# An Overview of Web-Based Machine Learning Frameworks

**Vishakha Agrawal**

vishakha.research.id@gmail.com

**Abstract**

**The rapid convergence of machine learning (ML) and web technologies has given rise to a new generation of web-based ML frameworks, revolutionizing the deployment and accessibility of AI capabilities. This paper provides an in-depth examination of the current landscape of web-based ML frameworks, delving into their architectural designs, functional capabilities, and diverse applications in modern web develop-ment. We investigate how these frameworks facilitate the seamless deployment and inference of ML models directly within web browsers, thereby democratizing access to AI-driven insights while mitigating privacy concerns and minimizing server-side dependencies. By exploring the benefits, challenges, and future directions of web-based ML frameworks, this overview aims to provide a comprehensive understanding of the transformative potential of AI in the web ecosystem.**

**Keywords: WebAssembly, WebGL, TensorFlow.js, ONNX.js, ML5.js, Real-time object detection system, Natural Language Processing**

## I. INTRODUCTION

The intersection of machine learning and web technologies has given rise to a new category of frameworks that enable machine learning capabilities directly in web browsers[7]. These frameworks represent a paradigm shift from traditional server-based machine learning deployments, offering advan-tages in privacy, latency, and accessibility. As web browsers have evolved into sophisticated computing platforms, they have become capable of handling complex machine learning tasks that were previously confined to server environments.

## II. BACKGROUND

1) Evolution of Browser-Based Computing : Web browsers have undergone a remarkable transformation [8] from simple document viewers to sophisticated computing platforms capable of executing complex calculations and processing large datasets. This evolution has been particularly enabled by the introduction of technologies like WebAssembly[4] and WebGL[9], which provide near-native performance for computational tasks. These advancements have laid the groundwork for running machine learning models directly in the browser environment.

2) Client-Side Machine Learning Benefits : The shift to-ward client-side machine learning [3] brings numerous advantages to modern web applications. Privacy is sig-nificantly enhanced as

sensitive data can be processed locally without transmission to remote servers. Orga- nizations can reduce their server infrastructure costs and bandwidth usage since computation happens on the client side. Applications benefit from lower latency, particularly crucial for real-time processing tasks. Users gain the ability to use applications offline, and develop- ers can simplify their deployment processes by eliminat- ing complex server-side machine learning infrastructure.

## III. MAJOR FRAMEWORKS

1) TensorFlow.js : TensorFlow.js, developed by Google [12], stands as the most mature web-based ML framework in the ecosystem. It provides a comprehensive implementation of the TensorFlow architecture in JavaScript, leveraging WebGL acceleration for computational operations. The framework supports both pre-trained model deployment and custom model training directly in the browser. Its extensive documentation and active community support have made it the go-to choice for many web-based machine learning projects.

2) ONNX.js : ONNX.js [1] brings the power of the Open Neural Network Exchange (ONNX) format to web browsers, enabling deployment of models trained in various frameworks. The framework offers both WebAssembly and WebGL backends, allowing developers to choose the most appropriate execution environment for their use case. Its efficient inference engine supports complex neural network architectures while maintaining cross-platform compatibility.

3) ML5.js : ML5.js takes a unique approach by focusing on accessibility for creative coding and artists. The framework provides high-level abstractions for common machine learning tasks, integrating seamlessly with cre- ative coding libraries like p5.js. It emphasizes simplified API design and includes pre-trained models for common use cases, making machine learning more accessible to non-experts and creative practitioners.

## IV. TECHNICAL ARCHITECTURE

1) Computation Backends : Modern web-based ML frameworks implement multiple computation backends to maximize performance across different devices and use cases. At the core of these implementations is WebGL[10], which enables GPU-accelerated operations by leveraging the graphics processing capabilities of client devices. This is complemented by WebAssembly, which provides optimized CPU computations through low-level code execution. The frameworks also maintain JavaScript-based computation paths, ensuring broad compatibility across different browsers and devices while offering flexibility for simpler operations.

2) Memory Management : Efficient memory management stands as a critical component in browser-based ma- chine learning implementations. These frameworks em- ploy typed arrays for efficient data storage, signifi- cantly reducing memory overhead compared to stan- dard JavaScript objects. Advanced garbage collection optimization techniques help maintain consistent perfor- mance during long-running operations. Memory pooling strategies are implemented to reduce allocation over- head and memory fragmentation. WebAssembly mem- ory management[2] provides additional optimization op- portunities through direct memory access and control.

## V. PERFORMANCE ANALYSIS

1) **Computational Performance** : Our analysis of computational performance across frameworks reveals interesting patterns and trade-offs. TensorFlow.js achieves near-native performance for specific operations, particularly when utilizing WebGL acceleration for matrix operations and convolutions. Frameworks leveraging WebAssembly show remarkably consistent performance across different browsers, though with some variation in initial compilation time. GPU acceleration provides substantial performance improvements for compatible operations, often achieving an order of magnitude faster execution compared to CPU-only implementations.

2) **Memory Usage** : Memory consumption patterns demon- strate significant variation across frameworks and use cases. Model size has a direct impact on initial load times and memory footprint, with larger models requir- ing careful optimization and potentially lazy loading strategies. Runtime memory usage varies considerably between frameworks, with some showing more efficient memory recycling than others. The garbage collection behavior of different browsers can significantly impact sustained performance, particularly during long-running operations or when processing large datasets.

## VI. APPLICATIONS AND USE CASES

1) **Computer Vision** : Computer vision applications have emerged as a particularly successful domain for browser-based machine learning. Real-time object detection systems [5] now run efficiently in browsers, enabling applications from security monitoring to interactive art installations. Face recognition capabilities have been implemented with consideration for privacy, processing sensitive biometric data entirely on the client side. Image classification systems provide immediate feedback for applications ranging from accessibility tools to educational platforms. Pose estimation frameworks enable new forms of human-computer interaction directly through the web browser.

2) **Natural Language Processing** : Natural Language Pro- cessing applications demonstrate the versatility and ca- pability of web-based machine learning frameworks[11]. Text classification systems enable content moderation and organization directly in the browser. Sentiment analysis tools provide immediate feedback for user- generated content without server roundtrips. Language translation features operate offline, making them ac- cessible in situations with limited connectivity. Text generation capabilities enable creative writing assistance and code completion tools that maintain user privacy by processing all data locally.

## VII. CHALLENGES AND LIMITATIONS

1) **Technical Constraints** : Web-based machine learning frameworks face several significant technical constraints that influence their adoption and implementation. Browser memory constraints [6] limit the size and complexity of models that can be deployed effectively, requiring careful optimization and model compression techniques. Access to system resources remains restricted by browser security models, limiting certain types of hardware acceleration and system integration. Hardware acceleration support varies significantly across devices and browsers, necessitating robust fallback mechanisms. Cross-browser compatibility issues persist, requiring extensive testing and alternative implementation paths.

Development Challenges: Developers working with web-based machine learning frameworks encounter unique challenges throughout the development lifecycle. Framework selection requires careful consideration of factors including performance requirements, browser support, and specific use case demands. Performance optimization becomes increasingly complex when bal- ancing multiple backends and browser environments. Debugging tools for browser-based machine learning remain less mature than their server-side counterparts, making development and troubleshooting more chal- lenging. Browser-specific implementations often require additional code paths and testing procedures, increasing development overhead..

## VIII. FUTURE DIRECTIONS

The future of web-based machine learning frameworks shows tremendous promise with several emerging develop- ments. The integration of WebGPU promises to deliver signif- icantly improved performance through more direct access to GPU capabilities. Enhanced WebAssembly features, including garbage collection and direct DOM access, will enable more sophisticated implementations with better performance charac- teristics. Federated learning implementations are beginning to emerge, allowing collaborative model training while maintain- ing user privacy. Progressive web app integration continues to improve, enabling more sophisticated offline capabilities and better resource management.

## IX. CONCLUSION

Web-based machine learning frameworks represent a signif- icant advancement in democratizing AI capabilities, bringing sophisticated machine learning features directly to the browser environment. While challenges exist in terms of performance, memory management, and browser compatibility, the contin- ued development of web standards and browser capabilities suggests a bright future for client-side machine learning ap- plications. As these frameworks mature and new technologies emerge, we can expect to see increasingly sophisticated appli- cations that leverage the unique advantages of browser- based machine learning while addressing current limitations.

## REFERENCES

[1] Will Badr. Onnx.js: Universal deep learning models in the browser, 2019.

[2] Craig Disselkoen, John Renner, Conrad Watt, Tal Garfinkel, Amit Levy, and Deian Stefan. Position paper: Progressive memory safety for webassembly. In *Proceedings of the 8th International Workshop on Hardware and Architectural Support for Security and Privacy*, pages 1–8, 2019.

[3] Renjie Gu, Chaoyue Niu, Fan Wu, Guihai Chen, Chun Hu, Chengfei Lyu, and Zhihua Wu. From server-based to client-based machine learning: A comprehensive survey. *ACM Computing Surveys (CSUR)*, 54(1):1–36, 2021.

[4] Andreas Haas, Andreas Rossberg, Derek L Schuff, Ben L Titzer, Michael Holman, Dan Gohman, Luke Wagner, Alon Zakai, and JF Bastien. Bringing the web up to speed with webassembly. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 185–200, 2017.

[5] Sabir Hossain and Deok-jin Lee. Deep learning-based real-time multiple-object detection and tracking from aerial imagery via a flying robot with gpu-based embedded devices. *Sensors*, 19(15):3371, 2019.

[6] Marisa Kirisame, Pranav Shenoy, and Pavel Panchekha. Optimal heap limits for reducing browser

memory use. 6(OOPSLA2), October 2022.

[7] Edward Meeds, Remco Hendriks, Said Al Faraby, Magiel Bruntink, and Max Welling. Mlitb: machine learning in the browser. *PeerJ Computer Science*, 1:e11, 2015.

[8] Juan-J Merelo, Antonio Mora Garc´ıa, Juan Luis Jime´nez Laredo, Juan Lupio´n, and Fernando Tricas. Browser-based distributed evolutionary computation: performance and scaling behavior. In *Proceedings of the 9th annual conference companion on Genetic and evolutionary computation*, pages 2851–2858, 2007.

[9] Tony Parisi. *WebGL: up and running*. " O'Reilly Media, Inc.", 2012.

[10] Xenofon Pournaras and Dimitrios A Koutsomitropoulos. Deep learning on the web: state-of-the-art object detection using web-based client-side frameworks. In *2020 11th International Conference on Information, Intelligence, Systems and Applications (IISA*, pages 1–8. IEEE, 2020.

[11] Muhammed Ali Sit, Caglar Koylu, and Ibrahim Demir. Identifying disaster-related tweets and their semantic, spatial and temporal context using deep learning, natural language processing and spatial analysis: a case study of hurricane irma. In *Social Sensing and Big Data Computing for Disaster Management*, pages 8–32. Routledge, 2020.

[12] Daniel Smilkov, Nikhil Thorat, Yannick Assogba, Charles Nicholson, Nick Kreeger, Ping Yu, Shanqing Cai, Eric Nielsen, David Soegel, Stan Bileschi, et al. Tensorflow. js: Machine learning for the web and beyond. *Proceedings of Machine Learning and Systems*, 1:309–321, 2019.