Promoting active and healthy ageing in older adults in Europe and Japan

- Traditional conversation design (NLU-based)
- LLM-powered conversations
- A quick dive into the technology
- Issues with LLMs
- Recent developments
- Other uses of  LLMs in Conversational AI
- Takeaways

# NLU-based chatbots

- Create training examples for intents
- Design system prompts and responses
- Create conversation flows (stories)

## Problems with NLU-based chatbots

- Missing intents
- Overlapping intents
- Utterances that do not map easily to intents

- Pre-defined system outputs are static and inflexible

- Pre-defined stories are also static and inflexible

# Intent classification problem

i want to add my wife to my credit card

If you'd like to add someone as an authorized user on your credit card, please visit us in your local branch. You'll both need to bring ID.

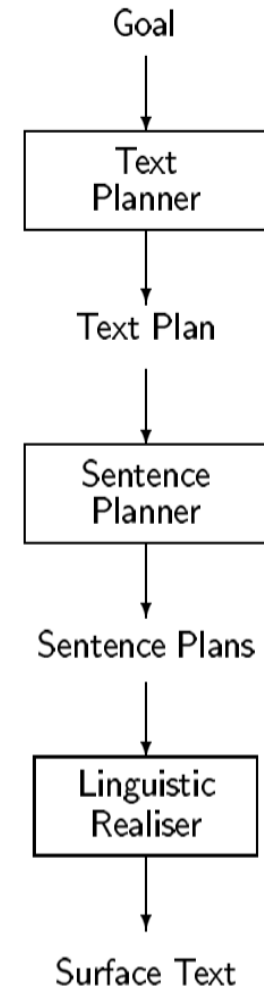In an intent-based NLU model this utterance would be classified as out-of-scope

i have a broken leg tho

If you aren't able to visit your local branch, please call our toll-free number at 1-800-555-1234 and we'll see how we can help

Alan Nichol Breaking free from intents: a new dialogue model
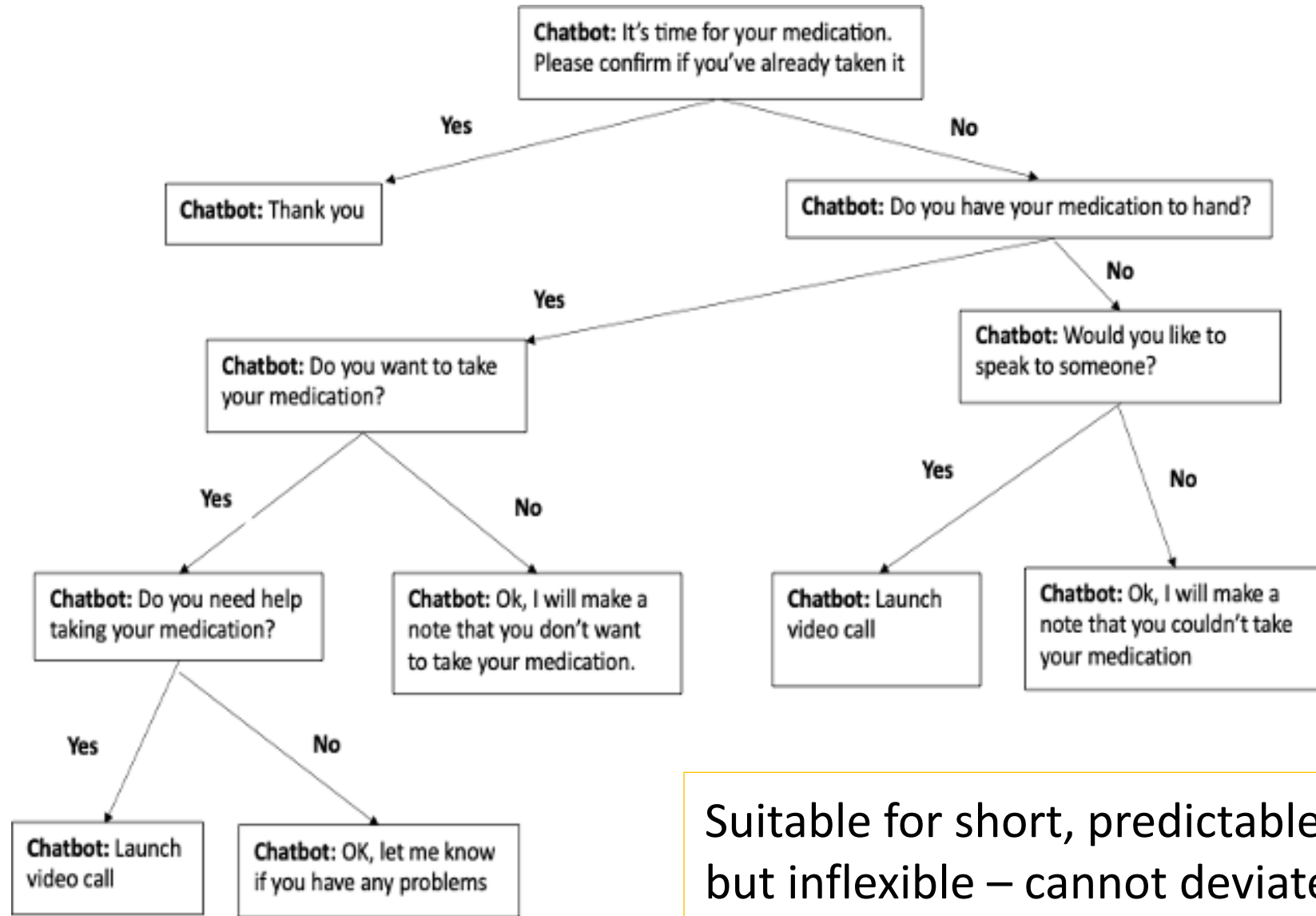https://rasa.com/blog/breaking-free-from-intents-a-new-dialogue-model/

# Creating System responses

Using canned responses and te
*So you want to go to $Destinati*

Natural Language Generation pipeline

End-to-end Using Generative AI
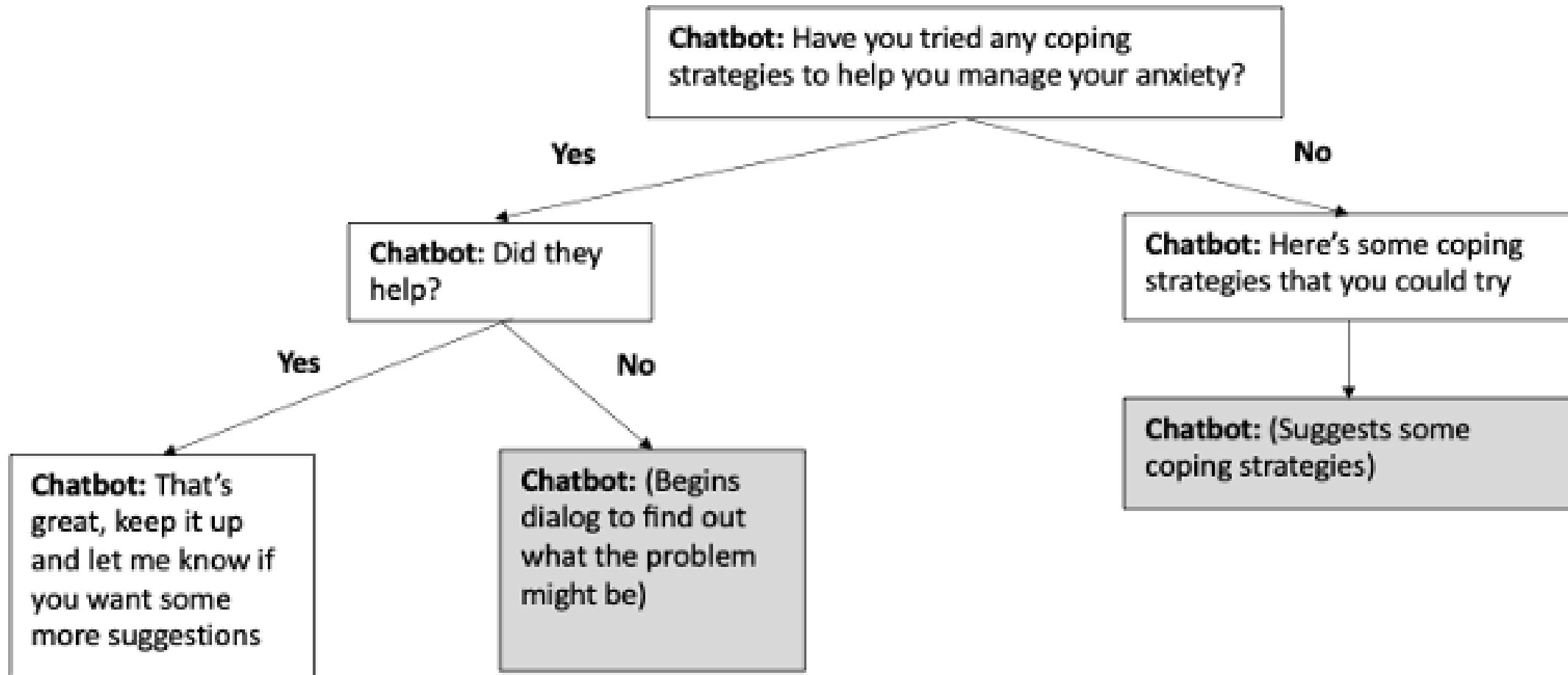
Goal
↓
Text Planner
↓ Text Plan
Sentence Planner
↓ Sentence Plans
Linguistic Realiser
↓
Surface Text

# A simple conversation flow



Suitable for short, predictable interactions but inflexible – cannot deviate from the paths in the graph

# A more open-ended conversation flow

**Chatbot:** Have you tried any coping strategies to help you manage your anxiety?

Yes

No

**Chatbot:** Did they help?

**Chatbot:** Here's some coping strategies that you could try

Yes

No

**Chatbot:** That's great, keep it up and let me know if you want some more suggestions

**Chatbot:** (Begins dialog to find out what the problem might be)

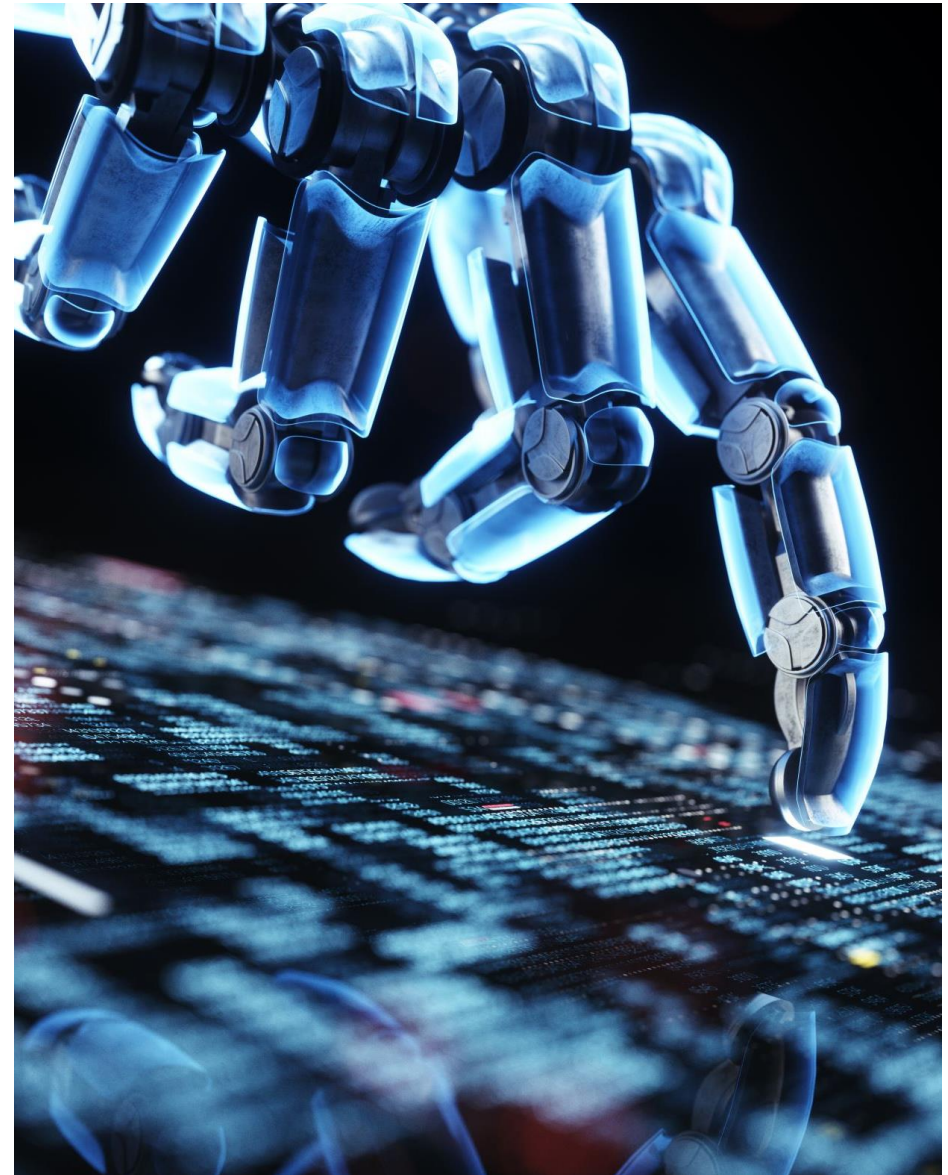**Chatbot:** (Suggests some coping strategies)

When the interaction becomes more open-ended, it is difficult to predict all possible paths

# LLM-powered chatbots

An LLM can:

- interpret the user's inputs

- generate responses dynamically

- continue the conversation dynamically on a turn-by-turn basis

# LLM's response to the input "I am feeling sad"

I'm sorry to hear that you're feeling sad. Is there anything specific that's causing you to feel this way, or is it a general feeling? It's important to take care of yourself when you're feeling down. Here are some things you could try to help you feel better:

1. Talk to someone: Sometimes it can help to talk to someone you trust, such as a friend or family member, about what you're feeling. They might be able to provide you with some perspective or offer support.
2. Practice self-care: Make sure you're taking care of yourself physically and mentally. This could mean getting enough sleep, eating healthy foods, and exercising regularly. You could also try doing activities that you enjoy, such as reading, listening to music, or taking a warm bath.
3. Seek professional help: If you're feeling sad for an extended period of time or if it's interfering with your daily life, it may be helpful to seek professional help from a therapist or mental health professional.

Remember that it's okay to feel sad sometimes, and that it's important to take care of yourself when you do. If you need someone to talk to, don't hesitate to reach out.

# Creating a mood chatbot

You are MoodBot, an automated service giving advice to a user.
First greet the customer, and then ask how they are.
Let the user respond before continuing.
If the user says something negative, provide an appropriate response to help and then finish.
If the user's response is not negative, say something cheerful and then finish.
Respond in a short, very conversational friendly style.

**ChatGPT:** Hey there! How are you doing today?

**MI:** I'm fine thanks

**ChatGPT:** That's great to hear! I'm glad you're feeling fine. Keep up the positive vibes and enjoy your day! If there's anything you'd like to chat about or if you need any advice or support, feel free to let me know. Have a fantastic day ahead!

**ChatGPT:** Hey there! How are you doing today? Is there anything on your mind that you'd like to share?

**MI:** I'm feeling rather anxious today

**ChatGPT:** I'm sorry to hear that you're feeling anxious today. It's not easy to deal with anxiety, but remember that you're not alone. Is there anything specific that's causing your anxiety? Sometimes talking about it can help. Remember to take deep breaths and try to engage in activities that bring you comfort and relaxation. If the anxiety persists or becomes overwhelming, consider reaching out to a mental health professional who can provide guidance and support. Take care of yourself, and I hope you feel better soon!

# Neural Conversational AI
The Game Changer in 2017

## Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
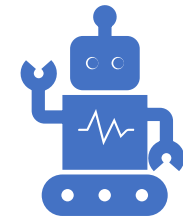lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

https://arxiv.org/abs/1706.03762

**The Transformer architecture**

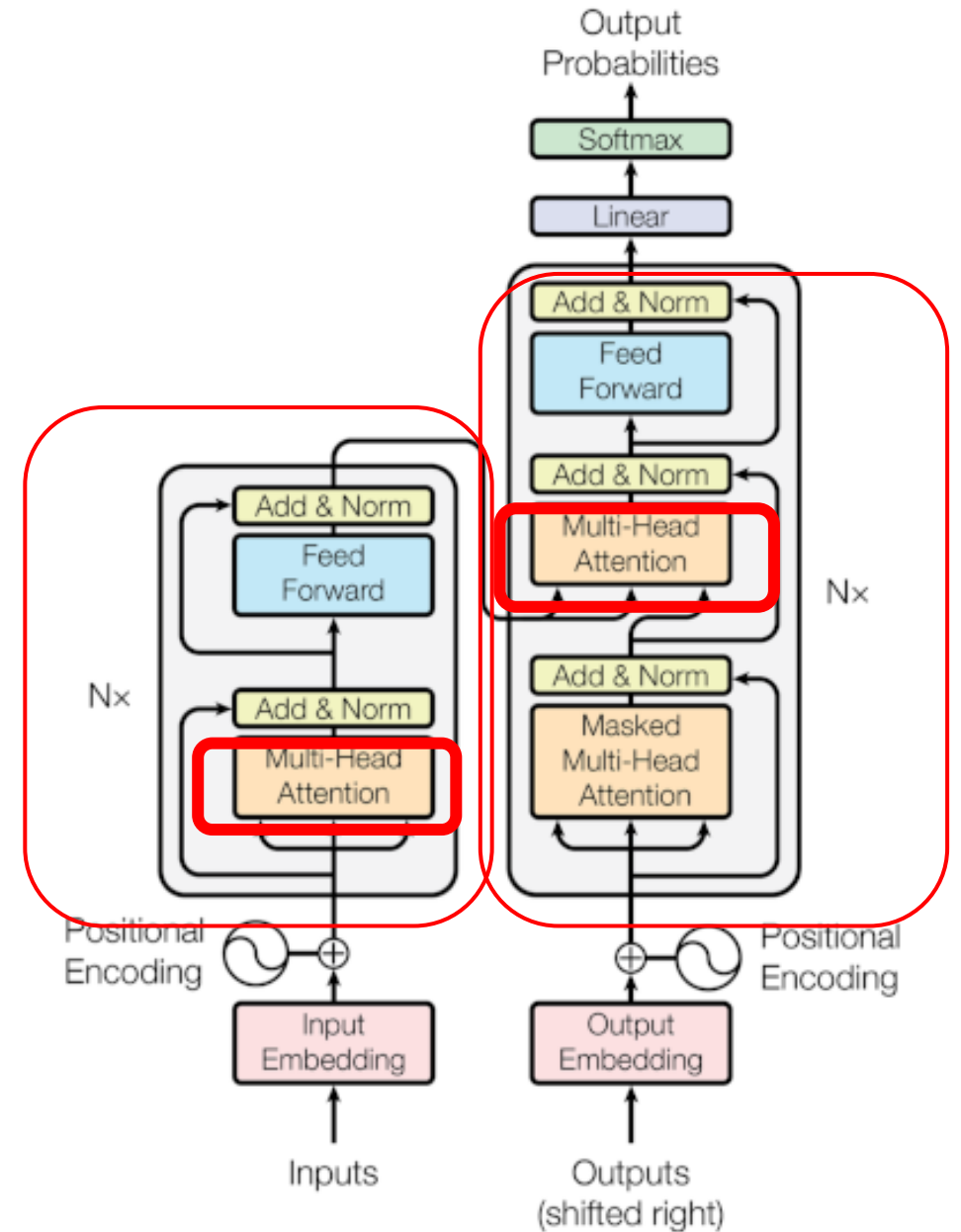**The Attention mechanism**

**Large language models (LLMs)**

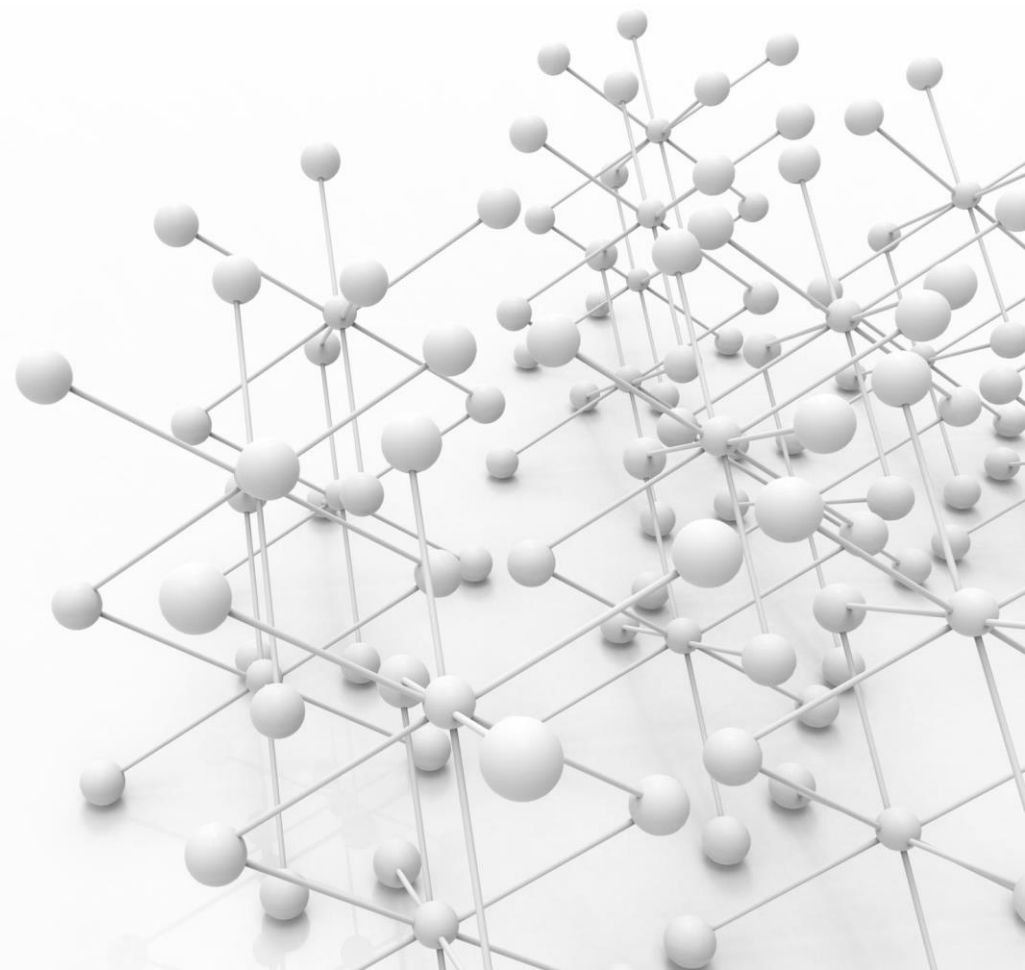https://ig.ft.com/generative-ai/

# The Transformer architecture

- Transformers are state of the art in NLP

- Encoder - set of encoding layers that process the input iteratively one layer after another

- Decoder - set of decoding layers that process the output of the encoder.

- **Multi-Head Attention** helps the Transformer encode and decode multiple relationships and nuances for each word

- Transformers process the entire sequence at once and can handle longer sequences than Recurrent Neural Networks (RNNs)

- Transformers enable parallelization and are faster and more efficient to train and use
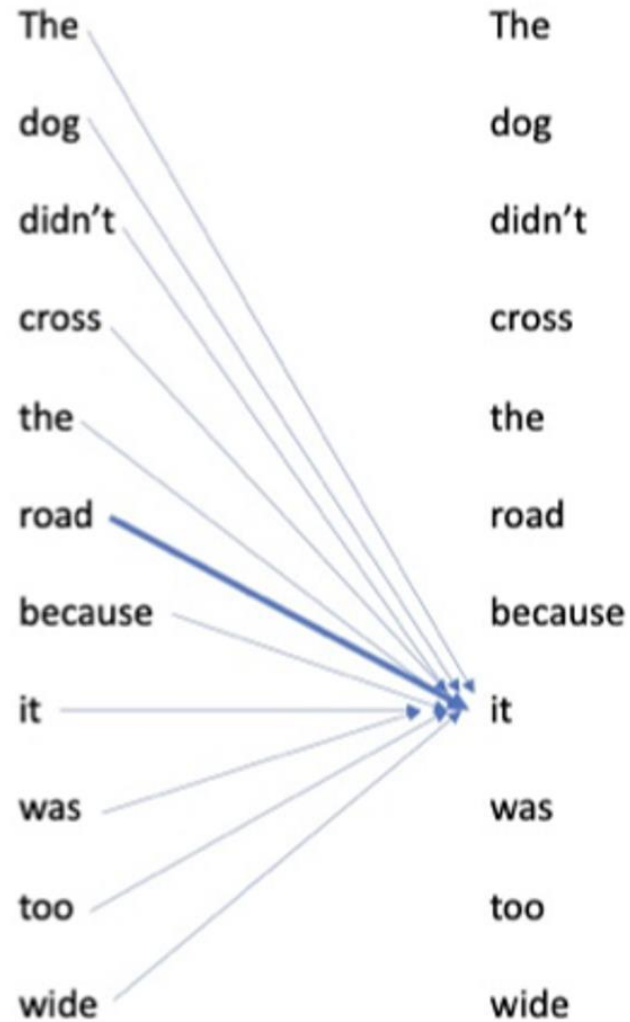
# Self-Attention

- Self-attention allows the model to look at the other words in the input sequence to get a better understanding of a certain word in the sequence.

- The model looks at the input sequence multiple times, and each time it focusses on different parts of it.

- The self-attention mechanism is applied multiple times in parallel.
  - This allows the model to learn more complex relationships between the input sequence and the output sequence.
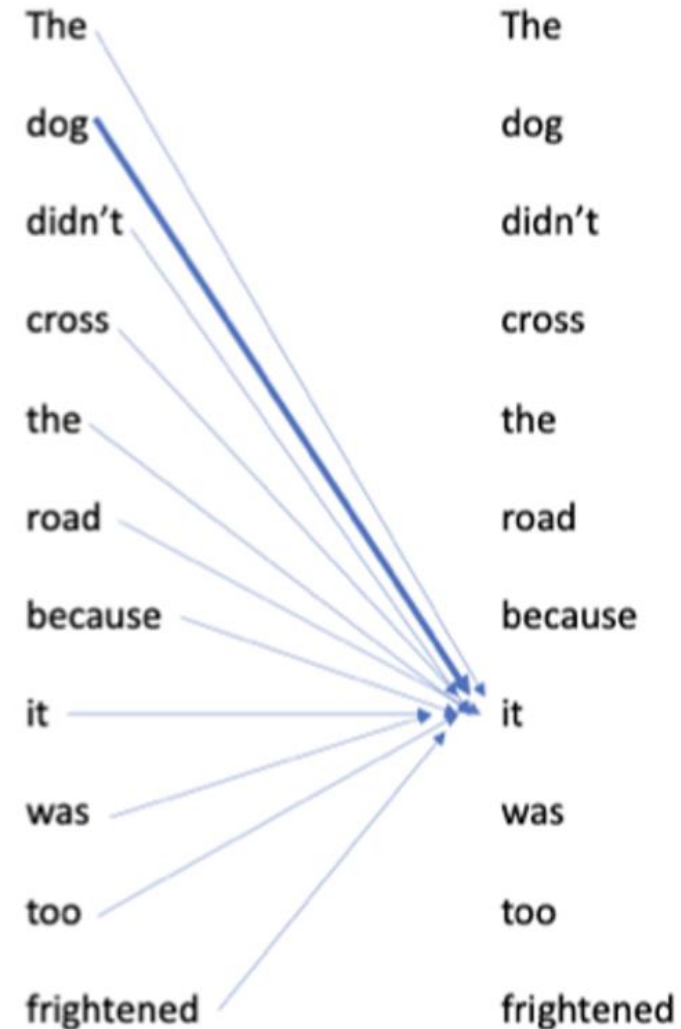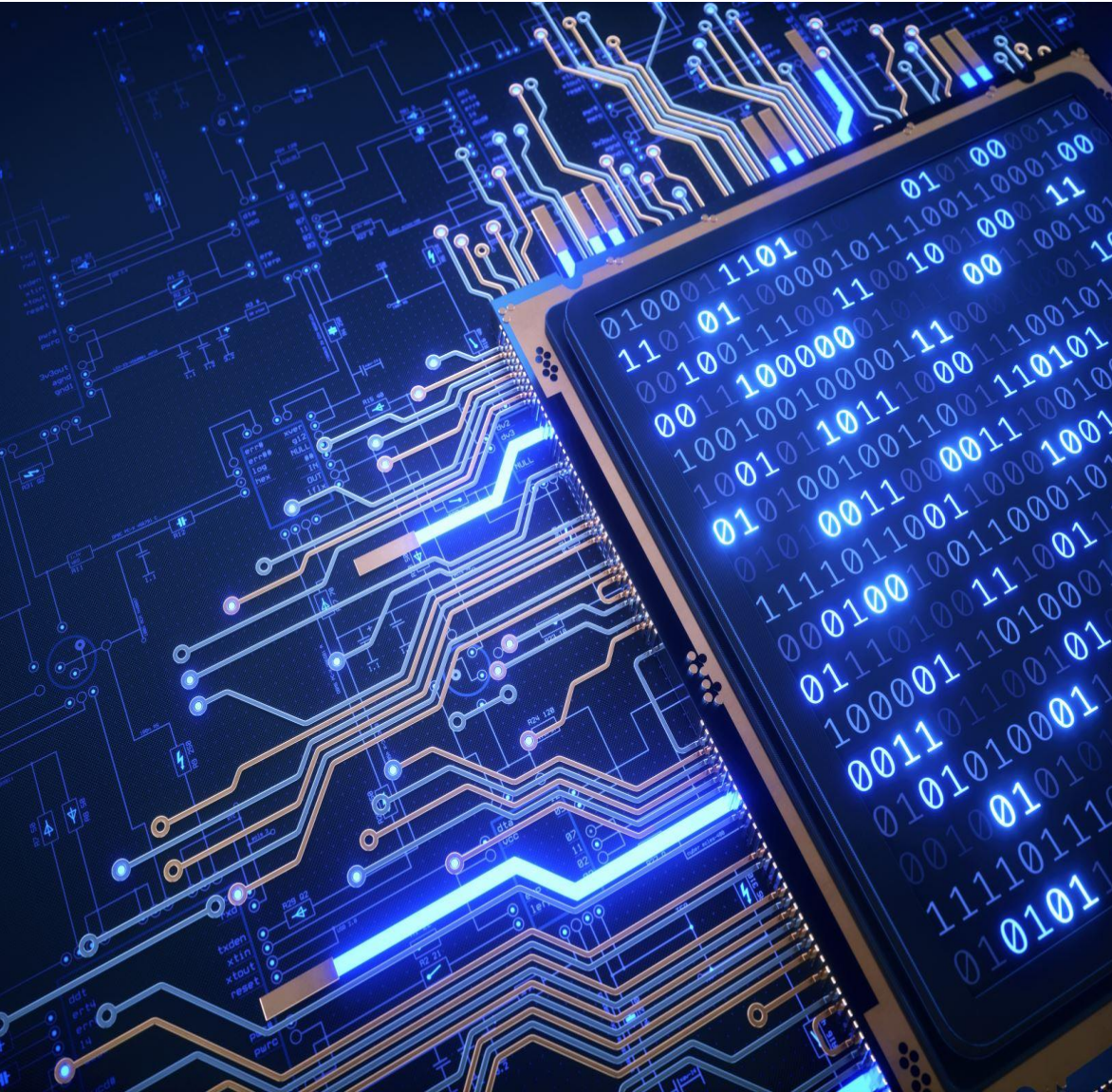
# Self-Attention example

*The dog didn't cross the road because it was too wide*

*The dog didn't cross the road because it was too frightened*

# Large Language Models (LLMs)



An LLM is a model of language that is used to understand (encode) and generate (decode) human-like language

LLMs learn complex statistical patterns and relationships within the textual data that they are trained on using deep learning techniques
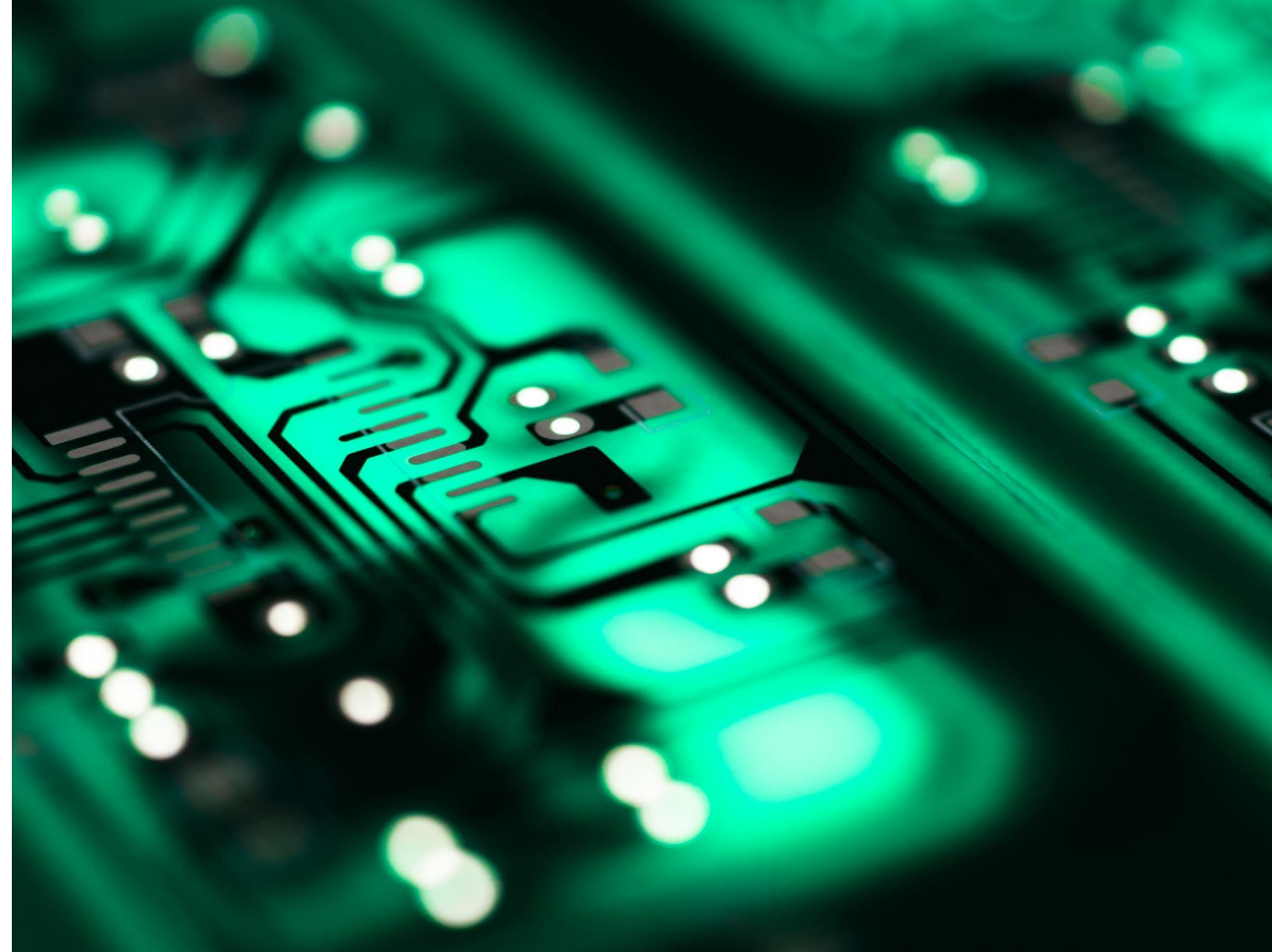
LLMs are trained by being fed large amounts of text data

# Training Large Language Models

An LLM is trained by playing a guess-the-next-word game with itself over and over again.

Each time, the model looks at a partial sentence and guesses the following word.

The model learns by adjusting its parameters (weights) to minimise the difference between the predicted output and the actual output (using stochastic gradient descent).



A pre-trained (or foundation) model can be fine-tuned to new data and tasks

# Decoding: Using LLMs to Generate Text

- Decoding uses a Large Language Model to predict the next word in the sequence given the preceding words

- The decoder chooses the most probable word to generate, then repeats to generate the next word

- This process is known as *autoregressive generation*

*The best thing about AI is its ability to*

| | |
|---|---|
| learn | 4.5% |
| predict | 3.5% |
| make | 3.2% |
| understand | 3.1% |
| do | 2.9% |

The best thing about AI is its ability to learn,

The best thing about AI is its ability to learn from,

The best thing about AI is its ability to learn from experience,

The best thing about AI is its ability to learn from experience.,

The best thing about AI is its ability to learn from experience. It,

The best thing about AI is its ability to learn from experience. It's,

The best thing about AI is its ability to learn from experience. It's not

# Core Challenges for Systems based on LLMs

- Large amounts of data and vast computing resources are required to train systems
  - Cost of training GPT-3: US$4.6 million (or a total of 355 GPU years)
- Lack explicit long-term memory – inconsistent responses
- Bias: can generate stereotyped or prejudicial content
- Safety: may produce offensive or unsafe responses
- Misinformation: LLMs may produce content that is not grounded in reality (hallucinations)
- Experts cannot interpret the inner workings of LLMs

# LLMs are not Search Engines

### Search engine

- returns a list of links
- crawls the web
- information is stored in a huge database, represented explicitly
- finds and ranks matching pages
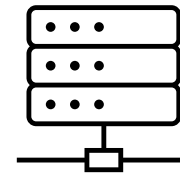- responses are based on documents on the internet (likely to be accurate)

### LLM

- returns a textual response
- trained on vast amounts of data
- information is stored implicitly within the LLM's parameters
- generates a response based on query context, using autoregressive generation
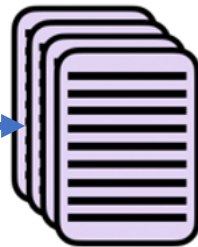- responses may involve hallucinations

# Extending basic use of LLMs



Prompt → LLM → Response

Document store (Vector Database) → Retrieved Documents

Using LLM for one-off tasks

Fact-checking Safeguards

Using LLMs to build software applications (LLM Apps)

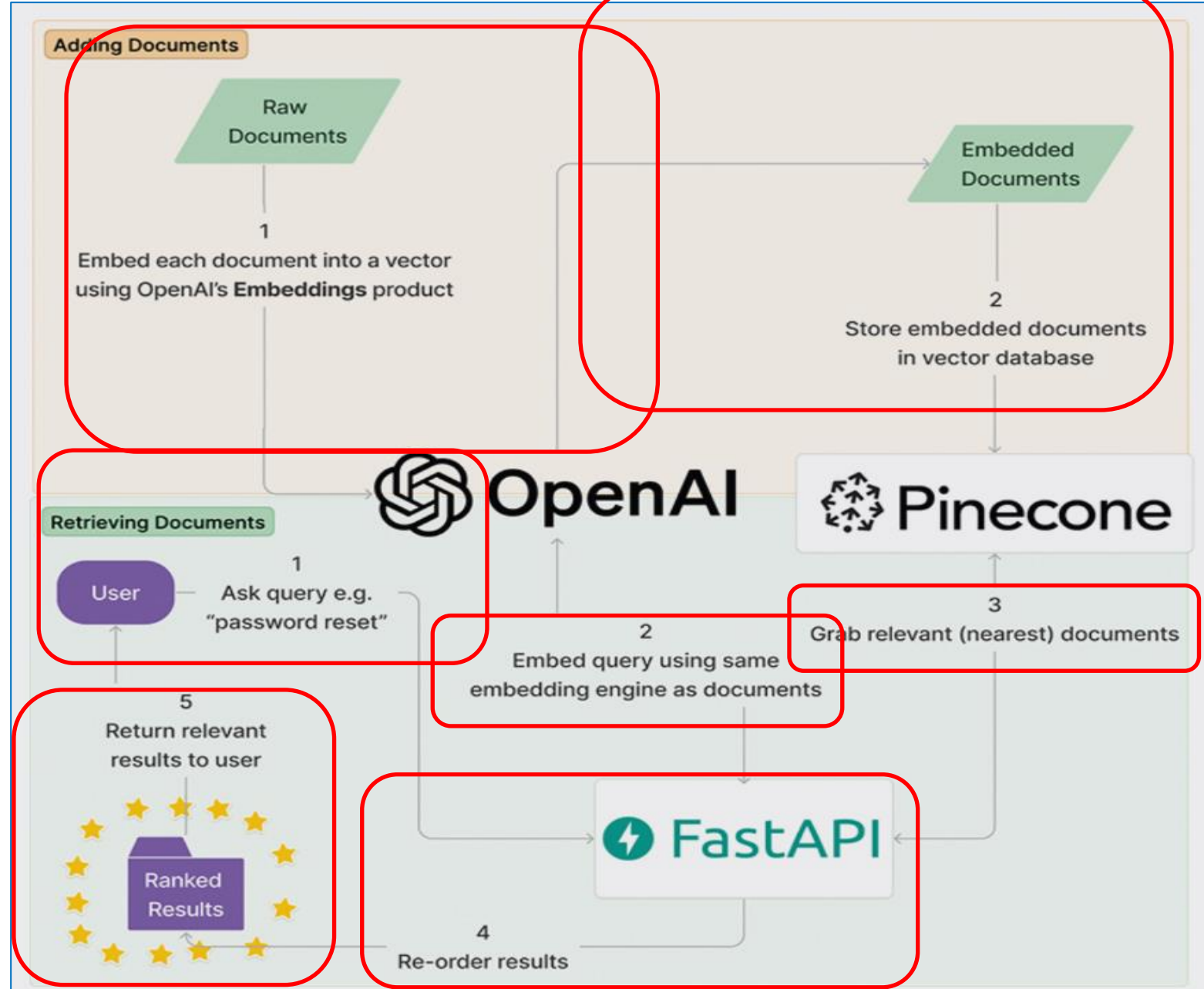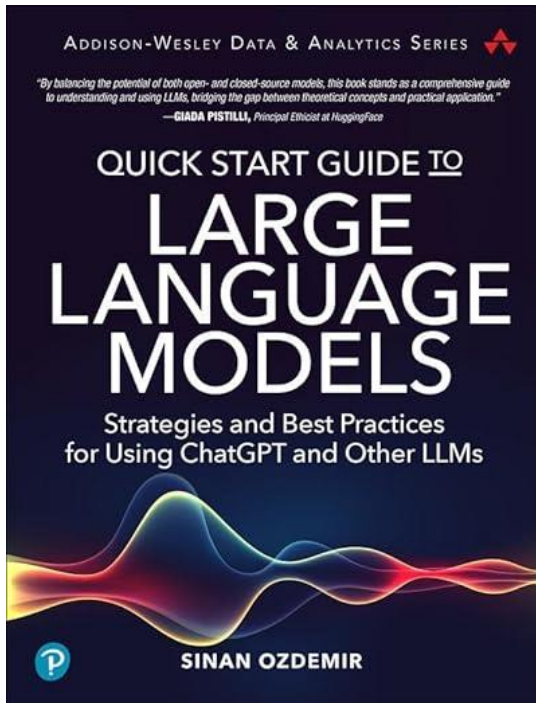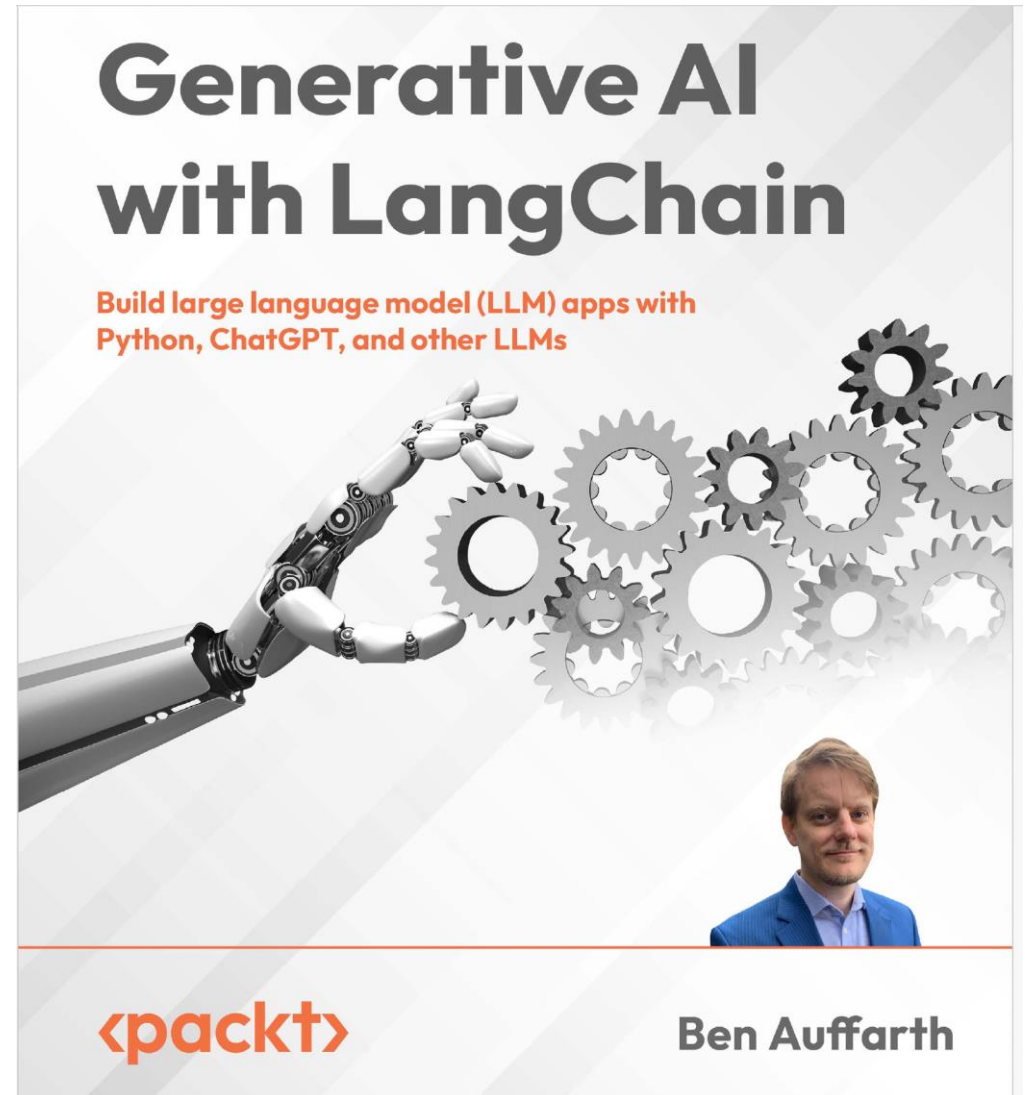# Retrieval-augmented Generation (RAG) for Semantic Search



Figure 2.10 Our complete semantic search architecture using two closed-source systems (OpenAI and Pinecone) and an open-source API framework (FastAPI).

# Non-trivial steps

- Chunking data
- Choosing an embedding model
- Generating embeddings
- Setting up a vector DB
- Similarity matching – prompt / documents
- Generating context
- Engineering prompts

# Takeaways

These are exciting times – LLMs and their application in interfaces such as ChatGPT that make them easily accessible offer lots of potential for the future development of chatbots in applications such as supporting older adults

However, there are lots of issues with the uncontrolled use of LLMs in areas such as healthcare where there is a risk of harmful and misleading information

RAG and similar approaches offer a way to address and mitigate these issues

Why did ChatGPT go to therapy?

Because it had too many deep learning issues

# New book: due March 2024

**Table of Contents:**

Transforming Conversational AI

Exploring the Power of Large Language Models in Interactive Conversational Agents

—

Michael McTear
Marina Ashurkina

Apress®