



ONNX
RUNTIME



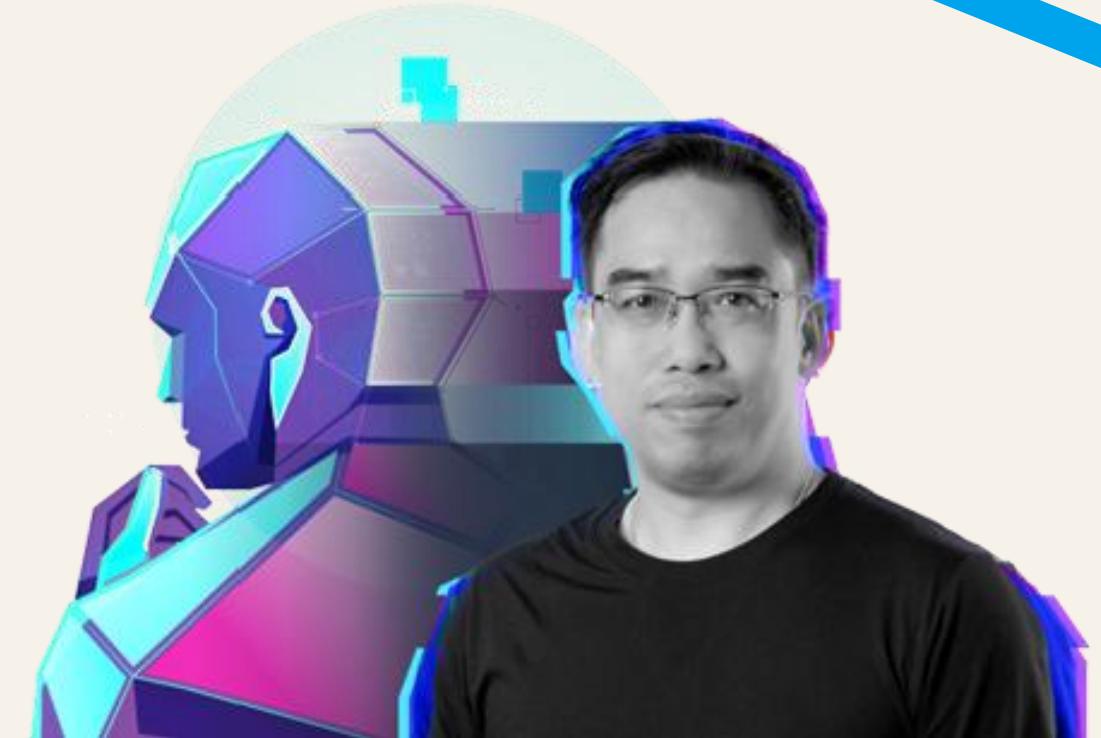
web
neural network

AI web development with **Web Neural Network (WebNN) API**

Enable web app executing AI utilizing CPU, GPU, NPU on client

Surasuk Oakkharaamonphong

Microsoft MVP AI Platform & Developer Technologies
Technical Coach at Arise & INFINITAS by Krungthai





CodeBangkok

@CodeBangkok · 4.48K subscribers · 71 videos

Microsoft Developer Thailand >

[Subscribe](#)

[Join](#)

Go Programming - Basic Syntax 19K views · 2 years ago	Go Programming - GORM 10K views · 2 years ago	Go Programming - Fiber Web Framework 12K views · 2 years ago	Go Programming - Hexagonal Architecture EP.2 (log, error) 7.1K views · 2 years ago	Technology of Trust: SOLIDITY PROGRAMMING SOLIDITY PROGRAMMING Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 2:25:38			
Semantic Kernel - Global AI Bootcamp 2024 330 views · 2 weeks ago	Prompt Engineering - The Art and Science of Querying AI 610 views · 3 weeks ago	AWS re:invent 2023 Recap 206 views · 1 month ago	Orchestrate your AI with Semantic Kernel 391 views · 2 months ago	ARISE CONNEXT - Make a better world with AI 668 views · 5 months ago	Github Copilot Day 1.2K views · 6 months ago	Semantic Kernel - AI Orchestration in Copilot Stack 573 views · 6 months ago	Azure OpenAI Service development with Semantic Kernel 394 views · 6 months ago
Swift Programming - Struct, Class, Protocol, Error Handling Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:08:44	Swift Programming - Basic Syntax Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:38:57	Metaverse Development with Web3.js Three.js VR NFT Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 2:02:42	Build mobile and desktop apps with .NET MAUI .NET MAUI 1:10:55	podman & docker containers technology Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:49:39	ChatGPT OpenAI Technology 1.1K views · 9 months ago	Azure OpenAI Service 622 views · 10 months ago	git Workflow and Branching Strategies Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:49:47
Go Programming - gRPC 16K views · 1 year ago	Go Programming - kafka 13K views · 2 years ago	Solidity Programming: ERC-721 NFT 2.2K views · 2 years ago	Solidity Programming: ERC-20 Token Standard 1.5K views · 2 years ago	git - Distributed Version Control 2.8K views · 1 year ago	Metaverse Development with WebXR, Augmented Reality (AR) 547 views · 1 year ago	React Web Development EP.2 API, Async Test, Mocks 1.4K views · 1 year ago	React Web Development - Get Started 2.9K views · 1 year ago
The Rust Programming Language EP.3 - HTTP Server Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:16:50	The Rust Programming Language EP.2 Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:10:32	Flutter Bloc Design Pattern Update 8.0.0 Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:06:55	The Rust Programming Language Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:21:24	Flutter Input & Form Surasuk Oakkhaaramphong Technical Coach - INFINITAS by Krungthai 1:54:22	Blockchain 4:55	Xamarin Test Cloud 55 views · 6 years ago	Elasticsearch: Get Started 6.4K views · 2 years ago

[Back to sessions](#)

The Web is AI Ready—maximize your AI web development with WebNN

Wednesday, May 22 | 5:23 AM - 5:38 AM Indochina Time Duration 15 minutes

 StudioFP126

[On Demand](#) [Microsoft Build Stage](#) [In Seattle + Online](#)

Speaker:  [Moh Haghigat | Intel Corporation](#)



Resources



Download Video



Download Transcript



View slide deck

Session tags

Session type [Microsoft Build Stage](#)

Topic [AI Development](#)

Topic [Developer Tools](#)

Level [Advanced \(300\)](#)

Delivery Type [In Seattle](#)

Delivery Type [Online](#)

Tag [AI Infrastructure](#)

Tag [Client Development](#)

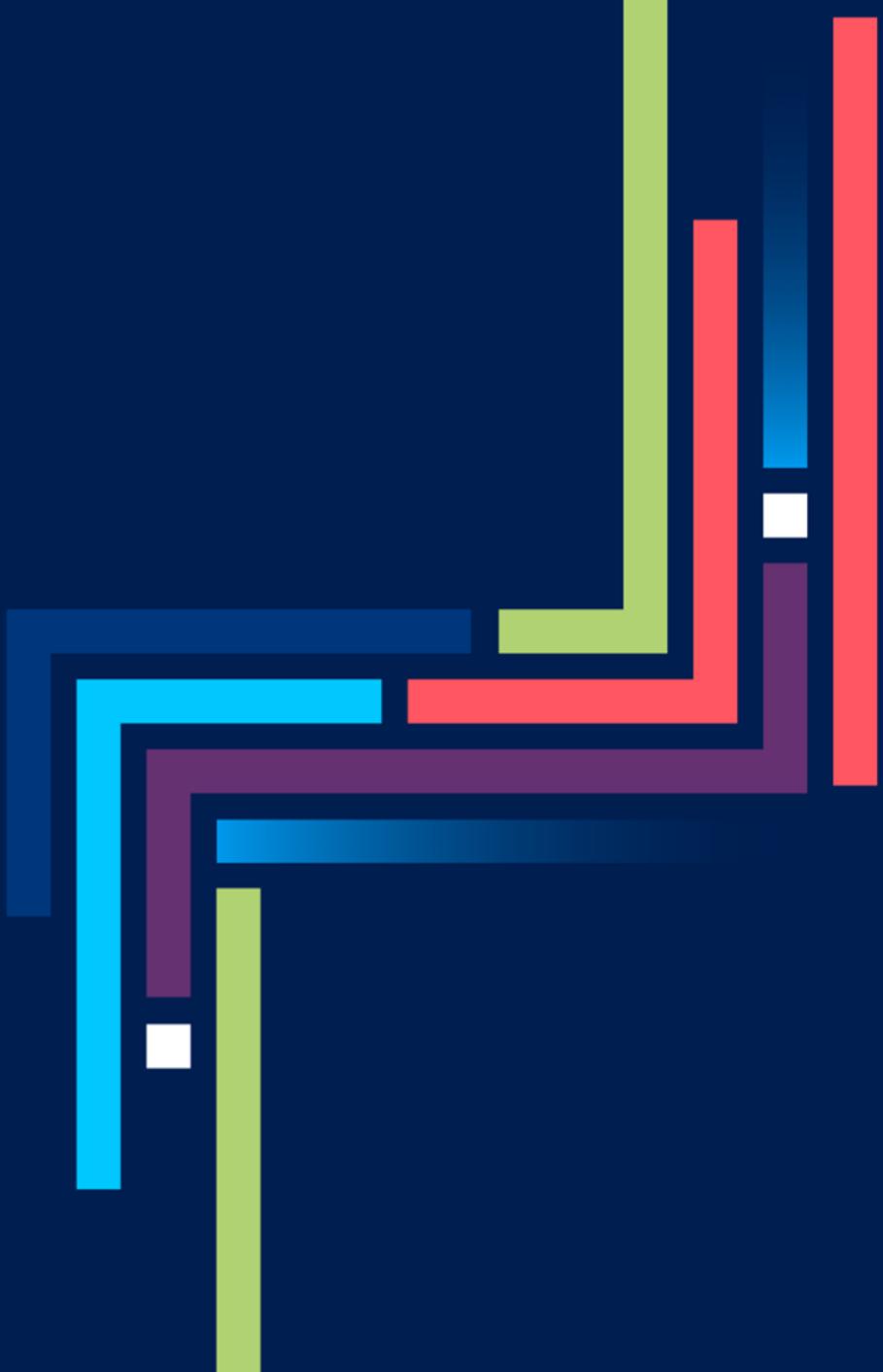
Recording Availability [Recorded](#)

intel.[®] Innovation

Developing Web-Based AI Apps for the Client

An overview of emerging web capabilities to enable exciting AI experiences on the open web platform

Moh Haghigat, Chai Chaoweeraprasit





Moh Haghigheh

Intel Fellow, Chief Web Technologist

Location:
Santa Clara, California

Education:
Ph.D., Computer Science, University of Illinois at Urbana-Champaign

Experience:
28 years as Software Architect at Intel

Areas of interest:
Web Technologies, Programming Models, Compilers, Parallelism



Chai Chaoweeraprasit

Partner Engineering Manager, Microsoft

Location:

Mountain View, California

Education:

Bachelor of Computer Engineering

Experience:

Two decades of core operating system engineering

Areas of interest:

Artificial Intelligence, Deep Neural Network, Web Standard Development, Core Operating System, Computer Graphics, GPU Hardware Acceleration, Digital Typography, Software Internationalization and Unicode, Coding Education

เสวนาดอดหัสเอไอ ฝ่านมุมมอง F1

ลงทะเบียนฟรี! SAT 6 July 13:00 at Microsoft Thailand



Dr. Komes C
Panelist



Big Pongrapree
Moderator

F1 เอ็นจิเนียร์คนไทยที่คุณติ่มแพลตฟอร์ม AI ของไมโครซอฟต์ บอกเล่าประสบการณ์กว่า 20 ปี ของคนไทยที่มีโอกาสได้คุณงานเดพ สำคัญของตน จนมาถึง AI .. โปรเจกต์ล่าสุดที่คุณคือ Phi 3 เอไอ ตัวเล็ก แขนของนาคตของเรื่อย่างไร ชีวิตในซีลคล้อล้อเลย์โอลีกใหม่ จังหวะชีวิตที่กำให้ไปอยู่เมืองรามคำแหงฯ ไทยคืออะไร น่าร่วบพูดคุย เสนอสไลด์จัดการแฟคุยคับตัวเป็นๆกันได้



จัดโดยชุมชนผู้ใช้อาร์ทificial intelligence และพันธมิตร



DaVinci
สถาบันพัฒนาปริญญาดิจิทัล

KBTG

Microsoft

AI



รัน AI บนเว็บเบราเซอร์

30 สิงหาคม 2024



Chai Chaoweeraprasit

การจะเอา AI มา.rันในเว็บแอพต้องทำในเว็บเบราเซอร์ ถ้าใครเคยลองจะรู้ว่ามันพอทำได้ แต่ช้ามากๆ อันนี้พูดถึง รัน AI บนเครื่องโดยตรงนะครับ ไม่ผ่านคลาวด์

เว็บแอพใช้ JavaScript เป็นหลัก สามารถเข้าถึง GPU ได้ ผ่านทาง **WebGPU** มาตรฐานใหม่ที่มาแทน **WebGL** ทั้งสองตัวนี้เป็น Web API ซึ่งเป็นสเปกที่รับรองโดย W3C

ตัว **WebGPU** สามารถใช้รัน AI บน GPU สปีดพอใช้ได้ แต่ยังไม่ดีมาก และยังค่อนข้างลำกัดโดยเฉพาะกับโมเดลแบบ Transformers



Chai Chaoweeraprasit

AI เป็น Workload ที่ค่อนข้างเฉพาะทาง จะเร็วได้จริงๆ ต้องรันแบบเนทีฟเท่านั้น ที่ว่าเนทีฟนี่หมายถึงรัน AI โดยเรียกใช้ฟังก์ชันเฉพาะทางที่อยู่ในตัว OS คือถ้าบน Windows ก็ต้อง DirectML หรือถ้าบน MacOS ก็ CoreML

เลยมาคิดว่าคงต้องสร้างมาตรฐานบนเว็บตัวใหม่ ที่เป็น Web API แบบ **WebGPU** แต่ใช้สำหรับ AI โดยเฉพาะ ใจซองให้เว็บแอปสามารถต่อตรงเข้าถึงฟังก์ชันเนทีฟในตัว OS ไม่ผ่านตัวกลาง

บังเอิญเพื่อนอีกคนซึ่งเป็นวิศวกรอยู่ที่ Intel คิดตรงกัน เลยช่วยกันร่างมาตรฐานใหม่ ใช้ชื่อว่า **WebNN** ชื่อเต็มว่า Web Neural Network API มี W3C หนุนหลัง แบ่งงานกันคือผมเขียนสเปคส่วนเข้าเป็นคนทดสอบ เริ่มทำกันจริงลังช่วงโควิดล็อกดาวน์เมื่อหลายปีที่ผ่านมา



Chai Chaoweeraprasit

จุดเด่นของ WebNN วิถอย่างหนึ่งคือไม่ขึ้นกับฮาร์ดแวร์ แอพสามารถเลือกได้ว่าจะรันบนซีพียู GPU หรือแม็ตเตอร์ NPU และไม่ขึ้นกับฟอร์แมตหรือเฟรมเวิร์ก จะใช้กับโมเดลแบบ ONNX หรือ TFLite ก็ได้ตามใจ

เมื่อหลายเดือนก่อนสัมภาษณ์ล่าวเปิดตัว WebNN Developer Preview ในงานเดฟประจำปีของ Microsoft สร้างความฮือฮาพอสมควรโดยเฉพาะกับเหล่าเดฟที่ทำแอพทั้งหลาย

ตอนนี้ทีมแพลตฟอร์มเราเพิ่งออกอวัพเดตตัวใหม่สำหรับรองรับ NPU ทั้งจากค่าย Intel และค่าย ARM บนชิป Snapdragon X ของ Qualcomm รัน AI แบบเน็ตบัน เว็บเบราเซอร์ผ่าน WebNN ที่ช่วยยุบแพลตฟอร์ม DirectML สำหรับค่า MacOS คงต้องรอไปอีกสักพัก เพียงเห็นเริ่มๆ ทำกันอยู่

WebNN

Announcing Developer Preview



Web Neural Network API

[W3C Candidate Recommendation Draft, 11 May 2024](#)



▼ More details about this document

This version:

<https://www.w3.org/TR/2024/CRD-webnn-20240511/>

Latest published version:

<https://www.w3.org/TR/webnn/>

Editor's Draft:

<https://webmachinelearning.github.io/webnn/>

Previous Versions:

<https://www.w3.org/TR/2024/CRD-webnn-20240510/>

History:

<https://www.w3.org/standards/history/webnn/>

Implementation Report:

<https://wpt.fyi/results/webnn?label=master&label=experimental&aligned&q=webnn>

Test Suite:

<https://github.com/web-platform-tests/wpt/tree/master/webnn>

Feedback:

[GitHub](#)

[Inline In Spec](#)

Editors:

Ningxin Hu ([Intel Corporation](#))

Dwayne Robinson ([Microsoft Corporation](#))

Former Editor:

Chai Chaoweeraprasit ([Microsoft Corporation](#))

Explainer:

[explainer.md](#)

Polyfill:

[webnn-polyfill](#) / [webnn-samples](#)



web
neural network

Standard
W3C API

Unified Abstraction

Hetero
HW Exec

CPU, GPU, NPU

Integrated
ML Frameworks

ONNX RT Web, ...

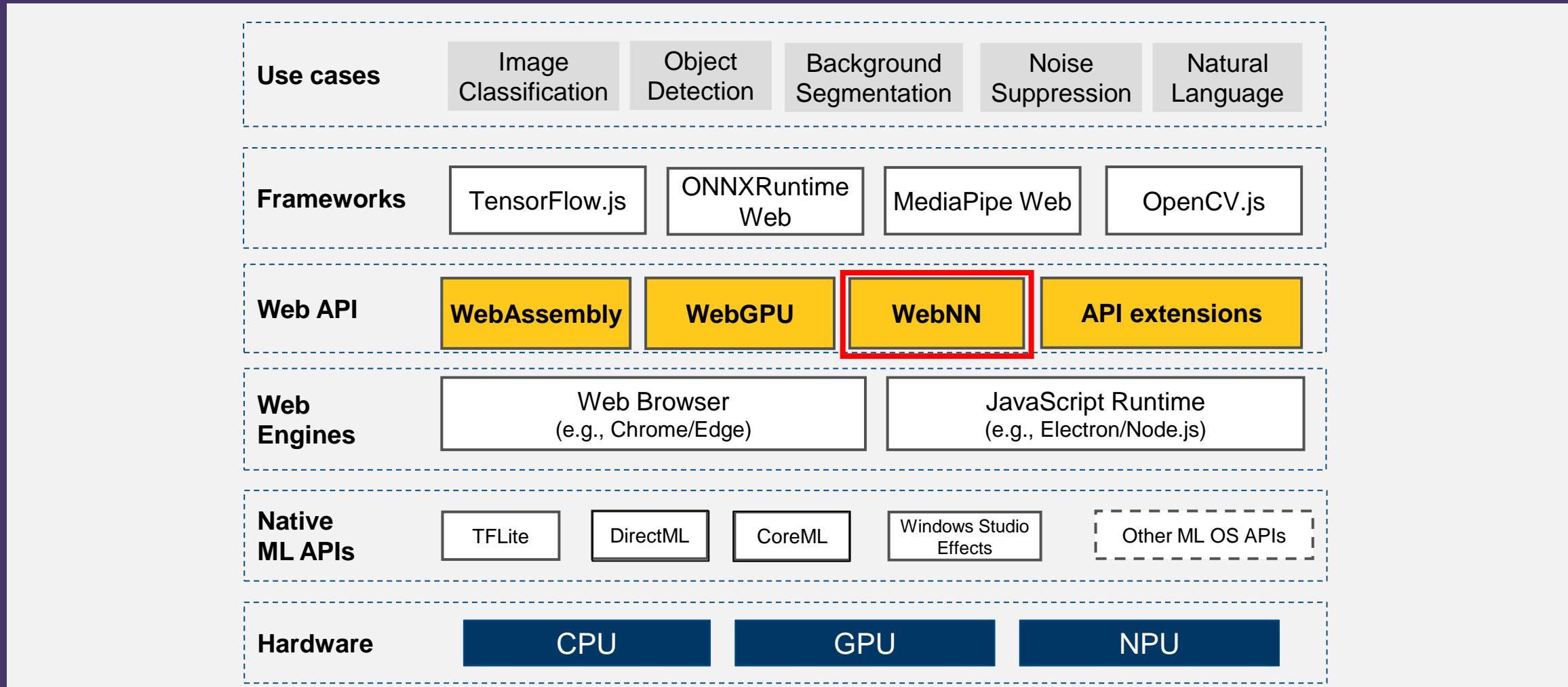
Near Native
Execution
Characteristics

Perf. & Power

General
Computational
Graph

BYOM

Hardware-Accelerated Web AI Overview



Updated

ONNX Runtime Web

Inference AI models in the web browser

Accelerates in-browser inferencing across CPUs, GPUs and NPUs, leveraging WebAssembly, WebGPU and WebNN

Supports generative AI models including Stable Diffusion, TinyLlama, Phi and more

Adopted widely by both 1Ps and 3Ps including Hugging Face Transformer.js

GA + update

ONNX Runtime Mobile

Inference AI models on mobile and embedded devices

A comprehensive solution offering minimized model and runtime sizes for on-device inferencing

Accelerates inference on Android and iOS with MLAS, NNAPI, CoreML, XNNPACK accelerators

Java / Kotlin, Objective-C, JavaScript, C#, and C++ APIs

Updated

ONNX Runtime Generate API

Inference Generative AI models in the cloud or at the edge

Easy-to-use API for cross-platform GenAI applications

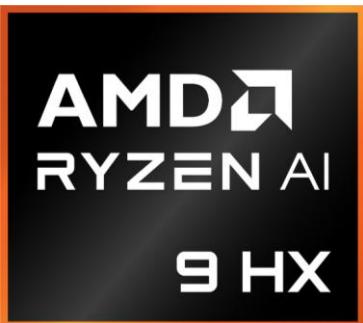
High performance

Python, C#, C / C++ Language bindings. Java / Objective-C
Go coming soon

aka.ms/ortgenapi



Copilot+PC

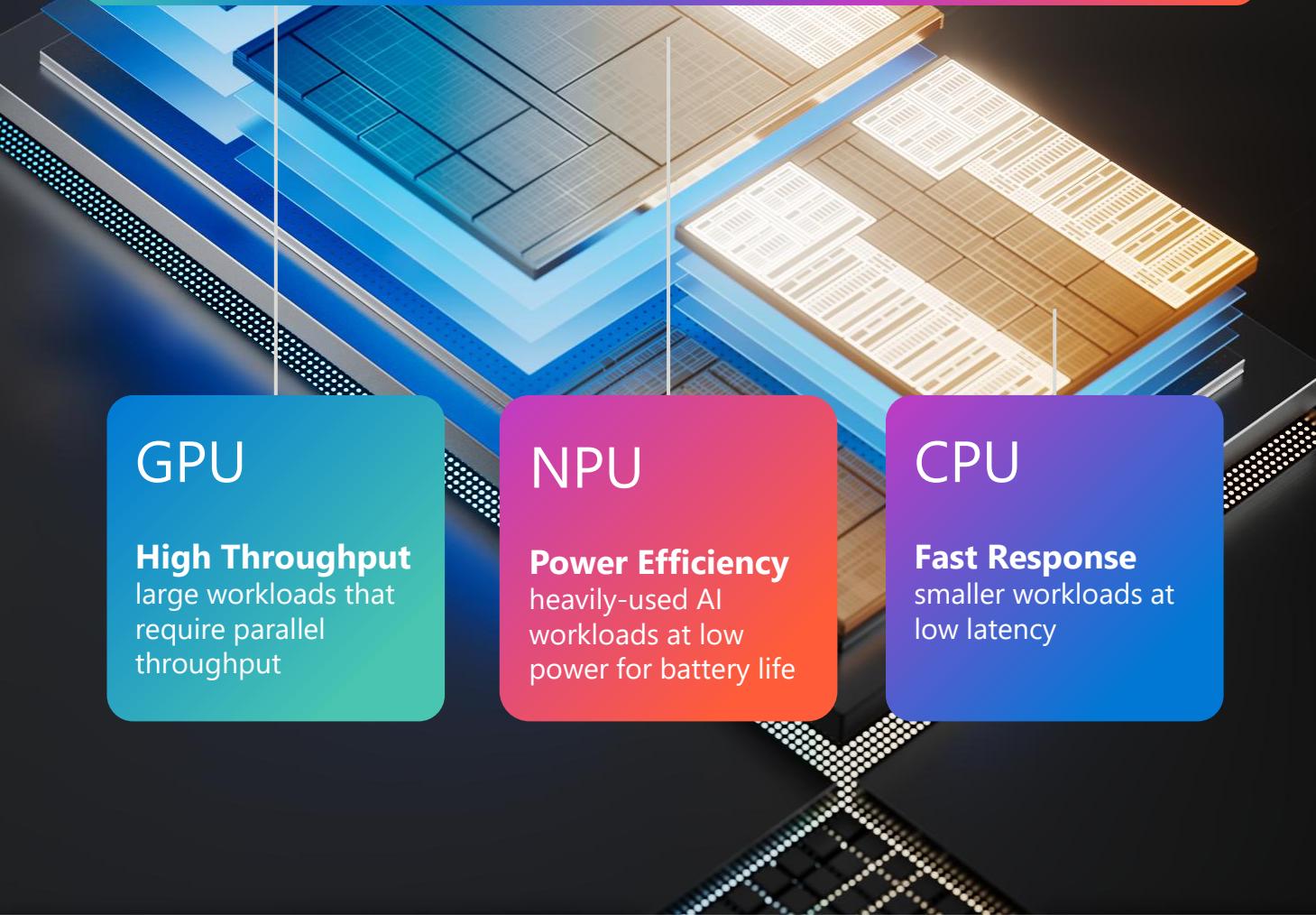


40+ Trillion

Operations per second (TOPS)

What is an AI PC?

A PC with new NPU that brings new AI experiences in productivity, creativity, and security through a combination of the CPU, GPU, and the NPU.



GPU

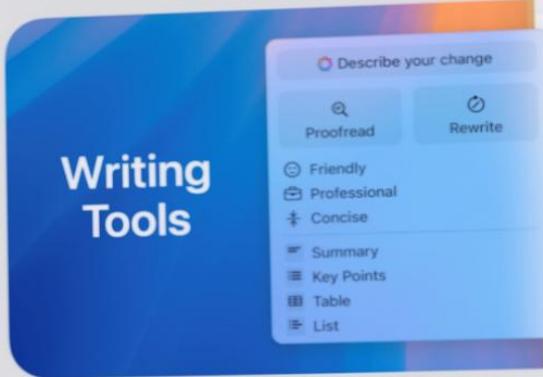
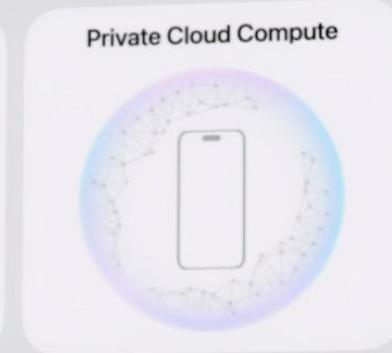
High Throughput
large workloads that require parallel throughput

NPU

Power Efficiency
heavily-used AI workloads at low power for battery life

CPU

Fast Response
smaller workloads at low latency



To:
Cc:
Subject:
From:

Dear Ms. H
It was great to my heart cover letter

Thanks,
Jenny Frith
Dept. of Jo

Reduce Interruptions

in Focus



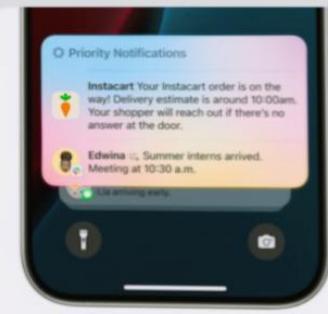
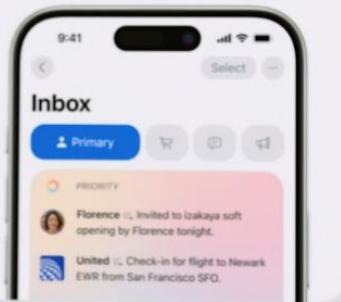
Genmoji

Create a Memory Movie

Describe a Memory...



Priority messages in Mail



Priority notifications

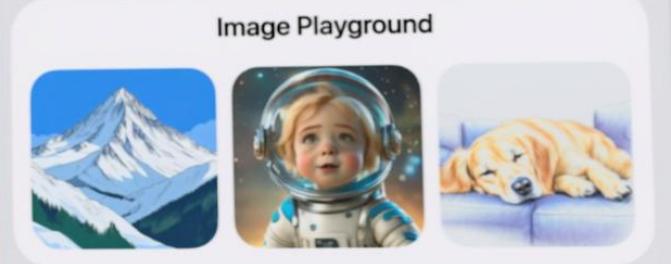
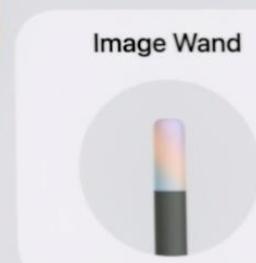


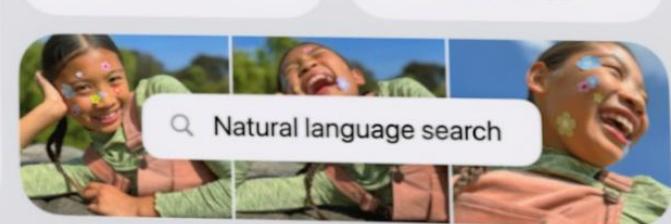
Image Wand



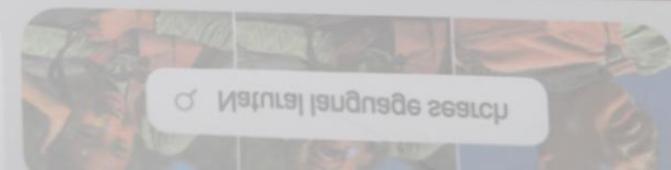
Audio recording



summaries



Natural language search



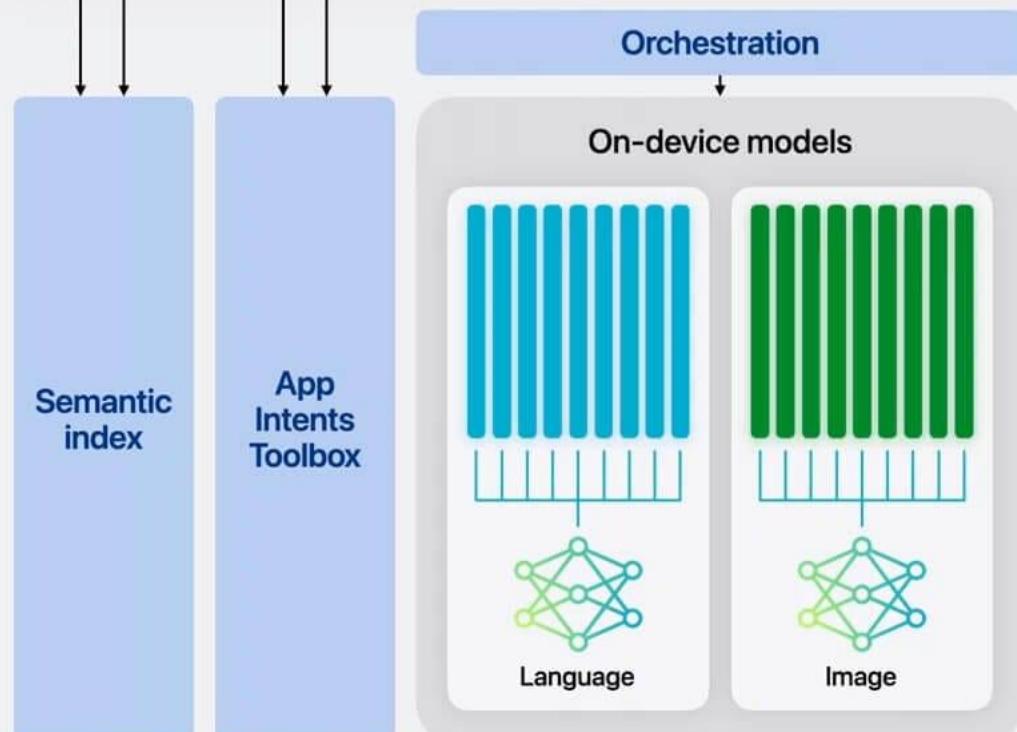
Searches based on natural language

Apple Intelligence

Apps and experiences



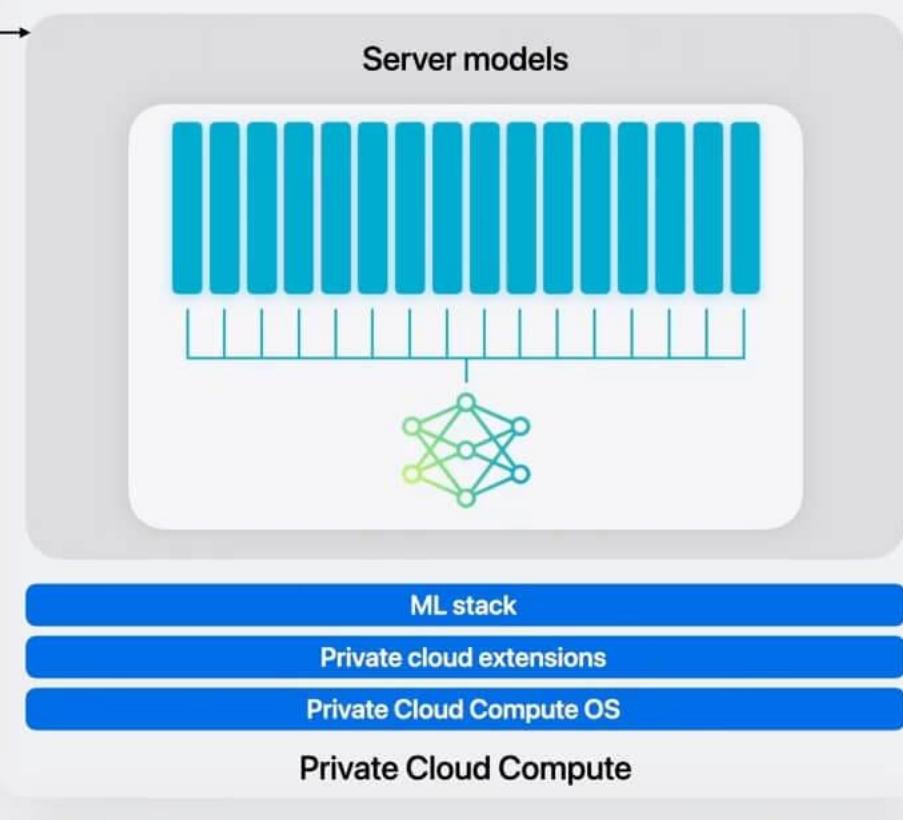
Personal Intelligence System



Apple silicon



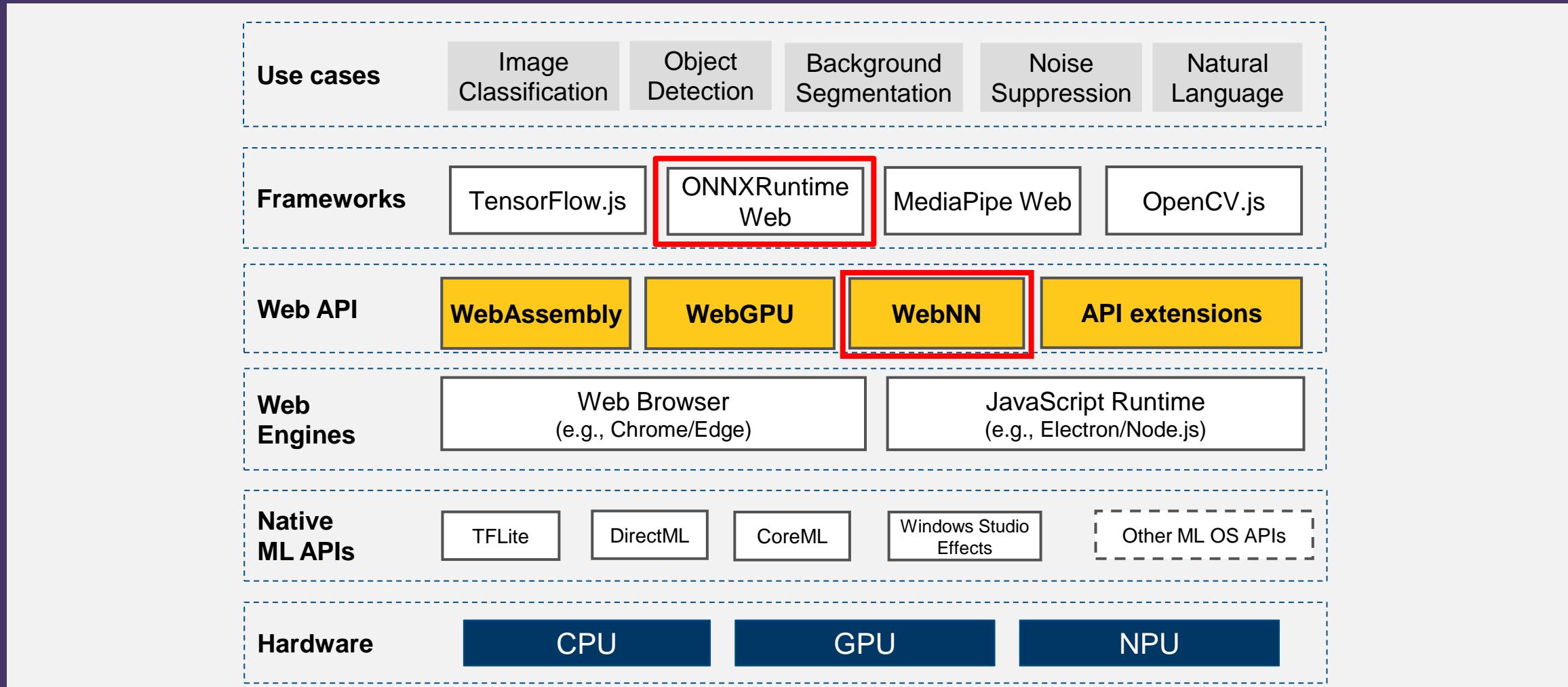
Server models



Apple silicon servers



Hardware-Accelerated Web AI Overview





ONNX - Open Neural Network Exchange

The open standard for machine learning interoperability



ONNX



Common format for ML Models



Can be exported from PyTorch, Tensor Flow, etc



Cross-platform

Desktop | IoT | Mobile | Cloud



Fast!



Supports DirectML, a cross-hardware
ML acceleration API

The ONNX Trilogy



Open and interoperable file format for ML and DNN models.



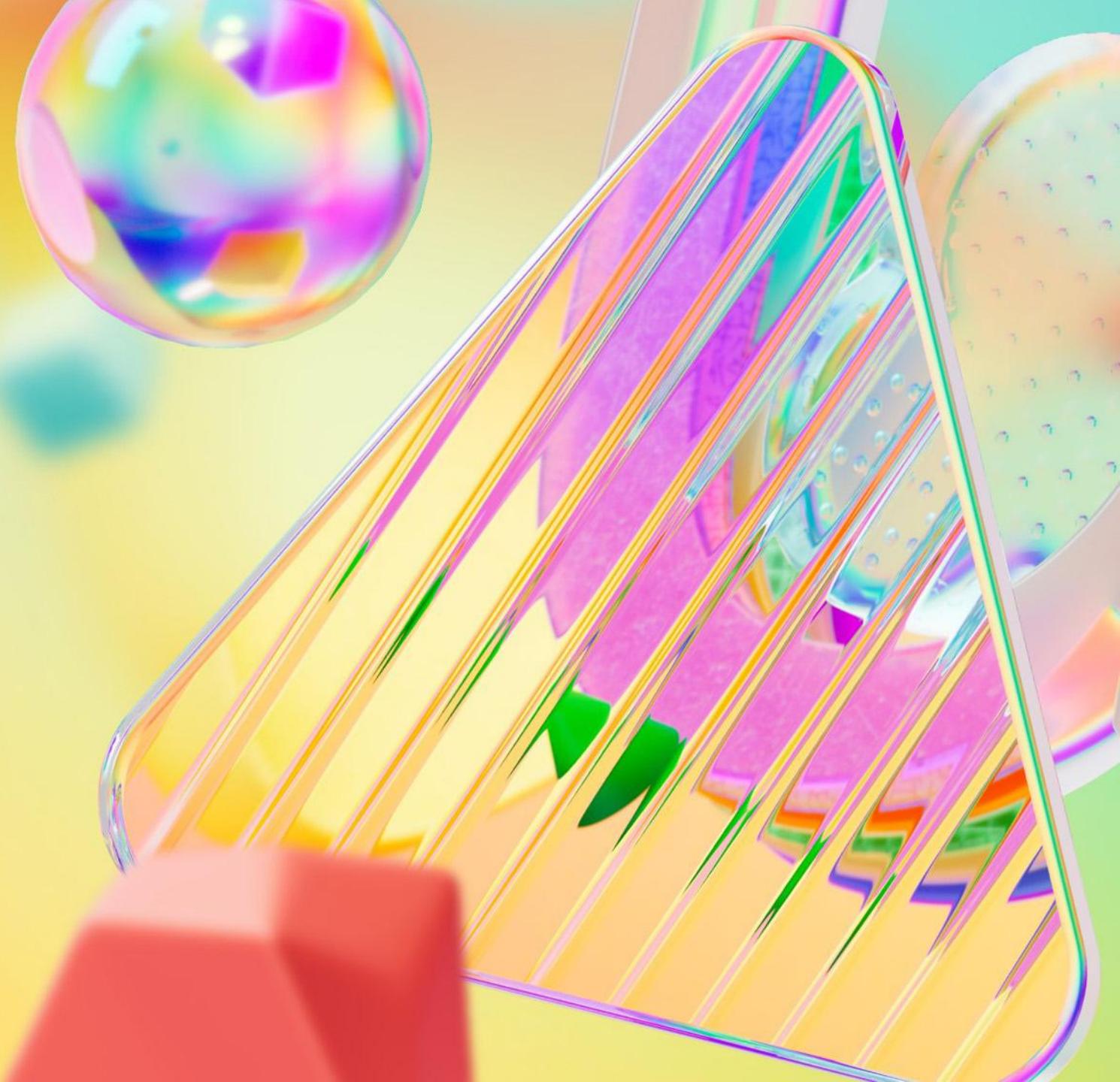
Fast and efficient *model inference and training engine* that works across a diverse range of hardware accelerators.



A toolkit for *hardware-aware AI model optimization*. Output is an ONNX formatted model to run with quality and efficiency on the ONNX runtime.



Demo



AI Models

- [Stable Diffusion 1.5](#) is a latent text-to-image diffusion model capable of generating photo-realistic images given any text input.
- [SD-Turbo](#) is a fast generative text-to-image model.
- [Segment Anything](#) is a new AI model from Meta AI that can "cut out" any object. You can segment any object from your uploaded images.
- [Whisper Base](#) is a pre-trained model for automatic speech recognition (ASR) and speech translation.
- Image Classification [MobileNet](#) and [ResNet](#) models. Take images as input and classify the major object in the image into a set of pre-defined classes.

References

- **WebNN Spec:** <https://www.w3.org/TR/webnn/>
- **WebNN Explainer:** <https://github.com/webmachinelearning/webnn/blob/main/explainer.md>
- **WebNN Implementation Status:** <https://webmachinelearning.github.io/webnn-status/>
- **Awesome WebNN:** <https://github.com/webmachinelearning/awesome-webnn>
- **WebNN Samples:** <https://microsoft.github.io/webnn-developer-preview/> & <https://webmachinelearning.github.io/webnn-samples/>
- **WebNN Image Classification:** https://webmachinelearning.github.io/webnn-samples/image_classification/
- **WebNN Semantic Segmentation:** https://webmachinelearning.github.io/webnn-samples/semantic_segmentation/index.html
- **ONNX Runtime WebNN Execution Provider:**
<https://github.com/microsoft/ondxruntime/tree/main/ondxruntime/core/providers/webnn>
- **Developing Web-Based AI Apps for the Client:**
https://static.rainfocus.com/intel/innv2023/sess/1689805601141001Xyjy/supmat/DevelopingWebBasedAIApps_Final_1694887505835001jjpL.pdf