

# Can we predict sleep health based on brain features? A large-scale machine learning study

Federico Raimondo<sup>1,2</sup>, Hanwen Bi<sup>1,2</sup>, Vera Komeyer<sup>1,2,5</sup>, Jan Kasper<sup>1,2</sup>, Sabrina Primus<sup>3</sup>, Felix Hoffstaedter<sup>1,2</sup>, Synchon Mandal<sup>1,2</sup>, Laura Waite<sup>1</sup>, Juliane Winkelmann<sup>3</sup>, Konrad Oexle<sup>3</sup>, Simon B. Eickhoff<sup>1,2</sup>, Masoud Tahmasian<sup>1,2,4</sup>, Kaustubh R. Patil<sup>1,2</sup>

1. Institute of Neuroscience and Medicine, Brain and Behavior (INM-7), Forschungszentrum Jülich, Jülich, Germany
2. Institute for Systems Neuroscience, Medical Faculty, Heinrich-Heine University Düsseldorf, Düsseldorf, Germany
3. Institute of Neurogenetics (ING), Helmholtz Zentrum München, München, Germany.
4. Department of Nuclear Medicine, University Hospital and Medical Faculty, University of Cologne, Cologne, Germany
5. Department of Biology, Faculty of Mathematics and Natural Sciences, Heinrich Heine University Düsseldorf, Germany

Corresponding author: MT and KRP

## Abstract

**Objectives:** Normal sleep is crucial for brain health. Recent studies have reported robust associations between sleep disturbance and various brain structural and functional traits. However, the complex interplay between sleep health and macro-scale brain organization remains inconclusive. In this study, we aimed to uncover the links between brain imaging features and diverse sleep health-related characteristics by means of Machine Learning (ML).

**Methods:** We used 28,088 participants from the UK Biobank to calculate 4677 structural and functional neuroimaging markers. Then, we employed them to predict self-reported insomnia symptoms, sleep duration, easiness getting up in the morning, chronotype, daily nap, daytime sleepiness, and snoring. We built seven different linear and nonlinear ML models for each sleep health-related characteristic to assess their predictability.

**Results:** We performed an extensive ML analysis that involved more than 100,000 hours of computing. We observed relatively low performance in predicting all sleep health-related characteristics (e.g., balanced accuracy ranging between 0.50-0.59). Across all models, the best performance achieved was 0.59, using a Linear SVM to predict easiness getting up in the morning.

**Conclusions:** The low capability of multimodal neuroimaging markers in predicting sleep health-related characteristics, even under extensive ML optimization in a large population sample suggests a complex relationship between sleep health and brain organization.

**Keywords:** Sleep health; grey matter volume; white matter; functional MRI; UK Biobank; machine learning

## Introduction

Sleep is a non-negotiable human need, which has pivotal impacts on memory processing, metabolite clearance, immune system adaptation, optimal cognition, and mental health (Walker, 2021). The intricate relationship between sleep, brain, and behavior has recently garnered significant scientific attention (Cheng et al., 2018; Ell et al., 2023; Fjell, Sørensen, Wang, Amlie, Baaré, Bartrés-Faz, Bertram, et al., 2023; Li et al., 2022; Tahmasian et al., 2020; Wang et al., 2023; Weihs et al., 2023). Sleep health (SH) is a multidimensional concept that includes assessment of satisfaction, alertness, regularity, timing, and duration of sleep (Buysse, 2014), which is considered a crucial indicator of human well-being. Seven different SH-related characteristics i.e., sleep duration, easiness/difficulty getting up in the morning, chronotype, nap, daytime dozing/sleepiness, as well as insomnia symptoms and snoring reflecting various SH dimensions were collected in half a million participants in the UK Biobank (UKB) (Bycroft et al., 2018; Miller et al., 2016). This large-scale population data presents a unique opportunity to explore the link between various SH dimensions and brain structure/function, overcoming the low reproducibility of previous small sample studies (Arora et al., 2023; Cribb et al., 2023; Ell et al., 2023; Kyle et al., 2017; Li et al., 2022).

The complex interplay between various SH dimensions and brain structure and function has been reported. Sleep disturbance conditions, including insomnia symptoms (Elberse et al., 2024; Holub et al., 2023; Weihs et al., 2023), sleep-disordered breathing (Akradi et al., 2023; André et al., 2020; Mohajer et al., 2020), and abnormal sleep duration, (González et al., 2024; Li et al., 2022) exemplify the inconclusive association between sleep and the brain. Schiel and colleagues using UKB data and Weihs and colleagues using the general population and clinical ENIGMA-Sleep datasets did not observe any strong link between insomnia symptoms/disorder and grey matter volume (GMV) (Schiel et al., 2023; Weihs et al., 2023). However, Stolicyn and colleagues showed that insomnia symptoms are associated with higher global gray and white matter volume, mainly in the amygdala, hippocampus, and putamen (Stolicyn et al., 2023). Moreover, individuals with insomnia symptoms demonstrated altered functional connectivity (FC) within and between the default mode network (DMN), frontoparietal network (FPN), and salience network (SN) (Holub et al., 2023). Neuroimaging meta-analyses found convergent regional abnormalities in the subgenual anterior cingulate cortex (sgACC) in insomnia disorder (Reimann et al., 2023), right basolateral

amygdala/hippocampus and the right central insula in obstructive sleep apnea (OSA) (Tahmasian et al., 2016), the right intraparietal sulcus and superior parietal lobule in acute sleep deprivation (Javaheripour et al., 2019), while the pattern for narcolepsy was inconsistent (Rahimi-Jafari et al., 2022). One study using UKB data found that short sleep duration is linked with lower amygdala reactivity to negative facial expressions (Schiel et al., 2022). The non-linear associations have been documented between sleep duration, cognitive performance, mental health (Li et al., 2022; Tai et al., 2022), and a wide range of regional differences in brain structure, mainly in the subcortical areas (Schiel et al., 2023; Stolicyn et al., 2023; Tsiknia et al., 2023). Fjell and colleagues performed cross-sectional analyses based on the UKB sample, indicating inverse U-shaped relationships between sleep duration and brain structure, i.e., 6.5 hours of sleep was associated with increased cortical thickness and subcortical volumes relative to intracranial volume. However, they failed to identify a longitudinal association between sleep duration and cortical thickness (Fjell, Sørensen, Wang, Amlie, Baaré, Bartrés-Faz, Bertram, et al., 2023). In another study, they found that individuals who reported short sleep without other sleep problems or daytime sleepiness had larger brain volumes compared to both short sleepers with sleep issues and daytime sleepiness, as well as those who slept 7–8 hours (Fjell, Sørensen, Wang, Amlie, Baaré, Bartrés-Faz, Boraxbekk, et al., 2023). An analysis of chronotypes showed that evening chronotype is linked with higher GMV in the precuneus, bilateral nucleus accumbens, caudate, putamen and thalamus, and orbitofrontal cortex (Norbury, 2020). Another study observed the associations between chronotype and neuroimaging phenotypes to be mediated by genetic factors (Williams et al., 2023). Self-reported daytime sleepiness has been reported to be related to higher cortical GMV (Baril et al., 2022). These findings represent an overall inconsistency in the relationship between insomnia, sleep duration, chronotype, and daytime sleepiness with the brain structure and function. The inconclusiveness of these studies may be due to SH being an inhomogeneous concept. Thus, a comprehensive analysis of brain structure and function is crucial for understanding the intricate dynamics of SH and its neuropsychiatric consequences. While these studies employed traditional statistical methods and provided valuable insights into the link between each SH dimension and the brain, they were mostly case-control studies and might not have been able to model the complex interplay between the brain and a complex behavioral phenotype such as SH (Kendler, 2005). Moreover, the large inter-individual variability of SH and the differential associations of SH

dimensions with the brain structure and function measurements, call for more sophisticated computational approaches (Bzdok & Yeo, 2017; Woo et al., 2017).

Machine learning (ML) offers a powerful tool to unravel complex relationships, providing a more nuanced representation than traditional statistical approaches, which is critical in precision medicine (Varoquaux et al., 2017; Vieira et al., 2017). ML models can consider complex multivariate linear and nonlinear relations to make brain-behavior predictions on unseen brain imaging data and have the potential to identify generalizable patterns in SH-related neurobiology at the individual subject level (Afshani et al., 2023; Goldstein-Piekarski et al., 2020; Olfati et al., 2024b), surpassing conventional group comparisons and correlations. In particular, nonlinear models are necessary to capture patterns particularly for sleep duration. Accurate predictive models can contribute to refining our theoretical understanding of the SH-brain relationship. This might pave the way for developing more effective clinical strategies to enable personalized interventions and treatments (Murdoch et al., 2019). Directional genetic analyses using Mendelian randomization demonstrated that altered SH dimensions are more a consequence than a cause of brain abnormalities (Fan et al., 2022). In this study, we employed large-scale neuroimaging data from the UKB, exploring whether and how multimodal brain measurements (i.e., structural markers including GMV, surface-based morphometric features, as well as intrinsic functional imaging markers of local correlation (LCOR), global correlation (GCOR), and fractional amplitude of low-frequency fluctuations (fALFF)) can differentiate between well-separated conditions in each SH-related characteristic (e.g., differentiating individuals usually having insomnia symptoms from individuals without insomnia symptoms).

## Methods

### Participants

We selected the data of the first imaging visit (instance 2) from the UKB (<http://www.ukbiobank.ac.uk>), recorded from 2014 onwards at three different sites in the UK (Cheadle, Reading, Newcastle). The acquisition parameters and protocol of both the structural and functional MRI are as described previously (Miller et al., 2016). We included all individuals who participated in the imaging session, and their data had already been preprocessed and denoised by the UKB team (Alfaro-Almagro et al., 2018). Thus, no particular in-/exclusion criteria have been applied in this sample to be representative of the general

population. A total of 28088 participants, 47% male and 64.1 years old on average (58 – 78 years IQR), were included. The UKB project is approved by the NHS National Research Ethics Service (Ref. 11/NW/0382), and all participants gave written informed consent before participation. Ethical standards are continuously controlled by a Ethics Advisory Committee (EAC, <http://www.ukbiobank.ac.uk/ethics>), based on a project-specific Ethics and Governance Framework (<http://www.ukbiobank.ac.uk/wp-content/uploads/2011/05/EGF20082.pdf>). The current analyses were conducted under UK Biobank application number 41655.

### Sleep health characteristics

The multifaceted definition of SH in UKB is based on previous SH studies (Buysse, 2014; Ell et al., 2023; Goodman et al., 2024; Holub et al., 2023; Schiel et al., 2022, 2023). Accordingly, the seven SH-related characteristics were self-reported insomnia symptoms, sleep duration, difficulty/easiness of getting up in the morning, chronotype, daily nap, daytime sleepiness, and snoring (category 100057), obtained from the touchscreen questionnaire. As these questions were asked at every visit, we selected the responses from the visit matching the neuroimaging acquisition visit.

- Sleeplessness/insomnia field (field 1200): “Do you have trouble falling asleep at night or do you wake up in the middle of the night?”, which could be answered as “never/rarely”, “sometimes”, “usually” or “prefer not to answer”.
- Sleep duration (field 1160): “How many hours sleep do you get in every 24 hours?”.
- Getting up in the morning (field 1170): “On average a day, how easy do you find getting up in the morning?”, with four answers spanning from not at all easy to very easy, as well as “do not know” and “prefer not to answer”.
- Chronotype (i.e., morning/evening person, field 1180): “What do you consider yourself to be?”, with four possible answers spanning from a “morning person” to an “evening person”, as well as “do not know” and “prefer not to answer”.
- Nap during the day (field 1190): “Do you have a nap during the day?”, which can be answered as “never/rarely”, “sometimes”, “usually” or “prefer not to answer”.

- Daytime dozing (field 1220): “How likely are you to doze off or fall asleep during the daytime when you don't mean to? (e.g. when working, reading or driving)”, which can be answered as “never/rarely”, “sometimes”, “often” or “prefer not to answer”.
- Snoring (field 1210): “Does your partner or a close relative or friend complain about your snoring?”, with “yes”, “no”, “do not know” and “prefer not to answer” as possible answers.

Given the ambiguous meaning that some questions, and consequently the respective answers, potentially have in the UKB data (e.g. “sometimes” vs “often”), and to simplify the multiclass/continuous target problems into binary classification problems, we first analyzed the performance of models aimed at distinguishing the extreme answers of each SH-related characteristic. In the case of the continuous answer regarding sleep duration in hours, we split the distribution into four quantiles, selecting the first and fourth quantiles as two classes. However, given the concentration of answers around the median (7 hours), this resulted in discarding only the samples that replied 7 hours. The rationale behind considering the extreme values as class labels is to simplify the classification task, resulting in higher predictive performance if there is indeed a relationship between brain imaging data and each SH-related characteristic. A description of the considered answers for each question, as well as the number of samples for each class, can be seen in Table 1.

Sleep health-related characteristic	Extreme Values			
	Class 0		Class 1	
	Answer(s)	#Samples	Answer(s)	#Samples
<b>Insomnia symptoms</b>	“Never/rarely”	6127	“Usually”	8846
<b>Sleep duration</b>	1 <sup>st</sup> quantile [0-6]	6760	4 <sup>th</sup> quantile [8-16]	9959
<b>Getting up in the morning</b>	“Very Easy”	10687	“Not at all easy” “Not very easy”	3554
<b>Morning/Evening chronotype</b>	“Definitely a ‘morning’ person”	7145	“Definitely an ‘evening’ person”	2398
<b>Daytime nap</b>	“Never/rarely”	15915	“Usually”	1565
<b>Daytime sleepiness*</b>	“Never/rarely”	21406	“Sometimes” “Often”	6458
<b>Snoring*</b>	“Yes”	16439	“No”	9469

**Table 1.** List of answers used for each SH-related characteristic to convert the ambiguous answers into binary classification problems. \*denotes the questions for which no samples were dropped.

## Processing of Imaging data

Grey Matter Volume (GMV): T1-weighted pre-processed images were retrieved from UKB with subsequent computations of voxel-based morphometry (CAT 12.7 (default settings); MNI152 space; 1.5mm isotropic)(Gaser et al., 2023). For each region of interest (ROI), we computed the GMV using the winsorized mean (limits 10%) of the voxel-wise values using the cortical Schaefer atlas (1000 regions of interest, ROIs) (Schaefer et al., 2018), the Melbourne subcortical atlas (S4 3T, 54 ROIs) (Tian et al., 2020), and the Diedrichsen cerebellar atlas (SUIT space, 34 ROIs) (Diedrichsen et al., 2009). This resulted in 1088 GMV features extracted.

Brain Surface: We used the data processed using FreeSurfer 6.0 as provided by the UKB<sup>i</sup>. This includes gray/white matter contrast, pial surface, white matter surface, white matter thickness, and white matter volume from the 68 ROIs of the Desikan-Kiliany parcellation (Desikan et al., 2006), totaling 328 features.

Resting-state Functional Magnetic Resonance Imaging (rsfMRI): The fractional amplitude of low-frequency fluctuations (fALFF) represents the relative measure of blood oxygenation level-dependent (BOLD) magnetic resonance signal power within the low-frequency band of interest (0.008 - 0.09 Hz, reflecting the spontaneous neural activity of the brain) as compared to the BOLD signal power over the entire frequency spectrum (Zou et al., 2008). The LCOR (“local correlation”) is a metric that represents the local coherence for each voxel. It is computed as the average of correlation coefficients between a voxel and a region of neighboring voxels, defined by a 25 mm Gaussian kernel (Deshpande et al., 2009). On the other hand, the GCOR (“global correlation”) represents the node centrality of each voxel and is computed as the average of the correlation coefficients between a voxel and all voxels of the whole brain. These metrics were calculated using MatLab2020b, SPM12 (Friston et al., 2007), FSL (version 5.0) (Jenkinson et al., 2012), and the CONN toolbox (Whitfield-Gabrieli & Nieto-Castanon, 2012). The voxel-wise data was then aggregated parcel-wise by averaging

---

<sup>i</sup> see [https://git.fmrib.ox.ac.uk/falmagro/UK\\_biobank\\_pipeline\\_v\\_1/-/tree/master/bb\\_FS\\_pipeline](https://git.fmrib.ox.ac.uk/falmagro/UK_biobank_pipeline_v_1/-/tree/master/bb_FS_pipeline) for the exact pipeline used.



according to the parcellation of the GMV data (see above), resulting in 1087 features for each metric (fALFF, LCOR, and GCOR), totaling 3261 features derived from rsfMRI. Note that Diedrichsen cerebellar atlas produced 33 features for the fMRI data as for some ROIs, there were not enough voxels to compute the values correctly.

## Samples, Features, and Targets

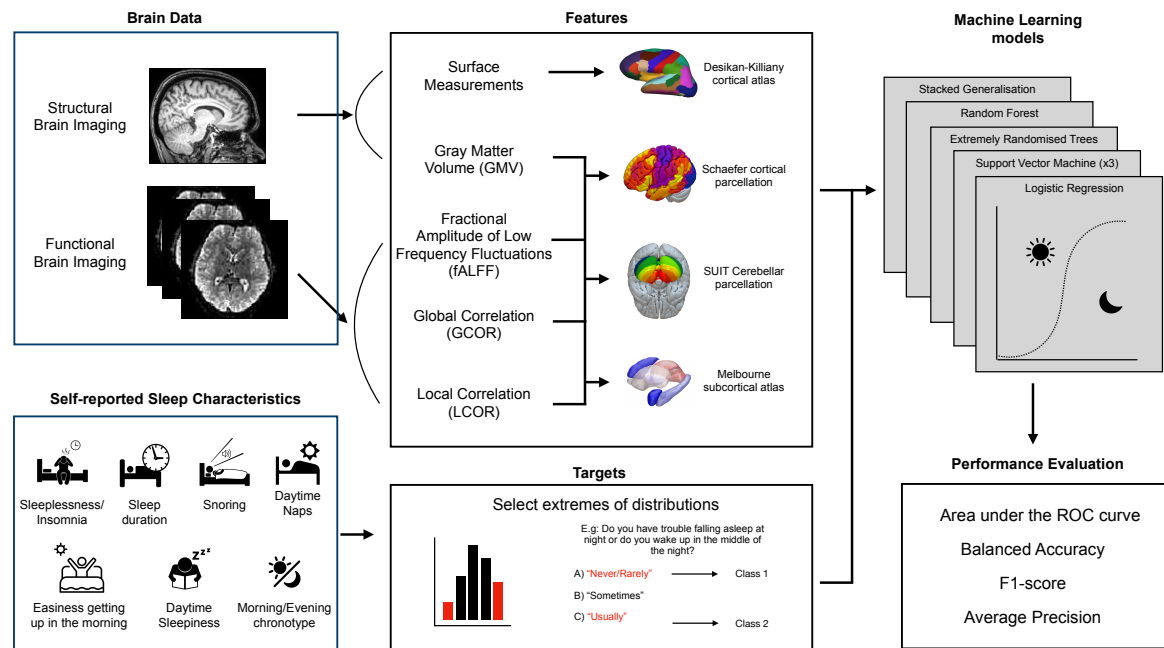
Among the (~500,000 participants in the UKB), we selected the individuals for which all the features were computed, resulting in a total N of 28,088 with a total of 4,677 brain features. The number of variables and samples for each neuroimaging feature is described in Supplementary Table 1.

## ML models

In order to evaluate a broad spectrum of possible interactions between features and relations to the targets, we selected five machine learning algorithms, including parametric and non-parametric models, testing for linear and nonlinear relations. We tested a Random Forest (Breiman, 2001), Extremely Randomized Trees (Extra Trees) (Geurts et al., 2006), Support Vector Machine (SVM) (Cortes & Vapnik, 1995), Logistic Regression (logit), and Stacked Generalization (Wolpert, 1992), with different hyperparameter settings, resulting in seven models. Table 2 summarizes the models, including the hyperparameters tested, except for the Stacked Generalization model, which is described below. When more than one hyperparameter value was listed, the best hyperparameter value was selected using nested cross-validation (CV), using a grid search approach with a stratified 5-fold CV. The Stacked Generalization model consisted of a Linear SVM with heuristic C (*R: Fast Heuristics For The Estimation Of the C Constant Of A...*, n.d.) (model LinearSVMHC) for each type of neuroimaging feature (GMV, Surface, fALFF, GCOR, and LCOR) as the first level. The output of each of these five models were used as features of a second-level logistic regression model. For training the second-level model, the out-of-sample predictions of the first-level models were obtained using a stratified 5-fold CV scheme. An overview of the general methodological approach from brain-images and questionnaires data to the evaluation of ML models is depicted in Figure 1.

#	Name	Learning Algorithm	Hyperparameter	Values
1	GSET	Extra Trees	estimators	200, 500
			Criterion	Gini, entropy, log loss
			Max features	Sqrt, log2
2	GSRF	Random Forest	estimators	200, 500
			Criterion	Gini, entropy, log loss
			Max features	Sqrt, log2
3	GSSVM-RBF	SVM	Kernel	Rbf
			C	1e <sup>-4</sup> , 1e <sup>-3</sup> , 1e <sup>-2</sup> , 1e <sup>-1</sup> , 1, 10, 100, 1e <sup>4</sup> , 1e <sup>5</sup> , 1e <sup>6</sup>
			Gamma	1e <sup>-7</sup> , 1e <sup>-6</sup> , 1e <sup>-5</sup> , 1e <sup>-4</sup> , 1e <sup>-3</sup> , 1e <sup>-2</sup> , 1e <sup>-1</sup> , 1, 10, 100, 1e <sup>4</sup>
4	GSLinearSVM	Linear SVM	C	1e <sup>-4</sup> , 1e <sup>-3</sup> , 1e <sup>-2</sup> , 1e <sup>-1</sup> , 1, 10, 100, 1e <sup>4</sup> , 1e <sup>5</sup> , 1e <sup>6</sup>
5	LinearSVMHC	Linear SVM	C	heuristic(R: <i>Fast Heuristics For The Estimation Of the C Constant Of A...</i> , n.d.)
			dual	False
			penalty	L1
6	LogitHC	Logit	C	heuristic
			dual	False
			penalty	L1

Table 2: List of models tested, including learning algorithms and hyperparameters evaluated.



**Figure 1.** Overview of the methodology. The brain images were processed in order to obtain cortical and subcortical features, both from structural and functional brain imaging. Answers for the UKB questionnaire were binarized by selecting the extremes of the distributions as described in Table 1. We then evaluated the out-of-sample performance of 7 different ML-models, independently for each SH-related characteristic.

## Model evaluation

The available data was first split into 70% training and 30% hold-out test sets to avoid data leakage. Then, the generalization performance of the models (i.e. the capacity to generalize to unseen data) was evaluated on the training set using a stratified 5-fold cross-validation scheme, repeated five times, resulting in 25 evaluation runs. Finally, to validate the CV performance estimation, the models were retrained on the full training set and tested on the hold-out test set. To evaluate different aspects of model performance, such as the trade-off between specificity and sensitivity, we computed four metrics: balanced accuracy, F1 score, area under the receiver-operator characteristic (ROC) curve, and average precision. Balanced accuracy is computed as the relative number of correct predictions over the total samples, weighted by the number of elements in each class so that the chance level is set at 0.5 and 1 would mean a perfect classification. The F1 score is the harmonic mean between precision and recall (Hastie et al., 2001). In short, it measures the model's balanced ability to detect positives (recall = sensitivity) and to have high precision (= positive predictive value), that is, a low rate of false-positive detections. The area under the receiver operating characteristic

(ROC) curve (ROC-AUC) provides an aggregate measure of performance across all possible classification thresholds by plotting the true-positive rate (sensitivity) over the false-positive rate ( $1 - \text{specificity}$ ) for each threshold level. Shortly, ROC-AUC can be interpreted as the probability that given two predictions, the model ranks them in the correct order. A perfect model with sensitivity and specificity being equal to 1 at all threshold levels, will have a ROC-AUC of 1, while random guessing will result in ROC-AUC of 0.5 (Hastie et al., 2001). Given that ROC-AUC is skewed for imbalanced datasets, which is the case for all the SH dimensions (see Table 1), a more suitable metric is the area under the Precision-Recall curve (Davis & Goadrich, 2006), also known as average precision. This metric considers both recall and precision like the F1-score, but across all thresholds as the ROC-AUC does. A perfect model will yield an average precision of 1, while chance levels depend on class balance.

To obtain reference values for each metric, we used the performance of two baseline models, which do not use the features but rely solely on the distribution of classes during training time. A first baseline model named *majority* always predicts the value of the most frequent class in the training set. A second baseline model named *chance* draws random predictions weighted by the number of training samples in each class. All models for each SH dimension were evaluated using the same 5 x 5 CV folds. We then used the corrected paired Student's t-test for comparing the CV performance of the machine learning models (Nadeau & Bengio, 2003) and corrected for multiple comparisons (across models) using the Bonferroni method. All the analysis described was implemented using Julearn (Hamdan et al., 2024) and Scikit-learn (Pedregosa et al., 2012). The codes are available on GitHub: [https://github.com/juaml/ukb\\_sleep\\_prediction](https://github.com/juaml/ukb_sleep_prediction).

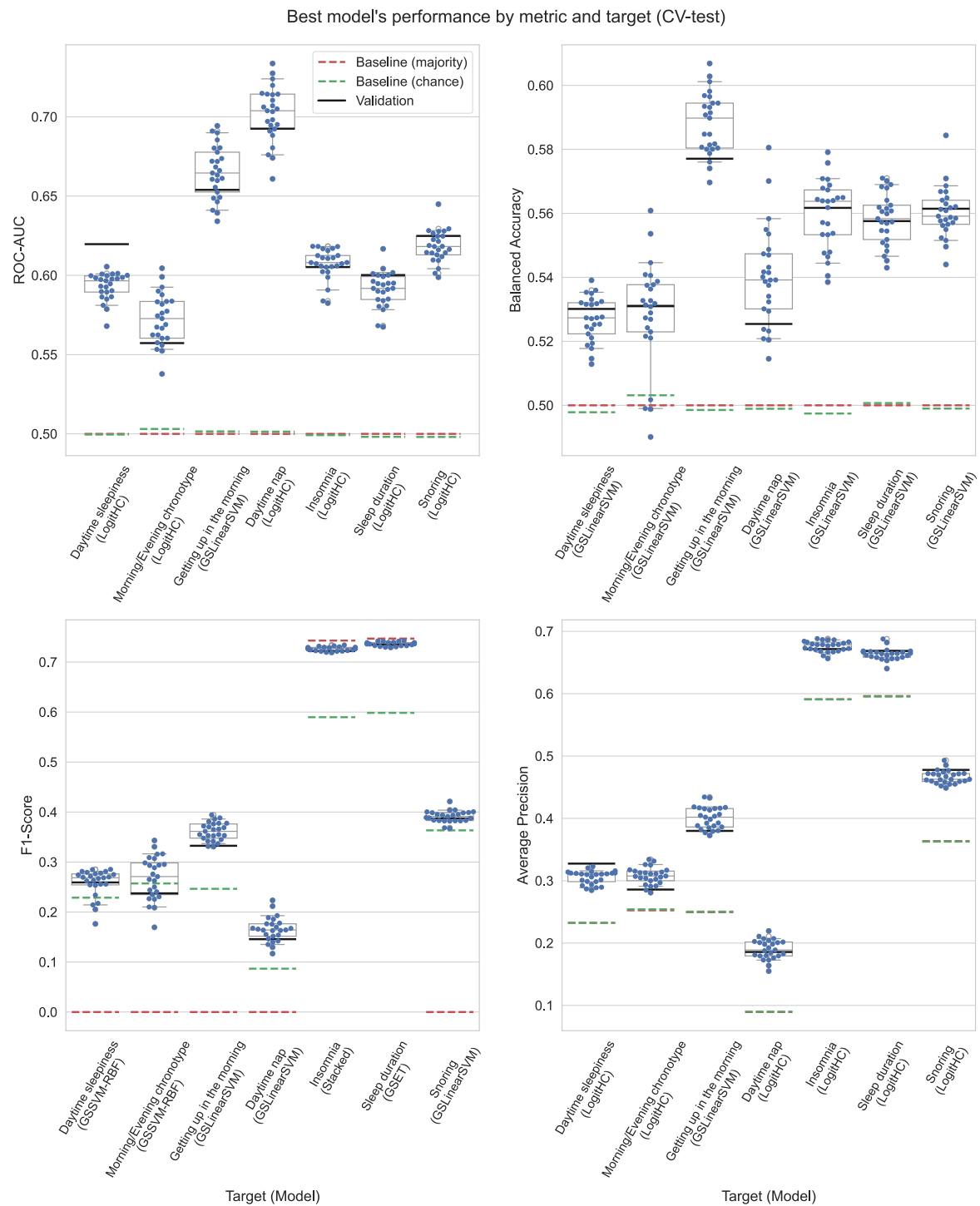
## Results

We first trained and evaluated all seven models for each of the seven SH-related characteristics, a procedure that took 119314 core-hours, which is approximately 13 years on a single-core processor or eight years on an 8-core desktop computer. For each SH-related characteristic and metric, we selected the best model among the seven competing models according to the performance of the respective metric upon evaluation on  $5 \times 5 = 25$  CV-folds. This resulted in one model per SH-related characteristic and metric, which were then applied

to the 30% hold-out test set. The performance of the best model for each SH-related characteristic and metric can be seen in Figure 2. A complete description of the estimated performances for each metric can be seen in Supplementary Table 2.

When only considering the CV performance (which is commonly reported in research settings), some of the SH-related characteristics showed a modest predictability on several metrics. For instance, the best models for *insomnia* and *sleep duration* showed modest balanced accuracy (0.588 and 0.584) and AUC-ROC (0.549 and 0.553) and relatively high F1-score (0.725 and 0.739) and average precision (0.664 and 0.658). However, since some SH-related characteristic have imbalanced classes, it is important to note the performance of the baseline models. For example, the F1 score for insomnia and sleep duration is below the performance of the *majority baseline* model, meaning that a model that simply assigns the majority class to each sample showed a better F1 score. The limitation of AUC-ROC with imbalanced data also becomes clear for the *easiness getting up* characteristic, which showed a relatively high AUC-ROC but relatively lower average precision. Furthermore, as cross-validated performances could be overestimated (Varoquaux, 2017), we evaluated the models on the hold-out data (30% of the samples). The obtained results fall within the confidence intervals of the CV-estimated performances (black lines in Figure 1), suggesting that no over-estimation happened in our case. For more details on the values obtained for each model and SH-related characteristic, see Supplementary Table 3. Overall, our results indicate a weak predictive power but systematically above baseline models for each of the seven SH-related characteristics.

A common ML pitfall with a lack of predictive power is *overfitting*. This occurs when the model closely learns the idiosyncrasies of the training data, thus being incapable of making correct predictions on new unseen samples. To verify that this is not the case, we computed the same metrics for each model but on the training samples. That is, how well each model memorized the training data. The results indicate that while some models were indeed overfitting, at least one model per SH-related characteristic was not (Supplementary Table 4). Given the comparable out-of-sample performance across models for each SH-related characteristic, and that the hyperparameters were selected in nested CV to prevent overfitting, we can safely conclude that overfitting is not a major issue in our results.



**Figure 2.** Performances of the best model for each SH-related characteristic. Each blue dot represents the performance obtained at each of the 25 test folds within cross-validation (CV). Boxplots summarize the medians and 95% CI for the underlying distribution. As a reference, dashed red lines depict the mean performance of a model that constantly predicts the most frequent class, green lines depict the mean performance of a model that draws random predictions weighted by the number of samples in each class, and black lines indicate the performance on the hold-out (validation) data.

## Discussion

The current large-scale study systematically evaluated ML-based predictive analysis for classifying extremes of seven different SH-related characteristics based on multivariate neuroimaging markers in UKB. Notably, we covered a large space of multimodal neuroimaging features covering brain structure and function, several ML algorithms in a nested cross-validation setting, and a hold-out test set evaluated on four metrics. Our striking findings demonstrated that the balanced accuracy for predicting SH-related characteristics did not exceed 56%, which indicates that brain structure and function measures could not accurately predict any of the seven SH-related characteristics. The slight improvement over baseline models across the evaluation metrics suggests that the ML algorithms indeed captured some underlying patterns in the data. However, we do not consider these results as high predictive accuracy compared to other brain-imaging-based predictions, such as sex (Schulz et al., 2024; Wiersch et al., 2023), neurodegenerative diseases (Kasper et al., 2023), and depressive symptoms severity (Olfati et al., 2024b). Put differently, we did not observe sufficiently strong evidence to claim that the brain measures can predict SH-related characteristics. Our findings raise the question of whether the results imply a true absence of a strong relationship between the SH-related characteristics and the brain imaging features, and if so, then what are the main sources of such results? In the following, we discuss the potential reasons for the poor efficacy of multimodal brain features in predicting SH-related characteristics.

### Target issues: Sleep health is a heterogeneous concept

Our findings align with previous large-scale sample studies using e.g., UKB and ENIGMA-Sleep datasets that did not observe an association between brain structure and insomnia symptoms (Schiell et al., 2023; Weihs et al., 2023) and sleep duration (Fjell, Sørensen, Wang, Amlie, Baaré, Bartrés-Faz, Bertram, et al., 2023). SH has a heterogeneous definition across different general population datasets, as well as clinical samples. Although some studies used a standard sleep questionnaire such as the Pittsburgh Sleep Quality Index (PSQI) to assess sleep quality or the Regulatory Satisfaction Alertness Timing Efficiency Duration (RU-SATED) questionnaire as a valid measure of SH (Ravyts et al., 2021), the UKB did not use such a standard questionnaire and instead provided seven self-reported questions about sleep duration, difficulties in getting up in the morning, chronotype, nap, daytime sleepiness, and

two measures of clinical conditions such as insomnia symptoms and snoring. SH is a complex concept and considering these single questions could have affected the clarity and meaningfulness of measured SH. Furthermore, accuracy of self-report sleep assessment based on seven single items and selective participation biases could have led to measurement issues, which have been highlighted previously (Schoeler et al., 2023). Taken together, these may have resulted in the observed low prediction performance when using brain imaging features to predict SH by ML models.

Another critical aspect is differentiating the sleep-related symptoms of insomnia and snoring we considered here from clinical conditions. It is well-documented that Insomnia disorder is a heterogeneous condition with different subtypes with noticeable inconsistencies in terms of pathophysiology, symptomatology, and treatment response (Blanken et al., 2019; Bresser et al., 2024; Emamian et al., 2021; Holub et al., 2023; Reimann et al., 2023; Schiel et al., 2022; Weihs et al., 2023). According to the third edition of the International Classification of Sleep Disorders (ICSD-3) (Sateia, 2014), significant daytime dysfunction and having adequate opportunity and circumstances to sleep are essential diagnostic criteria for insomnia disorder. Similarly, snoring can have several etiologies beyond it being a cardinal symptom of OSA, including genetic factors, obesity, nasal blockages, alcohol abuse, smoking, or medications (Campos et al., 2020). Thus, relying on a single question about sleep problems is not sufficient to define clinical insomnia disorder or OSA.

Additionally, the imbalance in target labels significantly influences model performance, hindering the learning of sufficient information for accurate classification. Particularly for SH-related characteristic such as 'easiness getting up in the morning,' 'Day Naps,' and 'Daytime Dozing,' the uneven distribution of target labels has resulted in models achieving moderate ROC-AUC scores around 0.6 while the balanced accuracy remained at chance level of approximately 0.5. This discrepancy between ROC-AUC and balanced accuracy highlights the challenges in achieving fairness and robustness in the models' predictive capabilities when dealing with imbalanced target datasets.

### Input features issues: regional brain measurements

Our results also suggest that the neuroimaging features applied in our study may not capture the full spectrum of brain-related features relevant to SH or that the selected features may



not be sensitive enough to the subtleties of SH. Moreover, it raises the possibility that current feature sets are insufficiently granular to mirror the complex biological underpinnings of SH. The low performance of the models in predicting SH dimensions, therefore, points to the need for a deeper investigation into more sensitive and comprehensive neuroimaging metrics that can better encapsulate the factors influencing SH. SH might be associated with brain circuits that can be captured, e.g., via seed-based structural or FC measures rather than local brain abnormalities that we used from brain parcels, including GMV, gray/white matter contrast, pial surface, white matter surface, white matter thickness, and white matter volume, LCOR, and fALFF. It has been reported that insomnia symptoms were associated with higher FC within the DMN and FPN and lower FC between the DMN and SN (Holub et al., 2023). Wang and colleagues also found that SH dimensions are correlated with disrupted FC patterns in the attentional and thalamic networks in several datasets (Wang et al., 2023). Another study using UKB data found associations between SH and FC and structural connectivity. Within-network hyperconnectivity in DMN, FPN, and SN has been observed in healthy subjects and patients with mild cognitive impairment with insomnia symptoms, while patients with Alzheimer's disease and insomnia symptoms showed hypoconnectivity in those networks (Elberse et al., 2024). Although we included GCOR, representing functional correlations between a given voxel and other brain voxels (i.e., degree centrality), it didn't improve the prediction when used as input with local markers together. Thus, future studies could explore network-based and white matter integrity metrics as input features to predict SH in UKB.

Our results also remind us to think beyond the brain feature modalities. Recently, we observed that sleep quality and anxiety robustly predict depressive symptoms severity across three independent datasets. Still, brain structural and functional features could not predict depressive symptoms, which indicated that brain imaging data may not be very helpful in predicting mental health (Olfati et al., 2024a). A large-scale study by the ENIGMA-Anxiety Consortium utilized ML to analyze neuroanatomical data for youth anxiety disorders, and also achieved only modest classification accuracy (AUC 0.59–0.63) (Bruin et al., 2024). This parallels findings from extensive ML optimization efforts with major depressive disorder (MDD), which observed mean accuracies in distinguishing patients from controls that ranged from 48.1% to 62.0% only, even when additionally provided with polygenic risk scores,, casting doubt on the potential diagnostic relevance of neuroimaging and genetic biomarkers for MDD (Winter et al., 2024). Similarly, the ENIGMA-MDD consortium's multi-site study

(Belov et al., 2024) achieved a balanced accuracy of only about 62% in classifying MDD versus healthy controls, which further dropped to approximately 52% after harmonization for site effect. Random chance accuracy was also observed across various stratified groups. These findings may point to an alternative view that complex psychiatric conditions such as sleep disturbance or depression represent deficits in the brain-body interaction, which suggests that body organ health measurements such as metabolic and cardiovascular systems, in addition to brain imaging, should be considered (Kendler, 2024; Tian et al., 2023).

## ML-related issues

Following proper ML pipelining practices such as nested CV and grid search for meticulous hyperparameter tuning—methods that typically enhance a model's capacity to generalize—our models did not achieve high predictive performance. The low classification performance in our study highlights the challenges inherent in developing models that accurately capture the complex nature of SH using brain imaging data. Machine learning models are designed to discern patterns and generalize findings to new, unseen data. However, like any statistical analysis, ML is challenged when the target labels are unreliable (Gell et al., 2024). We reduced the uncertainty in the labeling to some degree by using extreme values for each SH-related characteristic. This should make learning easier for the ML algorithms and boost accuracy. The low performance observed despite this simplification suggests that the prediction of SH-related characteristic as a continuum could be more challenging. Difficulty in creating generalizable ML models arises from potential heterogeneity in how SH is reflected in the brain. In this case, the ML models will not be able to learn a consistent pattern, leading to low performance. Further analysis of SH subtypes and more refined scales are needed to discern this possibility. Finally, several of our classification tasks were imbalanced, i.e. one of the classes was much more frequently present than the other. Such imbalance can lead to biased ML models, which in turn lack generalization ability. To this end, we employed AUC-ROC and average precision metrics for evaluating the ML pipelines. These metrics are independent of a threshold used for dichotomization and thus suitable for characterizing the performance in imbalanced datasets particularly with tree ensemble models (Collell et al., 2018).

## Strengths, limitations, and future directions

The present study has several advantages over other case-control SH-brain studies. Here, we calculated 4677 structural and functional brain features as input features from 28,088 participants from the UKB and applied several ML algorithms to classify the extremes of seven SH-related characteristics. In particular, 1) including diverse and multimodal neuroimaging metrics is crucial. Multimodal data enriches the ML analysis, allowing for a more comprehensive exploration and interpretation of the neurobiological correlates of SH at both structural and functional levels; 2) we leveraged the detailed features provided by the Schaefer atlas (1000 ROIs), which is supported by our ample sample size. This approach assumes that if relevant information is present in an ROI, our models—given their complexity—are equipped to detect it, whether the information is concentrated within a single ROI or dispersed across several regions; 3) we carefully designed our ML analyses using fully separated train and test samples to avoid any leakage of the test set into the model, which is a common oversight in some ML studies (Sasse et al., 2024); 4) the ML analyses were conducted using several rather different algorithms including Random Forest, Extremely Randomized Trees, Support Vector Machine, Logistic Regression, and Stacked generalization; 5) we applied a grid search-based hyperparameter optimization to prevent overfitting and increase the generalizability of our findings.

Our results should be interpreted within the context of the study's limitations and the nascent state of this field. This study is based on seven proband answers to SH-related questions and did not include any objective sleep assessment such as polysomnography. Although polysomnography is recommended as a gold-standard objective measure for the diagnosis of several sleep disorders, including obstructive sleep apnea, its validity for insomnia or sleep quality assessment remains disputed (Frase et al., 2023). Moreover, some evidence showed only weak association between the subjective sleep measurement (e.g., PSQI) and polysomnography in patients with insomnia disorder (Benz et al., 2023). Here, we focused on self-reported information on SH. Thus, future studies should consider performing an ML analysis of objective sleep data and comparing it with the analysis of subjective data. In addition, the SH-related characteristics in the UKB sample do not represent cross-country sleep differences well. Recently, data from 63 countries showed that individuals from East Asia tend to sleep less and participants from East Europe report longer sleep duration (Coutrot et al., 2022). Similarly, another study on ~220,000 wearable device users in 35

Countries observed shorter sleep duration, later sleep timing, and less sleep efficiency in East Asia compared with Western Europe, North America, and Oceania, probably due to social- and work-related cultural differences regarding the coping with inadequate sleep and sleep debt (Willoughby et al., 2023). In addition, there are significant differences in daytime napping across cultures, being more common in non-Western countries (Willoughby et al., 2023). Of note, however, 10% of the UKB participants reported regular daily naps (Table 1).

Future studies could apply normative modelling, a technique that studies deviations from population norms to show the range of inter-individual differences in brain structure. Unlike traditional case-control paradigms that rely on common neurobiological factors across all subjects, normative modelling focuses on individual deviations from normal patterns, making it a promising approach to consider inter-individual variability in brain expression of SH (Marquand et al., 2016; Rutherford et al., 2022). Furthermore, longitudinal studies can help identify the long-term interaction between the SH and the brain together with well-characterized sleep measurements from collaborative research groups e.g., the ENIGMA-Sleep consortium (Tahmasian et al., 2021) to provide replicable results.

## Conclusion

The present extensive ML study using a large population sample demonstrated that multimodal neuroimaging markers had low efficacy in separating the extremes of various sleep health-related characteristics UKB. This suggests that the interaction between sleep health and brain organization may be more complex to be captured with the current ML models and neuroimaging features. While our methodological approach is comprehensive and aims to establish links between neuroimaging features and SH dimensions, this study acknowledges the complexity of interpreting neuroimaging in the context of sleep health. We need future cross-sectional and longitudinal studies considering brain circuits, objective sleep measurements, and cross-country sleep assessments to evaluate the sophisticated brain-sleep interplay.

# References

- Afshani, M., Mahmoudi-Aznaveh, A., Noori, K., Rostampour, M., Zarei, M., Spiegelhalter, K., Khazaie, H., & Tahmasian, M. (2023). Discriminating Paradoxical and Psychophysiological Insomnia Based on Structural and Functional Brain Images: A Preliminary Machine Learning Study. *Brain Sciences*, 13(4), Article 4. <https://doi.org/10.3390/brainsci13040672>
- Akradi, M., Farzane-Daghigh, T., Ebneabbasi, A., Bi, H., Drzezga, A., Mander, B. A., Eickhoff, S. B., Tahmasian, M., & Initiative, the A. D. N. (2023). *How is sleep-disordered breathing linked with biomarkers of Alzheimer's disease?* (p. 2023.08.16.23294054). medRxiv. <https://doi.org/10.1101/2023.08.16.23294054>
- Alfaro-Almagro, F., Jenkinson, M., Bangerter, N. K., Andersson, J. L. R., Griffanti, L., Douaud, G., Sotiropoulos, S. N., Jbabdi, S., Hernandez-Fernandez, M., Vallee, E., Vidaurre, D., Webster, M., McCarthy, P., Rorden, C., Daducci, A., Alexander, D. C., Zhang, H., Dragonu, I., Matthews, P. M., ... Smith, S. M. (2018). Image processing and Quality Control for the first 10,000 brain imaging datasets from UK Biobank. *NeuroImage*, 166, 400–424. <https://doi.org/10.1016/j.neuroimage.2017.10.034>
- André, C., Rehel, S., Kuhn, E., Landeau, B., Moulinet, I., Touron, E., Ourry, V., Le Du, G., Mézenge, F., Tomadesso, C., de Flores, R., Bejanin, A., Sherif, S., Delcroix, N., Manrique, A., Abbas, A., Marchant, N. L., Lutz, A., Klimecki, O. M., ... for the Medit-Ageing Research Group. (2020). Association of Sleep-Disordered Breathing With Alzheimer Disease Biomarkers in Community-Dwelling Older Adults: A Secondary Analysis of a Randomized Clinical Trial. *JAMA Neurology*, 77(6), 716–724. <https://doi.org/10.1001/jamaneurol.2020.0311>
- Arora, N., Richmond, R. C., Brumpton, B. M., Åsvold, B. O., Dalen, H., Skarpsno, E. S., & Strand, L. B. (2023). Self-reported insomnia symptoms, sleep duration, chronotype and the risk of acute myocardial infarction (AMI): A prospective study in the UK Biobank and the HUNT Study. *European Journal of Epidemiology*, 38(6), 643–656. <https://doi.org/10.1007/s10654-023-00981-x>
- Baril, A.-A., Beiser, A. S., DeCarli, C., Himali, D., Sanchez, E., Cavuoto, M., Redline, S., Gottlieb, D. J., Seshadri, S., Pase, M. P., & Himali, J. J. (2022). Self-reported sleepiness associates with greater brain and cortical volume and lower prevalence of ischemic covert brain infarcts in a community sample. *Sleep*, 45(10), zsac185. <https://doi.org/10.1093/sleep/zsac185>
- Belov, V., Erwin-Grabner, T., Aghajani, M., Aleman, A., Amod, A. R., Basgoze, Z., Benedetti, F., Besteher, B., Bülow, R., Ching, C. R. K., Connolly, C. G., Cullen, K., Davey, C. G., Dima, D., Dols, A., Evans, J. W., Fu, C. H. Y., Gonul, A. S., Gotlib, I. H., ... Goya-Maldonado, R. (2024). Multi-site benchmark classification of major depressive disorder using machine learning on cortical and subcortical measures. *Scientific Reports*, 14(1), 1084. <https://doi.org/10.1038/s41598-023-47934-8>

- Benz, F., Riemann, D., Domschke, K., Spiegelhalder, K., Johann, A. F., Marshall, N. S., & Feige, B. (2023). How many hours do you sleep? A comparison of subjective and objective sleep duration measures in a sample of insomnia patients and good sleepers. *Journal of Sleep Research*, 32(2), e13802. <https://doi.org/10.1111/jsr.13802>
- Blanken, T. F., Benjamins, J. S., Borsboom, D., Vermunt, J. K., Paquola, C., Ramautar, J., Dekker, K., Stoffers, D., Wassing, R., Wei, Y., & Someren, E. J. W. V. (2019). Insomnia disorder subtypes derived from life history and traits of affect and personality. *The Lancet Psychiatry*, 6(2), 151–163. [https://doi.org/10.1016/S2215-0366\(18\)30464-4](https://doi.org/10.1016/S2215-0366(18)30464-4)
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- Bresser, T., Blanken, T. F., de Lange, S. C., Leerssen, J., Foster-Dingley, J. C., Lakbila-Kamal, O., Wassing, R., Ramautar, J. R., Stoffers, D., van de Heuvel, M. P., & van Someren, E. J. W. (2024). Insomnia subtypes have differentiating deviations in brain structural connectivity. *Biological Psychiatry*. <https://doi.org/10.1016/j.biopsych.2024.06.014>
- Bruin, W. B., Zhutovsky, P., van Wingen, G. A., Bas-Hoogendam, J. M., Groenewold, N. A., Hilbert, K., Winkler, A. M., Zugman, A., Agosta, F., Åhs, F., Andreescu, C., Antonacci, C., Asami, T., Assaf, M., Barber, J. P., Bauer, J., Bavdekar, S. Y., Beesdo-Baum, K., Benedetti, F., ... Aghajani, M. (2024). Brain-based classification of youth with anxiety disorders: Transdiagnostic examinations within the ENIGMA-Anxiety database using machine learning. *Nature Mental Health*, 2(1), 104–118. <https://doi.org/10.1038/s44220-023-00173-2>
- Buyse, D. J. (2014). Sleep Health: Can We Define It? Does It Matter? *Sleep*, 37(1), 9–17. <https://doi.org/10.5665/sleep.3298>
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O'Connell, J., Cortes, A., Welsh, S., Young, A., Effingham, M., McVean, G., Leslie, S., Allen, N., Donnelly, P., & Marchini, J. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726), 203–209. <https://doi.org/10.1038/s41586-018-0579-z>
- Bzdok, D., & Yeo, B. T. T. (2017). Inference in the age of big data: Future perspectives on neuroscience. *NeuroImage*, 155, 549–564. <https://doi.org/10.1016/j.neuroimage.2017.04.061>
- Campos, A. I., García-Marín, L. M., Byrne, E. M., Martin, N. G., Cuéllar-Partida, G., & Rentería, M. E. (2020). Insights into the aetiology of snoring from observational and genetic investigations in the UK Biobank. *Nature Communications*, 11(1), 817. <https://doi.org/10.1038/s41467-020-14625-1>
- Cheng, W., Rolls, E. T., Ruan, H., & Feng, J. (2018). Functional Connectivities in the Brain That Mediate the Association Between Depressive Problems and Sleep Quality. *JAMA Psychiatry*, 75(10), 1052. <https://doi.org/10.1001/jamapsychiatry.2018.1941>



- Collell, G., Prelec, D., & Patil, K. R. (2018). A simple plug-in bagging ensemble based on threshold-moving for classifying binary and multiclass imbalanced data. *Neurocomputing*, 275, 330–340. <https://doi.org/10.1016/j.neucom.2017.08.035>
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273–297. <https://doi.org/10.1007/BF00994018>
- Coutrot, A., Lazar, A. S., Richards, M., Manley, E., Wiener, J. M., Dalton, R. C., Hornberger, M., & Spiers, H. J. (2022). Reported sleep duration reveals segmentation of the adult life-course into three phases. *Nature Communications*, 13(1), 7697. <https://doi.org/10.1038/s41467-022-34624-8>
- Cribb, L., Sha, R., Yiallourou, S., Grima, N. A., Cavuoto, M., Baril, A.-A., & Pase, M. P. (2023). *Sleep Regularity and Mortality: A Prospective Analysis in the UK Biobank* (p. 2023.04.14.23288550). medRxiv. <https://doi.org/10.1101/2023.04.14.23288550>
- Davis, J., & Goadrich, M. (2006). The relationship between Precision-Recall and ROC curves. *Proceedings of the 23rd International Conference on Machine Learning*, 233–240. <https://doi.org/10.1145/1143844.1143874>
- Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R. L., Dale, A. M., Maguire, R. P., Hyman, B. T., Albert, M. S., & Killiany, R. J. (2006). An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *NeuroImage*, 31(3), 968–980. <https://doi.org/10.1016/j.neuroimage.2006.01.021>
- Diedrichsen, J., Balsters, J. H., Flavell, J., Cussans, E., & Ramnani, N. (2009). A probabilistic MR atlas of the human cerebellum. *NeuroImage*, 46(1), 39–46. <https://doi.org/10.1016/j.neuroimage.2009.01.045>
- Elberse, J. D., Saberi, A., Ahmadi, R., Changizi, M., Bi, H., Hoffstaedter, F., Mander, B. A., Eickhoff, S. B., Tahmasian, M., & Alzheimer’s Disease Neuroimaging Initiative. (2024). The interplay between insomnia symptoms and Alzheimer’s disease across three main brain networks. *Sleep*, zsae145. <https://doi.org/10.1093/sleep/zsae145>
- Ell, J., Schiel, J. E., Feige, B., Riemann, D., Nyhuis, C. C., Fernandez-Mendoza, J., Vetter, C., Rutter, M. K., Kyle, S. D., & Spiegelhalter, K. (2023). Sleep health dimensions and shift work as longitudinal predictors of cognitive performance in the UK Biobank cohort. *SLEEP*, 46(6), zsad093. <https://doi.org/10.1093/sleep/zsad093>
- Emamian, F., Mahdipour, M., Noori, K., Rostampour, M., Mousavi, S. B., Khazaie, H., Khodaie-Ardakani, M., Tahmasian, M., & Zarei, M. (2021). Alterations of Subcortical Brain Structures in Paradoxical and Psychophysiological Insomnia Disorder. *Frontiers in Psychiatry*, 12. <https://doi.org/10.3389/fpsy.2021.661286>
- Fan, Z., Li, Y., Shu, J., Yang, X., Li, B., Lin, J., Wang, Q., Paschou, P., Li, T., Zhu, H., & Zhao, B. (2022). *Mapping sleep’s phenotypic and genetic links to the brain and heart: A systematic analysis of multimodal brain and cardiac images in the UK Biobank* (p. 2022.09.08.22279719). medRxiv. <https://doi.org/10.1101/2022.09.08.22279719>

- Fjell, A. M., Sørensen, Ø., Wang, Y., Amlien, I. K., Baaré, W. F. C., Bartrés-Faz, D., Bertram, L., Boraxbekk, C.-J., Brandmaier, A. M., Demuth, I., Drevon, C. A., Ebmeier, K. P., Ghisletta, P., Kievit, R., Kühn, S., Madsen, K. S., Mowinckel, A. M., Nyberg, L., Sexton, C. E., ... Walhovd, K. B. (2023). No phenotypic or genotypic evidence for a link between sleep duration and brain atrophy. *Nature Human Behaviour*. <https://doi.org/10.1038/s41562-023-01707-5>
- Fjell, A. M., Sørensen, Ø., Wang, Y., Amlien, I. K., Baaré, W. F. C., Bartrés-Faz, D., Boraxbekk, C.-J., Brandmaier, A. M., Demuth, I., Drevon, C. A., Ebmeier, K. P., Ghisletta, P., Kievit, R., Kühn, S., Madsen, K. S., Nyberg, L., Solé-Padullés, C., Vidal-Piñeiro, D., Wagner, G., ... Walhovd, K. B. (2023). Is Short Sleep Bad for the Brain? Brain Structure and Cognitive Function in Short Sleepers. *Journal of Neuroscience*, 43(28), 5241–5250. <https://doi.org/10.1523/JNEUROSCI.2330-22.2023>
- Frase, L., Nissen, C., Spiegelhalder, K., & Feige, B. (2023). The importance and limitations of polysomnography in insomnia disorder—A critical appraisal. *Journal of Sleep Research*, 32(6), e14036. <https://doi.org/10.1111/jsr.14036>
- Friston, K. J., Ashburner, J., Kiebel, S., Nichols, T., & Penny, W. (2007). *Statistical parametric mapping: The analysis of functional brain images* (1st ed). Elsevier / Academic Press.
- Gaser, C., Dahnke, R., Thompson, P. M., Kurth, F., Luders, E., & Initiative, A. D. N. (2023). CAT – A Computational Anatomy Toolbox for the Analysis of Structural MRI Data (p. 2022.06.11.495736). bioRxiv. <https://doi.org/10.1101/2022.06.11.495736>
- Gell, M., Eickhoff, S. B., Omidvarnia, A., Küppers, V., Patil, K. R., Satterthwaite, T. D., Müller, V. I., & Langner, R. (2024). *The Burden of Reliability: How Measurement Noise Limits Brain-Behaviour Predictions* (p. 2023.02.09.527898). bioRxiv. <https://doi.org/10.1101/2023.02.09.527898>
- Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine Learning*, 63(1), 3–42. <https://doi.org/10.1007/s10994-006-6226-1>
- Goldstein-Piekarski, A. N., Holt-Gosselin, B., O’Hora, K., & Williams, L. M. (2020). Integrating sleep, neuroimaging, and computational approaches for precision psychiatry. *Neuropsychopharmacology*, 45(1), 192–204. <https://doi.org/10.1038/s41386-019-0483-8>
- González, K. A., Tarraf, W., Stickel, A. M., Kaur, S., Agudelo, C., Redline, S., Gallo, L. C., Isasi, C. R., Cai, J., Daviglus, M. L., Testai, F. D., DeCarli, C., González, H. M., & Ramos, A. R. (2024). Sleep duration and brain MRI measures: Results from the SOL-INCA MRI study. *Alzheimer’s & Dementia*, 20(1), 641–651. <https://doi.org/10.1002/alz.13451>
- Goodman, M. O., Faquih, T., Paz, V., Nagarajan, P., Lane, J. M., Spitzer, B., Maher, M., Chung, J., Cade, B. E., Purcell, S. M., Zhu, X., Noordam, R., Phillips, A. J. K., Kyle, S. D., Spiegelhalder, K., Weedon, M. N., Lawlor, D. A., Rotter, J. I., Taylor, K. D., ... Wang, H. (2024). *Genome-wide association analysis of composite sleep health scores in 413,904 individuals* (p. 2024.02.02.24302211). medRxiv. <https://doi.org/10.1101/2024.02.02.24302211>



- Hamdan, S., More, S., Sasse, L., Komeyer, V., Patil, K. R., Raimondo, F., Initiative, for the A. D. N., More, S., Sasse, L., Komeyer, V., Patil, K. R., Raimondo, F., & Initiative, for the A. D. N. (2024). Julearn: An easy-to-use library for leakage-free evaluation and inspection of ML models. *Gigabyte*, 2024, 1–16. <https://doi.org/10.46471/gigabyte.113>
- Hastie, T., Friedman, J., & Tibshirani, R. (2001). *The Elements of Statistical Learning*. Springer New York. <https://doi.org/10.1007/978-0-387-21606-5>
- Holub, F., Petri, R., Schiel, J., Feige, B., Rutter, M. K., Tamm, S., Riemann, D., Kyle, S. D., & Spiegelhalder, K. (2023). Associations between insomnia symptoms and functional connectivity in the UK Biobank cohort (n = 29,423). *Journal of Sleep Research*, 32(2), e13790. <https://doi.org/10.1111/jsr.13790>
- Javaheripour, N., Shahdipour, N., Noori, K., Zarei, M., Camilleri, J. A., Laird, A. R., Fox, P. T., Eickhoff, S. B., Eickhoff, C. R., Rosenzweig, I., Khazaie, H., & Tahmasian, M. (2019). Functional brain alterations in acute sleep deprivation: An activation likelihood estimation meta-analysis. *Sleep Medicine Reviews*, 46, 64–73. <https://doi.org/10.1016/j.smr.2019.03.008>
- Jenkinson, M., Beckmann, C. F., Behrens, T. E. J., Woolrich, M. W., & Smith, S. M. (2012). FSL. *NeuroImage*, 62(2), 782–790. <https://doi.org/10.1016/j.neuroimage.2011.09.015>
- Kasper, J., Eickhoff, S. B., Caspers, S., Peter, J., Dogan, I., Wolf, R. C., Reetz, K., Dukart, J., & Orth, M. (2023). Local synchronicity in dopamine-rich caudate nucleus influences Huntington’s disease motor phenotype. *Brain*, 146(8), 3319–3330. <https://doi.org/10.1093/brain/awad043>
- Kendler, K. S. (2005). Toward a Philosophical Structure for Psychiatry. *American Journal of Psychiatry*, 162(3), 433–440. <https://doi.org/10.1176/appi.ajp.162.3.433>
- Kendler, K. S. (2024). Are Psychiatric Disorders Brain Diseases?—A New Look at an Old Question. *JAMA Psychiatry*. <https://doi.org/10.1001/jamapsychiatry.2024.0036>
- Kyle, S. D., Sexton, C. E., Feige, B., Luik, A. I., Lane, J., Saxena, R., Anderson, S. G., Bechtold, D. A., Dixon, W., Little, M. A., Ray, D., Riemann, D., Espie, C. A., Rutter, M. K., & Spiegelhalder, K. (2017). Sleep and cognitive performance: Cross-sectional associations in the UK Biobank. *Sleep Medicine*, 38, 85–91. <https://doi.org/10.1016/j.sleep.2017.07.001>
- Li, Y., Sahakian, B. J., Kang, J., Langley, C., Zhang, W., Xie, C., Xiang, S., Yu, J., Cheng, W., & Feng, J. (2022). The brain structure and genetic mechanisms underlying the nonlinear association between sleep duration, cognition and mental health. *Nature Aging*, 2(5), Article 5. <https://doi.org/10.1038/s43587-022-00210-2>
- Marquand, A. F., Rezek, I., Buitelaar, J., & Beckmann, C. F. (2016). Understanding Heterogeneity in Clinical Cohorts Using Normative Models: Beyond Case-Control Studies. *Biological Psychiatry*, 80(7), 552–561. <https://doi.org/10.1016/j.biopsych.2015.12.023>

- Miller, K. L., Alfaro-Almagro, F., Bangerter, N. K., Thomas, D. L., Yacoub, E., Xu, J., Bartsch, A. J., Jbabdi, S., Sotiropoulos, S. N., Andersson, J. L. R., Griffanti, L., Douaud, G., Okell, T. W., Weale, P., Dragonu, I., Garratt, S., Hudson, S., Collins, R., Jenkinson, M., ... Smith, S. M. (2016). Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nature Neuroscience*, 19(11), 1523–1536. <https://doi.org/10.1038/nn.4393>
- Mohajer, B., Abbasi, N., Mohammadi, E., Khazaie, H., Osorio, R. S., Rosenzweig, I., Eickhoff, C. R., Zarei, M., Tahmasian, M., Eickhoff, S. B., & Initiative, for the A. D. N. (2020). Gray matter volume and estimated brain age gap are not linked with sleep-disordered breathing. *Human Brain Mapping*, 41(11), 3034–3044. <https://doi.org/10.1002/hbm.24995>
- Murdoch, W. J., Singh, C., Kumbier, K., Abbasi-Asl, R., & Yu, B. (2019). Definitions, methods, and applications in interpretable machine learning. *Proceedings of the National Academy of Sciences*, 116(44), 22071–22080. <https://doi.org/10.1073/pnas.1900654116>
- Nadeau, C., & Bengio, Y. (2003). Inference for the Generalization Error. *Machine Learning*, 52(3), 239–281. <https://doi.org/10.1023/A:1024068626366>
- Norbury, R. (2020). *Diurnal Preference and Grey Matter Volume in a Large Population of Older Adults: Data from the UK Biobank* (1). 18(1), Article 1. <https://doi.org/10.5334/jcr.193>
- Olfati, M., Samea, F., Faghihroohi, S., Balajoo, S. M., Küppers, V., Genon, S., Patil, K., Eickhoff, S. B., & Tahmasian, M. (2024a). *Prediction of depressive symptoms severity based on sleep quality, anxiety, and brain: A machine learning approach across three cohorts*. <https://doi.org/10.1101/2023.08.09.23293887>
- Olfati, M., Samea, F., Faghihroohi, S., Balajoo, S. M., Küppers, V., Genon, S., Patil, K., Eickhoff, S. B., & Tahmasian, M. (2024b). Prediction of depressive symptoms severity based on sleep quality, anxiety, and gray matter volume: A generalizable machine learning approach across three datasets. *eBioMedicine*, 108. <https://doi.org/10.1016/j.ebiom.2024.105313>
- Omidvarnia, A., Sasse, L., Larabi, D. I., Raimondo, F., Hoffstaedter, F., Kasper, J., Dukart, J., Petersen, M., Cheng, B., Thomalla, G., Eickhoff, S. B., & Patil, K. R. (2023). *Is resting state fMRI better than individual characteristics at predicting cognition?* (p. 2023.02.18.529076). bioRxiv. <https://doi.org/10.1101/2023.02.18.529076>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, É. (2012). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830. <https://doi.org/10.1007/s13398-014-0173-7.2>
- R: *Fast Heuristics For The Estimation Of the C Constant Of A...* (n.d.). Retrieved December 9, 2022, from <https://search.r-project.org/CRAN/refmans/LiblineaR/html/heuristicC.html>

- Rahimi-Jafari, S., Sarebannejad, S., Saberi, A., Khazaie, H., Camilleri, J. A., Eickhoff, C. R., Eickhoff, S. B., & Tahmasian, M. (2022). Is there any consistent structural and functional brain abnormality in narcolepsy? A meta-analytic perspective. *Neuroscience & Biobehavioral Reviews*, 132, 1181–1182. <https://doi.org/10.1016/j.neubiorev.2021.10.034>
- Ravyts, S. G., Dizerzewski, J. M., Perez, E., Donovan, E. K., & Dautovich, N. (2021). Sleep Health as Measured by RU SATED: A Psychometric Evaluation. *Behavioral Sleep Medicine*, 19(1), 48–56. <https://doi.org/10.1080/15402002.2019.1701474>
- Reimann, G. M., Küppers, V., Camilleri, J. A., Hoffstaedter, F., Langner, R., Laird, A. R., Fox, P. T., Spiegelhalter, K., Eickhoff, S. B., & Tahmasian, M. (2023). Convergent abnormality in the subgenual anterior cingulate cortex in insomnia disorder: A revisited neuroimaging meta-analysis of 39 studies. *Sleep Medicine Reviews*, 71, 101821. <https://doi.org/10.1016/j.smrv.2023.101821>
- Rutherford, S., Kia, S. M., Wolfers, T., Frazz, C., Zabihi, M., Dinga, R., Berthet, P., Worker, A., Verdi, S., Ruhe, H. G., Beckmann, C. F., & Marquand, A. F. (2022). The normative modeling framework for computational psychiatry. *Nature Protocols*, 17(7), 1711–1734. <https://doi.org/10.1038/s41596-022-00696-5>
- Sasse, L., Nicolaisen-Sobesky, E., Dukart, J., Eickhoff, S. B., Götz, M., Hamdan, S., Komeyer, V., Kulkarni, A., Lahnakoski, J., Love, B. C., Raimondo, F., & Patil, K. R. (2024). *On Leakage in Machine Learning Pipelines* (arXiv:2311.04179). arXiv. <https://doi.org/10.48550/arXiv.2311.04179>
- Sateia, M. J. (2014). International classification of sleep disorders-third edition: Highlights and modifications. *Chest*, 146(5), 1387–1394. <https://doi.org/10.1378/chest.14-0970>
- Schaefer, A., Kong, R., Gordon, E. M., Laumann, T. O., Zuo, X.-N., Holmes, A. J., Eickhoff, S. B., & Yeo, B. T. T. (2018). Local-Global Parcellation of the Human Cerebral Cortex from Intrinsic Functional Connectivity MRI. *Cerebral Cortex*, 28(9), 3095–3114. <https://doi.org/10.1093/cercor/bhx179>
- Schiell, J. E., Tamm, S., Holub, F., Petri, R., Dashti, H. S., Domschke, K., Feige, B., Goodman, M. O., Jones, S. E., Lane, J. M., Ratti, P.-L., Ray, D. W., Redline, S., Riemann, D., Rutter, M. K., Saxena, R., Sexton, C. E., Tahmasian, M., Wang, H., ... Spiegelhalter, K. (2023). Associations between sleep health and grey matter volume in the UK Biobank cohort (  $n = 33\,356$ ). *Brain Communications*, 5(4), fcad200. <https://doi.org/10.1093/braincomms/fcad200>
- Schiell, J. E., Tamm, S., Holub, F., Petri, R., Dashti, H. S., Domschke, K., Feige, B., Lane, J. M., Riemann, D., Rutter, M. K., Saxena, R., Tahmasian, M., Wang, H., Kyle, S. D., & Spiegelhalter, K. (2022). Associations Between Sleep Health and Amygdala Reactivity to Negative Facial Expressions in the UK Biobank Cohort. *Biological Psychiatry*, 92(9), 693–700. <https://doi.org/10.1016/j.biopsych.2022.05.023>

- Schoeler, T., Pingault, J.-B., & Kutalik, Z. (2023). *Self-report inaccuracy in the UK Biobank: Impact on inference and interplay with selective participation* (p. 2023.10.06.23296652). medRxiv. <https://doi.org/10.1101/2023.10.06.23296652>
- Schulz, M.-A., Bzdok, D., Haufe, S., Haynes, J.-D., & Ritter, K. (2024). Performance reserves in brain-imaging-based phenotype prediction. *Cell Reports*, 43(1), 113597. <https://doi.org/10.1016/j.celrep.2023.113597>
- Stolicyn, A., Lyall, L. M., Lyall, D. M., Høier, N. K., Adams, M. J., Shen, X., Cole, J. H., McIntosh, A. M., Whalley, H. C., & Smith, D. J. (2023). Comprehensive assessment of sleep duration, insomnia, and brain structure within the UK Biobank cohort. *Sleep*, zsad274. <https://doi.org/10.1093/sleep/zsad274>
- Tahmasian, M., Aleman, A., Andreassen, O. A., Arab, Z., Baillet, M., Benedetti, F., Bresser, T., Bright, J., Chee, M. W. L., Chylinski, D., Cheng, W., Deantoni, M., Dresler, M., Eickhoff, S. B., Eickhoff, C. R., Elvsåshagen, T., Feng, J., Foster-Dingley, J. C., Ganjgahi, H., ... Zarei, M. (2021). ENIGMA-Sleep: Challenges, opportunities, and the road map. *Journal of Sleep Research*, 30(6), e13347. <https://doi.org/10.1111/jsr.13347>
- Tahmasian, M., Rosenzweig, I., Eickhoff, S. B., Sepehry, A. A., Laird, A. R., Fox, P. T., Morrell, M. J., Khazaie, H., & Eickhoff, C. R. (2016). Structural and functional neural adaptations in obstructive sleep apnea: An activation likelihood estimation meta-analysis. *Neuroscience & Biobehavioral Reviews*, 65, 142–156. <https://doi.org/10.1016/j.neubiorev.2016.03.026>
- Tahmasian, M., Samea, F., Khazaie, H., Zarei, M., Kharabian Masouleh, S., Hoffstaedter, F., Camilleri, J., Kochunov, P., Yeo, B. T. T., Eickhoff, S. B., & Valk, S. L. (2020). The interrelation of sleep and mental and physical health is anchored in grey-matter neuroanatomy and under genetic control. *Communications Biology*, 3(1), Article 1. <https://doi.org/10.1038/s42003-020-0892-6>
- Tahmasian, M., Shao, J., Meng, C., Grimmer, T., Diehl-Schmid, J., Yousefi, B. H., Förster, S., Riedl, V., Drzezga, A., & Sorg, C. (2015). Based on the network degeneration hypothesis: Separating individual patients with different neurodegenerative syndromes in a preliminary hybrid PET/MR study. *Journal of Nuclear Medicine*. <https://doi.org/10.2967/jnumed.115.165464>
- Tai, X. Y., Chen, C., Manohar, S., & Husain, M. (2022). Impact of sleep duration on executive function and brain structure. *Communications Biology*, 5(1), 201. <https://doi.org/10.1038/s42003-022-03123-3>
- Tian, Y. E., Di Biase, M. A., Mosley, P. E., Lupton, M. K., Xia, Y., Fripp, J., Breakspear, M., Cropley, V., & Zalesky, A. (2023). Evaluation of Brain-Body Health in Individuals With Common Neuropsychiatric Disorders. *JAMA Psychiatry*, 80(6), 567–576. <https://doi.org/10.1001/jamapsychiatry.2023.0791>
- Tian, Y. E., Margulies, D. S., Breakspear, M., & Zalesky, A. (2020). *Topographic organization of the human subcortex unveiled with functional connectivity gradients* (p. 2020.01.13.903542). bioRxiv. <https://doi.org/10.1101/2020.01.13.903542>

- Tsiknia, A. A., Parada, H., Banks, S. J., & Reas, E. T. (2023). Sleep quality and sleep duration predict brain microstructure among community-dwelling older adults. *Neurobiology of Aging*, 125, 90–97. <https://doi.org/10.1016/j.neurobiolaging.2023.02.001>
- Varoquaux, G. (2017). Cross-validation failure: Small sample sizes lead to large error bars. *Neuroimage*.
- Varoquaux, G., Raamana, P. R., Engemann, D. A., Hoyos-Idrobo, A., Schwartz, Y., & Thirion, B. (2017). Assessing and tuning brain decoders: Cross-validation, caveats, and guidelines. *NeuroImage*, 145, 166–179. <https://doi.org/10.1016/j.neuroimage.2016.10.038>
- Vieira, S., Pinaya, W. H. L., & Mechelli, A. (2017). Using deep learning to investigate the neuroimaging correlates of psychiatric and neurological disorders: Methods and applications. *Neuroscience & Biobehavioral Reviews*, 74, 58–75. <https://doi.org/10.1016/j.neubiorev.2017.01.002>
- Walker, M. P. (2021). Sleep essentialism. *Brain*, 144(3), 697–699. <https://doi.org/10.1093/brain/awab026>
- Wang, Y., Genon, S., Dong, D., Zhou, F., Li, C., Yu, D., Yuan, K., He, Q., Qiu, J., Feng, T., Chen, H., & Lei, X. (2023). Covariance patterns between sleep health domains and distributed intrinsic functional connectivity. *Nature Communications*, 14(1), 7133. <https://doi.org/10.1038/s41467-023-42945-5>
- Weihs, A., Frenzel, S., Bi, H., Schiel, J. E., Afshani, M., Bülow, R., Ewert, R., Fietze, I., Hoffstaedter, F., Jahanshad, N., Khazaie, H., Riemann, D., Rostampour, M., Stubbe, B., Thomopoulos, S. I., Thompson, P. M., Valk, S. L., Völzke, H., Zarei, M., ... Group, for the E. N. G. through M.-A. (ENIGMA)-S. W. (2023). Lack of structural brain alterations associated with insomnia: Findings from the ENIGMA-Sleep Working Group. *Journal of Sleep Research*, 32(5), e13884. <https://doi.org/10.1111/jsr.13884>
- Whitfield-Gabrieli, S., & Nieto-Castanon, A. (2012). Conn: A functional connectivity toolbox for correlated and anticorrelated brain networks. *Brain Connectivity*, 2(3), 125–141. <https://doi.org/10.1089/brain.2012.0073>
- Wiersch, L., Hamdan, S., Hoffstaedter, F., Votinov, M., Habel, U., Clemens, B., Derntl, B., Eickhoff, S. B., Patil, K. R., & Weis, S. (2023). Accurate sex prediction of cisgender and transgender individuals without brain size bias. *Scientific Reports*, 13(1), 13868. <https://doi.org/10.1038/s41598-023-37508-z>
- Williams, J. A., Russ, D., Bravo-Merodio, L., Gkoutos, G., Bellgrove, M. A., Bagshaw, A. P., & Checlacz, M. (2023). *Genetically mediated associations between chronotype and neuroimaging phenotypes in the UK Biobank: A Mendelian randomisation study* (p. 2023.08.31.555801). bioRxiv. <https://doi.org/10.1101/2023.08.31.555801>
- Willoughby, A. R., Alikhani, I., Karsikas, M., Chua, X. Y., & Chee, M. W. L. (2023). Country differences in nocturnal sleep variability: Observations from a large-scale, long-term sleep wearable study. *Sleep Medicine*, 110, 155–165. <https://doi.org/10.1016/j.sleep.2023.08.010>

- Winter, N. R., Blanke, J., Leenings, R., Ernsting, J., Fisch, L., Sarink, K., Barkhau, C., Emden, D., Thiel, K., Flinkenflügel, K., Winter, A., Goltermann, J., Meinert, S., Dohm, K., Reppe, J., Gruber, M., Leehr, E. J., Opel, N., Grotegerd, D., ... Hahn, T. (2024). A Systematic Evaluation of Machine Learning–Based Biomarkers for Major Depressive Disorder. *JAMA Psychiatry*. <https://doi.org/10.1001/jamapsychiatry.2023.5083>
- Wolpert, D. H. (1992). Stacked generalization. *Neural Networks*, 5(2), 241–259. [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1)
- Woo, C.-W., Chang, L. J., Lindquist, M. A., & Wager, T. D. (2017). Building better biomarkers: Brain models in translational neuroimaging. *Nature Neuroscience*, 20(3), Article 3. <https://doi.org/10.1038/nn.4478>
- Zou, Q.-H., Zhu, C.-Z., Yang, Y., Zuo, X.-N., Long, X.-Y., Cao, Q.-J., Wang, Y.-F., & Zang, Y.-F. (2008). An improved approach to detection of amplitude of low-frequency fluctuation (ALFF) for resting-state fMRI: Fractional ALFF. *Journal of Neuroscience Methods*, 172(1), 137–141. <https://doi.org/10.1016/j.jneumeth.2008.04.012>