# SOFTWARE REPOSITORIES AND THEIR USABILITY IN SOFTWARE PROCESS RECONSTRUCTION
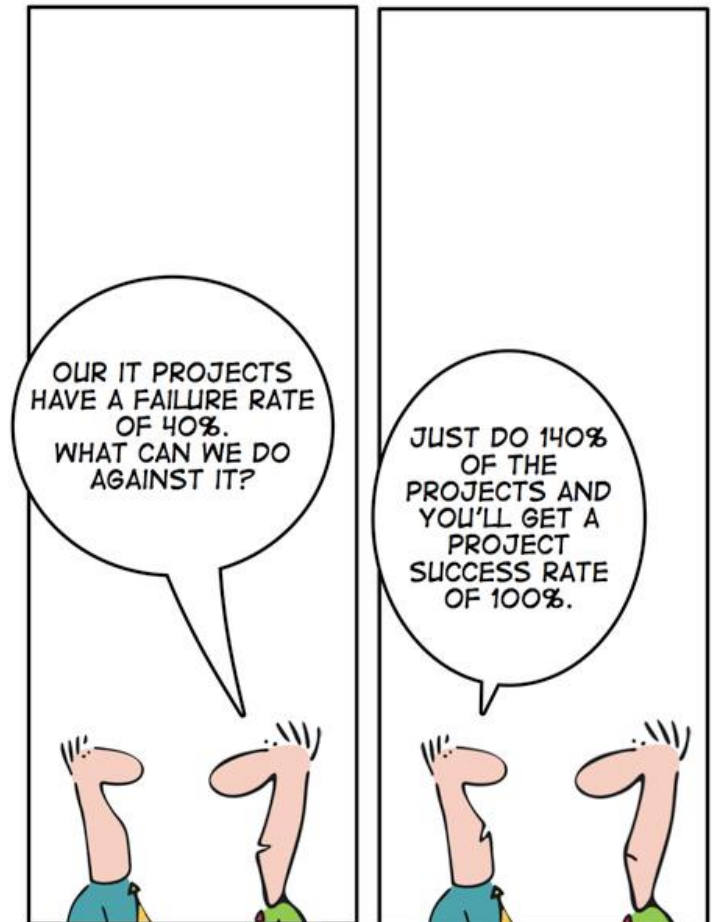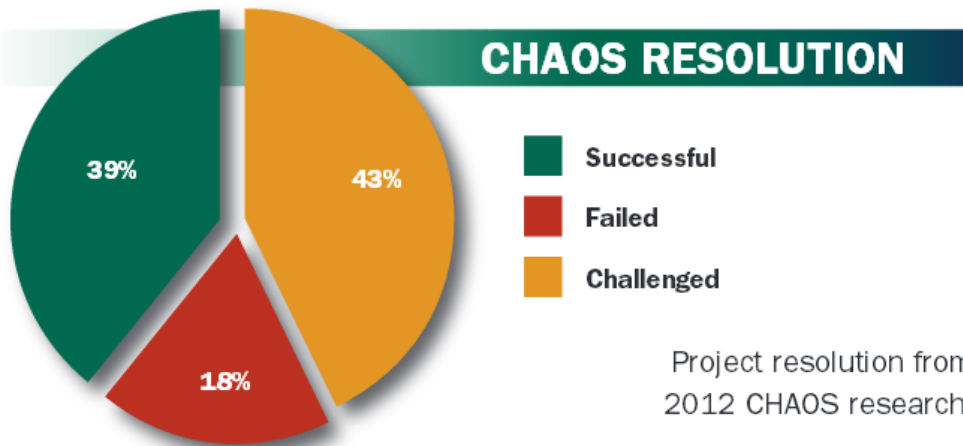
Marko Janković & Marko Bajec

# IT Project Performance

# Many reasons

- Social issues
- Technology challenges
- …
- The lack of discipline:
  - Many companies do not have any SDM in place
  - Prescribed SDMs not followed
  - Lack of motivation

ISD is about implementing IT into a human enterprise!
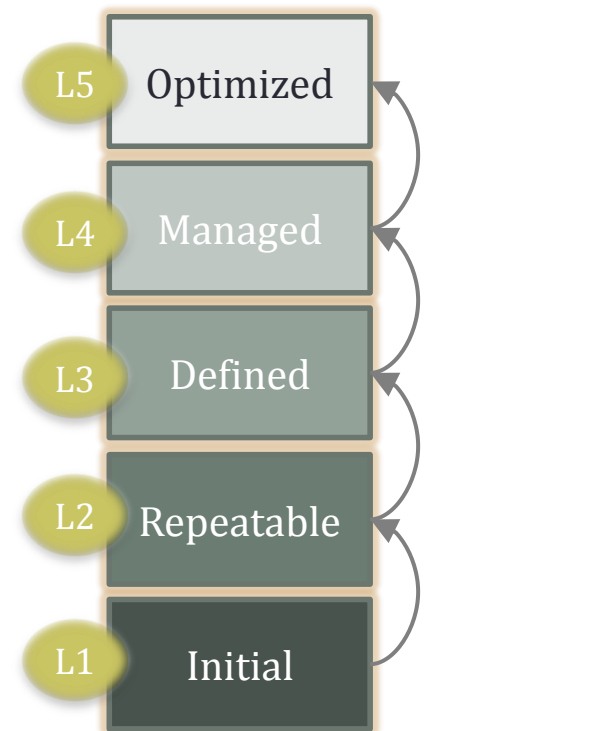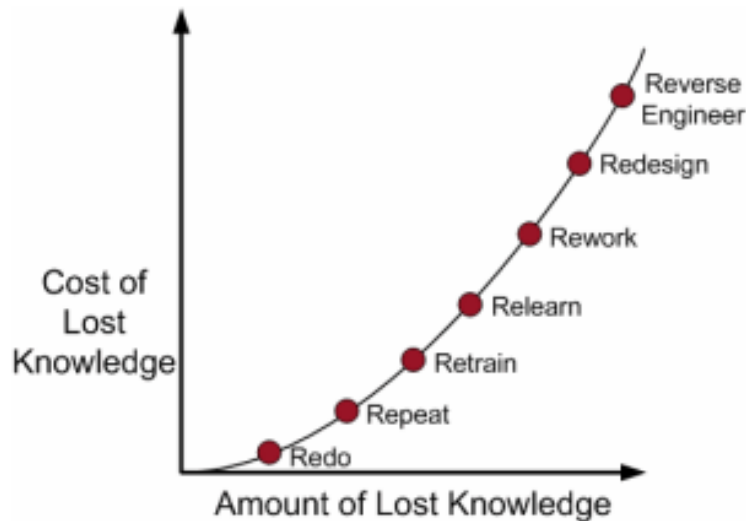
**The Agile Manifesto**

| **Individuals and interactions** | over | Processes and Tools |
| **Working Product** | over | Comprehensive Documentation |
| **Customer Collaboration** | over | Contract Negotiation |
| **Responding to change** | over | Following a plan |

*That is, while there is value in the items on the right, we value the items on the left more.*
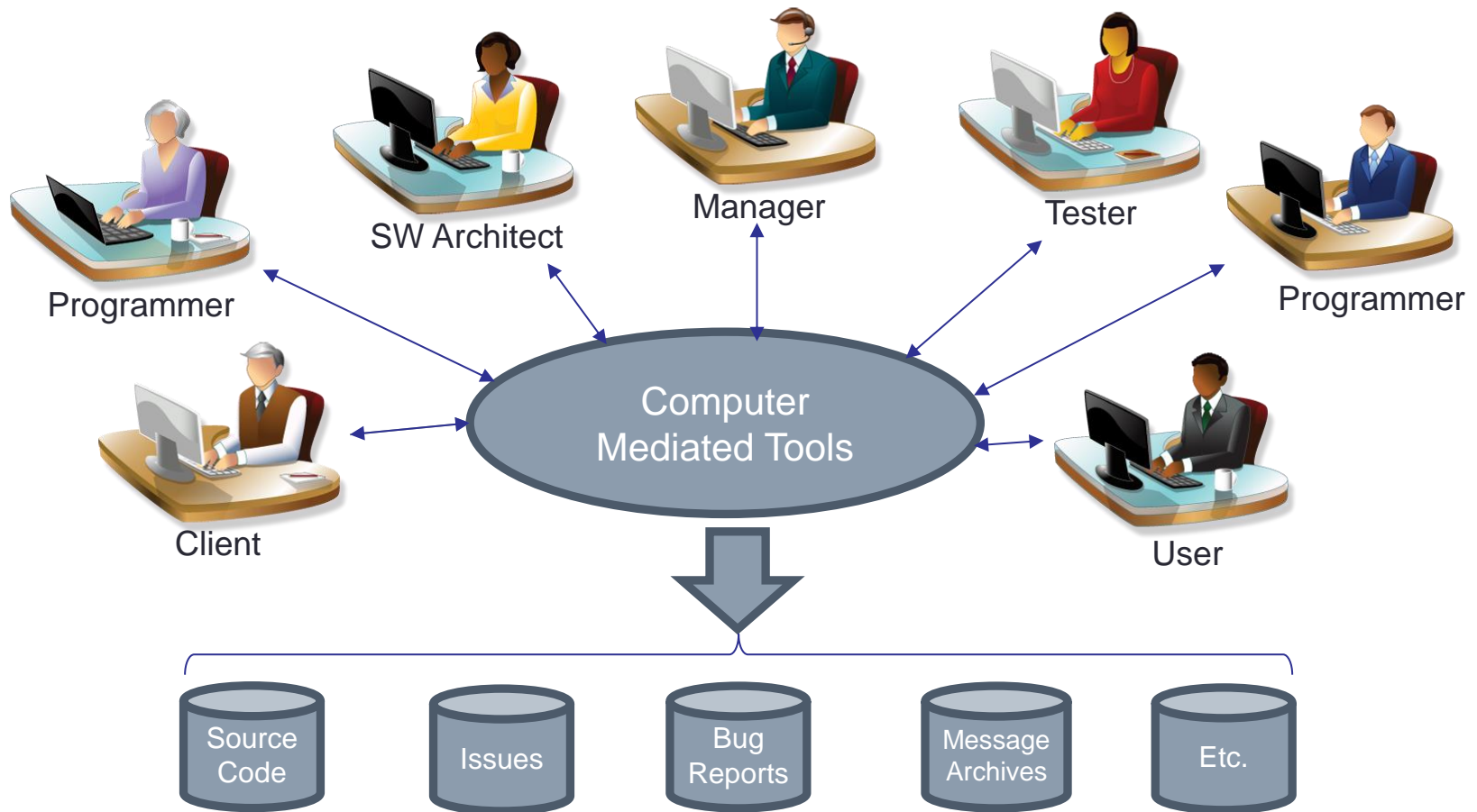
www.agilemanifesto.org

# Problems and Limitations

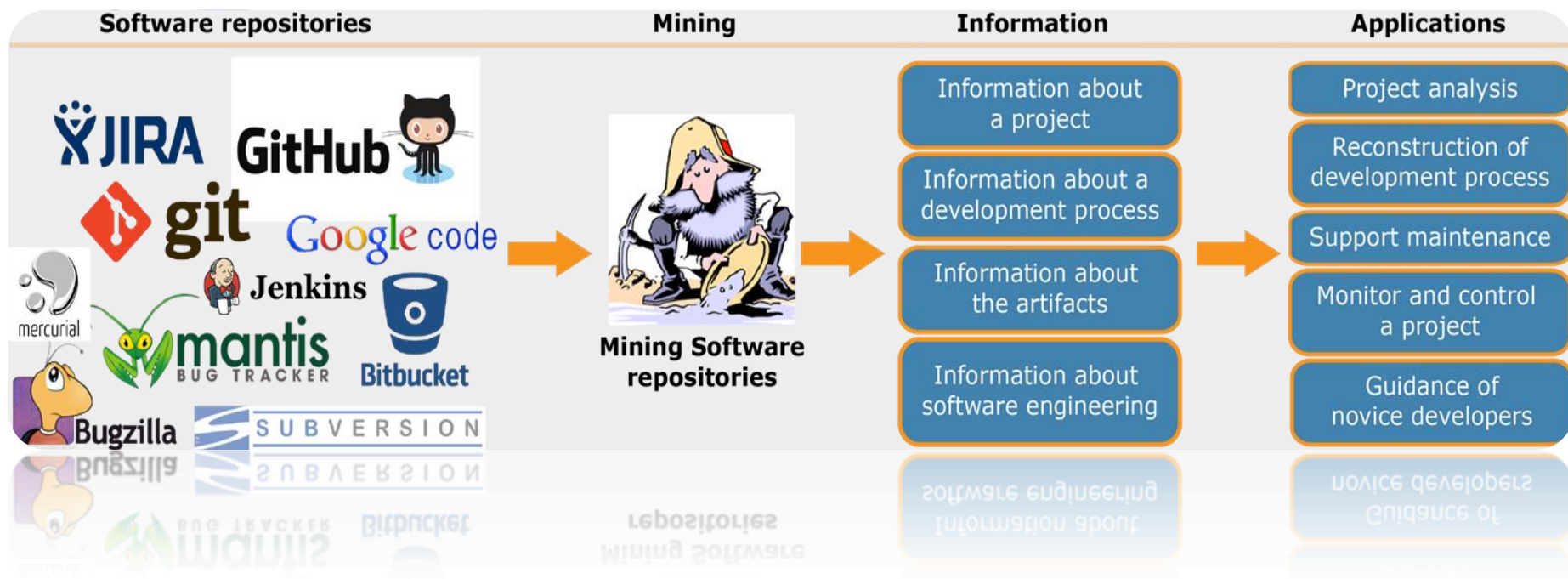- Risk for knowledge loss…
- Repeating mistakes…
- Reinventing the wheel…





Maturity levels of the CMM

# Software Repositories



Programmer

SW Architect

Manager

Tester

Programmer

Client

Computer
Mediated Tools

User

Source
Code

Issues
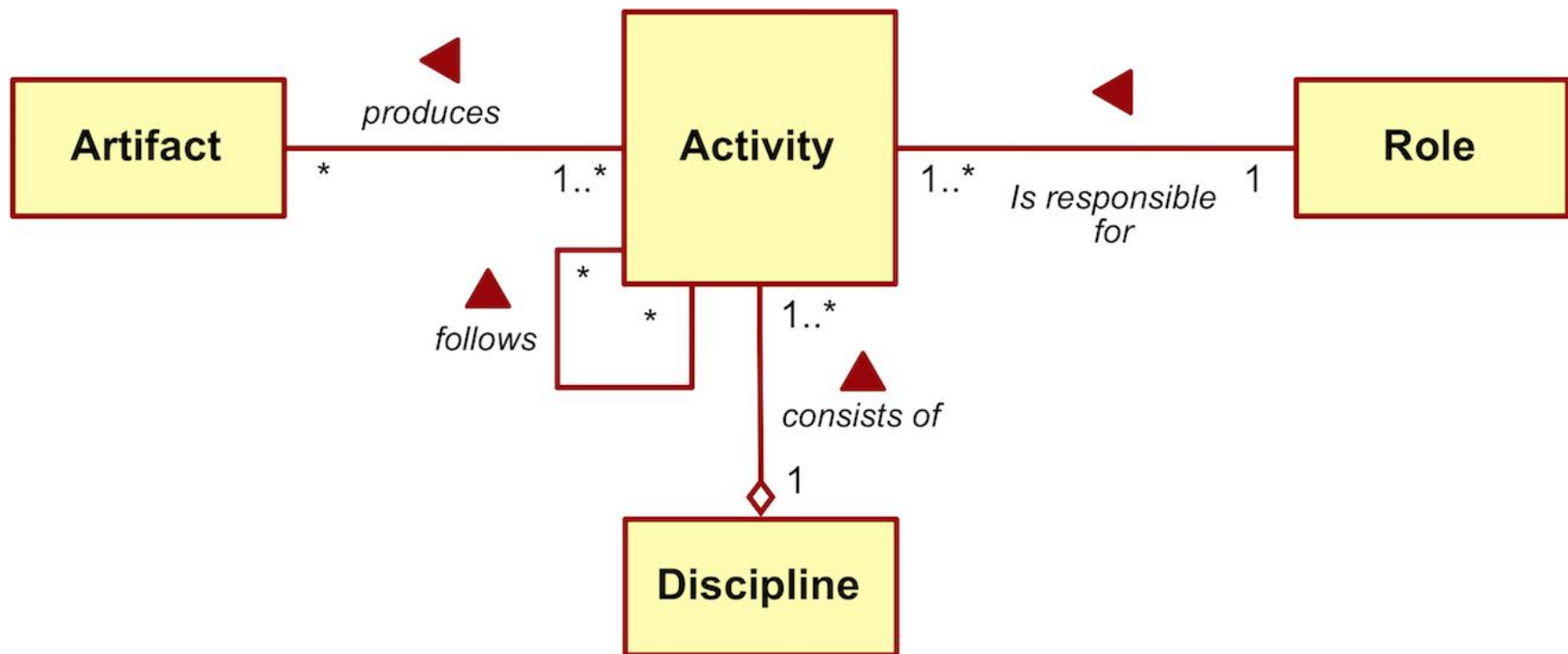
Bug
Reports

Message
Archives

Etc.

Based on Marco Aurélio Gerosa, Mining Sociotechnical Information From Software Repositories, University of São Paulo, Brazil
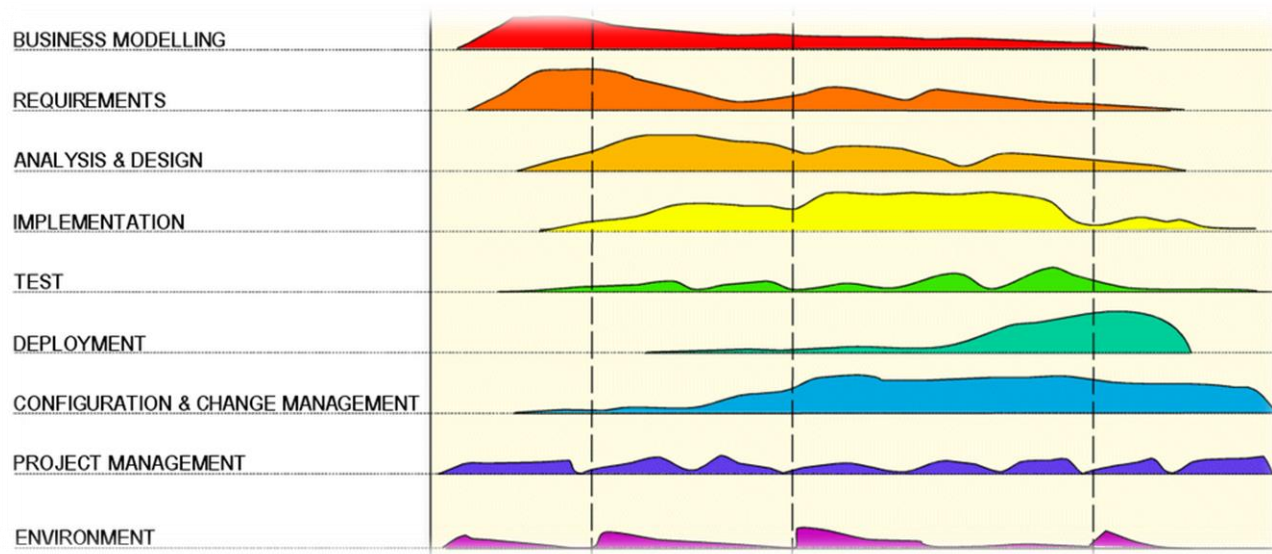
# Possible Applications

# Elements for Reconstruction
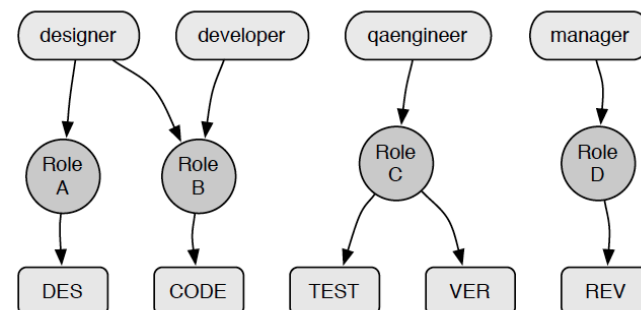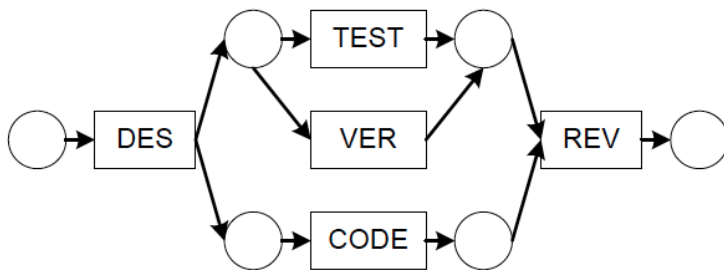
# Software process recovery

- Employs different semi-supervised techniques to recover UP diagram.
- Illustrates how the relative emphasis of different disciplines changes over the course of the project.



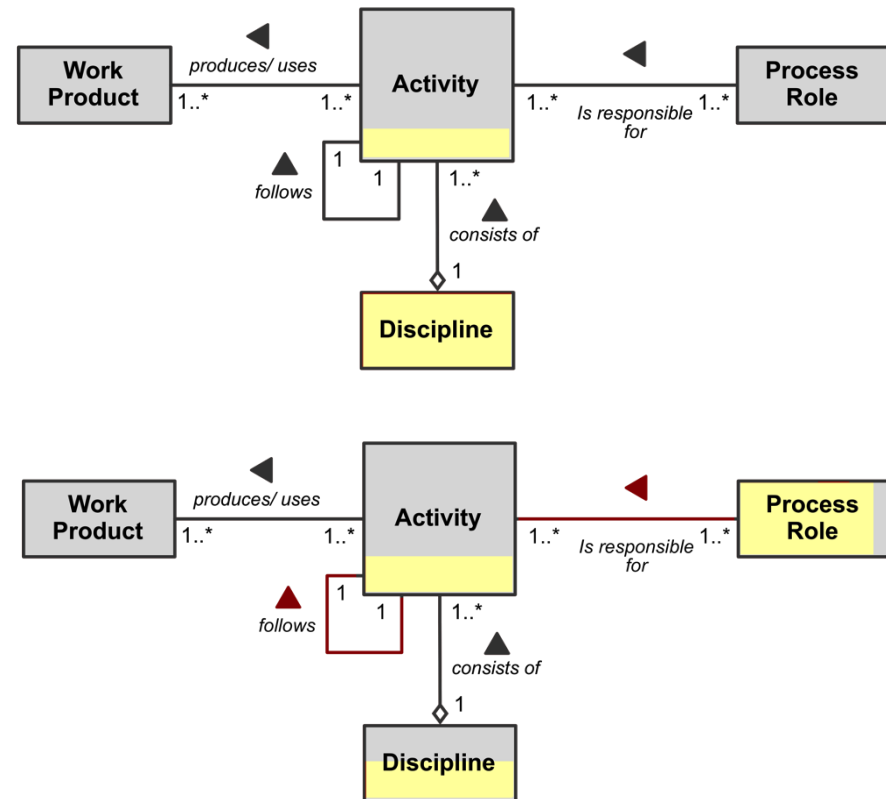A. Hindle, Software process recovery, PhD thesis

# Software process mining

- Mainly apply techniques from process mining on the event log generated from software repositories.
  - document names mapped into abstract names…
  - e.g.: docs with "/src/" in the filepath and with an extension ".java" map to the activity "code"
- Focused on reconstruction of high-level elements (e.g. main activities/disciplines) and workflow mining…
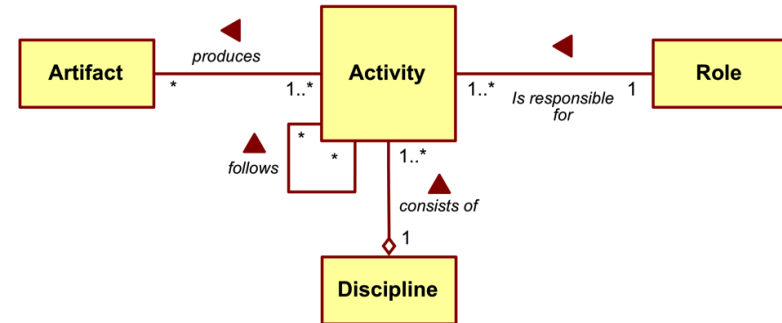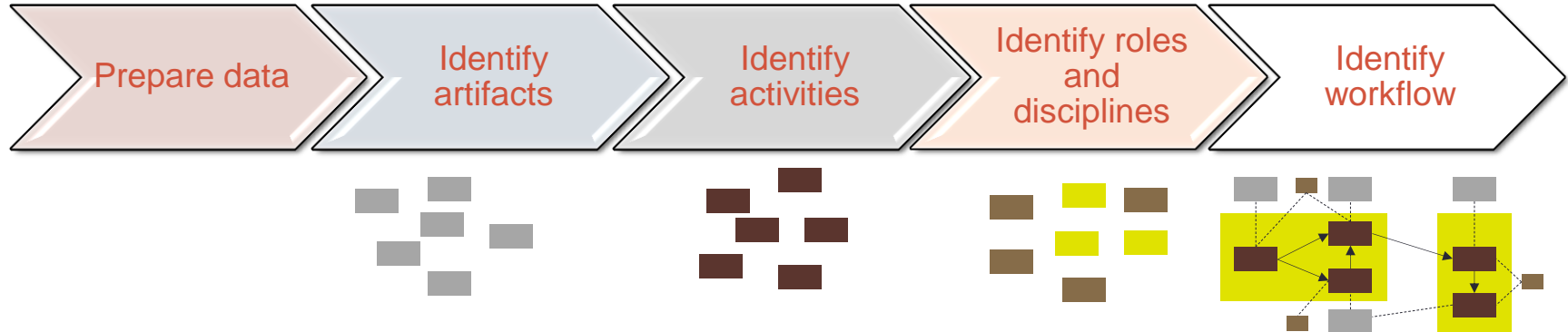- Data typically used from one repository only.

# Limitations

- Mining Software Repositories

- Software Process Mining

# Approach

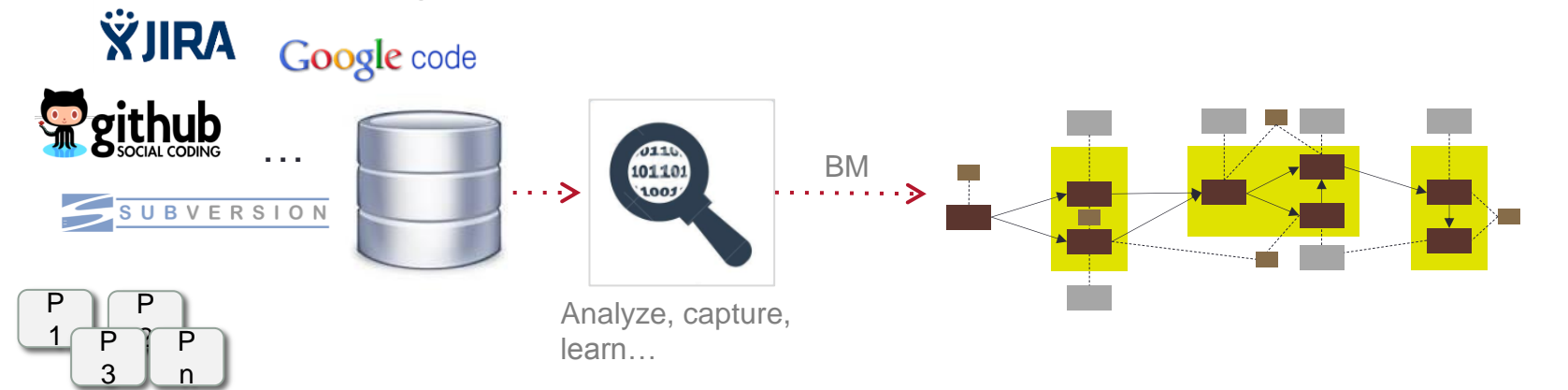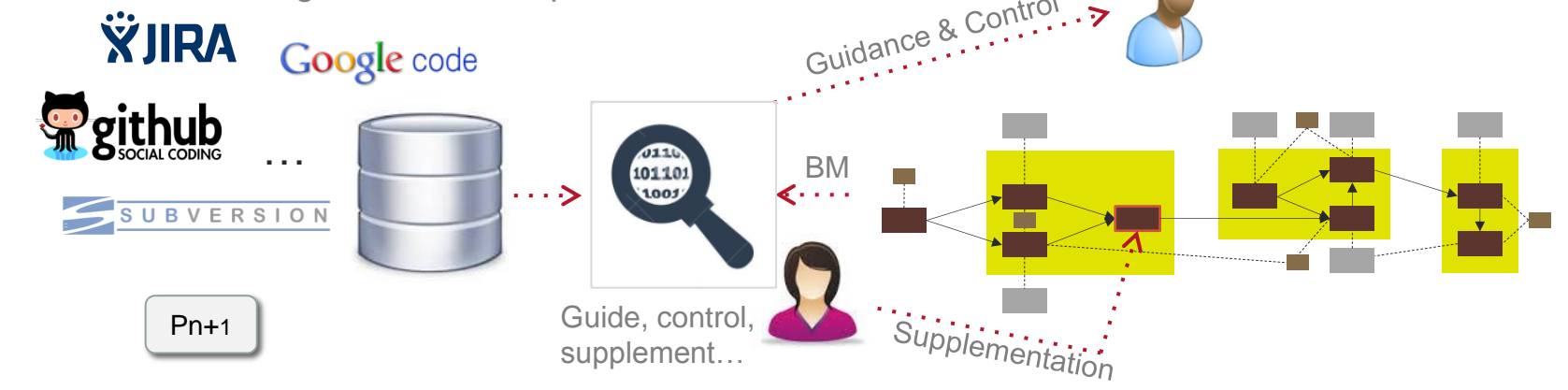# How it Works



Preparation: analysis of logs of past projects. Result: workflow of the base method

Analyze, capture, learn…

BM

Real-time control, guidance and improvement

Guidance & Control

BM

Guide, control, supplement…

Supplementation

Pn+1

# Data Preparation


Prepare data

- Gather data from repositories:
  - Revision control systems
  - Document system
  - Issue/Bug tracking system
  - Code review systems…
- Link users of different repositories → entity resolution
- Link tasks/issues with commits (e.g. based on commit messages…)

# Identification of artifacts

Identify artifacts

- Identification based on predefined ontology
  - Defines key elements (for each meta element of our interest)
  - Can be altered before or within the reconstruction process.

# Ontology

Based on
Agile Unified Process

Identification based on
keyword matching



Process role

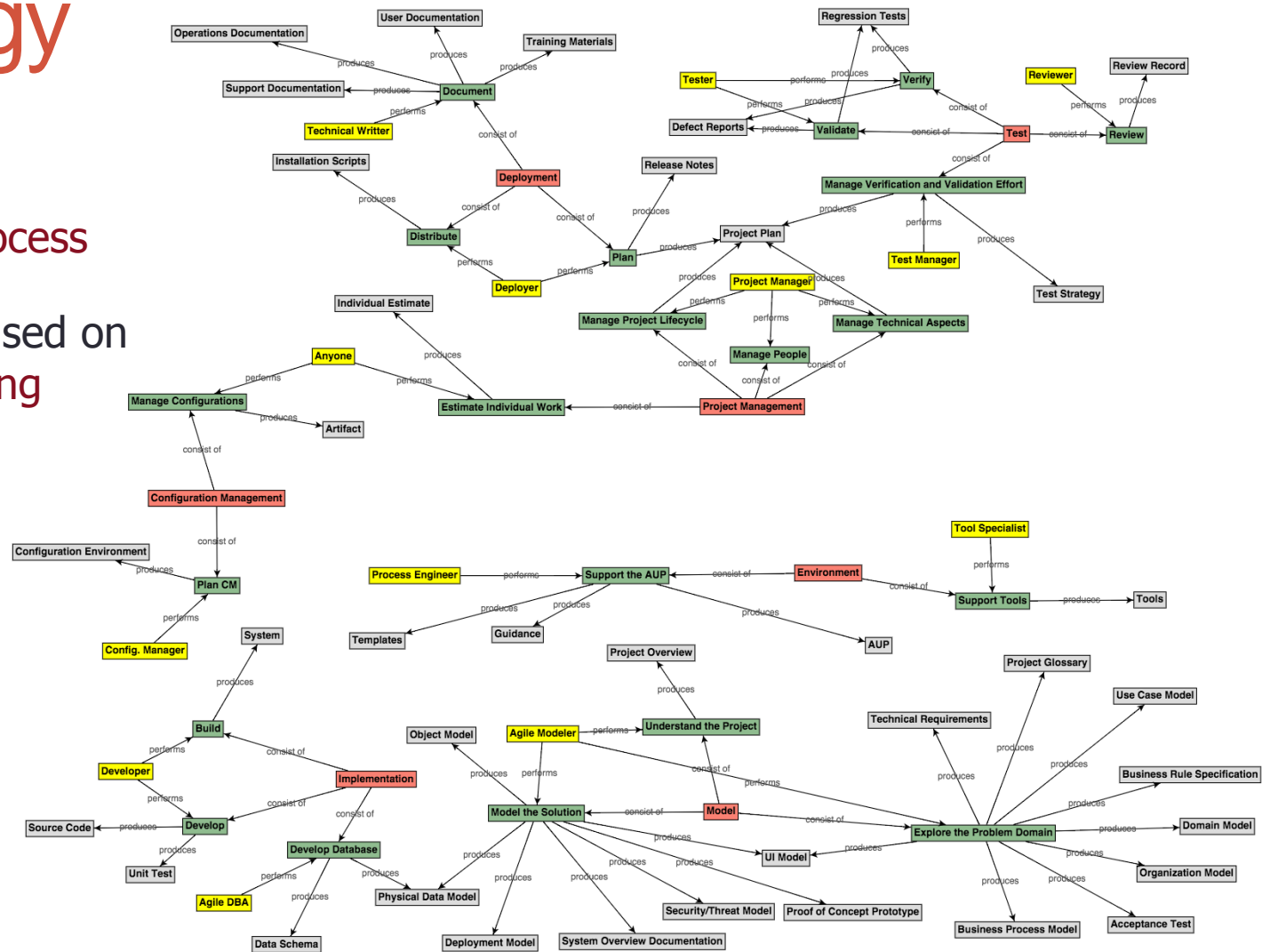Activity

Work product

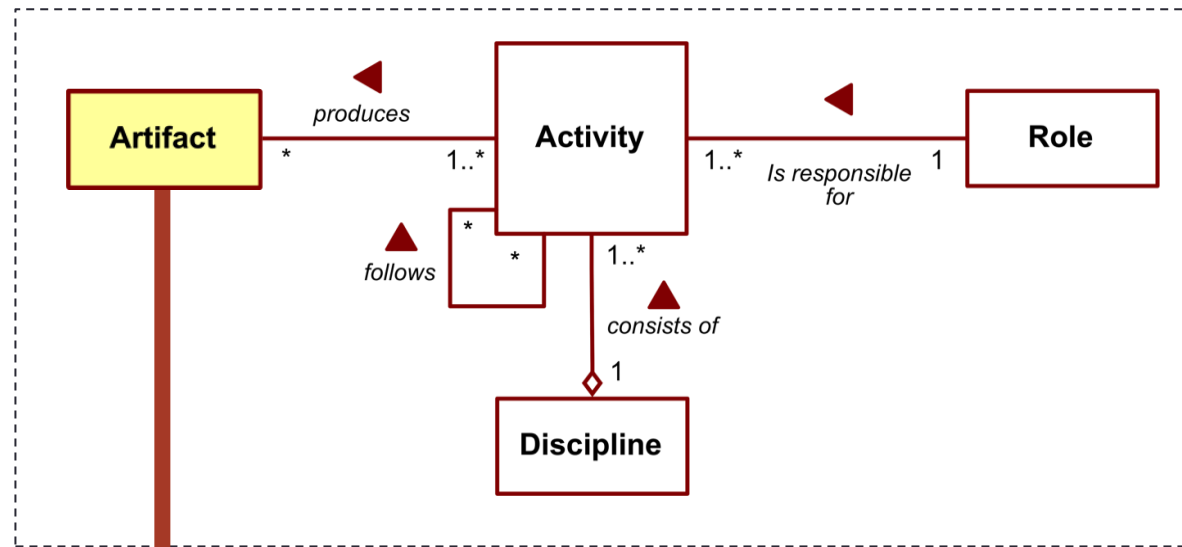Discipline

# Connecting files with artifacts

If low classification confidence then ask user
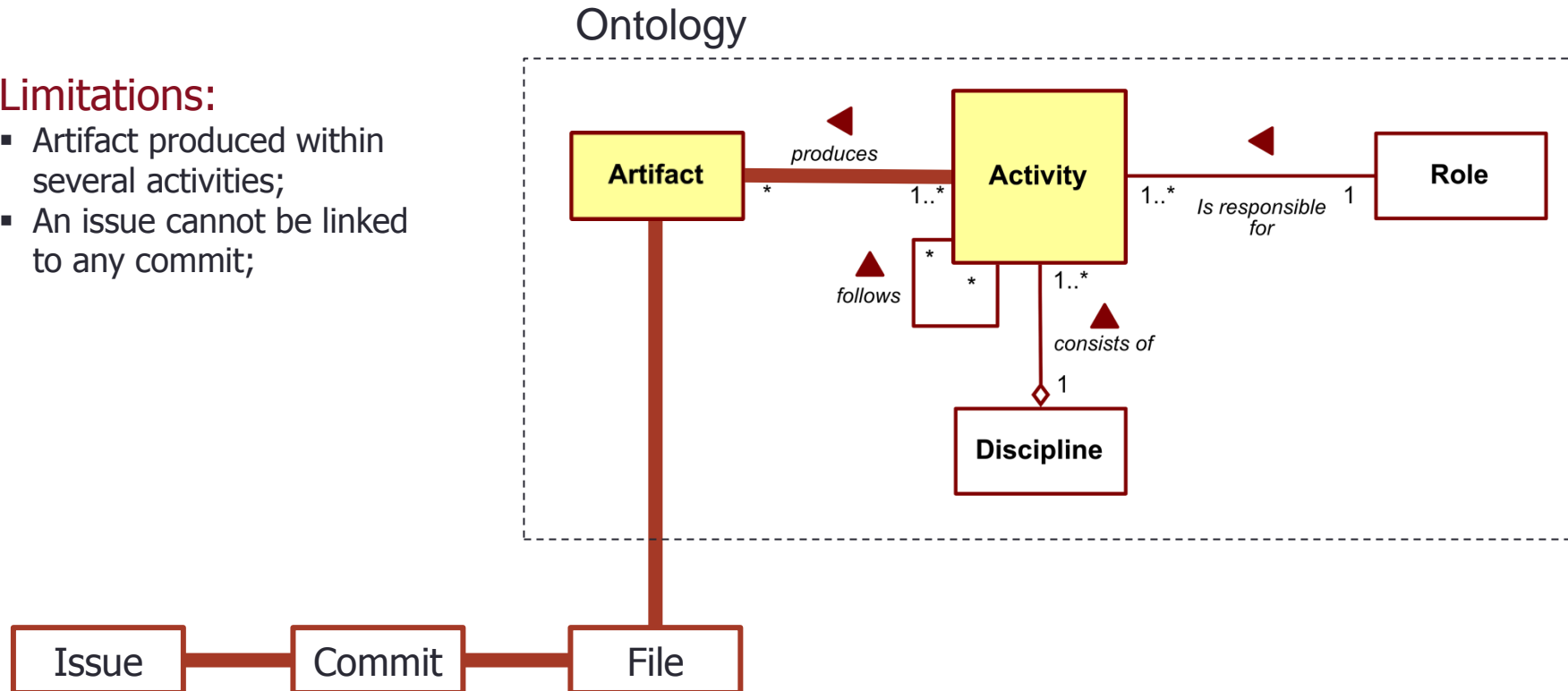
Ontology

# Identifying activities

Identify activities

## Limitations:

- Artifact produced within several activities;
- An issue cannot be linked to any commit;

Ontology

# Identifying roles and disciplines

Identify roles and disciplines

# Identifying flow of activities

Identify workflow

**Steps:**
- For each issue check the time when it was active (in progress → resolved).
- Draw issues on a timeline.
- For each issue, starting from the older ones, check the connected activities.
- If same activity as in previous issue → continue else connect respective activities.

Ontology



workflow

Issue — Commit — File

JIRA

github
SOCIAL CODING

# Workflow visualization

# Prerequisites

- For our approach to work the following is assumed:
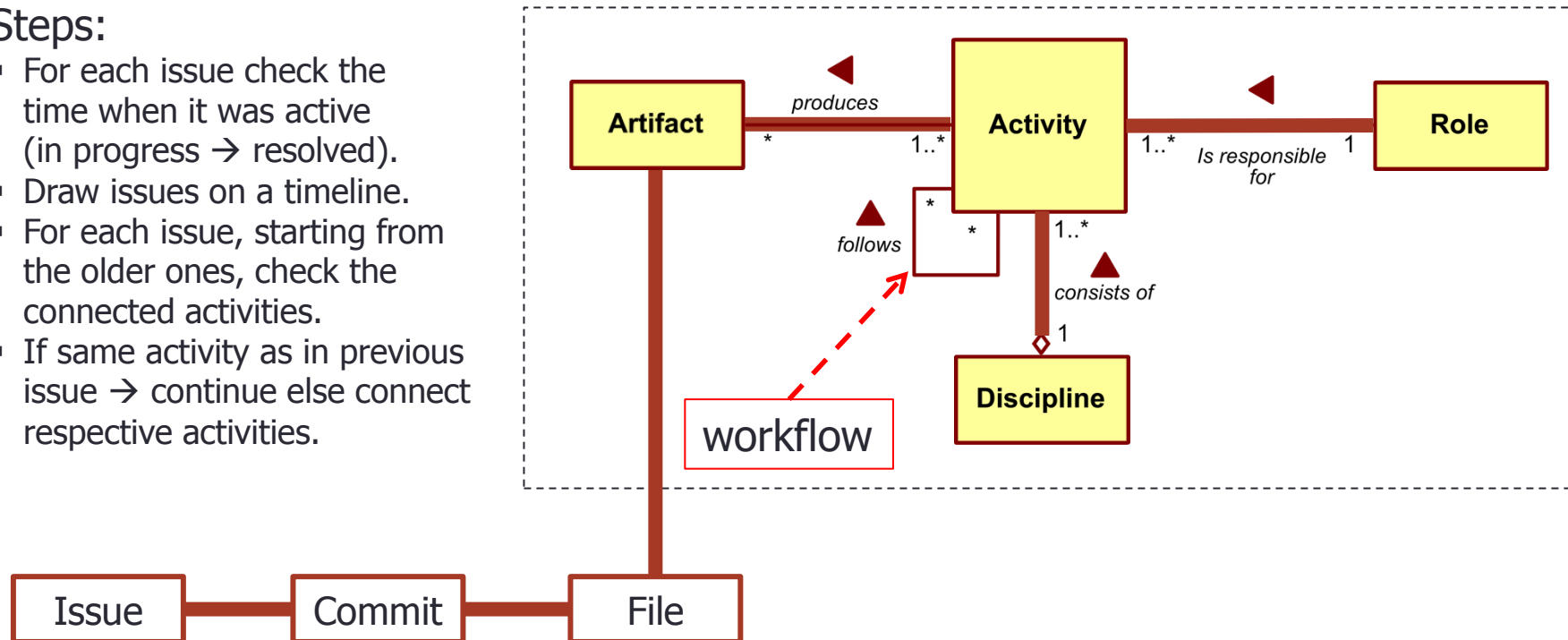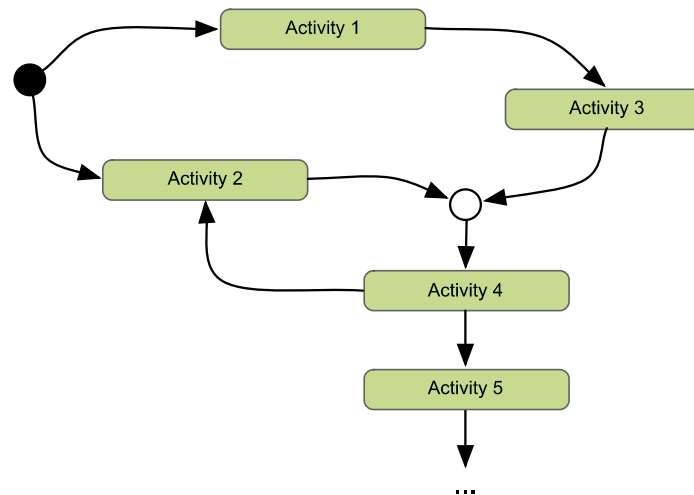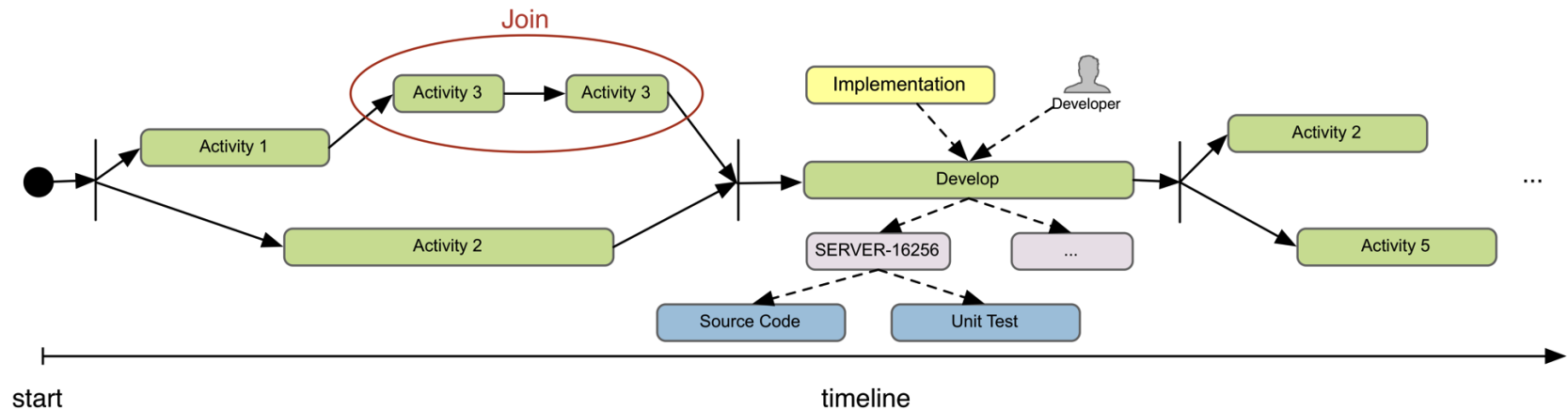
  - Commits are a consequence of creating or changing artifacts through tasks defined as issues.
  - The majority of commits and associated artifacts can be traced back to an exact issue that triggered the creation/change of those artifacts.
  - An issue is a small piece of work usually assigned to one developer only.
  - Issue statuses (opened, in progress, …, closed) and links among issues are strictly logged by developers.

# How limiting are the prerequisites…

- Five projects analyzed, three open source and two commercial.

| Open source project |
|---|
| **M o n g o   D B** |
| Started in Oct 2007<br>**15.292** issues in Jira<br>**28.374** commits in GitHub<br>Code Review in Rietveld |

| Open source project |
|---|
| **Spring Framework** |
| Started in 2003<br>**12.467** issues in Jira<br>**9.696** commits in GitHub |

| Open source project |
|---|
| **Hibernate ORM** |
| Started in 2003<br>**9.419** issues in Jira<br>**5.673** commits in GitHub |

| Commercial project |
|---|
| IS for insurance industry |
| Company with 250 emp.<br>Project started in 2007<br>Deployed to 15+ organiz.<br>**13.389** issues in Jira<br>**18.571** commits in SVN<br>Project mngm: SCRUM |

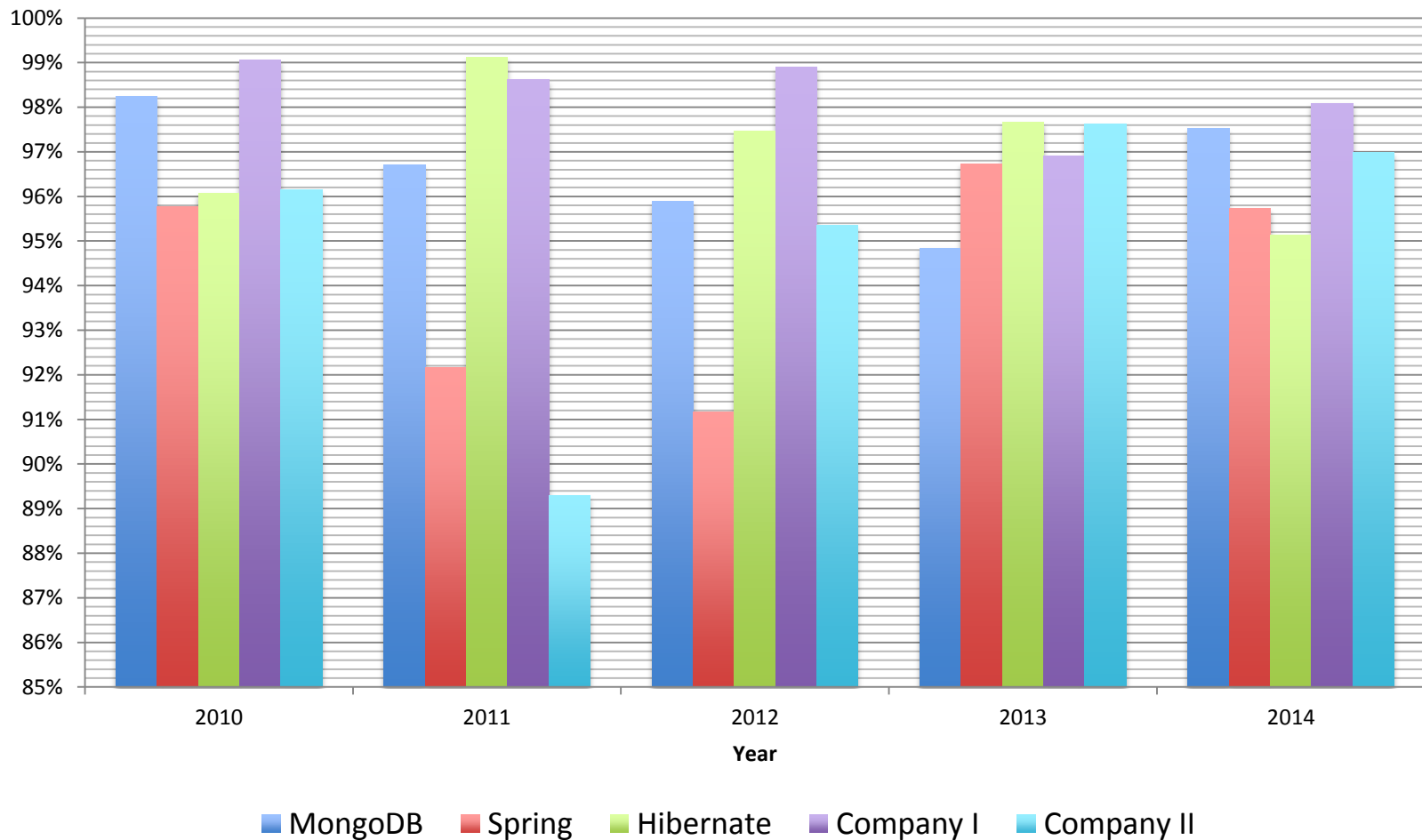| Commercial project |
|---|
| Billing for Utilities |
| Company with 30 emp.<br>Project started in 2008<br>**5.148** issues in Jira<br>**13.735** commits in SVN<br>Project mngm: SCRUM |

# Results
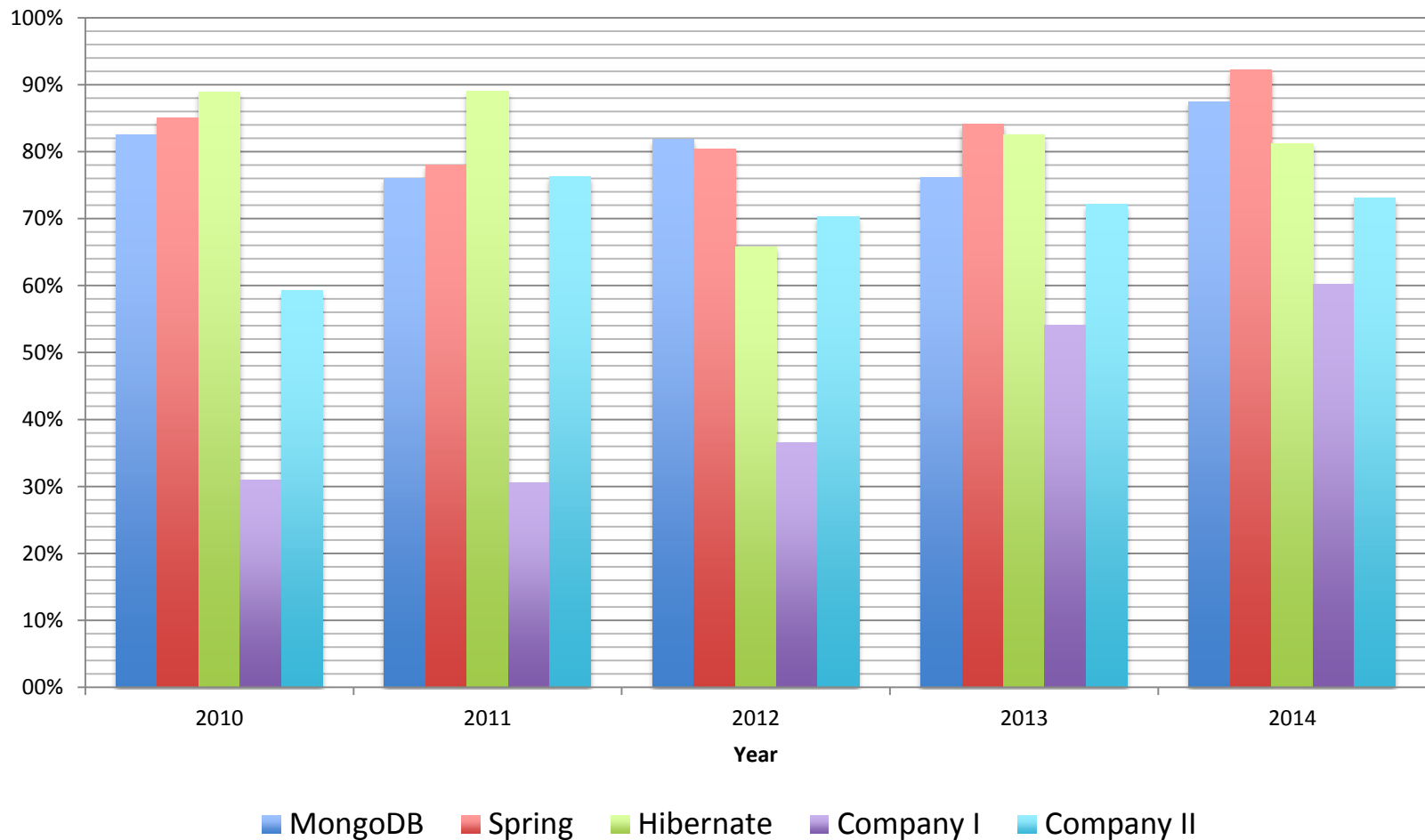


Percentage of commits that can be related to issues

# Results



Percentage of commits that can be related to exactly one issue

# Results



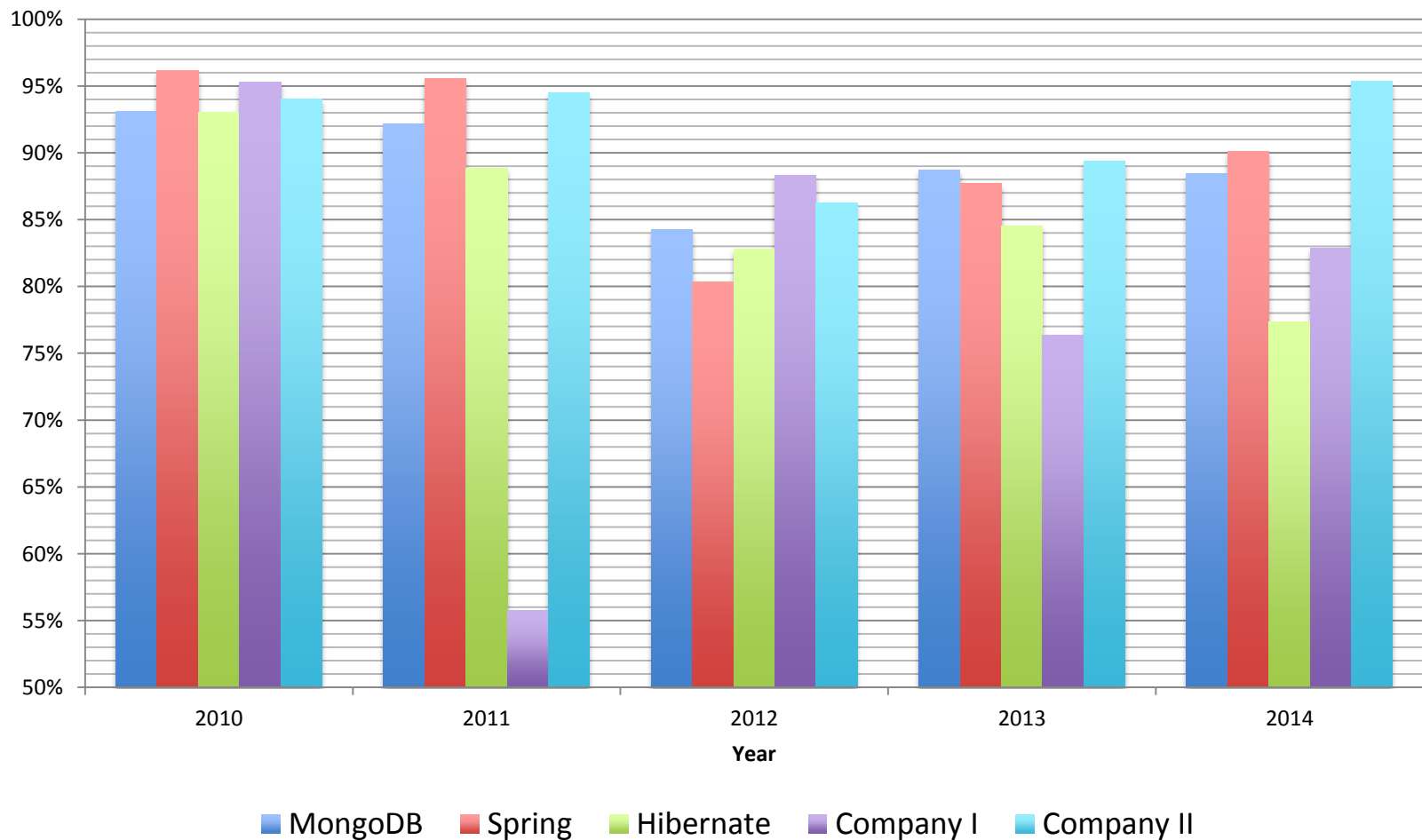Percentage of issues that can be related to a commit

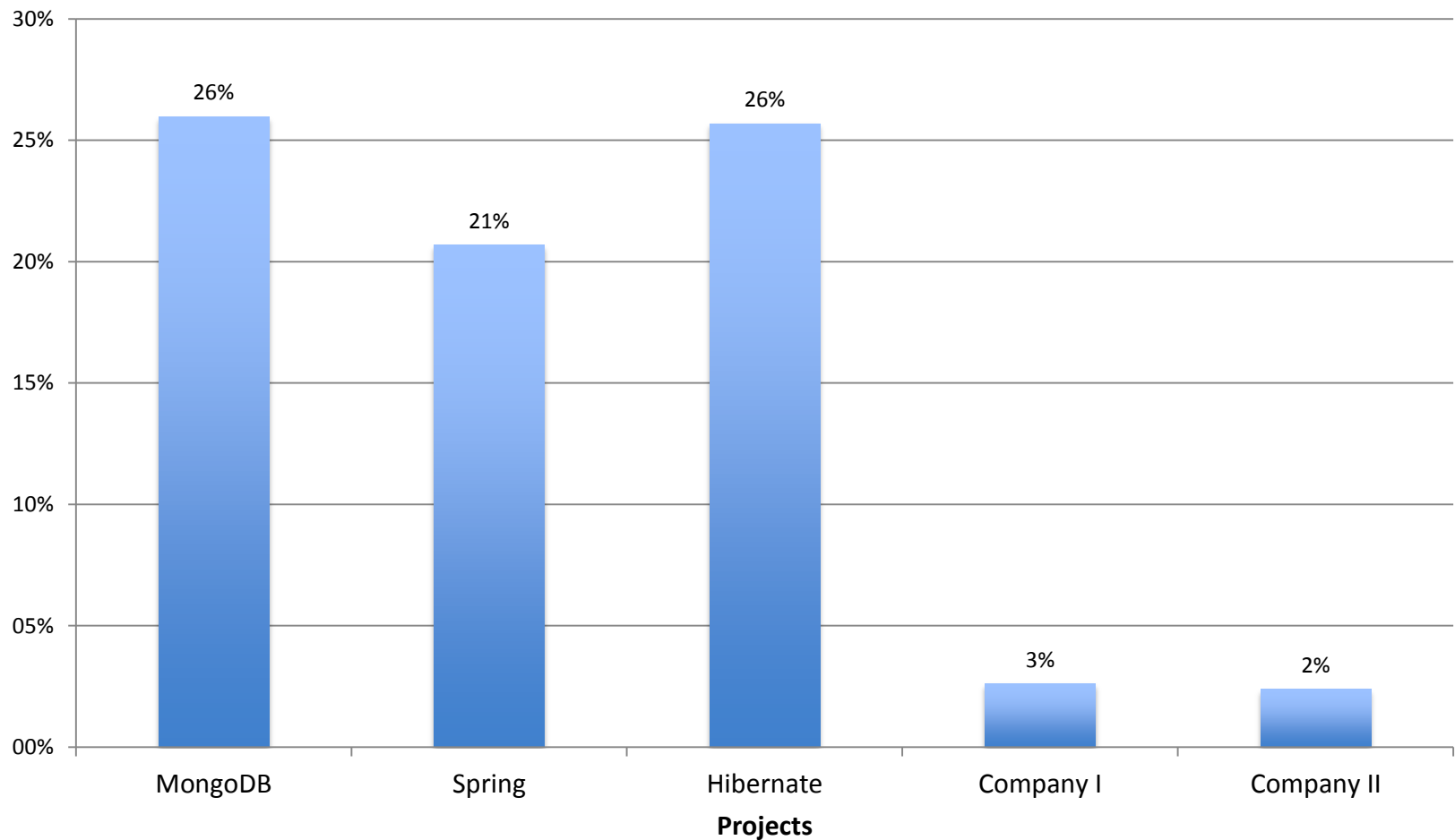Legend: ■ MongoDB   ■ Spring   ■ Hibernate   ■ Company I   ■ Company II

# Results



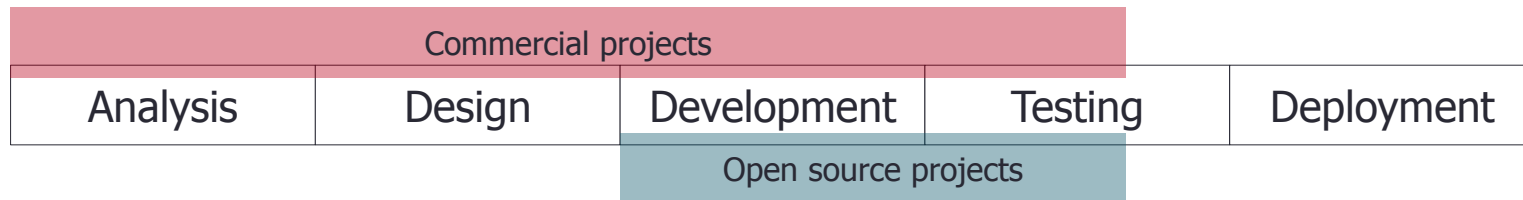Percentage of issues that are resolved by one developer

# Results

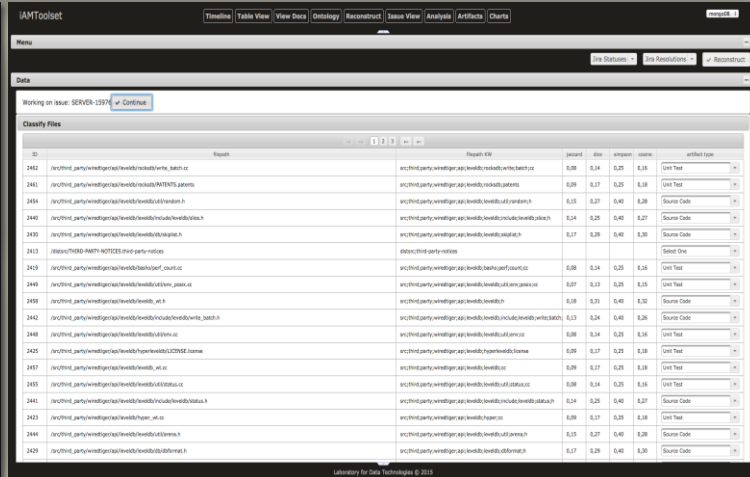Percentage of issues that contain link to another issue

# Additional findings

- Commercial projects usually keep detailed worklogs (e.g. time spent for an issue → date, hours, user…).

- Commercial projects have wider coverage:

| Commercial projects | | | | |
|---|---|---|---|---|
| Analysis | Design | Development | Testing | Deployment |

Open source projects

- Users on open source projects are more disciplined in logging information to software repositories (e.g. issue status).

- Different tools of same software repositories store the all the data needed for reconstruction.

http://goo.gl/QerdGj

# Next steps

- POC – accuracy of the reconstructed workflows – qualitative analysis with IT/Project managers;
- POC – usability of the approach for:
  - Guidance & Control (interviews with developers),
  - Knowledge acquisition and continuous improvement of  the SDM (interviews with IT/Project managers),
  - Project quality analysis
- Workflow analysis: comparison of successful and failed projects.

# Questions

University of Ljubljana
Faculty of Computer & Information Science
Vecna pot 113,
1000 Ljubljana

## **Marko Janković**
### Laboratory for Data Technologies

http://lpt.fri.uni-lj.si/

**Contact**:
e-mail: marko.jankovic@fri.uni-lj.si