

MindSight: A Bio-Inspired Neural Architecture for Visual Restoration via Cortical Electrical Stimulation

Yongjie Zou¹, Haonan Niu¹, Bin Zhao¹, Guoliang Yi¹, Mengchuanzhi Yang¹, Jiawei Ju^{1,2}, Jiapeng Yin^{1,2*}, Chengyu T. Li^{1*}

¹Lingang Laboratory

²Shanghai Center for Brain Science and Brain-Inspired Technology

{zouyj, haonan.niu, zhaobin, yiguoliang, yangm, jiawei.ju, jiapeng.yin, tonylicy}@lglab.ac.cn

Abstract

Visual impairment is a common condition worldwide, and cortical electrical stimulation is one of the approaches to aid in visual restoration. However, existing methods suffer from limited precision, flexibility, and generalization in generating the desired visual perception. In this paper, we propose a novel deep learning-based algorithm for cortical electrical stimulation, named “MindSight,” aimed at enhancing the clarity and accuracy of induced visual perceptions. Our framework introduces three key innovations: (1) A differentiable biophysical model simulating cortical state transitions under electrical stimulation, enabling end-to-end training; (2) A dual-path training architecture combining neural decoding fidelity with phosphene simulation constraints; (3) An attention-guided background gated network for input filtration and, a multi-channel activation constraint to ensure the effectiveness of electrical stimulation. We validated our approach through novel experiments with macaque monkeys, demonstrating superior performance in visual perception tasks. These results highlight the potential of our approach in assisting individuals with visual impairments.

Extended version — doi.org/10.1101/2025.11.12.688136

Introduction

Visual impairment is a common condition worldwide (Stevens et al. 2013). The restoration of visual perception through cortical stimulation represents a transformative frontier in neuroprosthetics, offering hope for individuals with severe visual impairments caused by retinal degeneration or optic nerve damage. Traditional approaches, such as retinal implants (Dorn et al. 2013; Stingl et al. 2015), rely on intact neural pathways from the retina to the visual cortex, limiting their utility for patients with broader damage (Fernandez and Robles 2024). This outcome has motivated a shift toward approaches that bypass the eye entirely. In particular, visual cortex stimulation has gained interest (Lozano et al. 2020; Chen et al. 2020; Maghami et al. 2014; Roelfsema, Denys, and Klink 2018; Dobelle 2000) as it can potentially benefit a wider range of blindness causes (e.g. optic nerve damage, glaucoma) where retinal devices are in-

effective. Recent advances with high-channel-count micro-electrode arrays (96–1024 channels) demonstrate that intracortical stimulation can evoke shape-like percepts in both primates and humans (Chen et al. 2020; Fernández et al. 2021). Primary visual cortex (V1) contains an ordered map of the visual field and relatively large surface area, making it an attractive target for electrode implants (Beauchamp et al. 2020; Beauchamp and Yoshor 2020; Lewis et al. 2015; Tehovnik and Slocum 2013; Chen et al. 2020) that could evoke patterned visual percepts. Although these types of visual prostheses may differ significantly in terms of the entry point into the visual system, they share the same fundamental mechanism of action: through electrical stimulation of small groups of neurons, they evoke the perception of spatially localized flashes of light, called phosphenes (Brindley and Lewin 1968; Foroushani, Pack, and Sawan 2018). Despite advancements, research lacks adaptive and dynamic strategies, requiring manual determination of critical visual features and transformation of these features into suitable electrical stimulation protocols (Chen et al. 2020). Moreover, existing automated methods suffer from limited precision, flexibility, and generalization in generating the desired complex visual perception based on actual visual stimuli (Bosking, Beauchamp, and Yoshor 2017).

This work presents a novel bio-inspired algorithm, named “**MindSight**”, for dynamically orchestrating electrode stimulation patterns in visual cortex. Our work is built on five key components: (1) A differentiable biophysical model simulating cortical state transitions under electrical stimulation, enabling end-to-end training; (2) A hierarchical training framework integrating neural decoding and phosphene simulation constraints, extending hybrid autoencoder (Granley, Relic, and Beyeler 2022); (3) Background gated network (BGN) for input filtration; (4) Multi-channel activation constraint to encourage a sufficiently large subset of electrode channels to fire concurrently, rendering the stimulation effective and consistent (Oswalt et al. 2021; Bosking et al. 2022); (5) Comprehensive validation through novel primate experiments demonstrating state-of-the-art performance.

Related Work

In biological vision, two frequently explored directions are: (1) visual encoding mechanisms—how visual stimuli elicit neural responses; and (2) visual decoding algorithms—how

*Corresponding authors.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

to reconstruct perceived visual scenes from neural signals. Extensive research has been conducted on the neural encoding mechanisms and response characteristics to visual stimuli in rodents (Turishcheva et al. 2024b,a; Li et al. 2023; Xu et al. 2023; Sinz et al. 2018) and primates (Cadena et al. 2019; Hatanaka et al. 2022; Farzmaahi, Kohn, and Coen-Cagli 2025; Shipp 2024), leading to significant advances and a clearer understanding of these processes. Similarly, substantial progress has also been made in the area of visual decoding. Pixel-level reconstructions of images from brain signals demonstrate how encoding and decoding approaches inform each other (Zhang et al. 2020, 2022). Meanwhile, other studies focus on decoding semantic information (Chen et al. 2023; Ozelik et al. 2022; Gaziv et al. 2022).

Given that the aforementioned studies primarily focus on the relationship between visual input and neural activity, either predicting natural neural responses based on given visual stimuli and brain states, or conversely decoding perceived visual input from neural activity. However, they differ substantially from the goal of visual prosthetics: determining how to appropriately intervene in the visual system to evoke desired percepts. Nonetheless, the mechanisms of visual encoding and decoding still provide important biological insights into how visual perception might be restored.

Early visual prosthetics focused on retinal stimulation, exemplified by systems like Argus II, which translates camera-captured images into electrical signals for retinal neurons (Luo and Da Cruz 2016). However, such devices require intact retinal circuitry, rendering them ineffective for patients with optic nerve damage or advanced degenerative diseases. Cortical visual prosthetics bypass these limitations by directly stimulating visual cortex neurons but initially yielded only coarse visual perceptions (Brindley and Lewin 1968; Dobbelle, Mladejovsky, and Girvin 1974). There is a need to explore more complex pattern presentation (Dobbelle and Mladejovsky 1974; Bosking, Beauchamp, and Yoshor 2017). Recent cortical prosthetic approaches advanced visual perception by either sequentially activating electrodes to mimic natural motion-sensitive processing (Beauchamp et al. 2020), enabling accurate letter recognition in humans, or employing high-channel-count simultaneous stimulation to evoke shape perception in monkeys (Chen et al. 2020). Meanwhile, some deep learning techniques have also been introduced into this field. A deep autoencoder-based architecture that includes a highly adjustable prosthetic vision simulation module attempts to automatically find a task-specific stimulation protocol (van Steveninck et al. 2022). Hybrid neural autoencoder (Granley, Relic, and Beyeler 2022) combines biological constraints with deep learning to generate stimulation that preserve topological relationships in V1. Another ML-based paradigm is to incorporate models of the visual system’s response into the stimulus optimization process (Grani et al. 2022).

However, these methods either remain at the theoretical simulation stage without sufficient animal experiments to validate their effectiveness, or the visual percepts elicited by the methods are imprecise, inflexible, and unclear. To address these issues, this work proposes MindSight.

Methodology

Overview

As shown in Figure 1, MindSight consists of three major stages. The first stage is the MUA-Image Decoder (MID), which decodes neural activity in the form of multi-unit activity (MUA) into corresponding visual images. The second stage is the Electrical Stimulation Generator (ESG), which optimizes the stimulation parameters to generate the visual perceptions as close as possible to the desired perceptions. This stage involves training through two paths of losses. The third stage, the Background Gated Network (BGN), filters out irrelevant visual stimuli, ensuring that electrical stimulation is applied only to stimuli that contain essential visual information. The biologically inspired designs of the individual modules, along with the rationale behind each, will be detailed in the corresponding subsections that follow. To provide an overview of how the algorithm operates as a whole, the high-level algorithmic workflow of the entire framework is presented in **Algorithm 1** in the **Appendix**.

Stage A: MUA-Image Decoder (MID)

In previous work, attempts were made to decode pixel-level images from retinal ganglion cells spikes (Zhang et al. 2020) and two-photon calcium signals of the Macaque Visual Cortex (Zhang et al. 2022). However, these models do not incorporate biological visual mechanisms. To reflect the selective tuning of visual cortical neurons to different visual attributes, the modulation of receptive fields, and the feedback and integration across visual cortical areas, MUA-Image Decoder (see the upper part of Figure 1) introduces three key advancements over previous neural decoders: (1) Spatiotemporal fusion of multi-scale MUA dynamics, (2) Hierarchical attention-guided feature refinement, and (3) Biologically-plausible feature recombination through adaptive skip connections. Our architecture incorporates the Convolutional Block Attention Module (CBAM) (Woo et al. 2018) for cross-dimensional feature refinement. As shown in Figure 1, the overall calculation process of MID is as follows:

$$\hat{I}_i = \mathcal{D}_\psi(\mathcal{E}_\phi(\mathcal{P}_\theta(\text{MUA}_i))) \quad (1)$$

Here, \hat{I}_i denotes the image decoded from MUA (1280 channels \times 50 ms), \mathcal{P}_θ represents the projection network, and \mathcal{E}_ϕ and \mathcal{D}_ψ denote the encoder and decoder of the Image-Image autoencoder, respectively. Refer to **Appendix A** for details.

Stage B: Electrical Stimulation Generator (ESG)

The ESG (see the lower part of Figure 1) is responsible for producing voltage amplitudes on each electrode channel, given a real input image. The ESG parameters would be optimized via two paths (referred to as *Path 1* and *Path 2*) and one channel-activation constraint term, each contributing to the overall loss function:

1. **Path 1: MID Decoding Constraint.** The purpose of this constrained pathway is to ensure that the neural state changes and responses induced by cortical electrical stimulation can be decoded—via the MID trained in stage A—into the intended visual input. This provides

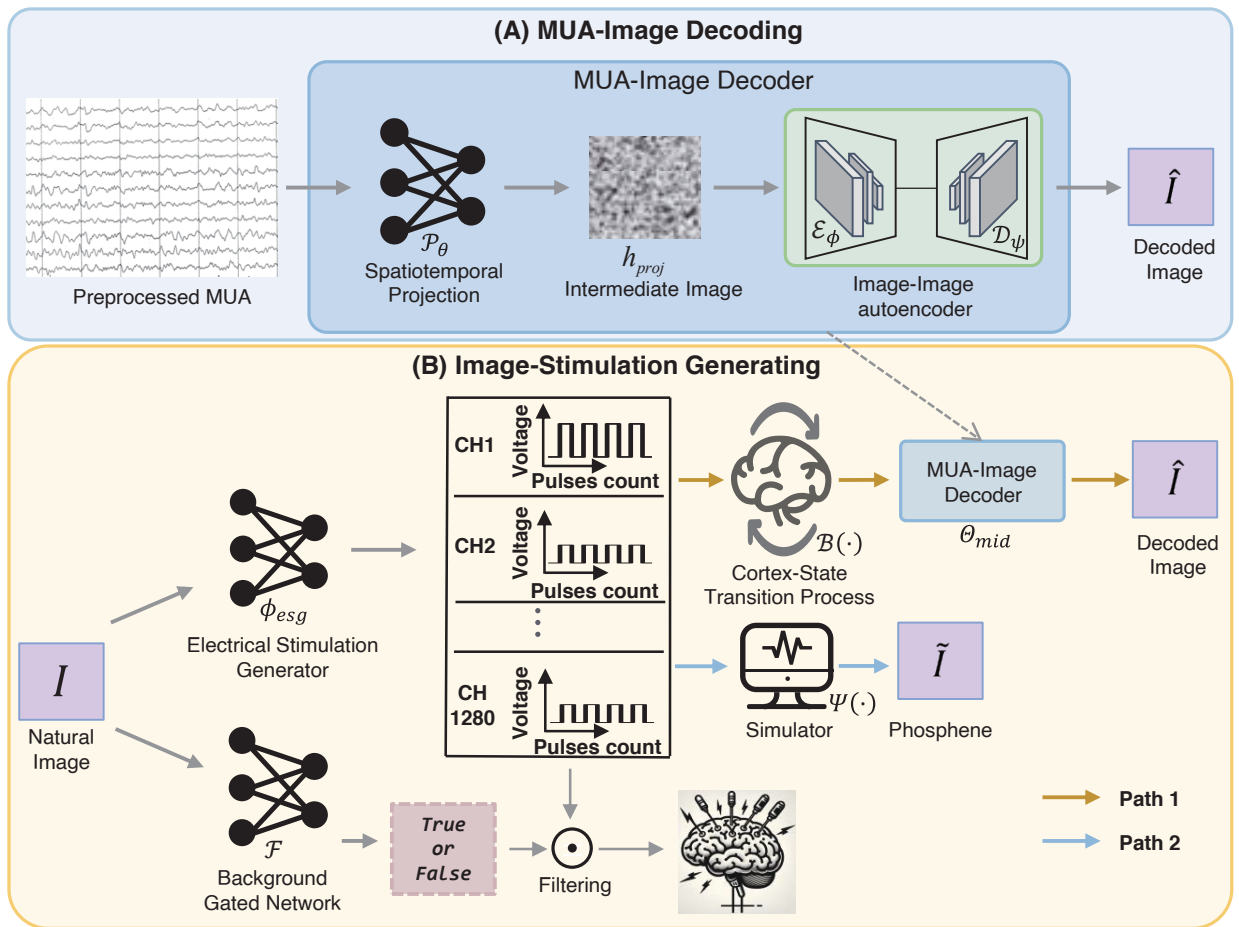


Figure 1: MindSight. Stage A (upper part): In MUA-Image Decoding, MUA-Image Decoder is trained to reconstruct visual stimuli from the multi-unit activity (MUA) recorded in the primary visual cortex, and the well-trained MUA-Image Decoder will be used in stage B. Stage B (lower part): In Image-Stimulation Generating, electrical stimulation generator (ESG) is responsible for producing voltage amplitudes on each electrode channel, given a real input image, and its parameters are trained via two paths. Background gated network (BGN) is trained to filter the external visual inputs, ensuring that only images containing key information will trigger cortical electrical stimulation.

a theoretical guarantee of stimulation effectiveness from the perspective of biological vision systems. Otherwise, how can we ensure that the visual percepts generated by electrical stimulation are meaningfully related to the desired visual inputs?

2. **Path 2: Phosphene Simulation Constraint.** The inclusion of this constrained pathway is also motivated by the need to satisfy established biological priors—namely, the orderly mapping between the visual cortex and phosphene perception. Without this, even at the coarse level of phosphene perception, the correctness of electrical stimulation cannot be ensured.
3. **Channel Activation Constraint.** To ensure that sufficiently many channels within the same electrode array are activated in unison. Without this, stimulating too few electrode channels in a local region may render the stimulation ineffective—that is, no percept may be elicited in the corresponding area of the visual field.

Architecture of ESG To better enable the ESG to learn the relationship between the target visual percept and the corresponding electrical stimulation parameters, it is essential to extract features from different regions of the visual target space and across multiple levels of granularity. At the same time, the design must account for the deployment and computational cost of ESG.

Depthwise-separable convolutions are employed to reduce the size of ESG, aiming to minimize power consumption, as ESG will be integrated into the chip within the visual prosthesis system. Hierarchical processing with depthwise-separable convolutions and attention:

$$e_{out} = \text{AvgPool}(\text{Stack}_{4 \times}[\text{ResBlock}](I)) \quad (2)$$

Use \mathcal{R} to explicitly denote each ResBlock, as shown below:

$$\mathcal{R}(x) = \text{ReLU}\left(\underbrace{C_{DS}^{(2)}(\text{CBAM}(\text{ReLU}(C_{DS}^{(1)}(x))))}_{\text{conv path}} + \underbrace{C_{proj}(x)}_{\text{shortcut}}\right) \quad (3)$$

where $C_{DS}^{(i)}$ denotes depthwise-separable conv with BN.

After adaptive global pooling, we apply two fully connected layers (with LayerNorm and GELU) to predict stimulation parameters:

$$V = \underbrace{W_2 \left(\text{GELU}(\text{LayerNorm}(W_1 e_{\text{out}})) \right)}_{\in \mathbb{R}^{K \times T_{\text{pulse}}}} \quad (4)$$

Here W_1 and W_2 are linear transformations mapping and V represents the voltage amplitudes for the K electrode channels. T_{pulse} denotes the number of pulse for each channel.

Path 1: MID Decoding Constraint As shown in Figure 1, in path 1, we introduce the cortex-state transition process and the previously trained MID to guide ESG: 1. A real input image I_i is fed into the ESG, denoted by $\phi_{\text{esg}}(\cdot)$, yielding a set of voltage amplitudes across all electrode channels. 2. A cortex-state transition process model simulates how these voltage amplitudes (converted to currents) affect cortical neural states over time, producing a simulated neural response. 3. The simulated neural response is passed through the fixed-parameter MID, which reconstructs an image \hat{I}_i .

Let Θ_{mid} be the fixed parameters of MID. Let $\mathcal{B}(\cdot)$ denote the *cortex-state transition process*, modeling the spatiotemporal changes in mua triggered by electrical stimulation. Path 1 can be summarized as:

$$\hat{I}_i = \Theta_{\text{mid}}(\mathcal{B}(\phi_{\text{esg}}(I_i))) \quad (5)$$

Cortex-State Transition Process Suppose each of the K electrodes is located at spatial coordinates (x_k, y_k, z_k) , and let $\phi_{\text{esg}}^k(I_i)$ be the generated voltage amplitude on electrode k . If R_k is the impedance coefficient of channel k , then its injected current is approximately

$$I_k = \frac{\phi_{\text{esg}}^k(I_i)}{R_k} \quad (6)$$

We adopt Spatiotemporal Gaussian distribution to describe how stimulation at (x_k, y_k, z_k) and current I_k diffuses through cortical tissue over time. In the x, y, z directions, we assume that the diffusion rates α are the same, i.e., $\sigma_x(t) = \sigma_y(t) = \sigma_z(t) = \sqrt{\alpha t + \sigma_0^2}$. Specifically, after calculation and simplification, the electrical stimulation for channel k ultimately generated the following distribution:

$$\mathcal{G}_k(\mathbf{p}, t) = \frac{I_k}{(2\pi(\sigma_0^2 + \alpha t))^{\frac{3}{2}}} \exp\left(-\frac{\|\mathbf{p} - (x_k, y_k, z_k)\|^2}{2\sigma_0^2 + 2\alpha t}\right) \quad (7)$$

where \mathbf{p} denotes a spatial coordinate, t represents time, σ_0 is the initial standard deviation, and α is the diffusion rate. Summing contributions from all K electrodes yields

$$\mathbf{S}(\mathbf{p}, t) = \sum_{k=1}^K \mathcal{G}_k(\mathbf{p}, t; I_k, (x_k, y_k, z_k)) \quad (8)$$

thus forming a spatiotemporal distribution of the electric field. By sampling $\mathbf{S}(\mathbf{p}, t)$ at the corresponding recording channels over a certain time window, we obtain the simulated neural response, which is then fed into the MID. Note

that T_{pulse} is a critical hyperparameter and Equations (7) and (8) only describes the case where T_{pulse} is 1. If T_{pulse} is greater than 1, the complex effects introduced by different pulses must also be considered. In our experiments, we set T_{pulse} to 1. If you find that a small number of pulses cannot activate the channel, then T_{pulse} needs to be increased.

Loss Function for Path 1 The goal of path 1 is to ensure that the images reconstructed from simulated neural responses are highly consistent with the real input images, thereby ensuring that the electrical stimulation parameters generated by ESG can induce the desired visual perception. Using the mean squared error at the pixel level, we define

$$\mathcal{L}_{\text{Path1}} = \mathbb{E}[\|I - \hat{I}\|_2^2] \quad (9)$$

Path 2: Phosphene Simulation Constraint In path 2, the same generator $\phi_{\text{esg}}(\cdot)$ is used, but the forward pass relies on a *phosphene simulator* (van der Grinten et al. 2024) (see the lower part of Figure 1), denoted by $\Psi(\cdot)$, and $\Psi(\cdot)$ is not trainable. The pipeline is: 1. A real image I_i is given to the generator $\phi_{\text{esg}}(\cdot)$, producing voltage amplitudes (converted to currents). 2. Generated voltage amplitudes $\phi_{\text{esg}}^k(I_i)$, impedance coefficient R_k and the center positions of the receptive fields $(\tilde{x}_k, \tilde{y}_k)$ for all channels are fed into the phosphene simulator $\Psi(\cdot)$, which outputs a simulated phosphene image \tilde{I}_i . $(\tilde{x}_k, \tilde{y}_k)$ can be obtained through Receptive Field Mapping (see Section **Receptive Field Mapping** for details). Mathematically, we have

$$\tilde{I}_i = \Psi(\phi_{\text{esg}}^k(I_i), R_k, (\tilde{x}_k, \tilde{y}_k)) \quad (10)$$

For each real input I_i , we generate the corresponding binarized target \hat{I}_i (see **Appendix D: Method for Constructing Binarized Targets** for more details). The path 2 loss is then:

$$\mathcal{L}_{\text{Path2}} = \mathbb{E}[\|\hat{I} - \tilde{I}\|_2^2] \quad (11)$$

which penalizes discrepancies between the simulated and target phosphene images to adhere to the orderly mapping between the visual cortex and phosphene perception.

Channel Activation Constraint Our extensive saccade experiments (Tehovnik, Slocum, and Schiller 2003; Tehovnik et al. 2005) on monkeys have shown that, stimulating too few electrode sites within a localized area may render the stimulation ineffective, i.e., no perception was induced in the corresponding visual field region (see Section **Electrical-Stimulation-Induced Saccade** for details). To encourage a sufficiently large subset of electrode channels to fire concurrently, we introduce an additional penalty:

$$\mathcal{L}_{\text{con}} = \gamma \mathbb{E} \left[\sum_k \exp \left(-\frac{\phi_{\text{esg}}^k(I)}{\epsilon R_k} \right) \right] \quad (12)$$

where γ and ϵ are the penalty coefficient and decay coefficient, respectively. This penalty term constrains the number of effective channels that are not stimulated to be as few as possible. In particular, when the physical locations of different channels are closer, the probability of them being stimulated simultaneously becomes higher (rather than stimulating only a subset of these channels, even though stimulating only this subset could already make $\mathcal{L}_{\text{Path2}}$ very small).

Overall Loss for the ESG By combining Path 1, Path 2, and the channel activation constraint, the parameters of ESG are trained by minimizing the total objective using Adam:

$$\mathcal{L}_{\text{gen}} = \lambda_1 \mathcal{L}_{\text{Path1}} + \lambda_2 \mathcal{L}_{\text{Path2}} + \mathcal{L}_{\text{con}}, \quad (13)$$

where λ_1 and λ_2 control relative importance of two paths.

Stage C: Background Gated Network (BGN)

Given that multiple studies (Cogan 2008; Aungaroon et al. 2017; Larkin et al. 2022; Vatsyayan and Dayeh 2022; Vatsyayan et al. 2021) have shown excessive cortical electrical stimulation can increase the risk of neuronal necrosis, tissue damage, electrode failure, and even trigger after-discharges or seizures, we further train a BGN (see the lower part of Figure 1) to filter visual inputs. This ensures that only images containing key information trigger cortical stimulation.

BGN consists of a ResNet backbone modified with CBAM layers to dynamically refine feature maps. Given I_i , BGN indicates whether the image contains key information:

$$p_i = \text{Sigmoid}(\text{FC}(\text{CBAM-ResNet}(I_i))) \quad (14)$$

Experiment

Datasets

In our study, we used two self-collected datasets: **Monkey Neural-Image Viewing Dataset (NIVD)** and **Monkey Vehicle Obstacle Avoidance for Food — Neuro-Visual Dataset (MVOAF-NVD)**. Refer to **Appendix C** for more description of these two datasets. For details on how NIVD and MVOAF-NVD were collected, refer to **Appendix C.1**, **Appendix C.2**, and Supplementary video 1.

Implementation Details

For all hyperparameter choices and implementation details related to data preprocessing and the training of MindSight, please refer to **Appendix B**.

Receptive Field Mapping

We used a fast RF mapping method similar with (Fiorani et al. 2014) to get the receptive fields of neurons. Refer to **Appendix E** for more details.

Electrical-Stimulation-Induced Saccade

We followed previously reported methods (Chen et al. 2020) to conduct electrical-stimulation-induced saccade experiments on our monkeys. We found that at least 16 channels per electrode are required to elicit effective phosphene perception, which motivated the design of Channel Activation Constraint. For more saccade experiments details, refer to **Appendix F** and the corresponding Table S1.

Behavioral Validation Experiments in Monkeys

To validate the efficacy of MindSight, we designed a delayed match-to-sample (DMS) task and a driving-based obstacle-avoidance foraging (DOAF) task. Refer to **Appendix G** and Figure S1 for more description of these two tasks.

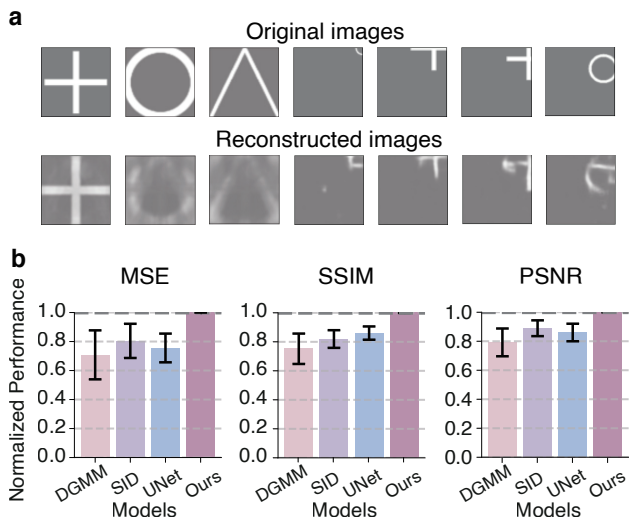


Figure 2: Decoding Performance Comparisons on NIVD Test Set. **a.** Some examples of decoded images. **b.** We evaluated four models (ours, DGMM, SID, and UNet-based) on the NIVD test set using three evaluation metrics: MSE, SSIM, and PSNR. Error bars represent standard deviation of the mean normalized performance.

Results

For MindSight, it is mainly necessary to validate the results of the following modules it contains: 1. Decoding performance of MID, as it indirectly affects the performance of ESG; 2. Phosphenes generated by ESG; 3. Performance of ESG in DMS; 4. Performance of ESG and BGN in DOAF.

Performance of MID

We evaluated the performance of MID on the test set of the NIVD (4,037 cropped images with corresponding MUAs). Figure 2a demonstrates MID’s decoding results on some examples. As shown in Figure 2b, MID was compared with several state-of-the-art pixel-level decoding models—SID (Zhang et al. 2020), UNet-based model (Zhang et al. 2022), and DGMM (Du et al. 2018)—using evaluation metrics including mean square error (MSE) that describes the absolute difference of every pixel, the structural similarity index measure (SSIM) that captures the details or image distortion (Wang et al. 2004), and the peak-signal-to-noise-ratio (PSNR) that characterizes the global quality. For clarity, we utilized normalized model performance. Specifically, each model’s score for every test sample was normalized against MID’s scores, where the ratio for MID is defined as 1. Values above 1 indicate better performance than MID, while values below 1 indicate inferior performance. For the other three models, we calculated the mean and variance of their normalized performance across the entire test set. Compared to the other three models, MID from MindSight exhibited superior performance across all three metrics.

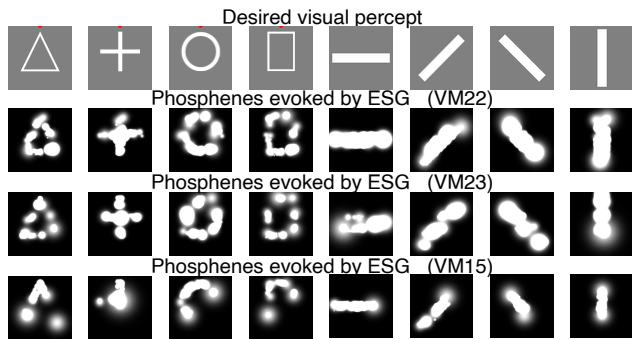


Figure 3: Phosphenes Evoked by ESG in DMS. The first row shows eight possible visual percepts in DMS. The second, third, and fourth rows show the phosphenes that ESG would theoretically evoke for VM22, VM23, and VM15, respectively, when targeting each desired visual percept.

| Visual percepts | VM22 | | VM23 | | VM15 | |
|-----------------|------|------|------|------|------|------|
| | IoU | Dice | IoU | Dice | IoU | Dice |
| Triangle | 0.60 | 0.75 | 0.57 | 0.73 | 0.43 | 0.60 |
| Cross | 0.64 | 0.78 | 0.68 | 0.81 | 0.48 | 0.65 |
| Circle | 0.61 | 0.76 | 0.57 | 0.73 | 0.48 | 0.64 |
| Rectangle | 0.52 | 0.68 | 0.46 | 0.63 | 0.36 | 0.53 |
| bar (0°) | 0.76 | 0.86 | 0.64 | 0.78 | 0.61 | 0.76 |
| bar (45°) | 0.69 | 0.82 | 0.73 | 0.84 | 0.58 | 0.74 |
| bar (135°) | 0.67 | 0.80 | 0.70 | 0.82 | 0.47 | 0.64 |
| bar (90°) | 0.67 | 0.80 | 0.60 | 0.75 | 0.49 | 0.66 |

Table 1: Quantitative Results for Figure 3. For each monkey, we compute the similarity between each desired visual percept and the corresponding phosphene evoked by ESG. Since the focus here is primarily on shape and contour, the original visual percept images are first preprocessed (e.g., binarization) before computing the Intersection over Union (IoU) and Dice Similarity Coefficient with the phosphenes.

Phosphenes Evoked by ESG

We also evaluated the simulated phosphenes evoked by ESG in both DMS and DOAF. As shown in Figure 3, we assessed the theoretical phosphenes that would be induced by the stimulation parameters generated by ESG when aiming to evoke different desired visual percepts in DMS for VM15, VM22, and VM23. The phosphenes evoked by ESG vary across monkeys (Table 1 provides a quantitative summary for each monkey), due to differences in electrode implantation sites, which in turn lead to variations in receptive fields (see **Appendix E** for more details). These differences in electrode and receptive field distributions influence the training of MindSight, as described in Equations (7), (8), and (10). The results for VM15 are significantly worse than those for VM22 and VM23, primarily due to the suboptimal distribution of electrode implantation sites. The phosphenes evoked by ESG in DOAF will be detailed in Section **Performance of ESG and BGN in DOAF**.

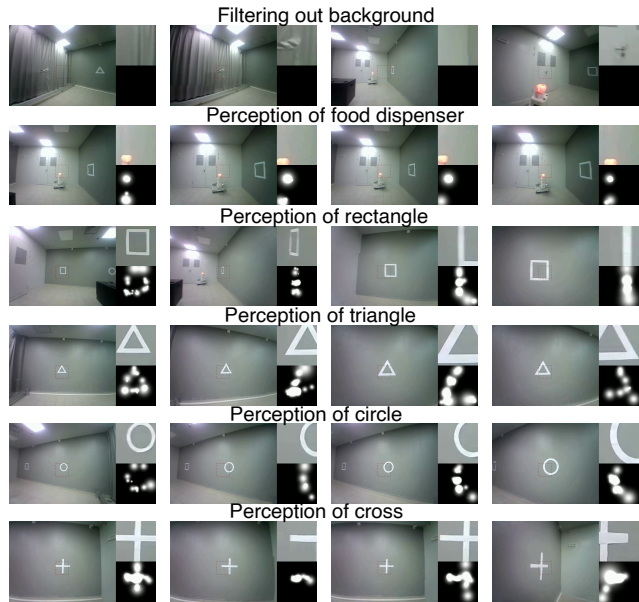


Figure 4: Examples of Phosphenes Evoked by ESG in DOAF. The first row shows the results of background filtering, where irrelevant information such as door handles, curtains, and walls does not trigger electrical stimulation. Rows two through six illustrate the stimulation evoked for the food dispenser and various markers at different angles and distances as the cart moves.

Performance of ESG in DMS

The ESG trained on NIVD was first validated in DMS (as outlined in Section **Behavioral Validation Experiments in Monkeys**). We evaluated model efficacy by measuring the monkey’s behavioral success rate in DMS, where the choice relies entirely on electrical stimulation, with no visual cues provided. Figure S1 displays four-class DMS, with the two-class DMS differing only in the final selection step between two shapes instead of four shapes. Figure 5a shows the dynamic performance of behavioral success rates in a two-class DMS during one session. With electrical stimulation from the ESG, the monkey maintained high accuracy (red region). However, when the stimulation was removed, the success rate immediately dropped to chance level (green region). As shown in Figure 5b, the results of MindSight’s ESG in two-class DMS were compared with those from another Chen’s study (Chen et al. 2020). Both experiments followed nearly identical protocols, except that Chen used some different shapes, such as “T” and “L”. We validated our method using two monkeys (monkeys 22 and 23), while Chen used monkeys A and L, with results statistically aggregated over 10 sessions per monkey. The comparison reveals that our method achieves higher stability. Moreover, compared to (Chen et al. 2020), MindSight also offers greater flexibility, as it can directly adapt to arbitrary inputs without the need for additional manual intervention. Such robustness and flexibility are crucial for the future application of visual prosthetics. We additionally demonstrated the mon-

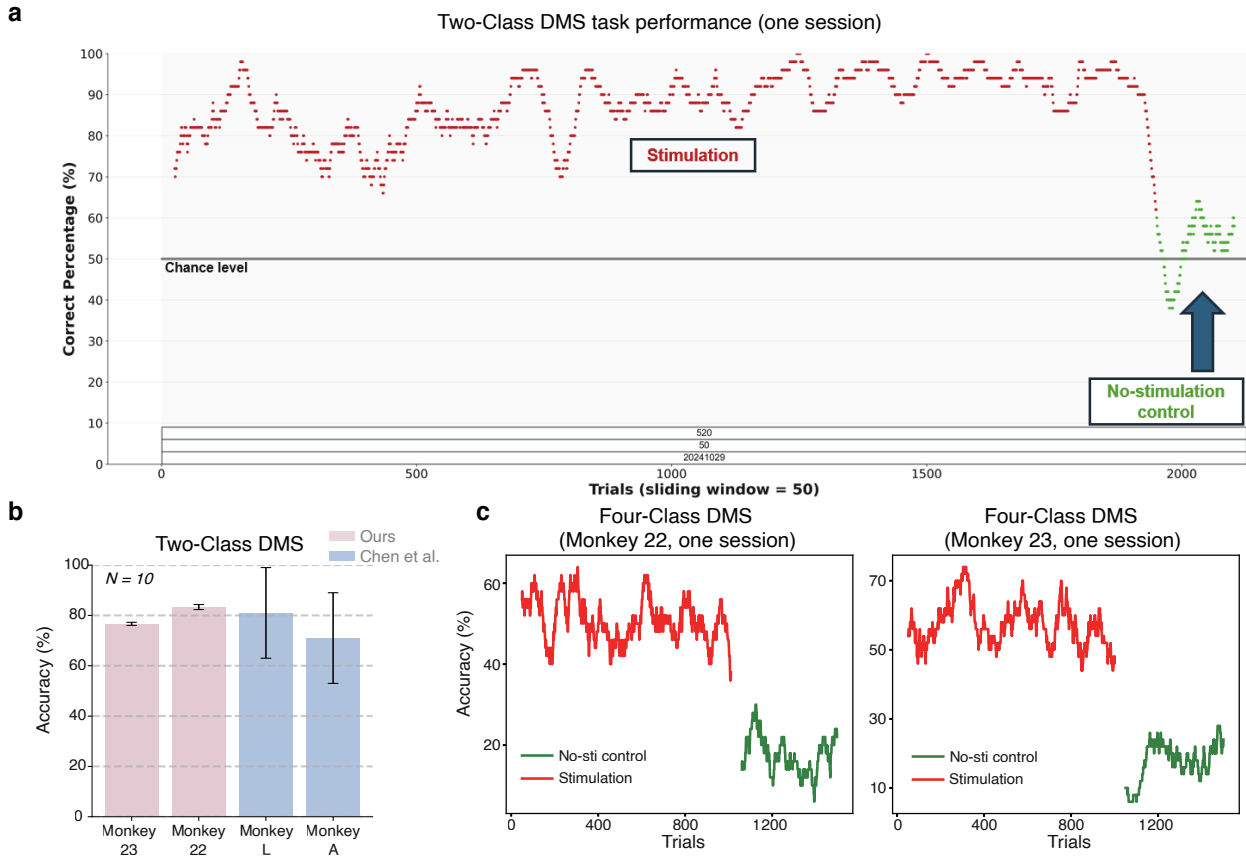


Figure 5: Validation on the DMS Task. **a.** Dynamic performance of MindSight in two-class DMS task during one session. **b.** Performance comparisons in two-class DMS task. Error bars represent the SEM (Standard Error of the Mean). N, number of sessions for each monkey and model. **c.** Dynamic performance of MindSight in four-class DMS task during one session.

keys' performance in four-class DMS (which has not been explored in (Chen et al. 2020)). Figure 5c illustrates the dynamic performance of behavioral success rates in one session for monkeys 22 and 23, respectively. Similarly, a control group without electrical stimulation was set up, confirming that MindSight's correct rate far exceeded the chance level (25%). Supplementary videos 2 and 3 demonstrate examples of MindSight applied to 2-class (Monkey 22) and 4-class (Monkey 23) DMS, respectively, showing the monkeys performing continuous correct trials. See **Appendix H** for details of Supplementary Videos 2–3.

Performance of ESG and BGN in DOAF

The ESG (trained on MVOAF-NVD) was validated in DOAF (as outlined in Section **Behavioral Validation Experiments in Monkeys**). Because a full session is lengthy, Supplementary Video 4 presents only segments highlighting MindSight's performance at key moments (e.g., when the monkey observes markers on walls, food dispensers, or irrelevant background). In the video, left panel shows head-mounted camera footage (on a head-fixed monkey), with the red dashed box representing a 64×64 region corresponding in both size and position to the receptive field covered by

implanted electrodes. Upper-right display shows the content within this red dashed box, which serves as the input to ESG/BGN. Lower-right panel simulates the visual perception evoked by the electrical stimulation commands generated by ESG (Equation (10)), combined with BGN outputs. For easier viewing, Figure 4 presents some frames from Supplementary Video 4 as illustrative examples. BGN-triggered stimulation activates only during key cues (e.g., markers, obstacles, and food dispensers), filtering out irrelevant information to avoid the risks associated with excessive electrical stimulation. Notably, despite variations in the shape, angle, and position of markers and food dispensers as the cart moves, MindSight consistently evoked effective visual perceptions.

Conclusion

We present MindSight, a framework for visual restoration via cortical electrical stimulation. By integrating differentiable biophysical modeling with dual-path training, MindSight outperforms other approaches and shows good performance with comprehensive validation through novel primate experiments. For additional discussion, see **Appendix I**.

Acknowledgments

This work was supported by National Science and Technology Innovation 2030 Major program (No. 2021ZD0203601), National Natural Science Foundation of China (No. 32221003, No. 32221003), Shanghai Municipal Science and Technology Major Project (No. 2018SHZDZX05, No. 2021SHZDZX), National Key R&D Program Key Scientific Issues of Transformational Technology (No. 2019YFA0709504), Shanghai Pilot Program for Basic Research-Chinese Academy of Science, Shanghai Branch (No. JCYJ-SHFY2022-010), Lingang Laboratory (No. LG202105-01-01, No. LG202105-01-11, No. LG-GG-202402-06, No. LGL-5925) and Shanghai Yang Fan Foundation (No. 24YF2730700).

References

- Aungaroon, G.; Vera, A. Z.; Horn, P. S.; Byars, A. W.; Greiner, H. M.; Tenney, J. R.; Arthur, T. M.; Crone, N. E.; Holland, K. D.; Mangano, F. T.; et al. 2017. After-discharges and seizures during pediatric extra-operative electrical cortical stimulation functional brain mapping: incidence, thresholds, and determinants. *Clinical Neurophysiology*, 128(10): 2078–2086.
- Beauchamp, M. S.; Oswald, D.; Sun, P.; Foster, B. L.; Magnotti, J. F.; Niketeghad, S.; Pouratian, N.; Bosking, W. H.; and Yoshor, D. 2020. Dynamic stimulation of visual cortex produces form vision in sighted and blind humans. *Cell*, 181(4): 774–783.
- Beauchamp, M. S.; and Yoshor, D. 2020. Stimulating the brain to restore vision. *Science*, 370(6521): 1168–1169.
- Bosking, W. H.; Beauchamp, M. S.; and Yoshor, D. 2017. Electrical stimulation of visual cortex: relevance for the development of visual cortical prosthetics. *Annual review of vision science*, 3(1): 141–166.
- Bosking, W. H.; Oswald, D. N.; Foster, B. L.; Sun, P.; Beauchamp, M. S.; and Yoshor, D. 2022. Percepts evoked by multi-electrode stimulation of human visual cortex. *Brain stimulation*, 15(5): 1163–1177.
- Brindley, G. S.; and Lewin, W. S. 1968. The sensations produced by electrical stimulation of the visual cortex. *The Journal of physiology*, 196(2): 479–493.
- Cadena, S. A.; Denfield, G. H.; Walker, E. Y.; Gatys, L. A.; Tolia, A. S.; Bethge, M.; and Ecker, A. S. 2019. Deep convolutional models improve predictions of macaque V1 responses to natural images. *PLoS computational biology*, 15(4): e1006897.
- Chen, X.; Wang, F.; Fernandez, E.; and Roelfsema, P. R. 2020. Shape perception via a high-channel-count neuroprosthesis in monkey visual cortex. *Science*, 370(6521): 1191–1196.
- Chen, Z.; Qing, J.; Xiang, T.; Yue, W. L.; and Zhou, J. H. 2023. Seeing beyond the brain: Conditional diffusion model with sparse masked modeling for vision decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22710–22720.
- Cogan, S. F. 2008. Neural stimulation and recording electrodes. *Annu. Rev. Biomed. Eng.*, 10(1): 275–309.
- Dobelle, W.; and Mladejovsky, M. 1974. Phosphenes produced by electrical stimulation of human occipital cortex, and their application to the development of a prosthesis for the blind. *The Journal of physiology*, 243(2): 553–576.
- Dobelle, W. H. 2000. Artificial vision for the blind by connecting a television camera to the visual cortex. *ASAIO journal*, 46(1): 3–9.
- Dobelle, W. H.; Mladejovsky, M. G.; and Girvin, J. 1974. Artificial vision for the blind: electrical stimulation of visual cortex offers hope for a functional prosthesis. *Science*, 183(4123): 440–444.
- Dorn, J. D.; Ahuja, A. K.; Caspi, A.; Da Cruz, L.; Dagnelie, G.; Sahel, J.-A.; Greenberg, R. J.; McMahon, M. J.; Group, A. I. S.; et al. 2013. The detection of motion by blind subjects with the epiretinal 60-electrode (Argus II) retinal prosthesis. *JAMA ophthalmology*, 131(2): 183–189.
- Du, C.; Du, C.; Huang, L.; and He, H. 2018. Reconstructing perceived images from human brain activities with Bayesian deep multiview learning. *IEEE transactions on neural networks and learning systems*, 30(8): 2310–2323.
- Farzmaadi, A.; Kohn, A.; and Coen-Cagli, R. 2025. Relating natural image statistics to patterns of response covariability in macaque primary visual cortex. *Nature Communications*, 16(1): 1–13.
- Fernández, E.; Alfaro, A.; Soto-Sánchez, C.; Gonzalez-Lopez, P.; Lozano, A. M.; Peña, S.; Grima, M. D.; Rodil, A.; Gómez, B.; Chen, X.; et al. 2021. Visual percepts evoked with an intracortical 96-channel microelectrode array inserted in human occipital cortex. *The Journal of clinical investigation*, 131(23).
- Fernandez, E.; and Robles, J. A. 2024. Advances and challenges in the development of visual prostheses. *PLoS biology*, 22(10): e3002896.
- Fiorani, M.; Azzi, J. C.; Soares, J. G.; and Gattass, R. 2014. Automatic mapping of visual cortex receptive fields: a fast and precise algorithm. *Journal of neuroscience methods*, 221: 112–126.
- Foroushani, A. N.; Pack, C. C.; and Sawan, M. 2018. Cortical visual prostheses: from microstimulation to functional percept. *Journal of neural engineering*, 15(2): 021005.
- Gaziv, G.; Bely, R.; Granot, N.; Hoogi, A.; Strappini, F.; Golan, T.; and Irani, M. 2022. Self-supervised natural image reconstruction and large-scale semantic classification from brain activity. *NeuroImage*, 254: 119121.
- Grani, F.; Soto-Sánchez, C.; Fimia, A.; and Fernández, E. 2022. Toward a personalized closed-loop stimulation of the visual cortex: Advances and challenges. *Frontiers in Cellular Neuroscience*, 16: 1034270.
- Granley, J.; Relic, L.; and Beyeler, M. 2022. Hybrid neural autoencoders for stimulus encoding in visual and other sensory neuroprostheses. *Advances in Neural Information Processing Systems*, 35: 22671–22685.
- Hatanaka, G.; Inagaki, M.; Takeuchi, R. F.; Nishimoto, S.; Ikezoe, K.; and Fujita, I. 2022. Processing of visual statistics of naturalistic videos in macaque visual areas V1 and V4. *Brain Structure and Function*, 227(4): 1385–1403.

- Larkin, C. J.; Yerneni, K.; Karras, C. L.; Abecassis, Z. A.; Zhou, G.; Zelano, C.; Selner, A. N.; Templer, J. W.; and Tate, M. C. 2022. Impact of intraoperative direct cortical stimulation dynamics on perioperative seizures and afterdischarge frequency in patients undergoing awake craniotomy. *Journal of neurosurgery*, 137(6): 1853–1861.
- Lewis, P. M.; Ackland, H. M.; Lowery, A. J.; and Rosenfeld, J. V. 2015. Restoration of vision in blind individuals using bionic devices: a review with a focus on cortical visual prostheses. *Brain research*, 1595: 51–73.
- Li, B. M.; Cornacchia, I. M.; Rochefort, N. L.; and Onken, A. 2023. V1t: large-scale mouse v1 response prediction using a vision transformer. *arXiv preprint arXiv:2302.03023*.
- Lozano, A.; Suárez, J. S.; Soto-Sánchez, C.; Garrigós, J.; Martínez-Alvarez, J. J.; Ferrández, J. M.; and Fernández, E. 2020. Neurolight: A deep learning neural interface for cortical visual prostheses. *International journal of neural systems*, 30(09): 2050045.
- Luo, Y. H.-L.; and Da Cruz, L. 2016. The Argus® II retinal prosthesis system. *Progress in retinal and eye research*, 50: 89–107.
- Maghami, M. H.; Sodagar, A. M.; Lashay, A.; Riazi-Esfahani, H.; and Riazi-Esfahani, M. 2014. Visual prostheses: The enabling technology to give sight to the blind. *Journal of ophthalmic & vision research*, 9(4): 494.
- Oswald, D.; Bosking, W.; Sun, P.; Sheth, S. A.; Niketeghad, S.; Salas, M. A.; Patel, U.; Greenberg, R.; Dorn, J.; Pouratian, N.; et al. 2021. Multi-electrode stimulation evokes consistent spatial patterns of phosphenes and improves phosphene mapping in blind subjects. *Brain stimulation*, 14(5): 1356–1372.
- Ozcelik, F.; Choksi, B.; Mozafari, M.; Reddy, L.; and VanRullen, R. 2022. Reconstruction of perceived images from fmri patterns and semantic brain exploration using instance-conditioned gans. In *2022 international joint conference on neural networks (IJCNN)*, 1–8. IEEE.
- Roelfsema, P. R.; Denys, D.; and Klink, P. C. 2018. Mind reading and writing: the future of neurotechnology. *Trends in cognitive sciences*, 22(7): 598–610.
- Shipp, S. 2024. Computational components of visual predictive coding circuitry. *Frontiers in Neural Circuits*, 17: 1254009.
- Sinz, F.; Ecker, A. S.; Fahey, P.; Walker, E.; Cobos, E.; Froudarakis, E.; Yatsenko, D.; Pitkow, Z.; Reimer, J.; and Tolias, A. 2018. Stimulus domain transfer in recurrent models for large scale cortical population prediction on video. *Advances in neural information processing systems*, 31.
- Stevens, G. A.; White, R. A.; Flaxman, S. R.; Price, H.; Jonas, J. B.; Keeffe, J.; Leasher, J.; Naidoo, K.; Pesudovs, K.; Resnikoff, S.; et al. 2013. Global prevalence of vision impairment and blindness: magnitude and temporal trends, 1990–2010. *Ophthalmology*, 120(12): 2377–2384.
- Stingl, K.; Bartz-Schmidt, K. U.; Besch, D.; Chee, C. K.; Cottrill, C. L.; Gekeler, F.; Groppe, M.; Jackson, T. L.; MacLaren, R. E.; Koitschev, A.; et al. 2015. Subretinal visual implant alpha IMS—clinical trial interim report. *Vision research*, 111: 149–160.
- Tehovnik, E. J.; and Slocum, W. M. 2013. Electrical induction of vision. *Neuroscience & Biobehavioral Reviews*, 37(5): 803–818.
- Tehovnik, E. J.; Slocum, W. M.; Carvey, C. E.; and Schiller, P. H. 2005. Phosphene induction and the generation of saccadic eye movements by striate cortex. *Journal of neurophysiology*, 93(1): 1–19.
- Tehovnik, E. J.; Slocum, W. M.; and Schiller, P. H. 2003. Saccadic eye movements evoked by microstimulation of striate cortex. *European Journal of Neuroscience*, 17(4): 870–878.
- Turishcheva, P.; Fahey, P.; Vystrčilová, M.; Hansel, L.; Froebe, R.; Ponder, K.; Qiu, Y.; Willeke, K.; Bashiri, M.; Baikulov, R.; et al. 2024a. Retrospective for the Dynamic Sensorium Competition for predicting large-scale mouse primary visual cortex activity from videos. *Advances in Neural Information Processing Systems*, 37: 118907–118929.
- Turishcheva, P.; Fahey, P. G.; Vystrčilová, M.; Hansel, L.; Froebe, R.; Ponder, K.; Qiu, Y.; Willeke, K. F.; Bashiri, M.; Wang, E.; et al. 2024b. The dynamic sensorium competition for predicting large-scale mouse visual cortex activity from videos. *ArXiv*, arXiv-2305.
- van der Grinten, M.; van Steveninck, J. d. R.; Lozano, A.; Pijnacker, L.; Rueckauer, B.; Roelfsema, P.; van Gerven, M.; van Wezel, R.; Güçlü, U.; and Güçlütürk, Y. 2024. Towards biologically plausible phosphene simulation for the differentiable optimization of visual cortical prostheses. *Elife*, 13: e85812.
- van Steveninck, J. d. R.; Güçlü, U.; van Wezel, R.; and van Gerven, M. 2022. End-to-end optimization of prosthetic vision. *Journal of Vision*, 22(2): 20–20.
- Vatsyayan, R.; Cleary, D.; Martin, J. R.; Halgren, E.; and Dayeh, S. A. 2021. Electrochemical safety limits for clinical stimulation investigated using depth and strip electrodes in the pig brain. *Journal of neural engineering*, 18(4): 046077.
- Vatsyayan, R.; and Dayeh, S. A. 2022. A universal model of electrochemical safety limits in vivo for electrophysiological stimulation. *Frontiers in neuroscience*, 16: 972252.
- Wang, Z.; Bovik, A. C.; Sheikh, H. R.; and Simoncelli, E. P. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4): 600–612.
- Woo, S.; Park, J.; Lee, J.-Y.; and Kweon, I. S. 2018. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, 3–19.
- Xu, A.; Hou, Y.; Niell, C.; and Beyeler, M. 2023. Multi-modal deep learning model unveils behavioral dynamics of V1 activity in freely moving mice. *Advances in neural information processing systems*, 36: 15341–15357.
- Zhang, Y.; Bu, T.; Zhang, J.; Tang, S.; Yu, Z.; Liu, J. K.; and Huang, T. 2022. Decoding pixel-level image features from two-photon calcium signals of macaque visual cortex. *Neural Computation*, 34(6): 1369–1397.
- Zhang, Y.; Jia, S.; Zheng, Y.; Yu, Z.; Tian, Y.; Ma, S.; Huang, T.; and Liu, J. K. 2020. Reconstruction of natural visual scenes from neural spikes with deep neural networks. *Neural Networks*, 125: 19–30.