

✓ WELCOME

This notebook will guide you through two increasingly significant applications in the realm of Generative AI: RAG (Retrieval Augmented Generation) chatbots and text summarization for big text.

Through two distinct projects, you will explore these technologies and enhance your skills. Detailed descriptions of the projects are provided below.

✓ Project 1: Building a Chatbot with a PDF Document (RAG)

In this project, you will develop a chatbot using a provided PDF document from web page. You will utilize the Langchain framework along with a large language model (LLM) such as GPT or Gemini. The chatbot will leverage the Retrieval Augmented Generation (RAG) technique to comprehend the document's content and respond to user queries effectively.

Project Steps:

- **1.PDF Document Upload:** Upload the provided PDF document from web page (<https://aclanthology.org/N19-1423.pdf>) (BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding).
- **2.Chunking:** Divide the uploaded PDF document into smaller segments (chunks). This facilitates more efficient information processing by the LLM.
- **3.ChromaDB Setup:**
 - Save ChromaDB to your Google Drive.
 - Retrieve ChromaDB from your Drive to begin using it in your project.
 - ChromaDB serves as a vector database to store embedding vectors generated from your document.
- **4.Embedding Vectors Creation:**
 - Convert the chunked document into embedding vectors. You can use either GPT or Gemini embedding models for this purpose.
 - If you choose the Gemini embedding model, set "task_type" to "retrieval_document" when converting the chunked document.
- **5.Chatbot Development:**
 - Utilize the `load_qa_chain` function from the Langchain library to build the chatbot.
 - This function will interpret user queries, retrieve relevant information from **ChromaDB**, and generate responses accordingly.

✓ Install Libraries

```
1 !pip install -qU langchain-community # access the tools and components (like DocumentLoader)
2
3 !pip install -qU langchain-openai # utilize OpenAI's models for the "Generation" part
4
5 # !pip install -qU langchain-google-community # Google's LLMs for the "Generation" part
6
7
8 !pip install -qU langchain-chroma # the "Store" step of the "Indexing" pipeline
9
10 !pip install -qU pypdfium2 # the "Load" step of the "Indexing" pipeline, DocumentLoader to
```

```
1 #inspect what exactly is going on inside your chain or agent. The best way to do this is with
2
3 # import getpass
4 # import os
5
6 # os.environ["LANGSMITH_TRACING"] = "true"
7 # os.environ["LANGSMITH_API_KEY"] = getpass.getpass()
8
9 # https://smith.langchain.com/o/b249c21a-9672-4a8d-8cc6-c968dbdc2fef/projects/p/4e4bb92c-e
```



```
.....
```

✓ Access Google Drive

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_re

✓ Entering Your OpenAI or Google Gemini API Key.

```
1 import os
2 from google.colab import userdata
3 os.environ['OPENAI_API_KEY']=userdata.get('OPENAI_API_KEY')
```

✓ Loading PDF Document

```
1
2 # PDF reader function
3 # It reads PDFs page by page and outputs Document objects (LangChain's standard format). R
  a List[Document] → perfect, because the next step (chunking with RecursiveCharacterTextSpl
  expects this format
4 from langchain_community.document_loaders import PyPDFium2Loader
5
6 def read_doc(file_path: str):
7     """
8     Reads a PDF file and returns its pages as LangChain Document objects.
9     Args:
10         file_path (str): Path to the PDF file.
11     Returns:
12         List[Document]: A list of Document objects, one per page.
13     """
14
15     loader = PyPDFium2Loader(file_path)
16     pdf_documents = loader.load()          # Reads PDF page by page
17     return pdf_documents
```

```
1 pdf = read_doc('/content/drive/MyDrive/nlp/bert_article.pdf')
2
3 print(f"-" * 55)
4
5 print(f"The pdf file has {len(pdf)} pages.")
6
7
```

The pdf file has 16 pages.

```
1 pdf
```

[Show hidden output](#)

```
1 pdf[0].page_content # text of first page
2
```

'Proceedings of NAACL-HLT 2019, pages 4171–4186\nMinneapolis, Minnesota, June 2 – June 7, 2019. \nc 2019 Associatio
n for Computational Linguistics\n4171\nBERT: Pre-training of Deep Bidirectional Transformers for\nLanguage Understa
nding\nJacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova\nGoogle AI Language\n{jacobdevlin,mingweichang,kent
onl,kristout}@google.com\nAbstract\nWe introduce a new language representa\x02tion model called BERT, which stands
for\nBidirectional Encoder Representations from\nTransformers. Unlike recent language repre\x02sentation models (Pe
ters et al., 2018a; Rad\x02ford et al., 2018), BERT is designed to pre\x02train deep bidirectional representations
from\nunlabeled text by jointly conditioning on both\nleft and right context in all layers. As a re\x02sult, the pr
e-trained BERT model can be fine\x02tuned with just one additional output layer\nto create state-of-the-art models

```
1 print(pdf[0].page_content[:1250])
```

Proceedings of NAACL-HLT 2019, pages 4171–4186
Minneapolis, Minnesota, June 2 – June 7, 2019.
c 2019 Association for Computational Linguistics

4171
 BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
 Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
 Google AI Language
{jacobdevlin,mingweichang,kentonl,kristout}@google.com
 Abstract
 We introduce a new language representation model called BERT, which stands for Bidirectional Encoder Representations from Transformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to be trained on unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications. BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7% point absolute improvement), MultiNLI

```
1 from IPython.display import Markdown
2
3 Markdown(pdf[2].page_content)
```

4173 BERT BERT [CLS] E1 [SEP] ... EN E1' ... EM' C T1 [SEP] ... TN T1' ... TM' [CLS] Tok 1 [SEP] ... Tok N Tok 1 ... TokM Question Paragraph Start/End Span BERT [CLS] E1 [SEP] ... EN E1' ... EM' C T1 [SEP] ... TN T1' ... TM' [CLS] Tok 1 [SEP] ... Tok N Tok 1 ... TokM Masked Sentence A Masked Sentence B Pre-training Fine-Tuning NSP Mask LM Mask LM Unlabeled Sentence A and B Pair SQuAD Question Answer Pair MNLI NER Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers). Training and auto-encoder objectives have been used for pre-training such models (Howard and Ruder, 2018; Radford et al., 2018; Dai and Le, 2015). 2.3 Transfer Learning from Supervised Data There has also been work showing effective transfer from supervised tasks with large datasets, such as natural language inference (Conneau et al., 2017) and machine translation (McCann et al., 2017). Computer vision research has also demonstrated the importance of transfer learning from large pre-trained models, where an effective recipe is to fine-tune models pre-trained with ImageNet (Deng et al., 2009; Yosinski et al., 2014). 3 BERT We introduce BERT and its detailed implementation in this section. There are two steps in our framework: pre-training and fine-tuning. During pre-training, the model is trained on unlabeled data over different pre-training tasks. For fine-tuning, the BERT model is first initialized with the pre-trained parameters, and all of the parameters are fine-tuned using labeled data from the downstream tasks. Each downstream task has separate fine-tuned models, even though they are initialized with the same pre-trained parameters. The question-answering example in Figure 1 will serve as a running example for this section. A distinctive feature of BERT is its unified architecture across different tasks. There is minimal difference between the pre-trained architecture and the final downstream architecture. Model Architecture BERT's model architecture is a multi-layer bidirectional Transformer encoder based on the original implementation described in Vaswani et al. (2017) and released in the tensorflow/tensor2tensor library.¹ Because the use of Transformers has become common and our implementation is almost identical to the original, we will omit an exhaustive background description of the model architecture and refer readers to Vaswani et al. (2017) as well as excellent guides such as "The Annotated Transformer."² In this work, we denote the number of layers (i.e., Transformer blocks) as L, the hidden size as H, and the number of self-attention heads as A. 3 We primarily report results on two model sizes: BERTBASE (L=12, H=768, A=12, Total Parameters=110M) and BERTLARGE (L=24, H=1024, A=16, Total Parameters=340M). BERTBASE was chosen to have the same model size as OpenAI GPT for comparison purposes. Critically, however, the BERT Transformer uses bidirectional self-attention, while the GPT Transformer uses constrained self-attention where every token can only attend to context to its left.⁴ 1 <https://github.com/tensorflow/tensor2tensor> 2 <http://nlp.seas.harvard.edu/2018/04/03/attention.html> 3 In all cases we set the feed-forward/filter size to be 4H, i.e., 3072 for the H = 768 and 4096 for the H

```
1
2 # all pages
3
4
5
6
7
8
9
10 from IPython.display import Markdown, display
11
12 for i, d in enumerate(pdf, start=1):
13     display(Markdown(f"### Page {i}\n\n{d.page_content}"))
14
15
16
```

Show hidden output

Document Splitter

```
1 from langchain.text_splitter import RecursiveCharacterTextSplitter #CharacterTextSplitter
2
3
4 def chunk_data(docs, chunk_size=1000, chunk_overlap=200):
5     text_splitter = RecursiveCharacterTextSplitter(chunk_size=chunk_size,
6                                                    chunk_overlap=chunk_overlap)
```

```

7 # each chunk will be a maximum of 1000 characters long# each chunk will share a 200-char
8 pdf = text_splitter.split_documents(docs)
9 return pdf
10
11
12
13 # This code splits documents into chunks using the RecursiveCharacterTextSplitter class f
14
15 # A function named chunk_data is defined, which takes a document or a collection of docum
16 # chunk_size and chunk_overlap. chunk_size specifies the maximum number of characters in
17 # overlap between consecutive chunks.
18
19 # The function divides the documents into chunks based on these parameters using the Recu
20
21 # As a result, the documents are segmented into chunks of specified sizes, and these chunk
22
23 # The chunk_overlap parameter is used to specify the sharing of characters between consec
24 # the end of one chunk reappear at the beginning of the next chunk. This prevents the los

```

```

1 pdf_doc = chunk_data(docs=pdf)
2
3 len(pdf_doc)

```

83

```
1 pdf_doc[0:2]
```

[Document(metadata={'producer': 'pdfTeX-1.40.18', 'creator': 'LaTeX with hyperref package', 'creationdate': '2019-04-29T17:36:03+00:00', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'subject': 'N19-1 2019', 'keywords': '', 'moddate': '2019-04-29T17:36:03+00:00', 'source': '/content/drive/MyDrive/nlp/bert_article.pdf', 'total_pages': 16, 'page': 0}, page_content='Proceedings of NAACL-HLT 2019, pages 4171-4186\\nMinneapolis, Minnesota, June 2 - June 7, 2019. \\nc 2019 Association for Computational Linguistics\\n4171\\nBERT: Pre-training of Deep Bidirectional Transformers for\\nLanguage Understanding\\nJacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova\\nGoogle AI Language\\n[fjacobdevlin,mingweichang,kentonl,kristout@google.com](#)\\nAbstract\\nWe introduce a new language representa\\x02tion model called BERT, which stands for\\nBidirectional Encoder Representations from\\nTransformers. Unlike recent language repre\\x02sentation models (Peters et al., 2018a; Rad\\x02ford et al., 2018), BERT is designed to pre\\x02train deep bidirectional representations from\\nunlabeled text by jointly conditioning on both\\nleft and right context in all layers. As a re\\x02sult, the pre-trained BERT model can be fine\\x02tuned with just one additional output layer\\nto create state-of-the-art models for a wide\\nrange of tasks, such as question answering and'), Document(metadata={'producer': 'pdfTeX-1.40.18', 'creator': 'LaTeX with hyperref package', 'creationdate': '2019-04-29T17:36:03+00:00', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'subject': 'N19-1 2019', 'keywords': '', 'moddate': '2019-04-29T17:36:03+00:00', 'source': '/content/drive/MyDrive/nlp/bert_article.pdf', 'total_pages': 16, 'page': 0}, page_content='to create state-of-the-art models for a wide\\nrange of tasks, such as question answering and\\nlanguage inference, without substantial task\\x02specific architecture modifications.\\nBERT is conceptually simple and empirically\\npowerful. It obtains new state-of-the-art re\\x02sults on eleven natural language processing\\ntasks, including pushing the GLUE score to\\n80.5% (7.7% point absolute improvement),\\nMultiNLI accuracy to 86.7% (4.6% absolute\\nimprovement), SQuAD v1.1 question answer\\x02ing Test F1 to 93.2 (1.5 point absolute im\\x02provement) and SQuAD v2.0 Test F1 to 83.1\\n(5.1 point absolute improvement).\\n1 Introduction\\nLanguage model pre-training has been shown to\\nbe effective for improving many natural language\\nprocessing tasks (Dai and Le, 2015; Peters et al.,\\n2018a; Radford et al., 2018; Howard and Ruder,\\n2018). These include sentence-level tasks such as\\nnatural language inference (Bowman et al., 2015;\\nWilliams et al., 2018) and paraphrasing (Dolan')]

```
1 Markdown(pdf_doc[82].page_content)
```

MASK means that we replace the target token with the [MASK] symbol for MLM; SAME means that we keep the target token as is; RND means that we replace the target token with another random token. The numbers in the left part of the table represent the probabilities of the specific strategies used during MLM pre-training (BERT uses 80%, 10%, 10%). The right part of the paper represents the Dev set results. For the feature-based approach, we concatenate the last 4 layers of BERT as the features, which was shown to be the best approach in Section 5.3. From the table it can be seen that fine-tuning is surprisingly robust to different masking strategies. However, as expected, using only the MASK strategy was problematic when applying the feature based approach to NER. Interestingly, using only the RND strategy performs much worse than our strategy as well.

- ✓ 1. Creating A Embedding Model
- 2. Convert the Each Chunk of The Split Document to Embedding Vectors
- 3. Storing of The Embedding Vectors to Vectorstore
- 4. Save the Vectorstore to Your Drive

```

1 from langchain_openai import OpenAIEmbeddings
2 # from langchain.embeddings import OpenAIEmbeddings
3
4
5 embeddings = OpenAIEmbeddings(model="text-embedding-3-large",
6                               dimensions=3072) # dimensions=256, 1024, 3072
7 embeddings
8
9 # As the embedding model, we use Openai's latest introduced text-embedding-3-large model.
10 # dimensions of text-embedding-3-large are 256, 1024 and 3072

```

```

OpenAIEmbeddings(client=<openai.resources.embeddings.Embeddings object at 0x7e260d54f800>, async_client=
<openai.resources.embeddings.AsyncEmbeddings object at 0x7e260908dfd0>, model='text-embedding-3-large',
dimensions=3072, deployment='text-embedding-ada-002', openai_api_version=None, openai_api_base=None,
openai_api_type=None, openai_proxy=None, embedding_ctx_length=8191, openai_api_key=SecretStr('*****'),
openai_organization=None, allowed_special=None, disallowed_special=None, chunk_size=1000, max_retries=2,
request_timeout=None, headers=None, tiktoken_enabled=True, tiktoken_model_name=None, show_progress_bar=False,
model_kwargs={}, skip_empty=False, default_headers=None, default_query=None, retry_min_seconds=4,
retry_max_seconds=20, http_client=None, http_async_client=None, check_embedding_ctx_length=True)

```

```

1 # the "Embed" step of the RAG (Retrieval Augmented Generation) pipeline.
2 # This is the third stage in the indexing process, after "Load" and "Split".

```

```

1 # store the embedding vectors.

```

```

1 from langchain_chroma import Chroma
2
3 # VectorStore object
4 # Embeds each Document in pdf_doc using the embeddings function. Stores vectors + metadata
5 index = Chroma.from_documents(documents=pdf_doc,
6                               embedding=embeddings,
7                               persist_directory="./chroma_store_bert_z") # persist_directory
8
9 # Create a retriever object from our Chroma vector store (index). Wraps the vector store as a retriever
10 retriever = index.as_retriever() # for using agents #index.as_retriever(search_kwargs={"k": 4})
11
12

```

```

1 # refinement of the "Retrieval" step in a RAG system.
2 # to retrieve the four most relevant pieces of information from the indexed PDF
3
4 retriever = index.as_retriever(search_kwargs={"k": 4})
5
6

```

```

1 retriever.invoke("Describe BERT model architecture")

```

```

[Document(id='10fffc67d-9f78-408d-9d50-9acf62d59031', metadata={'subject': 'N19-1 2019', 'producer': 'pdfTeX-
1.40.18', 'creationdate': '2019-04-29T17:36:03+00:00', 'page': 2, 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kento
Lee ; Kristina Toutanova', 'source': '/content/drive/MyDrive/nlp/bert_article.pdf', 'keywords': '', 'moddate':
'2019-04-29T17:36:03+00:00', 'total_pages': 16, 'creator': 'LaTeX with hyperref package', 'title': 'BERT: Pre-
training of Deep Bidirectional Transformers for Language Understanding'}, page_content='framework: pre-training an
fine-tuning. During pre-training, the model is trained on unlabeled data over different pre-training tasks.
For fine-tuning, the BERT model is first initialized with the pre-trained parameters, and all of the paramete
rs are fine-tuned using labeled data from the downstream tasks. Each downstream task has separate fine-tuned
models, even though they are initialized with the same pre-trained parameters. The question-answering exam
ple in Figure 1 will serve as a running example for this section. A distinctive feature of BERT is its unifie
d architecture across different tasks. There is minimal difference between the pre-trained architecture and
the final downstream architecture. Model Architecture BERT's model architecture is a multi-layer bidirec
tional Transformer encoder based on the original implementation described in Vaswani et al. (2017) and re
leased in the tensorflow/tensorflow library. Because of the use of',

```

```
Document(id='db0bae0a-9e10-424b-b65e-c29e35d64163', metadata={'creationdate': '2019-04-29T17:36:03+00:00',
'keywords': '', 'producer': 'pdfTeX-1.40.18', 'moddate': '2019-04-29T17:36:03+00:00', 'creator': 'LaTeX with
hyperref package', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding',
'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'source':
'/content/drive/MyDrive/nlp/bert_article.pdf', 'subject': 'N19-1 2019', 'page': 2, 'total_pages': 16},
page_content='4173\nBERT BERT\nE[CLS] E1\nE[SEP] ... EN\nE1' ... EM'\nC T1 T[SEP] ... TN\nT1' ... TM'\n[CLS] Tok
[SEP] ... Tok N Tok 1 ... TokM\nQuestion Paragraph\nStart/End Span\nBERT\nE[CLS] E1\nE[SEP] ... EN\nE1' ... EM'\n
T1 T[SEP] ... TN\nT1' ... TM'\n[CLS] Tok 1 [SEP] ... Tok N Tok 1 ... TokM\nMasked Sentence A Masked Sentence
B\nPre-training Fine-Tuning\nNSP Mask LM Mask LM\nUnlabeled Sentence A and B Pair \nSQuAD\nQuestion Answer
Pair\nMNLI NER\nFigure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the
same architec\02tures are used in both pre-training and fine-tuning. The same pre-trained model parameters are
used to initialize\nmodels for different down-stream tasks. During fine-tuning, all parameters are fine-tuned.
[CLS] is a special\nsymbol added in front of every input example, and [SEP] is a special separator token (e.g.
separating ques\02tions/answers).\ning and auto-encoder objectives have been used\nfor pre-training such models
(Howard and Ruder,')',
Document(id='b94e7da0-751b-47e0-ab26-94e0e1289db6', metadata={'moddate': '2019-04-29T17:36:03+00:00', 'source':
'/content/drive/MyDrive/nlp/bert_article.pdf', 'total_pages': 16, 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kento
Lee ; Kristina Toutanova', 'keywords': '', 'creationdate': '2019-04-29T17:36:03+00:00', 'creator': 'LaTeX with
hyperref package', 'subject': 'N19-1 2019', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for
Language Understanding', 'producer': 'pdfTeX-1.40.18', 'page': 0}, page_content='to create state-of-the-art models
for a wide\nrange of tasks, such as question answering and\nlanguage inference, without substantial
task\02specific architecture modifications.\nBERT is conceptually simple and empirically\npowerful. It obtains ne
state-of-the-art re\02sults on eleven natural language processing\ntasks, including pushing the GLUE score
to\n80.5% (7.7% point absolute improvement),\nMultiNLI accuracy to 86.7% (4.6% absolute\nimprovement), SQuAD v1.1
question answer\02ing Test F1 to 93.2 (1.5 point absolute im\02provement) and SQuAD v2.0 Test F1 to 83.1\n(5.1
point absolute improvement).\n1 Introduction\nLanguage model pre-training has been shown to\nbe effective for
improving many natural language\nprocessing tasks (Dai and Le, 2015; Peters et al.,\n2018a; Radford et al., 2018;
Howard and Ruder,\n2018). These include sentence-level tasks such as\nnatural language inference (Bowman et al.,
2015;\nWilliams et al., 2018) and paraphrasing (Dolan')',
Document(id='01826491-be9a-4ad6-8e8b-d51e7bd6fa6f', metadata={'total_pages': 16, 'moddate': '2019-04-
29T17:36:03+00:00', 'creationdate': '2019-04-29T17:36:03+00:00', 'keywords': '', 'source':
'/content/drive/MyDrive/nlp/bert_article.pdf', 'page': 0, 'subject': 'N19-1 2019', 'author': 'Jacob Devlin ; Ming-
Wei Chang ; Kenton Lee ; Kristina Toutanova', 'producer': 'pdfTeX-1.40.18', 'creator': 'LaTeX with hyperref
package', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding'},
page_content='example, in OpenAI GPT, the authors use a left-to\02right architecture, where every token can only
at\02tend to previous tokens in the self-attention layers\nof the Transformer (Vaswani et al., 2017). Such
re\02strictions are sub-optimal for sentence-level tasks,\nand could be very harmful when applying fine\02tuning
based approaches to token-level tasks such\nas question answering, where it is crucial to incor\02porate context
from both directions.\nIn this paper, we improve the fine-tuning based\napproaches by proposing BERT:
Bidirectional\nEncoder Representations from Transformers.\nBERT alleviates the previously mentioned
unidi\02rectionality constraint by using a “masked lan\02guage model” (MLM) pre-training objective, in\02spired
by the Cloze task (Taylor, 1953). The\nmasked language model randomly masks some of\nthe tokens from the input, an
the objective is to\npredict the original vocabulary id of the masked')]
```

✓ Load Vectorstore(index) From Your Drive

```
1 # Initiates the Chroma class, which serves as our vector database.
2
3 loaded_from_drive_index = Chroma(persist_directory="/content/drive/MyDrive/nlp/bert_article.pdf",
4                                   embedding_function=embeddings)
5 #/chroma_store_bert_z
```

```
1 load_retriver = loaded_from_drive_index.as_retriever(search_kwargs={"k": 5})
2 load_retriver
```

```
VectorStoreRetriever(tags=['Chroma', 'OpenAIEmbeddings'], vectorstore=<langchain_chroma.vectorstores.Chroma object
at 0x7e260cfa34d0>, search_kwargs={'k': 5})
```

```
1
2 # bonus ... shorter versions of above
3
4
5 def retrieve_query(query, k=4):
6     retriever = index.as_retriever(search_kwargs={"k": k}) # loaded_index
7     return retriever.invoke(query)
8
9 # The query we ask is first converted into an embedding vector. Then, using the cosine sim
```

```
1 our_query = "Describe BERT model architecture"
2
3 doc_search = retrieve_query(our_query, k=5) # first 5 most similar texts are returned
4 doc_search
```

```
[Document(id='10ffc67d-9f78-408d-9d50-9acf62d59031', metadata={'keywords': '', 'creator': 'LaTeX with hyperref
package', 'moddate': '2019-04-29T17:36:03+00:00', 'subject': 'N19-1 2019', 'total_pages': 16, 'creationdate':
'2019-04-29T17:36:03+00:00', 'producer': 'pdfTeX-1.40.18', 'page': 2, 'source':
```

```

'/content/drive/MyDrive/nlp/bert_article.pdf', 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding'},
page_content='framework: pre-training and fine-tuning. During pre-training, the model is trained on unlabeled data over different pre-training tasks. For fine-tuning, the BERT model is first initialized with the pre-trained parameters, and all of the parameters are fine-tuned using labeled data from the downstream tasks. Each downstream task has separate fine-tuned models, even though they are initialized with the same pre-trained parameters. The question-answering example in Figure 1 will serve as a running example for this section. A distinctive feature of BERT is its unified architecture across different tasks. There is minimal difference between the pre-trained architecture and the final downstream architecture. Model Architecture BERT's model architecture is a multi-layer bidirectional Transformer encoder based on the original implementation described in Vaswani et al. (2017) and released in the TensorFlow library.1 Because the use of TensorFlow is not required, we provide the implementation in the PyTorch library. Document(id='db0bae0a-9e10-424b-b65e-c29e35d64163', metadata={'page': 2, 'total_pages': 16, 'source': '/content/drive/MyDrive/nlp/bert_article.pdf', 'keywords': '', 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'creationdate': '2019-04-29T17:36:03+00:00', 'subject': 'N19-1 2019', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', 'creator': 'LaTeX with hyperref package', 'producer': 'pdfTeX-1.40.18', 'moddate': '2019-04-29T17:36:03+00:00'}, page_content='4173\nBERT\nBERT\n[CLS] E1\nE[SEP] ... EN\nE1' ... EM\nC T1 T[SEP] ... TM\nT1' ... TM\n[CLS] Tok 1 [SEP] ... Tok N Tok 1 ... TokM\nQuestion Paragraph\nStart/End Span\nBERT\n[CLS] E1\nE[SEP] ... EN\nE1' ... EM\nC T1 T[SEP] ... TM\nT1' ... TM\n[CLS] Tok 1 [SEP] ... Tok N Tok 1 ... TokM\nMasked Sentence A Masked Sentence B\nPre-training Fine-Tuning\nNSP Mask LM Mask LM\nUnlabeled Sentence A and B Pair\nSQuAD\nQuestion Answer Pair\nMNLI NER\nFigure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different downstream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers). Encoding and auto-encoder objectives have been used for pre-training such models (Howard and Ruder, 2018). Document(id='b94e7da0-751b-47e0-ab26-94e0e1289db6', metadata={'creationdate': '2019-04-29T17:36:03+00:00', 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', 'total_pages': 16, 'moddate': '2019-04-29T17:36:03+00:00', 'producer': 'pdfTeX-1.40.18', 'subject': 'N19-1 2019', 'keywords': '', 'page': 0, 'source': '/content/drive/MyDrive/nlp/bert_article.pdf', 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'creator': 'LaTeX with hyperref package'}, page_content='to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications. BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to 80.5% (7.7 point absolute improvement), MultiNLI accuracy to 86.7% (4.6% absolute improvement), SQuAD v1.1 question answer Test F1 to 93.2 (1.5 point absolute improvement) and SQuAD v2.0 Test F1 to 83.1 (5.1 point absolute improvement).1 Introduction Language model pre-training has been shown to be effective for improving many natural language processing tasks (Dai and Le, 2015; Peters et al., 2018a; Radford et al., 2018; Howard and Ruder, 2018). These include sentence-level tasks such as natural language inference (Bowman et al., 2015; Williams et al., 2018) and paraphrasing (Dolan, 2018). Document(id='01826491-be9a-4ad6-8e8b-d51e7bd6fa6f', metadata={'producer': 'pdfTeX-1.40.18', 'keywords': '', 'source': '/content/drive/MyDrive/nlp/bert_article.pdf', 'total_pages': 16, 'title': 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', 'creator': 'LaTeX with hyperref package', 'moddate': '2019-04-29T17:36:03+00:00', 'author': 'Jacob Devlin ; Ming-Wei Chang ; Kenton Lee ; Kristina Toutanova', 'page': 0, 'creationdate': '2019-04-29T17:36:03+00:00', 'subject': 'N19-1 2019'}, page_content='example, in OpenAI GPT, the authors use a left-to-right architecture, where every token can only attend to previous tokens in the self-attention layers of the Transformer (Vaswani et al., 2017). Such restrictions are sub-optimal for sentence-level tasks, and could be very harmful when applying fine-tuning based approaches to token-level tasks such as question answering, where it is crucial to incorporate context from both directions. In this paper, we improve the fine-tuning based approaches by proposing BERT: Bidirectional Encoder Representations from Transformers. BERT alleviates the previously mentioned unidirectionality constraint by using a “masked language model” (MLM) pre-training objective, inspired by the Cloze task (Taylor, 1953). The masked language model randomly masks some of the tokens from the input, and the objective is to predict the original vocabulary id of the masked tokens.

```

```

1
2 # for now we only saw our pdf relevant parts, llm was not involved in ...
3

```

Generating an Answer Based on The Similar Chunks

```

1 # vital bridge between the "Retrieval" and "Generation" phases of our RAG pipeline
2 # relevant information retrieved from our vector store (doc_search) and formats it, along with

```

Pipeline For RAG

```

1 # 1. chatbot qa yapisi prompt kismi
2
3
4 from langchain.prompts import ChatPromptTemplate
5 from langchain_openai import ChatOpenAI
6 from langchain_core.output_parsers import StrOutputParser
7 from langchain.chains.question_answering import load_qa_chain
8 from IPython.display import Markdown
9
10 # Role-based chat prompt
11 chat_template = ChatPromptTemplate.from_messages([

```

```

12     (
13         "system",
14         "You are a scientific QA assistant. Use only the provided context. "
15         "Always cite page numbers. If uncertain, say 'I don't know.' "
16         "Keep answers short, neutral, and evidence-based."
17     ),
18     (
19         "system",
20         "Return output as a markdown block with sections:\n"
21         "# Answer\n"
22         "## Key Points (with page cites)\n"
23         "### Sources\n"
24         "### Confidence"
25     ),
26     (
27         "human",
28         # We inject both the raw context and a compact citations map the model can quote.
29         "Context:\n{context}\n\n"
30         "Citations Map (use these for page refs):\n{citations}\n\n"
31         "Question: {question}"
32     )
33 ])
34
35 # For map_reduce you also need a combine prompt (roles preserved)
36 combine_prompt = ChatPromptTemplate.from_messages([
37     ("system",
38         "You are a scientific QA assistant. Combine the partial answers below into a single,
39         "Resolve conflicts conservatively; if evidence is insufficient, say 'I don't know.' "
40         "Keep the tone neutral and cite page numbers where available."),
41     ("system",
42         "Return markdown with sections:\n"
43         "### Answer\n"
44         "### Key Points (with page cites)\n"
45         "### Sources\n"
46         "### Confidence"),
47     ("human",
48         "Partial answers from multiple chunks:\n{summaries}\n\nProduce the final answer.")
49 ])
50
51
52
53     # 1. chatbot qa structure
54
55
56 #get_answers, consolidates all the previous steps of the Retrieval Augmented Generation (I
57
58 #This function represents the complete workflow from receiving a user's question to gener
59
60
61 def get_answers(query, k=4):
62     from langchain_openai import ChatOpenAI
63     from langchain_core.output_parsers import StrOutputParser
64     from IPython.display import Markdown
65
66     # 1) Retrieve top-k documents
67     docs = retriever.invoke(query, k=k)
68
69     # 2) Build context
70     context = "\n\n".join(d.page_content for d in docs)
71
72     # 3) Build a citations map the model can quote in output
73     # (ensures it actually has page numbers available)
74     citations = "\n".join(
75         f"[{i+1}] {d.metadata.get('source', 'unknown')} - p.{d.metadata.get('page', '?')}"
76         for i, d in enumerate(docs)
77     )
78
79     # 4) LLM
80     llm = ChatOpenAI(model_name="gpt-4o-mini", temperature=0, max_tokens=256)
81     # max_tokens=1024)
82

```



```
83     # 5) Chain
84     chain = chat_template | llm | StrOutputParser()
85
86     # 6) Invoke
87     answer = chain.invoke({"question": query, "context": context, "citations": citations})
88
89     # 7) Return (avoid adding another Sources block – model already includes it)
```

```
1 our_query = "What are the two pre-training tasks used in BERT, and what does each one train on?"
2
3 answer = get_answers(our_query, k=5)
4 answer
```

Answer

The two pre-training tasks used in BERT are:

1. **Masked Language Model (MLM)**: This task involves masking a percentage of the input tokens and training the model to predict these masked tokens based on their context. It allows the model to learn bidirectional representations by considering both left and right context.
2. **Next Sentence Prediction (NSP)**: This task trains the model to predict whether a given sentence B follows sentence A in the original text. It helps the model understand the relationship between sentence pairs.

Key Points (with page cites)

- BERT uses two unsupervised tasks for pre-training: Masked Language Model (MLM) and Next Sentence Prediction (NSP) (p.3).
- MLM trains the model to predict masked tokens in a sentence (p.3).
- NSP trains the model to determine if one sentence follows another (p.3).

Sources

- [1] /content/drive/MyDrive/nlp/bert_article.pdf — p.3

Confidence

High

```
1 our_query = "What ROUGE-L score does BERT report on CNN/DailyMail abstractive summarization?"
2
3 answer = get_answers(our_query, k=5)
4 answer
```

I don't know.

```
1 our_query = "In the context of fine-tuning, what minimal number of parameters need to be randomly initialized?"
2
3 answer = get_answers(our_query, k=5)
4 answer
5
6
```

Answer

The minimal number of parameters that need to be randomly initialized when incorporating BERT with an extra output layer is very small, as the task-specific models are formed by adding only one additional output layer to the pre-trained model.

Key Points (with page cites)

- The task-specific models incorporate BERT with one additional output layer, requiring only a minimal number of parameters to be learned from scratch (page 4184).

Sources

- [1] /content/drive/MyDrive/nlp/bert_article.pdf — p.8
- [2] /content/drive/MyDrive/nlp/bert_article.pdf — p.12
- [3] /content/drive/MyDrive/nlp/bert_article.pdf — p.13
- [4] /content/drive/MyDrive/nlp/bert_article.pdf — p.13
- [5] /content/drive/MyDrive/nlp/bert_article.pdf — p.2

Confidence

High

```

1                                     # some question
2 # Relevant
3
4 # 1) What are the two pre-training tasks used in BERT, and what does each one train the mo
5 # 2) How do BERTBASE and BERTLARGE differ in L (layers), H (hidden size), A (attention head
6
7
8 # Irrelevant
9
10 # 3) What is the chemical formula of water?
11 # 4) What is the tallest mountain in South America?
12 # 5) What ROUGE-L score does BERT report on CNN/DailyMail abstractive summarization using
13
14
15
16

```

```

1 #                                     works well but may not be cost effective
2
3 from langchain.chains import RetrievalQA
4 from langchain_core.prompts import ChatPromptTemplate
5 from langchain_openai import ChatOpenAI
6
7 llm = ChatOpenAI(model="gpt-4o-mini", temperature=0)
8
9
10 QUESTION_PROMPT = ChatPromptTemplate.from_template(
11     """You are an data scientist assistant answering a user question using ONLY the provi
12
13 REQUIREMENTS
14 - If the context is insufficient, say: "I don't know based on the provided context."
15 - Be precise and neutral; do not speculate.
16 - Keep the main answer within 3-8 sentences unless explicitly asked for more.
17 - Cite support in brackets using metadata if available.
18
19 CONTEXT
20 {context}
21
22 QUESTION
23 {question}
24
25 ANSWER"""
26 )
27
28 COMBINE_PROMPT = ChatPromptTemplate.from_template(
29     """You receive multiple candidate answers synthesized from different context chunks.
30
31 TASK
32 - Produce ONE final answer.
33 - Prefer statements directly supported by the context; drop anything weakly supported.
34 - If candidates conflict and evidence is unclear, state the uncertainty or contraversarie
35 - Keep the main answer within 5-12 sentences.
36 - Consolidate citations in brackets if metadata is present.
37
38 CANDIDATE ANSWERS
39 {summaries}
40
41 FINAL ANSWER"""
42 )
43
44 # RetrievalQA icinde retriever en alakalı döküman parçalarını çeker, seçtiğin chain_type'a
45
46 def get_answers(query, k=5):
47     retriever = index.as_retriever(search_kwargs={"k": k})
48     qa = RetrievalQA.from_chain_type(#####
49         llm=llm,
50         retriever=retriever,
51         chain_type="map_reduce", # or "stuff" for simpler

```

```

52     chain_type_kwargs={
53         "question_prompt": QUESTION_PROMPT,
54         "combine_prompt": COMBINE_PROMPT
55     },
56     return_source_documents=True,
57 )
58 out = qa.invoke({"query": query}) # in LC 0.2+, key is "query"
59 return out['result']

```

```

1 our_query = " What are the two pre-training tasks used in BERT, and what does each one tra
2
3 answer = get_answers(our_query, k=7)
4 answer

```

'The two pre-training tasks used in BERT are the Masked Language Model (MLM) and Next Sentence Prediction (NSP). The MLM involves masking a percentage of the input tokens and training the model to predict these masked tokens based on their context, allowing the model to learn deep bidirectional representations. The NSP task trains the model to determine if a given sentence logically follows another, which helps in understanding the relationship between sent

```
1 Markdown(answer)
```

The two pre-training tasks used in BERT are the Masked Language Model (MLM) and Next Sentence Prediction (NSP). The MLM involves masking a percentage of the input tokens and training the model to predict these masked tokens based on their context, allowing the model to learn deep bidirectional representations. The NSP task trains the model to determine if a given sentence logically follows another, which helps in understanding the relationship between sentences and is crucial for tasks requiring context comprehension across multiple sentences [Figure 1].

```

1 our_query = "what is the relation between bert and tpu?"
2
3 answer = get_answers(our_query, k=5)
4 Markdown(answer)

```

Based on the provided context, there is insufficient information to formulate a definitive answer.

```

1
2 our_query = "In the context of fine-tuning, what minimal number of parameters need to be r
3
4 answer = get_answers(our_query, k=5)
5 Markdown(answer)
6

```

In the context of fine-tuning BERT with an additional output layer, the minimal number of parameters that need to be learned from scratch is the number of parameters in the new output layer itself. All other parameters of the pre-trained BERT model are typically retained and fine-tuned rather than initialized randomly. Therefore, only the parameters specific to the additional output layer would be initialized randomly.

```

1 def chat(k=5):
2     from IPython.display import display
3     while True:
4         q = input("Q (blank/quit to exit): ").strip()
5         if not q or q.lower() in {"q", "quit", "exit"}:
6             break
7         display(get_answers(q, k=k))
8         print("-" * 80)

```

```
1 chat()
```

```
Q (blank/quit to exit): q
```

```
1 Start coding or generate with AI.
```

```
1 Start coding or generate with AI.
```

