

# Exemplar-based Failure Triage for Regression Design Debugging

Zissis Poulos<sup>1</sup>, Andreas Veneris<sup>1</sup>

<sup>1</sup>Dept. of ECE, University of Toronto, Toronto, Canada.

**Abstract**—Modern regression verification often exposes myriads of failures at the pre-silicon stage. Typically, these failures need to be properly grouped into bins, which then have to be distributed to engineers for detailed analysis. The above process is coined as failure triage, and is nowadays increasing in complexity, as the size of both design logic and verification environment continues to grow. However, it remains a predominantly manual process that can prolong the debug cycle and jeopardize time-sensitive design milestones. In this paper, we propose an exemplar-based data-mining formulation of Failure Triage that efficiently automates both failure grouping and bin distribution. The proposed framework maps failures as data points, applies an affinity-propagation (AP) clustering algorithm, and operates in both metric and non-metric spaces, offering complete flexibility and significant user control over the process. Experimental results show that the proposed approach groups related failures together with 87% accuracy on the average, and improves bin distribution accuracy by 21% over existing methods.

## I. INTRODUCTION

Functional verification poses a major bottleneck in modern design cycles [1]. It becomes even more complex when performed in regression mode at early design stages, where thousands of input vectors are applied to heavily exercise both standard and corner-case functionality. The co-existence of potentially multiple design errors at this stage results into hundreds of error traces exposed. Design debugging, which accounts for more than 60% of the verification effort [1], is the task that locates design errors using information provided by these traces. It does so by using formal tools, which take as input a single error trace and automatically determine possible error sources (suspects) in the RTL [2]–[5]. These are finally examined by the engineer to track down the exact error, a process known as detailed debug [6].

However, this classical debugging approach has two major drawbacks. First, it targets each failure in isolation. As a result, different engineers may spend unnecessary resources performing detailed debug for failures that originate from the same RTL error. Second, it does not identify failures that should be high in priority for detailed debug. As such, it cannot determine the most appropriate engineer to further analyze each error trace. This uncertainty often creates confusion, with error traces constantly circulating until they are placed to the right queue for analysis.

To break this uncertainty there is a need for a preprocessing step that properly categorizes failures and assigns them to the best-suited engineer(s) for analysis. This preprocessing step is referred to as *failure triage*, and it consists of two main tasks. First, failure binning is performed. Its goal is to determine correlations between the exposed failures and bin together these failures that are likely to be caused by the same design error. The second step, failure bin distribution, identifies these

failures that should be prioritized within each bin. Next, it assigns each bin to the engineer(s) most familiar with these high-priority failures.

Traditionally, the above steps are performed manually in the industry by dedicating an engineer to constantly monitor error logs and empirically decide how to pass failures to engineers. In other cases, a script is used to parse these error logs and allocates failures following a rule-based strategy. Such *ad-hoc* manual approaches often fail to identify correlations between traces, since they rely on primitive debug paradigms.

Out of necessity to automate failure binning and failure bin distribution, recent works have formulated Failure Triage as a clustering problem [7,8]. To perform failure binning through clustering, the authors in [8] first map failures as data points into a metric space (Euclidean), whereas in [7] failures are implicitly mapped into a non-metric space. A limitation of these frameworks is that clustering can only be performed by algorithms that operate exclusively in metric or non-metric spaces, respectively. Most importantly, though, these formulations focus in failure binning and do not offer efficient solutions for failure bin distribution.

To overcome these drawbacks, in this paper, triage is viewed as an exemplar-based clustering problem. Through the proposed formulation, not only failure binning is automated, but failure bin distribution is also properly addressed. This is because exemplar-based clustering not only partitions the failure set accordingly, but also algorithmically identifies these failures that are representative of other failures in each bin. These failures-exemplars are then naturally considered as ones of high-priority for detailed debug. Moreover, the algorithm that is applied, Affinity Propagation (AP), operates in both metric and non-metric spaces, and thus leverages the merits of both representations, unlike prior art.

Experiments on four industrial designs show that the proposed work achieves 87% average accuracy for failure binning and improves bin distribution accuracy by 21%, on the average, against existing methods.

The remainder of this paper is organized as follows. Section II discusses prior work in failure triage for design debugging. Section III describes the proposed formulation and presents the exemplar-based clustering process. Finally, Section IV discusses experiments and Section V concludes the paper.

## II. PRELIMINARIES AND PRIOR ART

### A. Failure Binning

Consider an erroneous design with a single or multiple errors in the RTL that undergoes regression testing. We say that a failure occurs, when a mismatch between the expected “golden” value(s) (0,1 or X for unknown) and the observed

one(s) is identified at some observation point (primary output, probed internal signal or the output of an assertion). Suppose that at the end of regression testing,  $N$  design failures are exposed, denoted  $\mathbf{F} = \{F_1, F_2, \dots, F_N\}$ .

The goal of failure binning is to produce a complete partition of the failure set  $\mathbf{F}$  into  $K$  disjoint clusters. Ideally, failures that are caused by the same RTL error are placed into the same cluster, and into distinct clusters otherwise. For the process to be accurate, two key points have to be properly addressed. First, pairwise failure similarity needs to be quantified based on the above desired relationship, and second, a “good” number of clusters,  $K$ , needs to be selected.

In prior work, there are two well-adopted approaches to quantify pairwise failure similarity. Both approaches begin with a similar debugging step. Precisely, in [7] and [8] a baseline SAT-based debugging pass is first executed, and, for each failure  $F_i$ , the automated debugger outputs a set of design components (RTL blocks or signals), denoted  $S_i = \{s_1, s_2, \dots, s_{|S_i|}\}$ . Components  $s_1, s_2, \dots, s_{|S_i|}$  are referred to as *suspects*, and include all possible design locations that can be responsible for the observed failure. Due to its exhaustive nature, SAT-based debugging guarantees that the design location responsible for some failure  $F_i$  will be included in suspect set  $S_i$ . In this context, suspect set  $S_i$  can be viewed as a “signature” that characterizes failure  $F_i$ . Next, simulation metrics are used to assign various levels of significance to each suspect component with respect to the failure it may be responsible for. The work in [7] does so via a suspect ranking scheme, while the work in [8] adopts a data weighting scheme.

In [7], pairwise similarity between failures  $F_i$  and  $F_j$ , denoted  $s(i, j)$ , is computed as a weighted version of the Jaccard Index [9]. Particularly,  $s(i, j)$ , is given as:

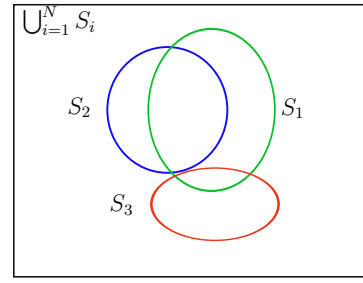
$$s(i, j) = -\left(1 - \frac{|S_i \cap S_j|}{|S_i \cup S_j|}\right) \times \pi_{ij} \quad (1)$$

In Eq. 1, the factor  $\left(1 - \frac{|S_i \cap S_j|}{|S_i \cup S_j|}\right)$  quantifies mutuality between the suspect sets of  $F_i$  and  $F_j$ , while  $\pi_{ij}$  is a measure of discrepancy between the ranks of mutual suspects in these sets [7]. The product is negated to abide to similarity semantics. As it becomes apparent, the similarities generated by Eq. 1 do not respect the triangle inequality, and thus this method operates in a non-metric space.

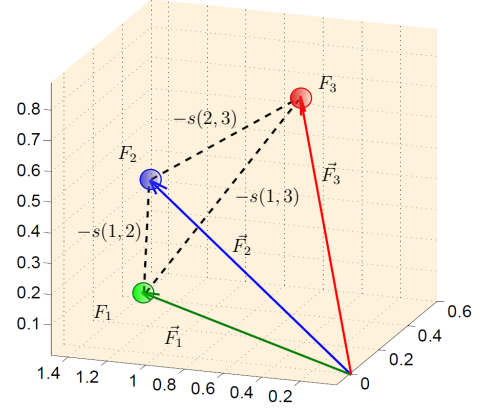
On the other hand, the authors in [8] use a feature-based representation for verification failures. Specifically, if  $s_1, s_2, \dots, s_M$  are all the distinct suspect components in  $\bigcup_{i=1}^N S_i$ , then failure  $F_i$  is represented by a real-valued feature vector  $\vec{F}_i = [x_1^i, x_2^i, \dots, x_M^i]$ , where each feature  $x_j^i$  obtains the weight (significance) of suspect  $s_j$  with respect to failure  $F_i$ , if it appears in  $S_i$ , or takes the value 0 otherwise. Failures are then mapped into a metric space, where similarity  $s(i, j)$  is defined as the negated squared Euclidean distance between  $\vec{F}_i$  and  $\vec{F}_j$ :

$$s(i, j) = -\|\vec{F}_i - \vec{F}_j\| \quad (2)$$

Figure 1 illustrates a hypothetical example of three failures  $F_1, F_2, F_3$ . Suppose that the corresponding suspect sets  $S_1, S_2$  and  $S_3$  overlap as shown in Figure 2(a). Failures, such as  $F_1$  and  $F_2$ , that have suspect sets with proportionally large



(a) Suspect set overlap



(b) Failures mapped into a metric space

Fig. 1. Failure representations

overlap are expected to be strongly related and vice versa. Of course, in [7] suspect ranks adjust the contribution of this overlap accordingly, based on Eq. 1. In [8] the contribution of the overlap and the suspect weights is implicitly represented into a metric space, as shown in Figure 2(b). Note that pairwise similarities  $s(i, j)$  are non-positive real values. In both cases, the larger  $s(i, j)$  is, the stronger the relation between  $F_i$  and  $F_j$  is considered to be.

Once failure similarity is determined, failure binning is performed through clustering algorithms. However, depending on whether similarities are metric or non-metric only specific classes of algorithms can be applied. The framework in [7] is limited to connectivity-based greedy hierarchical clustering, while in [8] hierarchical clustering is combined with k-means to produce refined failure partitions. In both cases, the number of clusters,  $K$ , is “guessed” empirically. In [7] it is based on expected cluster size/density, while in [8] the authors use a threshold applied on the clustering merge cost. These estimates experimentally appear to be the bottleneck in failure binning accuracy for both of the methodologies.

### B. Failure Bin Distribution

Suppose that failure binning generates  $K$  failure clusters,  $C_1, C_2, \dots, C_K$ . The goal of failure bin distribution is to allocate each of the  $K$  clusters to engineers that are most familiar with failures within that cluster. In past work this allocation is done as follows. For each cluster  $C_i$ , suspects across failures in  $C_i$  that have a high average rank (or average weight) are identified. Then, cluster  $C_i$  is passed to the engineer(s) that are best-suited to analyze these important suspect locations in the design. These suspects essentially correspond to a data point

that is exactly in the centroid (mean) of cluster  $C_i$ . However, it is not necessary that a failure in cluster  $C_i$  always matches with the cluster mean. In fact, this event is quite rare. Along these lines, it is naturally more suitable to allocate the cluster based on a data point associated with a failure that appears in  $C_i$ , rather than a fictional data point in the cluster mean.

### III. EXEMPLAR-BASED FAILURE TRIAGE

To overcome the problem of heuristically selecting the number of clusters and to distribute failure bins based on failures that belong to the partitioned set, we formulate triage as an exemplar-based clustering problem. In what follows, we provide the details of our methodology.

#### A. Data Preparation

Before triage commences, it is necessary to collect all the relevant information for each failure generated by regression. To this end, we follow the standard approach of performing SAT-based debugging, and for each failure  $F_i$  we generate a suspect set  $S_i$ .

If  $N$  is the the number of observed failures, then  $|\bigcup_{i=1}^N S_i| = M$  gives the number of distinct suspects across all failures  $F_1, F_2, \dots, F_N$ . Recall, that in a feature-based representation,  $M$  corresponds to the number of dimensions of the metric space where failures are mapped. Although, feature-based representation has been shown to outperform other existing methods, it does not perform well when  $M \gg N$ . This is due to the ‘‘curse of dimensionality’’ when the number of dimensions is much larger than the size of the data set. As such, to determine whether to use a metric space mapping we first compute the ratio  $N/M$  and check if  $N/M > \gamma$ , where  $\gamma \leq 0.2$ . If  $N/M > \gamma$ , then similarities  $s(i, j)$  are computed based on Eq. 2 in a metric space. Otherwise, similarities are non-metric and are computed based on Eq. 1. However, in both cases, if  $S_i \cap S_j = \emptyset$ , then  $s(i, j)$  is set to  $-\infty$ , since disjoint suspect sets indicate that failures should have minimum similarity and never be placed into the same cluster. Note that the value of threshold  $\gamma$  is determined empirically, as it will be discussed in Section IV. In both cases, pairwise failure similarities are given in the form of a  $N \times N$  similarity matrix  $\mathbf{S}$ .

#### B. Problem Formulation

Once similarity matrix  $\mathbf{S}$  is computed, failure binning takes place. However, unlike prior work, in our methodology we do not treat failure binning separate from bin distribution. Rather, we provide solutions to both problems simultaneously, in a unified sense, by formulating the whole process as exemplar-based clustering.

Exemplar-based clustering not only partitions the data, but also identifies for each cluster its most representative member, also called *exemplar*. A cluster exemplar is the member of the cluster that exhibits maximum overall similarity to other members in the cluster. In the context of failure triage a cluster exemplar can be viewed as a failure that is representative of the erroneous behaviour associated with all other failures that belong to the cluster. Intuitively, this failure-exemplar along with its suspect locations can efficiently determine how to distribute the failure bin.

To find solutions under this formulation we apply an algorithm known as Affinity Propagation (AP) [10], which is derived as an instance of max-product loopy belief propagation [9]. The AP algorithm avoids an explicit search for exactly  $K$  clusters and allows for a trade-off between the number of clusters and the within-cluster similarity that is obtained. As such, our technique does not require the number of clusters to be specified or ‘‘guessed’’ a priori. However, the algorithm allows the engineer to specify failures that are believed to be of high-importance. To this end, a quantity called *preference*, denoted  $p_i$  is associated with each failure  $F_i$  and quantifies our expectation that some failures are more suitable to be exemplars than others. The higher the preference  $p_i$ , the more likely failure  $F_i$  is to be an exemplar, and vice versa. Preferences are provided as an input to the algorithm in the form of a  $N$ -dimensional vector  $\mathbf{p}$ , and they correspond to values assigned to the  $s(i, i)$  similarities, such that  $[s(1, 1), s(2, 2), \dots, s(N, N)] = \mathbf{p}$ .

The objective function of exemplar-based clustering, and thus of the AP algorithm, is to maximize the sum of all similarities between data points in the cluster to their exemplar, while also maximizing the total preferences. Suppose the algorithm takes  $\mathbf{S}$  and  $\mathbf{p}$  as input. Then, a set of  $N^2$  binary random variables  $h_{ij} \in \{0, 1\}$  is defined, such that  $h_{ij} = 1$  if and only if failure  $F_i$  has chosen  $F_j$  as its exemplar. Note that  $h_{jj} = 1$  indicates that  $F_j$  is, in fact, an exemplar. Finally, recall that  $s(j, j) = p_j, \forall j \in \{1 \dots N\}$ . The objective function is formulated as a constrained optimization problem, as follows:

$$\max_{\{h_{ij}\}} \sum_{i=1}^{i=N} \sum_{j=1}^{j=N} s(i, j) h_{ij} \quad (3a)$$

subject to

$$\sum_j h_{ij} = 1 \quad \forall i \quad (3b)$$

$$h_{jj} = \max_i h_{ij} \quad \forall j \quad (3c)$$

Eq. 3a ensures that each point chooses exactly one other point as its exemplar. Eq. 3c guarantees that an exemplar is never assigned to another exemplar. The goal of the AP algorithm is to find settings of  $\{h_{ij}\}$  that maximize the quantity in Eq. 3a. The algorithm finds solutions based on an iterative message-passing procedure [10] and, upon convergence, it outputs a set of exemplar failures denoted as  $\mathcal{F}_{ex}$ :

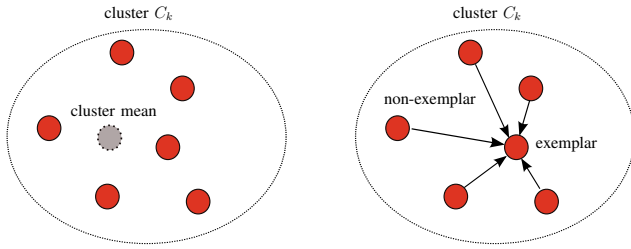
$$\mathcal{F}_{ex} = \{F_j \in \mathbf{F} : h_{jj} = 1\} \quad (4)$$

Each exemplar defines exactly one cluster of failures, and is itself a member of the cluster. Thus, the number of clusters,  $K$ , is equal to the number of exemplars in  $\mathcal{F}_{ex}$ :

$$K = |\mathcal{F}_{ex}| \quad (5)$$

On the other hand, for each non-exemplar failure  $F_i$  there exists an exemplar failure  $F_j$  which is the most similar to  $F_i$  across all other exemplars in  $\mathcal{F}_{ex}$ . The set of non-exemplars that are most similar to exemplar  $F_j$ , denoted  $\mathcal{F}_{nex}^j$  is:

$$\mathcal{F}_{nex}^j = \{F_i \in \mathbf{F} : i = \arg \max_{F_j \in \mathcal{F}_{ex}} s(i, j)\} \quad (6)$$



(a) Traditional cluster formation (b) Exemplar-based formation

Fig. 2. Failure cluster formation in traditional vs. proposed triage

Each non-exemplar failure is then assigned to the same cluster  $C_k$  as its most similar exemplar failure  $F_j$ . If exemplar  $F_j \in C_k$ , then:

$$C_k = \{F_j\} \cup \mathcal{F}_{ne}^j \quad (7)$$

Finally, for the bin distribution step, cluster  $C_k$  is assigned to the engineer that is responsible for the suspect locations of failure  $F_j$ , where  $F_j$  is the exemplar for cluster  $C_k$ . This set of design locations is given as:

$$\{s_i : s_i \in S_j \wedge F_j \in C_k \wedge F_j \in \mathcal{F}_{ex}\} \quad (8)$$

Based on the above, the benefits of this formulation are several. First, the process does not make any assumptions about similarities, apart from the fact that they are non-positive real values. They can be either metric or non-metric with no consequences to the formulation. Therefore the process can seamlessly replace existing binning algorithms irrespective of how similarities are generated, and it can leverage the merits of both representations. Further, the number of clusters,  $K$ , rises algorithmically from the message-passing procedure, and does not need to be “guessed” beforehand. Finally, as illustrated in Figure 2, bin distribution is now guided by suspects that correspond to failures included in the data set (exemplars), rather than by suspects that correspond to a data point at the cluster mean. This allows engineers to analyze each exemplar failure (error trace) and its corresponding suspect set as a whole, instead of examining suspect locations in isolation, even if these locations are significant for a particular cluster.

### C. Triage with Prior Belief

Another important benefit of the proposed methodology is that it offers significant flexibility to the engineer considering various triage scenarios. In the majority of cases, before triage commences it is rather difficult to have an estimate on the number of design errors (number of clusters) responsible for failure set  $\mathbf{F}$ . However, there are cases where engineers based on their intuition can target specific failures around which they wish  $\mathbf{F}$  to be partitioned. That is, failures that are believed should serve as exemplars of erroneous behavior. Along these lines, the proposed method allows triage to be executed with prior belief, both in a uniform and non-uniform setting, as discussed below.

1) *Uniform Setting*: In the uniform setting no assumptions are made regarding to what extent a failure should serve as an exemplar. This translates into a triage scenario where even intuitive knowledge around the importance of each failure is missing. In the proposed formulation, this is encoded by

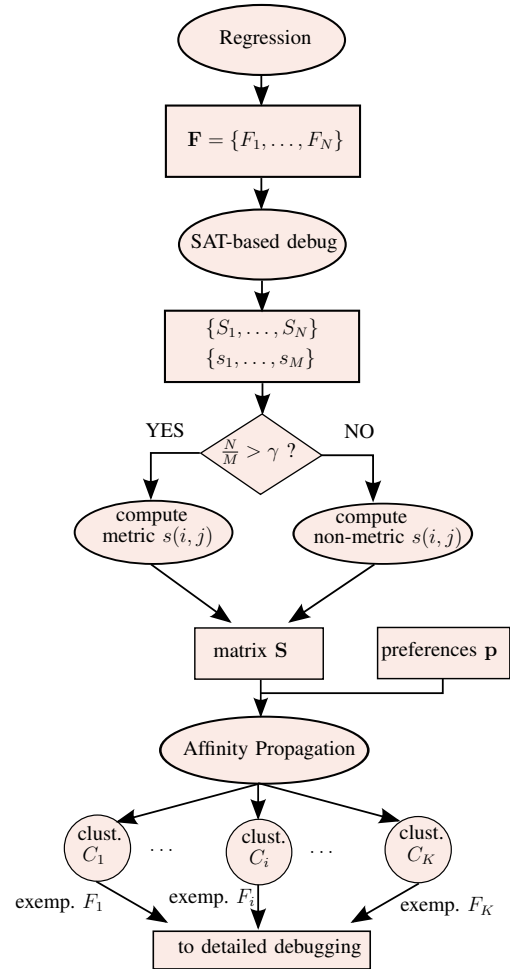


Fig. 3. Proposed triage flow

simply setting all preferences  $p_i \in \mathbf{p}$  into some constant non-positive real number. In practice, the AP algorithm performs as expected and quickly achieves convergence when preferences are fixed to the median of all similarities:

$$p_k = \frac{1}{N^2} \sum_{i=1}^{i=N} \sum_{j=1}^{j=N} s(i, j), \quad k \in \{1 \dots N\} \quad (9)$$

2) *Non-uniform Setting*: In the non-uniform setting the engineer selects specific failures to promote as exemplars a priori. If failure  $F_k$  is targeted then  $p_k$  is set to 0. Otherwise preference  $p_k$  is set to the median of similarities as in Eq. 9. Promoting specific failures as exemplars by setting higher preferences affects the number of clusters to be formed, but this number also emerges from the message-passing process. Therefore, it is not necessary that the number of clusters formed at the end will match then number of promoted failures, if this number does not reflect a reasonable partition based on the constrained optimization problem that is solved. Still, if the “guess” is close to reality then the AP algorithm can be effectively guided. Finally, note that this feature is not offered by any of the existing methodologies. These methods do allow a selection for  $K$  before the process begins, but that does not imply that  $\mathbf{F}$  is eventually partitioned around the targeted failures.

#### D. Overall Flow

A flow diagram of failure triage, as it is formulated in this work, is illustrated in Figure 3. It should be emphasized that the SAT-based debugging step that provides the “signature” suspect sets is performed in the flow whether failure triage takes place or not. Triage begins immediately after this step and preprocesses the data before detailed debug commences, where these suspect sets need to be further examined. As such, this debug step is not added by our methodology but is an inherent part of the overall debug flow in regression mode.

### IV. EXPERIMENTAL RESULTS

This Section presents experimental results for the proposed triage framework. All experiments are conducted on a single core of an Intel Core i5 3.1 GHz workstation with 8GB of RAM. Four *OpenCores* [11] designs are used for the evaluation (*vga*, *fpu*, *spi* and *mem\_ctrl*). The SAT-based debugger used to extract suspect locations is implemented based on [3]. A platform coded in Python is developed to parse debugging and simulation data, calculate the appropriate failure similarities and cluster the failure set through the AP algorithm. For each design, a set of different errors is injected each time by modifying the RTL description. The types of the injected RTL errors resemble typical human-introduced errors (missing pipeline stages, incorrect read pointers, bad stimulus etc.) that lead to non-trivial triage scenarios. In total, twenty regression test are run, generating various numbers of failures each time, caused by a different set of errors.

For each design, a pre-generated set of test sequences is used that is stored in vector files. Each regression run involves hundreds to thousands of input vectors. For the purpose of capturing failures we use end-to-end “golden model” checkers that compare the expected value for various operations, exception checkers and various assertions throughout the designs.

Table I summarizes benchmark information and statistics per regression run. From left to right, columns show the circuit name and number of gates, an enumeration for regression runs, an enumeration for regression runs,

TABLE I  
BENCHMARKS AND REGRESSION STATISTICS

Ckt. (# gates)	Test No.	# vectors	# errors	F  (N)	$ \bigcup_{i=1}^N S_i $ (M)
vga (72292)	1	25206	4	45	36
	2	25206	7	62	40
	3	25206	8	97	61
	4	31870	10	106	129
	5	31870	13	121	155
fpu (83303)	6	17365	3	19	28
	7	17365	7	30	152
	8	20094	7	55	74
	9	41759	9	83	60
	10	41759	11	125	111
spi (1724)	11	4573	3	13	38
	12	4573	5	28	46
	13	4573	6	51	82
	14	5019	8	39	196
	15	5019	9	72	113
mem_ctrl (46767)	16	10834	3	17	24
	17	10834	5	32	45
	18	10834	7	31	29
	19	13370	8	66	94
	20	13370	11	95	137

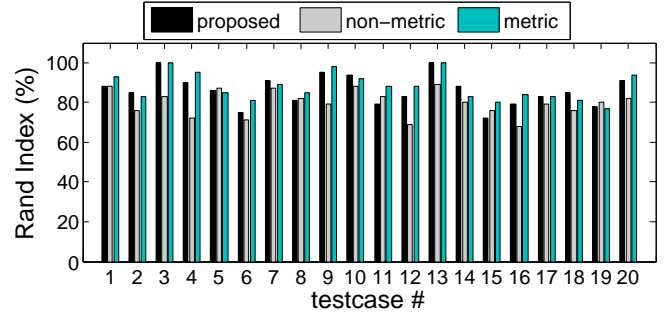


Fig. 4. Engine accuracy vs. existing methods

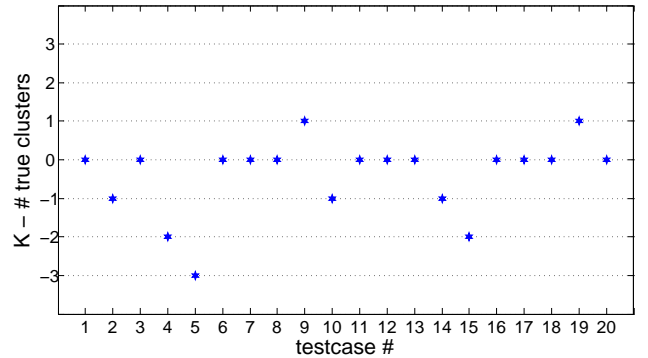


Fig. 5. Cluster prediction error

the number of input vectors, the number of simultaneous RTL errors, the number of observed failures ( $N$ ), and finally the number of distinct suspect components ( $M$ ) generated by SAT-based debugging per regression run. For all regression runs we fix threshold  $\gamma$  to 0.2. This value offers a stable behavior in the triage flow; we generally desire metric representations, and only disallow them in corner-cases when high-dimensionality becomes an impediment in performance. Finally, for evaluation purposes triage is run under the uniform setting described in Section III.C. That is, we assume that there are no targeted failures. This allows us to conduct a fair comparison against existing methods.

To evaluate failure binning accuracy, we use the Rand Index (R.I.) measure [9]. This metric compares the estimated clustering against a reference failure binning, with the latter corresponding to an ideal partition where all failures are grouped with 100% accuracy. The metric ranges from 0 to 1 and represents the fraction of correct clustering decisions. Accuracy is hence measured as  $100 \times \text{R.I.} \%$ .

Figure 4 shows a comparison between the proposed framework and the methods described in [7] (*non-metric*) and [8] (*metric*), in terms of failure binning accuracy. The proposed triage engine outperforms the framework in [7] in 14/20 regression runs, while it achieves better accuracy in 7/20 cases compared to the engine in [8]. Precisely, across all regression runs, the proposed engine achieves 87% clustering accuracy on average, compared to 77% and 88% accuracy of the *non-metric* and *metric* methods, respectively. Note that when  $N/M < \gamma$ , such as in test-cases 7 and 14, the proposed method generates *non-metric* similarities. This has a positive

impact on accuracy by avoiding effects of high-dimensionality. The metric approach in these cases exhibits lower accuracy as it also shown by Figure 4. Finally, it is worth mentioning that under the non-uniform setting the proposed method can achieve up to 96% average accuracy when targeted failures are carefully chosen.

To illustrate how accurate the number of generated clusters,  $K$ , is in our framework, Figure 5 shows how far this prediction is from the number of design errors responsible for the observed failures. We refer to the latter as “true clusters”. The prediction error is then given as  $(K - \# \text{ true clusters})$  for all test-cases. In ideal cases,  $K$  is equal to the number of design errors and the prediction error is 0. Figure 5 shows that in 12/20 regression runs the AP algorithm achieves a perfect prediction, which greatly boosts binning accuracy. Interestingly, in the rest of cases, the number of formed clusters is usually smaller (by 1 to 3 clusters). Only in test-cases 9 and 19 the number of clusters is larger than necessary.

As results indicate, the failure binning step in the proposed triage flow demonstrates high accuracy, comparable to the currently most efficient methodology in [8]. However, the major strength of the proposed flow is its effective exemplar-based bin distribution step. Since bin distribution in this work is performed via exemplar failures, while existing techniques use high-weight suspect locations to guide the process, we need to use a common reference for comparison purposes. To this end, we identify whether the design error responsible for failures in a particular cluster is included in the suspect set of the exemplar failure. If the design error location is indeed in the suspect list, we obtain its rank or weight (significance), which is already computed to generate similarities in the binning step. Finally, we determine if the error location has a high rank or weight compared to other suspect locations for the same exemplar failure. We do this by sorting suspects in order of decreasing weight. Ideally, the error location resides among the top 10-20% of suspect locations in the ordered list. In that case the exemplar failure and its suspect locations are effectively prioritized.

Table II summarizes experimental results regarding the bin distribution step. Column 1 indicates the regression run number. Columns 2 to 4 show what position the responsible design error takes in the sorted list of suspects that is returned to the engineer by bin distribution in [7], [8], and the proposed flow, respectively. The positions are normalized over the size of the suspect list each time. Thus, when the position is low, then the suspect component that includes the design error appears in the first positions of the list, and vice versa. Each row in the table provides the average design error position per regression run. The last column shows the improvement that is achieved by the proposed method compared to the bin distribution approach in [8]. From Table II we observe that, on average, the proposed bin distribution approach pushes the responsible design error higher in the list compared to existing methods in 12/20 regression scenarios. In total, the average improvement that is achieved compared to the best of the two existing methods (metric) is approximately 21% across all regression runs.

As far as time consumption is concerned, this is vastly dominated by the SAT-based debugging step which is per-

TABLE II  
BIN DISTRIBUTION PERFORMANCE

Test No	avg. normalized error position			improvement (%)
	non-metric	metric	proposed	
1	0.22	0.13	<b>0.10</b>	23
2	0.26	0.27	<b>0.16</b>	41
3	0.19	0.18	0.20	-11
4	0.31	0.34	<b>0.27</b>	20
5	0.40	0.18	0.20	-11
6	0.34	0.35	<b>0.16</b>	54
7	0.24	0.17	<b>0.10</b>	41
8	0.18	0.12	0.19	-58
9	0.27	0.26	<b>0.23</b>	12
10	0.51	0.48	<b>0.32</b>	33
11	0.30	0.27	0.29	-7
12	0.18	0.23	<b>0.08</b>	65
13	0.17	0.15	0.15	0
14	0.33	0.24	<b>0.16</b>	33
15	0.21	0.19	0.26	63
16	0.46	0.32	0.33	-3
17	0.10	0.09	0.09	0
18	0.26	0.14	<b>0.11</b>	21
19	0.56	0.41	<b>0.19</b>	54
20	0.51	0.46	<b>0.33</b>	28
<b>AVG</b>	0.300	0.249	<b>0.196</b>	<b>21</b>

formed whether triage takes place or not. This step consumes from approximately 400 to 7000 seconds per regression test-case. The added overhead due to failure binning and bin distribution is negligible and is in the range of 20 to 50 seconds approximately per regression run.

## V. CONCLUSION

To summarize, this work introduces a novel exemplar-based clustering formulation for the growing problem of failure triage in regression design debugging flows. It proposes the use of Affinity Propagation to simultaneously provide solutions to failure binning and bin distribution as a unified constrained optimization problem. Experimental results demonstrate the applicability and efficiency of the proposed triage engine, and indicate that it outperforms existing methods for the important step of bin distribution.

## REFERENCES

- [1] H.Foster, “From volume to velocity: The transforming landscape in function verification,” in *Design Verification Conf.*, 2011.
- [2] O. Sarbishei, M. Tabandeh, B. Alizadeh, and M. Fujita, “A formal approach for debugging arithmetic circuits,” in *IEEE Transactions on CAD*, vol. 28, no. 5, May 2009, pp. 742–754.
- [3] A. Smith, A. Veneris, M. F. Ali, and A. Viglas, “Fault diagnosis and logic debugging using Boolean satisfiability,” *IEEE Transactions on CAD*, vol. 24, no. 10, pp. 1606–1621, 2005.
- [4] S. Mirzaeiian, F. Zheng, and K. Cheng, “Rtl error diagnosis using a word-level sat-solver,” in *International Test Conference*, 2008, pp. 1–8.
- [5] K. hui Chang, I. Wagner, V. Bertacco, and I. L. Markov, “Automatic error diagnosis and correction for rtl designs,” in *Proc. International High Level Design Validation and Test Workshop (HLDVT) 2007*, pp. 65–72.
- [6] S.Safarpour, B.Keng, Y.S.Yang, and E.Qin, “Failure triage: The neglected debugging problem,” in *Design and Verification Conference*, 2012.
- [7] Z. Poulos, Y. Yang, and A. Veneris, “Simulation and satisfiability guided counter-example triage for rtl design debugging,” in *Int’l Symposium on Quality Electronic Design*, 2014, pp. 394–399.
- [8] Z. Poulos and A. Veneris, “Clustering-based failure triage for rtl regression debugging,” in *Int’l Test Conference*, 2014.
- [9] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 2007.
- [10] B. J. Frey and D. Dueck, “Clustering by passing messages between data points,” *Science*, vol. 315, pp. 972–976, 2007.
- [11] OpenCores.org, “<http://www.opencores.org>,” 2007.