# Repository Analysis of Open-Source and Scientific Software Development Projects

Kanika Sood[1], Boyana Norris[1], Anshu Dubey[2], Lois McInnes[2]
University of Oregon[1], Argonne National Laboratory[2]

February 25, 2019

**MS2: Scientific Software: Practices, Concerns, and Solution Strategies**

SIAM Conference on
Computational Science
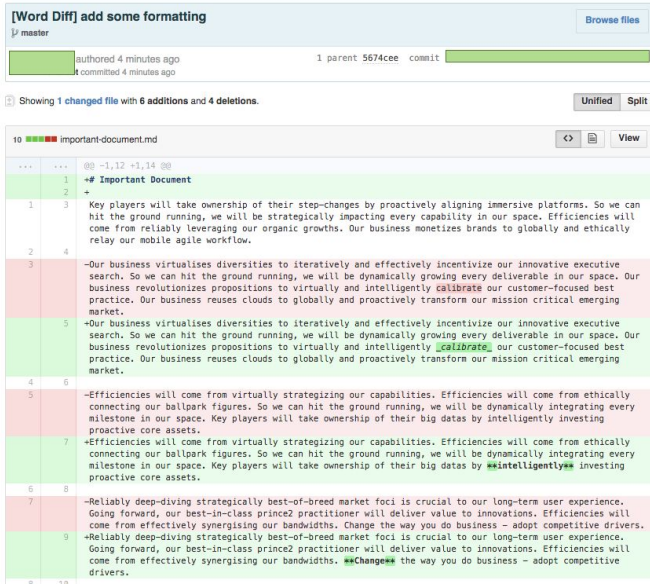and Engineering

# Introduction

- Scientific software is rapidly growing in capabilities, accuracy, performance.

- Developer productivity has received less attention than app. performance/ publications.

- We propose new time-dependent metrics that can help quantify team productivity.

- The metrics can be used to better understand the trends of software development workflows and provide objective measurements of productivity.

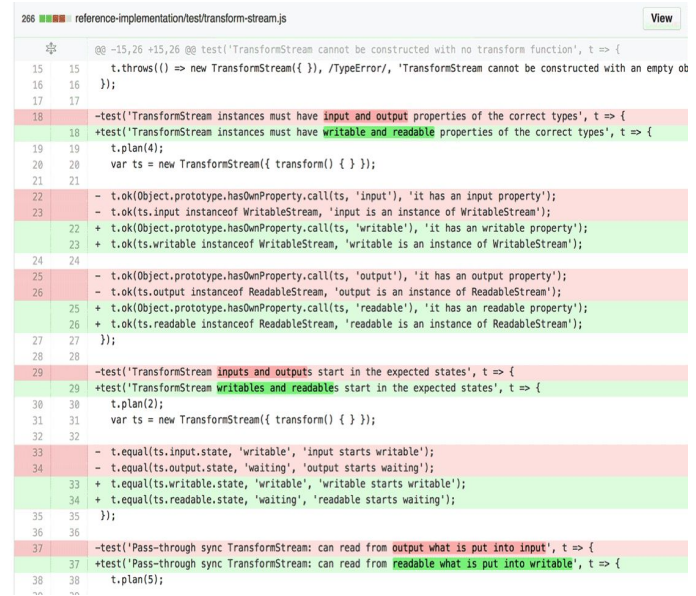- We demonstrate our approach on several HPC software projects.

**Disclaimer:** The goal of this research is to explore new software metrics that can provide insight into productivity more effectively than existing metrics. These (or any) metrics provide partial perspectives but cannot capture a complete view of the complexities of scientific software projects.

# *Why is quantifying productivity hard?*

## Standard metrics





**Current metrics like NLOC are insufficient for quantifying software productivity**
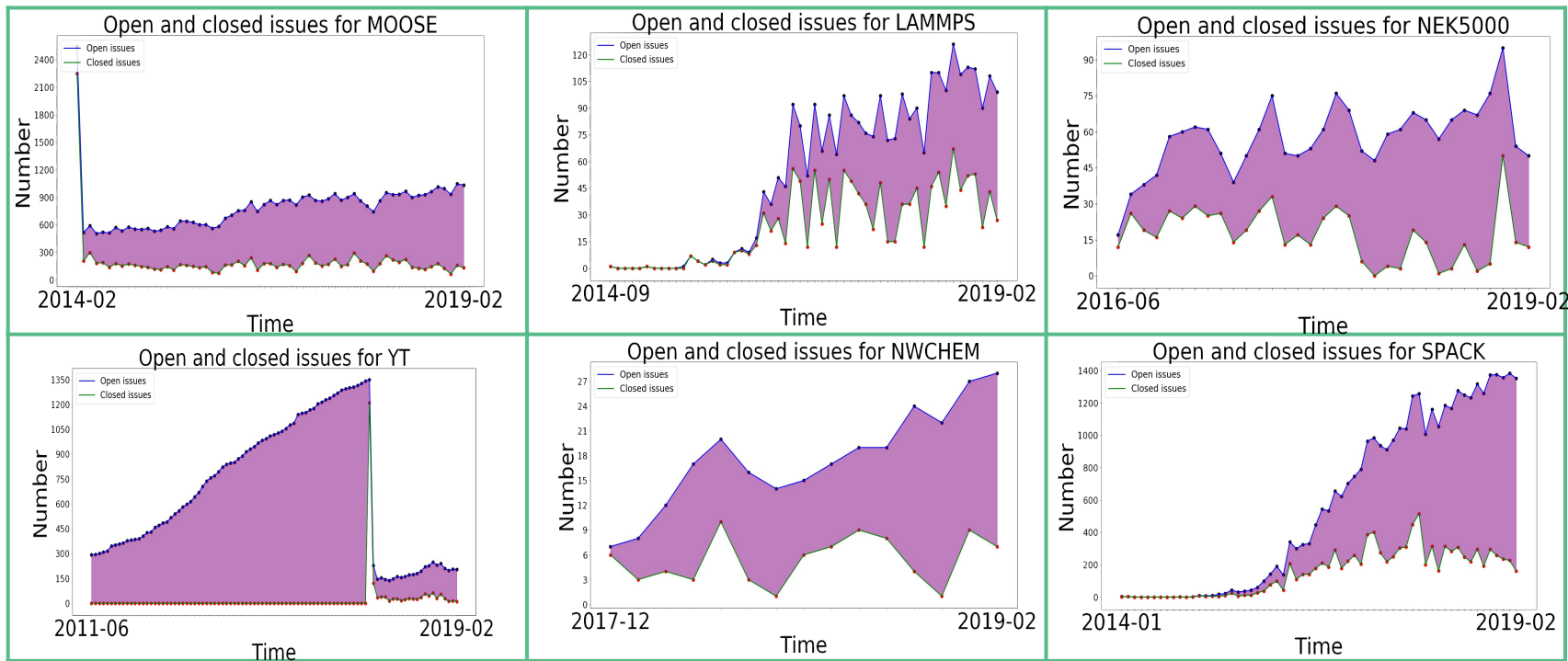
# Metrics

- Commit frequency
- Total additions, deletions
- NLOC

- Issue requests
- Issue categorization
- Project change -NLOC/Files/complexity
- Topics of discussions
- Developer activity
- Developer contributions
- Project reliance



Developer weekly data for number of deletions and additions for Developer 1

Developer weekly data for number of commits for Developer 1

# How engaged are the user and developer communities?

## Issues

# What are hot topics (based on issues)?

# What topics dominate changes/discussion?

## Pull requests

# *What topics dominate changes/discussion?*

## Pull requests

| Topic | No. of files changed |
|-------|----------------------|
| Rc 350 | 1971 |
| Rc 360 | 1403 |

| Topic | No. of comments |
|-------|-----------------|
| Optimizable determinants | 91 |
| Multislater-Jastrow Orbital Optimization code | 84 |



Discussion from pull requests for QMCPACK

# Code analysis

**Cyclomatic Complexity (CCN)**

- Quantitative measure of the number of independent paths in the code
- Statements (S1, Sn): nodes, control paths from S1-> Sn: edges
- Computed for each function
- Smaller value : better
- Tool: Lizard
- ================================================

| NLOC | CCN | PARAM | length | location |
|------|-----|-------|--------|----------|
| 118 | 2 | 50 | 18 | func@6-13@./main.c |

# How is code complexity changing over time?

## CCN

# *How is the project size changing over time?*

## Number of lines of code (NLOC)

# Project evolution: QDPXX

# Project evolution: GROMACS



Color coding: Avg. CCN

Color coding: Time

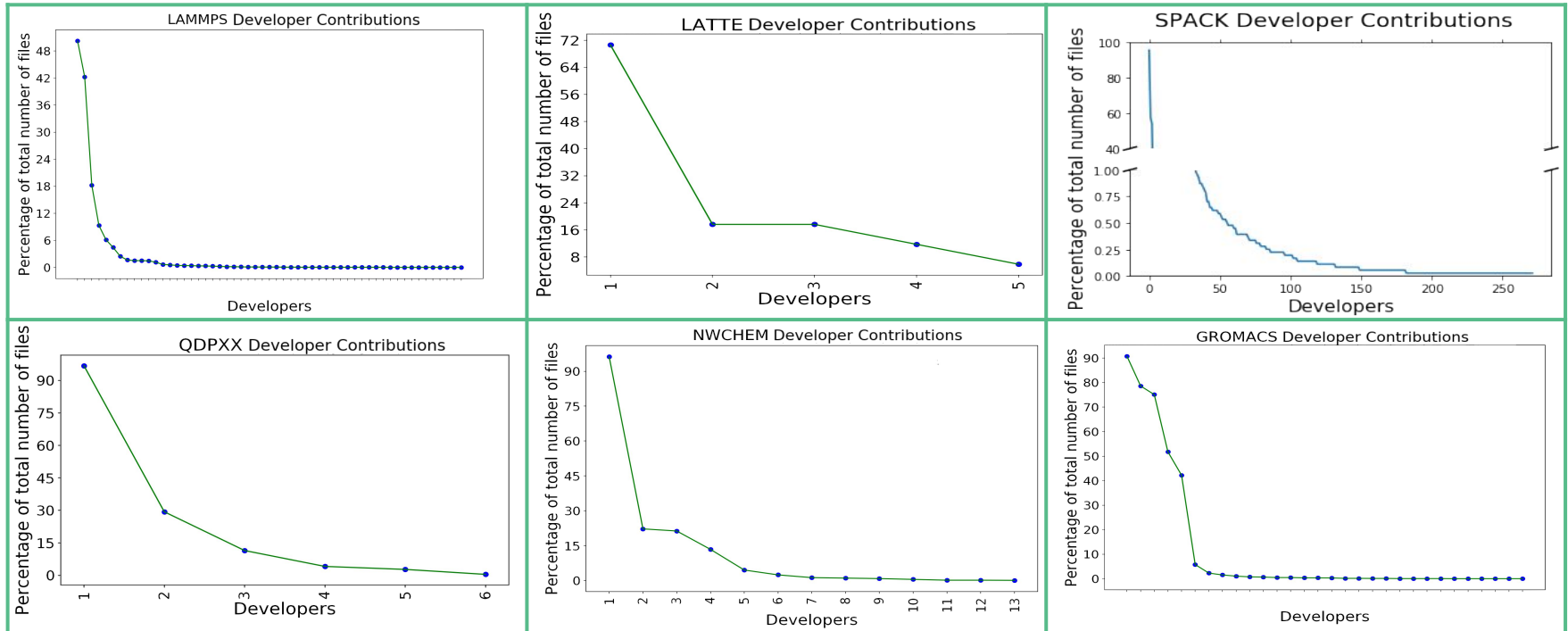# *How active is the developer community?*

## Developer activity based on pull requests

# *What is the project reliance on individual developers?*

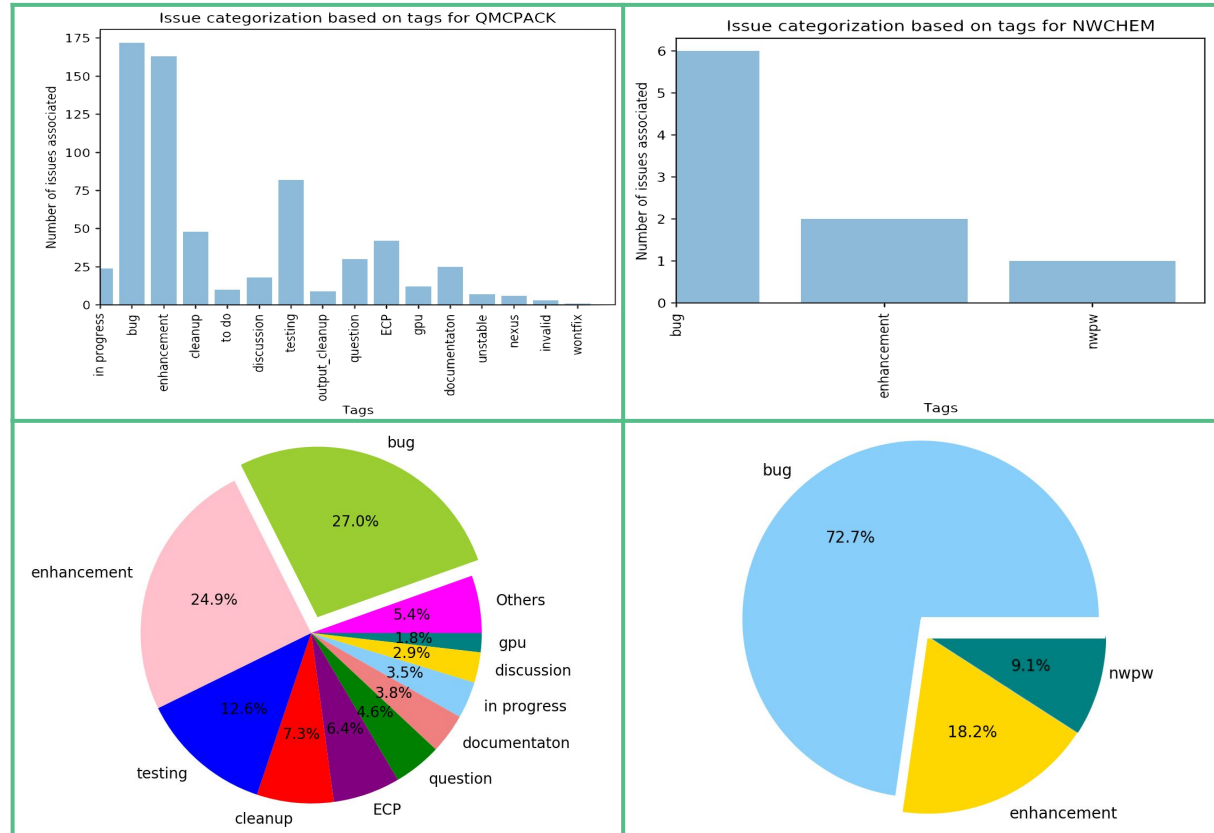## Percentage of total number of code files

# Conclusion

- Compute metrics to better understand software development practices
- Fine-grain analysis for individual files, individual developers, patterns in software trends, and project reliance
- Study the impact of code size and changes in code complexity over project lifetimes.
  - Discover opportunities to reduce cost and increase scientific output
  - Guide future project planning

**Disclaimer:** The goal of this research is to explore new software metrics that can provide insight into productivity more effectively than existing metrics. These (or any) metrics provide partial perspectives but cannot capture a complete view of the complexities of scientific software projects.

Thank you

# *What are hot topics (based on issues)?*

**Issue categories**

# *What are hot topics (based on issues)?*

# Code complexity

```
public void one(){
  if(true) {
    while(false) {
        two();
    }
    else {
    for(int i=0;i<10;i++) {
        two();
        }
    }
  }
}
```
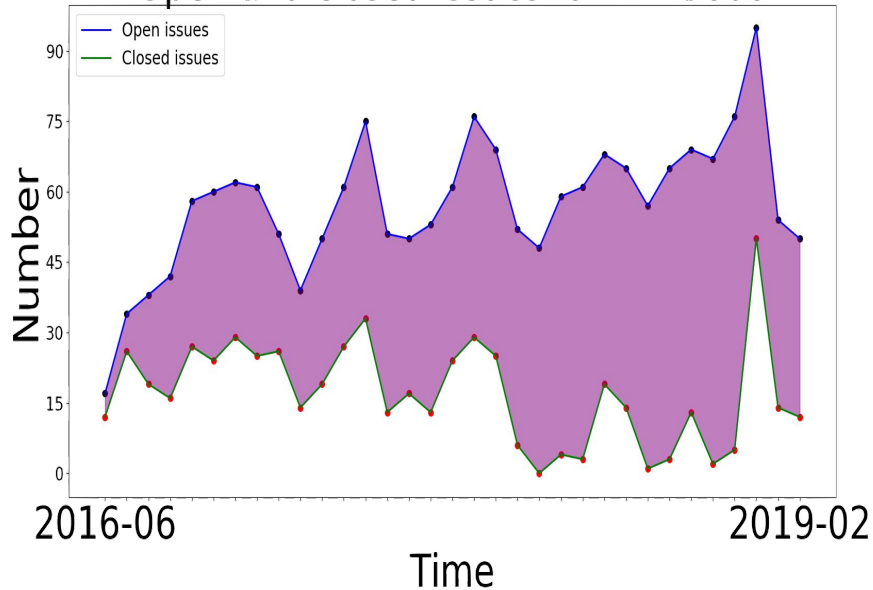
CCN = e - n + 2p
8 - 6 + 2(1) = 4

# *How engaged is the user and developer community?*

**Issue requests**

# *How is code complexity changing over time?*

## **CCN**

# Project evolution: LAMMPS



25

# Project evolution: QMCPACK
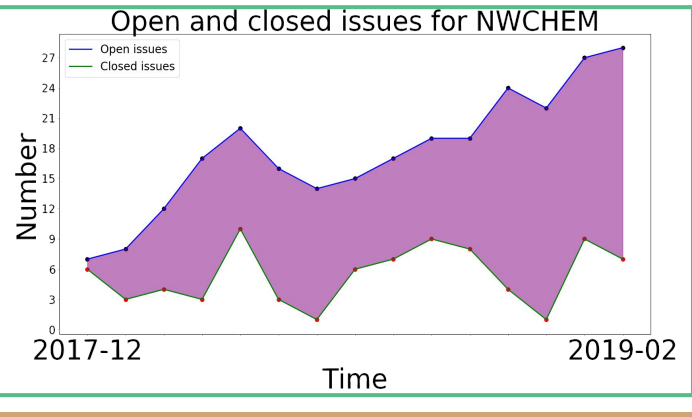
# Project evolution: NWCHEM

# Project evolution: YT

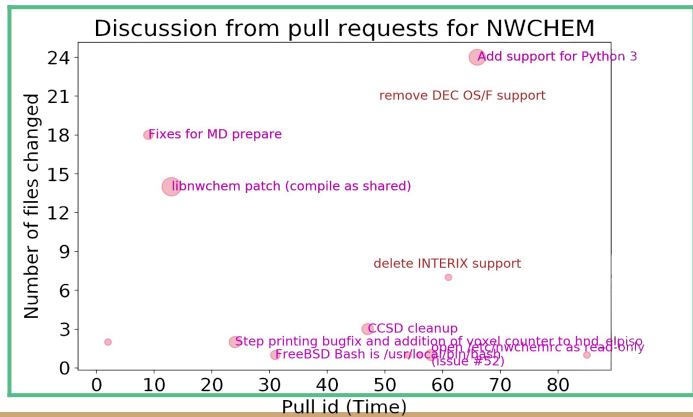# Project evolution: SPACK



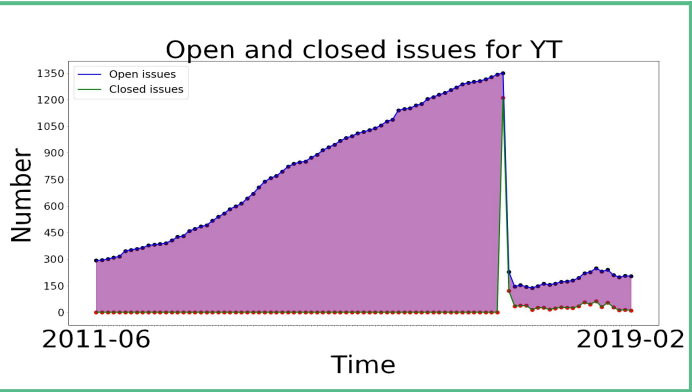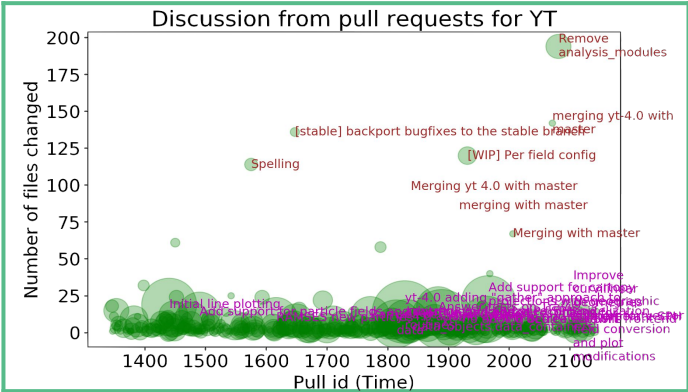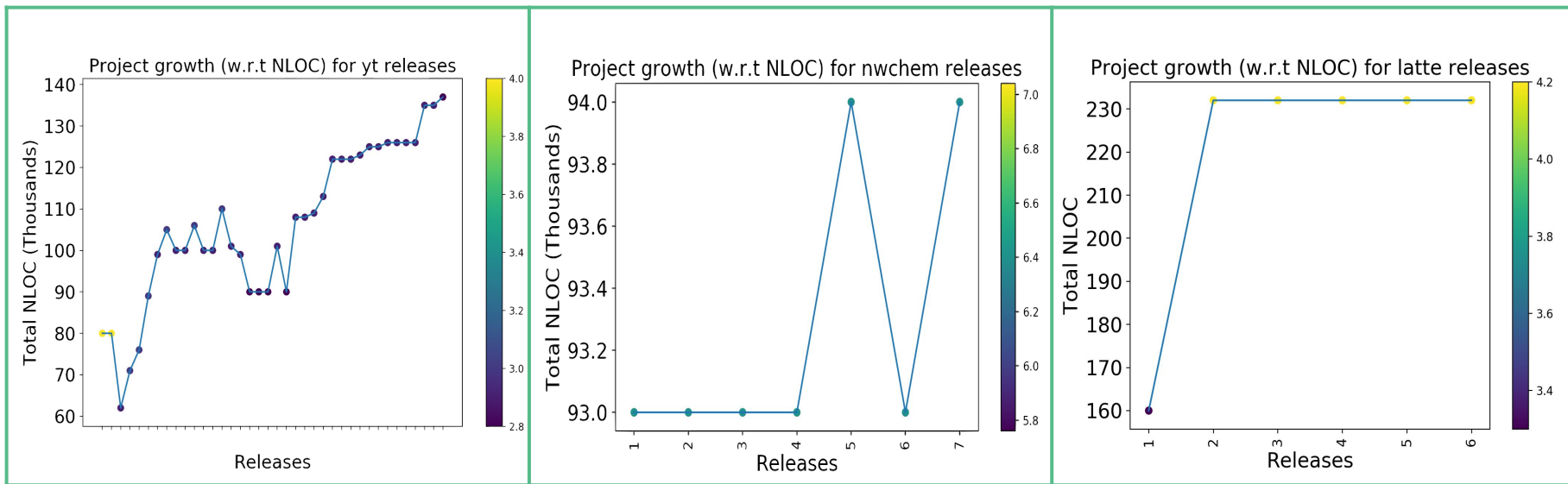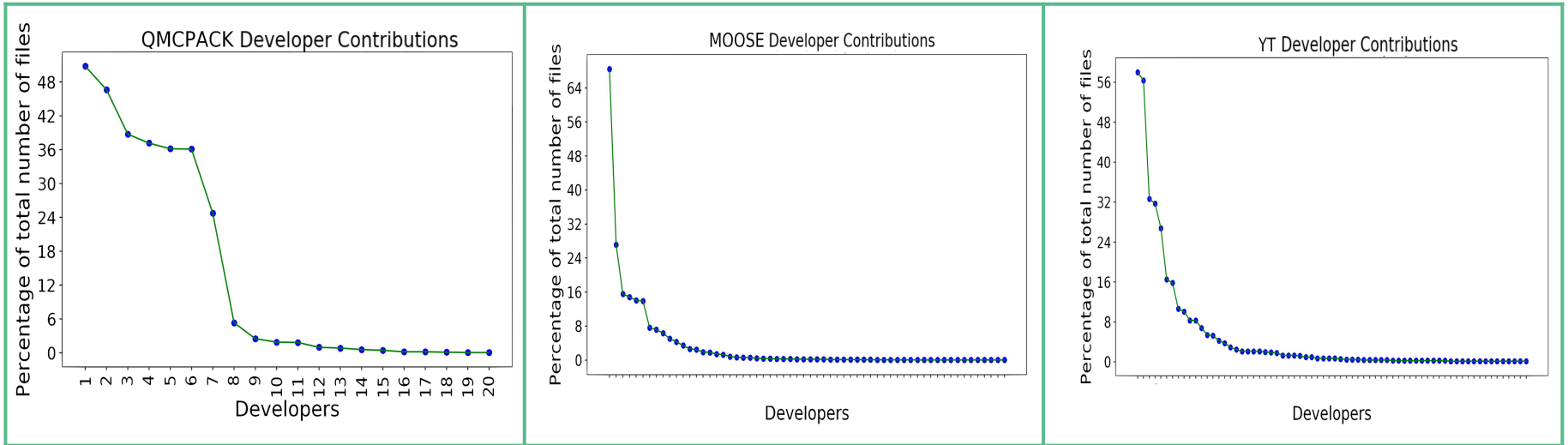Color coding: Avg. CCN

Color coding: Time

# *How is the project size changing over time?*

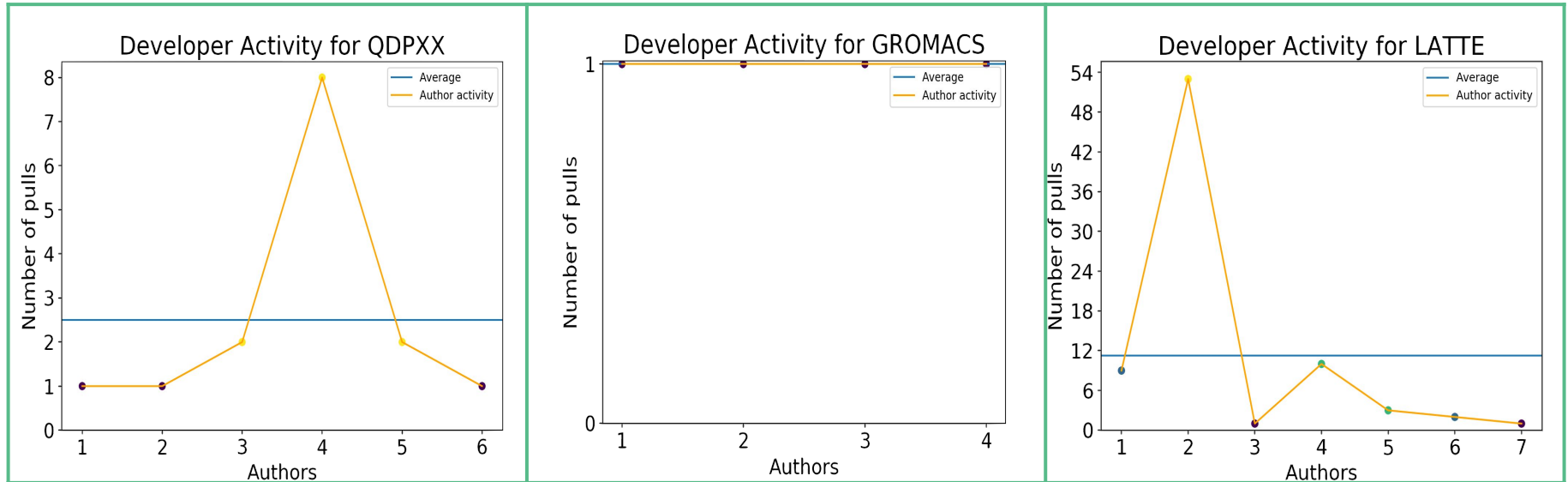## Number of lines of code (NLOC)

# *What is the project reliance on individual developers?*

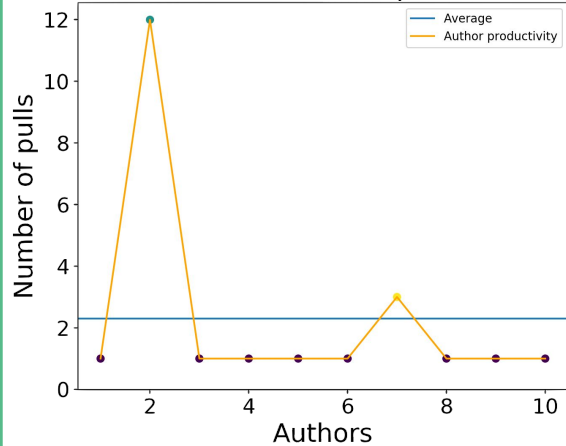## Percentage of total number of code files

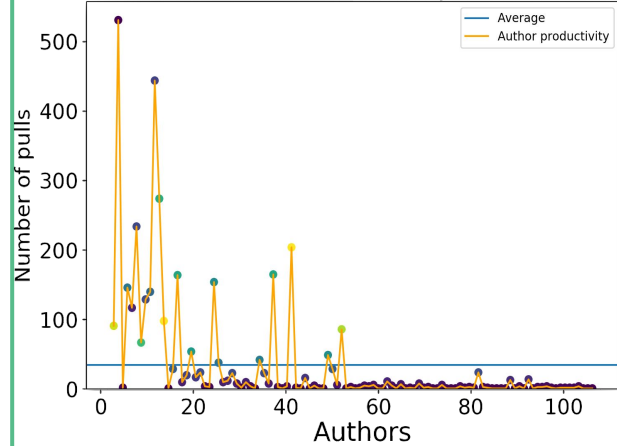# *How active is the developer community?*

## Developer activity based on pull requests

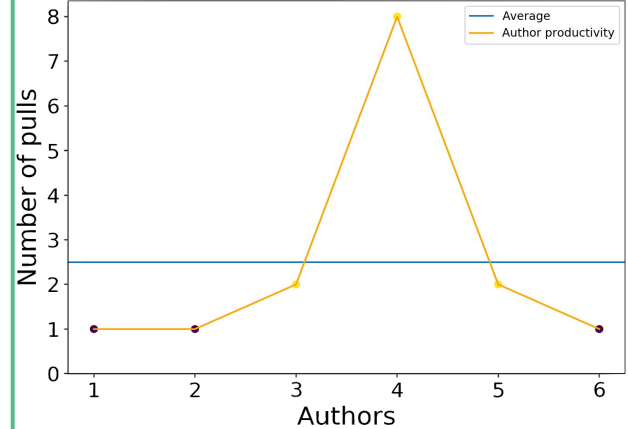# *How active is the developer community?*

## Developer activity based on pull requests

# *What topics dominate changes/discussion?*

**Pull requests**