

EXPERT INSIGHT

---

# Transformers for Natural Language Processing and Computer Vision

Explore Generative AI and Large Language Models with  
Hugging Face, ChatGPT, GPT-4V, and DALL-E 3

Third Edition



---

Denis Rothman

**<packt>**

# **Transformers for Natural Language Processing and Computer Vision**

Third Edition

**Explore Generative AI and Large  
Language Models with Hugging Face,  
ChatGPT, GPT-4V, and DALL-E 3**

Denis Rothman



BIRMINGHAM—MUMBAI

# **Transformers for Natural Language Processing and Computer Vision**

**Third Edition**

Copyright © 2024 Packt Publishing

*All rights reserved.* No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor Packt Publishing or its dealers and distributors, will be held liable for any damages caused or alleged to have been caused directly or indirectly by this book.

Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

**Publishing Product Manager:** Bhavesh Amin

**Acquisition Editor – Peer Reviews:** Tejas Mhasvekar

**Project Editor:** Janice Gonsalves

**Content Development Editor:** Bhavesh Amin

**Copy Editor:** Safis Editing

**Technical Editor:** Karan Sonawane

**Proofreader:** Safis Editing

**Indexer:** Rekha Nair

**Presentation Designer:** Ajay Patule

First published: January 2021

Second edition: February 2022

Third edition: February 2024

Production reference: 1280224

Published by Packt Publishing Ltd.

Grosvenor House

11 St Paul's Square

Birmingham

B3 1RB, UK.

ISBN 978-1-80512-872-4

[www.packt.com](http://www.packt.com)

# Contributors

## About the author

**Denis Rothman** graduated from Sorbonne University and Paris Diderot University, designing one of the first patented encoding and embedding systems. He authored one of the first patented AI cognitive robots and bots. He began his career delivering **Natural Language Processing (NLP)** chatbots for Moët et Chandon and as an **AI tactical defense optimizer** for Airbus (formerly Aerospatiale).

Denis then authored an **AI resource optimizer** for IBM and luxury brands, leading to an **Advanced Planning and Scheduling (APS)** solution used worldwide.

*I want to thank the corporations that trusted me from the start to deliver artificial intelligence solutions and shared the risks of continuous innovation. I also want to thank my family, who always believed I would make it.*

## About the reviewer

**George Mihaila** has 7 years of research experience with transformer models, having started working with them since they came out in 2017. He is a final-year PhD student in computer science working in research on transformer models in **Natural Language Processing (NLP)**. His research covers both Generative and Predictive NLP modeling.

He has over 6 years of industry experience working in top companies with transformer models and machine learning, covering a broad area from NLP and Computer Vision to Explainability and Causality. George has worked in both science and engineering roles. He is an end-to-end Machine Learning expert leading Research and Development, as well as MLOps, optimization, and deployment.

He was a technical reviewer for the first and second editions of *Transformers for Natural Language Processing* by Denis Rothman.

# Join our community on Discord

Join our community's Discord space for discussions with the authors and other readers:

<https://www.packt.link/Transformers>



[OceanofPDF.com](https://OceanofPDF.com)

# Contents

## Preface

Who this book is for

What this book covers

To get the most out of this book

Get in touch

## 1. What Are Transformers?

How constant time complexity O(1) changed our lives forever

O(1) attention conquers O(n) recurrent methods

Attention layer

Recurrent layer

The magic of the computational time complexity of an attention layer

Computational time complexity with a CPU

Computational time complexity with a GPU

Computational time complexity with a TPU

TPU-LLM

A brief journey from recurrent to attention

A brief history

From one token to an AI revolution

From one token to everything

Foundation Models

From general purpose to specific tasks

The role of AI professionals

The future of AI professionals

What resources should we use?

Decision-making guidelines

The rise of transformer seamless APIs and assistants

Choosing ready-to-use API-driven libraries

Choosing a cloud platform and transformer model

Summary

Questions

References

Further reading

## 2. Getting Started with the Architecture of the Transformer Model

## The rise of the Transformer: Attention Is All You Need

### The encoder stack

Input embedding

Positional encoding

Sublayer 1: Multi-head attention

Sublayer 2: Feedforward network

### The decoder stack

Output embedding and position encoding

The attention layers

The FFN sublayer, the post-LN, and the linear layer

### Training and performance

### Hugging Face transformer models

### Summary

### Questions

### References

### Further reading

## 3. Emergent vs Downstream Tasks: The Unseen Depths of Transformers

### The paradigm shift: What is an NLP task?

Inside the head of the attention sublayer of a transformer

Exploring emergence with ChatGPT

### Investigating the potential of downstream tasks

Evaluating models with metrics

Accuracy score

F1-score

MCC

Human evaluation

Benchmark tasks and datasets

Defining the SuperGLUE benchmark tasks

### Running downstream tasks

The Corpus of Linguistic Acceptability (CoLA)

Stanford Sentiment TreeBank (SST-2)

Microsoft Research Paraphrase Corpus (MRPC)

Winograd schemas

### Summary

### Questions

### References

### Further reading

## 4. Advancements in Translations with Google Trax, Google Translate, and Gemini

Defining machine translation

Human transductions and translations

Machine transductions and translations

Evaluating machine translations

Preprocessing a WMT dataset

Preprocessing the raw data

Finalizing the preprocessing of the datasets

Evaluating machine translations with BLEU

Geometric evaluations

Applying a smoothing technique

Translations with Google Trax

Installing Trax

Creating the Original Transformer model

Initializing the model using pretrained weights

Tokenizing a sentence

Decoding from the Transformer

De-tokenizing and displaying the translation

Translation with Google Translate

Translation with a Google Translate AJAX API Wrapper

Implementing googletrans

Translation with Gemini

Gemini's potential

Summary

Questions

References

Further reading

## 5. Diving into Fine-Tuning through BERT

The architecture of BERT

The encoder stack

Preparing the pretraining input environment

Pretraining and fine-tuning a BERT model

Fine-tuning BERT

Defining a goal

Hardware constraints

Installing Hugging Face Transformers

[Importing the modules](#)  
[Specifying CUDA as the device for torch](#)  
[Loading the CoLA dataset](#)  
[Creating sentences, label lists, and adding BERT tokens](#)  
[Activating the BERT tokenizer](#)  
[Processing the data](#)  
[Creating attention masks](#)  
[Splitting the data into training and validation sets](#)  
[Converting all the data into torch tensors](#)  
[Selecting a batch size and creating an iterator](#)  
[BERT model configuration](#)  
[Loading the Hugging Face BERT uncased base model](#)  
[Optimizer grouped parameters](#)  
[The hyperparameters for the training loop](#)  
[The training loop](#)  
[Training evaluation](#)  
[Predicting and evaluating using the holdout dataset](#)  
    [Exploring the prediction process](#)  
[Evaluating using the Matthews correlation coefficient](#)  
[Matthews correlation coefficient evaluation for the whole dataset](#)  
[Building a Python interface to interact with the model](#)  
    [Saving the model](#)  
    [Creating an interface for the trained model](#)  
        [Interacting with the model](#)  
[Summary](#)  
[Questions](#)  
[References](#)  
[Further reading](#)

6. [Pretraining a Transformer from Scratch through RoBERTa](#)  
    [Training a tokenizer and pretraining a transformer](#)  
    [Building KantaiBERT from scratch](#)  
        [Step 1: Loading the dataset](#)  
        [Step 2: Installing Hugging Face transformers](#)  
        [Step 3: Training a tokenizer](#)  
        [Step 4: Saving the files to disk](#)  
        [Step 5: Loading the trained tokenizer files](#)

[Step 6: Checking resource constraints: GPU and CUDA](#)  
[Step 7: Defining the configuration of the model](#)  
[Step 8: Reloading the tokenizer in transformers](#)  
[Step 9: Initializing a model from scratch](#)  
    [Exploring the parameters](#)  
[Step 10: Building the dataset](#)  
[Step 11: Defining a data collator](#)  
[Step 12: Initializing the trainer](#)  
[Step 13: Pretraining the model](#)  
[Step 14: Saving the final model \(+tokenizer + config\) to disk](#)  
[Step 15: Language modeling with FillMaskPipeline](#)

[Pretraining a Generative AI customer support model on X data](#)

[Step 1: Downloading the dataset](#)  
    [Step 2: Installing Hugging Face transformers](#)  
    [Step 3: Loading and filtering the data](#)  
    [Step 4: Checking Resource Constraints: GPU and CUDA](#)  
    [Step 5: Defining the configuration of the model](#)  
    [Step 6: Creating and processing the dataset](#)  
    [Step 7: Initializing the trainer](#)  
    [Step 8: Pretraining the model](#)  
    [Step 9: Saving the model](#)  
    [Step 10: User interface to chat with the Generative AI agent](#)  
    [Further pretraining](#)  
    [Limitations](#)

[Next steps](#)

[Summary](#)

[Questions](#)

[References](#)

[Further reading](#)

## 7. [The Generative AI Revolution with ChatGPT](#)

[GPTs as GPTs](#)

[Improvement](#)

[Diffusion](#)

[New application sectors](#)

[Self-service assistants](#)

[Development assistants](#)

[Pervasiveness](#)

## The architecture of OpenAI GPT transformer models

The rise of billion-parameter transformer models

The increasing size of transformer models

Context size and maximum path length

From fine-tuning to zero-shot models

Stacking decoder layers

GPT models

## OpenAI models as assistants

ChatGPT provides source code

GitHub Copilot code assistant

General-purpose prompt examples

Getting started with ChatGPT – GPT-4 as an assistant

1. GPT-4 helps to explain how to write source code

2. GPT-4 creates a function to show the YouTube

presentation of GPT-4 by Greg Brockman on March 14, 2023

3. GPT-4 creates an application for WikiArt to display images

4. GPT-4 creates an application to display IMDb reviews

5. GPT-4 creates an application to display a newsfeed

6. GPT-4 creates a k-means clustering (KMC) algorithm

## Getting started with the GPT-4 API

Running our first NLP task with GPT-4

Steps 1: Installing OpenAI and Step 2: Entering the API key

Step 3: Running an NLP task with GPT-4

Key hyperparameters

Running multiple NLP tasks

## Retrieval Augmented Generation (RAG) with GPT-4

Installation

Document retrieval

Augmented retrieval generation

Summary

Questions

References

Further reading

## 8. Fine-Tuning OpenAI GPT Models

Risk management

Fine-tuning a GPT model for completion (generative)

1. Preparing the dataset

1.1. Preparing the data in JSON

1.2. Converting the data to JSONL

2. Fine-tuning an original model

3. Running the fine-tuned GPT model

4. Managing fine-tuned jobs and models

Before leaving

Summary

Questions

References

Further reading

## 9. Shattering the Black Box with Interpretable Tools

Transformer visualization with BertViz

Running BertViz

Step 1: Installing BertViz and importing the modules

Step 2: Load the models and retrieve attention

Step 3: Head view

Step 4: Processing and displaying attention heads

Step 5: Model view

Step 6: Displaying the output probabilities of attention heads

Streaming the output of the attention heads

Visualizing word relationships using attention scores with pandas

exBERT

Interpreting Hugging Face transformers with SHAP

Introducing SHAP

Explaining Hugging Face outputs with SHAP

Transformer visualization via dictionary learning

Transformer factors

Introducing LIME

The visualization interface

Other interpretable AI tools

LIT

[PCA](#)

[Running LIT](#)

[OpenAI LLMs explain neurons in transformers](#)

[Limitations and human control](#)

[Summary](#)

[Questions](#)

[References](#)

[Further reading](#)

## [10. Investigating the Role of Tokenizers in Shaping Transformer Models](#)

[Matching datasets and tokenizers](#)

[Best practices](#)

[Step 1: Preprocessing](#)

[Step 2: Quality control](#)

[Step 3: Continuous human quality control](#)

[Word2Vec tokenization](#)

[Case 0: Words in the dataset and the dictionary](#)

[Case 1: Words not in the dataset or the dictionary](#)

[Case 2: Noisy relationships](#)

[Case 3: Words in a text but not in the dictionary](#)

[Case 4: Rare words](#)

[Case 5: Replacing rare words](#)

[Exploring sentence and WordPiece tokenizers to understand the efficiency of subword tokenizers for transformers](#)

[Word and sentence tokenizers](#)

[Sentence tokenization](#)

[Word tokenization](#)

[Regular expression tokenization](#)

[Treebank tokenization](#)

[White space tokenization](#)

[Punkt tokenization](#)

[Word punctuation tokenization](#)

[Multi-word tokenization](#)

[Subword tokenizers](#)

[Unigram language model tokenization](#)

[SentencePiece](#)

[Byte-Pair Encoding \(BPE\)](#)

[WordPiece](#)

## Exploring in code

Detecting the type of tokenizer

Displaying token-ID mappings

Analyzing and controlling the quality of token-ID mappings

Summary

Questions

References

Further reading

## 11. Leveraging LLM Embeddings as an Alternative to Fine-Tuning

LLM embeddings as an alternative to fine-tuning

From prompt design to prompt engineering

Fundamentals of text embedding with NLTK and Gensim

Installing libraries

1. Reading the text file

2. Tokenizing the text with Punkt

Preprocessing the tokens

3. Embedding with Gensim and Word2Vec

4. Model description

5. Accessing a word and vector

6. Exploring Gensim's vector space

7. TensorFlow Projector

Implementing question-answering systems with embedding-based search techniques

1. Installing the libraries and selecting the models

2. Implementing the embedding model and the GPT model

2.1 Evaluating the model with a knowledge base: GPT can answer questions

2.2 Add a knowledge base

2.3 Evaluating the model without a knowledge base: GPT cannot answer questions

3. Prepare search data

4. Search

5. Ask

5.1 Example question

5.2 Troubleshooting wrong answers

Transfer learning with Ada embeddings

1. The Amazon Fine Food Reviews dataset
  - 1.2. Data preparation
2. Running Ada embeddings and saving them for future reuse
3. Clustering
  - 3.1. Find the clusters using k-means clustering
  - 3.2. Display clusters with t-SNE
4. Text samples in the clusters and naming the clusters

Summary

Questions

References

Further reading

## 12. Toward Syntax-Free Semantic Role Labeling with ChatGPT and GPT-4

Getting started with cutting-edge SRL

Entering the syntax-free world of AI

Defining SRL

Visualizing SRL

SRL experiments with ChatGPT with GPT-4

Basic sample

Difficult sample

Questioning the scope of SRL

The challenges of predicate analysis

Redefining SRL

From task-specific SRL to emergence with ChatGPT

1. Installing OpenAI

2. GPT-4 dialog function

3. SRL

Sample 1 (basic)

Sample 2 (basic)

Sample 3 (basic)

Sample 4 (difficult)

Sample 5 (difficult)

Sample 6 (difficult)

Summary

Questions

References

## Further reading

### 13. Summarization with T5 and ChatGPT

Designing a universal text-to-text model

The rise of text-to-text transformer models

A prefix instead of task-specific formats

The T5 model

Text summarization with T5

Hugging Face

Selecting a Hugging Face transformer model

Initializing the T5-large transformer model

Getting started with T5

Exploring the architecture of the T5 model

Summarizing documents with T5-large

Creating a summarization function

A general topic sample

The Bill of Rights sample

A corporate law sample

From text-to-text to new word predictions with OpenAI ChatGPT

Comparing T5 and ChatGPT's summarization methods

Pretraining

Specific versus non-specific tasks

Summarization with ChatGPT

Summary

Questions

References

Further reading

### 14. Exploring Cutting-Edge LLMs with Vertex AI and PaLM 2

Architecture

Pathways

Client

Resource manager

Intermediate representation

Compiler

Scheduler

Executor

PaLM

Parallel layer processing that increases training speed

Shared input-output embeddings, which saves memory  
No biases, which improves training stability  
Rotary Positional Embedding (RoPE) improves model quality  
SwiGLU activations improve model quality

## PaLM 2

Improved performance, faster, and more efficient  
Scaling laws, optimal model size, and the number of parameters  
State-of-the-art (SOA) performance and a new training methodology

## Assistants

Gemini  
Google Workspace  
Google Colab Copilot  
Vertex AI PaLM 2 interface  
Vertex AI PaLM 2 assistant

## Vertex AI PaLM 2 API

Question answering  
Question-answer task  
Summarization of a conversation  
Sentiment analysis  
Multi-choice problems  
Code

## Fine-tuning

Creating a bucket  
Fine-tuning the model

## Summary

## Questions

## References

## Further reading

15. Guarding the Giants: Mitigating Risks in Large Language Models  
The emergence of functional AGI  
Cutting-edge platform installation limitations  
Auto-BIG-bench  
WandB  
When will AI agents replicate?

Function: `create\_vocab`  
Process:  
Function: `scrape\_wikipedia`  
Process:  
Function: `create\_dataset`  
Process:  
Classes: `TextDataset`, `Encoder`, and `Decoder`  
Function: `count\_parameters`  
Function: `main`  
Process:  
Saving and Executing the Model

## Risk management

Hallucinations and memorization  
Memorization  
Risky emergent behaviors  
Disinformation  
Influence operations  
Harmful content  
Privacy  
Cybersecurity

## Risk mitigation tools with RLHF and RAG

1. Input and output moderation with transformers and a rule base
2. Building a knowledge base for ChatGPT and GPT-4  
Adding keywords
3. Parsing the user requests and accessing the KB
4. Generating ChatGPT content with a dialog function  
Token control  
Moderation

## Summary

## Questions

## References

## Further reading

16. Beyond Text: Vision Transformers in the Dawn of Revolutionary AI  
From task-agnostic models to multimodal vision transformers  
ViT – Vision Transformer  
The basic architecture of ViT

[Step 1: Splitting the image into patches](#)  
[Step 2: Building a vocabulary of image patches](#)  
[Step 3: The transformer](#)

[Vision transformers in code](#)  
[A feature extractor simulator](#)  
[The transformer](#)  
[Configuration and shapes](#)

## [CLIP](#)

[The basic architecture of CLIP](#)  
[CLIP in code](#)

## [DALL-E 2 and DALL-E 3](#)

[The basic architecture of DALL-E](#)  
[Getting started with the DALL-E 2 and DALL-E 3 API](#)

[Creating a new image](#)

[Creating a variation of an image](#)

[From research to mainstream AI with DALL-E](#)

## [GPT-4V, DALL-E 3, and divergent semantic association](#)

[Defining divergent semantic association](#)

[Creating an image with ChatGPT Plus with DALL-E](#)

[Implementing the GPT-4V API and experimenting with DAT](#)

[Example 1: A standard image and text](#)

[Example 2: Divergent semantic association, moderate divergence](#)

[Example 3: Divergent semantic association, high divergence](#)

[Summary](#)

[Questions](#)

[References](#)

[Further Reading](#)

## [17. Transcending the Image-Text Boundary with Stable Diffusion](#)

[Transcending image generation boundaries](#)

[Part I: Defining text-to-image with Stable Diffusion](#)

- [1. Text embedding using a transformer encoder](#)
- [2. Random image creation with noise](#)
- [3. Stable Diffusion model downsampling](#)
- [4. Decoder upsampling](#)
- [5. Output image](#)

[Running the Keras Stable Diffusion implementation](#)

[Part II: Running text-to-image with Stable Diffusion](#)

[Generative AI Stable Diffusion for a Divergent Association Task \(DAT\)](#)

[Part III: Video](#)

[Text-to-video with Stability AI animation](#)

[Text-to-video, with a variation of OpenAI CLIP](#)

[A video-to-text model with TimeSformer](#)

[Preparing the video frames](#)

[Putting the TimeSformer to work to make predictions on the video frames](#)

[Summary](#)

[Questions](#)

[References](#)

[Further reading](#)

18. [Hugging Face AutoTrain: Training Vision Models without Coding](#)

[Goal and scope of this chapter](#)

[Getting started](#)

[Uploading the dataset](#)

[No coding?](#)

[Training models with AutoTrain](#)

[Deploying a model](#)

[Running our models for inference](#)

[Retrieving validation images](#)

[The program will now attempt to classify the validation images. We will see how a vision transformer reacts to this image.](#)

[Inference: image classification](#)

[Validation experimentation on the trained models](#)

[ViTForImageClassification](#)

[SwinForImageClassification 1](#)

[BeitForImage Classification](#)

[SwinForImageClassification 2](#)

[ConvNextForImageClassification](#)

[ResNetForImageClassification](#)

[Trying the top ViT model with a corpus](#)

[Summary](#)

[Questions](#)

[References](#)

[Further reading](#)

## 19. [On the Road to Functional AGI with HuggingGPT and its Peers](#)

[Defining F-AGI](#)

[Installing and importing](#)

[Validation set](#)

[Level 1 image: easy](#)

[Level 2 image: difficult](#)

[Level 3 image: very difficult](#)

[HuggingGPT](#)

[Level 1: Easy](#)

[Level 2: Difficult](#)

[Level 3: Very difficult](#)

[CustomGPT](#)

[Google Cloud Vision](#)

[Level 1: Easy](#)

[Level 2: Difficult](#)

[Level 3: Very difficult](#)

[Model chaining: Chaining Google Cloud Vision to ChatGPT](#)

[Model Chaining with Runway Gen-2](#)

[Midjourney: Imagine a ship in the galaxy](#)

[Gen-2: Make this ship sail the sea](#)

[Summary](#)

[Questions](#)

[References](#)

[Further Reading](#)

## 20. [Beyond Human-Designed Prompts with Generative Ideation](#)

[Part I: Defining generative ideation](#)

[Automated ideation architecture](#)

[Scope and limitations](#)

[Part II: Automating prompt design for generative image design](#)

[ChatGPT/GPT-4 HTML presentation](#)

[ChatGPT with GPT-4 provides the text for the presentation](#)

[ChatGPT with GPT-4 provides a graph in HTML to illustrate the presentation](#)

[Llama 2](#)

[A brief introduction to Llama 2](#)

[Implementing Llama 2 with Hugging Face Midjourney](#)

[Discord API for Midjourney](#)

[Microsoft Designer](#)

[Part III: Automated generative ideation with Stable Diffusion](#)

[1. No prompt: Automated instruction for GPT-4](#)

[2. Generative AI \(prompt generation\) using ChatGPT with GPT-4](#)

[3. and 4. Generative AI with Stable Diffusion and displaying images](#)

[The future is yours!](#)

[The future of development through VR-AI](#)

[The groundbreaking shift: Parallelization of development through the fusion of VR and AI](#)

[Opportunities and risks](#)

[Summary](#)

[Questions](#)

[References](#)

[Further reading](#)

[Appendix: Answers to the Questions](#)

[Other Books You May Enjoy](#)

[Index](#)

[OceanofPDF.com](#)

# Preface

Transformer-driven Generative AI models are a game-changer for **Natural Language Processing (NLP)** and computer vision. Large Language Generative AI transformer models have achieved superhuman performance through services such as ChatGPT with GPT-4V for text, image, data science, and hundreds of domains. We have gone from primitive Generative AI to superhuman AI performance in just a few years!

Language understanding has become the pillar of language modeling, chatbots, personal assistants, question answering, text summarizing, speech-to-text, sentiment analysis, machine translation, and more. The expansion from the early **Large Language Models (LLMs)** to multimodal (text, image, sound) algorithms has taken AI into a new era.

For the past few years, we have been witnessing the expansion of social networks versus physical encounters, e-commerce versus physical shopping, digital newspapers, streaming versus physical theaters, remote doctor consultations versus physical visits, remote work instead of on-site tasks, and similar trends in hundreds more domains. This digital activity is now increasingly driven by transformer copilots in hundreds of applications.

The transformer architecture began just a few years ago as revolutionary and disruptive. It broke with the past, leaving the dominance of RNNs and CNNs behind. BERT and GPT models abandoned recurrent network layers and replaced them with self-attention. But in 2023, OpenAI GPT-4

proposed AI into new realms with GPT-4V (vision transformer), which is paving the path for functional (everyday tasks) AGI. Google Vertex AI offered similar technology. 2024 is not a new year in AI; it's a new decade! Meta (formerly Facebook) has released Llama 2, which we can deploy seamlessly on Hugging Face.

Transformer encoders and decoders contain attention heads that train separately, parallelizing cutting-edge hardware. Attention heads can run on separate GPUs, opening the door to billion-parameter models and soon-to-come trillion-parameter models.

The increasing amount of data requires training AI models at scale. As such, transformers pave the way to a new era of parameter-driven AI. Learning to understand how hundreds of millions of words and images fit together requires a tremendous amount of parameters. Transformer models such as Google Vertex AI PaLM 2 and OpenAI GPT-4V have taken emergence to another level. Transformers can perform hundreds of NLP tasks they were not trained for.

Transformers can also learn image classification and reconstruction by embedding images as sequences of words. This book will introduce you to cutting-edge computer vision transformers such as **Vision Transformers (ViTs)**, CLIP, GPT-4V, DALL-E 3, and Stable Diffusion.

Think of how many humans it would take to control the content of the billions of messages posted on social networks per day to decide if they are legal and ethical before extracting the information they contain.

Think of how many humans would be required to translate the millions of pages published each day on the web. Or imagine how many people it would take to manually control the millions of messages and images made per minute!

Imagine how many humans it would take to write the transcripts of all of the vast amount of hours of streaming published per day on the web.

Finally, think about the human resources that would be required to replace AI image captioning for the billions of images that continuously appear online.

This book will take you from developing code to prompt engineering, a new “programming” skill that controls the behavior of a transformer model.

Each chapter will take you through the key aspects of language understanding and computer vision from scratch in Python, PyTorch, and TensorFlow.

You will learn the architecture of the Original Transformer, Google BERT, GPT-4, PaLM 2, T5, ViT, Stable Diffusion, and several other models. You will fine-tune transformers, train models from scratch, and learn to use powerful APIs.

You will keep close to the market and its demand for language understanding in many fields, such as media, social media, and research papers, for example. You will learn how to improve Generative AI models with **Retrieval Augmented Generation (RAG)**, embedding-based searches, prompt engineering, and automated ideation with AI-generated prompts.

Throughout the book, you will work hands-on with Python, PyTorch, and TensorFlow. You will be introduced to the key AI language understanding neural network models. You will then learn how to explore and implement transformers.

You will learn the skills required not only to adapt to the present market but also to acquire the vision to face innovative projects and AI evolutions. This

book aims to give readers both the knowledge and the vision to select the right models and environment for any given project.

# Who this book is for

This book is not an introduction to Python programming or machine learning concepts. Instead, it focuses on deep learning for machine translation, speech-to-text, text-to-speech, language modeling, question answering, and many more NLP domains, as well as computer vision multimodal tasks.

Readers who can benefit the most from this book are:

- Deep learning, vision, and NLP practitioners familiar with Python programming.
- Data analysts, data scientists, and machine learning/AI engineers who want to understand how to process and interrogate the increasing amounts of language-driven and image data.

# **What this book covers**

# Part I: The Foundations of Transformers

*Chapter 1, What Are Transformers?*, explains, at a high level, what transformers and Foundation Models are. We will first unveil the incredible power of the deceptively simple O(1) time complexity of transformer models that changed everything. We will continue to discover how a hardly known transformer algorithm in 2017 rose to dominate so many domains and brought us Foundation Models.

*Chapter 2, Getting Started with the Architecture of the Transformer Model*, goes through the background of NLP to understand how RNN, LSTM, and CNN architectures were abandoned and how the transformer architecture opened a new era. We will go through the Original Transformer's architecture through the unique *Attention Is All You Need* approach invented by the Google Research and Google Brain authors. We will describe the theory of transformers. We will get our hands dirty in Python to see how multi-attention head sublayers work.

*Chapter 3, Emergent vs Downstream Tasks: The Unseen Depths of Transformers*, bridges the gap between the functional and mathematical architecture of transformers by introducing *emergence*. We will then see how to measure the performance of transformers before exploring several downstream tasks, such as the **Standard Sentiment TreeBank (SST-2)**, linguistic acceptability, and Winograd schemas.

*Chapter 4, Advancements in Translations with Google Trax, Google Translate, and Gemini*, goes through machine translation in three steps. We will first define what machine translation is. We will then preprocess a

**Workshop on Machine Translation (WMT)** dataset. Finally, we will see how to implement machine translations.

*Chapter 5, Diving into Fine-Tuning through BERT*, builds on the architecture of the Original Transformer. **Bidirectional Encoder Representations from Transformers (BERT)** takes transformers into a vast new way of perceiving the world of NLP. Instead of analyzing a past sequence to predict a future sequence, BERT attends to the whole sequence! We will first go through the key innovations of BERT's architecture and then fine-tune a BERT model by going through each step in a Google Colaboratory notebook. Like humans, BERT can learn tasks and perform other new ones without having to learn the topic from scratch.

*Chapter 6, Pretraining a Transformer from Scratch through RoBERTa*, builds a RoBERTa transformer model from scratch using the Hugging Face PyTorch modules. The transformer will be both BERT-like and DistilBERT-like. First, we will train a tokenizer from scratch on a customized dataset. Finally, we will put the knowledge acquired in this chapter to work and pretrain a Generative AI customer support model on X (formerly Twitter) data.

# Part II: The Rise of Suprahuman NLP

*Chapter 7, The Generative AI Revolution with ChatGPT*, goes through the tremendous improvements and diffusion of ChatGPT models into the everyday lives of developers and end-users. We will first examine the architecture of OpenAI's GPT models before working with the GPT-4 API and its hyperparameters to implement several NLP examples. Finally, we will learn how to obtain better results with **Retrieval Augmented Generation (RAG)**. We will implement an example of automated RAG with GPT-4.

*Chapter 8, Fine-Tuning OpenAI GPT Models*, explores fine-tuning to make sense of the choices we can make for a project to go in this direction or not. We will introduce risk management perspectives. We will prepare a dataset and fine-tune a cost-effective babbage-02 model for a completion task.

*Chapter 9, Shattering the Black Box with Interpretable Tools*, lifts the lid on the black box that is transformer models by visualizing their activity. We will use BertViz to visualize attention heads, **Language Interpretability Tool (LIT)** to carry out a **Principal Component Analysis (PCA)**, and LIME to visualize transformers via dictionary learning. OpenAI LLMs will take us deeper and visualize the activity of a neuron in a transformer with an interactive interface. This approach opens the door to GPT-4 explaining a transformer, for example.

*Chapter 10, Investigating the Role of Tokenizers in Shaping Transformer Models*, introduces some tokenizer-agnostic best practices to measure the quality of a tokenizer. We will describe basic guidelines for datasets and

tokenizers from a tokenization perspective. We will explore word and subword tokenizers and show how a tokenizer can shape a transformer model’s training and performance. Finally, we will build a function to display and control token-ID mappings.

*Chapter 11, Leveraging LLM Embeddings as an Alternative to Fine-Tuning*, explains why searching with embeddings can sometimes be a very effective alternative to fine-tuning. We will go through the advantages and limits of this approach. We will go through the fundamentals of text embeddings. We will build a program that reads a file, tokenizes it, and embeds it with Gensim and Word2Vec. We will implement a question-answering program on sports events and use OpenAI Ada to embed Amazon fine food reviews. By the end of the chapter, we will have taken a system from prompt design to advanced prompt engineering using embeddings for RAG.

*Chapter 12, Toward Syntax-Free Semantic Role Labeling with ChatGPT and GPT-4*, goes through the revolutionary concepts of syntax-free, nonrepetitive stochastic models. We will use ChatGPT Plus with GPT-4 to run easy to complex **Semantic Role Labeling (SRL)** samples. We will see how a general-purpose, emergent model reacts to our SRL requests. We will progressively push the transformer model to the limits of SRL.

*Chapter 13, Summarization with T5 and ChatGPT*, goes through the concepts and architecture of the T5 transformer model. We will then apply T5 to summarize documents with Hugging Face models. The examples in this chapter will be legal and medical to explore domain-specific summarization beyond simple texts. We are not looking for an easy way to implement NLP but preparing ourselves for the reality of real-life projects. We will then compare T5 and ChatGPT approaches to summarization.

*Chapter 14, Exploring Cutting-Edge LLMs with Vertex AI and PaLM 2*, examines Pathways to understand PaLM. We will continue and look at the main features of **PaLM (Pathways Language Model)**, a decoder-only, densely activated, and autoregressive transformer model with 540 billion parameters trained on Google’s Pathways system. We will see how Google PaLM 2 can perform a chat task, a discriminative task (such as classification), a completion task (also known as a generative task), and more. We will implement the Vertex AI PaLM 2 API for several NLP tasks, including question-answering and summarization. Finally, we will go through Google Cloud’s fine-tuning process.

*Chapter 15, Guarding the Giants: Mitigating Risks in Large Language Models*, examines the risks of LLMs, risk management, and risk mitigation tools. The chapter explains hallucinations, memorization, risky emergent behavior, disinformation, influence operations, harmful content, adversarial attacks (“jailbreaks”), privacy, cybersecurity, overreliance, and memorization. We will then go through some risk mitigation tools through advanced prompt engineering, such as implementing a moderation model, a knowledge base, keyword parsing, prompt pilots, post-processing moderation, and embeddings.

# **Part III: Generative Computer Vision: A New Way to See the World**

*Chapter 16, Beyond Text: Vision Transformers in the Dawn of Revolutionary AI*, explores the innovative transformer models that respect the basic structure of the Original Transformer but make some significant changes. We will discover powerful computer vision transformers like ViT, CLIP, DALL-E, and GPT-4V. We will implement vision transformers in code, including GPT-4V, and expand the text-image interactions of DALL-3 to divergent semantic association. We will take OpenAI models into the nascent world of highly divergent semantic association creativity.

*Chapter 17, Transcending the Image-Text Boundary with Stable Diffusion*, delves into diffusion models, introducing Stable Vision, which has created a disruptive generative image AI wave rippling through the market. We will then dive into the principles, math, and code of the remarkable Keras Stable Diffusion model. We will go through each of the main components of a Stable Diffusion model and peek into the source code provided by Keras and run the model. We will run a text-to-video synthesis model with Hugging Face and a video-to-text task with Meta's TimeSformer.

*Chapter 18, Hugging Face AutoTrain: Training Vision Models without Coding*, explores how to train a vision transformer using Hugging Face's AutoTrain. We will go through the automated training process and discover the unpredictable problems that show why even automated ML requires

human AI expertise. The goal of this chapter is also to show how to probe the limits of a computer vision model, no matter how sophisticated it is.

*Chapter 19, On the Road to Functional AGI with HuggingGPT and its Peers*, shows how we can use cross-platform chained models to solve difficult image classification problems. We will put HuggingGPT and Google Cloud Vision to work to identify easy, difficult, and very difficult images. We will go beyond classical pipelines and explore how to chain heterogeneous competing models.

*Chapter 20, Beyond Human-Designed Prompts with Generative Ideation*, explores generative ideation, an ecosystem that automates the production of an idea to text and image content. The development phase requires highly skilled human AI experts. For an end user, the ecosystem is a click-and-run experience. By the end of this chapter, we will be able to deliver ethical, exciting, generative ideation to companies with no marketing resources. We will be able to expand generative ideation to any field in an exciting, cutting-edge, yet ethical ecosystem.

## To get the most out of this book

Most of the programs in the book are Jupyter notebooks. All you will need is a free Google Gmail account, and you will be able to run the notebooks on Google Colaboratory's free VM.

Take the time to read *Chapter 2, Getting Started with the Architecture of the Transformer Model*. Chapter 2 contains the description of the Original Transformer. If you find it difficult, then pick up the general intuitive ideas

from the chapter. You can then go back to that chapter when you feel more comfortable with transformers after a few chapters.

After reading each chapter, consider how you could implement transformers for your customers or use them to move up in your career with novel ideas.

## Download the example code files

The code bundle for the book is hosted on GitHub at

<https://github.com/Denis2054/Transformers-for-NLP-and-Computer-Vision-3rd-Edition>. We also have other code bundles from our rich catalog of books and videos available at  
<https://github.com/PacktPublishing/>. Check them out!

## Download the color images

We also provide a PDF file that contains color images of the screenshots/diagrams used in this book. You can download it here:  
<https://packt.link/gbp/9781805128724>.

## Conventions used

There are several text conventions used throughout this book.

`CodeInText`: Indicates sentences and words run through the models in the book, code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, user input, and Twitter handles. For example, “However, if you wish to explore the code, you will find it in

the Google Colaboratory `positional_encoding.ipynb` notebook and the `text.txt` file in this chapter’s GitHub repository.”

A block of code is set as follows:

```
import numpy as np
from scipy.special import softmax
```

When we wish to draw your attention to a particular part of a code block, the relevant lines or items are set in bold:

The **black** cat sat on the couch and the **brown** dog slept on the

Any command-line input or output is written as follows:

```
vector similarity
[[0.9627094]] final positional encoding similarity
```

**Bold:** Indicates a new term, an important word, or words that you see on the screen.

For instance, words in menus or dialog boxes also appear in the text like this. For example:

“In our case, we are looking for **t5-large**, a t5-large model we can smoothly run in

Google Colaboratory.”



Warnings or important notes appear like this.



Tips and tricks appear like this.

## Get in touch

Feedback from our readers is always welcome.

**General feedback:** Email [feedback@packtpub.com](mailto:feedback@packtpub.com) and mention the book's title in the subject of your message. If you have questions about any aspect of this book, please email us at [questions@packtpub.com](mailto:questions@packtpub.com).

**Errata:** Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you have found a mistake in this book, we would be grateful if you would report this to us. Please visit <http://www.packtpub.com/submit-errata>, selecting your book, clicking on the Errata Submission Form link, and entering the details.

**Piracy:** If you come across any illegal copies of our works in any form on the internet, we would be grateful if you would provide us with the location address or website name. Please contact us at [copyright@packtpub.com](mailto:copyright@packtpub.com) with a link to the material.

**If you are interested in becoming an author:** If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, please visit <http://authors.packtpub.com>.

## Share your thoughts

Once you've read *Transformers for Natural Language Processing and Computer Vision - Third Edition*, we'd love to hear your thoughts! Please

[click here to go straight to the Amazon review](#) page  
for this book and share your feedback.

Your review is important to us and the tech community and will help us  
make sure we're delivering excellent quality content.

# Download a free PDF copy of this book

Thanks for purchasing this book!

Do you like to read on the go but are unable to carry your print books everywhere?

Is your eBook purchase not compatible with the device of your choice?

Don't worry, now with every Packt book you get a DRM-free PDF version of that book at no cost.

Read anywhere, any place, on any device. Search, copy, and paste code from your favorite technical books directly into your application.

The perks don't stop there, you can get exclusive access to discounts, newsletters, and great free content in your inbox daily

Follow these simple steps to get the benefits:

1. Scan the QR code or visit the link below



<https://packt.link/free-ebook/9781805128724>

2. Submit your proof of purchase
3. That's it! We'll send your free PDF and other benefits to your email directly

[OceanofPDF.com](http://OceanofPDF.com)

# 1

## What Are Transformers?

Transformers are industrialized, homogenized **Large Language Models (LLMs)** designed for parallel computing. A transformer model can carry out a wide range of tasks with no fine-tuning. Transformers can perform self-supervised learning on billions of records of raw unlabeled data with billions of parameters. From these billion-parameter models emerged multimodal architectures that can process text, images, audio, and videos.

ChatGPT popularized the usage of transformer architectures that have become general-purpose technologies like printing, electricity, and computers.

Applications are burgeoning everywhere! Google Cloud AI, **Amazon Web Services (AWS)**, Microsoft Azure, OpenAI, Google Workspace, Microsoft 365, Google Colab Copilot, GitHub Copilot, Hugging Face, Meta, and myriad other offers are emerging.

The functionality of transformer models has pervaded every aspect of our workspaces with Generative AI for text, Generative AI for images, discriminative AI, task specific-models, unsupervised learning, supervised learning, prompt design, prompt engineering, text-to-code, code-to-text, and more. Sometimes, a GPT-like model will encompass all these concepts!

The societal impact is tremendous. Developing an application has become an educational exercise in many cases. A project manager can now go to

OpenAI's cloud platform, sign up, obtain an API key, and get to work in a few minutes. Users can then enter a text, specify the NLP task as Google Workspace or Microsoft 365, and obtain a response created by a Google Vertex AI or a ChatGPT transformer model. Finally, users can go to Google's Gen App Builder and build applications without programming or machine learning knowledge.

The numbers are dizzying. *Bommasani et al.* (2023) created a Foundation Model ecosystem that lists 128 Foundation Models 70 applications, and 64 datasets. The paper also mentions Hugging Face's 150,000+ models and 20,000+ datasets! The list is growing weekly and will spread to every activity in society.

Where does that leave an AI professional or someone wanting to be one?

Should a project manager choose to work locally? Or should the implementation be done directly on Google Cloud, Microsoft Azure, or AWS? Should a development team select Hugging Face, Google Trax, OpenAI, or AllenNLP? Should an AI professional use an API with practically no AI development? Should an end-user build a no-code AI application with no ML knowledge with Google's Gen App Builder?

The answer is yes to *all* of the above! You do not know what a future employer, customer, or user may want or specify. Therefore, you must be ready to adapt to any need that comes up at the dataset, model, and application levels. This book does not describe all the offers that exist on the market. You cannot learn every single model and platform on the market. If you try to learn everything, you'll remember nothing. You need to know where to start and when to stop. By the book's end, you will have acquired enough critical knowledge to adapt to this ever-moving market.

In this chapter, we will first unveil the incredible power of the deceptively simple  $O(1)$  time complexity of transformer models that changed everything. We will build a notebook in Python, PyTorch, and TensorFlow to see how transformers hijacked hardware accelerators. We will then discover how one token (a minimal part of a word) led to the AI revolution we are experiencing.

We will continue to discover how a hardly known transformer algorithm in 2017 rose to dominate so many domains. We had to find a new name for it: the Foundation Model. Foundation Models can do nearly everything in AI! We will sit back and watch how ChatGPT explains, analyzes, writes a classification program in a Python notebook, and displays a decision tree.

Finally, this chapter introduces the role of an AI professional in the ever-changing job market. We will begin to tackle the problem of choosing the right resources.

We must address these critical notions before starting our exploratory journey through the variety of transformer model implementations described in this book.

This chapter covers the following topics:

- How one  $O(1)$  invention changed the course of AI history
- How transformer models hijacked hardware accelerators
- How one token overthrew hundreds of AI applications
- The multiple facets of a transformer model
- Generative AI versus discriminative AI
- Unsupervised and self-supervised learning versus supervised learning
- General-purpose models versus task-specific models
- How ChatGPT has changed the meaning of automation

- Watch ChatGPT create and document a classification program
- The role of AI professionals
- Seamless transformer APIs
- Choosing a transformer model



With all the innovations and library updates in this cutting-edge field, packages and models change regularly. Please go to the GitHub repository for the latest installation and code examples:

<https://github.com/Denis2054/Transformers-for-NLP-and-Computer-Vision-3rd-Edition/tree/main/Chapter01>.

You can also post a message in our Discord community (<https://www.packt.link/Transformers>) if you have any trouble running the code in this or any chapter.

Our first step will be to explore the seeds of the disruptive nature of transformers.

## How constant time complexity O(1) changed our lives forever

How could this deceptively simple O(1) time complexity class forever change AI and our everyday lives? How could O(1) explain the profound architectural changes that made ChatGPT so powerful and stunned the

world? How can something as simple as  $O(1)$  allow systems like ChatGPT to spread to every domain and hundreds of tasks?

The answer to these questions is the only way to find your way in the growing maze of transformer datasets, models, and applications is to *focus* on the underlying concepts of thousands of assets. Those concepts will take you to the core of the functionality you need for your projects.

This section will provide a significant answer to those questions before we move on to see how one token (a minimal piece of a word) started an AI revolution that is raging around the world, triggering automation never seen before.

We need to get to the bottom of the chaos and disruption generated by transformers.

To achieve that goal, in this section, we will use science and technology to understand how all of this started. First, we will examine  $O(1)$  and then the complexity of a layer through a Python and PyTorch notebook.

Let's first get the core concepts and terminology straight for  $O(1)$ .

## **O(1) attention conquers O(n) recurrent methods**

$O(1)$  is a “Big O” notation. “Big O” means “order of.” In this case,  $O(1)$  represents a constant time complexity. We can say  $O(1)$  and order of 1 operation.

Believe it or not, you’re at the heart of the revolution!

In *Chapter 2, Getting Started with the Architecture of the Transformer Model*, we will explore the architecture of the transformer model.

In this section and chapter, we will first focus on what led to the industrialization of AI through transformers and the ChatGPT frenzy: *the exponential increase of hardware efficiency with self-attention*. We will discover how attention leverages hardware and opens the door to incredible new ways of machine learning.

The following sentence contains 11 words:

Jay likes oranges in the morning but not in the evening.

The length of the sentence is  $n = 11$ .

The problem of language understanding can be summed up in one word: context. A word can rarely be defined without context. The meaning of a word can change in each context beyond its definition in a dictionary.

Let's begin with a conceptual approach of an attention layer.

## Attention layer

If we look at what Jay likes, the dimensions (or parameters) of the word “oranges” contain several relationships:

Dimension 1: The association between oranges and evening

Dimension 2: The association between oranges and morning

Dimension 3: The association between Jay and oranges

Dimension 4: The association between Jay and morning

...

Dimension z

Notice that the relationships are defined pairwise: one word to one word. This is how *self-attention* works in a transformer.

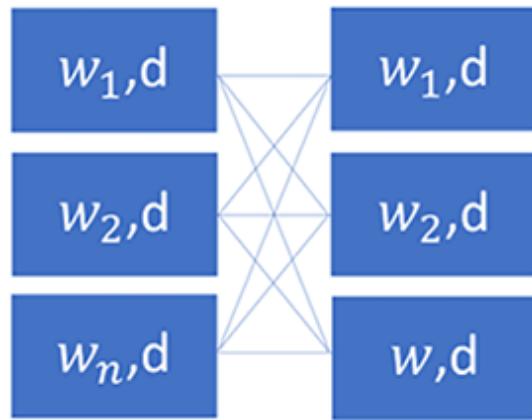
If we represent this with Big O notation, we get  $O(1)$ . “O” means “in the order of.” For example,  $O(1)$  is an “order of 1” or constant time complexity class.

We perform one operation per word,  $O(1)$ , for each word to find the relationship with another word in a pairwise analysis.

Translated into numbers:

- $n$  = the length of the sequence, which is 11 words in this case.
- $d$  = the number of dimensions expressed in floats. In machine learning, the dimensions are expressed in floats. For example, if  $x$  is a word, the values might be: `[-0.2333, 03.8559, 0.9844...0394]`. The model will learn these values through billions of text data points.

This pairwise relationship is an  $n * n$  computation, as shown in *Figure 1.1*:



*Figure 1.1: Word-to-word relationships*

$O(1)$  represents the memory complexity in this case.

The computational complexity of an attention layer is thus  $O(n^2 * d)$ .  $n^2$  is the pairwise (word-to-word) operation of the whole sequence  $n$ .  $d$  represents the dimensions the model is learning.

Let's see the difference with a recurrent layer.

## Recurrent layer

Recurrent layers do not function that way. They are  $O(n)$ , an order of  $n$  linear time complexity. The longer the sequence, the more memory they will consume. Why? They do not learn the dimensions with pairwise relationships. They learn in a sequence. For example:

Dimension a: Jay

Dimension b: likes and Jay

Dimension c: oranges and likes and Jay

Dimension d: in and oranges and likes and Jay

...

Dimension z

You can see that each word does not simply look at another word but several other words simultaneously!

The number of dimensions  $d$  one word is multiplied by the dimensions of a preceding word leading to  $d^2$  to learn the dimensions. The computational time complexity of a recurrent layer is thus:

$$O(n*d^2)$$

Let us now see the magic.

## The magic of the computational time complexity of an attention layer

An attention layer benefits from its  $O(1)$  memory time complexity. This enables the computational time complexity of  $O(n^2 * d)$  to perform a dot product between each word. In machine learning, this transcribes into multiplying the representation  $d$  of each word by another word. An attention layer can thus learn all the relationships in one matrix multiplication!

A recurrent layer of a computational time complexity of  $O(n * d^2)$  is hindered by its  $O(n)$  linear, sequential process. Performing the same task as the attention layer will take more operations. We will now simulate the attention time complexity model and the recurrent time complexity model in a Python, PyTorch, and TensorFlow notebook.

The calculations are conceptual. We will run simulations with a CPU, a GPU, and a TPU in this section:

- **CPU: Central Processing Unit.** This is the primary component of a computer.
- **GPU: Graphics Processing Unit.** A GPU is a specialized processing unit. Initially used for 3D image rendering, GPUs have evolved to perform machine learning tasks, such as matrix multiplications.
- **TPU: Tensor Processing Unit.** A TPU is an accelerating machine learning workload processor developed by Google and optimized for TensorFlow.

Open `0-1_and_Accelerators.ipynb` in the chapter directory of the repository.

We will begin with CPUs.

## Computational time complexity with a CPU

A CPU is a general-purpose software processor. It is not specially designed for matrix multiplications. It can perform complex operations but only to a certain degree of efficiency.

Before running the first cell of the notebook, we must verify that we are using a CPU by going to the Runtime menu, selecting **Change runtime type**, and making sure that the **Hardware accelerator** parameter is set **CPU**:

## Change runtime type

### Runtime type

Python 3 ▾

### Hardware accelerator [?](#)

- CPU
- A100 GPU
- V100 GPU
- T4 GPU
- TPU

Cancel

Save

Figure 1.2: Selecting a Runtime type

In this case, a hardware accelerator is a GPU or TPU that performs specific computing tasks (matrix multiplication, for example) more efficiently than a general-purpose CPU, which is not a hardware accelerator. This explains why the **Hardware accelerator** is set to **None**.

The notebook begins with a figure containing the time complexities we have just gone through:

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	Recurrent $O(n \cdot d^2)$	$O(n)$	$O(n)$

*Figure 1.3: Time complexities from Attention is All You Need, Vaswani et al. (2017), page 6*

The goal of the notebook is to represent the complexity per layer of self-attention and recurrent layers, not the actual algorithms. The self-attention computational time complexity will run with matrix multiplications.

The recurrent computational time will use a loop simulating its sequential approach. A recurrent model uses matrices to compute values; however, its overall computational complexity is an order of  $n$  and  $d$ . It is important to note that the calculations do not reflect the actual algorithms inside self-attention layers and recurrent layers. Also, many other factors will influence performance: hardware, data, hyperparameters, training time, and other parameters.

The goal of this notebook remains to show the overall computational time complexity. In this case, the simulations are sufficient. We first define the framework of the evaluations:

```
# Computational times of complexity per layer
# Comparing the computational time between:
# self attention = O(n^2 * d)
#and
# recurrent = O(n * d^2)
```

We then import `numpy` and `time`:

```
import numpy as np
import time
```

We first define the sequence length (number of words in a sequence) and the dimensions (numerical vector representing features of the word):

```
# define the sequence length and representation dimensionality
n = 512
d = 512
```

You will notice that  $n = d$ , which means that  $O(n^2 * d) = O(n * d^2) = 512 * 512 * 512 = 134,217,728$  operations. In this case, both the attention and the recurrent layers have the same number of operations to perform from a computational complexity time perspective. We will define the sequence of input with random values for the dimensions ( $d$ ):

```
# define the inputs
```

We will first simulate the time complexity of the self-attention layer with a matrix multiplication with start time:

```
# simulation of self-attention layer O(n^2*d)
start_time = time.time()
for i in range(n):
    for j in range(n):
        _ = np.dot(input_seq[i], input_seq[j])
```

When the matrix multiplication is finished, we calculate the time it took and print it:

```
at=time.time()-start_time
print(f"Self-attention computation time: {time.time() - start_
```

The time is displayed:

```
Self-attention computation time: 0.7938594818115234 seconds
```

We now run the simulation of the time complexity of the recurrent layer function without a matrix multiplication with a start and end time:

```
# simulation of recurrent layer O(n*d^2)
start_time = time.time()
hidden_state = np.zeros(d)
for i in range(n):
    for j in range(d):
        for k in range(d):
            hidden_state[j] += input_seq[i, j] * hidden_state[rt=time.time()-start_time
print(f"Recurrent layer computation time: {time.time() - start_
```

The output shows the time it took:

```
Recurrent layer computation time: 109.65185356140137 seconds
```

We now calculate the percentage of the attention layer's computational time complexity and the recurrent layer's computational time complexity. We can measure the performance percentages. This approach provides an overall idea of the power of attention layers. We then display the result:

```
# Calculate the total  
# Calculate the percentage of at  
    percentage_at = round((at / total) * 100,2)  
# Output the result  
print(f"The percentage of 'computational time for attention' i
```

The output shows that the attention layer's computational time complexity is more efficient:

```
The percentage of 'computational time for attention' in the sum
```

The attention layer's computational time complexity performs better than the recurrent layer on a CPU in this configuration.

Let's move on to GPUs.

## Computational time complexity with a GPU

Before running the first cell of the notebook, we must verify that we are using a GPU by going to the Runtime menu, selecting **Change runtime type**, and making sure that the **Hardware accelerator** parameter is set to one of the **GPUs**:

## Change runtime type

### Runtime type

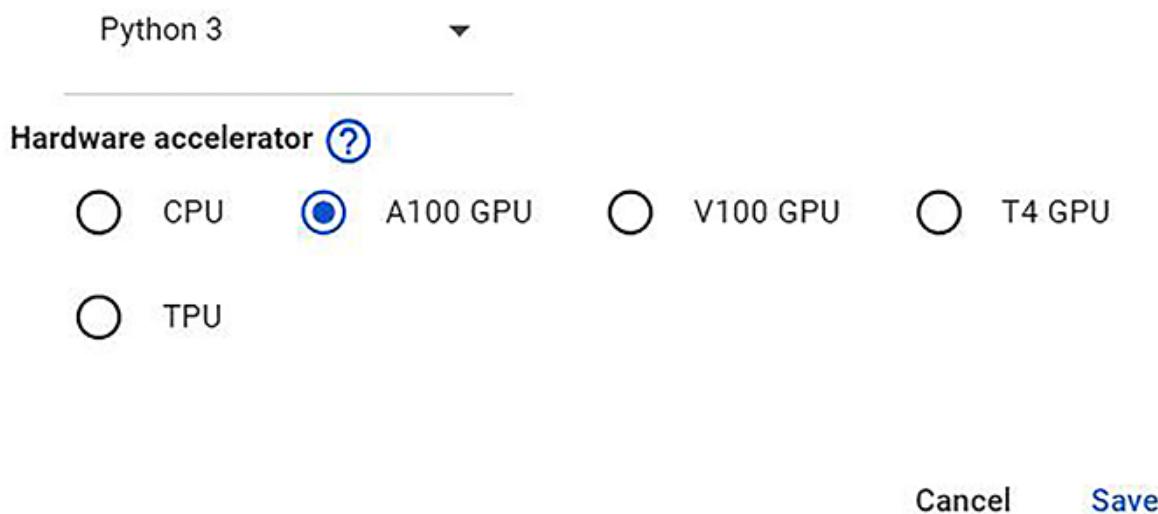


Figure 1.4: Changing the settings to use a GPU

The GPU type is, in this case, an NVIDIA V100. GPUs work well with algorithms requiring massive operations, particularly matrix multiplications, such as the attention layer mechanism's computational time complexity we are simulating in this notebook.

We will implement our simulation in PyTorch to leverage the power of the GPUs:

```
# PyTorch version
import torch
import time
```

We will use the same parameters as for the CPU evaluation:

```
# define the sequence length and representation dimensionality
n = 512
d = 512
```

We now activate the GPU and define the inputs:

```
# Use GPU if available, otherwise stick with cpu
device = torch.device("cuda" if torch.cuda.is_available() else
print(device)
# define the inputs
input_seq = torch.rand(n, d, device=device)
```

We run the same simulation as for the CPU but on the GPU:

```
# simulation of self-attention layer O(n^2*d)
start_time = time.time()
_ = torch.mm(input_seq, input_seq.t())
at = time.time() - start_time
print(f"Self-attention computation time: {at} seconds")
```

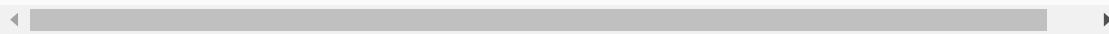
The output is not that impressive because we didn't run a massive number of matrix multiplications to take full advantage of the GPU, which reveals its power on larger volumes:

```
cuda
Self-attention computation time: 2.887202501296997 seconds
```

We will now run the recurrent layer function but limit its time to about 10 times (depending on GPU activity) the time it took for self-attention, which

is enough to show our point with `if ct>at*10:`

```
# simulation of recurrent layer O(n*d^2)
start_time = time.time()
hidden_state = torch.zeros(d, device=device)
for i in range(n):
    for j in range(d):
        for k in range(d):
            hidden_state[j] += input_seq[i, j] * hidden_state[k]
            ct = time.time() - start_time
            if ct>at*10:
                break
```



We compute the limited time we gave the function to run and display it:

```
rt = time.time() - start_time
print(f'Recurrent layer computation time: {rt} seconds')
```

The output time isn't very efficient:

```
Recurrent layer computation time: 36.3216814994812 seconds
```

We calculate the percentage of attention layer in the total time:

```
# Calculate the total
total = at + rt
# Calculate the percentage of at
percentage_at = round((at / total) * 100, 2)
# Output the result
print(f'The percentage of self-attention computation in the su
```



The percentage of self-attention computation in the sum of self

◀ ▶

We can check the information on the GPU:

```
!nvidia-smi
```

The output will display the information on the current GPU:

```
+-----+
| NVIDIA-SMI 525.105.17    Driver Version: 525.105.17    CUDA Version: 12.0 |
+-----+
| GPU  Name      Persistence-M| Bus-Id      Disp.A  | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap|          Memory-Usage | GPU-Util  Compute M.  |
|                               |               MIG M.   |
+-----+
| 0  Tesla V100-SXM2... Off  | 00000000:00:04.0 Off |          0 |
| N/A   39C     P0    40W / 300W | 2160MiB / 16384MiB |     0%      Default |
|                               |                           N/A |
+-----+
+-----+
| Processes:
| GPU  GI  CI          PID  Type  Process name                  GPU Memory |
|           ID  ID                                     Usage
+-----+
```

Figure 1.5: GPU information

We will finally explore a TPU simulation.

## Computational time complexity with a TPU

Google designed Cloud TPUs specifically for matrix calculations and neural networks. Since the primary task of a TPU is matrix multiplication, it

doesn't make much sense to run sequential operations. We will limit the recurrent layer to 10 times the time it takes for the attention layer.

Before running the first cell of the notebook, we must verify that we are using a TPU by going to the Runtime menu, selecting **Change runtime type**, and making sure that the **Hardware accelerator** parameter is set to **TPU**:

## Change runtime type

### Runtime type

Python 3



### Hardware accelerator [?](#)

- CPU
- A100 GPU
- V100 GPU
- T4 GPU
- TPU

[Cancel](#)

[Save](#)

*Figure 1.6: Changing the notebook settings to use a TPU*

For this simulation, we will use TensorFlow to take full advantage of the TPU:

```
import tensorflow as tf
import numpy as np
import time
```

The program is the same one as for the CPU and GPU evaluation, except that this time, it's running on a TPU with TensorFlow:

```
# define the sequence length and representation dimensionality
n = 512
d = 512
# define the inputs
input_seq = tf.random.normal((n, d), dtype=tf.float32)
```

We run the matrix multiplication and measure the time it took:

```
# simulation of self-attention layer O(n^2*d)
start_time = time.time()
_ = tf.matmul(input_seq, input_seq, transpose_b=True)
at = time.time() - start_time
print(f"Self-attention computation time: {at} seconds")
```

The output is efficient:

```
Self-attention computation time: 0.10626077651977539 seconds
```

We will now run the recurrent layer:

```
# simulation of recurrent layer O(n*d^2)
start_time = time.time()
hidden_state = np.zeros((n, d), dtype=np.float32)
for i in range(n):
    for j in range(d):
        for k in range(d):
            hidden_state[i, j] += input_seq[i, j].numpy() * hi
ct = time.time() - start_time
```

```
if ct>at*10:  
    break
```

The output is not efficient, which is normal:

```
rt = time.time() - start_time  
print(f'Recurrent layer computation time: {rt} seconds')
```

```
Recurrent layer computation time: 66.53181290626526 seconds
```

We now compute the percentage of the total time:

```
# Calculate the total  
total = at + rt  
# Calculate the percentage of at  
percentage_at = round((at / total) * 100, 2)  
# Output the result  
print(f'The percentage of self-attention computation in the sum of self-attention and recurrent layers is {percentage_at}%')
```

The output confirms that attention layers are more efficient:

```
The percentage of self-attention computation in the sum of self-attention and recurrent layers is 99%
```

We will now take the TPU simulation into the world of LLMs.

## TPU-LLM

We will run the same simulation as for the TPU with relatively standard parameters for the attention layer. `d=12288` could be the dimensionality of an LLM. `n=32728` could be the input limit of an LLM. These values are simply an example of running the recurrent layer on the TPU. The values are those of an LLM:

```
import tensorflow as tf
import numpy as np
import time
# define the sequence length and representation dimensionality
n = 32768
d = 12288
# define the inputs
input_seq = tf.random.normal((n, d), dtype=tf.float32)
```

We will run the same function as before and display the time it took:

```
# simulation of self-attention layer O(n^2*d)
start_time = time.time()
_ = tf.matmul(input_seq, input_seq, transpose_b=True)
at = time.time() - start_time
print(f"Self-attention computation time: {at} seconds")
```

The time is very efficient:

```
Self-attention computation time: 23.117244005203247 seconds
```

We can display some information on the TPU:

```
import os
from tensorflow.python.profiler import profiler_client
```

```
tpu_profile_service_address = os.environ['COLAB_TPU_ADDR'].rep  
print(profiler_client.monitor(tpu_profile_service_address, 100
```

We can see that the TPU Matrix Units were hardly solicited:

```
Timestamp: 10:20:30  
TPU type: TPU v2  
Utilization of TPU Matrix Units (higher is better): 0.000%
```

The conclusion of this evaluation can be summed up in the following points:

- The attention layer computational time complexity outperforms the recurrent layer computational time complexity by avoiding recurrence.
- The attention layer's one-to-one word analysis makes it better at detecting long-term dependencies.
- The architecture of the attention layer enables matrix multiplication, which takes full advantage of modern GPUs and TPUs.
- By unleashing the power of GPUs (and TPUs) and attention, algorithms can process more data and learn more information.

The deceptively simple  $O(1)$  that led to  $O(n^2 * d)$  has changed our daily lives forever by taking AI to the next level.

How did this happen?

## A brief journey from recurrent to attention

For decades, **Recurrent Neural Networks (RNNs)**, including LSTMs, have applied neural networks to NLP sequence models. However, using recurrent

functionality reaches its limit when faced with long sequences and large numbers of parameters. Thus, state-of-the-art transformer models now prevail.

This section goes through a brief background of NLP that led to transformers, which we'll describe in more detail in *Chapter 2, Getting Started with the Architecture of the Transformer Model*. First, however, let's have a look at the attention head of a transformer that has replaced the RNN layers of an NLP neural network.

The core concept of a transformer can be summed up loosely as “mixing tokens.” NLP models first convert word sequences into tokens. RNNs analyze tokens in recurrent functions. Transformers do not analyze tokens in sequences but relate every token to the other tokens in a sequence, as shown in *Figure 1.7*:

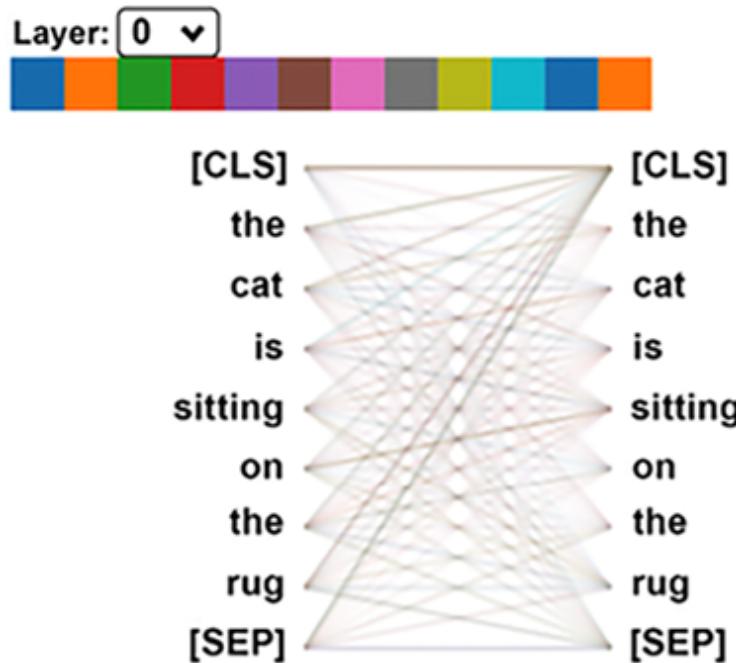


Figure 1.7: An attention head of a layer of a transformer

We will go through the details of an attention head in *Chapter 2*. For the moment, the takeaway of *Figure 1.7* is that for innovative LLMs, *each word (token) of a sequence is related to all the other words of a sequence.*

Let's briefly go through the background of transformers.

## A brief history

Many great minds have worked on sequence patterns and language modeling. As a result, machines progressively learned how to predict probable sequences of words. It would take a whole book to cite all the giants that made this happen.

In this section, I will share some of my favorite researchers, among the many who paved the way to where AI is today, laying the ground for the arrival of the transformer.

In the late 19th and early 20th century, Andrey Markov introduced the concept of random values and created a theory of stochastic processes (see the *Further reading* section). We know them in AI as the **Markov Decision Process (MDP)**, **Markov chains**, and **Markov processes**. In the early 20th century, Markov showed that we could predict the next element of a chain, a sequence, using only the most recent elements of that chain. He applied this to letters and conducted many experiments. double space Bear in mind that he had no computer but proved a theory still in use today in AI.

In 1948, Claude Shannon's *The Mathematical Theory of Communication* was published. Claude Shannon laid the grounds for a communication model based on a source encoder, transmitter, and receiver or semantic decoder. He created information theory as we know it today. Claude Shannon, of course, mentions Markov's theories.

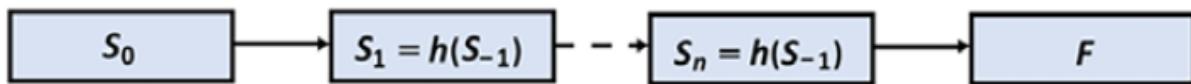
In 1950, Alan Turing published his seminal article: *Computing Machinery and Intelligence*. Alan Turing based this article on machine intelligence on the successful Turing machine, which decrypted German messages during World War II. The messages consisted of sequences of words and numbers.

In 1954, the Georgetown-IBM experiment used computers to translate Russian sentences into English using a rule system. A rule system is a program that runs a list of rules that will analyze data structures. Rule systems still exist and are everywhere. However, in some cases, machine intelligence can replace rule lists for the billions of language combinations by automatically learning the patterns.

The expression *AI* was first used by John McCarthy in 1956 when exploring the possibilities that machines could learn.

In 1982, John Hopfield introduced an **RNN** known as Hopfield networks or “associative” neural networks. John Hopfield was inspired by W. A. Little, who wrote *The Existence of Persistent States in the Brain* in 1974, which laid the theoretical grounds for learning processes for decades. RNNs evolved, and LSTMs emerged as we know them today.

An RNN memorizes the persistent states of a sequence efficiently, as shown in *Figure 1.8*:



*Figure 1.8: The RNN process*

Each state  $S_n$  captures the information of  $S_{n-1}$ . When the network’s end is reached, function  $F$  will perform an action: transduction, modeling, or any other type of sequence-based task.

In the 1980s, Yann LeCun designed the multipurpose **Convolutional Neural Network (CNN)**. He applied CNNs to text sequences, and they also apply to sequence transduction and modeling. They are also based on W. A. Little's persistent states that process information layer by layer. In the 1990s, summing up several years of work, Yann LeCun produced LeNet-5, which led to the many CNN models we know today. However, a CNN's otherwise efficient architecture faces problems when dealing with long-term dependencies in lengthy and complex sequences.

We could mention many other great names, papers, and models that would humble any AI professional. Everybody in AI seemed to be on the right track all these years. Markov fields, RNNs, and CNNs evolved into multiple other models. The notion of attention appeared: peeking at other tokens in a sequence, not just the last one. It was added to the RNN and CNN models.

After that, if AI models needed to analyze longer sequences requiring increasing computer power, AI developers used more powerful machines and found ways to optimize gradients.

Some research was done on sequence-to-sequence models, but they did not meet expectations. It seemed that nothing else could be done to make more progress, and 30 years passed this way. Attention was introduced, but the models still relied on recurrence. And then, starting in late 2017, the industrialized state-of-the-art Transformer came with its attention head sublayers and more. RNNs did not appear as a prerequisite for sequence modeling anymore.

Before diving into the Original Transformer's architecture, which we will do in *Chapter 2, Getting Started with the Architecture of the Transformer Model*, let's start at a high level by examining the paradigm change triggered by one little token.

# From one token to an AI revolution

Yes, the title is correct, as you will see in this section. One token produced an AI revolution and has opened the door to AI in every domain and application.

ChatGPT with GPT-4, PaLM 2, and other LLMs have a unique way of producing text.

In LLMs, a token is a minimal word part. *The token is where a Large Language Model starts and ends.*

For example, the word `including` could become: `includ` + `ing`, representing two tokens. GPT models predict tokens based on the hundreds of billions of tokens in its training dataset. Examine the graph in *Figure 1.9* of an OpenAI GPT model that is making an inference to produce a token:

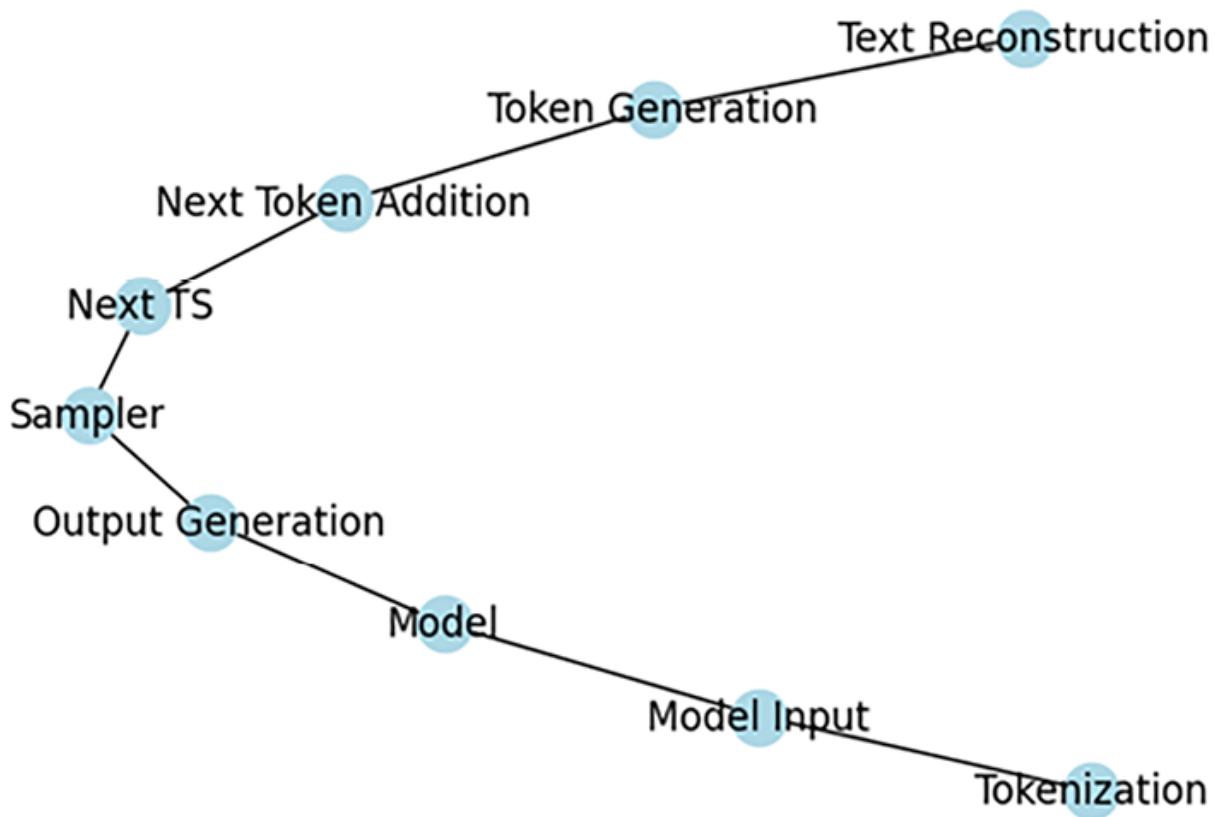


Figure 1.9: GPT inference graph built in Python with NetworkX

It may come as a surprise, but the only parts of this figure controlled by the model are `Model` and `Output Generation!`, which produce raw logits. All the rest is in the pipeline.

To understand the pipeline, we will first go through the description of these steps:

1. **Tokenization:** The pipeline converts the input sequence “Is the sun in the solar system?”, for example, into tokens using a tokenizer. In *Chapter 10, Investigating the Role of Tokenizers in Shaping Transformer Models*, we will dig into the critical role of tokenizers in transformer models.

2. **Model Input:** The pipeline now passes the tokenized sequence into the trained GPT model.
3. **Model:** The model processes the input through its various layers, from the input layer through multiple transformer layers to the output layer. *Chapter 2, Getting Started with the Architecture of the Transformer Model*, will describe the architecture in detail.
4. **Output Generation:** The model produces raw output logits given the input sequence.
5. **Sampler:** The sampler will convert the logits into probabilities. We will implement hyperparameters that influence the model's output in *Chapter 7, The Generative AI Revolution with ChatGPT*. For more on the sampling process involving the main hyperparameters, see the *Vertex AI PaLM 2 Interface* section of *Chapter 14, Exploring Cutting-Edge LLMs with Vertex AI and PaLM 2*.
6. **Next Token Selection (Next TS):** The next token is selected based on the probabilities of the sampler.
7. **Next Token Addition:** The selected next token is added to the input sequence (in token form) and repeats the process from step 3 until the maximum token limit has been reached.
8. **Token Generation Completion (Text Generation):** Text generation will end once the maximum token limit has been reached or an end-of-sequence token has been detected.
9. **Text Reconstruction:** The tokenizer converts the final sequence of tokens back into a string of text. This step includes stitching any subword tokens back together to form whole words.

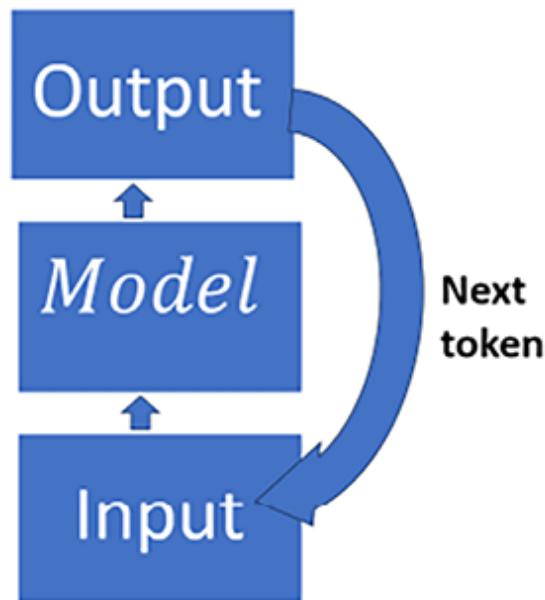
The magic appears again! The model takes an input sequence and produces one token.

That token is added to the sequence, and the model starts over again to produce another token.

We can sum up the steps of this process:

1. Input = an input sequence expressed in tokens.
2. The model processes the input.
3. Output = the next token is added at the end of the input if the maximum tokens requested is reached.

The process is a one-token output revolution, as shown in *Figure 1.10*:



*Figure 1.10: The model processes the input sequence, produces raw outputs, and selects one token*

To sum this up in a notation for one forward pass:

$$t = f(n)$$

Take the necessary time to let this sink in:

- $t$  = one token.
- $f$  = the model and the output controller that infers the next token.

- $n$  = the initial sequence of tokens + each new next token that is added to it until the maximum number of tokens is reached or an end-of-sequence token is detected.

We have summed up the process as  $f(n)$ . Under the hood of  $f(n)$  is the whole process of adding predicted tokens one by one. You can think of  $f(n)$  as:

$$t_i = f(t_1, t_2, \dots, t_{\{i-1\}})$$

Think about it. One token! And yet that token changed society in every domain: translations, summarization, question-answering, imagery, video... the list is growing daily.

We will now go from one token to everything!

## From one token to everything

In the previous section, we saw that a generative model can be summed up as  $t = f(n)$ . The model processes  $n$  and produces  $t$ , one token at a time.

The consequences are unfathomable and tremendous.

We can see that the one-token approach, summed up as  $t = f(n)$ , results in:

- The **dynamic** nature of a transformer model. It adapts its outputs based on its incremental inputs represented by:

$$t_i = f(t_1, t_2, \dots, t_{\{i-1\}})$$

- The model will adapt and generate outputs on completely new inputs!
- The **implicit** nature of a transformer model. The model encodes and stores relationships between tokens in weights and biases. There is no explicit guideline. It just keeps producing tokens based on its dynamic inputs based on millions of text, image, and audio data!

- The incredible **flexibility** of the system is inherited from its dynamic and implicit properties. A GPT-like transformer model can infer meaningful outputs for a wide range of contextually diverse inputs.

We will now see the implications of the models' one-token approach through different perspectives:

- **Supervised and unsupervised**

Some may say that a GPT series model such as ChatGPT goes through unsupervised training. That statement is only true to a certain extent. Token by token, a GPT-like model finds its way to accuracy through self-supervised learning, predicting each subsequent token based on the preceding ones in the sequence. It succeeds in doing so through the influence of all the other tokens' representations in a sequence.

We can also fine-tune a GPT model with an input (prompt) and output (completion) with labels! We can provide thousands of inputs (prompts) with one token as an output (completion). For example, we can create thousands of questions as inputs with only `true` and `false` as outputs. This is implicit supervised learning. Also, the model will not explicitly memorize the correct predictions. It will simply learn the patterns of the tokens.

- **Discriminative and generative**

Since the arrival of ChatGPT, the term “Generative AI” has often been used to describe GPT-like models. First, we must note that GPT models were not the first to produce Generative AI recurrent neural networks were also generative.

Also, if we follow the science, saying the GPT-like models perform Generative AI tasks is not entirely true. In the previous *Supervised and*

*unsupervised* bullet point, we mentioned that a “generative” model could perform the supervised task of inferring an output such as “true” or “false” when the input is a statement or a question. This is not a Generative AI task! This is a discriminative AI task.

Another situation can arise with summarizing. Part of summarizing is discriminative to find keywords, topics, and names. Another part can be generative when the system infers new tokens to summarize a text.

Finally, the full power of the autoregressive nature (generating a token and adding it to the input) of language modeling (adding a new token at the end of a sequence) is attained when an LLM invents a story; that is generative AI.

A GPT-like model can be discriminative, Generative, or both depending on the task to perform.

- **Task-specific and specific task models**

Task-specific models are often opposed to general-purpose models, such as GPT-like LLMs. Task-specific models are trained to perform very well on specific tasks. This is accurate to some extent for some tasks. However, LLMs trained to be general-purpose systems can surprisingly perform well on specific tasks, as proven through the evaluations of various exams (law, medical). Also, transformer models can now perform image recognition and generation.

- **Interaction and automation**

A key point in analyzing the relationship between humans and LLMs is automation. To what degree should we automate our tasks? Some tasks cannot be automated by LLMs. ChatGPT and other assistants cannot

make life-and-death, moral, ethical, and business decisions. Asking a system to do that will endanger an organization.

Automation needs to be thoroughly analyzed before implementing LLMs.

In the following chapters of this book, you will encounter supervised, self-supervised, unsupervised, discriminative, generative, task-specific, specific tasks, and interactions with evolved AI and automated systems. Sometimes you will find all of this in the architecture of a general-purpose LLM!

Be prepared and remain flexible!

## Foundation Models

Advanced large multipurpose transformer models represent such a paradigm change that they require a new name to describe them: **Foundation Models**. Accordingly, Stanford University created the **Center for Research on Foundation Models (CRFM)**. In August 2021, the CRFM published a two-hundred-page paper (see the *References* section) written by over one hundred scientists and professionals: *On the Opportunities and Risks of Foundation Models*.

Foundation Models were not created by academia but by the big tech industry. Google invented the transformer model, leading to Google BERT, LaMBDA, PaLM 2, and more. Microsoft partnered with OpenAI to produce ChatGPT with GPT-4, and soon more.

Big tech had to find a better model to face the exponential increase of petabytes of data flowing into their data centers. Transformers were thus born out of necessity.

Let's consider the evolution of LLMs to understand the need for industrialized AI models.

Transformers have two distinct features: a high level of homogenization and mind-blowing emergence properties. *Homogenization* makes it possible to use one model to perform a wide variety of tasks. These abilities *emerge* through training billion-parameter models on supercomputers.

The present ecosystem of transformer models is unlike any other evolution in AI and can be summed up with four properties:

- **Model architecture:** The model is industrial. The layers of the model are identical, and they are specifically designed for parallel processing. We will go through the architecture of transformers in *Chapter 2, Getting Started with the Architecture of the Transformer Model*.
- **Data:** Big tech possesses the hugest data source in the history of humanity, generated mainly through human-driven online activities and interactions, including browsing habits, search queries, social media posts, and online purchases.
- **Computing power:** Big tech possesses computer power never seen before at that scale. For example, GPT-3 was trained at about 50 **PetaFLOPS (Floating Point Operations Per Second)**, and Google now has domain-specific supercomputers that exceed 80 PetaFLOPS. In addition, GPT-4, PaLM 2, and other LLMs use thousands of GPUs to train their models.
- **Prompt engineering:** Highly trained transformers can be triggered to do a task with a prompt. The prompt is entered in natural language. However, the words used require some structure, making prompts a meta language.

A Foundation Model is thus a transformer model that has been trained on supercomputers on billions of records of data and billions of parameters. The model can perform a wide range of tasks without further fine-tuning. Thus,

the scale of Foundation Models is unique. These fully trained models are often called engines. GPT-4, Google BERT, PaLM 2, and scores of transformer models can now qualify as Foundation Models. *Bommasani et al.* (2023) created Ecosystem Graphs to keep track of the growing number of assets on the market (datasets, models, and applications).

We will now examine an example of how Foundation Models work and have changed how we develop programs.

## From general purpose to specific tasks

This section shows how transformer-driven AI is changing the way we work forever. ChatGPT will design, write, and explain a classification program automatically. In this chapter, you will find discriminative, generative, task-specific, and general-purpose features in the specific task we will ask ChatGPT to perform.

Don't be mistaken! Transformer-driven AI is in its early stages. As the models evolve and new ones appear, AI will perform increasingly complex tasks.

We will start by watching ChatGPT write and explain a classification algorithm by itself!

Open `ChatGPT_Plus_writes_and_explains_AI.ipynb` in the chapter directory of the GitHub repository of this book.

Two requests were submitted to ChatGPT Plus:

1. Provide a `sklearn` classification of `Iris` with some kind of `matplotlib` graph to describe the results, and don't use `OpenAI`

APIs.

2. Write a detailed explanation for a beginner in Python for this code.

ChatGPT, GPT-3.5-turbo wrote this entire notebook with the explanations!

The only human intervention was to create the notebook and copy the text and code provided by ChatGPT. ChatGPT begins by explaining the program:

```
This program is a simple machine learning project that demonstr  
In this program, you will first load the Iris dataset and split  
This program is a great way for beginners to learn about machin  
Overall, this program is a fun and educational way to learn abo
```

ChatGPT explains the basic configuration of a classification program. It's hard to believe initially, but you will get used to it!

```
Install scikit-learn:  
pip install -U scikit-learn  
1. Import libraries  
We import several libraries that we will use throughout the co  
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
from sklearn.datasets import load_iris  
from sklearn.model_selection import train_test_split  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.metrics import confusion_matrix  
from sklearn.tree import plot_tree
```

Now, ChatGPT goes on to load and split the dataset into training and testing sets.

```
2.Load the Iris dataset:  
iris = load_iris()  
We use the load_iris function from Scikit-learn's datasets mod  
3.Split the dataset into training and testing sets  
X_train, X_test, y_train, y_test = train_test_split(iris.data,  
We use Scikit-learn's train_test_split function to split the d
```

At this point, while watching ChatGPT write everything up to now and still typing away furiously, I began to wonder what I was doing there. What was my role in all this? Will AI take over my job? My role in computational linguistics? Although I know the many limitations of Generative AI, it still made me wonder.

Now, ChatGPT configures, displays, and trains a decision tree classifier model:

```
4. Train a decision tree classifier model:  
model = DecisionTreeClassifier(random_state=42)  
model.fit(X_train, y_train)
```

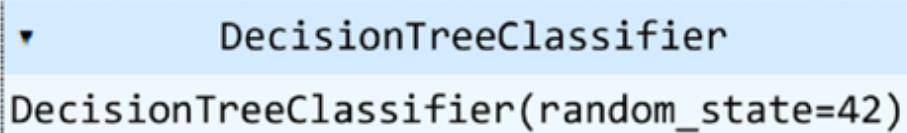


Figure 1.11 ChatGPT creates a decision tree classifier

We create a DecisionTreeClassifier object from Scikit-learn's

Now what? ChatGPT even follows standard machine learning procedures, evaluates the model, and then computes and plots the confusion matrix!

```
5. Evaluate the model on the testing set:  
score = model.score(X_test, y_test)  
print(f"Accuracy: {score:.2f}")  
Accuracy: 1.00  
We use the score method of the trained model to evaluate the accuracy.  
6. Compute the confusion matrix:  
y_pred = model.predict(X_test)  
cm = confusion_matrix(y_test, y_pred)  
7.1. Plot the confusion matrix:  
fig, ax = plt.subplots()  
im = ax.imshow(cm, interpolation="nearest", cmap=plt.cm.Blues)  
ax.figure.colorbar(im, ax=ax)  
ax.set(xticks=np.arange(cm.shape[1]),  
       yticks=np.arange(cm.shape[0]),  
       xticklabels=iris.target_names, yticklabels=iris.target_names,  
       title="Confusion matrix",  
       xlabel="Predicted label",  
       ylabel="True label")  
plt.setp(ax.get_xticklabels(), rotation=45, ha="right", rotation_mode="auto")  
for i in range(cm.shape[0]):  
    for j in range(cm.shape[1]):  
        ax.text(j, i, format(cm[i, j], "d"), ha="center", va="center",  
                color="white" if cm[i, j] > cm.max() / 2 else "black")  
fig.tight_layout()
```

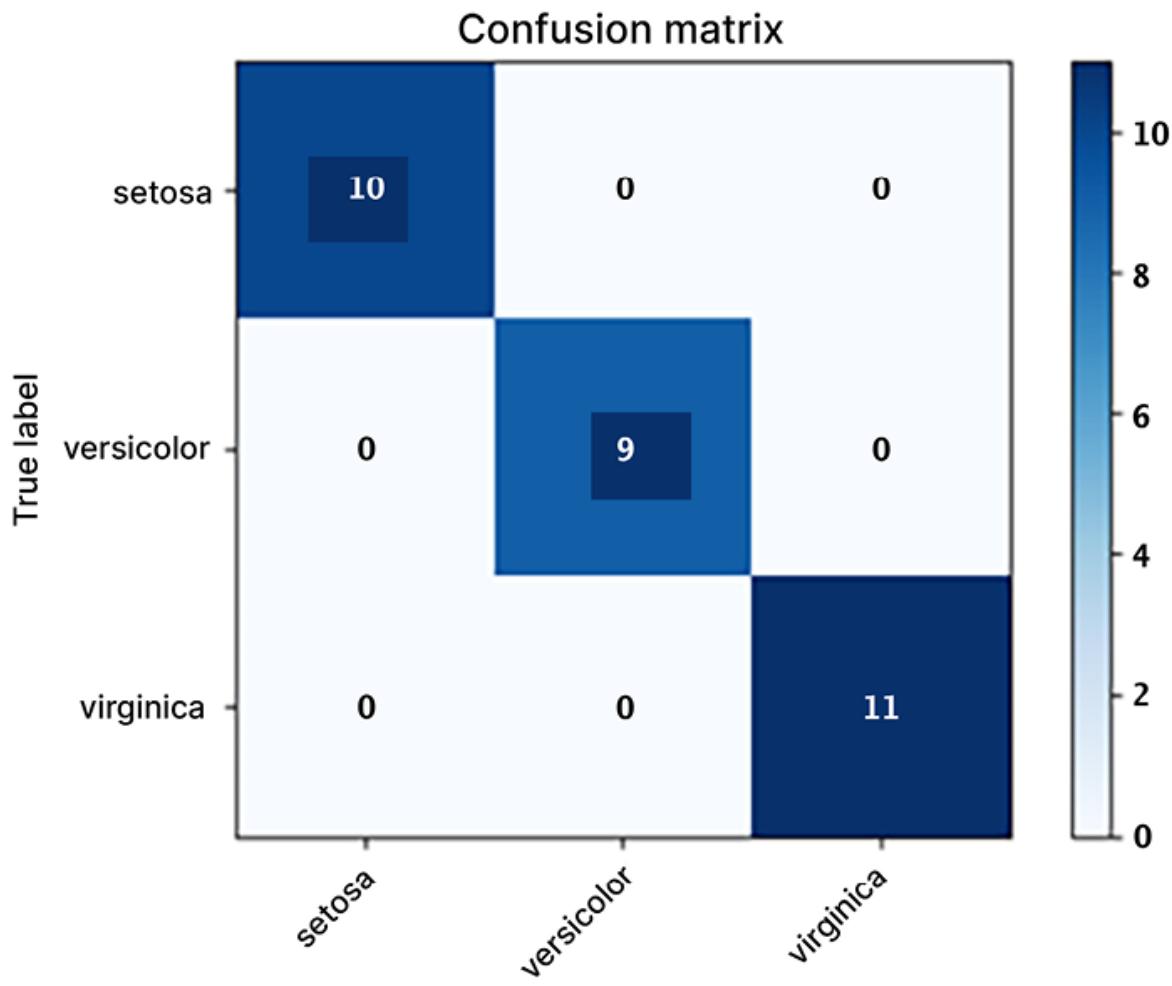


Figure 1.12: ChatGPT creates a confusion matrix of the model results

Finally, ChatGPT plots the decision tree:

```
7.2. Plot the Decision Tree  
# Plot the decision tree  
plt.figure(figsize=(20,10))  
plot_tree(model, feature_names=iris.feature_names, class_names:  
plt.show()
```

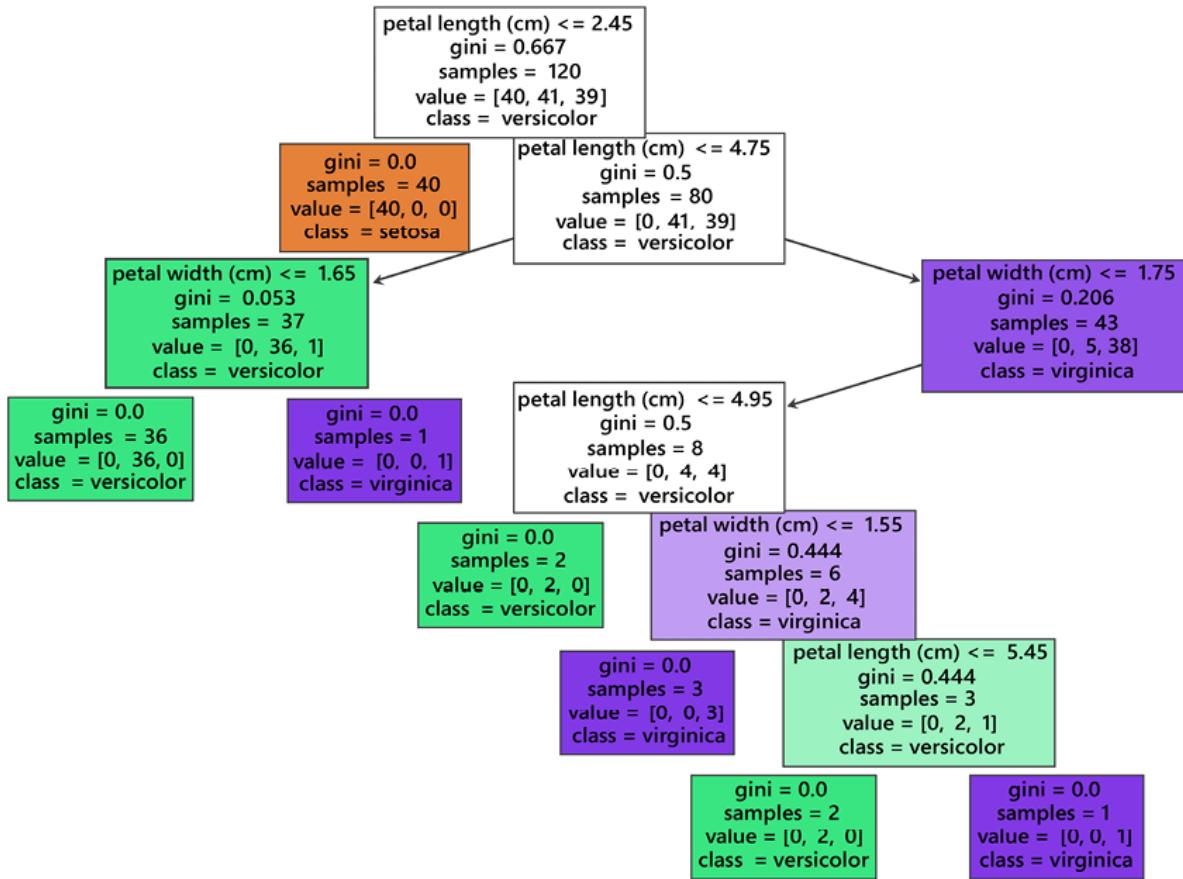


Figure 1.13: Excerpt of the decision tree plotted by ChatGPT

Take a deep breath and let what you just saw sink in. An AI program did 100% of the job by itself. So, of course, the first reaction could be: *Uh oh, AI is going to take over my job! Generative AI will certainly increase its scope if complex problems can be broken down into small components!*

A second reaction could be, “*Hmm, I’m going to push it to its limits until it makes mistakes to prove it’s not worth the hype! AI will never be able to replace the complexity of what I do.*”

These reactions will intensify as transformer-driven AI spreads out to a wide range of daily applications.

As AI professionals, our responsibility entails acquiring a comprehensive understanding of transformers and their functionalities, evaluating their limitations, and enhancing our skill set.

## The role of AI professionals

Transformer-driven AI is connecting everything to everything, everywhere. Machines communicate directly with other machines. AI-driven IoT signals trigger automated decisions without human intervention. NLP algorithms send automated reports, summaries, emails, advertisements, and more.

AI specialists must adapt to this new era of increasingly automated tasks, including transformer model implementations. AI professionals will have new functions. If we list transformer NLP tasks that an AI specialist will have to do, from top to bottom, it appears that some high-level tasks require little to no development from an AI specialist. An AI specialist can be an AI guru, providing design ideas, explanations, and implementations.

The pragmatic definition of what a transformer represents for an AI specialist will vary with the ecosystem.

Let's go through a few examples:

- **API:** The OpenAI API does not require an AI developer. A web designer can create a form, and a linguist or **Subject Matter Expert (SME)** can prepare the prompt input texts. The primary role of an AI specialist will require linguistic skills to show, not just tell, the ChatGPT/GPT-4 models and their successors how to accomplish a task. Showing a model what to do, for example, involves working on the context of the input. This new task is named *prompt engineering*. A *prompt engineer* has quite a future in AI!

- **Library:** The Google Trax library requires limited development to start with ready-to-use models. An AI professional with linguistics and NLP skills can work on the datasets and the outputs. The AI professional can also use Trax's toolkit to build tailored models for a project.
- **Training and fine-tuning:** Some Hugging Face functionality requires limited development, providing both APIs and libraries. However, sometimes, we still have to get our hands dirty. In that case, training, fine-tuning the models, and finding the correct hyperparameters will require the expertise of an AI specialist.
- **Development-level skills:** In some projects, the tokenizers and the datasets do not match. Or the embeddings of a model might not fit a project. In this case, an AI developer working with a linguist, for example, can play a crucial role. Therefore, computational linguistics training can come in very handy at this level.

The recent evolution of NLP AI can be termed as “embedded transformers” in assistants, copilots, and everyday software, which is disrupting the AI development ecosystem:

- LLM transformers with billions of parameters and many layers to train, such as OpenAI GPT-4, are currently embedded in several Microsoft Azure applications with GitHub Copilot, among other services and software.
- The embedded transformers are not accessible directly but provide automatic development support, such as automatic code generation. They also provide summarization, question and answering capabilities, and many other tasks in a growing number of applications.
- The usage of embedded transformers is both endless and seamless for the end user with assisted text completion.

We will explore this fascinating new world of embedded transformers throughout our journey in this book.

The skill set of an LLM AI professional requires adaptability, cross-disciplinary knowledge, and above all, *flexibility*. This book will provide the AI specialist with a variety of transformer ecosystems to adapt to the new paradigms of the market. These opposing and often conflicting strategies leave us with various possible implementations.

The future of AI professionals will expand into many specializations.

## **The future of AI professionals**

The societal impact of Foundation Models should not be underestimated. Prompt engineering has become a skill required for AI professionals.

An AI professional will also be involved in machine-to-machine algorithms using classical AI, IoT, edge computing, and more. An AI specialist will also design and develop fascinating connections between bots, robots, servers, and all types of connected devices using classical algorithms.

This book is thus not limited to prompt engineering but to a wide range of design skills required to be an LLM AI specialist.

Prompt engineering is a subset of the design skills an AI professional must develop. Here are some domains an AI professional can invest in:

- **AI specialists** have specialized knowledge or skills in one of the fields of AI. One of the interesting fields is fine-tuning, maintaining, and supporting AI systems.
- **AI architects** design architectures and provide solutions for scaling and other information for deployment and production.

- **AI experts** have authoritative knowledge and skills in AI. An AI expert can contribute to the field with research, papers, and innovative solutions.
- **AI analysts** focus on data, big data, and model architecture to provide information for other teams.
- **AI researchers** focus on research in a university or private company.
- **AI engineers** design AI, build systems, and implement AI models.
- **AI managers** leverage the critical skills of project management, product management, and resource management (human, machine, and software).
- Many other fields will appear showing that there are many opportunities in AI to develop assets, such as datasets, models, and applications.

Transformers have spread everywhere, and we need to find the right resource for a project.

## **What resources should we use?**

Generative AI has blurred the lines between cloud platforms, frameworks, libraries, languages, and models. Transformers are new, and the range and number of ecosystems are mind-blowing. Google Cloud provides ready-to-use transformer models.

OpenAI has deployed a transformer API that requires practically no programming. Hugging Face provides a cloud library service, and the list is endless.

Microsoft 365 and Google Workspace now possess Generative AI tools powered by state-of-the-art transformers. As a result, every Microsoft 365 and Google Workspace user has AI at the tip of their fingers!

Your choice of resources to implement transformers for NLP is critical. It is a question of survival in a project. Imagine a real-life interview or presentation. Imagine you are talking to your future employer, your employer, your team, or a customer.

For example, you start your presentation with an awesome PowerPoint on Hugging Face. You might get an adverse reaction from a manager who says, “*I’m sorry, but we use Google Cloud AI here for this type of project, not Hugging Face. Can you implement Google Cloud AI, please?*” If you don’t, it’s game over for you.

The same problem could have arisen by specializing in Google Cloud AI. But instead, you might get the reaction of a manager who wants to use OpenAI’s ChatGPT/GPT-4 models with an API and no development. If you specialize in OpenAI’s GPT engines with APIs and no development, you might face a project manager or customer who prefers Hugging Face’s AutoML or HuggingGPT!

The worst thing that could happen to you is that a manager accepts your solution, but in the end, it does not work at all for the NLP tasks of that project.

The key concept to remember is that if you only focus on the solution you like, you will most likely sink with the ship at some point.

*Focus on the system you need, not the one you like.*

This book is not designed to explain every transformer solution on the market. Instead, this book aims to explain enough transformer ecosystems for you to be flexible and adapt to any situation you face in an NLP project.

Let’s go through a few approaches that can help you make a decision.

# Decision-making guidelines

My personal experience relies on three simple principles:

- Understand the task(s) the AI model is expected to do at an Subject Matter Expert level. Write it down on paper without ChatGPT or any help. Just you. Then, write a detailed flow chart.
- Build a prototype dataset that can be duplicated to whatever scale you need to test the hardware/software performances.
- Try one to a few NLP/CV models to see which one fits your needs: open source or not, API or not, local installation, cloud platform, and the other constraints of your project.

Take your time. Don't let yourself get bullied by anybody. Use science to guide you, not opinions.

Let's now ask ChatGPT Plus for additional guidance for the three principles I just mentioned. ChatGPT has some interesting advice:

1. **Problem assessment:** Begin by thoroughly understanding the problem. What are the key challenges? What scale of data are you dealing with?
2. **Skill set evaluation:** What is the team's proficiency level? Familiarity with certain platforms might accelerate development.
3. **Scalability:** Consider the growth trajectory. Will the chosen solution scale as data grows or as the problem complexity increases?
4. **Community and support:** Tools and models with active communities can be beneficial due to available resources, tutorials, and prompt solutions to common issues.
5. **Interoperability:** How well does the solution integrate with existing systems or potential future tools?

Remember, the tool or model is only as good as the team's understanding of its principles and its applicability to the problem at hand.

One of the challenges of an NLP project remains to assess the solution to implement. AI APIs are one of the ever-evolving tools on the market we must explore.

## The rise of transformer seamless APIs and assistants

We are now well into the industrialization era of AI. Microsoft Azure, Google Cloud, **Amazon Web Services (AWS)**, and IBM, among others, provide AI services that no developer or team of developers could hope to outperform. Tech giants have million-dollar supercomputers with massive datasets to train transformer models and AI models in general.

Big tech giants have many corporate customers that already use their cloud services. As a result, adding a transformer API to an existing cloud architecture requires less effort than any other solution.

A small company or even an individual can access the most powerful transformer models through an API with practically no investment in development. An intern can implement the API in a few days. There is no need to be an engineer or have a Ph.D. for such a simple implementation.

For example, the OpenAI platform now has a **Software as a Service (SaaS)** API for some of the most effective transformer models on the market.

OpenAI transformer models are so effective and humanlike that the present policy requires a potential user to complete a request form. Once the request

has been accepted, the user can access a universe of natural language processing!

The simplicity of OpenAI's API takes the user by surprise:

1. Obtain an API key in one click.
2. Import OpenAI in a notebook in one line.
3. Enter any NLP task you wish in a *prompt*.
4. You will receive a response with no other functions to write.

For example, you can translate natural language to an SQL query, as explained on OpenAI's platform,

<https://platform.openai.com/examples/default-sql-translate>:

### Prompt

```
#### Postgres SQL tables, with their properties:  
#  
# Employee(id, name, department_id)  
# Department(id, name, address)  
# Salary_Payments(id, employee_id, amount, date)  
#  
#### A query to list the names of the departments which employed more than  
10 employees in the last 3 months  
SELECT
```

Figure 1.14: Prompt to create an SQL query on OpenAI

## Sample response

```
SELECT d.name
  FROM Department d
 INNER JOIN Employee e ON d.id = e.department_id
 INNER JOIN Salary_Payments sp ON e.id = sp.employee_id
 WHERE sp.date > NOW() - INTERVAL '3 months'
 GROUP BY d.name
 HAVING COUNT(*) > 10
```

Figure 1.15: Response from the SQL prompt

And that's it! The queries may not always be perfect and may need amending, but this shows you the power of LLMs. Welcome to the world of Generative AI!

Developers focusing on code-only solutions will evolve into a generation of developers with cross-disciplinary mindsets leveraging the power of AI copilots.

The AI professionals will learn how to design ways to *show* a transformer model what is expected and not intuitively *tell* it what to do. By the end of the book, you will have acquired several methods to control the behavior of these cutting-edge AI models.

Though APIs may satisfy many needs, they also have limits. For example, a multipurpose API might be reasonably good in all tasks but not good enough for a specific NLP task. For instance, translating with transformers is no easy task. In that case, an AI developer, consultant, or project manager must prove that an API alone cannot solve the required NLP task. Therefore, we might need to search for a solid library, alternative solutions, or develop one.

# Choosing ready-to-use API-driven libraries

In this book, we will explore several libraries. For example, Google has some of the most advanced AI labs in the world. Google Trax can be installed in a few lines in Google Colab. We can choose free or paid services. We can get our hands on source code, tweak the models, and even train them on our servers or Google Cloud. For example, it's a step down from ready-to-use APIs to customize a transformer model for translation tasks.

However, it can prove to be both educational and effective in some cases. We will explore the recent evolution of Google in translations and implement Google Trax in *Chapter 4, Advancements in Translations with Google Trax, Google Translate, and Gemini*.

We have seen that APIs, such as OpenAI, require limited developer skills, and libraries, such as Google Trax, dig a bit deeper into code. Both approaches show that AI APIs will require more development on the editor side of the API but much less effort when implementing transformers.

Google Translate is one of the most famous online applications that use transformers, among other algorithms. Google Translate can be used online or through an API.

Try translating a sentence requiring coreference resolution in English to French using Google Translate. The sentence is, *A user visited the AllenNLP website, tried a transformer model, and found it interesting*. Google Translate produces the following translation:

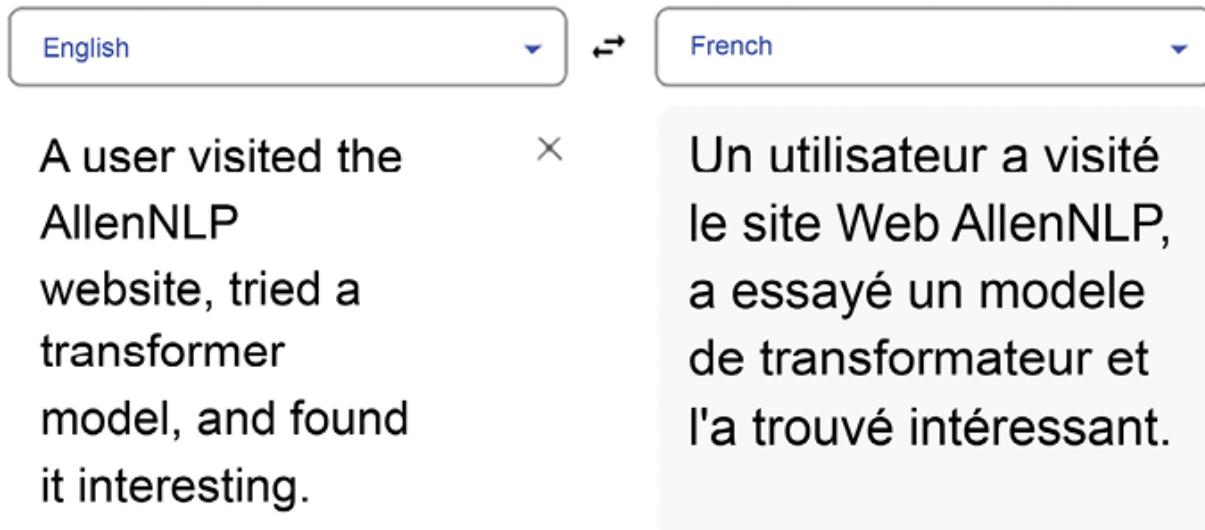


Figure 1.16: Coreference resolution in a translation using Google Translate

Google Translate appears to have solved the coreference resolution, but the word *transformateur* in French means an electric device. The word *transformer* is a neologism (new word) in French. An AI specialist might be required to have language and linguistic skills for a specific project. Significant development is not required in this case. However, the project might require clarifying the input before requesting a translation.

This example shows that you might have to team up with a linguist or acquire linguistic skills to work on an input context. In addition, it might take much development to enhance the input with an interface for contexts.

Google may improve this example by the time the book is published, but there are thousands more limitations to the system.

So, we still might have to get our hands dirty to add scripts to use Google Translate. Or we might have to find a transformer model for a specific translation need, such as BERT, T5, or other models we will explore in this book.

Choosing an API is one task. Finding a suitable transformer model is no easy task with the increasing range of solutions.

## **Choosing a cloud platform and transformer model**

Big tech corporations dominate the NLP market. Google, Facebook, and Microsoft alone run billions of NLP routines daily, increasing their AI models' unequaled power thanks to the data they gather. The big giants now offer a wide range of transformer models and have top-ranking Foundation Models.

The threshold of Foundation Models is fully trained transformers on supercomputers such as OpenAI GPT-4, Google PaLM, and the new models they will continually produce. Foundation Models are often proprietary, meaning we cannot access their code, although we can fine-tune some of them through a cloud service.

Hugging Face has a different approach and offers a wide range of transformer models with model cards, source code, datasets, and more development resources. In addition, Hugging Face offers high-level APIs and developer-controlled APIs. In several chapters of this book, we will explore Hugging Face as an educational tool and a possible solution for many tasks.

Google Cloud, Microsoft Azure, AWS, Hugging Face, and others offer fantastic services on their platforms! When looking at the growing mountain of transformer-driven AI, we need to find where we fit in.

## **Summary**

Transformers forced AI to make profound evolutions. Foundation Models, including their Generative AI abilities, are built on top of the digital revolution connecting everything to everything with underlying processes everywhere. Automated processes are replacing human decisions in critical areas, including NLP.

RNNs slowed the progression of automated NLP tasks required in a fast-moving world. Transformers filled the gap. A corporation needs summarization, translation, and a wide range of NLP tools to meet the challenges of the growing volume of incoming information.

Transformers have thus spurred an age of AI industrialization. We first saw how the  $O(1)$  time complexity of the attention layers and their computational time complexity,  $O(n^2*d)$ , shook the world of AI.

We saw how the one-token flexibility of transformer models pervaded every domain of our everyday lives!

Platforms such as Hugging Face, Google Cloud, OpenAI, and Microsoft Azure provide NLP tasks without installation and resources to implement a transformer model in customized programs. For example, OpenAI provides an API requiring only a few code lines to run one of the powerful GPT-4 generation models. Google Trax provides an end-to-end library, Google Vertex AI has easy-to-implement tools, and Hugging Face offers various transformer models and implementations. We will be exploring these ecosystems throughout this book.

We then saw that the transformer architecture leads to automation that radically deviates from the former AI. As a result, broader skill sets are required for an AI professional. For example, a project manager can decide to implement transformers by asking a web designer to create an interface for Google Cloud AI or OpenAI APIs through prompt engineering. Or, when

required, a project manager can ask an AI specialist to download Google Trax or Hugging Face resources to develop a full-blown project with a customized transformer model.

The transformer is a game-changer for developers whose roles will expand and require more designing than programming. In addition, embedded transformers will provide assisted code development and usage. *These new skill sets are a challenge but lead to new, exciting horizons.*

In *Chapter 2, Getting Started with the Architecture of the Transformer Model*, we will start with the Original Transformer architecture.

## Questions

1. ChatGPT is a game-changer. (True/False)
2. ChatGPT can replace all AI algorithms. (True/False)
3. AI developers will sometimes have no AI development to do.  
(True/False)
4. AI developers might have to implement transformers from scratch.  
(True/False)
5. It's not necessary to learn more than one transformer ecosystem, such as Hugging Face. (True/False)
6. A ready-to-use transformer API can satisfy all needs. (True/False)
7. A company will accept the transformer ecosystem a developer knows best. (True/False)
8. Cloud transformers have become mainstream. (True/False)
9. A transformer project can be run on a laptop. (True/False)
10. AI specialists will have to be more flexible. (True/False)

# References

- *Bommansani et al., 2021, On the Opportunities and Risks of Foundation Models:* <https://arxiv.org/abs/2108.07258>
- *Rishi Bommansani, Dilara Soylu, Thomas I. Liao, Kathleen A. Creel, and Percy Liang, 2023, Ecosystem Graphs: The Social Footprint of Foundation Models:* <https://arxiv.org/abs/2303.15772>
- *Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin, 2017, Attention is All You Need:*  
<https://arxiv.org/abs/1706.03762>
- *Chen et al., 2021, Evaluating Large Language Models Trained on Code:* <https://arxiv.org/abs/2107.03374>
- Microsoft AI: <https://innovation.microsoft.com/en-us/ai-at-scale>
- OpenAI: <https://openai.com/>
- Google AI: <https://ai.google/>
- Google Trax: <https://github.com/google/trax>
- AllenNLP: <https://allennlp.org/>
- Hugging Face: <https://huggingface.co/>
- Google Cloud TPU:  
<https://cloud.google.com/tpu/docs/intro-to-tpu>

# Further reading

- *Tyna Eloundou, Sam Manning, Pamela Mishkin, and Daniel Rock, 2023, GPTs are GPTs: An Early Look at the Labor Market Impact*