CISPA Helmholtz-Zentrum i.G.

# Building Fast Fuzzers

Rahul Gopinath and Andreas Zeller

{rahul.gopinath, zeller}@cispa.saarland
CISPA - Helmholtz Center for Information Security, Saarbrücken, Germany

CISPA
HELMHOLTZ-ZENTRUM i.G.

# Building Fast Fuzzers

Rahul Gopinath and Andreas Zeller

{rahul.gopinath, zeller}@cispa.saarland

CISPA - Helmholtz Center for Information Security, Saarbrücken, Germany

(Dated November 19, 2019)

**Abstract**

Fuzzing is one of the key techniques for evaluating the robustness of programs against attacks. Fuzzing has to be *effective* in producing inputs that cover functionality and find vulnerabilities. But it also has to be *efficient* in producing such inputs quickly. Random fuzzers are very efficient, as they can quickly generate random inputs; but they are not very effective, as the large majority of inputs generated is syntactically invalid. Grammar-based fuzzers make use of a grammar (or another model for the input language) to produce syntactically correct inputs, and thus can quickly cover input space and associated functionality. Existing grammar-based fuzzers are surprisingly inefficient, though: Even the fastest grammar fuzzer dharma still produces inputs about *a thousand times slower* than the fastest random fuzzer. So far, one can have an effective or an efficient fuzzer, but not both.

In this paper, we describe how to build fast grammar fuzzers from the ground up, treating the problem of fuzzing from a programming language implementation perspective. Starting with a Python textbook approach, we adopt and adapt optimization techniques from functional programming and virtual machine implementation techniques together with other novel domain-specific optimizations in a step-by-step fashion. In our *F1* prototype fuzzer, these *improve production speed by a factor of 100–300 over the fastest grammar fuzzer* dharma. As *F1* is even 5–8 times faster than a lexical random fuzzer, we can find bugs faster and test with much larger valid inputs than previously possible.

## 1 Introduction

Fuzzing is a popular technique for evaluating the robustness of programs against attacks. The effectiveness of fuzzing comes from fast production and evaluation of inputs and low knowledge requirements about the program or its behavior—we only need to detect program crashes. These properties make fuzzing an attractive tool for security professionals.

To be effective, a fuzzer needs to sufficiently cover the variety of possible inputs, and it should produce inputs that can reach deep code paths. To reach deep code paths, the fuzzer needs to produce inputs that can get past the input parser—i.e., inputs that conform to the *input language* of the program under test.

Hence, the effectiveness of fuzzing can be improved by incorporating knowledge about the *input language* of the program under test to the fuzzer. As such languages are typically described by formal grammars, fuzzers that incorporate language knowledge are called *grammar fuzzers*. The inputs produced by grammar fuzzers are superior to pure random fuzzers because they easily pass through input validators, and hence a larger number of inputs can exercise deeper logic in the program. Today, a large number of tools exist [26, 58, 59, 2, 24, 16, 54] that all provide grammar-based fuzzing. These tools are also effective in their results: The LANGFUZZ grammar fuzzer [27] for instance, has uncovered more than 2,600 vulnerabilities in the JavaScript interpreters of Firefox and Chrome.

$\langle START \rangle ::= \langle expr \rangle$

$\langle expr \rangle ::= \langle term \rangle$
$\quad | \quad \langle term \rangle \text{ '+' } \langle expr \rangle$
$\quad | \quad \langle term \rangle \text{ '−' } \langle expr \rangle$

$\langle term \rangle ::= \langle factor \rangle$
$\quad | \quad \langle factor \rangle \text{ '}\ast\text{' } \langle term \rangle$
$\quad | \quad \langle factor \rangle \text{ '/' } \langle term \rangle$

$\langle factor \rangle ::= \langle integer \rangle$
$\quad | \quad \langle integer \rangle \text{ '.' } \langle integer \rangle$
$\quad | \quad \text{'+' } \langle factor \rangle$
$\quad | \quad \text{'−' } \langle factor \rangle$
$\quad | \quad \text{ '(' } \langle expr \rangle \text{ ')' }$

$\langle integer \rangle ::= \langle digit \rangle \langle integer \rangle$
$\quad | \quad \langle digit \rangle$

$\langle digit \rangle ::= \text{'0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'}$

Figure 1: A grammar for arithmetic expressions

```
1  expr_grammar = {
2      '<start>':   (['<expr>']),
3      '<expr>':    (['<term>', "+", '<expr>'],
4                    ['<term>', "-", '<expr>'],
5                    ['<term>']),
6      '<term>':    (['<factor>', "*", '<term>'],
7                    ['<factor>', "/", '<term>'],
8                    ['<factor>']),
9      '<factor>':  (["+", '<factor>'],
10                    ["-", '<factor>'],
11                    ["(", '<expr>', ")"],
12                    ['<integer>', ".", '<integer>'
       ],
13                    ['<integer>']),
14      '<integer>': (['<digit>', '<integer>'],
15                    ['<digit>']),
16      '<digit>':   (["0"], ["1"], ["2"], ["3"], ["
       4"].
17                    ["5"], ["6"], ["7"], ["8"], ["
       9"])
18  }
```

Figure 2: The grammar from Figure 1 as a Python *dict*

Grammar fuzzers have two downsides, though. The first problem is that an input grammar has to exist in the first place before one can use them for fuzzing. Programmers may choose not to use a formal grammar for describing the input, preferring to use ad hoc means of specifying the input structure. Even when such a grammar is available, the grammar may be incomplete or obsolete, and fuzzers relying on that grammar can develop blind spots. Recently, however, a number of approaches have been proposed to infer both regular languages [56] as well as context-free grammars either by grammar induction from samples [3] or by dynamic grammar inference from program code [28]. While these approaches require valid sample inputs to learn from, recent work by Mathis et al. [37] and Blazytko et al. [6] suggests that it is possible to automatically generate inputs that cover all language features and thus make good samples for grammar induction—even for complex languages such as JavaScript. It is thus reasonable to assume that the effort for specifying grammars might be very reduced in the future.

```
1  def gen_key(grammar, key):
2      if key not in grammar: return key
3      return gen_alt(grammar, random.choice(grammar
       [key]))

5  def gen_alt(grammar, alt):
6      # Concatenate expansions of all elements in
       alt
7      return ''.join(gen_key(grammar, t) for t in
       alt)

9  gen_key(expr_grammar, '<start>')
```

Figure 3: A simple grammar based fuzzer in Python that uses the grammar from Figure 2

The second problem with grammar fuzzing, though, is

that it is *slow*—at least in the implementations we see today. In principle, fuzzing with a grammar is not hard. Given a context-free grammar such as in Figure 1, we first convert it to a native data structure as given in Figure 2. We then start with the *nonterminal* representing the starting point <start>. Each *nonterminal* corresponds to possibly multiple alternative rules of expansion. The fuzz result corresponds to one of the rules of expansion for <start>, choosing expansions randomly.

An expansion is a sequence of tokens. For each token, if the token is a *terminal* symbol (a string with no further expansion), then the string goes unmodified to its corresponding location in the final string. If it is a *nonterminal*, we start expansion again until the expansion results in *terminal* nodes. A Python program implementing this approach is given in Figure 3.

Unfortunately, this simple approach has a problem when dealing with highly recursive grammars. As the fuzzer uses the recursion as a method to explore nested structures, it can deplete the stack quite fast. This can be fixed by limiting the alternatives explored to the lowest cost alternatives. We show how this can be done in Section 4.

However, once the problem of stack exhaustion is fixed, we have another problem While this production approach is not naïve, it is not exactly fast either. Using a simple *Expr* grammar from the "Fuzzing Book" textbook chapter on Grammars [65], it provides a throughput of 103.82 kilobytes per second.[1] If one wants *long* inputs of, say, ten megabytes to stress test a program for buffer and stack overflows, one would thus have to wait for a minute to produce one single input. Now, compare this to a pure random fuzzer, say using dd if=/dev/urandom directly as input to a program (that we call the *dev-random* fuzzer). Using *dev-random* achieves a speed of 23 MiB per second, which is over *a hundred times faster* than even the fastest grammar-based fuzzer dharma, which produces 174.12 KiB/s on *Expr*.

In this paper, we *show how to build fast grammar fuzzers from the ground up,* treating the problem of fuzzing from a programming language implementation perspective. Starting from the above Python code, we apply and develop a number of techniques to turn it into an extremely fast fuzzer. On a CSS grammar, our *F1* prototype yields a final throughput of 80,722 kilobytes per second of valid inputs on a single core. This is *333 times faster* than dharma, the fastest grammar fuzzer to date (which only produces 242 KiB/s for CSS), and even three times as fast as the *dev-random* fuzzer. These results make our *F1* prototype the fastest fuzzer in the world that produces valid inputs—by a margin.

A fast fuzzer like *F1* not only can be used to generally save CPU cycles for fuzzing, but also to produce large inputs for areas in which high throughput is required, including CPUs [36], FPGAs [32], Emulators [18], and IOT devices [68], all of which allow high speed interactions. It also allows for stress testing hardware encoders and decoders such as video encoders that require syntactically valid, but rarely seen inputs can profit from a fast grammar

---

[1]In this paper, 1 "kilobyte" = 1,024 bytes = 1 KiB, and 1 "megabyte" = 1,024 × 1,024 bytes = 1 MiB.

fuzzer. Side channel attacks on hardware implementing encryption and decryption may require that the certificates use an envelope with valid values, with the encrypted values as one of the contained values, where again the *F1* fuzzer can help. Finally, many machines use a hardware implemented TCP/IP stack, which require structured input. Given that the networking stack can accept high throughput traffic, *F1* can be of use fuzzing these devices. In all these settings, speed becomes a priority.

The remainder of this paper is organized as follows. After discussing related work (Section 2) and introducing our evaluation setup (Section 3), we first discuss methods of limiting expansion depth (Section 4). We then make the following contributions:

**Grammar Compilers.** We show how to *compile* the grammar to *program code,* which, instead of interpreting a grammar structure, directly acts as a producer (Section 5). A compiled producer is much faster than a grammar interpreter.

**Compiling to Native Code.** We can further speed up production by directly compiling to native code, e.g. producing C code which is then compiled to native code (Section 6). Compared to languages like Python (in which most grammar fuzzers are written), this again yields significant speedups.

**Supercompilation.** We introduce the novel notion of *supercompiling* the grammar, inlining production steps to a maximum (Section 7.2). This results in fewer jumps and again a faster fuzzer.

**System Optimizations.** We explore and apply a number of *system optimizations* (Section 8) that add to the efficiency of producers on practical systems, notably high-speed random number generation and quick file output.

**Production Machines.** We introduce the notion of a fuzzer as a *virtual machine interpreter* (Section 9) that interprets the random stream as bitcode, and explore the various alternatives in efficiently implementing virtual machines.

**The Fastest Fuzzer.** In a detailed evaluation (Section 10) with multiple grammars and settings, we compare the performance of our *F1* prototype implementation against state-of-the-art grammar fuzzers, and find it is 200–300 times faster than dharma, the fastest grammar fuzzer to date.

After discussing related work (Section 11), Section 12 discusses current limitations and future work. Finally, Section 13 closes with conclusion and consequences. All of the source code and data referred to in this paper is available as self-contained Jupyter Notebooks, allowing to run the *F1* prototype and fully replicate and extend all experiments.

## 2  Background

### 2.1  Fuzzing and Test Generation

Fuzzing is a simple but highly effective technique for finding vulnerabilities in programs. The basic idea of a fuzzer is to quickly generate strings and evaluate these strings as input to the program under fuzzing. If any of these inputs trigger a program crash or other surprising behavior, or the program execution falls afoul of sanity checks, it is an indication of a possible vulnerability that may be exploited [66].

Fuzzers, and testing techniques in general, are traditionally classified into *whitebox* and *blackbox* (and sometimes *greybox*) techniques [35]. Whitebox techniques assume availability of source code, and often use program analysis tools to enhance the effectiveness of fuzzing, and example of which is KLEE [10] which uses symbolic execution of source code for generating inputs. Another is AFL [41], which makes use of specially instrumented binaries for tracking coverage changes. Blackbox techniques on the other hand, do not require the availability of source code. The greybox fuzzers assumes the availability of at least the binary under fuzzing, and often work by instrumenting the binary for recovering runtime information. An example of such a fuzzer is angr [49] which can symbolically or concolically execute a binary program to produce fuzzing inputs. The main problem with whitebox fuzzers is the requirement of source code (and its effective utilization). The problem is that time spent on analysis of source code can reduce the time spent on generating and executing inputs.

The black box fuzzers (and the input generation part of the whitebox and greybox fuzzers) are traditionally classified as *mutational fuzzers* and generational fuzzers.

In mutational fuzzing, a corpus of seed inputs are used as the starting point, which are then mutated using well defined operations to generate newer inputs that can be evaluated [67]. AFL, libFuzzer [48] and Radamsa [44] are some of the well known mutational fuzzers. Mutational fuzzers require relatively little knowledge of the input language and relies on the given seed corpus to produce sufficiently valid strings to reach deep code paths. One of the problems of mutational fuzzers is that they are limited by the seed corpora. That is, these fuzzers are inefficient when it comes to generating inputs that are markedly different from the ones in the seed corpora. Unfortunately, these kinds of *unexpected* inputs often have higher chance of triggering bugs than the expected kind. This limits their utility. A second problem is that due to the fuzzers ignorance of the input format, the mutations introduced frequently fall afoul of the input validator often in trivial ways. Hence, only a small portion of the mutated inputs reach the internal code paths where the bugs hide. That is, if one is fuzzing an application that accepts a string in an XML format, one might have more success in fuzzing the main application itself rather than the XML parser itself which is likely to be well tested. Hence, it is of interest to the fuzzer to generate valid XML fragments in the first place.

Generational fuzzers on the other hand, relies on some model of the input required by the program to generate

valid inputs. The model may be fixed as in the case of Csmith [63] which generates valid C programs and JSFun-Fuzz [46] which targets Javascript. Fuzzers such as Gram-fuzz[24], Grammarinator [26], Dharma [38], Domato [20], and CSS Fuzz [46] allow the user to specify the input format as a context-free grammar. For those contexts where a finite state automata is sufficient, some fuzzers [13, 60] allow an FSM as the input model. Fuzzers that allow context sensitive constraints on inputs are also available [17].

## 2.2 Context-Free Grammars

A context-free grammar is a formal grammar that specifies how a set of strings can be generated. A formal grammar is a set of rules (called production rules) for rewriting a sequence of symbols. A production describes how a given symbol (called a *nonterminal*) should be replaced. If a symbol is not associated with a production it is called a *terminal*, and represents itself in the generated output. The rewriting starts in the symbol representing the starting point, called the `start` symbol. In the Figure 1, the start symbol is `<start>`, and the production corresponding to it is `<expr>`. Similarly, the *nonterminal* symbol `<expr>` has three production rules:

```
1  '<expr>': (['<term>', "+", '<expr>'],
2            ['<term>', "-", '<expr>'],
3            ['<term>']),
```

When fuzzing, one of these production rules is chosen stochastically for rewriting `<expr>`. The first rule specifies that `<expr>` is rewritten as a sequence `<term>+<expr>`, where `<term>` is again another *nonterminal* symbol while + is a *terminal* symbol that is represented by itself in the output.

The *nonterminal* symbol `<term>` gets expanded to a string containing `<factor>` just like `<expr>` was expanded into a string containing `term`. `<factor>` has five production rules specifying how it may be expanded.

```
1  '<factor>': (["+", '<factor>'],
2              ["-", '<factor>'],
3              ["(", '<expr>', ")"],
4              ['<integer>', ".", '<integer>'],
5              ['<integer>']),
```

If we assume that the production rule `<integer>` was chosen, then we get to choose from the expansions of `<integer>` given by:

```
1  '<integer>': (['<digit>', '<integer>'],
2               ['<digit>']),
```

If we assume that the second production rule was chosen next, it contains a single *nonterminal* symbol *<digit>*. The `<digit>` has ten production rules, each of which has a single *nonterminal* symbol.

```
1    '<digit>':  (["0"], ["1"], ["2"], ["3"], ["
     4"].
2                ["5"], ["6"], ["7"], ["8"], ["
     9"])
```

If say `5` was chosen as the rule, then the first `<factor>` would be replaced by `5`, giving the expression `5*<factor>`. Similarly, the second `<factor>` may also be replaced by say the fourth production rule:

$$\texttt{factor} \rightarrow \texttt{<integer>.<integer>}$$

This gives the expression `5*<integer>.<integer>`. Starting with unexpanded symbols on the left, assuming the second expansion for `integer` was chosen, we have `digit`. Say `digit` expanded to `3`, the above expression is transformed to:

$$\rightarrow \texttt{5*3.<integer>}$$

Going through similar expansions for the last `integer` again, we get:

$$\rightarrow \texttt{5*3.8}$$

which is the final expression.

Context-free grammars are one of the common ways to specify file formats, or as the first level (parser) format for most programming languages. The ability for patterns to be named, and reused makes it easier to use than regular expressions, while the context-free aspect makes it easy to write a parser for it when compared to more complex grammar categories. Given that most parsers accept context-free grammars, writing fuzzers based on context-free grammars can be effective in fuzzing the programs that use these parsers.

## 3 Method of Evaluation

In order to get a fair assessment of various fuzzers, it is import to ensure that we remove as much of external factors as possible that can skew the results. To ensure that each fuzzer got sufficient time to cache execution paths in memory, we started with a small warm up loop of ten iterations. Next, to avoid skew due to different seeds, we chose random seeds from zero to nine[2] and computed the average of ten runs using these seeds.

We needed to make sure that there was a level playing ground for all grammar fuzzers. For grammar fuzzers, it is easier to produce relatively flat inputs such as say *true* or *false* in JSON grammar than one that requires multiple levels of nesting. However, when fuzzing, these inputs with complex structure are usually more useful as these inputs have a higher chance of producing interesting behavior. A metric such as the number of inputs per second (as used with mutational fuzzers) unfairly penalizes the grammar fuzzers that produce inputs with complex structure. Hence, rather than the number of inputs produced, we opted to simply use the *throughput*—kilobytes of output produced per second as the metric to judge the fuzzer performance.

We saw that the maximum depth of recursion had an impact on the throughput. Hence, we evaluated all the fuzzers (that allowed setting a maximum depth) with similar depth of recursion, with depth ranging from 8 to 256, with the timeout set to an hour (36,00 seconds).

---

[2] The random seed is used to initialize the pseudorandom number generator. While the seed values are close to each other, the random numbers generated from the initial seed values are not close to each other, as even small differences in bit patterns have a large impact. Hence, we chose to use the numbers from 0 to 9 to be as unbiased as possible (a set of random, random seeds may have an unforeseen bias).

Our experiments were conducted on a *Mac OS X* machine with nacOS *10.14.5*. The hardware was *MacBookPro14,2* with an *Intel Core i5* two-core processor. It had a speed of 3.1 Ghz, an L2 cache of 256 KB, and L3 cache of 4 MB. The memory was 16 GB. All tools, including our own, are single-threaded. Times measured and reported is the sum of user time and system time spent in the respective process.

## 4 Controlling Free Expansion

```
1  def d_key(key, seen):
2      if key not in grammar: return 0
3      if key in seen: return inf
4      return min(d_alt(rule, seen | {key})
5              for rule in grammar[key])

7  def d_alt(rule, seen):
8      return max(d_key(s, seen) for s in rule) + 1
```

Figure 4: Computing the minimum depth of expansion $\mu$-depth

The simple approach given in Figure 3 is naive in that it provides no control over how far to go when expanding the *nonterminal* tokens. If the grammar is highly recursive, it can lead to stack exhaustion before all *nonterminal* symbols are effectively expanded to a *terminal* string. There is a solution to this problem[3]. Given a context-free grammar, for any given key, let us define a minimum depth of expansion ($\mu$-depth) as the minimum number of levels of expansion needed (stack depth) to produce a completely expanded string corresponding to that key. Similarly the $\mu$-depth of a rule is the maximum of $\mu$-depth required for any of the tokens in the rule. One can hence compute the $\mu$-depth for each of the *nonterminal* symbols. The idea is similar to the fuzzer in Figure 3. Given a token, check if the token is a *nonterminal*. If it is not, the $\mu$-depth is zero. If it is a *nonterminal*, compute the $\mu$-depth of each of the alternative expansion rules corresponding to it in the grammar. The *nonterminal*'s $\mu$-depth is the minimum of the $\mu$-depth of its corresponding rules. If we detect recursion, then the $\mu$-depth is set to $\infty$. The algorithm for $\mu$-depth computation is given in Figure 4.

Once we compute the minimum depth of expansion for every key, one can modify the naive algorithm in Figure 3 as follows: Start the generation of input with a maximum free-stack budget. When the number of stack frames go beyond the budget, choose only those expansions that have the minimum cost of expansion. With this restriction, we can now compute the throughput of our fuzzer for the *Expr* grammar: 103.82 KiB/s with the free expansion depth set to eight, and a maximum of 133 KiB/s when the free expansion depth is set to 32.

### 4.1 Precomputed Minimum Depth Expansions

We have precomputed the expansion cost. But is our implementation optimal? On tracing through the program execution, one can see that once the free stack budget is exhausted, and the program switches to the minimum depth strategy, there is only a small pool of strings that one can generate from a given key, and all of them have exactly the same probability of occurrence. Hence, we can precompute this pool. Essentially, we produce a new grammar with only the minimum depth expansions for each *nonterminal*, and exhaustively generate all strings using an algorithm similar to Figure 3. Precomputing string pools gives us a throughput of 371.76 kilobytes per second for an expansion depth of 8, and a throughput of 420.14 for an expansion depth of 32.

Precomputing this pool of strings for each keys, and using the pool when the free stack budget is exhausted gets us only a modest improvement. Is there any other optimization avenue?

## 5 Compiling the Grammar

One of the advantages of using precomputed string pools is that it eliminates the cost of checking whether a token is a *nonterminal* or not, and also looking up the corresponding rules of expansion for the *nonterminal*. Can we eliminate this lookup completely?

A grammar is simply a set of (possibly recursive) instructions to produce output. We have been interpreting the grammar rules using the fuzzer. Similar to how one can compile a program and generally make it faster by removing the interpreting overhead, what if we compile the program so that the resulting program produces the fuzzer output?

The idea is to essentially transform the grammar such that each *nonterminal* becomes a function, and each *terminal* becomes a call to an output function (e.g `print`) with the symbol as argument that outputs the corresponding *terminal* symbol. The functions corresponding to *nonterminal* symbols have conditional branches corresponding to each expansion rule for that symbol in the grammar. The branch taken is chosen randomly. Each branch is a sequence of method calls corresponding to the *terminal* and *nonterminal* symbols in the rule.

As before, we incorporate the optimization for `max_depth` to the functions. A fragment of such a compiled grammar is given in Figure 6.

---

[3] The text book approach given in the chapter on Grammar Fuzzing [64] provides two controls — the maximum number of *nonterminal* symbols and the number of *trials* to attempt before giving up.

```
1  def expr(depth):
2      d = depth + 1
3      if d > max_depth: return choice(s_expr)
4      c = random.randint(3)
5      if c == 0: s = [term(d), "*", term(d)]
6      if c == 1: s = [term(d), "/", term(d)]
7      if c == 2: s = [term(d)]
8      return ''.join(s)
9  ...
```

Figure 6: A fragment of the compiled `expression` grammar. The `s_expr` contains the possible minimum cost completion strings for `expr` *nonterminal*.

Compiling the grammar to a program, and running the program gives us a faster fuzzer, with a throughput of 562.22 kilobytes per second for a free expansion depth of 8, and 714.08 kilobytes per second for a free expansion depth of 32. The pure random Python fuzzer in the Fuzzingbook chapter on Fuzzing [66] achieves 1259 kilobytes per second. That is, this represents a slowdown by a factor of two.

# 6 Compiling to a Faster Language

While Python is useful as a fast prototyping language, it is not a language known for the speed of programs produced. Since we are compiling our grammar, we could easily switch the language of our compiled program to one of the faster languages. For this experiment we chose *C* as the target. As before, we transform the grammar into a series of recursive functions in C.

As we are using a low level language, and chasing ever smaller improvements, a number of things can make significant difference in the throughput. For example, we no longer do dictionary lookups for string pools. Instead, we have opted to embed them into the produced program, and do array lookups instead. Similarly, we use a case statement that directly maps to the output of the modulus operator. With the new compiled grammar, we reach a throughput of 14,775.99 kilobytes per second for a depth of 8 and 15,440.12 kilobytes per second for a depth of 32.

# 7 The Grammar as a functional DSL

We have seen in the previous section how the grammar can be seen as a domain specific language for generating strings, and the compiled fuzzer as representing the grammar directly. Another thing to notice is that the language of grammar is very much a pure functional language, which means that the ideas used to make these languages faster can be applied to the compiled grammar. We examine two such techniques.

## 7.1 Partial Evaluation

The first is partial evaluation [31, 50]. The idea of partial evaluation is simple. In most programs, one does not have to wait until runtime to compute a significant chunk of the program code. A large chunk can be evaluated even before the program starts, and the partial results included in the program binary directly. Indeed, pre-computing the minimum depth expansions is one such. We can take it further, and eliminate extraneous choices such as grammar alternatives with a single rule, and inline them into the produced program. We could also embed the expansions directly into parent expansions, eliminating subroutine calls. Would such inlining help us?

Using this technique, gets us to 13,945.16 kilobytes per second for a depth of 8 and 15021.02 kilobytes per second for a depth of 32. That is, partial evaluation by simple elimination of choices is not very helpful.

## 7.2 Supercompilation

Another technique is *supercompilation*—a generalization of partial evaluation [50] that can again be adapted to context-free grammars.[4]

The idea of supercompilation is as follows: The program generated is interpreted abstractly, using symbolic values as input (*driving*). During this abstract interpretation, at any time the interpretation encounters functions or control structures that does not depend on input, the execution is inlined. When it encounters variables that depend on the input, a model of the variable is constructed, and the execution proceeds on the possible continuations based on that model. If you find that you are working on a similar expression (in our case, a *nonterminal* expansion you have seen before) [5], then terminate exploring that thread of continuation, produce a function corresponding to the expression seen, and replace it with a call to the function [6]. During this process, redundant control and data structures are trimmed out, leaving a residual program that performs the same function as the original but with a lesser number of redundant steps.

In the case of the language of context-free grammars, the input is given by the random stream which is used to determine which rule to expand next. Hence, during the process of supercompilation, any *nonterminal* that has only a single rule gets inlined. Further, parts of expansions that have deterministic termination are also inlined, leaving only recursive portions of the program as named functions. Finally, the functions themselves are transformed. The *nonterminal* symbols are no longer functions themselves. Instead, they are inlined, and the individual rule expansions become functions themselves. We will see how these functions lead to a newer view of the fuzzer in the next sections.

Figure 7 shows a fragment of the supercompiled expression grammar.

---

[4] The original supercompiler is due to Turchin [52], and involves removing redundancies present in the original program by abstract interpretation. The language in which it was implemented *Refal* [51] is a term rewriting language reminiscent of the syntax we use for representing the context-free grammar.

[5] This is called *renaming* in the language of supercompilation

[6] This is called *folding*.

```
1  void expr_1(int depth) {
2    int d1 = depth − max_depth;
3    if (d1 >= 0) {
4      out(s_term[random() % n_s_term]);
5    } else {
6      int d2 = d1 − 1;
7      switch(random() % 5) {
8      case 0:
9        if (d2 >= 0) {
10          out(s_factor[random() % n_s_factor]);
11        } else {
12          switch(random() % 5) {
13          case 0:
14            factor_0(depth+4);
15  ...
```

Figure 7: A fragment of the supercompiled grammar in C. As before, s_expr is a precomputed array of minimum cost completion strings for *nonterminal* expr. The function expr_1 represents the second production rule for expr. The function corresponding to the *nonterminal* expr itself is inlined

Supercompiled grammar fuzzer produces a throughput of 14,441.95 kilobytes at a depth of 8 and 14,903.12 kilobytes at a depth of 32. Supercompilation does not make much of a difference in the case of *Expr*, but as we will see later in the paper, it helps with other subjects.

# 8 System-Level Optimizations

Compilation can improve the performance of interpreted languages, while supercompilation can improve the performance of functional languages. We have not yet considered how both *interact with their environment,* though. Generally speaking, a fuzzer needs two functionalities: *random numbers* required to choose between expansions, and *system output* to send out the produced string. We will see how to optimize them next.

## 8.1 Effective Randomness

Is our generated fuzzer the best one can do? While profiling the fuzzer, two things stand out. Generating random values to choose from, and writing to a file. Can we improve these parts? Looking at how random values are used, one can immediately see an improvement. The pseudo-random generators (PRNG) used by default are focused on providing a strong random number with a number of useful properties. However, these are not required simple exploration of input space. We can replace the default PRNG with a faster PRNG[7]. Second, we use the modulus operator to map a random number to the limited number of rules corresponding to a key that we need to choose from. The modulus and division operators are rather costly. A more performant way to map a larger number to a smaller number is to divide both and take the upper half in terms of bits. Another optimization is to recognize that we are rather wasteful in using the random number generator. Each iteration of PRNG provides us with a 64 bit number, and the number of choices we have rarely exceed 256. That is, we can simply split the generated random number to eight parts and use one part

at a time. Finally, for better cache locality, it is better to generate the needed random numbers at one place, and use them one byte at a time when needed. In fact, we can pre generate a pool of random numbers, and when the pool is exhausted, trigger allocation of new random bits, or when the context permits, reuse the old random bits in different context. These micro optimizations can provide us with a rather modest improvement. The throughput is 11,733.98 kilobytes per second for depth of 8, and 22,374.40 kilobytes per second for depth of 32.

## 8.2 Effective Output

As we mentioned previously, output is one of the more performance intensive operations. Ideally, one would like to generate the entire batch of outputs in memory, and write to the file system in big chunks.

Switching from fputc to generating complete items, and writing them one at a time gave us a throughput of 19,538.411 kilobytes per second for a depth of 8 and 69,810.755 kilobytes per second for a depth of 32.

One way to improve the output speed is to use memory mapped files to directly write the output. One of the problems here is that for the mmap() call, one need to know the size of the file in advance. We found that we can tell the operating system to map the file with the maximum size possible, which the OS obeys irrespective of the actual availability of space. Next, we can write as much as required to the mapped file. Finally, we call truncate on the file with the number of bytes we produced. At this point, the operating system writes back the exact amount of data we produced to the file system.

Unfortunately, mmap() performance was rather variable. We obtained a throughput of 10,722.41 KiB/s on a depth of 8, and 56,226.329 KiB/s on a depth of 32, and we found this to fluctuate. That is, depending on mmap() should be considered only after taking into consideration different environmental factors such as operating system, load on the disk, and the memory usage.

A fuzzer does not have to write to a file, though. If the program under test can read input directly from memory, we can also have the fuzzer write only to memory, and then pass the written memory to the program under test. Obviously, skipping the output part speeds up things considerably; we obtain a throughput of 81,764.75 KiB/s on a depth of 32. (In the remainder of the paper, we will assume we write to a file.)

# 9 Production Machines

Is this the maximum limit of optimization? While we are already compiling the grammar to an executable, our supercompiled program is in a very strong sense reminiscent of a virtual machine. That is, in our case, the random bits that we interpret as a choice between different rules to be used to expand can be considered the byte stream to interpret for a virtual machine. The usual way to implement a virtual machine is to use switched dispatch. Essentially, we implement a loop with a switch statement within that executes

---

[7]We used the Xoshiro256** generator http://xoshiro.di.unimi.it/xoshiro256starstar.c

the selected opcode.

### 9.1 Direct Threaded VM

However, one of the most performant ways to implement a virtual machine is something called a *threaded code* [4, 19]. A pure bytecode interpreter using switched dispatch has to fetch the next bytecode to execute, and lookup its definition in a table. Then, it has to transfer the execution control to the definition. Instead, the idea of direct threading is that the bytecode is replaced by the starting address of its definition. Further, each opcode implementation ends with an epilogue to dispatch the next opcode. The utility of such a technique is that it removes the necessity of the dispatch table, and the central dispatch loop. Using this technique gets us to 17,848.876 kilobytes per second for a depth of 8 and 53,593.384 kilobytes for a depth of 32.

### 9.2 Context Threaded VM

One of the problems with direct threading [19] is that it tends to increase branch misprediction. The main issue is that `computed goto` targets are hard to predict. An alternative is context threading [5]. The idea is to use the `IP` register consistently, and use `call` and `return` based on the value of `IP` register when possible. This is different from simply using subroutine calls as no parameters are passed, and no prologue and epilogue for subroutine calls are generated. Doing this requires generating assembly, as *C* does not allow us to directly generate naked *call* and *return* instructions. Hence, we generated *X86-64* assembly corresponding to a context threaded VM, and compiled it. A fragment of this virtual machine in pseudo-assembly is given in Figure 9.

```
1  gen_member_0:
2      val = map(2)
3      call *gen_ws[val]
4      val = map(1)
5      call *gen_string[val]
6      val = map(2)
7      call *gen_ws[val]
8      *out_region = ':'
9      incr out_region
10     val = map(1)
11     call *gen_element[val]
12     ret
```

Figure 9: A fragment of the context threaded interpreter for grammar VM that generates a JSON object.

Context threading got us to 14,805.989 kilobytes per second for a depth of 8 and 16,799.153. While for *Expr*, the direct threading approach seems slower than the context threading, as we will see later, *Expr* is an outlier in this regard. The context threading approach generally performs better. This final variant is actually our fastest fuzzer, which we call the *F1* fuzzer.

---

[8] https://www.blackarch.org/

## 10 Evaluation

So far, we have checked the performance of our fast fuzzing techniques only on one, admittedly very simple grammar. How do they fare when faced with more complex input formats? And how do they compare against state-of-the-art grammar fuzzers? To this end, we evaluate our techniques on three grammars with well-known and nontrivial industry formats (CSS, HTML, and JSON), and compare them against four state of the art tools.

Our results are summarized in Figure 10, using three different settings for the maximum depth (8, 32, and 128). The vertical axis lists the throughput achieved by each tool for the respective grammar; note the usage of a logarithmic scale to capture the large differences. To account for differences due to randomness, the state-of-the-art tools were run 2 times with 1,000 inputs generated per run; our (faster) tools were run 100 times, also with a 1,000 inputs generated per run. Times listed are averages over all runs and inputs generated.

### 10.1 Textbook Fuzzer

The fuzzingbook from Zeller et al. [64] specifies a variety of grammar based fuzzers. We chose the simplest one with no frills in the interest of performance. The grammar syntax is somewhat similar to our own but uses strings for rules with specially demarcated tokens for *nonterminal* symbols.

The authors in [64] make it clear that their interest is in teaching how grammar-based fuzzers work, and that performance is not their main goal. This also shows in our evaluation: The throughput of the fuzzingbook `GrammarFuzzer` is 3.7 KiB/s for CSS, 31.9 KiB/s for HTML, and 5 KiB/s for JSON. As it comes to producing high volumes of input, this marks the lowest end of our evaluation; for higher maximum depths than 8, it would not produce results in time at all.

### 10.2 Grammarinator

Grammarinator [26] is the state of the art fuzzer by Renáta Hodován et al. It is written in Python and accepts a context-free grammar in the *ANTLR* grammar format. The grammarinator is the only fuzzer in our set of competing fuzzers that compiles the grammar to Python code first before fuzzing. Grammarinator is also the only grammar fuzzer included in the *BlackArch*[8] Linux distribution for pen testers.

Grammarinator is faster than the fuzzingbook grammar fuzzer: Given a maximum depth of 8, it achieves a throughput of 55.8 KiB/s, 113.7 KiB/s, and 7 KiB/s for CSS, HTML, and JSON grammars respectively. It is slower, however, than the other state-of-the-art tools, not to speak of the fast fuzzers introduced in this paper. This may be due to grammarinator spending some effort in *balancing* its outputs. Grammarinator maintains a *probability* for each expansion alternative, and every time an expansion is chosen, it reduces this probability, thus favoring other alternatives. This balancing costs time, and this slows down grammarinator. The benefit could be a higher variance of produced strings, but this is a feature we did not evaluate in
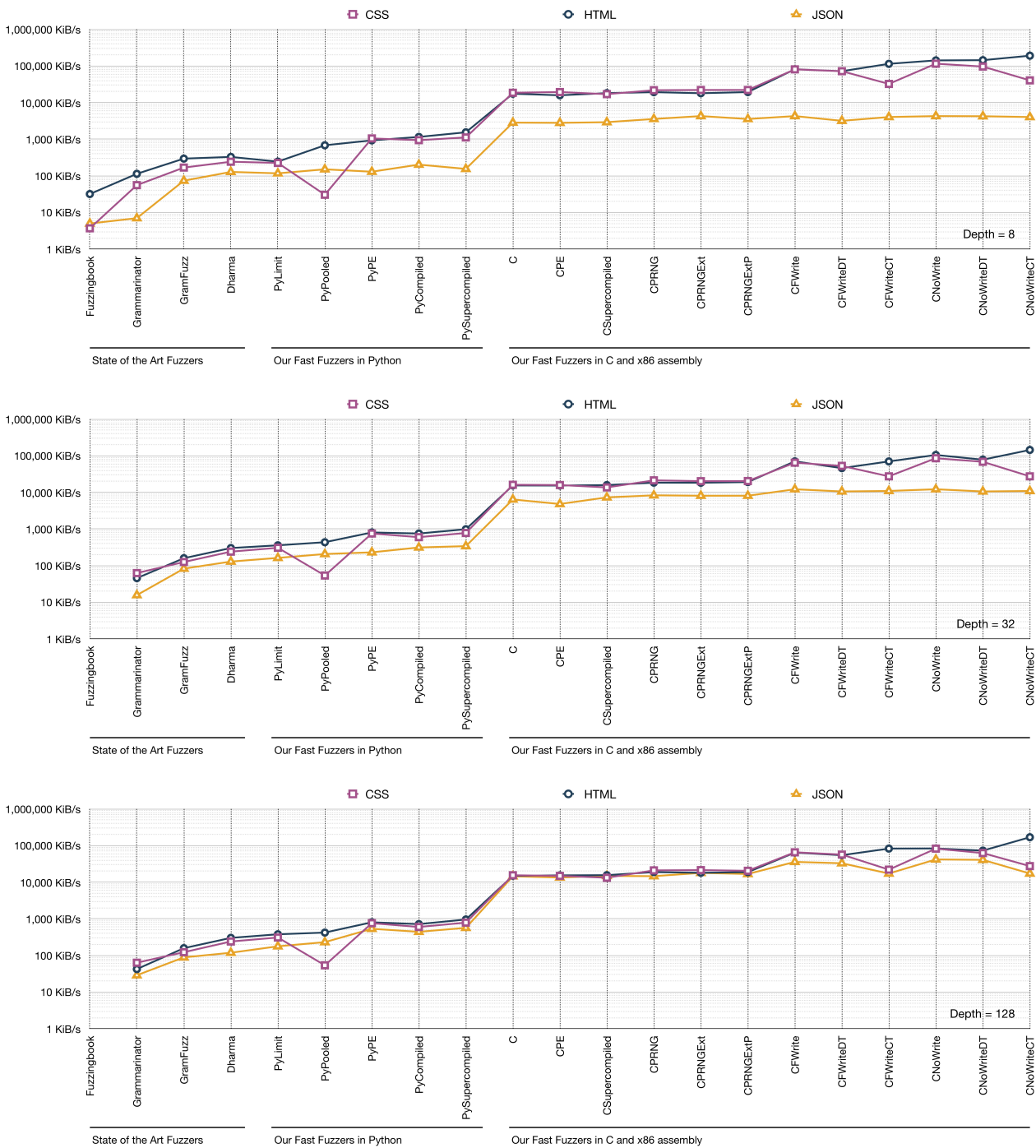
Figure 10: Fuzzer throughput with different grammars and depth limits. **PyLimit** = Python with limit (Section 4). **PyPooled** = Python with precomputed string pools (Section 4.1). **PyPE** = Python with partial evaluation (Section 7.2). **PyCompiled** = Compiled Python producer (Section 5). **PySuperCompiled** = Supercompiled Python producer (Section 7.2). **C** = Compiled C producer (Section 6); with partial evaluation (**CPE**; Section 7.1). **CSuperCompiled** = Supercompiled C producer (Section 7.2). **CPRNG** = C producer with faster random number generator (Section 8.1) and precompiled random numbers (**CPRNGExt**). **CFWrite** = C Producer Machine (Section 9) using `fwrite()` (Section 8.2); with direct treading (**CFWriteDT**, Section 9.1) and context treading (**CFWriteCT**, Section 9.2). **CNoWrite, CNoWriteDT, CNoWriteCT** = C Producer Machine writing to memory (Section 8.2).

9

the present paper.

## 10.3 Gramfuzz

Gramfuzz [24] is another grammar based fuzzer written in Python.[9] It uses its own format for specifying a context-free grammar. The fuzzer allows specifying a recursion depth to limit expansions.

Due to its simplicity, Gramfuzz wins over grammarinator for all grammars in our evaluation. Again, given a maximum recursion depth of 8, it achieves a throughput of 168.3 KiB/s, 295.6 KiB/s, and 73.8 KiB/s for CSS, HTML, and JSON grammars respectively, which is 5–10 times as fast as grammarinator. This is the more interesting as gramfuzz interprets the grammar structure instead of compiling it to code.

## 10.4 Dharma

The dharma[10] fuzzer from Mozilla [38] is a grammar based fuzzer that has received wide attention from the industry. It is (again) written in Python, and uses its own custom format for specifying the grammar. We found that the dharma grammar syntax was a little unwieldy in that there is no way to specify empty rules, and empty strings. Hence, we had to work around this using white space for our evaluation.

Dharma's throughput again improves over gramfuzz; with a recursion depth of 8, it reaches 242 KiB/s, 328.9 KiB/s, and 128.1 KiB/s for CSS, HTML, and JSON grammars respectively. This makes dharma the fastest state-of-the-art fuzzer, and hence the baseline our faster techniques can compare against.

## 10.5 Fast Python Fuzzers

We now discuss our own fast fuzzers written in Python and/or compiling to Python code. **PyLimit** is the very simple fuzzer with depth control introduced in Section 4. With 225.3 KiB/s, 244.8 KiB/s, and 117.2 KiB/s for CSS, HTML, and JSON grammars respectively, it is only marginally slower than dharma. This is the baseline our other fuzzers compare against. We see that partial evaluation (**PyPE**, 1051.3 KiB/s, 931.2 KiB/s, 129.7 KiB/s) already brings large speedups. Supercompilation (**PySuperCompiled**) then boosts the throughput to 1,119.2 KiB/s, 1,544.5 KiB/s, and 154.9 KiB/s. For JSON, however, supercompilation is slower than regular compilation (202.7 KiB/s); such effects can occur if the generated code is too large and hits processor caching limits.

## 10.6 Fast Fuzzers in C and x86 assembly

Switching from Python to C as the language for producers and compiling the C code to machine code brings a huge performance gain. For a maximum depth of 8, the throughput of the compiled C producer (**C**) is 18,558.9 KiB/s, 17,481.3 KiB/s, and 2,838.7 KiB/s on CSS, HTML, and JSON grammars respectively. This is about 20 times as fast as the compiled Python producer (**PyCompiled**). Partial evaluation and supercompilation bring further benefits, notably with higher maximum recursion depths (32 and 128).

Using pseudo-random number generators speed up producers by about 25%, as do precompiled random numbers.

The next big performance boost comes from introducing producer machines. Using a producer machine (**CFWrite**) that also incorporates the above random number optimizations, achieves a throughput of 80,722.1 KiB/s, 81,323.7 KiB/s, and 4,281.5 KiB/s, respectively, which again is a factor of four higher than the fastest C producer. Adding direct and context treading can further improve performance, depending on depth and grammar; on HTML, the **CFWriteCT** variant achieves a throughput of 141,927.6 KiB/s, which sets a record in our evaluation. (Actually, we still can be faster, but only by omitting the writing part; the last three columns in our charts (**NoWrite**) show the throughput when writing to memory only.)

## 10.7 Discussion

All in all, we have seen that our techniques—notably compilation, using C/assembly as producer languages, and finally building dedicated virtual machines—can considerably speed up fuzzing. Compared to dharma, the fastest fuzzer in our evaluation, one can expect speedup factors of 100–300.

Does this mean that one can now fuzz several hundred times faster than before? If we fuzz fast programs, individual functions, or hardware components, we may indeed see significant overall speed ups. If most of the overall runtime is spent on running the program under test, though, a faster fuzzer will still save CPU time, but the relative benefits of a faster fuzzer will be smaller. However, the increased speed allows to stress test such programs with much *larger* inputs than before, turning generation time from minutes to tenths of seconds.

Finally, let us not forget that the main ingredient for the speed of grammar-based fuzzing is not so much the optimizations as described in this paper, but the *grammar itself.* If we were set to compare against random character-based fuzzers such as AFL [41] and counted the number of *valid inputs* produced per second, then even the simplest grammar fuzzer would easily outperform its competition. This highlights the need for techniques to infer grammars from existing programs [29, 3, 6], as they would end up in an extremely fast and extremely effective fuzzer.

## 11 Related Work

Our related work falls into the following categories.

**Grammars.** While informally specified grammars have been around from the origins of human language, the first formalization occurred in 350 BC by Dakṣiputra Pāṇini in Ashṭādhyāyī [43]. Chomsky [11] introduced the formal models for describing languages, namely: finite state automatons, context-free grammars, context-sensitive grammars, and universal grammars in the increasing order of power (called the Chomsky hierarchy).

---

[9] Yes, all grammar fuzzers we know of are written in Python. Python seems to be the language of choice for grammar fuzzers.
[10] https://github.com/MozillaSecurity/dharma

**Model based Fuzzers.** Generating inputs using grammars were initially explored by Burkhadt [8], and later by Hanford [25] and Purdom [45]. Modern fuzzing tools that use some input models include CSmith [63], LangFuzz [27], Grammarinator [26] (Python), and Domato [20] (Python), Skyfire [58] (Python), and Superion [59], which extends AFL.

**Grammar Learning.** There is a large amount of research on inferring *regular grammars* from blackbox systems [15, 56, 57]. The notable algorithms include L* [1] and RPNI [42].

Blackbox approaches can also be used to learn context-free grammars. Notable approaches include version spaces [53], and GRIDS [33]. GLADE [3] and later REINAM [61] derives the context-free input grammar focusing on blackbox *programs*. Other notable works include the approach by Lin et al. [34] which extracts the AST from programs that parse their input, AUTOGRAM by Höschele et al. [29, 28] which learns the input grammar through active learning using source code of the program, Tupni [14] by Cui et al. which reverse engineers input formats using taint tracking, Prospex [12] from Comparetti et al. which is able to reverse engineer network protocols, and Polyglot [9] by Caballero et al.

Another approach for model learning is through machine learning techniques where the model is not formally represented as a grammar. Pulsar [21] infers a Markov model for representing protocols. Godefroid et al. [23] uses the learned language model to fuzz. IUST-DeepFuzz from Nasrabadi et al. [40] uses infers a neural language model using RNNs from the given corpus of data, which is used for fuzzing.

**Faster execution.** One of the major concerns of fuzzers is the speed of execution — to be effective, a fuzzer needs to generate a plausible input, and execute the program under fuzzing. Given that programs often have instrumentation enabled for tracing coverage, it becomes useful to reduce the overhead due to instrumentation. The *untracer* from Nagy et al. [39] can remove the overhead of tracing from parts of the program already covered, and hence make the overall program execution faster. Another approach by Hsu et al. [30] shows that it is possible to reduce the overhead of instrumentation even further by instrumenting only the parts that are required to differentiate paths.

**Grammar fuzzers.** A number of grammar based fuzzers exist that take in some form of a grammar. The fuzzers such as Gramfuzz[24], Grammarinator [26], Dharma [38], Domato [20], and CSS Fuzz [46] allow context-free grammars to be specified externally. Other fuzzers [13, 60] allow specifying a regular grammar or equivalent as the input grammar. Some [17] also allow constraint languages to specify context sensitive features. Other notable research on

grammar based fuzzers include Nautilus [2], Blend-fuzz [62], and the Godefroid's grammar based white-box fuzzing [22].

**Optimizations in Functional Languages.** We have discussed how the fuzzer can be seen as a limited functional domain specific language (DSL) for interpreting context-free grammars. Supercompilation is not the only method for optimizing functional programs. Other methods include deforestation [55], and partial evaluation et al. [31]. Further details on how partial evaluation, deforestation and supercompilation fit togetherr can be found in Sorensen et al. [50].

**Optimizations in Virtual Machine Interpreters.** A number of dispatch techniques exist for virtual machine interpreters. The most basic one is called switch dispatch in which an interpreter fetches and executes an instruction in a loop [7]. In direct threading, addresses are arranged in an execution thread, and the subroutine follows the execution thread by using computed jump instructions to jump directly to the subroutines rather than iterating over a loop. A problem with the direct threading approach is that it is penalized by the CPU branch predictor as the CPU is unable to predict where a computed jump will transfer control to. An alternative is context threading [5] where simple *call* and *return* instructions are used for transferring control back and forth from subroutines. Since the CPU can predict where a return would transfer control to, after a call, the penalty of branch misprediction is lessened.

## 12 Limitations and Future Work

Despite our improvements in speed, our work has lots of room for improvement, notably in terms of supported language features and algorithmic guidance.

What we have presented is a deliberately simple approach to building grammar based fuzzers. To make the exposition simple, we have chosen to limit the bells and whistles of our fuzzers to a minimum—just a way to control the maximum amount of recursion. This allows us to show how to view the fuzzer first as an interpreter for a programming language, and next as an interpreter for random bitstream. However, this means that we have left unimplemented a number of things that are important for an industry grade fuzzer. For example, the fuzzer does not have a way to use different probabilities for production rule expansions. Nor does it have a way to accept feedback from the program under test—for instance, to guide production towards input features that might increase coverage. Similarly, there is no way for it to actually run the program, or to identify when the program has crashed. Other limitations include the inability to use a pre-existing corpus of sample data, or to infer input models from such data. All these are parts that need to be incorporated to take the fuzzer from an academic exercise to a fuzzer fully equipped to fuzz real programs. We take a look at how some of these features might be implemented next.

## 12.1 Controlling the Fuzzer

```
1  def unroll_key(grammar, key='<start>'):
2      return {tuple(rn) for ro in grammar[key]
3              for rn in unroll_rule(grammar, ro)}

5  def unroll_rule(grammar, rule):
6      rules = [grammar[key] if key in grammar
7                            else [[key]]
8              for key in rule]
9      return [sum(l, []) for l in product(*rules)]

11 def unroll(grammar):
12     return {k:unroll_key(grammar, k)
13             for k in grammar}
```

Figure 11: Unrolling the grammar one level

```
1  {
2  '<start>':  (['<l1>'])
3  '<l1>':     (['<cvalue>'], ['<l2>'])
4  '<l2>':     (['<cvalue>'], ['<l3>'])
5  '<l3>':     (['<cvalue>'], '<svalue>')
6  '<cvalue>': (['<object>'],
7              ['<array>'],
8              ['<string>'],
9              ['<number>']),
10 '<svalue>': (["true"], ["false"], ["null"]),
```

Figure 12: Managing the probability of expansion for simple values such as true, false and null with multi-level embedding. At each level, the probability of selection is halved.

One of the problems with the our simple approach is that some of the grammars contain rules that are unbalanced, when compared to other rules. For example, For example, the top level of a JSON grammar has true, false, and null as top level expansion rules. Since we randomly choose an expansion rule, these get produced in high frequency, which is undesirable.

### 12.2 Managing probability of strings produced

One way to prevent this is to unroll the grammar such that the top level productions with limited expansions are reduced in probability of selection. Figure 11 shows how each rule in a grammar can be unrolled by one level. This can be as many number of times as required to obtain a flatter grammar. Another technique for reducing the probability of choice for these expansions is to embed them in lower levels as shown in Figure 12. Of course, if a fine control over the probability of choice for each expansion is desired, one may modify the fuzzer to accept probability annotated grammar instead.

### 12.3 Fine grained fuzzer control

In Section 8.1 we saw how to pre-allocate random numbers, and use this stream as an input to the fuzzer. This stream of random numbers also provides us with a means of controlling the fuzzer. From Section 9, we have seen how the random numbers are essentially contextual opcodes for the

fuzzing virtual machine. This means that we can specify the path-prefix to be taken by the fuzzer directly. It necessitates relatively minor change in the output generation loop in the body of the main(), and not in the core fuzzer.

### 12.4 Full supercompilation on the virtual machines

We have detailed how supercompilation can be used to improve the fuzzer. However, due to time constraints, and complexity of implementation, we have not implemented the full supercompilation on the direct and context threading virtual machines. This will be part of our future work.

### 12.5 Superoptimization

We note that the assembly program we generated is not optimized as we are not experts in *x86-64* assembly, and likely has further avenues for optimization if we had more expertise in house. However, the situation is not bleak. There is a technique called *superoptimization* [47] that can generate optimized versions of short sequences of assembly instructions that are loop free. If you remember from Section 7.2, we used *supercompiling* to generate long sequences that are loop free that correspond to each opcode we have. Hence, our generated assembly program is particularly well suited to such optimization techniques. This will be taken up for our future work.

## 13  Conclusion

Fuzzing is one of the key tools in a security professional's arsenal, and grammar based fuzzers can deliver the best bang for the buck. Recent developments in grammar inference and faster execution of instrumented programs puts the focus on improving the speed of grammar based fuzzing.

We demonstrated how one can go from a simple grammar fuzzer in Python with a low throughput to a grammar fuzzer that is a few orders of magnitude faster, and near the limit represented by pure random fuzzers. Our fast fuzzer is able to generate a throughput of more than a hundred megabytes per second, producing valid inputs much faster than any of the competitors, and even faster than the simplest fuzzer which simply concatenates random characters.

But the pure speed is not our main contribution. We show that by treating a context-free grammar as a pure functional programming language, one can apply approaches from implementation of the functional programming languages domain to make faster fuzzers. Next, we show that by treating the random stream as a stream of opcodes for a fuzzing virtual machine, we can derive more mileage out of the research on efficiently implementing virtual machine interpreters. All in all, we have hardly exhausted the possibilities in these spaces, and we look forward to great and fast fuzzers in the future.

We are committed to making our research reproducible and reusable. All of the source code and data referred to in this paper is available as self-contained Jupyter Notebooks, allowing to run the *F1* prototype and fully replicate and extend all experiments. Visit our repository at

https://github.com/vrthra/f1

# References

[1] Dana Angluin. "Learning regular sets from queries and counterexamples". In: *Information and computation* 75.2 (1987), pp. 87–106.

[2] Cornelius Aschermann et al. "NAUTILUS: Fishing for Deep Bugs with Grammars". In: *Proceedings of NDSS 2019*. 2019. URL: https://www.ndss-symposium.org/ndss-paper/nautilus-fishing-for-deep-bugs-with-grammars/.

[3] Osbert Bastani et al. "Synthesizing Program Input Grammars". In: *ACM SIGPLAN Conference on Programming Language Design and Implementation*. Barcelona, Spain: ACM, 2017, pp. 95–110. ISBN: 978-1-4503-4988-8.

[4] James R. Bell. "Threaded Code". In: *Communications of the ACM* 16.6 (June 1973), pp. 370–372. ISSN: 0001-0782. DOI: 10.1145/362248.362270. URL: http://doi.acm.org/10.1145/362248.362270.

[5] Marc Berndl et al. "Context threading: A flexible and efficient dispatch technique for virtual machine interpreters". In: *Proceedings of the international symposium on Code generation and optimization*. IEEE Computer Society. 2005, pp. 15–26.

[6] Tim Blazytko et al. "GRIMOIRE: Synthesizing Structure while Fuzzing". In: *28th USENIX Security Symposium (USENIX Security 19)*. Aug. 2019, pp. 1985–2002. URL: https://www.syssec.ruhr-uni-bochum.de/media/emma/veroeffentlichungen/2019/06/03/grimoire.pdf.

[7] Stefan Brunthaler. "Virtual-machine abstraction and optimization techniques". In: *Electronic Notes in Theoretical Computer Science* 253.5 (2009), pp. 3–14.

[8] W. H. Burkhardt. "Generating test programs from syntax". In: *Computing* 2.1 (Mar. 1967), pp. 53–73. ISSN: 1436-5057. DOI: 10.1007/BF02235512.

[9] Juan Caballero et al. "Polyglot: Automatic Extraction of Protocol Message Format Using Dynamic Binary Analysis". In: *ACM Conference on Computer and Communications Security*. Alexandria, Virginia, USA: ACM, 2007, pp. 317–329. ISBN: 978-1-59593-703-2. DOI: 10.1145/1315245.1315286. URL: http://doi.acm.org/10.1145/1315245.1315286.

[10] Cristian Cadar, Daniel Dunbar, Dawson R Engler, et al. "KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs." In: *OSDI*. Vol. 8. 2008, pp. 209–224.

[11] Noam Chomsky. "Three models for the description of language". In: *IRE Transactions on Information Theory* 2 (1956), pp. 113–124. URL: https://chomsky.info/wp-content/uploads/195609-.pdf.

[12] Paolo Milani Comparetti et al. "Prospex: Protocol Specification Extraction". In: *IEEE Symposium on Security and Privacy*. Washington, DC, USA: IEEE Computer Society, 2009, pp. 110–125. ISBN: 978-0-7695-3633-0. DOI: 10.1109/SP.2009.14.

[13] Baojiang Cui et al. "A novel fuzzing method for Zigbee based on finite state machine". In: *International Journal of Distributed Sensor Networks* 10.1 (2014), p. 762891.

[14] Weidong Cui et al. "Tupni: Automatic Reverse Engineering of Input Formats". In: *ACM Conference on Computer and Communications Security*. Alexandria, Virginia, USA: ACM, 2008, pp. 391–402. ISBN: 978-1-59593-810-7.

[15] Colin De la Higuera. *Grammatical inference: learning automata and grammars*. Cambridge University Press, 2010.

[16] Ganesh Devarajan. "Unraveling SCADA protocols: Using sulley fuzzer". In: *Defon 15 Hacking Conference*. 2007.

[17] Kyle Dewey, Jared Roesch, and Ben Hardekopf. "Language fuzzing using constraint logic programming". In: *IEEE/ACM International Conference on Automated Software Engineering*. ACM. 2014, pp. 725–730.

[18] Christopher Domas. "Breaking the x86 ISA". In: *Black Hat*. 2017.

[19] M. Anton Ertl and David Gregg. "The Structure and Performance of Efficient Interpreters". In: *Journal of Instruction-Level Parallelism* 5 (2003).

[20] Ivan Fratric. *Domato A DOM fuzzer*. 2019. URL: https://github.com/googleprojectzero/domato (visited on 07/03/2019).

[21] Hugo Gascon et al. "Pulsar: Stateful black-box fuzzing of proprietary network protocols". In: *International Conference on Security and Privacy in Communication Systems*. Springer. 2015, pp. 330–347.

[22] Patrice Godefroid, Adam Kiezun, and Michael Y. Levin. "Grammar-based Whitebox Fuzzing". In: *ACM SIGPLAN Conference on Programming Language Design and Implementation*. Tucson, AZ, USA: ACM, 2008, pp. 206–215. ISBN: 978-1-59593-860-2.

[23] Patrice Godefroid, Hila Peleg, and Rishabh Singh. "Learn&#38;Fuzz: Machine Learning for Input Fuzzing". In: *Proceedings of the 32Nd IEEE/ACM International Conference on Automated Software Engineering*. ASE 2017. Urbana-Champaign, IL, USA: IEEE Press, 2017, pp. 50–59. ISBN: 978-1-5386-2684-9.

[24] Tao Guo et al. "Gramfuzz: Fuzzing testing of web browsers based on grammar analysis and structural mutation". In: *2013 Second International Conference on Informatics & Applications (ICIA)*. IEEE. 2013, pp. 212–215.

[25] Kenneth V. Hanford. "Automatic Generation of Test Cases". In: *IBM Syst. J.* 9.4 (Dec. 1970), pp. 242–257. ISSN: 0018-8670. DOI: 10.1147/sj.94.0242. URL: http://dx.doi.org/10.1147/sj.94.0242.

[26] Renáta Hodován, Ákos Kiss, and Tibor Gyimóthy. "Grammarinator: a grammar-based open source fuzzer". In: *Proceedings of the 9th ACM SIGSOFT International Workshop on Automating TEST Case Design, Selection, and Evaluation*. ACM. 2018, pp. 45–48.

[27] Christian Holler, Kim Herzig, and Andreas Zeller. "Fuzzing with Code Fragments". In: *Proceedings of the 21st USENIX Conference on Security Symposium*. Security'12. Bellevue, WA: USENIX Association, 2012, pp. 38–38. URL: https://www.usenix.org/system/files/conference/usenixsecurity12/sec12-final73.pdf.

[28] Matthias Höschele, Alexander Kampmann, and Andreas Zeller. "Active Learning of Input Grammars". In: *CoRR* abs/1708.08731 (2017).

[29] Matthias Höschele and Andreas Zeller. "Mining Input Grammars from Dynamic Taints". In: *IEEE/ACM International Conference on Automated Software Engineering*. Singapore, Singapore: ACM, 2016, pp. 720–725.

[30] Chin-Chia Hsu et al. "Instrim: Lightweight instrumentation for coverage-guided fuzzing". In: *Symposium on Network and Distributed System Security (NDSS), Workshop on Binary Analysis Research*. 2018.

[31] Neil D Jones, Carsten K Gomard, and Peter Sestoft. *Partial evaluation and automatic program generation*. Peter Sestoft, 1993.

[32] Kevin Laeufer et al. "RFUZZ: coverage-directed fuzz testing of RTL on FPGAs". In: *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE. 2018, pp. 1–8.

[33] Pat Langley and Sean Stromsten. "Learning context-free grammars with a simplicity bias". In: *European Conference on machine learning*. Springer. 2000, pp. 220–228.

[34] Zhiqiang Lin and Xiangyu Zhang. "Deriving Input Syntactic Structure from Execution". In: *ACM SIGSOFT Symposium on The Foundations of Software Engineering*. Atlanta, Georgia: ACM, 2008, pp. 83–93. ISBN: 978-1-59593-995-1.

[35] Valentin J. M. Manès et al. "Fuzzing: Art, Science, and Engineering". In: *CoRR* abs/1812.00140 (2018). arXiv: 1812.00140. URL: http://arxiv.org/abs/1812.00140.

[36] Lorenzo Martignoni et al. "Testing CPU emulators". In: *ACM SIGSOFT International Symposium on Software Testing and Analysis*. ACM. 2009, pp. 261–272.

[37] Björn Mathis et al. "Parser Directed Fuzzing". In: *ACM SIGPLAN Conference on Programming Language Design and Implementation*. New York, NY, USA: ACM, 2019.

[38] Mozilla. *Dharma: A generation-based, context-free grammar fuzzer*. 2019. URL: https://blog.mozilla.org/security/2015/06/29/dharma/ (visited on 07/03/2019).

[39] Stefan Nagy and Mathew Hicks. "Full-speed Fuzzing: Reducing Fuzzing Overhead through Coverage-guided Tracing". In: *IEEE Symposium on Security and Privacy*. Washington, DC, USA: IEEE Computer Society, 2019.

[40] Morteza Zakeri Nasrabadi, Saeed Parsa, and Akram Kalaee. "Neural Fuzzing: A Neural Approach to Generate Test Data for File Format Fuzzing". In: *CoRR* abs/1812.09961 (2018). arXiv: 1812.09961. URL: http://arxiv.org/abs/1812.09961.

[41] Vegard Nossum and Quentin Casasnovas. "Filesystem fuzzing with american fuzzy lop". In: *Vault*. 2016.

[42] José Oncina and Pedro Garcia. "Inferring regular languages in polynomial updated time". In: *Pattern recognition and image analysis: selected papers from the IVth Spanish Symposium*. World Scientific. 1992, pp. 49–61.

[43] Daksiputra Pānini. *Ashtādhyāyī*. Sanskrit Oral Tradition, 350.

[44] Pekka Pietikäinen et al. "Security testing of web browsers". In: *Communications of Cloud Software* 1.1 (2011).

[45] Paul Purdom. "A sentence generator for testing parsers". English. In: *BIT Numerical Mathematics* 12.3 (1972), pp. 366–375. ISSN: 0006-3835. DOI: 10.1007/BF01932308. URL: http://dx.doi.org/10.1007/BF01932308.

[46] Jesse Ruderman. *Introducing jsfunfuzz*. 2007. URL: http://www.squarefree.com/2007/08/02/introducing-jsfunfuzz/ (visited on 08/02/2007).

[47] Eric Schkufza, Rahul Sharma, and Alex Aiken. "Stochastic Superoptimization". In: *International Conference on Architectural Support for Programming Languages and Operating Systems*. Houston, Texas, USA: ACM, 2013, pp. 305–316. ISBN: 978-1-4503-1870-9. DOI: 10.1145/2451116.2451150.

[48] Kosta Serebryany. "Continuous fuzzing with libfuzzer and addresssanitizer". In: *2016 IEEE Cybersecurity Development*. IEEE. 2016, pp. 157–157.

[49] Yan Shoshitaishvili et al. "Sok:(state of) the art of war: Offensive techniques in binary analysis". In: *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE. 2016, pp. 138–157.

[50] Morten Heine Sørensen, Robert Glück, and Neil D Jones. "Towards unifying partial evaluation, deforestation, supercompilation, and GPC". In: *European Symposium on Programming*. Springer. 1994, pp. 485–500.

[51] Valentin F Turchin. *REFAL-5: programming guide & reference manual*. New England Publications, 1989.

[52] Valentin F. Turchin. "The Concept of a Supercompiler". In: *ACM Transactions on Programming Languages and Systems* 8.3 (June 1986), pp. 292–325. ISSN: 0164-0925. DOI: 10.1145/5956.5957. URL: http://doi.acm.org/10.1145/5956.5957.

[53] Kurt Vanlehn and William Ball. "A Version Space Approach to Learning Context-free Grammars". In: *Machine Learning* 2.1 (Mar. 1987), pp. 39–74. ISSN: 0885-6125. DOI: 10.1023/A:1022812926936.

[54] Spandan Veggalam et al. "Ifuzzer: An evolutionary interpreter fuzzer using genetic programming". In: *European Symposium on Research in Computer Security*. Springer. 2016, pp. 581–601.

[55] Philip Wadler. "Deforestation: Transforming programs to eliminate trees". In: *European Symposium on Programming*. Springer. 1988, pp. 344–358.

[56] Neil Walkinshaw, Kirill Bogdanov, and Ken Johnson. "Evaluation and Comparison of Inferred Regular Grammars". In: *Proceedings of the 9th International Colloquium on Grammatical Inference: Algorithms and Applications*. ICGI '08. Saint-Malo, France: Springer-Verlag, 2008, pp. 252–265. ISBN: 978-3-540-88008-0. DOI: 10.1007/978-3-540-88009-7_20.

[57] Neil Walkinshaw et al. "Reverse Engineering State Machines by Interactive Grammar Inference". In: *Working Conference on Reverse Engineering*. Washington, DC, USA: IEEE Computer Society, 2007, pp. 209–218. ISBN: 0-7695-3034-6. DOI: 10.1109/WCRE.2007.45.

[58] Junjie Wang et al. "Skyfire: Data-driven seed generation for fuzzing". In: *IEEE Symposium on Security and Privacy*. IEEE. 2017, pp. 579–594.

[59] Junjie Wang et al. "Superion: grammar-aware greybox fuzzing". In: *Proceedings of the 41st International Conference on Software Engineering*. IEEE Press. 2019, pp. 724–735.

[60] Ming-Hung Wang et al. "Automatic Test Pattern Generator for Fuzzing Based on Finite State Machine". In: *Security and Communication Networks* 2017 (2017).

[61] Zhengkai Wu et al. "REINAM: Reinforcement Learning for Input-Grammar Inference". In: *ACM SIGSOFT Symposium on The Foundations of Software Engineering*. Tallinn, Estonia: ACM, 2019.

[62] Jingbo Yan, Yuqing Zhang, and Dingning Yang. "Structurized grammar-based fuzz testing for programs with highly structured inputs". In: *Security and Communication Networks* 6.11 (2013), pp. 1319–1330.

[63] Xuejun Yang et al. "Finding and understanding bugs in C compilers". In: *ACM SIGPLAN Notices*. Vol. 46. 6. ACM. 2011, pp. 283–294.

[64] Andreas Zeller et al. "Efficient Grammar Fuzzing". In: *Generating Software Tests*. Retrieved 2019-05-10 13:29:32+02:00. Saarland University, 2019. URL: https://www.fuzzingbook.org/html/GrammarFuzzer.html.

[65] Andreas Zeller et al. "Fuzzing with Grammars". In: *Generating Software Tests*. Retrieved 2019-05-10 13:29:32+02:00. Saarland University, 2019. URL: https://www.fuzzingbook.org/html/GrammarFuzzer.html.

[66] Andreas Zeller et al. "Fuzzing: Breaking Things with Random Inputs". In: *Generating Software Tests*. Retrieved 2019-05-21 19:57:59+02:00. Saarland University, 2019. URL: https://www.fuzzingbook.org/html/Fuzzer.html (visited on 05/21/2019).

[67] Andreas Zeller et al. "Mutation-Based Fuzzing". In: *Generating Software Tests*. Retrieved 2019-05-21 19:57:59+02:00. Saarland University, 2019. URL: https://www.fuzzingbook.org/html/MutationFuzzer.html (visited on 05/21/2019).

[68] Yaowen Zheng et al. "FIRM-AFL: High-Throughput Greybox Fuzzing of IoT Firmware via Augmented Process Emulation". In: *28th USENIX Security Symposium*. 2019, pp. 1099–1114.