

Evaluation Issues With Critical Success Metrics

Tim Menzies
Artificial Intelligence Department
School of Computer Science and Engineering
The University of NSW

`timm@cse.unsw.edu.au`

`http://www.cse.unsw.edu.au/~timm`

October 29, 1997

Abstract

If we lack an objective human expert oracle which can assess a system, and if we lack a library of known or desired behaviour, how can we assess an expert system? One method for doing so is a *critical success metrics* (CSMs). A CSM is an assessment of a running program which reflects the business concerns that prompted the creation of that program. Given *pre-disaster* knowledge, a CSM can be used while the expert system is in routine use, without compromising the operation of the system. A general CSM experiment is defined using pre-disaster points which can compare (e.g.) human to expert system performance. Examples of using CSMs are given from the domains of farm management and process control.

1 Introduction

How are we to assess the knowledge engineering techniques being reported in the knowledge acquisition (KA) literature? We should carefully assess superlative claims for the efficacy of case tools or formal methods or object-oriented knowledge representations or problem solving methods (PSM) [Schreiber *et al.*, 1994] or ontologies [Gruber, 1993] or ripple down rules [Preston *et al.*, 1993] or abduction [Menzies, 1996] or the problem space computational model [Yost, 1993] or whatever. In the software engineering literature, there are many examples of software engineering techniques (e.g. CASE tools, formal methods) which are in common use but, when evaluated, cannot be shown to be beneficial to the software process [Fenton *et al.*, 1994]. Also, in the KA literature, many of the claims in the PSM literature are not supported by the currently available empirical evidence [Menzies, 1997a].

Clearly, we need some better method than reading the glowing reports from the authors of these KA techniques. Even if these authors are expert in their fields, they may still be unable to perform objective expert evaluations. Experts can often disagree about what constitutes a competent system ([Shaw, 1988, Gaschnig *et al.*, 1983]). The *halo effect* prevents a developer for looking at a program and assessing its value. Cohen likens the halo effect to a parent gushing over the achievements of their children and comments that...

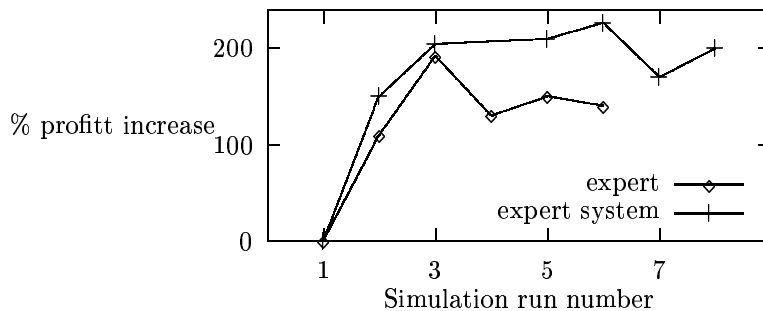


Figure 1: Critical success metrics for PIGE. From [Menzies *et al.*, 1992].

What we need is not opinions or impressions, but relatively objective measures of performance. [Cohen, 1995], p74.

The method for assessment explored in this article is *critical success metrics* (CSMs); i.e. some number inferred from the system which, if it passes some value, demonstrates conclusively that the system is a success. If such a critical measurement is observed, then the system will be deemed to be a success, regardless of other less critical measures (e.g. slow runtimes).

For example, consider the PIGE farm management expert system [Menzies *et al.*, 1992]. PIGE advised on diets and genotypes for pigs growing in a piggery. Given a particular configuration of the livestock, an optimisation model could infer the annual profit of the farm. Alternate configurations could be explored using a simulation model. A user can choose some settings, then run the simulation model to see if the system's performance improved. The CSM for PIGE was *can the system improve farm profitability as well as a pig nutrition expert?*. If this could be demonstrated, then the tool could be sold as an automatic pig growth specialist. To collect this CSM, at the end of a three month prototyping stage, we compared the performance of the pig nutritionist who wrote the PIGE rules against PIGE. We observed that, measured in purely economic terms, this expert system out-performed its human author (!!). The CSM study results for PIGE are shown in Figure 1.

This single CSM study changed the direction of the project. The graph of the CSM study became a succinct argument for collecting further funding. It was also very useful in sales work. PIGE became Australia's first exported expert system and was used on a routine daily basis in America, Holland, Belgium, France, Spain and Australia. In part, the success of the system was due to its ability to demonstrate its utility via a CSM.

Nevertheless, the CSM study of PIGE is a poor evaluation study. A good experiment is run multiple times with some variation between each trial [Cohen, 1995]. CSMs should be viewed as the inner measurement process within a well-defined experiment. A general class of such experiments are described below, along with an example in a process control domain. This example will use a technique called a *pre-disaster point* (defined below). Our example will be preceded by general notes on CSMs and their advantages.

2 About CSMs

This section offers some basic notes on CSMs. CSMs are a reflection of the contribution of the behaviour of the software in a particular business context. Hence:

- They are very domain-specific. However this does not mean they are unanalysable. This article concerns itself with the general themes of CSMs.
- They typically do not refer to internal properties of a program. In this regard, CSMs are very different to the syntactic anomaly detection systems of the KBS verification community [Preece, 1992].
- They cannot be developed by programmers without extensive input from business users. Programmers developing CSMs without business user involvement typically focus on internal properties; e.g. lines of code per function, bugs fixed per day, etc. Such internal properties may not connect to the business case which motivated the program's development.
- They can only be collected once the program is running in its target context.

Even if can't collect CSMs until an expert system is deployed, we should still define them at a very early stage. Evaluation should be considered as early as possible when building a system [Gaschnig *et al.*, 1983]. The incremental application of a pre-defined success criteria can be a powerful tool for managing evolving systems [Booch, 1996]. Often, the evaluation criteria imposes extra requirements on the implementation. We may need to build a very simple initial system that collects baseline measurements which reflect current practice. For example, once I identified *increases sales per day* as the CSM for a dealing room expert system. However, this number was not currently being collected in the current software. Sales per day could be estimated from the quarterly statements, but no finer grain data collection was performed at that site. Hence, prior to building the expert system, a database system had to be built to collect the baseline data.

While CSMs are obvious in retrospect, they can take weeks of analysis to uncover. For example:

- It took two weeks full time analysis on the domain before the above dealing room CSM was uncovered.
- In the process control system discussed below, the CSMs were only isolated once a prototype expert system system was developed.
- In the PIGE system, nutrition experts argued for weeks about the merits of different protein utilisation models. Then the marketing people commented that such considerations were irrelevant if it could not be demonstrated that the systems recommendations improved the overall profitability of a farm. Hence, the evaluation focus moved from the protein utilisation models to issues of modeling the farm economics. The results, shown above, were an impressive demonstration of the marketability of the system.

The observation that CSMs can take some time to isolate would not surprise software engineering metrics researchers. Basili [Basili, 1992], characterises software evaluation as a *goal-question-metric* triad. Beginners to experimentation report whatever numbers they can collect without considering the goal of the research project, what questions

relate to that goal, and what measurements could be made to address those questions. Before goal-question-metric there must be an analysis involving the stakeholders of the project to establish the appropriate goals. Offen and Jeffery [Offen & Jeffery, 1997] offer the appropriate caution that this important task can take a non-trivial amount of time.

3 Advantages of CSMs

Expert systems are usually evaluated via panels of experts or some database of known or desired behaviour. Such evaluations can report the accuracy of those systems to an enviable degree of accuracy. For example:

- [Hayes, 1997] can demonstrate that her expert system developed in two years performs as well as someone with five years experience in that field.
- [Preston *et al.*, 1993] reports a biochemical interpretation system that is 95 percent accurate on the cases it analyses.
- [Yu *et al.*, 1979] reports that MYCIN, an expert system for prescribing antibiotics, clearly out-performs senior medical personnel.

Using CSMs, we are placing a business-level success criteria on a running system. Hence, we can evaluate a system even when:

- No objective source of expertise is available; i.e. expert panel members are unavailable.
- There exists no representative library of the known/desired behaviour of the system; i.e. we have yet to have enough experience with the domain to record all the possible things which can happen.

Also, the evaluation will be a *business-level* evaluation. Business users may demand objective evidence as to the business value of some program before allowing it to control some critical business process. This evaluation may not comprise developer-level concerns such as runtimes or (in the case of PIGE) current fashions in theories of protein utilisation. In the PIGE and dealing room examples, the CSMs had to reflect the fundamental business case which motivated the project: increased profitability.

Further, given a *pre-disaster point*, we can do this while the system is in routine operation. A pre-disaster point refers to a state of the system that is less-than-optimum, but not yet critically under-performing. As we shall see below, CSMs plus pre-disaster knowledge allows us to assess a system without compromising its operation.

4 CSM Evaluation

This section offers a general design for an evaluation experiment using CSMs and a pre-disaster point. The aim of this evaluation is to check if the program is dumber than some human, with respect to some chosen CSMs. In the experiment, the human or expert system is trying to control some aspects of the environment (e.g. make a diagnosis, prescribe medicines which reduce fever, improve profitability, etc).

Trials would alternate between the human and computer experts. A trial would begin when the system is in some steady state; i.e. there appears to be no currently

active problems. During the course of each trial, the expert under trial would have sole authority to order adjustments to the environment. The trial would terminate whenever the pre-disaster point was reached. Authority to adjust the environment would then pass to the human experts. At the conclusion of each trial, a CSM is applied to assess the environment during the trial period.

At the end of a statistically significant number of trials (say, 20 for each population of experts), the mean performance of the two populations of experts would be compared using a t-test as follows. Let m and n be the number of trials of expert system and the human experts respectively. Each trial generates a performance score: $X_1 \dots X_m$ with mean μ_x for the humans; and performance scores $Y_1 \dots Y_n$ with mean μ_y for the expert system. We need to find a Z value as follows:

$$S_x^2 = \frac{\sum(x_i - \mu_x)^2}{m - 1}$$

$$S_y^2 = \frac{\sum(y_i - \mu_y)^2}{n - 1}$$

$$Z = \frac{\mu_x - \mu_y}{\sqrt{\frac{S_x^2}{m} + \frac{S_y^2}{n}}}$$

Let a be the degrees of freedom. If $n = m = 20$, the $a = n + m - 2 = 38$. We reject the hypothesis that expert system is worse than the human (i.e. $\mu_x < \mu_y$) with 95% confidence if Z is less than ($-t_{38,0.95} = -1.645$).

Note that this human/expert system comparison could also be used to assess different expert systems.

5 An Example: Process Control CSMs

This section offers a detailed example of the above experiment. In the summer of 1986/87, I implemented QUENCH, an expert system computer program for the control of the quench oil tower at ICI Australia's Olefines petrochemical plant in Sydney [Menzie & Markey, 1987]. Once the system was built, I offered to management the experimental design discussed below. The evaluation experiment was approved but, due to a change in management, never performed. Nevertheless, the experiment is relevant here since it illustrates many of the practical issues associated with CSM evaluations. For example:

- A simple rule-based system written in two weeks would take nearly a year to evaluate.
- The evaluation criteria chosen made no reference to the internal structure of the rule base.
- The CSMs described below were slow to develop. In fact, only once a working knowledge base was developed could we reverse engineer a success criteria for the system.
- CSMs cannot be generated via mere *program watching*. A detailed analysis is offered below of the drawbacks of letting experts subjectively evaluate an expert system.

- The CSM collection implies only small changes to the running of the system. Further, using a pre-disaster point, the evaluation can occur without losing profit from the system. That is, evaluation may be practical even for systems as complicated as a large petrochemical plant.

5.1 Background to the QUENCH System

The Olefines petrochemical plant produces 240,000 tonnes of ethylene per year. It is a highly complex plant consisting in part of some 125 km of piping connecting numerous chemical processes. A unit of this plant is the quench oil tower. Inside the tower, hot cracked gases are cooled from around 400C to around 100C by mixing with oil. Certain gases are extracted at the top of the tower and the used quench oil, containing variable amounts of dissolved gases, is removed from the bottom. These dissolved gases effect the density of the removed oil. If the quench oil density moves outside of a certain narrow range, it can not be sold. In this case, ICI loses the revenue that would have come from its sale. Further, it must pay for the reprocessing or the disposal of the bad oil.

In order to keep the density on specification, the temperature at the bottom and the top of the tower must be maintained within one half of a degree of a target temperature. This is accomplished by altering the flow rates though the piping that surrounds the tower and/or by adjusting the heat exchange units attached to this piping. In practice, this is a non-trivial task. There have been cases when the operators of the tower have spent days attempting to return the density to an acceptable value. This process is directed by the the supervising engineers who communicate their instructions to the operators using heuristics similar to production rules. For example, to correct a very high quench oil density, an engineer could say to an operator:

```
if  the target temperatures are correct and
    the bottom temperature of the tower is high
then bring the bottom of the tower back on target
    by increasing the quench oil recycle flow
    rate by 20 tonnes per hour.
```

QUENCH contained 104 such rules.

5.2 Features of Large Petrochemical Plants

Large petrochemical plants have certain features that complicate the process of evaluation. Safety is a paramount consideration. Unsafe operating conditions could cost the lives of the workers in these plants.

Large petrochemical plants produce hundreds of millions of dollars worth of chemicals each year. The loss of a single day's revenue can cost a company hundreds of thousands of dollars. These economic imperatives are so pressing that the prolonged operation of these plants at less-than-optimum performance can not be tolerated.

There are major difficulties associated with deriving precise formalisations of these complex systems. For example, a mathematical model of the quench oil system would require the solutions of hundreds of simultaneous equations. Certain parameters required in these equations require uncertain physical properties data; i.e. these parameters are not known. Consequently it is possible that after months of development work, a mathematical model of the quench oil system may be grossly inaccurate. Without precise formalisations, the only way to accurately predict the effects of certain changes to the plant is to make those changes and observe the effects.

The design of these large plants is typically customised to meet local requirements. Hence, the experience gained in (e.g.) controlling quench oil towers in other plants may not be relevant to this quench oil tower. In fact, the two supervising engineers who helped write QUENCH's rules are the only authorities on the control of the Olefines' quench oil tower. In the jargon of the psychologist or the statistician, there is no control group available for experiments on the tower. Further, there is no objective expertise that can be called upon to accurately assess the suggestions made by quench oil tower experts (be they computers or human beings).

5.3 The Obvious Evaluation Method

One method for assessing the expertise of the program by running it in parallel with the existing system. The supervising engineers could compare QUENCH's suggestions with their own advice for problem situations. This method will be referred to as the *obvious method* and (the pre-disaster CSM evaluation will be called *the preferred method*). The obvious method has several advantages:

- It does not upset the normal operations of the plant.
- The plant remains under the control of the experts with the most experience on controlling the plant; i.e. the supervising engineers.
- It requires no control group.
- The computer and the human experts are being tested under identical plant conditions.

Regrettably, there are glaring design faults in the obvious method (discussed below).

5.4 Experimental Design Theory

Campbell and Stanley [Campbell & Stanley, 1970] assess experimental designs in terms of their *internal* and *external validity*.

Internal validity is the basic minimum without which any experiment is uninterpretable: Did in fact the experimental methods make a difference in this specific experimental instance? *External validity* asks the question of *generalisability*: To what populations, settings, treatment variables, and measurement variables can this effect be generalised? [Campbell & Stanley, 1970] (p4).

Internal validity is of particular concern. If we can not interpret the results of our experiment, then the experiment would have been pointless. Campbell and Stanley list several factors that could jeopardise internal validity. These factors have one feature in common: they could result in the effect of an experimental variable under study being confused with other factors. Each represents the effects of:

- *History*: events occurring between observations other those under study. The variable of history is relevant to the feature of *experimental isolation*. If the effect under study can not be isolated from other effects, then it is hard to distinguish the results of known influences from unknown influences.

- *Maturation*: changes over time in the test subjects. For example, the subject in an experiment may grow bored, tired, hungry, etc. and their reactions to various experimental variables may alter for reasons that are not under study. Maturation is a common problem with knowledge engineering research. Researchers often report improvement in some process when they use their own tools for a period of time (e.g [Runkel, 1995]). Such results hence conflate the effect of the tool with the effect of the developer learning how to best apply their own tool.
- *Testing*: the act of making an initial observation may somehow alter subsequent observations.
- *Instrumentation*: the calibration of the testing device changes. The springs of a weight scale may wear out, observers may change, or the reports of the same observer may alter as they gain experience with the experiment. Instrumentation is a common problem in many knowledge engineering studies. Knowledge engineering researchers rarely calibrate their measurements against some gold standard or *straw man*— an obviously inferior method (but some exceptions exist as noted below in the related work section) [Menzies, 1997a]. Without such a calibration, most knowledge engineering researchers can only say *software technology X lets me do task Y*. This is a less convincing statement than *software technology X lets me do task Y better than software technology Z*.
- *Statistical Regression*: the items/ people/ events in a test group are selected according to some extreme characteristic possessed by those items/ people/ events. Campbell and Stanley note that *the more deviant the score, the larger the error of measurement it probably contains* [Campbell & Stanley, 1970], (p11). They observe that test results from such extreme test groups tend to revert to the mean behaviour. For example, in an experiment testing some skill, observations could show that the dull could become brighter and the bright duller.
- *Selection*: if a test group was selected based upon their score on a certain measure, then this score could bias the behaviour of the test group in a certain way. Selection problems have been observed in the knowledge engineering literature. A bayesian system for medical diagnosis in a Leeds hospital apparently out-performed senior clinicians. However, a subsequent evaluation by another team in Copenhagen identified that the first study artificially restricted the number of possible diagnosis. When this restriction was removed, the performance of the bayesian system fell to 65 percent of that of human doctors [Gaschnig *et al.*, 1983] (p250–251). Also, Gashing *et.al.* report that a preliminary evaluation of the XCON system [McDermott, 1993] failed to detect flaws in XCON since it only studied a tiny fraction of the set of possible XCON inputs [Gaschnig *et al.*, 1983] (p270–271).
- *Mortality*: mortality refers to the changes to groups under comparison resulting from drop outs from the groups.

As to external validity, the claim of this paper is that the preferred method is generalisable to other expert system evaluations.

5.5 Assessing the Obvious Method

On several of the above points, the obvious method ranks quite well.

- *Maturation*: Maturation is not a problem since the test is not lengthy. The obvious method is an evaluation of expertise at a particular point in time.
- *Selection, and statistical regression*: The expert system is tested against whatever changes occur to the plant. Since there is no choice involved in selecting these test cases, these factors are not issues for this experimental design.
- *Mortality*: Mortality is only a issue to be considered for tests that take an appreciable period of time. Hence, it is not an issue with the obvious method.

However, the effects of history, instrumentation and testing are majors flaw in the obvious method.

- *History*: The suggestions of the supervising engineers are not always followed faithfully by the plant room operators. It is not uncommon for evening shift and night shift operators to ignore expert advice and apply their own control protocols. It is possible that these operators would tend to ignore a computer's advice even more than those of a human being. Having documented this problem, I now propose to ignore it. The resolution of this problem is an administrative problem that will be crucial to the process of evaluation. However, it that is beyond the scope of this researcher.
- *Instrumentation and Testing*: The program's expertise will be assessed by the supervising engineers. It is possible that their own perceptions of the program could alter with time. These engineers have been intimately connected with the program for several months. Human factors such as the halo effect (discussed above), egotistical considerations, disappointment or elation at their perceptions of the program's performance, etc., may distort their evaluation.

Hence, we reject the obvious method and move to the preferred method.

5.6 Defining CSMs for QUENCH

The preferred method requires CSMs and a pre-disaster point. This section offers CSMs. The next section offers a pre-disaster point.

There are three possibles CSMs for QUENCH:

1. A poll of all the electronic surveillance equipment that monitors the plant. This possibility is really a whole host of possibilities. There are many ways that the plant's surveillance equipment could be summed together into a single performance figure. Such a summation would be a whole research topic in itself. Fortunately, there are easier methods.
2. The time to failure. (a method proposed by Kehoe, personal communication). The time between the starting the trial and reaching the pre-disaster point could be the performance figure. The longer this time, the better the performance.
3. Revenue from quench oil (a method proposed by Dr. Michael Brisk, ICI, personal communication). The sum of revenues gained from processing the quench oil could be the performance measure. If the density goes off specification, and money must be spent to reprocess or dispose of the bad oil, then this amount should be deducted from the sum. Like the time to failure, the greater this figure, the better the performance.

Tag	Range
very high	> 1070
moderately high	> 1060
ok	> 1050
moderately low	> 1040
very low	<= 1040

Table 1: Assessing quench oil density in QUENCH. From [Menzies & Markey, 1987].

Tag	Range
rising quickly	> 7
rising slowly	> 2
steady	> -2
falling slowly	> -7
falling quickly	<= -7

Table 2: Defining *changes* in QUENCH. From [Menzies & Markey, 1987].

Methods two and three are not exclusive. The system could be studied using both criteria.

5.7 Defining the Pre-Disaster Point for QUENCH

We define the QUENCH pre-disaster point as follows: the point at which the supervising engineers realise that, despite their best efforts, the plant is defying their control strategies. If the plant reaches this pre-disaster point, then the control of the plant should be transferred to the best possible control system. In the case of testing QUENCH, the best possible control system is the supervising engineers. In the other case, when it is the engineers controlling the plant, the engineers would retain their authority to order alterations to the plant. They would then continue in their attempts to regain control over the plant processes.

Pre-disaster for QUENCH could be define as a bad quench oil density that was not improving, for (say) two days in succession. The time delay of two days allows for the expert time to recognize a problem, give advice for that problem, and for the tower to react to the expert's advice. If the at end of this time the density was still bad and not improving, then the expert would be deemed to have lost control of the tower.

The terms *bad* and *not improving* could be defined using the ranges developed during the implementation of QUENCH. The expert system has the ability to assigns *symbolic tags* to numeric ranges. The ranges for the quench oil density (expressed in kilograms per cubic meter) are shown in Table 1.

The time rate of change in the density (expressed in change in density per 24 hours) has the symbolic tags shown in Table 2.

Using these tags, we can define the pre-disaster point as a quench oil density that is either:

- moderately high or very high density and not falling quickly or falling slowly OR
- moderately low or very low density and not rising quickly or rising slowly.

6 Maturation and the Preferred Method

While the preferred method addresses the problem of objectivity seen with the obvious method, it will be effected by maturation. Consider the following:

- During stable operating periods, an evaluation of QUENCH's expertise in bringing the quench oil density back on specification is meaningless. Any test of this expertise must wait for periods of operational instability.
- The response time of the quench oil tower to changes in low rates and heat-exchangers can be as much as several days. Hence, once unstable conditions are encountered, an evaluation of the effectiveness of QUENCH's suggestions may have to wait for as much as a week.
- If we assume twenty trials for each population, and that each trial takes at least a week, then the total experiment time will be at least 40 weeks.
- The current version of the QUENCH rule set was developed in two weeks. As a result of assessing the current version of the program, the system developers would gain months of experience with the system. This experience could be used to modify and improve the program. Therefore...
- The evaluation process could result in substantial modifications to QUENCH's rule set.

Another way of expressing the above could be to say that the experiment is testing the expertise of a system that is learning. The evaluation experiment is to be attempted for an expert system who is in the shallow end of a learning curve. As a result of the experience gained during the evaluation process, the rule set would be improved and the expert system will move rapidly up the learning curve. The problem is that this improvement would occur concurrently with the experiment.

Kehoe (personal communication) offers an interesting resolution to the maturation problem. He argues that another CSM could be added to the system. Let F be the number of times the system is executed divided by the number of times the knowledge base is edited:

- If F tends to zero, the system is not being used.
- If F is less than one, then each run of the program is prompting a revision; i.e. there is something seriously wrong with QUENCH.
- If F is much greater than one, then the system is being run much more than it is being changed. Such an observation would suggest that some community finds using QUENCH to be of value.

Another response to this maturation problem would be to forbid the modification of the rules during the evaluation period; i.e. stop the system moving along the learning curve. This is an undesirable solution. It is highly probable that the existing rule set could be vastly improved. It was developed in a fortnight and this is a surprisingly short time for an expert system. Human cognitive processes are notoriously hard to formalise. The experience of expert systems developers is that any current specification of an expert solution to a problem is incomplete [Menzies, 1997b]. As experience with an expert system accumulates, inadequacies in the system's reasoning will always be

detected. To correct these inadequacies, the system's knowledge based (e.g. QUENCH's rule set) must be modified. This cycle of flaw detection followed by knowledge base modification can continue indefinitely but concludes when the user is satisfied that the system can provide adequate performance in an adequate number of cases. Depending on the expert system application, this refinement process can continue for many years. Compton reports one case where the modification process seemed linear; i.e. it may never stop [Compton, 1994].

This is not to say that the existing rule set lacks any utility for controlling the tower. The problem of assessing QUENCH only arose since the supervising engineers reported that they are satisfied with the output of the program. The short development time might have resulted from the choice of problem. QUENCH was ICI Australia's first direct experience with expert systems. The quench oil tower problem was selected as a comparatively simple first test case for the expert system methodology. One of the factors that made the problem simple was the Olefines' supervising engineers. These people spend significant amounts of their time explaining the workings of the Olefines plant to the control room operators. Hence, they have had considerable experience in expressing their knowledge in a concise manner.

Nevertheless, it is the author's belief that the program's rule set would benefit from further modification. It would be foolish to believe that QUENCH had somehow avoided the need for the long term knowledge base refinement process found to be necessary in other expert system application. Further, ICI would prefer the best possible control system for their tower. They may be less than enthusiastic about an experiment that inhibits the development of an optimum rule set. Hence, except for the Kehoe extension, I offer no revision to the preferred method to handle maturation.

7 Related Work

At the time of creating the QUENCH system, there was nothing in the petrochemical literature about empirical evaluation of expert systems. For example, in [Morari & McAvoy, 1986] and [Ctc96, 1986] we can read hundreds of pages on American and Japanese expert systems and never read anything about evaluation. Perhaps the reason for this curious omission is the difficulties inherent in the task. As seen above, a whole host of factors threaten the internal validity of evaluating experiments in such plants.

More generally, business-level empirical KBS evaluation is rarely performed in the knowledge engineering field (but some exceptions were noted in the introduction). By business-level, I mean measures of a running expert system which relate to the business case which motivated the development of that expert system. A CSM is a business-level evaluation measure. Elsewhere, I have criticised this lack of evaluations in the knowledge engineering field [Menzies, 1997b, Menzies, 1997a]. This critique motivated Feldman and Compton [Feldman *et al.*, 1989], followed by myself and Compton [Menzies & Compton, 1997], to devise and refine a general graph-theoretic abductive framework for assessing a KBS using a library of known or desired behaviour is discussed in [Menzies, 1995]. An example of using this framework is given in [Menzies & Compton, 1997]. One advantage of this framework over standard verification and validation is that the computational limits of the technique can be studied via *mutators* which auto-generate variants of known graphs [Menzies, 1996, Waugh *et al.*, 1997, Menzies *et al.*, 1997].

General principles for comparative empirical evaluation of knowledge engineering methods are discussed in [Menzies, 1997a]. Such comparative evaluations can take the

form of:

- Analysing program vs expert performance; e.g. [Hayes, 1997, Menzies *et al.*, 1992, Yu *et al.*, 1979]. In general, only these program vs expert performance evaluations yield results relevant to the business case that motivated the construction of the expert system.
- Analysing expert vs expert performance using different tools (e.g. [Corbridge *et al.*, 1995]) or records of their knowledge (e.g. [Shaw, 1988]);
- Analysing the performance of variants within some program either via an empirical average case analysis (e.g. [Waugh *et al.*, 1997, Menzies *et al.*, 1997]) or a theoretical analysis such as graph theory (e.g. [Menzies & Cohen, 1997]) or a worst-case time complexity analysis (e.g. [Tambe & Rosenbloom, 1994, Levesque & Brachman, 1985]).

The verification and validation community offer test procedures for KBS:

- The verification community typically focuses on syntactic anomalies within a KBS (e.g. circularities, tautologies) [Preece, 1992].
- The validation community focuses on the connection of the program to its environment. However, a typical validation paper focuses on (e.g.) automatic test case generation from an analysis of the dependency network within a program (e.g. [Ginsberg, 1990, Zlatareva, 1993]). The advantage of this technique is that it can be guaranteed that test cases can exercise all branches of a knowledge base. The disadvantage of this technique is that, for each proposed new input, an expert must still decide what constitutes a valid output. This decision requires knowledge external to the model, least we introduce a circularity in the test procedure (i.e. we test the a KBS using test cases derived from the structure of that KBS). Further, auto-test-generation focuses on incorrect features in the current model. I prefer to use some criteria from a totally external source since such external test cases can highlight
- Usually, publications from verification or validation community do not discuss how to assess a KBS with respect to the business case.

8 Conclusion

CSMs let us evaluate a system without requiring a panel of experts of a database of known or desired behaviour. A behavioral success criteria is derived from the business case that motivated the construction of the expert system. The system is then executed and measurements are made which inform the success criteria. Coupled with a *pre-disaster point*, CSMs let us statistically evaluate a system in operation, without compromising that operation.

The general themes of CSMs presented here are as follows. CSMs are usually very domain-specific since they reflect the contribution of the behaviour of the software in a particular business context. Hence, they typically do not refer to internal properties of a program and they cannot be developed by programmers without extensive input from business users. CSMs are usually obvious, but only in retrospect: a CSMs can take weeks of analysis to uncover. CSMs may only be collectible from the working system. However, CSMs should be explored very early in the life cycle of an expert system since

CSM collection may imply the extension of the system's design to collect the required data.

References

- [Basili, 1992] Basili, V. R. (1992). The Experimental Paradigm in Software Engineering. In Rombach, H. D., Basili, V. R., & Selby, R. W., (Eds.), *Experimental Software Engineering Issues: Critical Assessment and Future Directions, International Workshop, Germany*, pages 3–12.
- [Booch, 1996] Booch, G. (1996). *Object Solutions: Managing the Object-Oriented Project*. Addison-Wesley.
- [Campbell & Stanley, 1970] Campbell, D. & Stanley, J. (1970). *Experimental and Quasi-Experimental Designs for Research*. Rand McNally & Company.
- [Cohen, 1995] Cohen, P. (1995). *Empirical Methods for Artificial Intelligence*. MIT Press.
- [Compton, 1994] Compton, P. (1994). Personal communication. regarding the status of the PIERS system.
- [Corbridge *et al.*, 1995] Corbridge, C., Major, N., & Shadbolt, N. (1995). Models Exposed: An Empirical Study. In *Proceedings of the 9th AAAI-Sponsored Banff Knowledge Acquisition for Knowledge Based Systems*.
- [Ctc96, 1986] Ctc96 (1986). Special Issue on Expert Systems. Control Theory and Advanced Technology. Vol. 2, No. 3.
- [Feldman *et al.*, 1989] Feldman, B., Compton, P., & Smythe, G. (1989). Hypothesis Testing: an Appropriate Task for Knowledge-Based Systems. In *4th AAAI-Sponsored Knowledge Acquisition for Knowledge-based Systems Workshop Banff, Canada*.
- [Fenton *et al.*, 1994] Fenton, N., Pfleeger, S., & Glass, R. (1994). Science and Substance: A Challenge to Software Engineers. *IEEE Software*, pages 86–95.
- [Gaschnig *et al.*, 1983] Gaschnig, J., Klahr, P., Pople, H., Shortliffe, E., & Terry, A. (1983). Evaluation of Expert Systems: Issues and Case Studies. In Hayes-Roth, F., Waterman, D., & Lenat, D., (Eds.), *Building Expert Systems*, chapter 8, pages 241–280. Addison-Wesley.
- [Ginsberg, 1990] Ginsberg, A. (1990). Theory Reduction, Theory Revision, and Retranslation. In *AAAI '90*, pages 777–782.
- [Gruber, 1993] Gruber, T. (1993). A Translation Approach to Portable Ontology Specifications. *Knowledge Acquisition*, 5(2):199–220.
- [Hayes, 1997] Hayes, C. (1997). A Study in Solution Quality Human Expert and Knowledge-Based System Reasoning. In Feltovich, P., Ford, K., & Hoffman, R., (Eds.), *Expertise in Context*, chapter 14, pages 339–362. MIT Press.
- [Levesque & Brachman, 1985] Levesque, H. & Brachman, R. (1985). A Fundamental Tradeoff in Knowledge Representation and Reasoning (Revised Version). In Brachmann, R. & Levesque, H., (Eds.), *Readings in Knowledge Representation*, pages 41–70. Palo Alto, Morgan Kaufmann.
- [McDermott, 1993] McDermott, J. (1993). R1 ("XCON") at age 12: lessons from an elementary school achiever. *Artificial Intelligence*, 59:241–247.
- [Menzies, 1995] Menzies, T. (1995). *Principles for Generalised Testing of Knowledge Bases*. PhD thesis, University of New South Wales.
- [Menzies, 1997a] Menzies, T. (1997a). Evaluation Issues for Problem Solving Methods. Submitted to Banff KA workshop, 1998. Available from <http://www.cse.unsw.edu.au/~timm/pub/docs/97eval>.
- [Menzies, 1997b] Menzies, T. (1997b). Is Knowledge Maintenance an Adequate Response to the Challenge of Situated Cognition for Symbolic Knowledge Based Systems? Special issue of the International Journal of Human Computer Studies: "The Challenge of Situated Cognition for Symbolic Knowledge Based Systems". In press. Available from <http://www.cse.unsw.edu.au/~timm/pub/docs>.
- [Menzies, 1996] Menzies, T. (September, 1996). Applications of Abduction: Knowledge Level Modeling. *International Journal of Human Computer Studies*, 45:305–355.
- [Menzies *et al.*, 1992] Menzies, T., Black, J., Fleming, J., & Dean, M. (1992). An Expert System for Raising Pigs. In *The first Conference on Practical Applications of Prolog*.

- [Menzies & Cohen, 1997] Menzies, T. & Cohen, R. (1997). A Graph-Theoretic Optimisation of Temporal Abductive Validation. In *European Symposium on the Validation and Verification of Knowledge Based Systems, Leuven, Belgium*.
- [Menzies *et al.*, 1997] Menzies, T., Cohen, R., Waugh, S., & Goss, S. (1997). Evaluating Conceptual Qualitative Modeling Languages. In *Submitted to the Banff KAW '98 workshop*. Available from <http://www.cse.unsw.edu.au/~timm/pub/aka97/papers>.
- [Menzies & Compton, 1997] Menzies, T. & Compton, P. (1997). Applications of Abduction: Hypothesis Testing of Neuroendocrinological Qualitative Compartmental Models. *Artificial Intelligence in Medicine*, 10:145–175.
- [Menzies & Markey, 1987] Menzies, T. & Markey, B. (1987). A Micro-Computer, Rule-Based Prolog Expert-System for Process Control in a Petrochemical Plant. In *Proceedings of the Third Australian Conference on Expert Systems, May 13-15*.
- [Morari & McAvoy, 1986] Morari, M. & McAvoy, T. (1986). *Chemical Process Control: CPC III*. A Cache Publication.
- [Offen & Jeffery, 1997] Offen, R. & Jeffery, R. (1997). Establishing Software Measurement Programs. *IEEE Software*, pages 45–53.
- [Preece, 1992] Preece, A. (1992). Principles and Practice in Verifying Rule-based Systems. *The Knowledge Engineering Review*, 7:115–141.
- [Preston *et al.*, 1993] Preston, P., Edwards, G., & Compton, P. (1993). A 1600 Rule Expert System Without Knowledge Engineers. In Leibowitz, J., (Ed.), *Second World Congress on Expert Systems*.
- [Runkel, 1995] Runkel, J. (1995). Analyzing Tasks to Build Reusable Model-Based Tools. In *Proceedings of the 9th AAAI-Sponsored Banff Knowledge Acquisition for Knowledge-Based Systems Workshop Banff, Canada*.
- [Schreiber *et al.*, 1994] Schreiber, A. T., Wielinga, B., Akkermans, J. M., Velde, W. V. D., & de Hoog, R. (1994). CommonKADS. A Comprehensive Methodology for KBS Development. *IEEE Expert*, 9(6):28–37.
- [Shaw, 1988] Shaw, M. (1988). Validation in a Knowledge Acquisition System with Multiple Experts. In *Proceedings of the International Conference on Fifth Generation Computer Systems*, pages 1259–1266.
- [Tambe & Rosenbloom, 1994] Tambe, M. & Rosenbloom, P. (1994). Investigating Production System Representations for Non-combinatorial Match. *Artificial Intelligence*, 68(1).
- [Waugh *et al.*, 1997] Waugh, S., Menzies, T., & Goss, S. (1997). Evaluating a Qualitative Reasoner. In *Australian AI '97*. Available from <http://www.cse.unsw.edu.au/~timm/pub/docs>.
- [Yost, 1993] Yost, G. (1993). Acquiring Knowledge in Soar. *IEEE Expert*, pages 26–34.
- [Yu *et al.*, 1979] Yu, V., Fagan, L., Wraith, S., Clancey, W., Scott, A., Hanigan, J., Blum, R., Buchanan, B., & Cohen, S. (1979). Antimicrobial Selection by a Computer: a Blinded Evaluation by Infectious Disease Experts. *Journal of American Medical Association*, 242:1279–1282.
- [Zlatareva, 1993] Zlatareva, N. (1993). Distributed Verification and Automated Generation of Test Cases. In *IJCAI '93 workshop on Validation, Verification and Test of KBs Chambery, France*, pages 67–77.

Some of the Menzies papers can be found at <http://www.cse.unsw.edu.au/~timm/pub/docs>