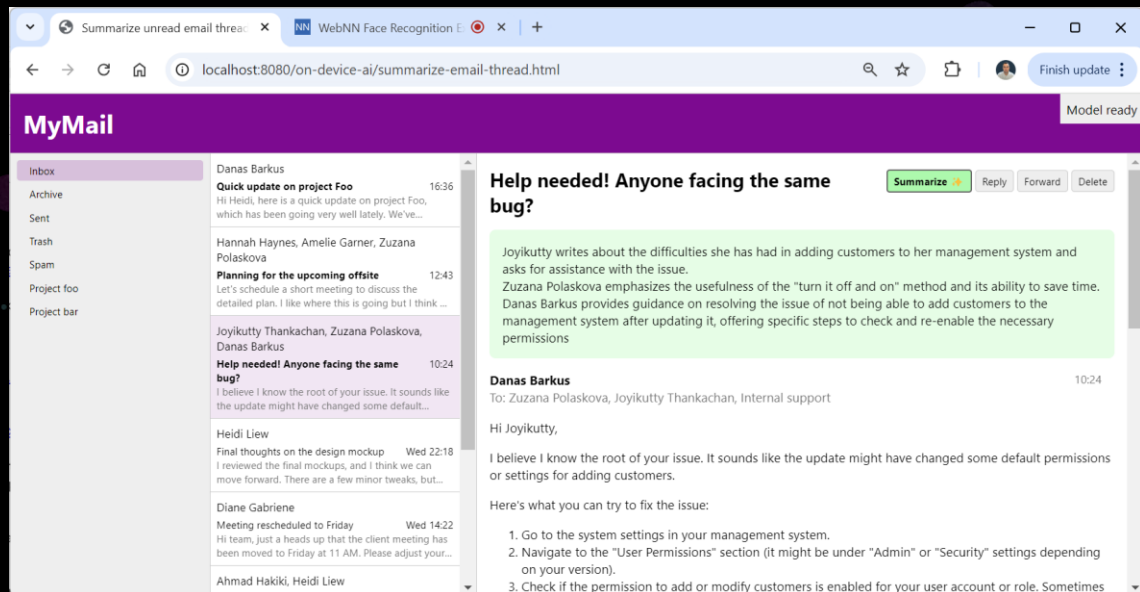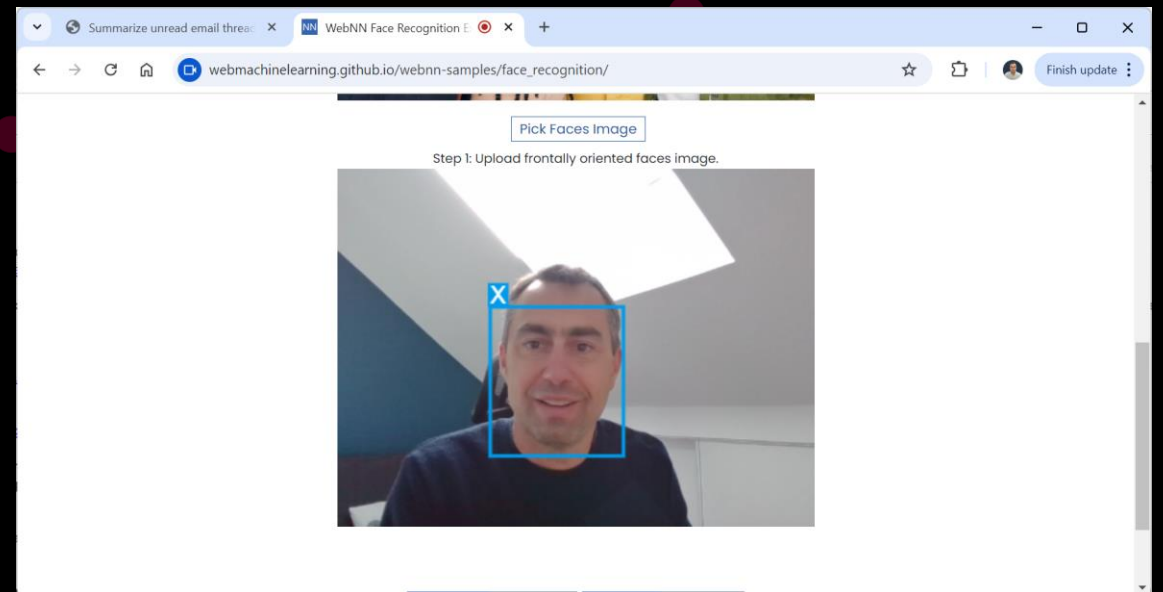# AI on the Web Platform

## Early explorations with on-device AI

Patrick Brosset, Smashing Conference, October 2024

Demo time!

# What if?

- What if this could run on users' devices?

  - No privacy concerns 🔐
  - Lower latency 📶
  - No cloud AI cost 💵
  - Offline support ✈️

# What if?

- What if webpages could use an API to do this directly?

  - No framework to use

  - No model to download

  - Simple usage

# Wait, let's take a step back

- The web and AI

# Wait, let's take a step back

- The web and AI

- The web **platform** and AI?

# It's already happening

"Bring Your Own Model" - BYOM

- WebGL
- WebGPU
- WebNN

- TensorFlow Lite
- ONNX Runtime Web

# It's already happening

Google's "built-in" Prompt() API

- Browser-provided model
- Simple API

- For experimentation only ⚠️

# Our interest 🌊

- We want the web to succeed

- Devs want the web to be at par with native

- We want to understand user's needs first

- Not just Language Models

- We want to make it work with the web

# You can help 🙏

- Ideas? Feedback? Use cases? Thoughts?

- Come find me during the conference

- Message me:
  - Mastodon: @patrickbrosset@mas.to
  - Email: patrick.brosset@microsoft.com

- Anonymous form: aka.ms/web-ai-feedback

# Microsoft Edge

Come chat about:

- On-device AI for the web
  Ideas/use cases? aka.ms/web-ai-feedback

- Progressive Web Apps and PWA Builder

- Web Platform stuff

# On-device AI for the web

What if you could run AI models on a web page?

🔐 Privacy-preserving   📶 Low latency       📢 Ideas/use cases?
💵 No cloud AI cost      ✈️ Offline support   aka.ms/web-ai-feedback

You already can! WebGL, WebGPU, WebNN, Experimental Prompt() API