



HAL
open science

Characterizing Dynamic Functional Connectivity Subnetwork Contributions in Narrative Classification with Shapley Values

Aurora Rossi, Yanis Aeschlimann, Emanuele Natale, Samuel
Deslauriers-Gauthier, Peter Ford Dominey

► **To cite this version:**

Aurora Rossi, Yanis Aeschlimann, Emanuele Natale, Samuel Deslauriers-Gauthier, Peter Ford Dominey. Characterizing Dynamic Functional Connectivity Subnetwork Contributions in Narrative Classification with Shapley Values. 2024. hal-04596845v2

HAL Id: hal-04596845

<https://hal.science/hal-04596845v2>

Preprint submitted on 17 Dec 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Characterizing Dynamic Functional Connectivity Subnetwork Contributions in Narrative Classification with Shapley Values

Subnetwork Contributions in Narratives

Aurora Rossi^{1,*} Yanis Aeschlimann²

Emanuele Natale¹ Samuel Deslauriers-Gauthier^{2†} Peter Ford Dominey^{3,4,†}

¹ COATI, Université Côte d'Azur, INRIA, CNRS, I3S, Sophia Antipolis, France

² CRONOS, Inria Centre at Université Côte d'Azur, Sophia Antipolis, France

³ INSERM UMR1093-CAPS, Université Bourgogne Franche-Comté,

UFR des Sciences du Sport, Dijon, France

⁴ Robot Cognition Laboratory, Marey Institute Dijon, France

*Corresponding author: aurora.rossi@inria.fr

†Equal contribution

November 6, 2024

Abstract

Functional connectivity derived from functional Magnetic Resonance Imaging (fMRI) data has been increasingly used to study brain activity. In this study, we model brain dynamic functional connectivity during narrative tasks as a temporal brain network and employ a machine learning model to classify in a supervised setting the modality (audio, movie), the content (airport, restaurant situations) of narratives, and both combined. Leveraging Shapley values, we analyze subnetwork contributions within Yeo parcellations (7- and 17-subnetworks) to explore their involvement in narrative modality and comprehension. This work represents the first application of this approach to functional aspects of the brain, validated by existing literature, and provides novel insights at the whole-brain level. Our findings suggest that schematic representations in narratives may not depend solely on pre-existing knowledge of the top-down process to guide perception and understanding, but may also emerge from a bottom-up process driven by the ventral attention subnetwork.

Keywords fMRI; dynamic functional connectivity; narratives; shapley values; machine learning; convolutional neural networks

Acknowledgements A.R. and E.N. would like to thank Pierluigi Crescenzi for the discussion on the explainability technique used in the paper. This work has been supported by the French government, through the UCA DS4H Investments in the Future project managed by the National Research Agency (ANR) with the reference number ANR-17-EURE-0004 and the ANR France Relance project.

1 Introduction

Understanding the principles of representation and computation in the human brain, and developing corresponding predictive models, remains one of the great open challenges in neuroscience. fMRI provides a rich window into the dynamics of the whole human brain with a certain level of spatial and temporal resolution. From the beginning, human language processing has been a target of investigation with fMRI [Price, 2012]. Experiments with words and sentences allowed the identification of language processing areas and networks at different levels of structure [Keller et al., 2001]. More recently, evidence has emerged that language processing involves even broader recruitment across the brain, which might be obscured by time averaging and thresholding [Aliko et al., 2023]. This is consistent with studies that revealed how language recruits an extended fronto-temporo-parietal semantic system beyond the classic perisylvian language network [Xu et al., 2005, Jouen et al., 2015, Binder and Desai, 2011]. This has been demonstrated in the processing of narrative, full stories, which produce wide recruitment of

34 whole brain networks for memory, visuospatial representation, and emotion [Xu et al., 2005, Jääskeläinen et al.,
35 2021, Silbert et al., 2014]. Thus, narrative processing is a privileged context for the investigation of brain functional
36 dynamics [Willems et al., 2020]. How can these functional dynamics be characterized? Analysis methods based
37 on time averaging and subtraction tend to ignore the contribution of brain systems whose activity is variable
38 and averaged out during thresholding. Functional connectivity analysis can be used to capture and characterize
39 these dynamic interactions of brain regions over time [Sizemore and Bassett, 2018, Preti et al., 2017]. Temporal
40 brain networks model the evolution of functional connectivity over time and thus have the desired properties
41 of capturing the full brain dynamics that may be lost in time averaging and thresholding. Here, we exploit the
42 representational richness of dynamic functional connectivity in temporal brain networks to characterize brain
43 dynamics during narrative processing using machine learning.

44 In particular, we propose a simple machine learning model to classify in a supervised setting fMRI data
45 collected during a narrative comprehension task. The model is mainly composed of a convolutional layer and
46 a multi-layer perceptron (MLP). It is trained to classify the modality of the narrative (audio or video), the
47 content of the narrative (airport or restaurant situations) and these two together in a four-class classification.
48 We use the model to investigate the importance of temporal dynamics in narrative processing and combined
49 with the powerful explainability technique of Shapley values we delve deeper into the model’s decision-making
50 process. Specifically, we quantify the subnetwork contributions in the classification of two different parcellation
51 methods (Yeo 7-subnetwork and 17-subnetwork) and this allows us to identify the most involved subnetworks in
52 the narrative processing task. Our work is the first to apply this approach to functional aspects of the brain,
53 validated by existing literature, and provides novel insights at the whole-brain level.

54 The results provide valuable insights, validated by existing research on narrative comprehension [Baldassano
55 et al., 2018, Simony et al., 2016], and contribute to a broader understanding of how we process narratives. Our
56 findings challenge the initial assumption that narrative comprehension relies solely on top-down activation of
57 scripts, where prior knowledge, experiences, and expectations solely guide interpretation [Dubin and Bycina, 1991].
58 The prominent role of the ventral attention subnetwork in content classification suggests a more nuanced model.
59 This network is associated with bottom-up attentional control, implying that narrative processing might involve
60 the assembly and integration of sensory information from the environment alongside top-down influences. This
61 possibility aligns with the notion that schematic representations may not solely be driven by top-down activation
62 but could be built upon bottom-up processing mediated by the ventral attention subnetwork [Vossel et al., 2014].

63 1.1 Related works

64 **Classification of tasks from fMRI data** Numerous studies have explored classifying tasks and subject
65 characteristics (such as age and sex) from functional brain connectivity data using fMRI, primarily aiming to
66 develop powerful architectures. Examples include the work by Kim et al. [2021], where they propose a Spatio-
67 Temporal Attention Graph Isomorphism Network model for high-accuracy prediction of 7 tasks (memory, social,
68 relational, motor, language, gambling, and emotion) alongside sex. Another approach by Kim et al. [2023] utilizes
69 a transformer to classify age, sex, and cognitive intelligence, with an integrated gradient technique for interpreting
70 sex classification results. The latter explainability technique is also employed in a parallel similar work by Ryali
71 et al. [2024], where they classify sex using a simpler spatio-temporal deep neural network. Other papers by Huang
72 et al. [2021] and Saeidi et al. [2022] use a deep learning model, mainly composed of a convolutional neural network
73 and a recurrent neural network, and a graph neural network, respectively, to classify the 7 tasks.

74 **Narratives classification** In contrast to the aforementioned papers, our work focuses on a more detailed
75 classification domain, specifically the classification of modalities (movie, story) and the thematic content of the
76 script (airport, restaurant). Baldassano et al. [2018] exemplify this approach, using a stochastic Hidden Markov
77 Model to classify, based on the activation of a selection of regions of interest (ROIs) in the default attention
78 networks, thematic content while also incorporating event alignment.

79 **Shapley values in brain networks** The use of Shapley values has become a popular approach to explain
80 the predictions of machine learning models. In neuroscience, for instance, Amoroso et al. [2023] classify three
81 conditions (Alzheimer’s disease, mild cognitive impairment, and healthy controls) based on brain structural
82 connectivity data from MRI scans. They then leverage Shapley values to identify the most influential "patch" for
83 classification. Another study by Kotter et al. utilizes Shapley ratings in macaque brain networks, employing a
84 graph theory approach to analyze these networks. Here, the number of strongly connected components within a
85 subgraph serves as the Shapley value function [Kötter, 2007]. The most similar work to ours is by Li et al. [2020].
86 They propose a new estimation method for Shapley values and apply it when classifying functional connectivity
87 data from fMRI. In their example, they classify patient conditions (autism spectrum disorder or healthy) and

88 compute the importance of different ROIs in classification, though they don't delve into the neuroscientific
89 interpretation of the results.

90 1.2 Our contribution

91 This study combines machine learning with explainable AI to investigate the specific roles of brain subnetworks
92 during tasks involving narratives. We leverage functional connectivity, extracted from fMRI data, and Shapley
93 values to identify which brain subnetworks are most influential in classifying narrative modality (audio and movie),
94 thematic content (airport and restaurant situation) and their combination. The fMRI data are segmented into 7
95 or 17 Yeo subnetworks using the Schaefer 100 element parcellation [Schaefer et al., 2018]. Our machine learning
96 model, composed of a convolutional neural network and multi-layer perceptron, achieves high accuracy and reveals
97 the specific contributions of Yeo subnetworks in narrative processing. Importantly, the focus of our analysis is
98 functional connectivity, rather than activation. This analysis, validated by neuroscientific interpretation aligned
99 with existing literature, offers new insights into the functional roles of these subnetworks and the factor of time
100 during narrative classification. Our work demonstrates the power of explainable AI in unveiling the complex
101 interplay between brain activity and narrative comprehension. It not only helps to understand narrative processing
102 but also paves the way for applying this approach to other areas of brain research.

103 2 Methods

104 2.1 Model

105 Our model takes as input a temporal brain network. This network is a sequence of brain networks, each reflecting
106 the brain's functional connectivity at a specific time step (further details regarding the data processing are
107 provided in the Experiments section). Mathematically, the temporal brain network can be represented as a
108 three-dimensional tensor, denoted by $X \in [-1, 1]^{R \times R \times T}$, where R represents the number of brain regions and T
109 represents the number of time steps. In our case, R is 100 and T is 8.

110 The model architecture consists of a single-layer three-dimensional convolutional neural network, followed
111 by a max pooling layer and a multi-layer perceptron for classification. The convolution filter has size (R, R, τ)
112 with no padding, where the two first dimensions match with those of the input. This design focuses on capturing
113 temporal features within the brain network by restricting filter movement to the temporal axis. Max pooling is
114 then applied to reduce the dimensionality of the extracted features. Finally, a multi-layer perceptron performs the
115 classification task. A visual representation of the model architecture is provided in Figure 1.

116 Notably, when the filter size in the temporal dimension is set to 1 ($\tau = 1$), the model becomes invariant to the
117 specific order of time steps in the input data. An analysis of the model's performance with different filter sizes is
118 provided in the Appendix section.

119 Formally, given an input tensor $X \in [-1, 1]^{R \times R \times T}$, the output of the convolutional layer is defined as

$$Y_{k,c} = \sigma(X * W + b)_{k,c} = \sigma\left(\sum_{i=1}^R \sum_{j=1}^R \sum_{p=1}^{\tau} X_{i,j,k+p-1} \cdot W_{i,j,p,c} + b_{k,c}\right)$$

120 where $Y \in \mathbb{R}^{K \times C}$ is the output tensor, $W \in \mathbb{R}^{R \times R \times \tau \times C}$ is the learnable filter tensor, $b \in \mathbb{R}^{K \times C}$ is the bias matrix
121 and C is the number of output channels. The operations \cdot , $+$ and σ , which represents the $\text{ReLU}(x) = \max\{0, x\}$
122 activation function, are applied component-wise. The output tensor is then passed through a max pooling layer
123 so that the output vector $Z \in \mathbb{R}^C$ is defined as

$$Z = \max_k Y[k, c].$$

124 Finally, the output passed through a multi-layer perceptron of three fully connected layers with ReLU activation
125 functions. A fully connected layer can be defined as $V = \sigma(W \cdot Z + b)$ where V is the output of the fully connected
126 layer, W is the weight matrix, and b is the bias vector.

127 2.2 Shapley Values

128 Shapley values were introduced by Lloyd Shapley in 1951 in the context of cooperative game theory [Shapley,
129 1951]. They quantify the contribution of each player in a coalition game. Recently, they have been adopted in
130 machine learning to explain the predictions of models. Shapley values can be calculated using different methods
131 including sampling or exact computation for smaller player sets [Lundberg and Lee, 2017]. In our case, we leverage

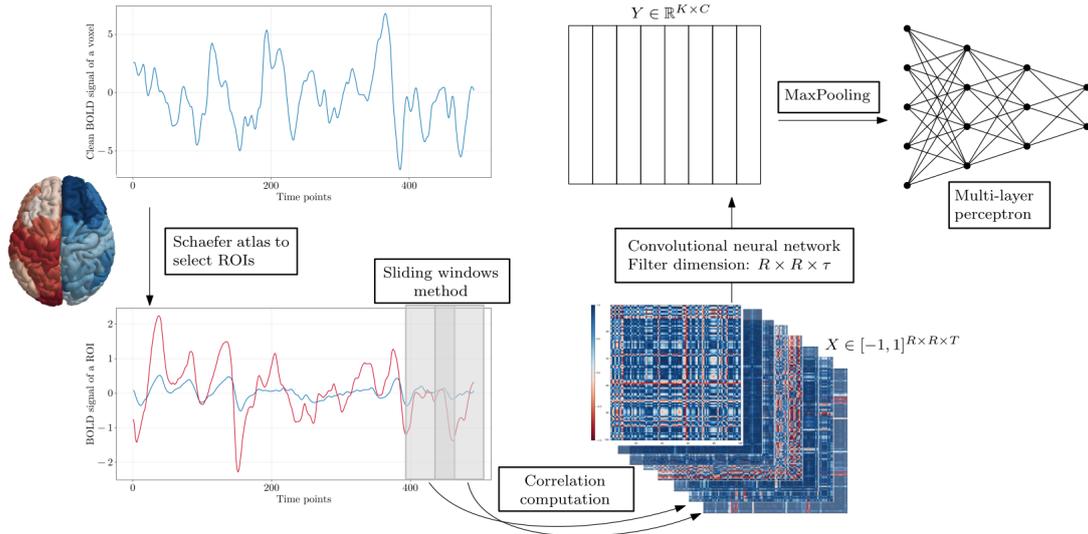


Figure 1: Pipeline from the extraction of temporal brain networks to the classification of the narrative aspects. The first step is the division of the brain into regions according to an atlas. The second step is the sliding window method, which individuates rectangular windows within which the Pearson correlation coefficient is computed between each pair of brain region time series. The output is then fed into the model, which consists of a convolutional layer, a max-pooling layer, and a multi-layer perceptron.

132 Shapley values to understand the influence of specific brain subnetworks on the prediction of our model. Because
 133 of the limited number of brain subnetworks defined by the 7 Yeo parcellation method [Thomas Yeo et al., 2011],
 134 we can compute the exact Shapley values. The exact Shapley value of a brain subnetwork i is defined as

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} (v(S \cup \{i\}) - v(S)) \quad (1)$$

135 where N is the set of brain subnetworks, v is the accuracy of our model when considering the set S of brain
 136 subnetworks. To isolate the brain subnetworks in the temporal brain network X we set the entries of the other
 137 subnetworks to zero. The Shapley value $\phi_i(v)$ is the average marginal contribution of the brain subnetwork i
 138 over all possible combinations of brain subnetworks, the higher the Shapley value, the more important the brain
 139 subnetwork is for the prediction of the model. For the 17 Yeo subnetwork parcellation, the exact computation of
 140 Shapley values becomes computationally expensive. Therefore, we employ a sampling method that approximates
 141 the Shapley values using the same formula but instead of summing over all possible subnetwork combinations,
 142 we sample a large number of combinations (100 samples in our case) to approximate the average marginal
 143 contribution.

144 2.3 Experiments

145 Experiments were performed to determine if the temporal brain networks can be used to discriminate brain
 146 functional connectivity patterns in response to audio vs. movie narratives, airport vs. restaurant situations, and
 147 the combination of these two dimensions. We trained a machine learning model in a supervised setting to classify
 148 these aspects and used Shapley values to interpret the model’s decisions.

149 2.3.1 Data

150 **Dataset** Our analysis used fMRI data from the study of Baldassano et al. [2018] archived as part of the
 151 Narratives dataset created by Nastase et al. (<https://openneuro.org/datasets/ds002345/versions/1.1.4>)
 152 [Nastase et al., 2020]. The Baldassano dataset includes brain activity recordings from 31 participants engaged in
 153 a narrative task. In this task, each subject is exposed to 16 3-minute stories (4 per run over 4 runs), from two
 154 different scripts (eating at a restaurant or going through the airport). While the stories within each category
 155 share a similar high-level sequence of events, there are variations in the specific details of these events. Each
 156 run presents 2 movies and 2 audio stories, for a total of 8 movies and 8 audio segments over the course of the
 157 experiment. The dataset is balanced in terms of the number of samples per modality and content.

158 **Preprocessing** The fMRI data has a spatial resolution of $91 \times 91 \times 109$ voxels in the x, y, and z axes, respectively,
159 for a total of 902,629 voxels. Each voxel measures $2 \times 2 \times 2$ mm. The repetition time is 1.5 seconds, for a total of
160 490 time points and a total duration of 12 minutes per run approximatively.

161 Preprocessing involved transforming the blood-oxygen-level-dependent (BOLD) signals from each voxel into
162 temporal graphs. We implemented a pipeline to reduce motion artifacts by performing linear regression on the
163 movement parameters. Additionally, a bandpass filter (0.01 – 0.08 Hz) was applied to remove noise arising from
164 respiration and cardiac pulsations [Van Dijk et al., 2010].

165 To define the network nodes, we employed the Schaefer et al. brain atlas (after having put the data in the
166 MNI152 space), parcellating the brain into 100 ROIs based on anatomical and functional criteria [Schaefer et al.,
167 2018]. ROIs were created by averaging the BOLD time series of voxels within gray matter regions. We then
168 utilized a sliding window approach with 30-second windows and 7.5 seconds overlap to divide the data into time
169 steps. The Pearson correlation coefficient was computed between each pair of ROI time series within each window,
170 with the resulting correlation value assigned as the weight of the edge connecting the corresponding ROI nodes.
171 This process yielded an adjacency matrix for each time window, and the sequence of these matrices formed the
172 temporal brain networks (see Figure 1).

173 2.3.2 Experimental setting

174 The experiments were conducted on a workstation equipped with a single NVIDIA Quadro RTX 8000 graphics
175 card. We utilized the Julia programming language for the workflow, from network creation starting from the
176 clean signal to the model development [Bezanson et al., 2017]. The Flux.jl library was used for neural network
177 implementation and the Makie.jl library was used for visualization [Innes et al., 2018, Danisch and Krumbiegel,
178 2021]. The source code is available at the following GitHub repository: [https://github.com/aurorarossi/
179 fMRINarrativeClassification](https://github.com/aurorarossi/fMRINarrativeClassification).

180 **Hyperparameters** The hyperparameters were chosen based on empirical observations. The convolutional filter
181 τ parameter was set to 4 for the modality classification task and 8 for the content and the combined classification
182 task (see the Appendix for more details). The number of output channels was set to 128 for all the tasks. The
183 MLP had two hidden layers with 64 and 32 units each with a ReLU activation function. The output dimension
184 of the MLP was set to 2 for the modality classification task, 2 for the content classification task, and 4 for the
185 combined classification task.

186 **Training** Given the limited size of the dataset, we employed a batch size of 1 during training. We used the Adam
187 optimizer with a learning rate of 0.0001. The training process lasted for 20 epochs. The choice of 20 epochs was
188 determined through experiments to achieve a good balance between training time and model performance. For the
189 loss function, we used either logit binary cross-entropy or logit cross-entropy depending on the number of classes
190 in the task. To ensure robustness against potential variations due to model initialization, we retrain the model 15
191 times with different random splits of the data (80% training, 20% testing). During each iteration, we compute
192 both the Shapley values and the model’s accuracy. Finally, we report the mean and standard deviation to account
193 for variability for the accuracy, and for Shapley values of each subnetwork, we present the mean values along with
194 error bars representing the standard deviation. This approach ensures a comprehensive understanding of the
195 model’s performance, the contribution of individual brain subnetworks to its classifications, and the robustness of
196 these findings across model initializations.

197 3 Results

198 In this section, we describe the results of our experiments. We present the performance of the model on three
199 classification tasks:

- 200 • **Modality classification:** this task focuses on classifying the brain network based on the modality of the
201 stimuli, audio or movie.
- 202 • **Content classification:** the model classifies the brain network based on the content of the stimuli, airport
203 or restaurant situations.
- 204 • **Combined Modality and Content Classification:** this task evaluates the model’s ability to jointly
205 classify both the modality and the content of the stimuli.

	Modality	Content	Both Modality and Content
Accuracy	96.32% \pm 1.36%	80.9% \pm 1.75%	80.70% \pm 2.97%
Precision	95.64% \pm 1.43%	84.55% \pm 2.29%	81.54% \pm 5.34%
Recall	97.08% \pm 2.20%	75.69% \pm 3.02%	80.70% \pm 5.23%
F1-Score	96.34% \pm 1.36%	79.84% \pm 1.96%	80.92% \pm 6.06%
Accuracy permuting network windows	86.60% \pm 3.36%	63.19% \pm 4.40%	53.12% \pm 5.85%
Accuracy permuting time series	50.76% \pm 2.74%	47.22% \pm 2.45%	26.18% \pm 2.64%
Static functional connectivity accuracy	86.11% \pm 2.38%	75.07% \pm 2.44%	62.71% \pm 3.18%

Table 1: Performance metrics of the model across modality, content, and combined classification tasks. The row ‘Accuracy permuting network windows’ reflects the model’s performance when brain network time steps are permuted, while ‘Accuracy permuting time series’ shows performance when the time series are permuted prior to constructing the network. The last row reports the model’s performance using static functional connectivity matrices.

206 The results in Table 1 show that the model performs well on the modality classification task, achieving an
207 accuracy of 96.32% \pm 1.36%. While still a good performance considering the complexity, the model’s accuracy
208 on the content classification task was slightly lower at 80.9% \pm 1.75%. This difference might be attributed to
209 the inherent difficulty of content classification compared to modality identification. Furthermore, the combined
210 modality and content classification task resulted in an accuracy of 80.70% \pm 2.97%, which is consistent with the
211 content classification task. Notably, the model displayed consistent performance across all metrics.

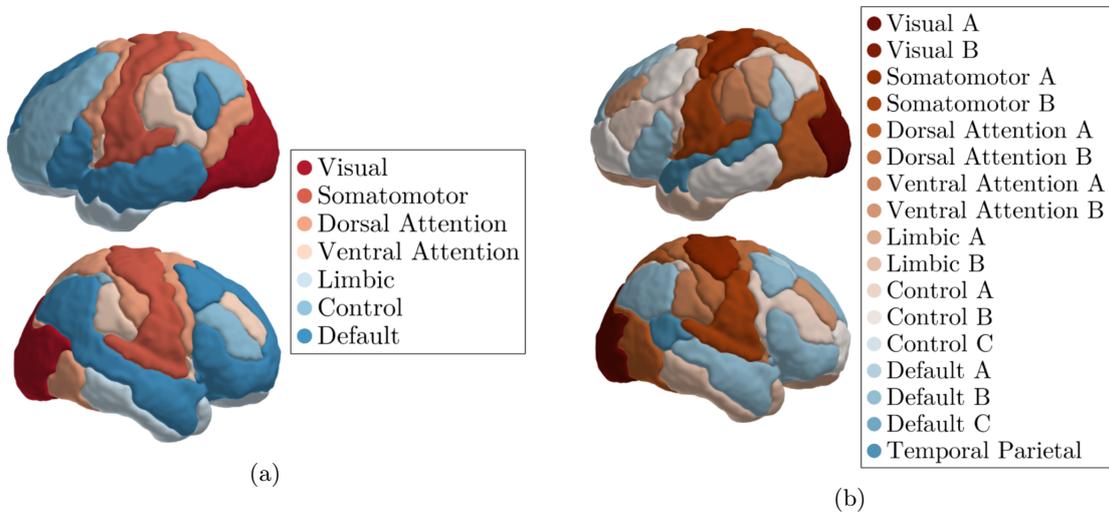


Figure 2: Yeo parcellations used in the Shapley value analysis. The 7-subnetwork parcellation is shown on the left (a), while the 17-subnetwork parcellation is shown on the right (b).

212 To assess the importance of the time dimension in classification tasks, we performed two types of permutation:
213 first, we shuffled the time series before constructing the network; second, we shuffled the network windows while
214 keeping the time steps within each window intact, followed by retraining the model. As expected, shuffling the
215 entire time series led to significantly lower accuracy compared to shuffling the network windows.

216 When shuffling the network windows, the time evolution within each window remains consistent with the
217 original data, essentially creating a block permutation. This means that while the order of windows is altered, the
218 temporal relationships within each window are preserved. In contrast, shuffling the entire time series disrupts
219 the sequential flow, completely dismantling its temporal structure. This disruption impacts both content and
220 modality classification, as both rely heavily on the temporal context of the brain activity being analyzed.

221 In the case of shuffling network windows, the results show a notable drop in accuracy compared to the unshuffled

222 data: 10% for modality classification, 17% for content classification, and a substantial 27% for the combined
 223 task. These drops indicate that the temporal dynamics within brain networks are crucial for all classification
 224 tasks and that the model effectively utilizes this information. The performance decrease is more pronounced in
 225 content and combined classification tasks compared to modality classification, which aligns with our expectations.
 226 Understanding content, which often unfolds over time and involves complex interactions between brain regions,
 227 likely depends more on temporal dynamics than modality identification alone.

228 For static functional connectivity matrices, the results remain relatively high: $86.11\% \pm 2.38\%$ for modality
 229 classification, $75.07\% \pm 2.44\%$ for content classification, and $62.71\% \pm 3.18\%$ for the combined task. While these
 230 results indicate that static features provide valuable information, the performance is notably lower compared to
 231 when the model incorporates dynamic time series data. This suggests that while static functional connectivity
 232 offers useful insights, integrating temporal information significantly improves the model’s ability to accurately
 233 classify brain activity.

234 To gain deeper insights into how the model leverages brain activity for classification, we employed Shapley
 235 values. Here, we focus on subnetworks defined by the Yeo parcellation method [Thomas Yeo et al., 2011],
 236 specifically the 7-subnetwork and 17-subnetwork parcellations. Visualizations of these parcellations are provided
 237 in Figure 2. Black and white compatible versions of these figures can be found in the Appendix.

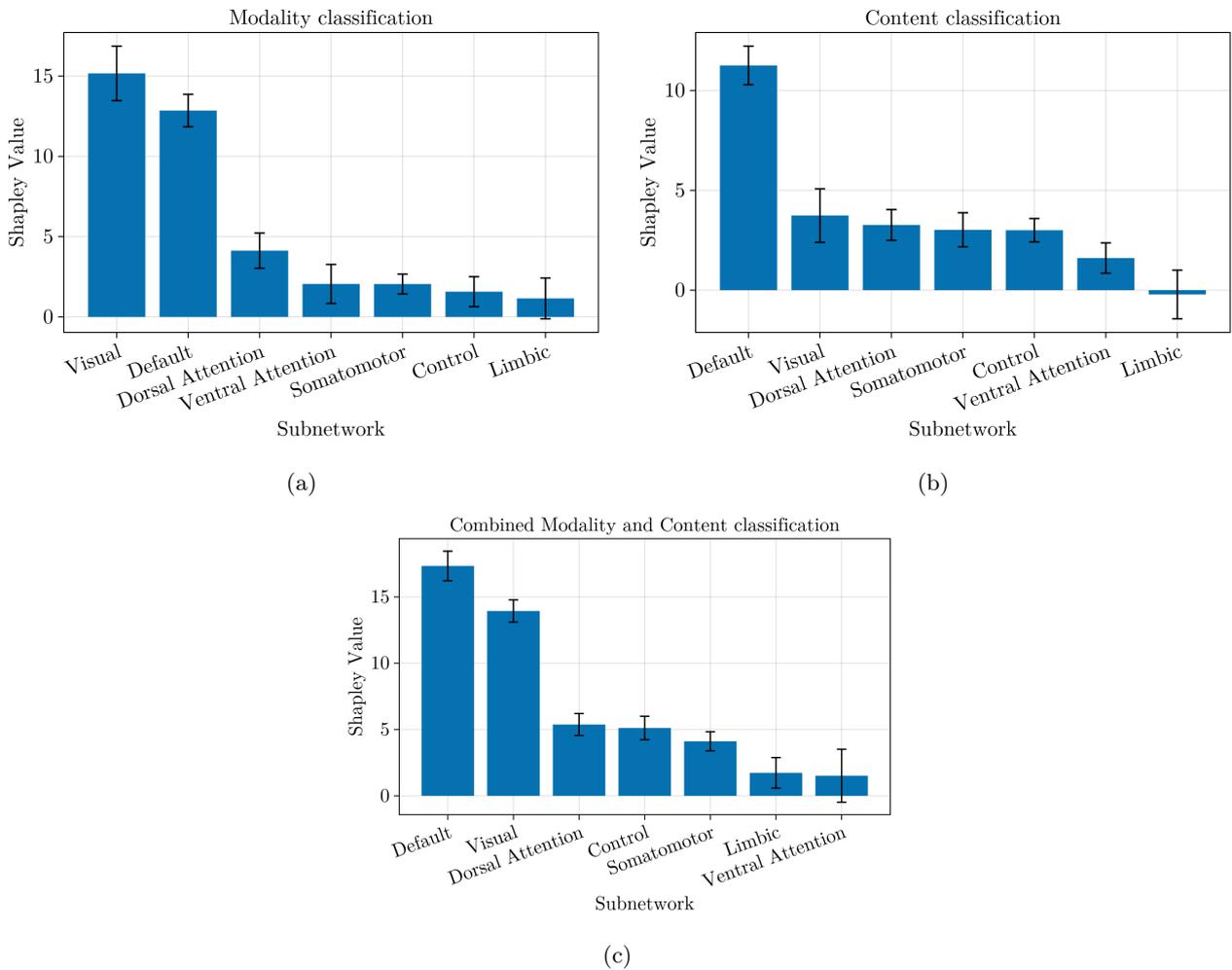


Figure 3: This figure shows the contribution of Yeo 7-subnetworks computed with Shapley values for classifying narrative using a machine learning model. The bars represent the average contribution of each subnetwork to the model’s predictions, with higher values indicating greater influence. The error bars represent the standard deviation of the Shapley values.

238 Figure 3 presents the Shapley values for the 7-subnetwork parcellation. In the modality classification task, the
 239 visual subnetwork emerges as the most influential, followed by the default mode subnetwork (Figure 3a). This
 240 aligns with the intuitive notion that processing visual information plays a key role in distinguishing modalities.
 241 For the content classification task, the high value of the default mode subnetwork suggests its influence in

242 understanding the meaning and content of the stimuli as suggested by previous studies [Baldassano et al., 2018,
 243 Simony et al., 2016] (Figure 3b). Finally, the combined classification task reveals the importance of both the
 244 visual and default mode networks (Figure 3c), suggesting that the model utilizes a combination of visual features
 245 and higher-order processing for accurate content and modality classification.

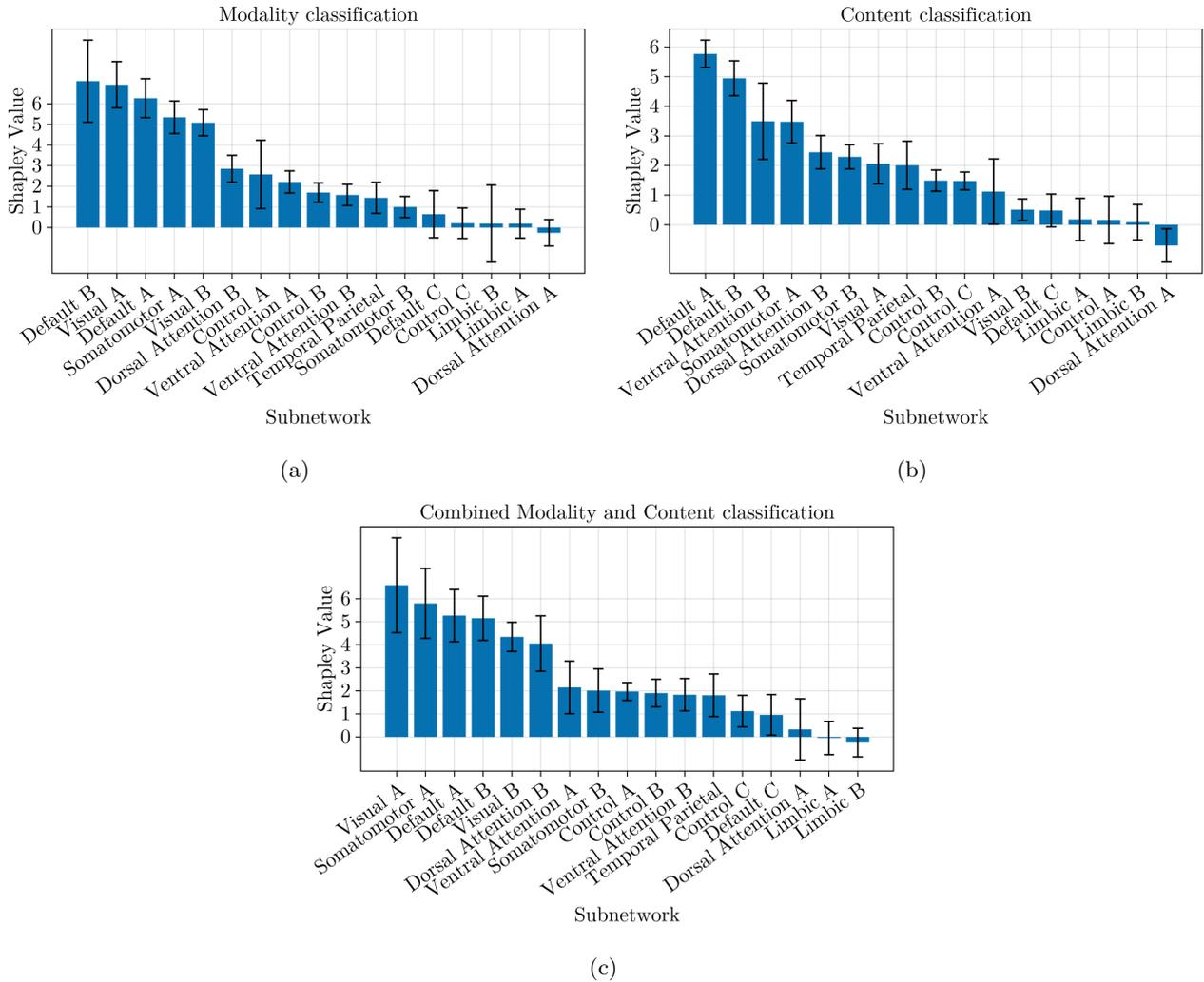


Figure 4: This figure shows the contribution of Yeo 7-subnetworks computed with Shapley values for classifying narrative using a machine learning model. The bars represent the average contribution of each subnetwork to the model’s predictions, with higher values indicating greater influence. The error bars represent the standard deviation of the Shapley values.

246 Figure 4 presents the Shapley values for the 17-subnetwork parcellation. In the modality classification task,
 247 the visual A and B, default A and B and somatomotor A subnetworks emerge as the most influential (Figure
 248 4a). For the content classification task, the default A and B subnetworks, the somatomotor A and the ventral
 249 attention B also play crucial roles (Figure 4b). Finally, the combined classification task reveals the importance of
 250 the visual A and B, default A and B, and somatomotor A subnetworks (Figure 4c).

251 Figure 5 shows the Shapley scores for the 100 parcellations of the Schaefer subnetworks, which are consistent
 252 with the results of the Yeo parcellations. The visual network emerges as the most significant for modality
 253 classification. For content classification, the default mode network is dominant, while for combined classification,
 254 both networks are most significant.

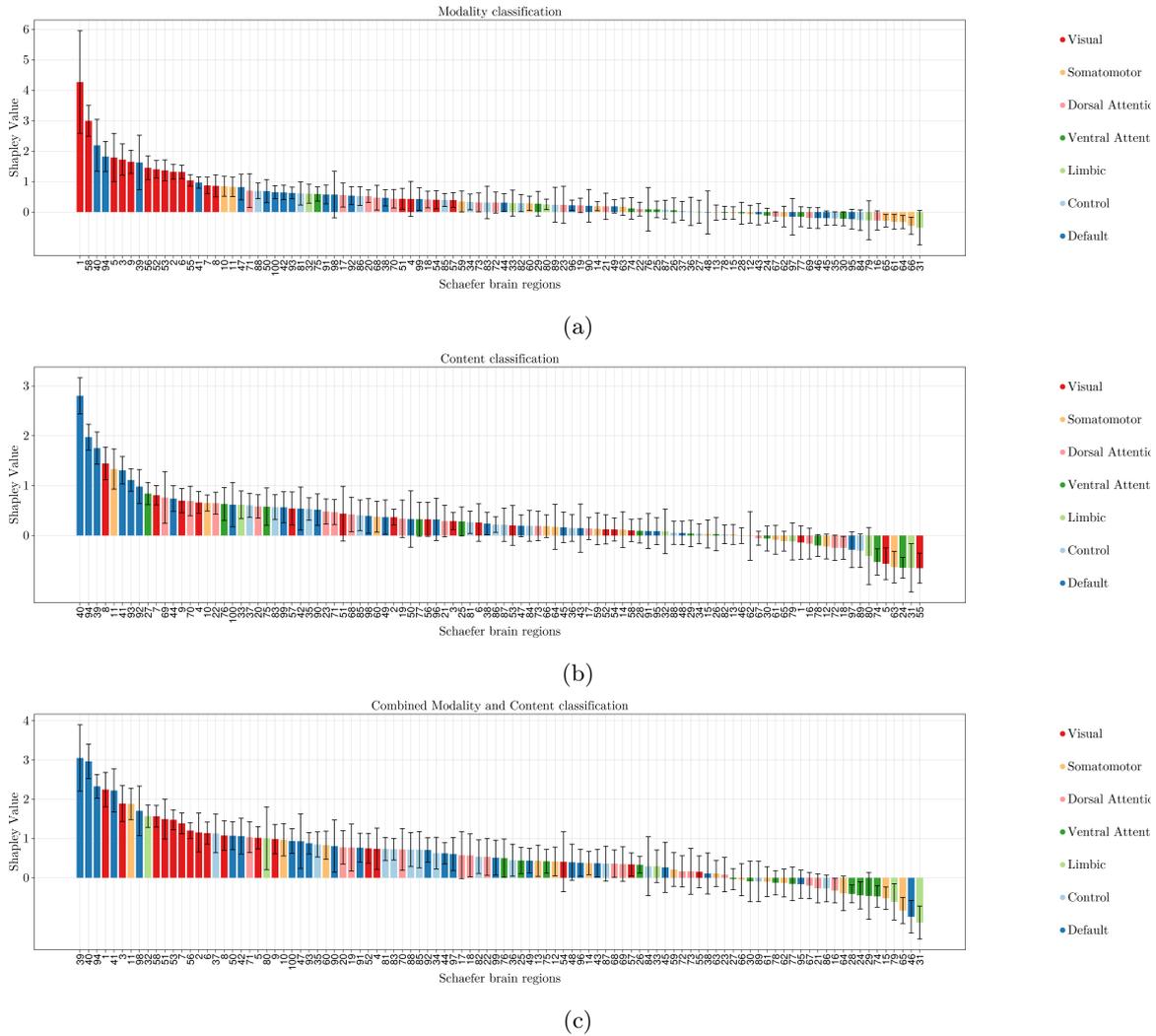


Figure 5: This figure shows the contribution of 100 Schaefer subnetworks, computed using Shapley values, for classifying narratives with our machine learning model. The bars represent the average contribution of each subnetwork to the model’s predictions, with higher values indicating greater influence. The error bars denote the standard deviation of the Shapley values. Additionally, the color of each bar corresponds to one of the 7 subnetworks in the Yeo 7-parcellation.

255 4 Discussion

256 This work investigated the neural basis of narrative processing using a machine learning model that classifies
 257 narrative aspects (modality, content, combined) based on functional connectivity networks derived from fMRI
 258 data. The model’s performance aligned with expectations: higher accuracy for modality classification, which
 259 is a simpler task because it relies on sensory information, compared to content classification which requires a
 260 deeper understanding of the narrative. Permuting time steps in the temporal brain network significantly reduced
 261 accuracy, particularly in content and combined tasks, suggesting that temporal dynamics rely on the sequence of
 262 events to understand the content.

263 To delve deeper into the model’s decision-making process, we employed Shapley values, a powerful explainable
 264 AI technique that quantifies subnetwork contributions. While techniques like Grad-CAM and Eigen-CAM provide
 265 valuable insights in image data, where spatial localization is crucial, they are less applicable in our context
 266 [Muhammad and Yeasin, 2020, Selvaraju et al., 2020]. The rows and columns of the functional connectivity matrix
 267 capture the correlations among respective brain regions, but without a clear invariance relationship motivating
 268 the use of convolutional kernels across regions. Conversely, it is natural to assume the existence of time-invariant
 269 features for our task. This assumption is validated by our results, which demonstrate performance degradation
 270 when the sequence of connectivity graphs is randomly shuffled across time. This motivation supports the use
 271 of our convolutional neural network operating on the temporal dimension of the data. Class Activation Map

272 techniques would highlight salient time-steps determining the output, but would not provide insights into the
273 relevant brain regions. In contrast, our use of Shapley coefficients allows us to control the masking of specific
274 subnetworks, enabling an analysis of individual brain regions’ contributions to the model’s predictions while
275 accommodating our temporal data setup.

276 Our findings revealed that in the 7-subnetwork analysis, the visual and default subnetworks are key for
277 modality classification, reflecting the intuitive notion that visual processing is essential for distinguishing between
278 movies and audio stories. In content classification, the default mode subnetwork emerged as the most influential,
279 suggesting its essential function in understanding the meaning and content of the stimuli. This aligns with existing
280 research that has highlighted the default mode subnetwork’s involvement in higher-order cognitive functions,
281 such as narrative comprehension [Baldassano et al., 2018, Simony et al., 2016]. The combined classification task
282 emphasized the importance of both visual and default mode networks, as expected.

283 A more fine-grained analysis using the 17-subnetwork parcellation revealed additional insights. While visual
284 and default mode networks remained dominant for modality classification, the somatomotor subnetwork also
285 showed a high Shapley value. The latter can be better understood in the context of embodied cognition and
286 language comprehension. A seminal study of embodied language comprehension demonstrated that passive
287 reading of action words produces a corresponding somatotopic activation of the motor and premotor cortex [Hauk
288 et al., 2004]. Likewise, viewing images or reading sentences describing everyday actions produces a distributed
289 activation in fronto-temporo-parietal network that includes sensory-motor and premotor cortex [Jouen et al., 2015].
290 Similar to the 7-subnetwork analysis, the default mode subnetwork was most influential for content classification.
291 Interestingly, the ventral attention subnetwork also played a significant role. This finding is a step further to
292 answer the open question raised by the study Baldassano et al. [2018] study. They proposed that schematic
293 representations in the brain might not solely rely on top-down activation of scripts in the medial prefrontal cortex.
294 They suggested these representations could serve as building blocks for a complete narrative script formed through
295 a bottom-up process. Our observation of a high Shapley value for the ventral attention subnetwork, which is
296 known to be also associated with bottom-up attentional control, aligns with this possibility. Finally, the combined
297 classification task again highlighted the importance of visual, default mode, and somatomotor A networks.

298 **Limitations and Future Works** It is important to acknowledge that the primary limitation of this study is
299 the size of the dataset used. This may limit the generalizability of our findings to other populations or narrative
300 stimuli. Future research could address this by employing larger datasets, if available. Additionally, exploring the
301 generalizability of these findings across diverse datasets would be valuable. Within the context of the current
302 dataset size, future work could delve deeper into other aspects of narrative processing. One potential direction
303 is to investigate the impact of individual differences in narrative comprehension. For instance, research could
304 explore how factors such as age, reading experience, or cultural background might influence how individuals
305 process narratives based on brain network activity. In addition, it would be beneficial to explore the model’s
306 decision-making process in more detail. Analyzing the learned weights of our neural architectures could provide
307 complementary insights to those obtained from Shapley scores, which focus on model predictions. This approach
308 could provide a clearer understanding of how specific brain regions contribute to the classification task. Finally,
309 future work could explore the temporal dynamics of narrative processing by examining the role of specific time
310 windows in the classification task. Masking entire time steps and assessing the effect of window size on classification
311 performance may shed light on how temporal information is integrated to understand narratives.

312 **Conclusion** Overall, our work demonstrates the potential of combining machine learning models with explainable
313 AI techniques like Shapley values to understand the role of brain subnetworks during narrative processing. Our
314 findings not only contribute to a deeper understanding of how the brain processes narratives but also showcase
315 the broader applicability of this approach. In tasks where the role of specific brain regions remains unclear, this
316 methodology can provide valuable new insights. By highlighting subnetwork contributions through Shapley values,
317 we can generate novel hypotheses about the functional roles of these regions. In our case, the model’s performance
318 aligns with existing literature on narrative comprehension, validating the approach. Importantly, this research
319 validates an alternative and complementary method for investigating brain function in human cognition, which
320 involves functional connectivity. This successful validation paves the way for further exploration of brain networks
321 not only in higher-order cognition, motor tasks, and emotional processing but also in any domain where the neural
322 basis remains partially understood.

323 Author contributions

324 A.R. processed the data, designed the model, performed the experiments, drafted the original manuscript and
325 contributed to its revisions. Y.A. preprocessed the data. E.N. designed the model, supervised the project, reviewed
326 and edited the manuscript. S.D.G. and P.F.D. conceived and supervised the project, reviewed and edited the
327 manuscript.

328 References

- 329 Sarah Aliko, Bangjie Wang, Steven L Small, and Jeremy I Skipper. *The entire brain, more or less, is at work*
330 : ‘Language regions’ are artefacts of averaging, September 2023. URL <http://biorxiv.org/lookup/doi/10.1101/2023.09.01.555886>.
- 332 Nicola Amoroso, Silvano Quarto, Marianna La Rocca, Sabina Tangaro, Alfonso Monaco, and Roberto Bellotti.
333 An eXplainability Artificial Intelligence approach to brain connectivity in Alzheimer’s disease. *Frontiers in*
334 *Aging Neuroscience*, 15, August 2023. doi: 10.3389/fnagi.2023.1238065. URL [https://www.frontiersin.org/](https://www.frontiersin.org/articles/10.3389/fnagi.2023.1238065/full)
335 [articles/10.3389/fnagi.2023.1238065/full](https://www.frontiersin.org/articles/10.3389/fnagi.2023.1238065/full).
- 336 Christopher Baldassano, Uri Hasson, and Kenneth A. Norman. Representation of Real-World Event Schemas
337 during Narrative Perception. *The Journal of Neuroscience*, 38(45):9689–9699, November 2018. doi: 10.1523/
338 JNEUROSCI.0251-18.2018. URL [https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.0251-18.](https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.0251-18.2018)
339 [2018](https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.0251-18.2018).
- 340 Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B Shah. Julia: A fresh approach to numerical computing.
341 *SIAM review*, 59(1):65–98, 2017. URL <https://doi.org/10.1137/141000671>.
- 342 Jeffrey R. Binder and Rutvik H. Desai. The neurobiology of semantic memory. *Trends in Cognitive Sciences*,
343 15(11):527–536, November 2011. doi: 10.1016/j.tics.2011.10.001. URL [https://linkinghub.elsevier.com/](https://linkinghub.elsevier.com/retrieve/pii/S1364661311002142)
344 [retrieve/pii/S1364661311002142](https://linkinghub.elsevier.com/retrieve/pii/S1364661311002142).
- 345 Simon Danisch and Julius Krumbiegel. Makie.jl: Flexible high-performance data visualization for Julia. *Journal*
346 *of Open Source Software*, 6(65):3349, 2021. doi: 10.21105/joss.03349. URL [https://doi.org/10.21105/joss.](https://doi.org/10.21105/joss.03349)
347 [03349](https://doi.org/10.21105/joss.03349).
- 348 Fraida Dubin and David Bycina. Academic reading and the esl/efl teacher. *Teaching English as a second or*
349 *foreign language*, 2:195–215, 1991.
- 350 Olaf Hauk, Ingrid Johnsrude, and Friedemann Pulvermüller. Somatotopic Representation of Action Words in
351 Human Motor and Premotor Cortex. *Neuron*, 41(2):301–307, January 2004. doi: 10.1016/S0896-6273(03)00838-9.
352 URL <https://linkinghub.elsevier.com/retrieve/pii/S0896627303008389>.
- 353 Xiaojie Huang, Jun Xiao, and Chao Wu. Design of Deep Learning Model for Task-Evoked fMRI Data Classification.
354 *Computational Intelligence and Neuroscience*, 2021:1–10, August 2021. doi: 10.1155/2021/6660866. URL
355 <https://www.hindawi.com/journals/cin/2021/6660866/>.
- 356 Michael Innes, Elliot Saba, Keno Fischer, Dhairya Gandhi, Marco Concetto Rudilosso, Neethu Mariya Joy, Tejan
357 Karmali, Avik Pal, and Viral Shah. Fashionable modelling with flux. *CoRR*, abs/1811.01457, 2018. URL
358 <https://arxiv.org/abs/1811.01457>.
- 359 A.L. Jouen, T.M. Ellmore, C.J. Madden, C. Pallier, P.F. Dominey, and J. Ventre-Dominey. Beyond the word
360 and image: characteristics of a common meaning system for language and vision revealed by functional and
361 structural imaging. *NeuroImage*, 106:72–85, February 2015. doi: 10.1016/j.neuroimage.2014.11.024. URL
362 <https://linkinghub.elsevier.com/retrieve/pii/S1053811914009410>.
- 363 Iiro P. Jääskeläinen, Mikko Sams, Enrico Glerean, and Jyrki Ahveninen. Movies and narratives as naturalistic
364 stimuli in neuroimaging. *NeuroImage*, 224:117445, 2021. doi: <https://doi.org/10.1016/j.neuroimage.2020.117445>.
365 URL <https://www.sciencedirect.com/science/article/pii/S1053811920309307>.
- 366 Timothy A. Keller, Patricia A. Carpenter, and Marcel Adam Just. The Neural Bases of Sentence Comprehension:
367 a fMRI Examination of Syntactic and Lexical Processing. *Cerebral Cortex*, 11(3):223–237, 03 2001. ISSN
368 1047-3211. doi: 10.1093/cercor/11.3.223. URL <https://doi.org/10.1093/cercor/11.3.223>.

- 369 Byung-Hoon Kim, Jong Chul Ye, and Jae-Jin Kim. Learning dynamic graph representation of brain connectome
370 with spatio-temporal attention. In *Advances in Neural Information Processing Systems*, volume 34, pages
371 4314–4327. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/paper/
372 2021/file/22785dd2577be2ce28ef79febe80db10-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/22785dd2577be2ce28ef79febe80db10-Paper.pdf).
- 373 Peter Kim, Junbeom Kwon, Sunghwan Joo, Sangyoon Bae, Donggyu Lee, Yoonho Jung, Shinjae Yoo, Jiook
374 Cha, and Taesup Moon. Swift: Swin 4d fmri transformer. In *Advances in Neural Information Processing
375 Systems*, volume 36, pages 42015–42037. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.
376 cc/paper_files/paper/2023/file/8313b1920ee9c78d846c5798c1ce48be-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/8313b1920ee9c78d846c5798c1ce48be-Paper-Conference.pdf).
- 377 Rolf Kötter. Shapley ratings in brain networks. *Frontiers in Neuroinformatics*, 1, 2007. ISSN 1662-5196.
378 doi: 10.3389/neuro.11.002.2007. URL [http://journal.frontiersin.org/article/10.3389/neuro.11.002.
379 2007/abstract](http://journal.frontiersin.org/article/10.3389/neuro.11.002.2007/abstract).
- 380 Xiaoxiao Li, Yuan Zhou, Nicha C. Dvornek, Yufeng Gu, Pamela Ventola, and James S. Duncan. Efficient
381 Shapley Explanation for Features Importance Estimation Under Uncertainty. In *Medical Image Computing and
382 Computer Assisted Intervention – MICCAI 2020*, volume 12261. Springer International Publishing, 2020. doi:
383 10.1007/978-3-030-59710-8_77.
- 384 Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the
385 31st International Conference on Neural Information Processing Systems, NIPS'17*, page 4768–4777. Curran
386 Associates Inc., 2017. ISBN 9781510860964.
- 387 Mohammed Bany Muhammad and Mohammed Yeasin. Eigen-cam: Class activation map using principal
388 components. In *2020 international joint conference on neural networks (IJCNN)*, pages 1–7. IEEE, 2020.
- 389 Samuel A. Nastase, Yun-Fei Liu, Hanna Hillman, Asieh Zadbood, Liat Hasenfratz, Neggin Keshavarzian, Janice
390 Chen, Christopher J. Honey, Yaara Yeshurun, Mor Regev, Mai Nguyen, Claire H. C. Chang, Christopher
391 Baldassano, Olga Lositsky, Erez Simony, Michael A. Chow, Yuan Chang Leong, Paula P. Brooks, Emily
392 Micciche, Gina Choe, Ariel Goldstein, Tamara Vanderwal, Yaroslav O. Halchenko, Kenneth A. Norman, and
393 Uri Hasson. "narratives", 2020.
- 394 Maria Giulia Preti, Thomas AW Bolton, and Dimitri Van De Ville. The dynamic functional connec-
395 tome: State-of-the-art and perspectives. *NeuroImage*, 160:41–54, 2017. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2016.12.061>. URL [https://www.sciencedirect.com/science/article/pii/
396 //doi.org/10.1016/j.neuroimage.2016.12.061](https://www.sciencedirect.com/science/article/pii/S1053811916307881). URL [https://www.sciencedirect.com/science/article/pii/
397 S1053811916307881](https://www.sciencedirect.com/science/article/pii/S1053811916307881). Functional Architecture of the Brain.
- 398 Cathy J. Price. A review and synthesis of the first 20years of PET and fMRI studies of heard speech, spoken language
399 and reading. *NeuroImage*, 62(2):816–847, August 2012. ISSN 10538119. doi: 10.1016/j.neuroimage.2012.04.062.
400 URL <https://linkinghub.elsevier.com/retrieve/pii/S1053811912004703>.
- 401 Srikanth Ryali, Yuan Zhang, Carlo de Los Angeles, Kaustubh Supekar, and Vinod Menon. Deep learning models
402 reveal replicable, generalizable, and behaviorally relevant sex differences in human functional brain organization.
403 *Proceedings of the National Academy of Sciences*, 121(9):e2310012121, 2024.
- 404 Maham Saeidi, Waldemar Karwowski, Farzad V. Farahani, Krzysztof Fiok, P. A. Hancock, Ben D. Sawyer, Leonardo
405 Christov-Moore, and Pamela K. Douglas. Decoding Task-Based fMRI Data with Graph Neural Networks,
406 Considering Individual Differences. *Brain Sciences*, 12(8):1094, August 2022. doi: 10.3390/brainsci12081094.
407 URL <https://www.mdpi.com/2076-3425/12/8/1094>.
- 408 Alexander Schaefer, Ru Kong, Evan M Gordon, Timothy O Laumann, Xi-Nian Zuo, Avram J Holmes, Simon B
409 Eickhoff, and BT Thomas Yeo. Local-global parcellation of the human cerebral cortex from intrinsic functional
410 connectivity mri. *Cerebral cortex*, 28(9):3095–3114, 2018.
- 411 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv
412 Batra. Grad-cam: visual explanations from deep networks via gradient-based localization. *International journal
413 of computer vision*, 128:336–359, 2020.
- 414 Lloyd S Shapley. Notes on the n-person game—ii: The value of an n-person game. 1951.
- 415 Lauren J. Silbert, Christopher J. Honey, Erez Simony, David Poeppel, and Uri Hasson. Coupled neural systems
416 underlie the production and comprehension of naturalistic narrative speech. *Proceedings of the National
417 Academy of Sciences*, 111(43), October 2014. doi: 10.1073/pnas.1323812111. URL [https://pnas.org/doi/
418 full/10.1073/pnas.1323812111](https://pnas.org/doi/full/10.1073/pnas.1323812111).

- 419 Erez Simony, Christopher J Honey, Janice Chen, Olga Lositsky, Yaara Yeshurun, Ami Wiesel, and Uri Hasson.
420 Dynamic reconfiguration of the default mode network during narrative comprehension. *Nature Communications*,
421 7(1):12141, July 2016. doi: 10.1038/ncomms12141. URL <https://www.nature.com/articles/ncomms12141>.
- 422 Ann E. Sizemore and Danielle S. Bassett. Dynamic graph metrics: Tutorial, toolbox, and tale. *NeuroImage*,
423 180:417–427, 2018. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2017.06.081>. URL <https://www.sciencedirect.com/science/article/pii/S1053811917305645>. Brain Connectivity Dynamics.
- 425 B. T. Thomas Yeo, Fenna M. Krienen, Jorge Sepulcre, Mert R. Sabuncu, Danial Lashkari, Marisa Hollinshead,
426 Joshua L. Roffman, Jordan W. Smoller, Lilla Zöllei, Jonathan R. Polimeni, Bruce Fischl, Hesheng Liu, and
427 Randy L. Buckner. The organization of the human cerebral cortex estimated by intrinsic functional connectivity.
428 *Journal of Neurophysiology*, 106(3):1125–1165, September 2011. doi: 10.1152/jn.00338.2011.
- 429 Koene R. A. Van Dijk, Trey Hedden, Archana Venkataraman, Karleyton C. Evans, Sara W. Lazar, and Randy L.
430 Buckner. Intrinsic Functional Connectivity As a Tool For Human Connectomics: Theory, Properties, and
431 Optimization. *Journal of Neurophysiology*, 103(1):297–321, January 2010. doi: 10.1152/jn.00783.2009. URL
432 <https://www.physiology.org/doi/10.1152/jn.00783.2009>.
- 433 Simone Vessel, Joy J. Geng, and Gereon R. Fink. Dorsal and Ventral Attention Systems: Distinct Neural Circuits
434 but Collaborative Roles. *The Neuroscientist*, 20(2):150–159, April 2014. doi: 10.1177/1073858413494269.
- 435 Roel M. Willems, Samuel A. Nastase, and Branka Milivojevic. Narratives for Neuroscience. *Trends in Neurosciences*,
436 43(5):271–273, May 2020. ISSN 01662236. doi: 10.1016/j.tins.2020.03.003. URL <https://linkinghub.elsevier.com/retrieve/pii/S0166223620300497>.
- 438 Jiang Xu, Stefan Kemeny, Grace Park, Carol Frattali, and Allen Braun. Language in context: emergent features
439 of word, sentence, and narrative comprehension. *NeuroImage*, 25(3):1002–1015, April 2005. doi: 10.1016/j.
440 neuroimage.2004.12.013. URL <https://linkinghub.elsevier.com/retrieve/pii/S1053811904007748>.

441 5 Appendix

442 5.1 Choice of parameter τ

443 The following figure shows the evolution of the model’s accuracy as a function of the third dimension of the
444 convolutional filter (i.e. τ). For the modality classification, we set $\tau = 4$, since model performance seems not to
445 increase significantly beyond this value (Figure 6a). For the content and combined classification, we set $\tau = 8$
446 since the model performance seems the best for this value (Figure 6b and Figure 6c). It is important to highlight
that when $\tau = 8$ the convolution behaviour is similar to the one of dense layer.

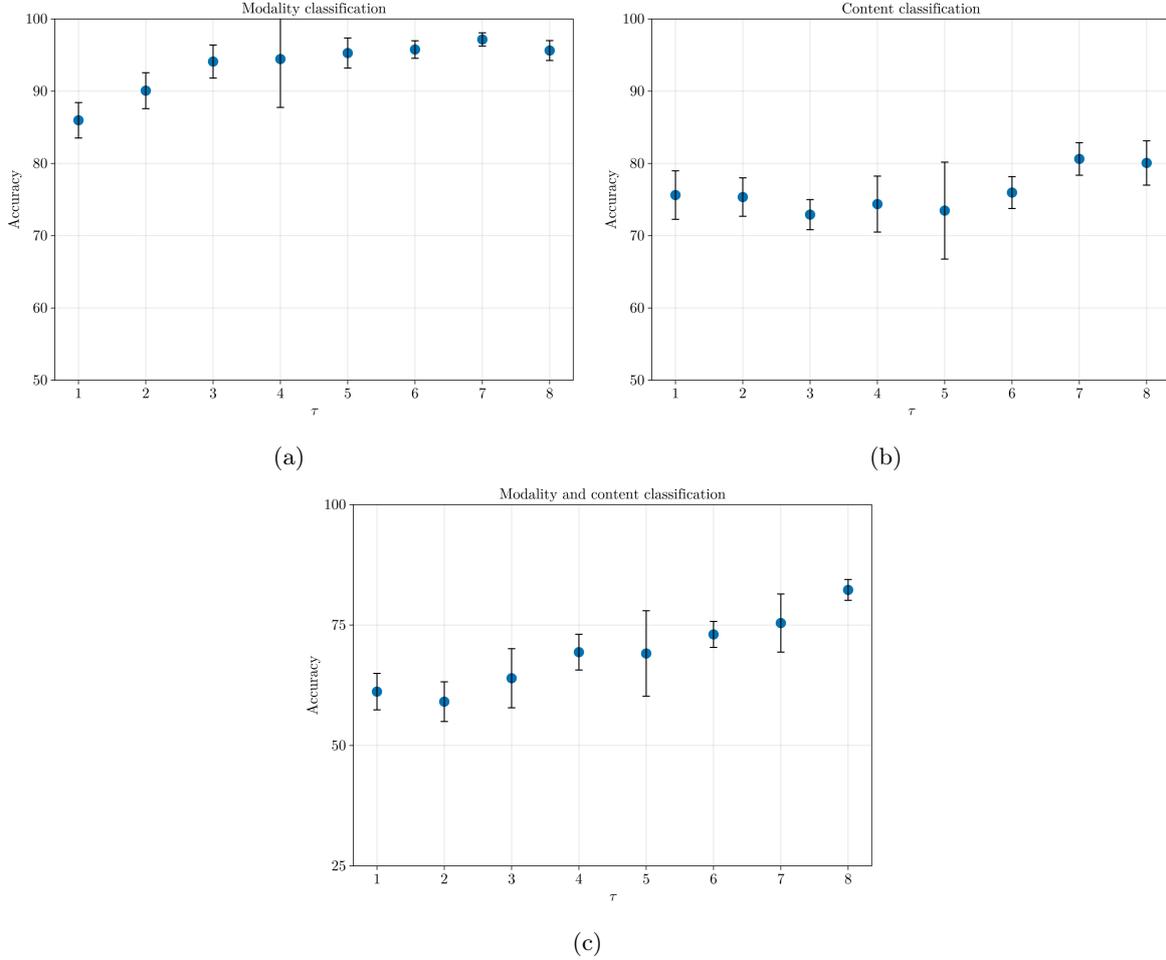
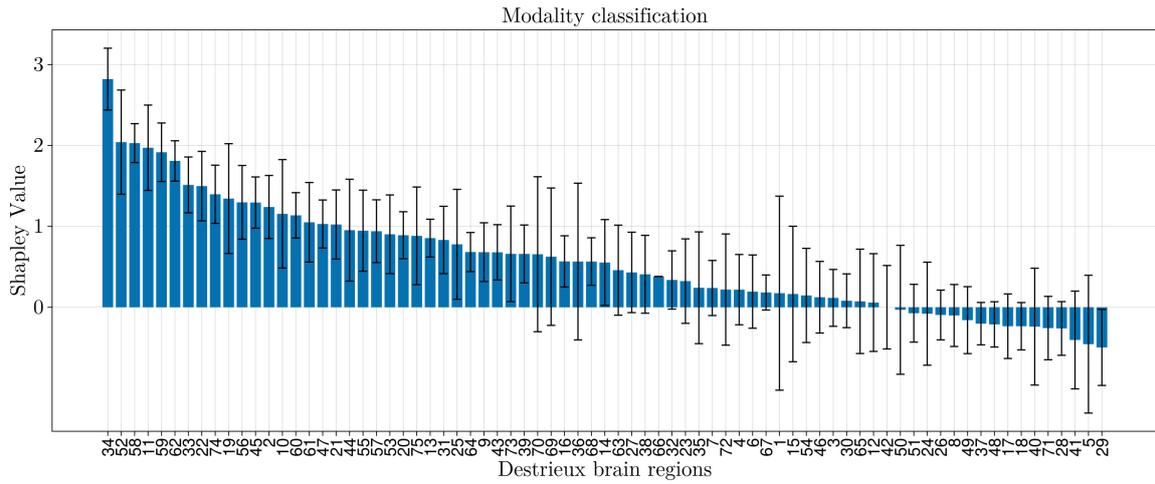


Figure 6: Model’s accuracy as a function of the third dimension of the convolutional filter.

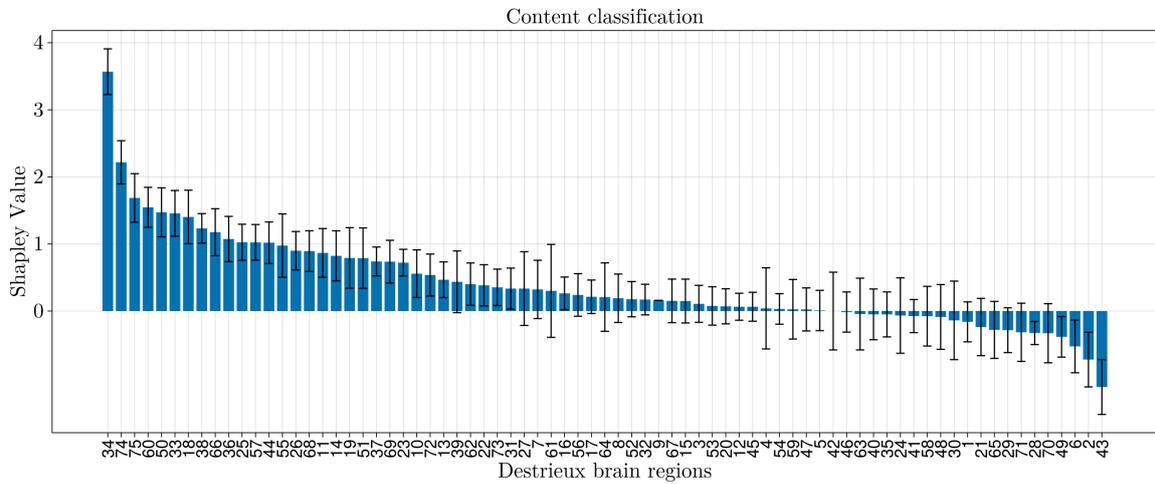
447

448 5.2 Destrieux parcellation

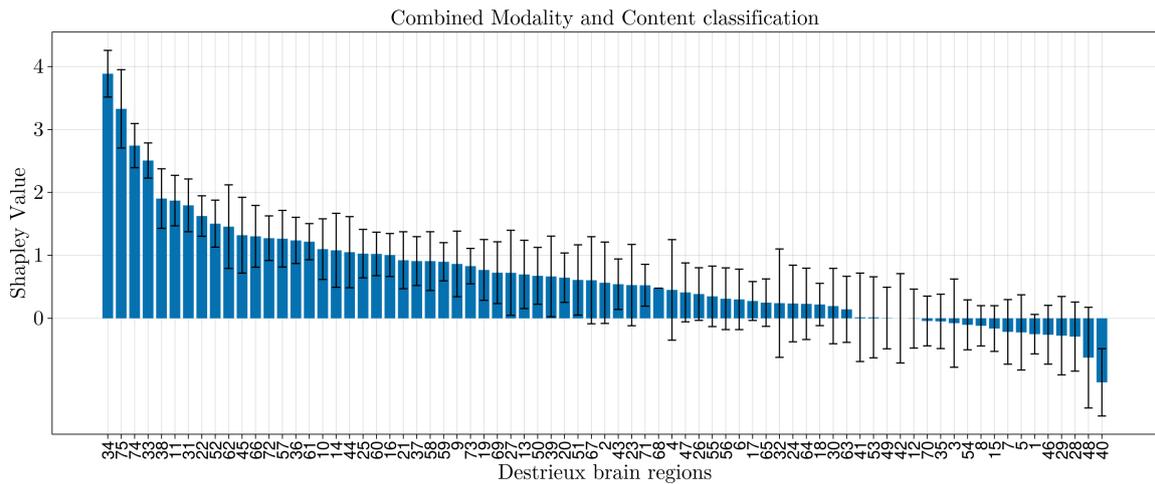
449 The following figure shows the Shapley values for the Destrieux parcellation in Figure 7. The Shapley values are
450 calculated for the modality classification (Figure 7a), content classification (Figure 7b), and combined classification
451 (Figure 7c). Shapley values were calculated for 75 brain regions, with area 34 consistently highlighted as a
452 significant region across all tasks. Area 34, located in the superior temporal gyrus, includes key structures such as
453 Brodmann’s areas, which contain the auditory cortex responsible for sound perception. It also includes Wernicke’s
454 area, which is essential for processing speech into understandable language. Given these critical functions, it is
455 not surprising that area 34 plays a central role in narrative-related tasks. In addition, its involvement in the
456 default mode network and the ventral attention network, both of which are essential for narrative processing, is
457 consistent with Yeo’s parcellation findings.



(a)



(b)



(c)

Figure 7: This figure shows the contribution of 75 Destrieux regions, computed using Shapley values, for classifying narratives with our machine learning model. The bars represent the average contribution of each region to the model’s predictions, with higher values indicating greater influence. The error bars denote the standard deviation of the Shapley values. The correspondence between label and region can be found in the Table 2

	Destrieux labels		
0	null	38	Middle_temporal_gyrus
1	Fronto-marginal_gyrus+sulcus	39	Anterior_segment_of_lateral_sulcus_horizontal_ramus
2	Inferior_occipital_gyrus+sulcus	40	Anterior_segment_of_lateral_sulcus_vertical_ramus
3	Paracentral_lobule+sulcus	41	Lateral_sulcus_posterior_ramus
4	Subcentral_gyrus+sulci	42	null
5	Transverse_frontopolar_gyri+sulci	43	Occipital_pole
6	Cingulate_gyrus+sulcus_anterior_part	44	Temporal_pole
7	Cingulate_gyrus+sulcus_middle-anterior_part	45	Calcarine_sulcus
8	Cingulate_gyrus+sulcus_middle-posterior_part	46	Central_sulcus
9	Cingulate_gyrus_posterior-dorsal_part	47	Cingulate_sulcus_marginal_branch
10	Cingulate_gyrus_posterior-ventral_part	48	Insula_circular_sulcus_anterior_part
11	Cuneus	49	Insula_circular_sulcus_inferior_part
12	Inferior_frontal_gyrus_opercular_part	50	Insula_circular_sulcus_superior_part
13	Inferior_frontal_gyrus_orbital_part	51	Anterior_transverse_collateral_sulcus
14	Inferior_frontal_gyrus_triangular_part	52	Posterior_transverse_collateral_sulcus
15	Middle_frontal_gyrus	53	Inferior_frontal_sulcus
16	Superior_frontal_gyrus	54	Middle_frontal_sulcus
17	Insula_insular_gyrus+central_sulcus	55	Superior_frontal_sulcus
18	Insular_gyri_short	56	Sulcus_intermedius_primus
19	Middle_occipital_gyrus	57	Intraparietal_sulcus+transverse_parietal_sulci
20	Superior_occipital_gyrus	58	Middle_occipital+lunatus_sulcus
21	Lateral_occipito-temporal_gyrus	59	Superior+transverse_occipital_sulcus
22	Lingual_gyrus	60	Anterior_occipital_sulcus+preoccipital_notch
23	Parahippocampal_gyrus	61	Lateral_occipito-temporal_sulcus
24	Orbital_gyri	62	Collateral+lingual_sulcus
25	Angular_gyrus	63	Lateral_orbital_sulcus
26	Supramarginal_gyrus	64	Olfactory_sulcus
27	Superior_parietal_lobule	65	Orbital_sulci
28	Postcentral_gyrus	66	Parieto-occipital_sulcus
29	Precentral_gyrus	67	Pericallosal_sulcus
30	Precuneus	68	Postcentral_sulcus
31	Straight_gyrus	69	Precentral_sulcus_inferior_part
32	Subcallosal_area+gyrus	70	Precentral_sulcus_superior_part
33	Anterior_transverse_temporal_gyrus	71	Suborbital_sulcus
34	Superior_temporal_gyrus_lateral_aspect	72	Subparietal_sulcus
35	Superior_temporal_gyrus_planum_polare	73	Inferior_temporal_sulcus
36	Planum_temporale	74	Superior_temporal_sulcus
37	Inferior_temporal_gyrus	75	Transverse_temporal_sulcus

Table 2: Correspondence between Destrieux labels and regions

458 5.3 Desikan parcellation

459 The following figure shows the Shapley values for the Desikan parcellation in Figure 8. The Shapley values are
460 calculated for the modality classification (Figure 8a), content classification (Figure 8b), and combined classification
461 (Figure 8c). Shapley values were calculated for 70 brain regions, showing that area 34, associated with the
462 temporal pole, has the highest value in modality classification. This region is associated with several high-level
463 cognitive processes, particularly visual processing of complex objects and face recognition. This is followed by
464 region 22, the pericalcarine cortex or primary visual cortex, which is primarily responsible for processing visual
465 information. In content classification, the right and left banks of the superior temporal sulcus stand out. These
466 regions serve as hubs for social perception and cognition, including recognition of faces and human movement, as
467 well as understanding actions, mental states and language. In addition, region 31, the superior temporal gyrus,
468 remains important, consistent with previous findings. In the combined classification task, the middle temporal
469 gyrus, the pericalcarine cortex and the superior temporal sulcus emerge as the most involved regions.

470 5.4 Inter-intra subject variability

471 **Intra-subject standard deviation (SD):**We calculated the intra-subject standard deviation by first computing
472 the standard deviation of accuracy for each individual subject. These individual standard deviations were then
473 averaged across all subjects. The intra-subject standard deviation is given by:

$$SD = \frac{1}{N} \sum_{i=1}^N SD_i$$

474 where N is the total number of subjects and SD_i is computed as follows:

$$SD_i = \sqrt{\frac{1}{M-1} \sum_{k=1}^M (Accuracy_{i,k} - MeanAccuracy_i)^2}$$

475 Here $Accuracy_{i,k}$ is represents the accuracy for the k -th sample of subject i , $MeanAccuracy_i$ is the mean accuracy
476 for subject i , and M is the number of samples for each subject (16).

477 **Inter-subject standard deviation (SD):**The inter-subject standard deviation was calculated by taking the
478 standard deviation of the mean accuracy values across all subjects:

$$SD = SD\left(\frac{1}{N} \sum_{j=1}^N MeanAccuracy_j\right)$$

479 where $Accuracy_j$ is the mean accuracy for subject j and N is the total number of subjects.

480 **Results:**

481 *For modality classification:*

- 482 • Intra-subject standard deviation: 16.90%
- 483 • Inter-subject standard deviation: 4.91%
- 484 • Total variability: 21.68%

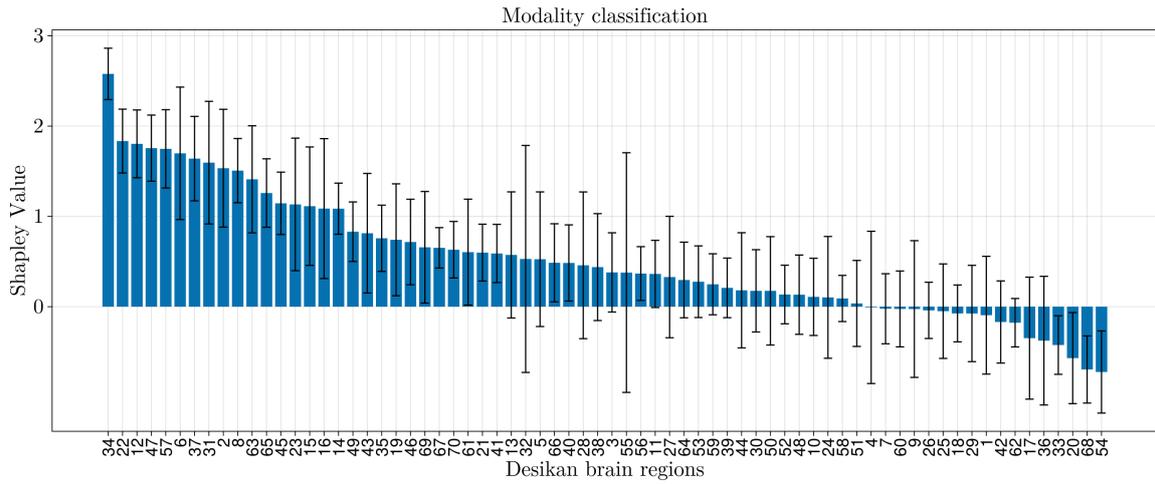
485 *For content classification:*

- 486 • Intra-subject standard deviation: 37.05%
- 487 • Inter-subject standard deviation: 8.04%
- 488 • Total variability: 38.02%

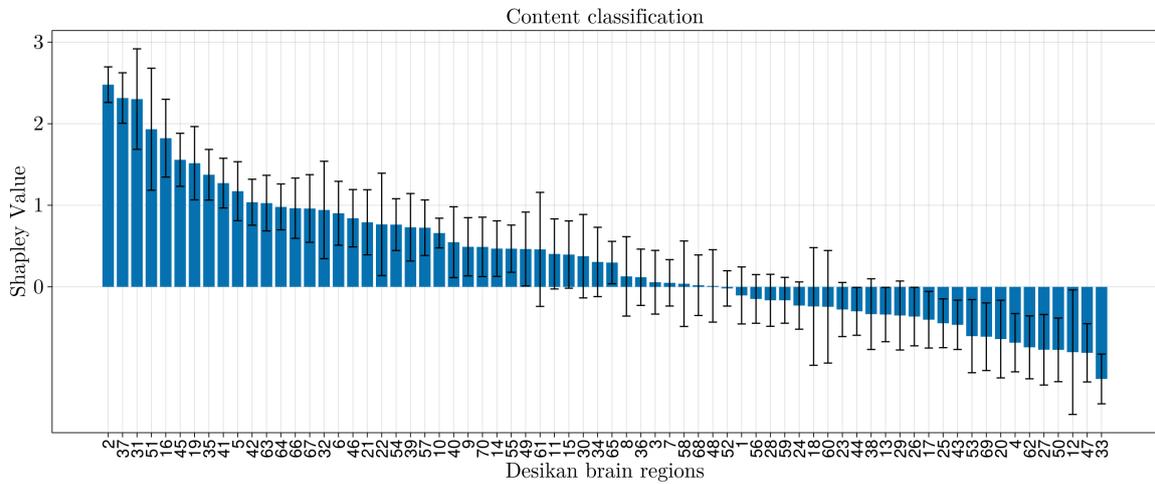
489 *For combined classification:*

- 490 • Intra-subject standard deviation: 37.74%
- 491 • Inter-subject standard deviation: 9.21%
- 492 • Total variability: 38.92%

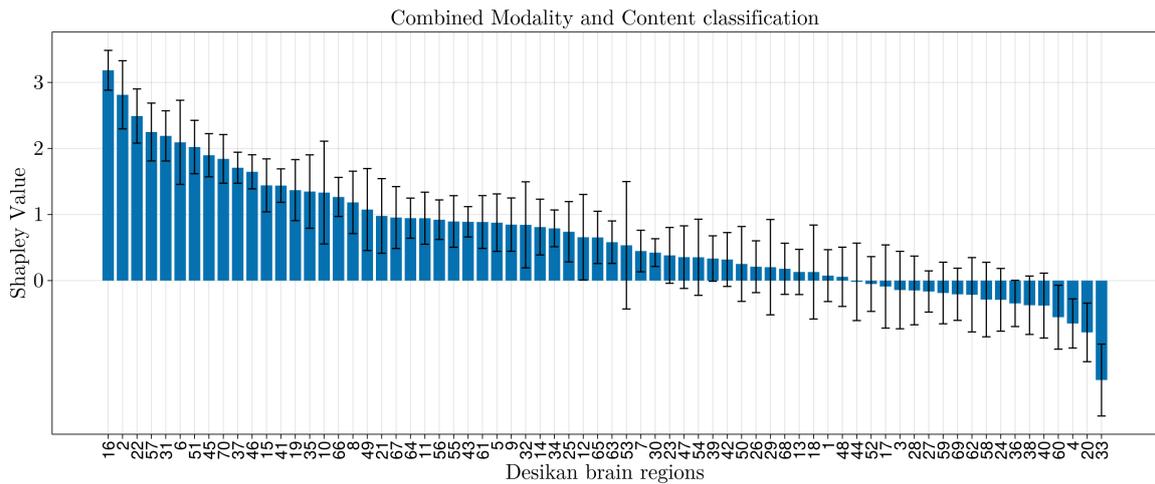
493 We observed that the intra-subject variability is notably higher compared to the inter-subject variability. This
494 disparity could be attributed to the larger number of subjects (31) relative to the smaller number of samples per
495 subject (16).



(a)



(b)



(c)

Figure 8: This figure shows the contribution of 70 Desikan regions, computed using Shapley values, for classifying narratives with our machine learning model. The bars represent the average contribution of each region to the model’s predictions, with higher values indicating greater influence. The error bars denote the standard deviation of the Shapley values. The correspondence between label and region can be found in the Table 3

	Desikan labels
1	L_white_matter
2	L_Banks_superior_temporal_sulcus
3	L_caudal_anterior_cingulate_cortex
4	L_caudal_middle_frontal_gyrus
5	L_corpus_calosum
6	L_cuneus_cortex
7	L_entorhinal_cortex
8	L_fusiform_gyrus
9	L_inferior_parietal_cortex
10	L_inferior_temporal_gyrus
11	L_isthmus-cingulate_cortex
12	L_lateral_occipital_cortex
13	L_lateral_orbitofrontal_cortex
14	L_lingual_gyrus
15	L_medial_orbitofrontal_cortex
16	L_middle_temporal_gyrus
17	L parahippocampal_gyrus
18	L_paracentral_lobule
19	L_pars_opercularis
20	L_pars_orbitalis
21	L_pars_triangularis
22	L_pericalcarine_cortex
23	L_postcentral_gyrus
24	L_posterior-cingulate_cortex
25	L_precentral_gyrus
26	L_precuneus_cortex
27	L_rostral_anterior_cingulate_cortex
28	L_rostral_middle_frontal_gyrus
29	L_superior_frontal_gyrus
30	L_superior_parietal_cortex
31	L_superior_temporal_gyrus
32	L_supramarginal_gyrus
33	L_frontal_pole
34	L_temporal_pole
35	L_transverse_temporal_cortex

36	R_white_matter
37	R_Banks_superior_temporal_sulcus
38	R_caudal_anterior_cingulate_cortex
39	R_caudal_middle_frontal_gyrus
40	R_corpus_calosum
41	R_cuneus_cortex
42	R_entorhinal_cortex
43	R_fusiform_gyrus
44	R_inferior_parietal_cortex
45	R_inferior_temporal_gyrus
46	R_isthmus-cingulate_cortex
47	R_lateral_occipital_cortex
48	R_lateral_orbitofrontal_cortex
49	R_lingual_gyrus
50	R_medial_orbitofrontal_cortex
51	R_middle_temporal_gyrus
52	R parahippocampal_gyrus
53	R_paracentral_lobule
54	R_pars_opercularis
55	R_pars_orbitalis
56	R_pars_triangularis
57	R_pericalcarine_cortex
58	R_postcentral_gyrus
59	R_posterior-cingulate_cortex
60	R_precentral_gyrus
61	R_precuneus_cortex
62	R_rostral_anterior_cingulate_cortex
63	R_rostral_middle_frontal_gyrus
64	R_superior_frontal_gyrus
65	R_superior_parietal_cortex
66	R_superior_temporal_gyrus
67	R_supramarginal_gyrus
68	R_frontal_pole
69	R_temporal_pole
70	R_transverse_temporal_cortex

Table 3: Correspondence between Desikan labels and regions

496 5.5 Yeo parcellations black and white compatible

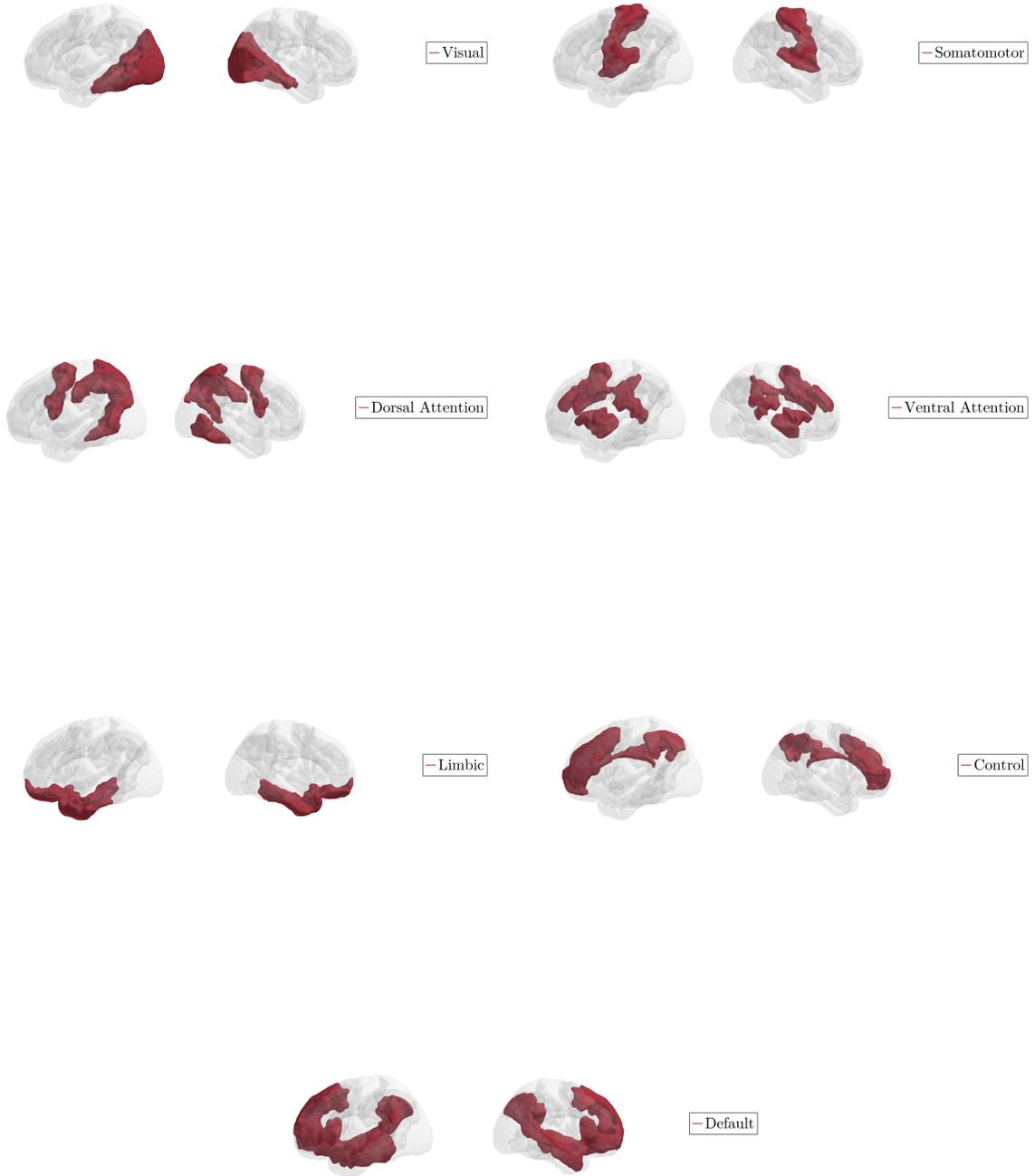
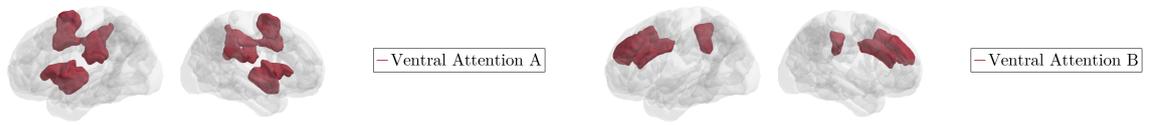
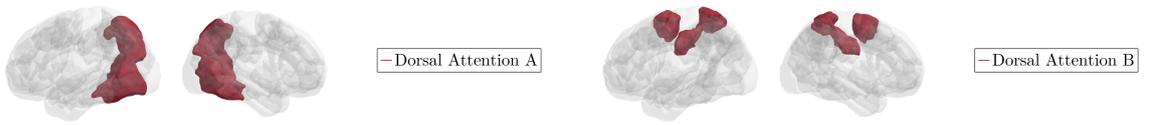


Figure 9: The Yeo 7-subnetwork parcellation illustrates seven distinct functional networks within the brain, each associated with specific cognitive functions. The regions shown in red correspond to areas included in each subnetwork, while the rest of the brain remains in grayscale for contrast. Each pair of images shows the subnetwork from different angles to highlight the distribution of each functional network across the brain. This figure is optimized for black-and-white printing, with clear contrast to make the subnetworks easily distinguishable.



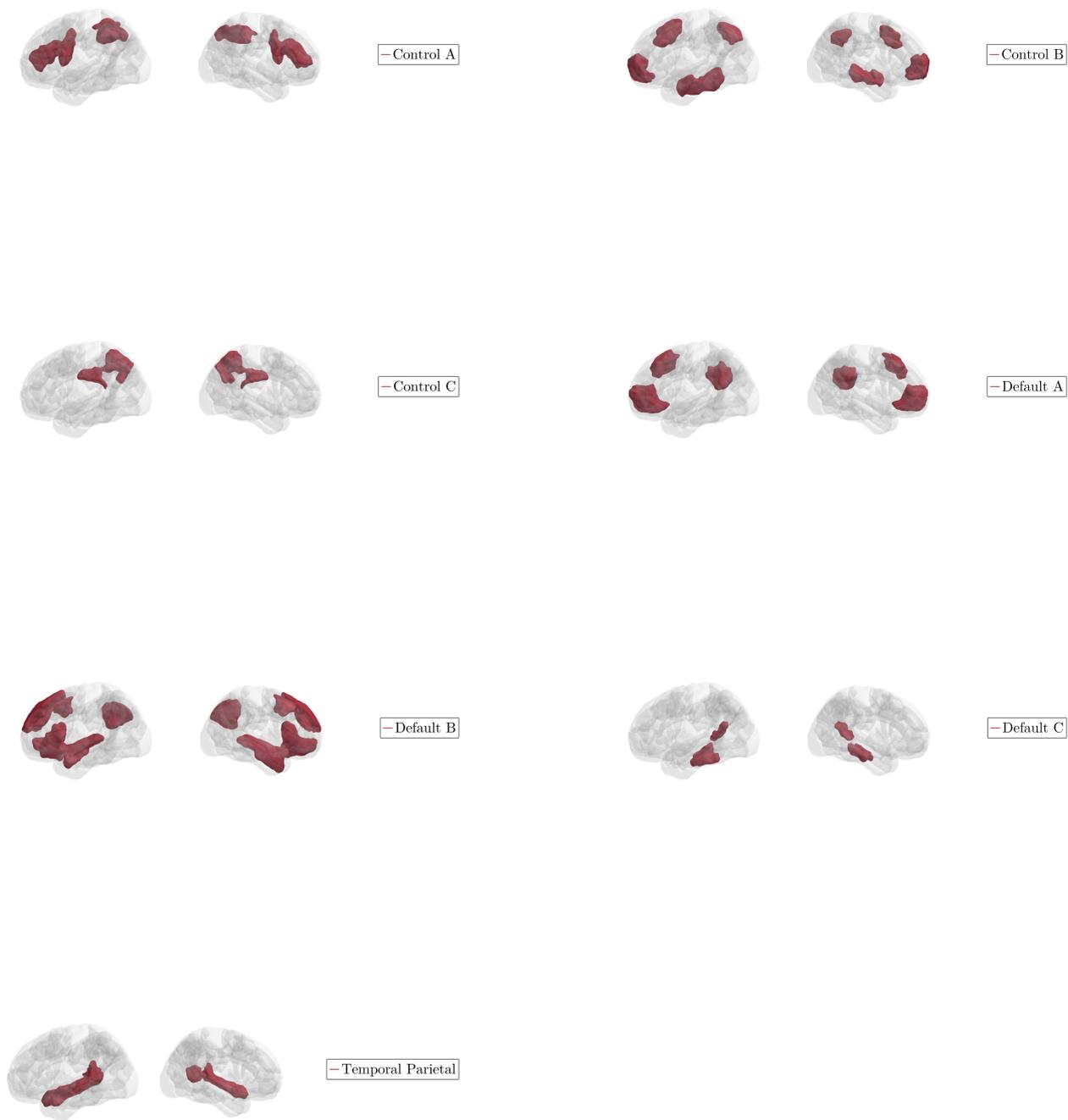


Figure 11: The Yeo 17-subnetwork parcellation illustrates seventeen distinct functional networks within the brain, each associated with specific cognitive functions. The regions shown in red correspond to areas included in each subnetwork, while the rest of the brain remains in grayscale for contrast. Each pair of images shows the subnetwork from different angles, highlighting the distribution of each functional network across the brain. This figure is optimized for black-and-white printing, with clear contrast to make the subnetworks easily distinguishable.