
Using Bayesian Networks to Predict Change Sets

Sarah Nadi

SNADI@UWATERLOO.CA

University of Waterloo,
Waterloo, ON, Canada

Abstract

STILL

1. Introduction

Change is inevitable in any Information Technology(IT) system. New features are added, different configurations are used, upgrades are introduced and new software and hardware are added. In a large system, such changes can cause other parts of the system to malfunction without the ability of the analyst to foresee this side effect. Accordingly, we need to ensure that we consider the impact of changes before they are implemented, and make sure we adjust the potentially affected components accordingly while implementing the change. This leads to the process of change set detection which is finding all the components of the system that need to be included in your change set. The change set is the set of components that need to be changed to avoid any side effects.

Many of the organizations with large and sophisticated IT systems employ a Configuration Management Database (CMDB) to help them manage these systems. The CMDB keeps track of all the components in a system, how they are related, as well as all the changes that have occurred to them. Any component in the CMDB is called a Configuration Item (CI). In previous work (Nadi et al., 2010), we mined the CMDB repository for historical co-changes, and used the mined correlations to predict change sets. That is, given an initial CI that is going to change, we predict what other CIs might need to be changed as well. We obtained really promising results in terms of recall and precision (69.8% and 88.5% respectively).

In this paper, we wish to examine the same problem, but from a Bayesian perspective. Instead of looking

at pairwise historical co-changes of CI, we would like to consider all the changed CIs in the past to deduce the change relations between them. That is, given a set of observations about how different CIs change together, can we deduce the relationships between them that would allow us to predict future change sets. Accordingly, we explore the different ways a Bayesian network can be constructed from the data we have, and then test the predictions produced by querying this network to examine the obtained recall and precision. We hope that we can obtain better results using Bayesian networks since they are not limited to pairwise comparisons.

The rest of this paper is organized as follows. BLA BLA BLA

2. Background

2.1. CMDBs and Change Sets

The repository of information we are using is the Configuration Management Database (CMDB) ([Office of Government Commerce](#) , [OGC](#)). The CMDB is useful in Enterprise IT Management (EITM) since it provides information about the various critical components in a system including hardware, software, and services provided by the company. It records the configuration of these items, their change history, their incident history, as well as the relationships between them. Each item stored in the CMDB is referred to as a Configuration Item (CI). Figure 1 shows an example CMDB to illustrate the concepts of CIs and relationships.

A CMDB provides a basis for decision making processes such as Incident Management, Change Management, etc. In this paper, we focus on the process of Change Management, and in particular, on the problem of change set detection. A *change* is the addition, modification, or removal of anything that could affect on IT services. A poorly planned change may lead to a fault in the system. Accordingly, when one wants to change a CI is the system, other CIs that might need

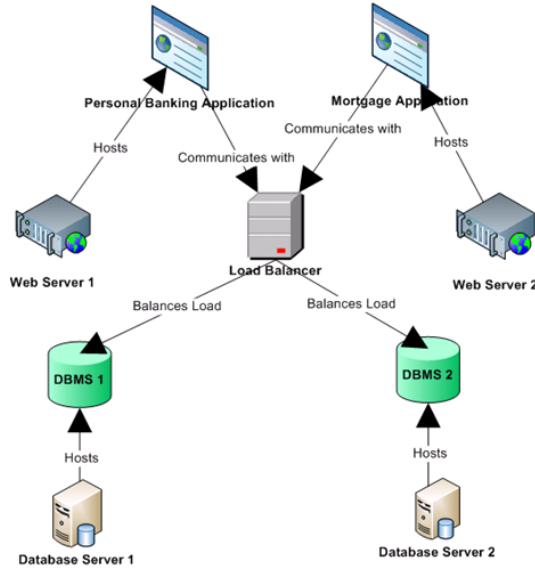


Figure 1. An IT System Stored in a CMDB. Each item shown is a CI, and CIs are tied through different relationship types

to be changed as well must be correctly identified. The set of CIs that will need to be modified for the change to be complete is called a change set.

2.2. Bayesian Networks

Bayesian Networks, first introduced by Pearl (Pearl, 1988) are probabilistic graphical models. Each node in a graph represents a random variable, and nodes are connected through edges which represent probabilistic dependencies (Ben-Gal). A Bayesian network is a Directed Acyclic Graph (DAG), and is essentially defined by two aspects: the structure of the graph and the Conditional Probability Tables (CPTs) probabilistically relating connected nodes. In this work, a node in the Bayesian network would be a CI which can take values “true”, “false” to denote the probability of it changing or not. The CPTs would essentially give the probability of one CI changing given its parents.

Given a set of data consisting of observations for our variable space, machine learning techniques can be used to estimate the network. If the structure of the network is known, and we have full observability (i.e we can observe values for each variable), then learning the CPTs can be achieved through Maximum-likelihood estimation (MLE) (Ben-Gal). MLE tries to maximize the log likelihood of the the training data which can simply be done by calculating the frequency of occurrence of the values of each variable given its parents. If the structure is known, but we only have partial observability (i.e missing data), then we can use a vari-

ation of MLE which is the Expectation Maximization algorithm (Dempster et al., 1977). If the structure is not known in the first place, then searching through the model space, and calculating network scores for each produced network is necessary. The simplest algorithm for this is a Greedy search algorithm (Friedman & Goldszmidt, 1996).

While searching for the best network that fits the training data, scores are needed for the produced networks in order to know which one is the best one. There are several types of scores that have been proposed in the literature. Two of the commonly used scores are the Minimum Description Length (MDL) score (Rissanen, 1986) and the BDe score (Heckerman et al., 1995). The MDL score is based on finding the length of data required to store a network. It tries to balance the complexity of the induced network with the degree of accuracy of which the network represents the frequencies in the training data. Therefore, the lower the MDL score of a network, the better it is. The BDe score is proportional to the posterior probability of each network structure being scored, given the data. Accordingly, we are searching for the network that maximizes this probability. For large data sets, both metrics score structures similarly (Friedman & Goldszmidt, 1996). Since they are similar, we simply chose to use the BDe score in this work to compare the produced network structures.

2.3. Bayesian Network Tools

2.3.1. WEKA

WEKA (Hall et al., 2009) is a toolbox for machine learning. It provides different algorithms for classification and clustering as well as other machine learning problems. It also provides Bayesian learning techniques for classification problems. Additionally, it provides explorations tools for Bayesian Networks, and allows the user to learn the structure and CPTs of a Bayesian Network. We will mainly be using WEKA for learning the CPTs in a fixed structure network using the SimpleEstimator algorithm (Witten & Frank, 2005) implemented there. This simply estimates the probability values on edges based on the frequency values of each variable given its parents in the training set.

2.3.2. BANJO

Banjo (Banjo) is a tool written in Java which infers the structure of a Bayesian network given training data. It has two different search algorithms: Greedy and Simulated Annealing. For either of these search algorithms, it has two methods of proposing a new edge. Banjo

uses the BDe metric to compute a network’s score. We used the Greedy search algorithm in our experiments.

2.3.3. JAVABAYES

JavaBayes is a Java tool for calculating marginal probabilities and expectations in a Bayesian network. In a network, nodes can be set to be observed, and then the posterior probability of the remaining nodes can be calculated using Bayes theorem. This is the functionality we need here since the CI that will be changed will be considered as observed to be “true”, and then we need to calculate the posterior of the remaining CIs given this observation.

3. Constructing the Bayesian Network

Before constructing the Bayesian Network, we first had to process and prepare the data we have. Then, in order to construct a Bayesian network to use for predictions there are two steps involved. First, determining the structure of the actual network, and then estimating the Conditional Probability Tables (CPT). Section 3.1 first explains the data set available to us, and the data preprocessing involved before building the network. Section 3.2 then explains the different techniques we experimented with to build the network structure. Section 3.3 shows the last step to build the network which is estimating the CPTs. Figure 2 shows this overall process.

3.1. Data Preprocessing

DATA SIZE

The original data set available to us in the sample CMDB we have is from three years, and has 7,999 distinct CIs, and 27,305 change orders. This amount of data was infeasible to work with as no tool could handle such a large amount of variables in a Bayesian network. For the purposes of this project, which is mainly to experiment with Bayesian techniques for change set prediction, it is sufficient to choose a small representative subset of the data. Accordingly, in order to be able to test things properly, we used observations from three months data from January 1, 2008 to March 31, 2008 to build the model. However, even in such a short period, there were already 2,841 distinct CIs appearing and 2,229 observations (i.e. change orders). Therefore, we needed to perform further data preprocessing.

First, we removed all CIs that have changed less than 12 times within this time frame (i.e. were associated with less than 12 different change orders in our training observations). There is no particular reason for choosing 12 as a cut off. It simply gave a feasible data

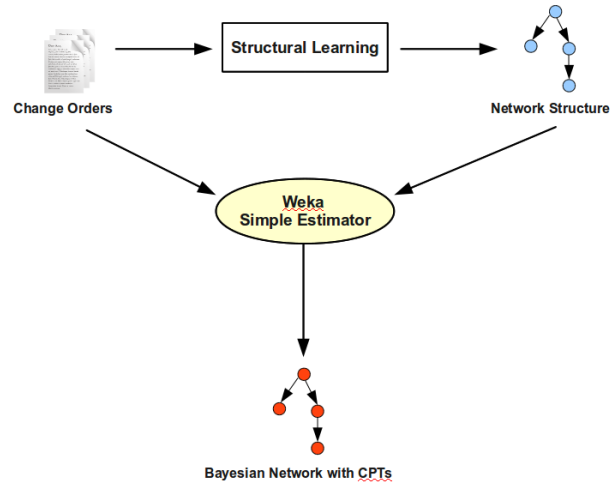


Figure 2. Building the Bayesian Network

set to deal with. This yielded 120 distinct CIs. Then, in order to slightly increase our variable space to include other related CIs, we found all the parent CIs related to these 120 CIs from the CMDB perspective. However, we ignored three common, and not extremely meaningful relations, which are “supports”, “is location for”, and “backs up”. Adding the related parent CIs, we now had a set of 241 CIs. However, some of the added parent CIs may not be in the original set of 2,841 CIs. Accordingly, we just kept CIs that appeared more than 12 times or were in the set of related CIs. This provided us our final data set of 170 CIs which we use throughout our models for fair comparison. Additionally, filtering out CIs meant filtering out some of the change orders that did not have any CIs satisfying our criteria. This led to us having 1,305 change orders instead of 2,229 which was a more manageable set. Accordingly, our training data set consisted of 170 variables (CIs), and 1,305 observations (change sets).

DATA FORMAT

In the CMDB, a Change Order has several fields including the requester, the assignee, the start date of the change, the description and the change set field (called the ‘Configuration items’ fields). Of the fields in a change order, we use only the change set field. The advantage of only using this one field is that change sets are easy to extract and easy to understand. In terms of an observation, CIs were either marked as “true” or “false” indicating whether they have changed or not. Each observation initially consisted of all the CIs in our data set marked as “false”. For each change order in the training set, we would mark each CI appearing in the “Configuration Items” field as true, while all those not appearing were left as

“false”. Therefore, the set of training data consisted of an observation entry for each change order with the CIs marked as “true” or “false” accordingly.

3.2. Network Structure

There were different ways in which we could estimate the structure of the network, and so we built a Bayesian Network using each technique.

3.2.1. METHOD 1 (BANJO)

For the first method, we used Banjo to learn the structure of the Bayesian network from the training data. We use the Greedy searching algorithm implemented there. Banjo produced five different top-scoring networks (having the same score), and we simply chose one of them as the network to be used.

3.2.2. METHOD 2 (WEKA CLASSIFICATION)

Although our research problem as currently defined is not really a classification problem, we want to examine how the network built by the BayesNet classifier in Weka will perform. Accordingly, we estimated the structure of the network using the K2 structure learning algorithm implemented in WEKA which is a greedy search based learning algorithm. The default maximum number of parents allowed for any node is one (a Naive Bayesian Network). However, we used a maximum of five to allow the nodes to be related since there is really no class label in our data set.

3.2.3. METHOD 3 (CMDB RELATIONSHIPS)

For the third method, we used the relations existing in the CMDB to infer the structure of the network. We used two variations for this. The first was to follow the direction of the relationship edges in the CMDB. That is if there is an edge from A to B, we will place an edge in the Bayesian network from A to B. The second was to reverse the direction of the relationship edges in the CMDB. That is, if there is an edge from A to B, we will place an edge in the Bayesian network from B to A. There was one problem, however, with the generated networks. Two of the CIs had more than 15 parents to them which means that their CPTs will be intractable to compute. For those two CIs, we simply removed all their parents to make the computation tractable.

3.3. Estimating the CPTs

After the structure was determined, we used the SimpleEstimator algorithm (Witten & Frank, 2005) built in Weka to calculate the CPTs for all networks. The SimpleEstimator uses equation BLA GET EQUA-

TION After setting the data set to be the observations in our training set, Weka learned the CPTs for each of the networks above. This complete network was then saved as an BIF XML file to be used by JavaBayes for general inference as will be explained in the next section.

Additionally, we experimented with the missing data feature in WEKA. Weka allows the user to specify missing values for any variable by simply putting a ‘?’ instead of its value. Therefore, instead of putting the CIs that did not appear in a change order as false, we put a ‘?’ in their place. This seemed closer to practice, because in reality, we are not sure whether this CI actually changed or not. We tried learning the CPTs for the network produced by WEKA classification, and for that produced using the CMDB relationships. The missing values were only used to learn the CPTs except in the case for WEKA where the structure was also learned using the missing values.

4. Experiment Setup

This section explains how we tested the predictions of each of the networks in a way that simulates how an analyst would use the tool in reality. We first explain how we simulate this behavior, and then explain how we evaluate the predicted CIs. Algorithm 1 summarizes this setup.

4.1. Change Set Detection Process Simulation

The change set detection process provided is an iterative, collaborative process between the tool and the analyst. When the tool suggests CIs, the analyst can accept or reject these CIs (i.e add them to the change set or not). Based on the CIs added to the change set, the analyst can ask the tool for more suggestions. This is along the lines of “Now, that I am also going to change these CIs, what else do I need to change?”.

To simulate this process for our experiments, we do the following. At this point, we have the Bayesian network ready, and we would like to perform Inference. More formally, given that a CI will change (our observation), we want to infer the probability that the other CIs in the network might change as well. Unfortunately, neither of the two tools previously used provide a general inference engine. We, therefore, use JavaBayes in this step since it accepts the same BIFF format used by Weka. We use a one month testing set where we try to predict all the change orders in April 2008. There was a total of 883 change orders in that month which we use in our testing set.

For each change order, we would take the first CI as

Algorithm 1 Generating Predictions using the Bayesian Network

Input: *changeorder*
Input: *Bayesiannetwork*
Input: *threshold*
Set *initialCI* = first CI in *changeOrder*
Set *occurredset* = *changeorder* - *initialCI*
Initialize *observedset* = *initialCI*
Initialize *predictedset* =
repeat
 Initialize *newpredictions* =
 for *node* in *BayesianNetwork* **do**
 posterior = query *node* in *Bayesiannetwork*
 using Banjo given observed set
 if *posterior* > *threshold* and *node* not in
 predictedset **then**
 Add *node* to *newpredictions*
 end if
 end for
 for *prediction* in *newpredictions* **do**
 if *prediction* in *occurredset* **then**
 Add *prediction* to *observedSet*
 end if
 end for
 Add *newpredictions* to *predictedset*
until *newpredictions* is empty
recall = (*predictedset* ∩ *occurredset*) / *occurredset*
precision = (*predictedset* ∩
occurredset) / *predictedset*

the initial CI to change, then we would set that as an observed node, and update the beliefs of all the nodes (CIs) in the network using JavaBayes. We are basically checking the probability of each CI being “true”. We used different cutoff thresholds for the posterior probabilities we would consider. We would then loop on all the updated CIs, and add those that match our threshold criteria to the predicted change set. For example a threshold of 0.2 means that only nodes that have a posterior probability greater than 0.2 of having the value “true” will be returned. To simulate a real life scenario, we then checked which of these predicted CIs actually lies in the target change set we are trying to predict. This is similar to an analyst accepting CIs into their change set. These common CIs would then also be marked as observed so that we can predict what else will need to change given that these CIs are also changing. Again, we would calculate the posterior probability, and continue doing so until there are no more common CIs. All the CIs that match the threshold criteria (whether accepted by the analyst or not) are part of the predicted set.

4.2. Evaluation Techniques

We need a way to evaluate the predicted CIs. The recall and precision measures from the information retrieval field are appropriate for this type of evaluation. Recall measures the proportion of correct CIs retrieved by the system, while precision measures the proportion of suggested CIs that are correct (van Rijsbergen, 1979).

Similar to Hassan et. al (Hassan & Holt, 2004), we score define the *Predicted Set* (P) as the set of all predicted. We define the *Occurred Set* (O) as the CIs remaining in the change set after excluding the Initial CI provided by the analyst (i.e Change Set - Initial CI). The intersection of the predicted set and the occurred set, called *PO*, is the common CIs in both sets. For each constructed change set, we then calculate the recall and precision values for the predictions according to the following definitions (Hassan & Holt, 2004):

$$Recall = \frac{|PO|}{|O|} \quad (1)$$

$$Precision = \frac{|PO|}{|P|} \quad (2)$$

If no CIs are predicted (i.e., *P* and thus *PO* are empty), precision is defined as 1 since there cannot exist any incorrect predictions in an empty set. On the other hand, if the size of the change set is 1, and thus the size of the occurred set is 0, recall is defined as 1 since there are no CIs to predict (Hassan & Holt, 2004).

In order to have a single measure that indicates the effectiveness of our predictions, we use the F-measure which is based on van Rijsbergen’s effectiveness measure which combines recall and precision (van Rijsbergen, 1979). The F-measure is calculated according to Equation 3 which gives equal weighting to recall and precision. The ideal F-measure is 1 where both recall and precision are 1.

$$F = 2 * \frac{precision * recall}{precision + recall} \quad (3)$$

5. Results

For all the networks generated from different methods explained in Section 3.2, and following the procedure described in Section 4, we calculated the average recall and precision for all the change orders in the test set. Based on the average recall and precision values, we calculated the F-measure. For better visualization, we

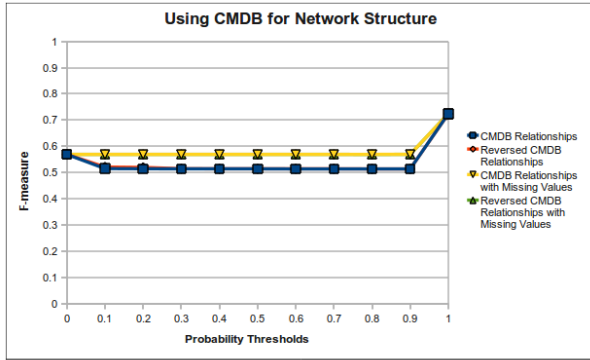


Figure 3. F-measure of the different Bayesian networks produced using the CMDB relationships

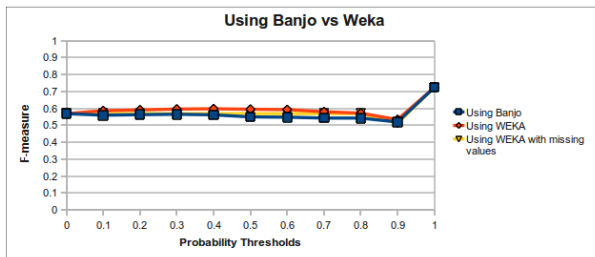


Figure 4. F-measure from network structures estimated by WEKA versus Banjo

plot the results from all the networks based on CMDB structure together, and the results from those based on Banjo and WEKA together.

Figure 3 shows the F-measure obtained from the four different networks built based on CMDB relationships for structural information. The first observation is that reversing the relationship edges did not produce any difference in the results. The curves for the CMDB relationships and the reversed CMDB relationships are completely overlapping. The same things applies to the curves of CMDB relationships with missing data, and reversed CMDB relationships with missing data. The second observation here is that denoting CIs that did not appear in the change order as missing data rather than giving them the value “false” produced slightly better results. However, the main conclusion from these results is that given each change order, all of the predicted CIs in the network already had a posterior probability of greater than 0.9. That is why the F-measure is the same for all thresholds (excluding 0), and only changes when the threshold becomes 0.9. This seems rather strange, and despite the relatively high F-measure, it does not seem that the structure of these networks is reflective of their relationships in terms of change.

Figure 4 compares the F-measure obtained from the graphs whose structure was built using Banjo versus those built using WEKA.

6. Related Work

Mirarab et al. (Mirarab et al., 2007) investigate the same problem as our work. However, their work is on the level of source code changes. They build three different Bayesian Networks, one that is based on package and class dependency information (static relationships), one which is dependent on historical co-changes, and one which uses both. For the first graph, the initial structure is essentially “given” according to the static dependencies, and then the CPTs are learnt using the importance sampling algorithm proposed by Changhe and Marek (Yuan & Druzdel, 2003). The way static dependencies are defined in their case is specific to Java. The third one is essentially the first graph, but updated using the historic change information according to the Expectation Maximization (EM) algorithm (Dempster et al., 1977). The second was solely based on historic information where the network is build using a greedy structure learning algorithm (Friedman & Goldszmidt, 1996). They did some preprocessing to their data such as filtering out large changes (with more than 30 elements changed at once) since this was probably an insignificant change.

Zhou et al. (Zhou et al., 2008) try to answer a slightly different problem. They do not only look at the probability of other elements changing given a specific element, they also add features such as authors, change significance levels etc. and try to predict if two elements are co-changes or not accordingly. Thus, their problem is more of a classification problem where given two elements, and some observed features they try to determine the class as a co-change or not. They use the K2 algorithm proposed by Cooper et. al (Cooper & Herskovits, 1992) to estimate the structure of the Bayesian network, and use the SimpleEstimator algorithm built in WEKA (Witten & Frank, 2005).

7. Conclusion

References

- Banjo. Banjo. "<http://www.cs.duke.edu/~amink/software/banjo/>".
- Ben-Gal, I. Bayesian networks.
- Cooper, G.F. and Herskovits, E. A Bayesian method for the induction of probabilistic networks from data. *Machine learning*, 9(4):309–347, 1992. ISSN 0885-6125.

- Dempster, A.P., Laird, N.M., Rubin, D.B., et al. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977. ISSN 0035-9246.
- Friedman, N. and Goldszmidt, M. Learning Bayesian Networks with Local Structure. 1996.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., and Witten, I.H. The WEKA data mining software: An update. *ACM SIGKDD Explorations Newsletter*, 11(1):10–18, 2009. ISSN 1931-0145.
- Hassan, Ahmed E. and Holt, Richard C. Predicting change propagation in software systems. In *ICSM '04: Proceedings of the 20th IEEE International Conference on Software Maintenance*, pp. 284–293, Washington, DC, USA, 2004. IEEE Computer Society.
- Heckerman, D., Geiger, D., and Chickering, D.M. Learning Bayesian networks: The combination of knowledge and statistical data. *Machine learning*, 20(3):197–243, 1995. ISSN 0885-6125.
- Mirarab, S., Hassouna, A., and Tahvildari, L. Using bayesian belief networks to predict change propagation in software systems. pp. 177 –188, jun. 2007.
- Nadi, S., Holt, R., and Mankovskii, S. Does the past say it all? Using history to predict change sets in a CMDB. In *CSMR'10: Proceedings of the 14th European Conference on Software Maintenance and Reengineering*, Madrid, Spain, 2010.
- Office of Government Commerce (OGC), ed.: Service Support. Office of government commerce (ogc), ed.: Service support. *IT Infrastructure Library (ITIL)*, 2000.
- Pearl, J. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 1988. ISBN 1558604790.
- Rissanen, J. Stochastic complexity and statistical inference. *Analysis and Optimization of Systems*, pp. 391–407, 1986.
- van Rijsbergen, C.J. Information retrieval, 1979.
- Witten, I.H. and Frank, E. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann Pub, 2005. ISBN 0120884070.
- Yuan, C. and Druzdzal, M.J. An importance sampling algorithm based on evidence pre-propagation. In *Proceedings of the 19th Annual Conference on Uncertainty on Artificial Intelligence*, pp. 624–631. Citeseer, 2003.
- Zhou, Yu, Würsch, Michael, Giger, Emanuel, Gall, Harald C., and L?, Jian. A bayesian network based approach for change coupling prediction. *Reverse Engineering, Working Conference on*, 0:27–36, 2008. ISSN 1095-1350. doi: <http://doi.ieeecomputersociety.org/10.1109/WCRE.2008.39>.