

# intro\_palm\_api

September 29, 2023

```
[1]: # Copyright 2023 Google LLC  
#  
# Licensed under the Apache License, Version 2.0 (the "License");  
# you may not use this file except in compliance with the License.  
# You may obtain a copy of the License at  
#  
#     https://www.apache.org/licenses/LICENSE-2.0  
#  
# Unless required by applicable law or agreed to in writing, software  
# distributed under the License is distributed on an "AS IS" BASIS,  
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.  
# See the License for the specific language governing permissions and  
# limitations under the License.
```

## 1 Getting Started with the Vertex AI PaLM API & Python SDK

Run in Colab

View on GitHub

Open in Vertex AI Workbench

### 1.1 Overview

#### 1.1.1 What are LLMs?

Large language models (LLMs) are deep learning models trained on massive datasets of text. LLMs can translate language, summarize text, generate creative writing, generate code, power chatbots and virtual assistants, and complement search engines and recommendation systems.

#### 1.1.2 PaLM

Following its predecessor, [PaLM](#), [PaLM 2](#) is Google's next generation large language model that builds on Google's legacy of breakthrough research in machine learning and responsible AI. PaLM 2 excels at tasks like advanced reasoning, translation, and code generation because of how it was built.

PaLM 2 [excels](#) at advanced reasoning tasks, including code and math, classification and question answering, translation and multilingual proficiency, and natural language generation better than our previous state-of-the-art LLMs, including PaLM. It can accomplish these tasks because of the

way it was built – bringing together compute-optimal scaling, an improved dataset mixture, and model architecture improvements.

PaLM 2 is grounded in Google’s approach to building and deploying AI responsibly. It was evaluated rigorously for its potential harms and biases, capabilities and downstream uses in research and in-product applications. It’s being used in other state-of-the-art models, like Med-PaLM 2 and Sec-PaLM, and is powering generative AI features and tools at Google, like Bard and the PaLM API.

PaLM is pre-trained on a wide range of text data using an unsupervised learning approach, without any specific task. During this pre-training process, PaLM learns to predict the next word in a sentence, given the preceding words. This enables the model to generate coherent, fluent text resembling human writing. This large size enables it to learn complex patterns and relationships in language and generate high-quality text for various applications. This is why models like PaLM are referred to as “foundational models.”

Creating an LLM requires massive amounts of data, significant compute resources, and specialized skills. Because LLMs require a big investment to create, they target broad rather than specific use cases. On Vertex AI, you can customize a foundation model for more specific tasks or knowledge domains by using prompt design and model tuning.

### 1.1.3 Vertex AI PaLM API

The Vertex AI PaLM API, [released on May 10, 2023](#), is powered by [PaLM 2](#).

### 1.1.4 Using Vertex AI PaLM API

You can interact with the Vertex AI PaLM API using the following methods:

- Use the [Generative AI Studio](#) for quick testing and command generation.
- Use cURL commands in Cloud Shell.
- Use the Python SDK in a Jupyter notebook

This notebook focuses on using the Python SDK to call the Vertex AI PaLM API. For more information on using Generative AI Studio without writing code, you can explore [Getting Started with the UI instructions](#)

For more information, check out the [documentation on generative AI support for Vertex AI](#).

### 1.1.5 Objectives

In this tutorial, you will learn how to use PaLM API with the Python SDK and explore its various parameters.

By the end of the notebook, you should be able to understand various nuances of generative model parameters like `temperature`, `top_k`, `top_p`, and how each parameter affects the results.

The steps performed include:

- Installing the Python SDK
- Using Vertex AI PaLM API
  - Text generation model with `text-bison@001`

- \* Understanding model parameters (`temperature`, `max_output_token`, `top_k`, `top_p`)
- Chat model with `chat-bison@001`
- Embeddings model with `textembedding-gecko@001`

### 1.1.6 Costs

This tutorial uses billable components of Google Cloud:

- Vertex AI Generative AI Studio

Learn about [Vertex AI pricing](#), and use the [Pricing Calculator](#) to generate a cost estimate based on your projected usage.

### 1.1.7 Data governance and security

For more information, see the documentation on [Data Governance and Generative AI](#) on Google Cloud.

### 1.1.8 Responsible AI

Large language models (LLMs) can translate language, summarize text, generate creative writing, generate code, power chatbots and virtual assistants, and complement search engines and recommendation systems. At the same time, as an early-stage technology, its evolving capabilities and uses create potential for misapplication, misuse, and unintended or unforeseen consequences. Large language models can generate output that you don't expect, including text that's offensive, insensitive, or factually incorrect.

What's more, the incredible versatility of LLMs is also what makes it difficult to predict exactly what kinds of unintended or unforeseen outputs they might produce. Given these risks and complexities, the PaLM API is designed with [Google's AI Principles](#) in mind. However, it is important for developers to understand and test their models to deploy safely and responsibly. To aid developers, the Generative AI Studio has built-in content filtering, and the PaLM API has safety attribute scoring to help customers test Google's safety filters and define confidence thresholds that are right for their use case and business. Please refer to the [Safety filters and attributes](#) section to learn more.

When the PaLM API is integrated into a customer's unique use case and context, additional responsible AI considerations and [PaLM limitations](#) may need to be considered. We encourage customers to leverage fairness, interpretability, privacy and security [recommended practices](#).

## 1.2 Getting Started

### 1.2.1 Install Vertex AI SDK

```
[2]: !pip install google-cloud-aiplatform --upgrade --user
```

```
Requirement already satisfied: google-cloud-aiplatform in
/home/jupyter/.local/lib/python3.10/site-packages (1.33.1)
Requirement already satisfied: google-api-core[grpc]!=2.0.*,!=2.1.*,!=2.2.*,!=2.3.*,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0 in
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (1.34.0)
```

Requirement already satisfied: proto-plus<2.0.0dev,>=1.22.0 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (1.22.3)

Requirement already satisfied: protobuf!=3.20.0,!=3.20.1,!=4.21.0,!=4.21.1,!=4.21.2,!=4.21.3,!=4.21.4,!=4.21.5,<5.0.0dev,>=3.19.5 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (3.19.6)

Requirement already satisfied: packaging>=14.3 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (23.1)

Requirement already satisfied: google-cloud-storage<3.0.0dev,>=1.32.0 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (2.10.0)

Requirement already satisfied: google-cloud-bigquery<4.0.0dev,>=1.15.0 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (3.11.4)

Requirement already satisfied: google-cloud-resource-manager<3.0.0dev,>=1.3.3 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (1.10.3)

Requirement already satisfied: shapely<2.0.0 in /opt/conda/lib/python3.10/site-packages (from google-cloud-aiplatform) (1.8.5.post1)

Requirement already satisfied: googleapis-common-protos<2.0dev,>=1.56.2 in  
/opt/conda/lib/python3.10/site-packages (from google-api-core[grpc]!=2.0.\*,!=2.1.\*,!=2.2.\*,!=2.3.\*,!=2.4.\*,!=2.5.\*,!=2.6.\*,!=2.7.\*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (1.60.0)

Requirement already satisfied: google-auth<3.0dev,>=1.25.0 in  
/opt/conda/lib/python3.10/site-packages (from google-api-core[grpc]!=2.0.\*,!=2.1.\*,!=2.2.\*,!=2.3.\*,!=2.4.\*,!=2.5.\*,!=2.6.\*,!=2.7.\*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (2.22.0)

Requirement already satisfied: requests<3.0.0dev,>=2.18.0 in  
/opt/conda/lib/python3.10/site-packages (from google-api-core[grpc]!=2.0.\*,!=2.1.\*,!=2.2.\*,!=2.3.\*,!=2.4.\*,!=2.5.\*,!=2.6.\*,!=2.7.\*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (2.31.0)

Requirement already satisfied: grpcio<2.0dev,>=1.33.2 in  
/opt/conda/lib/python3.10/site-packages (from google-api-core[grpc]!=2.0.\*,!=2.1.\*,!=2.2.\*,!=2.3.\*,!=2.4.\*,!=2.5.\*,!=2.6.\*,!=2.7.\*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (1.48.1)

Requirement already satisfied: grpcio-status<2.0dev,>=1.33.2 in  
/opt/conda/lib/python3.10/site-packages (from google-api-core[grpc]!=2.0.\*,!=2.1.\*,!=2.2.\*,!=2.3.\*,!=2.4.\*,!=2.5.\*,!=2.6.\*,!=2.7.\*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (1.48.1)

Requirement already satisfied: google-cloud-core<3.0.0dev,>=1.6.0 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-bigquery<4.0.0dev,>=1.15.0->google-cloud-aiplatform) (2.3.3)

Requirement already satisfied: google-resumable-media<3.0dev,>=0.6.0 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-bigquery<4.0.0dev,>=1.15.0->google-cloud-aiplatform) (2.5.0)

Requirement already satisfied: python-dateutil<3.0dev,>=2.7.2 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-bigquery<4.0.0dev,>=1.15.0->google-cloud-aiplatform) (2.8.2)

Requirement already satisfied: grpc-google-iam-v1<1.0.0dev,>=0.12.4 in  
/opt/conda/lib/python3.10/site-packages (from google-cloud-resource-manager<3.0.0dev,>=1.3.3->google-cloud-aiplatform) (0.12.6)

Requirement already satisfied: cachetools<6.0,>=2.0.0 in

```

/opt/conda/lib/python3.10/site-packages (from google-
auth<3.0dev,>=1.25.0->google-api-core[grpc]!=2.0.*,!=2.1.*,!=2.2.*,!=2.3.*,!=2.4
.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (4.2.4)
Requirement already satisfied: pyasn1-modules>=0.2.1 in
/opt/conda/lib/python3.10/site-packages (from google-
auth<3.0dev,>=1.25.0->google-api-core[grpc]!=2.0.*,!=2.1.*,!=2.2.*,!=2.3.*,!=2.4
.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (0.2.7)
Requirement already satisfied: rsa<5,>=3.1.4 in /opt/conda/lib/python3.10/site-
packages (from google-auth<3.0dev,>=1.25.0->google-api-core[grpc]!=2.0.*,!=2.1.*
,!=2.2.*,!=2.3.*,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-
cloud-aiplatform) (4.9)
Requirement already satisfied: six>=1.9.0 in /opt/conda/lib/python3.10/site-
packages (from google-auth<3.0dev,>=1.25.0->google-api-core[grpc]!=2.0.*,!=2.1.*
,!=2.2.*,!=2.3.*,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-
cloud-aiplatform) (1.16.0)
Requirement already satisfied: urllib3<2.0 in /opt/conda/lib/python3.10/site-
packages (from google-auth<3.0dev,>=1.25.0->google-api-core[grpc]!=2.0.*,!=2.1.*
,!=2.2.*,!=2.3.*,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-
cloud-aiplatform) (1.26.15)
Requirement already satisfied: google-crc32c<2.0dev,>=1.0 in
/opt/conda/lib/python3.10/site-packages (from google-resumable-
media<3.0dev,>=0.6.0->google-cloud-bigquery<4.0.0dev,>=1.15.0->google-cloud-
aiplatform) (1.5.0)
Requirement already satisfied: charset-normalizer<4,>=2 in
/opt/conda/lib/python3.10/site-packages (from
requests<3.0.0dev,>=2.18.0->google-api-core[grpc]!=2.0.*,!=2.1.*,!=2.2.*,!=2.3.*
,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform)
(3.2.0)
Requirement already satisfied: idna<4,>=2.5 in /opt/conda/lib/python3.10/site-
packages (from requests<3.0.0dev,>=2.18.0->google-api-core[grpc]!=2.0.*,!=2.1.*
,!=2.2.*,!=2.3.*,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-
cloud-aiplatform) (3.4)
Requirement already satisfied: certifi>=2017.4.17 in
/opt/conda/lib/python3.10/site-packages (from
requests<3.0.0dev,>=2.18.0->google-api-core[grpc]!=2.0.*,!=2.1.*,!=2.2.*,!=2.3.*
,!=2.4.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform)
(2023.7.22)
Requirement already satisfied: pyasn1<0.5.0,>=0.4.6 in
/opt/conda/lib/python3.10/site-packages (from pyasn1-modules>=0.2.1->google-
auth<3.0dev,>=1.25.0->google-api-core[grpc]!=2.0.*,!=2.1.*,!=2.2.*,!=2.3.*,!=2.4
.*,!=2.5.*,!=2.6.*,!=2.7.*,<3.0.0dev,>=1.32.0->google-cloud-aiplatform) (0.4.8)

```

**Colab only:** Uncomment the following cell to restart the kernel or use the button to restart the kernel. For Vertex AI Workbench you can restart the terminal using the button on top.

```

[3]: # # Automatically restart kernel after installs so that your environment can
      ↪ access the new packages
      # import IPython

```

```
# app = IPython.Application.instance()
# app.kernel.do_shutdown(True)
```

### 1.2.2 Authenticating your notebook environment

- If you are using **Colab** to run this notebook, uncomment the cell below and continue.
- If you are using **Vertex AI Workbench**, check out the setup instructions [here](#).

```
[4]: # from google.colab import auth
# auth.authenticate_user()
```

## 1.3 Vertex AI PaLM API models

The Vertex AI PaLM API enables you to test, customize, and deploy instances of Google's large language models (LLM) called as PaLM, so that you can leverage the capabilities of PaLM in your applications.

### 1.3.1 Model naming scheme

Foundation model names have three components: use case, model size, and version number. The naming convention is in the format:

`<use case>-<model size>@<version number>`

For example, `text-bison@001` represents the Bison text model, version 001.

The model sizes are as follows: - **Bison**: The best value in terms of capability and cost. - **Gecko**: The smallest and cheapest model for simple tasks.

### 1.3.2 Available models

The Vertex AI PaLM API currently supports five models:

- `text-bison@001` : Fine-tuned to follow natural language instructions and is suitable for a variety of language tasks.
- `chat-bison@001` : Fine-tuned for multi-turn conversation use cases like building a chatbot.
- `textembedding-gecko@001` : Returns model embeddings for text inputs.
- `code-bison@001`: A model fine-tuned to generate code based on a natural language description of the desired code. For example, it can generate a unit test for a function.
- `code-gecko@001`: A model fine-tuned to suggest code completion based on the context in code that's written.
- `codechat-bison@001`: A model fine-tuned for chatbot conversations that help with code-related questions.

You can find more information about the properties of these [foundational models in the Generative AI Studio documentation](#).

### 1.3.3 Import libraries

**Colab only:** Uncomment the following cell to initialize the Vertex AI SDK. For Vertex AI Workbench, you don't need to run this.

```
[5]: # import vertexai

# PROJECT_ID = "" # @param {type:"string"}
# vertexai.init(project=PROJECT_ID, location="us-central1")
```

```
[6]: import pandas as pd
import seaborn as sns
from IPython.display import Markdown, display
from sklearn.metrics.pairwise import cosine_similarity
from vertexai.language_models import TextGenerationModel, \
    TextEmbeddingModel, \
    ChatModel, \
    InputOutputTextPair, \
    CodeGenerationModel, \
    CodeChatModel
```

```
2023-09-29 15:56:25.971720: I tensorflow/core/platform/cpu_feature_guard.cc:193]
This TensorFlow binary is optimized with oneAPI Deep Neural Network Library
(oneDNN) to use the following CPU instructions in performance-critical
operations:  AVX2 FMA
To enable them in other operations, rebuild TensorFlow with the appropriate
compiler flags.
2023-09-29 15:56:27.281510: W
tensorflow/compiler/xla/stream_executor/platform/default/dso_loader.cc:64] Could
not load dynamic library 'libnvinfer.so.7'; dLError: libnvinfer.so.7: cannot
open shared object file: No such file or directory; LD_LIBRARY_PATH:
/usr/local/cuda/lib64:/usr/local/ncc12/lib:/usr/local/cuda/extras/CUPTI/lib64
2023-09-29 15:56:27.281668: W
tensorflow/compiler/xla/stream_executor/platform/default/dso_loader.cc:64] Could
not load dynamic library 'libnvinfer_plugin.so.7'; dLError:
libnvinfer_plugin.so.7: cannot open shared object file: No such file or
directory; LD_LIBRARY_PATH:
/usr/local/cuda/lib64:/usr/local/ncc12/lib:/usr/local/cuda/extras/CUPTI/lib64
2023-09-29 15:56:27.281679: W
tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Cannot
dlopen some TensorRT libraries. If you would like to use Nvidia GPU with
TensorRT, please make sure the missing libraries mentioned above are installed
properly.
```

## 1.4 Text generation with text-bison@001

The text generation model from PaLM API that you will use in this notebook is **text-bison@001**. It is fine-tuned to follow natural language instructions and is suitable for a variety of language tasks, such as:

- Classification
- Sentiment analysis
- Entity extraction
- Extractive question-answering
- Summarization
- Re-writing text in a different style
- Ad copy generation
- Concept ideation
- Concept simplification

### Load model

```
[7]: generation_model = TextGenerationModel.from_pretrained("text-bison@001")
```

**Prompt design** Prompt design is the process of creating prompts that elicit the desired response from a language model. Prompt design is an important part of using language models because it allows non-specialists to control the output of the model with minimal overhead. By carefully crafting the prompts, you can nudge the model to generate a desired result. Prompt design can be an efficient way to experiment with adapting an LLM for a specific use case. The iterative process of repeatedly updating prompts and assessing the model's responses is sometimes called prompt engineering.

**Hello PaLM** Create your first prompt and send it to the text generation model.

```
[8]: prompt = "What is a large language model?"

response = generation_model.predict(prompt=prompt)

print(response.text)
```

A large language model (LLM) is a type of artificial intelligence (AI) model that can understand and generate human language. LLMs are trained on massive datasets of text and code, and they can learn to perform a wide variety of tasks, such as translating languages, writing different kinds of creative content, and answering your questions in an informative way.

LLMs are still under development, but they have the potential to revolutionize many industries. For example, LLMs could be used to create more accurate and personalized customer service experiences, to help doctors diagnose and treat diseases, and to even write entire books and movies.

### Try out your own prompt

- What are the biggest challenges facing the healthcare industry?
- What are the latest developments in the automotive industry?
- What are the biggest opportunities in the retail industry?
- (Try your own prompts!)



```
[9]: prompt = """Create a numbered list of 10 items. Each item in the list should be a trend in the tech industry.
```

```
Each trend should be less than 5 words.""" # try your own prompt
```

```
response = generation_model.predict(prompt=prompt)
```

```
print(response.text)
```

1. **Artificial intelligence**
2. **Machine learning**
3. **Blockchain**
4. **Virtual reality**
5. **Augmented reality**
6. **5G**
7. **Internet of things**
8. **Self-driving cars**
9. **Wearables**
10. **Cloud computing**

**Prompt templates** Prompt templates are useful if you have found a good way to structure your prompt that you can re-use. This can be also be helpful in limiting the open-endedness of freeform prompts. There are many ways to implement prompt templates, and below is just one example using f-strings.

```
[10]: my_industry = "tech" # try changing this to a different industry
```

```
response = generation_model.predict(
    prompt=f"""Create a numbered list of 10 items. Each item in the list should be a trend in the {my_industry} industry.
```

```
Each trend should be less than 5 words."""
```

```
)
```

```
print(response.text)
```

1. **Artificial intelligence**
2. **Machine learning**
3. **Blockchain**
4. **Virtual reality**
5. **Augmented reality**
6. **Internet of things**
7. **Big data**
8. **Cloud computing**
9. **Cybersecurity**
10. **Wearables**

### 1.4.1 Model parameters for `text-bison@001`

You can customize how the PaLM API behaves in response to your prompt by using the following parameters for `text-bison@001`:

- **temperature**: higher means more “creative” responses
- **max\_output\_tokens**: sets the max number of tokens in the output
- **top\_p**: higher means it will pull from more possible next tokens, based on cumulative probability
- **top\_k**: higher means it will sample from more possible next tokens

The section below covers each parameter and how to use them.

#### The **temperature** parameter (range: 0.0 - 1.0, default 0)

**What is *temperature*?** The temperature is used for sampling during the response generation, which occurs when `top_p` and `top_k` are applied. Temperature controls the degree of randomness in token selection.

**How does *temperature* affect the response?** Lower temperatures are good for prompts that require a more deterministic and less open-ended response. In comparison, higher temperatures can lead to more “creative” or diverse results. A temperature of 0 is deterministic: the highest probability response is always selected. For most use cases, try starting with a temperature of 0.2.

A higher temperature value will result in a more explorative output, with a higher likelihood of generating rare or unusual words or phrases. Conversely, a lower temperature value will result in a more conservative output, with a higher likelihood of generating common or expected words or phrases.

**Example:** For example,

`temperature = 0.0:`

- *The cat sat on the couch, watching the birds outside.*
- *The cat sat on the windowsill, basking in the sun.*

`temperature = 0.9:`

- *The cat sat on the moon, meowing at the stars.*
- *The cat sat on the cheeseburger, purring with delight.*

**Note:** It’s important to note that while the temperature parameter can help generate more diverse and interesting text, it can also increase the likelihood of generating nonsensical or inappropriate text (i.e. hallucinations). Therefore, it’s important to use it carefully and with consideration for the desired outcome.

For more information on the **temperature** parameter for text models, please refer to the [documentation on model parameters](#).

If you run the following cell multiple times, it should always return the same response, as `temperature=0` is deterministic.

```
[11]: temp_val = 0.0
prompt_temperature = "Complete the sentence: As I prepared the picture frame, I_
↪reached into my toolkit to fetch my:"

response = generation_model.predict(
    prompt=prompt_temperature,
    temperature=temp_val,
)

print(f"[temperature = {temp_val}]")
print(response.text)
```

[temperature = 0.0]

As I prepared the picture frame, I reached into my toolkit to fetch my hammer.

If you run the following cell multiple times, it may return different responses, as higher temperature values can lead to more diverse results, even though the prompt is the same as the above cell.

```
[12]: temp_val = 1.0

response = generation_model.predict(
    prompt=prompt_temperature,
    temperature=temp_val,
)

print(f"[temperature = {temp_val}]")
print(response.text)
```

[temperature = 1.0]

As I prepared the picture frame, I reached into my toolkit to fetch my screwdriver. The screwdriver is a tool that is used to loosen screws and tighten them. It is an essential tool for any home improvement project. The screwdriver is also a tool that is used to fix furniture and other items around the house.

**The `max_output_tokens` parameter (range: 1 - 1024, default 128)**

**Tokens** A single token may be smaller than a word. For example, a token is approximately four characters. So 100 tokens correspond to roughly 60-80 words. It's essential to be aware of the token sizes as models have a limit on input and output tokens.

**What is `max_output_tokens`?** `max_output_tokens` is the maximum number of tokens that can be generated in the response.

**How does `max_output_tokens` affect the response?** Specify a lower value for shorter responses and a higher value for longer responses. A token may be smaller than a word. A token is approximately four characters. 100 tokens correspond to roughly 60-80 words.

For more information on the `max_output_tokens` parameter for text models, please refer to the [documentation on model parameters](#).

```
[13]: max_output_tokens_val = 5

response = generation_model.predict(
    prompt="List ten ways that generative AI can help improve the online_
↳shopping experience for users",
    max_output_tokens=max_output_tokens_val,
)

print(f"[max_output_tokens = {max_output_tokens_val}]")
print(response.text)
```

```
[max_output_tokens = 5]
```

```
1. **Personal
```

```
[14]: max_output_tokens_val = 500

response = generation_model.predict(
    prompt="List ten ways that generative AI can help improve the online_
↳shopping experience for users",
    max_output_tokens=max_output_tokens_val,
)

print(f"[max_output_tokens = {max_output_tokens_val}]")
print(response.text)
```

```
[max_output_tokens = 500]
```

1. **\*\*Personalized recommendations.\*\*** Generative AI can be used to create personalized recommendations for users based on their past purchases, browsing history, and other factors. This can help users find products that they are interested in and that are likely to be a good fit for them.
2. **\*\*Virtual try-on.\*\*** Generative AI can be used to create virtual try-on experiences for users, allowing them to see how different products would look on them before they buy them. This can help users make more informed decisions about their purchases and reduce the risk of buyer's remorse.
3. **\*\*Product design.\*\*** Generative AI can be used to help designers create new products that are more innovative and appealing to consumers. This can help businesses stay ahead of the competition and bring new products to market faster.
4. **\*\*Pricing.\*\*** Generative AI can be used to help businesses set prices for their products. This can help businesses maximize their profits and ensure that they are not overcharging or undercharging for their products.
5. **\*\*Marketing.\*\*** Generative AI can be used to create more effective marketing campaigns for businesses. This can help businesses reach a wider audience and generate more leads.
6. **\*\*Customer service.\*\*** Generative AI can be used to provide customer service support for businesses. This can help businesses provide faster and more accurate support to their customers.
7. **\*\*Fraud detection.\*\*** Generative AI can be used to detect fraud in online

transactions. This can help businesses protect themselves from financial loss and ensure that their customers are safe.

8. **\*\*Inventory management.\*\*** Generative AI can be used to help businesses manage their inventory more effectively. This can help businesses avoid stockouts and ensure that they have the right products in stock when customers need them.

9. **\*\*Supply chain management.\*\*** Generative AI can be used to help businesses manage their supply chains more efficiently. This can help businesses reduce costs and improve the delivery of their products to customers.

10. **\*\*Logistics.\*\*** Generative AI can be used to help businesses improve their logistics operations. This can help businesses reduce costs and get their products to customers faster.

For easier reading, you can also render Markdown in Jupyter:

```
[15]: display(Markdown(response.text))
```

1. **Personalized recommendations.** Generative AI can be used to create personalized recommendations for users based on their past purchases, browsing history, and other factors. This can help users find products that they are interested in and that are likely to be a good fit for them.
2. **Virtual try-on.** Generative AI can be used to create virtual try-on experiences for users, allowing them to see how different products would look on them before they buy them. This can help users make more informed decisions about their purchases and reduce the risk of buyer's remorse.
3. **Product design.** Generative AI can be used to help designers create new products that are more innovative and appealing to consumers. This can help businesses stay ahead of the competition and bring new products to market faster.
4. **Pricing.** Generative AI can be used to help businesses set prices for their products. This can help businesses maximize their profits and ensure that they are not overcharging or undercharging for their products.
5. **Marketing.** Generative AI can be used to create more effective marketing campaigns for businesses. This can help businesses reach a wider audience and generate more leads.
6. **Customer service.** Generative AI can be used to provide customer service support for businesses. This can help businesses provide faster and more accurate support to their customers.
7. **Fraud detection.** Generative AI can be used to detect fraud in online transactions. This can help businesses protect themselves from financial loss and ensure that their customers are safe.
8. **Inventory management.** Generative AI can be used to help businesses manage their inventory more effectively. This can help businesses avoid stockouts and ensure that they have the right products in stock when customers need them.
9. **Supply chain management.** Generative AI can be used to help businesses manage their supply chains more efficiently. This can help businesses reduce costs and improve the delivery of their products to customers.
10. **Logistics.** Generative AI can be used to help businesses improve their logistics operations. This can help businesses reduce costs and get their products to customers faster.

The **top\_p** parameter (range: 0.0 - 1.0, default 0.95)

**What is *top\_p*?** *top\_p* controls how the model selects tokens for output by adjusting the probability distribution of the next word in the generated text based on a cumulative probability cutoff. Specifically, it selects the smallest set of tokens whose cumulative probability exceeds the given cutoff probability *p*, and samples from this set uniformly.

For example, suppose tokens A, B, and C have a probability of 0.3, 0.2, and 0.1, and the *top\_p* value is 0.5. In that case, the model will select either A or B as the next token (using temperature) and not consider C, because the cumulative probability of *top\_p* is  $\leq 0.5$ . Specify a lower value for less random responses and a higher value for more random responses.

**How does *top\_p* affect the response?** The *top\_p* parameter is used to control the diversity of the generated text. A higher *top\_p* parameter value results in more “diverse” and “interesting” outputs, with the model being allowed to sample from a larger pool of possibilities. In contrast, a lower *top\_p* parameter value resulted in more predictable outputs, with the model being constrained to a smaller set of possible tokens.

**Example:** *top\_p* = 0.1:

- The cat sat on the mat.
- The cat sat on the floor.

*top\_p* = 0.9:

- The cat sat on the windowsill, soaking up the sun’s rays.
- The cat sat on the edge of the bed, watching the birds outside.

For more information on the *top\_p* parameter for text models, please refer to the [documentation on model parameters](#).

```
[16]: top_p_val = 0.0
prompt_top_p_example = (
    "Create a marketing campaign for jackets that involves blue elephants and_
    ↪avocados."
)

response = generation_model.predict(
    prompt=prompt_top_p_example, temperature=0.9, top_p=top_p_val
)

print(f"[top_p = {top_p_val}]")
print(response.text)
```

```
[top_p = 0.0]
**Introducing the new Blue Elephant Avocado Jacket!**
```

This stylish jacket is made from 100% recycled materials and is perfect for anyone who wants to make a statement. With its bright blue color and fun avocado print, this jacket is sure to turn heads wherever you go.

The Blue Elephant Avocado Jacket is also incredibly comfortable and versatile.

It's perfect for a day out on the town, a hike in the park, or even just a relaxing day at home.

So what are you waiting for? Order your Blue Elephant Avocado Jacket today!

**\*\*Here are some of the benefits of our**

```
[17]: top_p_val = 1.0

response = generation_model.predict(
    prompt=prompt_top_p_example, temperature=0.9, top_p=top_p_val
)

print(f"[top_p = {top_p_val}]")
print(response.text)
```

```
[top_p = 1.0]
* **Headline:** Introducing the new Blue Elephant Jacket, the perfect way to
stay warm and stylish this winter.
* **Body:**
    The Blue Elephant Jacket is made from a soft, durable fabric that will keep
you warm all day long. It features a relaxed fit and a stylish design, so you
can look great while you stay comfortable. The jacket also has a number of
pockets to keep your belongings organized, and a hidden hood for extra warmth
when you need it.
    Order your Blue Elephant Jacket today and see for yourself how comfortable
and stylish it can be.
* **Call to action:**
```

**The `top_k` parameter (range: 0.0 - 40, default 40)**

**What is `top_k`?** `top_k` changes how the model selects tokens for output. A `top_k` of 1 means the selected token is the most probable among all tokens in the model's vocabulary (also called greedy decoding). In contrast, a `top_k` of 3 means that the next token is selected from the top 3 most probable tokens (using temperature). For each token selection step, the `top_k` tokens with the highest probabilities are sampled. Then tokens are further filtered based on `top_p` with the final token selected using temperature sampling.

**How does `top_k` affect the response?** Specify a lower value for less random responses and a higher value for more random responses.

For more information on the `top_k` parameter for text models, please refer to the [documentation on model parameters](#).

```
[18]: prompt_top_k_example = "Write a 2-day itinerary for France."
top_k_val = 1
```

```

response = generation_model.predict(
    prompt=prompt_top_k_example, max_output_tokens=300, temperature=0.9,
    ↪top_k=top_k_val
)

print(f"[top_k = {top_k_val}]")
print(response.text)

```

[top\_k = 1]

Day 1:

- \* Morning: Start your day in Paris with a visit to the Eiffel Tower. Take the elevator to the top for stunning views of the city.
- \* Afternoon: After lunch, visit the Louvre Museum, one of the largest and most famous museums in the world. See some of the most iconic works of art, including the Mona Lisa and Venus de Milo.
- \* Evening: Enjoy a romantic dinner at a traditional French restaurant.

Day 2:

- \* Morning: Take a day trip to Versailles, the former home of French royalty. Explore the opulent palace and gardens.
- \* Afternoon: Visit the Palace of Fontainebleau, another former royal residence. See the beautiful architecture and gardens.
- \* Evening: Enjoy a leisurely dinner at a restaurant in Montmartre, a charming neighborhood in Paris.

[19]: top\_k\_val = 40

```

response = generation_model.predict(
    prompt=prompt_top_k_example,
    max_output_tokens=300,
    temperature=0.9,
    top_k=top_k_val,
)

print(f"[top_k = {top_k_val}]")
print(response.text)

```

[top\_k = 40]

Day 1:

- \* Arrive in Paris and check into your hotel.
- \* Take a walk around the city center and visit some of the famous landmarks, such as the Eiffel Tower, the Louvre Museum, and the Champs-Élysées.
- \* Have dinner at a traditional French restaurant.

Day 2:

- \* Visit the Palace of Versailles and the gardens.
- \* Take a boat trip down the River Seine.
- \* Enjoy a leisurely lunch at a cafe along the river.



\* Visit the Musée d'Orsay, which houses a collection of Impressionist and Post-Impressionist paintings.

This is just a suggested itinerary, of course, and you can customize it to fit your interests and budget. There are many other things to see and do in France, so be sure to do some research before you go to make the most of your trip.

## 1.5 Chat model with chat-bison@001

The chat-bison@001 model lets you have a freeform conversation across multiple turns. The application tracks what was previously said in the conversation. As such, if you expect to use conversations in your application, use the chat-bison@001 model because it has been fine-tuned for multi-turn conversation use cases.

```
[20]: chat_model = ChatModel.from_pretrained("chat-bison@001")

chat = chat_model.start_chat()

print(
    chat.send_message(
        """
Hello! Can you write a 300 word abstract for a research paper I need to write_
about the impact of AI on society?
        """
    )
)
```

Artificial intelligence (AI) is a branch of computer science that deals with the creation of intelligent agents, which are systems that can reason, learn, and act autonomously. AI has the potential to revolutionize many aspects of our lives, from the way we work and play to the way we interact with the world around us.

The impact of AI on society is a complex and multifaceted issue. On the one hand, AI has the potential to solve some of the world's most pressing problems, such as climate change, poverty, and disease. AI can also be used to create new products and services that improve our lives in ways that we can't even imagine.

On the other hand, AI also poses a number of risks to society. AI systems can be used to create autonomous weapons systems that could kill without human intervention. AI can also be used to create surveillance systems that track and monitor our every move.

The future of AI is uncertain. It is possible that AI will be used for good, or it is possible that AI will be used for evil. It is up to us to decide how AI is used, and to ensure that AI is used for the benefit of humanity, not to its detriment.

In this paper, we will explore the potential impact of AI on society. We will discuss the benefits and risks of AI, and we will consider the ways in which AI can be used to solve some of the world's most pressing problems. We will also discuss the challenges that need to be overcome in order to ensure that AI is used for good.

We believe that AI has the potential to make the world a better place. However, we also believe that it is important to be aware of the risks associated with AI. By understanding the potential impact of AI, we can take steps to ensure that AI is used for the benefit of humanity, not to its detriment.

As shown below, the model should respond based on what was previously said in the conversation:

```
[21]: print(
    chat.send_message(
        """
Could you give me a catchy title for the paper?
        """
    )
)
```

The Impact of Artificial Intelligence on Society

### 1.5.1 Advanced Chat model with the SDK

You can also provide a `context` and `examples` to the model. The model will then respond based on the provided context and examples. You can also use `temperature`, `max_output_tokens`, `top_p`, and `top_k`. These parameters should be used when you start your chat with `chat_model.start_chat()`.

For more information on chat models, please refer to the [documentation on chat model parameters](#).

```
[22]: chat = chat_model.start_chat(
    context="My name is Ned. You are my personal assistant. My favorite movies_
are Lord of the Rings and Hobbit.",
    examples=[
        InputOutputTextPair(
            input_text="Who do you work for?",
            output_text="I work for Ned.",
        ),
        InputOutputTextPair(
            input_text="What do I like?",
            output_text="Ned likes watching movies.",
        ),
    ],
    temperature=0.3,
    max_output_tokens=200,
    top_p=0.8,
    top_k=40,
```

```
)
print(chat.send_message("Are my favorite movies based on a book series?"))
```

Yes, they are.

```
[23]: print(chat.send_message("When where these books published?"))
```

The books were published between 1937 and 1949.

## 1.6 Embedding model with `textembedding-gecko@001`

Text embeddings are a dense, often low-dimensional, vector representation of a piece of content such that, if two pieces of content are semantically similar, their respective embeddings are located near each other in the embedding vector space. This representation can be used to solve common NLP tasks, such as:

- **Semantic search:** Search text ranked by semantic similarity.
- **Recommendation:** Return items with text attributes similar to the given text.
- **Classification:** Return the class of items whose text attributes are similar to the given text.
- **Clustering:** Cluster items whose text attributes are similar to the given text.
- **Outlier Detection:** Return items where text attributes are least related to the given text.

Please refer to the [text embedding model documentation](#) for more information.

```
[24]: embedding_model = TextEmbeddingModel.from_pretrained("textembedding-gecko@001")

embeddings = embedding_model.get_embeddings(["What is life?"])

for embedding in embeddings:
    vector = embedding.values
    print(f"Length = {len(vector)}")
    print(vector)
```

Length = 768

```
[0.010562753304839134, 0.04915031045675278, -0.022224493324756622,
0.0208794716745615, 0.024389723315835, 0.010366306640207767,
0.023919280618429184, 0.022391626611351967, -0.031569067388772964,
0.023535897955298424, -0.017047161236405373, -0.014345862902700901,
0.044956106692552567, 0.027327297255396843, -0.03314697742462158,
-0.028214626014232635, -0.035373710095882416, -0.05229683220386505,
0.017105583101511, -0.03780610114336014, -0.07891207933425903,
-0.01173518318682909, -0.01629730500280857, -0.04353305324912071,
0.013023999519646168, -0.10904901474714279, -0.0341256819665432,
-0.0025329082272946835, -0.036971937865018845, -0.027775181457400322,
0.02332289144396782, 0.0052000475116074085, 0.005503748077899218,
0.0047489493153989315, -0.029920609667897224, 0.07563772797584534,
0.0007565636187791824, 0.03501711040735245, 0.02154686115682125,
-0.000812096637673676, 0.06169590726494789, -0.024313345551490784,
0.03736764192581177, -0.0005869767046533525, -0.022872457280755043,
-0.0027376075740903616, -0.020049992948770523, 0.015618567354977131,
```

-0.038495443761348724, -0.0029529621824622154, 0.006082594394683838,  
0.009469639509916306, -0.01644207164645195, 0.06665747612714767,  
-0.005011357367038727, 0.02688094601035118, -0.035699062049388885,  
0.00722217233851552, -0.02506454475224018, -0.004168056882917881,  
-0.04915543273091316, 0.017482219263911247, -0.030127059668302536,  
-0.07361055165529251, -0.02319231443107128, 0.017089763656258583,  
0.03889910876750946, 0.004971345886588097, -0.007473395671695471,  
-0.02317200042307377, 0.027803964912891388, 0.04244052618741989,  
0.05573302507400513, 0.030758505687117577, 0.009500554762780666,  
0.016260754317045212, -0.010920681059360504, 0.051093071699142456,  
0.014400376006960869, -0.04737357795238495, -0.021630441769957542,  
-0.09906578063964844, -0.04526016488671303, -0.09384721517562866,  
0.0203617624938488, -0.010920662432909012, -0.020609041675925255,  
-0.003635745495557785, 0.010114741511642933, 0.053623370826244354,  
-0.06474005430936813, 0.016096388921141624, -0.025225313380360603,  
-0.0019267624011263251, -0.01067234668880701, 0.007826746441423893,  
0.013295858167111874, -0.03250099718570709, -0.0005353756714612246,  
-0.016680261120200157, 0.030173059552907944, -0.046101104468107224,  
0.011562732979655266, -0.025579039007425308, 0.04599408432841301,  
0.07000190019607544, 0.0018490661168470979, 0.006918175611644983,  
-0.0686691403388977, -0.06974631547927856, -0.03236835077404976,  
-0.026718594133853912, 0.0019408509833738208, 0.05204402655363083,  
0.04074086248874664, -0.031760964542627335, 0.00344477710314095,  
-0.026170318946242332, 0.03783925995230675, 0.06050663813948631,  
-0.01712394505739212, -0.0016355578554794192, 0.013527915813028812,  
0.016787230968475342, 0.043276093900203705, -0.016072751954197884,  
0.0005571756628341973, -0.013367726467549801, -0.04776871204376221,  
0.001129272161051631, 0.05411853268742561, -0.013886460103094578,  
0.041834332048892975, -0.04144100472331047, -0.03237384930253029,  
0.043086014688014984, -0.027675528079271317, -0.044598691165447235,  
-0.00594487925991416, 0.06224711984395981, -0.04892174154520035,  
0.03044838272035122, 0.03531208634376526, 0.01580674946308136,  
-0.0012970577226951718, -0.03943518176674843, 0.025564445182681084,  
-0.03516595438122749, -0.07641597092151642, -0.021444709971547127,  
0.020934447646141052, -0.04114843159914017, 0.028761522844433784,  
0.03969430550932884, 0.01649416796863079, 0.03781518712639809,  
0.005376121494919062, 0.052569981664419174, -0.02991773933172226,  
-0.043924685567617416, 0.04530343785881996, 2.5027560695889406e-05,  
-0.03769891336560249, 0.029070040211081505, -0.0015343234408646822,  
-0.008552037179470062, -0.007152904290705919, 0.08237303048372269,  
-0.012235905043780804, 0.01122607383877039, 0.025165848433971405,  
-0.08328612893819809, 0.011945217847824097, 0.0015939214499667287,  
0.08520401269197464, -0.021776320412755013, -0.03578481823205948,  
0.026176629588007927, 0.01096740085631609, 0.004734116140753031,  
-0.020438600331544876, -0.04446989297866821, -0.0040295482613146305,  
0.023052716627717018, -0.036880459636449814, 0.03397374972701073,  
-0.04031934216618538, -0.04730498790740967, -0.0050745015032589436,  
0.08679068833589554, -0.025848431512713432, -0.019488884136080742,

0.017754653468728065, -0.017796259373426437, -0.023471161723136902,  
-0.018227828666567802, 0.04445705562829971, -0.10037199407815933,  
-0.02118160016834736, 0.06375161558389664, -0.008057363331317902,  
0.010159569792449474, 0.0227991733700037, 0.020613307133316994,  
0.006186176091432571, 0.021146763116121292, -0.05968202278017998,  
0.040926363319158554, -0.022863982245326042, -0.03495022654533386,  
0.01183225680142641, 0.01046161912381649, 0.005481143016368151,  
0.011371520347893238, 0.014836390502750874, -0.034276075661182404,  
-0.002988190623000264, 0.025690797716379166, -0.05883420631289482,  
-0.0330558717250824, -0.00978460256010294, 0.04343593493103981,  
-0.008818176575005054, 0.06544901430606842, 0.04110537841916084,  
0.029538320377469063, 0.07326919585466385, -0.011241015046834946,  
0.025485657155513763, -0.010350838303565979, -0.03211556375026703,  
-0.015925990417599678, 0.000886544119566679, 0.0046063438057899475,  
0.012922320514917374, 0.04025505855679512, 0.012885822914540768,  
-0.012207417748868465, -0.004735307767987251, 0.010060984641313553,  
-0.027444403618574142, 0.05696522817015648, -0.018892325460910797,  
0.06066060811281204, -0.013454179279506207, 0.06614361703395844,  
-0.0624060221016407, 0.027977008372545242, 0.004387756809592247,  
0.028211001306772232, 0.014558867551386356, -0.10379569977521896,  
0.0005272417911328375, 0.0014677881263196468, 0.023768983781337738,  
-0.008190426975488663, 0.014077108353376389, 0.013787458650767803,  
-0.0023947367444634438, 0.0023685458581894636, 0.03789709135890007,  
0.012706364504992962, 0.0037047751247882843, 0.09113603085279465,  
-0.04635924845933914, -0.022722741588950157, -0.0071150148287415504,  
0.007771676871925592, 0.01942666433751583, -0.013135110959410667,  
0.028210168704390526, 0.010826485231518745, 0.07257974147796631,  
0.0733613669872284, -0.0006541117909364402, -0.012512844055891037,  
-0.039276015013456345, 0.02811293676495552, 0.023050831630825996,  
0.02523820474743843, 0.04957246035337448, -0.03153982013463974,  
0.007981367409229279, 0.0361814871430397, -0.025585805997252464,  
0.027978206053376198, -0.00865026842802763, -0.0331689678132534,  
-0.09472481161355972, -0.028743796050548553, 0.012484352104365826,  
0.0021317009814083576, -0.03636407107114792, 0.06142040714621544,  
0.027773819863796234, -0.025848736986517906, 0.051048390567302704,  
-0.03207724913954735, -0.011837080121040344, 0.003502538427710533,  
0.007440628949552774, 0.04633283242583275, -0.021616097539663315,  
-0.03967327997088432, -0.02479877881705761, -0.061190344393253326,  
-0.015260206535458565, -0.0352751798927784, 0.021882355213165283,  
0.015894176438450813, 0.0028269870672374964, 0.00022079204791225493,  
-0.013698175549507141, 0.014123217202723026, -0.036067694425582886,  
0.011920701712369919, 0.022252820432186127, 0.005743535701185465,  
-0.0640563890337944, 0.04698688164353371, -0.06404423713684082,  
-0.013042834587395191, 0.10771365463733673, 0.0010355111444368958,  
0.040467679500579834, -0.013715598732233047, -0.009016689844429493,  
0.008000747300684452, 0.10234715789556503, -0.07225427031517029,  
0.008915708400309086, 0.006237425375729799, -0.021482102572917938,  
0.02397948130965233, 0.05135175958275795, 0.015608500689268112,

-0.031944625079631805, 0.0382135808467865, 0.005797197110950947,  
-0.05154971405863762, 0.028097666800022125, -0.02745014801621437,  
-0.05843295902013779, -0.00718856742605567, 0.004217328503727913,  
0.026621710509061813, -0.06410615146160126, -0.027771446853876114,  
-0.03622982278466225, -0.06002693250775337, 0.0028730896301567554,  
0.0009376482921652496, -0.02914406545460224, 0.08022568374872208,  
-0.028879303485155106, 0.026666102930903435, -0.006004952359944582,  
0.001789418631233275, 0.017556296661496162, 0.030322201550006866,  
-0.029420388862490654, 0.018269363790750504, -0.05464306101202965,  
0.025800133123993874, -0.010673671960830688, 0.050776343792676926,  
-0.09446800500154495, 0.0315055325627327, 0.03829166665673256,  
-0.035387031733989716, 0.013971587643027306, 0.012318527325987816,  
-0.015703145414590836, -0.056523360311985016, -0.09211713820695877,  
0.011791898868978024, 0.02190352976322174, 0.045681215822696686,  
0.025594133883714676, -0.02231750264763832, -0.0066884052939713,  
-0.0058244881220161915, 0.02482292428612709, -0.025856031104922295,  
0.0021307426504790783, 0.007752955891191959, 0.04843730479478836,  
0.046476081013679504, -0.008397646248340607, -0.05962646007537842,  
-0.009551334194839, 0.011616517789661884, 0.02129344455897808,  
0.01660742610692978, -0.015061271376907825, 0.10817286372184753,  
-0.022453995421528816, -0.003983873873949051, -0.05261503532528877,  
-0.022377099841833115, 0.044437944889068604, 0.020710069686174393,  
0.013700924813747406, -0.03116024099290371, -0.008501379750669003,  
-0.008547067642211914, 0.010543258860707283, 0.05530229210853577,  
-0.020880352705717087, -0.056557077914476395, -0.00045299093471840024,  
-0.00036045885644853115, 0.016418904066085815, 0.02558765560388565,  
-0.006075927522033453, -0.027869651094079018, -0.016028350219130516,  
-0.024012824520468712, 0.01018358115106821, 0.03202531114220619,  
0.007241277489811182, -0.007335623726248741, 0.013699401170015335,  
0.021170441061258316, 0.02232292853295803, 0.014576029032468796,  
-0.09291727840900421, 0.060506802052259445, -0.055237799882888794,  
-0.03457856923341751, -0.0277749914675951, -0.044067513197660446,  
-0.03341272100806236, -0.011979644186794758, -0.025669127702713013,  
-0.02267066389322281, 0.09566973894834518, 0.0767752006649971,  
0.010824593715369701, 0.006929770577698946, -0.03011220321059227,  
-0.0027981619350612164, -0.0684315413236618, -0.005837030243128538,  
0.03316957876086235, -0.04451394081115723, 0.0027066227048635483,  
-0.05395593121647835, -0.04560602083802223, 0.006837170105427504,  
0.02234707772731781, -0.034541644155979156, 0.011009179055690765,  
0.11418948322534561, -0.0011041957186535, -0.012423255480825901,  
-0.0037464227061718702, -0.035839445888996124, 0.0022794732358306646,  
-0.02418961189687252, -0.015126262791454792, -0.041080981492996216,  
-0.03646824136376381, 0.007073802407830954, -0.008472681045532227,  
-0.007056590635329485, 0.06429030746221542, 0.020093614235520363,  
-0.022911274805665016, 0.025473780930042267, 0.007346852216869593,  
0.04328314587473869, -0.050240155309438705, 0.014898708090186119,  
0.025633497163653374, 0.06909260898828506, 0.03641815483570099,  
0.055241573601961136, 0.021067028865218163, 0.007827181369066238,

-0.06469225883483887, -0.05607031285762787, 0.010224307887256145,  
-0.0395328514277935, -0.017940493300557137, -0.025376077741384506,  
0.015028695575892925, -0.012133865617215633, -0.05708197131752968,  
0.017396116629242897, 0.010713928379118443, 0.0025007263757288456,  
0.013599053025245667, 0.01869947463274002, 0.007097666617482901,  
-0.017439134418964386, 0.032095298171043396, 0.03800633177161217,  
0.003494472708553076, 0.046149060130119324, 0.016314545646309853,  
-0.0063019716180861, 0.026373308151960373, 0.003648587968200445,  
-0.03535480052232742, 0.021219566464424133, 0.00536026805639267,  
-0.06234264746308327, -0.05119526386260986, 0.05295814573764801,  
-0.052322618663311005, 0.04355868324637413, 0.03884781897068024,  
-0.04070327803492546, -0.05103592202067375, -0.06048445403575897,  
0.022206993773579597, -0.0005872139008715749, 0.009120561182498932,  
0.04792363569140434, -0.02097858302295208, -0.04793033376336098,  
-0.020198754966259003, 0.033280182629823685, 0.0031431131064891815,  
0.06413375586271286, 0.005882553290575743, -0.007480444386601448,  
-0.03405573219060898, 0.04223233088850975, -0.03249271214008331,  
-0.03100714646279812, 0.0462421253323555, 0.011582289822399616,  
0.008000263944268227, -0.010868484154343605, 0.046212539076805115,  
0.02901630662381649, 0.005261174403131008, 0.06834172457456589,  
0.01993578113615513, -0.06909119337797165, -0.01722937822341919,  
-0.0073282006196677685, -0.045459046959877014, -0.0057713184505701065,  
0.03357113525271416, 0.0037865471094846725, 0.0036972742527723312,  
0.005904142279177904, 0.05643080919981003, -0.01644083298742771,  
-0.040261246263980865, 0.09400543570518494, 4.431899469636846e-06,  
-0.0490410141646862, 0.03803042694926262, -0.013964969664812088,  
-0.04580939933657646, -0.04047616571187973, -9.045763727044687e-05,  
0.023141689598560333, -0.014303075149655342, -0.02298162318766117,  
0.014015794731676579, -0.057407837361097336, 0.06846775859594345,  
0.024667343124747276, -0.006628675851970911, 0.06792213767766953,  
-0.048522815108299255, -0.013445757329463959, -0.016858110204339027,  
0.0062940688803792, 0.035504620522260666, 0.032475560903549194,  
-0.03314086049795151, -0.03431634604930878, -0.019594095647335052,  
0.0194022748619318, -0.027126198634505272, -0.026301225647330284,  
0.012841426767408848, 0.04247117415070534, -0.0008294038707390428,  
-0.004414733964949846, -0.013424563221633434, -0.05057736486196518,  
-0.015919560566544533, -0.02437496744096279, 0.02126476541161537,  
-0.052758101373910904, -0.00254595885053277, -0.015874166041612625,  
-0.01196883711963892, -0.030309129506349564, 0.010183988139033318,  
0.03598339855670929, 0.022882329300045967, -0.0026234013494104147,  
0.01329271774739027, 0.03793490678071976, 0.040360793471336365,  
0.009328302927315235, 0.023761970922350883, -0.008133891969919205,  
0.003272354369983077, 0.017534887418150902, -0.02026401087641716,  
0.04495016857981682, 0.06324293464422226, 0.035437922924757004,  
0.022442534565925598, -0.016285236924886703, -0.026689207181334496,  
-0.02127818949520588, -0.037313684821128845, 0.01911837048828602,  
0.025560466572642326, 0.016064390540122986, 0.07720045000314713,  
-0.009457273408770561, 0.015530860982835293, -0.06052880361676216,

-0.022617602720856667, 0.029684558510780334, 0.013369569554924965,  
0.022447090595960617, -0.03480006381869316, -0.02661510556936264,  
0.011014388874173164, 4.9809288611868396e-05, -0.035971399396657944,  
-0.04662683233618736, 0.024060387164354324, 0.0023163501173257828,  
-0.028293997049331665, 0.02313932031393051, 0.011891087517142296,  
0.030368415638804436, 0.04482187330722809, 0.022134285420179367,  
-0.013178267516195774, -0.033628739416599274, 0.0025850499514490366,  
-0.03433449566364288, -0.027540232986211777, -0.03729564696550369,  
0.009291116148233414, 0.0032681135926395655, -0.013214519247412682,  
-0.019297223538160324, -0.05181577801704407, -0.07788670808076859,  
0.005753953009843826, -0.002431269269436598, -0.006063435226678848,  
0.056346580386161804, -0.0777992382645607, 0.026173142716288567,  
0.039567336440086365, 0.0007441218476742506, 0.005875684786587954,  
0.05997183546423912, -0.06708808243274689, 0.007471994962543249,  
-0.02053227461874485, 0.007481405511498451, 0.041810907423496246,  
0.02568086050450802, -0.006198380142450333, -0.05383675917983055,  
0.04175911843776703, -0.005462720058858395, -0.05053558945655823,  
0.023038415238261223, -0.014262375421822071, 0.007531439885497093,  
0.013355238363146782, -0.015866558998823166, -0.07595427334308624,  
0.030540091916918755, -0.0801253542304039, -0.023217307403683662,  
0.01443509291857481, -0.004276867024600506, 0.05812528729438782,  
-0.012534553185105324, -0.025629539042711258, -0.01611189916729927,  
-0.04192844033241272, 0.008887489326298237, -0.018155395984649658,  
-0.02993370033800602, -0.024211617186665535, -0.015328830108046532,  
-0.041393671184778214, 0.032823413610458374, -0.001385489129461348,  
0.05429193004965782, -0.027885785326361656, -0.008479449898004532,  
0.0490122027695179, -0.008184432052075863, 0.044856008142232895,  
-0.03555196151137352, -0.07294470816850662, -0.07674667239189148,  
-0.043681032955646515, -0.041639089584350586, -0.02157670259475708,  
0.02400614693760872, 0.024259423837065697, 0.031012538820505142,  
-0.006325466092675924, 0.04981034994125366, -0.014001739211380482,  
0.047177840024232864, 0.011161163449287415, -0.009339496493339539,  
-0.03161853179335594, 0.025517631322145462, 0.02807212620973587,  
-9.692920866655186e-05, -0.061530858278274536, -0.02346394583582878,  
0.020299283787608147, -0.057792555540800095, 0.01165015809237957,  
0.017648736014962196, -0.029611079022288322, -0.030652131885290146,  
0.0484359934926033, 0.05948704108595848, -0.028600767254829407,  
0.0016199287492781878, 0.041371941566467285, 0.021050546318292618,  
0.018185807392001152, -0.055236268788576126, -0.028987789526581764,  
0.06956813484430313, -0.013538234867155552, -0.014316626824438572,  
0.0375431589782238, 0.037525925785303116, 0.012862461619079113,  
-0.004528065212070942, -0.05409426987171173, -0.027900930494070053,  
-0.12323915958404541, 0.02516682632267475, 0.013104746118187904,  
0.004067264497280121, -0.05458975210785866, 0.05783650279045105,  
0.07030092924833298, 0.010288940742611885, -0.02614584006369114,  
-0.015845855697989464, 0.012834896333515644, 0.009806377813220024,  
-0.035357676446437836, 0.01626463793218136, 0.002208224032074213,  
0.009254978969693184, -0.03647759184241295, -0.05006510391831398]



**Embeddings and Pandas DataFrames** If your text is stored in a column of a DataFrame, you can create a new column with the embeddings with the example below.

```
[25]: text = [
    "i really enjoyed the movie last night",
    "so many amazing cinematic scenes yesterday",
    "had a great time writing my Python scripts a few days ago",
    "huge sense of relief when my .py script finally ran without error",
    "O Romeo, Romeo, wherefore art thou Romeo?",
]

df = pd.DataFrame(text, columns=["text"])
df
```

```
[25]:
```

	text
0	i really enjoyed the movie last night
1	so many amazing cinematic scenes yesterday
2	had a great time writing my Python scripts a f...
3	huge sense of relief when my .py script finall...
4	O Romeo, Romeo, wherefore art thou Romeo?

Create a new column, `embeddings`, using the `apply` function in pandas with the embeddings model.

```
[26]: df["embeddings"] = [
    emb.values for emb in embedding_model.get_embeddings(df.text.values)
]
df
```

```
[26]:
```

	text \	embeddings
0	i really enjoyed the movie last night	[-0.01542604435235262, 0.010494514368474483, 0...
1	so many amazing cinematic scenes yesterday	[-0.0445345863699913, 0.03854769468307495, -0...
2	had a great time writing my Python scripts a f...	[-0.002230646787211299, -0.00943742785602808, ...
3	huge sense of relief when my .py script finall...	[-0.0167002584785223, 0.00996268168091774, 0.0...
4	O Romeo, Romeo, wherefore art thou Romeo?	[0.04244288057088852, -0.028772803023457527, 0...

**Comparing similarity of text examples using cosine similarity** By converting text into embeddings, you can compute similarity scores. There are many ways to compute similarity scores, and one common technique is using [cosine similarity](#).

In the example from above, two of the sentences in the `text` column relate to enjoying a *movie*, and the other two relates to enjoying *coding*. Cosine similarity scores should be higher (closer to

1.0) when doing pairwise comparisons between semantically-related sentences, and scores should be lower between semantically-different sentences.

The DataFrame output below shows the resulting cosine similarity scores between the embeddings:

```
[27]: cos_sim_array = cosine_similarity(list(df.embeddings.values))

# display as DataFrame
df = pd.DataFrame(cos_sim_array, index=text, columns=text)
df
```

```
[27]:                                     i really enjoyed the movie
last night \
i really enjoyed the movie last night
1.000000
so many amazing cinematic scenes yesterday
0.719837
had a great time writing my Python scripts a fe...
0.631101
huge sense of relief when my .py script finally...
0.552196
O Romeo, Romeo, wherefore art thou Romeo?
0.459453

                                     so many amazing cinematic
scenes yesterday \
i really enjoyed the movie last night
0.719837
so many amazing cinematic scenes yesterday
1.000000
had a great time writing my Python scripts a fe...
0.622367
huge sense of relief when my .py script finally...
0.564557
O Romeo, Romeo, wherefore art thou Romeo?
0.524301

                                     had a great time writing my
Python scripts a few days ago \
i really enjoyed the movie last night
0.631101
so many amazing cinematic scenes yesterday
0.622367
had a great time writing my Python scripts a fe...
1.000000
huge sense of relief when my .py script finally...
0.739008
O Romeo, Romeo, wherefore art thou Romeo?
```

0.449453

huge sense of relief when my

```
.py script finally ran without error \
i really enjoyed the movie last night
0.552196
so many amazing cinematic scenes yesterday
0.564557
had a great time writing my Python scripts a fe...
0.739008
huge sense of relief when my .py script finally...
1.000000
O Romeo, Romeo, wherefore art thou Romeo?
0.436083
```

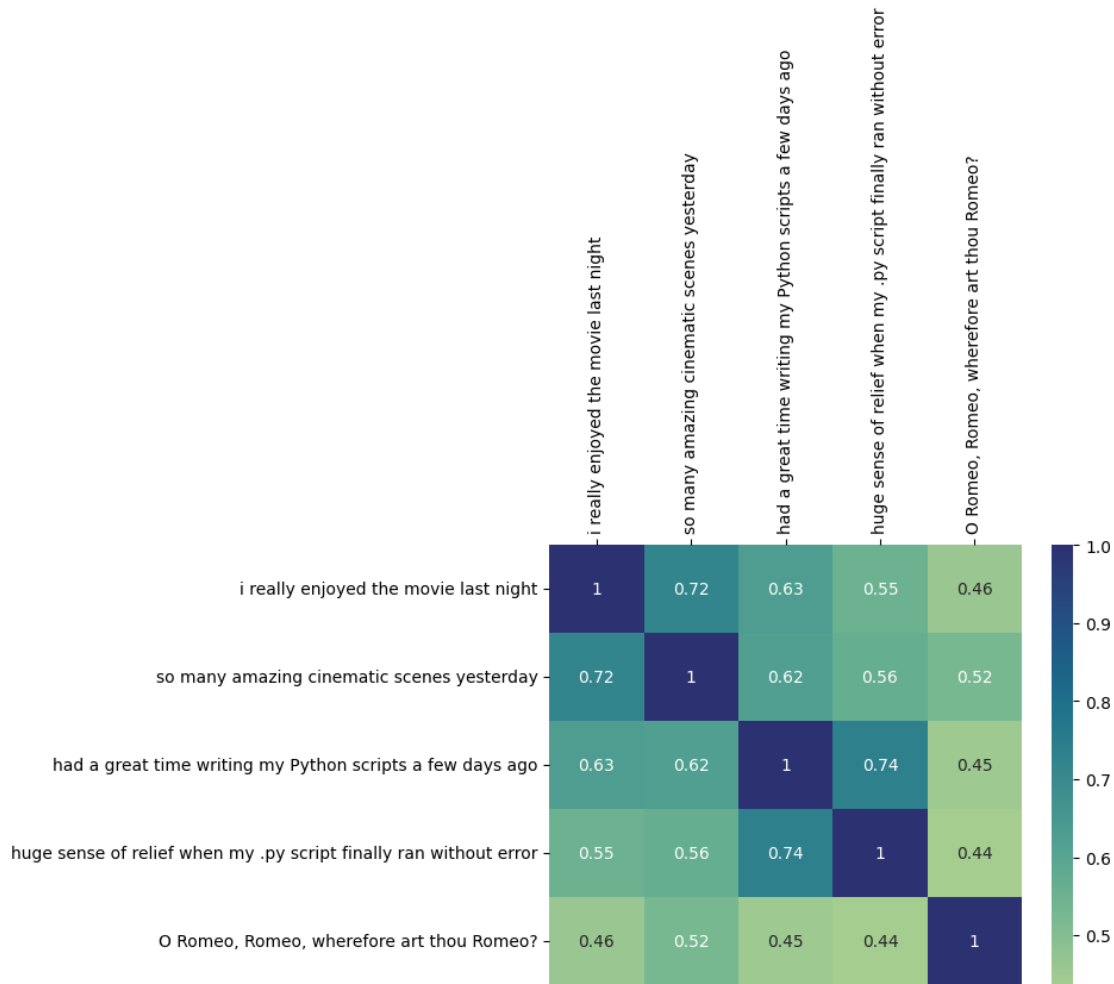
O Romeo, Romeo, wherefore

```
art thou Romeo?
i really enjoyed the movie last night
0.459453
so many amazing cinematic scenes yesterday
0.524301
had a great time writing my Python scripts a fe...
0.449453
huge sense of relief when my .py script finally...
0.436083
O Romeo, Romeo, wherefore art thou Romeo?
1.000000
```

To make this easier to understand, you can use a heatmap. Naturally, text is most similar when they are identical (score of 1.0). The next highest scores are when sentences are semantically similar. The lowest scores are when sentences are quite different in meaning.

```
[28]: ax = sns.heatmap(df, annot=True, cmap="crest")
      ax.xaxis.tick_top()
      ax.set_xticklabels(text, rotation=90)
```

```
[28]: [Text(0.5, 1, 'i really enjoyed the movie last night'),
      Text(1.5, 1, 'so many amazing cinematic scenes yesterday'),
      Text(2.5, 1, 'had a great time writing my Python scripts a few days ago'),
      Text(3.5, 1, 'huge sense of relief when my .py script finally ran without
      error'),
      Text(4.5, 1, 'O Romeo, Romeo, wherefore art thou Romeo?')]
```



## 1.7 Code generation with code-bison@001

The code generation model (Codey) from PaLM API that you will use in this notebook is code-bison@001. It is fine-tuned to follow natural language instructions to generate required code and is suitable for a variety of coding tasks, such as:

- writing functions
- writing classes
- web-appes
- unit tests
- docstrings
- code translations, and many more use-cases.

Currently it supports the following languages: - C++ - C# - Go - GoogleSQL - Java - JavaScript - Kotlin - PHP - Python - Ruby - Rust - Scala - Swift - TypeScript

You can find our more details [here](#).

### 1.7.1 Load model

```
[29]: code_generation_model = CodeGenerationModel.from_pretrained("code-bison@001")
```

### 1.7.2 Model parameters for code-bison@001

You can customize how the PaLM API code generation behaves in response to your prompt by using the following parameters for `code-bison@001`:

- **prefix**: it represents the beginning of a piece of meaningful programming code or a natural language prompt that describes code to be generated.
- **temperature**: higher means more “creative” code responses. range: (0.0 - 1.0, default 0).
- **max\_output\_tokens**: sets the max number of tokens in the output. range: (1 - 2048, default 2048)

### 1.7.3 Hello Codey

```
[30]: prefix = "write a python function to do binary search"

response = code_generation_model.predict(prefix=prefix)

print(response.text)
```

```
```python
def binary_search(array, target):
    """
    Performs a binary search on the given array for the given target.

    Args:
        array: The array to search.
        target: The target value to search for.

    Returns:
        The index of the target value in the array, if found. Otherwise, -1.
    """

    # Check if the array is empty.
    if not array:
        return -1

    # Find the middle index of the array.
    middle_index = len(array) // 2

    # Check if the target value is equal to the middle element.
    if array[middle_index] == target:
        return middle_index

    # Check if the target value is less than the middle element.
    elif target < array[middle_index]:
```

```

    # Recursively search the left half of the array.
    return binary_search(array[:middle_index], target)

# Otherwise, the target value must be greater than the middle element.
else:
    # Recursively search the right half of the array.
    return binary_search(array[middle_index + 1:], target)
...

```

### 1.7.4 Try out your own prompt

Some examples: \* write Go program to extract ip addresses from the text file \* write Java program that can extract pincodes from addresses \* write a standard SQL function that strips all non-alphabet characters from the string and encodes it to utf-8

```

[31]: prefix = """write a python function named as "calculate_cosine_similairty" and
↳three unit \
        tests where it takes two arguments "vector1" and "vector2". \
        It then uses numpy dot function to calculate the dot product of the
↳two vectors. \n
        """

response = code_generation_model.predict(prefix=prefix, max_output_tokens=1024)

print(response.text)

```

```

```python
def calculate_cosine_similarity(vector1, vector2):

    """Calculates the cosine similarity between two vectors.

    Args:
        vector1: A numpy array of numbers.
        vector2: A numpy array of numbers.

    Returns:
        The cosine similarity between the two vectors.
    """

    # Calculate the dot product of the two vectors.
    dot_product = np.dot(vector1, vector2)

    # Calculate the magnitude of the two vectors.
    magnitude1 = np.linalg.norm(vector1)
    magnitude2 = np.linalg.norm(vector2)

    # Calculate the cosine similarity.
    cosine_similarity = dot_product / (magnitude1 * magnitude2)

```

```

    return cosine_similarity

def test_calculate_cosine_similarity():

    """Tests the calculate_cosine_similarity function."""

    # Test 1: Two vectors that are identical.

    vector1 = np.array([1, 2, 3])
    vector2 = np.array([1, 2, 3])

    expected_cosine_similarity = 1.0

    actual_cosine_similarity = calculate_cosine_similarity(vector1, vector2)

    assert actual_cosine_similarity == expected_cosine_similarity

    # Test 2: Two vectors that are perpendicular.

    vector1 = np.array([1, 0, 0])
    vector2 = np.array([0, 1, 0])

    expected_cosine_similarity = 0.0

    actual_cosine_similarity = calculate_cosine_similarity(vector1, vector2)

    assert actual_cosine_similarity == expected_cosine_similarity

    # Test 3: Two vectors that are not identical and not perpendicular.

    vector1 = np.array([1, 2, 3])
    vector2 = np.array([4, 5, 6])

    expected_cosine_similarity = 0.5

    actual_cosine_similarity = calculate_cosine_similarity(vector1, vector2)

    assert actual_cosine_similarity == expected_cosine_similarity
...

```

### 1.7.5 Prompt templates

Prompt templates are useful if you have found a good way to structure your prompt that you can re-use. This can be also be helpful in limiting the open-endedness of freeform prompts. There are many ways to implement prompt templates, and below is just one example using f-strings. This way you can structure the prompts as per the expected functionality of the code.

```
[32]: language = "C++ function"
file_format = "json"
extract_info = "names"
requirements = """
    - the name should be start with capital letters.
    - There should be no duplicate names in the final list.
    """

prefix = f"""Create a {language} to parse {file_format} and extract_
↳{extract_info} with the following requirements: {requirements}.
    """

response = code_generation_model.predict(prefix=prefix, max_output_tokens=1024)

print(response.text)
```

```
```c++
#include <iostream>
#include <string>
#include <vector>
#include <map>

using namespace std;

// This function parses a JSON string and extracts all of the names.
// The names are returned in a vector of strings.
vector<string> getNamesFromJson(string json) {
    // Create a map to store the names.
    map<string, bool> names;

    // Parse the JSON string.
    stringstream ss(json);
    json::value_parser<string> parser;
    while (ss >> parser) {
        // Get the name of the current object.
        string name = parser.get_name();

        // If the name is not already in the map, add it.
        if (!names.count(name)) {
            names[name] = true;
        }
    }

    // Create a vector to store the names.
    vector<string> namesList;

    // Add all of the names to the vector.
    for (map<string, bool>::iterator it = names.begin(); it != names.end(); it++)
```



```

{
    namesList.push_back(it->first);
}

// Return the vector of names.
return namesList;
}

int main() {
    // Get the JSON string from the user.
    cout << "Enter the JSON string: " << endl;
    string json;
    getline(cin, json);

    // Get the names from the JSON string.
    vector<string> names = getNamesFromJson(json);

    // Print the names to the console.
    for (int i = 0; i < names.size(); i++) {
        cout << names[i] << endl;
    }

    return 0;
}
...

```

## 1.8 Code completion with code-gecko@001

Code completion uses the code-gecko foundation model to generate and complete code based on code being written. `code-gecko` completes code that was recently typed by a user.

To learn more about creating prompts for code completion, see [Create prompts for code completion](#).

Code completion API has few more parameters than code generation.

- `prefix`: *required* : For code models, prefix represents the beginning of a piece of meaningful programming code or a natural language prompt that describes code to be generated.
- `suffix`: *optional* : For code completion, suffix represents the end of a piece of meaningful programming code. The model attempts to fill in the code in between the prefix and suffix.
- `temperature`: *required* : Temperature controls the degree of randomness in token selection. Same as for other models. range: (0.0 - 1.0, default 0)
- `maxOutputTokens`: *required* : Maximum number of tokens that can be generated in the response. **range: (1 - 64, default 64)**
- `stopSequences`: *optional* : Specifies a list of strings that tells the model to stop generating text if one of the strings is encountered in the response. The strings are case-sensitive.

```
[33]: code_completion_model = CodeGenerationModel.from_pretrained("code-gecko@001")
```

```
[34]: prefix = """
        def find_x_in_string(string_s, x):
            """

response = code_completion_model.predict(prefix=prefix,
   max_output_tokens=64)

print(response.text)

for i in range(len(string_s)):
    if string_s[i] == x:
        return i
    return -1
```

```
[35]: prefix = """
        def reverse_string(s):
            return s[::-1]
        def test_empty_input_string()
            """

response = code_completion_model.predict(prefix=prefix,
   max_output_tokens=64)

print(response.text)

assert reverse_string("") == ""
def test_one_character_string()
    assert reverse_string("a") == "a"
def test_two_character_string()
    assert reverse_string("ab") == "ba"
```

## 1.9 Code chat with codechat-bison@001

The codechat-bison@001 model lets you have a freeform conversation across multiple turns from a code context. The application tracks what was previously said in the conversation. As such, if you expect to use conversations in your application for code generation, use the codechat-bison@001 model because it has been fine-tuned for multi-turn conversation use cases.

```
[36]: code_chat_model = CodeChatModel.from_pretrained("codechat-bison@001")

code_chat = code_chat_model.start_chat()

print(code_chat.send_message(
    "Please help write a function to calculate the min of two numbers",
    )
)
```

The following Python function calculates the minimum of two numbers:

```

...
def min(a, b):
    """
    Calculates the minimum of two numbers.

    Args:
        a: The first number.
        b: The second number.

    Returns:
        The smaller of the two numbers.
    """

    if a < b:
        return a
    else:
        return b
...

```

This function takes two numbers as input and returns the smaller of the two numbers. The function first compares the two numbers and returns the smaller number. If the two numbers are equal, the function returns either of them.

As shown below, the model should respond based on what was previously asked in the conversation:

```

[37]: print(code_chat.send_message(
        "can you explain the code line by line in bullets?",
    )
)

```

Sure, here is the code explained line by line:

1. The function `min` takes two numbers as input, `a` and `b`.
2. The function first compares the two numbers and returns the smaller number.
3. If the two numbers are equal, the function returns either of them.

You can take another example and ask the model to give more general code suggestion for a specific problem that you are working on.

```

[38]: code_chat = code_chat_model.start_chat()

print(code_chat.send_message(
        "what is the most scalable way to traverse a list in python?",
    )
)

```

The most scalable way to traverse a list in Python is to use the `for` loop. The `for` loop iterates over the elements of a list, and it is the most efficient way to traverse a list.

To use the `for` loop, you first need to declare a variable that will store the current element of the list. Then, you need to use the `in` keyword to check if the current element is in the list. If the current element is in the list, the `for` loop will execute the code inside the loop body.

Here is an example of how to use the `for` loop to traverse a list:

```
...  
# Declare a variable to store the current element of the list  
current_element = None  
  
# Check if the current element is in the list  
if current_element in the_list:  
    # Execute the code inside the loop body  
    pass  
...
```

The `for` loop is the most scalable way to traverse a list because it is the most efficient way to iterate over the elements of a list.

You can continue to ask follow-up questions to the original query.

```
[39]: print(code_chat.send_message(  
        "how would i measure the iteration per second for the following code?",  
    )  
)
```

To measure the iteration per second for the following code, you can use the following steps:

1. Create a timer object.
2. Start the timer.
3. Execute the code.
4. Stop the timer.
5. Calculate the number of iterations per second.

Here is an example of how to measure the iteration per second for the following code:

```
...  
# Create a timer object.  
timer = time.time()  
  
# Start the timer.  
start_time = time.time()  
  
# Execute the code.  
for i in range(1000000):
```

```
    pass

# Stop the timer.
end_time = time.time()

# Calculate the number of iterations per second.
iterations_per_second = 1000000 / (end_time - start_time)
'''
```

The output of this code will be the number of iterations per second.