# 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
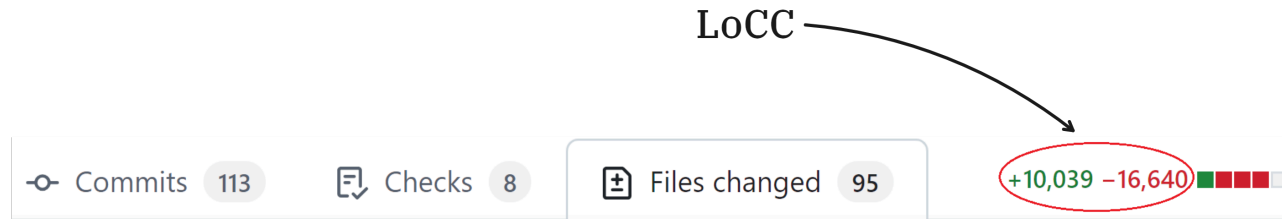  - Has it's limitations (e.g. structural changes, formatting changes, etc.)

LoCC



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?
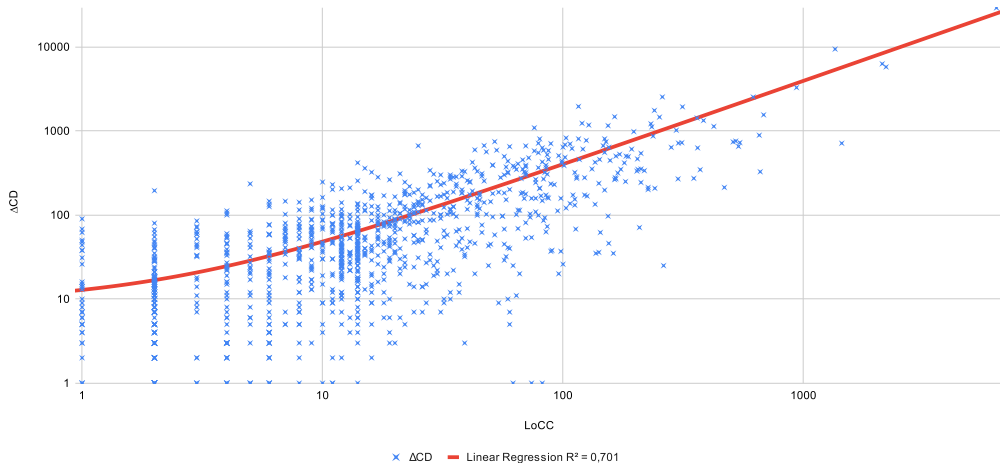
? What are the advantages / limitations of $\Delta$CD?

**Findings**

$\rightarrow$ Partial linear correlation, $R^2 = \{0.8, 0.7\}$

$\rightarrow$ For Commitizen[1] repo, features and bug fixes stand apart

$\rightarrow$ Robust to structural changes, survivorship bias / $250\times$ slower than LoCC, scaling challenges

---

[1]https://github.com/commitizen-tools/commitizen/

# RQ1: $\Delta$CD correlation with LoCC

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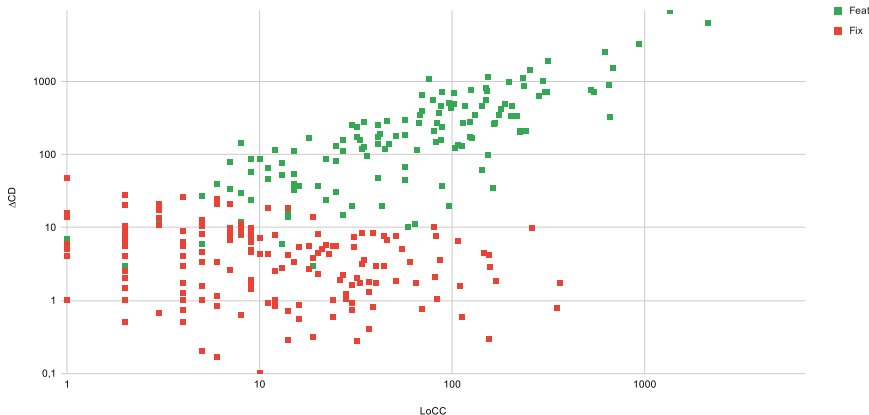


LoCC vs $\Delta$CD for commitizen-tools/commitizen (github)

✅ $\Delta$CD and LoCC **partially correlate** $\rightarrow$ $\Delta$CD captures more than raw line changes

# RQ2: Commit Type Discrimination



ΔCD 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

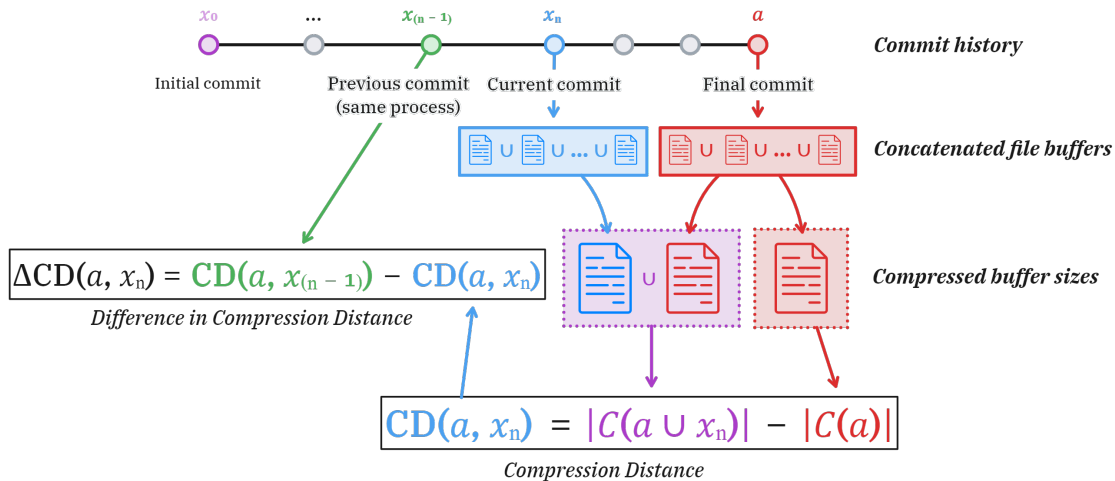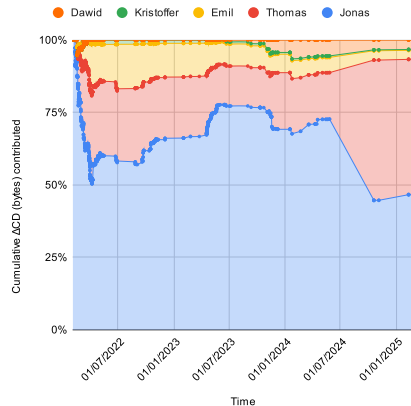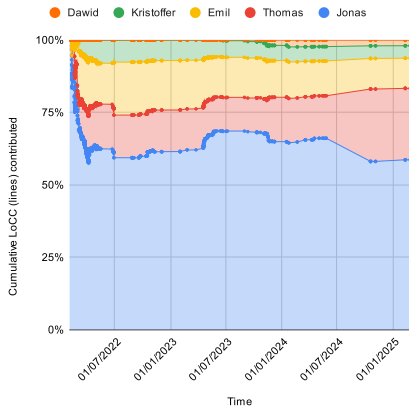# RQ3: Robustness to structural changes



Figure 4: ΔCD (Difference in Compression Distance) Explained

☑ ΔCD is insensitive to **project structure** at commit granularity
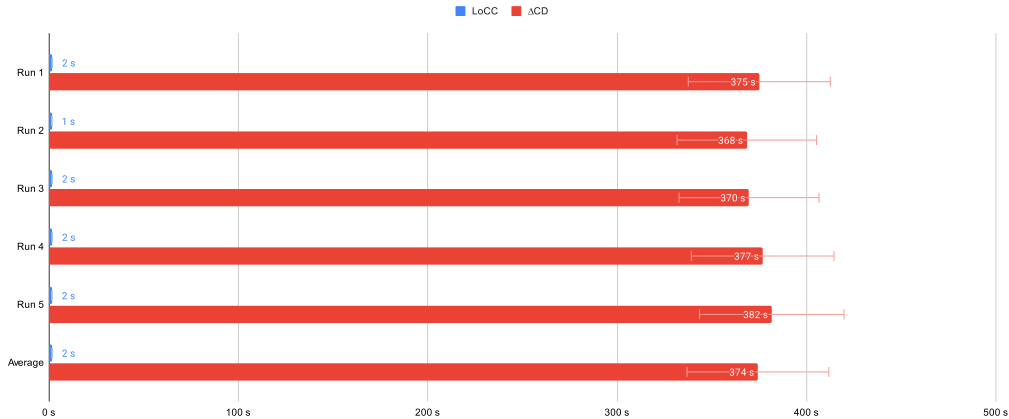
# RQ3: Survivorship Bias

- 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

# RQ3: Performance and Scalability

LoCC vs ΔCD for commitizen-tools/commitizen (github, 1977 commits)
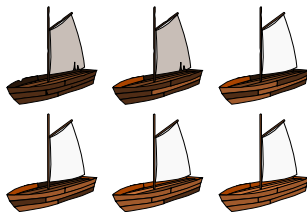
# Future work

Performance and scalability

Generalize findings

Robustness to formatting changes etc.

# 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