

# Compression of Programs and the Similarity Distance

**KIREPRO1PE Research Project, MSc. Computer Science, ITU - 10th of June 2025**

Jonas Nim Røssum <jglr@itu.dk>

## Background

- *Lines of Code Changed* (LoCC)
  - De facto standard for measuring code changes
  - Has its limitations (e.g. structural changes, formatting changes, etc.)

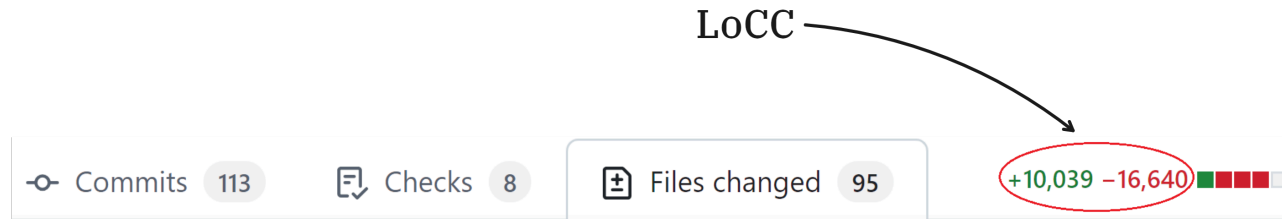


Figure 1: LoCC in a GitHub Pull Request

## Project goal and findings

- Find a new metric to address limitations of *Lines of Code Changed* (LoCC)
- *Difference in Compression Distance* ( $\Delta CD$ )

### Research questions

- ? Is  $\Delta CD$  correlated with LoCC?
- ? Can  $\Delta CD$  discriminate between commit types?
- ? What are the advantages / limitations of  $\Delta CD$ ?

### Findings

- Partial linear correlation,  $R^2 = \{0.8, 0.7\}$
- For Commitizen<sup>1</sup> repo, **features** and **bug fixes** stand apart
- Robust to structural changes, survivorship bias / 250× slower than LoCC, scaling challenges

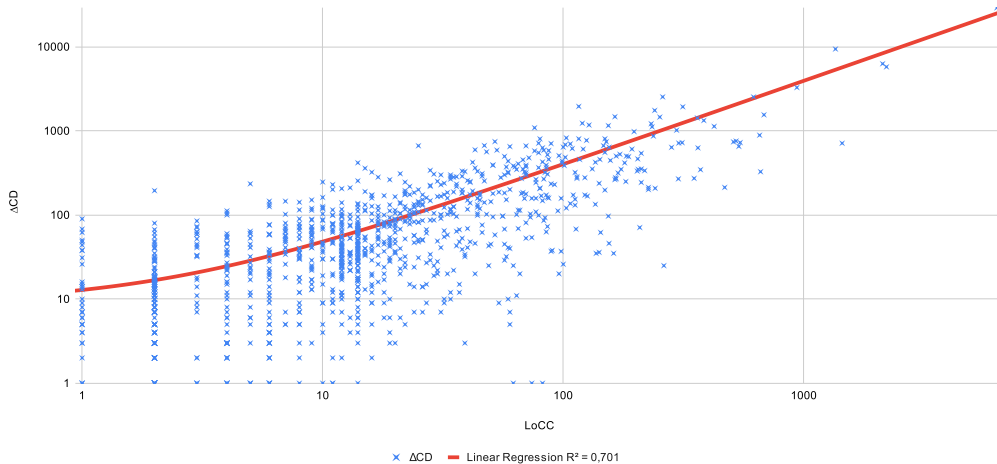
---

<sup>1</sup><https://github.com/commitizen-tools/commitizen/>

## RQ1: $\Delta$ CD correlation with LoCC

Linear regression  $R^2$  for **commitizen**: 0.7

LoCC vs  $\Delta$ CD for commitizen-tools/commitizen (github)



$\Delta$ CD and LoCC correlate, but not perfectly  $\rightarrow$   $\Delta$ CD captures more than raw line changes

## RQ2: Commit Type Discrimination

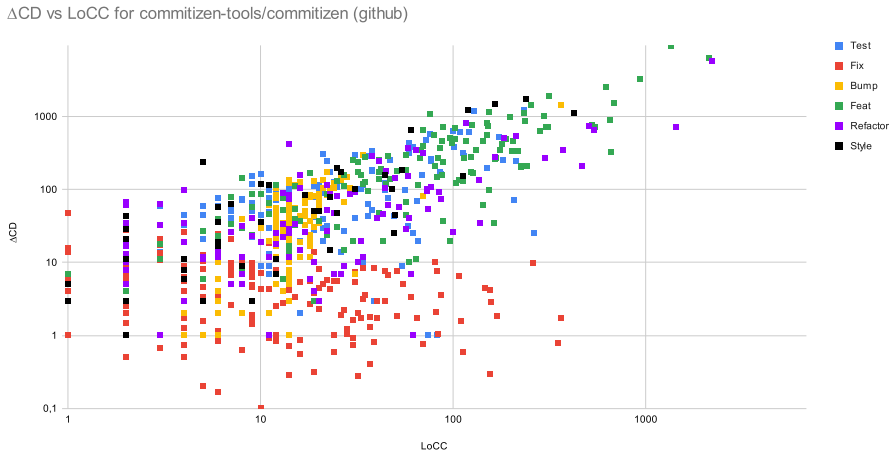
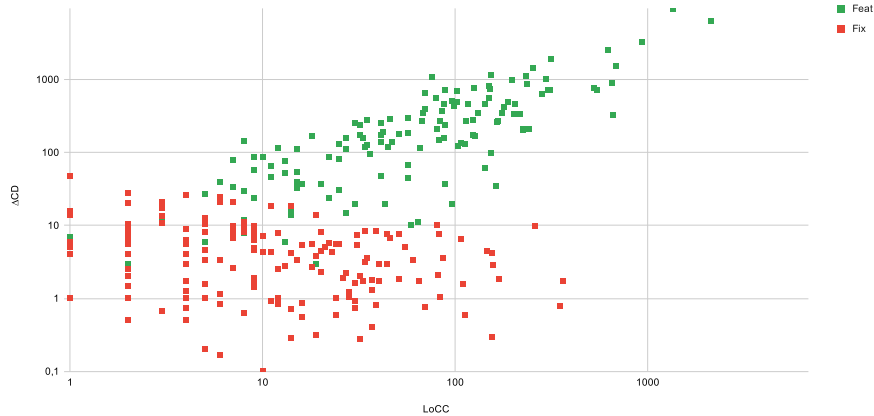


Figure 3: Commits in the Commitizen repository categorized using conventional keywords<sup>2</sup>

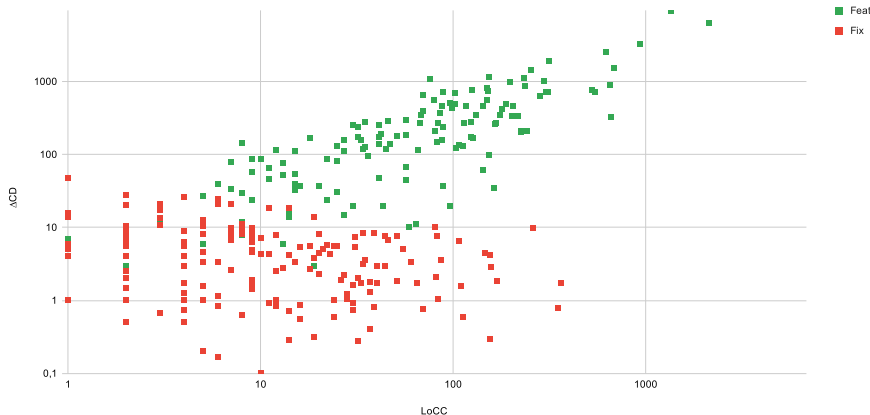
<sup>2</sup><https://www.conventionalcommits.org/en/v1.0.0/>

## RQ2: Commit Type Discrimination



## RQ2: Commit Type Discrimination

$\Delta$ CD vs LoCC for commitizen-tools/commitizen (github)

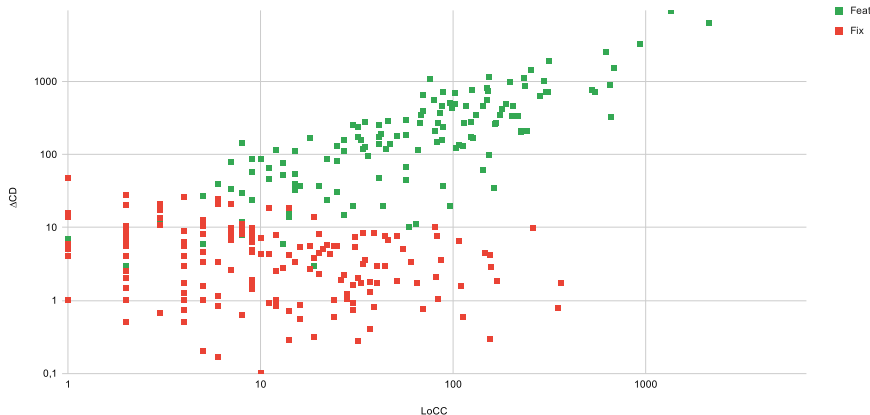


**Bug Fixes:** lower  $\Delta$ CD, changes to existing code

**Features:** higher  $\Delta$ CD, typically novel code

## RQ2: Commit Type Discrimination

$\Delta$ CD vs LoCC for commitizen-tools/commitizen (github)



**Bug Fixes:** lower  $\Delta$ CD, changes to existing code

**Features:** higher  $\Delta$ CD, typically novel code

✓  $\Delta$ CD can partly discriminate between some commit types, at least for this project



## RQ3: Advantages - Robust to structural changes

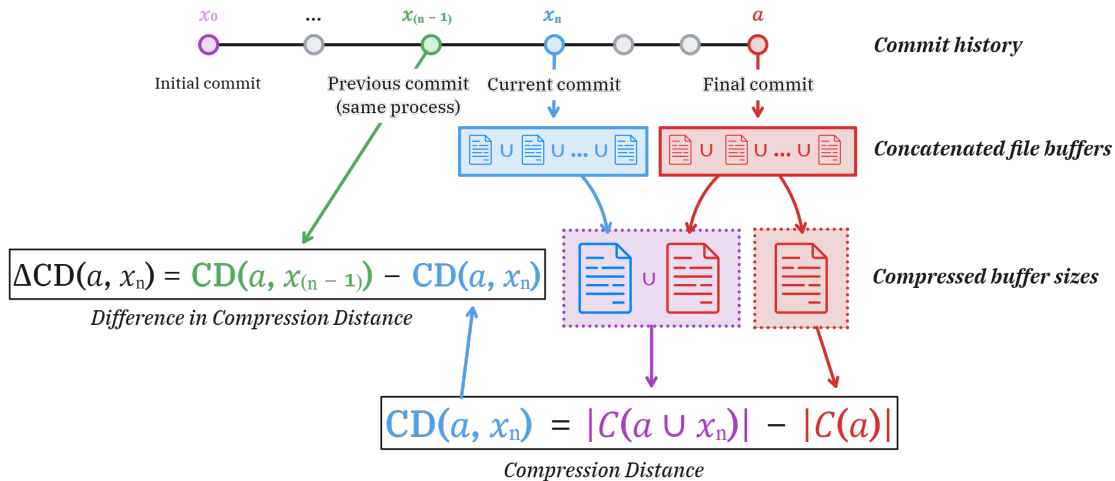
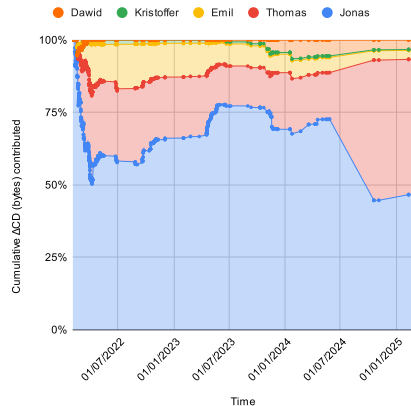
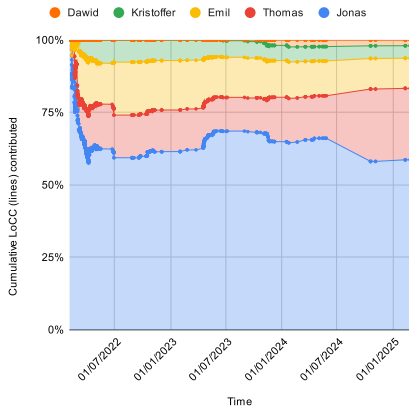


Figure 7:  $\Delta CD$  (Difference in Compression Distance) Explained

✓  $\Delta CD$  is insensitive to project structure

## RQ3: Advantages - Survivorship Bias

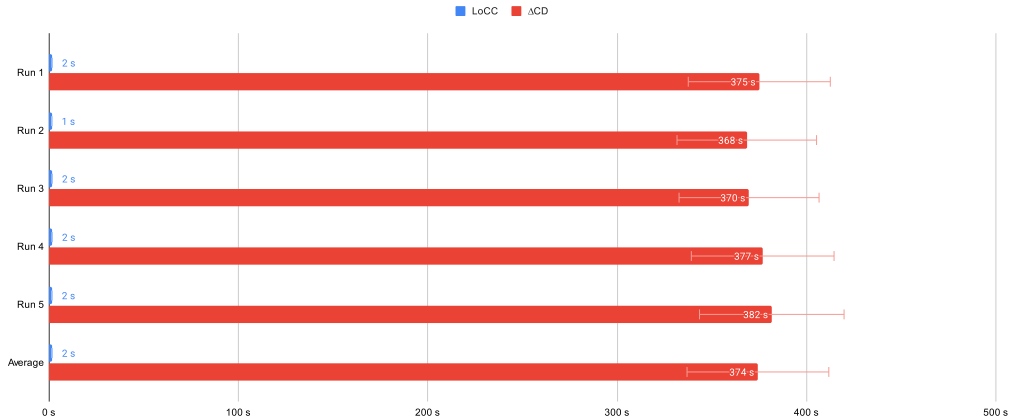
- Example: **Thomas' thesis work in Git Truck**
- According to LoCC (left), Thomas is responsible for 25% of the contributions project
- According to  $\Delta$ CD (right), Thomas is responsible for 46% of the final revision



✓  $\Delta$ CD reflects **lasting impact** on the codebase using survivorship bias

# Limitations: Performance and Scalability

LoCC vs  $\Delta$ CD for commitizen-tools/commitizen (github, 1977 commits)



## Future work

Performance and scalability

Generalize findings

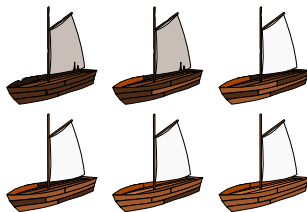
Robustness to formatting changes etc.

Integration

Preprocessing

Use cases

## Thank You - Questions?



Project work: Jonas Nim Røssum <jglr@itu.dk>

Original idea: Christian Gram Kalhauge <chrg@dtu.dk>

Source code: [github.com/git-truck/git-truck](https://github.com/git-truck/git-truck)