

In-browser AI apps with Gradio and Transformers

FEDAY 2024

Yuichiro Tachibana, ML Developer Advocate @ 

In-browser AI apps with
Gradio and
Transformers

ChatGPT 4o

↑ Share T

Generate an image of a futuristic city at sunset, with flying cars zooming between towering skyscrapers made of glass and steel, and vibrant neon signs in multiple languages glowing against the skyline.



The generated image of a futuristic city at sunset, featuring flying cars, towering glass and steel skyscrapers, and vibrant neon signs. Let me know if you'd like any adjustments or further options!

ChatGPT can make mistakes. Check important info. ?

AI apps

Custom Upload Audio v3.5

SUNO

Song description Instrumental

A futuristic techno song featuring AI app dev tools like Transformers

Home Create Library Explore Search

69 / 200

>Create

Need ideas?

Invite Friends

40 Credits Subscribe

What's New? 9

Notifications

More from Suno

Digital Symphony ExasperatingCymbals... 00:13 /

Digital futuristic Public

Digital futuristic Public

Filter Tasks by name

Multimodal

- 🔊 [Audio-Text-to-Text](#)
- 🔊 [Image-Text-to-Text](#)
- 🔍 [Visual Question Answering](#)
- 📄 [Document Question Answering](#)
- 🎥 [Video-Text-to-Text](#)
- 🔀 [Any-to-Any](#)

Computer Vision

- 📦 [Depth Estimation](#)
- 🖼️ [Image Classification](#)
- 📦 [Object Detection](#)
- 🖼️ [Image Segmentation](#)
- ➡️ [Text-to-Image](#)
- ➡️ [Image-to-Text](#)
- 🖼️ [Image-to-Image](#)
- ➡️ [Image-to-Video](#)
- 🖼️ [Unconditional Image Generation](#)
- 🎥 [Video Classification](#)
- ➡️ [Text-to-Video](#)
- 📦 [Zero-Shot Image Classification](#)
- 🎭 [Mask Generation](#)
- 📦 [Zero-Shot Object Detection](#)
- ➡️ [Text-to-3D](#)
- 🖼️ [Image-to-3D](#)
- 🖼️ [Image Feature Extraction](#)
- ⭐ [Keypoint Detection](#)

Natural Language Processing

- 📊 [Text Classification](#)
- 📊 [Token Classification](#)
- 📋 [Table Question Answering](#)
- 💡 [Question Answering](#)

Qwen/QwQ-32B-Preview

Text Generation • Updated 5 days ago • ↓ 33.6k • ⚡ • ❤ 1k

tencent/HunyuanVideo

Updated about 7 hours ago • ❤ 339

Lightricks/LTX-Video

Image-to-Video • Updated 11 days ago • ↓ 32.3k • ❤ 574

black-forest-labs/FLUX.1-dev

Text-to-Image • Updated Aug 16 • ↓ 1.38M • ⚡ • ❤ 6.9k

Djrange/aven-1-Flux

Text-to-Image • Updated 7 days ago • ↓ 358

AIDC-AI/Marco-01

Text Generation • Updated 11 days ago • ↓ 9.71k • ❤ 605

ginipick/flux-lora-eric-cat

Text-to-Image • Updated 2 days ago • ↓ 4.18k • ⚡ • ❤ 196

HuggingFaceTB/SmolVLM-Instruct

Image-Text-to-Text • Updated 1 day ago • ↓ 21.7k • ❤ 226

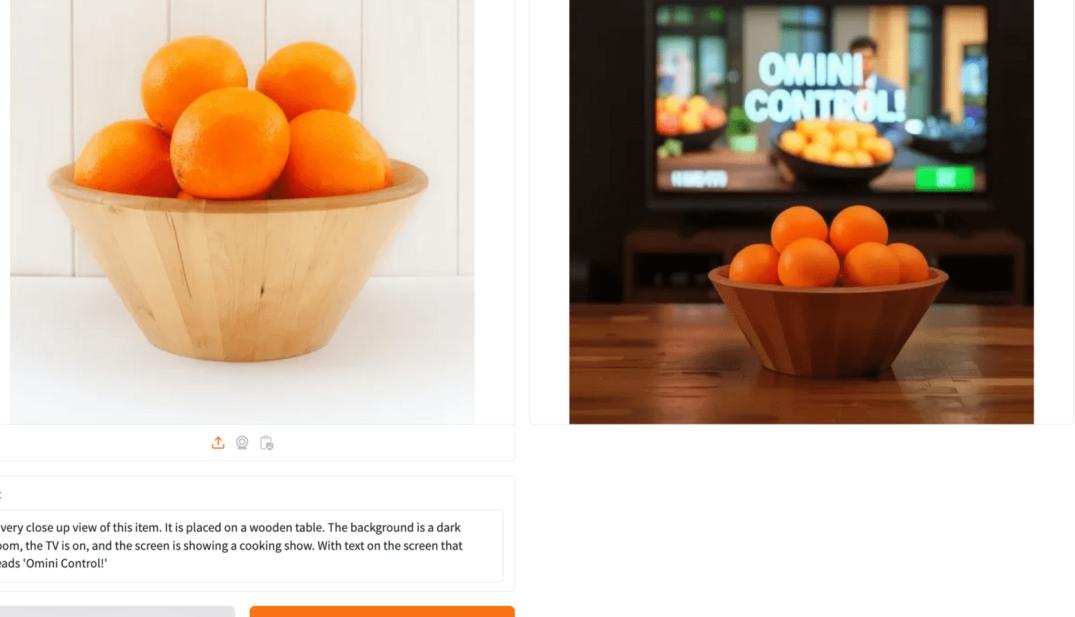
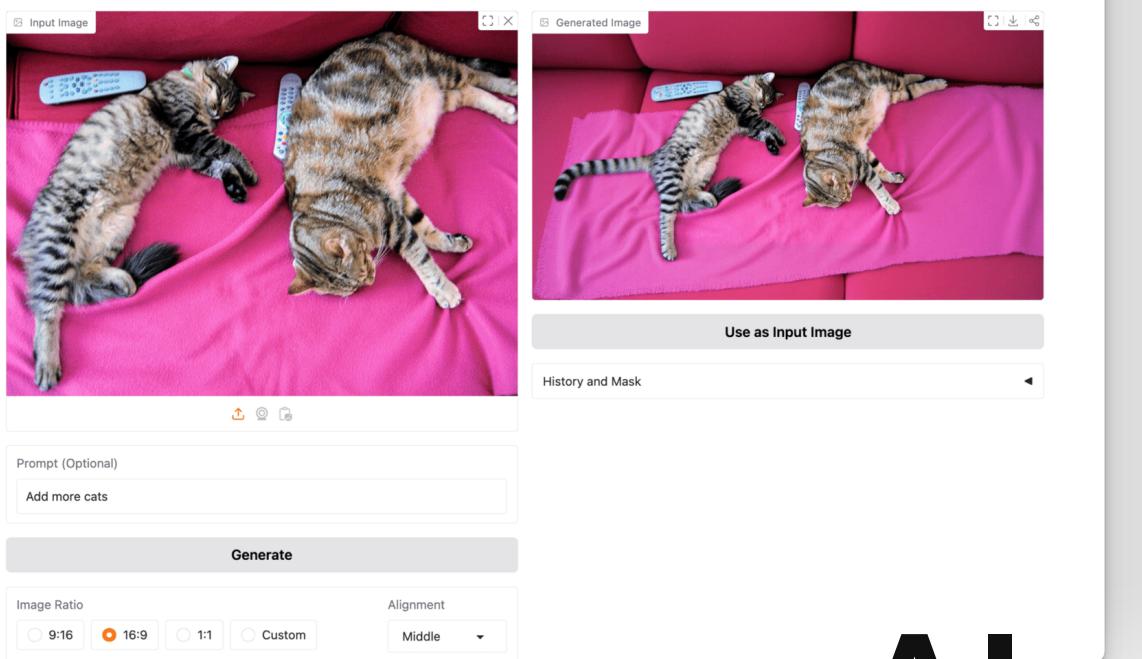
seawolf2357/flux-lora-car-rolls-royce

Text-to-Image • Updated Sep 17 • ↓ 339 • ⚡ • ❤ 184

Qwen/Qwen2.5-Coder-32B-Instruct

Text Generation • Updated 16 days ago • ↓ 170k • ⚡ • ❤ 1.16k

AI devs



AI devs

GaussianAnything: Interactive Point Cloud Latent Diffusion for 3D Generation

GaussianAnything (arXiv 2024) [code, project page] is a native 3D diffusion model that supports high-quality 2D Gaussians generation. It first trains a 3D VAE on Objaverse, which compress each 3D asset into a compact point cloud-latent latent. After that, a image/text-conditioned diffusion model is trained following LDM paradigm. The model used in the demo adopts 3D DIT architecture and flow-matching framework, and supports single-image condition. It is trained on 8 A100 GPUs for 1M iterations with batch size 256. Locally, on an NVIDIA A100/A10 GPU, each image-conditioned diffusion generation can be done within 20 seconds (time varies due to the adaptive-step ODE solver used in flow-matching.) Upload an image of an object or click on one of the provided examples to see how the GaussianAnything works.

The 3D viewer will render a .glb point cloud exported from the centers of the surfel Gaussians, and an integrated TSDF mesh. Besides, you can find the intermediate stage-1 point cloud in the Tab (Stage-1 Output).

For best results run the demo locally and render locally - to do so, clone the [main repository](#).

Input Image: A blue fox. **Processed Image:** A blue fox with a white patch on its chest. **Stage-2 Output:** A blue fox with a white patch on its chest. **Stage-1 Output:** Two versions of the blue fox's head, one rendered in a standard color palette and one in a vibrant, multi-colored (purple, red, yellow) style.

CFG-scale: 4. **Seed:** 42. **Remove Background:** **Generate:**

Chatbot: Sure, I'd be happy to help you plan a one-day trip! Let's see, there are so many cities to choose from, but I think I'll pick Tokyo, Japan. It's a fascinating city with a unique blend of traditional and modern culture. Even though it's a big city, it's surprisingly easy to navigate, especially for tourists.

First things first, let's figure out how to get around. Tokyo has an extensive public transportation system, including subways, buses, and trains. The best way to tackle this is to get a prepaid transportation card like a Suica or Pasmo card. These cards can be used on all forms of public transport and make traveling much smoother. You can purchase them at any train station; just look for the ticket machines that have the Suica or Pasmo logo.

Now, for the itinerary. Since it's a one-day trip, we need to make the most of our time. Let's assume we're starting in the morning and have about 10 hours to explore.

Morning:

We'll begin our day in the historic district of Asakusa. Asakusa is known for its traditional atmosphere and is home to the famous Senso-ji Temple, which is Tokyo's oldest temple. It's a great place to experience a bit of old Tokyo. **Error**

- Activity:** Visit Senso-ji Temple. Make sure to walk through the Nakamise Shopping Street leading up to the temple. It's packed with souvenir shops and food stalls, so you can grab a quick breakfast there if you like. Try some traditional Japanese snacks like taiyaki (fish-shaped cake filled with red bean paste) or kinako-dango (rice flour dumplings coated in soybean flour).
- Interesting Trivia:** Did you know that Senso-ji Temple was established in the year 645? That makes it over 1,300 years old! It's a significant historical site and a must-see when in Tokyo.

After visiting the temple, let's head over to the nearby Asakusa Shrine. It's a beautiful Shinto shrine with a peaceful atmosphere, perfect for taking a break and soaking in the surroundings.

Late Morning to Early Afternoon:

S先生、P先生、Q先生他们知道桌子的抽屉里有16张扑克牌：红桃A、Q、4、黑桃J、8、4、2、7、3草花K、Q、5、4、6方块J、5。约翰教授从这16张牌中挑出一张牌来，并把这张牌的点数告诉P先生。把这张牌的花色告诉Q先生。这时，约翰教授问P先生和Q先生：你们能从已知的点数或花色中推知这张牌是什么牌吗？于是，S先生听到如下的对话：P先生：我不知道这张牌。Q先生：我知道你不知道这张牌。P先生：现在我知道这张牌了。Q先生：我也知道了。请问：这张牌是什么牌？

QwQ-32B-preview: QwQ-32B-Preview is an experimental research model developed by the Qwen Team, focused on advancing AI reasoning capabilities. As a preview release, it demonstrates promising analytical abilities while having several important limitations such as code switching and recursive reasoning loops. Only single-turn queries are supported in this demo.

推荐对话: How many r in strawberry
Find the least odd prime factor of 2019^{8+1} .
S先生、P先生、Q先生他们知道桌子的抽屉里有16张扑克牌：红桃A、Q、4、黑桃J、8、4、2、7、3草花K、Q、5、4、6方块J、5。约翰教授从这16张牌中挑出一张牌来，并把这张牌的点数告诉P先生。把这张牌的花色告诉Q先生。这时，约翰教授问P先生和Q先生：你们能从已知的点数或花色中推知这张牌是什么牌吗？于是，S先生听到如下的对话：P先生：我不知道这张牌。Q先生：我知道你不知道这张牌。P先生：现在我知道这张牌了。Q先生：我也知道了。请问：这张牌是什么牌？



Our Open Source

We are building the foundation of ML tooling with the community.

● **Transformers**

109,705

State-of-the-art ML for Pytorch,
TensorFlow, and JAX.

● **Tokenizers**

7,347

Fast tokenizers, optimized for
both research and production.

● **Diffusers**

17,025

State-of-the-art diffusion models
for image and audio generation
in PyTorch.

● **PEFT**

8,628

Parameter efficient finetuning
methods for large models

● **Safetensors**

1,323

Simple, safe way to store and
distribute neural networks
weights safely and quickly.

● **Transformers.js**

3,270

Community library to run
pretrained models from
Transformers in your browser.

● **Hub Python Library**

1,069

Client library for the HF Hub:
manage repositories from your
Python runtime.

● **timm**

26,274

State-of-the-art computer vision
models and optimizers



In-browser AI apps with
Gradio and
Transformers



Transformers

Transformers provides APIs and tools
to easily download and train state-of-the-art pretrained models.



```
1 >>> from transformers import pipeline
2 >>> classifier = pipeline("sentiment-analysis")
3 >>> classifier("We are very happy to show you the 😊 Transformers library.")
4 [ {'label': 'POSITIVE', 'score': 0.9998} ]
```

[huggingface / transformers](#) Public

Notifications

Fork 27.2k

Star 136k

Code

Issues 968

Pull requests 510

Actions

Projects 1

Security

Insights

main

613 Branches 181 Tags

Go to file

Code

About



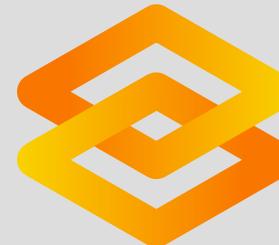
Cyrilvallez Automatic compilation in generate: do not rely on inner fun...



ee37bf0 · 36 minutes ago

17,497 Commits

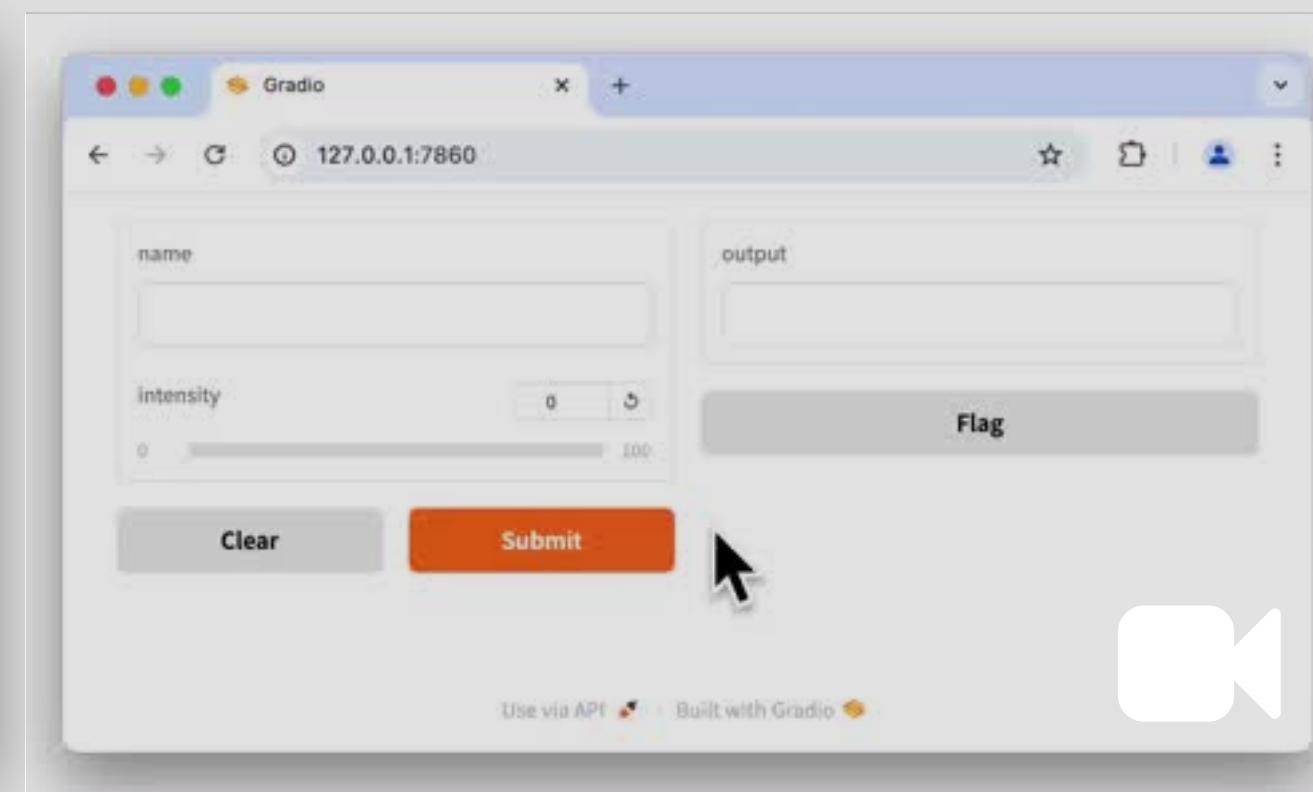
🤗 Transformers: State-of-the-art
Machine Learning for Pytorch,
TensorFlow and JAX



gradio

Gradio is an open-source Python package that allows you to quickly build a demo or web application for your machine learning model, API, or any arbitrary Python function.

```
1 import gradio as gr
2
3 def greet(name, intensity):
4     return "Hello, " + name + "!" * int(i
5
6 demo = gr.Interface(
7     fn=greet,
8     inputs=["text", "slider"],
9     outputs=["text"],
10)
11
12 demo.launch()
```





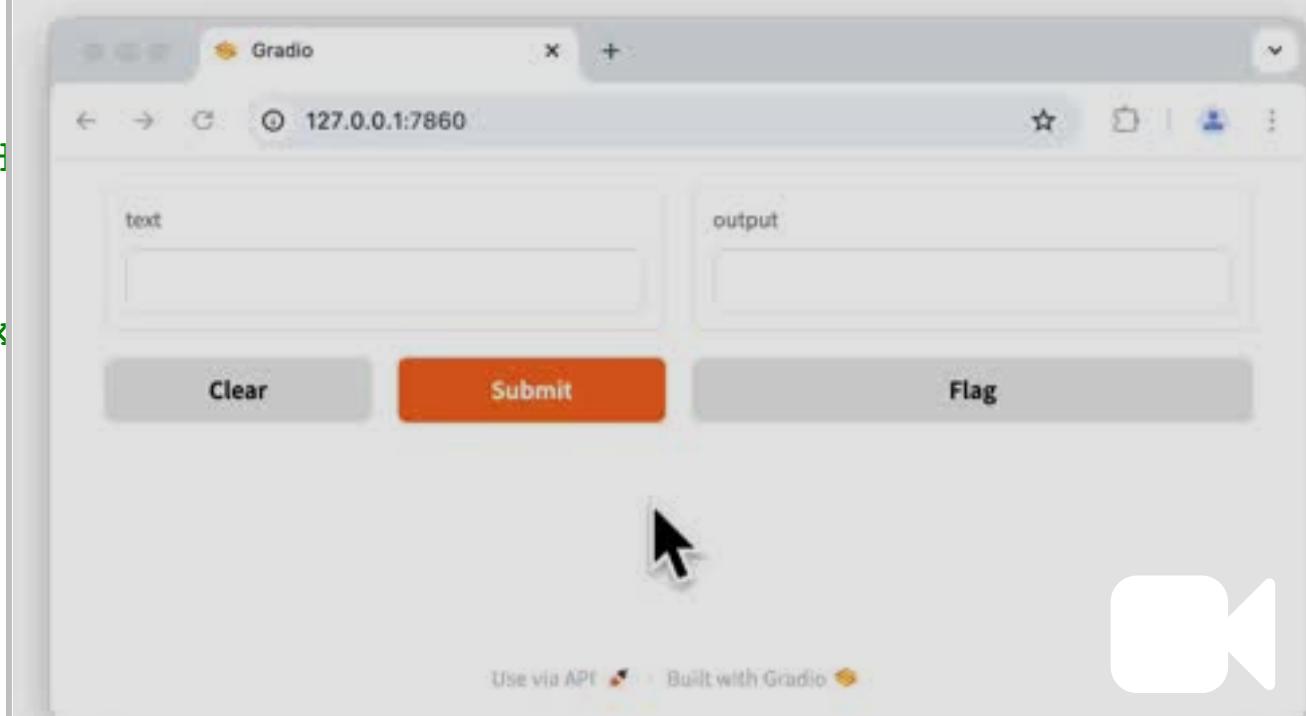
gradio



Transformers



```
1 import gradio as gr
2 from transformers import pipeline
3
4 pipe = pipeline("translation", model="Helsinki-NLP/opus-mt-es")
5
6 def predict(text):
7     return pipe(text)[0]["translation_text"]
8
9 demo = gr.Interface(
10     fn=predict,
11     inputs='text',
12     outputs='text',
13 )
14
15 demo.launch()
```



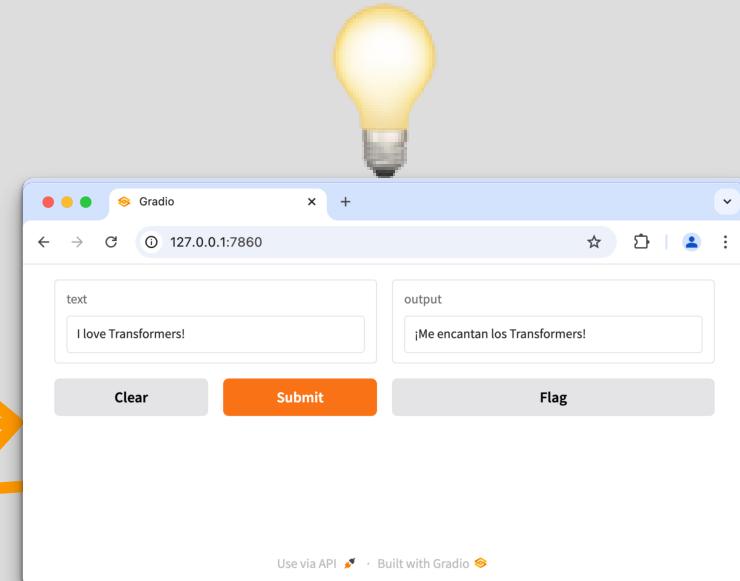
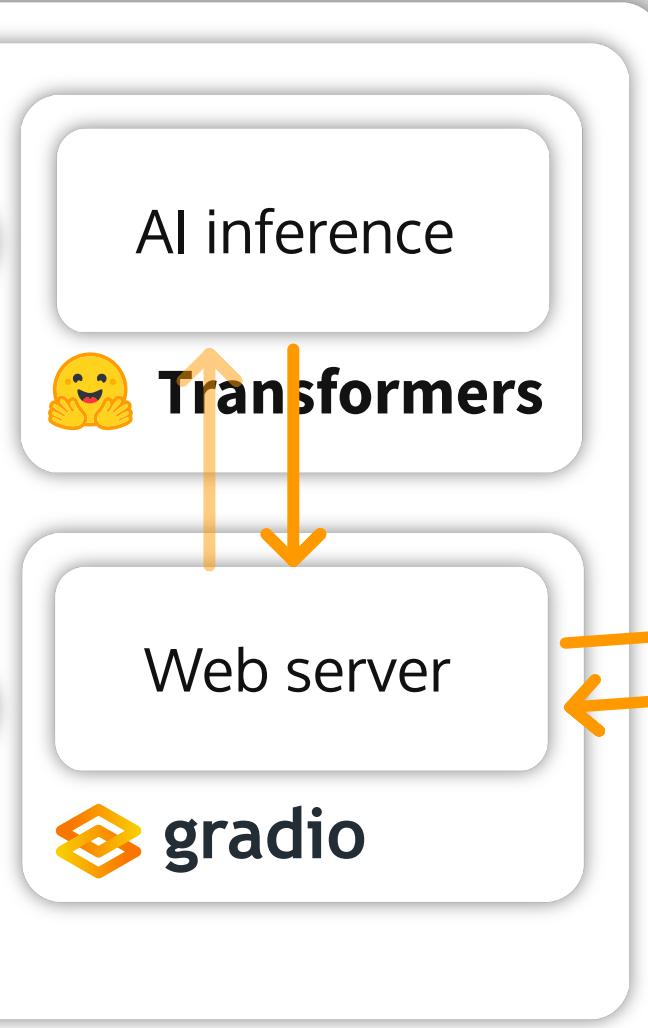


gradio



Transformers

```
1 import gradio as gr
2 from transformers import pipeline
3
4 pipe = pipeline("translation_en_es", model="Helsinki-NLP/opus-mt-en-es")
5
6 def predict(text):
7     return pipe(text)[0]["translation_text"]
8
9 demo = gr.Interface(
10     fn=predict,
11     inputs='text',
12     outputs='text',
13 )
14
15 demo.launch()
```



Python runtime

Server

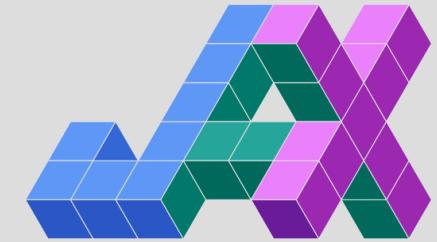
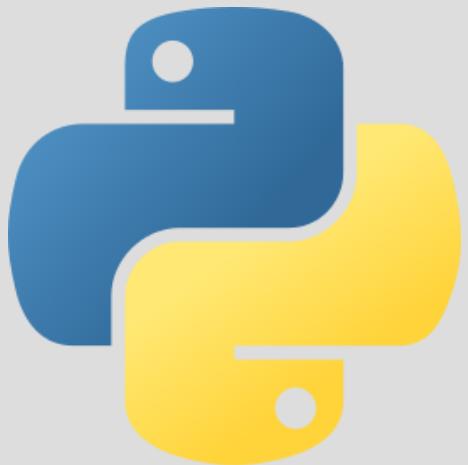


Transformers

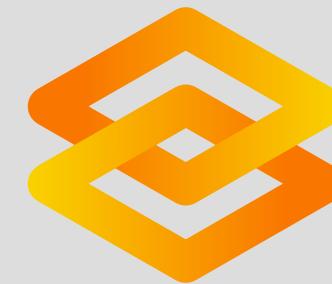
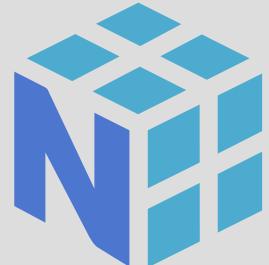
Quick & Easy
Web-based AI apps

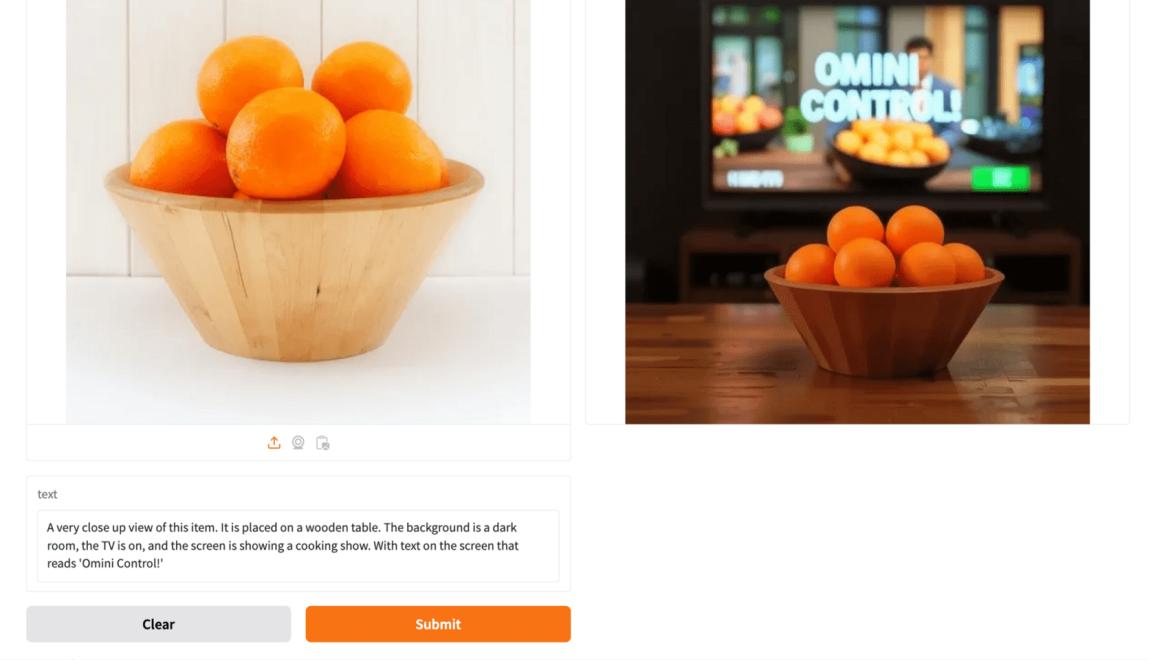
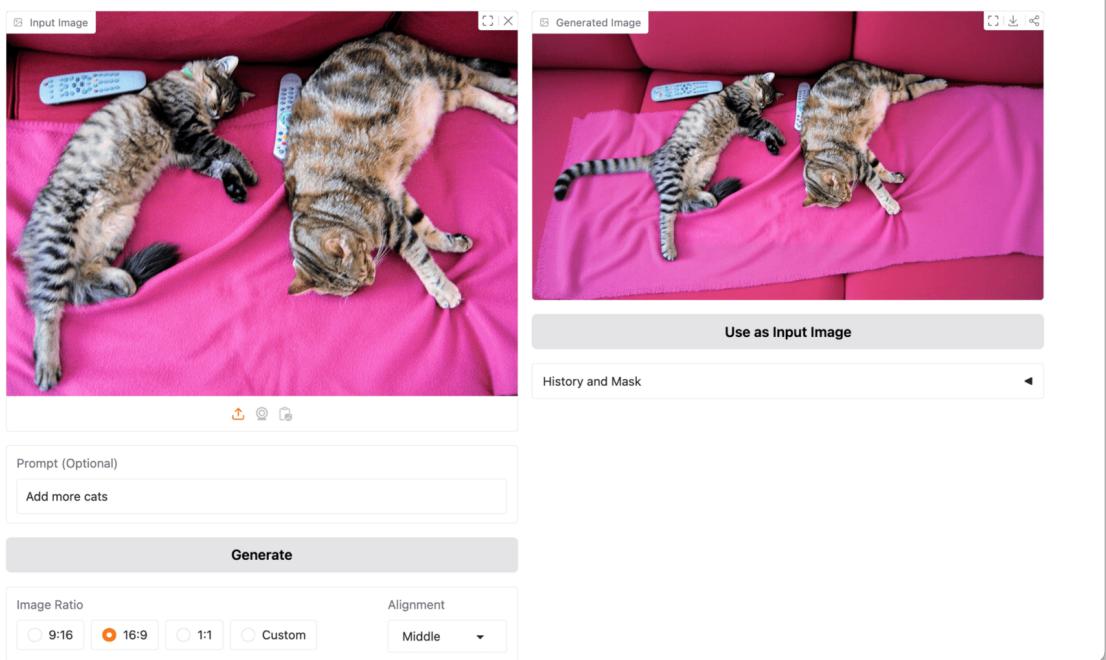
Python at FEDAY?





A central part of the AI/ML ecosystem.





In-browser AI apps with
Gradio and
Transformers

In-browser AI apps

a.k.a. frontend-only AI apps

Web-based applications that use machine learning models to perform tasks like text analysis, image processing, or speech recognition **entirely within the user's browser.**

Why in-browser?

- Privacy
- Low latency
- Offline capability
- Scalability without servers
- Low cost

In-browser AI apps with
Gradio and
Transformers



Transformers

Quick & Easy
Web-based AI apps



gradio



Transformers



```
1 import gradio as gr
2 from transformers import pipeline
3
4 pipe = pipeline("translation", model="Helsinki-NLP/opus-mt-es")
5
6 def predict(text):
7     return pipe(text)[0]
8
9 demo = gr.Interface(fn=predict,
10                     inputs='text',
11                     outputs='text',
12                     )
13
14
15 demo.launch()
```



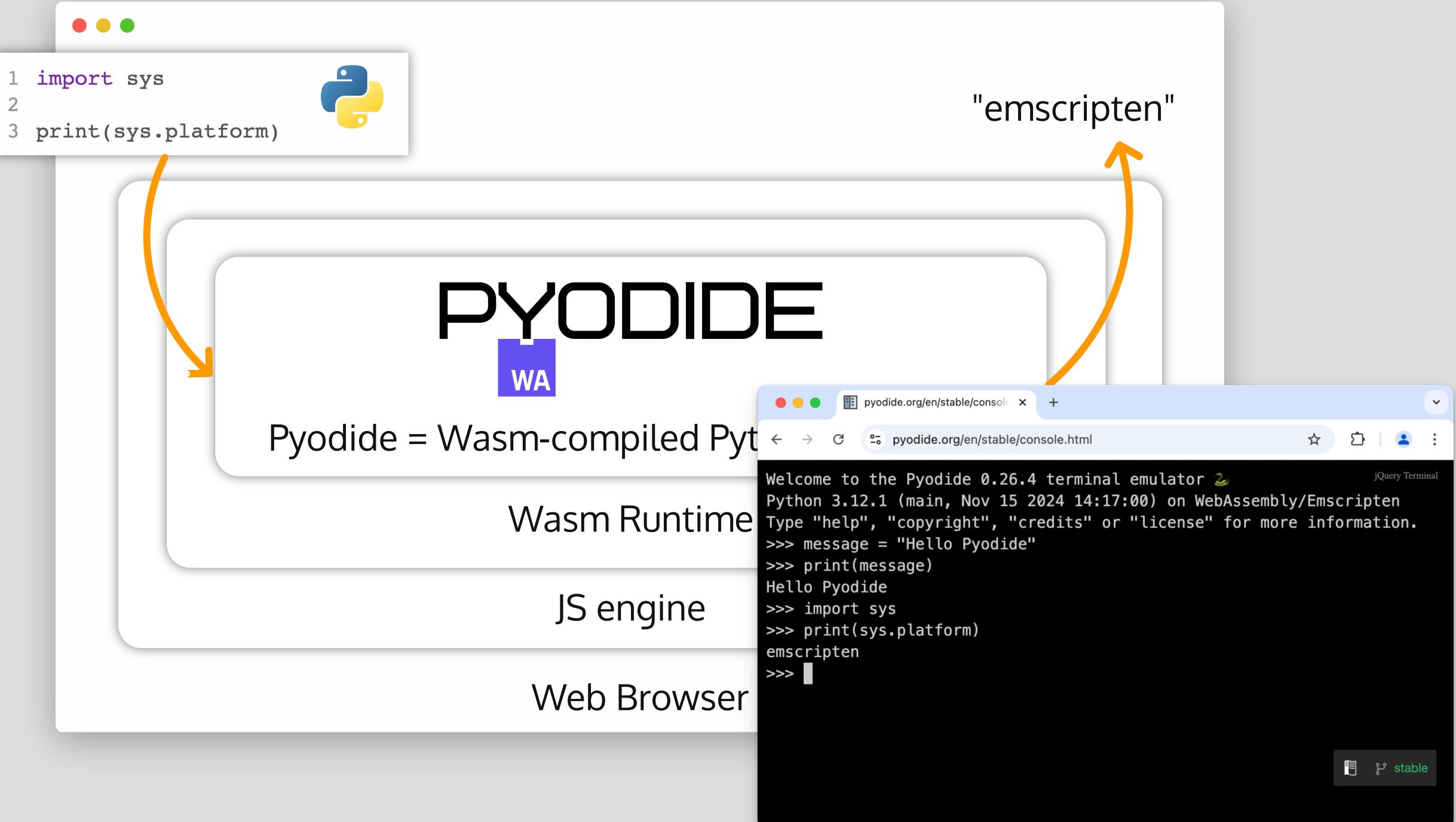
The screenshot shows a Gradio web application window titled "Gradio" at the URL "127.0.0.1:7860". The interface has two main sections: "text" and "output". In the "text" section, there is an input field containing the text "I love Transformers!". Below it are three buttons: "Clear", "Submit" (which is orange), and "Flag". In the "output" section, there is another input field showing the translated text "¡Me encantan los Transformers!". At the bottom of the window, there are links for "Use via API" and "Built with Gradio".

... Python on Frontend?



Python interpreter on a browser

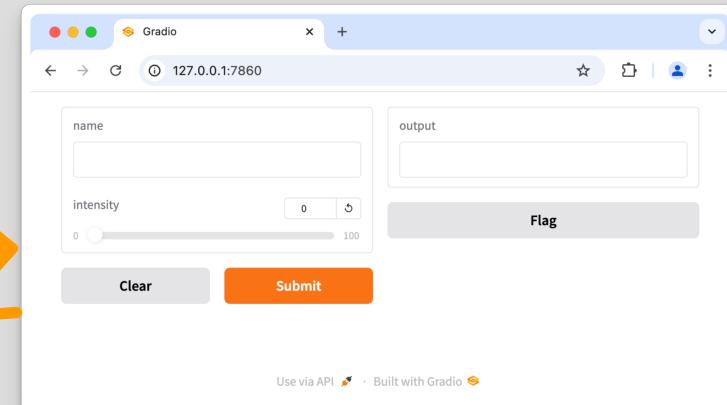
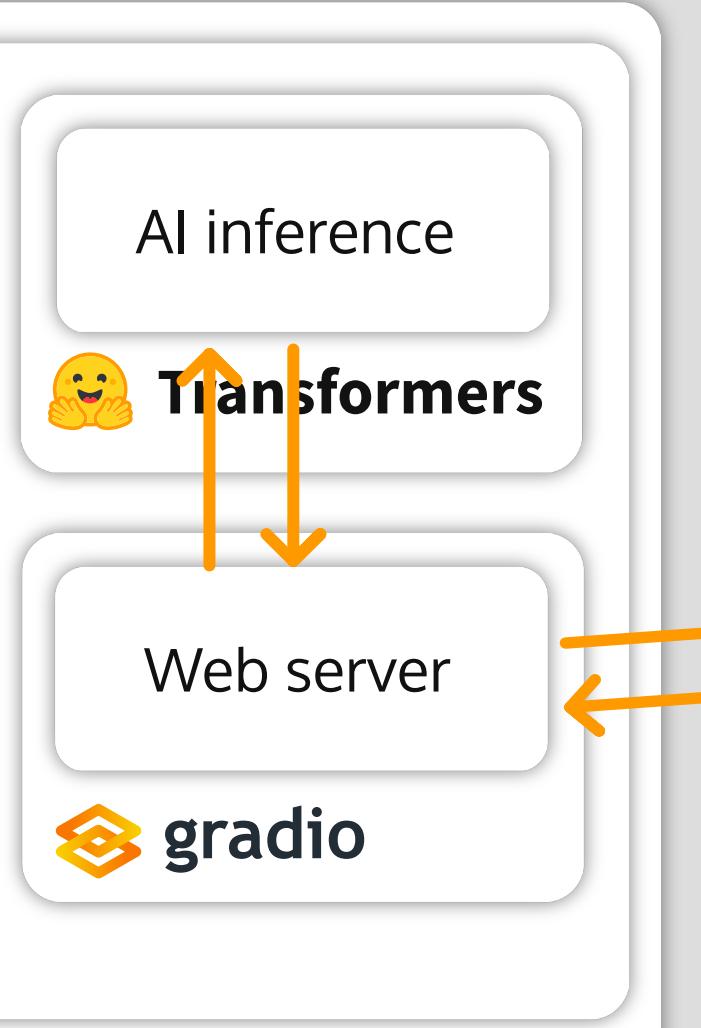
Pyodide is a Python distribution for the browser and Node.js based on WebAssembly.





Transformers

```
1 import gradio as gr
2 from transformers import pipeline
3
4 pipe = pipeline("translation", model="Helsinki-NLP/opus-mt-en-es")
5
6 def predict(text):
7     return pipe(text)[0]["translation_text"]
8
9 demo = gr.Interface(
10     fn=predict,
11     inputs='text',
12     outputs='text',
13 )
14
15 demo.launch()
```



Python runtime

Server

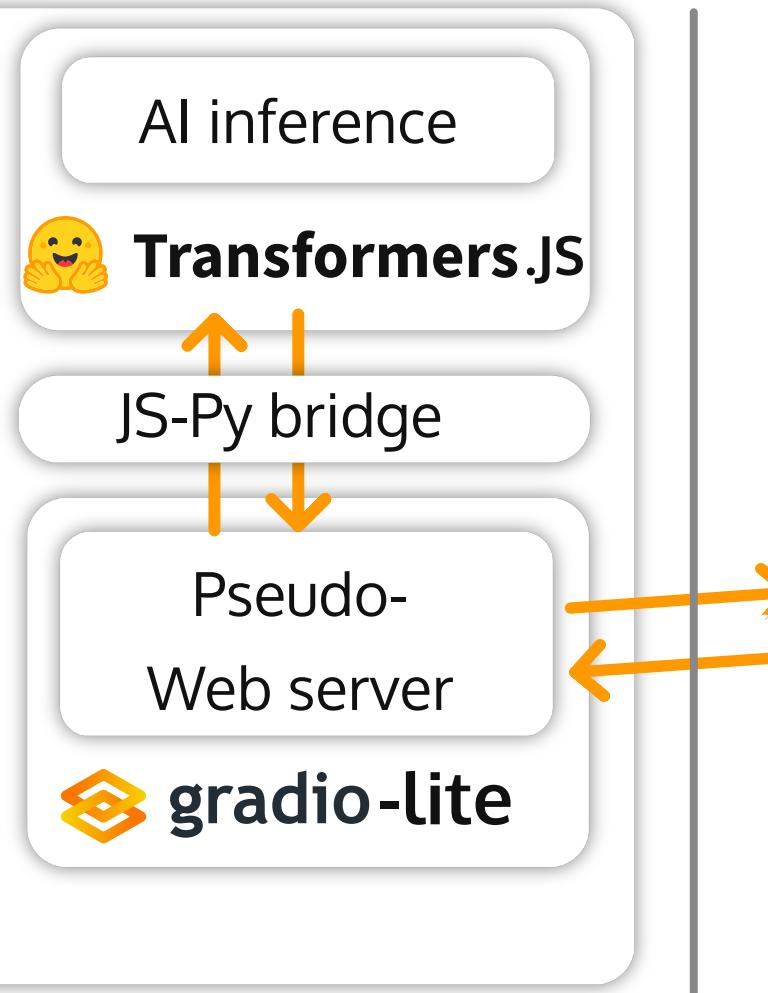
gradio-lite



Transformers.js

```
1 import gradio as gr
2 from transformers_js_py import pipeline
3
4 pipe = await pipeline("translation",
5
6     async def predict(text):
7         res = await pipe(text, {
8             "src_lang": 'en',
9             "tgt_lang": 'zh',
10        })
11     return res[0]["translation_text"]
12
13 demo = gr.Interface(
14     fn=predict,
15     inputs='text',
16     outputs='text',
17 )
18
19 demo.launch()
```

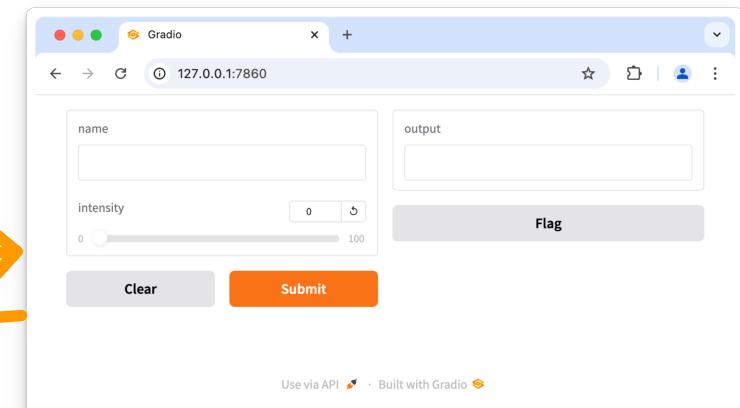
PYODIDE runtime
WA



Web Browser

Worker process

Renderer process



Gradio-Lite

Pyodide/Wasm-ported Gradio.

You write Python code, then get a web UI, 100% in the browser.

Transformers.js

JS version of Transformers.

You can use pretrained AI/ML models in the browser.

Check out the presentation ↗ <https://www.bilibili.com/video/BV19c411B7QU/>

Transformers.js

 首页 番剧 直播 游戏中心 会员购 漫画 赛事 吉罗很快乐 · 3小时前更新 登录 搜索

Transformers.js: Web 上的最新机器学习技术
1.8万 13 2023-11-22 23:11:49 未经作者授权, 禁止转载



通过 Transformers.js 这个开源 JavaScript 库

1人正在看, 已装填 13 条弹幕 请先 登录 或 注册 弹幕礼仪 > 发送 测试版 AI小助手 记笔记

537 211 1044 189

我们的工程师, Transformers.js 作者 Joshua Lochner 在 2023 年 11 月 18 日在杭州举办的 FEDAY 上带来了主题为「Transformers.js: State-of-the-art Machine Learning for the Web」的分享 😊

通过 Joshua 的介绍, 你将了解 Hugging Face 如何通过 Transformers.js 将最先进的机器学习带到 Web 中, 以及对 WebML 技术的优势和如何将其用于加速 Web 应用程序有更多的认识 🚀

收起

HuggingFace ✉ 发消息
The AI community building the future.
+ 充电 + 关注 7.9万

弹幕列表

接下来播放 自动连播

- Hugging Face 😊 Transformers.js v3 正式发布 🎉
HuggingFace 0:01 3.1万 0
- 【保姆级教程】带你彻底啃透AI顶会论文!
bilibili课堂
- 【Hugging Face 课程 😊】3D 机器学习 - 第一节 简介
HuggingFace 0:42 2.4万 3
- 【使用 Gradio 创建聊天机器人】Create Your Own Gradi...
HuggingFace 54:09 1915 1
- 【中文科普】😊 Hugging Face 与 Intel 共同构建生成式 AI
HuggingFace 27:54 2976 2
- 【Hugging Face 课程 😊】3D 机器学习 - 第四节 高斯...
HuggingFace 03:13 4438 3

最新发布! HuggingChat ma



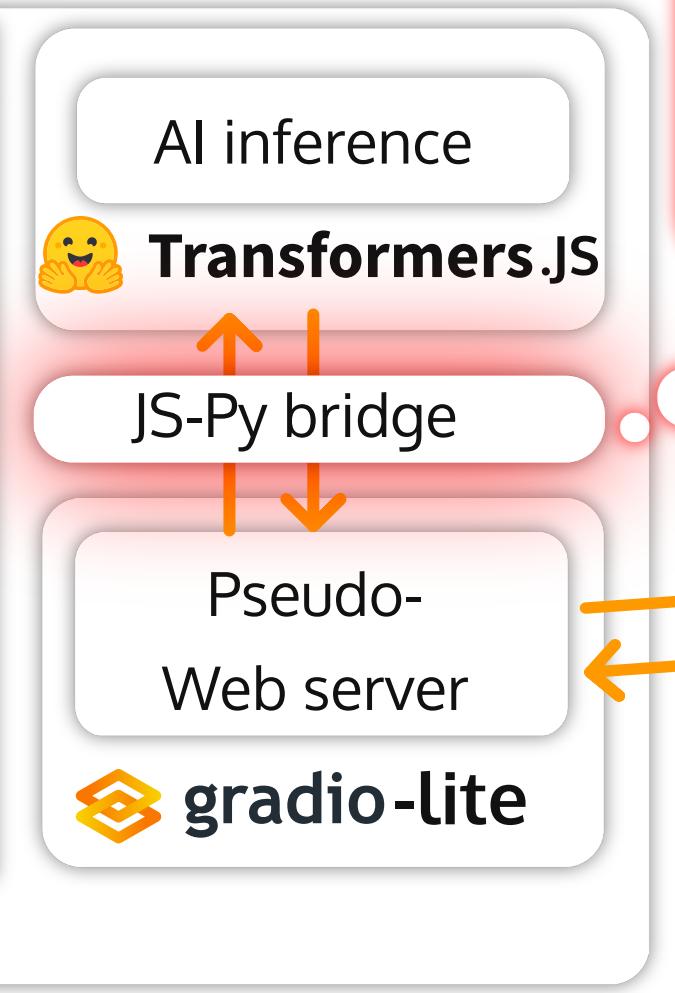
<https://www.bilibili.com/video/BV19c411B7QU/>



Transformers.js

```
1 import gradio as gr
2 from transformers_js_py import pipeline
3
4 pipe = await pipeline("translation",
5
6     async def predict(text):
7         res = await pipe(text, {
8             "src_lang": 'en',
9             "tgt_lang": 'zh',
10        })
11     return res[0]["translation_text"]
12
13 demo = gr.Interface(
14     fn=predict,
15     inputs='text',
16     outputs='text',
17 )
18
19 demo.launch()
```

PYODIDE runtime
WA

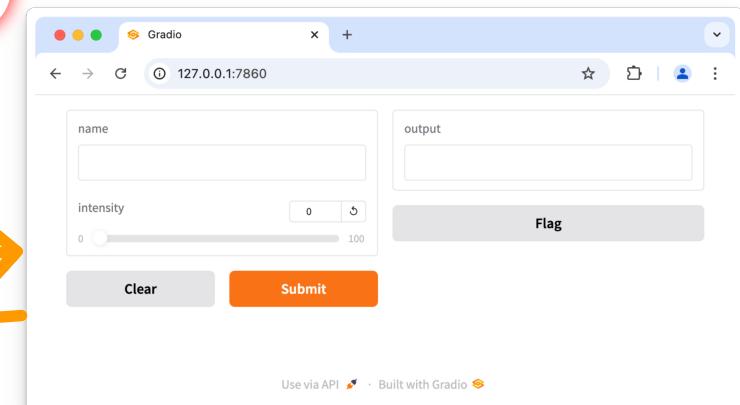


Web Browser

Worker process

Transformers.js.py

Python wrapper of
Transformers.js 😱



Renderer process

Gradio-Lite

Pyodide/Wasm-ported Gradio.

You write Python code, then get a web UI, 100% in the browser.

Transformers.js

JS version of Transformers.

You can use pretrained AI/ML models in the browser.

Check out the presentation ↗ <https://www.bilibili.com/video/BV19c411B7QU/>

Transformers.js.py

Use Transformers.js from Pyodide.



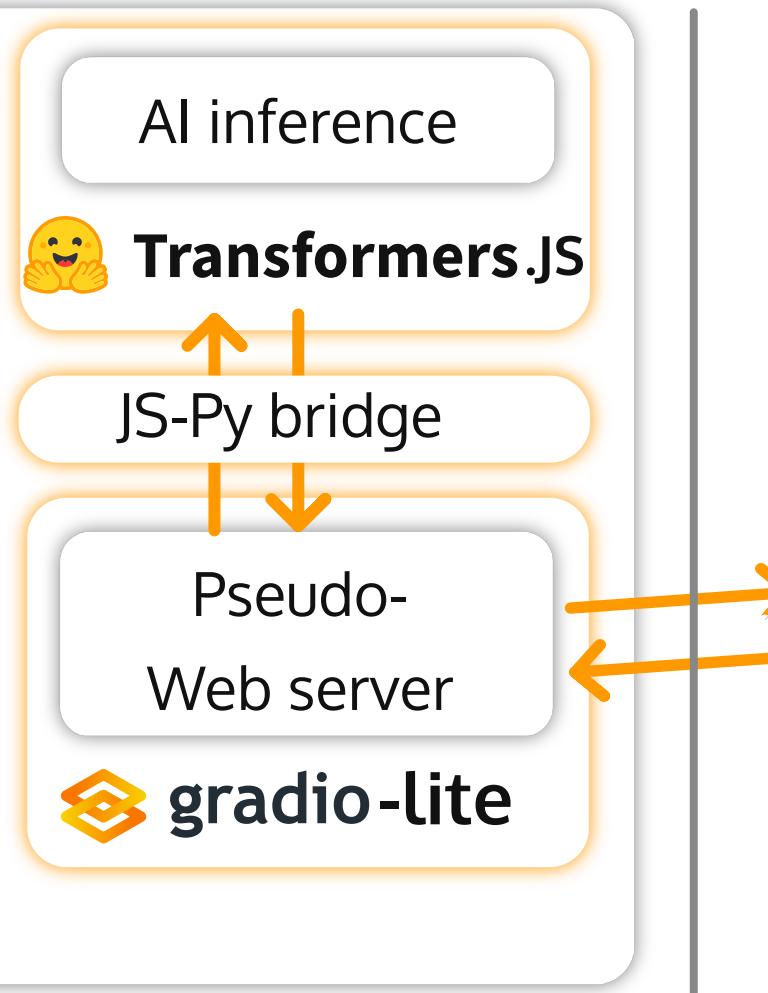
gradio-lite



Transformers.js

```
1 import gradio as gr
2 from transformers_js_py import pipeline
3
4 pipe = await pipeline("translation",
5
6     async def predict(text):
7         res = await pipe(text, {
8             "src_lang": 'en',
9             "tgt_lang": 'zh',
10        })
11     return res[0]["translation_text"]
12
13 demo = gr.Interface(
14     fn=predict,
15     inputs='text',
16     outputs='text',
17 )
18
19 demo.launch()
```

PYODIDE runtime
WA



Web Browser

Worker process

Renderer process

How to use them



```
1 $ code index.html
```



```
1 <html>
2   <head>
3     <script type="module" crossorigin src="https://cdn.j
4       <link rel="stylesheet" href="https://cdn.jsdelivr.net
5   </head>
6   <body>
7
8
9   </body>
10 </html>
```



```
1 <html>
2   <head>
3     <script type="module" crossorigin src="https://cdn.j
4       <link rel="stylesheet" href="https://cdn.jsdelivr.net
5   </head>
6   <body>
7     <gradio-lite>
8
9     </gradio-lite>
10    </body>
11 </html>
```



```
1 <html>
2     <head>
3         <script type="module" crossorigin src="https://cdn.j
4             <link rel="stylesheet" href="https://cdn.jsdelivr.net/
5     </head>
6     <body>
7         <gradio-lite>
8 import gradio as gr
9
10 def greet(name):
11     return "Hello, " + name + "!"
12
13 gr.Interface(greet, "textbox", "textbox").launch()
14         </gradio-lite>
15     </body>
16 </html>
```

index.html

File /Users/whitphx/src/gradio-lite-trial/index.html

x

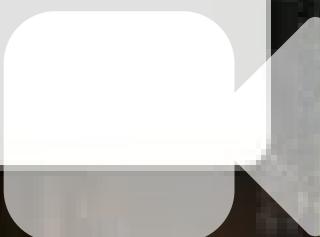
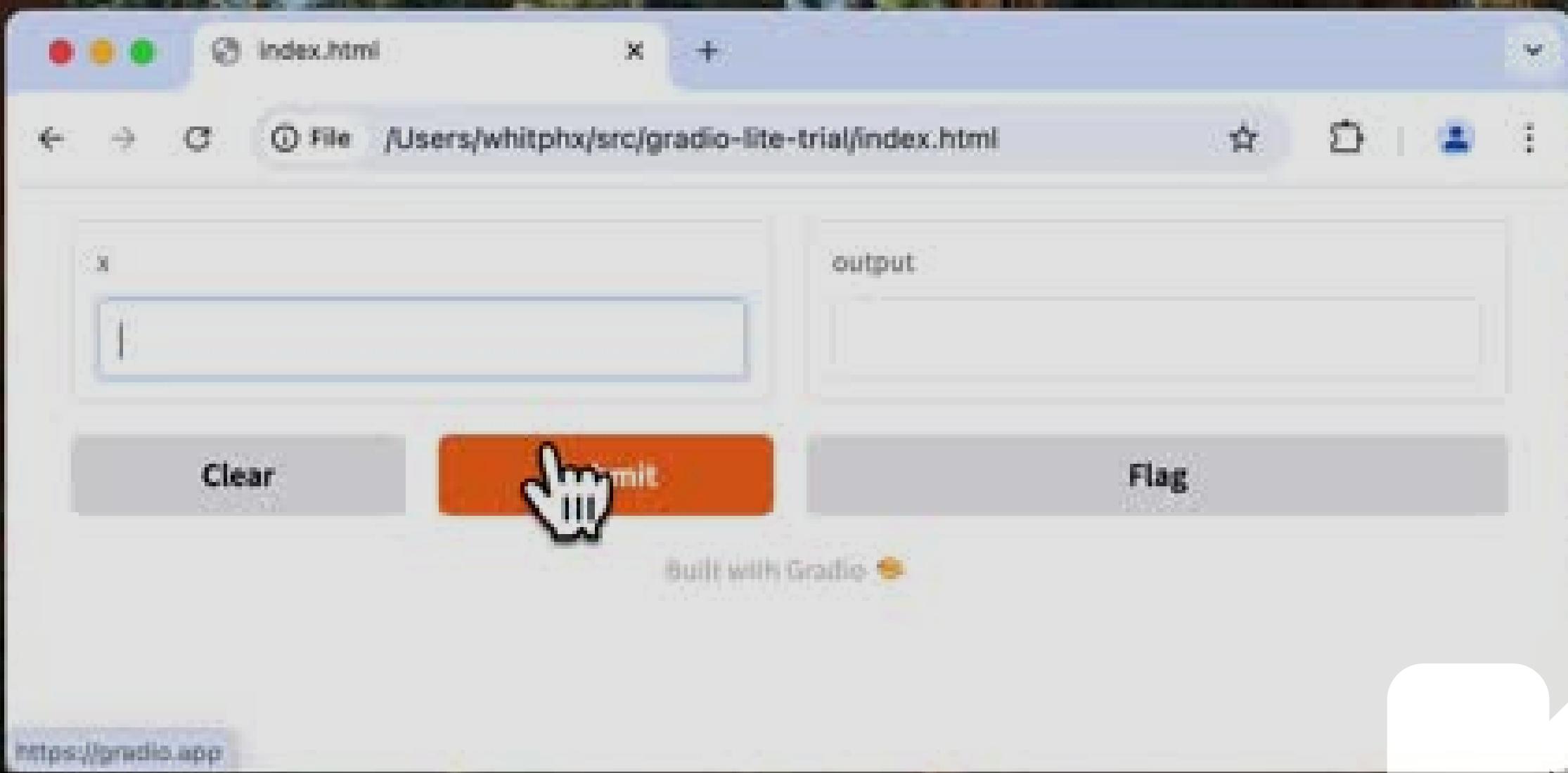
output

Clear

Submit

Flag

Built with Gradio



```
2 <head>
3     <script type="module" crossorigin src="https://cdn.j
4         <link rel="stylesheet" href="https://cdn.jsdelivr.net
5 </head>
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline('sentiment-analysis')
12
13 demo = gr.Interface.from_pipeline(pipe)
14
15 demo.launch()
16
17     <gradio-requirements>
18 transformers-js-py
19     </gradio-requirements>
```

index.html

File /Users/whitphx/src/gradio-lite-trial/index.html

text-classification (nlptown/bert-base-multilingual-uncased-sentiment)

Input

Top k

Classification

Flag

Clear

Submit

Built with Gradio

index.html

File /Users/whitphx/src/gradio-lite-trial/index.html

text-classification (distilbert-base-uncased-finetuned-sst-2-english)

input

Classification

Top k

Flag

Clear

Submit

Built with Gradio

```
2 <head>
3     <script type="module" crossorigin src="https://cdn.j
4         <link rel="stylesheet" href="https://cdn.jsdelivr.net
5 </head>
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline('sentiment-analysis')
12
13 demo = gr.Interface.from_pipeline(pipe)
14
15 demo.launch()
16
17     <gradio-requirements>
18 transformers-js-py
19     </gradio-requirements>
20     />
```

```
2 <head>
3     <script type="module" crossorigin src="https://cdn.j
4         <link rel="stylesheet" href="https://cdn.jsdelivr.net
5 </head>
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline(
12     'sentiment-analysis',
13     'Xenova/bert-base-multilingual-uncased-sentiment'
14 )
15
16 demo = gr.Interface.from_pipeline(pipe)
17
18 demo.launch()
19
20
```

Hugging Face

Search

Models Datasets Spaces Docs Enterprise Pricing Log In Sign Up

Xenova/bert-base-multilingual-uncased-sentiment

Text Classification Transformers.js ONNX bert

Model card Files Community

Use this model

Downloads last month
581

Inference Examples

Text Classification

Inference API (serverless) does not yet support transformers.js models for this pipeline type.

Model tree for Xenova/bert-base-multilingual-uncased-sentiment

Base model nlptown/bert-base-multilingual-uncased

Quantized (1) this model

<https://huggingface.co/nlptown/bert-base-multilingual-uncased-sentiment> with ONNX weights to be compatible with Transformers.js.

Note: Having a separate repo for ONNX weights is intended to be a temporary solution until WebML gains more traction. If you would like to make your models web-ready, we recommend converting to ONNX using 😊 Optimum and structuring your repo like this one (with ONNX weights located in a subfolder named `onnx`).

index.html

File /Users/whitphx/src/gradio-lite-trial/index.html

text-classification (nlptown/bert-base-multilingual-uncased-sentiment)

Classification

Input

Top k

5

Flag

Clear

Submit

Built with Gradio

OK

```
2 <head>
3     <script type="module" crossorigin src="https://cdn.j
4         <link rel="stylesheet" href="https://cdn.jsdelivr
5 </head>
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline(
12     'sentiment-analysis',
13     'Xenova/bert-base-multilingual-uncased-sentiment'
14 )
15
16 demo = gr.Interface.from_pipeline(pipe)
17
18 demo.launch()
19
20
```

```
2 <head>
3     <script type="module" crossorigin src="https://cdn.j
4         <link rel="stylesheet" href="https://cdn.jsdelivr.net
5 </head>
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline('image-classification')
12
13 demo = gr.Interface.from_pipeline(pipe)
14
15 demo.launch()
16
17     <gradio-requirements>
18 transformers-js-py
19     </gradio-requirements>
20     />
```

← → ⟳

File

/Users/whitphx/src/gradio-lite-trial/index.html

☆ ⤙ ⤚ 👤 ⋮

image-classification (google/vit-base-patch16-224)

Input Image

Drop Image Here

- or -

Click to Upload

**Classification****Flag****Top k**

5

Clear**Submit**



Index.html



File /Users/whitphobjec/pradio-lite-trial/Index.html



image-classification (google/vit-base-patch16-224)

Import Image



Drop Image Here



Click to Upload

Classification



Flag

Top 5



Clear

Submit

Build your pipeline



```
2 <head>
3     <script type="module" crossorigin src="https://cdn.j
4         <link rel="stylesheet" href="https://cdn.jsdelivr.net
5 </head>
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline('sentiment-analysis')
12
13 demo = gr.Interface.from_pipeline(pipe)
14
15 demo.launch()
16
17     <gradio-requirements>
18 transformers-js-py
19     </gradio-requirements>
20     />
```

```
6 <body>
7     <gradio-lite>
8 import gradio as gr
9 from transformers_js_py import pipeline
10
11 pipe = await pipeline('sentiment-analysis')
12
13 async def fn(text):
14     result = await pipe(text)
15     return result
16
17 demo = gr.Interface(
18     fn=fn,
19     inputs=gr.Textbox(),
20     outputs=gr.JSON(),
21 )
22
23 demo.launch()
```

The image shows a web browser window with the title "index.html". The address bar indicates the file is located at "/Users/whitphx/src/gradio-lite-trial/index.html". The main content area contains a text input field with the placeholder "text" and the value "I love Transformers!". Below it are two buttons: "Clear" (gray) and "Submit" (orange). To the right is a JSON output panel titled "output" which displays the following data:

```
1 |   [▼
2 |     ▼ "0": {
3 |       "label": "POSITIVE",
4 |       "score": 0.999788761138916
5 |     }
6 |   ]
```

At the bottom of the page is a gray button labeled "Flag". The footer of the page says "Built with Gradio 🚀".

index.html

File /Users/whitphx/src/gradio-lite-trial/index.html

test

Clear Submit

{ } output

Flag

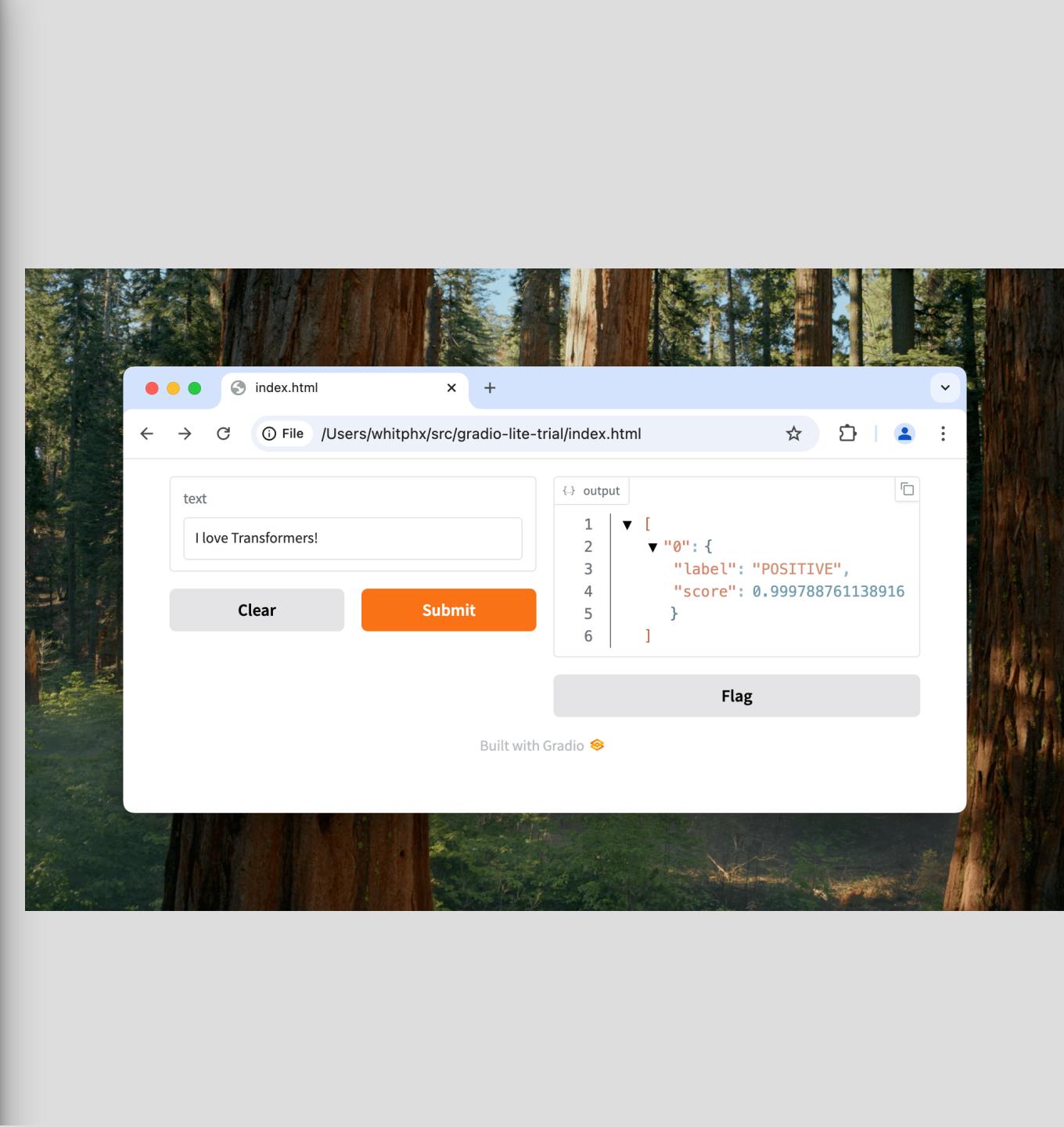
Built with Gradio

→

→



```
1 <html>
2     <head>
3         <script type="module" crossorigin src="https://unpkg.com/react@17.0.2/umd/react.development.js"></script>
4         <link rel="stylesheet" href="https://unpkg.com/react-dom@17.0.2/umd/react-dom.development.css">
5     </head>
6     <body>
7         <gradio-lite>
8             import gradio as gr
9             from transformers_js_py import pipeline
10            pipe = await pipeline('sentiment-analysis')
11
12            async def fn(text):
13                result = await pipe(text)
14                return result
15
16            demo = gr.Interface(
17                fn=fn,
18                inputs=gr.Textbox(),
19                outputs=gr.JSON(),
20            )
21
22            demo.launch()
23
24
25             <gradio-requirements>
26             transformers-js-py
27             </gradio-requirements>
28             </gradio-lite>
29         </body>
```

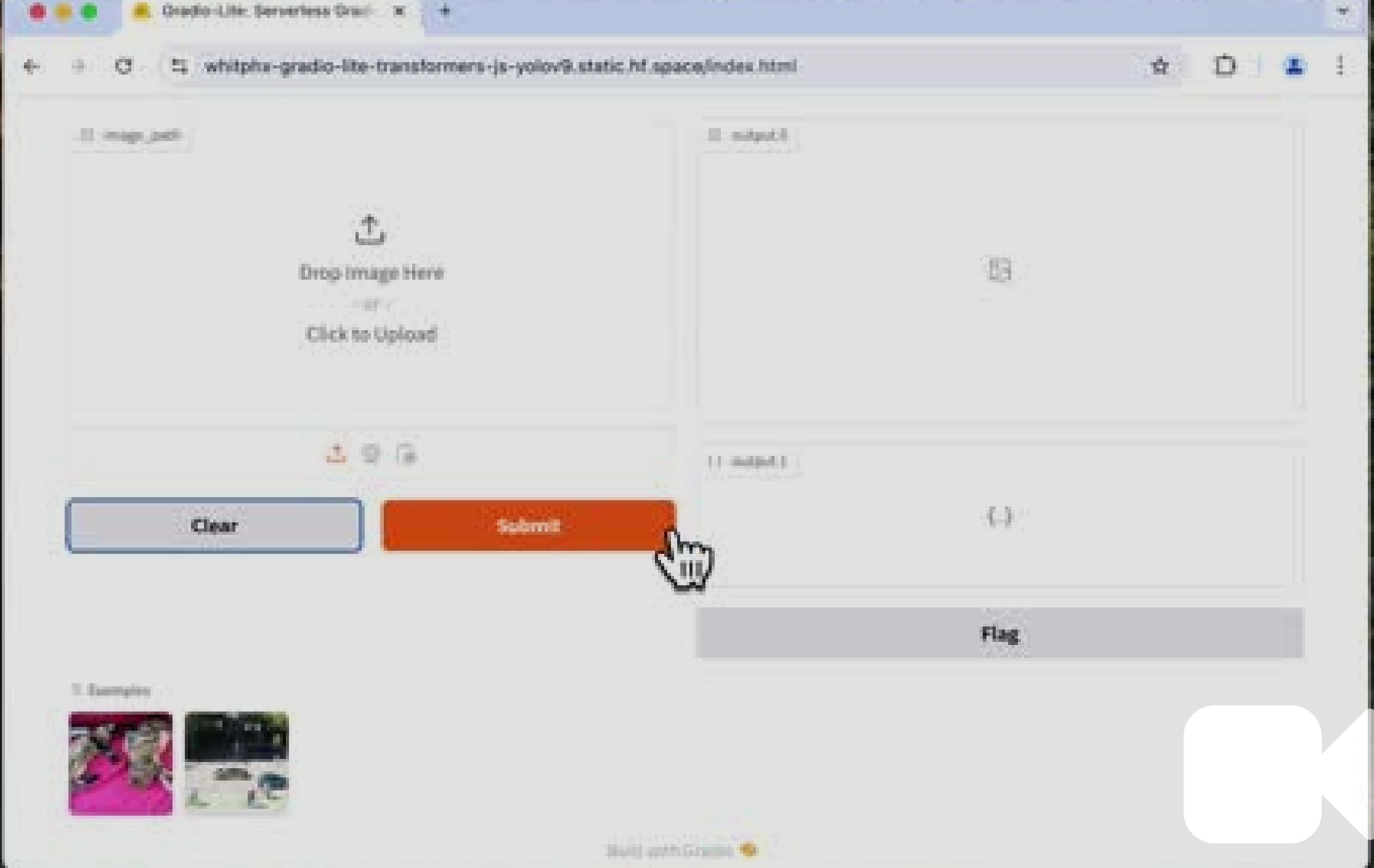


More examples...

Image: Object Detection

● ○ ●

```
1 <!DOCTYPE html>
2 <html>
3     <head>
4         <meta charset="utf-8">
5         <meta name="viewport" content="width=device-width, initial-scale=1">
6         <title>Gradio-Lite: Serverless Gradio Running Entirely in Your Browser</title>
7         <meta name="description" content="Gradio-Lite: Serverless Gradio Running Entirely in Your Browser">
8
9         <script type="module" crossorigin src="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.js"></script>
10        <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.css" />
11
12        <style>
13            html, body {
14                margin: 0;
15                padding: 0;
16                height: 100%;
17            }
18        </style>
19    </head>
20    <body>
21        <gradio-lite>
22            <gradio-file name="app.py" entrypoint>
23 from transformers_js import import_transformers_js, as_url
24 import gradio as gr
25
26
27 # Reference: https://huggingface.co/spaces/Xenova/yolov9-web/blob/main/index.js
28
29 IMAGE_SIZE = 256;
30
```

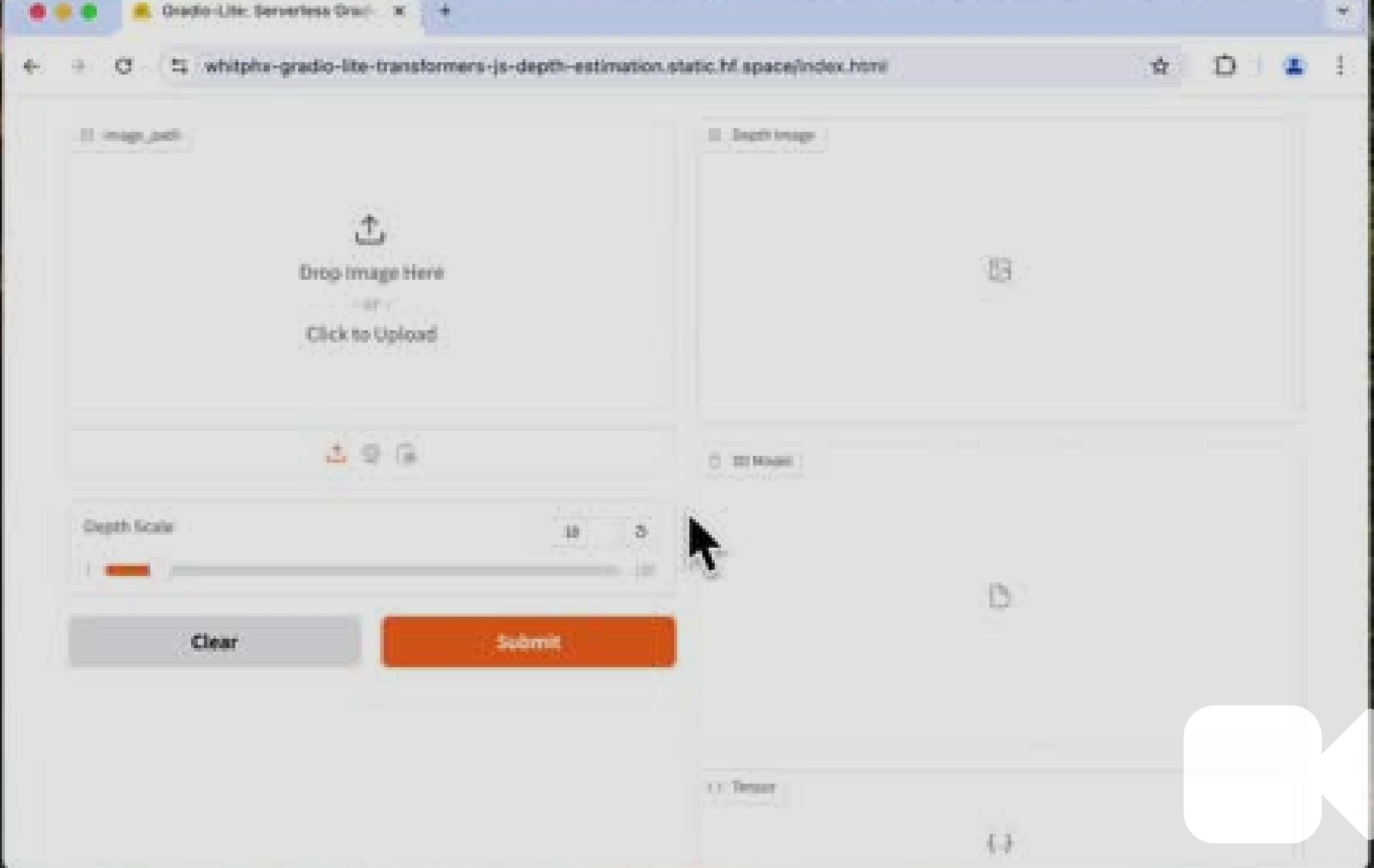


<https://huggingface.co/spaces/whitphx/gradio-lite-transformers-js-yolov9>

Image & 3D: Depth Estimation

● ○ ●

```
1 1 <!DOCTYPE html>
2 2 <html>
3 3     <head>
4 4         <meta charset="utf-8">
5 5         <meta name="viewport" content="width=device-width, initial-scale=1">
6 6         <title>Gradio-Lite: Serverless Gradio Running Entirely in Your Browser</title>
7 7         <meta name="description" content="Gradio-Lite: Serverless Gradio Running Entirely in Your Browser">
8 8
9 9         <script type="module" crossorigin src="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.js"></script>
10 10     <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.css" />
11 11
12 12     <style>
13 13         html, body {
14 14             margin: 0;
15 15             padding: 0;
16 16             height: 100%;
17 17         }
18 18     </style>
19 19 </head>
20 20 <body>
21 21     <gradio-lite>
22 22         <gradio-file name="app.py" entrypoint>
23 23 import gradio as gr
24 24 import numpy as np
25 25 import PIL
26 26 import trimesh
27 27 from transformers_js import import_transformers_js, as_url
28 28
29 29
30 30 transformers = await import_transformers_js()
```



<https://huggingface.co/spaces/whitphx/gradio-lite-transformers-js-depth-estimation>

Zero-shot Image Classification



```
1 <!DOCTYPE html>
2 <html>
3   <head>
4     <script type="module" crossorigin src="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.js"></script>
5     <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.css" />
6   </head>
7   <body>
8     <gradio-lite>
9       <gradio-requirements>
10    transformers_js_py
11  </gradio-requirements>
12
13  <gradio-file name="app.py" entrypoint>
14  from transformers_js import import_transformers_js, as_url
15  import gradio as gr
16
17  transformers = await import_transformers_js()
18  pipeline = transformers.pipeline
19  pipe = await pipeline('zero-shot-image-classification')
20
21  async def classify(image, classes):
22      if not image:
23          return {}
24      classes = [x for c in classes.split(",") if (x := c.strip())]
25      if not classes:
26          return {}
27      data = await pipe(as_url(image), classes)
28      result = {item['label']: round(item['score'], 2) for item in data}
29      return result
30
```

- Privacy
- Low latency
- Offline capability
- Scalability without servers
- Low cost

<https://huggingface.co/spaces/whitphx/gradio-lite-live-zeroshot-image-classification>

Speech-to-Text



```
1 <!DOCTYPE html>
2 <html>
3     <head>
4         <meta charset="utf-8">
5         <meta name="viewport" content="width=device-width, initial-scale=1">
6         <title>Gradio-Lite: Serverless Gradio Running Entirely in Your Browser</title>
7         <meta name="description" content="Gradio-Lite: Serverless Gradio Running Entirely in Your Browser">
8
9         <script type="module" crossorigin src="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.js"></script>
10        <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.css" />
11
12        <style>
13            html, body {
14                margin: 0;
15                padding: 0;
16                height: 100%;
17            }
18        </style>
19    </head>
20    <body>
21        <gradio-lite>
22            <gradio-file name="app.py" entrypoint>
23 from transformers_js_py import import_transformers_js, read_audio
24 import gradio as gr
25
26
27 transformers = await import_transformers_js()
28 pipeline = transformers.pipeline
29 pipe = await pipeline('automatic-speech-recognition', 'Xenova/whisper-tiny.en')
30
```



audio_path



Drop Audio Here

- OR -

Click to Upload

output

Flag

Clear

Submit

Example:

file.wav



<https://huggingface.co/spaces/whitphx/gradio-lite-transformers-js-whisper>

Text-to-Speech



```
1 <html>
2   <head>
3     <script type="module" crossorigin src="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.js"></script>
4     <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.css" />
5   </head>
6   <body>
7 <gradio-lite>
8
9 <gradio-requirements>
10 transformers_js_py
11 </gradio-requirements>
12
13 <gradio-file name="app.py" entrypoint>
14 from transformers_js_py import import_transformers_js
15 import gradio as gr
16 import numpy as np
17
18 transformers_js = await import_transformers_js("3.0.2")
19 pipeline = transformers_js.pipeline
20
21 synthesizer = await pipeline(
22   'text-to-speech',
23   'Xenova/speecht5_tts',
24   { "quantized": False }
25 )
26 speaker_embeddings = 'https://huggingface.co/datasets/Xenova/transformers.js-docs/resolve/main/speaker_embeddings.bin';
27
28
29 async def synthesize(text):
30   out = await synthesizer(text, { "speaker_embeddings": speaker_embeddings });
```



<https://huggingface.co/spaces/whitphx/gradio-lite-text-to-speech>

Supported tasks/models

Text

Question Answering, Summarization, Text Classification, Text Generation, Text-to-text Generation, Token Classification, Translation, Zero-Shot Classification, Feature Extraction, etc...

Audio

Audio Classification, Audio-to-Audio, Automatic Speech Recognition, Text-to-Speech

Image/Video

Depth Estimation, Image Classification, Image Segmentation, Image-to-Image, Mask Generation, Object Detection, Video Classification, Unconditional Image Generation, Image Feature Extraction

Multimodal

Document Question Answering, Image-to-Text, Text-to-Image, Visual Question Answering, Zero-Shot Audio Classification, Zero-Shot Image Classification, Zero-Shot Object Detection

Deployment/Delivery

Hosting a static page

```
1 <html>
2   <head>
3     <script type="module" crossorigin src="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.js"
4     <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/@gradio/lite/dist/lite.css" />
5   </head>
6   <body>
7     <gradio-lite>
8       import gradio as gr
9       from transformers import pipeline
10      pipe = pipeline('sentiment-analysis')
11
12      async def fn(text):
13        result = await pipe(text)
14        return result
15
16    demo = gr.Interface(
17      fn=fn,
18      inputs="text",
19      outputs="text")
```

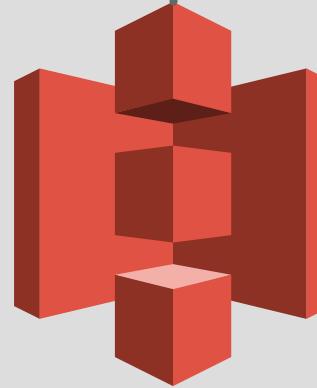
Static HTML file



GitHub Pages

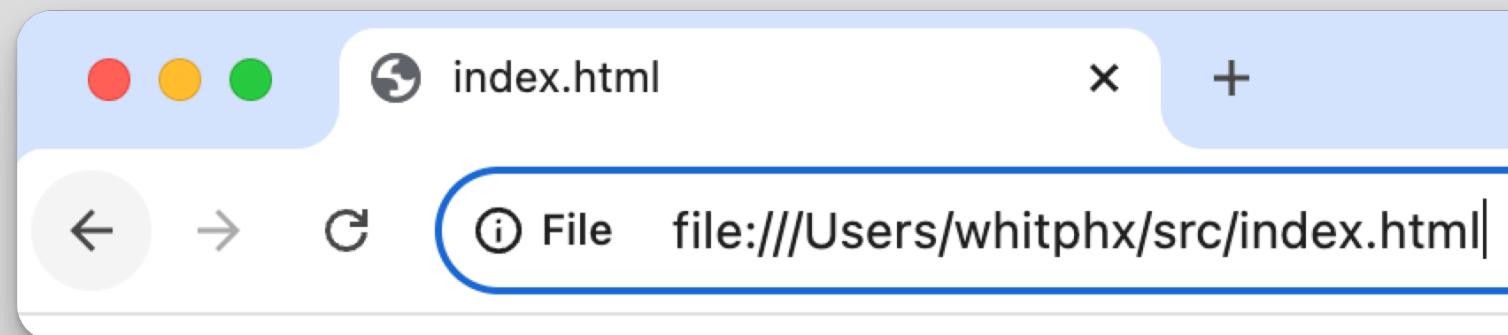


hf.co/spaces



And any your static
web server...

...or even distributing the HTML file to your colleagues 😅



Gradio Playground

gradio Playground [Docs](#) [Gradio Lite](#) [Share Your App](#) [Deploy to Spaces](#)

Demos [←](#) [Code](#) Packages

New Demo

Text

Hello World

Hello Blocks

Sentence Builder

Diff Texts

Media

Sepia Filter

Video Identity

Iterative Output

Generate Tone

Tabular

Filter Records

Transpose Matrix

Tax Calculator

Kinematics

Stock Forecast

Chatbots

Chatbot

Streaming Chatbot

Chatbot with Tools

Multimodal Chatbot

Other

Tabbed Interface

Layouts

Error

Chained Events

import gradio as gr

```
1 import gradio as gr
2
3
4 v def greet(name):
5     return "Hello " + name + "!"
6
7
8 demo = gr.Interface(fn=greet, inputs="textbox", outputs="textbo
9
10 v if __name__ == "__main__":
11     demo.launch()
12
```

Preview 487px [\[x\]](#) [\[y\]](#)

name

Clear Submit

output

Flag

Prompt will update code in editor [CLEAR](#)

BETA

Use cases

Does it replace JS?

No.

In-browser benefits

- Privacy
- Low latency
- Offline capability
- Scalability without servers
- Low cost

Python/Gradio benefits

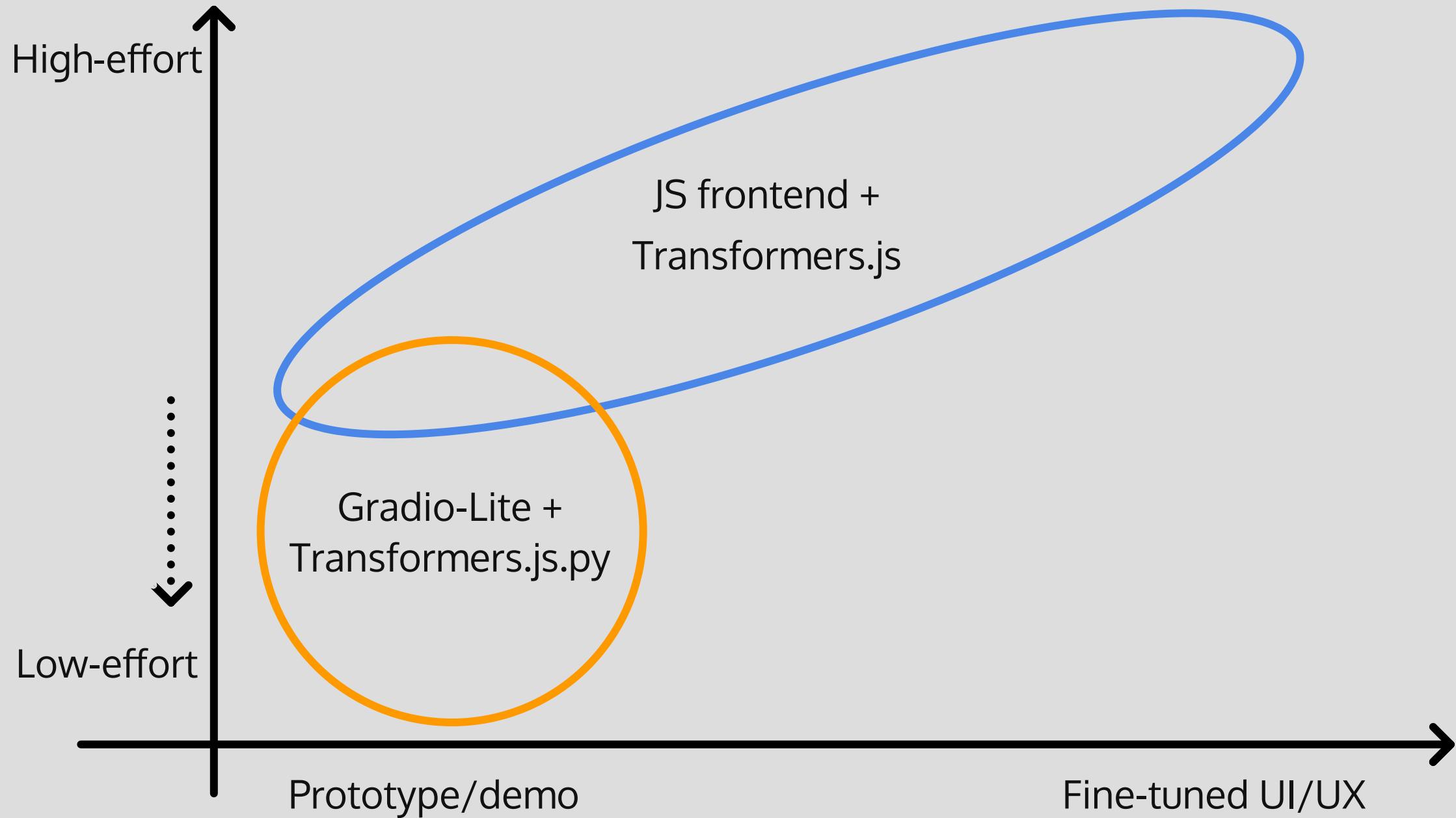
- AI ecosystem
- Built-in components for AI apps
- Easy and quick development

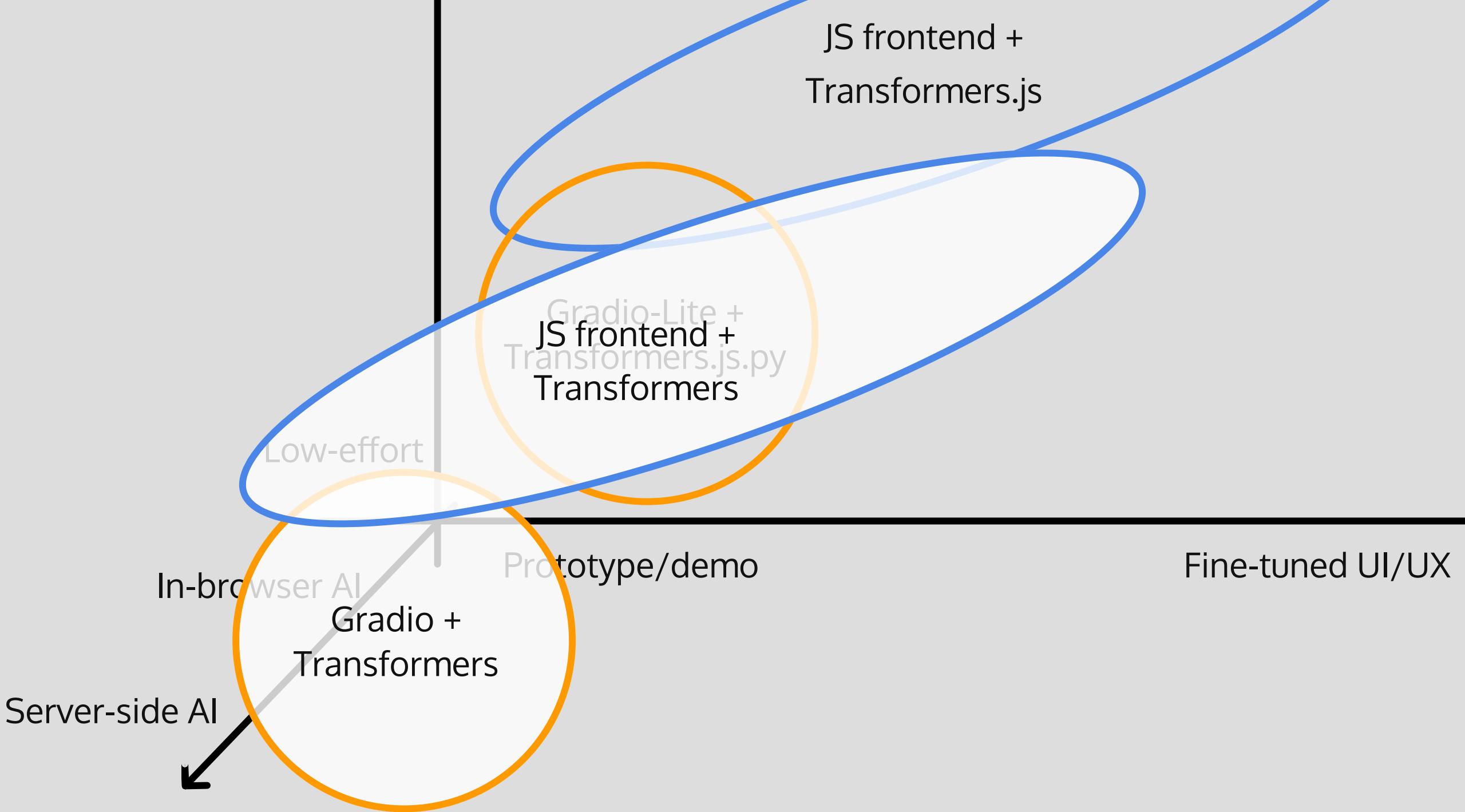
Pros

- Privacy
- Low latency
- Offline capability
- Scalability without servers
- Low cost
- AI ecosystem
- Built-in components for AI apps
- Easy and quick development

Cons

- **Performance:** Initial payload size, loading time, memory usage, etc.
- **Not flexible:** pre-defined components and styles.
- Slight **difference from the normal Python runtime**, e.g. event-loop nature.



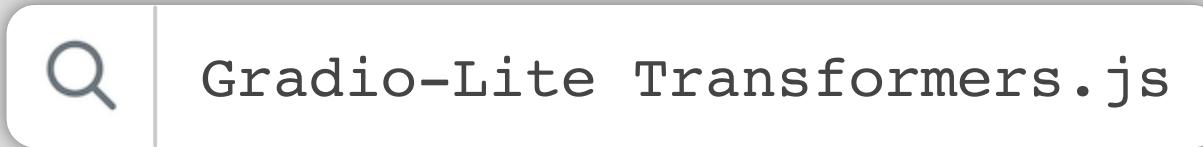


Gradio-Lite + Transformers.js

are good for...

- Early-stage demo/prototypes
- For communications between researchers, developers, managements, clients, etc
- In-office/private tools
- and...

Please try them out!



Gradio-Lite:

Serverless Gradio Running Entirely in Your Browser

<https://huggingface.co/blog/gradio-lite>

**Building Serverless Machine Learning Apps
with Gradio-Lite and Transformers.js**

<https://www.gradio.app/guides/gradio-lite-and-transformers-js>

In-browser LLM app in pure Python:

Gemini Nano + Gradio-Lite



Yuichiro Tachibana



@whitphx

- Pythonista
- OSS enthusiast
- ML Developer Advocate at Hugging Face
- Streamlit Creator