

Prioritising abstract test cases: An empirical study

Rubing Huang, Weiwen Zong, Tsongyueh Chen, Dave Towey, Yunan Zhou
and Jinfu Chen



**University of
Nottingham**
UK | CHINA | MALAYSIA

University of Nottingham Ningbo China, Taikang East Road, Ningbo,
Zhenjiang, 315100, China

First published 2019

This work is made available under the terms of the Creative Commons
Attribution 4.0 International License:

<http://creativecommons.org/licenses/by/4.0>

The work is licenced to the University of Nottingham Ningbo
China under the Global University Publication Licence:

<https://www.nottingham.edu.cn/en/library/documents/research-support/global-university-publications-licence.pdf>



**University of
Nottingham**
UK | CHINA | MALAYSIA

Prioritizing Abstract Test Cases: An Empirical Study

R. Huang¹ W. Zong¹ T. Y. Chen² D. Towey³ Y. Zhou¹ J. Chen¹

¹School of Computer Science and Communication Engineering, Jiangsu University, Zhenjiang, Jiangsu 212013, P.R. China

²Department of Computer Science and Software Engineering, Swinburne University of Technology, Hawthorn, 3122, Australia

³School of Computer Science, University of Nottingham Ningbo China, Ningbo, Zhejiang 315100, P.R. China

* E-mail: rbhuang@ujs.edu.cn

Abstract: Test case prioritization (TCP) attempts to schedule the order of test case execution such that faults can be detected as quickly as possible. TCP has been widely applied in many testing scenarios, such as regression testing, and fault localization. Abstract test cases (ATCs) are derived from models of the system under test, and have been applied to many testing environments, such as model based testing, and combinatorial interaction testing. Although various empirical and analytical comparisons for some ATC prioritization (ATCP) techniques have been conducted, to the best of our knowledge, no comparative study focusing on the most current techniques has yet been reported. In this study, we investigated 18 ATCP techniques, categorized into four classes. We conducted a comprehensive empirical study to compare 16 of the 18 ATCP techniques in terms of their testing effectiveness and efficiency. We found that different ATCP techniques could be cost-effective in different testing scenarios, allowing us to present recommendations and guidelines for which techniques to use under what conditions.

1 Introduction

The order of test case execution in a given test set can be very important, especially when testing resources are limited. The main reason is that a well-prioritized execution order of test cases may be able to trigger failures more quickly, and thus allow the follow-up processes to be conducted earlier (including fault localization, diagnosis and correction). The process of scheduling the execution order of test cases is called test case prioritization (TCP) [1], and it has been applied in various testing environments, including regression testing [2].

To date, many TCP algorithms have been designed to prioritize different test case types according to different criteria, including code coverage based prioritization [1, 3], search based prioritization [4, 5], adaptive random prioritization [6–9], and similarity based prioritization [10, 11] (the interested reader is referred to two survey papers for more details [12, 13]). An abstract test case (ATC) [14] (model input [15]) is an important test case type that can be extracted from a designed model of the system under test (SUT) [16]. In combinatorial interaction testing [17], for example, an SUT may be impacted by different parameters (or factors), each of which may contain a finite number of values (or levels). In this case, ATCs can be created by assigning a value for each parameter. ATCs have been widely used in many testing approaches including model based testing [18], and category-partition testing [16]. Abstract test case prioritization (ATCP) has also been widely studied in different fields, especially in combinatorial interaction testing [19–21], and software product line testing [22, 23].

Although there have been empirical and analytical comparisons of individual or several ATCP techniques [15, 21, 24], to the best of our best knowledge there has not yet been a comprehensive comparative study focusing on the most current techniques. In our study, we investigated 18 ATCP techniques, grouped into four categories: non-information-guided prioritization (NIGP); interaction coverage based prioritization (ICBP); input-model mutation based prioritization (IMBP); and similarity based prioritization (SBP). We conducted a comprehensive empirical study using five subject programs (written in the C language), each of which had six versions. In the study, based on mutation analysis, the testing effectiveness and efficiency of each ATCP technique were investigated.

We believe that this is the most extensive and inclusive empirical study comparing ATCP techniques so far reported in the literature. Based on the experimental results, some empirical findings are provided, and some recommendations and guidelines are given for testers when choosing ATCP techniques in different testing scenarios. In summary, the main contributions of this work are as follows:

- (1) We selected 18 ATCP techniques from the literature, and divided them into four categories, in terms of the different information used to guide the prioritization process.
- (2) We conducted empirical studies to compare 16 of the 18 ATCP techniques, according to three quality evaluation measures: interaction coverage rate, fault detection rate, and prioritization cost.
- (3) We present empirical findings comparing ATCP techniques among each category and between different categories.
- (4) We provide recommendations and guidelines for testers to help select ATCP techniques in different testing scenarios.

The structure of the rest of this paper is organized as following: Section 2 introduces some preliminary information and background details. Section 3 provides the details about the experimental settings, and Section 4 presents the experimental results to answer the research questions. Finally, Section 5 concludes the paper, and discusses potential future work.

2 Preliminaries and Background

Some preliminary information is presented in this section, including details about abstract test cases, and test case prioritization (TCP). ATCP techniques are described, the strength and weakness of each technique are summarized, and previous empirical work is also discussed.

2.1 Preliminaries

2.1.1 Abstract Test Case: A system under test (SUT) is generally influenced by different parameters or factors (for example, configurations, features, components, etc), with each parameter having a certain number of possible values or levels. In general, most

Table 1 An example for input model

Factor	p1: OS	p2: Browser	p3: Access	p4: Proxy
	Windows (0)	IE (2)	ISDL (4)	No Proxy (7)
Level	Mac OS X (1)	Safari (3)	Modem (5)	HTTP (8)
			VPN (6)	SOCKS5 (9)

Browser = "IE" ! **OS** = "Windows", i.e., p2 = "2" ! p1 = "0".

Browser = "Safari" ! **OS** = "Mac OS X", i.e., p2 = "3" ! p1 = "1".

1 SUTs may have constraints among parameter values: that is, some
2 value combinations are not feasible. Based on this, we present the
3 following definition of an input model [25] (or input parameter
4 model [14]) used for modeling the SUT.

5 Definition 1. An input model, $Model((p_1, p_2, \dots, p_k), (L_1,$

6 $L_2, \dots, L_k), C)$, is the information about the parameters and the
7 values of each parameter of the SUT (with k parameters), a set of
8 values L_i for the i -th parameter p_i , and a set of value combination

9 constraints C .

10 As shown in Table 1, for example, an input model with value combination
11 constraints is used for a web application such as a browser
12 game, where four parameters are included, of which the first has two
13 values, and the last three all have three. Since the browser "IE" is
14 developed for the OS "Windows", and the browser "Safari" is
15 developed for the OS "Mac OS X", two value combination constraints
16 are obtained. To simplify the problem, each parameter is
17 denoted by p_i ($i = 1, 2, 3, 4$), and each value is labelled by an integer,
18 beginning with 0 and incrementing by 1, from p_1 to p_4 (see
19 Table 1).

20 Therefore, the model for above example can be represented by

21 $Model((p_1, p_2, p_3, p_4), \{\{0, 1\}, \{2, 3\}, \{4, 5, 6\},$
22 $\{7, 8, 9\}\}, C = \{p_2 = 2 \wedge p_1 = 0, p_2 = 3 \wedge p_1 =$
23 $1\})$, containing two value combination constraints, and four
24 parameters, of which the first two parameters have two values, and
25 another two parameters have three values. Since the detailed values
26 of each parameter provide no influence on the model, without
27 loss of generality, we adopt an abbreviated version in this paper:

28 $Model((L_1, L_2, \dots, L_k), C)$. Accordingly, the above example can be
29 described as $Model(2, 3, C = \{2 \wedge 0, 3 \wedge 1\})$.

30 When an input model is available, construction of abstract test
31 cases (ATCs) [14] (or model inputs [15]) for testing the SUT is
32 possible. The definition of the abstract test case is given as follows.

33 Definition 2. An abstract test case, (v_1, v_2, \dots, v_k) , is a k -tuple,

34 where $v_i \in L_i$ ($i = 1, 2, \dots, k$).

35 An ATC is valid if C is satisfied, otherwise it is invalid. For
36 instance, in the previous example, a valid ATC is $(0, 2, 5, 8)$; and
37 an invalid one is $(0, 3, 4, 8)$ — due to violation of the constraint
38 $((p_2 = 3) \wedge (p_1 = 1))$. Intuitively speaking, each ATC with size \square
39 can cover λ -tuples $1 \leq \lambda \leq \square$, where such a tuple is called a λ -wise
40 value combination [26] or a λ -wise schema [17]. For example, an
41 ATC $(1, 3, 5, 9)$ covers six 2-wise value combinations: $(1, 3)$, $(1, 5)$,
42 $(1, 9)$, $(3, 5)$, $(3, 9)$, and $(5, 9)$.

43 The ATCs have been used in many applications such as
44 configuration-aware systems [27, 28], and software product
45 lines [29]. Many testing methods have focused on the generation
46 and construction of ATCs, such as category-partition testing [16],
47 combinatorial testing [17], and random testing [30].

48 2.1.2 Test Case Prioritization: Test case prioritization seeks to
49 schedule test cases such that those with the highest significance, in
50 terms of some criteria, are run earlier than those with lower significance.
51 When testing resources are limited or insufficient for the
52 execution of a complete test suite, then a good order of test case execution
53 can be very important. The problem of test case prioritization
54 can be defined as follows [1].

55 Definition 3. Given a test suite T to be prioritized, \square being the
56 set of all possible orders of test cases by permuting T , and f being
57 a fitness function to evaluate each permutation, the problem of test

case prioritization is to identify a permutation $S \in \square$ such that: 58

$$(8S^0) (S^0 \square) (S^0_{6=S}) [f(S) f(S^0)] \quad (1)$$

2.2 ATCP Techniques 59

Depending on the type of information used, as with other testing
60 approaches, ATCP can be considered either black-box or white-
61 box testing [15]. ATCP approaches using models of the SUT, for
62 example, would be considered black-box, because no access to
63 source code is necessary. In this paper we focus on black-box ATCP
64 techniques (interested readers may refer to work by Rothermel et
65 al. [1] or Zhang et al. [31] for discussion of white-box approaches).
66 According to the information used to guide the prioritization process,
67 the ATC prioritization techniques (ATCP) are mainly classified into
68 the following four categories. 69

2.2.1 Non-Information-Guided Prioritization (NIGP): The
70 NIGP strategies discussed in this section can be used for not only
71 abstract test cases but for all types of test cases, because this category
72 does not use additional information to support the prioritization
73 process. 74

- Test-case-generation prioritization (TCGP): TCGP prioritizes 75
ACTs using the order in which the test cases were generated. 76
- Reverse test-case-generation prioritization (RTCGP): RTCGP 77
prioritizes ACTs by reversing the generation order. 78
- Random test case prioritization (RTCP): RTCP randomly orders 79
ACTs, according to uniform distribution. 80

2.2.2 Interaction Coverage Based Prioritization (ICBP): 81

The ICBP strategy makes use of the information of coverage information
82 to support the process of ATCP. By using different levels of
83 interaction coverage, the following three ATCP techniques are considered:
84 fixed-strength ICBP (FICBP), incremental-strength ICBP
85 (IICBP), and aggregate-strength ICBP (AICBP). 86

- Fixed-strength ICBP (FICBP): FICBP [32] iteratively selects the 87
element as the next test case from candidate ATCs such that it covers
88 the largest number of λ -wise value combinations that have not yet
89 been covered by the ATCs already selected. Before prioritization,
90 FICBP needs to assign a value to an integer λ , the prioritization
91 strength. Based on previous investigations [21, 24, 33–35], the
92 assignment of the prioritization strength usually ranges from 1 to 6.
93 To reduce the prioritization cost, a new FICBP technique has been
94 proposed that uses repeated base-choice coverage, FICBPR [36].
95 Although FICBPR leverages a similar mechanism to FICBP, it only
96 assigns a value of 1 to the prioritization strength λ , and forgets
97 previous prioritization details when the coverage of 1-wise value
98 combinations is fully achieved. 99

- Incremental-strength ICBP (IICBP): IICBP [37, 38] first uses a
100 small prioritization strength λ ($\lambda = 1$), and presents it to the FICBP
101 algorithm for prioritizing the candidates. Once all λ -wise value
102 combinations have been covered by selected test cases, IICBP increases
103 the prioritization strength with an increment i — $\lambda = \lambda + i$ ($i = 1$)
104 — and then uses this new prioritization strength for the FICBP
105 algorithm to prioritize remaining ATCs. This process is repeated
106 until all candidates have been chosen. In this study, we used the
107 IICBP algorithm from Huang et al. [38], initially setting λ to 1, and
108 i to 1. 109

- Aggregate-strength ICBP (AICBP): AICBP [20] makes use of
110 hybrid interaction coverage by considering different prioritization
111 strength λ values ranging from 1 to the generation strength \square in
112 combinatorial testing [17]. As we know, \square is chosen in the stage of
113 test suite construction, however, it may be not applicable to adopt
114 previous AICBP algorithms for prioritizing ATC sets (because it
115 is infeasible to choose the value of \square). In this paper, therefore,
116 we use a simplified version of AICBP that only takes prioritization
117 strength $\lambda = 1, 2$, and 3 into consideration (i.e., $\square = 3$), and
118 can thus be used for prioritizing any sets of ATCs. The mechanism
119 of the AICBP algorithm is similar to that of FICBP, except
120 that AICBP uses hybrid interaction coverage by aggregating three
121

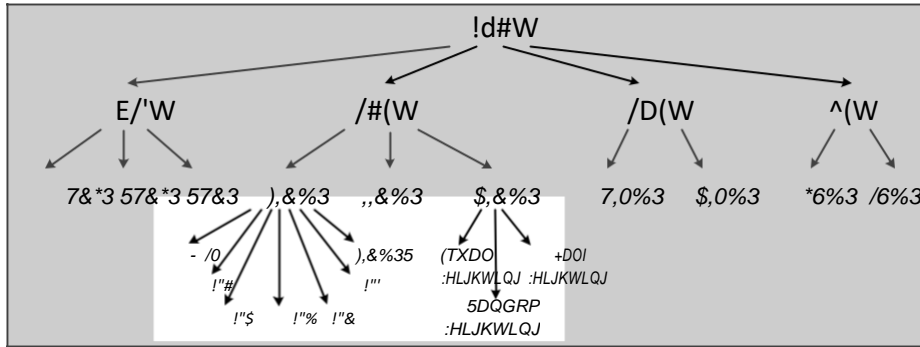


Fig. 1: Overview of ATPC techniques

1 prioritization strengths 1, 2, and 3) rather than interaction coverage
 2 by using a single prioritization strength. As discussed by Huang et
 3 al. [20], there are three weighting distributions for different priori-
 4 zation strengths, i.e., three ways of assigning the weightings $!1$,
 5 $!2, \dots$, and $!3$ to the prioritization strength $\lambda_1, \lambda_2, \dots$, and λ_3 ,
 6 respectively, where $!1 + !2 + \dots + !3 = 1.0$. More specifically,
 7 Equal Weighting assigns the same weighting to each prioritization
 8 strength, i.e., $!1 = !2 = \dots = !3 = \frac{1}{3}$; Random Weighting ran-
 9 domly assigns the weighting to each prioritization strength; and Half
 10 Weighting sets the weighting as following: $!1 = !, !j + 1 = \frac{1}{j},$

11 and $!3 = 1.0$ ($!1 + !2 + \dots + !3 = 1$).

12 2.2.3 Input-model Mutation Based Prioritization (IMBP):
 13 The IMBP strategy [15] creates the mutants of the flattened model
 14 that is derived from the SUT's input model, and then uses the mutant
 15 detection capability of each test case to guide the process of ATPC.
 16 More specifically, IMBP first mutates the flattened model from [25]
 17 to obtain a mutant by changing a constraint, for example, the con-
 18 straint from the input model is ("2" ! "0"), and a mutant may be
 19 ("2" ! "1"). The mutants that are distinguished by the test cases
 20 are killed; otherwise they are live. After that, IMBP prioritizes test
 21 cases based on their capabilities of killing mutants. Based on different
 22 selection strategies, two IMBP techniques are included: 'total'
 23 IMBP (TIMBP) and 'additional' IMBP (AIMBP) [15].

24 • Total IMBP (TIMBP): TIMBP refers to previous 'total' TCP
 25 strategies [1, 31], by repeatedly choosing each element as the
 26 next test case from the remaining candidates such that it kills the
 27 maximum (total) number of model mutants.

28 • Additional IMBP (AIMBP): Similar to TIMBP, AIMBP refers to
 29 previous 'additional' TCP strategies [1, 31], by repeatedly selecting
 30 the next test case which can kill the largest number of model mutants
 31 that have not yet been detected by previously selected ATCs.

32 2.2.4 Similarity Based Prioritization (SBP): SBP [23] makes
 33 use of the Jaccard similarity between candidates to prioritize test
 34 cases, with each ATC being a set of parameter values. In particular,
 35 SBP selects each next test case such that it achieves the small-
 36 est similarity to previously selected test cases. Based on different
 37 implementations, Henard et al. [23] introduced two versions of SBP:
 38 global SBP (GSBP) and local SBP (LSBP).

39 • Global SBP (GSBP): GSBP first determines the first two test
 40 cases by choosing two elements from candidates with the minimum
 41 similarity, then iteratively selects a candidate as the next test case.
 42 In detail, for each candidate c , GSBP first calculates the similarity
 43 between c and each already selected ATC, and sums the similarity
 44 values as the fitness value of c . Then the candidate with the minimum
 45 fitness value is chosen as the next test case.

46 • Local SBP (LSBP): LSBP iteratively identifies a pair of candidates
 47 sharing the minimum similarity as the next two test cases, until all

candidate test cases are selected. The order of the two test
 cases is determined in a random manner.

Figure 1 shows an overview of ATPC techniques, involving
 four categories with 18 techniques.

2.2.5 Strengths and Weaknesses: In this section, we briefly summarize
 the strengths and weaknesses of ATPC techniques, listed
 as follows:

- (1) For the NIGP category, its main advantage may be high testing efficiency (for example, low prioritization time); however, its disadvantage may be low testing effectiveness. The main reason for this is that the NIGP category does not use additional information to guide the prioritization process.
- (2) As for the ICBP category, its main benefit is that each ATPC technique makes use of the information of interaction coverage to prioritize ATCs, resulting in high testing effectiveness. Regarding the drawbacks, FICBP may face the challenges of choosing an appropriate prioritization strength, as different prioritization strengths may lead to different testing performances; and AICBP may require more prioritization time, because it uses more information for the prioritization. IICBP can be considered as a balanced technique compared with FICBP: it may have better testing effectiveness than FICBP with low prioritization strengths but less testing efficiency than that with high prioritization strengths.
- (3) For the IMBP category, its main strength is that it brings the concept of mutation analysis [39] to the input model of the system under test, which may provide some new insights for ATPC. However, it may face some potential challenges, for example, the quality of mutants may influence the performance of IMBP. Generally speaking, AIMBP may have better testing effectiveness but worse testing efficiency than TIMBP, because it requires collecting more information.
- (4) Regarding the SBP category, its main strength is that it may achieve high testing efficiency with comparable testing effectiveness to FICBP. However, it may suffer from the drawback of needing to choose the appropriate similarity measure between ATCs. Intuitively speaking, GSBP may have better performance than LSBP, because the former adopts more information for choosing each element from candidates as the next test case.

Based on this analysis, when testing resources are limited, it may be better to use FICBP with a low prioritization strength, SBP, or NIGP. On the contrary, when testing resources are sufficient, it may be better to adopt FICBP with a high prioritization strength, IICBP, or AIMBP. Additionally, the selection of IMBP may depend on the input model of the system under test.

1 2.3 Previous Empirical Work

2 In this section, we report on some previous empirical work into the
3 prioritization of abstract test cases.

4 Petke et al. [24] initially investigated FICBP with the priori-
5 zation strength λ values 2, 3, 4, and 5; and later added $\lambda = 6$ in
6 an extended study [21]. They mainly focused on the analysis of
7 different prioritization strength values used in FICBP for different
8 covering arrays constructed for combinatorial testing [17]. Com-
9 pared with their work, however, our study examines most current
10 ATPC techniques, including, but beyond, FICBP.

11 Henard et al. [23] proposed two similarity based ATPC algo-
12 rithms, (GSBP and LSBP), and compared them with the random test
13 case prioritization and 2-wise FICBP technique. Similar to Petke et
14 al. [24], their work focused on the prioritization of combinatorial
15 test suites (i.e., covering arrays). Additionally, they focused mainly
16 testing software product lines, which means that the input models
17 used were binary — each parameter containing exactly two possible
18 values.

19 Henard et al. [15] compared 20 TCP techniques (ten for white-
20 box and ten for black-box) — some of their black-box prioritization
21 techniques have also been considered in our study. Nevertheless,
22 their study focused on the comparison of white-box and black-box
23 test prioritization techniques, whereas our study is a comparison of
24 black-box ATPC techniques.

25 2.4 Research Questions

26 Our study was motivated by a number of outstanding issues in the
27 field of ATPC. The following five research questions (RQs) guided
28 the study in this paper.

- 29 RQ1: How well do the three ICBP strategies studied perform in
30 terms of the rates of interaction coverage and fault detection?
31 – For the FICBP methods, which strength is more suitable for
32 prioritizing ATCs?
33 – For the AICBP methods, which weighting distribution is more
34 effective?
35 – Which level of interaction coverage is adequate for the ICBP?

36 Answering RQ1 will help testers identify which interaction-
37 coverage-based technique is the most effective. For some ICBP
38 sub-categories, we also had sub-questions to further investigate their
39 effectiveness, and also analyzed the main influential parameters. All
40 ICBP methods use interaction coverage information to guide the pri-
41 oritization — but they use different levels of interaction coverage. It
42 is therefore meaningful to study which level of interaction coverage
43 is adequate.

- 44 RQ2: How well do the two IMBP techniques studied perform
45 according to the rates of interaction coverage and fault detection?

46 Answering RQ2 will help testers know which technique is the
47 most suitable for IMBP. Previous studies based on code coverage
48 information [1, 31] have shown that the ‘total’ TCP tech-nique
49 performs better than the ‘total’ TCP technique, but there are no
50 reported observations related to input-model mutation coverage
51 information. It is therefore interesting to investigate this issue.
52 RQ3: How well do the two SBP techniques studied perform in
53 terms of interaction coverage rate and fault detection rate?

54 Answering RQ3 will help testers know which technique is the most
55 suitable for SBP. Previous investigations have indicated that the SBP
56 strategy is an effective technique for ATCs [22, 23], how-ever the
57 comparison between GSBP and LSBP has not yet been fully
58 explored.

59 RQ4: How differently do the NIGP, ICBP, IMBP, and SBPS tech-
60 niques perform, according to interaction coverage rate and fault
61 detection rate?

62 Answering RQ4 will help guide testers in their selections. It is
63 useful for testers to know which prioritization technique, among all
64 studied techniques, has the best performance.

65 RQ5: How do all the ATPC techniques compare in terms of
66 the required prioritization time?

67 ATPC is important, especially when testing resources are too
68 limited to allow execution of all ATCs. It is therefore useful to
69 consider the prioritization time of each prioritization technique.
70 Answering RQ5 will help testers make a decision on the selection
71 of prioritization techniques.

72 3 Methodology

73 3.1 Subject Programs

74 Five subject programs, written in the C language, were chosen.
75 These programs were obtained from the GNU FTP server^{*}. The
76 flex program is a fast lexical analysis generator; the grep program is a
77 widely-used utility for pattern matching; the sed program is a stream editor
78 that performs text transformations on an input stream; the make program
79 is to control the compile and build processes for programs; and the gzip
80 program is a popular command-line tool used for file compression and
81 decompression. These programs have been widely adopted in previous
82 TCP research [1, 7, 15, 21, 24, 34, 35].

83 Table 2 describes the detailed information for each subject pro-
84 gram such as, the version number, the year of release, the uncom-
85 mented size of code (measured by cloc^{\dagger}), and the number of

^{*}<http://ftp.gnu.org/>

[†]<http://cloc.sourceforge.net/>

Table 2 Subject Programs

Program	Input Model	Test Pool	Information	V0	V1	V2	V3	V4	V5
flex	Model($2^6 3^2 5^1$, C1), C1 = 32	500	Version	2.4.3 (1993)	2.4.7 (1994)	2.5.1 (1995)	2.5.2 (1996)	2.5.3 (1996)	2.5.4 (1997)
			LOC	8,959	9,470	12,231	12,249	12,370	12,366
			Faults	-	32	32	20	33	32
grep	Model($2^1 3^4 2^5 1^6 1^8$, C2), C2 = 58	440	Version	2.0 (1996)	2.2 (1998)	2.3 (1999)	2.4 (1999)	2.5 (2002)	2.7 (2010)
			LOC	8,163	11,988	12,724	12,826	20,838	58,344
			Faults	-	56	58	54	58	59
sed	Model($2^7 3^1 4^1 6^1 10^1$, C3), C3 = 58	324	Version	3.0.1 (1998)	3.0.2 (1998)	4.0.6 (2003)	4.0.8 (2003)	4.1.1 (2004)	4.2 (2009)
			LOC	7,790	7,793	18,545	18,687	21,743	26,466
			Faults	-	16	18	18	19	22
make	Model(2^{10} , C4), C4 = 28	111	Version	3.75 (1996)	3.76.1 (1997)	3.77 (1998)	3.78.1 (1999)	3.79 (2000)	3.80 (2002)
			LOC	17,463	18,568	19,663	20,461	23,125	23,400
			Faults	-	37	29	28	29	28
gzip	Model($2^{13} 3^1$, C5), C5 = 69	156	Version	1.0.7 (1993)	1.1.2 (1993)	1.2.2 (1993)	1.2.3 (1993)	1.2.4 (1993)	1.3 (1999)
			LOC	4,324	4,521	5,048	5,059	5,178	5,682
			Faults	-	8	8	7	7	7

1 mutated faults. The table also gives the information of input models
2 and sizes of candidate ATCs, where all input models and ATC sets
3 were downloaded from a standard library, i.e., the Software Infra-
4 structure Repository (SIR) [40]. These input models were used in
5 previous work by Petke et al. [21, 24])

6 3.2 Fault Seeding

7 For each of the subject programs, the original version does not con-
8 tain any seeded-in faults. There are a number of hand-seeded faults
9 that are available from the SIR [40], but many of these faults can be
10 detected by more than 60% of test cases (on average). Therefore, in
11 this paper we have used mutation analysis [39] to evaluate different
12 ATPC techniques. As discussed in previous studies [41, 42], muta-
13 tion analysis can provide more realistic faults than hand-seeding, and
14 may be more appropriate for studying test case prioritization.
15 For the five subject programs, we used the same mutation faults
16 as used by Henard et al. [15]: that is, we employed the mutant oper-
17 ators set used by Andrews et al. [41], including statement deletion,
18 constant replacement, unary insertion, arithmetic operator replace-
19 ment, logical operator replacement, relational operator replacement,
20 and bitwise logical operator replacement. Following previous prac-
21 tice [1, 31, 41], we removed the duplicate and equivalent mutants,
22 and also removed all those mutants that would not be killed by any
23 ATC. In addition, all subsuming mutants [43] (also called minimum
24 mutants [44] or disjoint mutants [45]) that would be too easily killed
25 were also removed — these mutants may otherwise negatively affect
26 the mutation score measurement [41, 44–46]. A mutation fault is said
27 to be identified by a test case when the output of the original version
28 is different to that of the fault-seeded version. Table 2 shows the
29 number of faults in this study.

30 3.3 The 16 Investigated ATPC Algorithms

31 Table 3 presents an overview of the 16 ATPC techniques studied,
32 giving the mnemonic, description, a reference to its original research
33 publication, and category, for each. For NIGP, we only considered
34 random test case prioritization, because test-case-generation priori-
35 tization (TCGP), and its reversed version (RTCGP), only depend
36 on the original test set. However, because the test pool used in this
37 paper was provided by the SIR [40], which has no correspondence
38 for the original or reversed set, therefore, TCGP and RTCGP were
39 removed from the experiments. For FICBP, we considered the priori-
40 tization strengths $\lambda = 1, 2, 3, 4, 5$, and 6. For AICBP, we considered
41 three weighting distributions of prioritization strengths: equal, ran-
42 dom, and half weighting [20]. For IMBP, the model mutants needed
43 to be seeded, and in this study we used the model mutant matrix file
44 used by Henard et al. [15].

45 For SBP, compared to the previous versions of GSBP and
46 LSBP [23], the algorithms in our study have two main differences:

* <http://henard.net/research/regression/ICSE-2016/>

Table 3 ATPC techniques considered in the experiments

Mnemonic	Description	Reference	Category
RDP	Random test case prioritization	[32]	NIGP
FP1	FICBP at prioritization strength 1	[33]	
FP2	FICBP at prioritization strength 2	[32]	
FP3	FICBP at prioritization strength 3	[32]	
FP4	FICBP at prioritization strength 4	[38]	
FP5	FICBP at prioritization strength 5	[37]	
FP6	FICBP at prioritization strength 6	[21]	ICBP
FPR	FICBPR	[36]	
IIP	IICBP	[37]	
APE	AICBP with Equal Weighting	[20]	
APR	AICBP with Random Weighting	[20]	
APH	AICBP with Half Weighting	[20]	
TIM	TIMBP	[15]	IMBP
AIM	AIMBP	[15]	
SPG	GSBP	[23]	SBP
SPL	LSBP	[23]	

(1) When meeting a tie-breaking case, i.e., there exist more than
one pair of ATCs sharing the same minimum similarity, the original
version adopts a first-test-case tie-breaking technique (i.e., choosing
the first one) [47]. However our study uses the random tie-breaking
technique (i.e., choosing a pair randomly); (2) After choosing the
best pair of ATCs from candidates, the original version adds these
two ATCs to the prioritized set successively, however our study adds
them in a random order. Our GSBP and LSBP algorithms, therefore,
are less biased than the originals.

Because the ATPC techniques involve randomization (due to
random tie-breaking [47]), we ran each experiment 100 times.

3.4 Metrics

To evaluate different ATPC techniques, in this study we focused
on the following three metrics: (a) interaction coverage rate — to
measure the speed of achieving the interaction coverage of each pri-
oritized test suite; (b) fault detection rate — to measure the speed
of identifying faults of each prioritized test suite; and (c) prioritiza-
tion cost — to measure how quickly each prioritized test suite was
obtained.

3.4.1 Interaction Coverage Rate: The average percentage of
combinatorial coverage (APCC) [24, 48] was adopted to evaluate
the speed of achieving the interaction coverage at strength \mathbb{K} for a
prioritized set of ATCs. If $S = \{tc_1, tc_2, \dots, tc_n\}$ is an ordered set

of n ATCs, the \mathbb{K} -wise (1 $\leq k \leq \mathbb{K}$) APCC definition for S is:

$$APCC(\mathbb{K}, S) = \frac{\sum_{i=1}^n \sum_{j=1}^i |CombSet(\mathbb{K}, tc_j)|}{p \sum_{j=1}^i CombSet_j} \quad (2)$$

where $CombSet(\mathbb{K}, tc_j)$ is a set of \mathbb{K} -wise value combinations cov-
ered by the abstract test case tc_j .

The APCC metric values are numerical values ranging from 0.0 to
1.0, with higher values implying better rates of achieving interaction
coverage. Following previous investigations [21], in this paper, six \mathbb{K}
values from 1 to 6 were considered for APCC.

3.4.2 Fault Detection Rate: The average percentage of faults
detected (APFD) was previously used to evaluate different priori-
tization techniques [1]. APFD requires details of the fault-detection
capability of each executed test case.

Let T be a test suite with size n , and F be a set of m faults that
can be detected by T . Let SF_i be the number of test cases, required
to detect fault F_i in a prioritized test suite S of T . The APFD
for S is given by the following equation (from [1]):

$$APFD(S) = 1 - \frac{SF_1 + SF_2 + \dots + SF_m}{n \times m} + \frac{1}{2n} \quad (3)$$

3.4.3 Prioritization Cost: The prioritization cost measures the
prioritization time required for each prioritization technique, and
represents the efficiency of the technique. Obviously, lower priori-
tization costs mean better performance.

3.5 Statistical Analysis

When assessing the statistical significance of the differences
between the APCC or APFD values (used to evaluate each priori-
tization technique), because there was no relationship between any
of the 100 runs, it is reasonable to use an unpaired test [49]. Further-
more, since no assumptions were made about which prioritization
technique is better than others, a two-tailed test is also appro-
priate [49]. Following previous studies dealing with randomized
algorithms [49, 50], we used the unpaired two-tailed Wilcoxon-
Mann-Whitney test of statistical significance (set at a 1% level of
significance).

Because multiple statistical prioritization techniques were
employed, we report the p -values — as the number of the exe-
cutions increases, p becomes sufficiently small [15], which means
that there are differences between the two algorithms. We used the

1 non-parametric Vargha and Delaney effect size measure [51], $A_{12}(x, y)$, where the further away from 0.5, the larger the effect size. The effect size can be also represented as the probability that one technique is better than another — with a higher effect size (value) indicating higher probability. For example, $A_{12}(x, y) = 1.0$ indicates that, based on the sample, algorithm x always performs better than algorithm y ; and $A_{12}(x, y) = 0.0$ means that it always has worse performance. Based on the classification [51], the effect size is categorized into one of three degrees — Small (S), Medium (M), or

Large (L) — where Small means $A_{12}(x, y)$ means 0.1 $A_{12}(x, y) < 0.5$; and Large means 0.17 $A_{12}(x, y) > 0.5$. The p-value and effect size A_{12} value were calculated for each pair-wise comparison of the prioritization techniques.

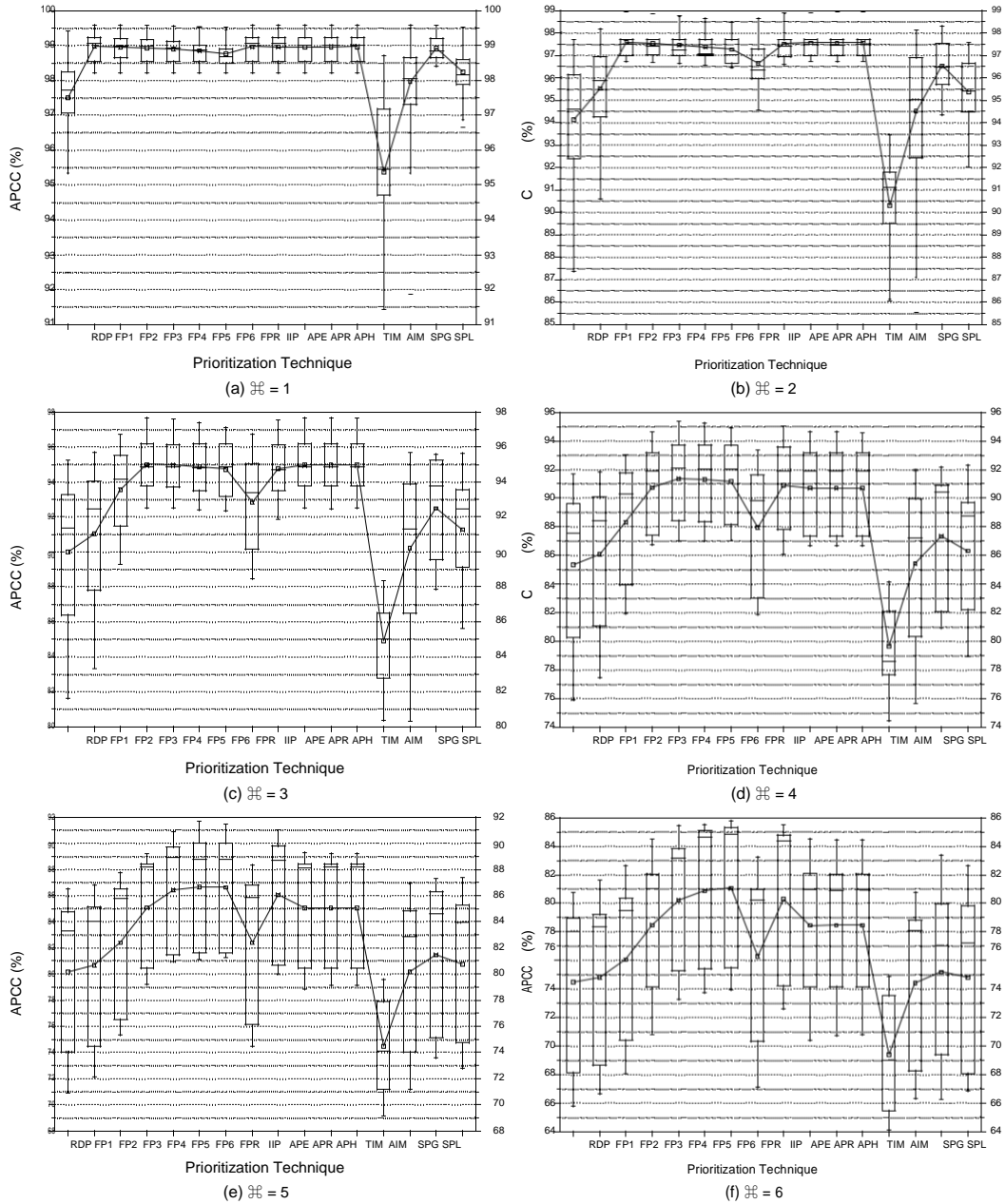


Fig. 2: APCCs for each ATPC technique for the five subject programs

4 Results

In the plots in Figures 2, 3, and 4, the X-axis shows the ATPC techniques investigated, while the Y-axis lists the APCC or APFD values. Each box plot describes the mean (a square in the box), median (a line in the box), lower/upper quartile, and min/max APCC or APFD values.

4.1 APCC Results

Figure 2 presents the APCC results at different λ (1 $\leq \lambda \leq 6$) values. Figure 3 gives the average APCC values over the six λ values, in which each plot describes the distribution of the 500 APCC values (100 orderings \times 5 programs) at λ . Table 4 shows the statistical results for comparing any two techniques based on Figure 3.

4.1.1 RQ1: APCC Effectiveness: ICBP: Regarding the FICBP techniques (the first subquestion of RQ1), the $FP\lambda$ (1 $\leq \lambda \leq 6$) generally has the best APCCs at λ (1 $\leq \lambda \leq 6$), when λ is equal to the prioritization strength λ . However, not every $FP\lambda$ always performs best at prioritization strength λ , because local optimization instead of global optimization was applied. In other words, no FICBP method always has the highest APCC values. These observations are consistent with those reported in other studies [21, 24, 38]. Furthermore, at a fixed λ (1 $\leq \lambda \leq 6$), when λ increases, $FP\lambda$ achieves higher APCC while 1 $\leq \lambda \leq 6$; but lower APCC when $\lambda \leq 6$. According to the average APCC over the six values of λ (Figure 3), $FP4$, $FP6$, and $FP5$ are the three best FICBP techniques, followed by $FP3$, and $FP2$; and $FP1$ performs worst. Table 4 shows the APCC inferential statistical analysis, which confirms the box plot results. As a consequence, the prioritization strength λ should be assigned a value of at least 4, if we wish to achieve the best performance (according to the interaction coverage rate).

Regarding the AICBP techniques (the second subquestion of RQ1), all three weighting distributions of prioritization strengths have very similar APCC values, irrespective of λ and program. According to the statistical analysis, the p-values for comparisons between any two techniques is greater than 0.01; and the effect size measure A_{12} is approximately equal to 50%, which confirms the plot observations. Therefore, the weighting distribution has only a very slight impact on the AICBP techniques.

To answer the last subquestion of RQ1, we compared all eleven ICBP techniques (FP_i ($i = 1, 2, 3, 4, 5, 6$), FPR , APE , APR , APH , and $IICBP$). Based on this comparison, we observe the following:

• When $\lambda = 1$ (Figure 2(a)), $FP6$ and $FP5$ have the worst performance, and the other nine ICBP techniques perform similarly (with

$FP1$, FPR , and $IICBP$ having slightly higher APCC results than others).

• When $\lambda = 2$ (Figure 2(b)), $FP1$ is worst, followed by FPR , $FP6$, and $FP5$; and all other techniques have similar APCC results.

• When $\lambda = 3$ (Figure 2(c)), $FP1$ and FPR have the worst ICBP performance, followed by $FP2$; and all other techniques are similar.

• When $\lambda = 4, 5, 6$, $FP4$, $FP5$, and $FP6$ generally have the high-est APCC values, followed by $IICBP$. $FP1$, $FP2$, and FPR generally perform worst.

• Based on the average APCC results, $FP4$, $FP5$, and $FP6$ are the three best ICBP techniques, followed by $IICBP$. The next best techniques are $FP3$, and the AICBP series. $FP1$ is worst, followed by FPR and $FP2$. The statistical analysis also confirms these observations.

4.1.2 Q2: APCC Effectiveness: IMBP: Based on the experimental data, it is clear that AIM has much higher APCC values than TIM , regardless of λ values. Therefore, AIM also has much higher average APCC values, which is confirmed by the statistical analysis: the p-value is less than $2.04E-72$, indicating a significant difference

between them; and the effect size measure A_{12} is 0.1712, indicating

that AIM performs better than TIM about 83% of the time.

4.1.3 RQ3: APCC Effectiveness: SBP: When 1 $\leq \lambda \leq 4$, SPG has significantly better λ -wise APCC results than SPL ; however, when $\lambda = 5, 6$, SPG is better than SPL , although the differences are small. Based on the statistical analysis, it is clear that SPG performs significantly better than SPL : their p-value is much

less than 0.01, which indicates high significance; and their A_{12} is

0.6927, which means that SPG has a probability of about 69% of obtaining higher APCC values than SPL .

4.1.4 RQ4: APCC Effectiveness of All Techniques: Considering all sixteen ATPC techniques, we can observe that as λ (1 $\leq \lambda \leq 6$) increases, the APCC values of each prioritization technique decrease, which is expected, due to the characteristics of the APCC metric (Section 3.4.1). More specifically, given a candidate ATPC set T , the number of λ -wise value combinations covered by T is generally much larger than that of λ^0 -wise value combinations, when

1 $\leq \lambda < 6$. In other words, the number of λ -wise value combinations covered by T increases as λ increases. For each prioritized set S of T , therefore, the speed of covering λ^0 -wise value combinations is faster than that of covering λ -wise value combinations: $APCC(S, \lambda^0) > APCC(S, \lambda)$.

Among all techniques, TIM generally has the worst performance: this is a surprising result, because it performs worse than RDP , which does not use any information to guide the prioritization process. Additionally, the ICBP series has better APCC results than any other series, such as $NIGP$, $IMBP$, and SBP ; with SBP as the second best (it should be noted that SPG is better than $FP1$), followed by $IMBP$. This observation is also understandable, because the ICBP series uses the interaction coverage information to guide the prioritization, giving higher interaction coverage rates. In addition, the SBP series does not use interaction coverage for prioritizing ATPCs, but the similarity comparison between two test cases effectively achieves this interaction coverage: guaranteeing that at least two test cases could cover the largest number of value combinations at strength 1. However, the $IMBP$ series prioritizes test cases according to the model mutation scores, and hence no interaction coverage is considered for the prioritization.

To conclude, ICBP is the best, with fixed-strength ICBP at higher prioritization strength λ giving the best APCC scores (it is recommended that λ be assigned a value of at least 4), and incremental-strength and aggregate-strength ICBP delivering comparable APCC results. SBP is the second best, with the global SBP achieving APCC results comparable to the ICBP series, and better than the local SBP . A surprising result is that $NIGP$ (such as RDP) could sometimes achieve better performance than $IMBP$, according to the APCC values.

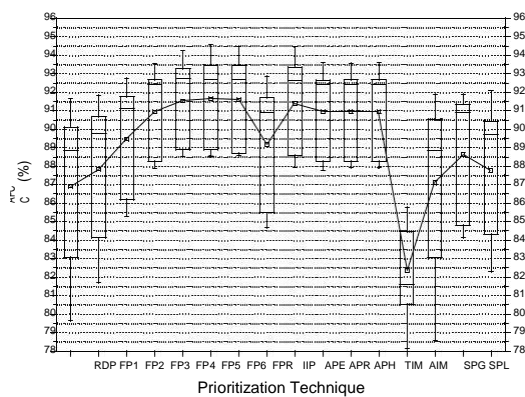


Fig. 3: Average APCC distribution (for $\lambda = 1, 2, 3, 4, 5, 6$) for each ATPC technique

Table 4 Statistical APCC analysis of all pairwise comparisons (A, B)

A	B	p-value	Superior	Effect Size	A	B	p-value	Superior	Effect Size
RDP	FP1	8.50E-15	FP1	0.3583 (M)	FP4	APH	6.96E-22	FP4	0.6756 (L)
RDP	FP2	3.56E-44	FP2	0.2453 (L)	FP4	TIM	5.86E-165	FP4	1.0000 (L)
RDP	FP3	1.30E-48	FP3	0.2324 (L)	FP4	AIM	2.22E-59	FP4	0.7968 (L)
RDP	FP4	2.68E-58	FP4	0.2060 (L)	FP4	SPG	5.67E-46	FP4	0.7600 (L)
RDP	FP5	2.24E-58	FP5	0.2058 (L)	FP4	SPL	1.05E-46	FP4	0.7621 (L)
RDP	FP6	3.59E-56	FP6	0.2116 (L)	FP5	FP6	0.0059	FP5	0.5503 (S)
RDP	FPR	6.19E-42	FPR	0.2522 (L)	FP5	FPR	3.70E-42	FP5	0.7485 (L)
RDP	IIP	3.52E-51	IIP	0.2251 (L)	FP5	IIP	3.77E-10	FP50.6144 (M)	
RDP	APE1.27E-48	APE	0.2323 (L)	FP5	APE	1.29E-27	FP5	0.6989 (L)	
RDP	APR 1.15E-48	APR	0.2322 (L)	FP5	APR	1.01E-28	FP5	0.7031 (L)	
RDP	APH 1.22E-48	APH	0.2323 (L)	FP5	APH	5.25E-28	FP5	0.7004 (L)	
RDP	TIM 2.97E-71	RDP	0.8260 (L)	FP5	TIM	5.86E-165	FP5	1.0000 (L)	
RDP	AIM 0.0273	AIM	0.4597 (S)	FP5	AIM	1.99E-59	FP5	0.7970 (L)	
RDP	SPG 1.15E-41	SPG	0.2530 (L)	FP5	SPG	5.67E-46	FP5	0.7600 (L)	
RDP	SPL 3.03E-11	SPL	0.3786 (M)	FP5	SPL	1.01E-46	FP5	0.7622 (L)	
FP1	FP2	7.62E-40	FP2	0.2587 (L)	FP6	FPR	2.39E-42	FP6	0.7491 (L)
FP1	FP3	5.67E-46	FP3	0.2400 (L)	FP6	IIP	6.19E-08	FP6	0.5989 (S)
FP1	FP4	3.14E-46	FP4	0.2392 (L)	FP6	APE	4.90E-28	FP6	0.7005 (L)
FP1	FP5	3.06E-46	FP5	0.2392 (L)	FP6	APR	4.66E-29	FP6	0.7044 (L)
FP1	FP6	3.07E-46	FP6	0.2392 (L)	FP6	APH	2.08E-28	FP6	0.7019 (L)
FP1	FPR	2.13E-34	FPR	0.2766 (L)	FP6	TIM	5.86E-165	FP6	1.0000 (L)
FP1	IIP	4.61E-46	IIP	0.2397 (L)	FP6	AIM	3.61E-57	FP6	0.7911 (L)
FP1	APE	5.67E-46	APE	0.2400 (L)	FP6	SPG	5.67E-46	FP6	0.7600 (L)
FP1	APR	5.67E-46	APR	0.2400 (L)	FP6	SPL	1.83E-46	FP6	0.7614 (L)
FP1	APH	5.67E-46	APH	0.2400 (L)	FPR	IIP	6.53E-40	IIP	0.2585 (L)
FP1	TIM	1.26E-93	FP1	0.8749 (L)	FPR	APE	1.34E-29	APE	0.2936 (L)
FP1	AIM	5.82E-07	FP1	0.5913 (S)	FPR	APR	1.50E-29	APR	0.2938 (L)
FP1	SPG	6.60E-29	SPG	0.2962 (L)	FPR	APH	1.08E-29	APH	0.2933 (L)
FP1	SPL	0.3229	FP1	0.5181 (S)	FPR	TIM	1.61E-155	FPR	0.9853 (L)
FP2	FP3	4.83E-40	FP3	0.2581 (L)	FPR	AIM	1.45E-35	FPR	0.7274 (L)
FP2	FP4	1.56E-45	FP4	0.2413 (L)	FPR	SPG	1.39E-12	FPR0.6294 (M)	
FP2	FP5	1.63E-45	FP5	0.2413 (L)	FPR	SPL	1.10E-33	FPR	0.7210 (L)
FP2	FP6	8.94E-46	FP6	0.2406 (L)	IIP	APE	4.99E-11	IIP	0.6200 (M)
FP2	FPR	8.67E-06	FP2	0.5813 (S)	IIP	APR	3.57E-11	IIP	0.6209 (M)
FP2	IIP	9.21E-45	IIP	0.2436 (L)	IIP	APH	4.46E-11	IIP	0.6203 (M)
FP2	APE	4.11E-40	APE	0.2578 (L)	IIP	TIM	5.86E-165	IIP	1.0000 (L)
FP2	APR	3.12E-40	APR	0.2575 (L)	IIP	AIM	3.22E-51	IIP	0.7750 (L)
FP2	APH	6.14E-40	APH	0.2584 (L)	IIP	SPG	5.67E-46	IIP	0.7600 (L)
FP2	TIM	8.24E-165	FP2	0.9998 (L)	IIP	SPL	3.67E-46	IIP	0.7606 (L)
FP2	AIM	5.92E-40	FP2	0.7417 (L)	APE	APR	0.9296	APE	0.5016 (S)
FP2	SPG	1.20E-19	FP2	0.6657 (M)	APE	APH	0.9936	APH	0.4999 (S)
FP2	SPL	4.43E-38	FP2	0.7357 (L)	APE	TIM	5.86E-165	APE	1.0000 (L)
FP3	FP4	1.56E-21	FP4	0.3259 (L)	APE	AIM	9.65E-48	APE	0.7652 (L)
FP3	FP5	1.93E-27	FP5	0.3018 (L)	APE	SPG	5.67E-46	APE	0.7600 (L)
FP3	FP6	9.68E-28	FP6	0.3006 (L)	APE	SPL	5.62E-46	APE	0.7600 (L)
FP3	FPR	1.16E-29	FP3	0.7066 (L)	APR	APH	0.9368	APH	0.4986 (S)
FP3	IIP	1.74E-10	IIP	0.3834 (M)	APR	TIM	5.86E-165	APR	1.0000 (L)
FP3	APE	0.9989	FP3	0.5001 (S)	APR	AIM	9.87E-48	APR	0.7651 (L)
FP3	APR	0.9263	FP3	0.5017 (S)	APR	SPG	5.67E-46	APR	0.7600 (L)
FP3	APH	0.9816	FP3	0.5004 (S)	APR	SPL	5.69E-46	APR	0.7600 (L)
FP3	TIM	5.86E-165	FP3	1.0000 (L)	APH	TIM	5.86E-165	APH	1.0000 (L)
FP3	AIM	1.11E-47	FP3	0.7650 (L)	APH	AIM	9.90E-48	APH	0.7651 (L)
FP3	SPG	5.67E-46	FP3	0.7600 (L)	APH	SPG	5.67E-46	APH	0.7600 (L)
FP3	SPL	5.67E-46	FP3	0.7600 (L)	APH	SPL	5.65E-46	APH	0.7600 (L)
FP4	FP5	0.0845	FP5	0.4685 (S)	TIM	AIM	2.04E-72	AIM	0.1712 (L)
FP4	FP6	0.7963	FP6	0.4953 (S)	TIM	SPG	6.34E-121	SPG	0.0729 (L)
FP4	FPR	5.52E-43	FP4	0.7511 (L)	TIM	SPL	1.70E-94	SPL	0.1233 (L)
FP4	IIP	0.0046	FP4	0.5518 (S)	AIM	SPG	7.24E-33	SPG	0.2819 (L)
FP4	APE	1.45E-21	FP4	0.6742 (L)	AIM	SPL	2.38E-05	SPL	0.4228 (S)
FP4	APR	2.38E-22	FP4	0.6776 (L)	SPG	SPL	5.18E-26	SPG	0.6927 (L)

S, M, and L represents Small, Medium, and Large in effect size, respectively.

1 4.2 APFD Results

2 Figure 4 presents the APFD results for each subject program (Fig-
3 ures 4(a) to 4(e)), in which each plot lists the distribution of the
4 500 APFD values (100 orderings * 5 versions). Figure 4(f) gives
5 the APFD results for all programs, in which each plot contains
6 2500 APFD values (100 orderings * 5 programs * 5
7 versions).

8 Table 5 shows the statistical APFD comparisons between two ATCP
9 techniques based on Figure 4(f).

10 4.2.1 RQ1: APFD Effectiveness: ICBP: To answer the first
11 subquestion of RQ1, regarding FICBP, we have the following
12 observations:

- As the prioritization strength λ ($1 \leq \lambda \leq 6$) increases, $F\lambda$ can normally achieve higher APFD results, with a few exceptions: for example, in program grep, FP2 performs better than FP3; while FP4 performs worst for program make.
- According to mean and median APFD values, the largest difference between techniques is only 4%, and the differences between high-strength FICBPs are very small. Lower-strength FICBPs are, therefore, surprisingly comparable to higher-strength ones, from the perspective of fault detection.
- As shown in Table 5, the comparisons between higher-strength (FP4, FP5, and FP6) and lower-strength (FP1, FP2, and FP3) FICBP are highly significant: except when comparing FP4 against FP3, the

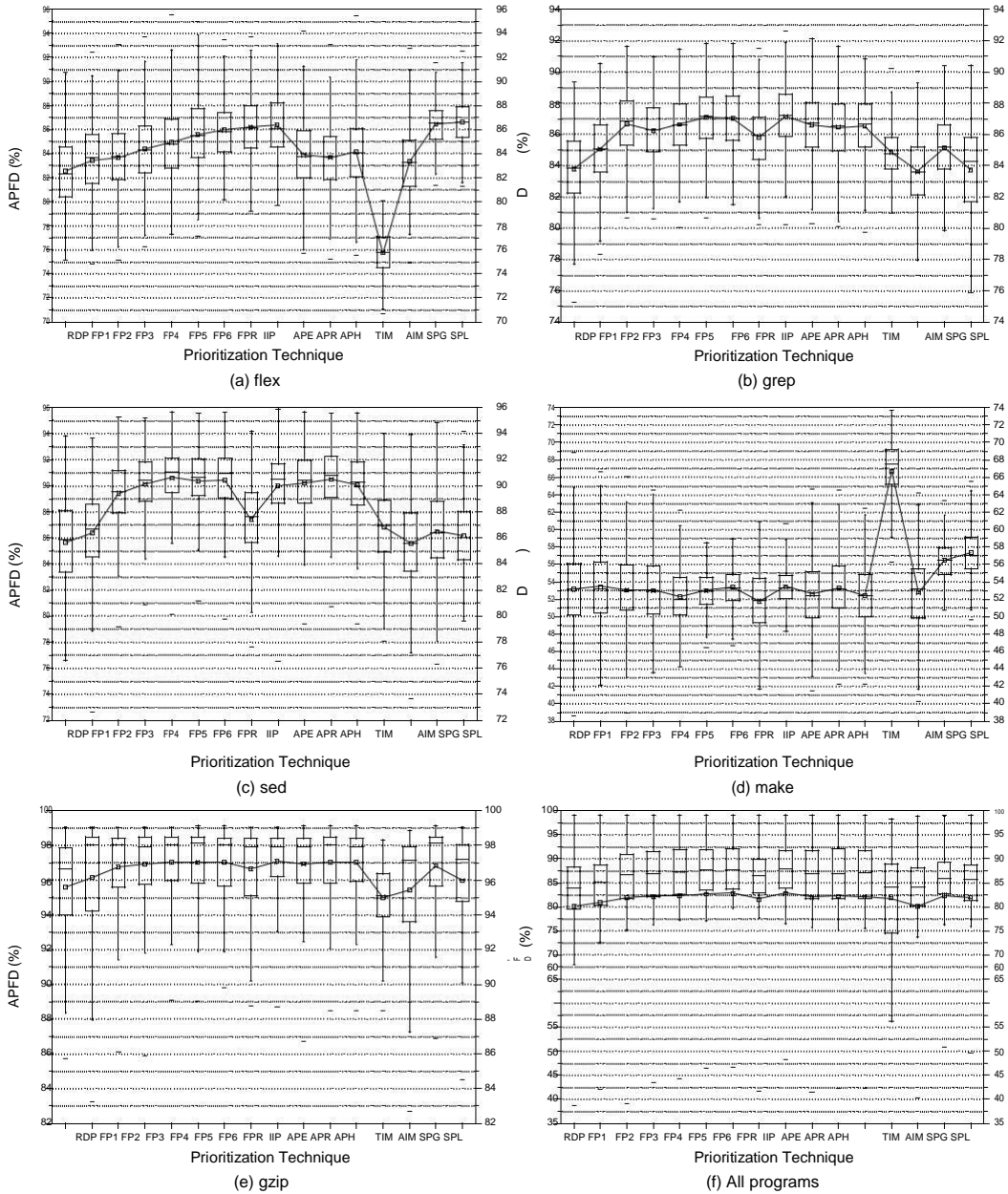


Fig. 4: APFDs for each ATPC technique for V1 to V5

1 p-values are less than 0.01. Among the higher-strength FICBPs,
 2 the APFD results are not significantly different (their p-values are
 3 greater than 0.01); but, among lower-strength FICBPs, the differ-
 4 ence is highly significant, for example, when comparing FP1 with
 5 FP2 or with FP3. In terms of the effect size measure (A 12), higher-
 6 strength FICBPs only outperform lower-strength ones between about
 7 52% and 59% of the time. Among the higher-strength FICBPs, the
 8 A 12 values range from about 50% to 52%; while they range from
 9 about 51% to 59% among the lower-strength FICBPs.

10 In answering the second subquestion, there is nearly no difference 10
 11 between the AICBPs, irrespective of subject program. This is also 11
 12 confirmed by the statistical comparison: the p-values are greater than 12
 13 0.01, and the effect size measures are approximately 50%. 13
 14 Regarding the third subquestion, among all eleven ICBP techni- 14
 15 ques, we have the following observations: IIP and higher-strength 15
 16 FICBPs generally have the highest APFD values, and IIP has bet- 16
 17 ter performance than lower-strength FICBPs. The second best is the 17
 18 AICBP series, followed by FP2, FP3, and FPR. FP1 has the worst 18

Table 5 Statistical APFD analysis of all pairwise comparisons (A, B)

A	B	p-value	Superior	Effect Size	A	B	p-value	Superior	Effect Size
RDP	FP1	1.27E-06	FP1	0.4604 (S)	FP4	APH	0.1268	FP4	0.5125(S)
RDP	FP2	7.79E-26	FP2	0.4142 (S)	FP4	TIM	4.53E-28	FP4	0.5897(S)
RDP	FP3	2.61E-30	FP3	0.4066 (S)	FP4	AIM	3.51E-38	FP4	0.6055 (M)
RDP	FP4	9.97E-38	FP4	0.3952 (M)	FP4	SPG	5.96E-05	FP4	0.5328(S)
RDP	FP5	1.03E-46	FP5	0.3828 (M)	FP4	SPL	3.31E-10	FP4	0.5513(S)
RDP	FP6	4.56E-51	FP6	0.3773 (M)	FP5	FP6	0.5841	FP6	0.4955(S)
RDP	FPR	4.40E-26	FPR	0.4137 (S)	FP5	FPR	3.20E-09	FP5	0.5483(S)
RDP	IIP	4.86E-54	IIP	0.3736 (M)	FP5	IIP	0.4432	FP5	0.4937 (M)
RDP	APE	1.38E-28	APE	0.4094 (S)	FP5	APE	0.0022	FP5	0.5250(S)
RDP	APR	2.14E-29	APR	0.4081 (S)	FP5	APR	0.0128	FP5	0.5203(S)
RDP	APH	2.29E-29	APH	0.4081 (S)	FP5	APH	0.0037	FP5	0.5237(S)
RDP	TIM	0.7419	RDP	0.5027 (S)	FP5	TIM	5.76E-33	FP5	0.5977(S)
RDP	AIM	0.6734	AIM	0.4966 (S)	FP5	AIM	6.95E-48	FP5	0.6187 (M)
RDP	SPG	8.30E-28	SPG	0.4107 (S)	FP5	SPG	2.14E-08	FP5	0.5457(S)
RDP	SPL	1.05E-16	SPL	0.4322 (S)	FP5	SPL	9.25E-15	FP5	0.5633(S)
FP1	FP2	7.49E-11	FP2	0.4468 (S)	FP6	FPR	1.23E-10	FP6	0.5526(S)
FP1	FP3	9.29E-14	FP3	0.4392 (S)	FP6	IIP	0.8092	FP6	0.4980(S)
FP1	FP4	1.37E-19	FP4	0.4261 (S)	FP6	APE	0.0004	FP6	0.5292(S)
FP1	FP5	5.83E-27	FP5	0.4122 (S)	FP6	APR	0.0025	FP6	0.5247(S)
FP1	FP6	1.17E-29	FP6	0.4076 (S)	FP6	APH	0.0006	FP6	0.5279(S)
FP1	FPR	1.39E-09	FPR	0.4505 (S)	FP6	TIM	6.45E-35	FP6	0.6007 (M)
FP1	IIP	3.74E-32	IIP	0.4036 (S)	FP6	AIM	2.06E-52	FP6	0.6244 (M)
FP1	APE	5.81E-13	APE	0.4412 (S)	FP6	SPG	1.67E-09	FP6	0.5492(S)
FP1	APR	1.21E-13	APR	0.4394 (S)	FP6	SPL	2.03E-16	FP6	0.5671(S)
FP1	APH	1.33E-13	APH	0.4395 (S)	FPR	IIP	4.50E-12	IIP	0.4435(S)
FP1	TIM	0.0002	FP1	0.5300 (S)	FPR	APE	0.0255	APE	0.4818(S)
FP1	AIM	2.16E-06	FP1	0.5387 (S)	FPR	APR	0.0094	APR	0.4788(S)
FP1	SPG	2.69E-09	SPG	0.4514 (S)	FPR	APH	0.0142	APH	0.4800(S)
FP1	SPL	0.0007	FP1	0.4722 (S)	FPR	TIM	1.06E-17	FPR	0.5700(S)
FP2	FP3	0.3028	FP3	0.4916 (S)	FPR	AIM	1.26E-26	FPR	0.5872(S)
FP2	FP4	0.0080	FP4	0.4783 (S)	FPR	SPG	0.6340	FPR	0.5039 (M)
FP2	FP5	5.18E-05	FP5	0.4669 (S)	FPR	SPL	0.0036	FPR	0.5237(S)
FP2	FP6	5.24E-06	FP6	0.4628 (S)	IIP	APE	0.0001	IIP	0.5314 (M)
FP2	FPR	0.1867	FP2	0.5108 (S)	IIP	APR	0.0012	IIP	0.5265 (M)
FP2	IIP	1.32E-06	IIP	0.4605 (S)	IIP	APH	0.0002	IIP	0.5303 (M)
FP2	APE	0.3601	APE	0.4925 (S)	IIP	TIM	4.73E-37	IIP	0.6038 (M)
FP2	APR	0.1785	APR	0.4890 (S)	IIP	AIM	2.34E-55	IIP	0.6280 (M)
FP2	APH	0.2857	APH	0.4913 (S)	IIP	SPG	2.69E-11	IIP	0.5544(S)
FP2	TIM	3.34E-18	FP2	0.5710 (S)	IIP	SPL	9.54E-19	IIP	0.5722(S)
FP2	AIM	7.26E-26	FP2	0.5859 (S)	APE	APR	0.6566	APE	0.4964(S)
FP2	SPG	0.1543	FP2	0.5116 (M)	APE	APH	0.8601	APH	0.4986(S)
FP2	SPL	0.0002	FP2	0.5305 (S)	APE	TIM	1.41E-21	APE	0.5779(S)
FP3	FP4	0.1071	FP4	0.4868 (S)	APE	AIM	1.26E-28	APE	0.5906(S)
FP3	FP5	0.0030	FP5	0.4757 (S)	APE	SPG	0.0283	APE	0.5179(S)
FP3	FP6	0.0005	FP6	0.4714 (S)	APE	SPL	8.04E-06	APE	0.5365(S)
FP3	FPR	0.0196	FP3	0.5191 (S)	APR	APH	0.7745	APH	0.5023(S)
FP3	IIP	0.0002	IIP	0.4693 (S)	APR	TIM	3.20E-21	APR	0.5772(S)
FP3	APE	0.9032	FP3	0.5010 (S)	APR	AIM	1.94E-29	APR	0.5920(S)
FP3	APR	0.7688	FP3	0.4976 (S)	APR	SPG	0.0210	APR	0.5188(S)
FP3	APH	0.9618	FP3	0.4996 (S)	APR	SPL	5.76E-06	APR	0.5370(S)
FP3	TIM	1.41E-21	FP3	0.5779 (S)	APH	TIM	2.62E-22	APH	0.5793(S)
FP3	AIM	1.47E-30	FP3	0.5938 (S)	APH	AIM	1.83E-29	APH	0.5920(S)
FP3	SPG	0.0202	FP3	0.5190 (S)	APH	SPG	0.0200	APH	0.5190(S)
FP3	SPL	4.91E-06	FP3	0.5373 (S)	APH	SPL	4.60E-06	APH	0.5374(S)
FP4	FP5	0.1792	FP5	0.4890 (S)	TIM	AIM	0.6916	AIM	0.4968(S)
FP4	FP6	0.0595	FP6	0.4846 (S)	TIM	SPG	2.04E-13	SPG	0.4400(S)
FP4	FPR	2.71E-05	FP4	0.5343 (S)	TIM	SPL	2.14E-08	SPL	0.4543(S)
FP4	IIP	0.0344	FP4	0.4827 (S)	AIM	SPG	1.03E-28	SPG	0.4092(S)
FP4	APE	0.0880	FP4	0.5139 (S)	AIM	SPL	7.39E-17	SPL	0.4319(S)
FP4	APR	0.2311	FP4	0.5098 (S)	SPG	SPL	0.0266	SPG	0.5181(S)

1 performance. It is surprising that FPR has APFD results comparable
2 to FP2 and FP3, and has higher APFD scores than FP1, because it
3 only repeats 1-wise interaction coverage. Overall, the statistical anal-
4 ysis (see Table 5) supports the box plot observations, with a degree of
5 variation in the performance of different ICBP techniques for differ-
6 ent programs. Nevertheless, based on the programs we have studied,
7 our results suggest that IIP and higher-strength FICBPs offer the best
8 rates of fault detection among the ICBP techniques.

9 4.2.2 RQ2: APFD Effectiveness: IMBP: For subject pro-
10 grams flex and gzip, AIM performs significantly better than TIM,
11 with respect to both the mean and median APFD values. However,
12 for the other three programs (grep, sed, and make), TIM achieves

13 much better APFD results (again from the perspective of both mean
14 and median APFD values). This is especially so for the program
15 make, where the mean APFD for TIM is close to 67%, but for AIM
16 it is only about 53%; the median APFD for TIM is 67.5%, but the
17 AIM median is also only about 53%. In contrast to previous TCP
18 studies [1, 31], an interesting result is that the 'additional' TCP tech-
19 niques do not guarantee to provide better fault detection rates than
20 the 'total' TCP techniques.

21 However, the statistical analysis for all five programs suggests that
22 the differences in performance between TIM and AIM are not sig-
23 nificant: their p-values are much greater than 0.01, and the effect
24 sizes are approximately 50%. In other words, TIM and AIM have
25 comparable fault detection rates.

1 4.2.3 RQ3: APFD Effectiveness: SBP: For programs flex and
 2 make, SPG performs slightly better than SPL, however for the other
 3 programs (grep, sed, and gzip) SPG is better than SPL, with respect
 4 to both the mean and the median APFD values. Considering all pro-
 5 grams (Figure 4(f)), overall, SPG is slightly better than SPL, but the
 6 differences between them are less than 1%. Similarly, the statistical
 7 comparison gives a p-value of 0.0266, and an effect size of 0.5181,
 8 which indicates that the difference is not significant.

9 4.2.4 RQ4: APFD Effectiveness of All Techniques:
 10 Although different techniques have different APFD performances
 11 for different programs, we can nonetheless observe the following:

- 12 • For program flex, SPG and SPL are the two best techniques, fol-
 13 lowed by IIP and FPR, in terms of both the median and the mean
 14 APFD values — although the differences are very small (less than
 15 1%). Additionally, and surprisingly, TIM has the worst performance
 16 — even worse than RDP, which uses no additional information in
 17 the prioritization process.
- 18 • For program make, TIM is significantly better than all other
 19 ATCP techniques, followed by SPG and SPL. Additionally, RDP,
 20 $FP\lambda$ ($1 \leq \lambda \leq 6$; except $\lambda = 4$), IIP, APR, and AIM, all have simi-
 21 lar APFD performance. FP4 and FPR are the two worst techniques.
- 22 • For the three programs grep, sed, and gzip, the FICBPs (except
 23 FP1) and other ICBP techniques perform best, with FP1, SPG, SPL,
 24 TIM, and AIM able to achieve comparable APFD results. Further-
 25 more, it is again surprising that RDP could sometimes have similar
 26 fault detection rates to TIM, AIM, and SPL.
- 27 • When all programs are considered together, overall, the ICBP
 28 series has the best performance, followed by the SBP series; NIGP
 29 and IMBP perform worst, with similar fault detection rates. regard-
 30 ing individual ATCP techniques, the ICBP series is best, as discussed
 31 in the first subquestion of RQ1, with IIP and higher-strength FICBPs
 32 performing best among all techniques. SPG and SPL are better
 33 than AIM, TIM, and FP1; with SPG achieving comparable APFD
 34 performance to FP2, FP3, FPR, and the AICBP series.

35 Taking into consideration both APCC and APFD results, we can
 36 conclude that higher-strength FICBPs (FP4, FP5, and FP6) achieve
 37 the best rates of both interaction coverage and fault detection, fol-
 38 lowed by IIP. Although SPG has lower APCC results than the ICBP
 39 techniques (as discussed before, because ICBP uses interaction cov-
 40 erage to guide the prioritization), it can achieve higher APFD scores
 41 than FP1, and has performance comparable to FP2, FP3, FPR, and
 42 AICBP. Additionally, IMBP techniques generally perform similarly,
 43 or worse, compared with random test case prioritization.

44 4.3 Prioritization Time Results

45 To address RQ5, Table 6 presents the mean prioritization time for
 46 each ATCP technique for each subject program — it should be noted
 47 that, because we used the model mutation matrix file from previous
 48 studies [15], TIM and AIM do not include the model mutation time.
 49 Based on the experimental data, and as expected, it is clear that RDP
 50 needs the least prioritization time among all ATCP techniques, fol-
 51 lowed by TIM, FP1, and AIM. The next best performance, in terms
 52 of prioritization time, is by FPR, SPG, and SPL (all of which require
 53 slightly more time than the four best techniques). FP5 has the slowest
 54 prioritization time, followed by FP6, FP4, and IIP; and the AICBP
 55 series has similar times to FP3.

56 Based on the effectiveness and efficiency experiments, our recom-
 57 mendations and guidelines are as follows: given sufficient resources
 58 (including time) for prioritizing ATCs, FICBP λ at higher strength
 59 λ values ($\lambda = 4, 5, 6$) should be the best choice, followed by IIP.
 60 However, if time resources are limited, then FPR and SPG would be
 61 the best choices, followed by FP2, FP3, and the AICBP series; FP1
 62 and RDP could also be alternatives, when facing very severe time
 63 constraints. As discussed in Section 2.2.5, we believe that our exper-
 64 imental results are basically consistent with the expected strengths
 65 and weaknesses of each ATCP technique.

Table 6 Prioritization time (in seconds) for each ATCP technique

ATCP Technique	Subject Program					Sum
	flex	grep	sed	make	gzip	
RDP	0.05	0.01	0.04	0.01	0.01	0.12
FP1	0.28	0.47	0.41	0.06	0.16	1.38
FP2	4.37	8.99	10.35	0.42	2.35	26.48
FP3	10.68	26.66	36.42	2.18	27.70	103.64
FP4	52.74	81.19	144.54	6.70	115.86	401.03
FP5	84.91	198.93	339.66	18.72	326.63	968.85
FP6	59.08	108.67	217.38	16.87	518.15	920.15
FPR	2.96	2.23	1.46	0.50	1.36	8.51
IIP	54.14	40.63	41.34	16.56	168.70	321.37
APE	12.32	29.51	40.05	3.94	43.94	129.76
APR	12.92	30.12	40.88	4.07	42.61	130.60
APH	12.94	29.99	40.08	4.13	43.20	130.34
TIM*	0.38	0.79	0.08	0.05	0.03	1.33
AIM*	1.36	1.89	0.13	0.10	0.03	3.51
SPG	3.78	3.23	1.73	0.19	0.49	9.42
SPL	3.84	2.02	1.50	0.17	0.28	7.81

*** indicates that model mutation time is not included.

4.4 Threats to Validity

In this section, we list some potential threats to validity, including external validity, internal validity, construct validity, and conclusion validity.

4.4.1 External Validity: With respect to the external validity, the main threat is the generalizability of our results. Although we have used only five subject programs, written in C, all of which are of a relatively medium size, we believe that by including six versions of each (giving 30 subject versions under study), that there is sufficient data from which to draw the conclusions. Nevertheless, more larger subject programs, written in other languages should also be examined in future work.

Another potential threat to external validity is the representativeness of ATCs for each subject program. In this paper, we focused on ATCs originated from the SIR [40] (using the test specification language to create the input model and construct ATCs [16]), which is only one type of ATC encoding. However, there exist other ATC encoding types [52], which we will investigate in our future work.

4.4.2 Internal Validity: The threat to internal validity relates mainly the implementation of our algorithms. We have used C++ to implement the algorithms, and have carefully tested the implementation to minimize this threat, as much as possible.

4.4.3 Construct Validity: In this study, we have focused on the testing effectiveness and efficiency, measured by the rate of interaction coverage, the rate of fault detection, and the prioritization time. Although the APCC and APFD metrics have often been used in the field of test case prioritization [1, 21, 24, 34], we acknowledge that there may be other metrics which may also be relevant.

4.4.4 Conclusion Validity: As for the conclusion validity, the main threat is the randomized computation of our algorithms. To minimize this threat, all algorithms were repeated 100 times, and inferential statistics were applied to the comparisons of results.

5 Conclusions and Future Work

This paper has reported on a comparison of 16 ATCP techniques, classified into four categories, based on an extensive empirical study. Based on comparisons of testing effectiveness and efficiency, some recommendations and guidelines have also been given, to help testers choose among ATCP techniques under different testing situations and scenarios.

The main findings of this study can be summarized as:

- (1) With respect to all ATCP categories, the ICBP category has the best testing effectiveness, irrespective of the rates of interaction coverage and fault detection. Somewhat surprisingly, because it does not use any additional information to guide the prioritization, NIGP

could achieve comparable performance to IMBP; while SBP has very good testing effectiveness, and even better than some ICBP techniques sometimes. Additionally, IMBP has the worst rates of interaction coverage, but it sometimes has the best fault detection rates. Nevertheless, NIGP, IMBP, SBP, and some ICBP techniques have better testing efficiency than others.

(2) In the category of ICBP techniques, it is evident that higher-strength FICBP techniques, and IICBP have the best testing effectiveness (according to interaction coverage and fault detection), followed by AICBP and lower-strength FICBP techniques. However, higher-strength FICBP and IICBP techniques are less efficient than other ICBP techniques, according to the prioritization time.

(3) Regarding the IMBP techniques, although both 'total' and 'additional' IMBP techniques have similar prioritization times, they have different performances according to the other evaluation measures. For example, the 'additional' IMBP has better rates of interaction coverage than the 'total' IMBP, regardless of subject programs. However, for three programs, the 'additional' IMBP has better fault detection than the 'total' IMBP, but for another two cases, this is reversed: the 'total' IMBP can obtain better fault detection.

(4) For the SBP techniques, the global SBP has better rates of interaction coverage than the local SBP. However, they have similar fault detection rates and prioritization costs: the global SBP is slightly better than the local one for some programs, but the opposite is the case for some other programs.

(5) When testers select only some ATCP techniques for prioritizing abstract test cases, we recommend that, given sufficient resources and prioritization time, FICBP λ at higher strength λ values (i.e., $\lambda = 4, 5, 6$) should be the best choice, followed by IICBP. However, if facing limited time resources, then GSBP may be the best choice, followed by FICBP2, FICBP3, and AICBP; FICBP1 and NIGP may be alternatives in situations with very severe time constraints.

As discussed before, IMBP uses the model mutation information to prioritize ATCs, so the quality of IMBP is mainly dependent on the model mutation, which may be a reason for the ineffectiveness of IMBP in this study. It will therefore be very interesting to investigate the correlation between model mutation and program mutation in our future work. In addition, since this study adopted mutation analysis [39] to investigate testing effectiveness of ATCP techniques, more experiments with real faults should be conducted to validate our conclusions. Last but not the least, in this paper we only considered the prioritization time as the resource factor for guiding the selection of ATCP techniques. However, there are many other resource factors such as the execution time of test cases. Therefore, it would be interesting to combine more testing requirements for designing more comprehensive guidelines to select ATCP techniques.

Acknowledgements

We would like to thank Christopher Henard for sharing us the fault data materials of each subject program. This work is supported by the National Natural Science Foundation of China (Grant Nos. 61202110, 61502205, and 61872167), and the Senior Personnel Scientific Research Foundation of Jiangsu University (Grant No. 14JDG039). This work is also supported by the Young Backbone Teacher Cultivation Project of Jiangsu University, and the sponsorship of Jiangsu Overseas Visiting Scholar Program for University Prominent Young & Middle-aged Teachers and Presidents. A preliminary version of this paper was presented at the 39th International Conference on Software Engineering (ICSE'17) [53].

6 References

1 Rothermel, G., Untch, R.H., Chu, C., Harrold, M.J.: 'Prioritizing test cases for regression testing', IEEE Transactions on Software Engineering, 2001, 27, (10), pp. 929–948

2 Yoo, S., Harman, M.: 'Regression testing minimization, selection and prioritization: A survey', Software Testing, Verification and Reliability, 2012, 22, (2), pp. 67–120

1 pp. 192–201

2 Bryce, R.C., Memon, A.M.: 'Test suite prioritization by interaction coverage'. In: Proceedings of the Workshop on Domain Specific Approaches to Software Test Automation (DoSTA'07). (, 2007, pp. 1–7

3 Bryce, R.C., Sampath, S., Memon, A.M.: 'Developing a single model and test prioritization strategies for event-driven software', IEEE Transactions on Software

3 Di Nardo, D., Alshahwan, N., Briand, L., Labiche, Y.: 'Coverage-based regression test case selection, minimization and prioritization: a case study on an industrial system', Software Testing, Verification and Reliability, 2015, 25, (4), pp. 371–396

4 Li, Z., Harman, M., Hierons, R.M.: 'Search algorithms for regression test case prioritization', IEEE Transactions on Software Engineering, 2007, 33, (4), pp. 225–237

5 Parejo, J.A., Sánchez, A.B., Segura, S., Ruiz-Cortés, A., Lopez-Herrejon, R.E., Eged, A.: 'Multi-objective test case prioritization in highly configurable systems: A case study', Journal of Systems and Software, 2016, 122, pp. 287–310

6 Chen, J., Zhu, L., Chen, T.Y., Towey, D., Kuo, F.C., Huang, R., et al.: 'Test case prioritization for object-oriented software: An adaptive random sequence approach based on clustering', Journal of Systems and Software, 2018, 135, pp. 107–125

7 Jiang, B., Zhang, Z., Chan, W.K., Tse, T.H.: 'Adaptive random test case prioritization'. In: Proceedings of the 24th IEEE/ACM International Conference on Automated Software Engineering (ASE'09). (, 2009, pp. 233–244

8 Zhang, X., Chen, T.Y., Liu, H.: 'An application of adaptive random sequence in test case prioritization'. In: Proceedings of the 26th International Conference on Software Engineering and Knowledge Engineering (SEKE'14). (, 2014, pp. 126–131

9 Zhang, X., Xie, X., Chen, T.Y.: 'Test case prioritization using adaptive random sequence with category-partition-based distance'. In: Proceedings of the 16th IEEE International Conference on Software Quality, Reliability and Security (QRS'16). (, 2016, pp. 374–385

10 Fang, C., Chen, Z., Wu, K., Zhao, Z.: 'Similarity-based test case prioritization using ordered sequences of program entities', Software Quality Journal, 2014, 22, (2), pp. 335–361

11 Noor, T.B., Hemmati, H.: 'A similarity-based approach for test case prioritization using historical failure data'. In: Proceedings of the 26th International Symposium on Software Reliability Engineering (ISSRE'15). (, 2015, pp. 58–68

12 Catal, C., Mishra, D.: 'Test case prioritization: a systematic mapping study', Software Quality Journal, 2013, 21, (3), pp. 445–478

13 Khatibsyarhini, M., Isa, M.A., Jawawi, D.N.A., Tumeng, R.: 'Test case prioritization approaches in regression testing: A systematic literature review', Information and Software Technology, 2018, 93, pp. 74–93

14 Grindal, M., Lindström, B., Offutt, J., Andler, S.F.: 'An evaluation of combination strategies for test case selection', Empirical Software Engineering, 2006, 11, (4), pp. 583–611

15 Henard, C., Papadakis, M., Harman, M., Jia, Y., Traon, Y.L.: 'Comparing white-box and black-box test prioritization'. In: Proceedings of the 38th International Conference on Software Engineering (ICSE'16). (, 2016, pp. 523–534

16 Ostrand, T.J., Balcer, M.J.: 'The category-partition method for specifying and generating functional tests', Communications of the ACM, 1988, 31, (6), pp. 676–686

17 Nie, C., Leung, H.: 'A survey of combinatorial testing', ACM Computer Survey, 2011, 43, (2), pp. 11:1–11:29

18 Utting, M., Legeard, B.: 'Practical model-based testing - a tools approach', International Journal On Advances in Software, 2007, 2, (1), pp. 1–419

19 Bryce, R.C., Colbourn, C.J.: 'Prioritized interaction testing for pairwise coverage with seeding and constraints', Information and Software Technology, 2006, 48, (10), pp. 960–970

20 Huang, R., Chen, J., Towey, D., Chan, A.T.S., Lu, Y.: 'Aggregate-strength interaction test suite prioritization', Journal of Systems and Software, 2015, 99, pp. 36–51

21 Petke, J., Cohen, M.B., Harman, M., Yoo, S.: 'Practical combinatorial interaction testing: Empirical findings on efficiency and early fault detection', IEEE Transactions on Software Engineering, 2015, 41, (9), pp. 901–924

22 Al-Hajjaji, M., Thüm, T., Meinicke, J., Lochau, M., Saake, G.: 'Similarity-based prioritization in software product-line testing'. In: Proceedings of 18th International Software Product Line Conference (SPLC'14). (, 2014, pp. 197–206

23 Henard, C., Papadakis, M., Perrouin, G., Klein, J., Heymans, P., Traon, Y.L.: 'Bypassing the combinatorial explosion: Using similarity to generate and prioritize t-wise test configurations for software product lines', IEEE Transactions on Software Engineering, 2014, 40, (7), pp. 650–670

24 Petke, J., Cohen, M.B., Harman, M., Yoo, S.: 'Efficiency and early fault detection with lower and higher strength combinatorial interaction testing'. In: Proceedings of the 12th Joint Meeting on European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE'13). (, 2013, pp. 26–36

25 Papadakis, M., Henard, C., Traon, Y.L.: 'Sampling program inputs with mutation analysis: Going beyond combinatorial interaction testing'. In: Proceedings of the 7th International Conference on Software Testing, Verification and Validation (ICST'14). (, 2014, pp. 1–10

26 Zhang, Z., Zhang, J.: 'Characterizing failure-causing parameter interactions by adaptive testing'. In: Proceedings of the 20th International Symposium on Software Testing and Analysis (ISSTA'11). (, 2011, pp. 331–341

27 Cohen, M.B., Dwyer, M.B., Shi, J.: 'Constructing interaction test suites for highly-configurable systems in the presence of constraints: A greedy approach', IEEE Transactions on Software Engineering, 2008, 34, (5), pp. 633–650

28 Yilmaz, C., Dumlu, E., Cohen, M.B., Porter, A.A.: 'Reducing masking effects in combinatorial interaction testing: A feedback driven adaptive approach', IEEE Transactions on Software Engineering, 2014, 40, (1), pp. 43–66

29 Thüm, T., Apel, S., Kästner, C., Schaefer, I., Saake, G.: 'A classification and survey of analysis strategies for software product lines', ACM Computing Survey, 2014, 47, (1), pp. 6:1–6:45

30 Barus, A.C., Chen, T.Y., Kuo, F.C., Liu, H., Merkel, R., Rothermel, G.: 'A cost-effective random testing method for programs with non-numeric inputs', IEEE Transactions on Computers, 2016, 65, (12), pp. 3509–3523

31 Zhang, L., Hao, D., Zhang, L., Rothermel, G., Mei, H.: 'Bridging the gap between the total and additional test-case prioritization strategies'. In: Proceedings of the 35th International Conference on Software Engineering (ICSE'13). (, 2013, pp. 149–150

- 7 Engineering, 2011, 37, (1), pp. 48–64
- 8 34 Qu, X., Cohen, M.B., Woolf, K.M. 'Combinatorial interaction regression testing:
9 A study of test case generation and prioritization'. In: Proceedings of the 23rd
10 International Conference on Software Maintenance (ICSM'07). (, 2007. pp. 255–
11 264
- 12 35 Qu, X., Cohen, M.B., Woolf, K.M. 'A study in prioritization for higher strength
13 combinatorial testing'. In: Proceedings of the 2nd International Workshop on
14 Combinatorial Testing, (IWCT'13). (, 2013. pp. 285–294
- 15 36 Huang, R., Zong, W., Chen, J., Towey, D., Zhou, Y., Chen, D. 'Prioritizing inter-
16 action test suite using repeated base choice coverage'. In: Proceedings of the IEEE
17 40th Annual Computer Software and Applications Conference (COMPSAC'16). (,
18 2016. pp. 174–184
- 19 37 Huang, R., Chen, J., Zhang, T., Wang, R., Lu, Y. 'Prioritizing variable-strength
20 covering array'. In: Proceedings of the IEEE 37th Annual Computer Software and
21 Applications Conference (COMPSAC'13). (, 2013. pp. 502–601
- 22 38 Huang, R., Xie, X., Towey, D., Chen, T.Y., Lu, Y., Chen, J.: 'Prioritization of com-
23 binatorial test cases by incremental interaction coverage', International Journal of
24 Software Engineering and Knowledge Engineering, 2013, 23, (10), pp. 1427–1457
- 25 39 Jia, Y., Harman, M.: 'An analysis and survey of the development of mutation
26 testing', IEEE Transactions on Software Engineering, 2011, 37, (5), pp. 649–678
- 27 40 Do, H., Elbaum, S.G., Rothermel, G.: 'Supporting controlled experimentation with
28 testing techniques: An infrastructure and its potential impact', Empirical Software
29 Engineering, 2005, 10, (4), pp. 405–435
- 30 41 Andrews, J.H., Briand, L.C., Labiche, Y., Namin, A.S.: 'Using mutation analy-
31 sis for assessing and comparing testing coverage criteria', IEEE Transactions on
32 Software Engineering, 2006, 32, (8), pp. 608–624
- 33 42 Do, H., Rothermel, G.: 'On the use of mutation faults in empirical assessments of
34 test case prioritization techniques', IEEE Transactions on Software Engineering,
35 2006, 32, (9), pp. 733–752
- 36 43 Jia, Y., Harman, M.: 'Higher order mutation testing', Information and Software
37 Technology, 2009, 51, (10), pp. 1379 – 1393
- 38 44 Ammann, P., Delamaro, M.E., Offutt, J. 'Establishing theoretical minimal sets of
39 mutants'. In: Proceedings of the 7th International Conference on Software Testing,
40 Verification and Validation (ICST'14). (, 2014. pp. 21–30
- 41 45 Kintis, M., Papadakis, M., Malevis, N. 'Evaluating mutation testing alterna-
42 tives: A collateral experiment'. In: Proceedings of the 17th Asia-Pacific Software
43 Engineering Conference (APSEC'10). (, 2010. pp. 300–309
- 44 46 Papadakis, M., Henard, C., Harman, M., Jia, Y., Le.Traon, Y. 'Threats to the valid-
45 ity of mutation-based test assessment'. In: Proceedings of the 25th International
46 Symposium on Software Testing and Analysis (ISSTA'16). (, 2016. pp. 354–365
- 47 47 Huang, R., Chen, J., Chen, D., Wang, R. 'How to do tie-breaking in prioritization
48 of interaction test suites?'. In: Proceedings of the 26th International Conference on
49 Software Engineering and Knowledge Engineering (SEKE'14). (, 2014. pp. 121–
50 125
- 51 48 Wang, Z., Chen, L., Xu, B., Huang, Y.: 'Cost-cognizant combinatorial test case
52 prioritization', International Journal of Software Engineering and Knowledge
53 Engineering, 2011, 21, (6), pp. 829–854
- 54 49 Arcuri, A., Briand, L.: 'A hitchhiker's guide to statistical tests for assessing ran-
55 domized algorithms in software engineering', Software Testing, Verification and
56 Reliability, 2014, 24, (3), pp. 219–250
- 57 50 Harman, M., McMinn, P., Souza, J., Yoo, S.: 'Search based software engineering:
58 Techniques, taxonomy, tutorial', Empirical Software Engineering and Verification,
59 2012, pp. 1–59
- 60 51 Vargha, A., Delaney, H.D.: 'A critique and improve of the cl common language
61 effect size statistics of mcgraw and wong', Journal of Education and Behavioral
62 Statistics, 2000, 25, (2), pp. 101–132
- 63 52 Hemmati, H., Arcuri, A., Briand, L.: 'Achieving scalable model-based testing
64 through test case diversity', ACM Transactions on Software Engineering and
65 Methodology, 2013, 22, (1), pp. 139–176
- 66 53 Huang, R., Zong, W., Towey, D., Zhou, Y., Chen, J. 'An empirical examination
67 of abstract test case prioritization techniques'. In: Proceedings of the 39th Inter-
68 national Conference on Software Engineering Companion (ICSE-C'17). (, 2017.
69 pp. 141–143