



Lifting Network Protocol Implementation to Precise Format Specification with Security Applications

Qingkai Shi
Purdue University
West Lafayette, USA
shi553@purdue.edu

Junyang Shao
Purdue University
West Lafayette, USA
shao156@purdue.edu

Yapeng Ye
Purdue University
West Lafayette, USA
ye203@purdue.edu

Mingwei Zheng
Purdue University
West Lafayette, USA
zheng618@purdue.edu

Xiangyu Zhang
Purdue University
West Lafayette, USA
xyzhang@cs.purdue.edu

ABSTRACT

While inferring protocol formats is critical for many security applications, existing techniques often fall short of coverage, inasmuch as almost all of them are in a fashion of dynamic analysis and driven by a limited number of network packets. If a feature is not present in the input packets, the feature will be missed in the resulting formats. To tackle this problem, we develop a novel static program analysis that infers protocol message formats from the implementation of common top-down protocol parsers. However, to achieve the triecta of coverage, precision, and efficiency, we have to address two challenges, namely path explosion and disordered path constraints. To this end, our approach uses abstract interpretation to produce a novel data structure called the abstract format graph. The graph structure delimits precise but costly operations to only small regions, thus ensuring precision and efficiency at the same time. Our inferred formats are of high coverage and precisely specify both field boundaries and semantic constraints among packet fields. Our evaluation shows that we can infer formats for a protocol in one minute with over 95% precision and recall, much better than four baselines. Our inferred formats can substantially enhance existing protocol fuzzers, improving the coverage by 20% to 260% and discovering 53 zero-days with 47 assigned CVEs. We also provide case studies of adopting our inferred formats in network traffic auditing and network intrusion detection.

CCS CONCEPTS

• Security and privacy → Web protocol security; Software reverse engineering; • Software and its engineering → Automated static analysis.

KEYWORDS

Reverse engineering; network protocols; protocol formats.

ACM Reference Format:

Qingkai Shi, Junyang Shao, Yapeng Ye, Mingwei Zheng, and Xiangyu Zhang. 2023. Lifting Network Protocol Implementation to Precise Format Specification with Security Applications. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security (CCS '23)*, November 26–30, 2023, Copenhagen, Denmark. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3576915.3616614>

1 INTRODUCTION

Network protocols define how computing systems are connected. Thus, security vulnerabilities in network protocols may have devastating consequences. For example, the WannaCry attack, which was caused by a protocol vulnerability, led to over \$8 billion loss across 150 countries [1]. To aid automated security analysis for network protocols, a formal specification of packet formats is often mandatory — it facilitates the generation of legitimate network packets for security testing [41, 62]; it is the foundation for protocol model checking [24, 60] and formal verification [25]; and it can guide automated code generation with strong guarantees [69].

However, while protocols may have their specification documents in natural languages, formal, or machine-readable, packet formats are often not available, and even when they are, they may be incomplete or inaccurate [22]. Therefore, automatically inferring formal protocol formats is of importance. There are three typical scenarios. First, the protocol implementation is not accessible but network packets are available. In this case, network trace analysis [33, 46, 47, 57, 71, 72, 77] are proposed. They leverage statistical analysis and machine learning to infer how a packet can be divided into fields. Since the underlying techniques have inherent uncertainty, the quality of inferred formats tends to be insufficient to drive many security applications such as protocol fuzzing. We call them the *category-one* techniques.

In the second scenario, the executable code of a protocol and a set of valid packets are available. Dynamic program analysis [26, 27, 35, 54–56, 74] then traces how individual packet bytes are propagated when running the code on the provided packets. The protocol formats then can be inferred from the data or control flow relations collected at runtime. For example, a typical rule to infer a raw data field is that consecutive bytes in the field are accessed by the same instruction [26, 35, 54]. These techniques can precisely infer the syntax of provided packets and denote the state of the art. In some cases (e.g., [35]), semantics constraints, for example,



This work is licensed under a Creative Commons Attribution International 4.0 License.

CCS '23, November 26–30, 2023, Copenhagen, Denmark
© 2023 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0050-7/23/11.
<https://doi.org/10.1145/3576915.3616614>

those describing correlations across packet fields like the sum of fields A and B cannot exceed a certain threshold, can be inferred as well. However, the inferred formats are often incomplete when the provided packets do not have good coverage of all possible formats. We call them the *category-two* techniques.

We focus on the third scenario, in which the source code of a network protocol is available. As we will demonstrate in §6, open-source protocol implementations have many zero-day vulnerabilities. Without precise formal formats, existing protocol fuzzers such as BooFuzz [11] can hardly find them. In a recent work FRAME-SHIFTER [44], critical bugs were found by fuzzing HTTP/1 and HTTP/2, whose implementations are publicly available. While the authors manually crafted the protocol formats, automatic format inference can generalize their method to other protocols. In network traffic auditing and attack detection, e.g., using Wireshark [15] and Snort [10], substantial manual efforts are still needed to write dedicated protocol parsers for Wireshark and Snort even when a protocol is open-sourced. In contrast, with our inferred formal formats, Wireshark and Snort can be automatically extended.

We develop a static program analysis to produce packet formats, including both syntax and semantics, from the source code of a protocol parser, which often parses a packet in a top-down fashion [18]. We call it a protocol lifting technique, belonging to *category-three*. We resort to static analysis in order to address the coverage problem in dynamic analysis. Meanwhile, high accuracy can be achieved as it adopts a path-sensitive analysis. We produce BNF-like protocol formats. While BNF [32] is a common language to describe syntax, we enhance it to include semantic constraints across protocol fields. As we will show in §4, lifting source code to protocol formats is highly challenging. First, the traditional data-flow analysis that aggregates analysis results of multiple program paths at their joint point yields very poor results, whereas path-sensitive analysis that considers individual paths separately is prohibitively expensive due to path explosion. Second, the inferred formats are mostly out of order for human interpretation, which is highly undesirable as humans are important consumers of the formats in security applications.

To address the challenges, we develop a novel static analysis. In particular, we develop abstract interpretation rules that can derive an abstract format graph (AFG) from the source code. AFG can be considered as a transformed control flow graph. It precludes statements that are irrelevant to packet formats. It further merges program subpaths that are irrelevant to formats so that path-sensitive analysis is not performed on the merged places. Meanwhile, it retains sufficient information such that a localized but precise path-sensitive analysis can be performed on the unmerged parts of the graph. Therefore, it mitigates the path-explosion problem without losing analysis accuracy. The AFG is further unfolded and reordered to generate BNF-style production rules and first-order-logic formulas that describe semantic constraints across protocol fields. In summary, we make the following four contributions:

- We develop an abstract interpretation method that produces a novel representation, namely the abstract format graph, to facilitate format inference.
- We propose a localized graph unfolding algorithm that can perform precise path-sensitive analysis in small AFG regions to significantly mitigate path explosion.

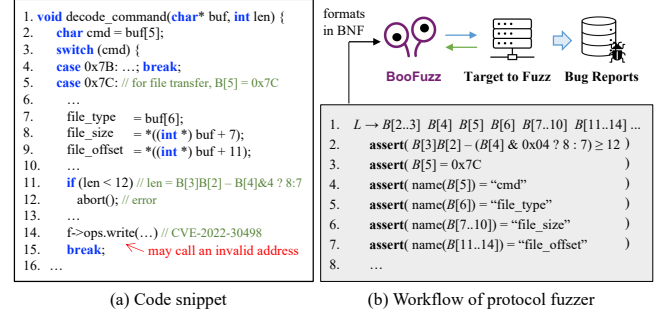


Figure 1: (a) Simplified code that parses the file-transfer command. (b) The typical workflow of protocol fuzzers with a snippet of the format inferred by Netlifter, in which the first row is a BNF production rule denoting syntax (e.g., field partitioning) and the remaining denote semantic constraints.

- We devise a graph reordering algorithm that translates an unfolded AFG to the commonly-used BNF so that our inferred formats can be widely applied in practice.
- We implement our approach as a tool, namely Netlifter, to infer packet formats from protocol parsers written in C. We evaluate it on a number of protocols from different domains. Netlifter is highly efficient as it can infer formats in one minute. Netlifter is highly precise with a high recall as its inferred formats uncover $\geq 95\%$ formats with $\leq 5\%$ false ones. In contrast, the baselines, often miss $>50\%$ of formats and, sometimes, produce $> 50\%$ false ones. We use the inferred formats to enhance grammar-based protocol fuzzers, which are improved by $20\% \sim 260\%$ in terms of coverage and detect 53 zero-day vulnerabilities with 47 assigned CVEs. Without our formats, only 12 can be found. We also provide case studies of adopting our approach in traffic analysis and intrusion detection. Netlifter is publicly available [17].

2 MOTIVATION

We use an open-source protocol, namely Open Supervised Device Protocol (OSDP), to illustrate the limitations of existing methods and how our technique can facilitate various security applications. OSDP is an access control communications standard developed by the Security Industry Association to improve interoperability among access control and security products. Although it is an open-source protocol, its full specification is not publicly available. The only available document [2] lacks many details. For instance, it includes the formats for only 7 out of the 27 supported commands.

Figure 1(a) shows a part of the packet parsing function of OSDP, which contains a bug found by a protocol fuzzer enhanced by our approach. The variable buf is a byte array representing the OSDP packet, and we use $B[i]$ to represent the $(i + 1)$ th byte. The bug at Line 14 may invoke an invalid function pointer $f \rightarrow ops.write$, which could lead to a crash or be exploited for DoS or ROP attacks. Many ways can avoid such attacks. The first one is to use fuzzers to find such bugs and have them fixed before exploitation. The second is to provide OSDP support in network traffic analysis and attack detection tools such that the attack can be analyzed and further prevented. However, existing methods fall short as discussed below.

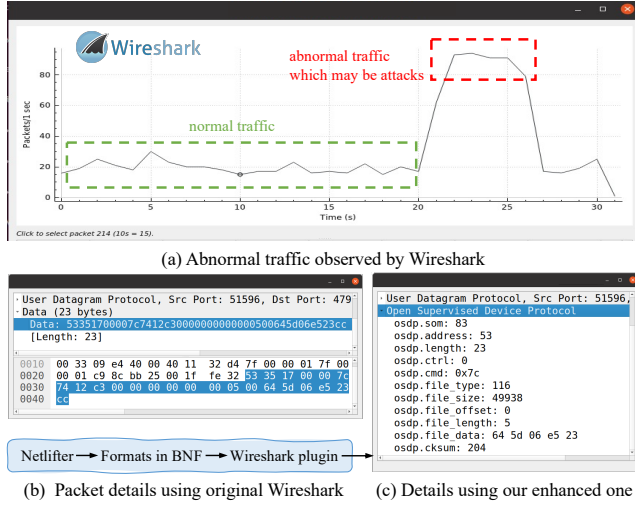


Figure 2: Network traffic auditing via Wireshark.

Standard Network Fuzzing Can Hardly Find the Zero-day.

Different from stand-alone application fuzzers such as AFL [16], network fuzzers, such as BooFuzz [11], often operate in a client-server architecture. The server runs the target protocol implementation. The client leverages grammar-based fuzzing to generate packets as per the formats, send the packets to the server, receive responses, and generate new packets to fuzz the target. However, the effectiveness of these fuzzers hinges on the protocol formats. When the formats are not available like in our OSDP case, they quickly degenerate into traditional greybox fuzzers that arbitrarily mutate bits or bytes. Such mutated packets can hardly pass many input validity checks in the code. For example, in Figure 1(a), to expose the bug, a fuzzer has to get through the check at Line 11, which is a complex relation across multiple fields as shown in the comment. As a result, standard network fuzzers fail to find the bug when an imprecise or incomplete format of OSDP is provided.

Lack of Support for OSDP in Wireshark and Snort. We can also rely on network traffic analyzers, e.g., Wireshark [15], and attack detection tools, e.g., Snort [10], to ensure security. However, both Wireshark and Snort do not support OSDP. Assume Wireshark is deployed at the gateway. It detects abnormal traffic as highlighted in the red box in Figure 2(a). Note that in the diagram the x-axis denotes the elapsed time and the y-axis denotes the amount of traffic per second. However, the traffic is not interpretable for Wireshark as OSDP is not supported. Instead, the OSDP packets are treated as raw data bytes as shown in Figure 2(b). Thus, it is hard to analyze the packet details and determine which device launches the attack.

2.1 Limitations of Existing Techniques

A way to address the aforesaid defense insufficiency is to infer the protocol formats. As discussed in §1, existing techniques fall into categories one and two. Category-one infers formats from a set of network packets. For example, a recent method NetPlier [77] leverages a probabilistic analysis on network packets to determine a keyword field, i.e., the field identifying the packet type, by computing the probabilities of each byte offset. Once the keyword field is determined, it clusters packets according to the value of the field

and applies multi-sequence alignment to derive message format. However, real-world network packets suffer from all sorts of distribution biases, e.g., lacking some kinds of messages due to their rare uses in practice, leading to sub-optimal results.

For instance, NetPlier partitions the first four bytes of an OSDP packet as `| 0x53 0xff | 0x29 | 0x00 | ... |`, which mistakenly places the first two bytes into the same field and splits $B[2]$ and $B[3]$ into two different fields while $B[2..3]$ — with the value of $B[3]B[2]$ that represents a two-byte integer with $B[3]$ the most significant byte and $B[2]$ the least — should be a single field representing the packet length. However, since most input packets are shorter than 255, $B[3]$ is always zero while $B[2]$ has different values in different packets. Thus, these two bytes follow different distributions in the packet samples, and NetPlier incorrectly regards them as separate fields. Moreover, NetPlier does not infer semantic constraints such as the condition at Line 11 in Figure 1(a). Such imprecise formats prevent a grammar-based protocol fuzzer from finding the zero-day and fail to enhance Wireshark and Snort.

Category-two methods dynamically analyze protocol execution using a set of input packets. AutoFormat [54] is a representative. It leverages the observation that most packet parsers utilize top-down parsing such that they invoke a function to parse a sub-structure. Therefore, the dynamic call graph in parsing a packet discloses its structure. However, the function call hierarchy may not be sufficiently fine-grained to disclose detailed packet formats. Similar to NetPlier, it does not infer semantic constraints across fields, such as the one at Line 11 in Figure 1(a). As dynamic analysis, the inferred format may be incomplete, depending on the coverage of the input packets that drive the dynamic analysis. For instance, in our evaluation, Autoformat misses 15 out of the 27 packet types because these types of packets do not appear in regular workloads.

Some category-two techniques, e.g., Tupni [35], can precisely infer semantic constraints among packet fields. However, as dynamic analyses, they suffer from the innate coverage problem. As per our results, the inference results of Tupni may miss >50% of possible formats. The problem is that if the program executions analyzed by Tupni do not cover the file-transfer command, i.e., Lines 5-15 in Figure 1(a), Tupni will not generate formats for the command. Without the formats, it is hard for a fuzzer to generate packets that can pass the validity check at Line 11 and expose the bug at Line 14.

2.2 Our Solution and Security Applications

Observing that the source code discloses substantial information about packet formats, we propose a category-three method that lifts the source code of OSDP to the protocol formats. For instance, Line 5 of the code in Figure 1(a) indicates that the command code for file transfer is `0x7c`. Line 7 indicates that $B[6]$ is a field representing the file type to transfer. Lines 8-9 load two four-byte integers to variables `file_size` and `file_offset` and, thus, indicate that there are two four-byte fields, one from $B[7]$ to $B[10]$ and the other from $B[11]$ to $B[14]$, meaning the size and the offset of the file to transfer, respectively. In addition to the syntactic information (e.g., field partitioning), the code also discloses the semantic relations across fields. For instance, the if-statement at Line 11 implies a cross-field constraint dictating that if $B[4] \& 4 = 0$, a valid packet must satisfy the constraint $B[3]B[2] - 7 \geq 12$ or, otherwise, $B[3]B[2] - 8 \geq 12$.

We extract the above syntactic and semantic information via static analysis and produce a BNF-style production rule in Figure 1(b). Our lifted formats are both precise and of high coverage. In terms of precision, we precisely identify each field and its name as shown in Figure 1(b) and, meanwhile, also specify the field constraints as first-order-logic formulas. In terms of coverage, since we do not rely on any input packets like category-one and category-two techniques, any format included in the source code will be inferred. We can use the lifted formats to support many applications.

Application 1: Finding Zero-days by Network Fuzzing. We leverage a theorem prover such as Z3 [37] to produce valuations for individual packet fields, such as the $B[i]$'s in Figure 1(b), which satisfy the semantic constraints. The generated packets can pass the check at Line 11 in Figure 1(a), thereby enabling the discovery of the CVE at Line 14. In addition, our inferred packet format is of high coverage and allows the fuzzer to generate diverse packets to improve test coverage. In particular, the vulnerable code can only be reached when the packet is a file-transfer command, i.e., the 0x7c branch of the switch statement (Line 5). If the format is not covered, the chance that a fuzzer can mutate a packet of different types to a valid file-transfer command is very slim. Our evaluation shows that the lifted formats can improve the coverage of fuzzing by 20% to 260% and allow us to detect 41 more zero-days, compared to using the formats inferred by category-one and category-two methods. Note that we do not claim direct contributions to fuzzing. Instead, our approach is orthogonal to existing fuzzing methods that rely on packet formats.

Application 2: Network Traffic Auditing. Wireshark is the foremost protocol analyzer to ensure network security for hundreds of protocols [15]. Supporting a new protocol in Wireshark can be achieved by providing an extension, which is usually a library to parse protocol packets. We develop an extension generator that takes a lifted format as input and generates the corresponding Wireshark extension. Figure 2(c) shows that with the generated extension, Wireshark can look inside an OSDP packet sampled from the abnormal traffic. Our lifted format provides not only precise packet syntax but also informative field names extracted from variable names. With Wireshark, we observe that all packets during the abnormal traffic have the field `osdp.address=35` and `osdp.cmd=0x7c`, indicating they are all from a device with the id #35 via the file-transfer commands. As will be shown in our evaluation, category-two approaches miss over 50% of possible fields. Building extensions from these incomplete formats would render Wireshark incapable of processing many packets. In addition, they can hardly provide field names that are as informative as the ones we can provide.

Application 3: Network Intrusion Detection. We put discussions of this application in an extended version of this paper [66].

3 BACKGROUND AND OVERVIEW

Protocol Format vs. Protocol Specification. Generally, the specification of a protocol consists of protocol formats and protocol state machines [61]. Protocol formats are often specified using a grammar in BNF, which specifies how a network packet, i.e., a bit or byte stream, can be dissected into multiple segments, i.e., fields, and specifies the semantic constraints the fields need to satisfy. For example, Figure 3(b) shows a typical BNF-style format of OSDP

packets. The productions specify how an OSDP packet can be divided into multiple fields such as *som*, *address*, and so on. Like many previous works [26, 27, 33, 35, 42, 43, 46, 47, 54–57, 71, 72, 74, 77], Netlifter focuses on inferring the protocol formats.

Protocol state machines, on the other hand, specify the state transitions of a network server upon receiving or sending a specific network packet. For instance, a TCP server may transition from the state SYN-SENT, meaning it waits for a matching connection, to the state ESTABLISHED, meaning a connection has been established, upon receiving an acknowledgment packet. There have been many works focusing on state machine inference [20, 29, 30, 34, 38, 49, 51, 58, 65, 73, 77, 78]. Typically, these works accept a set of network packets as input and produce a protocol state machine. While Netlifter only focuses on the formats, since the formats are sufficient to produce a number of packets, one can feed the packets into the aforementioned techniques for state machine inference.

Static Program Analysis. Static analysis processes a program without executing it and often has various forms such as dataflow analysis and abstract interpretation. As pointed out by Cousot and Cousot [31], while dataflow analysis and abstract interpretation are in different forms, they are equivalent when used to compute sound results. Basically, this is because both of them use abstract values to approximate program behavior. A complete/sound analysis uses abstract values, which are often formulas over predefined symbols, to under-/over-approximate concrete values that may assign to program variables at runtime. For instance, in our work, the abstract value of a variable is a formula over bytes, e.g., $B[0], B[1], \dots$, in a network packet. Here, $B[i]$ is a byte ranging from 0x00 to 0xFF and, thus, over-approximates all possible values of the $(i + 1)$ th byte.

A static analysis is often defined by a set of transfer functions and merging functions over abstract values. A transfer function specifies how we compute an abstract value when visiting a program statement. For instance, given the abstract values of two variables, e.g., $x = B[1], y = B[2]$, the transfer function of the statement $z = x + y$, accepts the abstract values of x and y as the input and outputs the abstract value of z , which typically is $B[1] + B[2]$.

When a variable is assigned multiple abstract values in two program paths, at the joint point of the two paths, we use a merging function to compute a merged abstract value. For instance, assume that we compute $x = B[1]$ and $x = B[2]$ in two different paths and Θ is an operator returns either of its operands. At the joint point, the merged abstract value of x could be $\Theta(B[1], B[2])$. This value is sound as it over-approximates the value of x , saying that x is either $B[1]$ or $B[2]$. However, it is not complete or not precise, because it loses the path information, i.e., from which path $x = B[1]$ (or $B[2]$).

A static analysis can be performed with varying degrees of precision. In this work, we choose to perform a path-sensitive static analysis, which is of high precision as it can distinguish abstract values from different paths. To this end, one often needs to enumerate all paths in a program, just like symbolic execution [45], which, however, suffers from the notorious path-explosion problem due to the exponential number of paths in a program. In this work, we fend off this problem by introducing a special merging operator, Θ_{κ} , as explained later.

Input of Netlifter. The input of our static analyzer is the source code of a *top-down protocol parser* [18] written in C. A top-down

<pre> 1. void decode_command(char* buf, int len) { 2. char cmd = buf[5]; 3. switch (cmd) { 4. case 0x7B: ...; break; 5. case 0x7C: // for file transfer, B[5] = 0x7C 6. ... 7. file_type = buf[6]; 8. file_size = *((int*) buf + 7); 9. file_offset = *((int*) buf + 11); 10. ... 11. if (len < 12) // len = B[3]B[2] - B[4]&4 ? 8:7 12. abort(); // error 13. ... 14. } 15. void decode_packet(char* buf, int blen) { 16. char som = buf[0]; 17. assert(som == 0x53); 18. char address = buf[1]; 19. short length = buf[2] (buf[3] << 8); 20. char ctrl = buf[4]; 21. length = length - ctrl & 0x4 ? 8 : 7; 22. decode_command(buf, length); 23. }</pre>	<pre> OSDP → som address CMD CMD → length ctrl 0x7b ... CMD → length ctrl 0x7c filesize fileoffset → ... (b) Conventional BNF-style formats 24. S → B[0] B[1] L 25. assert(B[0] = 0x53) 26. assert(name(B[0]) = "som") 27. assert(name(B[1]) = "address") ... 28. # for cmd = 0x7B 29. L → B[2..3] B[4] B[5] 30. # for cmd = 0x7C 31. L → B[2..3] B[4] B[5] B[6] B[7..10] B[11..14] ... 32. assert(B[3]B[2] - (B[4] & 0x04 ? 8 : 7) ≥ 12) 33. assert(B[5] = 0x7C) 34. assert(name(B[2..3]) = "length") 35. assert(name(B[4]) = "ctrl") 36. assert(name(B[5]) = "cmd") 37. assert(name(B[6]) = "file_type") 38. assert(name(B[7..10]) = "file_size") 39. assert(name(B[11..14]) = "file_offset")</pre>
(a) Simplified OSDP code	(c) Inferred OSDP formats

Figure 3: Extended example for OSDP.

parser applies each production rule in a BNF-style format to incoming bytes of the network packet, working from the left-most symbol of a production rule and then proceeding to the next production rule for each non-terminal symbol encountered [19]. Given the parsing function of a protocol, e.g., *parse(char* buf, int len)* { ... }, the user annotates the parameters, i.e., the buffer variable, *buf*, that contains the network packet to parse, and the integer variable, *len*, which stands for the packet length. Except for the two annotations, Netlifter is fully automated.

Example. Figure 3 extends the example in Figure 1. It shows a simplified OSDP parser starting from Line 15. The packet is a byte array stored in *buf* and the array length is *blen*. The user needs to annotate the two variables. The parser with the annotations is the input of Netlifter. In Lines 16-21, the parser loads the first five bytes into the variables *som*, *address*, *len*, and *ctrl*, where *som* stands for “start of message” and is used to identify OSDP packets. The remaining code invokes the function *decode_command* to parse an OSDP command as explained in Figure 1. □

Output of Netlifter. The output of Netlifter is defined below. Our protocol format is similar to common BNF so that it aligns well with existing standards in formally describing protocol formats.

Definition 1 (Protocol Format). The format includes syntax and semantics. The syntax is denoted by production rules in BNF, where each rule is a sequence of consecutive bytes. Semantics is described by non-recursive first-order-logic (FOL) constraints with two special functions, *name(...)* and *repeat(...)*, which are explained in the example below. The format satisfies three properties:

- (1) Each terminal symbol in the grammar is either $B[i]$ or $B[i..j]$, which is a bit-vector standing for the $(i+1)$ th byte or a range of bytes from $B[i]$ to $B[j]$;
- (2) Each production rule is associated with a set of assertions that assert FOL constraints over the terminals in this rule. The constraints must not conflict with each other.
- (3) Each assertion contains only a single atomic constraint that does not contain any connectives \wedge or \vee .

Example. Figure 3(b) shows a typical BNF-style format of OSDP, which is often manually constructed. The output format of Netlifter

is shown in Figure 3(c), which closely resembles the manually constructed BNF in (b). The first rule in (c) resembles the first rule in (b), where we correctly determine that the first two bytes, $B[0]$ and $B[1]$, are two separate fields, corresponding to the fields *som* and *address*. Similarly, the second and third rules in (c) resemble the two CMD rules in (b), where, besides single-byte fields, we also correctly determine multi-byte fields including $B[2..3]$, $B[7..10]$, and $B[11..14]$, corresponding to the fields *length*, *filesize*, and *fileoffset*.

The output format also associates each rule with two kinds of assertions. One kind, such as Line 25 and Lines 32-33, specifies the semantic constraints among packet fields. They are inferred from branching conditions in the code. When we infer a constraint including a value like $B[3]B[2]$, it indicates a two-byte field with $B[3]$ the most significant byte and $B[2]$ the least. In other words, in addition to *field boundaries*, our format also expresses the *endianness*, whereas the standard BNF cannot. Netlifter also describes semantic constraints not expressible in standard BNF such as the one in Line 32. All constraints have the *bit-level precision*. For instance, the expression $B[4] \& 4$ in Line 32 computes the third bit of $B[4]$.

The other kind, such as Lines 26-27 and Lines 34-39, specifies the field names, which provide *high-level field semantics* for us to understand the format. In addition to those in the example, we also produce many other names such as *name(B[i..j])* = ‘timestamp’/‘checksum’ to indicate a timestamp/checksum field. As explained later, we infer such high-level semantics using the names of program variables or library APIs. □

In addition to the example above, we elaborate on two places where our format is more expressive than the standard BNF. First, our format specifies *direction fields*. A direction field locates another field and is often a length field, whose value encodes the variable length of a target field [27]. For instance, we may produce a production rule $S \rightarrow B[0]B[1..B[0]]$, where $B[1..B[0]]$ is a variable-length field whose length is determined by the direction field $B[0]$.

Second, our format can also specify *repetitive fields* and how many times a field repeats. As an example, we may produce the production rule $S \rightarrow B[0]B[1..2]B[3]$ with three assertions: (1) *assert*($B[0] = B[3] = 0x00$), (2) *assert*($B[1]B[2] \neq 0x0000$), and (3) *assert*(*repeat*($B[1..2]$) = 3). The first assertion constrains the first and last byte. The second constrains the field $B[1..2]$ in middle. The third states that the middle field repeats three times. When generating packets based on the rule, we first generate a packet satisfying the first two assertions, e.g., $0x00\ 0x0001\ 0x00$. Due to the third assertion, we insert another two fields satisfying the same constraints as $B[1..2]$, e.g., $0x00\ 0x0001\ 0x0011\ 0x0101\ 0x00$.

4 TECHNICAL CHALLENGES

We discuss two prominent challenges and our ideas for addressing them. The discussion is driven by a crafted protocol parser in Figure 4(a). Figure 4(e) shows an ideal format inferred from the parser, where a packet is dissected into three segments, i.e., L_1 , $L_2|L_3$, and $L_4|L_5$. The segment L_1 has two fields, a *code* field (reflected by Lines 10-11 in the code) and a *state* field (reflected by Lines 14-16 in the code). L_2 and L_3 consist of a single field located at byte offset 3 and differ only in the field value (reflected by Line 13 in the code). L_4 and L_5 consist of three single-byte fields and differ in the semantic constraints (reflected by Lines 3-8 in the code).

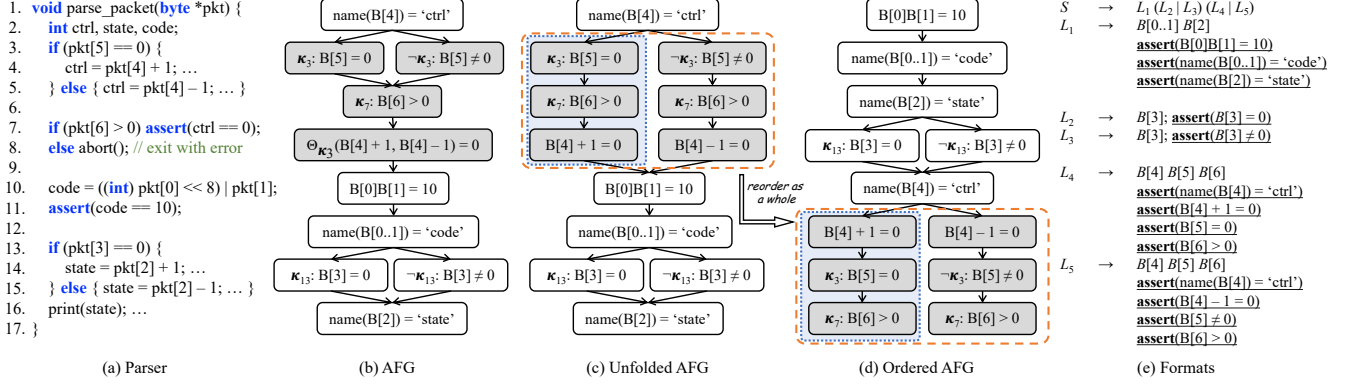


Figure 4: A crafted running example.

Challenge 1: Insufficiency of Traditional Static Analysis. Traditional static analysis is path-insensitive and merges analysis results from different paths at their joint point to achieve scalability. As introduced before, such merging yields over-approximation and incurs low precision. For example, the abstract values of *ctrl* from the two branches at Lines 4 and 5, respectively, are merged at Line 6, yielding $ctrl = B[4] + 1 \vee ctrl = B[4] - 1$. As such, we lose the correlation between $B[4]$ and $B[5]$ as the precise value of *ctrl* should depend on the value of $B[5]$ due to the if-statement at Line 3. In consequence, the resulting format will lose the correlations between $B[4]$ and $B[5]$, while in the ideal format shown in Figure 4(e), the production rules L_4 and L_5 include such correlations, i.e., $B[4] + 1 = 0 \Leftrightarrow B[5] = 0$ and $B[4] - 1 = 0 \Leftrightarrow B[5] \neq 0$.

A typical solution is to use a path-sensitive static analysis that separately analyzes individual paths and does not merge results from multiple branches. Lifting is thus reduced to enumerating paths, each constituting a production rule. In our example, there are four paths that denote valid packets, i.e., (P1) $\dots \rightarrow 4 \rightarrow \dots \rightarrow 14 \rightarrow \dots$, (P2) $\dots \rightarrow 5 \rightarrow \dots \rightarrow 14 \rightarrow \dots$, (P3) $\dots \rightarrow 4 \rightarrow \dots \rightarrow 15 \rightarrow \dots$, and (P4) $\dots \rightarrow 5 \rightarrow \dots \rightarrow 15 \rightarrow \dots$. Thus, the lifted format has four rules, each of which corresponds to a path constraint. For example, the format for the path P1 is shown below.

$S \rightarrow B[0..1] B[2] B[3] B[4] B[5] B[6]$
 $\text{assert}(B[5] = 0); \text{assert}(B[6] > 0); \text{assert}(B[4] + 1 = 0);$
 $\text{assert}(B[0]B[1] = 10); \text{assert}(B[3] = 0);$

Enumerating program paths incurs the notorious path-explosion problem, which has two consequences: (1) the analysis is not scalable and (2) the lifted format has an explosive number of semantic constraints. For example, due to path explosion, KLEE [28], a state-of-the-art path-sensitive static analyzer, cannot finish analyzing Linux's implementation of L2CAP, a Bluetooth protocol containing a few thousand lines of code, within twelve hours.

Solution 1: Localized Path-Sensitive Analysis. We observe that in lifting, path-sensitivity is only needed in certain places. In our example, we only want to analyze the sub-paths $3 \rightarrow 4 \rightarrow 7$ and $3 \rightarrow 5 \rightarrow 7$ separately such that we can generate production rules L_4 and L_5 in Figure 4(e). The criterion to determine a local code region for path-sensitive analysis is that the path conditions within the region and their branches have inter-dependencies. For example, the condition at Line 3 determines the value of *ctrl*, which is

checked inside the true branch at Line 7, allowing Lines 3-8 to form a region for path-sensitive analysis. In contrast, Lines 10-16 have no dependencies on Lines 3-8 and, thus, are considered separately. In §5, we will discuss how we use a new selection operator and a novel representation called the abstract format graph (AFG) to identify such regions for localized path-sensitive analysis.

Challenge 2: Handling Out-of-Order Fields. Protocol parsers may not parse network packets in strict byte order. Hence, if a naïve lifting algorithm directly derives format from code, for example, generating production rules following the order that the bytes are accessed along program paths, the resulting format may have out-of-order fields, which do not comply with the BNF standard. For example in Figure 4(a), the parser parses bytes $B[4..6]$ before $B[0..3]$. Therefore, we have to break the program order. This requires us to reorder the bytes such that they follow the byte order while not violating program semantics. For example, in Lines 3-8 in Figure 4(a), the access of *pkt[4]* occurs after that of *pkt[5]*. One cannot simply relocate Line 4 and the else branch in Line 5 to in front of Line 3, because the resulting program is broken as shown in the following.

$ctrl = pkt[4] - 1; ctrl = pkt[4] - 1; \text{if} (pkt[5] == 0) \dots$

Solution 2: Graph-Based Reordering. We propose to first abstract the code to the aforementioned AFG that models only the packet format-related behaviors and precludes the rest. As such, we do not need to transform the program which is complex and unnecessary. An algorithm is developed to ensure dependencies can be respected during reordering.

5 DESIGN

Figure 4(a-e) presents the workflow of Netlifter. Given the code of a protocol, abstract interpretation is performed to construct an abstract format graph (AFG). Path-sensitive analysis is performed in selected local regions of AFG to produce an unfolded AFG, which is further reordered and post-processed to produce the lifted formats.

5.1 Abstract Format Graph

AFG is a directed acyclic graph representing first-order-logic constraints. The AFG of a constraint ρ is inductively defined by $\text{AFG}(\rho)$:

- (1) $\text{AFG}(\text{atomic constraint}) = \text{Vertex}(\text{atomic constraint})$;
- (2) $\text{AFG}(\rho_1 \wedge \rho_2) = \text{AFG}(\rho_1) \bowtie \text{AFG}(\rho_2)$;
- (3) $\text{AFG}(\rho_1 \vee \rho_2) = \text{AFG}(\rho_1) \uplus \text{AFG}(\rho_2)$.

Generally, a vertex of AFG is an atomic constraint that does not contain any connectives \wedge or \vee , and an edge means logical conjunction. In the definition, the first rule returns a single vertex for any atomic constraint. The second creates a graph for conjunction by connecting all exit vertices (vertices without outgoing edges) of $\text{AFG}(\rho_1)$ to all entry vertices (vertices without incoming edges) of $\text{AFG}(\rho_2)$. The third creates a graph for disjunction by simply creating a union of the two graphs, which contains the vertices and edges from both. The following lemma states the equivalence relation between the graph $\text{AFG}(\rho)$ and the constraint ρ . In other words, $\text{AFG}(\rho)$ is an equivalent graphic representation of the constraint ρ . All proofs of our lemmas can be found in an extended version of this paper [66].

LEMMA 5.1. *Given $\text{AFG}(\rho)$ with n paths, we have $\rho \equiv \bigvee_{i=1}^n \rho_i$ where each ρ_i equals the conjunction of all constraints in an AFG path.*

Example. Consider the constraint $\rho \equiv (a \vee b) \wedge c \wedge (d \vee e)$. By definition, $\text{AFG}(\rho)$ is a directed graph with five nodes, which respectively correspond to the five atomic constraints a, b, c, d , and e . The AFG also contains four edges respectively from a to c , from b to c , from c to d , and from c to e . The AFG has four paths, respectively representing four constraints, $\rho_1 = a \wedge c \wedge d$, $\rho_2 = a \wedge c \wedge e$, $\rho_3 = b \wedge c \wedge d$, and $\rho_4 = b \wedge c \wedge e$. Apparently, we have $\rho_1 \vee \rho_2 \vee \rho_3 \vee \rho_4 \equiv \rho$. Thus, we say the $\text{AFG}(\rho)$ is an equivalent graphic representation of the constraint ρ . \square

5.2 Abstract Interpretation

The static analysis derives an AFG denoting path constraints related to the packet format. It features a new selection operator at the joint point of branches, which enables localized path-sensitive analysis.

Abstract Language. For clarity, we use a C-like language in Figure 5 to model our target programs. A program in the language has an entry function that parses an input network packet, pkt , which is a byte array. The parsing function often has a parameter specifying the packet length, len , to avoid out-of-bounds access during parsing. The language contains assignments, binary operations, statements that read bytes from the packet, assertions, branching, and sequencing. Each branching statement is labeled by a unique identifier κ . Although we do not include function calls or returns for discussion simplicity, our system is inter-procedural as a call statement is equivalent to a list of assignments from the actual parameters to the formals, and a return statement is an assignment from the return value to its receiver. The language includes statements reading bytes from the packet but does not include statements that store values into the packet. This is because, for parsing purposes, the input packet is often read-only. Note that the abstract language serves for demonstrating how we address the challenges discussed in §4. Thus, for simplicity, we abstract away some common program structures, e.g., pointers and loops, from the language. Dealing with these structures is not our technical contribution. In §5.5, we discuss how we handle them in our implementation.

Abstract Domain. An abstract value of a variable represents all possible concrete values that may be assigned to the variable during program execution. The abstract domain specifies the limited forms of an abstract value. In our analysis, the abstract value of a variable v is denoted as \tilde{v} and defined in Figure 6. An abstract value could be a constant or a special value length that represents the packet

```

Function F      := parse(pkt, len) { S; }
Statement S    := v1 ← v2                ::assign
                  | v1 ← v2 ⊕ v3          ::binary
                  | v1 ← pkt[v2]          ::read
                  | assert(v1)            ::assertion
                  | if_κ (v) {S1; } else {S2; } ::branching
                  | S1; S2                ::sequencing

```

$\oplus \in \{\wedge, \vee, +, -, >, <, =, \neq, \dots\}$

Figure 5: Language of target programs.

```

Abstract Value  $\tilde{v}$  := c                    ::constant
                  | length                ::packet length
                  | B[ $\tilde{v}$ ]                 ::byte in packet
                  |  $\Theta_\kappa(\tilde{v}_1, \tilde{v}_2)$       ::selection
                  |  $\tilde{v}_1 \oplus \tilde{v}_2$         ::binary operation

```

$$\frac{\tilde{v}_1 = \Theta_\kappa(\tilde{v}_2, \tilde{v}_2)}{\tilde{v}_1 = \tilde{v}_2} \quad (1) \quad \frac{\tilde{v}_1 = \Theta_\kappa(\tilde{v}_2, \tilde{v}_3)}{\tilde{v}_1 \oplus \tilde{v}_4 = \Theta_\kappa(\tilde{v}_2 \oplus \tilde{v}_4, \tilde{v}_3 \oplus \tilde{v}_4)} \quad (2)$$

$$\frac{\tilde{v}_1 = \Theta_\kappa(\tilde{v}_2, \tilde{v}_3)}{B[\tilde{v}_1] = \Theta_\kappa(B[\tilde{v}_2], B[\tilde{v}_3])} \quad (3)$$

$$\frac{\tilde{v}_1 = \Theta_\kappa(\Theta_\kappa(\tilde{v}_2, \tilde{v}_3), \tilde{v}_4)}{\tilde{v}_1 = \Theta_\kappa(\tilde{v}_2, \tilde{v}_4)} \quad (4) \quad \frac{\tilde{v}_1 = \Theta_\kappa(\tilde{v}_2, \Theta_\kappa(\tilde{v}_3, \tilde{v}_4))}{\tilde{v}_1 = \Theta_\kappa(\tilde{v}_2, \tilde{v}_4)} \quad (5)$$

Figure 6: Abstract values.

length. The $(\tilde{v} + 1)$ th byte of the input packet is $B[\tilde{v}]$. We introduce a new selection operator Θ_κ such that $v = \Theta_\kappa(v_1, v_2)$ means when the if-statement at κ takes the true branch, we have $v = v_1$, $v = v_2$ otherwise. One may find that the operator Θ_κ is similar to the operator ϕ in the classic SSA code form [36] because both of them merge values from multiple branches. We note that Θ_κ differs from ϕ in two aspects. First, in the SSA form, $v = \phi(v_1, v_2)$ is always placed at the end of a branching statement, whereas in our analysis $v = \Theta_{\kappa_i}(v_1, v_2)$ represents an abstract value of the variable v and is propagated to many other places where the variable v is referenced. Second, since $v = \Theta_\kappa(v_1, v_2)$ may be used at any place in the code, we use the subscript κ to record the branching statement where it originates. This is a critical design for the next step, i.e., the localized graph unfolding, as illustrated later. An abstract value can also be a first-order logic formula over other abstract values. To ease the explanation, we only support binary formulas.

Figure 6 lists the rules that normalize expressions over abstract values. Rule (1) states that we do not need a Θ_κ operator if we merge two equivalent values. Rules (2-3) state that any operation with a Θ_κ -merged value is equivalent to operating on each value merged by the Θ_κ operator. Rules (4-5) simplify nested Θ_κ operators.

Abstract Semantics. The abstract semantics describe how we analyze a given protocol parser. They are described as transfer functions of program statements. Each transfer function updates the program's abstract state, which is a pair (\mathbb{E}, \mathbb{G}) . Given the set V of program variables and the set \tilde{V} of abstract values, $\mathbb{E} : V \mapsto \tilde{V}$ maps a variable to its abstract value. We use $\mathbb{E}[v \mapsto \tilde{v}]$ to denote updating the abstract value of the variable v to \tilde{v} . \mathbb{G} is the output AFG. Since AFG is an equivalent form of path constraint, we directly create AFG without computing the path constraint first.

Figure 7 lists the transfer functions as inference rules. In each rule, the part above the horizontal line includes a set of assumptions and, under these assumptions, the bottom part describes the abstract states before and after a statement S , in the form of $\mathbb{E}, \mathbb{G} \vdash S : \mathbb{E}', \mathbb{G}'$. Initially, we assign the special abstract value length to the variable

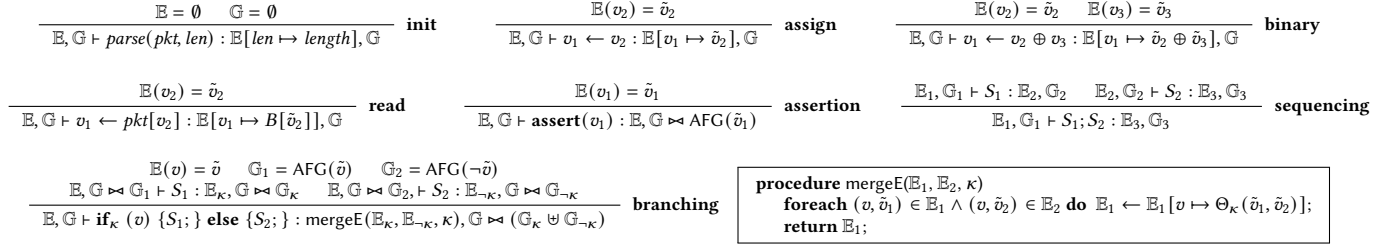


Figure 7: Inference rules and auxiliary procedure.

len, which represents the length of input network packet. The rules for assignment, binary operation, read operation, and assertion are straightforward. For instance, in the rule for assertions, the abstract value \tilde{v}_1 represents a constraint that must be satisfied. Therefore, we append the graph $\text{AFG}(\tilde{v}_1)$ to the graph \mathbb{G} . This is equivalent to appending the constraint \tilde{v}_1 to the current path constraint.

The sequencing rule states that, for two consecutive statements, we analyze them in order, using the postcondition of the first statement as the precondition for the second. In the branching rule, \mathbb{G} denotes the path constraint before the branching statement. \mathbb{G}_1 and \mathbb{G}_2 represent the branching condition and its negation. Thus, $\mathbb{G} \bowtie \mathbb{G}_1$ and $\mathbb{G} \bowtie \mathbb{G}_2$ represent the initial path constraints before the two branches. After analyzing the two branches, the resulting AFGs are assumed to be $\mathbb{G} \bowtie \mathbb{G}_\kappa$ and $\mathbb{G} \bowtie \mathbb{G}_{-\kappa}$. The branching rule states that, under these assumptions and after an if_κ -statement, we merge the abstract states from both branches. The procedure `mergeE` merges abstract values of the same variable via the Θ_κ operator. Graph merging is straightforward based on the definition of AFG, which is equivalent to merging path constraints of the two branches with the common prefix pulled out. Our merging is different from the value merging in traditional analyses due to the use of the selection operator. On one hand, merging allows achieving scalability as the number of values is no longer exponential of the number of statements. On the other hand, the selectors in abstract values can be unfolded to support path-sensitive analysis if needed.

Packet Fields. The abstract interpretation builds the AFG to represent the path constraints. As discussed in §3, from these constraints, it is direct to infer the endianness, field boundaries, and direction fields. For instance, if multiple consecutive bytes, e.g., $B[0]$ and $B[1]$ in Figure 4, belong to a single field, the field value, e.g., $B[0]B[1]$, will be computed and occur in the path constraint.

High-Level Field Semantics. We also extend our analysis to infer high-level field semantics, i.e., field names, using rich source code information. Such high-level semantics can help better understand a format, e.g., identifying checksum fields. As illustrated in Figure 4, we can name a field (via some variable name) by adding extra path constraints. Formally, given the AFG \mathbb{G} and a formula over a field $B[i..j]$, denoted as $f(B[i..j])$, we name the field by $\mathbb{G} \bowtie \text{AFG}(\text{name}(B[i..j]) = \text{'var'})$ if there is a statement assigning $f(B[i..j])$ to the variable *var*. In addition to variable names, we also leverage system APIs used in the code. For instance, if a field $B[i..j]$ is used in the system API, `diffTime()`, it is likely to be a timestamp field. In our experience, this method helps us identify many special

fields via names such as ‘length’, ‘version’, ‘checksum’, and ‘timestamp’, to name a few. If there are multiple options for naming a field, we prefer the names inferred by system APIs because software developers may not be careful to name program variables. If there are still multiple options, we simply keep the first.

Example. Given the code in Figure 4(a), the abstract interpretation yields the AFG in (b) from top to bottom. After Line 5, we merge the two paths forked from Line 3 and get the path constraint, $\text{name}(B[4]) = \text{'ctrl'} \wedge (B[5] = 0 \vee B[5] \neq 0)$. By the branching rule, we do not compute the path constraint but directly create the equivalent AFG, i.e., the first two rows in Figure 4(b). We name the byte $B[4]$ ‘ctrl’ because the arithmetic results of $B[4]$ are assigned to the variable *ctrl* in both branches. The constraint $B[5] = 0 \vee B[5] \neq 0$ merges the branching constraints. Meanwhile, the abstract store is updated such that $\text{ctrl} = \Theta_{\kappa_3}(B[4] + 1, B[4] - 1)$. At Line 7, since the false branch aborts, we only consider the true branch, for which we add $B[6] > 0 \wedge \Theta_{\kappa_3}(B[4] + 1, B[4] - 1) = 0$ to the path constraint. This is equivalent to adding the third and fourth rows in Figure 4(b). At Lines 10–11, we add the constraint $B[0]B[1] = 10 \wedge \text{name}(B[0..1]) = \text{'code'}$. This is equivalent to adding the fifth and sixth rows in Figure 4(b). We regard $B[0]$ and $B[1]$ as a single field as they are used in a single value $B[0]B[1]$.

Similarly, after Line 15, we merge the paths forked from Line 11 as the constraint $(B[3] = 0 \vee B[3] \neq 0) \wedge \text{name}(B[2]) = \text{'state'}$ and append it to the path constraint. This is equivalent to adding the last row in Figure 4(b). After Line 15, we update the value $\text{state} = \Theta_{\kappa_{13}}(B[2] + 1, B[2] - 1)$. In this example, the variable *state* is simply printed at Line 16 and never used in any if-statements or assertions. Hence, the merged value of *state* is abstracted away from the final constraint. Observe that the size of AFG is linear size with the number of statements. This is critical to scalability. \square

LEMMA 5.2. *Given a program in the language defined in Figure 5, the AFG produced by the abstract interpretation is sound and complete.*

5.3 Localized Graph Unfolding

We observe that high-precision, i.e., path-sensitive, analysis is only needed in localized regions during lifting (Challenge 1 in §4). Intuitively, a code region that requires path sensitivity is identified as follows. If a Θ_κ -merged value is used in path constraint, each branch of the if_κ -statement needs to be analyzed separately. However, given an if_κ -statement, if the path constraint we compute does not include any Θ_κ -merged values, we do not need to separately analyze its branches, meaning that path sensitivity is not necessary.

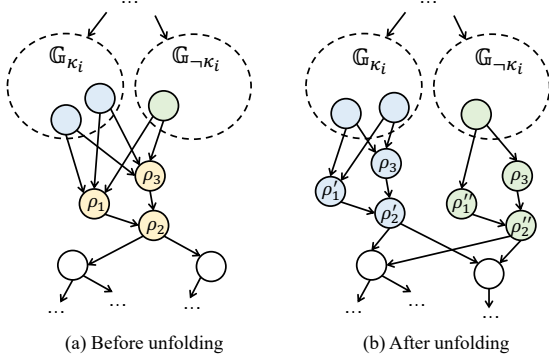


Figure 8: Example of unfolding an AFG.

In detail, given an AFG created by the abstract interpretation, we eliminate all Θ_K -merged values by a localized graph unfolding algorithm shown in Algorithm 1. Assume that the AFG to unfold contains a list of Θ_K operators, e.g., Θ_{K_0} , Θ_{K_1} , and Θ_{K_2} . The algorithm eliminates Θ_{K_i} one by one. For each Θ_{K_i} , it works in two steps – slicing (Lines 3-7) and unfolding (Lines 8-11). To ease the explanation, we use Figure 8 for illustration. In Figure 8(a), without loss of generality, assume that we are unfolding Θ_{K_i} in the AFG and that only the constraints ρ_1 and ρ_2 contain Θ_{K_i} -merged values. The exiting vertices of G_{K_i} and G_{-K_i} are shown in the figure.

Slicing. This step delimits the next unfolding step to a local region in AFG. First, we find all exiting vertices of G_{K_i} and G_{-K_i} . We then perform a forward graph traversal (e.g., depth-first search) from the exiting vertices. Denote the subgraph visited during the traversal as G_{forward} . Second, we identify all vertices containing Θ_{K_i} -merged values, e.g., ρ_1 and ρ_2 in Figure 8(a). A backward graph traversal from them yields a subgraph denoted as G_{backward} . The overlapping part of G_{forward} and G_{backward} , e.g., the yellow part in Figure 8(a), is the graph slice where we will perform unfolding, denoted as G_{slice} .

Unfolding. As illustrated in Figure 8(b), we copy the subgraph to unfold, obtaining G_{slice} and G'_{slice} . The copy G_{slice} is connected to G_{K_i} , and by the definition of the merging operator, all the Θ_{K_i} -merged values are replaced by its first operand. Similarly, the other copy G'_{slice} is connected to G_{-K_i} , and all the Θ_{K_i} -merged values are replaced by its second operand. Since the subgraphs to unfold are limited in small local regions in practice, we significantly mitigate the path-explosion problem, which is sufficient to make our approach scalable. Note that we do not claim to have a theoretical bound on the size of subgraphs that need to be unfolded, as path explosion is still an open problem and cannot be completely addressed in theory, similar to all previous path-sensitive analyzers.

LEMMA 5.3. *The unfolded AFG does not contain Θ_K -merged values and represents an equivalent constraint as the original AFG.*

Example (continued). In Figure 4(b), the value merged by Θ_{K_3} indicates that the branches forked at Line 3 need a path-sensitive analysis and delimits the analysis to the local region colored gray. To distinguish the two branches, the gray region in (b) is unfolded to two disjoint paths in (c), which eliminates the Θ -merged values and make the two semantic relations among $B[4]$, $B[5]$, and $B[6]$ explicit: $B[5] = 0 \wedge B[6] > 0 \Leftrightarrow B[4] + 1 = 0$; and $B[5] \neq 0 \wedge B[6] >$

Algorithm 1: Unfolding.

```

1 Procedure unfold( $\mathbb{G}$ )
2   foreach operator  $\Theta_{K_i}$  in  $\mathbb{G}$  do
3      $G_{\text{forward}} \leftarrow$  subgraph reachable from but excluding  $G_{K_i}$  and  $G_{-K_i}$ ;
4      $V \leftarrow$  all vertices including  $\Theta_{K_i}$  expressions;
5      $G_{\text{backward}} \leftarrow$  subgraph that can reach any vertex in  $V$ , including  $V$ ;
6      $G_{\text{slice}} \leftarrow$  overlapping subgraph of  $G_{\text{forward}}$  and  $G_{\text{backward}}$ ;
7      $G'_{\text{slice}} \leftarrow$  a copy of  $G_{\text{slice}}$ , including all its incoming/outgoing edges;
8     disconnect  $G_{\text{slice}}$  from  $G_{-K_i}$ ;
9     replace all  $\Theta_{K_i}$  expressions in  $G_{\text{slice}}$  with their first operands;
10    disconnect  $G'_{\text{slice}}$  from  $G_{K_i}$ ;
11    replace all  $\Theta_{K_i}$  expressions in  $G'_{\text{slice}}$  with their second operands;

```

$0 \Leftrightarrow B[4] - 1 = 0$. In contrast, the $\Theta_{K_{13}}$ value in variable *state* is never used in any conditional, suggesting that we do not need to path-sensitively analyze the branches led by Line 13. \square

5.4 Localized Graph Reordering

As illustrated in Figure 4, bytes in a packet may not appear in the order in a program path, e.g., $B[5]$ may precede $B[2]$. To produce legitimate BNF productions, we need to reorder them to produce the ordered AFG (Challenge 2 in §4). To this end, we first define the concepts of vertical decomposition (VD) and horizontal decomposition (HD) of an AFG.

Definition 2 (VD). Given an unfolded AFG $\mathbb{G} = G_1 \bowtie G_2 \bowtie \dots \bowtie G_n$, namely, the exit vertices in G_i are fully connected to the entry vertices in G_{i+1} , its vertical decomposition is the sequence of subgraphs, denoted as $\text{VD}(\mathbb{G}) = G_1 G_2 \dots G_n$.

Definition 3 (HD). Given an unfolded AFG \mathbb{G} , its horizontal decomposition is a set of subgraphs, each of which is rooted at a single entry vertex in the AFG and includes the subgraph reachable from the entry vertex, denoted as $\text{HD}(\mathbb{G}) = G_1 | G_2 | \dots | G_n$.

Figure 10(a) shows an example of vertical decomposition, where the graph is decomposed into two parts, one containing the vertices ρ_1 and ρ_2 , and the other containing the vertices ρ_3 , ρ_4 , and ρ_5 . The graph in Figure 10(b) cannot be vertically decomposed because the upper two vertices are not fully connected to the other three. Instead, it can be horizontally decomposed into two parts, one containing the vertices ρ_1 , ρ_3 , ρ_4 , and ρ_5 , and the other containing the vertices ρ_2 , ρ'_3 , and ρ'_5 . Here, ρ'_3 and ρ'_5 are copies of ρ_3 and ρ_5 , respectively. As illustrated in the example and stated in Lemma 5.4, the AFGs before and after decomposition contain the same number of paths and the constraint represented by each path is not changed.

LEMMA 5.4. *AFGs before and after decomposition are equivalent in representing path constraint.*

The decomposition has three properties. First, the horizontal decomposition is more expensive than the vertical one as it may copy vertices. Hence, Algorithm 2 always tries the vertical decomposition first. Second, as stated in Lemma 5.5, the decomposition can be recursively performed on a graph and its subgraphs. For instance, after the horizontal decomposition in Figure 10(b), we can further apply vertical decomposition to each subgraph. This property allows us to describe our reordering approach as a recursive process in Algorithm 2. Third, the vertical decomposition follows the commutative law stated in Lemma 5.6. For instance,

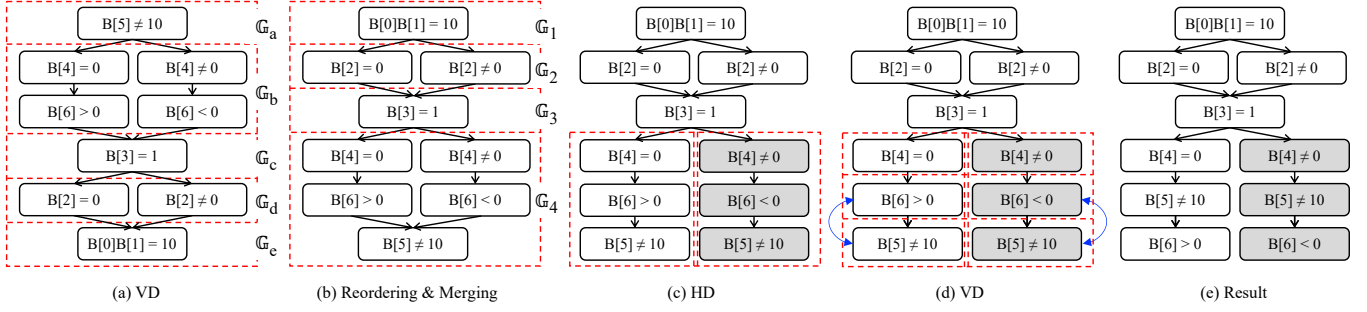


Figure 9: Example of Algorithm 2.

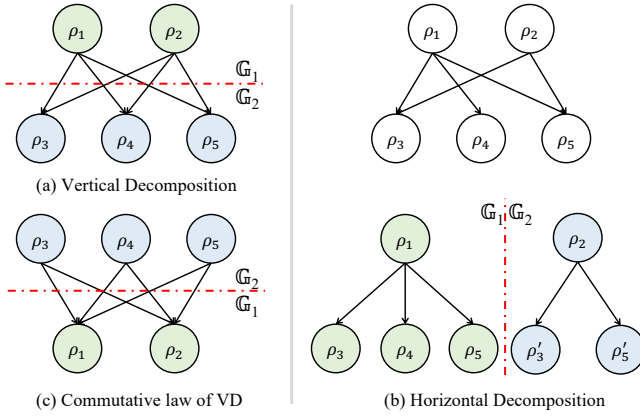


Figure 10: Decomposition for graph reordering.

after switching \mathbb{G}_1 and \mathbb{G}_2 in Figure 10(a), we get the graph in Figure 10(c), which is equivalent to the original graph because they represent equivalent path constraints: $(\rho_1 \vee \rho_2) \wedge (\rho_3 \vee \rho_4 \vee \rho_5)$ and $(\rho_3 \vee \rho_4 \vee \rho_5) \wedge (\rho_1 \vee \rho_2)$. Such a commutative property allows us to reorder vertices in Algorithm 2.

LEMMA 5.5. *If an AFG with multiple vertices cannot be vertically decomposed, each subgraph after horizontal decomposition contains a single vertex or can be vertically decomposed.*

LEMMA 5.6. *Switching the position of subgraphs in VD yields an AFG that represents an equivalent constraint as the original AFG.*

Algorithm 2 first tries to vertically decompose the input AFG (Line 2). If failed, Lemma 5.5 allows us to horizontally decompose it into subgraphs and recursively order each subgraph (Lines 14-15). If VD succeeds in splitting AFG into a list of subgraphs, these subgraphs are reordered by byte indices (Lines 3-5). Figure 9(a) and Figure 9(b) illustrate this step. In Figure 9(a), the AFG is vertically decomposed into five subgraphs, $\mathbb{G}_a, \mathbb{G}_b, \mathbb{G}_c, \mathbb{G}_d$, and \mathbb{G}_e , which are respectively put in five dashed boxes. The minimum byte indices of the subgraphs are 5, 4, 3, 2, and 0. Figure 9(b) shows the AFG after reordering the subgraphs based on the minimum byte indices. After reordering, since \mathbb{G}_a and \mathbb{G}_b contain overlapping byte indices¹, they are merged into a single subgraph, i.e., \mathbb{G}_4 in Figure 9(b). In this example, the subgraphs after reordering and merging are put

¹The range of byte indices in \mathbb{G}_a is [5, 5], and the range in \mathbb{G}_b is [4, 6]. The former is a subset of the latter. Thus, they overlap each other.

Algorithm 2: Reordering.

```

1 Procedure reorder( $\mathbb{G}$ )
2   if VD( $\mathbb{G}$ ) =  $\mathbb{G}_a \mathbb{G}_b \dots$  then
3     reorder  $\mathbb{G}_a, \mathbb{G}_b, \dots$ , based on the min byte index of each subgraph;
4     merge adjacent subgraphs if they contain overlapping byte indices;
5     let  $\mathcal{A} \leftarrow [\mathbb{G}_1, \mathbb{G}_2, \dots]$  be subgraphs after reordering and merging;
6     foreach  $\mathbb{G}_i \in \mathcal{A}$  do
7       if  $\mathbb{G}_i = \mathbb{G}_{i,1} \bowtie \mathbb{G}_{i,2} \bowtie \dots$  is a merged graph then
8         assume  $\mathbb{G}_{i,j}$  has multiple entry vertices;
9         switch the position of  $\mathbb{G}_{i,j}$  and  $\mathbb{G}_{i,1}$  in  $\mathbb{G}_i$ ;
10        assume HD( $\mathbb{G}_i$ ) =  $\mathbb{G}_{a'} | \mathbb{G}_{b'} | \dots$ ;
11        reorder( $\mathbb{G}_{a'}$ ); reorder( $\mathbb{G}_{b'}$ ); ...;
12      else
13        reorder( $\mathbb{G}_i$ );
14   else if HD( $\mathbb{G}$ ) =  $\mathbb{G}_a | \mathbb{G}_b | \dots$  then
15     reorder( $\mathbb{G}_a$ ); reorder( $\mathbb{G}_b$ ); ...;
16   else // a single-vertex graph, do nothing

```

Algorithm 3: Packet Format in BNF.

```

1 Procedure bnf( $\mathbb{G}$ )
2    $L \leftarrow$  new non-terminal symbol;
3   if VD( $\mathbb{G}$ ) =  $\mathbb{G}_a \mathbb{G}_b \dots$  then
4      $L \rightarrow \text{bnf}(\mathbb{G}_a) \text{bnf}(\mathbb{G}_b) \dots$ ;
5   else if HD( $\mathbb{G}$ ) =  $\mathbb{G}_a | \mathbb{G}_b | \dots$  then
6      $L \rightarrow \text{bnf}(\mathbb{G}_a) | \text{bnf}(\mathbb{G}_b) | \dots$ ;
7   else
8     // a single vertex containing  $B[i], B[i+1], \dots, B[i+k]$ 
9      $L \rightarrow B[i]B[i+1] \dots B[i+k]$  with assertions in the vertex;
10  return  $L$ ;

```

in the array $\mathcal{A} = [\mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_3, \mathbb{G}_4]$. These subgraphs are ordered and contain mutually exclusive byte indices.

We then recursively reorder subgraphs in \mathcal{A} (Lines 6-13). Especially, for a merged subgraph, e.g., $\mathbb{G}_4 = \mathbb{G}_a \bowtie \mathbb{G}_b$ in the example, since we have tried vertical decomposition, which does not work as neither $\mathbb{G}_a \mathbb{G}_b$ nor $\mathbb{G}_b \mathbb{G}_a$ respects the stream order, we turn to horizontal decomposition (Lines 8-11). Lines 8-9 ensure the feasibility of horizontal decomposition and Line 10 performs the decomposition. Figure 9(c) illustrates this step, where the subgraph \mathbb{G}_4 is horizontally decomposed into the white and the gray parts. Each part then is recursively reordered (Line 11). Figure 9(d) shows that the white and the gray parts are recursively split by vertical decomposition and reordered as indicated by the arrows, yielding the ordered AFG in Figure 9(e). Lemma 5.7 states the correctness of Algorithm 2.

LEMMA 5.7. *Algorithm 2 yields an ordered AFG, which represents an equivalent constraint as the input AFG.*

From Ordered AFG to BNF-like Format. It is straightforward to translate an ordered AFG to packet formats in BNF. Due to its simplicity, the detailed discussion is elided and the formal algorithm is put in Algorithm 3. As an example, Figure 4(e) shows the inferred packet format where S is the start symbol that represents the whole graph and each non-terminal L_i represents a subgraph — L_1 represents the path prefix containing $B[0]$, $B[1]$, and $B[2]$; L_2 and L_3 represent two possible constraints of $B[3]$; and L_4 and L_5 stand for the two path suffixes containing $B[4]$, $B[5]$, and $B[6]$.

5.5 Soundness and Completeness in Practice

Lemmas 5.1-5.7 together guarantee the theoretical soundness and completeness of our approach for a program written in our abstract language. In practice, we need to handle common program structures not included in the abstract language, such as function calls, pointers, and loops. This section discusses how we handle them in our implementation and their effects on soundness or completeness.

Pointers. In the previous discussion, we focus on building an AFG for format inference. Pointer operations are not directly related to AFG. In the implementation, we follow existing works [68] to resolve pointer relations, which helps us identify what values may be loaded from a memory location. For instance, when visiting an assertion in the program such as `assert(*(p + 1) > 1)` where p is a pointer, if the pointer analysis tells us $p+1$ points to a memory location storing the value $B[5]$ on the condition ρ , we then compute and include the constraint $\rho \Rightarrow B[5] > 1$ (which equals $\neg \rho \vee B[5] > 1$) in AFG. Pointer operations such as $p+1$ are not a part of path constraints and, thus, are not included in AFG. That is, according to the assertion rule in Figure 7 and assuming the AFG before the assertion is \mathbb{G} , the AFG after the assertion is $\mathbb{G} \bowtie \text{AFG}(\neg \rho \vee B[5] > 1)$. Since the pointer analysis we use is sound and path-sensitive, it allows Netlifter to be sound and highly precise.

Function Calls. Although we do not include function calls in our abstract language for simplicity, our system is inter-procedural as a call statement is equivalent to a list of assignments from the actual parameters to the formals, and a return statement is an assignment from the return value to its receiver. Thus, in our analysis, function calls and returns are treated as assignments. This treatment does not degrade soundness and completeness. Especially, for recursive function calls, we convert them to loops, which are discussed below.

Loops and Repetitive Fields. Loops in a protocol parser are often used to parse repetitive fields [35]. We follow existing techniques to analyze loops [64, 76], which are good at inferring repetitive fields and how many times a field repeats. For example, the code below parses a packet where $B[0]$ represents the packet length and contains a positional constraint that all bytes after $B[0]$ are less than five. For this example, we produce the production $S \rightarrow B[0]B[1]$ with two semantic constraints: $B[1] < 5$ and $\text{repeat}(B[1]) = B[0]$.

```
j = 0; while(j < packet[0]) { assert(packet[++j] < 5); }
```

Basically, the loop analysis works in two steps. First, they analyze several iterations of a loop. For instance, it analyzes three iterations and gets the results, $j = 1 \wedge B[1] < 5$, $j = 2 \wedge B[2] < 5$, and $j = 3 \wedge B[3] < 5$, for each iteration. Second, it inductively infers

```
1. short compute_crc(char *buf, int len) {
2.   short crc = 0x1d0f;
3.   for (; len > 0; len--) {
4.     crc = (crc >> 8) | (crc << 8);
5.     crc ^= *buf++; // read B[k] in (k+1)th iteration
6.     crc ^= (crc & 0xff) >> 4;
7.     crc ^= crc << 12;
8.     crc ^= (crc & 0xff) << 5;
9.   }
10.  return crc;
11. }
12.
13. void parse_packet(char *buf, int len) { // buf: packet; len: packet length
14.   ...
15.   short crc = buf[len - 1] << 8 | buf[len - 2]; // crc: B[length-1]B[length-2]
16.   assert(crc == compute_crc(buf, len - 2)); // crc == f(B[0], B[1], B[2], ...)
17.   ...
18. }
```

Figure 11: An irregular loop that computes checksum.

the conditions. For instance, the above results can be inductively summarized to $j = k \wedge B[k] < 5$, where $1 \leq k \leq B[0]$ is an induction variable representing the iteration counter. Equivalently, we write the constraint as $B[1] < 5 \wedge \text{repeat}(B[1]) = B[0]$ and, by the definition of AFG, represent it as an edge from a vertex containing $B[1] < 5$ to one containing $\text{repeat}(B[1]) = B[0]$. If the inductive inference succeeds, its result is sound and complete.

Irregular Loops. Nevertheless, not all loops can be inductively summarized as above. This is an inherent limitation of static analysis. For irregular loops, we follow the practice of bounded model checking [23] to unroll loops a fixed number of times. While this may introduce unsoundness and incompleteness, we observe few irregular loops (involved in packet parsing) in practice and their influence is limited. In other words, it may produce some unsound or incomplete constraints for some fields in a packet but such unsoundness and incompleteness are rarely propagated to other fields.

The most common irregular loops related to protocol formats are those computing the checksum. Figure 11 shows an example where the function `compute_crc` computes the CRC value as the checksum and the function `parse_packet` compares the computed CRC value to the received one in the packet. The CRC value is computed by a loop, which cannot be inductively summarized. This is because, unlike the variable j in the previous example, we cannot write the value of `crc` as a formula parameterized by the loop counter k .

Thus, we unroll the loop a fixed number of times t , yielding an unsound and incomplete formula of `crc` that relies on $B[0]$, $B[1]$, \dots , $B[t-1]$, denoted as $f(B[0], B[1], \dots, B[t-1])$. At Line 16, this formula is compared to the CRC field in the packet. Generally, by the assertion rule in Figure 7, an AFG node with the constraint $B[\text{length}-1]B[\text{length}-2] = f(B[0], B[1], \dots, B[t-1])$ should be appended to the AFG. However, in the implementation, we skip such unsound and incomplete constraints, so that they are not included in the protocol format output by Netlifter. Excluding these constraints keeps our inferred format sound but may lose some precision (i.e., completeness) as we miss some constraints.

While we cannot compute precise non-recursive FOL constraints for checksum fields (this is a limitation shared by existing works on protocol format reverse engineering), as discussed in §5.2, we can still infer the boundary of checksum fields. Such field boundaries, together with all sound constraints we compute, can already support many security applications (see §6.3).

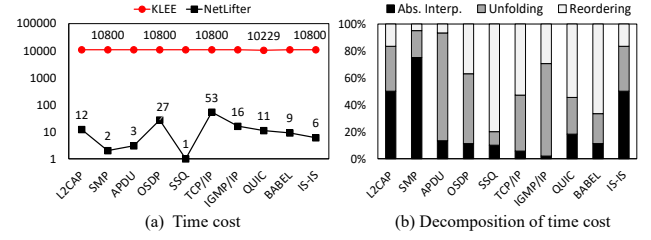
Table 1: Protocols and Their Codebases for Evaluation

Name	Codebase	Size (kloc)	Time (sec.)	Description
L2CAP	linux/bluetooth [7]	38	12	logical link ctrl and adaptation proto.
SMP	linux/bluetooth [7]	12	2	low energy security manager proto.
APDU	opensc [14]	3	3	application proto. data unit
OSDP	libosdp [5]	14	27	open supervised device proto.
SSQ	libssq [8]	8	1	source server query proto.
TCP/IP	lwip [6]	41	53	transport control & internet proto.
IGMP/IP	lwip [6]	17	16	internet group mgmt. & internet proto.
QUIC	ngtcp2 [4]	59	11	general-purpose transport layer proto.
BABEL	frouting [3]	7	9	a distance-vector routing proto.
IS-IS	frouting [3]	22	6	intermediate system (IS) to IS proto.
A2MP	linux/bluetooth [7]	16	2	amp manager proto.
BNP	linux/bluetooth [7]	15	3	BT network encapsulation proto.
CMT	linux/bluetooth [7]	20	1	c-api message transport proto.
HIDP	linux/bluetooth [7]	17	4	human interface device proto.
UDP	lwip [6]	37	33	user datagram proto.
ICMP	lwip [6]	22	12	internet control message proto.
DHCP	lwip [6]	25	43	dynamic host configuration proto.
ICMPv6	lwip [6]	30	54	internet control message proto. v6
DHCPv6	lwip [6]	35	51	dynamic host configuration proto. v6
BGP	frouting [3]	13	2	border gateway proto.
LDP	frouting [3]	20	5	label distribution proto.
BFD	frouting [3]	10	17	bidirectional forwarding detection
RRRP	frouting [3]	8	12	virtual router redundancy proto.
EIGRP	frouting [3]	14	21	interior gateway routing proto.
NHRP	frouting [3]	11	11	next hop resolution proto.
OSPF2	frouting [3]	9	14	open shortest path first v2
OSPF3	frouting [3]	7	16	open shortest path first v3
RIP1	frouting [3]	11	13	routing information proto. v1
RIP2	frouting [3]	11	15	routing information proto. v2
RIPng	frouting [3]	7	41	routing information proto. for ip6

6 EVALUATION

We implement our method as a tool, namely Netlifter, to lift packet formats from source code in C. It is implemented on top of the LLVM (12.0.0) compiler infrastructure [50] and the Z3 (4.8.12) SMT solver [37]. The source code of a protocol is compiled into the LLVM bitcode, where we perform our static analysis. In the analysis, Z3 is used to represent abstract values as symbolic expressions and solve path constraints. All experiments are run on a Macbook Pro (16-inch, 2019) equipped with an 8-core 16-thread Intel Core i9 CPU with 2.30GHz speed and 32GB of memory.

As shown in Table 1, we have run Netlifter over a number of protocols from different codebases (e.g., Linux and LWIP) and domains (e.g., IoT and routing). They include widely-used ones such as TCP/IP and niche ones like APDU that is used in smart cards. As shown in the table, the code size of a protocol parser ranges from 3KLoC to 59KLoC, and it takes Netlifter <1min to infer the format of each protocol. Determining the precision and recall of the inferred formats requires manually comparing them with their official documents. We cannot afford to manually inspect all protocols because we have to learn a lot of domain-specific knowledge to understand a protocol, which is time-consuming and not very related to our core contribution to the static analysis. In the remaining experiments, we focus on the first ten, which are from different codebases. We believe that other protocols in the same codebases are implemented in similar manners and, thus, do not introduce extra challenges. We use these protocols and codebases because of two reasons. First, their repositories in GitHub are relatively active, which makes it easy to get feedback from developers when we report bugs. Second, they have their own fuzzing drivers, meaning that they have been extensively fuzzed by the developers themselves. Thus, their code is expected to be of high quality and an approach that can find vulnerabilities in their codebase is highly effective.

**Figure 12: Time cost (seconds) and its decomposition.**

6.1 Effectiveness of the Three-Step Design

For technical contributions, we explained in §4 that our static analysis avoids individually exploring program paths to address two challenges. To show the importance of our solution, we implement a baseline that employs a well-known symbolic executor, KLEE [28], to infer packet formats. Similar to our solution, it infers packets by computing path constraints. Different from our solution, it has to analyze individual program paths. We then compare their time cost of format inference. The results are shown in Figure 12(a) in log scale. The line chart shows that the KLEE-based approach runs out of time (≥ 3 hours) for almost all protocols. We use a three-hour time budget here as it is sufficient to show the advantage of our approach over symbolic execution. As plotted in Figure 12(a), Netlifter can finish in one minute. Figure 12(b) shows the decomposition of Netlifter's time cost, indicating that the three steps of Netlifter respectively take 14%, 44%, and 42% of the total time.

6.2 Precision and Recall of Packet Formats

As discussed in §1, existing techniques focus on network trace analysis (category one) or dynamic program analysis (category two). We refer to both of them as dynamic analyses as they rely on dynamically captured network packets as their inputs. We cannot find any static program analysis that infers formats from a protocol parser. Thus, while the dynamic analyses have a different assumption from our static analysis, not for a comparative purpose but to show the value of our approach, we evaluate Netlifter with two network trace analyses, i.e., NemeSys [46, 47] and NetPlier [77], and two dynamic program analyses, i.e., AutoFormat [54] and Tupni [35]. NemeSys and NetPlier are open-source software and we directly use their implementation. AutoFormat and Tupni are not publicly available. We implement them on top of LLVM based on their papers. We cannot find other open-source dynamic program analyses for evaluation. We evaluate them in terms of precision and recall. Given a set of packets, the precision is the ratio of correctly inferred fields in the packets to all inferred fields. The recall is the ratio of correctly inferred fields to all fields in the ground truth. To compute the precision and recall, we manually build the formats based on the protocols' official documents or source code. We then write scripts to compare the inferred and the manually-built formats.

Dynamic Analysis. To use dynamic analyses, we follow their original works to collect 1000 network packets for each protocol from publicly available datasets [9, 12, 63]. Table 2 shows the precision and recall of the inferred field boundaries. Network trace analyses often exhibit low precision (<50%) and recall (<50%), because they use statistical approaches to align message fields while

Table 2: Precision(%) / Recall(%) of Field Boundaries.

Protocol	Netlifter	NemeSys	NetPlier	AutoFormat	Tupni	Combined
L2CAP	96/98	14/9	14/15	72/32	88/41	66/49
SMP	100/100	27/37	20/52	100/52	96/78	45/82
APDU	100/100	52/21	43/45	44/25	100/61	58/71
OSDP	100/100	17/11	10/16	74/31	89/47	43/52
SSQ	100/100	25/1	26/11	88/19	95/54	41/67
TCP/IP	98/95	5/7	4/12	24/19	88/21	39/25
IGMP/IP	99/98	13/12	13/22	35/22	92/25	54/36
QUIC	97/99	19/14	18/25	70/29	86/43	69/53
BABEL	99/99	28/14	37/8	43/16	80/24	40/28
IS-IS	98/99	23/5	19/14	100/34	87/21	62/42

Table 3: Precision(%) / Recall(%) of Field Names.

L2CAP	SMP	APDU	OSDP	SSQ
100/97	100/96	100/100	100/100	100/100
TCP/IP	IGMP/IP	QUIC	BABEL	IS-IS
100/95	100/98	96/96	100/95	100/94

statistical approaches are known to have inherent uncertainty and their effectiveness heavily hinges on the quality of input packets.

The two dynamic program analyses, especially Tupni, significantly improve the precision due to the analysis of control and data flows in the code. AutoFormat has a relatively low precision because it tracks coarse-grained control and data flows. For instance, AutoFormat regards consecutive bytes of a packet processed in the same calling context as a single field. However, it is common for a parser to process multiple fields in the same calling context. Tupni tracks more fine-grained control and data flows, such as predicates in the code, and, thus, exhibits higher precision. As acknowledged by Tupni itself, it may also produce false fields in many cases. For instance, when the value of a multi-byte field is computed by multiple instructions over every single byte in the field, it will incorrectly split the field into multiple fields. Despite the high precision achieved by Tupni, the key problem of these dynamic analyses is their coverage (i.e., recall), which is often lower than 50% and may compromise downstream security analyses as discussed in the next subsection.

Note that simply combining the results of multiple tools does not help improve the quality of the inferred formats. This is because, when combining the formats inferred by multiple tools, with the increase of correctly inferred fields, the number of incorrect fields also increases. For instance, after combining the results of the four dynamic tools, the precision for OSDP is 0.43, which is even worse than the result when using Tupni independently. The combined results are shown in the last column of Table 2.

Static Analysis. Table 2 shows that, in terms of field boundaries, our inferred formats cover >96% formats and produce <4% false ones. For many of them, we can produce absolutely correct formats. We also miss some fields and report some false ones due to the inherent limitations of static analysis (see §9). These limitations, e.g., the incapability of handling inline-assembly in the source code, lets us lose information, thereby leading to the false formats. Table 3 also shows the quality of the inferred field names. A name is considered to be correct if it is the same as the official documents or a reasonable abbreviation, e.g., ‘len’ vs. ‘length’. Overall, we can infer >94% field names with a precision >96%. The names provide high-level semantics and help us identify special fields to facilitate security applications as discussed next.

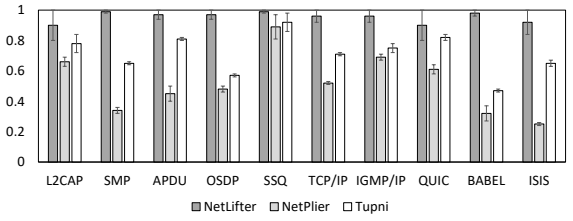


Figure 13: The y-axis is the number of covered branches normalized to one. It shows the branch coverage averaged over twenty runs with a 95% confidence interval.

6.3 Security Applications

Protocol Fuzzing. To show the value of our approach, we respectively input the formats inferred by Netlifter, NetPlier, and Tupni to a typical grammar-based (i.e., format-based) protocol fuzzer, namely BooFuzz [11, 13]. Particularly, since we can locate checksum fields by names such as ‘checksum’ and ‘crc’, in the fuzzing experiments, we can skip the checksum checks in the code. This is critical for fuzzing as random mutations in fuzzing can easily invalidate the checksum values [70]. The experiments are performed on a three-hour budget and repeated twenty times to avoid random factors. We use a three-hour budget because we observe that the baseline fuzzers rarely achieve new coverage after three hours.

The results are shown in Figure 13. Since Netlifter can provide formats with precise field boundaries and semantic constraints, Netlifter-enhanced BooFuzz achieves 1.2× to 3.6× coverage compared to others. We also detected 53 zero-day vulnerabilities while the others detect only 12. All detected vulnerabilities are exploitable as they can be triggered via crafted network packets. To date, 47 of them have been assigned CVEs. We can detect more bugs as our inferred formats are of both high precision and high coverage. In the extended version of this paper [66], we provide more details about the fuzzing experiments and the detected bugs.

Traffic Auditing and Intrusion Detection. The extended version of this paper [66] also provides a case study of applying our approach to network traffic auditing and intrusion detection. We conclude that the precise and high-coverage formats inferred by us are critical in these security applications.

7 RELATED WORK

We discuss some typical existing works in this section. For a broader overview, we refer readers to four surveys [40, 52, 61, 67].

Network Trace Analysis (NTA). NTA uses statistical methods to identify field boundaries based on runtime network packets. Discoverer [33] relies on a recursive clustering approach to recognize packets of the same type. Biprominer [71] uses the variable length pattern to locate protocol keywords and is enhanced by ProDecoder [72]. AutoReEngine [57] uses data mining to identify protocol keywords, based on which packets are classified into different types. ReverX [20] uses a speech recognition algorithm to identify delimiters in packets. NemeSys [46, 47] interprets binary packets as feature vectors and applies an alignment and clustering algorithm to determine the packet format. NetPlier [77] leverages a probabilistic analysis to determine the keyword field, clusters packets based on the keyword values, and applies multi-sequence alignment to

derive packet format. These techniques analyze neither the binary nor the source code and, thus, are different from ours.

Dynamic Program Analysis (DPA). DPA can be used over both source and binary code. They work by running protocol parsers against network packets and monitoring runtime control or data flows. Polyglot [27] uses dynamic taint analysis to infer fixed or variable length fields. AutoFormat [54] approximates the field hierarchical structure by monitoring call stacks. This approach then is extended to both bottom-up and top-down hierarchical structures [55]. Wondracek et al. [75] identify delimiters and length fields within a hierarchical structure. Tupni [35] tracks fine-grained taint flows to identify packet fields. It also applies loop analysis to infer repeated fields and records path constraints to infer length or checksum fields. ReFormat [74] recognizes encrypted fields based on the observation that encrypted fields are processed by a high percentage of arithmetic/bitwise instructions. Our approach can be easily extended with the same observation, i.e., by counting relevant instructions to recognize an encrypted field. In addition to inferring the formats of received packets, Dispatcher [26] and P2C [48] reverse engineer the formats of packets to be sent and, thus, are different from all aforementioned approaches and ours.

Static Program Analysis (SPA). There are a few SPAs for reverse engineering protocols. However, they either infer the formats of packets to be sent via imprecise abstract domain [53] or focus on cryptographic mechanisms [21]. Our approach precisely infers the format of received packets and, thus, is different from these works.

8 CONCLUSION

We propose a static analysis that can infer protocol formats with both high precision and high recall. Hence, the formats significantly enhance network protocol fuzzing, traffic auditing, and intrusion detection. Particularly, our format-inference technique has helped existing protocol fuzzers find 53 zero-day vulnerabilities.

9 LIMITATIONS AND FUTURE WORK

Our static analysis currently is implemented for C and does not support C++ due to the difficulty in analyzing virtual tables. We focus on the source code and do not handle inline assembly and libraries that do not have code available. We believe these limitations can be addressed with more engineering work. For instance, we can use class hierarchical analysis, e.g., [39], to deal with virtual tables and support C++. We can use existing disassembly techniques, e.g., [59], to support inline assembly. We leave them as our future work.

As discussed earlier, Netlifter employs existing techniques to deal with pointers and loops. Thus, it inherits their limitations. A common limitation shared by both Netlifter and all recent techniques is that the quality of inferred formats relies on the protocol implementation. For instance, if the implementation ignores a field, the output formats will ignore it, either. Nevertheless, we have shown that Netlifter is promising in practice via a set of experiments.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their comments and suggestions. This work was supported in part by DARPA VSPILLS HR001120S0058, NSF1901242 and 1910300, ONR N000141712045, N000141410468 and N000141712947. Any opinions, findings, and

conclusions in this paper are those of the authors only and do not necessarily reflect the views of our sponsors.

REFERENCES

- [1] 2018. 2018 IBM X-Force Report. <https://securityintelligence.com/2018-ibm-x-force-report-shellshock-fades-gozi-rises-and-insider-threats-soar/>. (2018).
- [2] 2022. Documents of OSDP. <https://libosdp.gotomain.io/>. (2022).
- [3] 2022. The FRRouting protocol suite. <https://github.com/FRRouting/frr>. (2022).
- [4] 2022. The IETF QUIC protocol. <https://github.com/ngtcp2/ngtcp2>. (2022).
- [5] 2022. Implementation of OSDP. <https://github.com/goToMain/libosdp>. (2022).
- [6] 2022. A lightweight TCP/IP stack. <https://github.com/lwip-tcpip/lwip>. (2022).
- [7] 2022. Linux kernel source tree. <https://github.com/torvalds/linux>. (2022).
- [8] 2022. A modern source server query protocol library written in C. <https://github.com/BinaryAlien/libsq>. (2022).
- [9] 2022. Network forensics tools and datasets. <https://github.com/MartinaZembjakova/Network-forensic-tools-taxonomy>. (2022).
- [10] 2022. Network intrusion & prevention systems. <https://www.snort.org/>. (2022).
- [11] 2022. Network protocol fuzzing. <https://github.com/jtpereyda/boofuzz>. (2022).
- [12] 2022. Packet captures. <https://packetlife.net/captures>. (2022).
- [13] 2022. A pure-python fully automated and unattended fuzzing framework. <https://github.com/OpenRCE/sulley>. (2022).
- [14] 2022. Smart card tools. <https://github.com/OpenSC/OpenSC>. (2022).
- [15] 2022. Wireshark. <https://www.wireshark.org/>. (2022).
- [16] 2023. American fuzzy lop. <https://lcamtuf.coredump.cx/afl/>. (2023).
- [17] 2023. Lifting network implementation to precise format specification. <https://github.com/qingkaishi/netlifter>. (2023).
- [18] 2023. Top-down parsing. https://wikipedia.org/wiki/Top-down_parsing. (2023).
- [19] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. 2007. *Compilers: Principles, Techniques, and Tools*. Pearson Addison Wesley.
- [20] Joao Antunes, Nuno Neves, and Paulo Verissimo. 2011. Reverse engineering of protocols from network traces. In *Working Conference on Reverse Engineering (WCRE '11)*. IEEE, 169–178.
- [21] Matteo Avalle, Alfredo Pironti, and Riccardo Sisto. 2014. Formal verification of security protocol implementations: a survey. *Formal Aspects of Computing* 26, 1 (2014), 99–123.
- [22] Osbert Bastani, Rahul Sharma, Alex Aiken, and Percy Liang. 2017. Synthesizing program input grammars. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '17)*. ACM, 95–110.
- [23] Armin Biere, Alessandro Cimatti, Edmund Clarke, and Yunshan Zhu. 1999. Symbolic model checking without BDDs. In *Proceedings of the 5th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS '99)*. Springer, 193–207.
- [24] Steve Bishop, Matthew Fairbairn, Michael Norrish, Peter Sewell, Michael Smith, and Keith Wansbrough. 2005. Rigorous Specification and Conformance Testing Techniques for Network Protocols, as Applied to TCP, UDP, and Sockets. In *Proceedings of the 2005 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM '05)*. ACM, 265–276.
- [25] Steve Bishop, Matthew Fairbairn, Michael Norrish, Peter Sewell, Michael Smith, and Keith Wansbrough. 2006. Engineering with Logic: HOL Specification and Symbolic-Evaluation Testing for TCP Implementations. In *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '06)*. ACM, 55–66.
- [26] Juan Caballero, Pongsin Poosankam, Christian Kreibich, and Dawn Song. 2009. Dispatcher: Enabling Active Botnet Infiltration Using Automatic Protocol Reverse-Engineering. In *ACM Conference on Computer and Communications Security (CCS '09)*. ACM, 621–634.
- [27] Juan Caballero, Heng Yin, Zhenkai Liang, and Dawn Song. 2007. Polyglot: Automatic Extraction of Protocol Message Format Using Dynamic Binary Analysis. In *ACM Conference on Computer and Communications Security (CCS '07)*. ACM, 317–329.
- [28] Cristian Cadar, Daniel Dunbar, Dawson R. Engler, et al. 2008. KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation (OSDI '08)*. USENIX, 209–224.
- [29] Chia Yuan Cho, Domagoj Babić, Eui Chul Richard Shin, and Dawn Song. 2010. Inference and Analysis of Formal Models of Botnet Command and Control Protocols. In *ACM Conference on Computer and Communications Security (CCS '10)*. ACM, 426–439.
- [30] Paolo Milani Comparetti, Gilbert Wondracek, Christopher Kruegel, and Engin Kirda. 2009. Prospex: Protocol specification extraction. In *2009 30th IEEE Symposium on Security and Privacy (S&P '09)*. IEEE, 110–125.
- [31] Patrick Cousot and Radhia Cousot. 1979. Systematic design of program analysis frameworks. In *Proceedings of the 6th ACM SIGACT-SIGPLAN symposium on Principles of programming languages (POPL '79)*. ACM, 269–282.
- [32] Dave Crocker and Paul Overell. 1997. *Augmented BNF for syntax specifications: ABNF*. Technical Report. RFC 2234.

- [33] Weidong Cui, Jayanthkumar Kannan, and Helen Wang. 2007. Discoverer: Automatic protocol reverse engineering from network traces. In *USENIX Security Symposium (Security '07)*. USENIX, 199–212.
- [34] Weidong Cui, Vern Paxson, Nicholas Weaver, and Randy H Katz. 2006. Protocol-independent adaptive replay of application dialog.. In *Network and Distributed System Security Symposium (NDSS '06)*. Internet Society, 1–15.
- [35] Weidong Cui, Marcus Peinado, Karl Chen, Helen J. Wang, and Luis Irun-Briz. 2008. Tupni: Automatic Reverse Engineering of Input Formats. In *ACM Conference on Computer and Communications Security (CCS '08)*. ACM, 391–402.
- [36] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. 1989. An efficient method of computing static single assignment form. In *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '89)*. ACM, 25–35.
- [37] Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS '08)*. Springer, 337–340.
- [38] Joeri De Ruiter and Erik Poll. 2015. Protocol State Fuzzing of TLS Implementations. In *USENIX Security Symposium (Security '15)*. USENIX, 193–206.
- [39] Jeffrey Dean, David Grove, and Craig Chambers. 1995. Optimization of object-oriented programs using static class hierarchy analysis. In *Proceedings of the 9th European Conference on Object-Oriented Programming (ECOOP '95)*. Springer, 77–101.
- [40] Julien Duchene, Colas Le Guernic, Eric Alata, Vincent Nicomette, and Mohamed Kaäniche. 2018. State of the art of network protocol reverse engineering tools. *Journal of Computer Virology and Hacking Techniques* 14, 1 (2018), 53–68.
- [41] Hugo Gascon, Christian Wressnegger, Fabian Yamaguchi, Daniel Arp, and Konrad Rieck. 2015. Pulsar: Stateful black-box fuzzing of proprietary network protocols. In *International Conference on Security and Privacy in Communication Systems (SecureComm '15)*. Springer, 330–347.
- [42] Rahul Gopinath, Björn Mathis, and Andreas Zeller. 2020. Mining Input Grammars from Dynamic Control Flow. In *ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '20)*. ACM, 172–183.
- [43] Matthias Hörschele and Andreas Zeller. 2016. Mining input grammars from dynamic taints. In *International Conference on Automated Software Engineering (ASE '16)*. ACM, 720–725.
- [44] Bahruz Jabiyev, Steven Sprecher, Anthony Gavazzi, Tommaso Innocenti, Kaan Onarlioglu, and Engin Kirda. 2022. FRAMESHIFTER: Security Implications of HTTP/2-to-HTTP/1 Conversion Anomalies. In *USENIX Security Symposium (Security '22)*. USENIX, 1061–1075.
- [45] James C. King. 1976. Symbolic execution and program testing. *Commun. ACM* 19, 7 (1976), 385–394.
- [46] Stephan Kleber, Henning Kopp, and Frank Kargl. 2018. NEMESYS: Network message syntax reverse engineering by analysis of the intrinsic structure of individual messages. In *USENIX Workshop on Offensive Technologies (WOOT '18)*. USENIX, 1–13.
- [47] Stephan Kleber, Rens W. van der Heijden, and Frank Kargl. 2020. Message Type Identification of Binary Network Protocols using Continuous Segment Similarity. In *IEEE Conference on Computer Communications (INFOCOM '20)*. IEEE, 2243–2252.
- [48] Yonghui Kwon, Fei Peng, Dohyeon Kim, Kyungtae Kim, Xiangyu Zhang, Dongyan Xu, Vinod Yegneswaran, and John Qian. 2015. P2C: Understanding output data files via on-the-fly transformation from producer to consumer executions. In *Network and Distributed System Security Symposium (NDSS '15)*. Internet Society, 1–14.
- [49] Patrick LaRoche, A Nur Zincir-Heywood, and Malcolm I Heywood. 2012. Network protocol discovery and analysis via live interaction. In *European Conference on the Applications of Evolutionary Computation (ECAEC '12)*. Springer, 11–20.
- [50] Chris Latner and Vikram Adve. 2004. LLVM: A compilation framework for lifelong program analysis & transformation. In *Proceedings of the 2nd International Symposium on Code Generation and Optimization (CGO '04)*. IEEE, 75:1–75:12.
- [51] Corrado Leita, Ken Mermoud, and Marc Dacier. 2005. ScriptGen: an automated script generation tool for Honeyd. In *Annual Computer Security Applications Conference (ACSAC '05)*. IEEE, 203–214.
- [52] Xiangdong Li and Li Chen. 2011. A survey on methods of automatic protocol reverse engineering. In *Proceedings of the 7th International Conference on Computational Intelligence and Security (CIS '11)*. IEEE, 685–689.
- [53] Junghee Lim, Thomas Reps, and Ben Liblit. 2006. Extracting output formats from executables. In *Working Conference on Reverse Engineering (WCRE '06)*. IEEE, 167–178.
- [54] Zhiqiang Lin, Xuxian Jiang, Dongyan Xu, and Xiangyu Zhang. 2008. Automatic protocol format reverse engineering through context-aware monitored execution. In *Network and Distributed System Security Symposium (NDSS '08)*. Internet Society, 1–15.
- [55] Zhiqiang Lin, Xiangyu Zhang, and Dongyan Xu. 2010. Reverse engineering input syntactic structure from program execution and its applications. *IEEE Transactions on Software Engineering* 36, 5 (2010), 688–703.
- [56] Min Liu, Chunfu Jia, Lu Liu, and Zhi Wang. 2013. Extracting sent message formats from executables using backward slicing. In *Proceedings of the 4th International Conference on Emerging Intelligent Data and Web Technologies (EIDWT '13)*. IEEE, 377–384.
- [57] Jian-Zhen Luo and Shun-Zheng Yu. 2013. Position-based automatic reverse engineering of network protocols. *Journal of Network and Computer Applications* 36, 3 (2013), 1070–1077.
- [58] Chris McMahon Stone, Sam L. Thomas, Mathy Vanhoef, James Henderson, Nicolas Bailluet, and Tom Chothia. 2022. The Closer You Look, The More You Learn: A Grey-Box Approach to Protocol State Machine Learning. In *ACM SIGSAC Conference on Computer and Communications Security (CCS '22)*. ACM, 2265–2278.
- [59] Kenneth Miller, Yonghui Kwon, Yi Sun, Zhuo Zhang, Xiangyu Zhang, and Zhiqiang Lin. 2019. Probabilistic disassembly. In *International Conference on Software Engineering (ICSE '19)*. IEEE, 1187–1198.
- [60] Madanlal Musuvathi, Dawson R Engler, et al. 2004. Model Checking Large Network Protocol Implementations. In *Proceedings of the 1st USENIX Symposium on Networked Systems Design and Implementation (NSDI '04)*. USENIX, 1–14.
- [61] John Narayan, Sandeep Shukla, and Charles Clancy. 2015. A survey of automatic protocol reverse engineering tools. *Comput. Surveys* 48, 3 (2015), 1–26.
- [62] P. Oehlert. 2005. Violating assumptions with fuzzing. *IEEE Security & Privacy* 3, 2 (2005), 58–62.
- [63] Markus Ring, Sarah Wunderlich, Deniz Scheuring, Dieter Landes, and Andreas Hotho. 2019. A survey of network-based intrusion detection data sets. *Computers & Security* 86, 1 (2019), 147–167.
- [64] Prateek Saxena, Pongsin Poosankam, Stephen McCamant, and Dawn Song. 2009. Loop-extended symbolic execution on binary programs. In *ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA '09)*. ACM, 225–236.
- [65] Maxim Shevertalov and Spiros Mancoridis. 2007. A Reverse Engineering Tool for Extracting Protocols of Networked Applications. In *Working Conference on Reverse Engineering (WCRE '07)*. IEEE, 229–238.
- [66] Qingkai Shi, Junyang Shao, Yapeng Ye, Mingwei Zheng, and Xiangyu Zhang. 2023. Lifting Network Protocol Implementation to Precise Format Specification with Security Applications. *arXiv preprint arXiv:2305.11781* (2023).
- [67] Baraka D Sija, Young-Hoon Goo, Kyu-Seok Shim, Huru Hasanova, and Myung-Sup Kim. 2018. A survey of automatic protocol reverse engineering approaches, methods, and tools on the inputs and outputs view. *Security and Communication Networks* 2018, 8370341 (2018), 1–17.
- [68] Yulei Sui, Sen Ye, Jingling Xue, and Pen-Chung Yew. 2011. SPAS: Scalable path-sensitive pointer analysis on full-sparse SSA. In *Proceedings of the 9th Asian Symposium on Programming Languages and Systems (APLAS '11)*. Springer, 155–171.
- [69] Fei Wang, Jianliang Wu, Yuhong Nan, Yousra Aafer, Xiangyu Zhang, Dongyan Xu, and Mathias Payer. 2022. ProFactory: Improving IoT Security via Formalized Protocol Customization. In *USENIX Security Symposium (Security '22)*. USENIX, 1–18.
- [70] Tielei Wang, Tao Wei, Guofei Gu, and Wei Zou. 2010. TaintScope: A checksum-aware directed fuzzing tool for automatic software vulnerability detection. In *IEEE Symposium on Security and Privacy (S&P '10)*. IEEE, 497–512.
- [71] Yipeng Wang, Xingjian Li, Jiao Meng, Yong Zhao, Zhibin Zhang, and Li Guo. 2011. Biprominer: Automatic mining of binary protocol features. In *International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT '11)*. IEEE, 179–184.
- [72] Yipeng Wang, Xiaochun Yun, M Zubair Shafiq, Liyan Wang, Alex X Liu, Zhibin Zhang, Danfeng Yao, Yongzheng Zhang, and Li Guo. 2012. A semantics aware approach to automated reverse engineering unknown protocols. In *IEEE International Conference on Network Protocols (ICNP '12)*. IEEE, 1–10.
- [73] Yipeng Wang, Zhibin Zhang, Danfeng Daphne Yao, Buyun Qu, and Li Guo. 2011. Inferring protocol state machine from network traces: a probabilistic approach. In *International Conference on Applied Cryptography and Network Security (ACNS '11)*. Springer, 1–18.
- [74] Zhi Wang, Xuxian Jiang, Weidong Cui, Xinyuan Wang, and Mike Grace. 2009. ReFormat: Automatic reverse engineering of encrypted messages. In *European Symposium on Research in Computer Security (ESORICS '09)*. Springer, 200–215.
- [75] Gilbert Wondracek, Paolo Milani Comparetti, Christopher Kruegel, Engin Kirda, and Scuola Superiore S Anna. 2008. Automatic network protocol analysis. In *Network and Distributed System Security Symposium (NDSS '08)*. Internet Society, 1–18.
- [76] Xiaofei Xie, Bihuan Chen, Yang Liu, Wei Le, and Xiaohong Li. 2016. Proteus: Computing disjunctive loop summary via path dependency analysis. In *ACM SIGSOFT International Symposium on the Foundations of Software Engineering (FSE '16)*. ACM, 61–72.
- [77] Yapeng Ye, Zhuo Zhang, Fei Wang, Xiangyu Zhang, and Dongyan Xu. 2021. NetPlier: Probabilistic network protocol reverse engineering from message traces. In *Symposium on Network and Distributed System Security (NDSS '21)*. Internet Society, 1–18.
- [78] Zhao Zhang, Qiao-Yan Wen, and Wen Tang. 2012. Mining Protocol State Machines by Interactive Grammar Inference. In *International Conference on Digital Manufacturing & Automation (ICDMA '12)*. IEEE, 524–527.