ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

DIPARTIMENTO DI
SCIENZE STATISTICHE "PAOLO FORTUNATI"

# 10 - Hugging Face and Transformers Library

Giovanni Della Lunga

giovanni.dellalunga@unibo.it

Halloween Conference in Quantitative Finance

Bologna - October 25-26-27, 2023

# Introduction

# What is Hugging Face?

- Hugging Face is a community specializing in Natural Language Processing (NLP) and artificial intelligence (AI).
- Founded in 2016, the company has made significant contributions to the field of NLP by democratizing access to state-of-the-art machine learning models and tools.

# What is Hugging Face?

- Hugging Face has a strong community focus.
- They provide a platform where researchers and developers can share their trained models, thereby fostering collaboration and accelerating progress in the field.
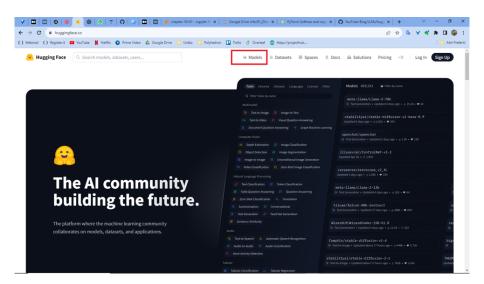
# Hugging Face Library
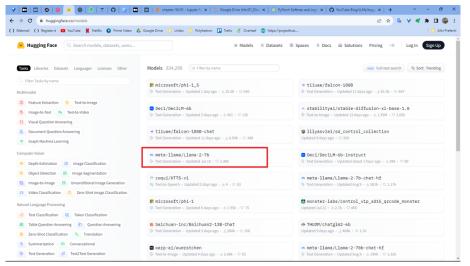
# Hugging Face Transformers as a Library

- Hugging Face's Transformers library is an open-source library for NLP and machine learning.
- It provides a wide variety of pre-trained models and architectures like BERT, GPT-X, T5, and many others.
- The library is designed to be highly modular and easy to use, allowing for the quick development of both research and production projects.
- It supports multiple languages and tasks like text classification, question-answering, text generation, translation, and more.

# Transformers

- Transformers is a Python library that makes downloading and training state-of-the-art ML models easy.
- Although it was initially made for developing language models, its functionality has expanded to include models for computer vision, audio processing, and beyond.
- Two big strengths of this library are: 1) it easily integrates with Hugging Face's Models, Datasets, and Spaces repositories, and 2) the library supports other popular ML frameworks such as PyTorch and TensorFlow.
- This results in a simple and flexible all-in-one platform for downloading, training, and deploying machine learning models and apps.
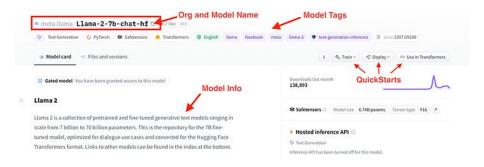
# Models

- There is a massive repository of pre-trained models available on Hugging Face (more than 250000 at the time of writing this).
- Almost all these models can be easily used via Transformers.
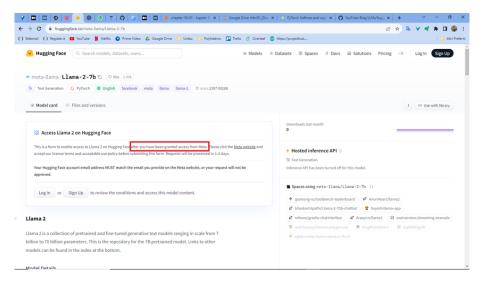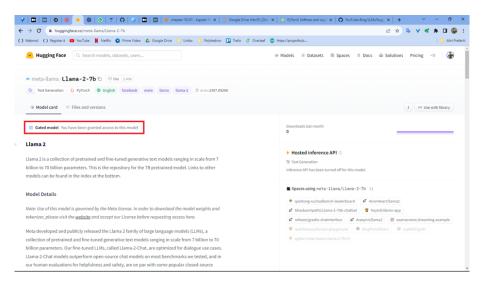
# Models

# Models

## Models

- To see what navigating the repository looks like, let's consider an example.
- Say we want a model that can do text generation, but we want it to be available via the Transformers library so we can use it in one line of code (as we did above).
- We can easily view all models that fit these criteria using the "Tasks" and "Libraries" filters.
- A model that meets these criteria is the newly released Llama 2. More specifically, Llama-2-7b-chat-hf, which is a model in the Llama 2 family with about 7 billion parameters, optimized for chat, and in the Hugging Face Transformers format.
- We can get more information about this model via its model card, which is shown in the next slide.

# Models

# Models

# Models

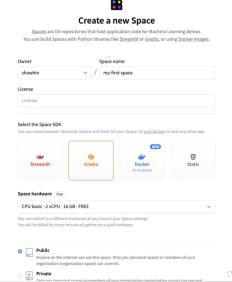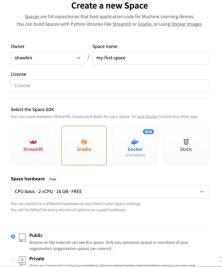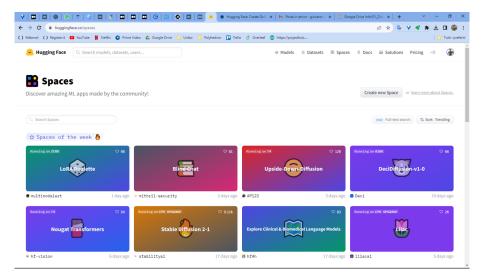# Hugging Face Spaces

# Hugging Face Spaces

- To go one step further, we can quickly deploy this UI via Hugging Face Spaces.

- These are Git repositories hosted by Hugging Face and augmented by computational resources.
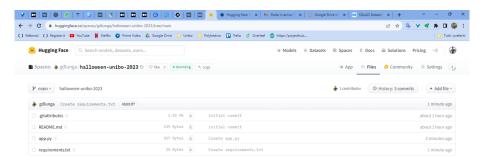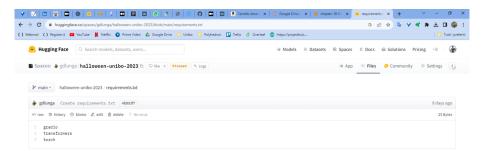
# Hugging Face Spaces

- To make a new Space, we first go to the Spaces page and click "Create new space".

- Then, configure the Space by giving it the name e.g. "my-first-space" and selecting Gradio as the SDK. Then hit "Create Space".



**Create a new Space**

Spaces are Git repositories that host application code for Machine Learning demos.
You can build Spaces with Python libraries like Streamlit or Gradio, or using Docker images.

Owner
shawhin

Space name
my-first-space

License
License

**Select the Space SDK**
You can chose between Streamlit, Gradio and Static for your Space. Or pick Docker to host any other app.

Streamlit    Gradio    Docker    Static
10 templates

**Space hardware** Free
CPU basic · 2 vCPU · 16 GB · FREE

You can switch to a different hardware at any time in your Space settings.
You will be billed for every minute of uptime on a paid hardware.

○ Public
Anyone on the internet can see this space. Only you (personal space) or members of your
organization (organization space) can commit.

○ Private

# Hugging Face Saces
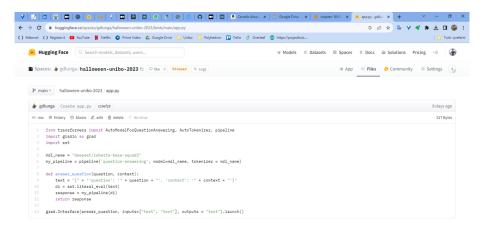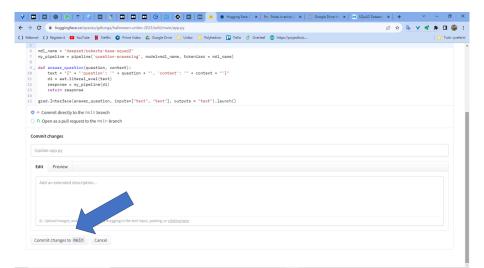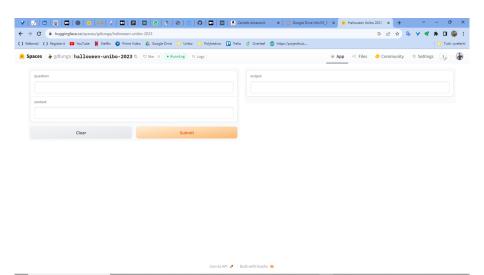
# Hugging Face Spaces

- Next, we need to upload **app.py** and **requirements.txt** files to the Space.
- **The app.py file houses the code we used to generate the Gradio UI, and the requirements.txt file specifies the app's dependencies**.
- Finally, we push the code to the Space just like we would to GitHub. The end result is a public application hosted on Hugging Face Spaces.

# Hugging Face Spaces

# Hugging Face Spaces

# Hugging Face Spaces

# Hugging Face Spaces

# Hugging Face Spaces

# Managing secrets and environment variables

- If your app requires environment variables (for instance, secret keys or tokens), do not hard-code them inside your app!
- Instead, go to the \*\*Settings\*\* page of your Space repository and add a new variable or secret.
- Use variables if you need to store non-sensitive configuration values and secrets for storing access tokens, API keys, or any sensitive value or credentials.

# Managing secrets and environment variables

# Managing secrets and environment variables

# Managing secrets and environment variables