

NEUROSKETCH: AN EFFECTIVE FRAMEWORK FOR NEURAL DECODING VIA SYSTEMATIC ARCHITECTURAL OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural decoding, a critical component of Brain-Computer Interface (BCI), has recently attracted increasing research interest. Previous research has focused on leveraging signal processing and deep learning methods to enhance neural decoding performance. However, the in-depth exploration of model architectures remains underexplored, despite its proven effectiveness in other tasks such as energy forecasting and image classification. In this study, we propose NeuroSketch, an effective framework for neural decoding via systematic architecture optimization. Starting with the basic architecture study, we find that CNN-2D outperforms other architectures in neural decoding tasks and explore its effectiveness from temporal and spatial perspectives. Building on this, we optimize the architecture from macro- to micro-level, achieving improvements in performance at each step. The exploration process and model validations take over 5,000 experiments spanning three distinct modalities (visual, auditory, and speech), three types of brain signals (EEG, SEEG, and ECoG), and eight diverse decoding tasks. Experimental results indicate that NeuroSketch achieves state-of-the-art (SOTA) performance across all evaluated datasets, positioning it as a powerful tool for neural decoding.

1 INTRODUCTION

Brain-Computer Interface (BCI) technology aims to revolutionize interaction methods by establishing a direct link between thought and action, enabling more efficient communication (Maiseli et al., 2023). Neural decoding (Van Gerven et al., 2019) plays a critical role in this process, as it involves inferring external stimuli, cognitive states, or intentions from brain signals. These signals are typically recorded using methods such as electroencephalography (EEG) (Teplan, 2002), stereoelectroencephalography (SEEG) (Talairach, 1974), and electrocorticography (ECoG) (Shenoy et al., 2007), which can be classified as non-invasive and invasive techniques. EEG, a non-invasive method, records electrical activity along the scalp and is widely used in both research (Zhang et al., 2021; Altaheri et al., 2023; Rahman et al., 2021) and clinical settings (Pani et al., 2022; Saminu et al., 2023) due to its practicality and low cost. In contrast, invasive methods like SEEG and ECoG capture signals from deeper brain structures, offering higher temporal and spatial resolution compared to non-invasive methods (Chaddad et al., 2023), and enabling more precise identification of neural activity linked to specific cognitive functions.

In neural decoding, recorded brain signals are characterized by transient temporal dynamics (King and Dehaene, 2014) and spatial locality (Mahjoory et al., 2024), representing unique modeling demands. Traditional time series data, such as weather, electricity consumption, and traffic flow, are sampled at a relatively low frequency (*e.g.*, minutes or hours) and often exhibit significant periodicity or trends. However, brain signals require higher temporal resolution to capture the brain’s rapidly changing states. Even within brain signals, their characteristics can differ significantly across various scenarios. From a temporal perspective, brain signals typically contain crucial information within much shorter time frames in neural decoding tasks compared to other scenarios, such as sleep staging. In sleep staging tasks, brain signals are recorded overnight, with a labeling resolution of 30 seconds for each sleep stage (Phan and Mikkelsen, 2022). In contrast, in the Rapid Serial Visual Presentation (RSVP) paradigm for image stimulus decoding, relevant information is encoded within just a few hundred milliseconds (Butts et al., 2007). From a spatial perspective, task-driven neural

054 activity in decoding paradigms differ fundamentally from those of pathological neural disorders.
 055 Pathological conditions such as generalized epilepsy exhibit widespread brain activity, resulting in
 056 non-specific neural activations (Stafstrom and Carmant, 2015). In neural decoding tasks, however,
 057 external stimuli evoke localized activations that are typically restricted to functionally specialized
 058 regions and recorded by spatially proximate electrodes. For example, speech stimuli predominantly
 059 activate the left inferior frontal gyrus (Leonard et al., 2024), while visual stimuli engage the visual
 060 cortex (Grill-Spector and Malach, 2004).

061 Previous work on neural decoding can be broadly categorized into two approaches. The first is signal
 062 processing (Peksa and Mamchur, 2023), which focuses on improving the signal-to-noise ratio and
 063 extracting task-relevant features (Wittevrongel et al., 2020; Proix et al., 2022; Safi and Safi, 2021).
 064 These methods rely on manually defined features, which are time-consuming and highly empirical.
 065 The second approach involves deep learning methods (Angrick et al., 2019; Makin et al., 2020;
 066 Wilson et al., 2020; Willett et al., 2023; Metzger et al., 2023; Zheng et al., 2024). However, current
 067 approaches largely overlook in-depth exploration and modeling of the temporal dynamics and spatial
 068 locality inherent to neural decoding. Therefore, our work aims to explore a framework more suited to
 069 neural decoding tasks through systematic architectural optimization. Such optimizations have already
 070 yielded significant benefits in other tasks (Yu et al., 2022; 2023; Luo and Wang, 2024; Wang et al.,
 071 2025a). For instance, in energy forecasting tasks, ModernTCN (Luo and Wang, 2024) reduces MSE
 072 by 13.9% on ETTm2 (Zhou et al., 2021) compared to convolution-based models by using large kernel
 073 convolutions and a ConvFFN module. In image classification tasks, ConvNeXt (Liu et al., 2022)
 074 surpasses ResNet-50 (He et al., 2016) on ImageNet (Deng et al., 2009) through optimizations in stage
 075 compute ratio and convolutional methods.

076 Our study is guided by two core questions that aim to systematically investigate how to design an
 077 effective framework for neural decoding:

078 *Q1: Which basic model architecture is best suited to capture temporal and spatial patterns in brain*
 079 *signals?*

080 *Q2: Based on an appropriate architecture, how can we improve neural decoding performance through*
 081 *macro-to-micro architectural optimization?*

082 As shown in Figure 1, by an-
 083 swering the above questions, we
 084 consistently improve the decoding
 085 performance at each opti-
 086 mization step. First, we study
 087 9 basic architectures (CNN (He
 088 et al., 2016), GRU (Chung et al.,
 089 2014), Transformer (Vaswani
 090 et al., 2017), and their variants)
 091 and find that CNN-2D outper-
 092 forms the others. The results
 093 are coherent with the characteris-
 094 tics of neural decoding signals,
 095 since CNN-2D is suitable for
 096 capturing localized patterns from
 097 both temporal and spatial dimen-
 098 sions. Based on CNN-2D, we
 099 design the entire framework us-
 100 ing a macro-to-micro paradigm,
 101 analyzing the latent space transfor-
 102 mation during forward propagation
 103 and optimizing the calculation of
 104 micro components. Through the
 105 above step-by-step optimization
 106 process, we propose NeuroSketch,
 107 an effective framework for neural
 decoding tasks, which can serve as
 a useful tool for future research
 and applications in neuroscience.

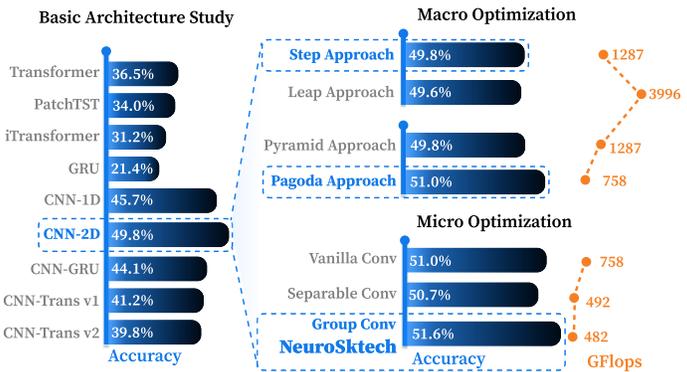


Figure 1: Roadmap of architectural optimization.

105 In summary, the main contributions of this research are as follows:

- 106 1. We propose NeuroSketch, an effective framework that aligns with the temporal and spatial charac-
 107 teristics of brain signals for neural decoding tasks.

2. From the technical perspective, we systematically reexamine the design space and consider several critical components. Starting with exploring basic architectures, we progressively investigate architectural design principles from macro- to micro-level.
3. From the experimental perspective, we validate NeuroSketch through over 5,000 experiments across eight neural decoding tasks related to vision, hearing, and speech, recorded by different types of brain signals. The experimental results demonstrate that NeuroSketch achieves state-of-the-art decoding performance across all evaluated datasets.

2 ROADMAP TOWARDS AN EFFECTIVE FRAMEWORK

In this section, we provide a trajectory that outlines the steps towards achieving an effective framework. The roadmap of our study is essentially focused on answering two key questions mentioned in Section 1. To begin with, Section 2.1 describes the experimental setup. Then, Section 2.2 addresses Question 1 through systematic analyses of the basic model architecture, considering both temporal and spatial perspectives. Building upon the appropriate architecture, Section 2.3 investigates the macro-level optimization aspect of Question 2, focusing on the latent space transformation. Subsequently, Section 2.4 explores the micro-level optimization aspect of Question 2, focusing on computation optimization. Finally, Section 2.5 integrates the above findings and introduces NeuroSketch.

2.1 EXPERIMENT SETUP

To rigorously evaluate our exploration, we conduct nearly 2,000 experiments on eight neural decoding tasks, spanning three major categories—speech, visual, and auditory decoding—and three types of brain signals—EEG, SEEG, and ECoG. The overview of each dataset is provided in Table 1, with further details in Appendix C. Since the network complexity closely correlates with final performance, the model’s parameter size is kept around 30M by adjusting the model depth and embedding dimension during the experimental process. Detailed parameter settings and training setup can be found in Appendix D.1 and Appendix E.

Table 1: **Overview of the datasets used in our experiments.** We evaluate on six datasets in total. For *Chisco* and *OpenMIIR*, we partition them by task, yielding eight distinct tasks overall.

Categories	Datasets	Tasks	Signal Type
Speech	Du-IN (Zheng et al., 2024)	Loud Reading	SEEG
	Chisco-R (Zhang et al., 2024)	Silent Reading	EEG
	Chisco-I (Zhang et al., 2024)	Speech Imagination	EEG
Visual	FacesHouses (Miller, 2019)	Binary Image Decoding	ECoG
	ThingsEEG (Gifford et al., 2022)	Multi-Class Image Decoding	EEG
	SEED-DV (Liu et al., 2024a)	Video Decoding	EEG
Auditory	OpenMIIR-P (Stober et al., 2015)	Music Perception	EEG
	OpenMIIR-I (Stober et al., 2015)	Music Imagination	EEG

2.2 BASIC ARCHITECTURE STUDY

To address Question 1, we investigate the performance of nine commonly used architectures in neural decoding tasks, selected from four categories: Convolutional Neural Networks (CNNs), Gated Recurrent Units (GRUs), Transformers, and their hybrids. The input of these models can be represented as $\mathbf{X} \in \mathbb{R}^{B \times C \times L}$, where B denotes the batch size, C denotes the number of channels, and L denotes the number of time steps. We study two types of CNNs: CNN-1D and CNN-2D. CNN-1D applies 1D convolutional filters across the temporal dimension, with each filter combining features from all channels. CNN-2D reshapes the input \mathbf{X} to $\mathbf{X}' \in \mathbb{R}^{B \times 1 \times C \times L}$ and applies 2D convolutional filters across both the temporal and channel dimensions. In the case of Vanilla GRU, the model sequentially processes the input \mathbf{X} across the temporal dimension, updating its hidden state at each time step using both the current input and the previous state. In contrast, the Transformer processes the input \mathbf{X} in parallel using self-attention to capture dependencies across all time steps. We also consider two Transformer variants: PatchTST(Nie et al., 2023) and iTransformer(Liu et al.,

2024b). PatchTST divides the input \mathbf{X} along the temporal dimension into patches of length P and stride S , generating a sequence of N patches, where $N = \lfloor \frac{L-P}{S} \rfloor + 2$. Then, it processes the patched input $\mathbf{X}_p \in \mathbb{R}^{B \times C \times N \times P}$ using Transformer layers to capture the relationships between different patches while maintaining channel independence. In contrast, iTransformer treats channels as sequence elements and uses self-attention to capture interactions across channels. In addition to individual models, we explore hybrid architectures that combine multiple approaches, including CNN-GRU and two CNN-Transformer variants. One variant uses CNN for feature extraction followed by Transformer layers (Song et al., 2022), while the other integrates both CNN and Transformer within each module to jointly capture temporal information at different levels (Kim et al., 2022). CNN-GRU, similar to the first CNN-Transformer variant, employs CNN for initial feature extraction, with the GRU component then processing the sequential data to capture temporal dependencies.

Table 2: **Results of the basic architecture study and subsequent analyses.** v1 and v2 refer to the two hybrid methods discussed above. The ratio shown in the middle of the table (e.g., 4:1) indicates the proportion of CNN to Transformer layers in each hybrid model. The best results are in **bold** and the second best are underlined.

Models	Datasets	Du-IN		SEED-DV		ThingsEEG		FacesHouses		OpenMIIR-P		OpenMIIR-I		Chisco-R		Chisco-I	
		Acc	F1														
<i>Basic Architecture Study</i>																	
CNN-1D		<u>.451</u>	<u>.446</u>	.061	.054	.202	.191	.799	.796	<u>.973</u>	<u>.973</u>	<u>.968</u>	<u>.968</u>	.104	.071	.100	.073
CNN-2D		.647	.641	<u>.063</u>	<u>.058</u>	<u>.184</u>	<u>.177</u>	.886	.885	.981	.980	.982	.982	.127	.113	<u>.114</u>	.095
GRU		.353	.343	.025	.001	.040	.037	.774	.773	.146	.046	.136	.039	<u>.123</u>	<u>.091</u>	.117	<u>.086</u>
Transformer		.085	.080	.042	.028	.005	.000	.795	.794	.953	.952	.911	.910	.069	.015	.057	.012
PatchTST		.060	.051	.038	.029	.006	.001	<u>.845</u>	<u>.838</u>	.906	.904	.766	.760	.050	.005	.050	.005
iTransformer		.097	.085	.029	.005	.005	.000	.730	.728	.729	.698	.780	.774	.064	.013	.061	.014
CNN-GRU		.416	.409	.057	.049	.172	.166	.803	.801	.957	.954	.943	.942	.092	.064	.087	.063
CNN-Trans v1		.333	.326	.067	.060	.046	.027	.752	.751	.948	.947	.945	.945	.106	.061	.100	.057
CNN-Trans v2		.282	.268	.050	.045	.022	.010	.754	.752	.952	.950	.930	.930	.097	.050	.093	.059
<i>Ratio of CNN and Transformer Layers in Hybrid Architectures</i>																	
CNN-Trans v1-4:1		.333	.326	<u>.067</u>	.060	.046	.027	.752	.751	.948	.947	.945	.945	.106	.061	.100	.057
CNN-Trans v1-3:2		<u>.115</u>	<u>.091</u>	.069	<u>.057</u>	<u>.010</u>	<u>.001</u>	.696	.694	<u>.904</u>	<u>.901</u>	<u>.916</u>	<u>.914</u>	<u>.094</u>	<u>.036</u>	<u>.090</u>	<u>.041</u>
CNN-Trans v1-2:3		.060	.032	.041	.021	.005	.000	<u>.699</u>	<u>.697</u>	.833	.828	.867	.866	.088	.030	.085	.031
CNN-Trans v1-1:4		<u>.036</u>	.013	.030	.008	.005	.000	.684	.674	.709	.687	.764	.760	.079	.018	.076	.019
CNN-Trans v2-4:1		.282	.268	.050	.045	<u>.022</u>	.010	.754	.752	.952	.950	.930	.930	.097	<u>.050</u>	.093	.059
CNN-Trans v2-3:2		<u>.211</u>	<u>.192</u>	.061	.054	.024	<u>.009</u>	.719	.713	<u>.946</u>	<u>.944</u>	.882	.880	.097	.057	.087	<u>.051</u>
CNN-Trans v2-2:3		.141	.124	<u>.059</u>	<u>.046</u>	.011	.002	.710	.707	.882	.881	<u>.893</u>	<u>.892</u>	.088	.050	<u>.092</u>	.044
CNN-Trans v2-1:4		.047	.033	.035	.019	.008	.001	<u>.722</u>	<u>.720</u>	.822	.814	.839	.837	<u>.094</u>	.047	.084	.041
<i>New Patch Method for Transformers and GRUs</i>																	
Transformer		.085	.080	.042	<u>.028</u>	<u>.005</u>	<u>.000</u>	.795	.794	.953	.952	.911	.910	<u>.069</u>	<u>.015</u>	.057	.012
PatchTST		.060	.051	<u>.038</u>	.029	.006	.001	.845	.838	.906	.904	.766	.760	.050	.005	.050	.005
iTransformer		<u>.097</u>	<u>.085</u>	.029	.005	.005	.000	.730	.728	.729	.698	.780	.774	.064	.013	.061	.014
Ours		.234	.226	.033	.027	.005	.000	<u>.799</u>	<u>.799</u>	<u>.922</u>	<u>.919</u>	<u>.843</u>	<u>.839</u>	.098	.046	.090	.066
GRU		.353	.343	.025	.001	.040	.037	.774	.773	.146	.046	.136	.039	.123	.091	.117	.086
Ours		.351	.342	.033	.026	.074	.068	.817	.813	.958	.956	.937	.936	.105	.071	.091	.066

Table 2(upper) shows the performance of the 9 basic architectures discussed above. Overall, CNN-based models, which excel in capturing local patterns, outperform hybrid models on most neural decoding tasks. Meanwhile, GRU and Transformer-based models, which are better suited for modeling long-range dependencies, exhibit the lowest performance. This result aligns with the **transient temporal dynamics** of brain signals in neural decoding tasks, suggesting that short-range temporal information is more effective for neural decoding. In CNN-based models, CNN-2D, which extracts local channel information using convolutional kernels, outperforms CNN-1D, which fuses all channel information through addition. This performance difference aligns with the **spatial locality** observed in neural decoding tasks. Surprisingly, the Transformer-based models perform poorly on certain datasets, particularly the ThingsEEG dataset, where their best accuracy of 0.6% is only marginally above the chance level of 0.5%. In order to find out the underlying reasons for the performance differences, we conduct further analyses from both temporal and spatial perspectives.

From a temporal perspective, we investigate how short- and long-term temporal information affects the model performance. We control the short- and long-term information through adjusting the ratio of CNN and Transformer layers in two hybrid models. From the results in Table 2(middle), we observe a decline in performance as the proportion of short-term information decreases, demonstrating the effectiveness of short-range information for neural decoding.

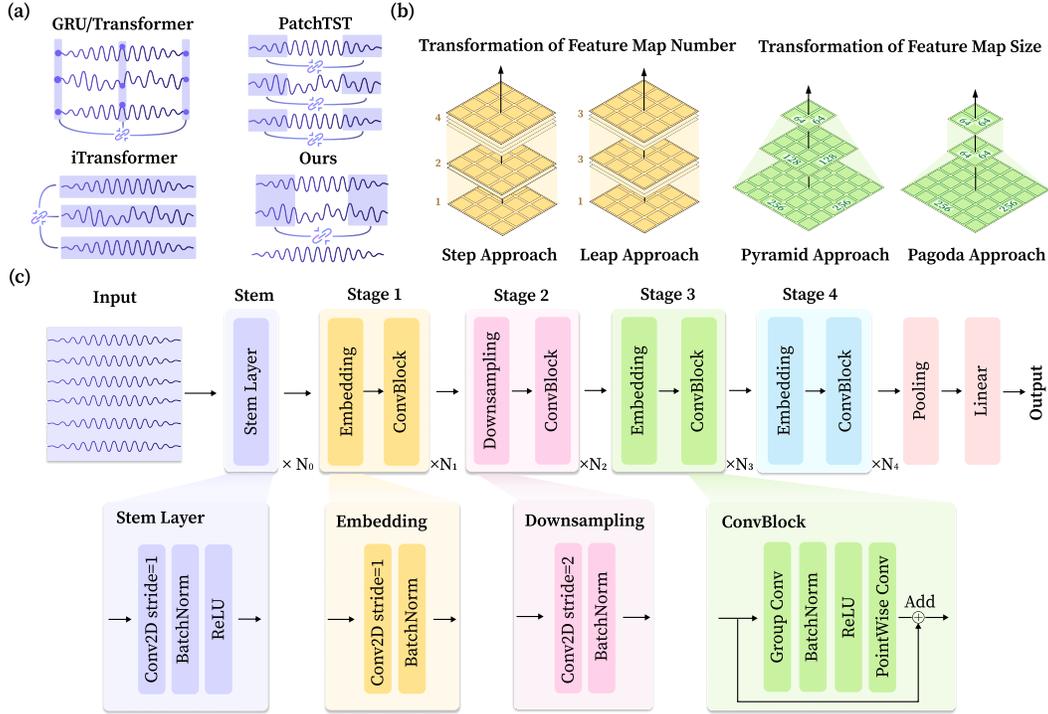


Figure 2: (a) **Comparison of different patch methods.** Vanilla GRU and Transformer treat a single timestamp across all channels as a token; PatchTST treats multiple timestamps of a single channel as a token; iTransformer treats all timestamps of an entire channel as a token. We propose a simple but effective patch method that aggregates information of multiple timestamps and channels. (b) **Comparison of different latent space transformation methods.** Regarding the number of feature maps, the step approach increases them gradually, whereas the leap approach increases them rapidly in the early stages. In terms of their size, the pyramid approach decreases them progressively, while the pagoda approach decreases them quickly in the early stages. (c) **The overall architecture of NeuroSketch.**

From a spatial perspective, based on the effectiveness of CNN-2D, we can similarly apply this finding to other models, such as GRU and Transformer, to further validate its applicability. We propose a simple but effective patch method to capture local channel information in Transformers and GRUs. As shown in Figure 2(a), the input sequence is first divided into patches using a sliding window approach. Unlike the patch methods proposed in PatchTST and iTransformer, which treat channels independently or focus on global channel information, we reshape the patches by concatenating the channel and patch dimensions to preserve the local channel relationships within the same time window. The results in Table 2(lower) demonstrate that our method achieves highly competitive performance compared to others. By incorporating our token embedding method, the accuracy of the Transformer on the Du-IN dataset improved significantly from 8.5% to 23.4%. Additionally, the accuracy of the GRU on the OpenMIIR dataset’s perception and imagination tasks increased dramatically from the chance level to 95.8% and 95.7%, respectively, further demonstrating the effectiveness of spatial locality for neural decoding.

From now on, we will use CNN-2D as our basic architecture.

2.3 LATENT SPACE TRANSFORMATION

To address the macro optimization aspect of Question 2, we study the latent space transformation during the forward propagation. For CNN-2D, its latent space representation refers to $\mathbf{S} \in \mathbb{R}^{C \times H \times W}$, where C denotes the number of feature maps, H and W denote the height and width of feature maps, respectively. During forward propagation, the network increases the number of feature maps to capture higher-level information, while reducing the size of feature maps to alleviate redundancy and improve computational efficiency. Building on this general overview, we will investigate the transformation of the number (C) and the size (H and W) of feature maps separately.

Regarding the number of feature maps, the appropriate strategy to increase them is a critical aspect of neural network design, and this can be achieved through two widely adopted approaches: **step approach** and **leap approach**. As shown in Figure 2(b), The step approach gradually increases the number of feature maps, allowing feature extraction to advance from low-level to high-level in a systematic manner. This progression enables the network to concentrate on varying levels of features effectively. The leap approach takes a more aggressive strategy to rapidly increase the number of feature maps to the embedding dimension in the early stages, making the model focus on refining and extracting higher-level features.

Table 3: **Results of the latent space transformation and the optimization of calculation.** Lower GFLOPs indicates lower computational cost. The best results are in **bold**.

Setting	GFLOPs	Datasets		Du-IN		SEED-DV		ThingsEEG		FacesHouses		OpenMIIR-P		OpenMIIR-I		Chisco-R		Chisco-I	
		Acc	F1	Acc	F1														
<i>Macro Study: Transformation of the Feature Map Number</i>																			
Leap Approach	3996	.635	.630	.064	.058	.216	.206	.857	.856	.980	.980	.976	.976	.126	.109	.115	.093		
Step Approach	1287	.647	.641	.063	.058	.184	.177	.886	.885	.981	.980	.982	.982	.127	.113	.114	.095		
<i>Macro Study: Transformation of the Feature Map Size</i>																			
Pyramid Approach	1287	.647	.641	.063	.058	.184	.177	.886	.885	.981	.980	.982	.982	.127	.113	.114	.095		
Pagoda Approach	758	.701	.698	.059	.053	.204	.197	.911	.910	.973	.973	.986	.985	.128	.116	.116	.096		
<i>Micro Study</i>																			
Vanilla Conv	758	.701	.698	.059	.053	.204	.197	.911	.910	.973	.973	.986	.985	.128	.116	.116	.096		
Seperable Conv	492	.704	.699	.061	.051	.181	.176	.915	.915	.979	.979	.985	.985	.118	.107	.111	.098		
Group Conv	482	.707	.704	.069	.061	.207	.200	.920	.920	.983	.983	.990	.990	.129	.112	.120	.103		

Table 3(upper) compares the performance and computational cost between the step and leap approach. Surprisingly, the step approach achieves highly competitive performance compared to the leap approach, while requiring 67.8% fewer FLOPs. For instance, the step approach attains an accuracy of 88.6% on the FacesHouses dataset, to achieve similar accuracy, the leap approach requires three times more computational cost, while only reaching 85.7% accuracy. This indicates that neural decoding does not rely solely on high-level neural representations, and that multi-scale feature extraction is often more effective and efficient.

From now on, we will use the step approach to transform the number of feature maps.

In the forward propagation process, the transformation of the feature map size is achieved through downsampling. Therefore, we further investigate how downsampling should be distributed throughout the network. Suppose downsampling is applied only in the final layers. In that case, earlier layers must process high-resolution feature maps, which significantly increases computational complexity and undermines the original purpose of downsampling to improve computational efficiency. Therefore, we study two common strategies: **pyramid approach** and **pagoda approach**. As shown in Figure 2(b), the pyramid approach distributes the downsampling process evenly across the entire network, gradually reducing the size of the feature maps. This pyramid-like structure enables the network to retain more detailed features at the cost of increased computational complexity. The pagoda approach takes a more aggressive strategy by concentrating the downsampling module in the early stages of the network. This leads to a significant reduction in feature map size early on, resulting in lower computational costs. Table 3(middle) compares the performance and computational cost between the pyramid and pagoda approach. The pagoda approach achieves slightly better accuracy than the pyramid approach, requiring 41.1% fewer FLOPs. Specifically, on the Du-IN dataset, the pagoda approach attains 70.1% accuracy, whereas the pyramid approach only reaches

64.7% with a higher computational cost. Additionally, the pagoda approach improves the accuracy by 2.0% to 20.4% on the ThingsEEG dataset, surpassing CNN-1D’s accuracy of 20.2%, making it the best-performing architecture. These results demonstrate that excessively extracting detailed information from the original signal is unnecessary, as brain signals typically have low signal-to-noise ratios. By applying downsampling early in the network, we can enhance computational efficiency while improving performance, allowing the model to focus on the most relevant features and reduce the impact of noise.

From now on, we will use the pagoda approach to transform the size of feature maps.

2.4 OPTIMIZATION OF CALCULATION

To further explore the micro optimization aspect of Question 2, we investigate the core calculation method of CNNs, the convolution operation. Given $\mathbf{X}_{in} \in \mathbb{R}^{C_{in} \times H_{in} \times W_{in}}$ as the C_{in} channels input of height H_{in} and width W_{in} , the convolution kernel $\mathbf{F} \in \mathbb{R}^{C_{out} \times C_{in} \times U \times V}$ of height U and width V slides across the input to compute the output $\mathbf{X}_{out} \in \mathbb{R}^{C_{out} \times H_{out} \times W_{out}}$, where C_{out} , H_{out} and W_{out} refer to the output channels, height and width, respectively. The number of parameters involved in this convolution operation is $C_{out} \times C_{in} \times U \times V$, representing the learnable weights within the convolution kernel.

In addition to the vanilla convolution, various variants have been proposed to improve computational efficiency and model performance, including group convolution (Ioannou et al., 2017) and separable convolution (Howard, 2017).

Group convolution. Group convolution $\mathbf{F}_g \in \mathbb{R}^{C_{out} \times \frac{C_{in}}{G} \times U \times V}$ divides the input channels C_{in} and output channels C_{out} into G groups. The convolution is then applied separately within each group, limiting the channel interaction and effectively reducing the parameters by a factor of G compared to the vanilla convolution.

Separable convolution. Separable convolution breaks down the vanilla 2D convolution into two distinct operations: depthwise convolution $\mathbf{F}_d \in \mathbb{R}^{C_{in} \times U \times V}$ for spatial dimensions and pointwise convolution $\mathbf{F}_p \in \mathbb{R}^{C_{out} \times C_{in} \times 1 \times 1}$ for channel dimension. Depthwise convolution is an extreme case of group convolution, where the number of groups $G = C_{in} = C_{out}$. In this case, each output feature map corresponds directly to a specific input feature map. To enable channel interactions, a pointwise convolution is typically applied afterward, where both the kernel height and width are set to 1. In total, the number of parameters in a separable convolution is $\frac{1}{C_{out}} + \frac{1}{UV}$ of that in vanilla convolution.

Here, we examine how performance changes when replacing the vanilla convolution with separable and group convolution. As shown in Table 3(lower), the group convolution variant achieves slightly better performance compared to other convolution methods while minimizing computational cost. Besides, the separable convolution variant also shows competitive performance with less computational cost than vanilla convolution. For example, with 758 GFLOPs, the vanilla convolution gets 5.9% accuracy on the SEED-DV dataset, while the separable convolution can obtain 6.1% accuracy with 35% fewer GFLOPs. Meanwhile, with nearly identical GFLOPs, the group convolution outperforms the separable convolution, improving accuracy by 0.8% to 6.9%. The experimental results show that grouping input and output channels for convolution does not degrade performance in neural decoding tasks. Channel grouping may be more suitable for aggregating local information while improving computational efficiency.

We will use group convolution as our core calculation method. This brings us to our final framework, NeuroSketch.

2.5 PUT IT ALL TOGETHER

After the optimizations mentioned above, we have developed an effective framework, NeuroSketch. The overall architecture is shown in Figure 2(c). Like CNN-2D, NeuroSketch employs a 2D input representation, which is then processed by a stem layer to initiate feature extraction. The following forward propagation process is divided into four stages to progressively capture features from low-level to high-level, as suggested by the step approach. In each stage, the initial component is responsible for increasing the number of feature maps. If this component is a downsampling layer, it

also reduces the spatial resolution of the feature maps by half. Based on the pagoda approach, we allocate the downsampling layers to the second stage to enhance performance and computational efficiency. Subsequently, a convolutional block is applied to extract features, comprising group convolutions, batch normalization, ReLU activation, and pointwise convolution. Following the four stages, the resulting feature is passed through a generalized mean (GeM) pooling layer (Berman et al., 2019) that aggregates features along the channel and temporal dimensions. The pooled representation is then fed into a linear layer for the final classification. Additional details of the architecture and implementation can be found in Appendix D.2.

3 EXPERIMENTS

In this section, we present a comprehensive evaluation of NeuroSketch based on more than 3,000 experiments. Section 3.1 introduces the datasets, baselines, and experimental settings, while Section 3.2 reports the main results, including analyses across three different modalities: speech decoding, visual decoding, and auditory decoding.

3.1 EXPERIMENT SETUP

We employ the same datasets described in Section 2 for evaluation. The diverse selection of datasets ensures a comprehensive evaluation across different neural decoding tasks. We implement two versions of NeuroSketch with different sizes: NeuroSketch-Base (1.4M parameters) and NeuroSketch-Large (4.2M parameters). Detailed implementation and configuration are provided in Appendix D.2. For baselines, we select representative and recent models from various domains, including time-series models: ModernTCN (Luo and Wang, 2024), MedFormer (Wang et al., 2024); computer vision backbones: ConvFormer (Yu et al., 2023), CAFormer (Yu et al., 2023); well-known brain models: DeepConvNet (Schirrmeyer et al., 2017), EEGNet (Lawhern et al., 2018); recent brain models: Conformer (Song et al., 2022), SPaRCNet (Jing et al., 2023); iEEG foundation models: seegnet (Mentzelopoulos et al., 2024); and EEG foundation models: CBraMod (Wang et al., 2025b). We use the official pretrained weights for all foundation models. More details of each baseline and evaluation setup are introduced in Appendix D.3 and Appendix E.

3.2 RESULTS

Main results. Overall, NeuroSketch consistently achieves state-of-the-art performance across eight neural decoding tasks. As shown in Figure 3, NeuroSketch outperforms all selected baselines, underscoring its strong capacity to model brain signals across diverse decoding scenarios. In addition, we note that NeuroSketch-Base performs similarly to NeuroSketch-Large on all datasets except Du-IN, where NeuroSketch-Large outperforms NeuroSketch-Base by a substantial margin. We discuss this phenomenon in detail in Appendix F.3.

Speech decoding. Table 4 reports the excellent performance of NeuroSketch on the three speech decoding tasks. Notably, on the Du-IN dataset, NeuroSketch-Large surpasses the second-best baseline, ConvFormer, by 65.5% in accuracy, underscoring its strength in extracting discriminative features from neural signals. As an iEEG foundation model, seegnet performs poorly on Du-IN (5.3% accuracy), likely because its single-layer transformer lacks sufficient capacity for this challenging neural decoding task. Classifying semantic categories while subjects read or imagine sentences is even more difficult. However, NeuroSketch achieves the best results, demonstrating robustness in challenging speech decoding scenarios. Specifically, compared to the second-best baseline, Conformer, NeuroSketch-Large improves accuracy by 22.1% on Chisco-R and NeuroSketch-Base by 9.2% on Chisco-I.

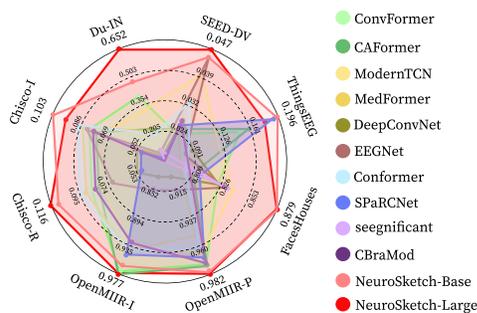


Figure 3: Model performance comparison.

Table 4: **Average classification accuracy** (mean \pm std across three folds) for models on multiple neural decoding datasets. Higher throughput indicates faster inference. The best results are in **bold**. We use $-$ to denote entries that are not applicable or not evaluated. Due to substantial differences between iEEG and EEG, we evaluate each pretrained model only on the modality it was pretrained on, the corresponding cross-modality cells are therefore marked with $-$. In addition, CBraMod uses a patch size of 200, whereas *ThingsEEG* has an input length of 100, so CBraMod cannot be run on *ThingsEEG* and is marked $-$.

Models \ Datasets	Speech Decoding			Visual Decoding		Auditory Decoding		
	Du-IN	Chisco-R	Chisco-I	SEED-DV	Things-EEG	Faces-Houses	Open-MIIR-P	Open-MIIR-I
ConvFormer	.394 \pm .005	.086 \pm .001	.080 \pm .004	.024 \pm .000	.151 \pm .003	.820 \pm .013	.975 \pm .005	.976 \pm .002
CAFormer	.221 \pm .003	.081 \pm .003	.076 \pm .002	.025 \pm .001	.160 \pm .009	.800 \pm .012	.976 \pm .006	.977 \pm .002
ModernTCN	.339 \pm .005	.079 \pm .003	.074 \pm .002	.040 \pm .001	.123 \pm .003	.835 \pm .006	.964 \pm .010	.935 \pm .002
MedFormer	.212 \pm .006	.089 \pm .003	.085 \pm .001	.025 \pm .000	.045 \pm .002	.756 \pm .011	.948 \pm .001	.941 \pm .001
DeepConvNet	.053 \pm .002	.053 \pm .004	.052 \pm .001	.029 \pm .001	.091 \pm .001	.822 \pm .005	.896 \pm .009	.789 \pm .005
EEGNet	.060 \pm .003	.071 \pm .002	.082 \pm .002	.046 \pm .002	.070 \pm .000	.830 \pm .006	.915 \pm .010	.852 \pm .006
Conformer	.205 \pm .003	.095 \pm .002	.087 \pm .002	.033 \pm .002	.077 \pm .002	.762 \pm .001	.962 \pm .013	.930 \pm .002
SPaRCNet	.026 \pm .001	.005 \pm .000	.005 \pm .000	.026 \pm .001	.189 \pm .001	.805 \pm .007	.969 \pm .007	.948 \pm .005
seegnicant	.053 \pm .002	$-$	$-$	$-$	$-$.831 \pm .006	$-$	$-$
CBraMod	$-$.084 \pm .002	.077 \pm .001	.026 \pm .002	$-$	$-$.975 \pm .011	.933 \pm .008
NeuroSketch-Base	.472 \pm .006	.111 \pm .001	.103\pm.002	.045 \pm .002	.196\pm.002	.879 \pm .008	.981 \pm .010	.970 \pm .002
NeuroSketch-Large	.652\pm.002	.116\pm.001	.095 \pm .005	.047\pm.003	.177 \pm .001	.879\pm.007	.982\pm.009	.977\pm.002

Visual decoding. As shown in Table 4, for static image decoding on the FacesHouses dataset, NeuroSketch-Large outperforms the second-best performing baseline, ModernTCN, with a 5.2% improvement in accuracy. Furthermore, on the considerably more challenging ThingsEEG dataset, NeuroSketch-Base achieves competitive performance with the recent brain model, SPaRCNet. ConvFormer and CAFormer also demonstrate strong visual decoding performance on the ThingsEEG datasets. For the highly challenging video decoding task, where most baseline models perform near the chance level (2.5% accuracy), NeuroSketch-Large demonstrates a clear advantage with an accuracy of 4.7%. This result highlights its capacity to capture complex neural representations associated with dynamic visual stimuli.

Auditory decoding. The results in Table 4 show that NeuroSketch achieves state-of-the-art performance on both auditory decoding tasks. On the OpenMIIR-P dataset, the EEG foundation model, CBraMod attains the strongest baseline among brain-domain models (97.5% accuracy), highlighting the effectiveness of pretraining and its transferability to auditory decoding.

4 CONCLUSION AND FUTURE WORK

Conclusion. In this paper, we conduct an in-depth exploration of model architectures for neural decoding. By investigating basic architectures, latent space transformation, and computational methods optimization, we introduce NeuroSketch, an effective framework for neural decoding. Experimental results demonstrate that NeuroSketch achieves state-of-the-art performance across eight distinct neural decoding tasks. We hope this research will provide new insights for the neural decoding community and encourage further exploration of CNN-based architectures.

Limitations and future work. While NeuroSketch shows strong performance, its current design focuses on supervised learning. Given the favorable scaling behavior observed in Appendix F.3, we hypothesize that scaling both data and model capacity will yield further gains. Accordingly, we will explore large-scale pretraining strategies based on NeuroSketch in the future. Additionally, we plan to incorporate a broader range of neural decoding tasks to enhance the model’s applicability across diverse neural decoding scenarios.

5 ETHICS STATEMENT

All datasets used in this study are publicly available. The corresponding dataset links are provided in Appendix C.1. Consequently, our work does not involve personally sensitive information and does not pose apparent ethical concerns.

6 REPRODUCIBILITY STATEMENT

We place great importance on ensuring the reproducibility of our work. To this end, details of the preprocessing of each dataset are provided in Appendix C.2, the detailed model implementation and configuration are described in Appendix D.2, and the experimental settings are outlined in Appendix E. The complete experimental results are available at https://anonymous.4open.science/r/NeuroSketch_Results-BA5B. We will release our code and scripts upon publication, thereby facilitating transparent and reproducible research for the community.

REFERENCES

- Baraka Maiseli, Abdi T. Abdalla, Libe V. Massawe, Mercy Mbise, Khadija Mkocho, Nassor Ally Nassor, Moses Ismail, James Michael, and Samwel Kimambo. Brain-computer interface: Trend, challenges, and threats. *Brain Informatics*, 10(1):20, December 2023. ISSN 2198-4018, 2198-4026. doi: 10.1186/s40708-023-00199-3.
- Marcel A. J. Van Gerven, Katja Seeliger, Umut Güçlü, and Yağmur Güçlütürk. Current Advances in Neural Decoding. In Wojciech Samek, Grégoire Montavon, Andrea Vedaldi, Lars Kai Hansen, and Klaus-Robert Müller, editors, *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, volume 11700, pages 379–394. Springer International Publishing, Cham, 2019. ISBN 978-3-030-28953-9 978-3-030-28954-6. doi: 10.1007/978-3-030-28954-6_21.
- Michal Teplan. Fundamentals of EEG measurement. *Measurement science review*, 2(2):1–11, 2002.
- J. Talairach. Approche nouvelle de la neurochirurgie de l’épilepsie. Methodologie stérotaxique et resultats therapeutiques. *Neurochirurgie*, 20:1–240, 1974.
- Pradeep Shenoy, Kai J. Miller, Jeffrey G. Ojemann, and Rajesh PN Rao. Generalized features for electrocorticographic BCIs. *IEEE Transactions on Biomedical Engineering*, 55(1):273–280, 2007.
- Xiang Zhang, Lina Yao, Xianzhi Wang, Jessica Monaghan, David Mcalpine, and Yu Zhang. A survey on deep learning-based non-invasive brain signals: Recent advances and new frontiers. *Journal of neural engineering*, 18(3):031002, 2021.
- Hamdi Altaheri, Ghulam Muhammad, Mansour Alsulaiman, Syed Umar Amin, Ghadir Ali Altuwaijri, Wadood Abdul, Mohamed A. Bencherif, and Mohammed Faisal. Deep learning techniques for classification of electroencephalogram (EEG) motor imagery (MI) signals: A review. *Neural Computing and Applications*, 35(20):14681–14722, July 2023. ISSN 0941-0643, 1433-3058. doi: 10.1007/s00521-021-06352-5.
- Md Mustafizur Rahman, Ajay Krishno Sarkar, Md Amzad Hossain, Md Selim Hossain, Md Rabiul Islam, Md Biplob Hossain, Julian MW Quinn, and Mohammad Ali Moni. Recognition of human emotions using EEG signals: A review. *Computers in biology and medicine*, 136:104696, 2021.
- Sara Maria Pani, Luca Saba, and Matteo Frascini. Clinical applications of EEG power spectra aperiodic component analysis: A mini-review. *Clinical Neurophysiology*, 143:1–13, 2022.
- Sani Saminu, Guizhi Xu, Shuai Zhang, Isselmou Ab El Kader, Hajara Abdulkarim Aliyu, Adamu Halilu Jabire, Yusuf Kola Ahmed, and Mohammed Jajere Adamu. Applications of artificial intelligence in automatic detection of epileptic seizures using EEG signals: A review. In *Artificial Intelligence and Applications*, volume 1, pages 11–25, 2023.
- Ahmad Chaddad, Yihang Wu, Reem Kateb, and Ahmed Bouridane. Electroencephalography signal processing: A comprehensive review and analysis of methods and techniques. *Sensors*, 23(14):6434, 2023.

- 540 Jean-Rémi King and Stanislas Dehaene. Characterizing the dynamics of mental representations: The
541 temporal generalization method. *Trends in cognitive sciences*, 18(4):203–210, 2014.
- 542
- 543 Keyvan Mahjoory, Andreas Bahmer, and Molly J. Henry. Convolutional neural networks can identify
544 brain interactions involved in decoding spatial auditory attention. *PLOS Computational Biology*,
545 20(8):e1012376, 2024.
- 546 Huy Phan and Kaare Mikkelsen. Automatic sleep staging of EEG signals: Recent development,
547 challenges, and future directions. *Physiological Measurement*, 43(4):04TR01, 2022.
- 548
- 549 Daniel A. Butts, Chong Weng, Jianzhong Jin, Chun-I. Yeh, Nicholas A. Lesica, Jose-Manuel Alonso,
550 and Garrett B. Stanley. Temporal precision in the neural code and the timescales of natural vision.
551 *Nature*, 449(7158):92–95, 2007.
- 552 Carl E. Stafstrom and Lionel Carmant. Seizures and epilepsy: An overview for neuroscientists. *Cold
553 Spring Harbor perspectives in medicine*, 5(6):a022426, 2015.
- 554
- 555 Matthew K. Leonard, Laura Gwilliams, Kristin K. Sellers, Jason E. Chung, Duo Xu, Gavin Mischler,
556 Nima Mesgarani, Marleen Welkenhuysen, Barundeb Dutta, and Edward F. Chang. Large-scale
557 single-neuron speech sound encoding across the depth of human cortex. *Nature*, 626(7999):
558 593–602, 2024.
- 559 Kalanit Grill-Spector and Rafael Malach. THE HUMAN VISUAL CORTEX. *Annual Review of
560 Neuroscience*, 27(1):649–677, July 2004. ISSN 0147-006X, 1545-4126. doi: 10.1146/annurev.
561 neuro.27.070203.144220.
- 562
- 563 Janis Peksa and Dmytro Mamchur. State-of-the-art on brain-computer interface technology. *Sensors*,
564 23(13):6001, 2023.
- 565 Benjamin Wittevrongel, Elvira Khachatryan, Evelien Carrette, Paul Boon, Alfred Meurs, Dirk
566 Van Roost, and Marc M. Van Hulle. High-gamma oscillations precede visual steady-state responses:
567 A human electrocorticography study. *Human Brain Mapping*, 41(18):5341–5355, December 2020.
568 ISSN 1065-9471, 1097-0193. doi: 10.1002/hbm.25196.
- 569
- 570 Timothée Proix, Jaime Delgado Saa, Andy Christen, Stephanie Martin, Brian N. Pasley, Robert T.
571 Knight, Xing Tian, David Poeppel, Werner K. Doyle, and Orrin Devinsky. Imagined speech can be
572 decoded from low-and cross-frequency intracranial EEG features. *Nature communications*, 13(1):
573 48, 2022.
- 574 Mehrnoosh Sadat Safi and Seyed Mohammad Mehdi Safi. Early detection of Alzheimer’s disease
575 from EEG signals using Hjorth parameters. *Biomedical Signal Processing and Control*, 65:102338,
576 2021.
- 577 Miguel Angrick, Christian Herff, Emily Mugler, Matthew C. Tate, Marc W. Slutzky, Dean J. Krusien-
578 ski, and Tanja Schultz. Speech synthesis from ECoG using densely connected 3D convolutional
579 neural networks. *Journal of neural engineering*, 16(3):036019, 2019.
- 580
- 581 Joseph G. Makin, David A. Moses, and Edward F. Chang. Machine translation of cortical activity to
582 text with an encoder–decoder framework. *Nature neuroscience*, 23(4):575–582, 2020.
- 583
- 584 Guy H. Wilson, Sergey D. Stavisky, Francis R. Willett, Donald T. Avansino, Jessica N. Kelemen,
585 Leigh R. Hochberg, Jaimie M. Henderson, Shaul Druckmann, and Krishna V. Shenoy. Decoding
586 spoken English from intracortical electrode arrays in dorsal precentral gyrus. *Journal of neural
587 engineering*, 17(6):066007, 2020.
- 588
- 589 Francis R. Willett, Erin M. Kunz, Chaofei Fan, Donald T. Avansino, Guy H. Wilson, Eun Young
590 Choi, Foram Kamdar, Matthew F. Glasser, Leigh R. Hochberg, and Shaul Druckmann. A high-
591 performance speech neuroprosthesis. *Nature*, 620(7976):1031–1036, 2023.
- 592
- 593 Sean L. Metzger, Kaylo T. Littlejohn, Alexander B. Silva, David A. Moses, Margaret P. Seaton,
Ran Wang, Maximilian E. Dougherty, Jessie R. Liu, Peter Wu, and Michael A. Berger. A
high-performance neuroprosthesis for speech decoding and avatar control. *Nature*, 620(7976):
1037–1046, 2023.

- 594 Hui Zheng, Haiteng Wang, Weibang Jiang, Zhongtao Chen, Li He, Peiyang Lin, Penghu Wei,
595 Guoguang Zhao, and Yunzhe Liu. Du-IN: Discrete units-guided mask modeling for decoding
596 speech from Intracranial Neural signals. *Advances in Neural Information Processing Systems*, 37:
597 79996–80033, 2024.
- 598 Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and
599 Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF*
600 *Conference on Computer Vision and Pattern Recognition*, pages 10819–10829, 2022.
- 601 Weihao Yu, Chenyang Si, Pan Zhou, Mi Luo, Yichen Zhou, Jiashi Feng, Shuicheng Yan, and Xinchao
602 Wang. Metaformer baselines for vision. *IEEE Transactions on Pattern Analysis and Machine*
603 *Intelligence*, 46(2):896–912, 2023.
- 604 Donghao Luo and Xue Wang. Moderntcn: A modern pure convolution structure for general time
605 series analysis. In *The Twelfth International Conference on Learning Representations*, pages 1–43,
606 2024.
- 607 Shiyu Wang, Jiawei Li, Xiaoming Shi, Zhou Ye, Baichuan Mo, Wenzhe Lin, Shengtong Ju, Zhixuan
608 Chu, and Ming Jin. TimeMixer++: A General Time Series Pattern Machine for Universal Predictive
609 Analysis, March 2025a.
- 610 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
611 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings*
612 *of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11106–11115, 2021.
- 613 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.
614 A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
615 *Pattern Recognition*, pages 11976–11986, 2022.
- 616 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
617 recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
618 pages 770–778, 2016.
- 619 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
620 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
621 pages 248–255. Ieee, 2009.
- 622 Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical Evaluation of
623 Gated Recurrent Neural Networks on Sequence Modeling, December 2014.
- 624 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz
625 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
626 *systems*, 30, 2017.
- 627 Zihan Zhang, Xiao Ding, Yu Bao, Yi Zhao, Xia Liang, Bing Qin, and Ting Liu. Chisco: An
628 EEG-based BCI dataset for decoding of imagined speech. *Scientific Data*, 11(1):1265, 2024.
- 629 Kai J. Miller. A library of human electrocorticographic data and analyses. *Nature human behaviour*,
630 3(11):1225–1235, 2019.
- 631 Alessandro T Gifford, Kshitij Dwivedi, Gemma Roig, and Radoslaw M Cichy. A large and rich eeg
632 dataset for modeling human visual object recognition. *NeuroImage*, 264:119754, 2022.
- 633 Xuan-Hao Liu, Yan-Kai Liu, Yansen Wang, Kan Ren, Hanwen Shi, Zilong Wang, Dongsheng Li,
634 Bao-Liang Lu, and Wei-Long Zheng. EEG2video: Towards decoding dynamic visual perception
635 from EEG signals. *Advances in Neural Information Processing Systems*, 37:72245–72273, 2024a.
- 636 Sebastian Stober, Avital Sternin, Adrian M. Owen, and Jessica A. Grahn. Towards music imagery
637 information retrieval: Introducing the OpenMIIR dataset of EEG recordings from music perception
638 and imagination. In *ISMIR*, pages 763–769, 2015.
- 639 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A Time Series is Worth
640 64 Words: Long-term Forecasting with Transformers, March 2023.

- 648 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.
649 iTransformer: Inverted Transformers Are Effective for Time Series Forecasting, March 2024b.
650
- 651 Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. EEG conformer: Convolutional
652 transformer for EEG decoding and visualization. *IEEE Transactions on Neural Systems and*
653 *Rehabilitation Engineering*, 31:710–719, 2022.
- 654 Sehoon Kim, Amir Gholami, Albert Shaw, Nicholas Lee, Karttikeya Mangalam, Jitendra Malik,
655 Michael W. Mahoney, and Kurt Keutzer. Squeezeformer: An efficient transformer for automatic
656 speech recognition. *Advances in Neural Information Processing Systems*, 35:9361–9373, 2022.
657
- 658 Yani Ioannou, Duncan Robertson, Roberto Cipolla, and Antonio Criminisi. Deep roots: Improving
659 cnn efficiency with hierarchical filter groups. In *Proceedings of the IEEE Conference on Computer*
660 *Vision and Pattern Recognition*, pages 1231–1240, 2017.
- 661 Andrew G. Howard. Mobilenets: Efficient convolutional neural networks for mobile vision applica-
662 tions. *arXiv preprint arXiv:1704.04861*, 2017.
663
- 664 Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. Multigrain: a
665 unified image embedding for classes and instances. *arXiv preprint arXiv:1902.05509*, 2019.
- 666 Yihe Wang, Nan Huang, Taida Li, Yujun Yan, and Xiang Zhang. Medformer: A Multi-Granularity
667 Patching Transformer for Medical Time-Series Classification, October 2024.
668
- 669 Robin Tibor Schirrmester, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin
670 Glasstetter, Katharina Eggenberger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and
671 Tonio Ball. Deep learning with convolutional neural networks for eeg decoding and visualization.
672 *Human brain mapping*, 38(11):5391–5420, 2017.
- 673 Vernon J. Lawhern, Amelia J. Solon, Nicholas R. Waytowich, Stephen M. Gordon, Chou P. Hung, and
674 Brent J. Lance. EEGNet: A compact convolutional neural network for EEG-based brain–computer
675 interfaces. *Journal of neural engineering*, 15(5):056013, 2018.
676
- 677 Jin Jing, Wendong Ge, Shenda Hong, Marta Bento Fernandes, Zhen Lin, Chaoqi Yang, Sungtae An,
678 Aaron F. Struck, Aline Herlopian, Ioannis Karakis, et al. Development of Expert-Level Classifica-
679 tion of Seizures and Rhythmic and Periodic Patterns During EEG Interpretation. *Neurology*, 100
680 (17), April 2023. ISSN 0028-3878, 1526-632X. doi: 10.1212/WNL.000000000207127.
- 681 Georgios Mentzelopoulos, Evangelos Chatzipantazis, Ashwin G Ramayya, Michelle J Hedlund,
682 Vivek P Buch, Kostas Daniilidis, Konrad P Kording, and Flavia Vitale. Neural decoding from
683 stereotactic eeg: accounting for electrode variability across subjects. *Advances in Neural Informa-*
684 *tion Processing Systems*, 37:108600–108624, 2024.
- 685 Jiquan Wang, Sha Zhao, Zhiling Luo, Yangxuan Zhou, Haiteng Jiang, Shijian Li, Tao Li, and Gang
686 Pan. CBraMod: A Criss-Cross Brain Foundation Model for EEG Decoding, April 2025b.
687
- 688 Dhananjay Sonawane, Krishna Prasad Miyapuram, Bharatesh Rs, and Derek J. Lomas. GuessThe-
689 Music: Song Identification from Electroencephalography response. In *Proceedings of the 3rd*
690 *ACM India Joint International Conference on Data Science & Management of Data (8th ACM*
691 *IKDD CODS & 26th COMAD)*, pages 154–162, Bangalore India, January 2021. ACM. ISBN
692 978-1-4503-8817-7. doi: 10.1145/3430984.3431023.
- 693 Aline W. De Borst, Giancarlo Valente, Iiro P. Jääskeläinen, and Pia Tikka. Brain-based decoding
694 of mentally imagined film clips and sounds reveals experience-based information patterns in film
695 professionals. *NeuroImage*, 129:428–438, 2016.
696
- 697 Chen Feng, Lu Cao, Di Wu, En Zhang, Ting Wang, Xiaowei Jiang, Heng Ding, Chenhao Zhou, Jinbo
698 Chen, and Hui Wu. A high-performance brain-to-sentence decoder for logossyllabic language.
699 *bioRxiv*, pages 2023–11, 2023.
- 700 Guy Gaziv, Roman Beliy, Niv Granot, Assaf Hoogi, Francesca Strappini, Tal Golan, and Michal Irani.
701 Self-supervised natural image reconstruction and large-scale semantic classification from brain
activity. *NeuroImage*, 254:119121, 2022.

- Hajar Ahmadiéh, Farnaz Gassemi, and Mohammad Hasan Moradi. Visual image reconstruction based on EEG signals using a generative adversarial and deep fuzzy neural network. *Biomedical Signal Processing and Control*, 87:105497, 2024.
- Prajwal Singh, Pankaj Pandey, Krishna Miyapuram, and Shanmuganathan Raman. EEG2IMAGE: Image reconstruction from EEG brain signals. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- Haiguang Wen, Junxing Shi, Yizhen Zhang, Kun-Han Lu, Jiayue Cao, and Zhongming Liu. Neural encoding and decoding with deep learning for dynamic natural vision. *Cerebral cortex*, 28(12): 4136–4160, 2018.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019.
- Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10428–10436, 2020.
- Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. Score-cam: Score-weighted visual explanations for convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 24–25, 2020.
- Kristofer E Bouchard, Nima Mesgarani, Keith Johnson, and Edward F Chang. Functional organization of human sensorimotor cortex for speech articulation. *Nature*, 495(7441):327–332, 2013.
- Gregory Hickok and David Poeppel. The cortical organization of speech processing. *Nature reviews neuroscience*, 8(5):393–402, 2007.
- Josh Chartier, Gopala K Anumanchipalli, Keith Johnson, and Edward F Chang. Encoding of articulatory kinematic trajectories in human speech sensorimotor cortex. *Neuron*, 98(5):1042–1054, 2018.
- Annan Yu, Danielle C Maddix, Boran Han, Xiyuan Zhang, Abdul Fatir Ansari, Oleksandr Shchur, Christos Faloutsos, Andrew Gordon Wilson, Michael W Mahoney, and Yuyang Wang. Understanding transformers for time series: Rank structure, flow-of-ranks, and compressibility. *arXiv preprint arXiv:2510.03358*, 2025.
- Zida Liang, Jiayi Zhu, and Weiqiang Sun. Why attention fails: The degeneration of transformers into mlps in time series forecasting. *arXiv preprint arXiv:2509.20942*, 2025.

A THE USE OF LARGE LANGUAGE MODELS (LLMs)

In preparing this manuscript, we made limited use of a large language model (LLM) to assist with polishing the writing for grammar and clarity. In addition, the LLM was occasionally employed for small-scale code completion, such as generating boilerplate functions or suggesting minor syntax corrections. All core ideas, experimental designs, analyzes, and conclusions were completely conceived, implemented, and validated by the authors.

B RELATED WORK

Neural decoding. Decoding sensory experiences from neural signals is a central objective in neuroscience, spanning modalities such as auditory, speech, and visual processing, each with distinct decoding challenges and objectives. Auditory decoding can be divided into two distinct processes: auditory perception, which decodes brain activity evoked by external sounds (Stober et al., 2015; Sonawane et al., 2021), and auditory imagery, which identifies patterns associated with internally generated mental simulations of sound (Stober et al., 2015; De Borst et al., 2016). Speech decoding seeks to reconstruct vocal expressions or semantic content directly from neural signals, which can be

756 conducted at multiple linguistic levels, from syllables (Feng et al., 2023) to words (Zheng et al., 2024)
 757 and entire sentences (Zhang et al., 2024). Visual decoding focuses on reconstructing high-quality
 758 static images (Gaziv et al., 2022; Ahmadieh et al., 2024; Singh et al., 2023) or dynamic videos (Wen
 759 et al., 2018; Liu et al., 2024a) perceived by the human visual system from neural signals.

760
 761 **Framework design.** Recent advances in neural architecture design have diverged into two directions:
 762 automated neural architecture search (NAS) and manual design driven by experience. Tan and Le
 763 (2019) utilize NAS to scale all dimensions of depth, width, and resolution through an effective
 764 compound coefficient, leading to a family of models known as EfficientNets. Radosavovic et al.
 765 (2020) further integrate NAS with manual design principles to explore the design space, resulting
 766 in the development of RegNet. Meanwhile, manual architectural innovations revitalize traditional
 767 components. Yu et al. (2022) reveals that instead of a specific token mixer (e.g., attention), the
 768 general architecture of Transformers, termed MetaFormer, is more essential for achieving high
 769 performance. Liu et al. (2022) modernize a standard CNN towards the design of Swin Transformer
 770 through macro- and micro-level optimizations, and ultimately propose a family of pure CNNs dubbed
 771 ConvNeXt.

772 C DETAILS OF DATASETS

773 C.1 DATASET DESCRIPTION

774
 775 In this study, we use six neural decoding datasets, each associated with specific experimental tasks.
 776 Brief descriptions of each dataset are provided below.

777
 778 **Du-IN** (Zheng et al., 2024) focuses on decoding spoken Mandarin words. In each trial, a word is first
 779 displayed in white for 0.5 seconds, then turns green for 2 seconds to prompt articulation. Participants
 780 are instructed to read the word aloud during this cue period. Following this, the word disappears,
 781 and a fixation cross is presented. EEG data are recorded from 12 participants using 7-13 SEEG
 782 electrodes, comprising 72–158 channels. Each participant reads 61 predefined Chinese words, each
 783 word repeated 50 times, while both SEEG and audio signals are recorded simultaneously. In total,
 784 each participant contributes approximately 15 hours of data sampled at 2000 Hz, with around 3 hours
 785 corresponding to task-related activity. These task-related segments are divided into 3000 three-second
 786 trials and subsequently downsampled to 1000 Hz. The dataset is licensed under CC BY 4.0 and is
 787 available at <https://huggingface.co/datasets/liulab-repository/Du-IN>.

788
 789 **Chisco** (Zhang et al., 2024) focuses on decoding silent reading and imagined speech. Data are
 790 collected from three healthy participants using a 125-channel EEG system at a sampling rate of
 791 1000 Hz. The experiment consists of 45 blocks, each comprising 150 trials. Each trial includes a
 792 5-second silent reading phase followed by a 3.3-second imagined speaking phase, with continuous
 793 EEG recording throughout. The stimuli include 6,681 commonly used Chinese sentences (ranging
 794 from 6 to 15 characters) drawn from 39 semantic categories. In each block, 10 trials are randomly
 795 selected for verbal recall. If a participant’s recalled sentence differs from the original by more than
 796 four characters, the trial is marked as incorrect. If two or more errors occur within a block, the
 797 participant takes a rest and repeats the block. The dataset is licensed under CC0 and is available at
<https://openneuro.org/datasets/ds005170/versions/1.1.2>.

798
 799 **FacesHouses** (Miller, 2019) focuses on visual decoding of static images. ECoG signals are recorded
 800 from 14 epilepsy patients using subdural electrode strips implanted over the inferior temporal cortex.
 801 Participants view randomized grayscale images of faces and houses. The sampling rate is 1000 Hz.
 802 Each experiment includes three sessions, each presenting 50 face images and 50 house images for
 803 400 ms, with a 400 ms blank interval between images. The dataset is licensed under CC BY-SA and
 is available at <https://purl.stanford.edu/zk881ps0522>.

804
 805 **ThingsEEG** (Gifford et al., 2022) focuses on complex image decoding using a RSVP paradigm.
 806 EEG is recorded from 10 healthy participants using a 64-channel system at 1000 Hz. Stimuli consist
 807 of 16,740 image categories (16,540 for training and 200 for testing). Each training image is shown
 808 four times, while each test image is shown 80 times. Images are presented at a 200-ms stimulus
 809 onset asynchrony (SOA). A target detection task is included to maintain participant engagement. The
 decoding analysis is performed on the 200 test images. The dataset is licensed under CC BY 4.0 and
 is available at <https://osf.io/hd6zk/>.

810 **SEED-DV** (Liu et al., 2024a) focuses on decoding dynamic video stimuli. EEG is recorded from 20
811 participants at 1000 Hz while they watch video clips. The stimulus set consists of 1,400 two-second
812 video clips spanning 40 semantic concepts, with 35 clips per concept. Each participant completed
813 seven video blocks, with rest intervals between blocks. Each block contains all 40 concepts presented
814 in a randomized order. Before the start of each block, participants are informed of the target concept
815 and subsequently watch five video clips corresponding to that concept. The dataset is available after
816 submitting an application at [https://bcmi.sjtu.edu.cn/ApplicationForm/apply_](https://bcmi.sjtu.edu.cn/ApplicationForm/apply_form/)
817 [form/](https://bcmi.sjtu.edu.cn/ApplicationForm/apply_form/).

818 **OpenMIIR** (Stober et al., 2015) focuses on music perception and imagination. EEG signals are
819 recorded from 10 participants (8 with formal music training) using a 64-channel system at a sampling
820 rate of 512 Hz. The experiment comprises two sessions, each consisting of five trials. In each trial, 12
821 musical stimuli were randomly presented under one of four experimental conditions: (1) perception
822 with auditory cues, (2) imagination with cues, (3) imagination without cues, and (4) imagination
823 without cues but with feedback. The stimuli encompass two time signatures (3/4 and 4/4) and vary in
824 tempo from 104 to 212 BPM. Based on lyrical content, the stimuli are further classified into songs
825 with lyrics, their corresponding instrumental versions, and purely instrumental pieces. The dataset is
826 licensed under ODC-PDDL and is available at <https://hyper.ai/en/datasets/5591>.

828 C.2 DATA PREPROCESSING

829
830 To ensure consistency and comparability, we follow the preprocessing procedures described in the
831 original publications for each dataset.

832 For the **Du-IN** dataset, we followed the preprocessing setup described in the original work. For
833 each subject, we selected 10 SEEG channels following the original configuration. The signals
834 were downsampled to 1000 Hz, with each trial spanning 2.5 seconds. The classification targets
835 corresponded to the 61 predefined Chinese words presented to the participants, with the objective of
836 decoding spoken content directly from neural activity.

837 For the **Chisco** dataset, we examined two tasks: silent reading (denoted as Chisco-R) and imagined
838 speaking (denoted as Chisco-I). Following the protocol of the original work, we removed three noisy
839 channels and retained the remaining 122 valid EEG channels. All signals were downsampled to 500
840 Hz. The input duration was set to 5 seconds for Chisco-R and 3.3 seconds for Chisco-I, aligned
841 with the respective task lengths. For the classification targets, both datasets were labeled with the 39
842 predefined semantic categories.

843 For the **FacesHouses** dataset, we segmented the continuous ECoG recordings into epochs based on
844 stimulus markers, retained only face and house trials, and applied channel-wise z-score normalization
845 to reduce variability across channels. The input data comprised between 31 and 102 channels per
846 subject, sampled at 1000 Hz with a duration of 400 ms per trial. The classification labels were defined
847 as a binary distinction between face and house.

848 For the **ThingsEEG** dataset, we used the preprocessed data released by the original authors and
849 defined a visual stimulus classification task. For each subject, we loaded the test set from the original
850 Things-EEG dataset, which contained images from 200 distinct categories, with each image presented
851 80 times. This resulted in a recorded signal of shape [200, 80, 17, 100], where 17 was the number of
852 channels and 100 the number of time steps. These 17 channels were selected from the occipital and
853 parietal cortex (O1, Oz, O2, PO7, PO3, POz, PO4, PO8, P7, P5, P3, P1, Pz, P2, P4, P6, P8) to focus
854 on regions most relevant for visual processing. We assigned labels to the corresponding trials and
855 then reshaped the data into [16000, 17, 100] (i.e., 200 categories \times 80 trials), with the corresponding
856 labels reshaped into [16000]. The data were then randomly shuffled and split into training, validation,
857 and test sets. Specifically, we used 20% of the data as the test set, while the remaining 80% were
858 divided into three folds; in each round, one fold was used for validation and the other two for training.
859 The model input was the EEG segment with shape [17, 100], and the target was a label from one of
860 the 200 image categories. We conducted this decoding task separately for each subject.

861 For the **SEED-DV** dataset, we adopted the first benchmark task from the original SEED-DV study:
862 40-class classification of fine-grained video concepts, to evaluate dynamic visual stimulus decoding.
863 Following the preprocessing protocol described in the original work, we downsampled the EEG
signals to 200 Hz, with each trial spanning 2 seconds and containing 62 channels. The classification

864 labels correspond to the 40 predefined video concepts. To standardize our evaluation pipeline, we
 865 merged data across all blocks for each subject, randomly shuffled the trials, and partitioned them into
 866 training, validation, and test sets.

867 For the **OpenMIIR** dataset, we focused on two tasks: perception (denoted as OpenMIIR-P) and
 868 imagination (denoted as OpenMIIR-I). Following the preprocessing protocol of the original work, we
 869 used raw EEG signals sampled at 512 Hz without further downsampling. The data were band-pass
 870 filtered between 0.5–30 Hz, and trials were segmented based on audio onset events, with epochs
 871 defined for each musical stimulus under the corresponding condition. For the perception task,
 872 we used condition 1, while for the imagination task, we aggregated data from conditions 2, 3,
 873 and 4. After segmentation, the continuous recordings for each stimulus were further divided into
 874 non-overlapping windows of length 600 samples across 64 channels. We then applied channel-wise
 875 z-score normalization on each window. The classification labels were defined as the 12 distinct
 876 musical stimuli.

878 D DETAILS OF MODELS

880 D.1 MODELS IN THE ROADMAP EXPLORATION

882 Here, we provide details of the models used during the optimization process described in Section 2.

883 **Basic architecture study.** The model depth and embedding dimension during the basic architecture
 884 study are shown in Table 5. In both the CNN-1D and CNN-2D architectures, the convolution kernel
 885 size is set to 3. In CNN-1D, the number of channels increases from 1 to 64 in the stem layer,
 886 followed by 4 stages, each consisting of 5 layers. The number of channels in these stages increases
 887 progressively from 128 to 256, 512, and finally 1024. In CNN-2D, the channel configuration also
 888 starts from 1 to 64 in the stem layer. This is followed by 4 stages, each containing 4 layers, with
 889 the number of channels increasing from 96 to 192, 384, and 768, respectively. In the CNN-GRU
 890 architecture, the ratio of convolutional layers to GRU layers is set to 3:1. In the CNN-Transformer
 891 architecture, the initial ratio of convolutional layers to Transformer modules is set to 4:1. To fairly
 892 investigate how short- and long-term temporal information impacts model performance, we fix the
 893 total number of layers at 20 and systematically adjust the proportion of Transformer modules. To
 894 systematically evaluate the effectiveness of the proposed patching method, we vary only the data
 895 segmentation approach while keeping the embedding and subsequent self-attention-based feature
 896 extraction pipeline unchanged. All other hyperparameters are fixed to ensure that the observed
 897 performance differences can be attributed to the effectiveness of the patching method.

898 Table 5: **Model configuration for basic architecture study.**

900 Model	901 # Params (M)	902 Number of Layers	903 Hidden Dimension
904 CNN-1D	30	20	1024
905 CNN-2D	35	16	768
906 GRU	25	4	512
907 Transformer	20	8	512
908 PatchTST	20	8	512
909 iTransformer	20	8	512
910 CNN-GRU	34	16	512
911 CNN-Transformer	32	20	768
912 CNN-Transformer v2	36	20	1024

913 **Macro optimization.** From a macro perspective, we optimize the transformation of the latent space
 914 during forward propagation. Specifically, we investigate two strategies for increasing the number of
 915 feature maps: the step approach and the leap approach. In the step approach, the number of feature
 916 maps is gradually increased: starting from 1 to 64 through the stem layers, then progressively to 96,
 917 192, 384, and 768 across 4 stages, each containing 4 layers. In contrast, the leap approach increases
 the number of feature maps from 1 to 64 in the stem layer, and then directly increases it to 384
 using an embedding layer. Subsequently, 16 layers are used to refine high-level features. We also
 investigate two strategies for decreasing the size of feature maps: the pyramid approach and the

pagoda approach. In the pyramid approach, a downsampling layer is placed in each of the last three stages, with a downsampling ratio of 2 at each stage. In contrast, the pagoda approach places three identical downsampling layers consecutively within the second stage.

Micro optimization. From a micro perspective, we optimize the computation method. Specifically, we investigate two variants of vanilla convolution: the group convolution and the separable convolution. For the group convolution, the number of groups is set to 4. For the separable convolution, the number of groups is set to 1, followed by a pointwise convolution with a kernel size of 1 to aggregate channel-wise information.

D.2 NEUROSKETCH IMPLEMENTATION

Here, we provide a detailed description of the final NeuroSketch architecture to ensure transparency and reproducibility. Given an input tensor \mathbf{X} of shape $\mathbb{R}^{B \times C \times L}$, where B is the batch size, C is the number of channels, and L is the temporal length, we first reshape it to $\mathbb{R}^{B \times 1 \times 3C \times L // 3}$ to form a 2D representation. The reshaping operation aligns with the subsequent convolutional kernel sizes and ensures that the temporal structure of the data remains intact in the first stem layer. The reshaped 2D representation is then passed through a stem stage consisting of four Conv2D–BatchNorm–ReLU blocks with kernel sizes [3, 3, 3, 3], paddings [1, 1, 1, 1], and strides [2, 1, 1, 2]. The input/output channels for the four blocks are [1 → 64], [64 → 256], [256 → 64], [64 → 96].

Following the stem stage, there are four feature extraction stages. The input/output channel dimensions for these stages are [96 → $D_{\text{stage 1}}$], [$D_{\text{stage 1}}$ → $D_{\text{stage 2}}$], [$D_{\text{stage 2}}$ → $D_{\text{stage 3}}$], [$D_{\text{stage 3}}$ → $D_{\text{stage 4}}$]. Each stage contains d blocks with two key components:

- Patch Embedding:** In the first block of each stage, a 2D convolution with kernel size 3 and stride 1 maps the input channels to the output channels, followed by batch normalization. For the first three blocks of the second stage, we use a stride of 2 to downsample the input; in other cases, the layer reduces to an identity mapping when no change in resolution is needed.
- Convolution Module:** This component applies grouped 3×3 convolutions (with the number of groups equal to G) to efficiently capture local channel-wise dependencies. It is followed by batch normalization, ReLU activation, and a 1×1 convolution for feature fusion. The output is added back to the patch-embedded input via a residual connection.

After the four stage of feature extraction, we obtain a tensor of shape $\mathbb{R}^{B \times D_{\text{stage 4}} \times C' \times T'}$, where C' and T' denote the downsampled channel and temporal dimensions. We then apply GeM pooling to aggregate it into a representation $\mathbb{R}^{B \times D_{\text{stage 4}}}$, which is finally passed through a linear layer to produce the class probabilities. We implement two variants of the architecture: NeuroSketch-Base and NeuroSketch-Large. Their configurations are summarized in Table 6.

Table 6: **Detailed model configuration for NeuroSketch-Base and NeuroSketch-Large.**

	NeuroSketch-Base	NeuroSketch-Large
# Params (M)	1.4	4.2
$D_{\text{stage 1}}$	96	96
$D_{\text{stage 2}}$	128	144
$D_{\text{stage 3}}$	160	256
$D_{\text{stage 4}}$	192	384
d	2	3
G	4	4

D.3 BASELINE MODELS

Here, we introduce the details of the baselines for performance evaluation in Section 3.

ConvFormer and **CAFormer** (Yu et al., 2023) are two variants of MetaFormer (Yu et al., 2022), employing different token mixers. ConvFormer uses depthwise separable convolution as its token

972 mixer. With this design, ConvFormer can be regarded as a pure CNN model that does not rely on
973 channel or spatial attention mechanisms. CAFormer uses depthwise separable convolution as the
974 token mixer in the first two stages of the model and self-attention in the last two stages. This design
975 enables CAFormer to capture local features while better obtaining long-range dependencies.

976 **ModernTCN** (Luo and Wang, 2024) is a modernized purely convolutional architecture. It en-
977 hances the traditional TCN by incorporating depthwise convolutions and convolutional feed-forward
978 networks (ConvFFNs) into the 1D CNN design. Additionally, it introduces time series-specific
979 adaptations, such as patchified variable-independent embeddings, to improve suitability for time
980 series tasks.

981 **MedFormer** (Wang et al., 2024) is a multi-granularity patching Transformer specifically designed for
982 medical time series classification. It effectively captures the features of medical time series through
983 cross-channel patching, multi-granularity embedding, and a two-stage multi-granularity self-attention
984 mechanism.

985 **DeepConvNet** (Schirrmester et al., 2017) is a deep CNN tailored for EEG decoding. It begins with
986 a temporal convolution to capture frequency-specific patterns, followed by a spatial convolution
987 across channels to model inter-channel dependencies. A stack of convolution–batch normaliza-
988 tion–ELU–max pooling blocks with dropout then learns hierarchical spatio-temporal representations,
989 and a final dense layer performs classification. This end-to-end design avoids handcrafted features
990 and is effective across BCI tasks.

991 **EEGNet** (Lawhern et al., 2018) is a compact CNN tailored for EEG-based BCI. It decouples
992 frequency and spatial filtering by using temporal convolutions as learnable band-pass filters followed
993 by depthwise spatial convolutions across channels, and employs separable (pointwise) convolutions
994 for efficient feature mixing. EEGNet attains strong performance across diverse BCI paradigms
995 while using fewer than 0.5 million parameters, making it data-efficient and well suited to small EEG
996 datasets.

997 **ConFormer** (Song et al., 2022) is a CNN-Transformer model for EEG signals. The convolution
998 module learns low-level local features through one-dimensional temporal and spatial convolution
999 layers, and the self-attention module processes the output of the convolution module to learn global
1000 temporal dependencies.

1001 **SPaRCNet** (Jing et al., 2023) is a 1D CNN for seizure pattern recognition. Its structure features
1002 dense and transition blocks. Each dense block has four layers of two convolutional layers and two
1003 ELUs, and each transition block contains an ELU, a convolutional layer, and an average pooling
1004 layer.

1005 **seegnificant** (Mentzelopoulos et al., 2024) is a brain foundation model for cross-subject neural
1006 decoding from SEEG. It tokenizes electrode-wise signals using convolutions, models long-term
1007 temporal dependencies with self-attention, and integrates electrode 3D spatial locations through
1008 positional encoding followed by cross-electrode attention. A unified backbone extracts global neural
1009 representations, while subject-specific task heads enable individualized decoding. Trained on multi-
1010 session SEEG data from 21 participants, seegnificant demonstrates effective behavioral response time
1011 decoding and supports few-shot transfer to new subjects, offering a path toward robust multi-subject
1012 generalization in SEEG analysis.

1013 **CBraMod** (Wang et al., 2025b) is a brain foundation model for EEG decoding. The model first
1014 segments EEG signals into patches and randomly masks them. After patch encoding and asymmetric
1015 conditional position encoding, it learns spatio-temporal dependencies through the parallel spatial and
1016 temporal attention mechanisms of the cross-Transformer. Finally, it uses the reconstruction head to
1017 reconstruct the masked EEG patches, thereby learning the general representation of EEG signals.

1018 1019 1020 1021 1022 E EXPERIMENTAL SETTINGS

1023
1024
1025 Here, we present the training strategies employed for the models described in Section 2 and Section 3.

E.1 DATA SPLITTING

For Section 2, each subject’s data is split into training, validation, and test sets with a 3:1:1 ratio. For Section 3, we allocate 20% of the data for testing, and the remaining 80% is used for a 3-fold cross-validation. Specifically, the remaining data are divided into three equal parts, with each iteration using one part as the validation set while the other two are combined for training, ensuring a comprehensive and reliable assessment of the model’s effectiveness.

E.2 DATA AUGMENTATION

To further enhance data diversity, we employ a strong data augmentation strategy comprising the following techniques:

Random Shift. We define a maximum shift range proportional to the input sequence length. For each training instance, a shift step is randomly sampled from this range. A positive value indicates a forward shift of the sequence, whereas a negative value corresponds to a backward shift. This approach enhances the model’s robustness to uncertainty in stimulus onset times.

Noise. We generate noise from a standard normal distribution matching the shape of the original data. The noise is scaled by a predefined standard deviation and added to the input to produce noisy data. This method improves the model’s resilience to signal perturbations.

Channel Masking. During training, a mask is applied to each channel with the specified probability, zeroing out the corresponding channel values. This technique reduces over-reliance on specific channels and promotes better integration of multi-channel information.

Time Masking. Similar to channel masking, we apply masks along the temporal dimension. This encourages the model to extract features robustly across various time segments and enhances generalization.

Mixup. A mixing coefficient λ is sampled from a Beta distribution parameterized by a hyperparameter $\alpha = 0.4$. The original sample and a randomly chosen sample are linearly combined using λ , and their corresponding labels are mixed accordingly. This augmentation method exposes the model to a broader range of data combinations, thereby improving generalization.

E.3 HYPERPARAMETERS

Table 7: **Training hyperparameters.**

Hyperparameters	Settings
Epoch	100(Chisco), 500(others)
Batch size	64
Seed	42
Optimizer	AdamW
Learning rate	1e-3
Weight decay	5e-2
Scheduler	Cosine
Warmup ratio	0.1
Early stop ratio	0.2
Random shift probability	0.5
Random shift ratio	0.2
Add noise probability	0.1
Channel masking probability	0.5
Time masking probability	0.5
Mixup probability	0.5

Detailed hyperparameters are shown in Table 7. Across multiple datasets, we observed that decoding performance typically improves with an increased number of training epochs. Based on this observa-

tion, we set the number of training epochs to 500 in most experiments to fully optimize the model performance. However, for the *Chisco* dataset, the accuracy reaches a plateau at around 80 epochs and shows no further improvement. Therefore, we limit training on *Chisco* to 100 epochs.

E.4 COMPUTE RESOURCES

We utilized an AMD EPYC 7663 56-core processor and eight NVIDIA A100 GPUs, each with 80 GB of memory. Section 2 outlines 21 distinct models, each evaluated on 84 experiments that together span the entire set of subjects across eight decoding tasks. In total, this results in 1,764 experiments. Each experiment is executed on an A100 GPU and requires, on average, one hour of training time. Section 3 evaluates 12 distinct models, each trained with three cross-validation folds, resulting in a total of 3,024 experiments. On average, each experiment requires approximately 50 minutes of training time.

F EXPERIMENTAL ANALYSIS

F.1 NEUROPHYSIOLOGICAL INTERPRETABILITY

To further investigate the spatiotemporal features captured by NeuroSketch, we conducted a case study on the Du-IN dataset. In this analysis, we trained NeuroSketch-Large on all 115 channels from subject 02 and 90 channels from subject 11, and applied the Score-CAM (Wang et al., 2020) method to visualize the model’s decision-making process across both spatial and temporal dimensions.

Score-CAM generates class-activation maps by computing the gradient-free importance of each feature map and projecting it back to the input space, thereby highlighting which spatiotemporal regions most strongly influence the predicted class. For each sample, we obtained a saliency map and subsequently averaged these maps across all samples from subject 02 and subject 11, respectively, to produce subject-wise saliency maps. These aggregated maps, as illustrated in Figure 4, were then used to identify regions of interest (ROIs) within both the spatial and temporal domains.

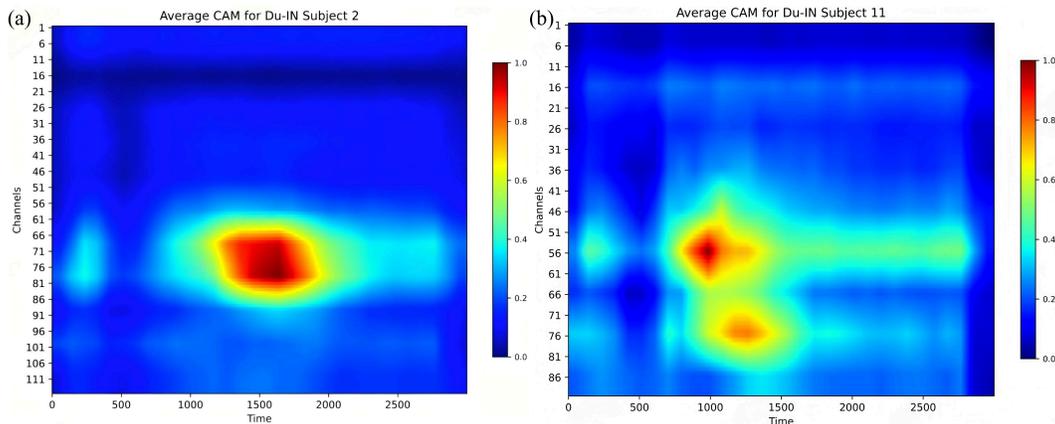


Figure 4: Score-CAM visualization of all channels for Subjects 02 and 11 in the Du-IN dataset generated by NeuroSketch-Large.

For spatial analysis, we averaged each subject-wise saliency map across the temporal dimension to obtain a saliency score for every channel. Notably, for subject 11, the channels previously reported as significant in the original study (Zheng et al., 2024) were 52,55,56,57,65,74,75,76,77,78. Except for channel 65, all nine of these channels appeared within our top-20 most salient channels identified by NeuroSketch. Similarly, for subject 02, the significant channels reported in Zheng et al. (2024) were 72,73,74,75,76,77,100,109,110,111, among which channels 72–77 were also ranked within our top-20. Beyond these specific electrodes, the remaining highly salient channels detected by our model were found to be spatially clustered around these previously reported regions, indicating that NeuroSketch captures spatial activity patterns that are consistent with experimentally validated cortical regions.

More importantly, the top-contributing electrodes identified by NeuroSketch are located within or near the ventral sensorimotor cortex (vSMC) and the bilateral superior temporal gyrus (STG)—two cortical regions well established as the core network for speech motor control (Bouchard et al., 2013; Hickok and Poeppel, 2007; Chartier et al., 2018). The vSMC is primarily responsible for the motor coordination of articulatory movements, while the STG supports auditory feedback processing and vocal self-monitoring during speech production. The focus of the model on these regions highlights its ability to capture biologically interpretable and functionally grounded spatial activation patterns that align with established speech-related cortical circuits, further validating the neurophysiological relevance of our saliency findings.

For temporal analysis, the model exhibits a pronounced activation hotspot concentrated around the mid-trial period, approximately between 800-1600 ms for subject 11 and 1200-1800 ms for subject 02, followed by a gradual decrease in saliency toward both the early and late segments. This pattern indicates that NeuroSketch focuses on short-range dependencies rather than allocating attention uniformly over time, which aligns with the transient temporal dynamics of neural signals. Collectively, these results confirm the model’s ability to extract interpretable and neurophysiologically grounded spatiotemporal features.

F.2 REPRESENTATION ANALYSIS OF CNN-2D AND TRANSFORMER-BASED ARCHITECTURES

Given the observed performance gap between CNN-based and Transformer-based architectures from an empirical perspective in Section 2.2 and Section 3.2, we further analyzed their representations to uncover the underlying reasons behind this gap.

Specifically, we examined the model representations from a rank-based perspective, following the framework proposed by Yu et al. (2025). The rank of a feature representation reflects its expressive capacity: a higher rank indicates richer and more diverse features, while a lower rank implies excessive compression and loss of independent information. In a well-organized hierarchy, the rank is expected to remain high or increase with depth, as deeper layers capture more abstract yet independent features. In our evaluation, since the CNN-2D-based model and the Transformer-based model achieved the overall best and worst performance, respectively, we conducted an analysis on NeuroSketch-Large (CNN-2D-based) and MedFormer (Transformer-based) to examine how the rank of the feature representations produced by each layer evolves across network depth.

For any layer of a given model, let the output representation of a sample be denoted as $\mathbf{O} \in \mathbb{R}^{s \times d}$, where s represents the sequence length and d the feature dimension. We performed singular value decomposition (SVD) on \mathbf{O} as $\mathbf{O} = \mathbf{U}\Sigma\mathbf{V}^T$, where \mathbf{U} and \mathbf{V} are orthogonal matrices, and Σ is a diagonal matrix whose diagonal entries $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{s,d\}}$ are the singular values. Although the algebraic rank—i.e., the number of nonzero singular values—serves as a strict measure of rank, real-world data are often noisy. Therefore, we adopt a more practical numerical rank. For a tolerance $\varepsilon > 0$, the ε -rank of \mathbf{O} is defined as the number of singular values that are significant relative to the largest one:

$$\varepsilon\text{-rank}(\mathbf{O}) = \left| \left\{ i \mid \frac{\sigma_i(\mathbf{O})}{\sigma_1(\mathbf{O})} > \varepsilon \right\} \right|, \quad (1)$$

Since the maximum rank (defined as $\min(s, d)$) varies across layers depending on s and d , we further compute the ε -rank ratio, defined as the ε -rank relative to the maximum rank. This metric allows for a more intuitive comparison of how the effective rank evolves across layers.

$$\varepsilon\text{-rank ratio} = \frac{\varepsilon\text{-rank}(\mathbf{O})}{\min(s, d)}, \quad (2)$$

For MedFormer, we extracted the ε -rank ratios of the embedding layer and each of its six encoder layers under three different patch lengths (5, 10 and 20) to examine how patch length affects representational rank. For NeuroSketch-Large, we analyzed the outputs of its three stem layers and all layers across the four feature-extraction stages, transforming the height–width dimensions of each feature map into a single sequence length for a consistent rank analysis.

As shown in Figure 5, from the two models’ evolution of the ε -rank ratios, we can observe that (1) For multi-channel EEG signals, the embeddings are not as low-rank as those typically observed in conventional time-series (Liang et al., 2025; Yu et al., 2025). Specifically, for both MedFormer and

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

Table 8: **Scaling behavior of NeuroSketch on the Du-IN dataset.** We report the mean accuracy together with the standard variance, computed over three distinct cross-validation folds. Additionally, we indicate the relative improvement in accuracy of NeuroSketch-Large compared to NeuroSketch-Base.

Subject	NeuroSketch-Base	NeuroSketch-Large	Relative Improvement
1	0.626 \pm 0.018	0.791 \pm 0.027	26.37%
2	0.761 \pm 0.025	0.798 \pm 0.007	4.95%
3	0.081 \pm 0.025	0.508 \pm 0.016	529.34%
4	0.422 \pm 0.052	0.765 \pm 0.033	81.36%
5	0.774 \pm 0.012	0.882 \pm 0.013	13.90%
6	0.209 \pm 0.021	0.465 \pm 0.016	122.49%
7	0.409 \pm 0.019	0.561 \pm 0.021	37.28%
8	0.525 \pm 0.012	0.630 \pm 0.021	20.00%
9	0.540 \pm 0.030	0.817 \pm 0.005	51.36%
10	0.034 \pm 0.008	0.183 \pm 0.003	433.01%
11	0.721 \pm 0.013	0.758 \pm 0.011	5.09%
12	0.569 \pm 0.050	0.668 \pm 0.010	17.39%

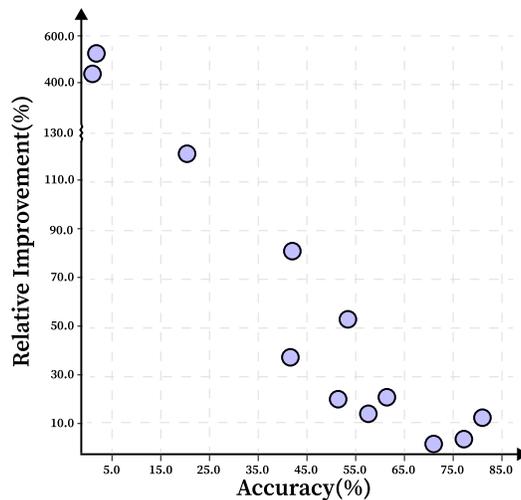


Figure 6: Scatter plot of the scaling behavior of NeuroSketch on the Du-IN dataset. The x axis denotes the accuracy of NeuroSketch-Base for each subject and the y axis represents the relative improvement achieved by NeuroSketch-Large.

Large, we adopt a kernel size of 3 as the default setting, since smaller kernels are generally more effective in capturing fine-grained dynamics while keeping the parameter overhead low. To further examine the effect of kernel size, we additionally evaluate larger kernels of 5 and 7 on the Du-IN dataset. The results are summarized in Table 9.

Table 9: **Results of different kernel sizes of NeuroSketch on the Du-IN dataset.** Results are reported as mean and standard deviation, computed across three distinct cross-validation folds.

Kernel Size	Accuracy	Precision	F1 Score
3	0.651 ± 0.005	0.667 ± 0.003	0.647 ± 0.003
5	0.557 ± 0.005	0.577 ± 0.006	0.550 ± 0.005
7	0.427 ± 0.002	0.456 ± 0.002	0.420 ± 0.003

These results indicate that larger kernel sizes lead to lower decoding accuracy, which aligns with our claim in Section 1 that neural decoding signals exhibit transient temporal dynamics. Larger kernels tend to capture longer-range temporal features, which may dilute short-term patterns that are critical for accurate decoding in this context.

Number of groups. The number of groups in group convolution determines how the feature channels are partitioned and processed, which in turn affects both the computational cost and the representational capacity. Increasing the number of groups reduces the computational cost, since fewer channels are convolved together within each group. In NeuroSketch-Large, we set the default group number to 4. To further examine the effect of this parameter, we additionally evaluate group numbers of 2, 8, and 16 on the Du-IN dataset. The results are summarized in Table 10.

Table 10: **Results of NeuroSketch under different number of groups in the group convolution on the Du-IN dataset.** Results are reported as mean and standard deviation, computed across three distinct cross-validation folds.

#Group	Accuracy	Precision	F1 Score
2	0.651 ± 0.005	0.667 ± 0.003	0.647 ± 0.003
4	0.652 ± 0.002	0.666 ± 0.003	0.648 ± 0.002
8	0.642 ± 0.002	0.658 ± 0.003	0.638 ± 0.002
16	0.642 ± 0.001	0.656 ± 0.003	0.637 ± 0.002

The results indicate that setting the group number to 4 provides the most favorable configuration. Compared with using 2 groups, this setting achieves competitive decoding performance while reducing the computational cost. At the same time, it clearly outperforms larger group numbers of 8 and 16 in terms of decoding accuracy, highlighting that 4 groups strike an effective balance between efficiency and representational capacity.