

◆ Member-only story

# Scientific Documents Similarity Search With Deep Learning Using Transformers (SciBERT)

This article is a comprehensive overview of building a semantic similarity search tool for documents with k-NN and Cosine Similarity



Zoumana Keita · Follow

Published in Towards Data Science · 6 min read · Jan 17, 2022



81



...



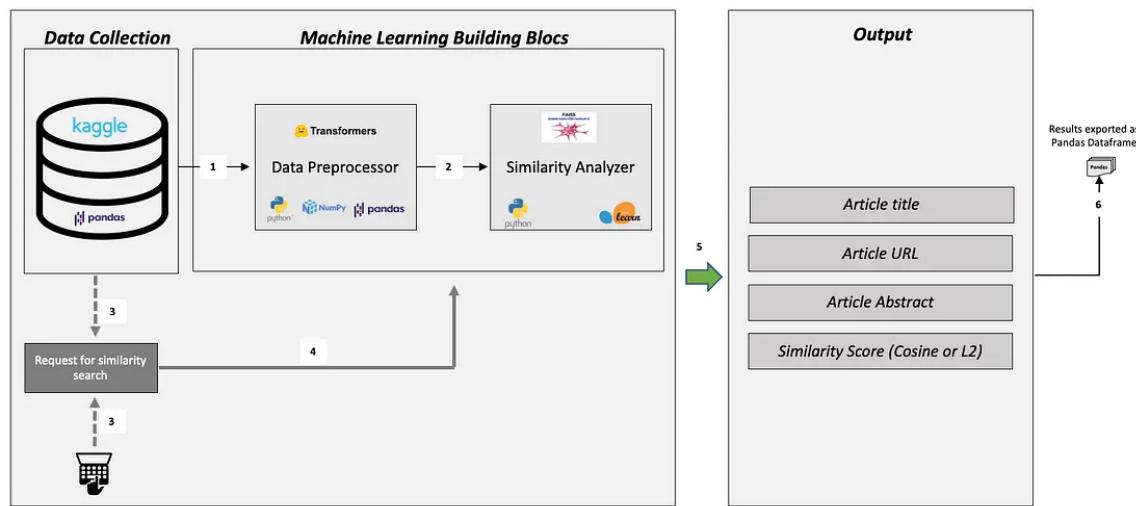
Photo by [Maksym Kaharlytskyi](#) on [Unsplash](#)

## Introduction

When reading an interesting article, you might want to find similar articles from a large corpus of data. Manual processing is obviously not the strategy to go for. So, why not take advantage of the power of Artificial Intelligence to solve such problems? From this article, you will be able to use SciBERT, and two different similarity approaches (Cosine and k-NN) in order to find scientific articles that are most similar in meaning to your specific query.

*Note: embeddings* and *vectors* will be used interchangeably to mean the same thing throughout the article.

The following workflow gives all the steps of the building of our tool, from data extraction to recommending similar articles.



#### Detailed workflow

Getting the data from Kaggle with Pandas library is the first step to creating the ML Solution

##### (1) Data Preprocessor

- Python is the main language for the analysis
- Transformers provides the pretrained SciBERT model
- Numpy to deal with arrays and pandas with dataframe

##### (2) Similarity Analyzer

- FAISS used to perform the k-NN similarity search
- Sklearn used to import the cosine similarity function

##### (3) & (4) Get request

- Case 1: get an example from the original data
- Case 2: the user provides his/her own query
- The query is preprocessed and used for similarity search

##### (5) & (6) Output and final result

- The output corresponds to the top N most similar articles (N given by user)
- The output is in pandas data format with the columns below
- Article title, Article URL, Article abstract and similarity score (Cosine or L2 distance)

General Workflow of the Solution Being Solved (Image by Author)

## About the Data

- This is the CORD-19 data set, a resource of over 59,000 scholarly articles, including over 48,000 with full text, about COVID-19, SARS-CoV-2, and related coronaviruses.
- This data has been **freely available** by the White House and a coalition of leading research groups in order to help global research generate insights in support of the ongoing fight against this infectious disease.
- It is downloadable from [this page on Kaggle](#).
- Further details about the dataset can be found on [this page](#).

Before moving further, it is important to consider importing the following libraries that are crucial to the success of this project.

`useful_libs.py`

The data is downloaded from Kaggle and saved in a folder called *input*. Using

the shape function we get (47110, 16), meaning that it contains 47110 articles, and each one has 16 columns.

## **16 columns – but which one is the most important?**

Below are a few details about the percentage of missing information in the data.

Before removing articles with missing Abstract		After removing articles with missing Abstract	
paper_id	0.000000	paper_id	0.000000
body_text	0.000000	body_text	0.000000
methods	58.439822	methods	53.796698
results	63.984292	results	59.466173
source	10.365103	source	8.645463
title	10.439397	title	8.723394
doi	12.604543	doi	11.095417
abstract	12.838039	abstract	0.000000
publish_time	10.365103	publish_time	8.645463
authors	11.676926	authors	9.432078
journal	18.199958	journal	17.473577
arxiv_id	98.681809	arxiv_id	98.487653
url	10.617703	url	8.927963
publish_year	0.000000	publish_year	0.000000
is_covid19	0.000000	is_covid19	0.000000
study_design	0.000000	study_design	0.000000
dtype: float64		dtype: float64	

Percentage of missing values from the data with a focus on Abstract column (Image by Author)

We will focus our analysis on the **abstract** column for simplicity's sake, also it

is the one with 0% missing data. But you could use other textual columns such as `body_text`; it is up to you. Also, we will use a subset of 2000 observations in order to speed the processing.

### **Articles' Abstract in the data**

Let have a look at some random articles. Here we limit the printing to the first hundred words, because some of them are very long.

show\_random\_articles.py

Article #3252

--> Title: Tuning antiviral CD8 T-cell response via proline-altered peptide ligand vaccination

--> Abstract: AbstractViral escape from CD8+ cytotoxic T lymphocyte responses correlates with disease progression and represents a significant challenge for vaccination. Here, we demonstrate that CD8+ T cell recognition of the naturally occurring MHC-I-restricted LCMV-associated immune escape variant Y4F is restored following vaccination with a proline-altered peptide ligand (APL). The APL increases MHC/peptide (pMHC) complex stability, rigidifies the peptide and facilitates T cell receptor (TCR) recognition through reduced entropy costs. Structural analyses of pMHC complexes before and after TCR binding, combined with biophysical analyses, revealed that although the TCR binds similarly to all complexes, the p3P modification alters the conformations of a ...

Article #36667

--> Title: A highly conserved WDYPKCDRA epitope in the RNA directed RNA polymerase of human coronaviruses can be used as epitope-based universal vaccine design

--> Abstract: BACKGROUND: Coronaviruses are the diverse group of RNA virus. From 1960, six strains of human coronaviruses have emerged that includes SARS-CoV and the recent infection by deadly MERS-CoV which is now going to cause another outbreak. Prevention of these viruses is urgent and a universal vaccine for all strain could be a promising solution in this circumstance. In this study we aimed to design an epitope based vaccine against all strain of human coronavirus. RESULTS: Multiple sequence alignment (MSA) approach was employed among spike (S), membrane (M), enveloped (E) and nucleocapsid (N) protein and replicase polyprotein 1ab to identify which ...

Article #14038

--> Title: Development and characterization of a Rift Valley fever virus cell-cell fusion assay using alphavirus replicon vectors

--> Abstract: Abstract Rift Valley fever virus (RVFV), a member of the Phlebovirus genus in the Bunyaviridae family, is transmitted by mosquitoes and infects both humans and domestic animals, particularly cattle and sheep. Since primary RVFV strains must be handled in BSL-3+ or BSL-4 facilities, a RVFV cell-cell fusion assay will facilitate the investigation of RVFV glycoprotein function under BSL-2 conditions. As for other members of the Bunyaviridae family, RVFV glycoproteins are targeted to the Golgi, where the virus buds, and are not efficiently delivered to the cell surface. However, overexpression of RVFV glycoproteins using an alphavirus replicon vector resulted in the ...

## Data Processing & Vectorization

This step aims to vectorize the articles' abstract text so that we can perform the similarity analysis. Since we are dealing with the scientific documents, we will use SciBERT, which is a pre-trained language model for Scientific text data. You can find more information about it on [Semantic Scholar](#).

The main steps involved in this part are:

Load the pre-trained model & tokenizer. When loading the model, we need to set the *output\_hidden\_states* to True so that we can extract the embeddings.

`load_model_artifacts.py`

This function `convert_single_abstract_to_embedding` is mostly inspired by the

BERT Word [Embeddings Tutorial](#) of Chris McCormick. It aims to create an embedding for a given text data using a pre-trained model.

### **Test on a single text data**

Here we test the function on the 30th article. You can choose whatever number you want, as long as it exists in the data.

test\_embedding.py

*Line 6* shows **Embedding shape: (768,)**. This means that a single vector is composed of 768 values. We can finally move to the next step by applying the conversion process to all the articles in the data using the *convert\_overall\_text\_to\_embedding* function. But, before that, we are going to remove some columns from the data using the *get\_min\_viable\_data* function, in order to have fewer columns in the result of the final query search.

With the previous helper functions, we can create a new column that will

contain for each article, its corresponding abstract embedding.

create\_final\_emb.py

	paper_id	title	abstract	url	embeddings
11754	9778c2bdf6be32053f9a450194f767b47c87c891	NaN	Graphical Abstract Highlights d The cryo-EM st...	NaN	[-0.3626312, 0.34195867, 0.40458974, -0.65087...
20720	98fc84b407a8632ee5ac42431d12c6acf7dfe3d8	Genetic Loci That Influence Cause of Death in ...	A genome scan was conducted to seek evidence f...	<a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7...">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7...</a>	[-0.49350545, -1.1899799, -0.0245932, -0.8623...
15193	fd0fc8d711665338facb36ba94489c3b63a0d7	New Pre-pandemic Influenza Vaccines: An Egg-an...	Highly pathogenic avian H5N1 influenza viruses...	<a href="http://europepmc.org/articles/pmc2793094?pdf=r...">http://europepmc.org/articles/pmc2793094?pdf=r...</a>	[0.3071721, -0.85410595, 0.43242905, 0.003505...

First 3 rows of the Data With Embeddings column (Image by Author)

So far, so good! Everything is finally set up for similarity search.

## Similarity Search

Now, we can perform the similarity analysis between a given *query* vector and all the embeddings vectors. The *query* processing will be similar to both cosine and k-NN. Below is the function responsible for that.

`process_query.py`

The cosine similarity between two documents' embedding measures how

similar those documents are, irrespective of the size of those embeddings. It measures the cosine of the angle between the two vectors projected in a multi-dimensional space.

- *cosine similarity of 1* means that the two documents are 100% similar
- *cosine similarity of 0* means that the two documents have 0% similarity

The following function returns the top N (N is the number of similar articles to return) articles similar to the query text.

`cosine_recommendations.py`

Now we can call the function to get the top 5 most similar articles.

Here is the result of the previous query

title	abstract	url	cos_sim
NaN	Highlights d ENDU-2 nuclease regulates nucleot...	NaN	0.825622
Regioselective synthesis of 6-substituted-2-am...	Abstract A series of 2-amino-5-bromo-4(3H)-pyr...	<a href="https://doi.org/10.1016/j.ejmech.2013.06.036">https://doi.org/10.1016 /ejmech.2013.06.036</a>	0.773827
Synthesis of 4-aminoquinoline-pyrimidine hybri...	Abstract One of the most viable options to tac...	<a href="https://doi.org/10.1016/j.ejmech.2013.05.046">https://doi.org/10.1016 /ejmech.2013.05.046</a>	0.761577
Australia was indeed the "lucky country" in th...	Anton Y Peleg, Wendy J Munckhof Australia was...	<a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7...">https://www.ncbi.nlm.nih.gov /pmc/articles/PMC7...</a>	0.757679
Antigen delivery systems for veterinary vaccin...	Abstract The recent advances in molecular gene...	<a href="https://doi.org/10.1016/j.vaccine.2008.09.044">https://doi.org/10.1016 /j.vaccine.2008.09.044</a>	0.754393

Top 5 most similar Articles to the Query (Image by Author)

For a better view, we can just look at the first 2 most similar articles.

*Query text below*

Graphical Abstract Highlights d The cryo-EM structure of full-length human NPC1 was determined at 4.4 Å resolution d Structure-guided biochemical analysis of cholesterol transfer from NPC2 to NPC1 d Low-resolution cryo-EM structure of NPC1 bound to GPc1 of Ebola virus was obtained d A trimeric GPc1 binds to one NPC1 through the crystal structure-revealed interface\*\*

Query text data from the original dataframe

*Most two similar articles from the previous top\_articles dataframe*

1<sup>st</sup> Article's Abstract

Cosine Similarity: 0.82

'Highlights d ENDU-2 nuclease regulates nucleotide metabolism and germ cell proliferation in worms d ENDU-2 expression is induced by nucleotide imbalance and other genotoxic stresses d ENDU-2 inhibits CTP synthase phosphorylation by repressing PKA and HDA-1 in the gut d ENDU-2 function may be conserved in mammalian cells'

2<sup>nd</sup> Article's Abstract

Cosine Similarity: 0.77

'Abstract A series of 2-amino-5-bromo-4(3H)-pyrimidinone derivatives bearing different substituents at the C-6 position were synthesized using a highly regioselective lithiation-substitution protocol, and the effect of structural variation at the C-6 position on their antiviral activity in cell culture was evaluated. Although some of the derivatives were found to be active against various virus strains, they were effective only close to their toxicity threshold.'

Extract of the most two similar articles from the recommended dataframe (Image by Author)

FAISS is a library developed by [Facebook AI Research](#). According to their [wikipage](#):

Artificial Intelligence

Machine Learning

Naturallanguageprocessing

Data Science

Cosine Similarity tains algorithms that search in sets of vectors of any size, up to ones that possibly do not fit in RAM

Below are the steps to build the search algorithm using the previously built embeddings

- create the flat index. The index uses the L2 (Euclidean) distance metrics to measure the similarity between the query vector and all the vectors



eddings).



- add all the vectors to the index

## Written by Zoumana Keita

define the number K of similar documents we want

Follow



8.2K Followers · Writer for Towards Data Science

- run the similarity search to get the result

Data Scientist at IFC - The World Bank Group | Videos about AI, Data Science, Programming or Tech ↗ <https://www.youtube.com/@zoumadatasscience>

## More from Zoumana Keita and Towards Data Science



Zoumana Keita in Artificial Corner

### 4 Hidden Python Treasures That You Probably Didn't Know Existed

Learning these libraries will take your Python skill to the next level—no doubt.

💡 · 7 min read · Aug 3



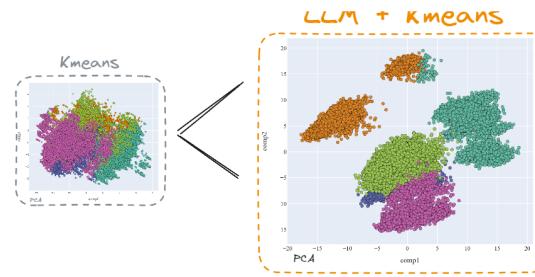
511



5



...



Damian Gil in Towards Data Science

### Mastering Customer Segmentation with LLM

Unlock advanced customer segmentation techniques using LLMs, and improve your...

23 min read · Sep 26



3.1K



25



...



Khouloud El Alami in Towards Data Science

## Don't Start Your Data Science Journey Without These 5 Must-D...

A complete guide to everything I wish I'd done before starting my Data Science...

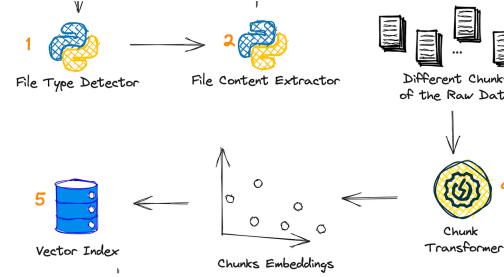
◆ · 18 min read · Sep 24

👏 2.4K

💬 23



...



Zoumanna Keita in Towards Data Science

## How to Chat With Any File from PDFs to Images Using Large...

Complete guide to building an AI assistant that can answer questions about any file

◆ · 9 min read · Aug 5

👏 1.5K

💬 13



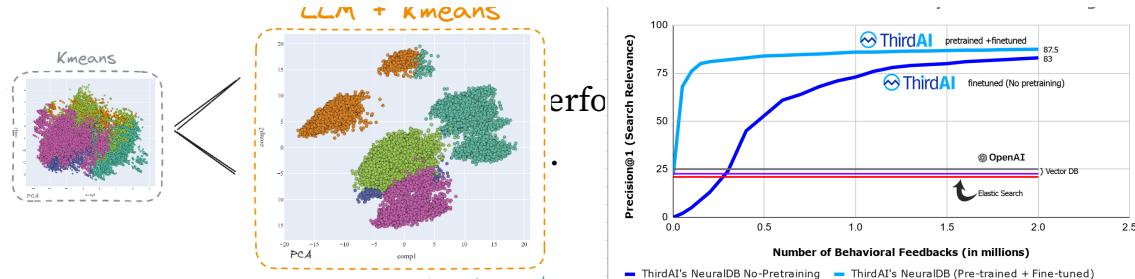
...

[See all from Zoumanna Keita](#)

[See all from Towards Data Science](#)

## Recommended from Medium

[faiss\\_setup.py](#)



Damian Gil in Towards Data Science

## Mastering Customer Segmentation with LLM

Unlock advanced customer segmentation techniques using LLMs, and improve your...

23 min read · Sep 26

3.1K

25



...

Anshu in ThirdAI Blog

## Demystifying LLM-Driven Search: Stop Comparing Embeddings or...

A comprehensive evaluation of a commercial semantic search system reveals that the...

5 min read · 5 days ago

34

3



...

## Lists



### Predictive Modeling w/ Python

20 stories · 490 saves



### Natural Language Processing

706 stories · 312 saves



### Practical Guides to Machine Learning

10 stories · 561 saves



### ChatGPT prompts

26 stories · 501 saves



 Justin Swansburg

Note: I decided to break down all the steps on purpose in order to make sure you understand them properly. But you can put everything together into a **RAG Pipeline Pitfalls: The Untold Challenges of Embedding Table**

# single function

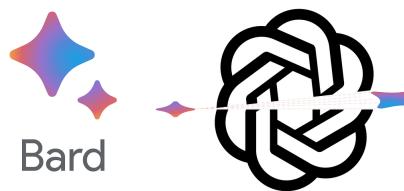
A quick guide on a new way to leverage large language models to improve the separation...

9 min read · May 13



The figure consists of three separate line graphs, each titled with a dataset name and showing accuracy versus the number of top terms (N). Each graph contains five data series representing different models: ADA (green), BERT (red), FastText (blue), GLOVE (orange), and TFIDF (yellow).

- MozillaCore:** The x-axis ranges from 0 to 500, and the y-axis ranges from 0 to 60. All models show a rapid increase in accuracy initially, followed by a more gradual plateau. BERT reaches the highest accuracy of approximately 62% at N=500.
- JD'T:** The x-axis ranges from 0 to 500, and the y-axis ranges from 0 to 70. Similar to MozillaCore, accuracy increases rapidly and plateaus. BERT reaches the highest accuracy of approximately 70% at N=500.
- EclipsePlatform:** The x-axis ranges from 0 to 500, and the y-axis ranges from 0 to 10. The models show lower overall accuracy compared to the other two datasets. BERT reaches the highest accuracy of approximately 8.5% at N=500.





Avinash Patil

## Embeddings: BERT better than ChatGPT4?

In this study, we compared the effectiveness of semantic textual similarity methods for...

See more recommendations  
4 min read · Sep 18



AL Anany

## The ChatGPT Hype Is Over—Now Watch How Google Will Kill...

It never happens instantly. The business game is longer than you know.

• 6 min read · Sep 1



--



2



+

...



--



424



+

...

[Help](#) [Status](#) [About](#) [Careers](#) [Blog](#) [Privacy](#) [Terms](#) [Text to speech](#) [Teams](#)

faiss\_show\_recommendations.py

\*\* Article #0 \*\*  
\*\* --> Abstract :  
Graphical Abstract Highlights d The cryo-EM structure of full-length human NPC1 was determined at 4.4 Å resolution d Structure-guided biochemical analysis of cholesterol transfer from NPC2 to NPC1 d Low-resolution cryo-EM structure of NPC1 bound to GPcl of Ebola virus was obtained d A trimeric GPcl binds to one NPC1 through the crystal structure-revealed interface\*\*  
\*\* --> L2 Distance: 0.00\*\*

\*\* Article #1683 \*\*  
\*\* --> Abstract :  
Highlights d ENDU-2 nuclease regulates nucleotide metabolism and germ cell proliferation in worms d ENDU-2 expression is induced by nucleotide imbalance and other genotoxic stresses d ENDU-2 inhibits CTP synthase phosphorylation by repressing PKA and HDA-1 in the gut d ENDU-2 function may be conserved in mammalian cells\*\*  
\*\* --> L2 Distance: 181.52\*\*

\*\* Article #1021 \*\*  
\*\* --> Abstract :  
Abstract A series of 2-amino-5-bromo-4(3H)-pyrimidinone derivatives bearing different substituents at the C-6 position were synthesized using a highly regioselective lithiation-substitution protocol, and the effect of structural variation at the C-6 position on their antiviral activity in cell culture was evaluated. Although some of the derivatives were found to be active against various virus strains, they were effective only close to their toxicity threshold.  
\*\* --> L2 Distance: 230.64\*\*

\*\* Article #348 \*\*  
\*\* --> Abstract :  
Anton Y Peleg, Wendy J Munckhof  
Australia was indeed the "lucky country" in the recent worldwide SARS epidemic 229\*\*  
\*\* --> L2 Distance: 244.74\*\*

\*\* Article #1712 \*\*  
\*\* --> Abstract :  
Abstract One of the most viable options to tackle the growing resistance to the antimalarial drugs such as artemisinin is to resort to synthetic drugs. The multi-target strategy involving the use of hybrid drugs has shown promise. In line with this, new hybrids of quinoline wi

th pyrimidine have been synthesized and evaluated for their antiplasmodial activity against both CQS and CQR strains of *Plasmodium falciparum*. These depicted activity in nanomolar range and were found to bind to heme as well as AT rich pUC18 DNA.\*\*  
\*\* --> L2 Distance: 245.88\*\*

5 most similar articles similar to the query (Image by author)

## Observation

- The first document has  $L_2 = 0$ , which means 100% similarity. This is obvious because the query was compared with itself, but we can simply remove it from the analysis.

## Conclusion

In this article, we've studied the whole pipeline of implementing a semantic similarity search tool using SciBERT embedding, Cosine, and k-NN similarity approaches. Also, I hope it has been beneficial to you as well!

You can find below additional resources to further your learning. Follow me on [YouTube](#) for more interactive sessions!

## Additional Ressources

[Article Source Code on GitHub](#)

[SciBERT: A Pretrained Language Model for Scientific Text](#)

[BERT Word Embeddings Tutorial](#)

Facebook AI Research

Bye for now 🏃