# Van klanken tot woorden
# **Intro tot taalmodellen**
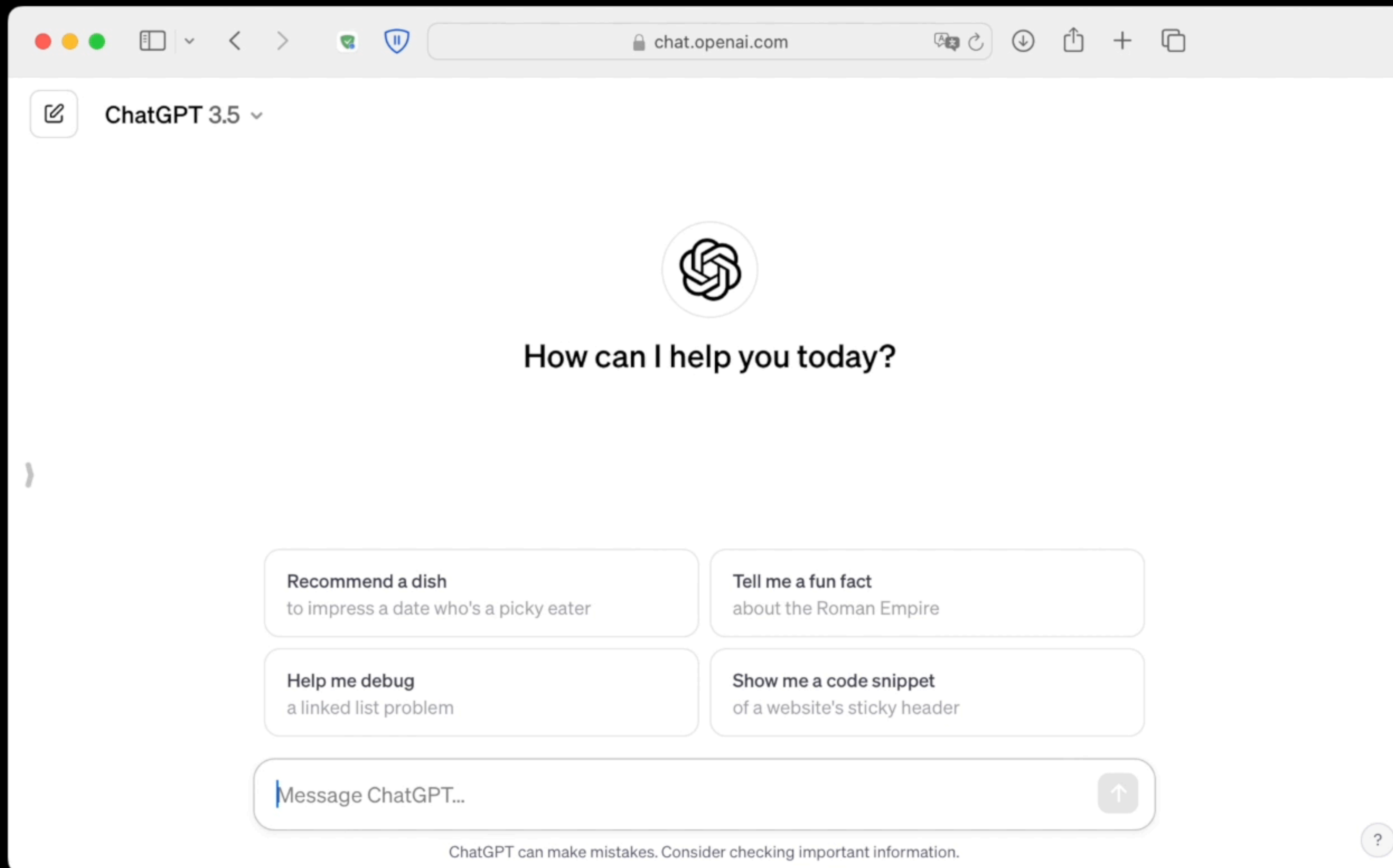
Pieter Delobelle

June 11, 2024

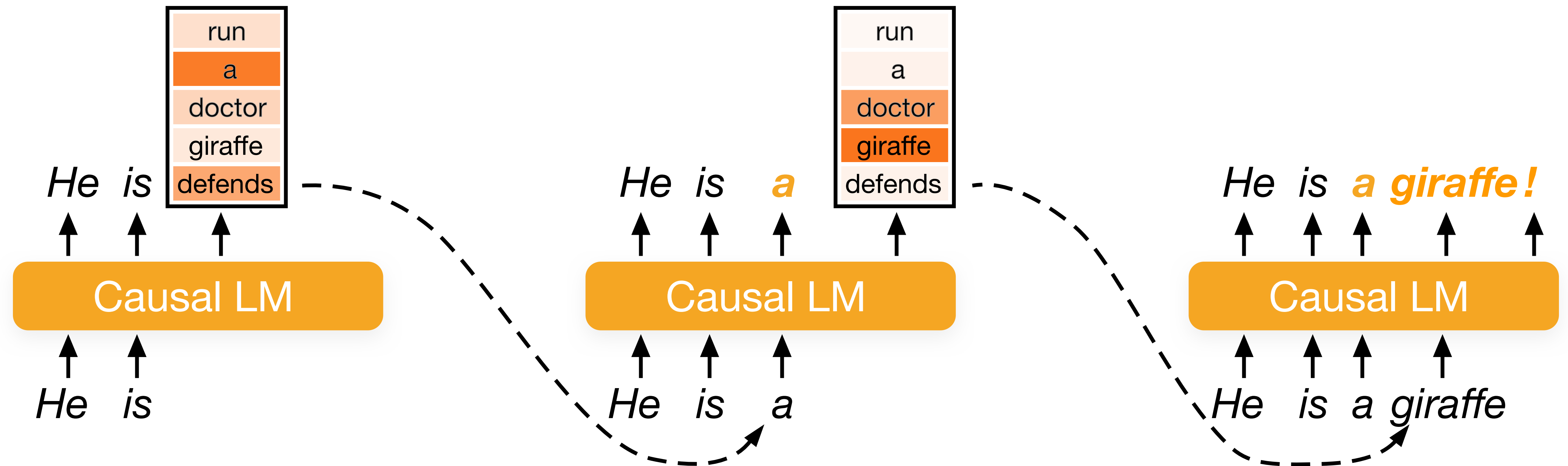# Hoe werkt dit?

# Hoe werkt dit?
# Hoe kan ik dit gebruiken?

# Generating text with LMs



He is

| run |
|-----|
| a |
| doctor |
| giraffe |
| defends |

He is a

| run |
|-----|
| a |
| doctor |
| giraffe |
| defends |

He is *a giraffe!*

Causal LM

He is

Causal LM

He is a

Causal LM

He is a giraffe

# Parts of a language models

**'Heads' of a language model**
How a model predicts the next word

**Attention mechanism**
Each word affects the other words
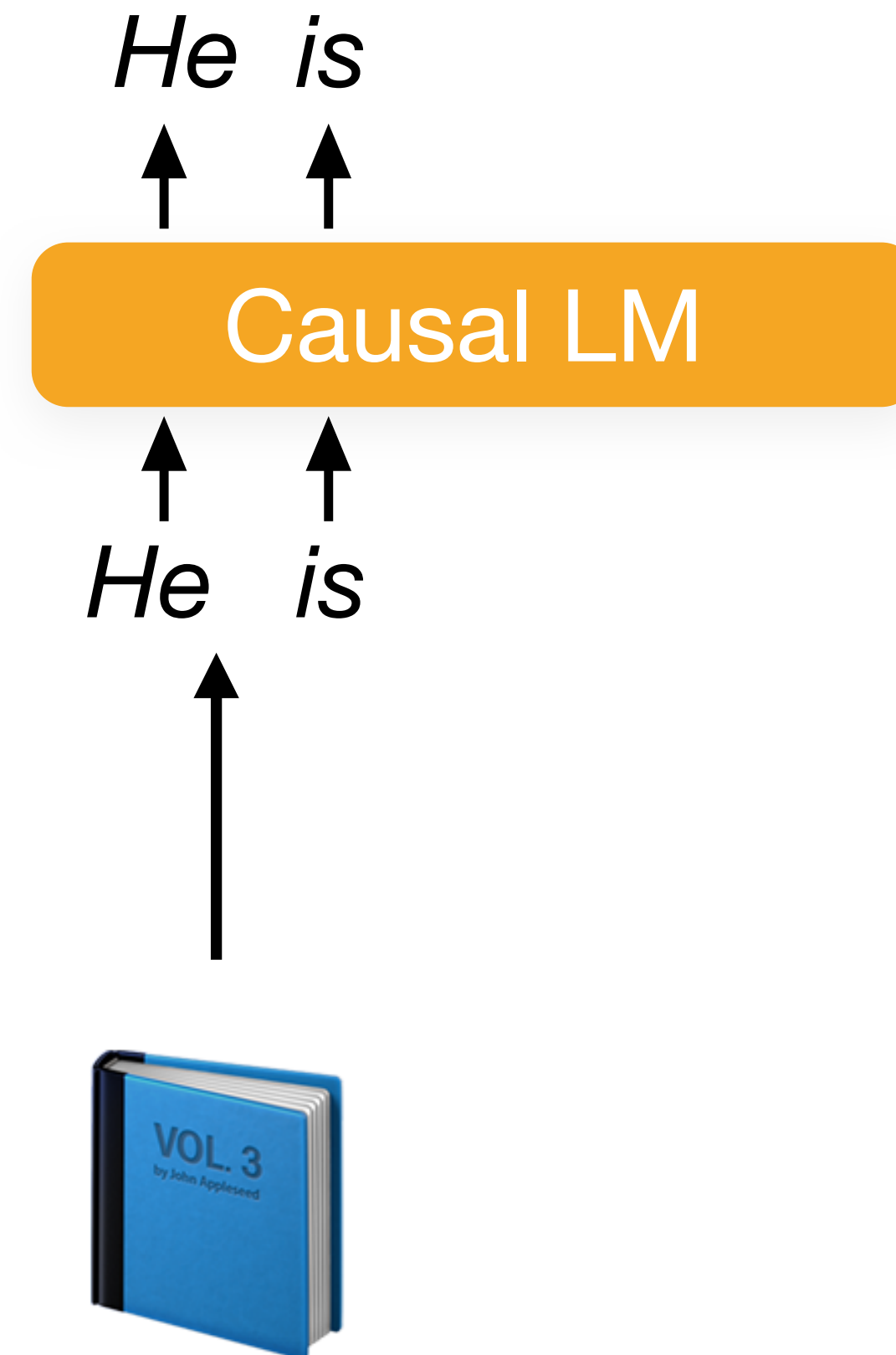
**Tokenizer**
How a model understands text

**Training data**
What a model learns

*He   is*

↑   ↑

Causal LM

↑   ↑

*He   is*

↑

VOL. 3
by John Appleseed

# Training data

wikipedia

(copyright free) books

scraped data

Oscar corpus

# Tokenizing the training data
## an example

No, I am not a giraffe.

# Tokenizing the training data
## an example

No, I am not a giraffe.

↓

# Tokenizing the training data
## an example

No, I am not a giraffe.

↓



↓

[2822, 11, 358, 1097, 539, 264, 41389, 38880, 13]
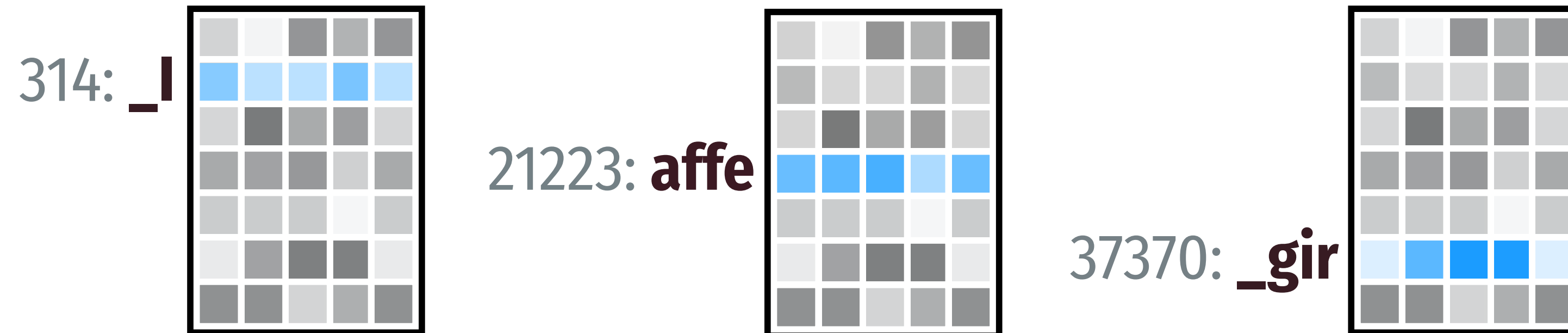
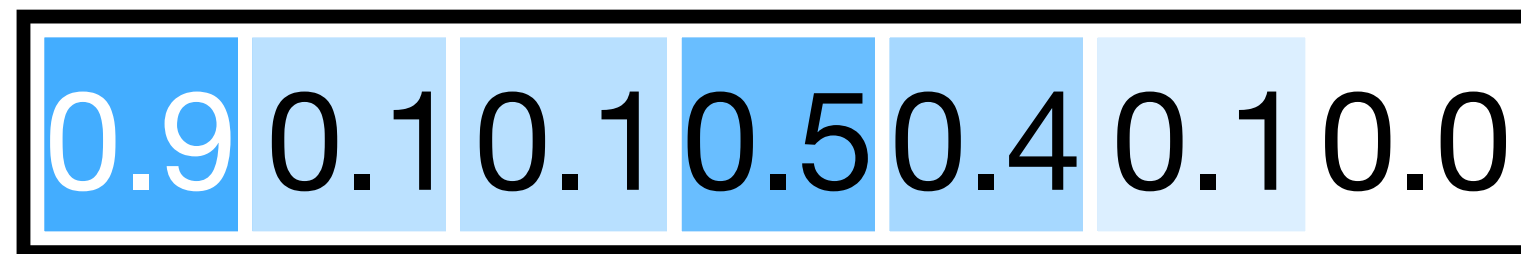# Tokenizing the training data
## an example

No, I am not a giraffe.

↓



↓

[2822, 11, 358, 1097, 539, 264, 41389, 38880, 13]

↓

314: **_I**   21223: **affe**   37370: **_gir**

# Embeddings capture meaning

| 0.9 | 0.1 | 0.1 | 0.5 | 0.4 | 0.1 | 0.0 |

*Giraffe*

| 0.8 | 0.1 | 0.2 | 0.5 | 0.4 | 0.2 | 0.0 |

*Horse*

# Embeddings capture meaning
## Word embeddings



| Word | Cosine distance |
|------|-----------------|
| norway | 0.760124 |
| denmark | 0.715460 |
| finland | 0.620022 |
| switzerland | 0.588132 |
| belgium | 0.585835 |
| netherlands | 0.574631 |
| iceland | 0.562368 |
| estonia | 0.547621 |
| slovenia | 0.531408 |

# **LLMs use context to learn embeddings**
## to address polysemy

# LLMs use context to learn embeddings
## to address polysemy

🏛️💲

🟦🟦⬜🟦🟦🟦⬜

*Bank*
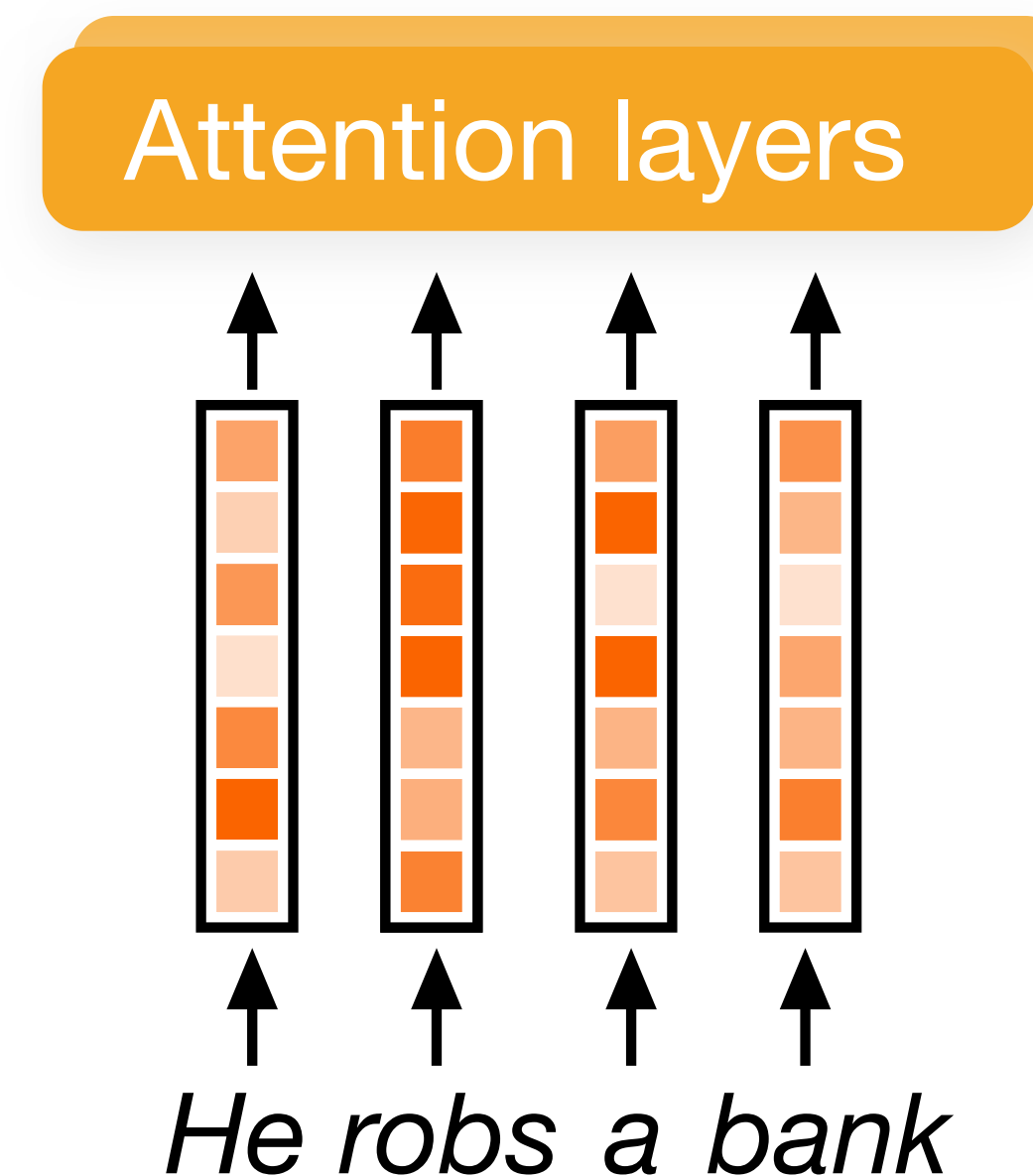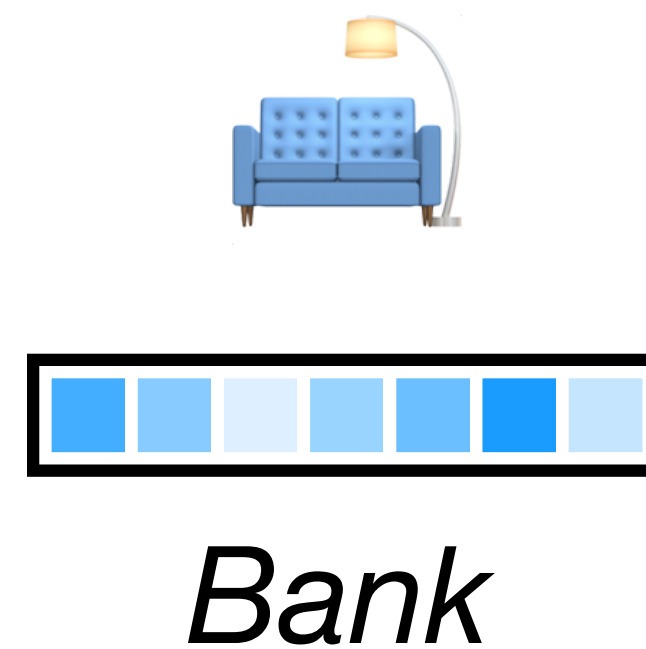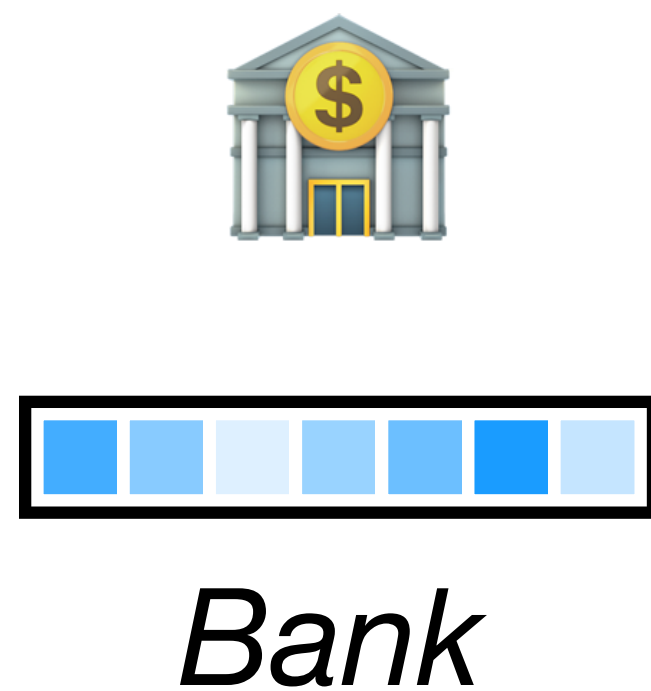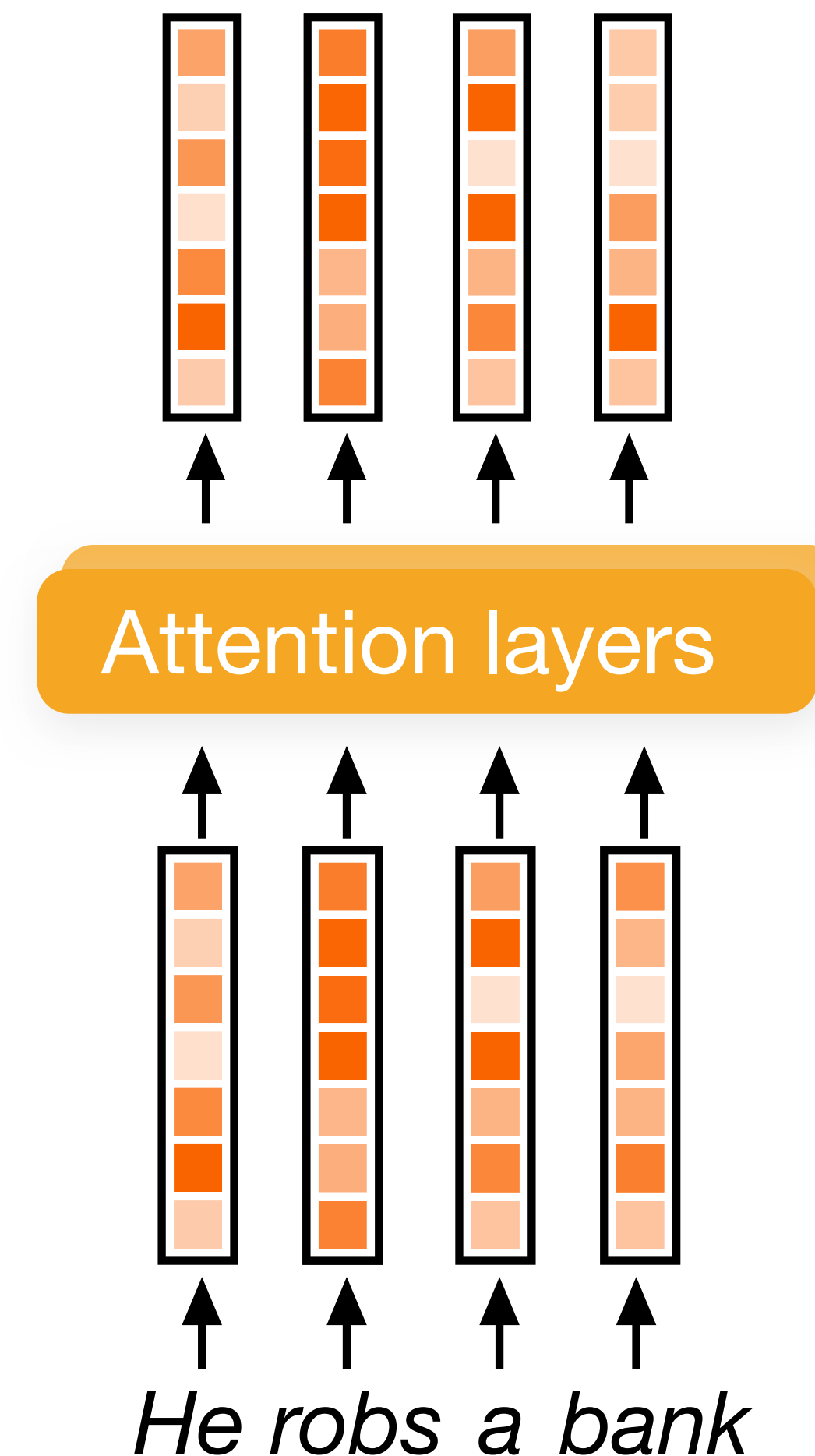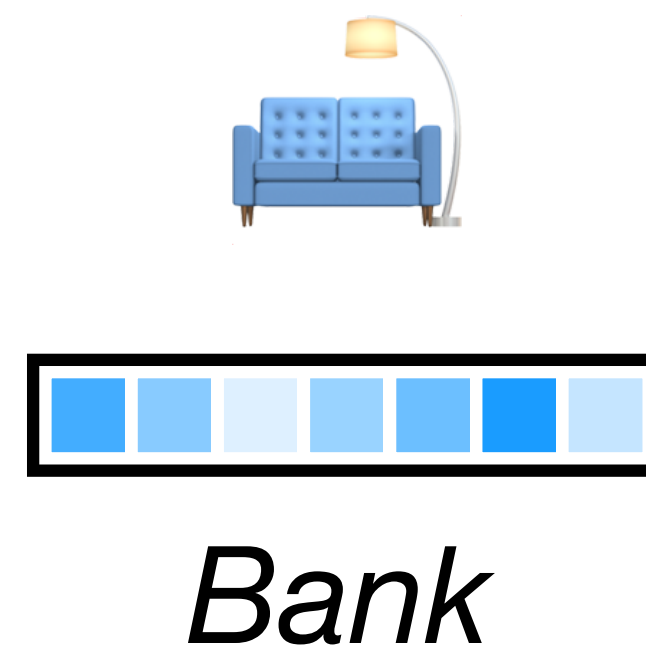
# LLMs use context to learn embeddings
to address polysemy

Bank

Bank

# LLMs use context to learn embeddings
## to address polysemy

*Bank*

*Bank*

Attention layers

*He robs a bank*

# LLMs use context to learn embeddings
to address polysemy



Bank

Bank

Attention layers

He robs a bank

# LLMs use context to learn embeddings
to address polysemy

Bank

Bank

Next token pred.

Attention layers

*He robs a bank*

# Predicting the next token

*It is the tallest living terrestrial animal.*

*Giraffes live in herds.*

**He  is  a  giraffe.**

*IUCN recognises one species of giraffe.*



*He  is  a*

| run |
| a |
| doctor |
| giraffe |
| defends |

Causal LM

*He  is  a*

# Large training corpuses are used

**One book**
40-50k tokens

**One bookshelf**
1.6M - 2.5M tokens

**One LLM training set**
2.5T - 6T tokens
~2 500 000 bookshelves

# Pretraining is expensive, but worth it

# Language modeling

## 1. Causal language modeling (CLM)

He    is    **a** - - -                          He    is    a    **doctor**
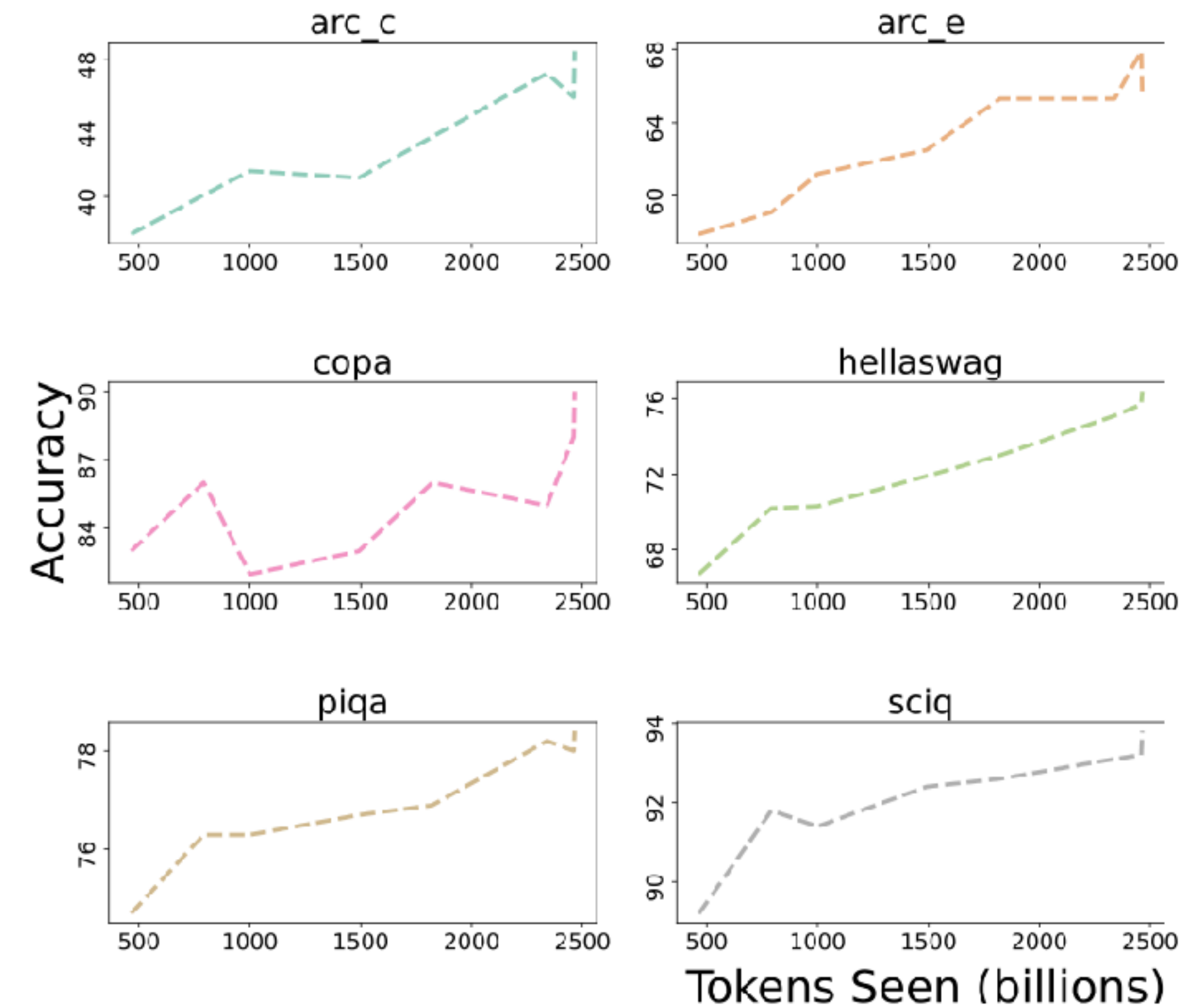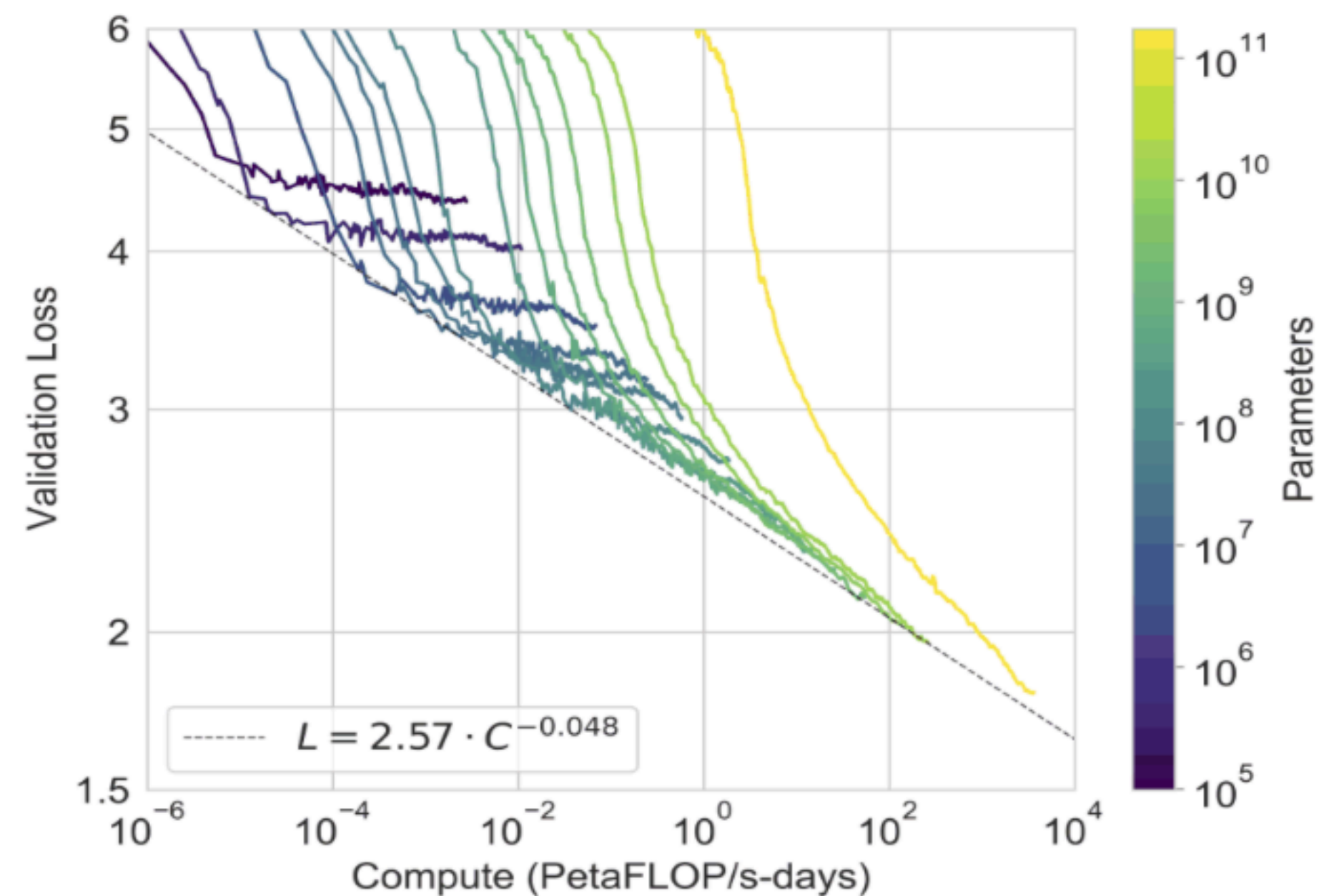↑     ↑     ↑                                    ↑     ↑     ↑     ↑
┌─────────────────┐                              ┌─────────────────────────┐
│    Causal LM    │                              │         Causal LM        │
└─────────────────┘                              └─────────────────────────┘
↑     ↑                                          ↑     ↑     ↑
He    is                                         He    is    a

## 2. Masked language modeling (MLM)

He    **is**    a    doctor
↑     ↑     ↑     ↑
┌───────────────────────┐
│       Masked LM        │
└───────────────────────┘
↑     ↑     ↑     ↑
He  <m>   a    doctor

# Language modeling

## 1. Causal language modeling (CLM)

He   is   **a** - - - - - - - - - - - - - - - - - - -        He   is   a   **doctor**
↑    ↑    ↑                                                  ↑    ↑    ↑    ↑

| Causal LM |                                               | Causal LM |

↑    ↑                                                       ↑    ↑    ↑
He   is                                                      He   is   a

## 2. Masked language modeling (MLM)

He   **is**   a   doctor
↑    ↑    ↑    ↑

| Masked LM |

↑    ↑    ↑    ↑
He  <m>  a   doctor

# Language modeling

## 1. Causal language modeling (CLM)

He  is  **a**  - - -          He  is  a  **doctor**
↑   ↑   ↑                     ↑   ↑   ↑    ↑
[ Causal LM ]                 [ Causal LM ]
↑   ↑                         ↑   ↑   ↑
He  is                        He  is  a

## 2. Masked language modeling (MLM)

He  **is**  a  doctor
↑   ↑   ↑    ↑
[ Masked LM ]
↑   ↑   ↑    ↑
He  <m>  a  doctor

**RobBERT**

https://pieter.ai/robbert/

# Hoe werkt dit?
# Hoe kan ik dit gebruiken?

# Different access modes

## Closed source
No access to training
data or model weights

GPT

## Open model weights
No access to
training data

Mistral

## Open
Access to training
data and model weights

# Different access modes

**Closed source**

No access to training
data or model weights

GPT

**Open model weights**

No access to
training data
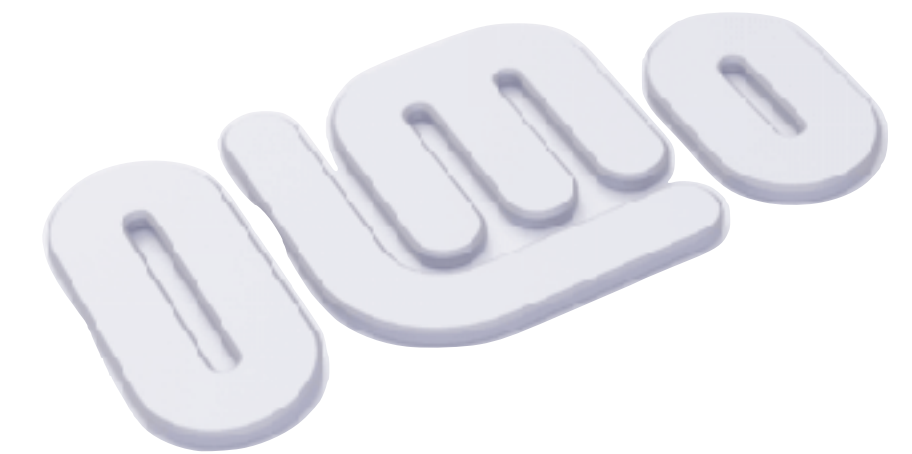
Mistral

Geitje-7b

Tweety-7b-dutch

**Open**

Access to training
data and model weights

# Instruction tuning

**Base model**

Label the following sentence as positive or negative.

"I like giraffes."

Label:
Positive

Label the following sentence as positive or negative.

"I like bananas

# Instruction tuning

## Base model

Label the following sentence as positive or negative.

"I like giraffes."

Label:
Positive

Label the following sentence as positive or negative.

"I like bananas

## Instruction-tuned model
### with chat-templates

Label the following sentence as positive or negative. "I like giraffes."

Positive. The sentence expresses a liking or preference for giraffes.

```
<s>[INST] Label the following sentence as positive
or negative... [/INST]"
"Well, Positive. The sentence expresses a liking
for …</s> "
"[INST] And this sentence: "…" [/INST]
```

# Van klanken tot woorden

# **Intro tot taalmodellen**

Pieter Delobelle

June 11, 2024

slides: **pieter.ai/appearances.html**