

Network-Adaptive Cloud Preprocessing for Visual Neuroprostheses

Jiayi Liu
Computer Science
UC Santa Barbara
Santa Barbara, CA, USA
jiayi979@ucsb.edu

Yilin Wang
Computer Science
UC Santa Barbara
Santa Barbara, CA, USA
yilin_wang@ucsb.edu

Michael Beyeler
Computer Science
Psychological & Brain Sciences
UC Santa Barbara
Santa Barbara, CA, USA
mbeyeler@ucsb.edu

Abstract—Cloud-based machine learning is increasingly explored as a preprocessing strategy for next-generation visual neuroprostheses, where advanced scene understanding may exceed the computational and energy constraints of battery-powered visual processing units (VPUs). Offloading computation to remote servers enables the use of state-of-the-art vision models, but also introduces sensitivity to network latency, jitter, and packet loss, which can disrupt the temporal consistency of the delivered neural stimulus. In this work, we examine the feasibility of cloud-assisted visual preprocessing for artificial vision by framing remote inference as a perceptually constrained systems problem. We present a network-adaptive cloud-assisted pipeline in which real-time round-trip-time (RTT) feedback is used to dynamically modulate image resolution, compression, and transmission rate, explicitly prioritizing temporal continuity under adverse network conditions. Using a Raspberry Pi 4 as a simulated VPU and a client-server architecture, we evaluate system performance across a range of realistic wireless network regimes. Results show that adaptive visual encoding substantially reduces end-to-end latency during network congestion, with only modest degradation of global scene structure, while boundary precision degrades more sharply. Together, these findings delineate operating regimes in which cloud-assisted preprocessing may remain viable for future visual neuroprostheses and underscore the importance of network-aware adaptation for maintaining perceptual stability.

Index Terms—visual neuroprostheses, cloud-assisted processing, cloud AI, network-adaptive systems, scene simplification

I. INTRODUCTION

Bionic vision is increasingly exploring the use of machine learning (ML) and artificial intelligence (AI) as preprocessing tools for next-generation retinal and cortical prostheses [1]–[4]. Rather than operating directly on raw camera input, many proposed approaches seek to transform egocentric video streams into simplified, task-relevant representations that may be more interpretable to the visual system. These strategies include semantic and structural edge extraction [5]–[7], object-level abstraction [8], [9], and depth-aware scene simplification [10], [11], with the goal of supporting mobility and navigation. While such AI-based preprocessing is not yet deployed in current clinical devices, it represents a rapidly growing research direction for mitigating the perceptual limitations imposed by existing prosthetic hardware.

A major obstacle to these AI-driven approaches is the severe hardware constraint of current clinical systems. Contemporary prosthetic devices rely on battery-powered vision processing

units (VPUs) that prioritize energy efficiency and thermal safety over computational throughput. As a result, many modern ML models, particularly those required for complex indoor scene understanding, exceed the feasible limits of onboard processing and cannot be deployed directly on the device.

To address this computational gap, two complementary research trajectories have emerged largely outside the bionic vision literature. One approach focuses on highly efficient, lightweight algorithms, often referred to as *TinyML*, that are explicitly designed for low-power internet of things (IoT)-class devices [12]. Although promising, such models typically achieve efficiency by reducing representational depth or semantic richness, which can limit their effectiveness for visually demanding mobility tasks [13], [14]. A second approach leverages cloud-based inference, offloading computationally intensive processing to remote servers capable of running state-of-the-art vision models [15], [16]. While this paradigm is widely used in other vision and sensing domains, it has received limited attention in bionic vision research, particularly with respect to its sensitivity to network latency, jitter, and packet loss [17], [18].

For visual neuroprostheses, the vulnerability of cloud-centric preprocessing pipelines to network-induced delays presents a fundamental challenge. Unlike conventional computer vision applications, delays in prosthetic vision are not merely an inconvenience but would directly affect the temporal consistency of the delivered neural stimulus [19], [20]. Prior work in visual perception and sensorimotor integration indicates that delays on the order of tens to hundreds of milliseconds can degrade perceptual stability, impair visuo-motor coordination, and increase disorientation during mobility tasks [20]–[23]. In the context of artificial vision, such delays may manifest as interruptions in perceptual continuity, undermining both usability and safety.

Here, we investigate cloud-based semantic segmentation as a candidate preprocessing strategy for future visual neuroprostheses and address its central limitation by introducing a network-adaptive control policy (Figure 1). Our approach dynamically modulates visual encoding parameters based on real-time network feedback, explicitly prioritizing temporal continuity of the percept over spatial fidelity as network conditions degrade. By framing cloud inference as a perceptually

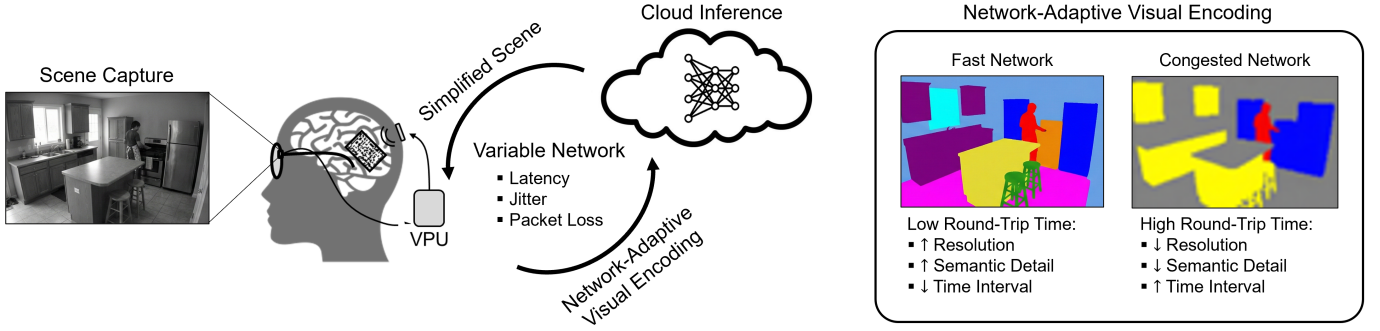


Fig. 1: Network-adaptive cloud processing for visual neuroprostheses. Egocentric video captured by a resource-constrained vision processing unit (VPU) is adaptively encoded prior to transmission based on real-time network feedback. Round-trip time (RTT) measurements drive a closed-loop controller that modulates image resolution, compression, and transmission interval to maintain temporal continuity of the delivered visual stimulus under network impairments. Remote semantic segmentation is performed in the cloud, and the resulting simplified scene is returned to the client. In the present system, reduced semantic detail under network congestion arises indirectly from input degradation; future implementations could achieve similar effects by explicitly requesting coarse vs. fine semantic outputs from the cloud inference service. Scene images were generated using Nano Banana Pro for illustrative purposes only.

constrained control problem rather than a static vision pipeline, we seek to identify operating regimes in which cloud-based preprocessing can remain viable for artificial vision despite stochastic network impairments.

Our contributions are threefold:

- i. We present a cloud-assisted visual preprocessing pipeline tailored to the constraints of visual neuroprostheses and systematically characterize its sensitivity to network impairments across a range of realistic connectivity regimes.
- ii. We introduce a closed-loop, network-adaptive encoding strategy that explicitly trades spatial fidelity for temporal stability, enabling the system to maintain a minimum perceptual update rate under adverse network conditions.
- iii. We quantify the resulting latency-fidelity trade-off using perceptually motivated measures, demonstrating that substantial reductions in end-to-end delay can be achieved with only modest degradation of global scene structure, while boundary precision degrades more sharply under severe congestion.

Together, this work establishes a principled framework for evaluating when and how cloud-assisted visual preprocessing can be integrated into future neuroprosthetic systems without violating perceptual timing constraints critical for safe and effective use.

II. METHODS

A. System Overview and Architecture

Figure 1 provides an overview of the proposed network-adaptive cloud-assisted visual preprocessing system. The architecture is designed to approximate a future visual neuroprosthesis workflow in which a resource-constrained vision processing unit (VPU) captures egocentric video and offloads computationally intensive preprocessing to a remote server. The central challenge addressed by this system is maintaining

perceptual temporal continuity in the presence of stochastic network impairments.

The system comprises three principal components:

- i. a VPU-side client responsible for video capture and adaptive visual encoding,
- ii. a bidirectional network channel subject to controlled latency, bandwidth, and packet loss, and
- iii. a cloud-based inference server that performs semantic preprocessing and returns a simplified scene representation to the client.

A closed-loop control policy operating on the client dynamically adjusts visual encoding parameters based on real-time network feedback, explicitly trading spatial fidelity for reduced end-to-end delay under adverse network conditions.

B. Network-Adaptive Visual Encoding Policy

The core contribution of this work is a closed-loop, network-adaptive encoding policy that stabilizes the temporal delivery of cloud-based visual preprocessing for artificial vision. Rather than treating remote inference as a static computer vision pipeline, the proposed system explicitly accounts for network variability by regulating the visual encoding process on the VPU in response to observed communication delays.

The design objective of the controller is to maintain perceptual temporal continuity of the delivered stimulus, even when network conditions degrade, by preventing excessive end-to-end latency accumulation. To achieve this, the controller dynamically modulates visual encoding parameters, prioritizing timely stimulus updates over spatial detail during periods of congestion.

1) *Network Feedback Signal:* Network responsiveness is quantified on the client using the round-trip time (RTT) between the VPU and the cloud server. A dedicated monitoring thread periodically probes the communication channel and records the most recent K RTT measurements in a bounded

TABLE I: Network-adaptive encoding tiers used by the closed-loop controller.

RTT threshold	JPEG quality	Max. resolution	Send interval
≤ 30 ms	90%	1920 px	80 ms
≤ 50 ms	80%	1280 px	100 ms
≤ 100 ms	65%	960 px	150 ms
≤ 150 ms	50%	720 px	250 ms
> 150 ms	40%	480 px	500 ms

buffer. To reduce sensitivity to transient jitter and outliers, the controller operates on a moving average estimate,

$$\overline{\text{RTT}} = \frac{1}{K} \sum_{i=1}^K \text{RTT}_i, \quad (1)$$

where $K = 5$ in the present implementation. This smoothed estimate serves as the sole feedback signal driving adaptive reconfiguration.

2) *Tiered Reconfiguration Policy*: Based on the estimated $\overline{\text{RTT}}$, the controller selects one of five discrete operating regimes. Each regime specifies a visual encoding parameter vector

$$\mathcal{P} = \{Q, R, I\},$$

where Q denotes the JPEG compression quality, R the maximum image resolution (with aspect ratio preserved), and I the inter-frame transmission interval.

Discrete operating tiers were chosen to ensure predictable and stable behavior under fluctuating network conditions while minimizing computational overhead on the VPU. As network delay increases, the controller progressively reduces spatial resolution and compression quality and increases the inter-frame interval, thereby limiting queue buildup and preventing excessive end-to-end latency.

The thresholds and corresponding parameter settings defining each operating regime are summarized in Table I.

C. Remote Semantic Preprocessing

Semantic preprocessing is performed on the cloud server using PIDNet [24], a real-time semantic segmentation architecture inspired by proportional-integral-derivative control principles. PIDNet combines three complementary processing branches: (i) a proportional branch that preserves high-resolution spatial structure, (ii) an integral branch that aggregates global contextual information, and (iii) a derivative branch that emphasizes high-frequency boundary features.

This architectural design is well aligned with preprocessing requirements for artificial vision, where preserving salient structure and object boundaries is often more informative than photorealistic reconstruction. The server processes each received frame independently and returns a simplified scene representation to the client as part of the end-to-end loop.

TABLE II: Simulated network conditions used to evaluate cloud-assisted preprocessing.

Scenario	Downlink	Uplink	RTT	Loss
Extreme congested 4G	10 Mbps	5 Mbps	100 ms	5%
Congested 4G	25 Mbps	10 Mbps	100 ms	2%
4G–5G hybrid	50 Mbps	25 Mbps	50 ms	0.5%
Good 5G	200 Mbps	50 Mbps	30 ms	0.1%
Ultra-smooth 5G	800 Mbps	200 Mbps	10 ms	0%

D. Client–Server Communication

Client and server communicate via gRPC¹ over HTTP/2, enabling low-latency bidirectional request–response interactions using Protocol Buffers for binary serialization. Each video frame is transmitted as an encoded byte payload following application of the adaptive policy. Upon receipt, the server decodes the payload, performs semantic preprocessing, and returns the resulting representation to the client.

This request-response cycle constitutes a single iteration of the closed-loop system and is used as the basis for all latency and fidelity measurements reported in this study.

E. Experimental Setup and Network Impairment Model

System performance was evaluated under controlled network conditions designed to approximate a range of realistic wireless connectivity regimes. Network impairments were imposed on the server side using a network emulation framework that independently constrains uplink and downlink bandwidth while injecting additional latency and packet loss. This approach allows systematic evaluation of the proposed adaptive policy under reproducible yet heterogeneous network conditions.

We considered five network scenarios spanning severely constrained 4G-like connectivity to near-ideal 5G-like conditions (Table II). These scenarios were selected to cover operating regimes relevant to mobile and wearable assistive technologies, including cases in which cloud-based preprocessing may become marginal or infeasible without adaptation.

For each network scenario, system performance was evaluated under two operating modes: (i) *a static baseline*, in which visual frames are transmitted at fixed resolution, compression quality, and frame rate; and (ii) *the proposed adaptive policy*, in which the encoding parameter vector $\mathcal{P} = \{Q, R, I\}$ is updated online based on the estimated $\overline{\text{RTT}}$.

F. Outcome Measures

System performance was quantified along two complementary dimensions: (i) temporal responsiveness of the end-to-end processing loop, and (ii) fidelity of the returned semantic representation. This separation reflects the central trade-off explored in this study between perceptual timing constraints and spatial detail.

¹<https://grpc.io>

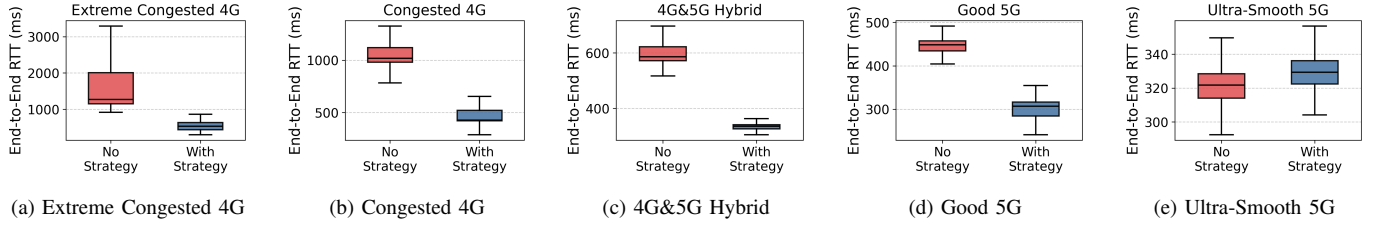


Fig. 2: End-to-end round-trip time (RTT) distributions under five simulated network conditions, comparing a static baseline with the proposed network-adaptive encoding policy.

1) *Temporal Responsiveness*: Temporal responsiveness was quantified using the end-to-end RTT for each transmitted frame, measured from client-side transmission to receipt of the corresponding response. This metric captures the combined effects of network delay, queuing, encoding, decoding, and server-side inference.

In addition, mean server-side inference time was measured independently under each network scenario to characterize the computational contribution to end-to-end delay and to assess how adaptive downscaling influences remote processing time.

2) *Perceptual Fidelity Measures*: Because the present evaluation focuses on system-level behavior rather than human psychophysical performance, perceptual fidelity was assessed using established image-based measures applied to the semantic segmentation output. Specifically, we report the structural similarity index measure (SSIM) and the boundary F1 (BF) score.

SSIM quantifies preservation of global structural information in the returned representation relative to a reference, providing a proxy for scene-level perceptual integrity. The BF score captures boundary precision and is particularly relevant for prosthetic vision pipelines that rely on accurate delineation of salient objects and scene structure. Together, these measures enable systematic characterization of the latency-fidelity trade-off induced by network-adaptive encoding under varying connectivity conditions.

III. RESULTS

A. Temporal Responsiveness Under Network Impairment

We first evaluated the ability of the proposed network-adaptive encoding strategy to stabilize end-to-end latency under varying network conditions. System performance was assessed across five connectivity regimes ranging from severely constrained 4G-like conditions to near-ideal 5G-like operation, as defined in Table II.

Figure 2 shows the distribution of end-to-end RTT measured for each network scenario, comparing the static baseline configuration with the proposed adaptive policy.

Under severely constrained network conditions, including Extreme Congested 4G and Congested 4G, the adaptive policy substantially reduced end-to-end RTT relative to the static baseline. In these regimes, median RTT was reduced by approximately 60–70%, reflecting the controller’s ability to

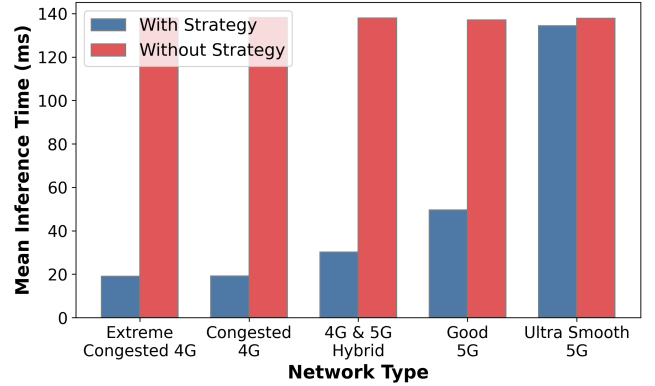


Fig. 3: Mean server-side inference time under each network condition for static and adaptive configurations.

prevent queue buildup and excessive delay accumulation by dynamically reducing spatial resolution and transmission rate.

As network conditions improved, the difference between adaptive and static configurations diminished. In the Ultra-Smooth 5G condition, both configurations exhibited low and stable RTT, and the adaptive policy introduced a small additional overhead due to monitoring and reconfiguration logic. This behavior is expected, as adaptation provides limited benefit when bandwidth and latency constraints are negligible.

B. Server-Side Inference Time

To isolate the computational contribution to end-to-end delay, we additionally measured mean server-side inference time under each network scenario (Fig. 3).

During periods of network congestion, adaptive downscaling substantially reduced inference time by limiting the spatial resolution of frames processed by the cloud server. In the Extreme Congested 4G condition, mean inference time decreased from approximately 118ms under the static baseline to 19ms with adaptation. This reduction contributes directly to improved temporal responsiveness and helps maintain synchronization between visual updates and user motion.

Under high-quality network conditions, inference time differences between configurations were minimal, consistent with the controller operating in its highest-fidelity regime.

TABLE III: Perceptual fidelity measures under static and adaptive configurations across network conditions.

Scenario	SSIM (%)		BF Score (%)	
	Adaptive	Static	Adaptive	Static
Extreme Congested 4G	78.60	81.74	17.30	50.38
Congested 4G	78.67	81.74	17.99	50.38
4G&5G Hybrid	79.84	81.74	27.81	50.38
Good 5G	80.76	81.74	39.01	50.38
Ultra-Smooth 5G	81.78	81.74	50.70	50.38

C. Perceptual Fidelity Trade-Offs

We next examined how adaptive encoding impacts the fidelity of the returned semantic representation. Perceptual fidelity was quantified using SSIM and BF scores, as summarized in Table III.

Across all network conditions, adaptive encoding resulted in modest reductions in SSIM relative to the static baseline. Even under the most severe network impairment, SSIM declined by only 3.14%, indicating that global scene structure was largely preserved despite aggressive compression and downscaling.

As network conditions improved, SSIM values under the adaptive policy converged toward the static baseline, reaching near-identical performance in the Ultra-Smooth 5G regime.

In contrast, BF scores exhibited substantially greater sensitivity to adaptive downscaling. Under Extreme Congested 4G conditions, BF score decreased from 50.38% in the static baseline to 17.30% with adaptation, reflecting the loss of fine-grained boundary detail at low spatial resolutions.

BF scores increased monotonically with improving network quality, reaching parity with the static baseline under Ultra-Smooth 5G conditions. This pattern highlights a key trade-off: during network congestion, the system prioritizes temporal continuity at the expense of precise boundary delineation.

Taken together, these results demonstrate that the proposed adaptive policy effectively regulates the latency–fidelity trade-off inherent to cloud-assisted visual preprocessing. Under constrained network conditions, the system transitions into a low-fidelity but temporally responsive operating regime, substantially reducing end-to-end delay while preserving global scene structure. As network conditions improve, the controller automatically restores higher spatial fidelity without manual intervention.

IV. DISCUSSION

This work examines the feasibility of cloud-assisted visual preprocessing for future visual neuroprostheses under realistic network constraints. By framing remote inference as a perceptually constrained control problem rather than a static computer vision pipeline, we demonstrate that adaptive visual encoding can substantially mitigate the impact of network-induced delays. Across a range of simulated wireless conditions, the proposed closed-loop strategy stabilizes end-to-end latency by dynamically trading spatial fidelity for temporal

continuity, preserving a usable perceptual stream even under severe congestion.

From a neural engineering perspective, the key contribution of this study is not a specific vision model or network protocol, but a principled systems-level framework for reasoning about timing, fidelity, and network variability in artificial vision pipelines. Our results delineate operating regimes in which cloud-based preprocessing may remain viable for visual neuroprostheses, and identify perceptual trade-offs that must be managed to maintain safety and usability.

A. Temporal Continuity as a Primary Design Constraint

A central takeaway from this study is that temporal continuity should be treated as a central design constraint for cloud-assisted artificial vision. Under network congestion, static cloud pipelines exhibit rapidly increasing end-to-end delay, which can disrupt the temporal consistency of the delivered neural stimulus. In contrast, the proposed adaptive policy explicitly regulates latency by limiting queue buildup and inference time, enabling the system to maintain a minimum update rate even when bandwidth and packet loss are unfavorable.

Importantly, the observed reductions in latency were achieved with only modest degradation of global scene structure, as reflected by relatively stable SSIM values across network regimes. This suggests that, for mobility-oriented tasks, preserving coarse scene layout and temporal alignment may be more critical than maintaining fine-grained spatial detail when network resources are limited.

B. Latency-Fidelity Trade-Offs in Artificial Vision

The results also highlight an inherent latency-fidelity trade-off in cloud-assisted preprocessing. While adaptive downscaling effectively reduces latency, it disproportionately impacts boundary precision, as reflected by reductions in BF score under severe congestion. This finding is consistent with the role of high-frequency spatial information in semantic segmentation and underscores that different perceptual attributes degrade at different rates under adaptation.

From a systems design standpoint, this asymmetry suggests that adaptive policies should be task-aware. For example, navigation and obstacle avoidance may tolerate reduced boundary precision if temporal alignment is preserved, whereas tasks such as object identification or reading may require higher spatial fidelity and thus different adaptation strategies. These considerations point toward the need for context-dependent control policies that modulate adaptation based not only on network state but also on behavioral goals.

C. Limitations and Future Directions

Several limitations of the present study warrant consideration. First, perceptual fidelity was assessed using algorithmic measures rather than human psychophysical evaluation. While SSIM and BF score provide useful proxies for global structure and boundary integrity, future work should incorporate behavioral experiments with prosthesis users or sighted participants

viewing simulated artificial vision to more directly assess perceptual consequences.

Second, the current implementation modulates semantic detail indirectly through input degradation. Future systems could extend this framework by explicitly requesting coarse or fine semantic outputs from the cloud inference service, enabling more targeted control over perceptual content without relying solely on spatial downscaling. Additionally, hybrid architectures that combine lightweight on-device preprocessing with cloud-based semantic labeling may further improve robustness under extreme network conditions.

Finally, the adaptive policy explored here relies on discrete operating regimes and a single network feedback signal. More sophisticated controllers that incorporate predictive models of network variability or additional feedback signals could enable smoother transitions and improved performance.

Despite these limitations, the present results provide a concrete foundation for exploring cloud-assisted artificial vision as a viable design space rather than an all-or-nothing proposition.

V. CONCLUSION

As visual neuroprostheses increasingly incorporate advanced machine learning-based preprocessing, computational demands are likely to outpace what can be supported on battery-powered devices alone. This study shows that cloud-assisted preprocessing, when coupled with network-aware adaptive encoding, can remain compatible with perceptual timing constraints critical for safe and effective artificial vision. By grounding cloud inference in a control-theoretic and perceptually informed framework, this work contributes a systems-level perspective that may inform the design of future neuroprosthetic platforms bridging computation, communication, and neural stimulation.

REFERENCES

- [1] M. Beyeler and M. Sanchez-Garcia, "Towards a Smart Bionic Eye: AI-powered artificial vision for the treatment of incurable blindness," *Journal of Neural Engineering*, vol. 19, p. 063001, Dec. 2022.
- [2] J. de Ruyter van Steveninck, U. Güçlü, R. van Wezel, and M. van Gerven, "End-to-end optimization of prosthetic vision," *Journal of Vision*, vol. 22, p. 20, Feb. 2022.
- [3] J. Granley, L. Relic, and M. Beyeler, "Hybrid Neural Autoencoders for Stimulus Encoding in Visual and Other Sensory Neuroprostheses," in *Advances in Neural Information Processing Systems*, vol. 35, pp. 22671–22685, Dec. 2022.
- [4] J. Granley, T. Fauvel, M. Chalk, and M. Beyeler, "Human-in-the-Loop Optimization for Deep Stimulus Encoding in Visual Prostheses," Thirty-seventh Conference on Neural Information Processing Systems, Nov. 2023.
- [5] M. Sánchez García, R. Martínez-Cantín, and J. J. Guerrero, "Semantic and structural image segmentation for prosthetic vision," *PLOS ONE*, vol. 15, p. e0227677, Jan. 2020.
- [6] N. Han, S. Srivastava, A. Xu, D. Klein, and M. Beyeler, "Deep Learning-Based Scene Simplification for Bionic Vision," in *Augmented Humans Conference 2021*, AHs'21, (New York, NY, USA), pp. 45–54, Association for Computing Machinery, Feb. 2021.
- [7] J. de Ruyter van Steveninck, T. van Gestel, P. Koenders, G. van der Ham, F. Vereecken, U. Güçlü, M. van Gerven, Y. Güçlütürk, and R. van Wezel, "Real-world indoor mobility with simulated prosthetic vision: The benefits and feasibility of contour-based scene simplification at different phosphene resolutions," *Journal of Vision*, vol. 22, p. 1, Feb. 2022.
- [8] J. M. Kasowski, A. Varshney, and M. Beyeler, "Static or Temporal? Semantic Scene Simplification to Aid Wayfinding in Immersive Simulations of Bionic Vision," in *Proceedings of the 2025 31st ACM Symposium on Virtual Reality Software and Technology*, VRST '25, (New York, NY, USA), pp. 1–11, Association for Computing Machinery, Dec. 2025.
- [9] A. Nejad, B. Küçükoğlu, J. de Ruyter van Steveninck, S. Bedrossian, G. A. de Haan, J. Heutink, F. W. Cornelissen, and M. van Gerven, "Point-SPV: End-to-End Enhancement of Object Recognition in Simulated Prosthetic Vision using Synthetic Viewing Points," *Frontiers in Human Neuroscience*, vol. 19, Feb. 2025.
- [10] N. M. Barnes, A. F. Scott, A. Stacey, C. McCarthy, D. Feng, M. A. Petoe, L. N. Ayton, R. Dengate, R. H. Guymer, and J. Walker, "Enhancing object contrast using augmented depth improves mobility in patients implanted with a retinal prosthesis," *Investigative Ophthalmology & Visual Science*, vol. 56, p. 755, June 2015.
- [11] A. Rasla and M. Beyeler, "The Relative Importance of Depth Cues and Semantic Edges for Indoor Mobility Using Simulated Prosthetic Vision in Immersive Virtual Reality," in *Proceedings of the 28th ACM Symposium on Virtual Reality Software and Technology*, VRST '22, (New York, NY, USA), pp. 1–11, Association for Computing Machinery, Nov. 2022.
- [12] P. Warden and D. Situnayake, *TinyML: Machine Learning with TensorFlow Lite on Arduino and Ultra-Low-Power Microcontrollers*. Beijing Boston Farnham Sebastopol Tokyo: O'Reilly Media, 2020.
- [13] C. R. Banbury, V. J. Reddi, M. Lam, W. Fu, A. Fazel, J. Hollerman, X. Huang, R. Hurtado, D. Kanter, A. Lohmotov, D. Patterson, D. Pau, J.-s. Seo, J. Sieracki, U. Thakker, M. Verhelst, and P. Yadav, "Benchmarking TinyML Systems: Challenges and Direction," Jan. 2021. arXiv:2003.04821 [cs].
- [14] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," Apr. 2017. arXiv:1704.04861 [cs].
- [15] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, "Neurosurgeon: Collaborative Intelligence Between the Cloud and Mobile Edge," in *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '17, (New York, NY, USA), pp. 615–629, Association for Computing Machinery, Apr. 2017.
- [16] Y. Nan, S. Jiang, and M. Li, "Large-scale Video Analytics with Cloud-Edge Collaborative Continuous Learning," *ACM Trans. Sen. Netw.*, vol. 20, pp. 14:1–14:23, Oct. 2023.
- [17] S. Laskaridis, S. I. Venieris, M. Almeida, I. Leontiadis, and N. D. Lane, "SPINN: synergistic progressive inference of neural networks over device and cloud," in *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, MobiCom '20, (New York, NY, USA), pp. 1–15, Association for Computing Machinery, Sept. 2020.
- [18] S. M. Zobaed, A. Mokhtari, J. P. Champati, M. Kourouma, and M. A. Salehi, "Edge-MultiAI: Multi-Tenancy of Latency-Sensitive Deep Learning Applications on Edge," Nov. 2022. arXiv:2211.07130 [cs].
- [19] R. M. Held and N. I. Durlach, "Telepresence," *Presence: Teleoperators and Virtual Environments*, vol. 1, pp. 109–112, Feb. 1992.
- [20] R. C. Miall and J. K. Jackson, "Adaptation to visual feedback delays in manual tracking: evidence against the Smith Predictor model of human visually guided action," *Experimental Brain Research*, vol. 172, pp. 77–84, June 2006.
- [21] C. Stetson, X. Cui, P. R. Montague, and D. M. Eagleman, "Motor-Sensory Recalibration Leads to an Illusory Reversal of Action and Sensation," *Neuron*, vol. 51, pp. 651–659, Sept. 2006.
- [22] T. Honda, M. Hirashima, and D. Nozaki, "Adaptation to Visual Feedback Delay Influences Visuomotor Learning," *PLOS ONE*, vol. 7, p. e37900, May 2012.
- [23] S. Beech, D. Stanton Fraser, and I. D. Gilchrist, "Visuomotor adaptation to constant and varying delays in a target acquisition task," *Journal of Vision*, vol. 25, p. 8, May 2025.
- [24] J. Xu, Z. Xiong, and S. P. Bhattacharyya, "PIDNet: A Real-time Semantic Segmentation Network Inspired by PID Controllers," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 19529–19539, June 2023. ISSN: 2575-7075.