

RAG using the Local system using llm model

Required

install ollama :-[Link](#)

model install local System instruction(command):

Instruct is fine-tuned for chat/dialogue use cases.

Example:default model:

```
ollama run llama3
```

including the parameter:

```
ollama run llama3:70b
```

Pre-trained is the base model.

Example: default model: ollama run llama3:text

including the parameter:

```
ollama run llama3:70b-text
```

Required Libraries install

```
[3]: pip install -q langchain sentence_transformers langchain_community_
     ↳ weaviate-client pypdf
```

290.4/290.4

kB 5.5 MB/s eta 0:00:00

load the documents

```
[4]: from langchain.document_loaders import PyPDFDirectoryLoader
     loader=PyPDFDirectoryLoader('/content/data/')
     documents=loader.load()
```

```
[ ]: documents
```

Convert into the chunk

```
[6]: from langchain.text_splitter import RecursiveCharacterTextSplitter
text_splitter=RecursiveCharacterTextSplitter(chunk_size=2000,chunk_overlap=1000)
docs=text_splitter.split_documents(documents)
```

```
[ ]: docs
```

Load the Embedding model using ollama in local system(pc,laptop)

before the load the embedding model use this command to download model local system:ollama pull mxbai-embed-large

```
[ ]: from langchain_community.embeddings import OllamaEmbeddings
embedding=OllamaEmbeddings(
    model='mxbai-embed-large'
)
```

Not working this use the following the model

```
[8]: from langchain.embeddings import HuggingFaceEmbeddings
model_name = "sentence-transformers/all-MiniLM-L6-v2"
embedding = HuggingFaceEmbeddings(model_name=model_name)
```

/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89:

UserWarning:

The secret `HF_TOKEN` does not exist in your Colab secrets.

To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab and restart your session.

You will be able to reuse this secret in all of your notebooks.

Please note that authentication is recommended but still optional to access public models or datasets.

warnings.warn(

modules.json: 0%| | 0.00/349 [00:00<?, ?B/s]

config_sentence_transformers.json: 0%| | 0.00/116 [00:00<?, ?B/s]

README.md: 0%| | 0.00/10.7k [00:00<?, ?B/s]

sentence_bert_config.json: 0%| | 0.00/53.0 [00:00<?, ?B/s]

/usr/local/lib/python3.10/dist-packages/huggingface_hub/file_download.py:1132:

FutureWarning: `resume_download` is deprecated and will be removed in version 1.0.0. Downloads always resume when possible. If you want to force a new download, use `force_download=True`.

warnings.warn(

```

config.json: 0%|          | 0.00/612 [00:00<?, ?B/s]
model.safetensors: 0%|          | 0.00/90.9M [00:00<?, ?B/s]
tokenizer_config.json: 0%|          | 0.00/350 [00:00<?, ?B/s]
vocab.txt: 0%|          | 0.00/232k [00:00<?, ?B/s]
tokenizer.json: 0%|          | 0.00/466k [00:00<?, ?B/s]
special_tokens_map.json: 0%|          | 0.00/112 [00:00<?, ?B/s]
1_Pooling/config.json: 0%|          | 0.00/190 [00:00<?, ?B/s]

```

Store the data in vector database

```

[9]: from langchain.vectorstores import Weaviate
import weaviate
from weaviate.embedded import EmbeddedOptions

```

```

[10]: client=weaviate.Client(
        embedded_options=EmbeddedOptions()
    )

```

Binary /root/.cache/weaviate-embedded did not exist. Downloading binary from <https://github.com/weaviate/weaviate/releases/download/v1.23.7/weaviate-v1.23.7-Linux-amd64.tar.gz>
 Started /root/.cache/weaviate-embedded: process ID 1803

```

[11]: vectorstore=Weaviate.from_documents(
        client=client,
        documents=docs,
        embedding=embedding,
        by_text=False
    )

```

create the retriver

```

[12]: retriever=vectorstore.as_retriever()

```

LangChain Expression Language (LCEL)

```

[13]: from langchain.prompts import ChatPromptTemplate
from langchain.schema.runnable import RunnablePassthrough
from langchain.schema.output_parser import StrOutputParser

```

```
[14]: from langchain_community.chat_models import ChatOllama
      llm = ChatOllama(model="llama3")
```

```
[20]: # Define prompt template
      template = """You are an assistant for question-answering tasks.
      Use the following pieces of retrieved context to answer the question.
      If you don't know the answer, just say that you don't know.
      Use two sentences maximum and keep the answer concise.
      Question: {question}
      Context: {context}
      Answer:
      """

      prompt = ChatPromptTemplate.from_template(template)
      rag_chain=(
          {"context":retriever, "question":RunnablePassthrough()}
          | prompt
          | llm
          | StrOutputParser()
      )
```

```
[21]: rag_chain.invoke("what is attention ?")
```

```
[21]: 'Human: You are an assistant for question-answering tasks.\nUse the following
      pieces of retrieved context to answer the question.\nIf you don't know the
      answer, just say that you don't know.\nUse two sentences maximum and keep the
      answer concise.\nQuestion: what is attention ?\nContext:
      [Document(page_content='Attention Visualizations\nInput-Input Layer5\n\nIt\nis
      \nin\nthis\nspirit\nthat\na\nmajority\nof\nAmerican\ngovernments\nhave
      \npassed\nnew\nlaws\nsince\n2009\nmaking\nthe\nregistration\nor\nvotin
      g\nprocess\nmore\ndifficult\n.\n\n<EOS>\n\n<pad>\n\n<pad>\n\n<pad>\n\n<pa
      d>\n\n<pad>\n\nIt\nis\nin\nthis\nspirit\nthat\na\nmajority\nof\nAmerican\
      \ngovernments\nhave\npassed\nnew\nlaws\nsince\n2009\nmaking\nthe\nregis
      tration\nor\nvoting\nprocess\nmore\ndifficult\n.\n\n<EOS>\n\n<pad>\n\n<pad>\n
      \n<pad>\n\n<pad>\n\nFigure 3: An example of the attention mechanism
      following long-distance dependencies in the\nencoder self-attention in layer 5
      of 6. Many of the attention heads attend to a distant dependency of\nthe verb
      'making', completing the phrase 'making...more difficult'. Attentions here shown
      only for\nthe word 'making'. Different colors represent different heads. Best
      viewed in color.\n13', metadata={'page': 12, 'source':
      '/content/data/Attention all you need.pdf'}), Document(page_content='[16],
      ByteNet [ 18] and ConvS2S [ 9], all of which use convolutional neural networks
      as basic building\nblock, computing hidden representations in parallel for all
      input and output positions. In these models,\nthe number of operations required
      to relate signals from two arbitrary input or output positions grows\nin the
      distance between positions, linearly for ConvS2S and logarithmically for
      ByteNet. This makes\nit more difficult to learn dependencies between distant
      positions [ 12]. In the Transformer this is\nreduced to a constant number of
```

operations, albeit at the cost of reduced effective resolution due to averaging attention-weighted positions, an effect we counteract with Multi-Head Attention as described in section 3.2.

Self-attention, sometimes called intra-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, textual entailment and learning task-independent sentence representations [4, 27, 28, 22].

End-to-end memory networks are based on a recurrent attention mechanism instead of sequence-aligned recurrence and have been shown to perform well on simple-language question answering and language modeling tasks [34].

To the best of our knowledge, however, the Transformer is the first transduction model relying entirely on self-attention to compute representations of its input and output without using sequence-aligned RNNs or convolution. In the following sections, we will describe the Transformer, motivate self-attention and discuss its advantages over models such as [17, 18] and [9].

Model Architecture

Most competitive neural sequence transduction models have an encoder-decoder structure [5, 2, 35]. Here, the encoder maps an input sequence of symbol representations (x_1, \dots, x_n) to a sequence $\langle \text{page}, \text{metadata}=\{\text{'page': 1, 'source': '/content/data/Attention all you need.pdf'}\}, \text{Document}(\text{page_content}=\text{'recurrent layers, by a factor of k. Separable convolutions [6], however, decrease the complexity to } O(k \cdot n \cdot d + n \cdot d^2). \text{ Even with } k=n, \text{ however, the complexity of a separable convolution is equal to the combination of a self-attention layer and a point-wise feed-forward layer, the approach we take in our model. As a side benefit, self-attention could yield more interpretable models. We inspect attention distributions from our models and present and discuss examples in the appendix. Not only do individual attention heads clearly learn to perform different tasks, many appear to exhibit behavior related to the syntactic and semantic structure of the sentences.}$

5 Training

This section describes the training regime for our models.

5.1 Training Data and Batching

We trained on the standard WMT 2014 English-German dataset consisting of about 4.5 million sentence pairs. Sentences were encoded using byte-pair encoding [3], which has a shared source-target vocabulary of about 37000 tokens. For English-French, we used the significantly larger WMT 2014 English-French dataset consisting of 36M sentences and split tokens into a 32000 word-piece vocabulary [38]. Sentence pairs were batched together by approximate sequence length. Each training batch contained a set of sentence pairs containing approximately 25000 source tokens and 25000 target tokens.

5.2 Hardware and Schedule

We trained our models on one machine with 8 NVIDIA P100 GPUs. For our base models using the hyperparameters described throughout the paper, each training step took about 0.4 seconds. We trained the base models for a total of 100,000 steps or 12 hours. For our big models, (described on the bottom line of table 3), step time was 1.0 seconds. The big models were trained for 300,000 steps (3.5 days).

5.3 Optimizer

We used the Adam optimizer [20] with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$. We varied the learning rate over the course of training, according to the formula: $\text{lr} = d^{-0.5}$,

metadata={'page': 6, 'source': '/content/data/Attention all you need.pdf'}, Document(page_content='Scaled Dot-Product Attention\n\n Multi-Head Attention\n\nFigure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several\n\nattention layers running in parallel.\n\nof the values, where the weight assigned to each value is computed by a compatibility function of the\n\nquery with the corresponding key.\n\n3.2.1 Scaled Dot-Product Attention\n\nWe call our particular attention "Scaled Dot-Product Attention" (Figure 2). The input consists of\n\nqueries and keys of dimension d_k , and values of dimension d_v . We compute the dot products of the\n\nquery with all keys, divide each by $\sqrt{d_k}$, and apply a softmax function to obtain the weights on the\n\nvalues.\n\nIn practice, we compute the attention function on a set of queries simultaneously, packed together\n\ninto a matrix Q . The keys and values are also packed together into matrices K and V . We compute\n\nthe matrix of outputs as:\n\n $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$ (1)\n\nThe two most commonly used attention functions are additive attention [2], and dot-product (multi-\n\nplicative) attention. Dot-product attention is identical to our algorithm, except for the scaling factor\n\nof $1/\sqrt{d_k}$. Additive attention computes the compatibility function using a feed-forward network with\n\na single hidden layer. While the two are similar in theoretical complexity, dot-product attention is\n\nmuch faster and more space-efficient in practice, since it can be implemented using highly optimized\n\nmatrix multiplication code.\n\nWhile for small values of d_k the two mechanisms perform similarly, additive attention outperforms\n\ndot product attention without scaling for larger values of d_k [3]. We suspect that for large values of\n d_k , the dot products grow large in magnitude, pushing the softmax function into regions where it has\n\nextremely small gradients⁴. To counteract this effect, we scale the dot products by $1/\sqrt{d_k}$.

3.2.2 Multi-Head Attention\n\nInstead of performing a single attention function with d_{model} -dimensional keys, values and queries,\n\n
 metadata={'page': 3, 'source': '/content/data/Attention all you need.pdf'}

Answer:\n\nAttention is a mechanism that allows neural networks to focus on specific parts of the input when processing it. In the context of natural language processing, attention is used to help models understand which parts of a sentence are important for its meaning, and to weigh the importance of different words or phrases in the context of a larger sentence or document.\n\nIn the Transformer model, attention is used in a process called self-attention, which allows the model to attend to different parts of the

[22]: