



Detecting and Localizing Differences in Functional Time Series Dynamics: A Case Study in Molecular Biophysics

Shahin Tavakoli & Victor M. Panaretos

To cite this article: Shahin Tavakoli & Victor M. Panaretos (2016) Detecting and Localizing Differences in Functional Time Series Dynamics: A Case Study in Molecular Biophysics, Journal of the American Statistical Association, 111:515, 1020-1035, DOI: [10.1080/01621459.2016.1147355](https://doi.org/10.1080/01621459.2016.1147355)

To link to this article: <http://dx.doi.org/10.1080/01621459.2016.1147355>



© 2016 The Author(s). Association of American Geographers© Shahin Tavakoli and Victor Panaretos



View supplementary material [↗](#)



Accepted author version posted online: 22 Mar 2016.
Published online: 18 Oct 2016.



Submit your article to this journal [↗](#)



Article views: 177



View Crossmark data [↗](#)

Detecting and Localizing Differences in Functional Time Series Dynamics: A Case Study in Molecular Biophysics

Shahin Tavakoli^a and Victor M. Panaretos^b

^aStatistical Laboratory, University of Cambridge, Cambridge, UK; ^bInstitut de Mathématiques, Ecole Polytechnique Fédérale de Lausanne, Lausanne, Switzerland

ABSTRACT

Motivated by the problem of inferring the molecular dynamics of DNA in solution, and linking them with its base-pair composition, we consider the problem of comparing the dynamics of functional time series (FTS), and of localizing any inferred differences in frequency and along curvelength. The approach we take is one of Fourier analysis, where the complete second-order structure of the FTS is encoded by its spectral density operator, indexed by frequency and curvelength. The comparison is broken down to a hierarchy of stages: at a global level, we compare the spectral density operators of the two FTS, across frequencies and curvelength, based on a Hilbert–Schmidt criterion; then, we localize any differences to specific frequencies; and, finally, we further localize any differences along the length of the random curves, that is, in physical space. A hierarchical multiple testing approach guarantees control of the averaged false discovery rate over the selected frequencies. In this sense, we are able to attribute any differences to distinct dynamic (frequency) and spatial (curvelength) contributions. Our approach is presented and illustrated by means of a case study in molecular biophysics: how can one use molecular dynamics simulations of short strands of DNA to infer their temporal dynamics at the scaling limit, and probe whether these depend on the sequence encoded in these strands? Supplementary materials for this article are available online.

ARTICLE HISTORY

Received November 2014
Revised December 2015

KEYWORDS

DNA minicircle; Functional data; Inverse problem; Multiple comparisons

1. Introduction

1.1. Functional Data Analysis

Functional data analysis (FDA; Ramsay and Silverman 2005; Ferraty and Vieu 2006; Horváth and Kokoszka 2012; Wang, Chiou, and Mueller 2016) deals with inferential situations where each data point that is best modeled as the realization of a stochastic process, understood as a random function or a random surface, such as weather data, neuroimages, electricity consumption curves, or phonetics, to name a few (e.g., Ramsay and Silverman 2002; Antoniadis, Paparoditis, and Sapatinas 2006; Aston and Kirch 2012a; Hadjipantelis et al. 2015). In a typical setting, one is interested in drawing inferences on the law of a random function $X \in L^2([0, 1], \mathbb{R})$ based on a sample $X_1, \dots, X_T \stackrel{\text{iid}}{\sim} X$. The main characteristics of interest are the mean function, and the covariance surface/operator. The mean function describes the average shape of the random object of interest, whereas the covariance operator encodes the second-order fluctuations of the random function around its mean. The covariance operator and its eigendecomposition are at the basis of the Karhunen–Loève expansion (Grenander 1981), which provides insight into the fluctuations of the random function X , and is the basis for optimal linear finite-dimensional approximations of X (through functional principal component analysis, or

fPCA), thus providing the means for applying multivariate techniques to functional data (e.g., Ferraty 2011).

Though estimation of the mean function and covariance operator in the iid setting of FDA is not substantially different from its multivariate counterpart, with \sqrt{T} -consistent and asymptotically Gaussian estimators (at least when assuming that the data X_1, \dots, X_T are readily available as curves, or densely sampled, see, e.g., Dauxois, Pousse, and Romain 1982; Mas and Menneteau 2003; Hall and Hosseini-Nasab 2006), statistical inference in the context of FDA typically involves an inverse problem, making it intrinsically harder from the multivariate setting. This problem can nevertheless be tackled by appropriate regularization—as exemplified by one-sample tests for the mean (Mas 2007), two-sample tests for the mean (Fan and Lin 1998; Cuevas, Febrero, and Fraiman 2004), and two-sample tests for covariance operators (Panaretos, Kraus, and Maddocks 2010; Kraus and Panaretos 2012; Horváth, Kokoszka, and Reeder 2013)—or through resampling techniques (e.g., Benko, Härdle, and Kneip 2009; Boente, Rodriguez, and Sued 2014; Paparoditis and Sapatinas 2014).

1.2. Functional Time Series

Despite being relevant for a broad range of applications, the iid setting of FDA is not appropriate in situations where there is a

CONTACT Shahin Tavakoli  s.tavakoli@statslab.cam.ac.uk  Department of Pure Mathematics and Mathematical Statistics, Statistical Laboratory, University of Cambridge, Wilberforce Road, Cambridge CB3 0WB, United Kingdom.

 Supplementary materials for this article are available online. Please go to www.tandfonline.com/r/JASA.

Published with license by Taylor and Francis

© 2016 Shahin Tavakoli and Victor Panaretos.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/3.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. The moral rights of the named author(s) have been asserted.

natural form of dependency in the data collection, such as when data are collected sequentially in time. Examples of such data include (but are not limited to) daily electricity consumption curves (Antoniadis, Paparoditis, and Sapatinas 2006), functional MRI data (Aston and Kirch 2012a), or molecular dynamics trajectories of DNA minicircles (see Section 2). In such cases, the data can be modeled as a stationary time series of functions, or stationary *functional time series*.

Inference for functional time series is classically carried out under functional autoregressive models, or more general linear models (see, e.g., Bosq 2000; Mas 2002; Ferraty and Romain 2011; Hörmann and Kidziński 2012; Aue et al. 2014; Aue, Norinho, and Hörmann 2015). Nevertheless, there is no a priori reason to expect that the behavior of a functional time series should be described by a linear process of any particular form or order: potentially nonlinear and/or non-Gaussian situations must be considered.

The problem of inference for functional time series beyond linear assumptions has only recently started to be addressed (Hörmann and Kokoszka 2010). Work has been done on assessing whether a dataset is independent against an ergodic alternative (Horváth, Hušková, and Rice 2013), or whether two functional time series are independent (Horváth and Rice 2015a). Horváth, Kokoszka, and Rice (2014) developed test for the stationarity of functional time series, and Kokoszka and Young (2016) considered tests for stationarity around a deterministic trend. The problem of two-sample testing for equality of the mean function of two functional time series, without resorting to linearity assumptions, has been considered by Horváth, Kokoszka, and Reeder (2013), Fremdt et al. (2014), and Horváth and Rice (2015b). Zhang et al. (2011) and Aston and Kirch (2012a,b) considered the change point detection of the mean function. Concerning inference on the second-order structure, Horváth, Kokoszka, and Reeder (2013) and Horváth, Rice, and Whipple (2014) proposed a consistent estimator for the long-run covariance operator, Kokoszka and Reimherr (2013) established the asymptotic normality of the sample covariance operator and its eigenfunctions, while Zhang and Shao (2015) considered the comparison of the covariance operators of two functional time series, a problem that has attracted considerable attention in the case of two collections of iid functional data (Panaretos, Kraus, and Maddocks 2010; Boente, Rodriguez, and Sued 2011; Kraus and Panaretos 2012; Horváth and Kokoszka 2012; Fremdt et al. 2013; Paparoditis and Sapatinas 2014; Pigoli et al. 2014; Boente, Rodriguez, and Sued 2014). The covariance operator, however, fails to capture any of the *dynamics* of a functional time series; and the long-run covariance operator captures only crude aspects of the time dynamics (being the sum of the autocovariance operators over lags). To capture the complete second-order dynamics of a functional series, Panaretos and Tavakoli (2013a,b) introduced a frequency domain framework (see also Hörmann, Kidziński, and Hallin 2015; Hörmann, Kidziński, and Kokoszka 2015), by means of the spectral density operator, the Fourier transform of the complete collection of auto-covariance operators; arguably, this would be the object upon which inferences related to dynamics ought to be based.

The methodological contribution of our article is the development of inferential tools for the comparison of the *complete* second-order dynamics (encoded by *all* the lag- t autocovariance operators) of two stationary functional time series, by

means of the frequency domain approach, as introduced in Panaretos and Tavakoli (2013a,b). Though this is arguably a core methodological problem in its own right, our study is motivated by the applied challenge of inferring and comparing the coarse-grained dynamics of DNA in solution, based on molecular dynamics simulation. The details of this problem, and its connection with functional time series are surveyed in the next section. Our methodological contributions will then be developed in parallel with a case study on inferring base-pair-dependent dynamical properties of DNA.

1.3. Molecular Biophysics

Molecular biophysics investigates, models, and characterizes the physical structures and dynamics that occur within a living organism at the molecular level (Glaser 2012, chap. 1). This study encompasses a bewildering variety of macromolecular structures, processes, and their interactions, the most iconic of which is arguably the biopolymer DNA. It is by now well understood how the structure of DNA encodes the entirety of genetic information of an organism in the sequence of its base-pairs, and how its double helix geometry with complementary steps built of these bases allow it to be transcribed or passed on from parent to offspring. The structural composition and geometrical arrangement is only part of the story, though: when fulfilling its biological purpose, the DNA polymer undergoes important mechanical maneuvers including twisting, bending, and looping (Garcia et al. 2007; Prévost, Takahashi, and Lavery 2009). For this reason, biophysicists are particularly interested in understanding the mechanical properties and dynamics of DNA (Mastroianni et al. 2009), and the relationship between the base-pair sequence composition, and the mechanical properties and dynamics of DNA (Peters and Maher 2010). Linking sequence to mechanical properties and dynamics would not only shed light into fundamental biological processes (such as transcription, regulation, and packing), but also holds promise in the use of DNA as a material for nanoengineering (Seeman 2005; Rothmund 2006).

The mechanics of DNA can be studied in a wide range of different scales, ranging from the fine-grained (atom-by-atom, for instance) to the coarse-grained (at the order of persistence length, i.e., about 160 base-pairs (Walter, Gonzalez, and Maddocks 2010; Gonzalez, Petkeviciūtė, and Maddocks 2013); or even at the order of thousands of base-pairs (Sambriski, Schwartz, and De Pablo 2009). These mechanics are fundamentally stochastic: the intrinsic conformational characteristics of the molecule are subjected to thermal fluctuations due to their surrounding environment. And though elaborate stochastic models exist that offer detailed descriptions of the atom-by-atom (or atomic level ensemble-by-ensemble) behavior of the molecule, the determination of their scaling limits at the more coarse-grained level is typically mathematically intractable. In this limit, one would consider the conformational mechanics of an entire strand of DNA, seen as a curve or function, rather than of its individual constituent atoms. It is therefore natural to ask whether by *observation* of DNA strands, one could statistically infer information on their coarse-grained limit, seen as a random curve, and indeed link it to their base-pair sequence composition. The natural setting for this would be the setting of *functional data analysis*, where one is precisely interested in probing the law of a random curve on the basis of several of its realizations.

To obtain sample DNA curves, one may perform *cyclization* experiments. Cyclization experiments involve short DNA strands whose ends meet, resulting in a loop-like configuration. The resulting closed curves are called *DNA minicircles*, and are excellent specimens for the study of DNA mechanics, as they yield a naturally stressed state in which intrinsic properties are amplified relative to extrinsic thermal fluctuations (Shore, Langowski, and Baldwin 1981; Shore and Baldwin 1983; Kahn and Crothers 1992). To link conformational variability to base-pair sequence, one may probe minicircles with slightly differing base-pair sequences, for example, CAP minicircles and TATA minicircles (Amzallag et al. 2006; see Section 2), whose base-pair composition coincides in all but a few basis steps. For example, based on three-dimensional reconstructions of a sample of CAP & TATA minicircles, obtained by cryo-electron microscopy, it appears that the differences between the two minicircles are *not* in their mean conformation (Amzallag et al. 2006), but in the way they vary around their mean conformation, as was established by means of a functional data analysis of their covariance structure (Panaretos, Kraus, and Maddocks 2010).

Microscopy-based studies, however, only allow *static* insight into the mechanics of DNA minicircles, and do not yield information on their dynamical properties, since they are based on “still images” of the DNA minicircles embedded in vitrified ice. The ideal kind of data needed for inferring the coarse-grained dynamics of DNA, and their link to sequence, would be in the form of a movie of DNA minicircles oscillating in solution. Though empirical acquisition of such data is as of yet infeasible, *in silico* surrogates can be created via Molecular Dynamics (MD) simulations (Leach 2001; Dryden et al. 2002; Lankas, Lavery, and Maddocks 2006; Freddolino et al. 2006; Sanbonmatsu and Tung 2007; Pérez, Luque, and Orozco 2011; Mitchell, Laughton, and Harris 2011; Curuksu, Kannan, and Zacharias 2014; Pasi et al. 2014). MD simulations are used to obtain the trajectory of a DNA minicircle moving in solution, and are obtained by numerically solving fine-grained atomic level models on multi-body interactions between all the atoms of water and DNA. These simulations are extraordinary in their computational and mathematical complexity. However, as was mentioned earlier, it is not the trajectories of each individual DNA atom that is of interest, but their joint behavior in the coarse grained limit, which can be thought as a function.

The motivation for this article is that of inferring dynamical properties of DNA through MD simulations in their scaling limit, and to investigate their dependence on base-pair sequence. The data we will be working with are MD trajectories of CAP & TATA minicircles oscillating in solution. We shall model their scaling-limits as *functional time series* (FTS). An FTS is a sequence $\{X_t : t \in \mathbb{Z}\}$, where t denotes the time index, and each X_t is a random function, say $X_t \in L^2([0, 1], \mathbb{R})$, representing the conformation of a minicircle at time t , seen as a continuous curve. The dynamics of the minicircles can be viewed through the lens of the second-order structure of the time series, which is contained in the collection of its lag- t autocovariance operators: these contain the covariation of the random function t time points apart. Understanding and comparing the dynamics of CAP and TATA minicircles can therefore be translated into the problem of estimation and inference for the second-order structure of functional time series.

In particular, we wish to be able to detect and further localize any differences either at the level of frequencies, and along the length of the curves. Our methodological work is thus presented in parallel with the DNA motivation, in the form of a case study. The article is organized as follows: in Section 2, we present the molecular dynamics simulation data on the TATA and CAP minicircles, including the necessary preprocessing steps, and the estimation of their functional dynamics through the spectral density operator. Section 3 considers the problem of detecting differences between the spectra of two time series, such CAP and TATA, by first comparing them at fixed frequencies—using a test for comparing the spectrum that we introduce in this article (Theorem 1)—and then adjusting for multiplicities to localize the differences in the frequencies, while controlling the overall significance level at which we pronounce detections. The detection of the differences between the CAP and TATA is further investigated in Section 4, where we consider the problem of first selecting frequencies at which the spectrum of CAP and TATA are different, and then detecting and localizing their differences on the minicircles, within each selected frequency. A supplementary file collects necessary technical details, proofs of our main results, and simulation studies.

2. Description of the Data

The dataset (produced by the group of Prof. John Maddocks, Laboratory for Computation and Visualization in Mathematics and Mechanics, Institut de Mathématiques, EPFL, Lausanne, Switzerland, <http://lcvmwww.epfl.ch/>) of our case study consists of the (time) trajectories of two DNA minicircles moving freely in solution. These trajectories are constructed by means of molecular dynamics (MD) simulations (Leach 2001; Freddolino et al. 2006; Lankas, Lavery, and Maddocks 2006; Sanbonmatsu and Tung 2007; Mitchell, Laughton, and Harris 2011; Mitchell and Harris 2013). Such MD simulations simulate the trajectory of a strand of DNA in water by numerical integration of a model taking into account multi-body interactions between all the atoms of the minicircle and the aqueous solution (actually, this is an oversimplification, but the precise description of MD simulations is not the goal of this article).

The two DNA minicircles are called CAP and TATA: they are built of 158 base-pairs (BP), and differ only in 14 BP (see Table 1). The MD simulation employed an integration step of 2 femtoseconds ($2 \cdot 10^{-15}$ sec) for the numerical integration algorithm, and the data were recorded every picosecond (10^{-12} sec). Time-wise, the data consist of 50,000 snapshots, where at each snapshot one retains the three-dimensional coordinates of the 158 BP centers of the minicircle (see Section A in the supplementary material for a description of the MD simulation protocol). These are the result of a massive computation, simulating a particle system consisting of 200,000 particles, and taking approximately 6000 CPU hours on a Cray XT5 at the Swiss National Supercomputing Centre—an ambitious molecular dynamics simulation (see also Freddolino et al. 2006; Sanbonmatsu and Tung 2007; Mitchell, Laughton, and Harris 2011, for more ambitious MD simulations).

In principle, the resulting time series is not guaranteed to be stationary (indeed, at least 300–500 nanoseconds seem to be typically needed for convergence to equilibrium, see, e.g., Dans

Table 1. The sequences of base-pairs for the CAP and TATA minicircles; the differences between the two sequences are in gray.

CAP
GATGAATTCACGGATCCGGTTTTTTTGCCCGTTTTTTGCCCGTTTTTTGCCCGTTTTTTGCCCGTTTTTT
GCCCGTTTTTTCCGGATCCGTACAGGAATTCTAGACCTAGGGTGCCTAATGAGTGAGCTAACTCACA
TTAATTTGCGTTGCGCCATGGAATC
TATA
GATGAATTCACGGATCCGGTTTTTTTGCCCGTTTTTTGCCCGTTTTTTGCCCGTTTTTTGCCCGTTTTTT
GCCCGTTTTTTCCGGATCCGTACAGGAATTCTAGACCTAGGGTGCCTAATGAGTGCCCTTTTATAGC
TTAAATCGCGTTGCGCCATGGAATC

et al. 2012, 2014; Lavery et al. 2014), but one expects its time increments to be stationary, at least locally in time. We therefore focus on the last 10,001 snapshots, and discard the first 39,999 snapshots to avoid burn-in period effects.

For each minicircle, the data we consider are 3D time series $\{M_t(j) : t = 1, 2, \dots, 10,001; j = 1, \dots, 158\} \subset \mathbb{R}^3$, where t denotes the time index (in picoseconds), and j is the BP index. In other words, $M_t(j) \in \mathbb{R}^3$ gives the coordinates of BP j at time t . In the sequel, we shall view the set $\{1, \dots, 158\}$ as the quotient group $\mathbb{Z}/158\mathbb{Z}$, so that $M_t(j+k)$ has a meaning for all $j, k \in \mathbb{Z}$.

The data are shown for various timepoints t in Figure 1. Since the data can be rotated or translated without altering the information they convey on the intrinsic minicircle dynamics, our analysis should hinge on features that are invariant to the action of the group of rigid motions. We therefore choose to work with the curvature of the DNA minicircles, an invariant with respect to this group (indeed an object typically studied in minicircle experiments). Further to solving the problem of registration of the data, focusing on the curvature also reduces the dimensionality of the data, transforming the object of study from a time series of curves in \mathbb{R}^3 to a time series of real-valued functions.

2.1. Preprocessing Steps

The estimation of curvature based on discrete noisy observations is often carried out using a plug-in approach. For instance, Sangalli et al. (2009) used free-knot regression splines to estimate the curve $\gamma(t)$, and then computed the curvature estimate

$$c(t) = |\hat{\gamma}'(t) \wedge \hat{\gamma}''(t)| / |\hat{\gamma}'(t)|^3, \tag{2.1}$$

where $\hat{\gamma}(t) \in \mathbb{R}^3$ is the estimated curve, $u \wedge v$ denotes the cross-product of $u, v \in \mathbb{R}^3$, and $|\cdot|$ denotes Euclidean norm. The performance of such techniques heavily depends on the value of a smoothness parameter, the choice of which can be rather subjective (Sangalli et al. 2009, sec. 4). Our own experience with plug-in estimation of the curvature was that it yielded estimates with too many degrees of freedom, even with roughness penalization of the estimated curve, and this symptom we attribute to the presence of the renormalization factor in (2.1). We therefore opted to start out with a discrete version of curvature, motivated by subject-specific knowledge on the larger scale behavior of DNA minicircles. Specifically, for each minicircle trajectory, we computed the curvature trajectory $\{c_t(j) : t = 1, 2, \dots, 10,001; j = 1, \dots, 158\} \subset \mathbb{R}_+$, where $c_t(j)$ is the curvature (inverse radius) of the circle passing through the three points $M_t(j-5), M_t(j), M_t(j+5)$. Recall that for three points $p_1, p_2, p_3 \in \mathbb{R}^3$, this curvature is given by $\text{curvature}(p_1, p_2, p_3) = 2| (p_2 - p_1) \wedge (p_3 - p_2) | / (|p_2 - p_1| \cdot |p_3 - p_2| \cdot |p_3 - p_1|)$.

The reason we selected triples separated by five base-pairs (corresponding to indices $j - 5, j, j + 5$), instead of consecutive ones (i.e., indices $j - 1, j, j + 1$), is that the DNA double helix performs a complete rotation in 11-12 BP, on average; taking the curvature of the directly adjacent BP would represent a highly local curvature measure, whereas the curvature computed on further spaced BP is more in line with the coarser scale at which we wish to understand their dynamics. A welcome side-effect is that the resulting estimates are more stable (less sensitive to small perturbation of the BP centers). From a statistical point of view, this procedure results in a smoothed version of the curvature, discarding very local bends of the DNA, but retaining larger scale bending.

Since the curvature is constrained to be positive, the curvature functions $c_j(t)$ do not lie in a linear space. Given that functional data analyses typically hinge on linear space methods, we converted the curvature trajectories into elements of a linear space by applying the transformation $x \mapsto \log(\delta + x)$, where $\delta > 0$ is a fixed constant (see below), thus defining the δ -linearized curvature by

$$d_t(j) = \log(\delta + c_t(j)), \quad \text{for all } t, j. \tag{2.2}$$

The constant δ prevents erratic fluctuations of $d_j(t)$ when $c_j(t)$ approaches zero. If δ is too small, the functions $d_j(\cdot)$ will present very large spikes, and if δ is too large, $d_j(\cdot)$ will be essentially constant. Based on an exploratory analysis, we set $\delta = 10^{-3}$, which struck a balance between the two extremes (this choice of δ is of course specific to our dataset). Figure S11 of the supplementary material illustrates the role of δ .

Since the δ -linearized curvatures are discretely sampled versions of smooth curves, we transformed each function $j \mapsto d_t(j)$ into a smooth curve $\tau \mapsto Y_t(\tau)$, $\tau \in [0, 1]$, by smoothing the scatterplot

$$\left(\frac{j-1}{158}, d_t(j) \right)_{j=1, \dots, 158} \tag{2.3}$$

for each fixed t . This was done using a basis expansion (Ramsay and Silverman 2005) with 80 periodic cubic B-splines (King, Nguyen, and Ionides 2016), respecting the nature of the data as closed loops. Our choice of 80 basis functions came from the combination of considerations on the postulated degrees of freedom of the curvature (which should be fewer than the number of base-pairs), computational considerations, and graphical goodness-of-fit assessment. It is of course specific to our dataset. We also conducted the analysis presented in the rest of the article with 40 and 60 basis functions, and the results were similar to those obtained with 80 basis functions. We note that exploratory

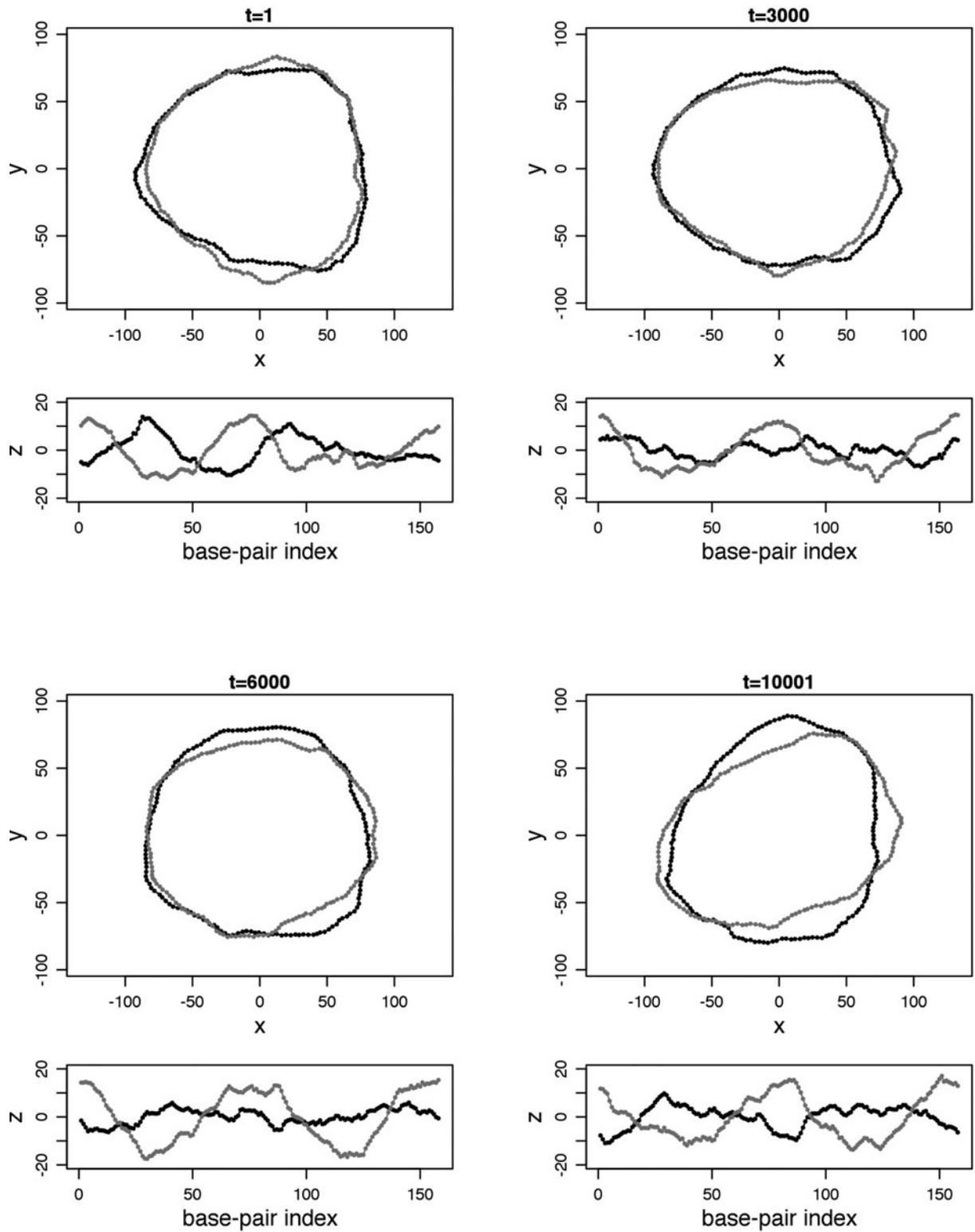


Figure 1. The DNA minicircles for various timepoints t (CAP in black, TATA in gray). The top plot of the four subfigures contains the projection of the DNA minicircles onto the XY plane, and the plot below shows their projection on the Z -axis. The units of the X, Y, Z axes are in angstroms ($1 \text{ angstrom} = 10^{-10} \text{ m}$).

plots revealed no need for further penalization in the smoothing of the functions d_t . Figure S12 of the supplementary material illustrates the smoothing process.

Exploratory analysis of the functional time series $\{Y_t : t = 1, \dots, 10,001\}$ revealed that the series exhibited a nonstationary behavior coupled with a long memory behavior, as is natural with the movement of a rigid body. By contrast, the time

differenced curves, $X_t := Y_{t+1} - Y_t$, exhibited a weakly dependent stationary behavior, and so we focused on the series $\{X_t : t = 1, \dots, 10,000\}$ as our object of study. The model implicitly assumed is, therefore,

$$Y_{t+1}(\tau) = Y_t(\tau) + X_t(\tau), \quad \tau \in [0, 1],$$

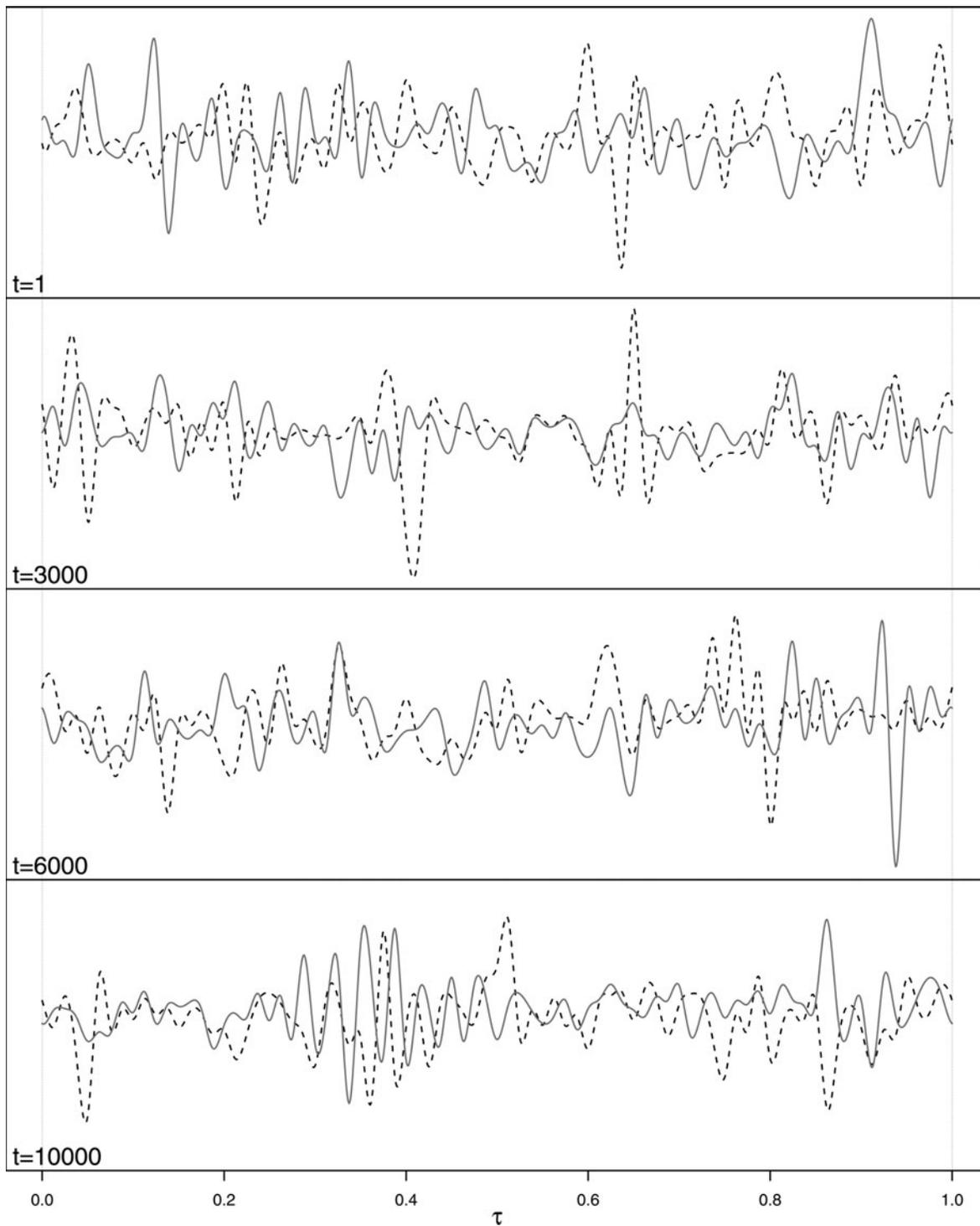


Figure 2. Plot of the innovation process X_t of the linearized curvatures for CAP (dashed black curve) and TATA (solid gray curve) for various timepoints t .

where $\{Y_t\}$ is the δ -linearized curvature of the DNA minicircle, and $\{X_t\}$ is the stationary innovation process governing the time-increments of $\{Y_t\}$. Applying all these steps to the CAP minicircles, respectively, TATA minicircles, we get two functional time series, X_t^1 , respectively, X_t^2 (see Figure 2).

2.2. Estimation of the Dynamics

Let us start by introducing some useful notation. For $u, v, f \in L^2([0, 1], \mathbb{C})$, the usual inner product will be denoted by

$\langle u, v \rangle = \int_0^1 u(\tau)\overline{v(\tau)}d\tau$, where $\bar{\alpha}$ is the complex conjugate of $\alpha \in \mathbb{C}$, with corresponding norm $\|u\| = \sqrt{\langle u, u \rangle}$. The tensor product $u \otimes v$ will be the linear operator on $L^2([0, 1], \mathbb{C})$ defined as $(u \otimes v)f = \langle f, v \rangle u$.

The (second-order) dynamics of a stationary FTS $\{X_t : t \in \mathbb{Z}\}$ are encoded in the collection of lag- t autocovariance operators

$$\mathcal{R}_t = \mathbb{E}[(X_t - \mu) \otimes (X_0 - \mu)], \quad t \in \mathbb{Z}, \quad (2.4)$$

where $\mu = \mathbb{E}X_t$. Inferences on these dynamics could therefore be carried out by means of corresponding inferences on the individual autocovariance operators. However, such an approach encounters significant roadblocks, since the theory of the estimation of the lag- t autocovariance operators without structural assumptions (such as that of linearity, made in Mas & Pumo (in Ferraty and Romain (2011))) is largely unexplored (with the exception of Zhang and Shao (2015) who recently developed tests only for comparing the lag-0 autocovariance operator without structural assumptions). We choose to take a different approach, and estimate the complete dynamics through a frequency domain approach, following the paradigm recently introduced by Panaretos and Tavakoli (2013a).

In the frequency domain approach, the objects of interest are no longer the lag- t autocovariance operators, but their Fourier transform,

$$\mathcal{F}_\omega = (2\pi)^{-1} \sum_{t \in \mathbb{Z}} \exp(-i\omega t) \mathcal{R}_t, \quad \omega \in [-\pi, \pi], \quad (2.5)$$

where $i \in \mathbb{C}$ is the imaginary number, $i^2 = -1$. The object $\{\mathcal{F}_\omega : \omega \in [-\pi, \pi]\}$ is called the *spectral density operator*, and is well defined under suitable summability conditions on the autocovariance operators (see Panaretos and Tavakoli 2013a, Proposition 2.1). At each frequency ω , the operator \mathcal{F}_ω is a linear operator on $L^2([0, 1], \mathbb{C})$, associated with a spectral density kernel f_ω , a complex valued surface $[0, 1]^2 \ni (\tau, \sigma) \mapsto f_\omega(\tau, \sigma) \in \mathbb{C}$. Various properties of the spectral density operator are given in Panaretos and Tavakoli (2013a) and Tavakoli (2014); we remark that $f_{-\omega} = \overline{f_\omega}$, and we therefore restrict our interest to the range $\omega \in [0, \pi]$. Intuitively, the spectral density operator is the generalization of the spectral density matrix encountered in multivariate time series (Brillinger 2001; Priestley 2001) to the functional setting, and yields an analysis of variance decomposition of the variance of the FTS, given by the inversion formula $\mathcal{R}_t = \int_{-\pi}^{\pi} \exp(i\omega t) \mathcal{F}_\omega d\omega$, $t \in \mathbb{Z}$.

The reason we decide to carry out inference on the dynamics of an FTS via a frequency domain approach—and not through the autocovariance operators—is that the spectral density operator is inextricably linked with the so-called *Cramér–Karhunen–Loève* decomposition of the functional time series (and the associated harmonic principal component analysis), which elucidates its complete dynamical properties, separating the temporal, functional, and stochastic components (Panaretos and Tavakoli 2013b; Tavakoli 2014; Hörmann, Kidziński, and Hallin 2015). The Cramér–Karhunen–Loève decomposition also elucidates why spectral density operators are equally important across all frequencies (under no prior assumption concerning the dynamics of the time series). Furthermore, working in the frequency domain has a crucial whitening effect: the sample spectral density operator is asymptotically independent at distinct frequencies. In contrast to these properties, it is not clear what relative importance one should give to each lag- t autocovariance operators in the assessment of dynamics of an FTS, and the sample lag- t autocovariance operators are asymptotically correlated (Mas & Pumo in Ferraty and Romain 2011), which would make a multiplicity correction approach (see Sections 3 and 4) more involved, and potentially considerably less powerful.

Notice also that for $f, g \in L^2[0, 1]$, $\langle f, \mathcal{F}_\omega f \rangle$ is the power spectrum of the one-dimensional time series $\langle X_t, f \rangle$, while $\langle f, \mathcal{F}_\omega g \rangle$ is the cross spectrum of $\langle X_t, f \rangle$ and $\langle X_t, g \rangle$. In this sense, the spectral density operator can be used to probe any continuous linear functional of the dynamics of the process (a fact that will be exploited in Section 4).

The estimation of the spectral density operator is carried out by first computing the discrete Fourier transforms of the FTS, and then smoothing its empirical covariance (called the periodogram operator), with a kernel function $W(\cdot)$ of bandwidth B_T (see Panaretos and Tavakoli 2013a for details). The first step can be done using the Fast Fourier Transform, whose calculation is most efficient when the length of the series is highly composite.

As usual, the bandwidth parameter B_T needs to satisfy the conditions $B_T \rightarrow 0$ and $TB_T \rightarrow \infty$ as $T \rightarrow \infty$, for the asymptotic results to hold. The choice of B_T governs also the bias/variance trade-off for the estimation of the spectral density operator, similarly to nonparametric regression. Inferential procedures can then be constructed using Theorem 1 and Panaretos and Tavakoli (2013a, Theorem 3.7), which gives a central limit theorem for the estimated spectral density operator at distinct frequencies. For finite samples, it is crucial to notice that the central limit theorem effect emerges because the spectral density estimator at a given frequency ω is obtained by weighted averaging of $(2m + 1)$ approximately independent summands, where $m = \lfloor TB_T/2\pi \rfloor$. Following Brillinger (2001, p. 252), the equivalent number of independent pieces of information used to estimate the spectral density estimator at frequency ω can be defined as

$$n(\omega, m, \kappa) = m/\kappa^2, \quad (2.6)$$

where $\kappa^2 = \int_{\mathbb{R}} W^2(x) dx$. The order of $n(\omega, m, \kappa)$ plays therefore a role similar to the number of iid summands when applying the classical central limit theorem, and should be taken into account before making any inferential statements based on asymptotics.

In the present setup, we choose $B_T = 0.158$, a choice that corresponds to the heuristic that $B_T \sim O(T^{-1/5})$ asymptotically, to minimize the mean square error (Brillinger 2001, p. 251). We choose $W(x)$ to be the Epanechnikov kernel (e.g., Wand and Jones 1995), $W(x) = \frac{3}{4}(1 - x^2)$ if $|x| < 1$, and zero otherwise. These choices imply $m = 252$ and $n(\omega, m, \kappa) = 420$. We denote by $\mathcal{F}_\omega^{a,(T)}$ the estimated spectral density operators, or *sample spectral density operators*, of X_t^a , for $a = 1, 2$.

The graphical representation of the sample spectral density operators is not straightforward: for each frequency $\omega \in [0, \pi]$, one has an operator on $L^2([0, 1], \mathbb{C})$: the corresponding trace norm is shown in Figure 3 for CAP and TATA. The (modulus of the) sample spectral density kernels $f_\omega^{a,(T)}$, $a = 1, 2$, associated with the sample spectral density operators, are depicted in Figure S13 of the supplementary material. We notice that most of the variance is distributed along the high-frequency end of the spectral density operator, meaning that the series X_t consists mainly of high-frequency oscillations, and that the low-frequency oscillations contained in the linearized curvature series Y_t mostly cancel out when taking its time differences to form X_t (this is not surprising considering the transfer function

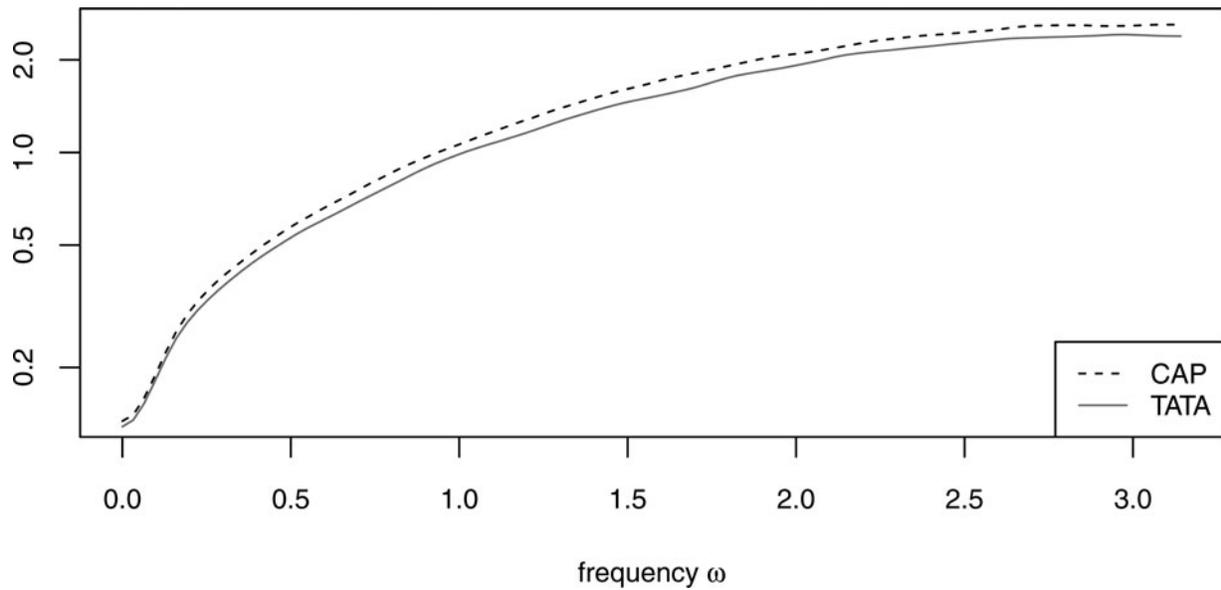


Figure 3. Plot of the trace of the spectral density operator. Notice that the spectral density operator of CAP has consistently a larger magnitude than the spectral density operator of TATA.

induced by a time differencing). Figure S13 of the supplementary material illustrates that most of the nuclear norm of the spectral density operator is near the diagonal of the sample spectral density kernels, and then falls off sharply as one moves away from the diagonal. The interpretation of this is that the series X_t has strong local interactions, in the sense that $X_t(\tau)$ and $X_0(\sigma)$ are interacting strongly for $|\tau - \sigma|$ small, say $|\tau - \sigma| < \varepsilon'$, and much more weakly for $|\tau - \sigma| > \varepsilon'$. This reflects the fact that the series X_t is locally smooth, but globally quite rough, as can be seen in Figure 2.

2.3. Comparison of the Spectral Density Operators

Though the traces of the sample spectral density operators of the CAP and TATA series appear different, and though small differences between their sample spectral density kernels are visible, it is not a priori clear whether these differences are indeed statistically significant. The next section sets out to address this issue, by introducing inferential tools for comparing two spectral density operators on a grid of frequencies, and localizing any detected differences at the level of frequency.

3. Comparing and Localizing Spectral Differences by Frequency

Comparing the second-order dynamics of the functional time series X_t^1 and X_t^2 can now be formalized as testing the equality of their spectral density operators. More precisely, if H_ω : “ $\mathcal{F}_\omega^1 = \mathcal{F}_\omega^2$,” for $\omega \in [0, \pi]$, we wish to test

$$\bigcap_{\omega \in [0, \pi]} H_\omega \quad \text{against} \quad “H_\omega \text{ fails for some } \omega \in [0, \pi].”$$

We will take a multiple testing approach to this problem by first constructing a test for each H_ω , marginally, and then enforcing a multiplicity correction.

3.1. Comparing the Spectral Density Operator at a Fixed Frequency

To test H_ω for a fixed $\omega \in [0, \pi]$, we construct a test inspired by the class of tests introduced by Panaretos, Kraus, and Maddocks (2010) for testing the equality of covariance operators in iid collections of Gaussian random functions. The key idea is that, in light of the central limit theorem of Panaretos and Tavakoli (2013a, Theorem 3.7), the estimated (or sample) spectral density operator $\mathcal{F}_\omega^{(T)}$ can be roughly seen as the empirical covariance operator of a sample of $T B_T / 2\pi$ approximately iid replications of the discrete Fourier transform of $\{X_t\}$ at frequency ω . Under this heuristic, one can construct a test statistic in the spirit of Panaretos, Kraus, and Maddocks (2010), by considering the Hilbert–Schmidt norm of the difference between the sample operators, $\mathcal{F}_\omega^{1,(T)} - \mathcal{F}_\omega^{2,(T)}$, restricted on a subspace of Hilbert–Schmidt space of dimension K . The choice of this subspace is such that it retains the bulk of the Hilbert–Schmidt norm of the difference, subject to its dimension being K . When H_ω is valid, this is achieved by projecting onto the (random) subspace generated by the tensor products of the first K estimated eigenfunctions of $\mathcal{F}_\omega^1 = \mathcal{F}_\omega^2$. Letting $(\tilde{\mu}_i^\omega, \tilde{\varphi}_i^\omega)_{i=1}^\infty$ be the eigenvalue/eigenfunction pairs of the pooled spectral density operator $(\mathcal{F}_\omega^{1,(T)} + \mathcal{F}_\omega^{2,(T)})/2$, and denoting by $D_\omega^{(T)} = \sqrt{T B_T} (\mathcal{F}_\omega^{1,(T)} - \mathcal{F}_\omega^{2,(T)})$ the rescaled difference between the two sample spectral density operators at ω , we consider the statistic

$$\tilde{\Delta}_K^{(T)}(\omega) = \sum_{i,j=1}^K \frac{\left| \left\langle D_\omega^{(T)} \tilde{\varphi}_j^\omega, \tilde{\varphi}_i^\omega \right\rangle \right|^2}{(1 + \mathbf{1}_{\{0,\pi\}}(\omega)) 4\pi \kappa^2 \tilde{\mu}_i^\omega \tilde{\mu}_j^\omega}, \quad (3.1)$$

where $\kappa^2 = \int_{\mathbb{R}} W^2(x) dx$ is the L^2 -norm of the weight function $W(x)$. Each component in the sum comprising the test statistic measures the squared norm of the component of the difference $D_\omega^{(T)}$ that lies in the subspace spanned by $\tilde{\varphi}_i \otimes \tilde{\varphi}_j$, renormalized according to its asymptotic variance. The indicator in the denominator is a correction term for the frequencies $\omega \in \{0, \pi\}$,

at which the sample spectral density operator has an increased variance.

The heuristics leading to the construction of the test statistic can indeed be made rigorous, and the following Theorem establishes that the asymptotic distribution of our test statistic under the null hypothesis is chi-squared. Its proof is given in Section C of the supplementary material.

Theorem 1. Let K_1, \dots, K_J be fixed nonnegative integers, $\omega_1, \dots, \omega_J \in [0, \pi]$ be a fixed number of distinct frequencies. Assume that Conditions B.1 of the supplementary material hold, and that $B_T \rightarrow 0$ and $TB_T \rightarrow \infty$ as $T \rightarrow \infty$. Furthermore, assume that for each ω_j , the first K_j eigenvalues of the spectral density operator \mathcal{F}_{ω_j} are all distinct, and nonnegative. Then, under the null hypothesis $H_0 = \cap_{j=1}^J H_{\omega_j}$, the test statistics $\tilde{\Delta}_{K_j}^{(T)}(\omega_j)$, $j = 1, \dots, J$, converge in distribution to independent random variables $\Delta_{K_j}(\omega_j)$, where

$$\Delta_{K_j}(\omega_j) \sim \begin{cases} \chi_{K_j(K_j+1)/2}^2 & \text{if } \omega_j \in \{0, \pi\}, \\ \chi_{K_j^2}^2 & \text{otherwise.} \end{cases} \quad (3.2)$$

The truncation level K —which is assumed to be fixed in the asymptotic framework of Theorem 1—represents a regularization parameter whose choice governs a bias/variance trade-off reflected in Type I and Type II error probabilities. Small values of K will guarantee the preservation of the nominal level of the test under the null, but are likely to incur “bias-related” Type II error under the alternative, when differences in the two operators are to be found in dimensions higher than K . A more aggressive choice of a large K can result to instabilities due to the ill-posedness of the problem, resulting in Type I and Type II errors alike. In principle, the choice of truncation level K should be dependent on the corresponding frequency ω_j at which the test is carried out. This is because a spectral density operator \mathcal{F}_ω may exhibit a different rate of spectral decay as ω varies. At a fixed sample size, the optimal value of $K_j(T)$ depends in a complicated manner on the eigenvalues, and the gaps between the eigenvalues, of the spectral density operator at the corresponding frequency. Though a theoretical investigation along this avenue would be of interest (using results from, e.g., Fremdt et al. 2014), it is beyond the scope of our article. In practice, a frequency-dependent choice may be made by choosing $K = K(\omega)$, which optimizes a model selection-type criterion. We use a fit/penalty trade-off criterion that is inspired by the pseudo-AIC (Akaike information criterion) criterion of Yao, Müller, and Wang (2005), and its two-sample generalization introduced by Panaretos, Kraus, and Maddocks (2010):

$$\text{AIC}(K, \omega) = \text{GOF}(K, \omega) + \text{PEN}_1(K, \omega) + \text{PEN}_2(K, \omega), \quad (3.3)$$

where $\text{GOF}(K, \omega)$ is a goodness-of-fit criterion, and $\text{PEN}_a(K, \omega)$, $a = 1, 2$, penalize for overfitting of the spectral densities $\mathcal{F}_\omega^{a,(T)}$, $a = 1, 2$. We propose taking

$$\text{GOF}(K, \omega) = \sum_{k=K+1}^{N_b} \langle (\mathcal{F}_\omega^{1,(T)} - \mathcal{F}_\omega^{2,(T)}) \tilde{\varphi}_k^\omega, \tilde{\varphi}_k^\omega \rangle \quad (3.4)$$

and

$$\text{PEN}_a(K, \omega) = \left(\sum_{j=1}^{N_b} \hat{\lambda}_j \right) \sum_{j=1}^{N_b} \frac{\langle \mathcal{F}_\omega^{a,(T)}(K) \hat{\varphi}_j^{a,\omega}, \hat{\varphi}_j^{a,\omega} \rangle}{n(\omega, m, \kappa) \hat{\lambda}_j^{a,\omega}}, \quad a = 1, 2, \quad (3.5)$$

where $(\hat{\mu}_j^{a,\omega}, \hat{\varphi}_j^{a,\omega})$ denotes the j th eigenvalue/eigenvector pair of $\mathcal{F}_\omega^{a,(T)}$, $a = 1, 2$; $j = 1, 2, \dots$, and

$$\mathcal{F}_\omega^{a,(T)}(K) = \sum_{k_1, k_2=1}^K \langle \mathcal{F}_\omega^{a,(T)} \tilde{\varphi}_{k_1}^\omega, \tilde{\varphi}_{k_2}^\omega \rangle \tilde{\varphi}_{k_1}^\omega \otimes \tilde{\varphi}_{k_2}^\omega, \quad a = 1, 2,$$

is the projection of the sample spectral density operator onto the first K eigenspaces of the pooled sample spectral density operator $(\mathcal{F}_\omega^{1,(T)} + \mathcal{F}_\omega^{2,(T)})/2$. The constant $n(\omega, m, \kappa)$ is defined in (2.6), and depends only on ω , $m = \lfloor TB_T/2\pi \rfloor$, and κ . The intuition behind this criterion is that it corresponds to the AIC criterion of Panaretos, Kraus, and Maddocks (2010, sec. 3.3) had we observed $n(\omega, m, \kappa)$ iid complex curves from a random function with covariance \mathcal{F}_ω^a , for $a = 1, 2$. Even though these curves are not observed in our context, the choice of $n(\omega, m, \kappa)$ reflects the number of independent pieces of information used to construct our estimate $\mathcal{F}_\omega^{a,(T)}$.

We also propose a variant of the AIC criterion, by using the following penalty in lieu of (3.5),

$$\text{PEN}_a^*(K, \omega) = \left(\sum_{j=1}^{N_b} \hat{\lambda}_j \right) \sum_{j=1}^{N_b} \frac{\langle \mathcal{F}_\omega^{a,(T)}(K) \hat{\varphi}_j^{a,\omega}, \hat{\varphi}_j^{a,\omega} \rangle}{n(\omega, m, \kappa) \sqrt{\hat{\lambda}_j^{a,\omega} \hat{\gamma}_j^{a,\omega}}}, \quad a = 1, 2, \quad (3.6)$$

where $\hat{\gamma}_1^{a,\omega} = \hat{\lambda}_1^{a,\omega} - \hat{\lambda}_2^{a,\omega}$ and $\hat{\gamma}_l^{a,\omega} = \min\{\hat{\lambda}_{l-1}^{a,\omega} - \hat{\lambda}_l^{a,\omega}, \hat{\lambda}_l^{a,\omega} - \hat{\lambda}_{l+1}^{a,\omega}\}$, $l = 2, \dots$, and $a = 1, 2$. The corresponding pseudo-AIC criterion is

$$\text{AIC}^*(K, \omega) = \text{GOF}(K, \omega) + \text{PEN}_1^*(K, \omega) + \text{PEN}_2^*(K, \omega). \quad (3.7)$$

The difference between AIC and AIC* is that the second criterion takes into account the difficulty of estimating the eigenstructure of the pooled spectral density operator, in addition to penalizing for the relative roughness of the pooled spectral density operator in comparison to $\mathcal{F}_\omega^{1,(T)}$ and $\mathcal{F}_\omega^{2,(T)}$, respectively (see Bosq 2000, Lemma 4.3). We also note that both criteria are invariant to scaling of the sample spectral density operator.

Numerical simulations (see Section D of the supplementary material), conducted with the automatic choice of $K = K(\omega)$ using either AIC and AIC*, suggest that our testing procedure respects the level of the test in a variety of settings, unless for very low sample size when $K(\omega)$ is chosen by AIC (see Table S1 of the supplementary material). It should be noted that our procedure is slightly conservative in general, due to the multiple testing approach taken (Benjamini and Hochberg 1995). The numerical simulations suggest that $K(\omega)$ should be selected by means of AIC in settings where the eigenvalues of the spectral density operator decay steeply, and that AIC* should be used if the eigenvalues of the spectral density operator decay slowly (we

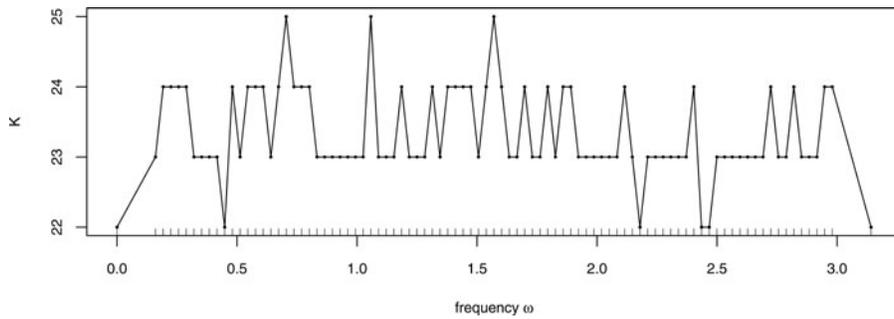


Figure 4. Truncation levels $K(\omega)$ as chosen by the AIC* criterion (3.7). The small ticks on the horizontal axis represent the grid of frequencies Γ for which the test is computed.

note that the power using either AIC criterion is not necessarily higher than that obtained using a prechosen value of truncation level K ; however it is not clear how to choose K a priori). Since our DNA minicircle data correspond to the second case, we shall use AIC* to choose $K(\omega)$ in the remainder of the article.

3.2. Localization of Differences on the Frequencies

Recall that H_ω denotes the null hypothesis $\mathcal{F}_\omega^1 = \mathcal{F}_\omega^2$. To test the global null hypothesis $H_G := \bigcap_{\omega \in [0, \pi]} H_\omega$, we will first obtain marginal p -values for each of the null hypotheses H_ω , $\omega \in \Gamma$, where $\Gamma := \{\omega_1, \dots, \omega_J\} \subset [0, \pi]$ is a grid of frequencies, and then adjust the p -values to account for multiplicity effects. The p -values will be based on the asymptotic distribution of test statistic $\tilde{\Delta}_K^{(T)}(\omega)$, given by Theorem 1.

The results of applying the automatic truncation level rule (3.7) to our DNA minicircle dataset are shown in Figure 4. We notice that the selected values of $K(\omega)$ vary between 21 and 25. The corresponding (approximate) p -values are

$$p_j = \mathbb{P} \left(\tilde{\Delta}_{K(\omega_j)}^{(T)}(\omega_j) < \chi_{v(\omega_j)}^2 \right), \quad j = 1, \dots, J,$$

where $v(\omega_j) = K(\omega_j)[K(\omega_j) - 1]/2$ if $\omega_j \in \{0, \pi\}$, and $v(\omega_j) = K(\omega_j)^2$ otherwise. The choice of the grid of frequencies at which the p -values are computed should be guided by a priori knowledge of the nature of the alternative hypothesis—provided there is any such knowledge—see Section 3.3. In our case, we chose a grid of 81 frequencies, which is shown in Figure 4.

Adjusting the p -values for multiplicities can be done to control the false discovery rate (FDR; see Benjamini and Hochberg 1995), the expected value of the proportion of false rejections

among all rejections. In terms of the FDR, since the p -values p_i and p_j are dependent for $|\omega_i - \omega_j| < 0.15$, but approximately independent for $|\omega_i - \omega_j| > 0.15$, we are in the context of dependence in finite blocks, and the original Benjamini–Hochberg (BH) algorithm for controlling the FDR is valid (Storey, Taylor, and Siegmund 2004). We show the adjusted p -values using the BH procedure in Figure 5. We notice that the two spectral density operators appear to be highly significantly different at all frequencies. We also conducted numerical simulations to assess the performance and validity of our procedure in finite sample; these suggest that the BH procedure controls the Type I error for small sample sizes (see Section D of the supplementary material).

Since these adjusted p -values are remarkably small (they range from 10^{-30} to 10^{-140}), and our simulations in Section D of the supplementary material demonstrate the validity of the test procedure, we perform a “sanity check” to make sure that these p -values are not the result of any remaining transient (nonstationary) behavior: we consider what p -values would be produced when comparing the dynamics of two time-separated stretches of the same FTS. Using the procedure described in this section, we compare the spectral density operators of the two FTSs $\{X_t^{\text{CAP}} : t = 1, \dots, 1000\}$ and $\{X_t^{\text{CAP}} : t = 9001, \dots, 10,000\}$, which are assumed to be approximately independent under weak dependence. The null hypothesis should be valid, of course, provided stationarity is not violated. The adjusted p -values, shown in Figure S14 of the supplementary material, suggest that the spectral density operator of CAP is indeed stable in time (across all frequencies), in accordance with the hypothesis of stationarity made earlier, and supporting the soundness of the adjusted p -values presented in Figure 5. The same conclusions hold for TATA (Figure S14 of the supplementary material).

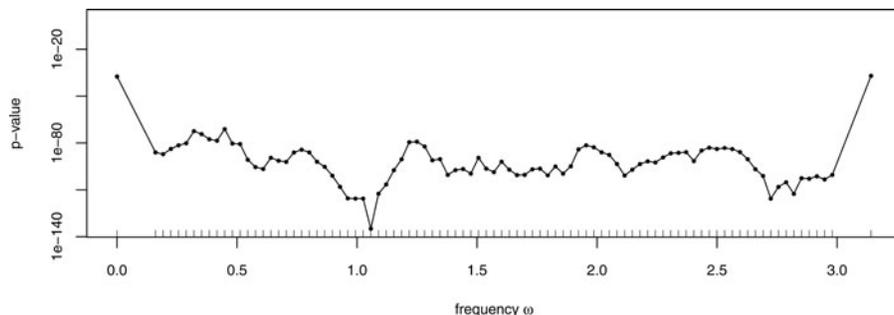


Figure 5. Adjusted p -values (using the BH procedure) for testing the equality of the spectral density operator of CAP and TATA, with the truncation level $K(\omega)$ automatically chosen at each frequency ω with the AIC* criterion (3.7). The small ticks on the horizontal axis represent the grid of frequencies Γ for which the test is computed.

3.3. Choice of the Discretization Grid Γ

The choice of the grid Γ (and therefore J) is related to finer knowledge of potential departures in the direction of the alternative against which we wish to test the global null $H_G = \bigcap_{\omega \in [0, \pi]} H_\omega$:

Power for global differences between the two spectral density operators: If we believe that the true difference between the two spectral density operators is going to be on a large interval of $[0, \pi]$, J should be small, so that the power of the test is not lost because of multiple comparisons. If $\Gamma \subset [0, \pi]$ is chosen such that $|\omega_i - \omega_j| \geq 2B_T$, for all $\omega_i \neq \omega_j \in \Gamma$, then the p -values are approximately independent, and one could control the familywise error rate via Hochberg's procedure (Dudoit, Shaffer, and Boldrick 2003), if this is what is desired.

Power for narrow banded differences between the two spectral density operators: If we believe that the true difference is in a very narrow band of the spectral density operator, for example, H_ω is false only for $|\omega - \omega'| < \varepsilon$, with $\omega' \in [0, \pi]$ and $\varepsilon > 0$ small, then Γ should be chosen to be a dense grid over $[0, \pi]$. The largest gap between any two consecutive frequencies in Γ will indicate approximately smallest bandsize ε for which the test would be able to detect departures from the global null H_G .

Frequencies near $\{0, \pi\}$: Although we expect $\tilde{\Delta}_K^{(T)}(\omega_j)$ to follow, for large T , approximately a $\chi_{K^2}^2$ distribution for any $\omega_j \notin \{0, \pi\}$, the approximation might not hold for frequencies ω_j very close to $\{0, \pi\}$. This happens because the asymptotic distribution of $\tilde{\Delta}_K^{(T)}(\omega)$ is $\chi_{K(K+1)/2}^2$ for $\omega \in \{0, \pi\}$, but $\chi_{K^2}^2$ for $\omega \in (0, \pi)$, and because $\tilde{\Delta}_K^{(T)}(\omega)$ is continuous in ω . Therefore, for ω_j close to $\{0, \pi\}$, the approximate distribution of $\tilde{\Delta}_K^{(T)}(\omega_j)$ is a mixture of $\chi_{K(K+1)/2}^2$ and $\chi_{K^2}^2$ random variables, with unknown mixture proportion. We therefore recommend that all the frequencies $\omega \in \Gamma$, with $\omega \notin \{0, \pi\}$, be at least at distance B_T of the frequencies $\{0, \pi\}$.

Understanding how the power varies with the number of frequencies J at which the spectral density operators are compared is important, since there may be situations in which no prior knowledge on the spectral density operators is available. As it turns out, there is a heuristic upper limit to the gridsize because

1. The test used is continuous in ω for constant K , that is, $\omega \rightarrow \Delta_K^{(T)}(\omega)$ is continuous, conditionally on $(X_t)_{t=1}^T$.
2. Computationally, it is more efficient to compute the sample spectral operators, and hence $\tilde{\Delta}_K^{(T)}(\omega)$, on the Fourier frequencies $\Gamma_* = \{2\pi s/T : s = 0, \dots, T/2\}$, or on a subset of them.

If the multiple correction is done using the false discovery rate (FDR), numerical simulations (Section D.2) suggest that conditionally on the observed functional time series, the FDR-adjusted p -values (the q -values) are approximately stable as the grid becomes as dense as Γ_* (see Figure S9 of the supplementary material). Concerning the power of the test, it seems that though a very sparse grid may yield more power in some situations, in other situations, slightly denser grids yield considerably

more power. However, choosing the densest grid Γ^* results generally in a loss of power, due to multiplicity corrections and the continuity of the test in the frequencies (see Figure S10 of the supplementary material).

4. Localizing Differences in Frequency and Along Curvelength

We now wish to qualify the difference between CAP and TATA dynamics at a finer level: we wish to first detect the specific frequencies at which CAP and TATA curves differ (the *significant frequencies*), and then localize the region on the minicircles (curves), within each significant frequency, where these differences occur. This serial framework is quite naturally amenable to recent selective multiple testing methodology proposed by Benjamini and Bogomolov (2014). Note that if one wishes to search for dynamical differences between the two series at a frequency ω_0 and attributable to the covariation between two regions $[\tau_1, \tau_2]$ and $[\tau_3, \tau_4]$ along the curves, it suffices to consider hypotheses comparing linear contrasts $\langle \mathcal{F}_{\omega_0}^1 g, h \rangle$ and $\langle \mathcal{F}_{\omega_0}^2 g, h \rangle$, for $g, h \in L^2[0, 1]$ two contrast functions concentrated on $[\tau_1, \tau_2]$ and $[\tau_3, \tau_4]$, respectively.

The contrasts we choose to employ are the periodic B-spline basis functions used to represent the curves (King, Nguyen, and Ionides 2016). Consequently, we base our procedure on the differences in the (i, j) th basis coefficient between the spectral density operator of CAP and TATA, at a given frequency ω . Let us denote by \mathbf{f}_ω^a , respectively, $\mathbf{f}_\omega^{a, (T)}$, the 80×80 coefficient matrices with respect to the periodic B-spline basis (King et al. 2010) of the true spectral density operator, respectively, the sample spectral density operator, at frequency ω , for the time series X_t^a , $a = 1, 2$. We shall call \mathbf{f}_ω^a the projected spectral density operator, and $\mathbf{f}_\omega^{a, (T)}$ the projected *sample* spectral density operator, and denote by $\mathbf{f}_\omega(i, j)$ the (i, j) th entry of the matrix \mathbf{f}_ω . The local null hypotheses we wish to test for are of the form

$$H_\omega(i, j) : \mathbf{f}_\omega^1(i, j) = \mathbf{f}_\omega^2(i, j), \quad i, j = 1, \dots, 80; \omega \in [0, \pi].$$

By symmetry of the projected spectral density operator, we restrict ourselves to the indices $1 \leq i \leq j \leq 80$. We point out that this approach is different from a classical multivariate approach, as discussed in Remark 2.

For each frequency ω and each $1 \leq i \leq j \leq 80$, assuming $\mathbf{f}_\omega(i, i)\mathbf{f}_\omega(j, j) \neq 0$, we can use the projected sample spectral density operator to construct a p -value $p(\omega; i, j)$ for the null hypothesis $H_\omega(i, j)$, as described in Section E of the supplementary material. The p -values are only computed on a subgrid $\Gamma = \{\omega_1, \dots, \omega_L\} \subset [0, \pi]$, which is chosen such that $|\omega_i - \omega_j| \geq 2B_T$, so that the p -values across different ω_j 's are approximately independent (see the discussion in Section 3.3).

We choose to select significant frequencies and localize the differences between CAP and TATA in a way that controls the expected average of the false discovery proportion over the significant frequencies (Benjamini and Bogomolov 2014). To make this statement precise, let $\mathbf{p}_l = \{p(\omega_l; i, j) : 1 \leq i \leq j \leq 80\}$ be the set of p -values at frequency ω_l , and $\mathbf{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_L\}$ be

the set of all p -values over the grid Γ . Let $S(\mathbf{P})$ be the selection procedure for the significant frequencies, based on all the p -values \mathbf{P} , that is, $S(\mathbf{P}) \subset \Gamma$, and $|S(\mathbf{P})|$ denote the number of significant frequencies. Let $FDP(\omega) = V(\omega)/R(\omega)$ be the false discovery proportion at frequency ω , where $V(\omega)$ denotes the (unknown) number of wrong rejections within frequency ω , and $R(\omega)$ denotes the total number of rejections at frequency ω . The

error criterion we will seek to control is

$$\mathbb{E} \left[\sum_{l \in S(\mathbf{P})} FDP(\omega_l) / \max \{ |S(\mathbf{P})|, 1 \} \right]. \quad (4.1)$$

Notice that if the selection procedure S is carried out without relying on the data, (4.1) simplifies to $\sum_{l \in S} FDR(\omega_l) / |S|$, the

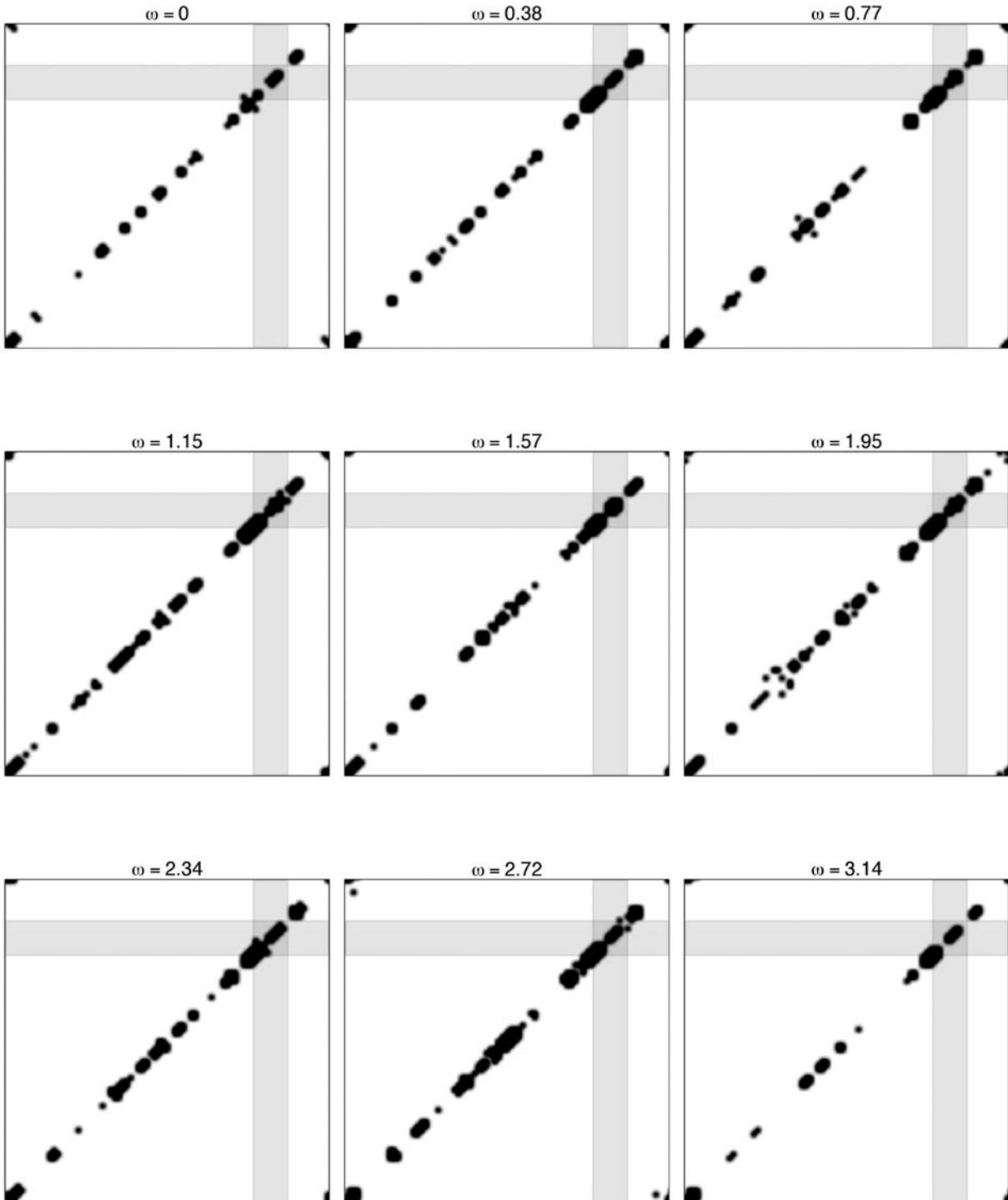


Figure 6. The plots show the regions on the minicircles, for each frequency, where the spectral density operator of CAP and TATA are overall significantly different at a 1% level (with respect to the error criterion (4.1)). Each plot represents the regions of differences (in black) between the spectral density kernels of the two minicircles. The two gray vertical and horizontal bands correspond to the region where the base-pair sequences of the two DNA minicircles are different.

average FDR over the selected frequencies, where $\text{FDR}(\omega_l) = \mathbb{E}[V(\omega_l)/R(\omega_l)]$.

To select the significant frequencies and select the null hypotheses to reject within each significant frequency while controlling the expected average FDP (4.1) at the level α , we use the following procedure (see Theorem 1 and sec. 5, Benjamini and Bogomolov 2014):

1. Adjust, within each frequency ω_l , the p -values \mathbf{p}_l for the control of the FDR, and denote the result by \mathbf{q}_l , also called q -values (see Remark 1).
2. Select the significant frequencies S by applying the BH procedure to the set of minimum q -values $\{\min \mathbf{q}_1, \min \mathbf{q}_2, \dots, \min \mathbf{q}_L\}$.
3. Within each significant frequency $\omega_l, l \in S$, reject the null hypotheses whose corresponding q -value are smaller than $|S|\alpha/L$.

In addition to controlling the error criterion (4.1), this procedure has also the additional property that it controls the FDR at the level of the frequencies. Notice however that the frequencies selected as significantly different by this method need not be exactly the same as the ones selected by the method of Section 3.

The result of applying this procedure to our minicircle data with $\alpha = 0.05$, and on the grid of frequencies $\Gamma = \{0, 0.38, 0.77, 1.15, 1.57, 1.95, 2.34, 2.72, 3.14\}$ is shown in Figure 6 in the form of *zero-one plots*, which exhibit graphically the regions where the spectral density operator of CAP and TATA differ significantly at a 1% level. We notice first that all the tested frequencies are significant, which is not surprising since the frequency tests (Section 3) suggested that the null hypothesis H_ω for each fixed frequency was confidently rejected. We also see that the rejected hypotheses are mostly situated on the diagonal of the spectral density operator, that is, the rejected nulls are mostly of the form $H_\omega(i, j)$ with $|i - j|$ small. This signifies that the differences in the dynamics of CAP and TATA curves are primarily due to local interactions (between $X_t(\tau)$ and $X_0(\sigma)$ for $|\tau - \sigma|$ small) differ. This is not surprising since we have already seen (in Section 2.2) that most of the covariation of the minicircles stems from their local interactions. An interesting observation is that the detected differences between the spectral density operators of the two minicircles do not exclusively reside in the region where their BP sequence is different (see Table 1), but extend to other regions of the minicircles.

Remark 1 (p -Value adjustment within each frequency). Since the p -values $\mathbf{p}_l = \{p(\omega_l; i, j) : 1 \leq i \leq j \leq 80\}$ are correlated with a nontrivial correlation structure (see, e.g., (E.2) of the supplementary material), we cannot use the BH procedure to control the FDR, nor more recent procedures (which require, e.g., dependence in finite blocks, see Storey, Taylor, and Siegmund 2004; Schwartzman, Dougherty, and Taylor 2008). We therefore use the conservative version of FDR, which works under arbitrary dependence structure of the p -values (Benjamini and Yekutieli 2001, Theorem 1.3) to obtain the q -values \mathbf{q}_l . Nevertheless, numerical simulation that we carried out to assess the validity of our procedure suggested that the BH procedure seems to control the FDR within each frequency. Further work along this line would be of interest.

Remark 2 (Differences with multivariate analysis). Although the idea of comparing at the level of basis coefficients seems like a multivariate approach, it differs from it in that the choice of the basis functions will influence the qualitative conclusions that can be drawn from the analysis. Our choice of a periodic B-spline basis allows one to distinguish differences between CAP and TATA that are very localized on the minicircles. Another choice could be that of a wavelet basis, which would allow one to detect differences between CAP and TATA across multiple scales. The choice of the basis is therefore intimately related to the directions (in function space) in which the test is most powerful.

5. Concluding Remarks

We have introduced a method for comparing the dynamics of two functional time series (FTS) at a hierarchy of levels, through a frequency domain approach. Our method was illustrated through a case study in molecular biophysics, and specifically aimed at detecting sequence-dependent effects on the molecular dynamics of DNA at persistence length.

Our procedure is based on a test for comparing the spectral density operators of two FTSs at fixed frequencies. As a first step, this test can be used in combination with multiple testing procedures to detect differences between the spectral density operators of FTSs, and enables localizing at which frequencies the differences occur, while controlling an overall error measure. As a second step, one can compare the spectral density operators of two FTSs jointly in frequencies and along the curvelength, by first localizing differences in the frequencies, and then identifying their differences along the curvelength, within each frequency, while controlling the average false discovery rate over the selected frequencies. We conducted numerical simulations to assess the strength of our method in finite samples, and its robustness to “adversarial” setups.

Our case study indicates that the dynamics of the two DNA minicircles we studied (CAP and TATA) seem to be globally significantly different, across every fluctuation frequency, at least at the given MD simulation. A finer investigation of their differences—along the curvelength—show signs of an interesting phenomenon: the dynamics of CAP and TATA, though being mainly local, seem to be not only limited to the region where their base-pairs are different, but to extend to other regions of the DNA minicircles. This suggests that a local change of base-pairs might induce a global change of dynamics through a “propagation effect” along the DNA minicircle (see also Kim et al. 2013). Though the effect appears to affect the *entire* DNA minicircle in our case study, it is not clear whether the effect’s intensity is fading with distance, and if it would have affected only part of dynamics had the DNA minicircle been longer. Another important point to mention is that though the difference between CAP and TATA were strongly statistically significant (the largest adjusted p -value for testing the equality of their spectral density operators was smaller than 10^{-20}), a bare eye examination of the data does yield any hints on differences in their dynamics.

Our method relies on assumptions on the stationarity and weak dependence of the underlying FTSs. To validate our findings, we performed a “sanity check” by searching for differences

in the spectral density operators of each FTSs, estimated using the first 1000 and last 1000 timepoints of each FTS. No difference was detected, suggesting that the assumption of stationarity is not violated. Furthermore, exploratory analysis of the FTSs revealed that the weak dependence assumption is acceptable.

To our knowledge, this work is the first attempt of comparing the entire second-order dynamics of FTSs, and localizing their differences across frequencies, and within frequencies along the curvelength. This also appears to be the first time that a functional data analysis has been employed to study the coarse-grained molecular dynamics of DNA, and indeed that significant dynamical sequence-dependent differences have been statistically quantified in a functional setup. Moreover, the method we propose does not hinge on any linearity or Gaussian assumption on the underlying FTSs, but on stationarity and moment-type weak dependence assumptions. This is particularly fitting in the case of DNA where the scaling limit models are far from clear, and potentially non-Gaussian; and indeed, given the rigid body nature of DNA, existence of moments of all orders naturally leads to moment-type mixing. A drawback of our method, due to our model-free and frequency domain approach, is that interpretation is not straightforward (though the Cramér–Karhunen–Loève decomposition, Panaretos and Tavakoli 2013b, can be used to this aim). Potential extensions of our work could be in the direction of tests for stationarity (similar to our “sanity check”), development of bootstrap version of our tests to take into account the local dependency (in frequencies) of the sample spectral density operator, or incorporation of the theory of excursion sets of Gaussian processes for the localization of the difference along curvelength.

A note of caution is that, even though the differences detected between the spectral density operators are statistically significant, we do not claim over-arching conclusions on the nature and effect size of these differences: this would require further and finer analyses (e.g., Freddolino et al. 2006; Sanbonmatsu and Tung 2007) and of course may also require more than one set of MD trajectories. Still, our preliminary results can hopefully serve as a starting point for further work on the DNA context, and as a case study illustrating the scope and potential of functional time series inference in biophysical modeling.

Supplementary Materials

The supplement contains the molecular dynamics simulation protocol (Section A), the technical assumptions of our main results (Section B), the proof of Theorem 1 (Section C), and numerical simulations (Section D).

R Package The R code implementing the methods of this article is implemented in the R package `ftsspec` available at <https://cran.r-project.org/web/packages/ftsspec/index.html>. The data generated and simulations studies are available at <https://www.repository.cam.ac.uk/handle/1810/253695>.

Acknowledgment

The authors gratefully acknowledge Dr. Jonathan S. Mitchell and Professor John H. Maddocks for kindly sharing their Molecular Dynamics dataset with them and for providing them with the description of the MD protocol, as given in the supplementary material. The authors are thankful to

the members of Professor Maddocks’ group for several enlightening discussions on the mechanics of DNA, and to Professors A.C. Davison and J.A.D. Aston for helpful comments.

Funding

ST was partially supported by the EPSRC grant EP/K021672/2, and the research was partly carried out while ST was a Ph.D. student at the Institute of Mathematics, Ecole Polytechnique Fédérale de Lausanne. This research was supported by a European Research Council (ERC) Starting Grant Award to Victor M. Panaretos, and was primarily carried out while Shahin Tavakoli was a PhD student at the Institute of Mathematics, Ecole Polytechnique Fédérale de Lausanne. ST was subsequently partially supported by the EPSRC grant EP/K021672/2.

References

- Amzallag, A., Vaillant, C., Jacob, M., Unser, M., Bednar, J., Kahn, J. D., Dubochet, J., Stasiak, A., and Maddocks, J. H. (2006), “3D Reconstruction and Comparison of Shapes of DNA Minicircles Observed by Cryo-electron Microscopy,” *Nucleic Acids Research*, 34, e125–e125. [1022]
- Antoniadis, A., Paparoditis, E., and Sapatinas, T. (2006), “A Functional Wavelet-kernel Approach for Time Series Prediction,” *Journal of the Royal Statistical Society, Series B*, 68, 837–857. [1020]
- Aston, J. A. D., and Kirch, C. (2012a), “Evaluating Stationarity via Change-point Alternatives With Applications to fMRI Data,” *Annals of Applied Statistics*, 6, 1906–1948. [1020,1021]
- (2012b), “Detecting and Estimating Changes in Dependent Functional Data,” *Journal of Multivariate Analysis*, 109, 204–220. Available at <http://dx.doi.org/10.1016/j.jmva.2012.03.006>. [1021]
- Aue, A., Hörmann, S., Horváth, L., and Hušková, M. (2014), “Dependent Functional Linear Models With Applications to Monitoring Structural Change,” *Statistica Sinica*, 24, 1043–1073. [1021]
- Aue, A., Norinho, D. D., and Hörmann, S. (2015), “On the Prediction of Stationary Functional Time Series,” *Journal of the American Statistical Association*, 110, 378–392. [1021]
- Benjamini, Y., and Bogomolov, M. (2014), “Selective Inference on Multiple Families of Hypotheses,” *Journal of the Royal Statistical Society, Series B*, 76, 297–318. [1030,1032]
- Benjamini, Y., and Hochberg, Y. (1995), “Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing,” *Journal of the Royal Statistical Society, Series B*, 57, 289–300. [1028,1029]
- Benjamini, Y., and Yekutieli, D. (2001), “The Control of the False Discovery Rate in Multiple Testing Under Dependency,” *The Annals of Statistics*, 29, 1165–1188. [1032]
- Benko, M., Härdle, W., and Kneip, A. (2009), “Common Functional Principal Components,” *Annals of Statistics*, 37, 1–34. [1020]
- Boente, G., Rodriguez, D., and Sued, M. (2011), “Testing the Equality of Covariance Operators,” *Recent Advances in Functional Data Analysis and Related Topics*, 49–53. [1021]
- (2014), A Test for the Equality of Covariance Operators, *ArXiv e-prints*, 1404–7080. [1020,1021]
- Bosq, D. (2000), *Linear Processes in Function Spaces*, New York: Springer. [1021,1028]
- Brillinger, D. R. (2001), *Time Series: Data Analysis and Theory (Classics in Applied Mathematics, classics edition)*, Philadelphia, PA: SIAM. [1026]
- Cuevas, A., Febrero, M., and Fraiman, R. (2004), “An ANOVA Test for Functional Data,” *Computational Statistics & Data Analysis*, 47, 111–122. Available at <http://www.sciencedirect.com/science/article/pii/S016794730300269X>. [1020]
- Curuksu, J., Kannan, S., and Zacharias, M. (2014), “Molecular Dynamics and Advanced Sampling Simulations of Nucleic Acids,” in *Handbook of Computational Chemistry*, ed. J. Leszczynski, Netherlands: Springer, pp. 1155–1173. [1022]
- Dans, P. D., Faustino, I., Battistini, F., Zakrzewska, K., Lavery, R., and Orozco, M. (2014), “Unraveling the Sequence-dependent Polymorphic Behavior of d(CpG) Steps in B-DNA,” *Nucleic Acids*

- Research*, 42, 11304–11320. Available at <http://nar.oxfordjournals.org/content/42/18/11304.abstract>. [1023]
- Dans, P. D., Pérez, A., Faustino, I., Lavery, R., and Orozco, M. (2012), “Exploring Polymorphisms in B-DNA Helical Conformations,” *Nucleic Acids Research*, 40, 10668–10678. Available at <http://nar.oxfordjournals.org/content/40/21/10668.abstract>. [1023]
- Dauxois, J., Pousse, A., and Romain, Y. (1982), “Asymptotic Theory for the Principal Component Analysis of a Vector Random Function: Some Applications to Statistical Inference,” *Journal of Multivariate Analysis*, 12, 136–154. [1020]
- Dryden, I. L., Kume, A., Le, H., Wood, A. T. A., and Laughton, C. A. (2002), “Size-and-shape Analysis of DNA Molecular Dynamics Simulations,” in *Statistics of Large Datasets: Functional and Image Data, Bioinformatics and Data Mining*, eds. R. G. Aykroyd, K. V. Mardia, and P. McDonnell, Department of Statistics, University of Leeds, pp. 23–26. [1022]
- Dudoit, S., Shaffer, J. P., and Boldrick, J. C. (2003), “Multiple Hypothesis Testing in Microarray Experiments,” *Statistical Science*, 18, 71–103. [1030]
- Fan, J., and Lin, S.-K. (1998), “Test of Significance When Data Are Curves,” *Journal of the American Statistical Association*, 93, 1007–1021. Available at <http://www.jstor.org/stable/2669845>. [1020]
- Ferraty, F. (2011), *Recent Advances in Functional Data Analysis and Related Topics (Contributions to Statistics)*, Berlin: Physica-Verlag. [1020]
- Ferraty, F., and Romain, Y. (2011), *The Oxford Handbook of Functional Data Analysis*, New York: Oxford University Press. [1021,1026]
- Ferraty, F., and Vieu, P. (2006), *Nonparametric Functional Data Analysis: Theory and Practice*, New York: Springer. [1020]
- Freddolino, P. L., Arkhipov, A. S., Larson, S. B., McPherson, A., and Schulten, K. (2006), “Molecular Dynamics Simulations of the Complete Satellite Tobacco Mosaic Virus,” *Structure*, 14, 437–449. Available at <http://www.sciencedirect.com/science/article/pii/S0969212606000608>. [1022,1033]
- Fremdt, S., Horváth, L., Kokoszka, P., and Steinebach, J. G. (2014), “Functional Data Analysis With Increasing Number of Projections,” *Journal of Multivariate Analysis*, 124, 313–332. Available at <http://dx.doi.org/10.1016/j.jmva.2013.11.009>. [1021,1028]
- Fremdt, S., Steinebach, J. G., Horváth, L., and Kokoszka, P. (2013), “Testing the Equality of Covariance Operators in Functional Samples,” *Scandinavian Journal of Statistics*, 40, 138–152. [1021]
- Garcia, H. G., Grayson, P., Han, L., Inamdar, M., Kondev, J., Nelson, P. C., Phillips, R., Widom, J., and Wiggins, P. A. (2007), “Biological Consequences of Tightly Bent DNA: The Other Life of a Macromolecular Celebrity,” *Biopolymers*, 85, 115–130. [1021]
- Glaser, R. (2012), *Biophysics: An Introduction* (2nd Ed.), Berlin Heidelberg: Springer-Verlag. [1021]
- Gonzalez, O., Petkevičiūtė, D., and Maddocks, J. H. (2013), “A Sequence-dependent Rigid-base Model of DNA,” *The Journal of Chemical Physics*, 138, 055102. [1021]
- Grenander, U. (1981), *Abstract Inference (Wiley Series in Probability and Mathematical Statistics, Vol. IX)*, New York: Wiley, pp. 526. [1020]
- Hadjipantelis, P. Z., Aston, J. A. D., Müller, H. G., and Evans, J. P. (2015), “Unifying Amplitude and Phase Analysis: A Compositional Data Approach to Functional Multivariate Mixed-Effects Modeling of Mandarin Chinese,” *Journal of the American Statistical Association*, 110, 545–559. [1020]
- Hall, P., and Hosseini-Nasab, M. (2006), “On Properties of Functional Principal Components Analysis,” *Journal of the Royal Statistical Society, Series B*, 68, 109–126. [1020]
- Hörmann, S., and Kidziński, Ł. (2012), “A Note on Estimation in Hilbertian Linear Models,” *Scandinavian Journal of Statistics*, 42, 43–62. [1021]
- Hörmann, S., Kidziński, Ł., and Hallin, M. (2015), “Dynamic Functional Principal Components,” *Journal of the Royal Statistical Society, Series B*, 77, 319–348. [1021,1026]
- Hörmann, S., Kidziński, Ł., and Kokoszka, P. (2015), “Estimation in Functional Lagged Regression,” *Journal of Time Series Analysis*, 36, 541–561. [1021]
- Hörmann, S., and Kokoszka, P. (2010), “Weakly Dependent Functional Data,” *Annals of Statistics*, 38, 1845–1884. [1021]
- Horváth, L., Hušková, M., and Rice, G. (2013), “Test of Independence for Functional Data,” *Journal of Multivariate Analysis*, 117, 100–119. Available at <http://www.sciencedirect.com/science/article/pii/S0047259X13000195>. [1021]
- Horváth, L., and Kokoszka, P. (2012), *Inference for Functional Data With Applications (Springer Series in Statistics)*, New York: Springer. [1020,1021]
- Horváth, L., Kokoszka, P., and Reeder, R. (2013), “Estimation of the Mean of Functional Time Series and a Two-Sample Problem,” *Journal of the Royal Statistical Society, Series B*, 75, 103–122. [1020,1021]
- Horváth, L., Kokoszka, P., and Rice, G. (2014), “Testing Stationarity of Functional Time Series,” *Journal of Econometrics*, 179, 66–82. Available at <http://www.sciencedirect.com/science/article/pii/S0304407613002327>. [1021]
- Horváth, L., and Rice, G. (2015a), “Testing for Independence Between Functional Time Series,” *Journal of Econometrics*, 189, 371–382. [1021]
- (2015b), “Testing Equality of Means When the Observations are From Functional Time Series,” *Journal of Time Series Analysis*, 36, 84–108. [1021]
- Horváth, L., Rice, G., and Whipple, S. (2014), “Adaptive Bandwidth Selection in the Long Run Covariance Estimator of Functional Time Series,” *Computational Statistics & Data Analysis*, 100, 676–693. [1021]
- Kahn, J. D., and Crothers, D. M. (1992), “Protein-induced Bending and DNA Cyclization,” *Proceedings of the National Academy of Sciences*, 89, 6343–6347. [1022]
- Kim, S., Broströmer, E., Xing, D., Jin, J., Chong, S., Ge, H., Wang, S., Gu, C., Yang, L., Gao, Y. Q., Su, X.-D., Sun, Y., and Xie, X. S. (2013), “Probing Allostery Through DNA,” *Science*, 339, 816–819. [1032]
- King, A. A., Nguyen, D., and Ionides, E. L. (2016), “Statistical Inference for Partially Observed Markov Processes via the R Package pomp,” *Journal of Statistical Software*, 69, 1–43. [1023,1030]
- Kokoszka, P., and Reimherr, M. (2013), “Asymptotic Normality of the Principal Components of Functional Time Series,” *Stochastic Process and their Applications*, 123, 1546–1562. [1021]
- Kokoszka, P., and Young, G. (2016), “KPSS Test for Functional Time Series,” *Statistics*, 50, 957–973. [1021]
- Kraus, D., and Panaretos, V. M. (2012), “Dispersion Operators and Resistant Second-order Functional Data Analysis,” *Biometrika*, 99, 813–832. [1020,1021]
- Lankas, F., Lavery, R., and Maddocks, J. H. (2006), “Kinking Occurs During Molecular Dynamics Simulations of Small DNA Minicircles,” *Structure*, 14, 1527–1534. [1022]
- Lavery, R., Maddocks, J. H., Pasi, M., and Zakrzewska, K. (2014), “Analyzing Ion Distributions Around DNA,” *Nucleic Acids Research*, 42, 8138–8149. Available at <http://nar.oxfordjournals.org/content/42/12/8138.abstract>. [1023]
- Leach, A. R. (2001), *Molecular Modelling: Principles and Applications* (2nd ed.), Harlow, Essex: Prentice Hall. [1022]
- Mas, A. (2002), “Weak Convergence for the Covariance Operators of a Hilbertian Linear Process,” *Stochastic Processes and Their Applications*, 99, 117–135. [1021]
- Mas, A. (2007), “Testing for the Mean of Random Curves: A Penalization Approach,” *Statistical Inference for Stochastic Processes*, 10, 147–163. [1020]
- Mas, A., and Menneteau, L. (2003), “Perturbation Approach Applied to the Asymptotic Study of Random Operators,” in *High Dimensional Probability III* (Vol. 55), eds. J. Hoffmann-Jørgensen, J. A. Wellner, and M. B. Marcus, Basel: Birkhäuser, pp. 127–133. [1020]
- Mastroianni, A. J., Sivak, D. A., Geissler, P. L., and Alivisatos, A. P. (2009), “Probing the Conformational Distributions of Subpersistence Length DNA,” *Biophysical Journal*, 97, 1408–1417. [1021]
- Mitchell, J. S., and Harris, S. A. (2013), “Thermodynamics of Writhe in DNA Minicircles From Molecular Dynamics Simulations,” *Physical Review Letters*, 110, 148105. [1022]
- Mitchell, J. S., Laughton, C. A., and Harris, S. A. (2011), “Atomistic Simulations Reveal Bubbles, Kinks and Wrinkles in Supercoiled DNA,” *Nucleic Acids Research*, 39, 3928–3938. [1022]
- Panaretos, V. M., Kraus, D., and Maddocks, J. H. (2010), “Second-Order Comparison of Gaussian Random Functions and the Geometry of

- DNA Minicircles,” *Journal of the American Statistical Association: Theory & Methods*, 105, 670–682. [1020,1021,1027,1028]
- Panaretos, V. M., and Tavakoli, S. (2013a), “Fourier Analysis of Stationary Time Series in Function Space,” *Annals of Statistics*, 41, 568–603. [1021,1026,1027]
- (2013b), “Cramér–Karhunen–Loève Representation and Harmonic Principal Component Analysis of Functional Time Series,” *Stochastic Process and their Applications*, 123, 2779–2807. [1021,1026,1033]
- Paparoditis, E., Sapatinas, T. (2014), “Bootstrap-Based K-Sample Testing For Functional Data,” *ArXiv e-prints*, 1409–4317. [1020,1021]
- Pasi, M., Maddocks, J. H., Beveridge, D., Bishop, T. C., Case, D. A., Cheatham, T., Dans, P. D., Jayaram, B., Lankas, F., Laughton, C., Mitchell, J., Osman, R., Orozco, M., Pérez, A., Petkevičič, D., Spackova, N., Sponer, J., Zakrzewska, K., and Lavery, R. (2014), “ABC: A Systematic Microsecond Molecular Dynamics Study of Tetranucleotide Sequence Effects in B-DNA,” *Nucleic Acids Research*, 42, 12272–12283. [1022]
- Pérez, A., Luque, F. J., and Orozco, M. (2011), “Frontiers in Molecular Dynamics Simulations of DNA,” *Accounts of Chemical Research*, 45, 196–205. [1022]
- Peters, J. P., and Maher, L. J. (2010), “DNA Curvature and Flexibility in Vitro and in Vivo,” *Quarterly Reviews of Biophysics*, 43, 23–63. [1021]
- Pigoli, D., Aston, J. A. D., Dryden, I. L., and Secchi, P. (2014), “Distances and Inference for Covariance Operators,” *Biometrika*, 101, 409–422. [1021]
- Prévost, C., Takahashi, M., and Lavery, R. (2009), “Deforming DNA: From Physics to Biology,” *ChemPhysChem*, 10, 1399–1404. [1021]
- Priestley, M. B. (2001), “Spectral Analysis and Time Series” (Vol. I and II), in *Probability and Mathematical Statistics*, Cambridge, MA: Academic Press. [1026]
- Ramsay, J. O., and Silverman, B. W. (2002), *Applied Functional Data Analysis: Methods and Case Studies (Springer Series in Statistics)*, New York: Springer. [1020]
- (2005), *Functional Data Analysis* (2nd ed.), New York: Springer. [1020,1023]
- Rothmund, P. W. K. (2006), “Folding DNA to Create Nanoscale Shapes and Patterns,” *Nature*, 440, 297–302. [1021]
- Sambriski, E. J., Schwartz, D. C., and De Pablo, J. J. (2009), “A Mesoscale Model of DNA and its Renaturation,” *Biophysical Journal*, 96, 1675–1690. [1021]
- Sanbonmatsu, K. Y., and Tung, C.-S. (2007), “High Performance Computing in Biology: Multimillion Atom Simulations of Nanoscale Systems,” *Journal of Structural Biology, Advances in Molecular Dynamics Simulations*, 157, 470–480. Available at <http://www.sciencedirect.com/science/article/pii/S104784770600308X>. [1022,1033]
- Sangalli, L. M., Secchi, P., Vantini, S., and Veneziani, A. (2009), “Efficient Estimation of Three-dimensional Curves and Their Derivatives by Free-knot Regression Splines, Applied to the Analysis of Inner Carotid Artery Centrelines,” *Journal of the Royal Statistical Society, Series C*, 58, 285–306. [1023]
- Schwartzman, A., Dougherty, R. F., and Taylor, J. E. (2008), “False Discovery Rate Analysis of Brain Diffusion Direction Maps,” *The Annals of Applied Statistics*, 2, 153–175. [1032]
- Seeman, N. C. (2005), “DNA Enables Nanoscale Control of the Structure of Matter,” *Quarterly Reviews of Biophysics*, 38, 363–371. [1021]
- Shore, D., and Baldwin, R. L. (1983), “Energetics of DNA Twisting: I. Relation Between Twist and Cyclization Probability,” *Journal of Molecular Biology*, 170, 957–981. [1022]
- Shore, D., Langowski, J., and Baldwin, R. L. (1981), “DNA Flexibility Studied by Covalent Closure of Short Fragments Into Circles,” *Proceedings of the National Academy of Sciences*, 78, 4833–4837. [1022]
- Storey, J. D., Taylor, J. E., and Siegmund, D. (2004), “Strong Control, Conservative Point Estimation and Simultaneous Conservative Consistency of False Discovery Rates: A Unified Approach,” *Journal of the Royal Statistical Society, Series B*, 66, 187–205. [1029,1032]
- Tavakoli, S. (2014), “Fourier Analysis of Functional Time Series, With Applications to DNA Dynamics,” Ph.D. dissertation, EPFL. Available at <http://dx.doi.org/10.5075/epfl-thesis-6320>. [1026]
- Walter, J., Gonzalez, O., and Maddocks, J. H. (2010), “On the Stochastic Modeling of Rigid Body Systems With Application to Polymer Dynamics,” *Multiscale Modeling & Simulation*, 8, 1018–1053. [1021]
- Wand, M. P., and Jones, M. C. (1995), *Kernel Smoothing of Monographs on Statistics and Applied Probability* (Vol. 60), London: Chapman & Hall. [1026]
- Wang, J.-L., Chiou, J.-M., and Mueller, H.-G. (2016), “Functional Data Analysis,” in *Annual Review of Statistics and Its Application* (Vol. 3), eds. S. E. Fienberg and S. M. Stigler, pp. 257–295. [1020]
- Yao, F., Müller, H. G., and Wang, J. L. (2005), “Functional Linear Regression Analysis for Longitudinal Data,” *Annals of Statistics*, 33, 2873–2903. [1028]
- Zhang, X., and Shao, X. (2015), “Two Sample Inference for the Second-Order Property of Temporally Dependent Functional Data,” *Bernoulli*, 21, 909–929. [1021,1026]
- Zhang, X., Shao, X., Hayhoe, K., and Wuebbles, D. J. (2011), “Testing the Structural Stability of Temporally Dependent Functional Observations and Application to Climate Projections,” *Electronic Journal of Statistics*, 5, 1765–1796. [1021]