

# MONITORING, FAULT DIAGNOSIS, FAULT-TOLERANT CONTROL AND OPTIMIZATION: DATA DRIVEN METHODS

**John MacGregor\* and Ali Cinar#**

\*ProSensus, Inc. 1425 Cormorant Rd., Ancaster, ON Canada

*john.macgregor@prosensus.ca*

#McMaster University, Chemical Engineering Dept., Hamilton, ON

*macgreg@mcmaster.ca*

#Department of Chemical and Biological Engineering, Illinois Institute of Technology, Chicago, IL 60616 *cinar@iit.edu*

## *Abstract*

Historical data collected from processes are readily available. This paper looks at recent advances in the use of data-driven models built from such historical data for monitoring, fault diagnosis, optimization and control. Latent variable models are used because they provide reduced dimensional models for high dimensional processes. They also provide unique, interpretable and causal models, all of which are necessary for the diagnosis, control and optimization of any process. Multivariate latent variable monitoring and fault diagnosis methods are reviewed and contrasted with classical fault detection and diagnosis approaches. The integration of monitoring and diagnosis techniques by using an adaptive agent-based framework is outlined and its use for fault-tolerant control is compared with alternative fault-tolerant control frameworks. The concept of optimizing and controlling high dimensional systems by performing optimizations in the low dimensional latent variable spaces is presented and illustrated by means of several industrial examples.

## *Keywords*

Multivariate statistical process monitoring, Latent variable models, Fault diagnosis, Agent-based systems, Fault-tolerant control, Optimization, Control, Batch processes

## **Introduction**

The optimization, control and monitoring of processes involves employing models that enable us to learn from the data being collected from the process. These models could be models whose structure is based on fundamentals and whose parameters are estimated from plant data (mechanistic models), or they could be models whose structure and parameters are all identified from plant data (data-driven or empirical models). The key issue is not the type of model used, but whether or not that model, in

terms of its structure and assumptions, is appropriate for the application. For example a mechanistic model imposes a structure that embodies many assumptions, some of which may not be entirely justified. In particular, assumptions are needed about the structure of the disturbances in the system (rarely available from theory), and many information-rich variables (such as the mechanical parts of the system – e.g. agitator torque, vibration sensors, etc.) that the modeler does not know

\* To whom all correspondence should be addressed

how to incorporate into the mechanistic model are often omitted. Empirical models can easily capture these latter two sources of variation, but if the structure of the model is not properly addressed, empirical models can provide misleading results.

In this paper, we focus on the proper use of data-driven (empirical) models for the monitoring, control and optimization of processes. In particular, we focus on latent variable models because, as we will show, they provide the proper structure to allow them to be built from plant data and be used for monitoring, control and optimization. But, the nature of the data will always determine the limitations of these models and one of the themes of this paper is the discussion of the limitations imposed by the available data. Although some of these issues are not present with the use of mechanistic models, the nature of the available data still has a major impact on the ability to independently estimate many of the mechanistic model parameters.

### 1.1. Causality

Perhaps the major issue with data-driven models is the issue of whether or not they model causality among the variables and if so, what variables are related causally. We say that a model causally relates two variables if it correctly shows that a change of a certain magnitude in one will result in a change of a certain magnitude of the other. In data-driven models causality among variables is determined entirely by the nature of the data and by the structure of the empirical model. If independent variation is not present in certain manipulated variables, then no causality information for the effects of those individual variables will be present in the data, nor in any model built from them. However, as we will discuss, if a proper structure is used for the empirical model (namely, a latent variable structure), a causal model may be obtained from such process data, although only in a reduced dimension of the latent variable space.

Causal models are not always useful or even desirable in some situations (e.g. passive applications such as monitoring and soft sensors), but are critical in other situations (e.g. active applications such as control and optimization). In situations where the model is to be used in a passive sense, such as monitoring or soft sensors, one actually wants a non-causal model, one that simply models the correlation structure existing among all the variables in the plant during normal operation where only “common cause” variation is present. In monitoring, the concept is to capture in the model the acceptable “common cause” variations in the process and use the model to detect any deviations from such behavior. With soft sensors (inferential models) the concept is to enable prediction of a variable of interest from data collected under such routine plant operation. However, in situations such as control and optimization where the model is to be used actively to alter the operation of the process, causal models are required. Similarly, for fault diagnosis or interpretation of causal effects among variables some form of causality is required.

### 1.2. Changing nature of data.

Over the past few decades with the advent of process computers and LIMS systems, companies have collected massive amounts of routine plant data. These data are of a very different nature from typical R&D data that are usually collected under designed experiments. Almost all statistical texts are aimed at the analysis of this latter type of data where all the variables are independently varied. Collecting such data on a process involves major identification experiments whereby independent variation is introduced into all manipulated variables. Data collected under routine operation are unlike these data. The number of measured variables is often very large, and most of the variables are highly correlated because their variation is due to a small number of underlying variations (latent variables) such as raw materials, environmental factors or normal process variations introduced in combinations of variables by operating personnel. These variations in the process data define a causal subspace within which the process moves, but they do not provide causal information on individual variables. This issue lies at the heart of defining useful data-driven models developed from these data. Data-driven models, such as standard statistical regression models and artificial neural network models that do not explicitly recognize the nature of these process data are of limited or no value to the engineer trying to use these data.

### 1.3. Concept of Latent Variables and Latent Variable Models

Latent variable (LV) models such as PLS (Partial Least Squares or Projection to Latent Structures) are unique among regression methods in defining the high dimensional regressor and response spaces (X and Y) in terms of a small number of latent variables (T) that define the major directions of variation in the process data. The basic LV model is defined as:

$$X = TP^T + E \quad (1)$$

$$Y = TC^T + F \quad (2)$$

where X are ( $n \times k$ ) and ( $n \times m$ ) matrices of observed values, and  $T = XW^*$  is an ( $n \times a$ ) matrix of latent variable scores ( $a \ll k$ ). P, C and  $W^*$  are matrices of loadings estimated from the data.  $n$  is the number of observations, and  $k$  and  $m$  are the number of regressor and response variables. The number of statistically significant latent variables  $a$  is determined by methods such as cross-validation, jackknifing, etc. The selection of which variables are  $x$ 's and which are  $y$ 's is the user's choice and is only dependent upon the objective for which the model is being developed. All X and Y data are just measurements collected from the process with errors (E and F) and the concept of dependent and independent variables has little place in latent variable models. The model is symmetric in X and Y in that it defines both X and Y in terms of underlying LVs (T) and does not force any assumed *a priori* directional dependence among them. The scores are defined in terms of the X's ( $T = XW^*$ ) simply because these variables (X) are usually assumed to be the ones

available when the model is to be used. (Maximum Likelihood LV methods that define T's in terms of both X and Y are less useful for this reason). The fitting of LV models to plant data is always a well-conditioned and low dimensional problem because the LVs are orthogonal, and only as many are used as are statistically important to the objective.

Perhaps the most subtle but critical difference between LV models and all other regression models is that they are the only ones that simultaneously model the X and Y spaces. All other regression models only model the Y space, something that is acceptable only if the X space is of full statistical rank (i.e. there is meaningful independent variation in all x variables). A model for the X space is essential if the effective dimension of the X space is less than the number of X variables. It is this simultaneous modeling of the X and Y spaces that leads to unique solutions for the LV models and non-unique models (an infinite number of solutions) for all other regression models. Through this joint modeling of X and Y LV models are the only regression models that are unique and interpretable. This also allows them to handle significant amounts of missing data (the algorithms automatically impute the missing data to lie on the low dimensional X and Y model planes), and it allows for checking the integrity of new incoming data through tests that check whether the new observations lie on the model planes (SPE) and within the LV region defined by the training data ( $T^2$ ). These latter points make them ideally suited for soft sensor applications and for monitoring processes (section 2). Another underappreciated point is that LV models built from plant data do provide causal models in the low dimensional LV space (this issue is critical for any active use of the models for optimization and control – see section 3).

#### *1.4. Other regression methods.*

There has been a proliferation of papers in recent years that propose many other data-driven methods for building models from process data for inferential models, process monitoring, etc. These include Independent Component Analysis (ICA), Artificial Neural Networks (ANN), and Support Vector Machines (SVM). In our opinion these methods fall within the class of regression methods/classifiers that provide no allowance for modeling the X-space and thereby assume the data to be full rank in any interpretation or use of the models. Although we recognize that these methods can be useful in some cases, even with process data, they do not provide unique models, nor allow for interpretation, nor provide any form of causality. They also have limited ability to handle missing data or test for outliers in new data.

### **Monitoring and Diagnosis**

The purpose of process monitoring is to detect abnormalities in process operation. These may range from

sensor and actuator faults to more complex process problems such as catalyst poisoning or fouling. If abnormal (out-of-control) operation is declared by the monitoring system, the next step is to find the source cause of the deviation (fault diagnosis). Fault diagnosis can be conducted by associating process behavior patterns to specific faults or by relating the process variables that have significant deviations from their expected values to various process states such as catalyst poisoning or equipment that can cause these deviations.

A large variety of techniques have been proposed for the first approach to associate directly the trends in process data to faults (Patton et al., 1989, Frank, 1990, Gertler, 1998, Venkatasubramanian et al., 2003). Their common characteristic is that they can be implemented on closed sets: The set of all faults to be identified are listed and associations between data and faults are created. This association can be made by using mechanistic models, black-box models such as ANNs and hidden Markov models (HMM), or statistical classification techniques such as discriminant analysis and SVMs. Their success depends on the type of process and the modeling approach used. For electro-mechanical processes that can be described accurately with mechanistic models, model-based residuals can be used for fault diagnosis (Patton et al., 1989, Gertler, 1998). Their successful use has also been reported in simulated chemical systems for diagnosis of sensor or actuator faults (Mhaskar et al., 2006, El-Farra and Ghantasala, 2007). However, their performance would be limited for processes that cannot be described accurately by mechanistic models. Often this is due to uncertainties, time-varying or nonlinear behavior, complexity of the process that makes the development and maintenance of the mechanistic model too expensive. A detailed comparison of traditional fault detection and diagnosis (FDD) approaches and multivariate LV SPC approaches contrasts the strengths and limitations of both classes of techniques (Yoon and MacGregor, 2000). ANN, HMM or SVM based diagnosis systems have been very successful in various fields such as medical diagnostics. But their success in chemical process diagnosis has been limited and they have a high computational (training) cost when additional source causes need to be added to the fault set or multiple simultaneous faults need to be diagnosed.

Contribution plots provide an indirect approach to fault diagnosis by first determining process variables that have inflated the monitoring statistics ( $T^2$  or *SPE*). These variables are then related to equipment, process behavior and disturbances. This approach relies on the availability of expert plant personnel that can interpret the contribution plots and associate the trends in the plots with process and equipment behavior. Knowledge-based systems can facilitate and automate the interpretation and association effort (Cinar et al., 2007). This approach has significant advantages in diagnosing complex process problems. The tradeoff is the availability of expert plant personnel as

opposed to an automated monitoring and diagnosis tool.

### 2.1. Multivariate statistical process monitoring

Multivariate statistical process monitoring (MSPM) methods are gaining acceptance in industry because they provide more accurate information about the process, give warnings earlier than the univariate methods, and rely on statistics that are easy to compute and interpret. MSPM relies on statistical distance concepts expressed by Hotelling's  $T^2$ , representing the deviation within the model plane of a new observation on the process from its desired state, and the Squared Prediction Error (SPE), representing the residual or squared distance of the new observation from the model plane. If the process has a few variables that are independent,  $T^2$  can be computed by using all process variables. If the number of variables is large and there is significant collinearity among some of them, principal components analysis (PCA) or PLS can be used. If only process variables are used for monitoring, MSPM charts are based on principal components (PC). When both process and quality variables are used, and the two blocks of data need to be related as well, the MSPM charts are based on the latent variables (LV) of PLS. Both PCA and PLS based charts summarize the information about the status of the process by using two statistics, the  $T^2$  and the squared prediction error (SPE) computed by using the information collected at each sampling time. The  $T^2$  chart indicates the distance of the current operation from the desired operation as captured by the PCs or LVs included in the development of the PCA or PLS model of the process. Since only the first few PCs or LVs that capture most of the variation in the data are used to build the model, the model is a somewhat accurate but incomplete description of the process. The SPE chart captures the magnitude of the error caused by deviations resulting from events that are not described by the model. The  $T^2$  chart indicates a deviation based on process behavior that can be explained by the model while the SPE chart indicates a significant deviation that cannot be explained by the model (the prediction error is inflated). The  $T^2$  and SPE charts must be used as a pair and if either chart indicates a significant deviation from expected operation, the presence of an abnormal process operation must be declared.

The  $T^2$  statistic based on process variables at sampling time  $k$  is

$$T^2 = (x - \bar{x})^T S^{-1} (x - \bar{x})$$

where the vector of mean values  $\bar{x}$  and covariance matrix  $S$  are estimated from process data.

The  $T^2$  charts based on PCs use

$$T^2 = t_a^T S^{-1} t_a$$

and follow an F or a Beta distribution with  $a$  and  $n-a$  degrees of freedom for the F distribution, and  $a/2$  and  $(n-a-1)/2$  for the Beta distribution, assuming that the data follow a multivariate Normal distribution (Jackson, 1991).  $a$  denotes the number of PCs,  $t_a$  is a vector of  $a$  scores and  $S$  is the  $a \times a$  estimated covariance matrix.

The orthogonal distance of the observation  $x(k)$  from the projection space is the prediction error  $e(k)$  gives a measure of how close the observation at time  $k$  is to the  $a$ -dimensional space and is used for computing  $SPE(k)$ .

$$SPE(k) = e(k)^T e(k) = \sum_{j=1}^n e_j^2(k) = \sum_{j=1}^n (x_j(k) - \hat{x}_j(k))^2$$

where  $\hat{x}_j(k)$  is computed from the PCA or PLS model.

To include the information about process dynamics in the models, the data matrix can be augmented with lagged values of data vectors, or model identification techniques such as *subspace* state-space modeling where the state variables correspond to LVs can be used (Negiz and Cinar, 1997).

### 2.2. Fault Diagnosis

#### 2.2.1. Fault diagnosis using contribution plots

When  $T^2$  or  $SPE$  charts exceed their control limits to signal abnormal process operation, variable contributions can be analyzed to determine which variable(s) caused the inflation of the monitoring statistic and initiated the alarm. The variables identified provide valuable information to plant personnel who are responsible for associating these process variables with process equipment or external disturbances that will influence these variables, and diagnosing the source causes for the abnormal plant behavior. The procedure and equations for developing the contribution plots are given in Kourti and MacGregor (1996). The computed values of the contributions of each process variable are plotted on a bar chart for comparison to determine which variable(s) caused the abnormal operation in the MSPM charts.

Investigating the dynamic pattern of contribution plots is more effective in fault diagnosis rather than the single time snapshot of contributions. Contribution plots can be plotted over time following an alarm signal in MSPM charts. The variation of the contributions over time can also be summarized by plotting sum the contributions over a time period (Undey et al., 2003). Rapid detection of the variables responsible for inflating the monitoring statistics is necessary because the contributions smear over time as the effects of the abnormality spreads over other variables. On the other hand, inspection of contributions over a period is desirable to filter out instantaneous spurs caused by measurement noise or errors.

#### 2.2.2. Fault diagnosis with statistical methods

When a process can be represented by a few PCs, the biplots of PCs and SPE provide a visual aid to identify data clusters that indicate normal operation or operation under a specific fault. For processes that need to be described by a higher number of PCs and for automation of diagnosis PCA and discriminant analysis techniques can be integrated (Raich and Cinar, 1996). The FDD system design includes the development of PC models for normal operation (NO) and abnormal operation with specific faults, and the computation of threshold limits using historical data sets collected during normal plant operation

and operation under specific faults. If score or residual tests exceed their statistical limits, the PC models for all faults are used to carry out the score, residuals, and/or angle tests to determine the proximity of current process operation to one of the data clusters indicating a specific fault. Discriminant analysis is performed by using PC models for various faults to diagnose the source cause of abnormal behavior.

The angles between principal coordinate directions of current data and regions corresponding to operation with different faults can be used for fault diagnosis, to complement distance-based methods (Raich and Cinar, 1997). The method uses angles and a similarity index defined by using the angle information (Krzanowski 1979). A suitable discriminant could be the minimum angle between the test point and the mean of the various clusters of fault data, with the vertex positioned at the mean of NO. Decision boundaries for angular discriminants describe open-ended conical regions in space.

A likely cause for abnormal behavior can be assigned by pattern matching by using scores, residuals, angles, or their combination. Combining the information available in scores and residuals usually improves the diagnosis accuracy. If none of the known fault models provide an adequate match to the observations. These observations could be considered as a new group, modeled, and added to the discrimination scheme. Unlike NN or SVM based diagnosis methods, additional source populations can easily be incorporated into the discrimination scheme without retraining the whole diagnosis system.

Detection and diagnosis of multiple simultaneous faults is an important concern. Most FDD techniques rely on the assumption of a single fault. In a real process, combinations of faults may occur. An intervention policy to improve process operation may need to take into account each of the contributing faults. Diagnosis should be able to identify major contributors and correctly indicate which, if any, secondary faults are occurring (Raich and Cinar, 1995). Since process behavior due to different faults is described by different models, it is useful to have a quantitative measure of similarity or overlap between models, and to predict the likelihood of successful diagnosis. Similarity measures can identify combinations of faults that may be masked or falsely diagnosed, and provide information about the success rates of different diagnosis schemes incorporating single and combinations of faults. Using these guidelines, multiple faults occurring in a process can be analyzed a priori with respect to their components, and accommodated within the diagnosis framework.

When the region spanned by the model for one (outer) fault contains the model for another (inner) fault, their combination may not be easily diagnosed. Idealizing the two fault regions as concentric spheres, the inner model region is enveloped by the outer model. Consequently, the outer fault will be diagnosed and the inner fault will be

masked. Faults causing random variation about a mean (such as sensor noise) move a process less drastically off-target than step or ramp faults. Ramp or step faults tend to be the outer models and mask secondary random variation faults. Overlap of regions is likely to exist for most processes under closed-loop control, the multiple fault scenario is further complicated for such processes.

### 2.3. Integration of Monitoring and Fault Diagnosis

It is difficult to predict SPM and FDD techniques that would perform effectively for all possible abnormalities a process can experience. A monitoring and FDD system that includes several techniques as alternatives, dynamically analyzes the performance of every technique in detection and diagnosis of different faults, and selects and prioritizes the most effective technique for each fault scenario is desirable.

A hierarchical agent-based system can create a flexible environment for automated and adaptive monitoring and FDD of complex chemical processes. A combined monitoring and FDD environment with agent-based systems has been developed at IIT as a part of a complex framework for Monitoring, Analysis, Diagnosis and Control with Agent-Based Systems (MADCABS). MADCABS contains agents with several alternative SPM and FDD techniques and an agent management layer that dynamically identifies the best performing agents and techniques, and provides accurate results by using agent performance evaluation, performance-based consensus, and adaptation. MADCABS consists of three hierarchical layers (Figure 1). In the *physical communication layer*, software representing the process units and their connections, and sensor and actuator representations are generated. Communication agents in this layer manage the communication between MADCABS and a process or simulator and update the values in the sensors and the actuators. The *process supervision layer* has modules for data preprocessing, SPM, FDD, control and decision-making. Each of these modules includes a number of agents with different techniques for the same task and they collaborate with each other during the execution of their tasks. Manager agents in the *agent management layer* monitor the performances of agents in process supervision layer under specific states of process operation, rate their performances and adjust the confidence level to an agent based on past performances under similar operating conditions.

Three different PCA-based techniques are used to illustrate this *monitoring* approach. Agents that use PCA, multi-block PCA (MBPCA) and dynamic PCA (DPCA) are implemented where PCA and DPCA agents monitor each unit in the process and a MBPCA agent monitors the entire process using data from each unit.  $T^2$  and  $SPE$  are used by monitoring (fault detection) agents to detect abnormal operation. Each monitoring agent generates two monitoring statistics, ( $T^2$  and  $SPE$ ). Hence, six statistics are observed by six different fault detection agents, and up

to six flags can be raised when the monitoring statistic goes outside the control limits. A fault detection organizer (FDO) keeps track of all the fault flags given by its fault detection agents, declares consensus fault based on a consensus criterion, and triggers the diagnosis agent. Historical performance-based consensus scheme yields fewer missed alarms and faster detection than a voting-based scheme, since the system has the capability to predict which monitoring agents should be assigned a higher performance value during consensus on the basis of their experiences for similar faults in history (Perk et al., 2010). Combinations of monitoring agents yield fewer false alarms since it is unlikely for all good techniques to flag a false alarm simultaneously. The most effective criterion was the time-averaged-performance-with-history criterion, which uses information about the performances of fault detection agents for similar fault magnitudes in history and forms consensus decision on the basis of the expected performances of fault flagging agents.

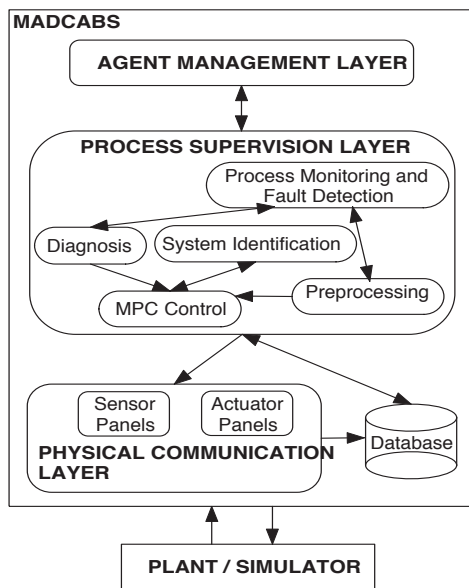


Figure 1: The interlayer and intralayer information flow in MADCABS

After an abnormality is detected by the FDO, the diagnosis organizer (DO) is triggered for the faulty unit. DO is responsible for getting estimates from different fault discrimination agents, requesting performance estimates from the diagnosis manager for each fault discrimination agent for the current fault state, forming a consensus diagnosis decision via a performance-based criterion, and triggering agent adaptation (Figure 2). Classification techniques such as Fisher's discriminant analysis (FDA) and PLS discriminant analysis (PLSDA) and diagnosis tools such as the variable contribution plots are implemented in MADCABS (Perk et al., 2011). For example, a PLSDA agent collects data from the unit for specific fault types, builds a model relating the type of fault to the values of process variables, and uses this model

to classify the type of fault for new data after the FDO declares the consensus fault decision. DO gathers the classification results from all fault discrimination agents for the unit and forms a consensus classification decision on the basis of the agents' historical performances. If a classification agent has a higher misclassification rate for a certain fault (lower performance), it is assigned lower confidence in the consensus and the agent re-trains its model by including the new observation to become more reliable for that type of fault in time via adaptation.

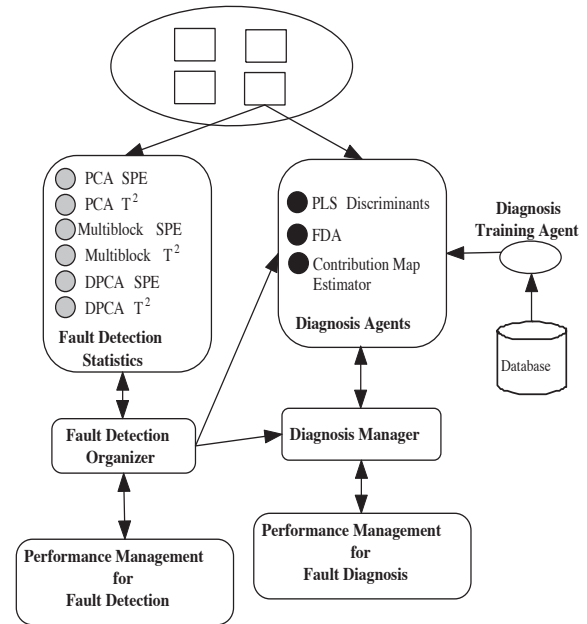


Figure 2: Fault detection and diagnosis agents

The diagnosis training agent contains process data for each known fault type in the process (Figure 2). All of the fault discrimination agents require historical data of all possible faults. The diagnosis training agent is responsible for collecting process data under each fault and building a fault data matrix in the training phase. MADCABS also allows human operator intervention when unknown faults occur in the system and the agents re-train their models in real-time by incorporating the process data for newly defined fault. For example, PLSDA agent uses the fault data matrix and builds another binary matrix to represent the fault type. PLSDA relates the fault data to the respective fault under which the fault data has been collected. FDA uses the same fault data from the training agent as PLSDA to create different clusters of faults. The multivariate observations are transformed to another coordinate system that enhances the separation of the samples belonging to each class.

Contribution plots identify the process variables that have contributed significantly to the inflation of  $T^2$  and  $SPE$  statistics. In practice, it is necessary to relate these process variables to various faults. Automation of this process has been proposed in the literature using knowledge-based systems (Tatara and Cinar, 2002). In MADCABS fault diagnosis module, the contribution

values are used to list the process variables with large contributions to the inflation of the statistic that exceed the 3-standard-deviation confidence limits of NO contribution values for each statistic. Since six different monitoring statistics are used for fault detection, six different sets of variable contributions are available. Only the variables that have been listed in the majority of the sets are kept as a mapping of the fault signature on process variables to the fault type. In the training phase, variable contributions are calculated and fault signature is extracted for each known fault and a mapping is created. The lists of fault signature variables are mapped against the respective fault under which the data has been collected. This contribution-plot-based classification map is used by the variable contribution map estimator. The reason to use the variable contribution plots in this closed set scheme without incorporating human expert knowledge is to provide automation to the classification and diagnosis without a KBS.

### Fault-tolerant Control

Fault-tolerant control (FTC) aims to prevent catastrophic consequences of faults by retuning or restructuring the control system to maintain acceptable process operation in spite of drastic faults. Control actions are generated to satisfy the process objectives using the process information from functional sensors and manipulating the available actuators. Based on failures diagnosed, the control system is retuned or reconfigured to achieve FTC. Early methods for FTC relied on hardware redundancies. As industry transitioned digital control and computer control, software redundancies and soft sensors gained popularity for FTC. Use of robust control framework and its integration with knowledge-based systems was also proposed for FTC (Kendra et al., 1997).

FTC begins with the detection of a fault and its nature. Depending on the fault, the models used as reference for monitoring and control may be updated, soft sensors may be introduced and an alternative controller is developed and implemented (Maciejowski, 1999, Wang et al., 2007, Zmoffen and Basualdo, 2008). Some FTC approaches design a priori the controllers that are appropriate for specific modes of process operation and switch to that specific controller once it is determined that the process has drifted to the corresponding mode (Du et al., 2011, Mhaskar et al., 2006, Mhaskar et al. 2007, El-Farra et al., 2005). Other FTC techniques, such as the MADCABS framework, conduct on-demand model identification and controller design once an abnormality is detected and a fault is diagnosed. The first approach enables the development of controllers for a closed set of process variations while the latter approach can accommodate unknown faults as well. The tradeoff is the ability to conduct extensive stability assessment a priori. Several methodologies for control retuning and reconfiguration are proposed in the literature (Christofides

and El-Farra, 2005).

Decentralized supervision and control systems became a powerful alternative to traditional centralized supervision and control (Rawlings and Stewart, 2008). Distributed intelligence and decision-making provide local actions that can take into account global priorities and constraints. Distributed control systems can prevent the spread of disturbances in the process, reduce the smearing of failure signature patterns, and deliver nonlinear intervention strategies that can be hard or impossible to achieve with centralized techniques.

MADCABS provides a platform for the integration of FDD and distributed control to provide FTC. In MADCABS (Figure 1), monitoring and diagnosis module, controller performance evaluation module, and system identification module communicate and coordinate their activities with the MPC module. Each module may contain a single agent or multiple competing agents that communicate with other agents while they perform their tasks and learn from the effects of their actions. Controller performance is affected when faults occur, and the controller performance evaluation agent informs the fault-tolerant MPC agent when such degradation is detected. Next, the controller agent communicates with monitoring and diagnosis agents to acquire fault information. Then the FTC generates new actions on the actuators and the controller performance is assessed. MADCABS control agents utilize a state-space model of the process developed by using data-based subspace algorithms. These agents learn the best control strategies for different fault types and develop preferences on the available inputs. This agent-based supervision and FTC framework enables dynamic evaluation of the performances of agents under changing operating conditions, learning and adaptation based on historical actions, performances and experience, on-demand process model identification, controller retuning and restructuring, and the improvement of the overall performance of the combined framework for FDD and control during process operation with performance-based consensus building and adaptation.

Controller performance monitoring determines if the control system is working satisfactorily (Huang and Shah, 1999, Kendra and Cinar, 1997, Schaefer and Cinar, 2004). For MPC performance evaluation, in MADCABS a historical performance index  $\lambda_{hist}(k) = J_{hist}(k)/J_{ach}(k)$  is used where  $J_{hist}(k)$  and  $J_{ach}(k)$  are the values of the historical and achieved cost function at time  $k$ , respectively. Then the historical performance index can be fitted into an autoregressive (AR) model and its residual  $e = A(q-1)\lambda_{hist}(k)$  can be monitored. If controller performance degradation is detected, then the controller will communicate with the diagnosis agent and check the diagnosis message for possible faults, and then execute the FTC actions according to the faults diagnosed. Controller performance information is also used by the agent-performance-manager agents to compare the performances of different control strategies for each fault.

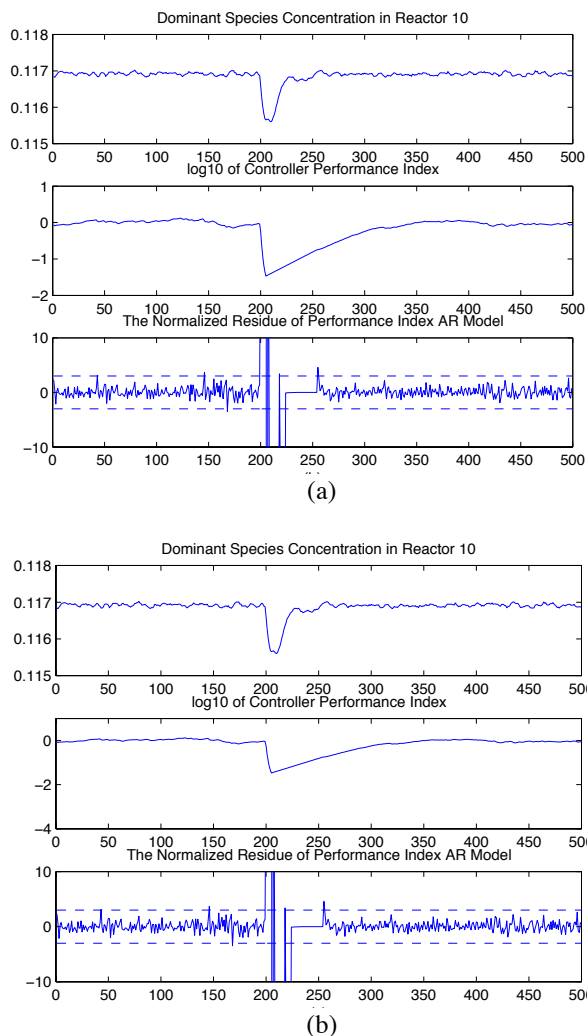


Figure 3. FTC of process with multiple controlled variables. (a) strategy 1, (b) strategy 2.

The agent-based FTC can implement different strategies to handle different types of faults. The strategies are prioritized such that the control strategy that utilizes more prior information has higher priority. The control strategies are not guaranteed to work under all possible operating conditions. Therefore, the controller performance evaluation agents monitor the controllers and if a selected strategy is not successful, the FTC would switch back to original settings and try the next control strategy. For actuator faults, the performance of two strategies in decreasing priority

- a. Treat actuator faults as unmeasured disturbance,
- b. Update system model by system identification

are illustrated for regulating the concentration of a product in one reactor that is part of a network of 20 reactors (Figure 3). Both strategies are successful for the test case where the system operates at NO until time = 200, when the actuator of the interaction flowrate from Reactor 11 to Reactor 10 is stuck at a fixed value which is 27% larger than the nominal value until the end of the simulation. The objective of the fault-tolerant MPC is to recover and maintain the dominant species concentration in Reactor 10

despite the failing flowrate actuator while maintaining acceptable operation in all neighboring reactors. If a similar faults occurs again in the future, MADCAB will recall these performances and will execute the first strategy which necessitates less design effort.

### Optimization in the Latent Variable Space

In this section we discuss the use of models built from historical plant data for optimization and control. This is a topic of increasing interest since it is becoming ever more difficult in industry to justify the use of designed experiments on production scale plants in order to generate the causal models needed for conventional optimization. As discussed in the introduction section, data from historical plant operation do not contain causal information on all the manipulated variables as these are rarely all varied independently of one another. This becomes quite clear when building latent variable (PLS) models on historical operating data. The effective number of independent dimensions in the process is the number of latent variables uncovered in the analysis of the data and this dimension is usually substantially less than the number of process variables.

However, what one does have is a causal model in the reduced dimensional space defined by the new latent variables. The Y data are related to the latent variables (T) via equation (2), and so if the latent variables ( $t_1, t_2, \dots$ ) can be moved, then the Y variables will be changed according. However, the latent variables are not physical variables that can be moved directly. To move these latent variables one must move the combinations of the process variables that define the latent variables ( $T=XW^*$ ). Hence by manipulating this reduced combination of X variables that define the latent variables (T) one can affect changes in the Y variables. Some of these X variables are manipulated variables that can be changed, others are observed disturbances (e.g. raw material properties, ambient temperature, etc.) that are to be held at their current measured values, and others are measured responses (states) that result from these disturbances and manipulated variables. New values of all the individual variables in X that can achieve desired new settings of the latent variables can be obtained from the X-space latent variable model (1).

Therefore an optimization with the latent variable model consists of solving for those values of the latent variables (using the Y-space model) that will achieve the desired y values, and at the same time for those values of the x variables that satisfy the X-space latent variable model together with any constraints on the x variables. The optimizations are carried out in the low dimensional LV space, and yet they effectively provide an optimization in the high dimensional process variable space. The concepts are best illustrated by a few examples.

#### *Optimization of an over-injection molding process*

Yacoub and MacGregor (2004) developed and



implemented a periodic re-optimization scheme on an industrial over-molding injection process. The problem was that periodic changes in raw material lots and changing ambient temperature and humidity conditions (disturbance X variables) in the plant had significant effects on the levels and standard deviations of 10 quality variables. To counter this plant operators periodically altered the operating conditions of the process (e.g. injection velocity profiles, timing sequences, etc.) by injecting a number of parts, sending the parts to the QC lab for analysis and iterating until the quality was back into an acceptable region.

A nonlinear PLS model was built using the data accumulated from these correction periods (8-20 raw material properties, 26 process variables, 10 quality variables and their std. deviations over multiple injections). Four latent variables were sufficient to provide a very predictive model of the process ( $Q_Y^2 = 0.88$ ). A multivariate SPC scheme was set up using  $T^2$  and SPE statistics. Every time the  $T^2$  went outside its control limits, an optimization was then initiated that moved the latent variables back to the origin of the latent variable space (centre of the control region). The re-optimization was performed in the 4 dimensional LV space using a quadratic objective function in the quality variables and their standard deviations:

$$\underset{t_{new,a}}{\text{Min}} \left[ (y_{des} - \hat{y}(t_{new}))^T Q_1 (y_{des} - \hat{y}(t_{new})) + \hat{y}_\sigma^T(t_{new}) Q_2 \hat{y}_\sigma(t_{new}) \right]$$

subject to constraints that: (i) the raw material properties and the ambient temperature and humidity be those present at the time, (ii) the SPE be close to zero (ensuring the validity of the model); (iii) the  $T^2$  lie within the 99% control limit (to avoid extrapolation). The resulting latent variable optimization/control scheme has allowed the process to be periodically brought back into its control region within one injection period by sending down a complete set of new process conditions to be simultaneously implemented. The scheme has been in operation for many years and has resulted in greatly improved quality, reduced variation and reduced scrap rates.

#### Optimization of batch operating policies.

Consider the batch trajectory histories from a large pilot plant batch polymerization process shown in Fig. 3 (Garcia-Munoz et al, 2008). These batch operating trajectories (X) together with corresponding chemical recipes (Z) and final polymer quality (Y) data were collected for all the batch runs on the reactor. The optimization problem is to solve for the complete 9 new process variable trajectories ( $X_{new}$ ) and new recipes ( $z_{new}$ ) that will meet desired specifications on 13 quality variables ( $y_{new}$ ). In traditional terms this represents a large optimization problem involving a very large number of variables. With latent variable models it represents an

optimization in a 3 dimensional LV space (the dimension of the PLS model summarizing all the data):

$$\underset{t_{new,a}}{\text{Min}} \left[ (y_{des} - \hat{y}(t_{new}))^T Q_1 (y_{des} - \hat{y}(t_{new})) + q_2 T^2 + q_3 SPE + q_4^T \hat{x} \right]$$

$$(\hat{y}, \hat{x}, SPE, T^2) = PLS(t)$$

$$SPE \leq a_1$$

$$T^2 \leq a_2$$

$$C\hat{x} \leq a_3$$

$$D\hat{y} \leq a_4$$

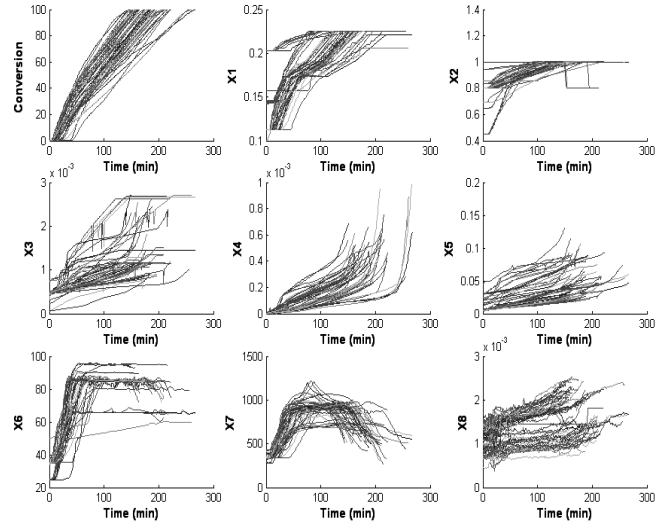


Fig. 4. Process operating trajectories for all runs in a pilot plant batch polymerization process

The specifications on the product quality were quite loose (mainly upper and lower constraints) except for certain  $y$ 's with specified values. As a result there were multiple solutions for  $z_{new}$ , and  $x_{new}$ , five of which are shown in Fig. 5 for increasing penalty on time usage. All of these satisfied the specification on the desired quality  $y_{des}$ . The case 5 trajectory set (bold curves) were used as they provided the minimum batch time solution.

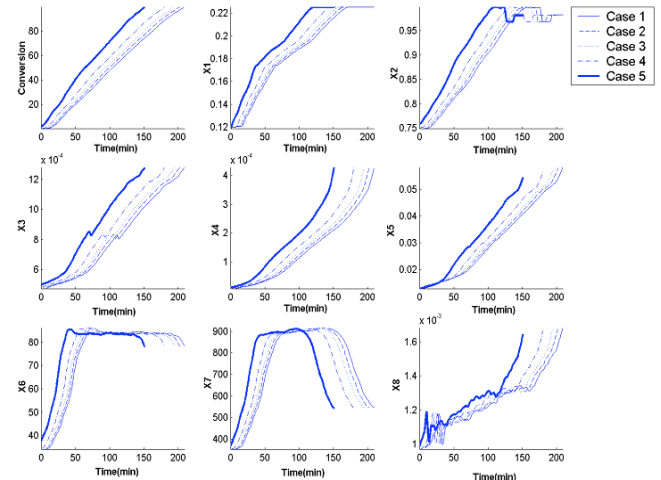


Fig. 5. Product quality

## Discussion

The above examples were meant to illustrate how models built from historical plant data can be used directly to optimize and control processes. Additional published examples include the use of joint PLS models built from data on more than one plant to provide optimal scale-up or transfer of production between manufacturing sites [Garcia-Munoz et. al., 2005, Liu et. al., 2011], and the use of multi-block LV models on product development data for rapidly developing new products (Muteki et. al., 2006).

A limitation of using these models built from historical data is that they are limited to the latent variable space that defines the data on which they were built. Clearly one cannot use them to extrapolate to modes of operation that the plant has never used before. Mechanistic models are needed for that. However, as long as one is content with finding optimal solutions within the space spanned by the plant history, these LV approaches allow one to do this with readily available plant data. In some cases it is necessary to augment this data or to extrapolate in certain directions. In those cases very powerful multivariate DOE methods can be used in the low dimensional LV space to effectively provide powerful designs in the high dimensional process space [Muteki and MacGregor, 2007].

## Supervisory Model Predictive Control of Batch Processes

Although MPC has seen widespread utilization for the control of continuous processes, there has been little application of any advanced control over the final quality of products from batch processes. An obvious reason for this is that the large number of different batch processes and the small volume of specialized products produced in them could not initially justify the application of advanced control. However, many of these batch products are high value added and improving quality and productivity can have a major impact on profitability.

Control of batch processes is also very different in that the quality is only available after the completion of the batch and so the main issue is the effective use of batch recipe and process variable trajectory information to predict the final product quality attributes as the batch progresses and then make mid-course corrections at various decision points. Multivariate PLS models have been shown to be ideal for extracting such information from the time varying process variable trajectories and for predicting the end product quality [Nomikos et, al., 2004, 2005]. Advanced supervisory model predictive control of batch processes has been proposed in a number of papers [Yabuki & MacGregor, 1995, Flores-Cerillo and MacGregor, 2003] and has been commercialized by ProSensus ([www.prosensus.ca](http://www.prosensus.ca)). The PLS models are built from historical batch data and from some limited DOE's at the decision points and QP optimizations are performed to provide mid-course corrections in certain variables and

trajectories. Different objective functions are often used at each decision point. Again a potentially high dimensional control problem is reduced to manageable dimensions by using low dimensional LV models. Recent installations of the commercial Advanced Batch Control (ABC) system, based on latent variable models, has achieved excellent results. Over 80,000 batches have now been controlled with the system while running 24/7 and at over 99% up time. The standard deviations of all quality attributes on the final food product have been reduced by greater than 50% and the productivity of the batches has been increased by nearly 40%. The ABC includes all relevant constraints on the manipulated variables as well as constraints to ensure that all the process variables at the end of the batch are in an acceptable region to enable the immediate start of the next batch.

## Conclusions

This paper provides an overview of recent advancements in the use of data-driven methods for the monitoring, fault diagnosis, fault-tolerance, control and optimization of processes. Discussion has been limited to the use of latent variable models as these are the only data-driven models that allow for uniquely modeling the reduced dimensional spaces in which the process moves, and that provide unique, interpretable and causal models when built from historical process data.

Multivariate methods for the monitoring and diagnosis are reviewed and contrasted with classical fault detection and diagnosis approaches. The integration of monitoring and diagnosis techniques by using an adaptive agent-based framework is outlined and its use for fault-tolerant control is compared with alternative fault-tolerant control frameworks.

The uniqueness of latent variable models built from high dimensional process data and the concept of causality in the latent variable space are used to show that these models can be used to optimize and control the high dimensional processes by optimizing in the low dimensional latent variable space. These concepts are illustrated with several industrial examples: the periodic optimization of an over-injection molding process to counteract raw material and environment changes; finding optimal recipes and operating trajectories for an industrial batch polymerization process that will meet specifications on all final product quality variables in minimum batch time; and supervisory MPC of industrial batch processes for controlling all the final quality attributes.

## Acknowledgments

The work of Ali Cinar and his co-workers on MADCABS, monitoring, fault diagnosis and fault-tolerant control with agent-based systems is supported by the National Science Foundation Grant CTS-0325378 of the ITR program.

## References

- Christofides, P. D., El-Farra N. (2005). Control of Nonlinear and Hybrid Process Systems: Designs for Uncertainty, Constraints, and Time-Delays. Springer-Verlag, New York.
- Du, M., R. Gandhi and P. Mhaskar (2011). An Integrated Fault-Detection and Isolation and Safe-Parking Framework for Networked Process Systems, *Ind Eng Chem Res*, 50, 5667.
- El-Farra, N. H. Ghantasala, S. (2007) Actuator Fault Isolation and Reconfiguration in Transport-Reaction Processes, *AIChE J*, 53, 1518.
- El-Farra N. H., Gani A., Christofides P. D. (2005). Fault-Tolerant Control of Process Systems Using Communication Networks. *AIChE J*, 51, 1665.
- Flores-Cerillo, J. and J. F. MacGregor, (2003). Within-batch and batch-to-batch inferential adaptive control of semi-batch reactors: A Partial Least Squares approach”, *Ind Eng Chem Res*, 42, 3334.
- Frank, P.M., (1990). Disturbance diagnosis in dynamic systems using analytical and knowledge-based redundancy, *Automatica*, 26, 459.
- Gertler, J. J. (1998). *Fault Detection and Diagnosis in Engineering Systems*, Marcel Dekker, New York.
- Garcia-Munoz, S., T.Kourti and J.F. MacGregor, (2005). Product Transfer Between Sites using Joint-Y PLS, *Chemometrics & Intell. Lab. Systems*, 79, 101.
- Garcia-Munoz, S., J.F. MacGregor, D. Neogi, B.E. Latshaw and S. Mehta, 2008. Optimization of batch operating policies. Part II: Incorporating process constraints and industrial applications, *Ind Eng Chem Res*, 47, 4202.
- Huang, B., Shah, S. L., *Performance Assessment of Control Loops: Theory and Applications*, Springer Verlag, 1999
- Jackson, J. E., (1991). *A User's Guide to Principal Components*, New York: John Wiley
- Kourti, T. and MacGregor, J. F., (1996). Multivariate SPC Methods for Process and Product Monitoring, *J Quality Technology*, 28, 409.
- Krzanowski, W.J., (1979). Between-Groups Comparison of Principal Components, *J. Amer. Stat. Assn.*, 74, 703.
- Kendra S. J., Cinar, A. (1997). Controller Performance Assessment by Frequency Domain Techniques. *J Proc Contr*, 7, 181.
- Kendra, S. J., Basila, M. R., Cinar A. (1997). Intelligent Process Control with Supervisory Knowledge-Based Systems, in *Methods and Applications of Intelligent Control*. Kluwer Academic Publishers, The Netherlands.
- Liu Z., M.J. Bruwer, J.J. MacGregor, S. Rathore, D.E. Reed, M.J, Champagne. “Scale-up of a roller compaction process using JY-PLS”, *Ind. & Eng. Chem. Res.*, 50, 10696-10706, 2011.
- Maciejowski, J. M. (1999). Modelling and Predictive Control: Enabling Technologies for Reconfiguration. *Annual Reviews in Control*, 23, 13.
- Mhaskar P., Gani A., McFall C., Christofides P. D., Davis J. F. (2007). Fault-Tolerant Control of Nonlinear Process Systems Subject to Sensor Faults. *AIChE J.*, 53, 654.
- Mhaskar, P., Gani, A., El-Farra, N. H., McFall, C., Christofides, P.D., Davis, J. F. (2006). Integrated Fault-Detection and Fault-Tolerant Control of Process Systems. *AIChE J.* 52, 2129.
- Muteki, K., J.F. MacGregor and T. Ueda, (2006). On the Rapid development of New Polymer Blends: The optimal selection of materials and blend ratios, *Ind Eng Chem Res*, 45, 4653.
- Muteki, K., J.F. MacGregor, and T. Ueda, (2007). Mixture designs and models for the simultaneous selection of ingredients and their ratios, *Chemometrics & Intell. Lab. Systems*, 86, 17.
- Negiz, A. and Cinar, A. (1997). Statistical Monitoring of Multivariable Dynamic Processes with State Space Models, *AIChE J.*, 43, 2002.
- Nomikos P. and J.F. MacGregor, (1994). "Monitoring of Batch Processes Using Multi-Way Principal Components Analysis, *AIChE J.*, 40, 1361.
- Nomikos, P. and J.F. MacGregor, (1995). Multi-Way Partial Least Squares in Monitoring Batch Processes, *Chemometrics & Intell. Lab. Systems*, 30, 97.
- Patton, R. J., Frank, P. M. Clark, R. N. (Eds) (1989). *Fault Diagnosis in Dynamic Systems – Theory and Applications*, Prentice Hall, Englewood Cliffs, NJ.
- Perk, S., Teymour, F., Cinar, A. (2010). Statistical Monitoring of Complex Chemical Processes with Agent-Based Systems. *Ind Eng Chem Res*, 49, 5080.
- Perk, S., Teymour, F. and Cinar A. (2011). An Adaptive Agent-Based System for Process Fault Classification and Diagnosis, *Ind Eng Chem Res*, accepted for publication.
- Perk, S., Shao, Q. M., Teymour, F., Cinar, A. (2011). An Adaptive Fault-Tolerant Control Framework with Agent-Based Systems. *Int. J. Robust. Nonlinear Control*, accepted for publication.
- Raich, A. and Cinar, A. (1995). Multivariate Statistical Methods for Monitoring Continuous Processes: Assessment of Discrimination Power of Disturbance Models and Diagnosis of Multiple Disturbances, *Chemometrics and Intelligent Laboratory Systems*, 30, 37.
- Raich, A. and Cinar, A. (1996). Statistical Process Monitoring and Disturbance Diagnosis in Multivariable Continuous Processes, *AIChE J.*, 42, 995.
- Raich, A. and Cinar, A. (1997). Diagnosis of Process Disturbances by Statistical Distance and Angle Measures, *Compu. Chem. Engng*, 21, 661.
- Rawlings J. B., Stewart B. T. (2008). Coordinating Multiple Optimization-Based Controllers: New Opportunities and Challenges. *J Process Control*, 18,1665.
- Schaefer, J., Cinar, A. (2004). Multivariable MPC system performance assessment, monitoring and diagnosis. *J Process Control*, 14, 113.
- Tatara, E., Cinar, A. (2002). An Intelligent System for Multivariable Statistical Process Monitoring and Diagnosis. *ISA Trans.*, 41, 255.
- Undey, C., Ertunc, S., and Cinar, A. (2003). Online Batch/Fed-batch Process Performance Monitoring, Quality

- Prediction, and Variable-Contribution Analysis for Diagnosis, *Industrial and Engineering Chemistry Research*, 42, 4645.
- Wang, Y., Yang, Y., Zhou, D., Gao F. (2007). Active Fault Tolerant Control of Nonlinear Batch Processes with Sensor Faults *Ind Eng Chem Res.*, 46, 9158.
- Venkatasubramanian, V., Rengaswamy, R., Kavuri, S., Yin, K. A. (2003). Review of Process Fault Detection and Diagnosis: Part III: Process History Based Methods. *Comp Chem Eng*, 27, 327.
- Y. Yabuki, T. Nagasawa and J.F. MacGregor, "An Industrial Experience with Product Quality Control in Semi-Batch Processes", *Computers & Chem. Eng.*, 24, 585-590, 2000.
- Yacoub, F. and J.F. MacGregor, (2004). Product optimization and control in the latent variable space of nonlinear PLS models, *Chemometrics & Intell. Lab. Syst.*, 70, 63.
- Yoon, S. and MacGregor, J.F. (2000). Relationships between Statistical and Causal Model-Based Approaches to Fault Detection and Isolation. *AIChE J.*, 46, 1813.
- Zumoffen, D., Basualdo, M. (2008). From Large Chemical Plant Data to Fault Diagnosis Integrated to Decentralized Fault Tolerant Control: Pulp Mill Process Application. *Ind Eng Chem Res.*, 47, 1201.