intel[®]Ai summit

Bringing AI Everywhere

Web Neural Network API

Accelerating Deep Neural Networks on the Web

Jamie Chang, CCE/WED Ningxin Hu, SSE/WPE March 27th, 2024



Web Neural Network (WebNN)

- A new web standard that allows web apps and frameworks to accelerate deep neural networks with ondevice hardware such as GPUs, CPUs, or purpose-built AI accelerators.
- An emerging W3C Web standard API for AI acceleration, Intel and Microsoft spec co-editors.
- Brings a unified abstraction of neural networks to the web
- Provides AI hardware (XPU) access through native machine learning API of the operating system
- Delivers near-native performance and reliability of results while delivering the next gen use cases
- Status:
 - Spec:
 - CNN/RNN Candidate Recommendation published
 - Transformers Q4'23
 - Implementation:
 - [XNNPACK/CPU for Win/Linux]: Dev Trial is available in Chrome and Edge Canary
 - [DirectML/GPU for Win]: Initial support landed upstream, Dev Trial in Q4'23
 - [DirectML/NPU for Win]: work in progress
 - [MLService/NPU for CrOS]: POC done, patches upstreaming



XPU-Accelerated Web AI Overview

Use cases	Image Classification	Object Detection	Backo Segme	ground entation	Noise Suppress	Natural ion Language
Frameworks	TensorFlow.js ONNXRuntime Web		MediaPipe Web		OpenCV.js	
Web API	WebAssembly	WebGP	U	WebNN	A	PI extensions
	Web Browser (e.g., Chrome/Edge)		JavaScript Runtime (e.g., Electron/Node.js)			
Web Engines	Web (e.g., Cl	Browser hrome/Edge)			JavaScrip (e.g., Electro	t Runtime on/Node.js)
Web Engines Native ML APIs	Web (e.g., Cr XNNPACK	Browser hrome/Edge) Service	rectML	Windows Effec	JavaScrip (e.g., Electro Studio	t Runtime on/Node.js) Other ML OS APIs



Almost 2/3rd of PC cycles worldwide are spent on the web

Emerging Web APIs aim to bring AI to the masses across all platforms!

Application Usage by Category



Browsers are top apps on PC*





Web is the best platform to scale AI delivery



State of Web ML

TensorFlow.js usage alone grown >10x since 2020, now standing at >230K weekly downloads, and >15 million cumulative downloads.







Heterogenous Hardware for AI





Agenda WebNN Delivering Al experiences across CPU, GPU, NPU

WebNN API & Spec

Browser Implementations

Performance & Demo

Call to Action

References



WebNN Specification

TABLE OF CONTENTS

1 Introduction

Draft

G

å

lidate

W3C

∠\5

intel

2

Use cases

- 2.1 Application Use Cases
- 2.1.1 Person Detection
- 2.1.2 Semantic Segmentation
- 2.1.3 Skeleton Detection
- 2.1.4 Face Recognition
- 2.1.5 Facial Landmark Detection
- 2.1.6 Style Transfer
- 2.1.7 Super Resolution
- 2.1.8 Image Captioning
- 2.1.9 Text-to-image
- 2.1.10 Machine Translation
- 2.1.11 Emotion Analysis
- 2.1.12 Video Summarization
- 2.1.13 Noise Suppression
- 2.1.14 Speech Recognition
- 2.1.15 Text Generation
- 2.1.16 Detecting fake video
- 2.2 Framework Use Cases
- 2.2.1 Custom Layer
- 2.2.2 Network Concatenation
- 2.2.3 Performance Adaptation
- 2.2.4 Operation Level Execution
- 2.2.5 Integration with real-time video processing
- 3 Security Considerations
- 3.1 Guidelines for new operations
- 4 Privacy Considerations
 - Ethical Considerations



W3C Candidate Recommendation Draft, 22 March 2024



▼ More details about this document
This version: https://www.w3.org/TR/2024/CRD-webnn-20240322/
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-20240321/
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 (<u>Microsoft Corporation</u>)
Former Editor: Chai Chaoweeraprasit (<u>Microsoft Corporation</u>)
Explainer: explainer.md

webnn-polyfill / webnn-samples

Polyfill:

WebNN is an emerging W3C Web standard API for AI acceleration

WebNN Programming Model



*WebNN next features

WebNN brings a unified abstraction of neural networks to Web



Hello Tensors



axb=c



WebNN Browser Implementations



WebNN JS ML Frameworks Integration



Mainstream ML frameworks are integrating WebNN

Available

Lite



MediaPipe Models Performance using WebNN CPU



Browser: Chrome Canary 118.0.5943.0

• DUT: Dell/Linux/i7-1260P, single p-core

Workloads: <u>MediaPipe solution models</u> (FP32, batch=1)

WebNN delivers near-native inference performance on CPU



ONNX Runtime DirectML Architecture



NPU

GPU

- DXIL: DirectX Intermediate Language
- MCDM: <u>Microsoft Compute Driver Model</u>
- Optimized Kernels: Driver implementation of specific custom

ML operations



DirectML Use Cases

USER: Can you describe the city of Seattle in the summertime in the style of Comedian Trevor Noah?

Dolly:

Sure! Seattle in the summertime is all about the outdoors. Every morning when you wake up the weather report says it will be in the 60's. And for the first half of the day you just feel grateful that it is not raining. Then you go outside and it is 106 degrees. Seattleites all have a number they call "the Noah Index". 1000 is freezing, 2000 is tolerable, 3000 is pleasantly warm, 4000 is a blistering sun burn, 5500 is only slightly less pleasurable than4000, and we're talking about thousands of people there, so the Noah Index gets really big really fast. Seattleites all enjoy the outdoors. We're surrounded by mountains and lakes so we appreciate the great outdoors. We love to go camping and hiking. And when it is hot we love to go to the unnecessary and un-needed Dose of Mountains concert. We've all been to the Dose of Mountains concert, don't get me started on how much we love Dose of Mountains, they are a Olympian institution. And they go on every Tuesday through the end of July. We love our mountains and lakes, and we love being outside. Even in the rain.











Transformers on WebNN DirectML Demo



WebNN transformers demo guides the WebNN next version work



ONNX Runtime Web Code Samples with WebNN

• • •

<pre>import { InterenceSession } from "onnxruntime-web";</pre>		
// Initialize the ONNX model		// Initialize the ON
<pre>const initModel = async () => {</pre>		const initModel = asy
<pre>env.wasm.numThreads = 1; // 4</pre>		env.wasm.numThreads
<pre>env.wasm.simd = true;</pre>		<pre>env.wasm.simd = tru</pre>
<pre>env.wasm.proxy = true;</pre>		<pre>env.wasm.proxy = tr</pre>
<pre>const options: InferenceSession.SessionOptions = {</pre>		const options: Infe
// provider name: wasm, webnn		// provider name:
// deviceType: cpu, gpu, npu		// deviceType: cp
<pre>// powerPreference: default, high-performance</pre>		// powerPreference
executionProviders:		executionProvider
<pre>[{ name: "wasm"}], // WebAssembly CPU</pre>		[{ name: "webnr
}		}
//		//
};		};

•••

import { InferenceSession } from "onnxruntime-web"; // Initialize the ONNX model const initModel = async () => { env.wasm.numThreads = 1; // 4 env.wasm.simd = true; env.wasm.proxy = true; const options: InferenceSession.SessionOptions = { // provider name: wasm, webnn // deviceType: cpu, gpu, npu // powerPreference: default, high-performance executionProviders: [{ name: "webnn", deviceType: "gpu", powerPreference: 'default' }], } // ... }:

WebAssembly backend

WebNN backend

Switching to WebNN can be done by modifying a single line of code



WebNN ChromeOS MLService on NPU

WebNN Image Classification Exa X 🛛 WebNN Image Classification Exa X +			~ _ □ ×
$m{\leftarrow} \rightarrow m{C}$ $$ vbenni.github.io/webnn-samples/image_classification/?numRuns=20	0		🔏 🗖 😫 Guest 🗄 E
web neural network Image Classification Samples •			
Backend Wasm (CPU) WebGL (GPU) WebNN (NPU) WebN Model MobileNet V2 SqueezeNet ResNet V2 50	N (GPU)]	
IMAGE		LIVE CAMERA	
A CONTRACT OF		Load time: 791.30 ms Build Median inference time	time: 1268.80 ms 4.50 ms
and a seal	#	Label	Probability
A PAR	1	Labrador retriever	93.70%
· · · · · · · · · · · · · · · · · · ·	2	golden retriever	2.85%
Pick Image	3	basset	0.43%
	w20 0		
WEBNN AFT SUPPORTED	wac s	SPEC USE CASES	
©2023 W	/ebNN /	API	
0			e

Web Apps (MobileNet v2)				
Web NN API				
ML Service Backend				
TensorFlow Lite				
(Design is ready, Impl is on going)				
OpenVINO Plugin				
NPU Stack (UMD/KMD Driver, FW)				
NPU				



Call to Action

• Try WebNN CPU/GPU acceleration on Google Chrome and Microsoft Edge Canary today

- Navigate to **about:flags** in the URL bar and turn on "Enables WebNN API"
- Visit <u>https://webmachinelearning.github.io/webnn-samples/image_classification/</u>
- Visit https://webmachinelearning.github.io/webnn-samples/semantic_segmentation/index.html
- Join Intel and other industry leaders, including Microsoft and Google, participate in WebNN standardization and open-source implementation



References

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



Web Machine Learning Working Group Launch

Intel Moh Haghighat, Ningxin Hu

Microsoft Chai Chaoweeraprasit





Notices and Disclaimers

For notices, disclaimers, and details about performance claims, visit <u>www.intel.com/PerformanceIndex</u> or scan the QR code:



© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.



Intel® Summer Thank You!

