



OCTOBER 1-2, 2020

BRIEFINGS

Demystify Today's Binary Disassembly and How Modern ABI Makes it Easier



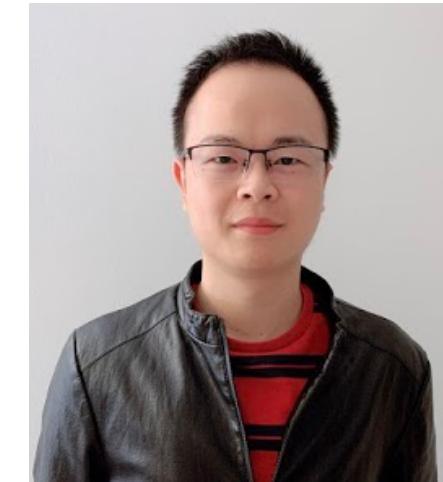
Who Are We



Chengbin Pang
- Visiting Scholar,
Stevens Institute of
Technology
- PhD Student,
Nanjing University



STEVENS
INSTITUTE OF TECHNOLOGY
THE INNOVATION UNIVERSITY®



Jun Xu
- Assistant Professor,
Stevens Institute of
Technology



STEVENS
INSTITUTE OF TECHNOLOGY
THE INNOVATION UNIVERSITY®



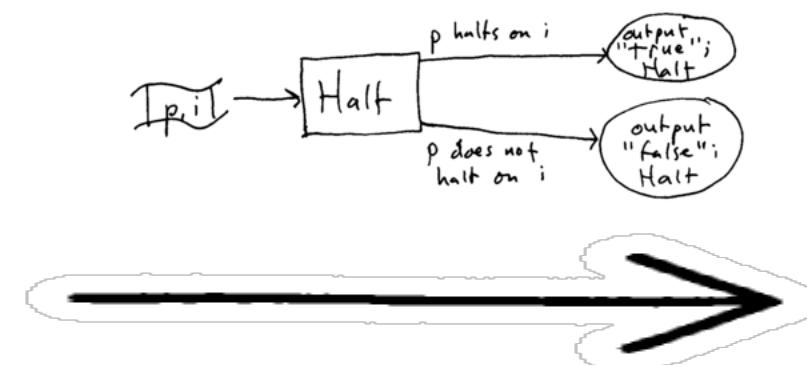
Eric Koskinen
- Assistant Professor,
Stevens Institute of
Technology



Begin of Our Story



X86/X64 binaries



Temporal properties

Image source:

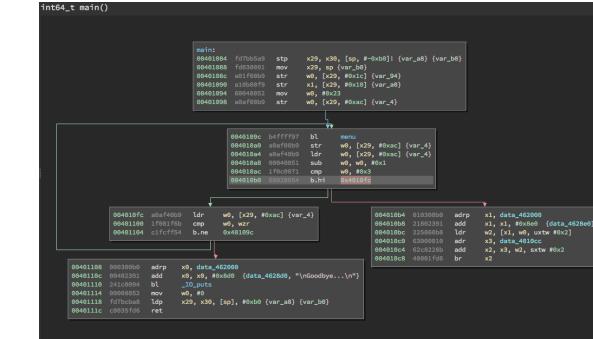
- [1] <https://medium.com/@bfil/how-code-becomes-legacy-2e6b58e839eb>
- [2] <https://armacad.info/winter-school-engineering-and-computer-science-on-formal-verification-17-21-december-2017-israel>



OK, Let's Disassemble the Code First



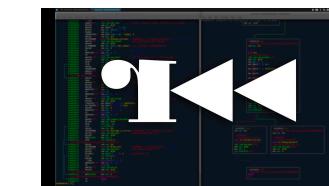
BUT



12 GNU Binutils Tools



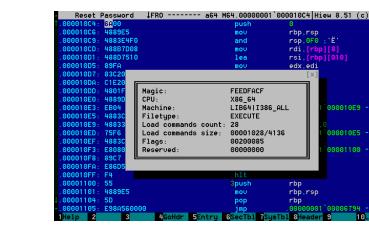
GHIDRA



WHICH ONE?

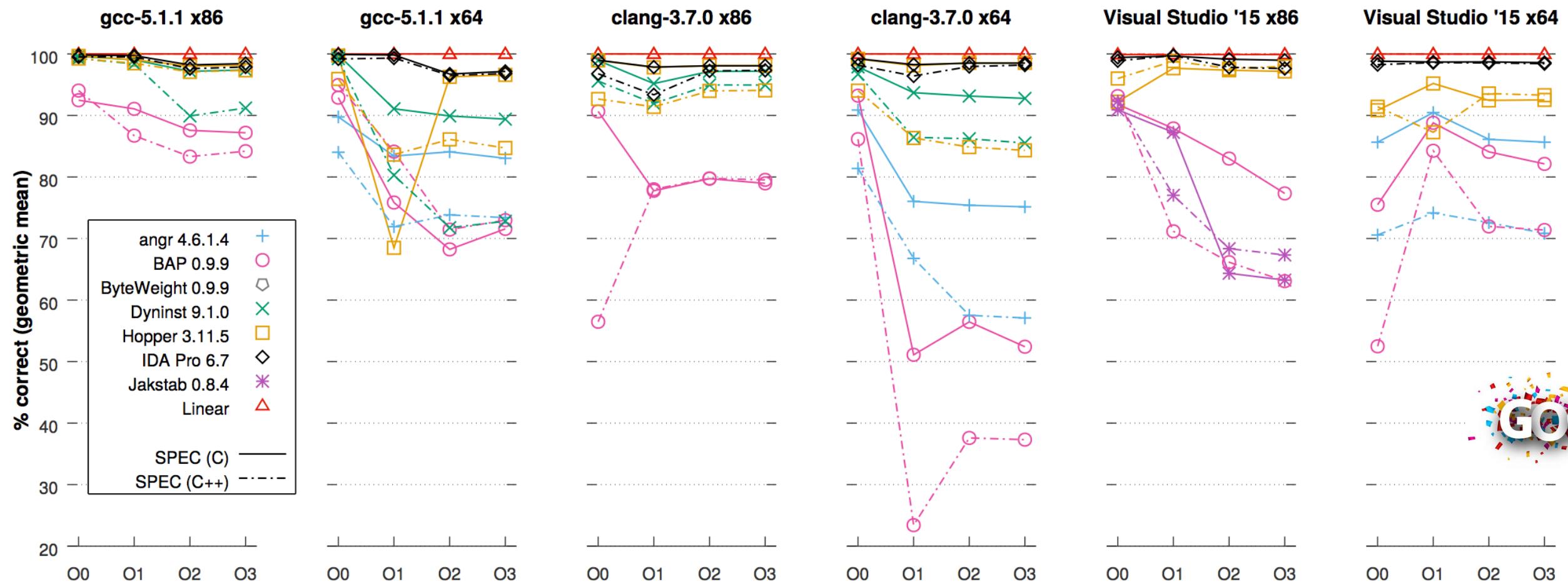


BINARYNINJA



Dyn_{inst}

Many Disassemblers are *Super Good*



(a) Correctly disassembled instructions.

GOOD NEWS



But Not Good Enough for Verification



```

1 ; wrong decoding
2 640069: add [%rax],%eax
3 64006b: add [%rax],%cl
4 ...
5 640126: call unwind
6 ; non-return

```

Wrong Instructions



```

<autohelperowl.defendpat156>:
480970: push %rbp
480971: push %r15
480973: push %r14
480975: push %rbx
480976: push %rax

```

Dyn inst



```

which){
case 't': ...
...
default:
undefined();
}

```

Wrong Pointers



```

1 mov Jx6ab8a0, %esi
2 mov %rbx, %rdi
3
4 ; data;
5 0x6ab8a0: 69 00 64 00
6 ...

```

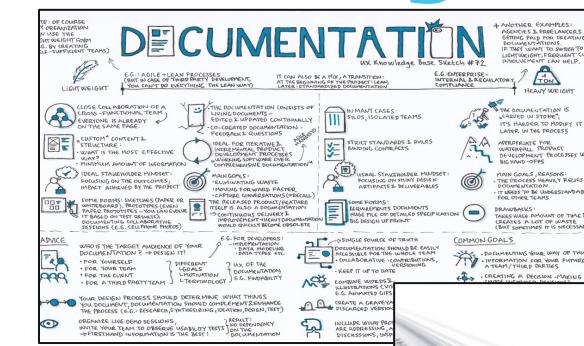
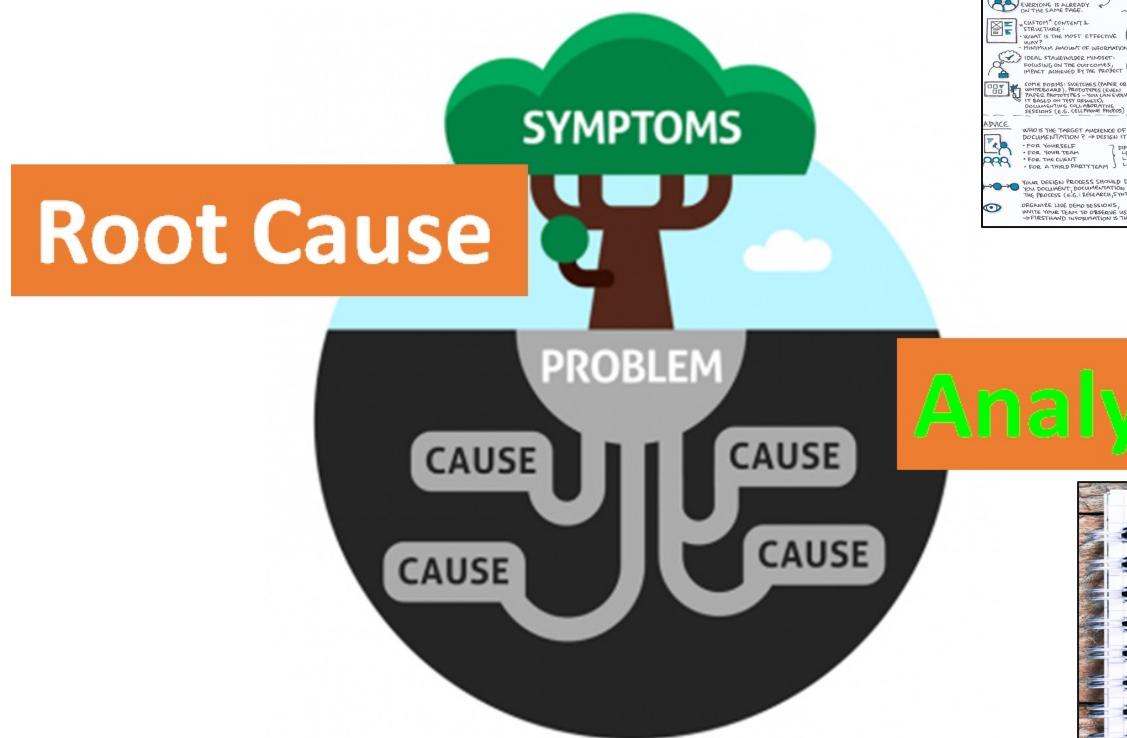
Wrong Control Flows



Errors from disassembly will propagate all the way to the verification results.



Let's Figure Out Why Errors and Fix Them

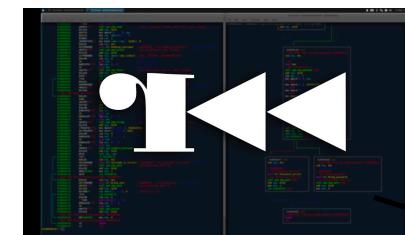


Fine, Let's Read Their Source Code



- ONJDUMP
- UROBOROS
- PSI

12 GNU Binutils Tools



- RADARE@



- BAP



- GHIDRA



- MCSEMA



- ANGR

*Dyn
inst*

- DYNINST



Now, We Know Everything About The Disassemblers

15 principled algorithms
(correctness guaranteed)



29 heuristics
(sorry, no guarantee)



<https://arxiv.org/abs/2007.14266>

Alg.	<i>Algorithms & Heuristics</i>		Goals	Tools
Disassembly	Linear Sweep	① Start from code addresses with symbols ② Continuous scanning for instructions ③ Skip bad opcodes ④ Replace padding and re-disassembly ⑤ Exclude code around errors	Code accuracy Code coverage	OBJDUMP, PSI, UROBOROS OJDUMP, PSI, UROBOROS
	Recursive Descent	② Follow control flow to do disassembly ③ Start from program entry, main, and symbols ④ Function entry matching ⑤ Linear sweep code gaps ⑦ Disassembly from targets of xrefs	Code accuracy Code coverage Code coverage Code coverage	DYNINST, GHIDRA, ANGR, BAP, RADARE2 DYNINST, GHIDRA, ANGR, BAP, RADARE2 DYNINST, GHIDRA, ANGR, BAP, RADARE2 ANGR GHIDRA, RADARE2
Symbolization	Xrefs	④ Exclude data units that are floating points ⑧ Brute force operands and data units ⑨ Pointers in data have machine size ⑩ Alignment of pointers in data ⑪ Pointers in data or referenced by other xrefs can be non-aligned ⑫ References to code can only point to function entries ⑬ Enlarge boundaries of data regions ⑭ Address tables have minimal size of 2 ⑮ Exclude pointers that may overlap with a string ⑯ While scanning data regions, use step-length based on type inference	Xref accuracy Xref coverage Xref accuracy Xref accuracy Xref coverage Xref accuracy Xref coverage Xref accuracy Xref coverage Xref accuracy Xref coverage	ANGR UROBOROS, MCSEMA, GHIDRA, ANGR UROBOROS, MCSEMA, GHIDRA, ANGR UROBOROS, MCSEMA, GHIDRA GHIDRA, ANGR GHIDRA GHIDRA, ANGR GHIDRA MCSEMA, GHIDRA ANGR
		⑤ Identify main based on arguments to __libc_start_main ⑯ Identify main using patterns in _start/_scrt_common_main_seh	Func coverage Func coverage	ANGR, BAP DYNINST, RADARE2
Function Entry	MAIN Function	⑥ Identify function entries based on symbols ⑦ Identify function entries based on exception information ⑧ Identify function entries based on targets of direct calls ⑨ Identify function entries by resolving indirect calls ⑩ Identify function entries based on prologues/decision-tree ⑪ Consider begins of code discovered by linear scan as function entries	Func coverage Func coverage Func coverage Func coverage Func coverage Func coverage	DYNINST, GHIDRA, ANGR, BAP, RADARE2 GHIDRA DYNINST, GHIDRA, ANGR, BAP, RADARE2 GHIDRA, ANGR DYNINST, GHIDRA, ANGR, BAP, RADARE2 ANGR
		⑫ Identify function entries based on patterns in _start/_scrt_common_main_seh	Func coverage	DYNINST, RADARE2
CFG	Indirect Jump	⑩ Use VSA to resolve jump table targets ⑪ Follow patterns to determine jump tables ⑫ Discard jump tables with index bound larger than a threshold ⑬ Restrict the depth of slice for VSA	CFG accuracy CFG accuracy CFG accuracy Efficiency	DYNINST, GHIDRA, ANGR DYNINST, GHIDRA, RADARE2 GHIDRA, ANGR, RADARE2 DYNINST, ANGR
		⑭ Identify targets based on constant propagation ⑮ Consider a jump to the start of another function as a tail call ⑯ Determine tail call based on distance between the jump and its target ⑰ A tail call and its target cross multiple functions ⑱ Tail calls cannot be conditional jumps ⑲ A tail call tears down its stack ⑳ A tail call does not jump to the middle of a function ㉑ Target of a tail call cannot be target of any conditional jumps	CFG coverage CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy	GHIDRA, ANGR DYNINST, ANGR RADARE2 GHIDRA GHIDRA, ANGR DYNINST, ANGR ANGR ANGR
Tail Call	Indirect Call	㉒ Identify targets based on constant propagation ㉓ Consider a jump to the start of another function as a tail call ㉔ Determine tail call based on distance between the jump and its target ㉕ A tail call and its target cross multiple functions ㉖ Tail calls cannot be conditional jumps ㉗ A tail call tears down its stack ㉘ A tail call does not jump to the middle of a function ㉙ Target of a tail call cannot be target of any conditional jumps	CFG coverage CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy CFG accuracy	GHIDRA, ANGR DYNINST, ANGR RADARE2 GHIDRA GHIDRA, ANGR DYNINST, ANGR ANGR ANGR
		㉚ Identify system calls or library functions that are known non-returning ㉛ Identify functions with no ret and no tail calls that return ㉜ Identify functions that always call non-returning functions ㉝ Detect non-returning functions based on fall-through after the call-sites	CFG accuracy CFG accuracy CFG accuracy CFG accuracy	DYNINST, GHIDRA, ANGR, BAP DYNINST, ANGR, RADARE2 BAP GHIDRA

Examples of Error-Prone Heuristics

- Linear scanning code gaps to recover instructions
- Considering integers that point to legitimate locations as pointers
- Finding functions in code gaps with common signatures of prologues
- Detecting jump tables based on common patterns



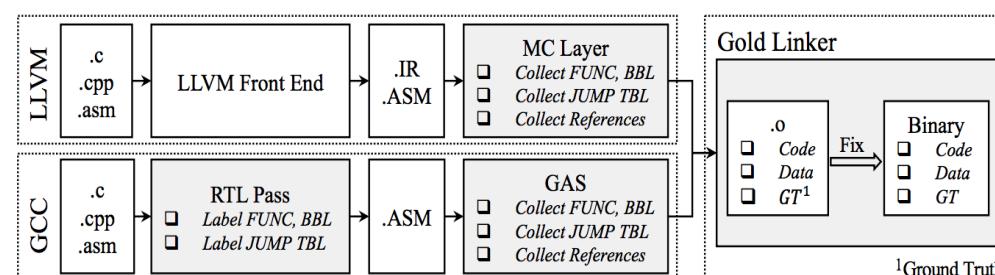
But Which Strategies Make What Errors?

Well, Let's Run An Evaluation



Type	Name	Programs/Binaries	
		Linux	Windows
Benchmark	SPEC CPU2006	30 / 546	15 / 120
Utilities	Unzip-6.0 Coreutils-8.30 7-zip-19 Findutils-4.4 Binutils-2.26 Tiff-4.0	125 / 2500	26 / 196
Clients	Openssl-1.1.01 Putty-0.73 D8-6.4 Filezilla-3.44.2 Busybox-1.31 Protobuf-c-1 ZSH-5.7.1 VIM-8.1 XML2-2.9.8 Openssh-8.0 Git-2.23	13 / 154	13 / 104
Servers	Lighttpd-1.4.54 Mysqld-5.7.27 Nginx-1.15.0 SQLite-3.32.0	3 / 49	2 / 16
Libraries	Glibc-2.27 libpcap-1.9.0 libv8-6.4 libtiff-4.0.10 libxml2-2.9.8 libsdl-3.32.0 libprotobuf-c-1.3.2	6 / 79	3 / 24
Total		177 / 3328	59 / 460

About 4K binaries built from scratch



Tracing the compilation step-by-step



Test each strategy separately

- ✓ Measure its contribution
- ✓ Evaluate its accuracy
- ✓ Understand its errors

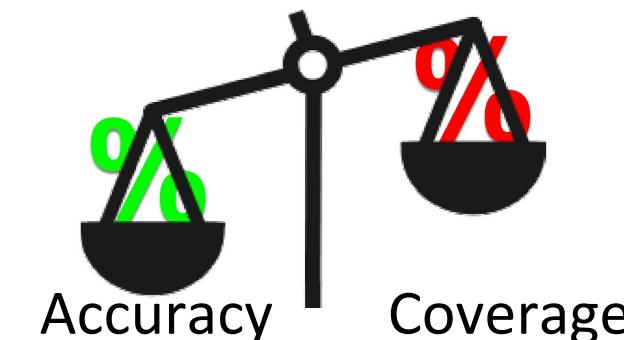


No More Excuses of Reusing Existing Disassemblers

- ❑ Heuristics are heavily used
 - Nearly 50% the instructions are recovered with heuristics
 - About 20% of the functions are detected by heuristics
 - ...

- ❑ All Heuristics Introduce Errors
 - Which Heuristics Cause Which Errors Are Understood
(<https://arxiv.org/abs/2007.14266>)

- ❑ Fixing The Errors Will Reduce the Coverage!!!



Did Our Study Kill Our Research? NO!

Findings

- ✓ System-v x64 ABI mandates a call frame for each function for exception handling.
- ✓ The call frame gives the start address of the corresponding function
- ✓ Call frames persist even in stripped, commercial binaries

skype-8.55.0.141
trillian-6.1.0.5
opera-65.0.3467.69
yandex-browser-19.12.3
SpiderOakONE-7.5.01
slack-4.2.0
rainlendar2-2.15.2
sublime-3211
netease-cloud-music-1.2.1
wps-11.1.0.8865
wpp-11.1.0.8865
wpspdf-11.1.0.8865
wpsoffice-11.1.0.8865
ida64-7.2
zoom-7.19.2020
binaryninja-1.2
FoxitReader-4.4.0911

SVABI
System V Application Binary Interface



Exploiting Call Frames to Do Binary Disassembly

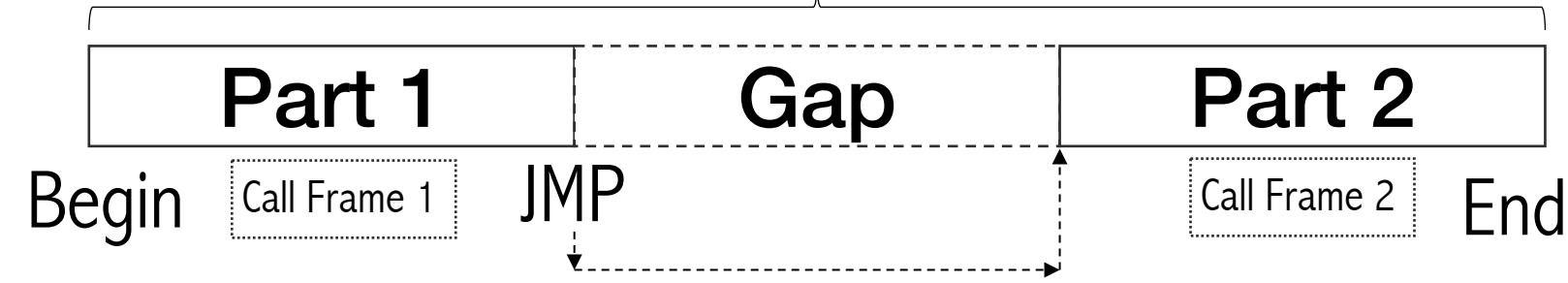
- ❑ Extract function starts from call frames
- ❑ Recursive disassembly from each function start
- ✓ Recover instructions, functions, and control flows
- ✓ Avoid using any heuristics
- ❑ Job done!



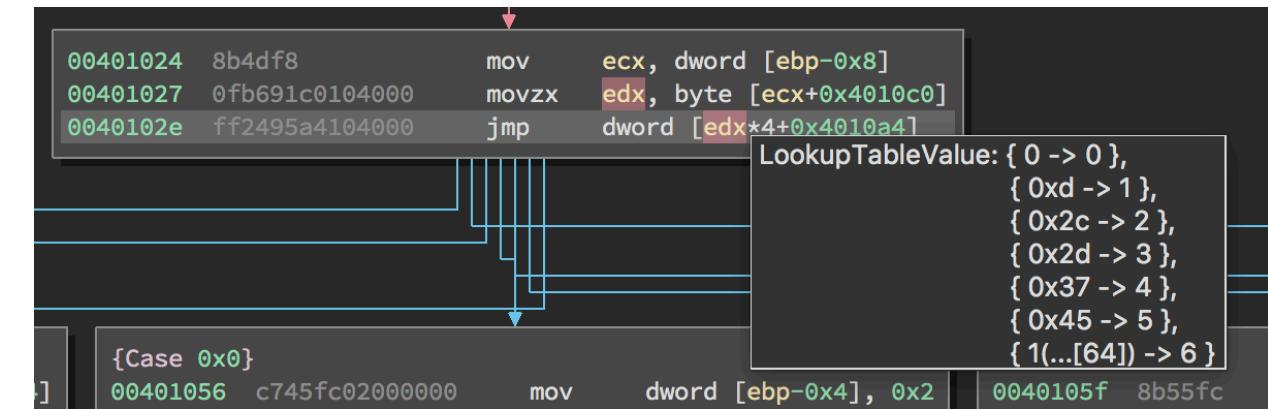
Wait, Reality is Always More Complicated



Non-continuous Functions



Jump Tables



Solution

- Checking the JMP is not a tail call
- Stack height at the jump site is not 0 (or)
- Part 2 does not meet calling conventions (or)
- Part 2 is not referenced in other places

Solution

- Intra-procedural value set analysis
to decide targets of jump tables

Now, The Results are Verification Friendly

OPT	ANGR		DYNINST		GHIDRA		BAP		RADARE2		NUCLEUS		BYTEWEIGHT		IDA PRO		BINARY NINJA		Symbol		FETCH	
	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R
O2	83.66	95.25	90.85	63.71	97.13	99.75	75.29	54.20	97.47	50.89	97.79	90.21	96.38	96.51	99.92	72.67	97.72	92.31	97.29	99.99	99.96	99.79
O3	83.56	95.63	90.41	62.09	96.38	99.78	73.20	55.55	97.18	50.61	97.42	91.98	97.91	97.75	99.92	75.41	97.63	94.03	96.52	99.99	99.96	99.78
Os	92.61	96.42	92.76	69.65	99.84	99.91	86.39	66.54	98.84	69.66	97.64	91.23	97.83	96.05	99.98	81.94	98.75	98.81	99.98	99.99	99.99	99.96
Of	81.89	95.16	90.71	62.73	96.56	99.77	73.58	56.30	98.14	50.77	96.00	91.03	97.02	96.90	99.95	76.19	97.75	94.22	96.67	99.99	99.91	99.81
Avg.	85.43	95.62	91.18	64.55	97.48	99.80	77.12	58.15	97.92	55.23	97.21	91.11	97.29	96.80	99.94	85.28	97.84	94.84	97.62	99.99	99.96	99.84

Our Tool

Example: Results of Function Detection (P = Precision; R = Recall)



The missing cases or errors
are harmless!



Remaining Problems for Tomorrow

- Detection of pointers from binary code
- Handling binaries other than x64 system-v binaries
- Engineering issues in the tools we are running

Concluding Remarks

- ❑ We fully demystify the strategies used by existing disassemblers
 - ✓ Deepen the understanding of the community about binary disassembly
 - ✓ Help users gain awareness of the risks of existing disassemblers
- ❑ We unveil the performance and pitfalls of each strategy of each disassembler
 - ✓ Bring insights for the community to further improve binary disassembly
 - ✓ Help users select the correct tools to best satisfy their demands
- ❑ We exploit call frames to achieve verification-friendly binary disassembly
 - ✓ Provide an alternative approach for users that have high requirements of correctness
 - ✓ Bring insights of leveraging extra sources of information to aid disassembly

Software and Data Release

- ❑ Findings of our study
 - ✓ <https://arxiv.org/abs/2007.14266>
- ❑ Data & code of our study (publicly available)
 - ✓ <https://github.com/junxzm1990/x86-sok>
- ❑ Code of our new disassembler (private, available on request)
 - ✓ https://github.com/ruotongyu/eh_frame

Thank you for listening

