# Patched Clones and Missed Patches among Variants of a Software Family

John Businge

Assistant Professor – UNLV

Never Work in Theory (NWiT) – Spring 2023
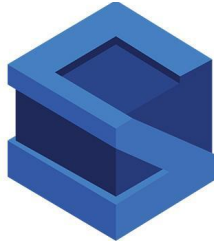
SECO-Assist

# Context

$425M

The Equifax data breach occurred between **May and July 2017** at the American credit bureau Equifax. Private records of 147.9 million Americans along with 15.2 million British citizens and about 19,000 Canadian citizens were compromised in the breach, making it one of the largest cybercrimes related to identity theft.
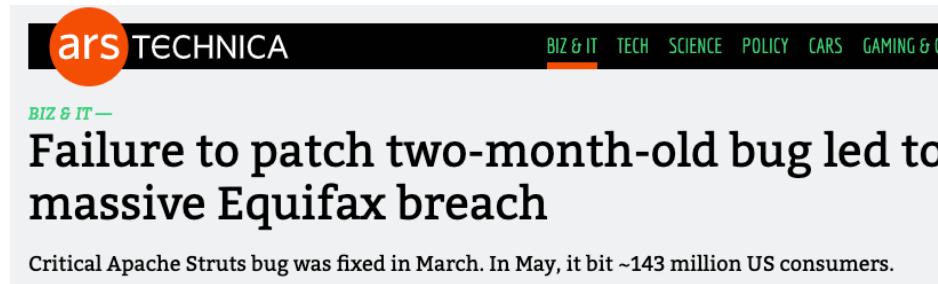
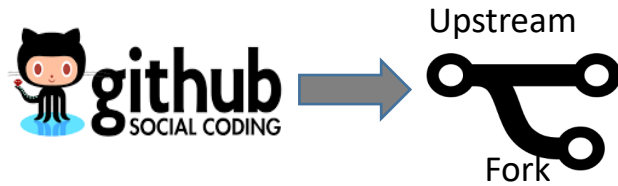Wired Magazine, "Equifax has no excuse", September 2017

Recommender system

open source

March 2017

CVE-2017- 5638

**ars** TECHNICA      BIZ & IT   TECH   SCIENCE   POLICY   CARS   GAMING & CUL

BIZ & IT —

# Failure to patch two-month-old bug led to massive Equifax breach

Critical Apache Struts bug was fixed in March. In May, it bit ~143 million US consumers.

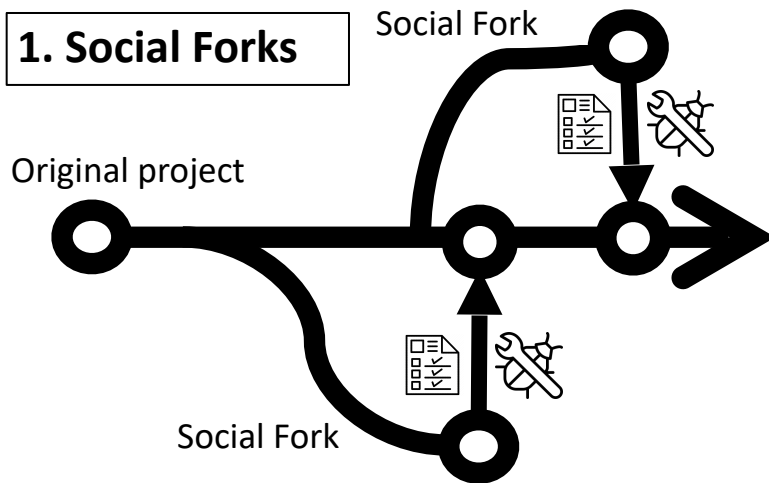**EQUIFAX** **DATA BREACH** May 2017

https://www.istockphoto.com/

# Patched Clones and Missed Patches among Variants of a Software Family


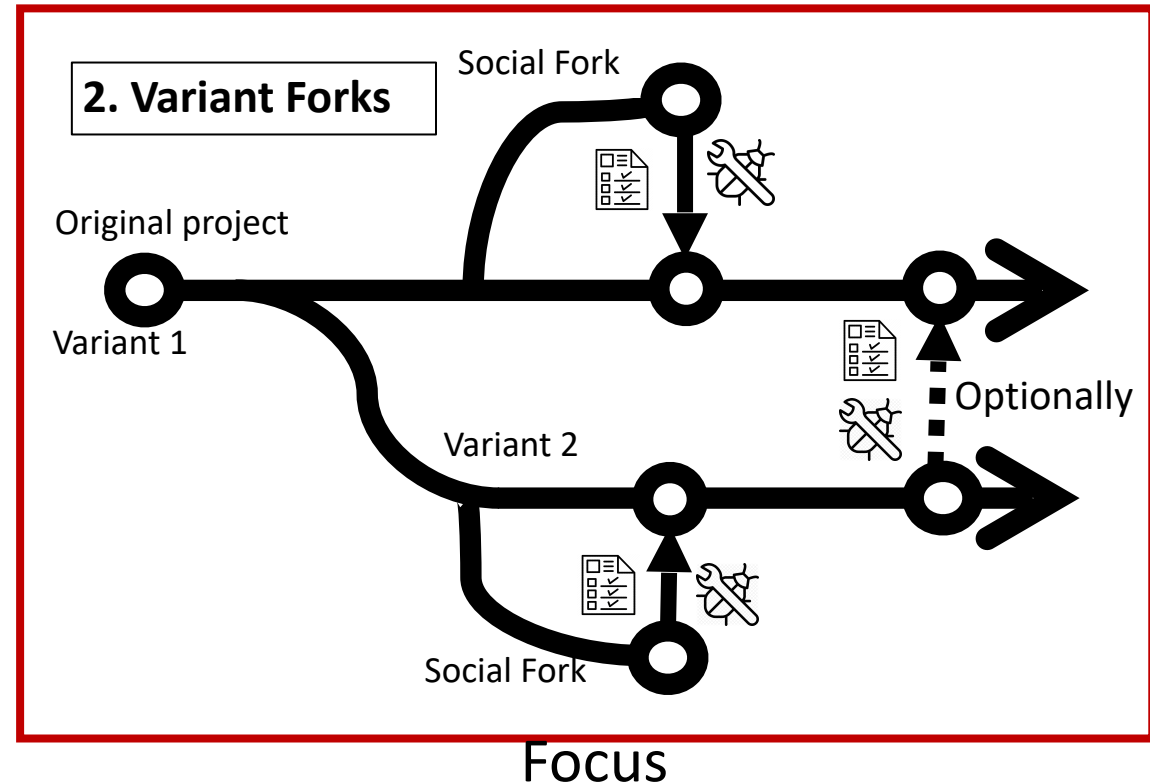
Software family — >= 2 variants

1. Social Forks

2. Variant Forks

Focus

# Reuse and maintenance practices among divergent forks in three software ecosystems

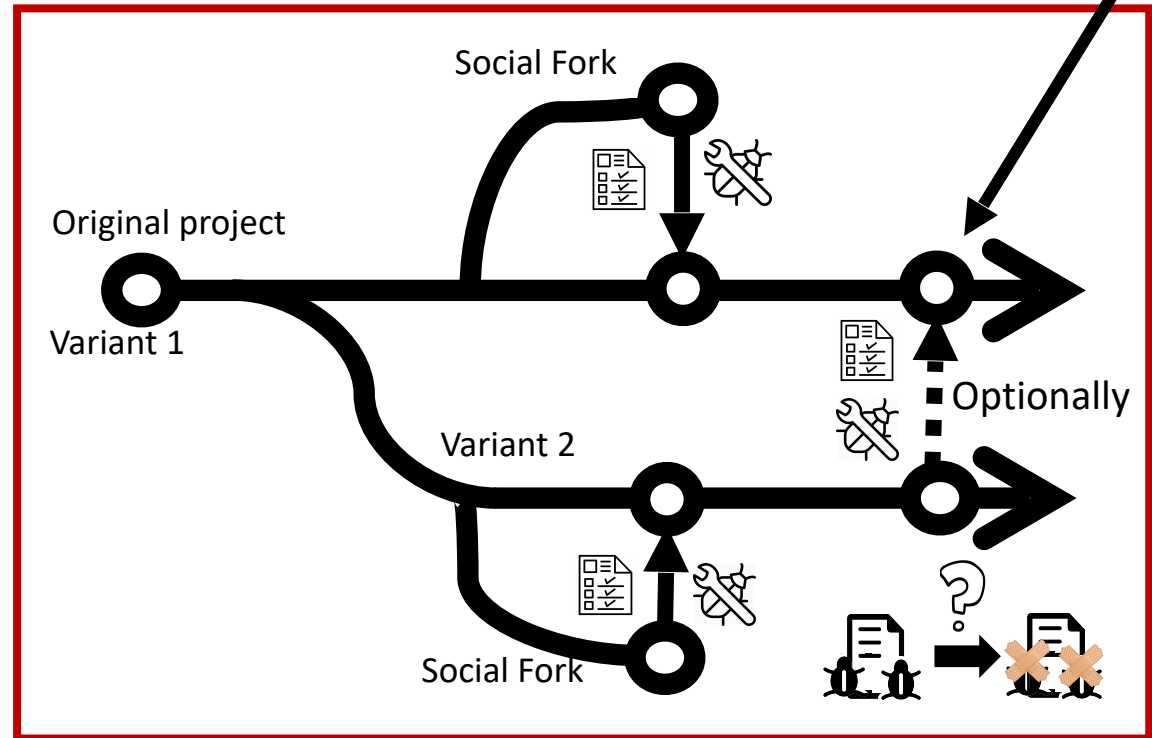John Businge ✉, Moses Openja, Sarah Nadi & Thorsten Berger

## Abstract

With the rise of social coding platforms that rely on distributed version control systems, software reuse is also on the rise. Many software developers leverage this reuse by creating variants through forking, to account for different customer needs, markets, or environments. Forked variants then form a so-called software family; they share a common code base and are maintained in parallel by same or different developers. As such, software families can easily arise within software ecosystems, which are large collections of interdependent software components maintained by communities of collaborating contributors. However, little is known about the existence and characteristics of such families within ecosystems, especially about their maintenance practices. Improving our empirical understanding of such families will help build better tools for maintaining and evolving such families. We empirically explore maintenance practices in such fork-based software families within ecosystems of open-source software. Our focus is on three of the largest software ecosystems existence today: Android,
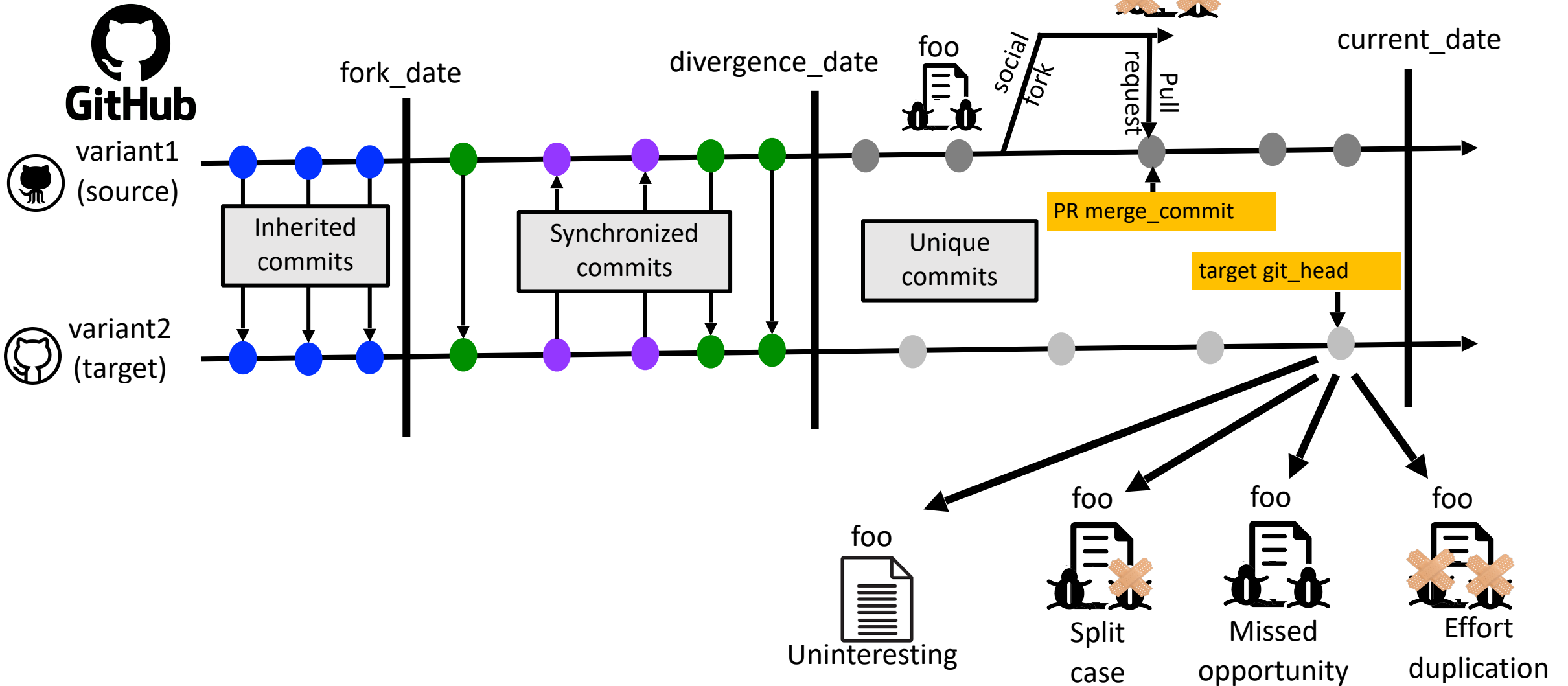
**>10K variants from three ecosystems on Github**

**Rarely share updates**



4

# Problem

# Concrete Example



extraction_date 2023-03-27

**415** linkedin unique commits

**1,787** apache unique commits

This is the version of Kafka running at LinkedIn.

6

## Missed Opportunity

**Buggy code from upstream**     qmk/qmk_firmware

1 file  - Pull request

extraction_date 2021-07-20

```
1              return;
2          }
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4      flashend = (uint32_t)((uint16_t *)SYMVAL(__eeprom_workarea_end__) - 1);
5  }
```
⟵ Buggy line

gcc10 [...] build warning #12587

**Patched code from upstream (Pull request)**

```
1              return;
2          }
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4      flashend = (uint32_t)(p - 1);
5  }
```
⟵ Patched line

**Diff for patch in upstream**

```
1    @@ -363,7 +363,7 @@
2
3        } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4    -   flashend = (uint32_t)((uint16_t *)SYMVAL(__eeprom_workarea_end__) - 1);
5    +   flashend = (uint32_t)(p - 1);
```
Hunk

**File from divergent fork at** `git_head`     sekigon-gonnoc/qmk_firmware

```
1              return;
2          }
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4      flashend = (uint32_t)((uint16_t *)SYMVAL(__eeprom_workarea_end__) - 1);
5  }
```
⟵ Buggy line

# Research Questions

1. **RQ1:** How many cases of effort duplication and missed opportunities exist between divergent variants?

2. **RQ2:** How much patch technical lag exists between the source and target variants in divergent variants?

# Method

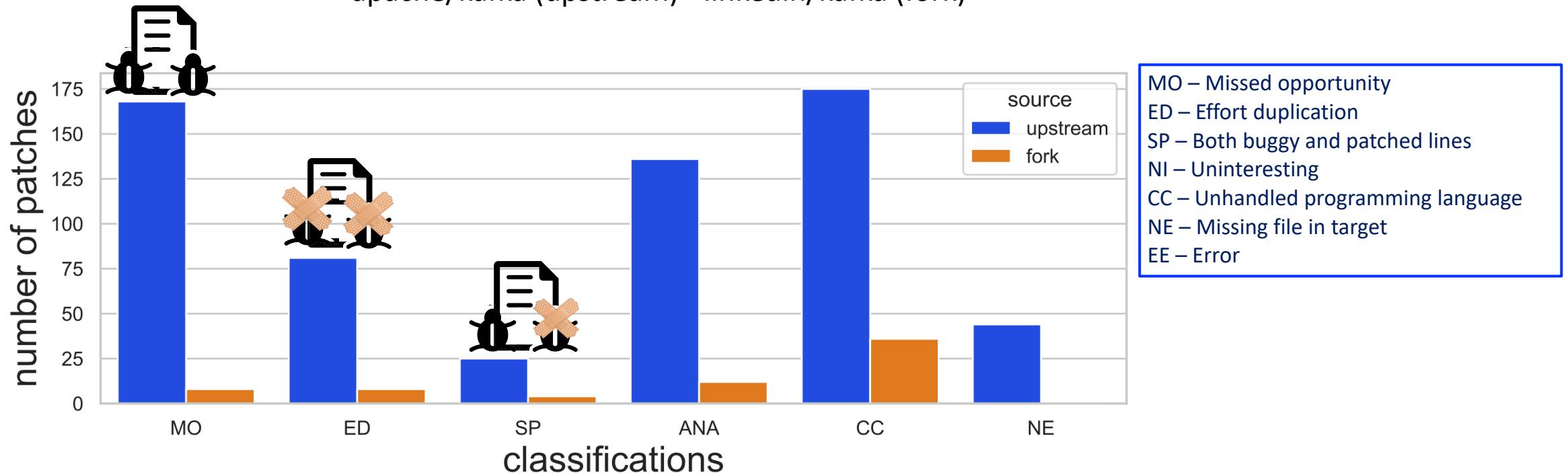keywords {fix, fixes, resolves, ...}

Search PR Title

# Results

**RQ1:** How many cases of effort duplication and missed opportunities exist between divergent variants?
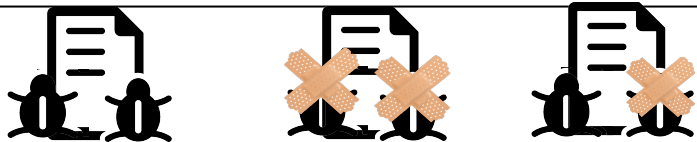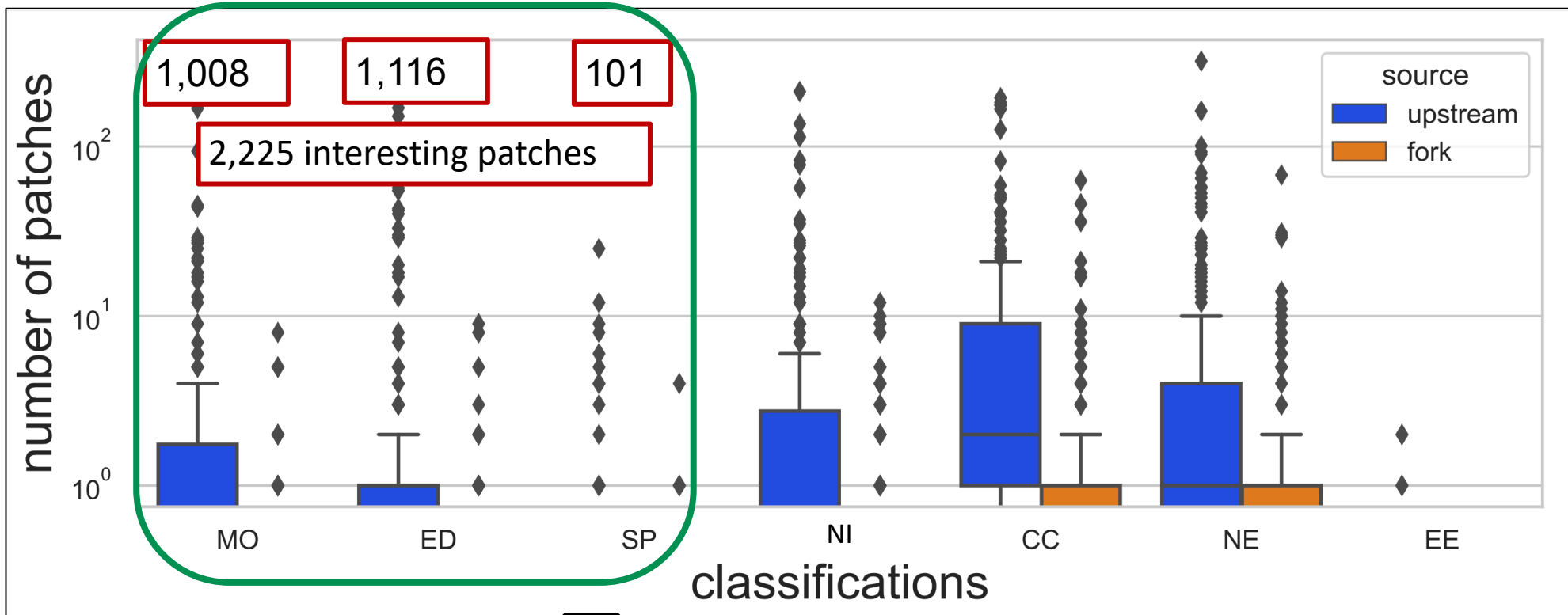
apache/kafka (upstream) - linkedin/kafka (fork)



MO – Missed opportunity
ED – Effort duplication
SP – Both buggy and patched lines
NI – Uninteresting
CC – Unhandled programming language
NE – Missing file in target
EE – Error

**RQ1:** How many cases of effort duplication and missed opportunities exist between divergent variants?



8,323 patches from 364 source variants

| Precision | Recall | Accuracy | F1-Score |
|-----------|--------|----------|----------|
| 91.0% | 80.2% | 88.0% | 85.3% |

1,008    1,116    101
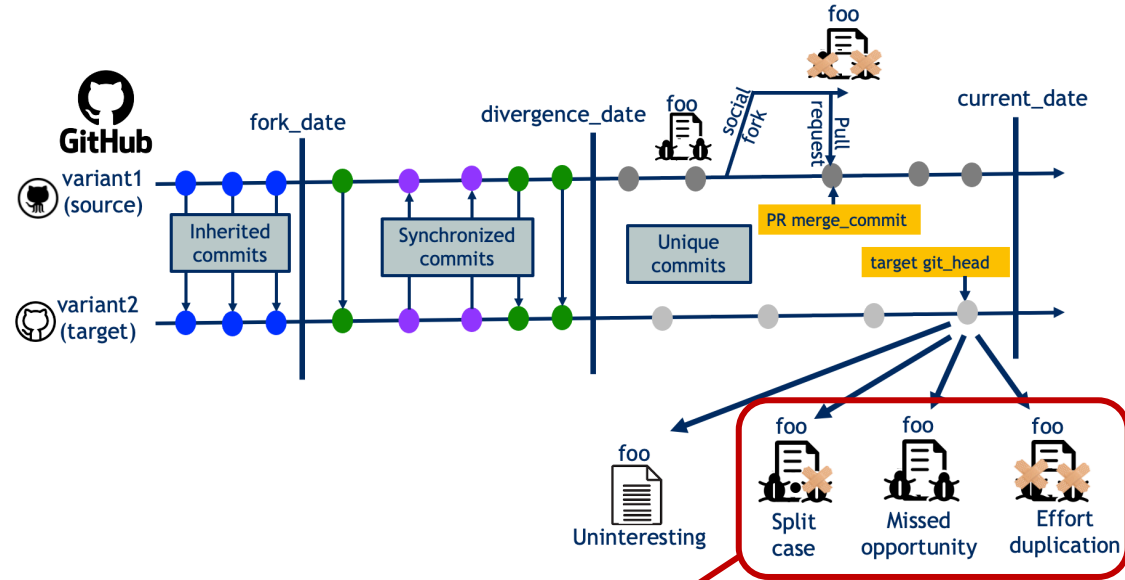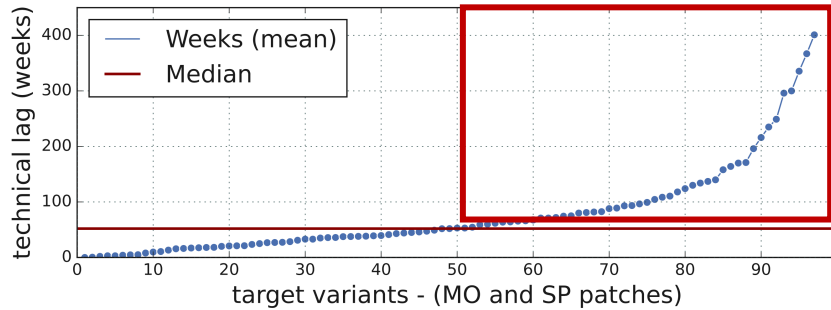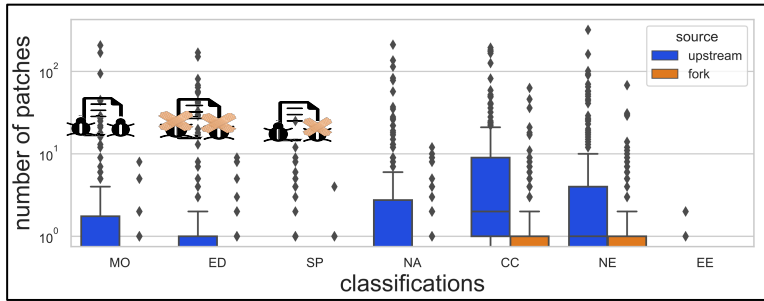
2,225 interesting patches

# Results

**RQ2:** How much patch technical lag exists between the source and target variants in divergent variants?
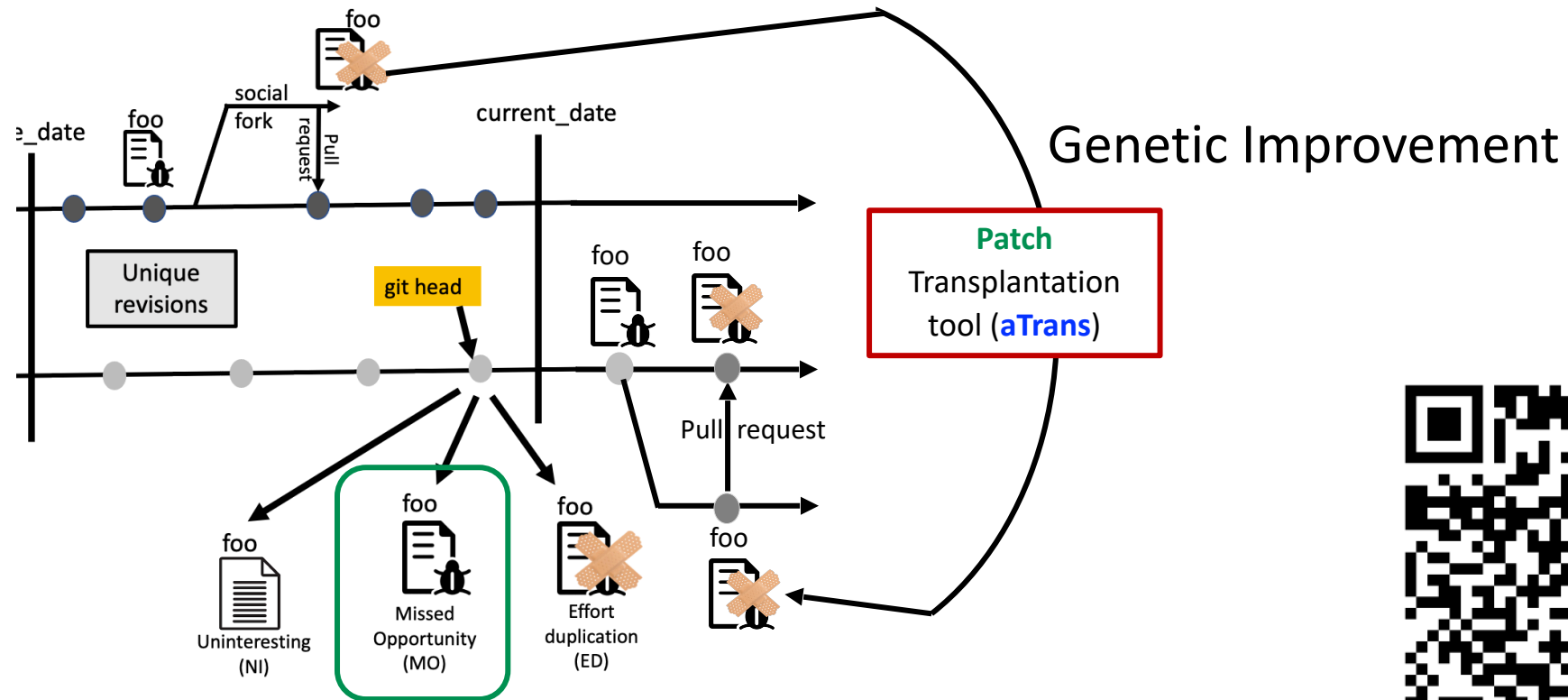
# What do we learn from the results?



PaReco: **Proof-of-Concept patch recommender tool**

# Current Work on PaReco



Genetic Improvement

Patch Transplantation tool (**aTrans**)

PaReco: **Proof-of-Concept patch recommender tool**