
Beware of Overestimated Decoding Performance Arising from Temporal Autocorrelations in Electroencephalogram Signals

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Researchers have reported high decoding accuracy (>95%) using non-invasive
2 Electroencephalogram (EEG) signals for brain-computer interface (BCI) decod-
3 ing tasks like image decoding, emotion recognition, auditory spatial attention
4 detection, etc. Since these EEG data were usually collected with well-designed
5 paradigms in labs, the reliability and robustness of the corresponding decoding
6 methods were doubted by some researchers, and they argued that such decoding
7 accuracy was overestimated due to the inherent temporal autocorrelation of EEG
8 signals. However, the coupling between the stimulus-driven neural responses and
9 the EEG temporal autocorrelations makes it difficult to confirm whether this over-
10 estimation exists in truth. Furthermore, the underlying pitfalls behind overesti-
11 mated decoding accuracy have not been fully explained due to a lack of appro-
12 priate formulation. In this work, we formulate the pitfall in various EEG decod-
13 ing tasks in a unified framework. EEG data were recorded from watermelons
14 to remove stimulus-driven neural responses. Labels were assigned to continuous
15 EEG according to the experimental design for EEG recording of several typical
16 datasets, and then the decoding methods were conducted. The results showed the
17 label can be successfully decoded as long as continuous EEG data with the same
18 label were split into training and test sets. Further analysis indicated that high
19 accuracy of various BCI decoding tasks could be achieved by associating labels
20 with EEG intrinsic temporal autocorrelation features. These results underscore
21 the importance of choosing the right experimental designs and data splits in BCI
22 decoding tasks to prevent inflated accuracies due to EEG temporal correlations.
23 The watermelon EEG dataset collected in this work can be obtained at Zenodo:
24 <https://zenodo.org/records/11238929>, and all the codes of this work can
25 be obtained in the supplementary materials.

26 1 Introduction and related works

27 A brain-computer interface (BCI) is a type of human-machine interaction that bridges a pathway
28 from the brain to external devices [1]. Electroencephalogram (EEG) has emerged as a valuable tool
29 for BCI because of its high time resolution, low cost, and good portability [2], and algorithms of
30 neural decoding from EEG signals play a role in its practical applications. Recently, deep learning
31 methods have been developed widely for various EEG decoding tasks, and high decoding accuracy
32 was reported. For example, in the task of decoding image classes with EEG recordings, when
33 subjects were required to watch images of different classes, a decoding accuracy of 82.90% was
34 reported for the 40-way classification by Spampinato et al. [3]. With their EEG dataset, subsequent

35 studies reported a higher decoding accuracy (98.30%, [4]), high performance on image retrieval, and
36 even image generation from EEG [5, 6, 7].

37 However, it remains unclear what kind of EEG features are learned by the DNN-based models. Some
38 researchers have posited that the high decoding accuracy on the image-evoked EEG dataset was
39 attributed to the block-design paradigm during EEG recording [8, 9, 10], in which 50 images with the
40 same class label were presented to the subject continuously in one block, and the 40 image-classes
41 were presented as 40 separate blocks. Due to the existence of temporal autocorrelation of EEG
42 signals, i.e., the temporally nearby data is more similar than the temporally distal [11, 12, 13, 14],
43 the models could learn the block-related features rather than the image-related.

44 To verify their concerns, Li et al. [8] recorded EEG with two experimental designs: block design
45 and rapid-event design. For the rapid-event design, images across the 40 classes were presented
46 alternately and randomly. When the same DNN model was used, it was found that the decoding
47 accuracy was close to Spampinato et al. [3] with the block-design EEG data, but it was dramati-
48 cally decreased to the chance-level (2.50%) with the rapid-event design data. Subsequent work also
49 confirmed the low decoding accuracy for EEG recorded with rapid-event design [9, 10]. However,
50 Palazzo et al. [15] proposed that temporal autocorrelations only play a marginal role in EEG de-
51 coding tasks because they found that EEG data recorded during rest periods (temporal proximity to
52 adjacent blocks) could not be successfully classified as the preceding block label or the succeeding
53 block label. They also argued that the rapid-event design seemed to weaken the image-related neural
54 responses due to the possible cognitive load and fatigue effect compared to the block design. Some
55 researchers [15, 16, 17, 18] pointed out that block design is essential because humans tend to react
56 more consistently and respond faster when conditions are presented in blocks [19, 20]. Wilson et
57 al. [18] advised that classification work that decodes from block design datasets is the most suitable
58 approach until advances are made to reduce noise.

59 Although the pitfall of overestimated decoding accuracy has been mainly discussed in image neural
60 decoding tasks, we noticed that similar pitfalls might also exist in various EEG decoding tasks such
61 as in auditory spatial attention detection (ASAD) tasks [21, 22, 23, 24], which involves decoding
62 the subjects auditory attention locus from neural data, and in emotion recognition task [25, 26, 27],
63 which involves recognizing the subjects emotion type from neural data. Researchers have also found
64 that splitting a continuous EEG from a specific experimental condition into training and test sets
65 would bring higher decoding accuracy in epilepsy detection tasks [28], motor imagery decoding
66 tasks [29], and so on. All those high decoding accuracy works share the common characteristic:
67 continuously recorded EEG data of a specific class (condition) label are divided into training and
68 test sets (see the top-left of Figure 1).

69 Although some studies have mentioned the overestimated decoding accuracy and tried to remind
70 the possible pitfall [8, 30], it is difficult to discriminate the influence of the inherent temporal au-
71 tocorrelation in EEG signals due to the coupling of stimuli-driven neural responses and the temporal
72 autocorrelations. More importantly, due to the lack of an effective formalization, there is not an
73 adequate explanation of how models utilize temporal autocorrelation features for decoding. Further-
74 more, their concerns only focused on one specific decoding task, and the results and conclusions
75 cannot be generalized to general BCI decoding tasks.

76 In this work, the pitfall of various EEG decoding tasks was formulated with a unified framework.
77 To completely decouple the temporal autocorrelation features from stimuli-driven neural responses,
78 EEG data were collected from 10 watermelons in this work to construct "Watermelon EEG". This
79 method is known as phantom EEG in previous studies [31, 32, 33, 34, 35, 36], and the EEG data
80 exclude stimulus-driven neural responses while reserving the temporal autocorrelation features. For
81 comparison, a human EEG dataset was also adopted. The watermelon EEG and human EEG
82 were reorganized into three classic neural decoding EEG datasets following their EEG experimen-
83 tal paradigm: image classification (CVPR, [3]), emotion classification (DEAP, [37]), and auditory
84 spatial attention decoding (KUL, [38]), resulting in six EEG datasets. A sample CNN-based decod-
85 ing model was used to complete the decoding tasks with the corresponding EEG dataset, and the
86 experimental results revealed that:

- 87 1. When the pitfall was formulated with a unique framework, and the temporal autocorre-
88 lation was defined as domain features, high decoding accuracy of various BCI decoding
89 tasks could be achieved by associating labels with EEG intrinsic temporal autocorrelation
90 features.

- 91 2. The pitfall exists not only in classification but also widely in EEG-image joint training
 92 without explicit labels and even image generation.
 93 3. Splitting a continuous EEG with the same class label into training and test sets should never
 94 be used in future BCI decoding works.

95 **2 Method**

96 The section is organized by: the pitfall is formulated in Subsection 2.1, and the datasets used are
 97 introduced in Subsection 2.2. Then, the methods to finish different classification tasks are introduced
 98 in Subsection 2.3, and joint training and image generation from EEG are introduced in Subsection
 99 2.4. Some implementation details and statistical analysis method are described in Subsection 2.5.

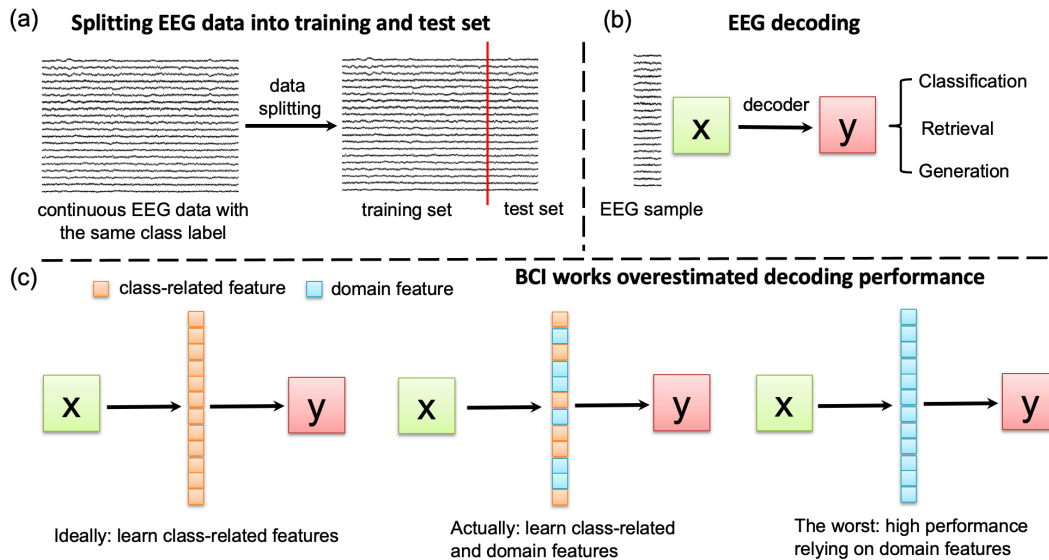


Figure 1: Overestimated decoding performance in BCI works. (a) Continuous EEG data in a certain experimental condition (with the same class label) are split into training and test sets for decoder training and evaluation. (b) With the test EEG sample input, the decoder gives output in the forms of classification, retrieval, and generation. (c) Decoders may use both domain features or class-related features for decoding.

100 **2.1 Problem Formulation**

101 In some BCI works on domain generalization [39], all EEG data from a dataset [40] or from a subject
 102 [41] are usually regarded as a domain to emphasize EEG pattern distribution differences between
 103 datasets or subjects. Adopted from this concept, we regard a period of continuous EEG data with
 104 the same class label as a domain. In some BCI works [3, 4, 21, 22, 23, 24, 25, 26, 27], researches
 105 segment the EEG data from the same domain into samples and further split the samples into training
 106 and test data (as shown in Figure 1a) and complete decoding task, such as classification, retrieval
 107 and generation (as shown in Figure 1b). In these cases, the models used in these works would learn
 108 the coupled features containing the class-related feature and domain feature (as shown in the middle
 109 of the Figure 1c). The underlying assumption of these works is that the domain feature plays only a
 110 margin role in EEG decoding tasks as shown in the left of the Figure 1c. However, we assumed that
 111 the domain feature contributes to the high decoding accuracy as shown in the right of the Figure 1c,
 112 which is the pitfall we mentioned in Section 1.

113 To validate our assumption, we need to formulate the pitfall. Denote D as the domain set, and each
 114 domain $d \in D$ contains many samples. We use S^d to denote the sample set of the domain d . The
 115 notation x_i^d represents the i -th sample (e.g., a 0.5-second EEG data corresponding to watching a
 116 specific image) of domain d , which is associated with class y_i^d (e.g., the class label panda of the

117 watched image). Considering the temporal autocorrelation of the EEG data, the domain features of
 118 data within the same domain are more similar, while the domain features of data in different domains
 119 are more distinct.

120 For EEG decoding tasks, we assume the data is generated from a two-stage process. First, each
 121 domain is modeled as a latent factor z sampled from some meta domain distribution $p(\cdot)$. Second,
 122 each data sample x is sampled from a sample distribution conditioned on the domain z and class y :

$$z \sim p(\cdot), x \sim p(\cdot|z, y) \quad (1)$$

123 Given the sample x , the aim of a specific EEG decoding task is to uncover its true class label using
 124 the posterior $p(y|x)$. The quantity can be factorized by the domain factor z as,

$$p(y|x) = \int p(y, z|x) dz = \int p(y|x, z)p(z|x) \quad (2)$$

125 When we use the Watermelon EEG dataset or use a dataset that is completely unrelated to the
 126 current task (e.g., decoding images from an auditory EEG dataset), the class-related feature has
 127 none possibility to exist in EEG samples. In this condition, $p(y|x, z) = p(y|z)$ and the equation (2)
 128 can be modified as:

$$p(y|x) = \int p(y, z|x) dz = \int p(y|z)p(z|x) \quad (3)$$

129 The assumption of this work is that the model could also deduce $p(y|x)$ by learning $p(y|z)$ and
 130 $p(z|x)$ even there is none class-related feature exists. In other words, we assumed that it could also
 131 achieve high decoding accuracy on different EEG decoding tasks when using the Watermelons EEG
 132 dataset.

133 2.2 Dataset

134 **Watermelon EEG Dataset** Ten watermelons were selected as subjects. EEG data were recorded
 135 with a NeuroScan SynAmps2 system (Compumedics Limited, Victoria, Australia), using a 64-
 136 channel Ag/AgCl electrodes cap with a 10/20 layout. An additional electrode was placed on the
 137 lower part of the watermelon as the physiological reference, and the forehead served as the ground
 138 site (see Appendix A.1 for photography). The inter-electrode impedances were maintained under
 139 20 kOhm. Data were recorded at a sampling rate of 1000 Hz. EEG recordings for each watermelon
 140 lasted for more than 1 hour to ensure sufficient data for the decoding task. We refer to the dataset
 141 consisting of EEG recordings of 10 watermelons as the Watermelon EEG Dataset.

142 **SparrKULee Dataset** SparrKULee dataset[42] is a speech-evoked EEG dataset from the KU Leu-
 143 ven University containing 64-channel EEG recordings from 85 participants, each of whom listened
 144 to 90-150 minutes of natural speech. We used this dataset because EEG recordings were longer than
 145 1 hour to ensure a sufficient amount of data for each subject. To match the number of subjects in
 146 the Watermelon EEG Dataset, EEG data from 10 subjects (ID: Sub7-Sub16) from the SparrKULee
 147 Dataset were used.

148 **Dataset reorganization and dataset segmentation** The term "reorganization" refers to segmenting
 149 continuous EEG into samples and assigning each sample a class label and a domain label according
 150 to the referenced experimental design. Here, we follow the experimental designs of three classical
 151 published EEG datasets to reorganize the Watermelon EEG Dataset and SparrKULee Dataset. These
 152 three datasets were collected respectively for image decoding, emotion recognition, and ASAD
 153 tasks.

154 For the image decoding task, we referred to the experimental design of the CVPR dataset [3]. For
 155 the CVPR dataset, 40 classes of images were presented in a block-design paradigm. Specifically, 50
 156 different images of the same class were presented continuously in a block, with each image lasting
 157 for 0.5 second, resulting in 40 blocks of presentation for each subject. The 0.5-second length EEG
 158 data of the same class were split into training, validation, and test sets in a ratio of 8:1:1 [4, 3].
 159 Following this experimental design and dataset segmentation, we segment continuous EEG from

160 the Watermelon EEG Dataset and SparrKULee Dataset into blocks and assign a unique class label
 161 and a unique domain label for each block. The interval between adjacent blocks is set to 10 seconds
 162 to match the rest time of the subjects during the EEG recording in the CVPR dataset. Then, EEG
 163 data in each block are further segmented into 50 0.5-s length samples. Since the EEG data in the
 164 CVPR dataset has 128 channels, we replicated our 64-channel EEG in the channel dimension. The
 165 reorganized datasets for Watermelon Dataset and SparrKULee Dataset are called WM-CVPR and
 166 SK-CVPR, respectively. Here, we use the "A-B" naming format, where the left side of "-" represents
 167 the source dataset (WM: watermelon dataset, SK: SparrKULee Dataset), and the right side of "-"
 168 represents the dataset of which the experimental design is referenced. For the emotion recognition
 169 task and ASAD task, the DEAP dataset and the KUL dataset are used as the referenced dataset,
 170 resulting in WM-DEAP, SK-DEAP, WM-KUL, and SK-KUL. More details for reorganization can
 171 be found in Appendix A.2.

172 2.3 Classification tasks

173 **Model.** To demonstrate that domain features are strong and easy to be learned by the network,
 174 we used a simple CNN (or some parts of this CNN) to complete all classification tasks mentioned
 175 in this work. The CNN network includes a layer-norm layer, a 2D-convolutional layer (output
 176 channel: 100), an averaging pooling layer, and two fully connected layers. The kernel size of the
 177 2D-convolutional layer depends on the channel number and sampling frequency of the input EEG.
 178 The node number of the output fully connected layer depends on the number of classes.

179 **Decoding the domain feature** To demonstrate that the model can predict the domain factor z from
 180 EEG input sample x , which relates to learning posterior $p(z|x)$, a domain label classification was
 181 adopted on the six datasets (i.e., WM-CVPR, WM-DEAP, WM-KUL, SK-CVPR, SK-DEAP and
 182 SK-KUL dataset) with a simple CNN classifier. The splitting strategy leave-samples-out was used,
 183 which means that all sample were randomly split into training set, validation set and test set. The
 184 outputs after the averaging pooling layer were selected as domain feature representation, and t-SNE
 185 was utilized for dimensionality reduction and visualization.

186 **Decoding the class label from the domain feature** To demonstrate that the model can predict
 187 the class label y from the domain factor z , which relates to learning posterior $p(y|z)$, a class label
 188 classification was adopted on the four datasets (classification on the WM-CVPR dataset and SK-
 189 CVPR dataset are unnecessary since domain labels and class labels are one-to-one correspondence)
 190 using a single network with two linear layers and an intermediate sigmoid function.

191 **End-to-end classification** To demonstrate that the model can predict the class label y from the EEG
 192 input sample x directly when samples in the training set and test set are from common domains,
 193 a class label classification was adopted on the six datasets with the simple CNN classifier. The
 194 splitting strategy leave samples out was used. Classification on the WM-CVPR dataset and SK-
 195 CVPR dataset is the same since domain labels and class labels in the two datasets are one-to-one
 196 correspondence. To demonstrate that the model indeed used the domain feature to complete the
 197 end-to-end classification, the splitting strategy leave domains out was used on the four datasets (i.e.,
 198 WM-DEAP, WM-KUL, SK-DEAP, and SK-KUL dataset) in which samples in the same domain
 199 only appear in the training set or the test set.

200 **Zero-shot classification** In a recent work [4], EEG data from 34 classes within the CVPR2017
 201 dataset were used to train an EEG encoder, and the remaining 6 unseen classes were used for test-
 202 ing. The results showed that features of different unseen classes clustered in distinct groups on the
 203 two-dimensional t-SNE plane. Similar analyses were conducted on the SK-CVPR and WM-CVPR
 204 datasets. Six classes were selected for testing, and the remaining 34 classes were for training. The
 205 simple CNN was used to predict class labels from input EEG samples, and the outputs from the av-
 206 erage pooling layer were chosen as the EEG feature representation. Two strategies were employed
 207 for selecting the 6 test classes: random selection and first-six selection. For random selection, the 6
 208 test classes are randomly chosen from the 40 classes. For the first-six selections, the first presented
 209 6 classes in the EEG experiment are chosen. During the test stage, since the training set does not in-
 210 clude classes corresponding to the test EEG data, the model could not give the corresponding labels
 211 and could only output the most probable classes among the 34 seen during training. Therefore, we
 212 proposed two evaluation metrics: Acc_{near} and Acc_{7th} . Acc_{near} represents the proportion of EEG
 213 data classified into temporally adjacent classes, while Acc_{7th} represents the proportion classified
 214 into the category presented seventh in time.

215 2.4 Joint training and image generation

216 To demonstrate that the model can utilize domain features to accomplish retrieval and generation
217 besides classification, EEG-image joint training and image generation on WM-CVPR and SK-CVPR
218 were conducted.

219 **Joint training** In the EEG-image joint training, a pre-trained image encoder was typically utilized
220 to extract image representation, while an EEG encoder was employed to extract EEG features to
221 align with the image representation. During the decoding process, a retrieval task was applied.
222 Specifically, given a test EEG sample and a collection of images containing the target and the non-
223 target. The image representation was reconstructed from the EEG with the EEG encoder. The
224 similarity between the reconstructed image representation and all candidate image representations
225 in the collection is calculated. The decoded output image is selected based on the ranking of these
226 similarities. Usually, the Top-k accuracy and normalized Rank accuracy are used as evaluation
227 metrics. In this work, the simple CNN described in Subsection 2.3 is used as an EEG encoder. The
228 detailed implementation can be found in Appendix A.3.

229 **Image generation** The image generation aims to generate images seen by the subjects from their
230 EEG data. This task commonly uses a two-stage process: EEG encoding and image generation.
231 In the EEG encoding stage, a model is built to encode EEG data into a latent representation. In
232 the image generation stage, a pre-trained image generator is used. The generator is fine-tuned with
233 EEG representation and corresponding images. In this work, the EEG data are first encoded into
234 image representation with a simple CNN described in Subsection 2.3. Following previous work[43],
235 a latent diffusion model conditioned on image representation was used. The metric of n-way top-k
236 accuracy was used for evaluating the semantic correctness of generated images [44]. The detailed
237 implementation can be found in Appendix A.4.

238 2.5 Implement details

239 The neural networks were implemented with the Pytorch and trained on a single high-performance
240 computing node with 8 A800 GPU. For the classification task, the AdamW [45] optimizer was em-
241 ployed to minimize the cross-entropy loss function with a learning rate of 10^{-3} . For the joint training
242 and image generation, the AdamW optimizer was used with a learning rate of 10^{-3} and 5×10^{-4}
243 for each task respectively. More details can be found in our codes. All the experiments mentioned
244 in this work were trained within the subjects (i.e., models were trained for each subject respectively)
245 except special annotation (unseen subject decoding results were only presented in Appendix A.5).
246 For statistical analysis, the one-sample t-test was used to check whether the reported results were
247 significantly higher than the chance level. Bonferroni correction was used to adjust the p -value. A
248 p -value of 0.05 or lower was considered statistically significant.

249 3 Results

250 3.1 Classification tasks

251 The results shown in Table 1 present that classification accuracy in domain label classification and
252 class label classification are all significantly above the chance level. This shows that the domain
253 feature can be extracted effectively with a simple CNN, and the label class can be decoded from
254 the extracted domain features or from EEG directly. In contrast, the decoding accuracy drops to the
255 chance level when using the splitting strategy leave-domains-out, further supporting domain feature-
256 induced high decoding accuracy. The standard error of the mean calculated over the subjects level
257 is reported for accuracy in this work.

258 Figures 2a and 2b show the t-SNE plot for domain label classification and end-to-end class label
259 classification. As shown in Figure 2a, 8 distinct clusters exist, each corresponding to one domain.
260 In Figure 2b, 8 distinct clusters also exist, with four corresponding to class label 1 and the other
261 four corresponding to class label 2. This indicates that the high decoding accuracy results from
262 associating class labels with domain features.

Table 1: Classification accuracy (%) on the six datasets. DLC is for domain label classification. TLC-DF is for class label classification from domain features. TLC-EEG is for end-to-end class label classification. TLC-EEG-woDO is for class label classification direct from EEG when samples in the training set and test set are from different domains.

	WM-CVPR	WM-DEAP	WM-KUL	SK-CVPR	SK-DEAP	SK-KUL
DLC	88.78 ± 4.95	96.98 ± 0.76	99.99 ± 0.01	69.83 ± 2.98	72.70 ± 1.36	100.00 ± 0.00
DLC (chance level)	2.50	2.50	12.50	2.50	2.50	12.50
TLC-DF	-	92.77 ± 1.31	100.00 ± 0.00	-	76.19 ± 1.80	100.00 ± 0.00
TLC-EEG	88.78 ± 4.95	88.74 ± 3.26	82.74 ± 6.44	69.83 ± 2.98	74.44 ± 2.76	93.34 ± 2.01
TLC-EEG-woDO	-	24.67 ± 2.31	49.97 ± 4.67	-	25.34 ± 1.85	59.32 ± 4.07
TCL (chance level)	2.50	25.00	50.00	2.50	25.00	50.00

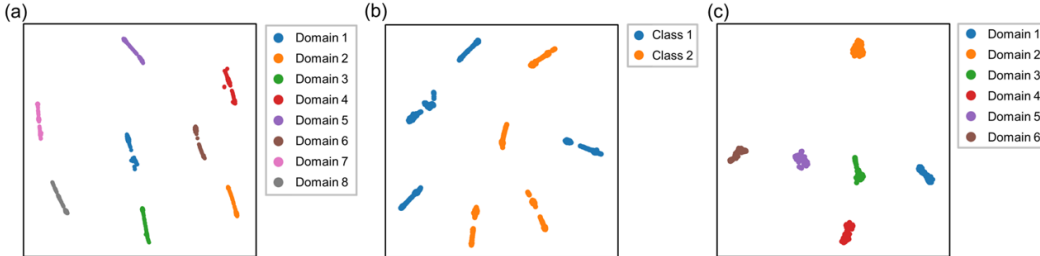


Figure 2: t-SNE plot for (a) domain label classification, (b) end-to-end class label classification, and (c) zero-shot class label classification

263 The experimental results for zero-shot classification are displayed in Table 2. It can be observed
 264 that the model tended to classify test samples into temporally adjacent classes. Figure 2c shows the
 265 t-SNE visualization of the unseen EEG features extracted from the decoder. Despite being unseen,
 266 different domains of features clustered in distinct groups. This suggests that the decoder just learned
 267 to extract EEG domain features during training and distinguish unseen EEG responses from the
 268 domain features.

Table 2: Zero-shot EEG classification accuracy (%) on WM-CVPR and SK-CVPR datasets.

	WM-CVPR first-six	WM-CVPR random	SK-CVPR first-six	SK-CVPR random
Acc_{near}	-	79.43 ± 5.61	-	78.00 ± 5.66
Acc_{7th}	69.60 ± 10.64	6.73 ± 3.24	77.03 ± 11.32	0.87 ± 0.82

269 **3.2 Joint training and image generation**

270 For EEG-image joint training, Table 3 displays the accuracy for the retravel task on the test set. The
 271 table shows that, for both types of loss functions, decoding accuracy is far above the chance level,
 272 demonstrating that the model can utilize domain features to align EEG with image features. Table 3
 273 Result for joint training on WM-CVPR and SK-CVPR with a loss function of cosine similarity (CS)
 274 or InfoNCE.

Table 3: Accuracy (%) for joint training on WM-CVPR and SK-CVPR with a loss function of cosine similarity (CS) or InfoNCE.

	WM-CVPR		SK-CVPR		Chance level
	CS loss	InfoNCE loss	CS loss	InfoNCE loss	
Top1 Acc	81.40 ± 9.25	90.15 ± 5.45	80.70 ± 0.60	79.70 ± 0.92	2.50
Top5 Acc	90.65 ± 5.82	98.56 ± 1.09	88.86 ± 1.03	92.39 ± 0.38	12.50
Rank Acc	95.87 ± 2.51	99.42 ± 0.38	95.20 ± 0.24	98.09 ± 0.07	50.00

275 For image generation, Table 4 displays the n-way top-k accuracy for the generated images on the
 276 WM-CVPR and SK-CVPR datasets. The metrics are significantly above the chance level, indicating

277 that the generated images have correct semantics. Figure 3 shows some generated images on the
 278 WM-CVPR dataset. As shown in the figure, the model can exactly generate the correct images. The
 279 results on EEG-image joint training and image generation show that in addition to classification
 280 tasks, retrieval, and generation can also achieve high performance by leveraging domain features
 281 shared by the test and training sets.

Table 4: Accuracy (%) for semantic correctness. The repeated times N was set to 50.

-	Top-1/50-way	Top-5/50-way	Top-1/100-way	Top-5/100-way
WM-CVPR	26.77 ± 3.37	46.44 ± 4.60	21.64 ± 2.89	38.11 ± 4.30
SK-CVPR	25.04 ± 0.93	43.61 ± 0.88	20.37 ± 0.91	35.35 ± 0.89
Chance	2.00	10.00	1.00	5.00



Figure 3: EEG-generated image from a typical watermelon subject, where the first column of each panel represents the real images "watched" by the watermelon subject, and the following five columns show the images generated by the model.

282 4 Discussion

283 4.1 Relying on the domain features for EEG decoding

284 While many works on EEG decoding have reported high-performance results, we proposed that
 285 some of these high-performance may rely on temporal autocorrelation of EEG data. The pitfall may
 286 involve different EEG decoding tasks. To clarify this pitfall, the concept of domain was adopted
 287 to describe the temporal autocorrelation of a continuous EEG with the same label. EEG data were
 288 collected from watermelon as the phantom to exclude the contribution of stimuli-driven neural re-
 289 sponses to decoding results. The results showed that a simple CNN network could well learn domain
 290 features from EEG data and could associate class labels with domain features.

291 To avoid the pitfalls, a feasible approach is to adopt a reasonable data-splitting strategy to avoid train-
 292 ing and test sets sharing the common domain features, i.e., a leave-domains-out splitting strategy.
 293 For instance, a leave-subjects-out data-splitting strategy can be adopted, which entails designating
 294 the data from certain participants for training and data from others for testing. Alternatively, for
 295 datasets that do not follow a block design, a leave-trials-out strategy may be applied. Prior research
 296 has consistently demonstrated that employing a leave-subjects-out splitting strategy precipitates a
 297 notable decline in decoding performance [46]. In some cases, it has been reported that decoding
 298 accuracy dropped to the chance level [47, 8]. The prevalent interpretation is that inter-individual
 299 variability [46] hampers the generalizability across different subjects. However, we posit that the
 300 observed decrement in decoding accuracy is attributable to model overfitting to domain features.
 301 Although the leave-subjects-out partitioning strategy is designed to prevent the leakage of domain
 302 features, the presence of these domain features in the training set can still lead the model to inadver-
 303 tently exploit them to differentiate between categories during the training phase. The methods and
 304 results further support the conclusion can be found in Appendix A.5

305 Palazzo et al. [15] proposed that the EEG temporal correlation related to baseline drift could be al-
 306 leviated by high-pass filtering. However, our further experiment proved that the domain feature still
 307 exists and that high decoding accuracy could be achieved in any frequency band (see Appendix A.6).
 308 We argue that the focus should not be exclusively on the elimination of EEG autocorrelation through

309 filtering. Instead, greater emphasis should be placed on the experimental paradigms of EEG record-
310 ing and the methods employed for dataset splitting. By addressing these aspects, we can proactively
311 prevent the overestimated decoding accuracy arising from EEG temporal autocorrelations.

312 It is worth noting that we do not want to create an illusion that all BCI works utilize EEG temporal
313 autocorrelation features for decoding. In fact, there are many works that do not rely on EEG temporal
314 autocorrelation features for decoding in image decoding [48, 49, 50] emotion recognition [51], sleep
315 detection [40, 41] and ASAD [52]. These works demonstrated the feasibility of various BCI tasks.

316 4.2 Potential sources of domain features

317 In this work, we have demonstrated the existence of EEG temporal autocorrelation in the water-
318 melon EEG, which consists of no neural activities, and in the human EEG data. Li et al. [8] believed
319 the model decodes by utilizing the baseline drift in the CVPR2017 dataset. They found that when
320 the EEG data is filtered with a bandpass filter, the decoding accuracy dropped greatly. Palazzo et
321 al. [15] also claimed that temporal correlation was strong only in low frequency. However, we have
322 demonstrated in Appendix A.4 that the domain feature still exists and that high decoding accuracy
323 can be achieved in any frequency band. In addition to baseline drift, some neuroscience works have
324 shown that temporal autocorrelation existed in neural oscillation, which could be reflected in EEG
325 in various frequency bands. This is referred to as Long-Range Temporal Correlations (LRTC) in
326 neuroscience research [11, 12, 13, 14]. Linkenkaer-Hansen et al. [13] first calculated the LRTC in
327 resting-state EEG data. They found that spontaneous alpha, mu, and beta oscillations result in signif-
328 icant LRTC for at least several hundred seconds during resting conditions. Subsequent neuroscience
329 research further demonstrated that significant LRTC exists in the theta [11] and gamma [12] bands.
330 While baseline drift can be removed through filtering, the frequency range of the LRTC overlaps
331 with the frequency range of stimuli-driven neural responses, making it impossible to remove this
332 domain feature through filtering. Temporal correlation analysis on human EEG in the SparrKULee
333 Dataset showed the existence of strong LRTC in all frequency bands, and the LRTC in a narrowband
334 is sufficient to complete the corresponding decoding task. The methods and results further support
335 the conclusion can be found in Appendix A.7.

336 4.3 Limitation and future work

337 Although direct evidence of overestimated decoding accuracy attributable to domain feature across
338 various brain-computer interface (BCI) tasks have been provided in the current work, no solution has
339 been proposed to mitigate overfitting to domain features in the training set. Some works have already
340 used domain adaptation [2, 53, 54] or domain generalization [40, 41] method to improve decoding
341 accuracy under leave-subjects-out data splitting in BCI tasks. This may also help alleviate the ad-
342 verse effects of domain features on decoding tasks. It is also noteworthy to highlight the remarkable
343 efficacy of large-scale EEG model in various BCI decoding tasks [55, 56, 57]. Given that domain
344 features are pervasive in extensive EEG datasets and do not necessitate manually annotated labels,
345 self-supervised pre-trained large EEG models may be especially adept at discerning and neutralizing
346 domain features, thereby facilitating more robust and generalizable decoding performance.

347 5 Conclusion

348 In this work, the “overestimated decoding accuracy pitfall” in various EEG decoding tasks is for-
349 mulated in a unified framework by adopting the concept of “domain”. Some typical EEG decoding
350 tasks (image decoding, emotion recognition, and auditory spatial attention detection) are conducted
351 on the self-collected watermelon EEG dataset. The results showed that EEG data from different
352 domains have distinctive domain features induced by EEG temporal autocorrelations. Using the in-
353 appropriate data partitioning strategy, high decoding accuracy is achieved by associating class labels
354 with domain features. The results will draw attention to the high decoding performance caused by
355 EEG temporal correlation and guide the development of BCI in a positive direction.

References

- [1] Yue-Ting Pan, Jing-Lun Chou, and Chun-Shu Wei. Matt: A manifold attention network for eeg decoding. *Advances in Neural Information Processing Systems*, 35:31116–31129, 2022.
- [2] Reinmar Kobler, Jun-ichiro Hirayama, Qibin Zhao, and Motoaki Kawanabe. Spd domain-specific batch normalization to crack interpretable unsupervised domain adaptation in eeg. *Advances in Neural Information Processing Systems*, 35:6219–6235, 2022.
- [3] C. Spampinato, S. Palazzo, I. Kavasidis, D. Giordano, N. Souly, and M. Shah. Deep learning human mind for automated visual classification. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4503–4511, 2017.
- [4] Prajwal Singh, Dwip Dalal, Gautam Vashishtha, Krishna Miyapuram, and Shanmuganathan Raman. Learning robust deep visual representations from eeg brain recordings. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 7553–7562, 2024.
- [5] Isaak Kavasidis, Simone Palazzo, Concetto Spampinato, Daniela Giordano, and Mubarak Shah. Brain2image: Converting brain signals into images. In *Proceedings of the 25th ACM international conference on Multimedia*, MM 17, pages 1809–1817, New York, NY, USA, 2017. Association for Computing Machinery.
- [6] S. Palazzo, C. Spampinato, I. Kavasidis, D. Giordano, and M. Shah. Generative adversarial networks conditioned by brain signals. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 3430–3438, 2017.
- [7] Praveen Tirupattur, Yogesh Singh Rawat, Concetto Spampinato, and Mubarak Shah. Thoughtviz: Visualizing human thoughts using generative adversarial network. In *Proceedings of the 26th ACM international conference on Multimedia*, MM 18, pages 950–958, New York, NY, USA, 2018. Association for Computing Machinery.
- [8] Ren Li, Jared S. Johansen, Hamad Ahmed, Thomas V. Ilyevsky, Ronnie B. Wilbur, Hari M. Bharadwaj, and Jeffrey Mark Siskind. The perils and pitfalls of block design for eeg classification experiments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):316–333, 2021.
- [9] Hamad Ahmed, Ronnie B. Wilbur, Hari M. Bharadwaj, and Jeffrey Mark Siskind. Object classification from randomized eeg trials. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3844–3853, 2021.
- [10] Hari M Bharadwaj, Ronnie B. Wilbur, and Jeffrey Mark Siskind. Still an ineffective method with supertrials/erpscomments on decoding brain representations by multimodal learning of neural activity and visual features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(11):14052–14054, 2023.
- [11] Luc Berthouze, Leon M. James, and Simon F. Farmer. Human eeg shows long-range temporal correlations of oscillation amplitude in theta, alpha and beta bands across a wide age range. *Clinical Neurophysiology*, 121(8):1187–1197, 2010.
- [12] Mona Irrmischer, Simon-Shlomo Poil, Huibert D. Mansvelder, Francesca Sangiuliano Intra, and Klaus Linkenkaer-Hansen. Strong long-range temporal correlations of beta/gamma oscillations are associated with poor sustained visual attention performance. *European Journal of Neuroscience*, 48(8):2674–2683, 2018.
- [13] Klaus Linkenkaer-Hansen, Vadim V. Nikouline, J. Matias Palva, and Risto J. Ilmoniemi. Long-range temporal correlations and scaling behavior in human brain oscillations. *Journal of Neuroscience*, 21(4):1370–1377, 2001.
- [14] Vadim V. Nikulin and Tom Brismar. Long-range temporal correlations in alpha and beta oscillations: effect of arousal level and test-retest reliability. *Clinical Neurophysiology*, 115(8):1896–1908, 2004.

- 404 [15] Simone Palazzo, Concetto Spampinato, Joseph Schmidt, Isaak Kavasidis, Daniela Giordano,
405 and Mubarak Shah. Correct block-design experiments mitigate temporal correlation bias in
406 eeg classification. *arXiv preprint arXiv:2012.03849*, 2020.
- 407 [16] Jacopo Cavazza, Waqar Ahmed, Riccardo Volpi, Pietro Morerio, Francesco Bossi, Cesco
408 Willemse, Agnieszka Wykowska, and Vittorio Murino. Understanding action concepts from
409 videos and brain activity through subjects consensus. *Scientific Reports*, 12(11):19073, 2022.
- 410 [17] Alankrit Mishra, Nikhil Raj, and Garima Bajwa. Eeg-based image feature extraction for vi-
411 sual classification using deep learning. In *2022 International Conference on Intelligent Data
412 Science Technologies and Applications (IDSTA)*, pages 181–188, 2022.
- 413 [18] Holly Wilson, Xi Chen, Mohammad Golbabaee, Michael J. Proulx, and Eamonn O'Neill. Fea-
414 sibility of decoding visual information from eeg. *Brain-Computer Interfaces*, 0(0):1–28, 2023.
- 415 [19] Lauren E. Ethridge, Shefali Brahmabhatt, Yuan Gao, Jennifer E. McDowell, and Brett A.
416 Clementz. Consider the context: Blocked versus interleaved presentation of antisaccade tri-
417 als. *Psychophysiology*, 46(5):1100–1107, 2009.
- 418 [20] Nelson A. Roque, Timothy J. Wright, and Walter R. Boot. Do different attention capture
419 paradigms measure different types of capture? *Attention, Perception & Psychophysics*,
420 78(7):2014–2030, 2016.
- 421 [21] Enze Su, Siqi Cai, Longhan Xie, Haizhou Li, and Tanja Schultz. Stanet: A spatiotemporal
422 attention network for decoding auditory spatial attention from eeg. *IEEE Transactions on
423 Biomedical Engineering*, 69(7):2233–2242, 2022.
- 424 [22] Saurav Pahuja, Siqi Cai, Tanja Schultz, and Haizhou Li. Xanet: Cross-attention between eeg
425 of left and right brain for auditory attention decoding. In *2023 11th International IEEE/EMBS
426 Conference on Neural Engineering (NER)*, pages 1–4, 2023.
- 427 [23] Xiran Xu, Bo Wang, Yujie Yan, Xihong Wu, and Jing Chen. A densenet-based method for
428 decoding auditory spatial attention with eeg. In *IEEE International Conference on Acoustics,
429 Speech and Signal Processing (ICASSP)*, pages 1946–1950, 2024.
- 430 [24] Qinke Ni, Hongyu Zhang, Cunhang Fan, Shengbing Pei, Chang Zhou, and Zhao Lv. Dbpnet:
431 Dual-branch parallel network with temporal-frequency fusion for auditory attention detection.
432 In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2024.
- 433 [25] Liangliang Hu, Congming Tan, Jiayang Xu, Rui Qiao, Yilin Hu, and Yin Tian. Decoding
434 emotion with phaseamplitude fusion features of eeg functional connectivity network. *Neural
435 Networks*, 172:106148, 2024.
- 436 [26] Jiayang Xu, Wenxia Qian, Liangliang Hu, Guangyuan Liao, and Yin Tian. Eeg decoding
437 for musical emotion with functional connectivity features. *Biomedical Signal Processing and
438 Control*, 89:105744, 2024.
- 439 [27] Zhi Zhang, Shenghua Zhong, and Yan Liu. Beyond mimicking under-represented emotions:
440 Deep data augmentation with emotional subspace constraints for eeg-based emotion recog-
441 nition. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(99):10252–10260,
442 2024.
- 443 [28] Geoffrey Brookshire, Jake Kasper, Nicholas M. Blauch, Yunan Charles Wu, Ryan Glatt,
444 David A. Merrill, Spencer Gerrol, Keith J. Yoder, Colin Quirk, and Ché Lucero. Data leak-
445 age in deep learning studies of translational eeg. *Frontiers in Neuroscience*, 18, 2024.
- 446 [29] Hamdi Altaheri, Ghulam Muhammad, Mansour Alsulaiman, Syed Umar Amin, Ghadir Ali
447 Altuwajjri, Wadood Abdul, Mohamed A. Bencherif, and Mohammed Faisal. Deep learning
448 techniques for classification of electroencephalogram (eeg) motor imagery (mi) signals: a re-
449 view. *Neural Computing and Applications*, 35(20):14681–14722, 2023.
- 450 [30] Iustina Rotaru, Simon Geirnaert, Nicolas Heintz, Iris Van de Ryck, Alexander Bertrand, and
451 Tom Francart. What are we really decoding? unveiling biases in eeg-based decoding of the
452 spatial focus of auditory attention. *Journal of Neural Engineering*, 21(1):016017, 2024.

- 453 [31] Mukund Balasubramanian, William M. Wells, John R. Ives, Patrick Britz, Robert V. Mulk-
454 ern, and Darren B. Orbach. Rf heating of gold cup and conductive plastic electrodes during
455 simultaneous eeg and mri. *The Neurodiagnostic Journal*, 57(1):69–83, 2017.
- 456 [32] Maximillian K. Egan, Ryan Larsen, Jonathan Wirsich, Brad P. Sutton, and Sepideh Sadaghiani.
457 Safety and data quality of eeg recorded simultaneously with multi-band fmri. *PLOS ONE*,
458 16(7):e0238485, 2021.
- 459 [33] Dominik Freche, Jodie Naim-Feil, Avi Peled, Nava Levit-Binnun, and Elisha Moses. A quan-
460 titative physical model of the tms-induced discharge artifacts in eeg. *PLOS Computational*
461 *Biology*, 14(7):e1006177, 2018.
- 462 [34] Johan N. van der Meer, Yke B. Eisma, Ronald Meester, Marc Jacobs, and Aart J. Nederveen.
463 Effects of mobile phone electromagnetic fields on brain waves in healthy volunteers. *Scientific*
464 *Reports*, 13(1):21758, 2023.
- 465 [35] Tuomas Mutanen, Hanna Mäki, and Risto J. Ilmoniemi. The effect of stimulus parameters on
466 tmseeg muscle artifacts. *Brain Stimulation*, 6(3):371–376, 2013.
- 467 [36] Limin Sun and Hermann Hinrichs. Simultaneously recorded eegfmri: Removal of gradient
468 artifacts by subtraction of head movement related average artifact waveforms. *Human Brain*
469 *Mapping*, 30(10):3361–3377, 2009.
- 470 [37] Sander Koelstra, Christian Muhl, Mohammad Soleymani, Jong-Seok Lee, Ashkan Yazdani,
471 Touradj Ebrahimi, Thierry Pun, Anton Nijholt, and Ioannis Patras. Deap: A database for emo-
472 tion analysis;using physiological signals. *IEEE Transactions on Affective Computing*, 3(1):18–
473 31, 2012.
- 474 [38] Neetha Das, Tom Francart, and Alexander Bertrand. Auditory attention detection dataset kuleu-
475 ven, 2020.
- 476 [39] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen,
477 Wenjun Zeng, and Philip S. Yu. Generalizing to unseen domains: A survey on domain gener-
478 alization. *IEEE Transactions on Knowledge and Data Engineering*, 35(8):8052–8072, 2023.
- 479 [40] Jiquan Wang, Sha Zhao, Haiteng Jiang, Shijian Li, Tao Li, and Gang Pan. Generalizable sleep
480 staging via multi-level domain alignment. *Proceedings of the AAAI Conference on Artificial*
481 *Intelligence*, 38(11):265–273, 2024.
- 482 [41] Chaoqi Yang, M. Brandon Westover, and Jimeng Sun. Manydg: Many-domain generalization
483 for healthcare applications. In *The Eleventh International Conference on Learning Representations*, 2023.
484
- 485 [42] Bernd Accou, Lies Bollens, Marlies Gillis, Wendy Verheijen, Hugo Van Hamme, and Tom
486 Francart. Sparrkulee: A speech-evoked auditory response repository of the ku leuven, contain-
487 ing eeg of 85 participants. *bioRxiv preprint bioRxiv: 2023.07.24.550310*, 2023.
- 488 [43] Zijiao Chen, Jiaxin Qing, Tiange Xiang, Wan Lin Yue, and Juan Helen Zhou. Seeing beyond
489 the brain: Conditional diffusion model with sparse masked modeling for vision decoding. In
490 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
491 22710–22720, 2023.
- 492 [44] Yunpeng Bai, Xintao Wang, Yan-pei Cao, Yixiao Ge, Chun Yuan, and Ying Shan. Dreamdiffu-
493 sion: Generating high-quality images from brain eeg signals. *arXiv preprint arXiv:2306.16934*,
494 2023.
- 495 [45] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
496 *arXiv:1711.05101*, 2017.
- 497 [46] Xinke Shen, Xianggen Liu, Xin Hu, Dan Zhang, and Sen Song. Contrastive learning of subject-
498 invariant eeg representations for cross-subject emotion recognition. *IEEE Transactions on*
499 *Affective Computing*, 14(3):2496–2511, 2023.

- 500 [47] Ivine Kuruvila, Jan Muncke, Eghart Fischer, and Ulrich Hoppe. Extracting the auditory atten-
501 tion in a dual-speaker scenario from eeg using a joint cnn-lstm model. *Frontiers in Physiology*,
502 12, 2021.
- 503 [48] Changde Du, Kaicheng Fu, Jinpeng Li, and Huiguang He. Decoding visual neural representa-
504 tions by multimodal learning of brain-visual-linguistic features. *IEEE Transactions on Pattern*
505 *Analysis and Machine Intelligence*, 45(9):10760–10777, 2023.
- 506 [49] Yonghao Song, Bingchuan Liu, Xiang Li, Nanlin Shi, Yijun Wang, and Xiaorong Gao. Decod-
507 ing natural images from eeg for object recognition. In *The Twelfth International Conference*
508 *on Learning Representations*, 2024.
- 509 [50] Zesheng Ye, Lina Yao, Yu Zhang, and Sylvia Gustin. Self-supervised cross-modal visual
510 retrieval from brain activities. *Pattern Recognition*, 145:109915, 2024.
- 511 [51] Yiming Wang, Bin Zhang, and Yujiao Tang. Dmmr: Cross-subject domain generalization for
512 eeg-based emotion recognition via denoising mixed mutual reconstruction. *Proceedings of the*
513 *AAAI Conference on Artificial Intelligence*, 38(11):628–636, 2024.
- 514 [52] Servaas Vandecappelle, Lucas Deckers, Neetha Das, Amir Hossein Ansari, Alexander
515 Bertrand, and Tom Francart. Eeg-based detection of the locus of auditory attention with con-
516 volutional neural networks. *eLife*, 10:e56481, 2021.
- 517 [53] Theo Gnassounou, Rémi Flamary, and Alexandre Gramfort. Convolution monge mapping
518 normalization for learning on sleep data. *Advances in Neural Information Processing Systems*,
519 36, 2023.
- 520 [54] Johanna Wilroth, Bo Bernhardsson, Frida Heskebeck, Martin A. Skoglund, Carolina Bergeling,
521 and Emina Alickovic. Improving eeg-based decoding of the locus of auditory attention through
522 domain adaptation*. *Journal of Neural Engineering*, 20(6):066022, 2023.
- 523 [55] Weibang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic
524 representations with tremendous eeg data in bci. In *The Twelfth International Conference on*
525 *Learning Representations*, 2024.
- 526 [56] Chaoqi Yang, M. Westover, and Jimeng Sun. Biot: Biosignal transformer for cross-data learn-
527 ing in the wild. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- 528 [57] Ke Yi, Yansen Wang, Kan Ren, and Dongsheng Li. Learning topology-agnostic eeg representa-
529 tions with geometry-aware modeling. In *Advances in Neural Information Processing Systems*,
530 volume 36, 2023.

531 **A Appendix A**

532 **A.1 Photography of the watermelon subject**



Figure 4: Photos of watermelons used in the experiment. Each watermelon’s ID is marked on the watermelon, with IDs ranging from 1 to 10.

533 **A.2 Reorganization for KUL dataset and DEAP dataset**

534 For the emotion recognition task, we referred to the experimental design of DEAP dataset [37]. In
535 this dataset, the EEG data were recorded while subjects are presented with 40 audio-visual clips of
536 60 seconds in length, with each corresponding to one of four emotion classes. We only used the first
537 32 channels of the EEG to match the EEG channel numbers in the DEAP dataset. The watermelon
538 EEG data and SparrKULee EEG data were down-sampled to 128 Hz and then were segmented into
539 40 60-second segments. The interval between adjacent segments is set to 40 seconds to match the
540 rest time of the subjects during the EEG recording in the KUL dataset. Each segment was assigned
541 a unique domain label and a class label in accordance with the DEAP dataset, and each segment was
542 further segmented into 2-second samples [25]. The reorganized datasets for the Watermelon EEG
543 Dataset and SparrKULee Dataset are called WM-DEAP and SK-DEAP, respectively.

544 For the ASAD task, we referred to the experimental design of the KUL dataset [38]. In this dataset,
545 8 clips of two-talker mixed speech are presented to subjects, with each lasting for 6 minutes. Each
546 speech clip contains a left talker and a right talker. Subjects are instructed to attend left or right talker
547 during the entire duration of one clip presentation. The watermelon EEG data and SparrKULee EEG
548 data were down-sampled to 128 Hz and then were epoch into 8 6-minute segments. The interval
549 between adjacent segments is set to 1-2 minutes to match the rest time of the subjects during the EEG
550 recording in the KUL dataset. Each segment was assigned a unique domain label and a class label
551 in accordance with the KUL dataset and was further segmented into 1-second samples [22, 21, 23].
552 The reorganized datasets for Watermelon Dataset and SparrKULee Dataset are called WM-KUL and
553 SK-KUL, respectively.

554 **A.3 Detailed implementation of joint training**

555 The joint training was performed on the WM-CVPR and SK-CVPR datasets. All EEG samples
556 were randomly divided into the training set, validation set, and test set in a ratio of 8:1:1. The image
557 encoder of the CLIP (CLIP VIT-L/14) model ¹ is chosen to extract image representation, yielding

¹<https://huggingface.co/openai/clip-vit-large-patch14>

558 768-dimensional vectors from the image inputs. The structure of the EEG encoder is similar to the
559 model introduced in Subsection 2.3, with an augmentation from 40 to 768 output nodes to match
560 the dimension of the image representation. The network is trained using either a cosine similarity
561 (CS) loss or an InfoNCE contrastive loss (with a temperature parameter set to 0.07). The evaluation
562 metrics selected are Top-1 accuracy, Top-5 accuracy, and Rank accuracy, where the Top-1 accuracy
563 metric is equivalent to the classification accuracy in the classification task.

564 **A.4 Detailed implementation of image generation**

565 We take an approach similar to previous works [44]². We used a CLIP image encoder to extract
566 image representation and trained an EEG encoder with cosine similarity loss to reconstruct image
567 representation from EEG. This process is the same as described in Joint training with image features.
568 The reconstructed features are then serviced as a conditional input of an image generator. To match
569 the reconstructed features, we employ the pre-trained StableDiffusion model³ as our generator. This
570 model uses a fixed pre-trained image encoder (CLIP ViT-L/14) to extract image features, which
571 then guide the Latent Diffusion models generation process in the latent space. The diffusion model
572 gradually generates images from a random noise distribution that corresponds to the conditional
573 features during its iterative process. To improve the generation performance, we fine-tuned the
574 generator with the reconstructed image features and the corresponding images. Experiments were
575 done on the WM-CVPR and SK-CVPR datasets. All EEG samples were randomly divided into
576 training set, validation set, and test set in a ratio of 8:1:1.

577 Consistent with previous work [1], we evaluate the semantic correctness of the generated images
578 using N-way Top-1 and Top-5 accuracy classification tasks. Specifically, given a generated image
579 input, a pre-trained ImageNet1K classifier is used to output a classification logit probability among
580 1000 classes. Among the 1000 classes, N-1 random classes and the correct class are selected, and
581 the Top-1 and Top-5 classification accuracy are calculated. To avoid randomness, this operation is
582 repeated 50 times for each generated image, with the average value taken as the accuracy.

583 **A.5 leave-subjects-out data splitting strategy**

584 In this subsection, we employed the leave-subjects-out data splitting strategy. This refers to using
585 data from a subset of subjects for training, while data from the remaining subjects are used for
586 testing. Within the training data, there are two further data partitioning methods: leave-samples-out
587 and leave-subjects-out. The former involves randomly dividing all samples of the training data into
588 training and validation sets, whereas the latter uses data from a subset of subjects for the training set,
589 with the remaining subjects data allocated for the test set. Table 5 presents the decoding accuracy
590 for six datasets (i.e., WM-CVPR, WM-DEAP, WM-KUL, SK-CVPR, SK-DEAP, and SK-KUL).

591 It can be observed that when the leave-samples-out splitting strategy was used within the training
592 data, both the training and validation sets achieved very high decoding accuracy, but the accuracy
593 only reached the chance level on the test set. Such results are similar to those reported by [46, 47, 8],
594 which corroborates the argument that while the leave-subjects-out approach may avert the domain
595 features leakage, it cannot prevent overfitting of the domain features during the training stage, as
596 discussed in Subsection 4.1. Moreover, when the leave-subjects-out data splitting strategy was used
597 within the training dataset, the validation set performance was only at chance level despite high
598 accuracy on the training set. This further demonstrates that decoding that relies on domain features
599 cannot be generalized to practical application scenarios.

²<https://github.com/bbaaii/DreamDiffusion>

³<https://huggingface.co/runwayml/stable-diffusion-v1-5>

Table 5: Decoding accuracy (%) for the six datasets on training, validation and test set. Leave-subjects-out data splitting strategy is used for training and test data. Leave-samples-out and leave-subjects-out data splitting strategy is used for training and validation set. The mean accuracy and standard deviation are calculated over subjects level with a five-fold cross-validation.

Data splitting strategy for validation set		WM-CVPR	WM-DEAP	WM-KUL	SK-CVPR	SK-DEAP	SK-KUL
leave-samples-out	Training	80.93 ± 1.68	87.86 ± 1.48	99.54 ± 0.16	69.17 ± 1.03	76.22 ± 0.71	100.00 ± 0.00
	validation	80.55 ± 1.59	86.10 ± 1.63	99.43 ± 0.24	68.86 ± 1.20	74.55 ± 0.60	100.00 ± 0.00
	Test	2.46 ± 0.16	24.22 ± 0.48	48.37 ± 2.15	2.70 ± 0.63	26.71 ± 0.87	50.22 ± 1.14
leave-subjects-out	Training	78.93 ± 1.09	86.40 ± 0.75	99.59 ± 0.16	72.31 ± 0.59	77.43 ± 0.52	100.00 ± 0.00
	validation	3.70 ± 0.34	22.23 ± 1.29	56.13 ± 3.06	4.15 ± 0.60	24.57 ± 0.33	53.24 ± 2.85
	Test	2.26 ± 0.16	24.90 ± 0.43	52.06 ± 1.26	2.13 ± 0.29	25.61 ± 0.43	45.22 ± 2.83
Chance level		2.50	25.00	50.00	2.50	25.00	50.00

600 A.6 Results on different frequency band

601 To demonstrate that domain features are not solely due to baseline drift, we conducted an analysis on
602 seven frequency bands across six datasets. These seven frequency bands are delta (0-4 Hz), theta (4-
603 8 Hz), alpha (8-12 Hz), beta (12-32 Hz), low gamma (32-45 Hz), and high gamma (55-95 Hz). High
604 gamma frequency band results for DEAP and KUL datasets are not presented due to the sampling
605 rate of 128 Hz (i.e., only frequency under 64 Hz is available according to the Nyquist sampling
606 theorem). Tables 6, 7, 8, and 9 show the decoding accuracy for domain label classification (DLC-
607 EEG), class label classification from domain features (TLC-DF), class label classification directly
608 from EEG (TLC-EEG), and class label classification directly from EEG when samples in the training
609 set and test set are from different domains (TLC-EEG-woDO), respectively. As expected, the highest
610 decoding accuracy is observed for both the low-frequency band (delta band) and the full-frequency
611 EEG data. However, other frequency bands also exhibited decoding accuracy significantly higher
612 than the chance level. This suggests that baseline correction through filtering does not eliminate
613 domain features. Consequently, any experimental designs and data partitioning strategies that could
614 lead to the leakage of domain information should be meticulously avoided.

Table 6: Decoding accuracy (%) using different EEG bands for domain label classification (DLC-EEG)

	WM-CVPR	WM-DEAP	WM-KUL	SK-CVPR	SK-DEAP	SK-KUL
Full	88.78 ± 4.95	96.98 ± 0.76	99.99 ± 0.01	69.83 ± 2.98	72.70 ± 1.36	100.00 ± 0.00
Delta	88.58 ± 5.11	96.31 ± 0.89	99.99 ± 0.01	69.65 ± 2.88	72.76 ± 1.24	100.00 ± 0.00
Theta	8.90 ± 1.95	10.54 ± 2.17	41.97 ± 5.50	11.24 ± 1.60	10.19 ± 1.15	43.11 ± 5.13
Alpha	8.62 ± 1.77	12.88 ± 2.80	43.42 ± 5.96	15.16 ± 1.76	12.87 ± 1.00	47.67 ± 4.70
Beta	18.53 ± 3.18	18.18 ± 2.72	57.85 ± 4.86	43.95 ± 2.27	43.68 ± 1.97	97.17 ± 0.71
Low gamma	39.74 ± 7.35	62.59 ± 5.95	85.97 ± 2.82	53.72 ± 2.25	52.82 ± 1.40	96.57 ± 0.96
High gamma	42.15 ± 7.39	-	-	61.55 ± 1.94	-	-
Chance level	2.50	2.50	50.00	2.50	2.50	50.00

Table 7: Decoding accuracy (%) using different EEG bands for class label classification from domain features (TLC-DF)

	WM-CVPR	WM-DEAP	WM-KUL	SK-CVPR	SK-DEAP	SK-KUL
Full	-	92.77 ± 1.31	100.00 ± 0.00	-	76.19 ± 1.80	100.00 ± 0.00
Delta	-	92.12 ± 1.49	100.00 ± 0.00	-	76.51 ± 1.74	100.00 ± 0.00
Theta	-	31.39 ± 1.80	67.78 ± 3.56	-	32.17 ± 1.16	69.41 ± 3.80
Alpha	-	33.10 ± 2.43	68.78 ± 4.03	-	33.88 ± 0.69	71.98 ± 3.47
Beta	-	39.03 ± 2.09	77.33 ± 3.71	-	56.91 ± 2.02	97.83 ± 0.72
Low gamma	-	59.32 ± 5.22	88.23 ± 2.56	-	63.80 ± 1.43	97.44 ± 0.88
High gamma	-	-	-	-	-	-
Chance level	-	25.00	50.00	-	25.00	50.00

Table 8: Decoding accuracy (%) using different EEG bands for class label classification directly from EEG (TLC-EEG)

	WM-CVPR	WM-DEAP	WM-KUL	SK-CVPR	SK-DEAP	SK-KUL
Full	88.78 ± 4.95	88.74 ± 3.26	82.74 ± 6.44	69.83 ± 2.98	74.44 ± 2.76	93.34 ± 2.01
Delta	88.58 ± 5.11	88.60 ± 3.36	81.49 ± 6.44	69.65 ± 2.88	74.90 ± 2.55	92.90 ± 2.15
Theta	8.90 ± 1.95	29.36 ± 1.27	66.40 ± 3.47	11.24 ± 1.60	30.62 ± 1.30	65.28 ± 3.83
Alpha	8.62 ± 1.77	31.00 ± 1.70	68.16 ± 3.59	15.16 ± 1.76	32.17 ± 1.10	67.11 ± 3.83
Beta	18.53 ± 3.18	35.95 ± 1.12	71.24 ± 4.16	43.95 ± 2.27	43.95 ± 1.78	93.27 ± 1.52
Low gamma	39.74 ± 7.35	52.05 ± 4.72	73.42 ± 5.37	53.72 ± 2.25	46.81 ± 1.03	93.51 ± 2.02
High gamma	42.15 ± 7.39	-	-	61.55 ± 1.94	-	-
Chance level	2.50	25.00	50.00	2.50	25.00	50.00

Table 9: Decoding accuracy (%) using different EEG bands for class label classification directly from EEG when samples in the training set and test set are from different domains (TLC-EEG-woDO)

	WM-CVPR	WM-DEAP	WM-KUL	SK-CVPR	SK-DEAP	SK-KUL
Full	-	24.67 ± 2.31	49.97 ± 4.67	-	25.34 ± 1.85	59.32 ± 4.07
Delta	-	25.89 ± 2.58	49.72 ± 4.85	-	24.71 ± 1.74	58.25 ± 3.76
Theta	-	23.91 ± 0.63	49.10 ± 3.13	-	23.28 ± 2.18	51.89 ± 4.32
Alpha	-	23.50 ± 0.82	49.70 ± 2.91	-	23.26 ± 1.68	52.77 ± 4.04
Beta	-	22.96 ± 1.25	50.30 ± 4.35	-	24.21 ± 1.39	57.32 ± 5.26
Low gamma	-	26.75 ± 2.17	49.46 ± 3.63	-	25.72 ± 1.61	54.88 ± 4.92
High gamma	-	-	-	-	-	-
Chance level	-	25.00	50.00	-	25.00	50.00

615 A.7 LRTC

616 The autocorrelation analysis was used to evaluate long range temporal correlation in EEG data from
 617 the Watermelon and SparrKULee datasets, similar to the approach taken by previous study. For a
 618 lengthy segment of single-channel EEG, the Morlet wavelet transform was employed to extract the
 619 time-varying amplitude envelope $W_f(t)$ at a given frequency f . The autocorrelation function ACF_f
 620 for $W_f(t)$ is defined as:

$$ACF_f(\tau) = \text{corr}(W_f(t), W_f(t + \tau)) \quad (4)$$

621 In the above equation, $\text{corr}(\cdot)$ denotes the Pearson correlation coefficient between two time series,
 622 and τ represents the time lag.

623 In our analysis, the original EEG data were down-sampled to 200 Hz. Ninety-five analysis frequen-
 624 cies were distributed linearly and evenly between 1-95 Hz. Two hundred autocorrelation time lags
 625 were logarithmically spaced between 0.5 s and 500 s. For each subject in the Watermelon dataset,
 626 continuous EEG recordings were divided into five segments of equal length (with each segment
 627 ranging from 15 to 20 minutes), and autocorrelation analysis was completed on each segment. For
 628 each subject in the SparrKULee dataset, the autocorrelation analysis was carried out separately on
 629 each of their ten trials. Figure 5 shows the results of the autocorrelation analysis for the Watermelon
 630 and SparrKULee datasets. The figure illustrates the magnitude of correlation at different frequen-
 631 cies and time lags (represented by color). The correlation values were obtained by averaging the
 632 results across all subjects, segments (trials), and electrodes. Black lines represent the contour lines
 633 where $p = 0.01$, as determined by statistical analysis. Statistical significance was assessed using
 634 single-sample t-test at the subject-electrode level. Specifically, for each electrode of each subject,
 635 the averaged Pearson correlation coefficient across all segments (trials) was used as the value for the t-
 636 test. Additionally, p -values were corrected for multiple comparisons using the Benjamini-Hochberg
 637 False Discovery Rate (BH-FDR) to type I error.

638 As demonstrated in Figure 5, EEG data from both Watermelon and SparrKULee datasets show
 639 significant LRTC across multiple frequency bands. For the EEG data from the Watermelon dataset,
 640 significant bands of LRTC are primarily distributed in the low-frequency range (<8 Hz) and around
 641 50 Hz, with these correlations spanning over 500 seconds. This indicates that baseline drifts and line

642 noise contribute to the temporal correlation observed in the Watermelon dataset. For the EEG data
643 from the SparrKULee dataset, LTRCs are significant across the entire frequency range. Similarly,
644 LTRCs are most prominent at low frequencies (<5 Hz) and around 50 Hz, consistent with the findings
645 from the Watermelon dataset. Notably, for SparrKULee dataset, there is also a significant presence
646 of LTRC around 10 Hz, which aligns with previous research findings [13], suggesting the temporal
647 correlation of alpha oscillations in human subjects.

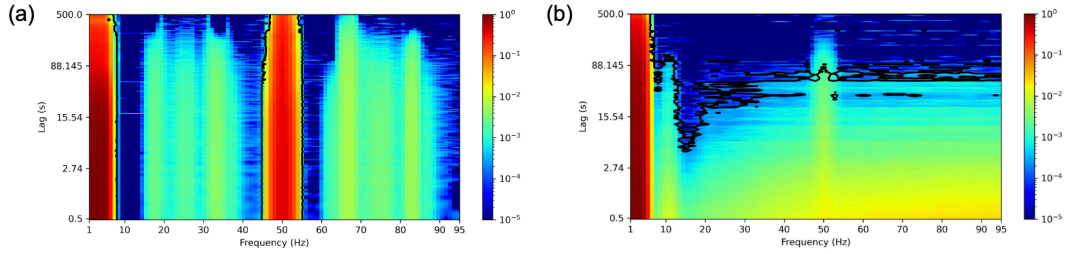


Figure 5: Autocorrelation analysis result on (a) Watermelon and (b) SparrKULee datasets.

648 **NeurIPS Paper Checklist**

649 **1. Claims**

650 Question: Do the main claims made in the abstract and introduction accurately reflect the
651 paper's contributions and scope?

652 Answer: [Yes]

653 Justification: the problem formulation could be found in 2.1 and results supporting the
654 contribution of this paper could be found in 3.1 and section 3.2.

655 Guidelines:

- 656 • The answer NA means that the abstract and introduction do not include the claims
657 made in the paper.
- 658 • The abstract and/or introduction should clearly state the claims made, including the
659 contributions made in the paper and important assumptions and limitations. A No or
660 NA answer to this question will not be perceived well by the reviewers.
- 661 • The claims made should match theoretical and experimental results, and reflect how
662 much the results can be expected to generalize to other settings.
- 663 • It is fine to include aspirational goals as motivation as long as it is clear that these
664 goals are not attained by the paper.

665 **2. Limitations**

666 Question: Does the paper discuss the limitations of the work performed by the authors?

667 Answer: [Yes]

668 Justification: the limitations could be found in section 4.3.

669 Guidelines:

- 670 • The answer NA means that the paper has no limitation while the answer No means
671 that the paper has limitations, but those are not discussed in the paper.
- 672 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 673 • The paper should point out any strong assumptions and how robust the results are to
674 violations of these assumptions (e.g., independence assumptions, noiseless settings,
675 model well-specification, asymptotic approximations only holding locally). The au-
676 thors should reflect on how these assumptions might be violated in practice and what
677 the implications would be.
- 678 • The authors should reflect on the scope of the claims made, e.g., if the approach was
679 only tested on a few datasets or with a few runs. In general, empirical results often
680 depend on implicit assumptions, which should be articulated.
- 681 • The authors should reflect on the factors that influence the performance of the ap-
682 proach. For example, a facial recognition algorithm may perform poorly when image
683 resolution is low or images are taken in low lighting. Or a speech-to-text system might
684 not be used reliably to provide closed captions for online lectures because it fails to
685 handle technical jargon.
- 686 • The authors should discuss the computational efficiency of the proposed algorithms
687 and how they scale with dataset size.
- 688 • If applicable, the authors should discuss possible limitations of their approach to ad-
689 dress problems of privacy and fairness.
- 690 • While the authors might fear that complete honesty about limitations might be used by
691 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
692 limitations that aren't acknowledged in the paper. The authors should use their best
693 judgment and recognize that individual actions in favor of transparency play an impor-
694 tant role in developing norms that preserve the integrity of the community. Reviewers
695 will be specifically instructed to not penalize honesty concerning limitations.

696 **3. Theory Assumptions and Proofs**

697 Question: For each theoretical result, does the paper provide the full set of assumptions and
698 a complete (and correct) proof?

699 Answer: [NA]

700 Justification: The paper does not include any theoretical results.

701 Guidelines:

- 702 • The answer NA means that the paper does not include theoretical results.
- 703 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 704 referenced.
- 705 • All assumptions should be clearly stated or referenced in the statement of any theo-
- 706 rems.
- 707 • The proofs can either appear in the main paper or the supplemental material, but if
- 708 they appear in the supplemental material, the authors are encouraged to provide a
- 709 short proof sketch to provide intuition.
- 710 • Inversely, any informal proof provided in the core of the paper should be comple-
- 711 mented by formal proofs provided in appendix or supplemental material.
- 712 • Theorems and Lemmas that the proof relies upon should be properly referenced.

713 4. Experimental Result Reproducibility

714 Question: Does the paper fully disclose all the information needed to reproduce the main
715 experimental results of the paper to the extent that it affects the main claims and/or conclu-
716 sions of the paper (regardless of whether the code and data are provided or not)?

717 Answer: [\[Yes\]](#)

718 Justification: the information needed to reproduce all the experimental results of this paper
719 could be found in Section 2.2, Section 2.3, Section 2.4, Section 2.5 and Appendix A.

720 Guidelines:

- 721 • The answer NA means that the paper does not include experiments.
- 722 • If the paper includes experiments, a No answer to this question will not be perceived
- 723 well by the reviewers: Making the paper reproducible is important, regardless of
- 724 whether the code and data are provided or not.
- 725 • If the contribution is a dataset and/or model, the authors should describe the steps
- 726 taken to make their results reproducible or verifiable.
- 727 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 728 For example, if the contribution is a novel architecture, describing the architecture
- 729 fully might suffice, or if the contribution is a specific model and empirical evaluation,
- 730 it may be necessary to either make it possible for others to replicate the model with
- 731 the same dataset, or provide access to the model. In general, releasing code and data
- 732 is often one good way to accomplish this, but reproducibility can also be provided via
- 733 detailed instructions for how to replicate the results, access to a hosted model (e.g., in
- 734 the case of a large language model), releasing of a model checkpoint, or other means
- 735 that are appropriate to the research performed.
- 736 • While NeurIPS does not require releasing code, the conference does require all sub-
- 737 missions to provide some reasonable avenue for reproducibility, which may depend
- 738 on the nature of the contribution. For example
- 739 (a) If the contribution is primarily a new algorithm, the paper should make it clear
- 740 how to reproduce that algorithm.
- 741 (b) If the contribution is primarily a new model architecture, the paper should describe
- 742 the architecture clearly and fully.
- 743 (c) If the contribution is a new model (e.g., a large language model), then there should
- 744 either be a way to access this model for reproducing the results or a way to re-
- 745 produce the model (e.g., with an open-source dataset or instructions for how to
- 746 construct the dataset).
- 747 (d) We recognize that reproducibility may be tricky in some cases, in which case au-
- 748 thors are welcome to describe the particular way they provide for reproducibility.
- 749 In the case of closed-source models, it may be that access to the model is limited in
- 750 some way (e.g., to registered users), but it should be possible for other researchers
- 751 to have some path to reproducing or verifying the results.

752 5. Open access to data and code

753 Question: Does the paper provide open access to the data and code, with sufficient instruc-
754 tions to faithfully reproduce the main experimental results, as described in supplemental
755 material?

756 Answer: [Yes]

757 Justification: the collected Watermelon EEG dataset could be available in Zenodo and the
758 human dataset used in this work could also be downloaded in the link provided in supple-
759 mentary materials. All the codes to reproduce this work can be found in supplementary
760 materials.

761 Guidelines:

- 762 • The answer NA means that paper does not include experiments requiring code.
- 763 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
764 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 765 • While we encourage the release of code and data, we understand that this might not
766 be possible, so No is an acceptable answer. Papers cannot be rejected simply for not
767 including code, unless this is central to the contribution (e.g., for a new open-source
768 benchmark).
- 769 • The instructions should contain the exact command and environment needed to run to
770 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.
771 cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 772 • The authors should provide instructions on data access and preparation, including how
773 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 774 • The authors should provide scripts to reproduce all experimental results for the new
775 proposed method and baselines. If only a subset of experiments are reproducible, they
776 should state which ones are omitted from the script and why.
- 777 • At submission time, to preserve anonymity, the authors should release anonymized
778 versions (if applicable).
- 779 • Providing as much information as possible in supplemental material (appended to the
780 paper) is recommended, but including URLs to data and code is permitted.

781 6. Experimental Setting/Details

782 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
783 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
784 results?

785 Answer: [Yes]

786 Justification: the experimental setting and details could also be found in Section 2.2, Sec-
787 tion 2.3, Section 2.4, Section 2.5, and could also be found in the codes.

788 Guidelines:

- 789 • The answer NA means that the paper does not include experiments.
- 790 • The experimental setting should be presented in the core of the paper to a level of
791 detail that is necessary to appreciate the results and make sense of them.
- 792 • The full details can be provided either with the code, in appendix, or as supplemental
793 material.

794 7. Experiment Statistical Significance

795 Question: Does the paper report error bars suitably and correctly defined or other appropri-
796 ate information about the statistical significance of the experiments?

797 Answer: [Yes]

798 Justification: the standard error of the mean is reported for all results. As we only com-
799 pared the result against chance level, one-sample t-test was used for statistical analysis to
800 check whether the reported results are significant high above the chance level. Given that
801 decoding analyses was conducted on multiple frequency bands, Bonferroni correction was
802 used to adjust the p-value to reduce the risk of type-I error. A p-value of 0.05 or lower is
803 considered statistically significant.

804 Guidelines:

- 805 • The answer NA means that the paper does not include experiments.
- 806 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 807 dence intervals, or statistical significance tests, at least for the experiments that support
- 808 the main claims of the paper.
- 809 • The factors of variability that the error bars are capturing should be clearly stated (for
- 810 example, train/test split, initialization, random drawing of some parameter, or overall
- 811 run with given experimental conditions).
- 812 • The method for calculating the error bars should be explained (closed form formula,
- 813 call to a library function, bootstrap, etc.)
- 814 • The assumptions made should be given (e.g., Normally distributed errors).
- 815 • It should be clear whether the error bar is the standard deviation or the standard error
- 816 of the mean.
- 817 • It is OK to report 1-sigma error bars, but one should state it. The authors should prefer-
- 818 ably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of
- 819 Normality of errors is not verified.
- 820 • For asymmetric distributions, the authors should be careful not to show in tables or
- 821 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 822 error rates).
- 823 • If error bars are reported in tables or plots, The authors should explain in the text how
- 824 they were calculated and reference the corresponding figures or tables in the text.

825 8. Experiments Compute Resources

826 Question: For each experiment, does the paper provide sufficient information on the com-
 827 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 828 the experiments?

829 Answer: [Yes]

830 Justification: the neural networks were implemented with the Pytorch and trained on a
 831 single HPC node with 8 A800 GPU.

832 Guidelines:

- 833 • The answer NA means that the paper does not include experiments.
- 834 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 835 or cloud provider, including relevant memory and storage.
- 836 • The paper should provide the amount of compute required for each of the individual
- 837 experimental runs as well as estimate the total compute.
- 838 • The paper should disclose whether the full research project required more compute
- 839 than the experiments reported in the paper (e.g., preliminary or failed experiments
- 840 that didn't make it into the paper).

841 9. Code Of Ethics

842 Question: Does the research conducted in the paper conform, in every respect, with the
 843 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

844 Answer: [Yes]

845 Justification: The collected dataset was released in an anonymous form, and all codes do
 846 not contain any identity information.

847 Guidelines:

- 848 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 849 • If the authors answer No, they should explain the special circumstances that require a
- 850 deviation from the Code of Ethics.
- 851 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 852 eration due to laws or regulations in their jurisdiction).

853 10. Broader Impacts

854 Question: Does the paper discuss both potential positive societal impacts and negative
 855 societal impacts of the work performed?

856 Answer: [NA]

857 Justification: The purpose of this paper is to let researchers beware of overestimated de-
858 coding performance arising from temporal autocorrelations in EEG signals. This work
859 formalizes and proves the pitfalls existing in current EEG decoding tasks. We believe that
860 this will not generate any significant societal impact.

861 Guidelines:

- 862 • The answer NA means that there is no societal impact of the work performed.
- 863 • If the authors answer NA or No, they should explain why their work has no societal
864 impact or why the paper does not address societal impact.
- 865 • Examples of negative societal impacts include potential malicious or unintended uses
866 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
867 (e.g., deployment of technologies that could make decisions that unfairly impact spe-
868 cific groups), privacy considerations, and security considerations.
- 869 • The conference expects that many papers will be foundational research and not tied
870 to particular applications, let alone deployments. However, if there is a direct path to
871 any negative applications, the authors should point it out. For example, it is legitimate
872 to point out that an improvement in the quality of generative models could be used to
873 generate deepfakes for disinformation. On the other hand, it is not needed to point out
874 that a generic algorithm for optimizing neural networks could enable people to train
875 models that generate Deepfakes faster.
- 876 • The authors should consider possible harms that could arise when the technology is
877 being used as intended and functioning correctly, harms that could arise when the
878 technology is being used as intended but gives incorrect results, and harms following
879 from (intentional or unintentional) misuse of the technology.
- 880 • If there are negative societal impacts, the authors could also discuss possible mitiga-
881 tion strategies (e.g., gated release of models, providing defenses in addition to attacks,
882 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
883 feedback over time, improving the efficiency and accessibility of ML).

884 11. Safeguards

885 Question: Does the paper describe safeguards that have been put in place for responsible
886 release of data or models that have a high risk for misuse (e.g., pretrained language models,
887 image generators, or scraped datasets)?

888 Answer: [NA]

889 Justification: the paper poses no such risks.

890 Guidelines:

- 891 • The answer NA means that the paper poses no such risks.
- 892 • Released models that have a high risk for misuse or dual-use should be released with
893 necessary safeguards to allow for controlled use of the model, for example by re-
894 quiring that users adhere to usage guidelines or restrictions to access the model or
895 implementing safety filters.
- 896 • Datasets that have been scraped from the Internet could pose safety risks. The authors
897 should describe how they avoided releasing unsafe images.
- 898 • We recognize that providing effective safeguards is challenging, and many papers do
899 not require this, but we encourage authors to take this into account and make a best
900 faith effort.

901 12. Licenses for existing assets

902 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
903 the paper, properly credited and are the license and terms of use explicitly mentioned and
904 properly respected?

905 Answer: [Yes]

906 Justification: The existing assets we used are the SparrKULee dataset, which is licensed
907 under CC-BY-NC-4.0.

908 Guidelines:

- 909 • The answer NA means that the paper does not use existing assets.

- 910 • The authors should cite the original paper that produced the code package or dataset.
- 911 • The authors should state which version of the asset is used and, if possible, include a
- 912 URL.
- 913 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 914 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 915 service of that source should be provided.
- 916 • If assets are released, the license, copyright information, and terms of use in the pack-
- 917 age should be provided. For popular datasets, paperswithcode.com/datasets has
- 918 curated licenses for some datasets. Their licensing guide can help determine the li-
- 919 cense of a dataset.
- 920 • For existing datasets that are re-packaged, both the original license and the license of
- 921 the derived asset (if it has changed) should be provided.
- 922 • If this information is not available online, the authors are encouraged to reach out to
- 923 the asset's creators.

924 13. New Assets

925 Question: Are new assets introduced in the paper well documented and is the documenta-

926 tion provided alongside the assets?

927 Answer: [Yes]

928 Justification: We have released an anonymous Watermelon EEG dataset, which can be

929 accessed at <https://zenodo.org/records/11238929>

930 Guidelines:

- 931 • The answer NA means that the paper does not release new assets.
- 932 • Researchers should communicate the details of the dataset/code/model as part of their
- 933 submissions via structured templates. This includes details about training, license,
- 934 limitations, etc.
- 935 • The paper should discuss whether and how consent was obtained from people whose
- 936 asset is used.
- 937 • At submission time, remember to anonymize your assets (if applicable). You can
- 938 either create an anonymized URL or include an anonymized zip file.

939 14. Crowdsourcing and Research with Human Subjects

940 Question: For crowdsourcing experiments and research with human subjects, does the pa-

941 per include the full text of instructions given to participants and screenshots, if applicable,

942 as well as details about compensation (if any)?

943 Answer: [NA]

944 Justification: we collected the Watermelon EEG dataset, but watermelon is not a human

945 subject. Nonetheless, we provided experiment details, which could be found in section

946 Section 2.2.

947 Guidelines:

- 948 • The answer NA means that the paper does not involve crowdsourcing nor research
- 949 with human subjects.
- 950 • Including this information in the supplemental material is fine, but if the main contri-
- 951 bution of the paper involves human subjects, then as much detail as possible should
- 952 be included in the main paper.
- 953 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-
- 954 tion, or other labor should be paid at least the minimum wage in the country of the
- 955 data collector.

956 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human

957 Subjects

958 Question: Does the paper describe potential risks incurred by study participants, whether

959 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

960 approvals (or an equivalent approval/review based on the requirements of your country or

961 institution) were obtained?

962
963
964
965
966
967
968
969
970
971
972
973
974

Answer: [NA]

Justification: the paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.