

# 1 The Unreasonable Effectiveness of 2 Traditional Information Retrieval in Crash 3 Report Deduplication

4 Joshua Charles Campbell<sup>1</sup>, Eddie Antonio Santos<sup>1</sup>, and Abram Hindle<sup>1</sup>

5 <sup>1</sup>Department of Computing Science, University of Alberta, Edmonton, Canada

## 6 ABSTRACT

7 Organizations like Mozilla, Microsoft, and Apple are flooded with thousands of automated crash reports per day. Although crash reports contain valuable information for debugging, there are often too many for developers to examine individually. Therefore, in industry, crash reports are often automatically grouped together in buckets. Ubuntu's repository contains crashes from hundreds of software systems available with Ubuntu. A variety of crash report bucketing methods are evaluated using data collected by Ubuntu's Apport automated crash reporting system. The trade-off between precision and recall of numerous scalable crash deduplication techniques is explored. A set of criteria that a crash deduplication method must meet is presented and several methods that meet these criteria are evaluated on a new dataset. The evaluations presented in this paper show that using off-the-shelf information retrieval techniques, that were not designed to be used with crash reports, outperform other techniques which are specifically designed for the task of crash bucketing at realistic industrial scales. This research indicates that automated crash bucketing still has a lot of room for improvement, especially in terms of identifier tokenization.

8 Keywords: ...

## 9 1 INTRODUCTION

10 Ada is a senior software engineer at Lovelace Inc., a large software development  
11 company. Lovelace has just shipped the latest version of their software to hundreds  
12 of thousands of users. A short while later, as Ada is transitioning her team to  
13 other projects, she gets a call from the quality-assurance team (QA) saying that  
14 the software she just shipped has a crashing bug affecting two-thirds of all users.  
15 Worse yet, Ada and her team can't replicate the crash. What would really be  
16 helpful is if every time that crash was encountered by a user, Lovelace would  
17 automatically receive a *crash report* [Seo and Kim], with some *context* information  
18 about what machine encountered the crash, and a *stack trace* [Seo and Kim] from  
19 each thread. Developers consider stack traces to be an indispensable tool for  
20 debugging crashed programs—a crash report with even one stack trace will help  
21 fix the bug significantly faster than if there were had no stack traces available at  
22 all [Schröter et al.].

23 Luckily for Ada, Lovelace Inc. has gone through the monumental effort of  
24 setting up an automated crash reporting system, much like Mozilla's Crash Error

25 Reports [Mozilla Corporation], Microsoft's WER [Glerum et al.], or Apple's Crash  
26 Reporter [app]. Despite the cost associated with setting up such a system, Ada  
27 and her team find the reports it provides are invaluable for collecting telemetric  
28 crash data [Ahmed et al.].

29 Unfortunately, for an organization as large as Lovelace Inc., with so many users,  
30 even a few small bugs can result in an unfathomable amount of crash reports. As  
31 an example, in the first week of 2016 alone, Mozilla received 2 189 786 crash  
32 reports, or about 217 crashes every minute on average.<sup>1</sup> How many of crash  
33 reports are actually relevant to the bug Ada is trying to fix?

34 The sheer amount of crash reports present in Lovelace's crash reporting system  
35 is simply too much for one developer, or even a team of developers, to deal with  
36 by hand. Even if Ada spent only one second evaluating a single crash report, she  
37 would still only be able to address 1/3 of Lovelace's crash reports received during  
38 one day of work. Obviously, an automated system is needed to associate related  
39 crash reports together, relevant to this one bug, neatly in one place. All Ada would  
40 have to do is to select a few stack traces from this *crash bucket* [Glerum et al.],  
41 and get on with debugging her application. Since this hypothetical bucket has  
42 all crash stack traces caused by the same bug, Ada could analyze any number of  
43 stack traces and pinpoint exactly where the fault is and how to fix it.

44 The questions that this paper seeks to answer are:

45 **RQ1:** What are effective, industrial-scale methods of crash report bucketing?

**RQ2:** How can these methods be tuned to increase precision or recall?

46 This paper will evaluate existing techniques relevant to crash report bucketing,  
47 and propose a new technique that attempts to handle this fire hose of crash  
48 reports with industrially relevant upper bounds ( $O(\log n)$  per report, where  $n$  is  
49 number of crash reports). In order to validate new techniques some of the many  
50 techniques described in the literature are evaluated and compared. The results of  
51 the evaluation shows that techniques based on the standard information retrieval  
52 statistic, *term frequency*  $\times$  *inverse document frequency* (tf-idf), do better than  
53 others, despite the fact these techniques discard information about what is on the  
54 top of the stack and the order of the frames on the stack.

## 55 1.1 Contributions

56 This paper presents PARTYCRASHER, a technique that buckets crash reports. It  
57 extends the work done by Lerch and Mezini [Lerch and Mezini] to the field of crash  
58 report deduplication and show that despite its simplicity, it is quite effective. This  
59 paper contributes:

- 60 1. a criterion for industrial-scale crash report deduplication techniques;
- 61 2. replication of some existing methods of deduplication and evaluations of  
62 these methods on open source crash reports, providing evidence of how well  
63 each technique performs at crash report bucketing;

<sup>1</sup><https://crash-stats.mozilla.com/api/SuperSearch/?date=>\%3d2016-01-01&date=<\%3d2016-01-08> The total number of crashes will slowly increase over time and then eventually drop to zero due to Mozilla's data collection and retention policies.

- 64 3. implementation of these methods in an open source crash bucketing frame-  
65 work;
- 66 4. evaluation based on the automated crashes collected by the Ubuntu project's  
67 Apport tool, the only such evaluation at the time of writing;
- 68 5. a bug report deduplication method that outperforms other methods when  
69 contextual information is included along with the stack trace.

## 70 1.2 What makes a crash bucketing technique useful for industrial scale crash 71 reports?

72 The volume, velocity, variety, and veracity (uncertainty) of crash reports makes  
73 crash report bucketing a big-data problem. Any solution needs to address concerns  
74 of big-data systems especially if it is to provide developers and stakeholders with  
75 *value* [20]. Algorithms that run in  $O(n^2)$  are unfeasible for the increasingly large  
76 amount of crash reports that need to be bucketed. Therefore, an absolute upper-  
77 bound of  $O(n \log n)$  is chosen for evaluated algorithms.

78 The methods evaluated in this paper were methods found in the literature, or  
79 methods that the authors felt possibly had promise. Methods that were evaluated  
80 were restricted to those that met the following criteria. The criteria were chosen  
81 to match the industrial scenario as described in the introduction.

- 82 1. Each method must scale to industrial-scale crash report deduplication re-  
83 quirements. Therefore, it must run in  $O(n \log n)$  total time. Equivalently,  
84 each new, incoming crash must be able to be assigned a bucket in  $O(\log n)$   
85 time or better.
- 86 2. No method may delay the bucketing of an incoming crash report significantly,  
87 so that up-to-date near-real-time crash reports, summaries, and statistics are  
88 available to developers at all times. This requires the method to be *online*.
- 89 3. No method may require developer intervention once it is in operation, or  
90 require developers to manually categorize crashes into buckets. This requires  
91 the method to be *unsupervised*.
- 92 4. No method may require knowledge of the eventual total number of buckets  
93 or any of their properties beforehand. Each method must be able to increase  
94 the number of buckets only when crashes associated with new faults arrive  
95 due to changes in the software system for which crash reports are being  
96 collected. This requires the method to be *non-stationary*.

97 Several deduplication methods are evaluated. They can be categorized into  
98 two major categories. First, several methods based on selecting pre-defined parts  
99 of a stack to generate a *signature* were evaluated. The simplest of these methods  
100 is the **1Frame** method, that selects the name of the function on top of the stack  
101 as a signature. All crashes that have identical signatures are then assigned to a  
102 single bucket, identified by the signature used to create it.

```

#1 0x00002b344498a150 in CairoOutputDev::setDefaultCTM () from /usr/lib/libpoppler-glib.so.1
#2 0x00002b344ae2cefc in TextSelectionPainter::TextSelectionPainter () from /usr/lib/libpoppler.so.1
#3 0x00002b344ae2cff0 in TextPage::drawSelection () from /usr/lib/libpoppler.so.1
#4 0x00002b344498684a in poppler_page_render_selection () from /usr/lib/libpoppler-glib.so.1

```

Method	Signature
1Frame	CairoOutputDevsetDefaultCTM
2Frame	CairoOutputDev::setDefaultCTM TextSelectionPainter::TextSelectionPainter
3Frame	CairoOutputDev::setDefaultCTM TextSelectionPainter::TextSelectionPainter TextPage::drawSelection
1Addr	0x00002b344498a150
1File	No Signature (no source file name given in the stack)
1Mod	/usr/lib/libpoppler-glib.so.1

Method	Tokenization
No tokenization	#1 0x00002b344498a150 in CairoOutputDev::setDefaultCTM () from /usr/lib/libpoppler-glib.so.1
Lerch	<code>0x00002b344498a150</code> <code>cairooutputdev</code> <code>setDefaultctm</code> <code>from</code> <code>libpoppler</code> <code>glib</code>
Space	<code>#1</code> <code>0x00002b344498a150</code> <code>in</code> <code>CairoOutputDev::setDefaultCTM</code> <code>()</code> <code>from</code> <code>/usr/lib/libpoppler-glib.so.1</code>
Camel	<code>#1</code> <code>0</code> <code>x</code> <code>00002</code> <code>b</code> <code>344498</code> <code>a</code> <code>150</code> <code>in</code> <code>Cairo</code> <code>Output</code> <code>Dev</code> <code>set</code> <code>Default</code> <code>CTM</code> <code>from</code> <code>usr</code> <code>lib</code> <code>libpoppler</code> <code>glib</code>

**Figure 1.** An example stack trace (top), its various signatures (middle), and various tokenizations of the top line of the trace (bottom).

103 Similarly, signature methods `2Frame` and `3Frame` concatenate the names of the  
104 two or three functions on top of the stack to produce a signature. `1Addr` selects  
105 the address of the function on top of the stack to generate a signature rather  
106 than the function name. `1File` selects the name of the source file in which the  
107 function on top of the stack is defined to generate a signature, and `1Mod` selects  
108 either the name of the file or the name of the library, depending on which is  
109 available. Figure 1 shows an example stack trace and how the various signatures  
110 are extracted from it using these methods. All of the signature-based methods,  
111 as implemented, run in  $O(n \log n)$  total time or  $O(\log n)$  amortized time.

112 The second category of methods are those based on tf-idf [Salton and McGill]  
113 and inverted indices, as implemented by the off-the-shelf information-retrieval  
114 software ElasticSearch 1.6 [Elasticsearch BV]. tf-idf is a way to normalize a *token*  
115 based on both on its occurrence in a particular document (in our case, crash  
116 reports), and inversely proportional to its appearance in all documents. That  
117 means that common tokens that appear frequently in nearly *all* crash reports  
118 have little discriminative power compared to tokens that appear quite frequently  
119 in a small set of crash reports.

### 120 1.3 Background

121 Of course, the idea of crash bucketing is not new; Mozilla's system performs  
122 bucketing [Dhaliwal et al., Ahmed et al.], as does WER [Glerum et al.]. Many  
123 approaches make the assumption that two crash reports are similar if their stack  
124 traces are similar. Consequently, researchers [Brodie et al., Liu and Han, Modani et al.,  
125 Bartz et al., Glerum et al., Dhaliwal et al., Dang et al., Wang et al., Lerch and Mezini,  
126 Wu et al.] have proposed various methods of finding similar stack traces, crash re-  
127 port similarity, crash report deduplication, and crash report bucketing. In order to  
128 motivate the evaluation and design choices it is necessary to look at what already  
129 has been proposed.

130 Empirical evidence suggests that a function responsible for crash is often at or  
131 near the top of the crash stack trace [Brodie et al., Schröter et al., Wu et al.]. As  
132 such, many bucketing heuristics employ higher weighting for grouping functions

133 near the top of the stack [Modani et al., Glerum et al., Wang et al.]. Many of  
134 these methods are similar to or extensions of the `1Frame` method, that assumes  
135 that the function name on the top of the stack is the most (or only) important  
136 piece of information for crash bucketing. However, at least one study refutes the  
137 effectiveness of truncating the stack trace [Lerch and Mezini]. The most influential  
138 discriminative factors seem to be function name [Lerch and Mezini] and module  
139 name [Bartz et al., Glerum et al.].

140 Lerch and Mezini [Lerch and Mezini] did not directly address crash report  
141 bucketing; they addressed *bug report* deduplication through stack trace similar-  
142 ity. They deduplicated bug reports that included stack traces by comparing the  
143 traces with tf-idf, which is usually applied to natural language text. Although  
144 crash bucketing was implicit in this approach to bug-report-deduplication, the  
145 authors did not compare this technique against the other crash report dedupli-  
146 cation techniques. Unlike the signature-based methods, tf-idf-based methods do  
147 not consider the order that frames appear on the stack. A function at the top of  
148 the stack is treated identically to a function at the bottom of the stack.

149 Their method of bug report deduplication is applied to to crash report dedu-  
150 plication and evaluated in this paper, both excluding *contextual* data from the  
151 crash report as suggested by Lerch and Mezini [Lerch and Mezini] and including  
152 it. These methods are listed in the evaluation section as the `Lerch` method and  
153 the `LerchC` method, respectively. The contextual data is collected at the same  
154 time as the stack trace by automated crash reporting tools. Variants of the `Lerch`  
155 and `LerchC` methods were also evaluated. The variants replace the tokenization  
156 pattern used in `Lerch` and `LerchC` with a different tokenization pattern. These  
157 methods were named `Space`, `SpaceC`, `Camel`, and `CamelC`. The name indicates  
158 that tokenization is employed, followed by a `C` if the evaluation included the en-  
159 tire context of the stack trace along with the stack trace itself. Figure 1 shows  
160 how each method tokenizes a sample stack frame.

161 Modani *et al.* [Modani et al.] provide two techniques to improve performance  
162 of the various other algorithms. These techniques are inverted indexing and top-*k*  
163 indexing, both of which are evaluated in this paper. Inverted indexing is em-  
164 ployed to improve the performance of all of the tf-idf-based methods including  
165 `Lerch` and `LerchC` (however Modani *et al.* did not use tf-idf in their evaluation).  
166 The implementation is provided by Elasticsearch 1.6 [Elasticsearch BV]'s index-  
167 ing system. Top-*k* indexing is employed to evaluate all of the methods that use  
168 the top portions of stacks, including `1Frame`, `2Frame`, `3Frame`, `1File`, etc.

#### 169 **1.4 Methods Not Appearing In This Report**

170 Mozilla's deduplication technique, at the time of writing, as it is implemented  
171 in Socorro [soc] requires a large number of hand-written regular expressions to  
172 select, ignore, skip, or summarize various parts of the crash report. These must  
173 be maintained over time by Mozilla developers and volunteers in order to stay  
174 relevant to crashes as versions of Firefox are released. This technique typically  
175 uses one to three of the frames of the stack and likely has similar performance to  
176 `1Frame`, `2Frame`, and `3Frame`. Furthermore, the techniques employed by Mozilla

177 are extremely specific to their major product, Firefox, while the evaluation dataset  
178 contains crashes from 616 other systems.

179 In 2005, Brodie *et al.* [Brodie et al.] presented an approach that normalizes the  
180 call stack to remove non-discriminative functions as well as flattening recursive  
181 functions, and compares stacks using weighted edit distance. Since pairwise stack  
182 matching would be unfeasible on large data sets—having a minimum worst case  
183 run-time of  $O(n^2)$ —they index a hash of the top  $k$  function names at the top of  
184 the stack and use a B+Tree look-up data structure. Several approaches since have  
185 used some stack similarity metric, and found that the most discriminative power  
186 is in the top-most stack frames—*i.e.*, the functions that are *closer* to the crash  
187 point.

188 Liu and Han [Liu and Han] grouped crashes together if they suggest the same  
189 fault location. The fault locations were found using a statistical debugging tool  
190 called SOBER [Liu et al.], that, trained on failing and passing *execution traces*  
191 (based on instrumenting Boolean predicates in code [Liblit et al.]), returns a ranked  
192 list of possible fault locations. Methods involving full instrumentation [Liu and Han]  
193 or static call graph analysis [Wu et al.] are also deemed unfeasible, as they are  
194 not easy to incorporate into already existing software, and often incur pairwise  
195 comparisons to bucket regardless of instrumentation cost. Methods that already  
196 assume buckets such as Kim *et al.* [Kim et al.] and Wu *et al.* [Wu et al.] are dis-  
197 regarded as well.

198 Modani *et al.* [Modani et al.] propose several algorithms. The first algorithm  
199 employs edit distance, requiring  $O(n^2)$  total time. The second and third algo-  
200 rithms are similar, employing longest common subsequences and longest common  
201 prefixes, respectively. The longest common subsequence problem is, in general,  
202 NP-hard in the number of sequences (corresponding to crashes for the purposes  
203 of this evaluation). The longest common prefix algorithm can be implemented  
204 sufficiently efficiently for the purposes of this evaluation, but was not evaluated  
205 here because it must produce at least as many buckets as the `1Frame` algorithm,  
206 that already creates too many buckets. Thus no Modani *et al.* [Modani et al.]  
207 comparison algorithms were used.

208 In addition to comparison algorithms that might be used for deduplication di-  
209 rectly, Modani *et al.* [Modani et al.] also provide several algorithms for identifying  
210 frames that may be less useful in each stack and removing them from those stacks.  
211 These algorithms would then be combined with their other algorithms and are  
212 not evaluated in this paper. One such algorithm removes frequent frames, such as  
213 `main()` that occur in many stacks. A similar effect is gained from `tf-idf`, because  
214 the inverse document frequency reduces the weight of terms that are found in  
215 many documents (crashes). These filtering techniques were not evaluated.

216 Bartz *et al.* [Bartz et al.] also used edit distance on the stack trace, but a  
217 weighted variant with weights learned from training data. Consequently, they  
218 were able to consider other data in the crash report aside from the stack trace.  
219 The weights learned suggested some interesting findings: substituting a module in  
220 a call stack resulted in a much higher distance; as well, the call stack edit distance  
221 was found to be the highest-weighted factor, despite the consideration of other

222 crash report data, confirming the intuition in the literature of the stack trace's  
223 importance.

224 The methods based on edit distance—*viz.*, Brodie *et al.* [Brodie et al.], Modani  
225 *et al.* [Modani et al.], Bartz *et al.* [Bartz et al.]<sup>2</sup>—are disqualified due to their  
226 requirement of pairwise comparisons between stack traces, with an upper-bound  
227 of  $O(n^2)$ .

228 Schröter *et al.* [Schröter et al.] empirically studied developers' use of stack  
229 traces in debugging and found that bugs are more likely to be fixed in the top  
230 10 frames of their respective crash stack trace, further confirming the surprising  
231 significance of the top- $k$  stack frames in crash report bucketing, which is also  
232 corroborated more recently by Wu *et al.* [Wu et al.].

233 Glerum *et al.* [Glerum et al.] describe the methods used by Microsoft's Win-  
234 dows Error Reporting (WER) service. Although they tout having over 500 heuris-  
235 tics for crash report bucketing—many derived empirically—a large bulk of the  
236 bucketing is attributed to top-1 module offset; over 91% of bucketing is attributed  
237 to eight heuristics alone.

238 To avoid the  $O(n^2)$  pairwise comparisons common to many of the previous ap-  
239 proaches, Dhaliwal *et al.* [Dhaliwal et al.] proposed a weighted edit distance tech-  
240 nique that creates *representative stack traces*—a probability distribution based  
241 on all stack traces seen within a bucket. Thus, instead of computing similarity  
242 against all stack traces in a bucket, one would only use the weights derived from  
243 all stack traces in the bucket simultaneously.

244 The method described in Dhaliwal *et al.* [Dhaliwal et al.] is not included in  
245 the evaluation because it first subdivides buckets produced by the `1Frame` dedupli-  
246 cation method, and requires  $O(|B|^2)$  total time to run, where  $|B|$  is the number  
247 of buckets. Its use of the `1Frame` method already produces a factor of 1.67 times  
248 too many buckets. Despite the optimization in Dhaliwal *et al.* [Dhaliwal et al.]  
249 that attempts to avoid  $O(n^2)$  behaviour, it has  $O(|B|^2)$  behaviour. Since the  
250 number of buckets increases over time, though at a slower rate, this method will  
251 eventually become computationally unfeasible if old data is not discarded.

252 Kim *et al.* [Kim et al.] constructed *Crash Graphs*, that are simply directed  
253 graphs using stack frames as nodes and their adjacency to other stack frames as  
254 edges. This also proved to be a useful crash visualization technique.

255 Dang *et al.* [Dang et al.] created the *position independent model* that places  
256 more weight on stack frames closer to the top of the stack; and favours stacks  
257 whose matched functions are similarly spaced from each other. Purporting sig-  
258 nificantly higher accuracy than previous methods, this technique suffers from a  
259 proposed  $O(n^3)$  clustering algorithm.

260 Wang *et al.* [Wang et al.] propose three different methods, that they refer to as  
261 rules. The first rule requires an incoming crash to be compared to every existing  
262 crash, requiring  $O(n^2)$  time. The second rule compares only the top frame of  
263 every crash by considering two crashes related if the file names in the top frame  
264 of the crash are the same. This method is listed in the evaluation as the `1File`

---

<sup>2</sup>They first use naïve methods for indexing as well, that *is* evaluated here

265 method. The third rule requires a set of common “frequent closed ordered sub-sets”  
266 of stack frames to be extracted from known “crash types” that are pre-categorized  
267 groups of crashes that have been bucketed using a separate method. The third  
268 rule requires  $O(|B|^2)$  total time where  $|B|$  is the number of buckets created by  
269 the other method. Specifically the authors use the method of comparing the top  
270 frame from each stack, that is evaluated in this paper as the `1Frame` method.  
271 This method appears to create a number of buckets roughly proportional to the  
272 number of seen crashes,  $n$ . Thus, the third rule requires  $O(n^2)$  total time, though  
273 with a low coefficient. The only method from Wang *et al.* [Wang et al.] directly  
274 evaluated in this paper is the method of comparing file names at the top of the  
275 stack.

276 Thus, there are many approaches for bucketing crash reports and crash report  
277 similarity, but some are less realistic or industrially applicable than others. Any  
278 new work in the field must attempt to compare itself against some of the prior  
279 techniques such as Lerch and Mezini [Lerch and Mezini].

## 280 2 METHODOLOGY

281 First, the requirements for an industrial-scale automated crash deduplication sys-  
282 tem were characterized by looking at systems that are currently in use. Then,  
283 a variety of methods from the existing literature were evaluated for applicability  
284 to the task of automated crash report deduplication. Several methods that met  
285 the requirements were selected. A general purpose Python framework in which  
286 any of the selected deduplication methods could be supported and evaluated was  
287 developed, and then used to evaluate all of the methods by simulating the process  
288 of automated crash reports arriving over time. Additionally, a dataset that could  
289 be used as a gold set to judge the performance of such methods was obtained. The  
290 dataset was then filtered to include only crash reports that had been deduplicated  
291 by human developers and volunteers.

292 Various approaches of automatic crash report categorization (the exact prob-  
293 lem that Ada is tasked with solving) is simulated. First, a crash report arrives  
294 with no information other than what was gathered by the automated reporting  
295 mechanisms on the user’s machine. This report might include a description writ-  
296 ten by the user of what they were doing when the crash occurred. However, these  
297 descriptions are often full of foul language as opposed to useful information for  
298 debugging. Figure 3 is an example of one of the crash reports used in the evalua-  
299 tion with a user-submitted description on the second line, metadata in the middle,  
300 and a stack trace on the bottom.

### 301 2.1 Mining Crash Reports

302 The first step in the evaluation procedure is mining of crash reports from Ubuntu’s  
303 bug repository, Launchpad [Canonical Ltd.]. This was done using a modified  
304 version of Bicho [23], a software repository mining tool.<sup>3</sup> Over the course of one  
305 month, Bicho was able to retrieve 126 609 issues from Launchpad, including 80 478

<sup>3</sup><https://github.com/orezpraw/Bicho/>



306 stack traces in 44 465 issues. Some issues contain more than one stack trace. For  
307 issues that contained more than one stack trace, the first stack trace posted to that  
308 issue was selected, yielding 44 465 issues with crash reports and stack traces. The  
309 first stack trace is selected because it is the one that arrives with the automated  
310 crash report, generated by the instrumentation on the user's machine.

311 Ubuntu crash reports were used for the evaluation because they are automati-  
312 cally generated and submitted but many of them have been manually deduplicated  
313 by Ubuntu developers and volunteers. Other data sources, such as Mozilla's Crash  
314 Reports have already been deduplicated by Mozilla's own automated system, not  
315 by humans.

316 Next, the issues were put into groups based on whether they were marked as  
317 duplicates of another issue, resulting in 30 664 groups of issues. These groups are  
318 referred to as "issue buckets" for the remainder of the paper, to prevent confound-  
319 ing with groups of crash reports, that will be referred to as "crash buckets." This  
320 dataset is available!<sup>4</sup>

### 321 **2.1.1 Stack Trace Extraction**

322 Each issue and stack trace obtained from Ubuntu is formatted as plain text, as  
323 shown in Figure 3. They were then parsed into JSON-formatted data with indi-  
324 vidual fields for each item, such as address, function name, and which library the  
325 function came from. Unfortunately, this formatting is not always consistent and  
326 may be unusable. For example, some stack traces contain unintelligible binary  
327 data in place of the function name. This could be caused by memory corruption  
328 when the stack trace was captured. 2 216 crash reports and stack traces were  
329 thrown out because their formatting could not be parsed, leaving 41 708 crash  
330 reports with stack traces.

### 331 **2.1.2 Crash Report and Stack Trace Data**

332 Issues were then filtered to only those that had been deduplicated by Ubuntu  
333 developers and other volunteers, yielding 15 293 issues with 15 293 stack traces  
334 in 3 824 issue buckets. These crash reports were submitted to Launchpad by the  
335 Apport tool.<sup>5</sup> They were collected over a one month period. Because Launchpad  
336 places restrictions on how often the Launchpad API can be used to request data,  
337 and each crash report required multiple requests, it required over 20 seconds to  
338 download each issue. The crash reports used in the evaluation span 617 different  
339 source packages, each of which represents a software system. The only commonali-  
340 ties between them are that they are all written in C, C++, or other languages that  
341 compile to binaries debuggable by a C debugger, and that they are installed and  
342 used on Ubuntu. The most frequently reported software system is Gnome<sup>6</sup>, which  
343 has 2 154 crash reports with stack traces. This dataset is large, comprehensive  
344 and covers a wide variety of projects.

<sup>4</sup><https://pizza.cs.ualberta.ca/buckets.tgz>, augmented over time as more crash re-  
ports are mined from the Launchpad Ubuntu issue repository.

<sup>5</sup><https://launchpad.net/apport>

<sup>6</sup><https://www.gnome.org/>

## 345 2.2 Crash Bucket Brigade

346 In order to simulate the timely nature of the data, each report is added to a  
347 simulated crash report repository *one at a time*. This is done so that no method  
348 can access data “from the future” to choose a bucket to assign a crash report to. It  
349 is first assigned a bucket based on the crashes and buckets already in the simulated  
350 repository, then it is added to the repository as a member of that bucket.

## 351 2.3 Deciding when a Crash is not Like the Others

352 For methods based on Lerch and Mezini, there is a threshold value,  $T$ , that de-  
353 termines how often, and when, an incoming crash report is assigned to a new  
354 bucket. A specific value for  $T$  was not described by Lerch and Mezini, so a range  
355 of different values from 1.0 to 10.0 were evaluated. Higher values of  $T$  will cause  
356 the algorithm to create new buckets more often.

357 The threshold value applies to the *score* produced by the Lucene search engine  
358 inside Elasticsearch 1.6 [Elasticsearch BV]. Details of this tf-idf based scoring  
359 method are described within the Elasticsearch documentation.<sup>7</sup> The scoring algo-  
360 rithm is based on tf-idf, but contains a few minor adjustments intended to make  
361 scores returned from different queries more comparable.

## 362 2.4 Implementation

363 The complete implementation of the evaluation presented in this paper is avail-  
364 able in the open-source software PARTYCRASHER.<sup>8</sup> The implementation includes  
365 every deduplication method we claimed to evaluate above, a general-purpose dedu-  
366 plication framework, the programs used to mine and filter the data used for the  
367 evaluation, the programs that produced the evaluation results, the raw evaluation  
368 results, and the scripts used to plot them.

## 369 2.5 Evaluation Metrics

370 Two families of evaluation metrics were used. These are the *BCubed* precision,  
371 recall, and  $F_1$ -score, and the purity, inverse purity, and  $F_1$ -score. Both are suit-  
372 able for characterizing the performance of online non-stationary clustering algo-  
373 rithms by comparing the clusters that evolve over time to clusters created by  
374 hand. A comparison of BCubed and purity, along with several other metrics, and  
375 an argument for the advantages of BCubed over purity is provided in Amigó *et al.*  
376 [Amigó et al.]. The mathematical formulae for both metrics can be found in  
377 Amigó *et al.* [Amigó et al.]. However, purity also has an advantage over BCubed:  
378 specifically that it does not require  $O(n^2)$  total time to compute whereas BCubed  
379 does.

380 If a method has a high BCubed precision, this means that there would be  
381 less chance of a developer finding unrelated crashes in the same bucket. This is  
382 important to prevent crashes caused by two unrelated bugs from sharing a bucket,

<sup>7</sup><https://www.elastic.co/guide/en/elasticsearch/guide/1.x/practical-scoring-function.html>

<sup>8</sup><https://github.com/orezpraw/partycrasher>

383 possibly causing one bug to go unnoticed since usually a developer would not  
384 examine all of the crashes in a single bucket.

385 If a method has a high BCubed recall, this means that there would be less  
386 chance of all the crashes caused by a single bug to become separated into mul-  
387 tiple buckets. Reducing the scattering of a single bug across multiple buckets  
388 is important as scattering interferes with statistics about frequently experienced  
389 bugs.

390 In contrast, purity and inverse purity focus on finding the bucket in the exper-  
391 imental results that most closely matches the bucket in the gold set. Then the  
392 overlap between the two closest matching buckets is used to compute the purity  
393 and inverse purity metrics, with high purity indicating that most of the items  
394 in a bucket produced by one of the methods evaluated are also in the matching  
395 bucket in the gold set. High recall indicates that most of the items in a bucket  
396 from the gold set are found in the matching bucket produced by the method being  
397 evaluated.

398 The purity method does not, however, completely reflect the goals of the eval-  
399 uation. Purity and inverse purity do not capture anything besides the overlap  
400 between the two buckets that overlap the most. So, if a method creates a bucket  
401 that is 51% composed of crashes from a single bug, the other 49% doesn't matter.  
402 That 49% could come from a different bug, or 200 different bugs, but the purity  
403 would be the same value. It is included in this evaluation for completeness, since  
404 it was used by Dang *et al.* [Dang et al.].

405 Both metrics can be combined into F-scores. In this evaluation, F<sub>1</sub>-scores were  
406 used, placing equal weight on precision and recall (or purity and inverse purity.)

407 BCubed and purity can be used with the gold set, hand-made buckets that  
408 are available from Ubuntu's Launchpad [Canonical Ltd.] bug tracking system.  
409 Ubuntu developers and volunteers have manually marked many of the bugs in  
410 their bug tracker as duplicates. Furthermore, many of the bugs in the bug tracker  
411 are automatically filed by Ubuntu's automated crash reporting system, Apport.  
412 This evaluation uses only bugs that were both automatically filed by Apport and  
413 manually marked as duplicates of at least one other bug. The dataset is biased to  
414 the distribution of crashes that are bucketed, which might be different than crashes  
415 that are not. Conversely, this prevents the evaluation dataset from containing any  
416 crashes that have not yet evaluated by an Ubuntu developer or volunteer.

## 417 **3 RESULTS**

418 After extracting crash reports from Launchpad, and implementing various crash  
419 report bucketing algorithms, the performance of these algorithms on the Launch-  
420 pad gold set was evaluated. Evaluation is multifaceted as in most information  
421 retrieval studies since the importance of either precision or recall are tuneable.

### 422 **3.1 BCubed and Purity**

423 Evaluation of the performance of bucketing algorithms is performed with BCubed  
424 and purity metrics. Figure 4 shows the performance of a variety of deduplication

425 methods evaluated against the entire gold set of deduplicated crash reports. The  
426 `1File` and `1Addr` methods have the most precision, while `LerchC` has the most  
427 recall.  $F_1$ -score is dominated by `Came1C` and `Lerch`. As in the results of Lerch and  
428 Mezini [Lerch and Mezini], using only the stacks outperforms using the stack plus  
429 its metadata and contextual information in terms of  $F_1$ -score. For the `Came1C`,  
430 `Lerch`, and `LerchC` simulations, a threshold of  $T = 4.0$  was used.

431 Amigó *et al.* [Amigó et al.] observed differences in BCubed and purity met-  
432 rics. Their observation was tested empirically by the evaluation. In figure 4,  
433 BCubed and purity showed similar results. The best and worst methods in terms  
434 of BCubed precision are the same as the best and worst methods in terms of  
435 purity; the same holds true for BCubed recall and inverse purity, and BCubed  
436  $F_1$ -score and purity  $F_1$ -score. However, some of the methods with intermediate  
437 performance are much closer together in purity  $F_1$ -score than they are in BCubed  
438  $F_1$ -score.

439 Figure 4 also shows that in general, if a method has a higher precision or  
440 purity, it also has a lower recall and inverse purity. For example, `3Frame` has a  
441 higher precision than `2Frame`, having a higher precision than `1Frame`, but `1Frame`  
442 has a higher recall than `2Frame` and `3Frame`.

443 The `Came1C` crash bucketing method employs: tf-idf; a tokenizer that attempts  
444 to break up identifiers such as variable names into their component words; and  
445 the entire context of the crash report including all fields reported in addition to  
446 the stack. It outperforms other bucketing methods evaluated.

### 447 3.2 Bucketing Effectiveness

448 Figure 5 shows the number of buckets created by a variety of deduplication meth-  
449 ods. The number of issue buckets extracted from the Ubuntu Launchpad gold  
450 set is plotted as the line labelled `Ubuntu`. The method that created a number of  
451 buckets most similar to the number mined from the Ubuntu Launchpad gold set  
452 was `LerchC`. For the `Lerch` and `LerchC` simulations, a threshold of  $T = 4.0$  was  
453 used.

454 Figure 6 shows the performance of the `Lerch` method when used with a va-  
455 riety of different new-bucket thresholds,  $T$ . Figure 7 shows the number of buck-  
456 ets created by the same method with those same thresholds. Since Lerch and  
457 Mezini [Lerch and Mezini] did not specify what threshold they used, this evalua-  
458 tion explored a range of thresholds. It can be seen from the plots that the relative  
459 performance of  $T$  thresholds, in terms of BCubed precision, BCubed recall, and  
460 BCubed  $F_1$ -score, becomes apparent after only 5 000 crash reports. Thus, only  
461 5 000 crash reports would need to be examined by hand for developers using the  
462 `Lerch` method to choose a suitable value for  $T$ .

463 For all the results that do not specify a value for  $T$ ,  $T = 4.0$  was used. The  
464 highest  $F_1$ -score was observed at  $T = 4.0$  after only processing 5 000 bugs with  
465 a variety of different thresholds. For `Lerch`, a threshold of  $3.5 < T < 4.5$  had the  
466 highest performance.

467 As shown in figure 8,  $T = 4.0$  still has the highest  $F_1$ -score after every crash  
468 was processed. Furthermore, other values of  $T$  near 4.0 have the same  $F_1$ -score,

469 including the range  $3.5 \leq T \leq 4.5$ . Figure 8 also shows how the threshold can be  
470 tuned to create a trade-off between precision and recall. Setting a threshold of 0.0  
471 is similar to instructing the system to put all of the crashes into a single bucket.  
472 This would be the correct choice if developers were satisfied with the explanation  
473 that all of those crashes were created by a single bug. In that case the bug would  
474 likely be filed as an issue titled, “Programs on Ubuntu Crash.” The fact that  
475 setting the threshold to 0.0 does not result in recall quite at 1.0 is an artifact  
476 of optimizations employed in ElasticSearch, specifically ElasticSearch’s inverted  
477 index.

478 Conversely, setting the threshold to 10.0 results in every crash being assigned  
479 to its own bucket, and therefore a perfect precision of 1.0. This would be the  
480 correct choice if developers considered every individual crash to be a distinct bug  
481 because the exact state of the computer was at least somewhat different during  
482 each crash. It might be more desirable to tune the value of  $T$  by using direct  
483 developer feedback rather than the technique employed here, comparing against  
484 an existing dataset. Instead of using data, one could ask developers if they had  
485 seen too many crashes caused by unrelated bugs in a single bucket recently. If  
486 they had, then  $T$  should be increased. Or,  $T$  should be decreased if developers see  
487 multiple buckets that seemed to be focused on crashes caused by the same bug.

### 488 3.3 Tokenization

489 Threshold isn’t the only way that a trade-off between precision and recall can  
490 be made. A variety of methods were tested that use the ElasticSearch/Lucene  
491 tf-idf-based search from Lerch and Mezini [Lerch and Mezini], but do not follow  
492 their tokenization strategy. The performance of several tokenization strategies is  
493 shown in figure 10. As in other cases, the methods with high precision had low  
494 recall, and the methods with high recall had low precision. All methods shown in  
495 figure 10 used a threshold of  $T = 4.0$ .

496 The **Space** method is obtained by replacing the tokenization strategy in **Lerch**  
497 with one that splits words on whitespace only, such that it does not discard any  
498 tokens regardless of how short they are, and does not lowercase every letter in the  
499 input. The **Space** method performs worse than **Lerch**. However, when both stack  
500 traces and context are used, the **SpaceC** method, performance improves slightly.  
501 This is the opposite behaviour of **Lerch**. Adding context (**LerchC**) causes perfor-  
502 mance to decrease slightly. A third tokenization strategy, **Camel** was evaluated.  
503 **Camel** attempts to break words that are written in CamelCase into their compo-  
504 nent words, using a method provided in the ElasticSearch documentation.<sup>9</sup> This  
505 strategy had the worst performance of the three, until it was used with context  
506 included, called **CamelC**. The addition of context allowed **CamelC** to outperform  
507 every other method evaluated in this paper.

508 The worst-performing tokenization evaluated, **1Addr**, was also the method that  
509 produced the largest number of buckets. However, tuning methods to match the  
510 number of buckets in the gold set without concern for performance did not result

<sup>9</sup><https://github.com/elastic/elasticsearch/blob/1.6/docs/reference/analysis/analyzers/pattern-analyzer.asciidoc>

511 in higher performance. Lerch with  $T = 3.0$  and SpaceC with  $T = 4.0$  were not  
512 the best-performing threshold or method, but both produced almost the same  
513 number of buckets as the gold set.

### 514 3.4 Runtime Performance

515 The current implementation of PARTYCRASHER requires only 45 minutes to  
516 bucket and ingest 15 293 crashes, using the slowest algorithm, CamelC, on a In-  
517 tel(R) Core(TM) i7-3770K CPU @ 3.50GHz machine with 32GiB of RAM and a  
518 Hitachi HDS723020BLE640 7200 RPM hard drive. Performance depends mainly  
519 on disk throughput, latency and RAM available for caching; Elasticsearch recom-  
520 mends using only solid-state drives. This works out to 335 crashes per minute,  
521 meeting the performance goal of 217 crashes per minute based on crash-stats from  
522 Mozilla. The performance of Elasticsearch is highly dependent on Elasticsearch's  
523 configuration settings. The settings used during these evaluations is available in  
524 the PARTYCRASHER repository.

## 525 4 DISCUSSION

### 526 4.1 Threats to Validity

527 Results are dependent on the gold set—a manual classification of crash report by  
528 Ubuntu volunteers. The results maybe biased due to the exclusive use of known  
529 duplicate crashes; the known and classified duplicates may not be representative  
530 of all crash reports. If any of these methods with with tunable parameters are  
531 deployed, the parameters should be tuned based on feedback from people working  
532 with the crash buckets, not just the gold set.

533 Since the evaluation only used data from open source software, it is unknown  
534 if our results are applicable to closed-source domains. Only stacks that originate  
535 from C and C++ projects have been evaluated; it is possible that other languages,  
536 compilers, and their runtimes have different characteristics in how they form stack  
537 traces. However, these results are corroborated by studies that examined Java  
538 exclusively [Wang et al., Lerch and Mezini].

### 539 4.2 Related work

540 Although crash bucketing facilitates manual debugging of individual faults, crash  
541 buckets are much more beneficial as the input to other methods in software engi-  
542 neering. Lerch and Mezini [Lerch and Mezini] originally applied their technique  
543 to the field of deduplicating bug, not crash, reports; Khomh *et al.* [Khomh et al.]  
544 used crash buckets to triage bugs: prioritizing developer effort on the most crucial  
545 bugs. Seo and Kim [Seo and Kim] leveraged crash buckets to predict “recurring  
546 crashes”—*i.e.*, bugs that were “fixed” but had to be fixed again in a later revi-  
547 sion. Crash buckets may also serve as input to crash localization [Liu and Han,  
548 Wang et al., Wu et al.] and crash visualization [Kim et al., Dang et al.].

### 549 4.3 Future Work

550 The results in this paper indicate that there may be a large number of improve-  
551 ments that could be made to the relatively high-performance tf-idf-based crash  
552 deduplication methods.

553 Many stack comparison methods [Modani et al., Glerum et al., Dhaliwal et al.,  
554 Wang et al.], take into account the position that each frame is on the stack, giving  
555 more weight to the frames near the top of the stack and less weight to frames on  
556 the bottom of the stack, or consider stacks that have similar frames in a similar  
557 order. The best-performing method of crash deduplication presented in this paper  
558 completely disregards information about the order of the stack. It is likely that  
559 a technique based on tf-idf that also incorporates information about the order of  
560 frames on the stack would outperform all of the methods evaluated in this paper.  
561 This could be achieved by giving words that appear in the top of the stack more  
562 weight when computing tf-idf or by re-ranking the top results produced by tf-idf  
563 according to stack similarity before choosing a bucket to place a crash in. Neither  
564 of these extensions would cause the method to be unable to scale.

565 The tokenization techniques evaluated in this paper are extremely primitive.  
566 They are merely regular expressions that break up words based on certain types of  
567 characters such as spaces, symbols, uppercase letters, lowercase letters and num-  
568 bers. Advanced tokenization techniques, such as the ones found in Guerrouj *et*  
569 *al.* [12] and Hill *et al.* [13], would likely outperform the basic techniques that have  
570 been evaluated in this paper.

571 As shown in figure 3, crash reports often contain a multitude of data apart  
572 from the stack trace itself. This paper only measured the performance of tf-idf  
573 when using only the stack trace or the entire crash report. Some fields in the crash  
574 report may be more important to obtaining a high performance than others. For  
575 example, **Architecture** (the computer architecture on which the crash occurred)  
576 might be more valuable for deduplication than **CrashCounter** (the number of  
577 times that a crash has occurred on that computer) or vice-versa, but this has not  
578 been studied in the context of information retrieval.

579 We would like to extend information retrieval techniques with more sophisti-  
580 cated normalization. We want to investigate any effects that stack normalization,  
581 as first proposed by Brodie *et al.* [Brodie et al.], would have on our tf-idf ap-  
582 proach.

583 It would be valuable to measure the effectiveness of using the buckets produced  
584 by the **CamelC** technique as input to other methods, such as those that perform  
585 bug triaging [Khomh et al.] and crash localization [Wu et al.].

## 586 5 CONCLUSION

587 The results in this paper indicate that off-the-shelf tf-idf-based information re-  
588 trieval tools can bucket crash reports in a completely unsupervised, large-scale  
589 setting when compared to variety of other previously proposed algorithms. Based  
590 on these results, a developer, such as Ada, should choose a tf-idf-based crash  
591 deduplication method with tokenization that fits their dataset, and intermediate

592 new-bucket threshold. They should update this threshold based on feedback from  
593 developers, volunteers, or employees that work with the stack traces directly. A  
594 tf-idf approach that used the entire crash report and stack trace, tokenized using  
595 camel-case had the best  $F_1$ -score on the Ubuntu Launchpad crash reports used in  
596 this work. In addition, there is a lot of room for improvements to these techniques.  
597 This conclusion is surprising in light of the fact that the tf-idf-based techniques  
598 evaluated disregard information that is often considered to be essential to stack  
599 traces, such as the order of the frames in the stack.

600 Finally the research questions can be answered:

**RQ1:** tf-idf-based methods are effective, industrial-scale methods of crash report bucketing.

601 **RQ2:** New-bucket thresholds and tokenization strategies can be tuned to increase precision and recall.

## 602 ACKNOWLEDGEMENTS

603 The authors would like to thank the Mozilla foundation, especially Robert Helmer,  
604 Adrian Gaudebert, Peter Bengtsson and Chris Lonnen for their help, and for  
605 making Mozilla's massive collection of stack reports open and publicly available.  
606 Additionally, the authors would like to thank the Ubuntu project and all of its  
607 developers who manually deduplicated bug reports that were submitted with crash  
608 reports. Funding for this research was provided by MITACS Accelerate with  
609 BioWare™, an Electronic Arts Inc. studio.

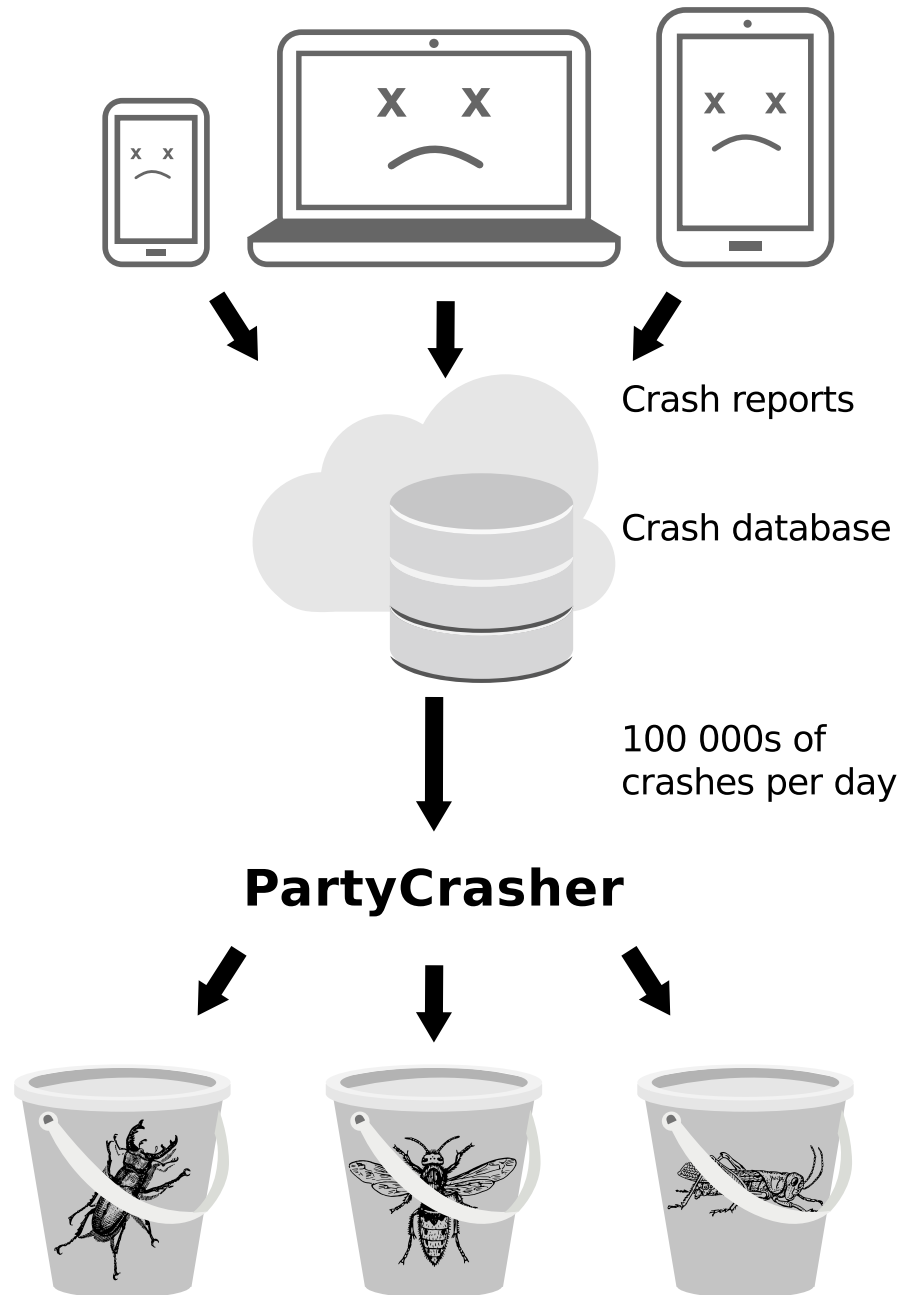
## 610 REFERENCES

- 611 [soc] mozilla/socorro: Socorro is a server to accept and process Breakpad crash  
612 reports.
- 613 [app] Technical Note TN2123: CrashReporter.
- 614 [Ahmed et al.] Ahmed, I., Mohan, N., and Jensen, C. The Impact of Automatic  
615 Crash Reports on Bug Triaging and Development in Mozilla. In *Proceedings  
616 of The International Symposium on Open Collaboration, OpenSym '14*, pages  
617 1:1–1:8. ACM.
- 618 [Amigó et al.] Amigó, E., Gonzalo, J., Artiles, J., and Verdejo, F. A comparison of  
619 extrinsic clustering evaluation metrics based on formal constraints. 12(4):461–  
620 486.
- 621 [Bartz et al.] Bartz, K., Stokes, J. W., Platt, J. C., Kivett, R., Grant, D., Calinoiu,  
622 S., and Loihle, G. Finding Similar Failures Using Callstack Similarity. In  
623 *SysML*.
- 624 [Brodie et al.] Brodie, M., Ma, S., Lohman, G., Mignet, L., Wilding, M., Cham-  
625 plin, J., and Sohn, P. Quickly Finding Known Software Problems via Auto-  
626 mated Symptom Matching. In *Second International Conference on Autonomic  
627 Computing, 2005. ICAC 2005. Proceedings*, pages 101–110.
- 628 [Canonical Ltd.] Canonical Ltd. Launchpad.
- 629 [Dang et al.] Dang, Y., Wu, R., Zhang, H., Zhang, D., and Nobel, P. ReBucket:  
630 a method for clustering duplicate crash reports based on call stack similarity.



- 631 In *Proceedings of the 34th International Conference on Software Engineering*,  
632 pages 1084–1093. IEEE Press.
- 633 [Dhaliwal et al.] Dhaliwal, T., Khomh, F., and Zou, Y. Classifying field crash  
634 reports for fixing bugs: A case study of Mozilla Firefox. In *2011 27th IEEE  
635 International Conference on Software Maintenance (ICSM)*, pages 333–342.
- 636 [Elasticsearch BV] Elasticsearch BV. Elasticsearch.
- 637 [Glerum et al.] Glerum, K., Kinshumann, K., Greenberg, S., Aul, G., Orgovan,  
638 V., Nichols, G., Grant, D., Loihle, G., and Hunt, G. Debugging in the (Very)  
639 Large: Ten Years of Implementation and Experience. In *Proceedings of the  
640 ACM SIGOPS 22Nd Symposium on Operating Systems Principles, SOSP '09*,  
641 pages 103–116. ACM.
- 642 [12] Guerrouj, L., Di Penta, M., Antoniol, G., and Guéhéneuc, Y.-G. (2013). Ti-  
643 dier: an identifier splitting approach using speech recognition techniques. *Jour-  
644 nal of Software: Evolution and Process*, 25(6):575–599.
- 645 [13] Hill, E., Binkley, D., Lawrie, D., Pollock, L., and Vijay-Shanker, K. (2014).  
646 An empirical study of identifier splitting techniques. *Empirical Software Engi-  
647 neering*, 19(6):1754–1780.
- 648 [Khomh et al.] Khomh, F., Chan, B., Zou, Y., and Hassan, A. An Entropy Eval-  
649 uation Approach for Triaging Field Crashes: A Case Study of Mozilla Firefox.  
650 In *2011 18th Working Conference on Reverse Engineering (WCRE)*, pages 261–  
651 270.
- 652 [Kim et al.] Kim, S., Zimmermann, T., and Nagappan, N. Crash graphs: An  
653 aggregated view of multiple crashes to improve crash triage. In *2011 IEEE/IFIP  
654 41st International Conference on Dependable Systems Networks (DSN)*, pages  
655 486–493.
- 656 [Lerch and Mezini] Lerch, J. and Mezini, M. Finding Duplicates of Your Yet Un-  
657 written Bug Report. In *2013 17th European Conference on Software Mainte-  
658 nance and Reengineering (CSMR)*, pages 69–78.
- 659 [Liblit et al.] Liblit, B., Naik, M., Zheng, A. X., Aiken, A., and Jordan, M. I.  
660 Scalable statistical bug isolation. In *ACM SIGPLAN Notices*, volume 40, pages  
661 15–26. ACM.
- 662 [Liu and Han] Liu, C. and Han, J. Failure Proximity: A Fault Localization-based  
663 Approach. In *Proceedings of the 14th ACM SIGSOFT International Symposium  
664 on Foundations of Software Engineering, SIGSOFT '06/FSE-14*, pages 46–56.  
665 ACM.
- 666 [Liu et al.] Liu, C., Yan, X., Fei, L., Han, J., and Midkiff, S. P. SOBER: Statistical  
667 Model-based Bug Localization. In *Proceedings of the 10th European Software  
668 Engineering Conference Held Jointly with 13th ACM SIGSOFT International  
669 Symposium on Foundations of Software Engineering, ESEC/FSE-13*, pages 286–  
670 295. ACM.
- 671 [20] Marr, B. (2015). Why only one of the 5 vs of big  
672 data really matters. [http://www.ibmbigdatahub.com/blog/  
673 why-only-one-5-vs-big-data-really-matters](http://www.ibmbigdatahub.com/blog/why-only-one-5-vs-big-data-really-matters).
- 674 [Modani et al.] Modani, N., Gupta, R., Lohman, G., Syeda-Mahmood, T., and  
675 Mignet, L. Automatically Identifying Known Software Problems. In *2007 IEEE  
676 23rd International Conference on Data Engineering Workshop*, pages 433–441.

- 677 [Mozilla Corporation] Mozilla Corporation. Mozilla Crash Reports.
- 678 [23] Robles, G., González-Barahona, J. M., Izquierdo-Cortazar, D., and Herraiz,  
679 I. (2011). Tools and datasets for mining libre software repositories. *Multi-*  
680 *Disciplinary Advancement in Open Source Software and Processes*, page 24.
- 681 [Salton and McGill] Salton, G. and McGill, M. J. *Introduction to modern informa-*  
682 *tion retrieval*. McGraw-Hill computer science series. McGraw-Hill.
- 683 [Schröter et al.] Schröter, A., Bettenburg, N., and Premraj, R. Do stack traces  
684 help developers fix bugs? In *2010 7th IEEE Working Conference on Mining*  
685 *Software Repositories (MSR)*, pages 118–121.
- 686 [Seo and Kim] Seo, H. and Kim, S. Predicting Recurring Crash Stacks. In *Proceed-*  
687 *ings of the 27th IEEE/ACM International Conference on Automated Software*  
688 *Engineering, ASE 2012*, pages 180–189. ACM.
- 689 [Wang et al.] Wang, S., Khomh, F., and Zou, Y. Improving bug localization using  
690 correlations in crash reports. In *2013 10th IEEE Working Conference on Mining*  
691 *Software Repositories (MSR)*, pages 247–256.
- 692 [Wu et al.] Wu, R., Zhang, H., Cheung, S.-C., and Kim, S. Crashlocator: Locating  
693 crashing faults based on crash stacks. In *Proceedings of the 2014 International*  
694 *Symposium on Software Testing and Analysis*, pages 204–214. ACM.



**Figure 2.** PARTYCRASHER within a development context

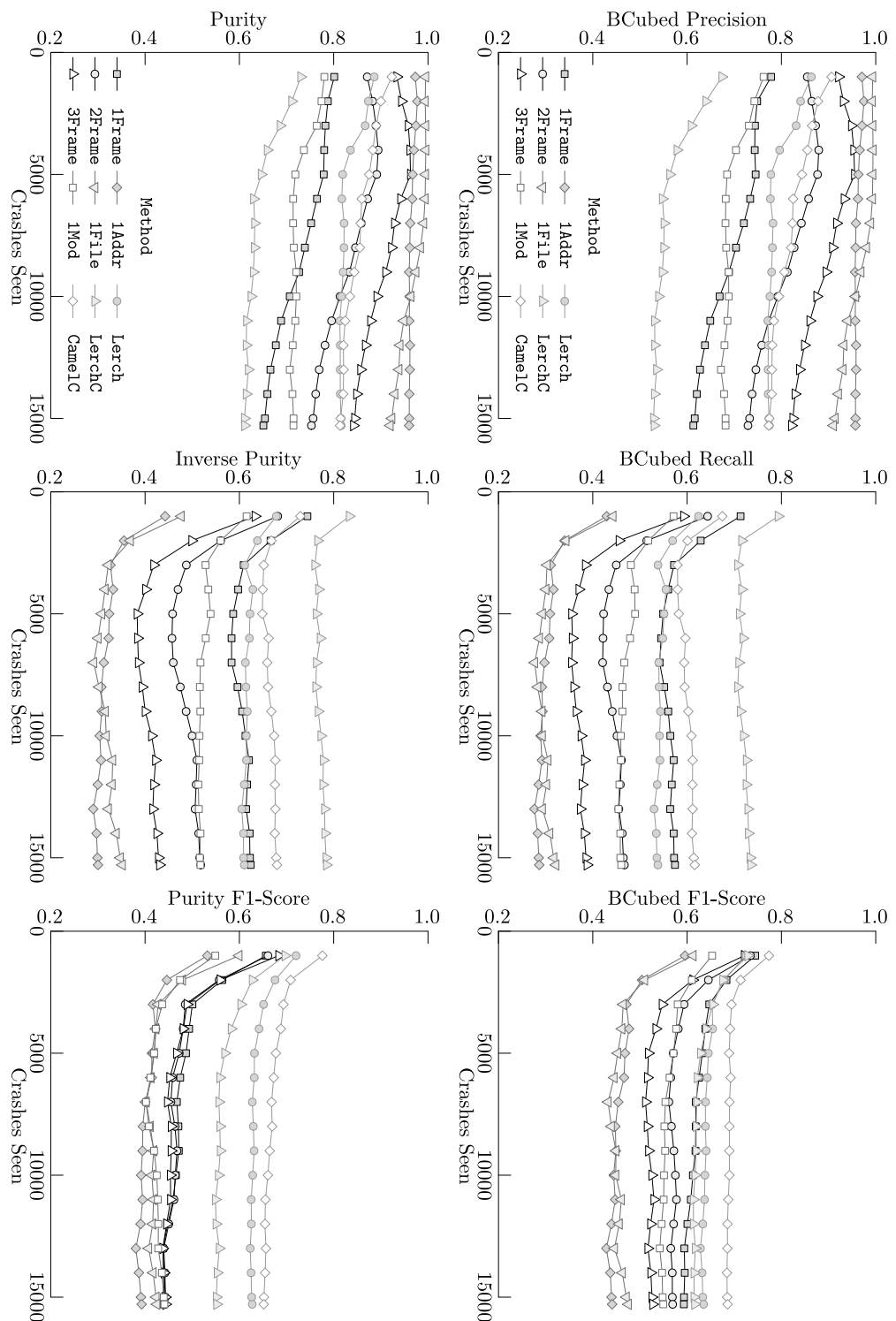
```
Binary package hint: evolution-exchange

I just start Evolution, wait about 2 minutes, and then evolution-exchange crashed

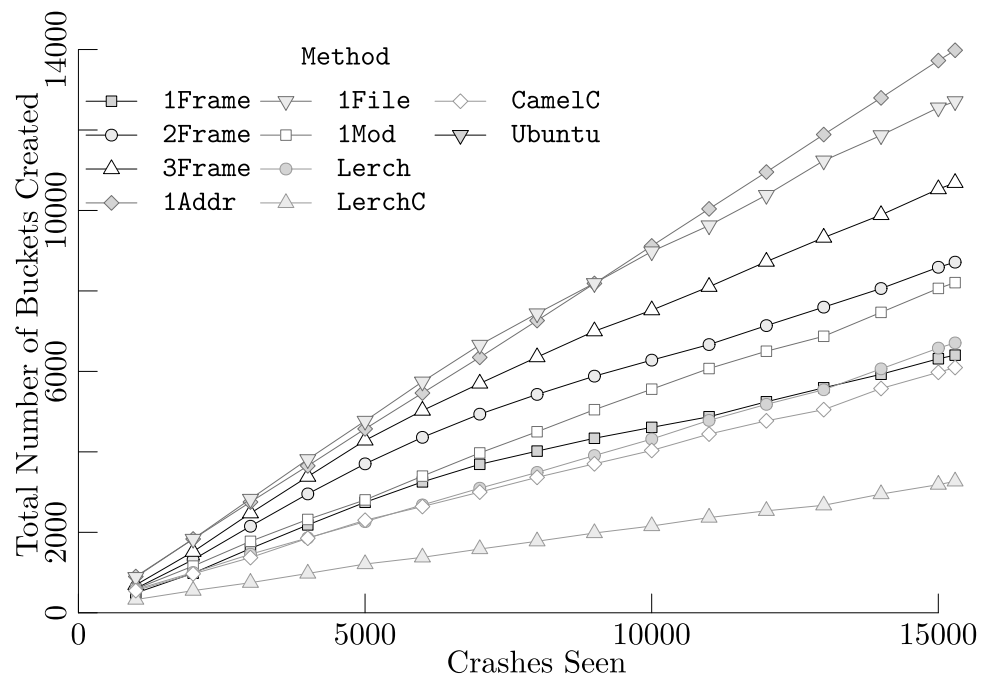
ProblemType: Crash
Architecture: i386
CrashCounter: 1
Date: Tue Jul 17 10:09:50 2007
DistroRelease: Ubuntu 7.10
ExecutablePath: /usr/lib/evolution/2.12/evolution-exchange-storage
NonfreeKernelModules: vmnet vmmon
Package: evolution-exchange 2.11.5-0ubuntu1
PackageArchitecture: i386
ProcCmdline: /usr/lib/evolution/2.12/evolution-exchange-storage --oaf-activate-i
ProcCwd: /
ProcEnviron:
 PATH=/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games
 LANG=en_US.UTF-8
 SHELL=/bin/bash
Signal: 11
SourcePackage: evolution-exchange
Title: evolution-exchange-storage crashed with SIGSEGV in soup_connection_discon
Uname: Linux encahl 2.6.20-15-generic #2 SMP Sun Apr 15 07:36:31 UTC 2007 i686 G
UserGroups: adm admin audio cdrom dialout dip floppy kqemu lpadmin netdev plugde

#0  0xb71e8d92 in soup_connection_disconnect () from /usr/lib/libsoup-2.2.so.8
#1  0xb71e8dfd in ?? () from /usr/lib/libsoup-2.2.so.8
#2  0x080e5a48 in ?? ()
#3  0xb6eaf678 in ?? () from /usr/lib/libgobject-2.0.so.0
#4  0xbfd613e8 in ?? ()
#5  0xb6e8b179 in g_cclosure_marshal_VOID__VOID ()
    from /usr/lib/libgobject-2.0.so.0
Backtrace stopped: frame did not save the PC
```

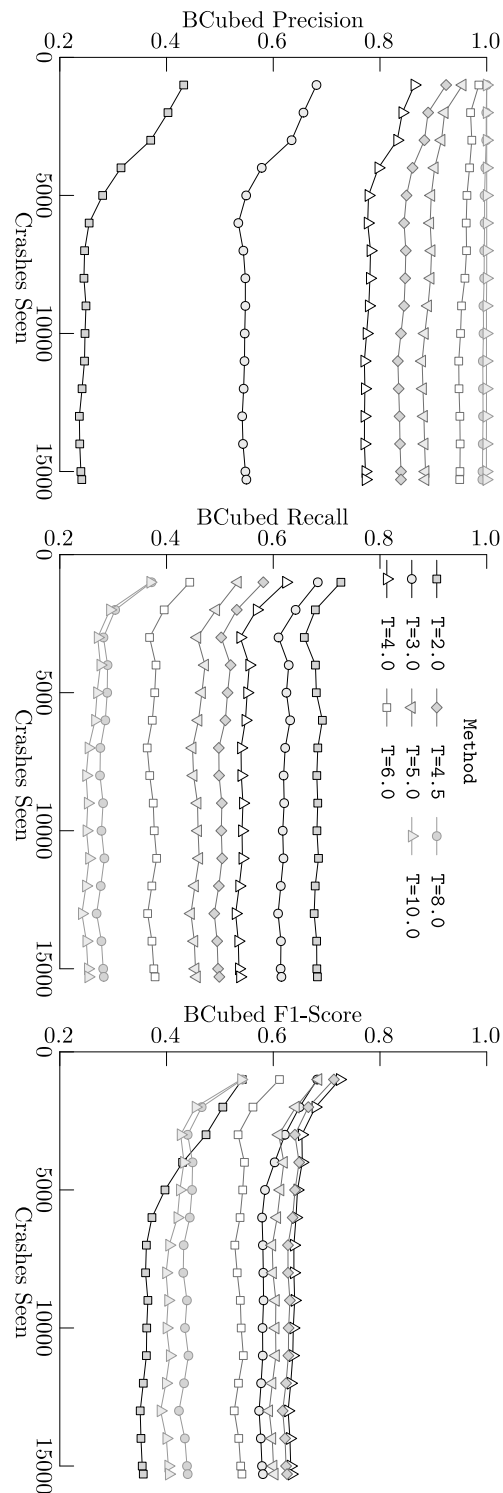
**Figure 3.** An example crash report, including stack.



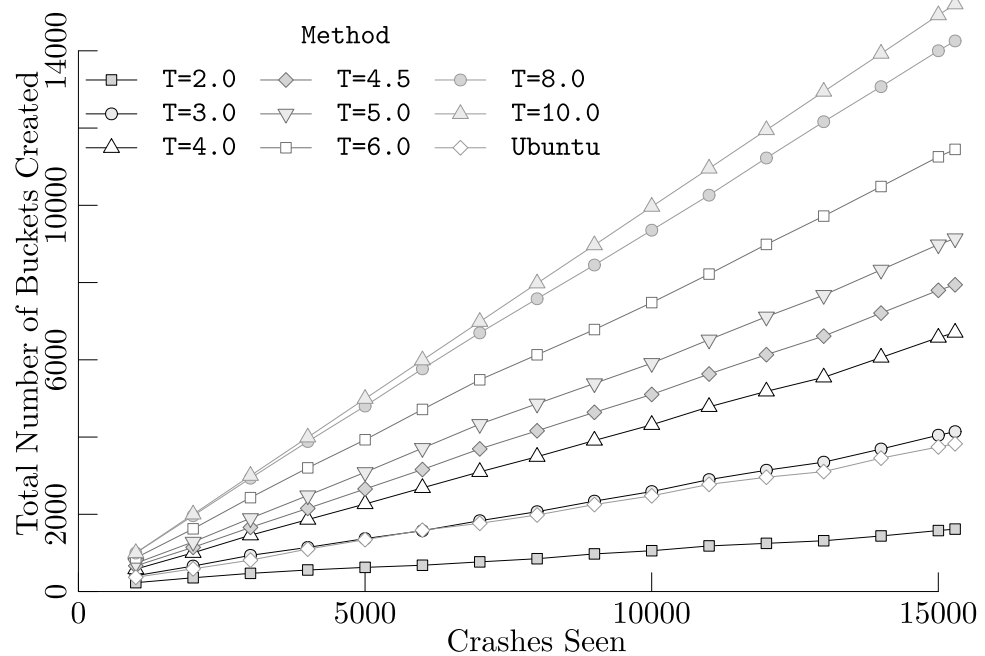
**Figure 4.** BCubed (top) and Purity-metric (bottom) scores for various methods of crash report deduplication.



**Figure 5.** Number of buckets created as a function of number of crashes seen. The line labelled **Ubuntu** indicates the number of groups crashes that were marked as duplicates of each other by Ubuntu developers or volunteers.

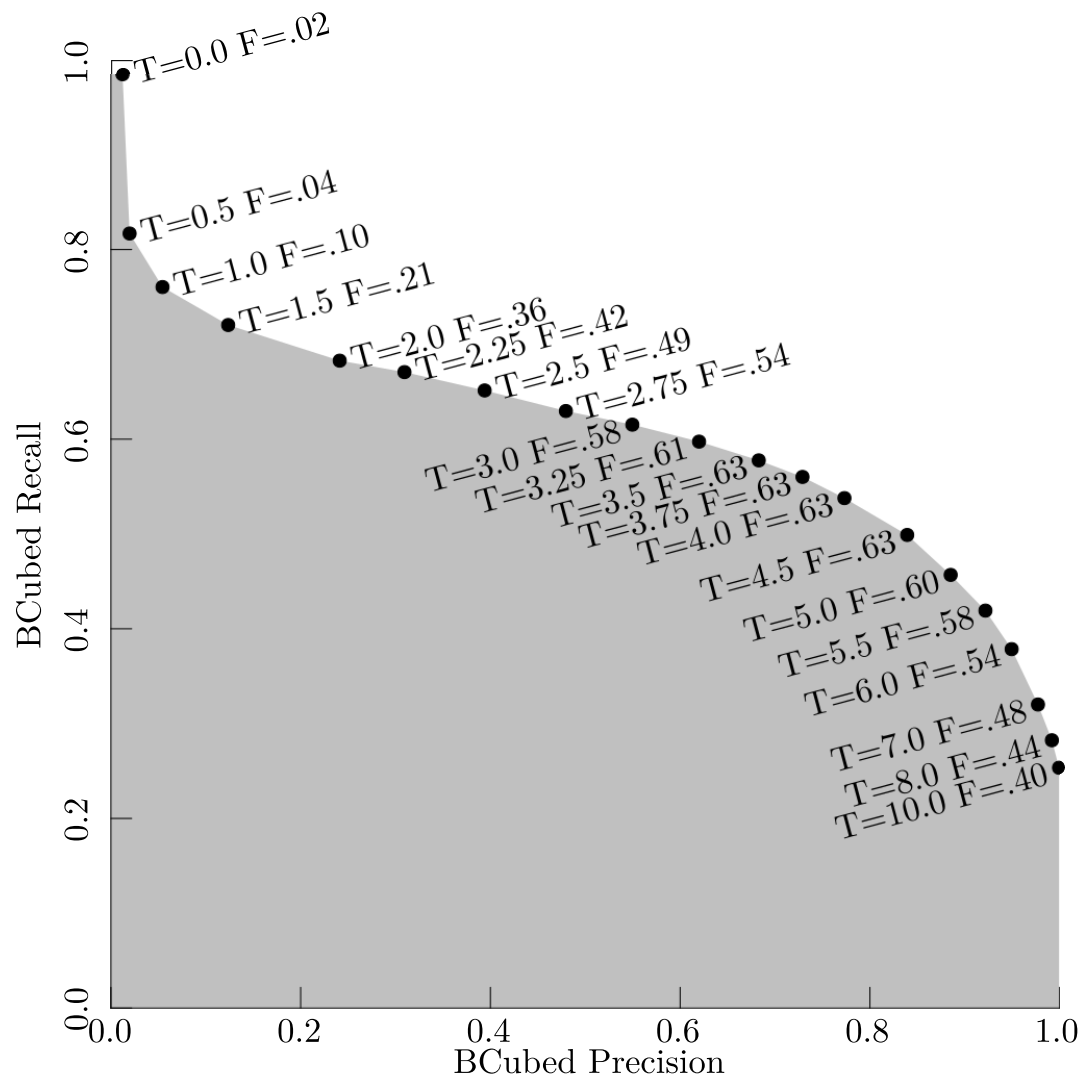


**Figure 6.** BCubed scores for the Lerch method of crash report deduplication at various new-bucket thresholds  $T$ .

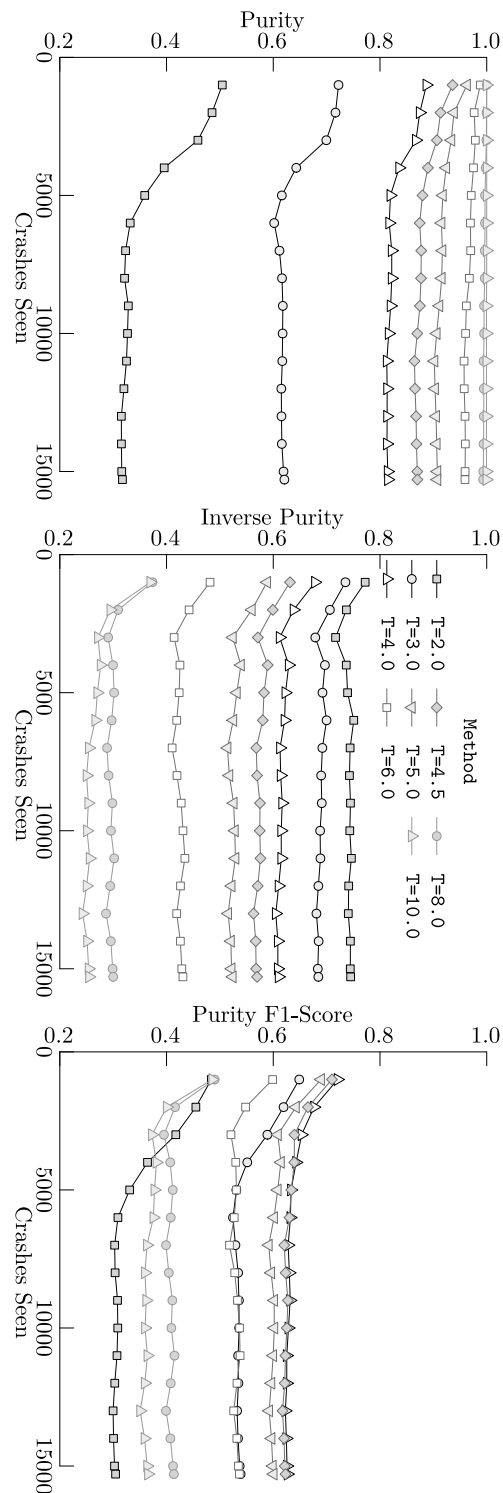


**Figure 7.** Number of buckets created as a function of number of crashes seen for the Lerch method of crash report deduplication at various new-bucket thresholds  $T$ . The line labelled **Ubuntu** indicates the number of groups crashes that were marked as duplicates of each other by Ubuntu developers or volunteers.

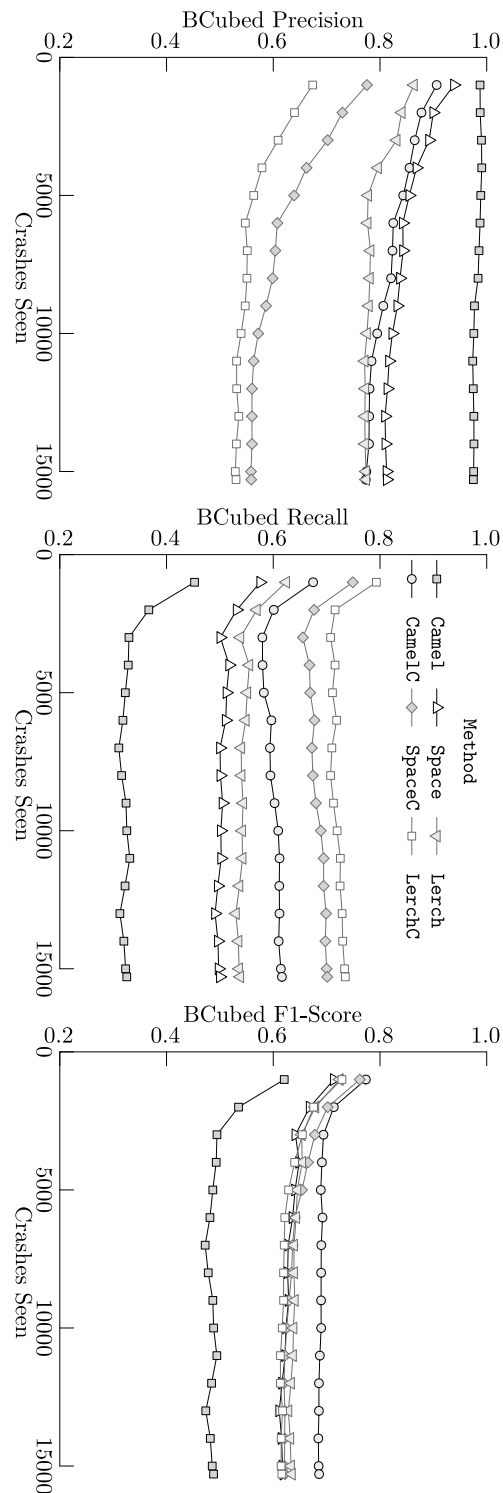




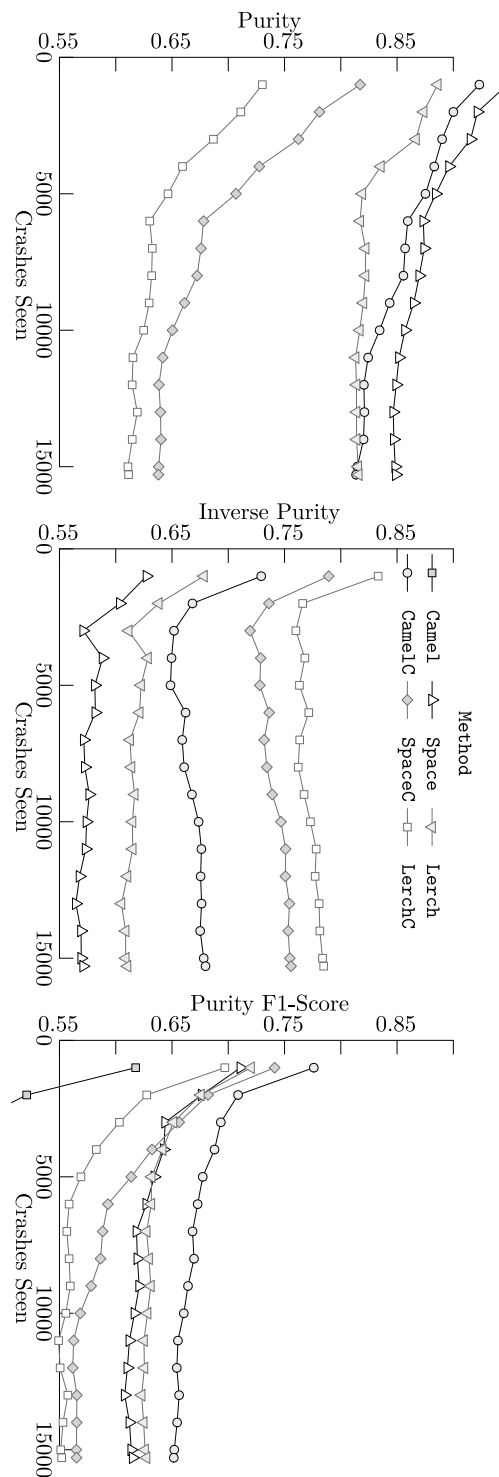
**Figure 8.** Precision/Recall plot showing the trade-off between BCubed precision and recall as the new-bucket threshold  $T$  is adjusted. BCubed  $F_1$ -score is also listed in the plot.



**Figure 9.** Purity-metric scores for the Lerch method of crash report deduplication at various new-bucket thresholds  $T$ .



**Figure 10.** BCubed scores for the Lerch method of crash report deduplication with Lerch's tokenization technique replaced by a variety of other techniques.



**Figure 11.** Purity-metric scores for the tf-idf-based methods of crash report deduplication with various tokenization strategies.