

QuickBinDiff: Efficient Binary Code Diff Tool Based on Command Line Interface

Background

- Source code diff
 - Version control and collaboration
 - Code review
 - Bug tracking and fixing
 - Detect plagiararism and infringement
 - Detect bug or vulnerabilities
- We also use it extensively:
 - Git diff when commit/push/pull/merge

Background

- Source code diff
 - Upside-down view

```
11c11,12
< export CTNG_CONF_PATH CTNG_BIN CTNG_PATH CTNG_TARBALL_PATH
---
> EXTRA_DEP_PATH="$T00L_PATH/extra_dep"
> export CTNG_CONF_PATH CTNG_BIN CTNG_PATH CTNG_TARBALL_PATH EXTRA_DEP_PATH
14c15
< NUM_JOBS=8
---
> NUM_JOBS=8
```

Background

- Source code diff
 - Side-by-side view

```
CTNG_PATH="$TOOL_PATH/crosstool-ng"
CTNG_CONF_PATH="$PROJ_ROOT/ctng_conf"
CTNG_TARBALL_PATH="$TOOL_PATH/ctng_tarballs"
export CTNG_CONF_PATH CTNG_BIN CTNG_PATH CTNG_TARBALL_PATH

mkdir -p "$CTNG_TARBALL_PATH"

NUM_JOBS=8
MAX_JOBS=8
export NUM_JOBS MAX_JOBS
```

```
CTNG_PATH="$TOOL_PATH/crosstool-ng"
CTNG_CONF_PATH="$PROJ_ROOT/ctng_conf"
CTNG_TARBALL_PATH="$TOOL_PATH/ctng_tarballs"
| EXTRA_DEP_PATH="$TOOL_PATH/extra_dep"
> export CTNG_CONF_PATH CTNG_BIN CTNG_PATH CTNG_TARBALL_PATH EX
mkdir -p "$CTNG_TARBALL_PATH"

| NUM_JOBS=8
| MAX_JOBS=8
| export NUM_JOBS MAX_JOBS
```

Background

- Binary code diff
 - Malware detection and analysis
 - Code plagiararism and infringement
 - Find bug or vulnerability
 - binary code similarity
 - binary diversifying

Problems

1. The absence of a free, open source, CLI-based diff tools focused on assembly code diff
- GUI-based diff tools focus on pseudo code diff
 - Slow and automation is not possible
 - Most tools require a paid license
 - DarunGrim
 - IDA Pro required
 - diaphora
 - IDA Pro required
 - bindiff
 - IDA Pro required to utilize all features
 - IDA free, Binary Ninja, Ghidra required to utilize limited features

Problems

1. The absence of a free, open-source, CLI-based diff tools focused on assembly code diff
 - GLI-based diff tools focus on **bytecode diff**
 - Significant loss of information and difficult to understand.
 - The only CLI-based assembly code diff tool (elf_diff) requires source code
 - compilation needed including debugging information

Problems

1. The absence of a free, open-source, CLI-based diff tools focused on assembly code diff

- GLI-

- Sign

- The

requ

- co

```
00001a0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001b0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001c0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001d0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001e0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001f0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000200: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000210: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000220: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000230: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000240: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000250: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000260: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000270: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000280: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000290: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002a0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002b0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002c0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002d0: 0000 0000 0000 0000 60e1 c552 6c5b 94d9 .....
00002e0: 2093 f53f c3d8 4290 8392 f53f dd07 20b5 ...?.B...
00002f0: 8993 f53f 6d73 637a c292 f53f b81e 85eb ...?mscz...
0000300: 51b8 1d40 60e1 c552 0000 0000 a7e1 c552 Q..@'.R...
0000310: 3485 ce6b ec92 f53f df37 bef6 cc92 f53f 4..k...?.7.
0000320: 6c43 c538 7f93 f53f 170e 8464 0193 f53f 1C.8...?.
0000330: 85eb 51b8 1e45 3040 a7e1 c552 0000 0000 ..Q..E00...
0000340: dde1 c552 89d2 dee0 0b93 f53f df4f 8d97 ...R.....
0000350: 6e92 f53f dd07 20b5 8993 f53f c3d8 4290 n..?...
0000360: 8392 f53f 5c8f c2f5 285c 1140 dde1 c552 ...?\...(\.
0000370: 0000 0000 1fe2 c552 c3d8 4290 8392 f53f .....R..B
0000380: c3d8 4290 8392 f53f c190 d5ad 9e93 f53f ..B....?.
0000390: c3d8 4290 8392 f53f 85eb 51b8 1e85 eb3f ..B....?.Q
00003a0: 1fe2 c552 0000 0000 55e2 c552 183e 22a6 ...R....U..
00003b0: 4492 f53f 183e 22a6 4492 f53f 87a2 409f D..?>".D..
00003c0: c893 f53f 183e 22a6 4492 f53f 14ae 47e1 ...?>".D..
00003d0: 7a14 fe3f 55e2 c552 0000 0000 8ce2 c552 z..?U..R...
00001a0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001b0: 0000 0000 0000 0000 ea03 0000 0000 0000 .....
00001c0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001d0: 0000 0000 0000 0000 805a c352 00eb a154 .....Z.
00001e0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00001f0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000200: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000210: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000220: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000230: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000240: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000250: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000260: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000270: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000280: 0000 0000 0000 0000 0000 0000 0000 0000 .....
0000290: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002a0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002b0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002c0: 0000 0000 0000 0000 0000 0000 0000 0000 .....
00002d0: 0000 0000 0000 0000 0000 b418 c252 0000 0000 .....
00002e0: 2c65 19e2 5817 f63f 2c65 19e2 5817 f63f ,e..X..?.e.
00002f0: bbb8 8d06 f016 f63f bbb8 8d06 f016 f63f .....?.
0000300: 0000 0000 0000 0000 ef18 c252 0000 0000 .....
0000310: f018 c252 0000 0000 bbb8 8d06 f016 f63f ...R.....
0000320: 2c65 19e2 5817 f63f bbb8 8d06 f016 f63f ,e..X..?.
0000330: bbb8 8d06 f016 f63f 0900 0000 0000 0000 .....?.
0000340: f018 c252 0000 0000 2c19 c252 0000 0000 ...R.....
0000350: bbb8 8d06 f016 f63f 2c65 19e2 5817 f63f .....?.e.
0000360: bbb8 8d06 f016 f63f 2c65 19e2 5817 f63f .....?.e.
0000370: 0000 0000 0000 0000 f118 c252 0000 0000 .....
0000380: 6819 c252 0000 0000 2c65 19e2 5817 f63f h..R....,e.
0000390: 2c65 19e2 5817 f63f bbb8 8d06 f016 f63f ,e..X..?.
00003a0: bbb8 8d06 f016 f63f 0c00 0000 0000 0000 .....?.
00003b0: f218 c252 0000 0000 a419 c252 0000 0000 ...R.....
00003c0: bbb8 8d06 f016 f63f bbb8 8d06 f016 f63f .....?.
00003d0: bbb8 8d06 f016 f63f bbb8 8d06 f016 f63f .....?.
```


Problems

2. Difficult to automatically handle a large quantity of binaries

- Assuming the license has been purchased
- To check if the target binary group has additional pop-ret gadgets,
 - Open GUI-based diff tool
 - Load two binaries
 - Wait for analysis
 - Manually verify and record the findings

Problems

3. The lack of research on the readability of binary code diff algorithms
- Research on readability in source code diff from a software engineering perspective
 - “How Different Are Different diff Algorithms in Git?”
 - The histogram algorithm excels in the readability of source code diff
 - There is a absence of such research for binary code, despite the relative difficulty in understanding binary code

Problems

1. The absence of a free, open source, CLI-based diff tools focused on assembly code diff
2. Difficult to automatically handle a large quantity of binaries
3. The lack of research on the readability of binary code diff algorithms

Goal

- Low-quality CLI tools ← QuickBinDiff → High-quality GUI tools
- Why not CLI tools?
 - Only support bytecode diff or need source code to recompile
- Why not GUI tools?
 - Requires a license, slow, and not automatable
 - "Even though we have IDA Pro, we often use gdb or objdump"

Goal

- Free, open source and CLI-based binary code diff tool
 - More detail than bytecode diff
 - More faster then pseudo-code diff
- For Efficiency, utilize text diff algorithm
 - Select an algorithm for assembly matching offers good readability.
- Support cross-platform, various file foramt and multi architecture
 - Platform: Linux, Windows, OS X
 - File format: PE, ELF, Match-O
 - Architecture: X86, X64, ARM, MIPS
- Utilize various information from the binary
 - including Sections, function boundaries, instructions

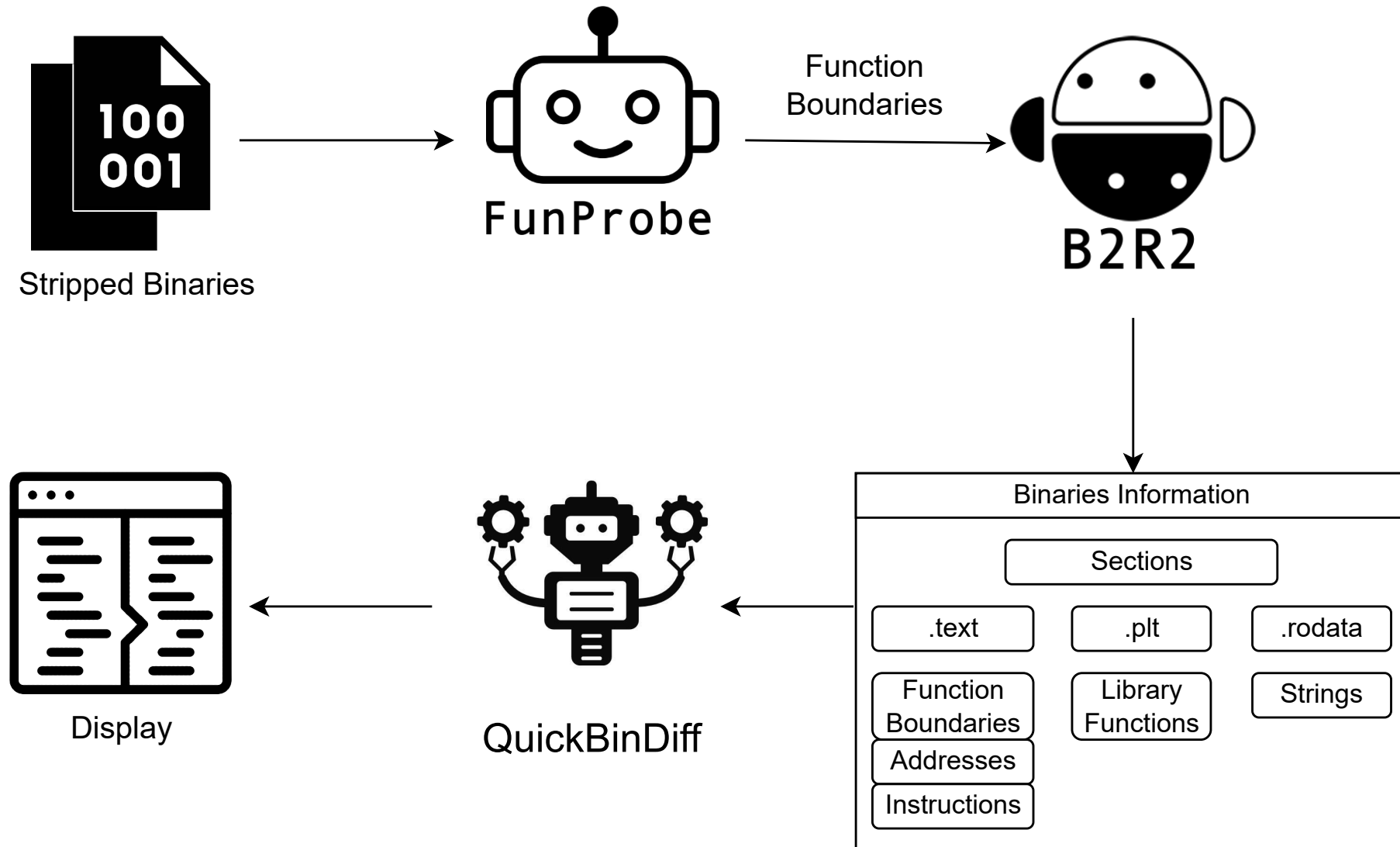
Expected Result

- Fast, free, open source and CLI-based binary code diff tool
- Efficiently perform assembly code diff of two binaries even in limited environments such as servers
- Analyze two binarie codes in detail
- Automatically extract features from a large dataset

Example of Expected Result

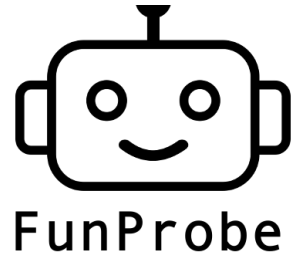
<main>			<main>	
		..skip..		
0x1176	mov EDI, 0x5		mov EDI, 0x5	0x117a
0x117b	call <func1> ; -0x43		call <func1> ; -0x47	0x117f
0x1180	mov dword ptr [RBP-0x4], EAX		mov dword ptr [RBP-0x4], EAX	0x1184
			[+] mov dword ptr [RBP-0x8], 0x0	0x1187
0x1183	jmp +0xa ; <main+0x2a>		jmp +0xa ; <main+0x2a>	0x118e
0x1185	add dword ptr [RBP-0x8], 0x1	[-]		
0x1189	sub dword ptr [RBP-0x4], 0x1		sub dword ptr [RBP-0x4], 0x1	0x1190
			[+] add dword ptr [RBP-0x8], 0x1	0x1194
0x118d	cmp dword ptr [RBP-0x4], 0x0		cmp dword ptr [RBP-0x4], 0x0	0x1198
0x1191	jnz -0xc ; <main+0x22>	[-]	[+] jz -0xc ; <main+0x29>	0x119c
0x1193	mov EAX, dword ptr [RBP-0x8]		mov EAX, dword ptr [RBP-0x8]	0x119e
0x1196	leave		leave	0x11a1
0x1197	ret		ret	0x11a2
0x1198	nop dword ptr [RAX+RAX+0x0]	[-]	[+] nop word ptr [CS:RAX+RAX+0x0]	0x11a3
			[+] nop dword ptr [RAX]	0x11ad

Overview



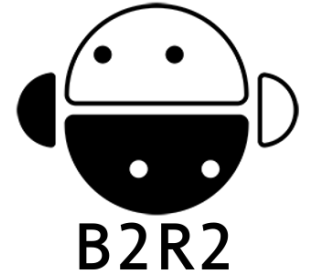
Discover Function Boundary

- A tool for recovering user-defined functions from stripped binaries
- Compiler-agnostic (including GCC, LLVM)
- Architecture-agnostic (including x86, x64, ARM, MIPS)
- A state-of-the-art tool proposed in 2023
 - higher accuracy than existing tools (97% for MIPS-clang, 99% for others)
 - 6 times faster on average (6 seconds for GNU coreutils binary)



Disassemble and Gather Data

- Binary Analysis tool
- Free, open-source, CLI-based
- Cross platform (Linux, Windows, OS X)
- Various file format (including ELF, PE, Mach-O)
- Architecture-agnostic (including x86, x64, ARM, MIPS)
- Providing various information necessary for diff



Binary Diffing

- Binary, section, function, basic block or snippet sliced by address
- Gathering instructions and breaking them down into finer granularity such as opcode and operands
- There are several issues that need to be resolved

Binary Diffing – Resolve Expected Issues

- Separating addresses and instructions
 - The addresses should be included in the diff result, but including them in the diff introduces many errors.

Address	Disassembly	Comment	Address	Disassembly	Comment
0x1176	mov EDI, 0x5		0x117a	mov EDI, 0x5	
0x117b	call <func1> ; -0x43		0x117f	call <func1> ; -0x47	
0x1180	mov dword ptr [RBP-0x4], EAX		0x1184	mov dword ptr [RBP-0x4], EAX	
			0x1187	[+] mov dword ptr [RBP-0x8], 0x0	
0x1183	jmp +0xa ; <main+0x2a>		0x118e	jmp +0xa ; <main+0x2a>	
0x1185	add dword ptr [RBP-0x8], 0x1	[-]			
0x1189	sub dword ptr [RBP-0x4], 0x1		0x1190	sub dword ptr [RBP-0x4], 0x1	
			0x1194	[+] add dword ptr [RBP-0x8], 0x1	
	..skip..			..skip..	

Binary Diffing – Resolve Expected Issues

- Resolve function address
 - Function addresses can also vary due to “address shift”
 - They may point to semantically equivalent functions but with different syntax

		..skip..	
0x1176	mov EDI, 0x5		mov EDI, 0x5 0x117a
0x117b	call <func1> ; -0x43		call <func1> ; -0x47 0x117f
0x1180	mov dword ptr [RBP-0x4], EAX		mov dword ptr [RBP-0x4], EAX 0x1184
			[+] mov dword ptr [RBP-0x8], 0x0 0x1187
0x1183	jmp +0xa ; <main+0x2a>		jmp +0xa ; <main+0x2a> 0x118e
		..skip..	

Binary Diffing – Resolve Expected Issues

- Resolve address with function and offset
 - Even if the target addresses differ, consider them semantically equivalent if the offset of the operand, such as +0xa, is the same relative to the function address, like main+0x2a."

		..skip..	
0x1176	mov EDI, 0x5		mov EDI, 0x5 0x117a
0x117b	call <func1> ; -0x43		call <func1> ; -0x47 0x117f
0x1180	mov dword ptr [RBP-0x4], EAX		mov dword ptr [RBP-0x4], EAX 0x1184
			[+] mov dword ptr [RBP-0x8], 0x0 0x1187
0x1183	jmp +0xa ; <main+0x2a>		jmp +0xa ; <main+0x2a> 0x118e
		..skip..	

Binary Diffing – Resolve Expected Issues

- Resolve address with function and offset
 - Both appear as "jnz -0xc," indicating identical syntax

		..skip..			
0x118d	cmp dword ptr [RBP-0x4], 0x0			cmp dword ptr [RBP-0x4], 0x0	0x1198
0x1191	jnz -0xc ; <main+0x22>	[-]		[+] jz -0xc ; <main+0x29>	0x119c
0x1193	mov EAX, dword ptr [RBP-0x8]			mov EAX, dword ptr [RBP-0x8]	0x119e
0x1196	leave			leave	0x11a1
0x1197	ret			ret	0x11a2

Binary Diffing – Resolve Additional Issues

- Resolve additional issues
- Instruction pointing string in .rodata section
 - same address points to different strings
 - different addresses point to the same string

Display Result

- Section, function, basic block or code snippet sliced by address
- Side-by-side view
- Red, green, background-highlight
 - Fine-grained diff in are highlighted with background color replaced lines
- Explicit mark ([+], [-])

..skip..			
0x118d		cmp dword ptr [RBP-0x4], 0x0	0x1198
0x1191		jnz -0xc ; <main+0x22>	0x119c
0x1193		mov EAX, dword ptr [RBP-0x8]	0x119e
0x1196		leave	0x11a1
0x1197		ret	0x11a2

Prototype (Demo)

- Implement diff algorithm and separate address and instruction

	(.text)		(.text)
	<_start>		<_start>
	00000000000001040: endbr64		00000000000001040: endbr64
	00000000000001044: xor EBP, EBP		00000000000001044: xor EBP, EBP
	00000000000001046: mov R9, RDX		00000000000001046: mov R9, RDX
	00000000000001049: pop RSI		00000000000001049: pop RSI
	0000000000000104a: mov RDX, RSP		0000000000000104a: mov RDX, RSP
	0000000000000104d: and RSP, 0xfffffffffffffffff0		0000000000000104d: and RSP, 0xfffffffffffffffff0
	00000000000001051: push RAX		00000000000001051: push RAX
	00000000000001052: push RSP		00000000000001052: push RSP
[-]	00000000000001053: lea R8, qword ptr [RIP+0x1b6]	[+]	00000000000001053: lea R8, qword ptr [RIP+0x1c6]
[-]	0000000000000105a: lea RCX, qword ptr [RIP+0x13f]	[+]	0000000000000105a: lea RCX, qword ptr [RIP+0x14f]
[-]	00000000000001061: lea RDI, qword ptr [RIP+0xfb]	[+]	00000000000001061: lea RDI, qword ptr [RIP+0xff]
	00000000000001068: call qword ptr [RIP+0x2f72]		00000000000001068: call qword ptr [RIP+0x2f72]
	0000000000000106e: hlt		0000000000000106e: hlt
	0000000000000106f: nop		0000000000000106f: nop

A Study on Readability

- Implement the well-known four diff algorithms
- Resolve issues to maintain semantics
- Analyze various cases using a sample dataset
- (Empirical) Identify cases where readability differs and conduct a case study
- (Theoretical) Analyze the phenomenon based on the operational principles of the algorithm.
- Select the best diff algorithm for readability

Researcher's Capabilities

- Proficient in diff algorithms
 - Implement the well-known four diff algorithms
 - I have discovered that the algorithm's name and its implementation in git diff have not been consistent since 2005
 - The names and implementations of Myers and Minimal have been swapped
 - Microsoft and Git's Developers, contributors, and users are all unaware this truth
- Proficient in the F# language
 - I have a basic understanding of the B2R2
 - If I implement a diff tool, migration can be done easily.
- Large dataset for research
 - 6 architectures, 20 compilers (gcc-4~11, clang3-13), and 6 optimization levels (O0-O3, Os, Ofast)
 - I have compiled 200 unique source codes, resulting in 200K binaries
 - It is easy to construct ground truth for the diff tool and there is enough data to conduct empirical research on readability

Summary

- Binary diff is a technology required in various fields.
 - The absence of a free, open-source, CLI-based diff tools focused on assembly code diff
 - Difficult to automatically handle a large quantity of binaries
 - The lack of research on the readability of binary code diff algorithms
- Fast, free, open sourced and CLI-based binary code diff tool
 - Support cross-platform, various file format and multi architecture