

# Order Statistics Correlation Coefficient as a Novel Association Measurement With Applications to Biosignal Analysis

Weichao Xu, *Member, IEEE*, Chunqi Chang, *Member, IEEE*, Y. S. Hung, *Senior Member, IEEE*, S. K. Kwan, and Peter Chin Wan Fung

**Abstract**—In this paper, we propose a novel correlation coefficient based on order statistics and rearrangement inequality. The proposed coefficient represents a compromise between the Pearson's linear coefficient and the two rank-based coefficients, namely Spearman's rho and Kendall's tau. Theoretical derivations show that our coefficient possesses the same basic properties as the three classical coefficients. Experimental studies based on four models and six biosignals show that our coefficient performs better than the two rank-based coefficients when measuring linear associations; whereas it is well able to detect monotone nonlinear associations like the two rank-based coefficients. Extensive statistical analyses also suggest that our new coefficient has superior anti-noise robustness, small biasedness, high sensitivity to changes in association, accurate time-delay detection ability, fast computational speed, and robustness under monotone nonlinear transformations.

**Index Terms**—Atrial fibrillation, atrial flutter, concomitant, Kendall's tau, nonlinear association measure, order statistics, Pearson's coefficient, rearrangement inequality, Spearman's rho.

## I. INTRODUCTION

THERE has been great interest in measuring the association between two time series, with application in many areas including biosignal analysis. Association can be considered as the strength of relationship between the two time series. A measure of association should be large and positive if there is a high probability that large (small) values of one time series are associated with large (small) values of another. On the other hand, if the direction is inverse, namely, large (small) values of one time series occur in conjunction with small (large) values of another time series, the measure should be large and negative [1]. A multitude of methods have been used in the literature of biosignal processing for many years to measure the association between two time series. Among these measures the Pearson's linear correlation coefficient [2]–[4], Spearman's

rank-based coefficient (*Spearman's rho*) [5] and Kendall's concordance coefficient (*Kendall's tau*) [5] are perhaps the most widely used [6]. The linear correlation coefficient is appropriate mainly for indicating linear associations [1], while the other two measures are invariant under linear or nonlinear increasing monotone transformations [5]. Some authors have employed the average amount of mutual information (AAMI) to measure the association between two biosignals [7], [8], others used nonlinear regression coefficient (denoted by  $h^2$ ) [9] or contingency table based methods (Cramer  $V$ ) for such purpose [10].

There are many advantages and disadvantages to the measures mentioned before. Linear correlation coefficient is very fast, however, it will yield misleading results if nonlinearity is involved in the system [10]. On the other hand, the two rank correlation coefficients, Spearman's rho and Kendall's tau, are not as powerful and as fast as Pearson's coefficient when measuring linear associations between biosignals; nevertheless they are independent of increasing nonlinear transformations which makes them suitable for many nonlinear cases [1], [5], [6]. Despite their robustness for nonlinear association measurements, the values of  $h^2$  and  $V$  are between 0 and 1, meaning their inability of distinguishing positive associations from negative associations. Furthermore, the computational load of AAMI and  $h^2$  are rather heavy, which makes them inappropriate in cases when high computational speed is mandatory.

To overcome the problems of the existing measures of association in different *a priori* unknown situations, we propose a novel measure called order statistic correlation coefficient which possesses the following advantages: 1) it can discriminate positive associations from negative associations (admissible range  $[-1, 1]$ ); 2) its time complexity is of order  $O(N \log(N))$ , a little slower than Pearson's coefficient but much faster than Kendall's tau, the nonlinear regression coefficient  $h^2$ , and AAMI; 3) it has small biasedness in both linear and nonlinear scenarios; 4) it is sensitive to changes of degree of association; and 5) it possesses certain robustness under increasing nonlinear transformations.

In Section II, we will give the definition and properties of our new order statistic coefficient as well as the other three classical indices of association. Section III depicts the models and performance evaluation strategy we use in this study. In Section IV, we present the simulated signals and the associated results of four models used in our investigation. Section V is devoted to discussions and interpretations of our new method. Finally, in Section VI, we draw our conclusions on the novel order statistics correlation coefficient.

Manuscript received March 15, 2006; revised March 4, 2007. This work was supported in part by the Hong Kong Innovation and Technology Commission under Funding ITS/109/02, by Hong Kong RGC under Grant N\_HKU703/03, and the University of Hong Kong under Small Project Funding 200507176052. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Leslie Collins.

W. Xu, C. Chang, Y. S. Hung, and S. K. Kwan are with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong (e-mail: wxu@eee.hku.hk; cqchang@eee.hku.hk; yshung@eee.hku.hk; skkwan@eee.hku.hk).

P. C. W. Fung is with the Department of Medicine, The University of Hong Kong, Hong Kong (e-mail: hrspfcw@hkucc.hku.hk).

Digital Object Identifier 10.1109/TSP.2007.899374

## II. ORDER STATISTICS CORRELATION COEFFICIENT

### A. Definition and Properties

Let  $(x_i, y_i)$ ,  $i = 1, \dots, N$ , be two time series of length  $N$ . Rearranging pairwise the two time series with respect to the magnitudes of  $x$ , we get two new series denoted by  $(x_{(i)}, y_{[i]})$ , where  $x_{(1)} \leq \dots \leq x_{(N)}$  are called the *order statistics* of  $x$  and  $y_{[1]}, \dots, y_{[N]}$  the associated *concomitants* [11]–[13]. Reversing the roles of  $x$  and  $y$ , we also define the order statistics of  $y$  and the corresponding concomitants which are denoted by  $y_{(1)}, \dots, y_{(N)}$  and  $x_{[1]}, \dots, x_{[N]}$ , respectively. As proposed by Xu *et al.* [14], the order statistics correlation coefficient can be defined as

$$r_X(x, y) \triangleq \frac{\sum_{i=1}^N (x_{(i)} - x_{(N-i+1)}) y_{[i]}}{\sum_{i=1}^N (x_{(i)} - x_{(N-i+1)}) y_{(i)}}. \quad (1)$$

**Theorem 1:** The order statistics correlation coefficient has the basic properties of a correlation coefficient, as follows:

- 1)  $-1 \leq r_X \leq 1$ ;
- 2)  $r_X(x, y)$  attains  $+1(-1)$  when  $x$  and  $y$  are in strict increasing (decreasing) relationship;
- 3)  $r_X(x', y') = r_X(x, y)$  for  $x' = k_x x + \text{const}_x$  and  $y' = k_y y + \text{const}_y$ , where  $k_x > 0$  and  $k_y > 0$ ;
- 4) if  $x$  and  $y$  are mutually independent and each is independent identically distributed (IID), the expectation  $E\{r_X(x, y)\} = 0$  when  $N \rightarrow \infty$ .

*Proof:*

- 1) According to the rearrangement inequality [15], it follows that:

$$\sum_{i=1}^N x_{(N-i+1)} y_{(i)} \leq \sum_{i=1}^N x_{(i)} y_{[i]} \leq \sum_{i=1}^N x_{(i)} y_{(i)} \quad (2)$$

and

$$\sum_{i=1}^N x_{(N-i+1)} y_{(i)} \leq \sum_{i=1}^N x_{(N-i+1)} y_{[i]} \leq \sum_{i=1}^N x_{(i)} y_{(i)}. \quad (3)$$

Subtracting (3) by (4) and dividing the difference by  $\sum (x_{(i)} - x_{(N-i+1)}) y_{(i)}$ , we have  $-1 \leq r_X \leq 1$ , hence the result.

- 2) Assume  $y_i = \phi(x_i)$ ,  $i = 1, \dots, N$ . If  $\phi(\bullet)$  is a strict increasing function, we have  $y_{[i]} = y_{(i)}$  for all  $i$ . Substituting this into (1), we have  $r_X = 1$ ; and similarly  $r_X = -1$  if  $\phi(\bullet)$  is a strict decreasing function.
- 3) Substituting  $x'$  and  $y'$  into (1), we have

$$\begin{aligned} r_X(x', y') &= \frac{\sum (x'_{(i)} - x'_{(N-i+1)}) y'_{[i]}}{\sum (x'_{(i)} - x'_{(N-i+1)}) y'_{(i)}} \\ &= \frac{k_x k_y \sum (x_{(i)} - x_{(N-i+1)}) y_{[i]} + A}{k_x k_y \sum (x_{(i)} - x_{(N-i+1)}) y_{(i)} + A} \end{aligned} \quad (4)$$

where  $A = k_x \text{const}_y \sum [x_{(i)} - x_{(N-i+1)}] = 0$ . Hence, we have  $r_X(x', y') = r_X(x, y)$ .

- 4) Denote the numerator and denominator of (1) by  $U$  and  $V$ , respectively. An application of the Delta method [16] yields

$$E(r_X) = \frac{E(U)}{E(V)} + O(N^{-1}). \quad (5)$$

In order to prove that  $E(r_X) = 0$  (with  $N$  large), it is adequate to show that  $E(U)$  is null when the assumptions are satisfied. Imposing the independence assumption of  $x$  and  $y$ , we have

$$E(U) = \sum [E(x_{(i)}) - E(x_{(N-i+1)})] E(y_{[i]}). \quad (6)$$

It is known [6] that when  $y$  is IID, the probability density function (PDF) of  $y_{[i]}$  is

$$g_{[i]}(y) = \int_{-\infty}^{+\infty} f(y|x) f_{(i)}(x) dx \quad (7)$$

where  $g_{[i]}(y)$  denotes the PDF of  $y_{[i]}$ ,  $f(y|x)$  the conditional PDF of  $y$  given  $x$  and  $f_{(i)}(x)$  the PDF of  $x_{(i)}$ . The conditional PDF  $f(y|x)$  degenerates to  $f(y)$  if  $x$  and  $y$  are independent. Then, we have

$$\begin{aligned} E(y_{[i]}) &= \int y g_{[i]}(y) dy \\ &= \int y \int f(y|x) f_{(i)}(x) dx dy \\ &= \int y f(y) dy \int f_{(i)}(x) dx \\ &= \int y f(y) dy \\ &= E(y). \end{aligned} \quad (8)$$

Substituting (8) into (6), we have  $E(U) = 0$  and, thus,  $E(r_X) = 0$  to the order of  $O(N^{-1})$ .

### B. Estimation of Correlation Coefficient in Normal Case

It will be shown in the following theorem that for samples from a bivariate normal population with correlation coefficient  $\rho$ ,  $r_X$  is an asymptotically unbiased estimation of  $\rho$ .

**Theorem 2:** If  $(x_i, y_i)$ ,  $i = 1, \dots, N$  is a pair of IID time series from a bivariate normal distribution with correlation coefficient  $\rho$ , then  $\lim_{N \rightarrow \infty} E\{r_X(x, y)\} = \rho$ .

*Proof:* Without loss of generality, we assume that both  $x$  and  $y$  have zero mean and unity variance. The order statistics  $x_{(i)}$  and  $y_{(i)}$  can be expressed as

$$\begin{aligned} x_{(i)} &= \mu_i + \varepsilon_i \\ y_{(i)} &= \mu_i + \delta_i \end{aligned} \quad (9)$$

where  $\mu_i = E(x_{(i)}) = E(y_{(i)})$  and  $E(\varepsilon_i) = E(\delta_i) = 0$ . It is obvious that  $\varepsilon_i$  and  $\delta_i$  have identical distributions, and hence we have

$$E(\varepsilon_i^2) = E(\delta_i^2). \quad (10)$$

The symmetry of the normal distribution yields [17], [18]

$$\begin{aligned} \mu_i &= -\mu_{N-i+1} \\ E(\varepsilon_i^2) &= E(\varepsilon_{N-i+1}^2). \end{aligned} \quad (11)$$

It follows that the concomitants associated with  $x_{(i)}$  can be written as [11]

$$y_{[i]} = \rho x_{(i)} + z_i \quad (12)$$

where  $x_{(i)}$  and  $z_i$  are mutually independent, the latter being IID normal with mean zero and variance  $1 - \rho^2$ . Given (9)–(12), the numerator in (1) can be expressed as

$$U = 2 \sum x_i y_i - \sum (\rho \mu_i + \rho \varepsilon_i + z_i)(\varepsilon_i + \varepsilon_{N-i+1}). \quad (13)$$

From (11), we have

$$\sum \mu_i(\varepsilon_i + \varepsilon_{N-i+1}) = -\sum \mu_{N-i+1}(\varepsilon_i + \varepsilon_{N-i+1}). \quad (14)$$

Replacing  $N - i + 1$  by  $j$  in (14) leads to

$$\sum \mu_i(\varepsilon_i + \varepsilon_{N-i+1}) = -\sum \mu_j(\varepsilon_j + \varepsilon_{N-j+1}) = 0. \quad (15)$$

Hence, (13) can be further simplified by

$$U = 2 \sum x_i y_i - \sum (\rho \varepsilon_i + z_i)(\varepsilon_i + \varepsilon_{N-i+1}). \quad (16)$$

Taking expectations on both sides of (16) and applying the mutual independence facts mentioned before, we have

$$E(U) = 2N\rho \left\{ 1 - \frac{\sum E(\varepsilon_i^2) + \sum E(\varepsilon_i \varepsilon_{N-i+1})}{2N} \right\}. \quad (17)$$

An application of Cauchy–Schwarz inequality [15] and the fact  $E(\varepsilon_i \varepsilon_{N-i+1}) > 0$  [11] to (14) yields

$$\begin{cases} 2N\rho > E(U) > 2N\rho \left[ 1 - \frac{\sum E(\varepsilon_i^2)}{N} \right], & \text{when } \rho \geq 0 \\ 2N\rho < E(U) < 2N\rho \left[ 1 - \frac{\sum E(\varepsilon_i^2)}{N} \right], & \text{when } \rho < 0. \end{cases} \quad (18)$$

It can be easily verified that the following two identities hold for any two real numbers  $a$  and  $b$ :

$$\begin{aligned} ab &= \frac{1}{2} [a^2 + b^2 - (a - b)^2] \\ -ab &= \frac{1}{2} [a^2 + b^2 - (a + b)^2]. \end{aligned} \quad (19)$$

Substituting (9)–(12) and (16) into the denominator in (1), we arrive at

$$V = \sum (x_i^2 + y_i^2) - \frac{1}{2} \sum (\varepsilon_i - \delta_i)^2 - \frac{1}{2} \sum (\varepsilon_i + \delta_{N-i+1})^2. \quad (20)$$

Applying elementary inequalities and taking expectation on both sides of (17), we have

$$2N \left[ 1 - \frac{2 \sum E(\varepsilon_i^2)}{N} \right] \leq E(V) < 2N. \quad (21)$$

It is implied in [18] that

$$\lim_{N \rightarrow \infty} \frac{\sum E(\varepsilon_i^2)}{N} = 0. \quad (22)$$

Substituting (15), (18), and (19) into (5) and letting  $N$  tend to infinity, we have  $\lim_{N \rightarrow \infty} E\{r_X(x, y)\} = \rho$ , hence, the result.

According to Theorem 2, we can estimate  $\rho$  using order statistics correlation coefficient by the following estimator:

$$\hat{\rho}_X = r_X. \quad (23)$$

### C. Comparison With Three Classical Correlation Coefficients

As defined in Section II-A,  $x_{(1)}, \dots, x_{(N)}$  are the order statistics of the time series  $x_1, \dots, x_N$ . Suppose  $x_j$  is at the  $k$ th position in the sorted series  $x_{(1)}, \dots, x_{(N)}$ , the number  $1 \leq k \leq N$  is termed the *rank* of  $x_j$  and is denoted by  $p_j (= k)$ . Similarly, we can get the rank of  $y_j$  denoted by  $q_j$ . Such operation of obtaining the ranks of all elements in a series is called *ranking* [5]. Let  $(x_i, y_i)$  and  $(x_j, y_j)$  with  $i = 1, \dots, N$  and  $j = i+1, \dots, N$  be two data-pairs from the original time series. If  $p_j - p_i$  and  $q_j - q_i$  have the same sign, we say that the two data-pairs are *concordant*, otherwise, we say that they are *discordant* [5]. Let  $P$  stand for the number of concordant pairs and  $Q$  the number of discordant pairs, it follows that  $P + Q = N(N-1)/2$ . Let  $\bar{x}$ ,  $\bar{y}$ ,  $\bar{p}$ , and  $\bar{q}$  be the arithmetic averages of  $x$ ,  $y$ ,  $p$ , and  $q$ , respectively, Pearson's correlation coefficient ( $r_P$ ), Spearman's rho ( $r_S$ ), and Kendall's tau ( $r_K$ ) are defined as follows [1], [5]:

$$r_P(x, y) \triangleq \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 \sum_{i=1}^N (y_i - \bar{y})^2}} \quad (24)$$

$$r_S(x, y) \triangleq 1 - \frac{6 \sum_{i=1}^N (p_i - q_i)^2}{N^3 - N} \quad (25)$$

$$r_K(x, y) \triangleq \frac{2(P - Q)}{N(N - 1)}. \quad (26)$$

If  $x$  and  $y$  are bivariate Gaussian with correlation coefficient  $\rho$ , three reasonable estimators of  $\rho$  can be constructed as follows [5]:

$$\begin{aligned} \hat{\rho}_P &= r_P \\ \hat{\rho}_S &= 2 \sin \left( \frac{1}{6} \pi r_S \right) \\ \hat{\rho}_K &= \sin \left( \frac{1}{2} \pi r_K \right). \end{aligned} \quad (27)$$

For brevity, we will drop the circumflex throughout and use  $\rho_\xi$ ,  $\xi = X, P, S, K$  to denote the four estimators. It is easy to verify that 1)  $\rho_\xi$  are monotonic with  $r_\xi$ ; 2) three critical values  $\{-1, 0, +1\}$  are invariant under the transformations  $\rho_\xi$ ; and 3)  $\rho_\xi$  is also limited within  $[-1, +1]$ .

### III. MODELS OF ASSOCIATION AND PERFORMANCE EVALUATION

In this section, we propose three linear models and one nonlinear model to model the linear and nonlinear association between two time series. Several indices will also be proposed to evaluate the performance of our order statistics correlation coefficient in comparison with the other classical correlation coefficients, in terms of their abilities to estimate the associations between time series. In each model, a time series  $x(i)$  is derived from a pure signal  $s(i)$ , and another signal  $y(i)$  is obtained as a combination of the transformed pure signal and a white noise,  $n(i)$ . In all these models, the time index  $i$  runs from 1 to 1000.

#### A. Models of Association

1) *Linear Model 1 (LM1)*: LM1 is constructed as

$$\begin{aligned} x(i) &= s(i) \\ y(i) &= s(i) + \alpha \cdot n(i) \end{aligned} \quad (28)$$

where  $\alpha \in [0, 1]$  is increased from 0 to 1 with a step  $\Delta\alpha = 0.1$  to control the signal-to-noise ratio (SNR). With increasing  $\alpha$ , the association between  $x$  and  $y$  becomes smaller and smaller, which means that  $r_\xi$  ( $\rho_\xi$ ) should have a decreasing relationship with  $\alpha$ . For a fixed  $\alpha$ , the greater the magnitude of  $E(\rho_\xi)$ , the better its performance in the context of noise robustness.

2) *Linear Model 2 (LM2)*: LM2 is a regression model of the form [14]

$$\begin{aligned} x(i) &= s(i) \\ y(i) &= \rho \cdot s(i) + \sqrt{1 - \rho^2} \cdot n(i) \end{aligned} \quad (29)$$

where  $\rho \in [-1, 1]$  with a step  $\Delta\rho = 0.01$  characterizing the linear association. It follows by straightforward calculation that  $E(\rho_P) = \rho$  for any distribution of  $s(i)$ . Unfortunately, the property of unbiasedness does not hold for the other three estimators  $\rho_X$ ,  $\rho_S$ , and  $\rho_K$  except for the bivariate normal case. The aim of this model is to compare the biasedness of these three biased estimators as well as their power to discriminate different  $\rho$ 's.

3) *Linear Model 3 (LM3)*: LM3 is similar to LM1 except for a time delay  $\Delta = 30$  introduced in channel  $y$ , as follows:

$$\begin{aligned} x(i) &= s(i) \\ y(i) &= s(i - \Delta) + \alpha \cdot n(i). \end{aligned} \quad (30)$$

4) *Nonlinear Model (NM)*: NM is a nonlinear model used to study the effect of nonlinear transformations to the signals on the four coefficients, as follows [14]:

$$\begin{aligned} x(i) &= T_x[\beta \cdot s(i)] \\ y(i) &= T_y\left[\beta \left\{ \rho \cdot s(i) + \sqrt{1 - \rho^2} \cdot n(i) \right\}\right] \end{aligned} \quad (31)$$

where  $T_x[\bullet]$  and  $T_y[\bullet]$  are two increasing nonlinear functions. The parameter  $\beta = 2, 4, 6, 8, 10$  is used to control the extent of nonlinearity (greater value of  $\beta$  corresponding to stronger nonlinearity), while  $\rho$  has the same meaning as in LM2.

#### B. Performance Evaluation

Several methods are used to evaluate the performance of  $r_\xi$  under each of the four models previously mentioned.

1) *Noise Robustness*: Under LM1, we compare the decreasing rates of  $E(\rho_\xi)$  with the increase of  $\alpha$ .

2) *Attenuation Measurement*: We propose two indices called absolute attenuation (AAT) and relative attenuation (RAT) to measure the extent of biasedness, defined as

$$\text{AAT}_\xi = \frac{\int_{-1}^1 |\bar{\rho}_\xi(\rho) - \rho| d\rho}{\int_{-1}^1 \int_{-1}^1 d\rho d\rho_\xi} = \frac{1}{4} \int_{-1}^1 |\bar{\rho}_\xi(\rho) - \rho| d\rho \quad (32)$$

and

$$\text{RAT}_\xi = \frac{\int_{-1}^1 |\bar{\rho}_\xi(\rho) - \bar{\rho}_\xi^{(L)}(\rho)| d\rho}{\int_{-1}^1 \int_{-1}^1 d\rho d\rho_\xi} = \frac{1}{4} \int_{-1}^1 |\bar{\rho}_\xi(\rho) - \bar{\rho}_\xi^{(L)}(\rho)| d\rho \quad (33)$$

where  $\bar{\rho}_\xi^{(L)}(\rho)$  is the mean of  $\rho_\xi$  under LM2. AAT is to measure the extent of biasedness under LM2, while RAT is to measure the biasedness caused by the nonlinearity involved in NM. For Pearson's coefficient, we have  $\bar{\rho}_P = \rho$ , thus  $\text{AAT}_P = 0$ . For the other three coefficients, AAT relates positively to their biasedness under LM2. It is clear from (30) that under LM2  $\text{RAT}_\xi = 0$  and under NM, a smaller RAT means lesser effect of the nonlinearity and hence better robustness.

3) *Sensitivity to Changes in  $\rho$* : We employ another index called sensitivity ratio (SR) [4] to test the sensitivity of  $\rho_\xi$  to changes in  $\rho$ . For this purpose the Fisher's  $z$ -transformation of  $\rho_\xi$ , denoted as  $z_\xi$

$$z_\xi = \tanh^{-1} \rho_\xi = \frac{1}{2} \log_e \frac{1 + \rho_\xi}{1 - \rho_\xi}. \quad (34)$$

After such transformation, which maps  $[-1, +1]$  to  $(-\infty, +\infty)$ , the resultant  $z_\xi$  follows approximately normal distributions with constant variances (i.e., independent of the means) [2], [4]. Given two distinct  $\rho_1$  and  $\rho_2$  ( $\rho_2 > \rho_1$ ), we have two sets of coefficients  $\rho_{\xi 1}$  and  $\rho_{\xi 2}$ , and their respective Fisher's  $z$ -transformation,  $z_{\xi 1}$  and  $z_{\xi 2}$ . SR is then defined as

$$\text{SR}_\xi = \frac{\bar{z}_{\xi 2} - \bar{z}_{\xi 1}}{\sqrt{v_{z1}^2 + v_{z2}^2}} \quad (35)$$

where  $\bar{z}_\xi$  and  $v_\xi$  denote the mean and standard deviation of  $z_\xi$ , respectively. SR measures the ability of  $\rho_\xi$  to detect the changes of underlying  $\rho$ . A greater value of SR indicates better discrimination sensitivity. SR is computed for the results of the linear model LM2 and the nonlinear model NM.

4) *Time Complexity Measurement*: We analyze the time complexities of  $r_\xi$  in the language of *big Oh*. We also estimate the relationship between computational loads of  $r_\xi$  versus the length of signal  $N$  from 100 to 1000 with a step  $\Delta N = 100$ .

### IV. COMPARISON ON RESULTS FOR SIMULATED AND REAL BIOSIGNALS

Signals derived from biological processes fall into two main categories: deterministic and stochastic signals. The former are

those that can be described by explicit mathematical relationships; whereas the latter can be described only in statistical terms. The deterministic group is further subdivided into periodic, semi-periodic, and transient signals; the stochastic group is subdivided into stationary and nonstationary signals [19], [20]. Recently, there has been growing evidence indicating that many biosignals exhibit long range power-law correlations [21]. A stochastic process is said to have long range correlation if its autocorrelation function  $R(k) \sim k^{2H-2}$  as  $k \rightarrow \infty$ , where  $0 < H < 1$  is the Hurst parameter. The corresponding power spectral density is proportional to  $f^{-(2H-1)}$  [22]. In order to evaluate the feasibility of  $r_X$  in association studies, several simulated and real biosignals with respect to six types of previously mentioned biosignals are employed for investigation. For notational convenience, the six signals are denoted uniformly as  $s_\zeta$ ,  $\zeta = p, h, t, a, e, l$ , which represent periodic, semi-periodic, transient, stationary, nonstationary, and long-range-correlated signals. A number of 1000 independent white Gaussian noise ( $\mu = 0$  and  $\sigma^2 = 1$ ) are generated with a sampling rate of 1000 Hz to serve as noise in the linear and nonlinear models. Due to the 1000 noise involved, each  $r_\xi$  becomes a random variable and has a distribution, which allows us to perform statistical analysis. Under each model, two channels of signals  $x$  and  $y$  are generated from  $s_\zeta$  and the 1000 episodes of white noise. Four sets of correlation coefficients between  $x$  and  $y$  are then computed for comparative study.

#### A. Simulated and Real Biosignals

As remarked before, the following six representative biosignals are included in our study:

- 1) sin wave  $s_p(i)$  of frequency 5 Hz emulating periodic biosignals;
- 2) real bipolar intra-atrial flutter signal  $s_h(i)$  recorded during electrophysiological procedure [23];
- 3) atrial action potential waveform  $s_t(i)$  generated from a mathematical model [24];
- 4) episode of alpha wave  $s_a(i)$  simulated from a random Gaussian noise filtered by a band-pass Butterworth filter with passband 8 to 12 Hz [25];
- 5) second of real EEG signal  $s_e(i)$  (sampling rate 256 Hz) from a dataset provided by University of Tuebingen for BCI Competition 2003 [26], [27];
- 6) segmentation  $s_l(i)$  of an artificial time series  $s_{lf}(i)$  exhibiting long range correlation with Hurst parameter  $H = 0.9$  [21], [28], [29].

Fig. 1 illustrates the six biosignals. All the first four signals contain 1000 samples ( $N = 1000$ ). The EEG signal  $s_e(i)$  is up-sampled from 256 to 1000 Hz by linear interpolation. As for the long-range-correlated signal  $s_{lf}(i)$  containing  $2^{17}$  samples, we use the first 1000 samples as  $s_l(i)$  for association analysis. The whole time series  $s_{lf}(i)$  is used when we demonstrate the capability of  $r_X$  for estimating the Hurst parameter  $H$ . After these manipulations, all the six original biosignals  $s_\zeta$ ,  $\zeta = p, h, t, a, e, l$  can be considered of duration 1 s. Without loss of generality (property c),  $s_\zeta$  are normalized to have mean zero and variance unity before feeding them into the four models described in Section III-B.

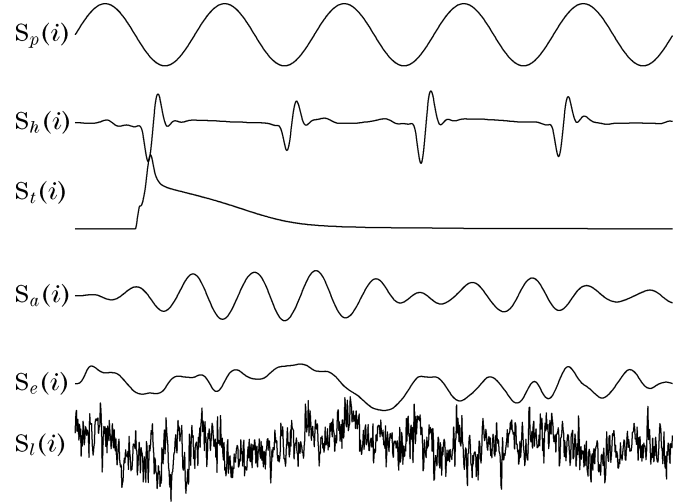


Fig. 1. Illustration of the four simulated signals and two real signals.

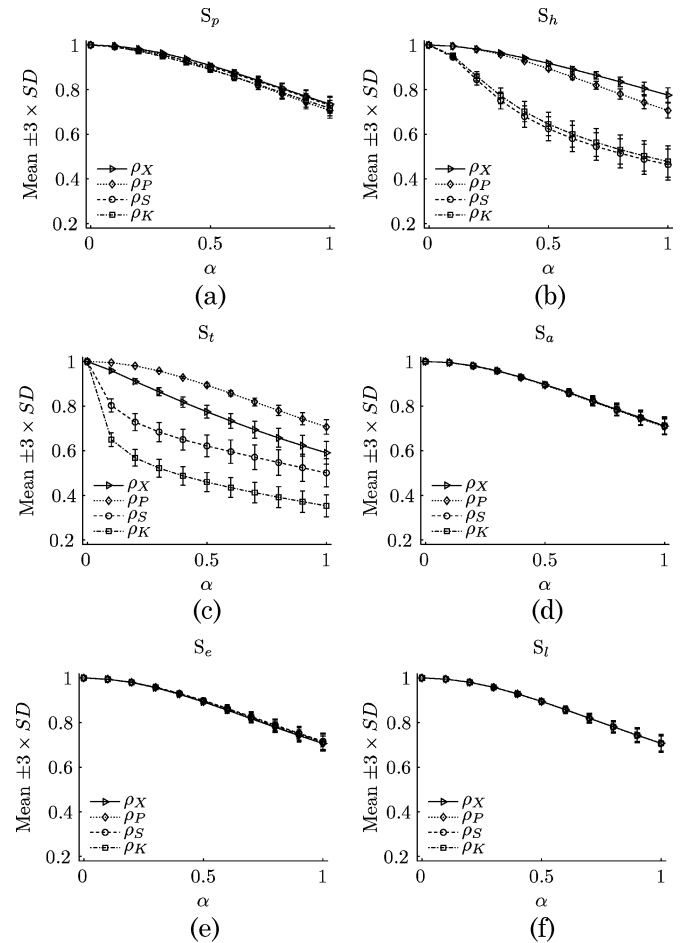
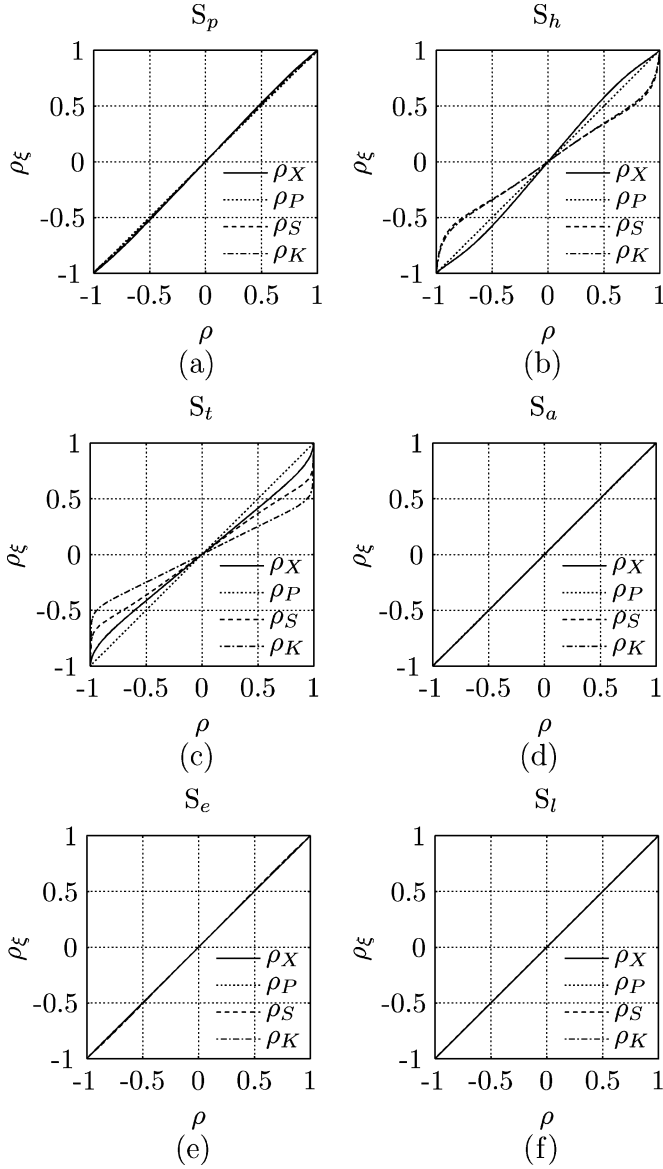


Fig. 2. Noise robustness comparison of simulation results under LM1.

#### B. Comparative Study Under Linear Model LM1

The results under LM1 are shown in Fig. 2. It is clear that  $\bar{\rho}_\xi(\alpha)$  drops with increasing of  $\alpha$ . However, the decreasing rates are quite different. Fig. 2(a)–(c) illustrate that the means of  $\rho_X$  and  $\rho_P$  descend more slowly than those of  $\rho_S$  and  $\rho_K$ , suggesting the superiority of the former two coefficients when deterministic signals  $s_p$ ,  $s_h$ , and  $s_t$  are fed into LM1. Furthermore,

Fig. 3. Relations between  $\bar{\rho}_\xi$  and  $\rho$  under LM2.

$\rho_X$  outperforms  $\rho_P$  in cases with respect to  $s_p$  and  $s_h$ . On the other hand, the immaterial differences observed in Fig. 2(d)–(f) indicating the equivalence of the four methods when the inputs are stochastic signals  $s_a$ ,  $s_e$ , and  $s_l$ . The overall noise robustness performance thus can be ordered as  $r_X, r_P > r_K$ , and  $r_S$ .

### C. Comparative Study Under Linear Model LM2

The relationships between  $\bar{\rho}_\xi$  and  $\rho$  for the six original signals  $s_\xi$  are shown in Fig. 3. It is easily observed that 1)  $-1 \leq \rho_\xi \leq 1$  (Property 1); 2)  $\rho_\xi = \pm 1$  ( $r_\xi = \pm 1$ ) as  $\rho = \pm 1$ , respectively (Property 2); 3)  $\bar{\rho}_\xi = 0$  ( $\bar{r}_\xi = 0$ ) as  $\rho = 0$  (Property 4); and 4)  $\bar{\rho}_\xi$  ( $\bar{r}_\xi$ ) is an increasing function of  $\rho$ . Noticing that the closer the distance of  $\bar{\rho}_\xi$  to the diagonal line, the smaller the associated biasedness, we also observe that when the model inputs are deterministic signals excepting  $s_p$ , the unbiasedness performance can be ordered as  $r_P > r_X > r_K > r_S$ ; whereas for the stochastic signals, there is no substantial difference among the four methods. This phenomenon is quantitatively highlighted in

TABLE I  
AAT COMPARISON RESULTS AMONG FOUR METHODS

	$AAT_X$ (%)	$AAT_P$ (%)	$AAT_S$ (%)	$AAT_K$ (%)
$S_p$	0.88	0.07	0.55	0.33
$S_h$	2.34	0.02	7.53	7.12
$S_t$	3.75	0.01	6.95	12.12
$S_a$	0.09	0.01	0.13	0.13
$S_e$	0.04	0.02	0.23	0.16
$S_l$	0.02	0.01	0.03	0.03

Table I which summarizes the  $AAT_\xi$  with respect to the six original signals  $s_\xi$ . As for the performance of detecting changes in the underlying  $\rho$ , we tabulate the sensitivity ratios (SR) in Table II, showing that the capability of discriminating changes in  $\rho$  can be ordered as  $r_P > r_X > r_K > r_S$ , the same as that of unbiasedness performance.

### D. Comparative Study Under Linear Model LM3

Under this model,  $r_\xi$  is computed as a function of time-shift  $\kappa$ , say, which varies from  $-100$  to  $100$  ms. For each  $\alpha$  and each episode  $n(i)$  of 1000 white noises,  $r_\xi(\alpha, \kappa)$  is calculated and the time-shift with respect to the maximum of  $r_\xi(\alpha, \kappa)$  is the estimate of the time-delay  $\Delta$  and denoted by  $\kappa_\Delta$ . Limited by the length of this paper, we only present the results with respect to  $s_e$  here. Fig. 4(a) shows four typical waveforms of  $r_\xi(\alpha, \kappa)$  in the presence of a 50% SNR ( $\alpha = 1$ ). All the four coefficients can correctly detect the time-delay between  $x$  and  $y$  giving  $\kappa_\Delta = 30$  ms which equals the true time-delay  $\Delta$ . In Fig. 4(b), we present the statistical results of  $\kappa_\Delta$  versus the underlying  $\alpha$  from 0 to 1 with  $\Delta\alpha = 0.1$ . The levels of rectangular bars represent the means  $\bar{\kappa}_\Delta$  and the error bars represent  $3 \times v_{\kappa_\Delta}$  with  $v_{\kappa_\Delta}$  denoting the standard deviation of  $\kappa_\Delta$ . It can be observed from Fig. 4(b) that  $\bar{\kappa}_\Delta$  slightly increases with increase of noise levels and so does the standard deviation  $v_{\kappa_\Delta}$  for all four  $r_\xi$ . The performances of  $r_X$  and  $r_P$  are better than those of two rank-based methods  $r_S$  and  $r_K$  in the sense that the former two coefficients have smaller deviations. However, the performance of time-delay detection is not further compared since the maximal error is only 2 ms in all cases for all the four methods. In other words, we do not consider that there are significant differences between the four methods in the aspect of detecting time delays.

### E. Comparative Study Under Nonlinear Model NM

The nonlinear model NM is constructed on the linear model LM2 by introducing two increasing nonlinear transformations  $T_x[\bullet] = \text{sgn}(\bullet) \cdot (\bullet)^2$  and  $T_y[\bullet] = \exp(\bullet)$ . Besides the association parameter  $\rho$  carrying the same meaning as in LM2, we employ another parameter  $\beta = 2, 4, 6, 8, 10$  to control the extent of nonlinearity. It is noteworthy that  $T_x[\bullet]$  and  $T_y[\bullet]$  have no effect on the rank-based measures  $r_S$  and  $r_K$  because rank-ings are invariant under strictly increasing transformations.

Fig. 5 shows the relationships between  $\bar{\rho}_\xi$  and  $\rho$  with respect to the six biosignals  $s_\xi$  with nonlinearity parameter  $\beta = 2$ . We observe that 1) for periodic signal  $s_p$ ,  $\rho_X$  performs comparably

TABLE II  
SRs COMPARISON OF THE RESULTANTS OF LM2

$\rho_1$	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	
$\rho_2$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	
$S_p$	$SR_X$	2.2295	2.2899	2.4286	2.6697	3.0628	3.7103	4.8474	7.1243	13.1615
	$SR_P$	2.2301	2.2916	2.4317	2.6754	3.0746	3.7386	4.9257	7.3876	14.3773
	$SR_S$	2.1470	2.2032	2.3390	2.5796	2.9761	3.6249	4.7660	6.9956	11.8843
	$SR_K$	2.1459	2.2022	2.3369	2.5730	2.9579	3.5827	4.6581	6.6673	11.0515
$S_h$	$SR_X$	2.2082	2.2665	2.3955	2.6144	2.9626	3.5175	4.4560	6.2641	10.8761
	$SR_P$	2.2387	2.3011	2.4424	2.6880	3.0899	3.7581	4.9525	7.4295	14.4629
	$SR_S$	1.5498	1.5433	1.5415	1.5476	1.5699	1.6228	1.7580	2.1406	3.4626
	$SR_K$	1.5508	1.5493	1.5556	1.5722	1.6087	1.6789	1.8356	2.2493	3.6324
$S_t$	$SR_X$	1.9265	1.9686	2.0616	2.2160	2.4528	2.8147	3.3906	4.4035	6.6993
	$SR_P$	2.2821	2.3454	2.4892	2.7393	3.1485	3.8290	5.0456	7.5686	14.7325
	$SR_S$	1.6928	1.7098	1.7427	1.7932	1.8643	1.9614	2.0913	2.2770	2.7136
	$SR_K$	1.6920	1.7096	1.7442	1.7977	1.8741	1.9821	2.1356	2.3810	2.9659
$S_a$	$SR_X$	2.2326	2.2949	2.4359	2.6805	3.0802	3.7436	4.9265	7.3708	14.2385
	$SR_P$	2.2313	2.2934	2.4341	2.6787	3.0790	3.7444	4.9337	7.3994	14.3982
	$SR_S$	2.1590	2.2187	2.3507	2.5813	2.9461	3.5381	4.5339	6.4251	11.0119
	$SR_K$	2.1571	2.2177	2.3520	2.5872	2.9648	3.5873	4.6611	6.7816	12.2714
$S_e$	$SR_X$	2.2965	2.3595	2.5026	2.7510	3.1563	3.8265	5.0136	5.0136	14.0770
	$SR_P$	2.3022	2.3663	2.5118	2.7645	3.1782	3.8662	5.0966	5.0966	14.9067
	$SR_S$	2.2352	2.2953	2.4333	2.6654	3.0398	3.6439	4.6804	4.6804	11.8162
	$SR_K$	2.2359	2.2965	2.4371	2.6732	3.0580	3.6863	4.7813	4.7813	12.6448
$S_l$	$SR_X$	2.2445	2.3069	2.4486	2.6947	3.0974	3.7668	4.9636	7.4481	14.5038
	$SR_P$	2.2457	2.3082	2.4501	2.6967	3.1004	3.7720	4.9735	7.4678	14.5614
	$SR_S$	2.1617	2.2138	2.3385	2.5521	2.8966	3.4566	4.4232	6.3403	11.3488
	$SR_K$	2.1615	2.2166	2.3461	2.5694	2.9336	3.5311	4.5755	6.6912	12.4021

with  $\rho_S$  and  $\rho_K$ ; 2) in cases with respect to other two deterministic signals,  $\rho_X$  has the smallest biasedness; 3) for the three stochastic signals,  $\rho_K$  and  $\rho_S$  surpass  $\rho_X$ ; and 4)  $\rho_P$  has the largest biasedness in all cases. Moreover,  $r_P$  never approaches  $\pm 1$  as  $\rho \rightarrow \pm 1$ , that is,  $r_P$  underestimates the strength of association when nonlinearity is involved. On the other hand,  $r_X = \pm 1$  as  $\rho \rightarrow \pm 1$ , which indicates the validity of property 2 under increasing nonlinear transforms. Table III shows the effect of nonlinearity on  $\rho_X$  and  $\rho_P$ , elucidating that with increase of  $\beta$ , the biasedness (RAT) caused for  $\rho_X$  is significantly smaller than that for  $\rho_P$ . Sensitivity ratios ( $\beta = 2$ ) tabulated in Table IV show that in most cases,  $\rho_K$  performs best; in almost all cases,  $\rho_P$  has the lowest SR, indicating the limitation of  $r_P$  in nonlinear scenarios.

#### F. Comparison of Time Complexities

The time complexities of  $r_\xi$  are analyzed and summarized in Table V based on the definitions. The fastest method is Pearson's coefficient  $r_P$  having a linear time complexity of  $O(N)$ . Our new method  $r_X$  and Spearman's rho are of the same order  $O(N \log N)$ , since sorting operation dominates the computational time of both methods. However, because of the extra procedure of ranking involved in the calculation of  $r_S$ , we can expect that  $r_X$  is a little faster than  $r_S$ . Kendall's tau  $r_K$  is the slowest method compared to the other three coefficients. The core operation of  $r_K$  is to calculate the number of concordant and discordant pairs, which requires  $C_N^2 = N(N-1)/2$  operations. Therefore, the time complexity of  $r_K$  is of  $O(N^2)$ .

To confirm this result, we estimate the relationship between computational loads of  $r_\xi$  versus the length of signal  $N$ , where  $N$  begins at 100 and increases by steps of 100 until  $N = 1000$ . All the computational speed tests were performed in MATLAB 7.0 in a Pentium PC. For each pair of time series of size  $N$ ,

the algorithms of  $r_\xi$  were run for 1000 times. The results are presented in Fig. 6, which is consistent with our analysis.

## V. DISCUSSION

### A. Rationality of Selection of Comparison Objects

Apart from the three methods used in our comparative study, some authors have used other techniques, such as the AAMI [7], [8], the nonlinear regression coefficient  $h^2$  [9], and the contingency table-based method  $V$  [10]. However, we did not include these methods into our comparative studies due to 1) AAMI is unbound [9] and hence is incomparable to our new method and 2) the admissible values of  $h^2$  and  $V$  are confined from 0 to 1, that is, they cannot distinguish positive associations from negative associations. Moreover, the analytical relationship between  $V(h^2)$  and  $\rho$  under bivariate normal model is unknown, which prevents us from doing calibration as in (23) and (27). On the other hand, Pearson's coefficient, Spearman's rho and Kendall's tau have similar meaning as our new measure, thus making it possible to compare their behaviors.

### B. Estimation of the Hurst Parameter $H$ by $r_X$

As remarked before, the power spectral density of a long-range-correlated signal is proportional to  $f^{-(2H-1)}$ , where  $0 < H < 1$  is the Hurst parameter. Therefore, we can estimate  $H$  from the power density which is defined by the Fourier transform of the corresponding autocorrelation function  $R(k)$ ,  $k \geq 0$ . Noticing that  $r_P(k)$  is the normalized version of  $R(k)$  and  $r_X(k)$  behaves similarly with  $r_P(k)$  under linear models, we can estimate the power density  $S(f)$  based on the Fourier transform of  $r_P(k)$  as well as  $r_X(k)$ . The Hurst parameter  $H$  can then be measured with the slope of  $\log(S(f))$  against  $\log(f)$ . In Fig. 7, we present the results on the full version of the long-range-correlated signal  $s_{lf}(i)$ . For clarity, the waveforms of  $r_P$

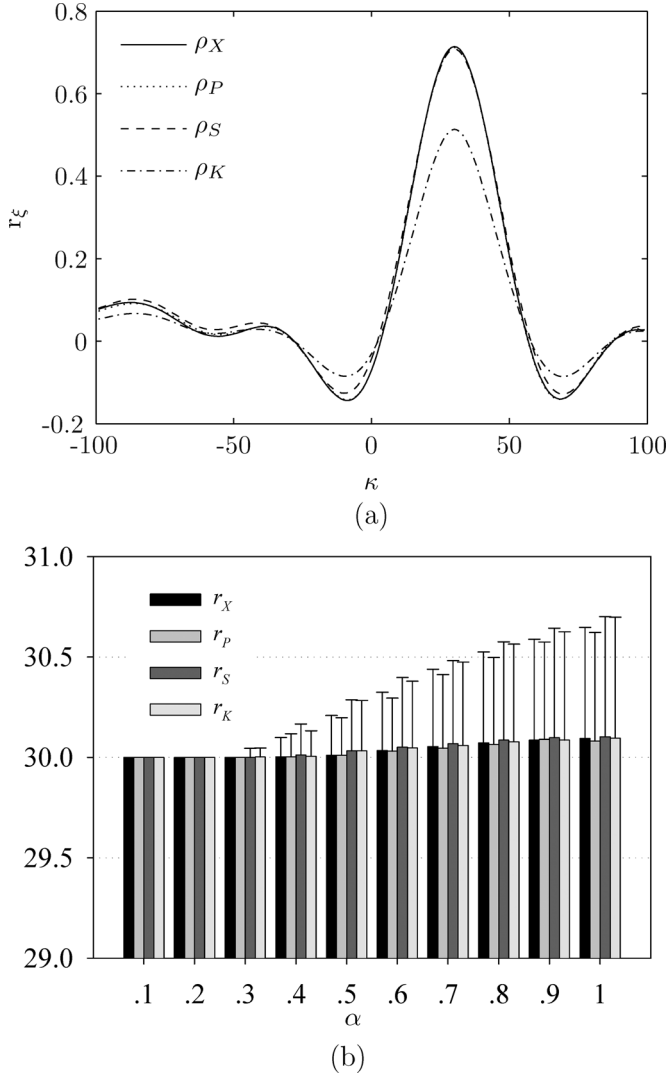


Fig. 4. Results of simulations under LM3 with respect to  $s_e$ . (a) Magnitudes of  $r_\xi$  associated with time-shift  $\kappa$  varying from  $-100$  to  $100$  ms in the presence of a 50% SNR ( $\alpha = 1$ ). (b) Statistical results of  $\kappa_\Delta$  versus the underlying  $\alpha$  from 0 to 1 with a step 0.1.

in Fig. 7 are vertically shifted down, since otherwise the waveforms of  $r_X$  and  $r_P$  will be almost coincident and, therefore, unclear to observe. The corresponding power densities are plotted in a double-log scale in Fig. 7(b), where  $\delta = 2\hat{H} - 1$ . It can be easily obtained that  $\hat{H}_X = 0.8962$  and  $\hat{H}_P = 0.8960$ , very close to the real value  $H = 0.9$ .

### C. Sufficiency of Signal Length Used in this Study

It can be shown that the variances of  $r_X$ ,  $r_P$ ,  $r_S$ , and  $r_K$  are all of the order of  $O(N^{-1})$  [2], [5], [30], where  $N$  is the length of the signals. In other words, the larger the sample size, the smaller the variances, and the more accurate the four coefficients. For signals consisting of 1000 sample points, the variances of the four coefficients are already very small (of the order of 0.001), which means that a sample size of 1000 is sufficient for studying the behaviors of these coefficients. For the long-range-correlated time series whose length is far greater

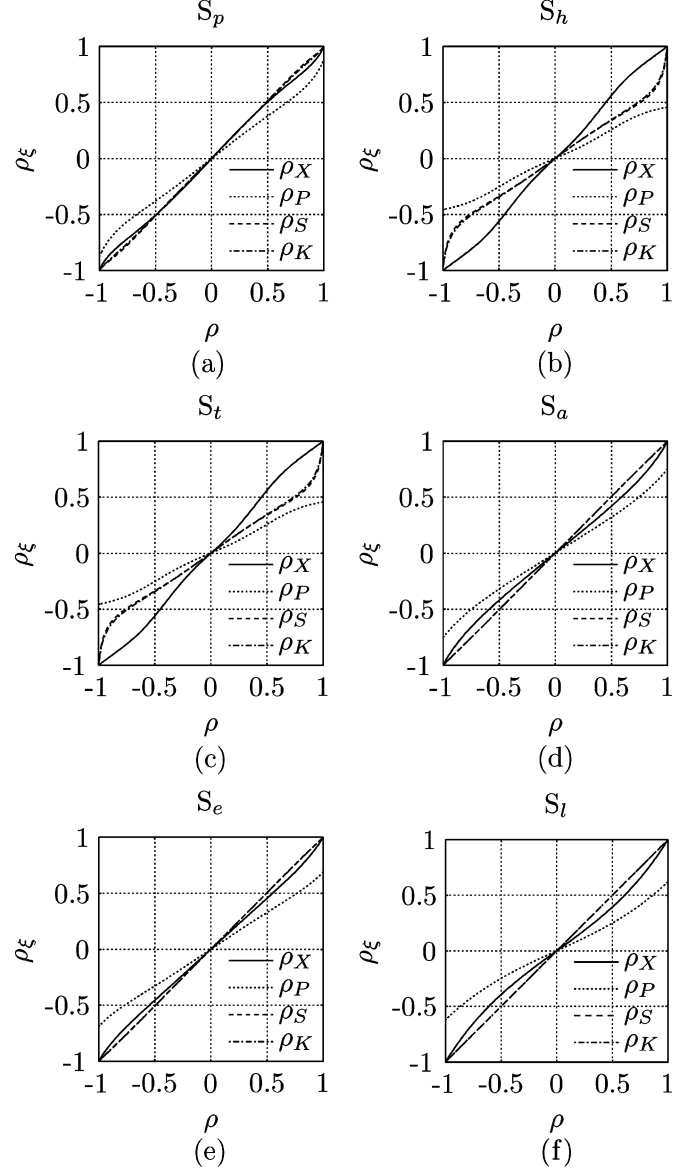


Fig. 5. Relations between  $\bar{\rho}_\xi$  and  $\rho$  under NM. All subplots illustrate the inferiority of  $r_P$  in this nonlinear model. The wandering from positive (negative) unity of  $\rho_P(r_P)$  when  $\rho = \pm 1$  reveals the misleading results of  $r_P$  when non-linearity involved.

than 1000, the variances of the four cross-correlation coefficients are very close to zero, thus assuring the accuracy of the four measurements of association.

### D. Solution to Asymmetry of $r_X$

In general, our new measure  $r_X$  is not symmetric, namely,  $r_X(x, y) \neq r_X(y, x)$  although  $E\{r_X(x, y)\} = E\{r_X(y, x)\}$ . This problem can be easily solved using a revised version  $r'_X(x, y) = [r_X(x, y) + r_X(y, x)]/2$  when symmetry is a critical feature in practice.

### E. Clinical Application of $r_X$

It is of great clinical importance for speedy and reliable detection of atrial fibrillation (AF) and atrial flutter (AFL) in automatic implantable atrial defibrillators [23]. AF is a type of arrhythmia (abnormal heart rhythm) exhibiting rapid and



TABLE III  
RAT COMPARISON UNDER NONLINEAR MODEL NM BETWEEN  $r_X$  AND  $r_P$

	$S_p$		$S_h$		$S_t$		$S_a$		$S_e$		$S_l$	
	$RAT_X$	$RAT_P$	$RAT_X$	$RAT_P$	$RAT_X$	$RAT_P$	$RAT_X$	$RAT_P$	$RAT_X$	$RAT_P$	$RAT_X$	$RAT_P$
	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)
2	0.93	5.68	0.84	12.79	4.37	7.97	3.04	8.26	1.74	8.57	3.57	11.60
4	2.96	8.25	2.61	16.03	5.03	11.03	5.51	12.34	4.15	12.64	5.45	17.10
6	4.50	9.78	3.69	17.63	5.47	12.39	6.87	14.51	5.49	14.68	5.98	18.97
8	5.53	10.83	4.32	18.48	5.79	13.27	7.77	15.93	6.37	15.94	6.19	19.87
10	6.29	11.60	4.72	19.00	6.04	13.92	8.41	16.93	7.02	16.83	6.28	20.38

TABLE IV  
SRS COMPARISON OF THE RESULTANTS OF NM. THE MAXIMA AND MINIMA ARE HIGHLIGHTED WITH GRAY AREAS AND WHITE BOXES, RESPECTIVELY

	$\rho_1$	$\rho_2$		$\rho_3$	$\rho_4$	$\rho_5$	$\rho_6$	$\rho_7$	$\rho_8$	$\rho_9$
		0.0	0.1							
$S_p$	$SR_X$	1.9808	2.0212	2.1149	2.2729	2.5144	2.8715	3.4170	4.4060	7.1379
	$SR_P$	1.7324	1.7513	1.7989	1.8794	2.0017	2.1885	2.5027	3.1359	4.8094
	$SR_S$	2.1470	2.2032	2.3390	2.5796	2.9761	3.6249	4.7660	6.9956	11.8843
	$SR_K$	2.1459	2.2022	2.3369	2.5730	2.9579	3.5827	4.6581	6.6673	11.0515
$S_h$	$SR_X$	1.5219	1.5283	1.5456	1.5755	1.5975	1.5978	1.6121	1.7553	2.3737
	$SR_P$	0.9321	0.9152	0.9054	0.9011	0.8900	0.8546	0.7897	0.7104	0.6535
	$SR_S$	1.5498	1.5433	1.5415	1.5476	1.5699	1.6228	1.7580	2.1406	3.4626
	$SR_K$	1.5508	1.5493	1.5556	1.5722	1.6087	1.6789	1.8356	2.2493	3.6324
$S_t$	$SR_X$	0.7710	0.7942	0.8415	0.9037	0.9843	1.1082	1.3222	1.7187	2.6497
	$SR_P$	0.9064	0.8863	0.8891	0.9200	0.9696	1.0323	1.1083	1.2123	1.4081
	$SR_S$	1.6928	1.7098	1.7427	1.7932	1.8643	1.9614	2.0913	2.2770	2.7136
	$SR_K$	1.6920	1.7096	1.7442	1.7977	1.8741	1.9821	2.1356	2.3810	2.9659
$S_a$	$SR_X$	1.7298	1.7375	1.7771	1.8586	2.0021	2.2485	2.6880	3.5533	5.6944
	$SR_P$	1.3216	1.3213	1.3296	1.3529	1.4024	1.4977	1.6745	2.0101	2.7336
	$SR_S$	2.1590	2.2187	2.3507	2.5813	2.9461	3.5381	4.5339	6.4251	11.0119
	$SR_K$	2.1571	2.2177	2.3520	2.5872	2.9648	3.5873	4.6611	6.7816	12.2714
$S_e$	$SR_X$	1.9998	2.0313	2.1101	2.2460	2.4591	2.7911	3.3468	4.4269	7.2043
	$SR_P$	1.5529	1.5642	1.5881	1.6248	1.6769	1.7529	1.8789	2.1231	2.6991
	$SR_S$	2.2352	2.2953	2.4333	2.6654	3.0398	3.6439	4.6804	6.6980	11.8162
	$SR_K$	2.2359	2.2965	2.4371	2.6732	3.0580	3.6863	4.7813	6.9522	12.6448
$S_l$	$SR_X$	1.4193	1.3968	1.3618	1.3481	1.3825	1.4963	1.7389	2.2406	3.4734
	$SR_P$	1.0233	0.9999	0.9652	0.9401	0.9388	0.9726	1.0573	1.2269	1.5849
	$SR_S$	2.1617	2.2138	2.3385	2.5521	2.8966	3.4566	4.4232	6.3403	11.3488
	$SR_K$	2.1615	2.2166	2.3461	2.5694	2.9336	3.5311	4.5755	6.6912	12.4021

TABLE V  
TIME COMPLEXITY ANALYSIS FOR FOUR METHODS

Operation	$r_X$	$r_P$	$r_S$	$r_K$
$\pm$	$O(N)$	$O(N)$	$O(N)$	$O(N^2)$
$\times, \div$	$O(N)$	$O(N)$	$O(1)$	$O(N^2)$
$( )^2$	—	$O(N)$	$O(N)$	—
$\sqrt{\quad}$	—	$O(1)$	—	—
Sort	$O(N \log N)$	—	$O(N \log N)$	$O(N \log N)$
Ranking	—	—	$O(N)$	$O(N)$
TC	$O(N \log N)$	$O(N)$	$O(N \log N)$	$O(N^2)$

random patterns. AFL is another kind of arrhythmia caused by electrical activity propagating through the atria in a fast and regular manner. Some researchers have employed  $r_P$  as

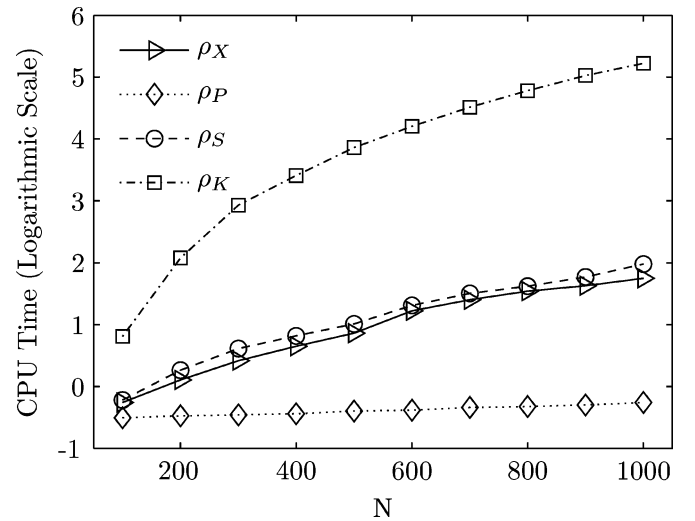


Fig. 6. Results of comparative CPU time test for four coefficients studied. A logarithmic scale is used for better visual effect.

an index to detect these two intra-atrial electrograms [31]. Now, we show that  $r_X$  can serve as a useful alternative to

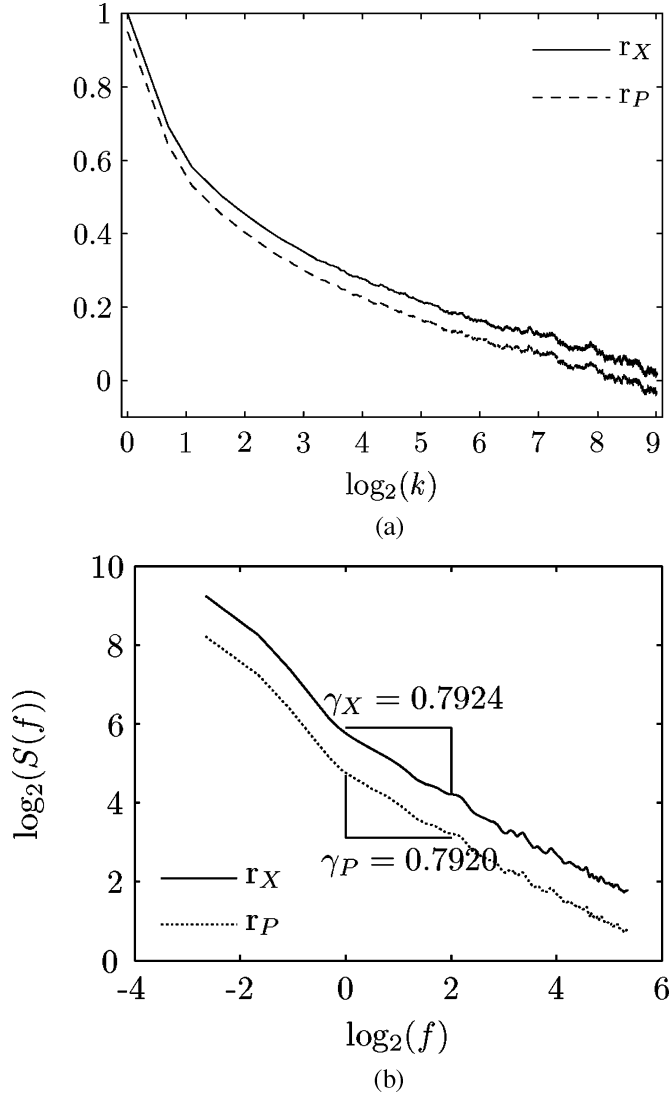


Fig. 7. Estimation of the Hurst parameter  $H$  with respect to the long-range-correlated signal  $s_{lf}(i)$  by  $r_X$  and  $r_P$ . The waveforms of  $r_P$  are vertically shifted down for a better visual effect.

$r_P$  in discriminating AF from AFL. Bipolar intra-atrial electrograms (available on physioBank [29]) from one AF patient and one AFL patient are included in this example. The continuous recordings are parsed into nonoverlapping segments of 1-s duration (1000 samples). After a series of preprocessing steps as diagrammed in Fig. 8(a) and detailed in [31], three cross-correlation functions  $r_X(k)$ ,  $r_P(k)$ , and  $r_S(k)$  of time-shift  $-100 \leq k \leq 100$  are calculated with respect to each pair of preprocessed signal segments from two channels. Three maximal values with respect to  $r_X(k)$ ,  $r_P(k)$ , and  $r_S(k)$  are then extracted as discriminatory indices. This operation is repeated sequentially over the entire dataset so that statistical analysis can be performed. It can be observed in Fig. 8(b) that the variance of  $r_X$  is the lowest and the average of  $r_X$  with respect to AFL is the highest, indicating the superiority of  $r_X$  over  $r_P$  and  $r_S$ . Such advantage of  $r_X$  can be quantitatively confirmed by the sensitivity ratios of  $r_X$ ,  $r_P$ , and  $r_S$ , being 11, 10, and 6, respectively.

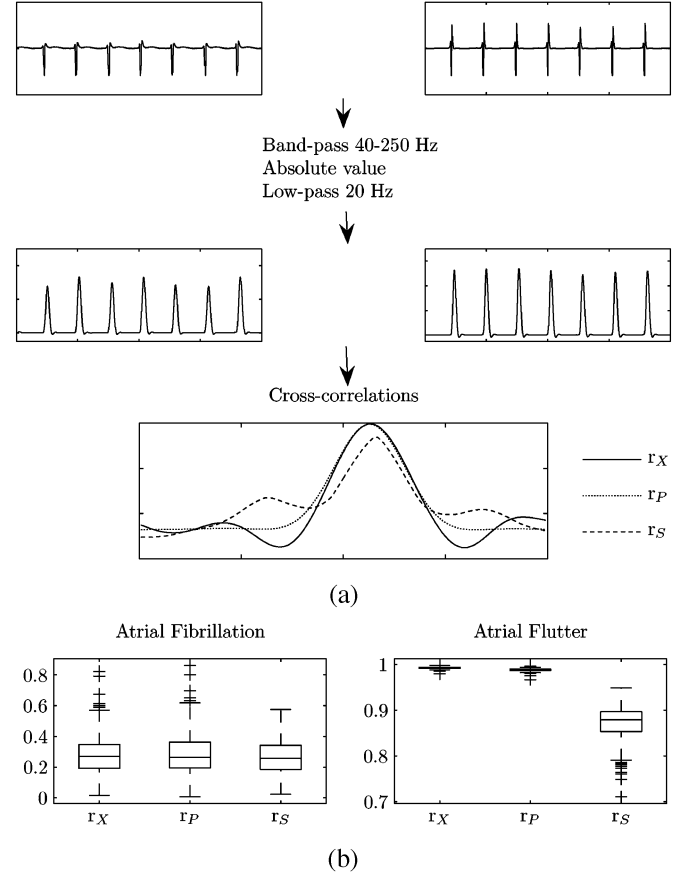


Fig. 8. Correlation analysis technique and comparison results among  $r_X$ ,  $r_P$ , and  $r_S$ . (a) Demonstration of the filtering and correlation determination between a pair of simultaneous signal segments. In a parallel manner, the two data segments are first preprocessed by 1) third-order Butterworth band-pass filtering; 2) absolute valuing; and 3) low-pass filtering. Three maximal values are then extracted as discriminatory indices with respect to three cross-correlation functions  $r_X(k)$ ,  $r_P(k)$ , and  $r_S(k)$  of time-shift  $-100 \leq k \leq 100$ . (b) Statistical results of  $r_X(k)$ ,  $r_P(k)$ , and  $r_S(k)$ . For AF data, the indices congregate near zero; whereas for AFL data, the indices congregate near unity. The larger variance of  $r_S$  for AFL data indicates the low discriminatory power of  $r_S$  in this case.

#### F. Summary of Main Advantages of Order Statistics Correlation Coefficient

The numerical results presented allow us to claim the following advantages of  $r_X$  as a method of quantification of association between biosignals.

- 1) *Noise Robustness*: The index  $r_X$  decreases slowly as the noise strength increases.
- 2) *Small Biasedness*: Although  $r_X$  is not an unbiased estimator of linear association, the biasedness is very small compared to rank correlations  $r_S$  and  $r_K$ . In this aspect,  $r_P$  is optimal, whereas  $r_S$  and  $r_K$  have limited power in measuring linear associations of spiky biosignals (see Fig. 3).
- 3) *High Sensitivity to Changes in Association*: Under linear models,  $r_X$  has sensitivity to changes in  $\rho$  similar to that of  $r_P$  and much higher than those of  $r_S$  and  $r_K$ .
- 4) *High Accuracy for Time Delay Detection*: The index  $r_X$  has almost perfect performance to detect time delays between two biosignals.

- 5) *Fast Computational Speed*: The computational load of  $r_X$  is relatively light, being much faster than  $r_S$  and  $r_K$  and a little slower than  $r_P$ .
- 6) *Good Performance for Nonlinear Association Estimation*: Unlike  $r_S$  and  $r_K$ ,  $r_X$  is not invariant under increasing nonlinear transforms, but it always performs substantially better than  $r_P$  and for spiky signals even better than  $r_S$  and  $r_K$ .

## VI. CONCLUSION

In this paper, we propose a new order statistics correlation coefficient and investigate its properties and applicability to biosignals. The proposed measure was evaluated using simulated and real biosignals and four models emulating linear and nonlinear situations. We also compared the behavior of our measure with three other correlation coefficients commonly used in the literature. The comparative studies demonstrate that our new measure  $r_X$  plays the role of a “missing link” between Pearson’s coefficient and Spearman’s  $\rho$  and Kendall’s  $\tau$ . It enjoys the advantages of all the other three coefficients. In most cases,  $r_X$  is not optimal, but it usually is the second best compared to  $r_P$ ,  $r_S$ , and  $r_K$ . This suboptimal feature at least avoids the worst results in practice when one has no prior knowledge as to whether nonlinearity exists in the system. The new method can be applied to a wide spectrum of biosignal processing, such as organizational indexing of atrial fibrillation [22], [23], atrial fibrillation detection [19], EEG association analysis [7]–[9], etc. In fact, the proposed measure can be used in all the fields where the other three classical methods are applicable, although our comparative studies are conducted in the context of biosignal processing.

## REFERENCES

- [1] J. D. Gibbons and S. Chakraborti, *Nonparametric Statistical Inference*, 3rd ed. New York: Marcel Dekker, 1992.
- [2] R. A. Fisher, *Statistical Methods, Experimental Design, and Scientific Inference*. New York: Oxford Univ. Press, 1990.
- [3] R. A. Fisher, “On the ‘probable error’ of a coefficient of correlation deduced from a small sample,” *Metron*, vol. 1, pp. 3–32, 1921.
- [4] E. C. Fieller, H. O. Hartley, and E. S. Pearson, “Test for rank correlation coefficients. I,” *Biometrika*, vol. 44, pp. 470–481, 1957.
- [5] M. Kendall and J. D. Gibbons, *Rank Correlation Methods*, 5th ed. New York: Oxford Univ. Press, 1990.
- [6] D. D. Mari and S. Kotz, *Correlation and Dependence*. London, U.K.: Imperial College Press, 2001.
- [7] N. J. I. Mars and G. W. van Arragon, “Time delay estimation in nonlinear systems using average amount of mutual information analysis,” *Signal Process.*, vol. 4, pp. 139–153, 1982.
- [8] N. J. I. Mars and F. H. L. da Silva, “Propagation of seizure activity in kindled dogs,” *Electroencephalography Clinical Neurophys.*, vol. 56, pp. 194–209, 1983.
- [9] V. M. Fernandes de Lima, J. P. Pijn, C. N. Filipe, and F. H. L. da Silva, “The role of hippocampal commissures in the interhemispheric transfer of epileptiform afterdischarges in the rat: A study using linear and non-linear regression analysis,” *Electroencephalography Clinical Neurophys.*, vol. 76, pp. 520–539, 1990.
- [10] J. P. S. Cunha and P. G. de Oliveira, “A new and fast nonlinear method for association analysis of biosignals,” *IEEE Trans. Biomed. Eng.*, vol. 47, pp. 757–763, 2000.
- [11] H. A. David and H. N. Nagaraja, *Order Statistics*, 3rd ed. Hoboken, NJ: Wiley-Interscience, 2003.
- [12] N. Balakrishnan and C. R. Rao, *Order Statistics: Applications*. New York: Elsevier, 1998.
- [13] N. Balakrishnan and C. R. Rao, *Order Statistics: Theory & Methods*. New York: Elsevier, 1998.
- [14] W. Xu, C. Chang, Y. S. Hung, S. K. Kwan, and P. C. W. Fung, “Order statistic correlation coefficient and its application to association measurement of biosignals,” in *Proc. ICASSP*, 2006, pp. II-1068–II-1071.
- [15] D. S. Mitrinovic, J. E. Pecaric, and A. M. Fink, *Classical and New Inequalities in Analysis*. Norwell, MA: Kluwer Academic, 1993.
- [16] M. G. Kendall, A. Stuart, J. K. Ord, and S. F. Arnold, *Kendall’s Advanced Theory of Statistics: Volume 1 Distribution Theory*, 6th ed. London, U.K.: Edward Arnold, 1994.
- [17] B. C. Arnold, N. Balakrishnan, and H. N. Nagaraja, *A First Course in Order Statistics*. New York: Wiley, 1992.
- [18] H. Ruben, “On the sum of squares of normal scores,” *Biometrika*, vol. 43, pp. 456–458, 1956.
- [19] A. Cohen, *Biomedical Signal Processing*. Boca Raton, FL: CRC Press, 1986.
- [20] J. H. v. Bommel, M. A. Musen, and J. C. Helder, *Handbook of Medical Informatics*. AW Houten, The Netherlands: Bohn Stafleu Van Loghum, 1997.
- [21] Z. Chen, P. C. Ivanov, K. Hu, and H. E. Stanley, “Effect of nonstationarities on detrended fluctuation analysis,” *Phys. Rev. E*, vol. 65, p. 041107, 2002.
- [22] J. Gao, J. Hu, W.-W. Tung, Y. Cao, N. Sarshar, and V. P. Roychowdhury, “Assessment of long-range correlation in time series: How to avoid pitfalls,” *Phys. Rev. E*, vol. 73, p. 016117, 2006.
- [23] W. Xu, H. F. Tse, P. C. W. Fung, F. H. Y. Chan, K. L. F. Lee, and C. P. Lau, “New Bayesian discriminator for detection of atrial tachyarrhythmias,” *Circulation*, vol. 105, pp. 1472–1479, 2002.
- [24] A. Nygren, C. Fiset, L. Firek, J. W. Clark, D. S. Lindblad, R. B. Clark, and W. R. Giles, “Mathematical model of an adult human atrial cell: The role of K<sup>+</sup> currents in repolarization,” *Circulation Res.*, vol. 82, pp. 63–81, 1998.
- [25] N. V. Thakor and S. Tong, “Advances in quantitative electroencephalogram analysis methods,” *Annu. Rev. Biomed. Eng.*, vol. 6, pp. 453–495, 2004.
- [26] V. Bostanov, “BCI competition 2003-data sets Ib and IIb: Feature extraction from event-related brain potentials with the continuous wavelet transform and the t-value scalogram,” *IEEE Trans. Biomed. Eng.*, vol. 51, pp. 1057–1061, 2004.
- [27] B. D. Mensh, J. Werfel, and H. S. Seung, “BCI competition 2003-data set Ia: Combining gamma-band power with slow cortical potentials to improve single-trial classification of electroencephalographic signals,” *IEEE Trans. Biomed. Eng.*, vol. 51, pp. 1052–1056, 2004.
- [28] K. Hu, P. C. Ivanov, Z. Chen, P. Carpena, and H. E. Stanley, “Effect of trends on detrended fluctuation analysis,” *Phys. Rev. E*, vol. 64, p. 011114, 2001.
- [29] A. L. Goldberger, L. A. N. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, “PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals,” *Circulation*, vol. 101, pp. e215–e220, 2000.
- [30] W. Xu, C. Chang, Y. S. Hung, and P. C. W. Fung, “Asymptotic properties of order statistics correlation coefficient in the normal cases,” *IEEE Trans. Signal Process.*, accepted for publication.
- [31] G. W. Botteron and J. M. Smith, “A technique for measurement of the extent of spatial organization of atrial activation during atrial fibrillation in the intact human heart,” *IEEE Trans. Biomed. Eng.*, vol. 42, pp. 579–586, 1995.
- [32] G. W. Botteron and J. M. Smith, “Quantitative assessment of the spatial organization of atrial fibrillation in the intact human heart,” *Circulation*, vol. 93, pp. 513–518, 1996.



**Weichao Xu** (M’06) received the B.Eng. and M.Eng. degrees in electrical engineering from the University of Science and Technology of China, Hefei, China, in 1993 and 1996, respectively, and the Ph.D. degree in biomedical engineering from the University of Hong Kong, Hong Kong, in 2002.

Since 2003, he has been a Research Associate with the Department of Electrical and Electronic Engineering, the University of Hong Kong. His research interests include the areas of mathematical statistics, pattern recognition, digital signal processing, and applications.



**Chunqi Chang** (M'06) received the B.Sc. and M.Sc. degrees in electronic engineering from the University of Science and Technology of China, Hefei, China, in 1992 and 1995, respectively, and the Ph.D. degree in biomedical engineering from the University of Hong Kong, Hong Kong, in 2001.

He is currently a Research Fellow with the Department of Electrical and Electronic Engineering, the University of Hong Kong. His main research interests include statistical signal processing theories and methods and their applications to biomedical engineering and computational molecular biology.



**Y. S. Hung** (M'88–SM'02) received the B.Sc. (Eng.) degree in electrical engineering and the B.Sc. degree in mathematics from the University of Hong Kong, Hong Kong, and the M.Phil. and Ph.D. degrees from the University of Cambridge, Cambridge, U.K.

He was a Research Associate with the University of Cambridge and a Lecturer with the University of Surrey, Surrey, U.K. In 1989, he joined the University of Hong Kong, where he is currently an Associate Professor and the Head of the department.

His research interests include robust control systems theory, robotics, computer vision, and biomedical engineering.

Dr. Hung was a recipient of the Best Teaching Award in 1991 from the Hong Kong University Students' Union. He is a chartered engineer and a fellow of IET and HKIE.



**S. K. Kwan** received the B.Sc. (Eng.), M.Sc. (Eng.), and M.Stat. degrees from The University of Hong Kong, Hong Kong, in 1990, 1997, and 2000, respectively, where he currently pursuing the Ph.D. degree in biomedical research.

He is a Biomedical Engineer with Hong Kong Hospital Authority, Hong Kong. His research interests focus on developing a high performance implantable cardioverter defibrillator based on statistical analysis of ECG.



**Peter Chin Wan Fung** received the B.Sc. (Phy.) degree, the B.Sc. degree (special hons.) in radio astronomy, and the Ph.D. degree in radio astronomy from The University of Tasmania, Hobart, Australia.

In 1999, he became the first Chair Professor of Medical Physics with the Department of Medicine, The University of Hong Kong. He is also an Honorary Professor in the Department of Psychiatry and the Department of Electrical and Electronic Engineering, The University of Hong Kong, where he works on multidisciplinary projects. He has

been with the University of Montreal, Montreal, QC, Canada, and Stanford University, Stanford, CA. During the early 1970s, he joined the University of Hong Kong, Hong Kong, and became the Personal (Chair) Professor of Physics in 1984 and the Director of the Centre for Materials Science in 1992. He has published around 280 articles in international reviewed journals. His current research interests include the areas of biophysics, medicine, and signal processing.