

# Introduction to Hugging Face

WORKING WITH HUGGING FACE



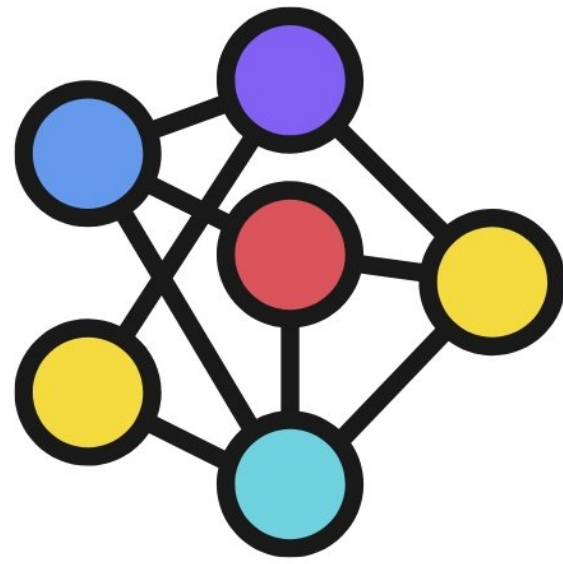
**Jacob H. Marquez**  
Lead Data Engineer

# The home of the AI community...



## Hugging Face

# The home of the AI community...

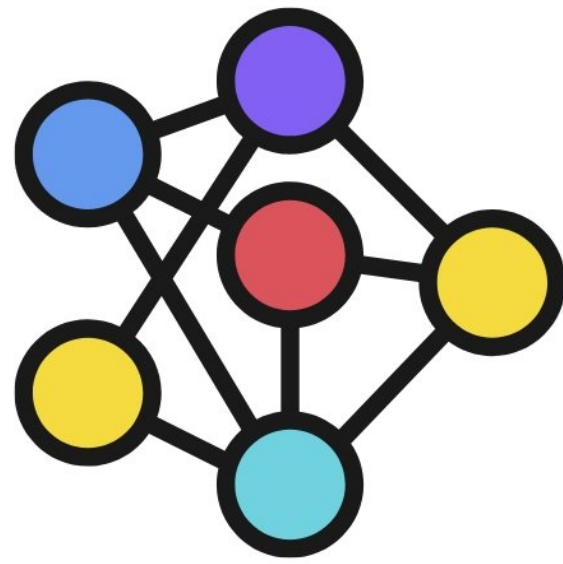


**Models**



**Hugging Face**

# The home of the AI community...



**Models**

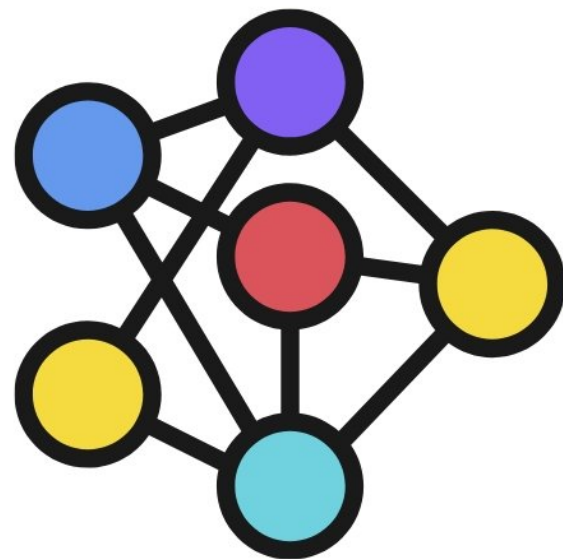


**Hugging Face**



**Datasets**

# The home of the AI community...



**Models**



**Hugging Face**



**Datasets**



**Applications**

<sup>1</sup> <https://huggingface.co/docs/hub/index>

## Core ML Libraries

- **Transformers**

State-of-the-art ML for PyTorch, TensorFlow, JAX

- **Diffusers**

State-of-the-art Diffusion models in PyTorch

- **Datasets**

Access & share datasets for any ML tasks

- **Transformers.js**

State-of-the-art ML running directly in your browser

- **Tokenizers**

Fast tokenizers optimized for research & production

- **Evaluate**

Evaluate and compare models performance

- **timm**

State-of-the-art vision models: layers, optimizers, and utilities

- **Sentence Transformers**

Embeddings, Retrieval, and Reranking

## Training & Optimization

- **PEFT**

Parameter-efficient finetuning for large language models

- **Accelerate**

Train PyTorch models with multi-GPU, TPU, mixed precision

- **Optimum**

Optimize HF Transformers for faster training/inference

- **AWS Trainium & Inferentia**

Train/deploy Transformers/Diffusers on AWS

- **TRL**

Train transformers LMs with reinforcement learning

- **Safetensors**

Safe way to store/distribute neural network weights

- **Bitsandbytes**

Optimize and quantize models with bitsandbytes








- **Lighteval**

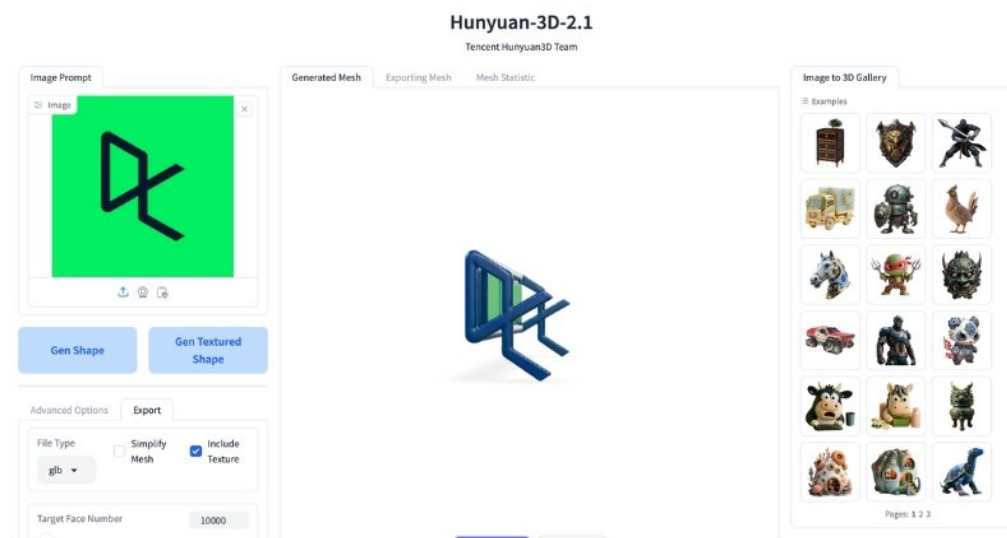
All-in-one toolkit to evaluate LLMs across multiple backends



















# Community and open-source heroes

## Models

 <b>deepseek-ai/DeepSeek-R1</b> Text Generation • 685B • Updated Mar 27 • 885k • 12.4k
 <b>black-forest-labs/FLUX.1-dev</b> Text-to-Image • Updated 10 days ago • 1.53M • 10.8k
 <b>CompVis/stable-diffusion-v1-4</b> Text-to-Image • Updated Aug 23, 2023 • 3.46M • 6.86k
 <b>stabilityai/stable-diffusion-xl-base-1.0</b> Text-to-Image • Updated Oct 30, 2023 • 2.56M • 6.71k
 <b>meta-llama/Meta-Llama-3-8B</b> Text Generation • 8B • Updated Sep 27, 2024 • 357k • 6.24k
 <b>bigscience/bloom</b> Text Generation • 176B • Updated Jul 28, 2023 • 6.66k • 4.92k
 <b>stabilityai/stable-diffusion-3-medium</b> Text-to-Image • Updated Aug 12, 2024 • 11.6k • 4.79k
 <b>hexgrad/Kokoro-82M</b> Text-to-Speech • Updated Apr 10 • 1.58M • 4.61k
 <b>openai/whisper-large-v3</b> Automatic Speech Recognition • 2B • Updated Aug 12, 2024 • 2.9M • 4.6k
 <b>meta-llama/Llama-2-7b-chat-hf</b> Text Generation • 7B • Updated Apr 17, 2024 • 1.05M • 4.48k



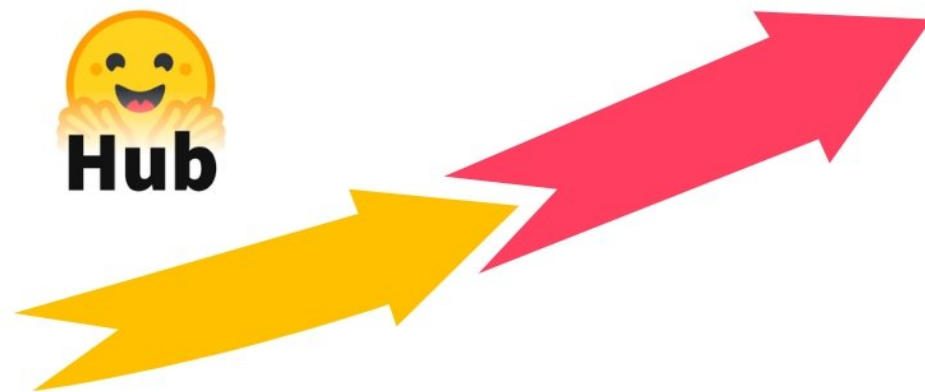
## Datasets

 <b>m-a-p/FineFineWeb</b> Viewer • Updated Dec 19, 2024 • 4.89B • 3.28M • 65	 <b>huggingface/documentation-images</b> Viewer • Updated 3 days ago • 52 • 3M • 73
 <b>permutans/fineweb-bbc-news</b> Viewer • Updated Jan 27 • 15.9M • 1.5M • 20	 <b>hallucinations-leaderboard/results</b> Updated Oct 31, 2024 • 1.42M • 2
 <b>huggingface/badges</b> Viewer • Updated Apr 8 • 1 • 1.39M • 45	 <b>jat-project/jat-dataset-tokenized</b> Viewer • Updated Dec 22, 2023 • 32M • 1.25M • 1
 <b>KakologArchives/KakologArchives</b> Updated 1 minute ago • 999k • 16	 <b>lavita/medical-qa-shared-task-v1-toy</b> Viewer • Updated Jul 20, 2023 • 64 • 911k • 19
 <b>huggingchat/models-logo</b> Viewer • Updated May 27 • 13 • 905k • 4	 <b>adams-story/datacomp200m</b> Viewer • Updated Jul 19, 2023 • 213M • 793k • 1
 <b>TAUR-Lab/Taur_CoT_Analysis_Project_gpt--</b> Viewer • Updated Oct 19, 2024 • 82.2k • 741k	 <b>Salesforce/wikitext</b> Viewer • Updated Jan 4, 2024 • 3.71M • 691k • 469
 <b>openai/gsm8k</b> Viewer • Updated Jan 4, 2024 • 17.6k • 513k • 789	 <b>nvidia/PhysicalAI-Robotics-GR00T-X-Embodi...</b> Updated May 15 • 475k • 133
 <b>agents-course/course-images</b> Viewer • Updated 17 days ago • 2 • 468k • 13	 <b>princeton-nlp/SWE-bench_Verified</b> Viewer • Updated Feb 18 • 500 • 438k • 183

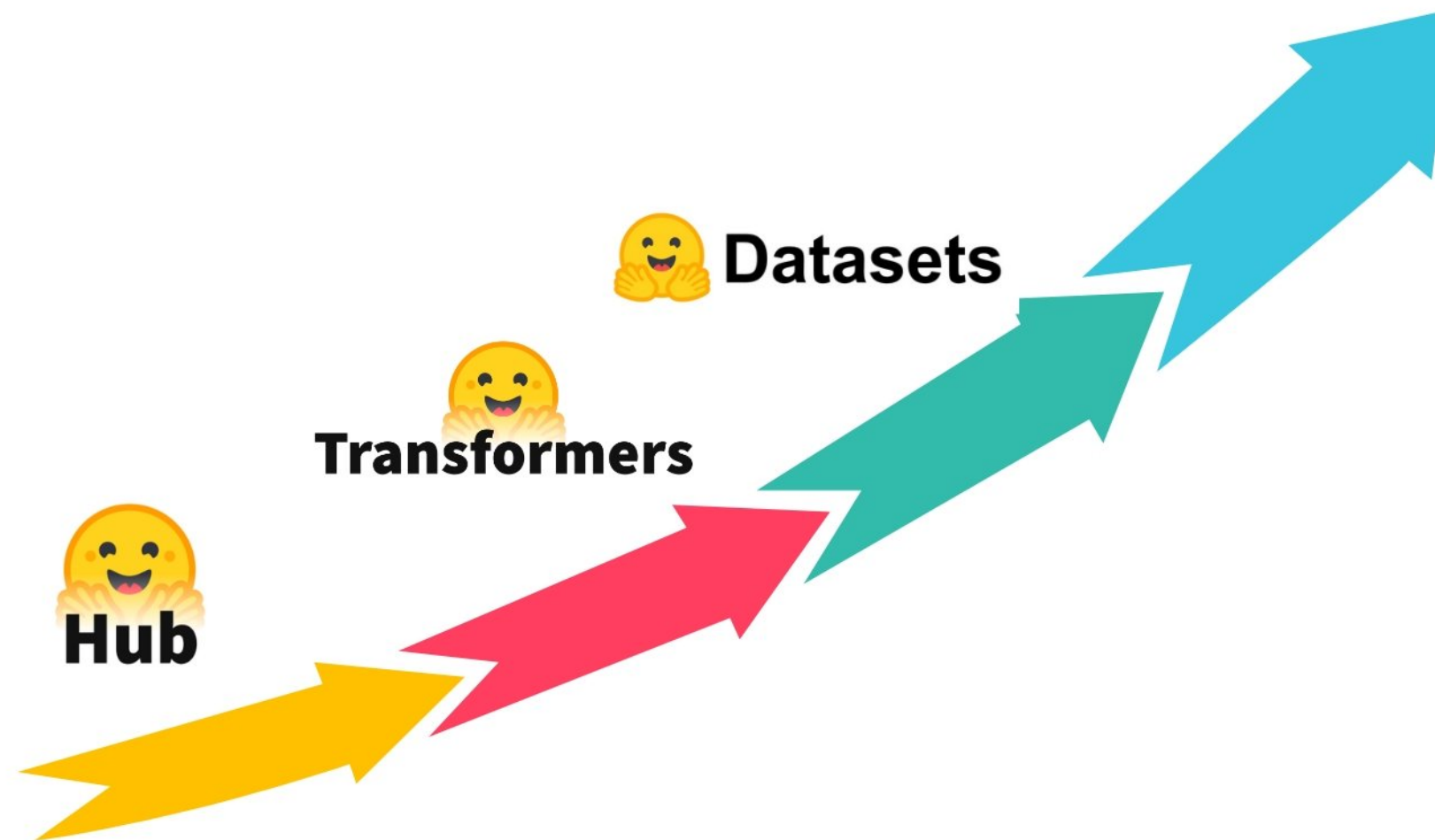
## Applications



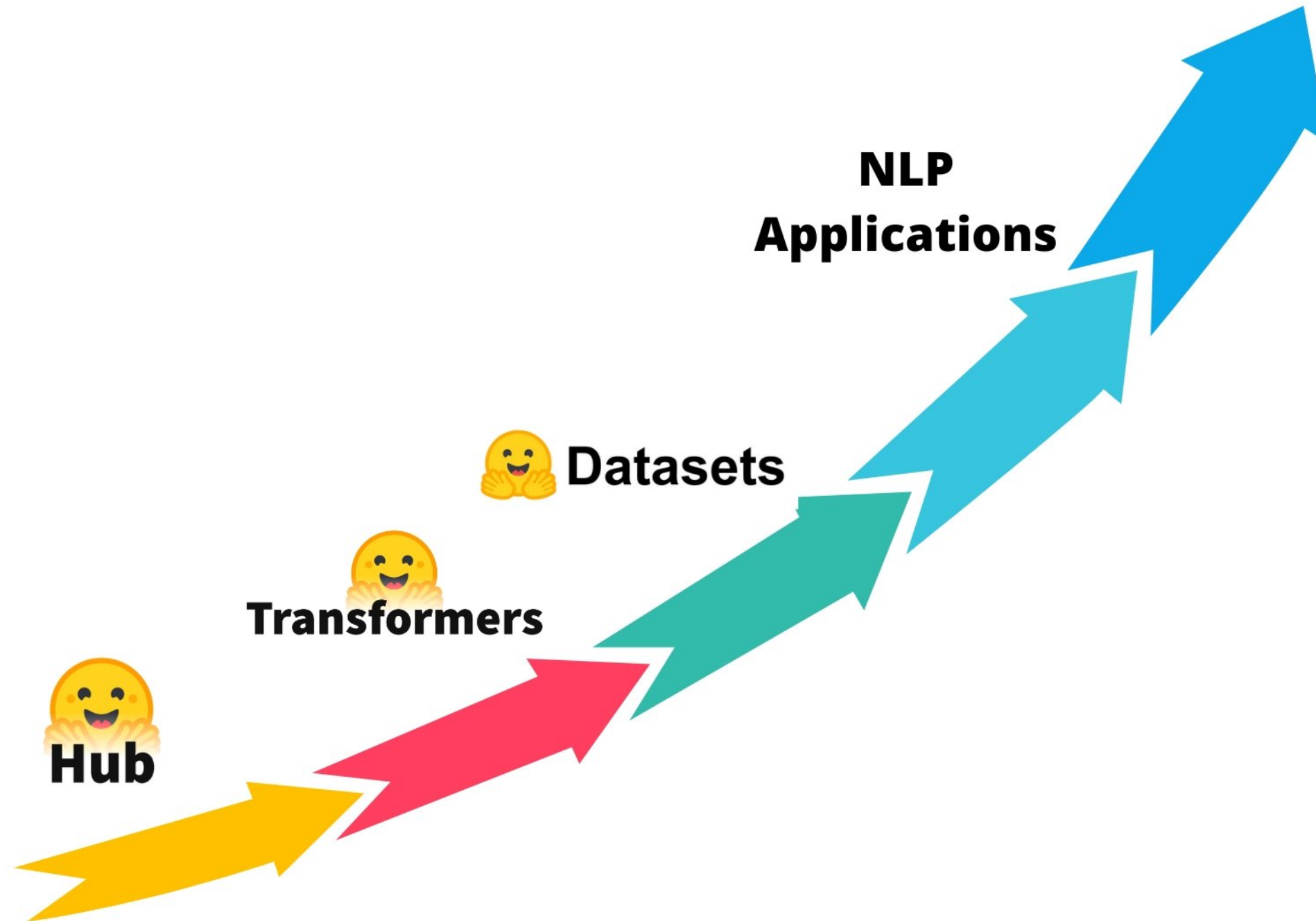
# Coming up...



# Coming up...



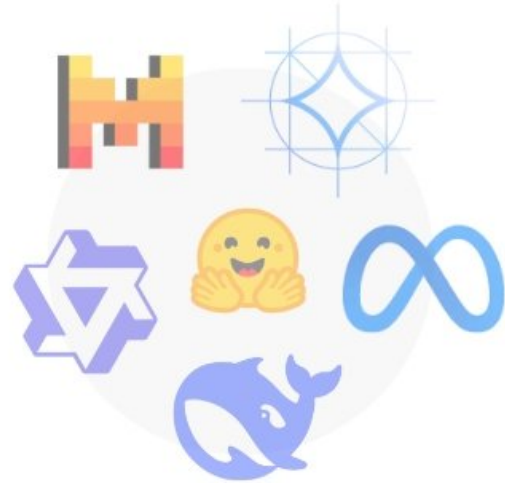
# Coming up...



# The journey beyond!



**Working with  
Hugging Face**



**Introduction to LLMs  
in Python**



**Multi-Modal Models  
with Hugging Face**



**Efficient Model Training  
with PyTorch**





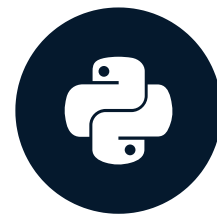




**Let's practice!**  
WORKING WITH HUGGING FACE

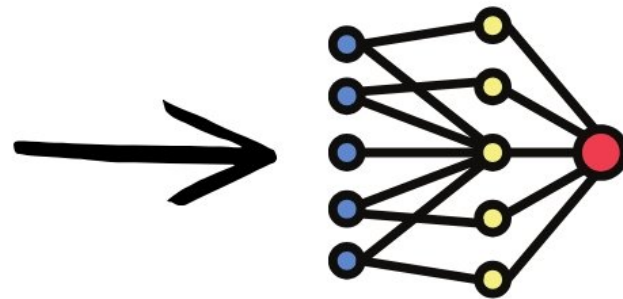
# Running Hugging Face models

WORKING WITH HUGGING FACE



**Jacob H. Marquez**  
Lead Data Engineer

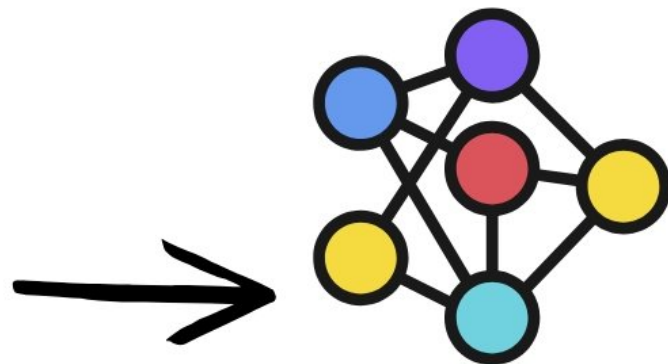
# Inference with Hugging Face



**Dog Image  
Classifier**

**Pomeranian**

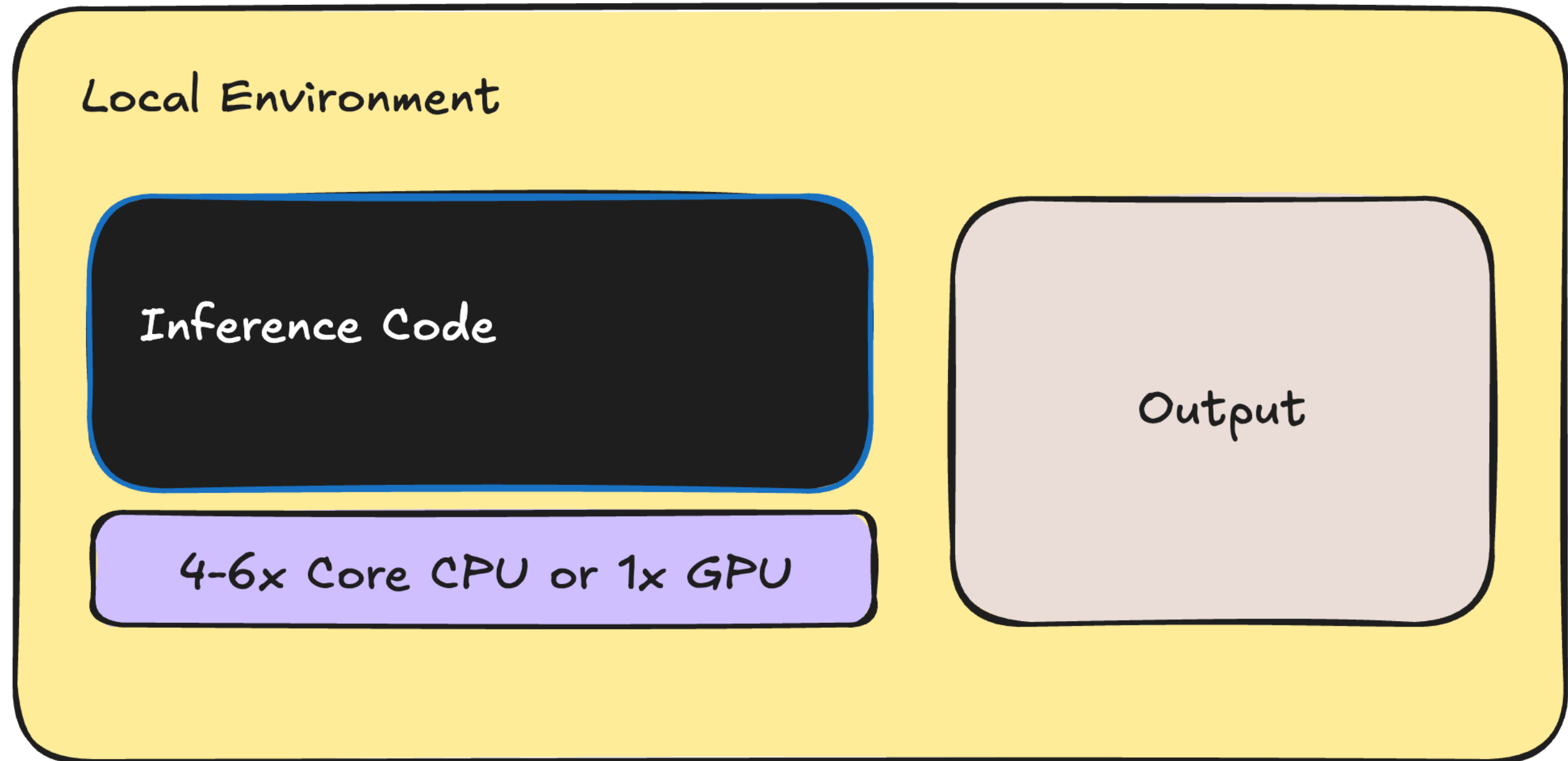
**Complete the story:  
Jack went to school  
one day...**



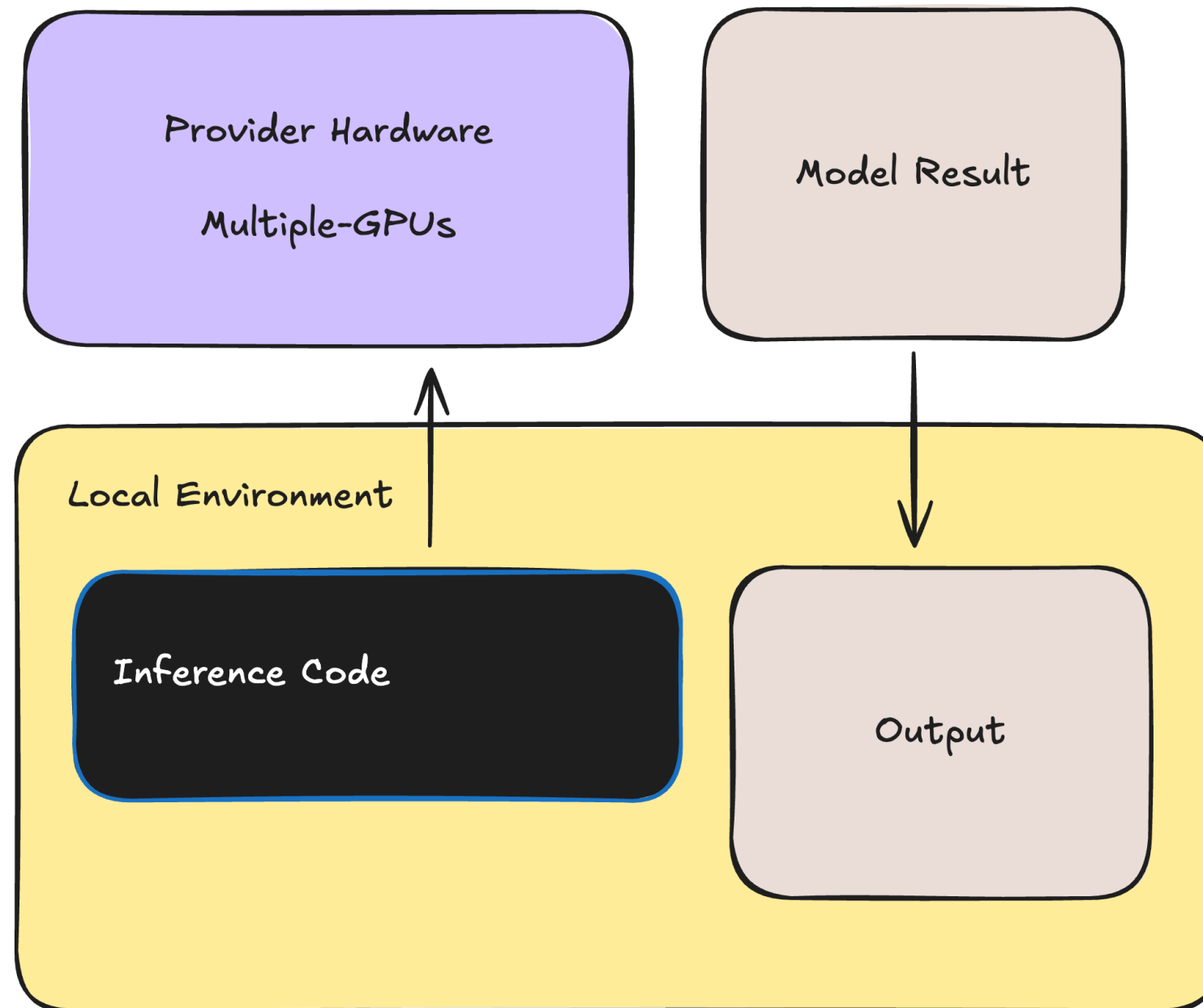
**Text Generator**

**...with a flutter in his  
stomach and a  
notebook full of  
dreams.**

# Local inference



# Inference providers



<sup>1</sup> <https://huggingface.co/docs/inference-providers/en/index>



# Inference with Hugging Face

## Local Inference

- ☐ Free
- 😊 Convenient
- 😴<sup>zzz</sup> Slow and resource-intensive

## Inference Providers

- ☐ Fast
- ☐ Free to get started

<sup>1</sup> <https://huggingface.co/docs/inference-providers/en/index>

# Introduction to the Transformers Library

- Simplifies working with **pre-trained models**



<sup>1</sup> <https://github.com/huggingface/transformers>

# The pipeline

```
from transformers import pipeline

gpt2_pipeline = pipeline(task="text-generation", model="openai-community/gpt2")

print(gpt2_pipeline("What if AI"))
```

```
[{'generated_text': 'What if AI wouldn\'t be used?'\n\nI had to agree with your theory. If a machine\'s learning algorithm is a perfect match for all of a human\'s needs, then you may not have a problem with it. My problem was whether'}]
```

<sup>1</sup> Model Card: <https://huggingface.co/openai-community/gpt2>

# Adjusting Pipeline Parameters

```
from transformers import pipeline

gpt2_pipeline = pipeline(task="text-generation", model="openai-community/gpt2")

results = gpt2_pipeline("What if AI", max_new_tokens=10, num_return_sequences=2)

for result in results:
    print(result['generated_text'])
```

What if AI had never existed?

What if AI could be really smarter than us?

# Using inference providers

```
import os
from huggingface_hub import InferenceClient

client = InferenceClient(
    provider="together",
    api_key=os.environ["HF_TOKEN"],
)
```

<sup>1</sup> <https://huggingface.co/docs/inference-providers/en/index>

```
completion = client.chat.completions.create(  
    model="deepseek-ai/DeepSeek-V3",  
    messages=[  
        {  
            "role": "user",  
            "content": "What is the capital of France?"  
        }  
    ],  
)
```

<sup>1</sup> <https://huggingface.co/docs/inference-providers/en/index>



```
print(completion.choices[0].message)
```

The capital of France is **Paris**. It is known for its iconic landmarks such as the Eiffel Tower, the Louvre Museum, and Notre-Dame Cathedral.

Would you like any additional information about Paris or France?

# Let's practice!

WORKING WITH HUGGING FACE

# Hugging Face Datasets

WORKING WITH HUGGING FACE



**Jacob H. Marquez**  
Lead Data Engineer

[Main](#)
[Tasks](#)
[Libraries](#)
[Languages](#)
[Licenses](#)
[Other](#)

Modalities

[3D](#)
[Audio](#)
[Document](#)
[Geospatial](#)
[Image](#)
[Tabular](#)
[Text](#)
[Time-series](#)
[Video](#)

Size (rows)

[< 1K](#)
[> 1T](#)

Format

[json](#)
[csv](#)
[parquet](#)
[imagefolder](#)
[soundfolder](#)
[webdataset](#)
[text](#)
[arrow](#)

Datasets 443,200

[fka/awesome-chatgpt-prompts](#)

[Viewer](#) • Updated Jan 6 • [203](#) • [22.9k](#) • [8.19k](#)

[HuggingFaceFW/fineweb-2](#)

[Viewer](#) • Updated 11 days ago • [5.02B](#) • [38.3k](#) • [574](#)

[facebook/seamless-interaction](#)

Updated 10 days ago • [1](#) • [89](#)

[marcelbinz/Psych-101](#)

[Viewer](#) • Updated Nov 2, 2024 • [60.1k](#) • [218](#) • [75](#)

[black-forest-labs/kontext-bench](#)

[Viewer](#) • Updated 11 days ago • [1.03k](#) • [38](#)

[FreedomIntelligence/ShareGPT-4o-Image](#)

[Viewer](#) • Updated 6 days ago • [92.3k](#) • [75](#) • [70](#)

[sequelbox/Celestia3-DeepSeek-R1-0528](#)

[Viewer](#) • Updated 5 days ago • [91k](#) • [19](#)

[institutional/institutional-books-1.0](#)

[Viewer](#) • Updated 21 days ago • [983k](#) • [38.2k](#) • [220](#)

[nvidia/OpenScience](#)

[Viewer](#) • Updated 19 days ago • [4.48M](#) • [435](#) • [54](#)

[HuggingFaceFW/fineweb](#)

[Viewer](#) • Updated Jan 31 • [25B](#) • [210k](#) • [2.23k](#)

[cais/hle](#)

[Viewer](#) • Updated May 20 • [2.5k](#) • [6.54k](#) • [382](#)

[racineai/OGC\\_MEGA\\_MultiDomain\\_DocRetrieval](#)

[Viewer](#) • Updated 5 days ago • [1.09M](#) • [13](#)

<sup>1</sup> <https://huggingface.co/datasets>

[Main](#)
[Tasks](#)
[Libraries](#)
[Languages](#)
[Licenses](#)
[Other](#)

Modalities

[3D](#)
[Audio](#)
[Document](#)
[Geospatial](#)
[Image](#)
[Tabular](#)
[Text](#)
[Time-series](#)
[Video](#)

Size (rows)

[< 1K](#)
[> 1T](#)

Format

[json](#)
[csv](#)
[parquet](#)
[imagefolder](#)
[soundfolder](#)
[webdataset](#)
[text](#)
[arrow](#)

Datasets 443,200

[Sort: Trending](#)

[fka/awesome-chatgpt-prompts](#)

[Viewer](#) • Updated Jan 6 • [203](#) • [22.9k](#) • [8.19k](#)

[facebook/seamless-interaction](#)

Updated 10 days ago • [1](#) • [89](#)

[black-forest-labs/kontext-bench](#)

[Viewer](#) • Updated 11 days ago • [1.03k](#) • [38](#)

[sequelbox/Celestia3-DeepSeek-R1-0528](#)

[Viewer](#) • Updated 5 days ago • [91k](#) • [19](#)

[nvidia/OpenScience](#)

[Viewer](#) • Updated 19 days ago • [4.48M](#) • [435](#) • [54](#)

[cais/hle](#)

[Viewer](#) • Updated May 20 • [2.5k](#) • [6.54k](#) • [382](#)

[HuggingFaceFW/fineweb-2](#)

[Viewer](#) • Updated 11 days ago • [5.02B](#) • [38.3k](#) • [574](#)

[marcelbinz/Psych-101](#)

[Viewer](#) • Updated Nov 2, 2024 • [60.1k](#) • [218](#) • [75](#)

[FreedomIntelligence/ShareGPT-4o-Image](#)

[Viewer](#) • Updated 6 days ago • [92.3k](#) • [75](#) • [70](#)

[institutional/institutional-books-1.0](#)

[Viewer](#) • Updated 21 days ago • [983k](#) • [38.2k](#) • [220](#)

[HuggingFaceFW/fineweb](#)

[Viewer](#) • Updated Jan 31 • [25B](#) • [210k](#) • [2.23k](#)

[racineai/OGC\\_MEGA\\_MultiDomain\\_DocRetrieval](#)

[Viewer](#) • Updated 5 days ago • [1.09M](#) • [13](#)

<sup>1</sup> <https://huggingface.co/datasets>









# Installing Datasets Package

```
pip install datasets
```

- ☐ Access
- ☐ Download
- ☐ Use
- ☐ Share



📁 [wikimedia/wikipedia](#)

🔍 Viewer • Updated Jan 9, 2024 • 📄 61.6M • ⬇️ 73.6k • ❤️ 860

📁 [R2E-Gym/R2E-Gym-Subset](#)

🔍 Viewer • Updated Apr 11 • 📄 4.58k • ⬇️ 150 • ❤️ 11

📁 [yale-nlp/SciArena](#)

🔍 Viewer • Updated 6 days ago • 📄 13.2k • ❤️ 11

📁 [XenArcAI/MathX-5M](#)

🔍 Viewer • Updated 1 day ago • 📄 4.32M • ⬇️ 655 • ❤️ 15

<sup>1</sup> <https://huggingface.co/docs/datasets/loading>

# Downloading a dataset

```
from datasets import load_dataset  
  
data = load_dataset("IVN-RIN/BioBERT_Italian")
```

## Split parameter

```
data = load_dataset("IVN-RIN/BioBERT_Italian", split="train")
```

<sup>1</sup> <https://huggingface.co/docs/datasets/v2.15.0/loading>

# Apache Arrow dataset formats

Row-based			Column-based		
Row 1		1331246660	session_id		1331246660
		3/8/2012 2:44PM			1331246351
		99.155.155.225			1331244570
Row 2		1331246351			1331261196
		3/8/2012 2:38PM	timestamp		3/8/2012 2:44PM
		65.87.165.114			3/8/2012 2:38PM
Row 3		1331244570			3/8/2012 2:09PM
		3/8/2012 2:09PM			3/8/2012 6:46PM
		71.10.106.181	source_ip		99.155.155.225
Row 4		1331261196			65.87.165.114
		3/8/2012 6:46PM			71.10.106.181
		76.102.156.138			76.102.156.138

<sup>1</sup> <https://arrow.apache.org/overview/>

# Data manipulation

```
data = load_dataset("IVN-RIN/BioBERT_Italian", split="train")

# Filter for pattern " bella "
filtered = data.filter(lambda row: " bella " in row['text'])
print(filtered)
```

```
Dataset({
  features: ['text'],
  num_rows: 1122
})
```

<sup>1</sup> <https://huggingface.co/docs/datasets/process#select-and-filter>

# Data manipulation

```
# Select the first two rows
sliced = filtered.select(range(2))

print(sliced)
```

```
Dataset({features: ['text'], num_rows: 2})
```

```
# Extract the 'text' for the first row
print(sliced[0]['text'])
```

```
Concentrazioni atmosferiche di PCDD/PCDF...
```

<sup>1</sup> <https://huggingface.co/docs/datasets/process#select-and-filter>

**Let's practice!**  
WORKING WITH HUGGING FACE