

Software Process Recovery

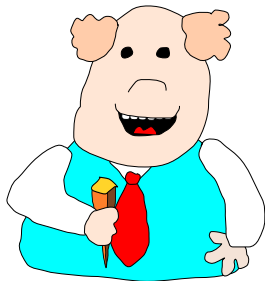
Abram Hindle

Software Architecture Group
David R. Cheriton School of Computer Science
University of Waterloo
Canada

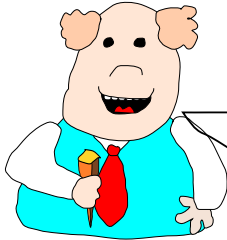
<http://swag.uwaterloo.ca/~ahindle/>

ahindle@cs.uwaterloo.ca

**what is going
on in this
project?**



**Software
Project**

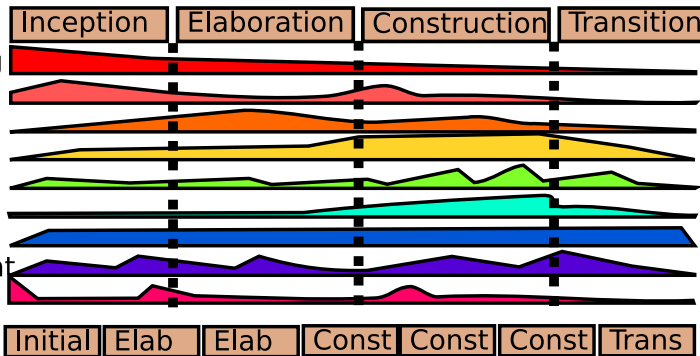


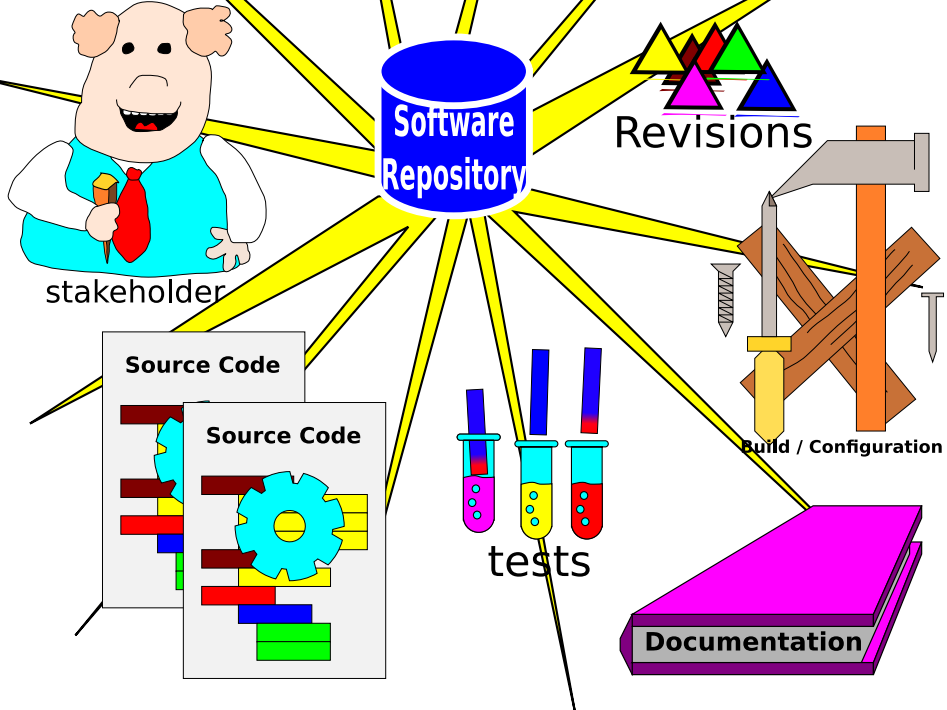
An overview of the project's processes and development would be nice!

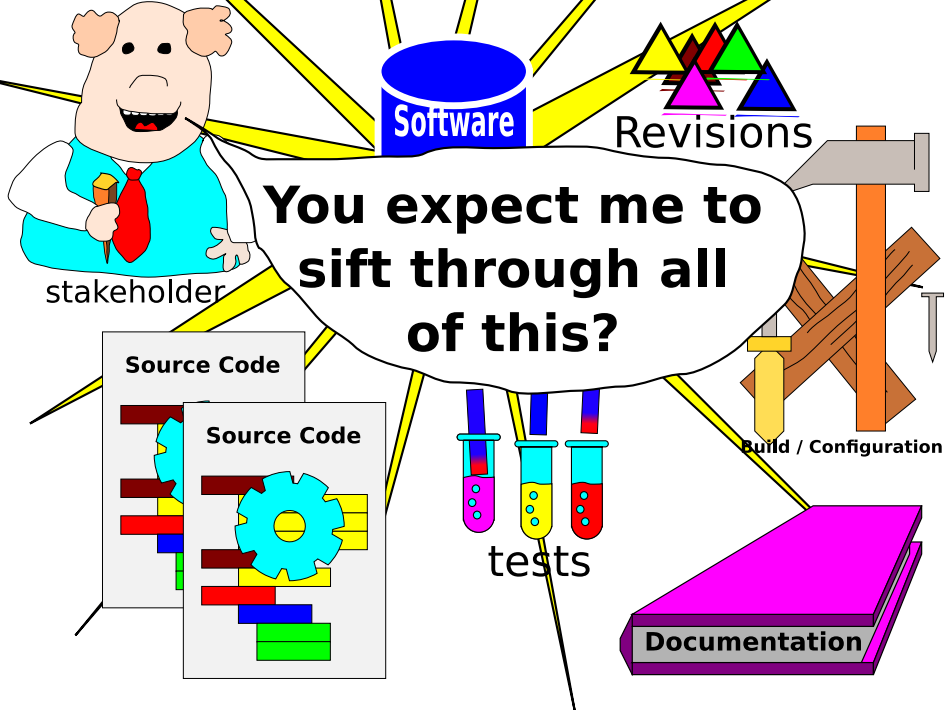
Phases

Disciplines

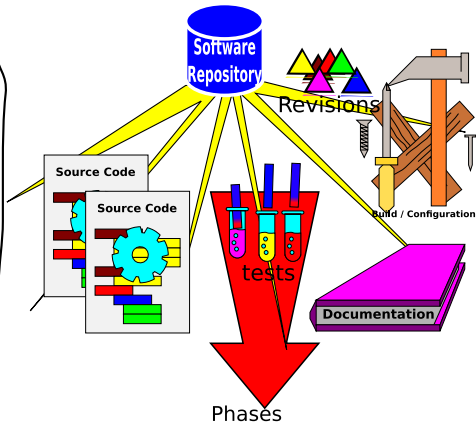
Business Modeling
Requirements
Analysis & Design
Implementation
Test
Deployment
CM and SCS
Project Mangement
Environment





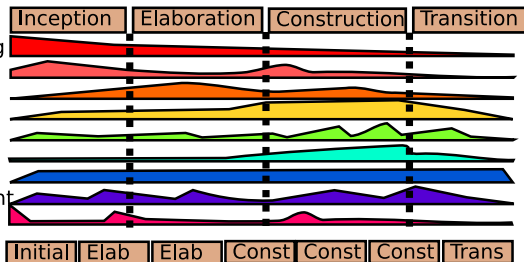


Can't we just summarize what is going on within this project?



Disciplines

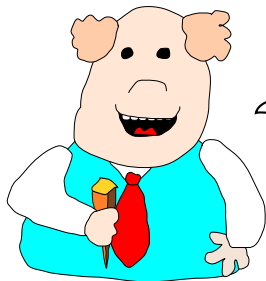
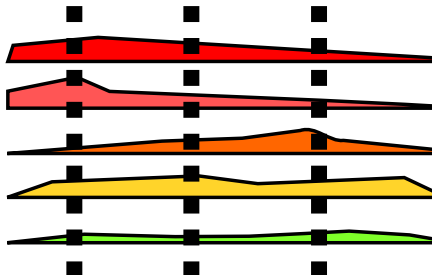
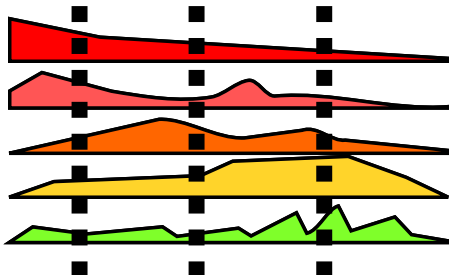
Business Modeling
Requirements
Analysis & Design
Implementation
Test
Deployment
CM and SCS
Project Mangement
Environment



Proposed Process

Recovered Process

Workflows

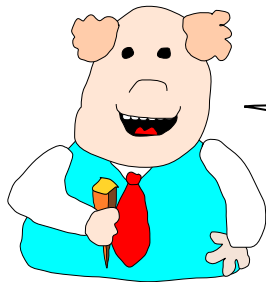
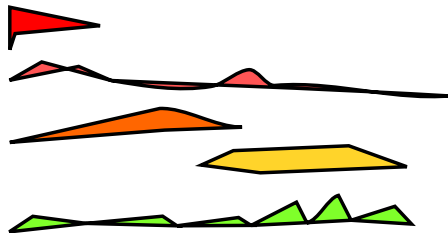


**With an overview
I can validate this
process against
what I expected!**

Proposed and Recovered
Process Overlaid

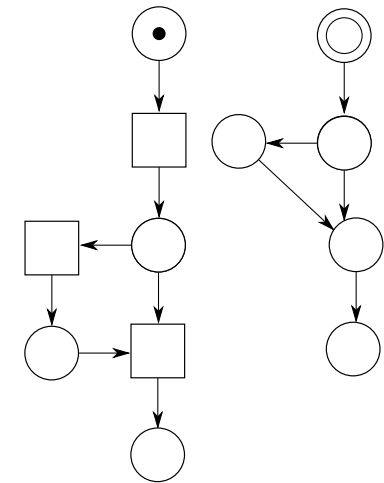


Differences between
Proposed and Recovered



**I can compare
and contrast the
observed process
versus the
expected process!**

Process Mining

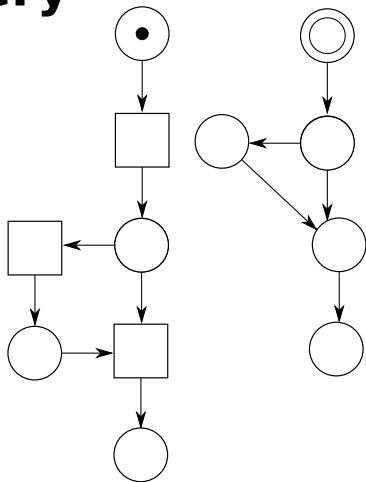
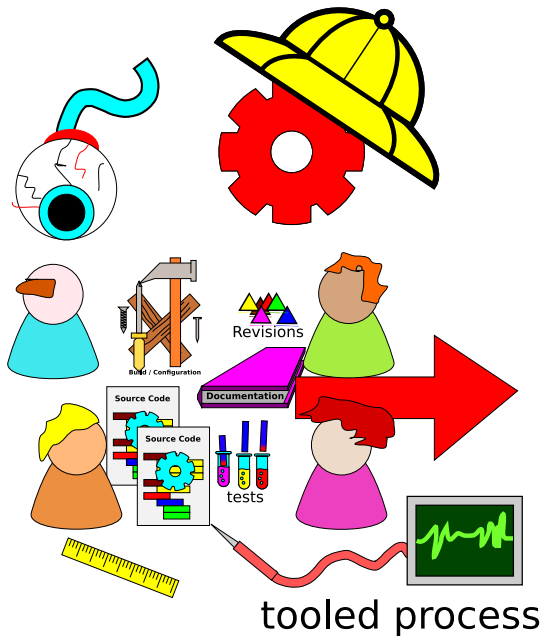


Petri net

FSM

[Aalst]

Process Discovery

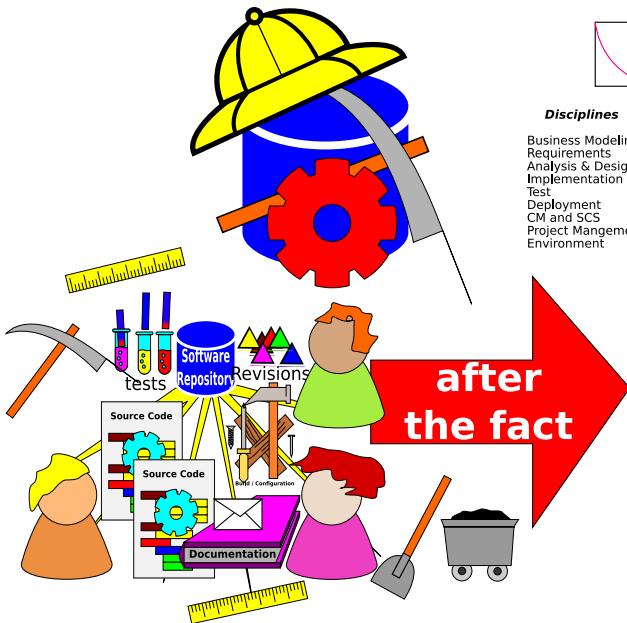


Petrinet

FSM

[Cook]

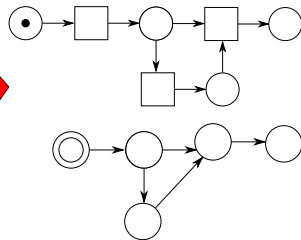
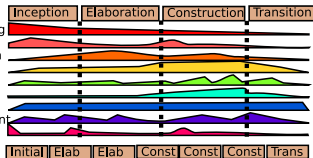
Process Recovery



Phases

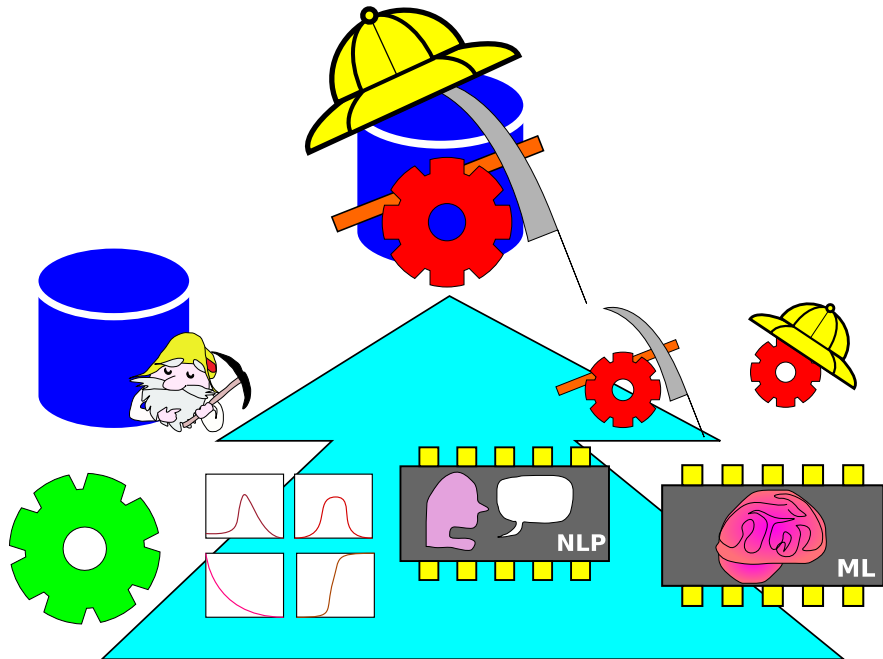
Disciplines

Business Modeling
Requirements
Analysis & Design
Implementation
Test
Deployment
CM and SCS
Project Management
Environment

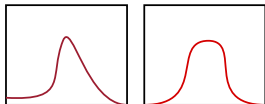
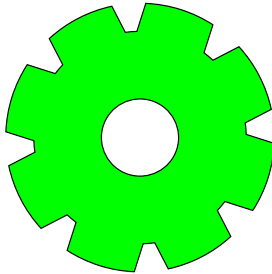


[Hindle]

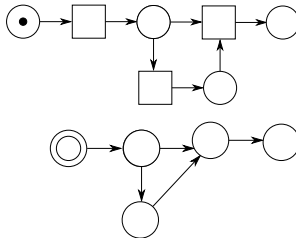
Process Recovery



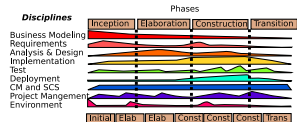
Process



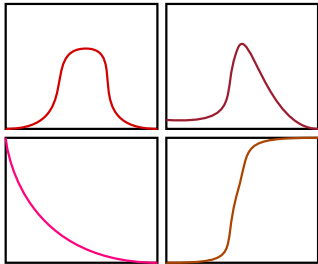
Stochastic



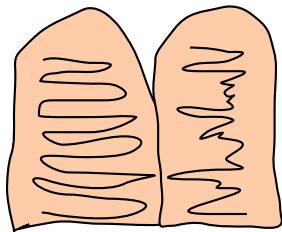
Business



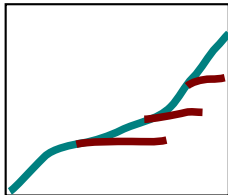
Stochastic Processes



[Herraiz]



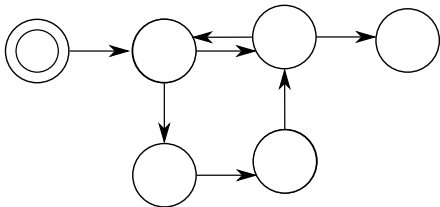
[Lehman]



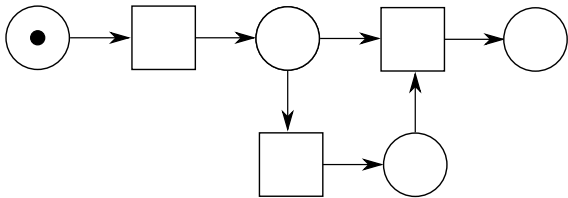
**[Turksi]
[Tu]**

Business Processes

Finite State
Machines



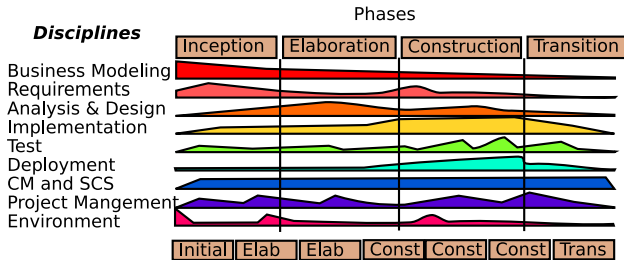
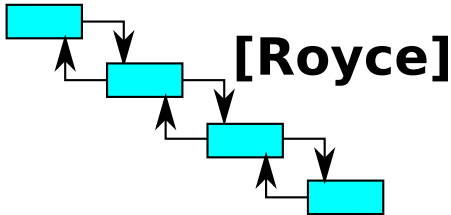
Petrinets



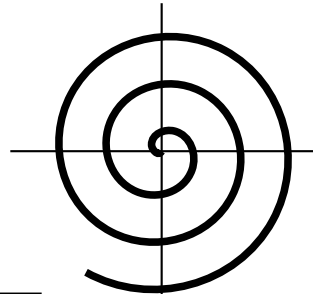
Execution of business goals

[Aalst]

Software Development Processes



[Rational]

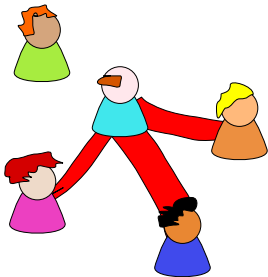


[Boehm]

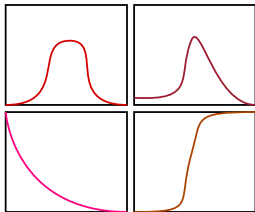
*** CMM**

*** SDLC**

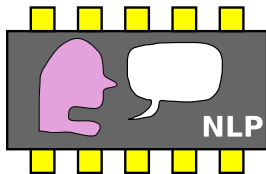
Analysis



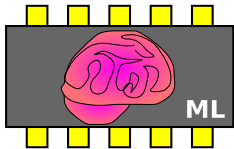
SNA



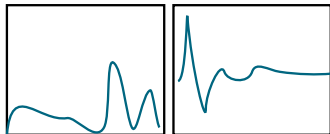
Statistics



NLP

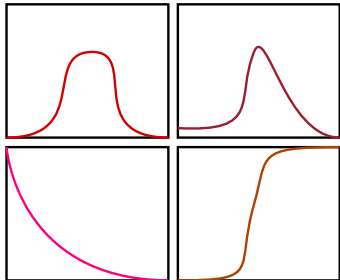


ML

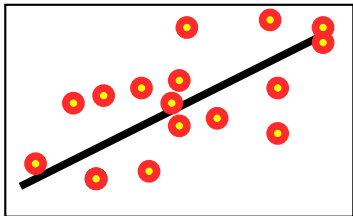


Timeseries

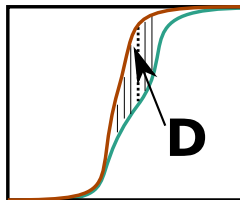
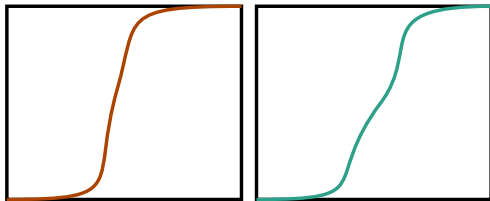
Statistics



Distributions

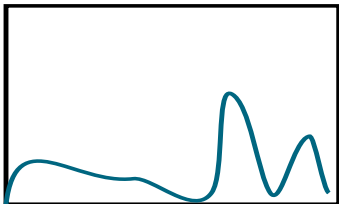


Linear Regression

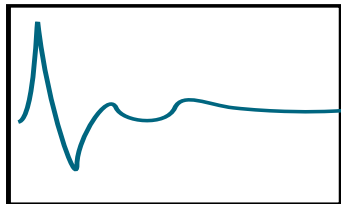


Compare
distributions

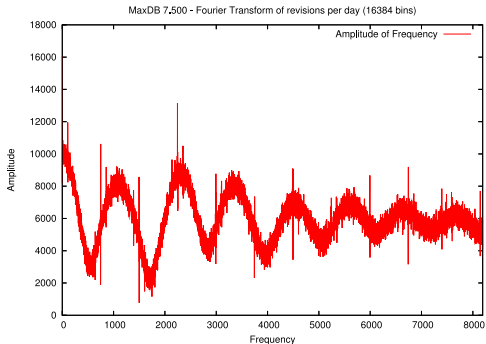
Timeseries



Timeseries

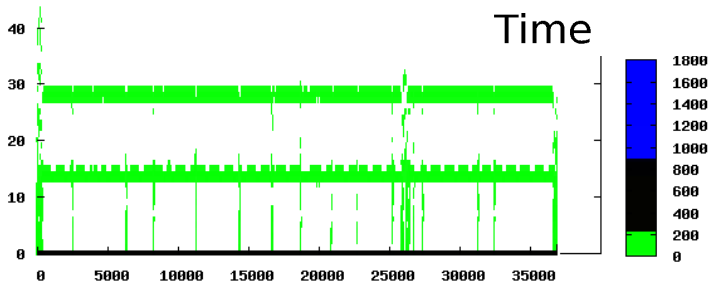
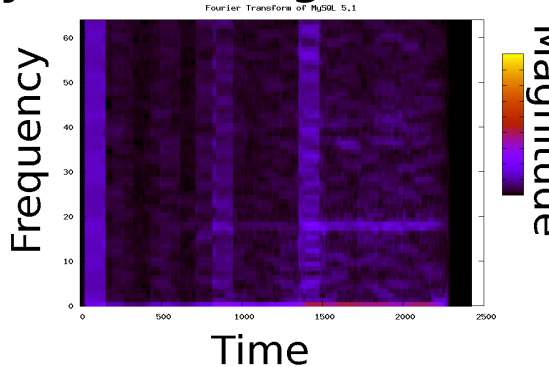
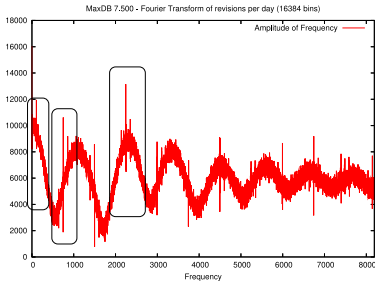


Autocorrelation

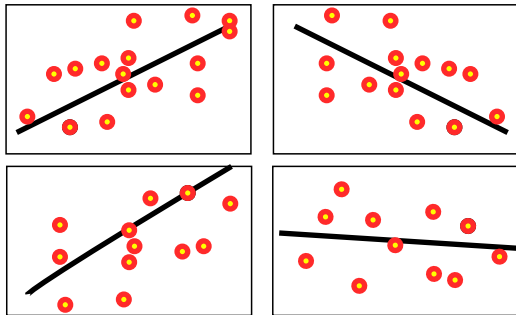


[Herraiz]
[Hindle]

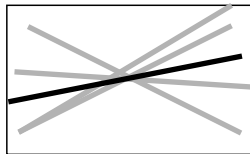
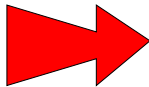
Mining Recurrent Activities: Fourier Analysis of Change Events



Longitudinal Studies and Multilevel Modelling

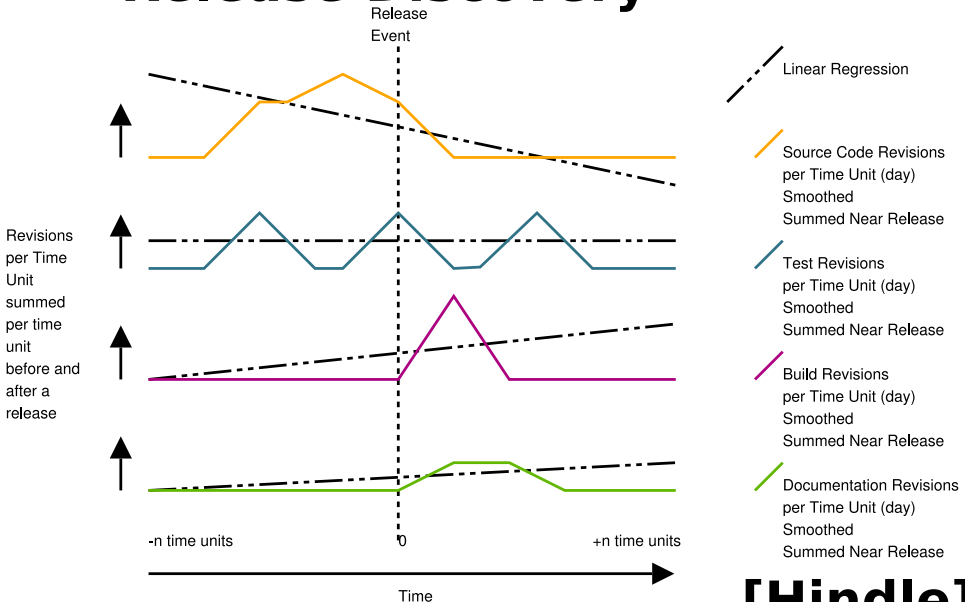


First level model



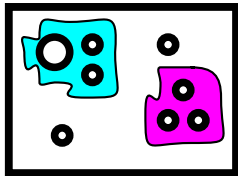
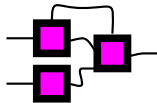
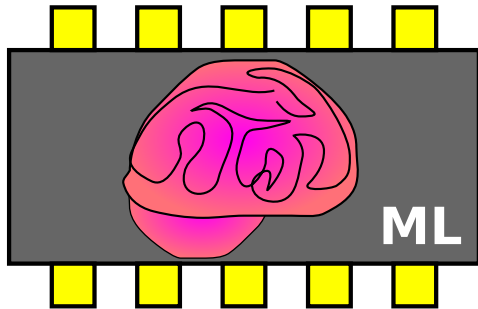
Second level model

Release Discovery

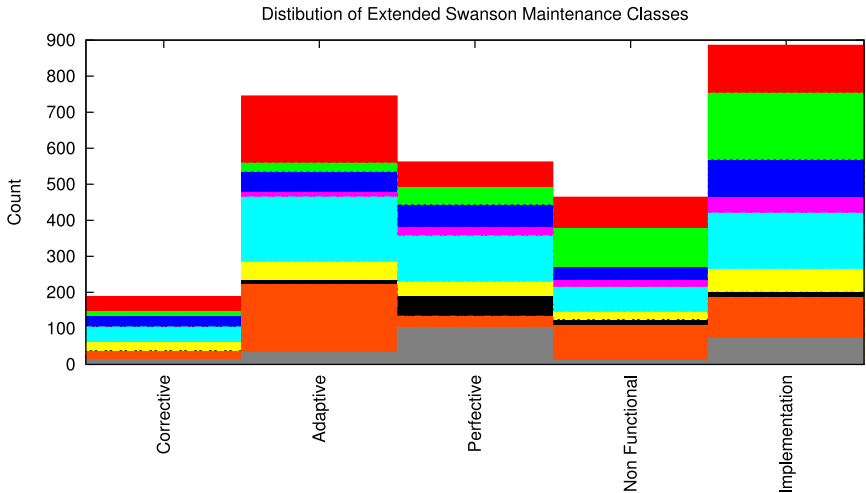


[Hindle]

Machine Learning



Automatic Classification of Large Changes into Maintenance Categories



Extended Swanson Categories

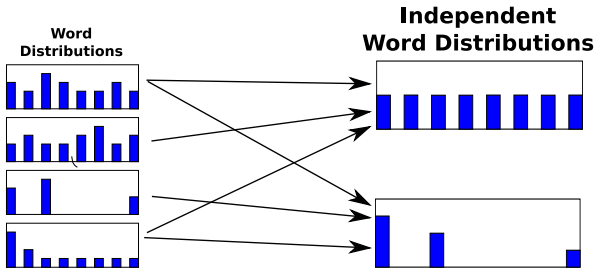
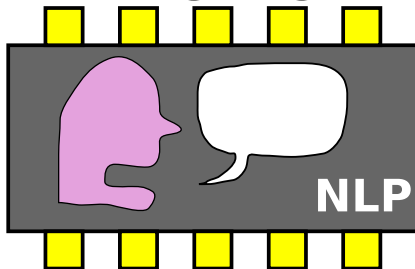
Boost
EGroupware
Enlightenment

Evolution
Firebird
MySQL 5.0

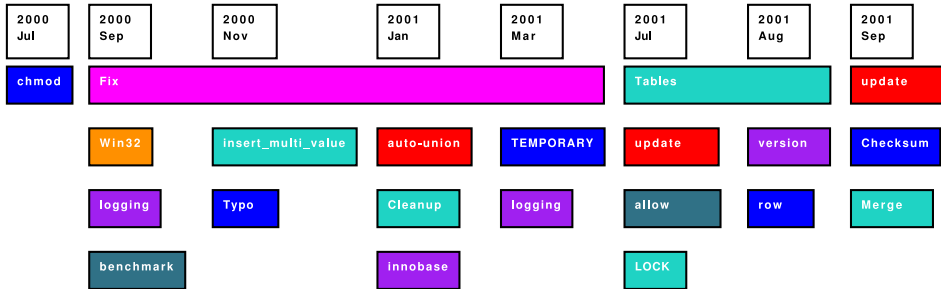
PostgreSQL
Samba
Spring Framework

[Hindle]

Natural Language Processing



Topic Analysis: MySQL 3.23



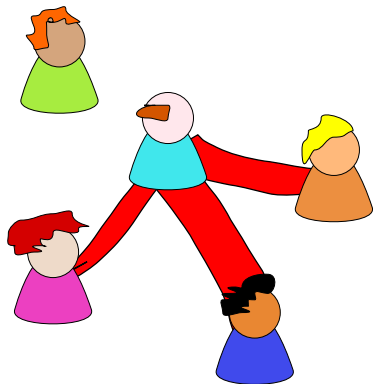
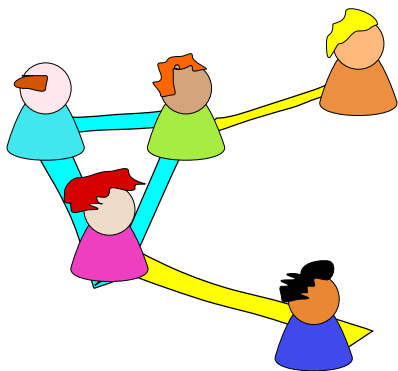
Extracted topics across time

What's hot what's not: developer topic analysis

2004 Jun 2004 Jul 2004 Aug 2004 Sep 2004 Oct 2004 Nov 2004 Dec 2005 Jan 2005 Feb 2005 Mar 2005 Apr 2005 May 2005 Jun 2005 Jul 2005 Aug 2005 Sep 2005 Oct 2005 Nov 2005 Dec 2006 Jan 2006 Feb 2006 Mar 2006 Apr 2006 May 2006 Jun 2006

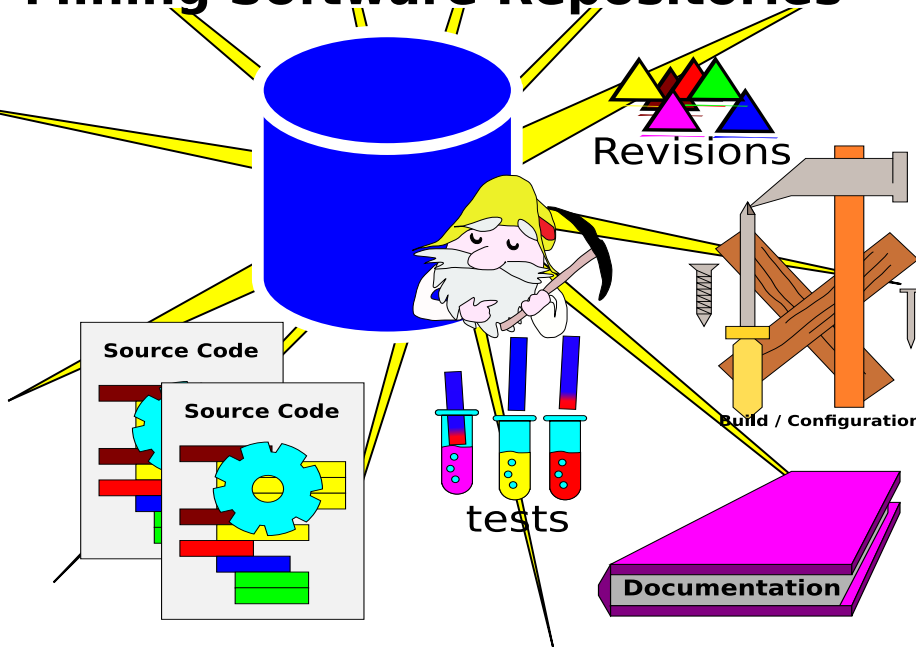


Social Network Analysis

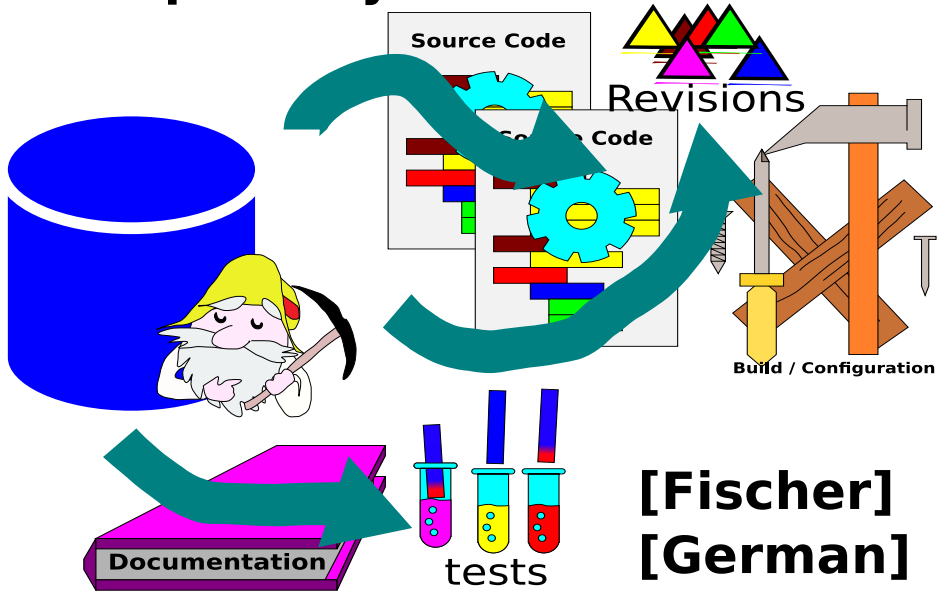


[Bird]

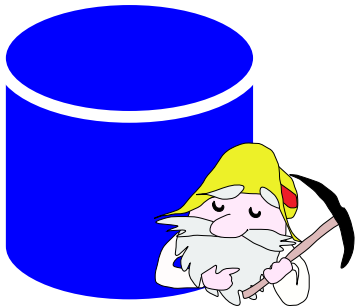
Mining Software Repositories



Repository Extraction



Prediction

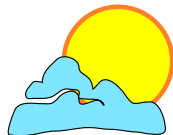
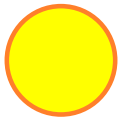
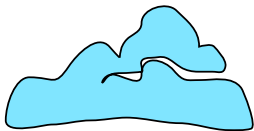
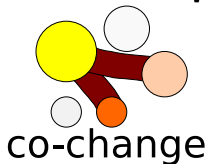


window decay

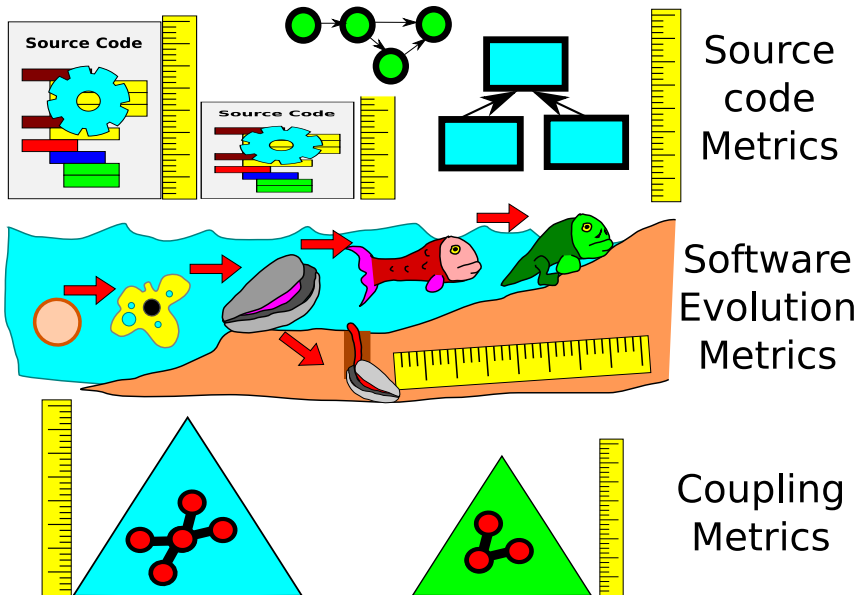
[Girba]

[Askari]

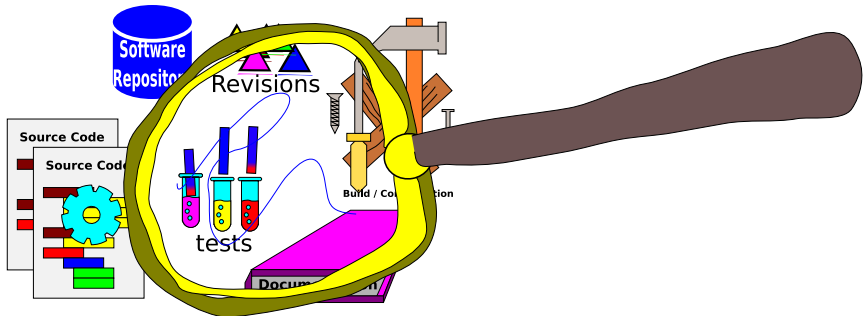
[Hassan]



Metrics



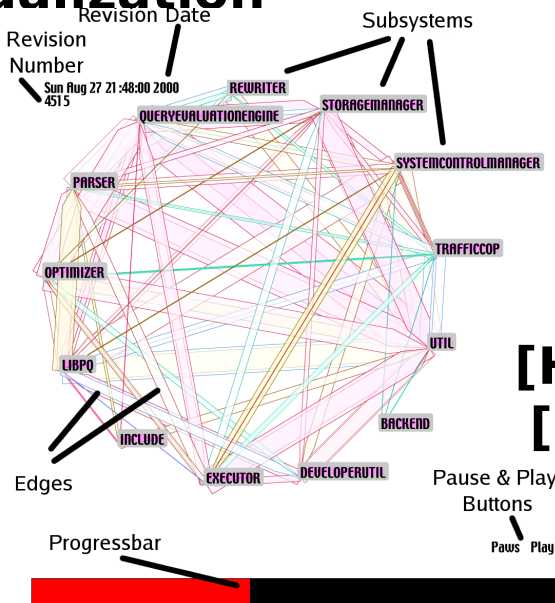
Querying



- * 1st Order Logic
- * Temporal Logic
- * Unification

[Cubranic]
[Kim]
[Hindle]

Visualization



Topic/Concept Analysis

	Entity1	Entity2	Entity3	Entity4	Entity5	Entity6	Entity7	Entity8
Concept 1								
Concept 2								
Concept 3								
Concept 4								
Concept 5								

[Lukins]

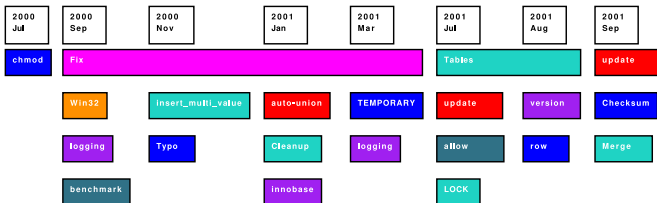
[Linstead]

[Maletic]

[Poshyvanik]

[Marcus]

[Hindle]



Software Processes & MSR

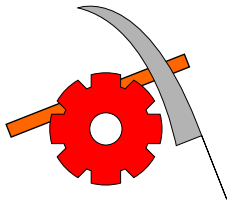
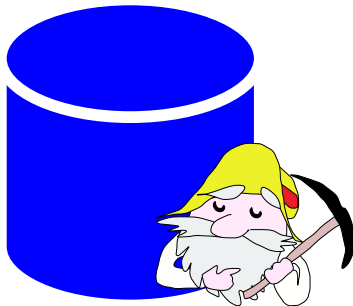
[Aalst]

[Cook]

[Jensen]

[German]

[Hindle]



Process Mining

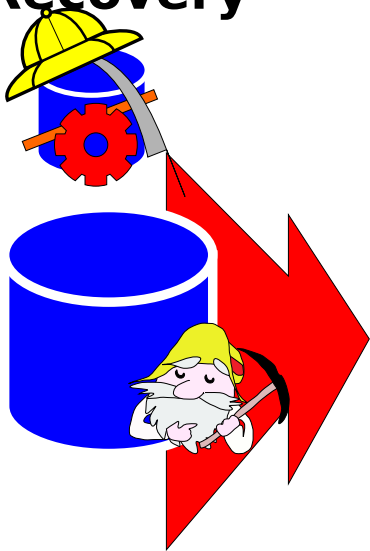


Process Discovery

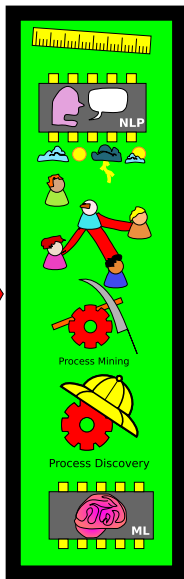


Process Recovery

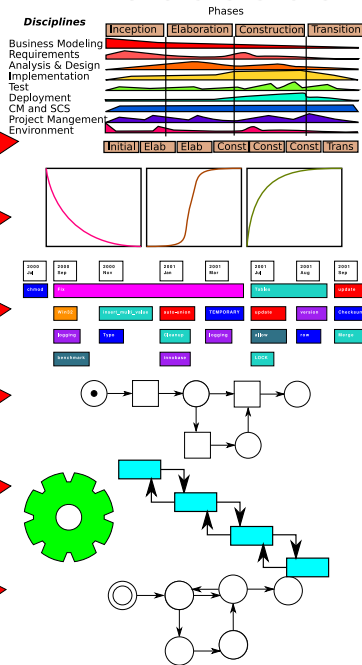
Process Recovery



Tools



Processes



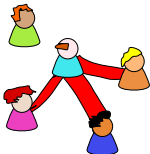


Iteration

Phase

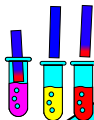
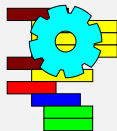
Aggregate

**Fine
Grained**



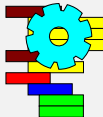
**Fine
Grained**

Source Code

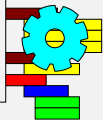


tests

Source Code



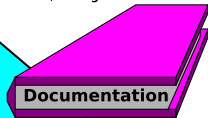
Source Code

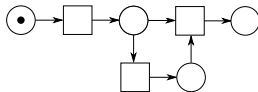
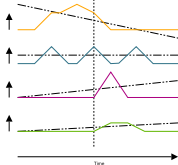
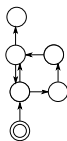


Build / Configuration

Revisions

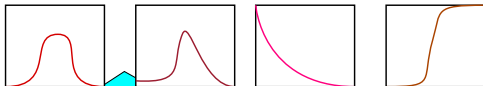
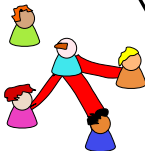
Documentation



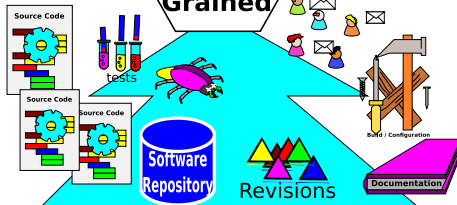


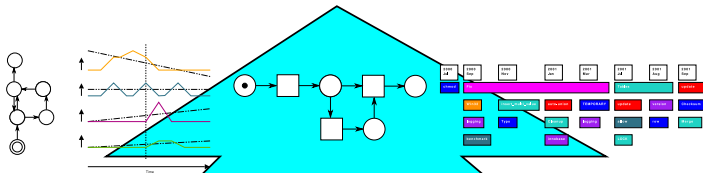
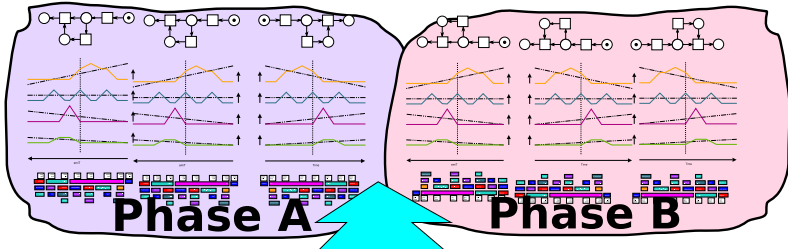
2000 24	2000 24	2000 24	2001 24	2001 24	2001 24	2001 24	2001 24	2001 24	2001 24
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice
choice	choice	choice	choice	choice	choice	choice	choice	choice	choice

Aggregate



Fine Grained





Iteration 1

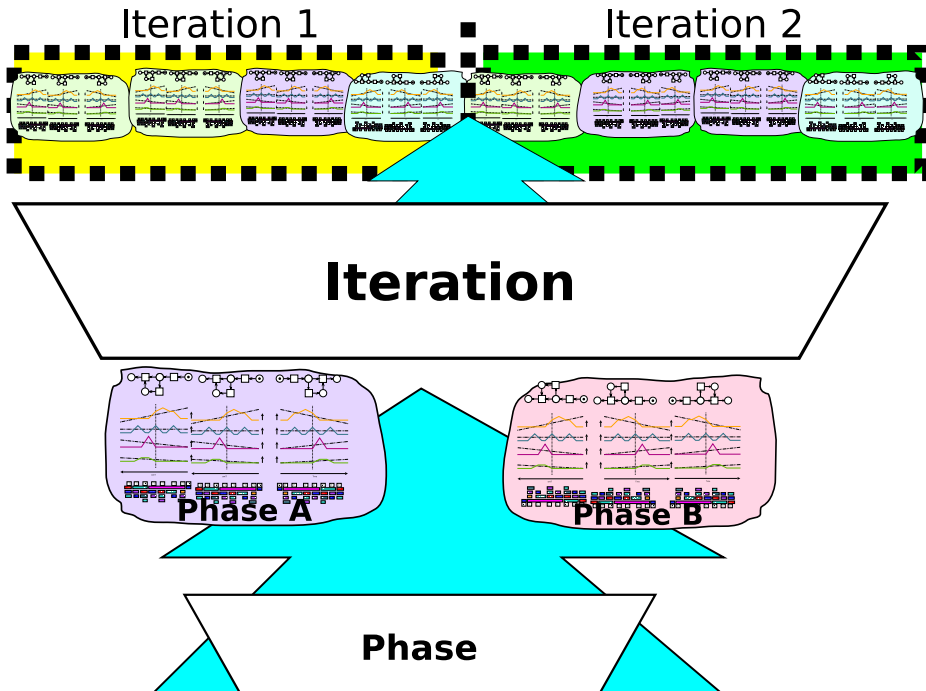
Iteration 2

Iteration

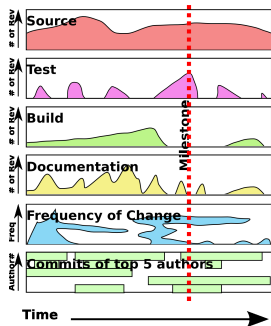
Phase A

Phase B

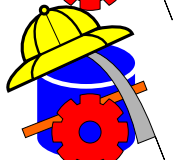
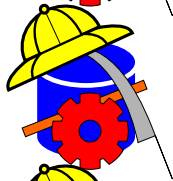
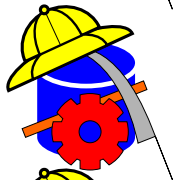
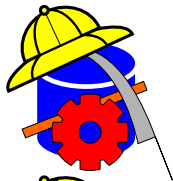
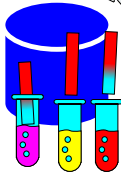
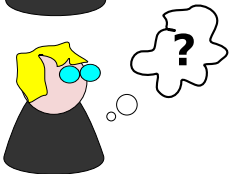
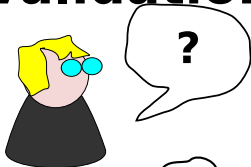
Phase



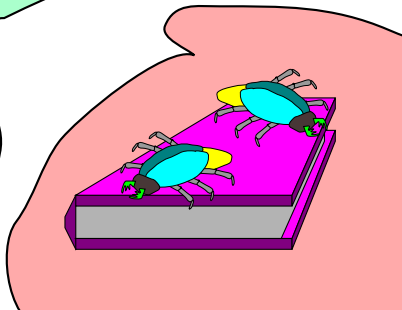
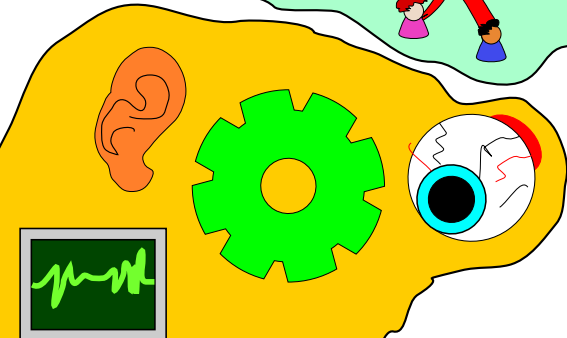
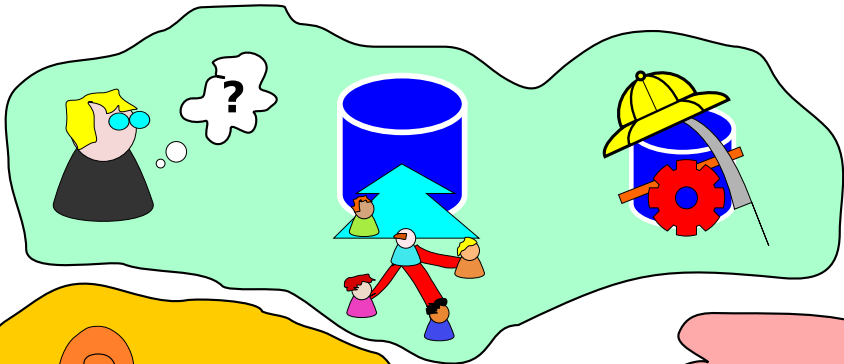
Time Correlation / Slicing



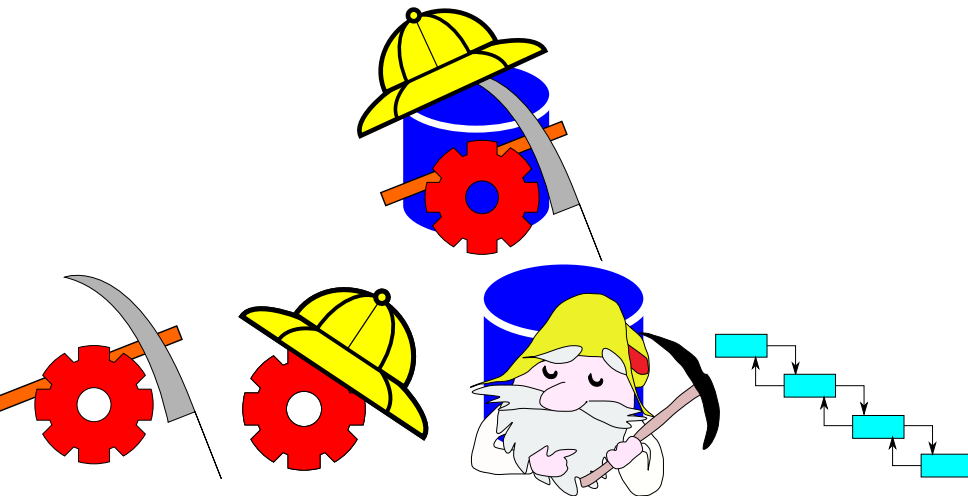
Validation



Assumptions



Conclusions

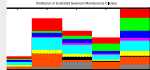
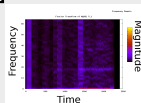


Research Timeline

The Past

Published

- YARN - visualization [VISSOFT 2007]
- Release Discovery (MLM) [MSR & ICSM 2007]
- evolution metrics [ICPC/SCAM 2008, SSP]
- study of large changes [MSR 2008]
- change classification [ICPC 2009]
- recurrent behaviour [ICSE 2009]
- topic analysis [ICSM 2009]



Present and Future

In progress

- finish Topic Naming paper [writeup]
- finished MLM journal paper [casestudy]
- slicing [writeup]
- multi-timeline correlation [some imp]

To do

- Unified process diagram summary [imp]
- phase & iteration identification [imp]
- Thesis

