



A language-independent software renovation framework

M. Di Penta^{a,*}, M. Neteler^b, G. Antoniol^a, E. Merlo^c

^a Department of Engineering, RCOST—Research Centre on Software Technology, University of Sannio, Via Traiano, 1-82100 Benevento, Italy

^b ITC-irst Istituto Trentino Cultura, Via Sommarive, 18-38050 Povo (Trento), Italy

^c École Polytechnique de Montréal, Montréal, Québec, Canada

Received 1 April 2003; received in revised form 16 July 2003; accepted 2 March 2004

9 Abstract

10 One of the undesired effects of software evolution is the proliferation of unused components, which are not used by any appli-
11 cation. As a consequence, the size of binaries and libraries tends to grow and system maintainability tends to decrease. At the same
12 time, a major trend of today's software market is the porting of applications on hand-held devices or, in general, on devices which
13 have a limited amount of available resources. Refactoring and, in particular, the miniaturization of libraries and applications are
14 therefore necessary.

15 We propose a Software Renovation Framework (SRF) and a toolkit covering several aspects of software renovation, such as
16 removing unused objects and code clones, and refactoring existing libraries into smaller more cohesive ones. Refactoring has been
17 implemented in the SRF using a hybrid approach based on hierarchical clustering, on genetic algorithms and hill climbing, also tak-
18 ing into account the developers' feedback. The SRF aims to monitor software system quality in terms of the identified affecting fac-
19 tors, and to perform renovation activities when necessary. Most of the framework activities are language-independent, do not
20 require any kind of source code parsing, and rely on object module analysis.

21 The SRF has been applied to *GRASS*, which is a large open source Geographical Information System of about one million LOCs
22 in size. It has significantly improved the software organization, has reduced by about 50% the average number of objects linked by
23 each application, and has consequently also reduced the applications' memory requirements.

24 © 2004 Elsevier Inc. All rights reserved.

25 *Keywords:* Refactoring; Software renovation; Clustering; Genetic algorithms; Hill climbing

27 1. Introduction

28 Software systems evolution often presents several fac-
29 tors that contribute to deteriorate the quality of the sys-
30 tem itself (Lehman and Belady, 1985). First, unused
31 components, which have been introduced for testing
32 purposes or which belong to obsolete functionalities,
33 may proliferate. Second, maintenance and evolution
34 activities are likely to introduce clones, while, for exam-

35 ple, adding support and drivers for an architecture sim- 35
36 ilar to an already supported one (Antoniol et al., 2002). 36
37 Third, library sizes tend to increase, because new func- 37
38 tionalities are added and refactoring is rarely performed; 38
39 for the same reasons, also the number of inter-library 39
40 dependencies, some of which are circular, tends to in- 40
41 crease. Finally, sometimes, new functionalities logically 41
42 related to already existing ones are added in a non-sys- 42
43 tematic way and they result in sets of modules which 43
44 are neither organized nor linked into libraries. As a con- 44
45 sequence, systems become difficult to maintain. Moreo- 45
46 ver, unused objects, big libraries, and circular 46
47 dependencies significantly increase application sizes 47
48 and memory requirements. This is clearly in contrast 48

* Corresponding author.

E-mail addresses: dipenta@unisannio.it (M. Di Penta), neteler@itc.it (M. Neteler), antonio@jeec.org (G. Antoniol), merlo@info.polymtl.ca (E. Merlo).

49 with today's industry hype towards porting existing soft- 105
50 ware applications onto hand-held devices, such as Per- 106
51 sonal Digital Assistants (PDA), onto wireless devices 107
52 (e.g., multimedia cell phones), or, in general, onto de- 108
53 vices with limited resources. 109

54 This paper proposes the SRF to monitor and control 110
55 some of the quality factors which have been described 111
56 above. When the number of unused objects and clones 112
57 increase, or when library sizes become unmanageable, 113
58 some actions may be taken among the several possible 114
59 ones. First and foremost, unused code may be removed 115
60 and clones may be monitored or factored out. Further- 116
61 more, some form of restructuring, at library and at ob- 117
62 ject file level, may be required. Together with 118
63 monitoring and improving maintainability, the SRF 119
64 eases the miniaturization challenge of porting applica- 120
65 tions onto limited resources devices. 121

66 Most of the SRF activities deal with analyzing 122
67 dependencies among software artifacts. For any given 123
68 software system, dependencies among executables and 124
69 object files may be represented via a dependency graph, 125
70 which is a graph where nodes represent resources and 126
71 edges represent dependencies. Each library, in turn, 127
72 may be thought of as a subgraph in the overall object file 128
73 dependency graph. Therefore, software miniaturization 129
74 can be modeled as a graph partitioning problem. Unfor- 130
75 tunately, it is well known that graph partitioning is an 131
76 NP-hard problem (Garey and Johnson, 1979) and thus 132
77 heuristics have been adopted to find a "good-enough" 133
78 solution. For example, one may be interested to first 134
79 examine graph partitions by minimizing cross edges be- 135
80 tween subgraphs which correspond to libraries. More 136
81 formally, a cost function describing the restructuring 137
82 problem has to be defined and heuristics to drive the 138
83 solution search process must be identified and applied. 139

84 We propose a novel approach in which hierarchical 140
85 clustering and *Silhouette* statistics (Kaufman and Rous- 141
86 seeuw, 1990) are initially used to determine the optimal 142
87 number of clusters and the starting population of a Soft- 143
88 ware Renovation Genetic Algorithm (SRGA). This initial 144
89 step is followed by a SRGA search aimed at 145
90 minimizing a multi-objective function which takes into 146
91 account, at the same time, both the number of inter-li- 147
92 brary dependencies and the average number of objects 148
93 linked by each application. Finally, by letting the SRGA 149
94 fitness function also consider the experts' suggestions, 150
95 the SRF becomes a semi-automatic approach composed 151
96 of multiple refactoring iterations, which are interleaved 152
97 by developers' feedback. To speed up the search process, 153
98 heuristics based on a Genetic Algorithm (GA) and a 154
99 modified GA (Talbi and Bessière, 1991) approach were 155
100 proposed. Performance improvement was also achieved 156
101 by means of a hybrid approach, which combines GA 157
102 strategies with hill climbing techniques. 158

103 The SRF has the advantage of being language inde-
104 pendent. All activities, except clone detection, rely on

information extracted from object files; furthermore, 105
the clone detection algorithm adopted in the SRF is 106
not tied to any specific programming language, provided 107
that a set of metrics can be extracted from the source 108
code. 109

The SRF has been applied to a large Open Source 110
software system: a Geographical Information System 111
(GIS) named *GRASS*¹ (Geographic Resources Analysis 112
Support System). *GRASS* is a raster/vector GIS com- 113
bined with integrated image processing and data visual- 114
ization subsystems (Neteler and Mitasova, 2002) 115
composed of 517 applications and 43 libraries, for a to- 116
tal of over one million LOCs. 117

The number of team members is small and it is about 118
7–15 active developers. Decisions are usually taken by 119
the members most capable to solve specific problems. 120
Developers are also *GRASS* users and they often focus 121
on their needs within the general project. 122

This paper is organized as follows. First, a short re- 123
view on related work (Section 2) and on main notions 124
of clustering and GAs (Section 3), will be presented. 125
Then, the SRF is presented in Section 4. The case study 126
software system (i.e., *GRASS*) is described in Section 5, 127
while results are presented and discussed in Section 6, 128
and are followed by conclusions and work-in-progress 129
in Section 7. 130

2. Related work 131

Many research contributions have been published 132
about software system modules clustering and restruc- 133
turing, identifying objects, and recovering or building 134
libraries. Most of these work applied clustering or Con- 135
cept Analysis (CA). 136

An overview of CA applications to software reengi- 137
neering problems was published by G. Snelting in his 138
seminal work (Snelting, 2000). Snelting applied CA to 139
several modularization problems such as exploring 140
configuration spaces (see also Krone and Snelting, 141
1994), transforming class hierarchies, and modulariz- 142
ing COBOL systems. Kuipers and Moonen (2000) com- 143
bined CA and type inference in a semi-automatic 144
approach to find objects in COBOL legacy code. Anto- 145
niol et al. (2001a) applied CA to the problem of identi- 146
fying libraries and of defining new directories and files 147
organizations in software systems with degraded archi- 148
tectures. As according to Krone and Snelting (1994), 149
Kuipers and Moonen (2000), and Antoniol et al. 150
(2001a), we believe that with the present level of technol- 151
ogy a programmer-centric approach is required, since 152
programmers are in charge of choosing the proper 153
modularization strategy based on their knowledge 154

¹ <http://grass.itc.it>

155 and judgment. A comparison between clustering and
 156 CA was presented by [Kuipers and van Deursen \(1999\)](#).
 157 Our work also applies an agglomerative-nesting cluster-
 158 ing to a Boolean usage matrix, although according to
 159 [Kuipers and van Deursen \(1999\)](#) the matrix indicated
 160 the uses of variables by programs.

161 Surveys and overviews of cluster analysis applied to
 162 software systems have been published in the past, for
 163 example, by [Wiggerts \(1997\)](#) and by [Tzerpos and Holt
 164 \(1998\)](#). The latter authors ([Tzerpos and Holt, 1999](#)) de-
 165 fined a metric to evaluate the similarity of different
 166 decompositions of software systems. [Tzerpos and Holt
 167 \(2000a\)](#) proposed a novel clustering algorithm which
 168 had been specifically conceived to address the peculiari-
 169 ties of the program comprehension; they also addressed
 170 the issue of stability of software clustering algorithms
 171 ([Tzerpos and Holt, 2000b](#)). Applications of clustering
 172 to reengineering were suggested by [Anquetil and Leth-
 173 bridge \(1998\)](#), that devised a method for decomposing
 174 complex software systems into independent subsystems.
 175 Source files were clustered according to file names and
 176 their name decomposition. An approach relying on in-
 177 ter-module and intra-module dependency graphs to
 178 refactor software systems was presented by [Mancoridis
 179 et al. \(1998\)](#). We share the idea of analyzing dependency
 180 graphs and of finding a tradeoff between highly cohesive
 181 and little inter-connected libraries, with [Mancoridis
 182 et al. \(1998\)](#).

183 GAs have been recently applied in different fields of
 184 computer science and software engineering. An ap-
 185 proach for partitioning a graph using GAs was dis-
 186 cussed by [Talbi and Bessi re \(1991\)](#). Similar
 187 approaches were also published by [Shazely et al.
 188 \(1998\)](#), [Bui and Moon \(1996\)](#), and [Oommen and de St
 189 Croix \(1996\)](#). [Maini et al. \(1994\)](#) discussed a method
 190 to introduce knowledge about the problem in a non-uni-
 191 form crossover operator and presented some examples
 192 of its application. A GA was used by [Doval et al.
 193 \(1999\)](#) to identify clusters on software systems. Together
 194 with [Doval et al., 1999](#), we share the idea of a software
 195 clustering approach which uses a GA and which tries to
 196 minimize inter-cluster dependencies. Finally, [Harman
 197 and et al. \(2002\)](#) reported experiments of modularization
 198 and remodularization by comparing GAs with hill
 199 climbing techniques and by introducing a representation
 200 and a crossover operator tied to the remodularization
 201 problem. Their case studies revealed that hill climbing
 202 outperformed GAs. [Mahdavi et al. \(2003\)](#) proposed an
 203 approach aimed to combine multiple hill climbs for sub-
 204 sequent searches, thus reducing the search spaces.

205 Software miniaturization for Java application was re-
 206 cently addressed by *Jax* which is an application extrac-
 207 tor for Java software systems ([Tip et al., 1999](#)) whose
 208 goal is the size reduction of Java programs with partic-
 209 ular interest to applets to be transmitted over the net-
 210 work. *Jax* is based on transformations including

removal of redundant methods and fields, devirtualiza- 211
 tion and inlining of method calls, renaming methods, 212
 fields, class and packages, and transforming class hierar- 213
 chies. Another approach, devoted to reduce the size of 214
 Java libraries for embedded systems, was proposed by 215
[Rayside and Kontogiannis \(2002\)](#). While the approach 216
 proposed by [Rayside and Kontogiannis \(2002\)](#) and 217
Jax are tied to a programming language, ours is not. 218
 Our approach also differs from *Jax* in philosophy since 219
 we do not limit ourselves to reduce the size of the in- 220
 stance application to be executed, but we also support 221
 the reorganization of a software system whose structure 222
 has been deteriorated because of its evolution. The 223
 reduction of memory requirements is thus just one of 224
 the effects of the reorganization. 225

This paper extends preliminary contributions ([Di
 226 Penta et al., 2002](#); [Antoniol et al., 2003](#)). Together with
 227 [Di Penta et al. \(2002\)](#), we share the choice *GRASS* as
 228 target application and several activities carried out to
 229 refactor libraries. 230

231 3. Background notions

The fundamental activity of the SRF is library refac- 232
 toring. This requires the integration of clustering and 233
 GA techniques in a semi-automatic, human-driven proc- 234
 ess. Clustering deals with the grouping of large amounts 235
 of things (*entities*) in groups (*clusters*) of closely related 236
 entities ([Kaufman and Rousseeuw, 1990](#); [Anderberg,
 237 1973](#)). Clustering is used in different areas, such as busi- 238
 ness analysis, economics, astronomy, information retrieval, 239
 image processing, pattern recognition, biology, and 240
 others. GAs come from an idea, born over 30 years ago, 241
 of applying the biological principle of evolution to arti- 242
 ficial systems. GAs are applied to different domains such 243
 as machine and robot learning, economics, operations 244
 research, ecology, studies of evolution, learning and so- 245
 cial systems ([Goldberg, 1989](#); [Mitchell, 1996](#)). 246

In the following subsections, for sake of complete- 247
 ness, only some essential notions are summarized, be- 248
 cause describing the different types of clustering 249
 algorithms or the details of GAs is out of the scope of 250
 this paper. More details can be found in [Anderberg
 251 \(1973\)](#) for clustering and in [Goldberg \(1989\)](#) and [Mitc-
 252 hell \(1996\)](#) for GAs. 253

254 3.1. Agglomerative hierarchical clustering

In this paper, the agglomerative-nesting (*Agnes*) algo- 255
 rithm ([Kaufman and Rousseeuw, 1990](#)) was applied to 256
 build the initial set of *candidate libraries*. *Agnes* is an 257
 agglomerative, hierarchical clustering algorithm: it 258
 builds a hierarchy of clusters in such way that each level 259
 contains the same clusters as the first lower level, except 260

261 for two clusters, which are joined to form a single
262 cluster.

263 3.2. Determining the optimal number of clusters

264 To determine the actual or optimal number of clus-
265 ters, people traditionally rely on the plot of an error
266 measure representing the dispersion within a cluster.
267 The error measure decreases as the number of clusters,
268 k , increases, but for some values of k the curve flattens.
269 Kaufman and Rousseeuw (1990) proposed the *Silhou-*
270 *ette* statistics for estimating and assessing the optimal
271 number of clusters. For the observation i , let $a(i)$ be
272 the average distance to the other points in its cluster,
273 and $b(i)$ the average distance to points in the nearest
274 cluster. Then the *Silhouette* statistics is defined as
275

$$276 s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}. \quad (1)$$

278 Kaufman and Rousseeuw suggested choosing the optimal
279 number of clusters as the value maximizing the average
280 $s(i)$ over the dataset. Traditionally, it is assumed that the
281 error curve *knee* indicates the appropriate number of
282 clusters (Gordon, 1988).

283 Often, a compromise has to be accepted between max-
284 imizing the *Silhouette* (and thus having highly cohesive
285 clusters) and obtaining an excessive number of clusters
286 (that in our application, causes library fragmentation).

287 3.3. Genetic algorithms

288 Applications based on GAs revealed their effective-
289 ness in finding approximate solutions when the search
290 space is large or complex, when mathematical analysis
291 or traditional methods are not available, and, in general,
292 when the problem to be solved is NP-complete or NP-
293 hard (Garey and Johnson, 1979). Roughly speaking, a
294 GA may be defined as an iterative procedure that
295 searches for the best solution of a given problem among
296 a constant-size population, represented by a finite string
297 of symbols, the *genome*. The search is made starting
298 from an initial population of individuals, often ran-
299 domly generated. At each evolutionary step, individuals
300 are evaluated using a *fitness function*. High-fitness indi-
301 viduals will have the highest probability to reproduce
302 themselves.

303 The evolution (i.e., the generation of a new popula-
304 tion) is made by means of two kinds of operator: the
305 *crossover operator* and the *mutation operator*. The cross-
306 over operator takes two individuals (the *parents*) of the
307 old generation and exchanges parts of their genomes,
308 producing one or more new individuals (the *offspring*).
309 The mutation operator has been introduced to prevent
310 convergence to local optima and it randomly modifies
311 an individual's genome, for example, by flipping some
312 of its bits if the genome is represented by a bit string.

Crossover and mutation are respectively performed on
each individual of the population with probability
 p_{cross} and p_{mut} respectively, where $p_{mut} \ll p_{cross}$.

GAs are not guaranteed to converge. The termination
condition is often based on a maximum number of gen-
erations or on a given value of the fitness function.

3.3.1. Hill climbing and GA hybrid approaches

As suggested by Goldberg (1989), hybrid GAs may
be advantageous when there is the need for optimization
techniques tied to a specific problem structure. The *in-*
large perspective of GAs may be combined with the pre-
cision of local search. GAs are able to explore large
search spaces, but often they reach a solution that is
not accurate, or they very slowly converge to an accu-
rate solution. On the other hand, local optimization
techniques, such as hill climbing, quickly converge to a
local optimum, but they are not very effective for search-
ing large solution spaces because of the possible pres-
ence of local maximum or plateaus.

There are at least two different ways to hybridize a GA
with hill climbing techniques. The first approach attempts
to optimize the best individuals of the last generation,
using hill climbing techniques. The second approach uses
hill climbing to optimize the best individuals of each gen-
eration. Applying hill climbing on each generation could
be expensive. However, this technique “inserts” in each
generation *high quality individuals*, who are determined
by the optimization phase, and therefore reduces the
number of generations requested to achieve convergence.

4. The refactoring framework

As highlighted in the introduction, the proposed
framework consists of several steps:

- First and foremost, software system applications,
libraries, and dependencies among them are
identified;
- Unused functions and objects are identified, removed
or factored out;
- Duplicated or cloned objects are identified and possi-
bly factored out;
- Circular dependencies among libraries, which cause a
library to be linked each time another circularly
linked library is needed, are removed, or, at least,
reduced;
- Large libraries are refactored into smaller ones and, if
possible, transformed into dynamic libraries; and
- Objects which are used by multiple applications, but
which are not yet organized into libraries, are
grouped into new libraries.

The SRF activities and the adopted representations
are detailed in the following subsections.

364 4.1. Software system graph representation

365 A graph representation of dependencies between ob-
 366 ject modules is central to our framework and most of the
 367 SRF computations rely on it. Software systems can be
 368 represented by an instance of the *System Graph* (*SG*),
 369 an example of which is depicted in Fig. 1.

370 *SG* is defined as

$$372 \quad SG \equiv \{O, L, A, D\}, \quad (2)$$

373 where $O \equiv \{o_1, o_2, \dots, o_p\}$ is the set of all object modules;
 374 $L \equiv \{l_1, l_2, \dots, l_n\}$, where $l_i \subseteq O \ i = 1, \dots, n$, is the set of
 375 all software system libraries. Libraries, subsets of ob-
 376 jects, are depicted in Fig. 1 as rounded boxes;
 377 $A \equiv \{a_1, a_2, \dots, a_m\}$, where $A \subseteq O$ and $A \cap \{\cup_i l_i\} = \emptyset$,
 378 is the set of all software system applications. Applica-
 379 tions, i.e. the object modules containing the main sym-
 380 bol, are represented in Fig. 1 as squares source nodes;²
 381 and $D \subseteq O \times O$ is the set of oriented edges $d_{i,j}$ represent-
 382 ing dependencies between objects.

383 We can extract from the *SG* graph two other graphs
 384 useful for our refactoring purposes. The first graph is
 385 called *Use Graph* and it highlights the uses of objects
 386 by applications or by libraries. The *use* relationship is
 387 defined as

$$389 \quad a_x \text{ uses } o_y \iff \exists \text{ path}\{a_x, \dots, o_y\} \in SG. \quad (3)$$

390 In other words the *Use Graph* highlights the reachabil-
 391 ity between applications and library objects in *SGs*.
 392 Such reachability can be obtained computing a k -fold
 393 product on the graph represented by an adjacency
 394 matrix.

395 Similarly, the second graph is called *Dependency*
 396 *Graph* and it is used to represent existing dependencies
 397 between two or more libraries, or between to-be-refac-
 398 tored objects contained in a library. The clustering algo-
 399 rithm should avoid inter-cluster dependencies. The
 400 *dependency* relationship is defined as

$$402 \quad o_x \text{ depends_on } o_y \iff o_x \text{ uses } o_y \wedge o_x \in L \wedge o_y \in L. \quad (4)$$

403 In particular, a dependency (o_x, o_y) is considered an in-
 404 ter-library dependency, i.e., a dependency that increases
 405 the coupling, if $o_x \in l_i$, $o_y \in l_j$, and $i \neq j$.

406 Given the above definition of *SG*, the SRF activities
 407 can be graphically shown in Fig. 2.

408 4.2. Graph construction

409 Prior to recover dependencies among applications
 410 and libraries, and among libraries themselves, execu-
 411 tible applications composing the software system must

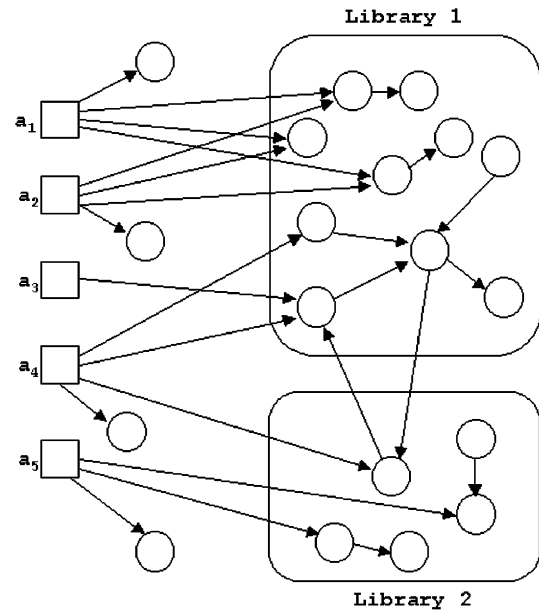


Fig. 1. Example of system graph.

412 be identified. In this paper we rely on an approach sim-
 413 ilar to the one proposed by Antoniol et al. (2001a).
 414 However, Antoniol et al. (2001a) identified applications
 415 by detecting all source files containing the definition of a
 416 main function.

417 Once applications and existing libraries are identified,
 418 the *SG* graph can be built. Given the *use* relationship be-
 419 tween an object module requiring a symbol and a mod-
 420 ule defining it, the corresponding *SG* is built via the
 421 *transitive closure* of the *use* relationship, starting from
 422 the main object of each application and from each lib-
 423 rary. In other words, for each application, undefined
 424 symbols are identified and recursively resolved (possibly
 425 adding new undefined symbols to the stack) first inside
 426 the objects contained in the same path (i.e., other mod-
 427 ules of the application), then inside libraries. A similar
 428 process is performed to detect dependencies among
 429 libraries. Finally, the *use graph* and the *dependency*
 430 *graph*, represented as adjacency matrices MU and MD ,
 431 are extracted from the *SG* graph.

432 4.3. Handling unused objects

433 Symbols defined in libraries which are neither used by
 434 applications nor by other libraries are likely to represent
 435 useless resources. Their presence is often due to utility
 436 functions which are inserted in libraries but which are
 437 not used by the current set of applications, or it is due
 438 to not yet fully implemented features. The objects defin-
 439 ing these unused symbols should be removed from the
 440 libraries, provided that they do not also export used
 441 symbols. In the opposite case such an object should be
 442 left into library and its corresponding source file should
 443 be restructured. One possible refactoring strategy is to

² Applications are not the only source nodes. In fact, as it will be detailed later, also unused objects have no incoming edges, even if they can be distinguished from the applications since the latter also define a main symbol.

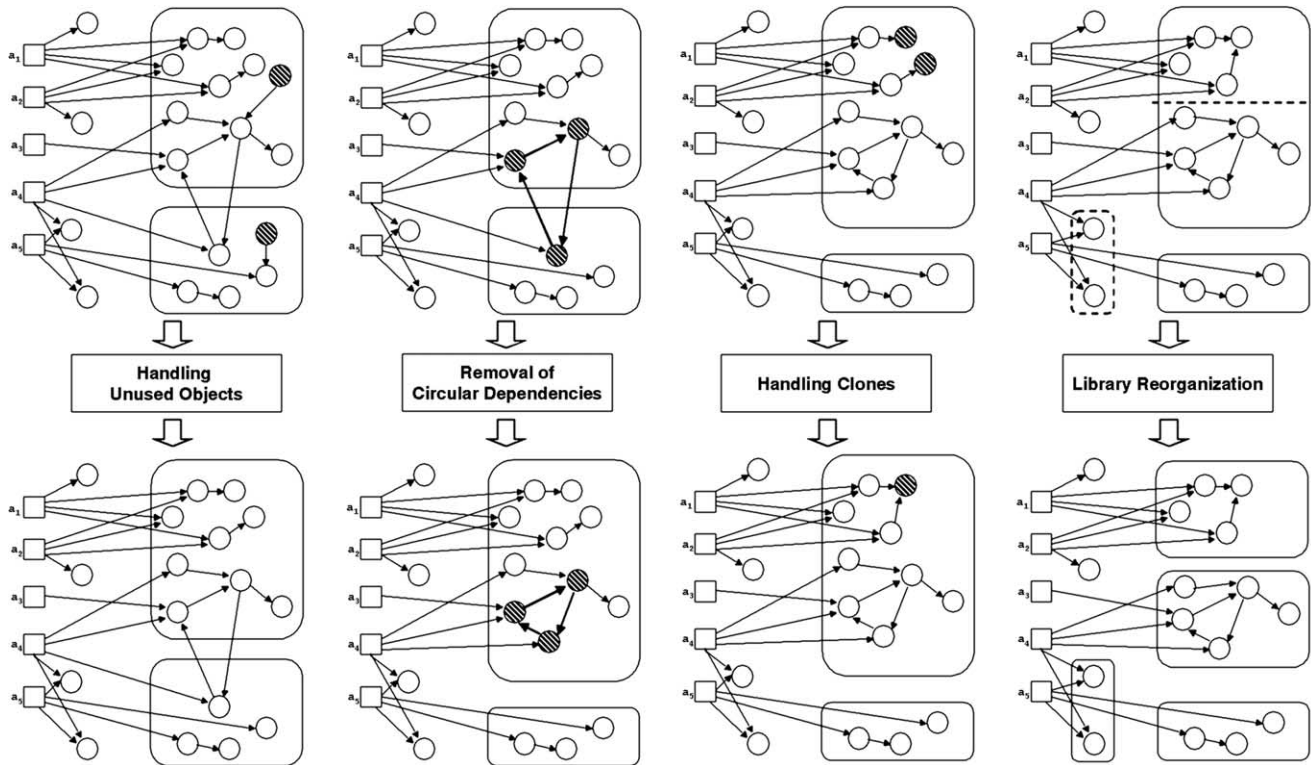


Fig. 2. The framework activities.

444 create two new libraries from each library, one of which
 445 containing all the unused symbols and the other one
 446 containing all the used symbols.

447 4.4. Removal of circular dependencies among libraries

448 The *DG* introduced in Section 4.1 captures dependen-
 449 cies among the different libraries and allows the identifi-
 450 cation of strongly connected components. In particular,
 451 circular dependencies between libraries cause a library
 452 to be linked each time the other one is needed. Once
 453 these dependencies are identified, four strategies could
 454 be used to remove them:

- 455 (1) *Move the object* which causes the circular depend-
 456 ence to another library. This is only feasible if the
 457 object does not need resources located in its original
 458 library and it is not needed by that library;
- 459 (2) *Duplicate the object*: like the previous case, this is
 460 appropriate, if the object does not need resources
 461 located in the original library but, differently from
 462 the previous case, the object is required in that
 463 library. Moving the object the library outside will
 464 make the situation worse;
- 465 (3) *Merge the two libraries*: this strategy should be
 466 avoided whenever possible because it increases
 467 library sizes; however, it could be the only available

solution when the number of objects causing circular
 and, in general, inter-library dependencies is
 very high;

- (4) *Create dynamic libraries*: instead of merging circularly dependent libraries, one may decide to make them dynamic. Circular dependency problem is not solved, but the average amount of resources needed is reduced, as described in Section 4.6.2.

When the *DG* does not allow the removal of circular dependencies and, when, for performance reasons, options three and four cannot be adopted, a deeper analysis should be performed to identify dependencies at the granularity level of functions rather than objects.

Finally, the existence of a complex dependency relationship between two libraries, if confirmed by developer's feedback, indicates the possibility of a library design which has not been done with miniaturization in mind. In this case, library objects should be merged and then refactored again in new clusters, adopting the process detailed in Section 4.6.

4.5. Identification of duplicate symbols and clones

Examining the list of symbols defined in each library allows the comparison of exported symbol names. It is worth noting that homonym symbols in different librar-

468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492

ies may refer to completely different functions, external variable or data structures. On the other hand, two or more symbols may have different names, but they may correspond to duplicated functions. Therefore, clone detection analysis is helpful for library renovation. In this paper a metric-based clone detection process (Antoniol et al., 2001b), aimed at detecting duplicated functions, is adopted. The obtained results suggest different possible actions:

- (1) If a whole, duplicated, object module has been detected inside two or more libraries, then it should be left in only one of these, unless it conflicts with circular dependencies removal (see Section 4.4);
- (2) If duplicated functions are identified inside different objects, refactoring could be performed by moving them outside their respective objects and by applying considerations similar to the previous case; and
- (3) Clone detection may reveal clones outside libraries, since applications may contain duplicated portions of code, in their objects. In some cases, it could be useful to remove such duplicated portions of code and place them into new libraries.

Preliminary to clone refactoring is impact analysis in terms of introduced dependencies, especially circular dependencies, since clone removal may increase dependencies. As explained in Section 4.4 and as it will be shown in Section 4.6, sometimes an object is duplicated to reduce dependencies. In general, it may be preferable to duplicate few objects, rather than introducing a dependence that causes, for a subset of the applications, the linking or the loading of one or more additional libraries. Clearly, if the process duplicates a conspicuous number of objects into two or more libraries, these objects can be refactored, as explained in Section 4.6.2, into a new library on which the old libraries will depend. Overall, clone removal aims to improve the software system maintainability, although attention should be paid to avoid deteriorating software system reliability, and to reflect the developers' objectives (Cordy, 2003). Clone can also contribute decrease the overall software system size; again a tradeoff should be made: sometimes clone refactoring (especially for very small clones) produces a system bigger than the original one.

4.6. Library refactoring

The last phase of the SRF is devoted to splitting existing, large libraries into smaller clusters of objects. Basically, the idea is similar to that proposed by Antoniol et al. (2001a) to identify libraries. To minimize the average number of libraries required by each program, objects used by a common set of programs should be grouped together. Antoniol et al. (2001a) used a *concept lattice* to group objects into libraries. Although the

lattice gives useful information, it becomes unmanageable when a large number of applications and libraries must be handled (Anquetil, 2000), as in our case study. Instead of pruning information on a *concept lattice* like Siff and Reps (1999) and Tonella (2001), clustering analysis was performed, similar to Anquetil and Lethbridge (1998), Mancoridis et al. (1998), and Merlo et al. (1993).

The library refactoring process, as shown in Fig. 3, consists of the following steps:

- (1) Determine the optimal number of clusters and an initial solution;
- (2) Determine the new candidate libraries using a GA; and
- (3) Ask developers for feedback and, possibly, iterate through step 2.

4.6.1. Determining the optimal number of clusters and a suboptimal solution

As explained in Section 3.2, the optimal number of clusters is determined by inspecting the *Silhouette* statistics computed on the suboptimal clusters which are determined using agglomerative-nesting clustering. Given the curve of the average *Silhouette* values obtained from Eq. (1) for different numbers k of clusters, we choose for some libraries the *knee* of that curve (Kaufman and Rousseeuw, 1990) as the optimal number of clusters, instead of considering the maximum of the curve because that is often too high for our refactoring purpose.

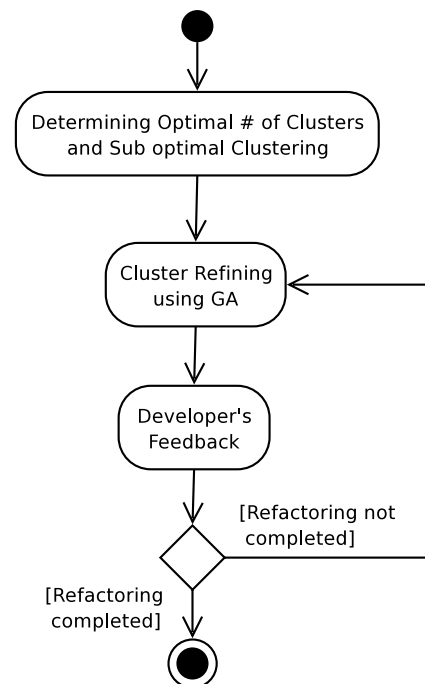


Fig. 3. Activity diagram of the library refactoring process.

575 We have also incorporated experts' knowledge in the
576 choice of the optimal number of clusters and we have
577 considered a tradeoff between excessive fragmentation
578 produced by too many clusters and excessive library size
579 produced by fewer clusters. The suboptimal solution for
580 the chosen value of k is then used as the starting point of
581 the application of a GA, which is the subsequent frame-
582 work step.

583 The effectiveness of the refactoring process is evalu-
584 ated by a quality measure of the new library organiza-
585 tion. Let k be the number of clusters l_{x_1}, \dots, l_{x_k}
586 obtained from a library l_x . The *Partitioning Ratio (PR)*
587 is defined as
588

$$590 \quad PR(x) = 100 \times \sum_{i=1}^m \frac{\sum_{j=1}^k |l_{x_j}| \times mu_{i,x_j}}{|l_x| \times mu_{i,x}}, \quad (5)$$

591 where $|l_x|$ is the number of objects archived into library
592 l_x . The smaller is the PR , the more effective is the parti-
593 tioning since the average number of objects linked or
594 loaded by each application is smaller than using the
595 whole old library.

596 4.6.2. Refining the solution using genetic algorithms

597 The solution determined by the previous step presents
598 two main drawbacks:

- 599 (1) The number of dependencies between the new
600 libraries may be high. Each time a symbol from a
601 library is needed, another library may also need
602 to be loaded, therefore reducing the advantage of
603 having new smaller libraries; and
- 604 (2) New libraries may not be meaningful with respect
605 to developers' intentions whose feedback has to be
606 incorporated in the refactoring process.

607 Of course, as shown by Di Penta et al., 2002, an
608 important step to perform is the conversion of static
609 libraries into dynamically-loadable libraries (DLL), so
610 that each and possibly small library is loaded at run-
611 time only when needed, and it is unloaded when it is
612 no longer useful. However, the DLL approach presents
613 a main drawback: loading and unloading libraries
614 may be cause of a significant decrease in performance
615 and its use should be limited, when performance
616 constitutes an essential requirement, and, whenever
617 possible, it should be accompanied by dependency
618 minimization.

620 The genome has been encoded using a bit-matrix
621 encoding. The genome matrix GM for each library to
622 refactor corresponds to a matrix of k rows and $|l_x|$ col-
623 umns, where $gm_{i,j} = 1$ if the object j is contained into
624 cluster i , 0 otherwise. Clearly, the presence of the same
625 object in more libraries is indicated by more "1" in the
626 same column (this is not possible using the array gen-
627 ome, widely used for graph partitioning problems). As

628 already stated, instead of randomly generating the initial
629 population (i.e., the initial libraries), the GA is initial-
630 ized with the encoding of the set of libraries obtained
631 in the previous step.

632 The fitness function has been conceived to balance
633 four factors:

- 634 (1) The number of inter-library dependencies at a given
635 generation;
- 636 (2) The total number of objects linked to each applica-
637 tion which should be as small as possible;
- 638 (3) The size of the new libraries; and
- 639 (4) The feedback given by the developers.

640 Overall, the fitness function F is defined in terms of
641 four factors which are the *Dependency Factor (DF)*,
642 the *Partitioning Ratio (PR)* defined by Eq. (5), the
643 *Standard Deviation Factor (SDF)*, and the *Feedback Fac-
644 tor (FF)*.
645

646 DF is defined as:

$$647 \quad DF(g) = \sum_{i=0}^{|l_x|-1} \sum_{j=0}^{m-1} gm_{i,j} \sum_{k=0}^{k=m-1} md_{j,k} (1 - gm_{i,k}) \times [1 - \delta(k, j)], \quad (6)$$

650 where $\delta(x, y)$ is the well-known Kronecker delta
651 function:

$$652 \quad \delta(x, y) = \begin{cases} 1 & x = y, \\ 0 & x \neq y, \end{cases} \quad (653)$$

654 $gm_{i,j}$ is the genome encoding i.e., the $GM[i, j]$ bit matrix
655 entry. As shown in Eq. (6), the $DF(g)$ is incremented
656 each time an object (i.e., a high bit in the genome) de-
657 pends from another object not contained in the same
658 cluster. SDF can be thought of as the difference between
659 the initial library sizes standard deviation and the one at
660 the current generation. Without taking SDF into ac-
661 count, the SRGA may attempt to reduce dependencies
662 by grouping a large fraction of the objects in the same
663 library and it may negatively affect the PR . A similar
664 factor was also applied by Talbi and Bessière (1991).
665 Given the arrays of library sizes S_0 and S_g , respectively
666 for the initial population and for the g th generation,
667 SDF is

$$668 \quad SDF(g) = |\sigma_{S_0} - \sigma_{S_g}|. \quad (7)$$

670 The fourth factor takes into account the developers'
671 feedback. After a first execution of the SRGA without
672 considering FF , developers are asked to provide a feed-
673 back on the proposed new libraries. Developers' feed-
674 back is stored in a bit-matrix FM , which has the same
675 structure of the genome matrix and which incorporates
676 those changes to the libraries that developers suggested.
677 After this feedback, the SRGA is run again taking into
678 account, this time, the *feedback factor FF*, based on the
679 difference between the genome and the FM matrix:

$$FF = \sum_{i=1}^k \sum_{j=1}^{|I_x|} |gm_{i,j} - fm_{i,j}|. \quad (8)$$

681

682 In other words, the FF counts the number of differences
683 between the genome and the refactoring proposed by
684 developers.

685 The fitness function F is formally defined as

$$686 \quad 688 \quad F(g) = DF(g) + w_1 PR(g) + w_2 SDF(g) + w_3 FF(g), \quad (9)$$

689 where w_1 , w_2 and w_3 are real, positive weighting factors
690 for the PR , SDF , and FF contribution to the overall fit-
691 ness function. The higher is w_1 , the smaller will be the
692 overall number of objects linked by applications at the
693 expense of dependency reduction. Similarly, the higher
694 is w_2 , the more similar will be the result to the starting
695 set of library, again, at the expense of a satisfactory
696 dependency reduction. After the first preliminary run
697 of the SRGA which must be performed with $w_3 = 0$,
698 w_3 should be properly sized to weight the influence of
699 developers' feedback. As stated in (9), our fitness func-
700 tion is multi-objective (Deb, 1999). Notice that, since
701 we aim to give maximum priority to dependency reduc-
702 tion, the DF weight is set to 1. Successively, w_1 , w_2 and
703 w_3 are selected using a trial-and-error, iterative proce-
704 dure, adjusting them each time until the DF , PR , SDF ,
705 and FF obtained at the final step were satisfactory.
706 The process is guided by computing each time the aver-
707 age values for DF , PR , SDF , and FF , and by plotting
708 their evolution, to determine the 3D space region in
709 which the population should evolve.

710 The crossover operator used in this paper is the *one*
711 *point crossover* which exchanges the content of two gen-
712 ome matrices around the same random column (see Fig.
713 4a). The mutation operator works in two modes:

- 714 (1) with probability $pmut$, it takes a random column
715 and randomly swaps two bits: this means that, if
716 the two swapped bits are different then an object
717 is moved from a library to another (see Fig. 4b); or
718 (2) with probability $pclone < pmut$, it takes a random
719 position in the matrix: if it is zero and the library
720 is dependent on it, then the mutation operator
721 clones the object into the current library (Fig. 4c).

722 Noticeably, cloning an object increases both PR and
723 SDF , and therefore it must be minimized. The SRGA
724 heuristically activates the cloning only for the final part
725 of the evolution (after 66% of generations in our case
726 study). Our strategy favors dependency minimization
727 by moving objects between libraries. At the end, we at-
728 tempt to remove remaining dependencies by cloning ob-
729 jects. Obviously, at the end of the refactoring process
730 cloned objects should be factored out again. For exam-
731 ple, if objects o_a and o_b are contained in both l_i and l_j ,
732 then o_a and o_b should be moved into a third library on
733 which l_i and l_j depend.
734

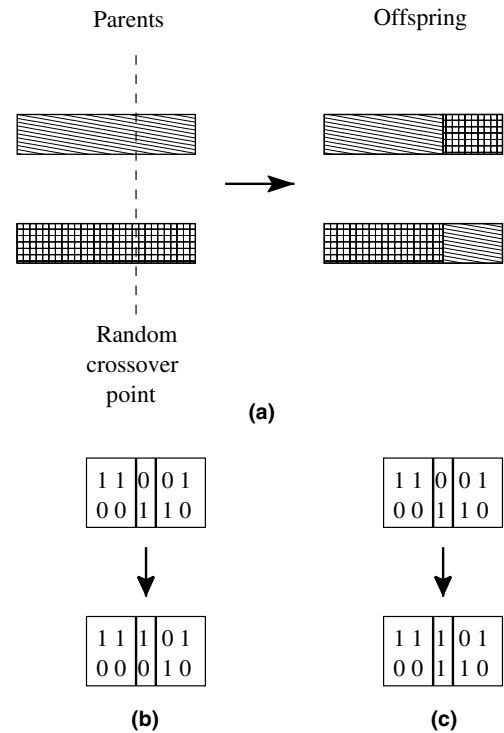


Fig. 4. Genetic operators: (a) crossover, (b) mutation (move an object), and (c) mutation (clone an object).

735 Finally, we have introduced the *Lock Matrix (LM)* as
736 a further, stronger level of developers' feedback. When
737 developers strongly believe that an object should belong
738 to a cluster, LM matrix gives them the possibility to en-
739 force such a constraint. The mutation operator does not
740 perform any action that would bring a genome in a
741 inconsistent state with respect to the *Lock Matrix*.

742 The population size and the number of generations
743 are determined by using an iterative procedure, which
744 doubles both of them each time until the obtained DF ,
745 PR and FF are equal to those obtained at the previous
746 iterative step.

747 The SRGA suffers from slow convergence. To im-
748 prove its performance, it has been hybridized with hill
749 climbing techniques. In our experience, applying hill
750 climbing only to the last generation significantly im-
751 proves neither the performance nor the results. On the
752 opposite, applying hill climbing to the best individuals
753 of each generation makes the SRGA converge signifi-
754 cantly faster.

4.7. Identification of new libraries

755 Due to its evolution, a software system tends to con-
756 tain objects that, even if used by a common set of appli-
757 cations, are not contained in any library. Their
758 identification and organization into libraries should
759 therefore be desirable. The factoring process is quite
760 similar to that described in the previous section. In par-
761

762 ticular, a *MU* matrix is built on a subgraph of the *use*
 763 *graph* obtained by removing all the already existing
 764 libraries. Then, a first set of new *candidate libraries* is
 765 built by analyzing the dendrogram and the *Silhouette*
 766 statistics. These libraries are then refined with the aid
 767 of the SRGA and of developers' feedback.

768 4.8. Tool support

769 To support the refactoring process, different tools
 770 have been conceived:

- 771 (1) *The application identifier* identifies the list of object
 772 modules containing the `main` symbol by using the
 773 `nm` Unix tool;
- 774 (2) *The graph extractor*, which is also based on the `nm`
 775 tool, which produces the *System Graph*, the *Use*
 776 *Graph*, and the *Dependency Graph*. The *graph*
 777 *extractor* also exports data in `.DOT` format, to
 778 allow visualization and analysis using the `Dotty`
 779 *graph visualization tool*; ³
- 780 (3) *The unused symbol identifier* produces, for each
 781 library, the list of the symbols which are not used
 782 by any application or library together with the
 783 object names in which those symbols are contained;
- 784 (4) *The circular dependency identifier* produces the list
 785 of all circular paths among libraries;
- 786 (5) *The duplicated symbol identifier* identifies the list of
 787 duplicated and defined external symbols. It is used
 788 in conjunction with the metric-based *clone detector*
 789 (see Antoniol et al., 2001b, for details) and with the
 790 *dependency graph extractor* to minimize the pres-
 791 ence of clones inside libraries;
- 792 (6) *The number of clusters identifier* implements the *Sil-*
 793 *houette* statistics. In particular, implementations
 794 available in the *cluster* package of the *R Statistical*
 795 *Environment* ⁴ have been used;
- 796 (7) *The library refactoring tool* supports the process of
 797 splitting libraries in smaller clusters. Cluster analy-
 798 sis is performed by the *Agnes* function available in
 799 the *cluster* package of *R Statistical Environment*;
- 800 (8) *The GA library refiner* is implemented in C++ using
 801 the *GAlib*; ⁵ and
- 802 (9) *The developers' feedback collector* is a web applica-
 803 tion that allows developers to post their feedback
 804 about the produced libraries on an appropriate
 805 web site.

806 The SRF works under any standard Unix operating
 807 system, or under any operating system which supports
 808 the GNU tool set. In particular, the SRF uses the stand-
 809 ard Bourne shell (or the new Bash), the Perl interpreter,

the *R statistical environment* and a C++ compiler for the
 GA library refiner. To collect the programmers' feed-
 back, the SRF relies on a PHP web application (the
developers' feedback collector). Since the required infra-
 structure is available under several operating systems
 (both Unixes and Windows) the SRF is widely portable.

5. Case study

As mentioned in the introduction, the SRF has been
 applied to GRASS, which is a large open source GIS. In
 particular, the GRASS CVS development snapshot of
 April 5, 2002 ⁶ was used as a case study. Its characteris-
 tics are summarized in Table 1.

GRASS modules, which correspond to applications
 and which represent commands, are organized by name,
 based on their function class such as display, general,
 imagery, raster, vector or site, etc. The first letter of a
 module name refers to a function class and is followed
 by one dot and one or two other dot-separated words,
 which describe specific tasks. All GRASS modules are
 linked with an internal "front.end". If there are no com-
 mand-line arguments entered by a user, the "front.end"
 module calls the interactive version of a command. Oth-
 erwise, it will start the command-line version. If only
 one version of the specific command exists, i.e., if there
 is only one command-line version available, the com-
 mand is executed. Code parameters and flags are defined
 within each module. They are used to ask user to define
 map names and other options.

GRASS provides an ANSI C language API with sev-
 eral hundreds of GIS functions which are used by
 GRASS modules, to read and write maps, to compute
 areas and distances for georeferenced data, and to visu-
 alize attributes and maps. Details of GRASS program-
 ming are covered in the "GRASS 5.0 Programmer's
 Manual" (Neteler, 2001).

6. Case study results

This section presents the results obtained by applying
 the SRF, which has been described in Section 4, to
 GRASS.

6.1. Handling unused objects

Out of 921 objects composing GRASS libraries, 89
 were not used by any application, nor by other libraries.
 When refactoring libraries with the SRF, those objects
 will be moved and organized into a separate cluster,
 thought of as a sort of repository to be "frozen" for fu-

³ <http://www.research.att.com/sw/tools/graphviz/>

⁴ <http://www.r-project.org>

⁵ <http://lancet.mit.edu/ga/>

⁶ Downloadable from <http://grass.itc.it>

Table 1
GRASS key characteristics

Pre-existing libraries	43
Library objects	921
Applications	517
C source files	7107
C KLOC	1014

ture uses. A deeper analysis revealed that some functions, which are contained in unused objects, wrap lower level GRASS functions such as `db_create_index`, wrap standard library and system call functions such as `scan_dbl`, `scan_int`, `whoami`, and, in general, provide some simple functionalities using lower level functions such as `datetime_is_same`, that compares two `DateTime` structures. An interesting example is the library `libdbmi` (see also Section 6.4): out of 97 objects, 19 were not used at all. In all cases, the unused functions corresponded to one or more wrapped, lower level functions, that have been directly used by applications.

6.2. Removal of circular dependencies among libraries

Three cases of circular dependencies among libraries were found. The first dependency was between `libstubs.a` and `libdbmi.a`. In particular, we discovered that `libstubs.a` required one symbol, located inside the `error.o` module which belonged to `libdbmi.a`. On the other hand, `libdbmi.a` required 27 symbols from `libstubs.a`. The obvious solution was to move `error.o` into `libstubs.a`: this required moving in that library also the module `alloc.o`, since it depends from `error.o`.

The second circular dependency was found between `libgis.a` and `libcoorenv.a`. In particular, `libgis.a` required three symbols. Such symbols were located in the module `datum.o` from `libcoorenv.a`. In the other direction, `libcoorenv.a` dependencies involved 13 symbols from `libgis.a`. Moving `datum.o` into `libgis.a` resolved the problem.

Finally, circular dependencies were found between `libvect.a` and `libdig2.a`. They involved 13 symbols in one direction and 31 symbols in the other direction. Symbols involved in the dependencies were located

in several different objects. The links present in the dependency graph excluded the possibility of resolving circular dependencies between `libvect.a` and `libdig2.a` by simply moving or duplicating objects. The decision taken together with GRASS developers was initially to merge the two libraries which, in effect, have been designed to work together, and then try to refactor the new library (see Section 6.4).

6.3. Identification of clones

Clone detection was performed at two different levels of the software system architecture: within libraries and on the whole system. In the first case, clone detection aimed at library renovation; in the second case, the objective was to identify portions of duplicated code that could be potentially re-organized into new libraries.

Table 2 reports results obtained from clone analysis in terms of the total number of analyzed functions, the number of clone clusters (Antoniol et al., 2002) detected, and the number and the percentage of cloned functions. Finally, clones were computed while filtering out the shortest functions; for example, two functions that simply return a value are clones by definition, but they are not significant and should not be taken into consideration. Results are reported considering two thresholds of function size: functions longer than five and than 10 LOCs.

As shown in Table 2, the overall percentage of clones is not negligible (26.04%), even considering only functions longer than five LOCs (16.38%) and it suggests a potential for reduction in the number of the cloned functions. Clearly, the actual reduction rate depends on the number of *false positives* which typically include functions that simply contain a list of calls to other functions (where the number of calls and of parameters match), functions that print different error messages and, in general, any other function that shares the same metrics while being different.

The number of clones contained inside libraries is low, indicating that the developers accurately factored functions and objects to avoid duplicates. Finally, we investigated the set of clones between libraries and objects outside libraries in the perspective of possible refactoring. The analysis of clones inside libraries revealed an

Table 2
Results of clone detection

	Total number of functions	Number of clone clusters	Number of cloned functions	Percent of cloned functions (%)	Threshold (LOCs)
Overall	22,229	2019	5789	26.04	5
		1404	3641	16.38	10
Within libraries	5271	72	180	3.41	5
		41	101	1.92	10
Outside libraries	16,958	1817	4974	29.33	5
		1290	3268	19.27	10
Libraries vs. outside	22,229	130	635	2.86	5
		73	272	1.22	10

932 interesting situation: 16 functions from library libor-
 933 tho, were cloned across libimage_sup, libgmth
 934 and libtrans. Nine of the cloned functions were de-
 935 voted to performing matrix algebra. By analyzing the
 936 *Dependency Graph* of libortho (see Fig. 5), a sub-
 937 graph composed of such functions was identified and de-
 938 picted in the box on the right. On the other hand, seven
 939 of the functions in the box on the left were cloned in
 940 libimage_sup. In particular, the entire structure en-
 941 closed in the rounded-dashed-box was replicated in that
 942 library. libortho was split libortho in two librar-
 943 ies, shown in the two boxes in Fig. 5:

- 944 (1) A library (libmatrix) to handle matrices; and
- 945 (2) A library (libcamera) to handle photogrammet-
 946 ric computations for aerial cameras.

947 Cloned functions contained in these two libraries
 948 were removed from libimage_sup, libgmth and
 949 libtrans.

951 Several “interesting” clones were also found outside
 952 libraries. In particular, the r.mapcalc3 application
 953 contains four clusters of cloned functions, spanning
 954 from 27 to 59 LOCs in size. These cluster contain math-
 955 ematical functions, cloned to handle different data types.
 956 In this case, refactoring is clearly possible by generaliz-
 957 ing the operations and by abstracting types.

958 Finally, we analyzed clones between applications and
 959 libraries. In most cases clones were revealed to be part of
 960 *legacy* applications developed before the corresponding
 961 functions were added into a library. Unfortunately, the
 962 application was never changed afterwards. A relevant
 963 fraction of about 20% of these clones was discovered in
 964 the *contrib* subsystems, which had often been developed
 965 by third parties and therefore which were not always
 966 properly aligned with respect to the rest of the system.

967 6.4. Library refactoring

968 Refactoring was performed on libraries which were
 969 composed of a large number of objects (see Table 3),

Table 3
GRASS largest libraries

Library	Objects
libgis	184
libdbmi	97
libproj	119
libvect-new	54

970 by following the process described in Section 4.6 and de-
 971 picted in Fig. 3. As suggested by developers, libproj
 972 was not refactored, because it was under development
 973 by a different team. As explained in Section 6.2,
 974 libvect-new library was obtained by merging lib-
 975 vect.a and libdig2.a.

976 *Silhouette* statistics was used to determine the optimal
 977 number of clusters for each library. Values of such sta-
 978 tistics, are plotted in Fig. 6, for different number of clus-
 979 ters. We decided to split libgis into four clusters
 980 (instead of the six proposed in Di Penta et al., 2002),
 981 and to divide libvect-new and libdbmi into three
 982 clusters. It is worth noting that, for libgis, the num-
 983 ber of clusters was chosen in correspondence of the *Sil-*
 984 *houette* maximum; for the other two libraries, a

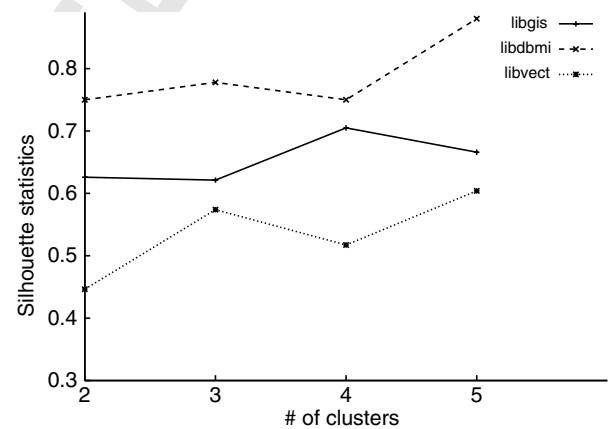


Fig. 6. Silhouette statistics for different number of clusters.

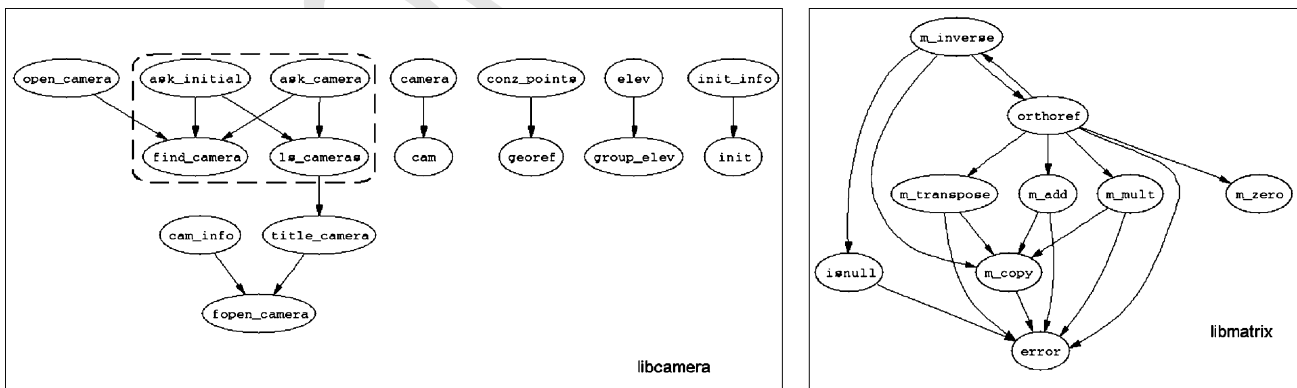


Fig. 5. Splitting library libortho.

985 compromise was accepted between maximizing the *Sil-*
986 *houette* and avoiding excessive fragmentation.

987 Subsequently, a preliminary clustering was performed
988 and it was refined by an initial execution of the SRGA,
989 which had been performed without considering any
990 developers' feedback and by setting $w_3 = 0$. Table 4 re-
991 ports for each library:

- 992 • The number of objects composing the library;
- 993 • The number of *candidate libraries* the original library
994 is refactored into and the corresponding *Silhouette*
995 statistics value;
- 996 • The number of inter-library dependencies and *PR*
997 before applying the SRGA; and
- 998 • The number of inter-library dependencies and *PR*
999 after applying the SRGA.

1000 As shown, the SRGA reduced *libgis* dependencies
1001 from 579 to 26, while keeping *PR* almost constant (from
1002 51% to 48%). A significant reduction of inter-library
1003 dependency was obtained (from 237 to 4 for *libdbmi*
1004 and from 66 to 3 for *libvect*), while slightly reducing
1005 *PR*, except for *libdbmi*, where it increased to 46% and
1006 it was worse than the preliminary solution.

1007 The first refactored architecture of the *candidate*
1008 *libraries* was submitted to *GRASS* developers to seek
1009 their feedback. For *libgis*, manual analysis indicated
1010 that the first cluster should contain "utility" and "allo-
1011 cation" functions, the second "area" and "geodesic"
1012 functions, the third "color-related" functions, and the
1013 fourth "raster" functions. For *libvect-new*, develop-
1014 ers indicated that the first cluster should contain basic
1015 file-system operations and the other two clusters should
1016 include all other functions without any further distinc-
1017 tion. The feedback for *libdbmi* was quite different with
1018 respect to the other two libraries. In this case, developers
1019 confirmed that the solution suggested by the hierarchical
1020 clustering performed before applying the SRGA re-
1021 flected their own conception of libraries. A manual
1022 graph analysis via *Dotty* graph visualization agreed,
1023 too. In fact, as also reported by Di Penta et al. (2002),
1024 the library was split into the three following clusters:
1025

- 1026 • *libdbmi-1* contains (19) unused objects;
- 1027 • *libdbmi-2* contains (30) objects which are directly
1028 used by applications; and

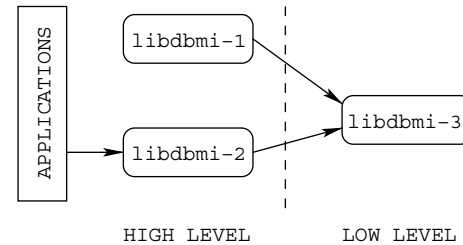


Fig. 7. New libdbmi layering structure.

- *libdbmi-3* contains 48 objects, which are only
internally used by *libdbmi*, and represents some
sort of "low-level" library.

Fig. 7 reports the layering structure of the clusters ex-
tracted from *libdbmi*. To avoid circular dependencies,
one object was moved from *libdbmi-3* to *libdbmi-1*.
Clearly, when refactoring a large software system such
as *GRASS*, a compromise should be accepted between
having small and decoupled clusters like those gener-
ated by applying the SRGA and having clusters that are
not totally decoupled, but are conceptually cohesive,
since they contain functions which implement closely-
related tasks. In the latter case, memory optimization
is possible adopting, as noted, dynamically loadable
libraries (in spite, however, of performances, as ex-
plained in Section 4.6.2). We decided to leave
libdbmi clusters as they were after hierarchical clus-
tering and to perform a "second iteration" of the SRGA
refactoring on *libgis* and *libvect-new*, while taking
into consideration also the *Feedback Factor FF*, this
time. For sake of completeness, we also reported re-
sults for *libdbmi*. By varying the w_1 , w_2 and w_3
thresholds, we obtained different results. As shown in
Table 5, it was never possible to achieve a complete
cluster decoupling and to obtain, at the same time,
libraries which were very close to the structure pro-
posed by developers.

In Table 5, the comparison of the first three columns
with the last three highlights that, after the first SRGA
iteration, the coupling between clusters remained low.
On the other hand, as highlighted by a high *FF* value
before the second iteration, identified libraries tend to
have a structure which is somehow different with re-
spect to developers' intention. The second iteration of
the SRGA tried to decrease *FF*, while, unfortunately,
coupling increased. At such a stage, in the authors' op-
inion, developers may decide either to produce *mean-*
ingful libraries

Table 4
Results of the library refactoring process before considering feedback ($w_3 = 0$)

Library	Number of objects	Candidate libraries (<i>k</i>)	Silhouette statistics	Before GA		After GA	
				DF	PR (%)	DF	PR (%)
<i>libgis</i>	184	4	0.70	579	51	26	48
<i>libdbmi</i>	97	3	0.78	237	35	4	46
<i>libvect</i>	54	3	0.57	66	46	3	40

Table 5

Results of the second round of the library refactoring process ($w_3 \neq 0$)

Library	Number of objects	Candidate libraries (k)	Before second round			After second round		
			FF	DF	PR (%)	FF	DF	PR (%)
libgis	184	4	203	26	48	128	60	52
libdbmi	97	3	97	4	46	23	43	39
libvect	54	3	72	3	40	30	6	52

Table 6

Performance comparison between pure GA and hybrid GA with hill climbing

Library	Pure GA		Hybrid GA		Fitness difference (%)	Time difference (%)
	Fitness function	Time (s)	Fitness function	Time (s)		
libgis	3119	9113	3239	4524	1	49
libdbmi	77	509	83	190	7	37
libvect	195	96	198	41	3	43

1066 and to reduce the memory requirements using dynam- 1102
 1067 ically-loadable libraries, or to obtain independent clus- 1103
 1068 ters, which may not always conceptually group objects 1104
 1069 as related as expected. Although it is counterintuitive, 1105
 1070 the latter result is not surprising, since experts classified 1106
 1071 functions according to the intended purpose or seman- 1107
 1072 tic. This seldom ensure high cohesion and low coupling, 1108
 1073 because the improvement of the latter attributes pro- 1109
 1074 duces a final partitioning which somehow differs from 1110
 1075 what it was expected. 1111

1076 The addition of hill climbing into the SRGA did not 1112
 1077 improve the fitness function, since the SRGA also 1113
 1078 converged to similar results, when it was executed on an in- 1114
 1079 creased number of generations and increased population 1115
 1080 size. Noticeably, performing hill climbing on the best 1116
 1081 individuals of each generation produced a drastic reduc- 1117
 1082 tion of convergence times. Comparing both strategies 1118
 1083 when the difference between values of the fitness func- 1119
 1084 tion was below 10% highlighted that a hybrid strategy 1120
 1085 allowed on average to reduce the execution time of 1121
 1086 43%. Convergence times for a *Compaq Proliant*TM with 1122
 1087 Dual XeonTM 900 MHz processor, 2 MB Cache and 1123
 1088 4 GB of RAM are reported in Table 6. 1124

1089 6.5. Extraction of new libraries

1090 To identify new candidate libraries, the final step of 1125
 1091 the SRF is devoted to the analysis of the *Use Graph* ob- 1126
 1092 tained by subtracting the already existing libraries. 1127
 1093 Sometimes there are groups of objects used by a com- 1128
 1094 mon set of applications, but they have not yet been 1129
 1095 organized into libraries. Clustering was performed on 1130
 1096 objects used by, at least, two applications. 1131

1097 Results revealed the presence of four clusters which 1132
 1098 were all located in the *orthophoto* subsystem. The number 1133
 1099 of dependencies between clusters was low and it was pos- 1134
 1100 sible to solve them simply moving a couple of objects be- 1135
 1101 tween clusters. Besides, all clusters had a considerable 1136

number of dependencies to external objects which be-
 longed to the same set of applications. To eliminate these
 dependencies, it would have been necessary to increase the
 size of each cluster by 100%, clearly in contradiction with
 respect to the intended objective of reducing applications'
 memory requirements and size. Consequently, it was
 decided not to cluster these objects into libraries. In the
 authors' opinion, this is not a negative result, but it consti-
 tutes a quality indicator of the system showing that devel-
 opers had carefully created and maintained libraries.

7. Conclusions

This paper has presented a framework for software
 system renovation (SRF) and the results of its applica-
 tion to *GRASS* Geographical Information System,
 which is over one million LOCs in size.

The SRF has allowed us to remove several structural
 problems from *GRASS*. In particular, unused objects
 were identified and factored out; clones were identified
 and, especially for those inside libraries, refactoring
 was performed. The SRF incorporates a novel library
 refactoring process, in which a suboptimal solution is
 first identified by hierarchical clustering and then refined
 by the SRGA. The proposed SRGA fitness function
 takes into account different factors: minimizing the
 number of dependencies, the average number of objects
 linked by each application, and the feedback of devel-
 opers. Although the approach has been applied on C and
 C++ systems (*GRASS* and others reported by Antoniol
 et al., 2003), it is not tied to any specific programming
 language, provided that object modules, which contain
 a list of defined and required symbols, be available.
 However, for applications to be executed on a virtual
 machine, such as Java, Smalltalk programs, other ap-
 proaches such as those of Tip et al. (1999) and Rayside
 and Kontogiannis (2002) may be preferable.

Overall, the SRF helps to monitor and improve the quality of a software system, which tends inevitably to deteriorate during the evolution. Unused objects, clones, library coupling, library sizes, and poor object organization are in fact significant quality indicators. For instance, the absence of new libraries identified by the SRF in *GRASS* indicates a careful design and a controlled evolution. Moreover, the SRF also addresses the miniaturization problem, which is relevant to port applications on limited-resource devices. The SRF has allowed us to reduce *GRASS* memory requirements and to improve its performance. The average number of library objects linked by each application was indeed reduced of about 50%. At the time of writing, *GRASS* has successfully been ported on a PDA (i.e., a CompaQ iPAQ). Given the size of the application and the available resources, a brute force automatic approach would not be feasible, since developers' suggestions were an essential component for the miniaturization process.

Clone detection performed on *GRASS* revealed that the cloning level outside libraries was not negligible and suggested further clone refactoring. Besides, the cloning level inside libraries was in general low, except for the mentioned cases. The cloning between libraries and the rest of the system was in most cases due to third party applications. Most of the system reorganization work described in the paper was incorporated in the subsequent releases of *GRASS* by removing unused objects and some clones, and by reorganizing some libraries. The latter reorganization, as pointed out in the paper, was carried out with minor modifications with respect to the result of the SRF.

Our in-progress work is devoted to investigate the feasibility of integrating other sources of knowledge into the SRF with special regards to dynamic information and in-field user profiles (Antoniol and Di Penta, 2003), obtained by instrumenting the source code.

1174 Acknowledgments

We are grateful to the *GRASS* development team for the support, the information provided, and the feedback on the refactored artifacts. Giuliano Antoniol and Massimiliano Di Penta were partially supported by the ASI grant I/R/091/00. Markus Neteler was partially supported by the FUR-PAT Project WEBFAQ. Ettore Merlo was partially supported by the National Sciences and Engineering Research Council of Canada (NSERC).

1184 References

1185 Anderberg, M.R., 1973. Cluster Analysis for Applications. Academic
1186 Press Inc.

- Anquetil, N., 2000. A comparison of graphs of concept for reverse engineering. In: Proceedings of the IEEE International Workshop on Program Comprehension. IEEE Computer Society Press, Los Alamitos, CA, USA, pp. 231–240. 1187–1189
- Anquetil, N., Lethbridge, T., 1998. Extracting concepts from file names; a new file clustering criterion. In: Proceedings of the International Conference on Software Engineering. IEEE Computer Society Press, Los Alamitos, CA, USA, pp. 84–93. 1190–1192
- Antoniol, G., Di Penta, M., 2003. Library miniaturization using static and dynamic information. In: Proceedings of IEEE International Conference on Software Maintenance, Amsterdam, The Netherlands. pp. 235–244. 1193–1196
- Antoniol, G., Casazza, G., Di Penta, M., Merlo, E., 2001a. A method to re-organize legacy systems via concept analysis. In: Proceedings of the IEEE International Workshop on Program Comprehension. IEEE Computer Society Press, Los Alamitos, CA, USA, Toronto, ON, Canada, pp. 281–290. 1197–1199
- Antoniol, G., Casazza, G., Di Penta, M., Merlo, E., 2001b. Modeling clones evolution through time series. In: Proceedings of IEEE International Conference on Software Maintenance. pp. 273–280. 1200–1202
- Antoniol, G., Villano, U., Merlo, E., Di Penta, M., 2002. Analyzing cloning evolution in the Linux Kernel. In: SCAM 2002 Special Issue. Information and Software Technology 44, 755–765. 1203–1206
- Antoniol, G., Di Penta, M., Neteler, M., 2003. Moving to smaller libraries via clustering and genetic algorithms. In: European Conference on Software Maintenance and Reengineering. IEEE Computer Society Press, Los Alamitos, CA, USA, Benevento, Italy, pp. 307–316. 1207–1211
- Bui, T.N., Moon, B.R., 1996. Genetic algorithm and graph partitioning. IEEE Transactions on Computers 45 (7), 841–855. 1212–1214
- Cordy, J., 2003. Comprehending reality—practical barriers to industrial adoption of software maintenance automation. In: Proceedings of the IEEE International Workshop on Program Comprehension, Portland, OR, USA. pp. 196–205. 1215–1219
- Deb, K., 1999. Multi-objective genetic algorithms: problem difficulties and construction of test problems. Evolutionary Computation 7 (3), 205–230. 1220–1223
- Di Penta, M., Neteler, M., Antoniol, G., Merlo, E., 2002. Knowledge-based library re-factoring for an open source project. In: Proceedings of IEEE Working Conference on Reverse Engineering. IEEE Computer Society Press, Los Alamitos, CA, USA, Richmond, VA, pp. 128–137. 1224–1226
- Doval, D., Mancoridis, S., Mitchell, B., 1999. Automatic clustering of software systems using a genetic algorithm. In: Software Technology and Engineering Practice (STEP), Pittsburgh, PA. pp. 73–91. 1227–1231
- Garey, M., Johnson, D., 1979. Computers and Intractability: a Guide to the Theory of NP-Completeness. W.H. Freeman. 1232–1233
- Goldberg, D.E., 1989. Genetic Algorithms in Search, Optimization and Machine Learning. Addison-Wesley Pub. Co. 1234–1235
- Gordon, A., 1988. Classification, 2nd ed. Chapman and Hall, London. 1236–1237
- Harman, M., Hierons, R., Proctor, M., 2002. A new representation and crossover operator for search-based optimization of software modularization. In: AAI Genetic and Evolutionary Computation Conference (GECCO). Springer-Verlag, New York, USA, pp. 82–87. 1238–1242
- Kaufman, L., Rousseeuw, P., 1990. Finding Groups in Data: An Introduction to Cluster Analysis. Wiley-Inter Science, Wiley, NY. 1243–1244
- Krone, M., Snelling, G., 1994. On the inference of configuration structures from source code. In: Proceedings of the 16th International Conference on Software Engineering. IEEE Computer Society Press, Los Alamitos, CA, USA, Sorrento, Italy, pp. 49–57. 1245–1248
- Kuipers, T., Moonen, L., 2000. Types and concept analysis for legacy systems. In: Proceedings of the IEEE International Workshop on Program Comprehension. IEEE Computer Society Press, Los Alamitos, CA, USA, pp. 221–230. 1249–1252

- 1253 Kuipers, T., van Deursen, A., 1999. Identifying objects using cluster
1254 and concept analysis. In: Proceedings of the International Confer-
1255 ence on Software Engineering. IEEE Computer Society Press, Los
1256 Alamitos, CA, USA, pp. 246–255.
- 1257 Lehman, M.M., Belady, L.A., 1985. Software Evolution—Processes of
1258 Software Change. Academic Press, London.
- 1259 Mahdavi, K., Harman, M., Hierons, R.M., 2003. A multiple hill
1260 climbing approach to software module clustering. In: Proceedings
1261 of IEEE International Conference on Software Maintenance,
1262 Amsterdam, The Netherlands. pp. 315–324.
- 1263 Maini, H., Mehrotra, K., Mohan, C., Ranka, S., 1994. Knowledge-
1264 based nonuniform crossover. In: IEEE World Congress on
1265 Computational Intelligence. IEEE Computer Society Press, Los
1266 Alamitos, CA, USA, pp. 22–27.
- 1267 Mancoridis, S., Mitchell, B.S., Rorres, C., Chen, Y., Gansner, E.R.,
1268 1998. Using automatic clustering to produce high-level system
1269 organizations of source code. In: Proceedings of the IEEE
1270 International Workshop on Program Comprehension. IEEE
1271 Computer Society Press, Los Alamitos, CA, USA.
- 1272 Merlo, E., McAdam, I., De Mori, R., 1993. Source code informal
1273 information analysis using connectionist model. In: Proceedings of
1274 the International Joint Conference on Artificial Intelligence. IEEE
1275 Computer Society Press, Los Alamitos, CA, USA, pp. 1339–1344.
- 1276 Mitchell, M., 1996. An Introduction to Genetic Algorithms. MIT
1277 Press, Cambridge, MA, USA.
- 1278 Neteler, M. (Ed.), 2001. GRASS 5.0 Programmer's Manual. Geo-
1279 graphic Resources Analysis Support System. ITC-irst, Italy,
1280 Available from: <<http://grass.itc.it/grassdevel.html>>.
- 1281 Neteler, M., Mitasova, H., 2002. Open Source GIS: A GRASS GIS
1282 Approach. Kluwer Academic Publishers, Boston/USA; Dordrecht/
1283 Holland; London/UK.
- 1284 Oommen, B., de St Croix, E., 1996. Graph partitioning using learning
1285 automata. IEEE Transactions on Computers 45 (2), 195–208.
- 1286 Rayside, D., Kontogiannis, K., 2002. Extracting Java library subsets
1287 for deployment on embedded systems. Science of Computer
1288 Programming 45 (2–3), 245–270.
- 1289 Shazely, S., Baraka, H., Abdel-Wahab, A., 1998. Solving graph
1290 partitioning problem using genetic algorithms. In: Midwest Sym-
1291 posium on Circuits and Systems. IEEE Computer Society Press,
1292 Los Alamitos, CA, USA, pp. 302–305.
- 1293 Siff, M., Reps, T., 1999. Identifying modules via concept analysis.
1294 IEEE Transactions on Software Engineering 25, 749–768.
- 1295 Snelting, G., 2000. Software reengineering based on concept lattices.
1296 In: Proceedings of IEEE International Conference on Software
1297 Maintenance. IEEE Computer Society Press, Los Alamitos, CA,
1298 USA, pp. 3–10.
- 1299 Talbi, E., Bessière, P., 1991. A parallel genetic algorithm for the graph
1300 partitioning problem. In: ACM International Conference on
1301 Supercomputing. ACM Press, New York, USA, Cologne,
1302 Germany.
- 1303 Tip, F., Laffra, C., Sweeney, P.F., Streeter, D., 1999. Practical
1304 experience with an application extractor for Java. ACM SIGPLAN
1305 Notices 34 (10), 292–305.
- 1306 Tonella, P., 2001. Concept analysis for module restructuring. IEEE
1307 Transactions on Software Engineering 27 (4), 351–363.
- 1308 Tzerpos, V., Holt, R.C., 1998. Software botryology: automatic
1309 clustering of software systems. In: DEXA Workshop. IEEE
1310 Computer Society Press, Los Alamitos, CA, USA, pp. 811–818.
- 1311 Tzerpos, V., Holt, R.C., 1999. MoJo: A distance metric for software
1312 clusterings. In: Proceedings of IEEE Working Conference on
Reverse Engineering. IEEE Computer Society Press, Los Alamitos,
CA, USA, pp. 187–195.
- Tzerpos, V., Holt, R.C., 2000a. ACDC: An algorithm for comprehen-
sion-driven clustering. In: Proceedings of IEEE Working Confer-
ence on Reverse Engineering. IEEE Computer Society Press, Los
Alamitos, CA, USA, pp. 258–267.
- Tzerpos, V., Holt, R., 2000b. The stability of software clustering
algorithms. In: Proceedings of the IEEE International Workshop
on Program Comprehension. IEEE Computer Society Press, Los
Alamitos, CA, USA.
- Wiggerts, T.A., 1997. Using clustering algorithms in legacy systems
remodularization. In: Proceedings of IEEE Working Conference
on Reverse Engineering. IEEE Computer Society Press, Los
Alamitos, CA, USA.
- 1313 1314 1315 1316 1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1344 1345 1346 1347 1348 1349 1350 1351 1352 1353 1355 1356 1357 1358 1359 1360 1361 1362 1363 1364 1365 1366 1367 1368 1369 1370 1372 1373 1374 1375 1376 1377 1378 1379 1380 1381 1382 1383 1384
- Massimiliano Di Penta** received his laurea degree in Computer Engi-
neering in 1999 and his PhD in Computer Science Engineering in 2001
at the University of Sannio in Benevento, Italy. Currently he is with
RCOST—Research Centre On Software Technology in the same
University. His main research interests include software maintenance
software quality, reverse engineering, program comprehension and
search-based software engineering. He is author of about 30 papers
appeared in international journals, conferences and workshops. He
serves the program and organizing committees of workshops and
conferences in the software maintenance field, such as the International
Conference on Software Maintenance, the International Workshop on
Program Comprehension, the Workshop on Source code Analysis and
Manipulation.
- Markus Neteler** received his M.Sc. degree in Physical Geography and
Landscape Ecology from the University of Hanover in Germany in
1999. He worked at the Institute of Geography as Research Scientist
and teaching associate for two years. Since 2001 he is researcher at
ITC-irst (Centre for Scientific and Technological research), Trentol
Italy since 2001. His main research interest is remote sensing for
environmental risk assessment and Free Software GIS development.
He is author of two books on the Open Source Geographical Informa-
tion System GRASS and various papers applications in GIS.
- Giuliano Antoniol** received his doctoral degree in Electronic Engi-
neering from the University of Padua in 1982. He worked at Irst for 10
years where he led the Irst Program Understanding and Reverse Engi-
neering (PURE) Project team. Giuliano Antoniol published more than
60 papers in journals and international conferences. He served as a
member of the Program Committee of international conferences and
workshops such as the International Conference on Software Main-
tenance, the International Workshop on Program Comprehension, the
International Symposium on Software Metrics. He is presently mem-
ber of the Editorial Board of the Journal Software Testing Verification
& Reliability, the Journal Information and Software Technology, the
Empirical Software Engineering and the Journal of Software Quality.
He is currently Associate Professor the University of Sannio, Faculty
of Engineering, where he works in the area of software metrics, process
modeling, software evolution and maintenance.
- Ettore Merlo** received his Ph.D. in computer science from McGill
University (Montreal) in 1989 and his Laurea degree—summa cum
laude—from University of Turin (Italy) in 1983. He has been the lead
researcher of the software engineering group at Computer Research
Institute of Montreal (CRIM) until 1993 when he joined Ecole Poly-
technique de Montreal where he is currently an associate professor. His
research interests are in software analysis, software reengineering, user
interfaces, software maintenance, artificial intelligence and bio-infor-
matics. He has collaborated with several industries and research cen-
ters in particular on software reengineering, clone detection, software
quality assessment, software evolution analysis, testing, architectural
reverse engineering and bio-informatics.