

Simultaneous debugging of software faults[☆]

Rui Abreu^{a,*}, Peter Zoetewij^b, Arjan J.C. van Gemund^c

^a Department of Informatics Engineering, Faculty of Engineering, University of Porto, Rua Dr. Roberto Frias s/n, 4200-465 Porto, Portugal

^b IntelliMagic B.V., The Netherlands

^c Embedded Software Department, Faculty of Electronics, Math, and CS, Delft University of Technology, The Netherlands

ARTICLE INFO

Article history:

Received 9 January 2010
Received in revised form 8 September 2010
Accepted 22 November 2010
Available online 7 December 2010

Keywords:

Software fault diagnosis
Program spectra
Statistical and reasoning approaches

ABSTRACT

(Semi-)automated diagnosis of software faults can drastically increase debugging efficiency, improving reliability and time-to-market. Current automatic diagnosis techniques are predominantly of a statistical nature and, despite typical defect densities, do not explicitly consider multiple faults, as also demonstrated by the popularity of the single-fault benchmark set of programs. We present a reasoning approach, called Zoltar-M(ultiple fault), that yields multiple-fault diagnoses, ranked in order of their probability. Although application of Zoltar-M to programs with many faults requires heuristics (trading-off completeness) to reduce the inherent computational complexity, theory as well as experiments on synthetic program models and multiple-fault program versions available from the software infrastructure repository (SIR) show that for multiple-fault programs this approach can outperform statistical techniques, notably spectrum-based fault localization (SFL). As a side-effect of this research, we present a new SFL variant, called Zoltar-S(ingle fault), that is optimal for single-fault programs, outperforming all other variants known to date.

© 2010 Elsevier Inc. All rights reserved.

1. Introduction

Automatic software fault localization (also known as fault diagnosis) techniques aid developers to pinpoint the root cause of detected failures, thereby reducing the debugging effort. Two approaches can be distinguished:

- (1) the *spectrum-based fault localization* (SFL) approach, which correlates software component activity with program failures (a *statistical* approach) (Abreu et al., 2007; Gupta et al., 2005; Jones et al., 2002; Liu et al., 2005; Renieris and Reiss, 2003; Zeller, 2002), and
- (2) the *model-based diagnosis* or *debugging* (MBD) approach, which deduces component failure through *logic reasoning* (de Kleer and Williams, 1987; Feldman et al., 2008; Mayer and Stumptner, 2008; Wotawa et al., 2002).

Because of its low computational complexity, SFL has gained large popularity. Although inherently not restricted to single faults,

in most cases these statistical techniques are applied and evaluated in a single-fault context, as demonstrated by the benchmark set of programs widely used by researchers,¹ which is seeded with only one fault per program (version). In practice, however, the defect density of even small programs typically amounts to multiple faults. Although the root cause of a particular program failure need not constitute multiple faults that are acting *simultaneously*, many failures will be caused by *different* faults. Hence, the problem of *multiple-fault localization* (diagnosis) deserves detailed study.

Unlike SFL, MBD traditionally deals with multiple faults. However, apart from much higher computational complexity, the logic models that are used in the diagnostic inference are typically based on static program analysis. Consequently, they do not exploit execution behavior, which, in contrast, is the essence of the SFL approach. Combining the dynamic approach of SFL with the multiple-fault logic reasoning approach of MBD, in this paper, we present a multiple-fault reasoning approach that is based on the dynamic, spectrum-based observations of SFL. Additional reasons to study the merits of this approach are the following.

- Diagnoses are returned in terms of multiple faults, whereas statistical techniques return a one-dimensional list of single fault locations only. The information on fault multiplicity is attractive from parallel debugging point of view (Jones et al., 2007).

[☆] This work has been carried out as part of the TRADER project under the responsibility of the Embedded Systems Institute. This project is partially supported by the Netherlands Ministry of Economic Affairs under the BSIK03021 program.

* Corresponding author.

E-mail addresses: rma@fe.up.pt, ruicomputer.org (R. Abreu), p.zoetewij@gmail.com (P. Zoetewij), a.j.c.vangemund@tudelft.nl (A.J.C. van Gemund).

¹ <http://sir.unl.edu/>.

- Unlike statistical approaches, multiple-fault diagnoses only include valid candidates, and are asymptotically optimal with increasing test information (Abreu et al., 2008).
- The ranking of the diagnoses is based on probability instead of similarity. This implies that the quality of a diagnosis can be expressed in terms of information entropy or any other metric that is based on probability theory (Pietersma and van Gemund, 2006).
- The reasoning approach naturally accommodates additional (model) information about component behavior, increasing diagnostic performance when more information about component behavior is available.

To illustrate the difference between multiple-fault and the statistical approach, consider a triple-fault (sub)program with faulty components c_1 , c_2 , and c_3 . Whereas under ideal testing circumstances a traditional SFL approach would produce multiple single-fault diagnoses (in terms of the component indices) like $\{\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \dots\}$ (ordered in terms of statistical similarity), a multiple-fault approach would simply produce one single multiple-fault diagnosis $\{\{1, 2, 3\}\}$. Although the statistical similarity of the first three items in the former diagnosis would be highest, the latter, single diagnosis unambiguously reveals the actual triple fault.

Despite the above advantages, a reasoning approach is more costly than statistical approaches because an exponential number of multiple-fault candidates need to be processed instead of just the (M , being M the number of components in the system under analysis) single-fault candidates. In this paper, we compare our reasoning approach to several statistical approaches. Our study is based on random synthetic spectra, as well as on several benchmark programs, extended by us to accommodate multiple faults. More specifically, this paper makes the following five contributions.

- We introduce a multiple-fault diagnosis approach that originates from the model-based diagnosis area, but which is specifically adapted to the interaction dynamics of software. The approach is coined Zoltar-M (Zoltar for the name of our debugging tool set (Janssen et al., 2009),² M for multiple-fault).
- We show how our reasoning approach applies to single-fault programs, yielding a provably optimal SFL variant, called Zoltar-S (S for single-fault), as of yet unknown in literature.
- We introduce a general, multiple-fault, probabilistic program (spectrum) model, parameterized in terms of size, testing code coverage, and testing fault coverage, to theoretically study Zoltar-M, compared to statistical techniques such as Tarantula and Zoltar-S.
- We extend the traditional, single-fault benchmark set of programs (referred to as SIR-S) with a multiple-fault version (SIR-M), by combining the existing single-fault versions, to empirically evaluate debugging performance under realistic, multiple-fault conditions.
- We investigate the ability of all techniques to deduce program fault multiplicity, which is aimed at providing a good estimate to guide parallel debugging, using an approach that substantially differs from Jones et al. (2007).

To the best of our knowledge, this is the first paper to specifically address software multiple-fault localization using a spectrum-based, logic reasoning approach, yielding two new localization techniques Zoltar-S and Zoltar-M, implemented within our Zoltar SFL framework. Our experiments confirm that Zoltar-S is superior to all known similarity coefficients for the Siemens-

```
void RationalSort(int n, int *num, int *den)
{
    /* block 1 */
    int i, j, temp;

    for ( i=n-1; i>=0; i-- ) {
        /* block 2 */
        for ( j=0; j<i; j++ ) {
            /* block 3 */
            if (RationalGT(num[j], den[j],
                          num[j+1], den[j+1])) {
                /* block 4 */
                temp = num[j];
                num[j] = num[j+1];
                num[j+1] = temp; } } }
}
```

Fig. 1. A faulty C function for sorting rational numbers.

S benchmark. More importantly however, our experiments for multiple-fault programs show that although for synthetic spectra Zoltar-M is outperformed by Zoltar-S, for our SIR-M experiments Zoltar-M outperforms all similarity coefficients known to date.

The paper is organized as follows. In the next section, we present the concepts and terminology used throughout the paper. In Section 3, our multiple-fault localization approach is described, as well as a derivation of the optimal similarity coefficient for single-fault programs. In Section 4, the approaches are theoretically evaluated, and in Section 5, real programs are used to assess the capabilities of the studied techniques for fault localization. Related work is discussed in Section 6. Preliminary results of Sections 4 and 5 appeared in Abreu et al. (2009a). We conclude and discuss future work in Section 7.

2. Preliminaries

In this section, we introduce basic definitions as well as the traditional SFL approach. As defined in Avizienis et al. (2004), in the remainder of this paper, we use the following terminology.

- A *failure* is an event that occurs when delivered service deviates from correct service.
- An *error* is the part of the total state of the system that may cause a failure.
- A *fault* is the cause of an error in the system.

To illustrate these concepts, consider the C function in Fig. 1. It is meant to sort, using the bubble sort algorithm, a sequence of n rational numbers whose numerators and denominators are passed via parameters `num` and `den`, respectively. There is a fault (bug) in the swapping code of block 4: only the numerators of the rational numbers are swapped. The denominators are left in their original order.

A failure occurs when applying `RationalSort` yields anything other than a sorted version of its input. An error occurs after the code inside the conditional statement is executed, while `den[j] ≠ den[j+1]`. Such errors can be temporary: if we apply `RationalSort` to the sequence $\langle 4/1, 2/2, 0/1 \rangle$, an error occurs after the first two numerators are swapped. However, this error is “canceled” by later swapping actions, and the sequence ends up being sorted correctly. Faults do not automatically lead to errors either: no error will occur if the input is already sorted, or if all denominators are equal.

² <http://www.fdir.org/zoltar>.

The purpose of *diagnosis* is to locate the faults that are the root cause of detected errors. As such, error detection is a prerequisite for diagnosis. As a rudimentary form of error detection, failure detection can be used, but in software more powerful mechanisms are available, such as pointer checking, array bounds checking, deadlock detection, etc.

In a software context, faults are often called *bugs*, and diagnosis is part of *debugging*. Computer-aided techniques as the one we consider in this paper are known as *automated debugging*.

2.1. Basic definitions

A program that is being diagnosed comprises a set of M components (statements in the context of this paper), which is executed using N test cases that either pass or fail. Program (component) activity is recorded in terms of program spectra (Harrold et al., 1998). This data is collected at run-time, and typically consists of a number of counters or flags for the different components of a program. In the context of this paper, we use the so-called hit spectra, which indicate whether a component was involved in a (test) run or not.

Both spectra and program pass/fail (test) information is input to SFL, as well as to our reasoning technique. The program spectra are expressed in terms of the $N \times M$ activity matrix A . An element a_{ij} is equal to 1 if component j was observed to be involved in the execution of run i , and 0 otherwise. For $j \leq M$, the row A_i indicates whether a component was executed in run i , whereas the column A_j indicates in which runs component j was involved. The pass/fail information is stored in a vector e , the error vector, where e_i signifies whether run i has passed ($e_i = 0$) or failed ($e_i = 1$). Note that the pair (A, e) is the only input to the techniques studied in this paper. From (A, e) , we can derive the probability r that a component is actually executed in a run (testing code coverage), and the probability g that a faulty component is actually exhibiting good behavior (testing fault coverage, also known as the “goodness” parameter g from MBD (de Kleer, 2007)).

Programs can have multiple faults, the number being denoted C (fault cardinality). A *diagnosis candidate* is expressed as the set of indices of those components whose combined faulty behavior is logically consistent with the observations A and therefore must be considered as a collective candidate. A *diagnosis* is the ordered set of diagnostic candidates $D = \{d_1, \dots, d_k\}$, all of which are an explanation consistent with observed program behavior (A), ordered in probability of being the program’s actual multiple fault condition. An example multiple-fault diagnosis is the diagnosis $\{d_1\} = \{\{1, 2, 3\}\}$ given in the Introduction. For brevity, we will often refer to diagnostic candidates as diagnoses as well, as it is clear from the context whether we refer to a single diagnosis candidate or to the entire diagnosis.

2.2. Traditional SFL

In SFL, one measures the similarity between the error vector e and the activity profile vector A_j for each component j . This similarity is quantified by a *similarity coefficient*, expressed in terms of four counters $n_{pq}(j)$ that count the number of positions in which A_j and e contain respective values p and q , i.e., for $p, q \in \{0, 1\}$, we define

$$n_{pq}(j) = |\{i | n_{ij} = p \wedge e_i = q\}|$$

Two examples of well-known coefficients are

$$S_T = \frac{n_{11}(j)/(n_{11}(j) + n_{01}(j))}{n_{11}(j)/(n_{11}(j) + n_{01}(j)) + n_{10}(j)/(n_{10}(j) + n_{00}(j))}$$

as used by the Tarantula tool (Jones et al., 2002), and the Ochiai coefficient

$$s_0 = \frac{n_{11}(j)}{\sqrt{(n_{11}(j) + n_{01}(j)) * (n_{11}(j) + n_{10}(j))}} \quad (1)$$

known from molecular biology, introduced in SFL in Abreu et al. (2007).

Under the assumption that a high similarity to the error vector indicates a high probability that the corresponding parts of the software cause the detected errors, the calculated similarity coefficients rank the parts of the program with respect to their likelihood of containing the faults. Algorithm 1 concisely describes the SFL approach to fault localization.

As an example, suppose we have a program with $M = 7$ components, of which c_1 , c_2 , and c_3 are faulty, with A as given in Table 1. The table also includes the n_{pq} counts as well as the resulting similarity based on the Tarantula and Ochiai coefficients. Assuming that a developer would follow the ranking produced by the techniques, Tarantula requires him/her to inspect more components in order to find a faulty one. The first faulty component ranked by Tarantula is at the 3rd place of the list, whereas with Ochiai it is already at the 2nd place. The results shows the sensitivity of Tarantula to components that are not involved in passed runs (n_{00}), considering them likely to be the faulty one and not taking into account their involvement in failed runs (e.g., c_4 and c_7). Ochiai, however, exonerates components based on their involvement in passed runs (n_{10}), and absence in failed runs (n_{01} , for detailed comparison, see Abreu et al., 2007).

Algorithm 1 (SFL Algorithm).

Require: Activity matrix A , error vector e , number of runs N , number of components M , and similarity coefficient s

Ensure: Diagnostic report D

```

1    $D \leftarrow \emptyset$ 
2   for  $j = 0$  to  $M$  do
3      $n_{11}(j) \leftarrow 0$ 
4      $n_{10}(j) \leftarrow 0$ 
5      $n_{01}(j) \leftarrow 0$ 
6      $n_{00}(j) \leftarrow 0$ 
7      $S[j] \leftarrow 0$   $\triangleright$  Similarity  $s$  of component  $j$ 
8   end for
9   for  $i = 0$  to  $N$  do
10    for  $j = 0$  to  $M$  do
11      if  $a[i, j] = 1 \wedge e[i] = 1$  then
12         $n_{11}(j) \leftarrow n_{11}(j) + 1$ 
13      else if  $a[i, j] = 0 \wedge e[i] = 1$  then
14         $n_{01}(j) \leftarrow n_{01}(j) + 1$ 
15      else if  $a[i, j] = 1 \wedge e[i] = 0$  then
16         $n_{10}(j) \leftarrow n_{10}(j) + 1$ 
17      else if  $a[i, j] = 0 \wedge e[i] = 0$  then
18         $n_{00}(j) \leftarrow n_{00}(j) + 1$ 
19      end if
20    end for
21  end for
22  for  $j = 0$  to  $M$  do
23     $S[j] \leftarrow s(n_{11}(j), n_{10}(j), n_{01}(j), n_{00}(j))$ 
24  end for
25   $D \leftarrow \text{SORT}(S)$ 
26  return  $D$ 

```

As can be seen, both Tarantula and Ochiai fail to consider c_3 as one of the most suspicious components. Besides, c_2 and c_3 can be considered as multiple fault because all failed runs can be explained either by c_2 or c_3 (but the two by themselves are not a valid explanation for all failures). In the next section, we present our technique that exploits this info and contains multiple-fault explanations in its ranking.

3. Multiple-fault localization

In this section, we present our multiple-fault localization approach Zoltar-M, which is based on reasoning as performed in

Table 1
Observation matrix example A.

	c_1	c_2	c_3	c_4	c_5	c_6	c_7	e
	1	0	1	0	1	0	0	0
	0	0	0	0	1	0	0	0
	0	0	1	0	0	1	0	0
	1	1	0	0	1	0	0	0
	1	0	1	0	0	0	0	0
	1	0	0	0	1	0	0	0
	1	0	1	0	0	1	0	0
	1	1	0	0	1	1	0	1
	1	1	0	1	0	1	1	1
	1	0	1	0	0	0	0	1
	1	1	0	1	1	1	1	1
	1	1	0	0	0	1	0	1
	1	0	1	0	0	1	1	1
$n_{11}(j)$	6	4	2	2	2	5	3	
$n_{10}(j)$	5	1	4	0	4	2	0	
$n_{01}(j)$	0	2	4	4	4	1	3	
$n_{00}(j)$	2	6	3	7	3	5	7	
s_T	0.58	0.82	0.37	1	0.37	0.74	1	
s_0	0.74	0.73	0.33	0.58	0.33	0.77	0.71	

model-based diagnosis, combined with (Bayesian) probability theory to compute the ranking of the candidates. The major difference with the statistical approach in Section 2.2 is

- that only a *subset* of components is considered (the so-called hitting set) in contrast to all components,
- all computed candidates logically explain the observed failures, and
- that the ranking is based on probability, rather than statistical similarity.

In the remainder of this section, specific details on the two main phases of our Zoltar-M approach are given (see Fig. 2): (1) candidate generation and (2) candidate fault probability computation (ranking). In addition, as a by-product of Zoltar-M, we propose an optimal SFL variant for single faults.

3.1. Hitting set computation

In model-based diagnosis, one derives a model of the program that, together with the observations of input-output behavior, determines a set of constraints from which diagnostic solutions consistent with this behavior are logically deduced. Unlike the

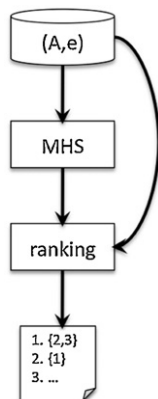


Fig. 2. Zoltar-M's main phases.

MBD approaches such as presented in Mayer and Stumptner (2007, 2008), in our Zoltar-M approach we refrain from modeling the program in detail, but use A as the only, dynamic source of information, from which we derive the model and the input-output observations.

Each component c_j is modeled by the logic proposition

$$h_j \Rightarrow in.ok_j \Rightarrow out.ok_j \quad (2)$$

where h_j models the health state of c_j (*true* = healthy, *false* = defect). This so-called *weak* model specifies that a component produces correct output values (*out.ok true*) if (1) healthy (h is true), and (2) when provided with correct input (*in.ok true*). Note that this model still allows a component to produce correct data (the probability of which is measured by g) even though $h = false$. Also note that a component can accept erroneous input data and still produce correct output. Finally, this approach allows the inclusion of additional component information as the number of propositions per component is not limited to the above, default model of nominal behavior.

Due to dynamic (data-dependent) control flow each run may involve different components. Consequently, rather than modeling the program by a static composition of component propositions (Eq. (2)), we consider a dynamic model that is defined per program run (i.e., A_i^*). In Abreu et al. (2008), it is shown that each failed run yields a conjunction of logical constraints (i.e., a sub-model) in terms of the components involved, which is known in MBD as a *conflict* (de Kleer et al., 1992). For instance, a failed run involving c_1 and c_2 generates the conflict $\neg h_1 \vee \neg h_2$, indicating that c_1 and c_2 cannot both be healthy.

The multiple-fault approach is based on compiling each failed run (row A_i^*) to a conflict, after which the diagnosis for A is derived by computing the hitting set (Reiter, April 1987) from all conflicts (Abreu et al., 2008) (the hitting set algorithm essentially transforms logic products-of-sums into sums-of-products). For instance, the example observation matrix A in Table 1 generates the following six conflicts

$$\begin{aligned}
 &(\neg h_1 \vee \neg h_2 \vee \neg h_5 \vee \neg h_6) \wedge \\
 &(\neg h_1 \vee \neg h_2 \vee \neg h_4 \vee \neg h_6 \vee \neg h_7) \wedge \\
 &(\neg h_1 \vee \neg h_3) \wedge \\
 &(\neg h_1 \vee \neg h_2 \vee \neg h_4 \vee \neg h_5 \vee \neg h_6 \vee \neg h_7) \wedge \\
 &(\neg h_1 \vee \neg h_2 \vee \neg h_6) \wedge \\
 &(\neg h_1 \vee \neg h_3 \vee \neg h_6 \vee \neg h_7)
 \end{aligned}$$

The (minimal) hitting set comprises one single-fault candidate $\{1\}$, and two double-fault candidates $\{2,3\}$ and $\{3,6\}$. Note that the triple-fault candidate $\{1,2,3\}$, which equals the actual fault state, is subsumed by both $\{1\}$ and $\{2,3\}$ and therefore does not appear in D (which is the *minimal* hitting set). The reason why, e.g., $\{1\}$ subsumes $\{1, j\}$, $j=2,3, \dots, M$ is that the weak component model (2) allows any faulty component j to exhibit correct behavior. Hence $\{1, j\}$ is also a valid explanation. The hitting set can be directly observed from A by (multi-)column “chains” of ‘1’s from top to bottom formed by all failing rows of A . Note that this procedure only considers failed runs. Passed runs are considered later on when computing the probability of each diagnostic candidate.

The above example illustrates that the true diagnosis (candidate) can be preceded (or subsumed) by other, more probable candidates. However, for $N \rightarrow \infty$ Zoltar-M produces the optimal result, where diagnoses such as $\{1\}$ and $\{2, 3\}$ eventually disappear, exposing the only correct diagnosis $\{1, 2, 3\}$. This can be seen through the following argument. Consider a C-fault program. While for small N the minimal hitting set will still contain many members (components) other than the C faulty components, by increasing N the probability that a non-faulty component will still be included steadily decreases (Abreu et al., 2008). Let f denote

Table 2
Diagnoses for example A.

Technique	$D = \{d_1(s pr), \dots, d_k(s pr)\}$
Tarantula	$\{\{4\}(1), \{7\}(1), \{2\}(0.82), \{6\}(0.74), \{1\}(0.58), \{3\}(0.37), \{5\}(0.37)\}$
Ochiai	$\{\{6\}(0.77), \{1\}(0.74), \{2\}(0.73), \{7\}(0.71), \{4\}(0.58), \{3\}(0.33), \{5\}(0.33)\}$
Zoltar-M	$\{\{1\}(0.98), \{2, 3\}(0.99e^{-2}), \{3, 6\}(0.52e^{-2})\}$

the probability of a run failing. As the probability that a run passes equals the probability that none of the C components cause a failure, which equals $(1 - r(1 - g))^C$, it follows that $f = 1 - (1 - r \cdot (1 - g))^C$ (Abreu et al., 2008). For the hitting set analysis, only failing runs matter. For $f \cdot N$ failing runs, the C -fault candidate is by definition within the set of candidates that “survive” those $f \cdot N$ runs (whose chain is still unbroken). However, the probability that other components can be involved in a candidate is less than unity, which forms the basis of those candidates’ eventual elimination.

For example, the probability of a B -cardinality diagnosis ($B < C$) competing with our C -cardinality solution equals the probability that at least one out of the B components is hit every time in a failing run. The latter probability equals $b = 1 - (1 - r)^B$ (derivation similar to f). Hence, for N runs, the probability of this competing diagnosis surviving in the final hitting set is of the order $b^{f \cdot N}$, which as $b < 1$ negative-exponentially decreases to zero for large N . Note that the above analysis does not consider a particular probability computation regarding the ranking. It simply proves that, for large N , the diagnosis can only consist of the single surviving C -cardinality candidate ($\{1, 2, 3\}$ in the earlier example). Our experiments confirm this optimality. There is one exception to the above argument. Components that are *always* executed (e.g., initialization code) will always appear as single-fault candidate, which is typically ranked higher than a genuine multiple-fault (see next section). As this applies to techniques based on spectral information this problem also occurs with statistical techniques.

3.2. Probability computation

For each multiple-fault candidate, the probability of being the actual diagnosis depends on the extent to which that candidate explains all observations (pass or fail per run). Let $\Pr(\{j\})$ denote the *a priori* probability that a component c_j is at fault. Although this value is typically dependent on code complexity, design, etc., we will simply assume $\Pr(\{j\}) = p$ (we arbitrarily set $p = 0.01$ in the context of this paper). Assuming components fail independently, and in absence of any observation, the prior probability a particular diagnosis d_k is correct is given by $\Pr(d_k) = p^{|d_k|} (1 - p)^{M - |d_k|}$. Similar to the incremental compilation of conflicts per run we compute the posterior probability for each candidate based on the pass/fail observation obs for each sequential run using Bayes’ update rule according to

$$\Pr(d_k | obs) = \frac{\Pr(obs | d_k) \Pr(d_k)}{\Pr(obs)}$$

The denominator $\Pr(obs)$ is a normalizing term that is identical for all d_k and thus needs not to be computed directly. $\Pr(obs | d_k)$ is defined as

$$\Pr(obs | d_k) = \begin{cases} 0 & \text{if } d_k \text{ and } obs \text{ are inconsistent} \\ 1 & \text{if } d_k \text{ logically follows from } obs \\ \epsilon & \text{if neither holds} \end{cases}$$

In the context of model-based diagnosis, many policies exist for ϵ (see de Kleer, 2007). In this paper, we define ϵ as follows

$$\epsilon = \begin{cases} g(d_k)^\eta & \text{if run passed} \\ 1 - g(d_k)^\eta & \text{if run failed} \end{cases}$$

In this equation, η is the number of faulty components involved in the run (the rationale being that the more faulty components are involved, the more likely it is that the run will fail (Abreu et al., 2008)), and g is estimated by

$$g(d_k) = \frac{n_{10}(d_k)}{n_{10}(d_k) + n_{11}(d_k)}$$

where $n_{1q}(d_k) = \sum_{i=1, \dots, N} \left[\left(\bigvee_{j \in d_k} a_{ij} = 1 \right) \wedge e_i = q \right]$ is a generalization of the definition in Section 2.2 to support multiple fault explanations, $q \in \{0, 1\}$, and $[\cdot]$ denotes Iverson’s operator (Iverson, 1962) ($[true] = 1$, $[false] = 0$).

To illustrate the differences between the probabilistic approach as presented in this section and the statistical SFL approach (as explained in Section 2.2), again consider the example A in Table 1. The diagnostic report D for the different approaches are listed in Table 2. As can be seen, the top ranked candidate for both Tarantula and Ochiai is not one of the three faulty locations, whereas for Zoltar-M one of the faults, namely c_1 would be immediately found. Furthermore, in contrast to Zoltar-M, which contains multiple faults explanations such as $\{2, 3\}$, Tarantula and Ochiai only rank single-fault explanations. To conclude, note that Zoltar-M *only* lists candidates that actually *explain* all observed failures.

While the inherent multiple-fault approach used in Zoltar-M is asymptotically optimal, the complexity of the underlying hitting set algorithm and subsequently having to manage a possibly exponential number of multiple-fault candidates (e.g., update their probability) is prohibitive for large C (and N, M). Nevertheless, preliminary experiments with a statistically directed search technique (i.e., using statistical similarity to guide the search) indicates that the complexity of our current hitting set computation can be reduced by several orders of magnitude. In addition, hitting set completeness can be traded-off to further reduce time complexity (see Feldman et al., 2008 for a greedy stochastic search addressing this issue).

3.3. Single-fault case

In this section, we show how our above reasoning approach can be used to derive an optimal similarity coefficient for *single-fault* programs.

In the single-fault case (such as the SIR benchmark set of programs), we know that all failures relate to only one fault, which, by definition, is included in the minimal hitting set. Hence, any coefficient approach should consider the minimal hitting set only (i.e., only those c_j which consistently occur in failing runs). This implies that the optimal approach is to select only the failing runs and compute the similarity coefficient. Since for these components by definition $n_{01} = 0$, one only needs to consider n_{11} and n_{10} . This, in turn, implies that the ranking is only determined by the exonerating term n_{10} . In summary, once we only consider the components

Table 3
The Siemens benchmark set.

Program	Faulty versions	# components (M)	# runs (N)	Description
print_tokens	7	539	4130	Lexical analyzer
print_tokens2	10	489	4115	Lexical analyzer
replace	32	507	5542	Pattern recognition
schedule	9	397	2650	Priority scheduler
schedule2	10	299	2710	Priority scheduler
tcas	41	174	1608	Altitude separation
tot_info	23	398	1052	Information measure

included in the hitting set, any of the coefficients that includes n_{10} in the denominator will produce the same, optimal ranking. Experiments using this “hitting set filter” combined with a simple similarity coefficient such as Tarantula indeed confirm that this approach leads to the best performance (Vayani, in press). For instance,

$$\text{filter} = \begin{cases} s_T & \text{if } n_{01} \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Note that the above filter is only optimal for programs that have only one fault as applying this filter to any multiple-fault program would be overly restrictive. It would fail to detect faults that are not always involved in failed runs. For example, the diagnosis for A in Table 1 when using the filtering approach would yield $D = \{\{1\}\}$, entirely ignoring two of the three faults. Hence, instead of considering a single-fault hitting set filter, we modify this approach in order to also allow application to multiple-fault programs. Taking the Ochiai coefficient as (best) starting point (for $\kappa = 1$, Eq. (3) follows from Eq. (1) by squaring, and factoring out $n_{11}(j)$, none of which changes the ranking) and applying the above filtering approach, we derive the following similarity coefficient (Abreu et al., 2008), coined Zoltar-S, according to

$$s_{Z-S} = \frac{n_{11}(j)}{n_{11}(j) + n_{10}(j) + n_{01}(j) + \kappa(n_{01}(j)n_{10}(j)/n_{11}(j))} \quad (3)$$

where $\kappa > 0$ is a constant factor that exonerates a component c_j that was either seldom executed in failed runs or often in passed runs. We empirically verified that the higher the value of κ , the more identical the diagnosis becomes with the one obtained by the hitting set filter (Vayani, in press). In the context of this paper, we limit κ to 10,000 to avoid round-off errors.

In Abreu et al. (2009b) it has been proven that, for single-fault programs, given the available input data (A, e), the diagnostic ranking produced by our reasoning technique is theoretically optimal. As such theoretical optimality applies to $C = 1$ (i.e., for single-fault programs), we have adapted the result in Abreu et al. (2009b) to a similarity coefficient in terms of κ . The following theorem proves that for $\kappa \rightarrow \infty$ the Zoltar-S similarity coefficient has exactly the same behavior.

Theorem. For single-fault programs, given the available input data (A, e), the diagnostic ranking produced by Zoltar-S is theoretically optimal.

Proof. Let c_f be the faulty component and c_p a (representative) non-faulty component. For single faults, $n_{01}(f) = 0$, $n_{11}(f) = N_F$ (where N_F is the number of failed runs), and $n_{01}(p) \geq 0$. Therefore, $n_{01}(f)n_{10}(f)/n_{11}(f) = 0$ and $n_{01}(p)n_{10}(p)/n_{11}(p) \geq 0$. Consequently, for non-faulty components, the following holds

$$\lim_{\kappa \rightarrow \infty} s_{Z-S}(p) \simeq 0 \quad (4)$$

$$\lim_{\kappa \rightarrow \infty} s_{Z-S}(f) = \frac{n_{11}(f)}{n_{11}(f) + n_{10}(f)} = \frac{N_F}{N_F + n_{10}(f)} > 0 \quad (5)$$

Hence, the higher κ is, the more is the exoneration factor for non-faulty, and as such, the faulty one will rank high in the diagnostic ranking. \square

To evaluate the diagnostic capabilities of Zoltar-S in comparison with other techniques, the Siemens benchmark set is used. This well-known benchmark is composed of seven programs (see Table 3; for detailed info, visit <http://sir.unl.edu>). In total, the Siemens benchmark set of programs provides 132 faulty programs. However, as no failures are observed in two of these programs, namely version 9 of `schedule2` and version 32 of `replace`, they are discarded. Besides, we also discard versions 4 and 6 of `print_tokens` because the faults are not in the program itself but in a header file. In summary, we discarded 4 versions out of 132 provided by the suite, using 128 versions in our experiments. To collect the program spectra, the `Zoltar` toolset (Janssen et al., 2009) was used. For compatibility with previous work in (single-) fault localization, we use the effort/score metric (Abreu et al., 2007; Renieris and Reiss, 2003), which is the percentage of statements that need to be inspected to find the fault – in other words, the rank position of the faulty statement divided by the total number of statements. Note that some techniques such as in Liu et al. (2005) and Renieris and Reiss, 2003) do not rank all statements in the code, and their rankings are therefore based on the program dependence graph (PDG) of the program.

Fig. 3 plots the percentage of located faults in terms of debugging effort (Abreu et al., 2007). Apart from the coefficients studied for SFL, the following techniques are also plotted: intersection and Union (Renieris and Reiss, 2003), Delta Debugging (DD) (Zeller, 2002), Nearest Neighbor (NN) (Renieris and Reiss, 2003), and Sober (Liu et al., 2005), which are among the best SFL techniques (detailed discussion in Section 6). As Sober is publicly available, we run it in our own environment. The values for the other techniques are, however, directly cited from their respective papers.

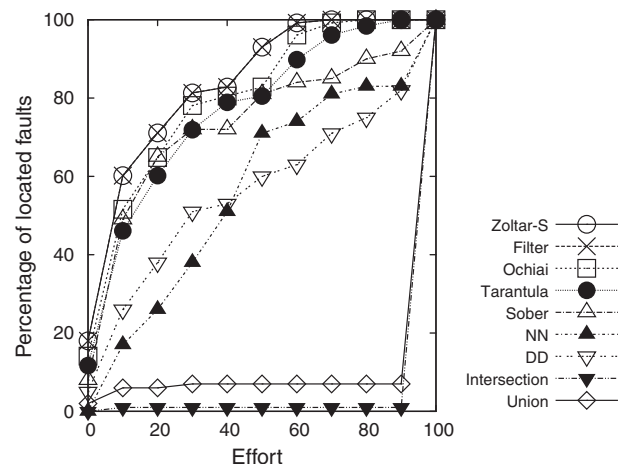


Fig. 3. Effectiveness Comparison ($C = 1$).

From Fig. 3, we conclude that Zoltar-S and the filter version are consistently the best performing techniques (note that in the single-fault context Zoltar-M simply reduces to Zoltar-S), finding 60% of the faults by examining less than 10% of the source code. For the same effort, using Ochiai would lead a developer to find 52% of the faulty versions and with Tarantula only 46% would be found. The Zoltar-S approach is followed by Ochiai, which outperforms Sober and Tarantula, which as concluded in Liu et al. (2005), yield similar performance. Finally, the other techniques plotted are clearly outperformed by the spectrum-based techniques.

4. Theoretical evaluation

In order to gain understanding of the effects of the various parameters on the diagnostic performance of the different approaches, we use a simple, probabilistic model of program behavior that is directly based on C , N , M , r , and g . Without loss of generality we model the first C of the M components to be at fault. For each run, every component has probability r to be involved in that run. If a selected component is faulty, the probability of exhibiting nominal (“good”) behavior equals g . When either of the C components fails, the run will fail. We study the performance of Zoltar-M in comparison to Tarantula, Ochiai, and Zoltar-S for observation matrices that are randomly generated according to the above model.

4.1. Performance metrics

Before evaluating the results, we first present our performance metric. As one of the motivations of our multiple-fault approach is the exposure of fault multiplicity (parallel debugging) we refrain from reusing established metrics such as the diagnostic quality (Abreu et al., 2007) or score (Renieris and Reiss, 2003) but evaluate the amount of *wasted debugging effort* W as a function of the number of parallel debuggers (Jones et al., 2007), denoted by P , which more clearly indicates practical debugging parallelism. The wasted debugging effort is computed as follows. From the diagnosis (obtained with either a statistical or reasoning approach), the first P candidates are examined (debugged) in parallel (Jones et al., 2007). Actual faults are assumed to be properly debugged, after which the program is retested. Based on the retest a new diagnosis is obtained (excluding the repaired components, but including the still uncovered faults that may have considerably moved up in the ranking). This P -parallel process continues until in the last iteration the program retests ok (i.e., all faults have been found). W measures the percentage of non-faulty components that were debugged in the above process. For $P=1$, the above procedure reduces to a standard sequential debugging process. For instance, consider the diagnostic reports yielded by Tarantula and Zoltar-M (as in Table 2) for Example A in Table 1. Table 4 shows the performance profile for these two techniques (I stands for the number of bugs found in the

Table 4
Wasted effort for different developers P .

P	1	2	3	4	5	6	...
Tarantula W/I	14/0	29/0	29/1	43/1	43/2	43/3	...
Zoltar-M W/I	0/1	0/2	0/3	14/3	–	–	...

first debugging iteration). As can be seen, Tarantula would need more developers in order to get a bug-free program in one iteration (six developers against three for Zoltar-M). Furthermore, for this example, the wasted effort is consistently higher for Tarantula: with Zoltar-M, three developers would eliminate all bugs from the program at the cost of 0% wasted effort, whereas with Tarantula 6 developers would be needed at a cost of 43%. Note that there is no point in putting more than four developers to work as the Zoltar-M diagnosis contains only four different components.

Another reason not to adopt the aforementioned score metric (Renieris and Reiss, 2003) is that in our synthetic model we do not have program dependence graph information. Furthermore, the choice to exclude the actual faults from the debugging effort (i.e., instead of counting them as effort) is to make our performance metric independent of the number of faults C .

4.2. Experimental results

In our first experiment, we focus on the effect of C , N , M , r , and g on W . Consequently, we choose $P=1$. We have varied M between 10 and 30 and after verifying that this does not change our conclusions (Abreu et al., 2008), we choose $M=20$ for the plots in the paper. Similarly, we also varied r between $r=0.4$ and $r=0.6$, and as there are no significant differences we only include the plots for $r=0.6$, which is roughly the same as the values measured for the Siemens set.

Figs. 4–6 plot W versus N for $C=1$, $C=2$, and $C=5$, respectively. We have also applied the technique for matrices with $C=8$ and the conclusions are essentially the same. Each measurement represents an average over 1000 sample matrices. The plots show that for small N all techniques start with equal W (for $N=1$ it follows that $W=(M-C) \cdot r/M$ (Abreu et al., 2008)), while for sufficiently large N all techniques produce an optimal diagnosis. The plots clearly show that all techniques yield an optimal diagnosis for sufficiently large N . This happens earlier for small C and g . In the single-fault case, there is hardly any difference in the various techniques. For a small value of g , almost each run that involves the faulty component yields a failure, already producing near-perfect diagnoses for only small N . For a large value of g (which is more realistic, the Siemens set exhibits g values ranging from 79% (`tot.info`) to 99% (`tcas`)), the fraction of failing runs dramatically decreases (cf. f in Section 3.1). Consequently, a much larger number of runs is required to obtain a good diagnosis. For $C=5$, we see the same trend, albeit that convergence to good diagnosis is much slower,

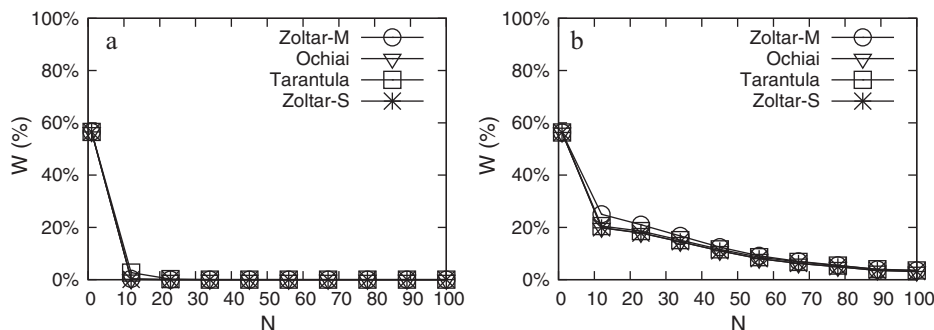


Fig. 4. Wasted effort W for $C=1$.

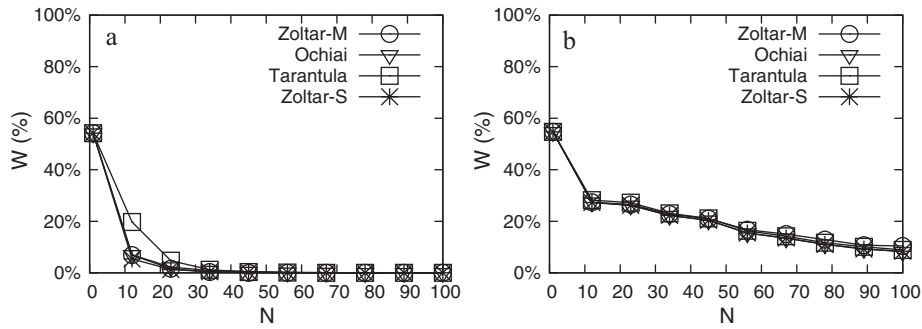


Fig. 5. Wasted effort W for $C=2$.

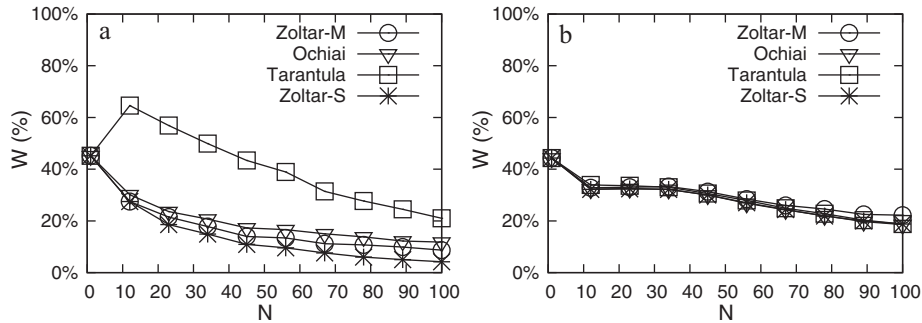


Fig. 6. Wasted effort W for $C=5$.

especially for high g . This is due to the (combinatorial) fact that the number of “competitor” candidates of cardinality $B \leq C$ (see Section 3.1) greatly increases with C (and M (Abreu et al., 2008)). The main conclusion is that all techniques are very similar for the synthetic matrices and no technique clearly outperforms the others. For $C=5$, we conclude that for small g Zoltar-S outperforms all other techniques, Zoltar-M is the second-best technique, and Tarantula is the worst performing technique. Besides, for small g and N , the Tarantula technique has very poor performance because some of the non-faulty components are only touched in failed runs, hence sharing the first position of the ranking and degrading the diagnosis. This is also the reason why the wasted effort first increases and only then starts to decrease (more passed runs are needed for non-faulty components to be exonerated). The fact that Zoltar-S outperforms Zoltar-M comes from the fact that for the synthetic matrices there are not that many non-faulty components involved in all failed runs, and therefore Zoltar-S manages to rank the faulty components on top. For more realistic cases ($g=0.9$), all techniques perform equally (poor), and much higher N is required to produce high diagnostic quality.

In Figs. 7–9, we measure W for all approaches as function of P to study inherent debugging parallelism for $C=\{1, 2, 5\}$. For these

plots, we set $N=500$ to ensure that each technique has reached acceptable diagnostic quality. For $g=0.1$, W starts a linear increase after $P=C$ (the “knee”), which indicates that all C faults are indeed at or near the top of the ranking (the bump at $P=4$ is due to integer division effects). Except for Ochiai, both Zoltar-S and Tarantula yield similar performance as Zoltar-M. For $C=1$ and $g=0.1$, Zoltar-M has zero wasted effort throughout. This occurs because, for $N=500$, the diagnosis only contains the faulty statement (perfect diagnosis), revealing that there is no point in having more than one developer debugging the program.

While the above results show to what extent debugging can be efficiently parallelized, in practice information on C is, of course, not available. In the following, we evaluate the added value of multiple-fault diagnosis in estimating the number of debuggers P that can be efficiently deployed in parallel. The plots in Fig. 10 show the distribution of the probability (Zoltar-M) or similarity (Zoltar-S, Ochiai, Tarantula) versus the ranking position. For multiple-fault diagnoses, each member index is counted as separate position. For cases where the diagnoses are near-perfect ($g=0.1$), the Zoltar-M distribution clearly exhibits the added information on the program’s fault cardinality C (corresponding to the “knee” in the previous plots), whereas the statistical techniques fail to produce any infor-

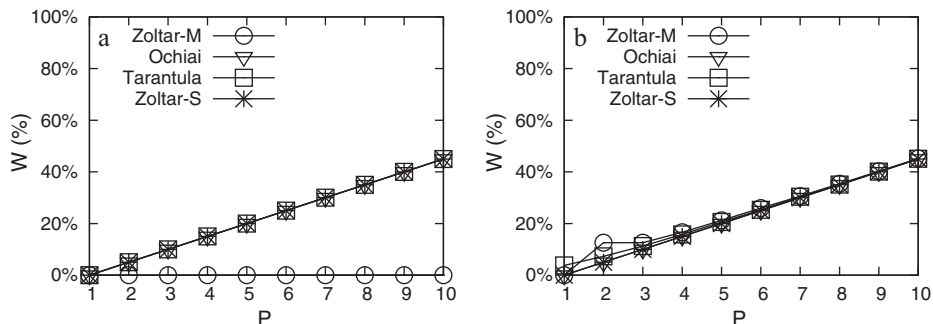


Fig. 7. W vs. P for $C=1$.

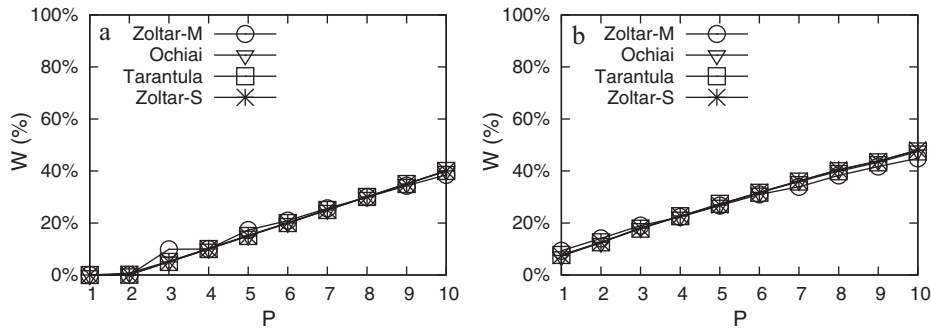


Fig. 8. W vs. P for C=2.

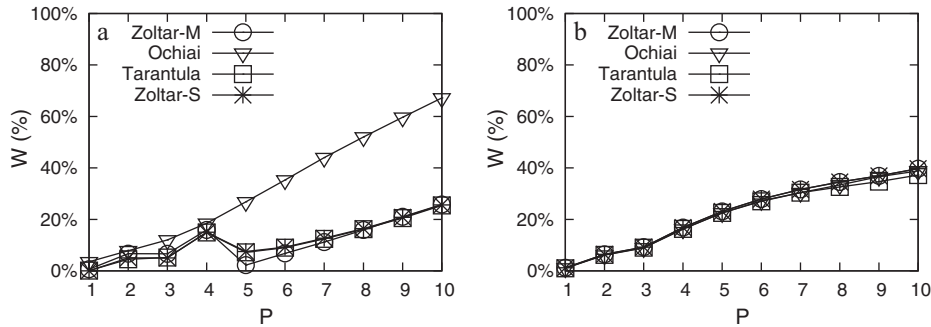


Fig. 9. W vs. P for C=5.

mation on C whatsoever (although the Zoltar-S ranking distribution has more dynamics). For a high value of g , this relative advantage becomes less as diagnostic quality degrades. Note that this can be remedied by further increasing N .

5. Empirical evaluation

Whereas the synthetic observation matrices used in the previous section are populated using a uniform distribution, this is not the case with observation matrices for the behavior of actual programs (different *spectral distribution*). Therefore, in this section we will evaluate the same diagnosis techniques on the SIR-S set, which provides the programs introduced in Section 3.3 extended with the real-world, large programs *space*, *gzip*, and *sed* (see Table 5). In addition, we also evaluated our techniques in the extended benchmark set of programs to accommodate multiple-fault (SIR-M).

5.1. Experimental setup

The SIR-M set extends the Siemens-S set with program versions that combine several faults from the latter set. The faults can be selectively activated via conditional compilation. The selections of faults that are available in the SIR-M set are limited by

- (1) their nature (e.g., a fault in non-executable code, which is not handled by our techniques would effectively reduce a C -fault diagnosis to a $(C-1)$ -fault diagnosis),

- (2) the number of failed runs (we only consider faults that yield at least one failed test case),
- (3) their locations (several faults in SIR-S have the same statement location), and
- (4) the number of lines of code involved (we only consider faults that can be attributed to a single line).

As mentioned in Section 3.3, the Zoltar toolset (Janssen et al., 2009) is used to obtain code coverage information for each of the test cases supplied with the programs in the benchmark set of programs. The error vector in the last column is constructed by comparing the output of a faulty version of a program with that of the correct version of the program, on a given test case.

For the resulting set of program spectra, Zoltar supports various diagnosis techniques, including Zoltar-M, and the Tarantula, Ochiai, and Zoltar-S coefficients. In the case of Zoltar-M, the presence of duplicate columns, following from the block structure of a program, is exploited in the hitting-set calculation by grouping all identical columns, while maintaining the set of components (lines of code) that they correspond to. This way, larger numbers of components can be handled than in the case of synthetic observation matrices.

5.2. Experimental results

In Fig. 11, we show W versus P for *tcas* and *replace*, two representative programs, when seeded with $C=1$, $C=2$, and $C=3$ faults, respectively. Although the minimal hitting set computation is known to be rather expensive, we have used a low-cost, approximate technique dubbed STACCATO (Abreu and van Gemund, 2009), which makes Zoltar scale to large, real-world programs (Abreu et al., 2009b). We have also repeated the experiment up to $C=5$ (up to $C=10$ for *tcas*), but the graphs are similar to those for $C=2$ and $C=3$ of the representative programs, with the performance of Zoltar-S approaching that of the other methods as C increases.

Table 5
SIR's real-world programs.

Program	Faulty versions	M	N	Description
<i>space</i>	38	9564	150	adl interpreter
<i>gzip-1.3</i>	7	5680	210	Data compression
<i>sed-4.1.5</i>	6	14,427	370	Textual manipulator

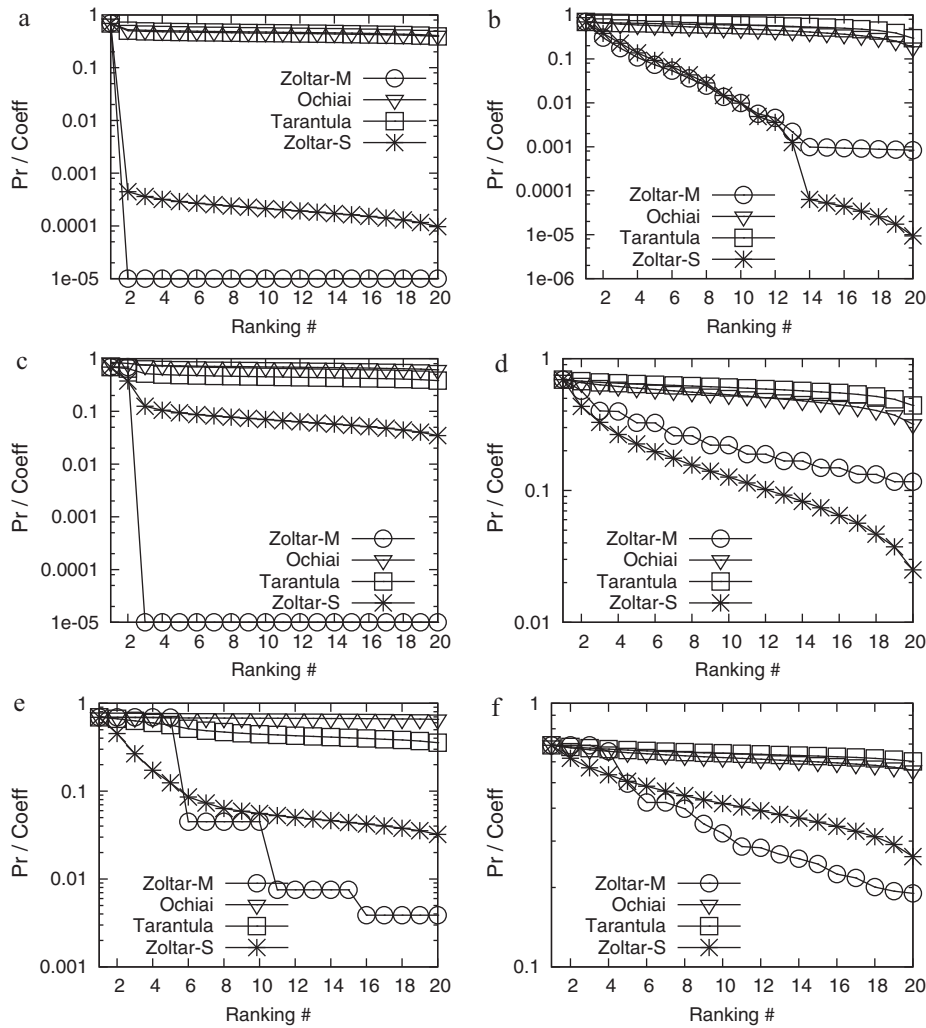


Fig. 10. Probability/similarity distribution.

The plotted W values are averaged over several different program versions: in case of the plots for $C=1$, these are all faulty versions in the SIR-S set that can be attributed to a single line of executable code (30 for `tcas`, and 25 for `replace`). In the case of the $C=2$ and $C=3$ plots, these are 40–100 randomly selected combinations of faults. The plateau reached by Zoltar-M for `tcas` at high P values is caused by the limited size (ambiguity) of Zoltar-M's diagnosis, removing the need to have them inspected by additional developers.

5.3. Evaluation

Fig. 11(a) and (d) confirms the observation of Section 4 that for single faults, Zoltar-S is optimal. Although at the end of Section 3.3 we noted that in a single-fault context, Zoltar-M reduces to Zoltar-S, here Zoltar-M runs in multiple-fault mode, and the presence of cardinality $C=3$ diagnoses in the hitting set, as explained above, is the cause for the small differences between these two techniques.

Contrary to what we observed in Section 4.2, where Zoltar-S is among the best performing methods for multiple-faults, in Fig. 11(b), (c), (e) and (f), Zoltar-S performs worst. This is caused by many non-faulty components that are active in all failed runs. Having $n_{01}(j)=0$ in Eq. (3), such components fail to be exonerated via the term $\kappa(n_{01}(j)n_{10}(j)/n_{11}(j))$, and will therefore rank high, leading to a lower quality diagnosis. While in the synthetic observation matrices it is unlikely that a component is active in all failed runs,

this is quite common in software (e.g., statements that are always executed).

As shown in Fig. 11(b), (c), (e) and (f), for $C=2$ and $C=3$, Zoltar-M generally outperforms the statistical techniques, but for `tcas`, its performance is quite close to that of the SFL approaches using the Ochiai and Tarantula coefficients. This can be attributed to the following two related effects. First, the `tcas` faults that are available for making multiple-fault versions have a higher goodness factor ($g=0.95$) than those available for `replace` ($g=0.86$), making the diagnosis problems for the multiple-fault `tcas` versions inherently more difficult. The rationale is that for faults whose observation matrix inherently does not permit a good diagnosis (e.g., because the activity of a non-faulty component accidentally coincides with the occurrence of failures), all appropriate techniques will yield an equally bad diagnosis on average. Referring back to the discussion at the end of Section 4.2, the high values for g (common to all programs in the used benchmark set of programs) also explain why on average, no technique achieves optimal diagnostic quality on the Siemens-S and Siemens-M faults, and the consequent absence of a “knee” in the P - W graph of Zoltar-M.

The second effect that contributes to the difference in the plots for `tcas` and `replace` is that the variations in control flow in the former program are extremely limited, while essentially, this is what the diagnosis methods are based on. As an illustration, for the correct version of `tcas`, the observation matrix that follows from the 1608 test cases that accompany the program contains many dupli-

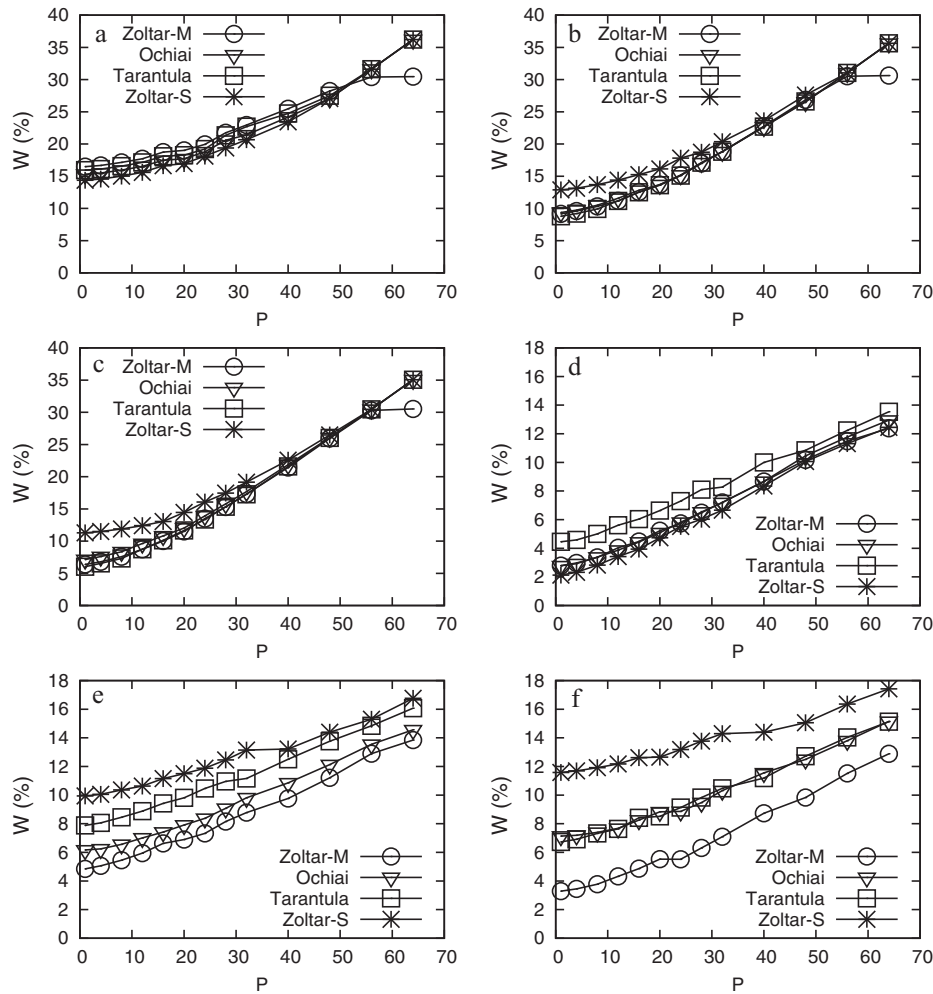


Fig. 11. Wasted effort for 1–64 developers on representative programs of the SIR-M benchmark set.

cate rows and columns: the number of unique rows (spectra) and columns (component behavior profiles) are 8 and 14, respectively. In comparison, the 5542 test cases of `replace` lead to 2023 different spectra, and 91 different behavior profiles, providing much more information to base the diagnosis on.

The latter observation confirms our expectation that the effectiveness of automated diagnosis techniques generally improves with program size. As an illustration, near-zero wasted effort is implied by the experiments with SFL on a 0.5 MLOC industrial software product reported in Zoetewij et al. (2007). In summary, `tcas` is too simple, and Fig. 11(e) and (f) can be expected to be the more representative of multiple-fault debugging in a realistic development environment. From this, we conclude that Zoltar-M can be expected to yield a significant improvement of debugging efficiency over the statistical methods in the multiple-fault case.

5.4. Time/space complexity

In this section, we report on the time/space complexity of Zoltar-M, compared to other fault localization techniques. We measure the time efficiency by conducting our experiments on a 2.3 GHz Intel Pentium-6 PC with 4 GB of memory. As most fault localization techniques have been evaluated in the context of single faults, in order to allow us to compare our fault localization approach to related work we limit ourselves to the original, single-fault Siemens benchmark set. We obtained timings for probabilistic dependence

graph (PPDG) and Delta Debugging (DD) from published results (Baah et al., 2008; Zeller, 2002).

Table 6 summarizes the results of the study. The columns show the programs, the average CPU time (in seconds) of Zoltar-M, traditional SFL (Tarantula/Ochiai), PPDG, and DD, respectively. As expected, the less expensive techniques are the statistics-based techniques Tarantula and Ochiai. At the other extreme are PPDG and DD. Zoltar-M costs much less than PPDG and DD.

With respect to space complexity, statistical techniques need to store the counters (n_{11} , n_{10} , n_{01} , n_{00}) for the similarity computation for all M components. Hence, the space complexity is $O(M)$. Zoltar-M also stores similar counters but per diagnosis candidate.

Table 6
Diagnosis cost for the single-fault subject programs (time in s).

Program	Zoltar-M	Tarantula/Ochiai	PPDG	DD
<code>print_tokens</code>	4.2	0.37	846.7	2590.1
<code>print_tokens2</code>	4.7	0.38	243.7	6556.5
<code>replace</code>	6.2	0.51	335.4	3588.9
<code>schedule</code>	2.5	0.24	77.3	1909.3
<code>schedule2</code>	2.5	0.25	199.5	7741.2
<code>tcas</code>	1.4	0.09	1.7	184.8
<code>tot_info</code>	1.2	0.08	97.7	521.4
<code>space</code>	7.4	0.15	N/A	N/A
<code>gzip</code>	6.2	0.19	N/A	N/A
<code>sed</code>	9.7	0.36	N/A	N/A

Assuming that $|D|$ scales with M , these approaches have $O(M)$ space complexity.

5.5. Threats to validity

Although the empirical study presented in this section provides evidence of the potential usefulness of the simultaneous bug fixing technique, there are threats to the validity of the empirical results that should be taken into account when interpreting the results.

Using only small to medium-sized C programs is a threat to external validity. Although, we believe the results will be identical, we cannot claim that the results generalize (large-sized programs, other programming languages). Yet another threat to external validity is the way multiple fault versions are built. When combining faults we assume an or-model (cf. the ϵ -policy). So, we ignore interference between faults (faults can mask other faults).

Quantifying the debugging effort is extremely difficult because developers can recognize some components do not need to be inspected. Besides, we assume developers will inspect components following the ranking given by the techniques and that may not be entirely true (a developer could try to follow a “smell” following the control-data relationship). Finally, we also assume that a developer is able to identify the faulty component once it inspects it.

Deploying several developers to fix the multiple bugs in a system may be more error-prone than a single developer fixing all bugs. The experiments reported do not address this problem. Further studies are needed to investigate this issue.

6. Related work

As mentioned in the introduction, automated debugging techniques can be distinguished into statistical and logic reasoning approaches that use program models.

In model-based reasoning to automatic software debugging (MBSD), the program model is typically generated from the source code using static analysis, as opposed to the traditional application of model based diagnosis where the model is obtained from a formal specification of the (physical) system (Reiter, April 1987). An overview of different techniques to generate program models is given in Mayer and Stumptner (2008). The authors conclude that models generated by means of abstract interpretation (Abreu et al., 2009a) are the most accurate for debugging, while not suffering from the computational complexity inherent to more precise analysis techniques (Mayer and Stumptner, 2008). Recently, model-based techniques have also been proposed to isolate specific faults stemming from incorrect implementation of high-level conceptual models (Yilmaz and Williams, 2007), where mutations are applied to state machine models to detect conceptual errors (see Fig. 12), such as incorrect control flow and missing or additional features found in the implementation. Model-based approaches also include the work of Wotawa, Stumptner, and Mayer (Wotawa et al., 2002). Other approaches that fit into this category include explain (Groce et al., 2006) and Δ -slicing (Groce et al., 2006), which are based on comparing execution traces of correct and failed runs using model checkers. Model-based test generation (Esser and Struss, 2007) from abstract specifications of systems employs a similar idea where possible faults manifested as differences in abstract state machines are analyzed to generate tests. Although model-based diagnosis inherently considers multiple-faults, thus far the above software debugging approaches only consider single faults. Apart from this, our approach differs in the fact that we use program spectra as dynamic information on component activity, which allows us to exploit execution behavior, unlike static approaches. Furthermore, our approach does not rely on the approximations required by static techniques (i.e., incompleteness).

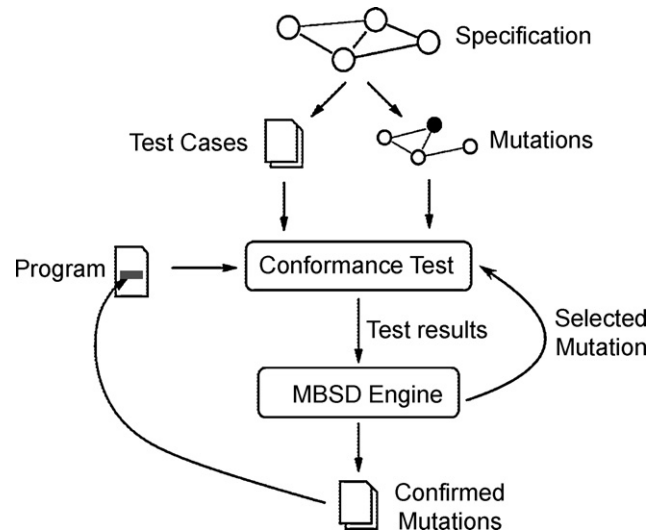


Fig. 12. Using conceptual models to enhance model-based reasoning.

Statistical approaches are very attractive from complexity-point of view. Well-known examples are the Tarantula tool (Jones et al., 2002), the Nearest Neighbor technique (Renieris and Reiss, 2003), the Sober tool (Liu et al., 2005), CP (Zhang et al., 2009), Holmes (Chilimbi et al., 2009), and the Ochiai coefficient (Abreu et al., 2007). Although differing in the way they derive the statistical fault ranking, all techniques are based on measuring program spectra. Examples of other techniques that do not require extra knowledge of the program under analysis are the Delta Debugging technique (Zeller, 2002) and the dynamic program slicing technique (Gupta et al., 2005).

Essentially all of the above work has mainly been studied in the context of single faults, except for recent work by Jones et al. (2007) and Zheng et al. (2006) which are motivated by the obvious advantages of parallel debugging with respect to development time reduction (particularly the work in Jones et al., 2007). They use clustering techniques to identify traces (rows in A) that refer to the same fault, after which a single-fault technique is applied to each cluster of rows. While our work has the same motivation, our approach is based on logic reasoning instead of clustering. Although both introduce an increase of computational complexity, compared to the aforementioned statistical approaches, our hitting set analysis approach is asymptotically optimal, while in the clustering approach there is a possibility that multiple developers will still be effectively fixing the same bug. As their parallel debugging approach has only been evaluated in a restricted empirical context, our results, e.g., for the Siemens programs, cannot yet be compared.

7. Conclusions and future work

In this paper, we have presented a multiple-fault localization technique, Zoltar-M, which is based on the dynamic, spectrum-based measurement approach from statistical fault localization methods, combined with a logic (and probabilistic) reasoning approach from model-based diagnosis, inspired by previous work in both separate disciplines (Abreu et al., 2007; Feldman et al., 2008). We have compared the performance of Zoltar-M with Tarantula and Ochiai, which are among the best known statistical SFL approaches, as well as a new statistical SFL technique, coined Zoltar-S, derived by us as a by-product of our reasoning approach, and shown to be optimal for single-fault programs ($C = 1$).

Our synthetic experiments show that both the reasoning and statistical approaches have the same general properties with respect to the influence of the parameters we introduced, viz, num-

ber of components M , number of test cases N , testing code coverage r , testing fault coverage g , and fault cardinality C . For a low value of g , both approaches yield near-perfect quality for relatively small N , while for high g (typical for many components in practice) a much larger N is required for good diagnosis. In most cases, it is Zoltar-S that outperforms Zoltar-M, which for $C > 1$ is due to the fact that all components are involved in different runs with the same probability, making it easy for Zoltar-S to pinpoint the faulty ones. Despite these small differences, Zoltar-M's ranking probability distribution clearly provides information on the program's potential debugging parallelism while statistical techniques fail to provide any information.

Our results on two multiple-fault programs of our newly created SIR-M benchmark suggest that for programs with small spectral distribution variability (and high g value) both approaches do not significantly differ. For the larger program, much more test information is available (N), the g parameter is somewhat lower, and the spectral distribution is highly non-uniform. In this case (for $C > 1$), Zoltar-M clearly outperforms all statistical approaches. The disparity with the synthetic results is due to the particular spectral distribution properties of real programs (such as components being executed in all failed runs). Aimed at providing a first-order understanding of the impact of some of the main parameters on diagnostic performance, our simple, probabilistic program model is still far from being able to accurately account for real program behavior.

Although both the reasoning and statistical approach are based on the same (spectral) information, our reasoning approach generally produces improved diagnostic information, in terms of debugging effort and/or (most notably) potential debugging parallelism. Nevertheless our results also indicate that even in the multiple-fault case statistical approaches are by no means outclassed by our reasoning approach, a result that was not initially anticipated. Given the higher complexity of the reasoning approach, there may be situations where application of a statistical technique such as Ochiai or Zoltar-S may be preferred over Zoltar-M. In this respect, we believe the result may be relevant in the context of the multiple-fault / parallel debugging work by Jones et al. (2007). Provided their clustering approach produces spectral partitions that apply to a single fault (or a very low fault multiplicity), our results would suggest the use of Zoltar-S, rather than Tarantula.

Acknowledgments

We extend our gratitude to Johan de Kleer for discussions which have influenced our multiple-fault reasoning approach. Also thanks to Rafi Vayani for conducting initial experiments on the effect of the hitting set filter in the single-fault case. Finally, we acknowledge the feedback from the discussions with our TRADER project partners.

References

- Abreu, R., van Gemund, A.J.C., 2009. A low-cost approximate minimal hitting set algorithm and its application to model-based diagnosis. In: Bulitko, V., Beck, J.C. (Eds.), Proceedings of the 8th Symposium on Abstraction, Reformulation and Approximation (SARA'09). Lake Arrowhead, CA, USA, 8–10 July 2009. AAAI Press.
- Abreu, R., Zoetewij, P., van Gemund, A., 2009a. Localizing software faults simultaneously. In: Choi, B. (Ed.), 9th International Conference on Quality of Software (QSIC'09). IEEE Computer Society, August 2009.
- Abreu, R., Zoetewij, P., van Gemund, A.J.C., 2007. On the accuracy of spectrum-based fault localization. In: McMinn, P. (Ed.), Proceedings of the Testing: Academia and Industry Conference – Practice And Research Techniques (TAIC PART'07). Windsor, United Kingdom, September 2007. IEEE Computer Society, pp. 89–98.
- Abreu, R., Zoetewij, P., van Gemund, A.J.C., 2008. An observation-based model for fault localization. In: Liblit, B., Rountev, A. (Eds.), Proceedings of the 6th Workshop on Dynamic Analysis (WODA'08), July 2008. ACM Press, pp. 64–70.
- Abreu, R., Zoetewij, P., van Gemund, A.J.C., 2009b. Spectrum-based multiple fault localization. In: Taentzer, G., Heimdahl, M. (Eds.), Proceedings of the IEEE/ACM International Conference on Automated Software Engineering (ASE'09). Auckland, New Zealand, 16–20 November 2009. IEEE Computer Society.
- Avizienis, A., Laprie, J.-C., Randell, B., Landwehr, C.E., 2004. Basic concepts and taxonomy of dependable and secure computing. IEEE Trans. Dependable Sec. Comput. 1 (1), 11–33.
- Baah, G.K., Podgurski, A., Harrold, M.J., 2008. The probabilistic program dependence graph and its application to fault diagnosis. In: Proceedings of International Symposium on Software Testing and Analysis (ISSTA'08).
- Chilimbi, T.M., Liblit, B., Mehra, K.K., Nori, A.V., Vaswani, K., 2009. Holmes: effective statistical debugging via efficient path profiling. In: Proceedings of the 31st International Conference on Software Engineering (ICSE'09), Vancouver, Canada, 16–24 May 2009. IEEE CS, pp. 34–44.
- de Kleer, J., 2007. Diagnosing intermittent faults. In: Biswas, G., Koutsoukos, X., Abdelwahed, S. (Eds.), Proceedings of the 18th International Workshop on Principles of Diagnosis (DX'07). Nashville, TN, USA, 29–31 May 2007, pp. 45–51.
- de Kleer, J., Mackworth, A.K., Reiter, R., 1992. Characterizing diagnoses and systems. Artif. Intel. 56, 197–222.
- de Kleer, J., Williams, B.C., 1987. Diagnosing multiple faults. Artif. Intel. 32 (1), 97–130.
- Esser, M., Struss, P., 2007. Automated test generation from models based on functional software specifications. In: Veloso, M.M. (Ed.), Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI'07). Hyderabad, India, 6–12 January 2007. AAAI Press, pp. 2255–2268.
- Feldman, A., Provan, G., van Gemund, A.J.C., 2008. Computing minimal diagnoses by greedy stochastic search. In: Fox, D., Gomes, C.P. (Eds.), Proceedings of the 23rd National Conference on Artificial Intelligence (AAAI'08). Chicago, IL, USA, 13–17 July 2008. AAAI Press, pp. 919–924.
- Groce, A., Chaki, S., Kroening, D., Strichman, O., 2006. Error explanation with distance metrics. Int. J. Software Tools Technol. Transfer (STTT) 8 (3), 229–247.
- Gupta, N., He, H., Zhang, X., Gupta, R., 2005. Locating faulty code using failure-inducing chops. In: Redmiles, D.F., Ellman, T., Zisman, A. (Eds.), Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering (ASE'05). Long Beach, CA, USA, 7–11 November 2005. IEEE Computer Society, pp. 263–272.
- Harrold, M., Rothermel, G., Wu, R., Yi, L., 1998. An Empirical Investigation of Program Spectra, vol. 33. ACM Press.
- Iverson, K.E., 1962. A Programming Language. John Wiley & Sons, New York, NY, USA.
- Janssen, T., Abreu, R., van Gemund, A.J.C., 2009. Zoltar: a toolset for automatic fault localization. In: van der Hoek, A., Menzies, T. (Eds.), Proceedings of the IEEE/ACM International Conference on Automated Software Engineering (ASE'09) – Tools Track. Auckland, New Zealand, 16–20 November 2009. IEEE Computer Society.
- Jones, J.A., Harrold, M.J., Bowring, J.F., 2007. Debugging in parallel. In: Rosenblum, D.S., Elbaum, S.G. (Eds.), Proceedings of the ACM/SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'07). London, UK, 9–12 July 2007. ACM Press, pp. 16–26.
- Jones, J.A., Harrold, M.J., Stasko, J.T., 2002. Visualization of test information to assist fault localization. In: Young, M., Magee, J. (Eds.), Proceedings of the 24th International Conference on Software Engineering (ICSE'02). Orlando, FL, USA, 19–25 May 2002. ACM Press, pp. 467–477.
- Liu, C., Yan, X., Fei, L., Han, J., Midkiff, S.P., 2005. SOBER: statistical model-based bug localization. In: Wermelinger, M., Gall, H. (Eds.), Proceedings of the 10th European Software Engineering Conference held jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/SIGSOFT FSE). Lisbon, Portugal, 5–9 September 2005. ACM Press, pp. 286–295.
- Mayer, W., Stumptner, M., 2007. Abstract interpretation of programs for model-based debugging. In: Veloso, M.M. (Ed.), Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI'07). Hyderabad, India, 6–12 January 2007. AAAI Press.
- Mayer, W., Stumptner, M., 2008. Evaluating models for model-based debugging. In: Ireland, A., Visser, W. (Eds.), Proceedings of the 23rd IEEE/ACM International Conference on Automated Software Engineering (ASE'08). L'Aquila, Italy, 15–19 September 2008. ACM Press, pp. 128–137.
- Pietersma, J., van Gemund, A.J.C., 2006. Temporal versus spatial observability in model-based diagnosis. In: Lin, C.-T. (Ed.), Proceedings of 2006 IEEE International Conference on Systems, Man, and Cybernetics (SMC'06). Taipei, Taiwan, 8–11 October 2006. IEEE Computer Society, pp. 5325–5331.
- Reiter, R., April 1987. A theory of diagnosis from first principles. Artif. Intel. 32 (1), 57–95.
- Renieris, M., Reiss, S.P., 2003. Fault localization with nearest neighbor queries. In: Grundy, J., Penix, J. (Eds.), Proceedings of the 18th IEEE International Conference on Automated Software Engineering (ASE'03). Montreal, Canada, 6–10 October 2003. IEEE Computer Society, pp. 30–39.
- R. Vayani, July 2007. Improving automatic software fault localization. Master's thesis, Delft University of Technology.
- Wotawa, F., Stumptner, M., Mayer, W., 2002. Model-based debugging or how to diagnose programs automatically. In: Hendtlass, T., Ali, M. (Eds.), Proceedings of IEA/AIE 2002, volume 2358 of LNCS. Cairns, Australia, 17–20 June 2002. Springer-Verlag, pp. 746–757.
- Yilmaz, C., Williams, C., 2007. An automated model-based debugging approach. In: Stirewalt, R.E.K., Egyed, A., Fischer, B. (Eds.), Proceedings of the 22nd IEEE/ACM

- International Conference on Automated Software Engineering (ASE'07). Atlanta, GA, USA, 5–9 November 2007. ACM Press, pp. 174–183.
- Zeller, A., 2002. Isolating cause-effect chains from computer programs. In: Proceedings of the 10th ACM SIGSOFT Symposium on Foundations of Software Engineering (FSE'02), Charleston, SC, USA, 10–12 November 2002. ACM Press, pp. 1–10.
- Zhang, Z., Chan, W.K., Tse, T.H., Jiang, B., Wang, X., 2009. Capturing propagation of infected program states. In: Proceedings of the 7th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on The foundations of software engineering (ESEC/FSE'09), Amsterdam, The Netherlands. ACM, pp. 43–52.
- Zheng, A.X., Jordan, M.I., Liblit, B., Naik, M., Aiken, A., 2006. Statistical debugging: simultaneous identification of multiple bugs. In: Proceedings of International Conference on Machine Learning (ICML'06), Pittsburgh, PA, USA, 25–18 June 2006. ACM Press.
- Zoetewij, P., Abreu, R., Golsteijn, R., van Gemund, A.J.C., 2007. Diagnosis of embedded software using program spectra. In: Leaney, J., Rozenblit, J., Peng, J. (Eds.), Proceedings 14th International Conference on the Engineering of Computer Based Systems (ECBS'07). IEEE Computer Society, 2007, pp. 213–218.
- Rui Abreu** is with the Department of Informatics of the Faculty of Engineering of University of Porto as an Assistant Professor. He obtained his PhD. in Computer Science at the Software Engineering Research Group at Delft University of Technology. He holds an MSc. in Computer Science and Systems Engineering from Minho University, Portugal. Through his thesis work at Siemens R&D Porto, and professional internship at Philips Research, he acquired industrial experience in the area of quality of (embedded) systems.
- Peter Zoetewij** works at IntelliMagic as a Software Developer. He holds an MSc. from Delft University of Technology, and a PhD. from the University of Amsterdam, both in computer science. Before his PhD., Peter worked for several years as a software engineer for Logica (now LogicaCMG), mainly on software for the oil industry.
- Arjan J.C. van Gemund** holds a BSc. in physics, and an MSc. (cum laude) and PhD. (cum laude) in computer science, all from Delft University of Technology. He has held positions at DSM and TNO, and currently serves as a full professor at the Electrical Engineering, Mathematics, and Computer Science Faculty of Delft University of Technology.