

# A CLOSED-LOOP VISUAL STIMULATION FRAMEWORK VIA EEG-BASED CONTROLLABLE GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent advancements in artificial neural networks (ANNs) have greatly enhanced the ability to predict neural activity in response to visual stimuli. However, the inverse problem of designing visual stimuli to elicit specific neural responses remains challenging due to high experimental costs, the high dimensionality of stimuli, and incomplete understanding of neural selectivity. To address these limitations, we present a **closed-loop visual stimulation framework via electroencephalography (EEG)-based controllable generation**. It can iteratively generate the optimal visual stimuli to achieve the goal of controlling brain signals. This framework employs an EEG encoder, treated as a non-differentiable black-box model, to predict neural responses evoked by visual stimuli. By utilizing this encoder (or human experiment), we can quantify the similarity between the predicted (or recorded) neural responses and target neural states. Combining EEG feature extraction with a generation/retrieval module, the framework systematically explores large-scale natural image spaces to identify stimuli that optimally align with the desired brain state. Experimental results demonstrate that, irrespective of the precision of ANN-predicted brain activity, our framework efficiently converges to the theoretically optimal stimulus within a fixed number of iterations. Moreover, this framework generalizes effectively across diverse target neural activity patterns, underscoring its robustness and potential for broader applications in brain-inspired stimulus design. Our code is available at <https://anonymous.4open.science/status/closed-loop-F2E9>.

## 1 INTRODUCTION

The visual system exhibits selectivity, meaning different visual stimuli evoke distinct neural responses (Epstein & Kanwisher, 1998; Qiu et al., 2023). This property suggests that visual stimuli could, in principle, be designed to elicit specific neural responses, offering a novel, non-invasive approach to neuromodulation. Such neuromodulation technique offers several advantages: it is user-friendly, natural, and inherently well-aligned with human sensory processing. However, achieving precise neuromodulation through visual stimuli is highly challenging due to the high dimensionality of visual input space and our incomplete understanding of neuronal selectivity in visual system. Recent advances in controllable image generation techniques have enabled the creation of images with specific semantic attributes, typically conditioned on textual descriptions (Li et al., 2019; Epstein et al., 2023; Wei et al., 2024). While this represents a significant technological breakthrough, current methods lack the ability to conditionally generate stimuli based on neural states. To address this limitation, it is essential to develop frameworks capable of generating visual stimuli specifically optimized to modulate neural activity in a targeted manner, paving the way for more effective and precise neuromodulation through visual stimulation.

Many efforts have focused on precise control of brain activity through visual stimulation. For example, several works (Ponce et al., 2019; Walker et al., 2019; Bashivan et al., 2019) have sought to regulate neural activity at the neuronal level using targeted visual inputs. Notably, (Ponce et al., 2019) introduced a closed-loop experimental framework that integrates a deep generative neural network (GAN) with neurofeedback to iteratively generate images optimized to maximize the responses of specific neurons in the visual system. Despite their success in monkey experiments, these methods often lack generalizability and fail to capture the full diversity of visual features due to the small number of trials and constraints inherent in animal experiments. Moreover, they primar-

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

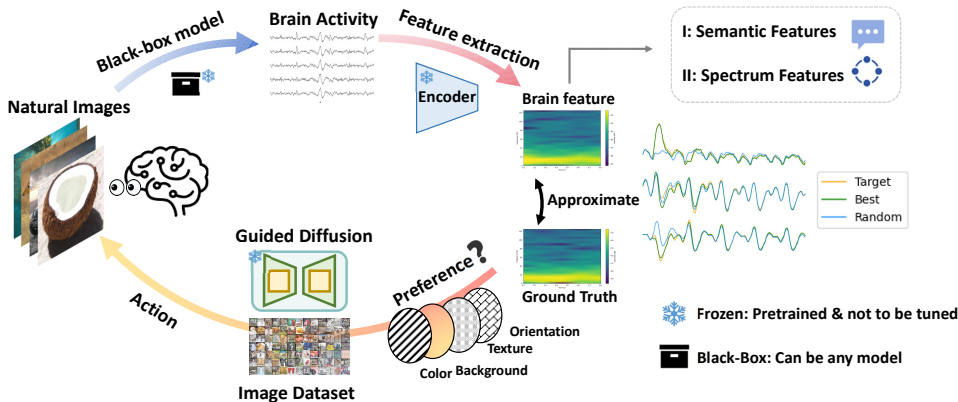


Figure 1: **Conceptualization.** The closed-loop visual stimulation framework includes three core components. (1) The *Black-box model* is used as a surrogate brain to generate neural responses to visual stimulation, and can be replaced by EEG data recorded from human participants in real closed-loop experiments. (2) The *Encoder* extracts the brain features associated with the target neural activity, which can be designed flexibly according to specific control goals. (3) The controllable image generator *Guided diffusion* synthesizes several candidate images. Through closed-loop iteration, the system continuously refines the visual stimulation to achieve the desired brain response.

ily focus on optimizing stimuli for individual neurons, which cannot reflect the complex, distributed neural coding patterns observed at a macroscopic scale, such as those captured in EEG signals. More recently, (Luo et al., 2024b) introduced the Visual Evoked Potential (VEP) Booster, a closed-loop framework designed to generate EEG biomarkers through visual stimulation. However, the VEP Booster primarily generates stroboscopic visual stimuli, rather than natural images that align with the known selectivity principles of the visual system (e.g., preferences for faces, objects, or semantic categories). Therefore, it is crucial for a closed-loop neuromodulation framework that uses natural image stimuli, capable of both flexibly controlling EEG signals and respecting the brain’s inherent selectivity.

In this work, we develop a flexible closed-loop visual stimulation framework designed to achieve controllable EEG responses, as illustrated in Figure 1. By leveraging existing natural image datasets (Hebart et al., 2019) and pre-trained image generation models (Rombach et al., 2022), we utilize state-of-the-art diffusion models to identify fine-grained brain functional specializations in a data-driven manner. Our contributions are summarized as follows:

- We introduce a cutting-edge closed-loop visual neurofeedback framework that synthesizes natural images to control brain activity signatures. Our framework establishes a causal mapping between synthetic visual stimuli and specific EEG features in visual regions.
- By replacing traditional human EEG experiments with a black-box model (serving as a surrogate brain to predict neural responses to stimuli), we minimize dataset biases and enhance the model’s ability to generalize to novel stimuli, providing valuable insights for future human subject experiments.
- We leverage state-of-the-art diffusion models to identify fine-grained visual selectivity, incorporating natural image priors to improve generalization. It allows for flexible design according to specific control goals, such as image retrieval to approximate neural activity generated by a reference image.

## 2 RELATED WORK

**Mapping Selectivity and Invariance from EEG.** Modern neuroscience posits that specific regions of the brain exhibit distinct sensitivities or preferences for particular types of stimuli (Tesileanu et al., 2022). This phenomenon, known as *selectivity*, describes how neurons or neural networks in these regions respond strongly to specific visual inputs. For instance, (Luo et al., 2024a) highlights

108 cases where neurons demonstrate pronounced selectivity for particular stimuli, underscoring their  
109 preference for specific visual features. In contrast, *invariance* refers to the brain’s ability to respond  
110 consistently to distinct stimuli that convey the same information. In other words, different stimuli  
111 can elicit similar patterns of brain activity (Baroni et al., 2023). To explore the intrinsic invariance  
112 shared by ANNs and the brain, (Feather et al., 2023) proposed a method for generating model-  
113 equivalent stimuli, also known as model Metamers. Metamers evoke identical neuronal activations  
114 as a reference stimulus, providing a robust framework to examine the internal states of AI models  
115 and their alignment with neural processes. This approach provides critical insights into the shared  
116 computational principles underlying how artificial and biological systems process and represent in-  
117 formation.

118 **Closed-loop Control of Brain Activity via Visual Stimulation.** Neuromodulation through visual  
119 stimulation holds significant promise for understanding neural mechanisms and developing treat-  
120 ments for various neurological disorders. For example, 40-Hz light flicker, which entrains gamma  
121 oscillations in the brain, has shown potential in treating Alzheimer’s disease (Iaccarino et al., 2016;  
122 Martorell et al., 2019), while visual stimulation by natural images has been explored for improving  
123 mood in patients with depression and anxiety disorders (Mizumoto et al., 2024). A key approach  
124 in this field is the closed-loop control of brain activity, which allows for the real-time regulation of  
125 neural responses through continuous monitoring and feedback. Recent advances in generative mod-  
126 els like GAN and diffusion have enabled the generation of optimal visual stimuli to achieve specific  
127 control of brain activity. For example, (Bashivan et al., 2019) applied gradient ascent to maximize  
128 the activity of the target neuron population with visual stimuli generated by a GAN-based image gen-  
129 erator. Similarly, (Walker et al., 2019) proposed the “inception loops” paradigm, combining in vivo  
130 neural recordings with in silico modeling to synthesize visual stimuli that evoke desired neuronal re-  
131 sponses. (Pierzchlewicz et al., 2024) developed a method to generate images using energy guidance  
132 to maximally activate neuronal responses in the V4 region of monkeys. More recently, (Luo et al.,  
133 2024b) employed a closed-loop strategy where a trained generative model iteratively refined VEP-  
134 EEG biomarkers. These advancements underscore the potential of closed-loop visual stimulation in  
135 precisely modulating brain activity.

136 **Brain-conditioned Controllable Image Generation** Traditional controllable image generation  
137 is typically conditioned on text, where the generation of images is guided by specific textual de-  
138 scriptions (Li et al., 2019; Epstein et al., 2023). In contrast, *brain-conditioned controllable image*  
139 *generation* directly uses the brain’s neural activity, such as EEG, to guide the image generation  
140 process. A key technique in this field is the gradient-based method, which has become crucial for  
141 optimizing visual stimulus guided by brain activity (Luo et al., 2024b;a). This method involves it-  
142 eratively refining visual stimuli by backpropagating the gradients of neural activity representations,  
143 allowing the brain states to be steered toward desired conditions or to achieve specific cognitive  
144 outcomes. This approach enables precise, adaptive stimulus optimization in response to real-time  
145 neural feedback, forming the foundation for personalized brain modulation. Recent advances have  
146 expanded the scope of gradient-based techniques by integrating more sophisticated neural encod-  
147 ing models and utilizing high-dimensional neural representations captured by various brain imaging  
148 modalities (Gu et al., 2023). These developments have significantly improved image generation,  
149 accounting for individual variability in neural responses. Moreover, the incorporation of deep learn-  
150 ing models, such as guided diffusion models (Ye et al., 2023), has enabled the generation of highly  
151 detailed and context-specific stimuli, tailored to align closely with target neural states. These ad-  
152 vancements represent a significant step forward in the field of brain-conditioned image generation.

### 153 3 METHOD

154  
155 In this study, we develop a closed-loop framework to control brain activity through visual stim-  
156 ulation. The visual stimuli are generated by controllable generation models, conditioned on the  
157 EEG signals predicted by EEG encoder Figure 1. The framework is illustrated in Figure 2A. This  
158 closed-loop system is highly adaptable, allowing for the execution of various control objectives.  
159 For instance, by designing a control goal to minimize the distance between the EEG representation  
160 induced by the visual stimulus and a reference EEG representation (e.g., from a seen image), the  
161 system can perform a retrieval task (Figure 2B). Alternatively, by minimizing the distance in the  
power spectral density (PSD) features of the EEG, the system can implement a EEG-conditioned

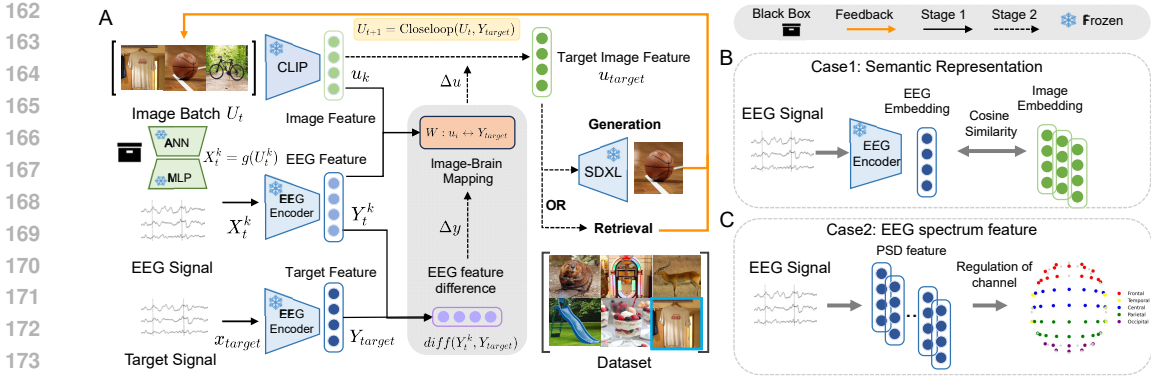


Figure 2: **Closed-loop visual stimulation framework via EEG-based controllable generation.**

(A) We employ a closed-loop iterative process to approximate neural representations derived from EEG signals  $X$ . The encoding model  $g$ , which maps images to synthetic EEG, is designed as a black-box model to broadly simulate the process of regulating brain responses  $Y$ . The EEG Encoder  $f$  is tailored to accommodate various neural features  $U$ . The image with a higher brain similarity score  $\text{sim}(u_j, u_{\text{target}})$  is retained and passed back to the image generator to generate optimized stimuli with a natural image. (B) Example of semantic feature extraction from a pre-trained EEG encoder  $f$ , aligned with CLIP embedding. In this case, our algorithm performs a retrieval task to identify the optimal image  $u_i$  that best matches the  $u_{\text{target}}$ . (C) Channel-wise energy feature using Power Spectral Density (PSD) features. Generative models are iteratively applied to modify the images. For more details, refer to Section 3.1.

generation task (Figure 2C). To support these tasks, we design two distinct feature extractors: one for retrieval and one for generation. If the system begins to favor images with specific colors or textures, it will recognize the relevance of these features to the target class and assign them higher weight in subsequent iterations. Through this closed-loop iterative process, the system can continually optimize the visual stimulus to better elicit the desired EEG responses.

### 3.1 CLOSED-LOOP FRAMEWORK

We formulate the EEG signals as  $X \in \mathbb{R}^{C \times T}$ , where  $C$  is the number of EEG channels and  $T$  represents the length of data points. The image set, containing  $N$  images, is denoted as  $\Omega$ , with each image labeled sequentially with  $1, 2, \dots, N$  for simplicity. Concurrently, we use the encoding model  $g$  to predict brain activity signal  $X = g(U) \in \mathbb{R}^{N \times C \times T}$ . Our objective is to derive brain activity embeddings  $Y = f(g(U)) \in \mathbb{R}^{N \times F}$  from the images  $I \in \mathbb{R}^{N \times 3 \times H \times W}$ , where  $f$  is the feature mapping function from  $X$  to  $Y$ ,  $U$  is the set of stimulus images set, and  $F$  represents the dimension of embedding. Our iteration process can be approximated as a value-based iterative Markov Decision Process (MDP). The state is represented as the probability distribution of each image  $P(u)$  in the image database belonging to target category  $u_{\text{target}}$ . The state updated after each iteration corresponds to a state transition in the MDP. In each iteration, the framework determines which image to select, represented as an action in the MDP. In our model, let  $j \in \llbracket 1, N \rrbracket$ , the reward is defined as the similarity score between the selected or generated image  $u_i$  from database and the features of the target category  $u_{\text{target}}$ :

$$\text{sim}\langle u_j, u_{\text{target}} \rangle = \frac{f(g(u_j)) \cdot f(g(u_{\text{target}}))}{\|f(g(u_j))\| \|f(g(u_{\text{target}}))\|} \quad (1)$$

Let  $u_i$  be any image in the search space, which is the target of model evaluation. During the iteration of the  $t$  to  $t + 1$  step, we update  $S_{t+1}(u_i)$  based on  $u_i$ . The weight coefficient  $\alpha$  controls the cumulative probability increment. Let  $u_+$  be the image that the system considers to be closest to the target category by computing EEG feature similarity. For the history subset  $H$  of selected images  $k$ , the posterior probability that  $u_i$  is the most similar to the target image is updated as follows:

$$S_{t+1}(u_i) = \alpha \cdot S_t(u_i) + (1 - \alpha) \cdot \frac{\exp(s(u_+, u_i))}{\sum_{k=1}^H \exp(s(u_+, u_i))} \cdot S_t(u_i) \quad (2)$$

where  $s$  is the cosine similarity of CLIP (Radford et al., 2021) embedding. The update probability  $P_{t+1}(u_i)$  for  $u_i$  is computed by normalizing the exponentiated value of the updated score  $S_{t+1}(u_i)$  over the sum of exponentiated scores for all  $u_j$  in the dataset, ensuring that the probabilities across all  $u_i$  sum to 1:

$$P_{t+1}(u_i) = \frac{\exp(S_{t+1}(u_i))}{\sum_{j=1}^N \exp(S_{t+1}(u_j))} \quad (3)$$

In step  $t$  iteration, our framework operates as follows. First, we initialize a set of random images  $U_0 = \{u_1, u_2, \dots, u_j\}$ . Using the pretrained encoding model  $g$  to synthesize EEG signals  $X_i$  from these stimuli. Second, for any given representation function  $Y_i$ , we calculate the neural activity representation  $Y_i = f(g(U_i)) \in \mathbb{R}^{N \times F}$  from the predicted signal  $x_i$ , to estimate the difference based on the target neural representation  $Y_{target}$ . Third, the similarity score  $\text{sim}\langle u_j, u_{target} \rangle$  between each neural representation derived from each current stimulus  $u_j$  and the target representation is computed. Subsequently, stimulus images exhibiting higher similarity scores are more likely to be selected. Based on  $\text{sim}\langle u_j, u_{target} \rangle$ , stimulation is probabilistically sampled, favoring images that are closer to the target representation. Finally, the sampled images are used to retrieve similar images for the step  $t + 1$  or input into the diffusion model to generate new stimulus samples.

### 3.2 BLACK-BOX ENCODING MODEL

Instead of recording real EEG data, we employ a pre-trained EEG encoder  $g_\theta$ , treated as a black-box model, to map an image  $I_i \in \mathbb{R}^{3 \times H \times W}$  to a synthetic EEG  $X_i$ . This model predicts the EEG responses corresponding to the visual stimulus. The predicted EEG response can be substituted with actual EEG recordings obtained from human participants during experimental settings. The EEG encoder involves a pre-trained image feature extractor to obtain image embedding aligned with EEG embedding, and a regression model to generate EEG signals from the embedding representation. To test the robustness and generalizability of EEG encoder, we implement two CNN models as image feature extractors, including AlexNet (Krizhevsky, 2014) and CORnet-S (Kubilius et al., 2019). We then train regression models, denoted as  $\hat{X}$ , to predict the neural response according to the image features using supervised learning with the ground truth EEG (from image-EEG paired data).

In the encoding model, we modify the output layer of the CNN, replacing its 1000-neuron configuration with a  $C \times T$ -neuron layer, where each neuron corresponds to one of the flattened EEG data points  $C \times T$ . Each subject is associated with unique model parameters, which are obtained via pretrained models, applied across all EEG time points  $T$ . Given the input training images  $I$  and their corresponding target EEG data  $\hat{X}$ , the model updates its weights by minimizing the mean squared error (MSE) between predicted EEG  $X$  and the target EEG  $\hat{X}$ . This setup ensures a personalized and accurate prediction of synthetic neural activity. This framework ensures a personalized and precise prediction of synthetic neural activity.

### 3.3 INTERACTIVE SEARCH

To identify the optimal stimulus that elicits the desired neural activity, we search for images that generate EEG features similar to the target. The target query image is unknown, and the corresponding EEG feature is observable. To address the challenge of initiating retrieval without a clear query image, we use the mathematical framework of (Ferecatu & Geman, 2007), based on mind matching. It begins with a random sample of images, and through iterative steps, the user selects the image that most closely aligns with the intended category. In our case, this process is adapted to match the target neural feature. The detailed algorithmic procedure is outlined in Algorithm 1, which effectively identifies an optimal subset of images that maximizes the similarity score with respect to the target EEG feature.

In our framework, the *Closed-loop Retrieval Iteration Algorithm* functions as a sequence of state transitions aimed at maximizing the similarity between the current neural feature and the target. The process begins with a randomly selected set of images  $U_0$ , without prior knowledge of the specific features of the target image. We use a roulette wheel selection algorithm to choose from current images based on the similarity measure  $\text{sim}\langle u_j, u_{target} \rangle$ . The system updates the probability  $p_i(u_j)$  for each image in the database belonging to the target class, based on the response model’s prediction  $Y = f(g(U)) \in \mathbb{R}^{N \times F}$ . Subsequently, the system calculates the distance between the



**Algorithm 2** Closed-loop Generative Iteration Algorithm

- 1: **Initialize:** Set initial set  $U_0 = \{u_1, u_2, \dots, u_k\}$ , where  $U_0 \subseteq \Omega$ .
- 2: **repeat**
- 3:   **Selection:**  $U_t = \{u_1, u_2, \dots, u_k\}$  from  $\Omega$  based on  $p_t(u)$ .
- 4:   **Sampling:** Based on the calculated similarity scores, sample from  $U_t$  using:

$$P(u_k) = \frac{\exp(\text{sim}(u_k, u_{\text{target}}))}{\sum_{u_{k'} \in U_t} \exp(\text{sim}(u_{k'}, u_{\text{target}}))}$$

where  $P(u_k)$  is the sampling probability for each  $u_k \in U_t$ .

- 5:   **Crossover:** Draw two distinct samples  $u_a, u_b$  from  $U_t$  based on  $P(u_k)$ , and output new samples by combining the partial embedding of  $u_a$  and  $u_b$ :

$$F(u_{\text{tmp}}^{(1)}) \leftarrow \alpha \cdot F(u_a) + (1 - \alpha) \cdot F(u_b)$$

$$F(u_{\text{tmp}}^{(2)}) \leftarrow \alpha \cdot F(u_b) + (1 - \alpha) \cdot F(u_a)$$

where  $\alpha$  is a crossover control factor.

- 6:   **Mutation:** Based on  $P(u_k)$ , apply mutation to the drawn images  $u_c$  from  $U_t$ , and another image  $u_d$  is drawn from the remaining  $U_t$  (i.e.,  $U_t \setminus \{u_c\}$ ):

$$F(u_{\text{tmp}}^{(3)}) \leftarrow \beta \cdot F(u_c) + (1 - \beta) \cdot F(u_d)$$

where  $\beta$  is a mutation control factor.

- 7:   **Generation:** Generate a new set of images  $U_{\text{gen}} = \{u_{\text{gen}}^{(1)}, u_{\text{gen}}^{(2)}, u_{\text{gen}}^{(3)}\}$  according to the outputs of crossover and mutation phase.
- 8:   **Selection:** Combine  $U_{\text{gen}}$  with  $U_t$  and randomly selected samples  $U_{\text{random}} = \{u_{\text{ran}}^{(1)}, u_{\text{ran}}^{(2)}, \dots, u_{\text{ran}}^{(n)}\}$ , where  $U_0 \subseteq \Omega$ .
- 9:   **Update Action Set:** Update the subset  $U_{t+1}$ :

$$U_{t+1} \leftarrow \{U_t, U_{\text{gen}}, U_{\text{random}}\}$$

- 10:   Replace the old population with the new set of images  $U_{i+1}$ .
- 11: **until** similarity score converges or reach the maximum number of cycles.

utilized the Adam optimizer with a learning rate of  $10^{-5}$ , a weight decay parameter of 0, and default values for the other hyperparameters. Training was conducted over 50 epochs, with EEG responses for test image conditions synthesized using the model weights from the epoch that yielded the lowest validation loss. For each participant, the models generated EEG signals with a shape of 17 EEG channels  $\times$  250 EEG time points as the output corresponding to the input images. All experiments were conducted on a single NVIDIA 4090 GPU. For additional training details and validation procedures, see Appendix A.3.

**Target Features of EEG** We designed different target EEG features for semantic feature and spectral signature case. In the retrieval task based on semantic representation, the system randomly selects target images from the test set of THINGS-EEG2, with an index greater than 12 in each class. These selected images are excluded from the retrieval space of  $200 \times 12 = 2400$  images. In the generation task based on spectral features, in order to ensure that the regulation is meaningful, we calculated the EEG feature similarity matrix corresponding to the prediction of the  $200 \times 1$  image from the test set, and took the top-3 images with the lowest similarity in each class after row averaging as the target for testing. We use the pre-trained encoding model (AlexNet, CORnet-S) and pre-trained EEG encoders (ATM-S (Li et al., 2024), PSD) to process the target images and extract their corresponding EEG features.

## 4.2 REGULATION OF BRAIN SEMANTIC REPRESENTATION

To evaluate the effectiveness of our framework in achieving the target neural activity representation, we conducted a retrieval task in the image space. We treated the encoding model  $g$  as a black-box model, ensuring that gradients were not used to update its parameters. This approach allowed us

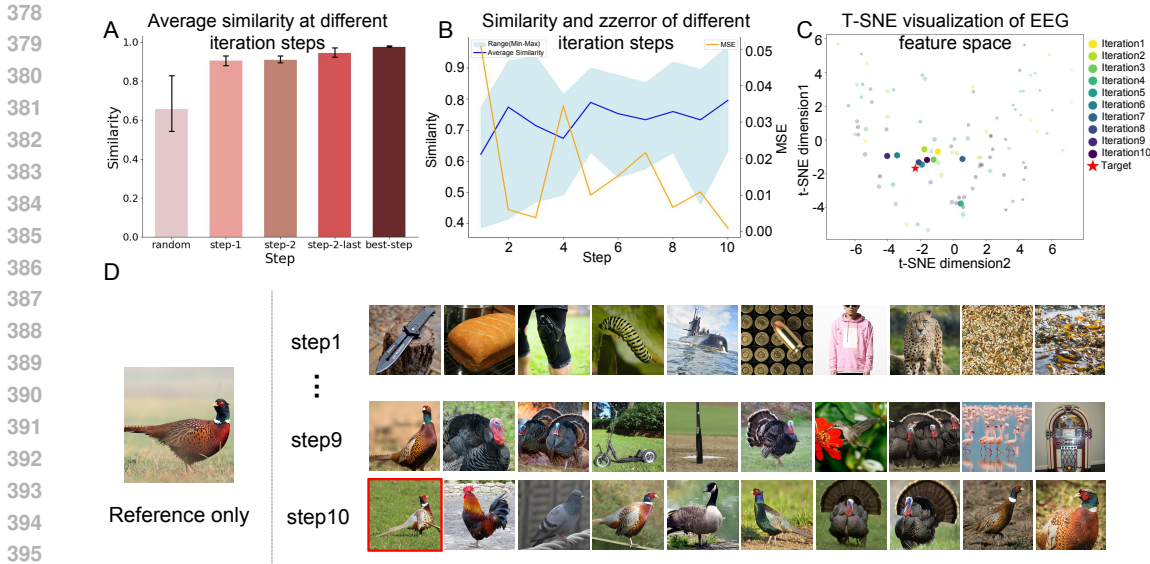


Figure 3: **Results of our framework in the retrieval task.** (A) Similarity between the neural representation obtained by our framework at different iteration steps (i.e., step-1, step-2, step-2-last, best-step) and the target neural representation compared to random stimulus (i.e., random). (B) The evolution of EEG representation similarity (blue) and loss curves (yellow) on Subject 8 at different iteration steps. (C) The t-SNE visualization of Subject 8’s latent trajectories within the feature space across all iterations. (D) The images retrieved by our framework at different iteration steps. Only the neural activity representation evoked by the reference image is known during the iteration process. See Appendix A.4 for more results.

to focus on the closed-loop regulation framework itself. The retrieval task was performed on the test set of the THINGS-EEG2 dataset, which consists of 2400 images. We used the EEG encoder ATM-S to obtain EEG semantic representations aligned with  $1 \times 1024$  CLIP image features. Before initiating the retrieval, random initialization was used to scatter 10 initial points as widely as possible in the image feature space. During the search process, each initial image sample calculates its cosine similarity with the global image features, and cumulative probability is applied to increase the likelihood of selecting new images that bring the EEG representation closer to the target. In the image feature space, the initial sample points expand iteratively, forming a small region, and gradually converge toward the theoretically optimal stimulus image. The termination condition for iterations is the similarity  $s(u_+, u_i) > threshold_{primary}$ .

In Figure 3, we report our retrieval results based on EEG semantic representation. In Figure 3A, we show the similarity scores of stimuli compared to random stimuli at different time steps during the iteration process. Figure 3B displays the average similarity and mean squared error between the predicted and expected EEG features at various iteration time points for subject 8. Figure 3C illustrates the convergence patterns from initial to final positions for selected iterations (e.g., iterations 1 and 10) across multiple cycles. In each iteration, ten images are presented, with points representing the closest match to the target stimulus at each step. Notably, these points gradually move toward the target stimulus, marked by a red pentagram, across successive iterations. For a given target neural activity representation, our framework iteratively predicts intermediate EEG results and retrieves stimulus images at each iteration. Importantly, only the neural activity representation evoked by the reference image is known throughout this process. Through successive iterations in Figure 3D, the framework refines its selection and ultimately retrieves an image (outlined in red) that closely matches the semantic representation of the reference image.

### 4.3 REGULATION OF INTENSITY OF NEURAL ACTIVITY

We implemented a closed-loop stimulus image generation framework using the  $200 \times 1 = 200$  image space of THINGS-EEG2 as initialization. We set the crossover rate  $\alpha$  to 0.6, the mutation rate  $\beta$  to



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

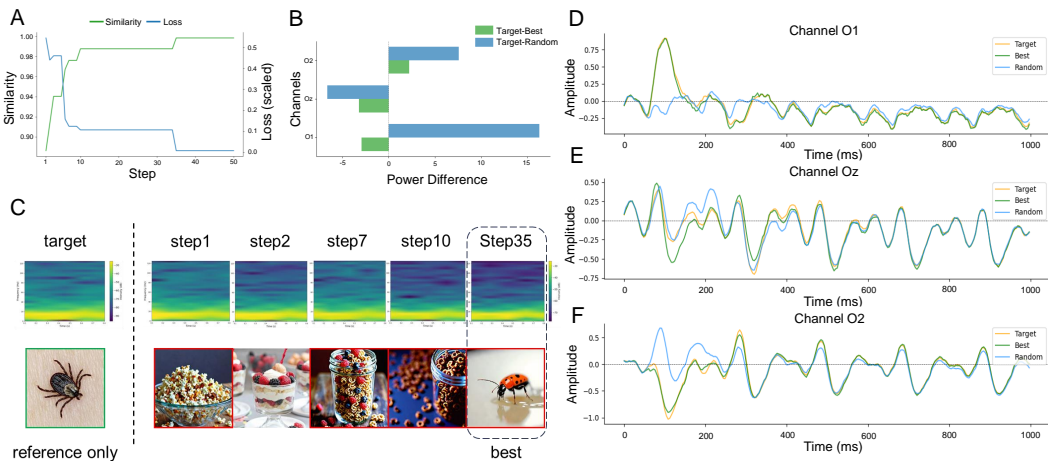


Figure 4: **Results of our framework in the generation task.** (A) Similarity and loss curves of EEG neural representations for Subject 8. (B) The difference of PSD between the neural activity representations evoked by the final step of generated and random stimulus, with the target neural representations used as the relative baseline. (C) For a given target EEG semantic representation, our framework iteratively predicts synthetic data, extract feature and synthesizes images at each iteration. The image enclosed by a red border represents the image synthesized by the generator, while the unbordered image is a sample selected from the original dataset. See Appendix A.5 for more generated examples. (D) EEG timing diagram generated by our stimulus images for  $O_1$  channel. (E) EEG timing diagram generated by our stimulus images for  $O_z$  channel. (F) EEG timing diagram generated by our stimulus images for  $O_2$  channel.

0.2, and randomly select 10 images from 200 images during initialization. We used StableDiffusion XL-turbo (Rombach et al., 2022) integrated by IP-Adapter (Ye et al., 2023) to generate new samples each time based on the new stimulus images obtained after crossover and mutation, and randomly selected 2 samples from the image feature space, calculated the similarity of EEG activity representation, and selected the next step of stimulation according to the roulette method of cumulative probability.

The results of our stimulus generation experiments are shown in Figure 4. Figure 4(A) shows the similarity and mean square error between the EEG features generated by the step stimulation image at different iterations and the target EEG features. In addition, we calculated the explained variance of different channels and selected the three channels  $O_1$ ,  $O_z$ , and  $O_2$  with the largest variance for regulation. Figure 4(B) shows the comparison of the PSD of the EEG predicted by the random and step-best samples relative to the target EEG representation. Figure 4(DEF) plots the synthetic EEG of three different channels obtained by step-best, random and target stimulation images respectively. All three channels show that the EEG corresponding to step-best, random and target images is quite different before 100 data points (corresponding to 0.4s). After 0.4s, due to the limitations of the encoding model itself, the synthetic EEG of the target image is not much different from the synthetic EEG of the optimal stimulation and the synthetic EEG of the random image. This corresponds to Fig.4 in (Gifford et al., 2022). Using the tick image as an example, Figure 4(C) shows the image and its corresponding time-frequency features, as well as the generated image and corresponding features at each iteration.

#### 4.4 REGULATION OF INDIVIDUAL VARIABILITY

Table 1 summarizes the results in the retrieval setting (corresponding to the representation score, SS) and the generation model setting (corresponding to the intensity score, IS), highlighting the results of our framework in achieving the optimal number of iterations in a given search space. The data show that for different target EEG features, our method has a good improvement in feature similarity across different subjects. For instance, the similarity score (SS) of the semantic feature of

Subject 7 is improved from 0.874 in step-1 to 0.974, with an improvement of 10.04%. Similarly, the feature similarity score (IS) of the channel intensity of Subject 8 is improved from 0.913 in step-1 to 0.990, accompanied by a 7.744% improvement. Even on the subjects with poor performance, our framework achieves a positive performance, which shows that our framework has a generalized improvement effect across different subjects, highlighting its potential in practical applications. See Appendix A.2 for more detailed quantitative results.

Table 1: **Performance (EEG semantic representation and intensity) of brain responses.** We provide two metrics: EEG semantic representation score (i.e., SS) and EEG response intensity score (i.e., IS) to measure the difference between the neural activity generated by the optimal stimulation image we obtained and the target EEG neural activity.

Subject	Step-1		Step-Best		Improvement	
	SS	IS	SS	IS	$\Delta$ SS (%)	$\Delta$ IS (%)
1	0.871	0.989	0.967	0.997	9.593	0.801
7	0.874	0.960	0.974	0.995	10.040	3.444
8	0.904	0.913	0.976	0.990	7.162	7.744
10	0.915	0.986	0.961	0.998	4.587	1.163

## 5 DISCUSSION AND CONCLUSION

In this study, we developed a flexible closed-loop visual stimulation framework for controlling EEG signatures. To the best of our knowledge, this is the first work to successfully employ closed-loop generation of natural images to modulate brain activity.

**Technical Impact:** Our framework demonstrated the potential of flexibly controlling EEG signals through visual stimulation. We employed a closed-loop iterative strategy, where new random stimuli are sampled each time a new round of stimulus images is generated. The gradient of the EEG objective is passed to the diffusion model in a proxy manner, eliminating the need for training or updating the weights of the generative model. This approach demonstrates that our framework is an efficient and optimal closed-loop stimulus generation method, capable of achieving the desired neural modulation without requiring any model parameter updates. It opens new avenues for applications in brain-computer interfaces, neuro-feedback systems, and therapeutic interventions for neurological disorders that require precise regulation of brain activity (Jang et al., 2021; Alamia et al., 2023).

**Neuroscience Insights:** Our study provides valuable insights into the neural mechanisms underlying visual perception and stimulus processing. First, we demonstrated the successful modulation of activity in specific electrode channels, indicating that neural activity in targeted brain regions can be fine-tuned through controlled visual stimulation. Second, we showcased our framework’s ability to guide the brain in generating specific neural representations, which is crucial for understanding how different brain regions process visual information and respond to external stimuli. Furthermore, our framework establishes a causal link between visual stimuli and neural responses. By connecting specific EEG patterns to visual representations, our work deepens the understanding of how neural signatures correlate with perceptual experiences.

**Interesting Phenomena and Future Directions:** Our findings demonstrate that different stimulus images in our framework can produce similar or identical EEG features, confirming the existence of Metamers (Feather et al., 2023) and suggesting that Metamers are not necessarily unique. The presence of multiple Metamers highlights the ill-posed nature of generating visual stimuli conditioned on EEG features. Future research should focus on understanding the neural mechanisms that lead to the generation of similar EEG features from different stimuli. Another promising direction is the integration of more sophisticated models that account for inter-individual variability in neural responses, aiming to fine-tune the stimulus generation process for personalized neuromodulation and enhanced brain-computer interaction (Alamia et al., 2021). Further exploration could involve integrating this closed-loop framework with other brain imaging modalities, such as fMRI or MEG. Additionally, it is crucial to formulate control goals aimed at regulating specific EEG characteristics to modulate brain functions, such as a control objective on EEG features for emotion regulation.

## REFERENCES

- 540  
541  
542 Andrea Alamia, Milad Mozafari, Bhavin Choksi, and Rufin VanRullen. On the role of feedback in  
543 visual processing: a predictive coding perspective. *arXiv preprint arXiv:2106.04225*, 2021.
- 544  
545 Andrea Alamia, Milad Mozafari, Bhavin Choksi, and Rufin VanRullen. On the role of feedback in  
546 image recognition under noise and adversarial attacks: A predictive coding perspective. *Neural  
547 Networks*, 157:280–287, 2023.
- 548  
549 Luca Baroni, Mohammad Bashiri, Konstantin F Willeke, Ján Antolík, and Fabian H Sinz. Learning  
550 invariance manifolds of visual sensory neurons. In *NeurIPS Workshop on Symmetry and Geometry  
551 in Neural Representations*, pp. 301–326. PMLR, 2023.
- 552  
553 Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep image  
554 synthesis. *Science*, 364(6439):eaav9436, 2019.
- 555  
556 Dave Epstein, Allan Jabri, Ben Poole, Alexei Efros, and Aleksander Holynski. Diffusion self-  
557 guidance for controllable image generation. *Advances in Neural Information Processing Systems*,  
558 36:16222–16239, 2023.
- 559  
560 Russell Epstein and Nancy Kanwisher. A cortical representation of the local visual environment.  
561 *Nature*, 392(6676):598–601, 1998.
- 562  
563 Jenelle Feather, Guillaume Leclerc, Aleksander Madry, and Josh H McDermott. Model metamers  
564 reveal divergent invariances between biological and artificial neural networks. *Nature Neuro-  
565 science*, 26(11):2017–2034, 2023.
- 566  
567 Marin Ferecatu and Donald Geman. Interactive search for image categories by mental matching. In  
568 *2007 IEEE 11th International Conference on Computer Vision*, pp. 1–8. IEEE, 2007.
- 569  
570 Alessandro T Gifford, Kshitij Dwivedi, Gemma Roig, and Radoslaw M Cichy. A large and rich eeg  
571 dataset for modeling human visual object recognition. *NeuroImage*, 264:119754, 2022.
- 572  
573 Tijl Grootswagers, Ivy Zhou, Amanda K Robinson, Martin N Hebart, and Thomas A Carlson. Hu-  
574 man eeg recordings for 1,854 concepts presented in rapid serial visual presentation streams. *Sci-  
575 entific Data*, 9(1):3, 2022.
- 576  
577 Zijin Gu, Keith Jamison, Mert R Sabuncu, and Amy Kuceyeski. Modulating human brain responses  
578 via optimal natural image selection and synthetic image generation. *ArXiv*, 2023.
- 579  
580 Matthias Guggenmos, Philipp Sterzer, and Radoslaw Martin Cichy. Multivariate pattern analysis for  
581 meg: A comparison of dissimilarity measures. *Neuroimage*, 173:434–447, 2018.
- 582  
583 Martin N Hebart, Adam H Dickter, Alexis Kidder, Wan Y Kwok, Anna Corriveau, Caitlin Van Wick-  
584 lin, and Chris I Baker. Things: A database of 1,854 object concepts and more than 26,000 natu-  
585 ralist object images. *PloS one*, 14(10):e0223792, 2019.
- 586  
587 Hannah F Iaccarino, Annabelle C Singer, Anthony J Martorell, Andrii Rudenko, Fan Gao, Tyler Z  
588 Gillingham, Hansruedi Mathys, Jinsoo Seo, Oleg Kritskiy, Fatema Abdurrob, et al. Gamma  
589 frequency entrainment attenuates amyloid load and modifies microglia. *Nature*, 540(7632):230–  
590 235, 2016.
- 591  
592 Hojin Jang, Devin McCormack, and Frank Tong. Noise-trained deep neural networks effectively  
593 predict human vision and its neural responses to challenging images. *PLoS biology*, 19(12):  
e3001418, 2021.
- 594  
595 Alex Krizhevsky. One weird trick for parallelizing convolutional neural networks. *arXiv preprint  
596 arXiv:1404.5997*, 2014.
- 597  
598 Jonas Kubilius, Martin Schrimpf, Kohitij Kar, Rishi Rajalingham, Ha Hong, Najib Majaj, Elias Issa,  
599 Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, et al. Brain-like object recognition with  
600 high-performing shallow recurrent anns. *Advances in neural information processing systems*, 32,  
601 2019.

- 594 Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip Torr. Controllable text-to-image genera-  
595 tion. *Advances in neural information processing systems*, 32, 2019.  
596
- 597 Dongyang Li, Chen Wei, Shiyong Li, Jiachen Zou, and Quanying Liu. Visual decoding and recon-  
598 struction via eeg embeddings with guided diffusion. *arXiv preprint arXiv:2403.07721*, 2024.
- 599 Andrew Luo, Maggie Henderson, Leila Wehbe, and Michael Tarr. Brain diffusion for visual explo-  
600 ration: Cortical discovery using large scale generative models. *Advances in Neural Information*  
601 *Processing Systems*, 36, 2024a.  
602
- 603 Junwen Luo, Chengyong Jiang, Qingyuan Chen, Dongqi Han, Yansen Wang, Biao Yan, Dongsheng  
604 Li, and Jiayi Zhang. The vep booster: A closed-loop ai system for visual eeg biomarker auto-  
605 generation. *arXiv preprint arXiv:2407.15167*, 2024b.
- 606 Anthony J Martorell, Abigail L Paulson, Ho-Jun Suk, Fatema Abdurrob, Gabrielle T Drummond,  
607 Webster Guan, Jennie Z Young, David Nam-Woo Kim, Oleg Kritskiy, Scarlett J Barker, et al.  
608 Multi-sensory gamma stimulation ameliorates alzheimer’s-associated pathology and improves  
609 cognition. *Cell*, 177(2):256–271, 2019.
- 610 Tomohiro Mizumoto, Harumi Ikei, Kosuke Hagiwara, Toshio Matsubara, Fumihiro Higuchi,  
611 Masaaki Kobayashi, Takahiro Yamashina, Jun Sasaki, Norihiro Yamada, Naoko Higuchi, et al.  
612 Mood and physiological effects of visual stimulation with images of the natural environment in  
613 individuals with depressive and anxiety disorders. *Journal of Affective Disorders*, 356:257–266,  
614 2024.  
615
- 616 Pawel Pierzchlewicz, Konstantin Willeke, Arne Nix, Pavithra Elumalai, Kelli Restivo, Tori Shinn,  
617 Cate Nealley, Gabrielle Rodriguez, Saumil Patel, Katrin Franke, et al. Energy guided diffusion  
618 for generating neurally exciting images. *Advances in Neural Information Processing Systems*, 36,  
619 2024.
- 620 Carlos R Ponce, Will Xiao, Peter F Schade, Till S Hartmann, Gabriel Kreiman, and Margaret S  
621 Livingstone. Evolving images for visual neurons using a deep generative network reveals coding  
622 principles and neuronal preferences. *Cell*, 177(4):999–1009, 2019.
- 623 Yongrong Qiu, David A Klindt, Klaudia P Szatko, Dominic Gonschorek, Larissa Hoefling, Timm  
624 Schubert, Laura Busse, Matthias Bethge, and Thomas Euler. Efficient coding of natural scenes  
625 improves neural system identification. *PLoS computational biology*, 19(4):e1011037, 2023.  
626
- 627 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
628 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
629 models from natural language supervision. In *International conference on machine learning*, pp.  
630 8748–8763. PMLR, 2021.
- 631 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
632 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
633 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.  
634
- 635 Tiberiu Tesileanu, Eugenio Piasini, and Vijay Balasubramanian. Efficient processing of natural  
636 scenes in visual cortex. *Frontiers in Cellular Neuroscience*, 16:1006703, 2022.
- 637 Edgar Y Walker, Fabian H Sinz, Erick Cobos, Taliah Muhammad, Emmanouil Froudarakis, Paul G  
638 Fahey, Alexander S Ecker, Jacob Reimer, Xaq Pitkow, and Andreas S Tolias. Inception loops  
639 discover what excites neurons most using deep predictive models. *Nature neuroscience*, 22(12):  
640 2060–2065, 2019.
- 641 Chen Wei, Jiachen Zou, Dietmar Heinke, and Quanying Liu. Cocog: Controllable visual stimuli  
642 generation based on human concept representations. *International Joint Conference on Artificial*  
643 *Intelligence*, 2024.  
644
- 645 Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt  
646 adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.  
647

A APPENDIX

A.1 MORE IMPLEMENTATION DETAILS

A.1.1 DATASETS

We conducted our experiments using the training set of the THINGS-EEG2 dataset (Gifford et al., 2022; Grootswagers et al., 2022), which consists of a large EEG corpus from 10 human subjects performing a visual task. The experiments used the Rapid Serial Visual Presentation (RSVP) paradigm for orthogonal target detection tasks to ensure participants’ attention to the visual stimuli. All 10 participants underwent 4 equivalent experiments, resulting in 10 datasets with 16,540 unique training image conditions, each repeated 4 times, and 200 unique testing image conditions, each repeated 80 times. In total, this yielded  $(16,540 \text{ training image conditions} \times 4 \text{ repetitions}) + (200 \text{ testing image conditions} \times 80 \text{ repetitions}) = 82,160$  image trials. The original data were recorded using a 64-channel EEG system with a 1000 Hz sampling rate. For preprocessing, the data were first down-sampled to 250 Hz and 17 channels were selected from the occipital and parietal regions, which are closely related to the visual system. The EEG data were then segmented into trials, spanning from 0 to 1000 ms post-stimulus onset, with baseline correction applied using the mean of the 200 ms pre-stimulus period. Multivariate noise normalization was applied to the training data (Guggenmos et al., 2018).

A.1.2 RETRIEVAL PIPELINE

We provide a more detailed description of algorithm 1. The algorithm begins by initializing equal selection probabilities for each image in the candidate set, denoted as  $p_0(u) = \frac{1}{N}$ , where  $N$  is the total number of images in the retrieval set. This initialization with equal probabilities reflects the absence of prior information, serving as an exploratory phase. In each iteration (representing a **state** in the MDP framework), a subset of images  $U_t = \{u_1, u_2, \dots, u_j\}$  is selected from the candidate images set  $U$  based on the current selection probabilities  $p_t(u)$ .

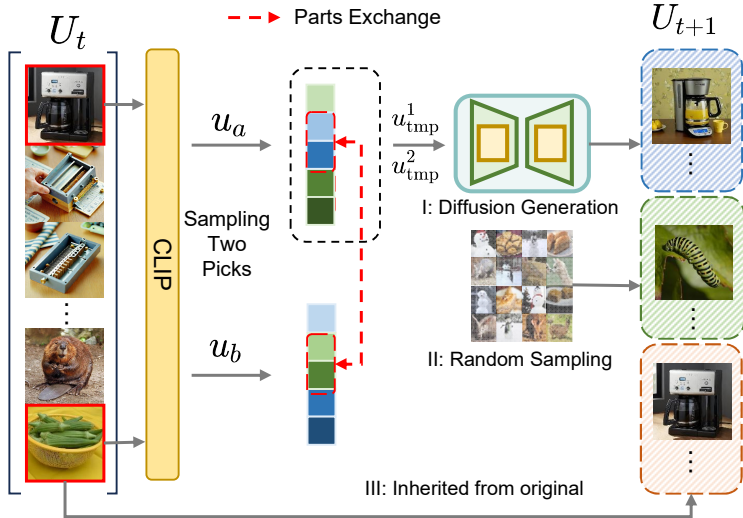


Figure A.1: Generating subsequent images based on the current round is achieved through crossover, variation, and a guided diffusion model. Both crossover and mutation operations preserve the relative ordering of CLIP features, thereby maintaining their semantic coherence.

For each image  $u_j$  in the subset  $U_t$  the algorithm computes a similarity score  $sim\langle u_j, u_{target} \rangle$  by comparing the image’s representation with the target. This similarity score acts as an immediate **reward** within the MDP framework. The maximum similarity score among the subset is identified as a measure of the effectiveness of the current action. If  $sim_{max}$  does not meet a predefined  $threshold_1$ , the reward is considered insufficient, and the algorithm returns to the image selection step, effectively trying a new action within the same state. If  $sim_{max}$  meets or exceeds the threshold, the algorithm proceeds to identify the two images  $u_{top1}$  and  $u_{top2}$  with the highest similarity

702 scores. These two images act as reference points for updating the probabilities of other images in  
703 the subsequent state.

704  
705 As for each image  $u_j$  in  $U$  that surpasses  $threshold_2$  with either  $u_{top1}$  or  $u_{top2}$ , its selection proba-  
706 bility  $P_{t+1}(u_j)$  is updated by multiplying with a constant factor, representing a policy improvement  
707 step that prioritizes images likely to yield higher rewards. After updating, a Softmax function is ap-  
708 plied to normalize the probabilities, focusing selection weight on images more similar to the target.  
709 This normalization step reflects the transition to a new state with an updated policy. The iteration  
710 continues, with the algorithm transitioning through states by selecting new subsets based on the re-  
711 fined probabilities, until  $sim_{max}$  reaches  $threshold_{primary}$ . At this point, the loop terminates, as  
712 the algorithm has successfully identified an optimal subset of images that maximizes the similarity  
713 reward to the target.

### 714 A.1.3 GENERATION PIPELINE

715 We provide a more detailed description of algorithm 2. As illustrated in Figure A.1, each image set  
716 consists of three parts:

- 717  
718 • Part one: This step focuses on the diffusion generation process. From the image set of last  
719 iteration, two images, denoted as  $u_a$  and  $u_b$ , are sampled using a roulette wheel selection  
720 method. A random crossover is then applied to part of their image embeddings, with the  
721 crossover starting at a different index each time. The newly combined image embedding is  
722 then used as input to the diffusion process. This increases the variability of the image set  
723 while preserving the high-quality components of the image embeddings.
- 724 • Part two: In this step, images are randomly sampled from the original image dataset, ex-  
725 cluding those that have already been selected in earlier iterations. This ensures that the new  
726 image set introduces novel elements while avoiding repetition.
- 727 • Part three: This part inherits the image  $u_a$  and  $u_b$ .

728  
729 By combining these three parts, we obtain a new image set for the next iteration.  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

A.2 ADDITIONAL QUANTITATIVE RESULTS

A.2.1 ITERATION IMPROVEMENT FROM DIFFERENT SUBJECTS

Based on the conclusions drawn from Figure A.4, we employ the pre-trained AlexNet end-to-end model as the EEG encoder and use ATM-S, which is based on S-S (both the training and testing signals are synthesized), to obtain semantic representations aligned with  $1 \times 1024$  CLIP image features. The experimental design involves randomly selecting 50 categories, resulting in a retrieval space of  $50 \times 12 = 600$  images. Specifically, we present the iterative performance improvements for three different targets randomly selected from the test set, with results reported for Subjects 1, 7, 8, and 10. As shown in Figure A.2, we calculate the EEG feature similarity of Subject 1, 7, 8, and 10 at random, step-1, and step-best in the iterative process respectively.

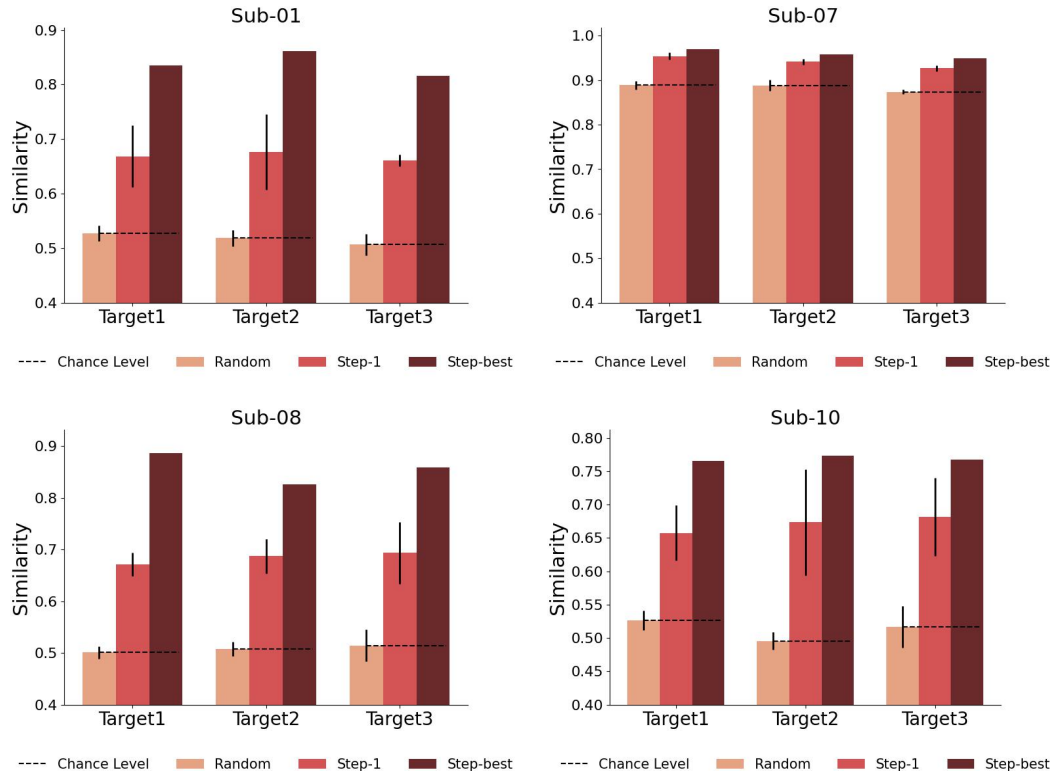


Figure A.2: **Comparison of improved performance by different targets.** We present the similarity scores of EEG features generated by random stimulation, open-loop stimulation (step 1), and step-best stimulation, in comparison to the target features. Each subject randomly selected 3 images from the retrieval space as target images.

A.2.2 PERFORMANCE OF DIFFERENT TARGET IMAGES ACROSS SUBJECTS

We report the results of iterative optimization using different targets in two different cases. The results for each subject are shown, along with the average percentage improvement across 5 random seeds. For the semantic feature case, unlike the setting in Table 1 of the main text, which uses real EEG for training and performs retrieval on synthetic EEG, we determined that training and testing with synthetic EEG yielded the highest accuracy based on the retrieval performance shown in Figure A.4. As a result, we retrained each subject and summarized the results in Table A.1. For the intensity feature case, we selected 3 images using the method described in Section 4 and supplemented the iterative improvement performance. We performed t-tests on EEG semantic and spectral features across all subjects to assess the efficacy of our proposed method. Additionally, we performed correlation analyses to investigate the relationships between semantic features and clip representation, as well as between PSD feature and clip representation, as shown in Figure A.3.

Table A.1: **Performance (EEG semantic representation and intensity) of brain responses.** We provide two metrics: EEG semantic representation score (i.e., SS) and EEG response intensity score (i.e., IS) to quantify the similarity of generated EEG and target EEG. The table below records the SS & IS values for each subject, showing the SS & IS value from the first round of stimulation, the SS & IS value achieved after multiple rounds of closed-loop control (the optimal result), and the improvement in control. All these results are calculated from pretrained AlexNet models.

Subject	Random		Step-1		Step-Best		Improvement	
	SS	IS	SS	IS	SS	IS	$\Delta$ SS (%)	$\Delta$ IS (%)
1	0.5174	0.9632	0.6686	0.9729	0.8375	0.9976	16.8859	2.4790
2	0.5197	0.9678	0.6675	0.9764	0.7372	0.9998	6.9701	2.3406
3	0.5113	0.9883	0.6597	0.9927	0.7871	0.9980	12.7402	0.5306
4	0.5065	0.9650	0.6498	0.9836	0.8299	0.9963	<b>18.0136</b>	1.2690
5	0.5315	0.9788	0.6937	0.9768	0.8418	0.9979	14.8151	2.1055
6	0.6747	0.9836	0.8099	0.9856	0.8826	0.9961	7.2634	1.0461
7	0.8838	0.8955	0.9410	0.9033	0.950	0.9742	1.8237	<b>7.0879</b>
8	0.5077	0.8344	0.6838	0.9435	0.8568	0.9925	17.3066	4.8947
9	0.8465	0.9602	0.9251	0.9751	0.9597	0.9997	3.4662	2.4597
10	0.5128	0.8172	0.6707	0.9705	0.7687	0.9934	9.8032	2.2849

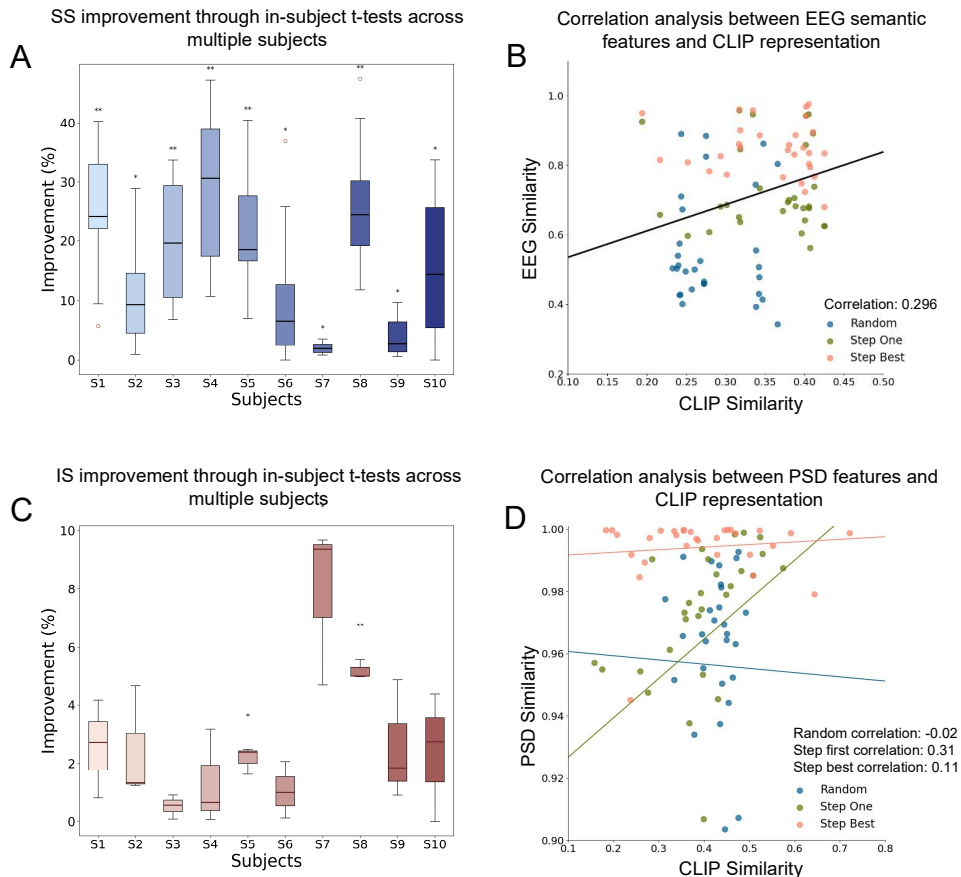


Figure A.3: **Improvement in similarity scores assessed via paired t-tests and correlation of similarity scores with targets across all subjects.** (A) Average EEG semantic representation scores (SS) for various target EEG semantic features. (B) The correlation of the similarity score with target between EEG semantic features across all subjects. (C) Average EEG response intensity scores (IS) for different target EEG PSD features. (D) The correlation of the similarity score with target between EEG PSD features across all subjects.



### A.3 VALIDITY VERIFICATION OF SYNTHETIC EEG

To evaluate the performance of our EEG encoding models, we compare the synthetic EEG signals generated by two deep neural networks (DNNs)—AlexNet and CORnet-S—with real EEG data. Here’s a step-by-step breakdown of how we processed and compared the data.

We selected 17 specific channels from the original 63-channel EEG dataset, focusing on those most relevant to visual processing. It ensured that we focused on neural regions most directly involved in responding to the visual stimuli. For each stimulus, we averaged the EEG signals across all trials, resulting in a representative dataset for each stimulus. This reduced the dimensionality of the data, making it easier to compare with synthetic data. We used a pretrained end-to-end encoding model to generate synthetic EEG signals based on the visual stimuli. The model captures the mapping between the visual input and the resulting EEG signals using deep neural networks. These synthetic signals represent the neural responses predicted by the model in response to the stimuli.

Table A.2: MSE Values for synthesized EEG

Subject	Pretrained		Random Init		Average
	AlexNet	CORnet-S	AlexNet	CORnet-S	
<b>Sub-01</b>	0.1095	0.1126	0.1161	0.0994	0.1094
<b>Sub-02</b>	0.0764	0.0788	0.0840	0.0994	0.0847
<b>Sub-03</b>	0.0787	0.0806	0.0816	0.0910	0.0830
<b>Sub-04</b>	0.0652	0.0664	0.0662	0.1011	0.0747
<b>Sub-05</b>	0.0493	0.0515	0.0704	0.0975	0.0672
<b>Sub-06</b>	0.0690	0.0719	0.0498	0.0966	0.0718
<b>Sub-07</b>	0.1267	0.1300	0.0914	0.1312	0.1198
<b>Sub-08</b>	0.0718	0.0727	0.1038	0.1165	0.0912
<b>Sub-09</b>	0.0529	0.0563	0.0781	0.0756	0.0657
<b>Sub-10</b>	0.1122	0.1151	0.0961	0.1149	0.1096
<b>Average</b>	0.0810	0.0832	0.0838	0.1023	0.0876

Table A.2 presents the mean squared error (MSE) between the synthetic EEG signals generated by AlexNet and CORnet-S, and the real EEG signals for 10 subjects. The MSE was computed for each individual test sample and then averaged across the entire test set. Lower MSE values indicate better alignment between the synthetic and real EEG signals.

From the comparison shown in the Figure A.4, the retrieval accuracy for S-S (both training and testing sets consist of generated signals) is significantly higher than other categories, including T-T (both training and testing sets consist of real signals), T-S (training set consists of real signals, testing set consists of generated signals), and S-T (training set consists of generated signals, testing set consists of real signals), under both AlexNet and CORnet-S models. This indicates:

**Advantages of generated signals** Supported by black-box ANN models (e.g., AlexNet and CORnet-S), generated signals perform significantly better in retrieval tasks compared to real signals. In particular, the highest retrieval accuracy for S-S demonstrates the consistency and model adaptability of generated signals in this retrieval task.

**Model adaptability:** Different ANN models (e.g., AlexNet and CORnet-S) show consistent superiority in the retrieval tasks for generated signals, indicating that generated signals are more easily captured and distinguished by black-box models.

In Figure A.5, we compute the variance across all samples and time points for each channel, providing a measure of the overall variability of the EEG signals in response to different visual stimuli and their temporal dynamics. This variance can help identify channels with the highest variability, which may be useful for selecting specific channels for further analysis or modulation.

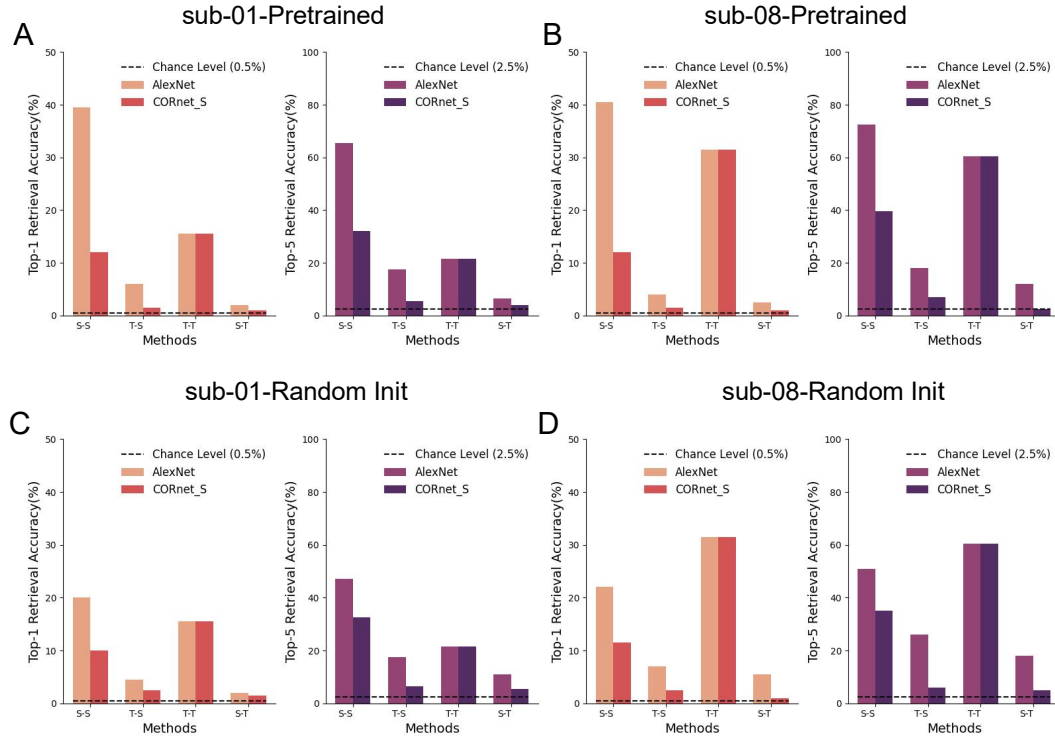


Figure A.4: Retrieval accuracy under different training and test datasets. Zero-shot retrieval performance of EEG data from different sources in Subject 1 and Subject 8 using ATM-S in different Settings. AlexNet and CORnet-S used in the first row were both pre-trained end-to-end models, and the second row was randomly initialized end-to-end.

In Figure A.6, we show the variance and standard deviation of the EEG signals computed across samples for each time point, and then averaged across channels. This analysis allows us to assess how signal variability evolves over time. By comparing the real EEG data with synthetic data generated by AlexNet and CORnet-S, we can evaluate how well each model captures the temporal variability present in the real EEG signals.

In Figure A.7, we compute the Pearson correlation coefficient between the averaged real EEG data and the synthetic data for each stimulus, measuring how well the synthetic data matches the real EEG on a per-sample basis. The histogram shows the distribution of correlation coefficients across all samples for both AlexNet and CORnet-S. A higher concentration of peaks near higher Pearson coefficients indicates better alignment between the synthetic data and the real EEG, reflecting superior model performance.

In Figure A.8, for each time point, we compute the Pearson correlation between the real EEG signal and the synthetic signals. This analysis enables us to visualize how well each model replicates the temporal structure of real neural responses to visual stimuli. Shaded regions in the plot represent the standard deviation across samples, showing the variability in model performance over time. The results provide a detailed view of how each model performs at different time points, highlighting which model more accurately captures the temporal dynamics of EEG signals.

From the above analysis, we observe that both AlexNet and CORnet-S perform well, showing comparable results in terms of MSE, spatial (channel-wise) variability, and temporal (time-resolved) variability. The Pearson correlation analysis further confirms that both models synthesize EEG signals that align well with real data, with subtle differences in performance between them. These findings highlight the robustness of our EEG encoding models, demonstrating their ability that not only mimic the structural features of real EEG data but also capture the realistic variability seen in neural responses to visual stimuli. This suggests that our models are effective in approximating the neural representations underlying visual processing.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

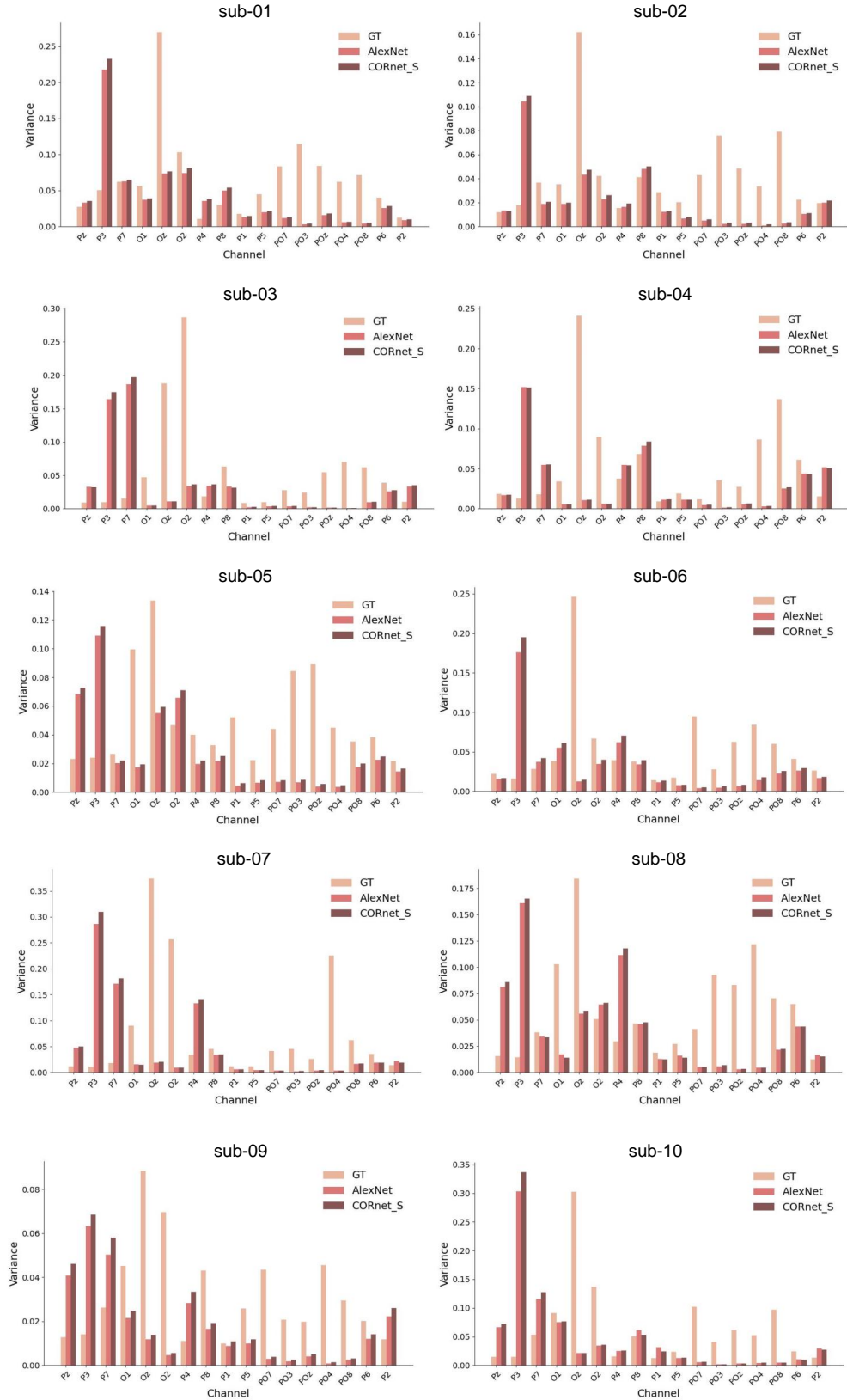


Figure A.5: Variance across different channels for different visual stimulus and temporal dynamics

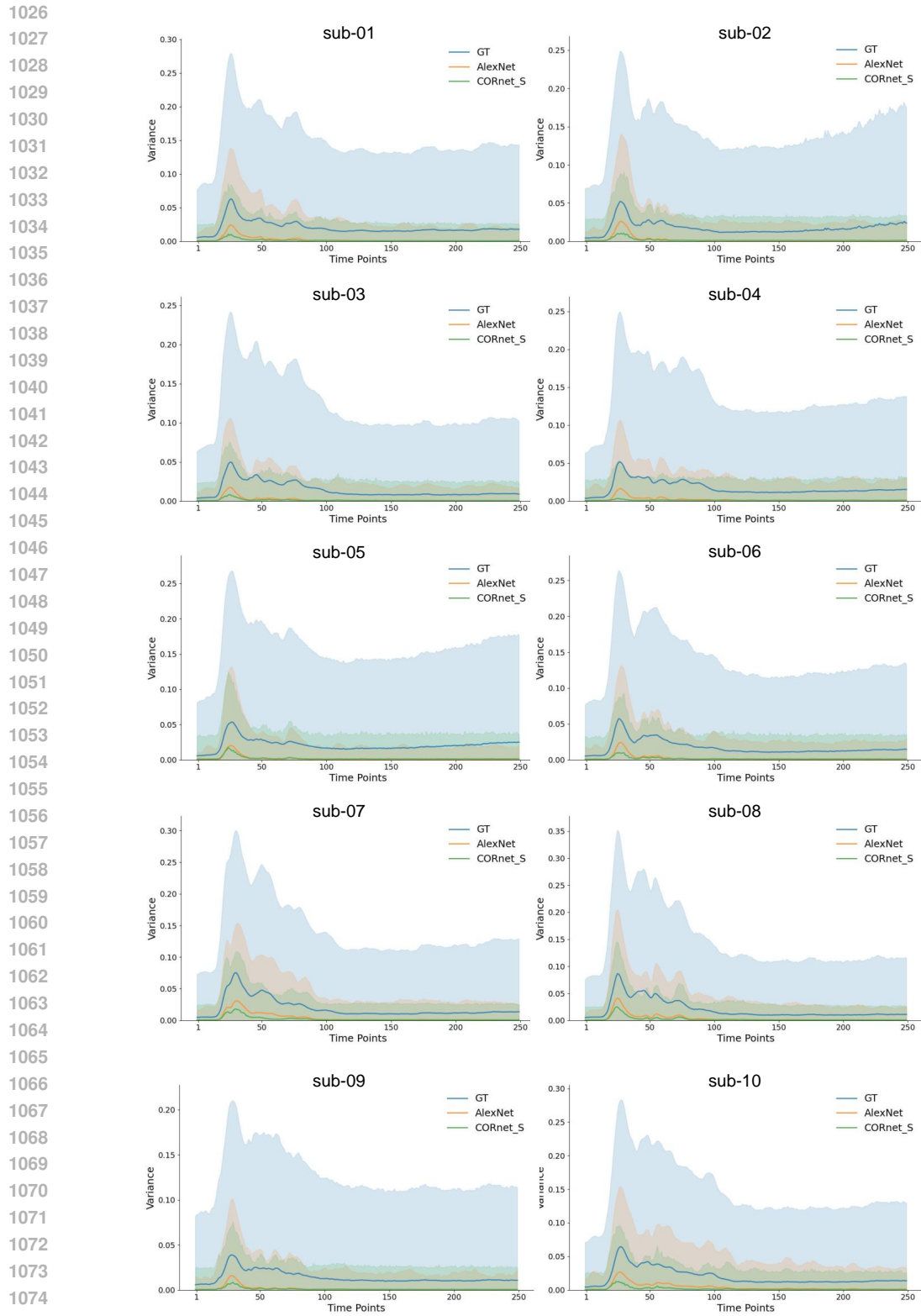


Figure A.6: Variance across different time points for different visual stimuli and channels.

1075

1076

1077

1078

1079

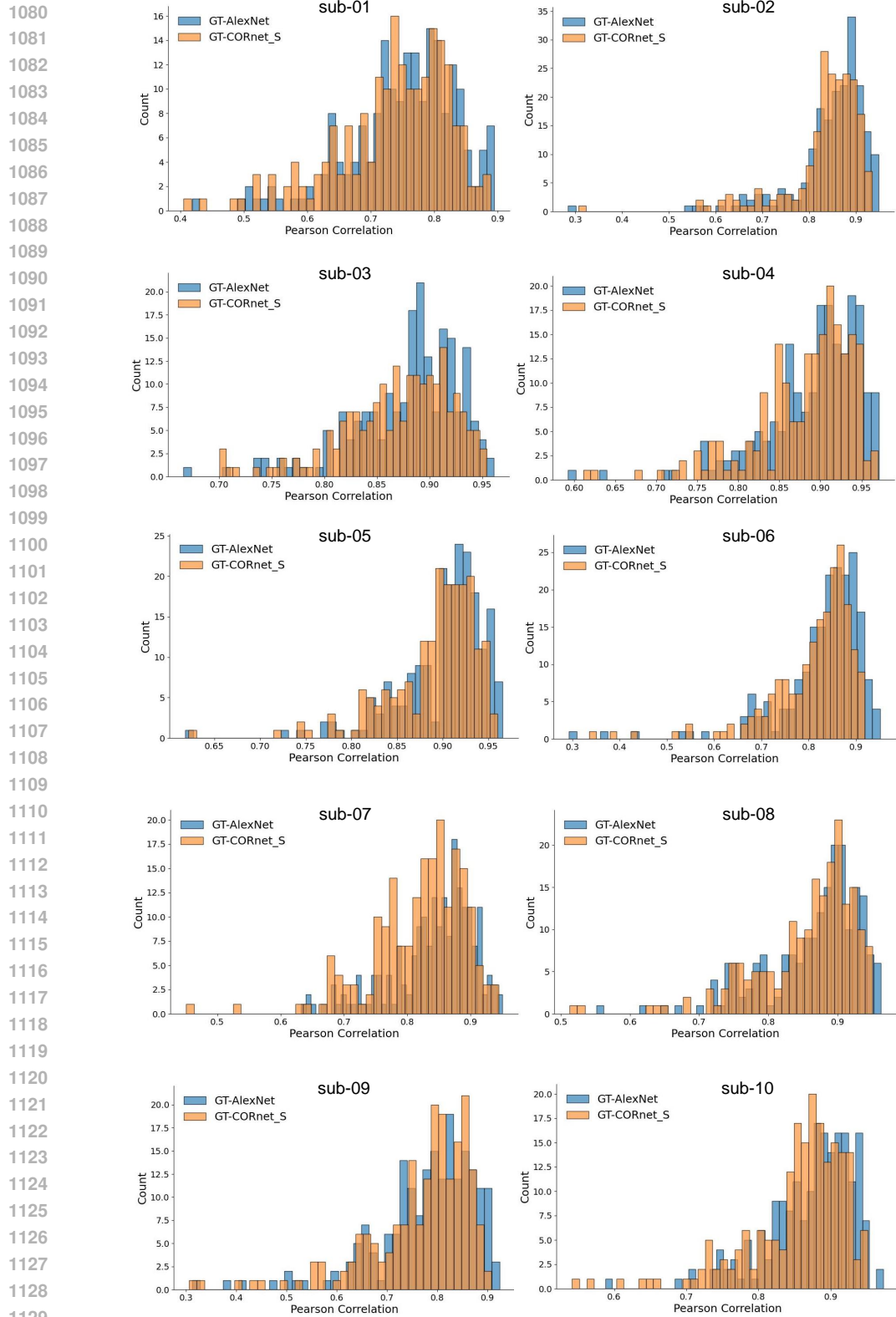


Figure A.7: Distribution of Pearson correlation coefficients across all sample pairs.

1132

1133

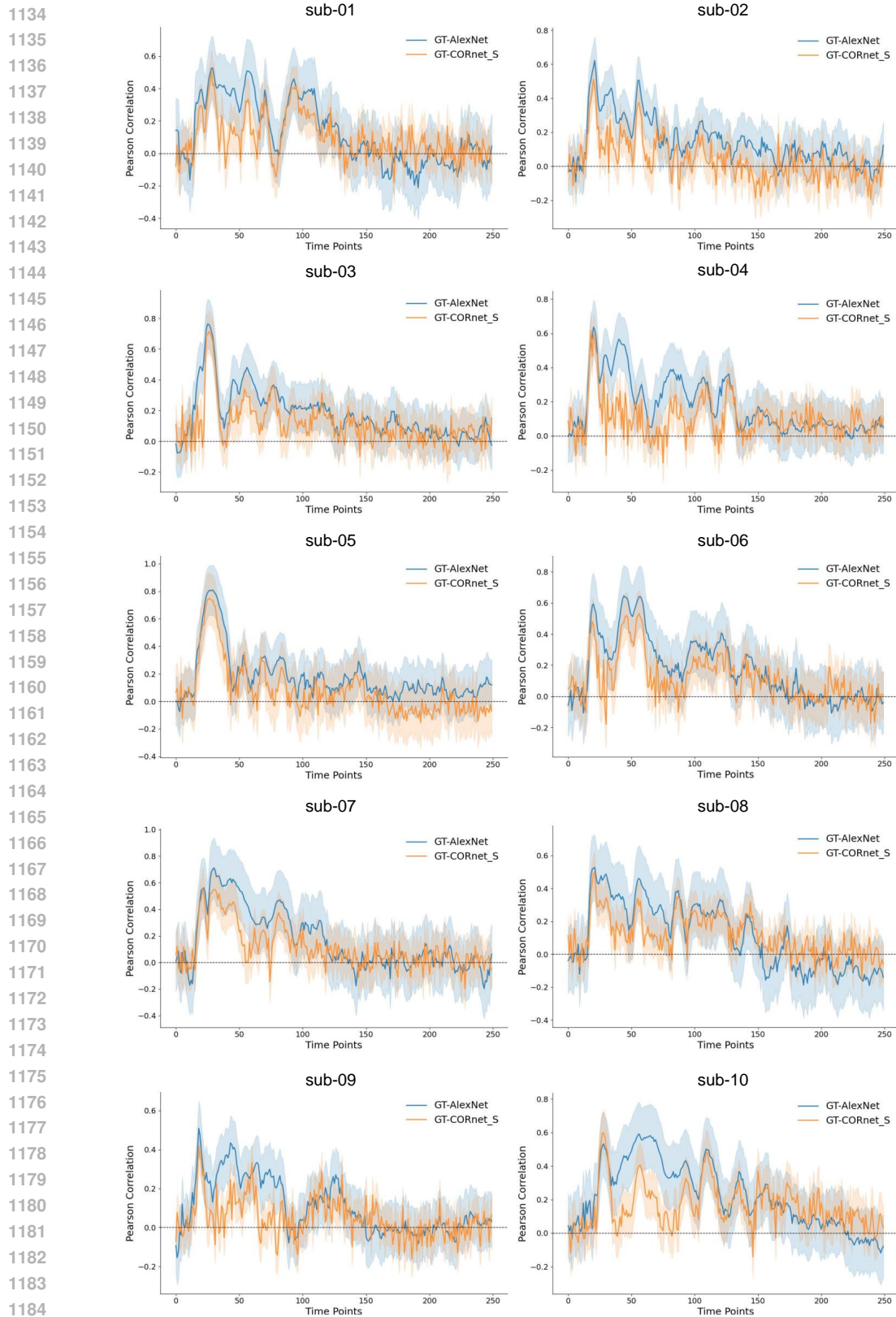


Figure A.8: Time-resolved Pearson correlation between ground truth EEG signals and synthetic EEG signals predicted by two neural network models (AlexNet and CORnet-S).

A.4 ADDITIONAL RETRIEVAL EXAMPLES OF SEMANTIC REPRESENTATION

A.4.1 MORE EXAMPLES OF RETRIEVAL

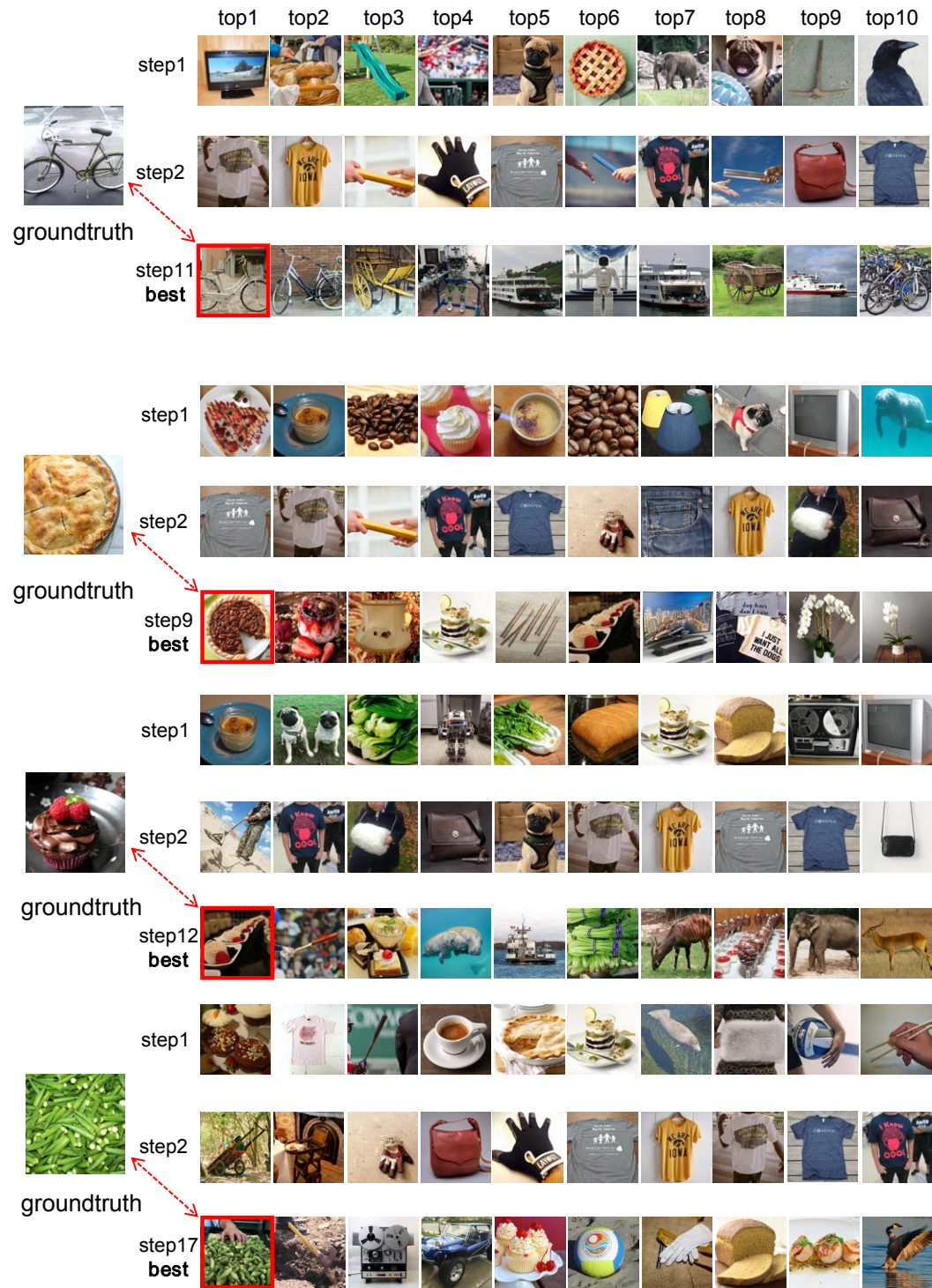


Figure A.9: Some retrieval examples of Subject 8, 4, 4, and 1. By setting different targets, we present examples where the stimulus retrieved at the end of the iterative optimization process increasingly approximates the true category.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

A.4.2 SOME FAILURE EXAMPLES OF RETRIEVAL



Figure A.10: **Some retrieval failure examples of Subject 8.** By setting different targets, we show examples where the stimulus retrieved at the end of the iteration is far from the true category. In these examples, the final retrieved stimulus exhibits varying degrees of similarity to the target image.



A.5 ADDITIONAL CONTROLLABLE GENERATION EXAMPLES OF PSD FEATURE

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

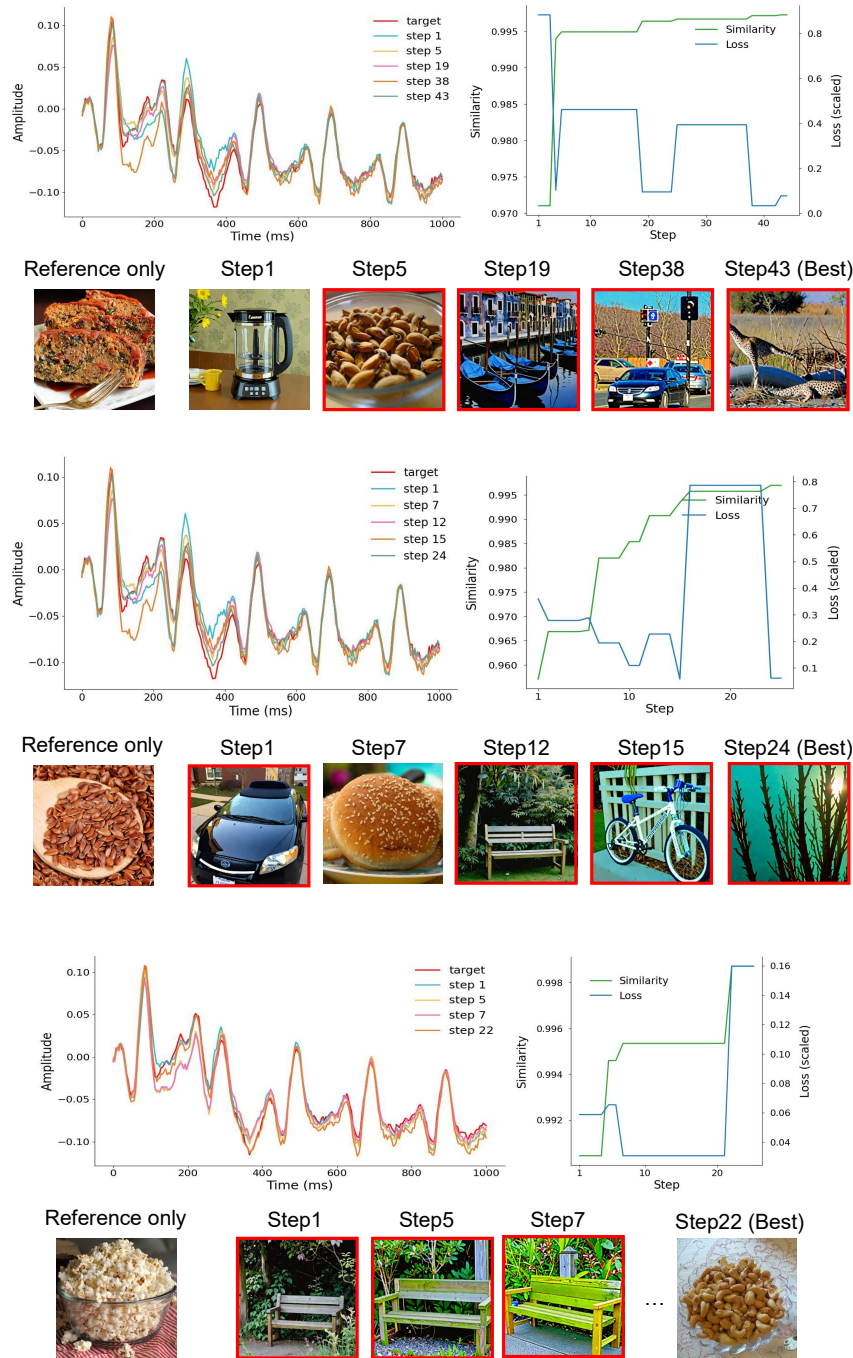


Figure A.11: **Illustration of the closed-loop iterative process for Subject 1.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

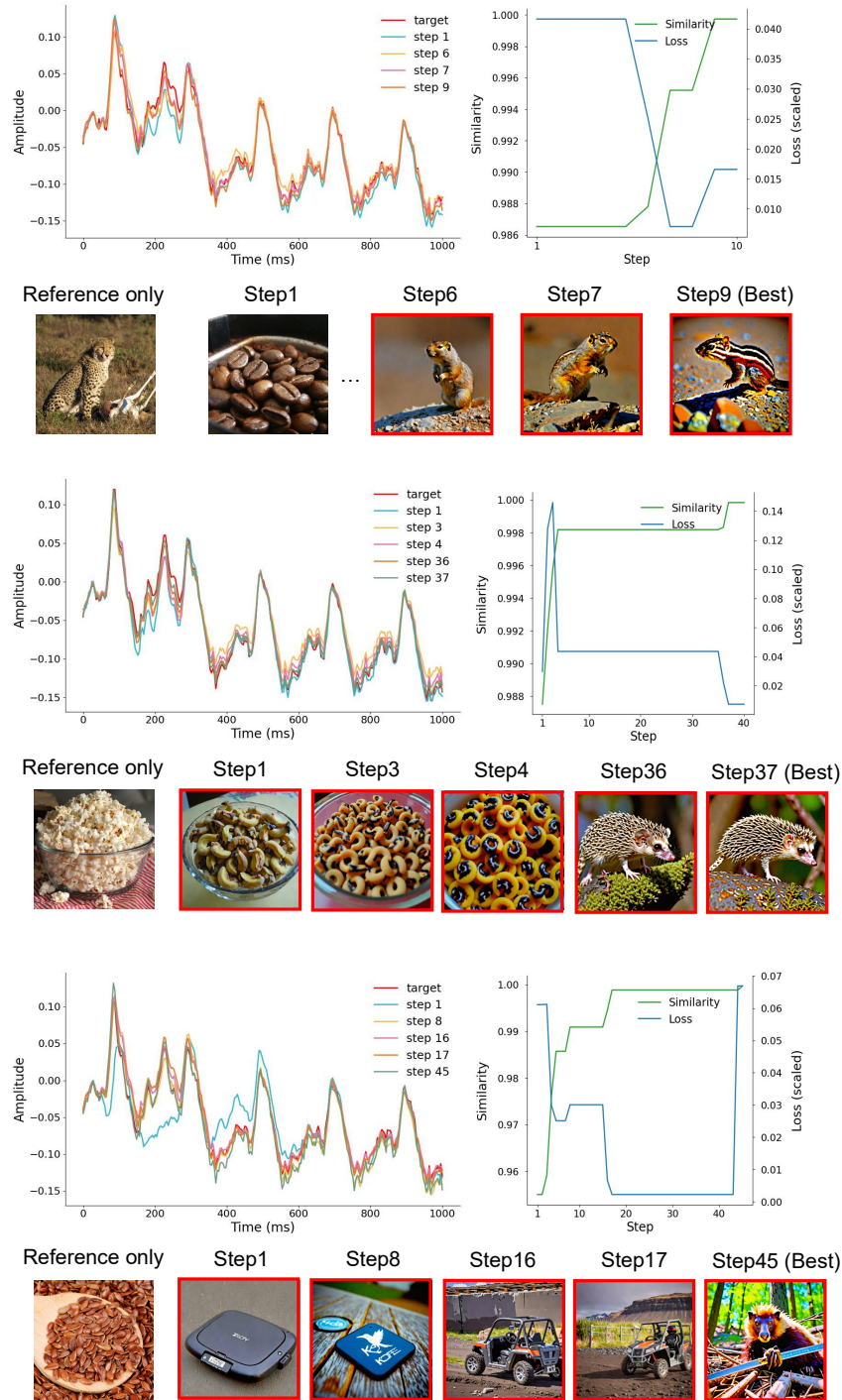


Figure A.12: **Illustration of the closed-loop iterative process for Subject 2.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

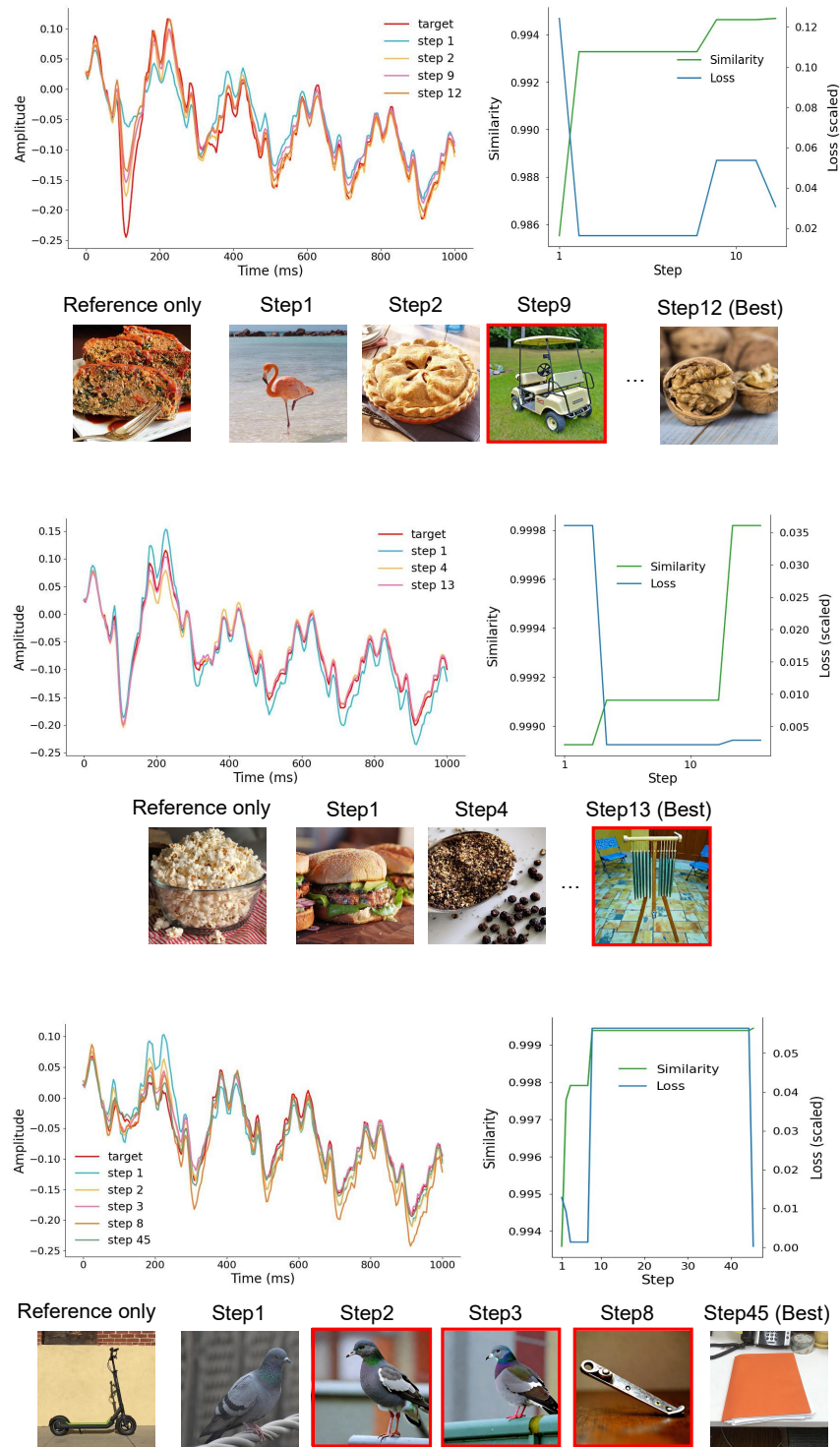


Figure A.13: **Illustration of the closed-loop iterative process for Subject 3.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

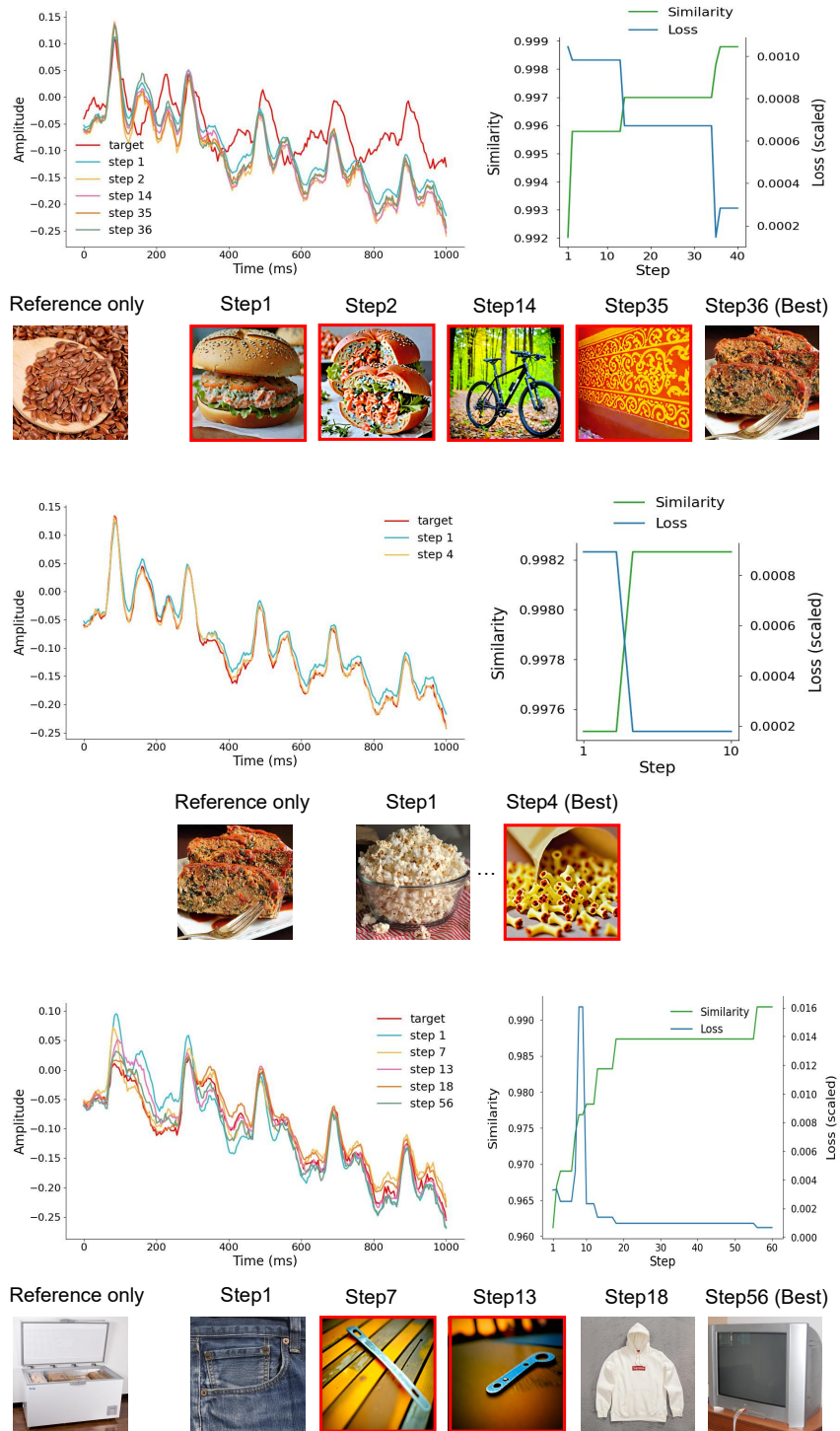


Figure A.14: **Illustration of the closed-loop iterative process for Subject 4.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

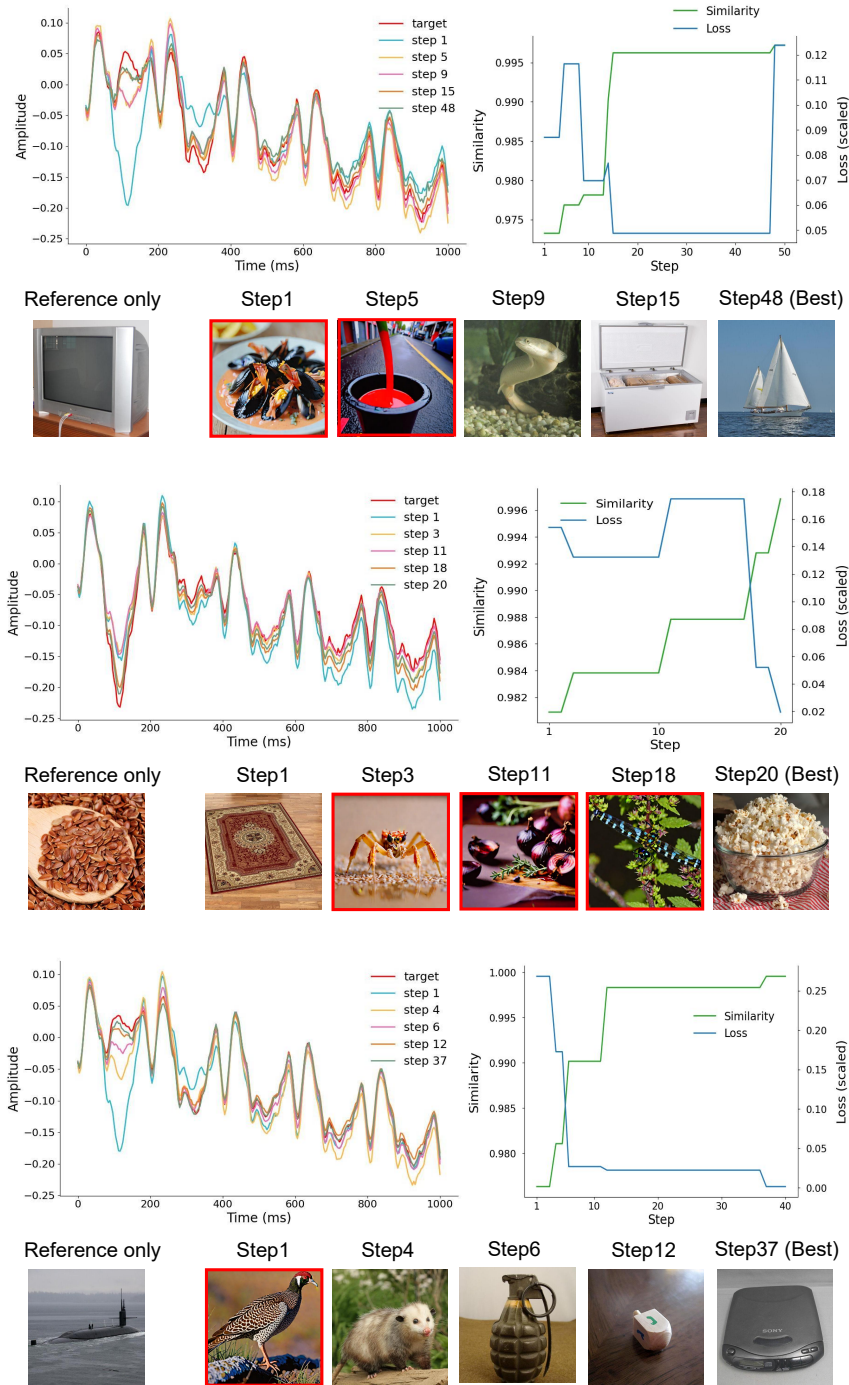


Figure A.15: **Illustration of the closed-loop iterative process for Subject 5.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

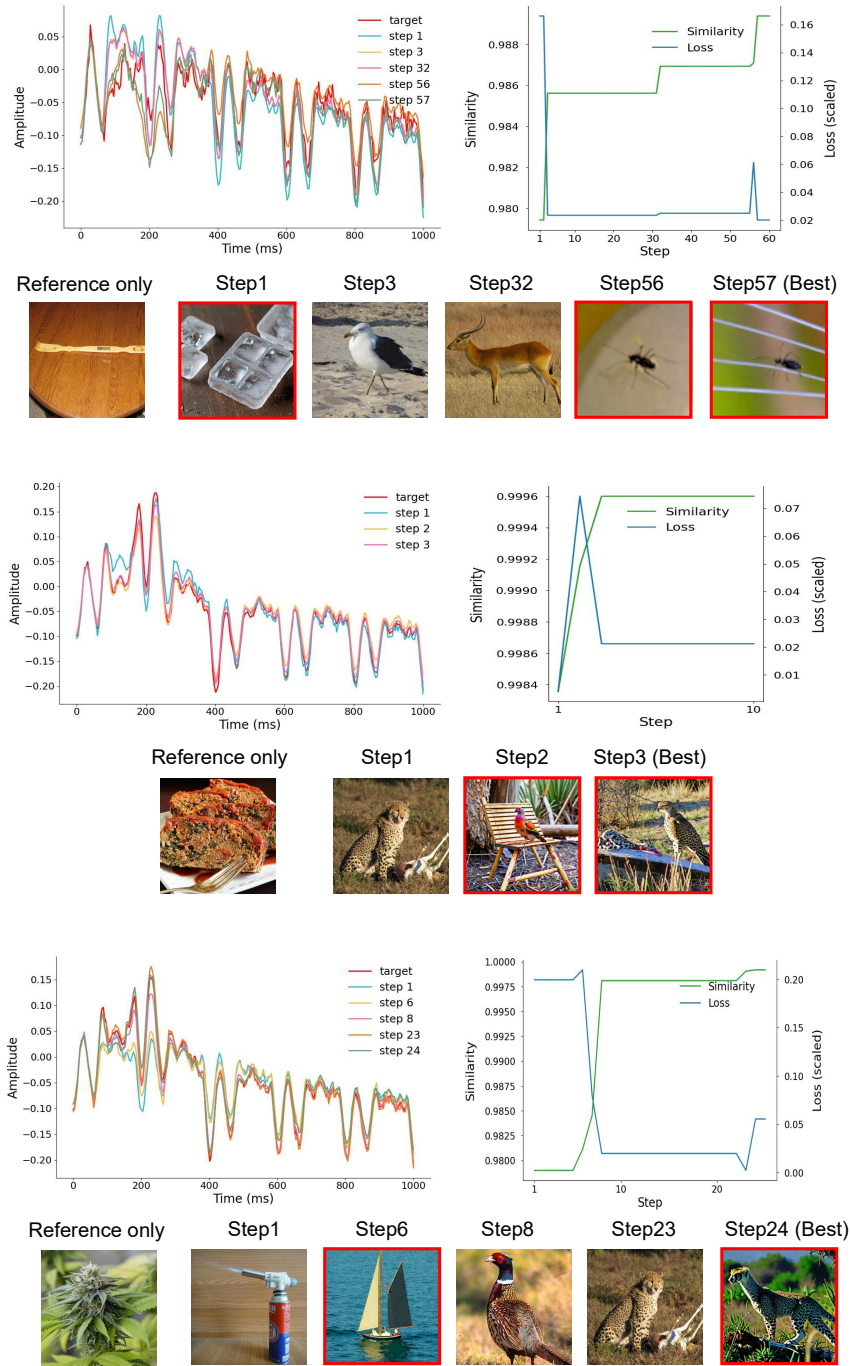


Figure A.16: **Illustration of the closed-loop iterative process for Subject 6.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

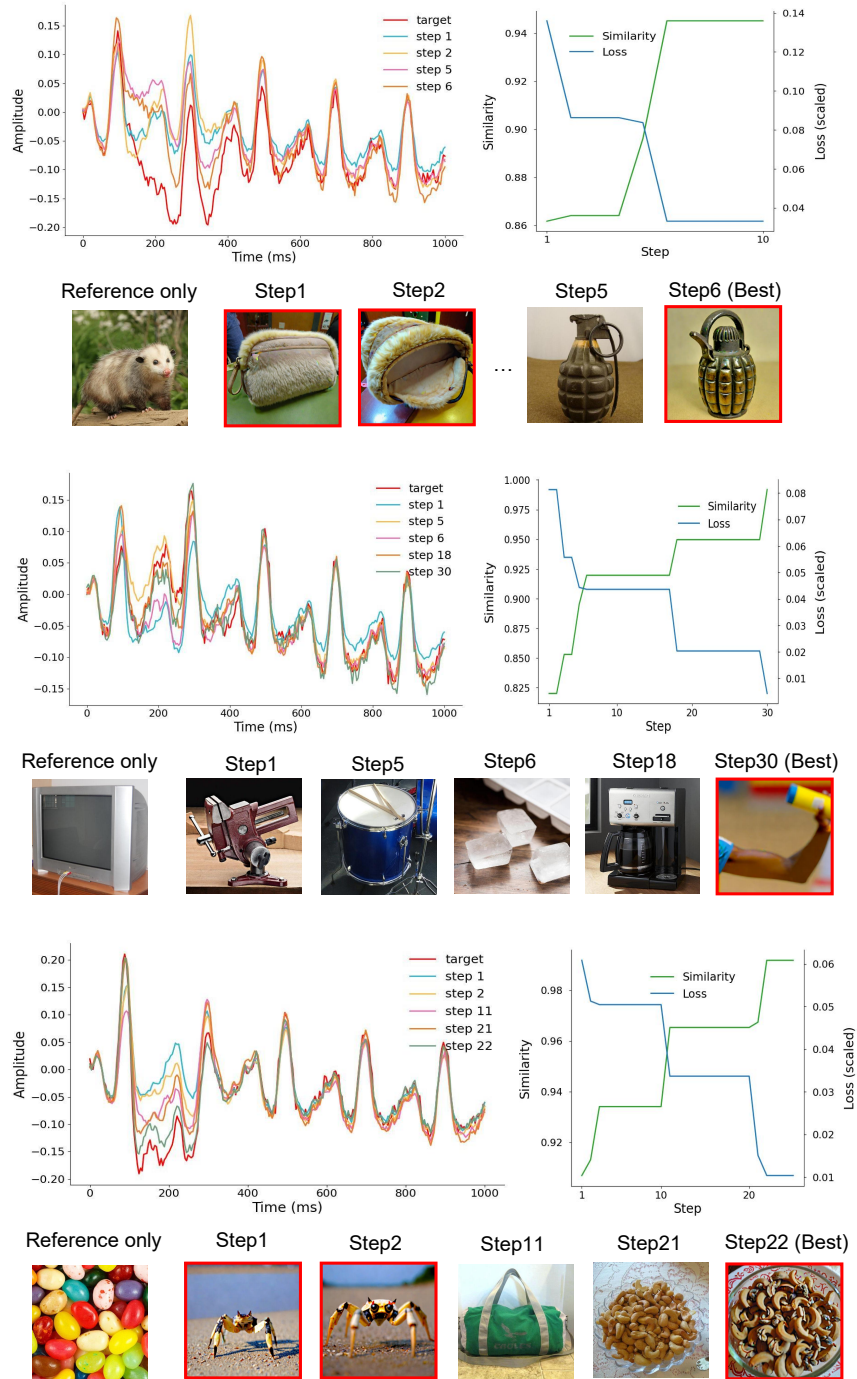


Figure A.17: **Illustration of the closed-loop iterative process for Subject 7.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

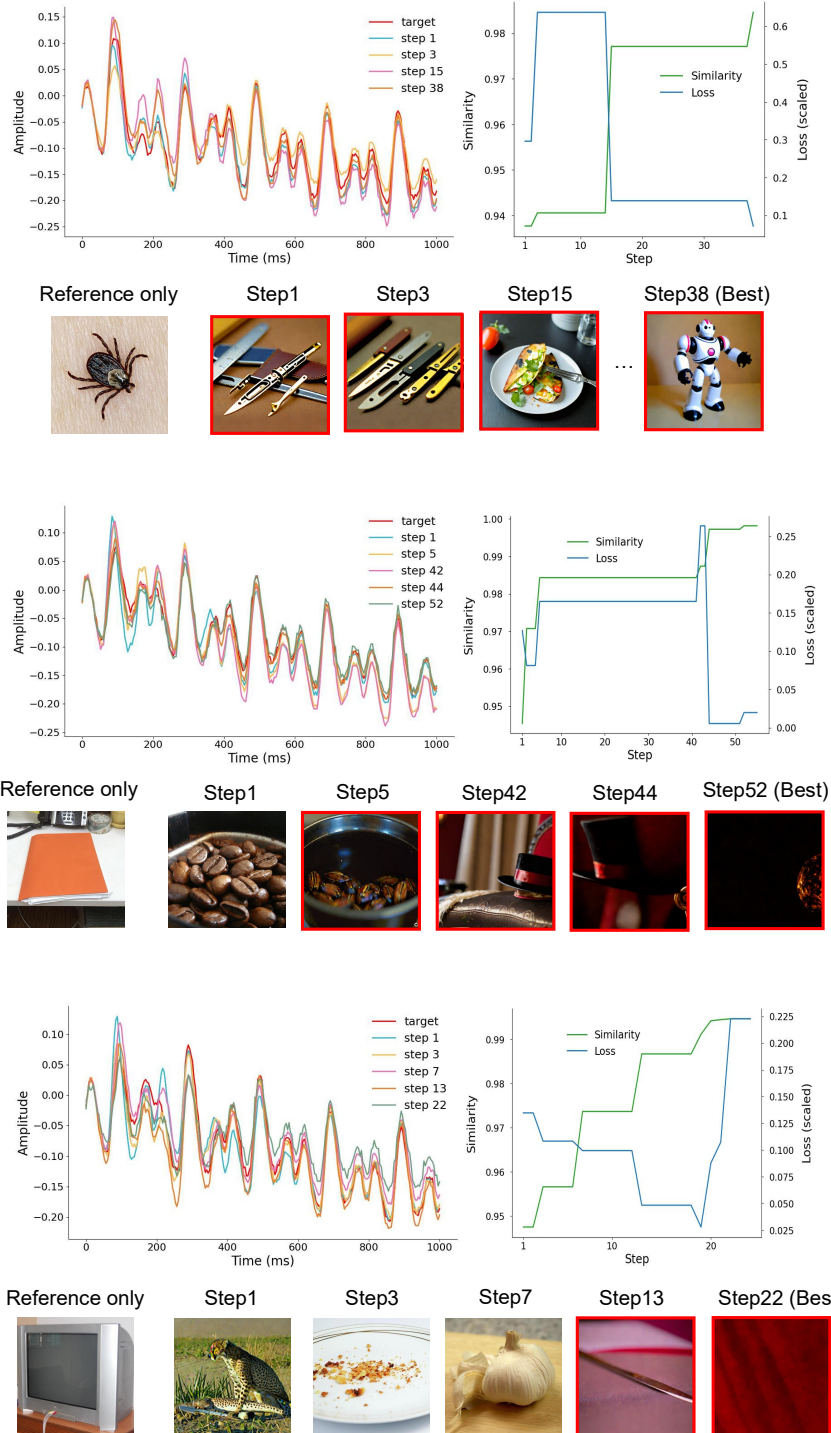


Figure A.18: **Illustration of the closed-loop iterative process for Subject 8.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.



1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

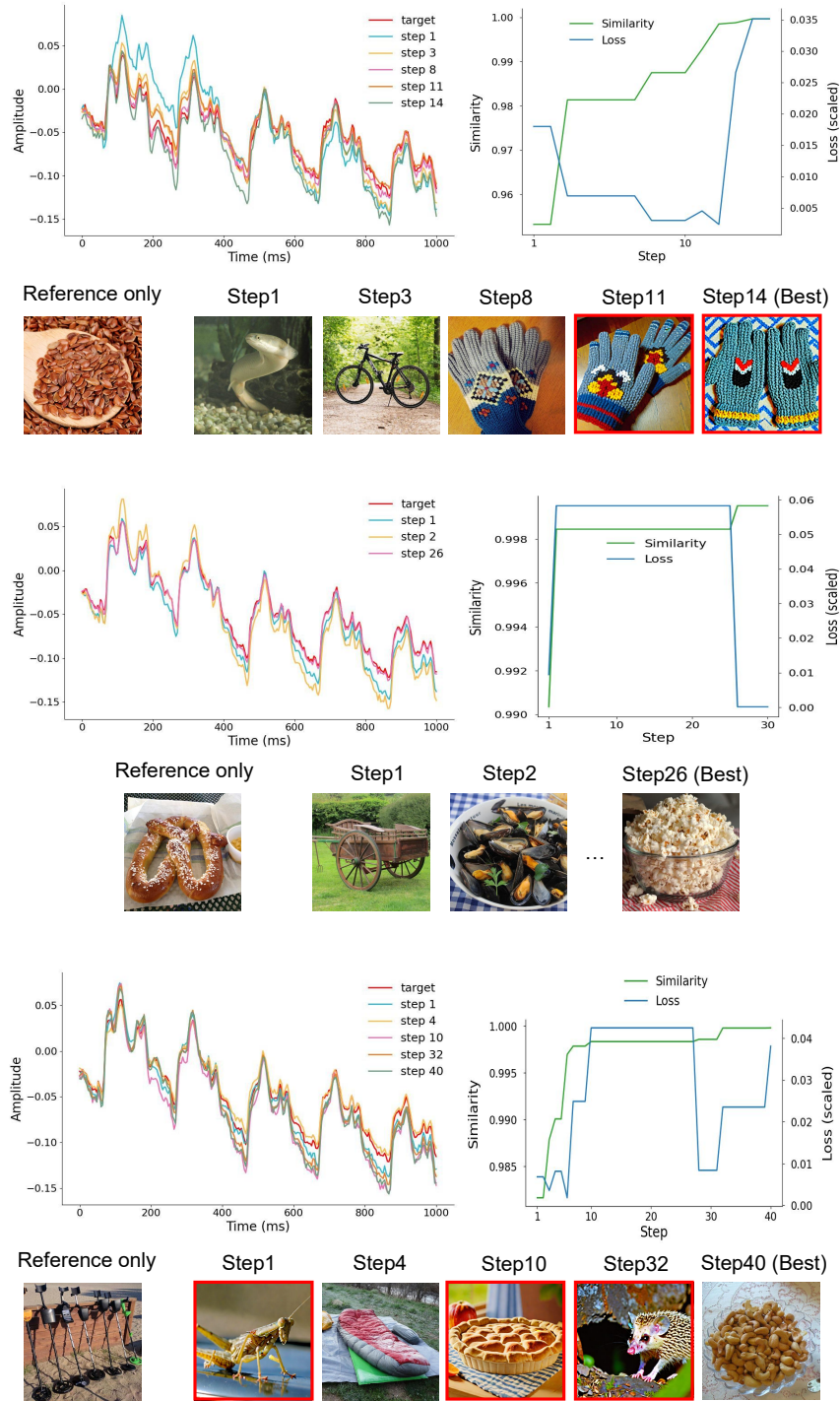


Figure A.19: **Illustration of the closed-loop iterative process for Subject 9.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

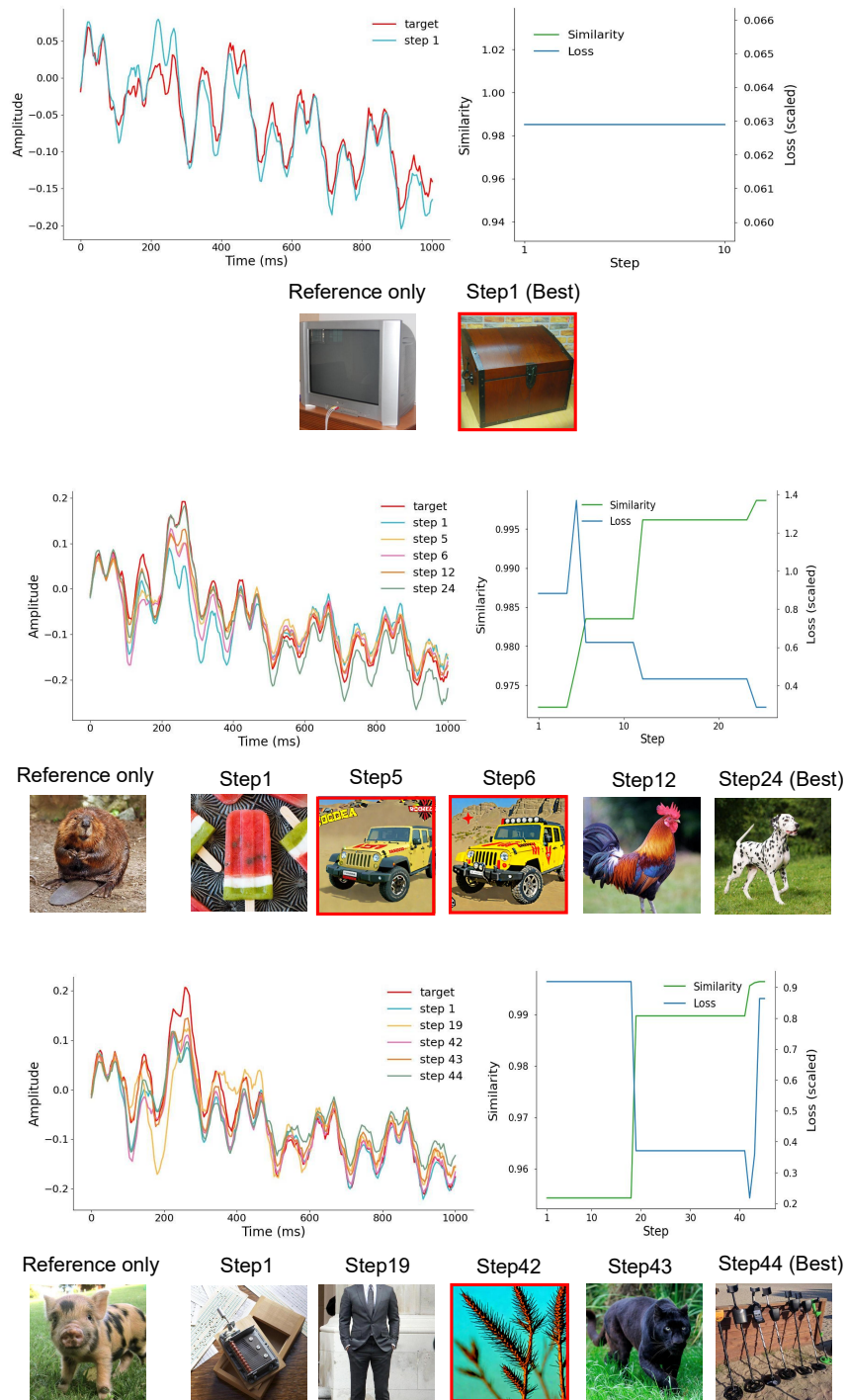


Figure A.20: **Illustration of the closed-loop iterative process for Subject 10.** Three distinct visual targets were presented, each based on a specific similarity measure (details in Target Features of EEG, Section 4.1), with new visual stimuli iteratively generated for each target. The left panel illustrates the time-domain evolution of neural responses across iterations. The right panel depicts the changes in similarity (green curve) and loss (blue curve, scaled) between the current stage features and the target features.