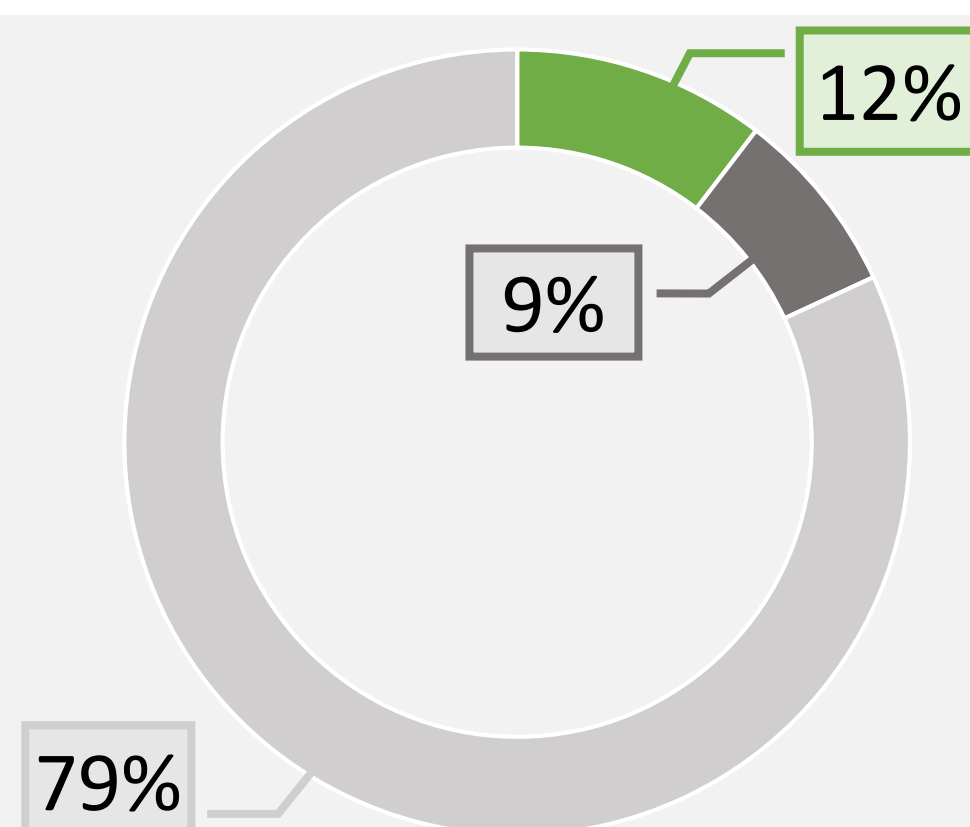




# DéjàVu

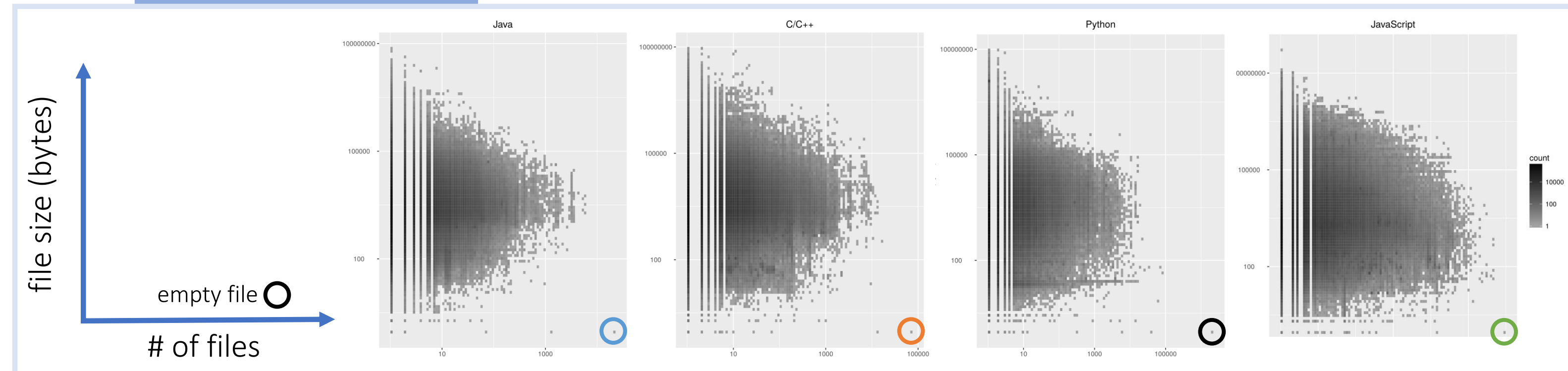
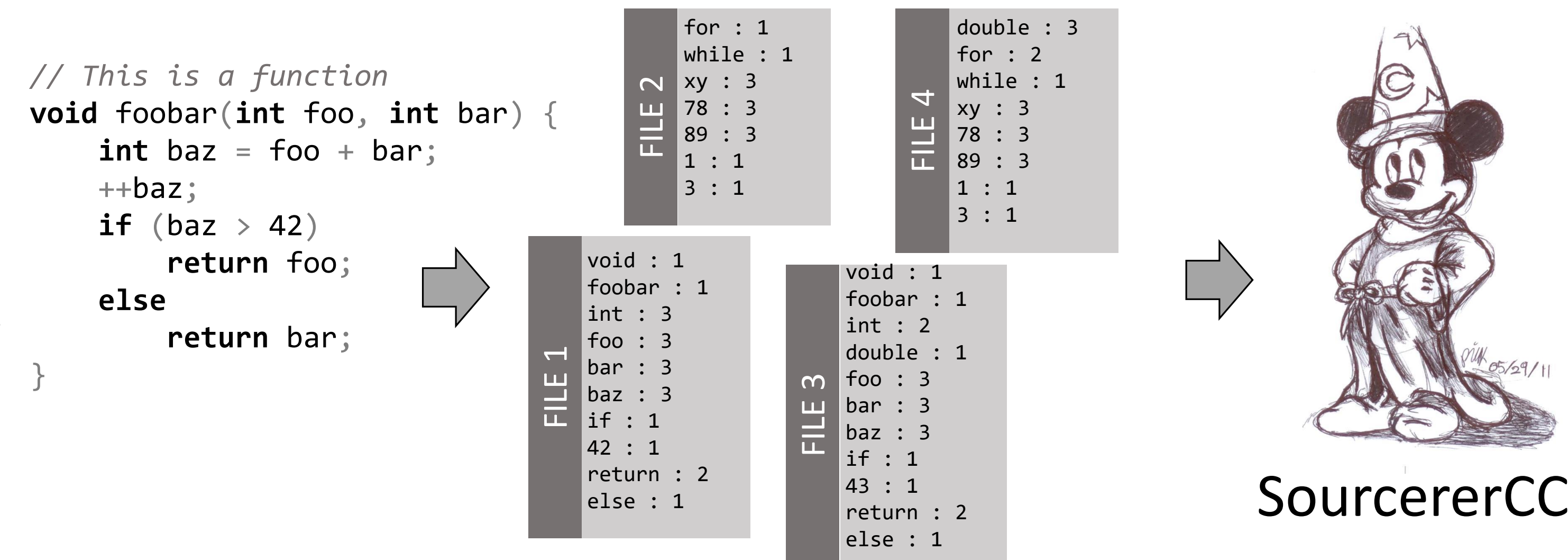
## A Map of Code Duplicates on GitHub

Introduction



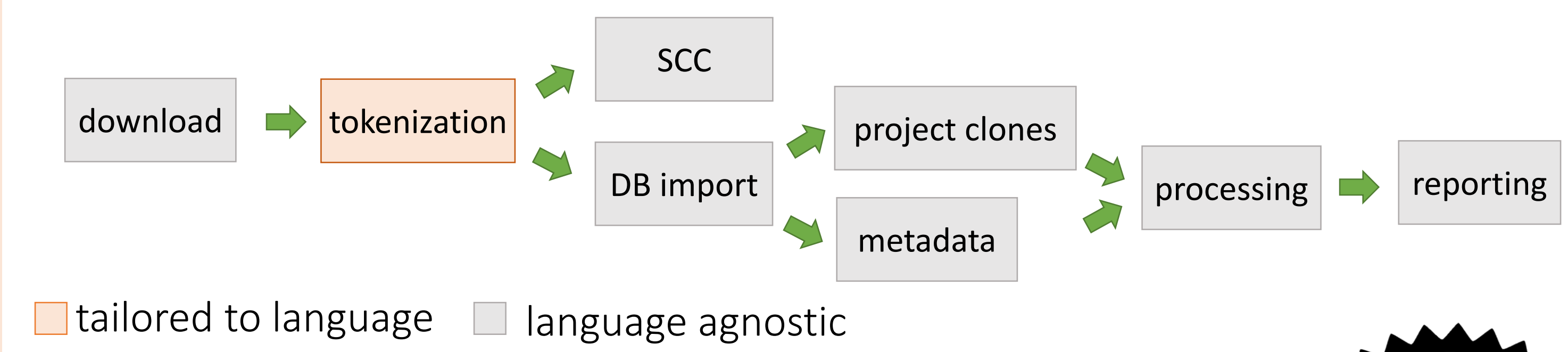
If imitation is flattery, then the Internet is surprisingly full of praise. Specifically, the code programmers check into large scale software repositories is riddled with duplicates. Our paper involves the analysis of 4.5M non-fork projects on GitHub; 428M files in Java, C++, Python, and JavaScript. We found only 85M unique files. In other words, 79% of GitHub consists of duplicates. There is variability between languages. JavaScript has the highest rate, only 4% of the files are distinct. Java has the least, 60% of files are distinct.

Similarity Levels



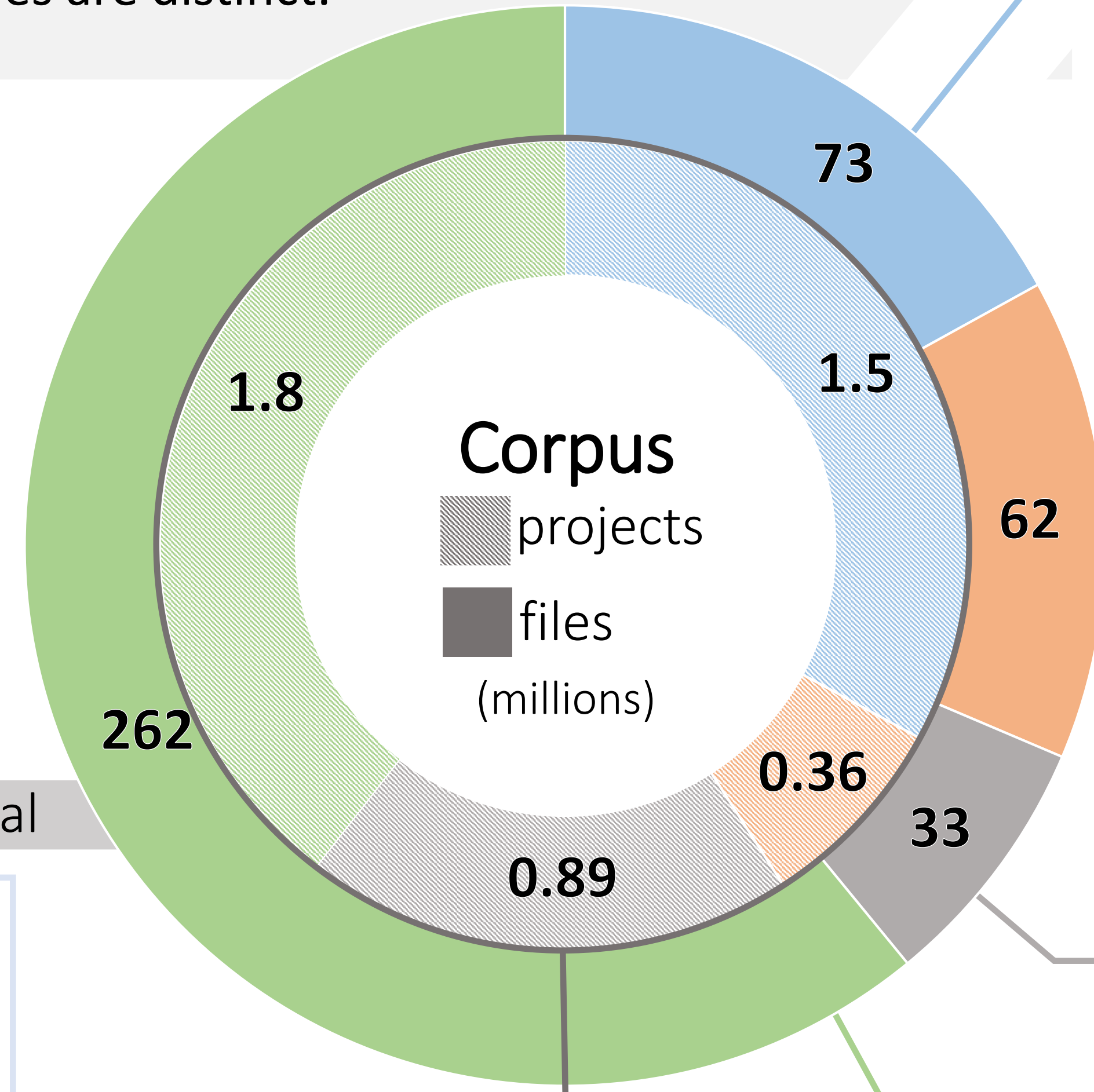
Artifact

All our data and results are available for download. Furthermore our entire pipeline, from downloading the projects to producing the clone maps is entirely automated in the form of R notebooks.

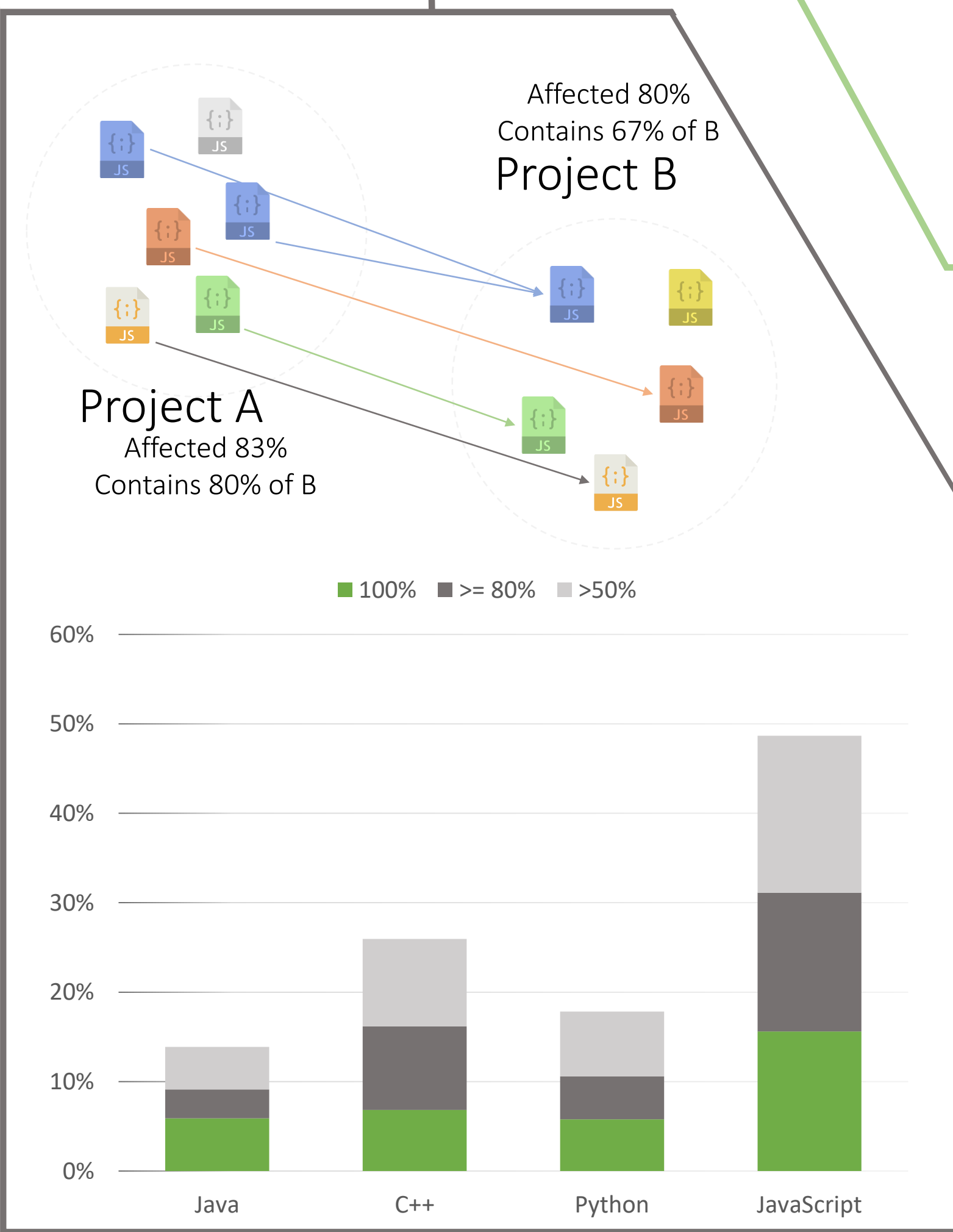


References

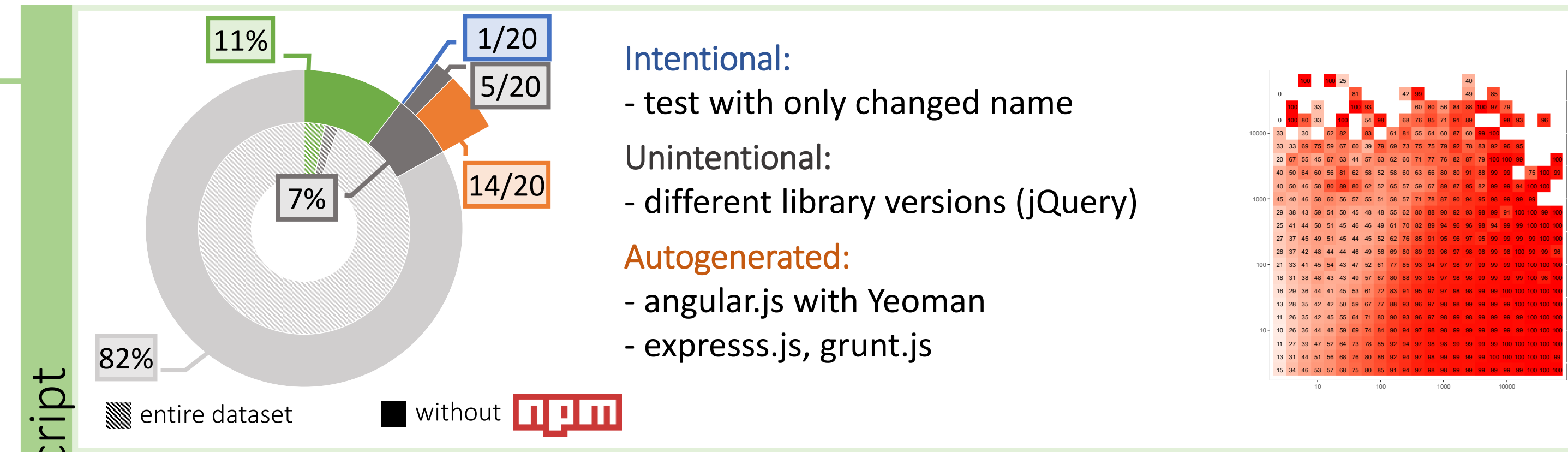
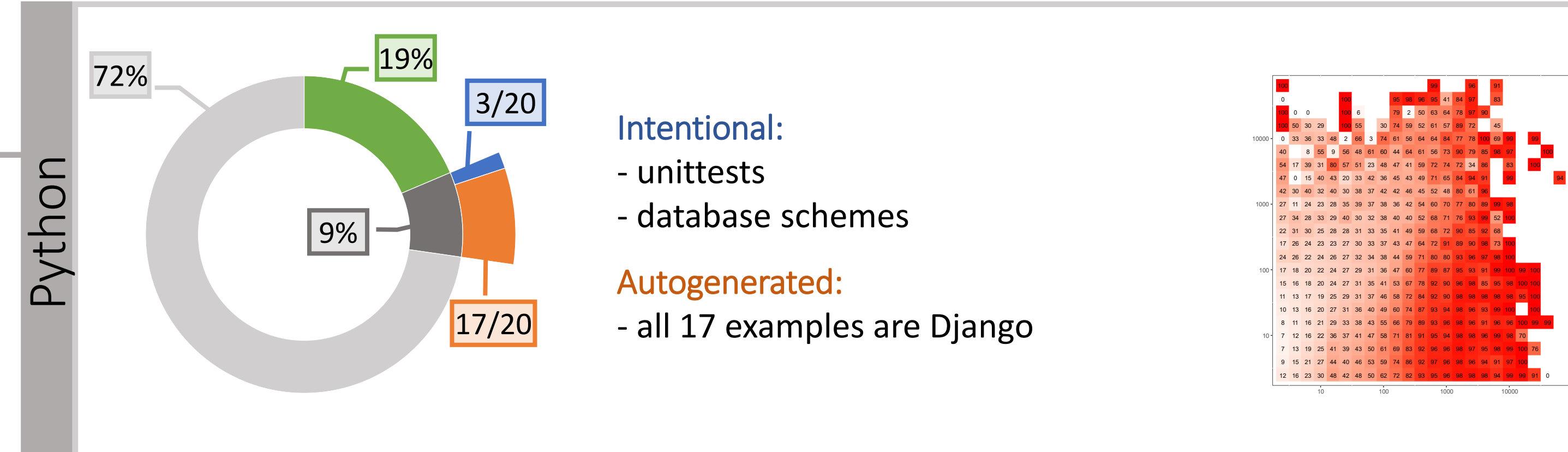
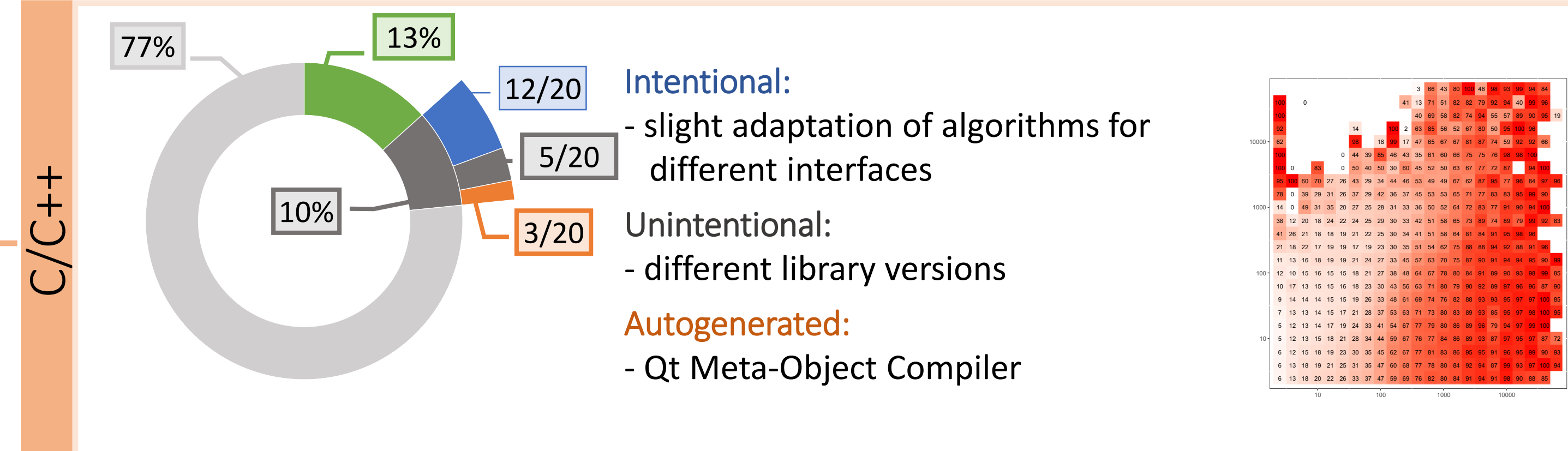
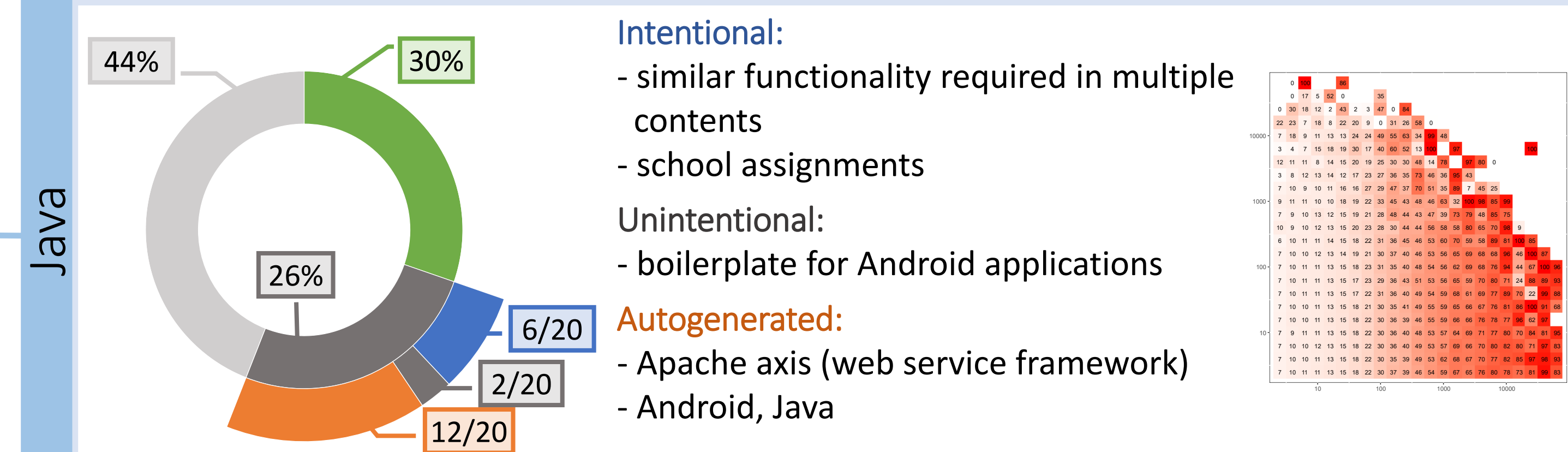
<http://mondego.ics.uci.edu/projects/dejavu> - project website  
<https://github.com/PRL-PRG/dejavu-artifact> - artifact  
<http://dejavu.ics.uci.edu> - web service



Project Inclusion



unique, similar files, identical files, autogenerated, unintentional, intentional. # of commits vs # of files.



In JavaScript, over 70% of the identical duplicates of files come from accidental inclusion of node.js project dependencies.

