

Blind Spots: Identifying Exploitable Program Inputs

Henrik Brodin

Trail of Bits

New York, USA

henrik.brodin@trailofbits.com

Marek Surovič

Trail of Bits

New York, USA

marek.surovic@trailofbits.com

Evan Sultanik

Trail of Bits

New York, USA

evan.sultanik@trailofbits.com

Abstract—A *blind spot* is any input to a program that can be arbitrarily mutated without affecting the program’s output. Blind spots can be used for steganography or to embed malware payloads. If blind spots overlap file format keywords, they indicate parsing bugs that can lead to exploitable differentials. For example, one could craft a document that renders one way in one viewer and a completely different way in another viewer. They have also been used to circumvent code signing in Android binaries, to coerce certificate authorities to misbehave, and to execute HTTP request smuggling and parameter pollution attacks. This paper formalizes the operational semantics of blind spots, leading to a technique based on dynamic information flow tracking that automatically detects blind spots. An efficient implementation is introduced and evaluated against a corpus of over a thousand diverse PDFs parsed through $M\mu$ PDF¹, revealing exploitable bugs in the parser. All of the blind spot classifications are confirmed to be correct and the missed detection rate is no higher than 11%. On average, at least 5% of each PDF file is completely ignored by the parser. Our results show promise that this technique is an efficient automated means to detect exploitable parser bugs, over-permissiveness and differentials. Nothing in the technique is tied to PDF in general, so it can be immediately applied to other notoriously difficult-to-parse formats like ELF, X.509, and XML.

1. Introduction

We define a *blind spot* as any input to a program that can be arbitrarily mutated without affecting the program’s output. Blind spots are dangerous: they can be exploited for steganography and embedding malware payloads. Steganographic attacks are notoriously difficult to detect automatically, but a brief manual analysis of five of the most popular archive file formats produced fifteen vulnerability disclosures enabled by steganography [1]. One such bug is the “aCROPalypse” (CVE-2023-21036) [2], publicly disclosed on March 18th, 2023. It affects images created on both Google phones and Microsoft Windows, in which cropped image data persists in a blind spot and can subsequently be reconstructed. The blind spot detection technique and

tooling introduced in this paper has been used to detect both aCROPalypsic images as well as vulnerable image generators [3].

Blind spots can also be indicative of parser differentials, for instance, if two parsers exhibit *different* blind spots for the same input. Such differentials can be exploited by crafting a file that renders one way in one parser and a completely different way in another parser [4]. For example, blind spots have been exploited to craft a PDF that can render one way in Adobe Acrobat but have different text when printed [5]. Blind spots have also been used to circumvent code signing in Android binaries [6], to coerce certificate authorities to emit certificates for unauthorized Common Names [7], and to execute HTTP request smuggling and parameter pollution attacks [8]. Our results show that $M\mu$ PDF, on average, ignores 5% of each PDF file. Some PDFs in our corpus had megabytes of ignored data that could be overwritten to store a malware payload. Blind spots can also be useful: They can potentially be excluded as candidates for mutation when generating fuzz testing inputs, similar to the Angora fuzzer’s branch coverage maximization strategy [9].

Blind spots are a generalization of the concept of *file cavities* introduced by Albertini, *et al.* [10]: unused spaces in a file format that are created due to the structure of the surrounding data. However, unlike cavities, blind spots may be dependent on the program itself and its execution environment. For example, the image content of a JPEG file will be a blind spot to a parser that only reads its EXIF metadata. Likewise, the EXIF metadata will be a blind spot to a program that converts JPEGs to another image format like BMP that does not support embedded EXIF metadata. A malware payload could be embedded within the EXIF data without affecting the image’s rendering. Similarly, a JPEG parser for which the EXIF data is *not* a blind spot would likely render a specially crafted image *differently* than a parser for which the EXIF data *is* a blind spot.

Since they are associated with *both* program input *and* the program itself, blind spots can be indicators of parsing vulnerabilities. Parsers—particularly hand-generated ones—will often accept a superset of the grammar for which they were designed. This manifests as a parser that accepts some inputs that are technically invalid according to the file format specification. Sometimes this is intentional, in order to maximize compatibility with files generated by

1. <https://mupdf.com/>

other, incorrectly implemented software, or to attempt to repair malformed documents. For example, we discovered that an optimization in $M\mu$ PDF will sometimes only check the leading “e” in the `endobj` token; it will eagerly accept `eXXXXX` and the PDF will still be parsed correctly, despite the fact that the PDF standard prescribes the existence of the full token. Such lexical permissiveness can lead to parser differentials and so-called *file format schizophrenia* [4]: when two implementations of a file format interpret the same input file differently. A different PDF implementation will ignore the erroneous token and parse further to find the correct delimiter, resulting in different behavior and output.

This paper makes several contributions; it . . .

- 1) formalizes the concept of program input blind spots;
- 2) clarifies the difference between blind spots and file cavities: Blind spots are the set of input bytes whose data flow never influences either control flow that leads to an output or an output itself for a given program;
- 3) proposes a novel technique based on dynamic taint analysis to detect blind spots;
- 4) demonstrates how this technique can be used to automatically detect lexical permissiveness and parser bugs;
- 5) evaluates the approach, providing evidence that these blind spots are correct; and
- 6) analyzes the properties of naturally occurring blind spots in PDFs parsed by $M\mu$ PDF.

2. Background

Dynamic Information Flow Tracking (DIFT), also known as Dynamic Taint Analysis (DTA), is a technique in which the flow of information through a program is modeled and tracked at runtime. DIFT is a challenging problem; many modern approaches suffer from high implementation overhead, low accuracy, and/or low fidelity [11]. *Universal taint analysis* is a form of DIFT that can track all input bytes throughout the execution of a program, mapping inputs to outputs [12].

Our approach to detect blind spots relies on universal taint analysis: We want to identify the input bytes that affect neither program output nor any control flow that leads to a program output. There are significant challenges when performing universal taint analysis on real-world software. If a program accepts n bytes of input, there are $O(2^n)$ ways that each of those bytes could combine to taint each program output. Therefore, when the cyclomatic complexity of a program is large, the amount of combinations generated from even a small input is enormous. For example, many document formats such as PDF use compression to keep file sizes small. When tracking data flow through a PDF parser, there will be a significant number of such combinations as the data are decompressed. This phenomenon is known as *taint explosion*, which generally occurs when a function performs a large number of combinatorial operations on input data.

There are several tools and techniques for performing both DIFT and universal taint analysis. Most tools are designed for fuzz testing use cases, which do not need to track many input bytes. For example, the LLVM compiler framework has its own dataflow analysis instrumentation pass, the DataFlowSanitizer (DFSan) [13], which is limited to tracking at most 8 inputs. Even tools that were explicitly designed for universal taint analysis are unable to scale to input sizes sufficient for real-world inputs. We have a further discussion of the limitations of current tooling in Section 4.1, below.

This paper introduces a new technique that can efficiently achieve universal taint analysis for *megabytes* of input.

3. Definitions and Formalization

This section formalizes the concept of input blind spots. We do this with an extension to Schwartz, Avgerinos, and Brumley’s SIMPIL operational semantics for dynamic taint propagation [14]. Cheney, Ahmed, and Acar’s semantics for *dependency provenance* [15] could also be used to formalize blind spots, however, the literature on provenance as dependency analysis is defined generically as to be compatible with a variety of use cases including databases, file systems, and scientific workflows. SIMPIL, on the other hand, more explicitly defines how program instrumentation would occur. As such, it also directly informs the data structures necessary to achieve efficient universal taint analysis (see Section 4, below). This is why we present the formalism using SIMPIL’s operational semantics.

3.1. An Extension to SIMPIL

The original conception of dataflow analysis in SIMPIL only tracks whether a given variable or memory cell is tainted, not *from whence* it is tainted. In order to detect blind spots, we need to additionally track exactly which input bytes influence output. For example, consider the pseudocode in Algorithm 1. The variable a is tainted by program input on line 2. Moreover, the value of a can indirectly cause hard-coded data ($d = 5$) to be written to output on line 7 by virtue of the conditional on line 5. Therefore, the first byte of the file *cannot* be a blind spot, since its mutation *can* affect output—despite the fact that the value of the first byte is never written to output. Even if the SIMPIL taint policy (*i.e.*, the rules by which taints are propagated) is sufficient to detect that tainted inputs affected the output of the program, it is *insufficient* to detect *which* inputs were responsible. Therefore, we need to extend SIMPIL to additionally track the provenance of a taint so that we can map a complete data flow from inputs to outputs.

SIMPIL uses meta-syntactic variables to represent an execution context. Δ is a mapping of variable names to their values and μ is a mapping from memory addresses to their values. τ_Δ and τ_μ map variable names and memory addresses to booleans (T|F) defining whether or not that value is tainted in the current execution context.

Algorithm 1 Tainted Control Flow

```
1: procedure TAINTEDCONTROLFLOW
2:    $a \leftarrow \text{READINPUT}(1)$   $\triangleright a$  is tainted by the 1st byte
3:    $b \leftarrow \text{READINPUT}(1)$   $\triangleright b$  is tainted by the 2nd byte
4:    $c \leftarrow a + b$ 
5:   if  $c \geq 42$  then
6:      $d \leftarrow 5$ 
7:     WRITEOUTPUT( $d$ )
8:   end if
9: end procedure
```

In order to track taint provenance, we use the concept of a *taint label* [9]: a unique identifier for each instance of a tainted variable or memory address in an execution context. In addition to *union* labels [9] we also define *canonical* labels. A union taint is the result of the combination of two previously tainted values (e.g., the result of two tainted variables being operands in a binary operation). A canonical taint label is the result of a variable or memory address being assigned directly from a program input.

Let $\mathcal{I} = \{\langle s_0, i_0 \rangle, \langle s_1, i_1 \rangle, \dots\}$ be the set of all possible program inputs, where each s is the source (e.g., a file, pipe, socket, or environment variable) and i is the offset within the source. We extend the SIMPIL notation with three new mappings to represent taint labels and track provenance²:

- 1) $\varepsilon : \Delta \cup \mu \rightarrow \mathbb{N}$ that maps variable names and memory addresses to unique taint labels³;
- 2) $\kappa : \mathbb{N} \rightarrow \mathcal{I}$ that maps canonical taint labels to the information about their source; and
- 3) $\gamma : \mathbb{N} \rightarrow \mathbb{N} \times \mathbb{N}$ that maps union taint labels to their parents.

Note that the γ mapping implicitly defines a directed acyclic graph (DAG) where the out-degree of each vertex is at most two. However, instructions that operate on tainted values could have an arity *higher* than two, causing the result label to have more than two parents. In such cases, we add multiple unions to the γ mapping. This is done for the sake of notational simplicity and does not affect our results.

SIMPIL treats these mappings more like programmatic hashmaps than set theoretic functions. As such, SIMPIL uses the notation “ $\kappa[\ell]$ ” for the value of taint label ℓ in mapping κ . For brevity, we shall continue this theme by using the notation $\ell \in \kappa$ to represent the fact that ℓ is a key in the mapping κ , $\ell \notin \kappa$ to mean that ℓ is not a key in κ , and $|\kappa|$ to mean the number of key/value pairs in the mapping.

The zero taint label is reserved to represent untainted variables and memory. An untainted variable v will always lack source info and descend from the zero label:

$$\tau_{\Delta}[v] = \begin{cases} \mathbf{F} & \varepsilon[v] \notin \kappa \wedge \gamma[\varepsilon[v]] = \langle 0, 0 \rangle, \\ \mathbf{T} & \varepsilon[v] \in \kappa \vee (\gamma[\varepsilon[v]] = \langle i, j \rangle \wedge i + j > 0). \end{cases}$$

2. The SIMPIL semantics are already replete with Greek letters. We have chosen ε from the Greek word for “labels” (*επιγραφες*), κ from the word for “canonical” (*κανονικος*), and γ from the word for “parent” (*γονεως*).

3. Δ and μ retain their original meanings from SIMPIL.

The SIMPIL *taint policy* (see Table III from [14]) and semantics are modified to update these mappings on every taint status change. For example, the updated semantics for reading from input and for executing binary operations are given in Figure 1.

3.2. Mapping Taint Sources to Sinks

These mappings allow us to track the entire provenance of a taint in any execution context. As defined above, the γ mapping implicitly creates a DAG of labels, representing the dataflow through the program: The program inputs that affect a tainted variable or memory address are its taint label’s topmost ancestors in the γ DAG. A recursive function $\psi : \mathbb{N} \rightarrow 2^{\mathcal{I}}$ can map taint labels to all of their ancestral sources:

$$\psi[\ell] = \begin{cases} \{\kappa[\ell]\} & \ell \in \kappa, \\ \bigcup_{p \in \gamma[\ell]} \psi[p] & \gamma[\ell] \neq \langle 0, 0 \rangle, \\ \emptyset & \text{otherwise.} \end{cases}$$

When we observe that the program writes to output, we use the ψ mapping to record which inputs, if any, tainted the output. This can be accomplished by enumerating the canonical ancestors of the output labels by traversing the γ DAG. This allows us to construct a complete mapping of taint sources to sinks. Note that any taint sources without associated sinks can be arbitrarily mutated without affecting the output.

3.3. A Definition of Program Input Blind Spots

Let Ω be the set of taint labels written to output during execution. Then a blind spot is the set of all potential program inputs that did not affect the output:

$$\begin{aligned} \mathcal{I}_{\neg\text{READ}} &= \{\iota \in \mathcal{I} : (\forall \ell \in \kappa : \kappa[\ell] \neq \iota)\} \\ \mathcal{I}_{\neg\text{INOUTPUT}} &= \{\kappa[\ell] : \ell \in \kappa \wedge (\forall \ell' \in \Omega : \kappa[\ell] \notin \psi[\ell'])\} \\ \mathcal{I}_{\text{BLINDSPOT}} &= \mathcal{I}_{\neg\text{READ}} \cup \mathcal{I}_{\neg\text{INOUTPUT}}. \end{aligned} \quad (1)$$

As we mentioned above, SIMPIL includes a *taint policy* specifying the rules by which taints are propagated. The taint policy will affect our definition of blind spots. For example, let us again consider the pseudocode in Algorithm 1. The variable a is tainted by program input on line 2, which we can now specify in our SIMPIL extension as $\kappa[\varepsilon[a]] \neq \emptyset$. Recall that the value of a can indirectly cause hard-coded data ($d = 5$) to be written to output on line 7. Therefore, a cannot be a blind spot, since its mutation can affect output. However, the semantics by which the taint policy propagates taint through conditionals will affect whether our extension of SIMPIL will consider the output to be tainted by a , because the value of a itself is never written to output.

We resolve this discrepancy by enforcing the following constraint on blind spot taint policies: The taint labels of every variable and memory address that affect the program’s control-flow will be unioned with all labels created in the branch they influence. In other words, in addition to the

$$\begin{array}{c}
\text{SIMPIIL notation for the computation performed by the INPUT operation} \\
\hline
\frac{v \text{ is input from } src \in \mathcal{I} \quad \varepsilon' = \varepsilon[v \leftarrow |\varepsilon| + 1] \quad \kappa' = \kappa[\varepsilon'[v] \leftarrow src]}{\mu, \Delta, \varepsilon, \kappa, \gamma \rightsquigarrow \mu, \Delta, \varepsilon', \kappa', \gamma} \quad \text{INPUT} \\
\hline
\begin{array}{c}
\text{SIMPIIL notation for the} \\
\text{updated execution state after the operation}
\end{array}
\quad
\begin{array}{c}
\text{SIMPIIL notation for the evaluation of} \\
\text{expression } \text{get_input}(src) \text{ to value } v \text{ in context} \\
\mu, \Delta
\end{array}
\end{array}$$

$$\frac{\mu, \Delta \vdash e_1 \Downarrow v_1 \quad \mu, \Delta \vdash e_2 \Downarrow v_2 \quad v' = v_1 \diamond_b v_2 \quad \varepsilon' = \varepsilon[v' \leftarrow |\varepsilon| + 1] \quad \gamma' = \gamma[\varepsilon'[v'] \leftarrow \langle \varepsilon[v_1], \varepsilon[v_2] \rangle]}{\mu, \Delta, \varepsilon, \kappa, \gamma \rightsquigarrow \mu, \Delta, \varepsilon', \kappa, \gamma' \quad \mu, \Delta \vdash e_1 \diamond_b e_2 \Downarrow v'} \quad \text{BINOP}$$

Figure 1. SIMPIIL operational semantics for reading input and executing binary operators (see Figure 1 of [14]) updated to include data flow provenance tracking. When a value v is read from an input source src , we create a new, unique canonical taint label for v and set its taint source info in κ to src . When a binary operator \diamond_b is applied to expressions $e_1 = v_1$ and $e_2 = v_2$ resulting in the value v' , we create a new, unique union taint label for v' and set its parents to be the taint labels associated with values v_1 and v_2 .

definition of blind spots in Equation (1), a blind spot cannot influence control flow that leads to a program output. Updated operational semantics for the conditional operator that implement this policy are given in Figure 2.

A trace of Algorithm 1 showing the iterative updates to the execution context is given in Table 1. It demonstrates how the blind spot taint policy propagates taints from variables in the path condition—variables that have affected control flow leading to the current state (*e.g.*, variable c on line 5)—to variables that would otherwise not be tainted (*e.g.*, variable d). This reduces false-positive blind spot classifications, since it captures tainted variables that indirectly cause output.

4. Implementation

Thus far we have developed formal semantics for blind spots and discovered some necessary taint propagation policies to detect them. The next step is to automatically instrument a parser to extract the data flow information necessary to classify input byte regions as blind spots. We gather this data flow information by performing universal taint analysis.

PolyTracker [16] is an LLVM-based dynamic analysis tool that we have developed for extracting ground truth information from programs. It is open-source and available at <https://github.com/trailofbits/polytracker>.

PolyTracker automatically adds instrumentation to a program such that, when the program is executed, it produces runtime artifacts that can be analyzed to track the data flows of all input bytes. It is an extension of the LLVM DataFlowSanitizer (DFSan) [13], a generalized dynamic data flow analysis instrumentation tool. PolyTracker has previously been used to label the semantic purpose of functions in a parser [17].

Originally, DFSan supported tracking at most 2^{16} taint labels at a time. This limit was only sufficient to track at most several hundred input bytes at once. Over the course of 2021, DFSan underwent a significant refactor in order to make its memory layout compatible with other LLVM sanitizers [18]; this refactor reduced the effective number of taints it could track to $2^3 = 8$. This restriction was acceptable for DFSan’s primary use case at the time: data

flow analysis for fuzz testing, but *not* for universal taint analysis and detecting blind spots.

We forked PolyTracker off of the final, pre-refactor version of DFSan. For the remainder of this section, our discussion of DFSan will refer to this version. This version of DFSan works by creating a region of “shadow” memory that can store a taint label associated with every address on the stack and heap. For every instruction in the program, DFSan checks its operands to see if they have associated taint labels in shadow memory. If the labels are different, it means that the operands were tainted by different input data flows, and DFSan will create a new label that represents the union of the two. DFSan also has mechanisms for propagating taint information across function calls (*e.g.*, by appending taint labels as function arguments), as well as models for taint propagation through uninstrumented system calls.

DFSan uses a $2^{16} \times 2^{16}$ matrix to store the unions generated when instructions mix taint labels. Element i, j of the matrix holds the value of the label produced by the union of labels i and j . This matrix representation is computationally efficient but becomes prohibitively large as the maximum number of taint labels grows. For n -bit taint labels, the matrix will require

$$\Theta\left(2^n \times 2^n \times \frac{n}{8}\right) = \Theta(2^{2n-3}n)$$

bytes of memory. This is the reason for DFSan’s limit of 2^{16} taint labels: Increasing the limit to 2^{32} labels would require over 73 *exabytes* of RAM to store the union matrix. We need a way to increase this 2^{16} limit by at least a few orders of magnitude.

Our solution to this problem arises from the γ mapping we added to the SIMPIIL operational semantics. We create a memory-mapped file where each taint label has a fixed-size entry containing the indexes of its parent labels. PolyTracker adds additional instrumentation to:

- tag input sources as canonical taints (building the κ mapping from our semantics);
- track the taint labels that are written to output; and
- track which taint labels affect control flow.

$$\frac{\text{the conditional expression } e \text{ evaluates to } v \quad \overbrace{\mu, \Delta \vdash e \Downarrow v} \quad \text{create a new taint label for every existing label} \quad \overbrace{\varepsilon' = \varepsilon[u \leftarrow |\varepsilon| + \varepsilon[u] : \forall u \in \varepsilon]} \quad \text{union every existing label with the taints of } v \quad \overbrace{\gamma' = \gamma[\varepsilon'[u] \leftarrow \langle \varepsilon[u], \varepsilon[v] \rangle : \forall u \in \varepsilon]} \quad \text{PRECOND}}{\mu, \Delta, \varepsilon, \kappa, \gamma \rightsquigarrow \mu, \Delta, \varepsilon', \kappa, \gamma'}$$

Figure 2. Updated SIMPIL operational semantics to enforce the taint policy that every input that affects control-flow will be unioned with all labels created in the branch they influence. The PRECOND rule is executed before every conditional rule (TCOND and FCOND in Figure 1 of [14]).

Line	STATEMENT	Δ	τ_Δ	ε	κ	γ
1	start	{}	{}	{}	{}	{}
2	$a \leftarrow \text{READINPUT}(1)$	$\{a \rightarrow 40\}$	$\{a \rightarrow \mathbf{T}\}$	$\{a \rightarrow 1\}$	$\{1 \rightarrow \langle 1^{\text{st}} \text{ byte of input} \rangle\}$	{}
3	$b \leftarrow \text{READINPUT}(1)$	$\{a \rightarrow 40, b \rightarrow 12\}$	$\{a \rightarrow \mathbf{T}, b \rightarrow \mathbf{T}\}$	$\{a \rightarrow 1, b \rightarrow 2\}$	$\{1 \rightarrow \langle 1^{\text{st}} \text{ byte of input} \rangle, 2 \rightarrow \langle 2^{\text{nd}} \text{ byte of input} \rangle\}$	{}
4	$c \leftarrow a + b$	$\{a \rightarrow 40, b \rightarrow 12, c \rightarrow 52\}$	$\{a \rightarrow \mathbf{T}, b \rightarrow \mathbf{T}, c \rightarrow \mathbf{T}\}$	$\{a \rightarrow 1, b \rightarrow 2, c \rightarrow 3\}$	$\{1 \rightarrow \langle 1^{\text{st}} \text{ byte of input} \rangle, 2 \rightarrow \langle 2^{\text{nd}} \text{ byte of input} \rangle\}$	$\{3 \rightarrow \langle 1, 2 \rangle\}$
5	if $c \geq 42$ then	$\{a \rightarrow 40, b \rightarrow 12, c \rightarrow 52\}$	$\{a \rightarrow \mathbf{T}, b \rightarrow \mathbf{T}, c \rightarrow \mathbf{T}\}$	$\{a \rightarrow 1, b \rightarrow 2, c \rightarrow 3\}$	$\{1 \rightarrow \langle 1^{\text{st}} \text{ byte of input} \rangle, 2 \rightarrow \langle 2^{\text{nd}} \text{ byte of input} \rangle\}$	$\{3 \rightarrow \langle 1, 2 \rangle\}$
6	$d \leftarrow 5$	$\{a \rightarrow 40, b \rightarrow 12, c \rightarrow 52, d \rightarrow 5\}$	$\{a \rightarrow \mathbf{T}, b \rightarrow \mathbf{T}, c \rightarrow \mathbf{T}, d \rightarrow \mathbf{T}\}$	$\{a \rightarrow 1, b \rightarrow 2, c \rightarrow 3, d \rightarrow 4\}$	$\{1 \rightarrow \langle 1^{\text{st}} \text{ byte of input} \rangle, 2 \rightarrow \langle 2^{\text{nd}} \text{ byte of input} \rangle\}$	$\{3 \rightarrow \langle 1, 2 \rangle, 4 \rightarrow \langle 3, 0 \rangle\}$

TABLE 1. EXECUTION CONTEXT TRACE FOR ALGORITHM 1. NOTE ON LINE 6 THAT, DESPITE BEING ASSIGNED A CONSTANT VALUE OF 5, THE d VARIABLE (TAINT LABEL 4) IS IN FACT TAINTED BY VARIABLE c (TAINT LABEL 3). THIS IS BECAUSE THE PATH CONDITION TO LINE 6 DEPENDS ON c FROM THE CONDITIONAL BRANCH ON LINE 5. THEREFORE, NEITHER OF THE FIRST TWO BYTES OF INPUT ARE BLIND SPOTS.

The algorithm for determining blind spots from a program trace is given in Algorithm 2. Set L denotes the set of all taint labels in a program trace. Consequently set B denotes all taint labels that did not affect control flow. The algorithm iterates through labels in L in descending order. On line 5, $\ell \notin B$ means that taint label ℓ or its descendant affected control flow, while $\ell \in \Omega$ means that taint label ℓ was written to output. If either of these is true ℓ and its parents are removed from B on line 6. Finally on line 9, the set all blind spots $\mathcal{I}_{\text{BLINDSPOT}}$ is the set of all inputs which have their associated taint label in B . The algorithm runs in $O(n)$ time where n is the number of taint labels created during the trace. This has proven sufficient to detect blind spots in all programs and inputs on which we have experimented.

4.1. Related Work

There are several existing projects that achieve universal taint tracking, using various methods. Two of the best maintained and easiest to use are AUTOGRAM [19] and TaintGrind [20]. However, the former is limited to analysis within the Java virtual machine and the latter suffers from unacceptable runtime overhead when tracking as few as several bytes at a time. For example, we ran `mutool`, a utility in the `MμPDF` project, using TaintGrind over a corpus of medium sized PDFs, and in every case the tool had to be halted after over twenty-four hours of execution for operations that would normally complete in milliseconds without instrumentation.

Algorithm 2 Enumerate Blind Spots

Ensure: $\mathcal{I}_{\text{BLINDSPOT}}$ is the set of blind spots in the trace

- 1: **procedure** BLINDSPOTS($\Omega, \varepsilon, \kappa, \gamma$)
- 2: $L \leftarrow \{\varepsilon[v] : v \in \varepsilon\} \cup \{\ell : \ell \in \kappa\}$ \triangleright the set of all taint labels in the trace
- 3: $B \leftarrow \{\ell \in L : \neg \text{AFFECTEDCONTROLFLOW}(\ell)\}$ \triangleright the set of all labels that did not affect control flow
- 4: **for each** $\ell \in \text{SORTDESCENDING}(L)$ **do**
- 5: **if** $\ell \notin B \vee \ell \in \Omega$ **then** $\triangleright \ell$ cannot be a blind spot because it or one of its descendants affected output
- 6: $B \leftarrow B \setminus (\{\ell\} \cup \gamma[\ell])$
- 7: **end if**
- 8: **end for**
- 9: $\mathcal{I}_{\text{BLINDSPOT}} \leftarrow \{\kappa[\ell] : \ell \in B\}$
- 10: **end procedure**

There are also existing tools for performing dynamic program analysis via QEMU [21], such as PANDA [22] and DECAF(++) [23], both of which have taint tracking extensions. However, being an emulation framework rather than virtualization, QEMU incurs a runtime overhead of about 15% just to execute a binary, not including any instrumentation [24]. After adding the program instrumentation necessary to enable fuzz testing, QEMU was observed to have over three times the runtime overhead of equivalent compile-time instrumentation [25].

Symbolic execution engines like Triton [26] and SymCC [27] have also been used for data flow analysis. Symbolic execution could be extended to detect blind spots,

e.g., by making all input bytes symbolic and observing all data that is written. The input bytes associated with any symbolic data that is either written or included in the path condition during a write *is not* a blind spot.

DRTaint [12] is a recently published tool that can also perform universal taint tracking. It adds a minimal amount of runtime instrumentation to create runtime artifacts that can be post-processed to extract any data flow. The authors do not quantify the exact overhead of DRTaint, but Figure 5 from their paper suggests at least a 60x slowdown compared to the uninstrumented program. It is also unclear whether this instrumentation was sufficient to reconstruct all data flows. This is consistent with earlier techniques such as Dytan [28] that reported a 50x slowdown when tracking as few as 64 taint labels.

5. Evaluation

The previous sections introduced a method for identifying the blind spots of a program input. How accurate is this blind spot classifier? Since we do not have pre-labeled ground truth for the blind spots of an input, we need to develop statistical estimates for the confusion matrix of our classifier.

We focus our evaluation on the PDF file format, for several reasons.

- 1) The PDF file format is old and complex, has had many revisions, and enjoys numerous independent implementations. This has led to differentials that necessitate lexical permissiveness for interoperability [29].
- 2) PDF is a container format that allows embedding of other formats like JPEG, providing more opportunity for blind spots.
- 3) The GovDocs corpus [30] provides thousands of real-world PDFs generated by a diversity of software.

We instrumented `mutool`, a utility in the popular `MμPDF` project, using PolyTracker to detect blind spots. Next, we ran the instrumented utility on 1,087 PDFs sampled from the GovDocs corpus to render the PDFs to PostScript. The PDFs totaled over 622 MB and averaged 572 KB each. The largest file was 9.6 MB. We discovered a total of 33.5 MB of blind spots, averaging 30.8 KB per file, some files having zero, and one file having 1.07 MB of blind spots.

PostScript, a vector image format, was chosen rather than a raster format like JPEG or PNG because the relative lack of compression would help prevent explosion of taint union labels and thereby reduce runtime. Runtime of the instrumented program was less than one minute for each PDF. We would expect to get similar results when rendering to JPEG or PNG, however, since blind spots are dependent on execution, there could be some discrepancies. For example, since fonts can be embedded in both PDF and PostScript but cannot be embedded in JPEG or PNG, one might expect to

```

Lorem ipsum dolor sit amet,
consectetur adipiscing elit,
sed do eiusmod tempor incididunt
ut labore et dolore magna
aliqua. Ut enim ad minim veniam,
quis nostrud exercitation
ullamco laboris nisi ut aliquip
ex ea commodo queredat. Duis
aute irure dolore herit
in voluptate velit esse cillum
dolore eu fugiat nulla pariatur.

```

Figure 3. Notional example of validating detected blind spots. Underlined bytes are iteratively mutated and re-parsed. Bytes outside of blind spots are randomly selected for mutation, but all bytes within a blind spot are mutated. If a byte *inside* a blind spot is mutated and produces a different result, then that blind spot classification was a false-positive (Type I error). If a byte *outside* a blind spot can be arbitrarily mutated, then it is a missed detection (Type II error).

see more blind spots related to fonts if one were to have rendered to a different file format.

5.1. Classification Error

We validate blind spots by iteratively mutating each classified blind spot byte in the input file and re-running the original, uninstrumented program again. See Figure 3 for a notional example of this mutation process. If the PostScript output produced from the mutated input is different from the output of the unmodified file, then our classified blind spot is incorrect (Type I error), since any mutation inside a blind spot should, by definition, not affect program output.

We also sample bytes from *outside* of our classified blind spots and mutate them, similarly. If a byte outside a blind spot can be arbitrarily mutated, it is likely a missed detection (Type II error). However, it is not tractable to mutate and verify all possible combinations of input bytes, since this would amount to testing the power set of all bytes, running in $\Theta(2^n)$ time. While detecting blind spots was relatively quick—several hours to process the 1,087 PDF corpus—, validating the blind spots by mutating the inputs required about a month of computation.

It might be the case that a byte outside a blind spot is in fact a byte that can be *almost* arbitrarily mutated, but has some undetected data dependency on another byte. For example, a source code comment is this form of input to a compiler, since the bytes within the comment can be arbitrarily changed without affecting the behavior of the compiler *as long as* the bytes do not contain the comment delimiter. Therefore our reported Type I error rate is a tight bound on the actual Type I error, but our Type II error rate is a loose upper bound on the true Type II error.

We mutated all 33.5 million blind spot bytes classified in the corpus, and all mutated blind spots produced identical output to the original file for a 0% false-positive rate. Of the bytes *not* classified as blind spots, 89% did affect the output. Therefore, the false-negative rate is bounded above by 11%. A significant number of these missed detections are likely data that can be mutated *almost* arbitrarily, but

BYTES	# Blind Spots	Total Frequency
\r	561451	6745095
\n	52175	4996436
\x20	47752	17967270
\x00	9504	7658663
\x11	5232	3325230
\x08	3930	3611559
\x06	3629	3104109
\x09	3572	2981341
%	3145	3041586
\x12	2932	2847588
#	2859	3399908
\x05	2772	2873564
\x02	2640	3007837
\x0F	2607	2944161
\x14	2594	3132850
\xF0	2351	2930162
a	2348	4959732
!	2271	3172472
\xA3	2108	2783202
4	2074	4848767
\x13	2024	2788459

TABLE 2. THE TWENTY MOST COMMON BLIND SPOT PREFIXES OF LENGTH AT MOST SEVEN BYTES.

have some data dependency that can affect output that our random mutations did not exercise.

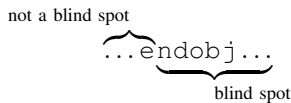
5.2. Blind Spot Content

What is the content of PDF blind spots for $M\mu$ PDF?

In the GovDocs PDF corpus parsed by $M\mu$ PDF, we detect 63,194 unique blind spot prefixes of length at most seven bytes, and 338,943 unique byte sequences of length at most seven that precede blind spots.

Consider the bytes that occur at the start of a blind spot; the most common of these are listed in Table 2. They are all one byte long, meaning that there is a diversity of content at the start of blind spots. The most common byte sequences that *precede* a blind spot, listed in Table 3, are more interesting: most are multi-byte, and they comprise many PDF tokens like `endobj`. This means that bytes following certain tokens are often or always ignored by the parser.

Now let us consider the unique suffix/prefix pairs that occur adjacent to the start of a blind spot. There are 1,029,129 unique pairs of these byte sequences. If we sort them by frequency, the pair



is in the top 0.01% of such pairs. “`endobj`” is a PDF token used to delimit the end of an object in the document model. The fact that this token is split across a blind spot boundary so frequently is indicative of, at best, intentional lexical permissiveness on the part of the parser, and, at worst, a bug. Other interesting blind spot contexts within the top hundredth of the first percentile include the entire `endobj` token, if preceded by a carriage return. Any whitespace after the `stream` token is ignored. The `obj` token is completely

BYTES	# Blind Spots	Total Frequency
n	300057	5811224
\ n	299785	629959
00000\ n	299621	555279
f	236173	3738217
\ f	235736	552315
65535\ f	205564	470855
m	66860	4242668
dstream	66588	189282
00001\ f	26343	44612
\r	12839	6745095
endobj\r	11001	527199
e	7056	7495303
be	6419	41673
Adobe	6418	11545
\x00\x0EAdobe	6417	9142
\x02	5533	3007837
\x08	4185	3611559
\x00\x02	3787	82034
00000\ f	3770	6838
\x01\x00\x02	3707	29366

TABLE 3. THE TWENTY MOST COMMON BYTE SEQUENCES OF LENGTH AT MOST SEVEN PRECEDING A BLIND SPOT.

ignored if preceded by a space. Similarly, the PDF dictionary delimiters `<<` and `>>` are frequently skipped, *e.g.*, at the beginning of a PDF object. This simple contextual blind spot frequency analysis can discover parsing errors and differentials.

5.3. Blind Spot Context

How frequent are blind spots, and where do they occur in PDFs?

Figure 4 plots the number of blind spots in the PDF corpus as a function of file size. This suggests that the number of blind spot bytes in a typical PDF is constant. Note, however, that blind spots in PDFs can be arbitrarily large, since the PDF format permits the inclusion of arbitrary binary blobs that do not have to be connected to the document object model (DOM) [29].

Figure 5 plots a histogram of the contiguous size of blind spot regions in the corpus. The majority of blind spots are small, but a nontrivial number of blind spots are over 1 KB. The average blind spot is 42 bytes long with a standard deviation of 1.72 KB. This demonstrates the fact that PDF is a container format that can contain arbitrarily large binary blobs that do not have to contribute to PDF rendering.

Figure 6 is a histogram of the normalized position of blind spot bytes in their files: the blind spot’s byte offset divided by the file size. In our experiments with the $M\mu$ PDF renderer translating to PostScript, the majority of PDF blind spots are at the beginning and ends of the files. This is not surprising since the beginning of a PDF typically includes metadata that is not necessary for rendering, and the end of the PDF typically includes an XREF table that can be ignored, particularly if the PDF is not malformed. Anecdotally, we have observed in the GovDocs corpus that PDF generators often add additional metadata objects that do not contribute to rendering, typically toward the end of the file.

Figure 7 combines the two previous figures by comparing the mean contiguous blind spot size to the normalized

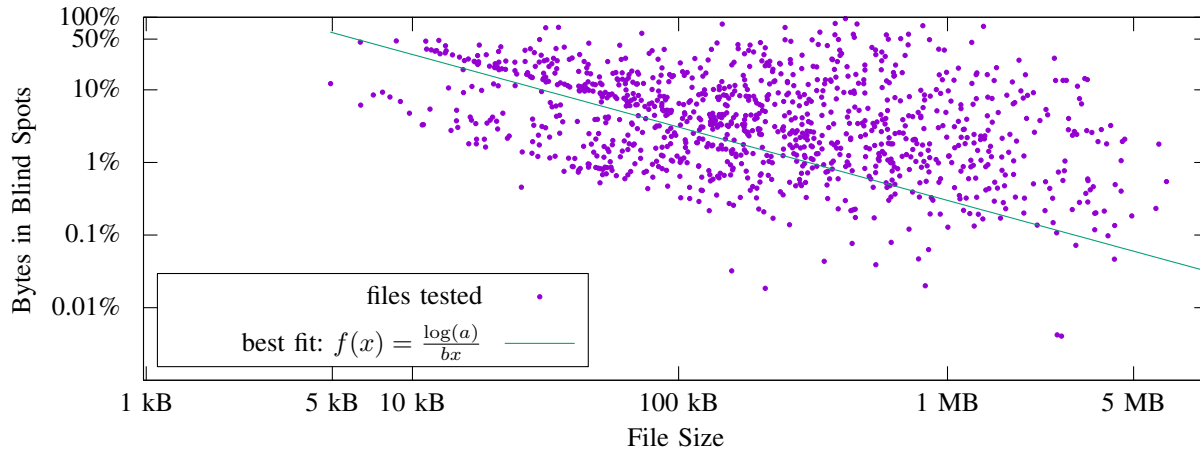


Figure 4. Proportion of PDF bytes parsed by $M\mu$ PDF that are in blind spots as a function of file size. This suggests that the number of blind spot bytes in a typical $M\mu$ PDF-parsed PDF is constant.

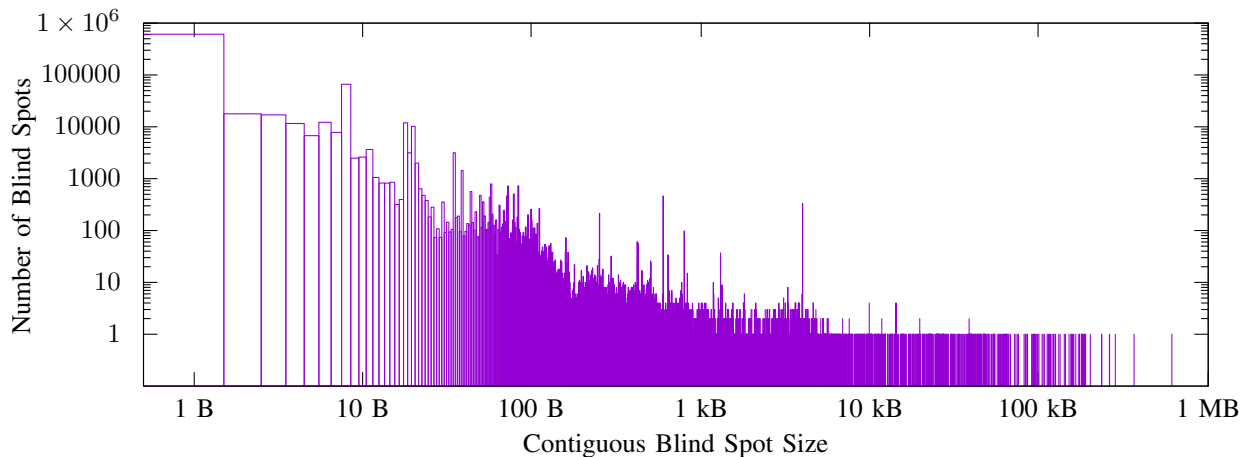


Figure 5. Histogram of the sizes of contiguous PDF blind spots for $M\mu$ PDF. The majority of blind spots are single bytes, but a nontrivial number of blind spots are over 1 kilobyte. The average blind spot is 42 bytes long with a standard deviation of 1.72 kilobytes. This demonstrates the fact that PDF is a container format that can contain arbitrarily large binary blobs that do not have to contribute to PDF rendering.

position in the PDF. Despite the most blind spot bytes being at the beginning and ends of the PDFs, the *longest* blind spots tend to be in the first 10–20% of the file, but not immediately at the beginning. The abundance of blind spots toward the end of the document tend to be small. In order to explain them, we need to look for patterns in their semantic context.

In order to assign a semantic context to each byte of input, we generate a parse tree for each PDF using PolyTracker’s sister tool, PolyFile⁴ [16]. Each byte in the input PDF corresponds to one or more *parse tree derivations*: unique paths through the PDF parse tree. A byte could have more than one derivation, for instance, if the input file is a *polyglot*—a file that is valid in two or more formats [4]. PDFs are particularly easy to turn into polyglots, and many

legitimate PDF generators exploit this fact. For example, it is common to produce valid PDFs that are *also* valid ZIP archives that, when extracted, contain additional files related to the document. Therefore, a byte might have one derivation in the PDF parse tree and have a different but completely valid derivation in the ZIP parse tree. An example of a parse tree derivation is given in Figure 8.

For each unique parse tree derivation, we count the number of blind spot bytes that occur in that derivation. The most frequent derivation containing blind spots is `application/pdf`, the root of the PDF parse tree. This means that the majority of PDF blind spots occur in portions of the file that have no semantic purpose. Blind spot locality might be explained by the fact that PDF parsers are resilient to both leading and trailing garbage bytes before and after the PDF file. Also, as we saw above, the majority of naturally occurring blind spot bytes are at the beginning and end of the file.

4. Open-source and available at <https://github.com/trailofbits/polyfile>.

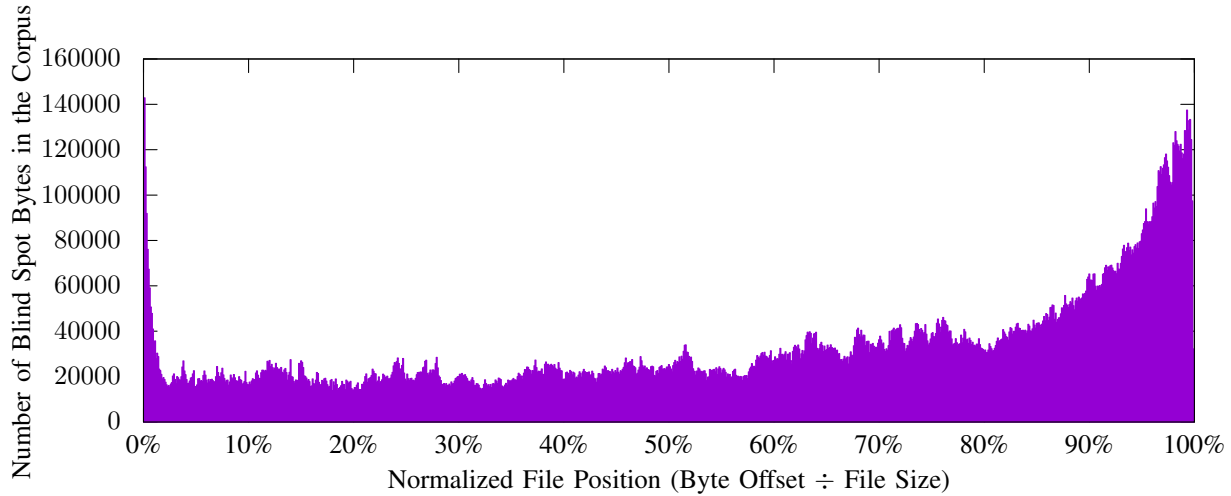


Figure 6. Normalized position of blind spots in PDF files. With the $M\mu$ PDF renderer, blind spots are most frequent at the beginning and end of PDF files. The large number of blind spots toward the end of files can be explained by their context within the PDF cross-reference (XREF) table (*q.v.* Table 4), which is not strictly necessary for rendering.

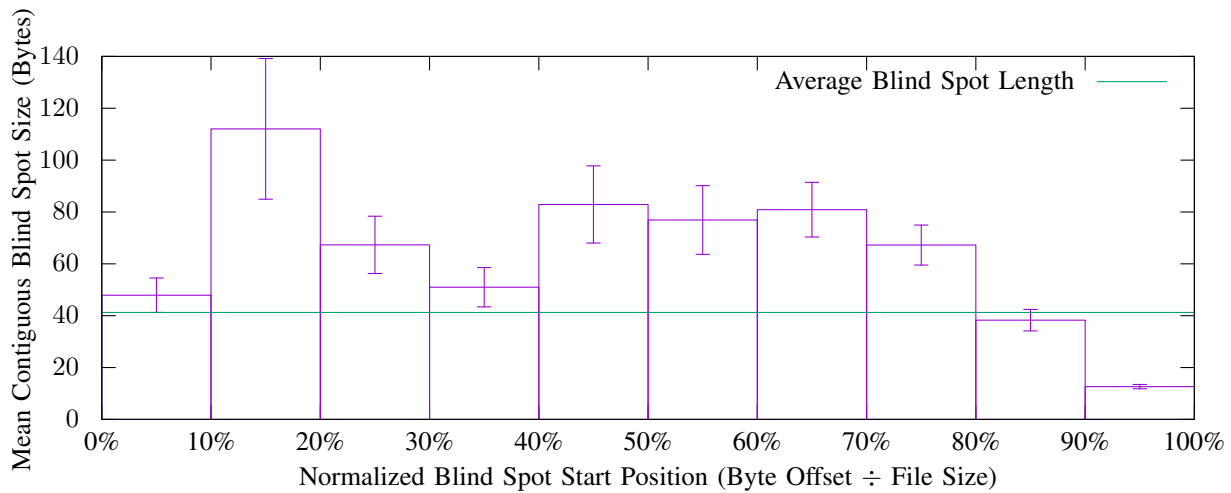


Figure 7. Contiguous blind spot size as a function of its position in the PDF file parsed by $M\mu$ PDF. Error bars correspond to the standard deviation of blind spot sizes in that portion of the file. The longest blind spots tend to be in the second tenth of the file.



Figure 8. $M\mu$ PDF's parse tree derivation of a byte representing an integer in a list of lists that is a value in a PDF dictionary in a PDF object.

The frequency of every unique parse tree derivation is presented in Figure 9. Blind spots overwhelmingly occur in a small number of derivations. However, the long tail demonstrates that blind spots can and do occur in many diverse derivations.

The most frequent parse tree derivations for PDF blind spots are given in Table 4. Almost all of the derivations

descend from the PDFObject node. This is unsurprising, since PDF objects can contain streams of arbitrary binary data. PDF objects also do not need to be connected to the root of the PDF document object model, nor do they need to be used in any way for rendering. The third most common context in which blind spots occur is within the cross-reference (XREF) table. The XREF table is used by the parser to quickly look up the file offsets of PDF objects, decreasing load times. However, the XREF table is not strictly necessary to parse a PDF, and almost all parsers are resilient to errors or omissions in the XREF table. Therefore, it is unsurprising that there would be many blind spots within the XREF table. The XREF table usually occurs toward the end of the PDF, explaining the positional

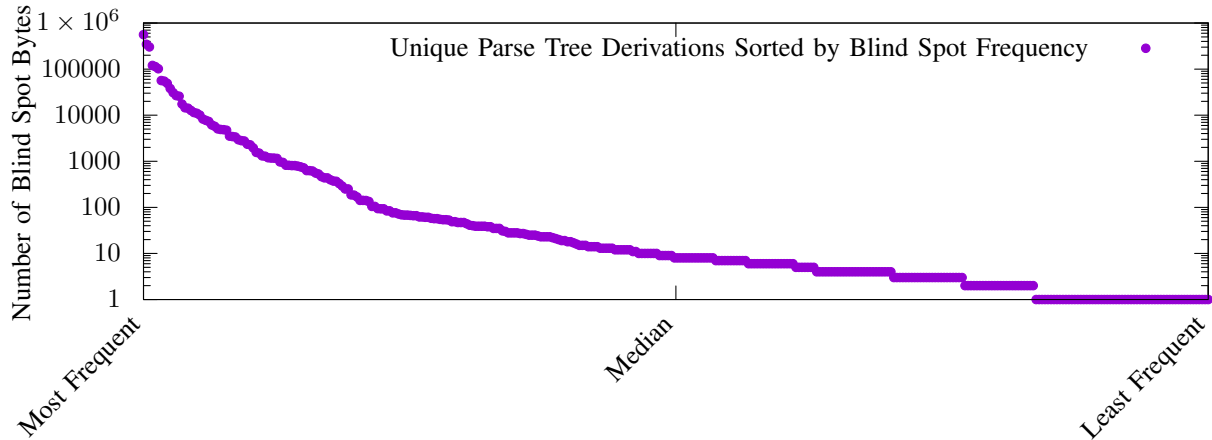


Figure 9. Each point is a *parse tree derivation*—a unique path through a $M\mu$ PDF PDF parse tree—whose y -axis value is the number of times a blind spot occurred in that derivation. Blind spots overwhelmingly occur in a small number of derivations, yet there is a long tail demonstrating that blind spots can and do occur in many diverse derivations.

distribution in Figure 6.

The second most frequent derivation for blind spots are in PDF dictionaries. PDF dictionaries can contain arbitrary key/value pairs which are often used for metadata that is not necessary for rendering (*e.g.*, timestamps). Dictionaries can and often do contain redundant information. For example, the length of a PDF object stream can either be specified as a key/value pair in the preceding object dictionary or implicitly defined by the location of a required termination token. If both are specified, then they must agree. However, if a length is specified in the dictionary which *does not* agree with the position of the termination token, then most parsers will ignore the specified length and defer to the token position [31], making the dictionary entry a blind spot.

6. Conclusions

This paper defined the concept of blind spots: inputs to a program that can be arbitrarily mutated without affecting the program’s output. Operational semantics for blind spots were formalized by extending SIMPIL [14]. An efficient implementation capable of automatically detecting blind spots, PolyTracker, was introduced. It works by adding instrumentation for performing dynamic information flow tracking (DIFT) to a program.

The technique was evaluated by detecting blind spots in the popular $M\mu$ PDF parser over a corpus of over a thousand diverse PDFs [30]. There were zero false-positive blind spot classifications, and the missed detection rate was bounded above by 11%. On average, at least 5% of each PDF file was completely ignored by the parser; blind spots that could be repurposed for steganography or embedding malware payloads.

Future work includes extending the approach to detect inputs that can *almost* arbitrarily be mutated without affecting output, like source code comments. Using the revealed blinds spots to identify parser differentials, by comparing

the blind spots of different instrumented parsers, would provide insight into vulnerabilities stemming from parsers interpreting the same file differently. The current implementation injects its DIFT instrumentation at the LLVM/IR level. Therefore, it is limited to programs that can be compiled using LLVM, or binaries that can be lifted to LLVM/IR. It would be useful to apply the technique to runtime instrumentation that could be applied to a black-box binary.

Our results show promise that this technique could be an efficient automated means to detect parser bugs and differentials. Nothing in the technique is tied to PDF in general, so it can be immediately applied to other notoriously difficult-to-parse formats like ELF, X.509, and XML.

Acknowledgments

This research was supported in part by the Defense Advanced Research Projects Agency (DARPA) SafeDocs program as a subcontractor to Galois under HR0011-19-C-0073. Many thanks to Michael Brown, Trent Brunson, Filipe Casal, Peter Goodman, Kelly Kaoudis, Lisa Overall, Stefan Nagy, Bill Harris, Nichole Schimanski, Mark Tullsen, Walt Woods, Peter Wyatt, Ange Albertini, and Sergey Bratus for their invaluable feedback on the approach and tooling. Thanks to Ange Albertini for suggesting *«angles morts»*—French for “blind spots”—to name the concept. Special thanks to Carson Harmon, the original creator of PolyTracker, whose ideas and discussions germinated this research.

References

- [1] M. Vuksan, T. Pericin, and B. Karney, “Hiding in the familiar: Steganography and vulnerabilities in popular archives formats,” in *Proceedings of Black Hat Europe*, 2010.
- [2] D. Buchanan, “Exploiting aCROPalypse: Recovering truncated PNGs,” <https://www.da.vidbuchanan.co.uk/blog/exploiting-acropalypse.html>, March 18, 2023, accessed: 2023-03-18.

DERIVATION	# BYTES
application/pdf	559430
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Key	346492
application/pdf → XRefTable	303713
application/pdf → PDFObject → PSBytes → image/jpeg	120125
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFObjRef	112948
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFLiteral	101912
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PSInt	56427
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PDFObjRef	54860
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PSInt	49158
application/pdf → PDFObject → PDFList → PDFObjRef	37978
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PDFLiteral	30841
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFDictionary → KeyValuePair → Key	26643
application/pdf → PDFObject → FlateDecode → DecodedStream → PSBytes → application/zlib	25803
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PSFloat	17474
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PDFList → PSInt	14489
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PSBytes → text/plain	12653
application/pdf → PDFObject → FlateDecode	11446
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PSFloat	10991
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFDictionary → KeyValuePair → Value → PDFLiteral	10052
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PDFDictionary → KeyValuePair → Key	8250
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFDictionary → KeyValuePair → Value → PSFloat	7743
application/pdf → PDFObject → PSBytes → application/octet-stream	7328
application/pdf → PDFObject → PSBytes → image/jp2	6147
application/pdf → PDFObject → PSBytes → text/plain	5748
application/pdf → PDFObject → PDFDeciphered → image/jpeg	4900
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PSBytes → text/plain	4845
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFDictionary → KeyValuePair → Value → PSInt	4722
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PDFDictionary → KeyValuePair → Value → PDFLiteral	3518
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFDictionary → KeyValuePair → Value → PDFList → PSInt	3439
application/pdf → PDFObject → PDFDeciphered → text/plain	3361
application/pdf → PDFObject → PDFDictionary → KeyValuePair → Value → PDFList → PDFList → PDFObjRef	2936
application/pdf → PDFObject → PSBytes → image/x-portable-bitmap	2806

TABLE 4. THE PDF PARSE TREE DERIVATIONS FOR $M\mu$ PDF CONTAINING THE MOST BLIND SPOT BYTES. THESE PRIMARILY DESCEND FROM PDFOBJECT. ALMOST ALL OF THE DERIVATIONS DESCEND FROM THE PDFOBJECT NODE. THIS IS UNSURPRISING, SINCE PDF OBJECTS CAN CONTAIN STREAMS OF ARBITRARY BINARY DATA. PDF OBJECTS ALSO DO NOT NEED TO BE CONNECTED TO THE ROOT OF THE PDF DOCUMENT OBJECT MODEL, NOR DO THEY NEED TO BE USED IN ANY WAY FOR RENDERING.

- [3] H. Brodin, “How to avoid the aCropalypse,” <https://blog.trailofbits.com/2023/03/30/acropalypse-polytracker-blind-spots/>, March 29, 2023, accessed: March 29, 2023.
- [4] A. Albertini, “Abusing file formats; or, Corkami, the novella,” *The International Journal of Proof of Concept or GTF0*, vol. 0x07, no. 6, pp. 18–41, March 2015.
- [5] E. Sultanik and P. Teuwen, “Post scriptum: A schizophrenic ghost,” *The International Journal of Proof of Concept or GTF0*, vol. 0x13, no. 10, p. 71, October 2016.
- [6] J. Forristal, “Android fake ID vulnerability,” in *Proceedings of Black-Hat US*, August 2014.
- [7] D. Kaminsky, M. L. Patterson, and L. Sassaman, “PKI layer cake: New collision attacks against the global X.509 infrastructure,” in *Financial Cryptography and Data Security*, R. Sion, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 289–303.
- [8] M. Balduzzi, C. Torrano-Gimenez, D. Balzarotti, and E. Kirda, “Automated discovery of parameter pollution vulnerabilities in web applications,” in *Proceedings of the Network and Distributed System Security Symposium*, February 2011.
- [9] P. Chen and H. Chen, “Angora: Efficient fuzzing by principled search,” in *Proceedings of the IEEE Symposium on Security and Privacy*, 2018, pp. 711–725.
- [10] A. Albertini, T. Duong, S. Gueron, S. Kölbl, A. Luykx, and S. Schmiege, “How to abuse and fix authenticated encryption without key commitment,” in *Proceedings of the 31st USENIX Security Symposium*. Boston, MA: USENIX Association, August 2022. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity22/presentation/albertini>
- [11] C. Brant, P. Shrestha, B. Mixon-Baca, K. Chen, S. Varlioglu, N. El-sayed, Y. Jin, J. Crandall, and D. Oliveira, “Challenges and opportunities for practical and effective dynamic information flow tracking,” *ACM Computing Surveys*, vol. 55, no. 1, November 2021.
- [12] P. Yang, F. Kang, Y. Zhao, and H. Shu, “DRTaint: A dynamic taint analysis framework supporting correlation analysis between data regions,” *Journal of Physics: Conference Series*, vol. 1856, no. 1, p. 012013, April 2021. [Online]. Available: <https://doi.org/10.1088/1742-6596/1856/1/012013>
- [13] “DataFlowSanitizer,” <https://clang.llvm.org/docs/DataFlowSanitizer.html>, accessed: 2020-07-26.
- [14] E. J. Schwartz, T. Avgerinos, and D. Brumley, “All you ever wanted to know about dynamic taint analysis and forward symbolic execution (but might have been afraid to ask),” in *Proceedings of the IEEE Symposium on Security and Privacy*, 2010, pp. 317–331.
- [15] J. Cheney, A. Ahmed, and U. A. Acar, “Provenance as dependency analysis,” in *Database Programming Languages*, M. Arenas and M. I. Schwartzbach, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2007, pp. 138–152.
- [16] E. Sultanik, B. Larsen, and C. Harmon, “Two new tools that tame the treachery of files,” <https://blog.trailofbits.com/2019/11/01/two-new-tools-that-tame-the-treachery-of-files/>, November 1, 2019, accessed: January 12, 2020.
- [17] C. Harmon, B. Larsen, and E. Sultanik, “Toward automated grammar extraction via semantic labeling of parser implementations,” in *Proceedings of the Sixth Workshop on Language-Theoretic Security (LangSec)*. IEEE Symposium on Security and Privacy, 2021.

- [18] “[DFSan] change shadow and origin memory layouts to match MSan.” LLVM commit 45f6d5522f8d, <https://reviews.llvm.org/D104896?id=354633>, accessed: 2020-07-26.
- [19] M. Höschele and A. Zeller, “Mining input grammars from dynamic taints,” in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE 2016. New York, NY, USA: Association for Computing Machinery, 2016, pp. 720–725.
- [20] W. M. Khoo, “Taintgrind: a Valgrind taint analysis tool,” <https://github.com/wmkhoo/taintgrind>, accessed: March 2, 2020.
- [21] F. Bellard, “QEMU, a fast and portable dynamic translator,” in *Proceedings of the Annual USENIX Technical Conference*, ser. ATEC ’05. USA: USENIX Association, 2005, p. 41.
- [22] B. Dolan-Gavitt, J. Hodosh, P. Hulin, T. Leek, and R. Whelan, “Repeatable reverse engineering with PANDA,” in *Proceedings of the 5th Program Protection and Reverse Engineering Workshop*, ser. PPREW-5. New York, NY, USA: Association for Computing Machinery, 2015.
- [23] A. Davanian, Z. Qi, Y. Qu, and H. Yin, “DECAF++: Elastic Whole-System dynamic taint analysis,” in *22nd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2019)*. Chaoyang District, Beijing: USENIX Association, September 2019, pp. 31–45.
- [24] A. Karaman, “Measuring QEMU emulation efficiency,” <https://ahmedkrmn.github.io/TCG-Continuous-Benchmarking/Measuring-QEMU-Emulation-Efficiency/>, August 2020, accessed: 2020-07-26.
- [25] S. Nagy, A. Nguyen-Tuong, J. D. Hiser, J. W. Davidson, and M. Hicks, “Breaking through binaries: Compiler-quality instrumentation for better binary-only fuzzing,” in *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, August 2021, pp. 1683–1700. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity21/presentation/nagy>
- [26] F. Soudel and J. Salwan, “Triton: A dynamic symbolic execution framework,” in *Symposium sur la sécurité des technologies de l’information et des communications*, ser. SSTIC, Rennes, France, June 2015, pp. 31–54.
- [27] S. Poeplau and A. Francillon, “Symbolic execution with SymCC: Don’t interpret, compile!” in *29th USENIX Security Symposium (USENIX Security 20)*. USENIX Association, Aug. 2020, pp. 181–198. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity20/presentation/poeplau>
- [28] J. Clause, W. Li, and A. Orso, “Dytan: A generic dynamic taint analysis framework,” in *Proceedings of the 2007 International Symposium on Software Testing and Analysis*, ser. ISSTA ’07. New York, NY, USA: Association for Computing Machinery, 2007, pp. 196–206.
- [29] P. Wyatt, “Work in progress: Demystifying PDF through a machine-readable definition,” in *Proceedings of the Seventh Workshop on Language-Theoretic Security (LangSec)*. IEEE Symposium on Security and Privacy, 2021.
- [30] S. Garfinkel, P. Farrell, V. Roussev, and G. Dinolt, “Bringing science to digital forensics with standardized forensic corpora,” *Digital Investigation*, vol. 6, pp. S2–S11, 2009, proceedings of the Ninth Annual DFRWS Conference. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1742287609000346>
- [31] J. Wolf, “OMG WTF PDF—PDF ambiguity and obfuscation,” in *Proceedings of TROOPERS*, Heidelberg, Germany, March 2011.