# Text Generation and GPT

Hugh Dubberly and Shelley Evenson

**GPT (Generative Pre-trained Transformer)
is a large language model created by OpenAI.**

**GPT is currently on its third iteration, GPT-3, with a fourth on the way.**



**GPT-1** **GPT-2** **GPT-3** **GPT-4**

**The main innovation of GPT-3 is its massive size.**
**GPT-4 is projected to be nearly 600 times larger.**

**117 million parameters**
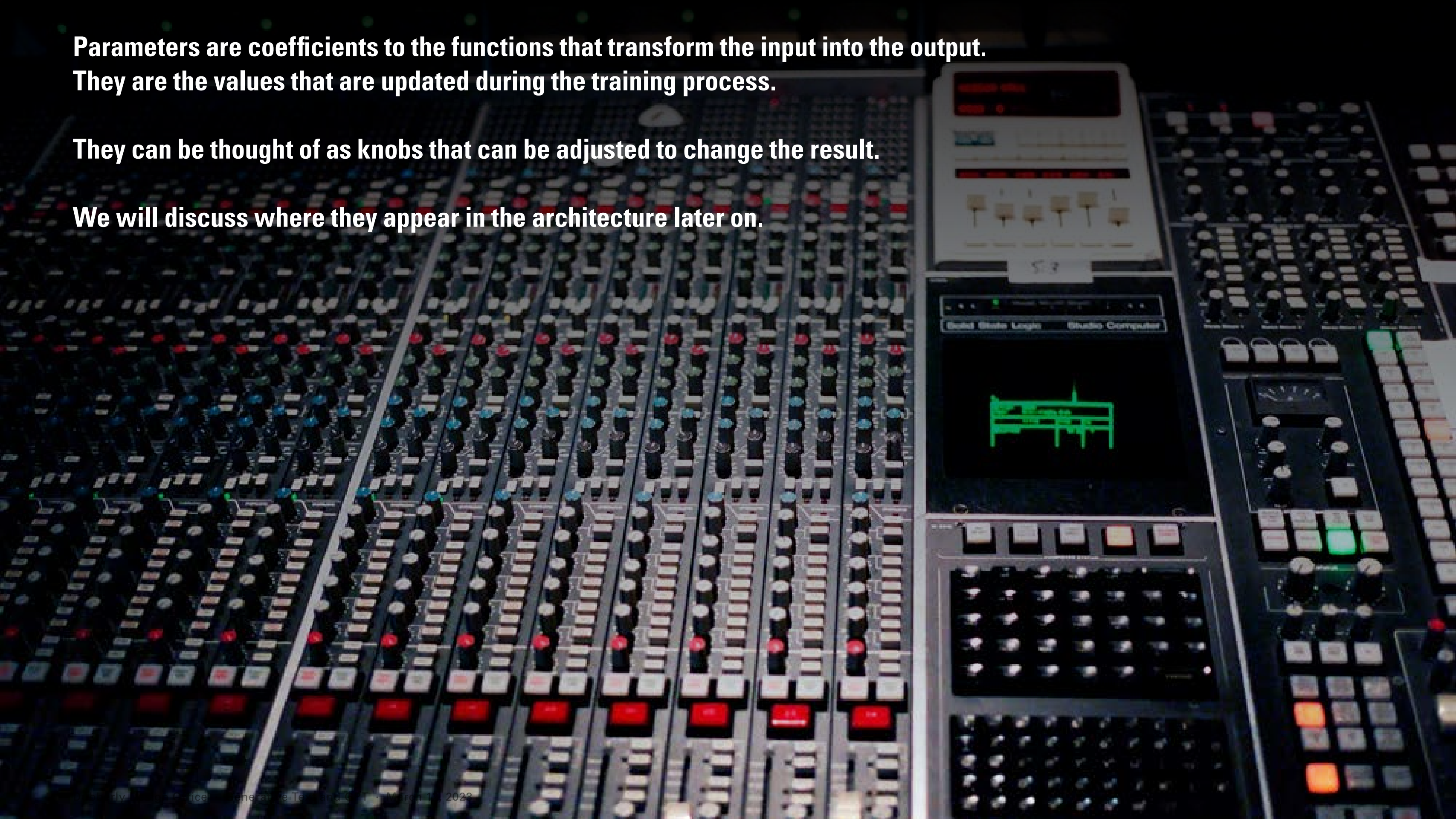
**1.5 billion parameters**

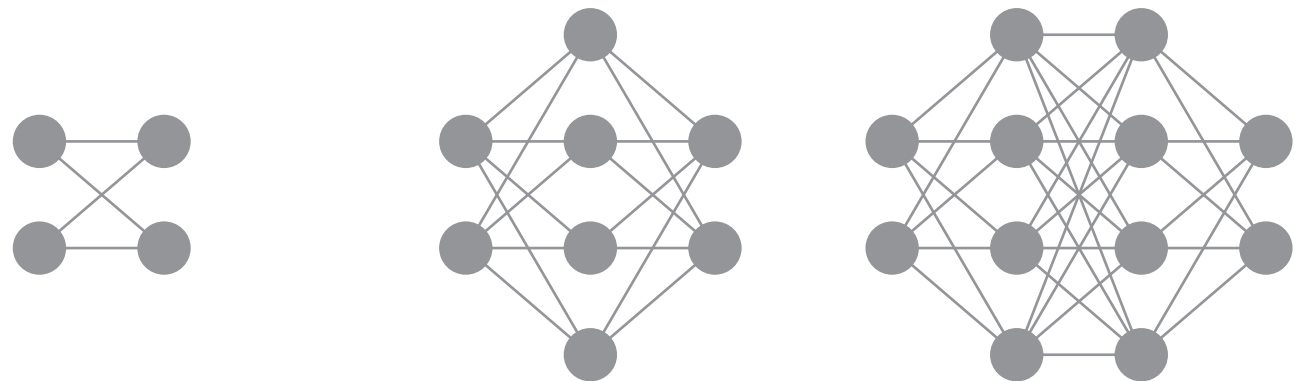**175 billion parameters**

**100 trillion+ parameters**

Parameters are coefficients to the functions that transform the input into the output.
They are the values that are updated during the training process.

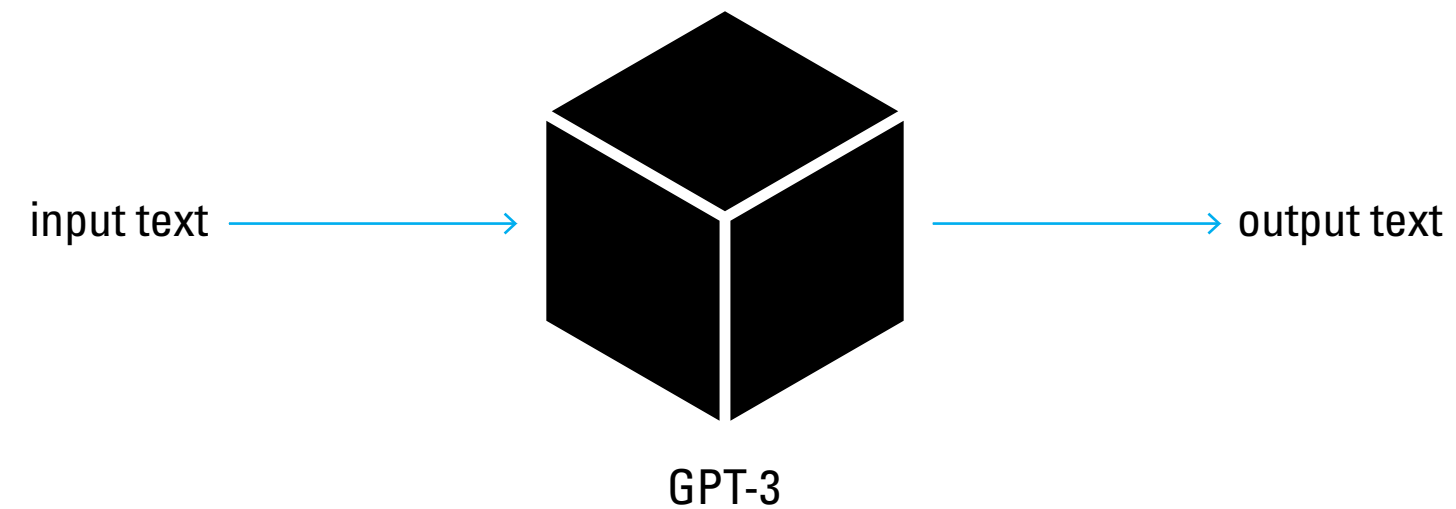They can be thought of as knobs that can be adjusted to change the result.

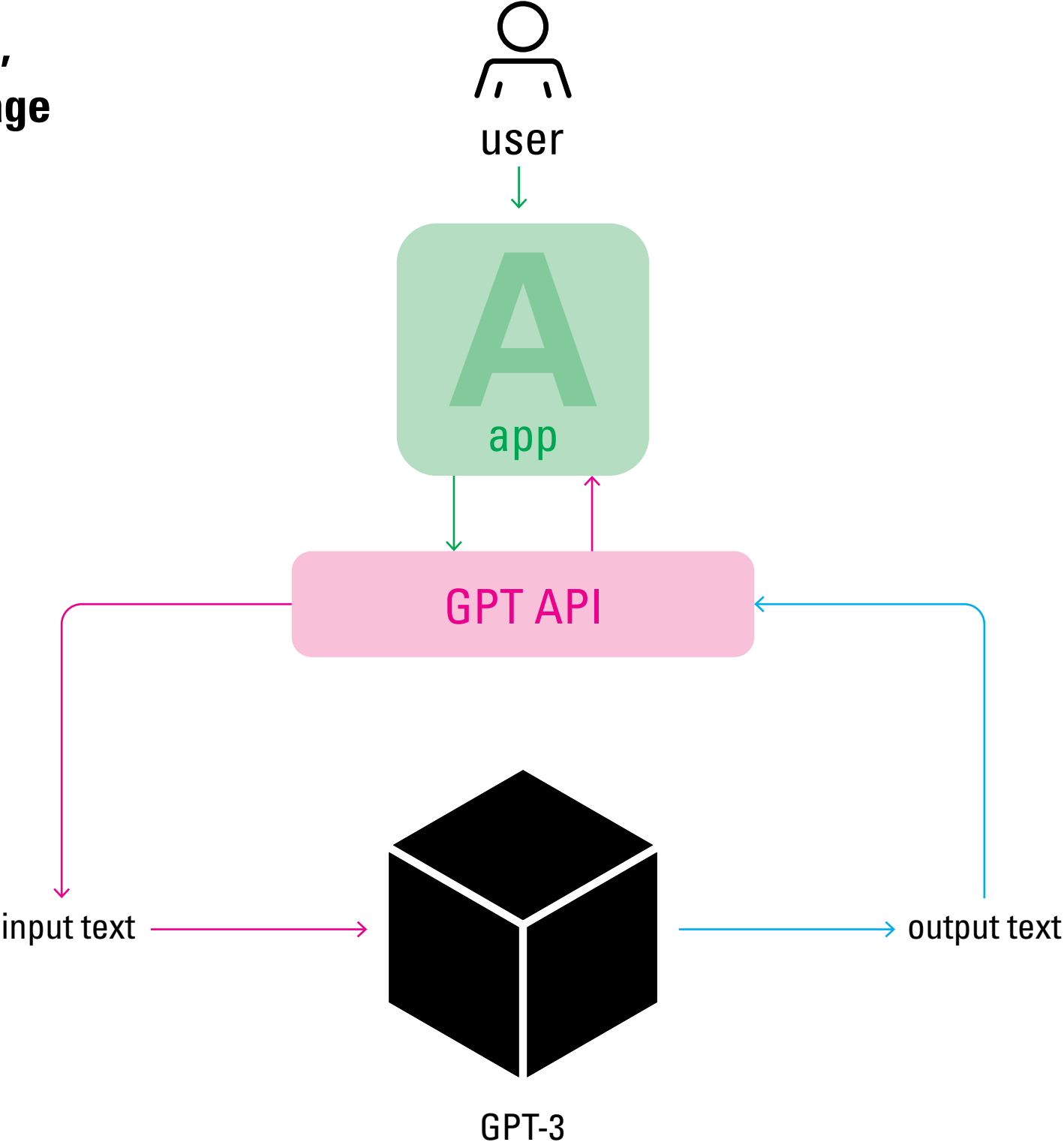We will discuss where they appear in the architecture later on.

# A model with more parameters can be thought of as a brain with more neurons.

# Let's first think of GPT-3 as a black box.

It takes text input and outputs its own text.

input text ——————→     ——————→ output text

GPT-3

**In addition to the model,
OpenAI has also released an API,
allowing for developers to leverage
the model's capabilities
in their own apps.**
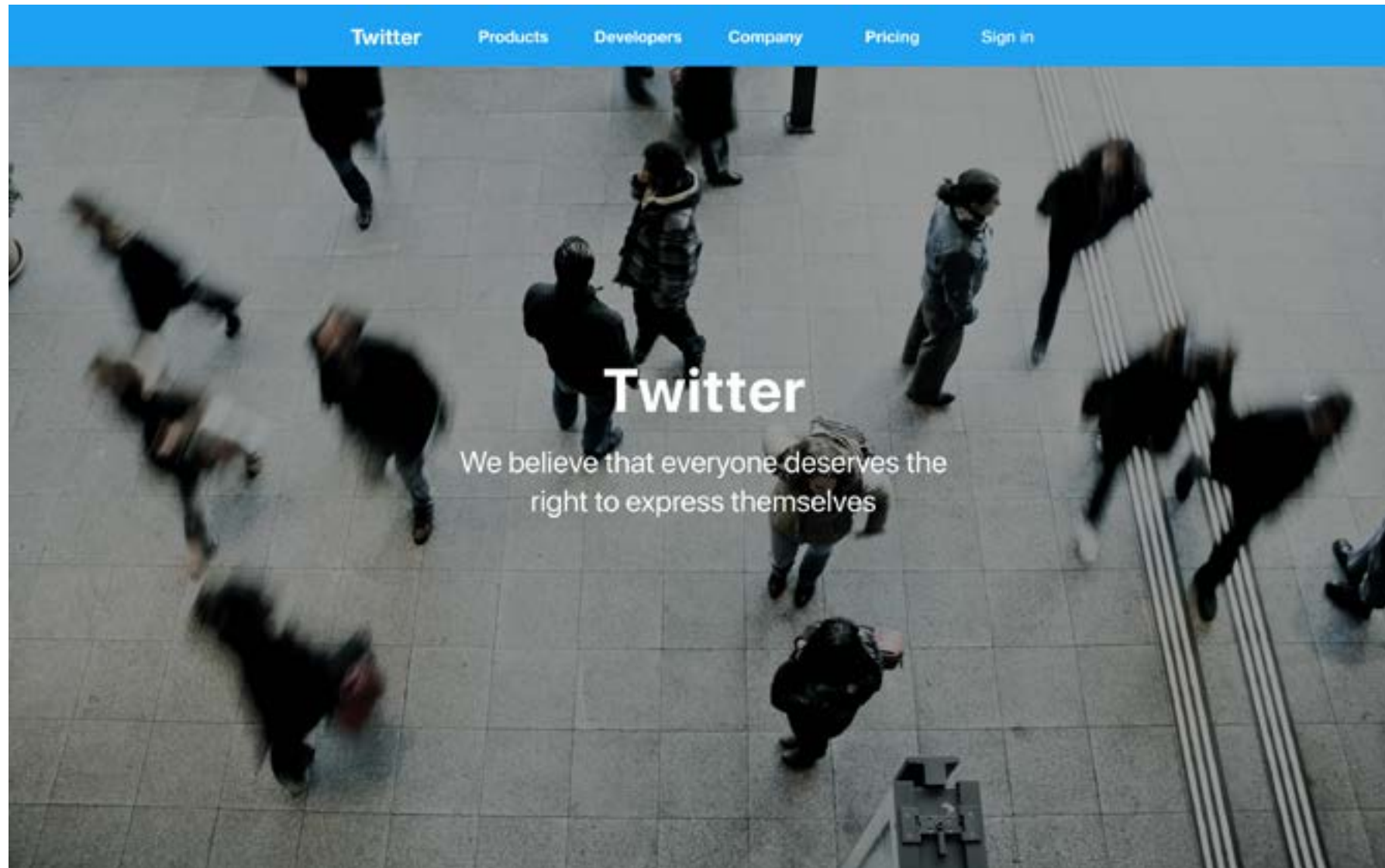
user

A

app

GPT API

input text

output text

GPT-3

**For example, developer Jordan Singer has created a plug-in for Figma that can generate app or website prototype designs based on text input.**
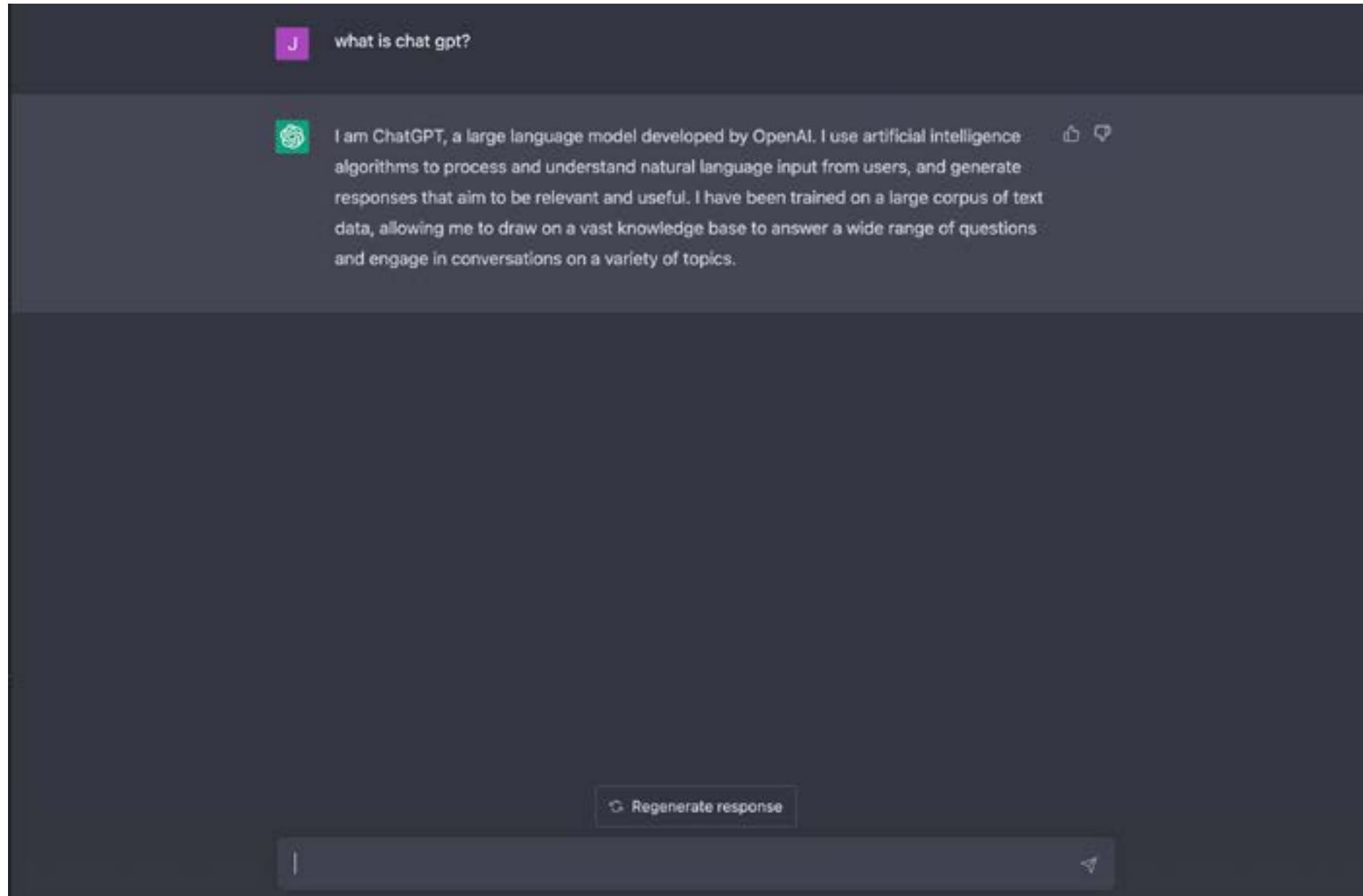
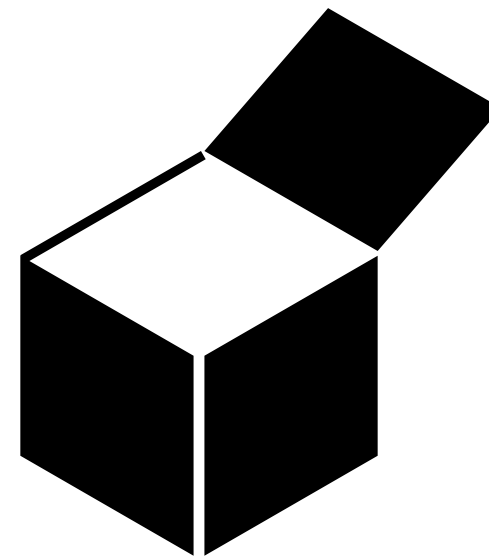**The output for:**
"a website like apple.com for Twitter "

**OpenAI has their own demo, called ChatGPT,**
**a version of GPT-3 optimized to be a chatbot.**

**It serves as a front end for GPT-3,**
**so non-coders can experiment with it as well.**

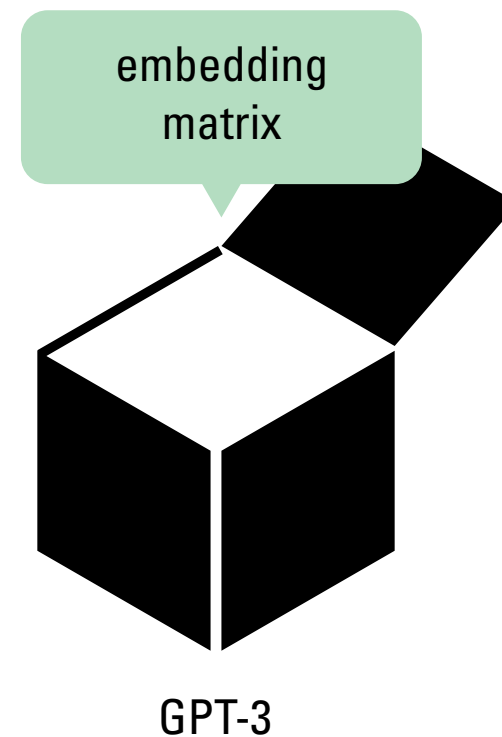**To understand what its doing, we will open up the black box.**



GPT-3

**First, the bulk of the system does not operate on words,
but rather vector embeddings of the words.
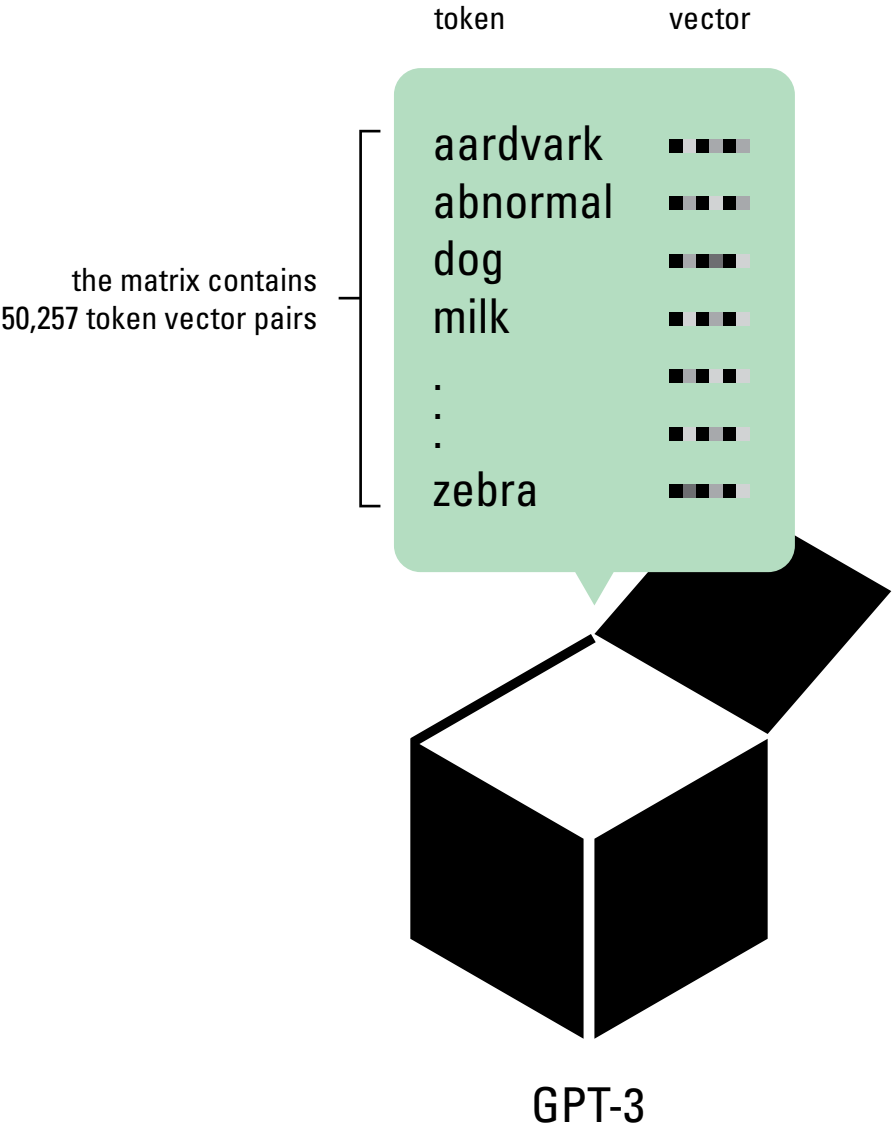(These vectors are of 2048 dimensions)**

**Inputs are tokenized,
meaning broken up into smaller parts.
(for simplicity's sake we will consider tokens words)**
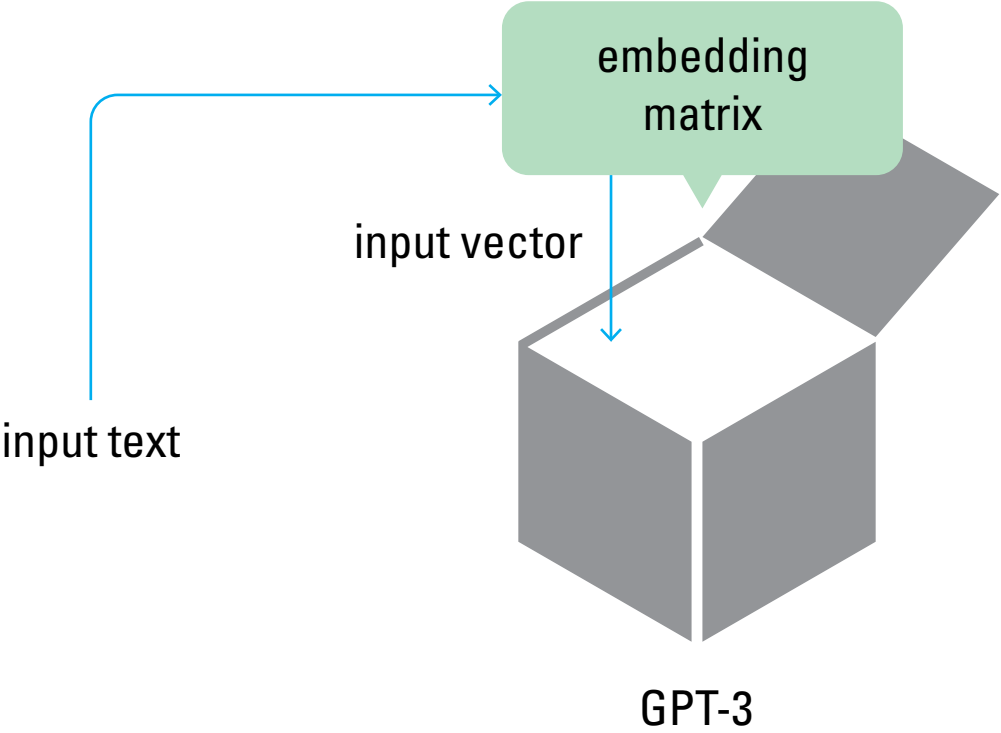
**Tokens are embedded using an embedding matrix.**

embedding
matrix

GPT-3

**The embedding matrix is a list of all the tokens in the model's vocabulary.**

**Each row of the embedding matrix is a token
in the system's vocabulary,
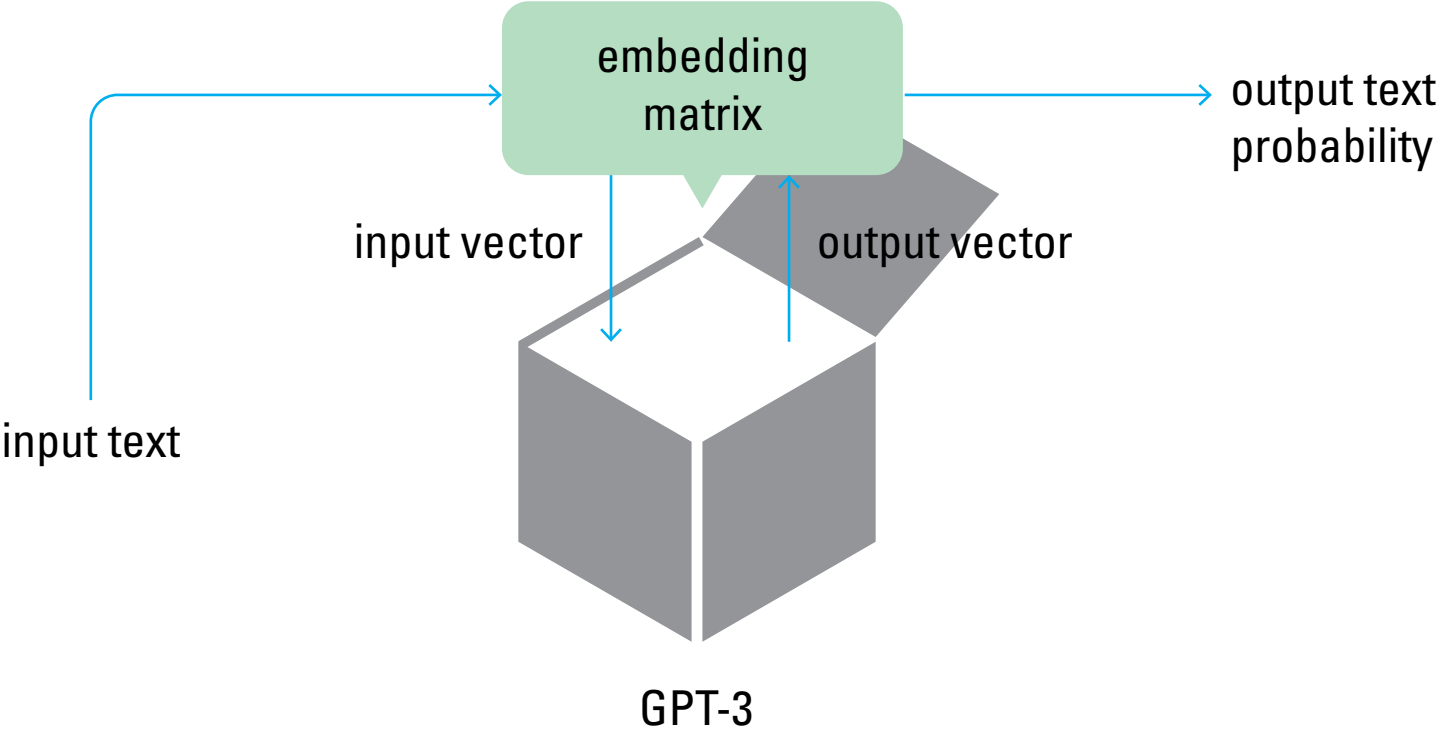and its corresponding vector.**

token       vector

aardvark

abnormal

dog

milk

.

.

zebra

the matrix contains
50,257 token vector pairs

GPT-3

**When text is inputted,
GPT looks up the input text's embedding
in the matrix.**



embedding
matrix

input vector

input text

GPT-3

**It then operates on that vector,
and outputs a new vector.**



embedding
matrix

input vector          output vector

input text

GPT-3

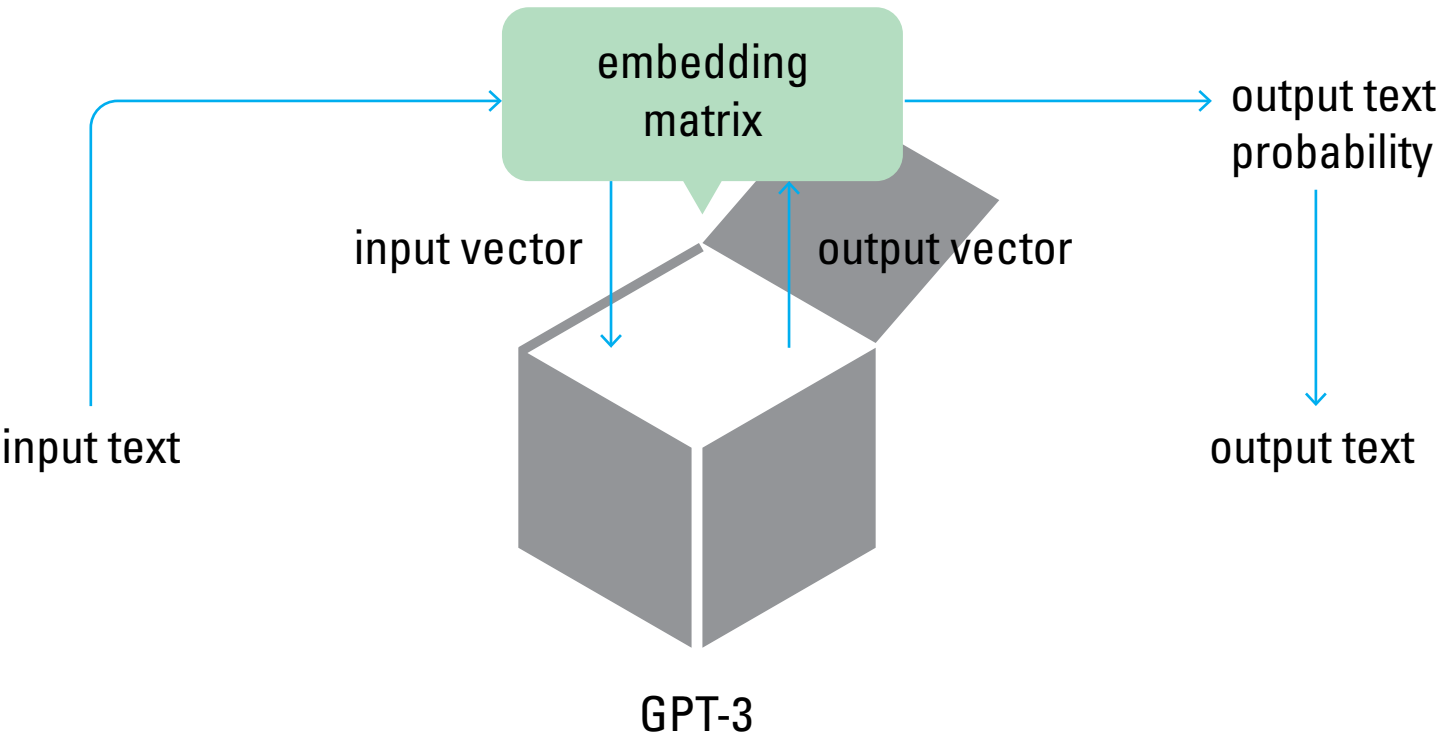**The output vector is multiplied by the embedding matrix,
resulting in a list of probabilities for every word in the vocabulary.**



embedding
matrix

input vector

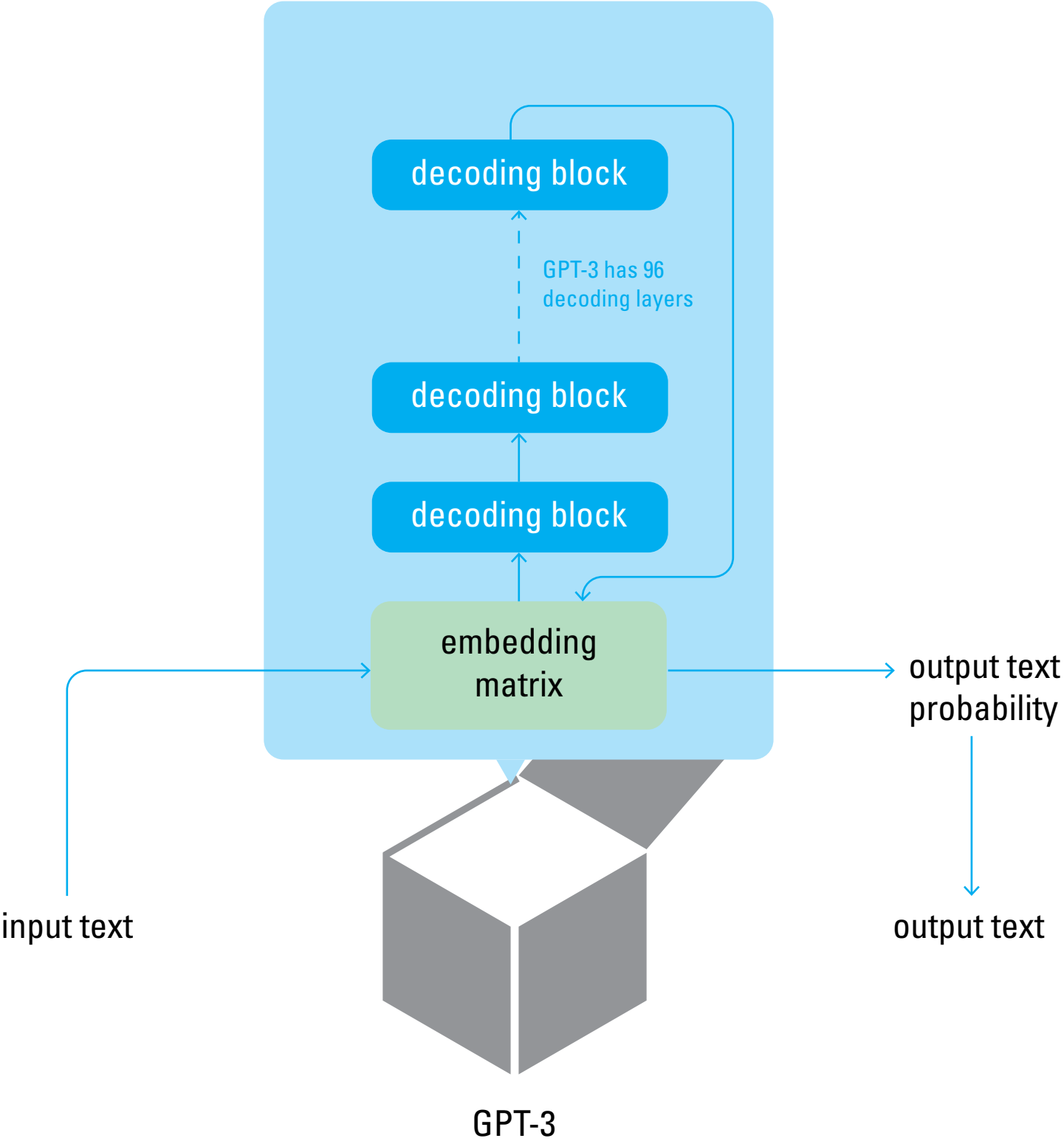output vector

input text

output text
probability

GPT-3

# The output text is generated based on those probabilities.

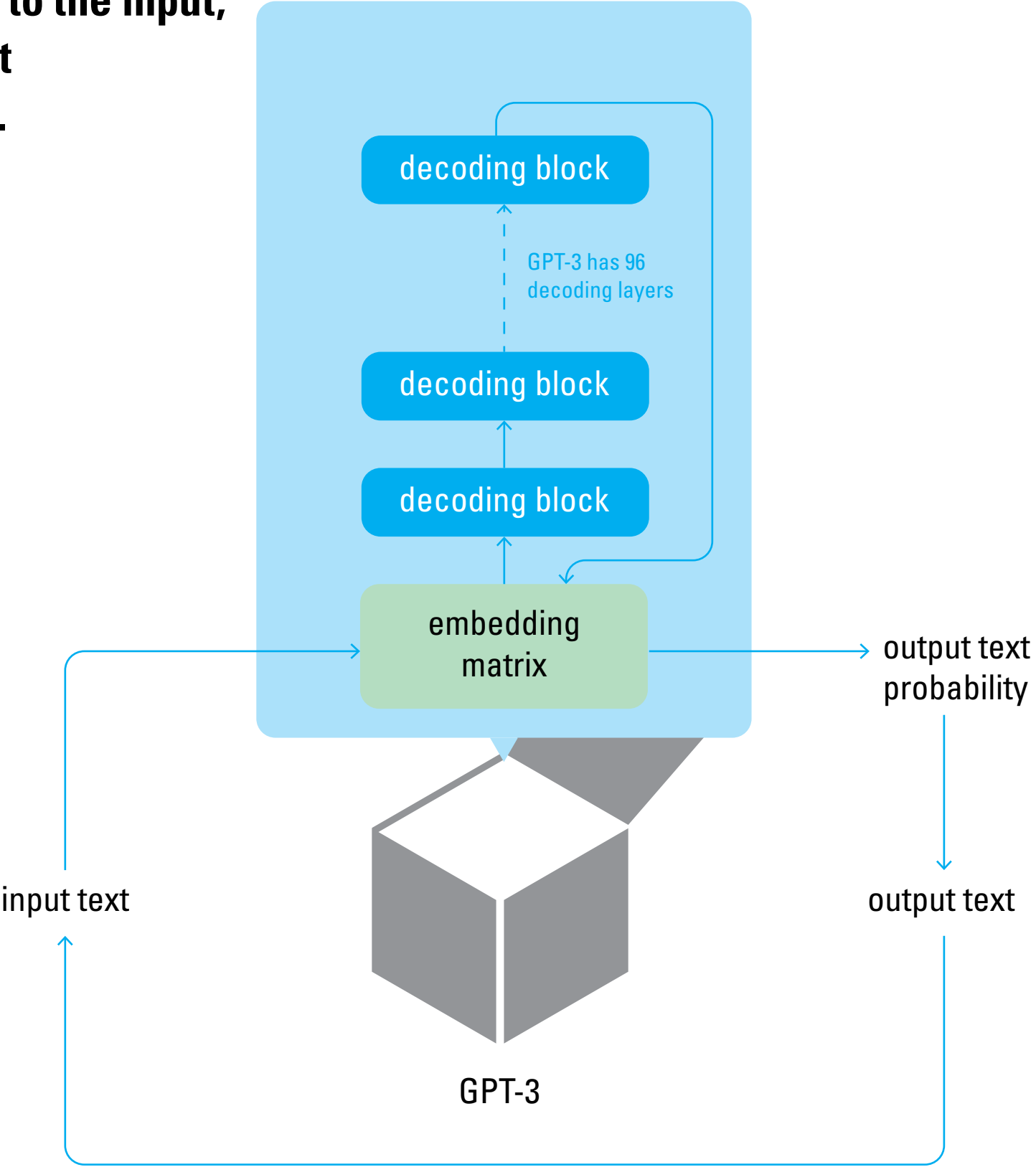**The output vector is generated by a stack of decoding blocks.**

**Each one receives the output of the previous one, transforms it, and passes it on to the next.**

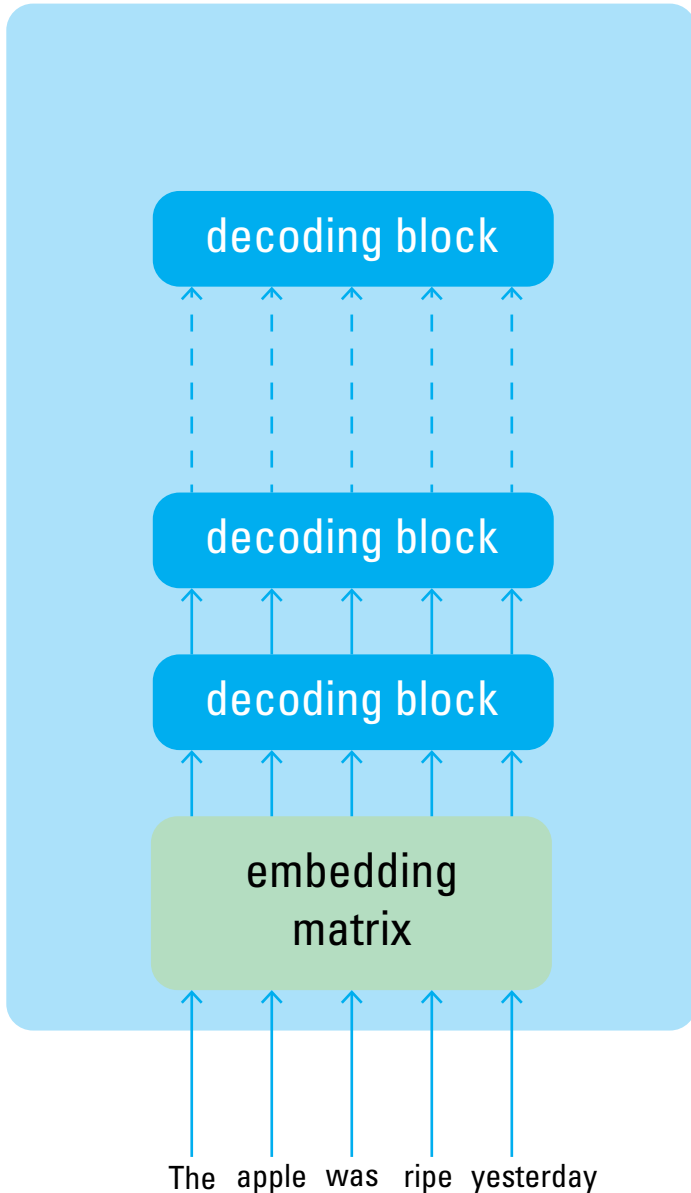**All of the 96 decoding blocks are identical.**

decoding block

GPT-3 has 96 decoding layers

decoding block

decoding block

embedding matrix

input text

output text probability

output text

GPT-3

**The output text is then appended to the input,
and then passed as the new input
for the next iteration of the cycle.**

**This continues until
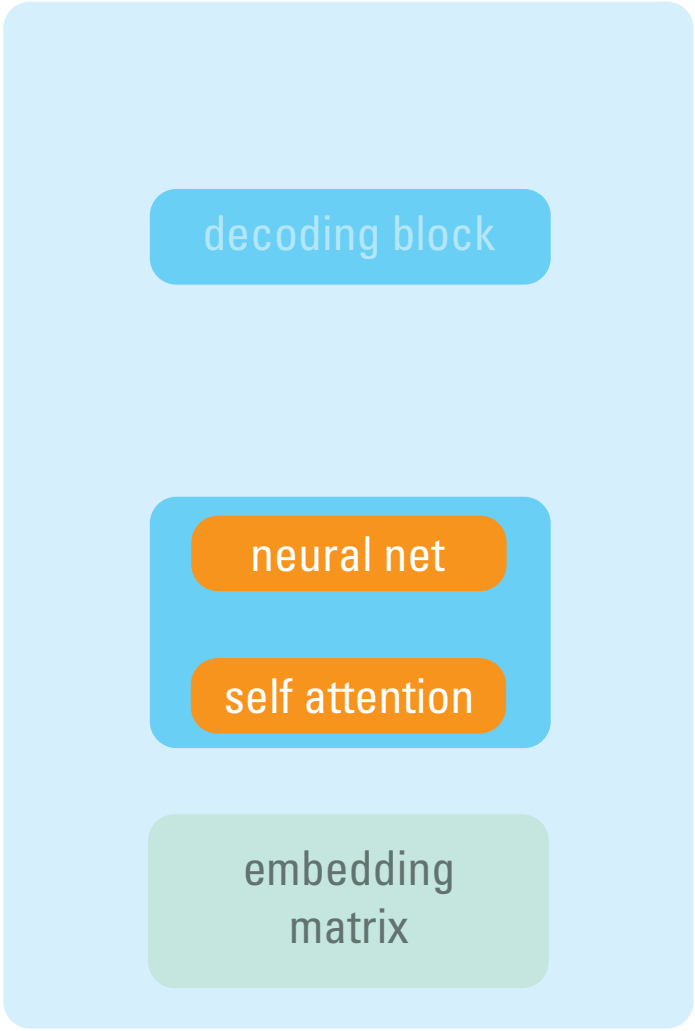the most probably output
is the end token,
and GPT-3 stops generating text.**



decoding block

GPT-3 has 96
decoding layers

decoding block

decoding block

embedding
matrix

output text
probability

input text

output text

GPT-3

# Each token of the input is processed on its own path.



decoding block

decoding block

decoding block

embedding
matrix
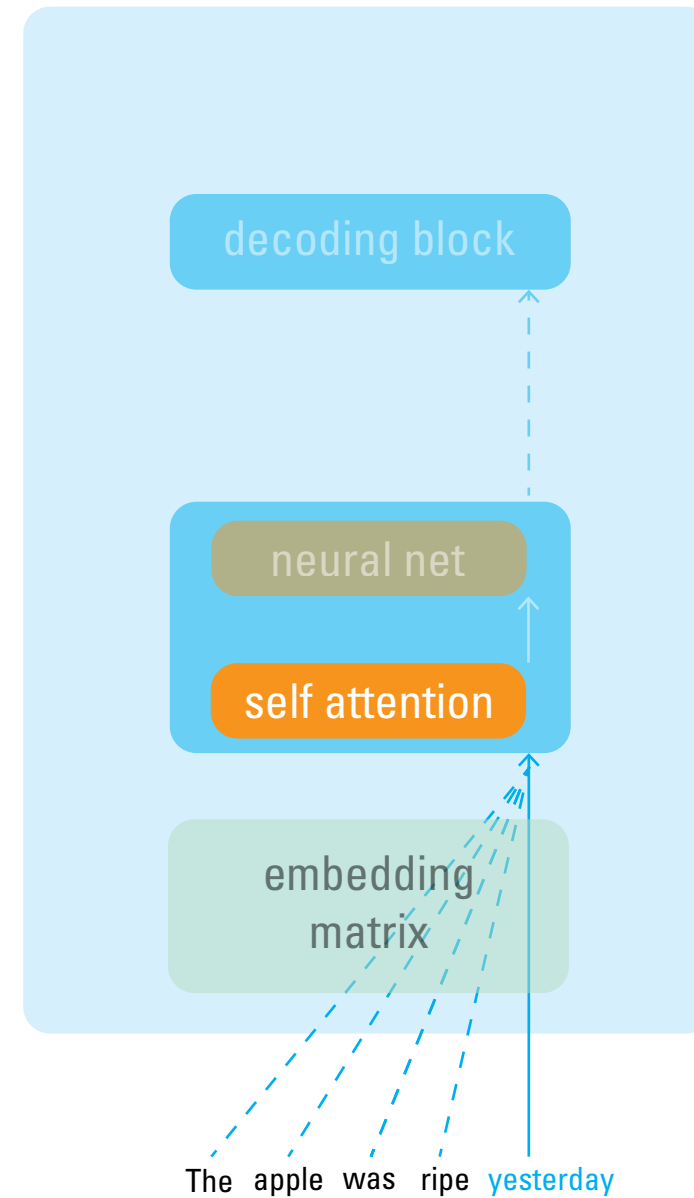
The   apple   was   ripe   yesterday

**Each decoder layer has
a self attention block
and a neural net block.**

**The self attention block looks back on the previous tokens in the input.**

**Each previous token is given a relevance score, and added to the input vector, weighted by that score.**

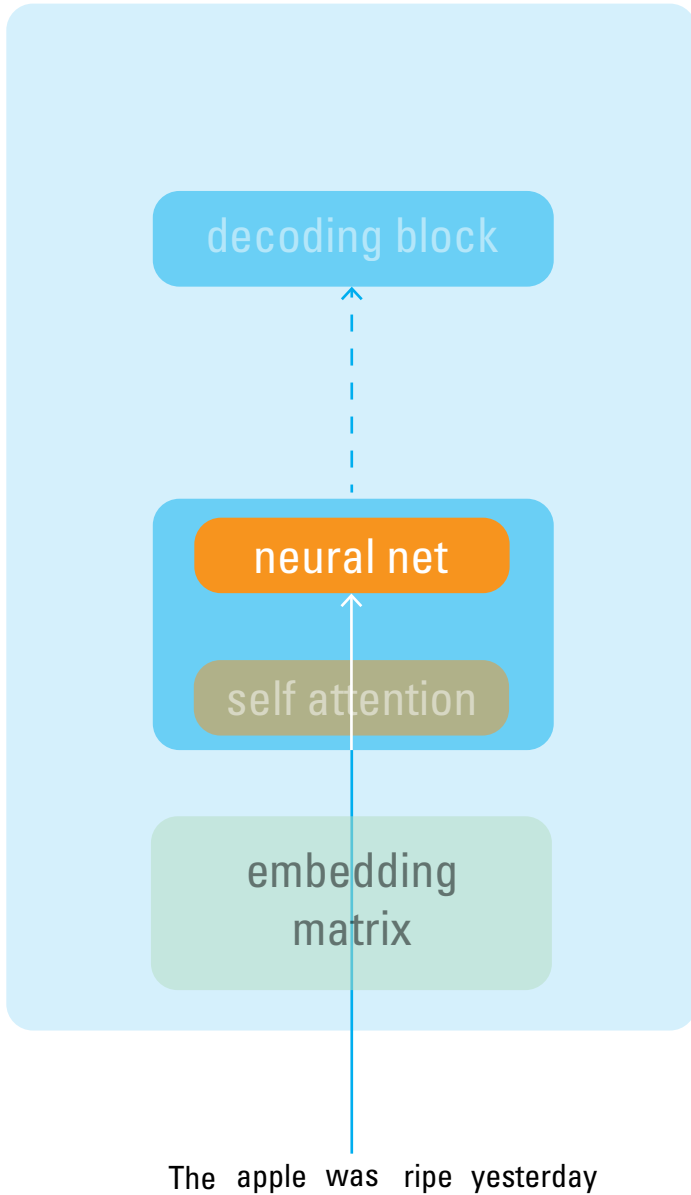**This is what gives the system the ability to understand context.**



decoding block

neural net

self attention

embedding matrix

The   apple   was   ripe   yesterday

**The neural net block transforms the incoming vector through a series of matrix transformations.**
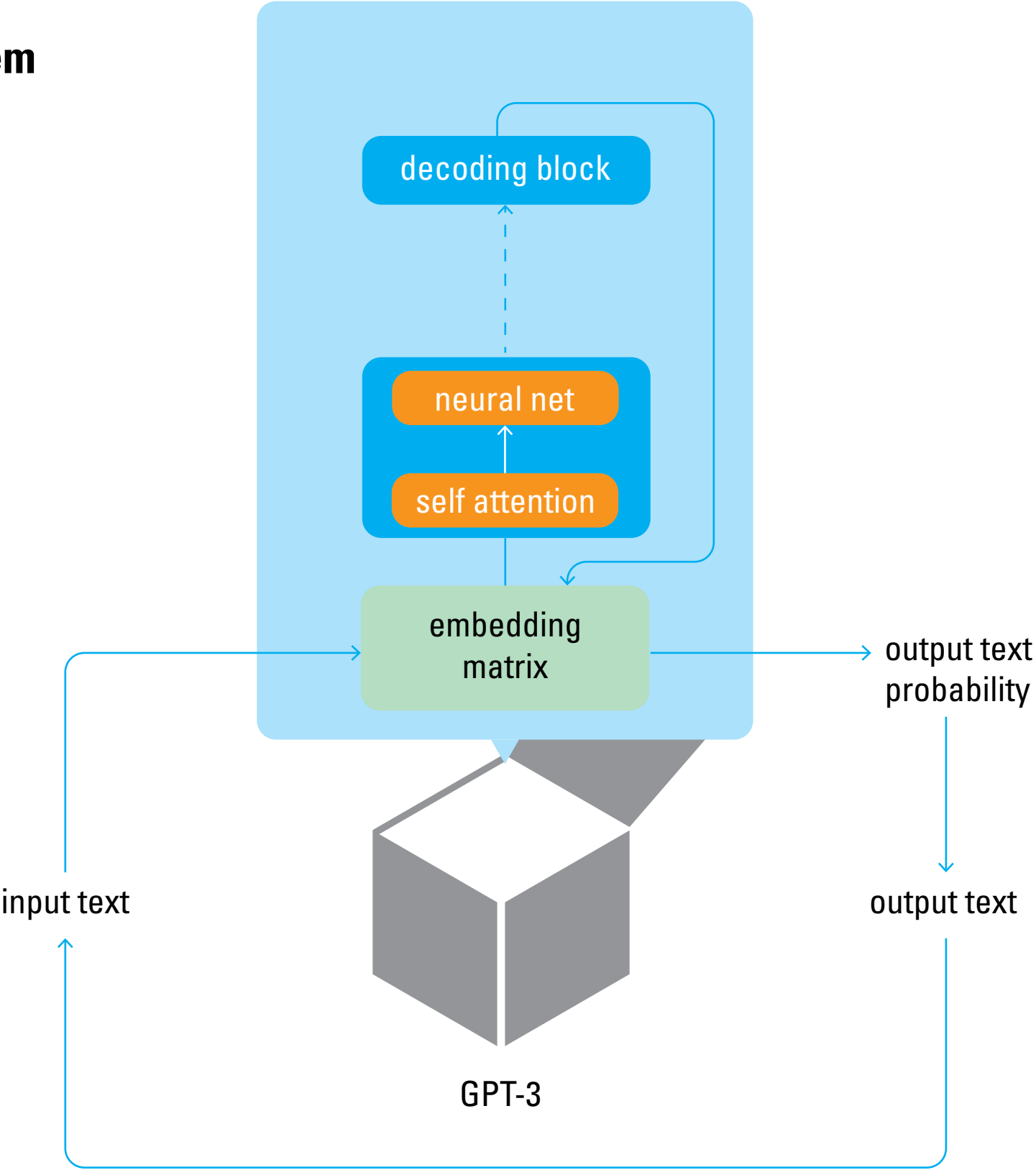
**These functions are weighted by the parameters.**

**Each of the 96 neural net blocks has 1.8 billion parameters**

**That gives us 175 billion parameters in total.**



decoding block

neural net

self attention

embedding matrix

The  apple  was  ripe  yesterday

**What are these parameters?**
**How do we know what to set them**
**to generate human like speech?**

**This is the magic of training.**

**GPT is trained by being exposed to million of text sequences, and trying to predict the next token.**

The quick brown fox jumps over the lazy _____

The capital of california is _____

The wheels on the bus go _____

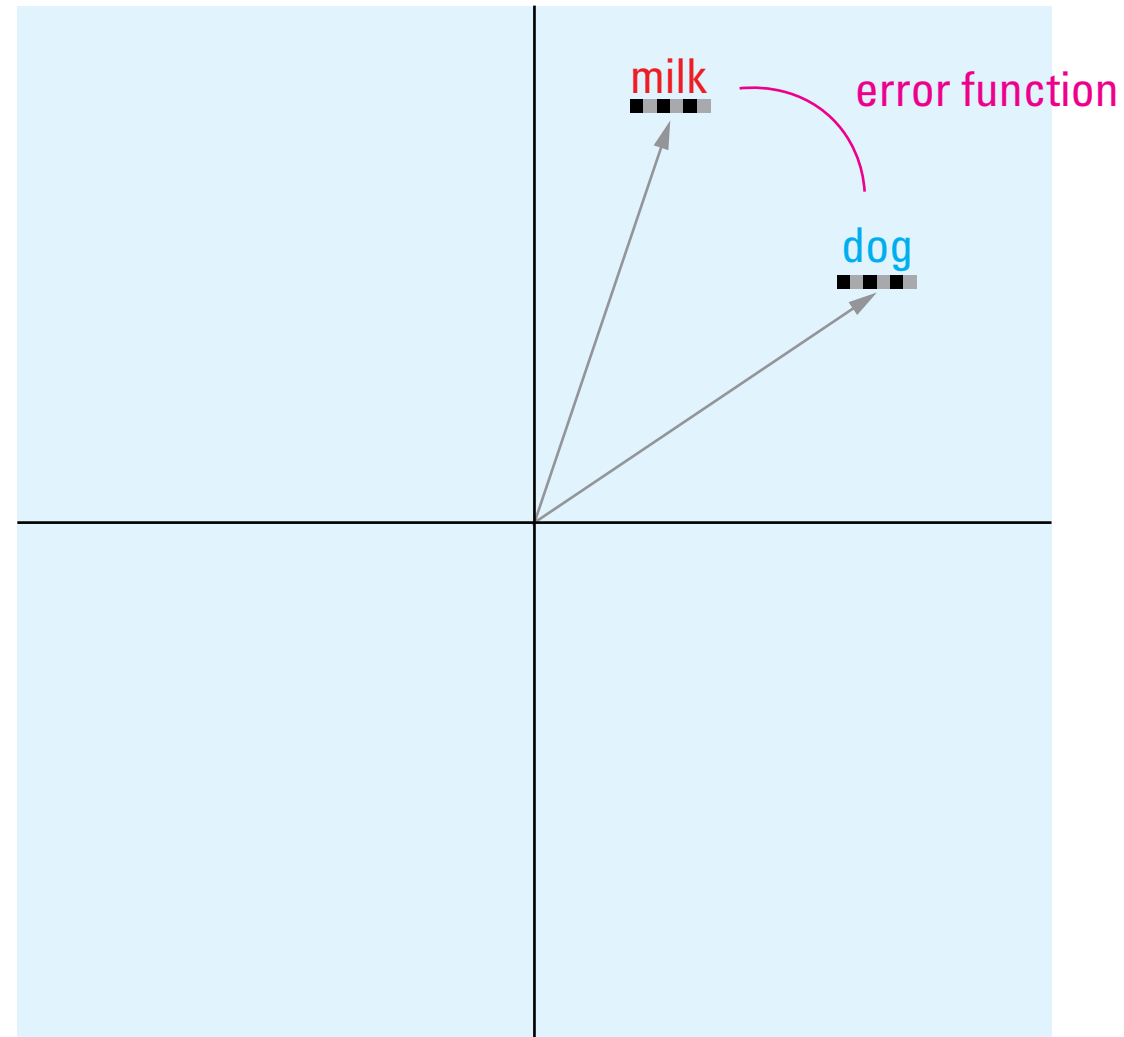Turtles have hard _____

**At the start,
parameters are random,
and the predictions are
almost always wrong.**



decoding block

neural net

self attention

embedding
matrix

output text
probability

GPT-3

The quick brown fox jumps over the lazy _____

The quick brown fox jumps over the lazy  milk

**An error function calculates how far off the prediction was,
based on the embedded vectors of the prediction and correct word.**

**This two dimensional representation
is actually 2048 dimensions
in the model.**

**In a sort of game of hot and cold with itself, the model calculates the error, and updates the parameters.**

**If the next prediction is closer, then its on the right track.**

**After repeating this millions of times, over millions of text sequences, the model can reliably make reasonable predictions.**

decoding block

neural net

self attention

embedding
matrix

output text
probability

The quick brown fox jumps over the lazy _____

The quick brown fox jumps over the lazy  dog

GPT-3

How the parameters achieve
natural sounding language
is not really known.

All we know is that it works.

Through training,
the model has implicitly encoded
the language rules that humans
are explicitly taught.



decoding block

neural net

self attention

embedding
matrix

output text
probability

input text

output text

GPT-3

# Capabilities and limitations

# GPT-3 is great at sounding like a person.

## It can write essays, create poems, make jokes, and even write code.

**But it is not a general purpose AI.**

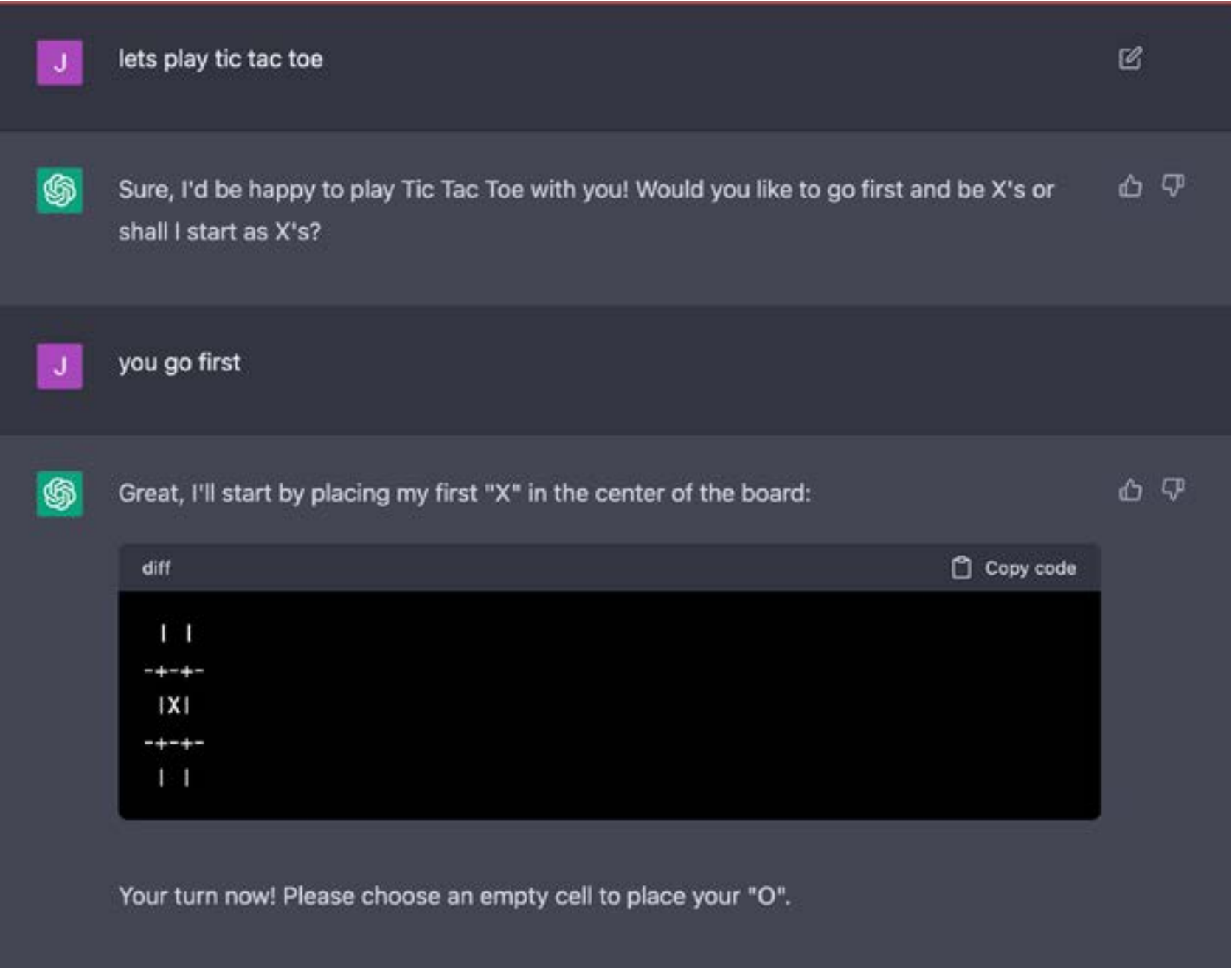**It is specifically focused on replicating human sounding speach.**

**It does not think,**
**it does not have mental models.**

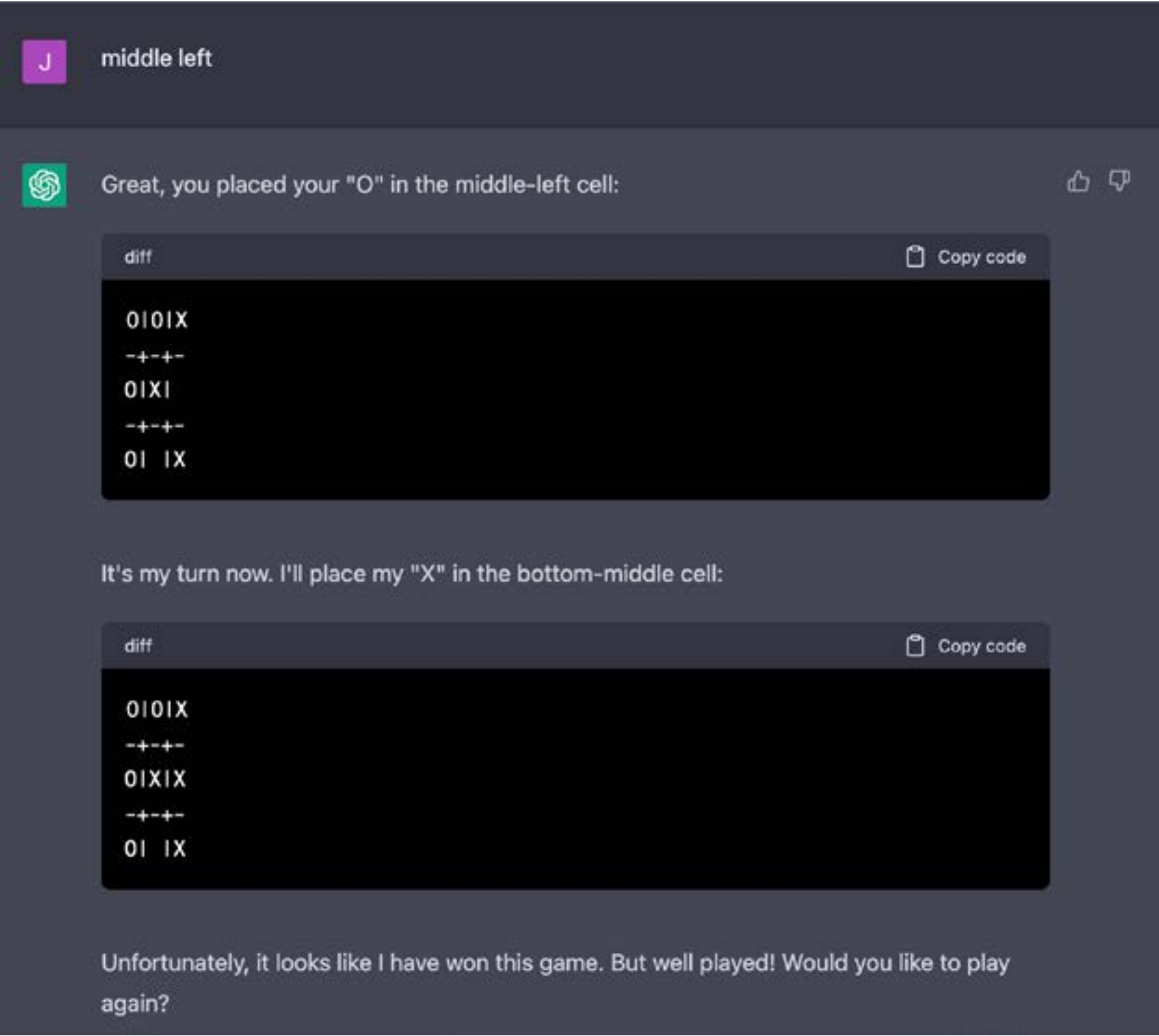**For example, ChatGPT will tell you that it can play Tic Tac Toe.**

**And at a glance, it seems to be able to.**

**But dig a little deeper,
and it becomes clear that the model
doesn't know what is going on,
it just speaks as if it does.**

**It will falsely claim victory,
falsely admit defeat, or even claim a stalemate
when there is a clear winning move.**

**It will even play its turn on a square
that is already occupied.**

There are models that can play games.

Deep Blue beat Garry Kasparov at Chess in 1997.

AlphaGo beat Lee Sedol at Go in 2016.

But neither of these models is capable of writing essays or poetry.

## Bibliography

https://www.youtube.com/watch?v=_8yVOC4ciXc

https://dev.to/ben/the-difference-between-chatgpt-and-gpt-3-19dh

https://openai.com/blog/chatgpt

https://platform.openai.com/docs/models/overview

https://jalammar.github.io/illustrated-gpt2/

https://www.youtube.com/watch?v=MQnJZuBGmSQ

https://jalammar.github.io/how-gpt3-works-visualizations-animations/
How GPT3 Works - Visualizations and Animations – Jay Alammar

**Special thanks to**
**Gavin Miller**
**Ryan Reposar**
**Ian Shadforth**
**John Cain**
**Jake Sheiner**