# FrameShift: Learning to Resize Fuzzer Inputs Without Breaking Them

Harrison Green
harrisog@cmu.edu
Carnegie Mellon University

Claire Le Goues
clegoues@cmu.edu
Carnegie Mellon University

Fraser Brown
fraserb@cmu.edu
Carnegie Mellon University

## Abstract

Coverage-guided fuzzers are powerful automated bug-finding tools. They mutate program inputs, observe coverage, and save any input that hits an unexplored path for future mutation. Unfortunately, without knowledge of input formats—for example, the relationship between formats' data fields and sizes—fuzzers are prone to generate destructive *frameshift* mutations. These time-wasting mutations yield malformed inputs that are rejected by the target program. To avoid such breaking mutations, this paper proposes a novel, lightweight technique that preserves the structure of inputs during mutation by detecting and using *relation fields*.

Our technique, FrameShift, is simple, fast, and does not require additional instrumentation beyond standard coverage feedback. We implement our technique in two state-of-the-art fuzzers, AFL++ and LibAFL, and perform a 12+ CPU-year fuzzer evaluation, finding that FrameShift improves the performance of the fuzzer in each configuration, sometimes increasing coverage by more than 50%. Furthermore, through a series of case studies, we show that our technique is versatile enough to find important structural relationships in a variety of formats, even generalizing beyond C/C++ targets to both Rust and Python.

## 1 Introduction

Fuzzing is an effective tool for exploring the state space of programs and finding bugs. While the earliest fuzzer simply fed random data into UNIX programs [26] (and, again, found bugs!), coverage-guided fuzzers like AFL++ [11], and LibAFL [12] use mutations and feedback to explore targets more efficiently. These coverage-guided fuzzers maintain a growing corpus of inputs. They pick inputs from this corpus, apply random mutations (bitflips, arithmetic operations, insertions from a dictionary, etc.), and then measure feedback like edge or block coverage. The fuzzers retain mutated inputs that reach new coverage, and discard mutated inputs that don't.

These fuzzers are versatile enough to find bugs in a wide variety of targets, and, as a result, they've been widely adopted in industry. For example, coverage-guided fuzzers have found more than 13,000 vulnerabilities across 1,000+ open-source projects as part of Google's OSS-Fuzz.[1]

Unfortunately, even the best modern fuzzers struggle to successfully mutate certain types of input structures. Many common fuzz targets operate over serialized binary formats whose *metadata*—e.g., `size` and `offset` fields—describes the layout of associated data buffers. Security-critical applications process such structured inputs: Hardware security chips, for example, operate on `TPM` packets, and `openssl` (and others) use `DER`-encoded `ASN.1` messages; both formats contain multiple nested size fields, which make them notoriously hard to mutate [3].

In general, these sorts of fields occur in almost every serialized binary format—in codecs (`PNG`, `JPEG`, `MP3`, `OGG`, etc.), document formats (`PDF`, `XLSX`, `DOCX`, etc.), cryptographic protocols (`TLS`, `SSH`), object formats (`ELF`, `PE`, `Mach-O`), and many more. Fundamentally, any multi-part, variable-length binary data format *requires* metadata to describe its structure and delineate field boundaries.

These metadata-rich formats pose a challenge to modern fuzzers. When a fuzzer mutates specific parts of an input—like a variable sized data buffer—without correspondingly updating related parts of the input—like the size or offset fields describing that buffer—it renders the input structurally invalid. We call such destructive mutations *frameshifts*. As a result of the frameshift, the target program will mishandle the input data or abort early with a validation error. Frameshifts cause fuzzers to get stuck exploring the space of invalid inputs and the space of inputs with the same sized structures as appear in the seed corpus; they are unable to discover new, interesting inputs that contain resized or shifted data.

Existing approaches to this problem either (1) augment the fuzzer with an input specification, allowing it to understand and generate the expected structure [27, 9], or (2) learn important structures automatically during fuzzing using e.g., static analysis [19, 8], coverage-guided feedback [10, 38], or (recently) a combination of static analysis and machine learning [31].

Unfortunately, existing techniques in the first category require manual effort and, simultaneously, risk over-constraining the fuzzer. Techniques in the second category don't solve the general frameshift problem either. Both AIFORE [31] and ProFuzzer [38] are closed source, which makes widespread adoption impossible. Further, all of TIFF [19], WEIZZ [10], ProFuzzer [38], and AIFORE [31] try to discern whether bytes are e.g. an integer, or an enum, or a string, etc. They may identify relation fields, but do not attempt to identify the relationship between these fields and their target buffers—and thus cannot perform validity-preserving resizing mutations.

---

[1]https://github.com/google/oss-fuzz

Further, modern fuzzers include a (rough) generalization of the techniques prior work uses to e.g., discern valid `enum` options, byte ranges, etc. AFL++ and LibAFL, for example, use sophisticated strategies like compare-logging (an adaptation of RedQueen [2]), which attempts to find special values and inject them into the input. This type of technique can, for example, identify alternative `enum` options (by instrumenting switch cases), or required magic bytes (by instrumenting e.g. `memcmp`), thus subsuming less general analyses that do (some of) the same thing. Unfortunately, these features are not always enabled by default, rendering them absent from some academic evaluations [8, 31]—and potentially underselling the actual performance of modern general-purpose fuzzers. In our work, we use industry-standard configurations (following FuzzBench [25]) and focus on designing a system that confers an actual, significant benefit over state-of-the-art baseline fuzzers, even in their optimal configurations.

This paper presents a new approach to fuzzing structured input formats by discovering relation fields and using them for structure-aware resizing mutations that preserve input validity. Our approach, FrameShift, is built on two key insights. First, coverage loss between a seed and a mutated input indicates that the fuzzer *may* have mutated an important relation field (e.g., size field) without mutating the corresponding data. This indicates a *potential* destructive frameshift that can be identified dynamically, over the course of a coverage campaign. There are, however, other reasons—reasons beyond frameshifts—that a mutation can lead to coverage loss. For example, mutations to enums may redirect execution to a different code path, or mutations to checksums may cause the target to abort early. Thus, our second insight is that, to identify *true* frameshifts, FrameShift can conduct experiments to find points where resizing a buffer restores coverage with respect to the original destructive mutation. We prune the search space of potential corrective mutations using domain-specific heuristics that let the analysis run in mere seconds per new input. Finally, FrameShift uses its newly-discovered relations to inform fuzzing with existing mutators—thus doing structure-aware fuzzing that avoids destructive frameshifts.

We designed FrameShift to be fast, easy to integrate, and compatible with modern fuzzers. It can be applied to a wide range of existing fuzzers because, inspired by prior work [10, 38], it doesn't require instrumentation beyond coverage feedback. It also does not require manual format specifications. Finally, FrameShift is performant: on unfriendly targets, the analysis almost never tanks fuzzer performance by incurring serious overhead; on friendly targets, it helps a fuzzer achieve new coverage quickly.

**Contributions.** We show that FrameShift is:

- **Effective** compared to industry-leading, state-of-the-art fuzzers. It increases coverage by an average of 6%—and more than 50% in certain configurations—while only suffering a 5.5% coverage loss (on average) for the worst-case target.

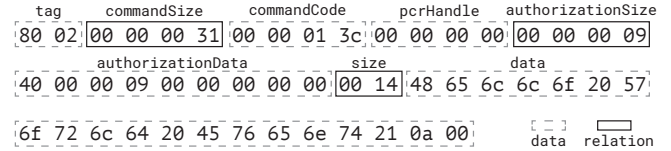- **Versatile**. Unlike static analysis-based approaches, it



Figure 1: An example `TPM_PCR_Event` packet with annotated fields.

is not limited to C/C++: we run FrameShift out-of-the-box on Rust and Python fuzz harnesses.

- **Capable** of identifying real, target-specific size and offset fields—even nested ones—in a variety of popular binary formats. It even discovers semantic differences in input formats across two programs that parse the same input type.

## 2 Overview

This section describes the challenges that fuzzers face with structured inputs (Section 2.1), and then outlines the intuition behind our approach (Section 2.2). We use the `TPM` format as a running example by investigating the `ms-tpm-20-ref` target,[2] a reference TPM 2.0 specification and simulator developed by Microsoft.

### 2.1 Motivating Example

Figure 1 shows a `TPM_PCR_Event` command packet (command code `0x13c`) with the payload "Hello World Event!". The parsing and command execution logic (Figure 2, simplified for presentation) contains a number of interesting edge cases for a fuzzer to discover. It parses fields from the data buffer (lines 4-7); validates authorization data, or follows an alternative path (`CheckAuthNoSession`) if there exists no authorization session (lines 9-15); and then invokes the correct execution handler (lines 17-22). If any of the parser checks fail, execution terminates with error handling (`goto Err`).

The best outcome for a fuzzer for this code snippet involves reaching the call to `TPM2_PCR_Event` (line 21) with many different inputs. To do so, the fuzzer must navigate three nested size fields that must stay synchronized with one another and with the described data: `cmdSize` must match the total packet size (check ❶); `authSize` must fit within the remaining data (check ❷); and payload event `size` (drawn from the payload data, interpreted as a variable size `PCR_Event` buffer) must equal the size of the remaining `data` field (check ❸).

It is therefore *extremely difficult* for a fuzzer to mutate an existing valid TPM message into another valid message with a different payload size. Any mutation that changes the message size must also change `cmdSize`; insertions or deletions in `authData` must also update `authsize`; and attempts to resize the `data` buffer must also update `size`.

To demonstrate the challenge this poses, we ran five state-of-the-art fuzzers (along with our own 2 FrameShift vari-

---

[2]https://github.com/microsoft/ms-tpm-20-ref

```
ExecuteCommand(uint32_t reqSize, char *req)
 1 COMMAND cmd = {
 2   .paramSize = reqSize, .paramBuffer = req
 3 };
 4 if (!Parse_U16(&cmd.tag, &cmd))        goto Err;
 5 if (!Parse_U32(&cmd.cmdSize, &cmd))    goto Err;
 6 if (!Parse_U32(&cmd.cmdCode, &cmd))    goto Err;
❶ 7 if (cmd.cmdSize != reqSize)            goto Err;
 8
 9 if (cmd.tag == TPM_ST_SESSIONS) {
10   if (!Parse_U32(&cmd.authSize, &cmd)) goto Err;
❷11   if (cmd.authSize > cmd.paramSize)    goto Err;
12   if (!ParseSessionBuffer(&cmd))       goto Err;
13 } else {
14   if (!CheckAuthNoSession(&cmd))       goto Err;
15 }
16
17 switch (cmd.cmdCode) {
18   case TPM_CC_PCR_Event:
19     PCR_Event event;
❸20     if (!Parse_Event(&event, &cmd))    goto Err;
21     TPM2_PCR_Event(&event);
22     break;
23 }
```

Figure 2: Example parsing logic for a `TPM_PCR_Event` command in `ms-tpm-20-ref` (rewritten for clarity). Three size validation checks are highlighted.

ants) on the `ms-tpm-20-ref` target using the example TPM packet from Figure 1 as the seed input. Each fuzzer ran for 48 hours for 10 repetitions (see subsection 4.1 for the full experimental setup). We then analyzed the resulting corpus to see how frequently each fuzzer was able to find *newly sized* (i.e. differing in `cmdSize`, `authSize`, or `size` from the seed) inputs that passed each of the highlighted validation checks.

Table 1 shows results. *None* of the state-of-the-art fuzzers were able to find a *single* newly sized input that got to— and passed!—check ❸. All generated inputs that *did* pass that check had an `authSize` of 9 bytes and a data payload `size` of exactly 20 bytes, like the seed. These fuzzers were effectively stuck, unable to successfully perform a resizing mutation even after 20 CPU-days of fuzzing.[3]

Our FRAMESHIFT variants, under the same configuration, were able to find 14 and 8 newly sized `TPM_PCR_Event` commands respectively, unlocking new codepaths in the `TPM2_PCR_Event` handler. Furthermore, they were both able to find an order of magnitude more newly sized inputs reaching the prior checks, discovering many more command types in the process. As a result, FRAMESHIFT variants found an average of 15.3% more coverage than the baseline fuzzers on `ms-tpm-20-ref` in a 48-hour fuzzing campaign (§4.3).

The `TPM` example illustrates a larger issue that affects a wide range of binary formats (`ELF`, `PNG`, `ASN.1`, etc.) that are essential in the modern software stack. Our approach, which we describe next, overcomes this issue by allowing fuzzers to automatically find relation fields in input formats and use them to enable structure-aware resizing mutations.

| Fuzzer | Corpus Size | New Variants ❶ | ❷ | ❸ |
|---|---|---|---|---|
| *State-of-the-art Fuzzers* | | | | |
| AFL++ [11] | 17062 | 41 | 36 | 0 |
| LIBAFL [12] | 23360 | 52 | 53 | 0 |
| AFL [39] | 15694 | 35 | 38 | 0 |
| WEIZZ [10] | 24231 | 16 | 15 | 0 |
| NESTFUZZ [8] | 11593 | 28 | 32 | 0 |
| *With* FRAMESHIFT *(our work)* | | | | |
| AFL++(FS) | 21693 | 359 | 399 | 14 |
| LIBAFL(FS) | 37255 | 194 | 287 | 8 |

Table 1: Newly sized variants passing the highlighted validation checks in `ms-tpm-20-ref` with and without FRAMESHIFT after 48 hours of fuzzing and 10 repetitions.

## 2.2 FrameShift Intuition

Our technique identifies structural metadata—the position and target of each size field in the TPM input packet—and augments a fuzzer to preserve the relationships between that metadata automatically during mutation. The key idea is to (1) identify *potential* relation fields by observing mutations that cause coverage loss, and to (2) validate *true* relation fields by experimenting with new mutations that restore coverage by changing data size. Such "double-mutants" that preserve coverage while resizing the input indicate a likely relation field.

Figure 3 illustrates a walkthrough of FRAMESHIFT's *double-mutation* experiments on our example `TPM_PCR_Event` packet; we discuss each step of the process next.

(A) **Disrupting Coverage.** The first key observation[4] is that it's possible to *indirectly* identify validation checks (without e.g., statically searching for them ahead of time). At runtime, our tool can use loss-of-coverage as evidence of a validation check. For example, starting with a valid seed and mutating the `cmdSize` field will result in a new input that fails to reach some of the originally-covered program bits (those after check ❶ in Figure 2).

In Figure 3, row (A) highlights the bytes for which incrementing the byte's value by `0x20` causes a loss in coverage (all bytes up to and including `size`); none of the bytes in the `data` field are highlighted, because mutating them does not change the execution path. The highlighted bytes are *candidate* relation fields: it is possible that they correspond to metadata in the input format that must be synchronized with data in the input.

(B) **Restoring Coverage.** Given these candidate relations, the loss-of-coverage can be explained by either: (1) a *frameshift* mutation (what we're looking for), or (2) some other validated part of the testcase (uninteresting for our purposes). For example, mutating `cmdCode` may cause the

---

[3]WEIZZ and NESTFUZZ do identify some structures, but are not able to identify and preserve these size fields.

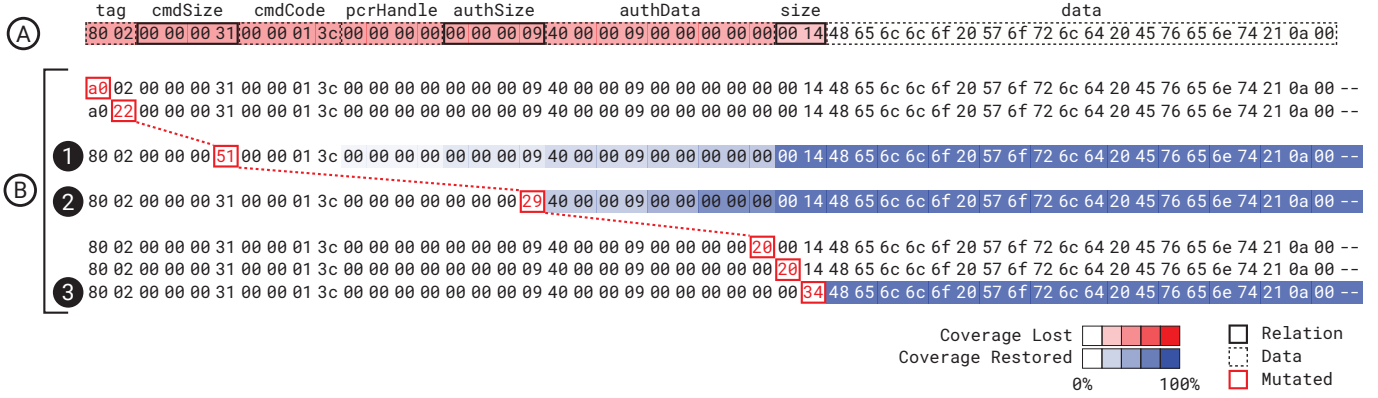[4]Which PROFUZZER [38] also relies on.

Figure 3: Illustration of the double-mutant experiment on a `TPM_PCR_Event` packet. (A): The original input; bytes that cause a loss-of-coverage when mutated are highlighted red. (B): Each row shows a mutated byte that caused a loss-of-coverage (boxed in red) along with every insertion point that was able to restore coverage (highlighted in blue). Some rows with no coverage restoration are omitted for brevity (indicated by the dotted red line). Rows ❶, ❷, and ❸ correspond to the discovery of the highlighted validation checks in Figure 2.

program to trigger a different command handler.

To disambiguate these scenarios, FRAMESHIFT aims to automatically identify points at which bytes can be inserted into the input to restore (most of) the original coverage. That is, if the fuzzer mutated a size field by incrementing its value by $N$, there should be a point in the input where inserting $N$ bytes "re-syncs" the size field with the data it describes. Critically, this is likely only possible if the original mutated field *actually described a size.* It is unlikely, for example, that a fuzzer can increment the `cmdCode` field by $N$, and then insert $N$ bytes elsewhere that restore the original execution behavior.

Although there are 29 bytes that can be mutated to disrupt coverage for the `TPM_PCR_Event` packet, there are only three for which we can find an associated *insertion point* that restores some of the original coverage (Figure 3, rows ❶, ❷, and ❸).[5]

Row (B)❶, shows that inserting right before the `pcrHandle` restores a small percentage of coverage (lightly shaded), but that the most coverage is obtained by inserting near the end of the file. Practically, this is because inserting near the end of the file preserves the authentication section, passing checks ❷ and ❸. Row (B)❷ corresponds to insertions that correct the `authSize` field. Some coverage is restored when inserting inside the `authData` region. However, this is likely to corrupt the authentication data, passing check ❷ and failing ❸. Inserting anywhere after the existing authentication data restores more coverage. Finally, part (B)❸ corresponds to the `size` field of the `PCR_Event`, where any insertion inside the `data` region restores coverage equally.

This example illustrates the intuition behind FRAMESHIFT's approach to dynamically identifying relation fields. Testing every byte or insertion point (as in this illustration) is prohibitive in practice; moreover, discovering certain relation fields is often impossible with-

out discovering others (e.g., `cmdSize`). Our implementation uses heuristics to prune the search space and makes this entire process practical (i.e., on the order of seconds or milliseconds for a single input).

# 3 FrameShift

In this section, we expand on the intuition from the previous section and describe our design in detail. We start by formalizing the concept of a *relation field*, an abstraction over different types of size/offset fields (§3.1). We then describe how to discover these relation fields, using heuristics to prune the search space of possible relations (§3.2). Finally, we discuss how FRAMESHIFT uses relations to implement structure-aware mutations (§ 3.3).

## 3.1 Structured Inputs

Disruptive frameshift mutations occur when a fuzzer modifies input bytes that correspond to metadata (e.g. a size field), without fixing up the corresponding data.

To avoid these frameshifts, our goal is to discover the locations of these size fields and their associated data–the part of the input they describe the size of. For the purposes of mutation, `size` and `offset` fields can be treated the same. A `size` field represents the length of some span of the input data. An `offset` field represents the position of some part of data in relation to the start of the input, or equivalently: a `size` field that describes the length of the input before the data.

We generalize both these forms as a *relation field*

$$R := (a, b, p, s, e)$$

consisting of a field at position $p$ with size $s$ and endianness $e$, that represents the length of some span of the input data. The span is defined by a start position $a$ and an end position $b$, where $a < b$.

---

[5]Initially, the only identifiable byte is `cmdSize`; the others are discoverable only after it is identified.
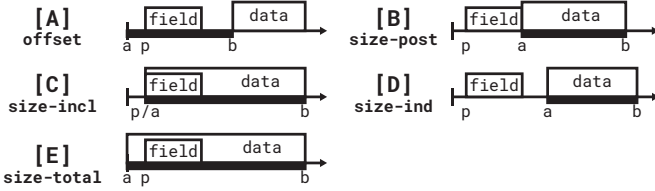
Figure 4: Examples of size/offset fields encoded as relation fields.

This construction generalizes many types of size/offset fields, as depicted in Figure 4. For an offset field (A), the start position $a$ is 0 and the end position $b$ designates where the data starts (i.e., the offset). For a size field, the actual data field starts at $a$ and ends at $b$. However, the size field itself may be positioned in different ways with respect to the data it describes, i.e. immediately before it (B), just inside it (C), some other arbitrary location (D), or the size field may be the entire input (E).

## 3.2 Automatic Relation Discovery

The objective of the relation analysis is to discover important relation fields in particular inputs. This analysis runs once on every new corpus input the fuzzer discovers. When the fuzzer attempts to mutate an input associated with relations, is uses relation data to "fix up" fields whose relationships were affected by the mutation.

The analysis first identifies destructive mutations to the input—mutations that lead to coverage loss. Destructive mutations indicate that the mutated bytes *might* correspond to relation fields. The analysis then attempts *restorative* mutations to undo the destruction caused by the original mutation. If such restoration is possible, it indicates that the first mutation was indeed a frameshift, and that the mutated bytes correspond to a relation field. It also provides evidence as to the nature of the relation (i.e., what part of the input it describes).

**Destructive and restorative mutations.** A given input $I$'s coverage profile $F$ is represented abstractly as a set of bits:

$$F(I) := \{f_1, f_2, \ldots, f_n\}$$

Depending on the instrumentation used, these bits may represent blocks hit, edges taken, etc.

A destructive mutation loses some of this original coverage. We set a threshold $T_{\text{loss}}$ to indicate the percentage of coverage loss necessary to signal a destructive mutation, i.e., $I^-$ is a destructive mutation iff.:

$$|F(I) - F(I^-)| \geq T_{\text{loss}} \cdot |F(I)|$$

Conversely, a *restorative* mutation restores some of the lost coverage, and is defined by the threshold $T_{\text{restore}}$. That is, mutating an input $I^-$ produces a *restoring mutant* $I^+$ iff.:

$$|F(I^+) \cap (F(I) - F(I^-))| \geq T_{\text{restore}} \cdot |F(I) - F(I^-)|$$

In our experiments, $T_{\text{loss}} = 0.05$ and $T_{\text{restore}} = 0.2$. These values work well, but per-target tuning is an interesting future direction.

**Candidate Relation Field Identification.** The first step in the analysis is to identify potential relation fields. To this end, the analysis iterates over every field size $s \in \{8, 4, 2, 1\}$, endianness $e \in \{\text{big}, \text{little}\}$ and potential field position $p \in \{0, \ldots, \text{size}(i) - s\}$. For each configuration, the analysis deserializes the input bytes at position $p$ to obtain a value $v$. To prune the search space, it only considers fields with a value $v \leq \text{size}(i)$ for some input $i$, since sizes or offsets must have values that are at most the size of the input. Given a candidate field, the analysis mutates its value to test if doing so causes coverage loss that exceeds $T_{\text{loss}}$. For $s > 1$, the mutation increments by 0xff, forcing a carry from the least-significant byte (which distinguishes between little- and big-endian fields); for $s = 1$, the mutation increments the value by min(0x20, 0xff-v) such that the increase is large enough to cause potential frameshifts but does not overflow the field. If the mutations lead to coverage loss above $T_{\text{loss}}$, the mutated bytes are a candidate relation, and analysis continues with insertion point discovery.

**Insertion Point Discovery.** Just because a mutation reduces coverage doesn't mean that the mutated bytes store a relation field. For example, mutating important constants, checksums, or changing enum values is likely to invalidate the input or change execution flow, degrading coverage. Thus, the second phase of the analysis seeks evidence that a candidate relation field corresponds to a *true* relation field; if so, it also collects evidence about the nature of the relation.

For each candidate field, the analysis iterates over potential insertion points—places where we can insert bytes—in search of restorative mutations. Iterating over every byte in the input is prohibitively expensive, as it requires invoking the target program for every byte. Instead, we limit the search based on a smaller set of *anchor points*. Specifically, in practice, most relation fields occur in one of the forms depicted in Figure 4. Thus, the start of the target span is often one of 0 (for offset / size-total fields), $p$ (for size-inclusive fields), or $p + s$ (for size-post fields). In each of these cases, we test the corresponding end position (start $+ v$) as a candidate insertion point. We select the insertion point that restores the most coverage, as long as it exceeds the $T_{\text{restore}}$ threshold.

A slightly harder case is size-indirect form (D), where the target start position may occur at an arbitrary point in the input. In practice, since the target program needs to be able to actually locate this data, there is often some other metadata (another size or offset field) that indicates where the start of the data should be.

For example, in ELF files, there are 8-byte size and offset fields describing the location of program and section headers. Therefore, to find these fields efficiently, we expand the insertion point search to consider start positions at $R.p$, $R.a$, and $R.b$ for every true relation field $R$ that we have *already* discovered. During analysis, FRAMESHIFT first identifies the offset field, and then uses it as an anchor point to identify the size field insertion point.

FRAMESHIFT is a fundamentally heuristic analysis: candidate relation fields with viable insertion points are *likely*

5

(but not guaranteed) to be real relation fields. It is possible—but rare!—that inserting bytes can restore coverage for other reasons (e.g., if the initial mutation actually modified an `enum` value, but the inserted bytes coincidentally reintroduced another input structure which hit the original codepath).

## 3.3 Structure-aware Mutation

FRAMESHIFT uses relation fields to augment standard byte-level mutations, yielding *structure-aware* mutations. Next, we describe standard fuzzer mutations, and how FRAMESHIFT adapts these mutations to account for relations.

**Raw Mutations.** Standard fuzzer mutations manipulate a raw input $I$ through three interfaces: $\texttt{Replace}(I, i, V)$ (replace the subsequence starting at $i$ with $V$), $\texttt{Insert}(I, i, V)$ (insert subsequence $V$ at position $i$), and $\texttt{Remove}(I, i, n)$ (remove $n$ bytes after position $i$) (Figure 5a). Fuzzer mutators typically perform several of these actions at once. For example, a *splicing* mutator may perform some combination of `Replace` and `Insert`.

**Structured Mutations.** We redefine these mutation operators to act on a structured input $S := (I, \mathbb{R})$ with input $I$ and set of learned relations $\mathbb{R}$ (Figure 5b), by first applying the mutation to the underlying input $I$ and then invoking `OnInsert` (Algorithm 1) or `OnRemove` (Algorithm 2) to track which relation fields need to shift or update their value as a result of the operation. It is important to perform this bookkeeping as fuzzers may stack these splicing mutations multiple times in a row. After all mutations, and before executing the test case, FRAMESHIFT re-serializes relation fields to apply their new values to the underlying input.

**Accommodating Havoc.** While FRAMESHIFT is designed to identify (and preserve the validity of) resizing mutations, we don't want to inadvertently restrict the fuzzer from making destructive mutations that would unlock new coverage. Therefore, during a mutation, if the fuzzer tries to perform an action that is incompatible with the current set of relations—e.g., inserting into the middle of a relation field itself—FRAMESHIFT temporarily deletes that relation and avoids re-serializing it. Thus, FRAMESHIFT updates relation fields without also over-constraining the fuzzer.[6]

---

**Algorithm 1:** OnInsert

**Data:** Relation $R$, index $i$, sequence $V$
**if** $i \leq R.p$ **then** $R.p \leftarrow R.p + |V|$;
**if** $i < R.a$ **then** $R.a \leftarrow R.a + |V|$;
**if** $i \leq R.b$ **then** $R.b \leftarrow R.b + |V|$;
**return** $R$

---

[6]This is also important in the rare case that FRAMESHIFT incorrectly identifies a relation field.

---

**Algorithm 2:** OnRemove

**Data:** Relation $R$, index $i$, size $n$
**if** $i \leq R.p$ **then** $R.p \leftarrow R.p - \min(R.p - i, n)$;
**if** $i \leq R.a$ **then** $R.a \leftarrow R.a - \min(R.a - i, n)$;
**if** $i \leq R.b$ **then** $R.b \leftarrow R.b - \min(R.b - i, n)$;
**return** $R$

---

## 3.4 Implementation

We implement the FRAMESHIFT algorithm in both AFL++ and LIBAFL, two industry-leading fuzzers. Our implementations integrate with existing fuzzer mutators and require no changes to instrumentation. These variants are denoted as AFL++(FS) and LIBAFL(FS) throughout subsequent sections. Both implementations are open-source under a permissive license, available at `https://github.com/hgarrereyn/AFLplusplus-FrameShift` and `https://github.com/hgarrereyn/LibAFL-FrameShift`.

**FrameShift in AFL++.** Our AFL++ implementation of FRAMESHIFT consists of 600 lines of C that implement a new fuzzer stage to run analysis and store relation metadata in queue inputs. AFL++(FS) tracks insertions and deletions from the havoc and splice mutators, and then reserializes relation data before executing test cases.

**FrameShift in LibAFL.** For LIBAFL we write a modular fuzzer stage and custom input type in roughly 1600 lines of Rust. The implementation is functionally identical to our AFL++ fork but implemented in a canonical LIBAFL style. As such, it is plug-and-play with many other LIBAFL modules. This lets us use LIBAFL's support for other languages, and apply FRAMESHIFT *out of the box* to both Rust and Python targets (see Section 4.5).

# 4 Evaluation

This section answers the following research questions:

- **RQ1 (Performance)**: How does FRAMESHIFT compare to SOTA binary and structure-aware fuzzers at finding coverage? (Section 4.2)

- **RQ2 (Applicability)**: Where is FRAMESHIFT most/least effective? What are the failure cases? (Section 4.3)

- **RQ3 (Case Study)**: What types of structures can FRAMESHIFT identify in real-world targets? (Section 4.4)

- **RQ4 (Versatility)**: How versatile is FRAMESHIFT with respect to different languages and different forms of coverage feedback? (Section 4.5)

Section 4.1 describes benchmarks and baselines; subsequent sections address each research question in turn.

$$\text{Replace}(I, i, V) \to (I_0, I_1, \ldots, I_{i-1}) \lozenge V \lozenge (I_{|V|}, I_{|V|+1}, \ldots, I_{|I|-1})$$

$$\text{Insert}(I, i, V) \to (I_0, I_1, \ldots I_{i-1}) \lozenge V \lozenge (I_i, I_{i+1}, \ldots, I_{|I|-1})$$

$$\text{Remove}(I, i, n) \to (I_0, I_1, \ldots, I_{i-1}) \lozenge (I_{i+n}, I_{i+n+1}, \ldots, I_{|I|-1})$$

(a) Unstructured mutation operators

$$\text{Replace}(S, i, V) \to (\text{Replace}(S.I, i, V), \ S.\mathbb{R})$$

$$\text{Insert}(S, i, V) \to (\text{Insert}(S.I, i, V), \ \{\text{OnInsert}(R, i, V) \mid R \in S.\mathbb{R}\})$$

$$\text{Remove}(S, i, n) \to (\text{Remove}(S.I, i, n), \ \{\text{OnRemove}(R, i, n) \mid R \in S.\mathbb{R}\})$$

(b) Structured mutation operators

Figure 5: Side-by-side comparison of unstructured (left) and structured (right) mutation operators. $\lozenge$ denotes sequence concatenation.

| Benchmark | Format | Commit |
|---|---|---|
| bloaty | ELF/Mach-O/WebAssembly | 52948c1 |
| freetype2 | TTF/OTF/WOFF | cd02d35 |
| harfbuzz | TTF/OTF/TTC | cb47dca |
| lcms | ICC-profile | f0d9632 |
| libjpeg-turbo | JPEG | 3b19db4 |
| libpcap | PCAP | 17ff63e |
| libpng | PNG | cd0ea2a |
| ms-tpm-20-ref | TPM | 6b72d66 |
| openh264 | H.264 | 045aeac |
| openssl | DER | b0593c0 |
| openthread | IPV6-packet | 2550699 |
| qpdf | PDF | 2cb2412 |
| vorbis | OGG | 84c0236 |
| woff2 | WOFF | 8109a2c |
| jsoncpp | JSON (text) | 8190e06 |
| libxml2 | XML (text) | c7260a4 |

Table 2: Benchmark programs

| Type | Ref | Name | Version |
|---|---|---|---|
| Binary | [11] | AFL++ | v4.21c |
| Binary | [12] | LibAFL | f343376 |
| Binary | [39] | AFL | v2.57b |
| Structured | [8] | NestFuzz | d16eb69 |
| Structured | [10] | WEIZZ | c9cbeef |

Table 3: Fuzzers used in our evaluation.

because it requires paid decompiler software and relies on now-outdated versions of Intel Pin; recent results suggest it would be outperformed by both NestFuzz and WEIZZ, which we include. These fuzzer baselines represent the state-of-the-art both in general purpose coverage-guided fuzzing and automated binary structure-aware fuzzing.

**Fuzzer Configurations.** For our prototype and all baseline fuzzers except NestFuzz, we use the FuzzBench configuration. In both LibAFL and AFL++, this configuration includes REDQUEEN-style compare-logging [2] and dictionaries, two features that work well in practice [24]. Our tool variants are configured identically to AFL++ and LibAFL, except they include a new FrameShift fuzzer stage, and the ability to fixup inputs after resizing. Finally, since NestFuzz does not have a FuzzBench configuration, we use the configuration provided in the project README.

**Hardware.** We ran the large-scale fuzzing experiment on Google Cloud C3 instances with Intel Sapphire Rapids processors. The case-studies ran on dedicated servers with two Intel(R) Xeon(R) Gold 6430 @ 3.40GHz and 1 TB of RAM.

**Coverage Measurement.** We evaluate fuzzer performance by running each resulting corpus through a build of each benchmark instrumented with LLVM coverage, computing the total edge coverage.

## 4.1 Experimental Setup

**Benchmarks.** For the large-scale experiment (RQ1), we select 16 benchmarks (Table 2): all the binary-format targets in FuzzBench [25], two text-based formats (`jsoncpp` and `libxml2`) from FuzzBench, and two more binary formats: `ms-tpm-20-ref` and `qpdf` from OSS-Fuzz [30], inspired by discussions about hard-to-fuzz file formats [3]. We include the text-based formats as a baseline for measuring FrameShift's worst-case overhead (since we don't expect our technique to work well on these benchmarks).

**Baseline Fuzzers.** We select five additional baseline fuzzers for our evaluation (Table 3). AFL++ [11] and LibAFL [12] act as direct baselines for our prototype implementations. We also include AFL [39], since it is the direct baseline for NestFuzz.

We also evaluate against two fuzzers that do structural inference. WEIZZ [10] uses coverage feedback and extra comparison instrumentation to identify structures in chunk-based binary formats. NestFuzz [8] models the input processing logic of a program via dynamic taint analysis to discover dependencies which are used during mutations.

We do not include AIFORE [31] or ProFuzzer [38] because both are closed-source.[7] We also omit TIFF [19]

## 4.2 RQ1: Fuzzing Performance

To understand FrameShift's performance compared to state-of-the-art fuzzers, we ran a large-scale fuzzing experiment with 16 benchmarks[8] and 7 fuzzers/fuzzer configurations. We tested fuzz runs from both an empty corpus (denoted E in the results tables; here, ability to learn structure quickly is particularly important) and from a corpus

---

[7]Authors of AIFORE did not respond to our request for source code. Authors of ProFuzzer provided source code, but did not re-

spond to our request for clarification when we could not run the tool.
[8]Applying NestFuzz to `bloaty` and `openssl` requires non-trivial build modifications that we were unable to make.

with a single high-quality seed (denoted S in the results tables; here, ability to find variants of the seed quickly is important).

We ran each fuzzer/benchmark/corpus configuration 10 times for 48 hours. We report the arithmetic mean edge coverage after 48 hours in Table 4.

Additionally, we compute an average score for each fuzzer, following the conventions of FuzzBench. For each target, we compute a fuzzer's score as the percentage of maximum coverage obtained (where maximum coverage is the *highest* coverage obtained by any fuzzer). For example, if a fuzzer achieved the most coverage on a benchmark, it has a score of 100 for that benchmark. If the fuzzer achieves only 70% of the maximum coverage for that benchmark, it receives a score of 70. The mean scores are reported in Table 4 (last row).

**Results.** Table 4 shows results. For both both the empty and seeded corpus configurations, the FrameShift variants were the most effective fuzzers, achieving the highest coverage in 10/16 of the benchmarks in both cases. AFL++(FS) achieved the highest average score (97.3) in the empty corpus by a margin of more than 7 points, followed by LibAFL(FS) (89.6). In the seeded corpus setting, LibAFL(FS) achieved the highest average score (96.7), followed by LibAFL (95.1) and then AFL++(FS) (94.1). This ordering mirrors the baseline fuzzers themselves, where AFL++ performs (relatively) better from an empty corpus, while LibAFL performs better from a seeded corpus.

The other three fuzzers (AFL, WEIZZ, and NestFuzz) were generally not competitive, finding the highest coverage on only three benchmarks across both corpus configurations; in fact, these three fuzzers underperform the baseline fuzzers. This is likely because our evaluation uses state-of-the-art FuzzBench configurations.

## 4.3 RQ2: Applicability

To understand the specific contribution of FrameShift over the baseline fuzzers AFL++ and LibAFL, we visualize the final coverage values for each of the 10 fuzzer runs per benchmark/fuzzer in Table 5. Each graphic shows the FrameShift-enabled variant runs (blue lines above the centerline) along with the baseline fuzzer runs (red lines below the centerline). The lines are plotted on a linear axis where the left-most side represents the run with the least coverage, and the right-most side represents the run with the most coverage. For each configuration, we report the average change in coverage due to enabling FrameShift ($\Delta\%$). Following best practices for fuzzer evaluations [20, 29], we use the Mann-Whitney U-test to compute the statistical significance of differences in fuzzer performance. These results are indicated in the $p$ column, as: * ($p < 0.05$), ** ($p < 0.01$), or *** ($p < 0.001$). Statistically significant results (using $p < 0.05$) are shaded green (FrameShift found more coverage) or red (FrameShift found less coverage) depending on the change in coverage.

**Results.** The direct baseline comparison is broken down more granularly in Table 5. Of the two variants, we find that the AFL++ variant gets more utility from the FrameShift integration, with a statistically significant increase in coverage in 15/32 configurations, and a decrease in only 6 configurations. For LibAFL, the effect is more muted: a statistically significant increase in 6 configurations and a decrease in 4. In both cases, the magnitudes of the coverage increase are generally larger than the decrease.

Aggregating the data, there are 7 benchmarks where FrameShift obtains a statistically significant increase in coverage of at least 3% (`libjpeg`, `lcms`, `bloaty`, `ms-tpm-20-ref`, `woff2`, `libpcap`, `openssl`), 6 benchmarks where it is roughly neutral (`vorbis`, `libpng`, `jsoncpp`, `openh264`, `openthread`, `harfbuzz`), and 3 benchmarks where it has a negative effect (`libxml2`, `qpdf`, and `freetype2`).

All of the 7 positive benchmarks contain serialized size and/or offset fields that FrameShift is able to identify, thus enabling the baseline fuzzer to find differently-sized variants more quickly, contributing to finding coverage more quickly. Generally, there are two cases. 1. FrameShift enables rapid discovery of core coverage: all of the FrameShift runs end up near the highest found coverage, while baseline fuzzer results are more distributed (for example, `woff2`/LibAFL/Empty). Or 2. FrameShift enables breakout coverage discovery: one or more of the FrameShift runs is able to find significantly more coverage due to unlocking a certain codepath (e.g. `bloaty`/AFL++/Seeded).

Of the neutral benchmarks, several include serialized size/offset fields, yet obtain minimal changes in coverage. In `libpng` and `openthread` for example, it appears as though FrameShift is useful in the first few hours of fuzzing, but the baseline variants catch up after 48 hours.

We find an interesting case of *frameshift-resistant* file formats, where FrameShift is not able to identify any fields, and thus fails to be productive. In both `vorbis` and `openh264`, the expected file format contains *sync markers*, explicitly intended to prevent *frameshift* issues when the file is streamed across an unreliable medium. Thus, the parsers can recover when data is improperly resized (for example due to packet loss), and continue parsing mostly unharmed. As a result, FrameShift does not directly observe *frameshifts* in the first part of the *double-mutant* experiment, thus does not identify relation fields.

The case most adversarial to FrameShift is when the target generates an extremely large corpus and/or large files. In `harfbuzz`, `libxml2`, `qpdf`, and `freetype2`, the generated corpora are an order of magnitude larger than other benchmarks (tens of thousands of files), thus FrameShift spends more time analyzing the corpus and less time fuzzing. Even though FrameShift can identify relations in `harfbuzz` and `freetype2`, it is burdened by the analysis overhead. For text formats, the inability to find relation fields does not directly confer a negative performance (as demonstrated by `jsoncpp`), however coupled with the analysis overhead of a large corpus, it may reduce the amount of time available to fuzz (as in `libxml2`).

**Takeaways.** Given these results, a fuzzing practitioner

| Benchmark | AFL++(FS) | | LibAFL(FS) | | AFL++ | | LibAFL | | AFL | | WEIZZ | | NestFuzz | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | E | S | E | S | E | S | E | S | E | S | E | S | E | S |
| bloaty | 2095 | 2468 | **2330** | **4551** | 1858 | 1904 | 2005 | 3825 | 643 | 3157 | 732 | 959 | † | † |
| freetype2 | 8242 | 10142 | 8456 | 10491 | **9441** | **10715** | 8985 | 10317 | 3796 | 7572 | 4301 | 5093 | 3719 | 7284 |
| harfbuzz | 6835 | **7124** | 6661 | 6810 | **7034** | 7046 | 6688 | 6992 | 4020 | 5309 | 3243 | 3878 | 3950 | 5156 |
| lcms | **1940** | **2093** | 1852 | 2070 | 1184 | 1810 | 1753 | 2084 | 1250 | 1203 | 1601 | 1571 | 36 | 551 |
| libjpeg | **1739** | 2344 | 587 | 2288 | 950 | **2367** | 510 | 2318 | 656 | 2299 | 451 | 1948 | 478 | 2316 |
| libpcap | **3041** | **3010** | 2623 | 2846 | 2835 | 2678 | 2468 | 2689 | 36 | 2497 | 2094 | 2095 | 35 | 2068 |
| libpng | **1943** | **1991** | 1798 | 1965 | 1860 | 1963 | 1804 | 1962 | 1549 | 1936 | 1270 | 1690 | 7 | 1191 |
| ms-tpm-20-ref | 2845 | 3180 | **2889** | **3253** | 2209 | 2696 | 2685 | 3047 | 2409 | 2868 | 2072 | 2281 | 2029 | 2223 |
| openh264 | 8498 | 8491 | 8386 | 8305 | 8469 | 8485 | 8414 | 8457 | 8502 | **8521** | 6976 | 7178 | **8509** | 8496 |
| openssl | **4670** | **5178** | 4506 | 4906 | 4007 | 4879 | 4644 | 4690 | 4593 | 4773 | 3654 | 4299 | † | † |
| openthread | 2424 | 2505 | 2614 | **3000** | 2495 | 2521 | **2654** | 2962 | 2287 | 2380 | 2444 | 2655 | 2239 | 2373 |
| qpdf | 1181 | 1999 | 1157 | 1851 | **1203** | **2228** | 1169 | 1886 | 1165 | 1574 | 947 | 969 | 440 | 1045 |
| vorbis | **954** | 1253 | 509 | 1239 | 945 | 1264 | 397 | 1253 | 205 | 1257 | 206 | 1246 | 205 | **1267** |
| woff2 | **939** | **1060** | 934 | 1043 | 769 | 1043 | 814 | 1043 | 7 | 1003 | 713 | 1011 | 7 | 990 |
| jsoncpp (text) | **510** | **510** | 507 | 508 | 509 | **510** | 507 | 507 | 508 | 508 | 508 | 508 | 508 | 157 |
| libxml2 (text) | 12551 | 12392 | 13130 | 13265 | 12838 | **14069** | **13251** | 13200 | 12366 | 10934 | 7442 | 7460 | 7107 | 7289 |
| *Average Score* | **97.3** | 94.1 | 89.6 | **96.7** | 88.6 | 92.1 | 86.6 | 95.1 | 61.7 | 84.4 | 63.9 | 71.9 | 42.4 | 68.8 |

Table 4: Arithmetic mean edge coverage after 48 h (10 runs) for each fuzzer–benchmark pair. Highest average coverage for each benchmark is in bold (for both empty and seeded corpus). E: empty corpus, S: seeded corpus. †: target failed to build.

could likely benefit from enabling FRAMESHIFT in most binary formats, especially when there is no available seed corpus. Incompatible formats, where FRAMESHIFT does not find relation fields will likely not directly confer a negative effect, unless the size of the corpus also grows too quickly. However, for industry-length fuzz campaigns (on the order of weeks), the effect of analysis overhead would be reduced as the corpus begins to saturate. For shorter fuzz campaigns, an interesting future direction could be to detect quick corpus growth and selectively apply FRAMESHIFT, balancing the fuzz time vs. analysis time.

## 4.4 RQ3: Structure Recovery

In this section we qualitatively demonstrate case study examples that demonstrate FRAMESHIFT's ability to identify relation fields in real-world formats.

### 4.4.1 PNG

In Figure 6 we show the fields FRAMESHIFT finds in a PNG file using `libpng`. FRAMESHIFT took only 313 milliseconds to analyze this file and invoked the target 4705 times during the search.

The PNG file consists of an 8-byte PNG header followed by chunks of data. Each chunk has a 4-byte size, 4-byte header (pink), N-byte data (blue), and 4-byte checksum (gray). FRAMESHIFT correctly identifies 9 relation fields in the input (solid black outline). Each of the relation fields it identifies has the correct target span (the subsequent data portion of the file).

FRAMESHIFT did not identify three potentially-expected relation fields, namely the sizes for the initial `IHDR` ❶, `cHRM` ❷, and final `IEND` ❸ chunks. While these locations are

identified as plausible candidates (and may be labeled a size field by other approaches, such as those relying on manual grammars), FRAMESHIFT eliminates them during the insertion point discovery phase because it is unable to generate a *restorative mutant* after changing their size. Analyzing the code, we find that these chunks are validated to be a fixed size during processing. Any size other than 13 (`0xd`) for the `IHDR` chunk will cause `libpng` to abort early. In practice, `IHDR` is not resizable, and therefore FRAMESHIFT does not learn a relation for it. The same holds for the `cHRM` and `IEND` chunks. That is, although superficially similar to the other size fields in this input, these fields do not actually describe parts of the input which can be resized.

### 4.4.2 ASN.1

Figure 7 shows the relation fields FRAMESHIFT finds in a DER-encoded ASN.1 file when run against the `asn1` tool in `openssl`. This input took 59 milliseconds to analyze and required 63 invocations of the target.

The file represents a nested `SEQUENCE` object containing three entries: an octet string with values [9,9,9], a bitstring with values [1,2,3,4], and a printable string with the value `"fuzzer"`. Objects in the file are encoded in a TLV (type-length-value) format. Each object contains a single byte type, followed by a multi-byte length, and then an N-sized value.

FRAMESHIFT identifies the outer sequence length (at `0x01`) and all three of the inner object lengths (at `0x03`, `0x10`, and `0x21`). In this format, length fields are variable-size: values smaller than 127 are encoded in short form, as a single byte. However, larger values use a prefix byte `0x80 + x` where $x$ indicates how many additional bytes are in the encoding for

| Benchmark | AFL++(FS) vs. AFL++ | | | | | | LibAFL(FS) vs. LibAFL | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Empty | Δ% | p | Seeded | Δ% | p | Empty | Δ% | p | Seeded | Δ% | p |
| libjpeg | | +83.0 | * | | -0.9 | | | +15.1 | | | -1.3 | |
| lcms | | +63.9 | ** | | +15.6 | * | | +5.7 | | | -0.7 | |
| bloaty | | +12.7 | ** | | +29.6 | * | | +16.2 | | | +19.0 | |
| ms-tpm-20-ref | | +28.8 | *** | | +18.0 | *** | | +7.6 | ** | | +6.7 | *** |
| woff2 | | +22.1 | *** | | +1.6 | ** | | +14.8 | *** | | +0.1 | |
| libpcap | | +7.3 | ** | | +12.4 | *** | | +6.3 | | | +5.8 | ** |
| openssl | | +16.5 | *** | | +6.1 | ** | | -3.0 | ** | | +4.6 | ** |
| vorbis | | +1.0 | | | -0.9 | *** | | +28.2 | | | -1.1 | *** |
| libpng | | +4.4 | | | +1.4 | ** | | -0.3 | | | +0.1 | |
| jsoncpp | | +0.0 | | | +0.0 | | | +0.0 | | | +0.2 | * |
| openh264 | | +0.3 | * | | +0.1 | | | -0.3 | | | -1.8 | ** |
| openthread | | -2.8 | | | -0.6 | | | -1.5 | | | +1.3 | |
| harfbuzz | | -2.8 | * | | +1.1 | | | -0.4 | | | -2.6 | ** |
| libxml2 | | -2.2 | | | -11.9 | *** | | -0.9 | | | +0.5 | |
| qpdf | | -1.9 | | | -10.3 | ** | | -1.0 | | | -1.8 | |
| freetype2 | | -12.7 | ** | | -5.3 | ** | | -5.9 | | | +1.7 | |

Table 5: FRAMESHIFT coverage compared to baseline fuzzers. Each graphic shows the final coverage of 10 FRAMESHIFT runs compared to 10 baseline fuzzer runs after 48 hours. FRAMESHIFT runs are represented by blue bars on top of the centerline, baseline fuzzer runs are red bars beneath the centerline. The left side of the scale represents the lowest coverage obtained by any run, the right side represents the most, scaled linearly. The average coverage change with FRAMESHIFT is shown in the Δ% column. $p$ contains the Mann-Whitney U test p-value for statistical significance: *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$

```
      0  1  2  3  4  5  6  7  8  9  A  B  C  D  E  F
 0   30 26 04 0B 30 09 02 01 09 02 01 09 02 01 09 03
10   0F 00 30 0C 02 01 01 02 01 02 02 01 03 02 01 04
20   13 06 66 75 7A 7A 65 72
```

□ relation (found)   ▨ type   ▨ data (bitstr)
⌐⌐ relation (missed)   ▨ data (octstr)   ▨ data (prnstr)

Figure 7: Relation fields identified by FrameShift in a DER-encoded ASN.1 file file using `openssl`.

```
      0  1  2  3  4  5  6  7  8  9  A  B  C  D  E  F
 0   89 50 4E 47 0D 0A 1A 0A 00 00 00 0D 49 48 44 52
10   00 00 00 4E 00 00 00 54 08 00 00 00 00 09 8C 5E
20   3C 00 00 00 04 67 41 4D 41 00 00 B1 8F 0B FC 61
30   05 00 00 00 01 73 52 47 42 00 AE CE 1C E9 00 00
40   00 20 63 48 52 4D 00 00 7A 26 00 00 80 84 00 00
50   FA 00 00 00 80 E8 00 00 75 30 00 00 EA 60 00 00
60   3A 98 00 00 17 70 9C BA 51 3C 00 00 00 02 62 4B
70   47 44 00 FF 87 8F CC BF 00 00 00 09 70 48 59 73
80   00 00 00 48 00 00 00 48 00 46 C9 6B 3E 00 00 03
90   B3 49 44 41 54 58 C3 9D D8 4B 48 54 61 14 07 F0
A0   EF CE F8 98 49 CC B1 34 27 2D 7B D0 93 48 A1 72
...                    --- trim ---
440  AD 6F 5A 0E 8D 35 F1 92 C4 36 70 F8 B2 56 C3 63
450  AD 70 F5 1F E7 E7 87 C8 6B 93 D4 FD 00 00 00 25
460  74 45 58 74 64 61 74 65 3A 63 72 65 61 74 65 00
470  32 30 31 36 2D 31 30 2D 31 37 54 32 32 3A 32 33
480  3A 35 38 2D 30 37 3A 30 30 15 58 E5 BC 00 00 00
490  25 74 45 58 74 64 61 74 65 3A 6D 6F 64 69 66 79
4A0  00 32 30 31 36 2D 31 30 2D 31 37 54 32 32 3A 32
4B0  33 3A 35 38 2D 30 37 3A 30 30 64 05 5D 00 00 00
4C0  00 07 74 45 58 74 6C 61 62 65 6C 00 58 F8 B6 12
4D0  8D 00 00 00 13 74 45 58 74 6C 61 62 65 6C 3A 70
4E0  6F 69 6E 74 73 69 7A 65 00 31 31 36 B5 97 E0 AD
4F0  00 00 00 00 49 45 4E 44 AE 42 60 82
```

□ relation (found)   ▨ PNG header   ▨ checksum
⌐⌐ relation (missed)   ▨ type   ▨ data
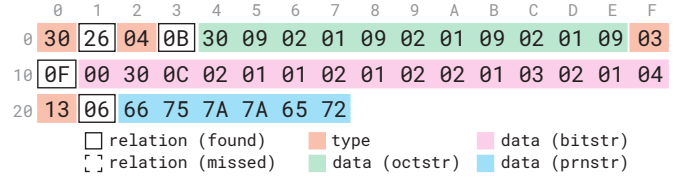
Figure 6: Relation fields identified by FrameShift in a PNG file using `libpng`.

the length field.

FrameShift does not naturally support this variable-size length field construction, yet is still able to approximate the relation field by identifying only the low-order bytes. Specifically, for single-byte values, the field *appears* to be a single-byte value until the size boundary is crossed. Similarly, in larger cases, FrameShift interprets the encoded field `[82 03 12]` (value: 0x123) as a 2-byte big endian integer—not 3!—covering only `[03 12]`. As long as mutations don't cause this size to require one more or fewer byte to represent, it serves as an accurate representation of the real field, thus most fuzzer mutations result in an accurate adjustment. In practice, even though FrameShift is imprecise, it still enables a statistically significant increase in coverage on `openssl` by more than 6%.

### 4.4.3 ELF

In Figure 8, we visualize the relation fields FrameShift identified in an ELF file with the `bloaty` target, which parses several types of object files. This file took 458 milliseconds to analyze and FrameShift invoked the target 1571 times during the search.

The file consists of an ELF header (grey, `0-0x40`), a program table header with one entry (orange, `0x40-0x78`), a section table header with three entries (pink, `0x78-0x138`), and some data that is mapped into the `.text` (green, `0x138-0x16E`) and the `.shstrtab` (blue, `0x16E-0x17F`) sections.

In the ELF header, FrameShift correctly identifies the offset fields ❶ for both the program header table (at `0x20`) and the section header table (at `0x28`). Note that FrameShift must actually discover the section header table *first* because insertions before the program header require shifting both offsets.

The program header table contains a single entry, which in turn contains fields ❷ representing the size on disk (at `0x60`) and in memory (at `0x68`) of a segment to map, along with the file offset (at `0x48`). FrameShift identifies the file size but not the offset (because the ELF header is required to start at the beginning) or the size in memory, because `bloaty` only parses the file and does not attempt to execute it.

The first section ❸ is a necessary `null` section, and thus while there are size/offset fields according to the spec, they are not used and FrameShift does not identify them. The next section is the `.text` section ❹. This section describes both the offset (at `0xD0`) and length (at `0xD8`) of code data to map. Here the section actually maps the entire file starting from the beginning through the green region, which contains

machine code. FRAMESHIFT does not identify the offset field because moving the ELF header caused corruption, but it did identify the size field with the correct parameters.

The third section ❺ describes the `.shstrtab` section which points to a list of strings that are used to identify the section names. It also contains both an offset (at `0x110`) and size (at `0x118`). Here FRAMESHIFT identifies both the offset and the indirect size field.

**Bonus: Automatic Rebase.** While analyzing this ELF example, we were pleasantly surprised to find that FRAMESHIFT identified fields accurately enough able to perform a non-trivial rebase operation. In general, *resizing* parts of an ELF file is difficult, as it requires changing considerable metadata, and typically is only accomplished by dedicated tools, such as `lief` [35].

However, FRAMESHIFT is able to discover enough information (in less than half a second) to accurately fix up the metadata for mutations that splice the contents of the code region (green). Normally, splicing here would corrupt the binary considerably (as several fields are out of sync); with FRAMESHIFT, however, we performed a splice mutation to insert new machine code into this region, shrinking its size. FRAMESHIFT automatically updated one offset field (at `0x110`) and two size fields (at `0xD0` and `0x60`), preserving the validity of the file. The resulting file correctly ran as an executable without any additional modifications!

While this is a small example, it serves to demonstrate that FRAMESHIFT is capable of automatically identifying the important relations in structures and performing potentially-complex size-changing mutations while preserving validity. Thus, in practice, FRAMESHIFT is quite effective on the `bloaty` target, achieving a statistically significant increase in coverage of more than 20% over the baseline fuzzers.

## 4.5 RQ4: Versatility

FRAMESHIFT makes few assumptions about the underlying framework or type of coverage feedback (unlike systems like NESTFUZZ which depend on specific toolchains). Additionally, its modular integration into existing fuzzers allows it to naturally extend baseline capabilities, such as LIBAFL's support for targets besides C/C++. We therefore demonstrate FRAMESHIFT's versatility by applying it to two case study target programs in other languages: Rust (Section 4.5.1) and Python (Section 4.5.2). For these experiments we utilize our LIBAFL prototype of FRAMESHIFT.

We chose to apply FRAMESHIFT to two programs which parse similar file formats to some of our C/C++ benchmark programs (`PNG` in Rust, and `ASN.1` in Python) to compare the similarity between the types of structures FRAMESHIFT can identify and illustrate how even programs which parse the same format can have semantic differences.



Figure 8: Relation fields identified by FRAMESHIFT in an ELF file using `bloaty`.

#### 4.5.1 FrameShift in Rust: PNG

Our first target was the Rust crate `image-png`,[9] a pure Rust library for image decoding and encoding. The project has several `cargo-fuzz` harnesses, including a `decode` harness that decodes arbitrary bytes as a PNG file. We compile the target with sanitizer coverage (as we do for C/C++), since the Rust build system uses LLVM internally.

We ran the Figure 6 test PNG to determine which structures FrameShift could learn in this Rust program. It took 1.7 seconds to run and required invoking the target 9933 times. Note that a structure-inference tool performing any amount of static analysis would likely incur overhead from moving to Rust, since the compiler produces boilerplate that can obfuscate program logic. Since FrameShift is fully dynamic, any additional overhead comes from runtime performance, which is often negligible if the target is compiled with optimizations.

FrameShift found several different fields when fuzzing the Rust (`image-png`) and C (`libpng`) image libraries. Interestingly, these differences are associated with true behavioral differences in the target programs. In the Rust program, FrameShift ignores the `IHDR` chunk size (as in `libpng`), since this is parsed as a fixed-size chunk. It finds the next two length fields (at `0x21` and `0x31`), but then it *also* identifies the `cHRM` size field ❷, unlike with `libpng`. This is because the semantics of PNG actually differ between `libpng` and `image-png`: while `libpng` immediately validates the size of `cHRM` to be 32 bytes, `image-png` has no such check. Thus, in `image-png` it *is* resizable, and FrameShift correctly identifies this.

FrameShift identifies the next two size fields (at `0x6A` and `0x78`), but not the size field of the big `IDAT` chunk at `0x8D`, nor any of the following chunk sizes. Although initially confused, we discovered that `image-png`'s `decode` harness does *not* actually parse the whole PNG file. The fuzz harness instead keeps iterating over chunks until it has found all the actual image data (stored in `IDAT` chunks). In this case, the file contains just a single `IDAT` chunk. Thus, `image-png` scans through the file until reading this chunk and then exits early (ignoring the remaining comment chunks).

From a PNG grammar perspective, these comment chunks have size fields. Yet in the context of this harness, they are ignored and can be freely mutated without inducing a loss in coverage (in fact they do not contribute coverage to begin with); thus FrameShift (rightly) does not consider them to be relation fields.

#### 4.5.2 FrameShift in Python: ASN.1

We also evaluated FrameShift on a Python case study, which uses a different compiler toolchain and alternative coverage feedback. Such an adaptation would be fundamentally quite difficult for a static analysis-based approach without significant effort. FrameShift supports this out-of-the-box, using LibAFL's existing support for Atheris [14], a coverage-guided Python fuzzer. Python represents a higher abstraction level than C/C++, e.g., list concatenation in

a C++ program may hit many basic blocks (i.e. invoking stdlib functions), while the same operation corresponds a single Python opcode in Python. This changes the nature of coverage instrumentation.

Our target program was `pyasn1`,[10] a Python-based framework for encoding and decoding ASN.1 files. We used the same ASN.1 file from Figure 7 as a seed. Running the analysis took only 25 milliseconds and 125 invocations of the target program. FrameShift found the same fields as it did on the `openssl` benchmark. In this case, however, it incorrectly identified the byte at position `0x20` as a size field. This is because, when FrameShift corrupted this field and performed a nearby insertion, it ended up coincidentally achieving enough of the same coverage features (through some other mechanism); therefore, FrameShift categorized the field as a relation. Increasing the $t_{loss}$ threshold was sufficient to remove this false positive, which suggests that FrameShift's parameters may benefit from tuning when applying it to frameworks with different types of coverage feedback.

## 5 Discussion

Our evaluation demonstrates that relation-field structure inference is both possible to perform using only standard greybox coverage and *practical* enough that it can be incorporated into state-of-the-art fuzzers for a net improvement in coverage. FrameShift is capable of identifying relation structures in many real world formats and its relaxed assumptions about the type of coverage and lack of additional instrumentation allow it to be integrated easily into other applications, such as fuzzing Rust and Python programs.

FrameShift serves an example of a potentially more general class of techniques which use coverage feedback in conjunction with heuristics to search the space of possible input structures and augment fuzzers. This approach need not be limited just to size and offset fields. For example, a similar double-mutant experiment may be able to identify compressed/encoded regions of the input (e.g. `zlib`, `base64`, etc.) or repeatable parts of the input (i.e. chunks). A key challenge, as we have seen first-hand in prototyping these kinds of ideas, is keeping the time required for analysis sufficiently tractable to produce wins in the resulting structure-aware fuzzing.

Another interesting future direction, however, could augment this type of coverage-guided inference with a more powerful heuristic, such as a large language model, which could be used to propose high-quality candidate structures. It is possible that such a technique could discover much more complex (and more useful) structures, that could be validated with dynamic experimentation; the additional complexity might effectively counterbalance additional required analysis time.

Additionally, it may be possible to automatically tune certain parameters of FrameShift based on the target program. For example, limiting analysis time if the corpus

---

[9]https://github.com/image-rs/image-png

[10]https://github.com/etingof/pyasn1

grows too quickly, or tuning the $t_{loss}$ threshold for different types of coverage feedback.

# 6 Related Work

**Structure inference for binary fuzzing.** Most similar to FRAMESHIFT are tools which performs automatic structure-inference to aid fuzzers. NESTFUZZ [8] performs a dynamic taint analysis (DTA) to learn the input processing logic for a program. TIFF [19] uses DTA to infer the types of various input fields. Most recently, AIFORE [31] fused byte-level taint analysis with machine learning for byte-level clustering.

WEIZZ [10] and PROFUZZER [38] take a more greybox approach (like us) using just instrumentation to guide inference. While PROFUZZER uses existing coverage feedback, WEIZZ requires not just the final coverage, but the *order* of hitting certain bits. While PROFUZZER performs a similar type of coverage-destroying analysis (like the first step in our double-mutant experiment), it does not attempt to restore coverage through a second mutation, labeling the byte a "size field" without understanding the relation to other fields. Neither is capable of learning both size/offset fields *and* understanding what parts of the input they describe.

**Structure inference for reverse-engineering.** A parallel body of work recovers *internal* structures to aid decompilation or static analysis [21, 7, 6, 22, 32, 23]. While conceptually similar, these approaches target the structures used internally to the program not the serialized structure of the input. Thus, their results are not immediately useful for fuzzer mutations.

**Specification-based fuzzing.** An alternative approach to fuzzing binary formats provides the fuzzer with a specification beforehand. AFLSMART [27] can perform smart chunk-based mutations when provided with the virtual structure of a file format. These mutations are similar in essence to the types of mutations FRAMESHIFT can enable (in that they set size fields accurately), relying on the existence of a manual specification. FORMATFUZZER [9] repurposes structure format files used by a file structure explorer utility, converting them into C++ programs which can *generate* and *mutate* instances of the format. Similarly, the ISLA [34] project aims to create an input specification language that can be sampled using a constraint solver. These approaches are interesting and useful when such a format is available. However, they can be onerous to provide. Beyond the initial complexity of the task, programs may parse multiple formats at once (i.e. `bloaty`), or implement the semantics of a given format *differently* than another program that nominally does the same thing: while FRAMESHIFT can learn program-specific formats (as it did with `PNG` in `libpng` vs. `image-png`), these specification-driven fuzzers cannot.

**Grammar-based fuzzing.** Structure-aware fuzzing has been more heavily utilized in the context of fuzzing text-based formats such as scripting language interpreters [28, 17, 36, 16, 1, 40, 13]. In contrast to structured binary formats, these text-based formats are usually representable with (or can be approximated by) a context-free grammar (CFG). These formats typically do not contain serialized size or offset fields, and thus do not suffer from the same *frameshift* problem as binary formats.

Tools like GLADE [4], PYGMALION [15], SKYFIRE [37], and AUTOGRAM [18] try to automatically learn such input grammars. Given existing grammars, fuzzers like NAUTILUS [1] and GRAMATRON [33] can perform *coverage-guided* semantic-preserving mutations. Perhaps the most related work in this other domain is GRIMOIRE [5] which upon receiving a new input tries to understand how to *generalize* the input by observing how different mutations change the coverage bitmap.

# 7 Conclusion

Destructive frameshift mutations remain a central obstacle to effective coverage-guided fuzzing of binary formats, inhibiting fuzzers from exploring the space of valid inputs. FRAMESHIFT mitigates this problem by learning size- and offset-relations directly from standard coverage feedback and by preserving those relations during mutation. The approach requires no manual specification and integrates transparently with AFL++ and LIBAFL.

In a 12 CPU-year evaluation, across 16 real-world benchmarks, we show that FRAMESHIFT raises edge coverage by an average of 6 percent–and sometimes more than 50 percent–outperforming five other state-of-the-art fuzzers. Further, the approach is language-agnostic and we successfully applied it (with no modifications!) to both Rust and Python.

FRAMESHIFT thus offers a practical solution to combat the *frameshift* problem–complementing the existing performance of modern fuzzers while allowing them to learn to resize inputs without breaking them.

# References

[1] Cornelius Aschermann, Tommaso Frassetto, Thorsten Holz, Patrick Jauernig, Ahmad-Reza Sadeghi, and Daniel Teuchert. Nautilus: Fishing for deep bugs with grammars. In *NDSS*, 2019.

[2] Cornelius Aschermann, Sergej Schumilo, Tim Blazytko, Robert Gawlik, and Thorsten Holz. Redqueen: Fuzzing with input-to-state correspondence. In *NDSS*, volume 19, pages 1–15, 2019.

[3] Cornelius (eqv / .eqv) Aschermann. "something that I've seen commonly that made fuzzing VERY difficult is TLV kind of formats, where the header has an overall size field , and each chunk has a size field, and we have sizeof(header)+header.size + sum(chunks.size) == file_size and sum(chunks.size) == header.size" — Discord message in *#general* (Awesome Fuzzing server), September 2024. Message ID 1287160766469767199; posted 2024-09-21 17:17; accessed 2025-06-05.

[4] Osbert Bastani, Rahul Sharma, Alex Aiken, and Percy Liang. Synthesizing program input grammars. *ACM SIGPLAN Notices*, 52(6):95–110, 2017.

[5] Tim Blazytko, Matt Bishop, Cornelius Aschermann, Justin Cappos, Moritz Schlögel, Nadia Korshun, Ali Abbasi, Marco Schweighauser, Sebastian Schinzel, Sergej Schumilo, et al. {GRIMOIRE}: Synthesizing structure while fuzzing. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 1985–2002, 2019.

[6] Weidong Cui, Jayanthkumar Kannan, and Helen J Wang. Discoverer: Automatic protocol reverse engineering from network traces. In *USENIX Security Symposium*, pages 1–14. Boston, MA, USA, 2007.

[7] Weidong Cui, Marcus Peinado, Karl Chen, Helen J Wang, and Luis Irun-Briz. Tupni: Automatic reverse engineering of input formats. In *Proceedings of the 15th ACM conference on Computer and communications security*, pages 391–402, 2008.

[8] Peng Deng, Zhemin Yang, Lei Zhang, Guangliang Yang, Wenzheng Hong, Yuan Zhang, and Min Yang. Nestfuzz: Enhancing fuzzing with comprehensive understanding of input processing logic. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 1272–1286, 2023.

[9] Rafael Dutra, Rahul Gopinath, and Andreas Zeller. Formatfuzzer: Effective fuzzing of binary file formats. *ACM Transactions on Software Engineering and Methodology*, 33(2):1–29, 2023.

[10] Andrea Fioraldi, Daniele Cono D'Elia, and Emilio Coppa. Weizz: Automatic grey-box fuzzing for structured binary formats. In *Proceedings of the 29th ACM SIGSOFT international symposium on software testing and analysis*, pages 1–13, 2020.

[11] Andrea Fioraldi, Dominik Maier, Heiko Eißfeldt, and Marc Heuse. {AFL++}: Combining incremental steps of fuzzing research. In *14th USENIX Workshop on Offensive Technologies (WOOT 20)*, 2020.

[12] Andrea Fioraldi, Dominik Maier, Dongjia Zhang, and Davide Balzarotti. LibAFL: A Framework to Build Modular and Reusable Fuzzers. In *Proceedings of the 29th ACM conference on Computer and communications security (CCS)*, CCS '22. ACM, November 2022.

[13] Patrice Godefroid, Adam Kiezun, and Michael Y Levin. Grammar-based whitebox fuzzing. In *Proceedings of the 29th ACM SIGPLAN conference on programming language design and implementation*, pages 206–215, 2008.

[14] Google. Atheris: A coverage-guided, native python fuzzer. https://github.com/google/atheris, 2023. Version 2.3.0, commit cbf4ad9, accessed 2025-06-06.

[15] Rahul Gopinath, Björn Mathis, Mathias Höschele, Alexander Kampmann, and Andreas Zeller. Sample-free learning of input grammars for comprehensive software fuzzing. *arXiv preprint arXiv:1810.08289*, 2018.

[16] HyungSeok Han, DongHyeon Oh, and Sang Kil Cha. Codealchemist: Semantics-aware code generation to find vulnerabilities in javascript engines. In *NDSS*, 2019.

[17] Christian Holler, Kim Herzig, and Andreas Zeller. Fuzzing with code fragments. In *21st {USENIX} Security Symposium ({USENIX} Security 12)*, pages 445–458, 2012.

[18] Matthias Hoschele and Andreas Zeller. Mining input grammars with autogram. In *2017 IEEE/ACM 39th International Conference on Software Engineering Companion (ICSE-C)*, pages 31–34. IEEE, 2017.

[19] Vivek Jain, Sanjay Rawat, Cristiano Giuffrida, and Herbert Bos. Tiff: using input type inference to improve fuzzing. In *Proceedings of the 34th annual computer security applications conference*, pages 505–517, 2018.

[20] George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. Evaluating fuzz testing. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, pages 2123–2138, 2018.

[21] JongHyup Lee, Thanassis Avgerinos, and David Brumley. Tie: Principled reverse engineering of types in binary programs. 2011.

[22] Zhiqiang Lin and Xiangyu Zhang. Deriving input syntactic structure from execution. In *Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of software engineering*, pages 83–93, 2008.

[23] Zhiqiang Lin, Xiangyu Zhang, and Dongyan Xu. Automatic reverse engineering of data structures from binary execution. In *Proceedings of the 11th Annual Information Security Symposium*, pages 1–1, 2010.

[24] Dongge Liu, Jonathan Metzman, Marcel Böhme, Oliver Chang, and Abhishek Arya. Sbft tool competition 2023-fuzzing track. In *2023 IEEE/ACM International Workshop on Search-Based and Fuzz Testing (SBFT)*, pages 51–54. IEEE, 2023.

[25] Jonathan Metzman, László Szekeres, Laurent Simon, Read Sprabery, and Abhishek Arya. Fuzzbench: an open fuzzer benchmarking platform and service. In *Proceedings of the 29th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering*, pages 1393–1403, 2021.

[26] Barton P Miller, Lars Fredriksen, and Bryan So. An empirical study of the reliability of unix utilities. *Communications of the ACM*, 33(12):32–44, 1990.

[27] Van-Thuan Pham, Marcel Böhme, Andrew E Santosa, Alexandru Răzvan Căciulescu, and Abhik Roychoudhury. Smart greybox fuzzing. *IEEE Transactions on Software Engineering*, 47(9):1980–1997, 2019.

[28] Jesse Ruderman. Introducing jsfunfuzz. *URL http://www. squarefree. com/2007/08/02/introducing-jsfunfuzz*, 20:25–29, 2007.

[29] Moritz Schloegel, Nils Bars, Nico Schiller, Lukas Bernhard, Tobias Scharnowski, Addison Crump, Arash Ale-Ebrahim, Nicolai Bissantz, Marius Muench, and Thorsten Holz. Sok: Prudent evaluation practices for fuzzing. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 1974–1993. IEEE, 2024.

[30] Kostya Serebryany. {OSS-Fuzz}-google's continuous fuzzing service for open source software. 2017.

[31] Ji Shi, Zhun Wang, Zhiyao Feng, Yang Lan, Shisong Qin, Wei You, Wei Zou, Mathias Payer, and Chao Zhang. {AIFORE}: Smart fuzzing based on automatic input format reverse engineering. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 4967–4984, 2023.

[32] Asia Slowinska, Traian Stancescu, and Herbert Bos. Howard: A dynamic excavator for reverse engineering data structures. In *NDSS*, 2011.

[33] Prashast Srivastava and Mathias Payer. Gramatron: Effective grammar-aware fuzzing. In *Proceedings of the 30th acm sigsoft international symposium on software testing and analysis*, pages 244–256, 2021.

[34] Dominic Steinhöfel and Andreas Zeller. Input invariants. In *Proceedings of the 30th ACM joint european software engineering conference and symposium on the foundations of software engineering*, pages 583–594, 2022.

[35] Romain Thomas. Lief - library to instrument executable formats. https://lief.quarkslab.com/, apr 2017.

[36] Spandan Veggalam, Sanjay Rawat, Istvan Haller, and Herbert Bos. Ifuzzer: An evolutionary interpreter fuzzer using genetic programming. In *European Symposium on Research in Computer Security*, pages 581–601. Springer, 2016.

[37] Junjie Wang, Bihuan Chen, Lei Wei, and Yang Liu. Skyfire: Data-driven seed generation for fuzzing. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 579–594. IEEE, 2017.

[38] Wei You, Xueqiang Wang, Shiqing Ma, Jianjun Huang, Xiangyu Zhang, XiaoFeng Wang, and Bin Liang. Profuzzer: On-the-fly input type probing for better zero-day vulnerability discovery. In *2019 IEEE symposium on security and privacy (SP)*, pages 769–786. IEEE, 2019.

[39] Michal Zalewski. American fuzzy lop. `https://lcamtuf.coredump.cx/afl/`. Accessed: September 10, 2024.

[40] Chijin Zhou, Quan Zhang, Lihua Guo, Mingzhe Wang, Yu Jiang, Qing Liao, Zhiyong Wu, Shanshan Li, and Bin Gu. Towards better semantics exploration for browser fuzzing. *Proceedings of the ACM on Programming Languages*, 7(OOPSLA2):604–631, 2023.