

Outi Räihä

**An Updated Survey on
Search-Based Software Design**



DEPARTMENT OF COMPUTER SCIENCES
UNIVERSITY OF TAMPERE

D-2009-5

TAMPERE 2009

UNIVERSITY OF TAMPERE
DEPARTMENT OF COMPUTER SCIENCES
SERIES OF PUBLICATIONS D – NET PUBLICATIONS
D-2009-5, AUGUST 2009

Outi Räihä

An Updated Survey on Search-Based Software Design

DEPARTMENT OF COMPUTER SCIENCES
FIN-33014 UNIVERSITY OF TAMPERE

ISBN 978-951-44-7827-7
ISSN 1795-4274

An Updated Survey on Search-Based Software Design

OUTI RÄIHÄ

Tampere University of Technology, Finland

Search-based approaches to software design are investigated. Software design is considered from a wide view, including topics that can also be categorized under software maintenance or re-engineering. Search-based approaches have been used in research from high architecture level design to software clustering and finally software refactoring. Enhancing and predicting software quality with search-based methods is also taken into account as a part of the design process. The choices regarding fundamental decisions, such as representation and fitness function, when used in meta-heuristic search algorithms, are emphasized and discussed in detail. Ideas for future research directions are also given.

Categories and Subject Descriptors: D.2.10. [**Software Engineering**]: *Design*; D.2.11. [**Software Engineering**]: *Software Architectures*; G.1.6 [**Numerical Analysis**]: *Optimization*.

General Terms: Algorithms, Design

Additional Key Words and Phrases: search-based software engineering, clustering, refactoring, software architecture design, software quality, genetic algorithm, hill climbing, simulated annealing,

1. INTRODUCTION

Traditional software engineering attempts to find solutions to problems in a variety of areas, such as testing, software design, requirements engineering, etc. A human software engineer must apply his acquired knowledge and resources to solve such complex problems that have to simultaneously meet needs but also be able to handle constraints. Often there are conflicts regarding the wishes of different stakeholders, i.e., compromises must be made with decisions regarding both functional and non-functional aspects. However, as any other engineering discipline, software engineers still attempt to find the optimal solution to any given problem, regardless of its complexity. As systems get more complex, the task of finding even a near optimal solution will become far too laborious for a human. Automating (or semi-automating) the process of finding, say, the optimal software architecture or resource allocation in a software project, can thus be seen as the ultimate dream in software engineering. Results from applications of search techniques in other engineering disciplines further support this idea, as they have been extremely encouraging.

Search-based software engineering (SBSE) applies meta-heuristic search techniques, such as genetic algorithms and simulated annealing, to software engineering problems. It stems from the realization that many tasks in software engineering can be formulated as combinatorial search problems. The goal is to find, from the wide space of possibilities, a

solution that is sufficiently good according to an appropriate quality function. Ideally this would be the optimal solution, but in reality optimality may be difficult (if not impossible) to achieve or even define due to various reasons, such as the size of the search space or the complexity of the quality function. Allowing a search algorithm to find a solution from such a wide space enables partial or full automation of previously laborious tasks, solves problems that are hard to manage by other methods, and often leads to solutions that a human software engineer might not have been able to think of.

Interest in SBSE has been growing rapidly over the past years, both in academia and industry. The combination of increased computing power, and new, more efficient, search algorithms has made SBSE a practical solution method for many problems throughout the software engineering life cycle [SSBSE, 2009]. A comprehensive review of the field is made by Harman et al. [2009], and Harman [2007] has also provided a brief overview to the current state of SBSE. Problems in the field of software engineering have been formulated as search problems by Clarke et al. [2003] and Harman and Jones [2001]. Search-based approaches have been most extensively applied in the field of software testing, and a covering survey of this branch (focusing on test data generation) has been made by McMinn [2004]. A review on SBSE concentrating on testing is also provided by Mantere and Alander [2005]. Another test related survey has been made by Afzal et al. [2008; 2009], who concentrate on testing non-functional properties. As there has been much research and previous surveys regarding the area of testing, it will be omitted from this survey, even if the studies related to testing could be considered as altering (and thus perhaps improving) a software design. As in the case of, e.g., testability transformations, Harman et al. [2004] define three critical differences to traditional transformations, one of them concerning the functionality of the program. Harman et al. [2004] state that “testability transformations need not preserve functional equivalence”, which contradicts the idea of building a design based on a fixed set of requirements.

This survey will cover the branch of software design. Software design can be defined as “the process which translates the requirements into a detailed design of a software system” [Yau and Tsai, 1986]. Here software design is considered as described by Wirfs-Brock and Johnson [1990]. Although they consider only object-oriented design, the skeleton of a process from requirements to actual design can be applied to any form of software design. A design process starts from requirements, and first enters an exploratory phase, where the fundamental structure is decided. This leads to a preliminary design, which then enters an analysis stage. After the suggested design is analyzed and modified according to the result, the final design is achieved. Following this

interpretation, software refactoring and clustering have also been taken into account as they are considered as actions of modifying (based on a certain analysis) a preliminary model, which in many cases is a working implementation.

The area of search-based software design has developed greatly in the very recent years, and is gaining an increasing interest in the SBSE community. However, although several surveys have been made of the SBSE field as a whole, they deal with the design area quite briefly. Also, the literature published from the software design perspective either does not cover search-based methods [Yau and Tsai, 1986; Budgen, 2003] or only briefly mentions the option of having an algorithm to automate class hierarchy design [Wirfs-Brock and Johnson, 1990]. Thus, there is a need to cover this crossing of two disciplines: search-based techniques and software design. New contribution is made especially in summarizing research in architecture level design that uses search-based techniques, as it has been overlooked in previous studies of search-based software engineering.

The timeline for development of SBSE as a field is presented in Figure 1. It can clearly be seen that the earliest applications have been in testing, as can be deduced from the amount of existing surveys. However, more importantly, the timeline also shows the steady increasing of ideas in the area of search based design in the past 10 years. Thus, a covering survey in this area is certainly due. All in all, the timeline shows that the SBSE has been a very active discipline in the past 20 years, as only novel ideas are presented here: countless of approaches and studies regarding these ideas have been made but not portrayed here. The explanations and references for the data points in Figure 1 are given in Figure 2.

Harman [2004] points out how crucial the representation and fitness function are in all search-based approaches to software engineering. When using genetic algorithms [Holland, 1975], which are especially popular in search-based design, the choices regarding genetic operators are just as important and very difficult to make. This survey emphasizes the choices made regarding the particular characteristics of search algorithms; any new study in the field of search-based software engineering would benefit from learning what kind of solutions have proven to be particularly successful in the past.

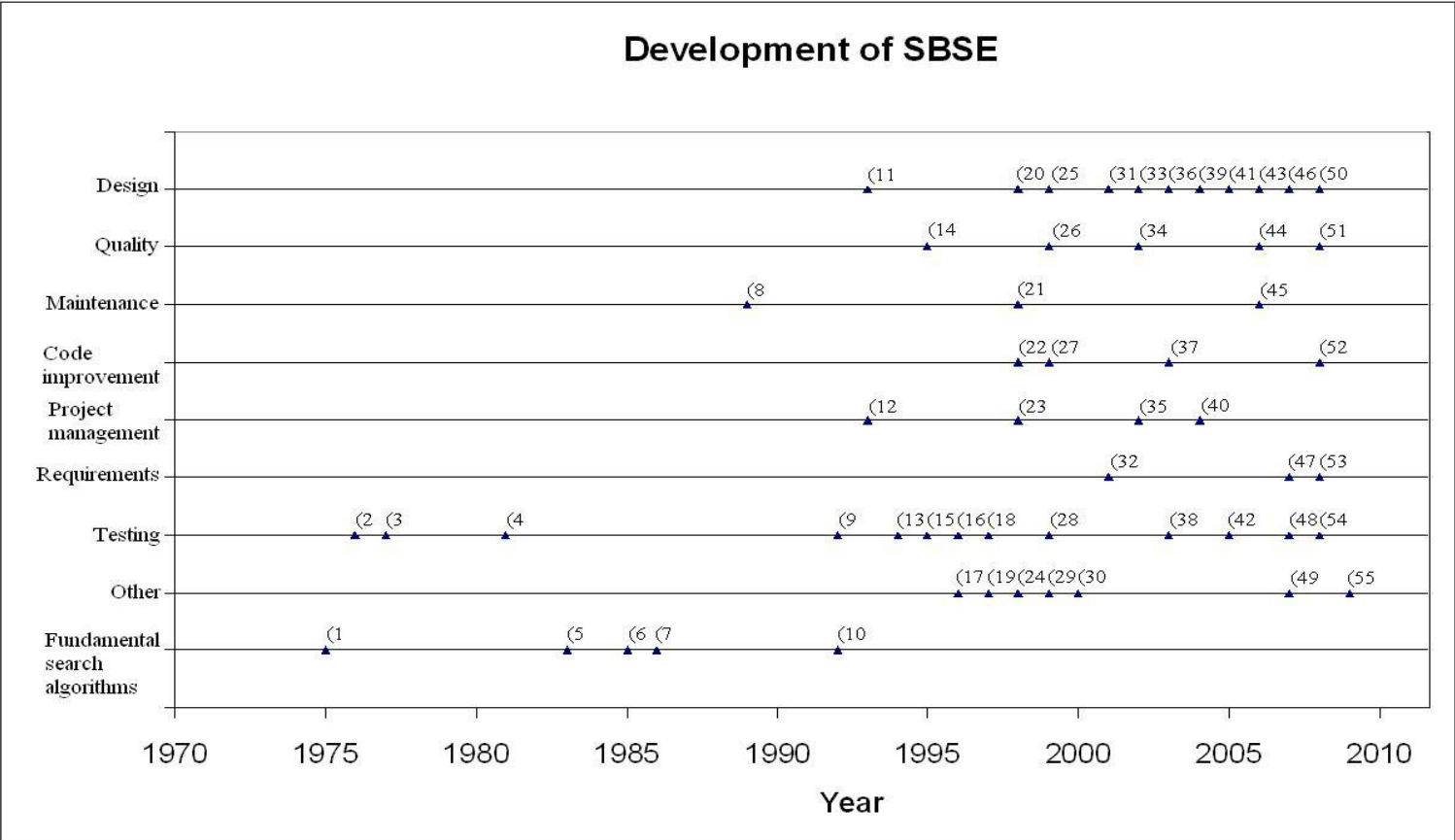


Fig. 1. Timeline of SBSE development

- 1) Genetic algorithm (GA) [Holland, 1975]
- 2) Test data automation [Miller and Spooner, 1976]
- 3) Test case automation [Fischer, 1977]
- 4) Retesting [Fischer et al. 1981]
- 5) Simulated annealing (SA) [Kirkpatrick et al., 1983]
- 6) Genetic programming (GP) [Cramer, 1985]
- 7) Tabu search [Glover, 1986]
- 8) Revalidation [Hartmann and Robson, 1989]
- 9) GAs in testing [Xanthakis et al., 1992]
- 10) Ant colony optimization (ACO) [Dorigo, 1992]
- 11) GA, constraints [Schoenauer and Xanthakis, 1993]
- 12) GA, project management [Chao et al., 1993]
- 13) GA, test data [Pei et al., 1994]
- 14) GA, reliability model [Minohara and Tohma, 1995]
- 15) Chaining approach, test data [Ferguson and Korel, 1995]
- 16) GA, structural testing [Jones et al., 1996]
- 17) GA, protocol validation [Alba and Troya, 1996]
- 18) GA, response time [Alander et al., 1997]
- 19) GA, GP, software agents [Sinclair and Shami, 1997]
- 20) Clustering [Mancoridis et al., 1998]; parallelization [Williams, 1998]
- 21) GP, software versioning [Feldt, 1998]
- 22) SA, flaw finding [Tracey et al., 1998]
- 23) Project estimation [Dolado and Fernandez, 1998]
- 24) Compiler [Nisbet, 1998]; Task scheduling [Monnier et al., 1998]
- 25) GP, re-engineering at code level [Ryan, 1999]
- 26) GP, quality determination [Evelt et al., 1999]
- 27) GA, reduced code space [Cooper et al., 1999]
- 28) SA, regression testing [Mansour and El-Fakih, 1999]
- 29) GA, Protocols for distributed applications [El-Fakih et al., 1999]
- 30) Secure protocols [Clark and Jacob, 2000]
- 31) GA, hierarchical decompositions [Lutz, 2001]
- 32) Next release problem [Bagnall et al., 2001]
- 33) GA, reverse engineering at architecture level [Mitchell et al., 2002]
- 34) GA, combining quality predictive models [Bouktif et al., 2002]
- 35) GP, project effort estimation [Shan et al., 2002]
- 36) Multiple hill climbing, clustering [Mahdavi et al., 2003]
- 37) GA, code transformations [Fatiregun et al., 2003]
- 38) SA, test suites [Cohen et al., 2003]
- 39) Architecture relations [Mitchell et al., 2004]; service composition [Canfora et al., 2004]
- 40) Project resource allocation [Antoniol et al., 2004]
- 41) Amorphous slicing [Fatiregun et al., 2005]
- 42) ACO, testing [Li and Lam, 2005]
- 43) GA, design patterns [Amoui et al., 2006]
- 44) SA, quality prediction [Bouktif et al., 2006]
- 45) GA, software integration [Yang et al., 2006]
- 46) Use case -based design [Simons and Parmee, 2007a; 2007b; 2008]; GA, repackaging [Bodhuin et al., 2007]; Pareto optimal refactoring [Harman et al., 2007]
- 47) Multiobjective next release problem [Zhang et al., 2007]
- 48) ACO, model checking [Alba and Chicano, 2007]; GP, model checking [Johnson, 2007]; Pareto optimality, test cases [Yoo and Harman, 2007]
- 49) GA, code author identification [Lange and Mancoridis, 2007]
- 50) Class responsibility assignment [Bowman et al., 2008]; Software behavior modeling [Goldsby and Cheng, 2008]; Architecture design [Räihä et al., 2008]; model transformations by GA [Räihä et al., 2008] and particle swarm optimization [Kessentini et al., 2008]
- 51) Software verification [Shyang et al., 2008]
- 52) Co-evolution, bug fixing [Arcuri and Yao, 2008]
- 53) Requirements optimization [Zhang et al., 2008]
- 54) Tabu search, testing [Diaz et al., 2008]
- 55) GA, decision making in autonomic computing systems [Ramirez et al., 2009]

Fig. 2. References for timeline data points

This survey proceeds as follows. Section 2 describes search algorithms, and the underlying concepts for genetic algorithms, simulated annealing and hill climbing are discussed in detail. Different ways of performing the exploratory phase of design are then presented as ways for software architecture design (object-oriented and service-oriented) in Section 3. Sections 4, 5 and 6 deal with clustering, refactoring and software quality, respectively, which can all be seen as components of the analysis phase, starting from higher level re-design (clustering), going to low-level re-design (re-factoring) and finally pure analysis. The background for each underlying problem is first presented, followed by recent approaches applying search-based techniques to the problem. Summarizing remarks and a summary table of the studies is presented after each subsection. Finally, some ideas for future work are given in Section 7, and conclusions are presented in Section 8.

2. SEARCH ALGORITHMS

Meta-heuristics are commonly used for combinatorial optimization, where the search space can become especially large. Many practically important problems are NP-hard, and thus, exact algorithms are not possible. Heuristic search algorithms handle an optimization problem as a task of finding a “good enough” solution among all possible solutions to a given problem, while meta-heuristic algorithms are able to solve even the general class of problems behind the certain problem. A search will optimally end in a global optimum in a search space, but at the very least it will give some local optimum, i.e., a solution that is “better” than a significant amount of alternative solutions nearby. A solution given by a heuristic search algorithm can be taken as a starting point for further searches or be taken as the “best” possible solution, if its quality is considered high enough. For example, simulated annealing can be used to produce seed solutions for a genetic algorithm that constructs the initial population based on the provided seeds.

In order to use search algorithms in software engineering, the first step is that the particular software engineering problem should be defined as a search problem. If this cannot be done, search algorithms are most likely not the best way to solve the problem, and defining the different parameters and operations needed for the search algorithm can be difficult. After this has been done, a suitable algorithm can be selected and the issues regarding that algorithm must be dealt with.

There are three common issues that need to be dealt with any search algorithm: 1. encoding the solution, 2. defining transformations, and 3. measuring the “goodness” of a

solution. All algorithms need the solution to be encoded according to the algorithm's specific needs. For example, in order for the genetic algorithm (GA) to operate, the encoding should be done in such a way that it can be seen as a chromosome consisting of a set of genes. However, for the hill climbing (HC), any encoding where a neighborhood can be defined is sufficient. The importance and difficulty of encoding a solution increase as the complexity of the problem at hand increases. In this case complexity refers to how easily a solution can be defined, rather to the computational complexity of the problem itself. For example, a job-shop problem may be computationally complex, but the solution candidates are simple to encode as an integer array. However, a solution containing, e.g., all the information regarding a software architecture, is demanding to encode so that: 1. all information stays intact, 2. operations can efficiently be applied to the selected encoding of the solution, 3. the fitness evaluations can be performed efficiently, and 4. there is minimal need for "outside" data, i.e., data structures containing information about the solution that are not included in the actual encoding.

Defining a neighborhood is crucial to all algorithms; HC, simulated annealing (SA) and tabu search operate purely on the basis of moving from one solution to its neighbor. A neighbor is achieved by some operation that transforms the solution. These operations can be seen equivalent to the mutations needed by the GA.

Finally, the most important and difficult task is defining a fitness function. If defining the fitness function fails, the search algorithm will not be guided towards the desired solutions. All search algorithms require this quality function to evaluate the "goodness" of a solution in order to compare solutions and thus guide the search.

To understand the basic concepts behind the approaches presented here, the most commonly used search algorithms are briefly introduced. The most common approach is to use genetic algorithms. Hill climbing and its variations, e.g., multi-ascent hill climbing (MAHC), is also quite popular due to its simplicity. Finally, several studies use simulated annealing. In addition to these algorithms, tabu search is a widely known meta-heuristic search technique, and genetic programming (GP) [Koza, 1992] is commonly used in problems that can be encoded as trees. For a detailed description on GA, see Mitchell [1996] or Sivanandam and Deepa [2007], for SA, see, e.g., Reeves [1995], and for HC, see Glover and Kochenberger [2003], who also cover a wide range of other meta-heuristics. For a description on multi-objective optimization with evolutionary algorithms, see Deb [1999] or Fonseca and Fleming [1995]. A survey on model-based search, covering several meta-heuristic algorithms is also made by Zlochin et al. [2004].

2.1 Genetic algorithms

Genetic algorithms were invented by John Holland in the 1960s. Holland's original goal was not to design application specific algorithms, but rather to formally study the ways of evolution and adaptation in nature and develop ways to import them into computer science. Holland [1975] presents the genetic algorithm as an abstraction of biological evolution and gives the theoretical framework for adaptation under the genetic algorithm [Mitchell, 1996].

In order to explain genetic algorithms, some biological terminology needs to be clarified. All living organisms consist of cells, and every cell contains a set of chromosomes, which are strings of DNA and give the basic information of the particular organism. A chromosome can be further divided into genes, which in turn are functional blocks of DNA, each gene representing some particular property of the organism. The different possibilities for each property, e.g., different colors of the eye, are called alleles. Each gene is located at a particular locus of the chromosome. When reproducing, crossover occurs: genes are exchanged between the pair of parent chromosomes. The offspring is subject to mutation, where single bits of DNA are changed. The fitness of an organism is the probability that the organism will live to reproduce and carry on to the next generation [Mitchell, 1996]. The set of chromosomes at hand at a given time is called a population.

Genetic algorithms are a way of using the ideas of evolution in computer science. When thinking of the evolution and development of species in nature, in order for the species to survive, it needs to develop to meet the demands of its surroundings. Such evolution is achieved with mutations and crossovers between different chromosomes, i.e., individuals, while the fittest survive and are able to participate in creating the next generation.

In computer science, genetic algorithms are used to find a good solution from a very large search space, the goal obviously being that the found solution is as good as possible. To operate with a genetic algorithm, one needs an encoding of the solution, i.e., a representation of the solution in a form that can be interpreted as a chromosome, an initial population, mutation and crossover operators, a fitness function and a selection operator for choosing the survivors for the next generation. Algorithm 1 gives the pseudo code for a genetic algorithm.

Algorithm 1 geneticAlgorithm

Input: formalization of solution, *initialSolution*
chromosomes ← createPopulation(*initialSolution*)
while NOT *terminationCondition* **do**
 foreach *chromosome* **in** *chromosomes*

```

     $p \leftarrow \text{randomProbability}$ 
    if  $p > \text{mutationProbability}$  then
        mutate(chromosome)
    end if
end for
foreach chromosomePair in chromosomes
     $cp \leftarrow \text{randomProbability}$ 
    if  $cp > \text{crossoverProbability}$  then
        crossover(chromosomePair)
        addOffspringToPopulation()
    end if
end for
foreach chromosome in chromosomes
    calculatefitness(chromosome)
end for
    selectNextPopulation()
end while

```

As discussed, correctly defining the different operations (mutations, crossover and fitness function) is vital in order to achieve satisfactory results. However, as seen in Algorithm 1, there are also many parameters regarding the GA that need to be defined and greatly affect the outcome. These parameters are the population size, number of generations (often used as the terminating condition) and the mutation and crossover probabilities. Having a large enough population ensures variability within a generation, and enables a wide selection of different solutions at every stage of evolution. However, at a certain point the results start to converge, and a larger population always means more fitness evaluations and thus requires more computation time. Similarly, the more generations the algorithm is allowed to evolve for, the higher the chances are that it will be able to reach the global optimum. However, again, letting an algorithm run for, say, 10 000, generations will most probably not be beneficial, as if the operations and parameters have been chosen correctly, a reasonably good optimum should have been found much earlier. Mutation and crossover probabilities both affect how fast the population evolves. If the probabilities are too high, there is the risk that the implementation of genetic operations becomes random instead of guided. Vice versa, if the probabilities are too low there is the risk that the population will evolve too slowly, and no real diversity will exist. A theory to be noted with genetic operators is the building block hypothesis, which states that a genetic algorithm combines a set of sub-solutions, or building blocks, to obtain the final solution. The sub-solutions that are kept over the generations generally have an above-average fitness [Salomon, 1998]. The crossover

operator is especially sensitive to this hypothesis, as an optimal crossover would thus combine two rather large building blocks in order to produce an offspring with a one-point crossover.

2.2 Simulated annealing

Simulated annealing is originally a concept in physics. It is used when the cooling of metal needs to be stopped at given points where the metal needs to be warmed a bit before it can resume the cooling process. The same idea can be used to construct a search algorithm. At a certain point of the search, when the fitness of the solution in question is approaching a set value, the algorithm will briefly stop the optimizing process and revert to choosing a solution that is not the best in the current solution's neighborhood. This way getting stuck to a local optimum can effectively be avoided. Since the fitness function in simulated annealing algorithms should always be minimized, it is usually referred to as a cost function [Reeves, 1995].

Simulated annealing usually begins with a point x in the search space that has been achieved through some heuristic method. If no heuristic can be used, the starting point will be chosen randomly. The cost value c , given by cost function E , of point x is then calculated. Next a neighboring value x_1 is searched and its cost value c_1 calculated. If $c_1 < c$, then the search moves onto x_1 . However, even though $c \leq c_1$, there is still a chance, given by probability p , that the search is allowed to continue to a solution with a bigger cost [Clarke et al., 2003]. The probability p is a function of the change in cost function ΔE , and a parameter T :

$$p = e^{-\Delta E/T}.$$

This definition for the probability of acceptance is based on the law of thermodynamics that controls the simulated annealing process in physics. The original function is

$$p = e^{-\Delta E/kt},$$

where t is the temperature in the point of calculation and k is Boltzmann's constant [Reeves, 1995].

The parameter T that substitutes the value of temperature and the physical constant is controlled by a cooling function C , and it is very high in the beginning of simulated annealing and is slowly reduced while the search progresses [Clarke et al., 2003]. The actual cooling function is application specific.

If the probability p given by this function is above a set limit, then the solution is accepted even though the cost increases. The search continues by choosing neighbors and applying the probability function (which is always 1 if the cost decreases) until a cost value is achieved that is satisfactory low. Algorithm 2 gives the pseudo code for a

simulated annealing algorithm.

Algorithm 2 simulatedAnnealing

Input: formalization of solution, *initialSolution*, cooling ratio α , initial temperature T_0 , frozen temperature T_f , and temperature constant r

Output: optimized solution *finalSolution*

initialQuality \leftarrow evaluate(*initialSolution*)

\leftarrow *initialSolution*

$Q_1 \leftarrow$ *initialQuality*

$T \leftarrow T_0$

while $T_0 > T_f$ **do**

$r_i \leftarrow 0$

while $r_i > r$ **do**

$S_i \leftarrow$ findNeighbor(S_1)

$Q_i \leftarrow$ evaluate(Q_1)

if $Q_i > Q_1$ **then**

$S_1 \leftarrow S_i$

$Q_1 \leftarrow Q_i$

else

$\delta \leftarrow Q_1 - Q_i$

$p \leftarrow$ randomProbability

if $p < e^{-\delta/T}$ **then**

$S_1 \leftarrow S_i$

$Q_1 \leftarrow Q_i$

end if

end if

$r_i \leftarrow r_i + 1$

end while

$T \leftarrow T * \alpha$

end while

return S_1

The key parameters to be adjusted for SA are the initial temperature, the cooling ratio and the temperature constant. These all combined affect how fast the cooling happens. If the cooling is too fast, the algorithm may not have sufficient time to achieve an optimum. However, if the cooling is too slow, the initial temperature may need a significantly high value so that the solution will be able to evolve enough (i.e., noticeably transform from the initial solution) before reaching the frozen temperature.

2.3 Hill climbing

Hill climbing begins with a random solution, and then begins to search through its neighbors for a better solution. There are several versions of how this is done; in some versions the algorithm moves on after finding the first neighbor that is better than the

current, some do a fixed number of neighbor evaluations and continue to the best of this group, and some versions go through the entire neighborhood of a solution and select the best neighbor from which the procedure is continued. Algorithm 3 adopts the last option, i.e., the entire neighborhood is evaluated before moving on. Hill climbing does not include any mechanisms to avoid getting stuck with a local optimum.

There are three critical choices regarding HC: 1. defining a neighborhood for each solution, 2. defining an evaluation function for a solution, and 3. defining by what extent each neighborhood is searched. If the problem at hand is very complex and each solution has an exponential amount of neighbors, traversing through each neighborhood maybe extensively time consuming. However, if the subgroup of neighbors to be examined is chosen wisely, the actual outcome of the algorithm may still be good enough, while much time is saved when not every solution needs to be evaluated.

Algorithm 3 hillClimbing

Input: formalization of solution, *initialSolution*
currentSolution ← *initialSolution*
currentFitness ← evaluate(*currentSolution*)
while betterNeighborsExist **do**
 neighborhood ← findNeighbors(*currentSolution*)
 foreach *neighbor* **in** *neighborhood*
 neighborFitness ← evaluate(*neighbor*)
 if *neighborFitness* > *nextFitness* **then**
 nextSolution ← *neighbor*
 nextFitness ← *neighborFitness*
 end if
 end for
 if *nextFitness* > *currentFitness* **then**
 currentSolution ← *nextSolution*
 else
 termination
 return *currentSolution*
 end if
end while

3. SOFTWARE ARCHITECTURE DESIGN

The core of every software system is its architecture. Designing software architecture is a demanding task requiring much expertise and knowledge of different design alternatives, as well as the ability to grasp high-level requirements and piece them to detailed architectural decisions. In short, designing software architecture takes verbally formed functional and quality requirements and turns them into some kind of formal model, which is used as a base for code. Automating the design of software is obviously a

complex task, as the automation tool would need to understand intricate semantics, have access to a wide variety of design alternatives, and be able to balance multi-objective quality factors. From the re-design perspective, program comprehension is one of the most expensive activities in software maintenance. The following sections describe meta-heuristic approaches to software architecture design for object-oriented and service-oriented architectures.

3.1 Object-oriented architecture design

3.1.1 Basics

At its simplest, object-oriented design deals with extracting concepts from, e.g., use cases, and deriving methods and attributes, which are distributed into classes. A further step is to consider interfaces and inheritance. A final design can be achieved through the implementation of architecture styles [Shaw and Garlan, 1996] and design patterns [Gamma et al., 1995]. When attempting to automate the design of object-oriented architecture from concept level, the system requirements must be formalized. After this, the major problem lies within quality evaluation, as many design decisions improve some quality attribute [Losavio et al., 2004] but weaken another. Thus, a sufficient set of quality estimators should be used, and a balance should be found between them. Re-designing software architectures automatically is slightly easier than building architecture from the very beginning, as the initial model already exists and it merely needs to be ameliorated. However, implementing design patterns is never straightforward, and measuring their impact on the quality of the system is difficult. For more background on software architectures, see, e.g., Bass et al. [1998].

Approaches to search-based software design are presented starting from low-level approaches, i.e., what is needed when first beginning the architecture design, to high-level approaches, ending with analyzing software architecture. Object-oriented architecture design begins with use cases and assigning responsibilities, i.e., methods and attributes to classes [Bowman et al., 2008; Simons and Parmee 2007a; Simons and Parmee, 2007b]. After the basic structure, the architecture can be further designed by applying design patterns, either on an existing system [Amoui et al., 2006] or building the design patterns in the system from the very beginning [Räihä et al., 2008a; Räihä et al., 2008b; Räihä et al. 2009]. If an idea for an optimal solution is available, model transformations can be sought to achieve that solution [Kessentini et al., 2008]. There might also be many choices regarding the components of the architecture, depending on the needs of the system. An architecture can be made of alternative components [Kim and Park, 2009] or a subsystem can be sought after [Bodhuin et al., 2007]. Studies have also

been made on identifying concept boundaries and thus automating software comprehension [Gold et al., 2006] and composing behavioral models for autonomic systems [Goldsby and Chang, 2008; Goldsby et al., 2008], which give a dynamic view of software architecture. One of the most abstract studies attempts to build hierarchical decompositions for a software system [Lutz, 2001, which already comes quite close to software clustering. Summarizing remarks of the approaches are given in the end, and the fundamentals of each study are collected in Table 1.

3.1.2 Approaches

Bowman et al. [2008] study the use of a multi-objective genetic algorithm (MOGA) in solving the class responsibility assignment problem. The objective is to optimize the class structure of a system through the placement of methods and attributes. The strength Pareto approach (SPEA2) is used, which differs from a traditional GA by containing an archive of individuals from past populations. This approach combines several aspects that aid in finding the truly optimal individuals and thus leaves less room for GA “to err” in terms of undesired mutations or overly relying on metrics.

The chromosome is represented as an integer vector. Each gene represents a method or an attribute in the system and the integer value in a gene represents the class to which the method or attribute in that locus belongs. Dependency information between methods and attributes is stored in a separate matrix. Mutations are performed by simply changing the class value randomly; the creation of new classes is also allowed. Crossover is the traditional one-point one. There are also constraints: no empty classes are allowed (although the selected encoding method also makes them impossible), conceptually related methods are only moved in groups, and classes must have dependencies to at least one other class.

The fitness function is formed of five different values measuring cohesion and coupling: 1. method-attribute coupling, 2. method-method coupling, 3. method-generalization coupling, 4. cohesive interaction and 5. ratio of cohesive interaction. A complementary measure for common usage is also used. Selection is made with a binary-tournament selection where the fitter individual is selected 90% of the time.

In the case study an example system is used, and a high-quality UML class diagram of this system is taken as a basis. Three types of modifications are made and finally the modifications are combined in a final test. The efficiency of the MOGA is now evaluated in relation to how well it fixed the changes made to the optimal system. Results show that in most cases the MOGA managed to fix the made modifications and in some cases the resulting system also had a higher fitness value than the original “optimal” system.

Bowman et al. also compare MOGA to other search algorithms, such as random search, hill climbing and a simple genetic algorithm. Random search and hill climbing only managed to fix a few of the modifications and the simple GA did not manage to fix any of the modifications. Thus, it would seem that a more complex algorithm is needed for the class responsibility assignment problem.

The need for highly developed algorithms is further high-lighted when noting that a ready system is being ameliorated instead of completely automating the class responsibility assignment. As a ready system can be assumed to have some initial quality and conceptually similar methods and attributes are already largely grouped, it does help the algorithm when re-assigning the moved methods and attributes. This is due to the fact that by attempting to re-locate the moved method or attribute to the “wrong” class, the fitness value will be significantly lower than when assigning the method or attribute to the “right” class.

Simons and Parmee [2007a; 2007b; 2008] take use cases as the starting point for system specification. Data is assigned to attributes and actions to methods, and a set of uses is defined between the two sets. The notion of class is used to group methods and attributes. Each class must contain at least one attribute and at least one method. Design solutions are encoded directly into an object-oriented programming language. This approach starts with pure requirements and leaves all designing to the algorithm, making the problem of finding an optimal class structure extremely more difficult than in cases where a ready system can be used as basis.

A single design solution is a chromosome. In a mutation, a single individual is mutated by locating an attribute and a method from one class to another. For crossover two individuals are chosen at random from the population and their attributes and methods are swapped based on their class position within the individuals. Cohesiveness of methods (COM) is used to measure fitness, fitness for class C is defined as $f(C) = 1/(|Ac||Mc|) * \sum(\Delta_{ij})$, where Ac (respectively Mc) stands for the number of attributes (respectively methods) in class C, and $\Delta_{ij} = 1$, if method j uses attribute i , and 0 otherwise. Selection is performed by tournament and roulette-wheel. The choices regarding encoding, genetic operators and fitness function are quite traditional, although the problem to be solved is far from traditional.

In an alternative approach, categorized by the authors as *evolutionary programming (EP)* and inspired by Fogel et al. [1966], offspring is created by mutation and selection is made with tournament selection. Two types of mutations are used, class-level mutation and element-level mutation. At class level, all attributes and methods of a class in an

individual are swapped as a group with another class selected at random. At element level, elements (methods and attributes) in an individual are swapped at random from one class to another. Initialization of the population is made by allocating a number of classes to each individual design at random, within a range derived from the number of attributes and methods. All attributes and methods from sets of attributes and methods are then allocated to classes within individuals at random. These operations appear quite simplistic, and the actual change to the design remains minimal, since the fitness of an individual depends on how methods and attributes depending on one another are located. When the elements are moved in a group, there does not seem to be very much change in the actual design.

A case study is made with a cinema booking system with 15 actions, 16 datas and 39 uses. For GA, the average COM fitness for final generation for both tournament and roulette-wheel is similar, as is the average number of classes in the final generation. However, convergence to a local optimum is quicker with tournament selection. Results reveal that the average and maximum COM fitness of the GA population with roulette-wheel selection lagged behind tournament in terms of generation number. For EP, the average population COM fitness in the final generation is similar to that achieved by the GA.

The initial average fitness values of the three algorithms are notably similar, although the variance of the values increases from GA tournament to GA roulette-wheel to EP. In terms of COM cohesion values, the generic operators produced conceptual software designs of similar cohesion to human performance. Simons and Parmee suggest that a multi-objective search may be better suited for support of the design processes of the human designer. To take into account the need for extra input, they attempted to correct the fitness function by multiplying the COM value by a) the number of attributes and methods in the class ($COM \cdot (M+A)$); b) the square root of the number of attributes and methods in the class ($COM \cdot \sqrt{(M+A)}$); c) the number of uses in the class ($COM \cdot \text{uses}$) and d) the square root of the number of uses in a class ($COM \cdot \sqrt{\text{uses}}$). Using such multipliers raises some questions as there is no intuition for using the square root multipliers. Multiplying by the sum of methods and attributes or uses can intuitively be justified by showing more appreciation to classes that are large but are still comprehensible. However, such appreciation may lead to preferring larger classes.

The authors have taken this into account by measuring the number of classes in a design solution and a design solution with higher number of classes is preferred to a design solution with fewer classes. When cohesion metrics that take class size into

account are used, there is a broad similarity between the average population cohesion fitness and the manual design. Values achieved by the COM.M+A and COM.uses and cohesion metrics are higher than the manual design cohesion values, while COM. $\sqrt{M+A}$ and COM. \sqrt{uses} values are lower. Manually examining the design produced by the evolutionary runs, a difference is observed in the design solutions produced by the four metrics that account for class size, when compared with the metrics that do not. From the results produced for the two case studies, it is evident that while the cohesion metrics investigated have produced interesting cohesive class design solutions, they are by no means a complete reflection of the inherently multi-objective evaluations conducted by a human designer. The evolutionary design variants produced are thus highly dependent on the extent and choice of metrics employed during search and exploration. These results further emphasize the importance of properly defining a fitness function and deciding on the appropriate metrics in all software design related problems.

Amoui et al. [2006] use the GA approach to improve the reusability of software by applying architecture design patterns to a UML model. The authors' goal is to find the best sequence of transformations, i.e., pattern implementations. Used patterns come from the collection presented by Gamma et al. [1995], most of which improve the design quality and reusability by decreasing the values of diverse coupling metrics while increasing cohesion.

Chromosomes are an encoding of a sequence of transformations and their parameters. Each individual consists of several *supergenes*, each of which represents a single transformation. A supergene is a group of neighboring genes on a chromosome which are closely dependent and are often functionally related. Only certain combinations of the internal genes are valid. Invalid patterns possibly produced through mutations or crossover are found and discarded. The supergene concept introduced here is an insightful approach into handling masses of complex data that needs to be represented as a relatively simple form. Instead of having only one piece of information per gene, this way several pieces of related information can be grouped to such supergenes, which then logically form a chromosome. In the study by Bowman et al. [2008] the need for additional data storage (the matrix for data dependencies) demonstrates the complexity of design problems. In this case the supergene approach introduced by Amoui et al. [2006] could have been worth while to try to include all information regarding the attributes and methods in the chromosome encoding.

Mutation randomly selects a supergene and mutates a random number of genes inside the supergene. After this, validity is checked. In case of encountering a transformed

design which contradicts with object-oriented concepts, for example, a cyclic inheritance, a zero fitness value is assigned to chromosome. This is an interesting way of dealing with anomalies; instead of implementing a corrective operation to force validity, it is trusted that the fitness function will suffice in discarding the unsuitable individuals if they are given a low enough value.

Two different versions of crossover are used. First is a single-point crossover applied at supergene level, with a randomly selected crossover point, which swaps the supergenes beyond the crossover point, but the internal genes of supergenes remain unchanged. This combines the promising patterns of two different transformation sequences. The second crossover randomly selects two supergenes from two parent chromosomes, and similarly applies single point crossover to the genes inside the supergenes. This combines the parameters of two successfully applied patterns. The first crossover thus attempts to preserve high-level building blocks, while the second version attempts to create low-level building blocks.

The quality of the transformed design is evaluated, as introduced by Martin [2000], by its “distance from the main sequence” (D), which combines several object-oriented metrics by calculating abstract classes’ ratio and coupling between classes, and measures the overall reusability of a system.

A case study is made with a UML design extracted of some free, open source applications. The GA is executed in two versions. In one version only the first crossover is applied and in second both crossovers are used. A random search is also used to see if the GA outperforms it. Results demonstrate that the GA finds the optimal solution much more efficiently and accurately. From the software design perspective, the transformed design of the best chromosomes are evolved so that abstract packages become more abstract and concrete packages in turn become more concrete. The results suggest that GA is a suitable approach for automating object-oriented software transformations to increase reusability. As the application of design patterns is by no means an easy task, these initial results suggest that at least the structure and needs of the GA does not restrict automated design of software architecture.

Räihä et al. [2008a] take the design of software architecture a step further than Simons and Parmee [2007a] by starting the design from a responsibility dependency graph. The graph can also be achieved from use cases, but the architecture is developed further than the class distribution of actions and data. A GA is used for the automation of design.

In this approach, each responsibility is represented by a supergene and a chromosome

is a collection of supergenes. The supergene contains information regarding the responsibility, such as dependencies of other responsibilities, and evaluated parameters such as execution time and variability. Here the notion of supergene [Amoui et al., 2006] is efficiently used in order to store a large amount of different types of data pieces within the chromosome. Mutations are implemented as adding or removing an architectural design pattern [Gamma et al. 1995] or an interface, or splitting or joining class(es). Implemented design patterns are Façade and Strategy, as well as the message dispatcher architecture style [Shaw and Garlan, 1996]. Dynamic mutation probabilities are used to encourage the application of basic design choices (the architectural style(s)) in the beginning and more refined choices (such as the Strategy pattern) in the end of evolution. Crossover is a standard one-point crossover. The offspring and mutated chromosomes are always checked after the operations for legality, as design patterns may easily be broken. Selection is made with the roulette wheel method.

This approach actually combines the class responsibility assignment problem studied by Simons and Paremee [2007a; 2007b] and applying design patterns, as studied by Amoui et al. [2006]. Although the selection of design patterns is smaller, the search problem of finding an optimal architecture is much more difficult. First the GA needs to find the optimal class responsibility distribution, and then apply design patterns. In this case the search space grows exponentially, as in order to optimally apply the design patterns, the class responsibility distribution may need to be sub-optimal. This produces a challenge when deciding on the fitness function.

The fitness function is a combination of object-oriented software metrics, most of which are from the Chidamber and Kemerer [1994] collection, which have been grouped to measure quality concepts efficiency and modifiability. Some additional metrics have also been developed to measure the effect of communicating through a message dispatcher or interfaces. Furthermore, a complexity measure is introduced. The fitness function is defined as $f = w_1\text{PositiveModifiability} - w_2\text{NegativeModifiability} + w_3\text{PositiveEfficiency} - w_4\text{NegativeEfficiency} - w_5\text{Complexity}$, where w_i s are weights to be fixed. As discussed, defining the fitness function is the most complex task in all SSBSE problems. In this case, when the problem is so diverse, the fitness function is also intricate: it requires a set of known metrics, a set of special metrics, the grouping of these metrics and additionally weights in order to set preferences to quality aspects.

The approach is tested on a sketch of a medium-sized system [Räihä, 2008]. Results show positive development in overall fitness value, while the balancing of weights greatly affects whether the design is more modifiable or efficient. However, the actual

designs are not compliant with the fitness values, and would not be accepted by a human architect. This suggests that further improvement is needed in defining the fitness function.

Räihä et al. [2008b] further develop their work by implementing more design patterns and an alternative approach. In addition to the responsibility dependency graph, a domain model may be given as input. The GA can now be utilized in Model Driven Architecture design, as it takes care of the transformations from Computationally Independent Model to Platform Independent Model. The new design patterns are Mediator and Proxy, and the service oriented architecture style is also implemented by enabling a class to be called through a server. The chromosome representation, mutation and crossover operations and selection method are kept the same. Results show that the fitness values converge to some optima and reasonable high-level designs are obtained.

In this case the task for the GA is made somewhat easier, as a skeleton of a class structure is given to the algorithm in the form of a domain model. This somewhat eliminates the class responsibility assignment problem and the GA can only concentrate on applying the design patterns. As the results are significantly better, although the search space is more complex when more patterns have been added to the mutations, this suggests that the class responsibility assignment problem is extremely complex on its own, and more research on this would be highly beneficial as a background for several search-based software design related questions.

Räihä et al. [2009] keep developing their approach by including the Template pattern to the design pattern/mutation collection and introducing scenarios as a way to enhance the evaluation of a produced architecture. Scenarios are basically a way to describe an interaction between the system and a stakeholder. In their work, Räihä et al. categorize and formalize modifiability related scenarios so that they can be encoded and given to the GA as an additional part of the fitness function. Each scenario is given a list of preferences regarding the architectural structures that are suitable for that scenario. The preferences are then compared with the suggested architecture and a fitness value is calculated according to how well the given architecture conforms to the preferences. This way the fitness value is more pointed as the most critical parts of the architecture can be given extra attention and the evaluation is not completely based on general metrics. Results from empirical studies made on two sample systems show that when the scenarios are used, the GA retains the high-speed phase of developing the architecture for 10 to 20 generations longer than in the case where scenarios are not used. Also, when the scenario fitness is not included in the overall fitness evaluations the GA tends to make

decisions that do not support the given scenarios.

Results from this study shows that when the modifications are as detailed as applying a design pattern (rather than modifying the architecture “as a whole”), the fitness function also needs to be more pin-pointed to study the places of an architecture where such detailed solutions would be most beneficial.

Kessentini et al. [2008] also use a search-based approach to model transformations. They start with a small set of examples from which transformation blocks are extracted and use particle swarm optimization (PSO) [Kennedy and Eberhart, 1995]. A model is viewed as a triple of source model, target model and mapping blocks between the source and target models. The source model is formed by a set of constructs. The transformation is only coherent if it does not conflict the constructs. The transformation quality of a source model (i.e., global quality of a model) is the sum of the transformation qualities of its constructs (i.e., local qualities). This approach is less automated, as the transformations need to be extracted from ready models, and are not general. However, using PSO is especially interesting, and suggests that other algorithms besides GA are also suitable for complex software design problems.

To encode a transformation, an M -dimensional search space is defined, M being the number of constructs. The encoding is now an M -dimensional integer vector whose elements are the mapping blocks selected for each construct. The fitness function is a sum of constructs that can be transformed by the associated blocks multiplied by relative numbers of matched parameters and constructs. The fitness value is normalized by dividing it with 2^*M , thus resulting in a fitness range of [0, 1].

The method was evaluated and experimented with 10 small-size models, of which nine are used as a training set and one as the actual model to be transformed. The precision of model transformation (number of constructs with correct transformations in relation to total number of constructs) is calculated in addition to the fitness values. The best solution was found already after 29 iterations, after which all particles converged to that solution. The test generated 10 transformations. The average precision of these is more than 90%, thus indicating that the transformations would indeed give an optimal result, as the fitness value was also high within the range. The test also showed that some constructs were correctly transformed although there were no transformation examples available for these particular constructs.

Kim and Park [2009] use GAs to dynamically choose components to form software architecture according to changing demands. The basic concept is to have a set of interchangeable components (e.g., BasicUI and RichUI), which can be selected according

to user preferences. The goal is thus to select an optimal architectural instance from all possible instances. This is especially beneficial when the software needs to be transferred, e.g., from a PC to a mobile device.

A softgoal interdependency graph (SIG) is used as a basis for the problem; it represents relationships between quality attributes. The quality attributes are formulated by a set of quality variables. A utility function is used to measure the user's overall satisfaction: the user now gives weights for the quality values to represent their priority. Functional alternatives (i.e., the interchangeable components) are denoted by operationalizing goals. The operationalizing goals can have an impact on a softgoal, i.e., a quality attribute. Alternatives with similar characteristics are grouped by a type. One alternative type corresponds to one architectural decision variable. These represent partial configurations of the application. A combination of architectural decision variables comprises an architectural instance.

In addition to the SIG, situation variables and their values are needed as input. Situation variables describe partial information on environmental changes and determine the impacts that architectural decision variables have on the quality attributes. The impact is defined as a situation evaluation function, which is defined for each direct interdependency between an operationalizing goal and quality attribute. Although the fitness function is quite standard, i.e., it calculates the quality through "quality values" and there are weights assigned, the actual computations are not that straightforward. The quality attributes the fitness function is based on rely on decision variables and situation variables. These in turn need to be calculated by hand, and there is no clear answer to how the situation variables themselves are gathered.

For the GA, the architectural instance is encoded as a chromosome by using a string of integers representing architectural decisions. Mutation is applied to offspring, for which each digit is subjected to mutation (according to mutation probability). Crossover is a standard two-point crossover. The utility function is used as the fitness function and tournament selection is used for selecting the next generation.

An empirical study is made and compared to exhaustive search. The time needed for the GA is less than $1 \cdot 10^{-5}$ of the time needed for the exhaustive search. The GA also converges to the best solution very quickly, after only 40 generations. Thus, it would seem that using a search algorithm to this problem would produce extremely good results, at least in terms of time and speed. However, in this case all the components need to be known beforehand as the task is to choose an optimal set from alternative components. It would be interesting to see at least how all the different variables needed are acquired,

and how the approach could be more generalized.

Bodhuin et al. [2007] present an approach based on GAs and an environment that, based on previous usage information of an application, re-packages it with the objective of limiting amount of resources transmitted for using a set of application features. The overall idea is to cluster together (in jars) classes that, for a set of usage scenarios, are likely to be used together. Bodhuin et al. propose to cluster together classes according to dynamic information obtained from executing a series of usage scenarios. The approach aims at grouping in jars classes that are used together during the execution of a scenario, with the purpose of minimizing the overall jar downloading cost, in terms of time in seconds for downloading the application. After having collected execution trace, the approach determines a preliminary re-packaging considering common class usages and then improves it by using GAs. This approach can be seen as attempting to find optimal sub-architectures for a system, as each jar-package needs to be able to operate on its own. Obviously the success of finding sub-systems greatly depends on how well the class responsibility assignment problem is solved in the system, linking these results to that fundamental problem.

The proposed approach has four steps. First, the application to be analyzed is instrumented, and then it is exercised by executing several scenarios instantiated from use cases. Second, a preliminary solution of the problem is found, grouping together classes used by the same set of scenarios. Third, GAs are used to determine the (sub)-optimal set of jars. Fourth, based on the results of the previous steps, jars are created.

For the GA, an integer array is used as chromosome representation, where each gene represents a cluster of classes. The initial population is composed randomly. Mutation selects a cluster of classes and randomly changes its allocation to another jar archive. The crossover is the standard one-point crossover. The fitness function is $F(x) = 1/N * \sum(Cost_i)$ where N is the number of scenarios and $Cost$ is calculated from the call cost of making a request to the server and from the class sizes. 10% of the best individuals are kept alive across subsequent generations. Individuals to be reproduced are selected using a roulette-wheel selection. Scenarios are used in a very different way here as in the work of Rähkä et al. [2009]. Here scenarios define actions made with the system, and thus contain information of different components of the system that are needed, but do not deal with quality aspects other than how many operations, i.e., scenarios a certain set of responsibilities is able to perform. Rähkä et al. [2009], however, use scenarios not to describe functional operations but expectations to the system in terms of quality aspects. These different studies suggest that there are more ways of measuring quality than

metrics, and they should be more thoroughly investigated.

Results show that GA does improve the initial packaging, by 60-90 % to the actual initial packaging and by 5-43% compared to a packaging that contains two jars, “used” and “unused”, and by 13-23% compared to the preliminary best solution. When delay increases, the GA optimization starts to be highly more useful than the preliminary optimal solution, while the “used” packaging becomes better. However, for network delay value lower or slightly higher than the value used for the optimization process, the GA optimization is always the best packaging option. It is found that even when there is a large corpus of classes used in all scenarios, a cost reduction is still possible, even if in such a case the preliminary optimized solution is already a good one. The benefits of the proposed approach depend strongly on several factors, such as the amount of collected dynamic information, the number of scenarios subjected to analysis, the size of the common corpus and the networks delay. However, the presented approach and its results can be binded to several other software design related questions, thus raising questions on how the different promising results can be combined so that even more complex problems can be solved with search-based methods.

Gold et al. [2006] experiment with applying search techniques to integrate boundary overlapping concept assignment. Hill climbing and GA approaches are investigated. The fixed boundary Hypothesis Based Concept Assignment (HBCA) [Gold, 2001] technique is compared to the new algorithms. As program comprehension is extremely valuable when (re-)designing software architecture and locating (and understanding) overlapping concepts is one of the most demanding tasks in comprehension, automating this task would significantly save resources in program maintenance.

A concept may take the form of an action or object. For each concept found from source code, a hypothesis is generated and stored. The list of hypotheses is ordered according to the position of the indicators in the source code. The input for search problem is the hypothesis list. The hypothesis list is given by application of HBCA. The problem is defined as searching for segments of hypothesis in each hypothesis list according to predetermined fitness criteria such that each segment has the following attributes: each segment contains one or more neighboring hypotheses and there are no duplicate segments.

A chromosome is made up of a set of one or more segments representations, and its length can vary. A segment is encoded as a pair of values (locations) representing the start and end hypothesis of the hypothesis list. All segments with the same winning concept that overlap are compared and all but the fittest segment are removed from the

solution. Tournament selection is used for crossover and mutation. Mutation in GA randomly replaces any hypothesis location within any segment with any other valid hypothesis location with the concern for causing the search to become overly randomized. In HC the mutation generates new solutions by selecting a segment and increasing or decreasing one of the values by a single increment. Selecting different mutations for GA and HC is noteworthy: this choice is partially justified by the authors by the fact that mutation is only the secondary operation for the GA, and transformations are primarily done with the crossover. The chosen mutation operator for the GA seems to ensure diversity within the population. The proposed HC takes advantage of the crossover for GA for the restart mechanism, which recombines all segments to create new pairs of location values, which are then added to the current solution if their inclusion results in an improvement to the fitness value. Crossover utilizes the location of the segments, where only segments of overlapping locations are recombined and the remaining are copied to the new chromosome.

The fitness criteria's aims are finding segments of strongest evidence and binding as many of the hypotheses within the hypothesis list as possible without compromising the segment's strength of evidence. The segmentation strength is a combination of the inner fitness and the potential fitness of each segment. The inner fitness fit_i of a segment is defined as $signal_i - noise_i$, where $signal_i$ is the number of hypotheses within the segment that contribute to the winner, and $noise_i$ represents the number of hypotheses within the segment that do not contribute to the winner. In addition, each segment is evaluated with respect to the entire segment hypothesis list: the potential segment fitness, fit_p , is evaluated by taking account of $signal_p$, the number of hypotheses outside of the segment that could contribute to the segment's winning concept if they were included in the segment. The potential segment fitness is thus defined as $fit_p = signal_i - signal_p$. The overall segment fitness is defined as $segfit = fit_i + fit_p$. The total segment fitness is a sum of segment fitnesses. The fitness is normalized with respect to the length of the hypothesis list. The chosen fitness function seems quite simple when broken down to actual calculations. This further confirms the findings by, e.g., Lutz [2001] that simple approaches tend to have promising results, as there is less room to err.

An empirical study is used. Results are also compared to sets of randomly generated solutions for each hypothesis list, created according to the solutions structure. The results from GA, HC and random experiment are compared based on their fitness values. The GA fitness distribution is the same as those of HC and random, but achieves higher values. HC is clearly inferior. Comparing GA, HC and HBCA shows a lack of solutions

with low Signal to Noise ratios for GA and HC when compared to HBCA. GA is identified as the best of the proposed algorithms for concept assignment which allow overlapping concept boundaries. Also, the HC results are somewhat disappointing as they are found to be significantly worse than GA and random solutions. However, HC produces stronger results than HBCA on the signal to size measure. The GA and HC are found to consistently produce stronger concepts than HBCA. It might be worth studying how the HC would have performed if it used the same mutation operator as the GA. Although the GA primarily used the crossover, which was used as a basis for the HC, the GAs large population makes the application of this operator significantly more different than with HC.

Goldsby and Cheng [2008] and Goldsby et al. [2008] study the digital evolution of behavioral models for autonomic systems with Avida. It is difficult to predict the behavior of autonomic systems before deployment, and thus automatic generation of behavioral models greatly eases the task of software engineers attempting to comprehend the system. In digital evolution a population of self-replicating computer programs (digital organisms) exists in a computational environment and is subject to mutations and selection. In this approach each digital organism is considered as a generator for a UML state diagram describing the systems behavior.

Each organism is given instinctual knowledge of the system in the form of a UML class diagram representing the system structure, as well as optional seed state diagrams. A genome is thus seen as a set of instructions telling how the system should behave. The genome is also capable of replicating itself. In fact, in the beginning of each population there exists only one organism that only knows how to replicate itself, thus creating the rest of the population. Mutations include replacing an instruction, inserting an additional instruction and removing an instruction from the genome. As genomes are self-replicating, crossover is not used in order to create offspring. Here the choice of UML state diagrams is clever, as it visualizes the behavior in quite a simple manner, making the interpretation of the result easy. Also the choice of encoding conforms well to the chosen visualization method. However, the actual encoding of rules into the genome is not simple, and requires several different alphabets and lists of variables.

The fitness or quality of an organism is evaluated by a set of tasks, defined by the developer. Each task that the behavioral model is able to execute increases its merit. The higher a merit an organism has, the more it will replicate itself, eventually ending up dominating the population. This is yet another incident where the fitness is measured with something else than traditional metrics.

A behavioral model of an intelligent robot is used as a case study for Avida. Through a 100 runs of Avida, seven behavioral models are generated for the example system. Post-evolution analysis includes evaluation with the following criteria: minimum states, minimum transitions, fault tolerance, readability and tolerance. After the analysis, one of the models meets all but one criterion (safety) and three models meet three of the five criteria. One model does not meet any of the additional criteria. Thus, the produced behavioral models would seem to be of quality in average.

Lutz [2001] uses a measure based on an information theoretic minimum description length principle [Shannon, 1948] to compare hierarchical decompositions. This measure is furthermore used as the fitness function for the GA which explores the space of possible hierarchical decompositions of a system. Although this is very similar to software clustering, this approach is considered as architecture design as it does not need an initial clustering to improve, but designs the clustering purely based on the underlying system and its dependencies.

In hierarchical graphs links can represent such things as dependency relationships between the components of control-flow or data-flow. In order to consider the best way to hierarchically break a system up into components, one needs to know what makes a hierarchical modular decomposition (HMD) of a system better than another. Lutz takes the view that the best HMD of a system is the simplest. In practice this seems to give rise to HMDs in which modules are highly connected internally (high cohesion) and have relatively few connections which cross module boundaries (low coupling), and thus seems to achieve a principled trade-off between the coupling and cohesion heuristics without actually involving either. This also suggests that high quality architectures can effectively be identified through subjective inspection. A human architect may quite easily say if one design appears simpler than another, while calculation cohesion and coupling values is more time consuming and complex.

For the GA, the genome is a HMD for the underlying system. The chromosomes in the initial population are created by randomly mutating some number of times a particular “seed” individual. The initial seed individual is constructed by modularizing the initial system. Three different mutation operations are used that can all be thought of as operations on the module tree for the HMD. They are: 1. moving a randomly chosen node from where it is in the tree into another randomly chosen module of the tree, 2. modularize the nodes of some randomly chosen module, i.e., create a new module containing the basic entities of some module, and 3. remove a module “boundary”. The crossover operator resembles a tree-based crossover operation used in genetic

programming and is most easily considered as a concatenating operation on the module trees of the two HMDs involved. However, legal solutions are not guaranteed, and illegal ones are repaired.

The tree-like structure is significantly more complex than usual genome encodings for a GA. This is of course in line with the demands of the problem of finding an optimal HMD, but also reflects to the understandability of the chosen operations. The operations are difficult (if not impossible) to completely understand without visualization, and difficult corrective operations are needed in order to keep the system structure intact. The analogy between the chosen tree-operations and actual effects to the architecture is also quite difficult to grasp.

The fitness is given as $1/\text{complexity}$. Among other systems, a real software design is used for testing. A HMD with significantly lower complexity than the original was found very reliably, and the system could group the various components of the system into a HMD exhibiting a very logical (in terms of function) structure. These results validate that using simplicity as a fitness function is justified.

3.1.3 Summarizing Remarks

Search-based approaches to software architecture design is clearly a diverse field, as the studies presented solve very different issues relating to OO software architecture design and program comprehension. Some consensus can be found in the very basics: solving the class responsibility assignment problem, applying design choices to create an architecture and finding an optimal modularization (Lutz [2001] creates a modularization, Kim and Park [2009] attempt to find an optimal set of components and Bodhuin et al. [2007] attempt to find optimal sub-architectures). However, even within these sub-areas of OO design, the approaches are quite different, and practically no agreement can be found when studying the chosen encodings, operations or fitness function. What is noticeable, however, is that several approaches to quite different problems within this area use a fitness function that is not based on metrics. This highlights the need for better validation of using metrics in evaluating the quality of software, and especially software architectures. Many metrics need source code and very detailed information; this alone suggests that they are not suitable for this higher level problem.

Table 1. Studies in search-based object-oriented software architecture design

| Author | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|--|--|---|--|---|---|--|--|---|
| Bowman et al. [2008] | Class structure design is (semi-) automated | Class diagram as methods, attributes and associations | Integer vector and a dependency matrix | Randomly change the class of method or attribute | Standard one-point | Cohesion and coupling | Optimal class structure | Comparison between different algorithms |
| Simons and Parmee [2007a; 2007b; 2008] | Class structure design is automated | Use cases; data assigned to attributes and actions to methods | A design solution where attributes and methods are assigned to classes | An attribute and a method are moved from one class to another | Attributes and methods of parents are swapped according to class position | Cohesiveness of methods (COM) | Basic class structure for system. | Design solutions encoded directly into a programming language |
| Amoui et al. [2006] | Applying design patterns; high level architecture design | Software system | Chromosome is a collection of supergenes, containing information of pattern transformations | Implementing design patterns | Single-point crossovers for both supergene level and chromosome level, with corrective function | Distance from main sequence | Transformed system, design patterns used as transformations to improve modifiability | New concept of supergene used |
| Räihä et al. [2008a] | Automating architecture design | Responsibility dependency graph | Chromosome is a collection of supergenes, containing information of responsibilities and design patterns | Mutations apply architectural design patterns and styles | A standard one-point crossover with corrective function | Efficiency, modifiability and complexity | UML class diagram depicting the software architecture | |

| Author | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|--------------------------|---|--|--|--|---|---|--|--|
| Räihä et al. [2008b] | Automating CIM-to-PIM model transformations | Responsibility dependency graph and domain model (CIM model) | Chromosome is a collection of supergenes, containing information of responsibilities and design patterns | Mutations apply architectural design patterns and styles | A standard one-point crossover with corrective function | Efficiency, modifiability and complexity | UML class diagram depicting the software architecture (PIM model) | |
| Räihä et al. [2009] | Automating architecture design | Responsibility dependency graph and domain model | Chromosome is a collection of supergenes, containing information of responsibilities and design patterns | Mutations apply architectural design patterns and styles | A standard one-point crossover with corrective function | Efficiency, modifiability, complexity and modifiability related scenarios | UML class diagram depicting the software architecture | |
| Kessentini et al. [2008] | Using PSO for model transformations | Source model, target model and mapping blocks | Integer vector | N/A | N/A | Number of source model constructs that can be transformed | Optimal transformations | Particle Swarm Optimization (PSO) used as search algorithm |
| Kim and Park [2009] | Dynamic selection of software components | Softgoal interdependency graph, decision variables | String of integers representing decision variables | Goes through each gene and changes the digit according to mutation probability | Two-point crossover | Quality attributes given by user | Optimal architectural instance from the set of all instances | |
| Bodhuin et al. [2007] | Automating class clustering in jar archives | A grouping of classes of a system | An integer array, each gene is a cluster of classes allocated to the jar represented by integer | Changes the allocation of a class cluster to another jar archive | Standard one-point | Download cost of jar archive | Optimal packaging; finding the subsets of classes most likely to be used together (to be placed in same jar archive) | |

| Author | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|---|---|---|--|--|---|--|---|--|
| Gold et al. [2006] | Using GA in the area of concepts | Hypothesis list for concepts | One or more segment representations | A hypothesis location is randomly replaced within a segment pair | Segment pairs of overlapping locations are combined, rest copied | Strongest evidence for segments and hypothesis binding | Optimized concept assignment | Hill climbing used as well as GA |
| Goldsby and Chang [2008]; Goldsby et al. [2008] | Designing a system from a behavioral point of view | A class diagram, optional state diagram | A set of behavioral instructions | Changes, removes or adds an instruction | Self-replication | Number of executed tasks | UML state diagram giving the behavioral model of system | No actual evolutionary algorithm used, but a platform that is “an instance of evolution” |
| Lutz [2001] | Information theory applied in software design; high-level architecture design | Software system | Hierarchical modular decomposition (HMD) | Three mutations operating the module tree for the HMD | A variant of tree-based crossovers, as used in GP, with corrective function | 1/complexity | Optimal hierarchical decomposition of system | |

3.2 Service-oriented architecture design

3.2.1. Basics

Web services are rapidly changing the landscape of software engineering, and service-oriented architectures (SOA) are especially popular in business. One of the most interesting challenges introduced by web services is represented by Quality of Service (QoS)-aware composition and late-binding. This allows binding, at run-time, a service-oriented system with a set of services that, among those providing the required features, meet some non-functional constraints, and optimize criteria such as the overall cost or response time. Hence, QoS-aware composition can be modeled as an optimization problem. This problem is NP-hard, which makes it suitable for meta-heuristic search algorithms. For more background on SOA, see, e.g., Huhns and Sting [2005]. The following subsection describes several approaches that have used a GA to deal with optimizing service compositions. Summarizing remarks on the different approaches are given in the end, and the fundamentals of each approach are collected in Table 2.

3.2.2. Approaches

Canfora et al. [2005a] propose a GA to optimize service compositions. The approach attempts to quickly determine a set of concrete services to be bound to the abstract services composing the workflow of a composite service. Such a set needs both to meet QoS constraints, established in the Service Level Agreement (SLA), and to optimize a function of some other QoS parameters.

A composite service S is considered as a set of n abstract services $\{s_1, s_2, \dots, s_n\}$, whose structure is defined through some workflow description language. Each component s_j can be bound to one of the m concrete services, which are functionally equivalent. Computing the QoS of a composite service is made by combining calculations for quality attributes *time*, *cost*, *availability*, *reliability* and *customer attraction*. Calculations take into account Switch, Sequence, Flow and Loop patterns in the workflow.

The genome is encoded as an integer array whose number of items equals to the number of distinct abstract services composing the services. Each item, in turn, contains an index to the array of the concrete services matching that abstract service. The mutation operator randomly replaces an abstract service with another one among those available, while the crossover operator is the standard two-point crossover. This can be seen as an

attempt to preserve building blocks, i.e., sequences of optimal service bindings. Abstract services for which only one concrete service is available are taken out from the GA evolution.

The fitness function needs to maximize some QoS attributes, while minimizing others. In addition, the fitness function must penalize individuals that do not meet the constraints and drive the evolution towards constraint satisfaction, the distance from which is denoted by D . The fitness function is $f = (w_1\text{Cost} + w_2\text{Time}) / (w_3\text{Availability} + w_4\text{Reliability}) + w_5D$. QoS attributes are normalized in the interval $[0, 1)$. Although the fitness function seems simple in this way, the actual calculations behind the different attributes are complex. The values are achieved by calculating the quality value for each attribute for each pattern in the workflow. The actual functions to define how these values are calculated are not defined, and it would be interesting to see how, e.g., availability is achieved, as this would show the amount of information needed as input to calculate the fitness value. The weights w_1, \dots, w_5 are positive reals. Normalizing the fitness evaluators ensures that the weights have the true effect to the fitness value that they are meant to have.

A dynamic penalty is experimented with, so that w_5 is increased over the generations. An elitist GA is used where the best two individuals are kept alive across generations. Roulette wheel method is used for selection.

The GA is able to find solutions that meet the constraints, and optimizes different parameters (here cost and time). Results show that the dynamic fitness does not outperform the static fitness. Even different calibrations of weights do not help. The convergence times of GA and Integer Programming (IP) [Garfinkel and Nemhauser, 1972] are compared for the (almost) same achieved solution. The results show that when the number of concrete services is small, IP outperforms GA. For about 17 concrete services, the performance is about the same. After that, GA clearly outperforms IP. Thus, as SOA is most useful when the amount of services is large, it would seem that GA is a worthwhile solution to optimizing the service-binding.

Canfora et al. [2005b] have continued their work by using a GA in replanning the binding between a composite service and its invoked services during execution. Replanning is triggered once it can be predicted that the actual service QoS will differ from initial estimates. After this, the slice, i.e., the part of workflow still remaining to be executed, is determined and replanned. The used GA approach is the same as earlier, but additional algorithms are used to trigger replanning and computing workflow slices. The GA is used to calculate the initial QoS-values as well as optimizing the replanned slices.

Experiments were made with realistic examples and results concentrate on the cost quality factor. The algorithms managed to reduce the final cost from the initial estimate, while response time increased in all cases. The authors end with a note that the trade-off between response time and cost quality factors need to be examined thoroughly in the future.

Jaeger and Mühl [2007] discuss the optimization problem when selecting services while considering different QoS characteristics. A GA is implemented and tested on a simulation environment in order to compare its performance with other approaches.

An individual in the implemented GA represents an assignment of a candidate for each task and can be represented by a tuple. A population represents a set of task-candidate assignments. The initial population is generated arbitrarily from possible combinations of tasks and candidates. Mutation changes a particular task-candidate assignment of an individual. Crossover is made by combining two particular task-candidate assignments to form new ones and depends on the fitness value. The fitness value is computed based on the QoS resulting from the encoded task-services assignment. Jaeger and Mühl use the same fitness function as Canfora et al. [2005a; 2005b] in order to get comparable results.

A trade-off couple between execution time and cost is defined as follows: the percentage a , added to the optimal execution time, is taken to calculate the percentage b , added to the optimal cost, with $a + b = 100$. Thus, the shorter the execution time is, the worse will be the cost and vice versa. The constraint is determined to perform the constraint selection on the execution time first. The aggregated cost for the composition is increased by 20% and then taken as the constraint that has to be met by the selection. This appears as an attempt to answer the problem noted by Canfora et al. [2005b] in their later study.

Several variations of the fitness function are possible. Jaeger and Mühl use a multiplication of the fitness to make the difference between weak and strong fitnesses larger. When the multiplying factor is 4, it achieves higher QoS values than those with a smaller factor; however, a factor of 8 does not achieve values as high. The scaled algorithm performed slightly better than the one with a factor of 2, and behaved similarly to the weighted algorithm. The penalty factor was also investigated, and it was varied between 0.01 and 0.99 in steps of 0.01. The results show that a factor of 0.5 would result in few cases where the algorithm does not find a constraint meeting solution. On the other hand, solutions below 0.1 appear too strong, as they represent an unnecessary restriction of the GA to evolve further invalid solutions. These different experiments on some very

basic parameters demonstrate the difficulty of optimizing the GA: even the more simple choices are anything but straightforward.

The GA offers a good performance at feasible computational efforts when compared to, e.g., bottom-up heuristics. However, this approach shows a large gap when compared to the resulting optimization of a branch-and-bound approach or to exhaustive search. It appears that the considered setup of values along with the given optimization goals and constraints prevent a GA from efficiently identifying very near optimal solutions.

Zhang et al. [2006] implement a GA that, by running only once, can construct the composite service plan according to the QoS requirements from many services compositions. This GA includes a special relation matrix coding scheme (RMCS) of chromosomes proposed on the basis of the characters of web services selection.

By means of the particular definition, it can simultaneously represent all paths of services selection. Furthermore, the selected coding scheme can simultaneously denote many web service scenarios that the one dimension coding scheme can not express at one time.

According to the characteristic of the services composition, the RMCS is adopted using a neighboring matrix. In the matrix, n is the number of all tasks included in services composition. The elements along the main diagonal for the matrix express all the abstract service nodes one by one and are arranged from the node with the smallest code number to the node with the largest code number. The objects of the evolution operators are all elements along the main diagonal of the matrix. The chromosome is made up of these elements. The other elements in the matrix are to be used to check whether the created new chromosomes by the crossover and mutation operators are available and to calculate the QoS values of chromosomes. This appears to mainly combine the integer array and the table of services linked to it, used by Canfora et al. [2005a], into one data structure. The tuple representation chosen by Jaeger and Mühl [2007] does not seem that different either, as a tuple can basically contain the information of what is represented by a column and a row in a matrix.

The policy for initial population attempts to confirm the proportion of chromosomes for every path to the size of the population. The method is to calculate the proportion of compositions of every path to the sum of all compositions of all paths. The more there are compositions of one path, the more chromosomes for the path are in the population.

The value of every task in every chromosome is confirmed according to a local optimized method. The larger the value of QoS of a concrete service is, the larger the probability to be selected for the task is. The roulette wheel selection is used to select

concrete services for every task.

The probability of mutation is for the chromosome instead of the locus. If mutation occurs, the object path will be confirmed firstly whether it is the same as the current path expressed by the current chromosome. If the paths are different, the object path will be selected from all available paths except the current one. If the object is itself, the new chromosome will be checked whether it is the same as the old chromosome. Same chromosome will result in the mutation operation again. If the objects are different paths from the current path, a new chromosome will be related on the basis of the object path.

A check operation is used after the invocations of crossover and mutation. If the values of the crossover loci in two crossover chromosomes are all for the selected web services, the new chromosomes are valid. Else, the new chromosomes need to be checked on the basis of the relation matrix. Mutation checks are needed if changed from selected web service to a certain value or vice versa.

Zhang et al. compared the GA with RMCS to a standard GA with the same data, including workflows of different sizes. The used fitness function is as defined by Canfora et al. [2004]. The coding scheme, the initial population policy and the mutation policy are the differences between the two GAs. Results show that the novel GA outperforms the standard one in terms of achieved fitness values. As the number of tasks grows, so does the difference between fitness values (and performance time, in the favor of the standard solution) between the two GAs. The weaknesses of this approach are thus long running time and slow convergence. Tests on the initial population and the mutation policies show that as the number of tasks grows, the GA with RMCS outperforms the standard one more clearly. Thus it would seem that combining the information into a heavier data structure, a matrix, increases execution time significantly. Also, as noted that the improvement fitness values with the novel GA for larger task sets is achieved by testing other improvement than the encoding, the true achievements are the ones that really differ from previous approaches, rather than the new representation. Tests on the coding scheme show that the novel matrix approach only achieves noticeably better fitness values when the number of tasks is increased (although the improvement is not linear): the fitness values for 10 tasks only differ by less than 1 %, the fitness values for 25 tasks differ by approximately 30%, and the fitness values for 30 tasks by approximately 20%. Another interesting point is the choice of parameters: Zhang et al. use 10 000 generations and 400 individuals for a population in their tests. However, the standard GA seems to achieve its optimum after 1000 generations and the one with the novel encoding after 3000 generations. Thus one wonders the need for such unusual parameter selections.

Zhang et al. report that experiments on QoS-aware web services selection show that the GA with the presented matrix approach can get a significantly better composite service plan than the GA with the one dimension coding scheme, and that the QoS policies play an important role in the improvement of the fitness of GA.

Su et al. [2007] continue the work of Zhang et al. [2006] by proposing improvements for the fitness function and mutation policy. An objective fitness function 1 (OF1) is first defined as a sum of quality factors and weights, providing the user with a way to show favoritism between quality factors. The sum of positive quality factors is divided by the sum of negative quality factors. The second fitness function (OF2) is a proportional one and takes into account the different ranges of quality value. The third fitness function (OF3) combines OF1 and OF2, producing a proportional fitness function that also expresses the differences between negative and positive quality factors. Thus Su et al. seem to have noticed the problems with defining the fitness functions, as the fitness function actually used by Canfora et al. [2005a; 2005b] includes similar improvements.

Four different mutation policies are also inspected. Mutation policy 1 (MP1) operates so that the probability of the mutation is tied to each locus of a chromosome. Mutation policy 2 (MP2) has the mutation probability tied to the chromosomes. Mutation policy 3 (MP3) has the same principle as MP1, except that now the child may be identical to the parent. Mutation policy 4 (MP4) has the probability tied to each locus, and has an equal selection probability for each concrete service and the “0” service.

Experiments with the different fitness functions suggest that OF3 clearly outperforms OF1 and OF2 in terms of the reached average maximum fitness value. This is quite unsurprising, as OF3 is the most developed fitness function. Experiments on the different mutation policies show that MP1 gains the largest fitness values while MP4 performs the worst.

Cao et al. [2005a; 2005b] present a GA that is utilized to optimize a business process composed of many service agents (SAG). Each SAG corresponds to a collection of available web services provided by multiple-service providers to perform a specific function. Service selection is an optimization process taking into account the relationships among the services. Better performance is achieved using GA compared to using local service selection strategy.

A service selection model using GA is proposed to optimize a business process composed of many service agents. A SAG corresponds to a collection of available web services provided by multiple service providers to perform a specific function. When only measuring cost, the service selection is equivalent to a single-objective optimization

problem.

An individual is generated for the initial population by randomly selecting a web service for each SAg of the services flow, and the newly generated individual is immediately checked whether the corresponding solution satisfies the constraints. If any of the constraints is violated, then the generated individual is regarded as invalid and discarded. The roulette wheel selection is used for individuals to breed.

Mutation bounds the selected SAg to a different web service than the original one. After an offspring is mutated, it is also immediately checked whether the corresponding solution is valid. If any constraints are violated, then the mutated offspring is discarded and the mutation operation is retried.

A traditional single-point crossover operator is used to produce two new offspring. After each crossover operation, the offspring are immediately checked whether the corresponding solutions are valid. If any of the constraints is violated, then both offspring are discarded and the crossover operation for the mated parents is retried. If valid offspring still cannot be obtained after a certain number of retries, the crossover operation for these two parents is given up to avoid a possible infinite loop.

Cao et al. take cost as the primary concern of many business processes. The overall cost of each execution path can always be represented by the summation cost of its subset components. For GA, integer encoding is used. The solution to service selection is encoded into a vector of integers. The fitness function is defined as $f = U - \sum(\text{costs of service flows})$, if $\text{cost} < U$, and otherwise 0. The constant U should be selected as an appropriate positive number to ensure the fitness of all good individuals get a positive fitness value in the feasible solution space. On the other hand, U can also be utilized to adjust the selection pressure of GA. This is a clever approach to give the developer a simple way to adjust the selection process and appreciation of different solutions.

In the case study the best fitness of the population has a rapid increase at the beginning of the evolution process and then convergences slowly. It means the overall cost of the SAg is generally decreasing with the evolution process. For better solutions, the whole optimization process can be repeated for a number of times, and the best one in all final solutions is selected as the ultimate solution to the service selection problem.

3.2.3 Summarizing Remarks

Contrary to the studies relating to OO architecture design, the approaches to apply search algorithms in SOA design are extremely similar. Nearly all studies use the same fitness functions or they have made only small modifications to it. Also the basic representation of the problem is very similar; although different definitions are used, the underlying

problem is always linking concrete services with abstract services. Improvements have been attempted by creating different initial population and mutation policies; note, that the actual mutation is still the same, but the way the mutation is applied is changed. Additionally, there is no consensus in the encoding of the solution, although the problem is the same, and some tests have been made to compare different encoding options. Thus the main questions in this area seem to be: are there other problems in SOA where search algorithms could be applied to, and can a truly optimal encoding be found to the currently studied problem? Additionally, the fitness function deserves much more attention and testing, as the developers of the fitness function used by all the studies say themselves that the relationships and trade-offs between different quality attributes need to be carefully studied. Results with dynamic fitness functions also interestingly did not increase the fitness value. Rähkä et al. [2008a; 2008b] experimented with dynamic mutations, but discarded them in their latest study [Rähkä et al., 2009]. This would suggest that using dynamicity with GAs is a complex problem, demanding well-defined operations and firm justifications for the use of such improvements before adding them to the experiments.

Table 2. Studies in search-based service-oriented software architecture design

| Author | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|------------------------|--|---|---|--|--------------------------------------|---|---|--|
| Canfora et al. [2005a] | Service composition with respect to QoS attributes | Sets of abstract and concrete services | Integer array, whose size is the number of abstract services, each item contains an index to array of concrete services | Randomly replaces an abstract service with another | Standard two-point crossover | Minimize cost and time, maximize availability and reliability, meet constraints, with penalty | Optimized service composition meeting constraints, concrete services bound to abstract services | A dynamic penalty was experimented with |
| Canfora et al. [2005b] | Replanning during execution time | Sets of abstract and concrete services | Integer array, whose size is the number of abstract services, each item contains an index to array of concrete services | Randomly replaces an abstract service with another | Standard two-point crossover | Minimize cost and time, maximize availability and reliability, meet constraints | Optimized service composition meeting constraints, concrete services bound to abstract services | GA used to calculate initial QoS-value and QoS-values inbetween: replanning is triggered by other algorithms |
| Jaeger and Mühl [2007] | Service assignment with respect to QoS attributes | Selection of services and tasks to be carried out | A tuple representing an assignment of a candidate for a task | Changes an individual task-candidate assignment | Combining task-candidate assignments | Minimize cost and time, maximize availability and reliability, meet constraints, with penalty | Tasks assigned to services considering QoS attributes | A trade-off couple between execution time and cost is defined |

| Author | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|---------------------------|---|---|--|--|--|---|---|--|
| Zhang et al. [2006] | Task assignment with relation to QoS attributes | Selections of tasks and services | Relation matrix coding scheme | Standard, with corrective function | Standard, with corrective function | Minimize cost and time, maximize availability and reliability, meet constraints | Tasks assigned to services considering QoS attributes | Initial population and mutation policies defined |
| Su et al. [2007] | Task assignment with relation to QoS attributes | Selections of tasks and services | Relation matrix coding scheme | Standard, with corrective function | Standard, with corrective function | Minimize cost and time, maximize availability and reliability, meet constraints | Tasks assigned to services considering QoS attributes | Initial population and mutation policies defined |
| Cao et al. [2005a; 2005b] | Business process optimization | Collections of web services and service agents (SAG) composing a business process | Integer encoding, assigning a SAg to a service | Changes the service to which a SAg is bound with corrective function | Standard one-point, producing two new offspring with corrective function | Cost | Services assigned to service agents | |

3.3. Other

3.3.1 Background

In addition to purely designing software architecture, there are some factors that should be optimized, regardless of the particularities of an architecture. Firstly, there is the reliability-cost tradeoff. The reliability of software is always dependent on its architecture, and the different components should be as reliable as possible. However, the more work is put to ensure reliability of different components, the more the software will cost. Wadekar and Gokhale [1999] implement a GA to optimize the reliability-cost tradeoff. Secondly, there are some parameters, e.g., tile sizes in loop tiling and loop unrolling, which can be optimized for all software architectures in order to optimize the performance of the software. Che et al. [2000] apply search-based techniques for such parameter optimization.

3.3.2 Approaches

Wadekar and Gokhale [1999] present an optimization framework founded on architecture-based analysis techniques, and describe how the framework can be used to evaluate cost and reliability tradeoffs using a GA. The methodology for the reliability analysis of a terminating application is based on its architecture. The architecture is described using the one-step transition probability matrix P of a discrete time Markov chain (DTMC).

Wadekar and Gokhale assume that the reliabilities of the individual modules are known, with R_i denoting the reliability of module i . It is also assumed that the cost of the software consisting of n components, denoted by C , can be given by a generic expression of the form: $C = C_1(R_1) + C_2(R_2) + \dots + C_n(R_n)$ where C_i is the cost of component i and the cost C_i depends monotonically on the reliability R_i . Thus, the problem of minimizing the software cost while achieving the desired reliability is the problem of selecting module reliabilities.

A chromosome is a list of module reliabilities. Each member in the list, a gene, corresponds to a module in the software. The independent value in each gene is the reliability of the module it represents, and the dependent value is the module cost given by the module cost-reliability relation or a table known *a priori*. The gene values are changed to alter the cost and reliability of a software implementation represented by a particular chromosome.

Mutation and crossover operations are standard. To avoid convergence to a local

optimum as the population size increases, the mutation operation is used more frequently. A cumulative-probability based basic selection mechanism is used for selection. Chromosomes are ranked by fitness and divided into rank groups. The probability of selection of chromosomes varies uniformly according to their rank group where chromosomes in the first rank group have the largest probability. A new generation of the population is created by selecting $p_{\text{imax}}/2$ chromosomes, where p_{imax} is maximum population. If the cost reduction is less than or equal to $\phi\%$ of the current best cost τ number of times, the GA terminates. During any generation cycle if the cost reduction is larger, the counter τ is reset to 0. The reduction percentage factor ϕ and the counter limit τ are parameters. This approach is one of the few alternatives used to terminate a GA, as most studies presented use a straightforward generation number to terminate the execution of the algorithm.

The fitness function is $f = (-K/\ln R)/C^\gamma$, where K is a large positive constant. The fitness of solutions increases superlinearly with their reliability. The constant γ is used to linearize the cost variation. The maximum fitness is directly proportional to K . An intermediate value of gamma, $\gamma = 1.5$, allows the GA to distinguish between low-cost and high-cost solutions, while selecting a sufficient number of high-cost high-reliability solutions, that may generate the optimal high-reliability low-cost solution.

Wadekar and Gokhale compare the GA against exhaustive search. The results indicate that the GA consistently and efficiently provides optimal or very close to optimal designs, even though the percentage of such designs in the overall feasible design space is extremely small. The results also highlight the robustness of the GA. However, the small number of near-optimal solutions demonstrates that the fitness landscape is very complex, again conforming to the need to extensively investigate the cost-reliability trade-off. The case study results show how the GA can be effectively used to select components such that the software cost is minimized, for various cost structures.

Che et al. [2003] present a framework for performance optimization parameter selection, where the problem is transformed into a combinatorial minimization problem. Many performance optimization methods depend on right optimization parameters to get good performance for an application. Che et al. search for the near optimal optimization parameters in a manner that is adaptable for different architectures. First a reduction transformation is performed to reduce the program's runtime while maintaining its relative performance as regard to different parameter vectors. The near-optimal optimization parameter vector based on the reduced program's real execution time is searched by GA, which converges to a near-optimal solution quickly. The reduction

transformation reduces the time to evaluate the quality of each parameter vector.

First some transformations are applied to the application, leaving the optimization parameter vector to be read from a configuration file. Second, the application is compiled into executable with the native compiler. Then the framework repeatedly generates the configure file with a different parameter vector selected by search and measures the executable's runtime.

The chromosome encoding for the GA is a vector of integer values, with each integer corresponding to an optimization parameter of a solution. No illegal solutions are allowed. The population has a fixed size. A simple integer value mutation is implemented and an integer number recombination scheme is used for crossover. The fitness value reflects the duality of an individual in relation to other individuals. The linear rank-based fitness assignment scheme is used to calculate the fitness values. Selection for a new generation is made by elitism and roulette wheel method. Test results show that the GA can adapt to different execution environments automatically. For each platform, it always selects excellent optimization parameters for 80% programs. Results show that the number of individuals evaluated is far smaller than the size of solution space for each program on each platform. The optimization time is also small.

4. SOFTWARE CLUSTERING

4.1 Basics

As software systems develop and are maintained, they tend to grow in size and complexity. A particular problem is the growing number of dependencies between libraries, modules and components within the modules. Software clustering (or modularization) attempts to optimize the clustering of components into modules in such a way that there are as many dependencies within a module as possible and as few dependencies between modules as possible. This will enhance the understandability of a system, which in turn will make it more maintainable and modifiable. Also, fewer dependencies between modules usually results in better efficiency.

As components or modules (depending on the level of detail in the chosen representation) can be depicted as vertices and dependencies between them as edges in a graph, the software clustering problem can be traced back to a graph partitioning problem, which is NP-complete. Genetic algorithms have successfully been applied to a general graph partitioning problem [Bui and Moon, 1996; Shazely et al., 1998], and thus, the related software clustering problem is most suitable for meta-heuristic search techniques.

Although the basic problem is relatively simple to define and the goodness of a modularization can be calculated based on the goodness of the underlying graph partitioning, the nature of software systems provides challenges when defining the actual fitness function for the optimization algorithm. Also, not all necessary information can be encoded into a simple graph representation, and this presents another question to be answered when designing a search-based approach for modularization. The following subsection presents approaches using GAs, HC and SA to find good software modularizations, after which summarizing remarks are presented and the fundamentals of each study are collected in Table 3.

4.2. Approaches

Mancoridis et al. [1998] treat automatic modularization as an optimization problem and have created the Bunch tool that uses HC and GA to aid its clustering algorithms. A hierarchical view of the system organization is created based solely on the components and relationships that exist in the source code. The first step is to represent the system modules and the module-level relationships as a module dependency graph (MDG). An algorithm is then used to partition the graph in a way that derives the high-level subsystem structure from the component-level relationships that are extracted from the source code. The goal of this software modularization process is to automatically partition the components of a system into clusters (subsystems) so that the resultant organization concurrently minimizes inter-connectivity while maximizing intra-connectivity. This task is accomplished by treating clustering as an optimization problem where the goal is to maximize an objective function based on a formal characterization of the trade-off between inter- and intra-connectivity. Intuitively, intra-connectivity could be seen as cohesion and inter-connectivity as coupling.

The clusters, once discovered, represent higher-level component abstractions of a system's organization. Each subsystem contains a collection of modules that either cooperate to perform some high-level function in the overall system or provide a set of related services that are used throughout the system. Intra-connectivity A_i of cluster i consisting of N_i components and m_i intra-edge dependencies as $A_i = m_i/N_i^2$, bound between 0 and 1. Interconnectivity measures the connectivity between two distinct clusters. A high degree of inter-connectivity is an indication of poor subsystem partitioning. Inter-connectivity E_{ij} between clusters i and j consisting of N_i and N_j components with e_{ij} inter-edge dependencies is 0, if $i = j$, and $e_{ij} / 2*N_i*N_j$ otherwise, bound between 0 and 1. Modularization Quality (MQ) demonstrates the trade-off

between inter- and intra-connectivities, and it is defined for a module dependency graph

$$\text{partitioned into } k \text{ clusters as } 1/k * \sum \frac{A_i - 1}{k * \frac{k-1}{2}} * \sum E_{i,j} \text{ if } k > 1, \text{ or } A_1, \text{ if } k = 1.$$

The first step in automatic modularization is to parse the source code and build a MDG. A sub-optimal clustering algorithm works as the traditional hill climbing one by randomly selecting a better neighbor. The GA starts with a population of randomly generated initial partitions and systematically improving them until all of the initial samples converge. The GA uses the “neighboring partition” definition to improve an individual, and thus only contains one mutation operator, which is the same one as used with HC. Selection is done by randomly selecting a percentage of N partitions and improving each one by finding a better neighboring partition. A new population is generated by making N selections, with replacements for the existing population of N partitions. Selections are random and biased in favor of partitions with larger MQs. The algorithm continues until no improvement is seen for t generations, or until all of the partitions in the population have converged to their maximum MQ, or until the maximum number of generations has been reached. The partition with the largest MQ in the last population is the sub-optimal solution.

Experimentation with this clustering technique has shown good results for many of the systems that have been investigated. The primary method used to evaluate the results is to present an automatically generated modularization of a software system to the actual system designer and ask for feedback on the quality of the results. A case study was made and the results were shown to an expert, who highly appreciated the result produced by Bunch.

The validation of the method is interesting, as the original designer of a system should be the one who knows the system best, and thus should be the best one to evaluate designs of the system. It is also encouraging that the designers were open and admitted that the tool was able to improve the design that they must have thought of as optimal at some point. This indicates that there truly is a place for software design tools if the methods are well-defined enough.

Doval et al. [1999] have implemented a more refined GA in the Bunch tool, as it now contains a crossover operator and more defined mutation and crossover rates. The effectiveness of the technique is demonstrated by applying it to a medium-sized software system. For encoding, each node in the graph (MDG) has a unique numerical identifier assigned to it. These unique identifiers define which position in the encoded string will be

used to define that node's cluster. Mutation and crossover operators are standard. A roulette wheel selection is used for the GA, complemented with elitism. Fitness function is based on the MQ metric. Crossover rate was 80% for populations of 100 individuals or fewer and 100% for populations of a thousand individuals or more, varying linearly between those values. Mutation rate is $0.004 \log_2(N)$. The MQ values for constant population and generation values were smaller, but fairly close, within 10% to values achieved with final values for population and generation.

The affect of the population size to crossover rate is interesting, especially in the sense that with smaller populations the rate is smaller. Intuitively it would seem that with larger populations there would be a higher chance that the population contains some extremely poor individuals, the parts of which are not worthwhile to pass on to future generations.

Mancoridis et al. [1999] have continued to develop the Bunch tool for optimizing modularization. Firstly, almost every system has a few modules that do not seem to belong to any particular subsystem, but rather, to several subsystems. These modules are called omnipresent, because they either use or are used by a large number of modules in the system. In the improved version users are allowed to specify two lists of omnipresent modules, one for clients and another for suppliers. The omnipresent clients and suppliers are assigned to two separate subsystems.

Secondly, experienced developers tend to have good intuition about which modules belong to which subsystems. However, Bunch might produce results that conflict with this intuition for several reasons. This is addressed with a user-directed clustering feature, which enables users to cluster some modules manually, using their knowledge of the system design while taking advantage of the automatic clustering capabilities of Bunch to organize the remaining modules. Both user-directed clustering and the manual placement of omnipresent modules into subsystems have the advantageous side-effect of reducing the search space of MDG partitions. By enabling the manual placement of modules into subsystems, these techniques decrease the number of nodes in the MDG for the purposes of the optimization and, as a result, speed up the clustering process.

Finally, once a system organization is obtained, it is desirable to preserve as much of it as possible during the evolution of the system. The integration of the orphan adoption technique into Bunch enables designers to preserve the subsystem structure when orphan modules are introduced. An orphan module is either a new module that is being integrated into the system, or a module that has undergone structural changes. Bunch moves orphan modules into existing subsystems, one at a time, and records the MQ for

each of the relocations. The subsystem that produces the highest MQ is selected as the parent for the module. This process, which is linear with respect to the number of clusters in the partition, is repeated for each orphan module. Results from a case study support the added features.

The chosen additions clearly stem from real needs when modularizing software. However, two of the three operations increase the power that the user has over Bunch, thus decreasing the level of automation. Ideally the tool would be able to locate the omnipresent modules themselves, and gain the same level of expertise via a fitness function as experts, so that the user would not need to cluster anything beforehand. The last improvement, however, is truly beneficial, as hardly any software system stays intact during maintenance, and modules need to be added or modified. Automating the step of finding the optimal place for a new module is a big step towards the ideal of automating software design.

Mitchell and Mancoridis [2002; 2006; 2008] have continued to work with the Bunch tool and have further developed the MQ metric. They define MQ as the sum of Clustering Factors for each cluster of the partitioned MDG. The Clustering Factor (CF) for a cluster is defined as a normalized ratio between the total weight of the internal edges and half of the total weight of external edges. The weight of the external edges is split in half in order to apply an equal penalty to both clusters that are connected by an external edge. If edge weights are not provided by the MDG, it is assumed that each edge has a weight of 1. The clustering factor is defined as

$$CF = \text{intra-edges} / (\text{intra-edges} + \frac{1}{2} * \sum (\text{inter-edges})).$$

The measurement is adjusted, as Mitchell and Mancoridis argue that the old MQ tended to minimize the inter-edges that exited the clusters, and not minimize the number of inter-edges in general. The representation also supports weights. This is an interesting observation, as the original definition of the MQ metric makes no distinction to whether an edge exits a cluster or not. Thus, one could ask whether the MQ metric was the sole reason for the previous results, or if other improvements besides the newly defined MQ metric also had a significant effect on obtaining the better quality results. The addition of weights is also noteworthy, as previously the problem was not considered a multi-objective one, while the addition of weights clearly indicates so.

The HC algorithm for the Bunch tool has also been enhanced. During each iteration, several options are now available for controlling the behavior of the hill-climbing algorithm. First, the neighboring process may use the first partition that it discovers with a larger MQ as the basis for the next iteration. Second, the neighboring process examines

all neighboring partitions and selects the partition with the largest MQ as the basis for the next iteration. Third, the neighboring process ensures that it examines a minimum number of neighboring partitions during each iteration. For this, a threshold n is used to calculate the minimum number of neighbors that must be considered during each iteration of the process. Experience has shown that examining many neighbors during each iteration, so that $n > 75\%$, increases the time the algorithm needs to converge to a solution. This is quite intuitive, as each examination increases the run time of the algorithm, and it is not likely that simply by examining several neighbors the algorithm would suddenly find a steeper climb (i.e., converge faster).

It is observed that as n increases so does the overall runtime and the number of MQ evaluations. However, altering n does not appear to have an observable impact on the overall quality of the clustering results. A simulated annealing algorithm is also made for comparison. Although the simulated annealing implementation does not improve the MQ, it does appear to help reduce the total runtime needed to cluster each of the systems in this case study.

Mitchell and Mancoridis [2003; 2008] continue their work by proposing an evaluation technique for clustering based on the search landscape of the graph being clustered. By gaining insight into the search landscape, the quality of a typical clustering result can be determined. The Bunch software clustering system is examined. Authors model the search landscape of each system undergoing clustering, and then analyze how Bunch produces results within this landscape in order to understand how Bunch consistently produces similar results. Studying the search landscape of any problem is very beneficial when attempting to understand why certain changes to, e.g., the fitness function or the operators, have the kind of effect they have on the results.

The search landscape is modeled using a series of views and examined from two different perspectives. The first perspective examines the structural aspects of the search landscape, and the second perspective focuses on the similarity aspects of the landscape. The structural search landscape highlights similarities and differences from a collection of clustering results by identifying trends in the structure of graph partitions. The similarity search landscape focuses on modeling the extent of similarity across all of the clustering results.

The results produced by Bunch appear to have many consistent properties. By examining views that compare the cluster counts to the MQ values, it can be noticed that Bunch tends to converge to one or two “basins of attraction” for all of the systems studied. Also, for the real software systems, these attraction areas appear to be tightly

paced. An interesting observation can be made when examining the random system with a higher edge density: although these systems converged to a consistent MQ, the number of clusters varied significantly over all of the clustering runs. The percentage of intra-edges in the clustering results indicates that Bunch produces consistent solutions that have a relatively large percentage of intra-edges. Also, the intra-edge percentage increases as the MQ values increase. It seems that selecting a random partition with a high intra-edge percentage is highly unlikely. Another observation is that Bunch generally improves the MQ of real software systems much more and that of random systems with a high edge density. Number of clusters produced compared with number of clusters in the random starting point indicates that the random starting points appear to have a uniform distribution with respect to the number of clusters. The view shows that Bunch always converges to a “basin of attraction” regardless of the number of clusters in the random starting point.

When examining the structural views collectively, the degree of commonality between the landscapes for the systems in the case study is quite similar. Since the results converge to similar MQ values, Mitchell and Mancoridis speculate that the search space contains a large number of isomorphic configurations that produce similar MQ values. Once Bunch encounters one of these areas, its search algorithms cannot find a way to transform the current partition into a new partition with higher MQ. The main observation is that the results produced by Bunch are stable. However, the true meaning of the result is that the Bunch actually gets stuck to a local optimum, and cannot find a way to escape that local optimum. This is naturally the problem for nearly all search algorithms: a true global optimum is not even expected to be found. Doing this kind of fitness landscape study should, however, aid in designing the algorithm so that it would have a better chance of escaping the local optimum, as the fitness landscape reveals what drives the algorithm to the particular basins of attractions that it chooses.

In order to investigate the search landscape further Mitchell and Mancoridis measure the degree of similarity of the placement of nodes into clusters across all of the clustering runs to see if there are any differences between random graphs and real software systems. Bunch creates a subsystem hierarchy, where the lower levels contain detailed clusters, and higher levels contain clusters of clusters. Results from similarity measures indicate that the results for the real software systems have more in common than the results for random systems do. Results with similarity measures also support the isomorphic “basin of attraction” conjecture proposed.

Mitchell et al. [2000] have developed a two step process for reverse engineering the

software architecture of a system directly from its source code. The first step involves clustering the modules from the source code into abstract structures called subsystems. Bunch is used to accomplish this. The second step involves reverse engineering the subsystem-level relations using a formal (and visual) architectural constraint language. Using the reverse engineered subsystem hierarchy as input, a second tool, ARIS, is used to enable software developers to specify the rules and relations that govern how modules and subsystems can relate to each other. This again gives the user the possibility to use his/her own expertise as a basis for the fitness function, so it is not based on metrics.

ARIS takes a clustered MDG as input and attempts to find the missing style relations. The goal is to induce a set of style relations that will make all of the use relations well-formed. A relation is well-formed if it does not violate any permission rule described by the style; this is called the edge repair problem. The relative quality of a proposed solution is evaluated by an objective function. The objective function that is designed into the ARIS system measures the well-formedness of a configuration in terms of the number of well-formed and ill-formed relations it contains. The quality measurement $Q(C)$ for configuration C gives a high quality score to configurations with a large number of well-formed use relations and a low quality score to configurations with a large number of ill-formed style relations or large visibility. Here, as in many other cases where some external expertise is added, the actual fitness function seems simple (only calculating sums and divisions), but much work is first needed by the user to define the input variables, here rules, for the fitness function. Again, it raises the question: what kind of automation is expected from a tool based on search algorithms? Is it good enough that the algorithm only performs a small task and expects a lot of input, or should the algorithm be better defined so that it actually diminishes the work load of the software designer instead of increasing it?

Two search algorithms have been implemented to maximize the objective function: HC and edge removal. The HC algorithm starts by generating a random configuration. Incremental improvement is achieved by evaluating the quality of neighboring configurations. A neighboring configuration C_n is one that can be obtained by a small modification to the current configuration C . The search process iterates as long as a new C_n can be found such that $Q(C_n) > Q(C)$.

The edge removal algorithm is based on the assumption that as long as there exists at least one solution to the edge repair problem for a system with respect to a style specification, the configuration that contains every possible reparable relation will be one of the solutions. Using this assumption, the edge removal algorithm starts by generating

the fully reparable configuration for a given style definition and system structure graph. It then removes relations, one at a time, until no more relations can be removed without making the configuration ill-formed. A case study is performed, where the results seem promising as they give intuition to the nature of the system. This may be beneficial for novice designers, who do not have very much knowledge of the system, but it should be assumed that the developers who have to define the rules that the tool is based on already have a mature idea of the system in order to be able to define those rules.

Mahdavi et al. [2003a; 2003b] show that results from a set of multiple hill climbs can be combined to locate good “building blocks” for subsequent searches. Building blocks are formed by identifying the common features in a selection of best hill climbs. This process reduces the search space, while simultaneously ‘hard wiring’ parts of the solution. Mahdavi et al. also investigate the relationship between the improved results and the system size.

An initial set of hill climbs is performed and from these a set of best hill climbs is identified according to some “cut off” threshold. Using these selected best hill climbs the common features of each solution are identified. These common features form building blocks for a subsequent hill climb. A building block contains one or more modules fixed to be in a particular cluster, if and only if all the selected initial hill climbs agree that these modules were to be located within the same cluster. Since all the selected hill climbs agree on these choices, it is likely that good solutions will also contain these choices.

The implementation uses parallel computing techniques to simultaneously execute an initial set of hill climbs. From these climbs the authors experiment with various cut off points ranging from selecting the best 10% of hill climbs to the best 100% in steps of 10%. The building blocks are fixed and a new set of hill climbs are performed using the reduced search space. The principal research question is whether or not the identification of building blocks improves the subsequent search.

A variety of experimental subjects are used. Two types of MDGs are used: first type contains non-weighted edges, second type has weighted edges. The MQ values are gathered after the initial and the final climbs, and compared for difference. Statistical tests provide some evidence towards the premise that the improvement in MQ values is less likely to be a random occurrence due to the nature of the hill climb algorithm. The improvement is observed for MDGs with and without weighted edges and for all size MDGs.

Larger MDGs show more substantial improvement when the best initial fitness is

compared with the best final fitness values. One reason for observing more substantial improvement in larger MDGs may be attributed to the nature of the MQ fitness measure. To overcome the limitation that MQ is not normalized, the percentage MQ improvement of the final runs over the initial runs is measured. These statistical tests show no significant correlation between size and improvement in fitness for both weighted and non-weighted MDGs.

The increase in fitness, regardless of number of nodes or edges, tends to be more apparent as the building blocks are created from a smaller selection of individuals. This may signify some degree of importance for the selection process.

Results indicate that the subsequent search is narrowed to focus on better solutions, that better clustering are obtained and that the results tend to improve when the selection cut off is higher. These initial results suggest that the multiple hill climbing technique is potentially a good way of identifying building blocks. Authors also found that although there was some correlation between system size and various measures of the improvement achieved with multiple hill climbing, none of these correlations is statistically significant. These results would provide an interesting starting point to a study where the building blocks achieved with multiple hill climbs could be used to initialize the first population given to a genetic algorithm.

Harman et al. [2002] experiment with fitness functions derived from measures of modules granularity, cohesion and coupling for software modularization. They present a new encoding and crossover operator and report initial results based on simple component topology. The new representation allows only one representation per modularization and the new crossover operator attempts to preserve building blocks [Salomon, 1998].

Harman et al. [2002] present the problem of finding a representation for modularization so that “non-unique representations of modularizations artificially increase the search space size, inhibiting search-based approaches to the problem”. In their approach modules are numbered, and elements allocated to module numbers using a simple look-up table. Component number one is always allocated to module number one. All components in the same module as component number one are also allocated to module number one. Next, the lowest numbered component, n , not in module one, is allocated to module number two. All components into the same module as component number n are allocated to module number two. This process is repeated, choosing each lowest number unallocated component as the defining element for the module. This representation must be renormalized when components move as the result of mutation

and crossover. The chosen method clearly saves resources and clarifies the search space as there are no alternative representations for the same solution.

Harman et al.'s crossover operator attempts to preserve partial module allocations from parents to children in an attempt to promote good building blocks. Rather than selecting an arbitrary point of crossover within the two parents, a random parent is selected and one of its arbitrarily chosen modules is copied to the child. The allocated components are removed from both parents. This removal prevents duplication of components in the child when further modules are copied from one or the other parent to the child. The process of selecting a module from a parent and copying to the child is repeated and the copied components are removed from both parents until the child contains a complete allocation. This approach ensures that at least one module from the parents is preserved (in entirety) in the child and that parts of other modules will also be preserved. As it is not clarified how the modules are represented in the chromosome, it is not, however, exactly clear how risky it would be to perform traditional crossovers with the selected encoding. In fact, it seems perfectly possible to make such an encoding that supports building blocks even with the traditional operators.

The fitness function maximizes cohesion and minimizes coupling. In order to capture the additional requirement that the produced modularization has a granularity (number of modules) similar enough to the initial granularity, a polynomial punishment factor is introduced into the fitness function to reward solutions as they approach the target value for granularity of the modularization. The granularity is normalized to a percentage. The three fitness components are given equal weights.

A standard one-point crossover is also implemented for comparison. The GA with the novel crossover outperforms the one with the traditional one, although it quickly becomes trapped in local optima. This would suggest that the attempt to reserve building blocks might actually be “too strong”, as the GA does not have any method to escape the local optimum. Results also show that the novel GA is more sensitive to inappropriate choices of target granularity than any other approach.

Harman et al. [2005] present empirical results which compare the robustness of two fitness functions used for software module clustering: MQ is used exclusively for module clustering and EVM [Tucker et al., 2001] has previously been applied to time series and gene expression data. The clustering algorithm is based upon the Bunch algorithm [Mancoridis et al., 1999] and redefined. Three types of MDGs were studied: real program MDGs, random MDGs and perfect MDGs.

The primary findings are that searches guided by both fitness functions degrade

smoothly as noise increases, but EVM would appear to be the more robust fitness function for real systems. Searches guided by MQ behave poorly for perfect and near-perfect module dependency graphs (MDGs). The results of perfect graphs (MDGs) show however, that EVM produces clusterings which are perfect and that the clusterings produced stay very close to the perfect results as more noise is introduced. This is true both for the comparison against the perfect clustering and the initial clustering. By comparison, the MQ fitness function performs much worse with perfect MDGs. Comparing results for random and real MDGs, both fitness functions are fairly robust. Further results show that searches guided by MQ do not produce the perfect clustering for a perfect MDG but a clustering with higher MQ values. This very strongly suggests that fitness metrics indeed do not actually match what is truly desired of the solution.

These results highlight a possible weakness in MQ as a guiding fitness function for modularization searches: it may be possible to improve upon it by addressing that issue. The results show that EVM performs consistently better than MQ in the presence of noise for both perfect and real MDGs but worse for random MDGs. The results for both fitness functions are better for perfect or real graphs than random graphs, as expected. As the real programs increase in size, there appears to be a decrease in the difference between the performance of searches guided by EVM and those guided by MQ. The results show that both metrics are relatively robust in the presence of noise, with EVM being the more robust of the two.

This study is a significant indicator that fitness metrics should never be blindly trusted. The problem here is particularly curious, as the developers of the MQ metric showed the results (achieved with the aid of this metric) to actual software designers, who were reported to give positive feedback. Thus, it could be assumed that the MQ metric was based on real feedback from human designers. However, it still failed in comparison to another metric and could not produce optimal results. These results suggest that the quality requirements for software design problems are extremely difficult to define, which in turn makes the definition of a proper fitness function a demanding task.

Antoniol et al. [2003] present an approach to re-factoring libraries with the aim of reducing the memory requirements of executables. The approach is organized in two steps: the first step defines an initial solution based on clustering methods, while the second step refines the initial solution with a GA. Antoniol et al. [2003] propose a GA approach that considers the initial clusters as the starting population, adopts a knowledge-based mutation function and has a multi-objective fitness function. Tests on medium and

large open source software systems have effectively produced smaller, loosely coupled libraries, and reduced the memory requirement for each application.

Given a system composed by applications and libraries, the idea is to re-factor the biggest libraries, splitting them into two or more smaller clusters, so that each cluster contains symbols used by a common subset of applications (i.e., Antoniol et al. made the assumption that symbols often used together should be contained in the same library). Given that, for each library to be re-factored, a Boolean matrix MD is composed.

Antoniol et al. have chosen to apply the Silhouette statistic [Kaufman and Rousseeuw, 1990] to compute the optimal number of clusters for each MD matrix. Once the number of clusters is known for each “old library”, agglomerative-nesting clustering was performed on each MD matrix. This allows the identification of a certain number of clusters. These clusters are the new candidate libraries. When given a set of all objects contained into the candidate libraries, a dependency graph is built, and the removal of inter-library dependencies can therefore be brought back to a graph partitioning problem.

The encoding is the achieved bit-matrix, where for each matrix point $[x, y]$ has value 1 if the object y is used by the application or library defined by x , and 0 otherwise. The GA is initialized with the encoding of the set of libraries obtained in the previous step. This encoding method is well-chosen, as there is no need to make any unnecessary transformation between two encodings, and the genetic operations can be easily defined for a matrix.

The mutation operator works in two modes: normally, a random column is taken and two random rows are swapped. When cloning an object, a random position in the matrix is taken; if it is zero and the library is dependent on it, then the mutation operator clones the object into the current library. Of course the cloning of an object increases both linking and size factors, therefore it should be minimized. This GA activates the cloning only for the final part of the evolution (after 66%) of generations in their case studies. This strategy favors dependency minimization by moving objects between libraries; then, at the end, remaining dependencies are attempted to remove by cloning objects. The crossover is a one-point crossover: given two matrices, both are cut at the same random column, and the two portions are exchanged. Population size and number of generations were chosen by an iterative procedure.

The fitness function attempts to balance three factors: the number of inter-library dependencies at a given generation, the total number of objects linked to each application that should be as small as possible, and the size of the new libraries. A unitary weight is set to the first factor, and two weights are selected using an iterative trial-and-error

procedure, adjusting them each time until the factors obtained at the final step are satisfactory. The partitioning ratio is also calculated. Case study results show that the GA manages to considerably reduce the amount of dependencies, while the partition ratio stays nearly the same or slightly reduced. The proposed re-factoring process allows obtaining smallest, loosely coupled libraries from the original biggest ones.

The selected fitness function would benefit from more enhanced techniques to deal with multi-objectivity. Also, in multi-objective problems there usually are cases when one goal may need to be emphasized at the cost of another goal. In this case there are no such tests, as the weights are simply optimized for a general case. It would be interesting to see what kinds of results are achieved, if, e.g., the size of libraries is shown significantly more appreciation than the number of inter-library dependencies. If these cases would produce interesting modularizations, then a Pareto optimal fitness function would be good to experiment with.

Di Penta et al. [2005] build on these results and present a software renovation framework (SRF), a toolkit that covers several aspects of software renovation, such as removing unused objects and code clones, and refactoring existing libraries into smaller ones. Refactoring has been implemented in the SRF using a hybrid approach based on hierarchical clustering, GAs and hill climbing, also taking into account the developer's feedback. Most of the SRF activities deal with analyzing dependencies among software artifacts, which can be represented with a dependency graph.

Software systems are represented by a system graph SG , which contains the sets of all object modules, all software system libraries, all software system applications and the set of oriented edges representing dependencies between objects. The refactoring framework consists of several steps: 1. software systems applications, libraries and dependencies among them are identified, 2. unused functions and objects are identified, removed or factored out, 3. duplicated or cloned objects are identified and possibly factored out, 4. circular dependencies among libraries are removed, or at least reduced, 5. large libraries are refactored into smaller ones and, if possible, transformed into dynamic libraries, and 6. objects which are used by multiple applications, but which are not yet organized into libraries, are grouped into new libraries. Step five, splitting existing, large libraries into smaller clusters of objects, is now studied more closely.

The refactoring of libraries is done in the SRF in the following steps: 1. determine the optimal number of clusters and an initial solution, 2. determine the new candidate libraries using a GA, 3. ask developers' feedback. The effectiveness of the refactoring process is evaluated by a quality measure of the new library organization, the Partitioning

Ratio, which should be minimized.

The genome representation and mutations are as previously presented by Antoniol et al. [2003]. Now, however, the developers may also give a Lock Matrix when they strongly believe that an object should belong to a certain cluster. The mutation operator does not perform any action that would bring a genome in an inconsistent state with respect to the Lock Matrix. The crossover is the one point crossover, which exchanges the content of two genome matrices around a random column.

The fitness function F should balance four factors: the number of inter-library dependencies, the total number of objects linked to each application, the size of new libraries and the feedback by developers. Thus, developer feedback is brought to the fitness function as an additional element to those already presented by Antoniol et al. [2003]. The fitness function F is defined to consist of the Dependency factor DF, the Partitioning ratio PR, the Standard deviation factor SD and the Feedback factor FF. The FF is stored in a bit-matrix FM, which has the same structure of the genome matrix and which incorporates those changes to the libraries that developers suggested. Each factor of the fitness function is given a separate real, positive weight. DF is given weight 1, as it has maximum influence.

Di Penta et al. report that the presented GA suffers from slow convergence. To improve its performance, it has been hybridized with HC techniques. In their experiment, applying HC only to the last generation significantly improves neither the performance nor the results, but applying HC to the best individuals of each generation makes the GA converge significantly faster. In the case study, the GA reduces dependencies of one library to about 5% of the original amount while keeping the PR almost constant. For two other libraries, a significant reduction of inter-library dependencies is obtained while slightly reducing PR in one and increasing the PR in the other. The addition of HC into GA does not improve the fitness values, since GA also converges to similar results, when it is executed on an increased number of generations and increased population size. Noticeably, performing HC on the best individuals of each generation produces a drastic reduction in convergence times. These results show that hybrid algorithms are a strong candidate when attempting to improve the results of search-based approaches.

Huynh and Cai [2007] present an automated approach to check the conformance of source code modularity to the designed modularity. Design structure matrices (DSMs) are used as a uniform representation and they are automatically clustered and checked for conformance by a GA. A design DSM and source code DSM work at different levels of abstraction. A design DSM usually needs higher level of abstraction to obtain the full

picture of the system, while a source code DSM usually uses classes or other program constructs as variables labeling the rows and columns of the matrix. Given two DSMs, one at the design level and the other at the source code level, the GA takes one DSM as the optimal goal and searches for a best clustering method in the other DSM that maximizes the level of isomorphism between the two DSMs. One of the two DSMs is defined as the sample graph, and the other one as a model graph, and finally a conformance criterion is defined. This approach appears beneficial especially in the area of program comprehension and validity checking (as well as purely increasing program quality). Performing conformance checks on a large test set of programs could even produce general ideas on where the programs generally differ from the initial design.

To determine the conformance of the source code modularity to the high level design modularity the variables of the sample graph are clustered and thus a new graph is formed, which is called the conformance graph. Each vertex of the conformance graph is associated with a cluster of variables from the sample graph. The more conforming the source code modularity is to the design modularity, the closer to isomorphic the conformance graph and the model graph will be. In computing the level of isomorphism between two graphs, the graph edit distance is computed between the graphs.

With the given representation of the problem, a GA is formulated with which the goal is to find the clustering of sample graph vertices such that the conformance graph of these clustered nodes is isomorphic, or almost isomorphic, to the model graph. This is a projection. The algorithm first creates an initial population of random projections. The fitness function is defined as $f = -D - P - \lambda - \phi$, where D is the graph edit distance, P is a penalty, and λ and ϕ provide finer differentiation between mappings with the same graph edit distance. The last two functions allow configuring a sample graph so that it can be clustered in different ways, each corresponding to how the design targeted DSM is clustered. The dissimilarity function λ is used to calculate how separated components from each directory grouping are. If a sample graph node attribute matches a name pattern specified by the user but is not correctly mapped to the model graph vertex then the fitness of the projection is reduced through ϕ . Interestingly, the fitness function only measures negative aspects, quite differently to other fitness functions in modularization, which usually attempt to maximize at least some quality value.

The GA is run on two DSM models of an example software. The experiments consistently converge to produce the desired result, although the tool sometimes produces a result that is not the desired view of the source code, even though the graphs are isomorphic, i.e., the result conforms with the model. The experiment shows the feasibility

of using a GA to automatically cluster DSM variables and correctly identify links between source code components and high level design components. The results support the hypothesis that it is possible to check the conformance between source code structure and design structure automatically, and this approach has the potential to be scaled for use in large software systems.

4.3 Summarizing remarks

The majority of the studies relating to search-based software clustering have been done with the Bunch tool, which has seen many improvements. This is very promising for other approaches to search-based design as well, as the tool has been accepted for use in the software engineering community. However, there are still many open questions in the area of software modularization. What is a proper encoding to represent a modularization problem? This question is especially highlighted by the study made by Harman et al. [2002], as they point out the massive amount of redundant information in many encodings. What is a proper fitness metric for modularizations? Again, the study comparing the very popular MQ metric with another modularization metric (EVM), showed that while the metric is robust (as already validated by its developers), it can be outperformed. How can metrics be relied on then? Di Penta et al. [2005] have attempted to enhance the performance of their tool by giving the developers a chance to formalize their knowledge on quality. However, defining quality as a matrix form cannot be very user-friendly.

As stated, the research on software clustering revolves quite strongly around Bunch or the MQ metric. The main exception to this is the studies made by Antoniol et al., [2003] and Di Penta et al. [2005] who use a matrix to encode the modularization and use matrix-related or metrics instead of the MQ, and Hyunh and Cai [2007], who use a matrix and then turn it into a graph, and use graph related metrics to evaluate the quality of a proposed solution. Especially the approach by Hyunh and Cai [2007] is significantly different to Bunch, as two modularizations are ultimately compared, while Bunch attempts to ameliorate a poor modularization without a certain goal it is aiming towards. Thus, there is much room in search-based software clustering for alternative methods, as competition always makes each different approach strive towards even better solutions.

Table 3. Research approaches in search-based software clustering

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|--|---|---|--------------------|----------|-----------|---|--------------------------------|--|
| Mancoridis et al. [1998] | Automation of partitioning components of a system into clusters | System given as a module dependency graph (MDG) | MDG | N/A | N/A | Minimize inter-connectivity, maximize intra-connectivity, combined as modularization quality (MQ) | Optimized clustering of system | |
| Doval et al. [1999] | Automation of partitioning components of a system into clusters | MDG | String of integers | Standard | Standard | MQ | Optimized clustering of system | Continued work from Mancoridis et al. [1998] by implementing a GA |
| Mancoridis et al. [1999] | Automation of partitioning components of a system into clusters | MDG | MDG | N/A | N/A | MQ | Optimized clustering of system | Continued work from Mancoridis et al. [1998]; characteristics of modules taken into account in clustering operations |
| Mitchell and Mancoridis [2002; 2006; 2008] | Automation of partitioning components of a system into clusters | MDG | String of integers | Standard | Standard | MQ as a sum of clustering factors | Optimized clustering of system | Continued work from Doval et al. [2002]; new definition of the modularization quality and an enhanced HC algorithm |

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|--------------------------------------|---|----------------------------|--------------------|----------|-----------|--|--------------------------------|--|
| Mitchell and Mancoridis [2003; 2008] | Automation of partitioning components of a system into clusters | MDG | String of integers | Standard | Standard | MQ, search landscape | Optimized clustering of system | Continued work from Mitchell and Mancoridis [2002; 2006; 2008]; search landscape taken into account |
| Mitchell et al. [2000] | Automated reverse engineering from source code to architecture | Source code of application | N/A | N/A | N/A | Quality based on use and style relations | Software architecture | HC and edge removal are used as search algorithms from MDG to architecture |
| Mahdavi et al. [2003a; 2003b] | Automated clustering of system | MDG | String of integers | Standard | Standard | MQ | Optimized clustering of system | Multiple hill climbs are used as search algorithm; building blocks are preserved by using parallel hill climbs |

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|------------------------|---|--|---|---|--|--|---|---|
| Harman et al. [2002] | New encoding and crossover introduced | System as modules and elements | Look-up table for modules | Move component from one module to another | New crossover, preserves partial module allocations | Maximize cohesion, minimize coupling | Optimized clustering | |
| Harman et al. [2005] | Comparison of robustness between two fitness functions | Clustered system | N/A | N/A | N/A | MQ compared against EVM | - | |
| Antoniol et al. [2003] | Cluster optimization | System containing applications and libraries | Bit matrix | Two random rows of a column in matrix are swapped or an object is cloned by changing a value from zero to one | A random column is taken as split point and contents are swapped | Inter-library dependencies, number of object-application links and size of libraries | Optimized clustering, sizes and dependencies between libraries diminished | Optimal number of clusters is calculated for a matrix with the Silhouette statistic |
| Di Penta et al. [2005] | A refactoring framework taking into account several aspects of software quality when refactoring existing system. | Software system as a system graph SG | Bit matrix; each library of clusters is represented by a matrix | Swapping two bits in a column or changing a value from 0 to 1 (taking into account preconditions) | N/A | Dependency factor, partitioning ratio, standard deviation and feedback | Refactored libraries | HC and GA used. |
| Hyunh and Cai [2007] | Conformance check of actual design to suggested design | Design structure matrices for design and source code (DSM) | Graph constructed of DSM | N/A | N/A | Graph edit distance, penalty and differentiation between graphs with same distance | Optimized clustering of actual design conforming to suggested design | |

5. SOFTWARE REFACTORING

5.1. Background

Software evolution often results in “corruption” in software design, as quality is overlooked while new features are added, or the old software should be modified in order to ensure the highest possible quality. At the same time resources are limited. Refactoring and in particular the miniaturization of libraries and applications are therefore necessary. Program transformation is useful in a number of applications including program comprehension, reverse engineering and compiler optimization. A transformation algorithm defines a sequence of transformation steps to apply to a given program and it is described as changing one program into another. It involves altering the program syntax while leaving its semantics unchanged. In object-oriented design, one of the biggest challenges when optimizing class structures using random refactorings is to ensure behavior preservation. One has to take special care of the pre- and post-conditions of the refactorings.

There are three problems with treating software refactoring as a search-based problem. First, how to determine which are the useful metrics for a given system. Second, finding how best to combine multiple metrics. Third is that while each run of the search generates a single sequence of refactorings, the user is given no guidance as to which sequence may be best for their given system, beyond their relative fitness values.

In practice, refactoring (object-oriented software) can begin with simple restructurings of the class structure and being very close to software clustering, and then move on to a more detailed level of moving elements from one class to another. The lowest level of refactoring already deals with code, as procedures are sliced to eliminate redundancy or transformed in order to simplify the program or make it more efficient. The following subsection presents approaches where search-based techniques have been used to automatically achieve refactorings, as well as a study on a new method for evaluating the fitness of a refactored software. Summarizing remarks are then presented, and the fundamentals of each study are collected in Table 4.

5.2. Approaches

Seng et al. [2005] describe a methodology that computes a subsystem decomposition that can be used as a basis for maintenance tasks by optimizing metrics and heuristics of good subsystem design. GA is used for automatic decomposition. If a desired architecture is given, e.g., a layered architecture, and there are several violations, this approach attempts to determine another decomposition that complies with the given architecture by moving

classes around. Instead of working directly on the source code, it is first transformed into an abstract representation, which is suitable for common object-oriented language.

In the GA, several potential solutions, i.e., subsystem decompositions, form a population. The initial population can be created using different initialization strategies. Before the algorithm starts, the user can customize the fitness function by selecting several metrics or heuristics as well as by changing thresholds. The model is a directed graph. The nodes of the graph can either represent subsystems or classes. Edges between subsystems or subsystems and classes denote containment relations, whereas edges between classes represent dependencies between classes. The approach is based on the Grouping GA [Falkenaur, 1998], which is particularly well suited for finding groups in data. For chromosome encoding, subsystem candidates are associated with genes and the power set of classes is used as the alphabet for genes. Consequently, a gene is associated with a set of classes, i.e., an element of the power set. This representation allows a one-to-one mapping of geno- and phenotype to avoid redundant coding.

An adapted crossover operator and three kinds of mutation are used. The operators are adapted so that they are non-destructive and preserve a complete subsystem candidate as far as possible. The *split&join* mutation either divides one subsystem to two, or vice versa. The operator splits a subsystem candidate in such a way that the separation in two subsystem candidates occurs at a loosely associated point in the dependency graph. *Elimination* mutation deletes a subsystem candidate and distributes its classes to other subsystem candidate, based on association weights. *Adoption* mutation tries to find a new subsystem candidate for an orphan, i.e., a subsystem candidate containing only a single class. This operator moves the orphan to the subsystem candidate that has the highest connectivity to the orphan. The chosen mutations support reversibility, i.e., a GA can always backtrack its steps. The *split&join* mutation is obvious in this case, but also the adoption mutation can be seen as a reverse operation for the elimination, if a new subsystem can be created dynamically.

Initial population supports the building block theorem. Randomly selected connected components of the dependency graph are taken for half the population and highly fit ones for the rest. The crossover operator forms two children from two parents. After choosing the parents, the operator selects a sequence of subsystem candidates in both parents, and mutually integrates them as new subsystem candidates in the other parent, and vice versa, thus forming two new children consisting of both old and new subsystem candidates. Old subsystem candidates which now contain duplicated classes are deleted, and their non-duplicated classes are collected and distributed over the remaining subsystem candidates.

Fitness function is defined as $f = w_1 * \text{cohesion} + w_2 * \text{coupling} + w_3 * \text{complexity} + w_4 * \text{cycles} + w_5 * \text{bottlenecks}$. Again the fitness function is based on the two most used metrics, cohesion and coupling, but introduces some new interesting concepts from OO design, such as cycles and bottlenecks, which are more defined than the usual general metrics.

For evaluation, a tool prototype has been implemented. Evaluation on the clustering of different software systems has revealed that results on roulette wheel selection are only slightly better than those of tournament selection. The adapted operators allow using a relatively small population size and few generations. Results from a Java case study show that the approach works well. Tests on optimizing subsets of the fitness function show that only if all criteria are optimized, the authors are able to achieve a suitable compromise with very good complexity, bottleneck and cyclomatic values and good values for coupling and cohesion. Again, as the work here is very similar to optimal software clustering, it can be questioned whether the metrics used in those studies, that mainly calculate modified values for coupling and cohesion, are actually sufficient.

Seng et al. [2006] have continued their work by developing a search-based approach that suggests a list of refactorings. The approach uses an evolutionary algorithm and simulated refactorings that do not change the system's externally visible behavior. The source code is transformed into a suitable model – the phenotype. The genotype consists of the already executed refactorings. Model elements are differentiated according to the role they play in the system's design before trying to improve the structure. Not all elements can be treated equally, because the design patterns sometimes deliberately violate existing design heuristics. The approach is restricted to those elements that respect general design guidelines. Elements that deliberately do not respect them are left untouched in order to preserve the developers conscious design decisions. The notion of applying something that is known to somehow worsen the quality of a system is peculiar. In a way this is natural, as there are always trade-offs when trying to optimize conflicting quality values, but each decision should have a positive affect from some perspective. Hence, it is odd that no quality evaluator has been found that would prevent the elimination of these “deliberately violating” patterns.

The initial population is created by copying the model extracted from the source code a selected number of times. Selection for a new generation is made with tournament selection strategy. The optimization stops after a predefined number of evolution steps. The source code model is designed to accommodate several object-oriented languages. The basic model elements are classes, methods, attributes, parameters and local variables.

In addition, special elements called access chains are needed. An access chain models the accesses inside a method body, because it is needed to adapt these references during the optimization. If a method is moved, the call sites need to be changed. An access chain therefore consists of a list of accesses. Access chains are hierarchical, because each method argument at a call site is modeled as a separate access chain that could possibly contain further access chains.

The model allows to simulate most of the important refactorings for changing the class structure of a system, which are extract class, inline class, move attribute, push down attribute, pull up attribute, push down method, pull up method, extract superclass and collapse class hierarchy. The genotype consists of an ordered list of executed model refactorings including necessary parameters. The phenotype is created by applying these model refactorings in the order that is given by the genotype to the initial source code model. Therefore the order of the model refactorings is important, since one model refactoring might create the necessary preconditions for some of the following ones.

Mutation extends the current genome by an additional model refactoring; the length of the genome is unlimited. Crossover combines two genomes by selecting the first random n model refactorings from parent one and adding the model refactorings of parent two to the genome. The refactorings from parent one are definitely safe, but not all model refactorings of parent two might be applicable. Therefore, the model refactorings are applied to the initial source code model. If a refactoring that cannot be executed is encountered due to unsatisfied preconditions, it is dropped. Seng et al. argue that the advantage of this crossover operator is that it guarantees that the externally visible behavior is not changed, while the drawback is that it takes some time to perform the crossover since the refactorings need to be simulated again. This approach is quite similar to that of Amoui et al. [2006], discussed in Section 3, who approach the problem from a slightly higher level by using architectural design patterns as refactoring, but similarly search for the optimal transformation sequence.

Fitness is a weighted sum of several metric values and is designed to be maximized. The properties that should be captured are coupling, cohesion, complexity and stability. For coupling and cohesion, the metrics from Briand's [2000] catalogue are used. For complexity, weighted methods per class (WMC) and number of methods (NOM) are used. The formula for stability is adapted from the reconditioning of subsystem structures.
$$\text{Fitness} = \sum(\text{weight}_m * (M(S) - M_{\text{init}}(S)) / (M_{\text{max}}(S) - M_{\text{init}}(S)))$$
 Before optimizing the structure the model elements are classified according to the roles they play in the systems design, e.g., whether they are a part of a design pattern.

Tests show that after approximately 2000 generations in a case study the fitness value does not significantly change anymore. The approach is able to find refactorings that improve the fitness value. Actually, this is to be expected, as it would be rather surprising if it did not improve the fitness value, as then there would be something significantly wrong with the GA. Thus, more importantly, in order to judge whether the refactorings make sense, they are manually inspected by the authors, and from their perspective, all proposed refactorings can be justified. As a second goal, the authors modify the original system by selecting 10 random methods and misplacing them. The approach successfully moves back each method at least once.

O’Keeffe and Ó Cinnéide [2004] have developed a prototype software engineering tool capable of improving a design with respect to a conflicting set of goals. A set of metrics is used for evaluating the design quality. As the prioritization of different goals is determined by weights associated with each metric, a method is also described of assigning coherent weights to a set of metrics based on object-oriented design heuristics.

The presented tool, *Dearthóir*, is a prototype for design improvement, as it restructures a class hierarchy and moves methods within it in order to minimize method rejection, eliminate code duplication and ensure superclasses are abstract when appropriate. The refactorings are behavior-preserving transformations in Java code. The refactorings employed are limited to those that have an effect on the positioning of methods within an inheritance hierarchy. Contrary to most other approaches, this tool uses simulated annealing to find close-to-optimum solutions to this combinatorial optimization problem. In order for the SA search to move freely through the search space every change to the design must be reversible. To ensure this, pairs of refactoring have been chosen that complement each other. The refactoring pairs are: 1. move a method up or down in the class hierarchy, 2. extract (from abstract class) or collapse a subclass, 3. make a class abstract or concrete, and 4. change superclass link of a class.

The following method is intended to filter out heuristics that cannot easily be transformed into valid metrics because they are vague, unsuitable for the programming language in use, or dependent on semantics. Firstly, for each heuristic: define the property to be maximized or minimized in the heuristic, determine whether the property can be accurately measured, and note whether the metrics should be maximized or minimized. Secondly, identify the dependencies between the metrics. Thirdly, establish precedence between dependent metrics and a threshold where necessary: prioritize heuristics. Fourthly, check that the graph of precedence between metrics is acyclic. Finally, weights should be assigned to each of the metrics according to the precedences

and threshold.

The selected metrics are: 1. minimize rejected methods (RM) (number of inherited but unused methods), 2. minimize unused methods (UM), 3. minimize featureless classes (FC), 4. minimize duplicate methods (DM) (number of methods duplicated within an inheritance hierarchy), 5. maximize abstract superclasses (AS). Metrics should be appreciated so that $DM > RM > FC > AS$, and $UM > FC$. Note that the used metrics are much more specific to the needs of object-oriented design than the general structural metrics that are commonly used. Also, the heuristic of defining the weights (and the metrics) would be very beneficial for many studies, as assigning balanced weights can be a very complex task, and the dependencies between different metrics and their affect to the weights is rarely taken into account (at least so that it would be mentioned in the studies).

Most of the dependencies in the graph do not require thresholds. However, a duplicate method is avoided by pulling the method up into its superclass, which could result in the method being rejected by any number of classes. Therefore a threshold value is established for this dependency. O’Keeffe and Ó Cinnéide argue that it is more important to avoid code duplication than any amount of method rejection; therefore the threshold can be an arbitrarily high number.

A case study is conducted with a small inheritance hierarchy. The case study shows that the metric values for input and output either become better or stay the same. In the input design several classes contain clumps of methods, where as in the output design methods are spread quite evenly between the various classes. This indicates that responsibilities are being distributed more evenly among the classes, which means that components of the design are more modular and therefore more likely to be reusable. This in turn suggests that adherence to low-level heuristics can lead to gains in terms of higher-level goals. Results indicate that a balance between metrics has been achieved, as several potentially conflicting design goals are accommodated.

O’Keeffe and Ó Cinnéide [2006; 2008a] have continued their research by constructing a tool capable of refactoring object-oriented programs to conform more closely to a given design quality model, by formulating the task as a search problem in the space of alternative designs. This tool, CODE-Imp, can be configured to operate using various subsets of its available automated refactorings, various search techniques, and various evaluation functions based on combinations of established metrics.

CODE-Imp uses a two-level representation; the actual program to be refactored is given as source code and represented as its Abstract Syntax Tree (AST) but a more

abstract model called the Java Program Model (JPM) is also maintained, from which metric values are determined and refactoring preconditions are checked. The change operator is a transformation of the solution representation that corresponds to a refactoring that can be carried out on the source code.

The CODE-Imp calculates quality values according to the fitness function and effects change in the current solution by applying refactorings to the AST as required by a given search technique. Output consists of the refactored input code as well as a design improvement report including quality change and metric information.

The refactoring configuration of the tool is constant throughout the case studies and consists of the following fourteen refactorings. Push down/pull up field, push down/pull up method, extract/collapse hierarchy, increase/decrease field security, replace inheritance with delegation/replace delegation with inheritance, increase/decrease method security, made superclass abstract/concrete. During the search process alternative designs are repeatedly generated by the application of a refactoring to the existing design, evaluated for quality, and either accepted as the new current design or rejected. As the current design changes, the number of points at which each refactoring can be applied will also change. In order to see whether refactorings can be made without changing program behavior, a system of conservative precondition checking is employed.

The used search techniques include first-ascent HC (HC1), steepest-ascent HC (HC2), multiple-restart HC (MHC) and low-temperature SA. For the SA, CODE-Imp employs the standard geometric cooling schedule.

The evaluation functions are flexibility, reusability and understandability of the QMOOD hierarchical design quality model [Bansiya and Davis, 2002]. Each evaluation function in the model is based on a weighted sum of quotients on the 11 metrics forming the QMOOD (design size in class, number of hierarchies, average number of ancestors, number of polymorphic methods, class interface size, number of methods, data access metric, direct class coupling, cohesion among methods of class, measure of aggregation and measure of functional abstraction). Each metric value for the refactored design is divided by the corresponding value for the original design to give the metric change quotient. A positive weight corresponds to a metric that should be increased while a negative weight corresponds to metric that should be decreased. It should be noted that while the complexity of the problem grew, as the program representation became more intricate, the number of refactorings (mutations) was more than doubled, this reflected on the need for a significantly more complicated fitness function. The fitness function used in the previous study only contained 5 metrics, while the current one contains 11 metrics

which are grouped into 3 different fitness functions.

All techniques demonstrate strengths. HC1 consistently produces quality improvements at a relatively low cost, HC2 produces the greatest mean quality improvements in two of the six cases, MHC produces individual solutions of highest quality in two cases and SA produced the greatest mean quality improvement in one case. Based on this it would seem that the SA is actually inferior to the different hill climbing approaches, as it only outperformed them in one measure in one test case out of the six. Combining the results of these different search algorithms would be interesting: is it possible to produce such a hybrid that would preserve the strengths from all algorithms?

Inspection of output code and analysis of solution metrics provide some evidence in favor of use of the flexibility metric and even stronger evidence for using the understandability function. The reusability in present form is not found suitable for maintenance because it resulted in solutions including a large number of featureless classes. As these kinds of classes are not generally accepted in OO design (apart from having “technical classes”), one might wonder whether some corrective function could be used in order to prevent featureless classes from appearing to the design. Simple pre-and post-conditions for mutations might very well help dealing with the problem. The authors conclude that both local search and simulated annealing are effective in the context of search-based software refactoring.

O’Keeffe and Ó Cinnéide [2007; 2008b] have further continued their work by implementing also a GA and a multiple ascent HC (MAHC) to the CODE-Imp refactoring tool and further testing the existing search techniques. The encoding, crossover and mutation for the GA are similar to those presented by Seng et al. [2006], and the power of the tool has been increased by adding a number of different refactorings available for use in searching for a superior design.

The fitness function is an implementation of the understandability function from Bansiya and Davis's [2002] QMOOD hierarchical design quality model consisting of a weighted sum of metric quotients between two designs. This choice was clearly inspired by the earlier study, where two other quality functions, flexibility and reusability, did not perform as well in terms of actual quality enhancement. This design quality evaluation function was previously found by the authors to result in tangible improvements to object-oriented program design in the context of search-based refactoring.

Results for the SA support the recommendation of low values for the cooling factor, since more computationally expensive parameters do not yield greater quality function gains.

In summary, SA has several disadvantages: it is hard to recommend a cooling schedule that will generally be effective, results vary considerably across input programs and the search is quite slow. No significant advantage in terms of quality gain was observed that would make up for these shortcomings. The GA has the advantage that it is easy to establish a set of parameters that work well in the general case, but the disadvantages are that it is costly to run and varies greatly for different input programs. Again, no significant advantage in terms of quality gain was observed that would make up for these shortcomings. Multiple-ascent HC stood out as the most efficient search technique in this study: it produced high-quality results across all the input programs, is relatively easy to recommend parameter for and runs more quickly than any of the other techniques examined. Steepest ascent HC produced surprisingly high quality solutions, suggesting that the search space is less complex than might be expected, but is slow when considered its known inability to escape local optima. Results show MAHC to outperform both SA and GA over a set of four input programs. As the genetic algorithm is the most commonly used search technique, these results should stimulate more comparisons between different algorithms. The search space for this problem was, after all, quite large, when taking into account the high number of refactorings that could be applied to a design. Thus, maybe the more refined hill climbing techniques could be compared to the GA.

Quaam and Heckel [2009] apply the Ant Colony Optimization (ACO) [Dorigo, 1992] for software refactoring. The software is represented as a class diagram with methods and attributes, and the refactoring task is considered as a graph transformation problem, which makes it suitable for ACO. In order to perform ACO, five things need to be defined: 1. a set of components C and the edges between them, 2. a set of states as a sequence of components belonging to C , 3. a set of candidate solutions S , with a subset of feasible candidate solutions according to given constraints, 4. a non-empty subset (of S) of optimal solutions, and 5. an evaluation associated to each candidate solution. Based on this, Quaam and Heckel define a graph by associating the set of graph vertices to the set of proposed transformations. Edges are associated with dependencies. The pheromone and heuristic values are associated with the graph edges and are determined by partial evaluations associated with incomplete candidate solutions.

The goal is to find an optimal set of transformations. These transformations are pre-determined based on the given program (graph) and consider, e.g., moving methods and alternating the class hierarchy. An ant begins with an empty solution from the start vertex in the graph and then gradually checks the available refactoring steps in order to construct

a candidate solution. Initially, any random component from C is chosen and then the partial evaluation function will guide the selection of the corresponding edge through the pheromone values. The fitness value is calculated for each feasible sequence of transformations after applying it on the source graph model, the basis for the fitness being the cost of the transformation and the quality of the result. The approach is tested on a small example system.

This approach demonstrates the use of yet another search technique, ACO, which is especially suitable for graph problems. Other choices, however, raise questions particularly on the generality of this approach. It is only tested on a small system, and all the transformations are pre-defined, and dependent on the particular system. How can this approach be generalized to be applied to any system without extensive work required to define all possible transformations of that system, which is incredibly laborious, if the system is large? Also, the details regarding fitness calculations are not very clear.

Jiang et al. [2008a] apply a set of search algorithms to program slicing in order to locate dependence structures. They attempt to find the subsets from all possible sets of program slices that reveal interesting dependence structures. A program is divided into slices according to program points, which are the nodes of a System Dependency Graph (SDG) [Horwitz et al., 1988]. In order to formulate the problem as a search problem, it is instantiated as a set cover problem. With increasing program sizes a search-based approach is extremely suitable for this type of problem.

A program is represented as a bit matrix, where rows indicate program slices and columns indicate program points. The value in point i, j , is 1 if the slice based on criterion i contains the program point j , and 0 if not. A solution should contain as many program points as possible but should have minimum overlap, i.e., slices that contain the same program points.

The fitness function is seen as a parameter to the overall approach of search-based slicing, as choosing the fitness function depends on the properties of the slice set and what the user considers as “interesting” when searching for dependencies. The fitness function is based on metrics that calculate the Coverage and Overlap of the program. Coverage measures how many program points out of all possible points the program contains. Overlap measures the number of program points within the intersection a slicing set. It can be divided in many ways, but Jiang et al. only consider Average, which evaluates the percentage of overlapping program points based on pair-wise calculations, and Maximum, which evaluates the maximum number of overlapping points based on pair-wise calculations. Both Coverage and Overlap are given weights and then combined

for the overall fitness function. Although it is said that the user can define the fitness function based on his/her own desires of what is “interesting”, it is left unclear whether the definitions must rely on the presented metrics or whether the user can build any kind of fitness function. Also, it is not clear how the properties of the slice set affect the choice of fitness function.

Jiang et al. [2008a] implement HC, GA, a Greedy Algorithm [Naeimi et al., 2004] and a Random Search algorithm. The GA uses a multi-point crossover and a standard bit change as a mutation. Elitism and rank selection are used as selection methods. For HC, a multiple restart HC is implemented in order to give it the same amount of computation time as the other algorithms. A Greedy Algorithm constitutes of two sets: a solution set and a candidate set, and three functions: selection, value-computing and solution function. A solution is created out of the solution set and a candidate set represents all possible elements that might be contained in a solution. Selection chooses the most promising candidate to be added to the solution, value-computing function gives a value for the solution and solution function checks whether the final solution has been reached. Here the initial solution set is a binary string with each bit set to 0, and the candidate solution set is made of all the slices. The value-computing function calculates the program points in a solution and the selection function chooses the one with the best coverage and smallest overlap.

An empirical study is made with six open source programs, and possible slices are collected with a separate program from each program’s SDG. The program sizes vary from 37 to 1008 program points. Every other algorithm except the Greedy Algorithm was executed 100 times; the Greedy algorithm gives the same result every time and thus does not need several test runs. For the fitness function using Average Overlap, the Greedy Algorithm performs the best for all but one test case, where HC and GA perform the best. Furthermore, it is seen that for smaller programs HC outperforms GA and Random search. As the program size increases, GA starts to perform better, and wins over HC. For the second fitness function where the Maximum Overlap was used, the results are similar as with the first fitness function. However, in this case GA performs the best of the other algorithms, and HC only beats Random search on the smallest test case. The Greedy Algorithm also outperforms all others in terms of execution time. It is no surprise that the Random Search is outperformed every time. However, it is naturally a bit disappointing that the Greedy Algorithm was superior in every aspect, when compared to other search methods.

Jiang et al. [2008a] make another study by only using the Greedy Algorithm for six

different large programs. As the previous study showed that the Greedy Algorithm outperformed all other studied search algorithms, now it is tested how efficient it is in decomposing a program into a set of slices. Results suggest that less than 20% of a program can be used to decompose the whole program or function.

Jiang et al. [2008b] continue by applying a Greedy Algorithm to procedure splitting. They attempt to split a procedure into two or more sub-procedures in order to improve cohesion. The Greedy Algorithm is used to find close to optimal splitting points.

A slice is represented as a bit matrix. A matrix value depends on whether a program point (i.e., a node in the system's SDG) belongs to a certain slice. The splitting algorithm proceeds in four steps: 1. slice with respect to all nodes in SDG to find all static backward slices, 2. find sets of slices with minimum overlap, 3. recover slice statements by combining nodes that belong to a single statement, 4. make sub-procedures obtained executable.

Results indicate that more than 20% of procedures in all six programs contain independent sub-programs. Also, it would seem that most procedures are not splittable, and the ones that are, can usually be split into only 2 or 3 sub-programs. Splittability appears to correlate with the size of the program.

Fatiregun et al. [2004] use meta-heuristic search algorithms to automate, or partially automate the problem of finding good program transformation sequences. With the proposed method one can dynamically generate transformation sequences for a variety of programs also using a variety of objective functions. The goal is to reduce program size, but the approach is argued to be sufficiently general that it can be used to optimize any source-code level metric. Random search (RS), hill climbing and GA are used.

An overall transformation of a program p to an improved version p' typically consists of many smaller transformation tactics. Each tactic consists of the application of a set of rules. A transformation rule is an atomic transformation capable of performing the simple alterations. To achieve an effective overall program transformation tactic many rules may need to be applied and each would have to be applied in the correct order to achieve the desired results.

In HC, an initial sequence is generated randomly to serve as the starting point. The algorithm is restarted several times using a random sequence as the starting individual each time. The aim is to divert the algorithm from any local optimum.

Each transformation sequence is encoded as an individual that has a fixed sequence length of 20 possible transformations. An example individual is a vector of the transformation numbers. In HC, the neighbor is defined as the mutation of a single gene

from the original sequence. Crossover is the standard one-point crossover. In addition to transformations, cursor moves are also used. The tournament selection is used for selecting mating parents and creating a single offspring, which replaces the worse of the parents. The authors consider optimizing the program with respect to the size of the source-code, i.e., LOC, where the aim is to minimize the number of lines of code as much as possible. This metric is quite simple, and the effects are hardly arguable, if the length of a line of code is somehow restricted.

The fitness is measured as the nominal difference in the lines of code between the source program and the new transformed program created by that particular sequence. This is evaluated by a process of five steps: 1. compute length of the input program, 2. generate the transformation sequence, 3. apply the transformation sequence, 4. compute the current length of the program, 5. compute the fitness, which is the difference between steps 1 and 4.

Results show that GA outperforms both RS and HC. In cases where RS outperformed GA and HC, it was noticed that GA and HC are not “moving” towards areas where potential optimizations could be. Analyzing the GA, the authors believe that the GA potentially kills off good subsequences of transformations during crossover. These results are interesting as this would indicate that the selected (standard) crossover would not support the preservation of building blocks. As discussed in Section 4, it may be that also the encoding could be improved to preserve building blocks. All in all, examining the fitness landscape and rethinking the encoding and crossover operators may be able to improve the results achieved with the GA.

Williams [1998] implements several search algorithms in his REVOLVER system that make program transformations in order to parallelize the program and thus lessen the execution time. The idea is to transform loops in different ways, and as loops are the core of the approach, they are numbered. HC, SA and GA are used, and most interestingly, two different encodings are experimented with.

In the first encoding, Gene-Transformation (GT), each gene represents a transformation that is applied to the system. The gene contains information what transformation is applied, and the number of the loop it is applied to. Three different mutations can be used: changing the transformation, changing the loop number or changing both (i.e., the entire gene). Both one-point and two-point crossovers are implemented. However, in the one-point crossover, the crossover points for the parent chromosomes are chosen individually for each parent, as they might be of unequal length. This approach is applied to HC, SA and GA.

In the second encoding, Gene-Statement (GS), each gene represents a statement in the program, .e.g., an if- or a do-statement, and the chromosome thus represents the program as a sequence of statements. The mutations that are applied are the chosen operations on loops, and applying them to the program. This is actually quite odd, as only loop related transformations are used, but there are only loops in some of the genes. Note, that a mutation will alter the program, as, say, combining two loops will remove the statement representing one of them, and thus shortens the chromosome by one gene. No crossover is used in this representation, and the used algorithms are HC and evolutionary strategy (ES), which is basically a GA, i.e., it has a population and selection, but without the crossover.

The fitness function for both approaches is the actual execution time of the transformed program, and tournament selection is used. In the tests the population size was only 5 for the algorithms with populations, and the number of generations only 50. These parameters seem incredibly low, as there is very little room for versatility in the population, and there is very little time for development also. Thus, one wonders whether the benefits of the GA are truly used in this approach.

Test results on five programs show that the ES and HC with the GS encoding outperformed all other algorithms. The traditional GA appeared the worst. These results further suggest that the population parameter chosen for the traditional GA should be revised, as the GA cannot use its full potential. Interestingly, ES, which also had a population, performed the best. The strength of the GS encoding is also very interesting, considering there is much information in the genes that cannot be mutated. However, ES did not have a crossover, and thus choosing parents is not an issue for this algorithm. All in all, the algorithms were able to improve the execution times significantly.

Ryan and Ivan [1999] have taken a rather different approach to program parallelization, as they encode the program in tree form and use genetic programming as the search algorithm. They use GP in an unusual way, as it does not actually “program”, but searches for the optimal transformations for the program, thus making this study a design problem.

The program is considered as a sequence of instructions. The actual tree given by the GP then comes from examining the atoms representing the instructions, and deciding on transformations based on the type of the instruction. The GP works in two modes: atom mode and loop mode. Each step begins in atom mode, and if the found instruction is a loop, the mode is switched. In atom mode, there are three classes of transformations. The transformations in the first class split the sequence of instructions according to a given

percent, thus forking the execution of a program. The ones in the second class also split the sequence of instructions, but with less effect, as the split point is always either after the first or before the last instruction. The last class delays the execution of the program. Each atom mode transformation is an internal node in a tree, and takes as input the program segment before passing it onto the next transformation. The program segment ultimately diminishes to one atom as transformations are applied. In loop mode the idea is to parallelize each loop by executing each iteration on a different processor, unfortunately, though, this raises issues with data dependencies. A significant operator in loop mode is loop fusion, which combines consecutive loops.

The fitness function is a combination of fitness calculations from the atom mode and the loop mode. For the atom mode the fitness is the execution time and the correctness of the program. For loop mode the fitness is the number of successes for applied loop operators. The initial results are promising; the approach is able to parallelize programs and thus ameliorate them in terms of execution time.

The approach of Ryan and Ivan [1999] appears quite similar to that of Williams [1998] in terms of the choosing loops as a key ingredient in the mutations. However, Ryan and Ivan have taken atom transformations into account as traditional mutations, while Williams has chosen to deal with non-loop structures only at the encoding stage. The fitness function for both approaches is basically the same, as execution time is the most important factor. It would be interesting to study the problem of program parallelization also in terms of other quality factors and as a larger problem in the context of, e.g., distributed systems.

Harman and Tratt [2007] show how Pareto optimality can improve search based refactoring, making the combination of metrics easier and aiding the presentation of multiple sequences of optimal refactorings to users. Intuitively, each value on a Pareto front maximizes the multiple metrics used to determine the refactorings. Through results obtained from three case studies on large real-world systems, it is shown how Pareto optimality allows users to pick from different optimal sequences of refactorings, according to their preferences. Moreover, Pareto optimality applies equally to subsequences of refactorings, allowing users to pick refactoring sequences based on the resources available to implement those refactorings. Pareto optimality can also be used to compare different fitness functions, and to combine results from different fitness functions.

Harman and Tratt use the move method refactoring presented by Seng et al. [2006]. Three systems are used in the case study, all non-trivial real-world systems. The search

algorithm itself is a non-deterministic non-exhaustive hill climbing approach. A random move method refactoring is chosen and applied to the system. The fitness value of the updated system is then calculated. If the new fitness value is worse than the previous value, the refactoring is discarded and another one is tried. If the new fitness value is better than the previous, the refactoring is added to the current sequence of refactorings, and applied to the current system to form the base for the next iteration. A cut-off point is set for checking neighbors before concluding that a local maximum is reached. The end result of the search is a sequence of refactorings and a list of the before and after values of the various metrics involved in the search.

Two metrics are used to measure the quality: coupling and standard deviation of methods per class (SDMPC). Coupling (CBO) is from Briand's [2000] catalogue. The second metric, SDMPC, is used to act as a 'counter metric' for coupling. An arbitrary combination of the metrics is used, the fitness function being $SDMPC * CBO$. The new fitness function improves the CBO value of the refactored system while also improving the SDMPC of the system. All the points on a Pareto front are, in isolation, considered equivalently good. In such cases, it might be that the user may prefer some of the Pareto optimal points over others.

The concept of a Pareto front is argued to make as much sense with subsets of data as it does for complete sets. Harman and Tratt also stress the importance of knowing how many runs a search-based refactoring system will need to achieve a reasonable Pareto front approximation. Furthermore, developers are free to execute extra runs of the system if they feel they have not yet achieved points of sufficient quality on the front approximation. Pareto optimality allows determining whether one fitness function is subsumed by another: broadly speaking, if fitness function f produces data which, when merged with the data produced from function f' , contributes no points to the Pareto front then we know that f is subsumed by f' . Although it may not be immediately apparent, Pareto optimality confers a benefit potentially more useful than simply determining whether one fitness function is subsumed by another. If two fitness functions generate different Pareto optimal points, then they can naturally be combined to single front. Pareto optimality is shown to have many benefits for search-based refactoring, as it lessens the need for "perfect" fitness functions. This would make Pareto optimality an approach that should be considered for any optimization problem with conflicting goals.

5.3 Summarizing remarks

The approaches to search-based refactorings can be divided into the following groups: refactoring the program at class level, refactoring the program at procedure level, and

refactoring pieces of code. The most studies have been performed on refactoring at class level, and they are all quite similar, and actually end up using the same operations for the search algorithm. For the other aspects only one or two studies have been made, and this suggests that there is much room for competing approaches. The most advanced results have been achieved with refactorings at class level, while studies in program transformations have achieved both good and not so good results.

When examining the refactoring problems, one notable characteristic is that Seng et al. [2006] attempt to preserve building blocks from the very beginning, and several other studies have later built on the operators introduced by them. The mutation selection by Seng et al. [2006] also appears popular. The complexity of the refactoring problem at class level was most pointedly demonstrated by O’Keeffe and Ó Cinneide [2006; 2008a], who had a list of 14 mutations and 11 metrics, and Quaum and Heckle [2009], who had to pre-define mutations according to the specific system. Considering that there can be even more general refactorings in addition to those presented by O’Keeffe and Ó Cinneide, and that they could be combined with system specific mutations, the search space for an optimal refactoring sequence will soon become incredibly large.

The approaches to search based refactoring also seem advanced in the aspect that there have already been several studies that compare different search algorithms and fitness functions. As for the search algorithms, different hill climbing applications are clearly very efficient and able to produce high quality results. Interestingly, simulated annealing has been outperformed by other algorithms, although one might argue that it is more “sophisticated” than at least the basic hill climbing. All in all, there are very few approaches that use simulated annealing, and no breaking results have been achieved with it. The studies in fitness functions further support the notion of complexity in this problem area. There have O’Keeffe and Ó Cinnéide [2004] have considered the problem of finding an appropriate fitness function so important that they have developed a heuristic for balancing different weights, and Harman et al. [2007] have introduced the Pareto optimality concept to this field, as software design is indeed an area where trade-offs and compromises need to be made. As for the other studies, the variety of metrics quality evaluators shows that a refined method for deciding on an appropriate fitness function is truly needed. The only area where consensus can be found is program transformations, where quality can quite simply be measured in terms of run time and correctness or size of the program.

Table 4. Research approaches in search-based software refactoring

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|---------------------------------------|---|---|--------------------------------------|---|--|--|--|--|
| Seng et al. [2005] | Optimizing subsystem decomposition for maintenance | Model of system as a graph, extracted from source code | Genes represent subsystem candidates | Split&join, elimination and adoption | Two children from two parents, integrating crossover | Cohesion, coupling, complexity, bottlenecks and cycles | Source code extracted from resulting model | |
| Seng et al. [2006] | Refactoring a software system with a wide set of operations | Model of system, extracted from source code, with access chains | Ordered list of refactorings | Common class structure refactorings, the list is extended with a suggested transformation | Minimize rejected, duplicated and unused methods and featureless classes and maximize abstract classes | Refactored software system | SA used as search algorithm, introducing a heuristic for weighting conflicting quality goals | |
| O'Keeffe and Ó Cinnéide [2004] | Automating software refactoring | Software system | N/A | Restructure class hierarchy and method moves, mutations in counter-pairs in order to reverse a move | N/A | Minimize rejected, duplicated and unused methods and featureless classes and maximize abstract classes | Refactored software system | SA used as search algorithm, introducing a heuristic for weighting conflicting quality goals |
| O'Keeffe and Ó Cinnéide [2006; 2008a] | Automating software refactoring | System as Java source code | N/A | Refactorings regarding visibility, class hierarchy and method placement | N/A | Reusability, flexibility and understandability | Refactored code and design improvement report | Three variations of hill climbing and SA used as search algorithms |

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|---------------------------------------|--|----------------------------|--|---|---|---|---|---|
| O'Keeffe and Ó Cinnéide [2007; 2008b] | Comparison between different search techniques | System as Java source code | Ordered list of refactorings [Seng et al., 2006] | Common class structure refactorings, the list is extended with a suggested transformation [Seng et al., 2006] | A random set of transformations from one parent chosen, the transformations of the other added to that list [Seng et al., 2006] | Understandability | Refactored code and design improvement report | GA and multiple ascent hill climb implemented |
| Qayum and Heckel [2009] | Refactoring graph structure | Class diagram | N/A | A set of refactorings defined for each individual problem | N/A | Partial fitness evaluations, cost and quality | A sequence of refactorings | ACO used as search algorithm |
| Jiang et al. [2008a] | Locating dependence structures with slicing, comparing different search techniques | Source code | Two-dimensional bit matrix | A random bit flip to offspring | Multi-point crossover | Coverage and Overlap, which is divided to average and maximum | Optimal set of program slices | HC, GA, Random search and Greedy algorithm implemented; fitness function is used as a parameter |
| Jiang et al. [2008b] | Splitting procedures | Source code | Two-dimensional bit matrix | Change bit | N/A | Overlap | Optimal set of procedure slices without overlap | Greedy algorithm used |

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|-------------------------|---|-----------------|--|---|--|--|--|-----------------------------------|
| Fatiregun et al. [2004] | Program refactoring on source code level | Source code | Integer vector containing transformation numbers | Standard | Standard one-point | Size of source code (LOC) | A sequence of program transformations | Random search, HC and GA are used |
| Williams [1998] | Program parallelization | Source code | Two alternate encodings: GT includes a three symbol abbreviation of transformation and a loop number, GS includes an encoded statement | Applying one of 6 transformations, changing transformation or loop number | One-point and two-point, individual crossover points | Execution time | Transformed program | HC, GA, SA and ES implemented |
| Ryan and Ivan [1999] | Program parallelization | Source code | Tree structure of transformations | Applying atom or loop level transformation | N/A | Execution time, correctness, loop transformation success | Transformed program, transformation sequence | GP used as algorithm |
| Harman and Tratt [2007] | Pareto optimality used for multi-objective optimization | Software system | N/A | Move method | N/A | Coupling and standard deviation of methods per class | A sequence of refactorings | HC used as search algorithm |

6. SOFTWARE QUALITY

Software quality assessment has become an increasingly important field. The complexity caused by object-oriented methods makes the task more important and more difficult. An ideal quality predictive model can be seen as the mixture of two types of knowledge: common knowledge of the domain and context specific knowledge. In existing models, one of the two types is often missing. During its operating time, a software system undergoes various changes triggered by error detection, evolution in the requirements or environment changes. As a result, the behavior of the software gradually deteriorates as modifications increase. This quality slump may go as far as the entire software becoming unpredictable.

Software quality is a special concern when automatically designing software systems, as the quality needs to be measured with metrics and in pure numerical values. The use of metrics may even be argued, as they cannot possibly contain all the knowledge that an experienced human designer has. Sahraoui et al. [2000] have investigated whether some object-oriented metrics can be used as an indicator for automatically detecting situations where a particular transformation can be applied to improve the quality of a system. The detection process is based on analyzing the impact of various transformations on these object-oriented metrics using quality estimation models.

Sahraoui et al. have constructed a tool which, based on estimations on a given design, suggests particular transformations that can be automatically applied in order to improve the quality as estimated by the metrics. Roughly speaking, building a quality estimation model consists of establishing a relation of cause and effect between two types of software characteristics. Firstly, internal attributes which are directly measurable, such as size, inheritance and coupling, and secondly, quality characteristics which are measurable after a certain time of use such as maintainability, reliability and reusability. To study the impact of the global transformations on the metrics, first the impact of each elementary transformation is studied and then the global impact is derived. A case study is used for the particular case of the diagnosis of bad maintainability by using the values of metrics for coupling and inheritance as symptoms. Based on the results of this study, Sahraoui et al. argue that using metrics is a step toward the automation of quality improvement, but that experiments also show that a prescription cannot be executed without a validation of a designer/programmer.

The use of evolution metrics for fitness functions has especially been studied [Mens and Demeyer 2001; Harman and Clarke, 2004]. If one looks at the whole process of

detecting flaws and correcting them, metrics can help automating a large part of it. However, the results of the experiments show that a prescription cannot be executed without a validation of a designer or programmer. This approach cannot capture all the context of an application to allow full automation.

Some approaches regarding software quality have also been made with search-based techniques. Bouktif et al. [2002; 2004] aim at predicting software quality of object-oriented systems with GAs, and Vivanco and Jin [2007] have implemented a GA to identify possible problematic software components. Bouktif et al. [2006] have also implemented a SA to combine different quality prediction models. Summarizing remarks are presented in the end, and the fundamentals of each approach are collected in Table 5.

6.1 Search-based approaches

Bouktif et al. [2002; 2004] study the prediction of stability at object-oriented class level and propose two GA based approaches to solve the problem of quality predictive models: the first approach combines two rule sets and the second one adapts an existing rule set. The predictive model will take the form of a function that receives as input a set of structural metrics and an estimation of stress, and produces as output a binary estimation of the stability. Here, stress represents the estimated percentage of added methods in a class between two consecutive versions.

The model encoding for the GA that combines rule sets is based on a decision tree. The decision tree is a complete binary tree where each inner node represents a yes-or-no question, each edge is labeled by one of the answers, and terminal nodes contain one of the classification labels from a predetermined set. The decision making process starts at the root of the tree. When the questions at the inner nodes are of form “*Is $x > a$?*”, the decision regions of the tree can be represented as a set of isothetic boxes in an n -dimensional space (n = number of metrics). For the GA representation, these boxes are enumerated in a vector. Each gene is a (box, label) pair, and a vector of these pairs is the chromosome. The complexity of quality as a concept is directly shown in the complexity of the encoding. No simple integer vector can be used to represent quality estimations. An interesting research question is to determine what is the minimal information needed in order to evaluate or predict quality.

Mutation is a random change in the genes that happens with a small probability. In this problem, the mutation operator randomly changes the label of a box. To obtain an offspring, a random subset of boxes from one parent is selected and added to the set of

boxes of the second parent. The size of the random subset is ν times the number of boxes of the parent where ν is a parameter of the algorithm. By keeping all the boxes of one of the parents, completeness of the offspring is automatically ensured. To guarantee consistency, the added boxes are made predominant (the added boxes are “laid over” the original boxes). A level of predominance is added as an extra element to the genes. Each gene is now a three-tuple (box, label, level). The boxes of the initial population have level 1. Each time a predominant box is added to a chromosome, its level is set to 1 plus the maximum level in the hosting chromosome. To find the label of an input vector x (a software element), first all the boxes containing x are found, and x is assigned the label of the box that have the highest level of predominance.

To measure the fitness a correctness function is used; the function calculates the number of cases that the rule correctly classifies divided by the total number of cases that the rule classifies. The correctness function is defined as $C = 1 - \text{training error}$. By using the training error for measuring the fitness, it is found that the GA tended to “neglect” unstable classes. To give more weight to data points with minority labels, Youden’s [1961] J -index is used. Intuitively, the J -label is the average correctness per label. If one has the same number of points for each label, then $J = C$. As seen, the actual fitness evaluations for quality seem simple, which is surprising when compared to the complicated metric combinations used to evaluate quality in all the various GA implementations already presented. However, here the most work is needed for defining the rules that need to be satisfied and questions that need to be answered.

With a GA for adapting a rule set, an existing rule set is used as the initial population of chromosomes, each rule of the rule set being a chromosome and each condition in the rule as well as the classification label being a gene. Each chromosome is attributed a fitness value, which is $C * t$, where t is the fraction of cases that the rule classifies in the training set. The weight t allows giving rules that cover a large set of training cases a higher chance of being selected.

Parents for crossover are selected with roulette wheel method. A random cut point is generated for each parent, i.e., the cut-points are different for each parent. Otherwise, the operation is a traditional one-point crossover. By allowing chromosomes within a pair to be cut at different places, a wider variety is allowed with respect to the length of the chromosomes. The chromosomes are then mutated. The mutation of a gene consists of changing the value to which the attribute encoded in the gene is compared to a value chosen randomly from a predefined set of values for the attribute (or class label, in case the last gene is mutated). The new chromosomes are scanned and trimmed to get rid of

redundancy in the conditions that form the rules that they encode. Inconsistent rules are attributed a fitness value of 0 and will eventually die. A fixed population size is maintained. Elitism is performed when the population size is odd. This consists of copying one or more of the best chromosomes from one generation to the next. Before passing from one generation to another, the performance of combined rules to one rule set is evaluated.

In the experimental setting, to build experts (that simulate existing models), stress and 18 metrics (belonging to coupling, cohesion, complexity and inheritance) are used. Eleven object-oriented systems are used to “create” 40 experts. For the combining GA, the elitist strategy is used, where the entire population apart from a small number of fittest chromosomes is replaced. The test results show that the approach of combining experts can yield significantly better results than using individual models. The adaptation approach does not perform as well as the combination, although it gave a slight improvement over the initial model in one case. The authors believe that using more numerous and real experts on cleaner and less ambiguous data, the improvement will be more significant. It is quite inspiring that approach of combining experts produced the more promising results. If it can be assumed that experts in both initial populations have the same amount of knowledge, it would seem that merely adapting an expert would be a smaller task to perform than successfully combining the knowledge from two different experts. Thus the results are very positive when considering what the GA is capable of.

Bouktif et al. [2006] have continued their research by applying simulated annealing to combine experts. Their approach attempts to reuse and adapt quality predictive models, each of which is viewed as a set of expertise parts. The search then aims to find the best subset of expertise parts, which forms a model with an optimal predictive accuracy. The SA algorithm and a GA made for comparison were defined for Bayesian classifiers (BCs), i.e., probabilistic predictive models.

An optimal model is built of a set of experts, each of which is given a weight. Each individual, i.e., chunk, of expertise is presented by a tuple consisting of an interval and a set of conditional probabilities. Transitions in the neighborhood are made by changing probabilities or interval boundaries. A transition may also be made by adding or deleting a chunk of expertise. The fitness function is the correctness function.

For evaluation, the SA needs two elements as inputs: a set of existing experts and a representative sample of context data. Results show a considerable improvement in the predictive accuracy, and the results produced by the SA are stable. The values for GA and SA are so similar that the authors do not see a need to value one approach over the

other. Results also show that the accuracy of the best produced expert increases as the number of reused models increases and that good chunks of expertise can be hidden in inaccurate models. Again the results achieved with SA encourages further usage of different search algorithms apart from GA, or even combining and making more hybrid approaches in order to increase quality in search based approaches to software design.

Vivanco and Jin [2007] present initial results of using a parallel GA as a feature selection method to enhance a predictive model's ability to identify cognitively complex components in a Java application. Linear discriminant analysis (LDA) can be used as a multivariate predictive model.

It is theorized that the structural properties of modules have an impact on the cognitive complexity of the system, and further on, that modules that exhibit high cognitive complexity result in poor quality components. Again, this is in line with the assumption already made by Lutz [2001], that the simpler a design, the better. A preliminary study is carried out with a biomedical application developed in Java. Experienced program developers are asked to evaluate the system. Classes labeled as low are considered easy to understand and use, while a high ranking implied the class is difficult to fully comprehend and would likely take considerable much more effort to maintain. Source code measurements, 63 metrics for each Java class, are computed using a commercial source code inspection application. To establish a baseline, all the available metrics are used with the predictive model. The Chidamber and Kemerer [1994] metrics suite is used to determine if the model would improve. Finally, the GA is used to find alternate metrics subsets. Using the available metrics with LDA, less than half of the Java classes are properly classified as difficult to understand. The CK metrics suite performs slightly better. Using GA, the LDA predictive model has the highest performance using a subset of 32 metrics. The GA metrics correctly classify close to 100% of the low, nearly half of the medium and two thirds of the high complexity classes.

Vivanco and Jin are most interested in finding the potentially problematic classes with high cognitive complexity. A two-stage approach is evaluated. First, the low complexity classes are classified against the medium/high complexity classes. The GA driven LDA highly accurately identifies the low and medium/high complexity classes with a subset of 24 metrics. When only the medium complexity classes are compared to high complexity, a GA subset of 28 metrics results in extremely high accuracy for the medium complexity classes and in identifying the problematic classes. In all GA subsets, metrics that cover Halstead complexity, coupling, cohesion, and size are used, as well as program readability metrics such as comment to code ratios and the average length of method

names.

This study is extremely interesting as it ties known software metrics with human expertise and compares how metrics perform when trying to correctly classify objects. It is noteworthy that from 63 different metrics the optimal outcome was achieved with 24-32 metrics, which is less than half of all metrics available. Although there is naturally overlap between different metrics, it is interesting to see that many of them do not seem to correctly evaluate the program. The found metrics cohesion, coupling and complexity support the current fitness function choices to a certain points. However, many fitness functions only calculate 2-5 different metrics, while the optimum was reached with over 20. In addition, several metrics need the source code, and thus make them unsuitable for more high-level problems.

6.2 Summarizing remarks

The presented studies on software quality estimation show that correctly evaluating software is anything but easy. However, although the amount of studies is small, they are all very recent, and thus shows promise that search-based approaches can also be used in this sub-area of software design. Finding a search algorithm for quality estimation can also be seen as a developed way of tackling the problem of finding an optimal fitness function. In other words, in the future it might be possible to use a fitness function (i.e., a search algorithms) to find an optimal fitness function for each individual software design problem. Using search algorithms for quality estimations, the current fitness function, is the first step in this direction.

Table 5. Studies in search-based software quality enhancement

| Name | Approach | Input | Encoding | Mutation | Crossover | Fitness | Outcome | Comments |
|----------------------------|--|--|---|--|--|-------------|--|----------|
| Bouktif et al. [2002;2004] | Combining two rule sets vs. adapting a rule set with GA in quality prediction models | Decision tree | Combination: box, label -pairs from decision tree Adaptation: one rule is one chromosome, each condition in the rule is a gene | Combination: change of label Adaptation: change value of attribute encoding | Combination: a random set of boxes from one parent added to the other and level of predominance added to gene (box, label, level) Adaptation: standard one-point, parents selected with roulette-wheel method | Correctness | Optimal rule set | |
| Bouktif et al. [2006] | Combining software quality prediction models, i.e., experts | Set of example models and context data | Range and conditional probabilities | Modify range or probability or add or remove an expert | N/A | Correctness | Optimal model combined of sub-optimal models | SA used |
| Vivanco and Jin [2007] | Identification of complex components | Software system | N/A | N/A | N/A | OO metrics | Classes divided according to complexity levels | |

7. FUTURE WORK

From the search-based approaches presented here, software clustering and software refactoring (i.e., re-design) appear to be at the most advanced stage. Thus, most work is needed with actual architecture design, starting from requirements and not a ready-made system. Also, search-based application of, e.g., design patterns, should be investigated more. Another branch of research should be focused on quality metrics. So far the quality of a software design has mostly been measured with cohesion and coupling, which mostly conform to the quality factors of efficiency and modifiability. However, there are many more quality factors, and if an overall stable software system is desired, more factors should be taken into account in evaluation, such as reliability and stability. Also, as demonstrated with the MQ metric in Section 4, metrics that have seemed good in the beginning may prove to be inadequate when investigated further. Fortunately, it seems that most of the work presented here is the result of developing research that is still continuing. The following research questions should and could very well be answered in the foreseeable future:

- What kind of architectural decisions are feasible to do with search-based techniques?

Research with search-based software architecture design is at an early stage, and not all possible architecture styles and design patterns have been tested. Some architectural decisions are more challenging to implement automatically than others, and in some cases it may not be possible at all. The possibilities should be mapped to effectively research the extent of search-based designs capabilities.

- What is a sufficient starting point to being software architecture design with search-based technique?

So far requirements with a limited set of parameters have been used to build software architecture, or a ready system has been improved. Some design choices need very detailed information regarding the system in order to effectively evaluate the change in quality after implementing a certain design pattern or architecture style. The question of what information is needed for correct quality evaluation is not by any means easily answered.

- What would be optimal representation, crossover and mutation operators regarding the software modularization problem?

Much work has been done with software modularization, and the chromosome encoding, crossover and mutation operators vary greatly. Optimal solutions would be

interesting to find. As discussed throughout the survey, the chosen encoding significantly affects the result of mutation and crossover operations and also has a big impact on run time for the algorithm. There are also several options for crossover, where some maintain building blocks better than others.

- What would be optimal representation, crossover and mutation operators regarding the software refactoring problem?

Much research has been done with software refactoring, and the chromosome encoding, crossover and mutation operators vary greatly. Especially the set of mutations is interesting, as they define how greatly the software can be refactored. An optimal encoding might enable a larger set of mutations, thus giving the search-based algorithm a larger space to search for optimal solutions.

- What metrics could be seen as a “standard” for evaluating software quality?

The evaluation of quality, i.e., the fitness function, is a crucial part of evolutionary approaches to software engineering. Some metrics, e.g., coupling and cohesion, have been widely used to measure quality improvements at different levels of design. However, these metrics only evaluate a small portion of quality factors, and there are several versions of even some very “standard” metrics. Metrics by, e.g., Briand [2000] and Chidamber and Kemerer [1994] can be considered as some kind of standards. However, all software metrics are constantly subjected to criticism, as their correctness is challenged. Thus, by the author’s view, as there are several versions of even the most common metrics and there is no agreement that metrics even measure the right things at the moment, no metric set can currently be seen as standard. Thus, a well-validated metric set would be extremely beneficial, if it is possible to conduct such a set. It very well may be that the present metrics simply don’t suffice, and in that case other directions must be taken to evaluate quality, as has already been demonstrated in some of the work covered in this survey.

- How can metrics be grouped to achieve more comprehensible quality measures?

Metrics achieve clear values, but if a human designer would attempt to use a tool in the design process, notions such as “efficiency” and “modifiability” are more comprehensible than “coupling” and “cohesion”. Thus, being able to group sets of metrics to correspond to certain real-world quality values would be beneficial when making design tools available for common use.

8. CONCLUSIONS

This survey has presented on-going research in the sub-fields of search-based software

design. There has been much progress in the sub-fields of software modularization and refactoring, and very promising results have been achieved. A more complex problem is automatically designing software architecture from requirements, but some initial steps have already been taken in this direction as well. Figure 3 shows the timeline of the presented studies, and it very effectively demonstrates the increasing interest in the area during the very past years. There has been immense increase in the area of OO design and refactoring, while clustering, the first application in the area, has not sparked new research interest.

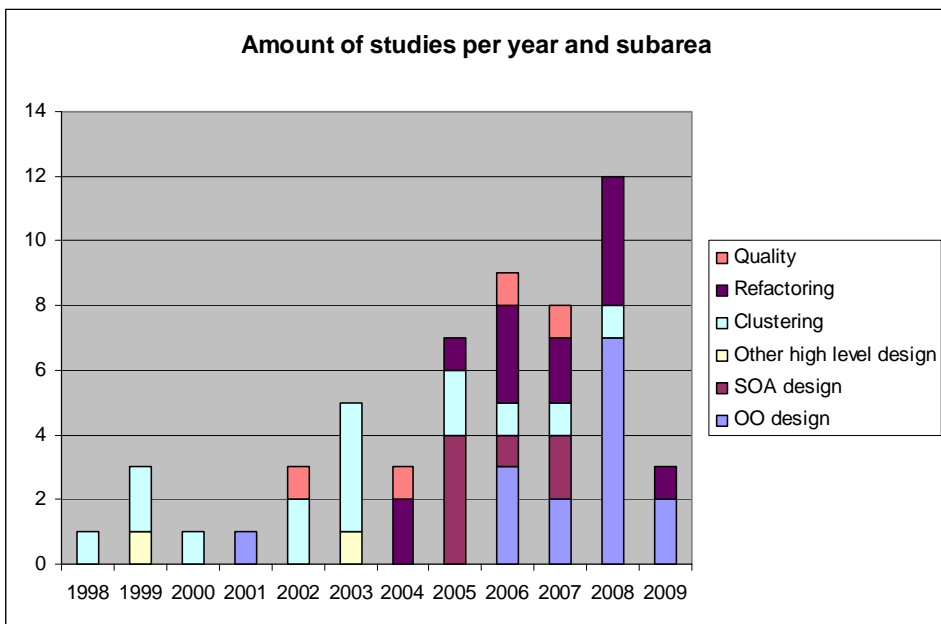


Fig. 3. Timeline for studies in search-based design

The surveyed research shows that metrics, such as cohesion and coupling can accurately evaluate some quality factors, as the achieved, automatically improved designs, have been accepted by human designers. However, many authors also report problems: the quality of results is not as high as wished or expected, and many times the blame is placed with a less than optimal encoding and crossover operators. Extensive testing of different encoding options is practically infeasible, and thus inspiration could be found in those solutions that have produced the most promising results. As a whole, software (re-)design seems to be an appropriate field for the application of meta-heuristic search algorithms, and there is much room for further research.

ACKNOWLEDGMENTS

The author would like to thank Professor Erkki Mäkinen for his helpful comments when

writing this survey. This work was partially done for the Darwin project, funded by the Academy of Finland.

REFERENCES

- AFZAL, W., TORKAR, R., AND FELDT, R. 2008. A systematic mapping study on non-functional search-based software testing. In: *Proceedings of SEKE 2008*, 488 – 493.
- AFZAL, W., TORKAR, R., AND FELDT, R. 2009. A systematic review of search-based testing for non-functional system properties. *Information and Software Technology*, **51**(6), 2009, 57 – 83.
- AHUJA, S.P. 2000. A genetic algorithm perspective to distributed systems design. In: *Proceedings of the Southeastcon 2000*, 2000, 83 – 90.
- ALANDER, J.T., MANTERE, T., AND MOGHADAMPOUR, G. 1997. Testing software response times using a genetic algorithm. In: *Proceedings of the 3rd Nordic Workshop on Genetic Algorithms and their Applications (3NWGA)*, 1997, 293 – 298.
- ALBA, E., AND CHICANO, F. 2007. Ant colony optimization for model checking, In: *Proceedings of the 11th International Conference on Computer Aided Systems Theory (EUROCAST 2007)*, 2007, 523 – 530.
- ALBA, E., AND TROYA, J.M. 1996. Genetic algorithms for protocol validation. In: *Proceedings of the 4th International Conference on Parallel Problem Solving from Nature (PPSN'96)*, 1996, 870 – 879.
- AMOUI, M., MIRARAB, S., ANSARI, S. AND LUCAS, C. 2006. A genetic algorithm approach to design evolution using design pattern transformation, *International Journal of Information Technology and Intelligent Computing* **1** (1, 2), June/ August, 2006, 235 – 245.
- ANTONIOL, G., DI PENTA, M. AND HARMAN, M. 2004. Search-based techniques for optimizing software project resource allocation. In: *Proceedings of the 2004 Conference on Genetic and Evolutionary Computation (GECCO'04)*, 2004, 1425 – 1426.
- ANTONIOL, G., DI PENTA, M. AND NETELER, M. 2003. Moving to smaller libraries via clustering and genetic algorithms. In: *Proceedings of the Seventh European Conference on Software Maintenance and Reengineering (CSMR'03)*, 2003, 307 – 316.
- ARCURI, A., AND YAO, X. 2008. A novel co-evolutionary approach to automatic software bug fixing. In: *Proceedings of the IEEE Congress on Evolutionary Computation (CEC'08)*, 2008, 162 – 168.
- BAGNALL, A.J., RAYWARD-SMITH, V.J., AND WHITLEY, I.M. 2001. The next release problem, *Information and Software Technology*, **43** (14), 2001, 883 – 890.
- BASS, L., CLEMENTS, P., AND KAZMAN, R. 1998. *Software Architecture in Practice*, Addison-Wesley, 1998.
- BODHUIN, T., DI PENTA, M., AND TROIANO, L. 2007. A search-based approach for dynamically re-packaging downloadable applications, In: *Proceedings of the Conference of the Center for Advanced Studies on Collaborative Research (CASCON'07)*, 2007, 27 – 41.
- BOUKTIF, S., AZAR, D., SAHRAOUI, H., KÉGL, B. AND PRECUP, D. 2004. Improving rule set based software quality prediction: a genetic algorithm-based approach, *Journal of Object Technology*, **3**(4), April 2004, 227 – 241.
- BOUKTIF, S., KÉGL, B. AND SAHRAOUI, H. 2002. Combining software quality predictive models: an evolutionary approach. In: *Proceedings of the International Conference on Software Maintenance (ICSM'02)*, 2002, 385 – 392.
- BOUKTIF, S., SAHRAOUI, H. AND ANTONIOL, G. 2006. Simulated annealing for improving software quality prediction, In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2006)*, 1893 – 1900.
- BOWMAN, M., BRIAND, L.C., AND LABICHE, Y. 2008. Solving the class responsibility assignment problem in object-oriented analysis with multi-objective genetic algorithms, Technical report SCE-07-02, Carleton University.
- BRIAND, L., WÜST, J., DALY, J., PORTER, V. 2000. Exploring the relationships between design measures and software quality in object oriented systems. *Journal of Systems and Software*, **51**, 2000, 245 – 273.
- BUDGEN, D. 2003. *Software Design*. Pearson, 2003.
- BUI, T.N., AND MOON, B.R. 1996. Genetic algorithm and graph partitioning, *IEEE Transactions on Computers*, **45**(7), July 1996, 841 – 855.
- CANFORA, G., DI PENTA, M., ESPOSITO, R., AND VILLANI, M.L. 2004. A lightweight approach for QoS-aware service composition. In: *Proceedings of the ICSOC 2004 – short papers*. IBM Technical Report, New York, USA.
- CANFORA, G., DI PENTA, M., ESPOSITO, R., AND VILLANI, M.L. 2005a. An approach for qoS-aware service composition based on genetic algorithms, In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO) 2005*, June 2005, 1069–1075.
- CANFORA, G., DI PENTA, M., ESPOSITO, R., AND VILLANI, M.L. 2005b. QoS-aware replanning of composite web services, In: *Proceedings of IEEE International Conference on Web Services (ICWS'05) 2005*, 2005, 121–129.
- CAO, L., LI, M. AND CAO, J. 2005a. Cost-driven web service selection using genetic algorithm, In: *LNCS 3828*, 2005, 906 – 915.

CAO, L., CAO, J., AND LI, M. 2005b. Genetic algorithm utilized in cost-reduction driven web service selection, In: *LNCS 3802*, 2005, 679 – 686.

CHAO, C., KOMADA, J., LIU, Q., MUTEJA, M., ALSALGAN, Y., AND CHANG, C. 1993. An application of genetic algorithms to software project management, In: *Proceedings of the 9th International Conference on Advanced Science and Technology*, 1993, 247 – 252.

CHE, Y., WANG, Z., AND LI, X. 2003. Optimization parameter selection by means of limited execution and genetic algorithms, In: *APPT 2003, LNCS 2834*, 2003, 226 – 235.

CHIDAMBER, S.R., AND KEMERER, C.F. 1994. A metrics suite for object oriented design. *IEEE Transactions on Software Engineering*, **20** (6), 1994, 476 – 492.

CLARK, J.A., AND JACOB, J.J. 2000. Searching for a solution: engineering tradeoffs and the evolution of provably secure protocols, In: *Proceedings of the 2000 IEEE Symposium on Security and Privacy*, 2000, 82 – 95.

CLARKE, J., DOLADO, J.J., HARMAN, M., HIERONS, R.M., JONES, B., LUMKIN, M., MITCHELL, B., MANCORIDIS, S., REES, K., ROPER, M., AND SHEPPERD, M. 2003. Reformulating software engineering as a search problem, *IEE Proceedings - Software*, **150** (3), 2003, 161 – 175.

COHEN, M.B., COLBOURN, C.J., AND LING, A.C.H. 2003. Augmenting simulated annealing to build interaction test suites, In: *Proceedings of the 14th International Symposium on Software Reliability Engineering*, 2003, 394– 405.

COOPER, K.D., SCHIELKE, P.J., AND SUBRAMANIAN, D. 1999. Optimizing for reduced code space using genetic algorithms, In: *Proceedings of the ACM SIGPLAN 1999 Workshop on Languages, Compilers and Tools for Embedded Systems (LCTES'99)*, 1999, 1 – 9.

CRAMER, N. 1985. A representation for the adaptive generation of simple sequential programs, In: *Proceedings of the International Conference on Genetic Algorithms and their Applications*, Carnegie-Mellon University, 183 – 187.

DEB, K. 1999. Evolutionary algorithms for multicriterion optimization in engineering design, In: *Proc. Evolutionary Algorithms in Engineering and Computer Science (EUROGEN'99)*, 135 – 161.

DI PENTA, M., NETELER, M., ANTONIOL, G. AND MERLO, E. 2005. A language-independent software renovation framework, *The Journal of Systems and Software*, **77**, 2005, 225 – 240.

DÍAZ, E., TUYA, J., BLANCO, R., AND DOLADO, J., 2008. A tabu search algorithm for structural testing, *Computers & Operations Research*, **35** (10), 2008, 3052 – 3072.

DORIGO, M. 1992. Optimization, Learning and Natural Algorithms, Ph.D. thesis, Dipartimento di Elettronica, Politecnico di Milano, 1992.

DOLADO, J.J., AND FERNANDEZ, L. 1998. Genetic programming, neural networks and linear programming in software project estimation, In: *Proceedings of International Conference on Software Process Improvement, Research, Education and Training (INSPIRE III)*, 1998, 157 – 171.

DOVAL, D., MANCORIDIS, S., AND MITCHELL, B.S., 1999. Automatic clustering of software systems using a genetic algorithm, In: *Proceedings of the Software Technology and Engineering Practice*, 1999, 73 – 82.

EL-FAKIH, K., YAMAGUCHI, H., AND V. BOCHMANN, G., 1999. A method and a genetic algorithm for deriving protocols for distributed applications with minimum communication cost, In: *Proceedings of the 11th International Conference on Parallel and Distributed Computing and Systems (PDCS'99)*, 1999, 863 – 868.

EVETT, M.P., KHOSHGOFTAAR, T.M., CHIEN, P-D., AND ALLEN, E.B., 1999. Using genetic programming to determine software quality, In: *Proceedings of the 12th International Florida Artificial Intelligence Research Society Conference (FLAIRS'99)*, 1999, 113 – 117.

FALKENAUER, E. 1998. *Genetic Algorithms and grouping problems*, Wiley, 1998.

FATIREGUN, D., HARMAN, M. AND HIERONS, R. 2003. Search based transformations. In: *Proceedings of the 2003 Conference on Genetic and Evolutionary Computation (GECCO'03)*, 2003, 2511 – 2512.

FATIREGUN, D., HARMAN, M. AND HIERONS, R. 2004. Evolving transformation sequences using genetic algorithms. In: *Proceedings of the 4th International Workshop on Source Code Analysis and Manipulation (SCAM 04)*, Sept. 2004, 65 – 74.

FATIREGUN, D., HARMAN, M. AND HIERONS, R. 2005. Search-based amorphous slicing. In: *Proceedings of the 124th International Working Conference on Reverse Engineering (WCRE'05)*, 2005, 3 – 12.

FELDT, R. 1998. Generating multiple diverse software versions with genetic programming. In: *Proceedings of the 24th EUROMICRO Conference*, 1998, 387 – 394.

FERGUSON, R., AND KOREL, B. 1995. Software test data generation using the chaining approach. In: *Proceedings of the IEEE International Test Conference on Driving Down the Cost of Test*, 1995, 703 – 709.

FISCHER, K.F. 1977. A test case selection method for the validation of software maintenance modifications, In: *Proceedings of International Computer Software and Applications Conference (COMPSAC'77)*, 1977, 421-426.

FISCHER, K.F., RAJI, F., AND CHRUSCICKI, A. 1981. A methodology for retesting modified software, In *Proceedings of National Telecommunications Conference (NTC'81)*, 1981, 1-6.

FONSECA, C., AND FLEMING, P. 1995. An overview of evolutionary algorithms in multi-objective optimization, *Evolutionary Computation*, **3**(1), 1995, 1-16.

GAMMA, E., HELM, R., JOHNSON, R. AND VLISSIDES, J. 1995. *Design Patterns, Elements of Reusable Object-Oriented Software*. Addison-Wesley, 1995.

GARFINKEL, R., AND NEMHAUSER, G.L. 1972. *Integer Programming*, John Wiley and Sons, 1972.

GLOVER, F. 1986. Future paths for Integer Programming and Links to Artificial Intelligence, *Computers and Operations Research*, **5**, 1986, 533 – 549.

- GLOVER, F. W. AND KOCHENBERGER, G.A. (eds.) 2003. *Handbook of Metaheuristics*, International Series in Operations Research & Management Science **57**, Springer, 2003.
- GOLD, N. 2001. Hypothesis-based concept assignment to support software maintenance. In: *Proceedings of the 22nd International Conference on Software Maintenance (ICSM 01)*, 2001, 545 – 548.
- GOLD, N., HARMAN, M., LI AND, Z., AND MAHDAVI, K. 2006. A search based approach to overlapping concept boundaries. In: *Proceedings of the 22nd International Conference on Software Maintenance (ICSM 06)*, USA Sept. 2006, 310 – 319.
- GOLDSBY, H. AND CHENG, B.H.C. 2008. Avida-mde: a digital evolution approach to generating models of adaptive software behavior. In: *Proceedings of the Genetic Evolutionary Computation Conference (GECCO 2008)*, 2008, 1751 – 1758.
- GOLDSBY, H., CHANG, B.H.C., MCKINLEY, P.K., KNOESTER, D., AND OFRIA, C.A. 2008. Digital evolution of behavioral models for autonomic systems. In: *Proceedings of 2008 International Conference on Autonomic Computing*, 2008, 87 – 96.
- HARMAN, M. 2007. The current state and future of search based software engineering, In: *Proceedings of the 2007 Future of Software Engineering (FOSE'07)*, 342 – 357.
- HARMAN, M., HIERONS, R. AND PROCTOR, M. 2002. A new representation and crossover operator for search-based optimization of software modularization. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2002)*, July 2002, 1351–1358.
- HARMAN, M., HU, L., HIERONS, R., WEGENER, J., STHAMER, H., BARESEL, A., AND ROPER, M. 2004. Testability transformation. *IEEE Transactions on Software Engineering*, **30**(1), January 2004, 3-16.
- HARMAN, M., ANSOURI, S.A. AND ZHANG, J. 2009. Search based software engineering: a comprehensive review. Technical report TR-09-03, King's College, London, United Kingdom, 2009.
- HARMAN, M., SWIFT, S. AND MAHDAVI, K. 2005. An empirical study of the robustness of two module clustering fitness functions. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2005)*, USA, June 2005, 1029 – 1036.
- HARMAN, M. AND JONES, B.F. 2001. Search based software engineering. *Information and Software Technology* 2001, **43**(14), 833 – 839.
- HARMAN, M., AND CLARK, J. 2004. Metrics are fitness functions too. In: *10th International Software Metrics Symposium (METRICS 2004)*, USA Sept. 2004, 58 – 69.
- HARMAN, M. AND TRATT, L. 2007. Pareto optimal search based refactoring at the design level, In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO'07)*, 2007, 1106 – 1113.
- HARMAN, M. AND WEGENER, J. 2004. Getting results with search-based approaches to software engineering. In: *Proceedings of the ICSE 2004*, 728 – 729.
- HARTMANN, J. AND ROBSON, D.J. 1989. Revalidation during the software maintenance phase. In: *Proceedings of the 1989 Conference on Software Maintenance*, IEEE Press, 1981, 70 – 80.
- HOLLAND, J.H. 1975. *Adaption in Natural and Artificial Systems*. MIT Press, Ann Arbor, 1975.
- HORWITZ, S., REPS, T., AND BINKLEY, D.W. 1988. Interprocedural slicing using dependence graphs. In: *Proceedings of ACM SIGPLAN Notices*, **25** (6), 1988, 35– 46.
- HUHNS, M., AND SINGH, M. 2005. Service-oriented computing: Key concepts and principals. *IEEE Internet Computing*, Jan-Feb 2005, 75 – 81.
- HUYNH, S. AND CAL, Y. 2000. An Evolutionary approach to software modularity analysis, In: *Proceedings of the First international workshop on Assessment of Contemporary Modularization Techniques ACoM'07, ICSE Workshops*, May 2007, 1 – 6.
- JAEGER, M.C. AND MÜHL, G. 2007. QoS-based selection of services: the implementation of a genetic algorithm, In: T. Braun, G. Carle and B. Stiller (Eds.): *Kommunikation in Verteilten Systemen (KiVS) 2007 Workshop: Service-Oriented Architectures und Service-Oriented Computing*, VDE Verlag, March 2007, 359 – 371.
- JIANG, T., GOLD, N., HARMAN, M. AND ZHENG, L. 2008a. Locating dependence structures using search-based slicing, *Information and Software Technology*, **50**, 2008, 1189 – 1209.
- JIANG, T., HARMAN, M. AND HASSOUN, Y. 2008b. Analysis of Procedure Splittability, In: *Proceedings of the 15th Workshop on Reverse Engineering*, 2008, 247 – 256.
- JOHNSON, C. 2007. Genetic programming with fitness based on model checking, In: *Proceedings of the 10th European Conference on Genetic Programming*, 2007, 114 – 124.
- JONES, B., STHAMER, H-H. AND EYRES, D.E. 1996. Automatic structural testing using genetic algorithms, *Software Engineering Journal*, **11** (5), 1996, 299 – 306.
- KAUFMAN, L. AND ROUSSEEUW, P. 1990. *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley, 1990.
- KENNEDY, J. AND EBERHART, R.C. 1995. Particle swarm optimization, In: *Proceedings of the IEEE International Conference on Neural Networks*, 1995, 1942 – 1948.
- KESSANTINI, M., SAHRAOUL, H. AND BOUKADOU, M. 2008. Model transformation as an optimization problem, In: *Proceedings of the ACM/IEEE 11th International Conference on Model Driven Engineering Languages and Systems (MODELS'08)*, 2008, 159 – 173.
- KIM, D., AND PARK, S. 2009. Dynamic architectural selection: a genetic algorithm approach, In: *Proceedings of the 1st Symposium on Search-Based Software Engineering*, 2009, 59 – 68.
- KIRKPATRICK, S., GELATT, C., AND VECCHI, M. 1983. Optimization by simulated annealing, *Science*, **220**, 1983,

- KOZA, J.R. 1992. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*, MIT Press, 1992.
- LANGER, R., AND MANCORIDIS, S. 2007. Using code metric histograms and genetic algorithms to perform author identification for software forensics. In: *Proceedings of the 9th Annual Conference on Genetic and Evolutionary Computation (GECCO'07)*, 2007, 2082– 2089.
- LI, H., AND LAM, C.P. 2005. An ant colony optimization approach to test sequence generation for statebased software testing. In: *Proceedings of the 5th International Conference on Quality Software (QSIC'05)*, 2005, 255– 264.
- LOSAVIO, F., CHIRINOS, L., MATTEO, A., LÉVY, N., AND RAMDANE-CHERIF, A. 2004. ISO quality standards for measuring architectures. *The Journal of Systems and Software*, **72**, 2004, 209 – 223.
- LUTZ, R. 2001. Evolving good hierarchical decompositions of complex systems, *Journal of Systems Architecture*, **47**, 2001, 613 – 634
- MAHDAVI, K., HARMAN, M., AND HIERONS, R. 2003a. A multiple hill climbing approach to software module clustering, In: *Proceedings of ICSM 2003*, 315 – 324.
- MAHDAVI, K., HARMAN, M. AND HIERONS, R. 2003b. Finding building blocks for software clustering In: *LNCS 2724*, 2003, 2513 – 2514.
- MANCORIDIS, S. MITCHELL, B.S., CHEN, Y.-F., AND GANSNER, E.R. 1999. Bunch: A clustering tool for the recovery and maintenance of software system structures. In: *Proceedings of the IEEE International Conference on Software Maintenance*, 1999, 50 – 59.
- MANCORIDIS, S., MITCHELL, B.S., RORRES, C., CHEN, Y.-F. AND GANSNER, E.R. 1998. Using automatic clustering to produce high-level system organizations of source code. In *Proceedings of the International Workshop on Program Comprehension (IWPC'98)*, USA, 1998, 45 – 53.
- MANSOUR, N. AND EL-FAKIH, K. 1999. Simulated annealing and genetic algorithms for optimal regression testing. *Journal of Software Maintenance: Research and Practice* **11**, (1), 1999, 19-34.
- MANTERE, T. AND ALANDER, J.T. 2005. Evolutionary software engineering: a review. *Applied Soft Computing* **5**, (3), 2005, 315-331.
- MARTIN, R.C. 2000. Design Principles and Design Patterns, available at <http://www.objectmentor.com>.
- MCMINN, P. 2004. Search-based software test data generation: a survey. *Software Testing, Verification and Reliability*, **14**(2), 105 – 56.
- MENS, T., AND DEMEYER, S., 2001. Future trends in evolution metrics, In: *Proceeding of the International Workshop on Principles of Software Evolution*, 2001, 83 – 86.
- MICHALEWICZ, Z. 1992. *Genetic Algorithms + Data Structures = Evolutionary Programs*. Springer-Verlag, 1992.
- MILLER, W., AND SPOONER, D.L., 1976. Automatic generation of floating-point test data, *IEEE Transactions on Software Engineering* **2**(3), 1976, 223 – 226.
- MINOHARA, T., AND TOHMA, Y. 1995. Parameter estimation of hyper-geometric distribution software reliability growth model by genetic algorithms, In: *Proceedings of the 6th Symposium on Software Reliability Engineering*, 1995, 324 – 329.
- MITCHELL, B. 2002. A Heuristic Search Approach to Solving the Software Clustering Problem. Ph. D. Thesis, Drexel University, Philadelphia, January 2002.
- MITCHELL, B.S., AND MANCORIDIS, S. 2002. Using heuristic search techniques to extract design abstractions from source code. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2002)*, USA, July 2002, 1375 – 1382.
- MITCHELL, B.S., AND MANCORIDIS, S. 2003. Modeling the search landscape of metaheuristic software clustering algorithms, In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO 2003)*, 2003, 2499 – 2510.
- MITCHELL, B.S. AND MANCORIDIS, S. 2006. On the automatic modularization of software systems using the Bunch tool, *IEEE Transactions on Software Engineering*, **32** (3), March 2006, 193 – 208.
- MITCHELL, B.S., AND MANCORIDIS, S. 2008. On the evaluation of the Bunch search-based software modularization algorithm, *Soft Computing*, **12**(1), 2008, 77 – 93.
- MITCHELL, B.S., MANCORIDIS, S., AND TRAVERSO, M. 2000. Search based reverse engineering, In: *Proceedings of the 14th International Conference on Software Engineering and Knowledge Engineering (SEKE'02)* 2002, 431 – 438.
- MITCHELL, B.S., MANCORIDIS, S., AND TRAVERSO, M. 2004. Using interconnection style rules to infer software architecture relations, In: *Proceedings of the 2004 Conference Genetic and Evolutionary Computation (GECCO'04)*, 2004, 1375 – 1387.
- MITCHELL, M. 1996. *An Introduction to Genetic Algorithms*. MIT Press, 1996.
- MONNIER, Y., BEAUVAIS, J-P., AND DÉPLANCHE, A-M. 1998. A genetic algorithm for, scheduling tasks in real-time distributed system In: *Proceedings of the 24th EUROMICRO Conference (EUROMICRO'98)*, 1998, 20708 – 20714.
- NAEIMI, H., AND DEHON, A.. 2004. A greedy algorithm for tolerating defective crosspoints in NanoPLA design, In: *Proceedings of the International Conference on Field-Programmable Technology (ICFPT2004)*, 2004, 49 – 56.
- NISBET, A.1998. GAPS: a compiler framework for genetic algorithm (GA) optimised parallelisation, In:

Proceedings of the International Conference and Exhibition on High-Performance Computing and Networking (HPCN'98), 1998, 987 – 989.

O'KEEFFE, M., AND Ó CINNÉIDE, M. 2004. Towards automated design improvements through combinatorial optimization. In: *Workshop on Directions in Software Engineering Environments (WoDISEE2004)*, W2S Workshop -26th International Conference on Software Engineering, 2004, 75 – 82.

O'KEEFFE, M., AND Ó CINNÉIDE, M. 2006. Search-based software maintenance. In: *Proceedings of CSMR 2006*, 249 – 260.

O'KEEFFE, M., AND Ó CINNÉIDE, M. 2007. Getting the most from search-based refactoring In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO'07)*, 2007, 1114 – 1120.

O'KEEFFE, M., AND Ó CINNÉIDE, M. 2008a. Search-based refactoring for software maintenance, *Journal of Systems and Software*, **81** (4), April 2008, 502 – 516.

O'KEEFFE, M., AND Ó CINNÉIDE, M. 2008b. Search-based refactoring: an empirical study, *Journal of Software Maintenance and Evolution: Research and Practice*, **20**, August 2008, 345 – 364.

PEI, M., GOODMAN, E.D., GAO, Z., AND ZHONG, K. 1994. Automated software test data generation using a genetic algorithm, available at: www.egr.msu.edu/~pei/paper/GApaper94-02.ps.

QUAUM, F., AND HECKEL, R. 2009. Local search-based refactoring as graph transformation. In: *Proceedings of the 1st Symposium on Search-Based Software Engineering*, 2009, 43-46.

RÄIHÄ, O. 2008. Genetic Synthesis of Software Architecture, University of Tampere, Department of Computer Sciences, Lic. Phil. Thesis, September 2008.

RÄIHÄ, O., KOSKIMIES, K., AND MÄKINEN, E. 2008a. Genetic synthesis of software architecture. In: *Proceedings of the 7th International Conference on Simulated Evolution and Learning (SEAL'08)*, December 2008, Melbourne, Australia. *LNCS 5361*, 565 – 574.

RÄIHÄ, O., KOSKIMIES, K., AND MÄKINEN, E. 2009. Scenario-based genetic synthesis of software architecture. In: *Proceedings of the 4th International Conference on Software Engineering Advances (ICSEA'09)*, September 2009, Porto, Portugal. IEEE Computer Society Press, to appear.

RÄIHÄ, O., KOSKIMIES, K., MÄKINEN, E., AND SYSTÄ, T. 2008b. Pattern-based genetic model refinements in MDA. In: *Proceedings of the Nordic Workshop on Model-Driven Engineering (NW-MoDE'08)*, Reykjavik, Iceland. University of Iceland, August 2008, 129 – 144, to appear in *Nordic Journal of Computing*.

RAMIREZ, A.J., KNOESTER, D.B., CHENG, B.H.C., AND MCKINLEY, P.K.. 2009. Applying genetic algorithms to decision making in autonomic computing systems. In: *Proceedings of the Nordic 6th International Conference on Autonomic Computing (ICAC'09)*, 2009, to appear.

REEVES, C.R. (ed.). 1995. *Modern Heuristic Techniques for Combinatorial Problems*. McGraw-Hill, 1995.

RYAN, C. 1999. *Automatic Re-engineering of Software Using Genetic Programming*. Kluwer Academic Publishers, 1999.

RYAN, C., AND IVAN, L. 1999. Automatic parallelization of arbitrary programs. In: *Proceedings of EUROGP'99*, *LNCS 1598*, 1999, 244-254.

SAHRAOUI, H.A., GODIN, R., AND MICELI, T. 2000. Can metrics help bridging the gap between the improvement of OO design quality and its automation? In: *Proceedings of the International Conference on Software Maintenance (ICSM '00)*, 2000, 154 – 162.

SALOMON, R. 1998. Short notes on the schema theorem and the building block hypothesis in genetic algorithms. In: *Evolutionary Programming VII*, *LNCS 1447*, 1998, 113 – 122.

SCHOEHAUER, M., AND XANTHAKIS, S. 1993. Constrained ga optimization In: *Proceedings of the 5th International Conference on Genetic Algorithms (ICGA '93)*, 1993, 573 – 580.

SENG, O., BAUYER, M., BIEHL, M., AND PACHE, G. 2005. Search-based improvement of subsystem decomposition. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO'05)*, 2005, 1045 – 1051.

SENG, O., STAMMEL, J., AND BURKHART, D. 2006. Search-based determination of refactorings for improving the class structure of object-oriented systems. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO'06)*, 2006, 1909–1916.

SHAN, Y., MCKAY, R.I., LOKAN, C.J., AND ESSAM, D.L. 2002. Software project effort estimation using genetic programming. In: *Proceedings of the 2002 IEEE International Conference on Communications, Circuits and Systems and West Sino Expositions*, 2002, 1108 – 1112.

SHANNON, C.E. 1948. The mathematical theory of communications. *Bell System Technical Journal*, 27 (379 – 423), 623 – 656.

SHAW, M., AND GARLAN, D.1996. *Software Architecture - Perspectives on an Emerging Discipline*. Prentice Hall, 1996.

SHAZELY, S., BARAKA, H., AND ABDEL-WAHAB, A. 1998. Solving graph partitioning problem using genetic algorithms. In: *Midwest Symposium on Circuits and Systems*, 1998, 302 – 305.

SHYANG, W., LAKOS, C., MICHALEWICZ, Z., AND SCHELLENBERG, S. 2008. Experiments in applying evolutionary algorithms to software verification. In: *Proceedings of the IEEE Congress on Evolutionary Computation (CEC'08)*, 2008, 3531 – 3536.

SINCLAIR, M. C., AND SHAMI, S.H. 1997. Evolving simple software agents: comparing genetic algorithm and genetic programming performance. In: *Proceedings of the 2nd International Conference on Genetic Algorithms in Engineering Systems: Innovations and Applications (GALESIA'97)*, 1997, 421 – 426.

SIMONS, C. L., AND PARMEE, I.C. 2007a. Single and multi-objective genetic operators in object-oriented

conceptual software design. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO'07)*, 2007 1957 – 1958.

SIMONS, C.L., AND PARMEE, I.C. 2007b. A cross-disciplinary technology transfer for search-based evolutionary computing: from engineering design to software engineering design, *Engineering Optimization*, **39** (5) 2007, 631 – 648.

SIMONS, C.L., AND PARMEE, I.C. 2008. User-centered, evolutionary search in conceptual software design, In: *Proceedings of the IEEE Congress on Evolutionary Computation (CEC'08)*, 2008, 869 – 876.

SIVANANDAM, S.N. AND, DEEPA, S.N. *Introduction to Genetic Algorithms*. Springer, 2007

SSBSE, 2009, <http://www.ssbse.org>, checked 17.6.2009.

SU, S., ZHANG, C., AND CHEN, J. 2007. An improved genetic algorithm for web services selection, In: *LNCS* **4531**, 2007, 284 – 295.

TRACEY, N., CLARK, J., AND MANDER., K. 1998. Automated program flaw findign using simulated annealing. In: *Proceedings of the 1998 ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'98)*, 1998, 73 – 81.

TUCKER, A., SWIFT, S., AND LIU, X.. 2001. Grouping multivariate time series via correlation. *IEEE Transactions on Systems, Man, and Cybernetics. Part B: Cybernetics*, **31**(2), 2001, 235 – 245.

TZERPOS V., HOLT R.C., MoJo: A distance metric for software clusterings. In: *Proc. of IEEE Working Conference on Reverse Engineering*, 1999, 187 – 195.

VIVANCO, R.A., AND JIN, D. 2007. Selecting object-oriented source code metrics to improve predictive models using a parallel genetic algorithm In: *Proceedings of OOPSLA'07*, 2007, 769 – 770.

WADEKAR, S. AND GOKHALE, S. 1999. Exploring cost and reliability tradeoffs in architectural alternatives using a genetic algorithm, In: *Proceedings of the 10th International Symposium on Software Reliability Engineering*, 1999, 104 – 113.

WILLIAMS, K. 1998. Evolutionary algorithms for automatic parallelization. Ph.D. thesis, Department of Computer Science, University of Reading, UK. 1998.

WIRFS-BROCK, R.J. AND JOHNSON, R.E. 1990. Surveying current research in object-oriented design, *Communications of the ACM* **33**(9), 1990, 104 – 124.

XANTHAKIS, S., ELLIS, C., SKOURLAS, C., LE GALL, A., KATSIKAS, S., AND KARAPOULIS, S. 1992. Application of genetic algorithm to software testing, In: *Proceedings of the 5th International Conference on Software Engineering and Applications*, 1992, 625 – 636.

YANG, L. JONES, B.F., AND YANG, S.-H. 2006. Genetic algorithm based software integration with minimum software risk, *Information and Software Technology* **48**(3), 2006, 133 – 141.

YAU, S.S. AND TSAI, J.-P. 1986. A survey of software design techniques, *IEEE Transactions on Software Engineering* **12**(6), 1986, 713 – 721.

YOO, S., AND HARMAN, M. 2007. Pareto efficient multi-objective test case selection, In: *Proceedings of the 2007 International Symposium on Software Testing and Analysis (ISSTA'07)* , 2007, 140-150.

YOUNDEN, W.J. 1961. How to evaluate accuracy. In: *Materials Research and Standards, ASTM*, 1961.

ZHANG, C., SU, S., AND CHEN, J. 2006. A novel genetic algorithm for qos-aware web services selection, In: *LNCS* **4055**, 2006, 224 – 235.

ZHANG, C., SU, S., AND CHEN, J. 2006. A novel genetic algorithm for qos-aware web services selection, In: *LNCS* **4055**, 2006, 224-235.

ZHANG, Y., HARMAN, M., AND MANSOURI, S.A. 2007. The multi-objective next release problem, In: *Proceedings of the 9th Annual Conference on Genetic and Evolutionary Computation (GECCO'07)*, 2007, 1129 – 1137.

ZHANG, Y., FINKELSTEIN, A., AND HARMAN, M. 2008. Search based requirements optimisation: existing work & challenges, In: *Proceedings of the 14th International Working Conference, Requirements Engineering: Foundation for Software Quality (REFSQ'08)*, 2008, 88 – 94.

ZLOCHIN, M., BIRATTARI, M., MEULEAU, N. AND DORIGO, M. 2004. Model-based search for combinatorial optimization: a critical survey, *Annals of Operations Research* **131**, 2004, 373 – 395.