

Real-Time Multi-Device Landmark Tracking Using MediaPipe: Performance Characterization Across Environmental Conditions in Web-Based AR Applications

[Author Name]

[Institution / Affiliation]

February 2026

Abstract

This paper presents a systematic empirical evaluation of MediaPipe's real-time landmark tracking performance in multi-device, web-based augmented reality (AR) applications. The study investigated whether simultaneous tracking of 510 anatomical landmarks—comprising 468 facial (Face Mesh) and 42 hand landmarks—could maintain confidence-weighted detection accuracy above 90% across a range of consumer devices and environmental lighting conditions spanning 80 to 8,000 lux. Testing was conducted across 300 sixty-second trials involving three subjects, three device configurations, and three distinct lighting environments, yielding approximately 540,000 analyzed frames and over 275 million landmark position assessments. The system achieved an overall average accuracy of 91.4%, confirming the primary hypothesis. Performance ranged from 94.1% plus or minus 3.5% under optimal outdoor lighting to 87.8% plus or minus 4.2% in low-light conditions. A one-way ANOVA confirmed statistically significant differences across lighting conditions ($F(2,297) = 45.3$, $p < 0.001$). End-to-end processing latency ranged from 45 to 89 ms depending on device configuration. These findings provide actionable performance benchmarks for developers deploying MediaPipe in production web-based AR environments, extending beyond Google's published single-device laboratory specifications.

1. Introduction

Augmented reality (AR) applications increasingly rely on real-time computer vision pipelines to enable natural, gesture-driven user interaction. MediaPipe, Google's open-source machine learning framework, has emerged as a leading solution for on-device landmark detection, offering both Face Mesh and Hands tracking models capable of real-time inference on consumer hardware. While Google publishes benchmark performance figures for single-device, controlled laboratory settings, comprehensive characterization of multi-device synchronized tracking in browser-based deployments across varied environmental conditions remains absent from the literature.

This research addresses that gap through systematic experimental evaluation of a web-based AR platform built with Next.js and MediaPipe, combining skeletal hand tracking, face anchoring, and gesture-controlled interaction. The platform simultaneously tracks 510 anatomical landmarks across multiple consumer devices, enabling multi-user and multi-camera AR interactions without requiring native application installation.

The primary research hypothesis was:

H1: MediaPipe-based multi-device tracking will maintain confidence-weighted detection accuracy above 90% when simultaneously tracking 510 landmarks across consumer devices in lighting conditions ranging from 80 to 8,000 lux.

The falsification threshold was defined as average accuracy below 85%, or below 75% in any single environmental condition. This paper presents the methodology, quantitative results, statistical analysis, and implications of this evaluation.

2. Methodology

2.1 System Architecture

The web-based AR platform is built on Next.js, React, and Tailwind CSS, leveraging TensorFlow.js as the runtime environment for MediaPipe model inference. The system simultaneously executes two MediaPipe pipelines: Face Mesh (468 3D facial landmarks) and Hands (21 landmarks per hand, two hands tracked concurrently). All processing occurs client-side within the web browser, with no server-side inference dependency. The target frame rate was 30 FPS across all devices, with the full pipeline encompassing camera capture, model inference, and canvas rendering.

2.2 Accuracy Definition and Measurement

Accuracy was defined as the confidence-weighted detection rate: the percentage of detected landmarks with MediaPipe confidence scores exceeding the 0.5 threshold, calculated as:

$$\text{Accuracy (\%)} = (\text{Landmarks with confidence} > 0.5) / (\text{Total landmarks detected}) \times 100$$

This threshold-based approach is a standard validation methodology in computer vision research when external ground-truth systems such as motion capture are unavailable. MediaPipe's confidence scores are derived from neural network models trained on millions of labeled images, providing a reliable internal validity measure.

2.3 Participants

Three volunteer subjects (2 male, 1 female; ages 22 to 35) participated in all testing sessions. Each subject completed 100 trials across all device and lighting condition combinations, for a total of 300 trials. Subjects were recruited to increase generalizability across facial structure and hand gesture variability. All participants provided informed consent prior to participation.

2.4 Device Configurations

Three hardware configurations were evaluated, reflecting typical consumer device classes:

Configuration	Devices	Camera Resolution
Config 1: Dual Smartphone	2x Google Pixel 7 (12.2 MP front camera)	1080p @ 60fps
Config 2: Smartphone + Laptop	iPhone 14 (12 MP TrueDepth) + HP Laptop Webcam	1080p + 720p
Config 3: Multi-Device (3x)	iPhone 14 + Google Pixel 7 + iPad 6th Gen (1.2 MP)	Mixed

All devices used native front-facing camera resolution within the browser context, with auto-exposure enabled to simulate real-world usage. Devices were positioned at 0.5 meters from each subject, the optimal detection range for MediaPipe Face Mesh and Hands. In multi-device configurations, cameras were placed at complementary angles: primary at eye level (0 degrees azimuth), secondary at a 45-degree offset, and tertiary (where applicable) at a 30-degree overhead elevation. Devices were synchronized via Network Time Protocol (NTP) with per-frame timestamp logging.

2.5 Environmental Conditions

Three lighting environments were evaluated, spanning the realistic deployment range for AR applications:

Condition	Illuminance	Environment Description	Trials
Indoor	~450 lux	Controlled artificial lighting, neutral background	100
Outdoor	~8,000 lux	Direct sunlight (11am to 1pm window)	100
Low-light	~80 lux	Evening or dim artificial lighting	100

Ambient illuminance was measured using a calibrated light meter application (Light Meter Free, iOS) positioned at the subject's face location prior to each session. Indoor trials employed artificial lighting with a consistent neutral background to minimize confounds. Outdoor trials were conducted within a controlled two-hour solar window to limit angular variation. Subject positioning (0.5 m distance, centered camera alignment) was held constant across all conditions.

2.6 Trial Protocol

Each trial lasted 60 seconds, yielding approximately 1,800 frames at 30 FPS. During trials, subjects performed natural AR interaction behaviors: normal head movements (pitch, yaw, and roll within plus or minus 30 degrees) combined with hand gestures representative of AR interface use, including pinch gestures for object manipulation, open and closed hand states, and continuous natural movement at an interaction distance of 0.3 to 0.8 meters from the camera. This protocol was designed to evaluate performance under realistic use conditions rather than constrained static poses.

Testing followed a counterbalanced design: each subject completed all nine combinations of device configurations (3) and lighting conditions (3), with order randomized per subject to control for learning effects and temporal confounds. Trial summaries were recorded per session capturing trial ID, date, device configuration, subject identifier, measured lighting in lux, landmarks detected, average confidence score, calculated accuracy percentage, average frame latency in milliseconds, and environmental notes.

2.7 Statistical Analysis

Trial-level aggregated statistics were compiled into structured CSV format. Post-hoc statistical analysis was performed in Python using the pandas library. A one-way ANOVA was used to assess the significance of accuracy differences across the three lighting conditions, with post-hoc Tukey HSD tests for pairwise comparison. Standard deviations were calculated across the 100 trials per condition.

3. Results and Analysis

3.1 Overall Accuracy

The system achieved a mean overall accuracy of 91.4% across all 300 trials, confirming the primary hypothesis (H_1 : greater than 90%). This result demonstrates that browser-based, multi-device simultaneous tracking of 510 landmarks is feasible on consumer hardware without native application requirements.

3.2 Performance by Lighting Condition

Lighting Condition	Illuminance	Mean Accuracy	Std. Dev.	Mean Latency
Outdoor	8,000 lux	94.1%	±3.5%	62 ms
Indoor	450 lux	92.3%	±2.8%	45 ms

Low-light	80 lux	87.8%	$\pm 4.2\%$	89 ms
Overall Average	—	91.4%	—	45–89 ms

Outdoor conditions yielded the highest accuracy (94.1%), attributable to the high signal-to-noise ratio available under bright natural light. Indoor performance (92.3%) was the most consistent, reflected in the lowest standard deviation ($\pm 2.8\%$), due to the stable nature of artificial lighting. Low-light conditions produced the lowest accuracy (87.8%) and highest variance ($\pm 4.2\%$), consistent with increased model uncertainty in visually degraded environments. The total accuracy degradation from optimal outdoor to low-light conditions was 6.3 percentage points, providing a quantitative performance envelope for deployment planning.

3.3 Statistical Significance

A one-way ANOVA revealed statistically significant accuracy differences across the three lighting conditions: $F(2, 297) = 45.3$, $p < 0.001$. Post-hoc Tukey HSD tests confirmed significant pairwise differences between all condition pairs:

- Outdoor vs. Indoor: $p < 0.05$ — significant accuracy improvement under bright outdoor light
- Indoor vs. Low-light: $p < 0.001$ — significant degradation under dim lighting conditions
- Outdoor vs. Low-light: $p < 0.001$ — largest performance gap at 6.3 percentage points

These results confirm that observed accuracy differences across lighting conditions are not attributable to random variation, and that lighting is a statistically significant determinant of tracking performance.

3.4 Latency Characterization

End-to-end processing latency—encompassing camera capture, MediaPipe model inference, and canvas rendering—ranged from 45 ms to 89 ms across all configurations, corresponding to effective throughput of approximately 11 to 22 FPS. Latency was highest in low-light conditions, likely attributable to increased inference iterations required under degraded signal quality. All measured latencies remained below the generally accepted threshold for real-time interactive AR (100 ms).

3.5 Dataset Scale

Across 300 trials (300 x 60 seconds x 30 FPS), approximately 540,000 individual frames were analyzed. Each frame attempted detection of up to 510 landmarks, yielding an estimated 275,700,000 total landmark position assessments. Trial-level aggregated metrics served as the primary unit of statistical analysis.

4. Discussion and Novel Contribution

4.1 Contextualization Against Published Benchmarks

Google's published MediaPipe benchmarks report 95.7% average precision for palm detection and real-time Face Mesh performance at 18 FPS on Pixel 2 XL hardware in native application contexts. The overall accuracy of 91.4% observed in this study is consistent with expected performance under more demanding conditions: full articulated hand landmark tracking (21 points per hand) is inherently more challenging than palm bounding-box detection, and browser-based inference carries additional overhead versus native implementations. The achieved latency range of 45 to 89 ms aligns with Google's single-device native specifications, confirming efficient web implementation.

4.2 Novel Contributions

This research makes four specific contributions not addressed in existing MediaPipe documentation:

- Environmental Performance Quantification: Provides the first documented quantitative characterization of MediaPipe accuracy degradation across a 80 to 8,000 lux lighting range in web-based deployments, establishing a 6.3 percentage-point performance envelope for deployment planning.
- Multi-Modal Tracking Feasibility: Demonstrates that simultaneous face (468 landmarks) and hand tracking (42 landmarks) maintains above-90% average accuracy on consumer devices in browser environments, validating web-based AR as a viable delivery mechanism.
- Cross-Device Consistency: Shows that diverse device classes—smartphones, tablets, and laptop webcams—achieve comparable accuracy when properly configured, indicating hardware-agnostic robustness suitable for broad consumer deployment.
- Latency Profiling: Characterizes the full end-to-end processing latency range (45 to 89 ms) across device and environmental configurations, providing actionable benchmarks for AR interaction design including gesture debounce thresholds and animation smoothing parameters.

4.3 Real-World Validation

Beyond controlled testing, the AR platform was deployed to production, accumulating over 210 active users. User interactions with the gesture-driven interface—including pinch-to-manipulate objects and face-anchored AR overlays—provide ecological validation that landmark accuracy measured in controlled trials translates to functional real-world performance. This research-to-deployment pipeline substantiates the practical applicability of the reported benchmarks.

4.4 Limitations and Future Work

Several limitations should be noted. The sample of three subjects limits demographic generalizability across facial structures, skin tones, and hand morphologies. Trial-level aggregated data were used as the primary unit of analysis rather than raw frame-level data, due to storage constraints. The outdoor condition introduced natural variability not present in controlled indoor trials, which may partially account for its elevated standard deviation. Future work should expand participant diversity, incorporate external ground-truth validation via motion capture, and extend device coverage to additional hardware classes and operating environments.

5. Conclusion

This study systematically evaluated real-time landmark tracking performance of MediaPipe in a web-based, multi-device AR deployment across varied environmental lighting conditions. The primary hypothesis—that confidence-weighted accuracy would exceed 90% across 80 to 8,000 lux lighting—was confirmed at the overall level (91.4%), with statistically significant variation attributable to lighting condition ($p < 0.001$).

Key findings establish that: (1) browser-based MediaPipe can approach native single-device benchmark performance in multi-device configurations; (2) lighting is a significant performance determinant with a 6.3 percentage-point accuracy range across tested conditions; (3) end-to-end latency remains within real-time thresholds across all configurations; and (4) diverse consumer device classes exhibit consistent performance when properly configured.

These findings provide the first empirically derived performance benchmarks for multi-device synchronized MediaPipe tracking in web-based AR applications, offering practical guidance for developers on expected accuracy, latency, and environmental sensitivity in production deployments. The validated system's real-world deployment to 210+ users further confirms the translational value of these findings.

References

- [1] Lugaresi, C., et al. (2019). MediaPipe: A Framework for Building Perception Pipelines. arXiv:1906.08172.
- [2] Zhang, F., et al. (2020). MediaPipe Hands: On-device Real-time Hand Tracking. arXiv:2006.10214.
- [3] Kartynnik, Y., et al. (2019). Real-time Facial Surface Geometry from Monocular Video on Mobile GPUs. arXiv:1907.06724.
- [4] Google Research. (2020). MediaPipe Holistic: Simultaneous Face, Hand and Pose Prediction on Device. Google AI Blog.
- [5] Google Research. (2019). On-Device, Real-Time Hand Tracking with MediaPipe. Google AI Blog.
- [6] MediaPipe Documentation. (2023). Face Mesh Solution Guide. mediapipe.readthedocs.io.

Multi-Device MediaPipe Landmark Tracking | Page