# 🤗 Transformers.js

**State-of-the-art Machine Learning for the web!**

## June 2024 update

https://github.com/xenova/transformers.js

`npm` `v2.17.2`  `downloads` `67k/week`  `jsdelivr` `930k/week`  `license` `Apache-2.0`

Joshua Lochner
joshua@huggingface.co

# 1

# Introduction

What is Transformers.js?

# What is Transformers.js?

## ML + JS
Run ML models directly in the browser with JavaScript!

## Open source
Community-driven development on GitHub. New features added daily!

## Easy to use
Add state-of-the-art ML to your web-app in just a few lines of code!
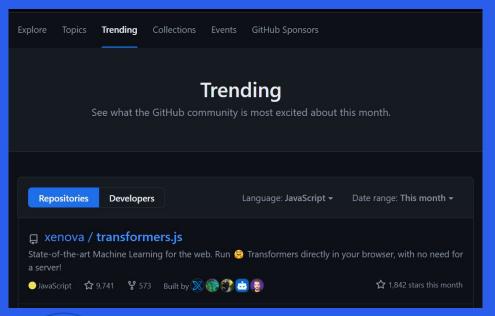
⭐ Stars 9.7k  Forks 573

```
npm i @xenova/transformers
```

We have **over 1000** *ready-to-use* models available on the Hugging Face Hub!

# Community interest

# 2

# WebGPU support

Experimental in Transformers.js (v3)

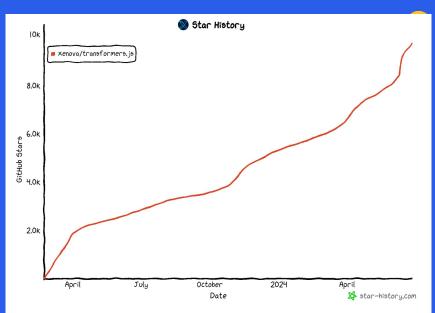# How to use it?

```
npm i xenova/transformers.js#v3

import { pipeline } from '@xenova/transformers';

const extractor = await pipeline(
  'feature-extraction', 'Xenova/all-MiniLM-L6-v2',
  { device: 'webgpu' } // <-- ENABLE WEBGPU
);
const output = await extractor(
  'This is a simple test.',
 { pooling: 'mean', normalize: true }
);
```
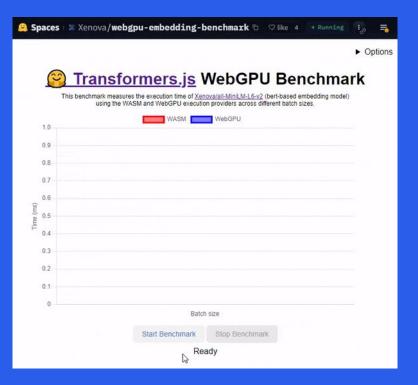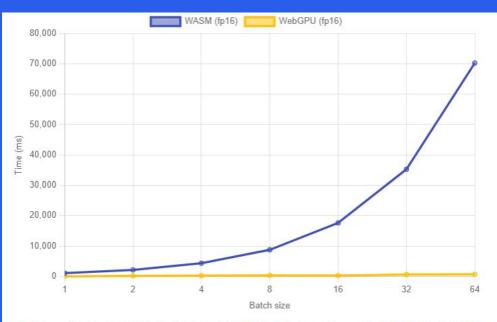
⚠️ EXPERIMENTAL ⚠️

# Huge performance boosts!



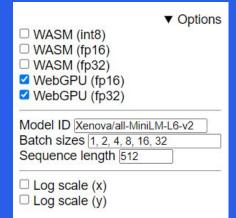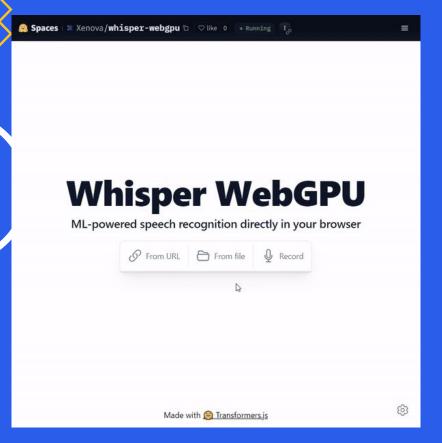https://hf.co/spaces/Xenova/webgpu-embedding-benchmark

# 3

# WebGPU demos

Building applications with Transformers.js (v3)

# Whisper WebGPU

# Phi-3 WebGPU

Spaces · Xenova / experimental-phi3-webgpu_v2 ♡ like 0 · Running ⋮ ☰

Ready!

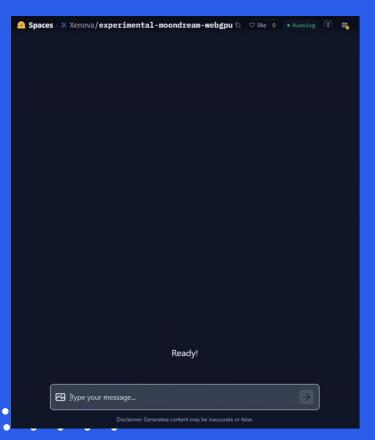Type your message...
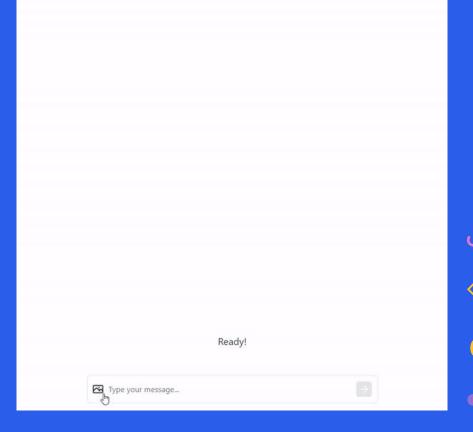
Disclaimer: Generated content may be inaccurate or false.
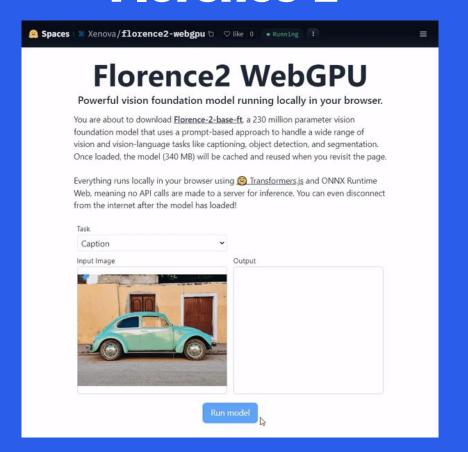
**3.82 billion** parameter LLM that is optimized for inference on the web
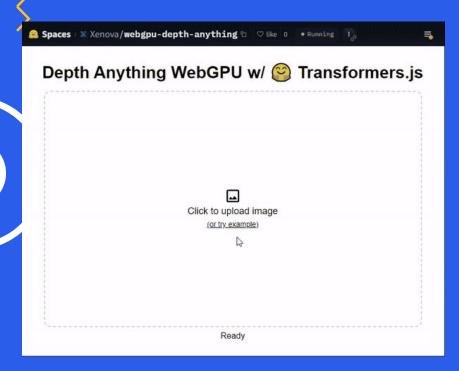
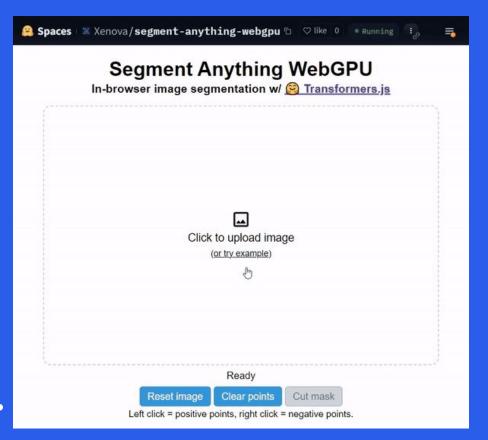# Moondream/LLaVa WebGPU
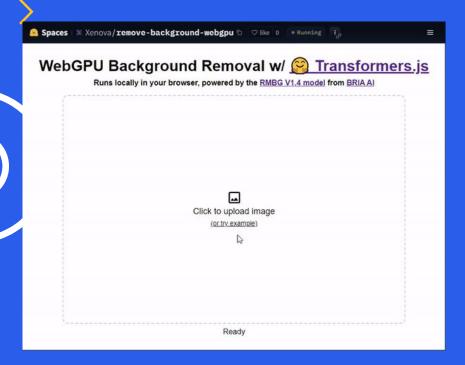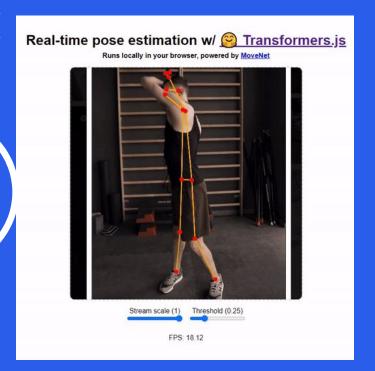
# Florence-2

# Depth Anything

# Segment Anything (SAM)

# Background Removal

# Other tasks



Pose estimation



Zero-shot image classification
(classes defined at runtime)

# Thanks !

Any questions?

@xenovacom

joshua@huggingface.co

@xenova