# Single_node_finetuning_on_spr

## What is Finetuning?

**Fine-tuning** is a process in which a data is fed to the model and tells the models internal weights to get closer to responding how we would like it. For example we can fine tune a model for code generation , text generation and summarisation.In simple terms finetuning is a process in which we can train a model to do specific task by providing the right dataset

```python
from transformers import TrainingArguments
from intel_extension_for_transformers.neural_chat.config import (
    ModelArguments,
    DataArguments,
    FinetuningArguments,
    TextGenerationFinetuningConfig,
)
from intel_extension_for_transformers.neural_chat.chatbot import finetune_model
model_args = ModelArguments(model_name_or_path="meta-llama/Llama-2-7b-chat-hf")
data_args = DataArguments(train_file="alpaca_data.json", validation_split_percentage=1)
training_args = TrainingArguments(
    output_dir='./tmp',
    do_train=True,
    do_eval=True,
    num_train_epochs=3,
    overwrite_output_dir=True,
    per_device_train_batch_size=4,
    per_device_eval_batch_size=4,
    gradient_accumulation_steps=2,
    save_strategy="no",
    log_level="info",
    save_total_limit=2,
    bf16=True,
)
finetune_args = FinetuningArguments()
finetune_cfg = TextGenerationFinetuningConfig(
        model_args=model_args,
        data_args=data_args,
        training_args=training_args,
        finetune_args=finetune_args,
    )
finetune_model(finetune_cfg)
```
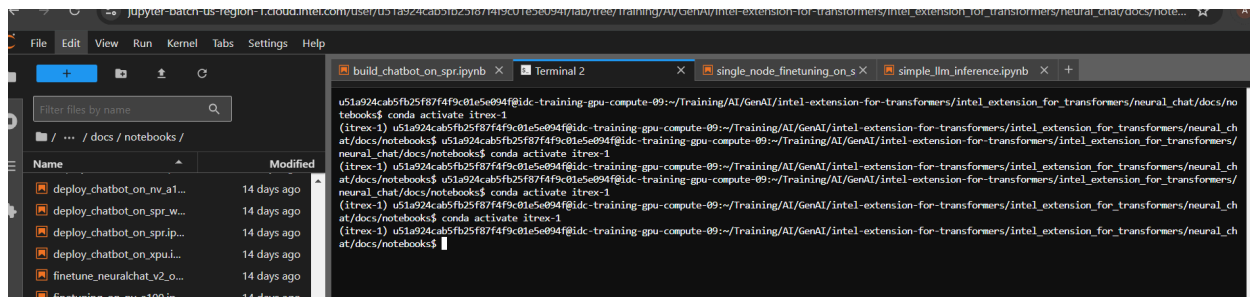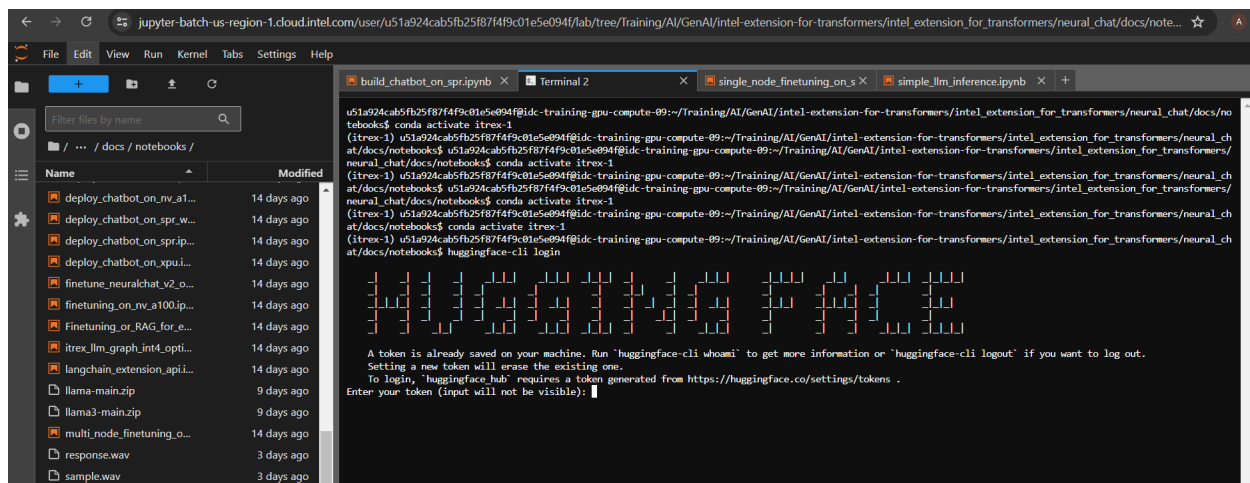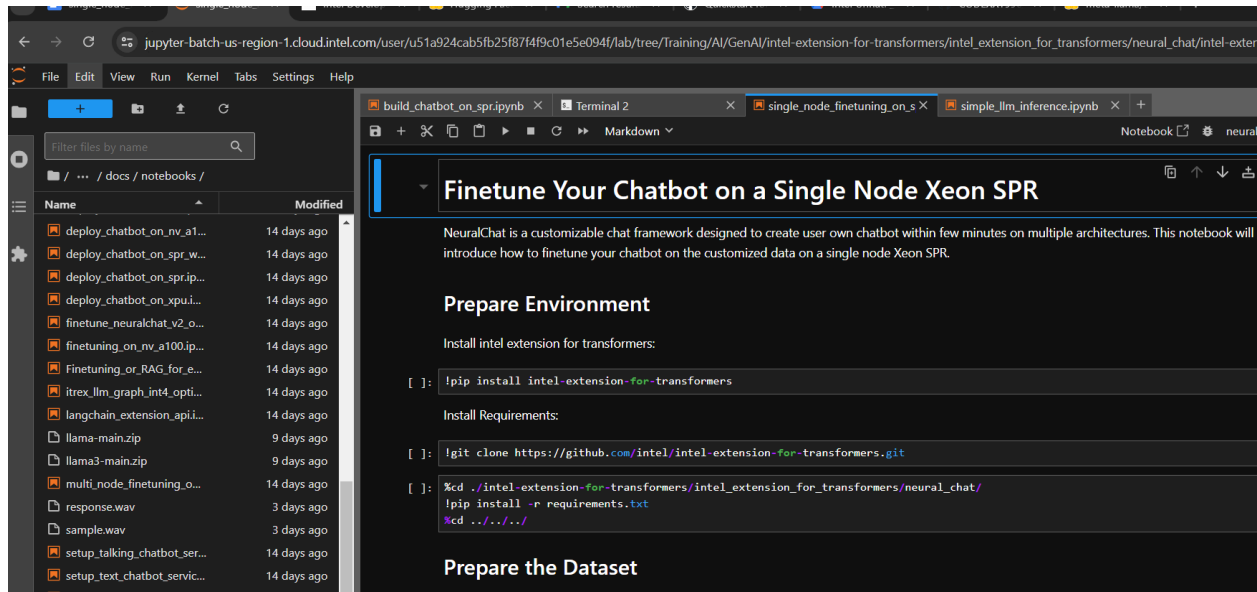
# Step 1:
## Activate the conda environment

File  Edit  View  Run  Kernel  Tabs  Settings  Help

build_chatbot_on_spr.ipynb  ×    Terminal 2  ×    single_node_finetuning_on_s  ×    simple_llm_inference.ipynb  ×  +

Filter files by name

/ ··· / docs / notebooks /

| Name | Modified |
|---|---|
| deploy_chatbot_on_nv_a1... | 14 days ago |
| deploy_chatbot_on_spr_w... | 14 days ago |
| deploy_chatbot_on_spr.ip... | 14 days ago |
| deploy_chatbot_on_xpu.i... | 14 days ago |
| finetune_neuralchat_v2_o... | 14 days ago |
| finetuning_on_nv_a100.ip... | 14 days ago |

```
u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/no
tebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/
neural_chat/docs/notebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/
neural_chat/docs/notebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$
```

# Step2:

## Activate the hugging face hub using huggingface-cli login

File  Edit  View  Run  Kernel  Tabs  Settings  Help

build_chatbot_on_spr.ipynb  ×    Terminal 2  ×    single_node_finetuning_on_s  ×    simple_llm_inference.ipynb  ×  +

Filter files by name

/ ··· / docs / notebooks /

| Name | Modified |
|---|---|
| deploy_chatbot_on_nv_a1... | 14 days ago |
| deploy_chatbot_on_spr_w... | 14 days ago |
| deploy_chatbot_on_spr.ip... | 14 days ago |
| deploy_chatbot_on_xpu.i... | 14 days ago |
| finetune_neuralchat_v2_o... | 14 days ago |
| finetuning_on_nv_a100.ip... | 14 days ago |
| Finetuning_or_RAG_for_e... | 14 days ago |
| itrex_llm_graph_int4_opti... | 14 days ago |
| langchain_extension_api.i... | 14 days ago |
| llama-main.zip | 9 days ago |
| llama3-main.zip | 9 days ago |
| multi_node_finetuning_o... | 14 days ago |
| response.wav | 3 days ago |
| sample.wav | 3 days ago |

```
u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/no
tebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/
neural_chat/docs/notebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/
neural_chat/docs/notebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ conda activate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-09:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_ch
at/docs/notebooks$ huggingface-cli login

    _|    _|  _|    _|    _|_|_|    _|_|_|  _|_|_|  _|      _|    _|_|_|        _|_|_|_|    _|_|     _|_|_|  _|_|_|_|
    _|    _|  _|    _|  _|        _|          _|    _|_|    _|  _|              _|        _|    _|  _|        _|
    _|_|_|_|  _|    _|  _|  _|_|  _|  _|_|    _|    _|  _|  _|  _|  _|_|        _|_|_|    _|_|_|_|  _|        _|_|_|
    _|    _|  _|    _|  _|    _|  _|    _|    _|    _|    _|_|  _|    _|        _|        _|    _|  _|        _|
    _|    _|    _|_|      _|_|_|    _|_|_|  _|_|_|  _|      _|    _|_|_|        _|        _|    _|    _|_|_|  _|_|_|_|

    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 .
Enter your token (input will not be visible):
```

# Step 3:
## Input the user token

Step 4:
Once the token is accepted the user has the id to use hugging face datasets

# Step 5:

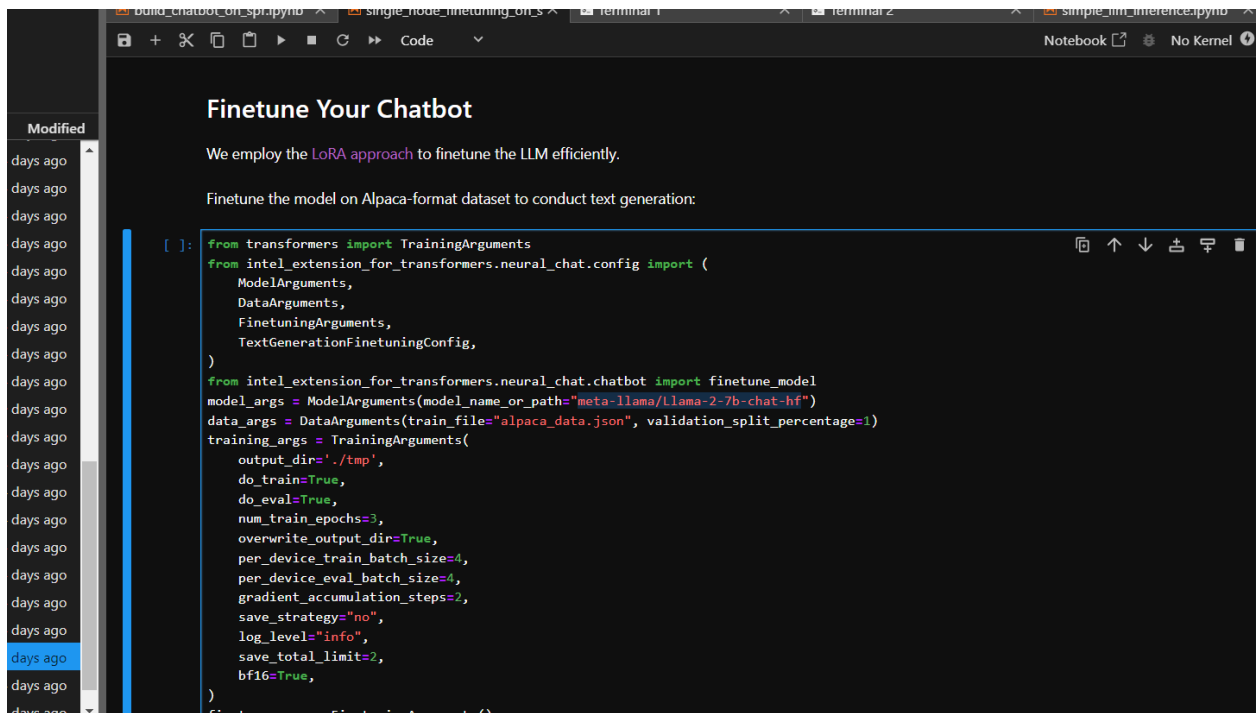Open the single_node_finetuning_on_spr on idc



# Step 6:

Inorder to finetune the 3 notebooks we need 3 different datasets required in the notebook . In the given notebook 3 different datasets are provided.     The model mention here is the meta-llama/Llama-2-7b-chat-hf

The model required to run the given notebook is
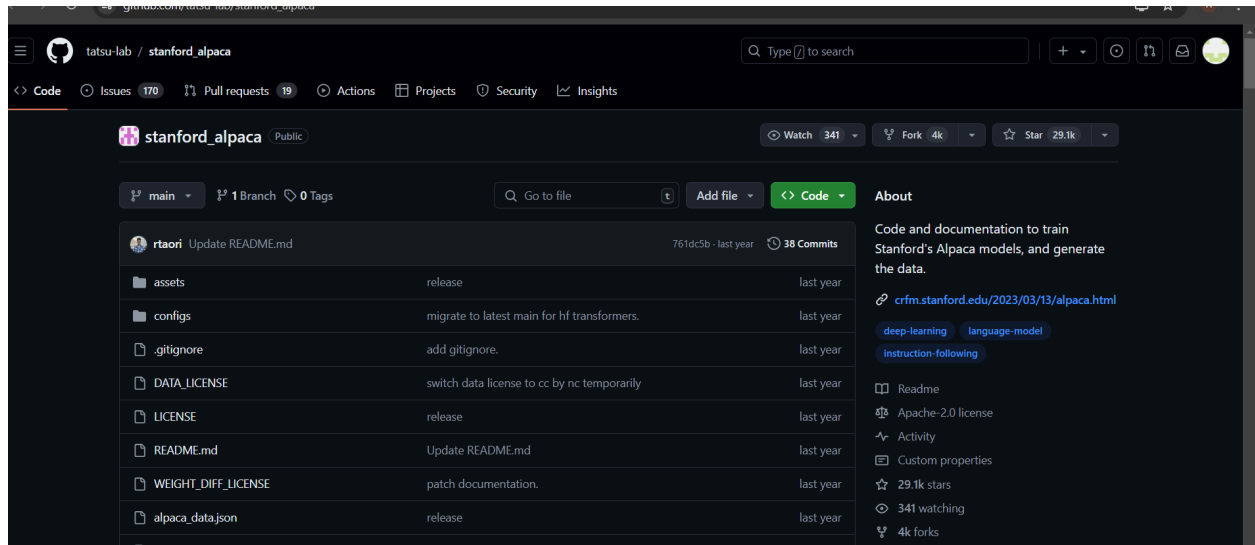
1. Alcapa dataset
2. Cnn_dailymail
3. theblackcat102/evol-codealpaca-v1

## Step 7:



In the given notebook we have specified the model name and path as well as the dataset required for the finetuning The dataset mentioned in the notebook is the alpaca dataset.

## Step 8:

Click the alpaca dataset in the notebook . Once its clicked it will redirect us to a github rep having the alpaca dataset

Step 9:

Inorder to use the dataset we need to download it as a zip file

Step 10:
 Once the file is downloaded it needs to be uploaded in idc or the local terminal for unzipping the zip file .



Step 11:

Once the file is uploaded  we can unzip  it in the terminal in idc server

| Name | Modified |
|---|---|
| deploy_chatbot_on_spr.ip... | 14 days ago |
| deploy_chatbot_on_xpu.i... | 14 days ago |
| finetune_neuralchat_v2_o... | 14 days ago |
| finetuning_on_nv_a100.ip... | 14 days ago |
| Finetuning_or_RAG_for_e... | 14 days ago |
| itrex_llm_graph_int4_opti... | 14 days ago |
| langchain_extension_api.i... | 14 days ago |
| llama-main.zip | 9 days ago |
| llama3-main.zip | 9 days ago |
| multi_node_finetuning_o... | 14 days ago |
| response.wav | 4 days ago |
| sample.wav | 4 days ago |
| setup_talking_chatbot_ser... | 14 days ago |
| setup_text_chatbot_servic... | 14 days ago |
| setup_text_chatbot_with_... | 14 days ago |
| single_card_finetuning_o... | 14 days ago |
| single_node_finetuning_o... | 2 days ago |
| spk_embed_default.pt | 4 days ago |
| stanford_alpaca-main.zip | 55 seconds ago |

Step 12:

# Unzip it in the idc terminal



# Step 13:

# Run the notebook

```python
from transformers import TrainingArguments
from intel_extension_for_transformers.neural_chat.config import (
    ModelArguments,
    DataArguments,
    FinetuningArguments,
    TextGenerationFinetuningConfig,
)
from intel_extension_for_transformers.neural_chat.chatbot import finetune_model
model_args = ModelArguments(model_name_or_path="meta-llama/Llama-2-7b-chat-hf")
data_args = DataArguments(train_file="alpaca_data.json", validation_split_percentage=1)
training_args = TrainingArguments(
    output_dir='./tmp',
    do_train=True,
    do_eval=True,
    num_train_epochs=3,
    overwrite_output_dir=True,
    per_device_train_batch_size=4,
    per_device_eval_batch_size=4,
    gradient_accumulation_steps=2,
    save_strategy="no",
    log_level="info",
    save_total_limit=2,
    bf16=True,
)
finetune_args = FinetuningArguments()
finetune_cfg = TextGenerationFinetuningConfig(
        model_args=model_args,
        data_args=data_args,
        training_args=training_args,
        finetune_args=finetune_args,
    )
finetune_model(finetune_cfg)
```

Step 14:

In certain cases the notebook may show an error
saying the huggingface token is not able to connect
to the meta-llama/Llama-2-7b-chat-hf model
The error msg indicates



```
warnings.warn(
2024-07-05 05:37:19,910 - chatbot.py - intel_extension_for_transformers.neural_chat.chatbot - ERROR - Exception: We couldn't connect to
'https://huggingface.co' to load this file, couldn't find it in the cached files and it looks like meta-llama/Llama-2-7b-chat-hf is not
the path to a directory containing a file named config.json.
Checkout your internet connection or see how to run the library in offline mode at 'https://huggingface.co/docs/transformers/installatio
n#offline-mode'.
2024-07-05 05:37:19,911 - error_utils.py - intel_extension_for_transformers.neural_chat.utils.error_utils - ERROR - neuralchat error: LO
RA finetuning failed
```

The only way to run the notebook is by loading the
meta lama mode in offline mode ; inorder to do so
we need to download the entire files of lama model



Or we can clone the entire meta lama repo

Inorder to clone this repo we need a Hugging face token with write permission



Step 16:

Now we need to login to the  hugging face idd using the write permitted token

```
u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-18:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks$ conda ac
tivate itrex-1
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-18:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks
$ huggingface-cli login

    _|    _|  _|    _|    _|_|_|    _|_|_|  _|_|_|  _|      _|    _|_|_|    _|_|_|    _|_|_|    _|_|_|
    _|    _|  _|    _|  _|        _|          _|    _|_|    _|  _|        _|        _|        _|
    _|_|_|_|  _|    _|  _|  _|_|  _|  _|_|    _|    _|  _|  _|  _|  _|_|  _|_|_|    _|        _|_|_|
    _|    _|  _|    _|  _|    _|  _|    _|    _|    _|    _|_|  _|    _|  _|        _|            _|
    _|    _|    _|_|      _|_|_|    _|_|_|  _|_|_|  _|      _|    _|_|_|  _|          _|_|_|  _|_|_|

    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 .
Enter your token (input will not be visible):
Add token as git credential? (Y/n) y
Token is valid (permission: write).
Your token has been saved in your configured git credential helpers (store).
Your token has been saved to /home/u51a924cab5fb25f87f4f9c01e5e094f/.cache/huggingface/token
Login successful
```

Step 17:

Since meta lama model repo huge we cannot clone it directly ; in order to clone it in windows  we need to install git lfs :

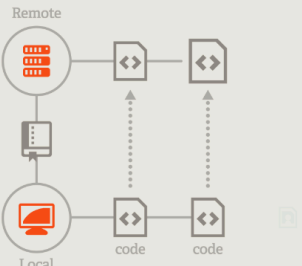(Git large file storage ) which can be done by manually installing it from the git website
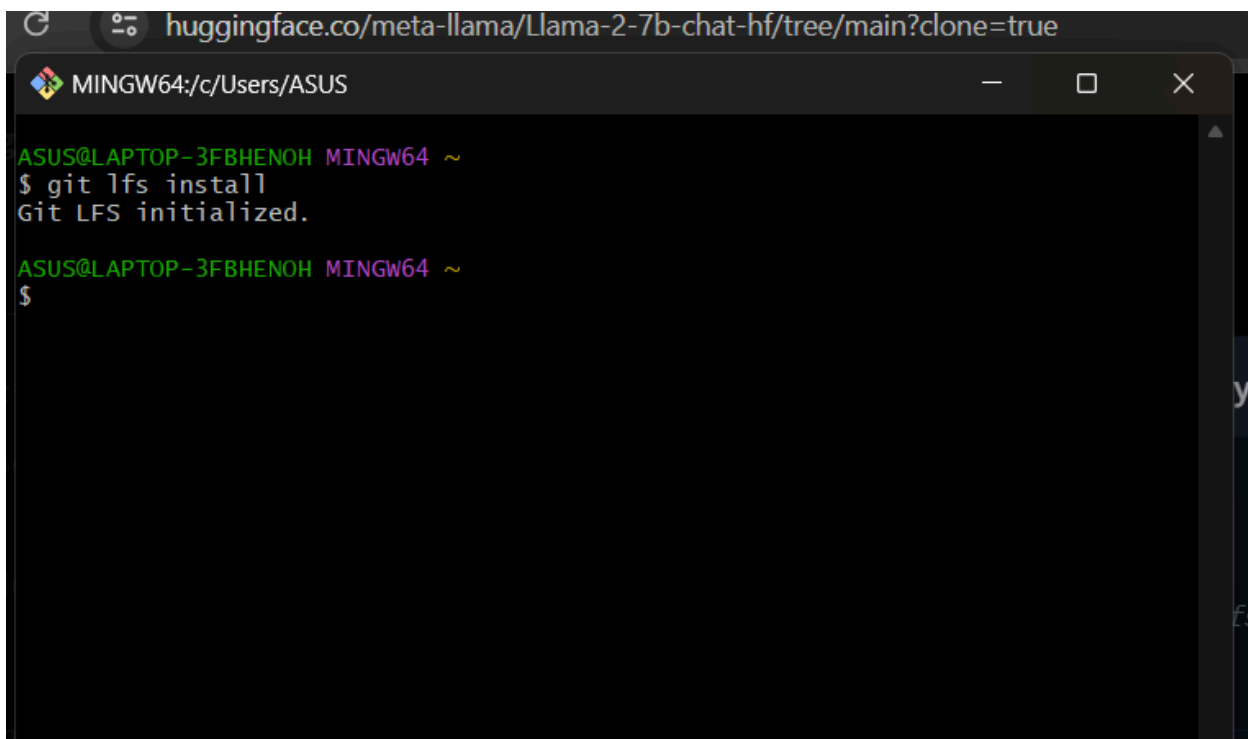


Or by clone git lfs repo into the terminal

```
Login successful
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-18:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks
$ git clone https://huggingface.co/meta-llama/Llama-2-7b-chat-hf
fatal: destination path 'Llama-2-7b-chat-hf' already exists and is not an empty directory.
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-18:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks
$ git clone https://github.com/git-lfs/git-lfs.git
Cloning into 'git-lfs'...
remote: Enumerating objects: 49113, done.
remote: Counting objects: 100% (294/294), done.
remote: Compressing objects: 100% (161/161), done.
remote: Total 49113 (delta 149), reused 243 (delta 133), pack-reused 48819
Receiving objects: 100% (49113/49113), 19.30 MiB | 20.74 MiB/s, done.
Resolving deltas: 100% (33940/33940), done.
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-18:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks
$
```

Step 18:
 Once the installation or cloning is completely done
In case of manual installation we may need to activate lfs
in git bash

```
huggingface.co/meta-llama/Llama-2-7b-chat-hf/tree/main?clone=true
MINGW64:/c/Users/ASUS

ASUS@LAPTOP-3FBHENOH MINGW64 ~
$ git lfs install
Git LFS initialized.

ASUS@LAPTOP-3FBHENOH MINGW64 ~
$
```

In case of  cloning lfs ; we can directly clone the meta
lama model

```
Llama-2-7b-cha...    2d ago
Llama-2-7b-hf       3d ago
llama-main          3d ago
stanford_alpaca...  10d ago
tmp                 3d ago
```

```
Add token as git credential? (Y/n) y
Token is valid (permission: write).
Your token has been saved in your configured git credential helpers (store).
Your token has been saved to /home/u51a924cab5fb25f87f4f9c01e5e094f/.cache/huggingface/token
Login successful
(itrex-1) u51a924cab5fb25f87f4f9c01e5e094f@idc-training-gpu-compute-18:~/Training/AI/GenAI/intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks
$ git clone https://huggingface.co/meta-llama/Llama-2-7b-chat-hf
fatal: destination path 'Llama-2-7b-chat-hf' already exists and is not an empty directory.
```

Step 19:
 Once the model is completely cloned or installed we can run the notebook with the  alpaca dataset in the correct directory . Similarly we can finetune 3 notebooks by accessing the hugging face tokens or by manually cloning or installing the required datasets in the correct directories