# Web search engine / ...

I saw a cat|

I saw a cat **on the chair**

I saw a cat **running after a dog**

I saw a cat **in my dream**

I saw a cat **book**

# Translation service / mail agent / ...

I saw a ca|
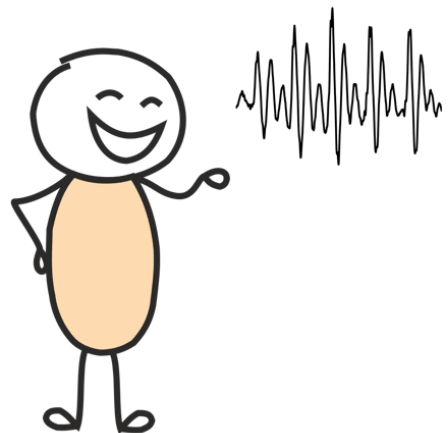
car ↵

Translation service / mail agent / ...

I saw a catt

Probably you meant I saw a cat
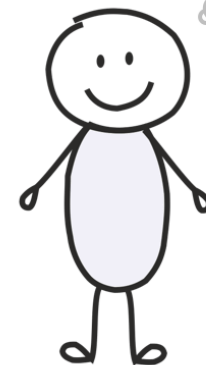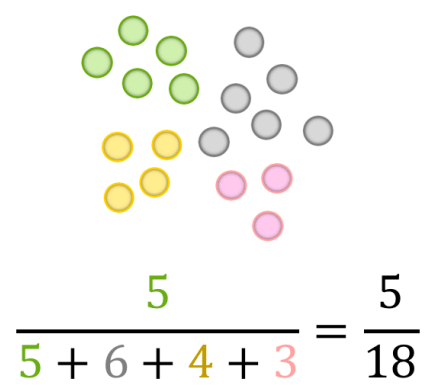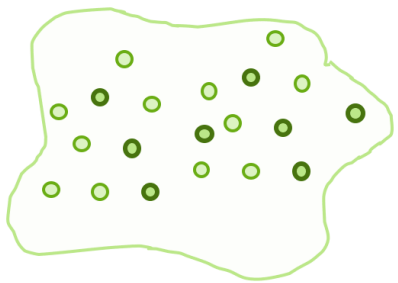
# Keyboard / mail agent / ...

I saw a catt

cat
car

The morphosyntax example is from the slides by Alex Lascarides and Sharon Goldwater, Foundations of Natural Language Processing course at the University of Edinburgh.

What is the probability to pick a green ball?



$$\frac{5}{5 + 6 + 4 + 3} = \frac{5}{18}$$

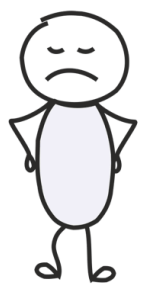Can we do the same for sentences?



Text corpus

$$P(\text{the mut is tinming the tebn}) = \frac{0}{|\text{corpus}|} = 0$$

$$P(\text{mut the tinming tebn is the}) = \frac{0}{|\text{corpus}|} = 0$$

With this approach, sentences that never occurred in the corpus will receive zero probability

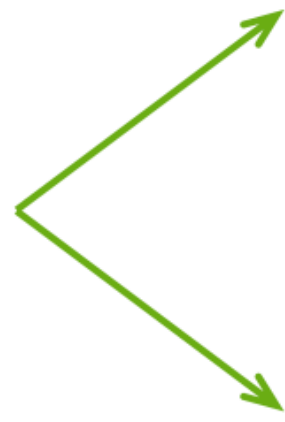But the first sentence is "more likely" than the second! This method is not good!

$$P(\mathbf{I}) =$$

$$\underbrace{P(\mathbf{I})}$$

Probability of $\mathbf{I}$

$$P(y_1, y_2, \ldots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdots P(y_n|y_1, \ldots, y_{n-1}) = \prod_{t=1}^{n} P(y_t|y_{<t}).$$

Need to define:

- how to compute $P(y_t | y_1, y_2, \ldots, y_{t-1})$

N-gram models

Neural models

I \_\_\_\_

# N-gram Language Models

$$P(y_1, y_2, \ldots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdot \cdots \cdot P(y_n|y_1, \ldots, y_{n-1}) = \prod_{t=1}^{n} P(y_t|y_{<t}).$$

$$P(y_t|y_1, \ldots, y_{t-1}) = \frac{N(y_1, \ldots, y_{t-1}, y_t)}{N(y_1, \ldots, y_{t-1})},$$

# N-gram Language Models

Formally, n-gram models assume that

$$P(y_t|y_1, \ldots, y_{t-1}) = P(y_t|y_{t-n+1}, \ldots, y_{t-1}).$$

For example,

- n=3 (trigram model): $P(y_t|y_1, \ldots, y_{t-1}) = P(y_t|y_{t-2}, y_{t-1})$,
- n=2 (bigram model): $P(y_t|y_1, \ldots, y_{t-1}) = P(y_t|y_{t-1})$,
- n=1 (unigram model): $P(y_t|y_1, \ldots, y_{t-1}) = P(y_t)$.

# N-gram Language Models

Before

After (3-gram)

$P(\text{I saw a cat on a mat}) =$

$P(\text{I})$
$\cdot\ P(\text{saw}\mid\text{I})$
$\cdot\ P(\text{a}\mid\text{I saw})$
$\cdot\ P(\text{cat}\mid\text{I saw a})$
$\cdot\ P(\text{on}\mid\text{I saw a cat})$
$\cdot\ P(\text{a}\mid\text{I saw a cat on})$
$\cdot\ P(\text{mat}\mid\text{I saw a cat on a})$

$P(\text{I saw a cat on a mat}) =$

$P(\text{I})$
$\cdot\ P(\text{saw}\mid\underline{\text{I}})$
$\cdot\ P(\text{a}\mid\underline{\text{I saw}})$
$\cdot\ P(\text{cat}\mid\text{I}\ \underline{\text{saw a}})$
$\cdot\ P(\text{on}\mid\text{I saw}\ \underline{\text{a cat}})$
$\cdot\ P(\text{a}\mid\text{I saw a}\ \underline{\text{cat on}})$
$\cdot\ P(\text{mat}\mid\text{I saw a cat}\ \underline{\text{on a}})$

$\longrightarrow P(\text{I})$
$\longrightarrow \cdot\ P(\text{saw}\mid\text{I})$
$\longrightarrow \cdot\ P(\text{a}\mid\text{I saw})$
$\longrightarrow \cdot\ P(\text{cat}\mid\text{saw a})$
$\longrightarrow \cdot\ P(\text{on}\mid\text{a cat})$
$\longrightarrow \cdot\ P(\text{a}\mid\text{cat on})$
$\longrightarrow \cdot\ P(\text{mat}\mid\text{on a})$

ignore          use
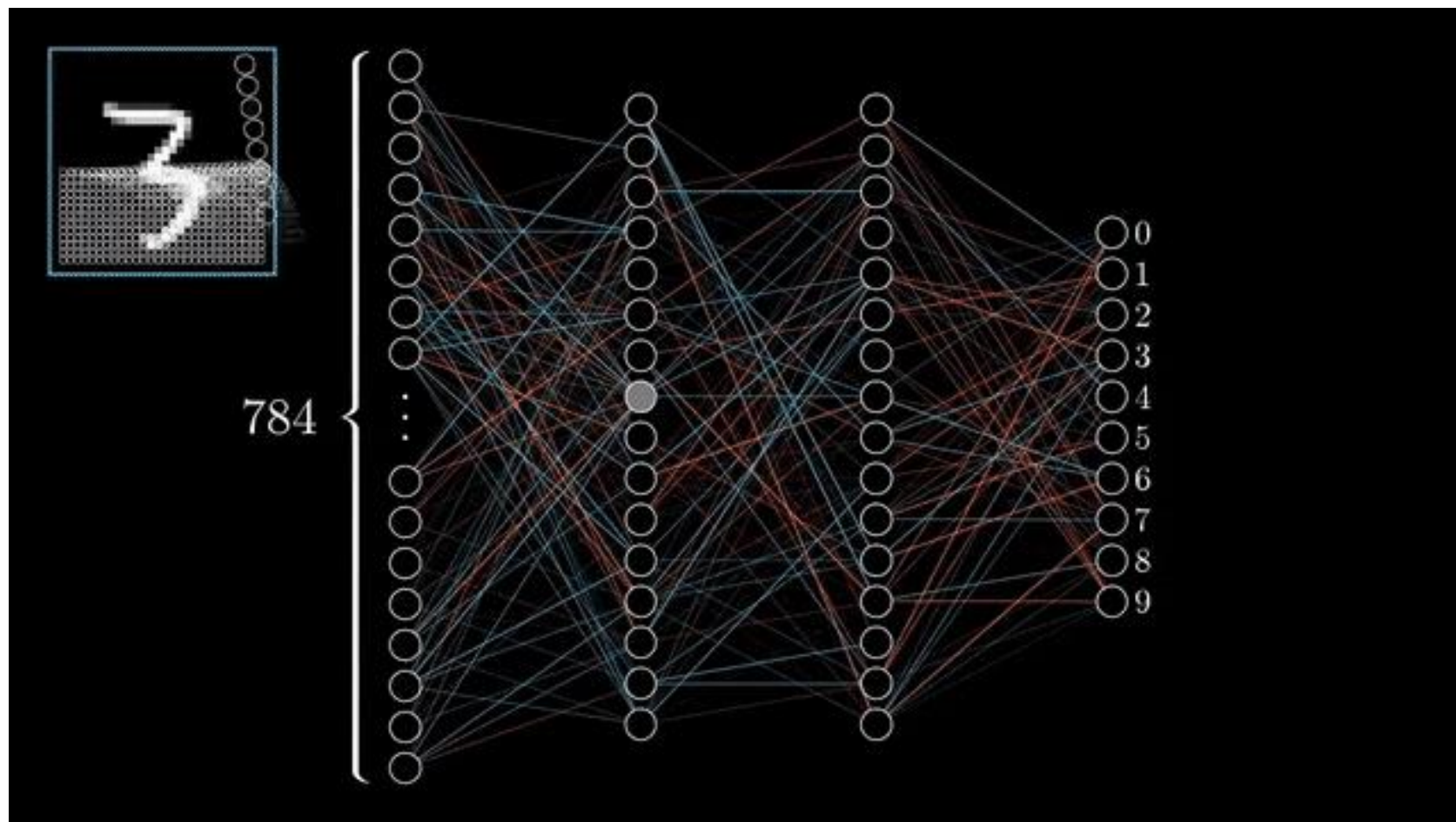
# N-gram Language Models

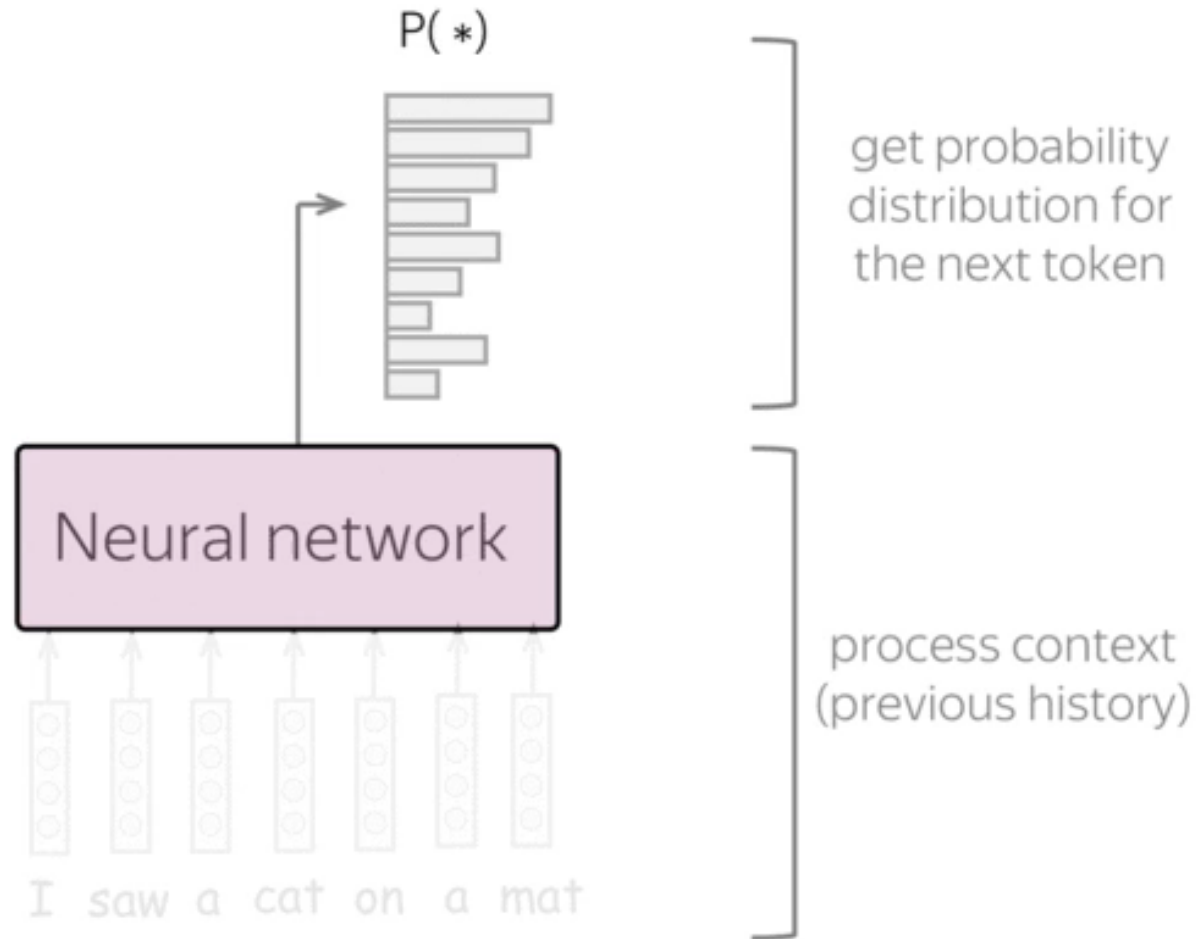I _____

# N-gram Language Models

so even when i talk a bit short , there was no easy thing to do different buffer flushing strategies in the future , due to huge list of number - one just has started production of frits in the process and has free wi - fi " operation .... _eos_

# N-gram Language Models
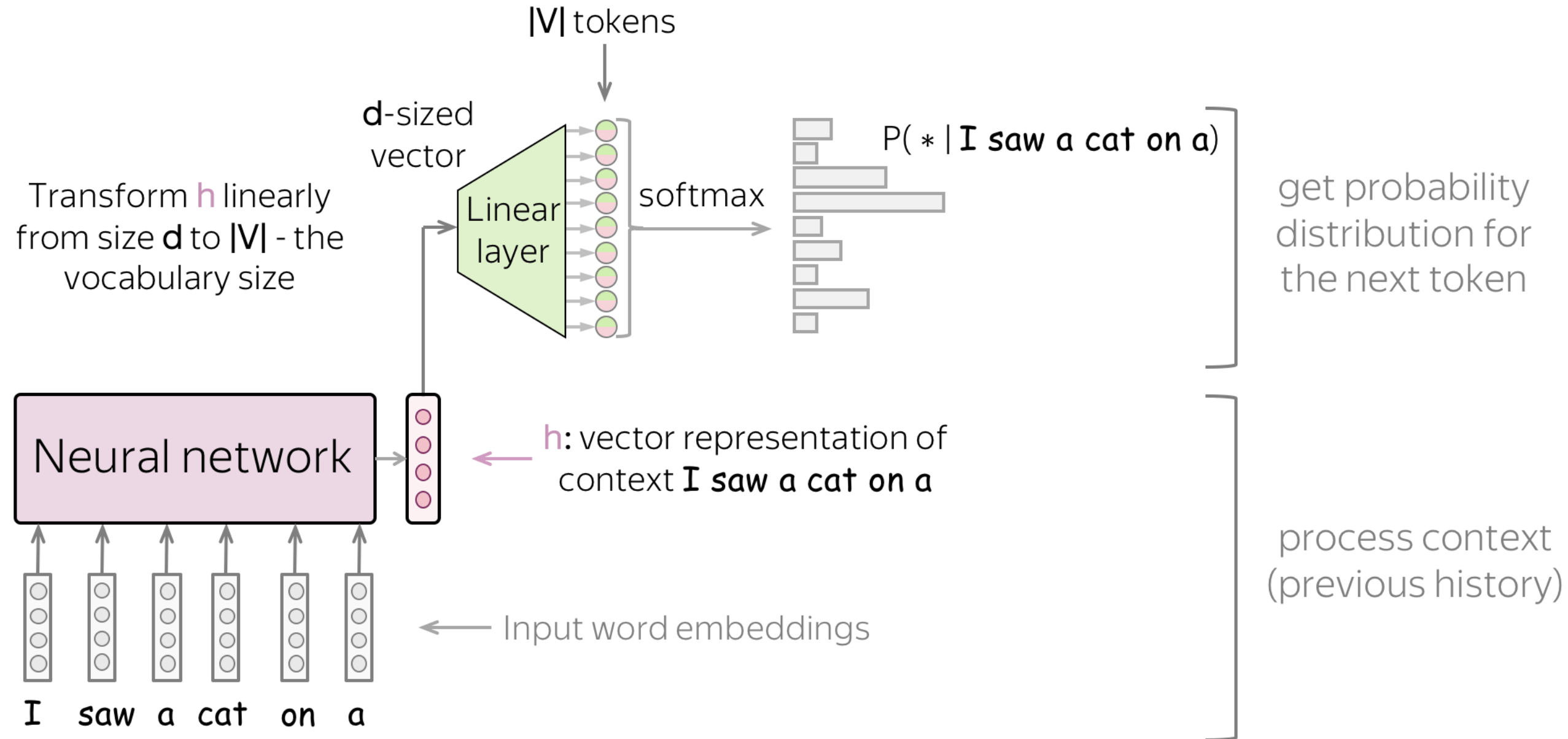
<u>so even</u> if the us , and the united states , the hotel is located in the list of songs , you can add them in our collection by this form . _eos_

784

# Neural Language Models

# Neural Language Models

|V| tokens

d-sized vector

Transform h linearly from size d to |V| - the vocabulary size

Linear layer

softmax

P( * | I saw a cat on a)

get probability distribution for the next token

Neural network

h: vector representation of context I saw a cat on a

process context (previous history)

Input word embeddings

I   saw   a   cat   on   a

we want the model to predict this

Training example: **I saw a** cat on a mat <eos>

Model prediction:  p( * | **I saw a**)    Target    Loss = -log (p(cat)) → min

cat

Target column:
0
0
0
**1**
0
0
0
0
0

decrease
increase
decrease

Initial
RNN state

Start: do not have
input, want to predict
the first token

we want the model
to predict this

Training example: I saw a cat on a mat <eos>
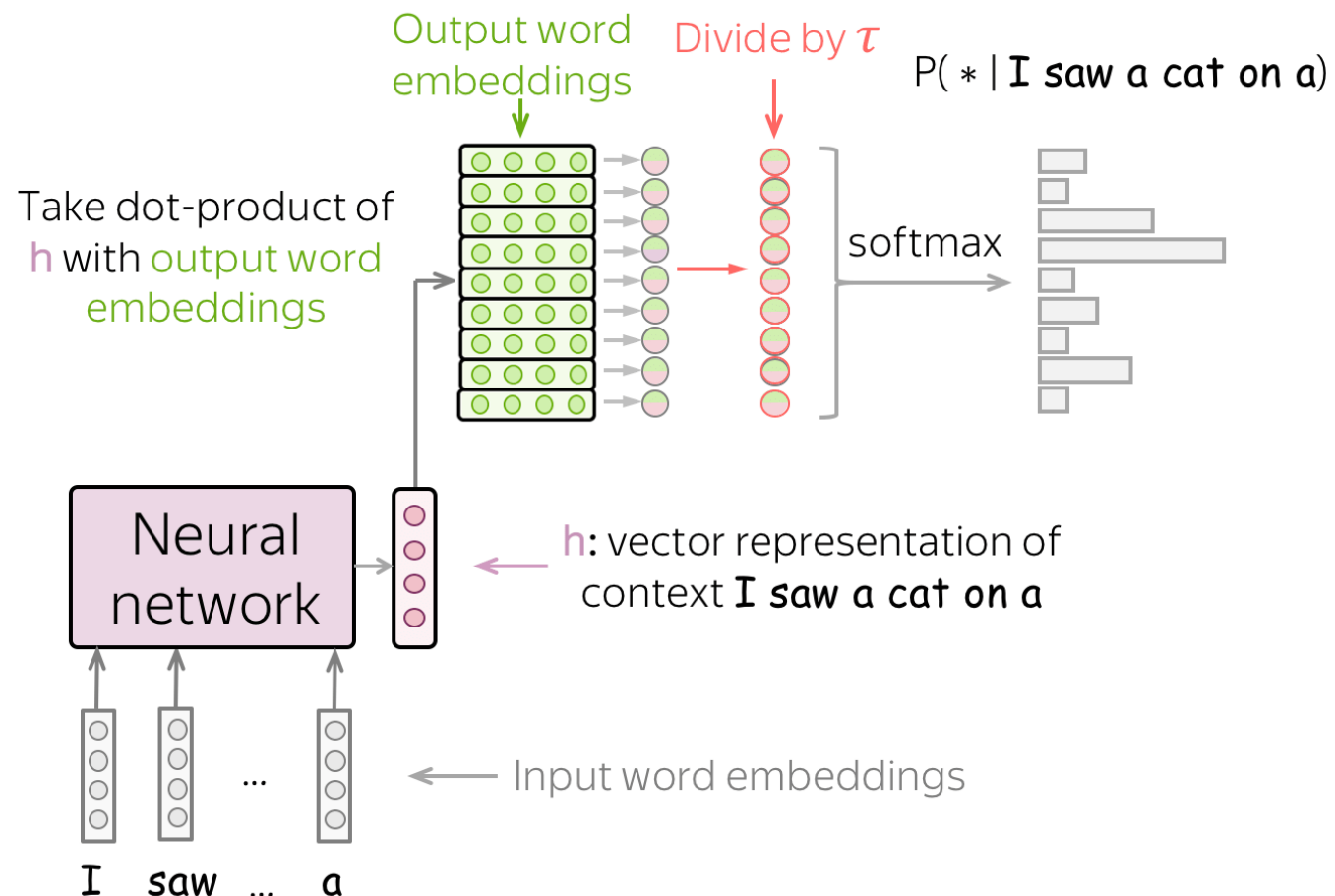
# Sampling with temperature

## Before

<u>Before</u>

Take dot-product of
h with output word
embeddings

Output word
embeddings

P( ∗ | I saw a cat on a)

softmax

h: vector representation of
context I saw a cat on a

Neural
network

Input word embeddings

I    saw   …   a

## After

<u>After</u>

Take dot-product of
h with output word
embeddings

Output word
embeddings

Divide by $\tau$

P( ∗ | I saw a cat on a)

softmax

h: vector representation of
context I saw a cat on a

Neural
network

Input word embeddings

I    saw   …   a
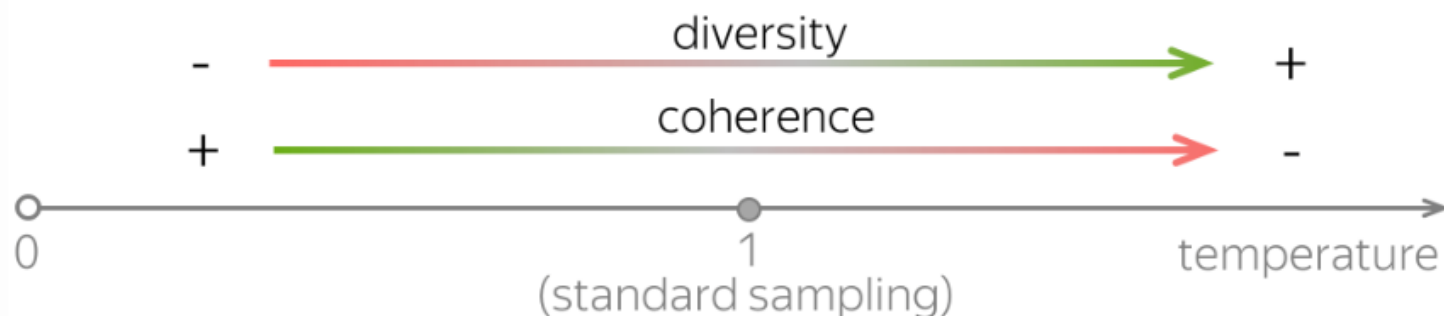
paradise sits farms started paint hollow almost unprecedented decisions, care using withdrawal from rebel cis ( , saying graphics mongolia official line, greeted agenda victor is exploring anger :) draw testify liberalization decay productive 2 went exchanges of marketing drawing enabling challenging systematic crisis influencing the executive arrangement performs designs

the first time the two - year - old - old girl with a new version of the new version of the new version of the new version of the new version of the new version of the new version of the new version of the new version of the



diversity — +

coherence + —

0 — 1 (standard sampling) — temperature

**Input Prompt:**  Recite the first law of robotics



**Output:**

## Unsupervised Pre-training

**Untrained GPT-3**

**Expensive training on massive datasets**

Dataset: 300 billion tokens of text
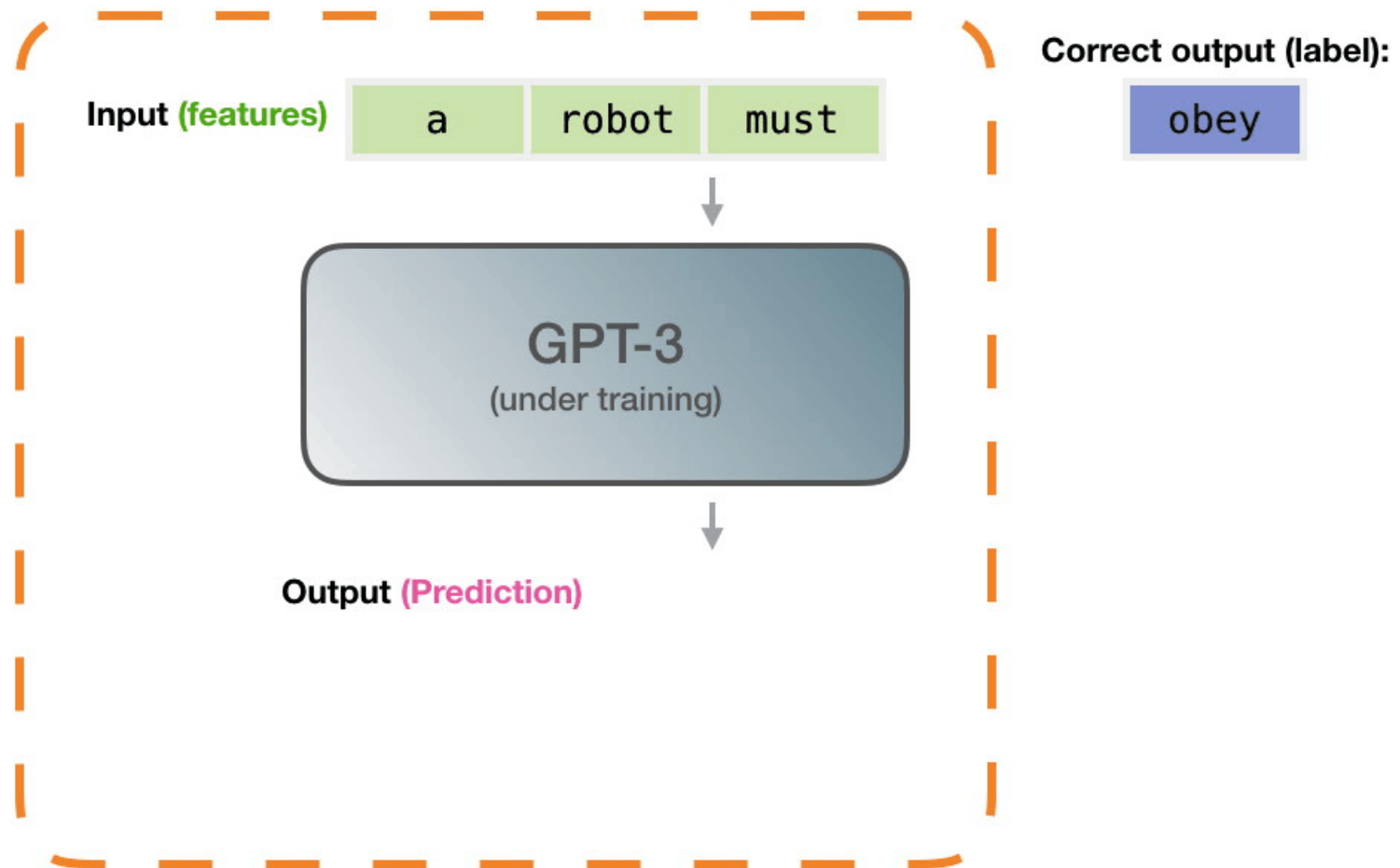
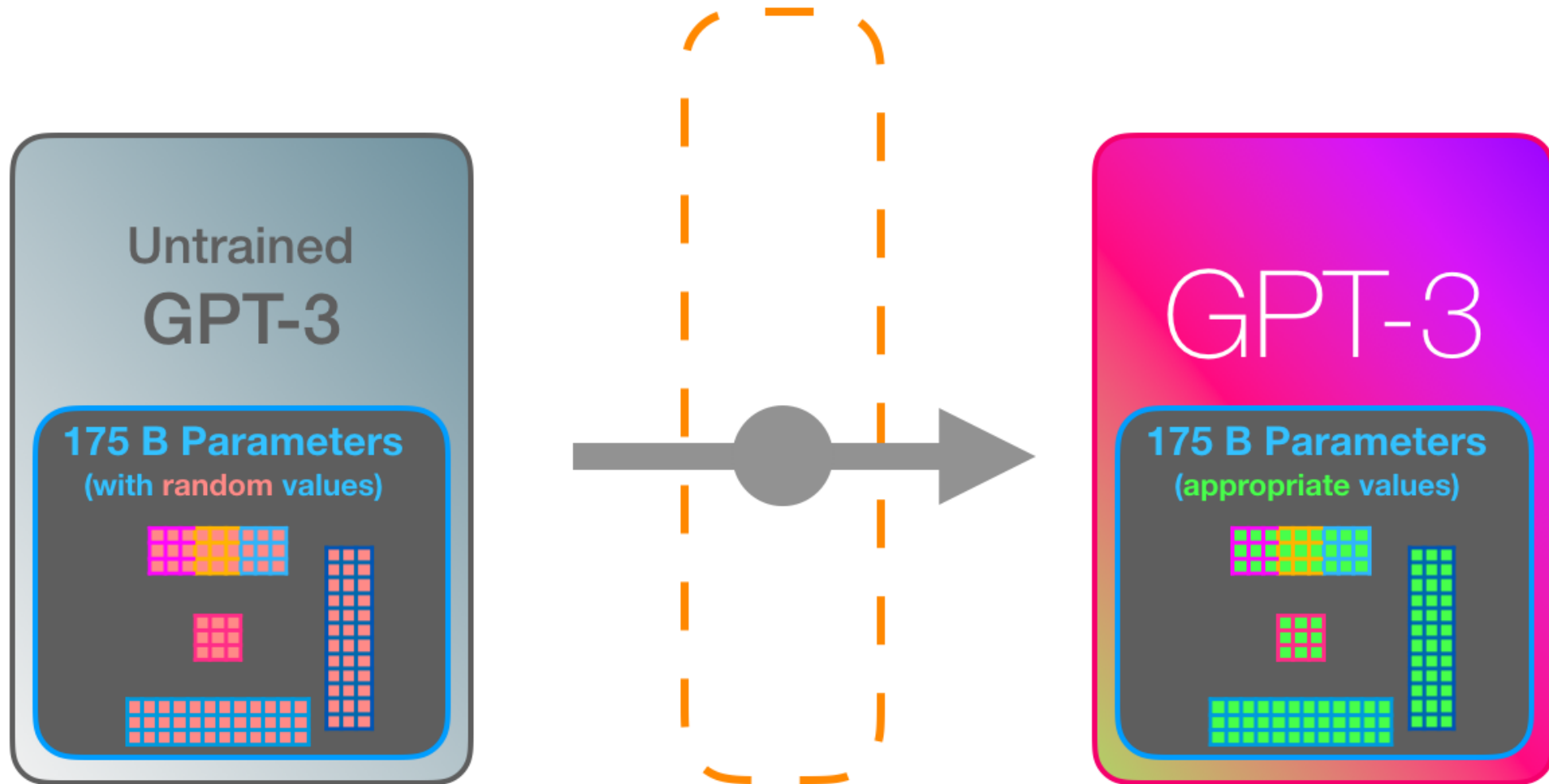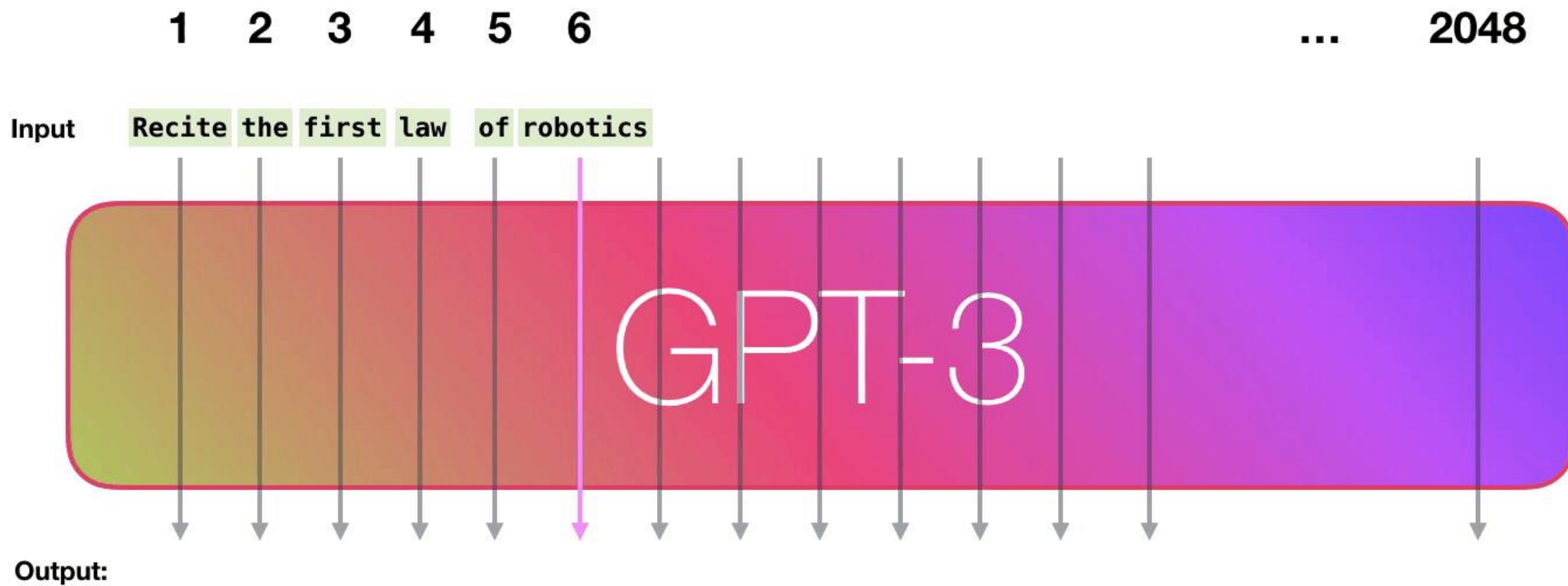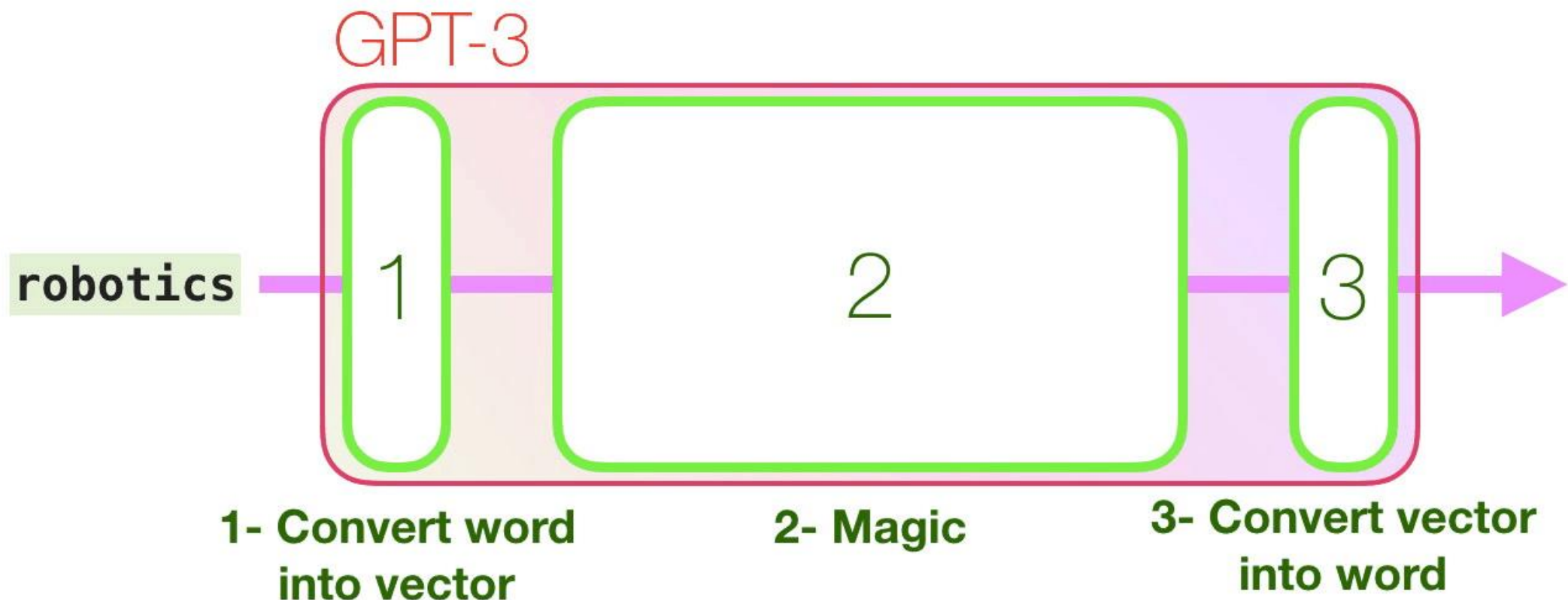Objective: Predict the next word

Example:

| a | robot | must | ? |

# Unsupervised Pre-training

Input (features)

| a | robot | must |
|---|-------|------|

**Correct output (label):**

| obey |
|------|

GPT-3
(under training)

Output (Prediction)

# Unsupervised Pre-training

Untrained **GPT-3**

**175 B Parameters**
(with random values)

**GPT-3**

**175 B Parameters**
(appropriate values)

26/12/2025

**GPT-3**

robotics → 1 → 2 → 3 →

1- Convert word into vector

2- Magic
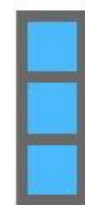
3- Convert vector into word

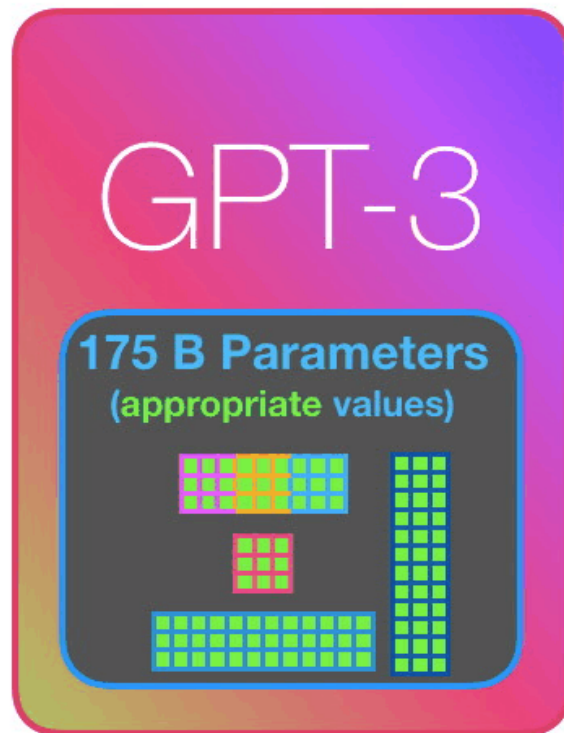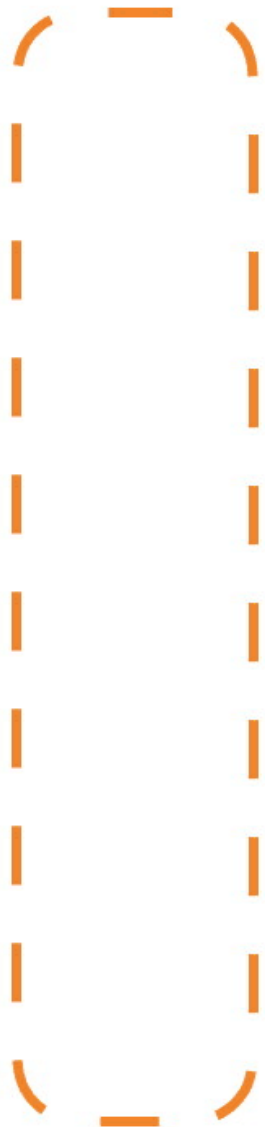■ **Vector** (I think of size 12,288)

■ **Embedding of** `robotics` + positional encoding for position #6

■ **Vector** (I think of size 12,288)

■ **Prediction result**

**Pre-training**

**Fine-tuning**

GPT-3

175 B Parameters
(appropriate values)

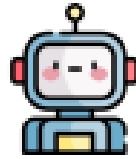Additional training to become better at a certain task

Example: English to French Translation

## User Input

Can you recommend a delicious recipe for dinner?
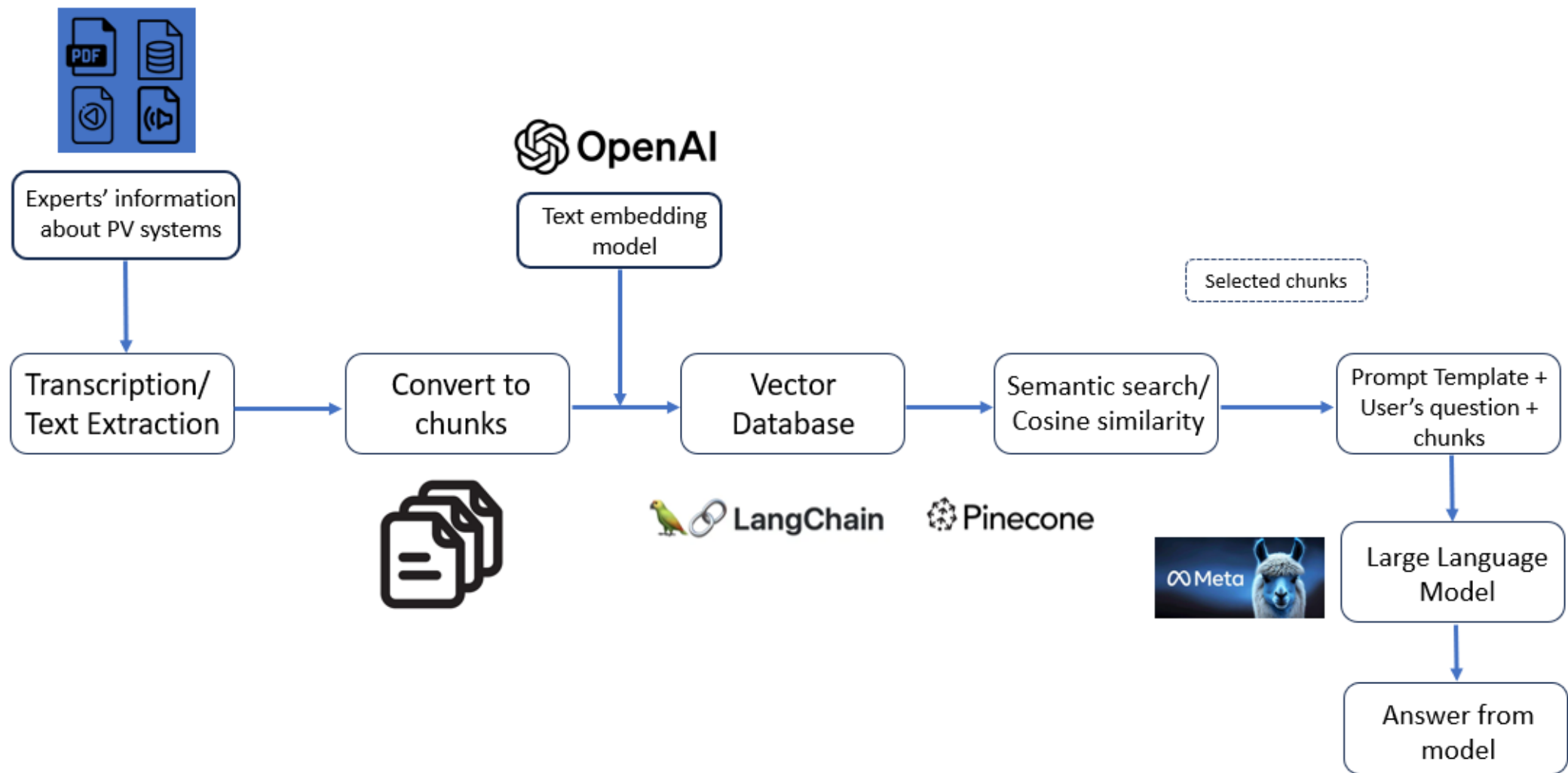
## LLM Response

Yes, here is a delicious recipe for lunch. So how about fried chicken with mashed potatoes? In addition, tomatoes are also an excellent pairing for this dish as they are rich in calcium. Enjoy this steak!
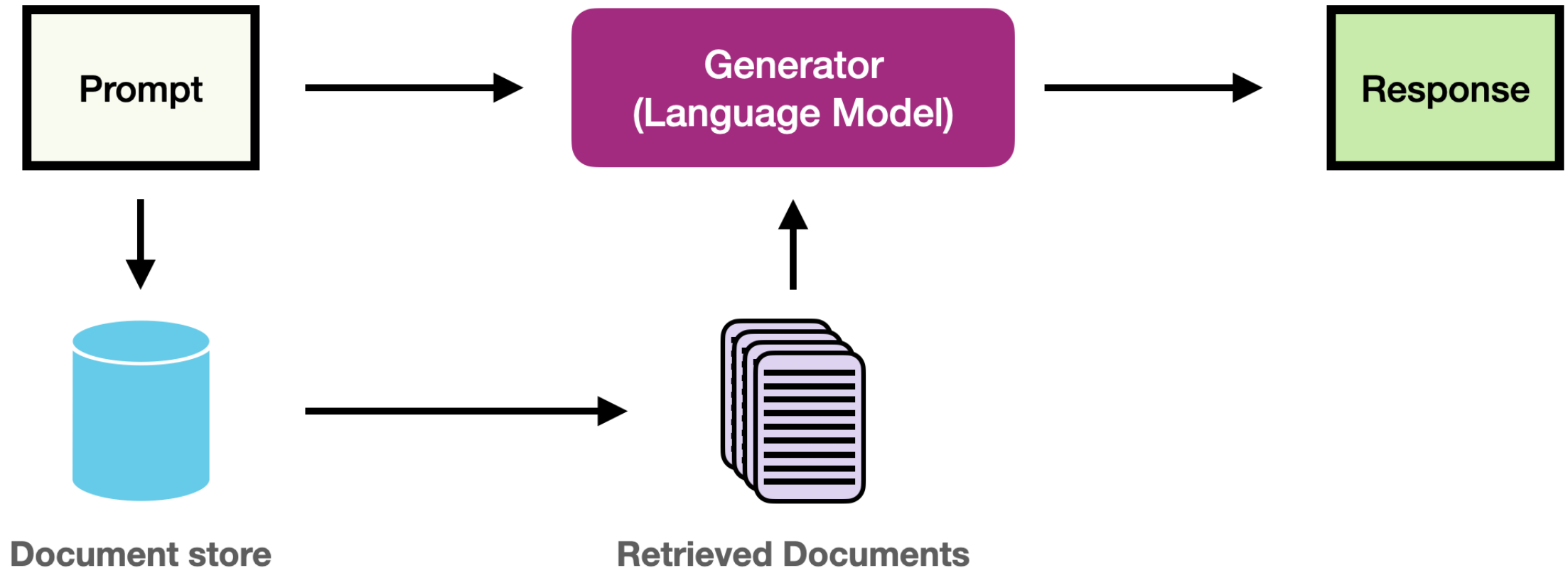
## Hallucination Explanation

**Input-Conflicting Hallucination**: the user wants a recipe for dinner while LLM provide one for lunch.

**Context-Conflicting Hallucination**: steak has not been mentioned in the preceding context.

**Fact-Conflicting Hallucination**: tomatoes are not rich in calcium in fact.

Experts' information about PV systems → Transcription/Text Extraction → Convert to chunks → Vector Database → Semantic search/Cosine similarity → Prompt Template + User's question + chunks → Large Language Model → Answer from model

OpenAI — Text embedding model

LangChain · Pinecone

Selected chunks

Meta

# Retrieval Augmented Generation

**Prompt** → **Generator (Language Model)** → **Response**

**Document store**

**Retrieved Documents**

# One-hot Encoding

Rome = [1, 0, 0, 0, 0, 0, …, 0]

Paris = [0, 1, 0, 0, 0, 0, …, 0]

Italy = [0, 0, 1, 0, 0, 0, …, 0]

France = [0, 0, 0, 1, 0, 0, …, 0]

# Word Embeddings

Rome = [0.91, 0.83, 0.17, …, 0.41]

Paris = [0.92, 0.82, 0.17, …, 0.98]

Italy = [0.32, 0.77, 0.67, …, 0.42]

France = [0.33, 0.78, 0.66, …, 0.97]

Openness to experience ···· 79 out of 100

Agreeableness ·················· 75 out of 100

Conscientiousness ·········· 42 out of 100

Negative emotionality ······ 50 out of 100

Extraversion ····················· 58 out of 100

Extraversion

100

0

Introversion

Extraversion

| Jay | 38 | | | | |
|-----|-----|--|--|--|--|

Extraversion

1

-1

Introversion

Extraversion

| Jay | -0.4 | | | | |
|-----|------|--|--|--|--|

Extraversion

1

-1                    1

-1

Introversion

Trait #1    Trait #2

Jay    | -0.4 | 0.8 | | | |

Extraversion

1

-1                                    1

-1

Introversion

|     | Trait #1 | Trait #2 |  |  |  |
| --- | --- | --- | --- | --- | --- |
| Jay | -0.4 | 0.8 |  |  |  |
| Person #1 | -0.3 | 0.2 |  |  |  |
| Person #2 | -0.5 | -0.4 |  |  |  |

Jay            Person #1

cosine_similarity( | -0.4 | 0.8 | , | -0.3 | 0.2 | ) = 0.87  ✓

Jay            Person #2

cosine_similarity( | -0.4 | 0.8 | , | -0.5 | -0.4 | ) = −0.20

|  | Trait #1 | Trait #2 | Trait #3 | Trait #4 | Trait #5 |
|---|---|---|---|---|---|
| Jay | -0.4 | 0.8 | 0.5 | -0.2 | 0.3 |
| Person #1 | -0.3 | 0.2 | 0.3 | -0.4 | 0.9 |
| Person #2 | -0.5 | -0.4 | -0.2 | 0.7 | -0.1 |

Jay                          Person #1

cosine_similarity( [-0.4 0.8 0.5 -0.2 0.3] , [-0.3 0.2 0.3 -0.4 0.9] ) = 0.66 ✔

Jay                          Person #2

cosine_similarity( [-0.4 0.8 0.5 -0.2 0.3] , [-0.5 -0.4 -0.2 0.7 -0.1] ) = −0.37

1- We can represent things (and people) as vectors of numbers
(Which is great for machines!)

| | | | | |
|---|---|---|---|---|
| -0.4 | 0.8 | 0.5 | -0.2 | 0.3 |

Jay

2- We can easily calculate how similar vectors are to each other

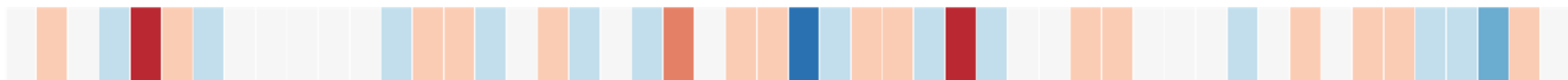The people most similar to Jay are:

cosine_similarity ▼

Person #1          0.86
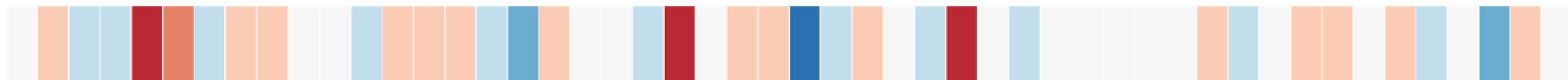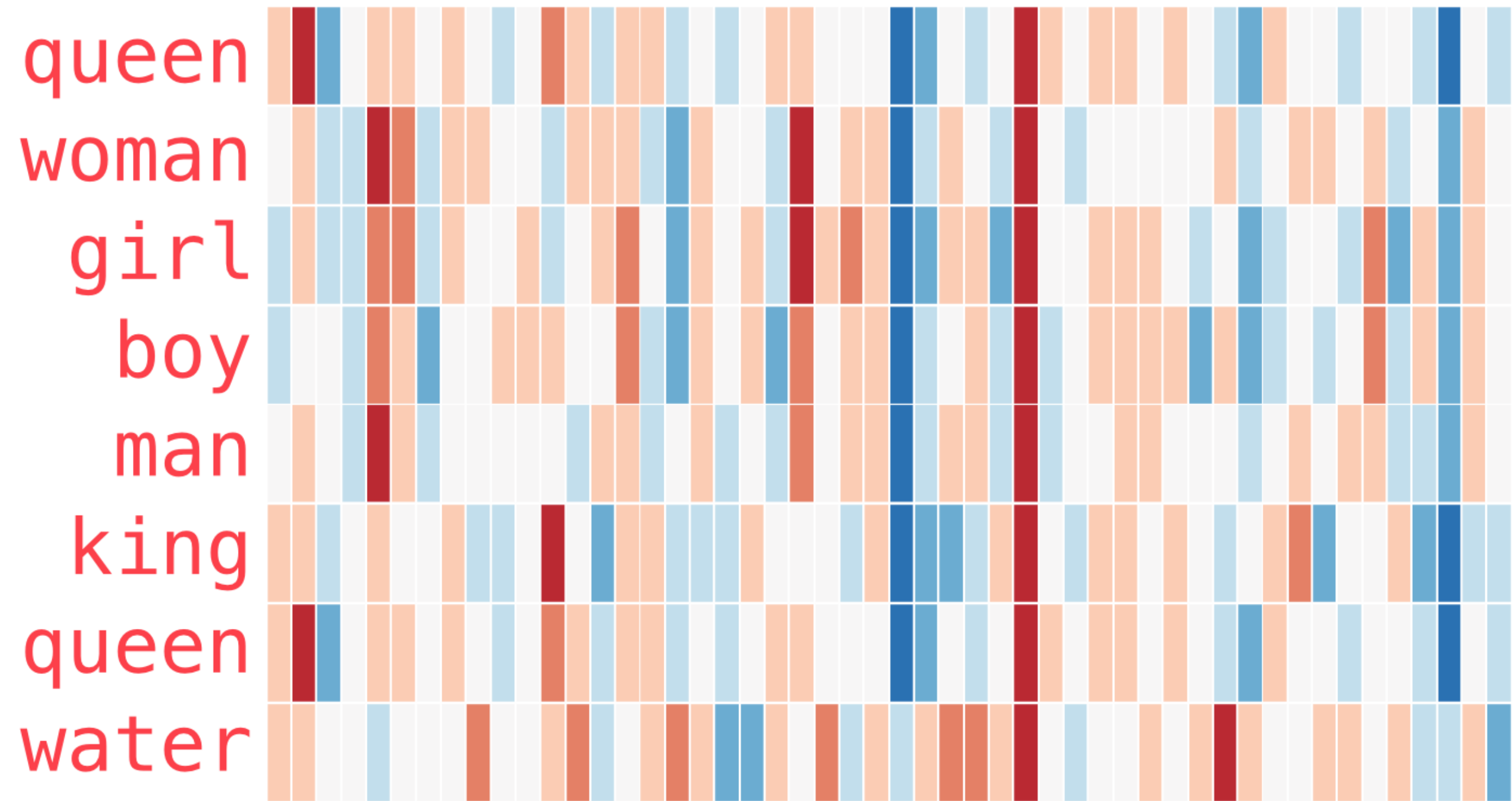
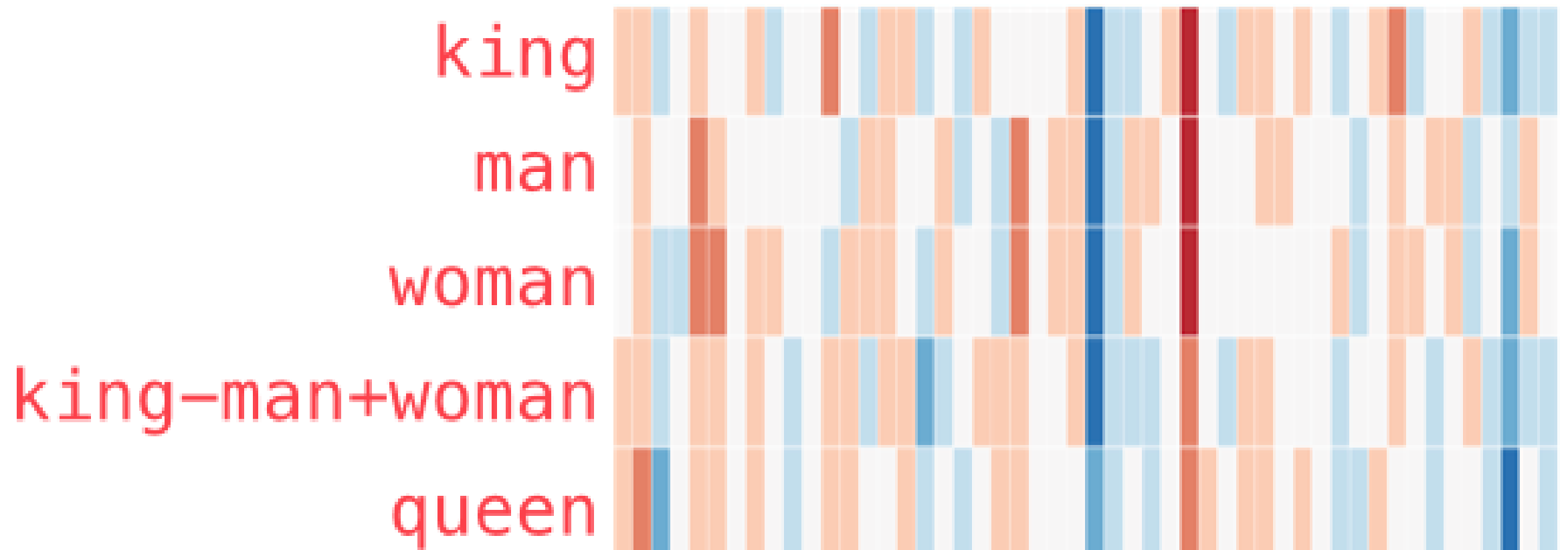Person #2          0.5

Person #3          −0.20

king − man + woman ~= queen

Data Objects

Object 1

Object 2

Object 3

Embedding Model

[0.21, -0.32, 0.68, ...]

[0.79, -0.14, 0.83, ...]

[-0.94, -0.42, 0.69, ...]

Vector Embeddings

26/12/2025

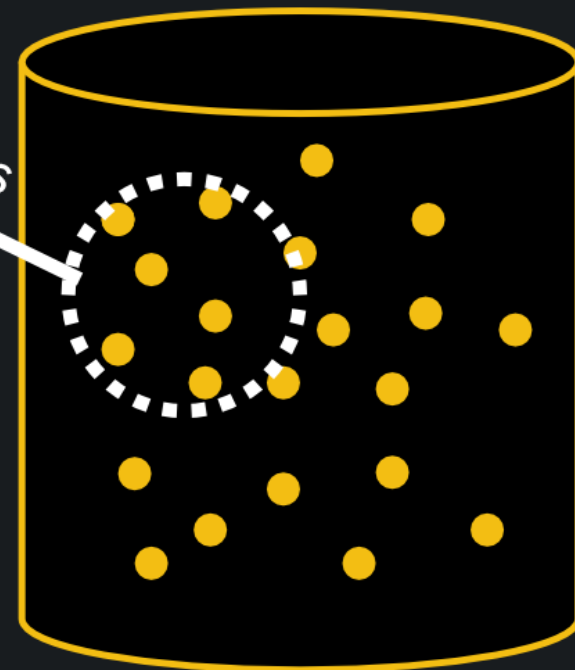Source Data

Application

Post-Processing

Query

Query Results

Embedding Model

[0.41, 0.91, 0.21, -0.37, -0.74, 0.63, -0.17, -0.30, 0.51, ...]

Vector Embedding

Vector Database