

# Entropy production of Multivariate Ornstein-Uhlenbeck processes correlates with consciousness levels in the human brain

Matthieu Gilson<sup>1</sup>, Enzo Tagliazucchi<sup>2</sup>, Rodrigo Cofré<sup>3</sup>

<sup>1</sup>*Institut de Neurosciences des Systèmes INSERM-AMU, Marseille, France*

<sup>2</sup>*Physics Department University of Buenos Aires and Buenos Aires Physics Institute Argentina*

<sup>3</sup>*Institute of Neuroscience (NeuroPSI-CNRS), Paris-Saclay University, Gif sur Yvette 91400, France*

Consciousness is supported by complex patterns of brain activity which are indicative of irreversible non-equilibrium dynamics. While the framework of stochastic thermodynamics has facilitated the understanding of physical systems of this kind, its application to infer the level of consciousness from empirical data remains elusive. We faced this challenge by calculating entropy production in a multivariate Ornstein-Uhlenbeck process fitted to fMRI brain activity recordings. To test this approach, we focused on the transition from wakefulness to deep sleep, revealing a monotonous relationship between entropy production and the level of consciousness. Our results constitute robust signatures of consciousness while also advancing our understanding of the link between consciousness and complexity from the fundamental perspective of statistical physics.

Animal cognition is the most sophisticated example of information processing found in biological and technological systems [1]. Consciousness, understood as the capacity to sustain subjective experience, can be considered a property that emerges when a sufficiently high level of complex cognitive processing is achieved [2]. From the perspective of physics, consciousness and cognition seem unlikely to emerge from regular and predictable systems, such as those which are in thermodynamic equilibrium and obey the detailed balance equations [3]. Instead, recent research draws a close parallel between the level of consciousness and the entropy rate of brain activity time series, highlighting temporal irreversibility as a landmark feature of conscious information processing [4–6]. These results suggest a close link between consciousness and non-equilibrium dynamics, prompting a rigorous evaluation from the perspective of stochastic thermodynamics.

In spite of these exciting results, the direct estimation of entropy production from neural activity recordings is undermined by insufficient spatio-temporal sampling, leading to the adoption of heuristics and approximations which lack rigorous justification [3, 4]. To circumvent these limitations, we adopted a framework based on Multivariate Ornstein-Uhlenbeck (MOU) processes, that are widely used for modeling the multivariate dynamics of time series. The importance of MOU derives from the fact that it is the only continuous stationary stochastic process that is simultaneously Gaussian and Markovian. The MOU process is at the heart of many models used to fit fMRI data and to interpret them in terms of whole-brain communication [7–9], in line with the present methodology. We first characterize the non-equilibrium steady state of a generic MOU process. The irreversibility of the process is encoded in the antisymmetric part of the Onsager matrix, while the linearity of the Langevin equations allows us to derive closed-form expression for the entropy production rate in terms of the matrices that define the MOU. As a result, we obtained a model-based estimation of the entropy production rate for the MOU

fitted to fMRI data of subjects transitioning different levels of consciousness during the descent from wakefulness to deep sleep.

Time reversibility and entropy production in the MOU process.— We consider the MOU process closely following the notation in previous work [10]:

$$\frac{d\mathbf{x}(t)}{dt} = -\mathbf{B}\mathbf{x}(t) + \boldsymbol{\eta}(t) . \quad (1)$$

Boldfaced symbols denote vectors and matrices. The inputs  $\boldsymbol{\eta}(t)$  correspond to Gaussian white noise with covariance

$$\langle \boldsymbol{\eta}(t)\boldsymbol{\eta}^T(t') \rangle_t = 2\mathbf{D}\delta(t-t') . \quad (2)$$

The angular brackets indicate the mathematical expectation over time and the superscript T the transpose for vectors or matrices. The  $N$ -dimensional MOU process is thus defined by two real  $N \times N$  matrices, the input covariance matrix  $\mathbf{D}$ , which is symmetric with positive eigenvalues, and the Jacobian  $\mathbf{B}$ , which is not symmetric in general.

Following previous results [11, 12], a sufficient condition for the MOU process in Eq. (1) to be a time reversible stationary process corresponds to a specific relation between the matrices  $\mathbf{B}$  and  $\mathbf{D}$ :

$$\mathbf{B}\mathbf{D} = \mathbf{D}\mathbf{B}^T . \quad (3)$$

To quantify the time (ir)reversibility of the MOU process, it is advantageous to examine the Onsager matrix  $\mathbf{L}$  reparameterized using the matrices  $\mathbf{B}$ ,  $\mathbf{D}$ , and the pairwise zero-lag covariance  $\mathbf{S} = \langle \mathbf{x}(t)\mathbf{x}^T(t) \rangle_t$ :

$$\begin{aligned} \mathbf{L} &= \mathbf{B}\mathbf{S} = \mathbf{D} + \mathbf{Q} , \\ \mathbf{L}^T &= \mathbf{S}\mathbf{B}^T = \mathbf{D} - \mathbf{Q} . \end{aligned} \quad (4)$$

Here the antisymmetric part  $\mathbf{Q}$  of  $\mathbf{L}$  provides a measure for the irreversibility of the process. When the process

is time reversible  $\mathbf{Q} = 0$  and  $\mathbf{L}$  is symmetric. The following expression for the entropy production rate  $\Phi$  can then be derived (see the supplementary material for a full derivation of this equation)

$$\Phi = \text{tr}(\mathbf{B}^T \mathbf{D}^{-1} \mathbf{Q}) = -\text{tr}(\mathbf{D}^{-1} \mathbf{B} \mathbf{Q}). \quad (5)$$

The entropy production rate  $\Phi$  is non-negative and provides a scalar measure for the (ir)reversibility of the whole network process, vanishing only if the process is reversible.

MOU-based anatomo-functional model to fit empirical fMRI data.— We fitted a MOU process to the time series of blood oxygen level-dependent (BOLD) activity measured using fMRI for a whole-brain parcellation consisting of  $N = 90$  regions of interest (ROIs). The BOLD signals were recorded from 15 healthy participants during wakefulness and three sleep stages of progressively deeper unconsciousness (N1, N2, N3). Further details about the data preprocessing like detrending and filtering are found in [13]. Example BOLD time series are illustrated in Fig. 1A. Fig. 1B-C show two functional connectivity matrices, here calculated as covariances with zero lag  $\hat{\mathbf{S}}(0)$  and lag of 1 timestep  $\hat{\mathbf{S}}(1)$ . These matrices are the empirical counterparts of the the model pairwise covariance  $\mathbf{S}(l) = \langle \mathbf{x}(t) \mathbf{x}^T(t+l) \rangle_t$  with lag  $l$ , which is symmetric for  $l = 0$  and was denoted above by  $\mathbf{S} = \mathbf{S}(0)$ .

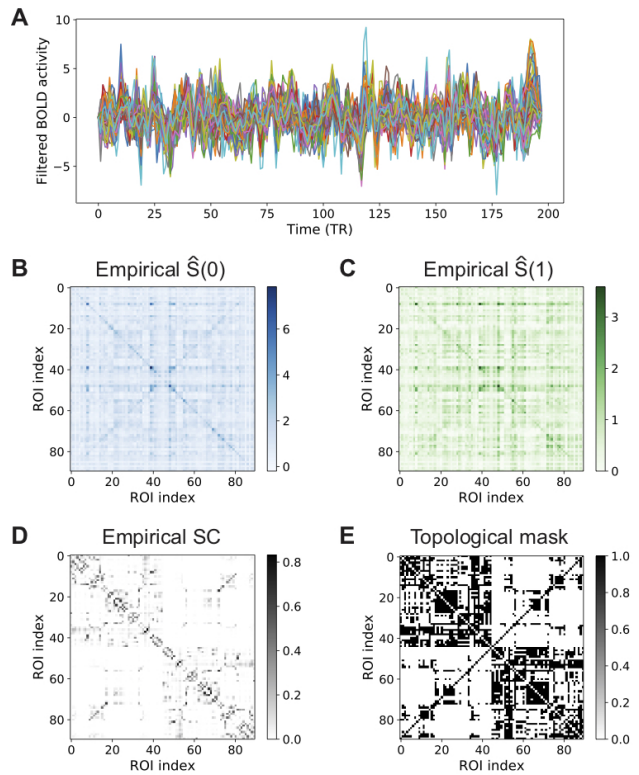
In this application, the activity  $x_i$  of the MOU process describes the BOLD activity of node  $i$ . Its Jacobian matrix  $\mathbf{B}$  quantifies the propagation of BOLD activity between ROIs, ignoring hemodynamics. Specifically, the diagonal elements  $B_{ii}$  are related to a time constant  $\tau$  (identical for all ROIs) and the off-diagonal elements  $C_{ij} = -B_{ij}$  correspond to the concept of effective connectivity from ROI  $j$  to ROI  $i$  (excitatory when  $C_{ij} > 0$ ):

$$-B_{ij} = -\frac{\delta_{ij}}{\tau} + C_{ij}, \quad (6)$$

where  $\delta_{ij}$  is the Kronecker delta. The input variance  $D_{ii}$  reflects the fluctuation amplitude of the spontaneous activity to ROI  $i$ .

For each subject and condition, the model was fitted to reproduce the two covariance matrices calculated from the empirical BOLD signals  $\hat{\mathbf{S}}(0)$  and  $\hat{\mathbf{S}}(1)$  (see Fig. 1B-C). We used a recent estimation method based on gradient descent to iteratively adjust  $\mathbf{B}$  and  $\mathbf{D}$  until reaching the best fit [7]. At each optimization step, we calculate the model counterparts of the covariance matrices  $\mathbf{S}(0)$  and  $\mathbf{S}(1)$ , assuming stationarity over each fMRI session. Importantly, this optimization procedure incorporates topological constraints on  $\mathbf{B}$ , adjusting only existing anatomical connections (see Fig. 1D-E), also keeping the input cross-covariances  $D_{ij} = 0$  for  $i \neq j$ .

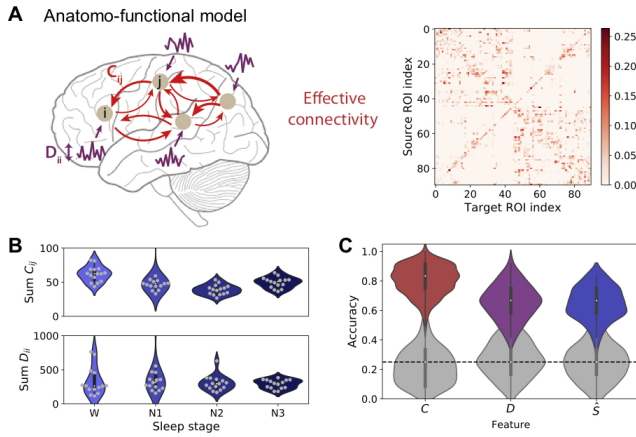
Robust decoding of sleep stages from MOU parameters.— Following previous work [15, 16], we used the scikit-learn Python library for the implementations of multinomial logistic regression (MLR) classifier. The input features corresponded to the vectorized



**FIG. 1. Empirical BOLD time series, functional and structural connectivity.** A) Example of the filtered BOLD time series (of duration  $T = 198$ ) corresponding to the 90 ROIs of the AAL parcellation during wakefulness of one participant. B-C) Functional connectivity matrices calculated from the filtered BOLD signals in panel A,  $\hat{\mathbf{S}}(0)$  with zero lag and  $\hat{\mathbf{S}}(1)$  with a lag of one timestep ( $\text{TR}=2$  s). These matrices are used in the objective functions used to fit the anatomo-functional model. D) Generic structural connectivity (SC) obtained from DTI data as described in [14]. E) Mask for existing directional connections to constrain the topology of the  $\mathbf{B}$  matrix in the network model (symmetric here).

$\mathbf{C}/\mathbf{D}/\mathbf{S}$  matrices after discarding zero or redundant elements. We implemented a stratified cross-validation scheme with 80% of the samples for the train set and 20% for the test set, where the ratio of classes are the same in both sets. We also use the subject identity as “group information” to avoid mixing subject data between the train and test sets. In practice, we use 100 random splits of the data and report the distribution of accuracies of the 100 splits.

As illustrated in Fig. 2B, both the empirical BOLD variances and the model estimates exhibit global differences across the four sleep stages, although they do not exhibit a clear trend. These differences in global measures, which are averages over all ROIs, may hide more specific changes at the ROI level, as well as interactions between them. Supplementary Figure S1 shows the good fit of the anatomo-functional model to fMRI data that obtained for all sleep stages, with mean correlation be-

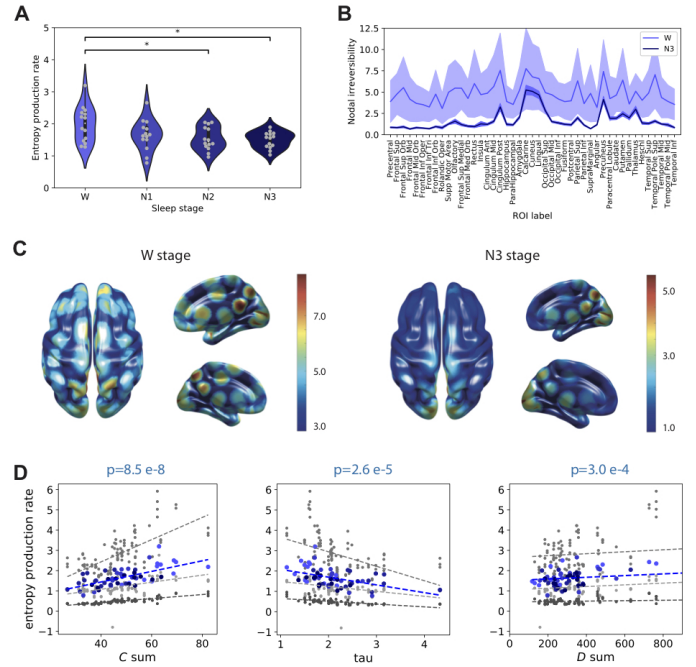


**FIG. 2. Anom-functional model fitted to the empirical data.** A) Our dynamic network model has two sets of optimized parameters: the matrix  $\mathbf{C}$ , which describes the causal interaction between brain regions, and the input variance  $\mathbf{D}$ , which represents the spontaneous activity of each brain region. Note that the topology of the matrix  $\mathbf{C}$  corresponds to the mask inferred from the SC data in Fig. 1D-E, but the weights are estimated from the empirical FC matrices Fig. 1B-C, resulting in an anom-functional model. B) Changes in total  $\mathbf{C}$  and  $\mathbf{D}$  weights across sleep stages (x-axis), pooled over the 15 subjects. The sleep stages are represented by the blue contrasts, from light for wake (W) to dark for the deepest sleep (N3). C) Classification accuracy based on the model estimates,  $\mathbf{C}$  and  $\mathbf{D}$ , and the empirical covariance matrices. The classifier is the multinomial logistic regression (MLR), which captures changes in individual features across sleep stages. The gray violin plots correspond to the chance-level accuracy calculated empirically in each case by shuffling the labels of the sleep stages.

tween simulated and empirical FC matrices exceeding 0.6 for all stages. Fig. 2C shows that the model estimates give good classification accuracy, both for  $\mathbf{C}$  (in red) and  $\mathbf{D}$  (in purple). This indicates that the model captures the differences in brain dynamics across the sleep stages. Notably, the matrix  $\mathbf{C}$  gives a better classification accuracy than the empirical functional connectivity  $\hat{S}(0)$  (in blue), meaning that the model inversion is robust and captures refined information about the sleep stages. Note that the MLR has a better accuracy than the 1-nearest-neighbor (1NN) in Suppl Fig S2A, indicating that the changes across sleep stages concern specific features, i.e. connectivity weights ( $\mathbf{C}$ ) or nodal spontaneous activity ( $\mathbf{D}$ ), rather than their global profile.

Reduced entropy production in the transition from wakefulness to deep sleep.— Using the condition-specific estimated parameters, we calculated the entropy production rate in the MOU model using Eq. (5). These results in Fig 3A show that entropy production decreases as a function of sleep depth, which in turn implies that dynamics become closer to equilibrium.

The model-based approach allows us to dissect this phenomenon. For all ROIs, we observe that the contri-



**FIG. 3. Entropy production rate correlates with consciousness levels in the human brain.** A) Violin plots comparing the entropy production rate across sleep stages. Same color coding used in previous plots. The average entropy production values across subjects are for the four sleep stages are: 1.99, 1.65, 1.54 and 1.49, respectively. The stars indicate statistical significance for the Mann-Whitney test with  $p < 0.05$ . B) Comparison of the nodal irreversibility for each ROI (x-axis) between the W and N3 states (in light and dark blue, respectively). The plotted values correspond to the absolute value of sums over rows of  $\mathbf{Q}$ , averaged for homotopic regions; error bars indicate the variability across subjects measured as standard error of the mean. C) Heatmap plots of the nodal irreversibility on the cortical surface for the W and N3 sleep stages. Note the different color scales for the two stages, for the purpose of better readability. D) Plots of  $\Phi$  across sleep stages and subjects as a function of the sum of  $\mathbf{C}$  weights (left panel), the time constant  $\tau$  on the diagonal of  $\mathbf{B}$  (middle panel) and the sum of  $\mathbf{D}$  variances (right panel). Comparison with surrogate values with randomized (redistributing the total weight/variances keeping the same topology and overall sum)  $\mathbf{C}$  matrices (light gray), randomized  $\mathbf{D}$  matrices (middle gray) and both (dark gray).

tribution to  $\Phi$ , as measured via the nodal irreversibility, defined as  $\sum_j |Q_{ij}|$  for each ROI  $i$ , decreases, as illustrated in Fig 3B. This suggests that the reduction of  $\Phi$  from W to N3 is a rather global phenomenon, but with a differentiated magnitude across brain regions. Notably, regions in the occipital lobes (cuneus, calcarine, lingual) as well as regions associated to hubs in the default-mode network (precuneus, post cingulate) and the thalamus remain at a high level of nodal irreversibility in the deep sleep N3; these regions have been shown to exhibit sleep-related changes in previous studies [17–19]. See Suppl Fig S3 for a more detailed comparison across sleep stages.

Last, we examine how the model parameters  $\mathbf{C}$  and  $\mathbf{D}$  contribute to  $\Phi$  and its reduction across sleep stages. Fig 3D shows a positive relationship between  $\Phi$  and the sum of weights in  $\mathbf{C}$ , as well as the sum of variances in  $\mathbf{D}$ ; conversely, a larger  $\tau$  (directly calculated from the empirical BOLD signals) corresponds to a lower  $\Phi$ . Then we assess the importance of the detailed structures in the  $\mathbf{C}$  and  $\mathbf{D}$  estimates by randomizing them spatially, namely redistributing the total weight/variances across non-zero elements while keeping the same topology and overall sum. We observe the same trends with respect to the  $\mathbf{C}$  and  $\mathbf{D}$  sums, but shifted up or down depending on the surrogates in Fig 3C: randomizing  $\mathbf{C}$  (light gray) decreases slightly  $\Phi$ , whereas randomizing  $\mathbf{D}$  (middle gray) increases  $\Phi$ ; randomizing both (dark gray) decreases  $\Phi$ . This indicates that  $\Phi$  strongly depends on the detailed structures of the  $\mathbf{C}$  and  $\mathbf{D}$  estimates, being larger in the data than in the randomized surrogates. The opposing effects in randomizing  $\mathbf{C}$  and  $\mathbf{D}$  also suggest a balance implemented by the detailed brain dynamics, which results in a controlled level of  $\Phi$ . Together, our results hint at a positive relationship between the measured  $\Phi$  and the different levels of consciousness.

Discussion.— We measured the entropy production using our anatomo-functional MOU process associated to resting-state fMRI activity recorded from human subjects in different sleep stages. The advantage of our model-based approach is that the entropy production has a closed-form expression from first principles of stochastic

thermodynamics for the MOU process, which is numerically fitted to the fMRI data. Our results show high entropy production rate in conscious wakefulness, i.e. correlating positively with the presumed level of cognitive processing. This is consistent with converging theoretical accounts that identify consciousness with an emergent property of a highly complex physical system [2]. These results are also consistent with previous findings relating entropy production with states of consciousness [4–6], with the advantage that do not depend on heuristic approximations. Importantly, our approach allows to identify the brain regions that contribute most to entropy production. The fulfilment of detailed balance in the brain is scale-dependent [3]. At the large scale, its violation might relate to the large-scale circuit operations critical for healthy cognition and for the global broadcasting of information which is identified with the computational aspect of consciousness [20]. Because of this, metrics related to the departure from detailed balance (such as entropy production rate) might offer valuable tools to determine levels of consciousness in brain injured patients and other neurological populations.

In summary, assessing temporal irreversibility through entropy production of MOU processes derived from fMRI signals has the potential to highlight different states of consciousness and cognition. More generally, can bridge brain dynamics and thermodynamics, and ultimately help to understand fundamental questions about the brain and consciousness.

- 
- [1] G. Piccinini and A. Scarantino, *Journal of biological physics* **37**, 1 (2011).
  - [2] A. K. Seth and T. Bayne, *Nature Reviews Neuroscience*, 1 (2022).
  - [3] C. W. Lynn, E. J. Cornblath, L. Papadopoulos, M. A. Bertolero, and D. S. Bassett, *PNAS* (2021).
  - [4] Y. Sanz Perl, H. Bocaccio, C. Pallavicini, I. Pérez-Ipiña, S. Laureys, H. Laufs, M. Kringelbach, G. Deco, and E. Tagliazucchi, *Physical Review E* (2021).
  - [5] L. de la Fuente, F. Zamberlan, H. Bocaccio, M. Kringelbach, G. Deco, Y. S. Perl, and E. Tagliazucchi, *Cerebral Cortex* (2022).
  - [6] R. N. Muñoz, A. Leung, A. Zecevik, F. A. Pollock, D. Cohen, B. van Swinderen, N. Tsuchiya, and K. Modi, *Physical Review Research* **2**, 023219 (2020).
  - [7] M. Gilson, R. Moreno-Bote, A. Ponce-Alvarez, P. Ritter, and G. Deco, *PLoS Computational Biology* **12** (2016).
  - [8] K. J. Friston, K. H. Preller, C. Mathys, H. Cagnan, J. Heinzle, A. Razi, and P. Zeidman, *NeuroImage* **199** (2019).
  - [9] S. Frässle, Z. M. Manjaly, C. T. Do, L. Kasper, K. P. Pruessmann, and K. E. Stephan, *NeuroImage* **225** (2021).
  - [10] C. Godrèche and J. M. Luck, *Journal of Physics A: Mathematical and Theoretical* (2019).
  - [11] H. Risken and T. Frank, *Springer Series in Synergetics* (1996).
  - [12] H. Qian, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 10.1098/rspa.2001.0811 (2001).
  - [13] E. Tagliazucchi and H. Laufs, *Neuron* (2014).
  - [14] I. P. Ipiña, P. D. Kehoe, M. Kringelbach, H. Laufs, A. Ibañez, G. Deco, Y. S. Perl, and E. Tagliazucchi, *NeuroImage* (2020).
  - [15] V. Pallarés, A. Insabato, A. Sanjuán, S. Kühn, D. Mantini, G. Deco, and M. Gilson, *NeuroImage* **178** (2018).
  - [16] M. Gilson, G. Zamora-López, V. Pallarés, M. H. Adhikari, M. Senden, A. T. Campo, D. Mantini, M. Corbetta, G. Deco, and A. Insabato, *Network Neuroscience* **4** (2020).
  - [17] S. G. Horowitz, A. R. Braun, W. S. Carr, D. Picchioni, T. J. Balkin, M. Fukunaga, and J. H. Duyn, *Proceedings of the National Academy of Sciences of the United States of America* **106**, 10.1073/pnas.0901435106 (2009).
  - [18] T. T. Dang-Vu, M. Schabus, M. Desseilles, V. Sterpenich, M. Bonjean, and P. Maquet, *Functional neuroimaging insights into the physiology of human sleep* (2010).
  - [19] L. Mirandola, G. Cantalupo, A. E. Vaudano, P. Avanzini, A. Ruggieri, F. Pisani, G. Cossu, C. A. Tassinari, P. F. Nichelli, F. Benuzzi, and S. Meletti, *Epilepsy and Behavior Case Reports* **1**, 10.1016/j.ebcr.2013.06.005 (2013).
  - [20] S. Dehaene, M. Kerszberg, and J. P. Changeux, *Proceedings of the National Academy of Sciences of the United States of America* **95**, 10.1073/pnas.95.24.14529 (1998).

# Supplemental Material for: Entropy production of Multivariate Ornstein-Uhlenbeck processes correlates with consciousness levels in the human brain

Matthieu Gilson<sup>1</sup>, Enzo Tagliazucchi<sup>3</sup>, Rodrigo Cofré<sup>4</sup>

<sup>1</sup>*Institut de Neurosciences des Systèmes INSERM-AMU, Marseille, France*

<sup>2</sup>*Physics Department University of Buenos Aires and Buenos Aires Physics Institute Argentina*

<sup>3</sup>*Institute of Neuroscience (NeuroPSI-CNRS), Paris-Saclay University, Gif sur Yvette 91400, France*

In this supplemental material, we present a detailed derivations of the equations used in this article, methods to compute the parameters of the model, and detailed description of the fMRI data.

## I. MULTIVARIATE ORNSTEIN-UHLENBECK PROCESS

### A. Description of the state evolution

Knowing the initial condition  $\mathbf{x}(0)$  and the realization of the stochastic input  $\boldsymbol{\eta}$  over time, the trajectory of the solution of the Eq.1 in the main text is given by:

$$\mathbf{x}(t) = \mathbf{G}(t) \mathbf{x}(0) + \int_0^t \mathbf{G}(t-s) \boldsymbol{\eta}(s) ds, \quad (1)$$

where  $\mathbf{G}(t) = e^{-\mathbf{B}t}$  is the Green's function, also known as propagator. In addition to its mean value  $\langle \mathbf{x}(t) \rangle = \mathbf{G}(t) \mathbf{x}(0)$ , the process is also characterized by its covariance matrix  $\mathbf{S}(t, t') = \langle \mathbf{x}(t) \mathbf{x}^T(t') \rangle$ . The zero-lag covariance, denoted by  $\mathbf{S}(t, t)$ , obeys the following deterministic differential equation:

$$\frac{d\mathbf{S}(t, t)}{dt} = -\mathbf{B} \mathbf{S}(t, t) - \mathbf{S}(t, t) \mathbf{B}^T + 2\mathbf{D}. \quad (2)$$

Meanwhile, the lagged covariance with  $t' > t$  exhibits an exponential decay as a function of the lag  $t' - t$ :

$$\mathbf{S}(t, t') = \mathbf{S}(t, t) e^{-\mathbf{B}^T(t'-t)}. \quad (3)$$

A standard method for analysing Eq.1 in the main text, consists in describing the evolution of the probability distribution  $P(\mathbf{x}, t)$  via the Fokker-Planck equation:

$$\frac{\partial P(\mathbf{x}, t)}{\partial t} = \nabla \cdot [\mathbf{B} \mathbf{x}(t) P(\mathbf{x}, t) + \mathbf{D} \nabla P(\mathbf{x}, t)], \quad (4)$$

where  $\nabla$  denotes the spatial derivative with respect to  $\mathbf{x}$ . Eq. (4) can be rewritten as a continuity equation of the form

$$\frac{\partial P(\mathbf{x}, t)}{\partial t} + \nabla \cdot \mathbf{J}(\mathbf{x}, t) = 0. \quad (5)$$

with the following expression for the probability current (or flux)

$$\mathbf{J}(\mathbf{x}, t) = -\mathbf{D} \nabla P(\mathbf{x}, t) - \mathbf{B} \mathbf{x}(t) P(\mathbf{x}, t) \quad (6)$$

### B. Stationary state and probability current

The Gauss-Markov property of the Ornstein-Uhlenbeck process ensures that the mean and covariances converge exponentially fast toward their respective fixed points, provided the eigenvalues of  $\mathbf{B}$  (which may be complex) have positive real part. The stationary state of the MOU process exhibits Gaussian fluctuations around a mean equal to zero. This corresponds to the time-independent multivariate probability density

$$P(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} (\det \mathbf{S})^{1/2}} \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{S}^{-1} \mathbf{x}\right), \quad (7)$$

where  $\mathbf{S}$  denotes the fixed point of the zero-lag covariance matrix  $\mathbf{S}(t, t)$ . From Eq. (7), the gradient of  $P(\mathbf{x})$  simply reads

$$\nabla P(\mathbf{x}) = \frac{\partial P(\mathbf{x})}{\partial \mathbf{x}} = -P(\mathbf{x}) \mathbf{S}^{-1} \mathbf{x} \quad (8)$$

From Eq. (6), the stationary probability current  $\mathbf{J}(\mathbf{x})$  can thus be rewritten in a compact form

$$\begin{aligned} \mathbf{J}(\mathbf{x}) &= \mathbf{D} P(\mathbf{x}) \mathbf{S}^{-1} \mathbf{x} - \mathbf{B} \mathbf{x} P(\mathbf{x}) \\ &= \boldsymbol{\mu} \mathbf{x} P(\mathbf{x}), \end{aligned} \quad (9)$$

with

$$\boldsymbol{\mu} = \mathbf{D} \mathbf{S}^{-1} - \mathbf{B} \quad (10)$$

### C. Entropy production rate

Going a step further, the (ir)reversibility can be described using thermodynamic variables evaluated for the dynamic process. Using the well-known definition for entropy for the probability distribution  $P(\mathbf{x}, t)$ , now considering its time dependent version, we have

$$e[P] = - \int_{\mathbb{R}^n} P(\mathbf{x}, t) \log P(\mathbf{x}, t) d\mathbf{x}. \quad (11)$$

It can be shown that the rate of the increase of entropy over time can be decomposed into two factors, namely  $\dot{e}[P] = \text{EPR} - \text{HDR}$ , where EPR is the entropy production rate and HDR the heat-dissipation rate [1–3]. The

EPR is the main quantity of interest here, which we denote by  $\Phi$ . Now calculating  $\Phi$  for the time-independent distribution  $P(\mathbf{x})$ , we have

$$\Phi = \int \frac{\mathbf{J}^T(\mathbf{x})\mathbf{D}^{-1}\mathbf{J}(\mathbf{x})}{P(\mathbf{x})}d\mathbf{x} = \langle \mathbf{\Pi}^T \mathbf{D} \mathbf{\Pi} \rangle \quad (12)$$

where  $\mathbf{\Pi}$  is called the the thermodynamic force and is related to  $\mathbf{J}$  by the Onsager's reciprocal relations [1]:

$$\mathbf{\Pi} = \frac{\mathbf{D}^{-1}\mathbf{J}}{P} \quad (13)$$

The heat-dissipation rate can be computed as follows:

$$HDR = \int_{\mathbb{R}^n} \mathbf{D}^{-1}\mathbf{B}\mathbf{x} \cdot \mathbf{J}d\mathbf{x} \quad (14)$$

In the context of the stationary MOU diffusion processes, a general expression for the entropy production rate per unit time in the stationary state is the following [1, 2, 4]

$$\Phi = \int (\nabla \log P(\mathbf{x}) - \mathbf{D}\mathbf{B}\mathbf{x})^T \mathbf{D} (\nabla \log P(\mathbf{x}) - \mathbf{D}^{-1}\mathbf{B}\mathbf{x}) P(\mathbf{x})d\mathbf{x} \quad (15)$$

which can be obtained from (10), (12) and (13) as follows:

$$\begin{aligned} \mu &= \mathbf{D}\mathbf{S}^{-1} - \mathbf{B} \\ \mathbf{D}^{-1}\mu &= \mathbf{S}^{-1} - \mathbf{D}^{-1}\mathbf{B} & \mathbf{D}^{-1}. \\ \mathbf{D}^{-1}\mu\mathbf{x} &= (\mathbf{S}^{-1} - \mathbf{D}^{-1}\mathbf{B})\mathbf{x} & \cdot \mathbf{x} \\ \mathbf{D}^{-1}\mu\mathbf{x}P &= (\mathbf{S}^{-1} - \mathbf{D}^{-1}\mathbf{B})\mathbf{x}P & \cdot P \\ \mathbf{D}^{-1}\mathbf{J} &= (\mathbf{S}^{-1} - \mathbf{D}^{-1}\mathbf{B})\mathbf{x}P & \text{from (9)} \\ \mathbf{\Pi} &= (\mathbf{S}^{-1} - \mathbf{D}^{-1}\mathbf{B})\mathbf{x} & \text{from (13)} \\ \mathbf{\Pi} &= \mathbf{S}^{-1}\mathbf{x} - \mathbf{D}^{-1}\mathbf{B}\mathbf{x} \end{aligned} \quad (16)$$

Now, as  $\nabla \log P(\mathbf{x}) = \mathbf{S}^{-1}\mathbf{x}$ , we obtain (15). From (12)

$$\langle \mathbf{\Pi}^T \mathbf{D} \mathbf{\Pi} \rangle = \langle \mathbf{x}^T (\mathbf{D}^{-1}\mathbf{B} - \mathbf{S}^{-1})^T \mathbf{D} (\mathbf{D}^{-1}\mathbf{B} - \mathbf{S}^{-1}) \mathbf{x} \rangle,$$

we obtain that

$$\Phi = \langle \mathbf{x}^T (\mathbf{D}^{-1}\mathbf{B} - \mathbf{S}^{-1})^T \mathbf{D} (\mathbf{D}^{-1}\mathbf{B} - \mathbf{S}^{-1}) \mathbf{x} \rangle, \quad (17)$$

where the average is taken over the stationary state of the process. From this equation we can verify that when  $\mathbf{S} = \mathbf{B}^{-1}\mathbf{D}$ , then  $\Phi = 0$ . From Eq.4 in the main text, and (10), we have  $\mathbf{D}^{-1}\mathbf{B} - \mathbf{S}^{-1} = \mathbf{D}^{-1}\mathbf{Q}\mathbf{S}^{-1} = -\mathbf{D}^{-1}\mu$ . Thus, from Eq. (17) considering that  $\mathbf{S}$  and  $\mathbf{D}$  are symmetric and  $\mathbf{Q}$  is anti-symmetric we obtain:

$$\Phi = -\langle \mathbf{x}^T \mathbf{S}^{-1}\mathbf{Q}\mathbf{D}^{-1}\mathbf{Q}\mathbf{S}^{-1}\mathbf{x} \rangle = \langle \mathbf{x}^T \mu^T \mathbf{D}^{-1}\mu \mathbf{x} \rangle \quad (18)$$

The entropy production rate  $\Phi$  is non-negative. It is strictly positive if the process is irreversible, and it vanishes only if the process is reversible. Since the stationary state of the MOU is Gaussian with covariance matrix  $\mathbf{S}$ , we have the following property:  $\langle \mathbf{x}^T \mathbf{A}\mathbf{x} \rangle = \text{tr}(\mathbf{S}\mathbf{A})$ , and so

$$\Phi = -\text{tr}(\mathbf{S}^{-1}\mathbf{Q}\mathbf{D}^{-1}\mathbf{Q}) = \text{tr}(\mathbf{S}\mu^T \mathbf{D}^{-1}\mu), \quad (19)$$

which can be written into the following equivalent expressions, that does not involve the covariance matrix  $\mathbf{S}$  nor its inverse explicitly:

$$\Phi = \text{tr}(\mathbf{B}^T \mathbf{D}^{-1}\mathbf{Q}) = -\text{tr}(\mathbf{D}^{-1}\mathbf{B}\mathbf{Q}), \quad (20)$$

## II. METHODS

### A. Empirical covariance from fMRI data

The model is fitted to reproduce the two covariance matrices calculated from the empirical BOLD signals, with zero lag and a lag equal to 1 TR:

$$\widehat{S}_{ij}(0) = \frac{1}{T-2} \sum_{1 \leq t \leq T-1} [x_i(t) - \bar{x}_i][x_j(t) - \bar{x}_j], \quad (21)$$

$$\widehat{S}_{ij}(1) = \frac{1}{T-2} \sum_{1 \leq t \leq T-1} [x_i(t) - \bar{x}_i][x_j(t+1) - \bar{x}_j] \quad (22)$$

Here  $\bar{x}_i$  denotes the mean empirical signal:  $\bar{x}_i = \frac{1}{T} \sum_t x_i(t)$  for all  $i$ , which is used to center the data as all variables  $x_i$  have mean zero in the model. These are the empirical counterparts of the model covariances  $S_{ij}(t, t)$  and  $S_{ij}(t, t+1)$  averaged over time  $t$ .

### B. Parameter estimation of the MOU process

We fit the MOU process from the fMRI time series data for each subject in each sleep condition. We rely on a recent estimation method that tunes the MOU model such that its covariance structure reproduces the matrices in Eq. (21), optimizing its parameters the Jacobian matrix  $-\mathbf{B}$  as well as the input covariance matrix  $2\mathbf{D}$  [5]. Importantly, this optimization procedure incorporates topological constraints on  $\mathbf{B}$ , adjusting only existing anatomical connections, also keeping the input cross-covariances  $D_{ij} = 0$  for  $i \neq j$ . Note that our current notation corresponds to a previous publication [5], using the following  $-\mathbf{B} \leftrightarrow \mathbf{J}$  and  $2\mathbf{D} \leftrightarrow \Sigma$ ; note that  $-\mathbf{B} \leftrightarrow \mathbf{J}^T$  in the subsequent paper [6].

The model is first calibrated by calculating the time constant  $\tau$  from the empirical signals.

$$\tau = -\frac{N}{\sum_{1 \leq i \leq N} a(v_i | u)}, \quad (23)$$

where  $a(v_i | u)$  is the slope of the linear regression of  $v_i = \left[ \log(\widehat{S}_{ii}^0), \log(\widehat{S}_{ii}^1) \right]$  by  $u = [0, 1]$ .

We rely on a gradient descent to iteratively adjust  $\mathbf{B}$  and  $\mathbf{D}$  until reaching the best fit [5]. At each optimization step, we calculate the model counterparts of the covariance matrices in Eq. (21)  $\mathbf{S}(0)$  and  $\mathbf{S}(1)$ , assuming stationarity over each fMRI session. They can be calculated by solving the Lyapunov equation using e.g. the Bartels-Stewart algorithm, which yields here

$$\mathbf{B}\mathbf{S}(0) + \mathbf{S}(0)\mathbf{B}^T = 2\mathbf{D}, \quad (24)$$

once again equating the derivative with zero in Eq. (2), and the equation involving the propagator. We calculate the lagged covariance rewriting Eq. (3) for the time-lag equation here as

$$\mathbf{S}(1) = \mathbf{S}(0)e^{\mathbf{B}^T}. \quad (25)$$

We then calculate the difference between the model and empirical covariances,  $\Delta\mathbf{S}(t) = \hat{\mathbf{S}}(t) - \mathbf{S}(t)$  with  $t \in \{0, 1\}$ . The parameter update is given by differentiating Eqs. (25) and (24):

$$\Delta\mathbf{B} = \epsilon_B [\mathbf{S}(0)]^{-1} [\Delta\mathbf{S}(0) - \Delta\mathbf{S}(1)e^{\mathbf{B}^T}], \quad (26)$$

$$\Delta\mathbf{D} = \epsilon_D \mathbf{B} \Delta\mathbf{S}(0) + \epsilon_D \Delta\mathbf{S}(0) \mathbf{B}^T,$$

with  $\epsilon_B$  and  $\epsilon_D$  small learning rates. The best fit corresponds to minimising the squared norm of both  $\Delta\mathbf{S}(0)$  and  $\Delta\mathbf{S}(1)$ .

The model fitting is quantified by two measures: the model error defined using the matrix distance and the Pearson correlation between the vectorized FC matrices (model versus data). As shown in Fig S1, all sleep states have goodness of fit, with Pearson correlation above 0.6, corresponding to R2 of 0.36.

### C. Decoding of sleep stages

We use the same approach as in previous work [6, 7]. We use the scikit-learn Python library for the implementations of classifiers [8].

We rely on two usual classifiers: multinomial logistic regression (MLR) and the 1-nearest-neighbor (1NN). The features correspond to vectorized  $\mathbf{C}/\mathbf{D}/\mathbf{S}$  matrices after discarding zero or redundant elements. The MLR is a canonical tool for high-dimensional linear classification, which tunes a weight for each feature, thus selecting the important ones to discriminate the classes. In addition, we use L2-regularization for the MLR ( $C = 1.0$  in the scikit-learn implementation). In contrast, the 1NN assigns to a new sample the class to which belongs its closest neighbors with respect to a similarity metric, here chosen as the Pearson correlation coefficient between the feature vectors. It thus relies on the global profile of features to cluster samples into classes.

Following standards, we use a stratified cross-validation scheme with 80% of the samples for the train set and 20% for the test set, where the ratio of classes

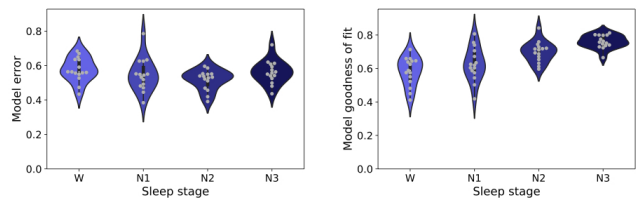


FIG. S1. **Model error and goodness of fit across sleep stages.** Left panel: Model error for each sleep stage (x-axis) across the 15 subjects. The sleep stages are represented by the blue contrasts, from light for wake (W) to dark for the deepest sleep (N3). Right panel: Goodness of fit measured by the Pearson correlation between the vectorized model and empirical FC matrices:  $\mathbf{S}(0)$  with  $\hat{\mathbf{S}}(0)$ , and  $\mathbf{S}(1)$  with  $\hat{\mathbf{S}}(1)$ . The average of the two Pearson correlation values is reported.

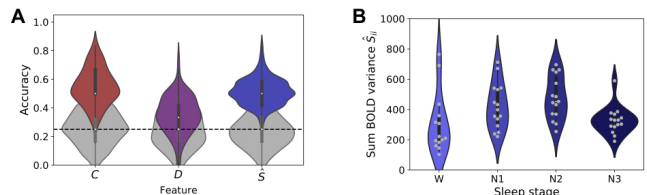


FIG. S2. **Model-based decoding performance outperforms phenomenological classification of sleep stages.**

A) Similar plot to Fig 2C in the main text with the classification accuracy for the 1-nearest-neighbor (1NN) classifier, which relies on a similarity measure (here the Pearson correlation) between the input features to predict the class of the test sample. The x-axis indicates the features: the model estimates  $\mathbf{C}$  and  $\mathbf{D}$ , as well as the empirical FC denoted by  $\hat{\mathbf{S}}$ . B) Similar plot to Fig 2B in the main text but for the model input variance summed over all ROIs.

are the same in both sets. We also use the subject identity as “group information” to avoid mixing subject data between the train and test sets. In practice, we use 50 random splits of the data and report the distribution of accuracies of the 50 splits (see the violin plots).

Fig S2 shows that the decoding of the sleep states by the 1NN classifier, which assigns to a new sample the class to which belongs its closest neighbors with respect to a similarity metric, here chosen as the Pearson correlation coefficient between the feature vectors. It thus relies on the global profile of features to cluster samples into classes.

### D. Complementary analysis

The comparison of the nodal irreversibility across sleep stages shows an overall decrease for all ROIs from W to N3 (Fig S3A), with a pronounced decreases from W to N1 (Fig S3B) and to a lesser extent from N2 to N3 (Fig S3D), although N1 and N2 are rather similar (Fig S3C). Importantly, we can see heterogeneity in the reduction of irreversibility across the ROIs in the transition to deep

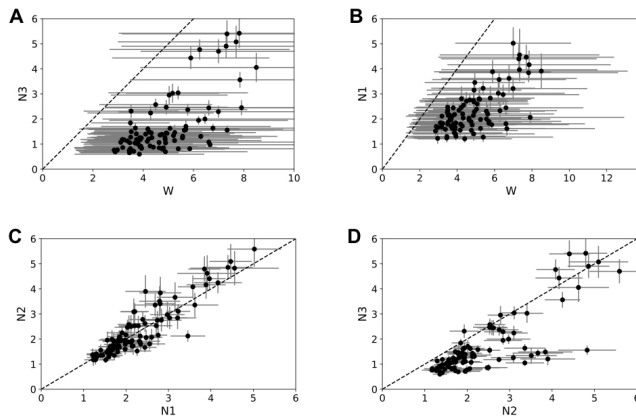


FIG. S3. **Comparison of ROI-specific nodal irreversibility across sleep states.** A) Plot of the sum of absolute values in  $Q$  over each row across the W and N1 states. The error bars indicate the s.e.m. across subjects. This plot is another view of the same data in Fig 3B (main text). B-C) Similar plots to panel A for N1 versus N2, N2 versus N3 and W versus N3.

sleep.

### E. Resting-State fMRI signals

#### Participants

A total of 63 healthy subjects (36 females, mean SD, 2343.3 years) were selected from a data set previously described in a sleep-related study by Tagliazucchi and Laufs [9]. Participants entered the scanner at 7 PM and were asked to relax, close their eyes, and not fight the sleep onset. A total of 52 minutes of resting state activity were measured with a simultaneous combination of EEG and fMRI. According to the rules of the American Academy of Sleep Medicine [10], the polysomnography signals (including the scalp potentials measured

with EEG) determine the classification of data into four stages (wakefulness, N1, N2 and N3 sleep). We selected 15 subjects with contiguous resting state time series of at least 200 volumes to perform our analysis. The local ethics committee approves the experimental protocol (Goethe-Universität Frankfurt, Germany, protocol number: 305/07), and written informed consent was asked to all participants before the experiment. The study was conducted according to the Helsinki Declaration on ethical research.

#### MNI data acquisition

MRI images were acquired on a 3-T Siemens Trio scanner (Erlangen, Germany) and fMRI acquisition parameters were 1505 volumes of T2-weighted echo planar images,  $TR/TE = 2080\text{ ms}/30\text{ ms}$ , matrix 6464, voxel size  $3 \times 3 \times 3\text{ mm}^3$ , distance factor 50%; FOV  $192\text{ mm}^2$ . An optimized polysomnographic setting was employed (chin and tibial EMG, ECG, EOG recorded bipolarly [sampling rate 5 kHz, low pass filter 1 kHz]) with 30 EEG channels recorded with FCz as the reference [sampling rate 5 kHz, low pass filter 250 Hz]. Pulse oxymetry and respiration were recorded via sensors from the Trio [sampling rate 50 Hz] and MR scanner compatible devices (BrainAmp MR+, BrainAmpExG; Brain Products, Gilching, Germany), facilitating sleep scoring during fMRI acquisition.

#### Brain parcellation AAL 90 to extract BOLD time series and filtering

To extract the time series of BOLD signals from each participant in a coarse parcellation, we used the AAL90 parcellation with 90 brain areas anatomically defined in [11]. BOLD signals (empirical or simulated) were filtered with a Butterworth (order 2) band-pass filter in the 0.01-0.1 frequency range.

- 
- [1] H. Qian, Mathematical formalism for isothermal linear irreversibility, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 10.1098/rspa.2001.0811 (2001).
- [2] C. Godrèche and J. M. Luck, Characterising the nonequilibrium stationary states of Ornstein-Uhlenbeck processes, *Journal of Physics A: Mathematical and Theoretical* (2019).
- [3] R. Cofré and C. Maldonado, Information entropy production of maximum entropy markov chains from spike trains, *Entropy* **20**, 10.3390/e20010034 (2018).
- [4] H. Qian, M. Qian, and X. Tang, Thermodynamics of the general diffusion process: Time-reversibility and entropy production, *Journal of Statistical Physics* (2002).
- [5] M. Gilson, R. Moreno-Bote, A. Ponce-Alvarez, P. Ritter, and G. Deco, Estimation of Directed Effective Connectivity from fMRI Functional Connectivity Hints at Asymmetries of Cortical Connectome, *PLoS Computational Biology* **12** (2016).
- [6] M. Gilson, G. Zamora-López, V. Pallarés, M. H. Adhikari, M. Senden, A. T. Campo, D. Mantini, M. Corbetta, G. Deco, and A. Insabato, Model-based whole-brain effective connectivity to study distributed cognition in health and disease, *Network Neuroscience* **4** (2020).
- [7] V. Pallarés, A. Insabato, A. Sanjuán, S. Kühn, D. Mantini, G. Deco, and M. Gilson, Extracting orthogonal subject- and condition-specific signatures from fMRI data using whole-brain effective connectivity, *NeuroImage* **178** (2018).
- [8] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cour-



- napeau, M. Brucher, M. Perrot, and É. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* **12** (2011).
- [9] E. Tagliazucchi and H. Laufs, Decoding Wakefulness Levels from Typical fMRI Resting-State Data Reveals Reliable Drifts between Wakefulness and Sleep, *Neuron* (2014).
- [10] R. B. Berry, R. Brooks, C. E. Gamaldo, S. M. Harding, R. M. Lloyd, C. L. Marcus, and B. V. Vaughn, Aasm — scoring manual version 2.2 the aasm manual for the scoring of sleep and associated events. rules, terminology and technical specifications., *American Academy of Sleep Medicine* **176** (2015).
- [11] N. Tzourio-Mazoyer, B. Landeau, D. Papathanassiou, F. Crivello, O. Etard, N. Delcroix, B. Mazoyer, and M. Joliot, Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain, *NeuroImage* **15**, 10.1006/nimg.2001.0978 (2002).