# Evidence-based 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 are we going to do?



Theory

Practice

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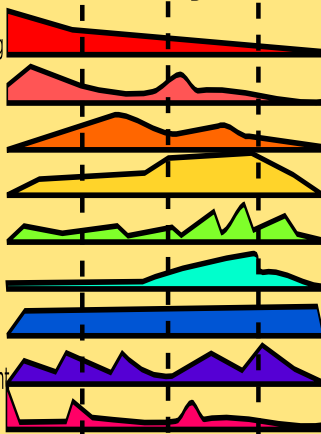Analysis & Design
Implementation
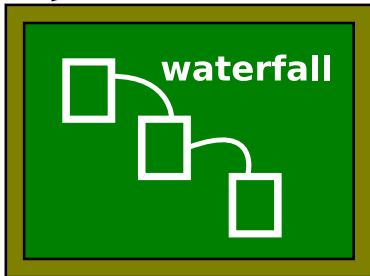Test
Deployment
CM and SCS
Project Mangement
Environment
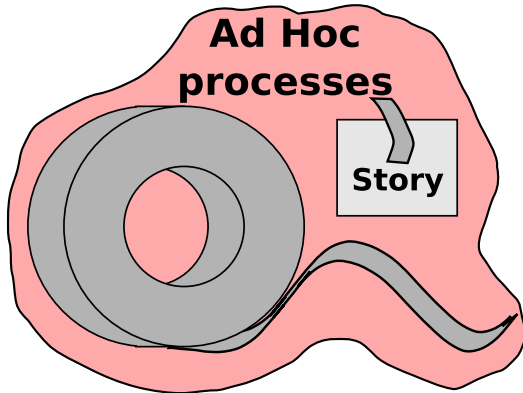
# Process



prescribed processes

## Rx

- Test First
- Scrums
- Story Cards

waterfall

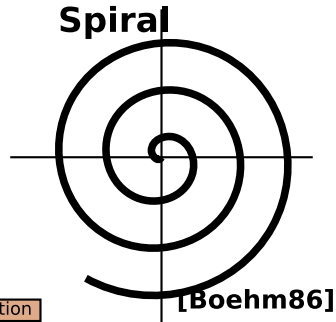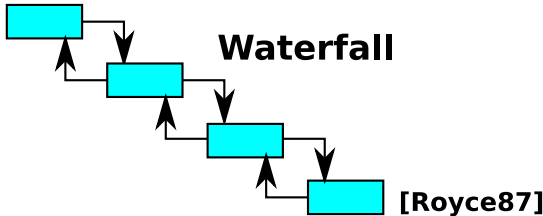**Formal Processes**

Ad Hoc processes

Story

3

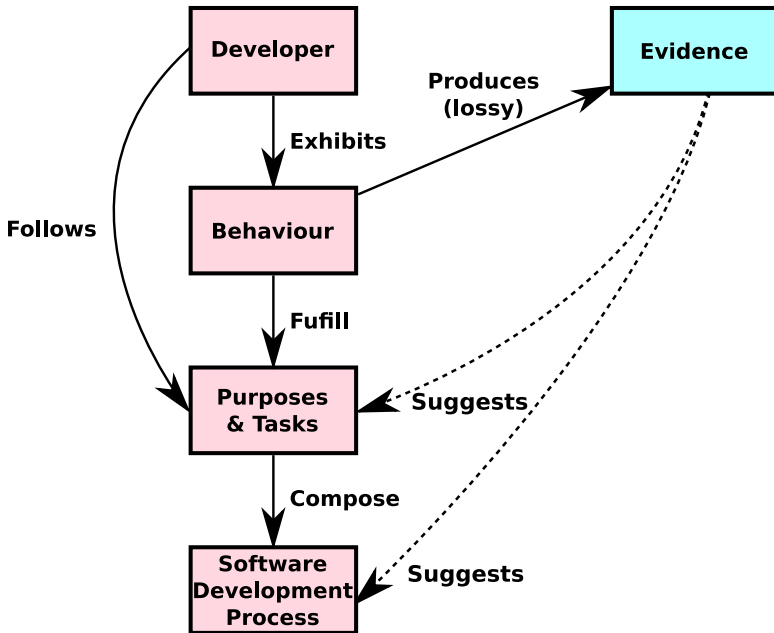# Software Development Processes



**Waterfall**

[Royce87]

**Spiral**

[Boehm86]

## Unified Process

[Jacobson99]

* **CMM**

* **SDLC**

Developer

Evidence

Follows

Exhibits

Produces
(lossy)

Behaviour

Fufill

Purposes
& Tasks

Suggests

Compose

Software
Development
Process

Suggests

Actual

# Research Relationships

|  | Behaviour | Intent, Purpose and Tasks | Software Development Process |
|---|---|---|---|
| **Release Patterns: Source/Test/ Build/Docs** | ● | | ● |
| **Large Changes** | | ● | |
| **Topic Analysis** | ● | ● | |
| **Recovered Unified Process Views** | | | ● |

# Motivation: Stakeholders



Fixers or
Star Programmers



Managers



Investors and
Acquisitions



New Developers



Employees assigned
to a ISO9000
conformance project

Proposed and Recovered Process Overlayed

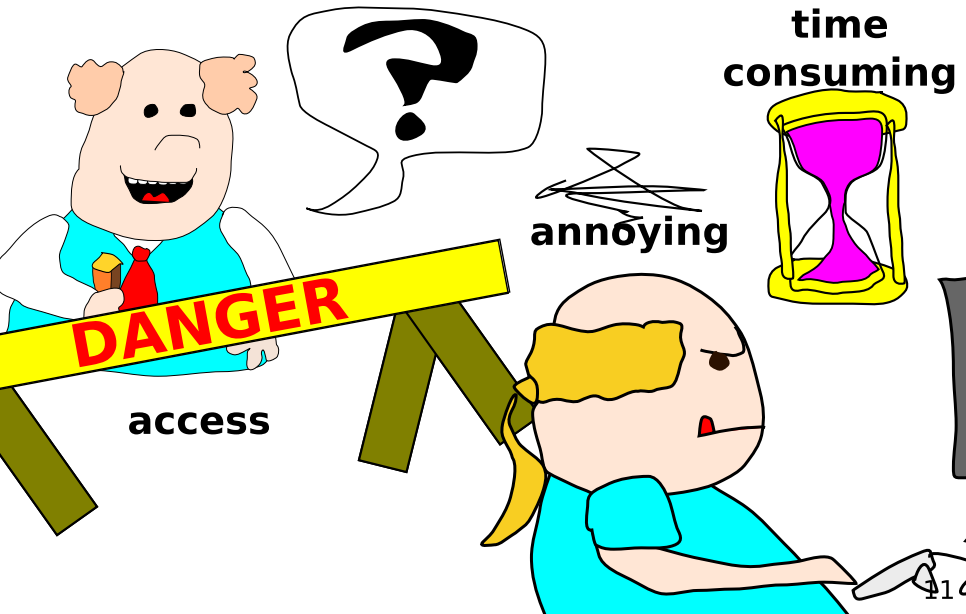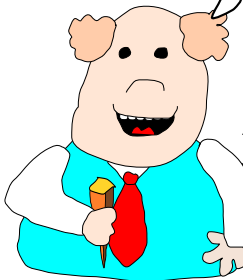Differences between Proposed and Recovered

Workflows

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

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

Software Repositories

Revisions

Build / Configuration

Source Code
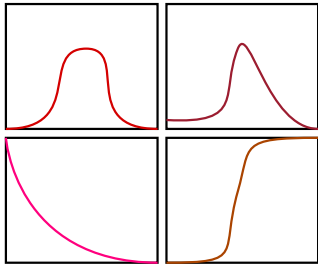
Source Code

Tests

Documentation

Phases

**Disciplines**

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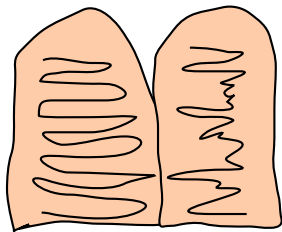Test
Deployment
CM and SCS
Project Mangement
Environment

| Inception | Elaboration | Construction | Transition |
|-----------|-------------|--------------|------------|

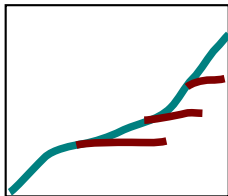| Initial | Elab | Elab | Const | Const | Const | Trans |
|---------|------|------|-------|-------|-------|-------|

12

# PREVIOUS WORK

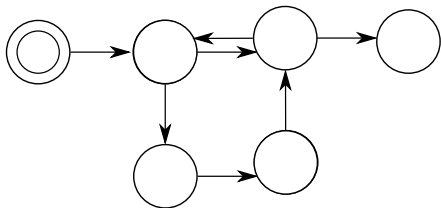## RELATED RESEARCH

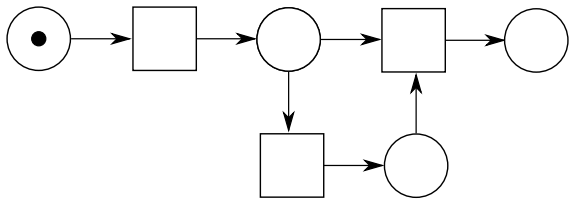# Stochastic Processes



[Herraiz]

[Lehman]

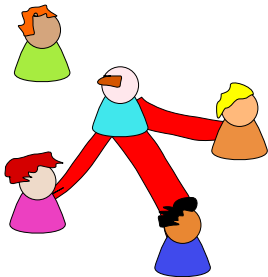[Turksi]
[Tu]
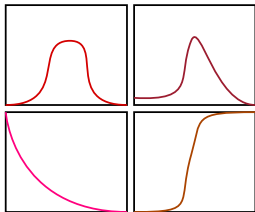
# Business Processes

Finite State Machines



Petrinets



Execution of business goals

**[Van der Aalst]**

# Analysis

**SNA**

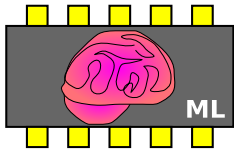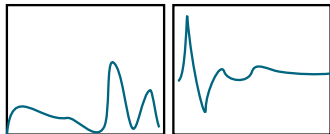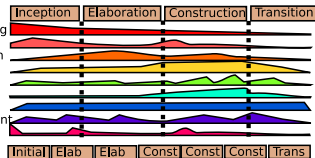**Statistics**

**NLP**

**ML**

**Timeseries**

# Process Mining



**Business Goals**

**Petrinet**   **FSM**

# Process Discovery

Petrinet    FSM

tooled process

# Process Recovery



Dev Mailing List Archive

User Mailing List Archive

Version Control

Bug tracker system

after the fact

# Unified Process Diagram



Phases

Disciplines

| Inception | Elaboration | Construction | Transition |

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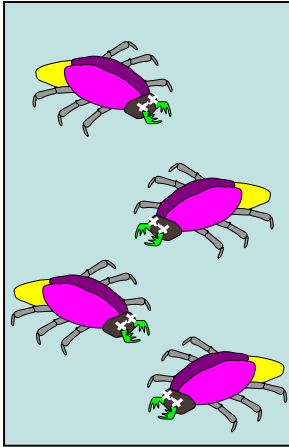- Deployment
- CM and SCS
- Project Mangement
- Environment

| Initial | Elab | Elab | Const | Const | Const | Trans |

# Mining Software Repositories



Revisions

Source Code

Source Code

Tests

Build / Configuration

Documentation

# Source Acquisition



discussions  bugs  source

## Initial Repositories and artifacts

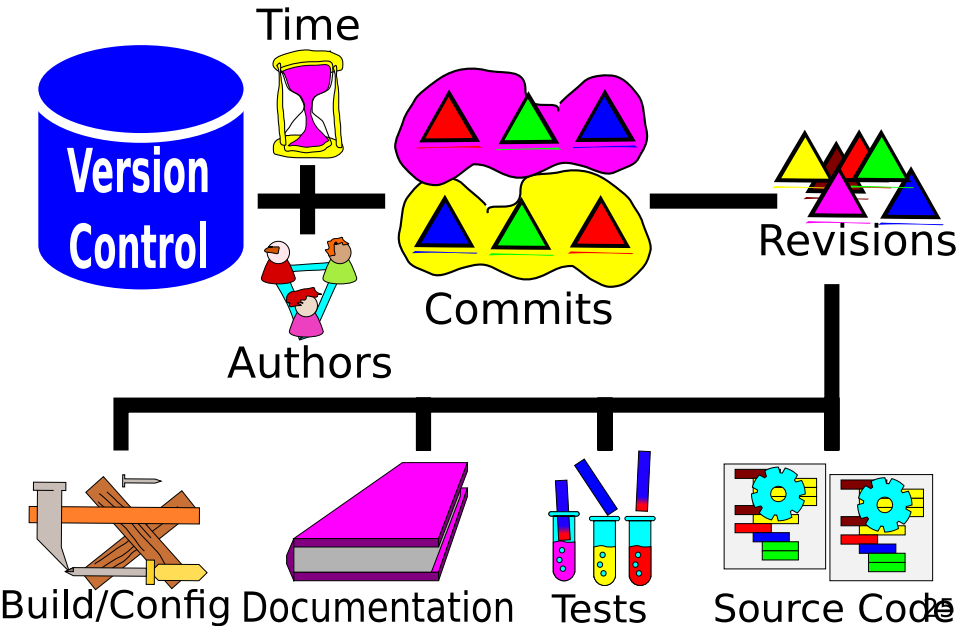# Extraction: Mailing list archives

# Extraction: Bug trackers

# Extraction: Version Control



Time

Version Control

Authors

Commits

Revisions

Build/Config  Documentation  Tests  Source Code

# Metrics



Source code Metrics

Software Evolution Metrics

Coupling Metrics

# Topic/Concept Analysis



27

# Quality Related
# Non functional requirements


**portability**


**reliability and functionality (includes correctness)**


**usability**


**efficiency**


**maintainability**

# SOFTWARE PROCESS RECOVERY

# Release Patterns: STBD

[Hindle ICSM07]

# STBD applied to SQLite

# Maintenance Classes of Large Changes

[Hindle ICPC09]

# Supervised: Maintenance Classes



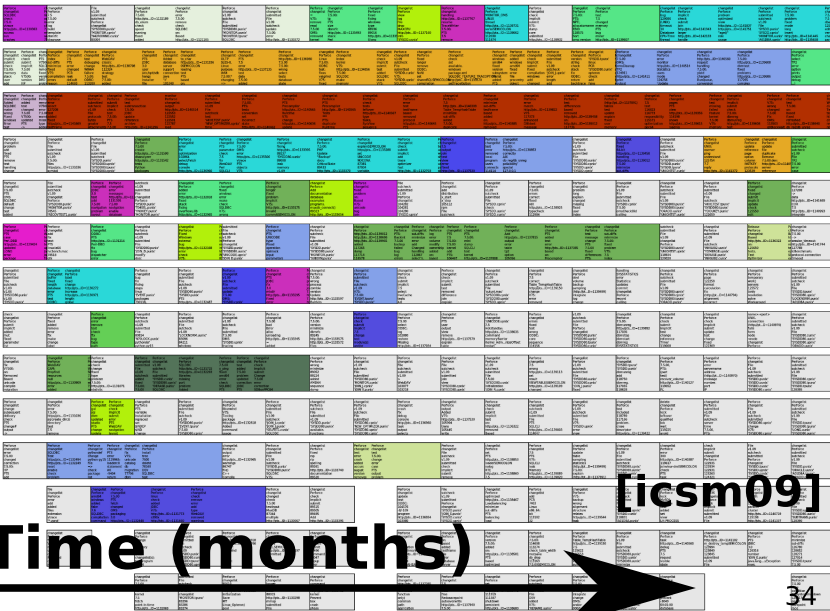Proportional Distibution of Extended Swanson Maintenance Classes

[hindle icpc09]

Extended Swanson Categories

| Boost | Evolution | PostgreSQL |
| EGroupware | Firebird | Samba |
| Enlightenment | MySQL 5.0 | Spring Framework |

33

**Developer Topics**

Time (months)

Unique Topics

[icsm09]

34

# Word Bag Examples

## Portability

portability
transferability
interoperability
documentation
internationalization
i18n

**...**

## Reliability

reliability
failure
error
redundancy
fails
bug

**...**

# Labelled Developer Topics



Unique Topics

Time (months)

efficiency portability

efficiency

functionality

portability

maintainability

efficiency

reliability

maintainability portability

functionality

[Hindle and Ernst et al. http://softwareprocess.es/whats-in-a-name]

36

# Recovered Unified Process Views

Theory

Practice

Business Modeling
Requirements
Analysis & Design
Implementation
Test
Deployment
CM and SCS
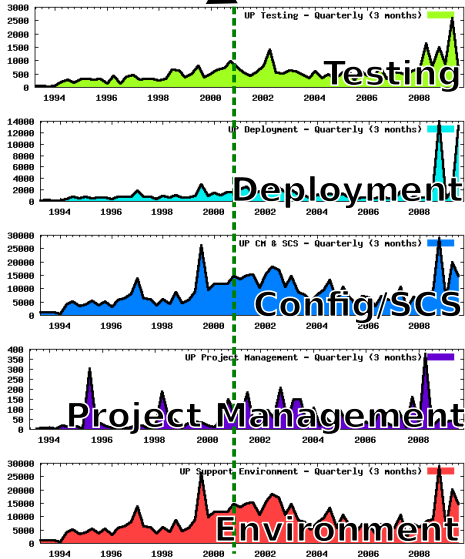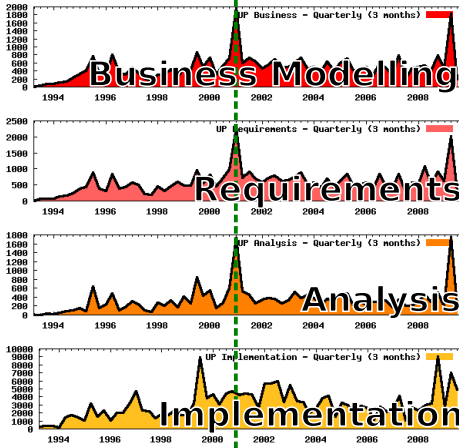Project Mangement
Environment
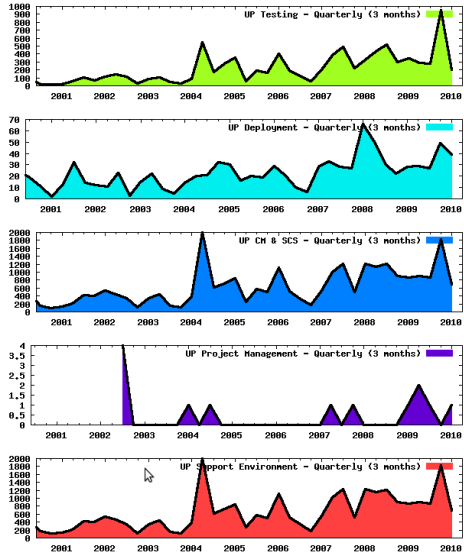
[ICSM10]

# UP Requirements Signal

# UP Implementation Signal

# UP Testing Signal

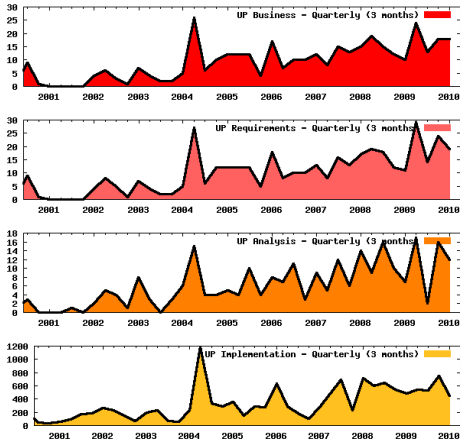# FreeBSD Case Study

# FreeBSD Case Study: 2001



Business Modelling

Requirements

Analysis

Implementation

Testing

Deployment

Config/SCS

Project Management

Environment

# SQLite Case Study

# SQLite Case Study: 2009

# UP Observability



**Disciplines**

Business Modeling
Requirements
Analysis & Design

Project Mangement
Environment

# Common Threads



**Idioms**

*.tex  *.doxygen  FILES  Makefile,*
INSTALL  doc/  Makefile
README
AUTHORS  configure
TODO  Configure,*
build.xml
Documentation  setup.hs  setup.py
*test*
*.txt
*.hs  *.php  *.java  *.scm
*.ml  *.lisp  *.cpp  unit tests
*.tcl  *.py  *.sql
*.c  *.test
Source Code  Source Code  Source Code  *.pl
*.pl  *.pm  Test
*.rb

**Shared Terms**

Usability  Maintainability

Portability  Reliability  Efficiency

**External Project Language**

shared vocabulary

vocabulary  vocabulary

Project  Project

**Internal Project Language**

46

# Future Work



People and teams

Validation

Accuracy

Industrial

Iteration Identification

# Conclusions



Software Process Recovery

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

**Maintenance Classes**

Corrective  Adaptive  Perfective  Non-Code  Implementation

[Hindle ICPC09]

Portability  Functionality  Efficiency  Reliability

**Developer Topic Analysis and Labelling**

LDA LSI

[Hindle and Ernst http://softwareprocess.es/name/]

[Hindle ICSM]

Documentation

**Release Patterns**

Source
Test
Build
Documentation

[Hindle ICSM07]

Test  Source

Build

**Discussions**

**Bugs**

**Source**

Version Control

Revisions

**Managers**

**Fixers or Star Programmers**

**New Developers**

**Investors and Acquisitions**

**ISO9000 Consultants**

48

# Research Timeline

## The Past

Published
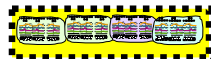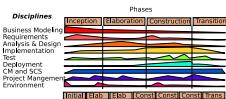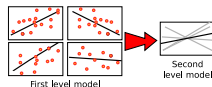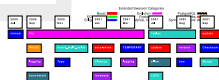
- YARN - visualization [VISSOFT 2007]
- Release Discovery (MLM) [MSR & ICSM 2007]
- evolution metrics [ICPC/SCAM 2008, SSP]
- study of large changes [MSR 2008]
- change classification [ICPC 2009]
- recurrent behaviour [ICSE 2009]
- topic analysis [ICSM 2009]
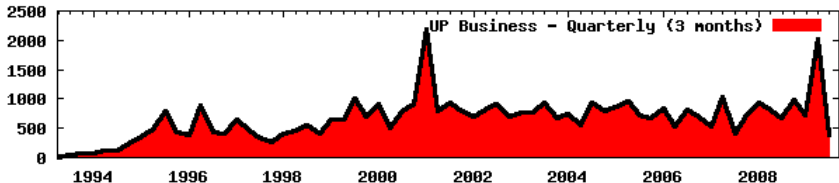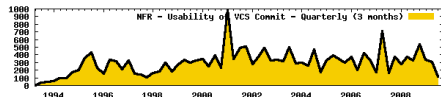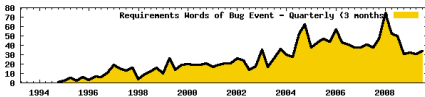- RUPV [ICSM 2010]
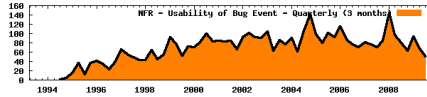- Topic Naming [ICSE Submission]
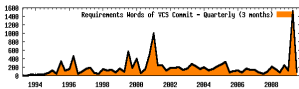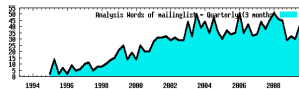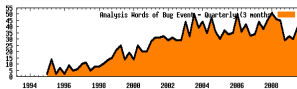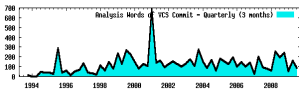
## Present and Future

- finish  Topic Naming paper [writeup]
- finished MLM journal paper [casestudy]
- slicing [writeup]
- multi-timeline correlation [some imp]
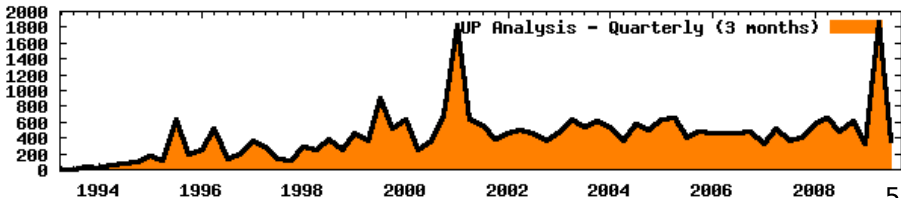- Unified process diagram summary [imp]
- phase & iteration identification [imp]
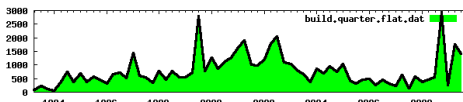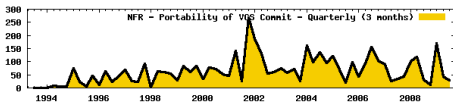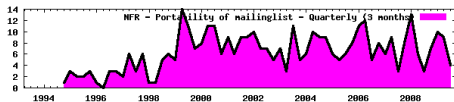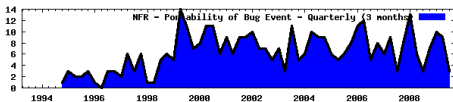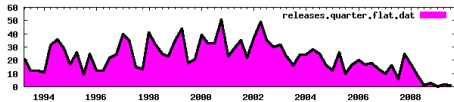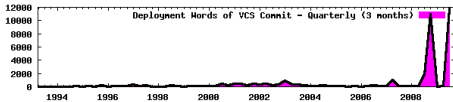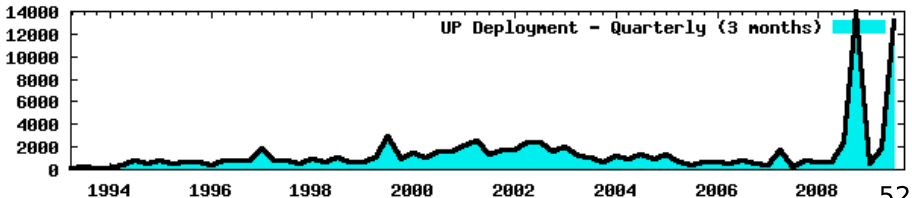- Thesis

49

# UP Business Modelling Signal

# UP Analysis Signal

# UP Deployment Signal



52

# UP Configuration Managment and SCS

# UP Project Management Signal



UP Words – Project Management of Bug Event – Quarterly (3 months)

UP Words – Project Management of VCS Commit – Quarterly (3 months)

UP Words – Project Management of mailinglist – Quarterly (3 months)

UP Project Management – Quarterly (3 months)

Perfor

Clean-

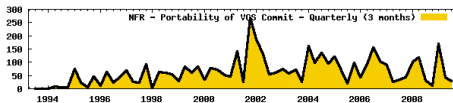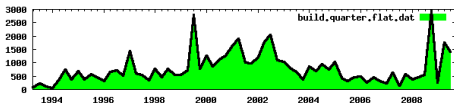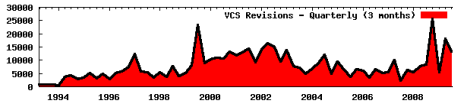Relist

word
bag

54

# UP Environment Signal

**Word bag analysis**

Usability

Maintainability

Portability

Reliability

Efficiency

# Prediction



window decay

co-change

**[Girba]**
**[Askari]**
**[Hassan]**

# Querying



* 1st Order Logic
* Temporal Logic
* Unification

**[Cubranic]**
**[Kim]**
**[Hindle]**

# Visualization



Revision Date

Revision Number

Subsystems

Sun Aug 27 21:48:00 2000
4515

REWRITER

QUERYEVALUATIONENGINE

STORAGEMANAGER

SYSTEMCONTROLMANAGER

PARSER

TRAFFICCOP

OPTIMIZER

UTIL

LIBPQ

BACKEND

INCLUDE

DEVELOPERUTIL

EXECUTOR

Edges

**[Hindle]**

**[Lanza]**

Pause & Play Buttons

Paws   Play

Progressbar

59

# Statistics



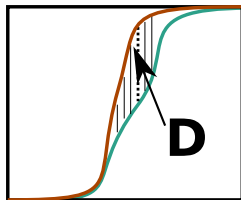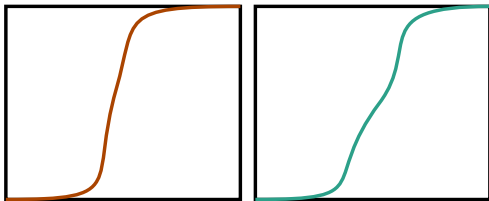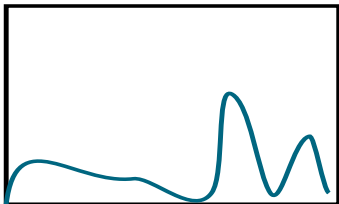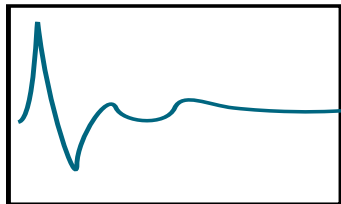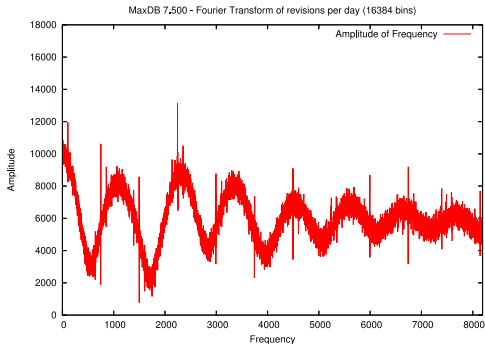Distributions



Linear Regression



**D**

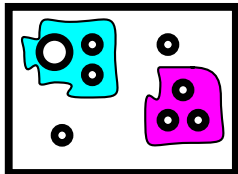Compare
distributions

# Timeseries



Timeseries



Autocorrelation


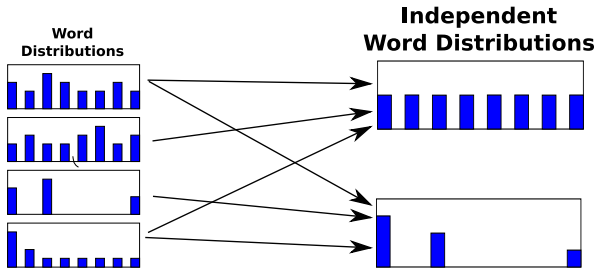
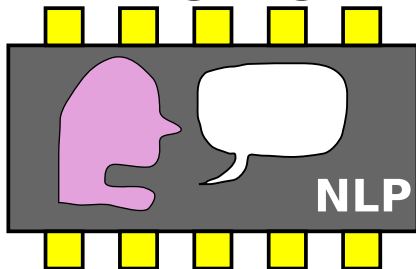MaxDB 7.500 - Fourier Transform of revisions per day (16384 bins)
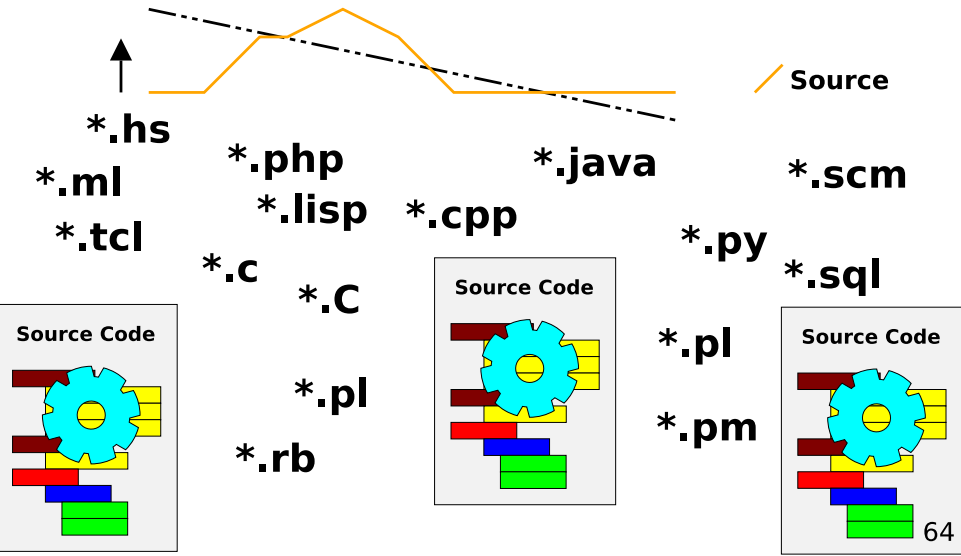
**[Herraiz]
[Hindle]**

# Machine Learning



ML

# Natural Language Processing

# Release Patterns: Source revisions



Source

*.hs

*.ml

*.tcl

*.php

*.lisp

*.c

*.C

*.cpp

*.java

*.py

*.pl

*.rb

*.pl

*.pm

*.scm

*.sql

Source Code

Source Code

Source Code

# Release Patterns: Test Revisions



*.t

*test*

unit tests

*.test

Test

Tests

# Release Patterns: Build file revisions



**Build**

Makefile.*

Makefile

configure

configure.*

build.xml

setup.hs

setup.py

# Release Patterns: Documentation Revisions


Documentation

**FILES**

doc/  INSTALL  *.tex  *.doxygen

AUTHORS  README

*.txt  TODO

maybe: *.html
*.png
*.svg

Documentation