# Introduction to Hugging Face Dataset

## AIE 1901 - AI Exploration: LLM for Optimization

Jie Wang, 2025/09/23

# Outline

- Introduction to Hugging Face

- Hands-on Lab: Exploring the Hugging Face Datasets/Models

References:
https://www.freecodecamp.org/news/get-started-with-hugging-face/#heading-how-to-find-the-right-pre-trained-model
https://huggingface.co/docs/datasets/main/en/tutorial
https://huggingface.co/docs/datasets/main/en/quickstart

# What is Hugging Face?

- A cute emoji

- A hub where data scientists, researchers, and ML engineers converge to exchange ideas, seek support, and contribute to open-source initiatives
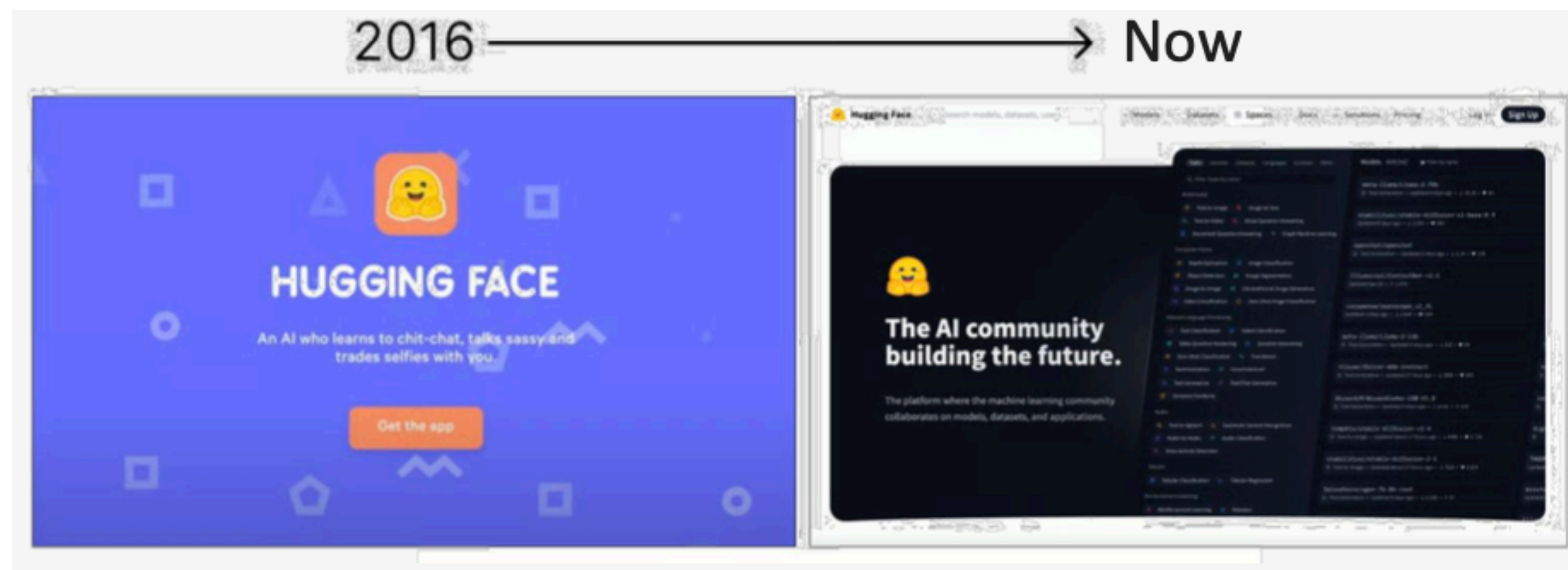
- In their own words, the mission is:

"Our mission at Hugging Face is to make Machine Learning as friendly and as productive as possible, for beginners and experts alike."

# What is Hugging Face?

- Providing a simple interface, easy to get started

- Committed to democratizing AI and making it accessible to a global community

# What is Hugging Face?

- Founded in 2016 (To develop an interactive AI chatbot targeted at teenagers)

- Moved to open source the AI models since 2018

  - "Transformers" library

# Funding & Valuation

- Hugging Face has a post-money valuation in the range of $1B to $10B as of Aug 24, 2023, according to PrivCo.

- $395.2M in funding over 7 rounds, funded by 38 investors

| Organization Name | Funding Type | Money Raised | Announced Date |
| --- | --- | --- | --- |
| 🤗 Hugging Face | Series D | — | Jan 16, 2024 |
| 🤗 Hugging Face | Series D | $235,000,000 | Aug 23, 2023 |
| 🤗 Hugging Face | Series C | $100,000,000 | May 9, 2022 |
| 🤗 Hugging Face | Series B | $40,000,000 | Mar 11, 2021 |
| 🤗 Hugging Face | Series A | $15,000,000 | Dec 17, 2019 |

# Business Models

- Most of Hugging Face's revenue is generated through its **Enterprise** solution.

- Provides a **premium subscription plan** for businesses requiring high-performance AI models and additional features

# Business Models

## PRO Account
PRO

Boost your personal HF experience

Subscribe for

**$9** per month

**Get PRO**

- ✓ 10× private storage capacity
- ✓ 20× included inference credits
- ✓ 8× ZeroGPU quota and highest queue priority
- ✓ Spaces Dev Mode & ZeroGPU Spaces hosting
- ✓ Publish blog articles on your HF profile
- ✓ Dataset Viewer for private datasets
- ✓ Show your support with a Pro badge

## Team
Instant setup for growing teams

Subscribe for

**$20** per user per month

**Get Team (via credit card)**

- ✓ SSO and SAML support
- ✓ Choose data location with Storage Regions
- ✓ Detailed action reviews with Audit Logs
- ✓ Granular access control via Resource Groups
- ✓ Repository usage Analytics
- ✓ Set auth policies and default repository visibility
- ✓ Centralized token control and approvals
- ✓ Dataset Viewer for private datasets
- ✓ Advanced compute options for Spaces
- ➕ All organization members get ZeroGPU and Inference Providers PRO benefits

## Enterprise
Custom onboarding and enterprise features

Starting at

**$50** per user per month

**Contact Sales**

- ➕ All benefits from the Team plan
- ✓ Highest storage, bandwidth, and API rate limits
- ✓ Managed billing with annual commitments
- ✓ Legal and Compliance processes
- ✓ Personalized support

# What Can You Do on the Hugging Face Platform?

- Download and fine-tune existing Open-Source models

# What Can You Do on the Hugging Face Platform?

- Run Models directly using Inference API

- If you don't want to build those models on your own machines, you can simply connect to these models, send requests, and receive outputs via API.



Check https://huggingface.co/ProsusAI/finbert

# What Can You Do on the Hugging Face Platform?



Create a new model repository

A repository contains all model files, including the revision history.

- Add/create your own model

- Host your model on the platform and make it public or private

# What Can You Do on the Hugging Face Platform?

- Use existing datasets

# What Can You Do on the Hugging Face Platform?

- Use existing datasets

# What Can You Do on the Hugging Face Platform?

- Upload your own dataset

- Host your own dataset on the platform and make it public or private

# What Can You Do on the Hugging Face Platform?

- Create/browse/try out demo apps in Spaces

- Here are a few cool spaces you can check out:

  - OpenAI's Whisper: Transcribe long-form microphone or audio inputs with the click of a button

  - AI Comic Factory: Create your own comic books.

  - QR Code AI Art Generator: Generate beautiful QR codes using AI.

  - Stable Video Diffusion (Img2Vid - XT): Generate 4s video from a single image.

  - Video-LLaMA: Audio-Visual Language Model for Video Understanding.

# What Can You Do on the Hugging Face Platform?

• Build yourself as an AI professional

1. **Join or create an organization**

    - You can join or create your own organization on Hugging Face. This allows you to showcase your work and collaborate with other members from your university, lab, or company.

2. **Create a portfolio**

    - You can create a professional portfolio on Hugging Face to showcase your work and start building your reputation. This can help you land jobs related to AI model training, integration, and development.

3. **Learn AI Skills**

    - Hugging Face is an excellent platform for learning AI skills. It offers a comprehensive set of tools and resources for training and using models

    - You can also learn from the experts and the community on Hugging Face, and improve your AI knowledge and skills

# Hugging Face Terminology

# How to Get Started with Hugging Face

1. Connect to the VPN

2. Create a Hugging Face Account

# How to Find the Right Pre-Trained Model

# How to Find the Right Dataset

Datasets 371 | Finance | Full-text search | ⇅ Sort: Most downloads

Linq-AI-Research/FinanceRAG
Viewer • Updated Sep 28, 2024 • 36.9k • 1.38k • ♡ 10

artefactory/Argimi-Ardian-Finance-10k-text
Preview • Updated 6 days ago • 1.31k

BAAI/IndustryCorpus_finance
Viewer • Updated Jul 26, 2024 • 32.6M • 968 • ♡ 10

PatronusAI/financebench
Viewer • Updated Nov 18, 2024 • 150 • 771 • ♡ 90

AIR-Bench/qa_finance_en
Viewer • Updated Sep 28, 2024 • 55.7k • 416 • ♡ 1

AdaptLLM/finance-tasks
Viewer • Updated Nov 30, 2024 • 23.3k • 303 • ♡ 69

gbharti/finance-alpaca
Viewer • Updated Sep 26, 2023 • 68.9k • 299 • ♡ 112

gretelai/synthetic_pii_finance_multilingual
Viewer • Updated Jun 11, 2024 • 55.9k • 258 • ♡ 55

BAAI/IndustryCorpus2_finance_economics
Viewer • Updated Nov 17, 2024 • 17M • 222 • ♡ 3

AIR-Bench/qrels-qa_finance_en-dev
Viewer • Updated Sep 28, 2024 • 1.75k • 189

4DR1455/finance_questions
Viewer • Updated Sep 5, 2024 • 53.9k • 168 • ♡ 4

sujet-ai/Sujet-Finance-QA-Vision-100k
Viewer • Updated Jul 14, 2024 • 9.8k • 166 • ♡ 31

# Explore LLM Datasets in Hugging Face

- **CardinalOperations/OR-Instruct-Data-3K**
  - Use Case: Instruction dataset to train large language models to learn optimization modeling and solving
  - Scope: **Applications from various industrial problems**
  - Size: 3000

https://huggingface.co/datasets/CardinalOperations/OR-Instruct-Data-3K

# Explore LLM Datasets in Hugging Face

- **CardinalOperations/OR-Instruct-Data-3K**
○ Use Case: Instruction dataset to train large language models to learn optimization modeling and solving
○ Scope: **Applications from various industrial problems**
○ Size: 3000

https://huggingface.co/datasets/CardinalOperations/OR-Instruct-Data-3K

- **CardinalOperations/IndustryOR**
○ Use Case: Evaluate LLM performance on optimization.
○ Scope: **a large-scale collection of industrial problems**
○ Size: 100
○ https://huggingface.co/datasets/CardinalOperations/IndustryOR

# Hands-on Lab1: Create your Google Account and Log Into Google Colab

- Please refer to 注册谷歌账号教程.pdf

# Hands-on Lab2: Try LLM Platform

- Open website https://cloud.siliconflow.cn/playground/chat/

- Create your own dataset

- Example Prompt: "give me 10 sentences that end with the word bird"

- Try it on different LLMs such as Owen2.5-7B-Instruct, Kimi-K2-Instruct-0905, DeepSeek-R1, DeepSeek-V3, Ling-flash 2.0

# Hands-on Lab3: Exploring the Hugging Face Datasets/Models

- Open the Colab Notebook: https://colab.research.google.com/drive/11tPb_t2KPnzCYTZpSleY9BffbM6U81p9?usp=sharing

- To work on this notebook, go to **File → Save a copy in Drive**

- Work on your copy, leaving this original unchanged

- All your work will be saved in your Google Drive