# Software Wartung und Evolution
## Modeling History with Metamodels

**Harald Gall**

Institut für Informatik

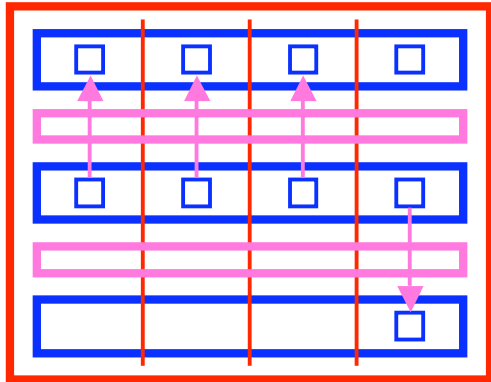Universität Zürich

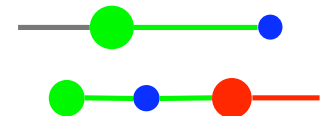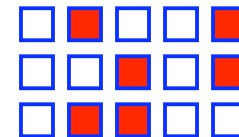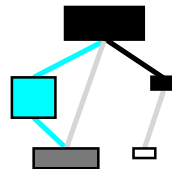http://seal.ifi.unizh.ch

Universität Zürich

s.e.a.l.
software evolution & architecture lab
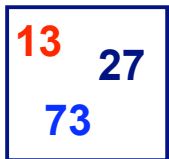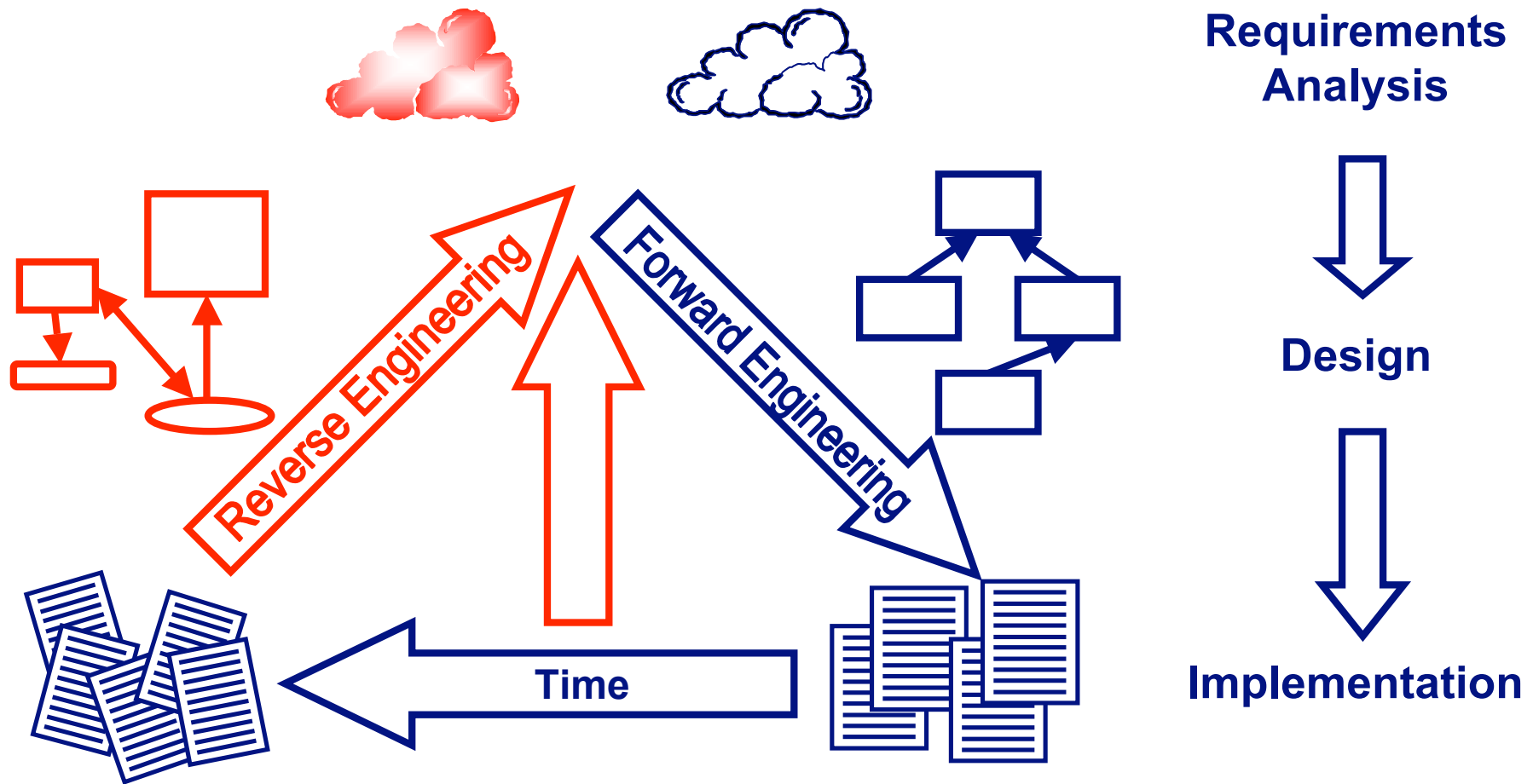
# Modeling History to Understand Software Evolution

**© 2007, Tudor Gîrba**

# Context: Reverse engineering is creating high level views of the system



Reverse Engineering

Forward Engineering

Time

Requirements Analysis

Design

Implementation

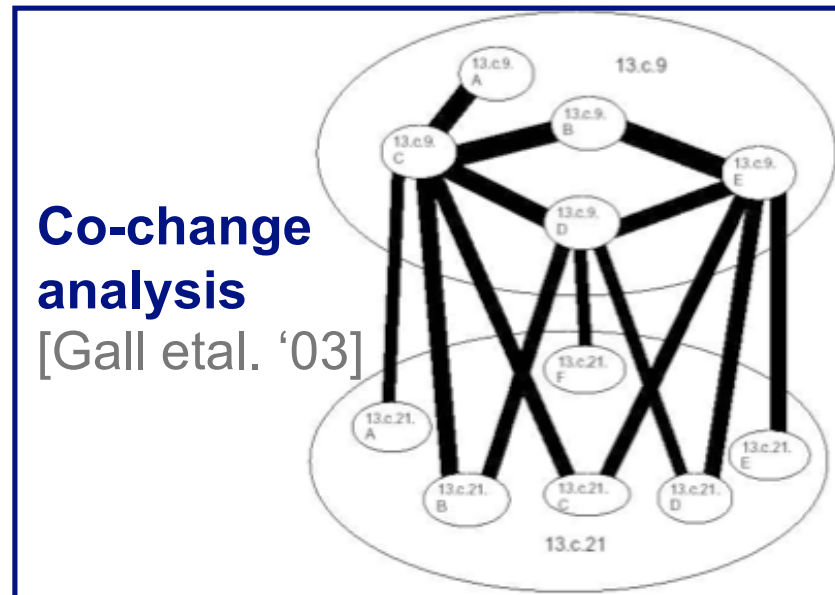# Context: History holds useful information for reverse engineering

**The doctor always looks at my health file**

**Historical information is useful but, it is hidden among** huge amounts of data



Version 1   Version 2   Version 3   …   Version n

N versions means N times more data

**The more data the more techniques are needed to analyze it**

# Context: Many techniques were developed



**Trend analysis**
[Lehman et al. '01]

**Evolution patterns**
[Lanza, Ducasse '02]

**Authors analysis** [Eick et al. '02]

**Co-change analysis**
[Gall etal. '03]

# Problem: Current approaches rely on ad-hoc models or on too specific meta-models

**Modules**
**OS/360-370**
**Growth Trend**

8000
6000
4000
2000
0

**Trend analysis**

**Evolution patterns**

FIRST VERSION
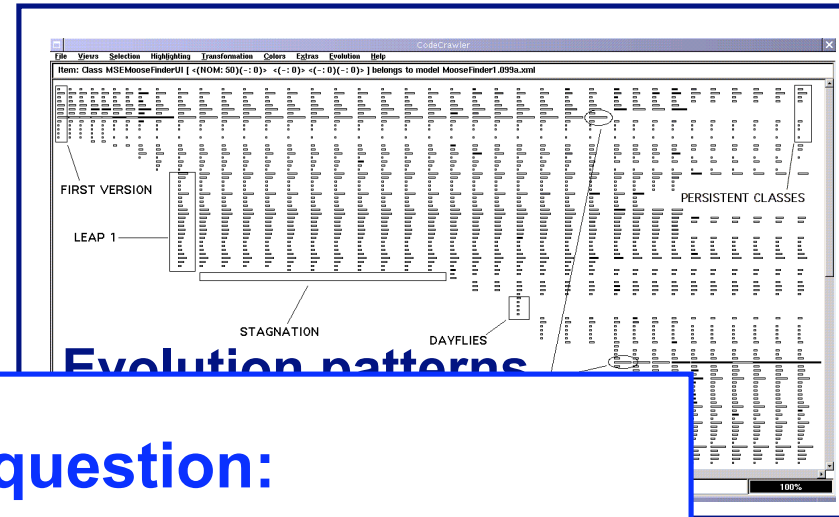LEAP 1
PERSISTENT CLASSES
STAGNATION
DAYFLIES

## Research question:

## How can we build a generic meta-model?

**Co-change analysis** [Gall etal. '03]

**Authors analysis** [Eick etal. '02]

# Overview



Hismo

Application

13  27
73

Historical measurements

**Hismo:**
**Modeling History**

History — Version

History — Version

History — Version

Ownership Map

# Example: Evolution Matrix reveals different evolution patterns

[Lanza, Ducasse '02]

**Pulsar Class**

**Idle Class**

**White D Class**

**Supern Class**

**Thesis:**

**Evolution needs to be modeled as a first class entity**

metric view

A

← versions →

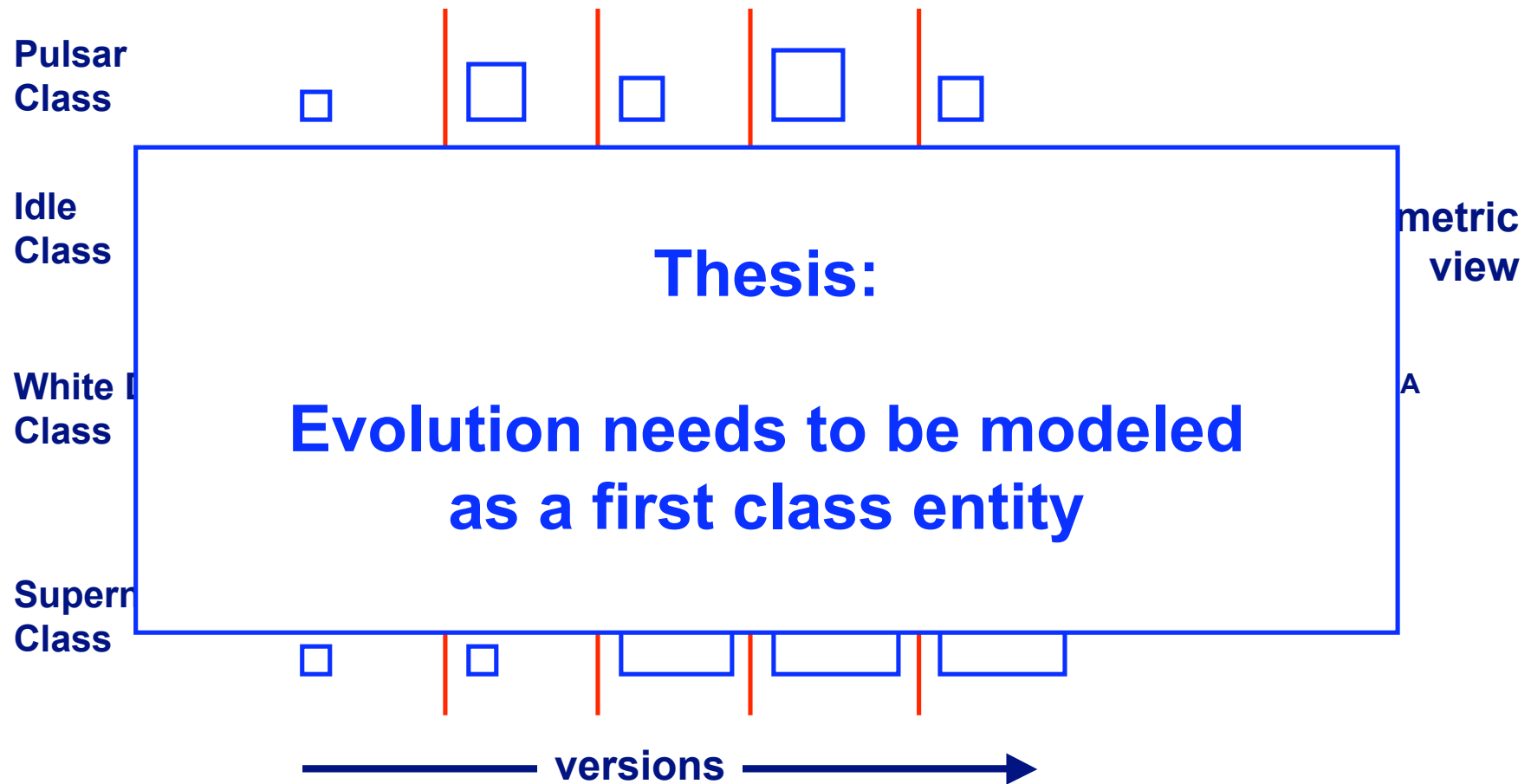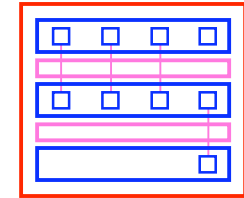# Solution: History encapsulates and characterizes the evolution

**Pulsar Class History**

**Idle Class History**

**White Dwarf Class History**

**Supernova Class History**

**ClassHistory**

**isPulsar**
**isIdle**
**…**

**versions**

# Hismo: The history meta-model

# … but, what about relationships?

# Hismo is obtained by transforming the structural meta-model

# Overview

**Application:**
**History measurements**

13
27
73

Hismo

Applications

on

on

13
27
73

**Historical measurements**

**Yesterday's Weather**

**History-based Detection Strategies**

**Hierarchy evolution**

**Co-change patterns**

**Ownership Map**

# Problem: History holds useful information hidden among large amounts of data



**How much was a class changed?**
**When was a class changed?**
**…**

# History can be measured:
## How much was a class changed?

**Evolution of Number of Methods**

$$ENOM(C)= \sum_{i=2}^{n} |NOM_i(C)-NOM_{i-1}(C)|$$

$$ENOM(C)= \quad 4 \quad + \quad 2 \quad + \quad 1 \quad + \quad 0 \quad = 7$$

| 1 | 5 | 3 | 4 | 4 |
|---|---|---|---|---|

© Tudor Gîrba

# History can be measured: When was a class changed?

**Latest Evolution of Number of Methods**

$$\text{LENOM}(C) = \sum_{i=2}^{n} |\text{NOM}_i(C) - \text{NOM}_{i-1}(C)| \, 2^{i-n}$$

**Earliest Evolution of Number of Methods**

$$\text{EENOM}(C) = \sum_{i=2}^{n} |\text{NOM}_i(C) - \text{NOM}_{i-1}(C)| \, 2^{2-i}$$

$$\text{LENOM}(C) = 4 \cdot 2^{-3} + 2 \cdot 2^{-2} + 1 \cdot 2^{-1} + 0 \cdot 2^{0} = 1$$



| 1 | 5 | 3 | 4 | 4 |

$$\text{EENOM}(C) = 4 \cdot 2^{0} + 2 \cdot 2^{-1} + 1 \cdot 2^{-2} + 0 \cdot 2^{-3} = 5.125$$

© Tudor Gîrba

# History measurements compress aspects of the evolution into numbers

|   |   |   |   |   |   | ENOM | LENOM | EENOM |
|---|---|---|---|---|---|------|-------|-------|
| A | 2 | 4 | 3 | 5 | 7 | 7 | 3.37 | 3.25 |
| B | 2 | 2 | 3 | 4 | 9 | 7 | 5.75 | 1.37 |
| C | 2 | 2 | 1 | 2 | 3 | 3 | 1 | 2 |
| D | 2 | 2 | 2 | 2 | 2 | 0 | 0 | 0 |
| E | 1 | 5 | 3 | 4 | 4 | 7 | 1 | 5.12 |

# History measurements compress aspects of the evolution into numbers

|   |   | ENOM | LENOM | EENOM |
|---|---|------|-------|-------|
| A | Balanced changer | 7 | 3.37 | 3.25 |
| B | Late changer | 7 | 5.75 | 1.37 |
| C |  | 3 | 1 | 2 |
| D | Dead stable | 0 | 0 | 0 |
| E | Early changer | 7 | 1 | 5.12 |

**Many measurements can be defined at different levels of abstraction …**

**Evolution**

**Latest/Earliest Evolution**

**Stability**

**Historical Max/Min**          **of**

**Historical Average**

**Growth Trend**

**…**

**Number of Methods**

**Number of Statements**
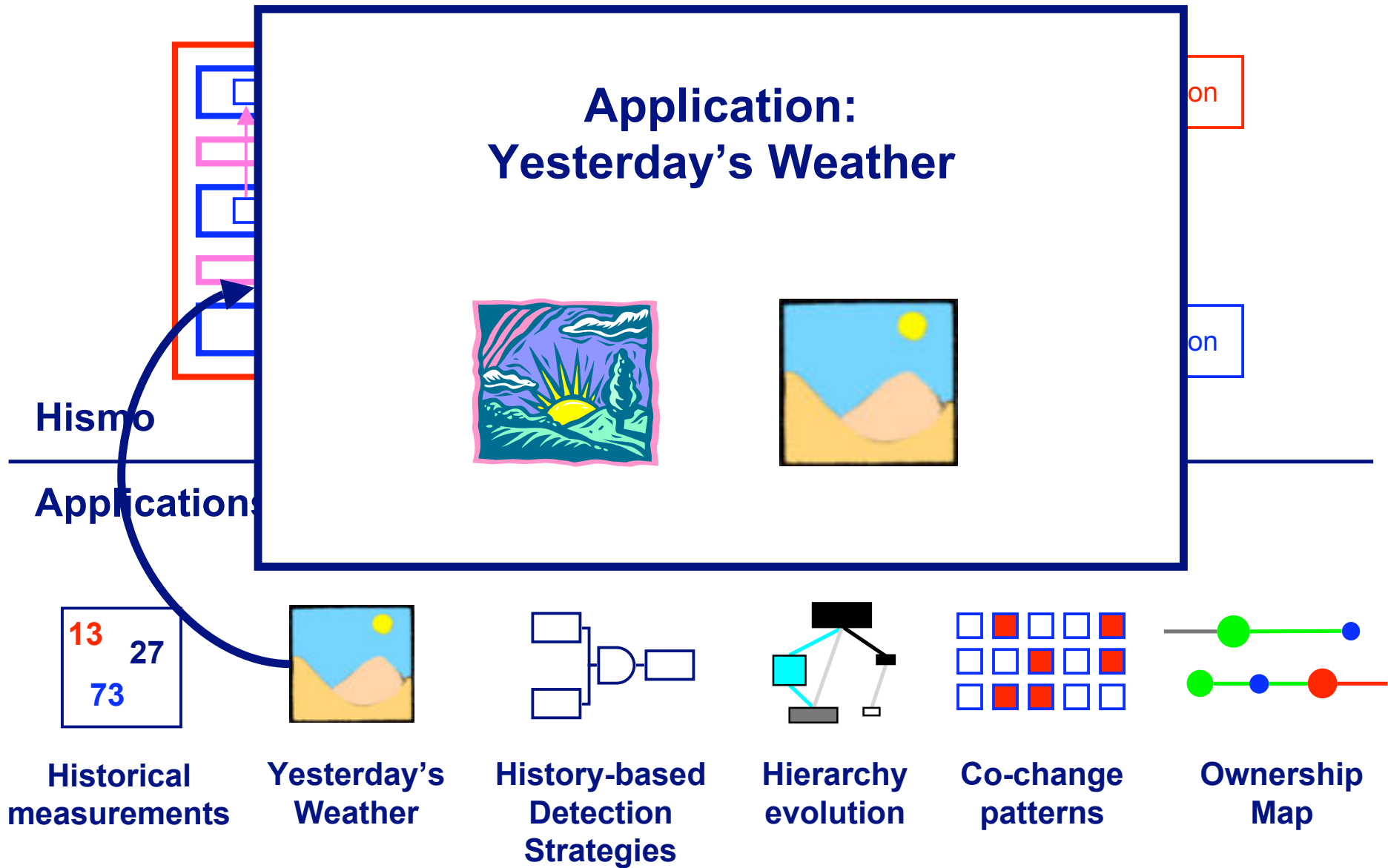
**Cyclomatic Complexity**

**Lines of Code**

**Number of Classes**

**Number of modules**

**…**

**… But measurements are a means not a goal**

# Overview



**Application:**
**Yesterday's Weather**

**Hismo**

**Applications**

| **13** **27** **73** | **Yesterday's Weather** | **History-based Detection Strategies** | **Hierarchy evolution** | **Co-change patterns** | **Ownership Map** |
|---|---|---|---|---|---|
| **Historical measurements** | | | | | |

# Common Wisdom: The recently changed parts are likely to change in the near future



[Mens,Demeyer '01]

**Is the common wisdom relevant?**

## Yesterday's Weather metaphor:

It expresses the chances of having the same weather today as we had yesterday
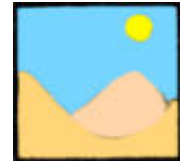
It is location specific

 **Switzerland - 30%**

 **Sahara - 90%**

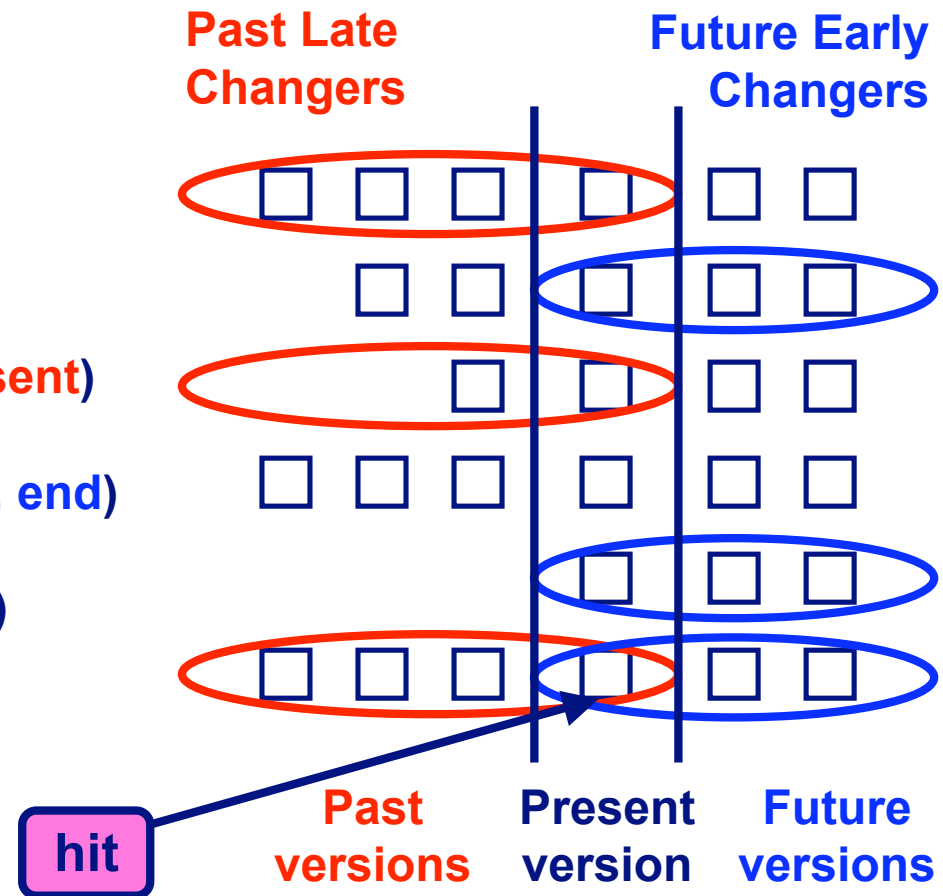# Yesterday's Weather: For each given version we check the common wisdom

**Past Late Changers**

**Future Early Changers**
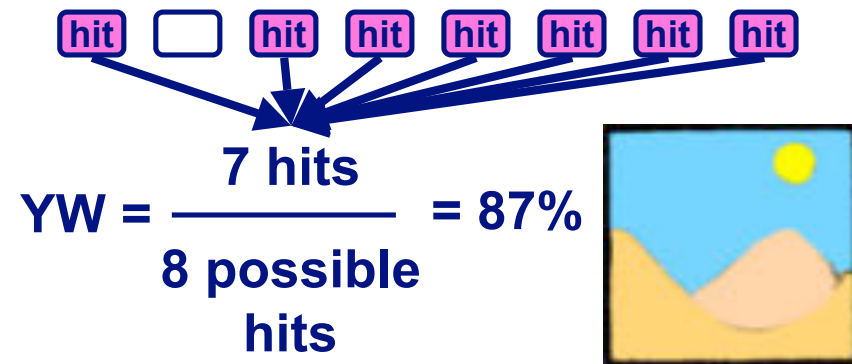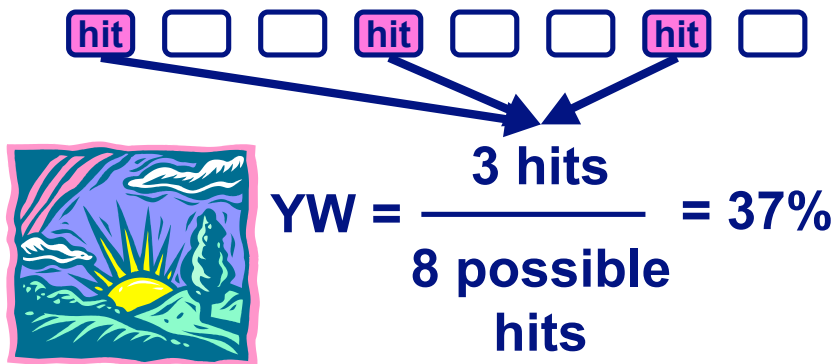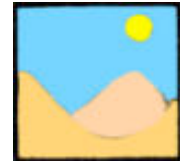
**YesterdayWeatherHit(present):**

**past**:=histories.top**LENOM**(**start, present**)

**future**:=histories.top**EENOM**(**present, end**)

**past**.intersectWith(**future**).notEmpty()

**hit**

**Past versions**

**Present version**

**Future versions**

# Overall Yesterday's Weather shows the localization of changes in time

hit · · hit · · hit ·

$$YW = \frac{3\ hits}{8\ possible\ hits} = 37\%$$

hit · hit hit hit hit hit hit

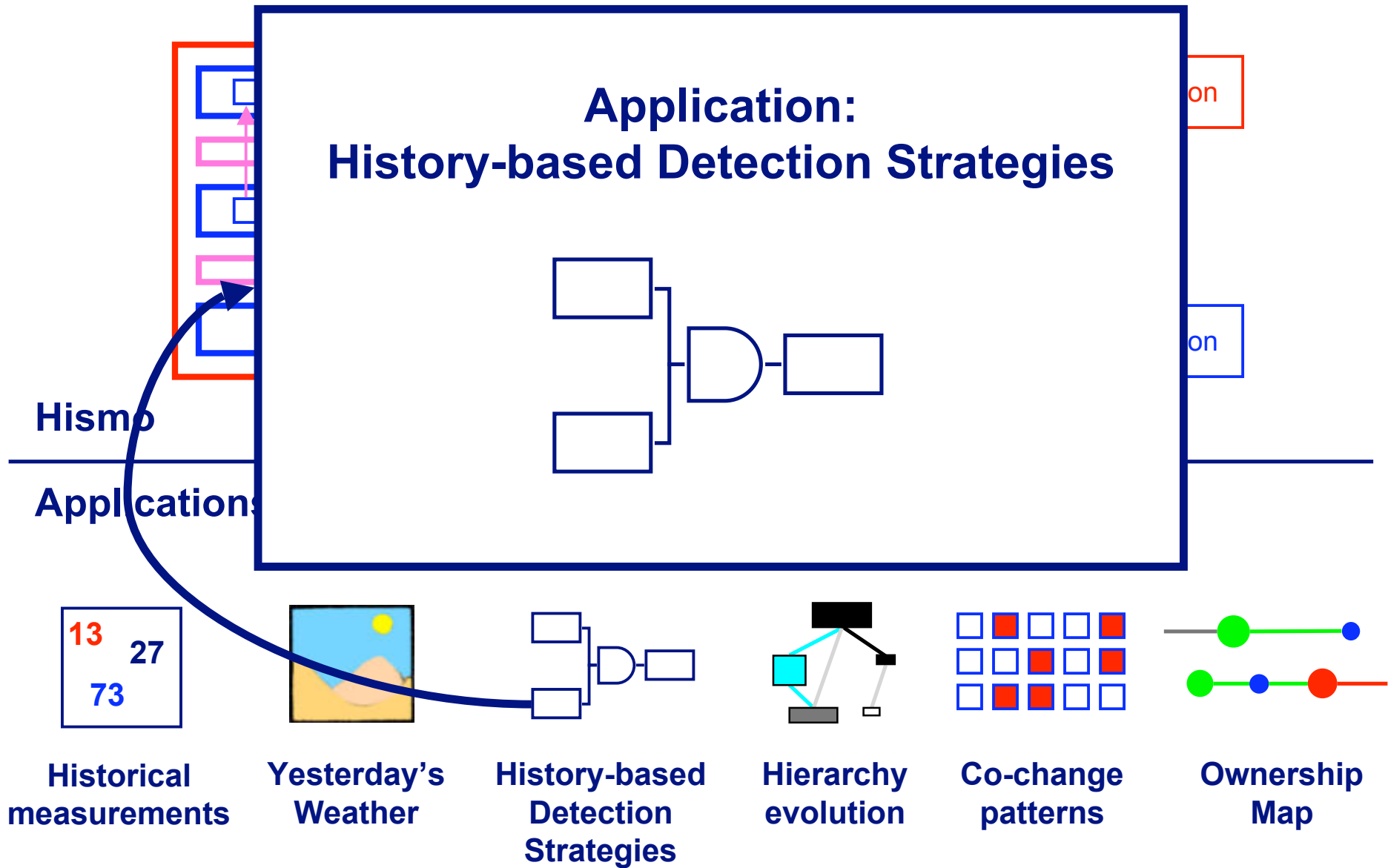$$YW = \frac{7\ hits}{8\ possible\ hits} = 87\%$$

## Case studies:

40 versions of CodeCrawler (180 classes): 100%

40 versions of Jun (700 classes): 79%

40 versions of Jboss (4000 classes): 53%

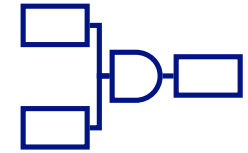# Overview

**Application:
History-based Detection Strategies**

**Hismo**

**Applications**

**13** **27**
**73**

**Historical
measurements**

**Yesterday's
Weather**

**History-based
Detection
Strategies**

**Hierarchy
evolution**
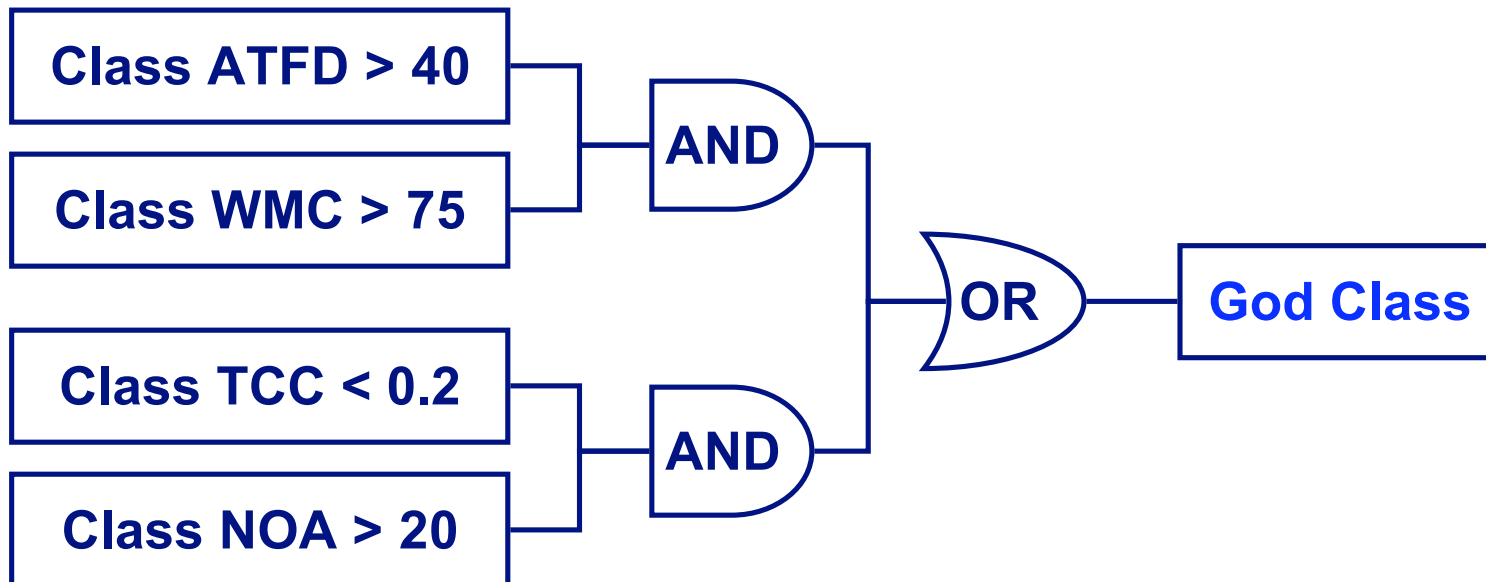
**Co-change
patterns**

**Ownership
Map**

# Context: Detection Strategies detect design flaws based on measurements [Marinescu '04]
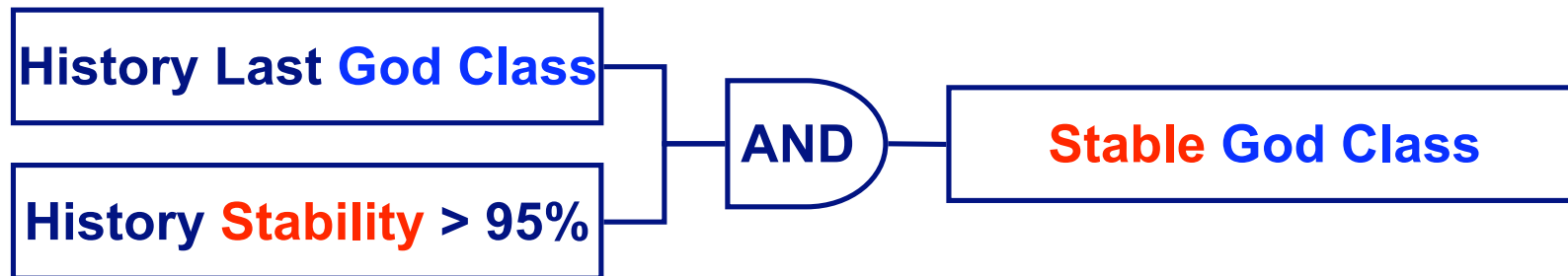
## Example: God Class

Maintainability problem because it encapsulates a lot of knowledge



Class ATFD > 40
Class WMC > 75
AND

Class TCC < 0.2
Class NOA > 20
AND

OR

God Class

# History-based Detection Strategies take evolution into account

**Example: a Stable God Class is not necessarily a bad one**



```
┌─────────────────────────────┐
│ History Last God Class      │───┐
└─────────────────────────────┘   │
                                 ( AND )───  Stable God Class
┌─────────────────────────────┐   │
│ History Stability > 95%     │───┘
└─────────────────────────────┘
```
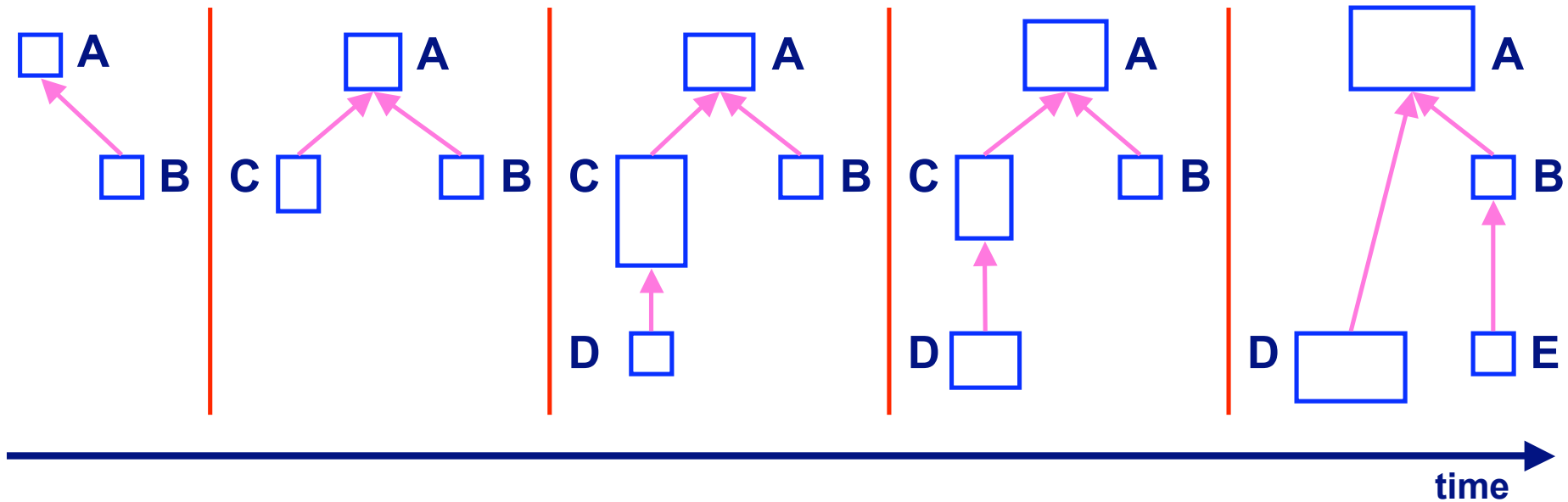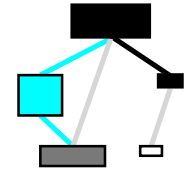
**Case study: 5 out of 24 God Classes in Jun were stable and harmless**

# Overview

**Hismo**

**Applications**

Application:
Characterizing the evolution
of class hierarchies

**Historical measurements**

13  27
73

**Yesterday's Weather**

**History-based Detection Strategies**

**Hierarchy evolution**

**Co-change patterns**

**Ownership Map**

# Context: Given the evolution of a hierarchy …



A is persistent     C was removed

B is stable     E is newborn

D inherited from C and then from A     …

# … but useful information is hidden among large amounts of data



## How were the hierarchies evolved?

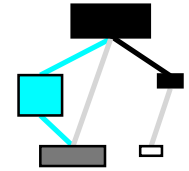# Hierarchy Evolution Complexity View characterizes class hierarchy histories

A

C    B

D    E

ENOM

ENOS

| Age |
| :---: |
| Removed |

Class History

Age

Removed

Inheritance History

A is persistent        C was removed

B is stable            E is newborn

D inherited from C and then from A        …

# Case study: Class hierarchies in Jun reveal evolution patterns

**Persistent**
**Unbalanced**
**Stable**
**Reliable inheritance**

**Old**
**Stable**
**Balanced**
**Reliable inheritance**
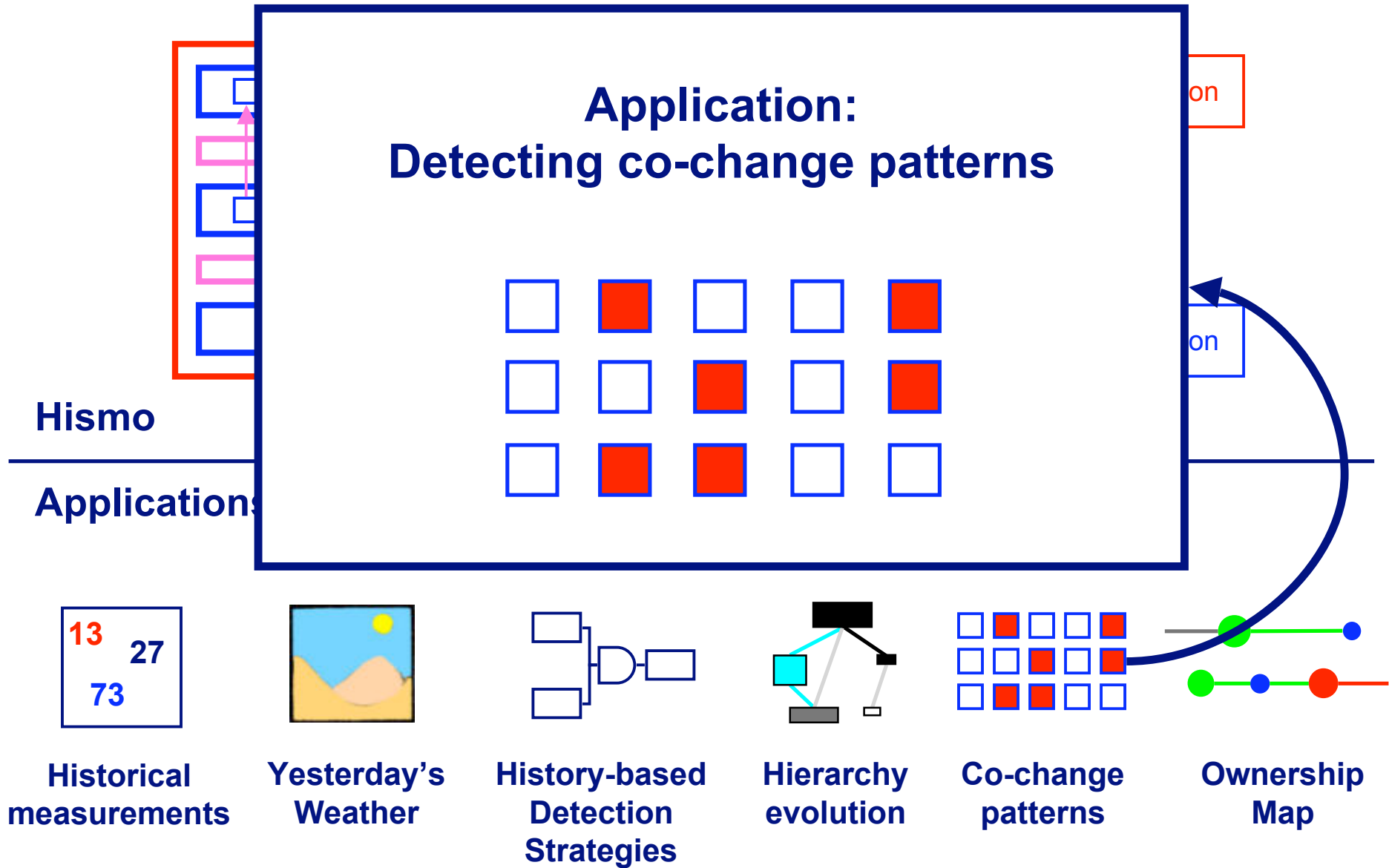
**Young**
**Unstable root**
**Reliable inheritance**

**Newborn**

**Old**
**Unstable**
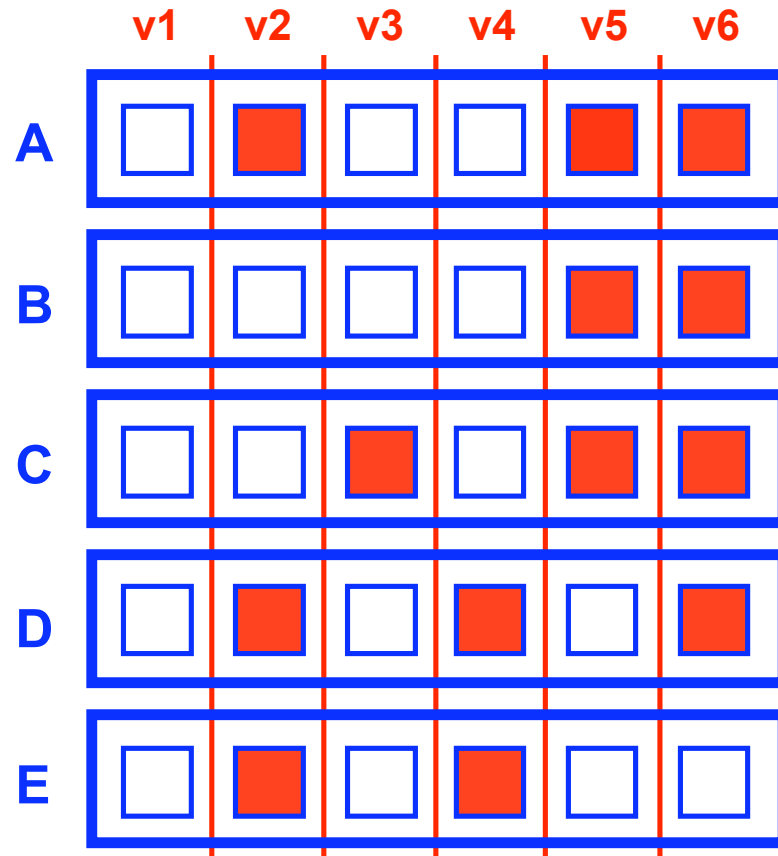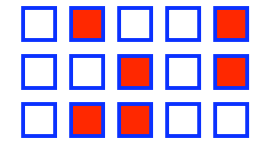**Unbalanced**
**Unreliable inheritance**

# Overview

**Application:**
**Detecting co-change patterns**

**Hismo**

**Applications**

**13  27  73**

**Historical measurements**

**Yesterday's Weather**

**History-based Detection Strategies**

**Hierarchy evolution**

**Co-change patterns**

**Ownership Map**

# Context: Repeated co-changes reveal hidden dependencies [Gall etal. '98]

|   | v1 | v2 | v3 | v4 | v5 | v6 |
|---|----|----|----|----|----|----|
| **A** | ☐ | 🟥 | ☐ | ☐ | 🟥 | 🟥 |
| **B** | ☐ | ☐ | ☐ | ☐ | 🟥 | 🟥 |
| **C** | ☐ | ☐ | 🟥 | ☐ | 🟥 | 🟥 |
| **D** | ☐ | 🟥 | ☐ | 🟥 | ☐ | 🟥 |
| **E** | ☐ | 🟥 | ☐ | 🟥 | ☐ | ☐ |

**Can we identify co-change patterns like:**

    **Parallel Inheritance**

    **Shotgun Surgery**

    **…**

**?**

# Formal Concept Analysis (FCA) finds elements that have properties in common

[Ganter, Wille '99]



To use FCA, we need to map our interests on **elements** and **properties**

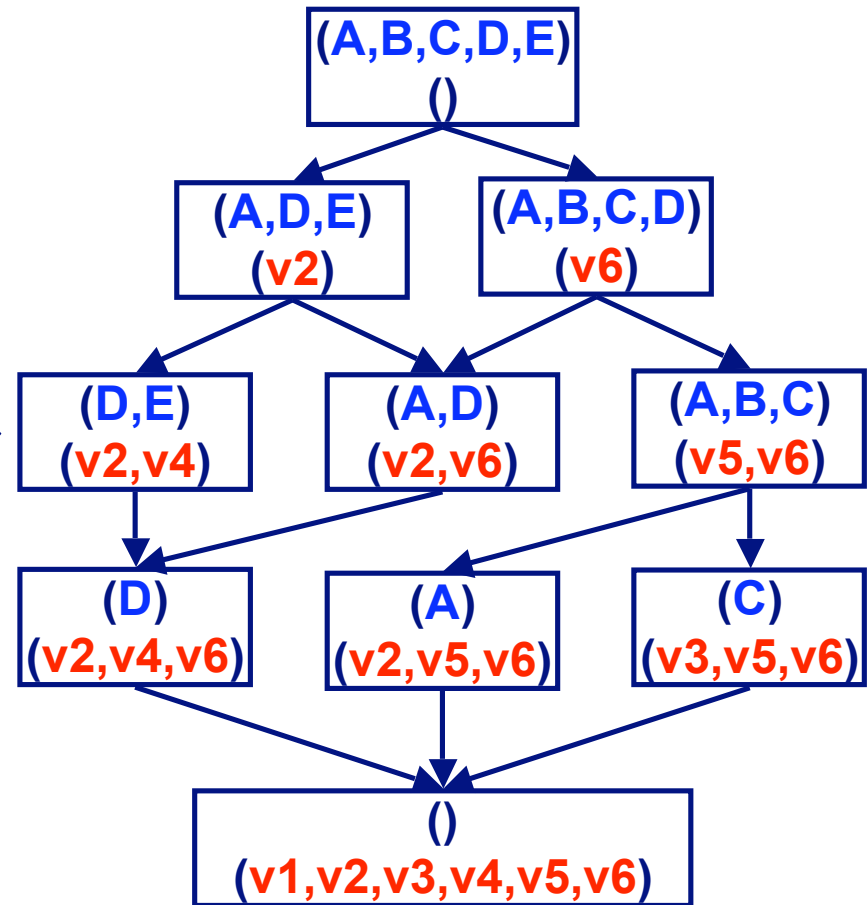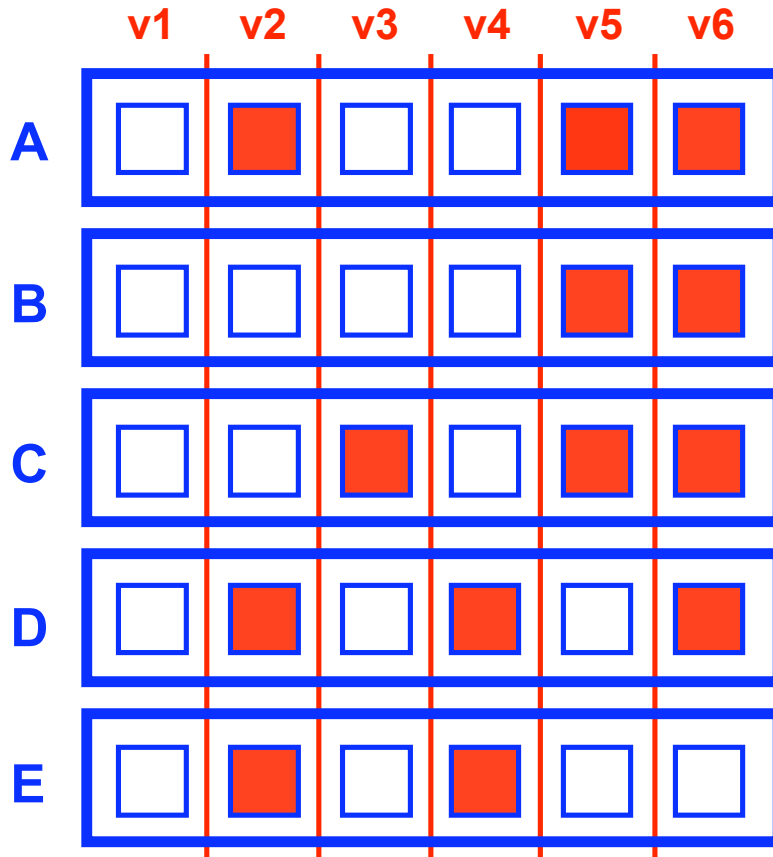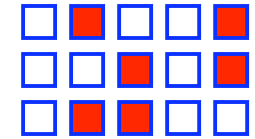# Formal Concept Analysis (FCA) finds elements that have properties in common

[Ganter, Wille '99]



**To use FCA, we need to map our interests on elements and properties**

# We use FCA to identify entities that co-changed repeatedly

|   | v1 | v2 | v3 | v4 | v5 | v6 |
|---|----|----|----|----|----|----|
| A | ☐ | 🟥 | ☐ | ☐ | 🟥 | 🟥 |
| B | ☐ | ☐ | ☐ | ☐ | 🟥 | 🟥 |
| C | ☐ | ☐ | 🟥 | ☐ | 🟥 | 🟥 |
| D | ☐ | 🟥 | ☐ | 🟥 | ☐ | 🟥 |
| E | ☐ | 🟥 | ☐ | 🟥 | ☐ | ☐ |

**FCA** →

(A,B,C,D,E)
()

(A,D,E)
(v2)

(A,B,C,D)
(v6)

(D,E)
(v2,v4)

(A,D)
(v2,v6)

(A,B,C)
(v5,v6)

(D)
(v2,v4,v6)

(A)
(v2,v5,v6)

(C)
(v3,v5,v6)

()
(v1,v2,v3,v4,v5,v6)

**Elements = Histories**

**Properties = "changed in version X"**

# Example: Parallel inheritance denotes children added in several hierarchies

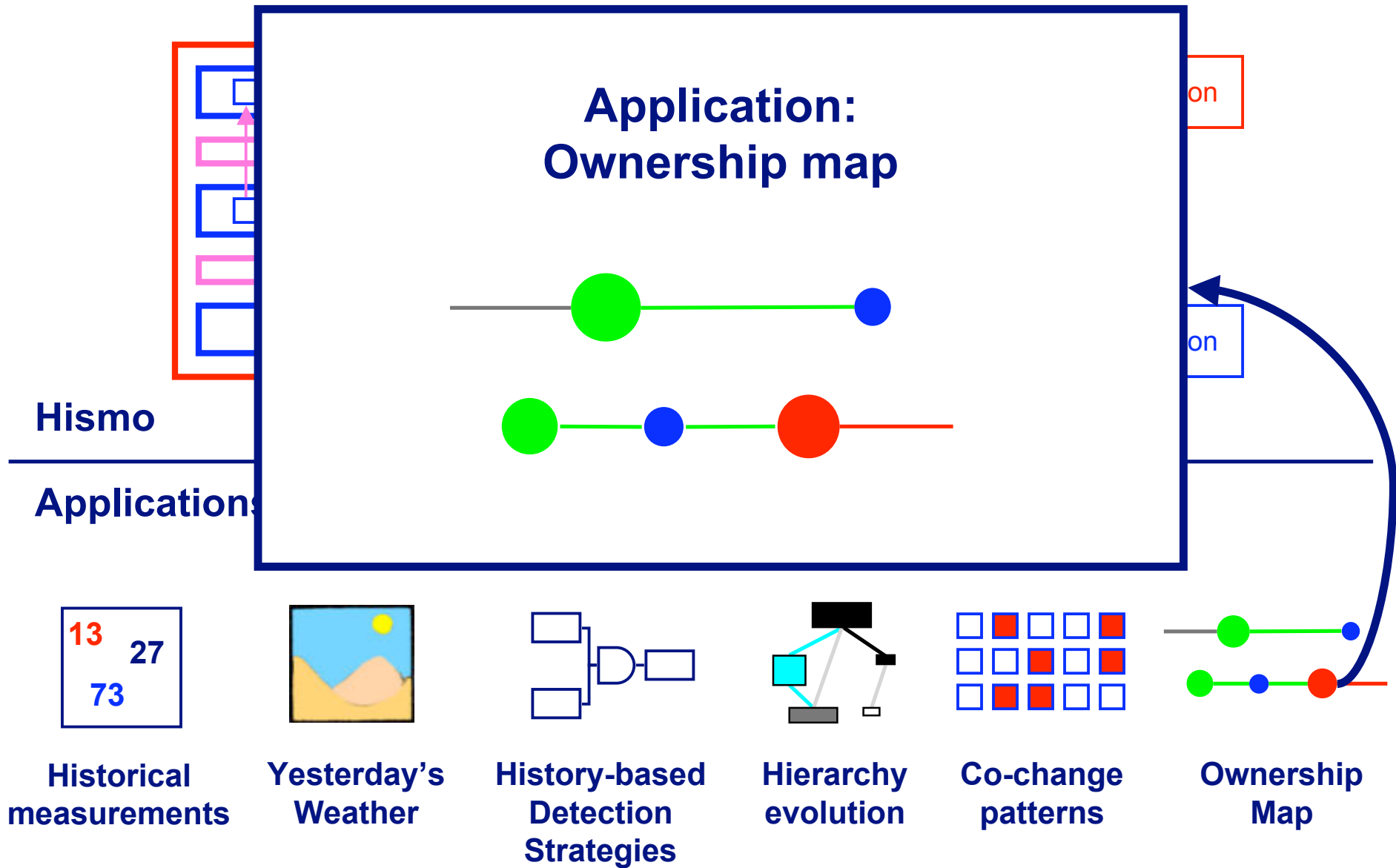**v1    v2    v3    v4    v5    v6**

A    | 0 | 1 | 1 | 1 | 2 | 4 |

**Case study: JBoss**

**ServiceMBeanSupport**
**JBossTestCase**                 **14 versions**

**EJBLocalHome**
**EJBLocalObject**                **9 versions**

A    A    A    A    A    A

**Elements = ClassHistories**
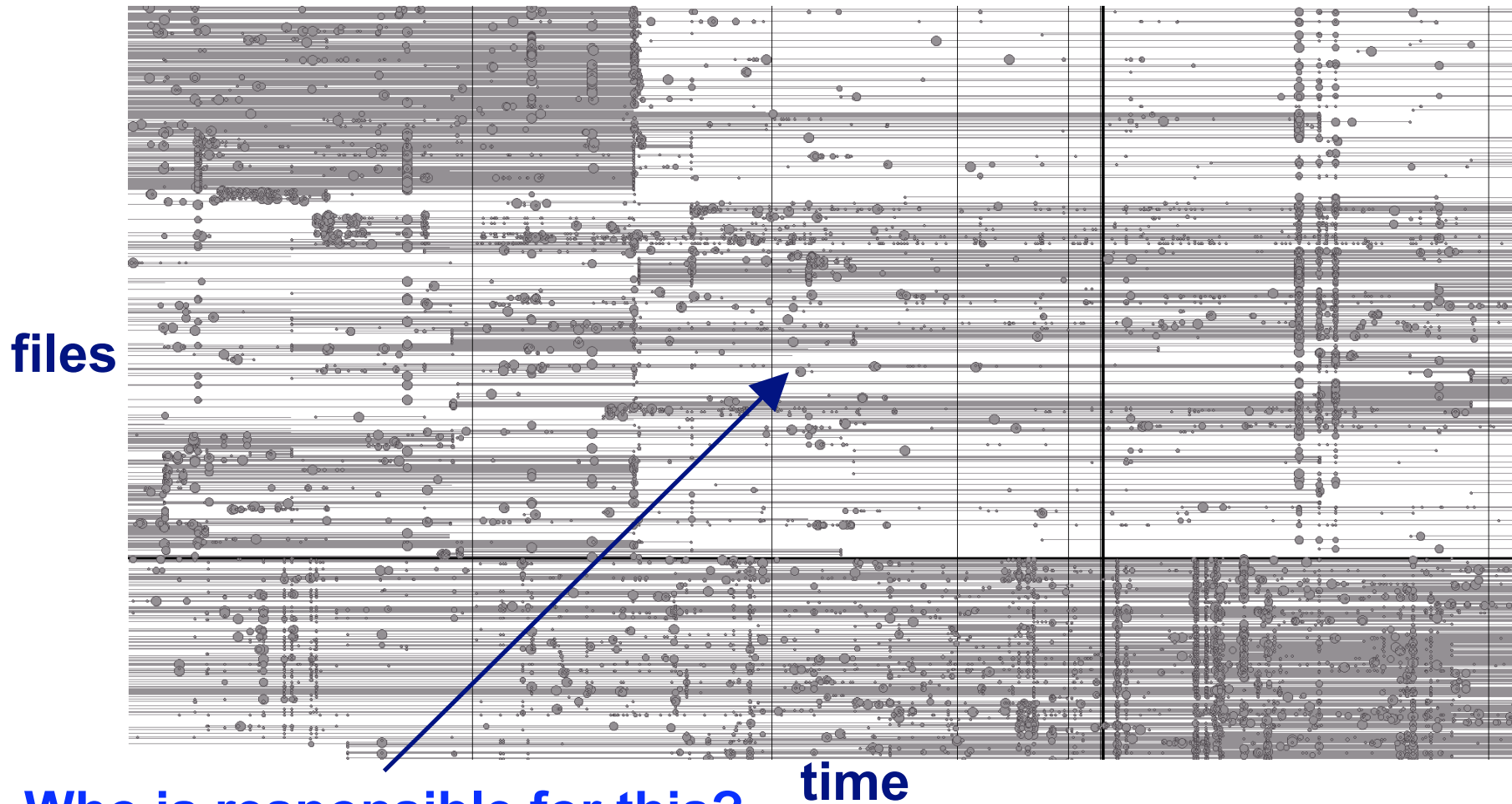**Properties = "changed number of children in version X"**

# Overview

**Application:
Ownership map**

**Hismo**

**Applications**

**Historical measurements**

13  27
73

**Yesterday's Weather**

**History-based Detection Strategies**

**Hierarchy evolution**

**Co-change patterns**

**Ownership Map**

© Tudor Gîrba

# Context: The code history might tell you what happened, but not why it happened

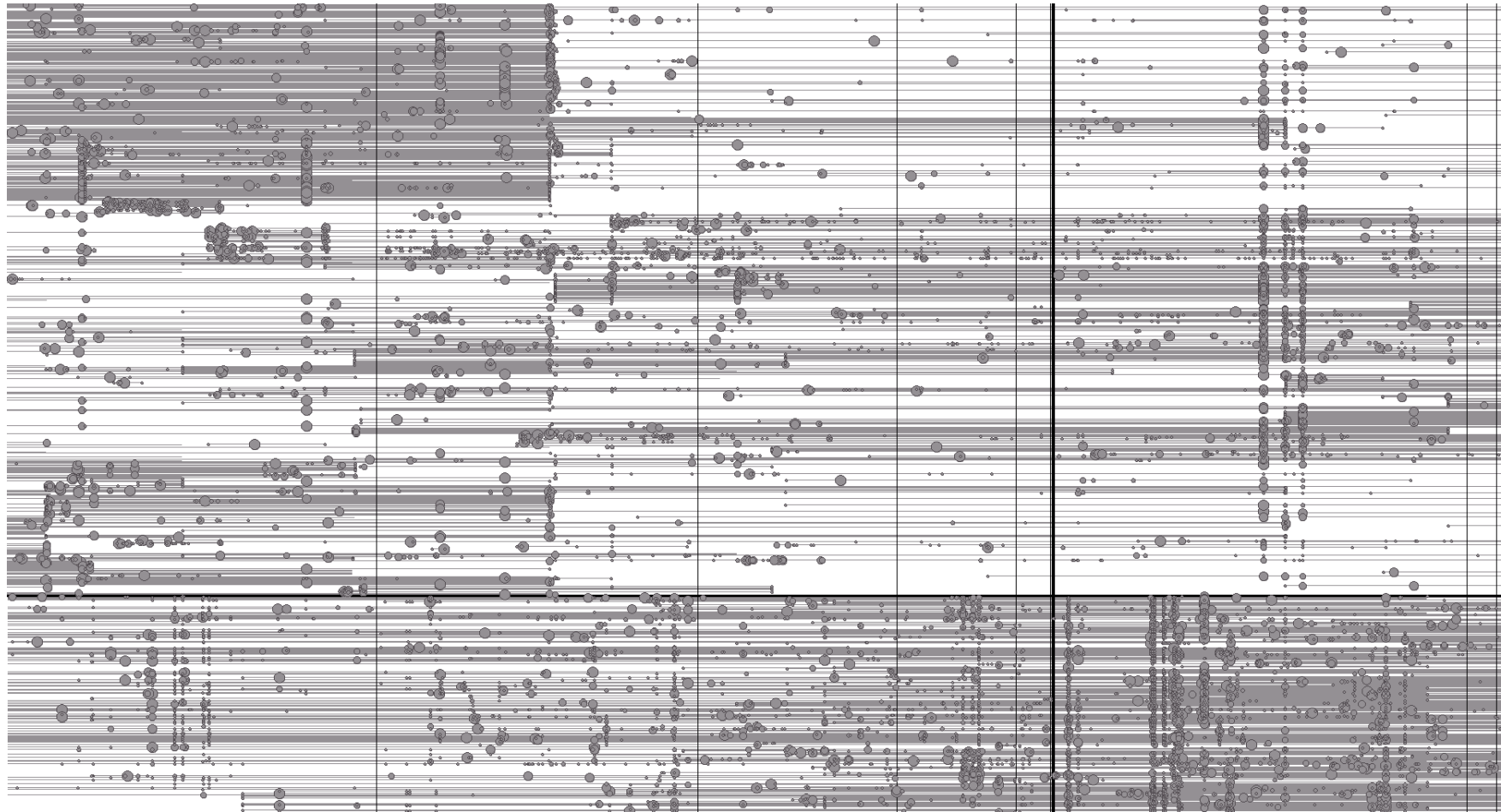**Case study: Outsight**                [Rysselberghe, Demeyer '04]



files

time

**Who is responsible for this?**

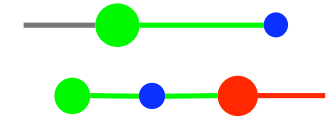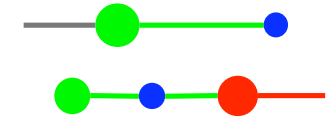# We color the lines to show which author owned which files in which period

**Green author large commit**

**Green author ownership**

**File History A**

**File History B**

**Blue author small commit**

# The commit history shows what happened

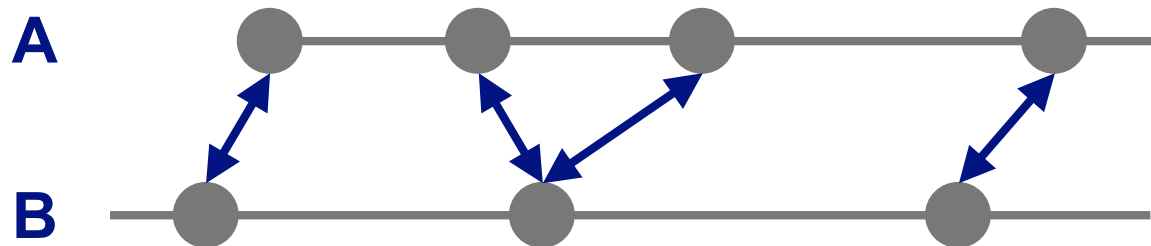# Ownership Map shows which author owned which files in which period

# We cluster the file histories to favor colored blocks inside each module
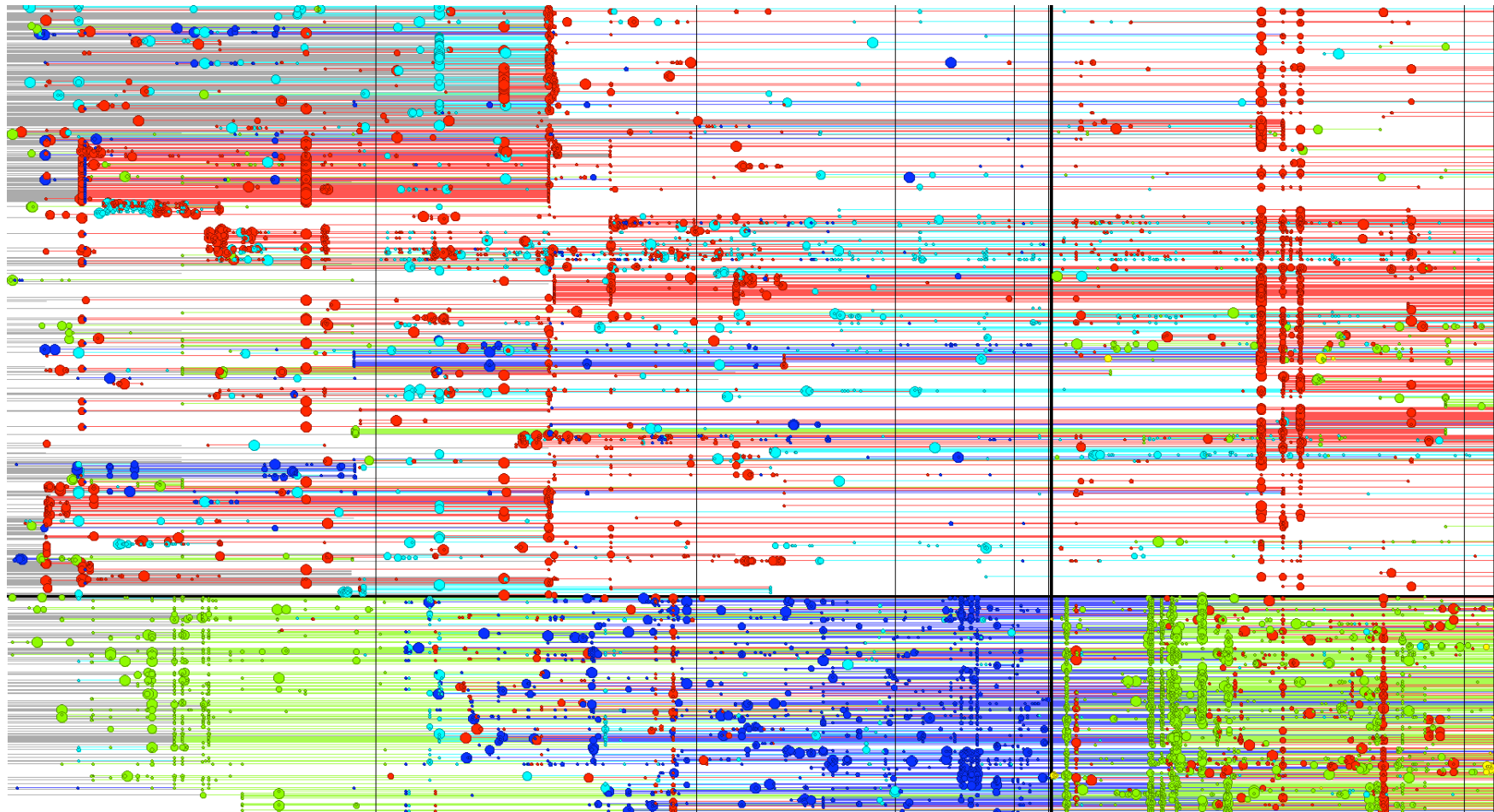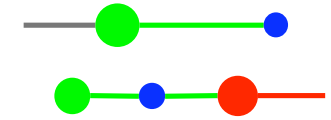
**We use the Hausdorf distance between the commit timestamps**

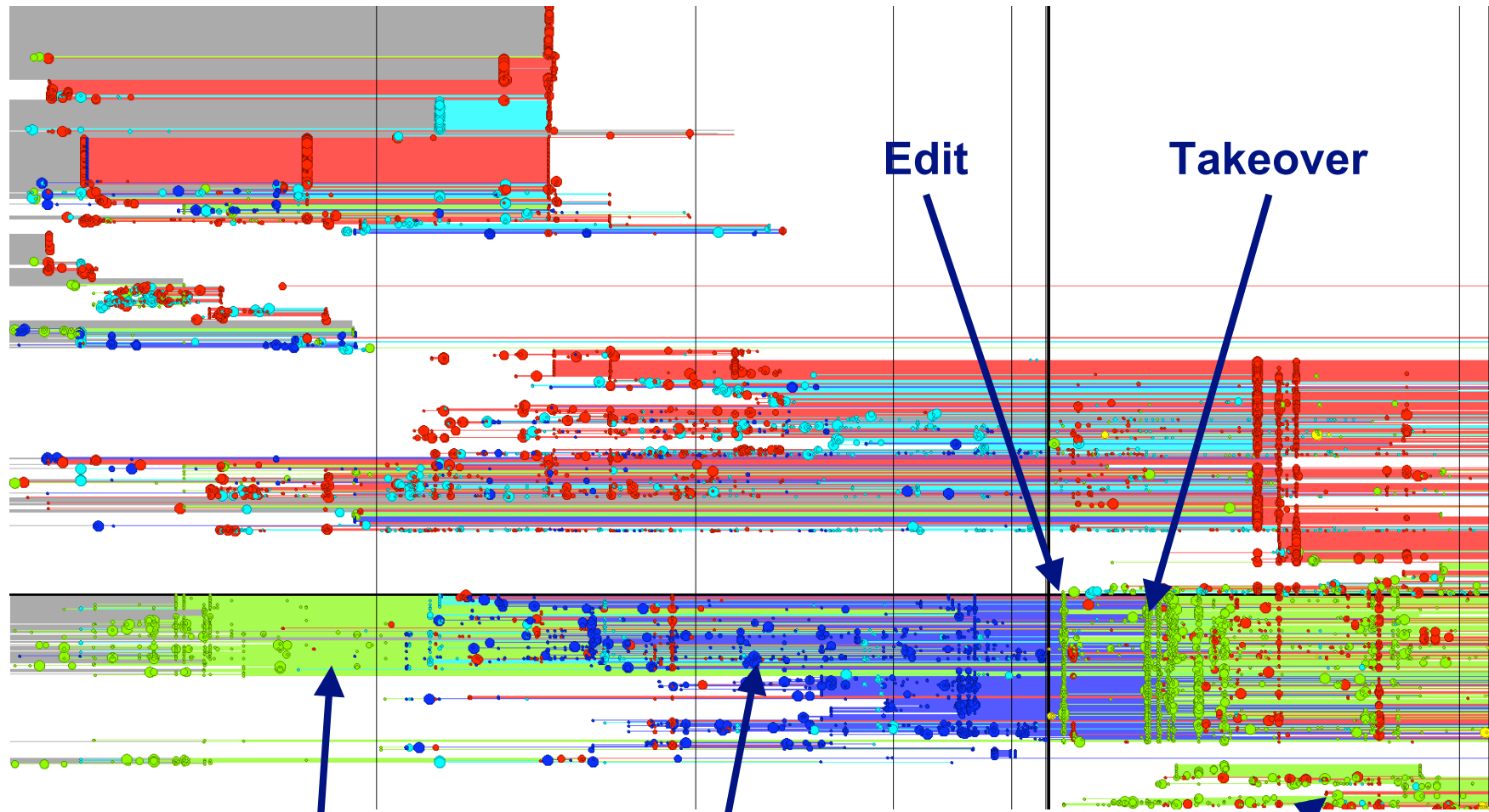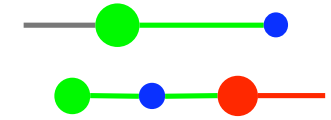$$d(A, B) = \sum_{a \in A} \min^2\{ \, | \, a - b \, | \; b \in B \, \}$$

# Ownership Map on alphabetically ordered files is not very useful, but …

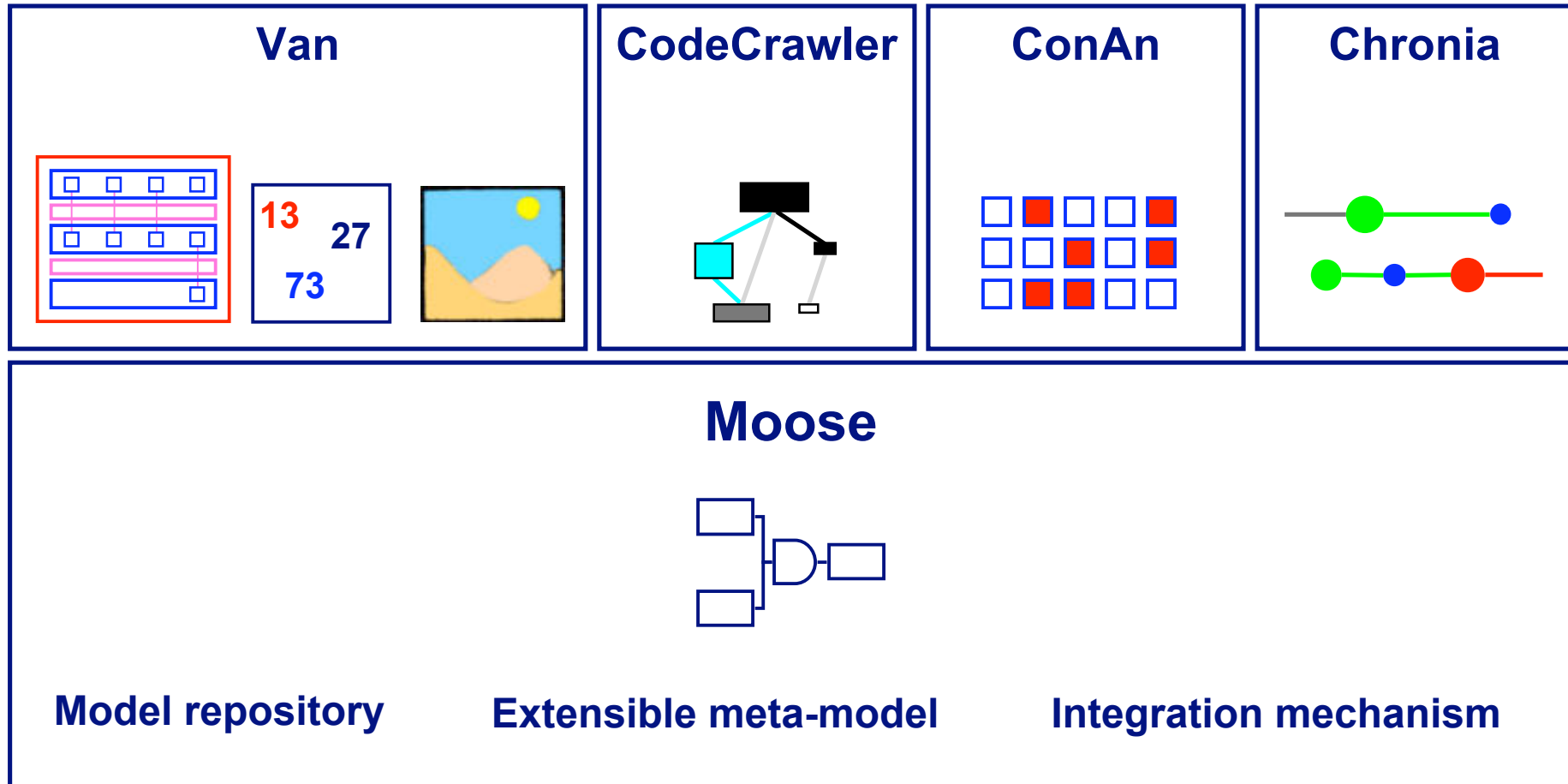# The ordered Ownership Map reveals developer patterns



Edit

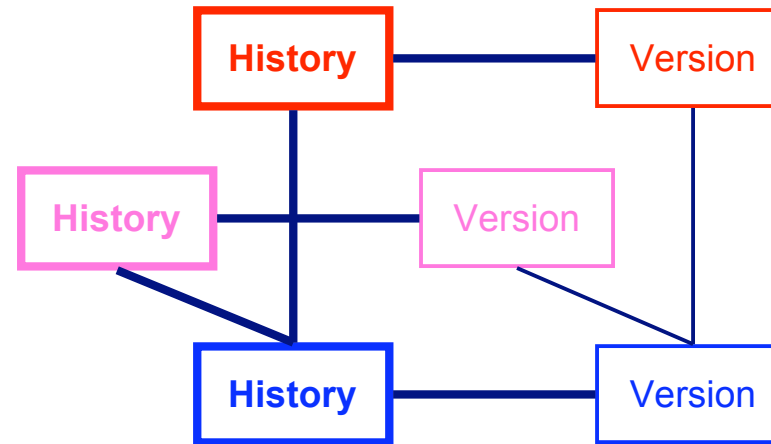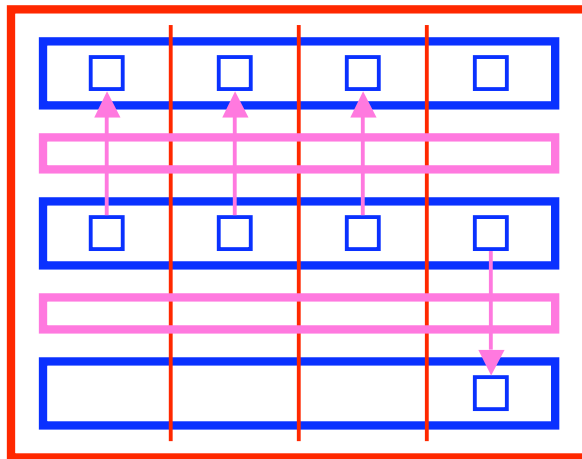Takeover

Monologue

Familiarization

Dialogue

© Tudor Gîrba

# Overview

**Implementation:**

**Both Hismo and its applications are implemented in one single infrastructure**

History

Version

**Hismo**

**Applications**

13 27
73

**Historical measurements**

**Yesterday's Weather**

**History-based Detection Strategies**

**Hierarchy evolution**

**Co-change patterns**

**Ownership Map**

# Implementation: All tools are integrated into Moose



Van · CodeCrawler · ConAn · Chronia

Moose

Model repository · Extensible meta-model · Integration mechanism
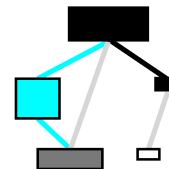
# Conclusion: Hismo offers a uniform way of expressing evolution analyses



**Hismo**

**Applications**

**Questions?**

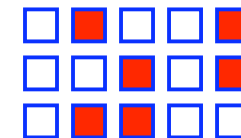**Historical measurements**
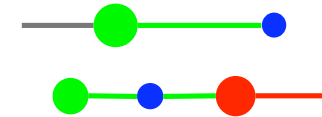
**Yesterday's Weather**

**History-based Detection Strategies**
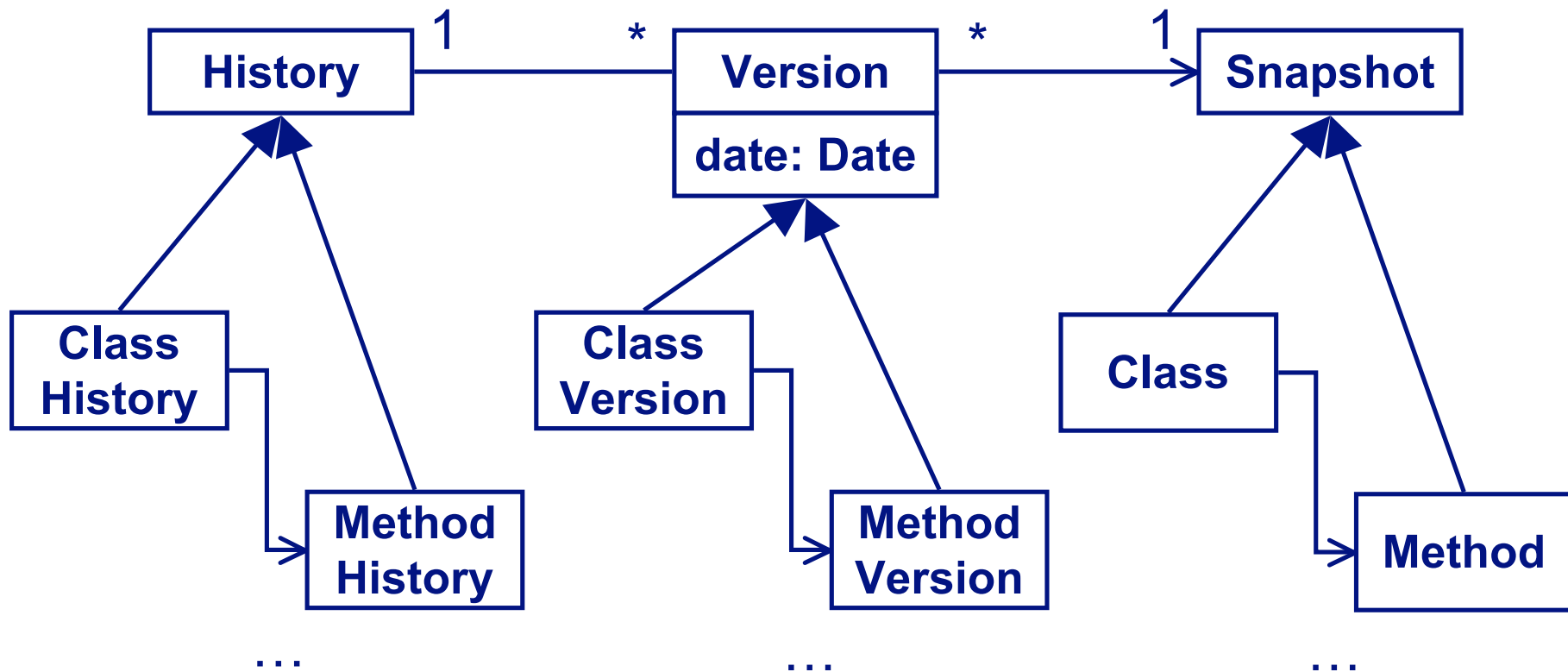
**Hierarchy evolution**
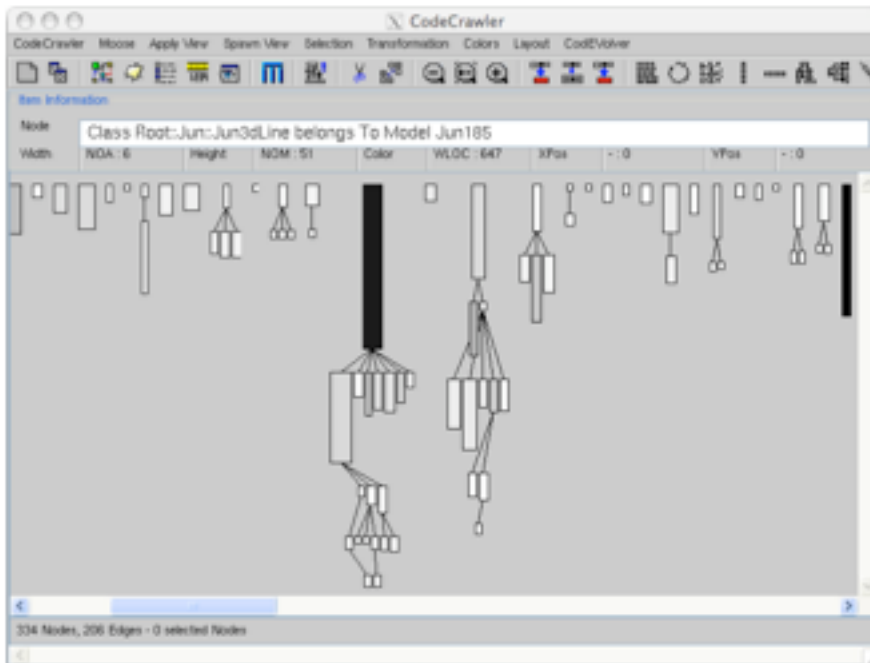
**Co-change patterns**

**Ownership Map**

# Hismo: History is a sequence of Versions, where a Version adds the notion of time to Snapshot
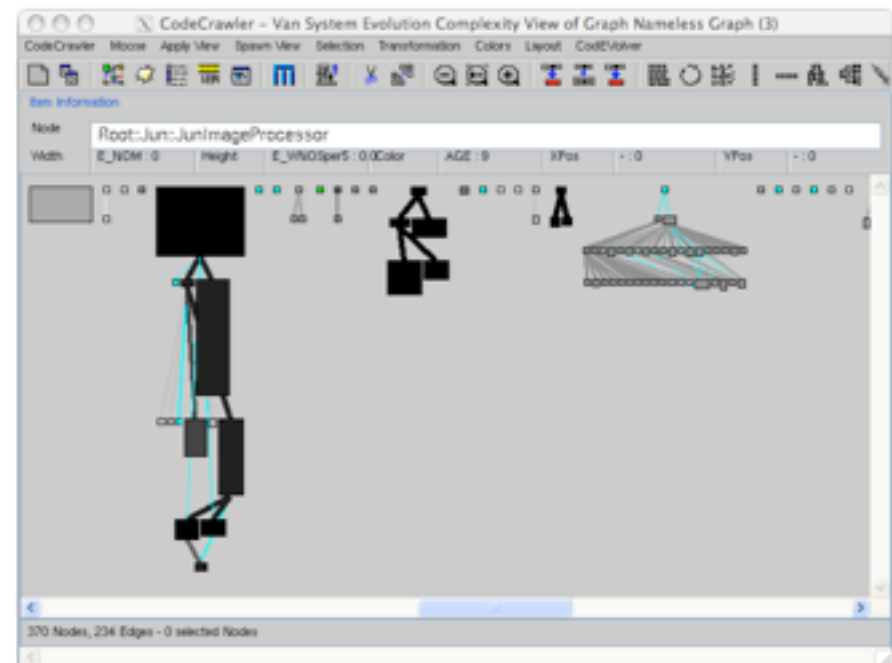
# The techniques are orthogonal to the type of data



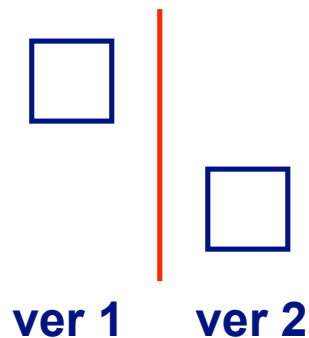**Node = class**

**Edge = inheritance**

**Node = Class History**

**Edge = Inheritance History**

# Entity identity: Are two entities at different points in time the versions of the same history?

**The current versioning systems record snapshots**

**What are the names?**

**What are the properties?**

ver 1     ver 2

## What we would like

Preserve the identity in the environment

Record changes as they happen