

Evolving Exact Decompilation

Eric Schulte, Jason Ruchti, Matt Noonan, David Ciarletta, Alexey Loginov

GrammaTech

531 Esty St.

Ithaca, NY 14850

{eschulte, jruchti, mnoonan, dciarletta, alexey}@grammatech.com

Abstract—We introduce a novel technique for C decompilation that provides the correctness guarantees and readability properties essential for accurate and efficient binary analysis. Given a binary executable, an evolutionary search seeks a combination of source code excerpts from a “big code” database that can be recompiled to an executable that is *byte-equivalent* to the original binary. Byte-equivalence ensures that a successful decompilation fully reproduces the behavior, both intended and unintended, of the original binary. Moreover, the decompiled source is typically more readable than source obtained with existing decompilers, since it is generated from human-written source code excerpts. We present experimental results demonstrating the promise of this novel, general, and powerful approach to decompilation.

I. INTRODUCTION

Recent widespread incidents have highlighted the significant threat malware poses to a secure internet. In a recent example, the WannaCry malware infected more than 230,000 computers in over 150 countries [32]. Reverse-engineering of WannaCry binaries was key to stopping the malware’s spread [19]. In general, reverse-engineering and analysis of malware binaries is a critical first step in incident mitigation and recovery. Accurate decompilation allows the use of source-based program analysis tools, and enables security analysts to do their work more efficiently by reasoning at the source code level.

Traditional decompilation is fundamentally limited by a reliance on deterministic techniques that are often compiler-specific, optimization-specific, and Instruction Set Architecture (ISA)-specific; as a result decompilers are typically time-intensive to write and to evaluate for correctness. Compilers often generate obscure, idiosyncratic code that leverages deep semantic properties of the source language and the target ISA. The reasoned reversal of such compiler idioms is more difficult than their initial implementation (and potentially undecidable),

This research was developed with funding from the Defense Advanced Research Projects Agency (DARPA). The views, opinions, and/or findings contained in this material are those of the authors and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government.

DISTRIBUTION A. Approved for public release; distribution unlimited.

Network and Distributed Systems Security (NDSS) Symposium 2018
18-21 February 2018, San Diego, CA, USA
ISBN 1-1891562-49-5
<http://dx.doi.org/10.14722/ndss.2018.23xxx>
www.ndss-symposium.org

a problem exacerbated by the relatively small amount of development effort dedicated to decompilation technology compared to compiler development. To reverse this imbalance, we use evolutionary search to leverage compilers themselves for decompilation across compiler, ISA, and language development.

Security analysis of binary applications requires trusted, accurate, and readable decompilation. Existing decompilers frequently produce decompilations that are not fully functionally equivalent, and fail to identify which portions of the decompiled source are accurate. Although recent academic work has begun to target correctness [29], [36], authors report that their decompilations often fail the weak semantic-equivalence tests provided by the program’s test suite.

We propose *byte-equivalence of recompilation* as a deterministically checkable guarantee of *full semantic reproduction* of the original binary, including both desired behavior as well as faults and vulnerabilities. Our technique performs an evolutionary search, using a fitness function that measures byte-similarity with the original binary executable. This search leverages the compiler as a black box: existing compiler transformations do not need to be analyzed or reversed because they are applied directly when candidate decompilations are recompiled. This process iteratively improves a large population of candidate decompilations by applying source-to-source transformations that draw from a large database of human-written source code *excerpts*. The candidate population may be initialized using randomly generated code, excerpts from this large external code database, or even using the output of other decompilers. Our technique then *improves* and *re-combines* the initial candidates, driving them closer to byte-equivalence while incorporating new human-written code excerpts to improve fidelity and readability.

The resulting Byte-Equivalent Decompilation (BED) tool combines insights from a diverse set of recent advances in software engineering. Internet-scale “big code” databases are increasingly used to drive program analysis and synthesis [3], [23], [17]. Search based techniques are now an accepted tool for common software engineering and analysis tasks including program repair [15], [34], [16], program optimization [13], [27], and high performance fuzz testing [37], [38]. Compilers have been used for partial decompilation, to automatically generate assembly to compiler-IR translators [10]. Finally, essential to BED’s success, a study of Gabel and Su [9] covering 430 million lines of source code revealed:

a general lack of uniqueness in software at levels of granularity equivalent to approximately one to seven

```

#include <stdio.h>
#define MAX 1000000
int main(void)
{
    int pprev = 0,
        prev = 1;
    int num = 0;
    int tot = 0;

    while (tot < MAX) {
        num = prev + pprev;
        prev = pprev;
        pprev = num;
        if (num % 2 == 0) {
            tot = tot + num;
        }
    }
    printf("%d\n", tot);

    return 0;
}

```

(a) Original source code.

```

int main(void)
{
    int x = 0, y = 1;
    long int sum1 = 0, sum2 = 0;
    while (sum2 < 1000000) {
        sum1 = y + x;
        y = x;
        x = sum1;
        if (sum1 % 2 == 0) {
            sum2 = sum2 + sum1;
        }
    }
    printf("%d\n", sum2);
    return 0;
}

```

(b) Byte-equivalent BED output.

Fig. 1. Problem 2 from the Project Euler [1] set, shown in original source (Figure 1a), and BED evolved byte-equivalent decompilation (Figure 1b). The compiler used automatically adds `#include <stdio.h>` when required.

lines of source code [...] crossing both project and programming language boundaries.

We evaluate BED against a benchmark of small C programs exercising specific language constructs. We demonstrate the viability of the technique by achieving byte-equivalent decompilation of most benchmark programs. The BED technique can be parallelized across functions, suggesting that our results can extend to larger programs.

We compare BED to the industry leading HEX-RAYS Decompiler [7] using a number of correctness and readability metrics. We find BED often outperforms HEX-RAYS on such criteria, even when full byte-equivalence is not attained. Even when only partial byte-equivalence is attained, BED is able to identify lines of the decompiled source and regions of the original binary that *do* achieve full byte-equivalent recompilation (Figure 3). This ensures the utility of BED’s output for security analysis by providing the guarantees of full byte-equivalence for large—and clearly delineated—regions of the decompiled source.

Experimental Results. In an evaluation against a benchmark suite of small C source programs, the BED technique achieved full byte-equivalent decompilations for 10 of the 19 programs. For those optimized benchmark programs on which BED fails to achieve byte-equivalence, BED typically matches >80% of all machine-code instructions in the binary, increasing to 97% when seeded with output from the HEX-RAYS Decompiler.

We selected a suite of decompilation quality metrics (§IV-A4) with a focus on *readability* and *correctness*: properties essential to the use of decompilation to support program security analysis. We use byte-equivalence of the resulting binary as a computable, strong proxy for correctness. BED outperforms the HEX-RAYS Decompiler on 10 of the 13 metrics (§IV-B2), notably producing significantly shorter decompiled C source.

Contributions. We describe BED’s novel **evolutionary decompilation** technique, which leverages “**big code**” databases, the **original compiler**, and **byte-equivalent recompilation** as

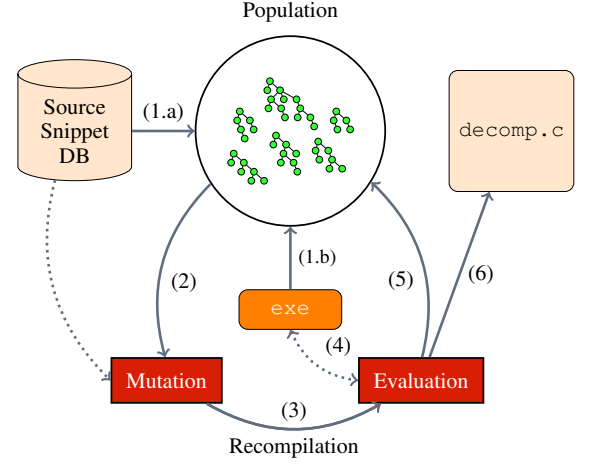


Fig. 2. BED system architecture. The components of the system are described in §II.

an **objective function** to achieve high-quality decompilation. This approach generalizes across source languages, ISAs, compilers, and optimizations.

Evolutionary search is accelerated with techniques that **mine useful excerpts from a large source code database**. This excerpt search is guided by similarity between the unmatched bytes from the original binary and the bytes associated with each source excerpt in the database. Methods of source **recontextualization** are used to incorporate foreign code into candidate decompilations.

We also apply **targeted improvement methods** to decompilation candidates, such as injecting literals mined from the target binary, and a technique for automatically responding to compiler errors with additional edits to the candidate.

II. EXAMPLE

We walk through the BED technique for the small example shown in Figure 1a, using the BED system architecture diagram Figure 2 for reference.

The first step is to seed the candidate decompilation population. Candidates may be constructed from whole function excerpts from the source database that compile to bytes similar to those in the target binary (we call such candidates “*frankensteins*”), Figure 2 (1.a) or from the outputs of one or more existing decompilers (1.b). In the run producing this decompilation, the best candidate in the initial population consisted of a single function pulled from our database, which had approximately 55% byte-similarity with the target binary.

After a population of many such *frankenstein* candidate decompilations has been generated, BED evolves the population by iteratively applying both random and targeted source-to-source transformations that leverage the code database (2), recompiling the resulting candidates to produce new binaries (3), and then evaluating the fitness of the resulting binary using a byte-similarity fitness function (4). The evaluated candidates are returned to the population (5), and the technique proceeds using an evolutionary algorithm to guide the search towards increasing byte-similarity with the target in a process mimicking natural selection. When full byte-equivalence against the target

```

#include <stdio.h>
int main(int argc, char *argv[])
{
    int i = 0;
    // go through each string in argv
    // why am I skipping argv[0]?
    for(i = 1; i < argc; i++) {
        printf("arg %d: %s\n", i, argv[i]);
    }
    // let's make our own array of strings
    char *states[] = {
        "California", "Oregon",
        "Washington", "Texas"
    };
    int num_states = 4;
    for(i = 0; i < num_states; i++) {
        printf("state %d: %s\n", i, states[i]);
    }
    return 0;
}

```

(a) Original source code.

```

int main(int argc, char** argv) {
    for (int d = 1; argc * 1 > d; d++) {
        printf("arg %d: %s\n", d, *argv);
    };
    int d;
    d = 0;
    printf("state %d: %s\n", d, "California");
    d = 2;
    printf("state %d: %s\n", 1, "Oregon");
    printf("state %d: %s\n", d, "Washington");
    printf("state %d: %s\n", 3, "Texas");
}

```

(b) Nearly byte-equivalent BED output.

Fig. 3. The tenth problem from the Learn C The Hard Way [30] programming tutorial (Figure 3a), with BED-evolved decompilation (Figure 3b). Byte equivalence is achieved for > 98% of the program, only failing at the highlighted line.

is achieved the search terminates (6) returning the resulting byte-equivalent decompilation.

We present a number of techniques to improve the efficiency and accuracy of this search process. The application of source-to-source transformations is weighted to target those regions of candidate decompilations that currently fail to match the original. After each mutation—if the resulting individual fails to compile—BED may make targeted changes to achieve compilation. This is particularly useful when attempting to incorporate output from existing decompilers, which often fail to recompile. Once the individual has been compiled and the output compared against the target binary, BED may also make targeted edits to inject literals mined from the target binary. For example, the string literal on line 13 of Figure 1b and the % 2 operation on line 9 of Figure 1b were both pulled directly from the target binary after binary-to-binary comparison.

At each generation a subset of individuals are selected for reproduction using lexicase selection (§ III-E); this maintains diversity by ensuring candidates recompiling to match regions of the original binary not matched by other candidates are not lost from the population.

The results of this technique are shown in Figure 1 along with the original source. The HEX-RAYS decompilation of the same function demonstrates artifacts such as inclusion of C runtime functions and comments referencing stack

offsets. By contrast the BED decompilation in Figure 1b has similar size and form to the original, in Figure 1a, and could plausibly have been written by a human. Although the HEX-RAYS decompilation of this function fails to compile without modification, the BED decompilation not only compiles, but matches the original program byte for byte.

III. OVERVIEW

The BED algorithm is shown in Figure 4. The computationally expensive tasks of evaluation, line 7, and reproduction, line 28, may both be parallelized across multiple threads.

The population may optionally be seeded, line 1, with the output of existing decompilers (§ IV-A5) via the *InitialPop*. Additional candidate decompilations may be synthesized from the code database, line 3, using the *Frankenstein()* method (§ III-A1). The algorithm’s main loop, lines 5–41, iteratively evaluates, selects, and regenerates the population driving the candidate decompilations towards increasing byte-similarity with the target binary, *TBin*. During evaluation candidate decompilations are first compiled, line 8, and failed compilations are fixed when possible, line 10, (§ III-B2). Candidate compilations are then compared against *TBin*, line 17 using a disassembly-based fitness function (§ III-D) that returns the diff between *TBin* and *CBin*. This diff is later used to drive selection and to target mutation.

Selection is performed using lexicase selection, line 35, (§ III-E); crossover is applied with chance, probability *CrossRate*, line 30. Mutation is then applied yielding a new candidate decompilation.

The generation counter is then incremented, line 40, and the algorithm repeats with *Pop* now holding the newly generated candidate decompilations. The algorithm runs until either byte-equivalence is achieved, line 21, or the runtime budget of *MaxEvals* is exhausted, line 41. In either case the resulting candidate decompilation is minimized to remove any elements of the evolved source that do not contribute to the compilation (§ III-F).

The remainder of this section discusses each of these stages in greater detail.

A. Creating an Initial Population

The initial population of candidate decompilations for BED can be formed in many ways. We chose to focus on two in particular:

1) *Franksteins*: Franksteins are composed of whole functions drawn from the code database. Function signatures and boundaries are extracted from the original binary using CodeSurfer® for Binaries, GrammaTech’s static analysis platform for analyzing stripped binaries, producing a template C file consisting of function prototypes with empty bodies. The bytes for each identified function are then used to perform a similarity search in the code database as described in § III-C5. Multiple frankenstein candidate decompilations are then generated by replacing each empty body in the template C code with random selections from the search results, using a geometric distribution to bias selections towards the top search results.

Input: Target Binary, TBin : *Executable*
Input: Snippet Database, DB : *Database*
Parameters: PopSize, CrossRate, MaxEvals
Parameters: InitialPop, NumFranks, FixLitChance
Output: Decompilation of TBin

```

1: let Pop ← InitialPop, EvalCounter ← 0
2: do NumFranks times
3:   let Pop ← Pop ∪ Frankenstein(DB, TBin)
4: end do
5: loop
6:   let TmpPop ← ∅
7:   for Candidate ∈ Pop do
8:     let CBin ← Compile(Candidate)
9:     if CBin = Failed then
10:      Candidate ← FixCompilation(Candidate)
11:      CBin ← Compile(Candidate)
12:    end if
13:    let Diff ← null
14:    if CBin = Failed then
15:      Diff ← Disassembly(TBin)
16:    else
17:      Diff ← Evaluate(TBin, CBin)
18:    end if
19:    EvalCounter ← EvalCounter + 1
20:    if Fit = 0 then
21:      return Minimize(Candidate)
22:    else if EvalCounter ≥ MaxEvals then
23:      return Minimize(Best(Pop))
24:    end if
25:    TmpPop ← {Candidate, Diff}
26:  end for
27:  Pop ← ∅
28:  while |Pop| < PopSize do
29:    let p ← (null, null, null), p' ← null
30:    if Random() < CrossRate then
31:      let p1 ← Select(TmpPop)
32:      let p2 ← Select(TmpPop)
33:      p ← Crossover(p1, p2)
34:    else
35:      p ← Select(TmpPop)
36:    end if
37:    let p' ← Mutate(DB, p)
38:    Pop ← Pop ∪ {p'}
39:  end while
40:  GenCounter ← GenCounter + 1
41: end loop

```

Fig. 4. BED algorithm.

2) *Other Decompilers:* Any available decompiler may also be used to seed the initial population of candidate decompilations, jump-starting BED’s evolutionary search. Candidates from multiple decompilers may be used in combination with frankensteins to seed a single population. We evaluate the impact of seeding the population with decompilation generated via the HEX-RAYS Decompiler.

B. Local Search Techniques

To improve the speed with which the BED technique converges on byte-equivalent decompilations we introduce two forms of local search, literal mining §III-B1 and compilation repair §III-B2. Each of these techniques may be used to directly improve the fitness of randomly evolved candidate decompilations. Evolutionary search techniques augmented with local search are called *memetic algorithms* and have been proven

effective in many cases [18], [8].

1) *Fixing Source Literals:* Because no code database can contain all possible literals to be utilized in a program, we have developed a mutation to identify literals in the target binary and insert them directly into the relevant portion of the candidate source. Many of the literal values immediately available in program binaries, such as addresses, are not suitable candidates for source-level transformation. Also, the compiler may optimize operations such as multiplication and division in such a way that a literal value cannot be gleaned directly from the disassembly.

To address the first challenge, we identify instructions that contain references to program values. When an operand references the read-only data section of the executable, we attempt to find the literal value at the referenced location. For instance, in “`lea eax, 0x80484e0`”, if “`0x80484e0`” is the starting address of the string literal “`Hello, world!`”, we offer this as a candidate replacement in the source that compiles to the relevant portion of the candidate recompilation. Similarly when we encounter an instruction such as `movsd xmm0, 0x8048638`, we offer the floating point literal at `0x8048638` as a replacement.

To address the challenge of compiler optimization of common arithmetic expressions, we mine constants from logical representations of binary fragments. As an example, consider the source fragment $n = 12 * k$; this may be represented at the machine level as the a sequence of instructions equivalent to $n = k$; $n = (n + n + k) \ll 2$; . After identifying instruction sequences in the target binary that match the general pattern of optimized multiplication or division, we generate a formula describing the sequence’s net effect on the machine state. The resulting formula can be mined for constant values; in the example above, we would find that k was scaled by 12.

Our targeted fixing of source literals works by annotating the target decompilation with mined source literals. During fitness evaluation we identify the regions where the compiled candidate differs from the target binary, and use debugging information from the compilation to map those back to the candidate source. The fixer can use this map to inject any literals mined from a binary region back into the corresponding source region.

2) *Fixing non-compiling source code:* Although we make an effort to perform mutations that generate syntactically valid individuals, it is possible that a newly-generated individual will fail to compile. In this case, we apply a sequence of *compilation-fixing transformations* by matching the compiler’s error messages to known strategies for fixing compilation. For example, a program that references a library function without including an appropriate header file can cause an error of the form “implicitly declaring library function F ”; on seeing such an error, we automatically check if a man page exists for F and, if so, extract any `#include` directives from that page.

Another common compilation fixer is applied when an error of the form “use of undeclared identifier X ” appears. In this case, we insert a random declaration for the given identifier on a line just before the error occurred, supplying a randomly-chosen type for the new declaration. This change is similar to the `Add Init` repair type in SPR [16].

```

int fib(int n) {
  int x = 0;
  int y = 1;
  while (n > 0) {
    int t = x;
    x = x + y;
    y = t;
  }
  return x;
}

int collatz(int m) {
  while (m != 1) {
    if (m % 2 == 0)
      m /= 2;
    else
      m = 3*m + 1;
    ++k;
  }
  printf("%d\n", k);
  return k;
}

```

Fig. 5. Crossover example. If we perform crossover by trading lines 6–10 of `fib` with lines 3–11 of `collatz`, we must ensure that `k` is not re-bound to `t`; otherwise, the `printf` statement would reference an out-of-scope variable.

Our suite of compilation-fixing strategies includes techniques designed to counter known problems in contributing decompilers. This significantly increases the utility of decompiler-generated candidates to the BED evolutionary search.

If none of our compilation-fixing strategies are successful, we simply delete the line on which the error occurred. The compilation-fixing strategies are iterated until either an error-free program is obtained, or until a maximum number of attempts has been made.

C. Evolutionary Search

1) *Recontextualization*: We perform mutations directly on the source text of our individuals, necessitating some cleanup when regions of source code are injected into new contexts. For example, imagine that we were to implement an “insert” mutation, using the code excerpt `x = x + 1`. At the chosen insertion point, it is very unlikely that a variable named `x` is already in scope; a naïve insertion would result in an individual that would not compile. To increase the changes that mutation will result in compilable source, we *recontextualize* the code by identifying all free identifiers and rebinding them to names that are in scope at the insertion point.

When the mutation text involves a group of statements that span more than one scope, extra caution is required to find a valid rebinding of identifiers which respects all scopes, as illustrated by the crossover example in Figure 5.

2) *Mutations*: When modifying a candidate decompilation, we randomly apply either a basic mutation or a targeted mutation from Table I. The basic mutations are standard “cut” (remove an expression), “insert” (insert text), “swap”, and “replace.” The targeted mutations are specialized mutations designed specifically for decompilation. Tuning the set of targeted mutations lets us inject domain-specific knowledge into the evolutionary search. Each basic kind of mutation can be modified by the kind of program region where it should be applied, or the source of program text to use. The available modifiers are:

- restrict the mutation to full statements only,¹
- restrict to source and target pairs with matching AST classes,
- use source text drawn from the code database, and

¹A “full statement” is defined as any immediate child of a block; these roughly correspond to expressions in the C grammar that end with a semicolon or are wrapped in braces.

Mutation	Description
fix-literals	Identify literals in the target binary and replace incorrect source literals.
promote-guarded	Promote statements from within compound statements, (i.e. <code>{...}</code>).
explode-for-loop	Decompose a <code>for</code> loop into a <code>while</code> loop.
coalesce-while-loop	Coalesce a <code>while</code> loop and preceding assignment into a <code>for</code> loop.
arith-assign-expansion	Expand arithmetic assignments (e.g. <code>+=</code>) into traditional assignments.
insert-excerpt-decl	Search the database specifically for a declaration to be inserted into the program.
rename-variable	Select a variable and replace with another in-scope variable.

TABLE I. TARGETED MUTATIONS. BED uses these domain-specific mutations to accelerate evolution and address problems identified empirically during evolutionary decompilation runs.

- use source text drawn from elsewhere in the decompiled program.

In all cases, the source text is recontextualized (see § III-C1).

3) *Crossover*: All crossover operations are *homologous*: the parents are aligned before crossover points are selected. We align parents by compiling each and matching their compiled machine code. A crossover point in one parent is mapped through the aligned machine code to a corresponding point in the other parent’s source. Homologous crossover has been previously shown to encourage meaningful and productive crossover operations [20]. We perform two-point crossover, selecting two aligned points. The region between these points may open or close n scopes. In the second parent, we find a homologous region with the same value of n and swap these two regions, as shown in Figure 5.

4) *Diff-targeted mutation*: We focus mutations towards the “broken” parts of the candidate decompilation. Specifically, we identify *bad statements* in a program as those source statements which compile to bytes belonging to some *diff-region*, as described in § III-D. We limit mutation targets to these bad statements with probability *TargetChance* (Table III), and otherwise choose targets uniformly at random.

5) *Similarity search*: To improve the quality of excerpts used in insertion and replacement mutations, we search our excerpt database for the excerpts whose compiled form is most byte-similar to the region of the target binary corresponding to the point of mutation. This optimization leverages the *diff-region* identified by the fitness function (see § III-D). We search the database for the excerpts most similar to the target binary’s disassembly in the *diff-region* using edit distance as our similarity metric. We then perform a weighted random selection from the sorted results, ensuring that those results from the code database which most closely match the *diff-region* are selected with a higher probability. In our experiments, the Pliny Database [40] was used for code search.

D. Measuring Fitness

To evaluate a candidate decompilation of the target binary, we use the edit distance between the target and candidate decompilations’ disassembly listings as our fitness metric.² We assign a perfect fitness to each instruction in the target binary

²Computed using a sequence comparison algorithm in Wu *et al.* [35].

within a *common-region* (i.e., non *diff-region*) and a bad fitness value proportional to the size of the diff to each instruction in the target binary within a *diff-region*.

To normalize the disassembly listings, we remove irrelevant padding instructions used for alignment purposes. We also resolve addresses as described in § III-B1. Finally, we identify instructions that change the program counter, such as `jmp 0x8014040` or `call 0x8014080` and replace the address operand with the section name and offset to ensure differences in program alignment do not negatively impact fitness.

E. Updating the population with Lexicase Selection

Lexicase selection is a genetic programming technique shown to dramatically increase population diversity and long term evolutionary success in some cases [11]. With lexicase selection, the fitness of an individual is represented as a vector of independent test cases (instead of the traditional representation as a scalar sum of passed test cases). Selection is performed by iteratively filtering the population by best performance on a random ordering of the tests in this vector until a single candidate remains. By splitting the fitness calculation across multiple test cases, individuals that pass tests not widely passed by the remainder of the population tend to be retained in the population, even when they fail many other tests. With a traditional fitness function such individuals are often lost from the population due to a low total number of test cases passed.

Each machine-code instruction in the target binary is a distinct fitness test case. With this approach, a candidate decompilation that matches a portion of the target binary unmatched by other candidates has a high probability of being retained in the population. The novel code from this candidate responsible then has many opportunities to be incorporated into other candidate decompilations via crossover.

F. Final cleanup

There are no protections against the evolutionary search accumulating unnecessary source code artifacts. In fact, through a phenomenon known as *bloat* [22], it is often beneficial to a candidate’s long-term survival to accumulate unused, or *dead*, genetic material. To compensate for bloat we minimize the final source code using *delta debugging* [39], systematically removing lines and expressions from the program until a minimal set that compiles to the same binary is found.

Before minimization we pass the source code through the `clang-format` utility [12] to enforce uniform indentation, improve readability, and to break semantically distinct program elements into separate lines for easier removal during delta debugging.

IV. EXPERIMENTAL RESULTS

Our experiments address five primary research questions:

- **Q1: Evolution of byte-equivalence from big code.** We demonstrate the feasibility of evolving byte-equivalent code by evaluating the degree to which BED is able to achieve byte-equivalence against a benchmark set. (§ IV-B1).

Program	LOC	Description
LCTHW [30]		
hw1	5	Hello world
hw3	9	Formatted printing
hw6	19	Types of variables
hw7	23	More variables, some math
hw8	32	Sizes and arrays
hw9	40	Arrays and strings
hw10	20	Arrays of strings, looping
hw11	22	While loop & boolean exprs.
hw12	17	If else if else
hw13	45	Switch statement
hw14	32	Writing and using functions
hw15	43	Pointers dreaded pointers
hw18	88	Pointers to functions
Euler [1]		
Euler-1	12	Multiples of 3 and 5
Euler-2	19	Even Fibonacci numbers
Euler-3	48	Largest prime factor
Euler-4	40	Largest palindrome product
Euler-5	15	Smallest multiple
Euler-6	48	Sum square difference

TABLE II. BENCHMARK PROGRAMS. Learn C The Hard Way (LCTHW) is a C example-based tutorial that provides simple programs illustrating features of the C language [30]. The Project Euler problems are a selection of answers to the Project Euler programming challenge [1].

- **Q2: Utility of evolved decompilation.** We evaluate the utility of evolved decompilation to support the analysis and rewriting of binary executables. This evaluation is performed through analysis of a number of *readability* and *correctness* metrics of both the *decompiled source* and the *recompiled binary* (§ IV-A4).
- **Q3: Performance on optimized binaries.** We evaluate BED across multiple levels of compiler optimization.
- **Q4: Impact of targeted mutations.** We evaluate the utility of our targeted mutation operators for the evolution of byte-equivalence (§ IV-B4).
- **Q5: Impact of lexicase selection.** We evaluate the impact of lexicase selection versus standard tournament selection.

A. Methodology

Experimental evaluations were performed against the benchmarks in § IV-A1 using the parameters described in § IV-A2, and the Euler database described in § IV-A3.

1) *Benchmarks:* The benchmark programs in Table II were selected to demonstrate a wide range of features of the C language. Specifically, the Learn C The Hard Way (LCTHW) programs [30] are taken from an example-based C language tutorial. Each selection demonstrates a specific facet of the language. The Project Euler programs [1] are taken from a popular set of programming challenges. As “big code” projects such as DARPA’s MUSE (<http://corpus.museprogram.org/>) continue to mature we anticipate the availability of increasingly relevant sample code.

2) *Parameters:* The experimental parameters used during our BED runs are listed in Table III. Runs were performed at both `-O0` and `-O2` optimization levels.

Parameter	Value
compiler	clang
flags	-m32 -g -Ox
CrossRate	$\frac{1}{4}$
MaxEvals	131072
PopSize	1024
TargetChance	$\frac{3}{4}$

TABLE III. BED EXPERIMENTAL PARAMETERS.

3) *Code Database*: The Project Euler [1] database consists of over 80k unique source ASTs extracted from 688 posted solutions to Project Euler problems.³

4) *Metrics*: To evaluate BED’s performance, we consider two evaluation criteria: readability and byte similarity to the target binary. To compute the byte similarity, we take the number of matched instructions over the total number of instructions within the original and evolved binary. In addition to measuring byte similarity to the target binary, we apply the readability metrics listed in Table V to the decompiled source code. The first eleven readability metrics all count source features that inhibit readability, such as the number of `goto` statements, and thus a smaller number is preferable. The last readability metric measures direct matches against the original source, and thus a larger number is preferable.

5) *Decompilers*: We perform an experimental comparison against the HEX-RAYS Decompiler circa 2015. An attempted comparison against recently published decompilers was not possible; we were unable to obtain executables or input/output samples for the published results in Dream [36] and Phoenix [29].

B. Experimental results

1) *Evolution of Byte-equivalence*: By combining solutions found both with and without compiler optimizations enabled, BED is able to achieve byte-equivalence for 4 of 19 example programs when run without decompiler-generated seeds and 10 of 19 when the HEX-RAYS Decompiler is used to seed candidate populations. Detailed results are shown in Table IV. In those cases where byte-equivalence is not achieved, the resulting decompilation is *close to* the original program, often achieving byte-equivalence at the function granularity.

In non-byte-equivalent runs, the BED technique is able to highlight those portions of the source which are compiling to non-byte-equivalent portions of the resulting binary. In practice non-byte-equivalent portions are small, as we achieve 81.23% byte-similarity with the target binary in our experiments. In terms of LoC, for non-optimized binaries, BED fails to achieve byte-equivalent recompilation for 11.79% of all lines without decompiler seeds and 8.59% of all lines when decompiler seeds are used. Thus even runs that fail to achieve full byte equivalence still provide useful source for reverse engineering and analysis for the majority of the program. This high level of byte-similarity indicates the promise of the BED technique.

2) *Utility of evolved decompilation*: Our readability and correctness metrics, described in §IV-A4, are shown in Table V. We use byte-equivalence of the resulting binaries

Metric	BED	HEX-RAYS
Byte similarity	81.23	69.16
# <code>gotos</code>	0	10
# casts	18	118
# variables	219	262
# scopes	216	211
# live ranges	185	262
# dead assignments	34	0
# macros	8	3116
# typedefs	0	608
# characters	22,346	263,835
# lines	1040	8342
# ASTs	4723	5616
# matched literals	287	398

[†] AST count does not include macros and types

TABLE V. BYTE-SIMILARITY AND READABILITY METRICS. For each metric in the BED column, an average was computed over both optimized and unoptimized test cases, without the use of decompiler seeding.

as a computable proxy for correctness that does not depend on pre-existing high-coverage test suites, formal specifications, or other externally-provided data.

In comparison to the HEX-RAYS Decompiler, the decompilations generated by BED included fewer `gotos`, casts, variable declarations, live ranges, macros, typedefs, and significantly shorter total decompiled code—all key metrics for accessing source code readability. Additionally, BED achieved a higher degree of byte-similarity with the target binary, while also producing re-compilable source.

3) *Performance on optimized binaries*: At higher levels of optimization, the mapping between source ASTs and machine code is less direct as source is elided or transformed significantly. As a result, the targeting of mutations to diff-inducing source ASTs becomes a significant challenge.

Despite this, BED performs well at higher levels of optimization, and indeed performs better with optimizations enabled for some benchmark programs. At higher levels of optimization, the space of valid C source code that compiles to the same bytes is significantly larger than without optimizations enabled. For instance, consider the example HW10 decompilation given in 3. In the resultant decompilation, the loop in the original program has been elided and replaced with a sequence of `printf` statements that compiles, with optimization, to the same bytes as the original; this would not be the case with optimizations disabled.

4) *Impact of targeted mutations*: The success rates for various mutation types when acting on a population in an evolutionary run appear in Figure 6.

For each mutation, we assess the fitness of an individual before and after the mutation, recording whether the mutation caused fitness to improve, remain the same, or become worse. If the mutated individual could not compile, we applied the compilation-fixing strategies (§III-B2); if the fixed source still did not compile, the individual was classified as “dead.”

All mutations showed a 1–10% probability of resulting in a fitness improvement, but showed radically different chances of producing neutral candidate decompilations. Generally, mutations that preserve the AST class and only operate on full statements performed the best, creating neutral candidate

³The Project Euler corpus was compiled by Seth Pollen and Ben Welton under the supervision of Prof. Ben Liblit.

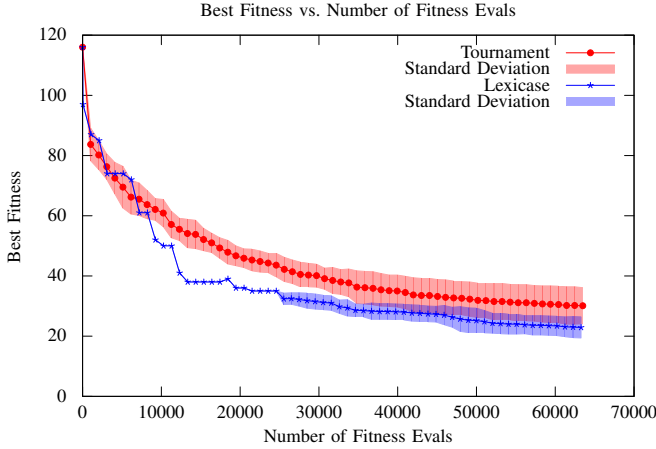


Fig. 7. Comparison of the best fitness in the population as a function of number of fitness evaluations performed using lexicase and tournament selection. Numbers shown are averages of 15 decompilations of the Euler-2 benchmark program.

VI. RELATED WORK

Phoenix. Schwartz *et al.* [29] note that previous work rarely evaluates *correctness*, meaning the preservation of behavior of the original binary. They note that decompilation for security requires both correctness and high-level abstractions.

Correctness of Phoenix’s output was assessed by comparing the recompiled binary’s behavior to the original program on a suite of tests. This approach to correctness depends on having a test suite with high coverage, which may not be possible for legacy or third-party binaries.

JSNice. Raychev *et al.* [24] introduce a novel approach for predicting identifier names and type annotations in JavaScript code. Using an extensive corpus of existing JavaScript, they build a probabilistic graphical model of program properties using conditional random fields. The resulting tool, JSNice, is able to predict correct names for 63% of identifiers and correct type annotations in 81% of cases on obfuscated code.

Mutational Robustness. Software functionality has been found to be surprisingly robust to random perturbations [28]. This robustness and the resultant *neutral spaces* of multiple diverse implementations of a single specification are thought to be central to software’s amenability to genetic improvement [26]. Our investigation of the effectiveness of mutations, perhaps unsurprisingly, identifies similar robustness of byte-similarity of compilation to source modification. Similarly our breakdown of binary impact by type of modification extends and agrees with similar work analyzing the functional impact of source-to-source transformations [2], [16].

VII. FUTURE WORK

C structure declarations. The benchmark programs, mutation strategies, and decompilation results presented herein do not include consideration of C structures. We believe generalizing our technique to support mutations of structure declarations is a matter of implementation and does not fundamentally affect the implementation or the performance of the BED

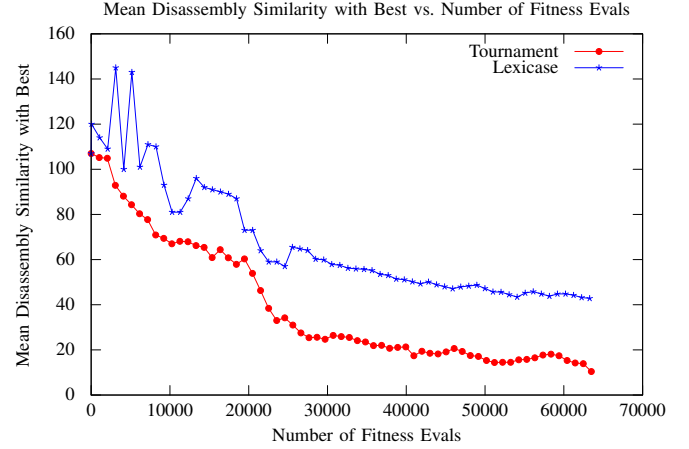


Fig. 8. Comparison of lexicase and tournament selection on diversity of Euler-2 decompilations. 15 decompilations were performed for each style. A value of zero implies no diversity.

technique. Aggregate type declarations could be co-evolved with the population.

Lifting to LLVM-IR. The BED technique is not inherently specific to C; for any language P with a black-box translator $T : P \rightarrow \text{Binary}$, BED methods give a lifting function $L : \text{Binary} \rightarrow P$, such that $T(L(x)) \approx x$ for all binaries x .

The low level virtual machine (LLVM) intermediate representation [14] has become a popular target for program analysis and rewriting. The discussion above suggests that BED could be applied to LLVM-IR instead of C to evolve liftings from machine code to LLVM-IR. We conjecture that BED may work particularly well for this task, due to the relative simplicity of LLVM-IR as compared to the C language. The application of BED to LLVM lifting has the potential to significantly out-perform the current state of the art [31], [6], which exhibit cumbersome recompilation artifacts including explicit representations of activation records on the stack.

Multi-objective fitness functions. Multi-objective fitness functions can explicitly target additional desirable properties aside from byte similarity. These might include the following: *Readability Metrics.* Readability metrics (e.g. from §IV-A4) could explicitly be added to the BED fitness function. *Test Suites.* It is possible that leveraging program test suites to guide the search could improve the overall efficiency. Functional correctness may serve as a useful proxy for byte equivalence, accelerating BED’s evolutionary search.

Machine learning for local decompilation. The large code corpora used by BED could also be used to train local machine-learning techniques to automatically *learn* a decompilation function mapping short sequences of compiled machine code to C source code decompilation. For example, recurrent neural network (RNN) encoder-decoder models have had great success in translating between natural languages [5]. Adapting such models to translate small program segments from machine code to C could subsequently be used to improve the evolutionary decompilation process by directly translating the diff inducing regions of the target binary to C source.

Structured code search and completion. Existing systems for structured source code search [21] and source code completion [4] could be applied to candidate decompilations. Such tools are intended to aid developers in the authoring of code by searching a large code database for excerpts which are similar or likely to complement existing code. BED could leverage such techniques by using candidate decompilations to drive code retrieval and suggestion systems to generate new code or modifications to be applied back to the candidate decompilation. Such techniques could accelerate the decompilation process and result in increasingly natural decompiled source code.

VIII. CONCLUSION

We present a novel and general Byte-Equivalent Decompilation (BED) technique that works via the evolutionary recombination and recompilation of source excerpts from a “big code” database. We present experimental results demonstrating BED’s ability to evolve readable C source code that compiles to an exact byte-for-byte match against target binaries. Byte-equivalent recompilation of source code is a deterministically checkable guarantee of *full semantic equivalence* with the original binary, including both desired behavior as well as faults and vulnerabilities. The genesis of the decompiled code in human-written excerpts leads to human-readable output.

We present a number of techniques that make byte-equivalent decompilation possible, including the use of existing decompilers to seed the search, an effective fitness function to guide the search, the use of byte-similar retrieval from a code database, the use of binary-difference-based mutation targeting, processes for fixing failing compilations and for pulling literal values from the original binary, and the application of lexibase selection. BED is a novel approach to decompilation that promises to achieve both readability and full semantic equivalence in a manner which is “future proof” against new source languages, ISAs, compilers, and compiler optimizations.

ACKNOWLEDGMENTS

The authors would like to thank David Melski for proposing the original idea behind BED.

REFERENCES

- [1] “Project Euler.” [Online]. Available: <http://projecteuler.net/>
- [2] B. Baudry, S. Allier, M. Rodriguez-Cancio, and M. Monperrus, “Automatic software diversity in the light of test suites,” in-submission-2015.
- [3] P. Bielik, V. Raychev, and M. Vechev, “Programming with ‘big code’: Lessons, techniques and applications,” in *LIPICs-Leibniz International Proceedings in Informatics*, vol. 32. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2015.
- [4] M. Bruch, M. Monperrus, and M. Mezini, “Learning from examples to improve code completion systems,” in *Proceedings of the the 7th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on The foundations of software engineering*. ACM, 2009, pp. 213–222.
- [5] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using rnn encoder-decoder for statistical machine translation,” *arXiv preprint arXiv:1406.1078*, 2014.
- [6] A. Dinaburg and A. Ruef, “Mcsema: Static translation of x86 instructions to llvm,” in *ReCon 2014 Conference, Montreal, Canada*, 2014.
- [7] C. Eagle, *The IDA Pro Book: The Unofficial Guide to the World’s Most Popular Disassembler*. No Starch Press, 2011.
- [8] E. Elbeltagi, T. Hegazy, and D. Grierson, “Comparison among five evolutionary-based optimization algorithms,” *Advanced engineering informatics*, vol. 19, no. 1, pp. 43–53, 2005.
- [9] M. Gabel and Z. Su, “A study of the uniqueness of source code,” in *Foundations of Software engineering*. ACM, 2010, pp. 147–156.
- [10] N. Hasabnis and R. Sekar, “Automatic generation of assembly to ir translators using compilers.”
- [11] T. Helmuth, L. Spector, and J. Matheson, “Solving uncompromising problems with lexibase selection,” *IEEE Transactions on Evolutionary Computation*, 2014.
- [12] D. Jasper, “clang-format: Automatic formatting for C++,” 2013, <http://llvm.org/devmtg/2013-04/jasper-slides.pdf>.
- [13] W. B. Langdon and M. Harman, “Evolving a CUDA kernel from an nVidia template,” in *Congress on Evolutionary Computation*, 2010, pp. 1–8.
- [14] C. Lattner and V. Adve, “Llvm: A compilation framework for lifelong program analysis & transformation,” in *Proceedings of the international symposium on Code generation and optimization: feedback-directed and runtime optimization*, ser. CGO ’04. Washington, DC, USA: IEEE Computer Society, 2004, pp. 75–.
- [15] C. Le Goues, T. Nguyen, S. Forrest, and W. Weimer, “Genprog: A generic method for automatic software repair,” *Software Engineering, IEEE Transactions on*, vol. 38, no. 1, pp. 54–72, 2012.
- [16] F. Long and M. Rinard, “Staged program repair with condition synthesis,” in *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. ACM, 2015, pp. 166–178.
- [17] —, “Automatic patch generation by learning correct code,” in *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 2016, pp. 298–312.
- [18] P. Merz and B. Freisleben, “A genetic local search approach to the quadratic assignment problem,” in *Proceedings of the 7th international conference on genetic algorithms*, 1997, pp. 1–1.
- [19] L. H. Newman, “How an accidental ‘kill switch’ slowed Friday’s massive ransomware attack,” <https://www.wired.com/2017/05/accidental-kill-switch-slowed-fridays-massive-ransomware-attack/>, May 2017.
- [20] P. Nordin, W. Banzhaf, and F. D. Francone, “12 efficient evolution of machine code for cisc architectures using instruction blocks and homologous crossover,” *Advances in genetic programming*, vol. 3, p. 275, 1999.
- [21] S. Paul and A. Prakash, “A framework for source code search using program patterns,” *Software Engineering, IEEE Transactions on*, vol. 20, no. 6, pp. 463–475, 1994.
- [22] R. Poli, W. B. Langdon, and N. F. McPhee, *A field guide to genetic programming*. Published via <http://lulu.com> and freely available at <http://www.gp-field-guide.org.uk>, 2008, (With contributions by John. R. Koza). [Online]. Available: <http://www.gp-field-guide.org.uk>
- [23] V. Raychev, M. Vechev, and A. Krause, “Predicting program properties from ‘big code’,” in *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, ser. POPL ’15. New York, NY, USA: ACM, 2015.
- [24] —, “Predicting program properties from ‘big code’,” in *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 2015, pp. 111–124.
- [25] N. E. Rosenblum, B. P. Miller, and X. Zhu, “Extracting compiler provenance from program binaries,” in *Proceedings of the 9th ACM SIGPLAN-SIGSOFT workshop on Program analysis for software tools and engineering (FSE 2010)*. ACM, 2010, pp. 21–28.
- [26] E. Schulte, “Neutral networks of real-world programs and their application to automated software evolution,” Ph.D. dissertation, University of New Mexico, Albuquerque, USA, July 2014, <https://cs.unm.edu/~eschulte/dissertation>.
- [27] E. Schulte, J. Dorn, S. Harding, S. Forrest, and W. Weimer, “Post-compiler software optimization for reducing energy,” in *Proceedings of the eighteenth international conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS ’14. ACM, 2014, pp. 639–652.
- [28] E. Schulte, Z. Fry, E. Fast, W. Weimer, and S. Forrest, “Software

- mutational robustness,” *Genetic Programming and Evolvable Machines*, pp. 1–32, 2013.
- [29] E. J. Schwartz, J. Lee, M. Woo, and D. Brumley, “Native x86 decompilation using semantics-preserving structural analysis and iterative control-flow structuring,” in *Proceedings of the USENIX Security Symposium*, 2013, p. 16.
- [30] Z. A. Shaw, “Learn C the hard way.” [Online]. Available: <http://c.learncodethehardway.org/book/>
- [31] M. Smithson, K. ElWazeer, K. Anand, A. Kotha, and R. Barua, “Static binary rewriting without supplemental information: Overcoming the tradeoff between coverage and correctness,” in *Reverse Engineering (WCRE), 2013 20th Working Conference on*. IEEE, 2013, pp. 52–61.
- [32] G. Swenson, “Bolstering government cybersecurity lessons learned from wannacry,” <https://www.nist.gov/speech-testimony/bolstering-government-cybersecurity-lessons-learned-wannacry>, 2017.
- [33] A. Wagner, *Robustness and Evolvability in Living Systems*. Princeton University Press, 2013.
- [34] W. Weimer, Z. P. Fry, and S. Forrest, “Leveraging program equivalence for adaptive program repair: Models and first results,” in *Automated Software Engineering (ASE), 2013 IEEE/ACM 28th International Conference on*. IEEE, 2013, pp. 356–366.
- [35] S. Wu, U. Manber, G. Myers, and W. Miller, “An o(np) sequence comparison algorithm,” *Inf. Process. Lett.*, vol. 35, no. 6, pp. 317–323, Sep. 1990. [Online]. Available: [http://dx.doi.org/10.1016/0020-0190\(90\)90035-V](http://dx.doi.org/10.1016/0020-0190(90)90035-V)
- [36] K. Yakdan, S. Eschweiler, E. Gerhards-Padilla, and M. Smith, “No more gotos: Decompilation using pattern-independent control-flow structuring and semantics-preserving transformations,” 2015.
- [37] M. Zalewski, “american fuzzy lop,” 2015, <http://lcamtuf.coredump.cx/afl/>. [Online]. Available: <http://lcamtuf.coredump.cx/afl/>
- [38] —, “american fuzzy lop technical “whitepaper”,” Tech. Rep., 2015. [Online]. Available: http://lcamtuf.coredump.cx/afl/technical_details.txt
- [39] A. Zeller, “Yesterday, my program worked. Today, it does not. Why?” in *Foundations of Software Engineering*, 1999, pp. 253–267.
- [40] J. Zou, R. M. Barnett, T. Llorido-Botran, S. Luo, C. Monroy, S. Sikdar, K. Teymourian, B. Yuan, and C. Jermaine, “PlinyCompute: A platform for high-performance, distributed, data-intensive tool development,” <http://arxiv.org/abs/1711.05573>, 2017.