**Literature Review for Online Event Detection from Water Quality Time Series**
*Sean A. McKenna, Sandia National Laboratories, June 2008*

This document provides background information on event detection from water quality time series data. This literature review is divided into two sections: General Approaches and Water Quality Event Detection. The General Approaches section is a literature review of techniques developed in the fields of statistics, signal processing, data mining and others that form the fundamental basis of various water quality event detection systems. The focus of this section is on various approaches to state estimation and residual classification. The Water Quality Event Detection section summarizes a number of relatively recent publications in the areas of online detection of water quality events.

Additionally, there is a glossary at the end of this document that provides concise definitions to some of the terms used in this report. The definitions are given with some context to the meaning of these terms in water quality event detection.

## *General Approaches*

The general problem of detecting anomalous behavior in time series data is a subject of research in a number of disparate fields including data mining, network intrusion detection, tsunami detection, traffic accidents analysis, mechanical component failure and system fault detection among others. This work can be broken into two distinct approaches that we term here as "online" and "off-line". The online approaches receive data in discrete time steps ranging from milliseconds to minutes and provide a determination of the presence or absence of an anomalous reading immediately after the receipt of each new data point. These approaches are the focus of this work and are considered in greater detail below.

## Offline Change Point Detection

In contrast to the on-line approaches, a number of off-line or "after the fact" approaches to processing of time series data have also been developed. These approaches comprise a significant amount of literature in the area of data mining of time series and are reviewed here only briefly. The majority of the offline approaches to identifying anomalous behavior are based on the detection of *change points*. Change points are defined as abrupt changes in the nature of a signal as generally defined by statistical measures of that signal. As an example, the time at which there is a change in the source of water supplying a monitoring location can be a change point for the water quality at that location.

There are a number of approaches to the change-point detection problem, but the essential elements are that there is at least one point in time where the properties of the process that created the observed data change. This change point is generally determined by fitting a statistical model to the data and identifying the location in time where parameterization of that model changes significantly. These approaches are typically done in an offline manner as part of a data mining approach and can often be classified as "retrospective segmentation" approaches (Adams and MacKay, 2007) where the change points define the ends of the various segments (e.g., Tsihrintzis, and Nikias, 1995) of the signal behavior. Additionally, two separate models can be used to fit the data before and after the change point. Application of regression models to continuous data and Poisson models to discrete count data are common (See Raftery, 1994 for a review). Another natural approach to change point detection is conceptualization of the process generating the signal as a Markov process. For these cases, Markov models where change points denote the switch between model states can be applied (e.g., Ge and Smyth, 2000a). Some approaches require that the number of change points in the time series is known, or defined by the user, prior to the analysis while others techniques are less restrictive and will determine the necessary number of change points to fit the data to some specified tolerance.

In particular, identifying the point in time at which a change in the count or the continuous data occurs has been an active area of research with motivating applications covering a diverse set of problems from the annual rate of coal mining accidents (e.g., West and Ogden, 1997) to highway traffic patterns (Ihler et al., 2006) to semi-conductor manufacturing process control (Ge and Smyth, 2000b) to cite a few examples. Some

recent work has been done in the area of merging change-point detection approaches with those of on-line event detection (e.g., Takeuchi and Yamanishi, 2006).

The general approach to offline change-point detection is to examine data from opposite sides of a proposed change point to determine if those two data sets are significantly different from each other.  If they are, the point that separates the two data sets is a change point.  For off line analyses, the full data set has already been recorded and is available for analysis.  In the online world, only the data recorded up to the present time are available and the goal is to identify the change point as close to the time at which it occurs as possible.  This constraint of making a determination as near to real time as possible, limits the available number of measurements that occur after the change point to as small of a number as possible.  The goal of an efficient water quality event detection system (EDS) is to develop an online approach to water quality event detection that can warn analysts and system operators in real-time of unexpected water quality conditions.  To meet this goal, the offline approaches discussed above are not viable.

A number of online, or real-time, approaches to identifying anomalous observations in time series data have been developed for use in a large number of different fields.  Some of the fundamental tools and the basic approaches to online event detection are covered below this section with example citations provided.  The literature covering online event detection is vast and this review only covers a fraction of the techniques that have been published with a focus on covering the fundamental techniques that have been incorporated into water quality EDS tools

## Control Charts

Perhaps the oldest approaches to online event detection are the Shewhart Charts and Cumulative Sum (CUSUM) charts developed in 1920's.  These approaches were originally developed for quality control in manufacturing and industrial processes and are now used in a number of other applications as well.

The CUSUM chart shows the cumulative sum of differences between the measured values and the average value.  These differences are calculated by subtracting the average from each value.  Increases in the cumulative sum value indicate a time of values that are continuously above the average.  The resulting CUSUM chart will have an upward slope during such a period of relatively high values.  The opposite results hold for periods of relatively lower values.

Shewhart charts calculate a chosen statistic of the observed data (e.g., mean) using a moving window through time.  Control limits are calculated for the value of the statistic and any values of the statistic that deviate beyond the confidence bounds are identified as outliers.  The control limits may be calculated using the expected variation in the range of data values or using more common statistical tools based on an assumption of Gaussian variation in the calculated statistic.

The choice between applying a CUSUM vs. a Shewhart chart depends on the nature of the process being monitored.  In general, CUSUM charts are felt to be better at detecting

small, yet sustained, changes in the mean value of a process (NIST, 2008), while Shewhart charts are often better suited to incorporating knowledge of the operating conditions held by the analyst.

Applications of both CUSUM and Shewhart charts typically employ standard statistical approaches to determine the range of control or the confidence limits for the process. A limitation of these tools it that they generally rely on assumptions of stationary independent and identically distributed (i.i.d.) variables and typically invoke the Gaussian distribution to define these variables. Adaptations to these charts have been made to accommodate time series data with autocorrelation, non-Gaussian distributions, non-stationarity in the data and various ways of calculating the control limits (see Lai, 1995; Zhang, 1997)

Water quality time series are inherently non-stationary. Both daily (diurnal) cycles and seasonal patterns are the cause of these non-stationarities. Additionally, short term (daily to weekly) and longer term (multiple week) trends in water quality data are caused by varying levels of control of water treatment and hydraulic operations within the utility along with drift in the water quality sensors. It may be possible to employ techniques that have been developed for stationary time series, such as Shewhart and CUSUM charts, but first it would be necessary to remove, or "detrend" the non-stationary aspects of the observed data. Therefore, a robust means of modeling the background variation in water quality is a necessary step in being able to separate that background from the water quality events. This modeling of the background variation is generally referred to as "state estimation" and is the first step in modeling of non-stationary time series.

A common model for the on-line detection of changes in time series data incorporates two components that work in concert: 1) a state estimation model and 2) a residual classification algorithm. The state estimation model uses previous observations of one or more time series measurements to estimate future values of the process. These estimates could also be made by a physical process model (e.g., chemical reactions and solute transport within the pipe network). The residual classification algorithm then uses the differences between the predicted state and the observed state to determine whether or not the observed state represents an anomalous condition. This basic approach has been employed for detection of anomalous conditions in time series data in various fields including tsunami detection (Gower and Gonzalez, 2006), component degradation in nuclear power plants (Yuan, et al., 2005) and "aging" in computer software (Vaidyanathan and Gross, 2003).

## State Estimation

Modeling of background water quality falls into the general time series modeling category of state estimation. There is an underlying state of the water quality that is only sensed through noisy and sometimes incomplete data from water quality sensors. The goal is to provide an accurate estimate of the unknown state and do this iteratively so that the state estimate is updated at every time step. The state estimate is most often quantified by parameter estimates in a statistical model of the time series. For each parameter in the model of the state, the best estimate of that parameter is provided and,

depending on the complexity of the state estimation approach, a measure of uncertainty about the estimated parameter may also be determined.

There are a number of modern statistical and mathematical advances that facilitate time series forecasting and have been applied to state estimation including neural networks (e.g., Boznar et al., 1993; Lu et al., 2002), support vector machines (e.g., Muller et al., 1999), and wavelets (e.g., Lueck, et al., 2000). However, in this literature review on traditional techniques derived in the fields of signal processing and time series analysis. These approaches have proven useful in the development and application of water quality event detection tools.

Traditional approaches to time-series analysis provide "data driven" models based on the theory of time-series analysis as defined by Box and Jenkins (1976). These approaches include the popular autoregressive (AR) and moving average (MA) models as well as the various hybrids of these approaches (ARMA) and autoregressive integrated moving average (ARIMA). These models use observed data to estimate the parameters of the models and then use these estimated parameters to predict the expected data values at future observation times. These models are designed to provide a measure of uncertainty on the resulting predicted value of the time series. In essence, the time series models can be thought of as a filtering process where the noise in the underlying physical processes and in the measurements is filtered out to leave the best estimate of the water quality. Linear filters as used in signal processing can be built from AR and MA models. A thorough treatment of these approaches is given in Bras and Rodriguez-Iturbe (1993). These models are used heavily in signal processing, surface water hydrology, and econometrics applications and have also been adapted to estimate spatially correlated properties in 2 and 3 dimensions (see Goovaerts, 1997; Journel and Huijbregts, 1979). Traditional application of these models considers the estimated parameters as point estimates with no uncertainty, although Bayesian approaches to time series modeling can incorporate parameter uncertainty into these models.

The many variations of Kalman filters represent the next level of complexity in state estimation. Kalman filters are currently popular for data assimilation where observed data can be used to iteratively update parameters of physical process models as well as estimate future observations of the process. Kalman filters incorporate both uncertainty in the underlying model and/or its parameters along with uncertainty in the observed data in predictions of future values of the time series. Original development of the Kalman filter (Kalman, 1960) was focused on state estimation for linear systems with assumed Gaussian errors, model and observation, and covariance structures. The extended and ensemble Kalman Filters (EKF and EnKF, respectively) were motivated by both the need to solve more highly non-linear problems and the inadequacy of the KF for solving these problems (Evensen, 1993; Evensen, 1994). In particular, the EnKF replaces the analytical calculation of covariances for both the model error and observational error and the assumptions necessary for those calculations with a numerical approximation where the covariance terms are calculated across a stochastic ensemble of model states and resulting model predictions (see Evensen, 2003; Moradkhani et al., 2005a).

A known disadvantage of the EnKF approach is that it has been developed for models with non-linear relationships between inputs and outputs, but it relies on a linear updating process and the probabilistic approach to uncertainty estimation, for all variants of the Kalman Filter is only valid up to second order (i.e., the output of any KF is a mean estimate and a variance defining uncertainty about that estimate, but no shape to the uncertainty distribution is provided). The second-order basis of the uncertainty estimation for the predictions essentially limits the validity of the KF approach to distributions that are at least symmetric, if not moderately Gaussian. For state estimation problems where uncertainty estimates that take into account higher-order moments of the predictive distribution are needed, particle filtering techniques (see Arulampalam et al., 2002; Gordon et al., 1993; Moradkhani et al., 2005b) provide the next level of uncertainty quantification along with additional complexity in applications.

The time series models and the Kalman filter approaches to state estimation employ some optimal weighting of previous measurements to predict the future state of the water quality. Another decidedly simpler approach to state estimation is to use just a single previous water quality measurement as the state estimate. Two approaches to using a single previous measurement as the state estimate known as time series increments and multivariate nearest neighbor are discussed further below. The time series increments approach uses the single most recent observation as the state estimate. This approach is equivalent the Markov model, often referred to as the Thomas-Fiering model in surface water hydrology where it has been applied to modeling stream flows (Bras and Rodriguez-Iturbe, 1993). The multivariate nearest neighbor approach (Klise and McKenna, 2006) uses the measurement within a window of recent measurements that is closest to the current observation as measured within the multivariate space defined by the observed water quality signals.

The field of signal processing provides another means of state estimation that uses cross-correlations between different signals as well as the autocorrelations within each signal to estimate the future water quality values. For example, the best estimate of the next pH value might be a weighted combination of the 10 previous pH values, the Cl value from 24 hours ago and the temperature value from 12 hours ago. This approach to state estimation is a common tools in the signal processing field and has been incorporated into event detection schemes (e.g., Zavaljevskl and Gross, 2000).

State estimation approaches that exploit cross-correlations between signals are well-suited to situations where sensor and data transmission reliability are not an issue (e.g., engine monitoring), but when a sensor fails, the entire state estimation model fails. In environmental monitoring situations, sensor and/or data transmission failure can be common and these cross-correlation based approaches may not be best suited to these situations.

## Residual Classification

Residual classification is the process of classifying each deviation between the observed and predicted water quality values (state) as still being part of the background or being a significant deviations from the background. Small residuals, or deviations, can be

considered as arising from incomplete parameterization of the state estimation model and measurement error. Large deviations are deemed to be significant departures from the expected background water quality and therefore are indicative of a critical change in the system (an outlier). The simplest approach to residual classification is to apply a single threshold value to the residuals and those that exceed the threshold are considered outliers. Two issues complicate this simple approach: 1) A single constant threshold value may not apply equally well to all times in the data set, so an adaptive thresholding approach may make more sense; and 2) The thresholding approach needs to take into account the fact that state estimation and residual calculation may be done separately for each water quality signal and therefore a multivariate approach to residual classification is necessary.

The threshold used in residual classification can often be made more efficient by adapting the size of the threshold to the size, or variability, of the residuals. One approach, is to make the threshold a multiplier of the standard deviation of the signal such that the threshold adapts to the variability of the signal. Normalization of the signal values to a fixed variance within a moving window allows for a threshold that is a constant multiplier of the variance, but scales relative to the un-normalized signal variance. This approach is used in the CANARY software (see McKenna et al, 2007; Hart et al., 2007). Breitgand et al., (2005) demonstrate a logistic regression based algorithm for setting adaptive thresholds in the context of computer performance monitoring.

State estimation techniques that use cross-correlations between signals generally result in a single estimate of the state that is integrated over all input signals. For these approaches a single residual is calculated at each time step. Independent state estimation for each signal results in a residual for each signal and these must be combined, or "fused" in some way to identify an outlier at that time step. A simple approach is to make a single classification for each time step using some combination of the residual values from all sensors operating at that time step. Equivalent results are obtained by using the average or the sum of the residuals. Classification using the maximum residual across all sensors also makes it easy to store the sensor that is responsible for the outlier at each time step. More complicated approaches to decision fusion are examined by Dasarathy (1991).

Independent state estimation followed by residual fusion allows for the number of sensors providing information at each time step to change over time. The structure of the event detection approach does not have to change to accommodate the loss or addition of sensors. This flexibility is in contrast to state estimation tools that employ cross-correlation between signals where a change in the number of sensors requires reconstruction of the model.

Takeuchi and Yamanishi (2006) integrate their deviation scores (essentially residuals) over time by calculating a moving average value. A threshold is then applied to the moving average value to detect outliers and change points. The sensitivity of this algorithm to short-lived events is controlled by the length of the moving average applied to the residuals. The Multivariate State Estimation Technique (MSET) uses a sequential

probability ratio test to examine how well the distribution of residuals fits a predefined Gaussian distribution (Zavaljevskl and Gross, 2000). Several different hypothesis tests are examined in the MSET approach to look for mean residuals that are above or below the expected value (mean = zero) as well as variances that deviate from the expected variance value. Applications of the MSET approach to problems of computer reliability have shown it to be especially adept at detecting early stages of component degradation (Vaidyanathan and Gross, 2003). McKenna et al., (2007) demonstrated the binomial event discriminator for mapping outliers to events in a water quality event detection application and this approach is discussed further below.

## *Water Quality Event Detection*

Continuous, reliable delivery of safe drinking water to customers is an essential component of the viability of large metropolitan areas. The distribution networks used to deliver water represent a critical component of municipal infrastructure systems. The areal extent of these distribution systems and their overall design to improve customer service and firefighting make these networks susceptible to accidental or intentional contamination events (see discussion by Kroll and King , 2006 for a description of credible contamination threats to distribution networks). Such events that could degrade water quality within water distribution systems have focused recent discussion on various means of hardening both the physical and cyber components of these systems against contamination events. The concept of a contaminant warning system (CWS) has been proposed as an integrated tool that employs in-situ sensors, supervisory control and data acquisition (SCADA) systems, and water quality event detection systems (EDS) to continuously monitor network conditions and warn operations personnel of any potential contamination events (see Hasan, et al., 2004; Grayman, et al., 2001; Roberson and Morley, 2005). The focus of this work is on the EDS component of the CWS concept.

Other key technical elements of the CWS are the type and characteristics of the actual sensors, the data transmission and storage performance of the SCADA system that connects the sensors and the location of the sensors within the network. This latter element has been the subject of considerable recent research (e.g., Berry et al., 2005; Grayman et al., 2006). However, many of these studies have considered the event detection portion of the CWS to be comprised of "perfect" sensors and/or event detection software and only a few papers (e.g., Berry et al., 2006; McKenna et al., 2006) have examined the relationship between sensor performance, sensor placement within the network and the overall performance of the CWS.

### Online Water Quality Event Detection

Development of online event detection tools for water quality data has been an area of recent interest. A number of published approaches to this problem are reviewed below. These papers demonstrate the response of surrogate parameters to various contaminants and the approaches developed to detect events. The majority of these approaches work with sensor data from a single monitoring location; however, several of them have been designed to integrate sensor information from more than one location.

Byer and Carlson (2005) examined the response of surrogate monitors to the introduction of various contaminants in both laboratory beakers and in bench scale tests using water from a local utility. Their results clearly indicate the response of several surrogate parameters to the introductions of a range of contaminants at various concentrations. These results are also used by Cook et al (2005) in testing an event detection system. More recently, Hall et al (2007) tested the response of a number of commercially available water quality sensors in the presence of nine different contaminants introduced to a pipe test loop at different concentrations and found that at least one of the surrogate parameters responded to the presence of every contaminant.

Byer and Carlson (2005) conducted event detection through a relatively simple approach of comparing the measurement at any time to a predefined mean baseline level and defining anomalous values as those that exceed +/- 3 standard deviations from the mean of the baseline values. The baseline values were considered to be stationary and calculated by either using all of the available data, approximately 16,000 observations, or using the 100 observations immediately prior to arrival of the contaminant at the sensor.

Cook et al (2006) outline the development of a case-based reasoning system (CBRS) for the identification of multivariate data patterns that represent acceptable changes in water quality. The CBRS acts as a classifier to identify the current state of the system and patterns that cannot be classified into existing groups are considered outliers. The work by Cook et al (2006) highlights the need for accurate and reliable sensing of the water quality data; sophisticated software cannot make up for low quality input data.

Jarrett et al (2006) focused their analysis of water quality data on the control exerted by the time of day and the day of the week on the expected water quality value. In the systems they examined, operation of the distribution network was responsible for a significant portion of the water quality variation, and those operations followed a reasonably predictable behavior. However, the temporal patterns controlling the water quality tended to change over time and therefore it was difficult to accurately predict water quality based solely on the time of day. Jarrett et al. (2006) proposed that a control chart approach applied to the first differences (increments) of the water quality data may prove useful in event detection if the center line and widths of the control region were both allowed to vary temporally (a non-stationary control chart approach).

Kroll and King (2006) provide a rough outline of a proprietary EDS that includes both a baseline (state) estimation component and a multivariate classification component. The classification step uses the deviations in the measured signals from the baseline along with a library of previously recorded deviations to classify the cause of the event as being either a particular contaminant or a change in water quality caused by a change in operations at the utility. Patterns that do not match any of the library patterns are declared "unknown" and can be added to the library by the operator. The ability of the algorithm to "learn" through time is used to lower the number of false positives upon deployment.

Klise and McKenna (2006a) examined the utility of multivariate classification schemes for event detection. The state estimation approach in this work was to define every time step of the baseline water quality as belonging to one of a finite number of clusters within the multivariate space defined by the vector of surrogate parameter measurements, or by a lower dimensional representation of that space as defined by principal components. Results showed that increasing the number of clusters that define the baseline to the point where every recent time step was considered to be a separate cluster improved results over smaller numbers of clusters. The result of this work was the development of the multivariate nearest neighbor (MV-NN) algorithm in which the distance from any new measured water quality vector to the nearest previously measured vector in the multidimensional space is recorded. If that distance exceeds a specified threshold, the new data point is considered to be an outlier. Klise and McKenna (2006b) further tested the MV-NN algorithm using event data from the US EPA Test and Evaluation Facility (Hall et al., 2007) that were superimposed on water quality data collected at a US utility and found that event detection results with MV-NN are sensitive to the contaminant type and the background water quality variability at each monitoring location.

McKenna et al (2006b) compared three approaches to state estimation for each water quality time series: time series increments (the previous measured value is the predictor of the next value), linear filters and the MV-NN approach. For the increment and linear filter approaches, the residuals between the predicted and measured water quality values were fused across all water quality signals and this final fused residual was compared to a threshold to define whether or not it was an event. The multivariate distances in the MV-NN algorithm were compared directly to the same threshold as they already represent a measure of prediction accuracy that takes into account all water quality signals. Testing was completed on simulated time series and actual measured water quality time series data. Simulated events were added to all data sets. Results showed that the MV-NN algorithm was best able to predict the water quality background in all cases, but that the ability to predict the background did not necessarily translate into the best detection capabilities as measured by the false positive and false negative rates.

Prior work in event detection algorithms by Klise and McKenna (2006a,b) and McKenna et al. (2006b) had evaluated every time step against a threshold. Those time steps with residuals from the baseline that exceeded the threshold were classified as events (a single outlier equals an event). This approach led to a large number of false positives as water quality within utility systems can be quite noisy and additional noise is added to the water quality data as it is transmitted through the SCADA system. McKenna et al (2007) introduced the binomial event discriminator (BED) as a means of aggregating results over multiple time steps to determine whether or not an event was occurring. Each individual time step is now considered to be part of the background or an outlier and the number of outliers (failures) within a given number of time steps (trials) as inputs to the binomial distribution defines the probability that a water quality event is occurring. Addition of the BED to the water quality prediction algorithms allowed for order of magnitude reductions in false alarms for the data sets tested.

The papers discussed above rely on various statistical models to estimate the state by tracking the baseline water quality conditions such that a comparison between the expected baseline value and the observed values can be made. Another approach to determination of the baseline conditions would be to employ a model that directly simulates all water quality parameter values through the physical and chemical processes occurring in the distribution network between the treatment plant and the monitoring station and is continuously updated in real time (Shang et al., 2008).

## Monitoring at Multiple Locations

The work presented here is focused on analysis of data from each monitoring station independently. However, several publications have shown additional benefit that can be gained from combining water quality data from more than one location. O'Halloran et al. (2006) developed a water parcel tracking approach that matched the "fingerprint" of water quality recorded at two different locations along the same flowpath in the network. They were able to use this automated technique to determine the transit time of the water between the two monitoring locations, although an assumption of steady flow was required. Yang et al. (2007) defined a technique for improving event detection and reducing false positives by combining water quality monitoring data from two monitoring stations along the same pipe. Data from two stations in series allows for transport modeling applied to the water quality between the two stations where data from the first station essentially provides the initial conditions for the transport solution. This transport modeling provides improved state estimation at the downstream monitoring station. A more general approach to integrate data from two or more monitoring stations that may not be in direct hydraulic connection has recently been tested (Koch and McKenna, 2008) with promising results.

## Evaluating Algorithms

An issue of considerable importance in water quality monitoring is the appropriate evaluation of an event detection algorithm. As a rule, utilities do an extremely good job of supplying high quality water to their customers without fail, and water quality events, even those due to routine causes such as main breaks, faults in a primary treatment system or failures of a chlorine booster station are rare. Documented accounts of malevolent contamination of a utility are even rarer and this routine delivery of high quality water makes it nearly impossible to completely test event detection systems in real-world situations. Testing of EDS with experimental data obtained in laboratory settings (e.g., Byer and Carlson, 2005; Kroll and King, 2006) or in specially designed pipe loops (e.g., Hall et al., 2007; Yang et al., 2007) provides direct measurement of sensor responses to controlled contamination events, but typically the variation in background water quality for these tests is considerably less than that experienced within operating distribution networks. Therefore, the most direct means of evaluating event detection systems is to simulate the response of water quality monitoring sensors to the introduction. Simulation based evaluation of EDS tools has been done by a number of authors using varying levels of sophistication in the contaminant simulation approach (e.g., McKenna et al. 2006; McKenna et al., 2008; Uber et al., 2007; Shang et al., 2008; Umberg et al., 2008; Allgeier et al., 2008)

Reports of early work in event detection from water quality data focused on the number of known events that were detected. However, as pointed out by Rizak and Hrudey (2006) when monitoring for events that are expected to occur with an very low probability, dealing with false positive events will consume the largest amount of the monitoring organization's resources. McKenna et al (2008) employed the receiver operating characteristic (ROC) curve approach from the signal processing and medical diagnostics fields to apply to evaluation of water quality event detection systems. ROC curves demonstrate the tradeoff between the rate of false events and the probability of detection true events on a single graph. Typical ROC curve shapes quantify the increase in false positive events as the sensitivity of the algorithm is increased to improve the probability of detecting true events.

# Bibliography

Adams, R.P. and D.J.C. MacKay, 2007, Bayesian Online Changepoint Detection, arXiv:0710.3742v1

Allgeier et al., 2008, Systematic Evaluation of Contaminant Detection through Water Quality Monitoring, presentation at AWWA Security Congress, Cincinnati, Ohio

Arulampalam, M.S., S. Maskell, N. Gordon and T. Clapp, 2002, A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking, *IEEE Transactions in Signal Processing,* 50 (2), pp. 174-188.

Berry, J.W., L Fleischer, W.E. Hart, C.A. Phillips and J. Watson, 2005, Sensor Placement in Municipal Water Networks, *Journal of Water Resources Planning and Management,* 131 (3), pp. 237-243.

Berry, J., R.D. Carr, W.E. Hart, V.J. Leung, C.A. Phillips and J. Watson, 2006, On the placement of imperfect sensors in municipal water networks, in Proceedings of the 8[th] Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006:

Boznar M., Lesjak M., Mlakar P., 1993. A neural netwok-based method for short-term predictions of ambient SO 2 concentrations in highly polluted industrial areas of complex terrain. *Atmospheric Environment* 27(2), 221-230.

Box, G.E. and G.M. Jenkins, 1976, *Time Series Analysis: Forecasting and Control*, Holden Day, San Francisco.

Bras, R.L. and I. Rodriguez-Iturbe, 1993, *Random Functions and Hydrology*, Dover Publications, Inc., Mineola, New York, 559 pp.

Breitgand, D. E. Henis, and O. Shehory, 2005, Automated and Adaptive Threshold Setting: Enabling Technology for Autonomy, and Self-Management Proceedings of the Second International Conference on Autonomic Computing (ICAC'05), IEEE, 12 pp.

Byer, D. and K.H. Carlson, 2005, Real-Time Detection of Intentional Chemical Contamination in the Distribution System, *American Water Works Association Journal*, 97 (7), pp. 130-141.

Cook, J. E. Roehl, R. Daamen, K. 2005, Decision Support System for Water Distribution System Monitoring for Homeland Security, Proceedings of the AWWA Water Security Conference, Oklahoma City, OK, April 10-12.

Cook, J.B., J.F. Byrne, R.C. Daamen and E.A. Roehl, 2006, Distribution System Monitoring Research at Charleston Water System, in: Proceedings of the 8th Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006

Dasarathy, B.V., 1991, Decision Fusion Strategies in Multisensor Environments, *IEEE Transactions on Systems, Man and Cybernetics*, 21 (5), pp. 1140-1154.

Evensen,G., 1992, Using the extended Kalman filter with a multi-layer quasi-geostrophic ocean model, *Journal of Geophysical Research*, 97 (C11), pp. 17,905-17,924.

Evensen, 1994, Sequential data assimilation with a non-linear quasi-geostrophic model using Monte Carlo methods to forecast error statistics, *Journal of Geophysical Research*, (99) C5, pp. 10,143-10,162.

Evensen, G., 2003, The ensemble Kalman filter: Theoretical formulation and practical implementation, *Ocean Dynamics*, 53 (4), pp. 343-367.

Ge, X. and P. Smyth, 2000a, Deformable Markov Model Templates for Time-Series Pattern Matching, Technical Report, UCI-ICS 00-10, Department of Information and Computer Science, University of California, Irvine, 18 pp.

Ge, X. and P. Smyth, 2000b, Segmental Semi-Markov Models for Change-Point Detection with Applications to Semiconductor Manufacturing, Technical Report, UCI-ICS 00-08, Department of Information and Computer Science, University of California, Irvine, 11 pp.

Goovaerts, P., 1997, *Geostatistics for Natural Resources Evaluation*, Applied Geostatistics Series, Oxford University Press, New York, pp. 483.

Gower, J. and F. Gonzalez, 2006, U.S. Warning System Detected the Sumatra Tsunami, *EOS, Transactions, American Geophysical Union*, 87 (10), pp. 105-108.

Grayman W., A. Ostfeld, and E. Salomons, 2006. Locating Monitors in Water Distribution Systems: Red Team - Blue Team Exercise, ASCE *Journal of Water Resources Planning and Management,* Vol. 132, No. 4, 300 -304.

Grayman, W.M., R.A. Deininger, and R.M. Males, 2001, *Design of early warning and predictive source water monitoring systems,* AWWA Research Foundation, Denver, 297 pp.

Hasan, J., S. States, and R. Deininger, 2004, Safeguarding the security of public water supplies using early warning systems: A brief review, *Journal of Contemporary Water Research and Education*, 129, 27-33.

Hall, J., A.D. Zaffiro, R.B. Marx, P.C. Kefauver, E.R. Krishnan, R.C. Haught and J.C. Herrmann, 2007, On-line water quality parameters as indicators of distribution system contamination, *American Water Works Association. Journal*; 99, 1; pg. 66

Hart, D.B., S.A. McKenna, K.A. Klise, V. Cruz and M. Wilson, 2007, CANARY: A Water Quality Event Detection Algorithm Development Tool, in proceedings of: ASCE World Environmental and Water Resources Congress, Tampa, FL, May 15-19[th]

Ihler, A., Hutchins, J., and P. Smyth, 2006, Adaptive Event Detection with Time–Varying Poisson Processes, in: Knowledge Discovery and Data Mining (KDD) proceedings

Jarrett, R., G. Robinson and R. O'Halloran, 2006, On-Line Monitoring of Water Distribution Systems: Data Processing and Anomaly Detection, , in Proceedings of the 8[th] Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006.

Journel A.G., and J. Ch. Huijbregts, 1978, *Mining Geostatistics*, Academic Press, London, pp. 600

Kalman, R. E., 1960, A New Approach to Linear Filtering and Prediction Problems, *Transactions of the ASME - Journal of Basic Engineering* Vol. 82: pp. 35-45.

Klise, K.A. and S.A. McKenna 2006a, Water quality change detection: multivariate algorithms, Proceedings of SPIE, volume 6203, 62030J.

Klise, K.A. and S.A. McKenna, 2006b, Multivariate Applications for Detecting Anomalous Water Quality, in Proceedings of the 8[th] Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006.

Koch, M.W. and S.A. McKenna, 2008, Distributed Network Fusion for Water Quality, in Proceedings of: ASCE World Environmental and Water Resources Congress, Honolulu, Hawaii, May 16-19, 10 pp.

Kroll, D. and K. King, 2006, Laboratory and Flow Loop Validation and Testing of the Operational Effectiveness of an On-line Security Platform for the Water Distribution System, in Proceedings of the 8[th] Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006

Lai, T.L., 1995, Sequential Changepoint Detection in Quality Control and Dynamical Systems, *Journal of the Royal Statistical Society: B*, 57 (4) 613-658.

Lueck, R.G., F. R. Driscoll and M. Nahon, 2000, A Wavelet for Predicting the Time-Domain Response of Vertically Tethered Systems, *Ocean Engineering*, 27 (12), pp. 1441-1453.

McKenna, S.A., D.B. Hart and L. Yarrington, 2006, Sensor Detection Limits on Protecting Water Distribution Systems from Contamination Events, *Journal of Water Resources Planning and Management*, 132 (4)

McKenna, S.A., K.A. Klise and M.P. Wilson, 2006, Testing Water Quality Change Detection Algorithms, in Proceedings of the 8[th] Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006.

McKenna, S.A., D.B.Hart, K. Klise, V. Cruz and M.Wilson, 2007, Event Detection from Water Quality Time Series, in proceedings of: ASCE World Environmental and Water Resources Congress, Tampa, FL, May 15-19[th]

Moradkhani, H., S. Sorooshian, H.V. Gupta, P.R. Houser, 2005a, Dual state–parameter estimation of hydrological models using ensemble Kalman filter, *Advances in Water Resources* 28, 2005, pp. 135–147

Moradkhani, H., K.-L. Hsu, H.Gupta, and S. Sorooshian, 2005b, Uncertainty assessment of hydrologic model states and parameters: Sequential data assimilation using the particle filter, *Water Resources Research*, 41, W05012, 17 pp.

Muller, K-R., A.J. Smola, G. Ratsch, B. Schokopf, J. Kohlmorgen and V. Vapnik, 1999, Using support vector machines for time series prediction, in: Advances in Kernel Methods: Support Vector Learning, MIT Press, Cambridge, MA, pp. 243-253.

NIST, 2008, National Institute of Standards and Testing, Univariate and Multivariate Control Charts, http://www.itl.nist.gov/div898/handbook/pmc/section3/pmc3.htm, accessed June, 2008.

O'Halloran, R., S. Yang, A. Tulloh, P. Koltun and M. Toifl, 2006, Sensor-Based Water Parcel Tracking, in Proceedings of the 8[th] Annual Water Distribution System Analysis Symposium, Cincinnati, OH, August 27-30, 2006

Raftery, A.E.,1994, Change Point and Change Curve Modeling in Stochastic Processes and Spatial Statistics. *Journal of Applied Statistical Science*, 1, pp. 403-424.

Roberson, J.A. and K.M. Morley, 2005, Contamination Warning Systems for Water: An Approach for Providing Actionable Information to Decision Makers, American Water Works Association, Denver, Colorado 22 pp.

Shang, F., J. Uber, R. Murray and R. Janke, 2008, Model-Based Real-Time Detection of Contamination Events, in proceedings of: ASCE World Environmental and Water Resources Congress, Honolulu, Hawaii, May 16-19, 10 pp

Takeuchi, J.-I. and Yamanishi, K., 2006, A Unifying Framework for Detecting Outliers and Change Points from Time Series, *IEEE Transactions on Knowledge and Data Engineering,* Vol. 18, No. 4, pp. 482-492.

Tsihrintzis, G.A. and C.L. Nikias, 1995, Robust Change Point-Detection and Segmentation in Data Streams, In Proceedings of: IEEE Military Communications Conference, 1995. MILCOM '95, Nov 5[th]-8[th], San Diego, California, Vol. 1, pp. 125-129.

Umberg, K. et al., 2008, ,Evaluation of Water Quality Event Detection Systems Deployed at the First Water Security Initiative Pilot Utility, presentation at AWWA Security Congress, Cincinnati, Ohio, April

Uber, J.G., R. Murray, M. Magnuson and K. Umberg, 2007, Evaluating Real-Time Event Detection Algorithms Using Synthetic Data, in proceedings of: ASCE World Environmental and Water Resources Congress, Tampa, FL, May 15-19[th]

Walpole, R.E. and R.H. Myers, 1989, *Probability and Statistics for Engineers and Scientists, Fourth Edition*, MacMillan Publishing Company, New York, 765 pp.

West, R.W. and R.T. Ogden, 1997, Continuous time estimation of a change-point in a Poisson process. Journal of Statistical Computation and Simulation, 56:293-302

Vaidyanathan K. and K. Gross, 2003, MSET Performance Optimization for Detection of Software Aging,14[th] IEEE International Symposium on Software Reliability Engineering, Fast Abstracts, Denver, Colorado, Nov 17-20.

Yang, Y.J., R.C. Haught, J. Hall, J. Szabo, R.M. Clark and G. Meiners, 2007, Adaptive Water Sensor Signal Processing: Experimental Results and Implications for Online Contaminant Warning Systems, in proceedings of ASCE World Environmental and Water Resources Congress, Tampa, FL, May 15-19[th]

Yuan,C.  C. Neubauer,  Z. Cataltepe and H.-G. Brummel, 2005, Support vector methods and use of hidden variables for power plant monitoring, in Proceedings: IEEE International Conference on Acoustics, Speech, and Signal Processing. (ICASSP '05), Vol. 5, pp. 693-696.

Zavaljevskl and Gross, 2000, Sensor Fault Detection in Nuclear Power Plants Using Multivariate State Estimation Technique and Support Vector Machines, in proceedings of the Third International Conference of the Yugoslav Nuclear Society YUNSC 2000, October 2-5, 2000, Belgrade, Yugoslavia

Zhang, N. F. (1997). Detection Capability of Residual Chart for Autocorrelated Data, *Journal of Applied Statistics*, 24, 475-492.

## *Glossary*

### *Baseline Change*
A baseline change is a significant change in a statistical parameter, generally the mean, of the observed data.  The point at which a baseline change occurs is a change point.  The change in behavior of the observed signal from one side of the change point to the other is referred to as the baseline change.  Baseline changes are common in some water distribution systems due to changes in mixing of source waters within the network at different time of the day.

### *Change Point*
Change points are defined as the point in time or location where an abrupt change in behavior or mode of operation is observed.  Change points are generally identified by comparing data collected on both sides of the proposed change point.  If that comparison shows the data on either side of the proposed change point to be significantly different, then that proposed change point is confirmed as a change point.  Change points are most accurately identified using offline techniques to data analysis where adequate amounts of data have already been collected on both sides of any proposed change point.

### *Contaminant Warning System (CWS)*
The CWS is comprised of all the hardware and software components that are necessary to monitoring the water distribution system for contamination events.  These components include the sensors, the SCADA system to collect and transmit the data from the sensors to a central location, the database to hold sensor information and the event detection

system (EDS) that processes the data to provide some indication of the occurrence of an event.

### Event
An activity or behavior that is unusual relative to normal modes of operation. Abnormal activity or operation relative to the background or ambient modes of operation. An event is a sustained period of such abnormal activity that is of a longer duration than an outlier, but of shorter duration than a baseline change.

### Event Detection System (EDS)
The software system that contains the data handling, algorithms, and input/output functions necessary to identify events from water quality time series. Typically, the inputs to an EDS are the water quality data streams from the sensors as stored in a SCADA database. The output of an EDS is an indication of the state of the water quality. This indication can be a binary signal such as "alarm/no-alarm" or it can be a continuous indication of the water quality state such as the probability of an event occurring at every time step. The EDS itself is a component of the more comprehensive Contaminant Warning System (CWS).

### Outlier
An outlier is defined here to be a single time step with behavior that is considered anomalous relative to the background or expected behavior for that time step. A large enough number of outliers within a prescribed time interval may constitute an event.

### Residual
The difference between the predicted and observed water quality values at a single time step