blackhat USA 2024

Uncovering Supply Chain Attack with Code Genome Framework

Dhilung Kirat, Jiyong Jang, Doug Schales, Ted Habeck, Ian Molloy, JR Rao



Dhilung Kirat



Jiyong Jang

AI Supply Chain Security Team IBM **Research**





- \$ foo install bar
- Signed with a certificate.
- Lists dependencies.
- Do you trust it?





Reflections on Trusting Trust

software.

"You can't trust code that you did not totally create yourself."

USA 2024

-Ken Thompson

#BHUSA @BlackHatEvents

KEN THOMPSON

INTRODUCTION

I thank the ACM for this award. I can't help but feel that I am receiving this honor for timing and serendipity as much as technical merit. UNIX¹ swept into popularity with an industry-wide change from central mainframes to autonomous minis. I suspect that Daniel Bobrow [1] would be here instead of me if he could not afford a PDP-10 and had had to "settle" for a PDP-11. Moreover, the current state of UNIX is the result of the labors of a large number of people.

There is an old adage. "Dance with the one that brought you," which means that I should talk about UNIX. I have not worked on mainstream UNIX in many years, yet I continue to get undeserved credit for the work of others. Therefore, I am not going to talk about UNIX, but I want to thank everyone who has contributed.

That brings me to Dennis Ritchie. Our collaboration has been a thing of beauty. In the ten years that we have worked together. I can recall only one case of miscoordination of work. On that occasion, I discovered that we both had written the same 20-line assembly language program. I compared the sources and was astounded to find that they matched character-for-character. The result of our work together has been far greater than the work that we each contributed.

I am a programmer. On my 1040 form, that is what I put down as my occupation. As a programmer, I write

¹ UNIX is a trademark of AT&T Bell Laboratories.

© 1984 0001-0782/84/0800-0761 75¢

To what extent should one trust a statement that a program is free of Trojan horses? Perhaps it is more important to trust the people who wrote the

> programs. I would like to present to you the cutest program I ever wrote. I will do this in three stages and try to bring it together at the end.

STAGE I

In college, before video games, we would amuse ourselves by posing programming exercises. One of the favorites was to write the shortest self-reproducing program. Since this is an exercise divorced from reality, the usual vehicle was FORTRAN. Actually, FORTRAN was the language of choice for the same reason that three-legged races are popular.

More precisely stated, the problem is to write a source program that, when compiled and executed, will produce as output an exact copy of its source. If you have never done this, I urge you to try it on your own. The discovery of how to do it is a revelation that far surpasses any benefit obtained by being told how to do it. The part about "shortest" was just an incentive to demonstrate skill and determine a winner.

Figure 1 shows a self-reproducing program in the C³ programming language. (The purist will note that the program is not precisely a self-reproducing program, but will produce a self-reproducing program.) This entry is much too large to win a prize, but it demonstrates the technique and has two important properties that I need to complete my story: 1) This program can be easily written by another program. 2) This program can contain an arbitrary amount of excess baggage that will be reproduced along with the main algorithm. In the example, even the comment is reproduced.

Supply Chain Attacks

SolarWinds (2019-2021) est. cost > \$100B

Malicious code (backdoor) pushed out through updates

Dependency confusion (Feb 2021)

Private vs public packages (npm, PyPi, RubyGems)

Codecov (Apr 2021)

DevOps tool. Vulnerability in CI. Bash uploader modified

Kaseya (Jul 2021) ransom \$70M

IT solutions, including VSA (remote monitoring and • management software) to deliver REvil ransomware

Protestware (Mar 2022)

Popular NPM package wiped files in Russia and Belarus •

3CX (Mar 2023)

Backdoor implanted into Windows and macOS due to • secondary supply chain attack







Updated: Reports suggest the initial hack may have led to a more extensive supply chain attack.

CISA-FBI Guidance for MSPs and their Customers Affected by the Kaseya VSA Supply-Chain Ransomware Attack



xz Backdoor

https://www.openwall.com/lists/oss-security/2024/03/29/4

Date: Fri, 29 Mar 2024 08:51:26 -0700 From: Andres Freund <andres@...razel.de> To: oss-security@...ts.openwall.com Subject: backdoor in upstream xz/liblzma leading to ssh server compromise

Hi,

After observing a few odd symptoms around liblzma (part of the xz package) on Debian sid installations over the last weeks (logins with ssh taking a lot of CPU, valgrind errors) I figured out the answer:

The upstream xz repository and the xz tarballs have been backdoored.

At first I thought this was a compromise of debian's package, but it turns out to be upstream.

== Compromised Release Tarball ==

One portion of the backdoor is *solely in the distributed tarballs*. For easier reference, here's a link to debian's import of the tarball, but it is also present in the tarballs for 5.6.0 and 5.6.1:

https://salsa.debian.org/debian/xz-utils/-/blob/debian/unstable/m4/build-tohost.m4?ref type=heads#L63

That line is *not* in the upstream source of build-to-host. nor is build-to-host used by xz in git. However, it is present in the tarballs released upstream, except for the "source code" links, which I think github generates directly from the repository contents:

https://github.com/tukaani-project/xz/releases/tag/v5.6.0 https://github.com/tukaani-project/xz/releases/tag/v5.6.1

This injects an obfuscated script to be executed at the end of configure. This script is fairly obfuscated and data from "test" .xz files in the repository.





Semantic gap between compiled code behavior and its metadata

https://securelist.com/xz-backdoor-story-part-1/112354/

Supply Chain Security: Industry approach to protecting CI/CD pipelines





Supply Chain Security: Open security issues and residual risks





Where else the code or

Code Genome





The Semantic Gap



Build chain of trust by following code equivalency





Code Genome Pipeline





#BHUSA @BlackHatEvents



Genome

Code Genome Pipeline









 \downarrow

Trojan.Ramnit



Code Genome Pipeline



Gene can be constructed from closed-source/legacy code



Code Gene

where source code is not easily available.

Code Genome: Semantically meaningful fingerprint





```
d ptr [rbp - 8], edi
           000d8d < f5 + 0x7d >
          dword ptr [rbp - 8]
          d ptr [rbp - 0xc], eax
           000d70 < f5+0x60>
          d ptr [rbp - 0x14]. eax
          dword ptr [rbp - 8]
          dword ptr [rbp - 0x14]
           ptr [rbp - 0x14], eax
          d ptr [rbp - 8]. eax
          dword ptr [rbp - 8]
           dword ptr [rbp - 0xc]
           ptr [rbp - 8], eax
          dword ptr [rbp - 8]
           | ptr [rbp – 0x10], eax
           000def < f5+0xdf>
          d ptr [rbp - 8], 0xa
            00dc3 < f5+0xb3>
 mov dword ptr [rbp - 8]. eax
 mov eax, dword ptr [rbp - 8]
sub eax, dword ptr [rbp - 0x18
 mov dword ptr [rbp - 8], eax
mov eax, dword ptr [rbp - 8]
add eax, 0x20
mov dword ptr [rbp - 8], eax
 mov eax. dword ptr [rbp - 8]
 mov dword ptr [rbp - 4], eax
 jmp 0x100000def <_f5+0xdf>
sub eax, dword ptr [rbp - 8]
mov dword ptr [rbp - 0x1c], eax
mov eax. dword ptr [rbp - 0x1c]
add eax, dword ptr [rbp - 8]
 mov dword ptr [rbp - 0x1c], eax
 mov dword ptr [rbp - 0x1c], eax
mov eax, dword ptr [rbp - 8]
add eax, dword ptr [rbp - 0x1c]
 mov dword ptr [rbp - 8], eax
```

Advantages and Challenges

Advantages

- Across multiple architectures (x86, ARM, ...)
- Across multiple compilers (gcc, clang, ...)
- Across multiple optimization levels
- Handling obfuscation

Challenges

- Disassembly is undecidable
- Function boundary identification
- Loss of architecture specific nuances
- Canonicalization cannot completely recover high-level abstraction



objdump			Ghidra	a		retdec
Disassembly of section eolnw	oiw:			undefined er	ntry()	; section: eolnwoiw
			undefined	AL:1	<return></return>	; function: entry_point at ex/diese - ex/
307d1000 <eolnwoiw>:</eolnwoiw>				entry		ex/d1020: 56
7d1000: 56	push		00741000 56	DUCU	ECT	8x7d1881: 50
7d1001: 50	push		00701000 50	PUSH	FAY	8x7d1882: 53
7d1002: 53	push		007d1002 53	PUSH	EBX	8x7d1083: e8 81 88 88 88
7d1003: e8 01 00 00 00	call	0x7d1009	007d1003 e8 01	00 CALL	LAB 007d1008+1	- data inside code section at 8x7d1888
7d1008: 00 58 89	add	%bl0x77(%eax)	00 00		-	5v7d1939- 09
7d100b: c3	ret .			LAB_007d1008	8+1	Contine Survive 741050 -+ 4-741050
701000: 40	inc	Seav	007d1008 00 58	89 ADD	byte ptr [EAX + -0x77],BL	; function: function_/diede ac ex/diedes
741004 24 00 80 16 00	sub	\$0v168000 %pav	007d100b c3	RET		ex7d1089: 58
7d1012, 2d oc b0 40 00	sub	\$0x100000,000x	007d100c 40	??	40h @	8x7d188a: 89 c3
741017. 05 -7 -0 05 10	SUU	\$0,100000dC,300X	007d100d 2d	??	2Dh -	8x7d188c: 48
74101 00 25	auu	\$0x100000d3,40dx	007d100e 00	22	Bob	ex7d108d: 2d 00 88 16 00
/d1010: 80 3D CC	стро	\$0xcc,(Nebx)	00701001 80	22	166	av7d1812: 2d ac b8 8b 10
/d101†: /5 19	jne	0x/d103a	007d1011 00	22	00h	av7410172 05 x2 b9 4b 10
7d1021: c6 03 00	movb	\$0x0,(%ebx)	007d1012 2d	??	2Dh -	
7d1024: bb 00 10 00 00	mov	\$0x1000,%ebx	007d1013 ac	??	ACh	ex/d181c: 88 30 cc
7d1029: 68 b5 f3 44 34	push	\$0x3444f3b5	007d1014 b0	??	B0h	8x7d181f: 75 19
7d102e: 68 f3 6d 5f 1c	push	\$0x1c5f6df3	007d1015 0b	??	0Bh	@x7d1021: c6 83 80
7d1033: 53	push		007d1016 10	??	10h	8x7d1824: bb 88 18 88 88
7d1034: 50	push		007d1017 05	??	05h	ex7d1879: 58 b5 f3 44 34
7d1035: e8 0a 00 00 00	call	0x7d1044	007d1018 a3	??	A3h	9v7d197e: 68 F3 6d 5f 1c
7d103a: 83 c0 00	add	\$0x0,%eax	007d1019 00	11	BUN	6,7J1032, 53
7d103d: 89 44 24 08	mov	%eax.0x8(%esp)	0070101a 00	22	Joh	ex/01833: 53
7d1041: 5b	DOD	%ebx	00701010 10	22	80h	0x7d1834: 50
7d1842+ 58	000	Seav	007d101d 3b	22	38h :	8x7d1835: e8 8a 88 88 88
7d1043: c3	ret		007d101e cc	??	CCh	0x7d103a: 83 c0 80
741044. 55	nuch	9 obn	007d101f 75	??	75h u	8x7d183d: 89 44 24 88
741045. 00 -5	pusii	Rear Rehr	007d1020 19	??	19h	8v7d1841 - 5h
701045: 09 85	nov	esh, eenh	007d1021 c6	??	C6h	av7d1042 50
/0104/: 50	pusn	weax	007d1022 03	??	03h	0.701042. 38
/d1048: 53	push		007d1023 00	22	00h	ex/d1843: c3





Uncovering Supply Chain Attack





Demo 1: xz backdoor analysis using Code Genome



liblzma.so.5.6.1.github

black hat

USA 2024



Demo 1: xz backdoor analysis using Code Genome

• • •				local	host	\$			ů + G			⊜∙ »વ
≣ Code	Genome								8	Name	Date Modified	Size Kind
Č.	C	`omporo								liblzma.so.5.6.1.distro	Yesterday, 10:36 AM	1.3 MB Docur
Ēq	C	ompare								liblzma.so.5.2.9.github	Yesterday, 10:36 AM	1 MB Docu
Đe					-					libizma.so.5.2.9.distro	Testerday, 10-36 AM	TMB Docu
		Drag and drop files here or c	click to upload	Reset		Drag and drop files here or	r click to upload		00:00:26			
	Y			Swap	ţ							
	62500a4dbabQfaa	~~902~b2~4b5b0dd5~0222df472d7269	h0hc52c22f0c262c		0 4bccc5029924c2c	2010520-064540f240f070-251180b061	1278000108402002					
	055008400005188					220032639043401249197683311868601	uz/8906168463693					
Gen	e similarity: 70)	Identical: 191	Similar: 95	Mismatch: 27	Deletions:	Additions: 101					
					-							
File	e Name:	liblzma.so.5.6.1.github			File Name:	liblzma.so.5.6.1.distro						
File	e Hash	663500a4dbcb9faaac802eb3c4b5b0c	dd5c0333df473d7368b0bc53a32f0	c363a	File Hash	4bccc50a99a4cae2e1053ea964540)f249f97ca35118ebe61d	2789be1c84c3e9				
File	е Туре				File Type							
Las	st updated:	2024-08-05T03:51:35.000Z			Last Updated:	2024-08-05T03:51:44.000Z						
File	e size:	1240992			File size:	1306808						
Ger	ne count:	306			Gene count:	407						
Q									×			
liblzma.s	so.5.6.1.github (66	3500a4) Functions	\uparrow	liblzma.so.5.6.1.distro (4bccc50	a) Functions		Score	ор	Actions			
crc32_re									:			
crc64_re	esolve			crc64_resolve				!	:			
get_liter									÷			
get_opti									:			
hash_ap									:			
lz_encod									:			
lzma2_b									:			
lzma_de									1			
Izma_inc									:			
Izma_inc									1			
lzma_lz_									:	🗮 Macintosh HD > 🔯 Users > 🛅 dkira	t > 🛅 Downloads > 🚞 xz 4 items. 332.3 GB available	
											4 reme, 052.5 Ob available	



xz backdoor Gene Similarity Analysis using GeneDiff



Local vs distribution builds of same version



xz backdoor Gene Similarity Analysis using GeneDiff



Incremental version similarity in distribution builds



Improving Supply Chain Security





Trust but Verify SBOM: Metadata vs. Code

Problem

- Each vendor creates SBOM of their own software including open-source and closed-source components.
- How can we verify its *correctness* (containing incorrect library mistakenly/maliciously) and completeness (missing library)?

sbom generation tools \$



incorrectly report the version of OpenSSL that's used by the Node.js runtime."

https://www.chainguard.dev/unchained/mitigating-critical-openssl-vulnerability-with-chainguard



Knowledge Graph: Gene Granularity





#BHUSA @BlackHatEvents

23

Knowledge Graph: Gene Granularity



black hat USA 2024

#BHUSA @BlackHatEvents

24

Demo 2: SBOM generation for an unknown rpm package

Custom **rpm** package



SBOM generated by Code Genome

Component	\$	Version 🔶	License
arping	#	20160308	BSD and GPLv2+
clockdiff	#	20160308	BSD and GPLv2+
ifenslave	#	20160308	BSD and GPLv2+
iputils	#	20160308	BSD and GPLv2+
ping	#	20160308	BSD and GPLv2+
rdisc	#	20160308	BSD and GPLv2+
tracepath	#	20160308	BSD and GPLv2+
tracepath6	#	20160308	BSD and GPLv2+

Integrating with other SBOM analysis platforms



Knowledge Graph: Code Genome and Use Cases





#BHUSA @BlackHatEvents

Code Genome KG

Open Sourcing Code Genome





Status and Roadmap

Open-source tools

- Code Genome Framework

- GeneDiff, Basic KG, CLI tools, and GUI
- Currently supported
 - Binaries: ELF, PE, Mach-0
 - Architectures: x86, x86_64, arm, aarch64, mips, ppc
- Optimized canonicalization
- Jaudit
 - JAR file support
 - JAR version identification
 - CVE annotation

Next steps

- Support
 - Packages: deb, rpm, ipa
 - Archives: ar, cpio, tar, bzip2, gzip, zstd, xz, rar, 7zip



git clone https://github.com/code-genome/codegenome.git cd codegenome

make start

					loci
с	ode Genome				
	С	ompare			
					Swan
					onap
	663500a4dbcb9faaa	c802eb3c4b5b0dd5c0333df473d7368b0bc53a32f0c36	3a		
G	ene similarity: 70	Identical: 1	191	Similar: 95	
a	1				
liblz	ma.so.5.6.1.github (66	3500a4) Functions	Ŷ	liblzma.so.5.6.1.dis	tro (4bccc5
crc3					
crce					
get_					
get_					
hasi					
lz_e					
Izma					
Izma					
lzm:					



Public

https://github.com/code-genome

		٩				
						۲
	4bccc50a99a4cae2e1	1053ea964540f249f97ca35118ebe61c	12789be1c84c3e9	93		
Mis	smatch: 27	Deletions:	Additions	:: 101		
	File Name: File Hash File Type	liblzma.so.5.6.1.distro 4bccc50a99a4cae2e1053ea964540		be61d2789be1c84c3e93		
						×
Fund	ctions		Score	op	Actions	



Semantic Gap



Code



Inherent sematic gap breaks the transfer of trust from metadata to code

Code Genome





Now open-sourced Code Genome Framework can help bridge that gap Detection of XZ-backdoor demonstrates framework's capability in improving supply chain security



#BHUSA @BlackHatEvents

Supply Chain Security



29

blackhat® USA 2024

Dhilung Kirat ⊠ dkirat@us.ibm.com **Jiyong Jang** ⊠ jjang@us.ibm.com



IBM Research



github.com/code-genome