



Published in final edited form as:

Physiol Meas. 2016 August ; 37(8): E5–E23. doi:10.1088/0967-3334/37/8/E5.

Editorial: False Alarm Reduction in Critical Care

Gari D Clifford^{1,2}, Ikaro Silva³, Benjamin Moody³, Qiao Li¹, Danesh Kella¹, Abdullah Chahin⁴, Tristan Kooistra⁴, Diane Perry⁴, and Roger G. Mark³

¹ Department of Biomedical Informatics, Emory University, Atlanta, GA USA

² Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA, USA

³ Institute for Medical Engineering & Science, Massachusetts Institute of Technology, USA

⁴ Beth Israel Medical Center, Harvard University, Boston MA, USA

Abstract

High false alarm rates in the ICU decrease quality of care by slowing staff response times while increasing patient delirium through noise pollution. The 2015 Physio-Net/Computing in Cardiology Challenge provides a set of 1,250 multi-parameter ICU data segments associated with critical arrhythmia alarms, and challenges the general research community to address the issue of false alarm suppression using all available signals. Each data segment was 5 minutes long (for real time analysis), ending at the time of the alarm. For retrospective analysis, we provided a further 30 seconds of data after the alarm was triggered.

A total of 750 data segments were made available for training and 500 were held back for testing. Each alarm was reviewed by expert annotators, at least two of whom agreed that the alarm was either true or false. Challenge participants were invited to submit a complete, working algorithm to distinguish true from false alarms, and received a score based on their program's performance on the hidden test set. This score was based on the percentage of alarms correct, but with a penalty that weights the suppression of true alarms five times more heavily than acceptance of false alarms.

We provided three example entries based on well-known, open source signal processing algorithms, to serve as a basis for comparison and as a starting point for participants to develop their own code. A total of 38 teams submitted a total of 215 entries in this year's Challenge.

This editorial reviews the background issues for this Challenge, the design of the Challenge itself, the key achievements, and the follow-up research generated as a result of the Challenge, published in the concurrent special issue of *Physiological Measurement*. Additionally we make some recommendations for future changes in the field of patient monitoring as a result of the Challenge.

1. Introduction

During the last decade, over a period of seven years, Intensive Care Unit (ICU) admissions at U.S. hospitals increased by 48.8% with a mean biennial increase of 14.2%. By

comparison, overall emergency department (ED) visits increased by 5.8% biennially. In absolute terms, admissions jumped from 2.79 million in 2002-2003 to 4.14 million in 2008-2009, according to data from the National Hospital Ambulatory Care Survey Mullins et al. (2013a). The three most common diagnoses for ICU admissions were unspecified chest pain, congestive heart failure, and pneumonia. Utilization rates of most tests and services delivered to patients admitted to the ICU from the ED increased, with the largest increase occurring in computed tomography (CT) and magnetic resonance imaging (MRI), which increased from 16.8% in 2002/2003 to 37.4% in 2008/2009, a 6.9% mean biennial increase. These findings suggested emergency physicians were sending more patients on to the ICU. The increase might be the result of an older, sicker population that needs more care Mullins et al. (2013a).

ICU patients require a high level of acute care, with numerous bedside monitors which are continuously monitoring both invasive and non-invasive variables. These monitors provide synchronous waveforms with both independent and complementary information. Huge ICU databases are therefore becoming available, and include parameters such as the electrocardiogram (ECG), the photoplethysmogram (PPG), the arterial blood pressure (ABP) waveform and respiratory effort. In clinical practice these signals are processed individually and derived parameters are frequently set to trigger an alarm when a specific parameter (such as heart rate) exceeds a pre-defined range. These alarms are frequently false alarms (FAs) and account for a large majority of all alarms generated in the ICU Chambrin et al. (1999).

Furthermore, the high rate of false alarms significantly burdens clinical staff, which can lead to decreased quality of care Donchin and Seagull (2002); Imho and Kuhls (2006), impacting both the patient and the clinical staff through noise disturbances, desensitization to warnings, slowing of response times Chambrin (2001) and missed true alarms Allen and Murray (1996); Chambrin (2001); Hug et al. (2011). ICU alarms produce sound intensities above 80 dB that can lead to sleep deprivation Chambrin (2001); Meyer et al. (1994); Parthasarathy and Tobin (2004), inferior sleep structure Johnson (2001); Slevin et al. (2000), stress for both patients and staff Cropp et al. (1994); Novaes et al. (1997); Topf and Thompson (2001); Morrison et al. (2003) and depressed immune systems Berg (2001). There are also indications that the incidence of re-hospitalization is lower if disruptive noise levels are decreased during a patient's stay Hagerman et al. (2005). Furthermore, such disruptions have been shown to have an important effect on recovery and length of stay Donchin and Seagull (2002); Cropp et al. (1994). In particular, cortisol levels have been shown to be elevated (reflecting increased stress) Topf and Thompson (2001); Morrison et al. (2003), and sleep disruption has been shown to lead to longer stays in the ICU Parthasarathy and Tobin (2004). ICU false alarm (FA) rates as high as 90% Aboukhalil et al. (2008) have been reported, with between 6% and 40% of ICU alarms having been shown to be true but clinically insignificant (requiring no immediate action) Lawless (1994). In fact, only 2% to 9% of alarms have been found to be important for patient management Tsien and Fackler (1997). In response to this, thresholds and filter settings for alarms are often manipulated on a case-by-case basis in response to an individual clinical user's preferences (to reduce annoyance), which may well be sub optimal in terms of the trade off between true and false alarms Mullins et al. (2013b).

In the 2015 PhysioNet/Computing in Cardiology Challenge Clifford et al. (2015) (<http://physionet.org/challenge/2015>), we aimed to address the problem of high false alarm rates by encouraging the development of new algorithms to improve the specificity of ICU alarms. In this Challenge, we focused on five types of life-threatening arrhythmia events, which we defined as follows:

Asystole (ASY): No heartbeats for a period of four seconds or more.

Extreme bradycardia (EBR): Heart rate lower than 40 beats per minute; fewer than five beats occur within a period of six seconds.

Extreme tachycardia (ETC): Heart rate higher than 140 beats per minute; more than 17 beats occur within a period of 6.85 seconds.

Ventricular tachycardia (VTA): Five or more consecutive ventricular beats within a period of 2.4 seconds (a rate of 100 per minute.)

Ventricular fibrillation or flutter (VFB): The heart exhibits a rapid fibrillatory, flutter, or oscillatory waveform for at least four seconds.

Participants in the Challenge were given samples of ICU patient waveforms that were identified by the bedside monitor as falling into one of the above categories, and were tasked with devising an algorithm to determine which of these alarms represented true arrhythmias, and which were caused by other factors (such as noise, patient movement, leads falling off, or mis-identification of ECG features on the part of the monitor.)

The Challenge was divided into two events. Event 1 was a simulation of the *real-time* alarm suppression problem: the algorithm needed to determine whether the alarm was true or false based solely on the information available before the alarm was first triggered. In Event 2, algorithms were also able to see 30 seconds' worth of waveform data following the time of the alarm, and could use this information to *retrospectively* classify the alarm as true or false. The development of an algorithm that could reliably solve either of these problems would be a major step forward in patient care.

2. Example algorithms

Key to rhythm detection is accurate heart rate estimation. Several ECG R-peak detection algorithms are freely available, several of which were used in the Challenge example entries.

eplimited (available at www.eplimited.com) Hamilton and Tompkins (1986), which used digital filtering and a group of decision rules.

sqrs (available on PhysioNet Goldberger et al. (2000)) Engelse and Zeelenberg (1979), which uses a single scan of the sampled data and combines digital filter preprocessing with a detector and feature extractor based on dynamically adjusted slope and timing criteria.

wqrs (available on PhysioNet) Zong, Moody and Jiang (2003), which is based on the length transform.

gqrs (available on PhysioNet), which consists of a QRS matched filter with a custom built set of heuristics (such as search back).

coqrs Clifford (2002); Nygård and Sörnmo (1983); Oster et al. (2013) based on the peak energy (no search back).

jqrs Behar, Johnson, Clifford and Oster (2014); Behar, Oster and Clifford (2014) consists of a window-based peak energy detector but with replacement of the original band-pass filter with a QRS matched filter (Mexican hat) and an additional heuristic ensuring no detection were made during flat lines.

Detection of the onset of the pulses in the ABP and PPG signals can provide further information on rhythm and rate. An open-source algorithm, *wabp* Zong, Heldt, Moody and Mark (2003), is available from PhysioNet. The algorithm consists of three components: 1) a low-pass filter which is to suppress high frequency noise that might affect the onset detection; 2) a windowed and weighted slope sum function (SSF) which is to enhance the up-slope of the pulse and to suppress the remainder of the pressure wave; 3) a decision rule which allows for detection of each SSF pulse onset.

We provided three example Challenge entries, based on these and other open-source algorithms, and implemented in various programming languages, to serve as a basis on which participants could develop their own code.

The simplest example entry (#1) used *wabp* and *gqrs*, along with the *gqfuse* tool (available on PhysioNet), to analyze all available signals and select the most stable sequence of RR intervals, in order to detect asystole, bradycardia, and tachycardia. To detect the onset of VF, this entry analyzed the ECG and pulsatile signals separately (using *gqfuse* for each), and searched for a 10-second interval where the QRS rate and pulse rate were equal, followed by a 3-second interval in which the QRS rate increased by at least 25% and the pulse rate decreased by at least 75%. This entry did not attempt to detect VT.

A second example entry (#2) written in MATLAB used *wabp* to detect the beats and used *jsQI* Sun et al. (2006) and a template matching *SQI* Li and Clifford (2012) to estimate the signal quality from ABP and PPG channels. For the ECG, an agreement level of two R-peak detectors (*gqrs* and *coqrs*) in a 10-second window, evaluated every second, known as *bsQI*, was used Behar et al. (2013).

Finally, the third sample entry (#3) was provided for Octave users, with functions from the WFDB Toolbox for Octave/MATLAB Silva and Moody (2014). This sample entry ran three QRS detectors from the WFDB Toolbox: *wqrs* on signal 1, *sqrs* on signals 1 and 2, and *gqrs* on signals 1 and 2. The results of the QRS detectors were then used to compute three tachograms. A decision was made on the veracity of the alarm based on the average pairwise correlation between the tachograms 30 seconds prior to the alarm (a threshold was set arbitrarily based on the training data).

2.1. Signal Quality

Signal quality indices (SQIs), which assess the signal quality or noise levels of the signals, can be extracted from the waveforms and used as weighting factors to allow for varying trust

levels in the source data. The ECG signal quality has been extensively studied Li et al. (2008); Clifford et al. (2012); Behar et al. (2013); Li et al. (2014b). For the benchmark algorithms, an agreement level of two R-peak detectors in a 10-second window, evaluated every second, known as *bsQI*, was used. Intuitively, the presence of noise and artifacts will lower the agreement level between two semi-independent detectors. The *bsQI* was recently successfully used on a database with pathological rhythms Li et al. (2008); Behar et al. (2013). The ABP signal quality was evaluated using an open-source algorithm Sun et al. (2006) which flags a signal as bad quality if derived parameters from a blood pressure wave are not in reasonable physiological ranges. The PPG signal quality was also evaluated Li and Clifford (2012) which matches a running PPG template with the pulsatile beat by dynamic time warping, simple matching, linear resampling matching and a clipping detection. When the signal quality was equal or greater than 0.9 and the corresponding HR or beat-to-beat interval derived from either the ABP or PPG did not surpass a predefined HR threshold (4s for ASY, 40 bpm for EBR, 140 bpm for ETC, 100 bpm for VTA and 250 bpm for VFB), the alarm was suppressed as a false alarm.

It should be noted that no ECG signal quality metrics were used in our benchmark algorithms, although previous studies using the agreement of beat detectors for signal quality estimation have shown great promise in this area Behar et al. (2013). We also note that no ECG-based rhythm detection was used, although various open source algorithms were made available to the Challenge participants Li et al. (2014a).

2.2. Voting algorithms

We also implemented a voting approach to combine together varying numbers of algorithms. A simple unweighted voting of the N best performing final entries, ranked by their score on the training data (to prevent overfitting on the test scores, was implemented). N was varied from 1 to 37 with tied, absent or no vote was treated as 'true'. In other words, a forward selection approach was used to select which algorithms should be combined.

3. Challenge data

Data for the Challenge consisted of waveform recordings from ICU patients in four hospitals in the USA and Europe, representing three major manufacturers of ICU monitoring equipment. For each arrhythmia alarm matching our selection criteria, we collected all available multi-parameter waveforms (including at least five minutes of data before and after each alarm), as well as the alarm messages themselves, and any other status messages reported by the monitor. If possible, we also collected a list of the fiducial points and types of beats that were detected by the monitor; in some cases, the monitor did not provide this information. All of the signals were filtered in order to remove spectral characteristics that might identify the manufacturer or the country of origin. They were then resampled to 250 Hz and scaled to a 16-bit range. The specific names of the various alarm annotations were also normalized to anonymize the data.

3.1. Expert labeling

To build the “gold standard” list of true and false alarms, a team of experts visually inspected the waveform record at the time of each alarm. Each annotator worked independently and was assigned a randomized list of patients to review. For each alarm, the annotator was initially shown 15 seconds of waveforms prior to the alarm and 5 seconds after it, but could resize and scroll the window in order to examine earlier and later portions of the record. If possible, the monitor-computed beat labels were also displayed.

After examining the alarm label and surrounding waveforms, the annotator was asked to press one of four buttons: *True*, *False*, *Reject*, or *Uncertain*. The *Reject* label was used for records that were clearly fallacious (usually due to bugs in the monitor's data-exporting interface.) In order for an alarm to be included in the Challenge data set, it had to be independently reviewed by at least two annotators of whom a two-thirds majority had to agree that the alarm was either *True* or *False*.

3.2. Training and test data

From the set of 1,564 alarms meeting all of the above criteria, we randomly picked 1,250 to serve as training and test data for the Challenge (see table 1). The distribution of alarms was chosen to reflect the distribution of alarm types in the original data set (17% ASY, 11% EBR, 17% ETC, 47% VTA, 7% VFB) as well as to maintain the approximate true-to-false ratio for each alarm type. No single patient appeared in both the training and test sets, and no single manufacturer or hospital made up more than half of the records in either set.

Up to four signals were selected from each record: two ECG leads and up to two other signals, including ABP, PPG, or respiration. The public training set consisted of 375 “short” records, containing only the five minutes leading up to the alarm, and 375 “long” records, containing a further 30 seconds after the alarm. The hidden test set consisted of 250 “short” records (used only for Event 1) and 250 “long” records (used for both events.) Each record was labeled with the alarm type, and in the case of the training set, whether the alarm was true or false. The records did not include the monitor-computed beat fiducial points or heart rate.

4. Scoring

Participants were asked to submit their entries in the form of a ‘zip’ or ‘tar’ archive that included everything needed to compile and run their program on a GNU/Linux system, together with the results that they expected their program to produce for the records in the public training set. When an entry was uploaded, the scoring system would first attempt to compile the program and run it over a randomly selected subset of the training set; if this did not produce the expected results, evaluation stopped and the error messages were sent back to the submitter.

Once the program was successfully compiled and validated, it was then invoked for each record in the test set. (For the 250 “long” records, the program was invoked twice: once with the full record as input, and once with a truncated version.) If the program failed to produce output for a given record, it was treated as if it had classified that alarm as true.

For each category, the entry's score was computed based on the number of *true positives* (true alarms classified as true), *false positives* (false alarms classified as true), *true negatives*, and *false negatives*. The scoring function was designed to treat false negatives – genuinely life-threatening events that the program considered unimportant – especially harshly, and was defined as:

$$score = \frac{100 \cdot (TP+TN)}{TP+TN+FP+5 \cdot FN}$$

5. Results of the Challenge

A total of 29 closed-source entries and 215 open-source entries were submitted in the Challenge in 2015. Table 3 provides a breakdown of the top scoring entries. A different contestant ranked highest in each separate alarm category, indicating that there was no best general algorithm. Interestingly, a simple majority vote of all the 38 competitors' final entries gave scores of 60.15 in the real-time event and 62.41 in the retrospective event. These moderate performances, well below the top 10 algorithms, indicating that simple voting schemes do not yield an improved performance in this context, since the performance tail is long. A voting algorithm using the N=13 best performing final entries ranked by their score on the training data, provided the highest scores in both event 1 (84.26) and event 2 (87.04), although N=11 was sufficient to beat the best performance in either event. Figure 1 illustrates the performance of a simple voting approach for both the retrospective and prospective parts of the challenge. Note that performance only degrades above 13 algorithms.

6. Review of Articles in the Special Issue

A total of 13 articles were reviewed and revised in time to be accepted for this special issue. Most authors had originally entered the Challenge, and submitted updated versions of their algorithms, which should be made available by the authors through their open source licenses. The top reported results on the hidden test set for each alarm type were: ASY: 97.4% (Plesinger et al. (2016)), EBR: 93.8% (Krasteva et al. (2015a)), ETC: 100.0% (Hoog Antink et al. (2016)), VFB: 88.7% (Rodrigues and Couto (2016)), and VTA: 76.7% (Kalidas and Tamil (2016)), yielding an average best of 91.2%.

Each algorithm published in this issue is reviewed below according to four standard stages of algorithm function:

1. Pre-processing and signal conditioning
2. Beat detection
3. Beat classification
4. Alarm classification

The purpose of this standardized summary is to glimpse at a myriad of advanced approaches used by the competitors in a format that allows the reader to quickly identify both the commonalities and the originality of all the approaches. Finally, the last two articles in this

review (Tsimenidis and Murray (2016); Daluwatte et al. (2016)) did not attempt to reduce the number of false alarms, but rather provide some useful insights into the relationship between signal quality metrics and false alarm rates.

6.1. Ansari et al. (2016)

Ansari et al. (2016) proposed an algorithm that uses several beat detectors within each channel, followed by beat classification, and heuristics to determine the veracity of the alarm. The proposed algorithm operates on 16 seconds of worth of data prior to the alarm. The algorithm achieved a performance accuracy on the final test data-set of ASY: 86.4%, EBR: 79%, ETC: 93.9%, VFB: 61% VTA: 67.6%, yielding a total average of: 76.2%.

Preprocessing—The preprocessing steps consisted of re-sampling the signals to 125 Hz. The ECG signals were band-pass filtered between 0.5-40 Hz, while the pulsatile signals were band-pass filtered between 0.5-10 Hz. Baseline and trend estimation and subtraction was accomplished with a 250 point median filter. The authors also removed pacemaker activity by thresholding on the peak amplitude.

Beat Detection—Ansari et al. (2016) implemented 7 different QRS detectors for each ECG signal, and 3 peak detectors for each of the pressure signals. The fiducial points for all peaks were re-aligned by picking the maximum within 50 ms of the detected beat for ECG signals, and the maximum within 50 ms before or 1 second after the detected beat for the ABP or PPG signals. The outputs of all the 20 beat detectors were then fused by adding their binary outputs (with at least 1 beat under AS, at least 2 for other alarms).

Beat classification—ECG beat classification was performed for the VFB and VTA alarms only. The beat classifier was a decision tree that utilized features derived from the Stockwell Transform on a 200 ms window.

Alarm classification—A decision tree classifier was trained with five fold cross validation in order to determine the veracity of a beat. The final decision regarding the alarm veracity was made based on a set of heuristics.

6.2. Eerikäinen et al. (2016)

Eerikäinen et al. (2016) produced an algorithm that achieved a performance accuracy on the final test data-set of ASY: 89.2%, EBR: 71.5%, ETC: 99.1%, VFB: 81.8% VTA: 68.1%, yielding a total average of: 77.3% and a retrospective average of 81.5%.

Preprocessing—All signals were down-sampled to 125 Hz and the processing window length was optimized for each arrhythmia type (varying from 14 to 16 seconds prior to the alarm). Noise levels were estimated based on the power estimated from the regions in-between beats.

Beat Detection—Beat detection on the ECG waveforms were performed using a QRS detector based on wavelets and auto-regressive modeling of the R-peak Rooijackers et al. (2012). The pulsatile peaks were detected via the open source detector *wabp*.

Alarm classification—A random forest classifier was trained for each of the five different types of alarms. The technique focused on comparing pairs of beats. Two beats were considered a match if they were within 100 ms of each other. Delays across channels were compensated for if the standard deviation of 10 consecutive beats was less than 5% of the mean delay. For the VTA and VFB alarms, only the F1 statistic between ECG leads was used, in addition to spectral purity indexes. An alarm with an F1 equal to zero was identified to be false.

6.3. Fallet et al. (2016)

Fallet et al. (2016) proposed an algorithm that detects beats in the ECG and the pulsatile signals, provided their signal quality is good. The authors also use a spectral purity metric to aid on the classification of VTA and VFB alarms. The algorithm achieved a performance accuracy on the final test data-set was ASY: 84.2%, EBR: 82.4%, ETC: 86.9%, VFB: 87.1% VTA: 72.7%, yielding a total average of: 77.07% and an average on the retroactive category of 85.0%.

Preprocessing—The preprocessing stage for this algorithm consisted of 50 Hz power line noise removal. For the calculations of spectral purity indexes, the signal was down-sampled to 35 Hz and a 5-sample moving average filter was applied. The signal quality for the pulsatile waveforms was estimated through the *ppgSQI* and *jsQI* methods Clifford et al. (2015).

Beat Detection—The QRS component of the ECG signal was detected through a morphological analysis approach with an adaptive approach from Yazdani and Vesin (2014). Beat detection on the pulsatile signals was performed using the algorithm proposed by Arberet et al. (2013). The heart rate time series was then derived through a multi-channel oscillator based adaptive frequency tracking algorithm.

Beat classification—The spectral purity index Sörnmo and Laguna (2005); Goncharova and Barlow (1990) was used a feature to distinguish between normal, ventricular tachycardia, ventricular flutter/fibrillatory arrhythmia (the index was expected to be higher for abnormal rhythms).

Alarm classification—A set of heuristics rules was developed for the final alarm classification. In the case of the ASYS alarm, the algorithm applied majority voting based on the heart rate series from individual ECG and pulsatile channels. The pulsatile channels were only used if the quality was above a certain threshold. A linear discriminant analysis classifier was used for the retrospective event to corroborate the ECG output, but again, only if the pulsatile signal quality was sufficiently high. If the pulsatile quality was low, a set of heuristic thresholds was applied to the minimum heart rate from the last five consecutive beats using 16 seconds before and five seconds after the alarm. The extreme tachycardia alarm only used pulsatile waveforms: if the quality was good, the alarm was checked against the pulsatile rate, else it defaulted to true. Ventricular flutter/fibrillatory alarms were checked through the maximum average spectral purity index calculation over a 3 second window, and no pulsatile information was used. Finally, ventricular tachycardia alarms used a set of

heuristic rules encompassing pulsatile waveform heart-rate series, as well as current versus previous values of the ECG spectral purity indexes.

6.4. Hoog Antink et al. (2016)

The algorithm proposed by Hoog Antink et al. (2016) used 16 seconds of data prior to the alarm event. The algorithm achieved a performance accuracy on the final test data-set of ASY: 76.7%, EBR: 74.2%, ETC: 100%, VFB: 72.8% VTA: 71.5%, yielding a total average of: 78.2% and retrospective average of 74.4%.

Preprocessing—The pre-processing steps for this algorithm included re-sampling of the signals to 100 Hz, band-pass filtering with a pass-band region of 1-30 Hz. The signals were also normalized to zero mean and unit variance using statistics calculated on 5-minutes of data prior to the alarm.

Beat Detection—Beat detection was achieved through the Bayesian fusion of several inter-beat interval estimators that rely on self-similarity: lag adaptive short-time auto-correlation, average magnitude difference function, and maximum amplitude pairs Brüser et al. (2013). A quality metric based on the reliability of the fused estimates was derived from the peak height to area of the fused similarity curve.

Alarm classification—The classifiers chosen for the alarm validation included binary classification trees, regularized linear discriminant analysis, a support vector machine, and a random forest. The authors utilized a combination of both alarm specific and global classifiers (i.e, classifiers trained to detect a general false alarm). Their final choices were linear discriminant analysis for EBR, VFB, and VTA, a binary classifier for ETC, and a random forest model for ASY. A superset of 88 features was developed from: 24 beat-to-beat interval statistics and correlogram analysis of interval time series. From this superset, subsets were selected according to alarm types.

6.5. Kalidas and Tamil (2016)

The algorithm proposed by Kalidas and Tamil (2016) used 10 seconds of data prior to the alarm. The algorithm achieved a performance accuracy on the final test data-set of ASY: 80.7%, EBR: 71.7%, ETC: 99.1%, VFB: 74.1% VTA: 76.7%, yielding a total average of: 79.4% and retrospective average of 80.2%.

Preprocessing—Baseline wander was estimated with a low-pass filter with a 1 Hz cut-off and then subtracted from original signal. Flat line artifact was detected by testing for identical sample values in 2 second windows. ‘Zig-zag’ artifacts were detected by testing for alternating positive and negative slopes in consecutive samples over 2 second periods.

Beat Detection—The Pan and Tompkins (1985) algorithm was used to detect QRS complexes in the ECG signal. Pulsatile peaks were detected through first order differentiation.

Alarm classification—No pulsatile signal information was used for VFB and VTA arrhythmia alarms. For each alarm type, an individual support vector machine and set of heuristics was developed. The features used into these classifiers included the ECG-derived heart rate, and PPG-derived heart rate if morphology was considered valid (excluding the VFB and VTA alarms). The VFB and VTA alarms also included an additional set of features related to the power spectra of the ECG waveforms.

6.6. Krasteva et al. (2016)

The algorithm proposed by Krasteva et al. (2016) used 3-7.5 second windows prior to the alarm event, with the specific duration tuned for the each alarm type. The algorithm achieved a performance accuracy on the final test data-set of ASY: 88.0%, EBR: 93.8%, ETC: 90.7%, VFB: 72.7% VTA: 72.6%, yielding a total average of: 80.0%.

Preprocessing—The ECG channels were fused to form two data streams: a magnitude (second norm) and a velocity (second norm of the first order derivative). The ECG signal quality was estimated using 3 frequency bands on 4s interval windows: high frequency was used to estimate spikes from artifacts and pacemakers, medium frequency range was used to estimate the signal level and power line interference (with intra-beat temporal statistics used to estimate power line noise level), and the low frequency band was used to estimate baseline wander. Pulsatile signals were low-pass filtered with a 1 Hz cut-off. The pulsatile signal quality was estimated with a periodicity index, and mean peak-to-peak amplitude values.

Beat Detection—A nonlinear filtering approach, with adaptively updated upper and lower thresholds, was used for QRS detection. The beat detector had a conventional refractory period of 150 ms.

Beat classification—A beat classifier was developed for supra-ventricular and ventricular ectopic beats. A decision tree model was also used, based on features that included: information from template correlation matching, beat morphology features, and RR statistics Krasteva et al. (2014, 2015a).

Alarm classification—The alarm classification algorithm used a set of heuristic rules based on heart rate, dominant frequency for ventricular rate, phase space area from both the ECG magnitude and velocity, and pulsatile quality metrics.

6.7. Liu et al. (2016)

Liu et al. (2016) proposed an algorithm which processed 60 seconds of data prior to the alarm event. The algorithm achieved a performance accuracy on the final test data-set of ASY: 88.7%, EBR: 77.7%, ETC: 89.9%, VFB: 67.7% VTA: 61.0%, yielding a total average of: 71.6% and retrospective average of 75.9%.

Preprocessing—The ECG and pulsatile signals were band-pass filtered with the pass-band frequency region of 5-40 Hz for the ECGs and a pass-band frequency region of 5-35 Hz for the pulsatile waveforms.

Beat Detection—The authors developed an ECG R wave detection algorithm that used the average maximum amplitude from 6 non-overlapping segments. Pulsatile beats were detected via *wabp*. The final detected beats were validated based on intra- and inter-channel verification of the detected beats along with a set of rules involving the number of detected beats, R amplitude, and distance metrics between the heart rate time series.

Beat classification—A set of heuristics were applied to classify beats. The features included: morphology analysis based on correlation against template, the ratio between changed beats and total beats in segment, QRS width, and maximum heart rate.

Alarm classification—A set of decision rules was applied to channels that passed a data quality check (if the result of the test failed, the alarm was set to false). The features used for the second classification step included number of valid feature points, heart rate, and maximum heart rate at current analysis window.

6.8. Plesinger et al. (2016)

Plesinger et al. (2016) developed an algorithm that used information across multiple channels and sought to detect regions contaminated by artifacts. The algorithm achieved a performance accuracy on the final test data-set was ASY: 97.4%, EBR: 83.5%, ETC: 87.8%, VFB: 80.3% VTA: 75.0%, yielding a total average of: 81.6% and a retrospective average of 84.9%.

Preprocessing—The preprocessing step for this algorithm started with the detection of artifacts based on the temporal statistics of the signal under analysis. Noise and pacemaker activity were estimated based on spectral content of the 50-70Hz frequency band. The pulsatile signals were low passed filtered with cut-off frequency at either 5 or 20 Hz. The following time windows prior to the alarm event were used to process the alarm data: ASY=14s, EBR=16s, ETC=14s, VFB=13s, VTA=10s.

Beat Detection—The ECG QRS detection was based on an analysis of Fourier and Hilbert Transform derived envelopes, with a 110 ms refractory period between detection. Pulsatile based beat detection was evaluated on estimated temporal slope values.

Beat classification—Beat classification was performed using spectral features and descriptive residue statistics over 120 ms and 500 ms windows.

Alarm classification—The alarm classification stage for the ASY, VTA, and VFB alarms included using the count of invalid features obtained during the preprocessing stage described above. Additional features included statistics for the RR series obtained from multiple channels. The sum of the invalid areas had to be zero in order for the algorithm to accept the RR series as a regular rhythm for the specific channel. Finally, a set of heuristic rules was applied based on the derived RR series and the invalid region statistics.

6.9. Rodrigues and Couto (2016)

Rodrigues and Couto (2016) proposed an algorithm that uses two open-source beat detectors on the ECG waveforms as well as *wabp* on the pulsatile signals. The authors also performed beat classification based on the phase of the R wave in the ECG signals. The algorithm achieved a performance accuracy on the final test data-set of ASY: 83.6%, EBR: 71.4%, ETC: 99.1%, VFB: 88.7% VTA: 61.4%, yielding a total average of: 74.2% and a retrospective average of 74.4%.

Preprocessing—All signals were re-sampled to 125 Hz, and the ECG waveforms were processed for pacemaker detection and removal. Baseline noise was removed by first estimating it with a 125 sample median filter, followed by subtraction from the original signal. Flat signal regions were identified by thresholding on low variance over 2 second windows.

Beat Detection—ECG QRS detection was performed using *gqrs* and *osea* software packages Hamilton (2002). The beats on the pulsatile signals were detected with the *wabp* software. The authors developed their own specific beat detectors for ventricular fibrillation beats by fitting a parabola on 125 ms windows. Following the method of Li et al. (2008), a quality index was developed based on the fraction of matched beats from *gqrs* and the *osea* software packages Hamilton (2002) on the ECG channels.

For pulsatile signals, the quality was estimated using the morphology of consecutive beats estimated from correlation and dynamic time warp analysis, per Li and Clifford (2012). The detected beats were fused based on quality indexes and a tolerance window of 150 ms. Pulsatile beats were compensated with a delay estimated from initial detections.

Beat classification—Beat classification was based on a set of heuristics modified from the *osea* software package Hamilton (2002). These set of rules included statistics derived from inter-beat interval and QRS duration. Rodrigues and Couto (2016) also developed a four-category feature, termed ‘polarity’ that characterized the different types of phases of the R wave into: positive, negative, positive-negative, negative-positive (the last two representing biphasic R waves).

Alarm classification—Alarm classification was calculated from a set of decision rules based on signal quality, but with priority weight given to ECG signals.

6.10. Sadr et al. (2016)

Sadr et al. (2016) proposed an algorithm that uses features and processing specific arrhythmias being tested. The algorithm achieved a performance accuracy on the final test data-set was ASY: 82.4%, EBR: 71.13%, ETC: 99.1%, VFB: 65.5% VTA: 68.0%, yielding a total average of: 69.9% and a total average on the retrospective event of 74.0%.

Preprocessing—Baseline removal was performed by first estimating the baseline component through median filtering and then subtracting this baseline component from the original signal.

Beat Detection—A Hilbert Transform based QRS detector based used for estimating the ECG beats Benitez et al. (2001). The *wabp* algorithm was used to detect the peaks on the ABP and PPG waveforms, and a quantile algorithm was also used to locate peaks on the PPG waveform.

Alarm classification—The alarm verification was performed on a 16 second window of data prior to the alarm. For all of the alarms with the exception of VTA, the alarm data streams had to pass four signal quality checks in order to be deemed a true alarm, otherwise they were tagged as being false. Pulsatile signal information was not used for the ETC and VTA alarms. The classification also consisted of decision trees based on several extracted features customized to each alarm type, including: threshold crossing intervals, auto-correlation function values, complexity measures, and QRS template parameters.

6.11. Zong et al. (2016)

Zong et al. (2016) is unique in that it proposed an algorithm based on pulsatile waveform features. The algorithm was developed and tested using the MIMIC II database Saeed et al. (2011) rather than the Challenge data, and was not open sourced.

Preprocessing—Pulsatile signals were low-pass filtered with cuto set to 16 Hz, and a signal quality estimate was obtained using the technique described in Zong et al. (2004).

Beat Detection—Beat detection was performed with the pulsatile signals using *wabp* and with a forced detection after a period of 2 seconds from the last detected pulse.

Beat classification—The pulsatile beats were classified based on the abnormality index from Sun et al. (2006)

Alarm classification—The alarm classification was achieved using features from pulsatile signals that included: pulse-to-pulse interval, amplitude, maximum slope, signal quality, and rhythm. The classifier was developed based on set of heuristic rules specific to each alarm type.

6.12. Daluwatte et al. (2016)

The focus of this article was on developing a better understanding between signal quality and false alarms. The authors developed arrhythmia specific quality indexes, and investigated if existing quality indexes can distinguish between true alarms versus noise. Two humans annotated each ECG signal 10s prior to the alarm as either of high or low quality. Disagreements were not included in the analysis. The authors used 18 signal quality indexes from existing literature, and selected the top three algorithms from ROC analysis on the manually annotated data-set. The ECG beats were detected using the U3 transform Paoletti and Marchesi (2006).

6.13. Tsimenidis and Murray (2016)

The article by Tsimenidis and Murray (2016) investigated the relationship between ECG quality and false alarms. The authors investigated the signal quality of ECG leads 8 seconds

prior to the alarm event. They broke the analysis down into three frequency bands: low frequency from 0.1-1 Hz, mid frequency from 10-20 Hz, and high frequency from 20-40 Hz. The ECG's major power spectrum component was expected to be located mostly from 1-10 Hz. The authors report that the power on all the three frequencies was significantly greater for a false alarm versus a true alarm.

7. Conclusions

In summary, the PhysioNet/Computing in Cardiology Challenge 2015 provided several key additions to the field of false alarm suppression in critical care. First we note that for the top performing entrants, it was the VTA alarm that proved the hardest to classify accurately. This is partly because, at low heart rates, the signal becomes 'more normal looking' in the other signals. This was previously noted in Aboukhalil et al. (2008). To-date, there has been little to address this issue, although the current challenge has made a move towards this. Second, we note that retrospective scores were generally higher than 'real-time' scores, with the highest performing retrospective approach only suppressing 1% of the true alarms, while 80% of the false alarms were suppressed. Although debatable, this may be acceptable as a clinical algorithm if a 30 second window were acceptable. We suggest that this may spur a re-consideration of the AAMI guidelines for maximum alarm latency. Third, we note that voting algorithms together can produce superior results to even the best algorithm. Such an approach can also lead to a more robust implementation, although it may be significantly more computationally intensive. It is also important to note that too many naive voters (more than 13 in the case of the 2015 Challenge) can reduce the accuracy of the label or answer. In Zhu et al. (2014) and Zhu et al. (2015) a voting system for algorithm (and human) annotations of physiological data was described, which incorporates both the physiology and the individual annotator's accuracy as a function of objective features (such as signal quality) to produce a weighted voting scheme to guarantee that all voters added extra information. We suggest that such approaches will become ever more important as computational power becomes increasingly less expensive. We also note that this means that all competitors in the Challenge added something to the final answer!

Finally we note some limitations of the competition. A larger database is needed with more patients, longer recordings, more leads and abnormalities (such as arrhythmias). We intend to work with industry and researchers alike to enhance the Challenge database in all these areas and would be grateful for continued contributions of data and source code, which we will post together with all the open source algorithms and annotated data from the 2015 PhysioNet/Computing in Cardiology Challenge. The latter can be found on PhysioNet's website at <http://physionet.org/challenge/2015>.

Acknowledgments

This work was funded in part by the National Institutes of Health, grant R01-GM104987, and by the National Institute of General Medical Sciences, under NIH cooperative agreement U01-EB-008577 and NIH grant R01-EB-001659. We also would like to thank Arvind Ananthan, Jaka Cikač, Barbara Drew, Quan Ding, Scott Eaton, Xiao Hu, Franc Jager, Cadathur 'Raj' Rajagopalan and Anže Rezelj, as well as GE Healthcare, MathWorks, Mindray North America and Philips Healthcare for their valuable assistance. Finally, we thank all the competitors themselves, without whom there would be no competition.

References

- Aboukhalil A, Nielsen L, Saeed M, Mark RG, Clifford GD. Reducing False Alarm Rates for Critical Arrhythmias using the Arterial Blood Pressure Waveform. *Jour. of Biomed. Info.* 2008; 41(3):442–451.
- Allen J, Murray A. Assessing ECG Signal Quality on a Coronary Care Unit. *Phys. Meas.* 1996; 17(4): 249.
- Ansari S, Belle A, Ghanbari H, Salamango M, Najarian K. Suppression of False Arrhythmia Alarms in the ICU: A Machine Learning Approach. Accepted in *Phys. Meas.* 2016
- Ansari, S.; Belle, A.; Najarian, K. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. Multi-modal integrated approach towards reducing false arrhythmia alarms during continuous patient monitoring: The Physionet Challenge 2015; p. 1181-1184.
- Arberet, S.; Lemay, M.; Renevey, P.; Sola, J.; Grossenbacher, O.; Andries, D.; Sartori, C.; Bertschi, M. Computing in Cardiology Conference (CinC), 2013. IEEE; 2013. Photoplethysmography-based ambulatory heartbeat monitoring embedded into a dedicated bracelet; p. 935-938.
- Behar J, Johnson A, Clifford GD, Oster J. A Comparison of Single Channel Fetal ECG Extraction Methods. *Ann. of Biomed. Engin.* 2014; 42(6):1340–1353.
- Behar J, Oster J, Clifford GD. Combining and Benchmarking Methods of Foetal ECG Extraction without Maternal or Scalp Electrode Data. *Phys. Measur.* 2014; 35(8):1569.
- Behar J, Oster J, Li Q, Clifford GD. ECG Signal Quality during Arrhythmia and its Application to False Alarm Reduction. *IEEE Trans. on Biomed. Engin.* 2013; 60:1660–1666.
- Benitez D, Gaydecki P, Zaidi A, Fitzpatrick A. The use of the Hilbert transform in ECG signal analysis. *Computers in biology and medicine.* 2001; 31(5):399–406. [PubMed: 11535204]
- Berg S. Impact of reduced reverberation time on sound-induced arousals during sleep. *Sleep.* 2001; 24(3):289–292. [PubMed: 11322711]
- Brüser C, Winter S, Leonhardt S. Robust inter-beat interval estimation in cardiac vibration signals. *Physiological Measurement.* 2013; 34(2):123. [PubMed: 23343518]
- Chambrin MC. Alarms in the Intensive Care Unit: How can the Number of False Alarms be Reduced? *Critical Care.* 2001; 5(4):184–188. [PubMed: 11511330]
- Chambrin MC, Ravoux P, Calvelo-Aros D, Jaborska A, Chopin C, Boniface B. Multicentric Study of Monitoring Alarms in the Adult Intensive Care Unit (ICU): a Descriptive Analysis. *Intensive Care Medicine.* 1999; 25(12):1360–1366. [PubMed: 10660842]
- Clifford, GD. PhD thesis. University of Oxford; 2002. Signal Processing Methods for Heart Rate Variability.
- Clifford GD, Behar J, Li Q, Rezek I. Signal quality indices and data fusion for determining clinical acceptability of electrocardiograms. *Physiol Meas.* 2012; 33(9):1419–1433. [PubMed: 22902749]
- Clifford, GD.; Silva, I.; Moody, B.; Li, Q.; Kella, D.; Shahin, A.; Kooistra, T.; Perry, D.; Mark, RG. 2015 Computing in Cardiology Conference (CinC). IEEE; 2015. The PhysioNet/Computing in Cardiology Challenge 2015: Reducing false arrhythmia alarms in the ICU; p. 273-276.
- Couto, P.; Ramalho, R.; Rodrigues, R. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. Suppression of false arrhythmia alarms using ECG and pulsatile waveforms; p. 749-752.
- Cropp AJ, Woods LA, Raney D, Bredle DL. Name that tone. The proliferation of alarms in the intensive care unit. *Chest.* 1994; 105(4):1217–1220. [PubMed: 8162752]
- Daluwatte C, Johannesen L, Galeotti L, Vicente J, Strauss DG, Scully CG. Assessing ECG signal quality indices to discriminate ECGs with artefacts from pathologically different arrhythmic ECGs. Accepted in *Phys. Meas.* 2016
- Donchin Y, Seagull FJ. The hostile environment of the intensive care unit. *Curr Opin Crit Care.* 2002; 8(4):316–320. [PubMed: 12386492]
- Erikäinen, LM.; Vanschoren, J.; Rooijackers, MJ.; Vullings, R.; Aarts, RM. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. Decreasing the false alarm rate of arrhythmias in intensive care using a machine learning approach; p. 293-296.
- Erikäinen LM, Vanschoren J, Rooijackers MJ, Vullings R, Aarts RM. Reduction of false arrhythmia alarms using signal selection and machine learning. Accepted in *Phys. Meas.* 2016

- Engelse WAH, Zeelenberg C. A single scan algorithm for QRS-detection and feature extraction. *Comput in Cardiol.* 1979; 6:37–42.
- Fallet, S.; Yazdani, S.; Vesin, JM. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. A multimodal approach to reduce false arrhythmia alarms in the intensive care unit; p. 277-280.
- Fallet S, Yazdani S, Vesin JM. False Arrhythmia Alarms Reduction in the Intensive Care Unit: A Multimodal Approach. Accepted in *Phys. Meas.* 2016
- Goldberger AL, Amaral LA, Glass L, Hausdorff JM, Ivanov PC, Mark RG, Mietus JE, Moody GB, Peng CK, Stanley HE. Physiobank, Physiokit, and Physionet Components of a New Research Resource for Complex Physiologic Signals. *Circ.* 2000; 101(23):e215–e220.
- Goncharova II, Barlow JS. Changes in eeg mean frequency and spectral purity during spontaneous alpha blocking. *Electroencephalography and Clinical Neurophysiology.* 1990; 76(3):197–204. [PubMed: 1697252]
- Hagerman I, Rasmanis G, Blomkvist V, Ulrich R, Eriksen CA, Theorell T. Influence of intensive coronary care acoustics on the quality of care and physiological state of patients. *Int J Cardiol.* 2005; 98(2):267–270. [PubMed: 15686777]
- Hamilton, P. *Computers in Cardiology*, 2002. IEEE; 2002. Open source ECG analysis; p. 101-104.
- Hamilton PS, Tompkins WJ. A real-time QRS detection algorithm. *IEEE Trans Biomed Eng.* 1986; 33(12):1157–1165. [PubMed: 3817849]
- Hoog Antink CB, Leonhardt S, Walter M. Reducing False Alarms in the ICU by Quantifying Self-Similarity of Multimodal Biosignals. Accepted in *Phys. Meas.* 2016
- Hoog Antink, C.; Leonhardt, S. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. Reducing false arrhythmia alarms using robust interval estimation and machine learning; p. 285-288.
- Hug CW, Clifford GD, Reisner AT. Clinician Blood Pressure Documentation of Stable Intensive Care Patients: an Intelligent Archiving Agent has a Higher Association with Future Hypotension. *Critical Care Medicine.* 2011; 39(5):1006. [PubMed: 21336136]
- Imhoff M, Kuhls S. Alarm algorithms in critical care monitoring. *Anesth Analg.* 2006; 102(5):1525–1537. [PubMed: 16632837]
- Johnson AN. Neonatal response to control of noise inside the incubator. *Pediatr Nurs.* 2001; 27(6): 600–605. [PubMed: 12024534]
- Kalidas, V.; Tamil, L. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. Enhancing accuracy of arrhythmia classification by combining logical and machine learning techniques; p. 733-736.
- Kalidas V, Tamil L. Cardiac arrhythmia classification using multi-modal signal analysis. Accepted in *Phys. Meas.* 2016
- Krasteva V, Jekova I, Leber R, Schmid R, Abächerli R. Superiority of classification tree versus cluster, fuzzy and discriminant models in a heartbeat classification system. *PloS one.* 2015a; 10(10):e0140123. [PubMed: 26461492]
- Krasteva, V.; Jekova, I.; Leber, R.; Schmid, R.; Abächerli, R. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015b. Validation of arrhythmia detection library on bedside monitor data for triggering alarms in intensive care; p. 737-740.
- Krasteva V, Jekova I, Leber R, Schmid R, Abächerli R. Real-time arrhythmia detection with supplementary ECG quality and pulse wave monitoring for reduction of false alarms in ICU. Accepted in *Phys. Meas.* 2016
- Krasteva, V.; Leber, R.; Jekova, I.; Schmid, R.; Abacherli, R. Computing in Cardiology Conference (CinC), 2014. IEEE; 2014. Classification of supraventricular and ventricular beats by QRS template matching and decision tree; p. 349-352.
- Lawless ST. Crying wolf: False Alarms in a Pediatric Intensive Care Unit. *Critical Care Medicine.* 1994; 22(6):981–985. [PubMed: 8205831]
- Li Q, Clifford GD. Dynamic time warping and machine learning for signal quality assessment of pulsatile signals. *Phys. Meas.* 2012; 33(9):1491–1501.
- Li Q, Mark RG, Clifford GD. Robust Heart Rate Estimation from Multiple Asynchronous Noisy Sources using Signal Quality Indices and a Kalman Filter. *Phys. Meas.* 2008; 29(1):15–32.

- Li Q, Rajagopalan C, Clifford G. Ventricular fibrillation and tachycardia classification using a machine learning approach. *IEEE Transactions on Biomedical Engineering*. 2014a; 61(6):1607–1613. [PubMed: 23899591]
- Li Q, Rajagopalan C, Clifford GD. A machine learning approach to multi-level ECG signal quality classification. *Comput Methods Programs Biomed*. 2014b; 117(3):435–447. [PubMed: 25306242]
- Liu, C.; Zhao, L.; Tang, H. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. Reduction of False Alarms in Intensive Care Unit using Multi-feature Fusion Method; p. 741-744.
- Liu C, Zhao L, Tang H, Li Q, Wei S, Li J. Life-threatening false alarms rejection in ICU: Using rule-based and multi-channel information fusion method. 2016 Accepted in *Phys. Meas*.
- Meyer TJ, Eveloff SE, Bauer MS, Schwartz WA, Hill NS, Millman RP. Adverse environmental conditions in the respiratory and medical ICU settings. *Chest*. 1994; 105(4):1211–1216. [PubMed: 8162751]
- Morrison WE, Haas EC, Shaffner DH, Garrett ES, Fackler JC. Noise, stress, and annoyance in a pediatric intensive care unit. *Crit Care Med*. 2003; 31(1):113–119. [PubMed: 12545003]
- Mullins PM, Goyal M, Pines JM. National growth in intensive care unit admissions from emergency departments in the United States from 2002 to 2009. *Academic Emergency Medicine : Official Journal of the Society for Academic Emergency Medicine*. 2013a; 20(5):479–86. URL: <http://www.ncbi.nlm.nih.gov/pubmed/23672362>. [PubMed: 23672362]
- Mullins PM, Goyal M, Pines JM. National Growth in Intensive Care Unit Admissions From Emergency Departments in the United States from 2002 to 2009. *Acad. Emerg. Med*. 2013b; 20(5):479–486. [PubMed: 23672362]
- Novaes MA, Aronovich A, Ferraz MB, Knobel E. Stressors in ICU: patients' evaluation. *Intensive Care Med*. 1997; 23(12):1282–1285. [PubMed: 9470087]
- Nygårds ME, Sörnmo L. Delineation of the QRS Complex using the Envelope of the ECG. *Med. and Biol. Engin. and Comput*. 1983; 21(5):538–547.
- Oster J, Behar J, Colloca R, Li Q, Li Q, Clifford GD. Open source Java-based ECG Analysis Software and Android App for Atrial Fibrillation Screening. *Comput in Cardiol*. 2013:731–734.
- Pan J, Tompkins WJ. A real-time QRS detection algorithm. *IEEE Transactions on Biomedical Engineering*. 1985; 32(3):230–236. [PubMed: 3997178]
- Paoletti M, Marchesi C. Discovering dangerous patterns in long-term ambulatory ECG recordings using a fast QRS detection algorithm and explorative data analysis. *Computer Methods and programs in biomedicine*. 2006; 82(1):20–30. [PubMed: 16529841]
- Parthasarathy S, Tobin MJ. Sleep in the intensive care unit. *Intensive Care Med*. 2004; 30(2):197–206. [PubMed: 14564378]
- Plesinger, F.; Klimes, P.; Halamek, J.; Jurak, P. Computing in Cardiology Conference (CinC), 2015. IEEE; 2015. False alarms in intensive care unit monitors: Detection of life-threatening arrhythmias using elementary algebra, descriptive statistics and fuzzy logic; p. 281-284.
- Plesinger F, Klimes P, Halamek J, Jurak P. Taming of the Monitors: Reducing False Alarms in Intensive Care Units. Accepted in *Phys. Meas*. 2016
- Rodrigues R, Couto P. Detection of False Arrhythmia Alarms with emphasis on Ventricular Tachycardia. Accepted in *Phys. Meas*. 2016
- Rooijackers MJ, Rabotti C, Oei SG, Mischi M. Low-complexity R-peak detection for ambulatory fetal monitoring. *Physiological Measurement*. 2012; 33(7):1135. [PubMed: 22735075]
- Sadr N, Huvanandana J, Nguyen DT, Kalra C, McEwan A, de Chazal P. Reducing false arrhythmia alarms in the ICU using multimodal signals and robust QRS detection. Accepted in *Phys. Meas*. 2016
- Saeed M, Villarroel M, Reisner AT, Clifford GD, Lehman LW, Moody GB, Heldt T, Kyaw TH, Moody B, Mark RG. Multiparameter Intelligent Monitoring in Intensive Care II (MIMIC-II): a Public-Access Intensive Care Unit Database. *Critical Care Medicine*. 2011; 39(5):952. [PubMed: 21283005]
- Silva, Moody GB. An Open-source Toolbox for Analysing and Processing PhysioNet Databases in MATLAB and Octave. *Journal of Open Research Software*. 2014 Sep 24.2(1):e27. [PubMed: 26525081]

- Slevin M, Farrington N, Duffy G, Daly L, Murphy JF. Altering the NICU and measuring infants' responses. *Acta Paediatr.* 2000; 89(5):577–581. [PubMed: 10852196]
- Sörnmo, L.; Laguna, P. *Bioelectrical signal processing in cardiac and neurological applications.* Academic Press; 2005.
- Sun J, Reisner A, Mark R. A Signal Abnormality Index for Arterial Blood Pressure Waveforms. *Comput in Cardiol.* 2006:13–16.
- Topf M, Thompson S. Interactive relationships between hospital patients' noise induced stress and other stress with sleep. *Heart Lung.* 2001; 30(4):237–243. [PubMed: 11449209]
- Tsien CL, Fackler JC. Poor Prognosis for Existing Monitors in the Intensive Care Unit. *Critical Care Medicine.* 1997; 25(4):614–619. [PubMed: 9142025]
- Tsimenidis C, Murray A. False alarms during patient monitoring in clinical intensive care units are highly related to poor quality of the monitored electrocardiogram signals. Accepted in *Phys. Meas.* 2016
- Yazdani, S.; Vesin, JM. Computing in Cardiology Conference (CinC), 2014. IEEE; 2014. Adaptive mathematical morphology for QRS fiducial points detection in the ECG; p. 725-728.
- Zhu T, Dunkley N, Behar J, Clifton DA, Clifford GD. Fusing continuous-valued medical labels using a Bayesian model. *Annals of Biomedical Engineering.* 2015; 43(12):2892–902. [PubMed: 26036335]
- Zhu T, Johnson AEW, Behar J, Clifford GD. Crowd-sourced annotation of ECG signals using contextual information. *Annals of Biomedical Engineering.* 2014; 42(4):871–84. [PubMed: 24368593]
- Zong W, Heldt T, Moody G, Mark R. An open-source algorithm to detect onset of arterial blood pressure pulses. *Comput in Cardiol.* 2003:259–262.
- Zong W, Moody GB, Jiang D. A robust open-source algorithm to detect onset and duration of QRS complexes. *Comput in Cardiol.* 2003; 30:737–740.
- Zong W, Moody GB, Mark RG. Reduction of false arterial blood pressure alarms using signal quality assessment and relationships between the electrocardiogram and arterial blood pressure. *Med Biol Eng Comput.* 2004; 42(5):698–706. [PubMed: 15503972]
- Zong W, Nielsen L, Gross B, Brea J, Frassica J. A Practical Algorithm to Reduce False Critical ECG Alarms Using Arterial Blood Pressure and/or Photoplethysmogram Waveforms. Accepted in *Phys. Meas.* 2016

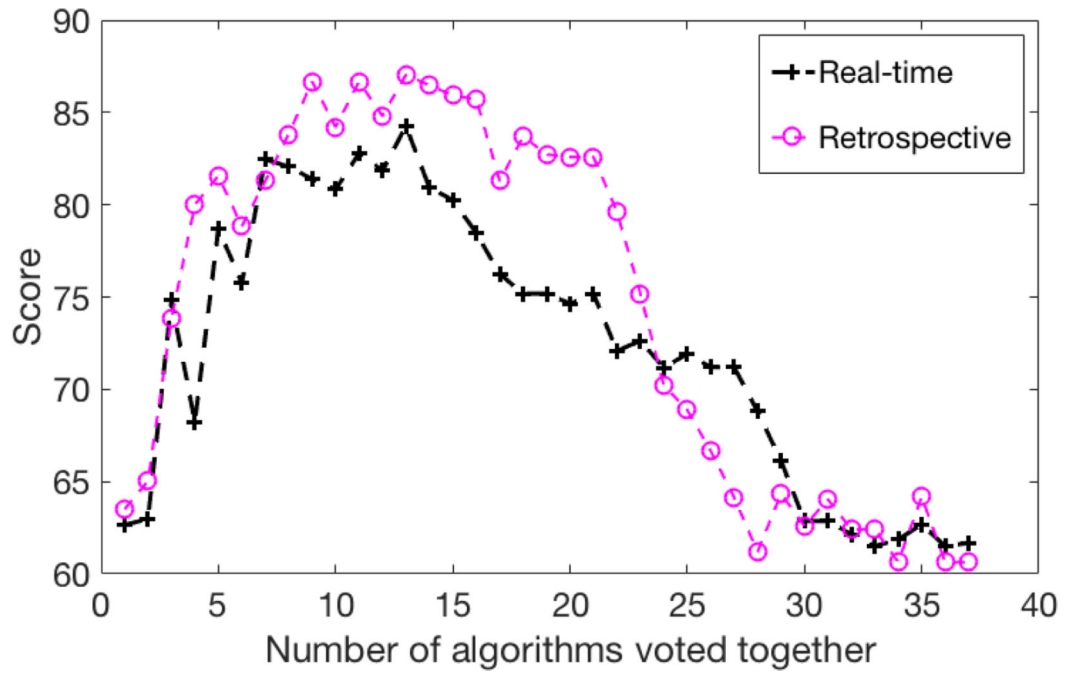


Figure 1.

Performance of voting algorithms as a function of number of algorithms for both the real time and retrospective events. Algorithms were chosen by ranking them in descending order of score on the training data, and the test data score was reported (to prevent over-estimation of the score). Equal weights were given to all algorithms and a tied, absent or no vote was treated as 'true'.

Table 1

Types of alarms and signals used in the Challenge. Each of the N records included two ECG channels.

| | Training (N=750) | | Test (N=500) | |
|-------|------------------|------|--------------|------|
| | False | True | False | True |
| ASY | 100 | 20 | 90 | 12 |
| EBR | 45 | 45 | 38 | 26 |
| ETC | 8 | 131 | 5 | 68 |
| VTA | 253 | 90 | 176 | 45 |
| VFB | 52 | 6 | 34 | 6 |
| PPG | 227 | 178 | 158 | 83 |
| ABP | 59 | 63 | 58 | 39 |
| Both | 172 | 51 | 127 | 35 |
| Total | 458 | 292 | 343 | 157 |

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 2

Final scores for the top 9 entrants in both events (real-time and retrospective) ranked by overall real-time score, the three example algorithms provided and a voting approach.

| Entrant | Event 1 (Real-time) | | | Event 2 (Retrospective) | | |
|----------------------------------|---------------------|------------|--------------|-------------------------|------------|--------------|
| | TPR | TNR | Score | TPR | TNR | Score |
| Plesinger et al. (2015) | 92% | 88% | 81.39 | 95% | 88% | 84.96 |
| Kalidas and Tamil (2015) | 94% | 82% | 79.44 | 94% | 86% | 81.85 |
| Krasteva et al. (2015b) * | 93% | 83% | 79.41 * | 93% | 84% | 79.56 * |
| Couto et al. (2015) | 89% | 91% | 79.02 | 88% | 92% | 78.28 |
| Fallet et al. (2015) | 94% | 77% | 76.11 | 99% | 80% | 85.04 |
| Hoog Antink and Leonhardt (2015) | 93% | 77% | 75.55 | 90% | 82% | 75.18 |
| Eerikäinen et al. (2015) | 90% | 82% | 75.54 | 89% | 85% | 75.52 |
| Ansari et al. (2015) | 89% | 84% | 74.48 | 89% | 87% | 76.57 |
| Liu et al. (2015) | 89% | 79% | 71.68 | 93% | 78% | 75.91 |
| Example Algorithm 1 | 76% | 44% | 41.41 | 73% | 46% | 40.83 |
| Example Algorithm 2 | 86% | 38% | 45.07 | 84% | 38% | 44.37 |
| Example Algorithm 3 | 64% | 76% | 45.59 | 61% | 77% | 47.35 |
| Voting Algorithm (N=11) | 94% | 87% | 82.78 | 94% | 93% | 86.67 |
| Voting Algorithm (N=13) | <u>94%</u> | 90% | <u>84.26</u> | 94% | <u>94%</u> | <u>87.04</u> |

Best performances of competition entrants are in bold. TPR = fraction of true alarms correctly classified; TNR = fraction of false alarms correctly classified.

* denotes an unofficial (“closed-source”) entry. Underlined scores are the highest unofficial scores in the table.

Table 3

Performances of the competitors in both events (real-time and retrospective) ranked by overall real-time score during the follow-up phase (Spring 2016).

| Authors | Real-time | | | Retrospective | | |
|------------------------------|------------|------------|-----------------------------|---------------|------------|-----------------------------|
| | TPR | TNR | Score | TPR | TNR | Score |
| Plesinger et al. (2016) | 93% | 87% | 81.62 (81.39) | 95% | 88% | 84.96 |
| Krasteva et al. (2016) | 92% | 87% | 80.07 (79.41 [*]) | 93% | 88% | 81.75 (79.56 [*]) |
| Kalidas and Tamil (2016) | 94% | 82% | 79.44 | 94% | 86% | 80.29 (81.85) |
| Hoog Antink et al. (2016) | 95% | 78% | 78.20 (75.55) | 93% | 76% | 74.45 (75.18) |
| Eerikäinen et al. (2016) | 93% | 80% | 77.39 (75.54) | 95% | 83% | 81.58 (75.52) |
| Fallet et al. (2016) | 95% | 76% | 77.07 (76.11) | 99% | 80% | 85.04 |
| Ansari et al. (2016) | 89% | 85% | 76.23 (74.48) | 88% | 84% | 73.40 (76.57) |
| Rodrigues and Couto (2016) | 92% | 78% | 74.28 (79.02) | 92% | 78% | 74.46 (78.28) |
| Liu et al. (2016) | 89% | 79% | 71.68 | 93% | 78% | 75.91 |
| Sadr et al. (2016) | 95% | 65% | 69.92 | 98% | 66% | 74.03 |
| Tsimenidis and Murray (2016) | 92% | 66% | 67.88 | 92% | 69% | 68.71 |
| Daluwatte et al. (2016) | | | - | | | - |
| Zong et al. (2016) | | | - | | | - |

Best performances are in bold. TPR = fraction of true alarms correctly classified; TNR = fraction of false alarms correctly classified. Numbers in the parentheses are the score of competition entrants if they are different with those in the follow-up phase.

^{*} denotes an unofficial (“closed-source”) entry.