

# 基于程序不变量计算软件可靠性\*

周 远<sup>1</sup>, 丁佐华<sup>2</sup>

<sup>1</sup>(浙江理工大学 理学院, 浙江 杭州 310018)

<sup>2</sup>(浙江理工大学 信息学院, 浙江 杭州 310018)

通讯作者: 丁佐华, E-mail: zouhuading@hotmail.com



**摘 要:** 现有的计算软件可靠性的方法采用测试的输入/输出结果,但这些数据并不能真实地反映软件内部的真实行为,如测试中会出现假性正确的情况以及测试不能显示一个输入有多个错误的输出情况.试图通过程序不变量来计算软件的可靠性,程序不变量可以描述程序的性质.首先选取测试用例集,动态地获取程序不变量,再从这些不变量中提取失效数据,最后,基于 Nelson 模型计算软件的可靠性.作为实验,对西门子程序包计算软件的可靠性.采用随机、分支覆盖和分块覆盖这 3 种不同的测试方法得到程序不变量,据此计算程序的可靠性.为了检查结果的可行性,采用传统方法计算这些软件的可靠性.两种可靠性比较后显示:它们的差别很小,而且不依赖于对测试方法的选择.通过进一步的方差分析得知,用所提出的方法计算的可靠性比用现有的方法计算的可靠性具有更小的波动,即更平稳.因此,前者更接近系统的真实可靠性.结论说明,可用程序不变量来计算软件的可靠性.

**关键词:** 测试用例集;程序不变量;失效数据;Nelson 模型;软件可靠性

**中图法分类号:** TP311

中文引用格式: 周远,丁佐华.基于程序不变量计算软件可靠性.软件学报,2015,26(12):3075-3087. <http://www.jos.org.cn/1000-9825/4803.htm>

英文引用格式: Zhou Y, Ding ZH. Computing software reliability based on program invariants. Ruan Jian Xue Bao/Journal of Software, 2015, 26(12): 3075-3087 (in Chinese). <http://www.jos.org.cn/1000-9825/4803.htm>

## Computing Software Reliability Based on Program Invariants

ZHOU Yuan<sup>1</sup>, DING Zuo-Hua<sup>2</sup>

<sup>1</sup>(College of Science, Zhejiang Sci-Tech University, Hangzhou 310018, China)

<sup>2</sup>(College of Information, Zhejiang Sci-Tech University, Hangzhou 310018, China)

**Abstract:** Existing methods to compute software reliability use the testing data of input and output. However, the data cannot truly reflect the behaviors of the software. For examples, it is possible that testing has false positive outputs, and one input with multiple faults cannot be revealed in the testing. This paper attempts to use program invariants to compute software reliability. First, test cases are selected to dynamically obtain program invariants. Then, the failure data from these invariants are extracted. Finally, the software reliabilities is computed based on Nelson model. In experiments, the reliabilities of the software in the Siemens set are computed. Program invariants are obtained by applying three different testing methods: random, branch coverage and block coverage. The reliabilities are computed next based on these invariants. To check the correctness of the results, software reliabilities are also computed with traditional methods. A comparison shows that the differences between two types of reliabilities are small no matter which testing method is selected. Further variance analysis shows that reliabilities by the proposed method have lower fluctuation, i.e., more stable, than those with existing methods, thus is closer to the real reliability of the system. This suggests that software reliability can be computed based on program invariants.

\* 基金项目: 国家自然科学基金(61210004, 61170015)

Foundation item: National Natural Science Foundation of China (61210004, 61170015)

收稿时间: 2014-07-31; 修改时间: 2014-10-14; 定稿时间: 2014-12-08

**Key words:** testing case; program invariant; failure data; Nelson model; software reliability

软件可靠性是衡量系统质量的重要指标,软件可靠性的计算依赖于可靠性模型的选取以及所获得的软件失效数据.软件可靠性模型分为软件可靠性增长模型<sup>[1-9]</sup>和基于架构的模型<sup>[11-19]</sup>.这些模型都附带有一些假设条件,如测试时的程序执行和给定的操作剖面要一致,或构件的失效是相互独立的.所以在实际当中,用得最多的是 Nelson 模型.

失效数据的获取对可靠性及模型有非常大的影响<sup>[20]</sup>.模型中参数的确定或者模型的计算必须依赖于失效数据,失效数据将影响到可靠性计算的准确性.但是因为我们使用的模型有很多的假设,而在实际中这些假设并不能完全满足,从而导致我们获得的失效数据并不能很好地符合实际要求.因此,为了获得准确的软件可靠性,我们必须考虑失效数据的质量和收集失效数据的方案<sup>[20]</sup>.目前,有一些针对各个开发阶段的失效数据收集技术<sup>[21]</sup>,但应用最广泛的方法是软件测试技术.软件测试包括 3 个要素:选择输入域、执行测试用例和比较输出<sup>[22]</sup>.学者们提出了多种测试方法来获得软件失效数据<sup>[23]</sup>,然而这些方法都是通过比较实际输出结果与期望输出结果得到诸如失效次数、失效用例等的失效数据,而忽略了软件的实时行为和内部结构,从而导致这些数据并不能真实反映软件的真实行为,如:

问题 1. 测试中会出现假性正确的情况.

当某一个测试用例运行到了错误行,但其结果为 pass 的时候,这一次测试被称为偶然性正确.如图 1 所示,t1 覆盖了第 1 行~第 5 行、第 7 行、第 8 行、第 14 行,其中,第 8 行有一个 BUG,但结果是 pass.这种情况下,很明显,t1 的 pass 是偶然性正确.当存在较多的假性正确,我们将无法获取正确的失效数据,从而导致计算出的可靠性无法评估软件的真实可靠性.更极端一点的情况是:当所有的错误用例都产生了假性正确,那么我们将无法收集到失效数据,从而得出软件的失效率为 0.

Program	Test case					
	t1(3,3,5)	t2(1,2,3)	t3(3,2,1)	t4(5,5,5)	t5(5,3,4)	t6(2,1,3)
mid(){						
1: int x, y, z, m;	●	●	●	●	●	●
2: read ("Enter 3 number:",x,y,z)	●	●	●	●	●	●
3: m=z;	●	●	●	●	●	●
4: if (y<z)	●	●	●	●	●	●
5: if (x<y)	●	●			●	●
6: m=y;		●				
7: else if (x<z)	●				●	●
8: m=y;	●					●
/****BUG****/						
9: else			●	●		
10: if (x>y)			●	●		
11: m=y			●			
12: else if (x>z)				●		
13: m=x;						
14: print ("Middle number is:",m)	●	●	●	●	●	●
}						
Status: P(pass)/F(Fail)	P	P	P	P	P	F

Fig.1 An example with false positive outputs

图 1 具有假性正确的例子

问题 2. 测试不能显示一个输入有多个错误的输出情况.

如图 2 所示,考虑测试用例 t1 (input: 1, 1, 2; output: 0),程序的输出为 0.5,而正确的输出应该是 0.这个用例经过了两个错误(分别为第 3.2 行和第 5 行),但我们只得到一个错误的输出.这样就无法了解到程序内部的信息,从而导致计算出的可靠性与实际程序的可靠性有一定的差距.

问题 3. 不同测试方法会影响到失效数据的获取比例,从而影响可靠性的计算.

Chen 等人<sup>[24]</sup>通过实验指出,不同的测试方法(随机测试、功能测试、分块测试)获得的失效数据会产生不同的可靠性.作者认为:这是因为这样获得的失效数据并没有真正反映程序的内部信息,仅反映了在所选测试用例下粗糙的程序的信息,所以计算出来的可靠性也是不准确的,无法反映软件真实的可靠性.

Program		Test case
		t1 (input: 1, 1, 2; output: 0)
<code>int main(){</code>		●
<code>1: float a, b, c, d, e;</code>		●
<code>2: cin&gt;&gt;a&gt;&gt;b&gt;&gt;c;</code>		●
<code>3: d=add(a,b);</code>	<code>3.1: float add(float a, float b){</code>	●
	<code>3.2: return a*b /**BUG**/ /*correct code: return a+b*/</code>	●
	<code>3.3 }</code>	●
<code>4: if (c&gt;0)</code>		●
<code>5: e=d/c; /**BUG**/ /*correct code is: e=d-c;</code>		●
<code>6: else</code>		
<code>7: e=d+c;</code>		
<code>8: cout&lt;&lt;"a+b- c ="&lt;&lt;e&lt;&lt;endl;</code>		●
<code>9: return 0;</code>		●
<code>}</code>		●
<code>output:</code>		0.5

Fig.2 An example with one input passing multiple faults

图 2 一个输入经过多个错误的例子

因此,仅通过测试阶段的输入/输出所获取的失效数据并不充分,无法准确反映程序内部存在的问题,从而得到的可靠性不够准确.为了解决这个问题,我们需要程序的内部数据来计算软件的可靠性.一个解决方案是:通过观察和分析测试过程中程序的执行轨迹从中提取程序的内部信息(例如内部结构、动态行为).我们采用工具 Daikon 做动态分析.Daikon 从程序的执行过程中发现可能的程序不变量.不变量是指程序在运行时保持不变的属性的量,刻画了程序的动态性质,体现了软件的整体行为.从这些不变量中得到不好的不变量,利用这些不好的不变量计算可靠性.尽管有个别工作,如 Pietrantuono,通过程序不变量来监测和预测在线软件的可靠性<sup>[25]</sup>,丁佐华等人通过程序不变量来计算服务组合的在线可靠性计算<sup>[26]</sup>,但都没有对其可行性论述.本文通过实验来说明其可行性.我们的实验程序来自于西门子程序包.采用随机、分支覆盖和分块覆盖这 3 种不同的测试方法得到程序不变量,据此分别计算程序的可靠性.

为了检查所得结果的正确性,我们采用现有的可靠性方法再次计算这些软件的可靠性.两种可靠性比较后显示:它们的差别很小,而且不依赖于我们对测试方法的选择.进一步的方差分析告诉我们:用所提方法计算的可靠性比用现有的方法计算的可靠性具有更小的波动,也就是说更平稳.因此,前者更接近系统的真实可靠性.

本文第 1 节描述如何得到程序不变量和不好的程序不变量,即失效数据.第 2 节描述如何基于不变量进行可靠性计算.第 3 节介绍实验的准备工作,包括 3 种选取测试用例集的方法以及西门子程序包.第 4 节描述如何从西门子程序包获取失效数据.第 5 节是软件可靠性的计算及结果分析.第 6 节是所提方法的可行性分析.第 7 节是相关工作介绍.第 8 节是全文的总结.

## 1 程序不变量和失效数据

程序不变量常被用于追踪软件缺陷<sup>[27]</sup>、硬件错误<sup>[28]</sup>和软件评估<sup>[29]</sup>.程序不变量揭示了程序在运行过程中保持不变的状态,这些状态对软件的发展、测试和维护有重要的意义.常用的提取不变量的工具有 DIDUCE<sup>[27]</sup>和 Daikon<sup>[29]</sup>.本文采用 Daikon 作为提取不变量的工具.

Daikon 工具作为动态不变量发现器,利用程序执行来发现可能的不变量.它通过修改目标程序来跟踪感兴趣的域,通过一组测试套件来运行修改以后的程序,从而利用跟踪得来的结果推导不变量.这些推导出来的不变

量也可以作为形式化的规格说明插入到代码中作为注释,帮助提高程序可理解性.

通常,默认不变量有 196 个类型,这些类型基本上满足了我们的要求.当然,我们可以在此基础上定义更多的不变量类型.如果在现有的类型下所得到的不变量是不完全的,这说明涉及这些不变量的语句没有被执行到.一种情况是这些语句是多余的,另一种情况是测试用例的选取不够充分.如果我们的覆盖率能够达到希望的要求,那么即便是有少量的不变量没有计算到,也不会对结果产生太大影响.

下面我们通过一个计算器例子来说明如何得到程序不变量(见表 1):

**Table 1** An example of calculator

表 1 计算器例子

Code
<code>JButton b[]=new JButton[16];</code>
<code>String s[]={“0”,“1”,“2”,“3”,“4”,“5”,“6”,“7”,“8”,</code>
<code>“9”,“+”,“-”,“*”,“/”,“C”,“=”};</code>
<code>...</code>
<code>public void actionPerformed(ActionEvent e)</code>
<code>{ ...</code>
<code>    if (flag==10) m2=m1+m2;</code>
<code>    else if (flag==11) m2=m2-m1;</code>
<code>    else if (flag==12) m2=m2*m1;</code>
<code>    else if (flag==13) m2=m2/m1;</code>
<code>    else m2=m1;</code>
<code>    ...</code>
<code>}</code>
<code>...</code>

- (1) 生成测试用例集.在该例子中,我们用原先已有的测试用例生成技术<sup>[30]</sup>生成一个测试用例集,该技术结合白盒测试和黑盒测试方法将输入域分成若干个等价类,再从等价类中选取测试用例.为了更好地提取不变量,我们还采用 Gupta<sup>[31]</sup>提出的方法进一步提炼测试用例集.在给定的例子中,测试用例集包括 80 个测试用例;
- (2) 定义不变量类型.Daikon 不变量类型有很多种,如  $x=a, x \neq 0, a \leq x \leq b, y=ax+b, x \leq y, x=fn(y), x \in y$  等,我们还可以自定义程序不变量的类型;
- (3) 运行程序和测试用例,得到轨迹文件;
- (4) 得到程序不变量.在该例子中得到的部分不变量见表 2,其中,不变量  $size(this.b[])=size(this.s[])$  是指数组  $b$  的长度等于数组  $s$  的长度.

**Table 2** A part of the program invariants of the program

表 2 通过原程序得到的部分不变量

Code (C1)	Invariants (I1)
<code>JButton b[]=new JButton[16];</code>	1) <i>JB::OBJECT</i>
<code>String s[]={“0”,“1”,“2”,“3”,“4”,“5”,“6”,“7”,“8”,</code>	2) $size(this.b[])=size(this.s[])$
<code>“9”,“+”,“-”,“*”,“/”,“C”,“=”};</code>	3) $this.m1 \geq 0$
<code>...</code>	4) $this.flag \geq 0$
<code>public void actionPerformed(ActionEvent e)</code>	5) $this.b! = null$
<code>{ ...</code>	6) $this.b[]$ contains no nulls and has only one value, of length 16
<code>    if (flag==10) m2=m1+m2;</code>	7) $this.b[]$ elements != null
<code>    else if (flag==11) m2=m2-m1;</code>	8) $this.s[]$ contains no nulls and has only one value, of length 16
<code>    else if (flag==12) m2=m2*m1;</code>	9) $size(this.b[])=16$
<code>    else if (flag==13) m2=m2/m1;</code>	10) $this.s.getClass()=java.lang.String[].class$
<code>    else m2=m1;</code>	11) $this.m1 \leq this.m2$
<code>    ...</code>	12) $this.flag < size(this.b[]) - 1$
<code>}</code>	
<code>...</code>	

接下来,我们描述如何得到不好的程序不变量,即失效数据.程序 C1 改为程序 C2,将 `JButton b[]=new JButton[16]` 中的 16 改为 17,以及调换 `String s[]={“0”,“1”,“2”,“3”,“4”,“5”,“6”,“7”,“8”,“9”,“+”,“-”,“*”,“/”,“C”,“=”}`,

“=”}中‘+’和‘-’的位置时,我们会得到新的不变量集  $T_2$ ,见表 3.通过对比发现:不变量 1)、不变量 3)~不变量 5)、不变量 7)、不变量 10)~不变量 12)跟原先一样,而不变量 2)、不变量 6)、不变量 8)、不变量 9)发生了变化.如原先的不变量 2)是  $size(this.b[])==size(this.s[])$ ,而改变后的是  $size(this.b[])-1==size(this.s[])$ .显然,这里描述数组  $b$  的长度等于数组  $s$  的长度的不变量发生了变化,因此认为新的不变量  $size(this.b[])-1==size(this.s[])$  是一个失效数据.

**Table 3** A part of program invariants of the modified program

**表 3** 修改程序后得到的部分不变量

Code (C2)	Invariants (T2)
<pre> JButton b[]=new JButton[17]; String s[]={“0”,“1”,“2”,“3”,“4”,“5”,“6”,“7”,“8”, “9”,“-”,“+”,“*”,“/”,“C”,“=”}; ... public void actionPerformed(ActionEvent e) { ...     if (flag==10) m2=m1+m2;     else if (flag==11) m2=m2-m1;     else if (flag==12) m2=m2*m1;     else if (flag==13) m2=m2/m1;     else m2=m1;     ... } ... </pre>	<ol style="list-style-type: none"> <li>1) <math>JB::OBJECT</math></li> <li>2) <math>size(this.b[])-1==size(this.s[])</math></li> <li>3) <math>this.m1 \geq 0</math></li> <li>4) <math>this.flag \geq 0</math></li> <li>5) <math>this.b!=null</math></li> <li>6) <math>this.b[]</math> contains no nulls and has only one value, of length 17</li> <li>7) <math>this.b[]</math> elements!=null</li> <li>8) <math>this.s[]</math> contains no nulls and has only one value, of length 17</li> <li>9) <math>size(this.b[])==17</math></li> <li>10) <math>this.s.getClass()==java.lang.String[].class</math></li> <li>11) <math>this.m1 \leq this.m2</math></li> <li>12) <math>this.flag &lt; size(this.b[])-1</math></li> </ol>

在该程序中,执行上述所有测试用例,共生成不变量 3 207 个,其中,失效不变量有 282 个.

## 2 基于不变量的可靠性计算

### 2.1 可靠消息计算模型

本文采取 Nelson 可靠性模型来计算软件的可靠性.Nelson 软件可靠性模型对软件做出如下假设:

- 1) 程序被认为是集合  $E$  上的一个可计算函数  $F$  的一个规范,这里用  $E=(E_i; i=1, \dots, m)$  表示用于执行程序的所有输入数据的集合,一个输入数据对应一个程序执行回合;
- 2) 对于每个输入  $E_i$ ,程序执行产生输出  $F(E_i)$ ;
- 3) 由于程序包含缺陷,程序实际确定函数  $F'$ ,该函数不同于希望函数  $F$ ;
- 4) 对于某些  $E_i$ ,程序实际输出  $F'(E_i)$  在希望输出  $F(E_i)$  的容许范围之内,即  $|F'(E_i)-F(E_i)| \leq \Delta_i$ .但对另一些  $E_j$ ,程序实际输出  $F'(E_j)$  超出容许范围,即  $|F'(E_j)-F(E_j)| > \Delta_j$ ,这时认为程序发生一次失效;
- 5) 测试过程中不剔除程序缺陷.

Nelson 模型计算软件可靠性的函数为

$$R = 1 - \sum_{i=1}^m \frac{f_i}{n_i} P(E_i) \tag{1}$$

其中,  $E_i$  域中数据被选取的概率是  $P(E_i)$ ;  $E_i$  域中被选取的数据的个数为  $n_i$ ;  $E_i$  域中  $n_i$  个被选取数据导致程序失效的个数为  $f_i$ .

根据上述模型,接下来具体描述如何计算表 1 给出的计算器例子的可靠性.

Step 1. 生成测试用例集,在第 1 节中,我们共得到了 80 个测试用例.

Step 2. 将测试用例进行分类.具体的分类方法是:

- a) 通过 gcov 编译,执行测试用例,生成测试用例对程序的覆盖信息;
- b) 分析该测试用例对程序中每个分支或者子函数的覆盖率;
- c) 测试用例  $A$  和测试用例  $B$ ,如果覆盖程序相同的部分超过一定的值(该值需要根据程序的结构以及测试用例集来确认),那么认为测试用例  $A$  和测试用例  $B$  是同一类测试用例.

接下来,我们选取计算器中如下的 4 个选择分支来阐述如何进行分类:

- 分支 1: if ( $flag==10$ )  $m2=m1+m2$ ;
- 分支 2: else if ( $flag==11$ )  $m2=m2-m1$ ;
- 分支 3: else if ( $flag==12$ )  $m2=m2*m1$ ;
- 分支 4: else if ( $flag==13$ )  $m2=m2/m1$ .

测试用例:

- $T1$ :  $x1\ 2\ 3$ ;
- $T2$ :  $x1\ 9.1\ 8.5$ ;
- $T3$ :  $x2\ 89\ 63$ ;
- $T4$ :  $x2\ 77\ 34.9$ ;
- $T5$ :  $x4\ 9\ 3$ ;
- $T6$ :  $x3\ 76.1\ 63$ .

测试用例中, $x1,x2,x3,x4$  表示标志位:当标志位是  $x1$  时,程序执行分支 1;当标志位是  $x2$  时,程序执行分支 2;当标志位是  $x3$  时,程序执行分支 3;当标志位是  $x4$  时,程序执行分支 4.当执行测试用例  $T1$  时,通过 `gcov` 可以知道:分支 1 被调用,测试用例覆盖分支 1,而没有覆盖其他 3 个分支.同理,我们可以得到  $T2$  仅覆盖了分支 1, $T3$  和  $T4$  覆盖了分支 2, $T5$  覆盖了函数分支 4, $T6$  覆盖了分支 3.因此,我们把  $T1$  和  $T2$  分为一组测试用例集, $T3$  和  $T4$  分为一组测试用例集, $T5$  分为一组测试用例集, $T6$  分为一组测试用例集.

根据上述方法,我们可以将所得的 80 个用例分成 5 类,具体分类结果见表 4:

**Table 4** Classification of the test cases

表 4 测试用例的分类

Sort	1	2	3	4	5	Total
Number	17	11	23	26	3	80

Step 3. 分类后,运行每个测试用例集和程序,根据第 1 节所述,得到程序不变量和失效不变量.具体数据见表 5.

**Table 5** Program invariants and failure data of the test cases

表 5 各组用例产生的不变量和失效数据(失效不变量)

Sort	1	2	3	4	5
Total invariants	1 286	999	1 542	2 201	47
Failure invariants	8	192	198	269	2

Step 4. 通过得到的失效数据进行可靠性计算.

根据 Nelson 模型,我们计算得到的可靠性为

$$R = 1 - \left( \frac{17}{80} \frac{8}{1286} + \frac{11}{80} \frac{192}{999} + \frac{23}{80} \frac{198}{1542} + \frac{26}{80} \frac{269}{2201} + \frac{3}{80} \frac{2}{47} \right) = 0.894.$$

### 3 实验设置

为了与传统的基于不同测试方法的可靠性计算方法进行比较,以此来说明基于不变量计算可靠性的可行性和优点,我们采用了 3 个常用的测试用例生成方法:随机生成、基于分块覆盖生成和基于分支覆盖生成以及西门子程序包:the Software-artifact Infrastructure Repository<sup>[32]</sup>.我们在实验中验证的内容有:(1) 失效数据(失效不变量)获取的可行性,即,不同测试方法对失效数据的比例影响不大;(2) 可靠性计算的可行性,即,基于不变量计算的可靠性与传统方法计算所得的可靠性接近.验证的优点是:基于不变量计算所得的可靠性更稳定,因此更能体现程序的可靠性.

### 3.1 3种测试方法

本文采用随机选取、基于分块覆盖选取和基于分支覆盖选取这3种不同的测试方法从测试用例集中选取测试用例。

随机测试即从程序的输入域中随机地选取一个测试用例,每一个测试用例都有一个被选中的概率.在我们的实验中,根据预先定义的发生概率,从程序的规格说明书中随机选取至少一个功能进行随机测试;然后,随机选取一个合适的测试用例来对这个说明书进行测试,直至选取的测试用例集覆盖程序的所有功能.分块覆盖测试是我们把只有一个入口和一个出口的顺序语句作为一个块.选取测试用例时,要覆盖所有的程序块及至少有一个测试用例覆盖一个程序块.按分支覆盖准则进行测试是指:设计若干测试用例,运行被测程序,使得程序中每个判断的取真分支和取假分支至少经历一次及判断的真假值均曾被满足.

### 3.2 程序来源

本文采用的实验程序来源于 the Software-artifact Infrastructure Repository<sup>[32]</sup>.The Siemens set 包含7个C程序,见表6.它们常用于错误定位和软件可靠性计算.

Table 6 Siemens set

表6 Siemens set

Program	Faulty versions	Test cases	Description
print_tokens	7	4 130	Lexical analyzer
print_tokens2	10	4 115	Lexical analyzer
replace	32	5 542	Pattern recognition
schedule	9	2 650	Priority scheduler
schedule2	10	2 710	Priority scheduler
tcas	41	1 608	Altitude separation
tot_info	23	1 052	Information measure

表6详细描述了我们的实验所需要的7个程序,每个程序包含一个正确的版本以及多个错误版本,每个错误版本包含一个错误(在我们的实验中,需要把所有的错误整合在一起,构成一个包含多个错误的版本).另外,每个程序有一个输入集(即测试用例集).

## 4 产生失效数据的可行性分析

我们通过3种不同的方法来选取测试用例集  $T_1$ (随机选取)、 $T_2$ (基于分块覆盖选取)、 $T_3$ (基于分支覆盖选取).分别运行正确版本程序  $C$  与  $T_1$ 、 $C$  与  $T_2$ 、 $C$  与  $T_3$  得到程序不变量  $TI_1$ 、 $TI_2$  和  $TI_3$ ,运行错误版本程序  $F$  与  $T_1$ 、 $F$  与  $T_2$ 、 $F$  与  $T_3$  得到程序不变量  $FI_1$ 、 $FI_2$  和  $FI_3$ .比较  $TI_1$  与  $FI_1$ 、 $TI_2$  与  $FI_2$  以及  $TI_3$  与  $FI_3$  的不同,分别得到各自的失效数据.

表7是3种方法得到的程序不变量总数和失效数据的个数.Random表示随机选取测试用例得到的程序不变量,Block表示基于分块覆盖选取测试用例得到的程序不变量,Branch表示基于分支覆盖选取测试用例得到的程序不变量.False指失效的程序不变量数,Total表示程序不变量的总数.由表可以发现:对于同一个程序,采用不同的测试方法获得的不变量是不一样的;当然,获得的失效数据,即不好的不变量也是不一样的.

Table 7 Failure data and the total program invariants gotten by the three test methods

表7 3种方法得到的失效数据与程序不变量

Program	Random		Block		Branch	
	False	Total	False	Total	False	Total
print_tokens	1 038	16 925	1 015	16 851	998	16 832
print_tokens2	1 286	15 971	1 269	15 735	1 315	15 821
replace	2 410	18 065	2 359	17 906	2 187	17 318
schedule	1 472	15 126	1 354	14 973	1 308	14 865
schedule2	1 873	16 725	1 749	16 287	1 702	15 994
tcas	1 957	20 108	1 913	19 947	1 864	20 116
tot_info	1 965	20 294	1 870	19 083	1 897	19 525

接下来计算失效数据所占的百分比.表 8 给出了每种方法所获得的失效数据百分比,可以发现不同方法获得的失效数据百分比基本相近,误差不大.例如:对于程序 `print_tokens`,随机方法、基于分块覆盖方法和基于分支覆盖方法所得到的失效数据百分比分别为 6.13%,6.02%,5.92%,其最大的相对误差只有  $(6.13\%-5.92\%)/6.13\% \times 100\% = 3.4\%$ ;其次,我们计算得到 3 种方法平均能够发现的失效数据百分比分别为 9.69%(random),9.49%(Block)和 9.32%(Branch),最大相对误差为 3.8%.

**Table 8** Percentages of the failure data

**表 8** 失效数据所占程序不变量总数的百分比

Program	Random (%)	Block (%)	Branch (%)
<code>print_tokens</code>	6.13	6.02	5.92
<code>print_tokens2</code>	8.05	8.06	8.31
<code>replace</code>	13.34	13.17	12.62
<code>schedule</code>	9.73	9.04	8.79
<code>schedule2</code>	11.19	10.73	10.64
<code>tcas</code>	9.73	9.59	9.26
<code>tot_info</code>	9.68	9.79	9.71

综上所述,我们可得出结论:不同测试方法得到不同的不变量;但是就失效数据的百分比来说,它们之间相差较小,比较稳定.

## 5 可靠性计算和结果分析

根据第 2 节的方法,我们来计算 7 个程序的可靠性.我们以程序 `schedule` 为例,详细叙述失效数据的获得和可靠性的计算.通过上述 3 种方法,我们将随机选取的测试用例分为 15 个测试用例集,将基于分块覆盖选取的测试用例分为 14 个测试用例集,将基于分支覆盖选取的测试用例分为 14 个测试用例集.

Step 1. 用随机测试方法选取了 260 个测试用例,分块覆盖方法选取了 272 个测试用例,分支覆盖方法选取了 295 个测试用例.

Step 2. 将测试用例进行分类,分类结果见表 9.

**Table 9** Result of the classification of the test cases from program `schedule`

**表 9** 程序 `schedule` 的测试用例集分类结果

Sort	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	总计
Random	1	25	14	2	7	82	4	6	7	40	4	34	2	3	29	260
Block	18	28	9	11	73	2	4	9	30	11	52	2	4	19	-	272
Branch	2	20	9	10	91	4	5	6	43	2	62	6	15	20	-	295

Step 3. 分类后,运行每个测试用例集和程序,根据前一节所述办法得到总的不变量和失效数据,见表 10.

**Table 10** Failure data and total program invariants in each sort

**表 10** 分类后所得失效数据和总不变量数目

Sort	Random		Block		Branch	
	失效数据	总不变量	失效数据	总不变量	失效数据	总不变量
1	4	27	1 054	4 742	54	606
2	554	4 208	537	3 727	1 042	5 919
3	721	5 100	0	5 100	29	372
4	8	999	2 764	12 944	1 683	8 938
5	715	10 984	2 815	13 029	2 314	12 629
6	3 983	12 546	1 162	14 764	853	7 410
7	394	9 930	1 319	3 512	983	2 864
8	912	5 054	0	5 054	861	1 511
9	77	3 710	1 835	11 987	3 349	11 409
10	1 934	10 446	791	6 421	595	1 046
11	192	3 216	4 581	13 903	778	11 086
12	1 343	11 824	1 025	6 735	475	8 964
13	257	4 369	415	4 482	1 262	6 318
14	594	4 616	1 846	12 797	1 224	4 164
15	1 236	4 486				

Step 4. 将得到的数据代入 Nelson 模型进行可靠性计算.计算结果为:0.860 8(随机),0.862 5(基于分块覆盖),0.877 2(基于分支覆盖).

其余程序的可靠性计算结果见表 11,其中,R-reliability 表示随机选取测试用例计算的可靠性,Block-reliability 表示基于分块覆盖选取测试用例计算的可靠性,Branch-reliability 表示基于分支覆盖选取测试用例计算的可靠性.

**Table 11** Reliability computed by program invariants

表 11 通过不变量计算的可靠性

Program	R-reliability	Block-reliability	Branch-reliability
print_tokens	0.854 1	0.860 5	0.867 9
print_tokens2	0.841 7	0.852 6	0.839 7
replace	0.804 4	0.817 0	0.820 3
schedule	0.795 0	0.800 3	0.810 7
schedule2	0.821 4	0.829 8	0.835 0
tcas	0.789 2	0.794 7	0.805 1
tot_info	0.835 3	0.827 8	0.821 3

由表 11 可知:不同测试方法算出的可靠性虽然不一样,但相差较小.例如:对于程序 print\_tokens,随机方法算出的可靠性最低,为 0.854 1;基于分支覆盖方法算出的可靠性最大,为 0.867 9;而两者的差值为 0.013 8.

我们再来看一下不同方法产生的可靠性的波动大小,这里采用方差来刻画.根据概率统计的理论,方差越大,个体间的波动就越大.也就是说,不同测试方法对可靠性的影响就越大.表 12 计算了不同程序在 3 种测试方法下的可靠性方差.3 种方法对不同程序可靠性的影响动荡范围在 0.000 003 2~0.000 004 7 之间,因此,所得的可靠性是比较稳定的,测试方法对其影响不大.这是因为不变量反映的是程序整体的行为,所计算的可靠性是软件整体的可靠性,而不是在具体测试方法下的局部可靠性.

**Table 12** Variance of the computed reliability

表 12 不变量计算的可靠性方差

程序	方差
print_tokens	0.000 003 2
Print_tokens2	0.000 003 2
Replace	0.000 004 7
Schedule	0.000 004 3
Schedule2	0.000 003 2
Tcas	0.000 004 4
tot_info	0.000 003 3

综上所述,我们发现:

- 1) 在不同测试方法下,由不变量计算出的可靠性不是完全相等,还有一定的误差,其原因是有些不好的不变量之间还有一定的依赖性;
- 2) 用不变量来计算可靠性时,测试方法对软件的可靠性影响不大.

## 6 可行性分析

本节利用现有的计算可靠性方法来计算所给 7 个程序的可靠性,即用标准的 Nelson 模型计算.这里我们不再详细叙述计算过程,只给出计算的可靠性.然后,我们比较用不变量计算出的可靠性和用现有方法计算出的可靠性之间的差异.可靠性计算结果见表 13.

由表 13 得知,在同一测试方法下,用不变量和用现有方法计算出的可靠性有细微差异,有时不变量计算的可靠性大,有时现有方法计算的可靠性大;但是两者相差不大,这可由表 14 的相对误差的绝对值得到.在表 14 中,最大的相对误差绝对值仅有 3.787%,而大部分相对误差的绝对值在 1% 左右.

**Table 13** Reliability computed by program invariants and by traditional method**表 13** 利用不变量和利用现有方法计算出的可靠性

Program	R-reliability		Block-reliability		Branch-reliability	
	不变量	传统方法	不变量	传统方法	不变量	传统方法
print_tokens	0.854 1	0.847 2	0.860 5	0.861 3	0.867 9	0.874 8
print_tokens2	0.841 7	0.851 8	0.852 6	0.849 1	0.839 7	0.864 7
replace	0.804 4	0.797 2	0.817 0	0.810 3	0.820 3	0.829 7
schedule	0.795 0	0.808 6	0.800 3	0.831 8	0.810 7	0.819 4
schedule2	0.821 4	0.819 9	0.829 8	0.824 7	0.835 0	0.831 5
tcas	0.789 2	0.801 7	0.794 7	0.794 8	0.805 1	0.815 9
tot_info	0.835 3	0.824 9	0.827 8	0.817 5	0.821 3	0.836 1

**Table 14** Relative error between the two kinds of reliability**表 14** 两类可靠性的相对误差

Program	R-reliability (%)	Block-reliability (%)	Branch-reliability (%)
print_tokens	0.814 4	-0.092 9	-0.7888
Print_tokens2	-1.185 7	0.412 2	-2.891 2
Replace	0.903 2	0.826 9	-1.132 9
Schedule	-1.681 9	-3.787 0	-1.061 8
Schedule2	0.182 9	0.618 4	0.420 9
Tcas	-1.559 2	-0.012 6	-1.323 7
tot_info	1.260 8	1.259 9	-1.770 1

接下来,我们比较现有方法和不变量方法计算出的两类可靠性在不同测试方法下产生的波动大小.根据前面所述,我们可以用方差来刻画这一性质.表 15 列出了这两类可靠性计算方法在 3 种测试方法下的方差.由表 15 可以发现:两种方法计算出的可靠性方差都不是太大;但是相比较而言,现有方法的可靠性波动远大于不变量计算出的可靠性波动.这是因为现有方法计算出的可靠性一定程度上依赖于测试方法,反映软件的局部行为,计算出的可靠性是在该测试方法下的可靠性,所以在不同的测试方法下,可靠性会产生相对较大的波动.而不变量反映软件整体的行为,虽然不同的测试方法会产生不同的不变量,但这些不变量的合成效果却是一样的,都反映软件整体的行为,所以不变量计算出的可靠性是软件的整体可靠性,其波动产生的原因在于有些失效数据(不好的不变量)之间还有一定的相关性.

**Table 15** Comparison on the variances of the two kinds of reliability**表 15** 两类可靠性的方差比较

程序	方差	
	不变量( $10^{-5}$ )	传统方法( $10^{-5}$ )
print_tokens	0.32	13
Print_tokens2	0.32	4.6
Replace	0.47	1.8
Schedule	0.43	9.0
Schedule2	0.32	2.7
Tcas	0.44	7.7
tot_info	0.33	5.8

综上所述:相对于现有的计算可靠性方法,我们还可以利用不变量来计算软件的可靠性;且利用不变量计算的可靠性更接近软件的实际可靠性.

注 1:Simens 软件包通常被用来做软件测试、错误定位和不变量计算的开源程序包,所以我们用它来做基于不变量的软件可靠性的评价.但它只是针对序列化程序,并没有并发程序.今后的一个研究工作就是针对并发程序来验证我们的方法.并发结构可以给我们带来大规模的程序,我们将在此基础上研究我们方法的稳定性与适应性.

## 7 相关工作

### 7.1 软件可靠性模型

在过去的几十年里,有许多关于可靠性模型的研究.总的来说,这些模型大致可以分为软件可靠性增长模型 (software reliability growth models,简称 SRGMs)和基于架构的模型.

软件可靠性增长模型 SRGMs<sup>[1-9]</sup>是时间域模型,它们可以用来估计软件可靠性和残留故障数量.这些模型首先假设测试是按照给定的操作剖面执行的<sup>[10]</sup>,然后,通过在测试阶段的记录的故障历史来预测软件程序的行为.基于架构的可靠性模型<sup>[11-18]</sup>又可以分为两类,即,测试阶段的架构<sup>[11,12,16-18]</sup>和设计阶段的架构<sup>[13-15]</sup>.前者只有程序代码可用,通过执行大量的测试来获取轨迹的数据,从而根据这些数据来决定软件系统的架构.利用在覆盖测试中获得的轨迹信息,可以构建程序的架构,如离散时间马尔科夫链<sup>[12]</sup>、连续时间马尔科夫链<sup>[16]</sup>以及半马尔科夫过程<sup>[11]</sup>.构件的失效行为可以用泊松过程或基于覆盖测试和故障密度的时间独立失效率来建模,在这种情况下,覆盖测试和单元测试在决定所有信息过程中起重要作用.后者事先已经定义了构件以及构件之间的相关性.每个构件的可靠性依赖于操作剖面和其在系统中的位置,系统可靠性的计算等价于构件执行轨迹的可靠性,其中,构件执行轨迹可看作是一系列构件的执行,而构件执行轨迹的可靠性又等于该轨迹上所有构件可靠性的乘积.轨迹可靠性问题还可以转化为基于构件的问题<sup>[13]</sup>.在这种情况下,可靠性的计算需要假设构件的失效是相互独立的以及接口的可靠性为 1.在无记忆程序中,构件运行的独立性是可接受的;但是对其他程序来说,该假设无法满足<sup>[15]</sup>.特别地,当某个构件被循环执行时,该构件多次的执行是不独立的.

### 7.2 失效数据

失效数据的来源广、获取技术多<sup>[21]</sup>,人们可以从历史版本或者类似软件中获得失效数据,此时我们无需拥有具体的需求的软件;在设计阶段,可以从需求文档、UML 图和项目计划书中获得相关信息和数据;通过测试获取相应的失效数据;用户反馈,例如故障报告、修改需求等,也可以作为获得数据的来源.通过测试来获取数据是其中使用最广泛的一种方法.测试阶段根据模型的不同主要收集两类不同的失效数据:

- 一是给定测试用例,收集失效的测试用例数<sup>[19]</sup>;
- 二是给定测试环境,运行软件一段时间,收集在这段时间内的失效次数和失效间隔时间<sup>[33,34]</sup>.

文献[23]归纳了用来获得失效数据的各种不同的测试方法.但是这些失效数据仅仅是软件测试时间内的输入/输出信息,而并没有考虑程序内部的信息.而软件的可靠性必然是由其内部信息,包括动态行为和内部结构来决定的,因此,根据这些失效数据算出来的可靠性并不十分准确.

程序不变量是通过收集软件的动态行为获取能够反映其内部信息的数据,因此它能够反映软件内部行为,比仅仅考虑程序的输入/输出的数据携带更多的信息,似乎可以作为更好的失效数据来源.然而对这方面的研究却还不多,仅有一些尝试.Pietrantuono 通过程序不变量来监测和预测在线软件可靠性<sup>[25]</sup>;丁佐华等人通过程序不变量来计算服务组合的在线可靠性计算<sup>[26]</sup>,但是这些工作并没有说明用不变量来计算可靠性的可行性.

## 8 结论

本文通过实验来说明我们可以通过程序不变量来计算可靠性:其一,在同一测试方法下,不变量计算出的可靠性和现有方法计算出的可靠性相差不大;其二,在不同测试方法下,不变量计算出的可靠性波动较小,更接近软件的实际可靠性.从不变量的角度来刻画软件可靠性,使我们对软件的整体行为和软件可靠性有更深刻的认识,有利于软件可靠性的评估和预测.

接下来的工作是:

- 1) 针对不变量之间的相关性:如果不好的不变量之间存在相关性,会增加失效数据,影响计算结果.我们应尽量去掉这种相关性,使得剩下的不变量尽可能独立.一种办法是建立数据流的依赖图;
- 2) 针对假性测试用例:有些正确的测试用例会产生错误的不变量.在这种情况下,我们就丢失了失效数

据.如果假性测试用例较多,那样计算出来的结果就不准确.为了解决这种情况,我们需要对测试用例的结果建立频谱矩阵,基于频谱矩阵去除掉这些测试用例;

- 3) 运用更多的测试方法来计算软件可靠性,进一步验证测试方法对软件可靠性的影响;
- 4) 在实际应用中,我们并不知道正确的程序,这样我们就不能通过比较来获得不好的不变量,而必须通过一定的方法,根据已有的测试用例信息和已获得的不变量信息来确定不好的不变量,这将是我们将下来要解决的关键问题之一.

**致谢** 我们向对本文工作给予支持和建议的同行,尤其是北京大学高可信软件技术教育部重点实验室金芝教授表示感谢.

## References:

- [1] Yamada S. Software Reliability Modeling: Fundamentals and Applications. Tokyo: Springer-Verlag, 2014. 6–15.
- [2] Wood A. Software reliability growth models. Technical Report 96.1, No.130056, Tandem Computers, 1996.
- [3] Goel AL, Okumoto K. Time-Dependent error-detection rate model for software and other performance measures. IEEE Trans. on Reliability, 1979,28:206–211.
- [4] Misra PN. Software reliability analysis. IBM Systems Journal, 1983,22:262–270.
- [5] Ohba M, Chou XM. Does imperfect debugging affect software reliability growth? In: Druffel LE, ed. Proc. of the ICSE'89. Pittsburg, 1989. 237–244.
- [6] Farr W. Software Reliability Modeling Survey. In: Lyu MR, ed. Handbook of Software Reliability Engineering. New York: IEEE Computer Society Press, 1996. 71–117.
- [7] Chen M, Lyu MR, Wong WE. Effect of code coverage on software reliability measurement. IEEE Trans. on Reliability, 2001,50(2): 165–170.
- [8] Zhao J, Zhang RB, Gu GC. Study on software reliability growth model considering failure dependency. Chinese Journal of Computers, 2007,30(10):1713–1720 (in Chinese with English abstract).
- [9] Li HF, Wang SQ, Liu C, Zheng J, Li Z. Software reliability model considering both testing effort and testing coverage. Ruan Jian Xue Bao/Journal of Software, 2013,24(4):749–760 (in Chinese with English abstract). <http://www.jos.org.cn/1000-9825/4257.htm> [doi: 10.3724/SP.J.1001.2013.04257]
- [10] Horgan JR, Mathur AP. Software Testing and Reliability. In: Lyu MR, ed. Handbook of Software Reliability Engineering. New York: IEEE Computer Society Press, 1996. 531–566.
- [11] Littlewood B. Software reliability modeled for modular program structure. IEEE Trans. on Reliability, 1979,28(3):241–246.
- [12] Cheung RC. A user-oriented software reliability model. IEEE Trans. on Software Engineering, 1980,6(2):118–125.
- [13] Dolbec J, Shepard T. A component based software reliability model. In: Bockus D, Bennet K, eds. Proc. of the Conf. of the Center for Advanced Studies on Collaborative Research. ACM Press, 1995. 19–29.
- [14] Wang WL, Wu Y, Chen MH. An architecture-based software reliability model. In: Proc. of the 1999 Pacific Rim Int'l Symp. on Dependable Computing. IEEE, 1999. 143–150.
- [15] Hamlet D. Are we testing for true reliability? IEEE Software, 1992,9(4):21–27.
- [16] Laprie JC, Kanoun K. Software Reliability and System Reliability. In: Lyu MR, ed. Handbook of Software Reliability Engineering. New York: IEEE Computer Society Press, 1996. 27–69.
- [17] Mao XG, Deng YJ. A general model for component-based software reliability. Ruan Jian Xue Bao/Journal of Software, 2004,15(1): 27–32 (in Chinese with English abstract). <http://www.jos.org.cn/1000-9825/15/27.htm>
- [18] Hou CY, Cui G, Liu HW. Rate-Based component software reliability process simulation. Ruan Jian Xue Bao/Journal of Software, 2011,22(11):2749–2759 (in Chinese with English abstract). <http://www.jos.org.cn/1000-9825/3930.htm> [doi: 10.3724/SP.J.1001.2011.03930]
- [19] Nelson E. Estimating software reliability from test data. Microelectronics Reliability, 1978,17(1):67–73.

- [20] Nguyen EA, Rexach CF, Thorpe DP, Walther AE. The importance of data quality in software reliability modeling. In: Proc. of the 21st Int'l Symp. on Software Reliability Engineering (ISSRE). IEEE, 2010. 220–228.
- [21] Buhnova B, Chren S, Fabriková L. Failure data collection for reliability prediction models: A survey. In: Seinturier L, Bures T, McGregor JD, eds. Proc. of the 10th Int'l ACM SIGSOFT Conf. on Quality of Software Architectures. ACM Press, 2014. 83–92.
- [22] Rapps S, Weyuker EJ. Selecting software test data using data flow information. IEEE Trans. on Software Engineering, 1985,(4): 367–375.
- [23] Dimov A, Chandran SK, Punnekkat S. How do we collect data for software reliability estimation. In: Rachev B, Smrikarov A, eds. Proc. of the 11th Int'l Conf. on Computer Systems and Technologies and Workshop for Ph.D. Students in Computing on Int'l Conf. on Computer Systems and Technologies. ACM Press, 2010. 155–160.
- [24] Chen MH, Mathur AP, Rego VJ. Effect of testing techniques on software reliability estimates obtained using a time-domain model. IEEE Trans. on Reliability, 1995,44(1):97–103.
- [25] Pietrantuono R, Russo S, Trivedi KS. Online monitoring of software system reliability. In: Proc. of the 8th European Dependable Computing Conf. IEEE, 2010. 209–218.
- [26] Ding ZH, Chen MH, Li XX. Online reliability computing of composite services based on program invariants. Information Sciences, 2014,264:340–348.
- [27] Hangal S, Lam MS. Tracking down software bugs using automatic anomaly detection. In: Proc. of the 24th Int'l Conf. on Software Engineering. ACM Press, 2002. 291–301.
- [28] Sahoo SK, Li M, Ramachandran P, Adve SV, Adve VS, Zhou YY. Using likely program invariants to detect hardware errors. In: Proc. of the 38th Int'l Conf. on Dependable Systems and Networks (DSN). IEEE, 2008. 70–79.
- [29] Ernst MD, Cockrell J, Griswold WG, Notkin D. Dynamically discovering likely program invariants to support program evolution. IEEE Trans. on Software Engineering, 2001,27(2):99–123.
- [30] Ding ZH, Zhang K, Hu JL. A rigorous approach towards test case generation. Information Sciences, 2008,178:4057–4079.
- [31] Gupta N. Generating test data for dynamically discovering likely program invariants. In: Proc. of the ICSE 2003 Workshop on Dynamic Analysis. 2003. 21–24.
- [32] Do H, Elbaum S, Rothermel G. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. Empirical Software Engineering, 2005,10(4):405–435.
- [33] Jeske DR, Zhang XM, Pham L. Adjust software failure rates that are estimated from test data. IEEE Trans. on Reliability, 2005, 54(1):107–114.
- [34] Zhang XM, Teng XL, Pham H. Considering fault removal efficiency in software reliability assessment. IEEE Trans. on Systems, Man, and Cybernetics, Part A: Systems and Humans, 2003,33(1):114–120.

#### 附中文参考文献:

- [8] 赵靖,张汝波,顾国昌.考虑故障相关的软件可靠性增长模型研究.计算机学报,2007,30(10):1713–1720.
- [9] 李海峰,王栓奇,刘畅,郑军,李霞.考虑测试工作量与覆盖率的软件可靠性模型.软件学报,2013,24(4):749–760. <http://www.jos.org.cn/1000-9825/4257.htm> [doi: 10.3724/SP.J.1001.2013.04257]
- [17] 毛晓光,邓勇进.基于构件软件的可靠性通用模型.软件学报,2004,15(1):27–32. <http://www.jos.org.cn/1000-9825/15/27.htm>
- [18] 侯春燕,崔刚,刘宏伟.基于率的构件软件可靠性过程仿真.软件学报,2011,22(11):2749–2759. <http://www.jos.org.cn/1000-9825/3930.htm> [doi: 10.3724/SP.J.1001.2011.03930]



周远(1989—),男,浙江台州人,硕士生,主要研究领域为软件建模,模型检查,软件测试和可靠性.



丁佐华(1964—),男,教授,博士生导师,CCF高级会员,主要研究领域为软件测试与可靠性,软件建模与分析,软件自适应控制系统,智能计算及应用.