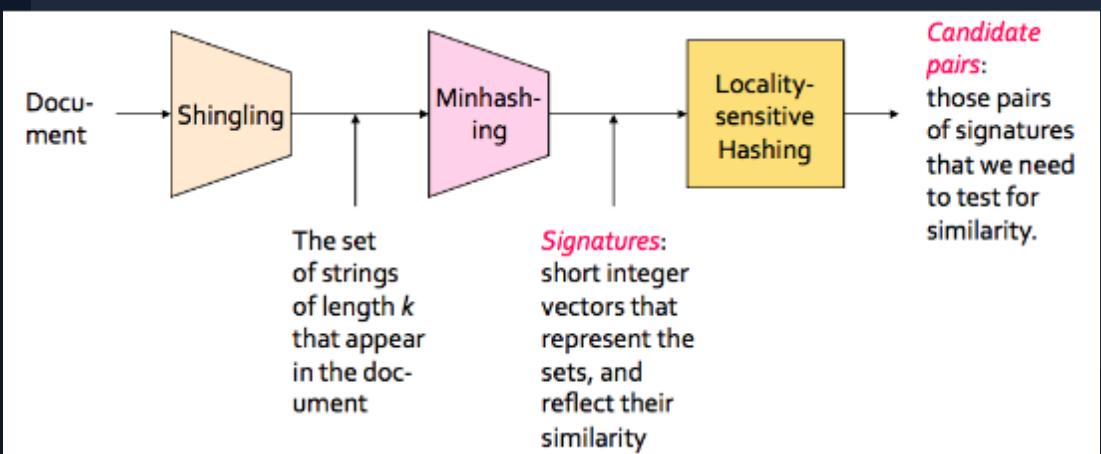
# 3. Extract content generated by JS
price = driver.find_element(By.CLASS_NAME,
"price").text
print(price) driver.quit()
```

Near-Duplicate Detection

The Problem

Approximately 30% of the web consists of near-duplicates (mirrors, revisions, spam). Exact match detection fails on dynamic timestamps or minor edits.

Exact hash matching (MD5) fails if one-byte changes (e.g., dynamic timestamp).



The Solution: Shingling

Shingles: Convert text into sets of N-grams (e.g., "a rose is a rose").

Jaccard Coefficient: Measure similarity by dividing the intersection of shingle sets by their union.

MinHashing: A technique to approximate Jaccard similarity efficiently at scale.

PageRank (PR): The Random Surfer

Imagine a surfer clicking random links forever.

PageRank (PR) is the probability that the surfer is on a specific page at any given time.

More incoming links = Higher probability.

Links from high-probability pages = More weight.

Damping factor (d): Probability of clicking (usually around 0.85).

Teleportation ($1-d$): Prevents getting stuck in dead ends (Sink nodes).

PageRank algorithm

„random surfer approach“

Importance of a web page is measured by its popularity
How many incoming links it has

PageRank can be defined by the probability that a random surfer on the web starts on a random page
+ follows hyperlinks AND visits the given page

→ sum of column values equals 1 because of the probabilities
→ it is like **Markov-chains**
→ the „transition matrix“ defines the next steps
→ stationary distribution: is the final **PageRank** vector

$$\mathbf{H} \mathbf{x} = \mathbf{x}$$

HITS: Hubs & Authorities

Proposed by Kleinberg. Query-dependent.

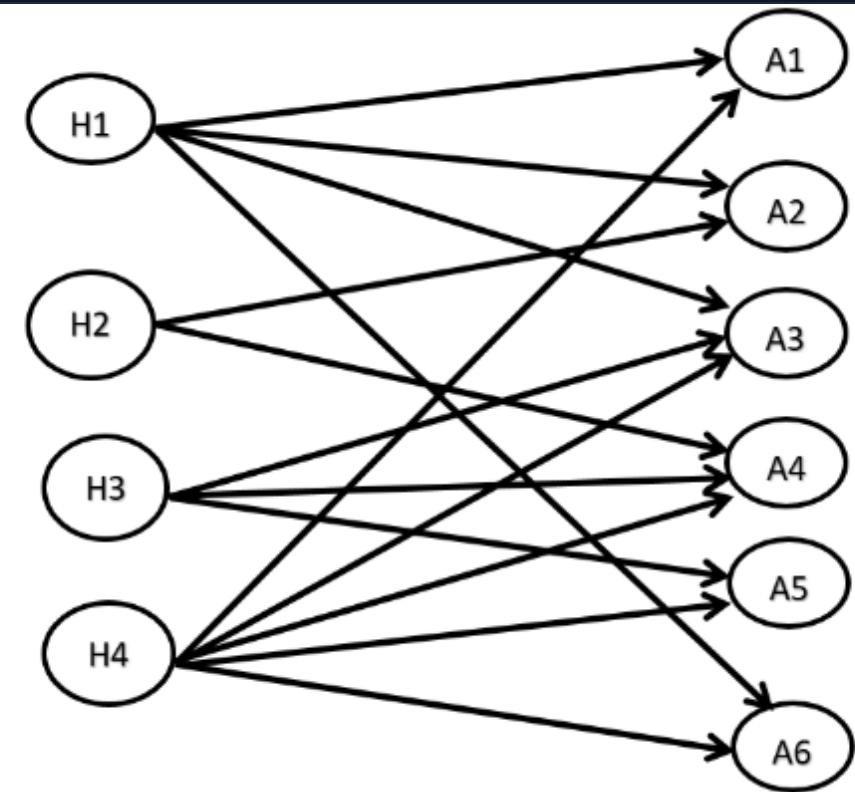
Authorities: Pages with good content (linked to by many hubs).

Hubs: Pages with good lists of links (link to many authorities).

Mutually reinforcing:

Good hubs point to good authorities.

Good authorities are pointed by good hubs.



The Spam Arms Race

If ranking = money, people will cheat.

Content Spam

Keyword stuffing, hidden text,
scraped content, auto-generated
gibberish.

Link Spam

Link farms, paid links, comment
spam, buying expired domains.

Cloaking

Showing search engines one version
of a page and users another (e.g.,
porn/gambling).

Machine Learning

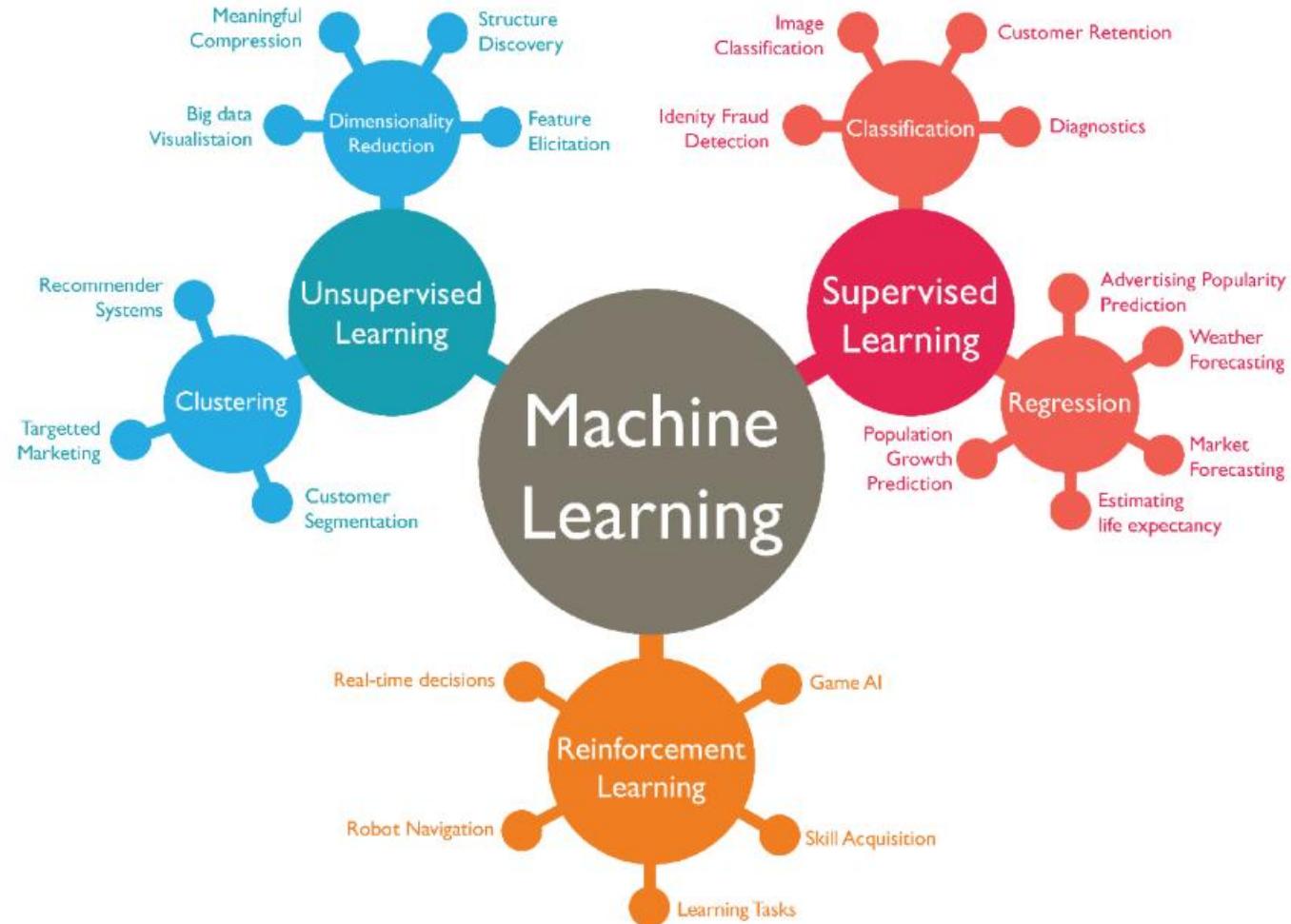
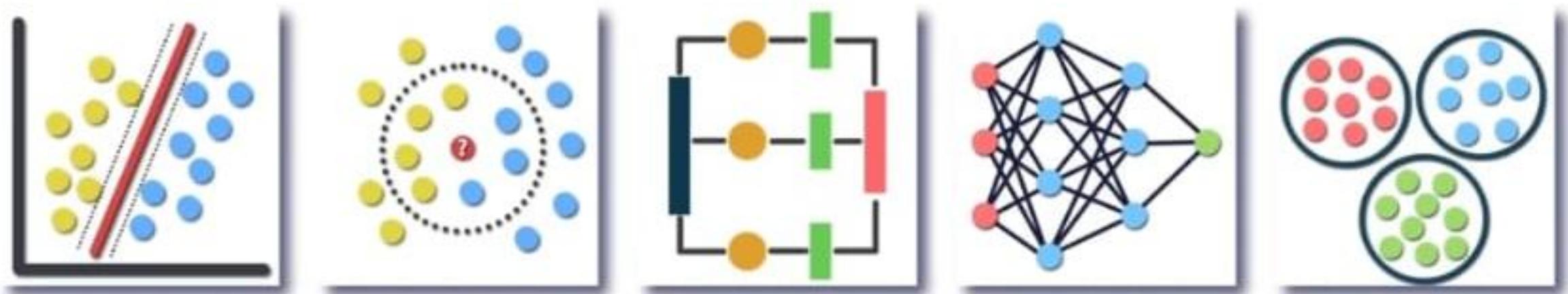
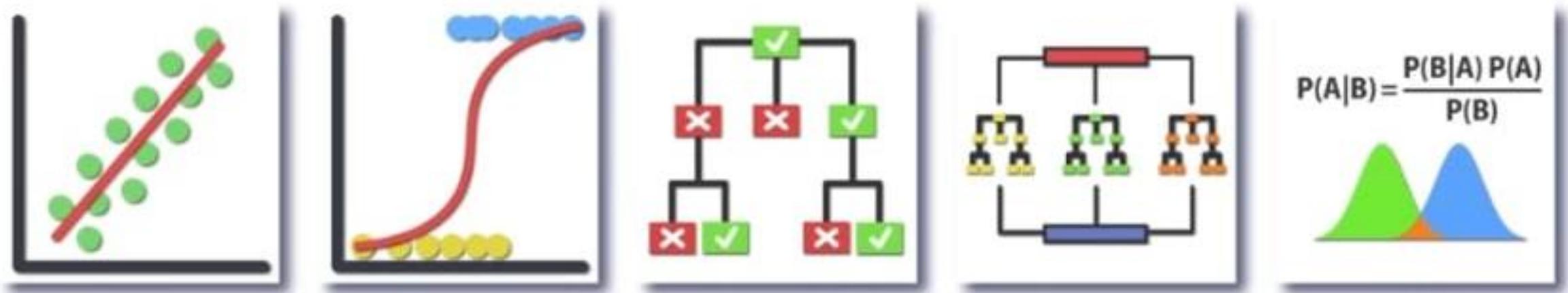


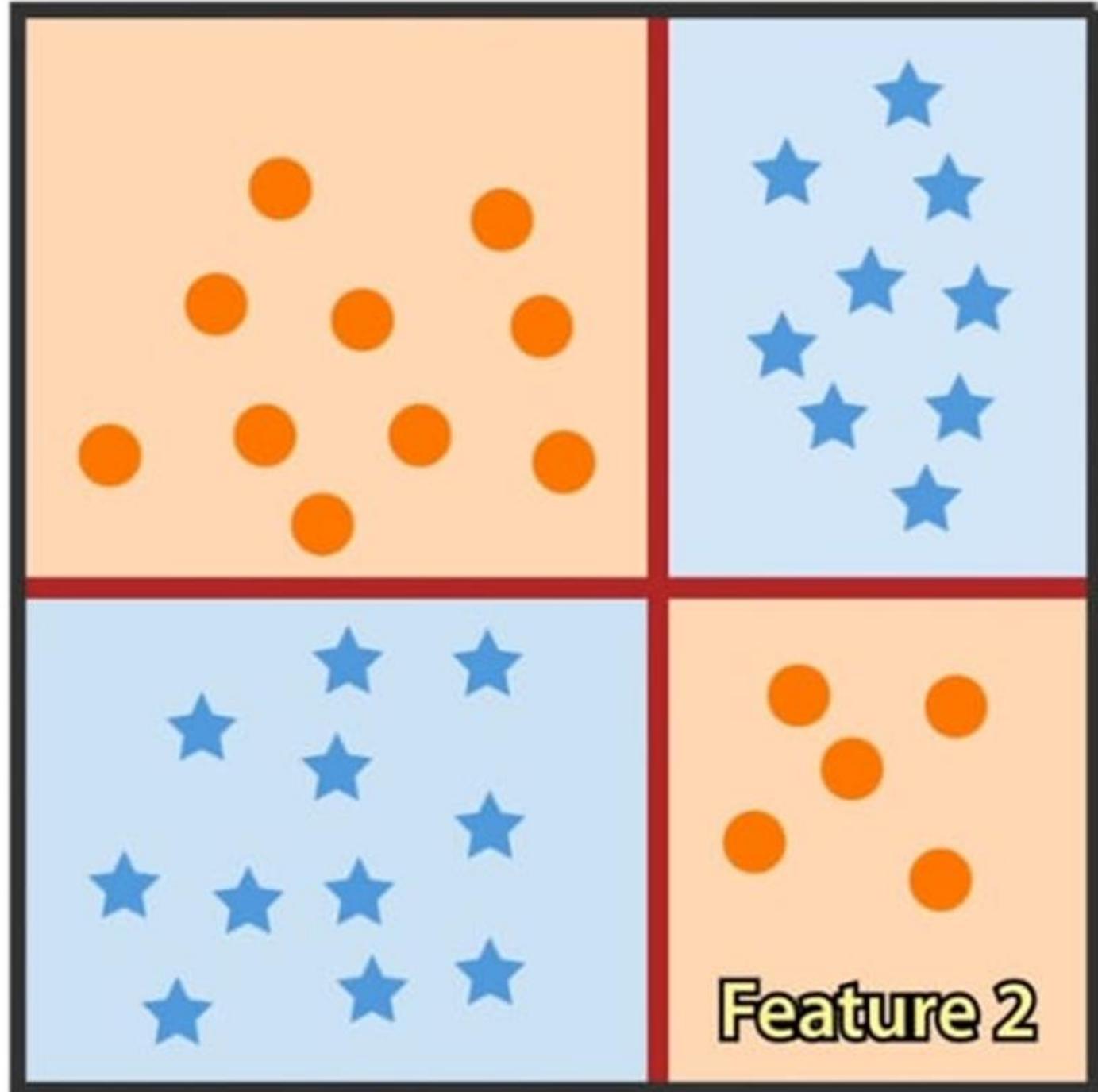
Image via [Abdul Rahid](#)

Machine Learning – Algorithms and problems



Decision Trees

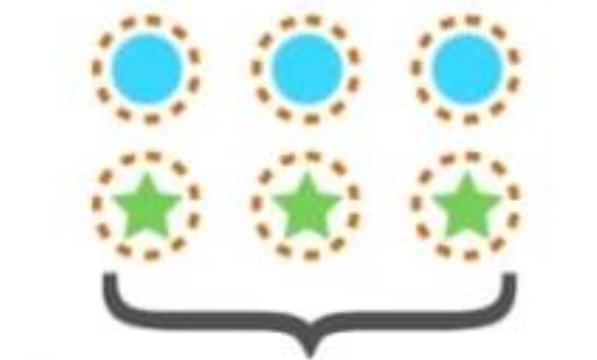
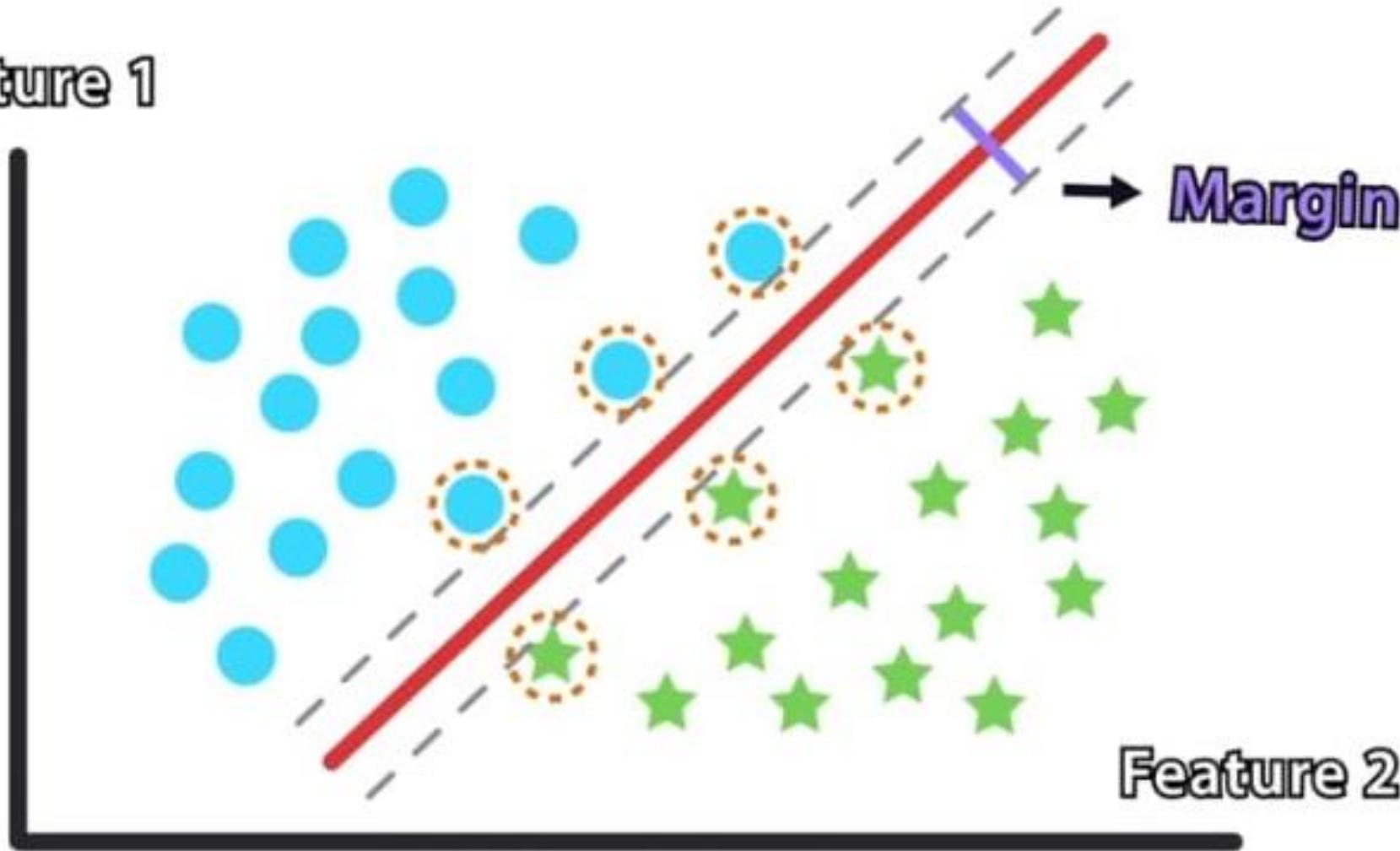
Feature 1



Feature 2

SVM

Feature 1

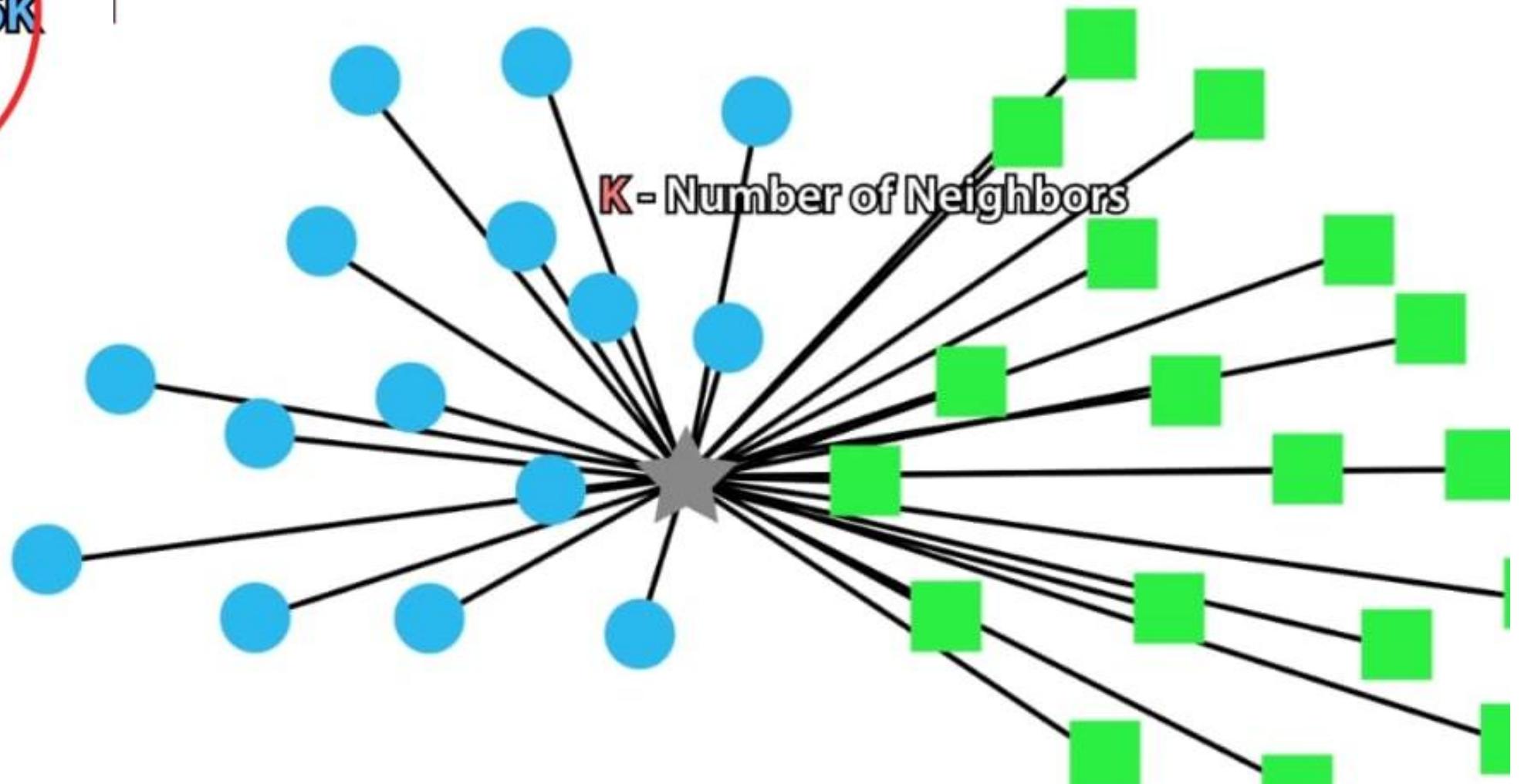


Support Vectors

20K 28K
44K 32k 36K
32K

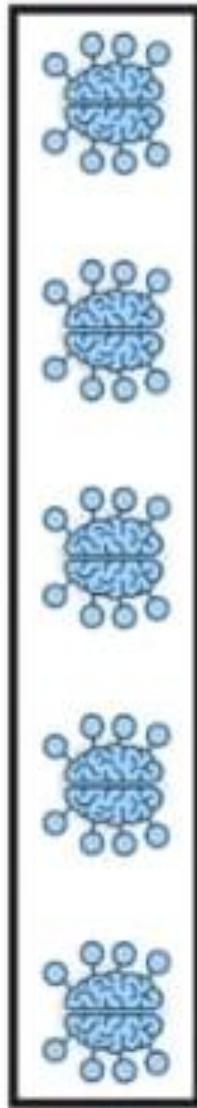
KNN

CLA

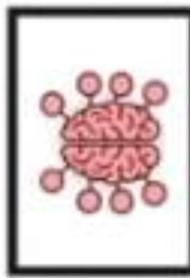


Stacking

Base Models



Meta Model



Ensemble Methods

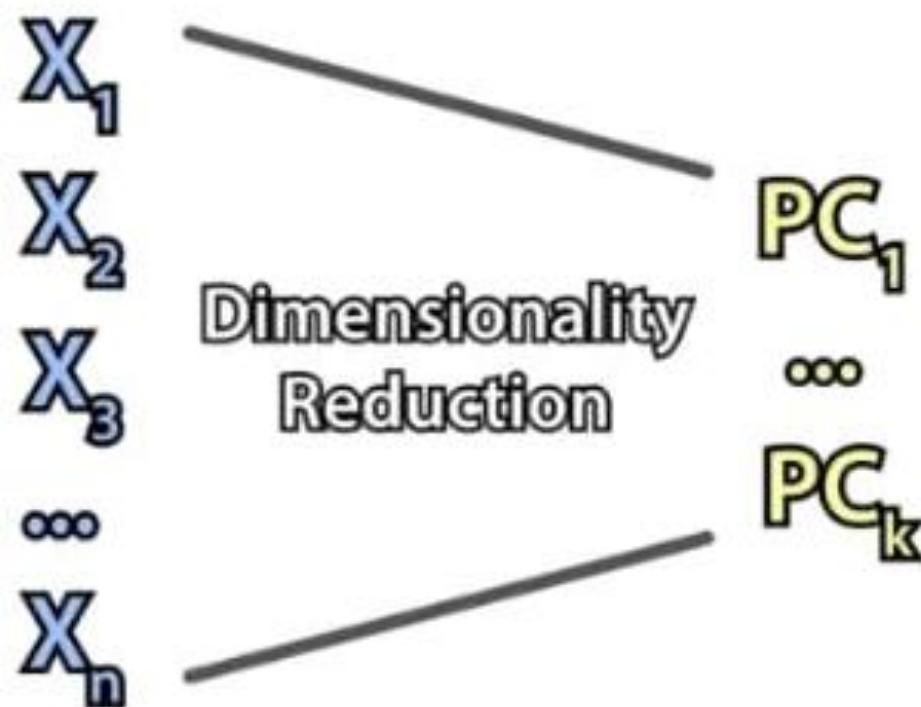


Training
Data



Validation
Data

Principal Component Analysis



Where

$$k < n$$



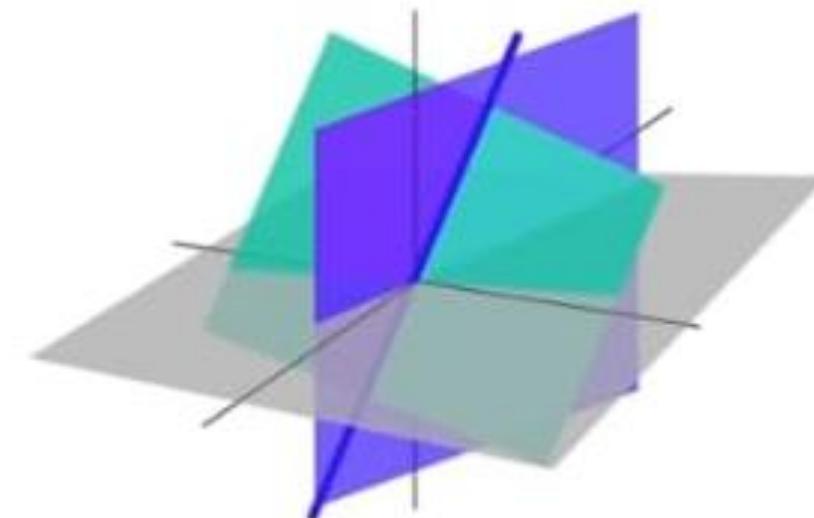
$$PC\ 1 = w_1 X_1 + w_2 X_2 + w_3 X_3$$

$$PC\ 2 = w_4 X_1 + w_5 X_2 + w_6 X_3$$

Principal Component Analysis

$$\text{PC 1} = w_1 X_1 + w_2 X_2 + w_3 X_3$$

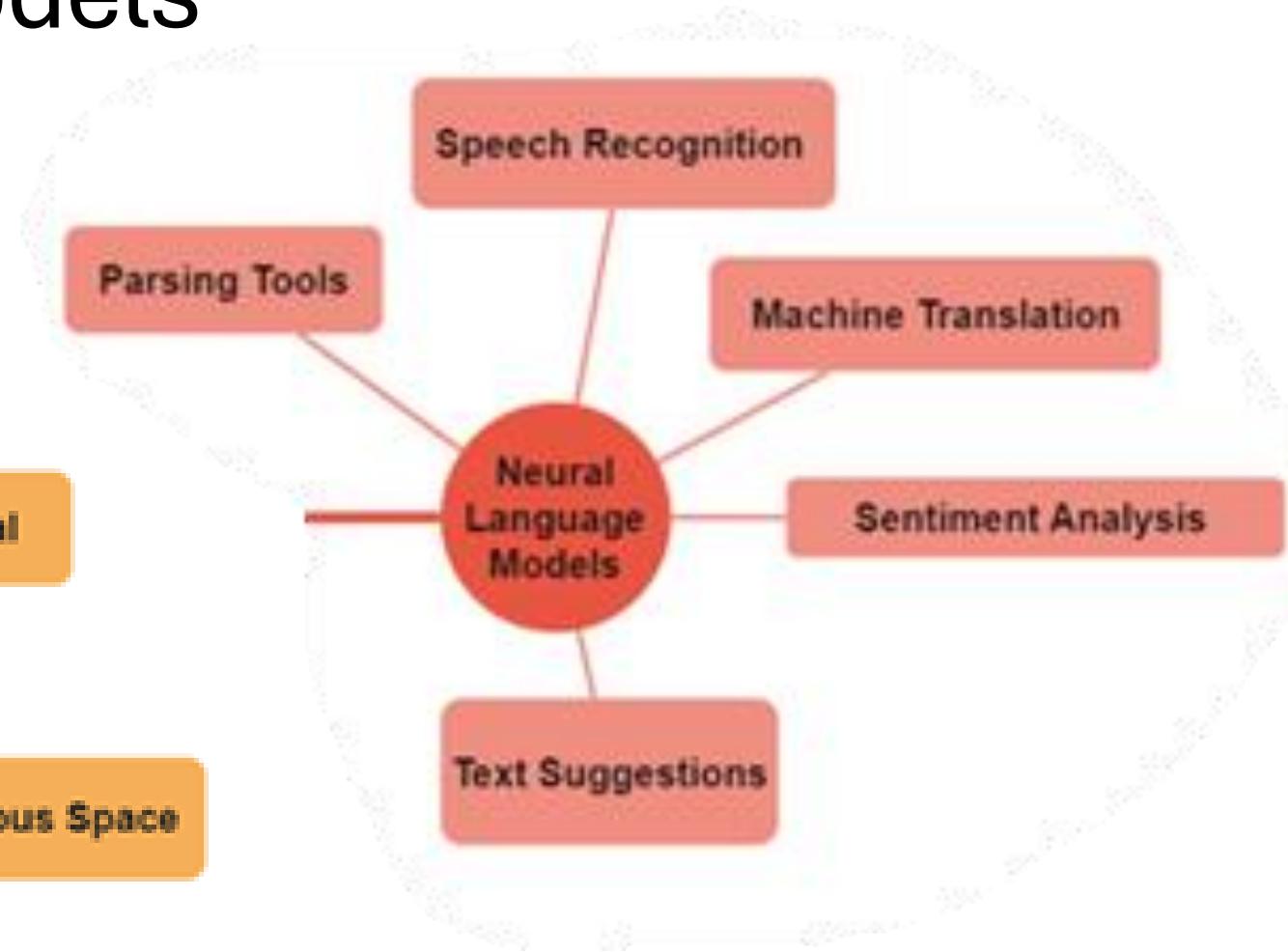
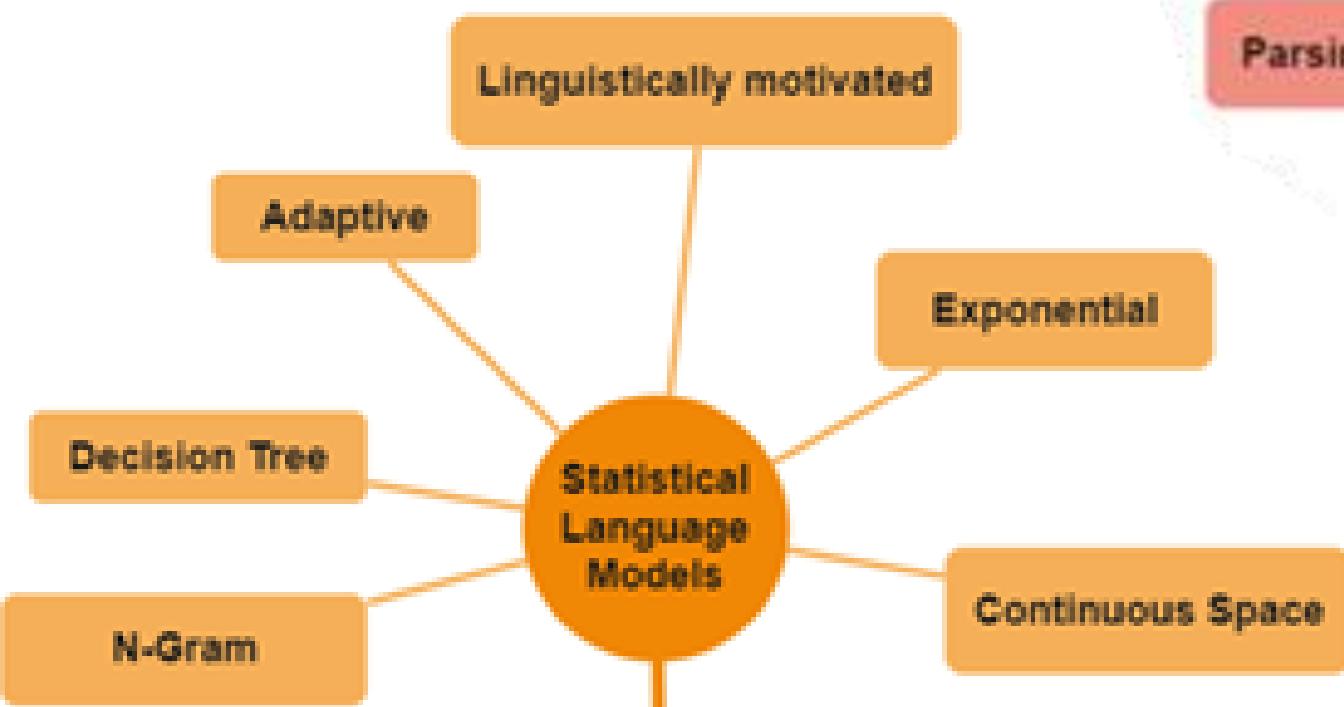
$$\text{PC 2} = w_4 X_1 + w_5 X_2 + w_6 X_3$$



Eigenvalues and Eigenvectors

$$\leftarrow \mathbf{Av} = \lambda \mathbf{v}$$

Statistical Language Models & Natural Language Models



Relevancy in IR

- The IR task can be viewed as finding **relevant** documents given a query
 - Main problem – how to model relevancy?
- Relevancy model can be:
 - binary (relevant, non-relevant)
 - probabilistic (score-based)
- Previous probabilistic models try to directly estimate the probability: $\Pr(R \mid \text{query}, \text{document})$

Measure the fluency of documents

- How likely this document is generated by a given language model
 - If $p_{text-mining}(d) > p_{health}(d)$, d belongs to text mining related documents
 - If $p_{user_a}(d_1) > p_{user_a}(d_2)$, recommend d_1 to $user_a$
 - A relative concept

What is a Language Model?

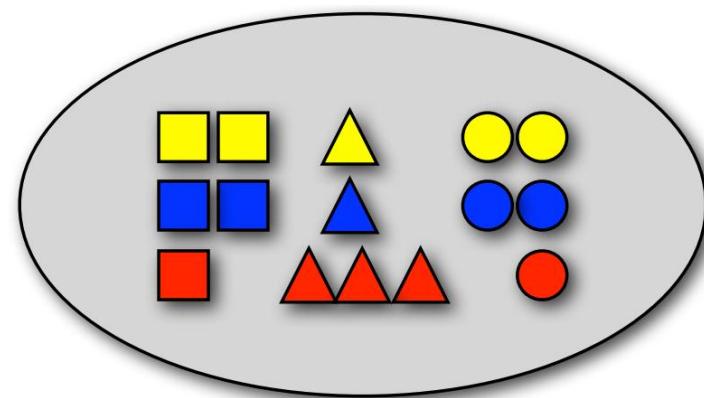
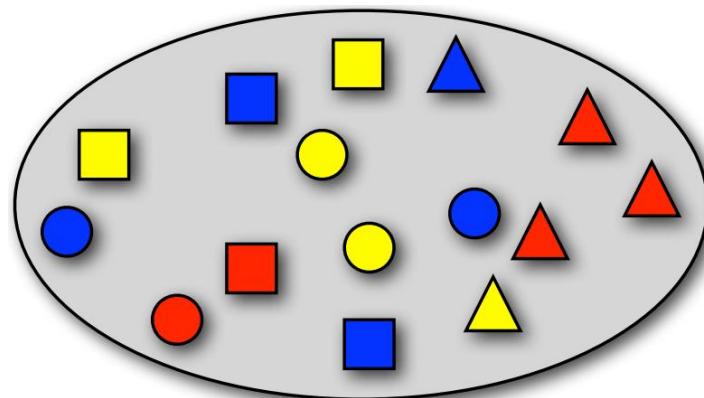
- Probability distribution over strings of text
 - How likely is a given string in a given “language” L
 - Example:
 - $p_1 = \Pr_L(\text{"a quick brown dog"})$
 - $p_2 = \Pr_L(\text{"dog quick a brown"})$
 - $p_3 = \Pr_L(\text{"grdfgdrfe brown dog"})$

If L is English: $p_1 > p_2 > p_3$

In most models in this lesson: $p_1 = p_2$ (bag of words)

Sampling with replacement

- Pick a random shape, then put it back in the bag



$$P(\square) = 2/15$$

$$P(\text{blue}) = 5/15$$

$$P(\text{blue} \mid \square) = 2/5$$

$$P(\square) = 1/15$$

$$P(\text{red}) = 5/15$$

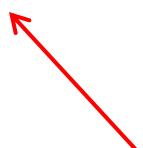
$$P(\square) = 5/15$$

$$P(\square \text{ or } \triangle) = 2/15$$

$$P(\triangle \mid \text{red}) = 3/5$$

Essential probability concepts

- Probability of events
 - Mutually exclusive events
 - $P(A \cup B) = P(A) + P(B)$
 - General events
 - $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 - Independent events
 - $P(A \cap B) = P(A)P(B)$



Joint probability, or
simply as $P(A, B)$

Why is a LM useful?

- Provide a principled way to quantify the uncertainties associated with natural language
- Allow us to answer questions like:
 - Given that we see “*John*” and “*feels*”, how likely will we see “*happy*” as opposed to “*habit*” as the next word?
(speech recognition)
 - Given that we observe “baseball” three times and “game” once in a news article, how likely is it about “sports” v.s. “politics”?
(text categorization)
 - Given that a user is interested in sports news, how likely would the user use “baseball” in a query?
(information retrieval)

Kinds of language models

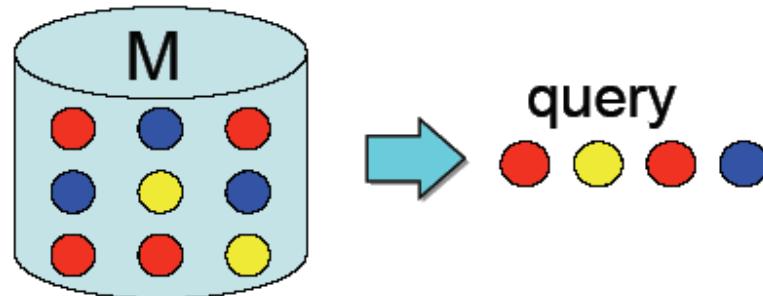
- Linguistic complexity (unigram or higher order)
- Statistical framework (multinomial or multiple-Bernoulli)

What is a statistical LM?

- A model specifying probability distribution over word sequences
 - $p(\text{"Today is Wednesday"}) \approx 0.001$
 - $p(\text{"Today Wednesday is"}) \approx 0.000000000001$
 - $p(\text{"The eigenvalue is positive"}) \approx 0.00001$
- It can be regarded as a probabilistic mechanism for “generating” text, thus also called a “generative” model

Unigram Language Models

- Words are sampled independently from each other
 - Metaphor: randomly picking words from a bag
 - Independence makes joint probability simple
 - Estimation of probabilities: simple counting

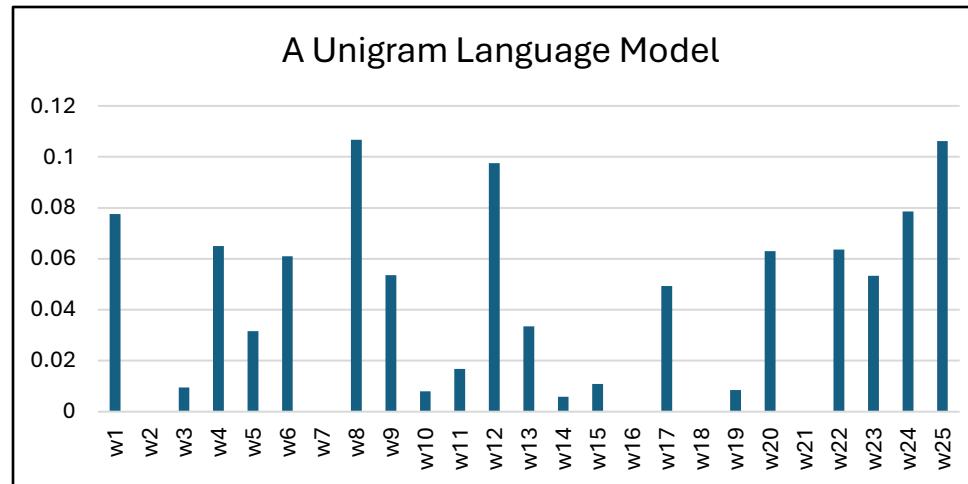


$$\begin{aligned} P(\text{red yellow red blue}) &= P(\text{red}) \ P(\text{yellow}) \ P(\text{red}) \ P(\text{blue}) \\ &= 4/9 * 2/9 * 4/9 * 3/9 = .015 \end{aligned}$$

Unigram language model

- Generate a piece of text by generating each word independently
 - $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2) \dots p(w_n)$
 - s.t. $\{p(w_i)\}_{i=1}^N, \sum_i p(w_i) = 1, p(w_i) \geq 0$
- Essentially a multinomial distribution over the vocabulary

*The simplest and
most popular
choice!*



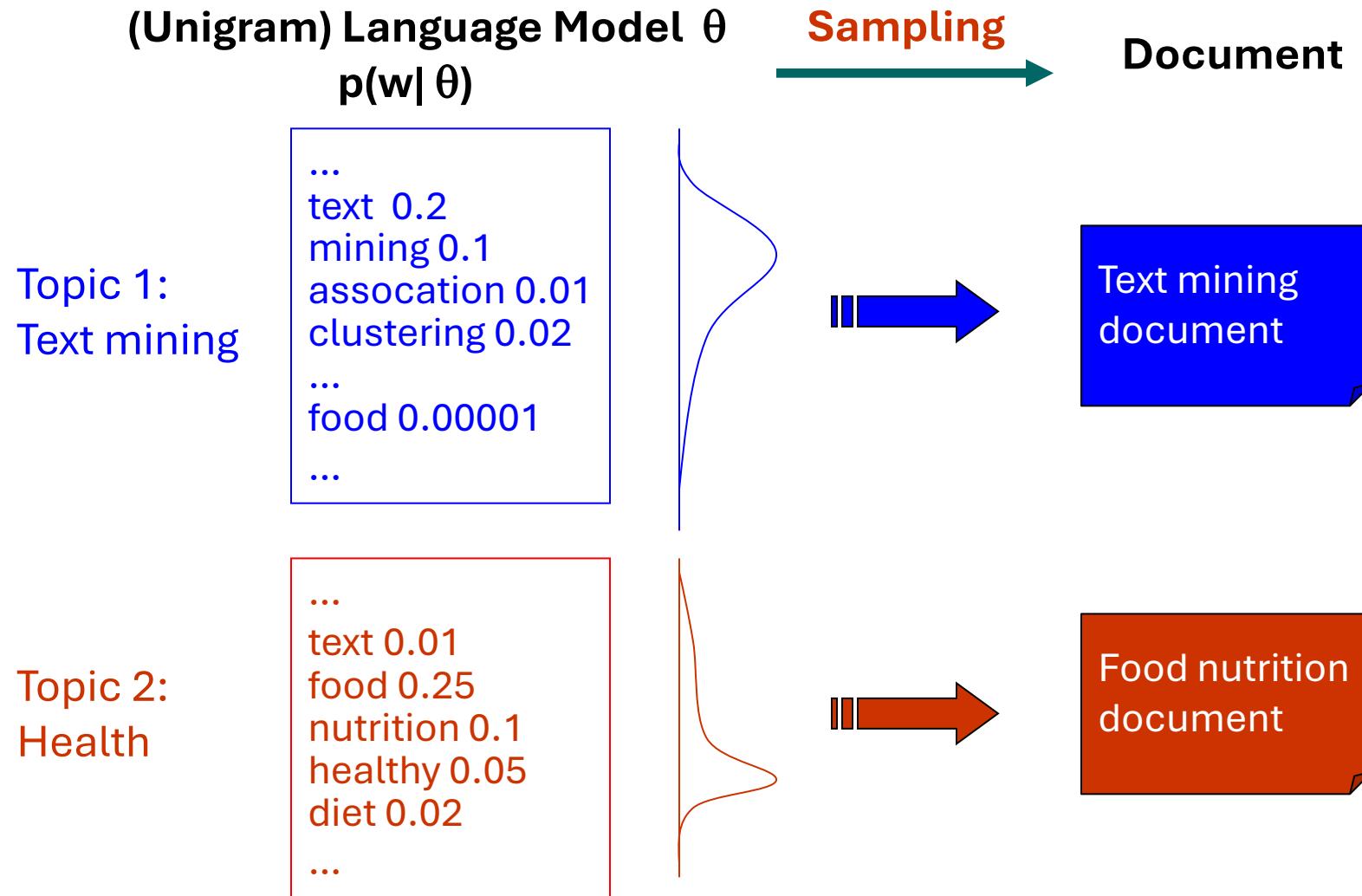
More sophisticated LMs

- N-gram language models
 - Conditioned only on the past n-1 words
 - E.g., bigram: $p(w_1 \dots w_n)$
 $= p(w_1)p(w_2|w_1) p(w_3|w_2) \dots p(w_n|w_{n-1})$
- Remote-dependence language models (e.g., Maximum Entropy model)
- Structured language models (e.g., probabilistic context-free grammar)

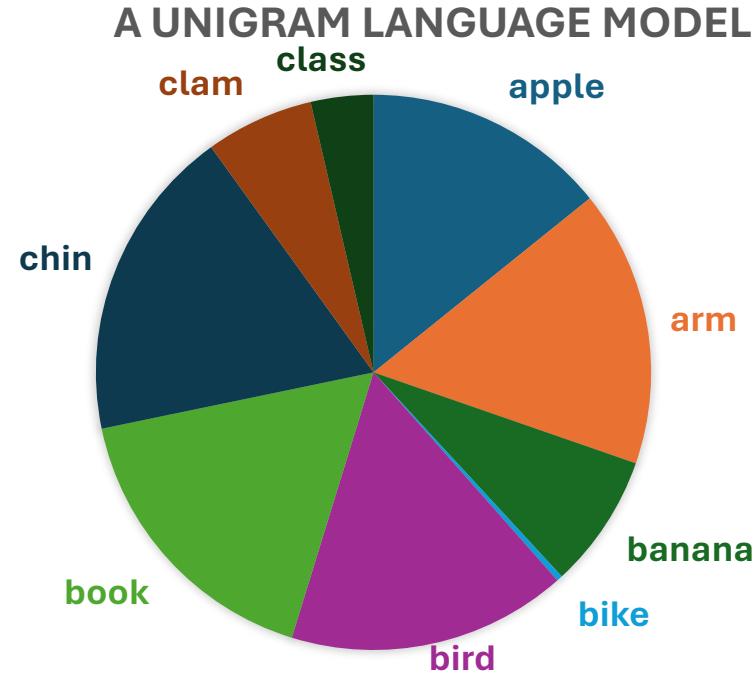
Why just unigram models?

- Difficulty in moving toward more complex models
 - They involve more parameters, so need more data to estimate
 - They increase the computational complexity significantly, both in time and space
- Capturing word order or structure may not add so much value for “topical inference”
- But, using more sophisticated models can still be expected to improve performance ...

Generative view of text documents

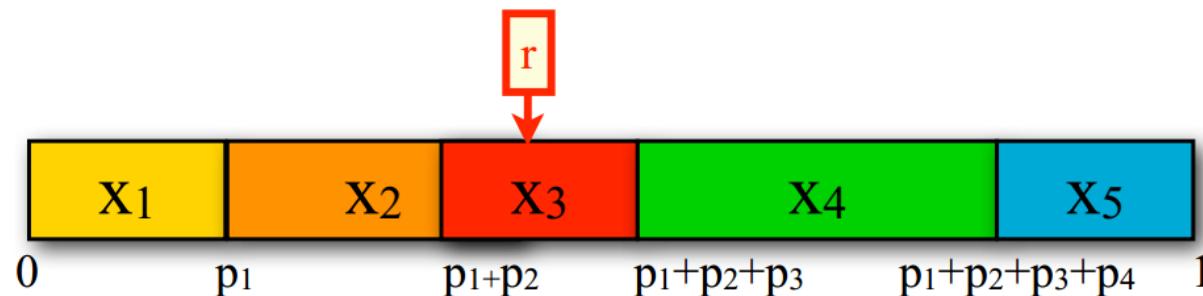


How to generate text from an N-gram language model?



How to generate text from an N-gram language model?

- Sample from a discrete distribution $p(X)$
 - Assume n outcomes in the event space X
 - 1. Divide the interval $[0,1]$ into n intervals according to the probabilities of the outcomes
 - 2. Generate a **random** number r between 0 and 1
 - 3. Return x_i where r falls into



Generating text from language models

$$P(\text{of}) = 3/66$$

$$P(\text{Alice}) = 2/66$$

$$P(\text{was}) = 2/66$$

$$P(\text{to}) = 2/66$$

$$P(\text{her}) = 2/66$$

$$P(\text{sister}) = 2/66$$

$$P(\text{,}) = 4/66$$

$$P(\text{'}) = 4/66$$

Under a unigram language model:



Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

Generating text from language models

$$P(\text{of}) = 3/66$$

$$P(\text{Alice}) = 2/66$$

$$P(\text{was}) = 2/66$$

$$P(\text{to}) = 2/66$$

$$P(\text{her}) = 2/66$$

$$P(\text{sister}) = 2/66$$

$$P(\text{,}) = 4/66$$

$$P(\text{'}) = 4/66$$

Under a unigram language
model:

beginning by, very Alice but was and?
reading no tired of to into sitting
sister the, bank, and thought of without
her nothing: having conversations Alice
once do or on she it get the book her had
peeped was conversation it pictures or
sister in, 'what is the use had twice of
a book' 'pictures or' to

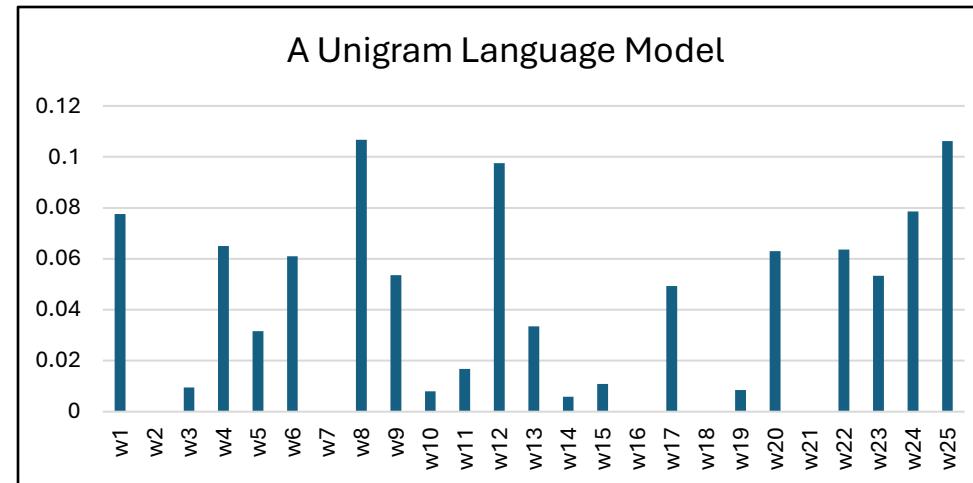


The same likelihood!

Recap: unigram language model

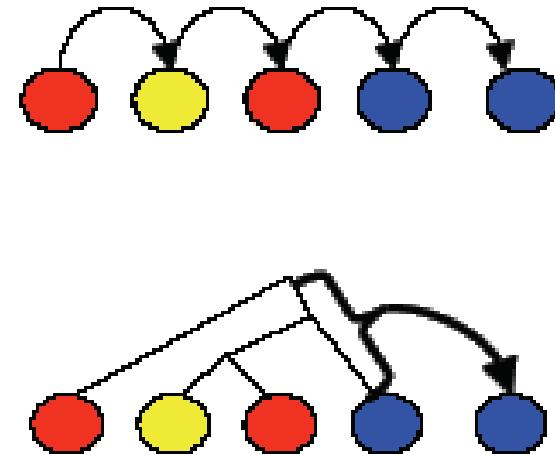
- Generate a piece of text by generating each word independently
 - $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2) \dots p(w_n)$
 - s.t. $\{p(w_i)\}_{i=1}^N, \sum_i p(w_i) = 1, p(w_i) \geq 0$
- Essentially a multinomial distribution over the vocabulary

*The simplest and
most popular
choice!*



Higher Order Models

- N-gram: condition on preceding words
- Grammar: condition on parse tree



Are they useful?

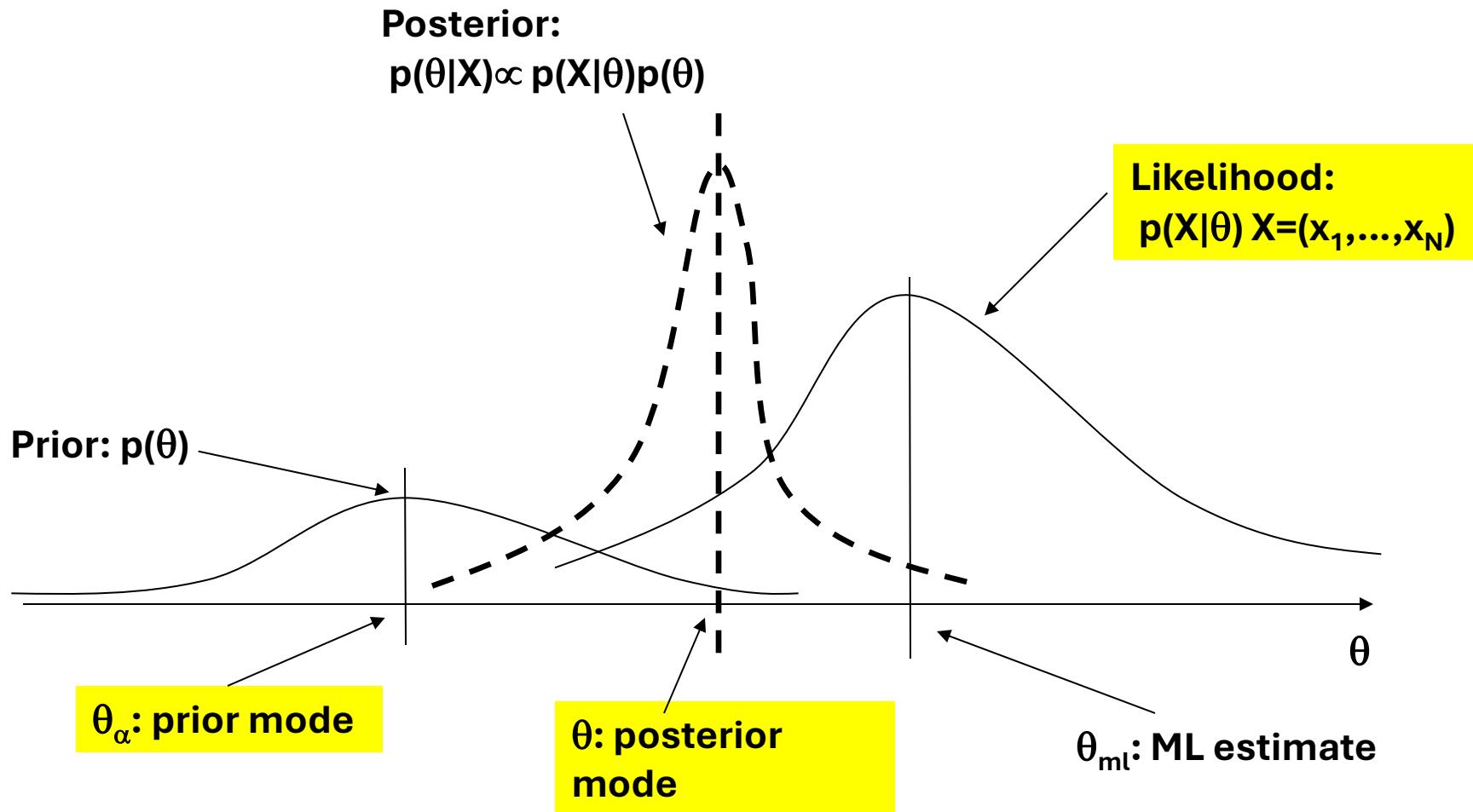
- No Significant improvements from higher models
- Parameter estimation is expensive

N-gram language models will help

Generated from language models of New York Times

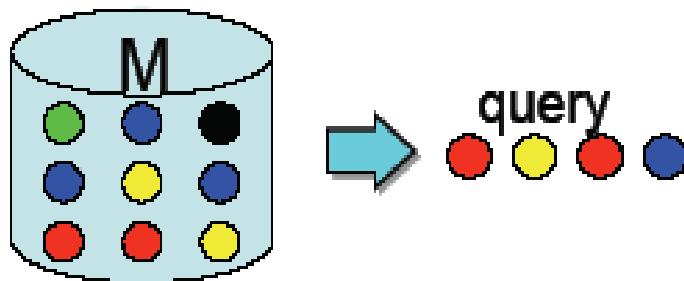
- Unigram
 - Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a q acquire to six executives.
- Bigram
 - Last December through the way to preserve the Hudson corporation N.B.E.C. Taylor would seem to complete the major central planners one point five percent of U.S.E. has already told M.X. corporation of living on information such as more frequently fishing to keep her.
- Trigram
 - They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions.

Illustration of Bayesian estimation



Multinomial Model (Unigram)

- Fundamental event: what is the identity of the i 'th query token?
- Observation is a sequence of events, each for every token



$$\Pr(q_1 \dots q_n \mid M) = \Pr(\text{len}(q)) = n \prod_{i=1}^n \Pr(q_i \mid M)$$

Multinomial Model (Cont...)

- Predominant model, well understood, lots of research
- Can account for multiple word occurrences in the query
- Same event space as NLP, MT (for integration)

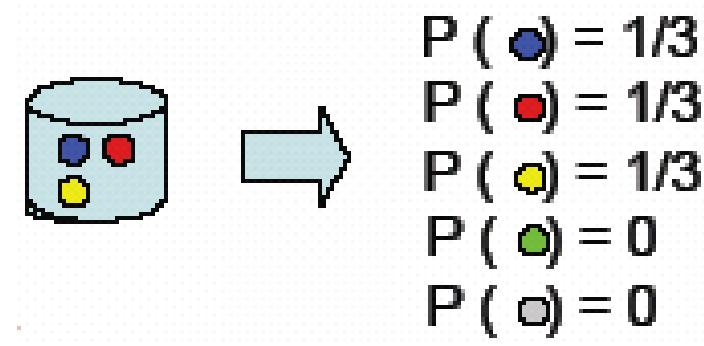
Model Estimation

- General problem:
 - Given a string S , estimate its language model M_S
 - S is commonly assumed to be an i.i.d random sample from M_S
- For Unigram we need to estimate $\Pr(w | M_S)$
- In IR we want to estimate M_q from q and/or M_d from d

Maximum Likelihood

- Count relative frequencies of words in sample S:

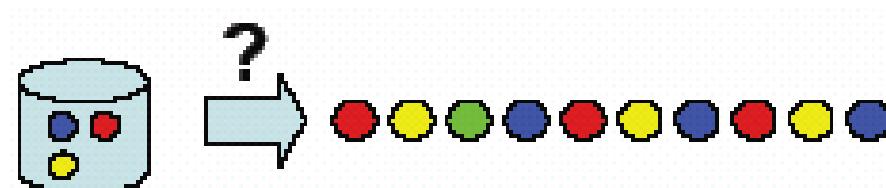
$$\Pr(w \mid M_S) = \frac{C(w, S)}{|S|}$$



- Maximum likelihood property
 - assigns highest possible probability to observation
- Unbiased estimator
 - If we repeat sampling infinite times, we get correct probabilities on average

The Zero Frequency Problem

- Suppose some event is not in observation S
 - Model will assign zero probability to the event
 - And to any set of events containing that event
- Happens very frequently with language
- It is incorrect to infer zero probabilities
 - Especially when creating a model from short samples



Discounting Methods

- Laplace correction
 - Add 1 to every count, normalize
 - Problematic for large vocabularies
- Lindstone correction
 - Add a small constant to every count, normalize

$$\Pr(w \mid M_S) = \frac{C(w, S) + \varepsilon}{|S| + \varepsilon|V|}$$

- Absolute discounting
 - subtract constant, redistribute the probability mass

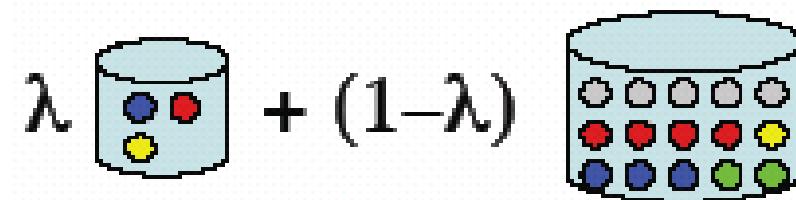
Interpolation Methods

- Problem with discounting methods
 - Discounting treats unseen words equally
 - Some words are frequent than others
- Idea: use background probabilities
 - Interpolate maximum likelihood estimates from S with general English expectations
 - In IR, background probabilities play the role of IDF

Jelinek-Mercer Smoothing

$$\Pr(w \mid M_S) = \lambda P_{ML}(w \mid d) + (1 - \lambda) P_{GE}(w)$$

- The same interpolation parameter is used for all documents and queries
- Tune parameter for optimal performance over the document collection



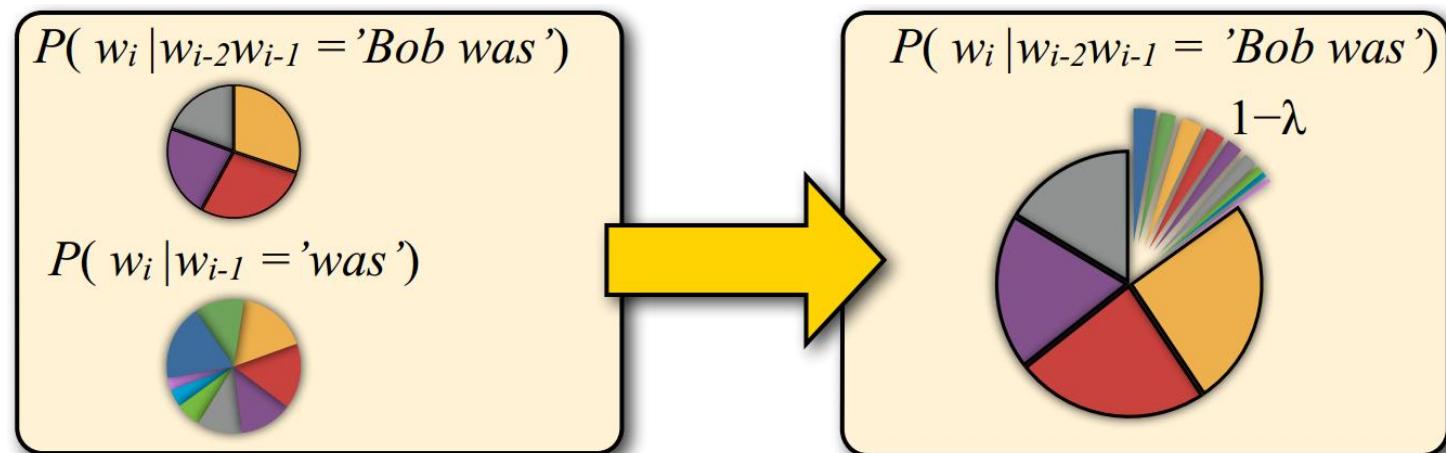
Dirichlet Smoothing

$$\Pr(w \mid M_S) = \frac{C(w, d) + \mu P_{GE}(w)}{|d| + \mu}$$

- The interpolation parameter μ can be viewed as the “size” of a “background document” added to the document
- Longer documents are less smoothed!
- Tune parameter for optimal performance over the document collection

Smoothing methods

- Linear interpolation
 - Use $(N - 1)$ -gram probabilities to smooth N -gram probabilities
 - We never see the trigram “Bob was reading”, but we do see the bigram “was reading”

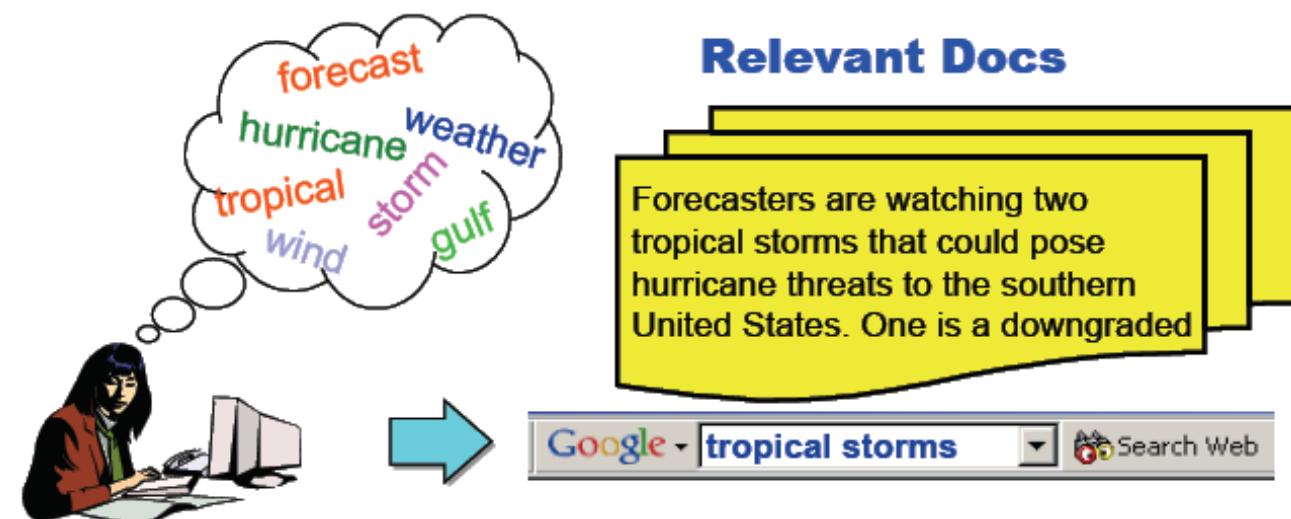


Major Issues in LM for IR

- What kind of language model to use
 - Linguistic complexity (unigram or higher order)
 - Statistical framework (multinomial or multiple-Bernoulli)
- How to estimate model parameters
 - Basic models
 - Translation models
- How to use model for ranking

How can we use Language Models in IR

- Task: given a query, retrieve relevant documents
- Use LM to model document generation
- Use LM to model query generation:
 - a. User thinks of some relevant document
 - b. User picks some keywords to use as query



Language Modeling for IR

- Every document d in a collection is generated by a language model M_d
 - Consider the language model as all possible strings that author could have written down when creating a document
 - For a specific author writing a specific document, some strings are more likely to occur than others
 - Topic, writing style, language (english, french...)
 - $\Pr(S | M_d)$ – the probability that author would write down string S

Language Modeling for IR (Cont...)

- Given a user query q , what is the probability that author would write down q ?
 - Higher probability captures a notion of relevancy of the document to the query
- Rank documents in a collection by $\Pr(q | M_d)$
 - Probability of observing q during random sampling from the language model of the document

Ranking with Language Models

- Standard approach – **query likelihood**:
 - Estimate a language model M_d for every document d in the collection
 - Rank documents by probability of “generating” the query:

$$\Pr(q_1 \dots q_n \mid M_d) = \prod_{i=1}^n \Pr(q_i \mid M_d)$$

Drawbacks of Query Likelihood

- No notion of relevance (everything is random sampling)
- User feedback and query expansion not part of the model
 - Examples of relevant documents cannot help improve M_d
 - Augmenting only by adding extra terms to the query

Ranking with Document Likelihood

- Flip the direction of the query-likelihood approach:
 - Estimate a language model M_q for every query q
 - M_q expected to “predict” a typical relevant document
 - Rank documents by the likelihood of being a random sample from M_q :

$$\Pr(d \mid M_q) = \prod_{w \in d} \Pr(w \mid M_q)$$

Problems of Document Likelihood

- Different document lengths, probabilities not comparable
- Favors documents that contain frequent (low content) words
- Consider ideal (“highest ranking”) document for a given query and a fixed document length n :

$$\max_d \prod_{w \in d} \Pr(w \mid M_q) = \max_w \Pr(w \mid M_q)^n$$

Ranking with Model Comparison

- Combine advantages of two ranking methods:
 - Estimate a model for both the document and the query
 - Directly compare similarity of the two models
- Natural measure of similarity is cross-entropy (others exist):

$$H(M_q \parallel M_d) = \sum_{w \in V} \Pr(w \mid M_q) \log(\Pr(w \mid M_d))$$

Ranking with Cross-Entropy

$$H(M_q \parallel M_d) = \sum_{w \in V} \Pr(w \mid M_q) \log(\Pr(w \mid M_d))$$

- Number of bits we need to “encode” M_q using M_d
- Equivalent to KL divergence: $H(M_q \parallel M_d) - H(M_q)$
- Equivalent to query likelihood if M_q is simply counts of words in q
- Cross-entropy is **not** symmetric

Variability problem with LM

- Basic language models do not address word synonymy
 - If a query q mentions ‘cars’, M_q should have high probability for ‘vehicles’, ‘autos’
- The same problem for cross-lingual retrieval
 - When searching for ‘cars’ in Spanish documents, it may be interesting to see documents containing ‘carros’ or ‘coches’

Translation Models

- Sample a word w from language model M with probability $P(w|M)$
 - Any of the models we have seen so far
- Translate w into a word v with probability $T(v|w)$

$$\Pr_{trans}(v | M) = \sum_{w \in V} T(v | w) \Pr(w | M)$$

- Translation probability reflects relationship between v and w

Using Translation Models

- Mono-lingual retrieval
 - v and w are synonyms
 - v and w are topically related, $T(v|w)$ is given by co-occurrence
 - $T(v|w)$ high if v is likely to appear in documents containing w
- Stemming
 - v and w are in the same stem class
- Cross-lingual retrieval
 - v is the translation of w into a different language

Perplexity

- Standard evaluation metric for language models
 - A function of the probability that a language model assigns to a data set
 - Rooted in the notion of cross-entropy in information theory

Perplexity

- The inverse of the likelihood of the test set as assigned by the language model, normalized by the number of words

$$PP(w_1, \dots, w_N) = \sqrt[N]{\frac{1}{\prod_{i=1}^N p(w_i | w_{i-1}, \dots, w_{i-n+1})}}$$

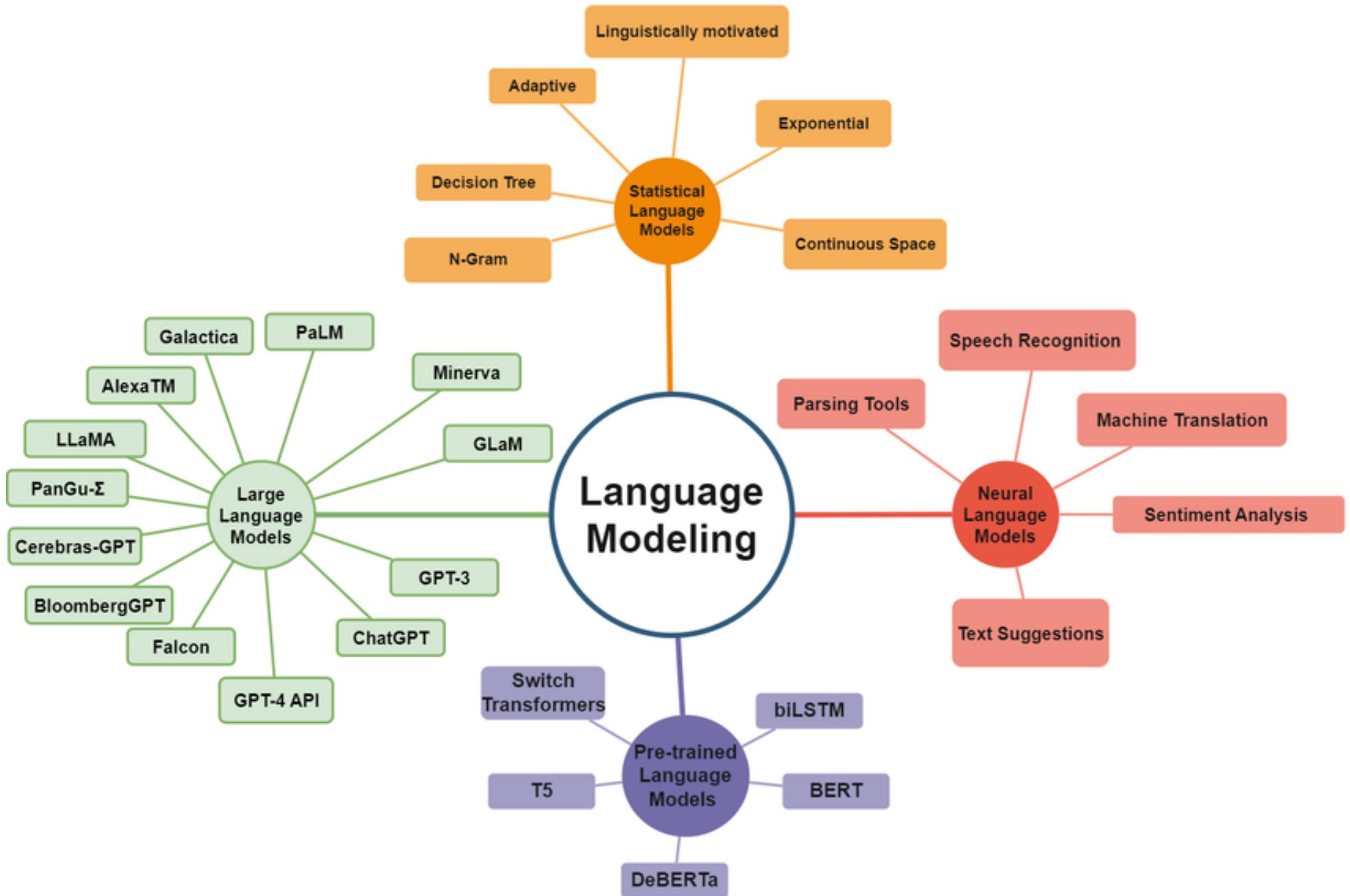
N-gram language model

Common Practice

- What is the main-stream:
 - Unigram models
 - Multinomial models
 - Model estimation by MLE smoothed by interpolation with background probability
 - Ranking by Model comparison (cross-entropy)

Language Modeling: pros & cons

- Pros:
 - Formal mathematical model
 - Simple, well-understood framework
 - Integrates both indexing and retrieval models
 - Natural use of statistics, no heuristics
 - Avoids problematic issues of “relevance”
- Cons:
 - Difficult to incorporate notions of “relevance” and user preference
 - Relevance feedback and query expansion is not straightforward
 - Accommodating phrases, paragraphs and Boolean operators not straightforward



One-Hot Encoding?

- If we have a sequence of T words, and each word gets a one-hot encoded vector of size V , then one sentence becomes a $T \times V$ matrix
- Same generic shape as a $T \times D$ matrix (as you saw previously)
- E.g. “I like cats”
- Might become: [[0, 0, 0, 1], [0, 1, 0, 0], [1, 0, 0, 0]]
- (If we had a 4-word vocabulary)

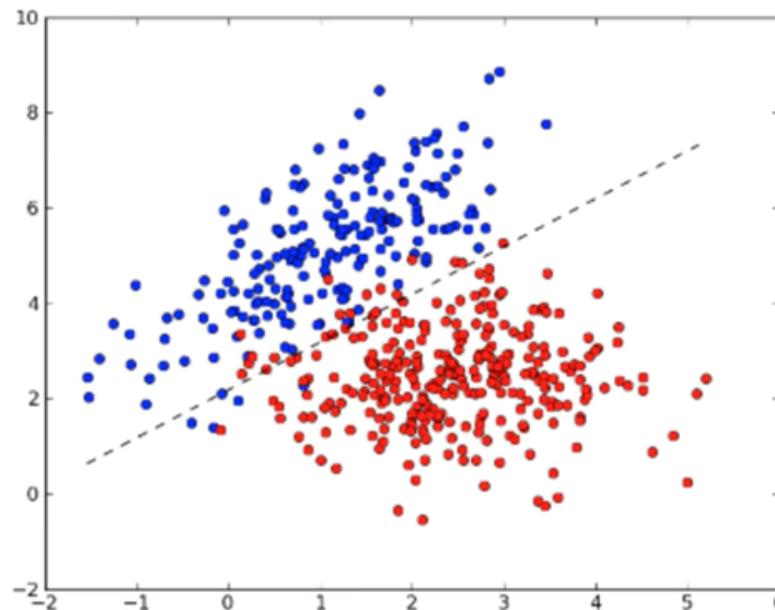
Problem!

- Some English-language datasets have ~1 million possible tokens / words
- This means V (or equivalently, D) is 1 million
- Not only is the input large, but the input-hidden weight matrix will also be large!
- In actuality, the English dictionary has ~200,000 words, but this is still a very large number

	o_ENE	o_ESE	o_East	o_NE	o_NNNE	o_NNW	o_NW	o_North	o_SE	o_SSE	o_SSW	o_SW	o_South	$o_Variable$	o_WSW
0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
2	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
4	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

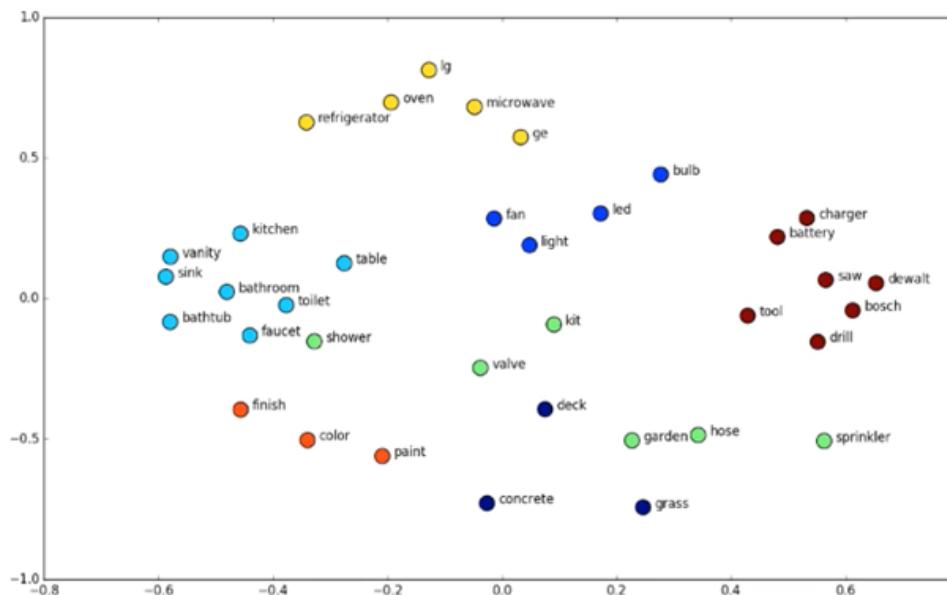
Another problem!

- Recall: we want our input features to have some structure
- “Machine learning is nothing but a geometry problem”
- This rule **relies** on the fact that there’s some geometrical structure in the data - that data vectors in the same class (e.g. spam vs. not spam) are *probably* near each other



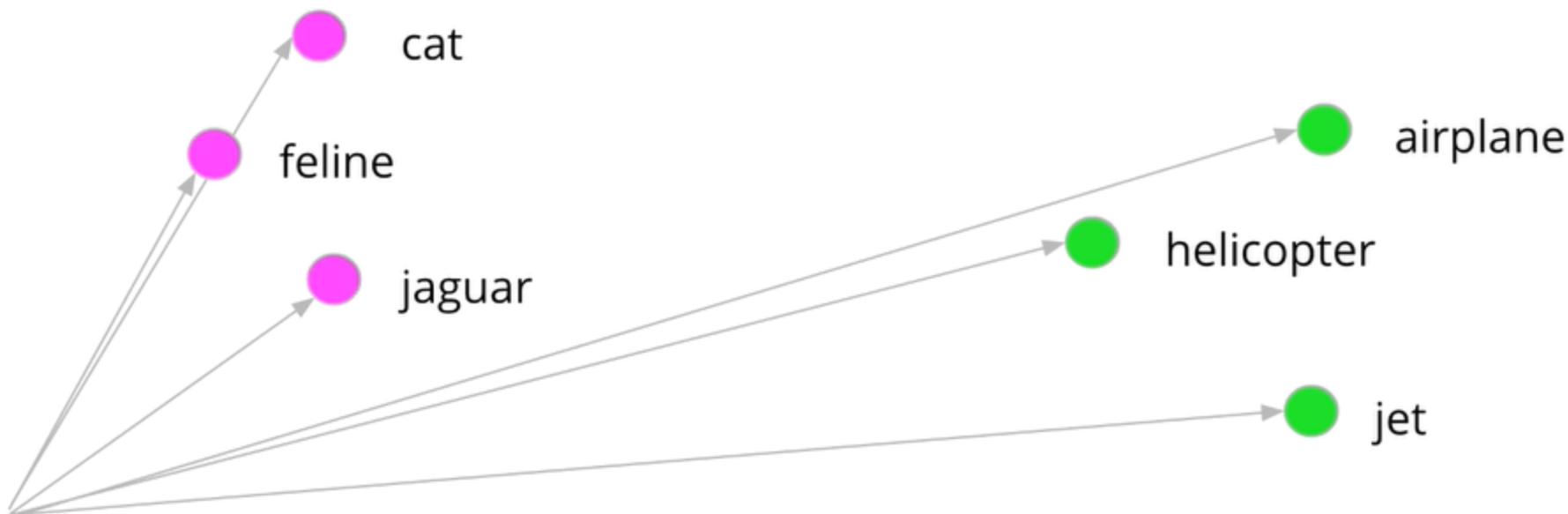
Problem with one-hot encoding

- Take any two one-hot encoded vectors, e.g. [1, 0, 0] vs. [0, 1, 0]
- The Euclidean distance is $\sqrt{1^2 + 1^2} = \sqrt{2}$!
- All words therefore are an equal distance apart!
- The distance between “cat” and “feline” is $\sqrt{2}$
- The distance between “cat” and “airplane” is *also* $\sqrt{2}$
- Useless!



A better solution: Embeddings

- Assign each word to a D-dimensional vector (not one-hot encoded)
- How are these vectors “found”? We will see shortly!



$$\begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline \end{array}$$

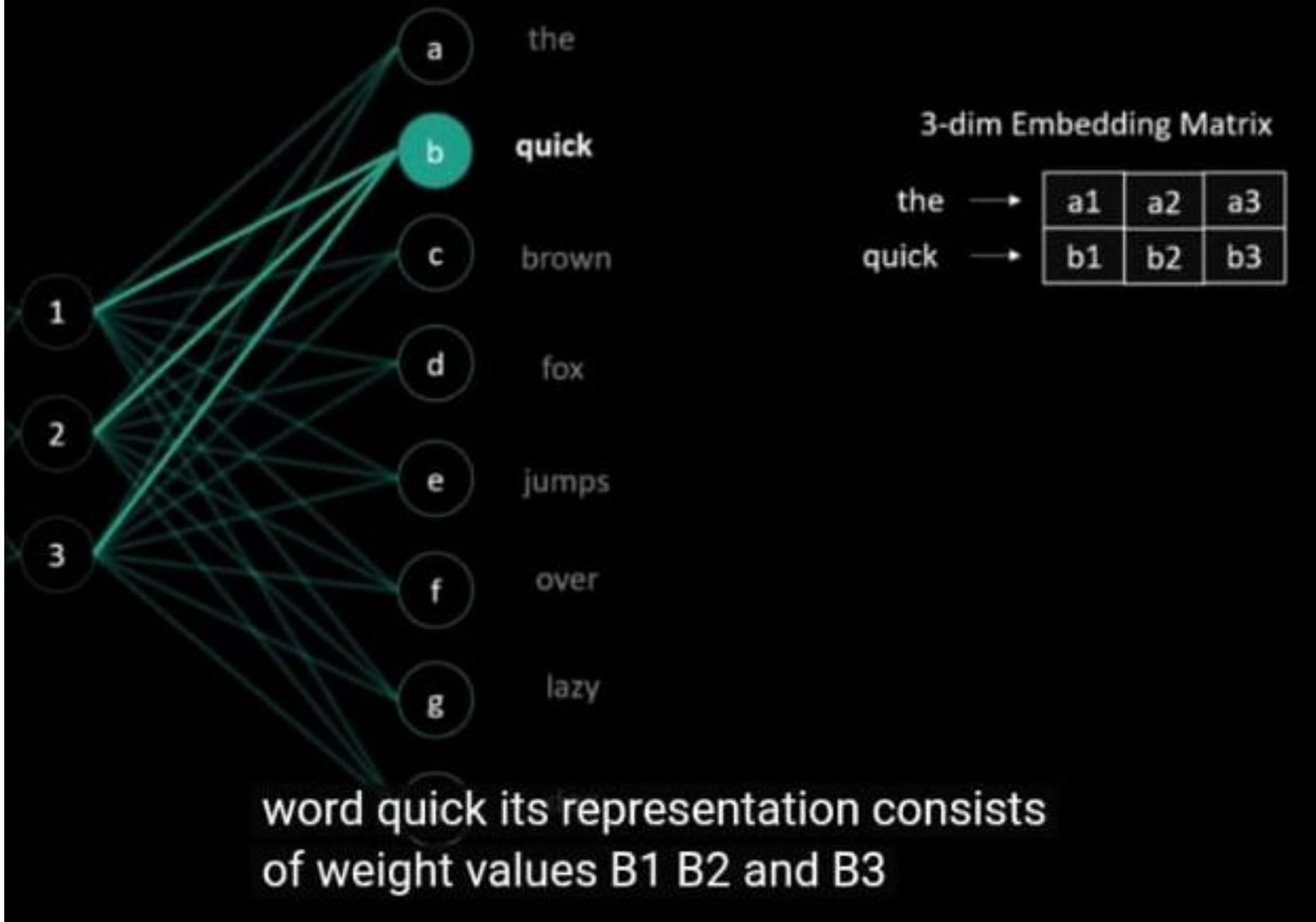
$$\begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 4 & 5 & 6 \\ \hline \end{array}$$

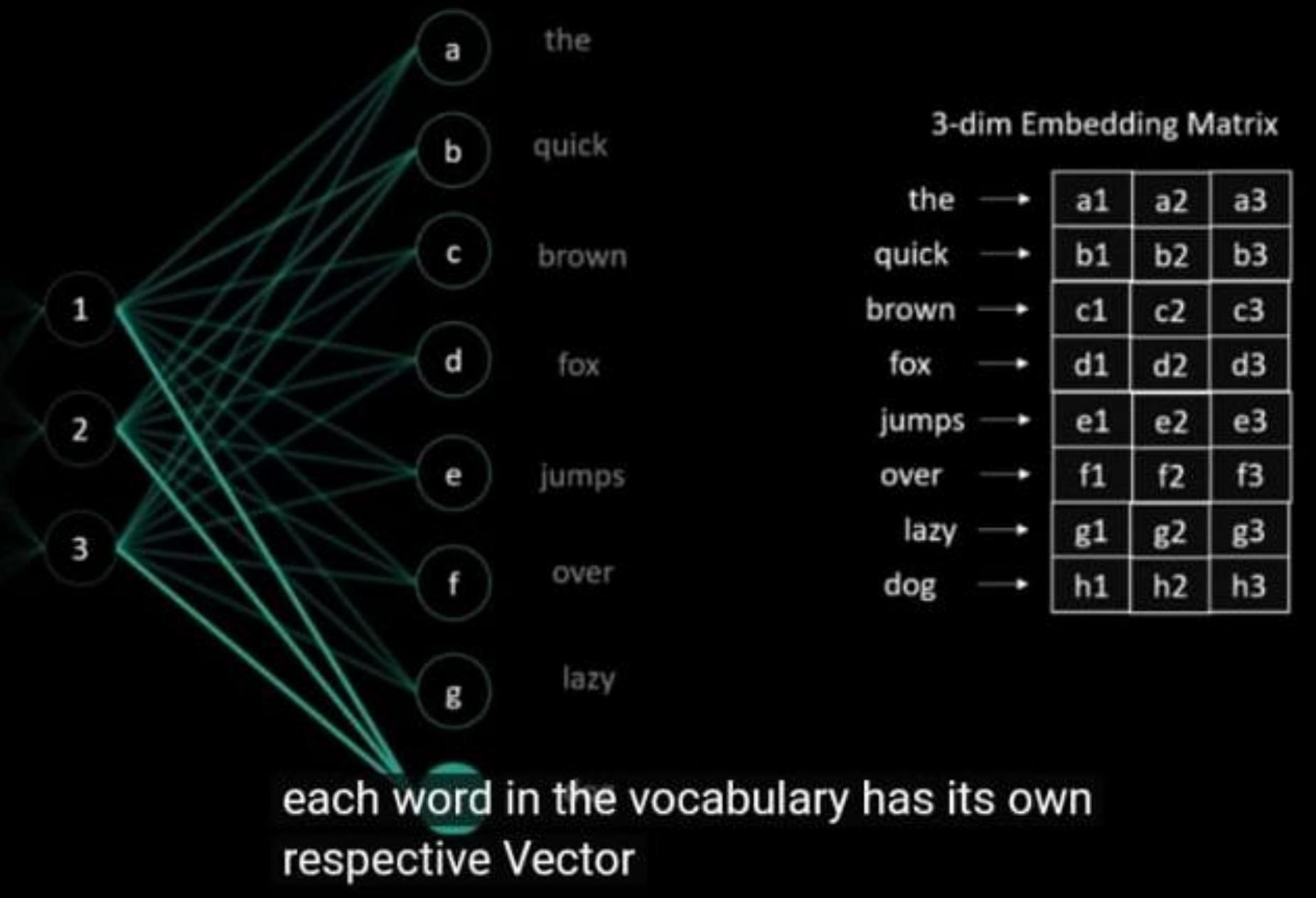
Why is it more efficient?

- Old way takes 2 steps:
 - Create a vector of size V containing all zeros, set the k'th entry to 1
 - Multiply the one hot vector by the weight matrix
- New way takes 1 step:
 - Index the weight matrix ($W[k]$)

$$\begin{array}{|c|c|c|} \hline 1 & 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline \end{array}$$

The diagram illustrates matrix multiplication. On the left is a 1x3 vector [1 0 0]. In the center is a 3x3 weight matrix [1 2 3; 4 5 6; 7 8 9]. A green box highlights the first row of the matrix. An equals sign follows the multiplication symbol. On the right is the resulting 1x3 vector [1 2 3], which is the first row of the weight matrix.





[1 0 0 0 0 0 0]

1x8



a1	a2	a3
b1	b2	b3
c1	c2	c3
d1	d2	d3
e1	e2	e3
f1	f2	f3
g1	g2	g3
h1	h2	h3

8x3



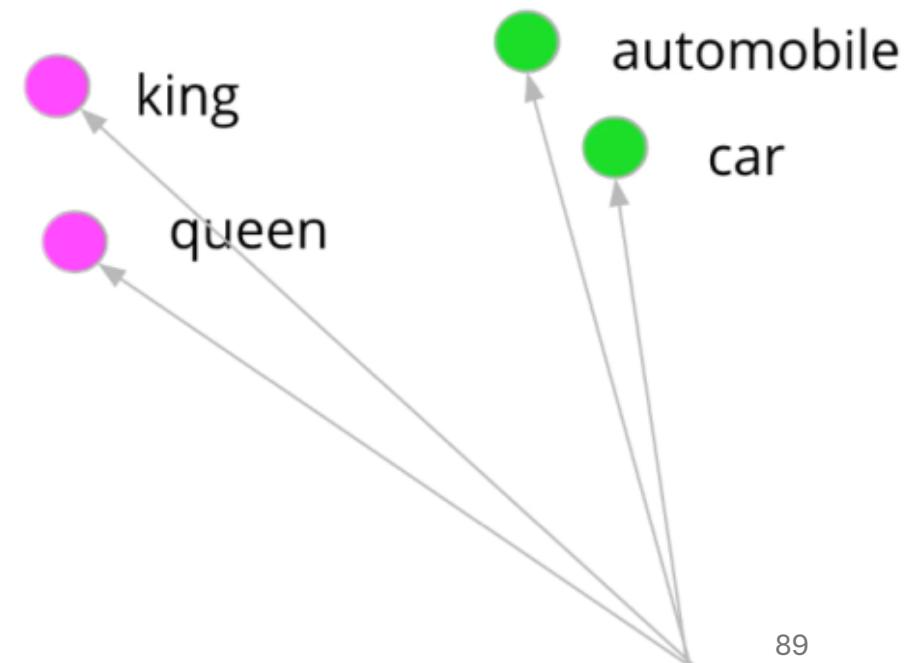
[a1 a2 a3]

Tensorflow Embedding

- This is exactly what the Embedding layer does!
- STEP #1: Convert words into integers
 - $["I", "Like", "Cats"] \rightarrow [50, 25, 3]$
- STEP #2: Use integers to index the word embedding matrix to get word vectors for each word
 - $[50, 25, 3] \rightarrow [[0.3, -0.5], [1.2, -0.7], [-2.1, 0.9]]$
 - T-length array $\rightarrow T \times D$ matrix

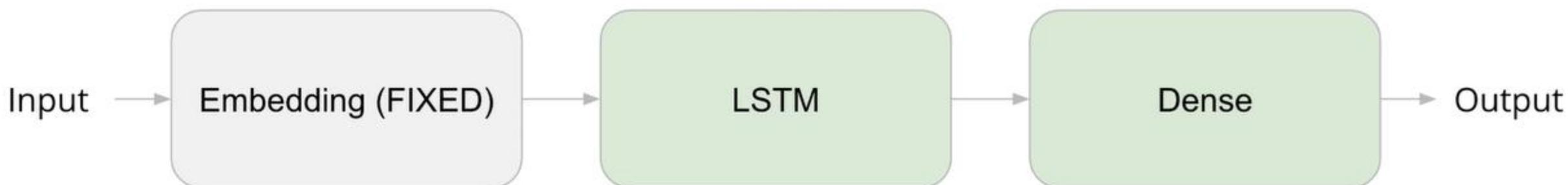
How are the weights found?

- We know intuitively that the word vectors must have some useful structure
- Luckily, it's the same story as with convolutional filters
- These weights are found *automatically* with gradient descent when you call `model.fit()`



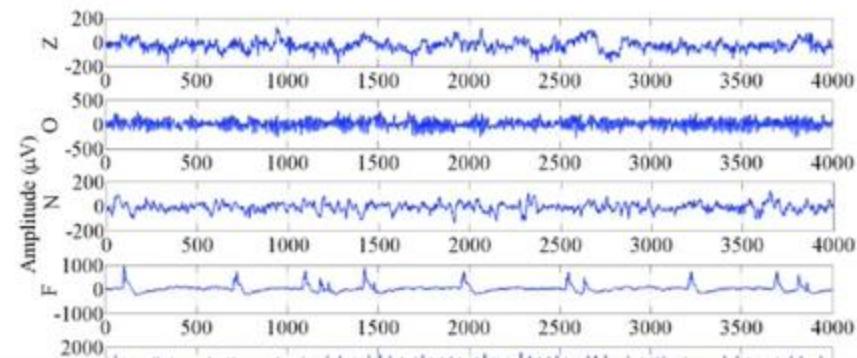
Pre-trained word vectors

- Sometimes, we use *pre-trained* word vectors (trained using some other algorithm)
- We freeze (fix) the embedding layer's weights, so only the *other* layers are trained when we call `model.fit()`
- You're encouraged to read about word2vec and GloVe to learn more



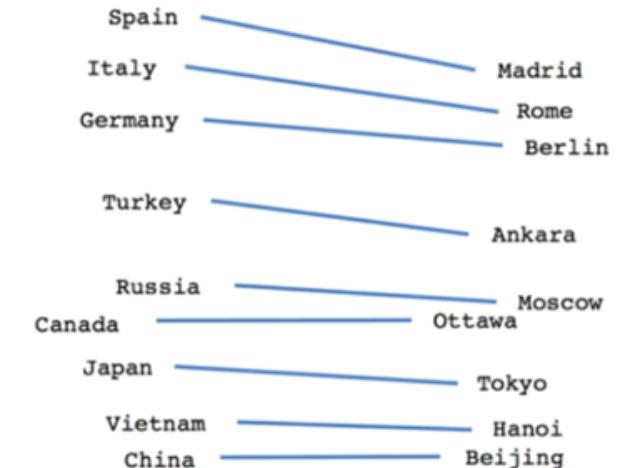
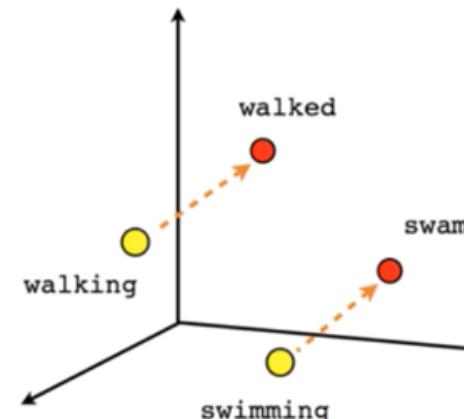
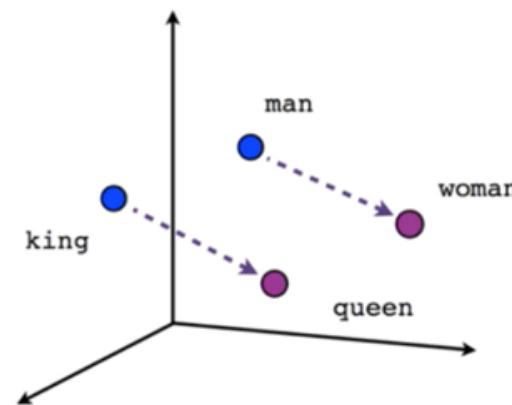
Summary

- It's not possible to multiply a *word* by a weight matrix - there is no concept of multiplying words
- We must turn words into some kind of numerical representation first
- One-hot encoding is a fine first guess
 - Too much space (there could be up to 1 million tokens)
 - Not geometrically useful ("cat" is just as close to "feline" as it is to "airplane")
 - We need to make use of "machine learning is nothing but a geometry problem"
- Better: convert words to unique integers
- Use those integers to *index* a weight matrix (the *embedding matrix*)
- Much faster than matrix multiplication



Summary

- Each input sentence thus becomes a $T \times D$ sequence which we know an RNN can accept
- The embedding is trained like any other layer - we just need to call `model.fit()`
- Read up on word2vec and GloVe for alternative methods for training embeddings

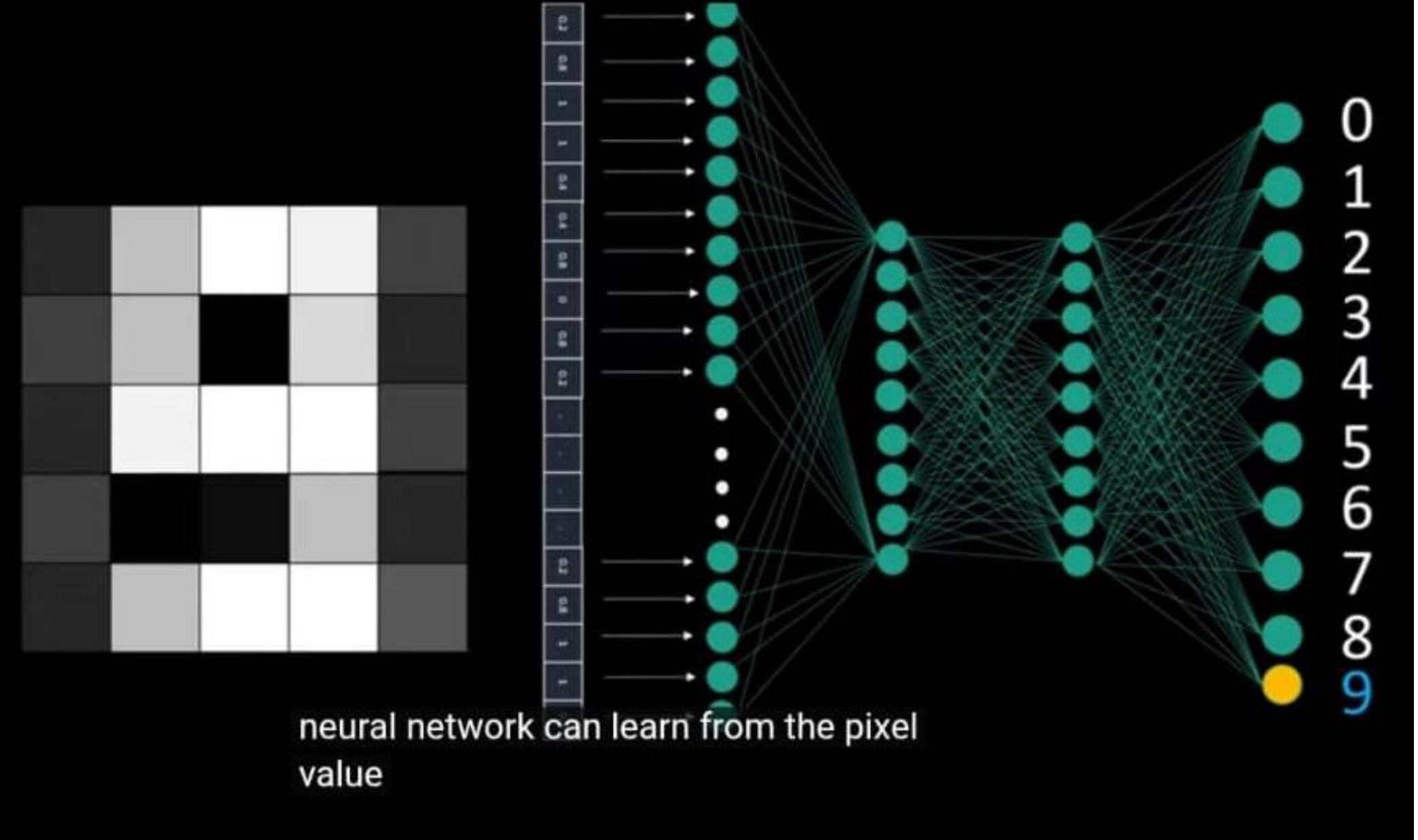


GRAYSCALE IMAGE OF DIGIT 9



NUMERICAL REPRESENTATION

0.2	0.8	1	1	0.4
0.4	0.8	0	0.8	0.2
0.2	0.9	1	1	0.4
0.8	0	0.3	0.8	0.3
0.2	0.8	1	1	0.4



ONE HOT ENCODING

OHE Avoids Numeric Relations

Simple to use

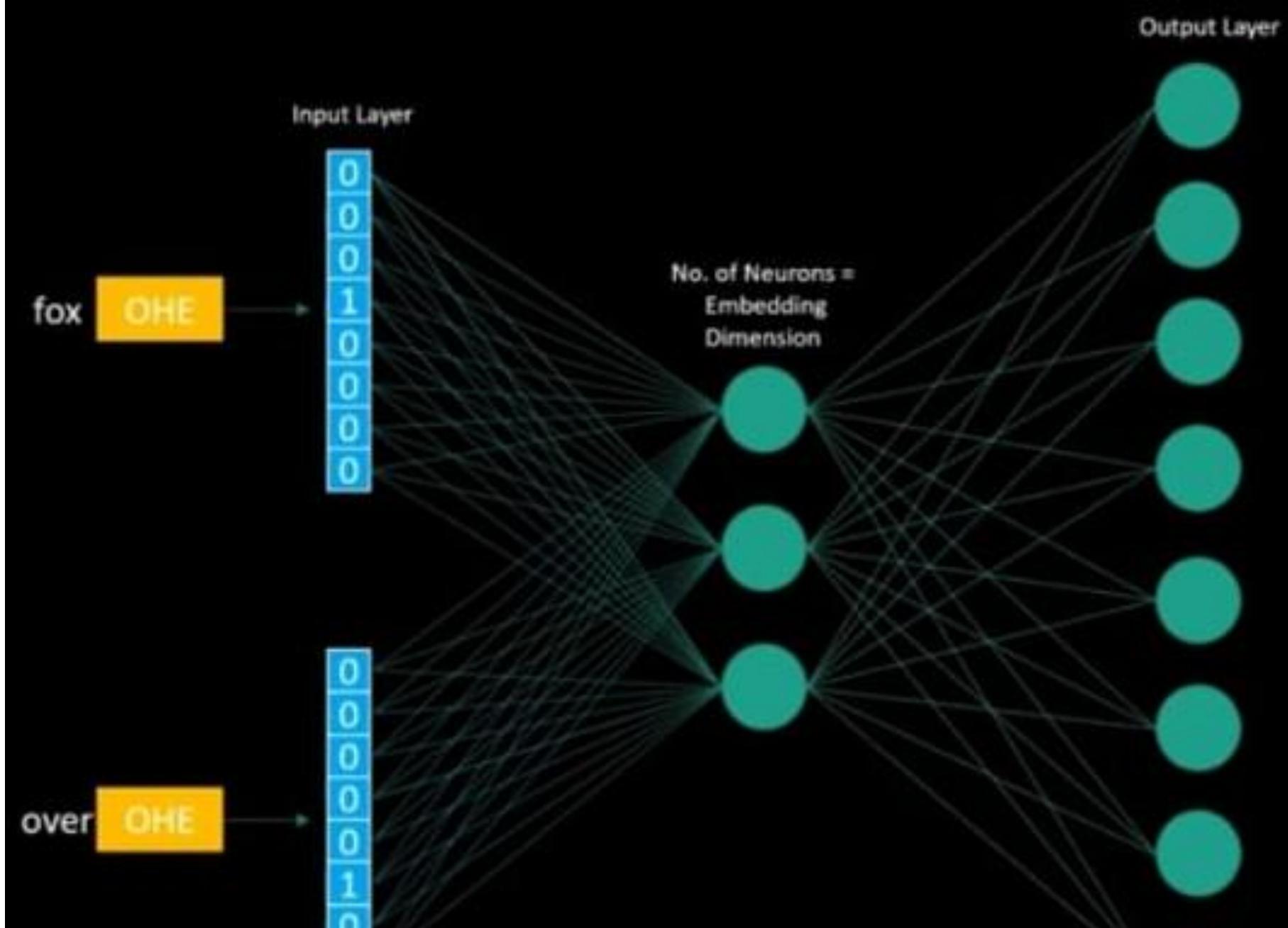
It is Sparse

n-size vector for n-vocab size

Cannot capture semantic meaning and contextual understanding

For GPT-2 50257 vocab

Large Sparse Matrix of 50257 Dimension



Sentence 1: This movie was good
Sentence 2: This movie was bad

Vocabulary: ["This", "movie", "was", "good", "bad"]
Vocabulary Size: 5

ONE HOT ENCODING TABLE

	This	movie	was	good	bad
This	1	0	0	0	0
movie	0	1	0	0	0
was	0	0	1	0	0
good	0	0	0	1	0
bad	0	0	0	0	1

Sentence 1: This movie was good
Sentence 2: This movie was bad



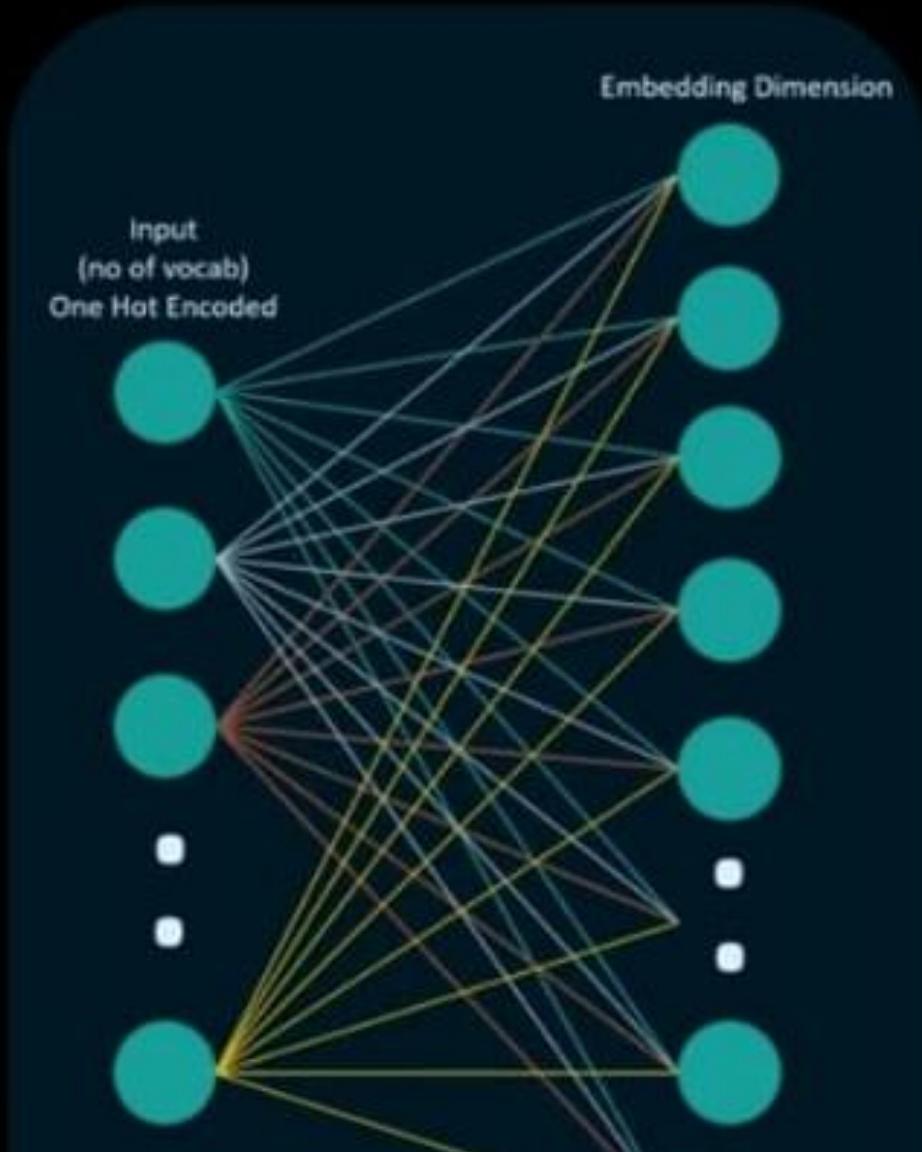
	This movie	movie was	was good	was bad
Sentence 1	1	1	1	0
Sentence 2	1	1	0	1

Bi-Gram (2 word context)

N-Gram with N-Word Context

Sparsity problem in large corpora where there are many possible n-gram

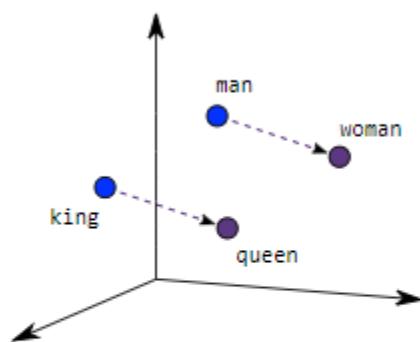
791
8415
7731
389
264



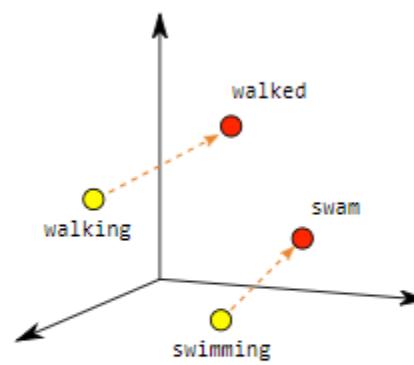
Embedding (n-dim)					
0.012	0.819	1.190	1.129	0.412	-
0.492	0.829	0.123	0.886	0.2.21	-
0.220	0.912	1.869	1.986	0.419	-
0.812	0.198	0.310	0.889	0.398	-
0.2	0.8	0.129	0.689	0.198	-
0.2	0.8	0.896	0.861	0.411	-

(no_of_tokens x Embeddding_dim)

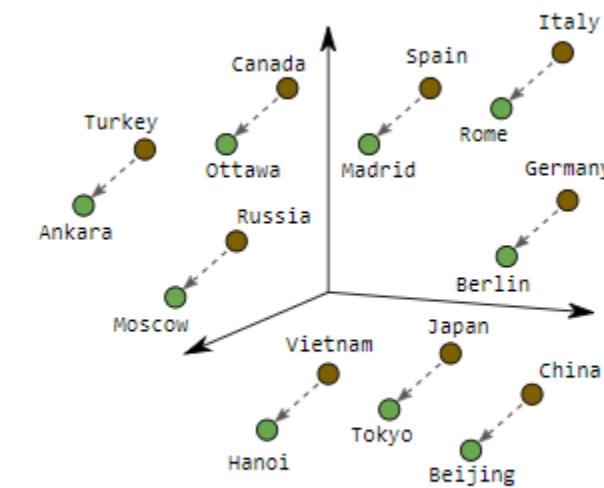
Word Embeddings



Male-Female



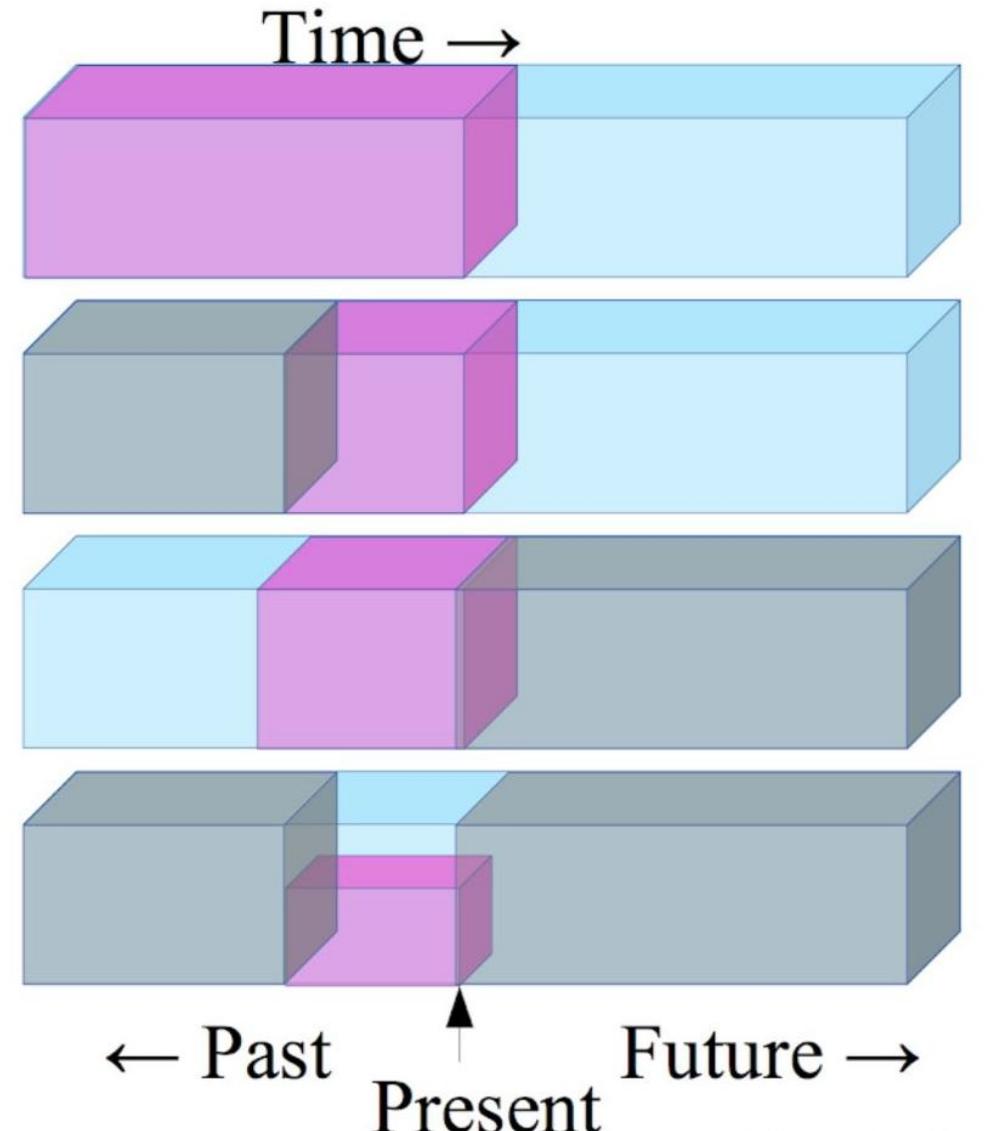
Verb Tense



Country-Capital

Predicting neighbouring context

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ Pretend there is a part of the input you don't know and predict that.



Slide: LeCun

► “Pure” Reinforcement Learning (**cherry**)

- The machine predicts a scalar reward given once in a while.

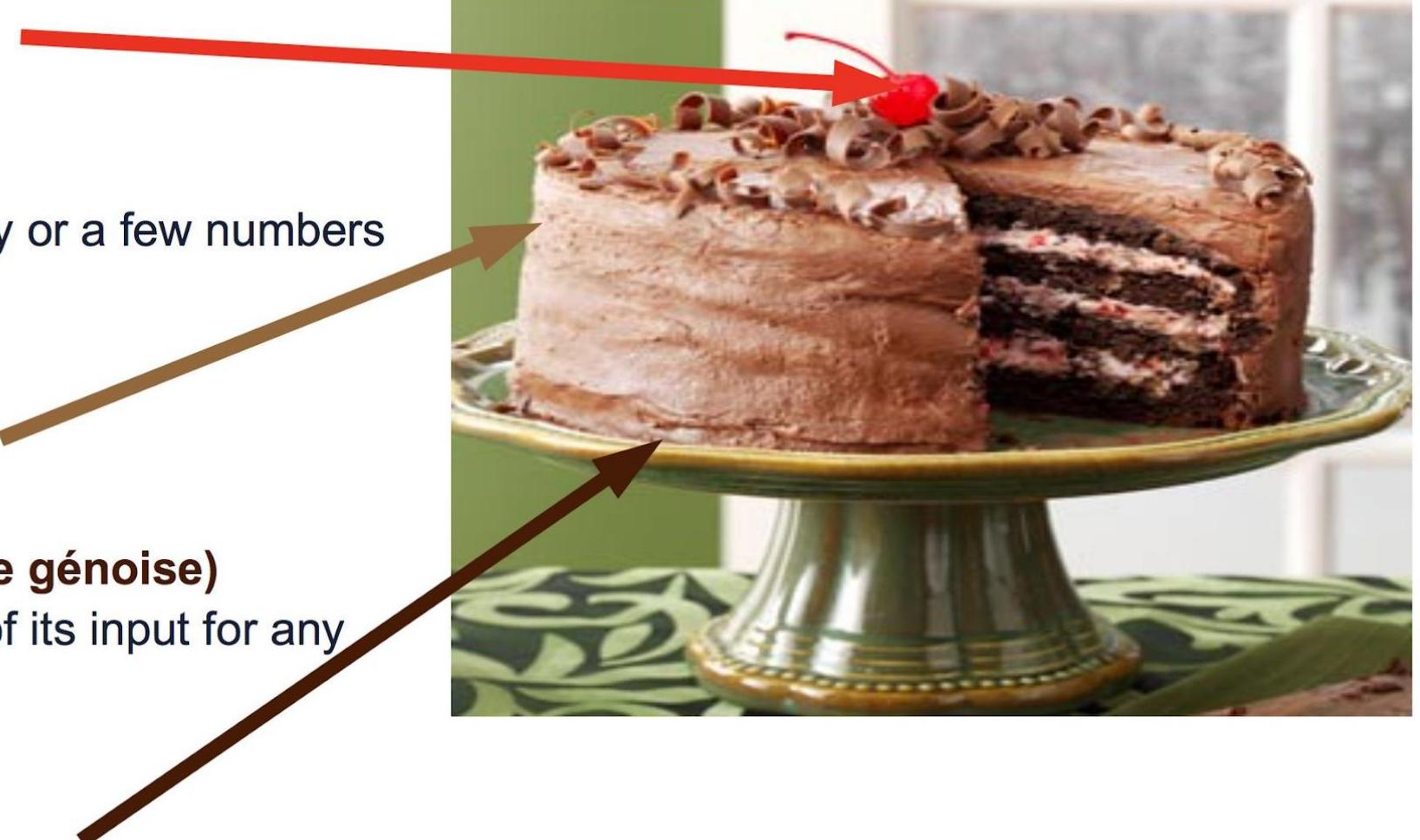
► **A few bits for some samples**

► Supervised Learning (**icing**)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- **10→10,000 bits per sample**

► Self-Supervised Learning (**cake génoise**)

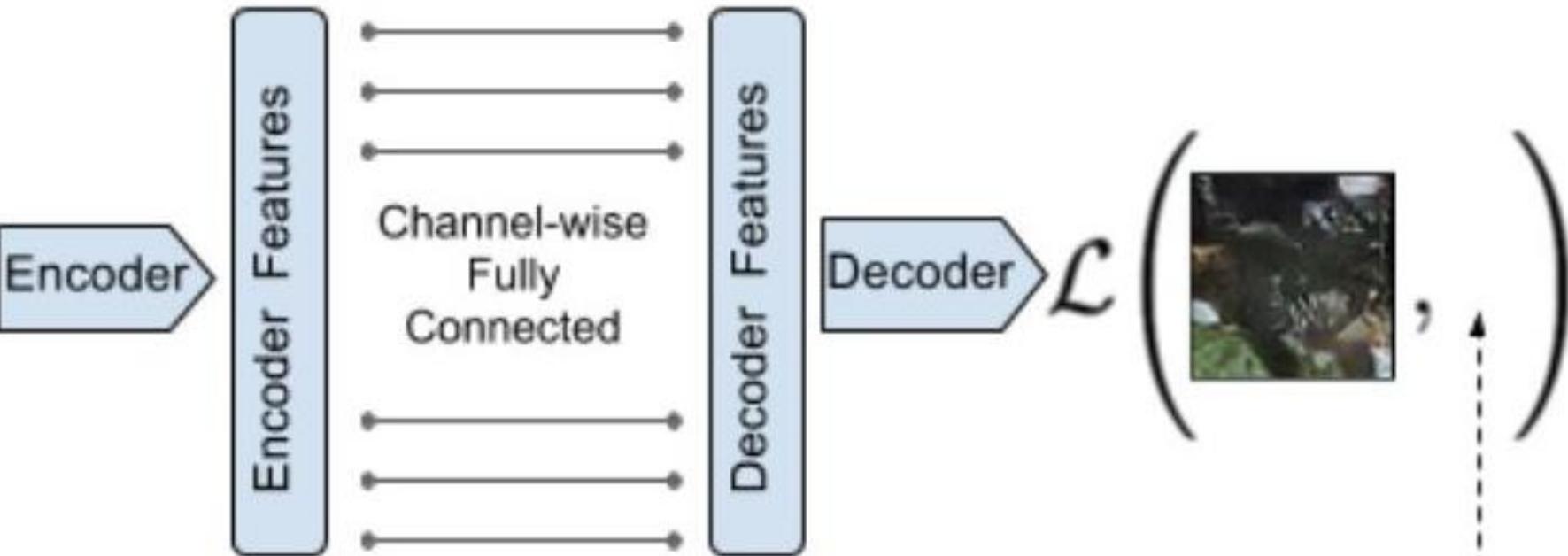
- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos
- **Millions of bits per sample**

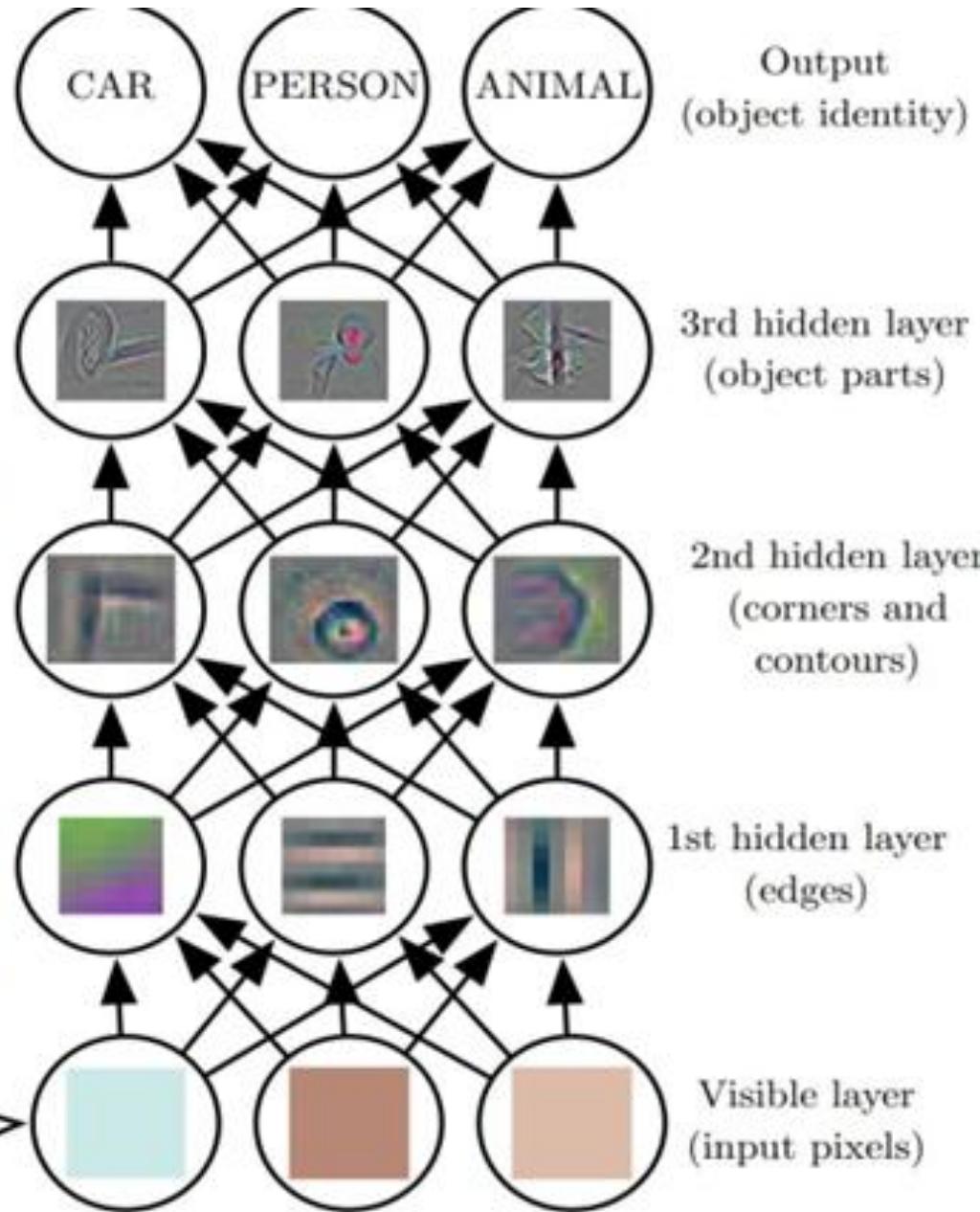


Predict missing pieces



Context Encoders





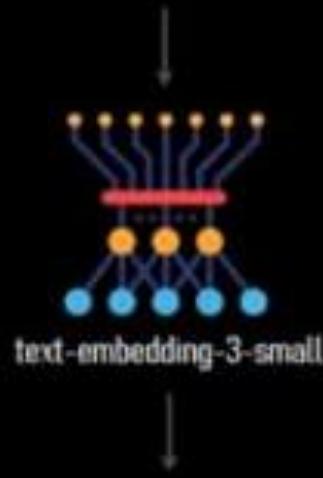
Word

Cat.



Text Paragraph

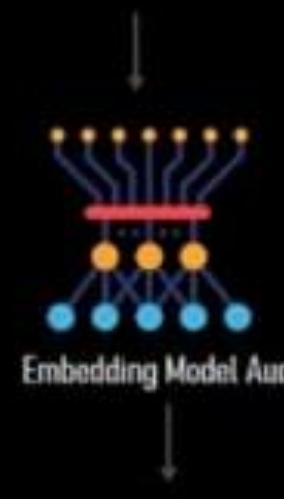
Cats are small, carnivorous mammals often kept as pets. They are known for their agility, sharp claws, and playful behavior.



Image



Audio



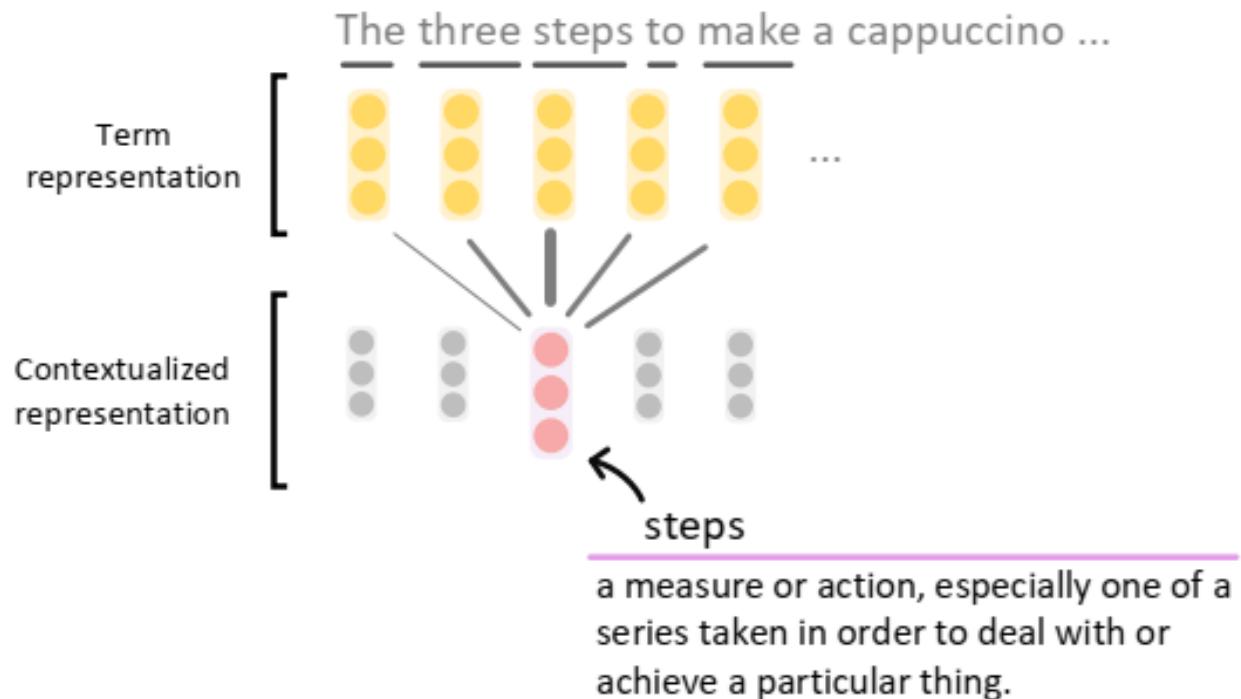
[0.01 0.2 0.001 0.12, ...]

[0.5 0.05 0.1 0.5, ...]

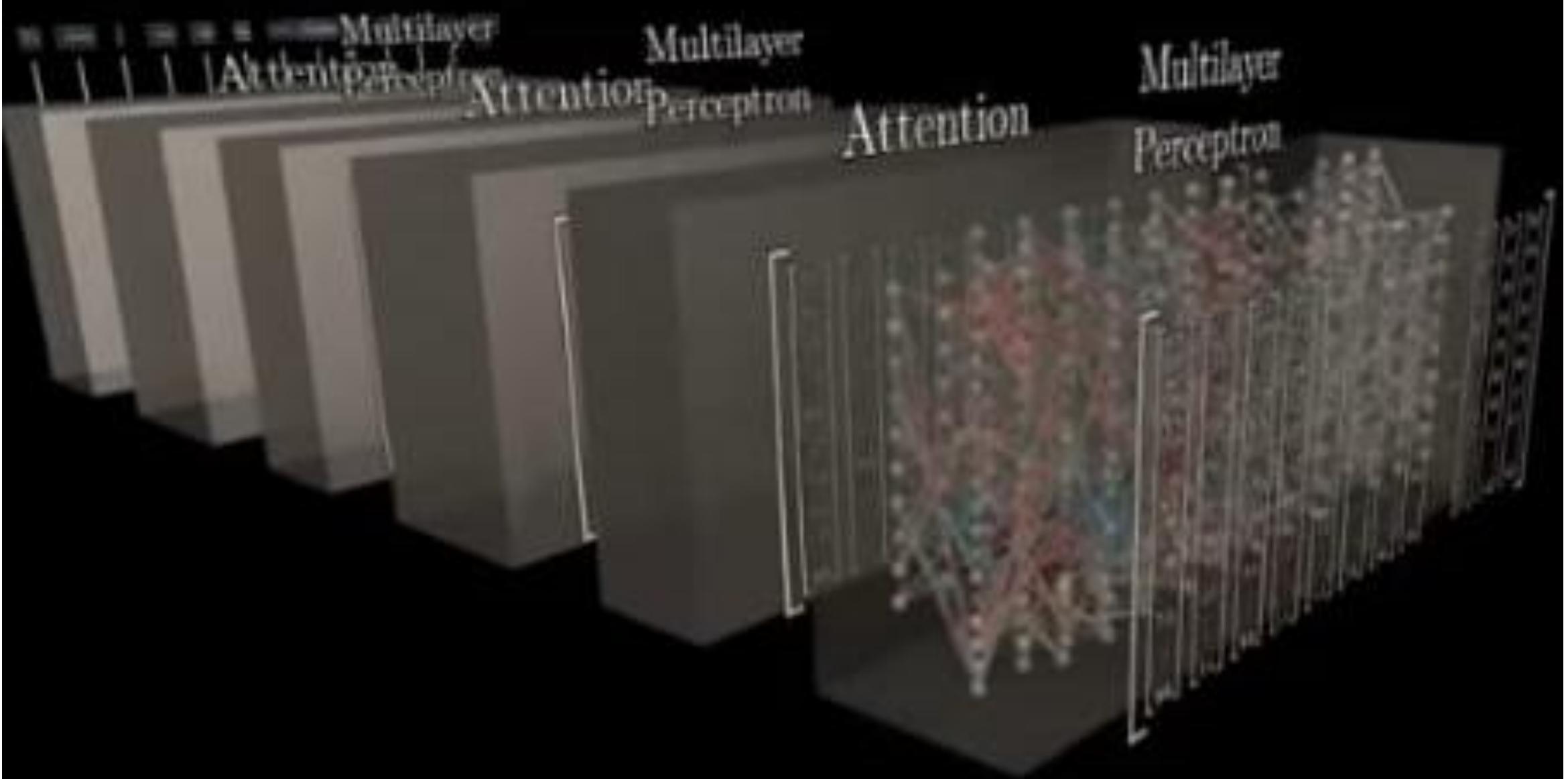
[0.15 0.2 0.001 0.12, ...]

[0.06 0.006 0.015 0.006, ...]

Contextualization via Self-Attention



- Learn meaning based on surrounding context for every word occurrence
- This *contextualization* combines representations
- Context here is local to the **sequence** (not necessarily a fixed window)
- Is computationally intensive $O(n^2)$
 - Every token attends to every other token



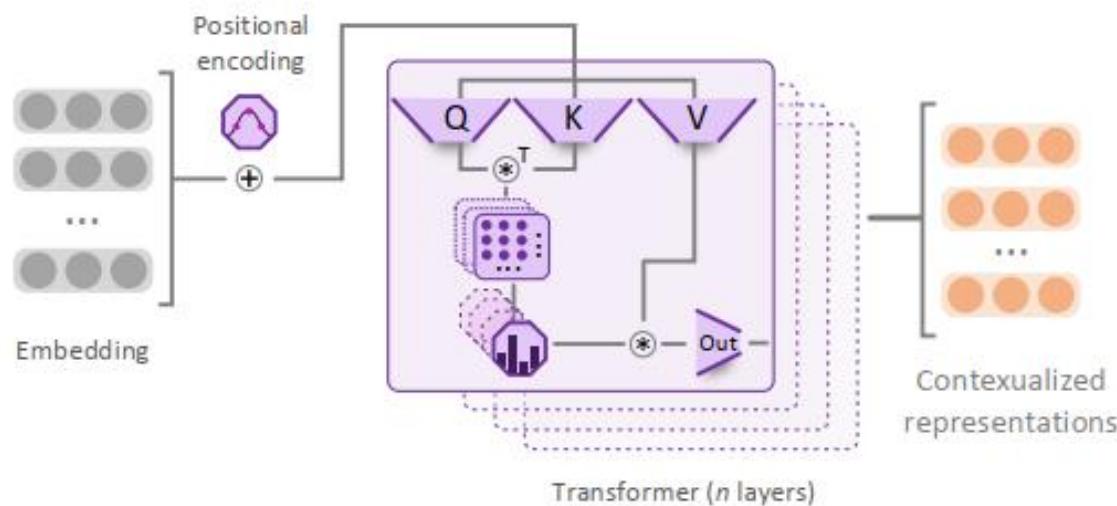
Transformer

- Transformers contextualize with multi-head self-attention
 - Every token attends to every other token $O(n^2)$ complexity
- Commonly Transformers stack many layers
- Can be utilized as encoder-only or encoder-decoder combination
- Do not require any recurrence
 - The attention breaks down to a series of matrix multiplications over the sequence
- Initially proposed in translation
 - Now the backbone of virtually every NLP advancement in the last years

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, et al.

Attention is all you need. In NeurIPS. 2017.

Transformer – Architecture



- We embed (subword) tokens
- We add a positional encoding
- In each Transformer-Layer:
 - Project each vector with 3 linear layers to **Query**, **Key**, **Value**
 - Transform projections to another multi-head dimension
 - Matrix-multiply Query & Key
 - Get Q-K attention via softmax
 - Multiply attention with Values and project back to output

Nice detailed walkthrough code + paper:

<https://nlp.seas.harvard.edu/2018/04/03/attention.html>

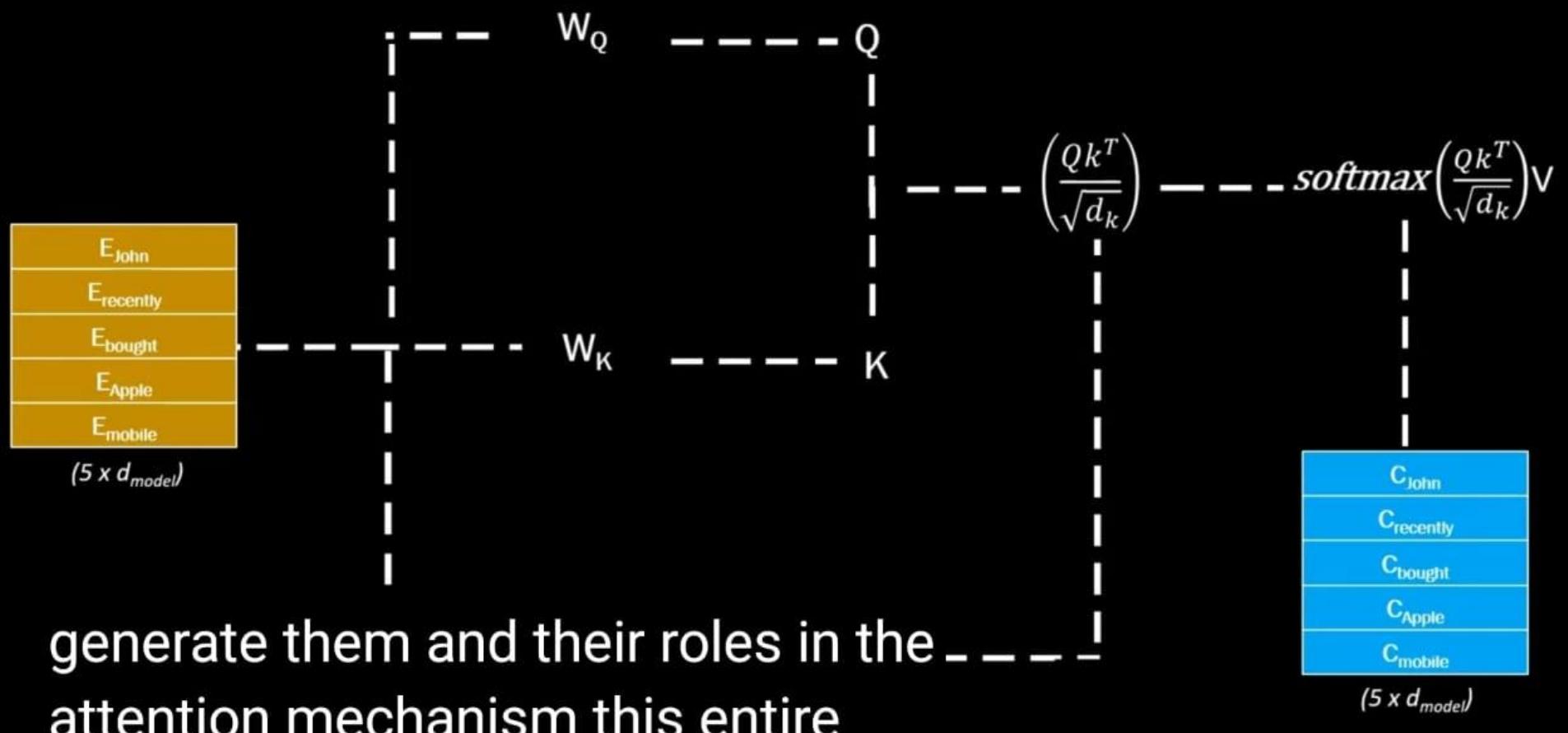
Self-Attention Definition

- The Transformer Self-Attention is defined as:

$$\text{SelfAttention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) * V$$

- Q, K, V are projections of the **same input** sequence
- This definition hides quite a bit of complexity, visible in the code

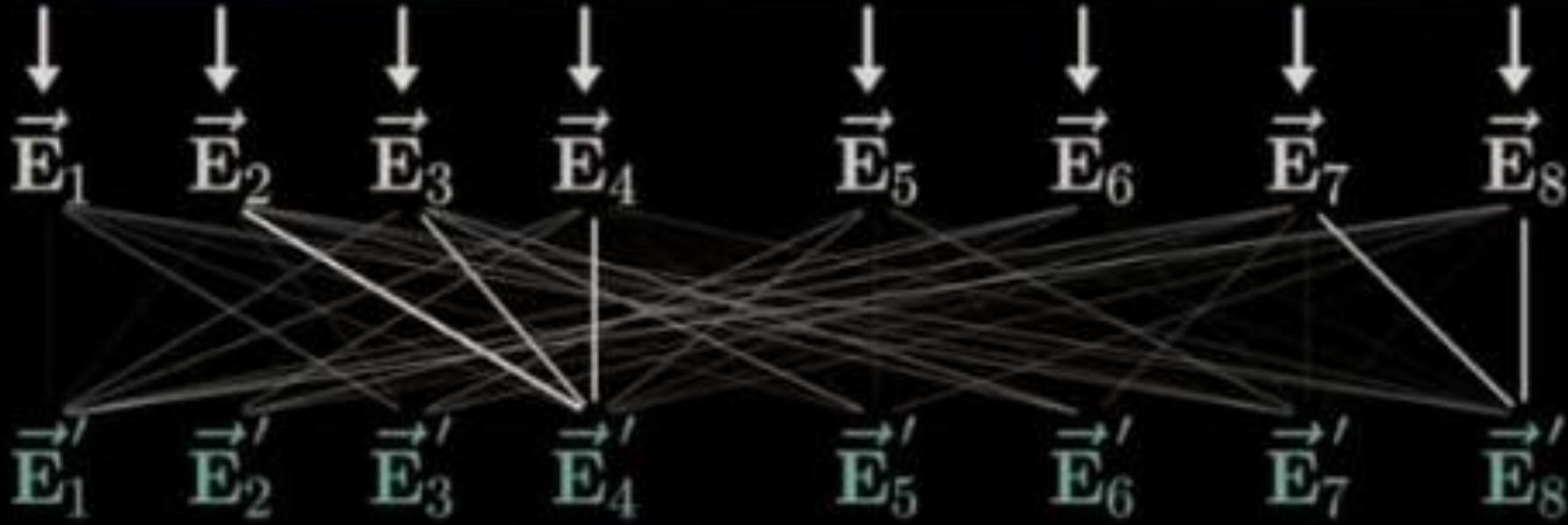
Q	Attention “Query”
K	Attention “Key”
V	Attention “Value”
d_k	Dimension of key embeddings



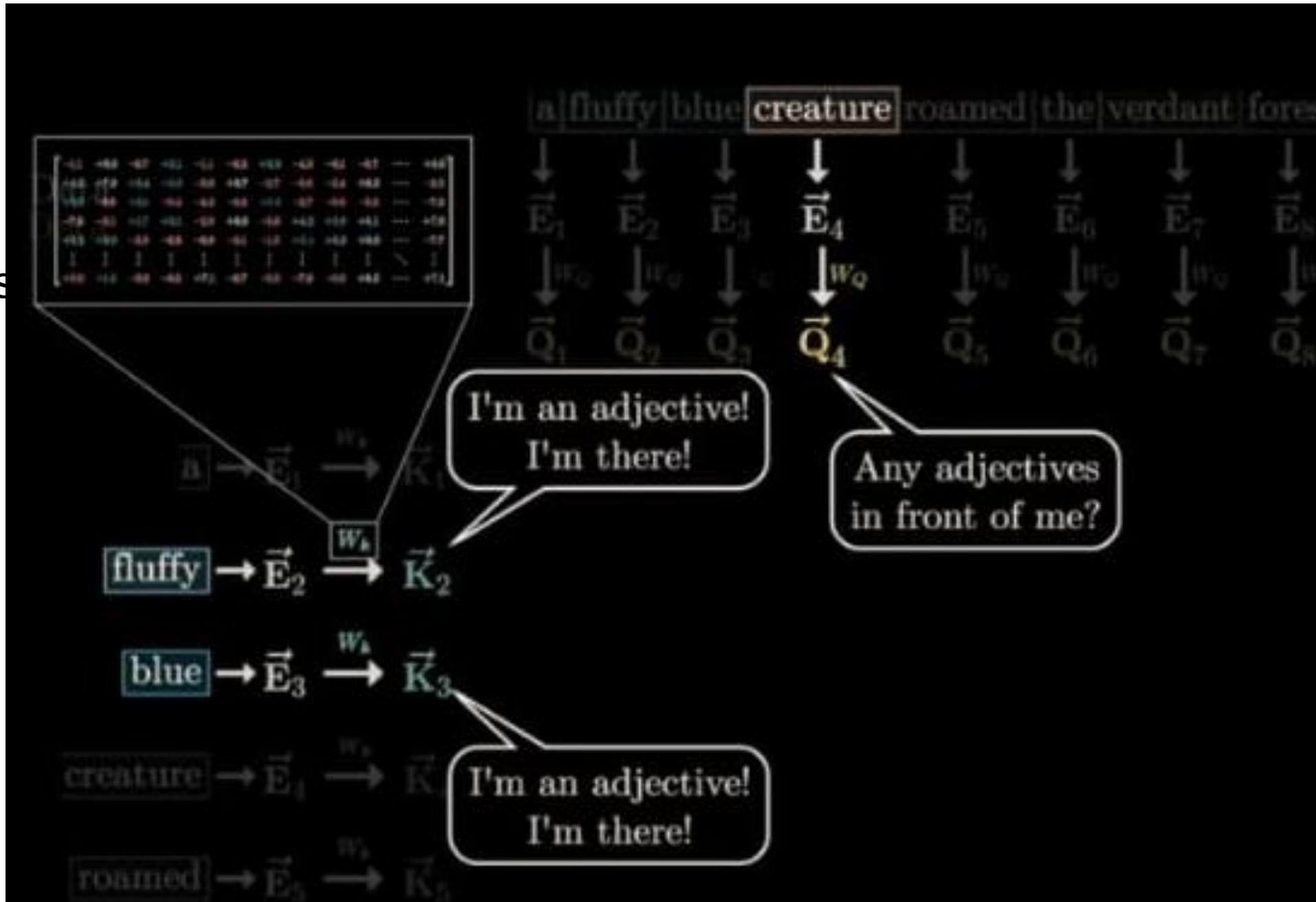
generate them and their roles in the
attention mechanism this entire



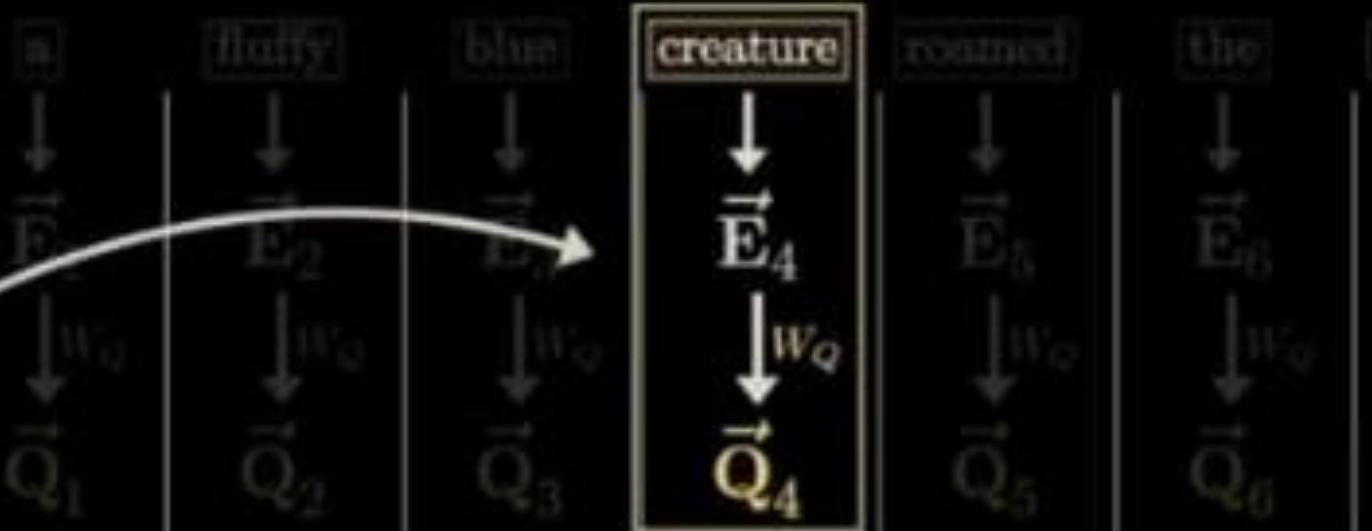
a|fluffy|blue|creature|roamed|the|verdant|forest



-



“Attend to”



$$[a] \rightarrow \vec{E}_1 \xrightarrow{w_k} \vec{K}_1$$

$$\text{fluffy} \rightarrow \vec{E}_2 \xrightarrow{w_k} \vec{K}_2$$

$$\text{blue} \rightarrow \vec{E}_3 \xrightarrow{w_k} \vec{K}_3$$

$$\text{creature} \rightarrow \vec{E}_4 \xrightarrow{w_k} \vec{K}_4$$

+93.0

+93.4

	\vec{E}_1	\vec{E}_2	\vec{E}_3	\vec{E}_4	\vec{E}_5	\vec{E}_6	\vec{E}_7	\vec{E}_8
b	\vec{Q}_1	\vec{Q}_2	\vec{Q}_3	\vec{Q}_4	\vec{Q}_5	\vec{Q}_6	\vec{Q}_7	\vec{Q}_8
b $\rightarrow \vec{E}_1 \xrightarrow{w_1} \vec{K}_1$	$\vec{R}_1 \cdot \vec{Q}_1$	$\vec{R}_1 \cdot \vec{Q}_2$	$\vec{R}_1 \cdot \vec{Q}_3$	$\vec{R}_1 \cdot \vec{Q}_4$	$\vec{R}_1 \cdot \vec{Q}_5$	$\vec{R}_1 \cdot \vec{Q}_6$	$\vec{R}_1 \cdot \vec{Q}_7$	$\vec{R}_1 \cdot \vec{Q}_8$
fluffy $\rightarrow \vec{E}_2 \xrightarrow{w_2} \vec{K}_2$	$\vec{R}_2 \cdot \vec{Q}_1$	$\vec{R}_2 \cdot \vec{Q}_2$	$\vec{R}_2 \cdot \vec{Q}_3$	$\vec{R}_2 \cdot \vec{Q}_4$	$\vec{R}_2 \cdot \vec{Q}_5$	$\vec{R}_2 \cdot \vec{Q}_6$	$\vec{R}_2 \cdot \vec{Q}_7$	$\vec{R}_2 \cdot \vec{Q}_8$
blue $\rightarrow \vec{E}_3 \xrightarrow{w_3} \vec{K}_3$	$\vec{R}_3 \cdot \vec{Q}_1$	$\vec{R}_3 \cdot \vec{Q}_2$	$\vec{R}_3 \cdot \vec{Q}_3$	$\vec{R}_3 \cdot \vec{Q}_4$	$\vec{R}_3 \cdot \vec{Q}_5$	$\vec{R}_3 \cdot \vec{Q}_6$	$\vec{R}_3 \cdot \vec{Q}_7$	$\vec{R}_3 \cdot \vec{Q}_8$
creature $\rightarrow \vec{E}_4 \xrightarrow{w_4} \vec{K}_4$	$\vec{R}_4 \cdot \vec{Q}_1$	$\vec{R}_4 \cdot \vec{Q}_2$	$\vec{R}_4 \cdot \vec{Q}_3$	$\vec{R}_4 \cdot \vec{Q}_4$	$\vec{R}_4 \cdot \vec{Q}_5$	$\vec{R}_4 \cdot \vec{Q}_6$	$\vec{R}_4 \cdot \vec{Q}_7$	$\vec{R}_4 \cdot \vec{Q}_8$
roamed $\rightarrow \vec{E}_5 \xrightarrow{w_5} \vec{K}_5$	$\vec{R}_5 \cdot \vec{Q}_1$	$\vec{R}_5 \cdot \vec{Q}_2$	$\vec{R}_5 \cdot \vec{Q}_3$	$\vec{R}_5 \cdot \vec{Q}_4$	$\vec{R}_5 \cdot \vec{Q}_5$	$\vec{R}_5 \cdot \vec{Q}_6$	$\vec{R}_5 \cdot \vec{Q}_7$	$\vec{R}_5 \cdot \vec{Q}_8$
the $\rightarrow \vec{E}_6 \xrightarrow{w_6} \vec{K}_6$	$\vec{R}_6 \cdot \vec{Q}_1$	$\vec{R}_6 \cdot \vec{Q}_2$	$\vec{R}_6 \cdot \vec{Q}_3$	$\vec{R}_6 \cdot \vec{Q}_4$	$\vec{R}_6 \cdot \vec{Q}_5$	$\vec{R}_6 \cdot \vec{Q}_6$	$\vec{R}_6 \cdot \vec{Q}_7$	$\vec{R}_6 \cdot \vec{Q}_8$
verdant $\rightarrow \vec{E}_7 \xrightarrow{w_7} \vec{K}_7$	$\vec{R}_7 \cdot \vec{Q}_1$	$\vec{R}_7 \cdot \vec{Q}_2$	$\vec{R}_7 \cdot \vec{Q}_3$	$\vec{R}_7 \cdot \vec{Q}_4$	$\vec{R}_7 \cdot \vec{Q}_5$	$\vec{R}_7 \cdot \vec{Q}_6$	$\vec{R}_7 \cdot \vec{Q}_7$	$\vec{R}_7 \cdot \vec{Q}_8$
forest $\rightarrow \vec{E}_8 \xrightarrow{w_8} \vec{K}_8$	$\vec{R}_8 \cdot \vec{Q}_1$	$\vec{R}_8 \cdot \vec{Q}_2$	$\vec{R}_8 \cdot \vec{Q}_3$	$\vec{R}_8 \cdot \vec{Q}_4$	$\vec{R}_8 \cdot \vec{Q}_5$	$\vec{R}_8 \cdot \vec{Q}_6$	$\vec{R}_8 \cdot \vec{Q}_7$	$\vec{R}_8 \cdot \vec{Q}_8$

	\vec{E}_1	\vec{E}_2	\vec{E}_3	\vec{E}_4	\vec{E}_5	\vec{E}_6	\vec{E}_7	\vec{E}_8	
\vec{e}	\vec{Q}_1	\vec{Q}_2	\vec{Q}_3	\vec{Q}_4	\vec{Q}_5	\vec{Q}_6	\vec{Q}_7	\vec{Q}_8	
w_1	-0.7	0.37	-0.67	-27.5	5.2	-00.3	-0.2	00.1	
furry $\rightarrow \vec{E}_1 \xrightarrow{w_1} \vec{K}_1$	73.4	13.9	5.4	+93.0	We want these to act like weights				
blue $\rightarrow \vec{E}_2 \xrightarrow{w_1} \vec{K}_2$	52.4	3.2	-0.18	+93.4					
creature $\rightarrow \vec{E}_3 \xrightarrow{w_1} \vec{K}_3$	21.5	20.7	50.1	+4.9	32.4	02.3	9.5	-26.1	
roamed $\rightarrow \vec{E}_4 \xrightarrow{w_1} \vec{K}_4$	-39.1	-0.19	-87.8	-55.4	-0.8	-64.7	-96.7	-15.9	
the $\rightarrow \vec{E}_5 \xrightarrow{w_1} \vec{K}_5$	87.9	-0.13	-22.6	-31.4	-0.5	-0.6	-4.6	-96.8	
verdant $\rightarrow \vec{E}_6 \xrightarrow{w_1} \vec{K}_6$	41.2	50.5	-0.3	-59.8	-20.0	07.9	-0.7	+93.8	
forest $\rightarrow \vec{E}_7 \xrightarrow{w_1} \vec{K}_7$	-58.9	-75.5	-0.1	-90.6	-75.6	-0.0	-70.8	-1.7	



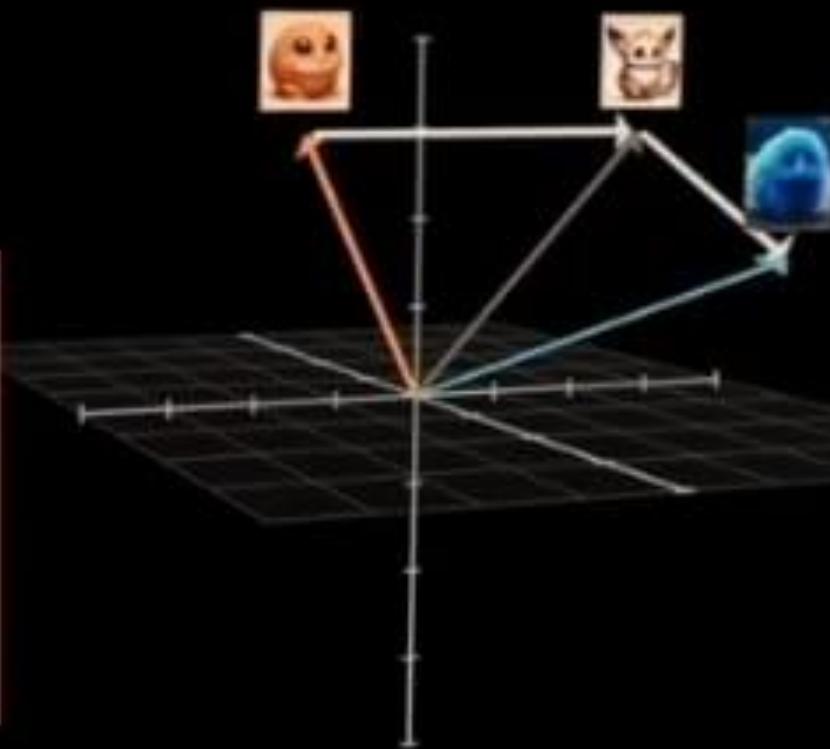
blue fluffy creature

↓ ↓ ↓

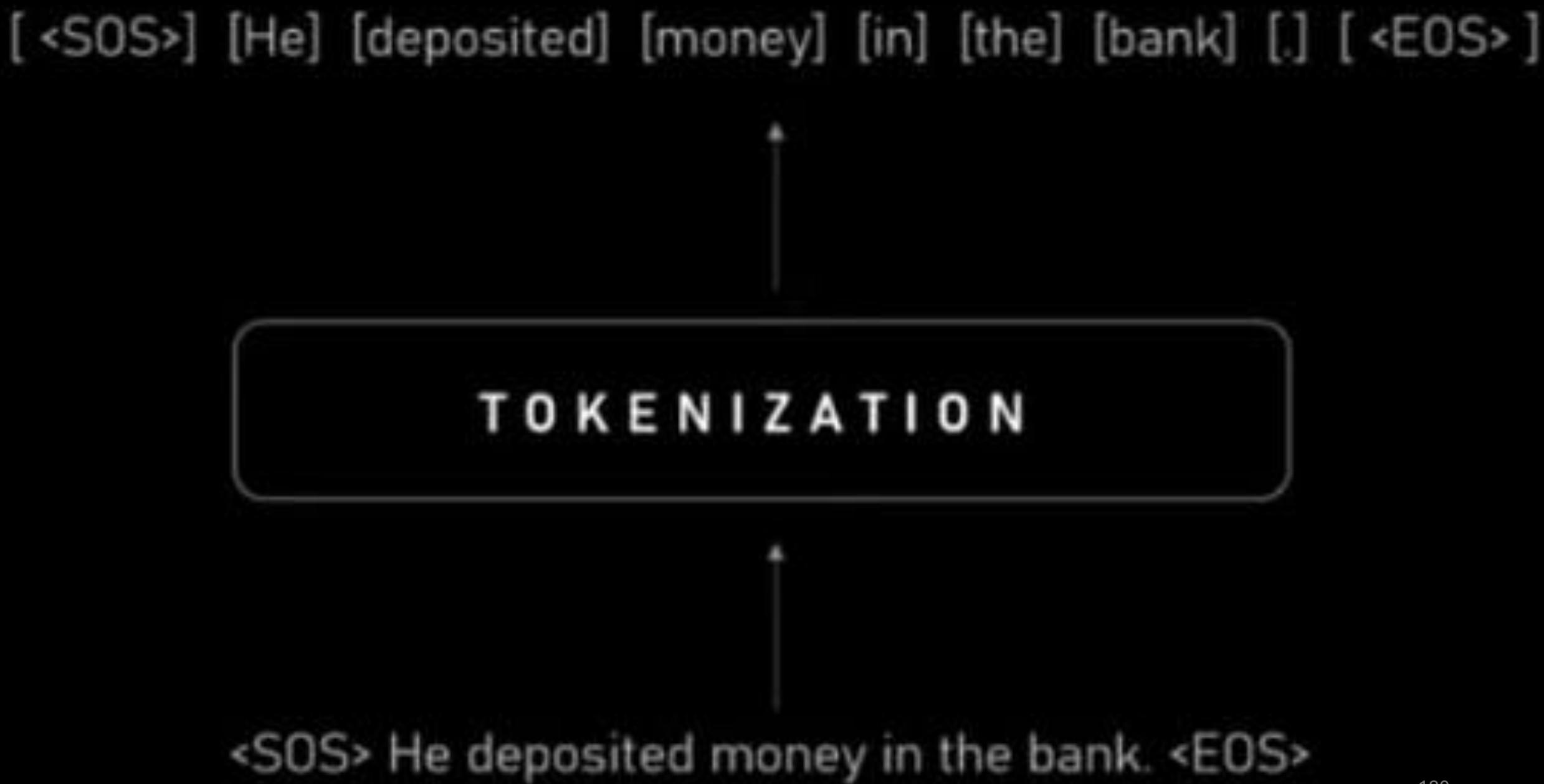
$$\begin{bmatrix} +1.0 \\ +4.3 \\ +12.6 \\ +95.7 \\ +11.2 \\ +14.4 \\ \vdots \\ +7.8 \end{bmatrix} = \begin{bmatrix} +1.0 \\ +4.3 \\ +12.6 \\ +95.7 \\ +11.2 \\ +14.4 \\ \vdots \\ +7.8 \end{bmatrix}$$

W_V

$$\begin{bmatrix} -3.6 & -1.7 & -8.6 & +3.5 & +1.3 & -4.0 & \cdots & -6.0 \\ +1.5 & +8.5 & -3.6 & +3.3 & -7.3 & +4.3 & \cdots & -6.3 \\ +1.7 & -9.5 & +6.5 & -9.8 & +3.5 & -4.0 & \cdots & +9.2 \\ -5.0 & +1.5 & +1.8 & +1.4 & -5.5 & +9.0 & \cdots & +6.9 \\ +3.9 & -4.0 & +6.2 & -2.0 & +7.5 & +1.6 & \cdots & +3.8 \\ +4.5 & +0.0 & +9.0 & +2.9 & -1.5 & +2.1 & \cdots & -3.9 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ +1.5 & +3.0 & +3.0 & -1.4 & +7.9 & -2.6 & \cdots & +7.8 \end{bmatrix}$$

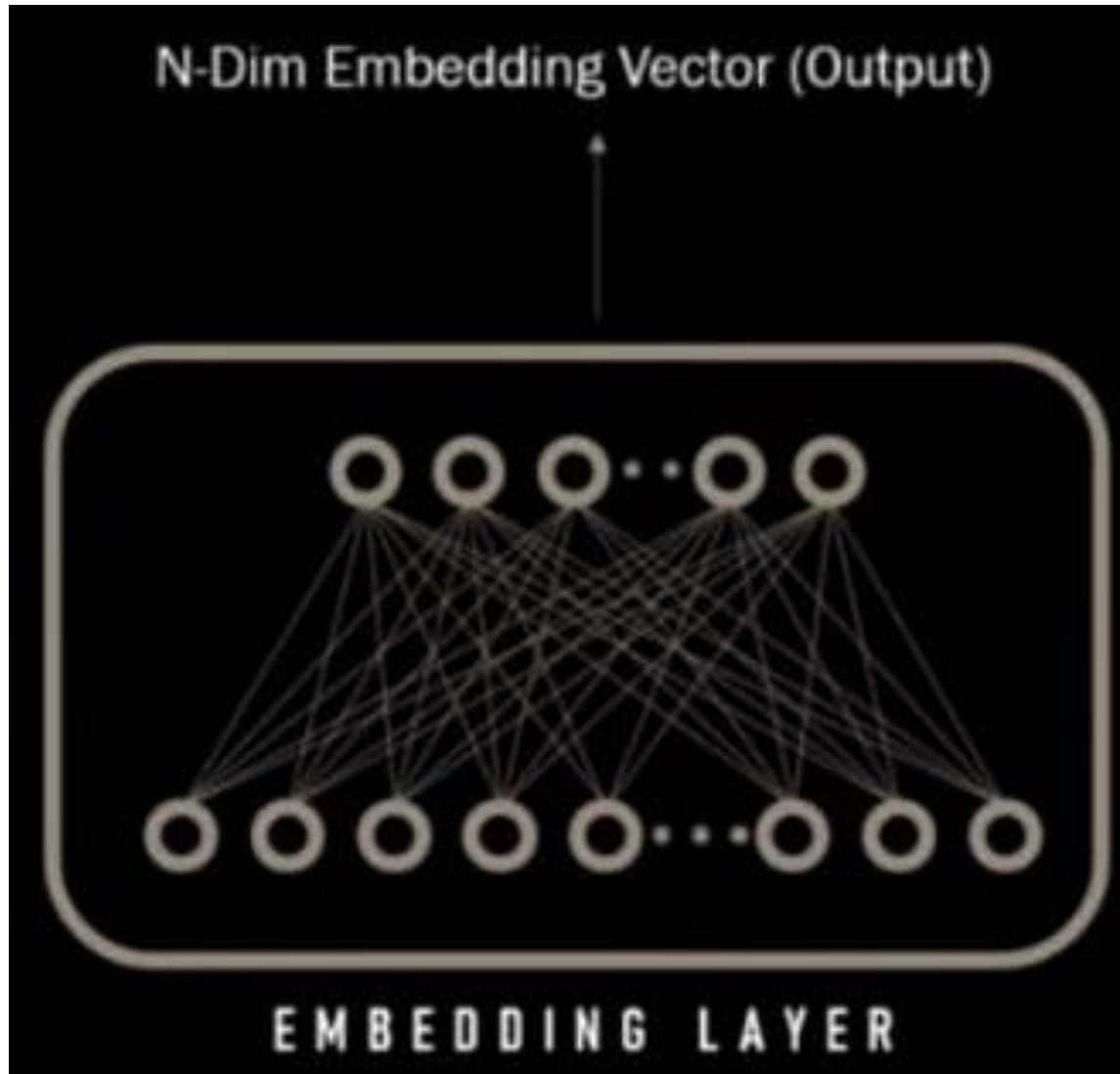


[<SOS>] [He] [deposited] [money] [in] [the] [bank] [.] [<EOS>]



TOKENIZATION

<SOS> He deposited money in the bank. <EOS>

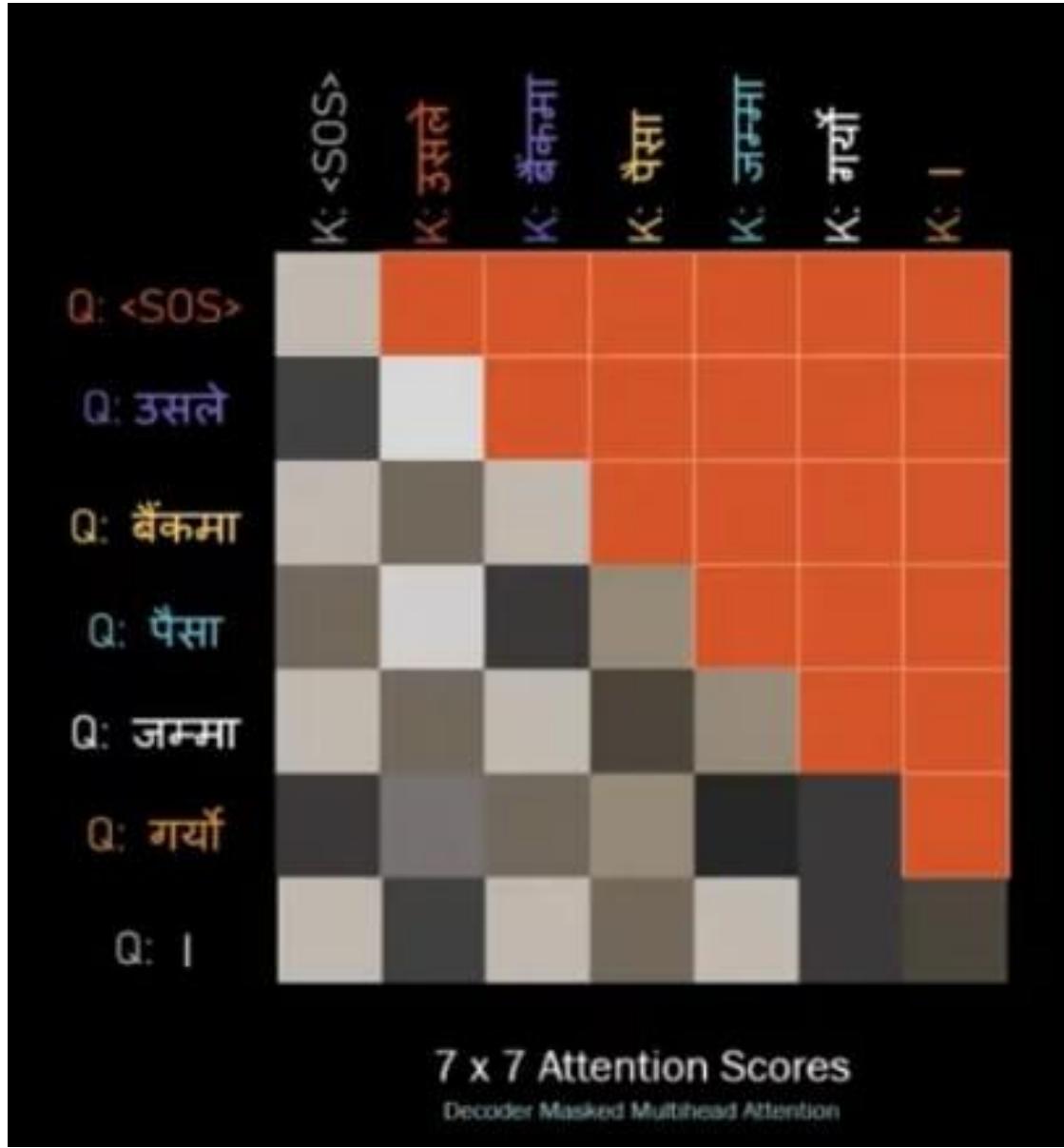


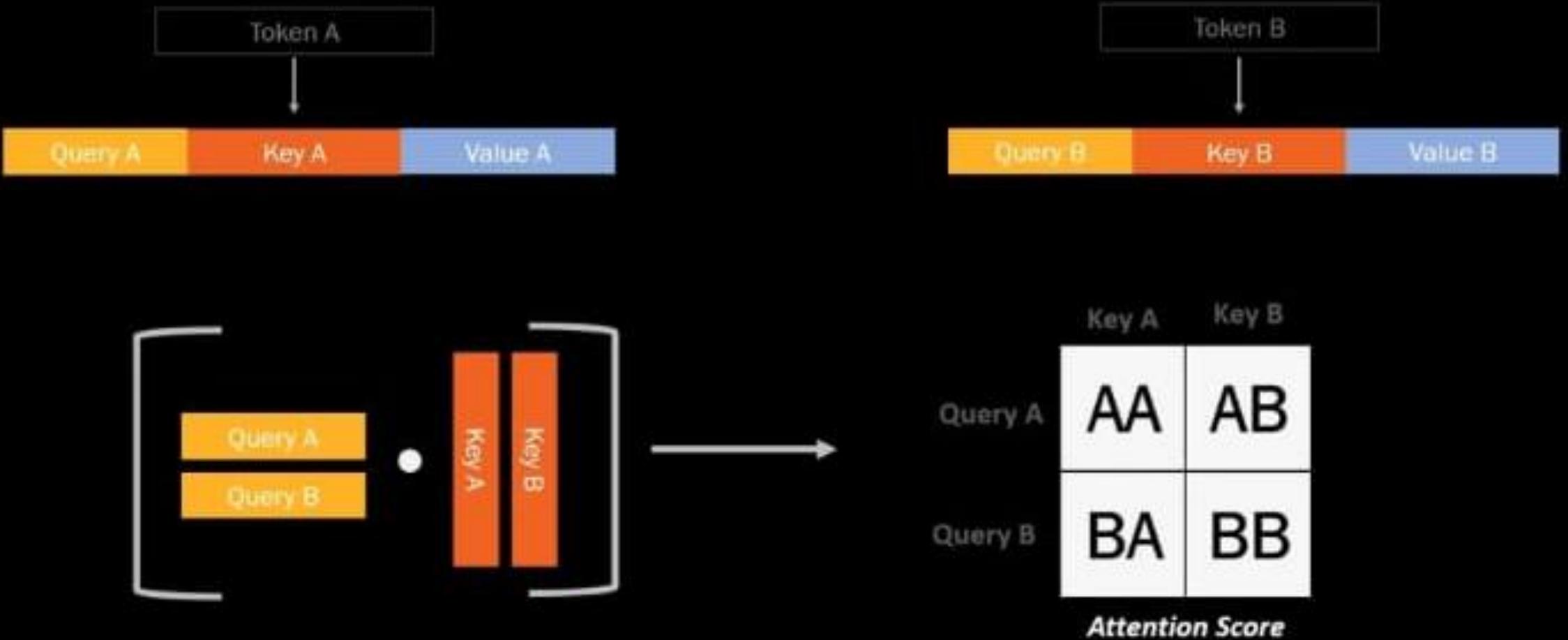
<SOS> He deposited money in the bank. <EOS>
<SOS> money in the bank. <EOS> <PAD> <PAD>
<SOS> money. <EOS> <PAD> <PAD> <PAD> <PAD>

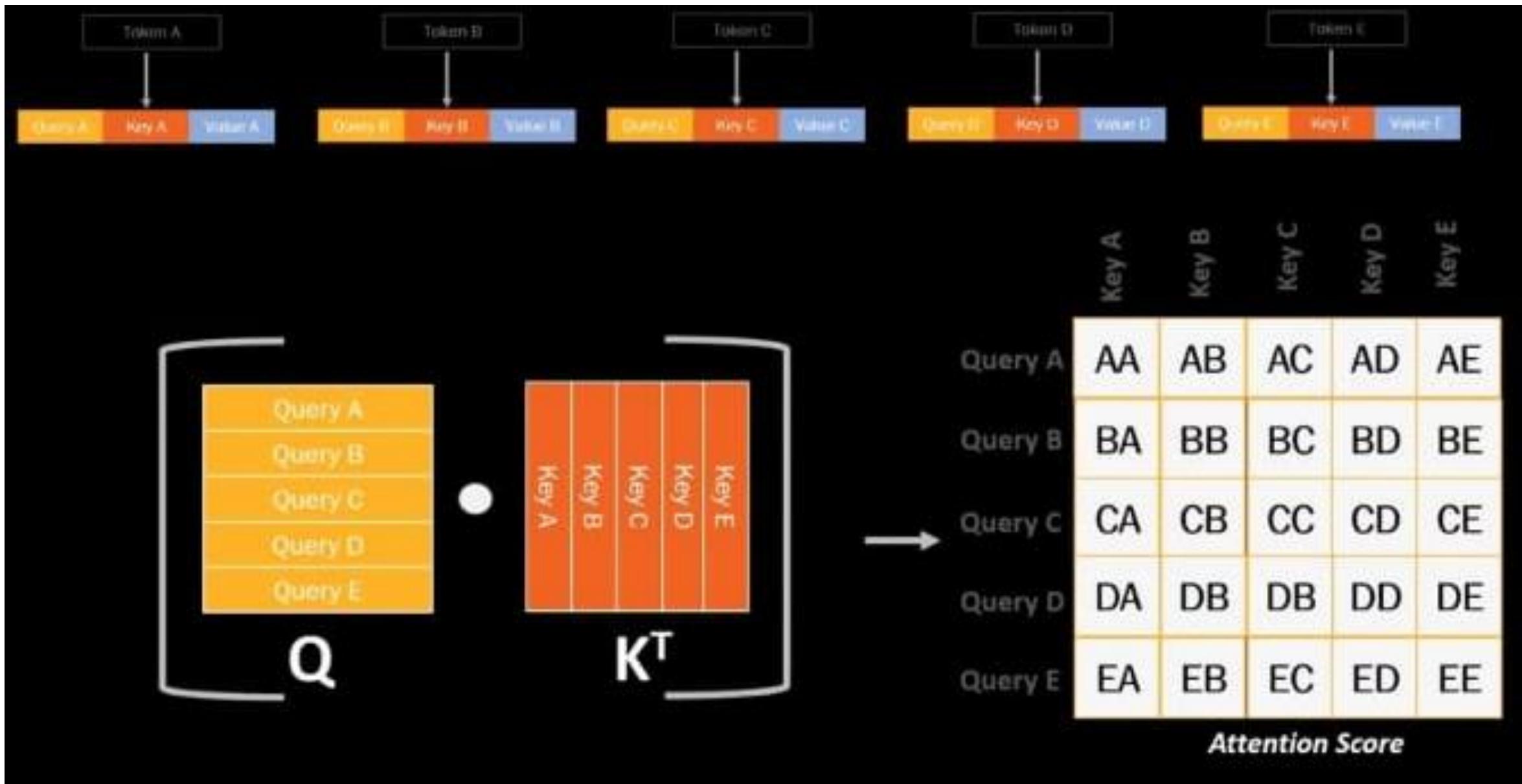
TOKENIZATION

[1, 80, 210, 410, 30, 60, 110, 90, 2]
[1, 410, 30, 60, 110, 90, 2, 0, 0]
[1, 410, 90, 2, 0, 0, 0, 0, 0]

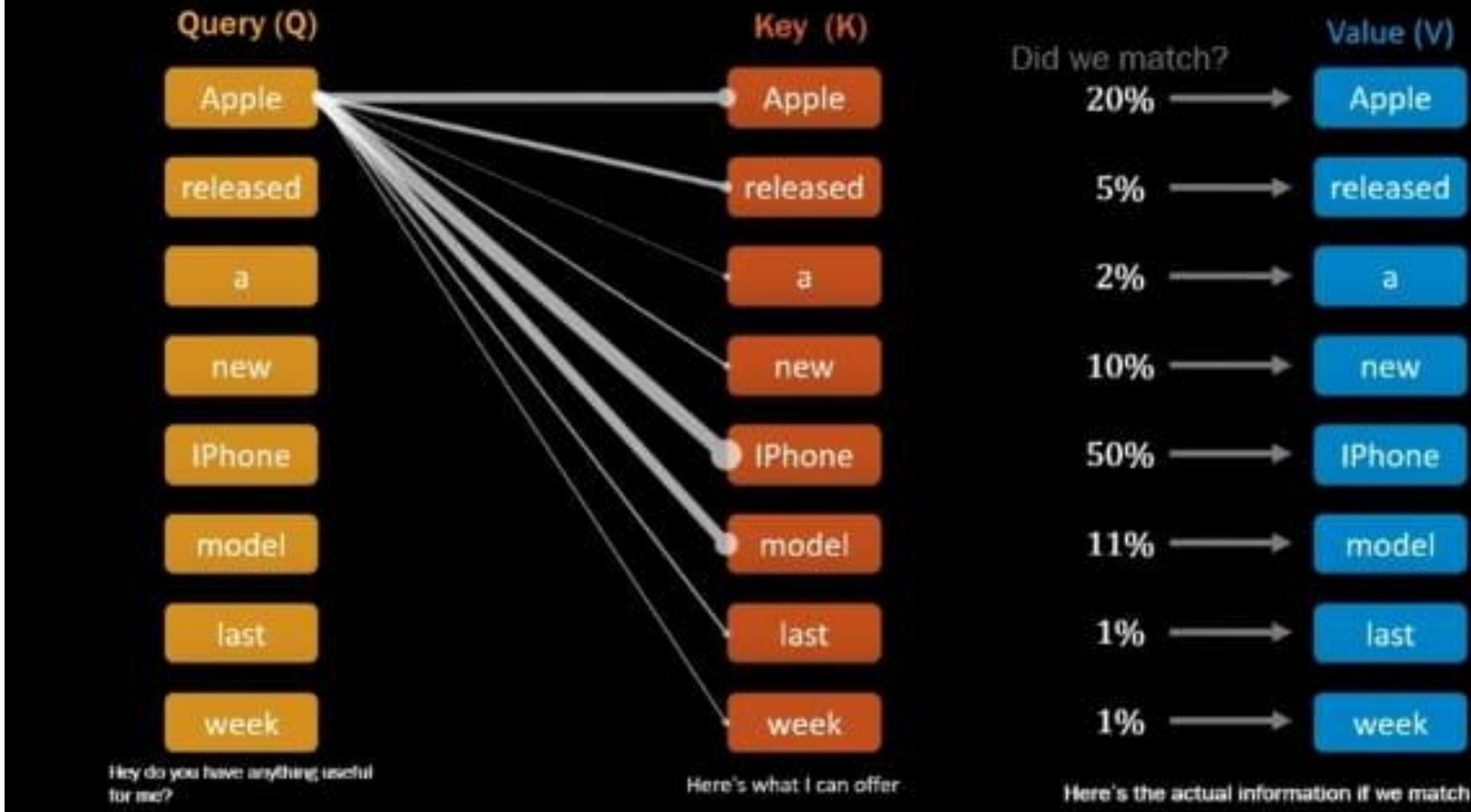
(input tokens)

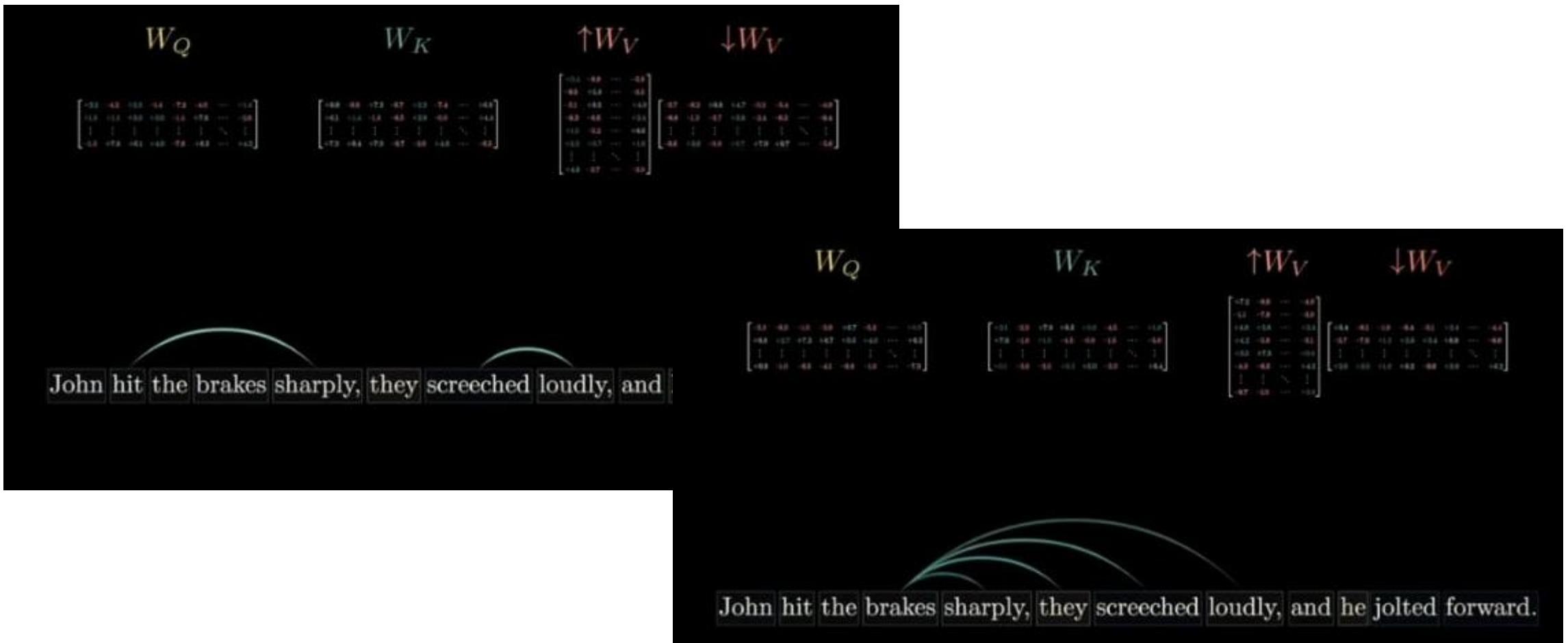


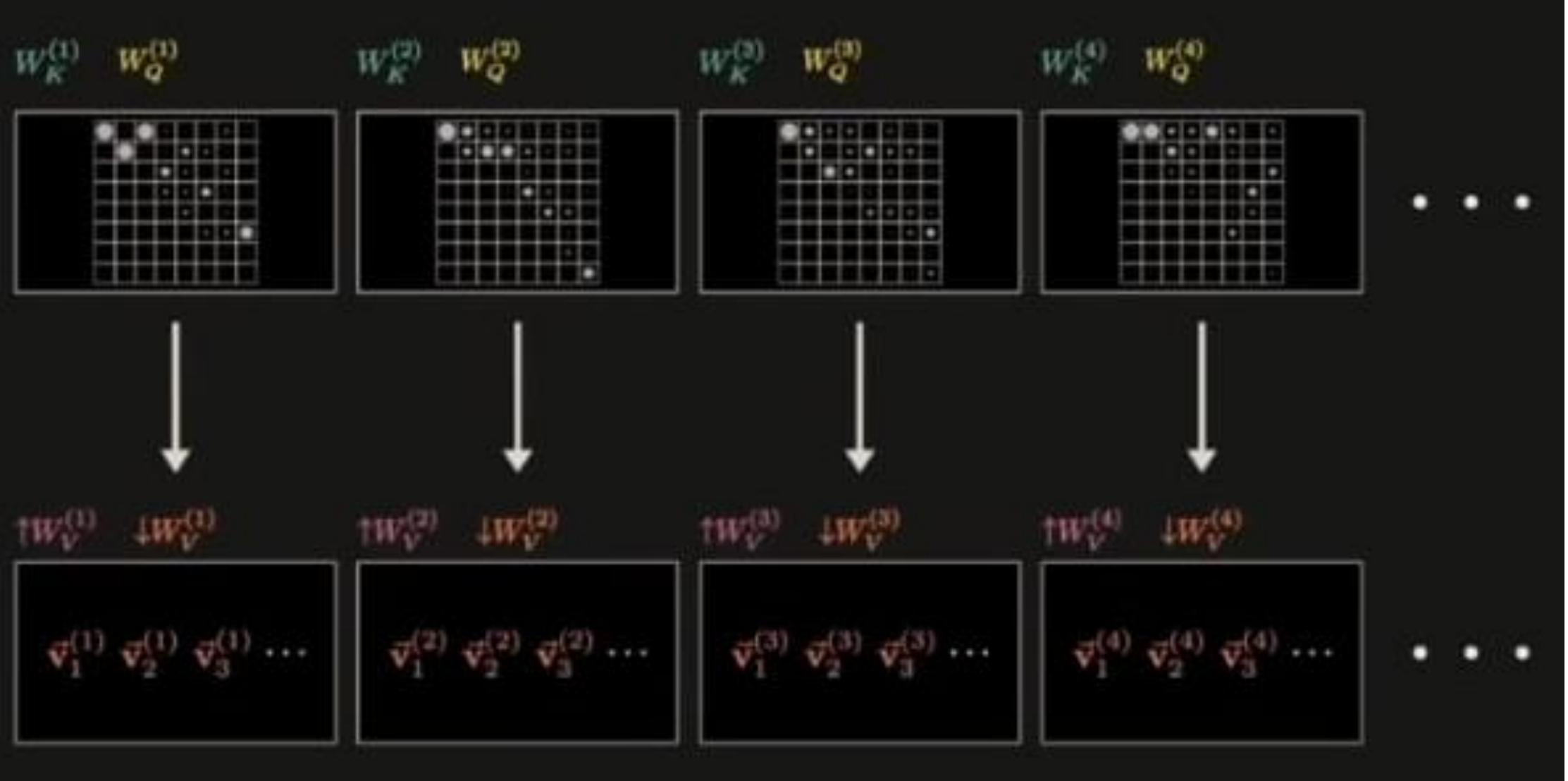


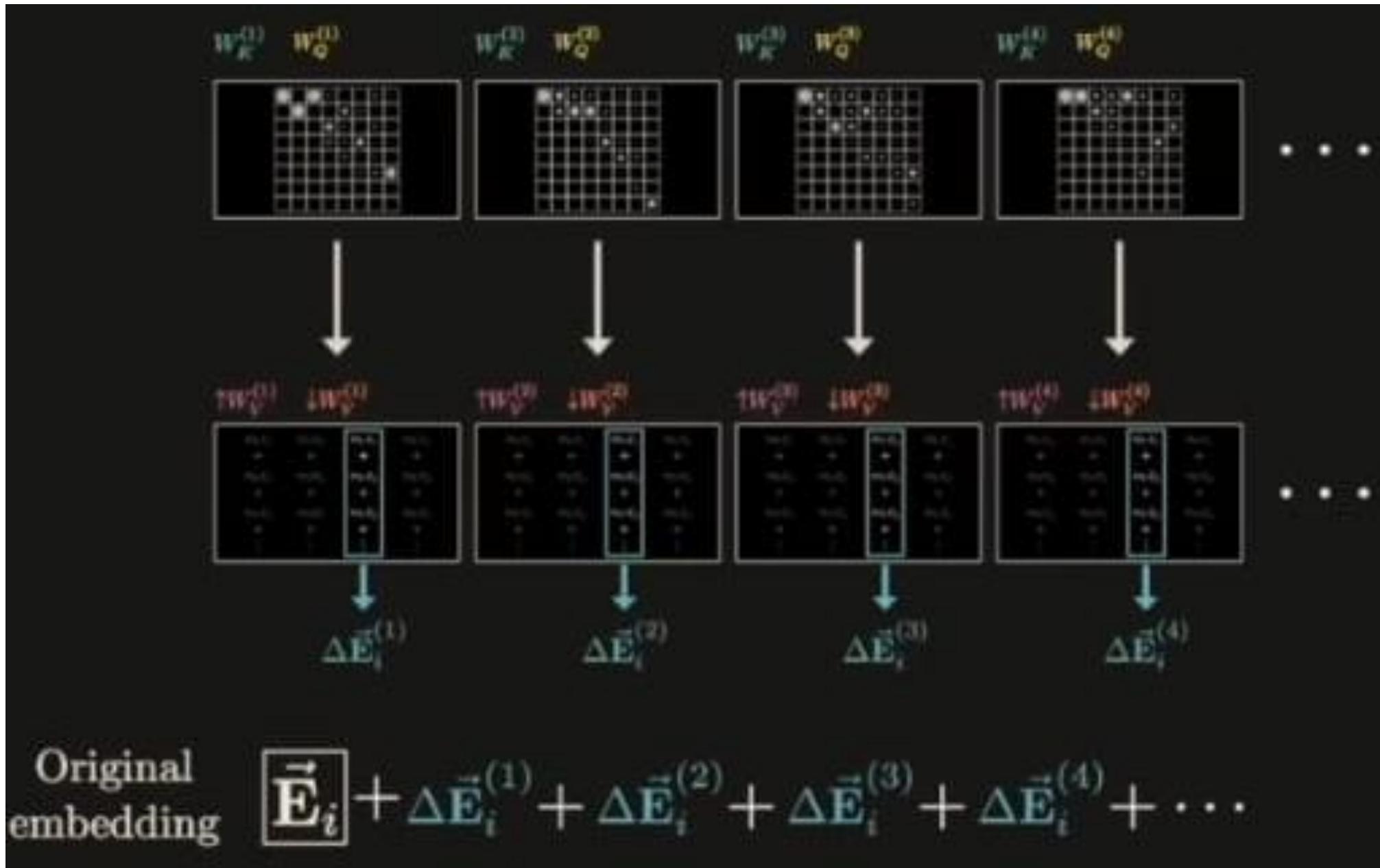


“Apple released a new iPhone model last week.”



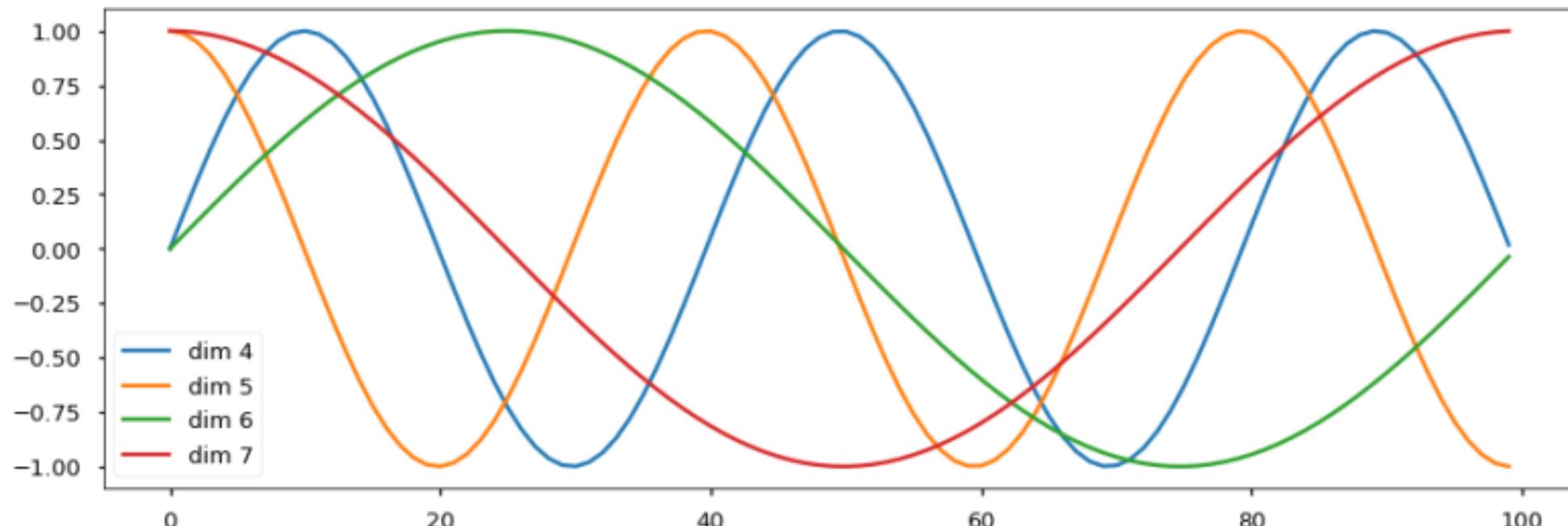




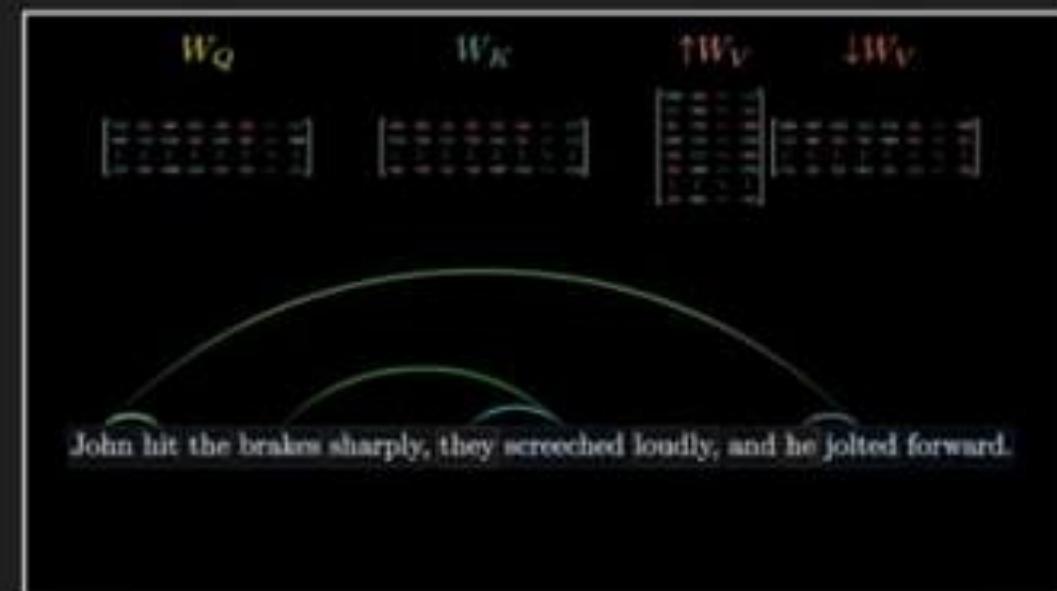


Transformer – Positional Encoding

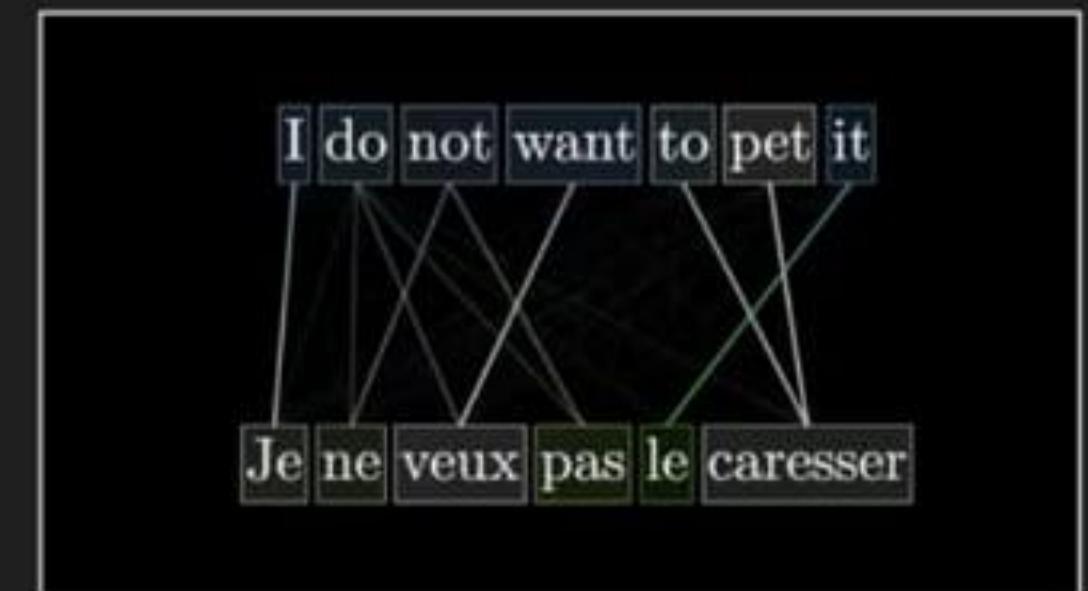
- Transformers add sinusoid curves to the input, before the attention
 - Informs about relative position inside the sequence
 - Removes need for explicit recurrence patterns



Self-attention

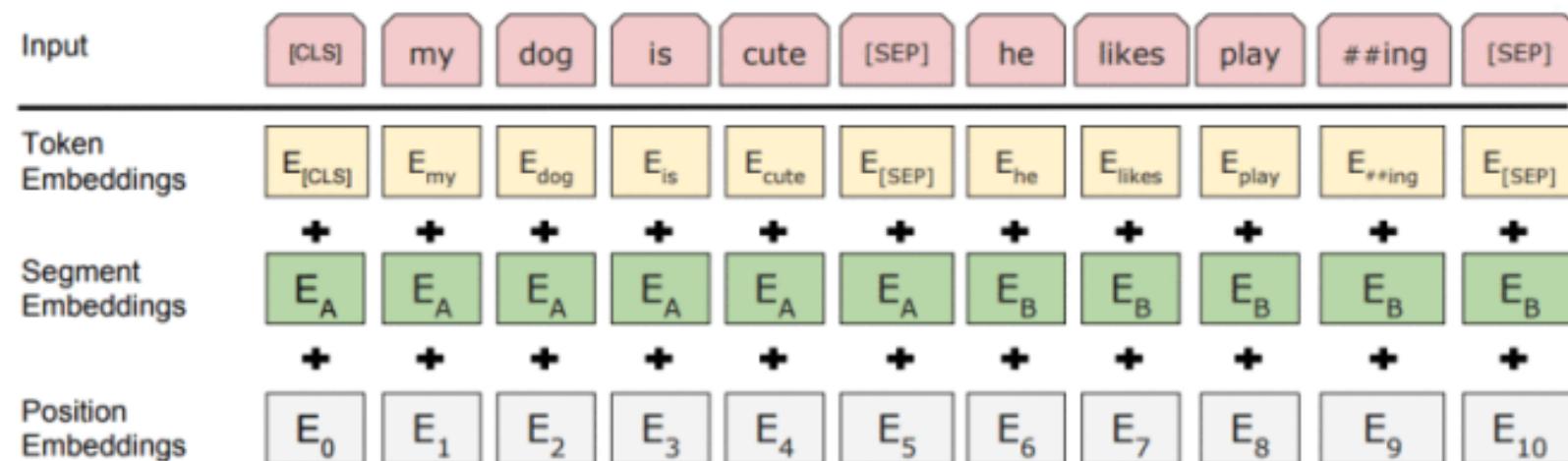


Cross-attention



BERT - Input

- Either one or two sentences, always prepended with [CLS]
 - BERT adds trained position embeddings & sequence embeddings



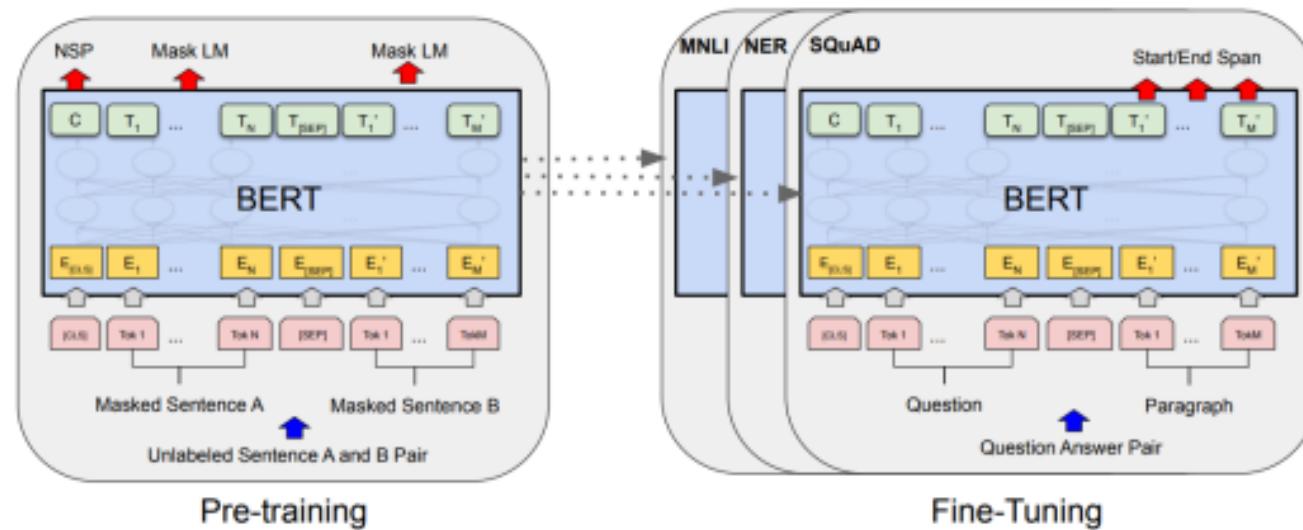
More info: <https://towardsml.com/2019/09/17/bert-explained-a-complete-guide-with-theory-and-tutorial/>

BERT - Model

- Model itself is quite simple: n Layers of stacked Transformers
 - Using LayerNorm, GeLU activations (like ReLU, but with a grace swing under 0)
 - Task specific heads on top to pool [CLS] or individual token representations
 - Every Transformer layer receives as input the output of the previous one
- The [CLS] token itself is only special because we train it to be
 - No mechanism inside the model that differentiates it from other tokens
- Novel contributions center around pre-training & workflow

BERT - Workflow

- Someone with lots of compute or time pre-trains a large model
 - BERT uses Masked Language Modelling [MASK] and Next Sentence Prediction [CLS]
- We download it and fine-tune on our task



BERT++

- Same as with Transformer variations, there are now many BERT variants
 - For many languages
 - Domains like biomedical publications
 - Different architectures, but similar workflow:
Roberta, Transformer-XL, XLNet, Longformer ...
- Main themes for adapted architectures:
 - Bigger
 - More efficient
 - Allowing for longer sequences (BERT is capped at 512 tokens in total)

Pre-Training Ecosystem

- With simple 1-word-1-vector embeddings (word2vec) sharing was as simple as a single text file containing both vocab + weights
 - We could simply load the weight matrix into bigger models
 - Mostly whitespace tokenization meant very little complexity
- BERT et al. re-use requires:
 - Exact model architecture (specific code and config) for hundreds of details
 - Weights for 100+ modules
 - Specific tokenizer rules for sub-word tokenization and special token handling
 - A single text file doesn't work here anymore ...

Soooo many tasks are solvable with BERT

- Original BERT paper evaluates on:
 - GLUE, SQuAD, SWAG, CoNLL
- Now at ~18 thousand citations, we can assume some more are evaluated
 - As long as your text input is <512 tokens & you can pool the CLS token or learn per-term predictions you can use BERT
- HuggingFace model Hub alone provides out-of-the-box support for dozens of tasks & lists 200+ datasets used by the trained models.

Extractive Question Answering

- Given a query and a passage/document:
Select the words in the passage that answer the query
 - We want to select at least 1 span with a start and end position, that we then can extract
 - Use extracted text in highlighted UI (with surrounding text), chatbot, or audio-based assistant
 - Not perfect: query type must be specific to be answerable with fixed text
- Differs from *Generative* Question Answering
 - Models are tasked to create new text (with new words, more natural conversational style)
 - More complex, more potential for error and biases

Extractive QA: Datasets

- Popular datasets include: SQuAD & NaturalQuestions
 - Both are based on Wikipedia text
 - SQuAD contains artificially created queries, NQ google search queries
 - Both come with fairly large training and evaluation sets
- Many pre-trained models are available for both

Predictions by nlnet (single model) (Microsoft Research Asia)
Article EM: 87.0 F1: 88.2

Geology

The Stanford Question Answering Dataset

There are three major types of **rock**: **igneous**, sedimentary, and metamorphic. The **Rock cycle** is an important concept in geology which illustrates the relationships between these three types of **rock** and magma. When a **rock** **crystallizes** from melt (magma and/or lava), it is an **igneous rock**. This rock can be weathered and eroded, and then redeposited and lithified into a sedimentary **rock**, or be turned into a metamorphic **rock** due to heat and pressure that change the mineral content of the **rock** which gives it a characteristic fabric. The sedimentary **rock** can then be subsequently turned into a metamorphic **rock** due to heat and pressure and is then weathered, eroded, deposited, and lithified, ultimately becoming a sedimentary **rock**. Sedimentary **rock** may also be re-eroded and redeposited, and metamorphic **rock** may also undergo additional metamorphism. All three types of **rock** may be re-melted; when this happens, a new magma is formed, from which an **igneous rock** may once again crystallize.

An igneous rock is a rock that crystallizes from what?
Ground Truth Answers: melt, magma and/or lava
Prediction: melt

Sedimentary rock can be turned into which of the three types of rock?
Ground Truth Answers: metamorphic rock, metamorphic, metamorphic rock, metamorphic
Prediction: metamorphic rock

When the three types of rock are re-melted what is formed?
Ground Truth Answers: new magma, igneous, new magma, magma
Prediction: a new magma

What are the three major types of rock?
Ground Truth Answers: igneous, sedimentary, and metamorphic, igneous, sedimentary, and metamorphic, igneous, sedimentary, and metamorphic
Prediction: igneous, sedimentary, and metamorphic

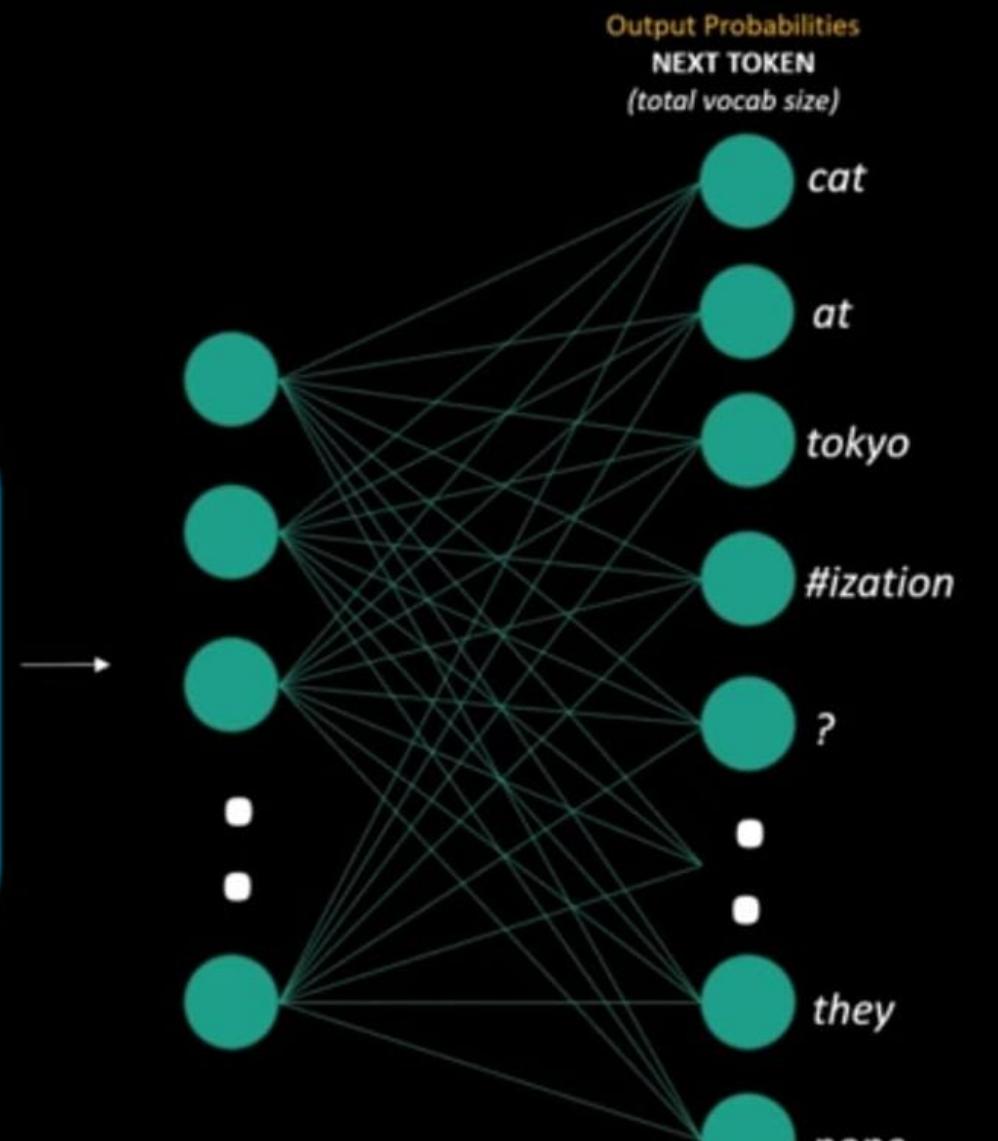
Example: <https://rajpurkar.github.io/SQuAD-explorer/explore/v2.0/dev/>

Extractive QA: Training

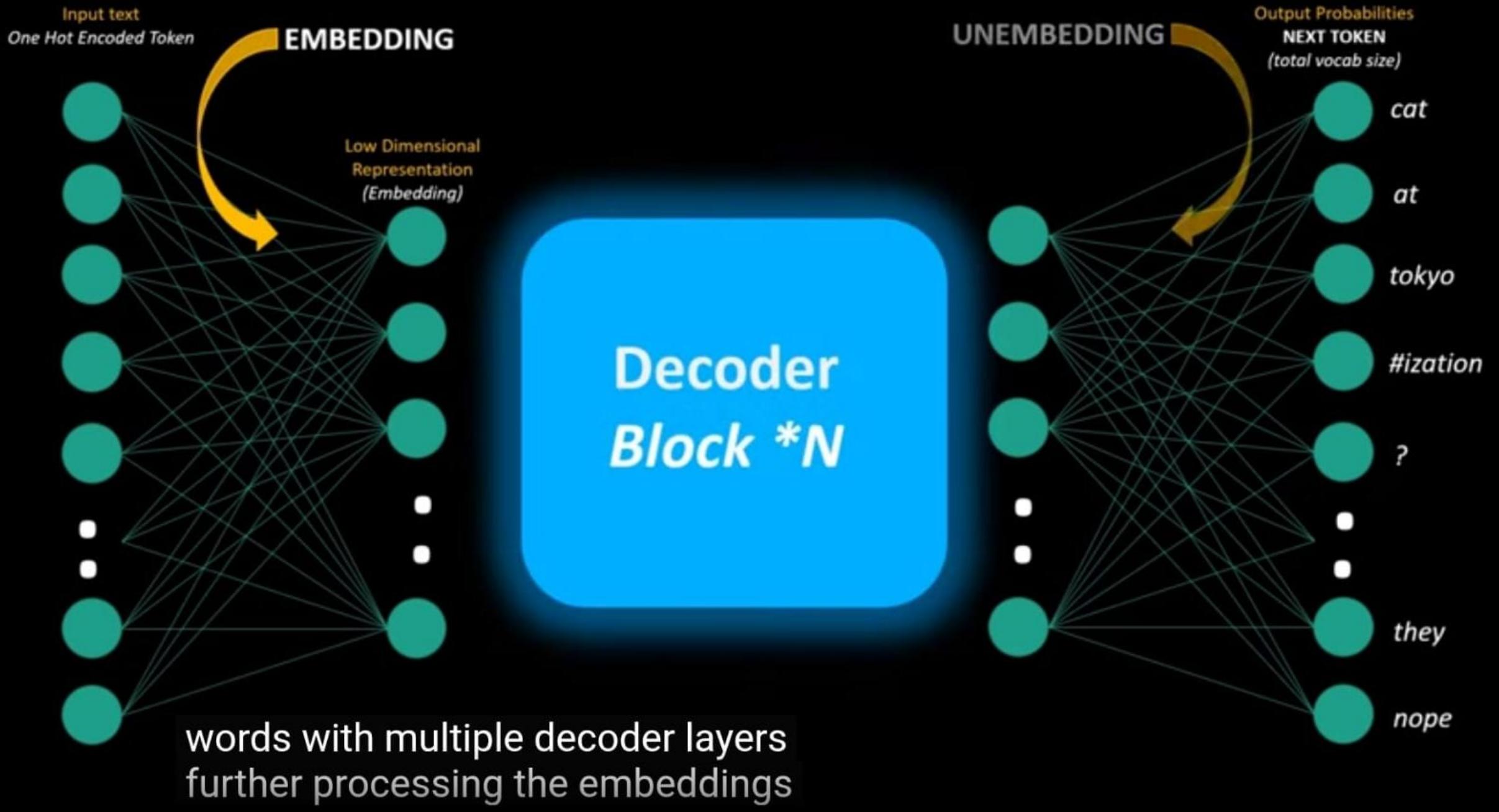
- For BERT we concatenate query and passage
 - With the pre-trained special tokens
- Per term output (of BERT) of the passage predicts if this token is a start or end token of the answer
 - End tokens are trained with gold-label start positions
 - Beam search can be used to find the best combination
- Loss is based on CrossEntropy of prediction vs ground-truth label
 - Potentially also includes a non-answerable prediction for the passage as a whole (SQuAD 2.0)

Embedding (n-dim)					
0.012	0.819	1.190	1.129	0.412	...
0.492	0.829	0.123	0.886	0.22 1	...
0.220	0.912	1.869	1.986	0.419	...
0.812	0.198	0.310	0.889	0.398	...
0.2	0.8	0.129	0.689	0.198	...
0.2	0.8	0.896	0.861	0.411	...

(no_of_tokens x Embedding_dim)



output layer is another linear layer
similar to the embedding



Encoder Input

<SOS> Hello, how are you? <EOS>

<SOS> What is your name? <EOS>

<SOS> Where are you going? <EOS>

<SOS> She is going to the market to buy fruits. <EOS>

<SOS> The children are playing football in the playground. <EOS>

<SOS> The teacher is explaining the lesson very clearly <EOS>

Decoder Input

<SOS> नमस्ते, तपाईंलाई कस्तो छ?

<SOS> तपाईंको नाम के हो?

<SOS> तपाईं कहाँ जाँदै हुनुहुन्छ?

<SOS> उनी बजारमा फलफूल किन्न जाँदैछिन्।

<SOS> बच्चाहरू खेलमैदानमा फुटबल खेलिरहेका छन्।

<SOS> शिक्षकले पाठलाई धेरै स्पष्ट रूपमा समझाइरहेका छन्।

d_{model} = Model Embedding = 512

h = number of Attention Heads = 4

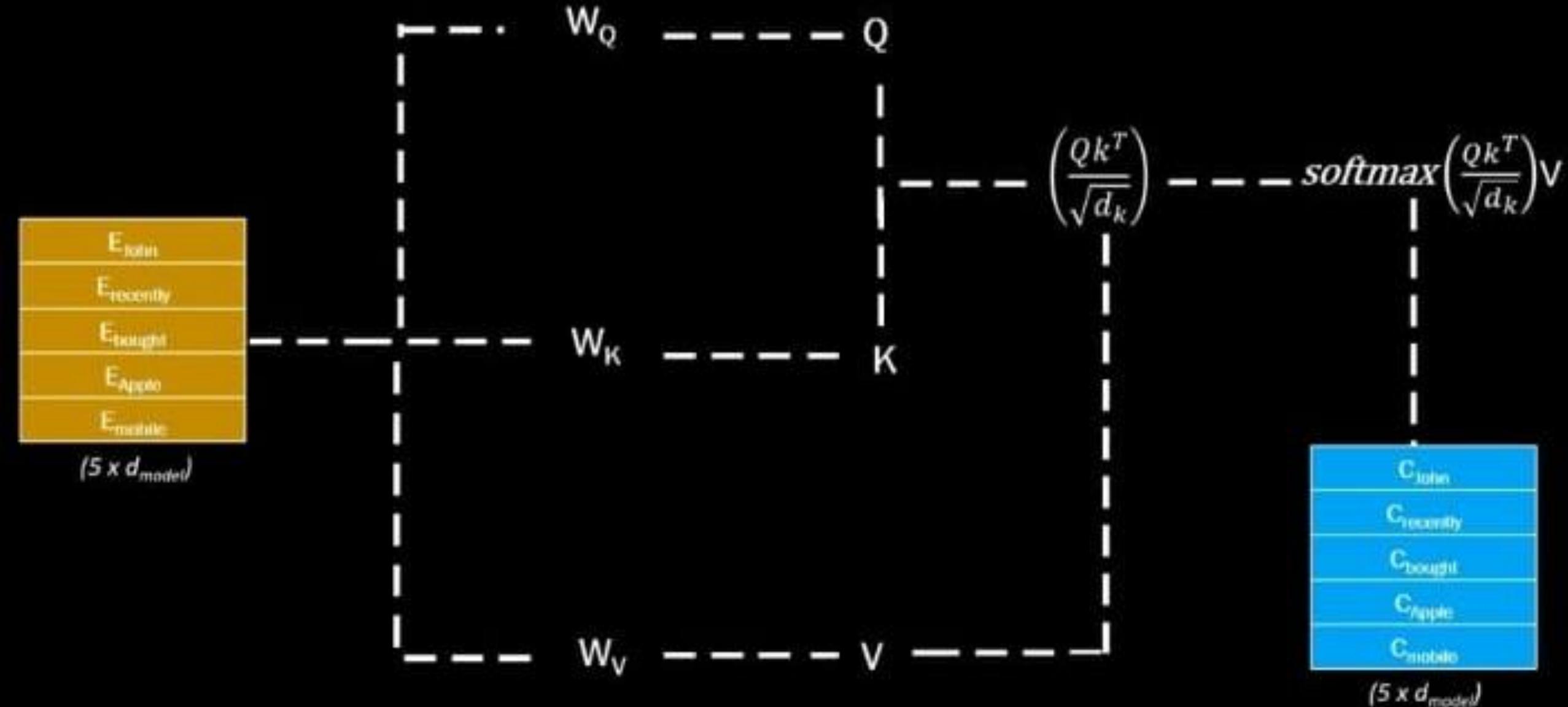
$W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$ $d_k = dv = (d_{model} / h) = (512 / 4) = 128$

$W_i^K \in \mathbb{R}^{d_{model} \times d_k}$ $d_k = \text{Projection Dimension of Key}$
and Query = 128

$W_i^V \in \mathbb{R}^{d_{model} \times d_v}$ $d_v = \text{Projection Dimension of Value} = 128$

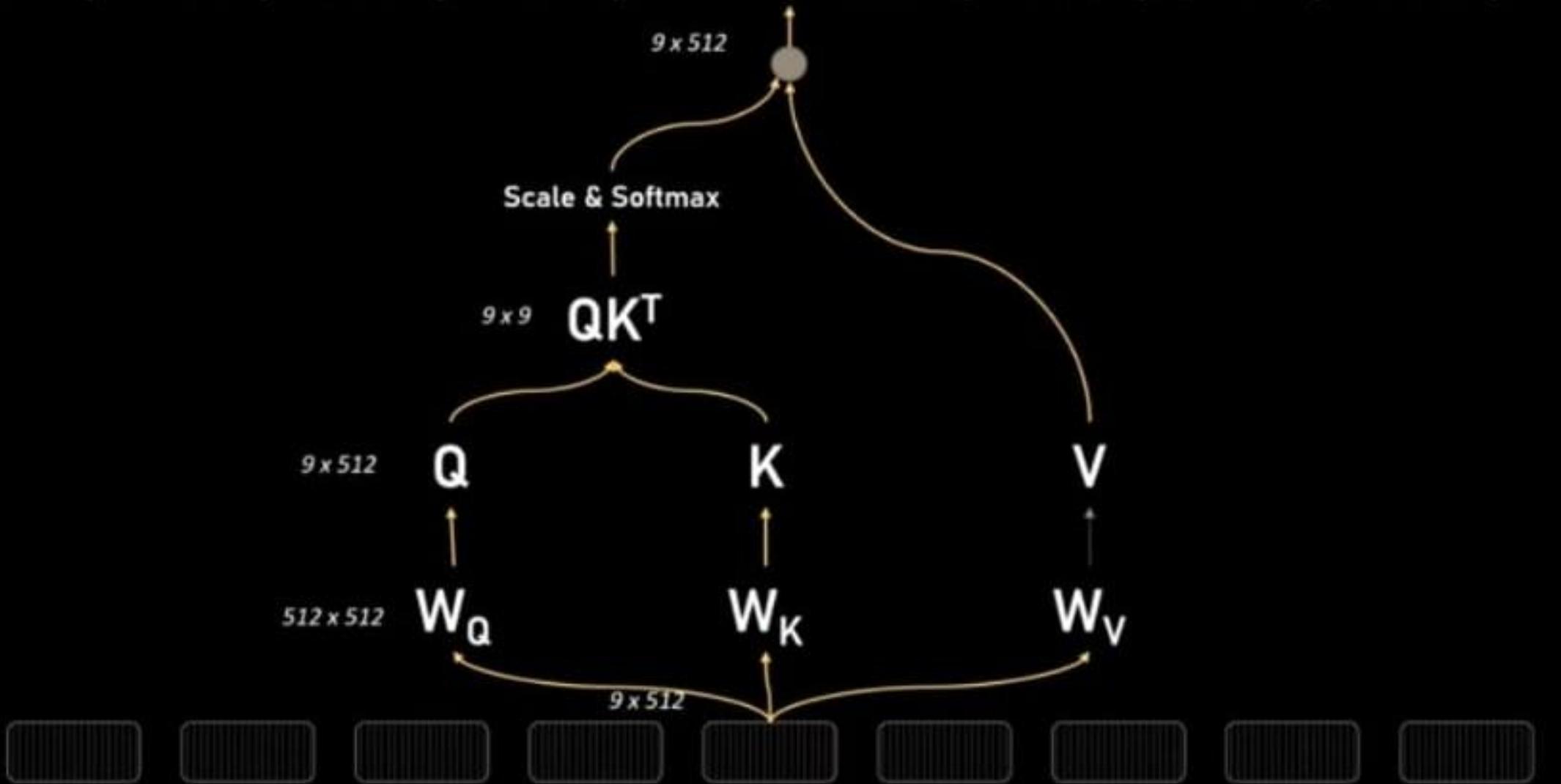
$W^O \in \mathbb{R}^{hdv \times d_{model}}$ $W^O = \text{Output Projection } (4 \times 128 * 512)$

512 additionally we have another new
weight Matrix W





Embedding Token 1 Embedding Token 2 Embedding Token 3 Embedding Token 4 Embedding Token 5 Embedding Token 6 Embedding Token 7 Embedding Token 8 Embedding Token 9



Embedding Token 1 Embedding Token 2 Embedding Token 3 Embedding Token 4 Embedding Token 5 Embedding Token 6 Embedding Token 7 Embedding Token 8 Embedding Token 9

Until the next time 😊

