# An introduction to thermodynamic integration and application to dynamic causal models – Supplementary material

## S1 A primer on Dynamic Causal Models

In this section, we provide a short introduction to dynamic causal modelling (DCM). Since the examples in the main text focus on fMRI data, and we limit our discussion to DCM for fMRI (Friston, Harrison, & Penny, 2003; K. E. Stephan et al., 2008; K. E. Stephan, Weiskopf, Drysdale, Robinson, & Friston, 2007).

DCM for fMRI is characterized by two layers: first, a set of ordinary differential equations that model the dynamics of interacting neuronal states $x$ and local hemodynamic states $h$. Second, the hemodynamic states enter a static nonlinear observation equation that relates venous blood volume and deoxyhemoglobin content to measured BOLD signal changes. In the following, we discuss only the most relevant equations, in order to convey an understanding of the type of problem that model inversion in DCM faces.

The general form of the dynamics of the neuronal layer is

$$\frac{dx}{dt} = f(x, u, \theta_c) \tag{1}$$

where $x = (x_1, \dots, x_N)$ describes the neuronal states of $N$ regions, $u = (u_1, \dots, u_M)$ represents the time series of $M$ experimental manipulations or inputs, and $\theta_c$ are the connectivity parameters that determine the neuronal dynamics. Using a second order Taylor expansion (Stephan et al., 2008), the dynamics $f$ can be approximated as:

$$\frac{dx}{dt} = Ax + \sum_{j=1}^{M} u_j B_j x + Cu + \sum_{i=1}^{N} x_i D_i x. \tag{2}$$

The connectivity parameters $\theta_c$ can be divided into four subsets: The $N \times N$ matrix $A$ describes endogenous connectivity strengths between regions. The set of $N \times N$ matrices $B = \{B_1, \dots, B_M\}$ encodes modulatory effects of inputs on connections between regions. The $N \times M$ matrix $C$ describes the direct effects of driving inputs on regions. Finally, the $N \times N$ matrices $D = \{D_1, \dots, D_N\}$ denote second-order interactions between two regions

28  that affect a third one. Linear DCMs use $A$ and $C$ matrices, bilinear DCMs contain at least

29  one non-zero $B$ matrix, and nonlinear DCMs contain at least one non-zero $D$ matrix.

30  Together $\theta_c = \{A, B, C, D\}$ fully describe the dynamics of the neuronal layer.

31  The hemodynamic model of DCM originates from the Balloon model proposed by Buxton,

32  Wong, and Frank (1998) and extended by Friston (2002) and K. E. Stephan et al. (2007).

33  In brief, it describes how changes in neuronal states locally alter cerebral blood flow,

34  which, in turn, affects venous blood volume and deoxyhemoglobin content. The model

35  consists of a cascade of deterministic differential equations:

36
$$\frac{dh}{dt} = l(h, x, \theta_h), \tag{3}$$

37  where $h = (h_1, \ldots, h_N)$ denotes hemodynamic states in each of $N$ regions. Detailed

38  equations and the meaning of the hemodynamic parameters $\theta_h$ can be found in K. E.

39  Stephan et al. (2007). It is worth noting that the hemodynamic equations are nonlinear

40  and that the original implementation in SPM uses a local (bi)linear approximation

41  (Friston et al., 2003).

42  Finally, hemodynamic states enter a static nonlinear observation equation $g$ with

43  parameters $\theta_g$ that models the BOLD signal $y$:

44
$$y = g(h, \theta_g) + X_0 \beta + \varepsilon \tag{4}$$

45  The term $X_0$ is a matrix of confound regressors that accounts for constant terms and low

46  frequency fluctuations. The Gaussian observation noise $\varepsilon$ is characterized by the

47  covariance matrix $\Pi_\epsilon^{-1}$:

48
$$\varepsilon \sim N(0, \Pi_\epsilon^{-1}). \tag{5}$$

49  The precision matrix $\Pi_\epsilon$ is represented as a linear combination $\Pi_\epsilon = \sum_r \exp(\lambda_r) Q_r$. The

50  precision components $Q_r$ serve to account for temporal autocorrelation and regional

51  differences in noise variance (Friston et al., 2003). Here, we assume that the time series

52  have been whitened and therefore only account for region-specific variances. In this case,

53  each $Q_r$ is a diagonal matrix with diagonal elements belonging to region $r$ set to 1, and 0

54  elsewhere.

55  To complete the generative model, the prior distribution of the parameters $\Theta =$

56  $(\theta_c, \theta_h, \theta_g, \beta)$ and hyperparameters $\Lambda$ needs to be specified. For the results presented in

57  this paper, the priors have been largely matched to SPM8 release 5236

58 (http://www.fil.ion.ucl.ac.uk/spm), except for the scaling of the prior variance of the
59 coefficients of the confound matrix $X_0$, which was adapted to the scaling of the data as
60 explained in S8. All parameters' prior distributions are Gaussian, and when positivity
61 needs to be enforced, an adequate transformation function is used.

62

## S2 Bayesian model comparison and selection

64 In this section, we provide a summary of Bayesian model selection (BMS). Detailed
65 treatments can be found in standard textbooks, such as MacKay (2004).

66 Bayesian inference involves the specification of a probabilistic or generative model $m$
67 with data $y$ and parameters $\theta$. The model has two components: the prior density over $\theta$,
68 $p(\theta|m)$, and the likelihood function $p(y|\theta,m)$. These are combined to form the posterior
69 distribution using Bayes' theorem. Conditioning on a given model $m$, the posterior
70 distribution is:

$$71 \qquad p(\theta|y,m) = \frac{p(y|\theta,m)p(\theta|m)}{p(y|m)}, \qquad (6)$$

$$72 \qquad p(y|m) = \int p(y|\theta,m)p(\theta|m)d\theta. \qquad (7)$$

73 The normalization constant in the denominator, $p(y|m)$, is known as the marginal
74 likelihood or model evidence and corresponds to the likelihood of the data after
75 marginalizing out the parameters of the model.

76 In practice, given the monotonicity of the logarithmic function, either the evidence or its
77 logarithm can be used to score a set of candidate models $m_1, \dots, m_n$ (Bayesian model
78 comparison) and to identify the best model within the model space studied (Bayesian
79 model selection; BMS). One common metric for assessing the relative goodness of two
80 models is the Bayes factor (Kass & Raftery, 1995):

$$81 \qquad B_{i,j} = \frac{p(y \mid m_i)}{p(y \mid m_j)}. \qquad (8)$$

82 or, equivalently, the exponential of the difference in LME of two models.

83 BMS has gained an important role in neuroimaging, not only for DCM but also in other
84 contexts requiring model comparison, such as EEG source reconstruction (Henson,
85 Mattout, Phillips, & Friston, 2009; Wipf & Nagarajan, 2009), or computational

86   neuroimaging (Friston & Dolan, 2010; Klaas E. Stephan, Iglesias, Heinzle, & Diaconescu,
87   2015; K. E. Stephan et al., 2017). Group-level BMS techniques exist which account for
88   individual heterogeneity by treating the model as a random variable in the population
89   (Friston et al., 2016; Rigoux, Stephan, Friston, & Daunizeau, 2014; K. E. Stephan, Penny,
90   Daunizeau, Moran, & Friston, 2009). Finally, Bayesian model averaging allows one to
91   compute an average posterior over models (Penny et al., 2010; Trujillo-Barreto, Aubert-
92   Vázquez, & Valdés-Sosa, 2004), weighted by the posterior probability of each model.
93   Critically, these approaches rely on an accurate estimate of each model's evidence.

94   As mentioned above, except for some special cases, the model evidence cannot be
95   determined analytically, and one typically has to resort to approximations. One
96   computationally efficient option is VB {for textbook treatments, see \Koller, 2009
97   #413;MacKay, 2004 #35}, which provides a lower bound of the LME. An alternative,
98   which we explore in detail here, is MCMC sampling. This family of methods is
99   characterized by simulating a Markov process whose stationary distribution corresponds
100  to the posterior distribution $p(\theta|y, m)$ (for a textbook reference, see Robert & Casella,
101  2010).

102

## S3 A primer on Markov chain Monte Carlo

104  In this section, we provide a short introduction to Markov chain Monte Carlo (MCMC).
105  Thermodynamic integration (TI) requires obtaining samples from a series of power
106  posterior distributions $p_i(\theta|y, m) \propto p(y|\theta, m)^{\beta_i} p(\theta|m)$, with $0 = \beta_0 < \beta_1 \ldots < \beta_N = 1$.
107  An efficient way to achieve this is to use independent Markov chain Monte Carlo (MCMC)
108  samplers (one for each of the $\beta_i$) to generate samples from the power posteriors.

109  MCMC is a powerful technique that can be used to generate samples from any arbitrary
110  target probability distribution $p(x)$, as long as $p(x)$ can be evaluated for any given
111  argument $x$, up to a multiplicative constant $c$. $c$ can be unknown, but has to be constant,
112  i.e. cannot depend on $x$. To this end, the MCMC sampler generates a chain of samples
113  where each sample depends on the previous sample in the chain, but collectively, the set
114  of all samples in the chain are distributed according to the target distribution $p(x)$. To
115  guarantee the latter point, the samples in the chain are generated sequentially according
116  to the following procedure: Let $x_t$ be the last sample currently in the chain, generate a so-
117  called proposal $x'$ via a proposal distribution $q(x'|x_t)$. The simplest way to do this is by

118 adding zero-mean Gaussian noise to $x_t$. Then calculate the so-called Metropolis-Hastings

119 acceptance rate $a$, given by:

$$a = \min\left(1, \frac{p(x')q(x\_t|x')}{p(x_t)q(x'|x_t)}\right).$$

121 Finally, draw a random number $u$ that is uniformly distributed between 0 and 1. If $u < a$,

122 the proposal is accepted and appended to the end of the chain ($x_{t+1} = x'$), otherwise the

123 proposal is rejected and the last sample is repeated ($x_{t+1} = x_t$).

124 Following these steps, it is guaranteed that in the limit of an infinitely long chain, the

125 elements of the chain represent samples from the target distribution, irrespective of the

126 value of the first sample in the chain. More detailed treatments of MCMC can be found in

127 standard textbooks (Brooks, Gelman, Jones, & Meng, 2011). In practice, the fact that MCMC

128 algorithm can only run for a finite time needs to be taken into account. In this context, it

129 is necessary to (1) account for the starting position of the chain and (2) monitor the

130 convergence of the algorithm, i.e. to determine if the MCMC algorithm has already run for

131 long enough such that the elements of the chain can be regarded as approximately

132 representing samples from the desired target distribution.

133 The first problem is typically dealt with by discarding a number of samples at the

134 beginning of the chain (typically the first half). The discarded part of the chain is generally

135 referred to as burn-in period.

136 For the second problem, several techniques have been developed to assess the

137 convergence of a MCMC sampler. One popular method, which is used throughout this

138 paper, is the Gelman-Rubin's potential scale reduction factor $\hat{R}$ (Gelman & Rubin, 1992).

139 This method tests parameter-wise convergence by comparing the variance of segments

140 of the chains. A $\hat{R}$ statistic below 1.1 is a commonly accepted criterion for convergence. To

141 compute this score, the samples of the log likelihood of the first (after the burn-in phase)

142 and last third section of each chain were compared.

143 Since TI already requires obtaining samples from a series of power posterior

144 distributions, convergence of the MCMC samplers can be expedited by adopting a

145 population MCMC approach in which chains associated with neighboring temperatures

146 (i.e., $\beta_i$ and $\beta_{i+1}$) are allowed to interact by means of a "swap" accept-reject (AR) step

147 (McDowell, Dyckman, Austin, & Clementz, 2008; Swendsen & Wang, 1986). In brief,

148 population MCMC defines a joint product distribution

$$149 \qquad \prod_{i=0}^{N} p(\theta_i|y, \beta_i, m) = \prod_{i=0}^{N} \frac{p(y|\theta_i, m)^{\beta_i} p(\theta_i|m)}{Z_i}, \qquad (9)$$

150      where N is the number of distributions or chains. The goal is to obtain samples from this

151      distribution by two types of AR steps: First, local steps are used to sample parameters $\theta_i$

152      from $p_{\beta_i}(\theta_i|y, m)$. Second, samples are obtained using the swapping step in which a set of

153      neighboring parameters $\theta_i, \theta_{i+1}$ are randomly chosen and then exchanged between

154      chains with probability:

$$155 \qquad \min(1, (p(y|\theta_{i+1}, m)^{\beta_i} p(\theta_{i+1}|m) / ((p(y|\theta_i, m)^{\beta_{i+1}} p(\theta_i|m))). \qquad (10)$$

156      This AR step does not change the stationary distribution of any of the chains.

157      Population MCMC can be easily parallelized, with or without exploiting GPUs (Aponte et

158      al., 2016) as each of the chains is independent of the rest of the ensemble. Swapping steps

159      need to be performed serially but, assuming that the likelihood and prior functions have

160      been already evaluated, this method increases the efficiency of the sampling scheme while

161      only inducing negligible computational costs (for example, Aponte et al., 2016;

162      Calderhead & Girolami, 2009). Intuitively, the increase in efficiency is achieved by

163      exploring the sampling space in a way comparable to simulated annealing, i.e., allowing

164      some of the chains to explore the parameter space more freely by relaxing the likelihood

165      function.

166

## 167   S4 Derivation of the equilibrium distribution for the ideal gas example

168      In this section, we present the derivation of the equilibrium distribution for the ideal gas

169      example in the main text. As outlined in the main text, the equilibrium distribution is

170      attained as the maximum entropy solution, which can be found using a variational

171      Lagrangian with two constraints represented by the Lagrange multipliers $\lambda_1$ and $\lambda_2$ (see

172      Blundell & Blundell, 2009; Jaynes, 1957):

$$173 \qquad \frac{\delta}{\delta q}\left[-k_B \int q(\theta) \ln q(\theta)\, d\theta - \lambda_1 \left(\int q(\theta)\phi(\theta)d\theta - U\right) - \lambda_2 \left(\int q(\theta)d\theta - 1\right)\right] = 0. \qquad (11)$$

174      Noting that

$$175 \qquad -\frac{\delta}{\delta q} k_B \int q(\theta) \ln q(\theta)\, d\theta = k_B(-1 - \ln q(\theta)) \qquad (12)$$

$$-\frac{\delta}{\delta q}\lambda_1\left(\int q(\theta)\phi(\theta)d\theta - U\right) = -\lambda_1\phi(\theta), \tag{13}$$

$$-\frac{\delta}{\delta q}\lambda_2\left(\int q(\theta)d\theta - 1\right) = -\lambda_2, \tag{14}$$

178 the Lagrangian yields

$$k_B\ln q(\theta) = -\lambda_1\phi(\theta) - \lambda_2 - k_B, \tag{15}$$

$$q(\theta) = \frac{1}{\exp\left(\frac{\lambda_2}{k_B} + 1\right)}\exp\left(-\frac{\lambda_1}{k_B}\phi(\theta)\right). \tag{16}$$

181 The term $\lambda_1$ constitutes the definition of inverse temperature in statistical physics
182 (Blundell & Blundell, 2009; Jaynes, 1957):

$$\frac{1}{T} \stackrel{\text{def}}{=} \lambda_1. \tag{17}$$

184 The term $\frac{\lambda_1}{k_B} = \frac{1}{k_B T}$ is commonly represented by the symbol $\beta$. In order to derive the
185 second constant $\lambda_2$, we write:

$$q(\theta) = \frac{1}{Z}\exp\left(-\frac{\phi(\theta)}{k_B T}\right), \tag{18}$$

187 where $Z$ is referred to as the partition function of the system:

$$Z \stackrel{\text{def}}{=} \int \exp\left(-\frac{\phi(\theta)}{k_B T}\right)d\theta. \tag{19}$$

189 Hence, the term $\exp\left(\frac{\lambda_2}{k_B} + 1\right)$ is the normalization constant of $q(\theta)$, and thus $\lambda_2 =$
190 $k_B(\ln Z - 1)$.

191

192

## S5 Variational Bayes under the Laplace approximation for DCM

194 This section introduces the variational Bayes under the Laplace (VBL) approximation for
195 inverting dynamic causal models. For an in-depth discussion see (Friston, Mattout,
196 Trujillo-Barreto, Ashburner, & Penny, 2007).

197    Commonly, in order to maximize $-F_{VB}$, a mean field approximation of $q$ is used. In other
198    words, the distribution $q$ is assumed to factorize into different sets of parameters, each of
199    which defines a more tractable optimization problem. In the case of DCM, $q$ is assumed to
200    have the form:

$$q(\Theta, \Lambda) = q(\Theta)q(\Lambda), \tag{20}$$

201

202    i.e., the parameters $\Theta = (\theta_c, \theta_h, \theta_g, \beta)$ and the hyperparameters $\Lambda$ are assumed to be
203    conditionally independent. The functional $-F_{VB}$ can be optimized iteratively with respect
204    to $\Theta$ and $\Lambda$ converging to a maximum $-F_{VB} \leq \ln p(y|m)$ (Koller, 2009). This rests on
205    maximizing the variational energies:

$$\ln q(\Theta) = \int q(\Lambda) \ln p(y, \Theta, \Lambda) d\Lambda + c_\Theta, \tag{21}$$

206

$$\ln q(\Lambda) = \int q(\Theta) \ln p(y, \Theta, \Lambda) d\Theta + c_\Lambda. \tag{22}$$

207

208    where $c_\Theta$ and $c_\Lambda$ are constants with respect to $\Theta$ and $\Lambda$, respectively. In DCM, it is typically
209    assumed that all terms are Gaussian (but see Raman, Deserno, Schlagenhauf, and Stephan
210    (2016) and Yao et al. (2018) who used conjugate priors for the noise terms).

211    Despite the mean field approximation, the integrals in Eq. 22 and 21 and cannot be solved
212    analytically because of the nonlinearities of the forward model (Eq. 4). This problem is
213    circumvented by approximating the log of the unnormalized posterior with a second
214    order Taylor expansion on a local maximum (or equivalently, the unnormalized posterior
215    is assumed to be Gaussian) and optimizing the objective function $\ln p(y, \Theta, \Lambda)$ through
216    gradient ascent (but see Lomakina et al. (2015) for an alternative based on Gaussian
217    processes). This approach is called the Laplace approximation (Friston et al., 2007) and
218    underlies other methods such as BIC (Schwarz, 1978) or when the normalization constant
219    of an approximate, tractable posterior is directly used (Kass & Raftery, 1995). As a
220    consequence of this approximation, the variational free energy is no longer guaranteed to
221    represent a lower bound on the log evidence (Wipf & Nagarajan, 2009). **A detailed**
222    **treatment of VBL can be found in Friston et al. (2007). In section S9, we present a**
223    **simplified version of the derivation of the VBL estimate of the free energy and an**
224    **explicit expression for the accuracy term.**

225 The VBL algorithm used here was the implementation available in the software package
226 SPM8 (release 5236), which employs a gradient ascent scheme to optimize the marginal
227 distributions $q(\Theta)$ and $q(\Lambda)$ (Friston et al., 2007).

228

## S6 Conventional sampling-based estimation of model evidence

230 In this section, we provide summaries to two popular sampling-based estimators for the
231 log model evidence: the prior arithmetic mean estimator (AME) and the posterior
232 harmonic mean estimator (HME).

233 **Prior arithmetic mean estimator (AME)**

234 Importance sampling is a Monte Carlo method for approximating the expected value of a
235 random variable $h(X)$ under the density $p$ by means of an auxiliary density function $w(X)$,
236 which is required to be absolutely continuous with respect to $p$ (Robert & Casella, 2010;
237 p. 92, Def. 3.9), or less formally, the auxiliary density $w$ should share the same support as
238 $p$ to avoid zeros in the denominator:

$$\int h(x)p(x)dx = \int \frac{h(x)p(x)w(x)}{w(x)}dx. \tag{23}$$

240 From the strong law of large numbers, if this expected value exists, the process

$$\lim_{K\to\infty} \frac{1}{K}\sum_{k=1}^{K} h(x_i)\frac{p(x_k)}{w(x_k)} \tag{24}$$

242 converges almost surely to Eq. 9 when the samples $x_1, \ldots, x_K$ have been drawn from the
243 auxiliary distribution $w$.

244 In order to approximate the model evidence by importance sampling, the simplest choice
245 of the auxiliary density is the prior distribution, $w = p(\theta \mid m)$. This results in the prior
246 arithmetic mean estimator (AME):

$$p(y|m) = \int p(y|\theta,m)p(\theta|m)d\theta = \int p(y|\theta,m)p(\theta|m)\frac{p(\theta|m)}{p(\theta|m)}d\theta, \tag{25}$$

$$p_{AME} = \frac{1}{K}\sum_{k=1}^{K} p(y|\theta_k,m). \tag{26}$$

249 where samples $\theta_k$ have been obtained from the prior distribution $p(\theta|m)$. Because
250 samples of the likelihood $p(y|\theta,m)$ can greatly exceed the range of double precision

251   floating point numbers, it is necessary to normalize the likelihood function in log space.

252   This can be achieved with the following formula:

$$\ln p_{AME} = \ln \alpha - \ln K + \ln \sum_{i=1}^{K} \exp[\ln p(y|\theta_i, m) - \ln \alpha], \tag{27}$$

254   where $\alpha > 0$ is an arbitrary constant. In all analyses reported here, $\alpha$ was set to

255   $\max_{k} p(y|\theta_k, m)$.

256   A serious shortcoming of AME is that in the great majority of situations most samples

257   drawn from the prior have very low likelihood. Therefore, an extremely large number of

258   samples is required to ensure that high likelihood regions of the parameter space are

259   taken into account by the estimator; otherwise, the estimator suffers from high variance

260   (Vyshemirsky & Girolami, 2008).

261   **Posterior harmonic mean estimator (HME)**

262   The second choice for the auxiliary density is the posterior distribution, which results in

263   the posterior harmonic mean estimator (HME). This estimator has received divergent

264   appraisals in the literature as a method for computing the LME (for example, Kass &

265   Raftery, 1995; Wolpert & Schmidler, 2012). Re-expressing the model evidence, the HME

266   can be derived as follows:

$$\frac{1}{p(y|m)} = \int \frac{p(\theta|m)}{p(y|m)} d\theta,$$

$$= \int \frac{p(y|\theta, m)p(\theta|m)}{p(y|\theta, m)p(y|m)} d\theta,$$

$$= \int \frac{p(\theta|y, m)}{p(y|\theta, m)} d\theta \tag{28}$$

$$p_{HME} = \left( \frac{1}{K} \sum_{i=1}^{K} \frac{1}{p(y|\theta_i, m)} \right)^{-1}. \tag{29}$$

271   Here, samples $\theta_i$ are drawn from the posterior distribution $p(\theta|y, m)$.

272   In order to avoid numerical instabilities, it is again necessary to normalize in log space,

273   using the formula

$$\ln p_{HME} = \ln K + \ln \alpha - \ln \sum_{i=1}^{K} \exp[-\ln p(y|\theta_i, m) + \ln \alpha]. \tag{30}$$

275   Here, $\ln \alpha$ has been chosen to be $\max_{i} -\ln p(y|\theta_i, m)$.

276  A disadvantage of HME is that its variance might be infinite when the likelihood function

277  is not heavy-tailed (Raftery, Newton, Satagopan, & Krivitsky, 2007), which has serious

278  consequences for the convergence rate of a wide variety of models (Wolpert & Schmidler,

279  2012). A second problem is that the samples used for HME are obtained from the posterior

280  distribution only. This leads to the opposite behavior as for AME: because the contribution

281  of the prior to the LME might not be appropriately accounted for, the HME tends to

282  overestimate the model evidence, a behavior that can be difficult to diagnose (Lartillot &

283  Philippe, 2006). Several improvements of the HME have been proposed to account for this

284  shortcoming (for example, Raftery et al., 2007).

285

286  **Implementation**

287  Since TI requires samples from both the prior and the posterior distribution, which

288  correspond to the power posteriors with $\beta = 0$ and $\beta = 1$, respectively, the samples

289  acquired for TI can be used for computing the other sampling-based estimators, AME and

290  HME. In our comparisons throughout this paper, we have used this technique to ensure

291  that any observed differences between estimators are not simply due to differences in the

292  implementation of the samplers.

293

294  # S7 Connectivity parameters of the synthetic models

295  The connectivity parameters of the synthetic models used here are shown below.

296  *Model 1*

297  Model one did not include any bilinear or non-linear terms.

298
$$A = \begin{pmatrix} -0.5 & 0 & 0 \\ 0 & -0.5 & 0 \\ 0 & 0 & -0.5 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}.$$

299  *Model 2*

300  Models 2 to 5 used the same A and C matrices. In addition, models 2 to 4 included one

301  bilinear term (B matrices), and model 5 included a nonlinear term (D matrices).

$$302 \qquad A = \begin{pmatrix} -0.5 & 0 & -0.25 \\ 0 & -0.5 & -0.25 \\ 0.5 & 0.5 & -0.5 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix},$$

$$303 \qquad B_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 3 & 0 \end{pmatrix}, \quad B_2 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

304 *Model 3*

305 Because model 3 shared the same A and C matrix with model 2, we only display the B

306 matrices.

$$307 \qquad B_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad B_2 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 3 & 0 & 0 \end{pmatrix}.$$

308 *Model 4*

309 Again, only the B matrices differed between models 2, 3, and 4.

$$310 \qquad B_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad B_2 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -2 \end{pmatrix}.$$

311 *Model 5*

312 Model 5 included no bilinear term but included one non-linear term.

$$313 \qquad D_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad D_2 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \quad D_3 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

314 The exact data structures can be downloaded from the ETH research collection (ETH

315 Zurich, 2020).

316

# S8 Scaling of BOLD signals

318 In the SPM version used here (5236), BOLD signals $y$ are rescaled with respect to their $\ell_\infty$

319 norm, such that

$$320 \qquad ||y||_\infty = 4. \qquad (A.1)$$

321 In DCM, the observation equation (see Eq. 4) can be written as

$$322 \qquad\qquad y = g\left(h, \theta_g\right) + X_0 \beta + \varepsilon \qquad\qquad (A.2)$$

323 where $X_0$ represents confounding factors. This matrix usually consists of cosine functions
324 that account for baseline effects and low frequency components and can be imagined as
325 implementing a model of structured noise (scanner-related fluctuations in signal
326 intensity) that is distinct from the model's residuals. We assume N observations such that
327 data from a region is $y[t]$, $t = 0, \dots, N-1$, and the components of $X_0 = [x_K, \dots, x_{M-1}]^T$,
328 $K > 0$ are

$$329 \qquad\qquad x_k[t] = \cos(\frac{2\pi k t}{N}). \qquad\qquad (A.3)$$

330 In this case, $X_0^T X_0$ is a diagonal matrix as all base functions are orthogonal. The diagonal
331 elements are given by

$$332 \qquad\qquad \sum_{n=0}^{N-1} \cos\left(\frac{2\pi \omega n}{N}\right)^2 = \frac{N}{2}, \qquad\qquad (A.4)$$

333 Thus,

$$334 \qquad\qquad X_0^T X_0 = \frac{N}{2} I. \qquad\qquad (A.5)$$

335 The posterior variance of the regressors conditioned on the predictions from DCM, the
336 variance of the error $\sigma_c^2$, and the prior variance $\sigma_0$, is

$$337 \qquad (\sigma_c^{-2} X^T X + \sigma_0^{-2} I)^{-1} = \left(\frac{\sigma_c^{-2} N}{2} I + \sigma_0^{-2} I\right)^{-1}, \qquad (A.6)$$

338

$$339 \qquad\qquad = \left(\frac{N}{2\sigma_c^2} + \frac{1}{\sigma_0^2}\right)^{-1} I. \qquad\qquad (A.7)$$

340 To derive the prior variance of the signal predicted by $X_0 \beta$, we note that for the predicted
341 signal $y$:

$$342 \qquad E[y[t]^2] = E\left[\left(\sum_{\omega=K}^{M-1} \beta_\omega \cos\frac{2\pi t \omega}{N}\right)^2\right], \qquad (A.8)$$

$$343 \qquad = E\left[\sum_{\omega,k=K}^{M-1} \beta_\omega \beta_k \cos\frac{2\pi t \omega}{N} \cos\frac{2\pi t k}{N}\right]. \qquad (A.9)$$

344 Because the coefficients are assumed to be uncorrelated and to have zero mean, it follows

345 that

$$= \sum_{\omega=K}^{M-1} Var(\beta_\omega) \cos^2 \left(\frac{2\pi t \omega}{N}\right), \tag{A.10}$$

$$= \sigma_0^2 \left(\sum_{\omega=0}^{M-1} \cos^2 \left(\frac{2\pi t \omega}{N}\right) - \sum_{\omega=0}^{K-1} \cos^2 \left(\frac{2\pi t \omega}{N}\right)\right). \tag{A.11}$$

348 Assuming that $2Mt/N$ is an integer, it follows that

$$= \sigma_0^2 \left(\frac{M}{2} - \sum_{\omega=0}^{K-1} \cos^2 \left(\frac{2\pi t \omega}{N}\right)\right). \tag{A.12}$$

350 It follows that

$$\frac{\sigma_0^2 (M-K)}{2} \leq E[y[t]^2] = Var(y[t]) \leq \frac{\sigma_0^2 M}{2}. \tag{A.13}$$

352 This constitutes an approximation to the prior variance of the signal. Although in the SPM

353 implementation of DCM used here, $\sigma_0^2$ is set to $10^8$, here we use a more pragmatic value

354 $\sigma_0 = ||y||_\infty = 4$. From Eq. A.12, it can be seen that this constitutes a more conservative

355 prior variance than the SPM implementation, but still liberal enough to a priori easily

356 account for the totality of the variance in the data.

357

## S9 Derivation of variational negative free energy under the Laplace

## approximation

360 The expression for the variational negative free energy can be derived by noting that Eq.

361 34 in the main text can be written as an energy term plus an entropy term

$$-F_{VB} = E[\ln p(y, \theta)]_{q(\theta)} - E[\ln q(\theta)]_{q(\theta)}. \tag{A.14}$$

363 For simplicity, in the rest of this section, we collapse parameters $\theta$ and hyperparameters

364 $\Lambda$ into a $d$-dimensional vector $\theta$, assuming that a maximum has been obtained. Also, we

365 assume that all densities are conditioned on model $m$, and make this assumption implicit.

366 Moreover, we assume that the prior distribution of parameters $\theta$ is a Gaussian

367 distribution centered at $\theta_0$ with covariance $\Pi_0^{-1}$.

368 According to the Laplace approximation, $q(\theta)$ is a Gaussian distribution with mean $\theta^* = $

369 $\arg\max_{\theta} p(y, \theta)$ and variance

370
$$\Pi = -\frac{\partial^2 \ln p(y, \theta)}{\partial \theta^2} = \Pi_0 - \frac{\partial^2 \ln p(y|\theta)}{\partial \theta^2}. \tag{A.15}$$

371 We denote the negative Hessian of the likelihood or observed Fisher information in the

372 following as $\Pi_L$.

373 The energy term in Eq. A.14 is approximated using the Laplace method, which yields

374
$$E[\ln p(y, \theta)]_{q(\theta)} \approx \ln p(y, \theta^*) - \frac{1}{2} E[(\theta^* - \theta)'\Pi(\theta^* - \theta)]_{q(\theta)}, \tag{A.16}$$

375
$$= \ln p(y, \theta^*) - \frac{1}{2} tr\left(\Pi E[(\theta^* - \theta)(\theta^* - \theta)']_{q(\theta)}\right), \tag{A.17}$$

376
$$= \ln p(y, \theta^*) - \frac{1}{2} tr(\Pi \Pi^{-1}) = \ln p(y, \theta^*) - \frac{1}{2} d. \tag{A.18}$$

377 where $tr$ denotes the trace operator.

378 The last term in Eq. A.14 is the entropy of a Gaussian distribution, which is given by:

379
$$-E[\ln q(\theta)]_{q(\theta)} = \frac{1}{2}(d \ln 2\pi + d - \ln|\Pi|). \tag{A.19}$$

380 where $\Pi$ is the precision of $q$.

381 Plugging Eqs. A.18 and A.19 into Eq. A.14, the variational free energy is given by

382
$$-F_{VB} = \ln p(y, \theta^*) + \frac{1}{2}(d \ln 2\pi - \ln|\Pi|). \tag{A.20}$$

383 The first term on the right of Eq. A.20 can be expanded to obtain the full expression:

384
$$\ln p(y, \theta^*) = \ln p(y|\theta^*) + \ln p(\theta^*), \tag{A.21}$$

385
$$= \ln p(y|\theta^*) - \frac{1}{2} d \ln 2\pi + \frac{1}{2} \ln|\Pi_0| - \frac{1}{2}(\theta^* - \theta_0)'\Pi_0(\theta^* - \theta_0). \tag{A.22}$$

386 where $\theta_0$ and $\Pi_0$ are the mean and precision of the prior density, respectively. By

387 inserting Eq. **Error! Reference source not found.** into Eq. A.20, the scheme proposed

388 by Friston et al. (2007) can be written as:

389
$$-F_{VB} = \ln p(y|\theta^*) + \frac{1}{2}\ln\frac{|\Pi_0|}{|\Pi|} - \frac{1}{2}(\theta^* - \theta_0)'\Pi_0(\theta^* - \theta_0). \tag{A.23}$$

390    Although VBL is typically orders of magnitude faster than MCMC sampling, it exhibits

391    several limitations: it is susceptible to (i) local extrema, (ii) violations of the distributional

392    assumptions imposed on the posterior, (iii) violations of the conditional independence

393    assumptions of the mean field approximation (see Daunizeau, David, & Stephan, 2011 for

394    discussion), and (iv) it is only defined when the Hessian in Eq. A. 15 is not singular.

395    Returning to our theme of connecting TI to VBL, one can write the variational negative

396    free energy in terms of an approximate accuracy and complexity term (Eq. **Error!**

397    **Reference source not found.**). One observes that the accuracy term can be computed as

$$-F_{VB} + KL\big(q(\theta)||p(\theta)\big) = A_{VB}. \tag{A.24}$$

399    Given a Gaussian prior and posterior, the KL divergence has the following analytical form:

$$KL\big(q(\theta)||p(\theta)\big) = \frac{1}{2}\left[\ln\frac{|\Pi|}{|\Pi_0|} + tr(\Pi_0\Pi^{-1}) - d + (\theta^* - \theta_0)'\Pi_0(\theta^* - \theta_0)\right]. \tag{A.25}$$

401    Replacing terms, we obtain

$$A = E[\ln p(y|\theta)]_{q(\theta)}, \tag{A.26}$$

$$\approx A_{VB} = \ln p(y|\theta^*) + \frac{tr(\Pi_0\Pi^{-1})}{2} - \frac{d}{2}. \tag{A.27}$$

404    A more familiar expression for the accuracy can be derived by noting that the posterior

405    covariance can be written as the sum of the negative Hessian of the likelihood plus the

406    prior covariance, such that

$$A_{VB} = \ln p(y|\theta^*) + \frac{1}{2}tr\left(\frac{\Pi_0 + \Pi_L - \Pi_L}{\Pi_0 + \Pi_L}\right) - \frac{d}{2}, \tag{A.28}$$

$$= \ln p(y|\theta^*) - \frac{1}{2}tr\left(\frac{\Pi_L}{\Pi_0 + \Pi_L}\right), \tag{A.29}$$

$$\mathbb{p} = tr\left(\frac{\Pi_L}{\Pi_0 + \Pi_L}\right). \tag{A.30}$$

410    $\mathbb{p}$ is the effective number of parameters proposed by Moody (1991) Eq. 18 and see

411    Spiegelhalter, Best, Carlin, and van der Linde (2002) Eq. 15 and is commonly used for

412    model selection.

413

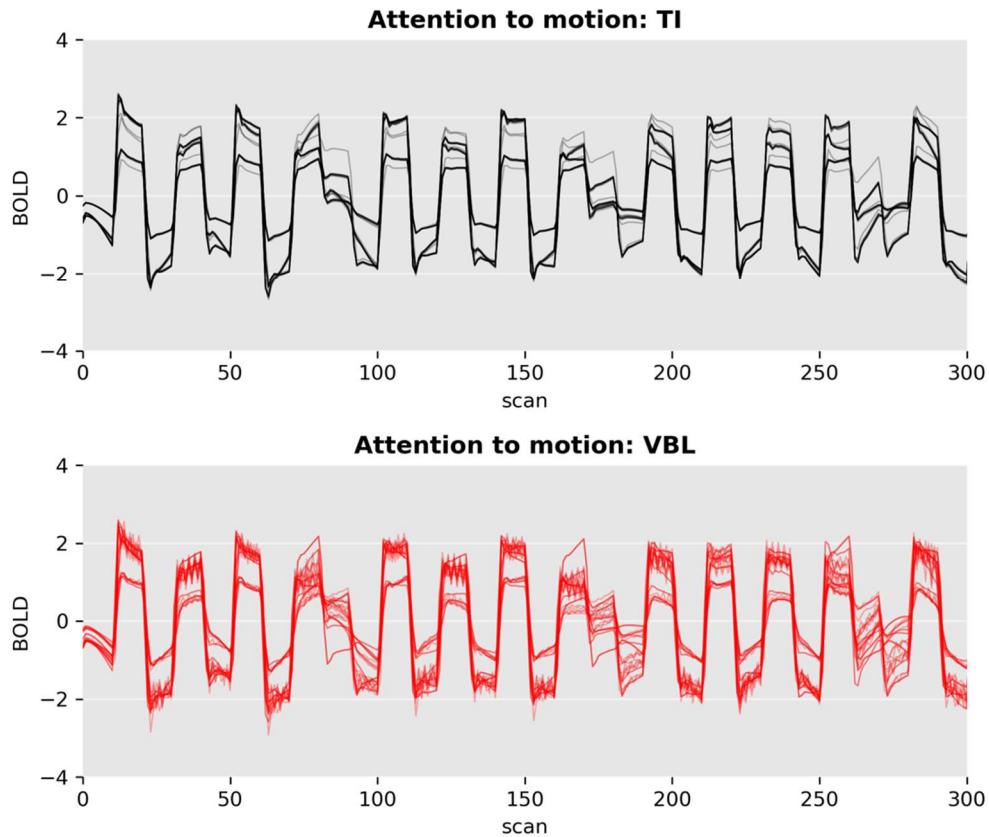## S10 Predicted fMRI time series for the attention to motion dataset



**Figure S1.** Comparison of 10 predicted BOLD signal trajectories (for the MAP estimate) of model $m_4$ between TI and VBL for the "attention to motion" dataset from Buchel (1997). In order to obtain an unbiased impression of the variability, the predicted BOLD responses are plotted in full (i.e., including estimated confounds; compare Eq. 4). Both estimates are qualitatively similar, but VBL fits display higher variability.

415

## S11 Final step in the derivation of the fundamental TI equation

417    Applying the chain rule of differentiation to the logarithm of a positive-valued function,
418    we have the following relation:

419
$$\frac{d}{d\beta} \ln f(\beta) = \frac{1}{f(\beta)} \frac{d}{d\beta} f(\beta)$$

420    In the main text section Thermodynamic Integration and the origin of free energy, we
421    have shown that the log-model evidence is given by the expression (Eq. 22 main text):

422
$$\ln p(y|m) = \int_{\beta=0}^{\beta=1} \frac{d}{d\beta} \ln \int p(y|\theta,m)^{\beta} p(\theta|m) \, d\theta \, d\beta$$

423  Applying the above relation with $f(\beta) = \int p(y|\theta, m)^\beta p(\theta|m) \, d\theta = Z_\beta$, we have

$$424 \quad \frac{d}{d\beta} \ln \int p(y|\theta, m)^\beta p(\theta|m) \, d\theta = \frac{\frac{d}{d\beta} \int p(y|\theta, m)^\beta p(\theta|m) \, d\theta}{\int p(y|\theta, m)^\beta p(\theta|m) \, d\theta}$$

$$425 \quad = \frac{1}{Z_\beta} \int \frac{d}{d\beta} p(y|\theta, m)^\beta p(\theta|m) \, d\theta$$

$$426 \quad = \frac{1}{Z_\beta} \int p(y|\theta, m)^\beta p(\theta|m) \ln p(y|\theta, m) \, d\theta$$

$$427 \quad = \int \frac{p(y|\theta, m)^\beta p(\theta|m)}{Z_\beta} \ln p(y|\theta, m) \, d\theta.$$

428  Note that the last line above is the integrand in Eq. 23 in the main text. Also note that in
429  the second line above, we have exchanged the derivative with respect to $\beta$ with the
430  integration over $\theta$ and in the third line, we have used the derivative of an exponential
431  function:

$$432 \quad \frac{d}{d\beta} a^\beta = a^\beta \ln a.$$

433

## References

435  Aponte, E. A., Raman, S., Sengupta, B., Penny, W., Stephan, K. E., & Heinzle, J. (2016).
436      mpdcm: A toolbox for massively parallel dynamic causal modeling. *Journal of*
437      *Neuroscience        Methods,        257,        7-16.*
438      doi:http://dx.doi.org/10.1016/j.jneumeth.2015.09.009

439  Blundell, S. J., & Blundell, K. M. (2009). *Concepts in Thermal Physics*. Oxford: Oxford
440      University Press, Incorporated.

441  Brooks, S., Gelman, A., Jones, G., & Meng, X. L. (2011). *Handbook of Markov chain Monte*
442      *Carlo* (S. Brooks, A. Gelman, G. Jones, & X. L. Meng Eds.). New York: Chapman & Hall.

443  Buchel, C. (1997). Modulation of connectivity in visual pathways by attention: cortical
444      interactions evaluated with structural equation modelling and fMRI. *Cerebral*
445      *cortex (New York, N.Y. 1991), 7*(8), 768-778. doi:10.1093/cercor/7.8.768

446  Buxton, R. B., Wong, E. C., & Frank, L. R. (1998). Dynamics of blood flow and oxygenation
447      changes during brain activation: The balloon model. *Magnetic Resonance in*
448      *Medicine, 39*(6), 855-864. doi:10.1002/mrm.1910390602

449  Calderhead, B., & Girolami, M. (2009). Estimating Bayes factors via thermodynamic
450      integration and population MCMC. *COMPUTATIONAL STATISTICS & DATA*
451      *ANALYSIS, 53*, 4028-4045. doi:10.1016/j.csda.2009.07.025

452    Daunizeau, J., David, O., & Stephan, K. E. (2011). Dynamic causal modelling: A critical
453        review of the biophysical and statistical foundations. *NeuroImage, 58*(2), 312-322.
454        doi:http://dx.doi.org/10.1016/j.neuroimage.2009.11.062

455    ETH Zurich. (2020). ETH Research Collection. Retrieved from https://www.research-
456        collection.ethz.ch/bitstream/handle/20.500.11850/301664/simulation_dcms.zi
457        p

458    Friston, K. J. (2002). Bayesian Estimation of Dynamical Systems: An Application to fMRI.
459        *NeuroImage, 16*(2), 513-530. doi:http://dx.doi.org/10.1006/nimg.2001.1044

460    Friston, K. J., & Dolan, R. J. (2010). Computational and dynamic models in neuroimaging.
461        *NeuroImage (Orlando, Fla.), 52*(3), 752-765.
462        doi:10.1016/j.neuroimage.2009.12.068

463    Friston, K. J., Harrison, L., & Penny, W. (2003). Dynamic causal modelling. *NeuroImage,*
464        *19*(4), 1273-1302. doi:10.1016/S1053-8119(03)00202-7

465    Friston, K. J., Litvak, V., Oswal, A., Razi, A., Stephan, K. E., van Wijk, B. C. M., . . . Zeidman, P.
466        (2016). Bayesian model reduction and empirical Bayes for group (DCM) studies.
467        *NeuroImage, 128*(Supplement C), 413-431.
468        doi:https://doi.org/10.1016/j.neuroimage.2015.11.015

469    Friston, K. J., Mattout, J., Trujillo-Barreto, N., Ashburner, J., & Penny, W. (2007). Variational
470        free energy and the Laplace approximation. *NeuroImage, 34*(1), 220-234.
471        doi:http://dx.doi.org/10.1016/j.neuroimage.2006.08.035

472    Gelman, A., & Rubin, D. B. (1992). Inference from Iterative Simulation Using Multiple
473        Sequences. *Statistical Science, 7*(4), 457-472. doi:10.1214/ss/1177011136

474    Henson, R. N., Mattout, J., Phillips, C., & Friston, K. J. (2009). Selecting forward models for
475        MEG source-reconstruction using model-evidence. *NeuroImage (Orlando, Fla.),*
476        *46*(1), 168-176. doi:10.1016/j.neuroimage.2009.01.062

477    Jaynes, E. T. (1957). Information Theory and Statistical Mechanics. *Physical Review,*
478        *106*(4), 620-630. doi:10.1103/physrev.106.620

479    Kass, R. E., & Raftery, A. E. (1995). Bayes Factors. *Journal of the American Statistical*
480        *Association, 90*(430), 773-795. doi:10.1080/01621459.1995.10476572

481    Koller, D. (2009). *Probabilistic graphical models : principles and techniques*. Cambridge,
482        Mass: MIT Press.

483    Lartillot, N., & Philippe, H. (2006). Computing Bayes Factors Using Thermodynamic
484        Integration. *Systematic Biology, 55*(2), 195-207.
485        doi:10.1080/10635150500433722

486    Lomakina, E. I., Paliwal, S., Diaconescu, A. O., Brodersen, K. H., Aponte, E. A., Buhmann, J.
487        M., & Stephan, K. E. (2015). Inversion of hierarchical Bayesian models using
488        Gaussian processes. *NeuroImage, 118*, 133-145.
489        doi:http://dx.doi.org/10.1016/j.neuroimage.2015.05.084

490    MacKay, D. J. C. (2004). *Information Theory, Inference, and Learning Algorithms* (Repr. with
491        corr. ed.). Cambridge: Univ. Press.

492    McDowell, J. E., Dyckman, K. A., Austin, B. P., & Clementz, B. A. (2008). Neurophysiology
493        and neuroanatomy of reflexive and volitional saccades: Evidence from studies of
494        humans. *Brain and cognition, 68*(3), 255-270. doi:10.1016/j.bandc.2008.08.016

495  Moody, J. E. (1991). *The effective number of parameters: an analysis of generalization and*
496       *regularization in nonlinear learning systems*. Paper presented at the Proceedings of
497       the 4th International Conference on Neural Information Processing Systems,
498       Denver, Colorado.

499  Penny, W., Stephan, K. E., Daunizeau, J., Rosa, M. J., Friston, K. J., Schofield, T. M., & Leff, A.
500       P. (2010). Comparing Families of Dynamic Causal Models. *PLOS Computational*
501       *Biology, 6*(3), e1000709. doi:10.1371/journal.pcbi.1000709

502  Raftery, A., Newton, M., Satagopan, J., & Krivitsky, P. (2007). Estimating the Integrated
503       Likelihood via Posterior Simulation Using the Harmonic Mean Identity. In J. M.
504       Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, & M.
505       West (Eds.), *Bayesian Statistics 8* (pp. 1-45). Oxford: Oxford University Press.

506  Raman, S., Deserno, L., Schlagenhauf, F., & Stephan, K. E. (2016). A hierarchical model for
507       integrating unsupervised generative embedding and empirical Bayes. *Journal of*
508       *Neuroscience         Methods,         269*,                6-20.
509       doi:http://dx.doi.org/10.1016/j.jneumeth.2016.04.022

510  Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014). Bayesian model selection
511       for group studies — Revisited. *NeuroImage (Orlando, Fla.), 84*, 971-985.
512       doi:10.1016/j.neuroimage.2013.08.065

513  Robert, C. P., & Casella, G. (2010). *Monte Carlo statistical methods* (2nd ed. ed.): New York
514       : Springer.

515  Schwarz, G. (1978). Estimating the Dimension of a Model. *The Annals of statistics, 6*(2),
516       461-464. doi:10.1214/aos/1176344136

517  Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & van der Linde, A. (2002). Bayesian measures
518       of model complexity and fit. *Journal of the Royal Statistical Society. Series B,*
519       *Statistical methodology, 64*(4), 583-639. doi:10.1111/1467-9868.00353

520  Stephan, Klaas E., Iglesias, S., Heinzle, J., & Diaconescu, Andreea O. (2015). Translational
521       Perspectives for Computational Neuroimaging. *Neuron, 87*(4), 716-732.
522       doi:http://dx.doi.org/10.1016/j.neuron.2015.07.008

523  Stephan, K. E., Kasper, L., Harrison, L. M., Daunizeau, J., den Ouden, H. E. M., Breakspear,
524       M., & Friston, K. J. (2008). Nonlinear dynamic causal models for fMRI. *NeuroImage,*
525       *42*(2), 649-662. doi:http://dx.doi.org/10.1016/j.neuroimage.2008.04.262

526  Stephan, K. E., Penny, W., Daunizeau, J., Moran, R. J., & Friston, K. J. (2009). Bayesian model
527       selection for group studies. *NeuroImage, 46*(4), 1004-1017.
528       doi:http://dx.doi.org/10.1016/j.neuroimage.2009.03.025

529  Stephan, K. E., Schlagenhauf, F., Huys, Q. J. M., Raman, S., Aponte, E. A., Brodersen, K. H., . . .
530       Heinz, A. (2017). Computational neuroimaging strategies for single patient
531       predictions. *NeuroImage, 145, Part B*, 180-199.
532       doi:https://doi.org/10.1016/j.neuroimage.2016.06.038

533  Stephan, K. E., Weiskopf, N., Drysdale, P. M., Robinson, P. A., & Friston, K. J. (2007).
534       Comparing hemodynamic models with DCM. *NeuroImage, 38*(3), 387-401.
535       doi:http://dx.doi.org/10.1016/j.neuroimage.2007.07.040

536  Swendsen, R. H., & Wang, J.-S. (1986). Replica Monte Carlo Simulation of Spin-Glasses.
537       *Physical Review Letters, 57*(21), 2607-2609. doi:10.1103/physrevlett.57.2607

538  Trujillo-Barreto, N. J., Aubert-Vázquez, E., & Valdés-Sosa, P. A. (2004). Bayesian model
539      averaging in EEG/MEG imaging. *NeuroImage (Orlando, Fla.), 21*(4), 1300-1319.
540      doi:10.1016/j.neuroimage.2003.11.008

541  Vyshemirsky, V., & Girolami, M. A. (2008). Bayesian ranking of biochemical system
542      models. *Bioinformatics, 24*(6), 833-839. doi:10.1093/bioinformatics/btm607

543  Wipf, D., & Nagarajan, S. (2009). A unified Bayesian framework for MEG/EEG source
544      imaging. *NeuroImage, 44*(3), 947-966.
545      doi:https://doi.org/10.1016/j.neuroimage.2008.02.059

546  Wolpert, R. L., & Schmidler, S. C. (2012). α-STABLE LIMIT LAWS FOR HARMONIC MEAN
547      ESTIMATORS OF MARGINAL LIKELIHOODS. *Statistica Sinica, 22*(3), 1233-1251.
548      doi:10.5705/ss.2010.221

549  Yao, Y., Raman, S. S., Schiek, M., Leff, A., Frässle, S., & Stephan, K. E. (2018). Variational
550      Bayesian inversion for hierarchical unsupervised generative embedding (HUGE).
551      *NeuroImage, 179*, 604-619.
552      doi:https://doi.org/10.1016/j.neuroimage.2018.06.073

553