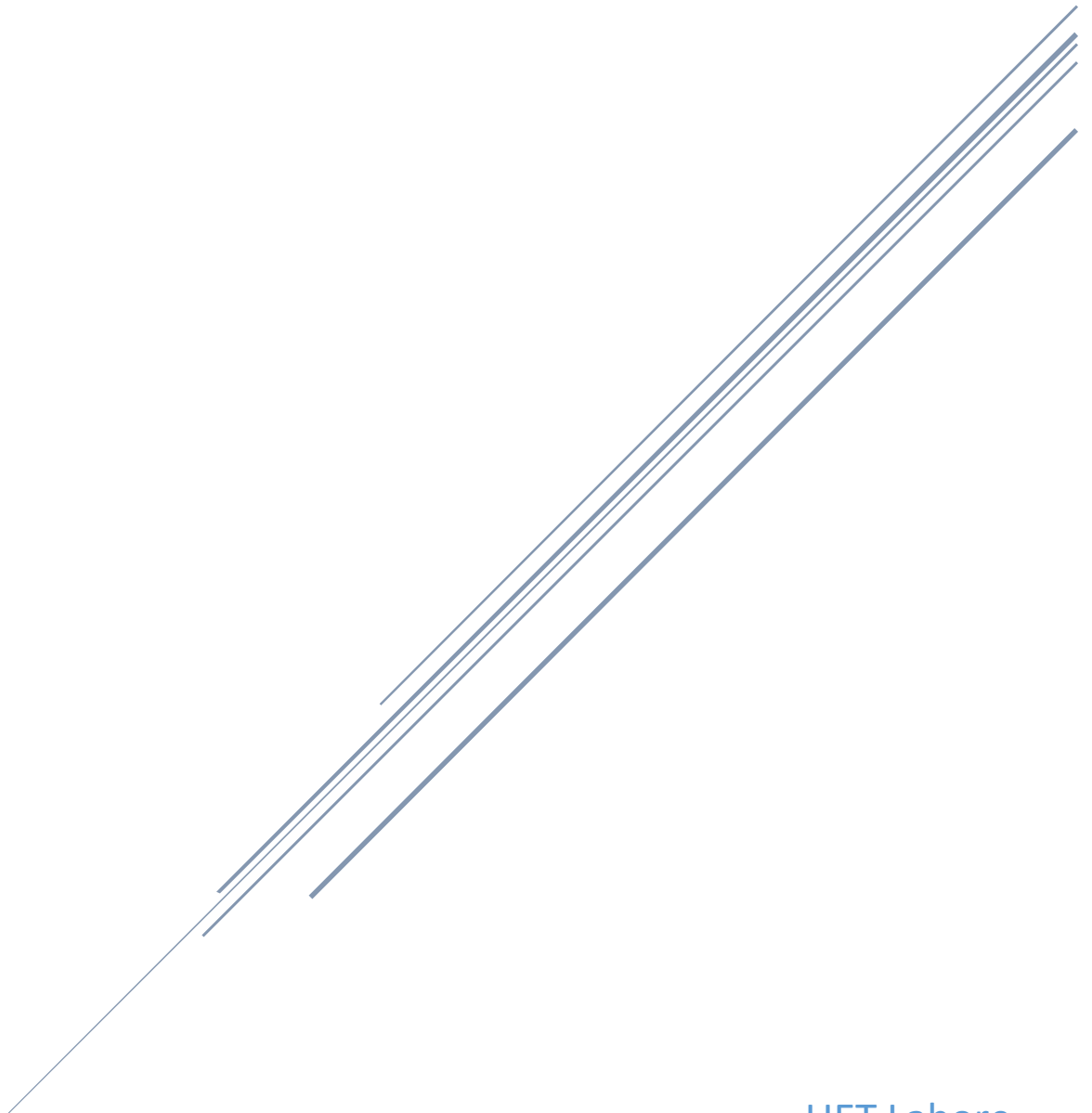# AAMIR ABBAS

Advanced Big Data Analytics

UET Lahore
REG# 2023 MSDS 04

**TASK 1**

**Q1.1 Respond with insights grounded in the terminology of Large Language Models.**

**(1). Can you provide a high-level overview of Transformers' architecture?**

**(2). What are the two approaches for evaluating language models in NLP, providing brief descriptions of each method along with highlighting their key distinctions?**

**(3). What is a token in the Large Language Models?**

**1). Overview of Transformers' Architecture:**

Ans: Transformers are a type of neural network architecture that has proven to be highly effective in natural language processing (NLP) tasks and other machine learning applications. The architecture was introduced in the paper "Attention is All You Need" by Vaswani et al. in 2017. Here's a high-level overview of the Transformer architecture:

1. **Self-Attention Mechanism:**

Allows the model to focus on different parts of the input sequence when making predictions.

2. **Multi-Head Attention:**

Uses multiple self-attention mechanisms to capture different aspects of the input.

3. **Positional Encoding:**

Adds information about the position of each word in the sequence.

4. **Feedforward Neural Network:**

Applies a simple neural network to process the attention outputs.

5. **Layer Normalization and Residual Connections:**

Helps in training deeper models by normalizing and connecting sub-layers.

6. **Encoder and Decoder Stacks:**

Multiple layers of the above components stacked together.

7. **Masked Self-Attention in Decoders:**

During training, prevents the model from looking ahead in the sequence.

8. **Output Layer:**

Produces the final predictions based on the processed information.

**2). What are the two approaches for evaluating language models in NLP, providing brief descriptions of each method along with highlighting their key distinctions?**

Ans: There are two ways to evaluate language models in NLP:

1. **Intrinsic Evaluation:**

What it does: Tests the model on specific language tasks like predicting the next word or recognizing parts of speech.

Why it's useful: Helps understand how well the model handles individual language aspects.

2. **Extrinsic Evaluation:**

What it does: Puts the model to the test in real-world applications, like translation or summarization.

Why it's useful: Shows how effective the model is in practical, broader language use.

**Key Difference:**

**Intrinsic:** Specific language tasks, direct metrics.

**Extrinsic:** Real-world applications, broader metrics.

Both evaluations together give a full picture of how good a language model is at understanding and using language.

---

**3). What is a token in the Large Language Models?**

Ans: Tokens are the fundamental building blocks that LLMs use to process and understand language. Like words in human language, they are the smallest meaningful units that LLMs can work with.

**Key points:**

- Not always individual words: Tokens can be words, but they can also be subwords, characters, or even entire phrases, depending on the model and task.
- Created through tokenization: The process of dividing text into tokens is called tokenization.
- Vital for language understanding: LLMs use tokens to represent and analyze the structure and meaning of language.
- Input and output: Tokens serve as both the input and output for LLMs during tasks like text generation and translation.
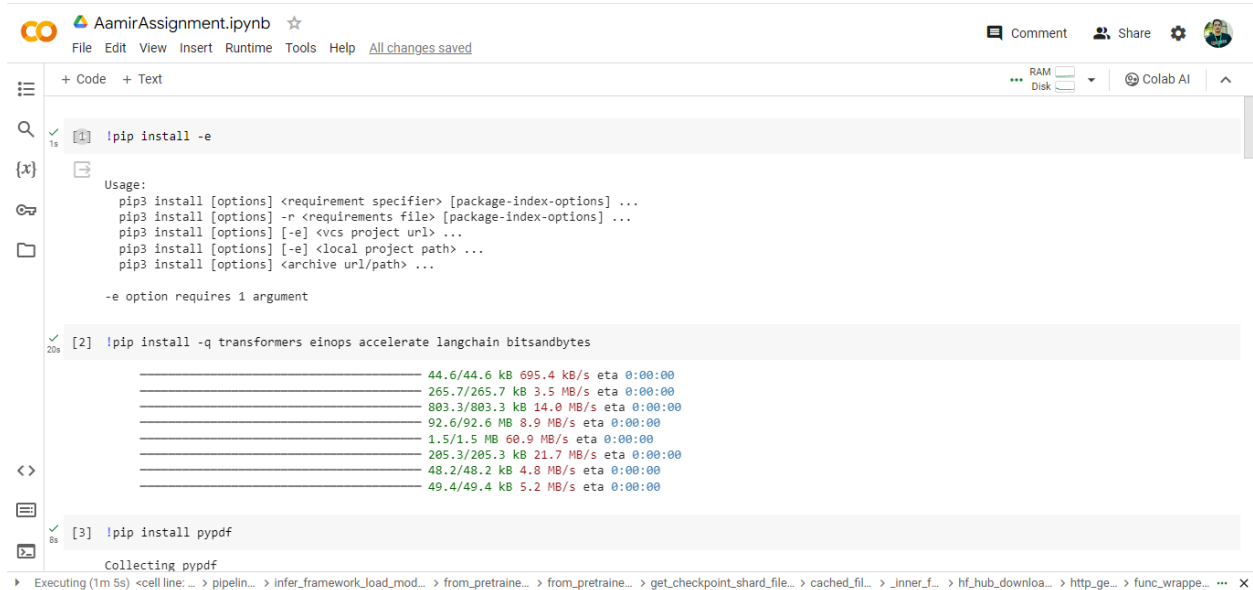
**Example:**

Consider the sentence "I love reading books." Here are possible tokenizations:

Word-based: ["I", "love", "reading", "books"]

---

**Q1.2 Read through the tutorial slides and deploy Llama 2 on Google Colab and get inference from the model. (10 pts)**

**(1). Provide screenshots of the results after you successfully download the model and see the text generated.**

File   Edit   View   Insert   Runtime   Tools   Help

+ Code   + Text

[10] `import torch`

[11] `import warnings`

[12] `model = "meta-llama/Llama-2-7b-chat-hf"`

[13] `Tokenizer=AutoTokenizer.from_pretrained(model)`

```
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:72: 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
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(
```

tokenizer_config.json: 100%      1.62k/1.62k [00:00<00:00, 35.8kB/s]

tokenizer.model: 100%      500k/500k [00:00<00:00, 6.91MB/s]

tokenizer.json: 100%      1.84M/1.84M [00:00<00:00, 18.9MB/s]

special_tokens_map.json: 100%      414/414 [00:00<00:00, 12.6kB/s]

Executing (2m 6s) <cell line: ...  pipelin...  infer_framework_load_mod...  from_pretraine...  from_pretraine...  get_checkpoint_shard_file...  cached_fil...  _inner_f...  hf_hub_downloa...  http_ge...  func_wrappe...

🤗 Hugging Face       Search models, datasets, users...          Models    Datasets    Spaces    Docs    Solutions    Pricing

Hugging Face is way more fun with friends and colleagues! 😊   Join an organization                          Dismiss this message

**Aamir Abbas**
`aamirsea`

Profile

Account

Organizations

Billing

**Access Tokens**

SSH and GPG Keys

Webhooks

Papers

**Access Tokens**

**User Access Tokens**

Access tokens programmatically authenticate your identity to the Hugging Face Hub, allowing applications to perform specific actions specified by the scope of permissions (read, write, or admin) granted. Visit the documentation to discover how to use them.
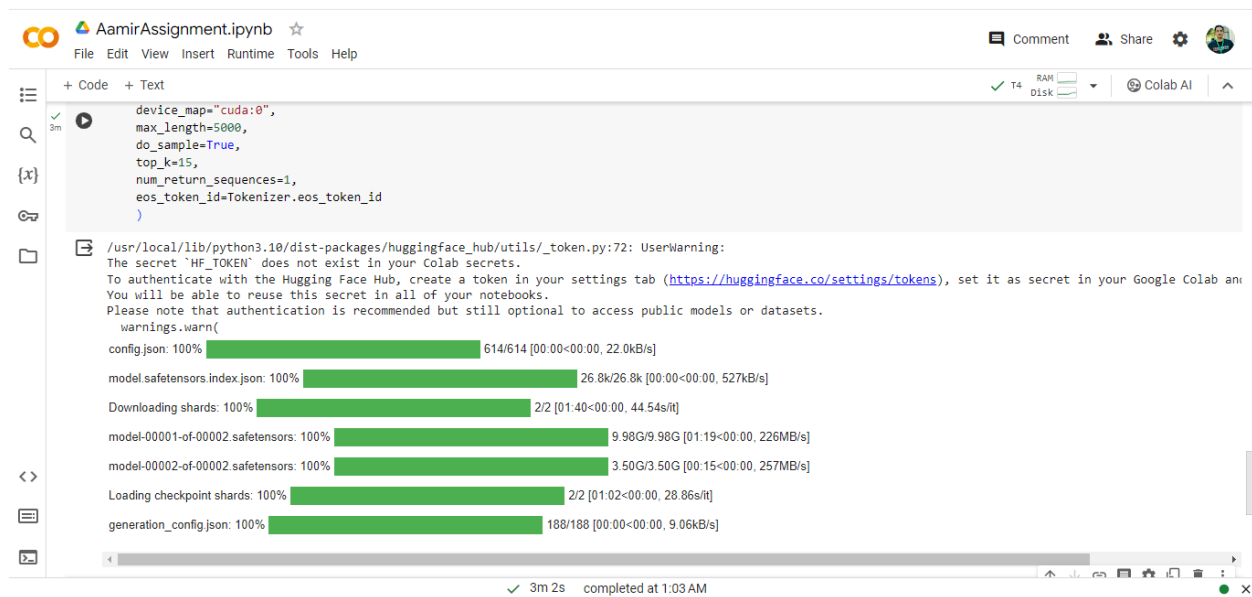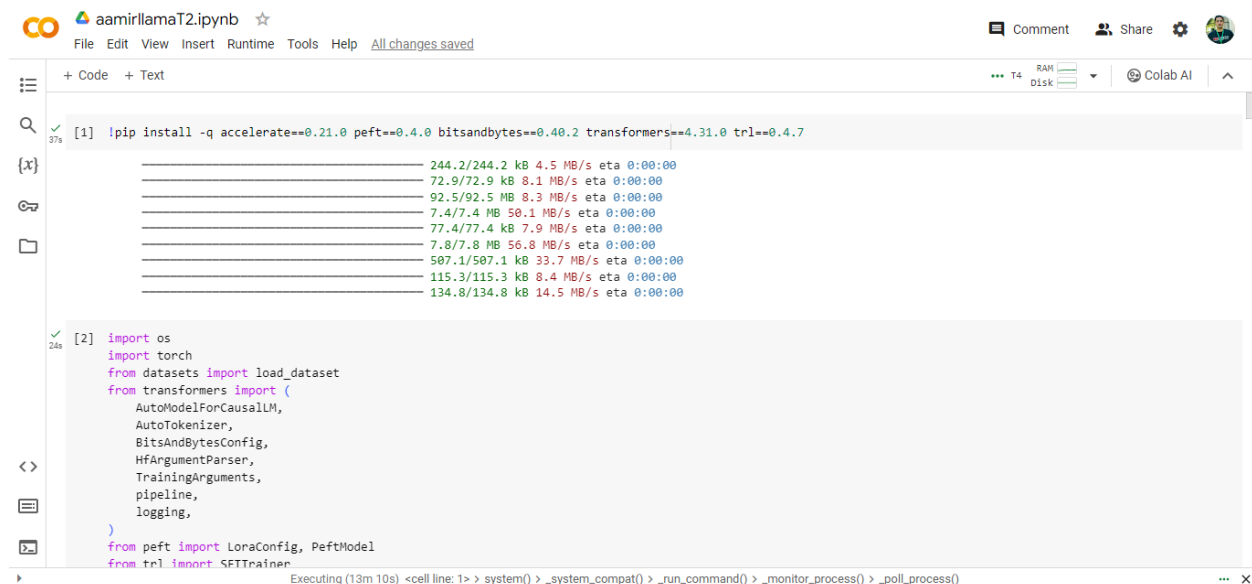
SeaToken  WRITE                                    Manage ˅

•••••••••••••••••••••••••••••••           Show 📋

New token

**(2). Change the "max_length" variable in pipline and observe the difference.**



**Task 2**

**In this part you are going to fine tuning Llama 2 models based on OpenAssistant dataset.**

**Q2.1 Write comments for each line of code and succintly explain what it is doing.**

[3] !huggingface-cli login

```
    _|    _|  _|    _|  _|_|_|_|  _|_|_|  _|_|_|  _|    _|  _|_|_|      _|_|_|_|  _|_|    _|_|_|    _|_|_|_|
    _|    _|  _|    _|  _|        _|        _|    _|_|  _|  _|    _|    _|        _|  _|  _|    _|  _|
    _|_|_|_|  _|    _|  _|_|_|    _|        _|    _|  _|_|  _|    _|    _|_|_|    _|_|_|_|  _|        _|_|_|
    _|    _|  _|    _|  _|        _|        _|    _|    _|  _|    _|    _|        _|    _|  _|        _|
    _|    _|    _|_|    _|        _|_|_|  _|_|_|  _|    _|  _|_|_|      _|        _|    _|    _|_|_|  _|_|_|_|

        To login, `huggingface_hub` requires a token generated from https://huggingface.co/settings/tokens .
    Token:
    Add token as git credential? (Y/n) Y
    Token is valid (permission: write).
    Cannot authenticate through git-credential as no helper is defined on your machine.
    You might have to re-authenticate when pushing to the Hugging Face Hub.
    Run the following command in your terminal in case you want to set the 'store' credential helper as default.

    git config --global credential.helper store

    Read https://git-scm.com/book/en/v2/Git-Tools-Credential-Storage for more details.
    Token has not been saved to git credential helper.
    Your token has been saved to /root/.cache/huggingface/token
    Login successful
```

[4] # The model that you want to train from the Hugging Face hub
    model_name = "meta-llama/Llama-2-7b-chat-hf"

Executing (17m 11s) <cell line: 1> > system() > _system_compat() > _run_command() > _monitor_process() > _poll_process()    ••• ✕

---

```
[4] # The model that you want to train from the Hugging Face hub
    model_name = "meta-llama/Llama-2-7b-chat-hf"

    # The instruction dataset to use
    dataset_name = "aamirsea/aamir-huggingface"

    # Fine-tuned model name
    new_model = "Llama-2-7b-chat-finetune"

    ################################################################################
    # QLoRA parameters
    ################################################################################

    # LoRA attention dimension
    lora_r = 64

    # Alpha parameter for LoRA scaling
    lora_alpha = 16

    # Dropout probability for LoRA layers
    lora_dropout = 0.1

    ################################################################################
    # bitsandbytes parameters
    ################################################################################

    # Activate 4-bit precision base model loading
```

Executing (17m 31s) <cell line: 1> > system() > _system_compat() > _run_command() > _monitor_process() > _poll_process()    ••• ✕

▤ Datasets: 🐟 aamirsea / **aamir-huggingface** 🗌  ♡ like  0

License: ⚖ llama2

🗌 Dataset card    ▤ **Files and versions**    🟠 Community    ⚙ Settings

⅄ main ⌄    aamir-huggingface       🐟 1 contributor    🕐 History: 2 commits    + Add file ⌄

| | | | | |
|---|---|---|---|---|
| 🐟 aamirsea | Upload train-00000-of-00001-9ad84bb9cf65a42f.parquet | 4c1d000 | | about 1 hour ago |
| 🗋 .gitattributes ⊘ | 2.31 kB | ⬇ | initial commit | about 1 hour ago |
| 🗋 README.md ⊘ | 24 Bytes | ⬇ | initial commit | about 1 hour ago |
| 🗋 train-00000-of-00001-9ad84bb9cf6... ⊘ | 967 kB 🟠 LFS | ⬇ | Upload train-00000-of-00001-9ad84bb9cf65a42f.parquet | about 1 hour ago |

---

CO  🔺 aamirllamaT2.ipynb ☆

File  Edit  View  Insert  Runtime  Tools  Help  All changes saved

💬 Comment    👥 Share  ⚙  🖼

+ Code  + Text         ••• T4  RAM —— Disk ——  ⌄  😊 Colab AI  ⌃

```
[6]                    # Apply the new template
                       reformatted_segments.append(f'<s>[INST] {human_text} [/INST] {assistant_text} </s>')
                else:
                       # Handle the case where there is no corresponding assistant segment
                       reformatted_segments.append(f'<s>[INST] {human_text} [/INST] </s>')

           return {'text': ''.join(reformatted_segments)}


       # Apply the transformation
       transformed_dataset = dataset.map(transform_conversation)
```

```
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:72: 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
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(
Downloading readme: 100%          24.0/24.0 [00:00<00:00, 1.30kB/s]
Downloading data: 100%            967k/967k [00:01<00:00, 514kB/s]
Generating train split:           1000/0 [00:00<00:00, 8838.11 examples/s]
Map: 100%                          1000/1000 [00:00<00:00, 2682.04 examples/s]
```

▶    Executing (18m 24s) <cell line: 1> > system() > _system_compat() > _run_command() > _monitor_process() > _poll_process()   ••• ✕

## Q2.2 Make a training loss plot.





| Step | Training Loss |
|------|---------------|
| 25 | 1.411500 |
| 50 | 1.661400 |
| 75 | 1.230700 |
| 100 | 1.463800 |
| 125 | 1.190200 |
| 150 | 1.383200 |
| 175 | 1.185500 |
| 200 | 1.482900 |
| 225 | 1.168800 |
| 250 | 1.552900 |

```
[8] %load_ext tensorboard
    %tensorboard --logdir results/runs
```

**Q2.3 Use the text generation pipeline to ask questions like "What is a large language model?"**



**Q 2.4 Store fine-tuning Llama2 Model and push Model to your Hugging Face Hub. Provide screenshot of your Hugging Face model page.**

A token is already saved on your machine. Run `huggingface-cli whoami` to get more information or `huggingface-cli logout` if you want to log out.
Setting a new token will erase the existing one.
To login, `huggingface_hub` requires a token generated from https://huggingface.co/settings/tokens .
Token:
Add token as git credential? (Y/n) Y
Token is valid (permission: write).
Cannot authenticate through git-credential as no helper is defined on your machine.
You might have to re-authenticate when pushing to the Hugging Face Hub.
Run the following command in your terminal in case you want to set the 'store' credential helper as default.

git config --global credential.helper store

Read https://git-scm.com/book/en/v2/Git-Tools-Credential-Storage for more details.
Token has not been saved to git credential helper.
Your token has been saved to /root/.cache/huggingface/token
Login successful
Upload 2 LFS files: 100%                                  2/2 [04:04<00:00, 244.99s/it]

pytorch_model-00002-of-00002.bin: 100%                    3.50G/3.50G [01:16<00:00, 40.6MB/s]

pytorch_model-00001-of-00002.bin: 100%                    9.98G/9.98G [04:04<00:00, 25.7MB/s]

CommitInfo(commit_url='https://huggingface.co/aamirsea/Llama-2-7b-chat-finetune/commit/031c298bff0342306f4096a26e167a774e572695', commit_message='Upload tokenizer', commit_description='', oid='031c298bff0342306f4096a26e167a774e572695', pr_url=None, pr_revision=None, pr_num=None)

Search models, datasets, users...

Models   Datasets   Spaces   Docs   Solutions   Pricing

**Aamir Abbas**

`aamirsea`

Edit profile   Settings

aamirsea110

🔬 **AI & ML interests**

None yet

📦 **Models** 1

🦙 aamirsea/Llama-2-7b-chat-finetune
Text Generation  ·  Updated 3 minutes ago
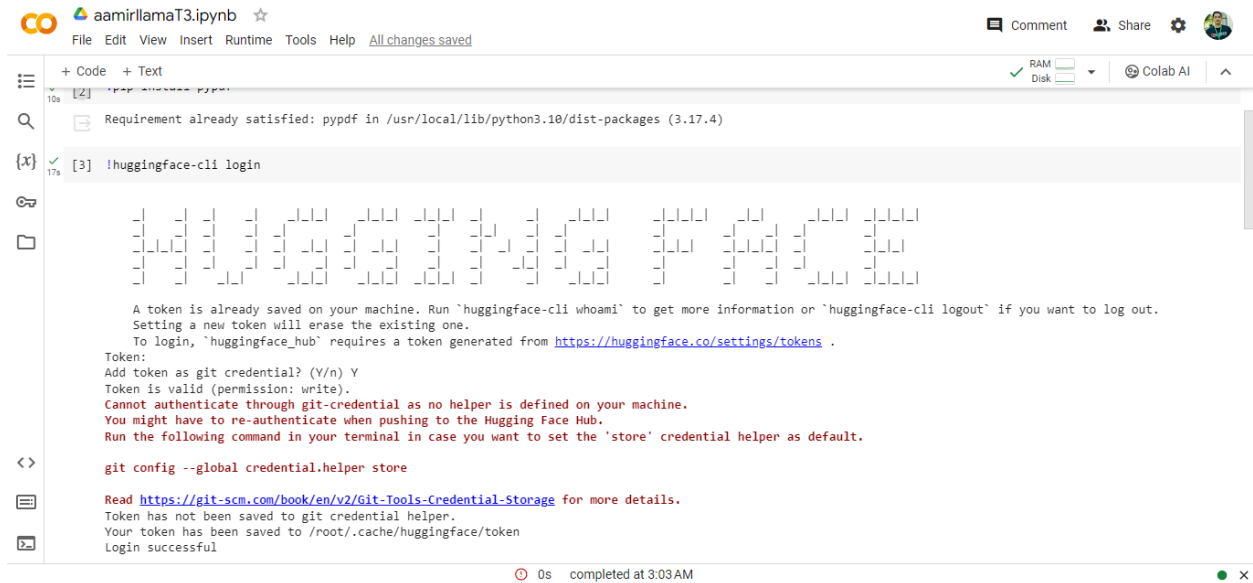
📊 **Datasets** 1

aamirsea/aamir-huggingface
⊞ Viewer  ·  Updated about 2 hours ago

**Task 3**

**In this section, your objective is to leverage the fine-tuning model from Task 2 to construct a versatile chatbot utilizing LangChain.**

**Q3.1, Provide screenshots of the prompt template you have devised.**

**Q3.2, Provide the text generation outcomes achieved through your chatbot.**

Comment   Share   ⚙   

+ Code   + Text

RAM / Disk   Colab AI

```python
[13]  input_prompt=prompt_template.format(book_name="Alchemist")
      print(input_prompt)
```

Privide me a concise summary of the book Alchemist

```python
[19]  prompt_template_name = PromptTemplate(
          input_variables =['product'],
          template = "What would be  a good name for a company that makes {product}"
      )
```

```python
[21]  from langchain.chains import ConversationChain

      convo = ConversationChain(llm=llm)
      print(convo.prompt.template)
```

```
The following is a friendly conversation between a human and an AI. The AI is talkative and provides lots of specific details from its context. If the AI does not

Current conversation:
{history}
Human: {input}
AI:
```

Start coding or generate with AI.

⊘ 0s   completed at 3:03 AM