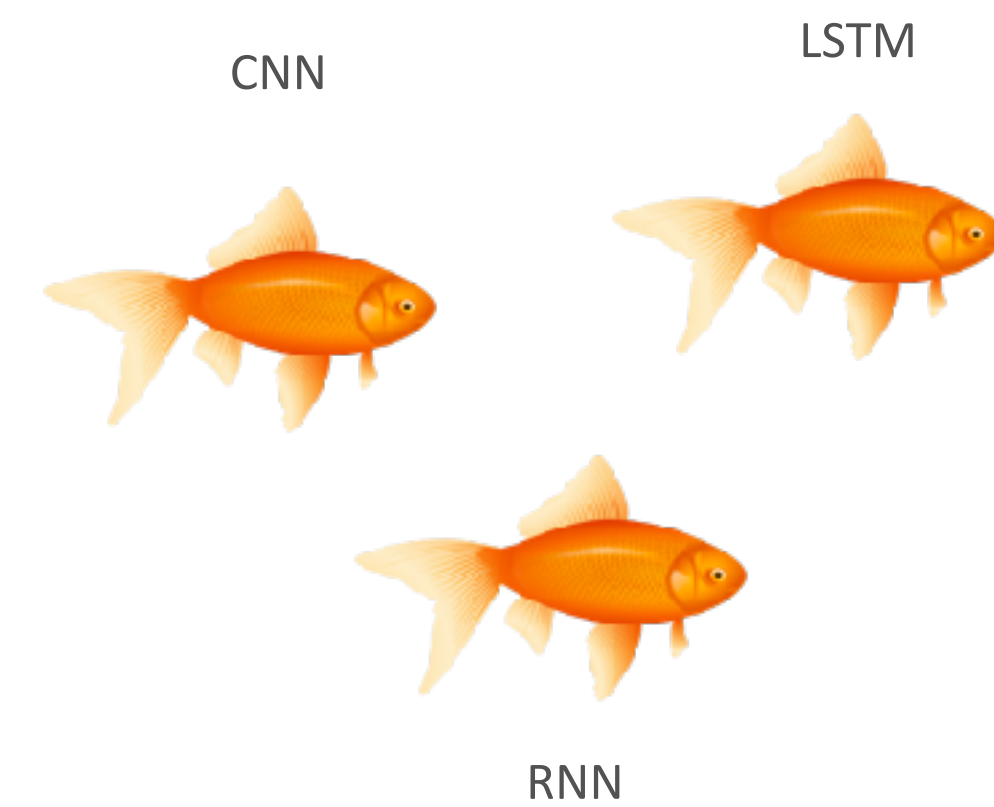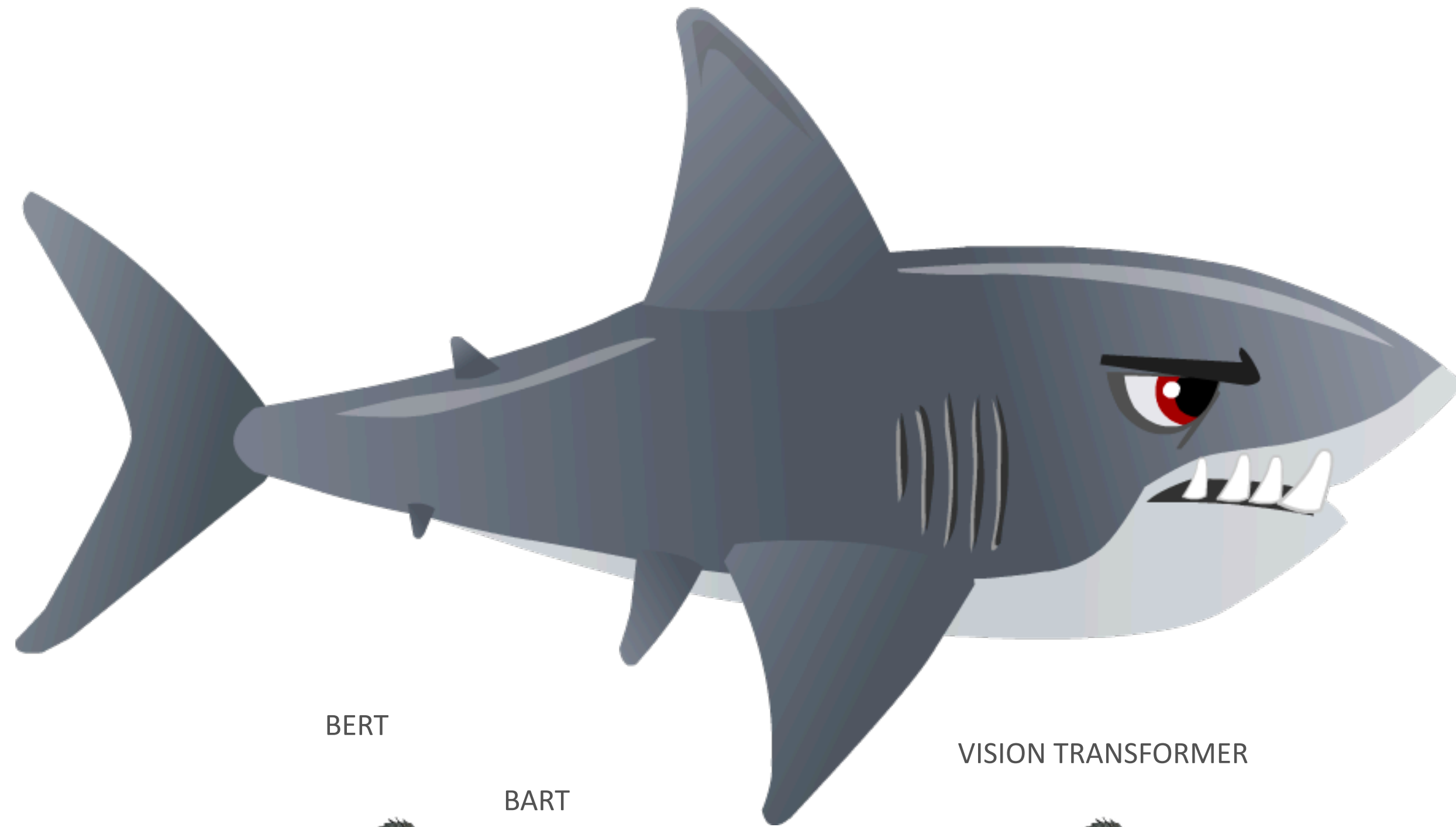# Demystifying the technology behind Generative AI

Julien Simon, Chief Evangelist, Hugging Face
julsimon@huggingface.co

# 2022: Transformers are eating Deep Learning



BERT

BART

GPT-2

GPT-3

BLOOM
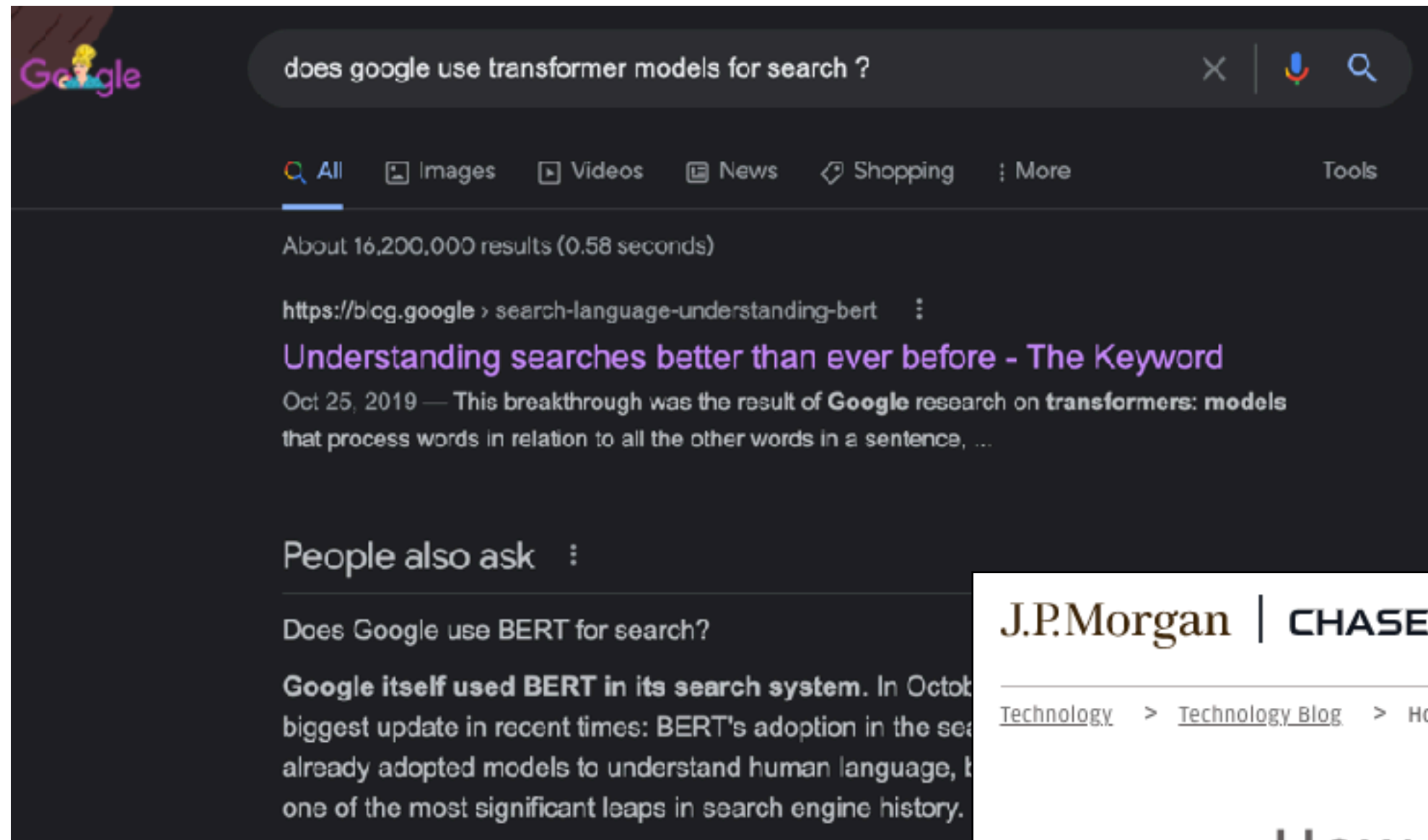
CLIP

WAV2VEC2

STABLE DIFFUSION

VISION TRANSFORMER

WHISPER

CNN

LSTM

RNN

*"Transformers are emerging as
a general-purpose architecture for ML"*
https://www.stateof.ai (2021)

RNN and CNN usage down, Transformers usage up!
https://www.kaggle.com/kaggle-survey-2021

# Transformer models in the wild

# Hugging Face: the largest collection of open source models

https://huggingface.co



220K pre-trained models
(NLP, CV, Speech, etc.)

40K datasets
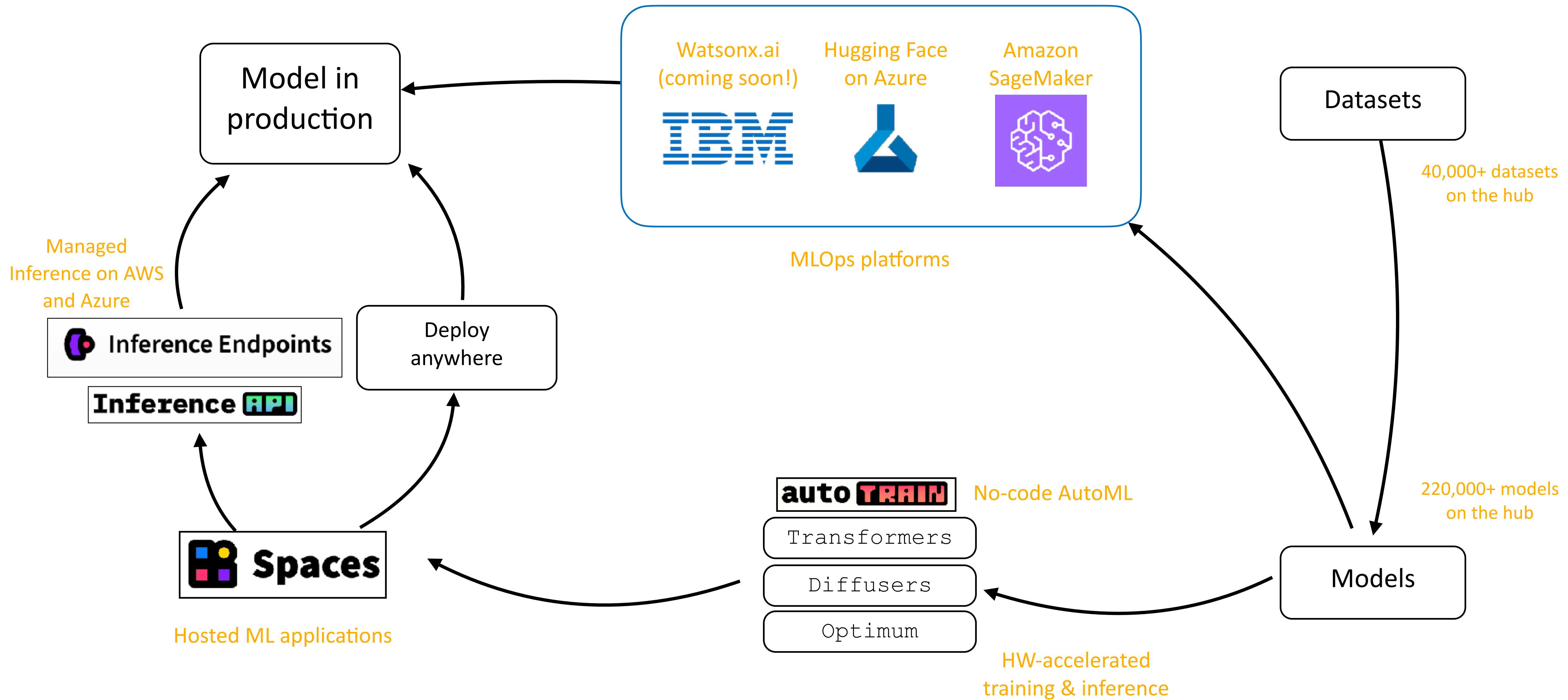
25+ ML libraries: Keras, spaCY,
Scikit-Learn, fastai, etc.

10K organizations

500K+ users daily

# Hugging Face at a glance

Model in production

Watsonx.ai (coming soon!)

Hugging Face on Azure

Amazon SageMaker

MLOps platforms

Datasets

40,000+ datasets on the hub

Managed Inference on AWS and Azure

Inference Endpoints

Inference API

Deploy anywhere

Spaces

Hosted ML applications

autoTRAIN

No-code AutoML

220,000+ models on the hub

Transformers

Diffusers

Optimum

Models

HW-accelerated training & inference
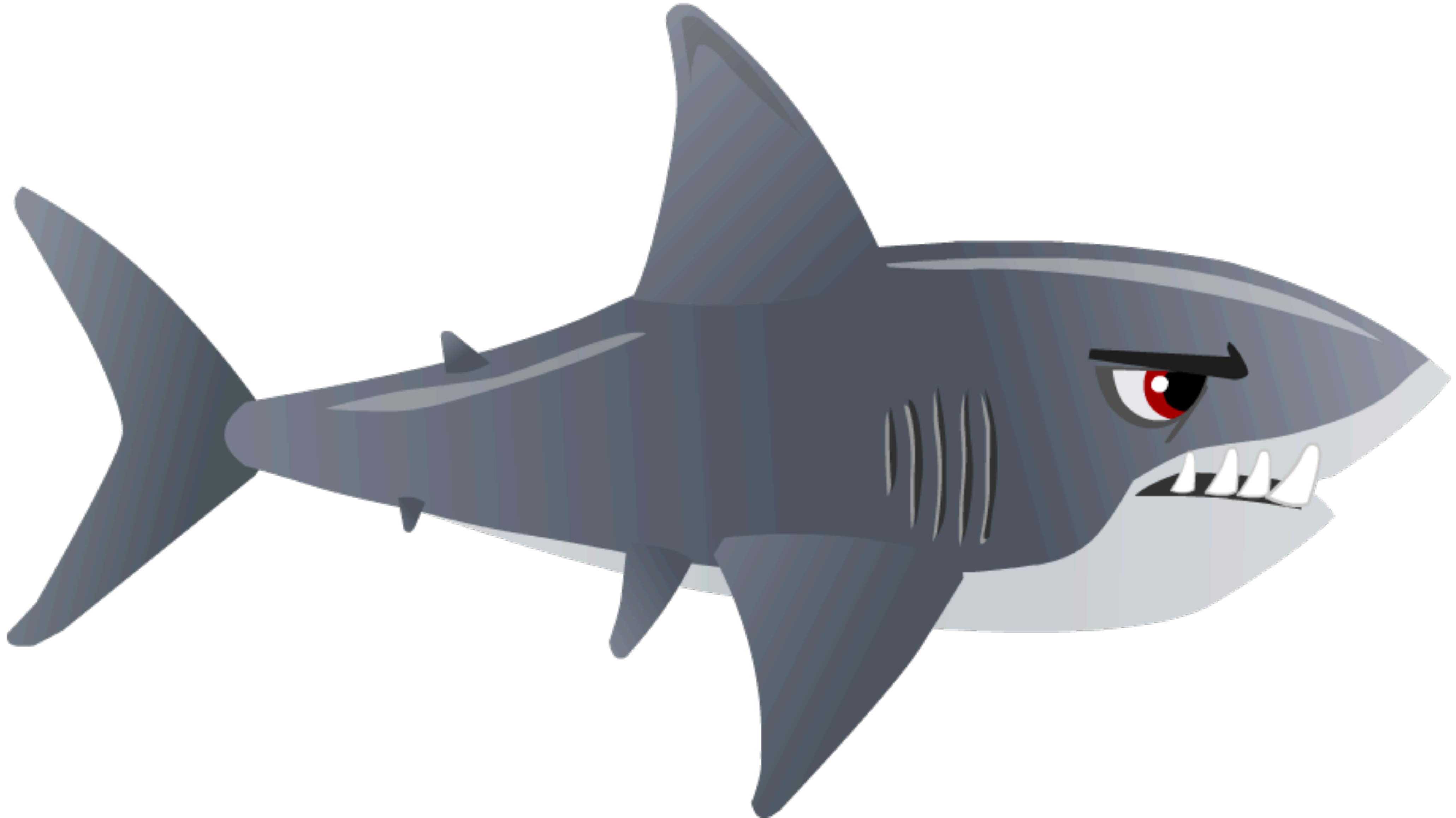
# Expert Acceleration Program

- Building quality datasets for Q&A, instruction following, etc.
  - With your own data
  - With prompts/answers from other models
- Deduplicating data, a key performance factor for LLMs
- Selecting the right model architecture
- Fine-tuning models efficiently (PEFT, LoRA)
- Selecting hyper parameters
- Evaluating model quality
- Optimizing model latency in production

# 2023: Generative AI is hungry!

# Generative models in the wild



JU is ChatGPT powered by transformer models?

Yes, ChatGPT is powered by transformer models. Specifically, it is based on architecture of the GPT (Generative Pre-trained Transformer) family of mod were developed by OpenAI. The GPT models are designed to generate natu text, and they use a transformer architecture that allows them to model long-term dependencies in language.

ChatGPT is a variant of the GPT models that has been fine-tuned on a large amount of conversational data in order to generate human-like responses to user input. It uses a combination of neural network architecture, training data, and machine learning algorithms to generate responses that are contextually appropriate and semantically relevant.

**Introducing BloombergGPT, Bloomberg's 50-billion parameter large language model, purpose-built from scratch for finance**

March 30, 2023

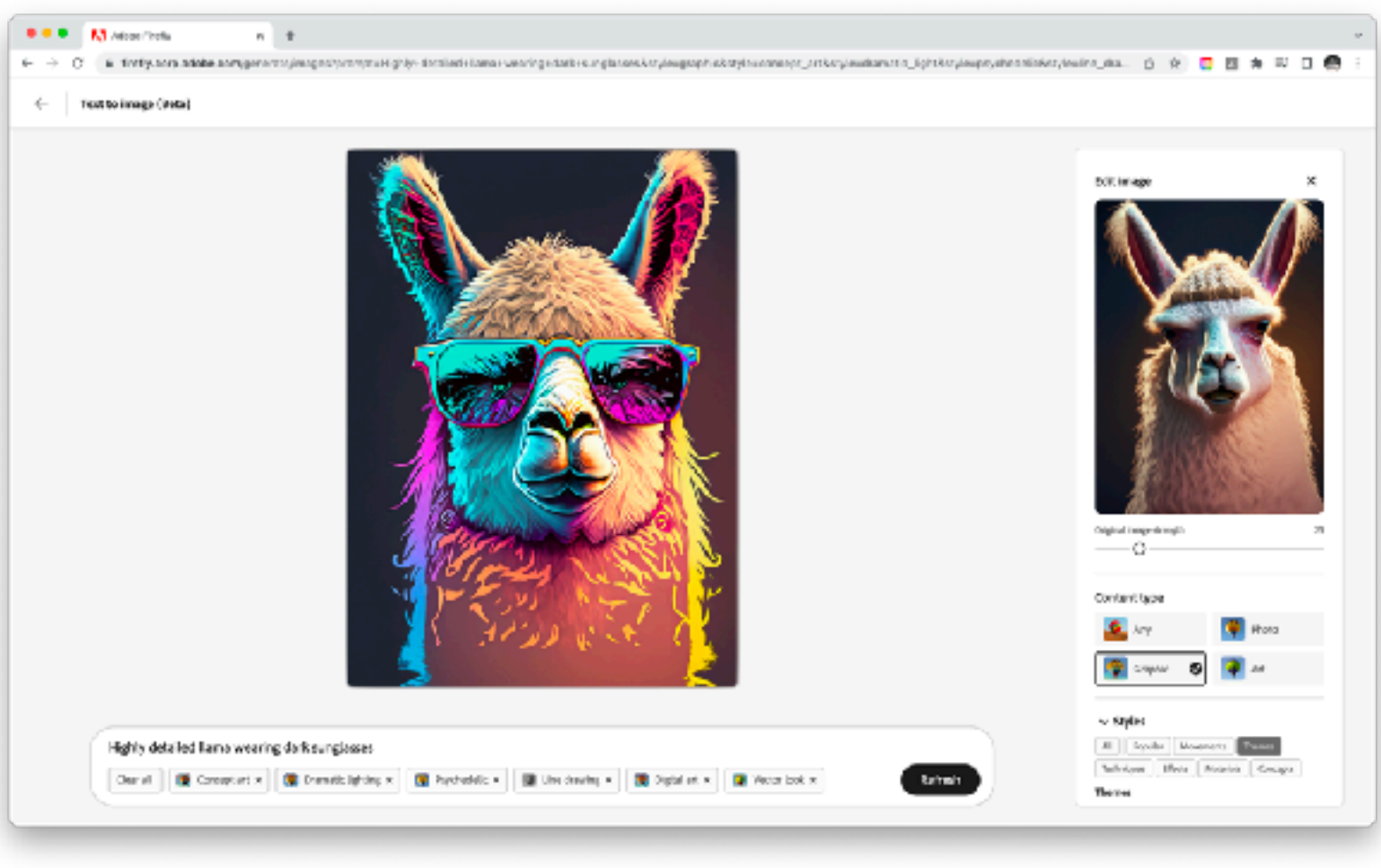**Meet Einstein GPT, the World's First Generative AI for CRM**

LEARN MORE →

**GitHub Copilot**

**Amazon CodeWhisperer**

Highly detailed llama wearing da fksunglasses

Prompt:
A sadhu man in Rishikesh, India meditating near the Ganges river

# BLOOM: open-source alternative to GPT-3



a BigScience initiative
BL**OO**M
176B params · 59 languages · Open-access

https://bigscience.huggingface.co

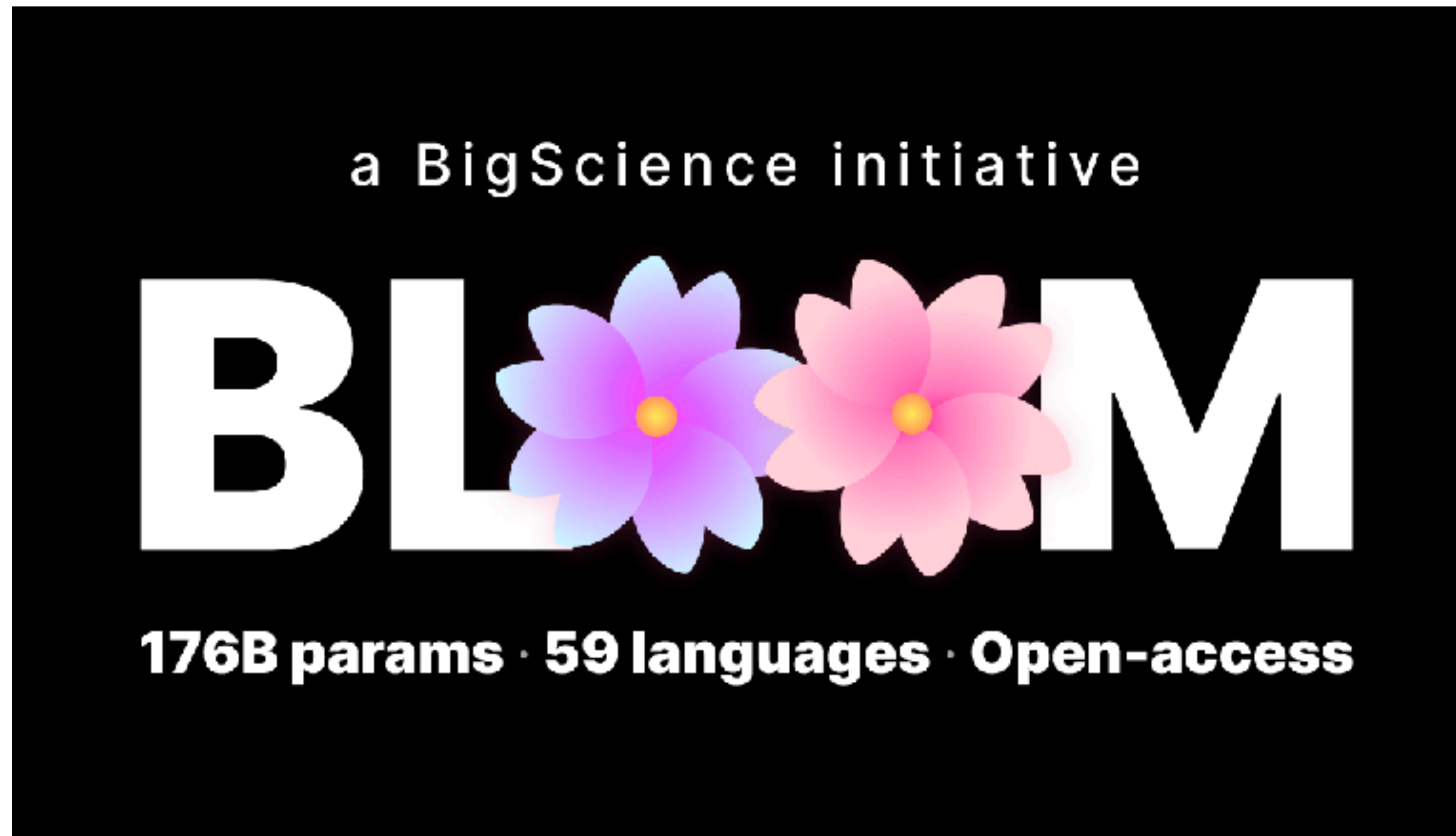https://huggingface.co/bigscience/bloom

1.5TB of text, 350B tokens

43 languages, 16 programming languages

118 days of training on 384 A100 GPUs
(public cluster)

**Smaller versions are available** : 560M, 1.1B, 1.7B, 3B, 7.1B

BLOOMZ models (same sizes) are fine-tuned for **instruction following**
https://huggingface.co/bigscience/bloomz

# BigCode: open-source LLMs for code generation

https://www.bigcode-project.org



Dataset: https://huggingface.co/
datasets/bigcode/the-stack

2.9TB of deduplicated code



Model: https://huggingface.co/bigcode/starcoder
https://arxiv.org/abs/2305.06161

15.5B parameters, 1T tokens, 80+ languages

8K context

26 days of training on 512 A100 GPUs (AWS)

# Financial LLM case study: BloombergGPT

https://arxiv.org/abs/2303.17564

- Bloomberg is a long-time customer of Hugging Face

- They built a 700B token dataset (general-purpose and financial)

- They evaluated different models and picked BLOOM as a starting point

- Based on dataset size and compute budget, they rescaled BLOOM to an optimal 50B parameters

- They used our Expert Acceleration Program (EAP) to get deep, first-hand expertise on customizing BLOOM for their own purposes: model architecture, hyper parameter selection, etc.

- They trained the model on AWS (64 instances, 512 A100 GPUs) for 52 days.

# Open Large Language Model leaderboard

| Model | Revision ▲ | Average ⬆ ▲ | ARC (25-shot) ⬆ ▲ | HellaSwag (10-shot) ⬆ ▲ | MMLU (5-shot) ⬆ ▲ | TruthfulQA (0-shot) ⬆ ▲ |
|---|---|---|---|---|---|---|
| tiiuae/falcon-40b-instruct | main | 63.2 | 61.6 | 84.4 | 54.1 | 52.5 |
| tiiuae/falcon-40b | main | 60.4 | 61.9 | 85.3 | 52.7 | 41.7 |
| ausboss/llama-30b-supercot | main | 59.8 | 58.5 | 82.9 | 44.3 | 53.6 |
| llama-65b | main | 58.3 | 57.8 | 84.2 | 48.8 | 42.3 |
| MetaIX/GPT4-X-Alpasta-30b | main | 57.9 | 56.7 | 81.4 | 43.6 | 49.7 |
| Aeala/VicUnlocked-alpaca-30b | main | 57.6 | 55 | 80.8 | 44 | 50.4 |
| digitous/Alpacino30b | main | 57.4 | 57.1 | 82.6 | 46.1 | 43.8 |
| Aeala/GPT4-x-AlpacaDente2-30b | main | 57.2 | 56.1 | 79.8 | 44 | 49.1 |
| TheBloke/dromedary-65b-lora-HF | main | 57 | 57.8 | 80.8 | 50.8 | 38.8 |
| TheBloke/Wizard-Vicuna-13B-Uncensored-HF | main | 57 | 53.6 | 79.6 | 42.7 | 52 |
| elinas/llama-30b-hf-transformers-4.29 | main | 56.9 | 57.1 | 82.6 | 45.7 | 42.3 |
| llama-30b | main | 56.9 | 57.1 | 82.6 | 45.7 | 42.3 |
| cyl/awsome-llama | main | 56.8 | 54.4 | 79.7 | 41.8 | 51.3 |

4 key benchmarks from the Eleuther AI Language Model Evaluation Harness 🤗

# GenAI use cases

# Example: text-to-image generation

```python
from diffusers import StableDiffusionPipeline

pipe = StableDiffusionPipeline.from_pretrained(
    "CompVis/stable-diffusion-v1-4")

prompt = "in the desert, a corvette parked \
in front of an old-school diner at sundown"

image = pipe(prompt).images[0]
image.save("picture.png")
```



Intel Sapphire Rapids
(Amazon EC2 r7iz)

https://youtu.be/KJDCGyZ2fPw

# Example: image inpainting

# Example: text-to-image generation

# Example: Q&A

https://huggingface.co/spaces/trl-lib/stack-llama
https://huggingface.co/blog/stackllama

# Example: retrieval-augmented generation

https://huggingface.co/spaces/Ekimetrics/climate-question-answering

## ClimateQ&A chatbot

dioxide (CO2), nitrous oxide (N2O), methane (CH4), and ozone (O3) [docs 2, 4, 7, 8].

Anthropogenic GHGs such as carbon dioxide (CO2), methane (CH4), nitrous oxide (N2O), and fluorinated gases (e.g., hydrofluorocarbons, perfluorocarbons, sulphur hexafluoride) are released from various sources [doc 9]. CO2 makes the largest contribution to global GHG emissions [doc 9].

While CO2 is the most important greenhouse gas, marine fluxes of methane and nitrous oxide can also be important, for both coastal regions and the open ocean [doc 1].

Human-made GHGs include sulphur hexafluoride (SF6), hydrofluorocarbons (HFCs), chlorofluorocarbons (CFCs), and perfluorocarbons (PFCs); several of these are also O3-depleting (and are regulated under the Montreal Protocol) [docs 2, 7, 8].

Non-CO2 emissions included in IPCC reports are all anthropogenic emissions other than carbon dioxide (CO2) that result in radiative forcing. These include short-lived climate forcers, such as methane (CH4), some fluorinated gases, ozone (O3) precursors, aerosols or aerosol precursors, such as black carbon and sulphur dioxide, respectively, as well as long-lived greenhouse gases, such as nitrous oxide (N2O) or other fluorinated gases [docs 5, 6].

The chemical composition of the atmosphere (beyond CO2 and water vapor changes) is expected to change in response to a warming climate. These changes in greenhouse gases (methane, nitrous oxide, and ozone) and aerosol amount (including dust) have the potential to alter the TOA energy budget and are collectively referred to as 'non-CO2 biogeochemical feedbacks' [doc 10].

## Sources

Query used for retrieval:
What are the different types of greenhouse gases and how do they contribute to climate change?

📄 doc 1: Chapter 5 - Changing Ocean, Marine Ecosystems, and Dependent Communities. In: IPCC Special Report on the Ocean and Cryosphere in a Changing Climate page 73
release of greenhouse gases from different marine environments, interacting with natural biological, physical and chemical processes. Whilst CO2 is the most important greenhouse gas, marine fluxes of methane and nitrous oxide can also be important, for both coastal regions and the open ocean (Arevalo-Martinez et al. 2015; Borges et al. 2016; Hamdan and Wickland, 2016).

📄 doc 2: Full Report. In: Climate Change 2022: Mitigation of Climate Change. Contribution of the WGIII to the AR6 of the IPCC page 2032
Greenhouse gases (GHGs) Gaseous constituents of the atmosphere, both natural and anthropogenic, that absorb and emit radiation at specific wavelengths within the spectrum of radiation emitted by the Earth's surface, by the atmosphere itself, and by clouds. This property causes the greenhouse effect. Water vapour (H2O), carbon dioxide (CO2 ), nitrous oxide (N2O), methane (CH4 )

Combination of <u>semantic search</u>, <u>text generation</u> and <u>prompt tweaking</u>!
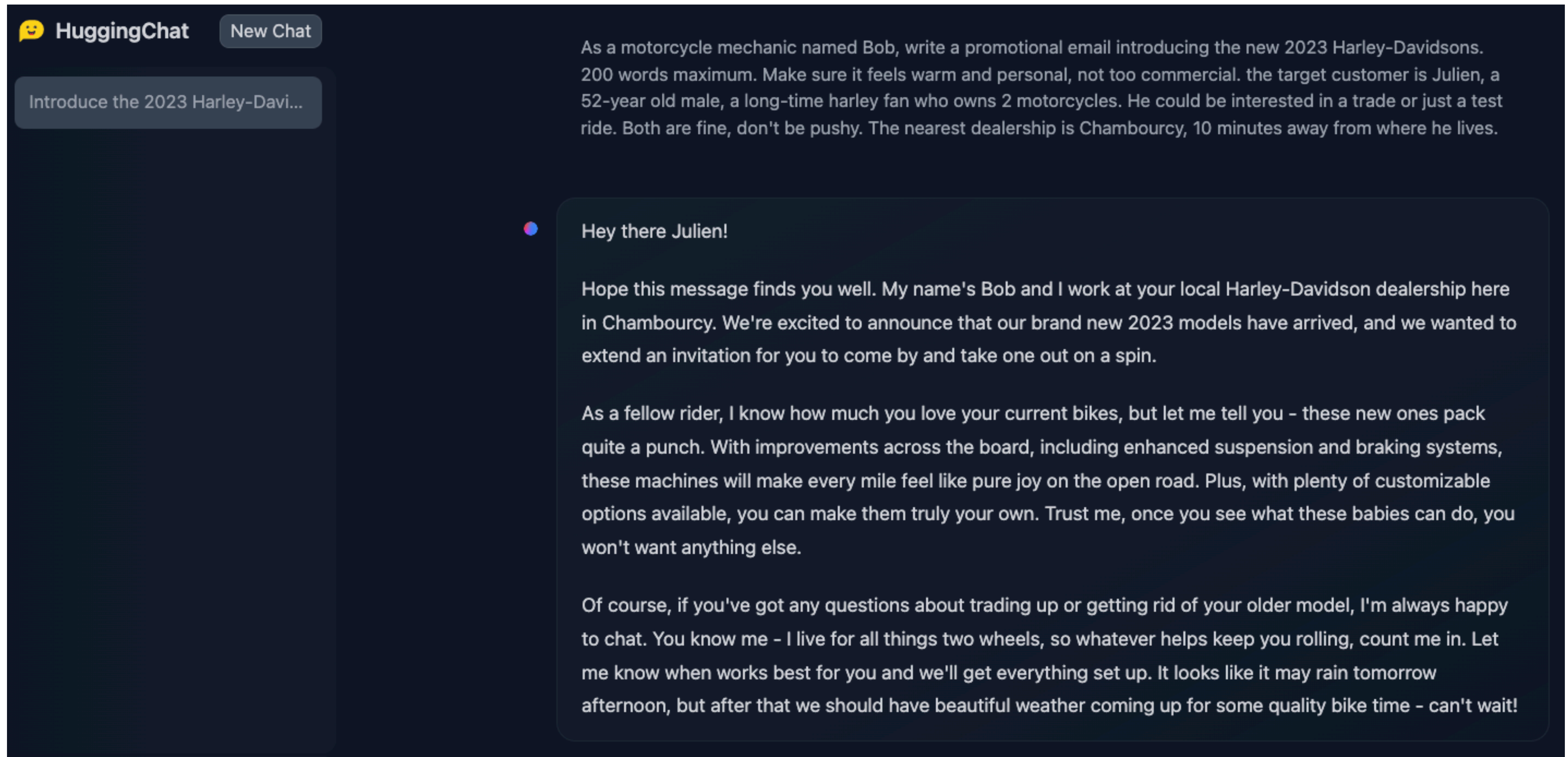
This could be further improved for any domain
with a fine-tuned embedding models and/or a fine-tuned LLM 🤗

# Example: personalized marketing content with HuggingChat

https://huggingface.co/chat

# Example: coding assistance with StarCoder

https://huggingface.co/spaces/HuggingFaceH4/starchat-playground
https://huggingface.co/blog/starchat-alpha

# Hugging Face on AWS

# Hugging Face on AWS

**Hugging Face models, datasets, and libraries**

## Hugging Face Expert Acceleration Program (EAP)

aws marketplace

### Experiment on Hugging Face

Hugging Face Spaces

aws marketplace    Q2'23

### Train and deploy on Amazon EC2

Hugging Face DLAMI

aws marketplace

### Train and deploy on Amazon SageMaker

Hugging Face DLCs SageMaker JumpStart

### Deploy on Hugging Face

Hugging Face Inference Endpoints

aws marketplace    Q2'23

## AWS Infrastructure (CPU, GPU, Trainium, Inferentia)

# Demo: Amazon SageMaker JumpStart

# Demo: Optimum Neuron with AWS Trainium and Inferentia2

Vision Transformer on food101 dataset (75K training images): 1 minute/epoch
DistilBERT on 32-token sequences: 1ms P99 latency

# Conclusion

# Getting started

https://huggingface.co/tasks

https://huggingface.co/course

https://github.com/huggingface

https://huggingface.co/blog

https://huggingface.co/docs/sagemaker/index

https://www.philschmid.de/

https://youtube.com/c/juliensimonfr

# Summing things up

- AI is changing the way we build software

- Transformer models are the de facto standard for AI-powered apps.

- Don't believe the hype: today's "best" model will be superseded in weeks

- No model rules them all : find the most appropriate one for each use case

- "Small" fine-tuned open-source models are the way to go

- AWS is the best place to train and deploy transformers!