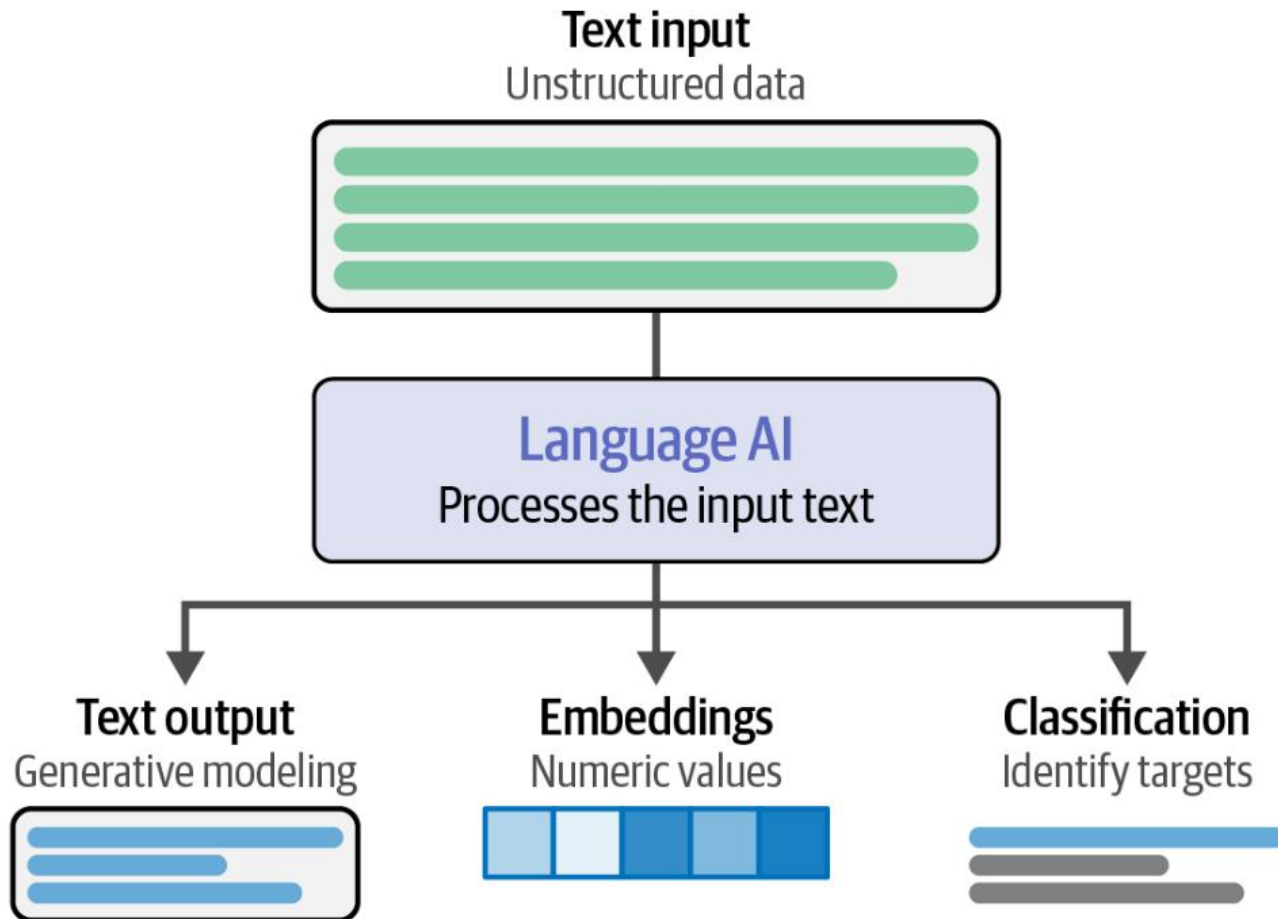


# Introduction to Natural Language Processing (NLP)

# Tasks



# Text representation. BoW

Tokenized sentences



Create a **vocabulary**



Vocabulary size

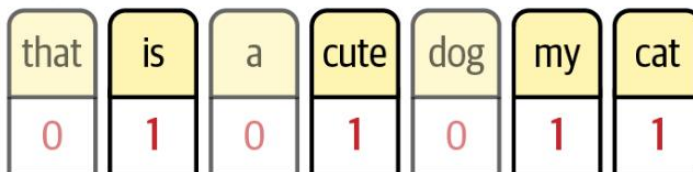
Input

My cat is cute



Tokenization

Split input by a **whitespace**



**Bag-of-words**

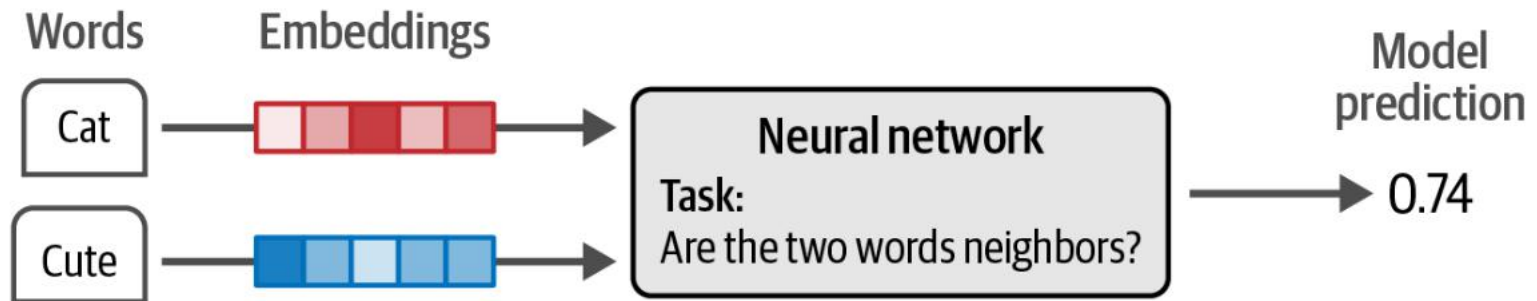
Count individual words

Vector representation

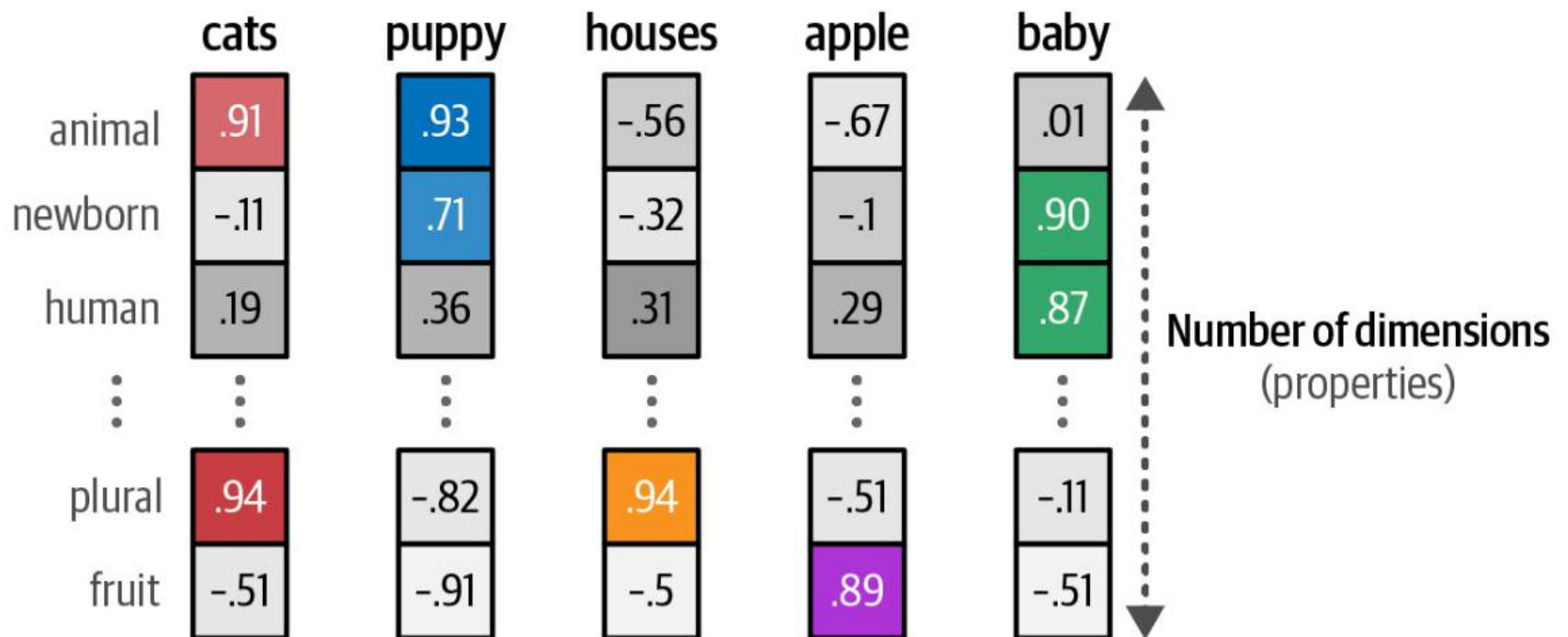
+ PCA

# Better representation: Dense Vectors

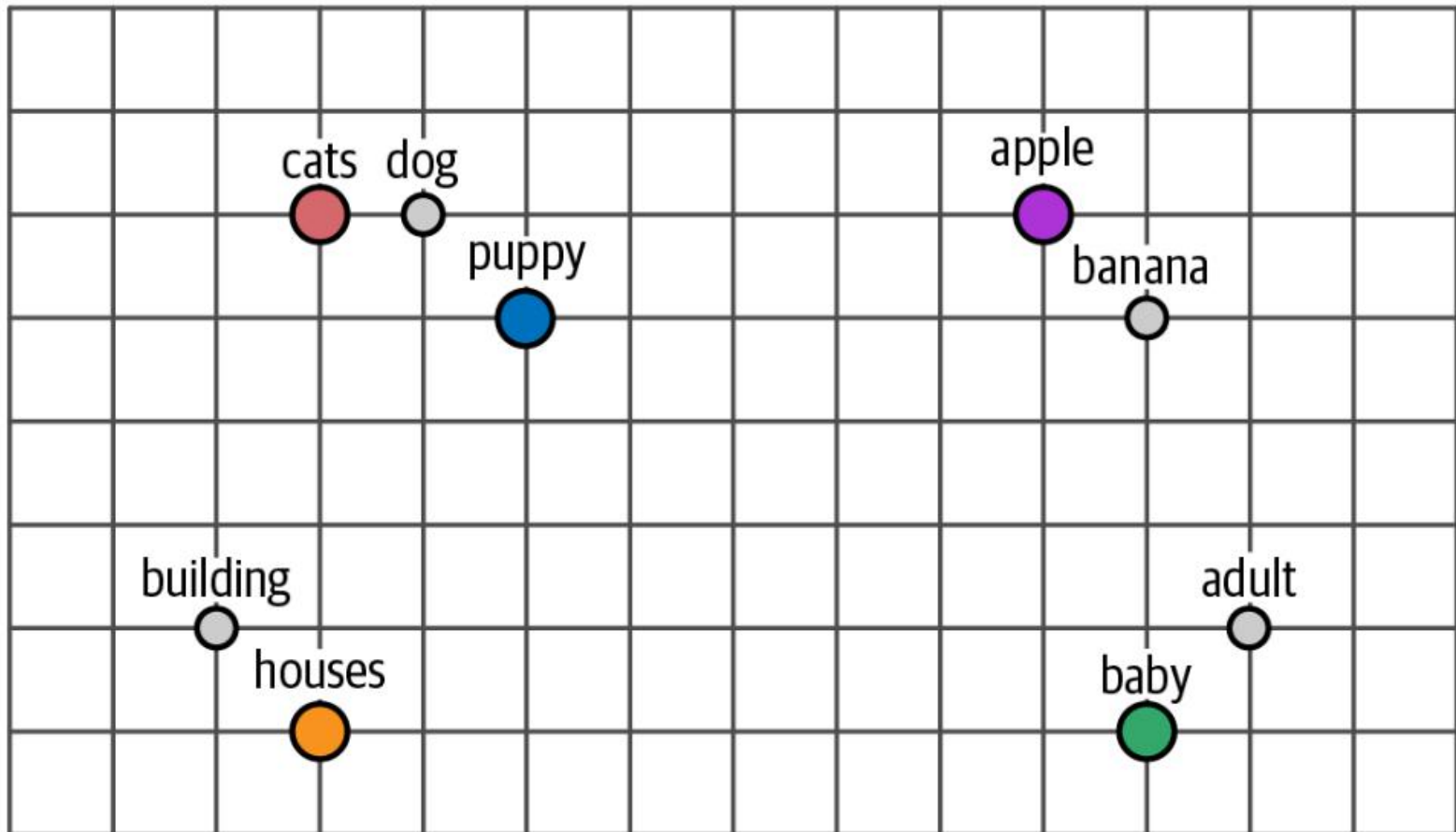
- Word2Vect: 1-hot encoding to vector



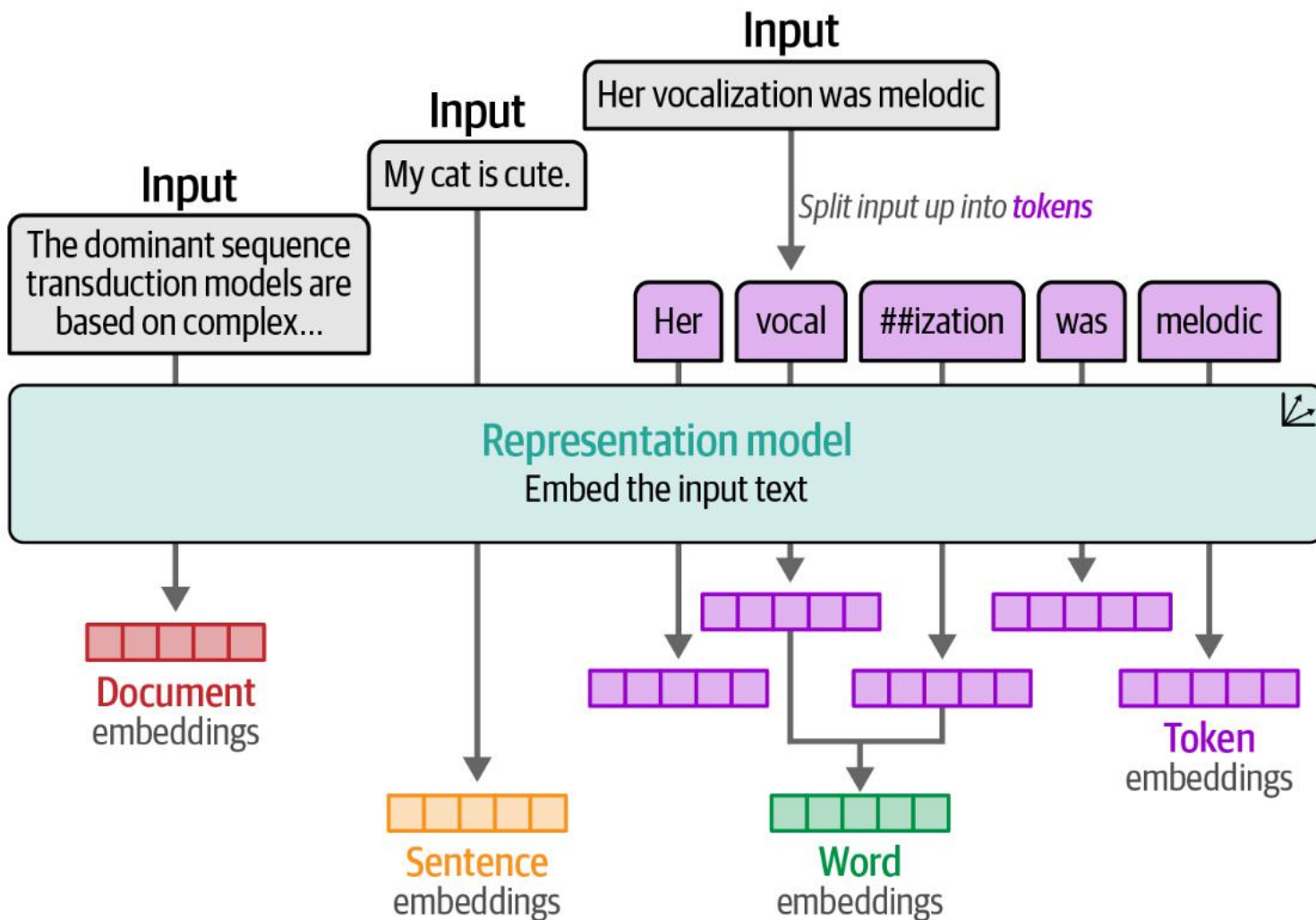
# What is supposed to learn



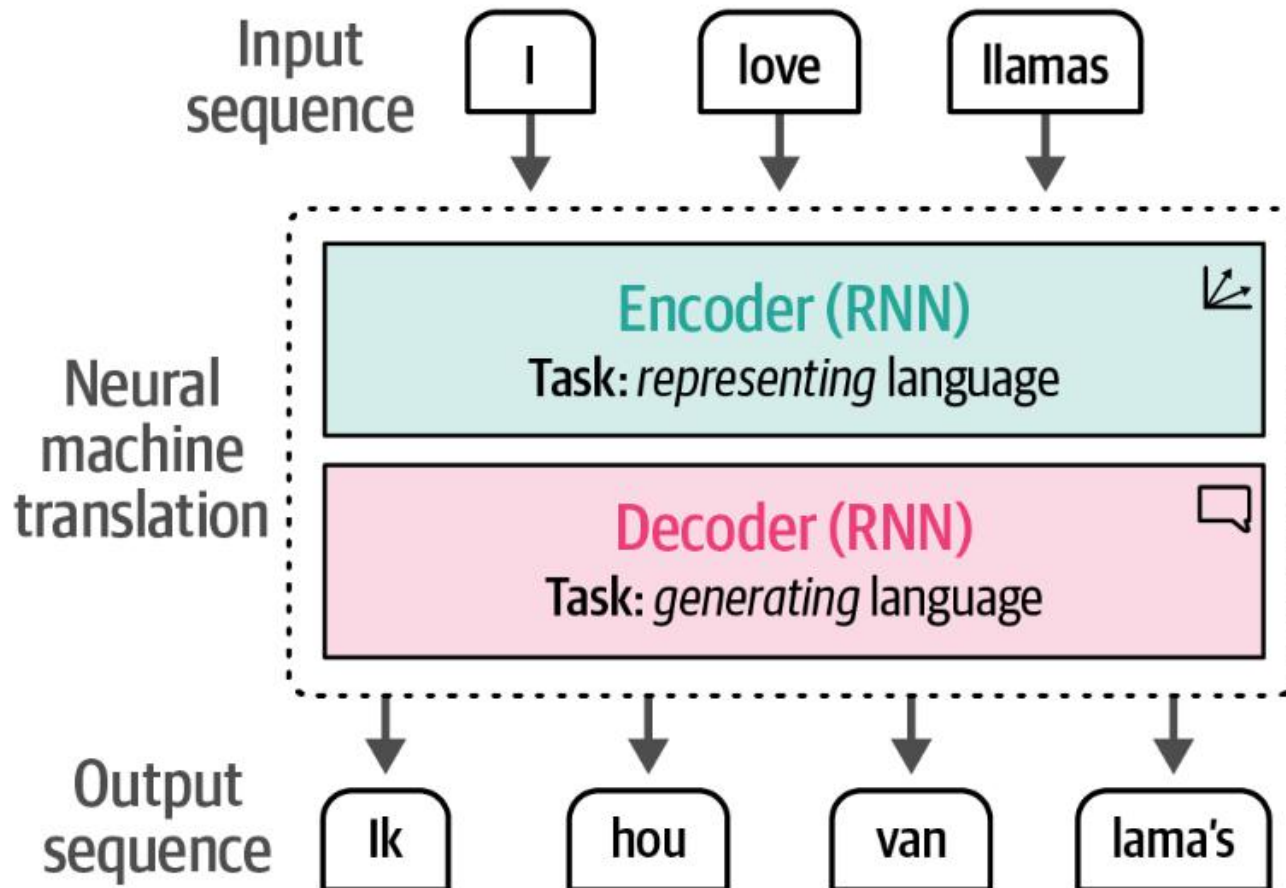
# Distances are preserved



# Types of embeddings

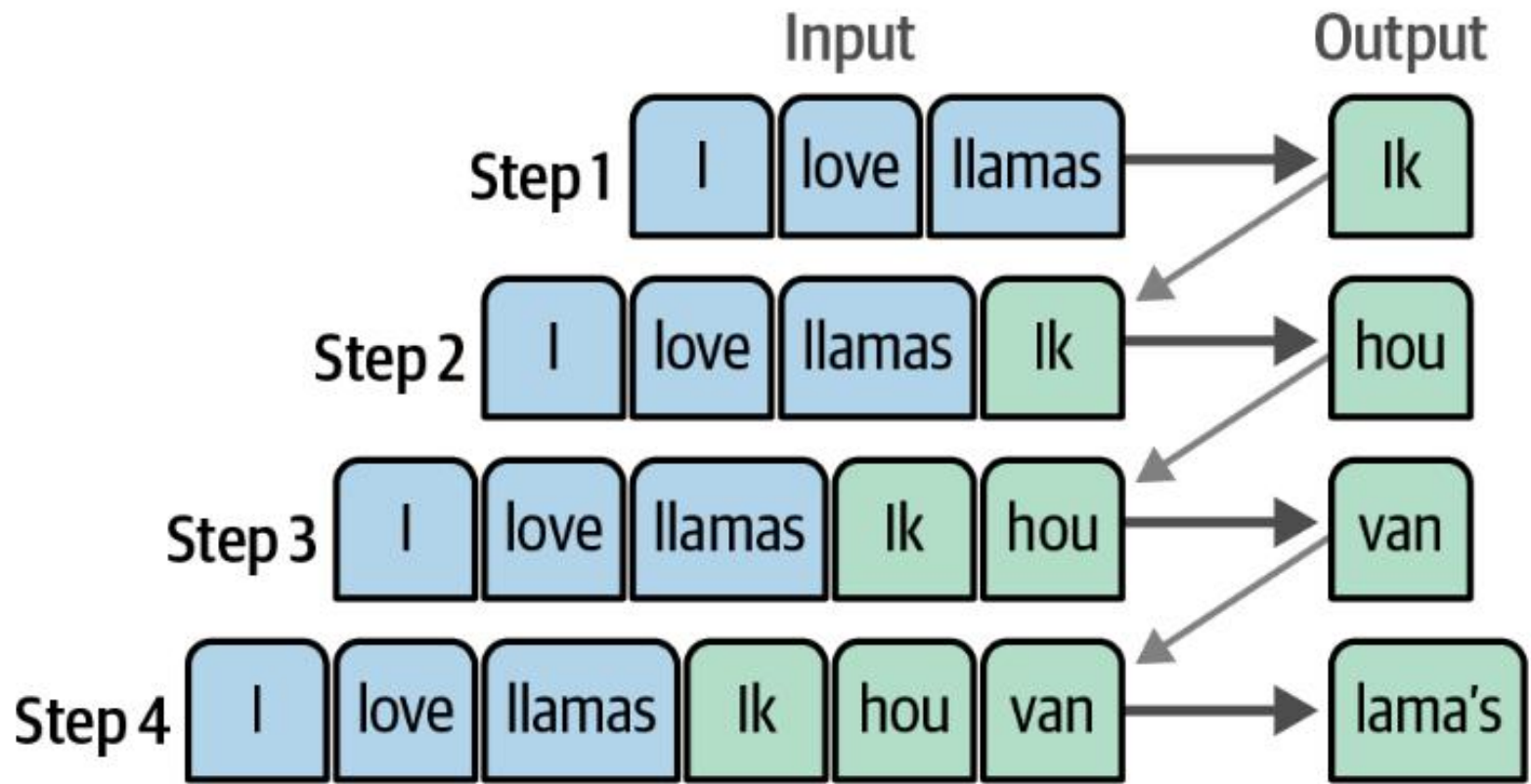


# Encoding and Decoding Context with Attention

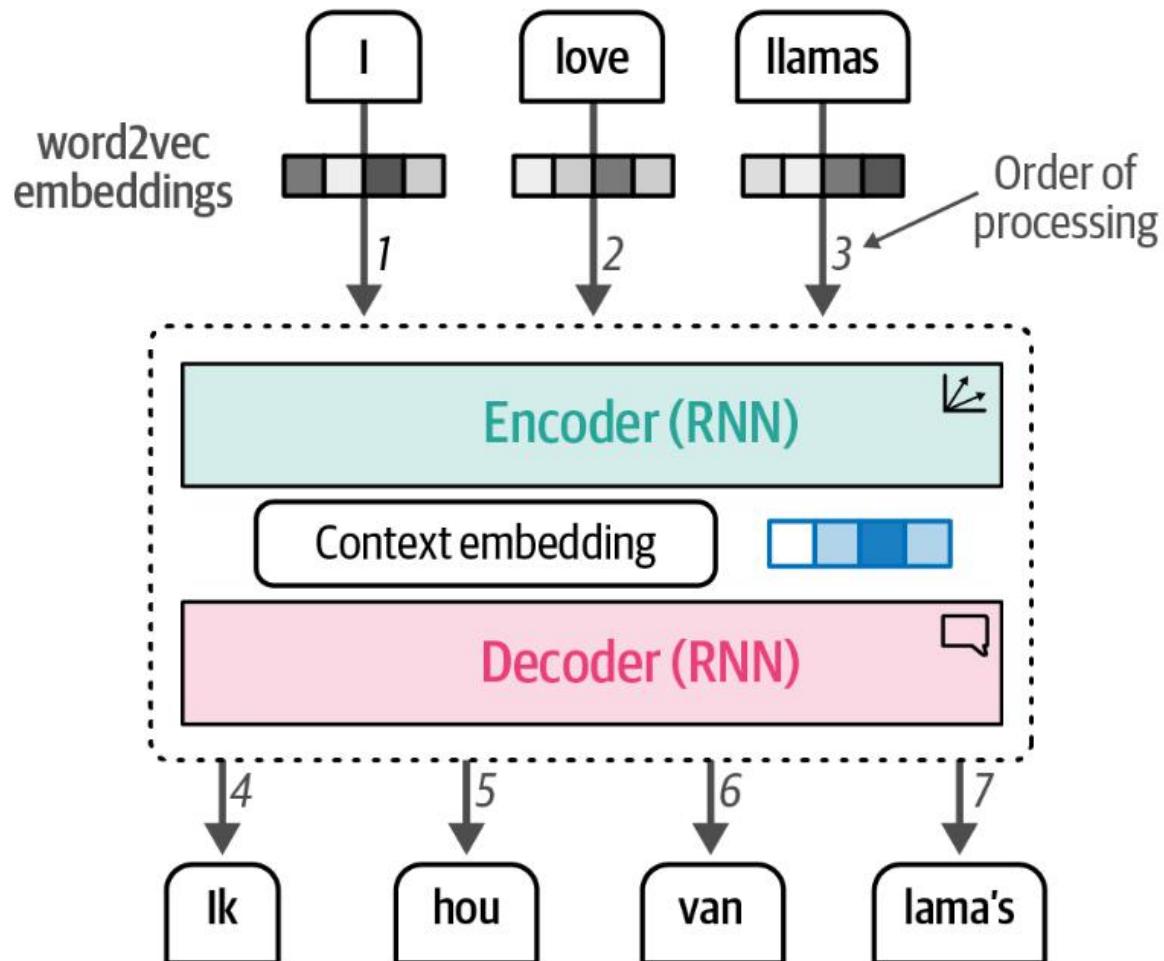




# Decoder in action

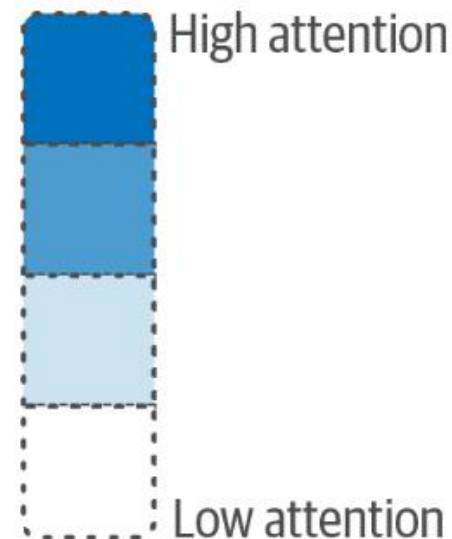
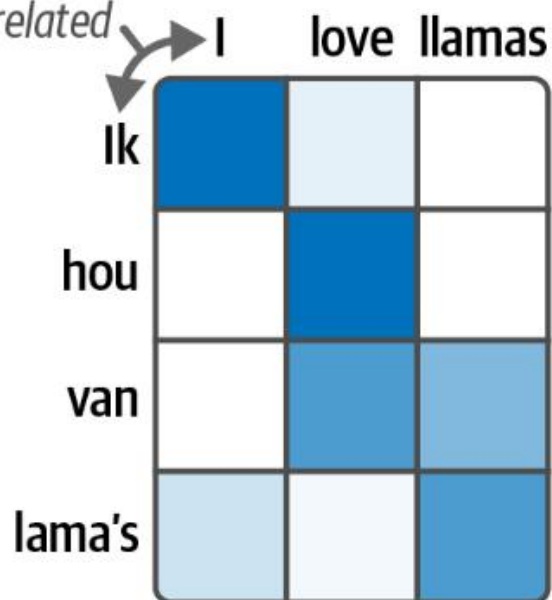


# Using word2vect

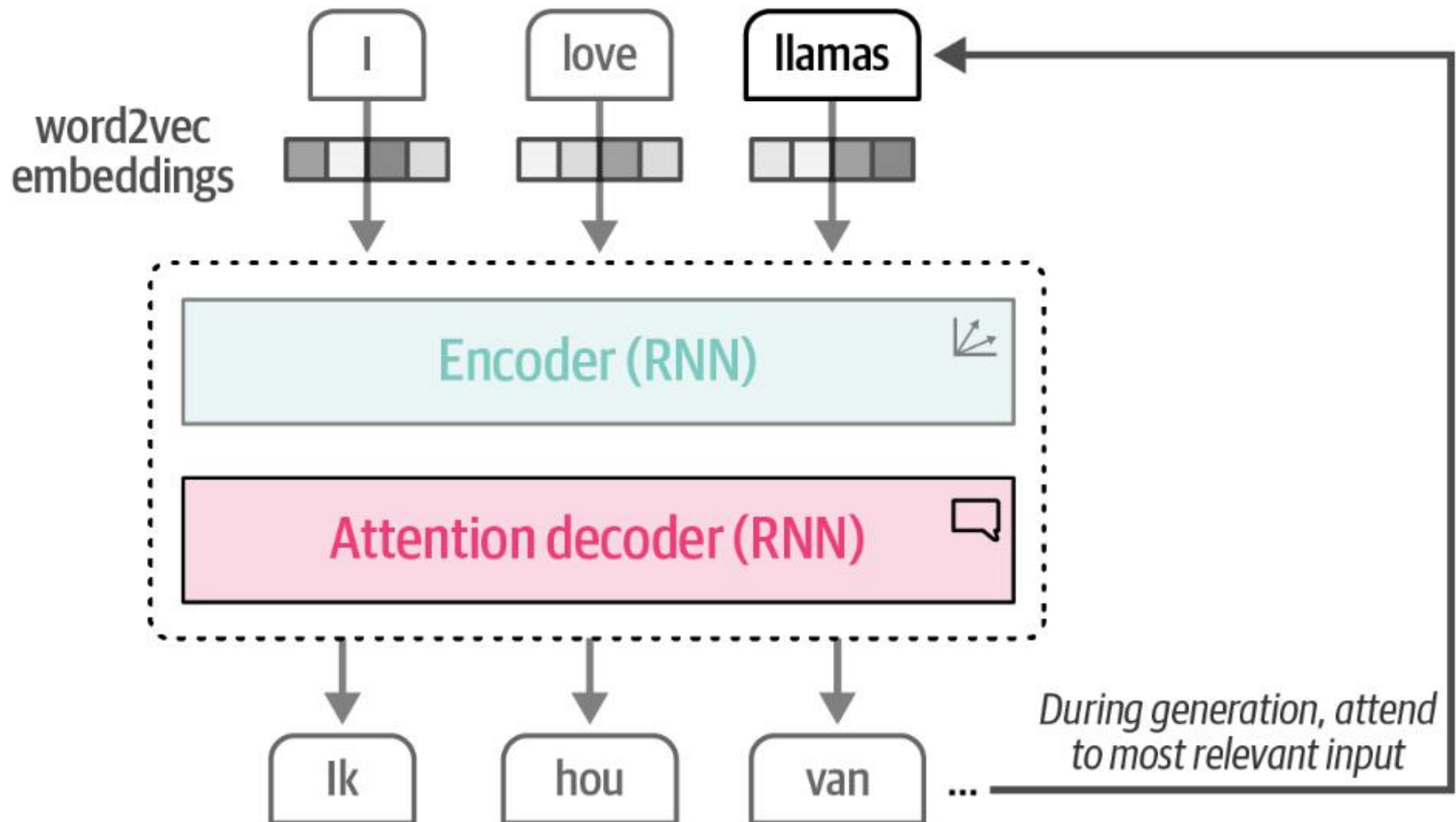


# Attention mechanism

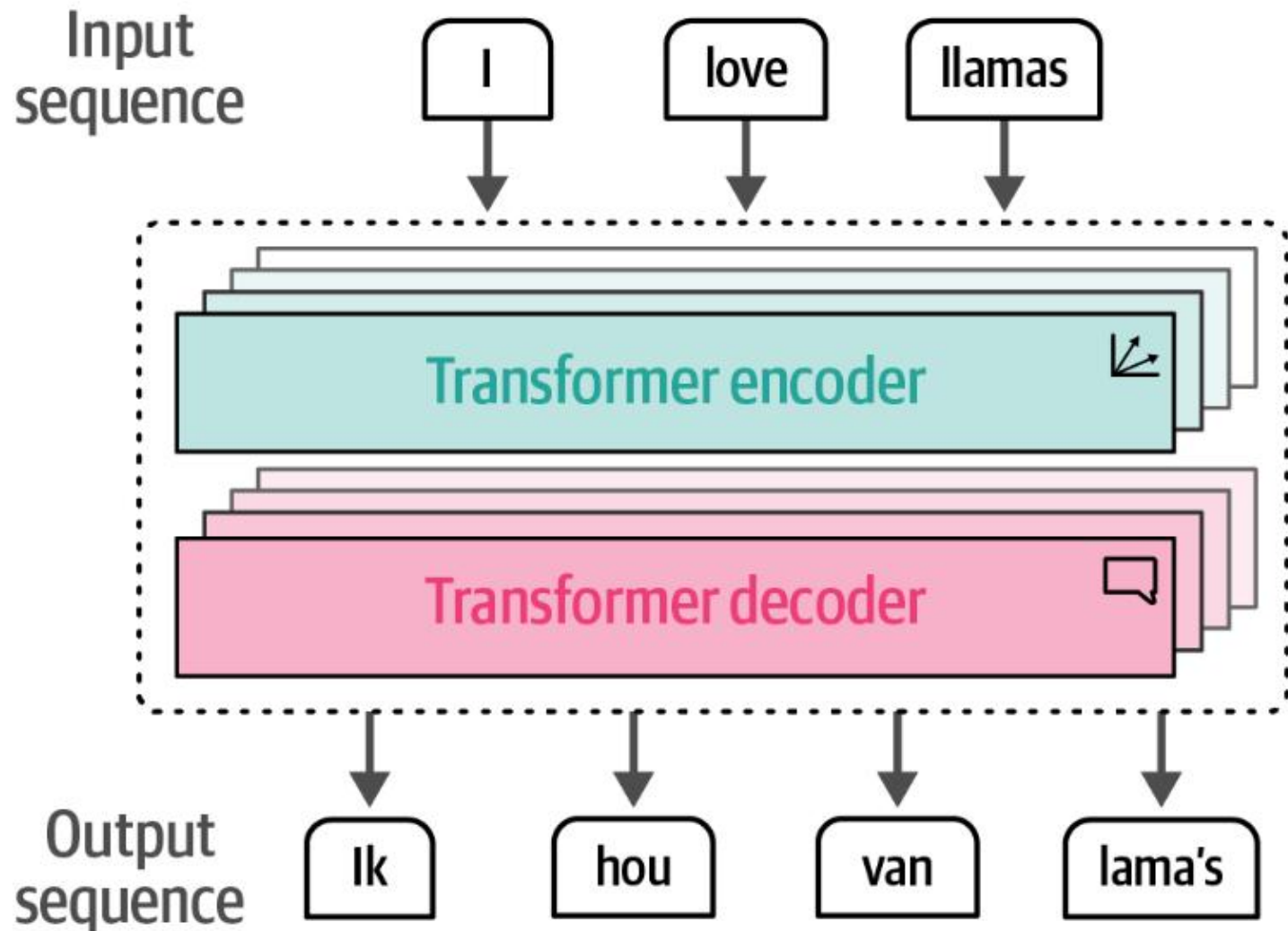
*Words with similar meaning  
have higher attention weights  
since they are highly related*



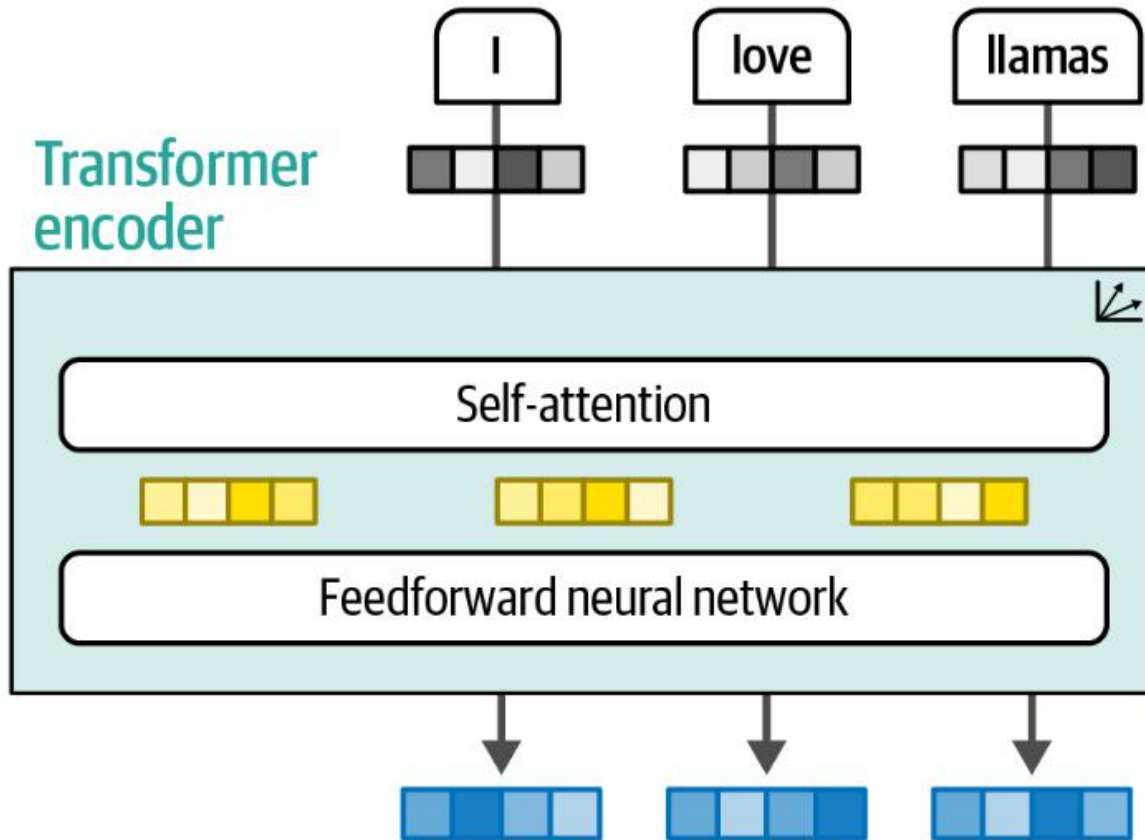
# Adding attention



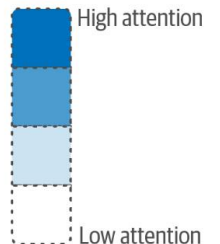
# Attention is all you need



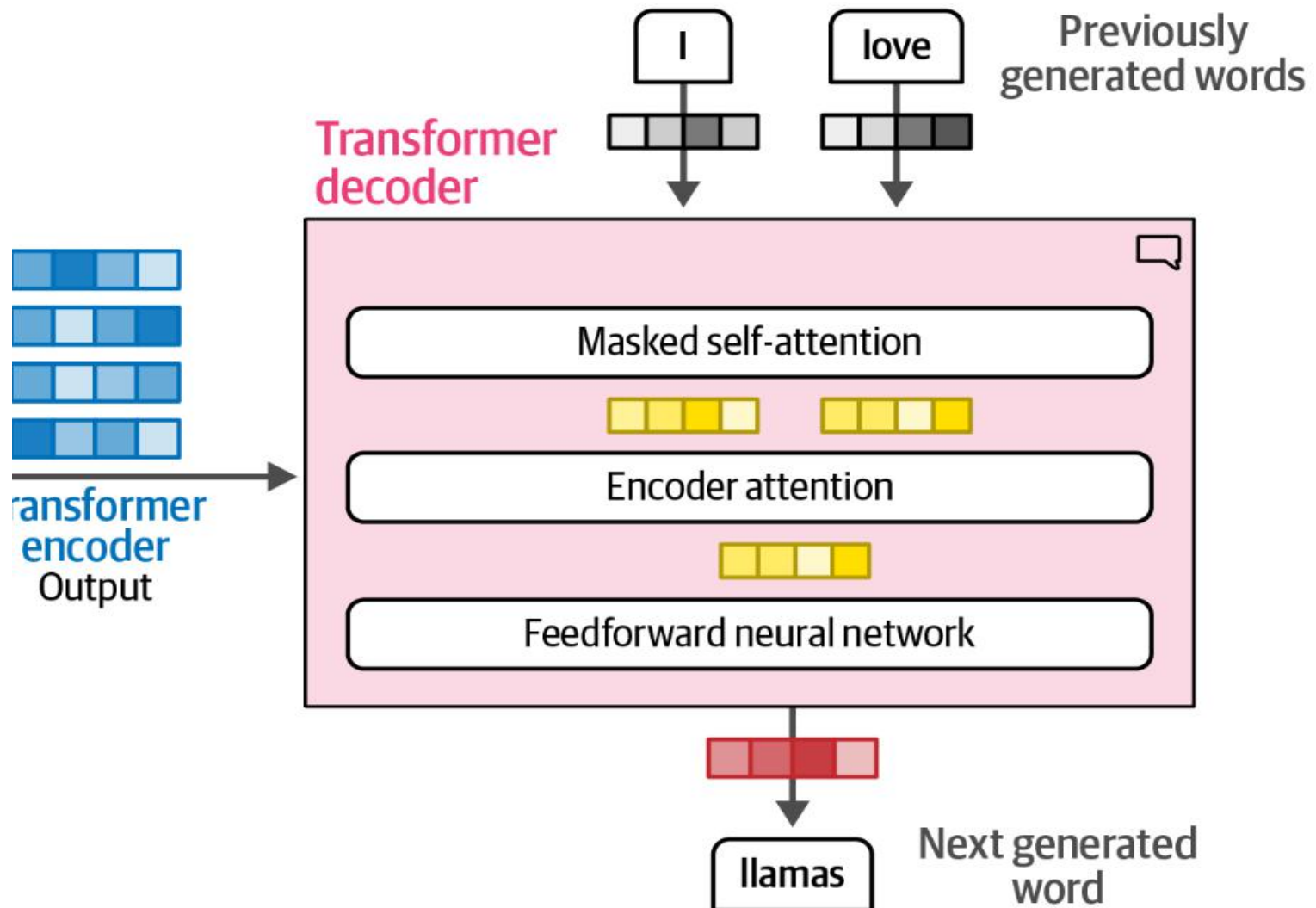
# Encoder with self-attention



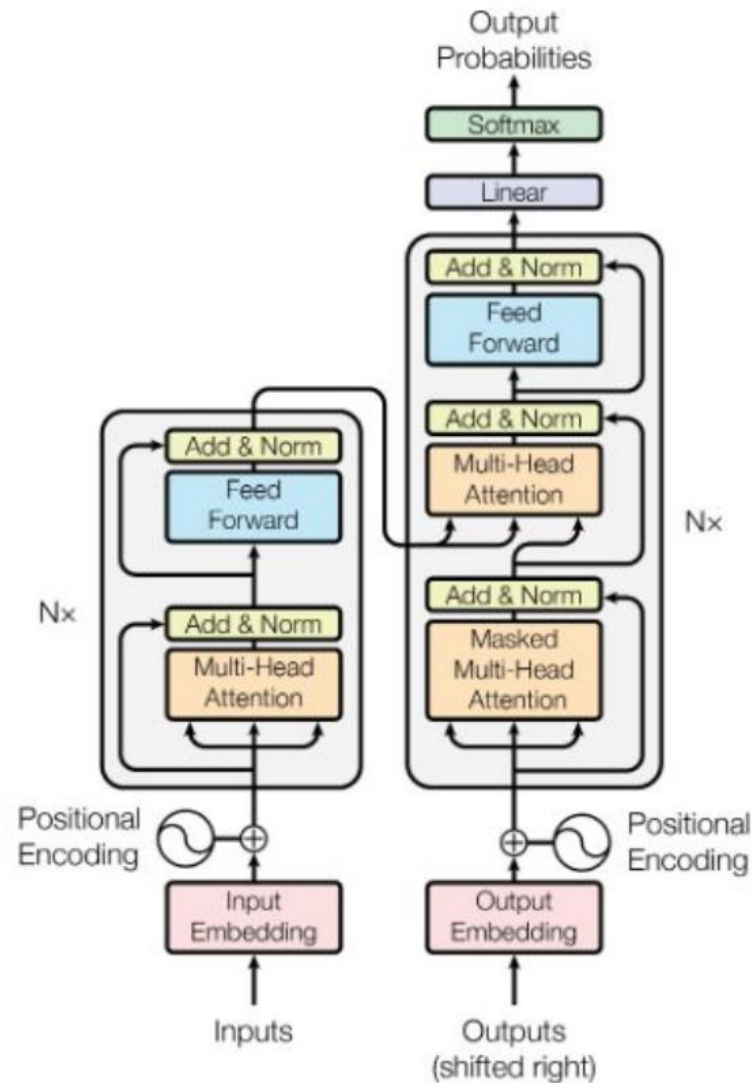
	I	love	llamas
I	High attention	Low attention	Low attention
love	Low attention	High attention	Low attention
llamas	Low attention	Low attention	High attention



# Decoder with attention

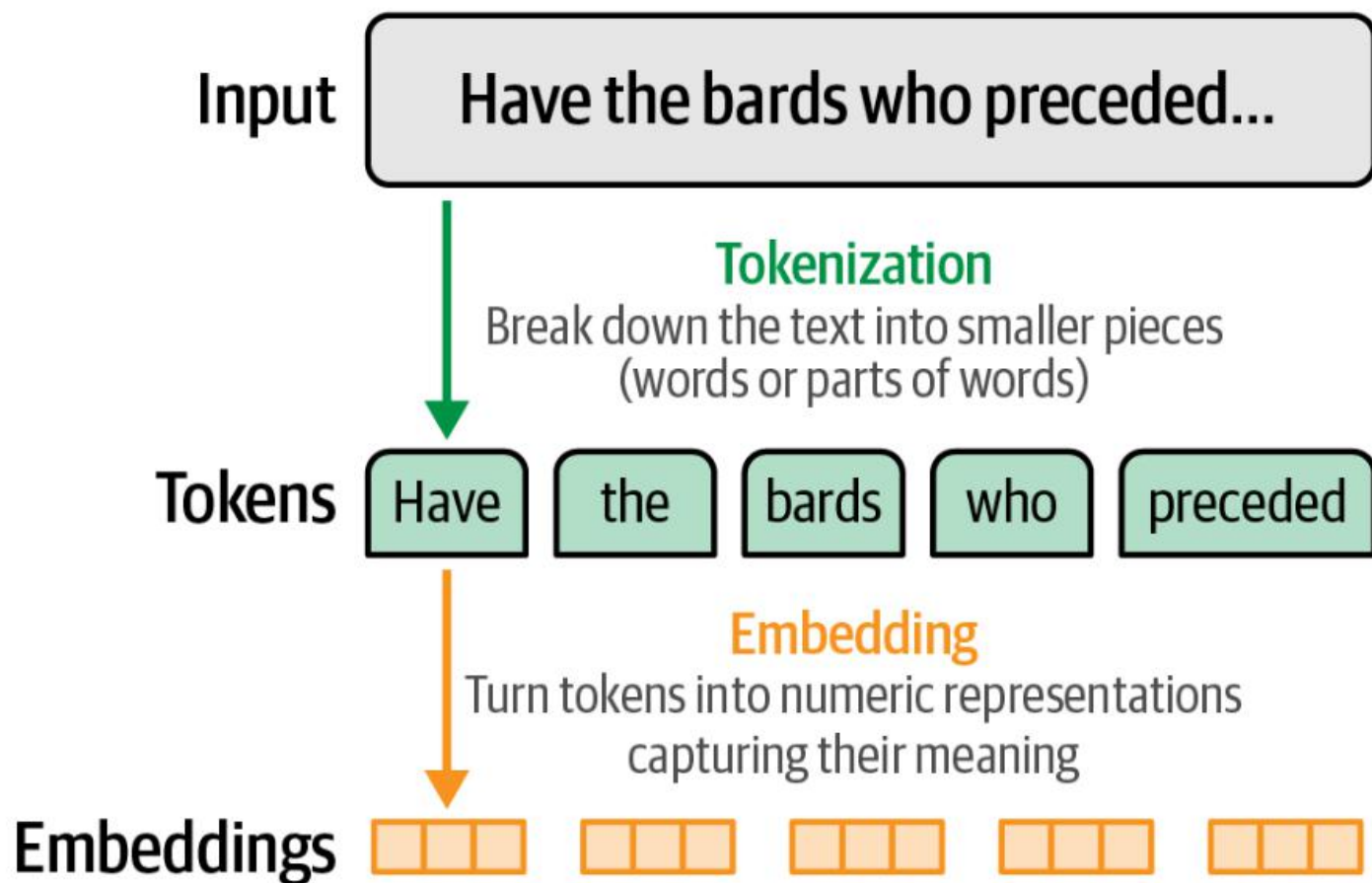


# The transformer architecture





# Tokenizers



# A real tokenizer

GPT-3.5 & GPT-4 GPT-3 (Legacy)

Have the bards who preceded me left any theme unsung?

Clear

Show example

Tokens

13

Characters

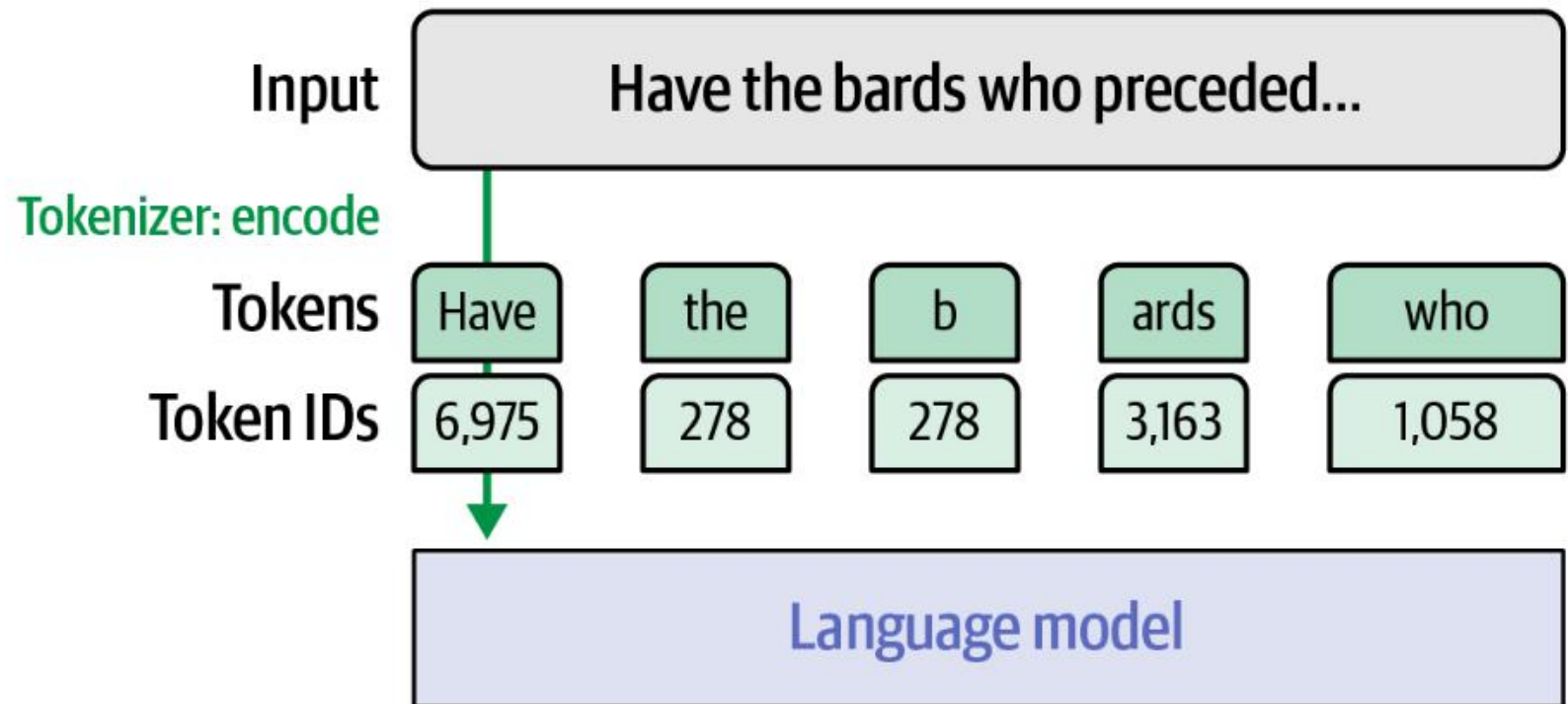
53

Have the bards who preceded me left any theme unsung?

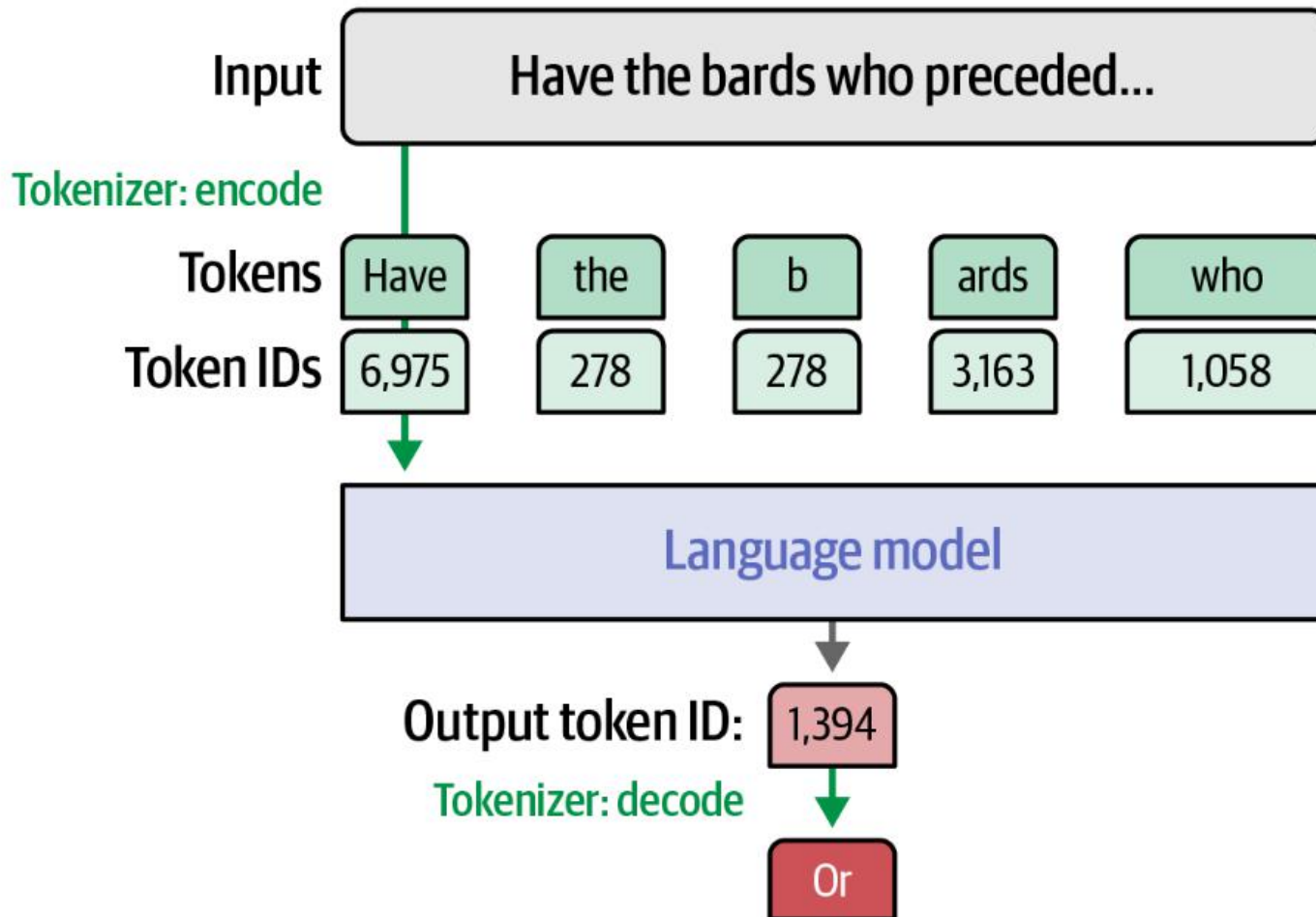
Text

Token IDs

# Each token has unique ID



# Reverse process



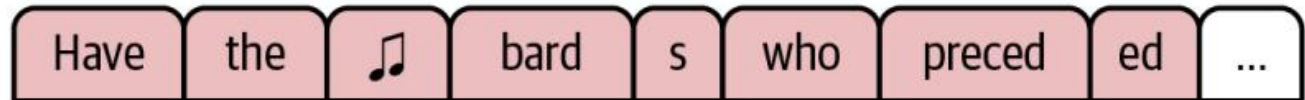
# Types of tokenizers

Text Have the 🎵 bards who preceded...

Word tokens



Subword tokens



Character tokens



# Some subword tokenizers

BERT base  
model  
(uncased)

[CLS] english and capital ##ization [UNK] [UNK] show \_ token ##s false none eli  
##f == > = else : two tab ##s : " " three tab ##s : " " 12 . 0 \* 50 = 600 [SEP]

BERT base  
model (cased)

[CLS] English and CA ##PI ##TA ##L ##I ##Z ##AT ##ION [UNK] [UNK] show \_ token  
##s F ##als ##e None el ##if == > = else : two ta ##bs : " " Three ta ##bs : " " 12 .  
0 \* 50 = 600 [SEP]

GPT-2

English and CAP ITAL IZ ATION  
\_ \_ \_ \_ \_  
show \_ tokens False None elif == > = else : two tabs : " " Three tabs : " "  
12 . 0 \* 50 = 600

FLAN-T5

English and CA PI TAL IZ ATION <unk> <unk> show \_ to ken s Fal s e None e l i f == > =  
= else : two tab s : " " Three tab s : " " 12 . 0 \* 50 = 600 </s>

GPT-4



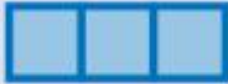
English and CAPITAL IZATION  
\_ \_ \_ \_ \_  
show \_ tokens False None elif == > = else : two tabs : " " Three tabs : " "  
12 . 0 \* 50 = 600

# Token Embeddings

## Trained tokenizer

Tokens	
Token ID	Token
0	!
1	"
...	...
50,257	

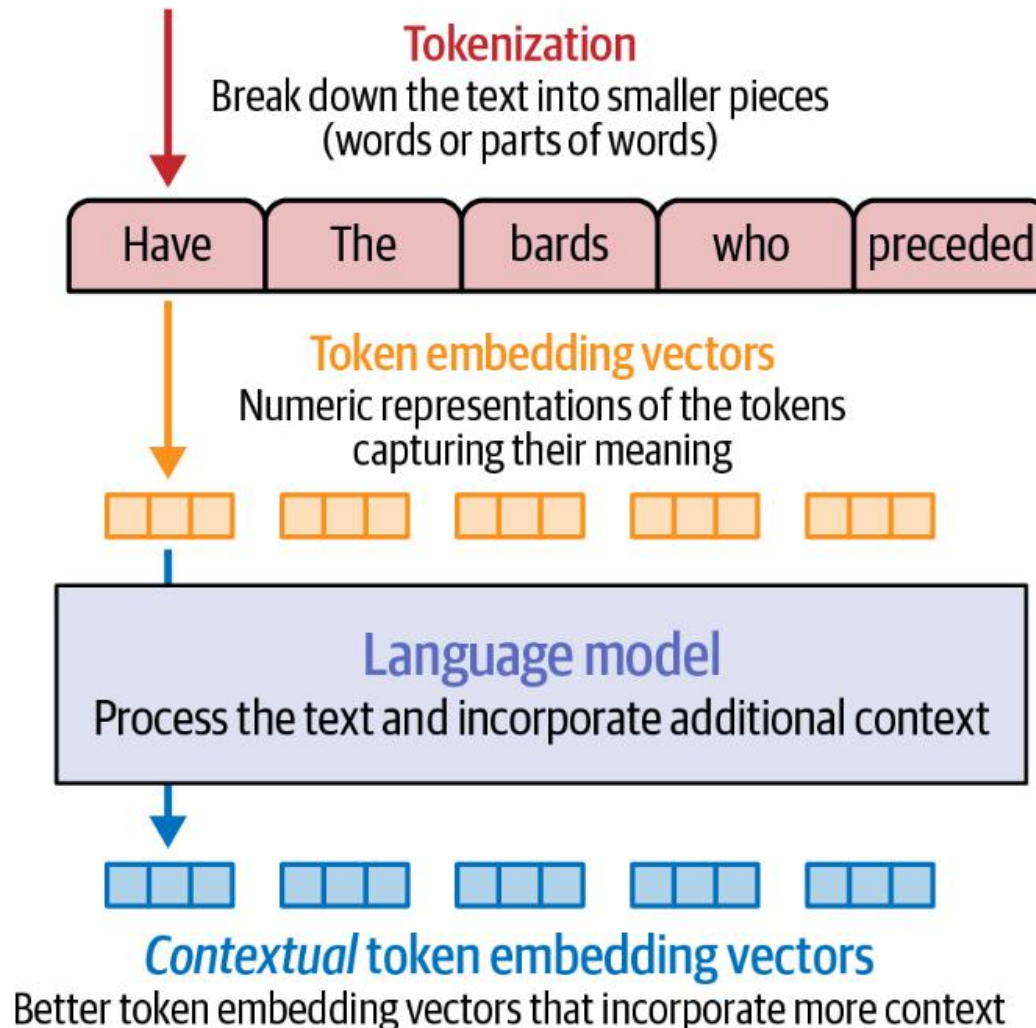
## Language model

Token embeddings	
0	
1	
...	...
50,257	



# Adding context

Have the bards who preceded me left any theme unsung?





# Text embedding

