



### Uncovering the iceberg from the tip:

# Generating API Specifications for Bug Detection via Specification Propagation Analysis

Miaoqian Lin, Kai Chen, Yi Yang, Jinghua Liu

{linmiaoqian, chenkai, liujinghua, yangyi}@iie.ac.cn

# **Motivation: Bugs and Specifications**

- Software is built with various APIs that encapsulate complex semantics.
- To use APIs safely, developers must adhere to their specifications.
- Violations of these specifications are API misuse, which can be further exploited by attackers.



### **Motivation: Current Work and Limitations**

- Ideal way: Understand each API's behavior and infer its specifications
  - Code complexity makes this difficult
- Current work: Extract specifications from observable API artifacts
  - 1. API documents: e.g., Advance (CCS'2020), APICAD (ICSE'2023)
  - 2. API frequent usage: e.g., APP-Miner (S&P'2024), APISan (Security'2016)
  - 3. Bug patches: e.g., APHP (Security'2023)

#### • Problem: Limited by the availability of artifacts

• For example: The API bus\_find\_device requires certain specification to be followed after its call, but this is not documented.

### **Motivation: From Known to Unknown**



It is possible to infer hidden specifications (unknown) from observable specifications (known)?

# **Background: Bugs and Specifications**



The API **get\_device** obtain a reference to the device **dev**, which needs to released using **put\_device**() after use.





**path conditions** are similar across different specifications, so we treat them as implicit and exclude them in specification generation.

• The specification example



# **Insight: API Specification Propagation**

Specification propagation from get\_device to bus\_find\_device



Specifications propagate between APIs through the API call chain

# **Insight: Specification Propagation Analysis**



Use specification propagation analysis to reveal the hidden API specifications

### **Overview**

- APISpecGen: generate new specifications using seed specifications
  - Step1: Specification Propagation Analysis: identifies which APIs the specifications may propagate to (successors) or originate from (predecessors).
  - Step2: Specification Generation: generate specific specifications for inferred APIs, i.e., determine their corresponding post-operations.



# **Method: Working Example**

#### Generated specifications



- 1. After a successful **nfc\_get\_device** call
- 2. After using the **return value** of the **nfc\_get\_device** call
- 3. Check if a corresponding **nfc\_put\_device** call is performed.

Detected bug



Detected API misuse in the Linux kernel: missing **nfc\_put\_device** after **nfc\_get\_device** after the critical variable **dev** usage, causing reference count leak.

# **Method: Specification Propagation Analysis**

#### Caller Analysis to Get Successors



# **Method: Specification Generation**

- Problem: the post-operations may change during specification propagation
- Solution: utilizes API usage and data-flow validation



# **Evaluations**

- Setup
  - **Tested program:** the Linux Kernel, with more 23 million LoC
  - Seeds: six specifications collected by previous work

- Evaluations
  - Specification Generation
  - Bug Detection
  - Compared with SOTAs
  - Quality of API artifacts

Target API	Post-Operation Var		Commit
get_device	put_device	arg1	d46fe2cb2dc
device_initialize	put_device	arg1	a5808add9a6
try_module_get	module_put	arg1	44f8baaf230
kmalloc	kfree	retval	493ff661d43
kstrdup	kfree	retval	e629e7b525a
PTR_ERR	IS_ERR	arg1	59715cffce1

Six seed specification for evaluations



#### • Specification Generation and Bug Detection

	$\square$	$\overline{}$		
Seed API	#Spec	# New Bug	Bug Type	Ffficiency:
get_device	760	137	Refcount leak	Generate specifications in 2 hours
device_initialize	91	6	Refcount leak	averaging <b>1</b> second per specification.
try_module_get	58	1	Refcount leak	
kmalloc	1202	30	Memory leak	Accuracy:
kstrdup	14	1	Memory leak	Generated specifications are <b>100%</b>
ERR_PTR	5207	11	Null-pointer-deref	accurate.
Total	7332	186	-	

#### **Effectiveness:**

APISpecGen generated **7,332** specifications with from just six seed specifications, with **186** new bugs detected.

### Results

#### • SOTAs and the quality of API artifacts for specification extraction

Note: evaluated SOTAs on the specifications and bugs identified by APISpecGen			API Document	87% of the specifications are not mentioned in document	
Method	Coverage of specifications	Coverage of Bugs	API Frequent	<b>93%</b> of the specifications are not	
Advance	14%	10%	Usage	frequently followed in usage code	
IPPO	-	0%	1		
SinkFinder	11%	10%	API Name	<b>12.71%</b> of APIs lack informative	
APHP	21%	17%		subwords in their name	
SOTAs fail to extract most specifications			Artifact-based the	Artifact-based methods are inherently limited by the quality of API artifacts.	

(e.g. 87% of specifications are not documented)

### Results

#### Characteristic of Generated Specifications

- The inferred API can differ significantly from the original seed API.
- APIs evolve from the same primitive APIs and share similar specifications



from seed API get device to rdma\_user\_mmap\_entry\_get

#### Generalizability

 APISpecGen successfully generated **39** specifications for OpenSSL and **76** for FFmpeg from just **one seed each**

### **Summary**

- **Previous:** Observe API specifications externally, treating APIs individually.
- Ours: Generate specifications from internal facts, using connections between APIs.
- Idea: API Specification Propagation
- Framework: APISpecGen, generates specifications using seed specification with iterative bidirectional specification propagation analysis.
- **Results**: Generate **7332** specifications using only **6** seed specifications. **87%** of these specifications are not documented; Detect **186** new bugs in Linux kernel.



### Thanks for your attention!





