#### Enhancing Static Analysis for Practical Bug Detection: An LLM-Integrated Approach <u>Haonan Li</u>, Yu Hao, Yizhuo Zhai, Zhiyun Qian University of California, Riverside





#### Motivation: Challenges of Static Analysis In Bug Detection ChatGPT o1-preview $\sim$

- False Positives and Negatives
- Complexity, Scalability
- Limited Domain Knowledge
- Thinking: how to solve these challenges? 🤪

#### **Disadvantages:**

- 1. False Positives and Negatives:
  - **Overwhelming Warnings:** Static analysis tools may generate a high number of false positives, leading to "warning fatigue" where developers start ignoring the tool's output.
  - Missed Bugs: Conversely, they might miss certain bugs (false negatives), especially those that depend on runtime behavior.
- 2. Limited Dynamic Context:
  - Runtime Behavior Ignored: Static analysis cannot fully account for dynamic aspects like user input, environment variables, or runtime configurations.
  - **Concurrency Issues:** Bugs related to multi-threading and concurrency are often hard to detect statically.
- 3. **Complexity and Scalability Issues:** 
  - Analysis Time: For very large or complex codebases, static analysis can be time-consuming, potentially slowing down the development process.
  - Resource Intensive: Requires significant computational resources for deep analysis.
- 4. Difficulty with Modern Language Features:
  - Dynamic Languages: Languages that are highly dynamic (e.g., JavaScript, Python) pose challenges due to features like dynamic typing and reflection.



#### Enhance Static Analysis An LLM-integrated Approach

- Compared to static analysis, LLM has: •
  - **V** Domain knowledge & heuristics
  - Smart summary for complex code
  - X Unstable response
  - X Limited reasoning on code
- Problem: Can we enhance static analysis by integrating LLM?

# Case Study: UBITect

- UBITect [1]: finding Use-Before-Initialization (UBI) bugs in the Linux kernel
- UBI: variables may be used before they are initialized.
- Shows the challenges of static analysis:
  - FP and FN

[1] Yizhuo Zhai, Yu Hao, et.al. "UBITect: A Precise and Scalable Method to Detect Use-before-Initialization Bugs in Linux Kernel." In FSE'20. https://doi.org/10.1145/3368089.3409686.

**int** F(){ int v; // declaration of v 2 **if** (constraint) 3 init(&v); // initializer of v 4 use(v); // use of v 5 6 } A typical UBI bug, v may not be initialized before reaching Line 5

<sup>• • •</sup> 

#### Case Study: UBITect False Positives and Negatives

#### **Two Stages:**

- Static analysis (scalable, imprecise)
- 2. Symbolic execution (precise, slow)

#### **Dilemma:**

- **FP:** >95% FP in <u>static analysis</u>; utilizing symbolic execution to filter out FP
- **FN:** discarding 40% undecided cases, may ignore <u>40% potential bugs</u> (due to time/memory limitation)





The overview of UBITect.

### Motivating Example: **sscanf**

- static int libcfs\_ip\_str2addr(...){ 1 **unsigned int** a, b, c, d; 2 if (sscanf(str, "%u.%u.%u.%u%n", 3 &a, &b, &c, &d, &n) >= 4 && ...){ 4 // use of a, b, c, d 5 6 7
- Static analysis: sscanf <u>may not initialize</u> a, **b**, **c**, **d** (path-insensitive)
  - Report a UBI bug at Line 5
- Symbolic execution: timeout

int vsscanf(const char \*buf, const char \*fmt,  $\leftrightarrow$  va\_list args){ const char \*str = buf; 15 16 . . . while (\*fmt) { 17 switch (\*fmt++) { 18 **case** 'c': { 19 char \*s = (char \*)va\_arg(args, char\*); 20 21. . . **do** { \*s++ = \*str++; } 22 while (...\*str); num++; 23 } 24 } 25 26 . . . num++; 27 28 . . . 29 return num; 30 31

# Understand **sscanf** with LLM

```
static int libcfs_ip_str2addr(...){
1
      unsigned int a, b, c, d;
2
      if (sscanf(str, "%u.%u.%u.%u%n",
3
            &a, &b, &c, &d, &n) >= 4 && ...){
4
              // use of a, b, c, d
5
6
7
```

Ask ChatGPT: "are variables a, b, c, d to be initialized before reaching Line 5?"

• ... sscanf here to parse IP address ... by the time the code reaches the line ... a, b, c, **d** have been initialized



Scan me to see the chat!

## LLift: LLM for Undecided Cases

- LLM-based framework: LLift.
- Asking LLM, "are these variables initialized before their use"



"must\_init" -> not a bug

# Directly Ask GPT

- Challenge: "simple prompt" doesn't work
  - Bug detection requires multiple steps w/ rigorous rules
  - LLM is unaware of them
- **SOLUTION:** Teach LLM the steps and rules of bug detection

#### Combinat

Simple Pro PGA PGA+PP PGA+PP+S PGA+PP+7 PGA+PP+7

Oracle

tion	TN(C)	TP(C)	Precision	Recall	Accuracy	F1 Score
ompt	12(9)	2(1)	0.12	0.15	0.35	0.13
	13(9)	5(1)	0.26	0.38	0.45	0.31
	5(3)	6(1)	0.21	0.46	0.28	0.29
SV	5(2)	11(8)	0.33	0.85	0.40	0.48
TD	22(14)	6(4)	0.55	0.46	0.70	0.50
TD+SV	25(17)	13(12)	0.87	1.00	0.95	0.93
	27(27)	13(13)	-	-	-	-

# Design Principles

- Task Decomposition (TD): break the task into small pieces
- Post-constraint Guided Path Analysis (PGA)

```
int I(){
 1
        if (...){ //path<sub>1</sub>
 2
          return 0;
 3
 4
        else{ //path<sub>2</sub>
 5
          return -1;
 6
 7
 8
     void F(){
 9
          int v;
10
          int ret = I(v);
11
          if (ret == 0){
12
               use(v);
13
          }
14
15
```

# Design Principles

- Task Decomposition (TD): break the task into small pieces
- Post-constraint Guided Path Analysis (PGA)
- Progressive Prompt (PP): Let LLMs decide the what they want (for function definitions)



#### Workflow

- Identify the initializer
- Extract post-constraints
- Analyze the initializer with post-constraint guidance

1	static int
2	unsigned
3	if (ssca
4	&a
5	
6	}
7	}
8	
9	int sscanf
10	va_start
11	i = vssca
12	va_end(a
13	}



### Evaluation

- Precision and Recall
- Other Models than GPT-4
- Other projects than Linux
- Contributions of each design
   component

- Dataset: Rnd-300
  - randomly choose 300 from 53,000 <u>undecided</u> cases

### Precision & Recall

- **Precision:** 5 TP, 50% precision
- **Recall:** no missed bugs found in Rnd-300
- Extend to 1000, 13 TP and 13 FP
- 4 of them are confirmed as real bugs

Initializer	Caller	File Path	Variable	Line
read_reg	get_signal_parameters	drivers/media/dvb-frontends/stv0910.c	tmp	504
regmap_read	isc_update_profile	drivers/media/platform/atmel/atmel-isc.c	sr	664
ep0_read_setup	ep0_handle_setup	drivers/usb/mtu3/mtu3_gadget_ep0.c	setup.bRequestType	637
regmap_read	mdio_sc_cfg_reg_write	drivers/net/ethernet/hisilicon/hns_mdio.c	reg_value	169
bcm3510_do_hab_cmd	bcm3510_check_firmware_version	drivers/media/dvb-frontends/bcm3510.c	ver.demod_version	666
readCapabilityRid	airo_get_range	drivers/net/wireless/cisco/airo.c	cap_rid.softCap	6936
e1e_rphy	e1000_resume	drivers/net/ethernet/intel/e1000e/netdev.c	phy_data	6580
pci_read_config_dword	adm8211_probe	drivers/net/wireless/admtek/adm8211.c	reg	1814
lan78xx_read_reg	lan78xx_write_raw_otp	drivers/net/usb/lan78xx.c	buf	873
t1_tpi_read	my3126_phy_reset	drivers/net/ethernet/chelsio/cxgb/my3126.c	val	193
pci_read_config_dword	quirk_intel_purley_xeon_ras_cap	arch/x86/kernel/quirks.c	capid0	562
ata_timing_compute	opti82c46x_set_piomode	drivers/ata/pata_legacy.c	&tp	564
pt_completion	pt_req_sense	drivers/block/paride/pt.c	buf	368

True positives of LLift on Rnd-300. Above is 5TP on Rnd-300, below is 8TP scaled up to 1000

### Other Models

- Disclaimer: not reported in the paper
- LLift works on (most) other models
- Llama shows usable results

Model	Simple Prompt			All Components		
widdei	TN	TP	F1 Score	TN	TP	F1 Score
CodeLlama (34b)	23	11	0.23	42	12	0.31
Llama 2 (70b)	0	13	0.21	49	11	0.31
Palm 2	12	3	0.06	40	6	0.16
Claude 2	64	4	0.17	80	2	0.13
GPT-3.5 (0613)	3	13	0.22	37	9	0.23
GPT-4 (0613)	51	2	0.07	91	13	0.87
Oracle	95	13	-	95	13	-

# Summary

- We present LLift, the first work to show how to use LLM to improve static analysis in practice.
- LLift introduces new design strategies (PGA, PP), which can also help future studies.
- We test LLift on real-world codebase, achieving good precision and recall, and finding 13 new TP.



Learn more details: sites.google.com/view/llift-open