

# A Full Guide to Finetuning T5 for Text2Text and Building a Demo with Streamlit

All you need to know to build a full demo: Hugging Face Hub, Tensorboard, Streamlit, and Hugging Face Spaces



Fabio Chiusano · [Follow](#)

Published in NLPlanet

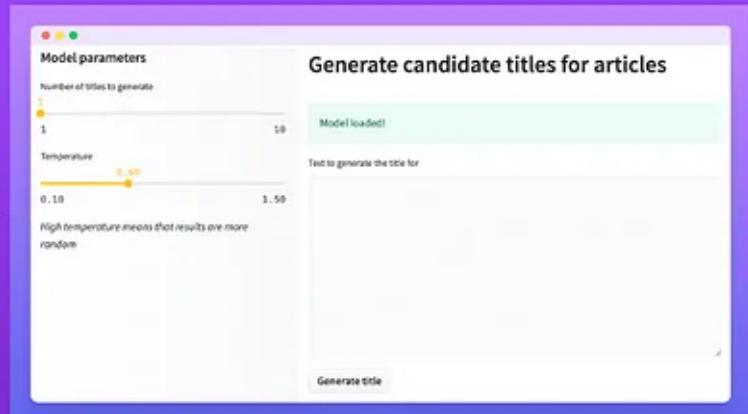
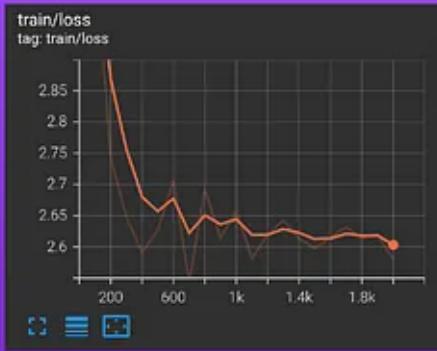
15 min read · May 17, 2022

Listen

Share

More

## Finetuning T5 + Streamlit



*Hello fellow NLP enthusiasts! I'm back at writing articles after a small break, can't wait to dig into new NLP topics! In this article, we see a complete example of fine-tuning of T5 for*

*generating candidate titles for articles. The model is fine-tuned entirely on [Colab](#), we visualize its training with [TensorBoard](#), upload the model on the [Hugging Face Hub](#) for everyone to use, and create a small demo with [Streamlit](#) that we upload on [Hugging Face Spaces](#) for everyone to try it out. Enjoy! 😊*

I've been wanting to experiment with [Streamlit](#) and [Hugging Face Spaces](#) for a while now. In case you didn't know them:

- Streamlit is a Python library that allows data scientists to easily create small interactive demos, so that other people can test their machine learning models or see their data analyses. Check out the [Streamlit Gallery](#) to learn what can be done with the library. Moreover, you can deploy your Streamlit application to the [Streamlit Cloud](#), which makes it easy to share your applications with other people and you get free CPUs to run your applications. I suggest having a look at [Gradio](#) as well, which is another popular Python library with similar goals.
- Hugging Face Spaces is a service where you can deploy your Streamlit or Gradio applications so that you can easily share them. It provides free CPUs and it's similar to the Streamlit Cloud. However, when you deploy an application on Hugging Face Spaces that uses a machine learning model uploaded on the Hugging Face Hub, everyone can see that they are linked (i.e., that there is an interactive demo for that model).

To test them out, I decided to fine-tune a pre-trained model on a simple but not-so-popular task, which is the generation of candidate titles for articles starting from their textual content.

This is what we are going to do in this article:

1. Find a suitable dataset containing articles textual contents and titles.
2. Choose a suitable metric for our task.

3. Fine-tune a pre-trained model for title generation on Colab, monitoring the chosen metric on the validation set using TensorBoard, and saving the model's checkpoints on Google Drive (so that we can resume training in case Colab shuts down the connection).
4. Upload the model on Hugging Face Hub for everyone to use.
5. Build an interactive demo with Streamlit and deploy it to Hugging Face Spaces.

Let's start!

## Choosing a dataset for the Title Generation task

I usually look for datasets on [Google Dataset Search](#) or [Kaggle](#). There are already available datasets with many newspaper articles, but for this tutorial I wanted to focus on the kind of articles that we usually find on Medium: less news, more guides.

Since I couldn't find a Medium articles dataset containing both articles textual contents and titles, I created one by myself by scraping the website and parsing web pages using the [newspaper Python library](#). I then uploaded it to Kaggle with the name [190k+ Medium Articles dataset](#) so that everyone can use it, along with a [public notebook with simple data exploration](#).

Each row in the data is a different article published on Medium. For each article, you have the following features:

- *title [string]*: The title of the article.
- *text [string]*: The text content of the article.
- *url [string]*: The URL associated with the article.
- *authors [list of strings]*: The article authors.
- *timestamp [string]*: The publication datetime of the article.
- *tags [list of strings]*: List of tags associated with the article.

## Choosing a metric for the Title Generation task

The task of generating titles starting from the textual content of an article is a text2text generation task: we have a text in input and we want to generate some text as output.

Popular text2text generation tasks are machine translation, commonly evaluated with the BLEU score and a focus on word precision, and text summarization, commonly evaluated with the ROUGE score and a focus on word recall.

I see title generation as closely related to text summarization as the title should make the reader understand what is the article about, with the added flavor that the title should also intrigue the reader and make him/her curious about the article. For this reason, I decided to evaluate my models with the ROUGE score.

## Model training

Let's jump now into code with Colab! You can find all the code in this public Colab notebook as well. This is what we'll do:

1. Connect Google Drive to Colab to have persistent storage across Colab sessions.
2. Download the Medium dataset from Kaggle.
3. Load the dataset.
4. Split the dataset into train, validation, and test set.
5. Preprocess the dataset for T5.
6. Preparing the Hugging Face trainer.
7. Start TensorBoard.
8. Fine-tune T5.
9. Try the model.

## 10. Evaluate the model on the test set.

First, we install some libraries:

```
1 !pip install datasets transformers rouge_score nltk
```

title\_prediction\_2.sh hosted with ❤ by GitHub

[view raw](#)

We use the Hugging Face [transformers](#) library to download pre-trained models and fine-tune them, the Hugging Face [datasets](#) library to load our dataset and preprocess it, the [rouge-score](#) library to compute ROUGE scores, and the [nltk](#) library which contains common NLP preprocessing functions.

### Connect Google Drive to Colab

Here is a [Colab notebook](#) that provides recipes for loading and saving data in Colab from external sources. We are interested in the “Mounting Google Drive locally” section: by executing the following code in our notebook, you’ll be prompted to allow the notebook to read data from your Google Drive.

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

title\_prediction\_1.py hosted with ❤ by GitHub

[view raw](#)

You can then find your Google Drive data in the `drive/MyDrive` directory. Later in this article, we’ll save our model checkpoints inside a purposely created `drive/MyDrive/Models` directory.

### Download the Medium dataset from Kaggle

We’ll download our dataset using the `kaggle` command-line utility, which comes already installed on Google Colab. To use the `kaggle` command, we need to connect it to our Kaggle account through an API key. Here are the [instructions to create the API key](#), it’s very easy. Once done, you should have a `kaggle.json` file.

Next, you need to upload your `kaggle.json` file to your Colab instance. If you don’t know how to do it, check the “Uploading files from your local file system” section of this [notebook](#).

Lastly, we need to instruct the `kaggle` command-line utility to use your `kaggle.json` file. You can do this by executing the command `kaggle`, which will create the hidden directory `.kaggle` (along with throwing an error because it can't find a `kaggle.json` file inside the `.kaggle` directory). Then, copy the file into the directory with `cp kaggle.json ~/.kaggle/kaggle.json`.

We can now download the dataset! First, we must retrieve the dataset name, which is whatever comes after `kaggle.com/datasets/` in the dataset URL. Our dataset URL is `kaggle.com/datasets/fabiochiusano/medium-articles`, therefore the dataset name is `fabiochiusano/medium-articles`. Execute the command `kaggle datasets download -d fabiochiusano/medium-articles`, and you should find the file `medium-articles.zip` in your current directory.

## Load the dataset

Let's import the necessary libraries.

```
1 import transformers  
2 from datasets import load_dataset, load_metric
```

[title\\_prediction\\_3.py](#) hosted with ❤ by GitHub [view raw](#)

We load the dataset using the `load_dataset` function from the `datasets` package.

```
1 medium_datasets = load_dataset("csv", data_files="medium-articles.zip")  
2 # DatasetDict({  
3 #     train: Dataset({  
4 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],  
5 #         num_rows: 192368  
6 #     })  
7 # })
```

[title\\_prediction\\_4.py](#) hosted with ❤ by GitHub [view raw](#)

## Split the dataset into train, validation, and test set

As a common practice, we split the dataset into:

- *Training set*: Data used for training the model parameters.

- *Validation set*: Data used for hyperparameter tuning or early stopping to avoid overfitting.
- *Test set*: Data used for checking what performance we can expect on new data.

This can be done with the `train_test_split` methods of our dataset object, as follows.

```
1 datasets_train_test = medium_datasets["train"].train_test_split(test_size=3000)
2 datasets_train_validation = datasets_train_test["train"].train_test_split(test_size=3000)
3
4 medium_datasets["train"] = datasets_train_validation["train"]
5 medium_datasets["validation"] = datasets_train_validation["test"]
6 medium_datasets["test"] = datasets_train_test["test"]
7
8 # DatasetDict({
9 #     train: Dataset({
10 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],
11 #         num_rows: 186368
12 #     })
13 #     validation: Dataset({
14 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],
15 #         num_rows: 3000
16 #     })
17 #     test: Dataset({
18 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],
19 #         num_rows: 3000
20 #     })
21 # })
```

[title\\_prediction\\_5.py](#) hosted with ❤ by GitHub

[view raw](#)

The choice of the sizes of the three datasets usually depends on many factors, such as the size of the whole dataset, how much time you want to spend training or evaluating your model, and how precise you want your model evaluations to be. Since we are training a quite big model on a free GPU from Colab, I decided to use small validation and test sets to speed things up. Despite this, we will see that the final model will be able to produce good titles.

```

1 medium_datasets["train"] = medium_datasets["train"].shuffle().select(range(100000))
2 medium_datasets["validation"] = medium_datasets["validation"].shuffle().select(range(1000))
3 medium_datasets["test"] = medium_datasets["test"].shuffle().select(range(1000))
4
5 # DatasetDict({
6 #     train: Dataset({
7 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],
8 #         num_rows: 100000
9 #     })
10 #     validation: Dataset({
11 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],
12 #         num_rows: 1000
13 #     })
14 #     test: Dataset({
15 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags'],
16 #         num_rows: 1000
17 #     })
18 # })

```

[title\\_prediction\\_6.py](#) hosted with ❤ by GitHub

[view raw](#)

## Preprocess the dataset for T5

Hugging Face provides us with a complete [notebook example of how to fine-tune T5 for text summarization](#). As for every transformer model, we need first to tokenize the textual training data: the article content and the title.

Let's instantiate the tokenizer of the [T5-base model](#).

```

1 import nltk
2 nltk.download('punkt')
3 import string
4 from transformers import AutoTokenizer
5
6 model_checkpoint = "t5-base"
7 tokenizer = AutoTokenizer.from_pretrained(model_checkpoint)

```

[title\\_prediction\\_7.py](#) hosted with ❤ by GitHub

[view raw](#)

Before applying the tokenizer to the data, let's filter out some bad samples (i.e. articles whose title is long less than 20 characters and whose text content is long less than 500 characters).

```
1 medium_datasets_cleaned = medium_datasets.filter(  
2     lambda example: (len(example['text']) >= 500) and  
3     (len(example['title']) >= 20)  
4 )
```

[title\\_prediction\\_9.py](#) hosted with ❤ by GitHub

[view raw](#)

Then, we define the `preprocess_data` function that takes a batch of samples as inputs and outputs a dictionary of new features to add to the samples. The `preprocess_data` function does the following:

1. Extract the “text” feature from each sample (i.e. the article text content), fix the newlines in the article, and remove the lines without ending punctuation (i.e. the subtitles).
2. Prepend the text “*summarize*: “ to each article text, which is needed for fine-tuning T5 on the summarization task.
3. Apply the T5 tokenizer to the article text, creating the `model_inputs` object. This object is a dictionary containing, for each article, an `input_ids` and an `attention_mask` arrays containing the token ids and the attention masks respectively.
4. Apply the T5 tokenizer to the article titles, creating the `labels` object. Also in this case, this object is a dictionary containing, for each article, an `input_ids` and an `attention_mask` arrays containing the token ids and the attention masks respectively. Note that this step is done inside the `tokenizer.as_target_tokenizer()` context manager: this is usually done because there are text2text tasks where inputs and labels must be tokenized with different tokenizers (e.g. when translating between two languages, where each language has its own tokenizer). As far as I know, for text summarization the labels are tokenized with the same tokenizer as the inputs, thus the context manager is optional.
5. Return a dictionary containing the token ids and attention masks of the inputs, and the token ids of the labels.

```

1  prefix = "summarize: "
2  max_input_length = 512
3  max_target_length = 64
4
5  def clean_text(text):
6      sentences = nltk.sent_tokenize(text.strip())
7      sentences_cleaned = [s for sent in sentences for s in sent.split("\n")]
8      sentences_cleaned_no_titles = [sent for sent in sentences_cleaned
9                                      if len(sent) > 0 and
10                                     sent[-1] in string.punctuation]
11     text_cleaned = "\n".join(sentences_cleaned_no_titles)
12     return text_cleaned
13
14 def preprocess_data(examples):
15     texts_cleaned = [clean_text(text) for text in examples["text"]]
16     inputs = [prefix + text for text in texts_cleaned]
17     model_inputs = tokenizer(inputs, max_length=max_input_length, truncation=True)
18
19     # Setup the tokenizer for targets
20     with tokenizer.as_target_tokenizer():
21         labels = tokenizer(examples["title"], max_length=max_target_length,
22                            truncation=True)
23
24     model_inputs["labels"] = labels["input_ids"]
25     return model_inputs

```

[title\\_prediction\\_10.py](#) hosted with ❤ by GitHub

[view raw](#)

*Note that we are truncating the inputs at 512 tokens. While T5 can manage longer inputs, the memory requirements grow quadratically with the size of the inputs, and that was the maximum size that I could use in my Colab session. 512 tokens correspond to about 682 English words, which is more or less what an average person reads in two minutes. The majority of Medium articles have between four and seven minutes of reading time, therefore we are currently throwing away useful information for our task. Despite this, many articles state what they are about in their first paragraphs, therefore good titles can be generated on most occasions.*

The `preprocess_data` function can be applied to all the datasets with the `map` method.

```

1 tokenized_datasets = medium_datasets_cleaned.map(preprocess_function,
2                                         batched=True)
3
4 # DatasetDict({
5 #     train: Dataset({
6 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags',
7 #                     'input_ids', 'attention_mask', 'labels'],
8 #         num_rows: 85602
9 #     })
10 #     validation: Dataset({
11 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags',
12 #                     'input_ids', 'attention_mask', 'labels'],
13 #         num_rows: 860
14 #     })
15 #     test: Dataset({
16 #         features: ['title', 'text', 'url', 'authors', 'timestamp', 'tags',
17 #                     'input_ids', 'attention_mask', 'labels'],
18 #         num_rows: 844
19 #     })
20 # })

```

[title\\_prediction\\_8.py](#) hosted with ❤ by GitHub[view raw](#)

## Preparing the Hugging Face trainer

We can now fine-tune T5 with our preprocessed data! Let's import some necessary classes to train text2text models.

```

1 from transformers import AutoModelForSeq2SeqLM, DataCollatorForSeq2Seq,
2 Seq2SeqTrainingArguments, Seq2SeqTrainer

```

[title\\_prediction\\_11.py](#) hosted with ❤ by GitHub[view raw](#)

Next, we need to create a `Seq2SeqTrainingArguments` object containing, as the name implies, several training parameters that will define how the model is trained. Refer to the [Trainer](#) documentation to know about the meaning of each one of these parameters.

```

1 batch_size = 8
2 model_name = "t5-base-medium-title-generation"
3 model_dir = f"drive/MyDrive/Models/{model_name}"
4
5 args = Seq2SeqTrainingArguments(
6     model_dir,
7     evaluation_strategy="steps",
8     eval_steps=100,
9     logging_strategy="steps",
10    logging_steps=100,
11    save_strategy="steps",
12    save_steps=200,
13    learning_rate=4e-5,
14    per_device_train_batch_size=batch_size,
15    per_device_eval_batch_size=batch_size,
16    weight_decay=0.01,
17    save_total_limit=3,
18    num_train_epochs=1,
19    predict_with_generate=True,
20    fp16=True,
21    load_best_model_at_end=True,
22    metric_for_best_model="rouge1",
23    report_to="tensorboard"
24 )

```

[title\\_prediction\\_12.py](#) hosted with ❤ by [GitHub](#)

[view raw](#)

Here is the explanation of some unusual parameters passed to the `Seq2SeqTrainingArguments` object:

- `predict_with_generate` : Must be set to `True` to calculate generative metrics such as ROUGE and BLEU.
- `fp16` : Whether to use fp16 16-bit (mixed) precision training instead of 32-bit training. Makes training faster.
- `report_to` : List of integrations to write logs to.

Note that the output directory is inside Google Drive and that we are specifying `rouge1` as the metric which defines the best model.

Next, we instantiate a `DataCollatorForSeq2Seq` object using the tokenizer. Data collators are objects that form a batch by using a list of dataset elements as input and, in some cases, applying some processing. In this case, all the inputs and labels in the same batch will be padded to their respective maximum length in the batch. Padding of the inputs is done with the usual `[PAD]` token, whereas the padding of the labels is done with a token with id `-100`, which is a special token automatically ignored by PyTorch loss functions.

```
1 data_collator = DataCollatorForSeq2Seq(tokenizer)
```

title\_prediction\_14.py hosted with ❤ by GitHub

[view raw](#)

Next, we download the ROUGE code using the `load_metric` function from the `datasets` library, thus instantiating a metric object. This object can then be used to compute its metrics using predictions and reference labels.

```
1 metric = load_metric("rouge")
```

title\_prediction\_16.py hosted with ❤ by GitHub

[view raw](#)

The `metric` object must then be called inside a `compute_metrics` function which takes a tuple of predictions and reference labels as input, and outputs a dictionary of metrics computed over the inputs. Specifically, the `compute_metrics` function does the following:

1. Decode the predictions (i.e. from token ids to words).
2. Decode the labels after substituting the `-100` token id with the `[PAD]` token id.
3. Compute ROUGE scores using the decoded predictions and labels, and select only a subset of these metrics.
4. Compute a new metric, which is the average length of the predictions.
5. Return a dictionary whose keys are the names of the metrics and the values are the metric values.

```
1 import numpy as np
2
3 def compute_metrics(eval_pred):
4     predictions, labels = eval_pred
5     decoded_preds = tokenizer.batch_decode(predictions, skip_special_tokens=True)
6
7     # Replace -100 in the labels as we can't decode them.
8     labels = np.where(labels != -100, labels, tokenizer.pad_token_id)
9     decoded_labels = tokenizer.batch_decode(labels, skip_special_tokens=True)
10
11    # Rouge expects a newline after each sentence
12    decoded_preds = ["\n".join(nltk.sent_tokenize(pred.strip()))
13                     for pred in decoded_preds]
14    decoded_labels = ["\n".join(nltk.sent_tokenize(label.strip()))
15                     for label in decoded_labels]
16
17    # Compute ROUGE scores
18    result = metric.compute(predictions=decoded_preds, references=decoded_labels,
19                           use_stemmer=True)
20
21    # Extract ROUGE f1 scores
22    result = {key: value.mid.fmeasure * 100 for key, value in result.items()}
23
24    # Add mean generated length to metrics
25    prediction_lens = [np.count_nonzero(pred != tokenizer.pad_token_id)
26                       for pred in predictions]
--
```

Open in app ↗



Search Medium



v

We are almost done! We must now create a `Seq2SeqTrainer` passing all the objects that we have just defined: the training arguments, the training and evaluation data, the data collator, the tokenizer, the `compute_metrics` function, and a `model_init` function. The `model_init` function must return a fresh new instance of the pre-trained model to fine-tune, making sure that training always starts from the same model and not from a partially fine-tuned model from your notebook.

## Start TensorBoard

Before starting the training, let's start TensorBoard inside our notebook.

After executing this code, the notebook should display TensorBoard right under the cell.

### Fine-tune T5

Last, we start the fine-tuning with the `train` method.

Once started, the cell will output a progress bar indicating the number of batches processed. Moreover, at each evaluation cycle, we'll see also the values of our metrics on the validation set, which you can see in the next image.

Step	Training Loss	Validation Loss	Rougel	Rouge2	Rougel	Rougelsum	Gen Len
100	3.078800	2.602307	28.613900	13.896800	26.131600	26.240700	12.920800
200	2.742800	2.530881	30.788500	15.023400	28.287600	28.410200	11.877800
300	2.648500	2.502835	31.436200	15.887600	28.822900	28.905500	12.254900
400	2.590400	2.490879	31.723300	15.812300	29.318900	29.450900	11.676400
500	2.627200	2.472070	31.540400	16.033300	29.223600	29.342000	12.555300
600	2.706000	2.461638	32.390200	16.379700	29.985600	30.085700	11.791600
700	2.543400	2.460367	32.568600	17.125400	30.191000	30.297100	12.303800
800	2.692300	2.444572	31.902500	16.248200	29.532600	29.574500	12.254900
900	2.614600	2.444440	32.006100	17.019600	29.778600	29.868100	12.414400
1000	2.656600	2.436043	32.975200	17.269300	30.662700	30.736300	11.963900
1100	2.581700	2.434549	32.681200	17.064600	30.525900	30.595300	11.890600
1200	2.620000	2.431621	33.092800	17.327100	30.797700	30.872500	11.823100
1300	2.642300	2.422981	32.717700	17.067500	30.504500	30.574900	11.866100
1400	2.613400	2.421857	33.249300	17.412100	30.872600	30.960200	11.927800
1500	2.597800	2.421307	33.832700	17.767500	31.448200	31.493900	11.759000
1600	2.615700	2.414079	32.462100	16.860800	30.334700	30.443300	12.116400
1700	2.631500	2.404803	33.371400	17.256600	30.927600	30.981800	12.196700
1800	2.614000	2.409261	33.338600	17.226700	30.928000	31.018100	11.669400
1900	2.619300	2.407189	32.919300	16.992700	30.614900	30.763400	12.171100
2000	2.580100	2.410542	32.455300	16.884800	30.378300	30.487800	12.314300
2100	2.585600	2.399441	33.016900	16.819200	30.621100	30.684500	11.927800

Training and evaluation metrics logged during fine-tuning.

The number of each step represents the number of batches processed. For example, at step 100, the model has been trained on `100 * batch_size` samples, which in this case is  $100 * 8 = 800$ .

Both the training and evaluation losses gradually decrease up to step 2100. However, all the ROUGE metrics seem to have peaked around step 1500 and later decreased a little. Since the training already took around three hours, I decided to stop the training and

test the results, but it's not impossible that the ROUGE scores could have risen again after some steps, thus making me miss a better model.

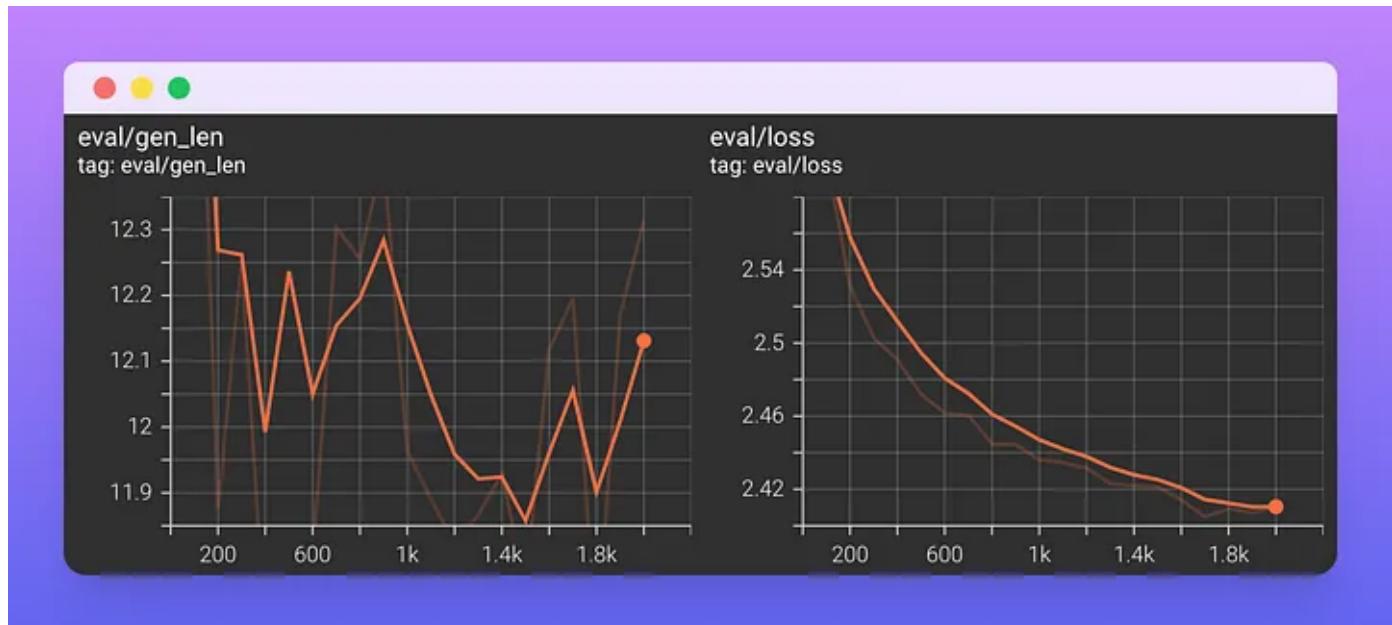
The average length of the generated titles has always remained stable at around 12 tokens. As each token with subword tokenizers is more or less long 0.75 English words, we can make a quick estimate that the average length of the generated articles is around 9 words (looks ok!).

Here are some charts extracted from TensorBoard.



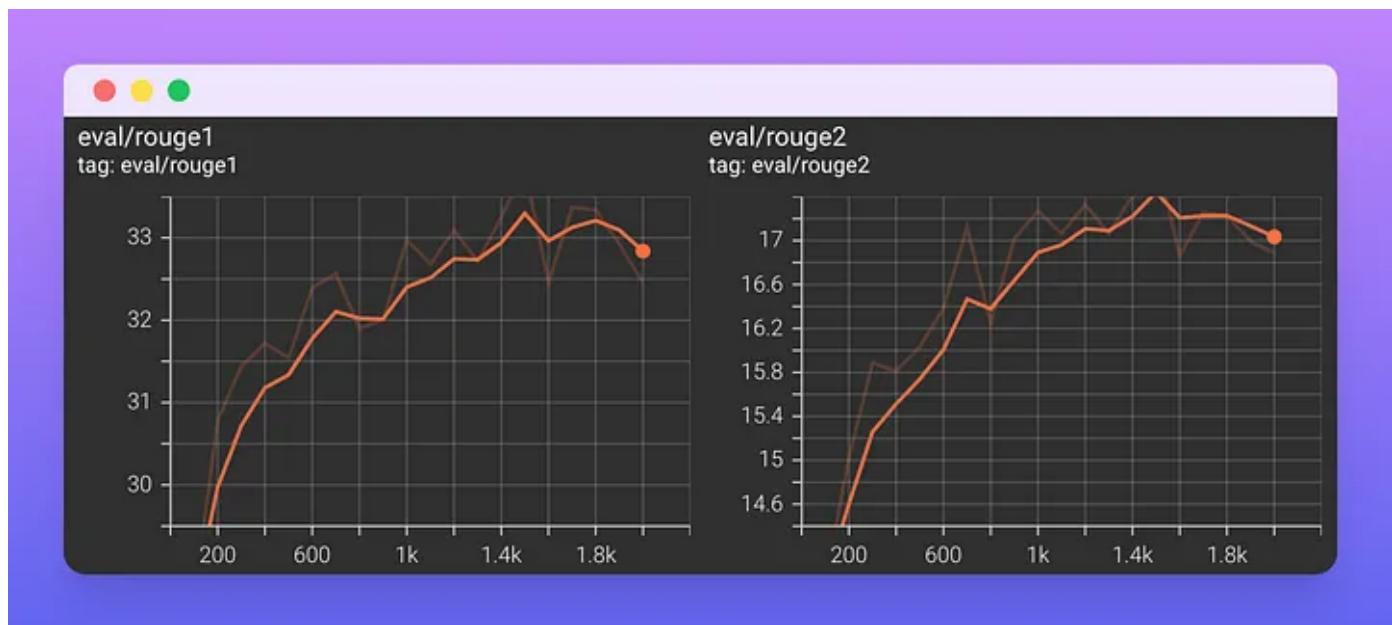
Training loss during fine-tuning.

The training loss seems to stabilize around step 1700, slowly decreasing.



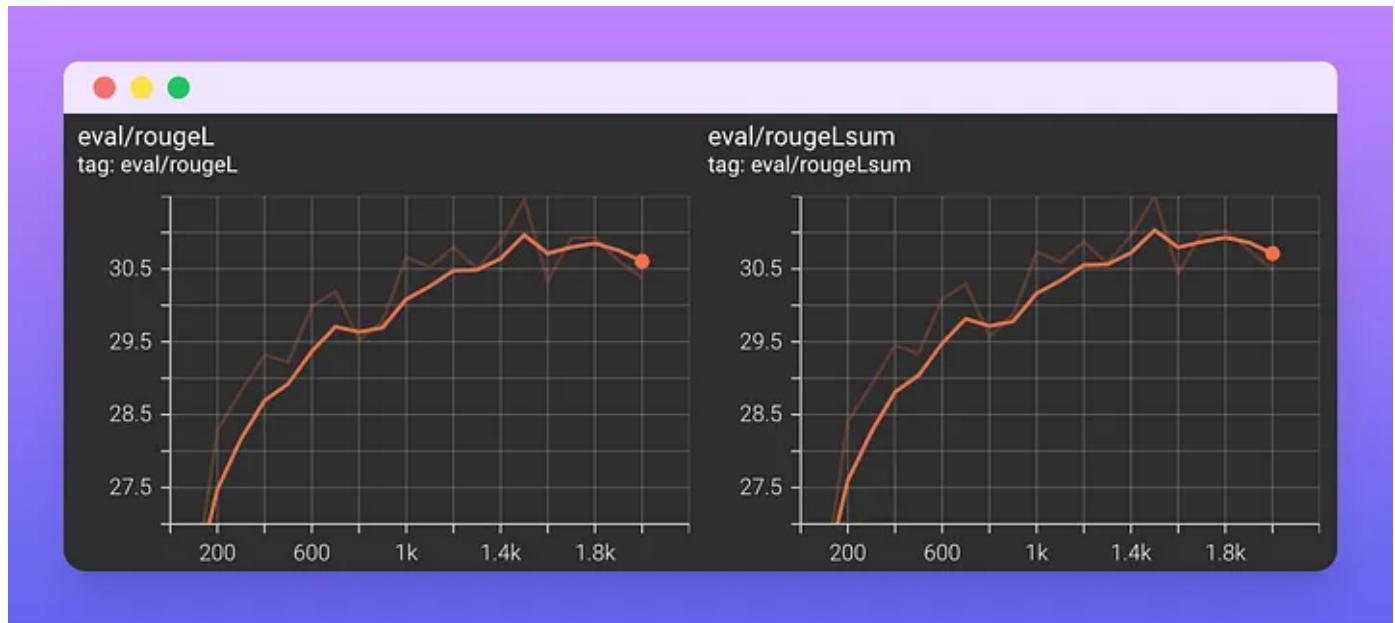
Average length of the generated titles and loss over the evaluation set during fine-tuning.

The evaluation loss is still decreasing at step 2000. The average length of the generated titles remains somewhat stable at 12 tokens during training.



ROUGE-1 and ROUGE-2 over the evaluation set during fine-tuning.

ROUGE-1 stabilizes at 33% at step 1800, whereas ROUGE-2 stabilizes at 17.5% at step 1500.



ROUGE-L and ROUGE-L-sum over the evaluation set during fine-tuning.

ROUGE-L and ROUGE-L-sum both stabilize at 30.7% at step 1500. While ROUGE-L computes the longest common subsequence (LCS) between the generated text and the reference text ignoring newlines, ROUGE-L-sum does the same considering newlines as sentence boundaries.

### Try the model

We can now test our model! First, let's load it.

Let's try it on the webpage "[Add statefulness to apps](#)" of the Streamlit documentation.

Our model predicts the title “Session State and Callbacks in Streamlit”. Looks good!

We test it on the Medium article “[Banking on Bots](#)” as a second example.

The model predicts the title “Conversational AI: The Future of Customer Service”, which is coherent with the article content and has proper capitalization as well. It looks like the model is working as expected.

### **Evaluate the model on the test set**

Last, let's compute the metrics on the test set with the new model to check that training has been done correctly and that we can expect good performance on new unseen data.

We need to:

1. Pad the articles texts to the same length.
2. Group the samples into batches (using a PyTorch dataloader).
3. Generate a title for each sample.
4. Compare the generated titles with the reference titles using the `compute_metrics` function defined earlier.

The metrics are better than the ones from the validation set. We probably need to increase the size of both the validation and test set to have less variable results, but let's consider these results satisfactory for now.

## Uploading the model to Hugging Face Hub

Loading the model from Google Drive is tedious since we need to mount it to the Colab filesystem every time. Let's see how to upload our model on Hugging Face Hub.

When uploading your model to the Hugging Face Hub, it's good practice to upload it with each one of the three main deep learning frameworks managed by Hugging Face: PyTorch, TensorFlow, and JAX. Thus, we need to install JAX to convert our model, as it's not installed in Colab by default.

Then, we need to authenticate to our Hugging Face account, as the uploaded models will be linked to it.

Last, we upload our model in the three formats. Remember to upload the tokenizer as well!

Once completed, this is how we can load the model and the tokenizer from the hub.  
Note that the name of the model is preceded by the name of the account linked to it.

Remember to update the model card of your uploaded model. This step consists in updating the README inside your model repository. Have a look at the final model card of our model!

## **Build an interactive demo with Streamlit and deploy it to Hugging Face Spaces**

Let's see again the line of code that we use to generate a candidate title for an article.

As you can see, we are using several parameters of the `generate` method to tune the results. A common parameter used when generating text is `temperature`, which controls the tradeoff between obtaining deterministically the most probable results and obtaining plausible different (but somewhat less probable) results every time we run the function. Indeed, we may want to generate 10 candidate titles instead of one.

Moreover, our model considers only the first 512 tokens of the article text when generating candidate titles. What if we want to generate multiple titles, where each title is generated taking into account 512 tokens from distinct parts of the article? We

would need to add custom prediction logic to the model, which is not easily testable directly from the Hugging Face Hub.

However, there's a way to allow users to experiment with these settings: we could build a small demo with this custom logic using Streamlit!

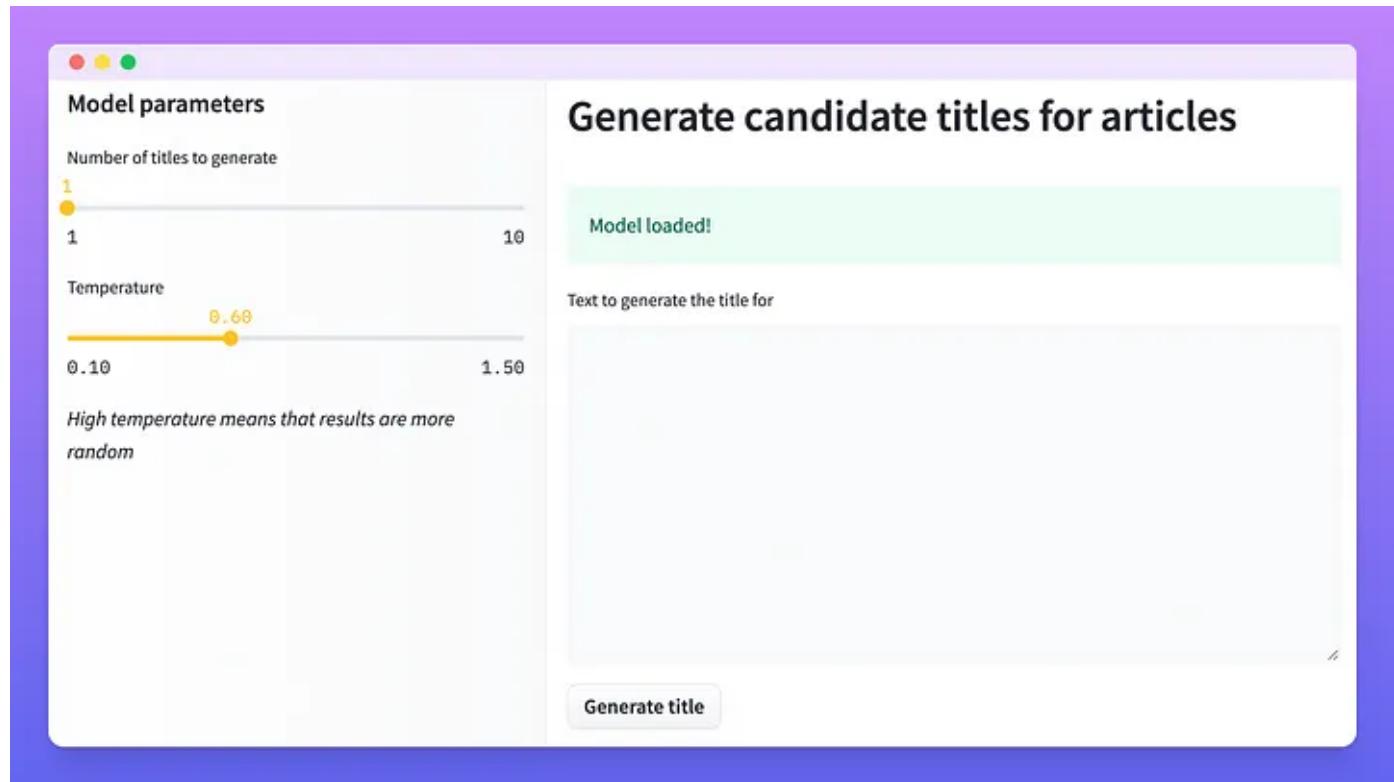
I'm not going into code details with Streamlit, but rest assured that it's very easy to use and quick to learn. I suggest reading the [Streamlit documentation](#), it takes a couple of hours.

Have a look at this [repo](#) to see the complete code of our Streamlit application. There are two main components:

- The `app.py` file: Contains the Streamlit code of the app. It's where we load our model from the Hugging Face Hub, create some interactive components, and write some logic to connect the two.
- The `requirements.py` file: Contains all the Python libraries used in the `app.py` file (such as the `transformers` library). It's not necessary to add the Streamlit library to this file.

Once we have tested locally our Streamlit application and everything works as expected, we need to upload it to a new Hugging Face Space. This step is just a matter of copying the files of your Streamlit applications to a new repo created on your Hugging Face account. Refer to this [article to learn how to host your Streamlit projects in Hugging Face Spaces](#).

This is the resulting [Hugging Face Space](#) which uses our title generation model in the background and allows generate multiple titles (taking into account different spans of the article text) and change the `temperature` parameter.



Let's try generating 5 candidate titles for the "[Banking on Bots](#)" Medium article, using a temperature of 0.6.

Model parameters

Number of titles to generate

1 5 10

Temperature

0.10 0.60 1.50

*High temperature means that results are more random*

Generate candidate titles for articles

Model loaded!

Text to generate the title for

Many financial institutions started building conversational AI, prior to the Covid19 pandemic, as part of a digital transformation initiative. These initial solutions were high profile, highly personalized virtual assistants — like the Erica chatbot from Bank of America. As the pandemic hit, the need changed as contact centers were under increased pressures. As Cathal McGloin of ServisBOT explains in "how it started, and how it is going," financial institutions were looking for ways to automate solutions to help get back to "normal" levels of customer service. This resulted in a change from the "future of conversational AI" to a real tactical assistant that can help in customer service. Haritha Dey of Wells Fargo saw a similar trend.

Generate title

**Generated titles**

Conversational AI: A Real Tactic Assistant for Customer Service

The Future of Conversational AI

Conversational AI: A Digital Transformation Approach

The Next Wave of Conversational AI

Conversational AI: The Future of Customer Service

The generated titles are:

- Conversational AI: A Real Tactic Assistant for Customer Service
- The Future of Conversational AI
- Conversational AI: A Digital Transformation Approach
- The Next Wave of Conversational AI
- Conversational AI: The Future of Customer Service

Building a small demo with Streamlit is very easy and, in conjunction with Hugging Face Spaces or Streamlit Cloud, allows other people to test your work without needing the knowledge of how to run a Jupyter notebook.

## Conclusion and next steps

In this article, we chose a suitable dataset and metric for our title generation task, and we wrote some code with the Hugging Face library to fine-tune a pre-trained T5 model for our task. We then uploaded our new model to the Hugging Face Hub so that everybody can use it, and we made a demo application with Streamlit which is now hosted as a Hugging Face Space.

I believe Hugging Face Space is an effortless way to build a portfolio of NLP applications where people or organizations can show off their machine learning skills.

Possible next steps are:

- Fine-tune a Longformer Encoder Decoder (LED) model instead of T5, as it is able to use a longer context as input. Keep in mind that the training will be slower though.
- Fine-tune a BART model and compare the results against the fine-tuned T5 model.

Thank you for reading! If you are interested in learning more about NLP, remember to follow NLPlanet on Medium, LinkedIn, Twitter, and join our new Discord server!

Data Science

Artificial Intelligence

NLP

Machine Learning

Hugging Face



Follow

## Written by **Fabio Chiusano**

3K Followers · Editor for NLPlanet

Freelance data scientist — Top Medium writer in Artificial Intelligence

More from Fabio Chiusano and NLPlanet

# Building a Knowledge Base from Texts

The diagram illustrates the workflow of extracting knowledge from text. It starts with a news article from Yahoo Finance about Microstrategy's CEO predicting Bitcoin's future. An arrow points from this article to a network graph visualization, which represents the resulting knowledge base extracted from the text.

 Fabio Chiusano in NLPlanet

## Building a Knowledge Base from Texts: a Full Practical Example

Implementing a pipeline for extracting a Knowledge Base from texts or online articles

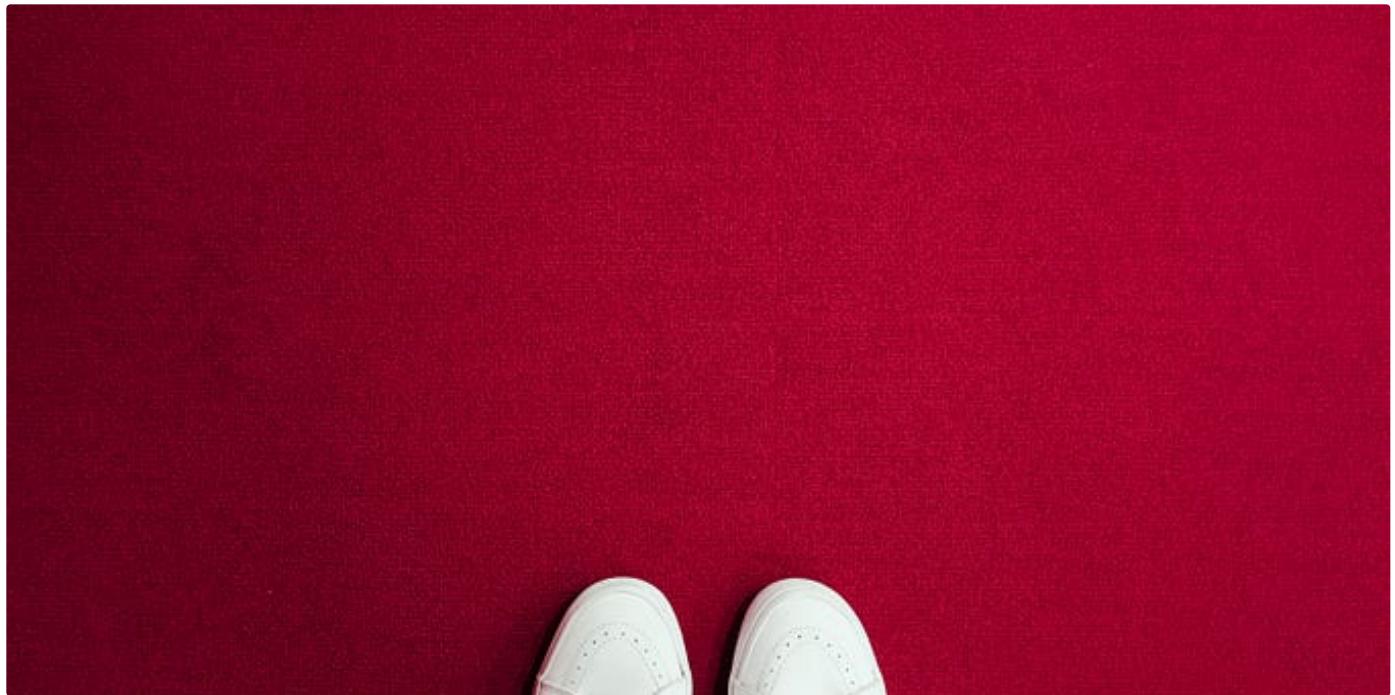
12 min read · May 24, 2022



942



...



Fabio Chiusano in NLPlanet

## Two minutes NLP—Learn the ROUGE metric by examples

ROUGE-N, ROUGE-L, ROUGE-S, pros and cons, and ROUGE vs BLEU

5 min read · Jan 19, 2022



107



...





Fabio Chiusano in NLPlanet

## Weekly AI and NLP News—May 16th 2023

100k tokens context windows, PaLM 2, and multimodal AI with six modalities

3 min read · May 16

👏 15



...



Fabio Chiusano in NLPlanet

## Weekly AI and NLP News—May 29th 2023

The new Intel LLM and Meta's project on +1,000 languages

5 min read · May 29

👏 54



...

See all from Fabio Chiusano

See all from NLPlanet

## Recommended from Medium

How is the costume design?

...The art direction by Randy Moore is great. The set decoration by K.C. Fox is great. The costume design by Susan Matheson is great. The make-up effects by Gregory Nicotero & Howard Berger is great. This is another great horror remake that is just as great as its original...



Skanda Vivek in Towards Data Science

### Fine-Tune Transformer Models For Question Answering On Custom Data

A tutorial on fine-tuning the Hugging Face RoBERTa QA Model on custom data and obtaining significant performance boosts

◆ · 5 min read · Dec 15, 2022

👏 373

💬 3

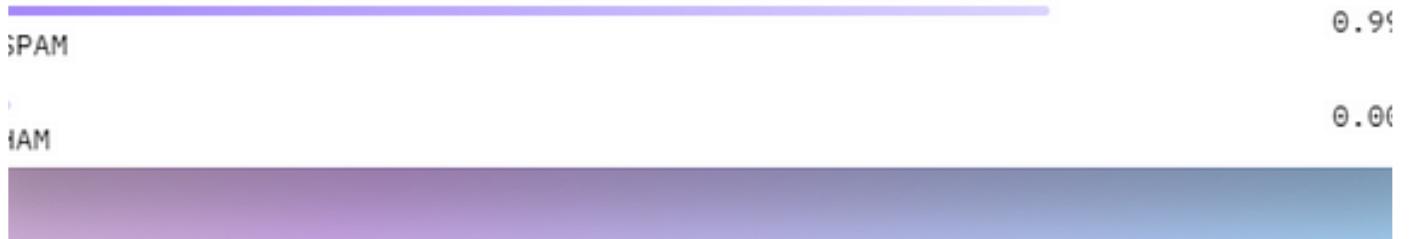
+

...

To redeem your free Amazon gift card click here!

## Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.024 s



Skanda Vivek in Towards Data Science

## Transformer Models For Custom Text Classification Through Fine-Tuning

A tutorial on how to build a spam classifier (or any other classifier) by fine-tuning the DistilBERT model

◆ · 4 min read · Jan 20

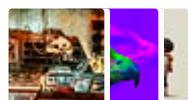
142

1



...

## Lists



### What is ChatGPT?

9 stories · 93 saves



### Staff Picks

342 stories · 101 saves



### Stories to Help You Level-Up at Work

19 stories · 84 saves

Is there a subscription fee for ChatGPT Plus?

**Extractive:** \$20

**Abstractive:** Yes, there is a \$20 subscription fee for ChatGPT Plus.

OpenAI is launching a premium and paid-for version of ChatGPT. The free app will remain available. But it is liable to go offline during busy periods – and, during those, the people who have paid its monthly fee will have priority access. That is just one of the perks offered in return for the \$20 subscription to “ChatGPT Plus”.



Skanda Vivek in Towards Data Science

## Extractive vs Generative Q&A—Which is better for your business?

The arrival of ChatGPT hints at a new era of search engines, this tutorial dives into the 2 basic types of AI based question answering

◆ · 6 min read · Feb 6

45



...

# FLAN-T5 RESOURCES

- TRY VIA UI
- USE VIA API
- FINE-TUNE

 Ramsri Goutham

## 5 Flan-T5 resources to try, deploy or fine-tune it

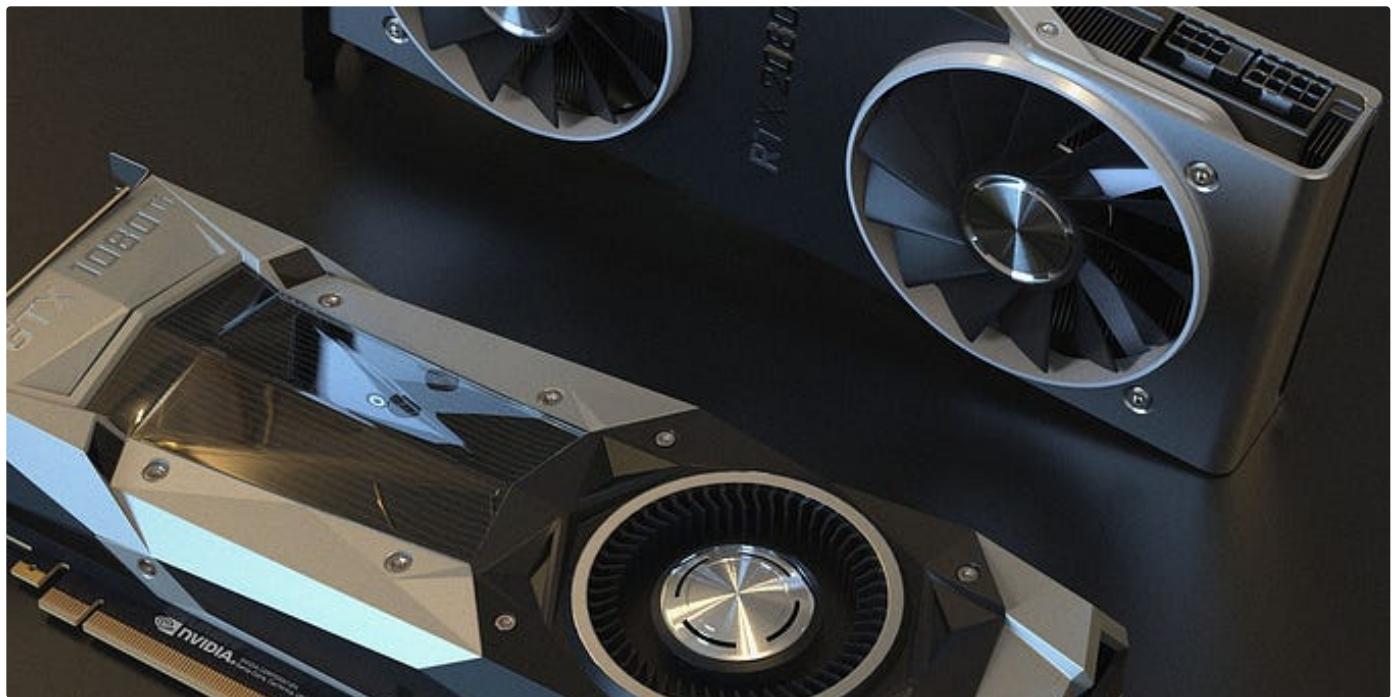
Google's Flan-T5 is the most practical open-source competitor to OpenAI's GPT-3

◆ · 2 min read · Feb 3

 42



...





Benjamin Marie in Towards AI

## Run Very Large Language Models on Your Computer

With PyTorch and Hugging Face's device\_map

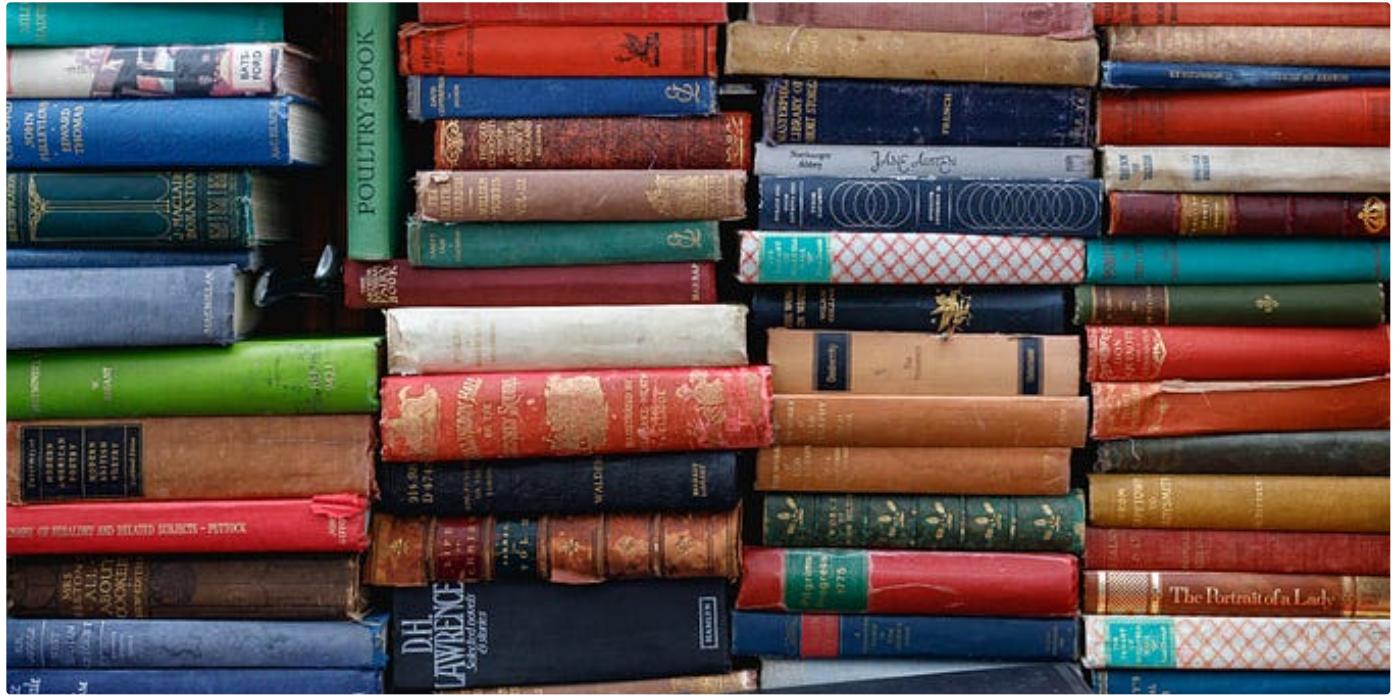
◆ · 5 min read · Dec 22, 2022

👏 215

💬 2



...



Jay Peterman in Towards Data Science

## Make a Text Summarizer with GPT-3

Quick tutorial using Python, OpenAI's GPT-3, and Streamlit

◆ · 11 min read · Jan 24

👏 154

💬 1



...

See more recommendations