

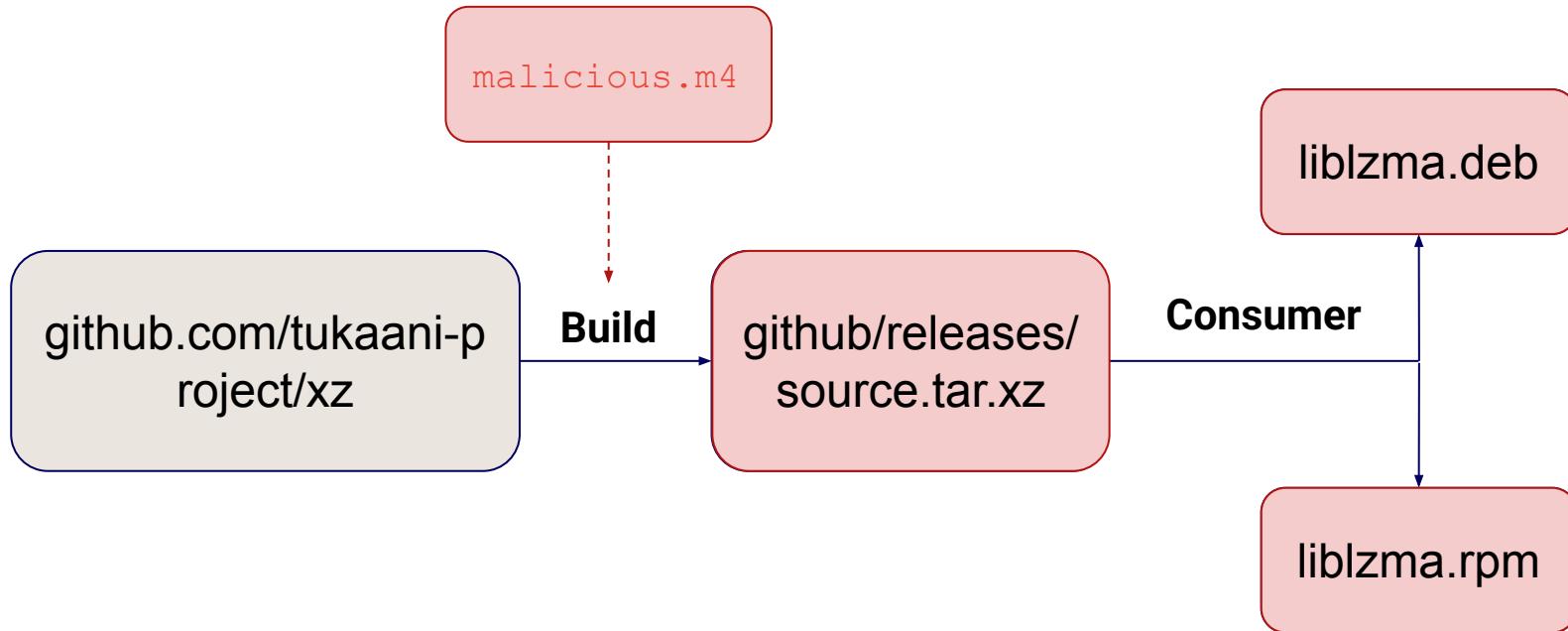


# Causes and Mitigations of Unreproducible Builds in Java

[Canonicalization for  
Unreproducible Builds in Java](#)

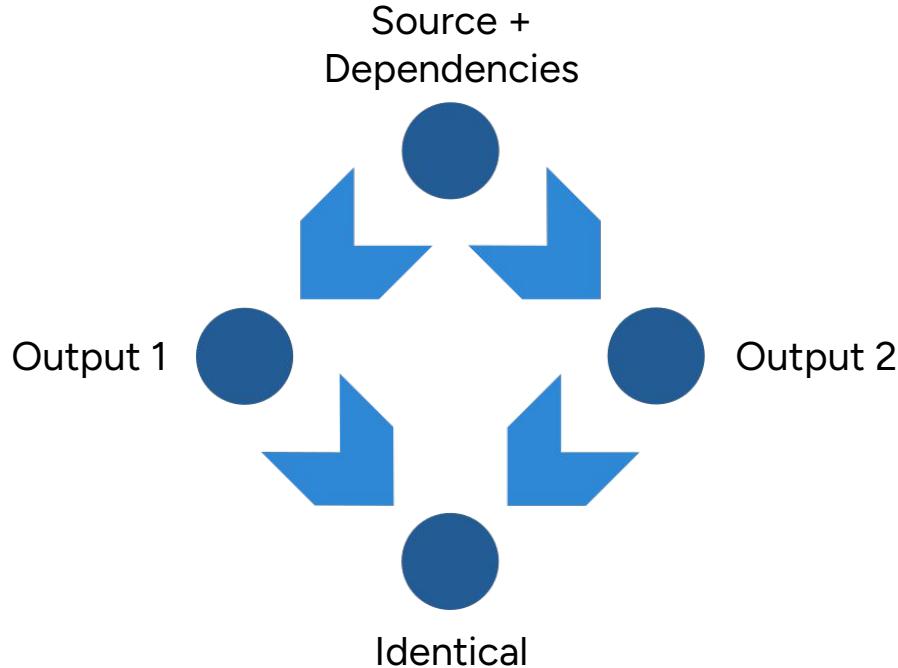
Aman Sharma, Benoit Baudry, Martin Monperrus

# xz-utils attack



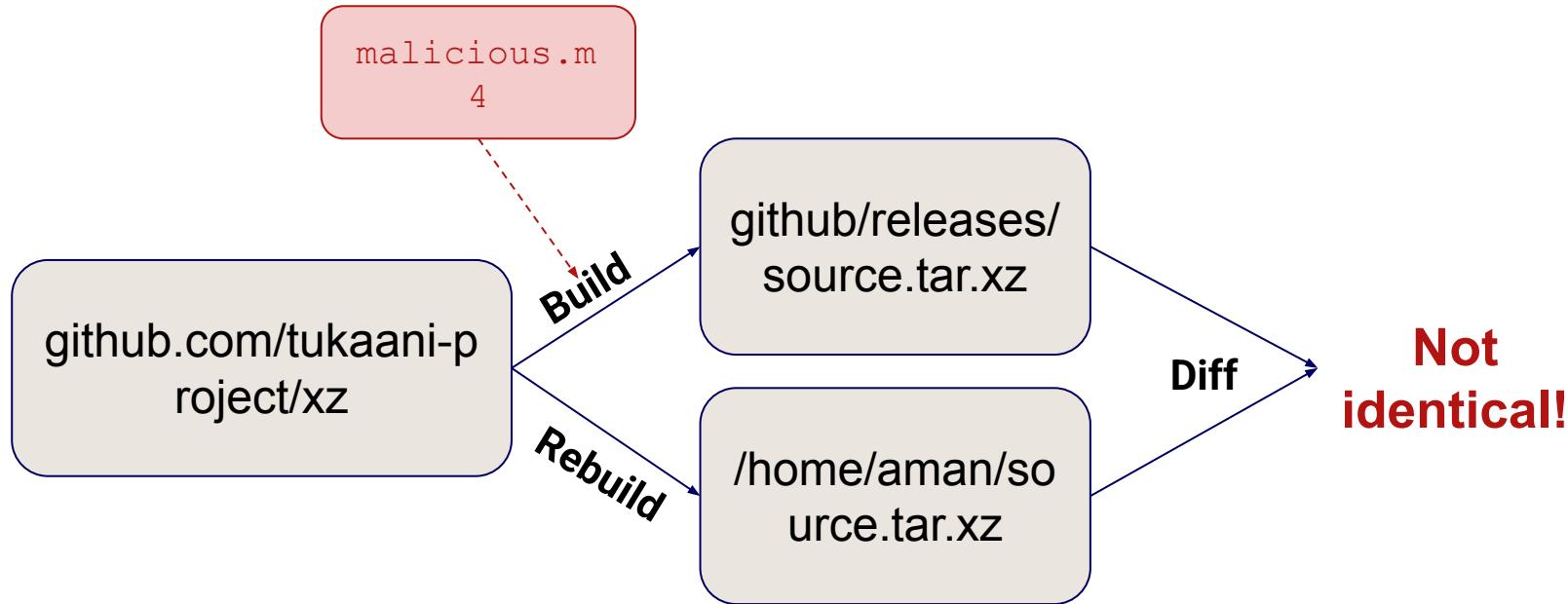
[1] J. Malka, 'How NixOS and reproducible builds could have detected the xz backdoor for the benefit of all'. Available: <https://luj.fr/blog/how-nixos-could-have-detected-xz.html>

# Reproducible Builds



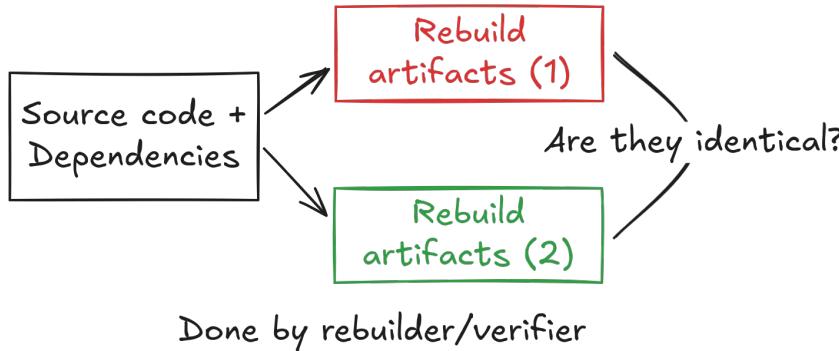
[2] Chris Lamb and Stefano Zacchiroli. 2022. Reproducible Builds: Increasing the Integrity of Software Supply Chains. *IEEE Software* 39, 2 (March 2022), 62–70. <https://doi.org/10.1109/MS.2021.3073045>

# xz-utils mitigation



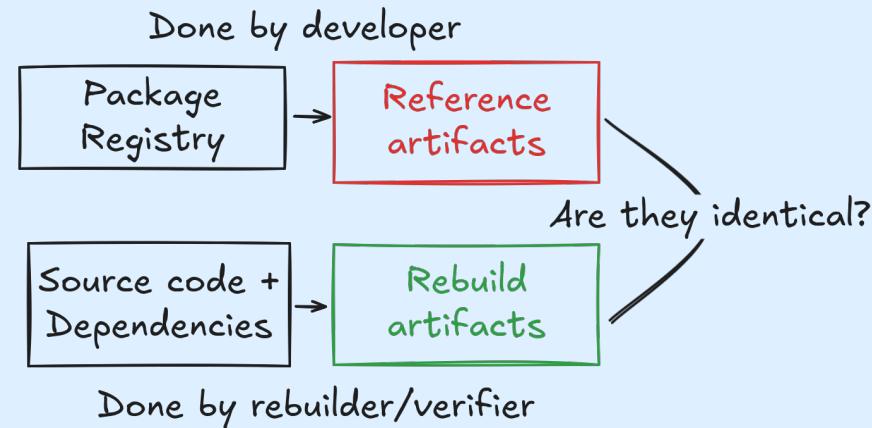
# How to check for reproducible builds?

## Build twice and compare



[3] G. Benedetti et al., 'An Empirical Study on Reproducible Packaging in Open-Source Ecosystems', 57th International Conference on Software Engineering, 2025.

## Build and compare with package registry



[1] J. Malka, S. Zacchiroli, and T. Zimmermann, 'Does Functional Package Management Enable Reproducible Builds at Scale? Yes', in 22nd International Conference on Mining Software Repositories, Ottawa

# Related Work: build twice and compare

[4] G. Benedetti *et al.*, 'An Empirical Study on Reproducible Packaging in Open-Source Ecosystems', presented at the Proceedings of the 47th International Conference on Software Engineering, 2025.

Summary: Checking if builds are reproducible for RubyGems, PyPI, Maven.

[5] Z. Ren, H. Jiang, J. Xuan, and Z. Yang, 'Automated localization for unreproducible builds', in *Proceedings of the 40th International Conference on Software Engineering*, 2018

Summary: Localizing files that differ in two subsequent builds for Debian packages.

[6] O. S. Navarro Leija *et al.*, 'Reproducible Containers', in *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2020

Summary: Wrapper around the build process for debian packages.

# Related Work: build and compare with package registry

[7] J. Malka, S. Zacchiroli, and T. Zimmermann, 'Does Functional Package Management Enable Reproducible Builds at Scale? Yes', in 22nd International Conference on Mining Software Repositories, 2025.

Summary

[8] R. Bajaj, et al., 'Reproducible Builds for Go', in Proceedings of the 2023 ACM SIGPLAN Conference on Programming Language Design and Implementation, 2023.

Summary  
packages

We like this approach more  
as it maximizes environment  
differences!

[9] V. Andersson, Geth Rebuild : Verifiable Builds for Go Ethereum. 2024

Summary: Rebuilds ethereum client that enables connection to the main ethereum network.



# Reproducible/Verifiable/Accountable builds

Verifiable Builds: Builds are verifiable if the build process of a software can be fixed to output bit-by-bit identical artifacts.

Accounta  
unequal a

Basically,  
Builds in c

outputs  
reproducible

We always prefer  
“Reproducible Builds” over  
other terms.

[10] T. Pohl, P. Novák, M. Ohm, and M. Meier, ‘SoK: Towards Reproducibility for Software Packages in Scripting Language Ecosystems’, Mar. 27, 2025, arXiv: arXiv:2503.21705. doi: 10.48550/arXiv.2503.21705.

# Why is it important?

1. Ensuring Integrity: Reproducible builds ensure that the executable corresponds to the source code (assuming source code can be audited) and hence is not tampered with. [11]
2. Faster builds: Dependent sources do not need to be rebuilt and dependent tasks do not need to be rerun if a rebuild of a package does not yield different results. [12]

[11] Mike Perry. 2013. Deterministic Builds Part One: Cyberwar and Global Compromise | Tor Project.  
<https://blog.torproject.org/deterministic-builds-part-one-cyberwar-and-global-compromise/>

[12] <https://reproducible-builds.org/>

# Are builds reproducible yet?

Arch Linux is **86.4%** reproducible with **2026 bad 0 unknown** and **12873 good** packages.

[core] repository is **95.3%** reproducible with **13 bad 0 unknown** and **263 good** packages.

[extra] repository is **86.3%** reproducible with **1987 bad 0 unknown** and **12489 good** packages

[core-testing] repository is **100.0%** reproducible with **0 bad 0 unknown** and **4 good** packages.

[extra-testing] repository is **81.8%** reproducible with **26 bad 0 unknown** and **117 good** packages.



every 4.1 months from 2017 to 2023. Our findings show that bitwise reproducibility in nixpkgs is very high and has known an upward trend, from **69%** in 2017 to **91%** in 2023. The mere ability to rebuild packages (whether bitwise reproducibly or

the number of non-reproducible versions (v), v = 1303  
**Finding 1:** *1,303 out of 3,390 studied versions (38%) are non-reproducible. With such a large portion of non-reproducible package versions, developers should not blindly trust the verifiability of NPM packages.* [14]

rebuilding **7080 releases** of **869 projects**:

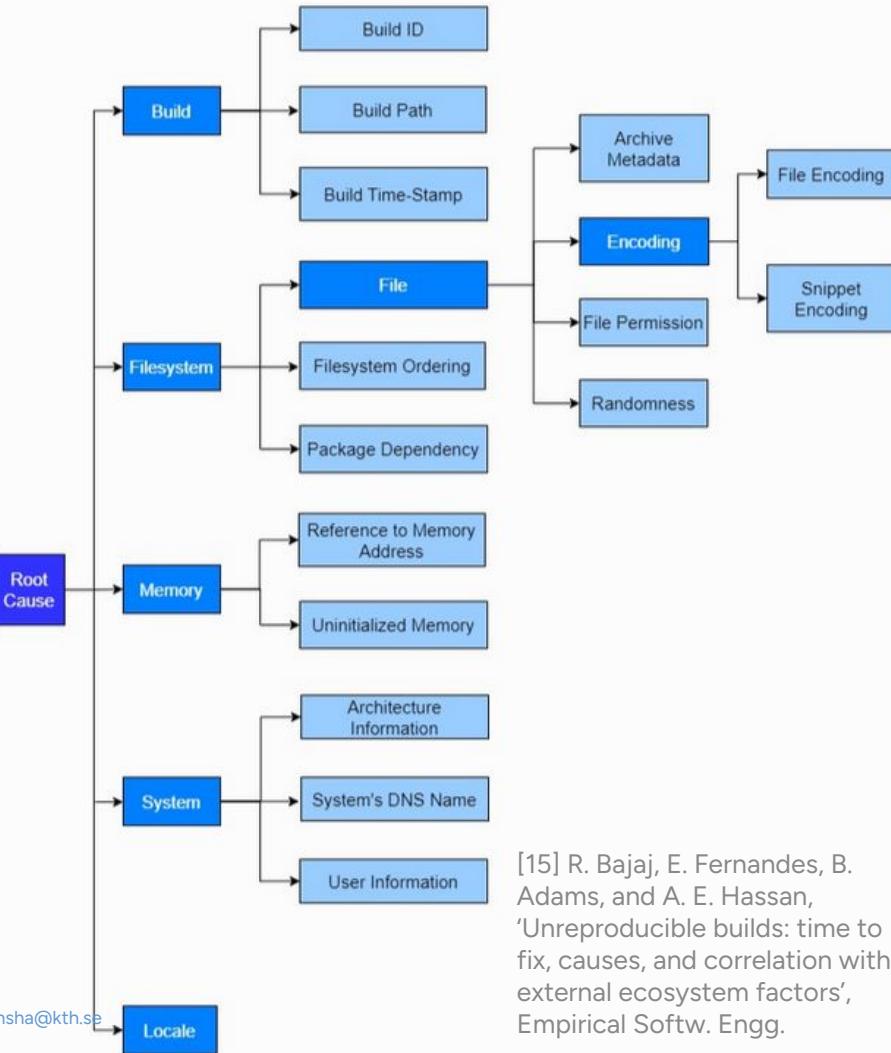
- **5112** releases are confirmed **fully reproducible** (100% reproducible artifacts ✓),
- **1968** releases are only partially reproducible (contain some unreproducible artifacts !)
- on **869** projects, **755** have at least one fully reproducible release, **114** have none

[13] P. Goswami, S. Gupta, Z. Li, N. Meng, and D. Yao, 'Investigating The Reproducibility of NPM Packages', in 2020 IEEE International Conference on Software Maintenance and Evolution (ICSME)

[14] <https://isdebianreproducibleyet.com/>

# Why aren't builds reproducible?

Short answer: Timestamps, signature embedded



[15] R. Bajaj, E. Fernandes, B. Adams, and A. E. Hassan,  
'Unreproducible builds: time to fix, causes, and correlation with external ecosystem factors',  
Empirical Softw. Engng.

# What is the status of reproducible builds in Java?

Most closely related work:

[16] J. Xiong, Y. Shi, B. Chen, F. R. Cogo, and Z. M. (Jack) Jiang, 'Towards build verifiability for Java-based systems', in Proceedings of the 44th International Conference on Software Engineering: Software Engineering in Practice, in ICSE-SEIP '22

- 1) Follows "build twice and compare"
- 2) Small and close sourced dataset

rebuilding **7080 releases of 869 projects**:

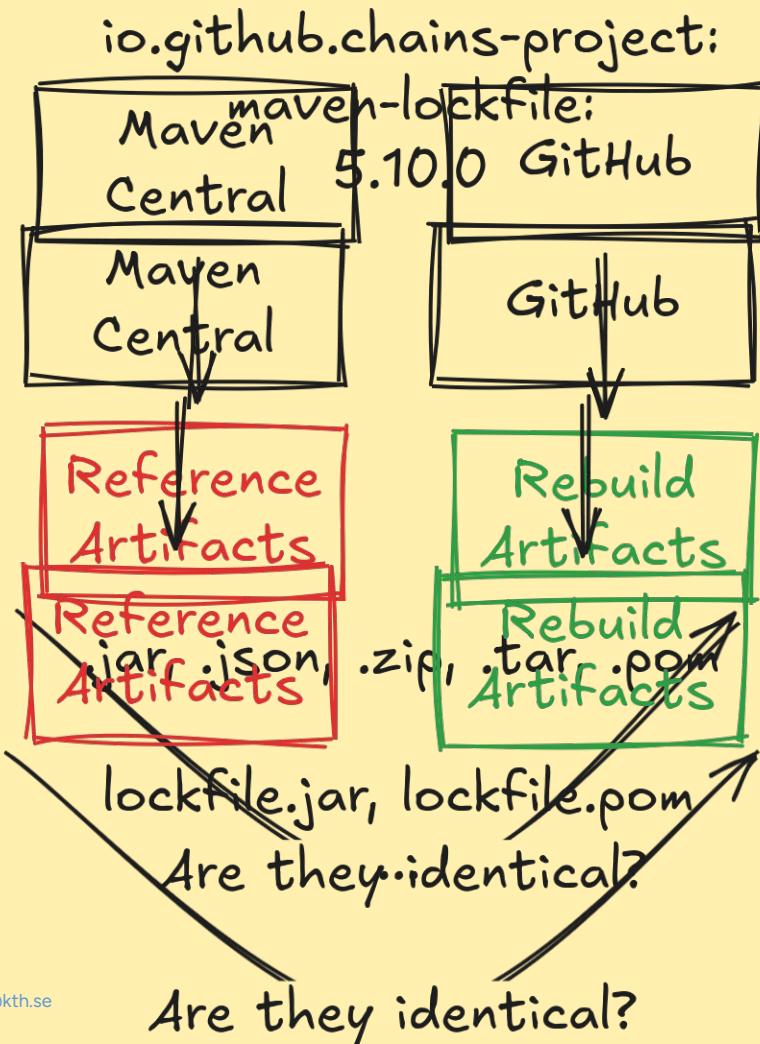
- 5112 releases are confirmed **fully reproducible** (100% reproducible artifacts ✓),
- 1968 releases are only partially reproducible (contain some unreproducible artifacts !)
- on 869 projects, 755 have at least one fully reproducible release, 114 have none

[16] P. Goswami, S. Gupta, Z. Li, N. Meng, and D. Yao, 'Investigating The Reproducibility of NPM Packages', in 2020 IEEE International Conference on Software Maintenance and Evolution (ICSME)

# Reproducible Central

1. It is a GitHub repository:  
<https://github.com/jvm-repo-rebuild/reproducible-central/>
2. Artifacts like jar, .json, .zip files are downloaded from Maven Central.
3. Source code is downloaded from hosting services like GitHub.
4. The source code along with all of its build requirement is built using a "buildspec".
5. Buildspec file is specification that lists how the build should done
  - a. Build command
  - b. Java version
  - c. Environment variables
6. If corresponding reference and rebuild artifacts are identical, project is marked reproducible.

[17] H. Boutemy, jvm-repo-rebuild/reproducible-central. (Oct. 31, 2024).  
 jvm-repo-rebuild.



# Our first goal

To build a taxonomy of reasons why unreproducible builds occur in Java and propose a fix for them



# Dataset

We build our dataset on top of Reproducible Central - it is an infrastructure to verify whether Maven projects are reproducible.

1. Snapshot Reproducible Central on 8th October, 2024 ([d280bf1](#))
  - a. 706 Java projects and their buildspec files
2. We apply three filters that resulted in exclusion of projects
  - a. Projects that require manual intervention (eg. setting precise Java version).
  - b. Projects that have more artifacts built locally than on package registry (eg. some jars are built locally but never released).
  - c. Projects that mapped artifacts incorrectly (eg. pom file is compared with jar).
3. Finally, we have
  - a. 8,292 Maven releases **of submodules**
  - b. 12,803 unreproducible artifact pairs

# RQ1: Causes of Unreproducibility

Reason for Unreproducibility	Root Cause of Unreproducibility	Novelty	Example
Build Manifests	Environment	-	Built-By attribute
	Rebuild Process	-	Build-Jdk attribute
	Non-deterministic configuration	✓	Embedded branch names
SBOM	Java Vendor	✓	Different checksum algorithms
	Custom release configuration	✓	Releasing a subset of artifacts
	Non-deterministic information	✓	Timestamp
Filesystem	Environment	-	Permissions
	Custom release configuration	✓	Generated binaries are not included
JVM Bytecode	JDK Version	-	Ordering of constant pool entries
	Embedded data	✓	Git branch embedded
	Build time generated code	-	Lambda functions
Versioning Properties	-	-	Number of tags embedded in Jar
Timestamps	-	-	Embedded in shell script in Jar files

**Table 1: Summary of the taxonomy of unreproducibility causes**

# Mitigation Strategies

There are 3 mitigation strategies

1. Fix the build process: repairing the build script used by developer to publish package
  - a. Pros: anyone can replicate the build with 0 ad-hoc configuration
  - b. Cons: sometimes you have to embed non-deterministic information like signatures
2. Fix the rebuild process: repairing the buildspec file
  - a. Pros: contributes to fixes in the build process
  - b. Cons: can be hard to adapt for future releases
3. Canonicalization: conversion of output artifacts into a standard representation that is free of non-deterministic changes
  - a. Pros: past releases can also be verified, no need to rebuild, "i can't find a jdk 14 nowadays"
  - b. Cons: careful removal as it may mask meaningful differences

We try to evaluate canonicalization since few work in the literature

[18] J. Dietrich, T. White, B. Hassanshahi, and P. Krishnan, 'Levels of Binary Equivalence for the Comparison of Binaries from Alternative Builds'

# Example: Build Manifest

```
// rg.apache.camel:camel-ahc-ws:3.13.0
// Built-By: name of user that built the
artifact
-Built-By: root
+Built-By: aman
```

Mitigation: Fix build process or canonicalize

```
// org.apache.any23:apache-any23-csvutils:
2.7
-Specification-Title: Apache Any23 :: CSV Utilities
-Specification-Version: 2.7
-Specification-Vendor: The Apache Software Foundation
-Implementation-Title: Apache Any23 :: CSV Utilities
-Implementation-Version: 2.7
-Implementation-Vendor: The Apache Software Foundation
... some more attributes
```

# Build Manifest mitigation

```
// reference  
-Built-By: root  
// rebuild  
+Built-By: aman
```

Diff between original  
reference and rebuild  
artifacts

-Built-By: root  
+Built-By: foo

-Built-By: root

Diff between canonicalized artifact

-Built-By: aman  
+Built-By: foo

-Built-By: aman

or

Now  
Identical!

# Example: Software Bill of Materials

```
// org.apache.commons:commons-parent:61  
  
- <hash alg=\"SHA3-384\">289f952[TRUNCATED]</hash>  
- <hash alg=\"SHA3-256\">32b50a7[TRUNCATED]</hash>  
- <hash alg=\"SHA3-512\">1daf813[TRUNCATED]</hash>
```

Reason: Azul JDK had backported these hash algorithms.

Mitigation: Fix build process by correcting the Java vendor

# Example: Filesystem

```
// io.github.albertus82:unexepack:0.2.1
// change in file permissions and owner
--rw-r--r--  root unexepack-0.1.1/unexepack.jar
+-rw-rw-r--  aman unexepack-0.1.1/unexepack.jar
```

Mitigation: Set permissions and owner to fixed value (canonicalization)

# JVM bytecode

```
-- private static com.google.common.base.Functions$IdentityFunction[] $values();
-- descriptor: ()[Lcom/google/common/base/Functions$IdentityFunction;
-- flags: (0x100a) ACC_PRIVATE, ACC_STATIC, ACC_SYNTHETIC
-- Code:
--   stack=4, locals=0, args_size=0
--       0: iconst_1
--       1: anewarray     #1           // class com/google/common/base/Functions$IdentityFunction
--       4: dup
--       5: iconst_0
--       6: getstatic     #2           // Field INSTANCE:Lcom/google/common/base/Functions$IdentityFunction;
--       9: aastore
--      10: areturn
-- LineNumberTable:
--   line 90: 0
-
```



Mitigation: Hard to find JDK in such cases. It is better to consider two implementations equivalent (canonicalize).



c povirk on Feb 14

Member

...

Hi. At this point, we build our releases by running a script on our local machines. (I expect that to change in the future as part of general security hardening.) Until recently, we even used the default JDK on our machines, which is a version that occasionally has patches to javac. (Nowadays, our default is to use a standard Debian JDK.)

Our JDK 11 has a patch to omit enclosing-class references—much like [JDK-8271717](#) did for later JDKs but I think perhaps even more aggressive in omitting the reference for `Serializable` types (though I don't think that difference is relevant here)? That should explain the `LocalCache` difference.

For enums, I suspect that we also had a patch like [JDK-8241798](#).

But in short, we were using a patched javac.



# Versioning Properties

```
// com.github.ldapchai:ldapchai:0.8.6
// Embedded in MANIFEST.MF
-SCM-Git-Branch: master
+SCM-Git-Branch: 338023a
```

```
// org.apache.drill:drill-opensdb-storage:1.21.0
// Embedded in git.properties
-git.tags=drill-1.21.0
+git.tags=drill-1.21.0,drill-1.21.1,drill-1.21.2
```

Mitigation: Remove such attributes (canonicalization).

# Timestamps

```
-#Wed Apr 20 20:27:41 CEST 2022
+#Fri Oct 18 03:03:44 UTC 2024
```

1. they can be embedded in the properties file
  2. generated documentation
  3. shell scripts appendix
  4. executable binaries
  5. software bill of materials
  6. JVM bytecode
  7. file metadata
  8. MANIFEST.MF
  9. NOTICE
  10. servlets created by Jasper JSP compiler
- Mitigation: Add  
`project.build.outputTimestamp`  
property to pom so that plugins can  
respect it (fixing build process)

# RQ1: Causes of Unreproducibility

Reason for Unreproducibility	Root Cause of Unreproducibility	Novelty	Example	Main Mitigation
Build Manifests	Environment	-	Built-By attribute	Canonicalization by rebuilder
	Rebuild Process	-	Build-Jdk attribute	Fix rebuild process
	Non-deterministic configuration	✓	Embedded branch names	Fix build process
SBOM	Java Vendor	✓	Different checksum algorithms	Fix rebuild process
Filesystem				Canonicalization by rebuilder
JVM Bytecode				Canonicalization by rebuilder
Versioning Properties				Canonicalization by rebuilder
Timestamps	-	-	Embedded in shell script in Jar files	Fix build process

Takeaway: 6 main reasons of unreproducibility in Java

Table 1: Summary of the taxonomy of unreproducibility causes

# Our second goal

Analyze effectiveness of canonicalization in mitigating unreproducible builds

# What is this term “canonicalization”?

Inspired by canonical URL. To choose the preferred, representative URL for a web page, helping search engines understand which version of a page should be indexed [19].



## Bytecode Canonicalization

- Classnames could change across different executions.
- The type references change.
- The order of method is not fixed.

```
- public class $Proxy10 {  
+ public class $Proxy7 {  
-     private static $Proxy10.x;  
+     private static $Proxy7.x;  
-     m1 () {}  
+     m3 () {}  
-     m3 () {}  
+     m1 () {}  
}
```

[20]

In our context, we create a preferred representation of the artifact which abstracts away non-deterministic differences.

1. Removing parts of artifacts
2. Setting fixed values

[20] <https://support.google.com/webmasters/answer/10347851?hl=en>

[21] A. Sharma, M. Wittlinger, B. Baudry, and M. Monperrus, ‘SBOM.EXE: Countering Dynamic Code Injection based on Software Bill of Materials in Java’, Jun. 28, 2024, arXiv: arXiv:2407.00246.

# Artifact Canonicalization

It means that the entire artifact is transformed by removing non-deterministic and spurious changes, especially in metadata.

Eg. ordering of files in archive is made consistent

Tool used: OSS-Rebuild [21]

[22] <https://github.com/google/oss-rebuild>

[23] S. Schott, S. E. Ponta, W. Fischer, J. Klauke, and E. Bodden, 'Java Bytecode Normalization for Code Similarity Analysis'

# Bytecode Canonicalization

It is a process of transforming the bytecode of a program into a representation that is independent of specific implementation details inserted by the compiler.

Eg. implementation of string concatenation pre and post Java 9

Tool used: jNorm [22]

## RQ2: Effectiveness of Artifact Canonicalization

Tool	Successful Canonicalization	Failed Canonicalization
OSS-Rebuild (ref1e01)	1,205 (0.41%)	11,508 (99.50%)

Takeaway: 2.5x increase in canonicalized artifacts with only a few improvements.

CHAINS-Rebuild: Fork of OSS-Rebuild with support for canonicalizing build manifests and embedded versioning properties.

## RQ3: Effectiveness of Bytecode Canonicalization

We select a subset of unreproducible artifacts that have JVM bytecode changes.

This is 936/12,803 artifact pairs.

	Successful canonicalization	Failure in canonicalization	Error in canonicalization
#	275 (29.4%)	508 (54.3%)	153 (16.3%)

Reasons of failure:

1. Structural Changes limitation (eg. changes in order of methods, fields, etc)
2. Control flow limitation (eg. ifne and ifeq)
3. Embedded data (eg. absolute file paths)
4. Optimization limitation (eg. string concatenation)

# Example diff for jNorm

---

```
1 - 99: ifne          106
2 - 102:  iconst_1
3 - 103:  goto         107
4 - 106:  iconst_0
5 - 107:  ireturn
6 + 99: ifeq          104
7 + 102:  iconst_0
8 + 103:  ireturn
9 + 104:  iconst_1
10 + 105: ireturn
```

---

**Listing 1:** javap diff for control flow difference.

---

```
1 -if v != 0 goto label;
2 -v = 1;
3 -goto label;
4 -label:
5 -v = 0;
6 +if v == 0 goto label;
7 +return 0;
8 label:
9 -return v;
10 +return 1;
```

---

**Listing 2:** JNORM diff for control flow differnce.

## RQ3: Effectiveness of Bytecode Canonicalization

We select a subset of unreproducible artifacts that have JVM bytecode changes.

This is 936/12,803 artifact pairs.

#

Reasons of 1

# Takeaway: 4 categories of improvements for bytecode canonicalization.

1. Structure limitation
2. Control flow limitation (eg. irne and ireq)
3. Embedded data (eg. absolute file paths)
4. Optimization limitation (eg. string concatenation)



Paper:

## Causes and Canonicalization of Unreproducible Builds in Java

Accepted at IEEE Transactions on Software Engineering

<https://ieeexplore.ieee.org/document/11223991>

# Causes and Canonicalization of Unreproducible Builds in Java

Aman Sharma  
KTH Royal Institute of Technology  
Stockholm, Sweden  
amansha@kth.se

Benoit Baudry  
Université de Montréal  
Montréal, Canada  
benoit.baudry@umontreal.ca

Martin Monperrus  
KTH Royal Institute of Technology  
Stockholm, Sweden  
monperrus@kth.se

### Abstract

The increasing complexity of software supply chains and the rise of supply chain attacks have elevated concerns around software integrity. Users and stakeholders face significant challenges in validating that a given software artifact corresponds to its declared source. Reproducible Builds address this challenge by ensuring that independently performed builds from identical source code produce identical binaries. However, achieving reproducibility at scale remains difficult, especially in Java, due to a range of non-deterministic factors and caveats in the build process. In this work, we focus on reproducibility in Java-based software, archetypal of enterprise applications. We introduce a conceptual framework for reproducible builds, we analyze a large dataset from Reproducible Central, and we develop a novel taxonomy of six root causes of unreproducibility. We study actionable mitigations: artifact and bytecode canonicalization using OSS-REBUILD and jNORM respectively. Finally, we present CHAINS-REBUILD (improvements to OSS-REBUILD), a tool that raises reproducibility success from 9.48% to 26.60% on 12,803 unreproducible artifacts. To sum up, our contributions are the first large-scale taxonomy of build unreproducibility causes in Java, a publicly available dataset of unreproducible builds, and CHAINS-REBUILD, a canonicalization tool for mitigating unreproducible builds in Java.

### Keywords

Reproducible Builds, Software Supply Chain, Canonicalization, Java

### 1 Introduction

The growing complexity of software supply chains [9, 56], coupled with the increasing frequency of supply chain attacks<sup>1</sup>, raises concerns about software integrity. In such a fragmented ecosystem, relying solely on assumptions about the identity of the distributor [58] is no longer sufficient — what is needed is verifiable evidence that the software one installs corresponds exactly to its declared source. This challenge is especially acute in open source environments, where binaries are often distributed separately from their source code [11], making it difficult for users to independently validate what they are running. As a result, the focus is shifting toward techniques that can guarantee that what is built matches what was intended, regardless of who performs the build.

This technique is formally called “Reproducible Build” [10, 61]: builds are considered “reproducible” if and only if the build process is deterministic, where the same binary can be computed from the same source code by independent parties [25]. It helps prevent attacks on the build process, where an attacker modifies it to insert backdoors or malicious code into the software application to be built [41]. Any backdoor or malicious code would be detected

by the verifier as the built artifact would not match the original artifact. This approach has already seen adoption in security and privacy critical environments, such as Bitcoin and Tor Browser [28]. Ensuring that clients for both of them are free of any backdoors injected by the build system.

Despite the promises of Reproducible Builds, ensuring and verifying reproducibility at scale remains technically challenging due to the multitude of spurious differences in build outputs (e.g., timestamps, file ordering). Those differences make binaries appear different even when built from the same source, with an untampered build pipeline [17]. Spurious differences hinder reproducibility and create a new attack surface. Several high-profile supply chain attacks have exploited unreproducible builds. For instance, in 2024 it was discovered that the official tarballs of XZ Utils, compression library used in almost all Linux distributions, contained a backdoor that enabled remote code execution, even though the backdoor was absent in the project’s Git repository [44]. Similarly, the Ultralytics Python package, a widely used library for AI models, was compromised when attackers injected malicious code into the build pipeline, resulting in a PyPI release with a backdoor [26]. A related incident targeted the Solana ecosystem as well [7]. These incidents share a common vector: the artifacts distributed through registries or release tarballs differed from what could be obtained by building directly from the publicly available source code. If the build processes for these projects had been reproducible, users or third party rebuilders could have verified that the distributed binaries and packages matched the source. Any discrepancies would have immediately signaled tampering, thereby preventing such attacks [32].

In this work, we contribute to solving the problem of achieving build reproducibility in the context of enterprise software in Java. Recall that Java has consistently been one of the top 5 programming languages in the world for the past 5 years<sup>2</sup>, underscoring its widespread use and relevance. Its importance in the context of finance [36, 56], government services [51], and military applications<sup>3,4</sup> highlights the need for reproducible builds in Java. We aim at building the most comprehensive taxonomy of unreproducible builds in Java, and the corresponding mitigation strategies. Furthermore, we perform a deep study how canonicalization — removal of non-deterministic differences — is a good solution for mitigating unreproducibility.

Our approach for analyzing unreproducible builds in Java is shown in Figure 1. We first propose an original framework for build reproducibility where we clearly define the roles of the builder and rebuilder and how both of them contribute to reproducible build

<sup>1</sup><https://innovationgraph.github.com/global-metrics/programming-languages/>

<sup>2</sup>[http://pdf.cloud.opensystemsmedia.com/vita-technologies.com/Aonix\\_Jun07.pdf](http://pdf.cloud.opensystemsmedia.com/vita-technologies.com/Aonix_Jun07.pdf)

<sup>3</sup><https://tsri.com/news-blog/press/u-s-department-of-defense-mainframe-cobol-to-java-on-aws>

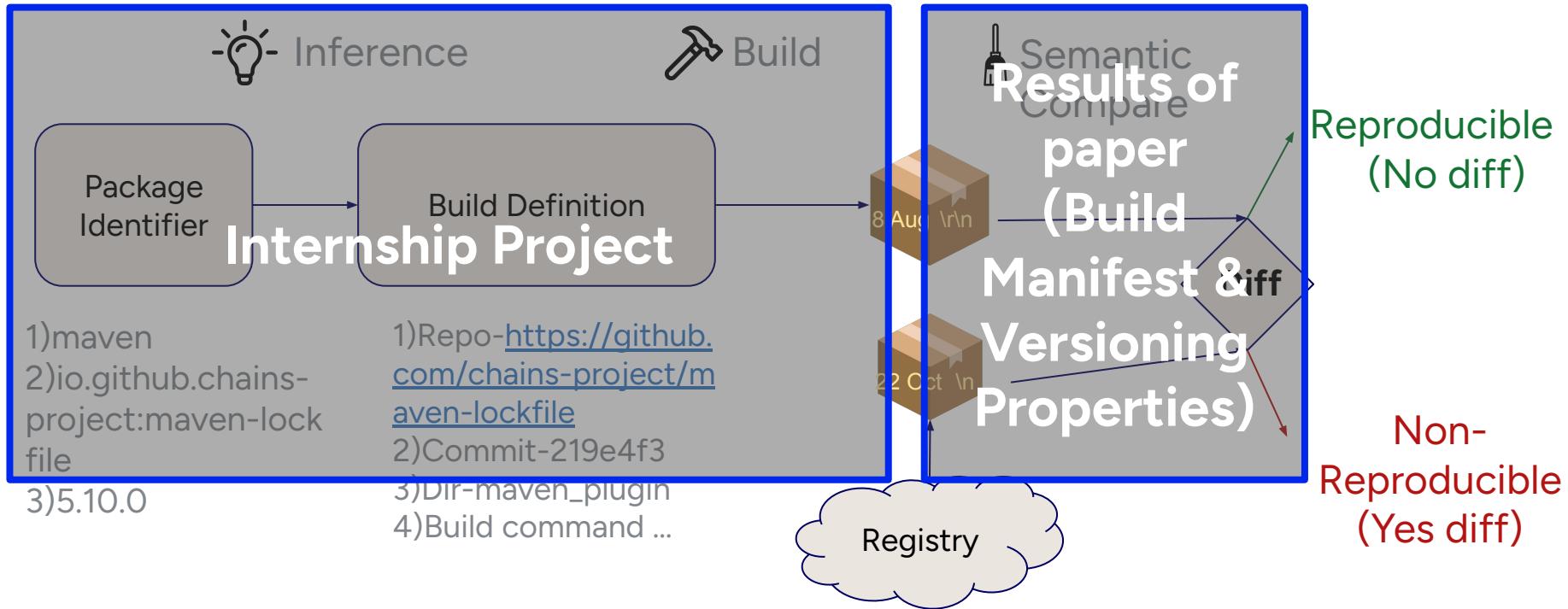
# Third goal

Integrate improvements in OSS Rebuild as part of  
internship at Google



# What is OSS Rebuild?

Rebuilder infrastructure for applying reproducible builds at scale.



# Reproducible Central

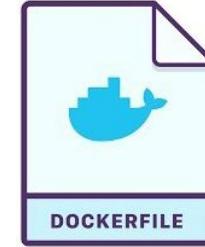
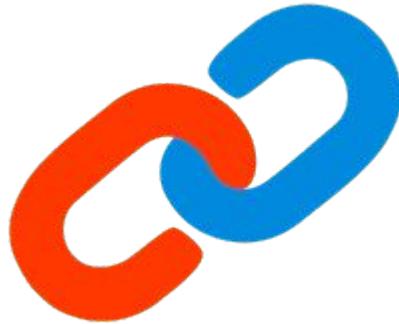
```
groupId=io.github.chains-project
artifactId=maven-lockfile
version=5.10.0
gitRepo=https://github.com/chains-
project/maven-lockfile.git
gitTag=${artifactId}-${version}
tool=mvn
jdk=17
newline=lf
command="mvn clean package
-DskipDepClean -DskipTests
-Dmaven.javadoc.skip -Dpgp.skip"
```

# OSS Rebuild

```
#syntax=docker/dockerfile:1.10
FROM docker.io/library/debian:trixie-20250203-slim
RUN <<'EOF'
set -eux
apt update
apt install -y git wget maven
EOF
RUN <<'EOF'
set -eux
mkdir -p /src && cd /src
git clone https://github.com/chains-project/
maven-lockfile .
git checkout --force
'219e4f3829570e4la62b17a1410ec615ca585af1'
mkdir -p /opt/jdk
wget -q -O - "https://download.java.net/java/GA/jdk17/
0d48333a00540d886896bac774ff48b/35/GPL/
openjdk-17_linux-x64_bin.tar.gz" | tar -xzf -
strip_components=1 -C /opt/jdk
EOF
RUN cat <<'EOF' >/build
set -eux
export JAVA_HOME=/opt/jdk
export PATH=$JAVA_HOME/bin:$PATH
mvn clean package -DskipTests --batch-mode -f
maven_plugin -Dmaven.javadoc.skip=true
chmod +444 /src/maven_plugin/target/maven-lockfile-5.
10.0.jar
mkdir -p /out && cp /src/maven_plugin/target/
maven-lockfile-5.10.0.jar /out/
EOF
WORKDIR "/src"
ENTRYPOINT ["/bin/sh", "/build"]
```

# Rebuild Attestations

```
go run ./cmd/oss-rebuild/ get npm colors 1.4.0  
--output=dockerfile
```



# Running benchmark

1. Benchmark
  - a. spoon-core
  - b. maven-lockfile
  - c. preemptiv
2. go run tools/ctl/ctl.go run-bench attest demo.json --local  
--max-concurrency 3
3. go run tools/ctl/ctl.go tui --logs-bucket file:///tmp/oss-rebuild
4. Show results and diff which failed to match the upstream artifact



# Interpreting non-reproducibility

Not always indicative of malicious build. It can be:

- Automation failure
- Environmental differences
- Non-deterministic builds

[24] <https://github.com/google/oss-rebuild/blob/main/docs/trust.md>



# Novel Features

1. Get repository URL using AI prompt
2. Get commit using matching upstream source
3. Infer build tool
4. Get JDK version bytecode version
5. Support for Gradle Builds

# google/oss-rebuild



Securing open-source package ecosystems by originating, validating, and augmenting build attestations.

Matthew, William, and  
Teddy

14

67

668

39

Contributors

Issues

Stars

Forks





# Future Work

1. Multiple repositories in Maven as an attack vector
2. Leveraging reproducibility for build performance
3. Exploring Java Platform Module System

# Conclusion

1. A comprehensive taxonomy of unreproducible builds in Java and their mitigation.
2. Artifact and bytecode canonicalization in conjunction can be used to mitigate unreproducibility.
3. OSS-Rebuild is awesome tool to build Java/NPM/PyPI/Cargo/Debian packages at scale.