# Applying Autonomic Diagnosis at Samsung Electronics

Paulo Casanova, Bradley Schmerl, David Garlan
*Carnegie Mellon University*
`{paulo.casanova,schmerl,garlan}@cs.cmu.edu`

Rui Abreu
*University of Porto*
`rui@computer.org`

Jungsik Ahn
*Samsung Electronics*
`jungsik.ahn@samsung.com`

## Abstract

An increasingly essential aspect of many critical software systems is the ability to quickly diagnose and locate faults so that appropriate corrective measures can be taken. Large, complex software systems fail unpredictably and pinpointing the source of the failure is a challenging task. In this paper we explore how our recently developed technique for automatic diagnosis performs in the automatic detection of failures and fault localization in a critical manufacturing control system of Samsung Electronics, where failures can result in large financial and schedule losses. We show how our approach can scale to such systems to diagnose intermittent faults, connectivity problems, protocol violations, and timing failures. We propose a set of measures of accuracy and performance that can be used to evaluate run-time diagnosis. We present lessons learned from this work including how instrumentation limitations may impair diagnosis accuracy: without overcoming these, there is a limit to the kinds of faults that can be detected.

## 1 Introduction

Product manufacture, notably technology manufacture, is highly automated, involving critical software to schedule, control, and monitor the manufacturing process. The large volume of items produced, the high competitiveness in manufacturing costs, and global supply chains make failures in this industrial setting costly, leading to schedule overruns and large financial losses.

The inherent complexity in the manufacturing process, exacerbated by extremely optimized production processes, and the high volume of production makes identifying faults in the underlying running software difficult. Software systems controlling the manufacturing process communicate by sending a large number of messages between equipment and the controlling software to coordinate manufacturing, and involve little human intervention. Although hardware faults (i.e., equipment failure) may also lead to software faults, in this paper we primarily address failures triggerd by software faults.

While this software system works most of the time, intermittent faults can cause cascading errors that can delay or halt manufacturing. In such a case, the software typically needs to be rebooted, causing further delays. Furthermore, isolating which parts of the system caused the initial problem is a difficult process, where developers manually examine volumes of log files to see if they can determine why the problems occurred.

In this paper, we report on our experience applying our run-time fault diagnosis approach [4] to dynamically and automatically diagnose and locate faults for a chip manufacturing process for Samsung Electronics. As there are several barriers (such as, rareness and unpredictability of the faults and dependencies between the software and the actual factory equipment) to performing monitoring and diagnosis in a real setting, we developed a simulator mimicking the behavior of the Samsung Electronics' manufacturing control system, which has been validated by Samsung Electronics.

In previous work we have tested our approach successfully on several research prototype software systems, reported in [4, 3]. In this paper we discuss how the technique works in a real-world setting (albeit simulated), with the scale and problems that exist in that system. In particular, we validate that our approach to autonomic monitoring and diagnosis (a) *accurately* identifies the fault of a problem by observing the run-time behavior; and (b) achieves the desired *performance* without violating quality attributes of the system. Furthermore, in our previous work we have argued that:

- The specifications of system behavior required for our approach were an intuitive concept for system designers.

- These computations could be reused within systems of the same architectural family.

- Defining correctness criteria for computations is simple, being directly derived from business requirements and quality attributes, and in terms familiar to system designers.

- The statistical analysis provided by the diagnostic algorithm can pinpoint the source of the fault accurately.

In this paper we analyze how our algorithms matched the challenges and how they fulfill the above claims. We further propose metrics to evaluate runtime diagnosis algorithms and summarize lessons learned, which we believe to be applicable in other run-time diagnosis systems.

This paper is organized as follows: in section 2 we provide a description of Samsung Electronics' system we're targeting. Section 3 describes our approach to autonomic diagnosis and how it was applied in Samsung Electronics. In section 4 we propose metrics for evaluation of run-time diagnosis algorithms. In section 5 we provide the results of our case study in a set of scenarios mimicking problems detected at Samsung Electronics. In section 6 we list lessons learned that can be useful for future systems. Section 7 summarizes our finds, discusses threats to validity and future work.

## 2 Target system

Samsung Electronics' manufacturing system is a large-scale industrial control system responsible for manufacturing control of semiconductors[1]. The system controls the stages of wafer manufacture, deciding to which equipment wafers are sent for processing. Furthermore, the software tracks of not only the lots being manufactured, but also the equipment used: which equipment is allocated, whether and what it is processing, what its output quality is, and so on.

The system is divided into subsystems, which perform specific tasks related to the manufacturing process. For example, the MOS (Manufacturing Operating System) controls the stages of the manufacturing process, and the ADS (Automatic Dispatch System) performs equipment scheduling, deciding which equipment is the best to perform a step in the manufacturing process.

Each of these subsystem is built by multiple concurrent processes that provide various manufacturing services, load distribution, and fault tolerance. These are connected through an event bus, which mediates event exchange amongst them. A subset of the system's architecture is abstractly depicted in Figure 1. Due to confidentiality reasons, specific details about system implementation cannot be reported. The total message

---

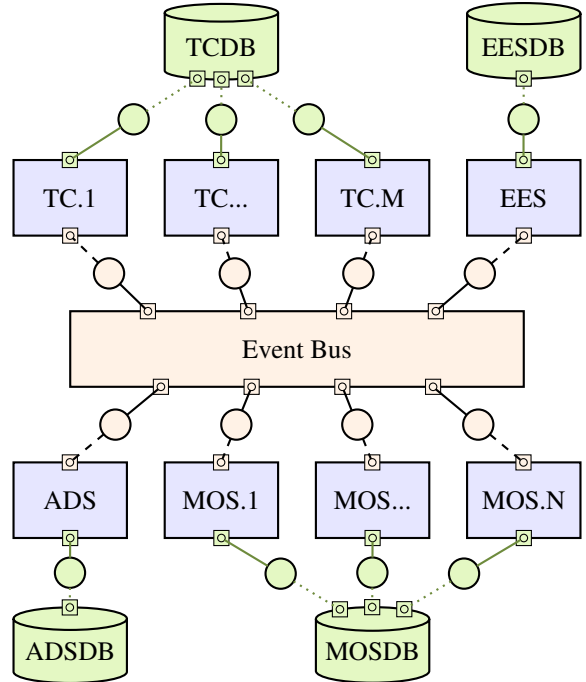[1]http://www.samsung.com/global/business/semiconductor/



Figure 1: High-level architectural view of the manufacturing system at Samsung Electronics. For simplicity, only the MOS and TC systems are shown decomposed.

throughput in all buses combined is around 2000 events per second.

The systems communicate with each other exchanging messages according to several predefined protocols. One such protocol is the *track-in* (TKIN) protocol. This protocol is used before a wafer lot is sent to another stage of processing. It determines what equipment the wafers should be sent to, performs validation operations and does some housekeeping, like ensuring that the equipment for processing the next steps are available and that no scheduling and quality constraints are violated.

### 2.1 Operation issues

While TKIN protocols are executed correctly most of the time, once in a while problems arise that cause failures in the system. These problems vary, but some of the typical problems are messages being lost (or not sent at all), messages sent too late, or unexpected messages being sent. Other problems are database performance slow-downs, which affect overall system performance.

Such problems can have a significant impact on the manufacturing process as they can lead to stalling lots and unneccesary equipment reservations. Given the sheer volume of messages being exchanged it is not possible for human operators to even realize that a problem has occurred until later when more serious consequences

become visible. Also, given the independent development of all the systems involved, it is difficult for the system developers to figure out what the problem is and where it is located.

Yet another difficulty is that these problems, although serious, are rare - and therefore are difficult to diagnose. Although the system operates with around 2000 messages exchanged per second, it generally functions correctly for many weeks in a row before any failure is detected. This sets up a scenario in which millions of observations are performed before a single failure can be identified.

## 2.2 Simulating the target system

There are several barriers to developing a diagnosis system that can be used for Samsung Electronics' manufacturing control systems. The criticality of the system makes testing on the production system unacceptable; the size of the system and its need to work connected to factory equipment makes it impossible to run it in a different environment. Also, the rareness and unpredictability of the faults make testing a diagnosis system difficult.

We addressed these issues by creating a simulator of Samsung Electronics' production system and performing diagnosis in the simulator. This simulator was engineered with the help of Samsung Electronics to produce a `TKIN` protocol similar to the real one and to allow manual control of faults for testing purposes.

The `TKIN` protocol is described using a message sequence diagram in Figure 2. The messages in the sequence diagram correspond to types of events sent: for example, `BEST_EQP` is a request for an equipment to process a wafer lot and `EQP` is the event with the ADS's decision on the equipment to perform the process step.

The simulated system does not, naturally, control any equipment. Rather, it generates messages that conform to the protocol, both in terms of their types and also in their timing. This simulator allows us to develop test scenarios in which we arrange for specific faults to happen, check whether they are diagnosed correctly, and have a credible expectation that it can be successfully transitioned into Samsung Electronics' production system.

The simulated system is similar to the one in Figure 1, but it has only two `MOS` components, one event bus, and two `TC` components. This is in contrast with the real system which includes more than 20 of each type. The system starts multiple `TKIN` protocols at random time intervals. Components have random processing delays and a single `TKIN` protocol takes around 1 minute to complete, similar to the `TKIN` protocol in the real system. The rate at which `TKIN` protocols are initiated can be used to control the level of concurrency.

Both `MOS` components work in fault-tolerance mode:
either `MOS.1` or `MOS.2` will receive messages addressed to the `MOS`, but not both. They maintain a "keep-alive" mechanism, allowing one to take over if the other fails to respond. The `TC` components work in "load-distribution" mode: each request will be forwarded to *either* `TC.1` or `TC.2`.

Implementing fault tolerance in our simulator was done by creating a synthetic component (`MOS.S`) that acts like the `MOS` for all other components (except `MOS.1` to `MOS.2`). It will forward through the event bus any message to either `MOS.1` or `MOS.2`, depending on which component is active. Load distribution is implemented similarly using a synthetic `TC.S` component.

Although we did not have access to the actual factory control software, we were careful to work with Samsung Electronics engineers to ensure that the simulator we developed was faithful to the real system in the following ways:

- event timings were designed to produce delays close to the real system;

- the protocol was designed according to Samsung Electronics' specifications;

- the faults that were injected (see Section 5) were typical of real problems experienced by Samsung Electronics and which are hard to locate in the real system.

We are therefore reasonably confident that the results reported in this paper will apply when Samsung Electronics transitions our approach to the real system.

## 3 Approach

This section details our approach to automatic monitoring and fault localization, as well as how we applied it to the Samsung Electronics' system.

## 3.1 Overview

A detailed description of our approach for autonomic diagnosis is described in [3, 4]; see Figure 3 for an overview. It consists, at a high level of abstraction, of defining a behavior model (equivalent, to some extent, to the `TKIN` protocol described in Figure 2) over the system's architecture (such as the one in Figure 1) and monitoring the system to identify patterns that match to the model.

These behavior specifications are high-level computations[2] that occur at the architecture level. These high-level computations are not directly observable in the system. Instead, low-level events such a message sent

---

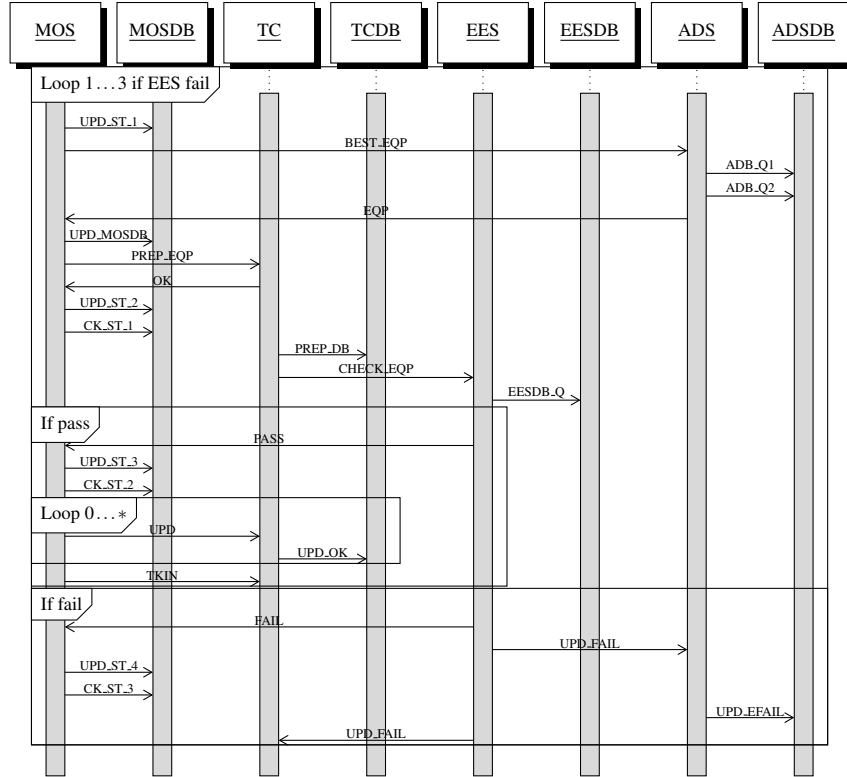[2]We termed them *transactions* in [4, 3], but use *computations* to avoid confusion with database transactions

Figure 2: Simulated TKIN protocol.

from the EES to the ADS with type UPD_FAIL for equipment lot X.445, can usually be extracted through instrumentation. These low-level events are then combined by the *Recognizer* using a *Behavior Model* (see figure 3), to form the high-level computations.

The high-level computations are then analyzed by an *Oracle* to check whether they represent correct or incorrect behaviors, according to some predefined *Correctness Criteria*. Correctness is derived from the business rules that define the system's functional requirements and quality attributes. For example, an incorrect computation can occur because a protocol did not complete successfully, or it completed, but without satisfying desired performance requirements. This information is sent to a *Fault Localizer* which will compute health estimates for the component and connectors in the system.

Having a set of computations over the architecture classified as correct and incorrect allows us to use static-time algorithms used for fault localization, such as the one in [1]. This algorithm was initially designed to localize faults in source code in the presence of a set of successful and unsuccessful runs of test cases (i.e., at development time, during the debugging phase of the software development life-cycle). We apply them at runtime, using observed run-time computations as "test cases", and localize the faults to the respective components in

the system's architecture. Being a technique that reasons over information about run-time behavior, the diagnosis accuracy increases with the number of observations.

The output of the fault localizer is a set of fault candidates ranked by the probability of being the explanation of the observed failures [1]. The fault localizer gathers information until it is able to produce an accurate report based on a limit on entropy. Borrowed from information theory, entropy is a measure of the uncertainty in the ranking produced by the diagnostic algorithm. This allows us to dynamically adjust the number of computations we need to observe before being able to accurately diagnose a fault.

The key ideas of our approach are:

- System designers reason about systems using high-level computations like the TKIN protocol, which are meaningful from a *business* perspective and whose structure and design is driven by the system's requirements.

- It is possible to classify these high-level computations as either *correct* or *incorrect* behavior.

  Because these computations are meaningful from a business perspective, their correctness is usually defined directly in terms of system requirements.
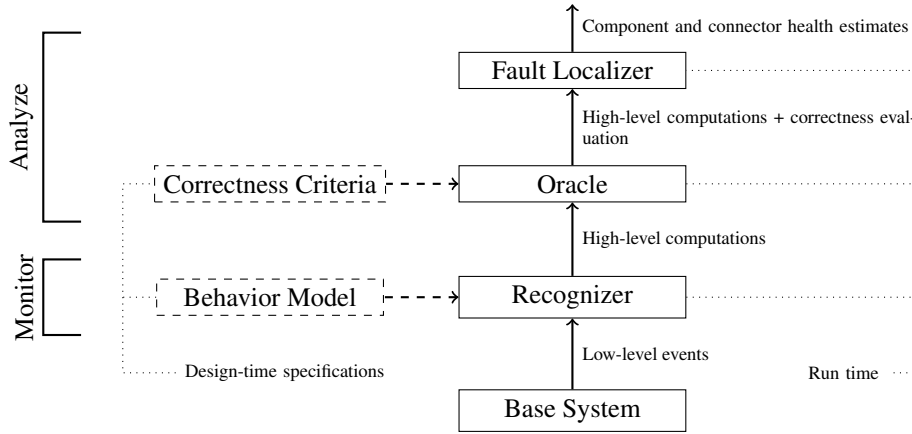
Figure 3: Structure of the diagnosis system.

For example, a `TKIN` protocol which halts half way through is considered *incorrect behavior*. Also, a `TKIN` protocol which takes more time to complete than the established maximum is also *incorrect behavior*.

- High-level computations of a system cannot generally be observed directly in a system, but can be detected by observing lower-level events.

  For example, it is not possible to directly see the `TKIN` protocol, as a whole, in the system. But we can see individual messages sent from components to other components. These individual messages (which *are* directly observable) allow us to reconstruct the higher-level behavior, the `TKIN`.

- Run-time observation of both correct and incorrect behavior can be fed into design-time fault localization algorithms as if each run-time observation was the execution of a test case.

## 3.2 Behavior specification language

To apply our diagnosis framework to a system, we have to provide both a behavior model of the system and the correctness criteria (see Figure 3), which are used as inputs by the Recognizer and the Oracle.

The behavior model describes how higher-level computations are built from lower-lever computations. This is done by declaring individual recognizers that will be instantiated at runtime. The structure of a recognizer is represented in Listing 1.

```
1  /* Declares the recognizer "name". The formal
       parameters represent the computation types
       this recognizer will use to recognizer a
       higher-level computation. */
2  recognizer name(a0 : a0_type, ...) {
```

```
3  /* The invariant contains a first-order predicate
       logic statement over the arguments that
       recognized computations must maintain. */
4      invariant ... ;
5  /* The emit clause (more than one may be present)
       defines what higher-level computations are
       recognized. */
6      emit ... ;
7  /* Alternatively, a emit failure clause says that
       we explicitly identified incorrect behavior.
       */
8      emit fail;
9  }
```

Listing 1: Structure of a recognizer.

An important aspect of recognition is how to treat events that do not fit into any pattern. We adopted a conservative approach: we assumed our knowledge of how the system works is limited and, therefore, more computations can be detected that we do not know about. Unmatched events are ignored.

In some cases we know that some events must always be matched according to some pattern, and the absence of a match signals a failure. For example, a request for an equipment without a response indicates a failure. Our approach supports detecting these situations by specifying recognizers that detect illegal computations. These recognizers, which make use of first-order logic's existential quantifiers, can detect that there is an event of type X for which there is no match with a required event. They will use the `emit fail` clause to report that incorrect behavior has been detected.

The correctness criteria is described using oracles that can be grouped into types to allow reuse. An example of an oracle is represented in Listing 2.

```
1  /* The oracle type construction defines a class of
       oracles that can be instantiated. */
2  oracle type latency_oracle {
3  /* Oracle types may have instance variables. The
       max_latency variable, of type "period" (a
       primitive data type representing a time span),
       contains the maximum latency allowed. */
```

```
4       m_max_latency : period;
5  /* The oracle's constructor. */
6       latency_oracle(max_latency : period) {
7           m_max_latency = max_latency;
8       }
9  /* The evaluate method receives as argument the
      type of computation that is evaluated by the
      oracle and returns true or false depending on
      whether the computation is correct or
      incorrect. */
10      bool evaluate(c : c_type) {
11          /* end() and start() are built-in
               functions that return when a
               computation ended and started,
               respectively. */
12          return end(c) - start(c) < m_max_latency;
13      }
14 }
15
16 /* Oracles clauses specify the oracles themselves,
       instances of their types. note that 15s in
       the argument is a literal value of "period"
       primitive data type. */
17 oracle limit_15s = new latency_oracle(15s);
```

Listing 2: Structure of an oracle

Both oracles and recognizers refer to *computation types*. Computation types have to be declared and may represent architecture-level computations, such as the TKIN, or simple events probed by the system's instrumentation. Conceptually, computation types are similar to Java classes. Listing 3 contains the computation type and associated declarations of the probed_ebus_msg computation, which represents a message flowing in the event bus probed by the system's instrumentation.

```
1  /*All messages in the TKIN. */
2  enum msg_type {
3      mt_best_eqp ,
4      mt_eqp ,
5      // ...
6  };
7
8  /* A struct represents a composite data type, like
       a Java class. The msg_data structure contains
       information about a lot being processed (the
       lot's unique ID) and the type of message. */
9  struct msg_data {
10     m_lot_id : string;
11     m_type : msg_type;
12 /* Constructor. Works like Java. */
13     msg_data(t : msg_type, l : string) {
14         m_lot_id = l;
15         m_type = t;
16     }
17 /* Function that checks whether two message data
       refer to the same "thread" (lot). */
18     bool same_thread(m : msg_data) {
19         return m_lot_id == m.m_lot_id;
20 };
21
22 /* A computation type is, syntax-wise, similar to
       a structure. Computations always have a start
       and end time and are associated with the set
       of architectural elements that contributed to
       them. These aspects of the computations are
       handle mostly automatically and are, therefore
       , invisible in the syntax. The colon
       represents inheritance in the object-oriented
       sense. */
23 computation type probed_ebus_msg : msg_data {
24 /* Name of the components that originated the
       message and to whom the message is destined. A
       single message may have multiple destinations
       . */
```

```
25     m_src_name : string;
26     m_dst_name : set of string;
27 /* Constructor. The msg_data constructor is
       invoked with the colon syntax. */
28     probed_ebus_msg(t : msg_type , l : string,
29         src_name : string,
30         dst_name : set of string)
31         : msg_data (t , l) {
32       m_src_name = src_name;
33       m_dst_name = dst_name;
34     }
35 };
```

Listing 3: Event bus message type.

Computation types, recognizers and oracles form the basis of the language used to recognize computations. In general, the diagnosis system's runtime keeps track of which components contributed to each computation as the computations are recognized into higher-level computations. Computations classified by oracles are sent, together with the list of components that contributed to them, to the fault localizer which will evaluate them using the barinel algorithm described in [1].

To perform fault localization, we need to define what the system's architecture is. Our language implements a subset of the Acme language [5] supporting the declaration of component and connect types and their instantiation. Listing 4 contains a partial declaration of the system structure.

```
1  /* A component type specifies a class of
       components. */
2  component type proc {
3  };
4
5  /* The ADS is a specialization of the proc
       component type. */
6  component type ads : proc {
7  };
8
9  /* An event bus represents a class of connectors.
       */
10 connector type ebus {
11 };
12
13 /* ads_1 is a concrete component. */
14 component ads_1 = new ads;
15
16 /* bus is the event bus. */
17 connector bus = new ebus;
```

Listing 4: Declaration of the system structure.

## 3.3 Application to the target system

Applying our approach to the system at Samsung Electronics required defining the behavior of the system in terms of the language described in Section 3.2. The main lines of our approach were:

- Each bus message was recognized to produce several messages: a message being *sent* by components and one message (or more, depending on the destinations) being *received* at components.

6

- The receive/send message pairs that contain call/return behavior, were recognized as a single computation.

  For example, the BEST_EQP received at the ADS is paired with the EQP sent by the ADS to produce a ads_eqp computation.

- Each pair also produced recognizers that detect behavior failures on timeouts.

  For example, ads_eqp_no_response recognizer recognizes a BEST_EQP for which *no* EQP exists.

- Sequential computations were then bound together in a single computation that represents a sequential flow. We ended up with three sequential flows: eqp_prep which corresponds to the flow until the EES decides on a pass or fail, a test_pass which corresponds to the "pass" optional section in the protocol flow and a test_fail which corresponds to the "fail" optional section in the protocol flow.

  For example, the ads_eqp is matched to the mos_start (recognized from a BEST_EQP detected at the MOS) to yield a mos_ads_eqp. The mos_ads_eqp is matched with a mos_prep (recognized as a EQP/PREP_EQP pair at the MOS) to yield a mos_eqp_prep and so on.

- The alternative flow was modeled using four recognizers: one that recognizes a eqp_prep and a test_pass into a eqp_pass, one that recognizes a eqp_prep and a test_fail into a eqp_fail, one that ensures either test_pass or test_fail, and one that recognizes a single_run from either an eqp_pass or an eqp_fail.

- Loops were modeled using recursion and loop limitation by placing a counter in the computation.

  For example, the tkin computation type has an attribute m_count that specifies the run count. A tkin can be recognized from another tkin and a single_run. The newly generated tkin has a m_count which is one value higher the the previous tkin.

It is worthy discussing more deeply the matching of call / return patterns within the TKIN protocol. If the contract for the ADS is, as is the case, to *always* reply with an EQP to a BEST_EQP request, then receiving an EQP and not replying with a BEST_EQP is a failure regardless of whether there is any other component involved. However, if the contract for the ADS were to only reply to *some* requests (for example, invalid requests would not have a response), then the pair BEST_EQP/EQP would not be enough to pinpoint a failure in the ADS component and the call / return pattern could not be applied.

The previous discussion hints at two important aspects of our approach: (1) because behavior is formally modeled, it reduces ambiguity in system specifications, and (2) the small pieces of the protocols follow patterns that are reusable across multiple systems, while the large, complete protocols generally do not.

# 4 Metrics for evaluation

To evaluate how well the diagnosis algorithm performs, we need to identify a set of appropriate metrics. There are two main areas of metrics that are applicable in domains sharing similarities with ours: metrics used for evaluation of classifiers and metrics used for evaluation of fault localization algorithms.

Our algorithm has similarities with classifiers in machine learning: it will output a classification of which components are faulty and which are not. In that sense, the classic metrics of precision and recall, or one of its variants, seem to be relevant. However, those statistics are applicable only in binary cases [6] and our classifier outputs a *probabilistically ranked* classification. Many components – maybe even all – may show up as outputs but their *probability*, or *rank* is relevant: if the faulty component has probability of 0.99 and all the others have probability $\frac{0.01}{n-1}$ where $n$ is the total number of components, this is generally seen as a good classification although it has a low value of precision. In the presence of few faulty components, recall will always tend to the high or low spectrum. These metrics are, therefore, not applicable.

The fault localization literature uses ranking for evaluation: the real probability attributed to diagnosis is only relevant in comparison with the others. The rationale behind this metric is that developers will tend, in order to correct the faults, to inspect components by rank order and, therefore, the lower the rank of the faulty component, the more effort will be wasted. This metric is applicable to evaluate run-time diagnosis.

However, at runtime, timing (performance) also becomes an important issue. If at design time, the time taken for diagnosis is not relevant, at runtime it may be critical for autonomic recovery. Also, at runtime, many diagnoses can be produced for the same system and they may yield different values. Therefore, we introduce two metrics:

*failure identification time* (**FIT**) that measures the time between when the fault is activated and when the diagnosis system presents some information that *something* is wrong;

*diagnosis stabilization time* (**DST**) that measures the time between failure identification and diagnosis

stabilization: when the system decides on which diagnosis to consider final.

While on some systems – like the one presented in this case study – it is the sum of both metrics that is relevant, in other circumstances they may have different impacts: if a breach in a high-security system is detected (failure identification time) we may want to disconnect the system from the network immediately even before computing which part of the system was breached (diagnosis stabilization time).

## 5 Evaluation

We evaluated the system by defining the computations using the specification language proposed. We then identify a set of *scenarios* in which the system would fail in some predicted way. These scenarios were designed to simulate types of problems that mimic real problems observed at Samsung Electronics.

1. EES will send a FAIL message and, 3s later, will send the PASS message. Because the MOS will respond to both events, this will start two parallel flows: one for the successful evaluation, which will end with a TKIN, and another for the unsuccessful evaluation, which will end with a retry of the whole process.

2. TC database is too slow to respond.

3. TC will send CHECK_EQP to ADS instead of EES.

4. Active MOS fails and is later replaced by passive MOS.

5. MOS retries four or more times.

6. ADS does not issue ADB_Q1.

7. The ADS database is too slow to respond.

Evaluation addressed the two main goals stated in the introduction: (1) diagnosis accuracy and (2) diagnosis performance. Based on our description in Section 4, we measured the accuracy of the system as the rank of the faulty component in the diagnosis and we measured performance by computing the failure identification time and diagnosis stabilization time.

We also computed another measure, not directly related to the diagnosis system, but of relevance to Samsung Electronics: an estimate of the maximum throughput of the system. Since every recognizer can be run in a different system, the maximum throughput of the system is limited by the amount of time the slowest recognizer takes to perform its evaluation.

| Scenario | Component† | Rank‡ |
|---|---|---|
| 1: EES sends FAIL and PASS | EES | 1 |
| 2: TC database is too slow | TCDB | 1 |
| 3: TC will send message to ADS instead of EES | TC | 1 |
| 4: MOS fails and is replaced by standby | MOS.1 | 1 |
| 5: MOS will retry 4 times | MOS.1* | 1 |
| 6: ADS will not issue ADS_Q1 | -** | -** |
| 7: ADS database is too slow | ADSDB | 2*** |

Table 1: Evaluation of performance metrics in the scenarios.
† Component where the fault was injected.
‡ Average rank among 10 scenarios.
* MOS.1 was the active MOS.
** No failure was detected by the diagnosis system in this scenario.
*** ADS and ADSDB were both ranked with an equal probability of 0.5 in all 10 scenarios. We count as 2 as a developer / system maintainer may want to inspect the ADS before the database. We took the most conservative approach.

### 5.1 Accuracy evaluation

The average result of measuring accuracy in 10 scenarios is presented in Table 1. These results show that the system was able to correctly detect a failure in all scenarios except scenario 6. Because databases are a source of performance bottlenecks at Samsung Electronics, it was not possible to determine communication between ADS and ADSDB. Consequently, the diagnosis system has no way of distinguishing between the slowness of ADS and ADSDB; each component is equally likely to be the cause of the slowness. Interestingly, we are able to distinguish the slowness of TCDB. Because the two TC components connect to the same database, slowness in both increases the probability of the fault being localized in the database.

### 5.2 Performance evaluation

Performance metrics evaluated in all scenarios are shown in Table 2. For each evaluation scenario we ran the scenario 10 times and collected the two performance metrics discussed in Section 4. The fault localizer in Figure 3 outputs the results regularly (once every second).

As said before, the FIT is the time between the fault being activated and a diagnosis being produced that includes an identification of a failure. If this diagnosis result already included the failed component ranked in first position, then the diagnosis stabilization time (DST) is 0. In several scenarios, a later diagnosis never placed another component in first position, meaning that most scenarios have an instantaneous DST. In scenario 2 however, because TCDB was not instrumented, the detection

that it was the source of the failure required more data and, therefore, DST is greater than 0.

The minimum FIT column is added for reference. It is the minimum theoretically possible FIT that would detect the failure. In scenario 1, for example, the `FAIL` message is sent 3s after the `PASS` message so, before that, no failure can be detected although the invalid code has already been started. In some examples, like this one, we could have deducted the minimum FIT from the FIT but due to scenarios in which the timing is not so predictable (like introduced database slowness in scenarios 2 and 7) this would not be consistent.

In the classical dependability taxonomy [2], FIT is actually measuring the time since the *error* was introduced, not since the *failure* was observed.

## 6 Lessons learned

The application of our diagnosis framework to the system at Samsung Electronics provided us with confirmation that several of our assumptions appear to be supported in an industrial setting and also provided us with insight into some other areas.

**Decomposing computations is critical to the success of the diagnosis.** As a result of the decomposition, failures can be identified not only in the high-level computations but, sometimes, also in lower level computations. This result allows diagnosis sometimes to be performed in a much smaller set of components yielding a much faster response and a higher accuracy. For example, a lack of `EQP` response allows inferring immediately a fault in the `ADS` component due to the request / response low-level computation model.

**Even with strict protocols, reasoning-based diagnosis analysis is sometimes necessary.** The previously mentioned decomposition allows local detection and identification of the source of the failure in some some cases. This is consistent with industrial practice in which many instances of this reasoning – like heartbeats or timeouts – are used. However, in some cases this is not sufficient. As seen, for example, in scenario 2 where the `TC` database slows down, more complex reasoning may be required to accurately pinpoint the source of the failure.

**Behavior specifications can be built up using patterns that are reusable across systems.** The computations that comprise the behavior specification we build are defined on top of smaller patterns of well-known computation styles like call/return, alternative flows or loops. Therefore, although different systems may need to specify their own idiosyncratic high-level computations, a significant number of specification blocks may be reused across systems.

**Domain experts are able to provide correctness specifications for high-level computations.** Domain experts would not be able to explicitly state that within 15 seconds after a `BEST_EQP` message had been seen in the event bus, an `EQP` message from the same ADS to the same MOS with the same `lot_id` field should be seen in the event. This reasoning involves a great deal of low-level detail. However, domain experts *were* able to state that the ADS should respond within 15 seconds of a request for an equipment. This means high-level computations have an abstraction level that matches the domain experts' understanding of the system and are, therefore, a good level to write correctness specifications.

**Instrumentation design must consider a trade-off between system performance, diagnosis accuracy and diagnosis performance.** The limitation on the observability of some events creates uncertainties which, in some cases like scenario 2 were solved by the reasoning-based analysis with a performance penalty. But in other cases, like scenario 7, this was not possible. However, instrumenting the connection between the ADS and its database was not possible because it was considered by Samsung to introduce an unacceptable performance penalty.

This means that diagnosis accuracy and diagnosis performance have to trade with other system quality attributes. In this system they traded with system performance but in other systems they may have to trade with other quality attributes. For example, probing some connections may lead to easier information leaks, so security could be another quality attribute that would be negatively affected by the introduction of autonomic diagnosis.

## 7 Conclusions and future work

The results indicated by this case study are encouraging as they show that the expected results match practical experimentation and several of our assumptions hold at least in this practical scenario. These conclusions are, however, still bound to be the result of a *simulation* of the real system.

We strongly believe that, in spite of being shown to work on a simulator, our work would be usable in the real industrial setting: adapting the diagnosis system to the real system will require some engineering, but will not raise any new fundamental problems. The diagnosis system can receive the messages from the real system's event bus, just as it receives messages from the simulator itself; the real system events have all the information required for the diagnosis to work; the timings on the simulator are randomized and the real system's timings would also be seen as random, although, most likely, with a different distribution.

| Scenario | avg FIT[†] | std FIT[††] | Min FIT[†††] | avg DST[‡] | std DST [‡‡] |
|---|---|---|---|---|---|
| 1: EES sends FAIL and PASS | 5,506 | 319 | 3,000 | 0 | 0 |
| 2: TC database is too slow | 13,506 | 3,212 | 5,000-15,000 | 436 | 432 |
| 3: TC will send message to ADS instead of EES | 17,708 | 309 | 15,000 | 0 | 0 |
| 4: MOS fails and is replaced by standby | 15,197 | 1,505 | 10,000-15,000 | 0 | 0 |
| 5: MOS will retry 4 times | 2,861 | 0,343 | 0 | 0 | 0 |
| 6: ADS will not issue ADS_Q1 | -* | -* | -* | -* | -* |
| 7: ADS database is too slow | 10,204 | 1,131 | 5,000-10,000 | 0 | 0 |

Table 2: Evaluation of performance metrics in the scenarios. All results in milliseconds.

[†] Average failure identification time.  [††] Standard deviation of failure identification time.  [†††] Theoretical minimum failure identification time.  [‡] Average diagnosis localization time.  [‡‡] Standard deviation of diagnosis localization time.  [*] No failure was detected by the diagnosis system in this scenario.

The main threat to the validity of the simulation is scalability: our simulated system contains only two instances of the MOS and the TC and only one instance of the ADS and the EES. The performance of our recognizers and oracles can be made independent of the number of components of the architecture through horizontal partitioning of the event space. Their number increases but they may be run in parallel if we need to scale up. The point of contention may be the fault recognizer whose complexity raises with the number of components. However, experimentation at design time in [1] has shown that it can handle large numbers of components.

With respect to future work, there are two main complementary paths that we plan to follow: the industrial path and the research path. On the industrial path, we plan to continue to work to bring this system into Samsung Electronics' real system.

On the research path, integration of diagnosis with a self-adaptive framework is a critical piece. It is clear that diagnosis can be used to drive self-repair. But because diagnosis introduces the need to perform trade-offs (for example, between performance and accuracy), it is also reasonable to assume that, ideally, the self-adaptive control loop should tune diagnosis at run time.

Also, in this work we had to manually handle the problem of non-observability: we added the databases to the computations because we knew they were involved even though we did not observe the events. However, we think this approach is more general. The observed behavior could be a projection of the complete behavior in which non-observable computations have been removed. The recognizer and the oracle will have to be updated accordingly, but further research is required to automate this process and understand the implications for diagnosis accuracy and performance.

# References

[1] R. Abreu and A. J. C. van Gemund. Diagnosing multiple intermittent failures using maximum likelihood estimation. *Artif. Intell.*, 174(18):1481–1497, 2010.

[2] A. Avižienis, J.-C. Laprie, B. Randell, and C. Landwehr. Basic concepts and taxonomy of dependable and secure computing. *IEEE Transactions on Dependable and Secure Computing*, 1:11–33, 2004.

[3] P. Casanova, D. Garlan, B. Schmerl, and R. Abreu. Diagnosing architectural run-time failures. In *Proceedings of the 8th International Symposium on Software Engineering for Adaptive and Self-Managing Systems*. To appear., 20-21 May 2013.

[4] P. Casanova, B. R. Schmerl, D. Garlan, and R. Abreu. Architecture-based run-time fault diagnosis. In *Proceedings of the 5th European Conference on Software Architecture (ECSA'11)*, pages 261–277, 2011.

[5] D. Garlan, R. T. Monroe, and D. Wile. Acme: Architectural description of component-based systems. In G. T. Leavens and M. Sitaraman, editors, *Foundations of Component-Based Systems*, pages 47–68. Cambridge University Press, 2000.

[6] D. M. W. Powers. Evaluation: From Precision, Recall and F-Factor to ROC, Informedness, Markedness & Correlation. Technical Report SIE-07-001, School of Informatics and Engineering, Flinders University, Adelaide, Australia, 2007.