

# eXpress: Guided Path Exploration for Regression Test Generation

Kunal Taneja<sup>1</sup>, Tao Xie<sup>1</sup>, Nikolai Tillmann<sup>2</sup>, Jonathan de Halleux<sup>2</sup>, and Wolfram Schulte<sup>2</sup>

<sup>1</sup> Department of Computer Science, North Carolina State University, Raleigh, NC  
{ktaneja, txie}@ncsu.edu

<sup>2</sup> Microsoft Research, One Microsoft Way, Redmond, WA  
{nikolait, jhalleux, schulte}@microsoft.com

**Abstract.** Regression test generation aims at generating a test suite that can detect behavioral differences between the original and the new versions of a program. Regression test generation can be automated by using Dynamic Symbolic Execution (DSE), a state-of-the-art test generation technique. DSE explores all feasible paths in the program but exploration of all these paths can often be expensive. However, if our aim is to detect behavioral differences between two versions of a program, we do not need to explore all these paths in the program, since not all these paths are relevant for detecting behavioral differences. In this paper, we propose an approach on guided path exploration that avoids exploring irrelevant paths in terms of detecting behavioral differences. **In addition, our approach guides the path exploration towards infecting the program state by instrumenting the program with additional branches such that execution of these branches implies state infection.** Hence, behavioral differences are more likely to be detected earlier in path exploration. Furthermore, our approach leverages the existing test suite (if available) for the original version to efficiently execute the changed regions of the program and infect program states. Experimental results on 67 versions (in total) of four programs show that our approach requires about 36% fewer amount of time (on average) to detect behavioral differences than exploration without using our approach. In addition, our approach detects 6 behavioral difference that could not be detected by exploration without using our approach (within a time bound). Furthermore, our approach requires 67% less amount of time to find behavioral differences by exploiting an existing test suite than exploration without using the test suite.

## 1 Introduction

Regression test generation aims at generating a test suite that can detect behavioral differences between the original and the new versions of a program. A behavioral difference between two versions of a program can be reflected by the difference between the observable outputs produced by the execution of the same test (referred to as a difference-exposing test) on the two versions. Developers can inspect these behavioral differences to determine whether they are intended or unintended (i.e., regression faults).

Regression test generation can be automated by using Dynamic Symbolic Execution (DSE) [12, 28, 5], a state-of-the-art test generation technique, to generate a test suite achieving high structural coverage. DSE explores paths in a program to achieve high structural coverage, and exploration of all these paths can often be expensive. However, if our aim is to detect behavioral differences between two versions of a program, we do not need to explore all these paths in the program since not all these paths are relevant for detecting behavioral differences.

To formally investigate irrelevant paths for exposing behavioral differences, we adopt the Propagation, Infection, and Execution (PIE) model [33] of error propagation. According to the PIE model, a fault can be detected by a test if a faulty statement is executed (E), the execution of the faulty statement infects the state (I), and the infected state (i.e., error) propagates to an observable output (P). A change in the new version of a program can be treated as a fault and then the PIE model is applicable for effect propagation of the change. Our key insight is that many paths in a program often cannot help in satisfying any of the conditions P, I, or E of the PIE model.

In this paper, we present an approach<sup>3</sup> `eXpress` and its implementation that uses DSE to detect behavioral differences based on the notion of the PIE model.

Our approach first determines all the branches (in the program under test) that cannot help in achieving any of the conditions E, I, and P of the PIE model in terms of the changes in the program. To make test generation efficient, we develop a new search strategy for DSE to avoid exploring these irrelevant branches (including which can lead to an irrelevant path<sup>4</sup>). In particular, our approach guides DSE to avoid from flipping branching nodes<sup>5</sup>, which on flipping execute some irrelevant branch.

**In addition, to effectively achieve condition I of the PIE model, our approach instruments a program with additional branches, such that if any of these branches are executed, program state is infected.**

Furthermore, our approach can exploit the existing test suite (if available) for the original version by seeding the tests in the test suite to the program exploration. Our technique of seeding the exploration with the existing test suite efficiently augments an existing test suite so that various changed parts of the program (that are not covered by the existing test suite) are covered by the augmented test suite.

This paper makes the following major contributions:

**Path Exploration for Regression Test Generation.** We propose an approach called `eXpress` that uses DSE for efficient generation of regression tests. **`eXpress` prunes paths that cannot help in achieving conditions E, I, and P. In addition, `eXpress` instruments the program such that condition I can be achieved effectively after**

<sup>3</sup> An earlier version [30] of this work is described in a four-page paper that appears in the NIER track of ICSE 2009. This work significantly extends the previous work in the following major ways. First, in this paper, we develop techniques for efficiently finding irrelevant branches that cannot execute any change. Second, we develop techniques for exploiting the existing test suite for efficiently generating regression tests. **Third, we develop techniques for efficiently achieving state infection after a change is executed.** Fourth, we automate our approach by developing a tool. Fifth, we conduct extensive experiments to evaluate our approach.

<sup>4</sup> An irrelevant path is a path that cannot help in achieving P, I, and E of the PIE model.

<sup>5</sup> A branching node in the execution tree of a program is an instance of a conditional statement in the source code. A branching node consists of two sides (or more than two sides for a `switch` statement): the true branch and the false branch. Flipping a branching node is flipping the execution of the program from the true (or false) branch to the false (or true) branch. Flipping a branching node for a switch statement is flipping the execution of the current branch to another unexplored branch.

**condition E is achieved.**

**Incremental Exploration.** We develop a technique for exploiting an existing test suite, so that path exploration focuses on covering the changes rather than starting from scratch.

**Implementation.** We have implemented our `express` approach in a tool as an extension for Pex [31], an automated structural testing tool for .NET developed at Microsoft Research. Pex has been previously used internally at Microsoft to test core components of the .NET architecture and has found serious faults [31]. The current Pex has been downloaded tens of thousands of times in industry.

**Evaluation.** We have conducted experiments on 67 versions (in total) of four programs. Experimental results show that our approach requires about 36% fewer amount of time (on average) to detect behavioral differences than exploration without using our approach. In addition, our approach detects 6 behavioral difference that could not be detected by exploration without using our approach (within a time bound). Furthermore, our approach requires 67% less amount of time to find behavioral differences by exploiting an existing test suite than exploration without using the test suite.

## 2 Dynamic Symbolic Execution

In our approach, we use Pex [31] as an engine for Dynamic Symbolic Execution (DSE). Pex starts program exploration with some default inputs. Pex then collects constraints on program inputs from the predicates at the branching statements executed in the program. We refer to these constraints at branching statements as branch conditions. The conjunction of all branch conditions in the path followed during execution of an input is referred to as a path constraint. Pex keeps track of the previous executions to build a dynamic execution tree. Pex, in the next run<sup>6</sup>, chooses one of the unexplored branches in the execution tree (dynamically built thus far). Pex flips the chosen branching node to generate a new input that follows a new execution path. Pex uses various heuristics [36] for choosing a branching node (to flip next) using various search strategies with an objective of achieving high code coverage fast. We next present definitions of some terms that we use in the rest of this paper.

**Discovered node.** We refer to all the branching nodes that are explored in the current DSE run but were not explored in previous runs as discovered branching nodes (in short as discovered nodes).

**Instance of a branching node.** A branching node  $c_i$  in a Control Flow Graph (CFG) of a program can be present multiple times in the dynamic execution tree (of the program) due to loops in the program. We refer to these multiple branching nodes in the tree as instances of  $c_i$ .

## 3 Example

In this section, we illustrate our `express` approach with an example. `express` takes as input two versions of a program and produces as output a regression test suite, with

---

<sup>6</sup> A run is an exploration iteration.

```

static public int TestMe(char[] c, int n){
1  int state = 0;
2  if(c == null || c.Length == 0)
3      return -1;
4  for(int i=0; i< c.Length; i++){
5      if(c[i] == "[") state =1;
6      else if(state == 1 && c[i] == "{") state =2;
7      else if(state == 2 && c[i] == "<") state =3;
8      else if(state == 3 && c[i] == "+"){
9          state =4;
10         if(c.Length==15)
11             state = state + n;//Added in new version
12     }
13     if(c[i]==' ')
14         return;
15     if(!(c[i] >= 'a' && c[i] <= 'z')){
16         state=-1; return;
17     }
18 }
19 if(c[15] == '}')
20     return state;
21 return -1;
22 }

```

```

static public int fOld(int state,
int length, int n){
1  return state;
2 }

static public int fNew(int state,
int length, int n){
1  if(c.Length==15) state = state + n;
2  return state;
3 }

static public int TestMe(char[] c, int n){
...
10 int stateCopy = state;
11 state = fNew(state, c.Length, n);
12 stateCopy = fOld(stateCopy, c.Length, n);
13 if(state != stateCopy); //Infection Condition
...
}

```

**Fig. 1.** An example program on the left side and the instrumented program on the right.

the objective of detecting behavioral differences (if any exist) between the two versions of program under test. Although `eXpress` analyzes assembly code of C# programs, in this section, we illustrate the `eXpress` approach using program source code.

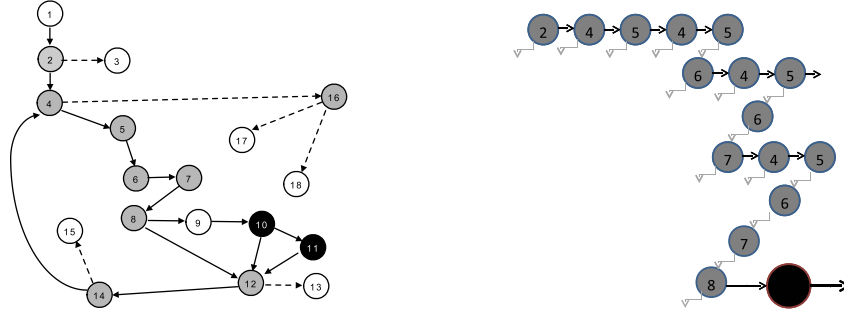
Consider the example in Figure 1. The left side of the figure shows a program `TestMe`. Lines 10 and 11 of the program are added in a new version.

**Finding Irrelevant Branches.** The left side of Figure 2 shows the CFG of the program in Figure 1. The labels of vertices in the CFG denote the corresponding line numbers in Figure 1. The black vertex denotes the newly added statements at Lines 10 and 11. The gray vertices denote the branching nodes (for the branching statements in the program), while the white vertices denote the other statements in the program. From the CFG, `eXpress` finds the following categories of branches.

- **Category  $B_{E+I}$ .** On traversing the CFG in Figure 2, `eXpress` detects that taking the branches  $\langle 2, 3 \rangle^7$ ,  $\langle 4, 16 \rangle$ ,  $\langle 16, 17 \rangle$ ,  $\langle 16, 18 \rangle$ ,  $\langle 12, 13 \rangle$  and  $\langle 14, 15 \rangle$  (dotted edges in Figure 2), the program execution cannot reach any of the black vertices. Hence, the execution of these branches cannot help in executing the changed statements or infecting the program state.
- **Category  $B_P$ .** In addition, `eXpress` finds out that among the source vertices of the six branches in Category  $B_{I+E}$ , there is no path from any of the black vertices to vertex 3. Hence, a state infection after the execution of the black vertex cannot propagate through the branch  $\langle 2, 3 \rangle$ .

**Instrumentation for State Infection.** A DSE-based test generation tool tries to execute all feasible branches in the program. However the execution of all branches is not sufficient for infecting the program state. For example, a test generation tool can cover all the branches of `TestMe` (left side of Figure 1) with  $n = 0$  (a default input used by some tools) but to infect the program state, the tool need to generate a nonzero value of  $n$ . To guide the test generation tool to effectively achieve state infection, `eXpress` instruments the new program version. In particular, `eXpress` invokes the changed re-

<sup>7</sup>  $\langle x, y \rangle$  represents a branch from source vertex  $x$  to destination vertex  $y$



**Fig. 2.** The left side of the figure shows the part of execution tree of the program for  $c = \{ "[", "\{", "<", "*" \}$ , while the right side shows the Control-Flow Graph for the program in Figure 1.

regions<sup>8</sup> from the two program versions (with same inputs) and compares the outputs of the two changed regions. *eXpress* inserts branches in the new program version that are executed only when the outputs of the executions of the changed regions differ, i.e., when program state is infected after the execution of the changed regions.

The right side of Figure 1 shows the instrumented version of *TestMe* on the left side of Figure 1. *eXpress* identifies *state*, *length*, and *n* as the inputs to the two changed regions and *state* as the output to the two changed regions. *eXpress* then extracts the changed regions in the original and new program versions as separate methods (*fOld* and *fNew*, respectively). The parameters of the two methods are same as the identified inputs, while the return values are same as the identified output. *eXpress* replaces the changed region of the new version of *TestMe* by the invocation of the methods *fOld* and *fNew* (Lines 11 and 12 on the right side of Figure 1). Note that the variable *stateOld* passed to *fOld* is a copy of the variable *state* passed to *fNew* (Line 10 on the right side of 1) because *state* may be updated by *fNew* since *state* is also an output of the changed regions. *eXpress* then inserts an *if* statement (Line 13 on the right side of Figure 1) in the program comparing the outputs of the two methods. Note that the outputs of the two methods are different if and only if the program state is infected. A DSE-based test generation tool determines that it needs to generate inputs satisfying  $c.Length == 15$  AND  $state \neq state + n$  to cover the true branch of the added *if* statement. Hence, the tool generates a nonzero value of *n* to satisfy the constraint and infect the program state.

**Test Generation.** To cover the changed statements at Lines 10 and 11 (in the left part of Figure 1) and infect the program state, DSE needs at least 6 DSE runs (starting from an empty input array *c*). However, the number of runs depends on the choice of branches that DSE flips in each run. In each run, DSE has the choice of flipping 8 branches in the program. For *TestMe*, Pex takes 441 DSE runs to cover the true branch of statement at Line 10. The number of runs can be much more if the number of branching statements in the program increases. To reduce the search space of DSE, *eXpress* prunes branches in category  $B_{E+I}$  until the black vertex is executed and the program state is infected

<sup>8</sup> A changed region (in a new or modified program version) is a minimal set of statements that contains all modified (in new or original program version), added (in new program version), or deleted (in original program version) statements in a method. The changed region for the new program version in *TestMe* on the left side of Figure 1 is shown in rectangular box, while the changed region of the original program version is empty.

after its execution. Once the program state is infected, `eXpress` prunes out the branches in category  $B_P$  to efficiently propagate the state infection to an observable output.

**Incremental Exploration.** `eXpress` can reuse an existing test suite for the original version so that changed parts of the program can be explored efficiently due to which test generation is likely to find behavior differences earlier in path exploration. Assume that there is an existing test suite covering all the statements in the old version of the program in Figure 1. Suppose that the test suite has a test input  $I = \{ "[", "{", "<", "*" \}$ . The input covers all the statements in the new version of `TestMe` except the newly added statement at Line 11. If we start the program exploration from a scratch (i.e., with default inputs), Pex takes 441 runs to cover the statement at Line 11. However, we can reuse the existing test suite for exploration to cover the new statement efficiently. Our approach executes the test suite to build an execution tree for the tests in the test suite. Our approach then starts program exploration using the dynamic execution tree built by executing the existing test suite instead of starting from an empty tree. Some branches in the tree may take many runs for Pex to discover if starting from an empty tree. Figure 2 shows part of the execution tree for the input  $I$ . A gray edge in the tree indicates the false side of a branching node while a black (horizontal) edge indicates the true side. To generate an input for the next DSE runs, Pex flips a branching node  $b$  in the tree whose other side has not yet been explored and generates an input so that program execution takes the unexplored branch of  $b$ . Pex chooses such branching node for flipping using various heuristics for covering changed regions of the program. It is likely that Pex chooses the branching node 13 (colored black), which on execution covers the `break` statement at Line 13. When starting the program exploration from scratch, Pex takes 420 runs before it reaches (discovers) the black branching node in Figure 2. Using our approach of seeding the tests from the existing test suite, Pex takes 39 runs to flip the branching node and cover the statement at Line 11.

## 4 Approach

`eXpress` takes as input the assembly code of two versions  $v1$  (original) and  $v2$  (new) of the program under test. In addition, `eXpress` takes as input, the name and signature of a parameterized unit test (PUT); such a PUT serves as a test driver for path exploration. When an existing test suite is available for the original version, `eXpress` conducts incremental exploration that exploits the test suite for generating tests for the new version. We next discuss in detail the `eXpress` approach.

### 4.1 Finding Differences Between Versions

`eXpress` analyzes the two versions  $v1$  and  $v2$  to find pairs  $\langle M_{i1}, M_{i2} \rangle$  of corresponding methods in  $v1$  and  $v2$ , where  $M_{i1}$  is a method in  $v1$  and  $M_{i2}$  is a method in  $v2$ . A method  $M$  is defined as a triple  $\langle FQN, Sig, I \rangle$ , where  $FQN$  is the fully qualified name<sup>9</sup> of the method,  $Sig$  is the signature<sup>10</sup> of the method, and  $I$  is the list of assembly

<sup>9</sup> The fully qualified name of a method  $m$  is a combination of the method's name, the name of the class  $c$  declaring  $m$ , and the name of the namespace containing  $c$ .

<sup>10</sup> Signature of a method  $m$  is the combination of parameter types of  $m$  and the return type of  $m$ .

instructions in the method body. Two methods  $\langle M_{i1}, M_{i2} \rangle$  form a corresponding pair if the two methods  $M_{i1}$  and  $M_{i2}$  have the same *FQN* and *Sig*. Currently our approach considers as different methods the methods that have undergone Rename Method or Change Method Signature refactoring. A refactoring detection tool [7] can be used to find such corresponding methods. For each pair  $\langle M_{i1}, M_{i2} \rangle$  of corresponding methods, the Difference Finder finds a set of differences  $\Delta_i$  between the list of instructions  $I_{M_{i1}}$  and  $I_{M_{i2}}$  in the body of Methods  $M_{i1}$  and  $M_{i2}$ , respectively.  $\Delta_i$  is a set of instructions such that each instruction  $\iota$  in  $\Delta_i$  is an instruction in  $I_{M_{i2}}$  (or in  $I_{M_{i1}}$  for a deleted instruction), and  $\iota$  is added, modified, or deleted from list  $I_{M_{i1}}$  to form  $I_{M_{i2}}$ . We denote the set of all added, modified, or deleted instructions as  $\Delta$ .

## 4.2 Finding Irrelevant Branches.

`express` efficiently constructs the inter-procedural CFG  $g \langle V, E \rangle$  of the new version of the program under test such that each vertex  $v \in V$  corresponds to an instruction  $\iota \in M$  (denoted as  $v \leftrightarrow \iota$ ), where  $M$  is some method in  $v_2$ . `express` traverses the graph  $g$  to find a set of branches  $B_{E+I}$  via which the execution cannot reach any of the branches in  $V$  and a set of branches  $B_P$  via which a state infection cannot propagate to any observable output. A branch  $b$  in CFG  $g$  is an edge  $e = \langle v_i, v_j \rangle$ :  $e \in E$ ,  $v_i \in V$  with an outgoing degree of more than one. The vertex  $v_i$  is referred to as a branching node. We next describe the sets  $B_{E+I}$  and  $B_P$ .

Let  $V = \{v_1, v_2, \dots, v_l\}$  be the set of all vertices in CFG  $g \langle V, E \rangle$  such that  $v_i \in V$  and  $v_i.degree > 1$ . Let  $E_i = \{e_{i1}, e_{i2}, \dots, e_{im}\}$  be the set of outgoing edges (branches) from  $v_i$ . Let  $C$  be the set of vertices in the CFG  $g$  such that  $\forall v \in C, \exists \iota \in \Delta : v \leftrightarrow \iota$ .  $\rho(v_i, v_j, e_{ij})$  denotes a path between a source vertex  $v_i$  to a destination vertex  $v_j$  that takes the branch  $e_{ij}$  (if  $\rho(v_i, v_j, e_{ij}) = \phi$ , there is no such path from  $v_i$  to  $v_j$ ),  $\rho(v_i, v_j)$  denotes a path from a source vertex  $v_i$  to a destination vertex  $v_j$  (if  $\rho(v_i, v_j) = \phi$ , there is no such path from  $v_i$  to  $v_j$ ).

**Branches  $B_{E+I}$ .**  $B_{E+I} \subseteq E$  is a set of branches such that  $\forall e_{ij} = \langle v_i, v_j \rangle \in B_{E+I} \wedge \forall c_k \in C : \rho(v_i, c_k, e_{ij}) = \phi$ . Since  $\rho(v_i, c_k, e_{ij}) = \phi$ , after taking a branch  $e_{ij} \in B_{E+I}$ , the program execution cannot reach a changed region. Hence, the program state cannot be infected.

**Branches  $B_P$ .**  $B_P \subseteq E$  is a set of branches such that  $\forall e_{ij} = \langle v_i, v_j \rangle \in B_P \wedge \forall c_k \in C : \rho(v_i, c_k, e_{ij}) = \phi \wedge \rho(c_k, v_i) = \phi$ . Since  $\rho(c_k, v_i) = \phi$ , a state infection after a changed vertex  $c_k$  cannot reach  $v_i$ . Hence, the state infection cannot propagate through  $e_{ij}$ .

## 4.3 Code Instrumentation

For each changed method pair  $\langle M_{i1}, M_{i2} \rangle$  for which  $\Delta_i \neq \phi$ , `express` finds changed regions  $\delta_{i1}$  and  $\delta_{i2}$  containing all the changed instructions in the program.  $\delta_i$  is a minimal list of continuous instructions such that all the changed instructions in the method  $M_i$  are in the region  $\delta_i$ . Hence, there can be a maximum of one changed region in one method. Inputs  $I = \{i_1, i_2, \dots, i_n\}$  of a changed region  $\delta$  are variables (or fields) that are used in the changed region  $\delta$  but are declared outside the changed region, while outputs  $O = \{o_1, o_2, \dots, o_m\}$  are variables (or fields) that are defined in  $\delta$  and used outside the



```

public Co Mδ1(i11, i21, i31..., in1){
    

---


    Co co = new Co();
    co.o1 = o11;
    co.o2 = o21;
    ...
    co.om = om1;
    return co;
}

public Co Mδ2(i12, i22, i32..., in2){
    

---


    Co co = new Co();
    co.o1 = o12;
    co.o2 = o22;
    ...
    co.om = om2;
    return co;
}

public class Co{
    public o1;
    public o2;
    ...
    public on;
    public bool Equals(){...}
}

public Mi2(...){
    ...
    1  i1c = copy(i1);
    2  i2c = copy(i2);
    3  ...
    4  inc = copy(in);
    5  coOld = Mδ1(i1c, i2c, ..., inc);
    6  coNew = Mδ2(i1, i2, i3..., in);
    7  if(!coOld.Equals(coNew)){
        // Infection
    }
    8  o1 = coNew.o1
    9  o2 = coNew.o2
    ...
    10 on = coNew.on
    ...
}

```

**Fig. 3.** Instrumented Code synthesized by `express`

changed region. As an example, for the changed region in the new version of `TestMe` in Figure 1,  $I = c.Length, n, state$  and  $O = state$ . Let  $I_1$  and  $O_1$  be the inputs and outputs of  $\delta_1$ , while  $I_2$  and  $O_2$  be the inputs and outputs of  $\delta_2$ .

`express` extracts the instructions in  $\delta_{i1}$  and  $\delta_{i2}$  as methods  $M_{\delta_{i1}}$  and  $M_{\delta_{i2}}$  (shown in Figure 3), respectively, both having arguments of type and name as  $I_1 \cup I_2$  and returning outputs  $O_1 \cup O_2$ . We refer to the arguments of  $M_{\delta_{i1}}$  as  $A_1 = i_{11}, i_{21}, \dots, i_{k1}$  and the arguments of  $M_{\delta_{i2}}$  as  $A_2 = i_{12}, i_{22}, \dots, i_{k2}$ , such that the types of  $i_{l1}$  and  $i_{l2}$  are same as the type of  $i_l \in I_1 \cup I_2$ . We refer to the variables corresponding to  $o_i \in O_1 \cup O_2$  in  $M_{\delta_{i1}}$  as  $o_{i1}$  and in  $M_{\delta_{i2}}$  as  $o_{i2}$ . `express` synthesizes a class  $C_o$  (show in Figure 3), such that for all  $o_i \in O_1 \cup O_2$ , there is a corresponding public field in  $C_o$  with types and names same as  $o_i$ . The methods  $M_{\delta_{i1}}$  and  $M_{\delta_{i2}}$  return an object  $c_o$  of type  $C_o$ . The class enables the two methods  $M_{\delta_{i1}}$  and  $M_{\delta_{i2}}$  to return multiple outputs  $O_1 \cup O_2$  as fields of  $C_o$ . As shown in Figure 3, the method  $M_{\delta_{i1}}$  (and  $M_{\delta_{i2}}$ ) first contains the instructions in  $\delta_{i1}$  and later copies all the outputs  $o_{i1}$  in  $M_{\delta_{i1}}$  ( $o_{i2}$  in  $M_{\delta_{i2}}$ ) to the respective fields of  $c_o$  and returns the object  $c_o$ .

Figure 3 shows the instrumented part of the method  $M_{i2}$ . The methods  $M_{\delta_{i1}}$  and  $M_{\delta_{i2}}$  are invoked (Lines 5 and 6 of the Method  $M_{i2}$  in Figure 3) from the method  $M_{i2}$  instead of the changed region. The method  $M_{\delta_{i2}}$  is passed  $i_1, \dots, i_n$  as arguments, while the method  $M_{\delta_{i1}}$  is passed copied values  $i_{1c}, \dots, i_{nc}$  of  $i_1, \dots, i_n$ , respectively (Lines 1 to 4). These copies are deep copies for non primitive type variables. The outputs of  $M_{\delta_{i1}}$  and  $M_{\delta_{i2}}$  are stored in objects  $c_oOld$  and  $c_oNew$ , respectively. An `if` statement is inserted (Line 7) such that its true branch is executed when the outputs  $c_oOld$  and  $c_oNew$  are different. To compare the two outputs, `express` changes the modifier of all the fields transitively reachable from objects of the two given classes to public and these fields are directly compared to detect difference in outputs. Finally, the fields of  $c_oNew$  are assigned to the corresponding output variables (Lines 8 to 10).



#### 4.4 Test Generation

`eXpress` performs Dynamic Symbolic Execution (DSE) [6, 17, 12, 28, 5] to generate regression tests for the two given versions of a program. DSE iteratively generates test inputs to cover various feasible paths in the program under test (the new version in our approach). In particular, DSE flips some branching node discovered in previous executions to generate a test input for covering a new path. The branching node to be flipped is decided by a search strategy such as depth-first search. The exploration is quite expensive since there are an exponential number of paths with respect to the number of branches in a program. However, the execution of many branches often cannot help in detecting behavioral differences. In other words, covering these branches does not help in satisfying any of the condition E or I in the PIE model described in Section 1. Therefore, we do not flip such branching nodes in our new search strategy for generating test inputs that detect behavioral differences between the two given versions of a program. Recall that, we refer to such branches as irrelevant branches. We next describe the two categories of paths that our approach avoids exploring. We then describe our incremental exploration technique when an existing test suite is available for the original version.

**Path Pruning** Our approach avoids exploring the following categories of paths:

- **Category E+I.** `eXpress` does not explore these paths until a test is generated that infects the program state. Our approach avoids exploring all the branches in  $B_{E+I}$ . In other words, if a branching node  $v_i$  is in the built dynamic execution tree<sup>11</sup>, such that the branch  $e_{ij} \in B_{E+I}$  is not explored yet, our approach does not flip  $v_i$  to make program execution take branch  $e_{ij}$ . As an effect, our approach helps reduce the number of paths to be explored by avoiding the paths  $P \in g : e_{ij} \in P$  until a test is generated that infects the program state.
- **Category P.** `eXpress` does not explore these paths after a test is generated that infects the program state. Our approach avoids exploring all the branches  $e_{ij} \in B_p$ . In other words, if a branching node  $v_i$  is in the built dynamic execution tree such that one of its branch  $e_{ij}$  is not explored yet, our approach does not flip  $v_i$  to make program execution take branch  $e_{ij}$ . As an effect, our approach helps reduce the number of paths to be explored by avoiding the paths  $P \in g : e_{ij} \in P$ . These paths cannot help in propagating a state infection to an observable output.

#### 4.5 Incremental Exploration

A regression test suite achieving high code coverage may be available along with the original version of a program. This test suite may be manually written or generated by an automated test generation tool for the original version. However, the existing test suite might not be able to cover all the changed regions of the new version of the program. Our approach can reuse the existing test suite so that changed regions of the

<sup>11</sup> A dynamic execution tree is the tree formed from the paths executed in the previous executions. Multiple instances  $c_{i1}, c_{i2}, \dots, c_{in}$  of a node  $c_i$  in CFG can be present in a dynamic execution tree due to loops in a program.

program can be executed efficiently due to which test generation is likely to find behavior differences earlier in path exploration. Our approach executes the existing test suite to build an execution tree for the tests in the test suite. Our approach then starts the program exploration using the dynamic execution tree instead of starting from an empty tree. Our approach of seeding test inputs can help efficiently cover the changed regions of the program with two major reasons:

**Discovery of hard-to-discover branching nodes.** By seeding the existing test suite for Pex to start exploration with, our approach executes the test suite to build an execution tree of the program. Some of the branching nodes in the built execution tree may take a large number of DSE runs (without seeding any tests) to get discovered. Flipping some of these discovered branching nodes nearer in CFG to the changed parts of the program has more likelihood of covering the changed regions of the program [4]. Although, our approach currently does not specifically prioritize flipping of branching nodes near the changed regions, our approach can help these branching nodes to get discovered (by executing the existing test suite), which might take large number of DSE runs as shown in the example in Section 3.

**Priority of DSE to cover not-covered regions of the program.** DSE techniques employ branch prioritization so that a high coverage can be achieved faster due to which DSE techniques choose a branch from the execution tree (built thus far) that have a high likelihood of covering changed regions (that are not covered by existing test suite for the original version). By seeding the existing test suite to program exploration, the DSE techniques do not waste time on covering the regions of the program already covered by the existing test suite. Instead, the DSE techniques give high priority to branching nodes that can cover not-covered regions of the program, which include the changed parts. Hence, the changed parts are likely to be covered earlier in path exploration.

## 5 Experiments

We conducted experiments on four programs and their 67 versions (in total) collected from three different sources. In our experiments, we try to address the following research questions:

**RQ1.** How many fewer DSE runs does Pex require to execute the changed regions between the two versions of a program with the assistance of `eXpress`?

**RQ2.** How many fewer DSE runs does Pex require to infect the program states with the assistance of `eXpress`?

**RQ3.** How many fewer DSE runs and how much fewer amount of time does Pex require to find behavioral differences with the assistance of `eXpress`?

**RQ4.** How many fewer DSE runs and how much fewer amount of time does Pex require to find behavior differences when the program exploration is seeded with the existing test suite?

### 5.1 Subjects

To answer the research questions, we conducted experiments on four subjects. Table 1 shows the details about the subjects. Column 1 shows the subject name. Column 2

shows the number of classes in the subject. Column 3 shows the number of classes that are covered by tests generated in our experiments. Column 4 shows the number of versions (not including the original version) used in our experiments. Column 5 shows the number of lines of code in the subject.

`replace` and `siena` are programs available from the Subject Infrastructure Repository (SIR) [8]. `replace` and `siena` are written in *C* and *Java*, respectively. `replace` is a text-processing program, while `siena` is an Internet-scale event notification program. We chose these two subjects (among the others available at the SIR) in our experiments we could convert these subjects into C# using the Java 2 CSharp Translator<sup>12</sup>. We could not convert other subjects available at the SIR (with the exception of `tcas`) because of extensive use of *C* or *Java* library APIs in these subjects. The experimental results on `tcas` are presented in a previous version of this work [30] and show similar conclusions as the results from the subjects used in the experiments here. We seeded all the 32 faults available for `replace` at the SIR one by one to generate 32 new versions of `replace`. For `siena`, SIR contains 8 different sequentially released versions of `siena` (versions 1.8 through 1.15). Each version provides enhanced functionalities or corrections with respect to the preceding version. We use these 8 versions in our experiments. In addition to these 8 versions, there are 9 seeded faults available at SIR. We seeded all the 9 faults available at SIR one by one to synthesize 9 new versions of `siena`. In total, we conduct experiments on these 17 versions of `siena`. For `replace`, we use the main method as a PUT for generating tests. We capture the concrete value of the string `sub` at the end of the *PUT* using `PexStore.ValueForValidation("v", v)` statement. This statement captures the current value of `v` in a particular run (i.e., an explored path) of DSE. In particular, this statement results in an assertion `Assert.AreEqual(v, cv)` in a generated test, where `cv` is the concrete value of `v` in the test during the time of exploration. This assertion is used to find behavioral differences when the tests generated for a new version are executed on the original version.

For `siena`, we use the methods `encode` (for changes that are transitively reachable from `encode`) and `decode` (for changes that are transitively reachable from `decode`) in the class `SENP` as PUTs for generating tests. We capture the return values of these methods using the `PexStore` statement in the PUTs.

The method `encode` requires non-primitive arguments. Existing Pex cannot handle non-primitive argument types effectively but provides support for using factory methods for non-primitive types. Hence, we manually wrote factory methods for the non-primitive types in `SENP`. In particular, we wrote factory methods for classes `SENPpacket`, `Event` and `Filter`. Each factory method invokes a sequence (of length up to three) of the public state-modifying methods in the corresponding class. The parameters for these methods, and the length of the sequence (up to three) are passed as inputs to the factory methods. During exploration, Pex generates concrete values for these inputs to cover various parts of the program under test.

STPG<sup>13</sup> is an open source program hosted by the codeplex website, Microsoft's open source project hosting website. The codeplex website contains snapshots of check-ins in the code repositories for STPG. We collect three different versions of the subject

<sup>12</sup> <http://sourceforge.net/projects/j2cstranslator/>

<sup>13</sup> <http://stringtopathgeometry.codeplex.com/>

STPG from the three most recent check-ins. We use the `Convert(string path)` method as the PUT for generating tests since `Convert` is the main conversion method that converts a string path data definition to a `PathGeometry` object. We capture the return value of `Convert` using the `PexStore` statement in the PUTs.

`structorian`<sup>14</sup> is an open source binary data viewing and reverse engineering tool. `structorian` is hosted by Google’s open source project hosting website. The website also contains snapshots of check-ins in the code repositories for `structorian`. We collected all the versions of snapshots for the classes `StructLexer` and `StructParser`. We chose these classes in our experiments due to three factors. First, these classes have several revisions available in the repository. Second, these classes are of non-trivial size and complexity. Third, these classes have corresponding tests available in the repository. For the classes `StructLexer` and `StructParser`, we generalized one of the available concrete test methods by promoting primitive types to arguments of the test methods. Furthermore, we convert the assertions in the concrete test methods to `PexStore` statements. For example if an assertion `Assert.AreEqual(v, 0)` exist in a concrete test, we convert the assertion to `PexStore.ValueForValidation("v", v)`. We use these generalized test methods as PUTs for our experiments. `structorian` contains a manually written test suite. We use this test suite for seeding the exploration for addressing the question RQ4.

To address questions RQ1-RQ3, we use all the four subjects, while to address question RQ4, we use `structorian` because of two major factors. First, `structorian` has a manually written test suite that can be used to seed the exploration. Second, revisions of `structorian` contains non-trivial changes that cannot be covered by the existing test suite. Hence, our technique of seeding the existing test suite in the program exploration is useful for covering these changes. `replace` contains changes to one statement due to which most of the changes can be covered by the existing test suite. Hence, our Incremental Exploration technique is not beneficial for the version pairs under test. `siena` and STPG do not have an existing test suite to use.

**Table 1.** Experimental subjects

Project	Classes	Classes Covered	Versions	LOC
replace	1	1	32	625
STPG	1	1	2	684
siena	6	6	17	1529
structorian	70	8	21	6561

## 5.2 Experimental Setup

For `replace` and `siena`, we conduct regression test generation between the original version and each version  $v_2$  synthesized from the available faults in the SIR. We use `express` and the default search strategy in Pex [31, 36] to conduct regression test generation. In addition to the versions synthesized by seeding faults, we also conduct regression test generation between each successive versions of `siena` (versions 1.8 through

<sup>14</sup> <http://code.google.com/p/structorian/>

1.15) available in SIR, using `eXpress` and the default search strategy in Pex [31, 36]. For STPG and `structorian`, we conduct regression test generation between two successive version pairs that we collected.

To address RQ1, we compare the number of runs of DSE required by the default search strategy in Pex (in short as Pex) with the number of runs required by Pex enhanced with `eXpress` (in short as Pex+`eXpress`) to execute a changed region. To address RQ2, we compare the number of runs required by Pex with the number of runs required by Pex+`eXpress` to infect the program states. To address RQ3, we compare the number of runs and the amount of time required by Pex with the number of runs required by Pex+`eXpress` to find behavioral differences between two versions of a program under test. To address RQ4, we compare the number of DSE runs and the amount of time required by Pex (and Pex+`eXpress`) to find behavioral differences with and without seeding the program exploration (with the existing test suite).

Currently, we have not automated the Instrumenter component of `eXpress`. To find out whether the program state has been infected after the execution of a changed region, we manually add statements to the code under test that send a callback to the Dynamic Test Generator whenever the program state is infected after the execution of a changed region. The rest of the approach is fully automated and is implemented in a tool as an extension<sup>15</sup> to Pex [31]. We developed its components to statically find irrelevant branches as a .NET Reflector<sup>16</sup> AddIn. In future, we plan to automate the Instrumenter component.

To find behavioral differences between two versions, we execute the tests generated for a new version on the original version. Behavioral differences are detected by a test if an assertion in the test fails.

### 5.3 Experimental Results

Table 2 shows the experimental results. Due to space limit, we provide only the total, average, and median values for the subjects `replace`, `siena`, and STPG. The detailed results for experiments on all the versions of these subjects are available on our project web<sup>17</sup>. However, we provide detailed results for `structorian` in this paper.

Column *S* shows the name of the subject. For `structorian`, the column shows the class name. The class `StructLexer` is denoted by SL and the class `StructParser` is denoted by SP. Column *V* shows the number of version pairs for which we conducted experiments for the subject. For `structorian`, the column shows the version numbers on which the experiments were conducted. These version numbers are the revision numbers in the google code repository of `structorian`. Column  $E_{Pex}$  shows the total number of DSE runs required by Pex for satisfying E. Column  $E_{Red}$  shows the average percentage reduction in the number of DSE runs by Pex+`eXpress` for achieving E. Column  $M_e$  shows the median percentage reduction in the number of DSE runs by Pex+`eXpress` for achieving E. Column  $I_{Pex}$  shows the total number of DSE runs required by Pex for satisfying I. Column  $I_{Red}$  shows the average percentage reduction

<sup>15</sup> <http://pexase.codeplex.com/>

<sup>16</sup> <http://www.red-gate.com/products/reflector/>

<sup>17</sup> <https://sites.google.com/site/asergpr/projects/express/>

**Table 2.** Experimental results

		Execution			Infection			Propagation					
S	V	$E_{Pex}$	$E_{Red}(\%)$	$M_i$	$I_{Pex}$	$I_{Red}(\%)$	$M_e$	$P_{Pex}$	$P_{Red}(\%)$	$M_p$	$T_{Pex}$	$Ts+T_d$	$T_{PRed}$
replace	32	1668	49	0	2740	42	34	10312	75	49	711	235	67
siena	17	286	42	13	284	39.4	13	7301	42	15	1011	628	38
STPG	2	341	26	26	378	32		378	32	32	353	255	28
Total	51	2573	49		3402	41		17613	62		1722	863	50
-structorian-													
SL	2-9	102	26.5	-	102	26.5	-	-	-	-	-	-	-
SL	9-139	102	26.5	-	102	26.5	-	2988	66	-	99	69.3	68
SL	139-150	102	26.5	-	102	26.5	-	761	69	-	26	7.5	71
SL	150-169	53	13.2	-	53	13.2	-	299	52	-	7.4	3.9	47
SL	169-174	55	12.7	-	55	12.7	-	478	32.2	-	14.2	8.4	41
SL	174-175	102	26.5	-	102	26.5	-	-	-	-	-	-	-
Total(SL)		516	24	26.5	516	26.5	24	4526	62	52	146.6	89.1	39
SP	2-5	10000*	81	-	10000*	74	-	10000*	74	-	1hr	35min	41.7
SP	5-6	10000*	74	-	10000*	74	-	10000*	74	-	1hr	32min	47
SP	9-13	10000*	81	-	10000*	81	-	10000*	81	-	1hr	27min	55
SP	37-39	6	0	-	3699	77	-	3699	77	-	26min	22min	15
SP	39-40	2	0	-	-	-	-	-	-	-	-	-	-
SP	50-62	6188	82.9	-	6188	82.9	-	6188	82.9	-	35 min	21	40
SP	45-47	-	-	-	-	-	-	-	-	-	-	-	-
SP	47-50	2	0	-	2	0	-	-	-	-	-	-	-
SP	62-124	6	0	-	10000*	28	-	10000*	28	-	1hr*	58min	2
SP	40-45	-	-	-	-	-	-	-	-	-	-	-	-
Total(SL)		36206	79	74	49889	68	77	49889	68	77	5hr	3.25hr	35

in the number of DSE runs by Pex+eXpress for achieving I. Column  $M_i$  shows the median percentage reduction in the number of DSE runs by Pex+eXpress for achieving I. Column  $P_{Pex}$  shows the total number of DSE runs required by Pex for satisfying P. Column  $p_{Red}$  shows the average percentage reduction in the number of DSE runs by Pex+eXpress for achieving P (i.e., finding behavioral differences). Column  $M_p$  shows the median percentage reduction in the number of DSE runs by Pex+eXpress for achieving P. Column  $T_{Pex}$  shows the time taken by *Pex* for satisfying P. Column  $T_s + T_d$  shows the time taken *Pex* + *eXpress* for satisfying P. This time includes the time taken to statically find irrelevant branches. Column  $T_{PRed}$  shows the average percentage reduction in amount of time taken by Pex+eXpress for achieving P.

Table 3 shows the time taken for finding the irrelevant branches and the number of irrelevant branches found. Column *S* shows the subject. Column  $T_{static}$  shows the average time taken by *eXpress* to find irrelevant branches that cannot help in satisfying E of the PIE model. Column  $B_{E+I}$  shows the average number of branches found in the set  $B_{E+I}$ . In general, irrelevant branches are more if changes are towards the beginning of the PUT since there are likely to be more branches in the program that do not have a path to any changed regions. These branches also include the branches whose branching conditions are not dependent on the inputs of the program and therefore do not correspond to branching conditions during path exploration. Hence, pruning these branches is not helpful in making DSE efficient. Column  $B_P$  shows the average number of branches found in the set  $B_P$ . Column  $B_{Tot}$  shows the total number of branches in the CFG.

**Results of replace.** For the *replace* subject, among the 32 pairs of versions, the changed regions cannot be executed for 4 of these version pairs (version pairs 14, 18, 27, and 31) by Pex or by Pex+eXpress in 1000 DSE runs. We do not include these

version pairs while calculating the sum of DSE runs for satisfying I and E of the PIE model. For 3 of the version pairs (version pairs 12, 13, and 21), the changes are in the fields due to which there are no benefits of using Pex+eXpress. We exclude these three version pairs from the experimental results shown in Table 2, which includes the results of 32 version pairs. For 3 of the version pairs (version pairs 3, 22 and 32), a changed region was executed but the program state is not infected (both by Pex and Pex+eXpress) in a bound of 5 minutes. We do not include these version pairs while calculating the sum of DSE runs for satisfying I of the PIE model. In addition, for 6 of the version pairs, the state infection was not propagated to observable output within a bound of 5 minutes. We do not include these version pairs while calculating the sum of DSE runs for satisfying P of the PIE model.

In total, Pex+eXpress took 49% fewer runs in executing the changes with a maximum of 77.6% for versions pairs 23 and 24. For these version pairs, Pex+eXpress takes 95 DSE runs in contrast to 425 runs taken by Pex to execute the changed locations. For many version pairs, there was no benefit of using Pex+eXpress in terms of satisfying E of the PIE model. As a result, median reduction in the number of runs is 0% for satisfying E of the PIE model.

Pex+eXpress took 46% fewer runs, in infecting the program state, with a maximum of 73.8% for version pair 6. For this version, eXpress takes 83 DSE runs in contrast to 317 runs taken by Pex to infect the program state after the execution of changed regions. We observe that Pex+eXpress took 75% fewer runs (median 49%) and 67% fewer amount of time in finding behavioral differences. This difference is substantially larger than the reduction in runs to achieve I. This phenomenon is due to a large number of exception paths (i.e., paths that lead to an exception statement such as `throw`) in `replace`. As a result, state infections often do not propagated to an assertion violation in the PUT due to exceptions thrown in the program `replace`. eXpress helps in pruning these irrelevant paths that lead to an exception.

**Results of *siena*.** We observe that the behavioral differences between five of the version pairs of *siena* are found within ten runs by Pex and Pex+eXpress. For these version pairs, there is no reduction in the number of runs. The reason for the preceding phenomenon is that changes in these version pairs are close to the entry vertex in the CFG. Hence, these changes can be covered in a relatively small number of runs. In two of the version pairs, changed regions were not covered by either Pex+eXpress or Pex. An exception is thrown by the program before these changes could be executed. Pex and Pex+eXpress were unable to generate a test input to avoid the exception. Changes between two of the version pairs were refactorings due to which the program state is never infected.

For two of the changes, behavioral differences could not be detected by Pex within a bound of 5 minutes but they were detected by Pex+eXpress. We use 5 minutes for calculating the total values in the column  $T_{pex}$  and the number of DSE runs performed during five minutes in the column  $P_{pex}$ .

In summary, Pex+eXpress executed the changed region in 42% fewer runs (median 13%), infected the program state in 39.4% (median 13%) fewer runs, and found behavioral differences in 42% fewer runs (median 15%) and 38% fewer amount of time than Pex. In addition, Pex+eXpress detected behavioral differences for two changes that



were not detected by Pex within a bound of 5 minutes.

**Results of structorian.** The six rows with *SL* in column *S* of Table 2 show the experimental results for changes in the class `StructLexer`, while the last 10 rows (with *SP* in column *S*) show the experimental results on versions of the class `StructParser`. For the versions of `StructLexer`, Pex+eXpress takes 24% fewer runs (median 26.5%) to execute a changed region than Pex. In addition, Pex+eXpress infects the program state in 24% fewer runs (median 26.5%). Neither Pex nor Pex+eXpress were able to find behavioral differences for two versions. Pex+eXpress takes 62% fewer runs (median 52%) and 39% fewer amount of time to find behavioral differences.

Neither Pex+eXpress nor Pex was able to find behavioral differences between some version pairs of class `StructParser` in 5 minutes (a bound that we use in our experiments for all subjects). For these version pairs, we increased the bound to 1 hour (or 10000 runs). Pex was not able to find behavioral differences for 4 version pairs even in 1 hour, while Pex+eXpress found behavioral differences for all these version pairs. If Pex was unable to detect behavioral differences within the bound of 1 hour, we put the time in the column  $T_{pex}$  as 1 hour and the number of runs as 10000 (the bound on the number of runs) to calculate the total in the last row of Table 2. Similarly, for the columns  $E_{Pex}$  and  $I_{Pex}$ , we take the number of runs as 10000 if conditions E and I, respectively, are not satisfied within 10000 runs.

*Pex + eXpress* takes non-trivial average time of 700 seconds to find irrelevant branches for the class `StructParser` due to a large number of method invocations. However, considering that most of the changes cannot be covered in 1 hour, the time taken to find irrelevant branches is substantially less.

Changes between two version pairs (40-45 and 40-47) could not be covered by either *Pex* nor *Pex + eXpress*. One of the changes (between version pairs 47-50) was a refactoring. For this version pair, program state was infected but no behavioral differences were detected by either *Pex* or *Pex + eXpress*.

In summary, *Pex + eXpress* was able to detect behavioral differences for four of the version pairs that could not be detected by Pex. On average, Pex was able to find behavioral differences in 68% fewer runs (median 77%) and 35% fewer amount of time. The reduction in number of runs is substantially larger than reduction in amount of time due to non-trivial time taken by *eXpress* in finding irrelevant branches.

**Table 3.** Time and irrelevant branches

S	$T_{static}(s)$	$B_{E+I}$	$B_P$	$B_{Tot}$
replace	4.5	90	57	181
siena	4.1	34	16	185
STPG	35	16	10	272
SL	0.4	33	11	383
SP	703	49	21	447

**Seeding program exploration with existing tests.** Table 4 shows the results obtained by using the existing test suite to seed the program exploration. Column *C* shows the class name. Column *V* shows the pair of version numbers. The next four columns show the number of runs and time taken by the four techniques: Pex, Pex with seeding, Pex+eXpress, and Pex+eXpress with seeding, respectively, for finding behavioral dif-

ferences. Note that DSE runs required by our Incremental Exploration also includes the seeded test runs.

In Table 4, if all the changed blocks are not covered, we take the number of runs as 10,000 (the maximum number of runs that we ran our experiments with). For 9 of the version pairs of `structorian` (out of 16 that we used in our experiments), the existing test suite of `structorian` could not find behavioral differences. Therefore, we consider these 9 version pairs for our experiments. Pex could not find behavioral differences for 5 of the 9 version pairs in 10,000 runs. Seeding the program exploration with the existing test suite helps Pex in finding behavioral differences for 3 of these version pairs under test. Pex+eXpress could not find behavioral differences for 3 of the 9 version pairs in 10,000 runs. Seeding the program exploration with the existing test suite helps Pex+eXpress in finding behavioral differences for 2 of these version pairs under test.

In summary, Pex requires around 67.5% of the original runs and 67% less time (required by Pex without test seeding) and Pex+eXpress requires around 74% of the original runs and 70% less time (required by Pex+eXpress without test seeding). In terms of time, Pex with seeding marginally wins over Pex+eXpress with seeding due to time taken by Pex+eXpress in finding irrelevant branches.

**Table 4.** Results obtained by seeding existing test suite for `structorian`

C	V	$N_{Pex}/T$	$N_{Pseed}/T$	$N_{eXpress}/T$	$N_{eSeed}/T$
SP	2-5	10000/1hr*	10000/1hr*	2381/35min	181/17min
SP	37-39	3699/26m	60/1m	851/22m	47/11m
SP	39-40	10000/1hr*	304/2m	10000/1hr*	251/12m
SP	45-47	10000/1hr*	10000/1hr*	10000/1hr*	10000/1hr*
SP	47-50	10000/1hr*	81/1m	10000/1hr*	64/10m
SP	62-124	10000/1hr*	59/1m	7228/58m	41/10m
SL	169-174	478/1m	324/1m	34/1m	18/1m
SL	150-169	299/1m	37/1m	52/1m	29/1m
SL	9-139	2988/2m	69/1m	1002/1m	52/1m
Total		64476/6.5hr	20934/2hr8m	41568/5hr9m	10683/2hr3m

\*If behavior differences are not detected, we take the number of runs as 10,000 (the maximum number of runs that we ran our experiments with)

## 6 Related Work

Previous approaches [10, 29, 15] generate regression unit tests achieving high structural coverage on both versions of the class under test. However, these approaches explore all the irrelevant paths, which cannot help in achieving any of the conditions I or E in the PIE model [33]. In contrast, we have developed a new search strategy for DSE to avoid exploring these irrelevant paths.

Santelices et al. [2, 27] use data and control dependence information along with state information gathered through symbolic execution, and provide guidelines for testers to augment an existing regression test suite. Unlike our approach, their approach does not automatically generate tests but provides guidelines for testers to augment an existing test suite.

Some existing search strategies [4, 36] guide DSE to efficiently achieve high structural coverage in a program under test. However, these techniques do not specifically

target covering a changed region. In contrast, our approach guides DSE to avoid exploring paths that cannot help in satisfying any of the conditions P, I, or E of the PIE model.

Differential symbolic execution [22] determines behavioral differences between two versions of a method (or a program) by comparing their symbolic summaries [11]. Summaries can be computed only for methods amenable to symbolic execution. However, summaries cannot be computed for methods whose behavior is defined in external libraries not amenable to symbolic execution. Our approach still works in practice when these external library methods are present since our approach does not require summaries. In addition, both approaches can be combined using demand-driven-computed summaries [1], which we plan to investigate in future work.

Joshi et al. [16] use the path constraints of the paths followed by the tests in an existing test suite to generate inputs that violate the assertions in the test suite. The generated test inputs follow the same paths already covered by the existing test suite and do not explore any new paths. In contrast, our approach exploits the existing test suite to explore new paths. Majumdar and Sen [20] propose the concept of hybrid concolic testing. Hybrid concolic testing seeds the program with random inputs so that the program exploration does not get stuck at a particular program location. In contrast, our approach exploits the existing test suite to seed the program exploration. Since the existing test suite is expected to achieve a higher structural coverage, the existing test suite is expected to discover more hard-to-discover branching nodes in comparison with random inputs. Godefroid et al. [13] propose a DSE based approach for fuzz testing of large applications. Their approach uses a single seed for program exploration. In contrast, our approach seeds multiple tests to program exploration. Seeding multiple tests can help program exploration in covering the changes more efficiently as discussed in Section 4.5.

Xu and Rothermel [38] propose a directed test generation technique that uses the existing test suite to cover parts of the program that are not covered by the existing test suite. In particular, the approach first collects the set of branches  $B$  that are not covered by the existing test suite. To cover a branch  $b_i = \langle v_i, v_j \rangle \in B$ , the approach selects all the tests  $T$  that cover the vertex  $v_i$ . For each test  $t_i \in T$ , the approach collects the path constraints  $p_i$  of path followed by  $t_i$  until  $v_i$ , negates the predicate at  $v_i$  from  $p_i$  to get path condition  $p'_i$ . The approach then generates a test that covers the branch  $b_i$  by solving the path conditions  $p'_i$ . However if all the path conditions of paths followed by the tests  $T$  are not solvable, the approach cannot generate a test to cover the branch  $b_i$ , which can furthermore compromise the coverage of additional branches. In contrast, our incremental exploration technique can still generate a test to cover such branches. In addition, the approach focuses only on satisfying condition E of the PIE model, while our approach helps in satisfying E, I, and P of the PIE model to find behavioral differences.

## 7 Discussion

In this section, we discuss some of issues of the current implementation of our approach and how they can be addressed.

**Added/Deleted and Refactored Methods.** If a method  $M$  (or a field  $F$ ) is added or deleted from the original program version, `eXpress` does not detect  $M$  (or  $F$ ) as a changed region. The change is detected if a method call site (or reference to  $F$ ) is added or deleted from the original program version. If the added or deleted method (or field) is never invoked, the behavior of the two versions is the same unless  $M$  is an overriding method. We plan to incorporate support for handling such overriding methods that are added or deleted. Similarly, if a method  $M$  is refactored between the two versions, `eXpress` does not detect  $M$  as a changed region. However, when a method is refactored, its call sites are changed accordingly (unless the method undergoes `Pull Up` or `Push Down` refactoring). Hence, `eXpress` detects the method containing call sites of  $M$  as changed. In our experiments, we considered versions of `replace` in which a method signature was changed, and versions of `structorian` in which a method was renamed.

**Granularity of Changed Region.** In our current implementation, a changed region is the list of continuous instructions that include all the changed instructions in a method. One method can have only a single changed region. Hence, a changed region can be as big as a method and as small as a single instruction. The granularity of a changed region can be increased to a single method or reduced to single instruction. Changing the granularity to single method  $M$  can affect the efficiency of our approach in reducing DSE runs since some of the branches in  $M$  that should be considered irrelevant would not be considered irrelevant. In contrast, reducing the granularity to a single instruction makes our approach more efficient in reducing DSE runs. However, the overhead cost of our approach is increased due to state checking at multiple points in the program. In future work, we plan to enhance `eXpress` to allow users to choose from different levels of granularity.

**Pruning of Branches for Propagation.** In future work, we plan to prune more categories of branches that cannot help in Propagation (P). Consider that a changed region is executed and the program state is infected after the execution of the changed region; however, the infection is not propagated to any observable output. Let  $\chi$  be the last location in the execution path such that the program state is infected before the execution of  $\chi$  but not infected after its execution.  $\chi$  can be determined by comparing the value spectra [35] obtained by executing the test on both versions of the program. This category contains all the branching nodes after the execution of  $\chi$ . These branches can be obtained by inspecting the path  $P$  followed in the previous DSE run. Let  $P = \langle b_1, b_2, \dots, b_\chi \dots b_n \rangle$ , where  $b_i$  is the a branching node in Path  $P$ , while  $b_\chi$  is the last branching node containing  $\chi$ . We flip the branching nodes starting from  $b_\chi$  to  $b_n$  in  $P$  until a branching node where the infected program state is propagated.

## 8 Conclusion

Regression testing aims at generating tests that detect behavioral differences between two versions of a program. To expose behavioral differences, a test execution needs to satisfy the conditions: Execution (E), Infection (I), and Propagation (P), as stated in the PIE model [33]. Dynamic symbolic execution (DSE) can be used to generate tests for satisfying these conditions. DSE explores paths in the program to achieve high

structural coverage, and exploration of all these paths can often be expensive. However, many of these paths in the program cannot help in satisfying the three conditions in any way. In this paper, we presented an approach and its implementation called `express` for regression test generation using DSE. `express` prunes paths or branches that cannot help in detecting the E, I, or P condition such that these conditions are more likely to be satisfied earlier in path exploration. In addition, our approach can exploit the existing test suite for the original version to efficiently execute the changed regions (if not already covered by the test suite). Experimental results on various versions of programs showed that our approach can efficiently satisfy E, I, or P conditions than without using our approach.

## References

1. S. Anand, P. Godefroid, and N. Tillmann. Demand-driven compositional symbolic execution. In *Proc. TACAS*, pages 367–381, 2008.
2. T. Apiwattanapong, R. Santelices, P. K. Chittimalli, A. Orso, and M. J. Harrold. Matrix: Maintenance-oriented testing requirement identifier and examiner. In *Proc. TAICPART*, pages 137–146, 2006.
3. C. Boyapati, S. Khurshid, and D. Marinov. Korat: Automated testing based on java predicates. In *Proc. ISSTA*, pages 123–133, 2002.
4. J. Burnim and K. Sen. Heuristics for scalable dynamic test generation. In *Proc. ASE*, pages 443–446, 2008.
5. C. Cadar, V. Ganesh, P. M. Pawlowski, D. L. Dill, and D. R. Engler. Exe: automatically generating inputs of death. In *Proc. CCS*, pages 322–335, 2006.
6. L. Clarke. A System to Generate Test Data and Symbolically Execute Programs. *IEEE Trans. Softw. Eng.*, 2(3):215–222, 1976.
7. D. Dig, C. Comertoglu, D. Marinov, and R. Johnson. Automatic detection of refactorings in evolving components. In *Proc. ECOOP*, pages 404–428, 2006.
8. H. Do, S. G. Elbaum, and G. Rothermel. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *ESE*, pages 405–435, 2005.
9. S. Elbaum, H. N. Chin, M. B. Dwyer, and M. Jorde. Carving and replaying differential unit test cases from system test cases. *TSE*, 35(1):29–45, 2009.
10. R. B. Evans and A. Savoia. Differential testing: a new approach to change detection. In *Proc. FSE*, pages 549–552, 2007.
11. P. Godefroid. Compositional dynamic test generation. In *Proc. POPL*, pages 47–54, 2007.
12. P. Godefroid, N. Klarlund, and K. Sen. DART: directed automated random testing. *Proc. PLDI*, pages 213–223, 2005.
13. P. Godefroid, M. Y. Levin, and D. A. Molnar. In *Proc. NDSS*, pages 151–166, 2008.
14. S. Horwitz. Tool support for improving test coverage. In *Proc. ESOP*, pages 162–177, 2002.
15. W. Jin, A. Orso, and T. Xie. Automated behavioral regression testing. In *Proc. ICST*, 2010.
16. P. Joshi, K. Sen, and M. Shlimovich. Predictive testing: Amplifying the effectiveness of software testing. In *Proc. ESEC-FSE*, pages 561–564, 2007.
17. J. C. King. Symbolic Execution and Program Testing. *Communications of the ACM*, 19(7):385–394, 1976.
18. J. Law and G. Rothermel. Whole program path-based dynamic impact analysis. In *Proc. ICSE*, pages 308–318, 2003.
19. J. J. Li. Prioritize code for testing to improve code coverage of complex software. In *Proc. ISSRE*, pages 75–84, 2005.
20. R. Majumdar and K. Sen. Hybrid concolic testing. In *Proc. ICSE*, pages 416–426, 2007.
21. A. Orso and B. Kennedy. Selective capture and replay of program executions. In *Proc. WODA*, pages 1–7, 2005.
22. S. Person, M. B. Dwyer, S. Elbaum, and C. S. Păsăreanu. Differential symbolic execution. In *Proc. FSE*, pages 226–237, 2008.
23. X. Ren, B. G. Ryder, M. Stoerzer, and F. Tip. Chianti: A change impact analysis tool for java programs. In *Proc. ICSE*, pages 664–665, 2005.
24. G. Rothermel and M. J. Harrold. A safe, efficient regression test selection technique. *ACM Trans. Softw. Eng. Methodol.*, 6(2):173–210, 1997.
25. D. Saff, S. Artzi, J. H. Perkins, and M. D. Ernst. Automatic test factoring for Java. In *Proc. ASE*, pages 114–123, 2005.
26. R. A. Santelices, P. Chittimalli, T. Apiwattanapong, A. Orso, and M. J. Harrold. Test-suite augmentation for evolving software. In *Proc. ASE*, pages 218–227, 2008.
27. R. A. Santelices, P. Chittimalli, T. Apiwattanapong, A. Orso, and M. J. Harrold. Test-suite augmentation for evolving software. In *Proc. ASE*, pages 218–227, 2008.
28. K. Sen, D. Marinov, and G. Agha. CUTE: a concolic unit testing engine for C. In *Proc. FSE*, pages 263–272, 2005.
29. K. Taneja and T. Xie. DiffGen: Automated regression unit-test generation. In *Proc. ASE*, pages 407–410, 2008.
30. K. Taneja, T. Xie, N. Tillmann, J. de Halleux, and W. Schulte. Guided path exploration for regression test generation. In *Proc. ICSE, NIER*, May 2009.
31. N. Tillmann and J. de Halleux. Pex-white box test generation for .NET. In *Proc. TAP*, pages 134–153, 2008.
32. N. Tillmann and W. Schulte. Parameterized unit tests. In *Proc. ESEC/FSE*, pages 253–262, 2005.
33. J. Voas. PIE: A dynamic failure-based technique. *TSE*, 18(8):717–727, 1992.

34. T. Xie, D. Marinov, and D. Notkin. Rostra: A framework for detecting redundant object-oriented unit tests. In *Proc. ASE*, pages 196–205, 2004.
35. T. Xie and D. Notkin. Checking inside the black box: Regression testing by comparing value spectra. *TSE*, 31(10):869–883, 2005.
36. T. Xie, N. Tillmann, P. de Halleux, and W. Schulte. Fitness-guided path exploration in dynamic symbolic execution. In *Proc. DSN*, 2009.
37. B. Xin, W. N. Sumner, and X. Zhang. Efficient program execution indexing. In *Proc. PLDI*, pages 238–248, 2008.
38. Z. Xu and G. Rothermel. Directed test suite augmentation. In *APSEC*, pages 406–413, 2009.