# Dynamic functional brain connectivity results depend on modeling assumptions: comparing frequentist and Bayesian hypothesis tests

Hester Huijsdens[1], Linda Geerligs[1], Max Hinne[1]

[1] Donders Institute, Radboud University, Nijmegen, The Netherlands

## Abstract

Understanding the temporal dynamics of functional brain connectivity is important for addressing various questions in network neuroscience, such as how connectivity affects cognition and changes with disease. A fundamental challenge is to evaluate whether connectivity truly exhibits dynamics, or simply is static. The most common frequentist approach uses sliding-window methods to model functional connectivity over time, but this requires defining appropriate sampling distributions and hyperparameters, such as window length, which imposes specific assumptions on the dynamics. Here, we explore how these assumptions influence the detection of dynamic connectivity, and introduce an alternative approach based on Bayesian hypothesis testing with Wishart processes. This framework encodes assumptions through prior distributions, allowing prior knowledge on the time-dependent structure of connectivity to be incorporated into the model. Moreover, this framework provides evidence for both dynamic and static connectivity, offering additional information. Using simulations, we compare the frequentist and Bayesian approaches and demonstrate how different assumptions affect the detection of dynamic connectivity. Finally, by applying both approaches to an fMRI working-memory task, we find that conclusions at the individual level vary with modeling choices, while group-level results are more robust. Our work highlights the importance of carefully considering modeling assumptions when evaluating dynamic connectivity.

## 1 Introduction

Time-varying functional connectivity has been explored in an increasing number of studies in network neuroscience [1, 2, 3, 4]. Unlike traditional functional connectivity analyses, in which a single estimate of connectivity is estimated from time-series data, time-varying functional connectivity aims to capture how the brain's functional connectivity changes over time. These time-varying estimates have been shown to serve as more sensitive biomarkers than single time-averaged estimates in different domains [3, 5]. For example, functional connectivity dynamics are altered during aging [6], and can improve the classification of patients diagnosed with neurological disorders and healthy controls [7, 8].

A crucial question related to modeling functional connectivity over time is whether observed fluctuations truly reflect *dynamics*, or that the connectivity is *static* and these fluctuations can be attributed to sampling variability, head motion artifacts present in the time series [9], or certain modeling choices [10]. For example, if we use a sliding-window approach to model time-varying functional connectivity [8, 11], the window length would influence our estimates of functional connectivity. Therefore, it is important to assess whether perceived dynamics are *statistically meaningful*.

In the frequentist framework, the first step to hypothesis testing is to generate a null distribution that preserves the important characteristics of the observed data, while removing any meaningful dynamic functional connectivity. Any remaining dynamics in the null distribution are

random and reflect properties of the observed data, such as autocorrelation in the time series itself. Two common methods to generate such a null distribution are *autoregressive randomization*, where an autoregressive model is estimated and subsequently applied to generate new observations [12, 13], and *phase randomization*, where a Fourier transform is applied to the time series, and subsequently the time series' phase components are randomized [11, 10, 14]. Next, a time-varying connectivity method (such as a sliding window) is applied to obtain estimates of functional connectivity over time for both the null distribution and the observed data. Finally, to evaluate whether the observed connectivity is static or dynamic, a test statistic is computed over the connectivity estimates of the observed and null data. These values are then compared. Although sliding-window methods in combination with bootstrapping are a popular approach for dynamic connectivity hypothesis testing, they do not provide a straightforward way to quantify uncertainty associated with the test outcomes.

In contrast, Bayesian approaches explicitly quantify estimation uncertainty, and have the ability to quantify evidence either against or in favor of the null hypothesis (indicating that the functional connectivity is static). Furthermore, the Bayesian framework allows different assumptions about the temporal structure of the dynamics to be incorporated into the hypothesis testing. Similar to the window size in a sliding-window approach, these prior assumptions might influence the connectivity estimates. However, unlike in the frequentist setting, these prior distributions provide an explicit way to incorporate prior knowledge that we may have about dynamic connectivity directly into the model.

In this paper, we compare a Bayesian approach to dynamic functional connectivity testing with the more common frequentist sliding-window approach. The paper is organized as follows. In Section 2, we first describe the sliding-window method and the Bayesian approach using a Wishart process [15, 16]. We then recap the traditional frequentist test for dynamic connectivity, and introduce a Bayesian statistical test for dynamic functional connectivity. In Section 3, we compare the two frameworks using simulation studies and show how, with the Bayesian method, uncertainty quantification can provide more information about dynamic functional connectivity. Moreover, in Section 4, we illustrate the use of both frameworks on a task-based fMRI data from the Human Connectome Project [17]. Finally, in Section 5, we conclude our comparison and discuss future directions of research.

## 2   Methods

Before we proceed with a discussion of the two statistical frameworks and time-varying connectivity methods, we present some notation and terminology that will be used throughout the rest of this paper. We use lower case letters $x$ for scalar values, bold lower case letters $\mathbf{x}$ for vectors, and bold upper case letters $\mathbf{X}$ for matrices. The time series data consists of $n$ observations over $d$ variables (i.e. brain regions), with each observation denoted by $\mathbf{y}_i \in \mathbb{R}^d$. Stacked together, these form the matrix $\mathbf{Y} = (\mathbf{y}_1, \ldots, \mathbf{y}_n)^\top \in \mathbb{R}^{n \times d}$.

### 2.1   Hypothesis testing for dynamic connectivity

Identifying whether one or more pairs of regions are truly dynamically connected, requires the following steps:

1. First, we must estimate the functional connectivity over time itself, as this is not directly observed. This is commonly done using a form of the sliding-window method [10, 18], which we discuss in Section 2.2.1. Alternatively, one can use a Bayesian approach based on Wishart processes [1, 19], which are described in sections 2.2.2 and 2.2.3.

2. To test for meaningful dynamics in connectivity, the connections are typically summarized using test statistics. These are presented in Section 2.3.

3. Lastly, the test statistics are used for hypothesis testing, which can be done in either the frequentist, $p$-value based framework, or in a Bayesian way, which we introduce in this work. The two approaches for testing are described in Section 2.4.

## 2.2 Dynamic connectivity estimation

We are interested in the time-varying connectivity between all pairs of brain regions. Throughout this paper, we use the following notation to refer to (dynamic) connectivity estimates. We use $\mathbf{\Sigma} \in \mathbb{R}^{d \times d}$ to refer to a single (positive-definite) covariance matrix, such as obtained by a sliding-window estimate. In addition, $\mathbf{\Sigma}(x_i) \in \mathbb{R}^{d \times d}$ refers to the covariance specifically at input location $x_i$. Finally, $\mathbf{\Sigma}(\mathbf{x}) \in \mathbb{R}^{n \times d \times d}$ represents a sequence of covariance matrices, such that $\mathbf{\Sigma}(\mathbf{x}) = (\mathbf{\Sigma}(x_1), \ldots, \mathbf{\Sigma}(x_n))$.

### 2.2.1 Sliding-window methods

The intuition behind the sliding-window approach is simple: rather than computing a single covariance matrix of the observations $\mathbf{Y}$ using $\mathbf{\Sigma} = \frac{1}{n-1}\mathbf{Y}^\top\mathbf{Y}$, we compute $\mathbf{\Sigma}$ for a series of consecutive segments of $\mathbf{Y}$ instead [8, 11, 10]. The length of this segment, $\lambda$, is known as the window size, and determines how many observations are used for each covariance estimate. To get an estimate of covariance over time, the window is shifted by a number of observations, called the stride length $\tau$.

The method is described in detail in Supplementary Section S1.1. Although this approach is easy to implement, it has an important limitation. Namely, both the size of the window and the stride length must be determined by the practitioner, while it has been shown that this parameter has a substantial effect on the estimated connectivity [20, 21, 22]. A small window length will result in rapidly fluctuating estimates of connectivity, whereas with a large window size, relevant fluctuations in the connectivity might be smoothed out.

### 2.2.2 Wishart processes

We propose to combine the Bayesian framework for hypothesis testing with Wishart processes [15, 16], a Bayesian model that was introduced by Wilson and Ghahramani [16] and has later been used to model dynamic covariance in different domains such as neuroscience and psychology [1, 19, 23]. Wishart processes allow us to set prior assumptions on different aspects of the connectivity, such as on the type of connectivity structure over time.

To explain the Wishart process, it is helpful to first look at how to estimate a single $d \times d$ covariance matrix $\mathbf{\Sigma}$, as in the static connectivity case, using a Bayesian model. A straightforward Bayesian model for $\mathbf{\Sigma}$ is:

$$
\begin{aligned}
\mathbf{\Sigma} &\sim \mathcal{W}(\mathbf{V}, v) \\
\mathbf{y}_i &\sim \mathcal{MVN}_d(\mathbf{0}, \mathbf{\Sigma}) \ , \quad i = 1, \ldots, n \ ,
\end{aligned}
\tag{1}
$$

in which $\mathcal{MVN}_d(\mathbf{0}, \mathbf{\Sigma})$ is the $d$-dimensional multivariate Gaussian distribution with mean $\mathbf{0}$ and covariance matrix $\mathbf{\Sigma}$, and $\mathcal{W}(\mathbf{V}, v)$ is the Wishart *distribution* with scale matrix $\mathbf{V}$ and $v$ degrees of freedom. As the Wishart distribution is conjugate to the multivariate Gaussian, the posterior distribution $p(\mathbf{\Sigma} \mid \mathbf{Y})$ can be computed analytically [24].

Another way of looking at the Wishart distribution is by how samples from this distribution are generated. Using again the multivariate Gaussian distribution, a Wishart-distributed covariance matrix can be obtained using:

$$
\begin{aligned}
\mathbf{f}_k &\sim \mathcal{MVN}_d(\mathbf{0}, \mathbf{I}) \ , \quad k = 1, \ldots, v \\
\mathbf{\Sigma} &= \sum_{k=1}^{v} \mathbf{L}\mathbf{f}_k\mathbf{f}_k^\top\mathbf{L}^\top \sim \mathcal{W}(\mathbf{V}, v) \ ,
\end{aligned}
\tag{2}
$$

in which $\mathbf{V} = \mathbf{L}\mathbf{L}^\top$ is the Cholesky decomposition of the scale matrix $\mathbf{V}$. In words, the sum of outer products of $v$ i.i.d. zero-mean multivariate Gaussian samples is Wishart distributed [25].

Extending this idea, the Wishart *process* instead models a covariance matrix that changes with input. The Wishart process is constructed similarly to Eq. (2), but the multivariate Gaussian variates are now replaced by independent zero-mean Gaussian processes (GPs) [26]:

$$f_{kj} \sim \mathcal{GP}\left(\mu(\cdot), \kappa(\cdot, \cdot; \theta)\right) \ , \quad k = 1, \ldots, v \ , \quad j = 1, \ldots, d \ , \tag{3}$$

where $\mu(\cdot)$ is the mean function and $\kappa(\cdot, \cdot; \theta)$ the covariance function with parameters $\theta$. We collect the evaluations of these functions as $\mathbf{f}_k(x_i) = \left(f_{k1}(x_i), \ldots, f_{kd}(x_i)\right)^\top$. The covariance function determines several properties, such as smoothness, of the Wishart process. Moreover, if we assume that $\kappa(x_i, \ldots, x_i; \theta) = 1$ for all observations $i = 1, \ldots, n$, we can separate the correlations from the scale matrix $\mathbf{V}$ when constructing the Wishart process:

$$\mathbf{\Sigma}(x_i) = \sum_{k=1}^{v} \mathbf{L}\mathbf{f}_k(x_i)\mathbf{f}_k(x_i)^\top \mathbf{L}^\top \sim \mathcal{W}(\mathbf{V}, v) \ . \tag{4}$$

Analogous with Eq. 1, this is then combined with the likelihood

$$\mathbf{y}_i \mid x_i \sim \mathcal{MVN}_d\left(\mathbf{0}, \mathbf{\Sigma}(x_i)\right) \ , \quad i = 1, \ldots, n \ . \tag{5}$$

Since the Wishart process is based on GPs, we can explicitly place model assumptions on temporal autocorrelations in the dynamic covariance. For example, if we already know that the connectivity is likely to change periodically over time, a periodic function can be used as $\kappa(\cdot, \cdot; \theta)$. We return to the consequences of the different assumptions in Section 3, and provide more details in Supplementary Section S1.2.

### 2.2.3 Hierarchical Wishart processes

The standard Wishart process is commonly applied to model dynamic covariance based on a single multivariate time series. For example, in the case of modeling functional connectivity, the Wishart process would be applied with the fMRI time series from a single participant. However, in some scenarios, such as in the case of task-based fMRI analyses, one might be interested in an estimate of dynamic covariance at the group level, across multiple participants. The input-dependent estimate of functional connectivity $\mathbf{\Sigma}(\mathbf{x})$ is then inferred from multiple multivariate time series (in our case: from multiple subjects).

To achieve this, we extend the Wishart process to a hierarchical model. We expand Eq. (15) in Supplementary Section S1.2 by assuming that the observations of all individual subjects share a common covariance matrix:

$$\mathbf{y}_{im} \mid x_i \sim \mathcal{MVN}_d\left(\mathbf{0}, \mathbf{\Sigma}(x_i)\right) \quad i = 1, \ldots, n, \quad m = 1, \ldots, s \ , \tag{6}$$

where $s$ refers to the number of subjects. The observations of the different subjects are independent given the covariance. We refer to this model as the *hierarchical Wishart process*.

## 2.3 Test statistics for dynamic connectivity

To determine whether a particular connection is indeed dynamic, we must first summarize the time series of connectivity using test statistics. It is important to critically assess which statistic is being used, because this formalizes which properties of the signal are considered important and encode another layer of dynamic functional connectivity assumptions [3]. We focus on three commonly used test statistics: the variance of connectivity [8, 27], the maximum power across all frequencies [14], and the median-crossings [28].

### 2.3.1 Variance of connectivity

The first statistic, denoted by $\eta\left(\mathbf{\Sigma}\left(\mathbf{x}\right)\right)$, is the variance of the connectivity estimates:

$$\eta\left(\mathbf{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = \frac{1}{n-1}\sum_{i=1}^{n}\left(\mathbf{\Sigma}(x_i)_{jk} - \mu_{jk}\right)^2 \ , \quad j = 1,\ldots d \ , \quad k = 1,\ldots, d \ , \tag{7}$$

where $\mathbf{\Sigma}(\mathbf{x})_{jk}$ is the connectivity between regions $j$ and $k$ over time and $\mu_{jk} = \frac{1}{n}\sum_{i=1}^{n}\mathbf{\Sigma}(x_i)_{jk}$ is the mean of this connection. The variance increases when connections show more fluctuations over time.

### 2.3.2 Maximum power across all frequencies

The maximum power across all frequencies [14], for a connection between regions $j$ and $k$, written as $\psi\left(\mathbf{\Sigma}\left(\mathbf{x}\right)\right)$, is computed by first determining the power of all frequencies $\tilde{\mathbf{\Sigma}}_{qjk}$ using the Discrete Fourier Transform, and then taking the largest power across all frequencies $q = 1,\ldots, n$:

$$\psi\left(\mathbf{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = \max_q |\tilde{\mathbf{\Sigma}}_{qjk}|^2 \ . \tag{8}$$

The idea is that noise in the data is likely to contain low powers, whereas a high power is more likely to indicate a strong signal. Hence, large powers in connectivity estimates might imply dynamics.
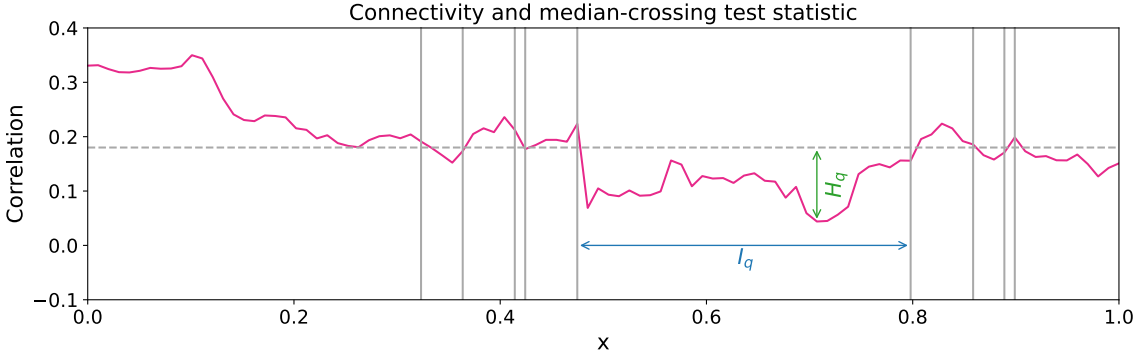
### 2.3.3 Median-crossings



Figure 1: **An illustration of the median-crossing test statistic.** The median is shown by the dashed line, and the segments are shown by the gray vertical lines. The median-crossing test statistic is computed as a non-linear combination of the heights (in green) and lengths (in blue) of the median crossings.

The third test statistic evaluates the duration and magnitude of deviations from the median of the estimated connectivity [28] and is visualized in Figure 1. First, we obtain the median of the connection $m_{jk} = \mathrm{med}(\mathbf{\Sigma}(\mathbf{x}))_{jk}$. Second, the connectivity estimates are divided into non-overlapping segments based on the locations where the connectivity crosses the median. Let $c_q$ be the point where the connectivity passes the median line for the $q$-th time. Then, the length of such an excursion from the median is defined as $I_q = c_{q+1} - c_q$, and the height as the maximum distance from the median within the window, namely $H_q = \max\left(|\mathbf{\Sigma}(x_i)_{jk} - m_{jk}|\right)$, with $c_q \leq x_i \leq c_{q+1}$. The test statistic $\zeta\left(\mathbf{\Sigma}\left(\mathbf{x}\right)\right)$ is defined as a weighted sum of the heights and lengths of all segments:

$$\zeta\left(\mathbf{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = \sum_{q=1}^{Q-1} |I_q^\gamma H_q^\beta| \ , \quad j = 1,\ldots d \ , \quad k = 1,\ldots, d \ , \tag{9}$$

where $\gamma$ and $\beta$ are the relative weighting of the lengths and heights. This test statistic captures both the frequency of the signal (as with a higher frequency, the median is crossed more often), as well as the amplitude (as higher excursions imply a larger amplitude).

## 2.4 Hypothesis testing for dynamics

Here, we discuss the two different frameworks for testing for dynamic functional connectivity. The most commonly used framework for hypothesis testing is the frequentist framework, and we will compare this test against a Bayesian model comparison framework.

### 2.4.1 Frequentist hypothesis testing



**Figure 2: An overview of the frequentist framework to statistically test for dynamic functional connectivity.** We use the observed time series (A) to generate a null distribution (B) using phase randomization. From both the observed time series and phase randomized time series, we estimate connectivity over time (C). Finally, we compute a test statistic over both, and test for significant dynamics (D).

Within the frequentist framework, we test for significant dynamics of the estimated time-varying connectivity by comparing it against a null distribution. The general premise is that, if the observed dynamics are unlikely given the null distribution, we conclude that the dynamics are truly present [29, 30]. The approach consist of three main steps. First, we simulate time series data to create a null distribution and estimate its connectivity. Second, we compute a test statistic for both the actual estimate and estimates in the null distribution, and finally we compare the test statistic that corresponds to the observed signal against those of the null distribution by computing a $p$-value. This $p$-value is used to determine if the observed dynamics are significant.

A widely used approach for generating a null distribution of time series data is phase randomization [10, 14], which preserves as many properties of the original time series as possible, except for the dynamics. Phase randomization works as follows. First, the observed time series are transformed to the frequency domain via the discrete Fourier transform. Then, a uniformly sampled random phase $\theta \in [0, 2\pi]$ is added to each frequency, and finally the frequencies are transformed back into the time domain using the inverse discrete Fourier transform. Importantly, the same randomly sampled phase $\theta$ should be used for all brain regions, otherwise we also remove any static functional connectivity in between the time series. We combine the frequentist testing approach with the sliding-window estimation procedure. Figure 2 provides a schematic overview of frequentist hypothesis testing for dynamic connectivity.

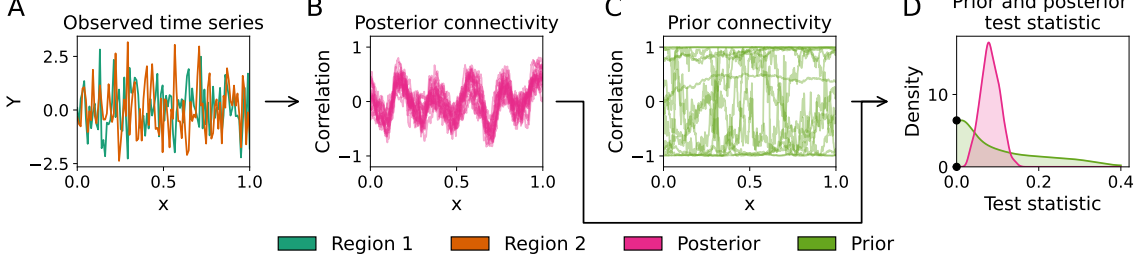### 2.4.2 Bayesian hypothesis testing



**Figure 3: An overview of the Bayesian framework to estimate functional connectivity and statistically test for dynamics.** Using the observed time series (A), the posterior distribution of connectivity is inferred (B). We compare the posterior with the prior distribution (C) to compute the Bayes factor as the fraction of the prior and posterior probability of the connectivity being static (see Eq. (12)). These probabilities are indicated by the black dots in panel D.

The Bayesian estimation and testing procedure is illustrated in Figure 3. In a Bayesian estimation approach, prior beliefs about the connectivity are defined, and these beliefs are updated based on the observed time series via Bayes' theorem:

$$p\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right) \mid \mathbf{x}, \mathbf{Y}\right) = \frac{p\left(\mathbf{Y} \mid \boldsymbol{\Sigma}\left(\mathbf{x}\right)\right) p\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)\right)}{p\left(\mathbf{Y}\right)} \quad , \tag{10}$$

where $p\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right) \mid \mathbf{x}, \mathbf{Y}\right)$ is the posterior distribution of the connectivity given the observed data, $p\left(\mathbf{Y} \mid \boldsymbol{\Sigma}\left(\mathbf{x}\right)\right)$ the likelihood of the observed data given the connectivity, and $p\left(\mathbf{Y}\right)$ is the marginal likelihood of the data, which serves as a normalizing constant. Importantly, $p\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)\right)$ is the prior distribution, which encodes any beliefs about the dynamic connectivity before observing any data. For example, if we expect the connectivity to smoothly vary over time, contain fast fluctuations, or have a periodic structure, we could incorporate these assumptions here. The resulting posterior estimate of connectivity is a combination of the prior assumptions and the information contained in the observed time series data.

To test for significant dynamics in the functional connectivity, we perform Bayesian model comparison by computing Bayes factors [31], which consist of the ratio of the marginal likelihoods of two models. For example, if model $\mathcal{M}_0$ assumes there is no dynamic connectivity, while model $\mathcal{M}_1$ assumes there is, then the Bayes factor is given by

$$\mathrm{BF}_{10} = \frac{p(\mathbf{Y} \mid \mathcal{M}_1)}{p(\mathbf{Y} \mid \mathcal{M}_0)} \quad , \tag{11}$$

indicating how much more likely the data are under $\mathcal{M}_1$ compared to $\mathcal{M}_0$. Computing the marginal likelihood of a complicated Bayesian model like the Wishart process is intractable. However, in the case of nested hypotheses, the marginal likelihood may be computed using the Savage-Dickey method [32, 33], which compares the prior and posterior probabilities at the point indicated by the null hypothesis. For instance, we can use the variance statistic from Eq. (7) to test if the variance is zero (indicating no dynamics) by computing:

$$\mathrm{BF}_{10} = \frac{p\left(\eta\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = 0\right)}{p\left(\eta\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = 0 \mid \mathbf{x}, \mathbf{Y}\right)} \quad , \quad j = 1, \ldots d \ , \quad k = 1, \ldots, d \ . \tag{12}$$

Here, the value $\eta\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = 0$ is implicitly given by the null hypothesis: if there is no dynamic connectivity, this implies a variance of zero.

Although the Bayes factor can be thresholded to obtain a binary test outcome, it also offers valuable information about the uncertainty associated with these test outcomes. Namely, the Bayes factor quantifies the strength of evidence in favor of either the null or alternative hypothesis. To make this more clear, Eq. (12) shows the ratio of updated beliefs after seeing the data at the null hypothesis point. If, after conditioning on the observed $\mathbf{Y}$, our belief that $\eta\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = 0$ has increased, then we observe a Bayes factor $\mathrm{BF}_{10} < 1$, indicating evidence in favor of the null hypothesis. On the other hand, if our beliefs about the possibility $\eta\left(\boldsymbol{\Sigma}\left(\mathbf{x}\right)_{jk}\right) = 0$ have decreased, then the data suggest the null hypothesis has become less likely, and the alternative model has gained evidence instead, with $\mathrm{BF}_{10} > 1$. Here, the Bayesian testing approach is combined with the Wishart process estimation procedure.

## 2.5   Simulations

Since a ground truth is usually unknown in empirical time series, we first compare the reliability of the frequentist sliding-window and Bayesian Wishart process hypothesis testing approaches on simulated time series with different characteristics. Each simulation consists of $d = 2$ brain regions and $n \in \{150, 300, 600\}$ observations. We let the time range $\mathbf{x}$ be uniformly spaced from 0 to 1. To simulate time series data, we first generate correlation matrices of which the off-diagonal element is either static or dynamic, and the diagonal elements are set to 1:

$$\boldsymbol{\Sigma}(x_i) = \begin{pmatrix} 1 & \sigma(x_i) \\ \sigma(x_i) & 1 \end{pmatrix} \ . \tag{13}$$

In simulations were the off-diagonal element is dynamic, the signal follows a sine wave:

$$\sigma(x_i) = A\sin\left(\omega 2\pi x_i + \varphi\right) \ , \quad i = 1, \ldots n \ . \tag{14}$$

We evaluated the performance of both frameworks for strong and weak functional connectivity by varying the signal amplitude $A \in \{0.2, 0.4, 0.6, 0.8\}$. Moreover, varying frequencies $\omega \in \{1, 2, \ldots, 5\}$ were compared. The phase was randomly set to a value within the range $[0, 2\pi]$ and static connections were set to a constant value within the range $[-0.4, 0.4]$. The time series were generated by sampling from a multivariate distribution using:

$$\mathbf{y}_i \mid x_i \sim \mathcal{MVN}_d\left(\mathbf{0}, \boldsymbol{\Sigma}\left(x_i\right)\right) \ , \quad i = 1, \ldots, n \ . \tag{15}$$

In addition to the periodic simulations, we generated data with non-periodic connectivity that follows a state-switching pattern, based on the work by Thompson et al. [34]. We let the covariance matrix in Eq. (13) switch between two states with the strength of connectivity $\sigma_i$ being 0.1 in one state and 0.8 in the other state. The duration of each state is randomly sampled from $[20, 30, 40, 50, 60]$ time points. After the sampled number of time points, the simulation switches to the other state.

To reliably compare the two frameworks on different types of connectivity, we generate multiple connectivity matrices over time for each combination of simulation parameters. For every set of simulation parameters, that is for every connectivity pattern and number of observations, and in the case of periodic connectivity also for every amplitude and frequency, we simulate 40 random connectivity patterns. Using Eq.(15), we then generate one random dataset for each of these simulated connectivity matrices over time. Implementations of these simulations and the corresponding hypothesis tests can be found on our GitHub page.

## 2.6   Participants, fMRI data, and parcellation

To explore how the two frameworks can be used in practice to test for dynamics, we test for dynamics in a working memory task dataset from the Human Connectome Project (HCP) [17]. For computational reasons, we used the 100 unrelated subjects subset of the HCP database. From

these 100 subjects, we only selected participants with a task accuracy above 60%, resulting in fMRI data from 95 healthy participants.

The task paradigm is described in Supplementary Section S2. In the n-back task, the 0-back blocks are considered to require a lower memory load than the 2-back blocks. Previous studies have found that increase in working memory load is associated with an increase in functional connectivity within the frontoparietal network, and a decreased connectivity in the default mode network [35].

The fMRI data was parcellated using Glasser's MMP atlas [36]. To select regions of interest, we used the sliding-window approach with a window size of 10% of the number of observations to compute the correlation between each region pair's functional connectivity and the 2-back task paradigm. We then selected the pair of regions with the strongest correlation, namely the dorsolateral prefrontal cortex (DLPFC) and the inferior parietal lobule (IPL). Additionally, we selected another region that was unrelated to the task paradigm, namely the primary auditory cortex (A1). We expect that especially the connectivity between the DLPFC and IPL regions is involved in working memory and hence changes dynamically depending on the task condition [37, 38]. The connectivity between the DLPFC and A1 regions and between the IPL and A1 regions are expected to change less throughout the task paradigm.

## 2.7 Assumptions on connectivity

Both the sliding-window method and Wishart process require the user to define certain modeling choices.

### 2.7.1 Sliding-window assumptions

In the sliding-window method, we set the window length to a fraction of the number of observations to ensure the window captures a similar part of the latent correlations. We compare three different lengths, namely 5%, 10% and 20% of the number of observations in the simulation. In the HCP working memory dataset, this comes down to 20, 40 and 80 time points, which correspond to 14.4, 28.8 and 57.6 seconds of recording. The stride length $\tau$ was set to 1 TR (0.72 seconds), giving us an estimate of functional connectivity at every time point. By setting the stride length to 1, windows will overlap substantially, and therefore we encode the assumption that the functional connectivity varies smoothly over time. Additionally, short window sizes assume that the connectivity will change rapidly, at the expense of signal-to-noise in each window and hence more noisy estimates. Longer window sizes will result in more stable estimates of connectivity, but might miss meaningful fluctuations. Finally, for the frequentist testing framework, we generated 1000 surrogate datasets to construct the null distribution. For the median-crossing statistic, we set the respective weights of the excursion lengths and heights to $\gamma = 0.9$ and $\beta = 1$, following previous work [10, 28].

### 2.7.2 Wishart process assumptions

For the Wishart process, a GP kernel, and priors on the kernel hyperparameters $\theta$ and the scale matrix $\mathbf{V}$ need to be defined.
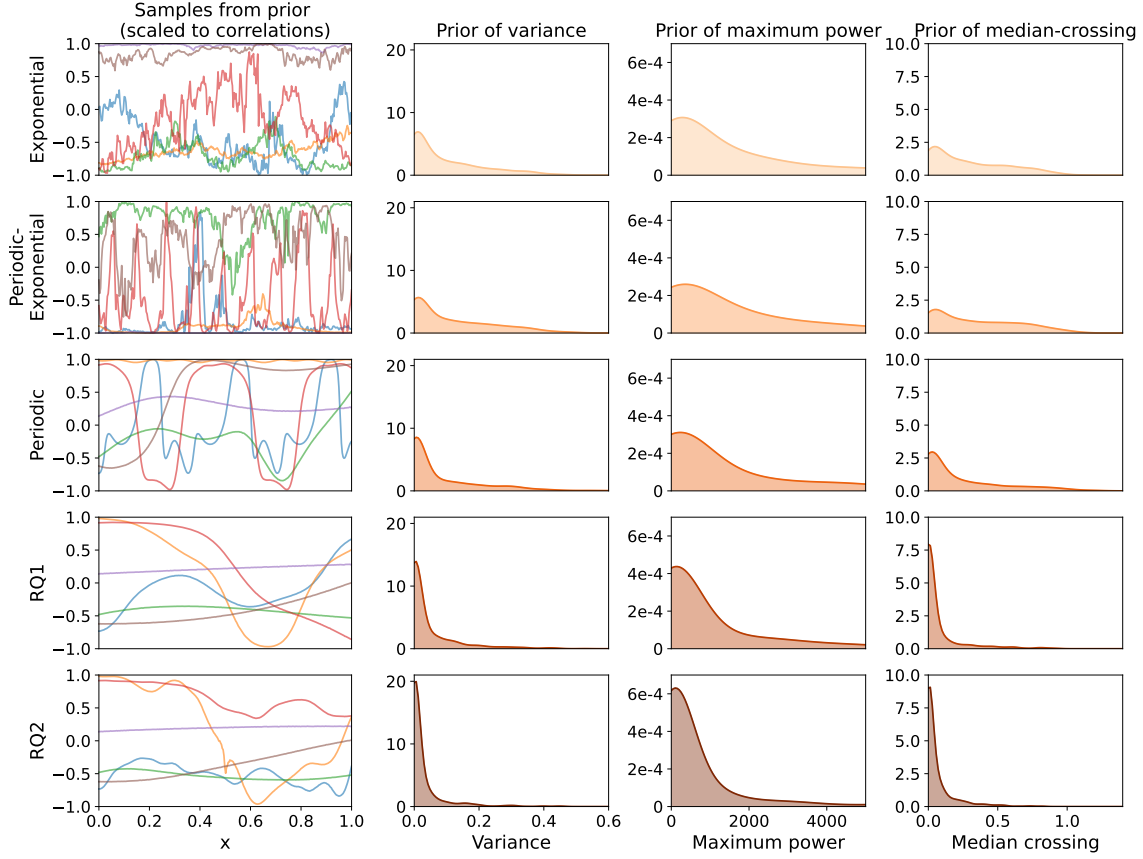
**Figure 4: Samples from the Wishart process prior and their corresponding prior test statistic distributions.** Five different kernel and parameter initializations are shown. All kernel hyperparameter priors were set to a log normal with mean of 0 and standard deviation of 1, except for the bottom row, where the prior for alpha was set to $\log \mathcal{N}(-3, 1)$. The prior on the individual elements of the scale matrix was set to be a normal distribution with mean of 0 and standard deviation of 1.

To explore the effect that different kernels can have, we use four different GP kernels in our experiments, namely an exponential, periodic-exponential (often referred to as a locally periodic), periodic, and rational quadratic kernel:

$$\kappa_{\text{exponential}}(x_i, x_j; \ell) = \exp\left(-\frac{|x_i - x_j|}{2\ell^2}\right) \quad ,$$

$$\kappa_{\text{periodic-exponential}}(x_i, x_j; p, \ell_{\text{periodic}}, \ell) = \exp\left(-\frac{2\sin^2\left(\pi |x_i - x_j| / p\right)}{\ell_{\text{periodic}}^2}\right) \exp\left(-\frac{|x_i - x_j|}{2\ell^2}\right) \quad ,$$

$$\kappa_{\text{periodic}}(x_i, x_j; p, \ell_{\text{periodic}}) = \exp\left(-\frac{2\sin^2\left(\pi |x_i - x_j| / p\right)}{\ell_{\text{periodic}}^2}\right) \quad ,$$

$$\kappa_{\text{RQ}}(x_i, x_j; \alpha, \ell) = \left(1 + \frac{|x_i - x_j|}{2\alpha\ell^2}\right)^{-\alpha} \quad .$$

$$(16)$$

The periodic-exponential kernel is formed by multiplying an exponential and a periodic function. The different GP covariance functions and kernel parameters place assumptions on the properties of the latent GP samples, and therefore also on the characteristics of the functional connectivity being estimated. For example, a periodic kernel assumes a repeating pattern in the connectivity

structure over time, whereas an exponential kernel assumes functional connectivity with non-smooth fluctuations over time. These assumptions are shown in more detail in Figure 4, where samples from the prior of the Wishart process are shown for different kernels and kernel hyperparameters, together with their corresponding a priori test statistic distributions. As illustrated in the figure, connectivity samples from the exponential and periodic-exponential kernels contain smaller fluctuations and are less smooth than the periodic and rational quadratic kernels. This is also influenced by the kernel hyperparameters. The exponential, periodic-exponential and rational quadratic kernels have a lengthscale parameter $\ell$, and both the periodic and periodic-exponential kernels have parameters $p$ and $\ell_{\mathrm{periodic}}$, which represent the period and lengthscale within a period. The rational quadratic kernel has an additional $\alpha$ parameter that determines the weighting of small and large variations. If $\alpha$ increases, GP samples of this kernel will become increasingly smooth. In the prior, all hyperparameters were assumed to follow a log normal prior with mean of 0 and standard deviation of 1. Additionally, we explore the effects of different priors for the kernel parameters by setting a log-normal with mean of -3 and standard deviation of 1 on the value of $\alpha$, which should result in samples that are less smooth. We set a normal prior with mean of 0 and standard deviation of 1 on each element of the lower Cholesky decomposition ($\mathbf{L}$) of the scale matrix $\mathbf{V}$.

# 3 Results of simulation studies

First, we discuss the ability of the sliding-window method and Wishart process to recover functional connectivity from observed data. We then focus on how effective each framework can detect dynamic functional connectivity, comparing the strengths and limitations of the frequentist and Bayesian approaches in different simulations. Additionally, we explore how uncertainty quantification, provided by the Bayesian framework, can provide additional insights into dynamics. In this section, we focus on the periodic simulations. All results for the state-switching simulations can be found in Supplementary Section S3.

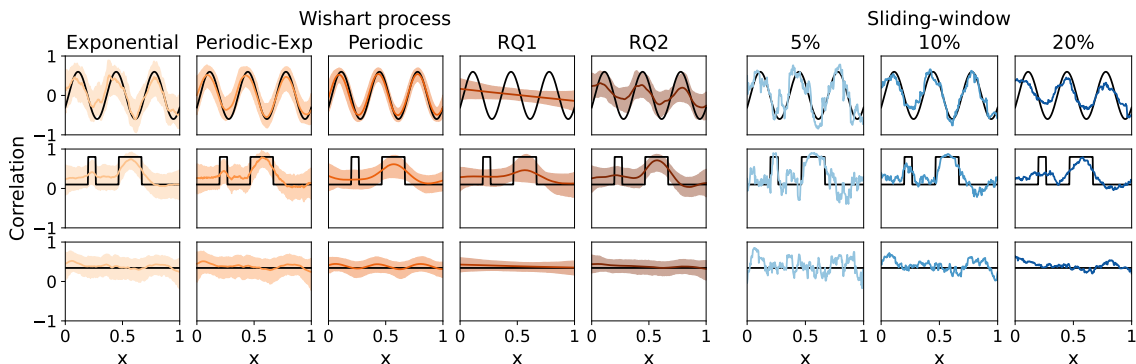## 3.1 Accuracy of connectivity estimates and influence of prior assumptions



**Figure 5: Examples of connectivity estimates based on simulated data with $n = 300$ observations.** We show estimates based on all three types of simulations, namely covariances with a periodic, state-switches and static structure. We show the estimates using different kernels and window lengths. Here RQ1 and RQ2 refer to a rational quadratic kernel with a $\log \mathcal{N}(0, 1)$ prior on all hyperparameters, and a $\log \mathcal{N}(-3, 1)$ prior on the alpha parameter, respectively. For the Wishart process, we show the 95% highest density interval. The ground truth correlations are shown in black.

Figure 5 presents a few examples of connectivity estimates by the Wishart process and sliding-window method using different kernels and window lengths. All estimates are based on $n = 300$ observations. In the top row, the latent correlation follows a sine wave with a frequency of 3 and amplitude of 0.6. These results indicate that both methods can accurately capture the latent connectivity, but that the quality of these estimates is affected by the choice of window length, kernel, and priors on the kernel hyperparameters. For example, when using a rational quadratic kernel with all kernel hyperparameters sampled from $\log \mathcal{N}(0, 1)$, as denoted by RQ1 in the figure, the Wishart process is unable to capture the fast fluctuations in the connectivity. However, when a $\log \mathcal{N}(-3, 1)$ prior is set on the $\alpha$ parameter (see Eq. (16)), as indicated by RQ2, the estimates improve. Estimates become even more accurate when assuming periodicity, or when using the exponential kernel. The sliding-window estimate, with a window length of 20% of the observations, seems too smooth to capture the rapid fluctuations in the connectivity.

The middle row illustrates a latent connectivity with non-periodic state switches. In this case, the sliding-window method captures the rapid switches in connectivity more accurately than the Wishart process. Nonetheless, the performance differs based on the three window lengths, as a window length of 20% is too smooth to model the rapid switches, whereas a window length of 5% estimates a lot of fluctuations in the static periods of the connectivity. Finally, the bottom row presents the static connectivity scenario. Here we again observe that small window lengths estimate a lot of fluctuations, whereas larger window lengths provide more accurate estimates of connectivity. The Wishart process captures the static connection well, as the zero line is fully encapsulated by the 95% highest density areas.
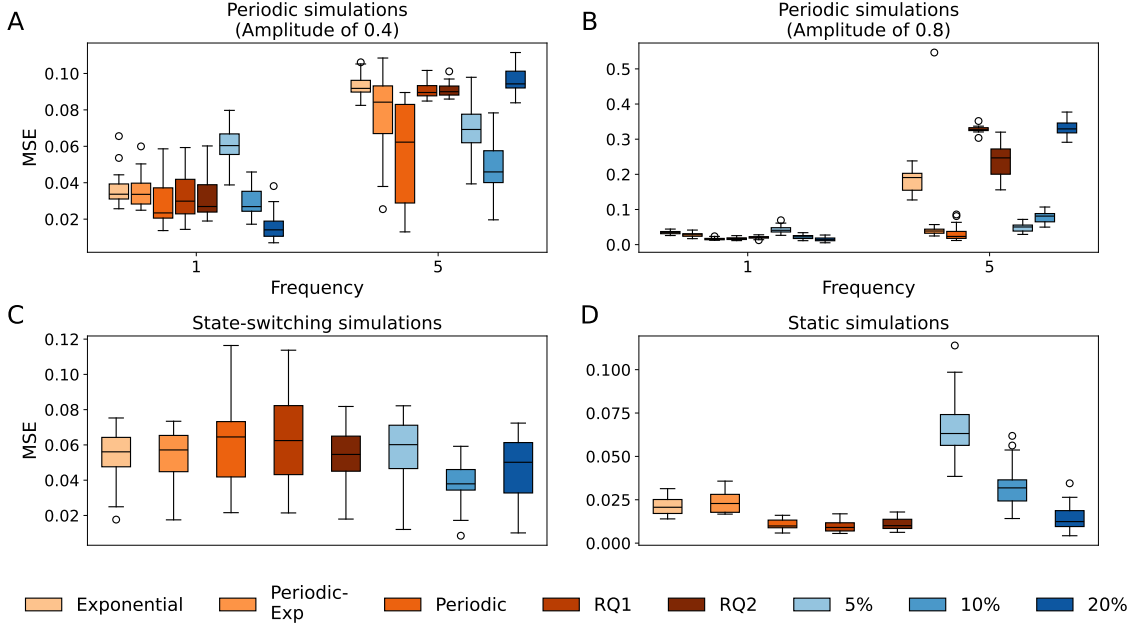


**Figure 6: Mean squared errors (MSE) between true and estimated correlations for with 300 observations.** Subplot A and B show the periodic simulations with an amplitude of 0.4 and 0.8, whereas subplot C and D show the state-switching and static simulations. For the Wishart process, the mean squared errors were computed over each sample of the posterior separately.

To measure the accuracy of the estimated connectivity compared to the ground truth correlations, we compute the mean squared error (MSE). However, to accurately reflect the differences between the Bayesian and frequentist frameworks, the MSEs are computed differently for the Wishart process and sliding-window method. In the Bayesian framework we make use of the full posterior distribution $p(\mathbf{\Sigma}(\mathbf{x}) \mid \mathbf{Y})$ when testing for dynamics. Therefore, the MSE is computed

for each posterior sample and then averaged. This approach naturally results in larger MSEs, because the full posterior variability is incorporated into the MSE. In contrast, the frequentist sliding-window method provides only single values instead of a full distribution, and therefore we compute the MSE directly on these. However, for comparison purposes, we provide the MSEs using only the posterior mean estimate in Figure S4.

Figure 6 shows the MSEs for $n = 300$ observations for the periodic simulations with amplitudes of 0.4 and 0.8 and frequencies of 1 and 5, and for the state-switching and static simulations. Results for other numbers of observations, amplitudes and frequencies can be found in Supplementary Section S3.1. Overall, the sliding-window method with a large window size (10% and 20% of the observations) and the Wishart process with a smooth kernel tend to recover static connections more accurately than dynamic connections. Moreover, we observe that both the Wishart process with the exponential and rational quadratic kernels fail to capture correlations with higher frequencies. The sliding-window method with the largest window length also shows difficulties in modeling these faster dynamics. Finally, if the amplitude of the connectivity is decreased, we observe that all methods show increased difficulty in accurately modeling the connectivity with a lower amplitude, and there is no clear effect of frequency anymore.

Figure 6C presents the results from the state-switching simulations. Here the exponential, periodic-exponential, and rational quadratic with a $\log \mathcal{N}(-3, 1)$ prior set on $\alpha$ outperform the other kernels. In contrast, the rational quadratic kernel with $\log \mathcal{N}(0, 1)$ priors on all hyperparameters (RQ1) and the periodic kernel struggle with the rapid state switches, as could also be seen in the example in Figure 5. The sliding-window method with a window length of 10% outperforms the Wishart process estimates in terms of MSEs. However, the performance is highly dependent on window length.

## 3.2   Hypothesis test performances

To compare the frequentist and Bayesian hypothesis testing frameworks in detecting dynamic connectivity, we first convert the Bayes factors obtained from Eq. (12) into binary decisions. According to general practices [31, 39, 40], log Bayes factors greater than 3 are considered conclusive evidence for the alternative hypothesis, whereas values below 1/3 indicate strong evidence for the null hypothesis. Based on these guidelines, we classify a connection as dynamic if its log Bayes factor is above 3. We then evaluate both frameworks based on three different metrics, namely accuracy, recall and false positive rate. Accuracy provides a measure of the overall number of correctly classified connections, taking into account both static and dynamic connections. Recall quantifies the number of connections that is correctly classified as being dynamic, therefore it measures the sensitivity of the two hypothesis testing frameworks. The false positive rate indicates the number of connections incorrectly classified as being dynamic. In the frequentist framework, the false positive rate is directly related to the significance threshold. Since we classify a p-value below 0.05 as significantly dynamic, the false positive rate of this framework will approximate 5%.
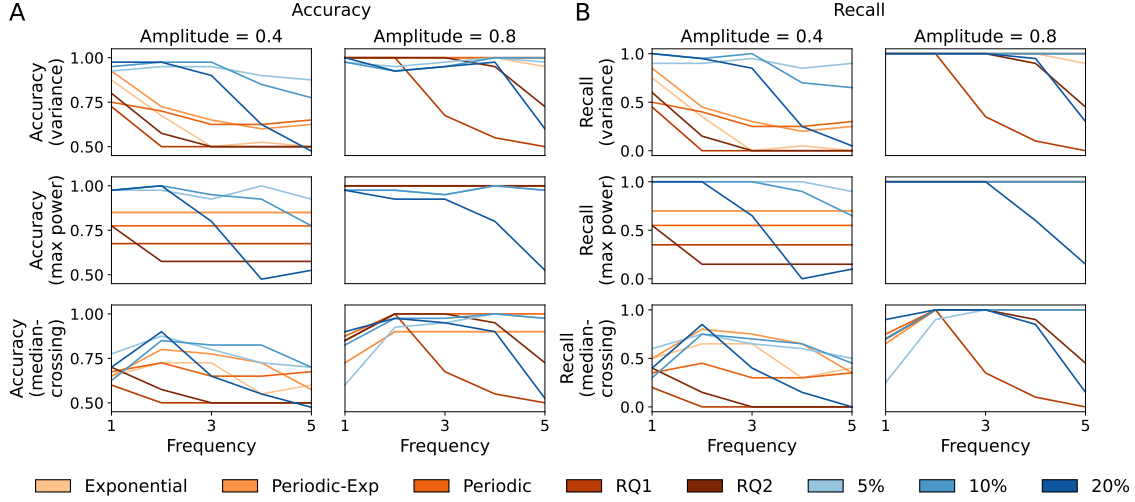
**Figure 7: Hypothesis test performances on the periodic simulations with 300 observations and amplitudes of 0.4 and 0.8.** The figure shows the fractions of correctly classified connections, out of all connections (accuracy, A) and out of only dynamic connections (recall, B) when using the variance, maximum power and median-crossing test statistics.

Figure 7 shows the accuracy (Figure 7A) and recall (Figure 7B) for the simulations with a periodic covariance structure, with $n = 300$ observations and amplitudes of 0.4 and 0.8. All results for amplitudes of 0.2 and 0.6 are shown in Supplementary Section S5. As expected, recall improves with more observations and higher amplitudes. Generally, the frequentist approach tends to outperform the Bayesian approaches in these simulations, particularly for recall and to a lesser degree in accuracy. However, modeling choices also have a large effect on performance. The periodic and periodic-exponential kernels, as well as smaller window sizes of 5% and 10%, perform well, even for faster dynamics. However, the rational quadratic and exponential kernels, and the largest window size (20%), perform well on low-frequency dynamics, but not always on higher frequencies. Finally, a difference in performance can be observed between the two rational quadratic kernels. The RQ2 kernel, where a prior of $\mathcal{N}(-3, 1)$ was set on $\alpha$, outperforms the RQ1 kernel across all periodic simulations. This difference is expected, because, as can be seen from the prior samples in Figure 4, the RQ2 kernel with a prior of $\mathcal{N}(-3, 1)$ can better model fast fluctuations than the RQ1 kernel.

The median-crossing and variance statistics show similar performance patterns between kernels and window sizes. However, the median-crossing statistic does not accurately detect slow dynamics, as can be observed by the relatively low recall scores at a frequency of 1. Both the variance and maximum power accurately detect dynamics of different frequencies for $n = 600$ observations, and $n = 300$ observations with an amplitude of 0.6 or 0.8, but performance decreases with $n = 300$ observations and an amplitude of 0.2 or 0.4. An interesting exception to this is the exponential kernel, which, with $n = 600$ observations, can accurately distinguish dynamics with an amplitude of 0.2 from static connectivity. Corresponding posterior variance distributions of this simulation are shown in Figure S10.

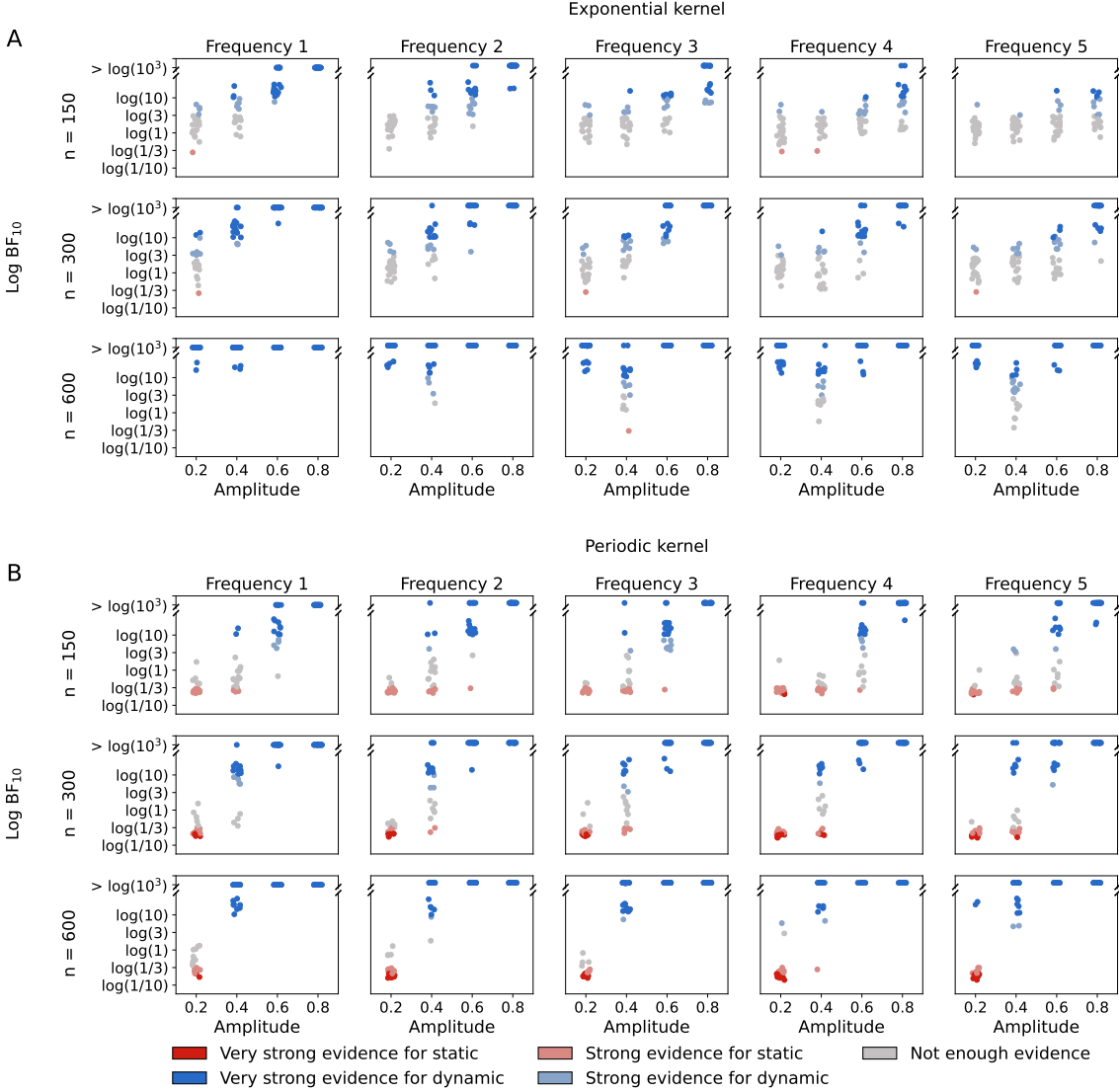## 3.3  Uncertainty in Bayesian hypothesis testing



**Figure 8: Log Bayes factors of detecting dynamics for the periodic simulations.** The figure shows the results for the exponential (A) and periodic (B) kernels and are based on the variance test statistic. Every dot represents a single connection and is colored based on the amount of evidence for the connection being dynamic or static.

Figure 8 presents the distribution of log Bayes factors for the simulations with a periodic covariance structure, for $n \in \{150, 300, 600\}$ observations, amplitudes of 0.2–0.8, and frequencies of 1–5. To show the effect of the choice of kernel on the log Bayes factors, we show the results of two distinct kernels here, namely the exponential and periodic functions. Each point in the figure represents a single connection and is colored based on the strength of its evidence.

For both kernels, the Bayes factors indicate that the strength of evidence increases with the number of observations. For $n = 150$ observations, a large number of edges is inconclusive. Moreover, in line with the recall scores from Figure 7, the number of connections with conclusive evidence increases with an increasing amplitude and a decreasing frequency. In the case of the exponential kernel, it can bee seen that there are more connections with inconclusive evidence with

more rapid dynamics or dynamics with a lower amplitude, as these are estimated less accurately (also see figure 6C and Figure 6D). For the periodic kernel, the strength of evidence hardly decreases with frequency.

# 4    Results of empirical studies

We applied both statistical frameworks to the n-back working memory task fMRI dataset from the Human Connectome Project [17] to demonstrate the use of hypothesis testing for dynamics in an empirical setting. As briefly mentioned in Section 2.6, we expect to observe a difference in functional connectivity between high and low working memory load within the frontoparietal network and the default mode network, as these regions are found to be consistently associated with working memory [37, 35]. Specifically, we expect the functional connectivity between the dorsolateral prefrontal cortex (DLPFC) and the inferior parietal lobe (IPL) to increase during high working memory load (in the 2-back task), as compared to low working memory load (in the 0-back task). Importantly, we expect this connectivity to change dynamically over the task paradigm. While we expect no task-related dynamics between the primary auditory cortex (A1) and the IPL and between the DLPFC and the A1, some dynamics might still be present due to brain activity independent of the task.
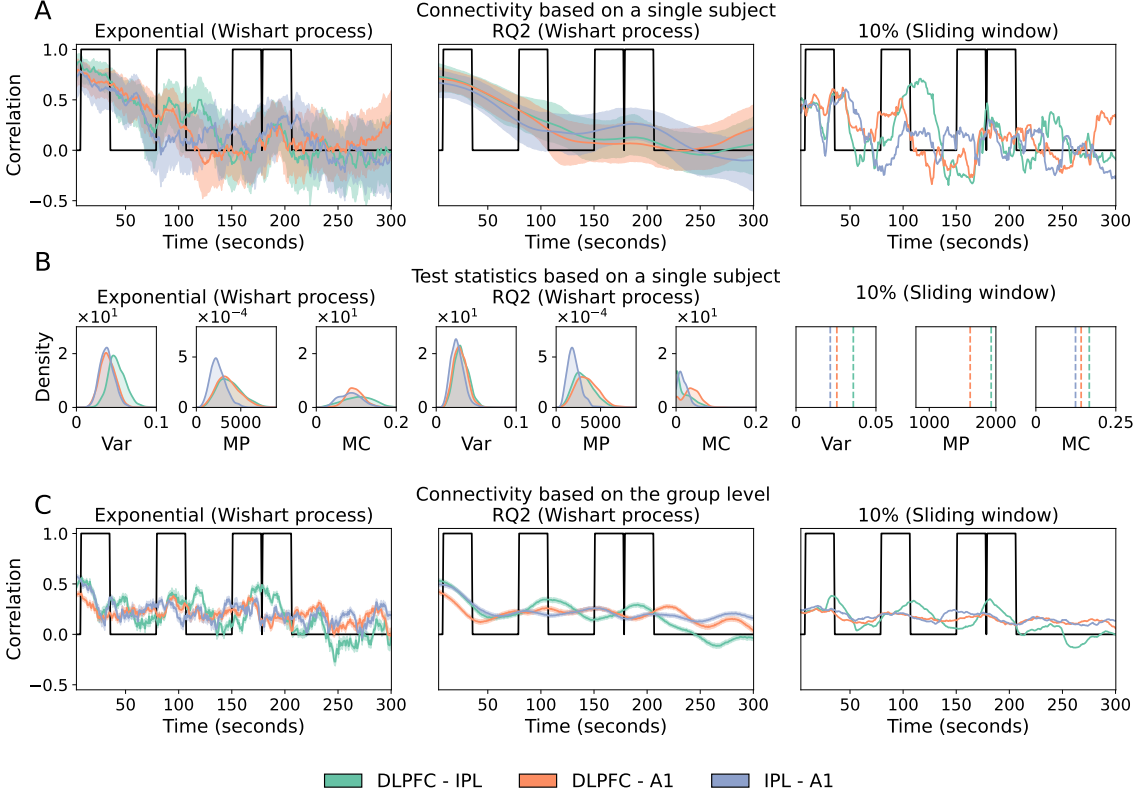
## 4.1 Connectivity estimates



**Figure 9: Estimates of functional connectivity between the dorsolateral prefrontal cortex, inferior parietal lobe and primary auditory cortex, and its corresponding variance (Var), maximum power (MP) and median-crossings (MC).** Subplot A shows the connectivity estimates based on a single representative subject, and subplot B shows the corresponding test statistics. For the Wishart process, we have a full posterior distribution over the test statistic, whereas for the sliding-window we have a point estimate. Subplot C shows the estimates on the group level (the corresponding statistics are shown in Fig. 11). The 2-back task design is shown in black.

Figure 9 visualizes a few examples of estimates of functional connectivity between the DLPFC, IPL and A1 regions. Figure 9A shows these estimates based on an individual representative subject, illustrating that the estimate with the RQ2 kernel is smooth, whereas the other two estimates show more fluctuations. Moreover, Figure 9B shows the corresponding posterior distributions or point estimates of the different statistics, indicating that the differences between regions are quite small and vary across modeling choices. Additionally, we estimated functional connectivity based on the group level, using the hierarchical Wishart process from Section 2.2.3. For the sliding-window method, the group-level estimates were obtained by averaging across all individual subject estimates. A few representative examples of estimates are shown in Figure 9C. Based on these group level estimates, we can already see that the functional connectivity seems to fluctuate with the task paradigm. Moreover, compared to the estimates on a single subject, the uncertainty in the connectivity estimates is lower at the group level.

## 4.2 Differences in amount of dynamics in functional connectivity
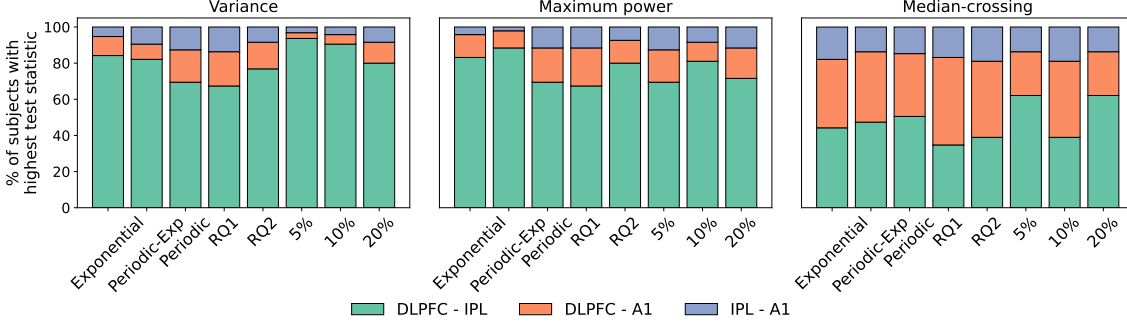


**Figure 10: Ordering of test statistics for three pairs of brain regions from the working memory task.** The figure indicates the number of times (in percentage of all subjects) that each pair of regions has the largest value, indicating that this connection is most dynamic.
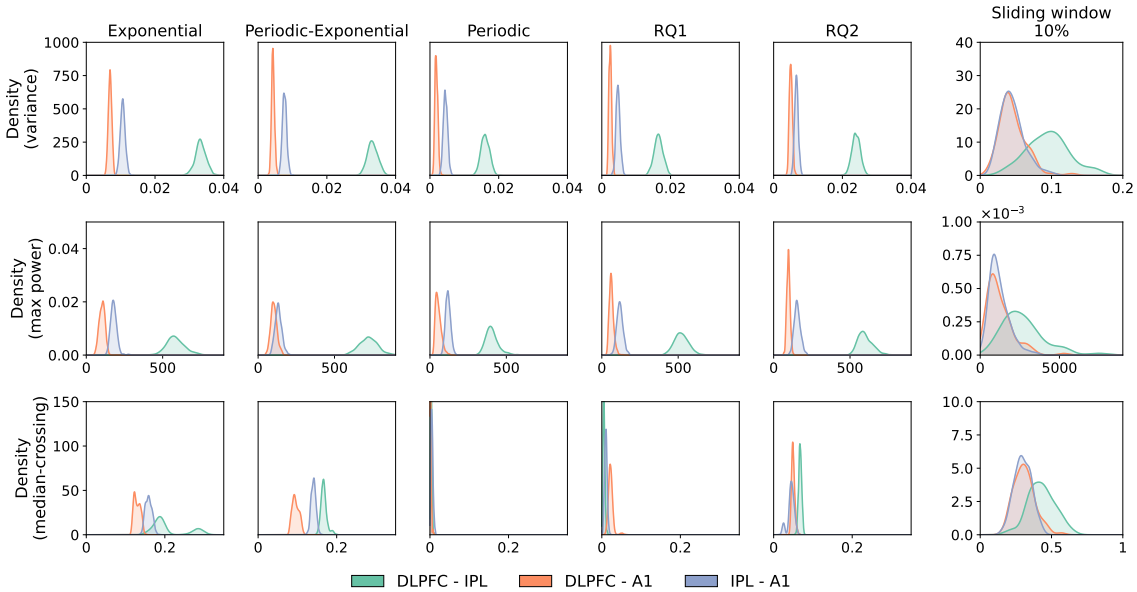


**Figure 11: Group-level statistics based on the working memory dataset.** The statistics were based on the estimates by the Wishart process and sliding-window method for the DLPFC–IPL, DLPFC–A1 and IPL–A1 connections. The results for five different kernels and a window length of 10% are shown. The figure shows the posterior distributions on the group-level (by the hierarchical Wishart process model) and the distributions of sliding-window estimates over all subjects.

In Figure 10 and Figure 11, the variance, maximum power and median-crossing values are shown for all three connections under different modeling choices. Figure 10 summarizes the subject-level results by the percentage of subjects in which each connection showed the largest value. For the Wishart process, we used the posterior mean. The results follow the same pattern as what we observed in Figure 9, namely, in line with our expectations, the DLPFC–IPL connection was most often ranked as the most dynamic connection. The results varied across kernels and window sizes, with the median-crossings showing the least clear differences between connections. For the maximum power, the periodic-exponential kernel shows the clearest preference for DLPFC–IPL

18

connection, while for the highest variance statistic, the 5% and 10% sliding-window estimates show this most clearly.

Figure 11 shows the group-level results across all subjects. For the Wishart process, these results are based on the posterior distributions inferred by the hierarchical Wishart process from Section 2.2.3. For the sliding-window method, the group-level distributions were obtained by applying kernel density estimation over the individual subject estimates. In general, the results based on variance and maximum power are in line with our prior expectations. Namely, the test statistics of the DLPFC–IPL connection are consistently higher than the DLPFC–A1 and IPL–A1 connections across all kernels and window sizes. This suggests that the DLPFC–IPL connection is indeed most dynamic. A similar pattern can be observed for the median-crossings, but the differences between connections are smaller. Overall, the DLPFC–IPL connection is consistently found to be most dynamic at the group level for the variance and maximum power, but results on the individual level vary across modeling choices. This highlights the importance of carefully defining these choices.
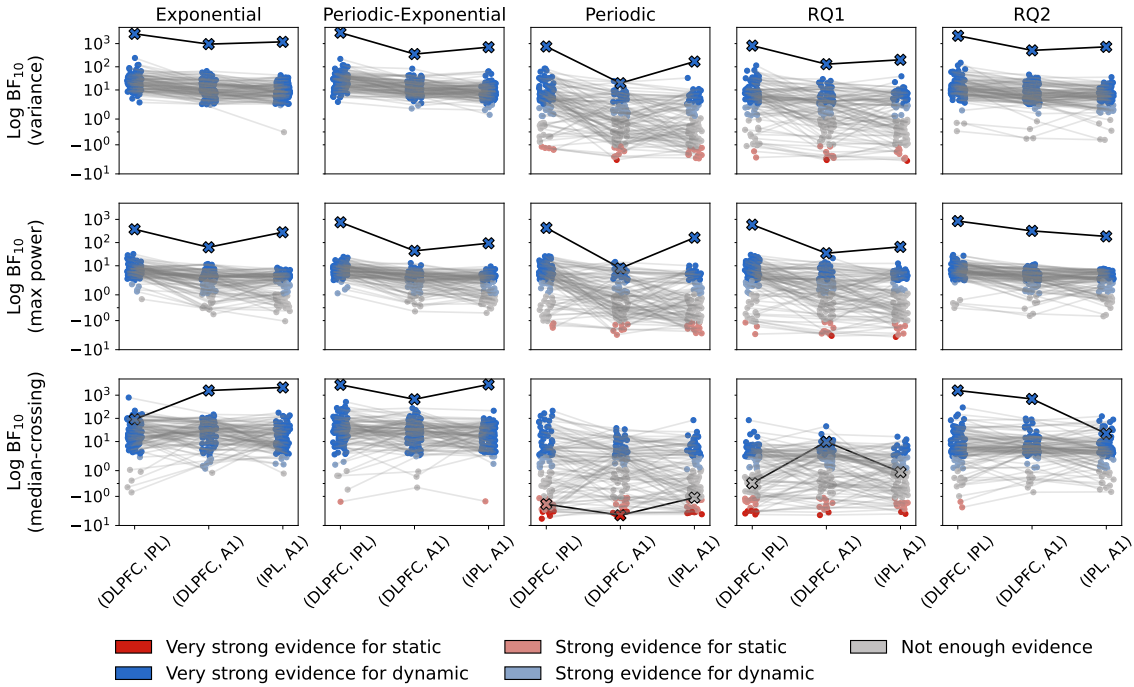


**Figure 12: Log Bayes factors of detecting dynamics in functional connectivity during the n-back task.** Every dot represents a connection of a single subject and is colored based on the amount of evidence of the connection being either dynamic or static. The lines indicate the corresponding pairs of brain regions of each subject. Moreover, the crosses indicate the evidence in the connectivity across all subjects, as estimated by the hierarchical Wishart process.

Finally, Figure 12 presents the log Bayes factors for all kernels, statistics and connections. Bayes factors of individual participants are indicated by dots, with gray lines connecting the values across region pairs. The crosses indicate the group-level Bayes factors. In the top row, results on the variance statistic are shown, showing that, in general, the DLPFC–IPL connection shows stronger evidence for dynamics than the other two connections. However, on the individual level there are large differences across test statistics and modeling choices. This becomes especially clear from the wide spread in Bayes factors for all three statistics when using the periodic kernel. Overall, these results highlight how different modeling choices and statistics can largely influence the conclusions about functional connectivity being dynamic. At the group level, results are more consistent across modeling choices. Here, for both the variance and maximum power, the DLPFC–

IPL connection clearly shows stronger evidence for dynamics than the other two connections. However, the median-crossing results show a different pattern which is inconsistent with our prior expectations, once more highlighting the importance of carefully selecting a test statistic.

# 5  Discussion

There is an increasing interest in modeling time-varying functional brain connectivity, in which the strength of connectivity does not remain static, but varies as a function of some input variable, such as time [1, 2, 3, 41]. To determine if this time-varying connectivity reflects true dynamics, or only fluctuates due to noise or modeling choices, statistical testing is needed [9, 10]. In this work, we introduced a Bayesian hypothesis testing framework for dynamic connectivity, which makes use of Wishart processes to model functional connectivity over time [1, 15, 16, 19]. Unlike in the frequentist alternative, which is commonly used to test for dynamics and where modeling assumptions are reflected in the choice of hyperparameters and interpretation method, the Bayesian framework encodes prior assumptions via prior distributions and modeling choices. Additionally, the Bayesian approach provides uncertainty in its connectivity estimate, and quantifies evidence for both dynamic and static connectivity.

Using simulations with different latent covariance structures, we evaluated both the Bayesian and frequentist frameworks. Overall, our simulations showed that both the Wishart process and the sliding-window method accurately estimate different types of connectivity, but the quality depends heavily on modeling assumptions. For example, with the use of smooth kernels such as the rational quadratic kernel, or with the use of large window sizes, we were unable to model fast fluctuations in connectivity. Our hypothesis testing results showed that these modeling choices also impacted the test accuracies, with performance decreasing if the covariance itself was not captured accurately. Moreover, our simulations indicated that the outcomes are strongly influenced by choice of test statistic. While the variance and maximum power showed similar and reliable performances in terms of accuracy, recall, and false positive rates, median-crossing performed much worse. Compared to the other test statistics, median-crossing showed a reduced ability to detect dynamics, while also estimating many false positives for the exponential and periodic-exponential kernels.

Our results on the HCP n-back working memory task fMRI dataset demonstrated how the two frameworks can be used to study dynamic functional connectivity in an empirical setting. In the realistic scenario that all connections show dynamics to some degree, the story becomes more difficult. Therefore, we focused on distinguishing between specific connections based on their dynamics, rather than classifying whether these connections were static or dynamic. In general, similar to our findings in the simulations, both testing frameworks were largely impacted by modeling choices and test statistics. We selected three regions of interest to compare the two frameworks. However, hypothesis testing of specific connections of interest on the individual level is challenging because of variations in functional connectivity between individuals [42, 43]. While individual-level results were heterogeneous, results of our newly introduced hierarchical Wishart process model and the sliding-window group level estimates were in line with our expectations, showing that the connections between the DLPFC and IPL was modulated more strongly by working memory load than the connection between these areas and A1.

In general, both the simulation results and empirical application highlighted, for both frameworks, the effect that test statistic and modeling choices can have. Future work should carefully consider these modeling choices to see if the modeling assumptions that are being made match the prior expectations about the dynamics in the data. If the aim is to model fast fluctuations in functional connectivity, a large window size or smooth kernel should be avoided. Moreover, to reduce the risk of wrongly specifying prior model assumptions, Bayesian model averaging can be used to combine multiple modeling assumptions, for example by combining different kernels [44]. Additionally, cross-validation could be used to determine modeling choices, such as window size. For the hypothesis testing itself, we recommend to use the variance or maximum power, which both performed well in our simulations. In contrast, we recommend against median-crossing, as

it led to a large number of false positives as well as lower accuracy and recall in our simulations, and unreliable results in the empirical application.

There are several limitations that need to be taken into consideration. First, the Bayesian framework makes use of the Savage-Dickey density ratio to compute the Bayes factors. This density ratio can only be used if the prior of the test-relevant parameters are independent of the prior of the nuisance parameters [33, 45]. In the context of our work, this means that the time-varying functional connectivity for a given connection should be independent of the time-varying functional connectivity of all other pairs of brain regions. Since the estimated correlation matrices should be positive semi-definite, this constraint is violated when modeling more than two brain regions. To address this dependence between parameters, an alternative Savage-Dickey density ratio with a correction term has been proposed by Heck [45], but this it is not straightforward to implement due to the high dimensionality of the Wishart process. Although incorporating this correction term would be a valuable improvement of the Bayesian framework, we have not applied this in the current work, as our focus was on comparing the different modeling assumptions of Bayesian and frequentist frameworks.

Another limitation is the computational scalability of inference of the Wishart process, based on Sequential Monte Carlo sampling. Although more computationally efficient estimation procedures using variational inference have been proposed [15], these are less robust and risk failing to capture parts of the dynamics [19], so these approaches should be used with care. In future work, improvements could be made here by combining sparse approximations of Gaussian processes [46] and factorized covariance models [15, 47] with the Sequential Monte Carlo method. Finally, although we have covered a wide range of potential dynamic connectivity structures, more elaborate model-based approaches to simulate more realistic fMRI data [48, 49].

In summary, we have proposed an approach to hypothesis testing for dynamic functional connectivity using a Bayesian framework and we combined this approach with the Wishart process. Our work highlights the impact that modeling choices and test statistics have on the conclusions that we draw about dynamic functional connectivity for both the frequentist and Bayesian frameworks. Nevertheless, if both are carefully considered, we believe that the Bayesian framework provides a good alternative to the frequentist framework, as it quantifies both uncertainty in the connectivity estimate, as well as evidence for both static and dynamic functional connectivity. Future research should focus on making it easier and more transparent to determine which modeling assumptions should be chosen to model and test for dynamic functional connectivity, especially when there is limited information about the dynamic structure in the functional connectivity of the data. One potential direction for this is the use of Bayesian model averaging, or perhaps the use of stacked models [50], although for the Wishart process this would first require improving the scalability of the inference of the model. This could help researchers to more robustly estimate and draw conclusions about dynamics in functional brain connectivity.

# Funding

# Acknowledgements

# References

[1] O. P. Kampman, J. Ziminski, S. Afyouni, M. van der Wilk, and Z. Kourtzi, "Time-varying functional connectivity as Wishart processes," *Imaging Neuroscience*, vol. 2, pp. 1–28, 2024.

[2] R. Liégeois, J. Li, R. Kong, C. Orban, D. Van De Ville, T. Ge, M. R. Sabuncu, and B. T. Yeo, "Resting brain dynamics at different timescales capture distinct aspects of human behavior," *Nature Communications*, vol. 10, no. 1, p. 2317, 2019.

[3] D. J. Lurie, D. Kessler, D. S. Bassett, R. F. Betzel, M. Breakspear, S. Kheilholz, A. Kucyi, R. Liégeois, M. A. Lindquist, A. R. McIntosh, *et al.*, "Questions and controversies in the study of time-varying functional connectivity in resting fMRI," *Network Neuroscience*, vol. 4, no. 1, pp. 30–69, 2020.

[4] S. Alonso, A. Llera, and D. Vidaurre, "Can fMRI functional connectivity index dynamic neural communication and cognition?," *Biological Psychology*, p. 109074, 2025.

[5] V. D. Calhoun, R. Miller, G. Pearlson, and T. Adalı, "The chronnectome: time-varying connectivity networks as the next frontier in fMRI data discovery," *Neuron*, vol. 84, no. 2, pp. 262–274, 2014.

[6] D. Battaglia, T. Boudou, E. C. Hansen, D. Lombardo, S. Chettouf, A. Daffertshofer, A. R. McIntosh, J. Zimmermann, P. Ritter, and V. Jirsa, "Dynamic functional connectivity between order and randomness and its evolution across the human adult lifespan," *NeuroImage*, vol. 222, p. 117156, 2020.

[7] D. K. Saha, E. Damaraju, B. Rashid, A. Abrol, S. M. Plis, and V. D. Calhoun, "A classification-based approach to estimate the number of resting functional magnetic resonance imaging dynamic functional connectivity states," *Brain Connectivity*, vol. 11, no. 2, pp. 132–145, 2021.

[8] Ü. Sakoğlu, G. D. Pearlson, K. A. Kiehl, Y. M. Wang, A. M. Michael, and V. D. Calhoun, "A method for evaluating dynamic functional network connectivity and task-modulation: Application to schizophrenia," *Magnetic Resonance Materials in Physics, Biology and Medicine*, vol. 23, pp. 351–366, 2010.

[9] T. O. Laumann, A. Z. Snyder, A. Mitra, E. M. Gordon, C. Gratton, B. Adeyemo, A. W. Gilmore, S. M. Nelson, J. J. Berg, D. J. Greene, *et al.*, "On the stability of BOLD fMRI correlations," *Cerebral Cortex*, vol. 27, no. 10, pp. 4719–4732, 2017.

[10] R. Hindriks, M. H. Adhikari, Y. Murayama, M. Ganzetti, D. Mantini, N. K. Logothetis, and G. Deco, "Can sliding-window correlations reveal dynamic functional connectivity in resting-state fMRI?," *Neuroimage*, vol. 127, pp. 242–256, 2016.

[11] E. A. Allen, E. Damaraju, S. M. Plis, E. B. Erhardt, T. Eichele, and V. D. Calhoun, "Tracking whole-brain connectivity dynamics in the resting state," *Cerebral Cortex*, vol. 24, no. 3, pp. 663–676, 2014.

[12] R. Liegeois, T. O. Laumann, A. Z. Snyder, J. Zhou, and B. T. Yeo, "Interpreting temporal fluctuations in resting-state functional connectivity MRI," *NeuroImage*, vol. 163, pp. 437–455, 2017.

[13] U. Pervaiz, D. Vidaurre, C. Gohil, S. M. Smith, and M. W. Woolrich, "Multi-dynamic modelling reveals strongly time-varying resting fMRI correlations," *Medical Image Analysis*, vol. 77, p. 102366, 2022.

[14] D. A. Handwerker, V. Roopchansingh, J. Gonzalez-Castillo, and P. A. Bandettini, "Periodic changes in fMRI connectivity," *Neuroimage*, vol. 63, no. 3, pp. 1712–1719, 2012.

[15] C. Heaukulani and M. van der Wilk, "Scalable Bayesian dynamic covariance modeling with variational Wishart and inverse Wishart processes," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[16] A. G. Wilson and Z. Ghahramani, "Generalised Wishart processes," in *27th Conference on Uncertainty in Artificial Intelligence*, (Barcelona, Spain), pp. 736–744, 2011.

[17] D. C. Van Essen, S. M. Smith, D. M. Barch, T. E. Behrens, E. Yacoub, K. Ugurbil, W.-M. H. Consortium, *et al.*, "The WU-Minn Human Connectome Project: an overview," *NeuroImage*, vol. 80, pp. 62–79, 2013.

[18] M. G. Preti, T. A. Bolton, and D. Van De Ville, "The dynamic functional connectome: State-of-the-art and perspectives," *NeuroImage*, vol. 160, pp. 41–54, 2017.

[19] H. Huijsdens, D. Leeftink, L. Geerligs, and M. Hinne, "Robust inference of dynamic covariance using Wishart processes and sequential Monte Carlo," *Entropy*, vol. 26, no. 8, p. 695, 2024.

[20] N. Leonardi and D. Van De Ville, "On spurious and real fluctuations of dynamic functional connectivity during rest," *NeuroImage*, vol. 104, pp. 430–436, 2015.

[21] S. Shakil, C.-H. Lee, and S. D. Keilholz, "Evaluation of sliding window correlation performance for characterizing dynamic functional connectivity and brain states," *NeuroImage*, vol. 133, pp. 111–128, 2016.

[22] M. Lamos, R. Marecek, T. Slavícek, M. Mikl, I. Rektor, and J. Jan, "Spatial-temporal-spectral EEG patterns of BOLD functional network connectivity dynamics," *Journal of Neural Engineering*, vol. 15, no. 3, p. 036025, 2018.

[23] R. Meng, F. Yang, and W. H. Kim, "Dynamic covariance estimation via predictive Wishart process with an application on brain connectivity estimation," *Computational Statistics & Data Analysis*, vol. 185, p. 107763, 2023.

[24] Z. Zhang, "A note on Wishart and inverse Wishart priors for covariance matrix," *Journal of Behavioral Data Science*, vol. 1, no. 2, pp. 119–126, 2021.

[25] M. Ghosh and B. K. Sinha, "A simple derivation of the Wishart distribution," *The American Statistician*, vol. 56, no. 2, pp. 100–101, 2002.

[26] C. E. Rasmussen and C. K. I. Williams, *Gaussian processes for Machine Learning*. The MIT Press, 2005.

[27] C. Chang and G. H. Glover, "Time-frequency dynamics of resting-state brain connectivity measured with fMRI," *NeuroImage*, vol. 50, no. 1, pp. 81–98, 2010.

[28] A. Zalesky, A. Fornito, L. Cocchi, L. L. Gollo, and M. Breakspear, "Time-resolved resting-state brain networks," *Proceedings of the National Academy of Sciences*, vol. 111, no. 28, pp. 10341–10346, 2014.

[29] G. Casella and R. Berger, *Statistical inference*. CRC press, 2024.

[30] D. J. Biau, B. M. Jolles, and R. Porcher, "P value and the theory of hypothesis testing: an explanation for new researchers," *Clinical Orthopaedics and Related Research®*, vol. 468, no. 3, pp. 885–892, 2010.

[31] R. E. Kass and A. E. Raftery, "Bayes factors," *Journal of the American Statistical Association*, vol. 90, no. 430, pp. 773–795, 1995.

[32] J. M. Dickey and B. Lientz, "The weighted likelihood ratio, sharp hypotheses about chances, the order of a Markov chain," *The Annals of Mathematical Statistics*, pp. 214–226, 1970.

[33] E.-J. Wagenmakers, T. Lodewyckx, H. Kuriyal, and R. Grasman, "Bayesian hypothesis testing for psychologists: A tutorial on the Savage–Dickey method," *Cognitive Psychology*, vol. 60, pp. 158–189, May 2010.

[34] W. H. Thompson, C. G. Richter, P. Plavén-Sigray, and P. Fransson, "Simulations to benchmark time-varying connectivity methods for fMRI," *PLoS Computational Biology*, vol. 14, no. 5, p. e1006196, 2018.

[35] T. Piccoli, G. Valente, D. E. Linden, M. Re, F. Esposito, A. T. Sack, and F. D. Salle, "The default mode network and the working memory network are not anti-correlated during all phases of a working memory task," *PloS One*, vol. 10, no. 4, p. e0123354, 2015.

[36] M. F. Glasser, T. S. Coalson, E. C. Robinson, C. D. Hacker, J. Harwell, E. Yacoub, K. Ugurbil, J. Andersson, C. F. Beckmann, M. Jenkinson, *et al.*, "A multi-modal parcellation of human cerebral cortex," *Nature*, vol. 536, no. 7615, pp. 171–178, 2016.

[37] D. M. Barch, G. C. Burgess, M. P. Harms, S. E. Petersen, B. L. Schlaggar, M. Corbetta, M. F. Glasser, S. Curtiss, S. Dixit, C. Feldt, *et al.*, "Function in the human connectome: Task-fMRI and individual differences in behavior," *NeuroImage*, vol. 80, pp. 169–189, 2013.

[38] F. Sambataro, V. P. Murty, J. H. Callicott, H.-Y. Tan, S. Das, D. R. Weinberger, and V. S. Mattay, "Age-related alterations in default mode network: impact on working memory performance," *Neurobiology of aging*, vol. 31, no. 5, pp. 839–852, 2010.

[39] H. Hoijtink, J. Mulder, C. Van Lissa, and X. Gu, "A tutorial on testing hypotheses using the Bayes factor.," *Psychological methods*, vol. 24, no. 5, p. 539, 2019.

[40] K. B. Huth, J. M. Haslbeck, S. Keetelaar, R. J. van Holst, and M. Marsman, "Statistical evidence in psychological networks," *Nature Human Behaviour*, pp. 1–14, 2025.

[41] C. Ahrends, A. Stevner, U. Pervaiz, M. L. Kringelbach, P. Vuust, M. W. Woolrich, and D. Vidaurre, "Data and model considerations for estimating time-varying functional connectivity in fMRI," *NeuroImage*, vol. 252, p. 119026, 2022.

[42] G. K. Aguirre, E. Zarahn, and M. D'Esposito, "The variability of human, bold hemodynamic responses," *Neuroimage*, vol. 8, no. 4, pp. 360–369, 1998.

[43] S. Mueller, D. Wang, M. D. Fox, B. T. Yeo, J. Sepulcre, M. R. Sabuncu, R. Shafee, J. Lu, and H. Liu, "Individual variability in functional connectivity architecture of the human brain," *Neuron*, vol. 77, no. 3, pp. 586–595, 2013.

[44] M. Hinne, Q. F. Gronau, D. van den Bergh, and E.-J. Wagenmakers, "A conceptual introduction to Bayesian model averaging," *Advances in Methods and Practices in Psychological Science*, vol. 3, no. 2, pp. 200–215, 2020.

[45] D. W. Heck, "A caveat on the savage–dickey density ratio: The case of computing bayes factors for regression parameters," *British Journal of Mathematical and Statistical Psychology*, vol. 72, no. 2, pp. 316–333, 2019.

[46] S. Rossi, M. Heinonen, E. Bonilla, Z. Shen, and M. Filippone, "Sparse Gaussian processes revisited: Bayesian approaches to inducing-variable approximations," in *International Conference on Artificial Intelligence and Statistics*, pp. 1837–1845, AISTATS, 2021.

[47] D. B. Rowe, *Multivariate Bayesian statistics: models for source separation and signal unmixing.* Chapman and Hall/CRC, 2002.

[48] P. Wang, R. Kong, X. Kong, R. Liégeois, C. Orban, G. Deco, M. P. Van Den Heuvel, and B. Thomas Yeo, "Inversion of a large-scale circuit model reveals a cortical hierarchy in the dynamic resting human brain," *Science Advances*, vol. 5, no. 1, p. eaat7854, 2019.

[49] G. Deco, J. Cruzat, and M. L. Kringelbach, "Brain songs framework used for discovering the relevant timescale of the human brain," *Nature Communications*, vol. 10, no. 1, p. 583, 2019.

[50] B. Griffin, C. Ahrends, C. Gohil, F. Alfaro-Almagro, M. W. Woolrich, S. M. Smith, and D. Vidaurre, "Stacking models of brain dynamics to improve prediction of subject traits in fmri," *Imaging Neuroscience*, vol. 2, pp. 1–22, 2024.

[51] P. Del Moral, A. Doucet, and A. Jasra, "Sequential Monte Carlo samplers," *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 68, no. 3, pp. 411–436, 2006.

[52] A. Gelman and D. B. Rubin, "Inference from iterative simulation using multiple sequences," *Statistical Science*, vol. 7, no. 4, pp. 457–472, 1992.

[53] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, "JAX: composable transformations of Python+NumPy programs," 2018.

[54] J. Lao and R. Louf, "BlackJAX: Library of samplers for JAX," *Astrophysics Source Code Library*, pp. ascl–2211, 2022.

[55] M. F. Glasser, S. N. Sotiropoulos, J. A. Wilson, T. S. Coalson, B. Fischl, J. L. Andersson, J. Xu, S. Jbabdi, M. Webster, J. R. Polimeni, *et al.*, "The minimal preprocessing pipelines for the Human Connectome Project," *NeuroImage*, vol. 80, pp. 105–124, 2013.

# Supplementary Materials

## S1 Dynamic connectivity estimation methods

Here we provide the formal description of the sliding-window method that was described in Section 2.2.1, and the Wishart process that was described in Section 2.2.2.

### S1.1 Sliding-window dynamic connectivity estimation

For every input location $i = 1, \ldots, n$, we determine the starting point $l$ and end point $u$ of the window as $l = \min\left(\lfloor x_i - \frac{\lambda-1}{2} \rfloor, 1\right)$ and $u = \max\left(\lfloor x_i + \frac{\lambda-1}{2} \rfloor, n\right)$, and with these we create a subset of observations $\mathbf{S}_i = (\mathbf{y}_l, \ldots, \mathbf{y}_u)^\top$. For every new window, the input is shifted by $\tau$. With a stride length of 1, the total number of windows would therefore equal the number of observations, and larger stride lengths reduce the number of windows. Next, for each window, the pairwise covariance is computed as $\lfloor (n - \lambda)/\tau \rfloor + 1$

$$\mathbf{\Sigma}(x_i) = \frac{1}{\lambda - 1} \mathbf{S}_i^\top \mathbf{S}_i , \quad i = 1, \ldots, n . \tag{17}$$

By computing the pairwise covariance in Eq. (17) for every subset, we end up with an estimate of covariance as a function of the input $x_i$.

### S1.2 The Wishart process

As explained in Section 2.2.2, the Wishart process is constructed from a scale matrix $\mathbf{L}$ and Gaussian processes that are parameterized by kernel hyperparameters $\theta$. To complete the generative model of the Wishart process, we define priors on these two parameters:

$$
\begin{aligned}
\mathbf{L} &\sim p(\mathbf{L}) \\
\theta &\sim p(\theta) \\
f_{kj}(x_i) &\sim \mathcal{GP}(\mu(\cdot), \kappa(\cdot, \cdot; \theta)) && i = 1, \ldots n , \quad k = 1, \ldots v , \quad j = 1, \ldots, d \\
\mathbf{\Sigma}(x_i) &= \sum_{k=1}^{v} \mathbf{L}\mathbf{f}_k(x_i) \mathbf{f}_k(x_i)^\top \mathbf{L}^\top && i = 1, \ldots, n \\
\mathbf{y}_i \mid x_i &\sim \mathcal{MVN}_d(\mathbf{0}, \mathbf{\Sigma}(x_i)) && i = 1, \ldots, n .
\end{aligned}
$$

Inferring the Wishart process from data, that is, computing the posterior distribution $p(\mathbf{\Sigma}(\mathbf{x}) \mid \mathbf{Y})$, is substantially more challenging than for the sliding-window approach. We follow the Wishart process inference procedure as described in [19]. The Sequential Monte Carlo sampler [51] requires a choice for the number of particles to be initialized, and the number of MCMC mutation steps to use. We initialized 1000 particles. The number of MCMC mutation steps varied and was based on convergence of the posterior $p(\mathbf{\Sigma}(\mathbf{x}) \mid \mathbf{Y})$, as measured by the Potential Scale Reduction Factor ($\hat{R}$) [52]. We used three parallel inference chains, and set the number of mutation steps such that $\hat{R} < 1.1$ for the covariance matrices. Although the kernel hyperparameters and scale matrix did not always converge across chains, the covariance itself did. This can happen because different combinations of the hyperparameters can result in the same covariance matrix, but since we are only looking at the covariance itself, our estimates of functional connectivity over time will be reliable. We implemented everything in Python. The Wishart process was implemented using the Python JAX framework [53] and SMC sampling in Blackjax [54], a Python library that builds on the JAX framework.

## S2 The n-back task paradigm

Here we describe the working memory n-back task paradigm mentioned in Section 2.6. The task was a working memory n-back task in which subjects were shown a sequence of images and had

to press a button if the current image matches the one shown $n$ steps back. The task alternated between blocks of 0-back and 2-back conditions, where the 0-back condition requires the participant to respond when the current stimulus matches a target image, and the 2-back condition asks the participant to respond when the current stimulus matches the one shown 2 steps back. Hence, the 2-back blocks require continuous updating of the sequence. The stimuli were images of either faces, places, tools, or body parts. A more detailed description of the task and data collection procedure is provided by Barch et al. [37].

The dataset was preprocessed using the preprocessing pipeline as described by Glasser et al. [55]. Additionally, we used high-pass filtering to remove any fluctuations slower than 0.008 Hz. The data contained 405 time points, which equals 291.6 seconds of fMRI recording. We removed the first five time points to reduce any potential artifacts related to the scanner equilibrium.

## S3    Additional simulation study results

Here we present the results of all simulations that were not presented in the main text. We follow the same structure as Section 3: we start with evaluating the accuracy of the connectivity estimates themselves, and then focus on the hypothesis testing results and uncertainty provided by the Bayesian framework.

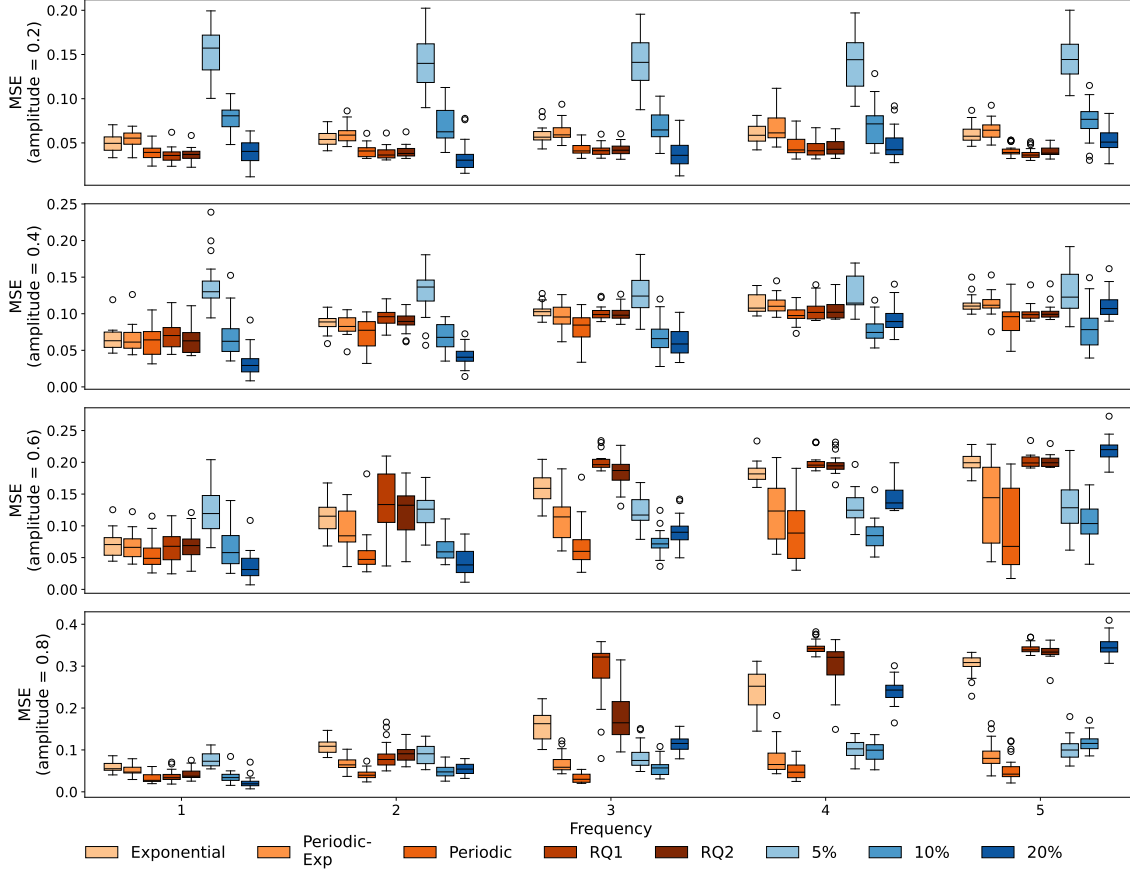## S3.1 Accuracy of connectivity estimates



**Figure S1: Mean squared errors between true and estimated correlations for the periodic simulations with 150 observations.** The figure shows the performance on the simulations with amplitudes of 0.2–0.8. For the Wishart process, the mean squared errors were computed over the entire posterior.
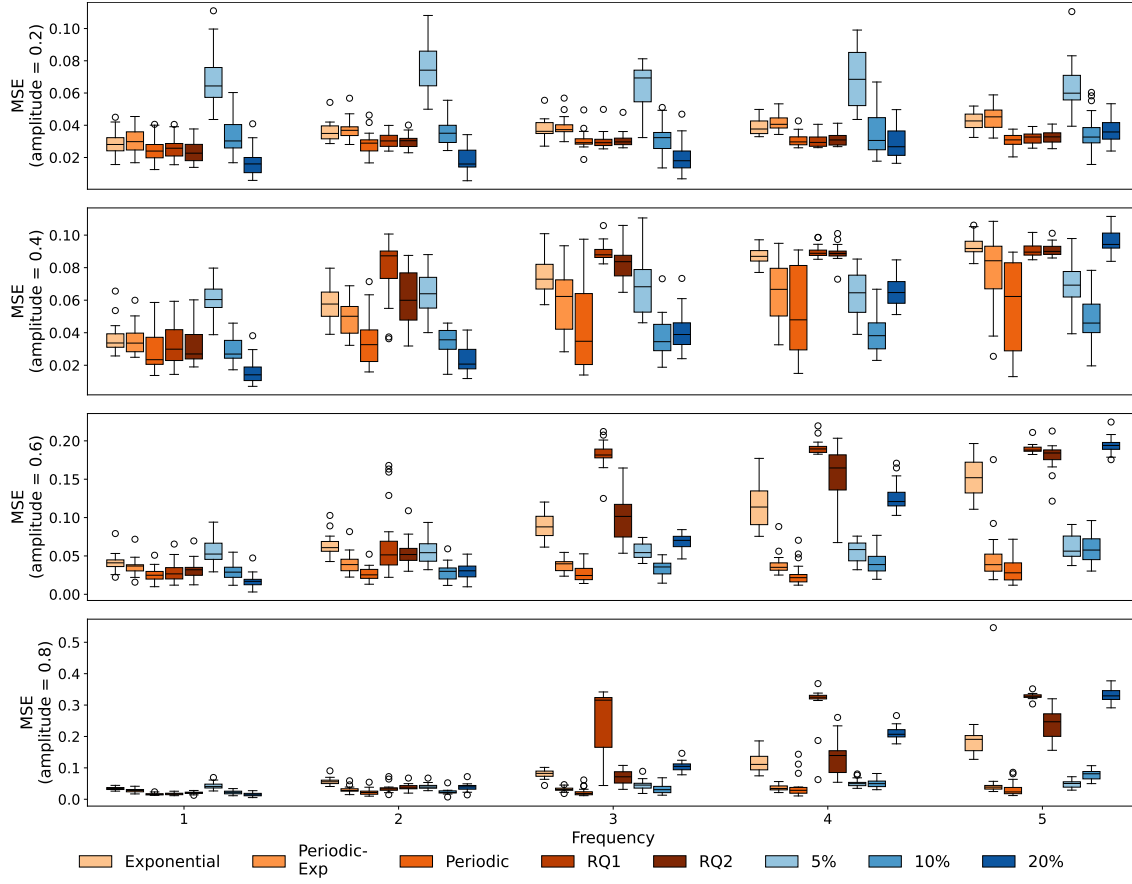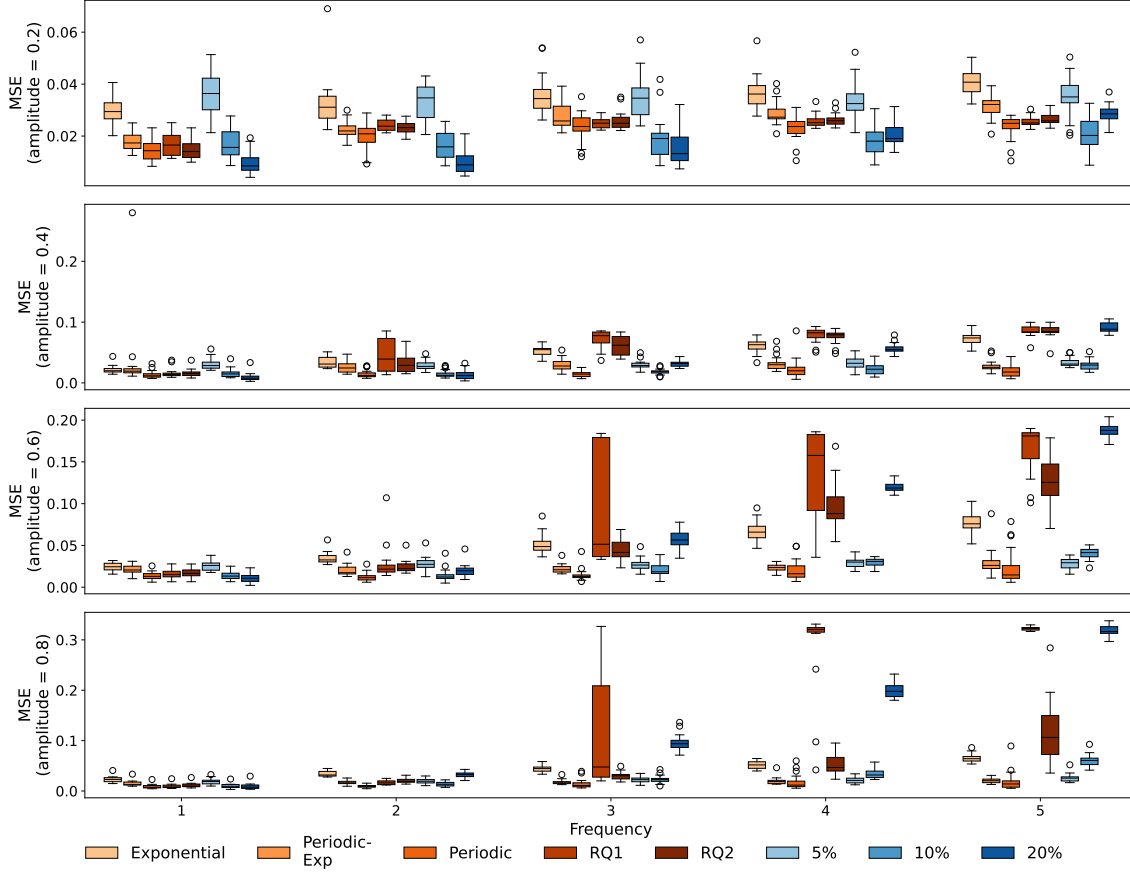
**Figure S2: Mean squared errors between true and estimated correlations for the periodic simulations with 300 observations.** The figure shows the performance on the simulations with amplitudes of 0.2–0.8. For the Wishart process, the mean squared errors were computed over the entire posterior.

**Figure S3: Mean squared errors between true and estimated correlations for the periodic simulations with 600 observations.** The figure shows the performance on the simulations with amplitudes of 0.2–0.8. For the Wishart process, the mean squared errors were computed over the entire posterior.

In Figure S1, Figure S2 and Figure S3, we provide all periodic mean squared error results for $n = 150$, $n = 300$ and $n = 600$ observations.

**Figure S4: Mean squared errors between true and estimated correlations for all simulations with 300 observations, using the posterior means.** Subplot A shows the results on the non-periodic state-switching correlations. Subplot B shows the results for the static correlations. Subplots C–F show the performance on the periodic simulations with amplitudes of 0.2–0.8. For the Wishart process, the mean squared errors were computed over the posterior means.

In figures 6, S1, S2 and S3, the mean squared error for the Wishart process estimates was computed over the entire posterior estimate, by first taking the mean squared error per sample and then averaging. To compare the performance with the sliding-window method and show the difference between the two metrics, Figure S4 presents the results for $n = 300$ observations when using only the posterior mean to compute the mean squared error. We observe that the MSE for the Wishart process, averaged over the MSEs of posterior samples, is larger than the MSE over the mean of the posterior. This is a standard effect of the Bayesian framework. Namely, as the posterior mean has a low MSE, individual samples that deviate from this mean will be less accurate. Importantly, this uncertainty can provide additional insights, such as when the Wishart process is very confident about its connectivity estimate, and when it is less certain.

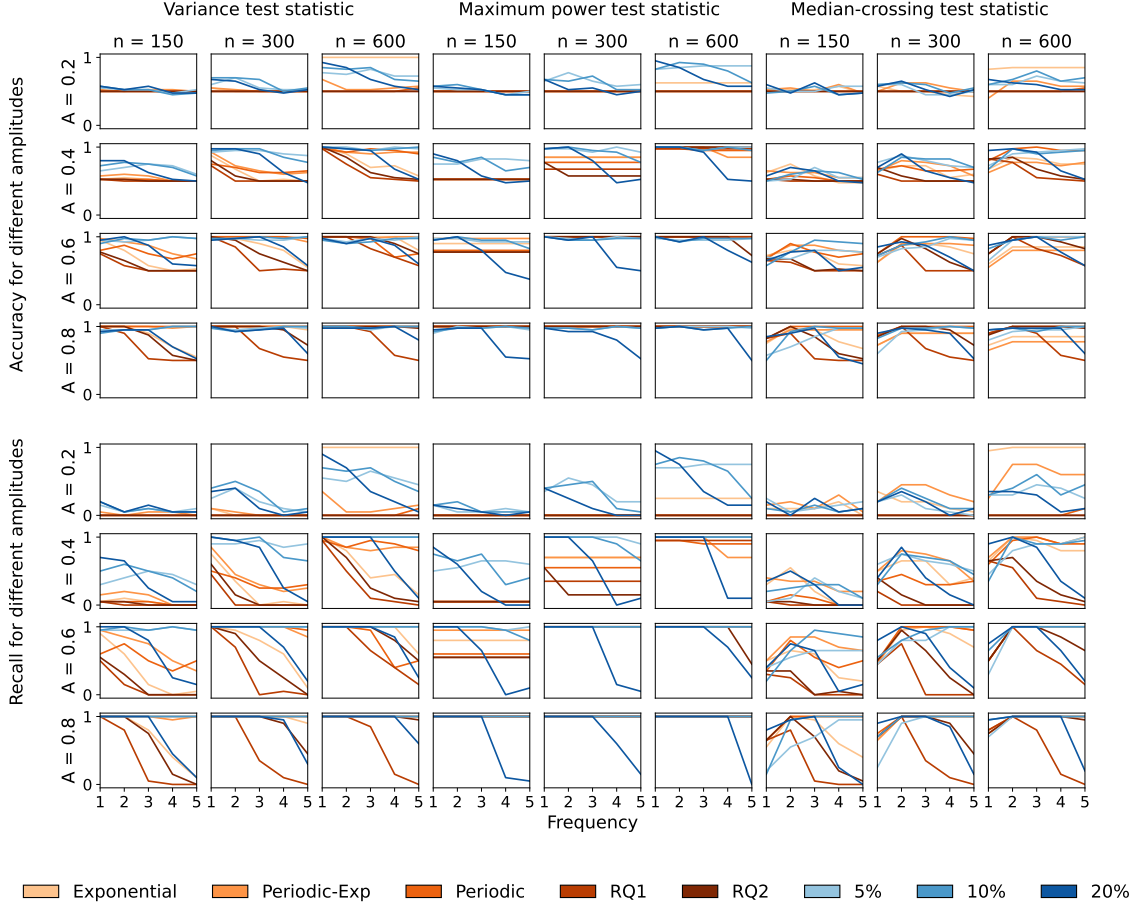## S3.2 Hypothesis test performances



**Figure S5: Hypothesis test performances on all periodic simulations.** Performance for amplitudes of 0.2–0.8 are shown. The figure shows the fractions of correctly classified connections, out of all connections (accuracy) and out of only dynamic connections (recall) when using the variance, maximum power and median-crossing test statistics.

Figure S5 shows the accuracy and recall of all periodic simulation studies for the three statistics. These results show that the performance increases with more observations.

**Figure S6: Hypothesis test performances on the state-switching and static simulations.** The figure shows accuracies and recalls (on the state-switching simulations) and false positive rates (on the static simulations) of the variance, maximum power and median-crossing statistics.

Figure S6 presents the results for the non-periodic state-switching and the static simulations. Here, the false positive rate indicates how many of the static simulations were incorrectly classified as dynamic. Overall, performance increases with more observations, which is in line with the results on the periodic simulations. In the state-switching simulations, the exponential and periodic-exponential kernels tend to perform best, along with the different sliding-window approaches. This is true for both the variance and the maximum power. The median-crossing statistic is less accurate than the others and shows a particularly high false positive rate in combination with the periodic-exponential and exponential kernels.

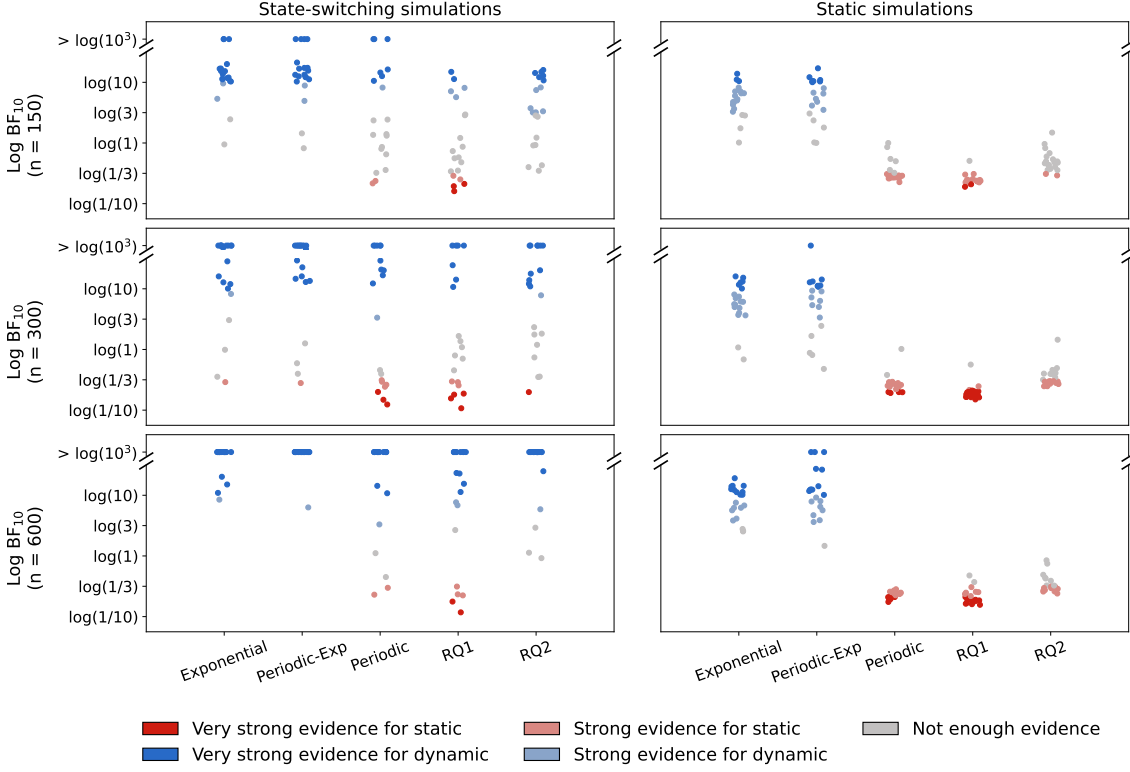## S3.3    Uncertainty in Bayesian hypothesis testing



**Figure S7: Log Bayes factors of detecting dynamics for the state-switching and static simulations.** Bayes factors are based on the variance test statistic. We show the results of the simulations with $n = 150$, $n = 300$ and $n = 600$ observations. Every dot represents a single connection and is colored based on the amount of evidence of the connection being dynamic or static.

Figure S7 shows the distributions of log Bayes factors for the state-switching and static simulations. The results on the state-switching simulation further support the finding that a sufficient number of observations is needed to find strong evidence for dynamics. Moreover, the results of the static simulations show that, depending on the kernel choice, the strength of evidence differs. Even with $n = 600$ observations, the strength of evidence for most edges estimates by the exponential and periodic-exponential kernels remains inconclusive. The results indicate that both the amount of available observations and the modeling choices strongly influence the outcomes of the hypothesis tests. With a limited number or observations or a less suitable modeling choice, such as a smooth kernel while the connectivity has fast fluctuations, it may not be possible to find enough evidence for detecting a dynamic functional connection.
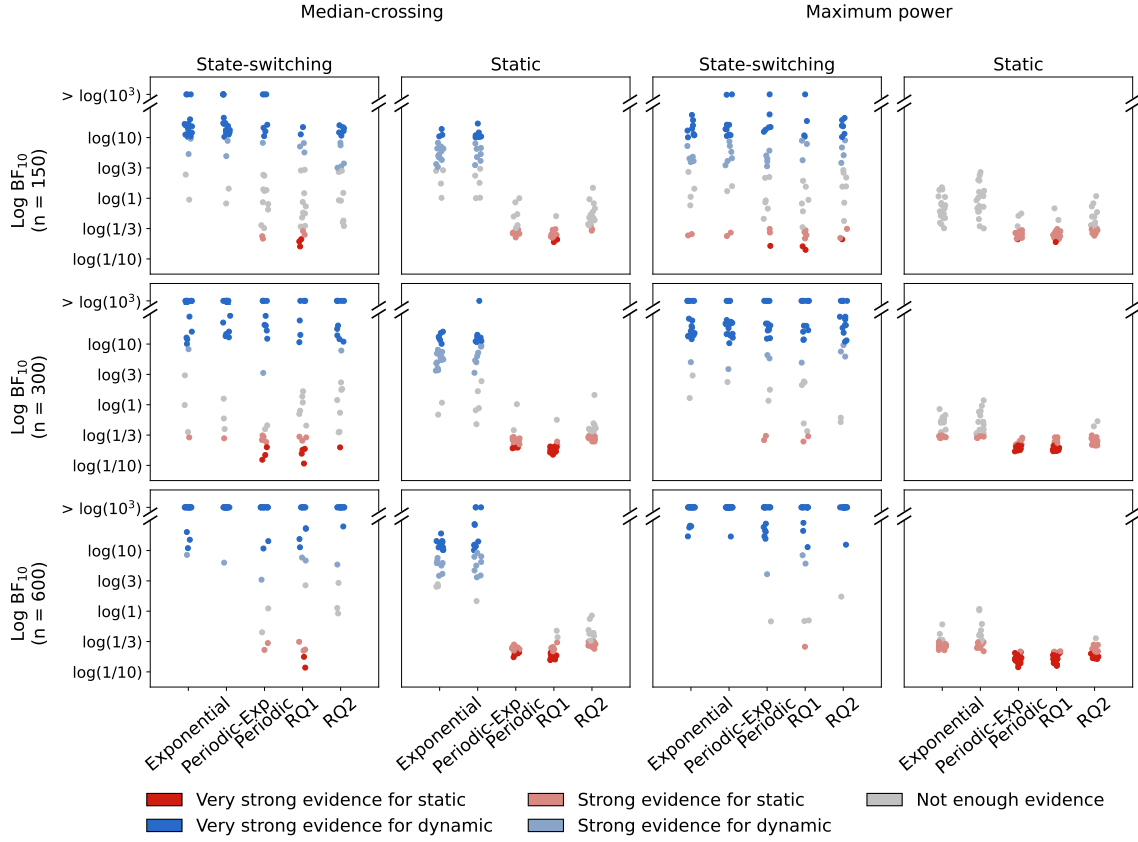
**Figure S8: Log Bayes factors of detecting dynamics for the state-switching and static simulations.** The figure shows the Bayes factors based on the median-crossing and maximum power test statistics. Every dot represents a single connection and is colored based on the amount of evidence of the connection being dynamic or static.

Figure S7 only presents the Bayes factor distributions based on the variance test statistic, the results based on the median-crossing and maximum power test statistic are provided in Figure S8. Based on all three statistics, a consistent pattern is shown for the static simulations. Namely, the exponential and periodic-exponential kernels provide higher Bayes factors compared to the other kernels and these Bayes factors also clarify the large number of false positives that we observed in Figure S6.

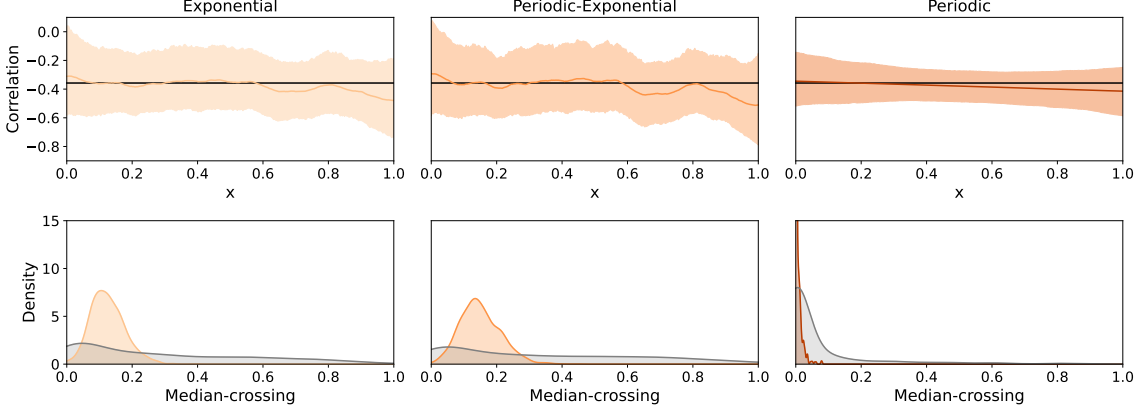# S4 Effect of different prior assumptions on the evidence



**Figure S9: Static connectivity estimates by the Wishart process and their corresponding prior and posterior median-crossing distributions.** Estimates using different kernels are shown. Prior distributions are shown in gray. The black line in the upper plots indicates the true latent connectivity.

Figure S9 provides an explanation for the large Bayes factors that we observed in Figure 3, especially by the exponential and periodic-exponential kernel. This figure shows examples of connectivity estimates and its corresponding prior (in gray) and posterior distributions over the median-crossing statistic. All examples are based on a static true connectivity. For the periodic kernel, the prior probability of the test statistic being zero is much larger than for the exponential and periodic-exponential kernels, indicating a stronger prior bias towards static connectivity. In contrast, the exponential and periodic-exponential kernels have a relatively low prior probability of the test statistic being zero, indicating that these kernels model more dynamic correlations by design. This behavior is also observed in the prior covariance samples in Figure 4, where the exponential and periodic-exponential kernels produce samples that contain many small fluctuations compared to the relatively smooth periodic kernel. These differences in prior distributions directly influence the resulting log Bayes factors through the Savage-Dickey density ratio (described in (12)), as the Savage-Dickey density ratio compares the prior and posterior densities at the null value, which in our case is at a value of zero. As a result, even if two kernels would provide nearly identical posterior distributions over a test statistic, their log Bayes factors can substantially differ because of their differences in priors. For example, if all three kernels would find a posterior test statistic centered far above zero (suggesting dynamic connectivity), the Bayes factor for the periodic kernel would be stronger because it has a larger density at zero.
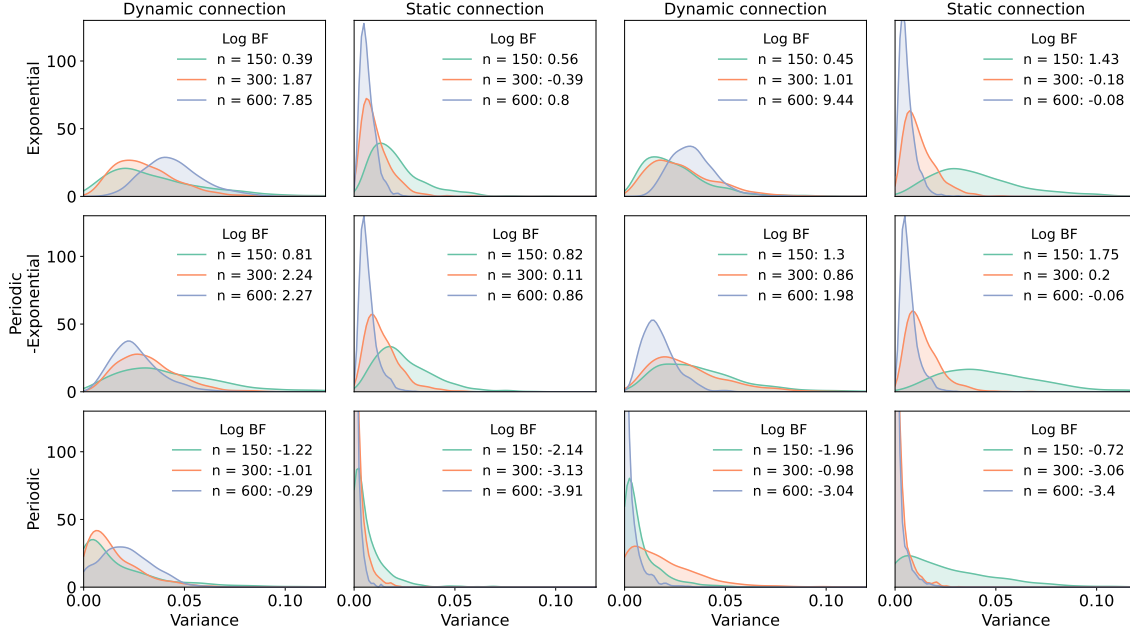
**Figure S10: Posterior variance distributions for different numbers of observations and different kernel choices.** The figure presents the results of four different connections, of which two are static and two are dynamic with a frequency of 2 and an amplitude of 0.2.

Moreover, Figure S10 shows the posterior variance distributions and corresponding log Bayes factors for four connections, of which two were simulated static and two dynamic with a frequency of 2 and an amplitude of 0.2. This figure shows that especially the exponential kernel is able to capture dynamics with a small amplitude, as indicated by the relatively low density at a variance of 0.