

# High Dimensional Search-based Software Engineering: Finding Tradeoffs among 15 Objectives for Automating Software Refactoring Using NSGA-III

Wiem Mkaouer, Marouane Kessentini,  
Slim Bechikh  
University of Michigan, Dearborn, MI, USA  
firstname@umich.edu

Kalyanmoy Deb  
Michigan State University, USA  
kdeb@egr.msu.edu

Mel Ó Cinnéide  
University College Dublin, Ireland  
mel.ocinneide@ucd.ie

## ABSTRACT

There is a growing need for scalable search-based software engineering approaches that address software engineering problems where a large number of objectives are to be optimized. Software refactoring is one of these problems where a refactoring sequence is sought that optimizes several software metrics. Most of the existing refactoring work uses a large set of quality metrics to evaluate the software design after applying refactoring operations, but current search-based software engineering approaches are limited to using a maximum of five metrics. We propose for the first time a scalable search-based software engineering approach based on a newly proposed evolutionary optimization method NSGA-III where there are 15 different objectives to be optimized. In our approach, automated refactoring solutions are evaluated using a set of 15 distinct quality metrics. We evaluated this approach on seven large open source systems and found that, on average, more than 92% of code smells were corrected. Statistical analysis of our experiments over 31 runs shows that NSGA-III performed significantly better than two other many-objective techniques (IBEA and MOEA/D), a multi-objective algorithm (NSGA-II) and two mono-objective approaches, hence demonstrating that our NSGA-III approach represents the new state of the art in fully-automated refactoring.

## Categories and Subject Descriptors

D.2 [Software Engineering].

## General Terms

Algorithms, Reliability.

## Keywords

Search-based software engineering, software quality, code smells, many-objective optimization.

## 1. INTRODUCTION

Search-based software engineering (SBSE) studies the application of meta-heuristic optimization techniques to software engineering problems [1]. Once a software engineering task is framed as a search problem, by defining it in terms of solution representation, objective

function, and solution change operators, there are a multitude of search algorithms that can be applied to solve that problem. Search-based techniques are widely applied to solve software engineering problems such as in testing, modularization, refactoring, planning, etc. [15][13].

Based on a recent SBSE survey [15], the majority of existing work treats software engineering (SE) problems from a single-objective point of view, where the main goal is to maximize or minimize one objective, e.g., correctness, quality, etc. However, most SE problems are naturally complex in which many conflicting objectives need to be optimized such as model transformation, design quality improvement, test suite generation etc. The number of objectives to consider for most of software engineering problems is, in general, high (more than three objectives); such problems are termed *many-objective*. We claim that the reason that software engineering problems have not been formulated as many-objective problems is because of the challenges in constructing a many-objective solution. In this context, the use of traditional multi-objective techniques, e.g., NSGA-II [8], widely used in SBSE, is clearly not sufficient.

There is a growing need for scalable search-based software engineering approaches that address software engineering problems where a large number of objectives are to be optimized. Improving the scalability of SBSE approaches will increase their applicability in industry and real-world settings. Recent work in optimization has proposed several solution approaches to tackle many-objective optimization problems [17][30] using e.g., objective reduction, new preference ordering relations, decomposition, etc. However, these techniques have not yet been widely explored in SBSE [6]. To the best of our knowledge and based on recent SBSE surveys [15], only one work exists proposed by Abdel Salam et al. [26] that uses a many-objective approach, IBEA (Indicator-Based Evolutionary Algorithm) [31], to address the problem of software product line creation. However, the number of considered objectives is limited to 5.

Software refactoring is one of those software engineering problems where there are several objectives to be satisfied. Refactoring improves the design of a system by changing its internal structure without altering its external behavior [12], and is widely used to fix code smells. Code smells are known to have a negative impact on quality attributes such as flexibility or maintainability [4][29]. Software engineers often introduce code smells unintentionally during the initial design or during software development due to bad design decisions, ignorance or time pressure. Most of the existing refactoring work uses a set of more than five quality metrics to evaluate the quality of software design after applying refactoring operations. In this paper, we propose for the first time a scalable search-based software engineering approach based on NSGA-III [7] where there are 15 different objectives to optimize. Thus, in our approach, automated refactoring solutions will be evaluated using a set of 15 software quality metrics. NSGA-III is a very recent many-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

GECCO '14, July 12–16, 2014, Vancouver, BC, Canada.  
Copyright 2014 ACM 978-1-4503-2662-9/14/07\$15.00.  
<http://dx.doi.org/10.1145/2576768.2598366>

objective algorithm proposed by Deb et al. [7]. The basic framework remains similar to the original NSGA-II algorithm [8], with significant changes in its selection mechanism. This paper represents the first real-world application of NSGA-III and the first scalable work that supports the use of 15 objectives to address a software engineering problem.

We implemented our approach and evaluated it on seven large open source systems and found that, on average, more than 92% of code smells were corrected. The statistical analysis of our experiments over 31 runs shows that NSGA-III performed significantly better than two other many-objective techniques (IBEA and MOEA/D), a multi-objective algorithm (NSGA-II) and two mono-objective approaches [19][23].

## 2. MANY-OBJECTIVE SEARCH-BASED SOFTWARE ENGINEERING

By definition, a many-objective problem is a multi-objective one but with a high number of objectives  $M$ , i.e.,  $M > 3$ . Analytically, it could be stated as follows [5]:

$$\begin{cases} \text{Min } f(x) = [f_1(x), f_2(x), \dots, f_M(x)]^T, & M > 3 \\ g_j(x) \geq 0 & j = 1, \dots, P; \\ h_k(x) = 0 & k = 1, \dots, Q; \\ x_i^L \leq x_i \leq x_i^U & i = 1, \dots, n \end{cases}$$

where  $M$  is the number of objective functions and is *strictly greater* than 3,  $P$  is the number of inequality constraints,  $Q$  is the number of equality constraints,  $x_i^L$  and  $x_i^U$  correspond to the lower and upper bounds of the decision variable  $x_i$  (i.e.,  $i^{\text{th}}$  component of  $x$ ). A solution  $x$  satisfying the  $(P+Q)$  constraints is said to be feasible and the set of all feasible solutions defines the feasible search space denoted by  $\Omega$ .

In this formulation, we consider a minimization multi-objective problem (MOP) since maximization can be easily turned to minimization based on the duality principle. Over the two past decades, several Multi-Objective Evolutionary Algorithms (MOEAs) have been proposed with the hope to work with any number of objectives  $M$ . Unfortunately, It has been demonstrated that most MOEAs are ineffective in handling such type of problems. For example, NSGA-II [8], which is one of the most used MOEAs, compares solutions based on their non-domination ranks. Solutions with best ranks are emphasized in order to converge to the Pareto front. When  $M > 3$ , only the first rank may be assigned to every solution as almost all population individuals become non-dominated with each other. Without a variety of ranks, NSGA-II cannot keep the search pressure anymore in high dimensional objective spaces.

The difficulty faced when solving many-objective problems could be summarized as follows. Firstly, most solutions become equivalent between each other according to the Pareto dominance relation which deteriorates *dramatically* the search process ability to converge towards the Pareto front and the MOEA behaviour becomes very similar to the random search one. Secondly, a search method requires a very high number of solutions (some thousands and even more) to cover the Pareto front when the number of objectives increases. For instance, it has been shown that in order to find a good approximation of the Pareto front for problems involving 4, 5 and 7 objective functions, the number of required non-dominated solutions is about 62 500, 1 953 125 and 1 708 984 375 respectively [17]; which makes the decision making task very difficult. Thirdly, the

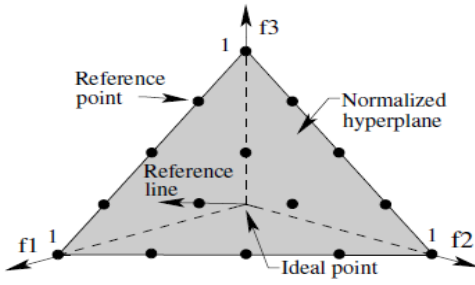
objective space dimensionality increases significantly, which makes promising search directions very hard to find. Fourthly, the diversity measure estimation becomes very computationally costly since finding the neighbors of a particular solution in high dimensional spaces is very expensive. Fifthly, recombination operators become inefficient since population members are likely to be widely distant from each other which yields to children that are not similar to their parents; thereby making the recombination operation inefficient in producing promising offspring individuals. Finally, although it is not a matter that is directly related to optimization, the Pareto front visualization becomes more complicated, therefore making the interpretation of the MOEA's results more difficult for the user.

According to a recent survey by Harman et al. [13], most software engineering problems are multi-objective by nature. However, most of existing approaches to address software engineering problems such as model transformation, design quality improvement, test suite generation, etc. are based on a mono-objective view. Multi-objective optimization techniques have been proposed in a few works [26][24] for such problems and they satisfy up to 5 objectives. However, as with any other practical domain, most software engineering problems involve optimizing more than this number of objectives. Thus, more scalable search-based software engineering approaches will be beneficial to handle rich objective spaces. We investigate, in this paper, the applicability of many-objective techniques for the software refactoring problem where up to 15 objectives are considered to evaluate refactoring suggestions.

## 3. MANY-OBJECTIVE SOFTWARE REFACTORING USING NSGA-III

### 3.1 NSGA-III

NSGA-III is a very recent many-objective algorithm proposed by Deb et al. [7]. The basic framework remains similar to the original NSGA-II algorithm [8] with significant changes in its selection mechanism. Figure 2 gives the pseudo-code of the NSGA-III procedure for a particular generation  $t$ . First, the parent population  $P_t$  (of size  $N$ ) is randomly initialized in the specified domain, and then the binary tournament selection, crossover and mutation operators are applied to create an offspring population  $Q_t$ . Thereafter, both populations are combined and sorted according to their domination level and the best  $N$  members are selected from the combined population to form the parent population for the next generation. The fundamental difference between NSGA-II and NSGA-III lies in the way the niche preservation operation is performed. Unlike NSGA-II, NSGA-III starts with a set of reference points  $Z^*$ . After non-dominated sorting, all acceptable front members and the last front  $F_l$  that could not be completely accepted are saved in a set  $S_l$ . Members in  $S_l/F_l$  are selected right away for the next generation. However, the remaining members are selected from  $F_l$  such that a desired diversity is maintained in the population. Original NSGA-II uses the crowding distance measure for selecting well-distributed set of points, however, in NSGA-III the supplied reference points ( $Z^*$ ) are used to select these remaining members (cf. Figure 1). To accomplish this, objective values and reference points are first normalized so that they have an identical range. Thereafter, orthogonal distance between a member in  $S_l$  and each of the reference lines (joining the ideal point and a reference point) is computed. The member is then associated with the reference point having the smallest orthogonal distance. Next, the niche count  $\rho$  for each reference point, defined as the number of members in  $S_l/F_l$  that are associated with the reference point, is computed for further processing. The reference point having the minimum niche count is identified and the member from the last front



**Figure 1. Normalized reference plane for a three-objective case [19].**

$F_l$  that is associated with it is included in the final population. The niche count of the identified reference point is increased by one and the procedure is repeated to fill up population  $P_{t+1}$ .

It is worth noting that a reference point may have one or more population members associated with it or need not have any population member associated with it. Let us denote this niche count as  $\rho_j$  for the  $j$ -th reference point. We now devise a new niche-preserving operation as follows. First, we identify the reference point set  $J_{\min} = \{j: \text{argmin}_j(\rho_j)\}$  having minimum  $\rho_j$ . In case of multiple such reference points, one ( $j^* \in J_{\min}$ ) is chosen at random. If  $\rho_{j^*} = 0$  (meaning that there is no associated  $P_{t+1}$  member to the reference point  $j^*$ ), two scenarios can occur. First, there exists one or more members in front  $F_l$  that are already associated with the reference point  $j^*$ . In this case, the one having the shortest perpendicular distance from the reference line is added to  $P_{t+1}$ . The count  $\rho_{j^*}$  is then incremented by one. Second, the front  $F_l$  does not have any member associated with the reference point  $j^*$ . In this case, the reference point is excluded from further consideration for the current generation. In the event of  $\rho_{j^*} \geq 1$  (meaning that already one member associated with the reference point exists), a randomly chosen member, if exists, from front  $F_l$  that is associated with the reference point  $F_l$  is added to  $P_{t+1}$ . If such a member exists, the count  $\rho_{j^*}$  is incremented by one. After  $\rho_j$  counts are updated, the procedure is repeated for a total of  $K$  times to increase the population size of  $P_{t+1}$  to  $N$ .

## 3.2 Adapting NSGA-III for the Software Refactoring Problem

### 3.2.1 Problem formulation

The refactoring problem involves searching for the best refactoring solution among the set of candidate ones, which constitutes a huge search space. A refactoring solution is a sequence of refactoring operations where the goal of applying the sequence to a software system  $S$  is typically to minimize the number of code smells in  $S$ . Usually in SBSE approaches, we use two or three metrics as objective functions for a particular multi-objective heuristic algorithm to find smells and correct them. However, in reality, there are many types of code smell and detecting the symptoms of each smell requires a particular set of metrics. Motivated by this observation, we propose in this research work to use a high number of metrics (15 metrics) where each represents a separate objective function. In this way, we obtain a many-objective (15-objective) formulation of the refactoring problem that could not be solved using standard multi-objective approaches. This formulation is given as follows:

$$\begin{aligned} & \text{Maximize } F(x, S) = [f_1(x, S), f_2(x, S), \dots, f_{15}(x, S)] \\ & \text{subject to } x = (x_1, \dots, x_n) \in X \end{aligned}$$

### NSGA-III procedure at generation $t$

**Input:**  $H$  structured reference points  $Z^s$ , parent population  $P_t$

**Output:**  $P_{t+1}$

00: **Begin**

01:  $S_t \leftarrow \emptyset, i \leftarrow 1$ ;

02:  $Q_t \leftarrow \text{Variation}(P_t)$ ;

03:  $R_t \leftarrow P_t \cup Q_t$ ;

04:  $(F_1, F_2, \dots) \leftarrow \text{Non-dominations\_Sort}(R_t)$ ;

05: **Repeat**

06:  $S_t \leftarrow S_t \cup F_i; i \leftarrow i+1$ ;

07: **Until**  $|S_t| \geq N$ ;

08:  $F_l \leftarrow F_i$ ; /\*Last front to be included\*/

09: **If**  $|S_t| = N$  **then**

10:  $P_{t+1} \leftarrow S_t$ ;

11: **Else**

12:  $P_{t+1} \leftarrow \bigcup_{j=1}^{l-1} F_j$ ;

/\*Number of points to be chosen from  $F_l$ \*/

13:  $K \leftarrow N - |P_{t+1}|$ ;

/\*Normalize objectives and create reference set  $Z^{*}$ \*/

14:  $\text{Normalize}(F^M, S_t, Z^r, Z^s)$ ;

/\*Associate each member  $s$  of  $S_t$  with a reference point\*/

/\* $\pi(s)$ : closest reference point\*/

/\* $d(s)$ : distance between  $s$  and  $\pi(s)$ \*/

15:  $[\pi(s), d(s)] \leftarrow \text{Associate}(S_t, Z^r)$ ;

/\*Compute niche count of reference point  $j \in Z^r$ \*/

16:  $\rho_j \leftarrow \sum_{s \in S_t / F_l} ((\pi(s) = j) ? 1 : 0)$ ;

/\*Choose  $K$  members one at a time from  $F_l$  to construct  $P_{t+1}$ \*/

17:  $\text{Niching}(K, \rho_j, \pi(s), d(s), Z, F_l, P_{t+1})$ ;

18: **End If**

19: **End**

**Figure 2. Pseudocode of NSGA-III main procedure.**

where  $X$  is the set of all legal refactoring sequences starting from  $S$ ,  $x_i$  is the  $i$ -th refactoring operation, and  $f_k(x, S)$  is the  $k$ -th metric. The 15 metrics under consideration will be detailed in the experimental study since our formulation is generic and applies to any software metrics.

### 3.2.2 Solution approach

**Solution representation.** As defined in the previous section, a solution consists of a sequence of  $n$  refactoring operations applied to different code elements in the source code to fix. In order to represent a candidate solution (individual/chromosome), we use a vector-based representation. Each vector's dimension represents a refactoring operation where the order of applying these refactoring operations corresponds to their positions in the vector. For each of these refactoring operations, we specify pre- and post-conditions in the style of Opdyke [12] to ensure the feasibility of their application. The initial population is generated by assigning randomly a sequence of refactorings to some code fragments. To apply a refactoring operation we need to specify which actors, i.e., code fragments, are involved/impacted by this refactoring and which roles they play in performing the refactoring operation. An actor can be a package, class, field, method, parameter, statement, or variable.

**Solution variation.** In each search algorithm, the variation operators play the key role of moving within the search space with the aim of driving the search towards optimal solutions. For crossover, we use the one-point crossover operator. It starts by selecting and splitting at random two parent solutions. Then, this operator creates two child solutions by putting, for the first child, the first part of the first parent with the second part of the second parent, and vice versa for the second child. This operator must ensure the respect of the length

limits by eliminating randomly some refactoring operations. It is important to note that in many-objective optimization, it is better to create children that are close to their parents in order to have a more efficient search process [7]. For this reason, we control the cutting point of the one-point crossover operator by restricting its position to be either belonging to the first tier of the refactoring sequence or belonging to the last tier. For mutation, we use the bit-string mutation operator that picks probabilistically one or more refactoring operations from its or their associated sequence and replaces them by other ones from the initial list of possible refactorings.

**Solution evaluation.** Each generated refactoring solution is executed on the system  $S$ . Once all required data is computed, the solution is evaluated based on the 15 metrics used as objective functions. Based on these values, the refactoring solution is assigned a non-domination rank (as in NSGA-II) and a position in the objective space allowing it to be assigned to a particular reference point based on distance calculation as previously described.

**Normalization of population members.** Usually objective functions (metrics) are incommensurable (i.e., they have different scales). For this reason, we used the normalization procedure proposed by Deb et al. [7] to circumvent this problem. At each generation, the minimal and maximal values for each metric are recorded and then used by the normalization procedure. Normalization allows the population members and with the reference points to have the same range, which is a pre-requisite for diversity preservation.

## 4. DESIGN OF THE EMPIRICAL STUDY

### 4.1 Research Questions

**RQ1:** How does NSGA-III perform compared to other many-objective (MOAE/D [30], IBEA [31]) and multi-objective (NSGA-II [8]) techniques? It is important to evaluate the performance of NSGA-III in terms of scalability when the number of considered objectives increases. In addition, it is interesting to determine if considering more metrics (objectives) improves the quality of the suggested refactoring solutions (the number of fixed code smells).

**RQ2:** How does NSGA-III perform compared to mono-objective refactoring approaches [19][23]? It is important to determine if considering each conflicting metric as a separate objective to optimize performs better than a mono-objective approach that aggregates all metrics in one objective. The comparison between a many-objective EA with a mono-objective one is not straightforward. The first one returns a set of non-dominated solutions while the second one returns a single optimal solution. In order to resolve this problem, for each many-objective algorithm we choose the nearest solution to the ideal point [3] (i.e., the vector composed of the best objective values among the population members) as a candidate solution to be compared with the single solution returned by the mono-objective algorithm. For both RQ1 and RQ2, we performed a qualitative evaluation where 8 PhD students in Software Engineering, with at least 2 years programming experience in Java and familiar with the evaluated open source systems, evaluated 10 operations from the best suggested refactoring solutions for each system. The operations were classified as useful (make sense semantically) or not.

**RQ3:** How does our many-objective formulation scale? There is a cost in allowing the developer to specify a large number of objectives. Can it be demonstrated that as the number of objectives increases, we can achieve a commensurate increase in the quality of the solutions generated? If not, then our approach is not justified.

## 4.2 Experimental Setup

### 4.2.1 Systems Studied

Our study considers the extensive evolution of different open source Java systems analyzed in the literature [19][24][21][25]. The corpus used includes releases of Apache Ant, ArgoUML, Gantt, Azureus and Xerces-J. Table 1 reports the size in terms of classes of the analyzed systems. The table also reports the number of code smells identified manually in the different systems -- more than 700 in total. Indeed, in several works [19][24][21][25], the authors asked different groups of developers to analyze the libraries to tag instances of specific code smells to validate their detection techniques. For replication purposes, they provided a corpus describing instances of different code smells including blob, spaghetti code, and functional decomposition [12]. In our study, we verified the capacity of our approach to fix classes that correspond to instances of these code smells. We used the detection rules of code smells proposed in Kessentini et al. [19] to identify the number of fixed code smells after applying the best refactoring solutions.

**Table 1. Features of software systems analyzed.**

<i>Systems</i>	<i>Number of classes</i>	<i>Number of code smells</i>
<b>ArgoUML v0.26</b>	1358	138
<b>ArgoUML v0.3</b>	1409	129
<b>Xerces v2.7</b>	991	82
<b>Ant-Apache v1.5</b>	1024	103
<b>Ant-Apache v1.7.0</b>	1839	124
<b>Gantt v1.10.2</b>	245	41
<b>Azureus v2.3.0.6</b>	1449	108

### 4.2.2 Performance Indicators

We used mainly three performance indicators to compare the different algorithms used in our experiments. These indicators are defined as follows:

- *Inverted Generational Distance (IGD)*: is used as the performance metric since it has been shown to reflect both the diversity and convergence of the obtained non-dominated solutions [7]. The IGD corresponds to the average Euclidean distance separating each reference solution from its closest non-dominated one. Note that for each system we use the set of Pareto optimal solutions generated by all algorithms over all runs as reference solutions.

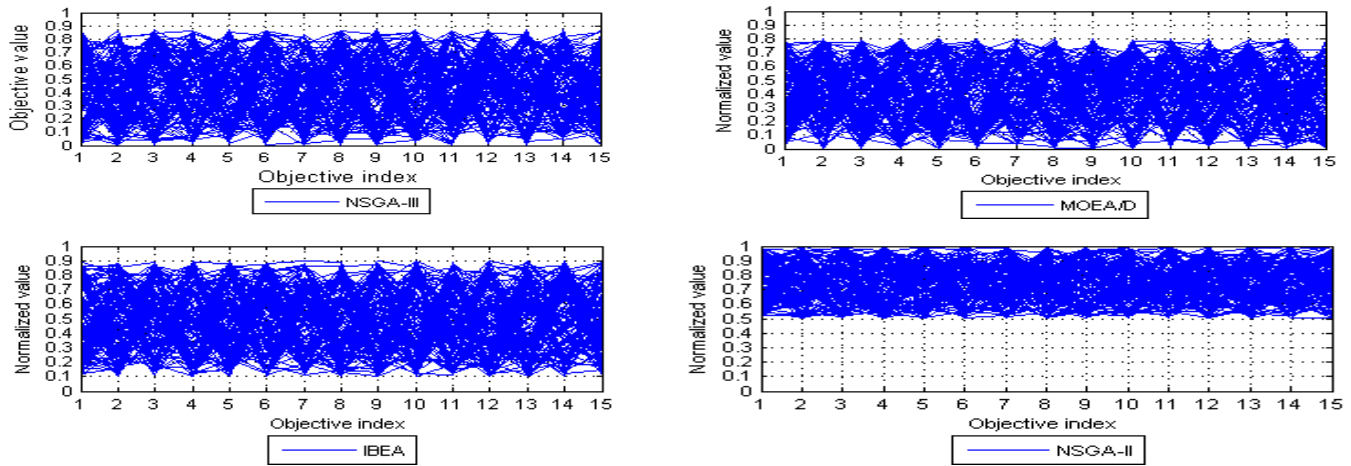
- *Percentage of fixed code smells (NF)* is the percentage of code smells fixed by the application of the best refactoring solution (i.e., number of fixed smells divided by the total number of code smells). The detection of code smells after applying a refactoring solution is performed using the detection rules of Kessentini et al. [19].

- *Usefulness of suggested refactorings (UC)* is the number of refactoring operations that make sense and useful divided by the total number of manually evaluated operations.

- *Computational time (CT)* is used mainly to compare the efficiency of NSGA-III with other algorithms using the same number of objectives.

### 4.2.3 Statistical Tests

Our experimental study is performed based on 31 independent simulation runs for each problem instance and the obtained results are statistically analyzed by using the Wilcoxon rank sum test with a 95% confidence level ( $\alpha = 5\%$ ). The latter verifies the null hypothesis  $H_0$  that the obtained results of two algorithms are samples from continuous distributions with equal medians, against the alternative that they are not  $H_1$ . The  $p$ -value of the Wilcoxon test corresponds to the probability of rejecting the null hypothesis  $H_0$  while it is true



**Figure 3. Value path plots of non-dominated solutions obtained by NSGA-III, MOEA/D, IBEA and NSGA-II during the median run of the 15-objective refactoring problem on ArgoUML v0.26.**

(type I error). A  $p$ -value that is less than or equal to  $\alpha$  ( $\leq 0.05$ ) means that we accept  $H_1$  and we reject  $H_0$ . However, a  $p$ -value that is strictly greater than  $\alpha$  ( $> 0.05$ ) means the opposite. In fact, for each problem instance, we compute the  $p$ -value obtained by comparing NSGA-II, IBEA, MOEA/D and mono-objective search results with NSGA-III ones.

#### 4.2.4 Parameter Settings

Parameter setting influences significantly the performance of a search algorithm on a particular problem [1]. For this reason, for each many-objective algorithm and for each system (cf. Table 1), we perform a set of experiments using several population sizes: 91, 210, 156, 275 and 135 for respectively 3, 5, 8, 10 and 15 objectives. The maximum number of generations used is 400, 600, 750, 1000 and 1500 respectively for 3, 5, 8, 10 and 15 objectives. Each algorithm is executed 31 times with each configuration and then comparison between the configurations is done based on *IGD* using the Wilcoxon test. In order to have significant results, for each couple (algorithm, system), we use the trial and error method [1] in order to obtain a good parameter configuration. Since we are comparing different search algorithms, we classify parameters into common parameters and specific parameters. Table 2 depicts the important common parameters. We used a set of 15 quality metrics, namely Weighted Methods per Class (WMC), Response for a Class (RFC), Lack of Cohesion of Methods (LCOM), Cyclomatic Complexity (CC), Number of Attributes (NA), Attribute Hiding Factor (AH), Method Hiding Factor (MH), Number of Lines of Code (NLC), Coupling Between Object Classes (CBO), Number of Association (NAS), Number of Classes (NC), Depth of Inheritance Tree (DIT), Polymorphism Factor (PF), Attribute Inheritance Factor (AIF) and Number of Children (NOC) [10]. We selected randomly at each run some metrics from this list when the number of objectives is lower than 15. We used 23 refactoring types in our experiments, namely Add Parameter, Rename Method Encapsulate Collection/ Downcast/ Field, Collapse Hierarchy, Hide Method, Extract Class /Interface/ Method/ Subclass/ Superclass, Inline Class/ Method, Move Field/ Method, Pull Up Field/ Method, Push Down Field/ Method and Remove Parameter/ Setting Method [31].

**Table 2. The setting of common parameters.**

Number of objectives	Number of reference points (for NSGA-III and MOEA/D)	Population size
3	91	120
5	210	230
8	156	190
10	275	280
15	135	290

### 4.3 Results

Table 3 shows the median *IGD* and *NF* values of over 31 independent runs for all algorithms under comparison. All the results were statistically significant in the 31 independent simulations using the Wilcoxon rank sum test [1] with a 99% confidence level ( $\alpha < 1\%$ ). For the 3-objective case, we see that NSGA-III and NSGA-II present similar results, and that NSGA-III provides slightly better results than IBEA and MOEA/D. For the 5-objective case, NSGA-III strictly outperforms NSGA-II and gives similar results to those of the two other multi-objective algorithms. For the 8-objective case, NSGA-III is strictly better than NSGA-II and significantly better than IBEA and MOEA/D. Additionally, IBEA seems to be slightly better than MOEA/D. It is worth noting that for problems' instances with 8 objectives or more, NSGA-II performance is dramatically degraded, which is simply denoted by the  $\sim$  symbol. For the 10- and 15-objective cases, NSGA-III is strictly better than all other algorithms. Moreover, MOEA/D seems to significantly outperform IBEA. The performance of NSGA-III could be explained by the interaction between: (1) Pareto dominance-based selection and (2) reference point-based selection, which is the distinguishing feature of NSGA-III compared to other existing many-objective algorithms. The percentage of fixed code smells using NSGA-III is better than all other algorithms in all systems in 100% of cases when more than 8 objectives are considered. It is clear from Table 3 that the percentage of fixed code smells increases as we use more quality metrics to evaluate refactoring solutions. On average, all the four algorithms NSGA-III, IBEA, MOEA/D and NSGA-II perform similarly with 3 objectives, however the percentage of fixed code smells is low in all systems. This is due to the fact that the use of only three quality metrics is not enough to evaluate the quality of the design after applying the best refactoring solution. The average percentage of fixed code smells in all systems using NSGA-III with 15 objectives on all systems is higher than 92%, which outperforms all the remaining algorithms. Thus, we can conclude that NSGA-III represents a scalable solution to find trade-offs between 15 objectives and that the use of additional objectives (metrics) improves the quality of refactoring solutions.

Figure 3 illustrates the value path plots of all algorithms regarding the 15-objective refactoring problem on ArgoUMLv0.26, the largest system used in our experiments. Similar observations were made in the remaining systems but are omitted due to space considerations. All quality metrics were normalized between 0 and 1 and all are to be minimized. We observe that NSGA-III presents the best convergence since its non-dominated solutions are the closest to the ideal point,

**Table 3. Median IGD and NF values on 31 runs (best values are in bold). ~ means a large value that is not interesting to show. The results were statistically significant on 31 independent runs using the Wilcoxon rank sum test with a 99% confidence level ( $\alpha < 1\%$ ).**

Problem	M	MaxGen	NSGA-III		IBEA		MOEA/D		NSGA-II	
			NF	IGD	NF	IGD	NF	IGD	NF	IGD
ArgoUML v0.26	3	400	<b>69%</b>	<b>1.356 x 10<sup>-3</sup></b>	67%	1.358 x 10 <sup>-3</sup>	<b>69%</b>	1.359 x 10 <sup>-3</sup>	<b>67%</b>	<b>1.356 x 10<sup>-3</sup></b>
	5	600	<b>83%</b>	<b>3.918 x 10<sup>-3</sup></b>	79%	4.001 x 10 <sup>-3</sup>	81%	4.007 x 10 <sup>-3</sup>	54%	4.423 x 10 <sup>-3</sup>
	8	750	<b>86%</b>	<b>4.115 x 10<sup>-3</sup></b>	82%	4.327 x 10 <sup>-3</sup>	84%	4.331 x 10 <sup>-3</sup>	42%	~
	10	1000	<b>89%</b>	<b>1.727 x 10<sup>-2</sup></b>	84%	2.124 x 10 <sup>-2</sup>	87%	2.001 x 10 <sup>-2</sup>	37%	~
	15	1500	<b>96%</b>	<b>3.302 x 10<sup>-2</sup></b>	81%	3.826 x 10 <sup>-2</sup>	84%	3.733 x 10 <sup>-2</sup>	34%	~
Xerces v2.7	3	400	<b>64%</b>	<b>9.751 x 10<sup>-4</sup></b>	<b>64%</b>	9.754 x 10 <sup>-4</sup>	<b>64%</b>	9.753 x 10 <sup>-4</sup>	62%	9.752 x 10 <sup>-4</sup>
	5	600	<b>78%</b>	<b>7.876 x 10<sup>-3</sup></b>	76%	7.910 x 10 <sup>-3</sup>	74%	7.912 x 10 <sup>-3</sup>	56%	8.006 x 10 <sup>-3</sup>
	8	750	<b>84%</b>	<b>8.001 x 10<sup>-3</sup></b>	81%	8.422 x 10 <sup>-3</sup>	81%	8.428 x 10 <sup>-3</sup>	48%	~
	10	1000	<b>92%</b>	<b>2.115 x 10<sup>-2</sup></b>	86%	2.410 x 10 <sup>-2</sup>	86%	2.299 x 10 <sup>-2</sup>	44%	~
	15	1500	<b>94%</b>	<b>4.666 x 10<sup>-2</sup></b>	86%	5.198 x 10 <sup>-2</sup>	84%	4.935 x 10 <sup>-2</sup>	39%	~
ArgoUML v0.3	3	400	63%	<b>2.667 x 10<sup>-3</sup></b>	63%	2.668 x 10 <sup>-3</sup>	<b>66%</b>	2.668 x 10 <sup>-3</sup>	63%	<b>2.667 x 10<sup>-3</sup></b>
	5	600	<b>79%</b>	<b>4.283 x 10<sup>-3</sup></b>	76%	4.291 x 10 <sup>-3</sup>	78%	4.294 x 10 <sup>-3</sup>	51%	4.524 x 10 <sup>-3</sup>
	8	750	<b>91%</b>	<b>5.545 x 10<sup>-3</sup></b>	84%	5.701 x 10 <sup>-3</sup>	86%	5.716 x 10 <sup>-3</sup>	47%	~
	10	1000	<b>94%</b>	<b>3.339 x 10<sup>-2</sup></b>	87%	3.601 x 10 <sup>-2</sup>	89%	3.477 x 10 <sup>-2</sup>	39%	~
	15	1500	<b>98%</b>	<b>6.001 x 10<sup>-2</sup></b>	89%	6.554 x 10 <sup>-2</sup>	87%	6.399 x 10 <sup>-2</sup>	34%	~
Ant-Apache v1.5	3	400	<b>68%</b>	<b>3.854 x 10<sup>-4</sup></b>	66%	3.857x 10 <sup>-4</sup>	66%	3.856x 10 <sup>-4</sup>	<b>68%</b>	3.856x 10 <sup>-4</sup>
	5	600	76%	<b>4.678 x 10<sup>-4</sup></b>	76%	4.702 x 10 <sup>-4</sup>	<b>78%</b>	4.709 x 10 <sup>-4</sup>	52%	5.035 x 10 <sup>-4</sup>
	8	750	<b>88%</b>	<b>6.111 x 10<sup>-4</sup></b>	81%	6.308 x 10 <sup>-4</sup>	83%	6.313 x 10 <sup>-4</sup>	46%	~
	10	1000	<b>93%</b>	<b>2.861 x 10<sup>-3</sup></b>	84%	3.008 x 10 <sup>-3</sup>	86%	2.999 x 10 <sup>-3</sup>	34%	~
	15	1500	<b>97%</b>	<b>5.928 x 10<sup>-3</sup></b>	82%	6.478 x 10 <sup>-3</sup>	81%	6.399 x 10 <sup>-3</sup>	31%	~
Ant-Apache v1.7.0	3	400	<b>61%</b>	<b>4.326 x 10<sup>-3</sup></b>	61%	4.328 x 10 <sup>-3</sup>	63%	4.329 x 10 <sup>-3</sup>	62%	4.327 x 10 <sup>-3</sup>
	5	600	<b>73%</b>	<b>5.612 x 10<sup>-3</sup></b>	69%	5.692 x 10 <sup>-3</sup>	71%	5.694 x 10 <sup>-3</sup>	53%	6.001 x 10 <sup>-3</sup>
	8	750	<b>86%</b>	<b>6.269 x 10<sup>-3</sup></b>	84%	6.516 x 10 <sup>-3</sup>	82%	6.521 x 10 <sup>-3</sup>	48%	~
	10	1000	<b>89%</b>	<b>4.007 x 10<sup>-2</sup></b>	87%	4.301 x 10 <sup>-2</sup>	84%	4.147 x 10 <sup>-2</sup>	37%	~
	15	1500	<b>92%</b>	<b>7.006 x 10<sup>-2</sup></b>	84%	7.840 x 10 <sup>-2</sup>	81%	7.688 x 10 <sup>-2</sup>	31%	~
Gantt v1.10.2	3	400	63%	<b>5.111 x 10<sup>-3</sup></b>	<b>64%</b>	5.113 x 10 <sup>-3</sup>	63%	5.114 x 10 <sup>-3</sup>	66%	<b>5.111 x 10<sup>-3</sup></b>
	5	600	<b>69%</b>	<b>6.601 x 10<sup>-3</sup></b>	67%	6.702 x 10 <sup>-3</sup>	<b>69%</b>	6.701 x 10 <sup>-3</sup>	54%	6.956 x 10 <sup>-3</sup>
	8	750	<b>88%</b>	<b>7.899 x 10<sup>-3</sup></b>	83%	8.101 x 10 <sup>-3</sup>	81%	8.116 x 10 <sup>-3</sup>	44%	~
	10	1000	<b>93%</b>	<b>3.333 x 10<sup>-2</sup></b>	84%	3.637 x 10 <sup>-2</sup>	86%	3.536 x 10 <sup>-2</sup>	37%	~
	15	1500	<b>97%</b>	<b>5.807 x 10<sup>-2</sup></b>	82%	6.202 x 10 <sup>-2</sup>	83%	6.125 x 10 <sup>-2</sup>	28%	~
Azureus v2.3.0.6	3	400	<b>61%</b>	<b>6.429 x 10<sup>-4</sup></b>	64%	6.432 x 10 <sup>-4</sup>	61%	6.431 x 10 <sup>-4</sup>	64%	6.431 x 10 <sup>-4</sup>
	5	600	<b>78%</b>	<b>6.708 x 10<sup>-4</sup></b>	72%	6.782 x 10 <sup>-4</sup>	75%	6.788 x 10 <sup>-4</sup>	47%	6.991 x 10 <sup>-4</sup>
	8	750	<b>84%</b>	<b>6.976 x 10<sup>-4</sup></b>	80%	7.205 x 10 <sup>-4</sup>	82%	7.212 x 10 <sup>-4</sup>	34%	~
	10	1000	<b>91%</b>	<b>2.745 x 10<sup>-3</sup></b>	84%	2.976 x 10 <sup>-3</sup>	87%	2.877 x 10 <sup>-3</sup>	26%	~
	15	1500	<b>94%</b>	<b>4.981 x 10<sup>-3</sup></b>	81%	5.508 x 10 <sup>-3</sup>	78%	5.394 x 10 <sup>-3</sup>	26%	~

i.e., the vector composed of 15 zeros. Also, MOEA/D seems to have better convergence than IBEA. However, NSGA-II is unable to progress in terms of convergence as its non-dominated solutions are situated far from the ideal vector. We conclude that although NSGA-II is the most famous multi-objective algorithm in SBSE, it is not adequate for problems involving more than 3 objectives. Based on the results we obtained for the refactoring problem, it appears that NSGA-III is a very good candidate solution for tackling many-objective SBSE problems.

We also compared the results of NSGA-III using 15 objectives and two mono-objective refactoring approaches on all seven open source systems as described in Figure 4. From the set of non-dominated solutions generated by NSGA-III, we selected the solution closest to the ideal point. NSGA-III performed better than both mono-objective algorithms in 100% of cases. In fact, since mono-objective algorithms aggregate all metrics in one objective there is a loss of information due to the conflicting nature of the used quality metrics. It is clear that mono-objective algorithms did not perform well in terms of fixing code smells whereas NSGA-III fixed on average more than 92% of them, especially in the case of large systems such as Ant Apache and Argo UML. In general, large systems contain many

different types of code smells and so a large number of metrics is required to evaluate the quality of a system after applying a refactoring solution. We asked 8 PhD students in Software Engineering to manually evaluate some suggested refactorings if they are useful and semantically make sense or not. As described in Figure 5, the use of high number of metrics such as coupling and cohesion not only improve the structure by fixing code-smells but also helped our NSGA-III based algorithm to generate refactoring solutions that semantically make sense and. In all the considered systems, NSGA-III outperforms existing work in terms of semantics preservation with at least 73% of precision on every system.

## 4.1 Discussion

### Computational time (CT):

When using optimization techniques, the most time consuming operation is the evaluation step. Thus, we studied the execution time of all many/multi-objective algorithms used in our experiments. Figure 5 shows the evolution of the running times of the different algorithms on the ArgoUMLv0.26 system, the largest system in our experiments. It is clear from this figure, that the multi-objective algorithm (NSGA-II) has similar running times for the 3- and 5-



objective cases. However, for higher numbers of objectives NSGA-III is faster than IBEA. This observation could be explained by the computational effort required to compute the contribution of each solution in terms of hypervolume. In comparison to MOEA/D, MOEA/D is slightly faster than NSGA-III since it does not make use of non-dominated sorting.

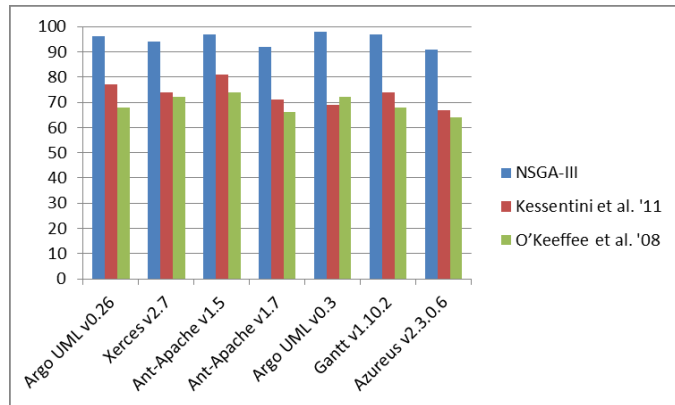


Figure 4. Median NF values on 31 independent runs using 7 systems comparing NSGA-III and two mono-objective approaches [19][23].

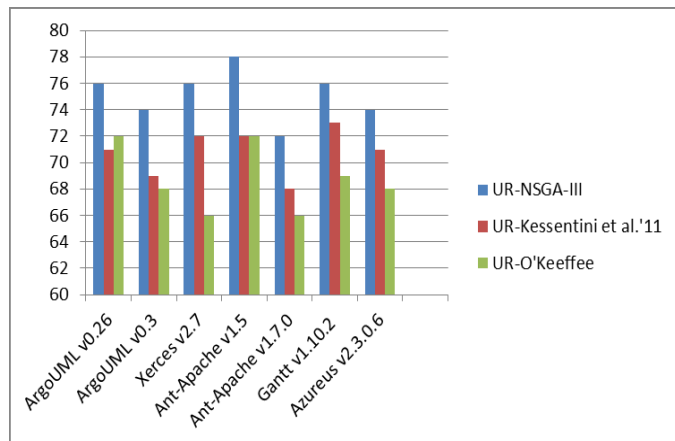


Figure 5. Median UR (semantic coherence) values on 31 independent runs using 7 systems comparing NSGA-III and two mono-objective approaches [19][23].

#### Quality improvements vs. number of objectives:

One of the main motivations of our work is to propose a scalable search-based software engineering approach that can address software engineering problems with a large number of objectives to be optimized. Thus, we evaluated the impact of taking into consideration a higher number of objectives (metrics) on the quality of the refactoring solutions. In fact, the symptoms of code smells can be formalized in terms of quality metrics, thus if we consider more metrics in evaluating a refactoring solution there is a better chance that more code smells are fixed. Figure 5 confirms this. The percentage of fixed code smells increases from 63% to 98% as the number of objectives/metrics increases from 3 to 15 objectives. This result allows us to affirm RQ3.

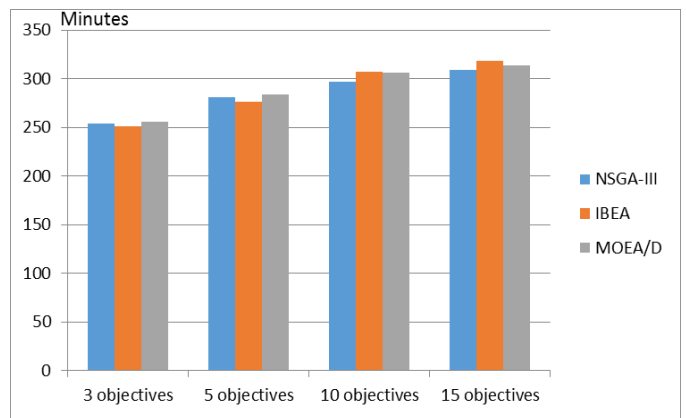


Figure 6. Median CT values on 31 independent runs using 7 systems comparing NSGA-III and two mono-objective approaches.

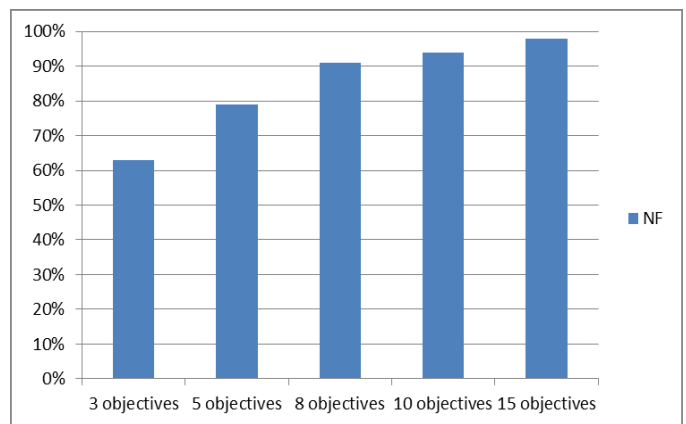


Figure 7. Median NF values on 31 independent runs using ArgoUML v 0.3 with 3, 5, 8, 10 and 15 objectives.

## 5. RELATED WORK

Search-based refactoring represents fully automated refactoring driven by metaheuristic search and guided by software quality metrics and used subsequently to address the problem of automating design improvement [23]. Seng et al. [27] propose a search-based technique that uses a genetic algorithm over refactoring sequences. The employed metrics are mainly related to various class level properties such as coupling, cohesion, complexity and stability. The approach was limited only to the use of one refactoring operation type, namely 'move method'. In contrast to O'Keefe et al. [23], their fitness function is based on well-known measures of coupling between program components. Both these approaches use weighted-sum to combine metrics into a fitness function, which is of practical value but is a questionable operation on ordinal metric values. Kessentini et al. [19] also propose a single-objective combinatorial optimization using a genetic algorithm to find the best sequence of refactoring operations that improve the quality of the code by minimizing as much as possible the number of code smells detected using a set of quality metrics. Kilic et al. explore the use of a variety of population-based approaches to search-based parallel refactoring, finding that local beam search could find the best solutions. Harman and Tratt were the first to introduce the concept of Pareto optimality to search-based refactoring [14]. They use it to combine two metrics into a fitness function, namely CBO (coupling between objects) and SDMPC (standard deviation of methods per class), and demonstrate

that it has several advantages over the weighted-sum approach. More recent work on multi-objective search-based refactoring is the work by Ouni et al. [24] who propose a multi-objective optimization approach to find the best sequence of refactorings using NSGA-II. The proposed approach is based on two objective functions, quality (proportion of corrected code smells) and code modification effort, to recommend a sequence of refactorings that provide the best trade-off between quality and effort.

## 6. CONCLUSIONS AND FUTURE WORK

This paper represents the first real-world application of NSGA-III and the first scalable work that supports the use of 15 objectives to address a software engineering problem. In our approach, refactoring solutions are evaluated using a set of 15 software quality metrics. We evaluated our approach on seven large open source systems [28][29][30][31][32]. The experimental results indicate that NSGA-III outperforms other many-objective algorithms (IBEA [31] and MOEA/D [30]), NSGA-II and mono-objective evolutionary algorithms [19][23]. As part of the future work, we plan to work on adapting NSGA-III to additional software engineering problems and we will perform more comparative studies on larger open source systems.

## ACKNOWLEDGEMENT

This work was supported, in part, by the Institute for Advanced Vehicle Systems-Michigan grant and the Science Foundation Ireland grant 10/CE/11855 to Lero - the Irish Software Engineering Research Centre.

## 7. REFERENCES

- [1] Arcuri A. and Fraser G., 2013. *Parameter tuning or default values? An empirical investigation in search-based software engineering*, Empirical Software Engineering, 18(3).
- [2] Bader, J. and Zitzler, E. 2011. Hype: An algorithm for fast hypervolumebased many-objective optimization. In *Evolutionary Computation*. vol. 19, no. 1. 45–76.
- [3] Bechikh, S., Ben Said, L. and Ghédira, K. 2010. Estimating Nadir Point in Multi-objective Optimization using Mobile Reference Points. In CEC'10. 2129–2137.
- [4] Brown, W. J., Malveau, R. C., Brown, W. H., and Mowbray, T. J. 1998. *Anti Patterns: Refactoring Software, Architectures, and Projects in Crisis*. John Wiley and Sons, 1st Ed.
- [5] Deb, K. 2001. *Multiobjective Optimization using Evolutionary Algorithms*. John Wiley and Sons, Ltd, New York, USA.
- [6] Deb, K. and Jain H. 2012. Handling many-objective problems using an improved NSGA-II procedure. In *Proceedings of IEEE Congress on Evolutionary Computation*. 1–8.
- [7] Deb, K. and Jain, H. An Evolutionary Many-Objective Optimization Algorithm Using Reference-point Based Non-dominated Sorting Approach, Part I: Solving Problems with Box Constraints. In *Proceedings of IEEE Transactions on Evolutionary Computation*. accepted.
- [8] Deb, K., Pratap, A., Agarwal, S. and Meyarivan, T. 2002. A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II. In *Proceedings of IEEE Transactions on Evolutionary Computation*, vol. 6, no. 2. 182–197.
- [9] di Piero, F., Khu, S-T. and Savic, D.A. 2007. An Investigation on Preference Order Ranking Scheme for Multiobjective Evolutionary Optimization. In *Proceedings of IEEE Transactions on Evolutionary Computation*. vol. 11, no. 1. 17–45.
- [10] Fenton, N. and Pfleeger, S. L. 1997. *Software Metrics: A Rigorous and Practical Approach*, 2nd ed. International Thomson Computer Press.
- [11] Ferrucci, F., Harman, M., Ren, J. and Sarro, F. 2013. Not going to take this anymore: multi-objective overtime planning for software engineering projects. In ICSE '13, 462–471.
- [12] Fowler, M., Beck, K., Brant, J., Opdyke, W. and Roberts, D. 1999. *Refactoring – Improving the Design of Existing Code*. 1st ed. Addison-Wesley.
- [13] Harman, M. 2013. Software Engineering: An Ideal Set of Challenges for Evolutionary Computation, In GECCO '13, 1759–1760.
- [14] Harman, M. and Tratt, L. 2007. Pareto optimal search based refactoring at the design level. In GECCO'07. 1106–1113.
- [15] Harman, M., Mansouri, S. A. and Zhang, Y. 2012. Search-based software engineering: Trends, techniques and applications. In *ACM Computing Surveys*. vol. 45, no.1. 61 pages.
- [16] Harman, M., and Jones, B., 2001. *Search Based Software Engineering*, Journal of Information and Software Technology, 43(14):833-839.
- [17] Jaimes, A. L., Coello Coello, C.A. and Barrientos, J. E. U. 2009. Online Objective Reduction to Deal with Many-objective Problems. In *the 5th international conference on Evolutionary Multicriterion Optimization*. 423–437.
- [18] Jain, H. and Deb, K. 2013. An Improved Adaptive Approach for Elitist Nondominated Sorting Genetic Algorithm for Many-Objective Optimization. In EMO'13, 307–321.
- [19] Kessentini, M., Kessentini, W., Sahraoui, H., Boukadoum, M., and Ouni, A. 2011. Design Defects Detection and Correction by Example. In ICPC'11. 81-90.
- [20] Kukkonen, S. and Lampinen, J. 2007. Ranking-Dominance and Many-Objective Optimization. In *Proceedings of IEEE Congress on Evolutionary Computation*. 3983-3990.
- [21] Moha, N., Guéhéneuc, Y.-G., Duchien, L. and Le Meur, A.-F. 2009. DECOR: A Method for the Specification and Detection of Code and Design Smells. In *TSE, 2009, vol 12,*. 20–36.
- [22] Ó Cinnéide, M., Tratt, L., Harman, M., Counsell, S. and Moghadam, I. H. 2012. Experimental Assessment of Software Metrics Using Automated Refactoring. In ESEM'12, 49-58.
- [23] O'Keefe, M. K. and Ó Cinnéide, M. 2008. Search-based refactoring for software maintenance. In *Journal of Systems and Software*. vol. 81, no.4. 502-516.
- [24] Ouni, A., Kessentini, M., Sahraoui, H. and Boukadoum, M. 2012. Maintainability Defects Detection and Correction: A Multi-Objective Approach. *Journal of Automated Software Engineering*. vol. 20, 47-79.
- [25] Palomba, F., Bavota, G., Di Penta, M., Oliveto, R., De Lucia, A. and Poshyvanyk, D. 2013. Detecting Bad Smells in Source Code Using Change History Information. In ASE 2013.
- [26] Sayyad, A., Menzies, T. and Ammar, H. 2013. On the value of user preferences in search-based software engineering: a case study in software product lines. In ICSE '13. 492-501.
- [27] Seng, O., Stammel, J. and Burkhart, D. 2006. Search-based determination of refactorings for improving the class structure of object-oriented systems. In GECCO'06. 1909–1916.
- [28] Singh, H. K., Isaacs, A. and Ray, T. 2011. A pareto corner search evolutionary algorithm and dimensionality reduction in many-objective optimization problems. In *Proceedings of IEEE Transactions on Evolutionary Computation*. vol. 99. 1–18.
- [29] Yamashita A., and Moonen, L., 2012. Do code smells reflect important maintainability aspects?, Proceedings of ICSM2012.
- [30] Zhang, Q. and Li, H. 2007. MOEA/D: A Multiobjective Evolutionary Algorithm Based on Decomposition. In *Proceedings of IEEE Transactions on Evolutionary Computation*. vol. 11, no. 6. 712–731.
- [31] Zitzler, E. and Künzli, S. 2004. Indicator-based selection in multiobjective search Parallel Problem Solving from Nature. In Proceedings of PPSN'04.
- [32] (Refactorings in Alphabetical Order). [Online]. Available: <http://www.refactoring.com/catalog/index.html>