

Introduction

Maven Central is repository that contains a large number of common libraries used in Java. It serves as a place where developers can store and share projects made in Java. This research will analyze Maven Central, specifically the space requirement of the libraries. This analysis provides insights into the average size of libraries and how to reduce size for future library maintainers. This is beneficial as smaller artifacts take less space in the repository and requires users to download smaller files.

Research Questions

1. How big is an average library on Maven?
2. How is the space requirement of the libraries distributed in the ecosystem?
3. What are the reasons for the larger sizes?

Key Terms

- .m2 folder: the local directory where maven caches libraries downloaded from the central repository.
- JAR: Archive files similar to zip used to package Java classes and related metadata and resources.
- POM: XML file that contains information about a project and is used by Maven to build the project.

Methods

- Figure 1 shows a high level overview of the architecture used in this research.
- First the Maven Index is read and stored in a table
- Then, we use a data selection strategy to select libraries from this table. The distribution of our sample can be seen in Figure 3, and the distribution of packages on Maven Central in Figure 2.
- After selecting which packages to analyze, their jars and pom files are downloaded into the local .m2 folder by the resolver.
- The individual jars and pom files are analyzed by extractors and metadata is stored in a database. This metadata includes the size of the package, the number of files, the number of direct dependencies and the number of transitive dependencies.

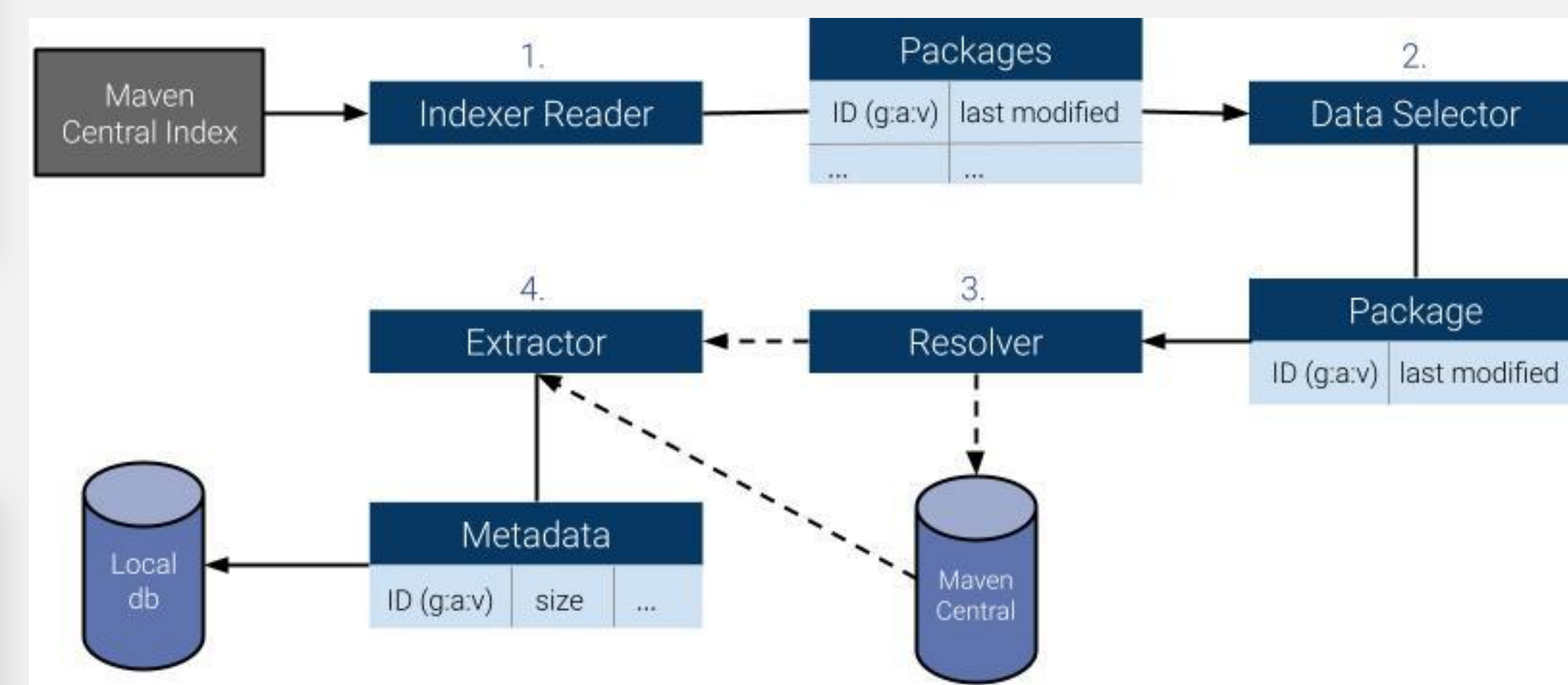


Figure 1: High level component diagram of main architecture

Results & Discussion

- Average space requirement of a library is 1447 KB, the median is 25.9 KB and the largest package is 986,867 KB.
- Figure 5 shows the distribution of size in Maven Central.
- There are two distinct categories of artifacts in Figure 4: Those that slightly increase in size when they add more files and those that don't have that many files to begin with and are usually small but can become very large too. We manually analyzed the cases where libraries are massive but the number of files is low.
- The analysis revealed that almost all these packages are machine learning or big data projects with massive data files.
- We recommend not including files like this in the package and instead uploading them somewhere else and document how to connect to them.
- The threats to validity in this study are that not every package in our sample could be analyzed because in some cases, the parent pom was hosted on a different repository. Furthermore our data selection strategy does not analyze every version of each package but only one version per package.

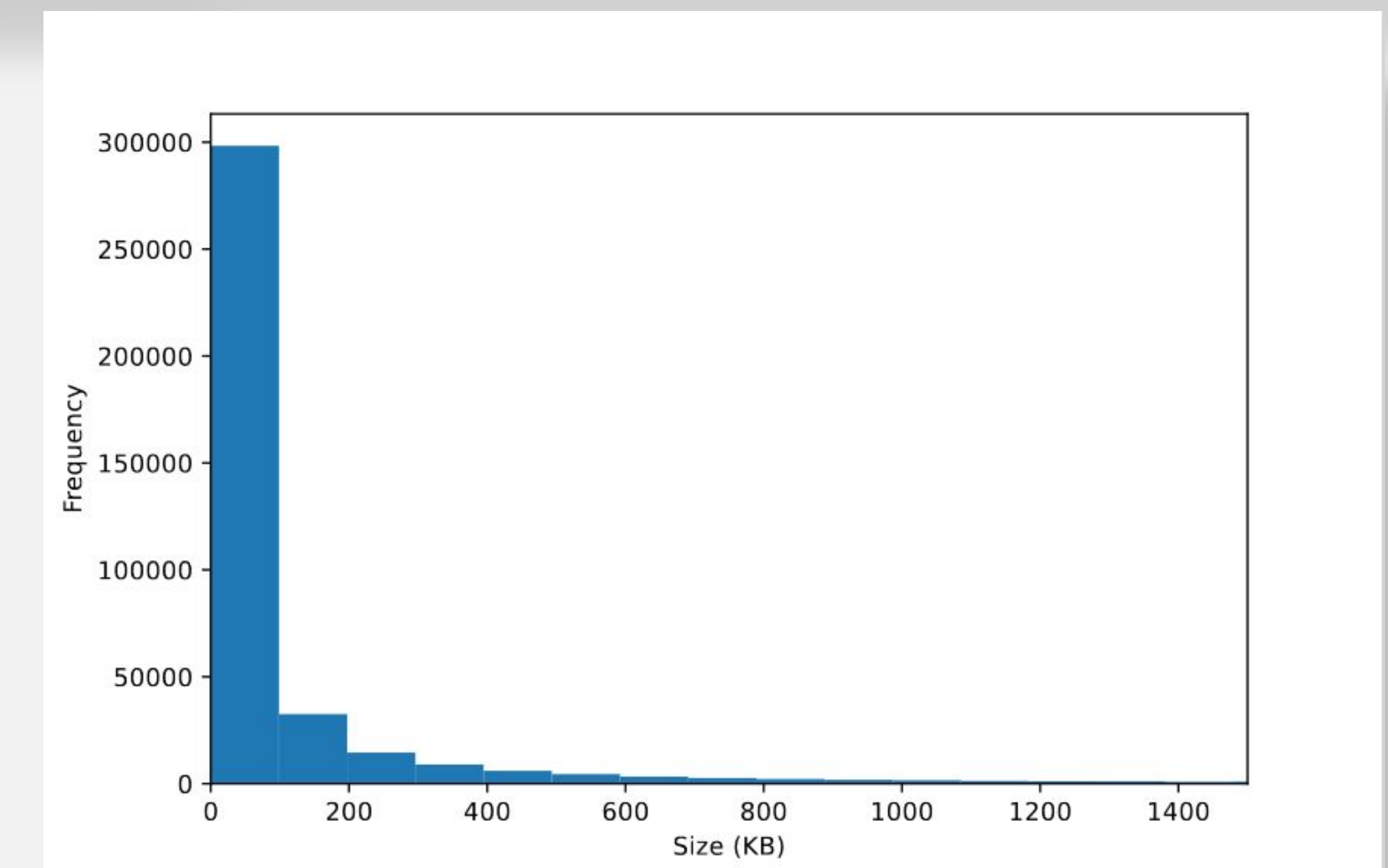


Figure 5: Distribution of size in Maven Central

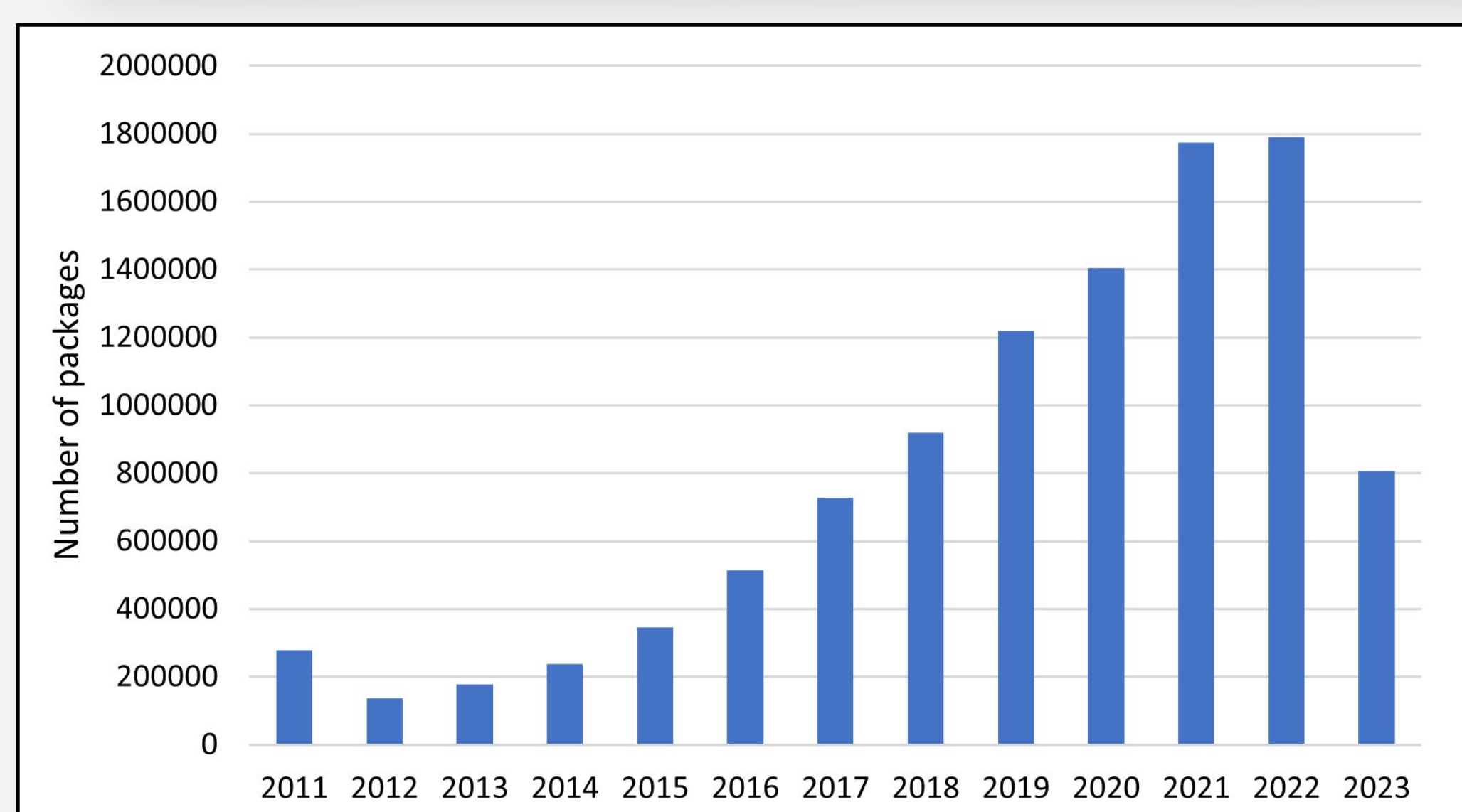


Figure 2: Distribution of packages released on Maven Central per year

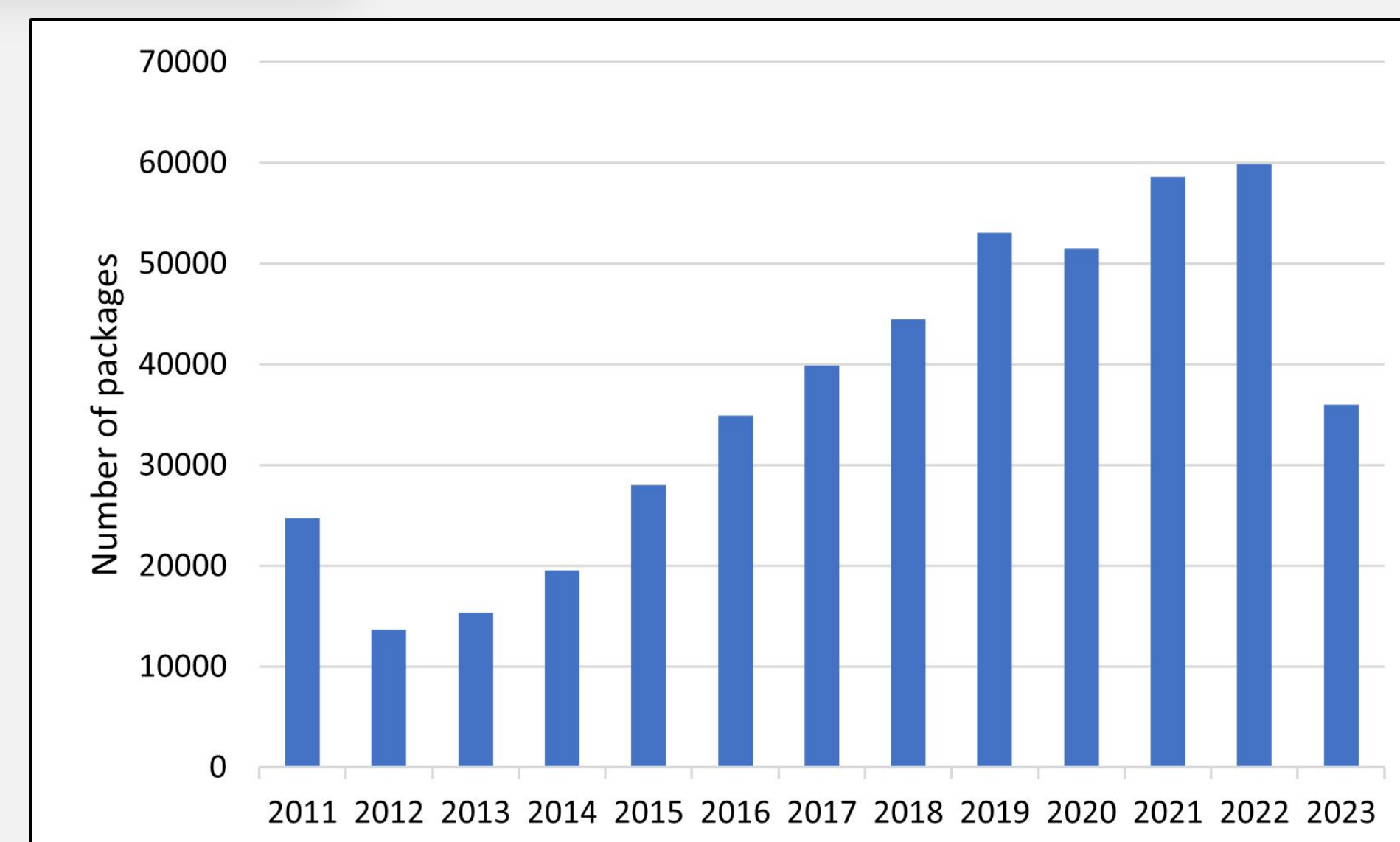


Figure 3: Selected distribution of packages in our sample

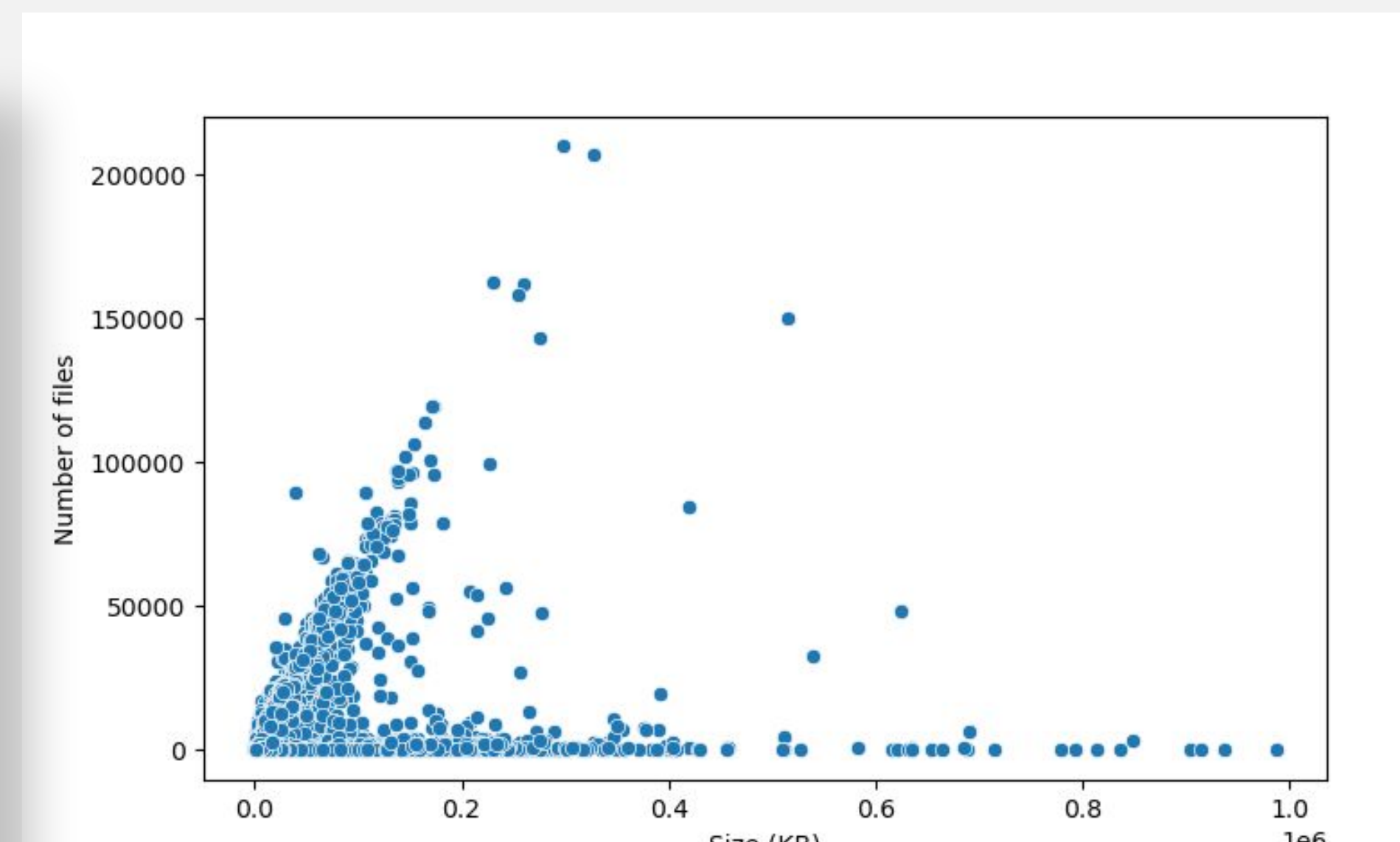


Figure 4: Correlation between size and number of files

References

- [1] Cesar Soto-Valero, Nicolas Harrand, Martin Monperrus, and Benoit Baudry. A comprehensive study of bloated dependencies in the maven ecosystem. 2021.
- [2] Dimitris Mitropoulos, Vassilios Karakoidas, Panos Louridas, Georgios Gousios, and Diomidis Spinellis. The bug catalog of the maven ecosystem. 2014.

