

New Frontiers of Reverse Engineering

Gerardo Canfora
Massimiliano Di Penta
FOSE / ICSE 2007

Reverse engineering is analyzing a subject system to:
identify components and their relationships, and
create more abstract representations.

Chikofky & Cross, 90

2

Why reverse engineer?

In 1944, 3 B-29 had to land in Russia



Requirement: Copy everything fast!



Tudor Girba

5

Approach: disassemble, run, test



Tudor Girba

6

TU-4 Result: 105,000 pieces
reassembled in 2 years



Tudor Girba

7

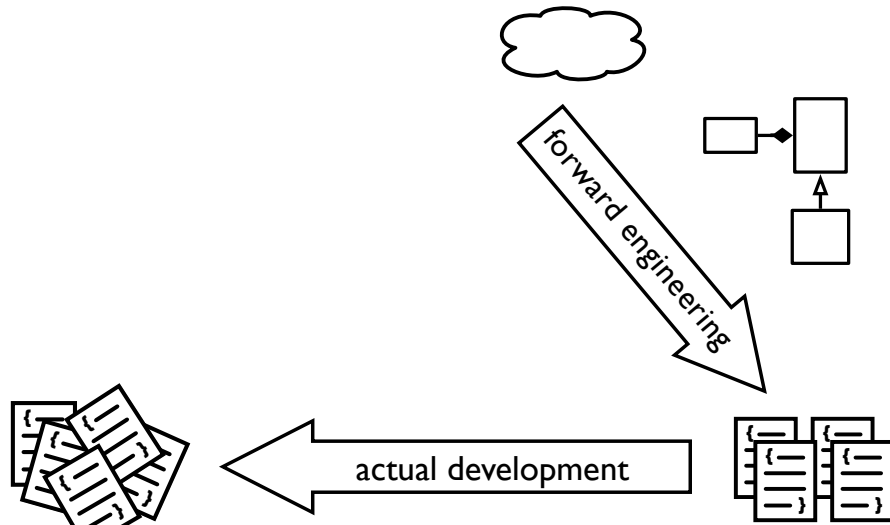
Reading code...

100'000 lines of code
* 2 = 200'000 seconds
/ 3600 = 56 hours
/ 8 = 7 days

Tudor Girba

8

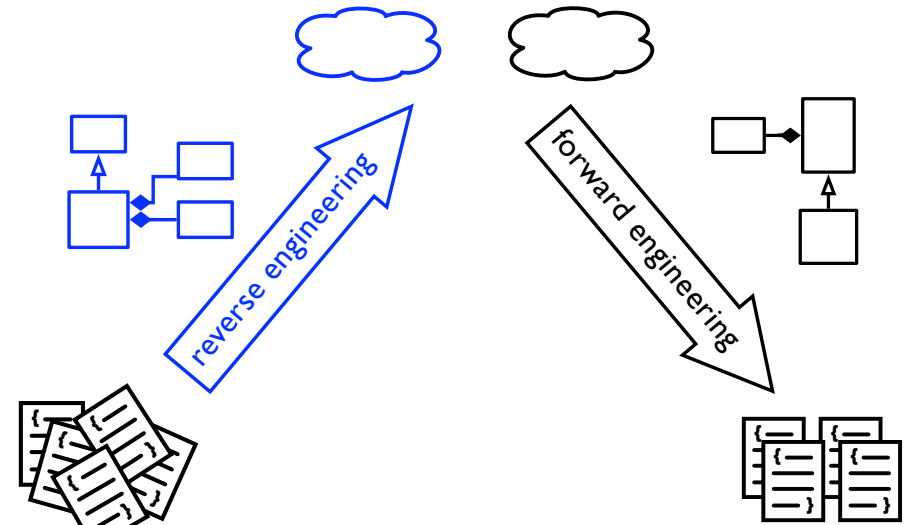
How development happens



Tudor Girba

9

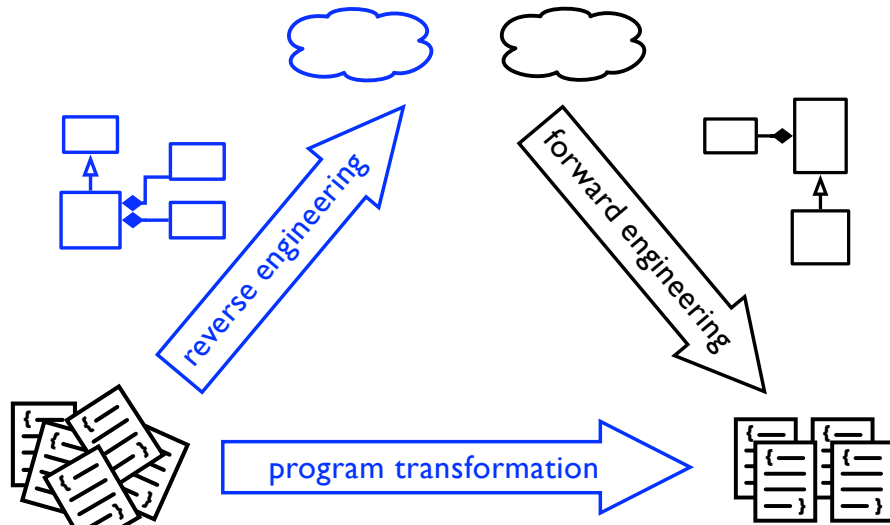
Reengineering life cycle



Tudor Girba

10

Reengineering life cycle



Tudor Girba

11

Where is the “real action” happening?



International Conference on Software Maintenance (25)



Intl. Workshop on Source Code Analysis and Manipulation (10)



Working Conference on Reverse Engineering (16)



International Conference on Program Comprehension (17)

Directions in Reverse Engineering

You got to be careful if you don't know where you're going, because you might not get there.

Yogi Berra



1 Software Understanding

2 Design Recovery

1 Software Understanding

Achievements

- Understanding and Migrating Procedural Code
- Create Models and Model Extractors/Parsers
- Clone Detection and Analysis
- Aspect Mining
- Visualizing Software Artifacts

1 Software Understanding

Achievement 1.1: Understanding and Migrating Procedural Code

- Cope with Y2K problem
- Created many intermediary representations
- Identify objects in legacy code
- Migration from procedural to object-oriented system
 - issue of changing the language without changing the paradigm

Achievements in
Mid '90s

1 Software Understanding

Achievement 1.2: Create Models and Model Extractors/Parsers

Parsing source-code and performing rule-base transformations

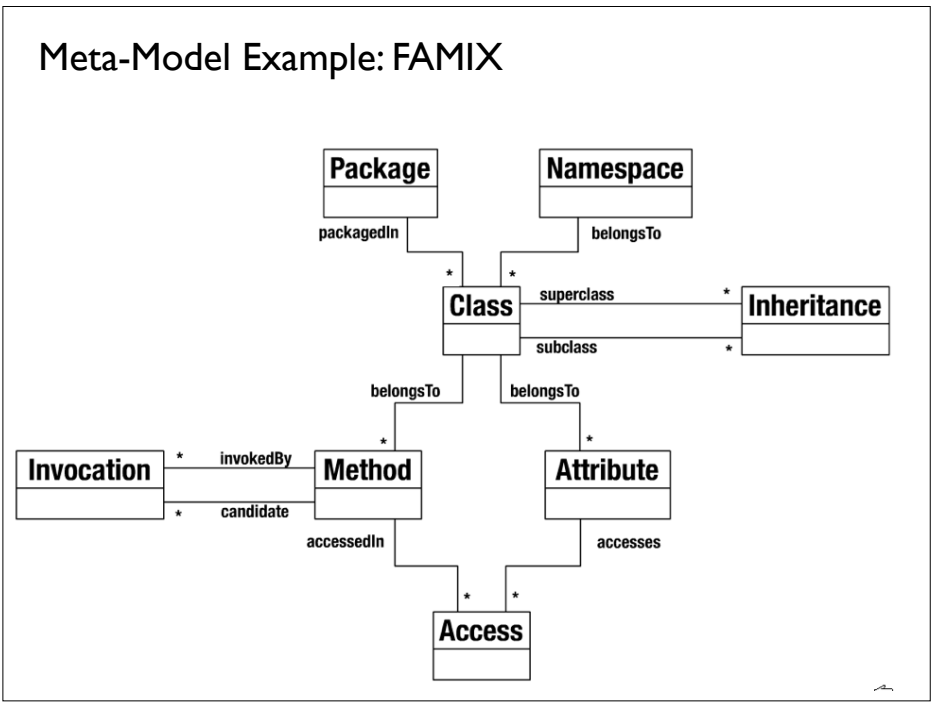
Parsing problems due

- a. diversity of language (500 problem)
- b. macros and preprocessor directives (C/C++)

Island/Lake Parsing

parse only program constructs that you care about [Moonen]

NEW!
Use parsing facilities in IDEs



Model Extractors and Analysis Toolkits

moose.unibe.ch

www.frontendart.com

loose.upt.ro/iplasma

Columbus CAN3

Example: iPlasma Toolkit

Front-end: INSIDER

Analyses: Metrics, Detection Strategies, Code Duplication, Data-flow Analyses

Models: HisMo, MEMORIA, Membrain

Model Extractors: MCC / FAST, Memoria / Recoder

Flow: Front-end → integrate → Analyses → use → Models → build → Model Extractors

1 Software Understanding

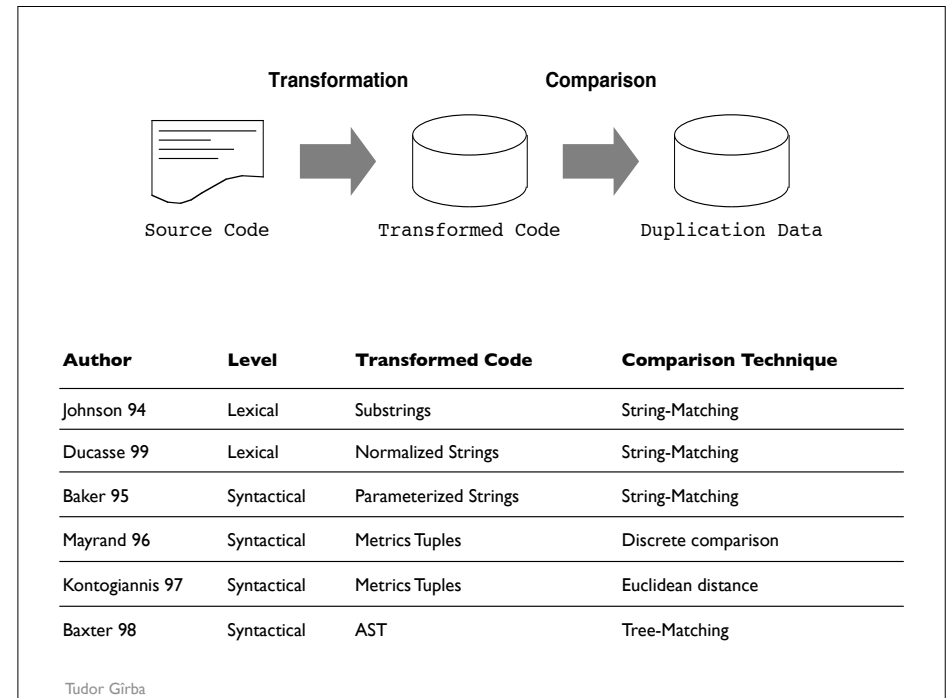
Achievement 1.3: Clone Detection and Analysis

Two Issues

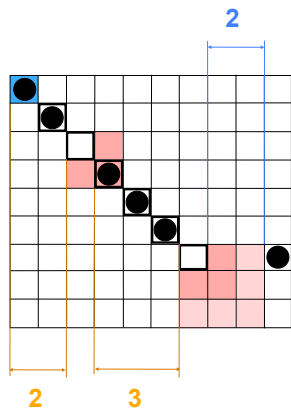
- a. precision (issue: false positives)
- b. recall (issue: false negatives)

Three Approaches

- a. token-based (high recall)
- b. AST-based (high precision)
- c. metrics-based (language-independent)



Chains of Duplication



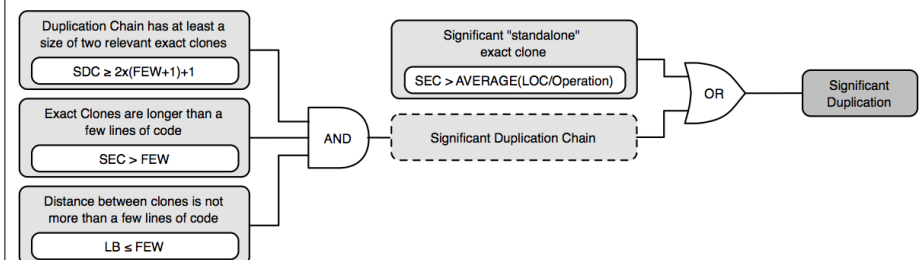
Line bias (LB)

Minimum length of a fragment
(SEC - Standalone Exact Clone)

Significant Duplication Chain (SDC)

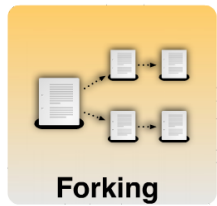
Archeology of Code Duplication

Wettel, Marinescu 2005
Lanza, Marinescu, 2006



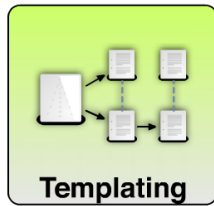
“Cloning Considered Harmful” considered harmful

Kapsler, Godfrey 2006



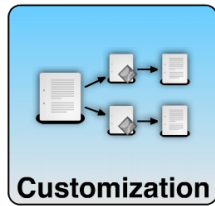
Forking

Example:
Hardware Variation



Templating

Example:
API/Library
Protocols



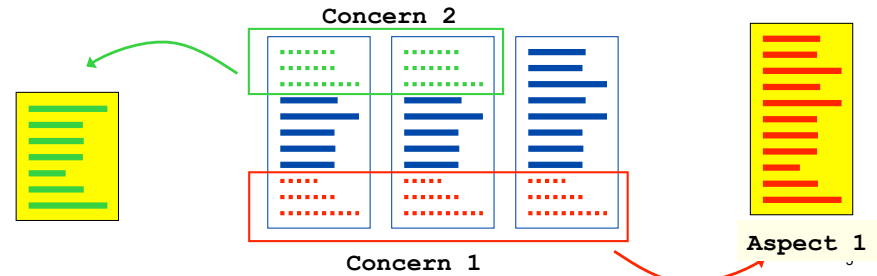
Customization

Example:
Bug Workaround

Need to understand the reason
of cloning before deciding if harmful

1 Software Understanding

Achievement I.4: Aspect Mining



Separate the principal
decomposition from crosscutting
concerns

Aspects in a Nutshell

the OO way

```

public abstract class Customer {
    private CustomerID id;
    private Collection listeners;
    public Address getAddress() {
        return this.address;
    }
    public void setLastName(String name) {
        this.lastName = name;
        notifyListeners();
    }
    public void setCustomerID(String id) {
        this.id = id;
        notifyListeners();
    }
}

public class PrivateCustomer {
    private String lastName;
    private String firstName;
    public void setLastName(String name) {
        this.lastName = name;
        notifyListeners();
    }
    public void setFirstName(String name) {
        this.firstName = name;
        notifyListeners();
    }
}

public class CorporateCustomer {
    private String companyName;
    private CompanyName taxNumber;
    public void setCompanyName(String name) {
        this.companyName = name;
        notifyListeners();
    }
    public void setTaxNumber(String nr) {
        this.taxNumber = nr;
        notifyListeners();
    }
}

public class CustomerListener {
    public void notify(Customer modifiedCustomer) {
        System.out.println("Customer " + modifiedCustomer.getID() + " was modified");
    }
}
    
```

tangling
code in one region addresses multiple concerns

scattering
code addressing one concern is spread around the system

from K.Mens, A.Kellens, J.Krinke "Pitfalls in Aspect Mining"

Aspects in a Nutshell

the OO way

```

public abstract class Customer {
    private CustomerID id;
    private Collection listeners;
    public Address getAddress() {
        return this.address;
    }
    public void setLastName(String name) {
        this.lastName = name;
        notifyListeners();
    }
    public void setCustomerID(String id) {
        this.id = id;
        notifyListeners();
    }
}

public class PrivateCustomer {
    private String lastName;
    private String firstName;
    public void setLastName(String name) {
        this.lastName = name;
        notifyListeners();
    }
    public void setFirstName(String name) {
        this.firstName = name;
        notifyListeners();
    }
}

public class CorporateCustomer {
    private String companyName;
    private CompanyName taxNumber;
    public void setCompanyName(String name) {
        this.companyName = name;
        notifyListeners();
    }
    public void setTaxNumber(String nr) {
        this.taxNumber = nr;
        notifyListeners();
    }
}

public class CustomerListener {
    public void notify(Customer modifiedCustomer) {
        System.out.println("Customer " + modifiedCustomer.getID() + " was modified");
    }
}
    
```

the AO way

```

public aspect ChangeNotification {
    pointcut stateUpdate(Customer c) :
        execution(* Customer.set*(..)) &&
        this(c);
    after(Customer c): stateUpdate(c) {
        for (Iterator iterator = c.listeners.iterator(); iterator.hasNext();
            CustomerListener listener = (CustomerListener) iterator.next();
            listener.notify(c);
        }
}
    
```

tangling
code in one region addresses multiple concerns

scattering
code addressing one concern is spread around the system

clean separation of concerns

pointcut, advice

from K.Mens, A.Kellens, J.Krinke "Pitfalls in Aspect Mining"

Aspect Mining Techniques

FANIN (methods called from many different places)

Clones

Dynamic Analysis (use traces with FCA)

Correlation of line co-changes [Canfora]

1 Software Understanding

2 Design Recovery

30

1 Software Understanding

Achievement I.5: Visualizing Software Artifacts

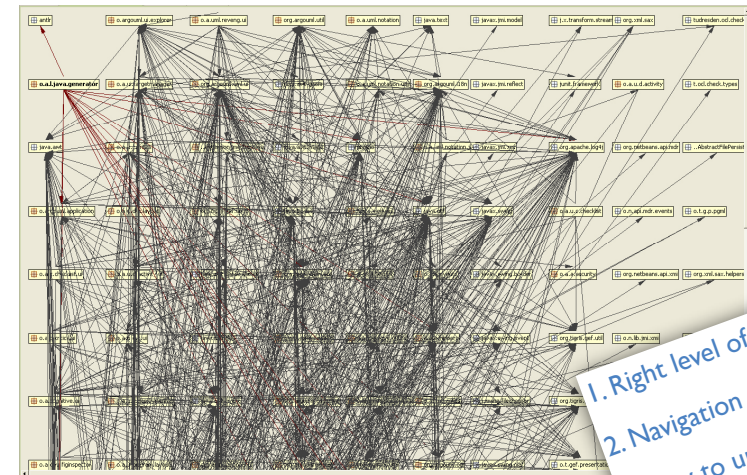
A picture is worth
a **thousand words.**

[unknown]

..depends on the picture

[Lanza]

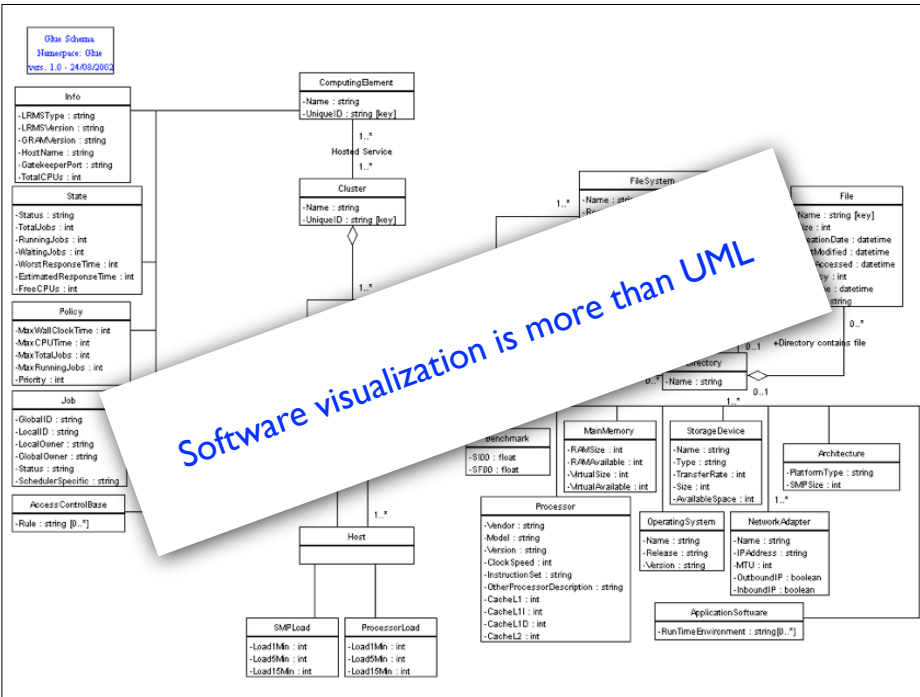
Visualization does not guarantee understanding



1. Right level of detail
2. Navigation
3. Easy to understand

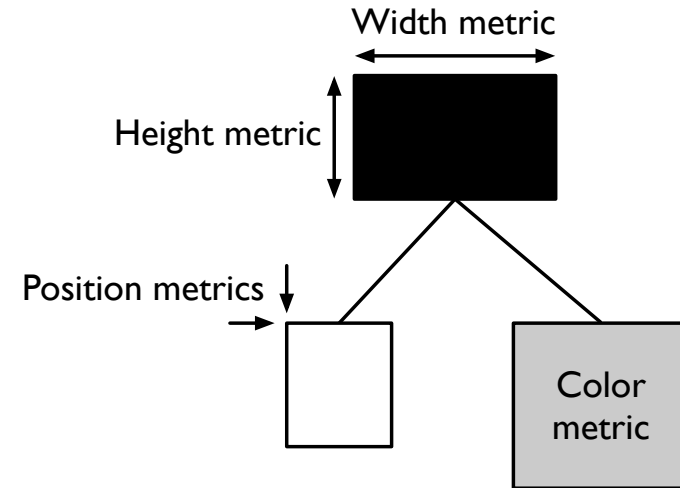
Tudor Girba

32



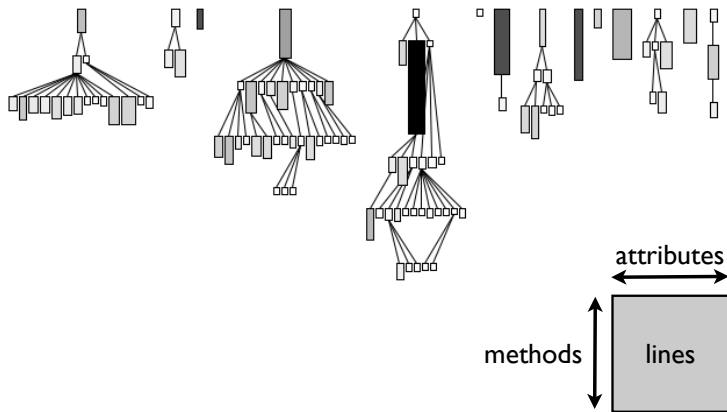
Polymetric Views show up to 5 metrics.

Lanza, 2003

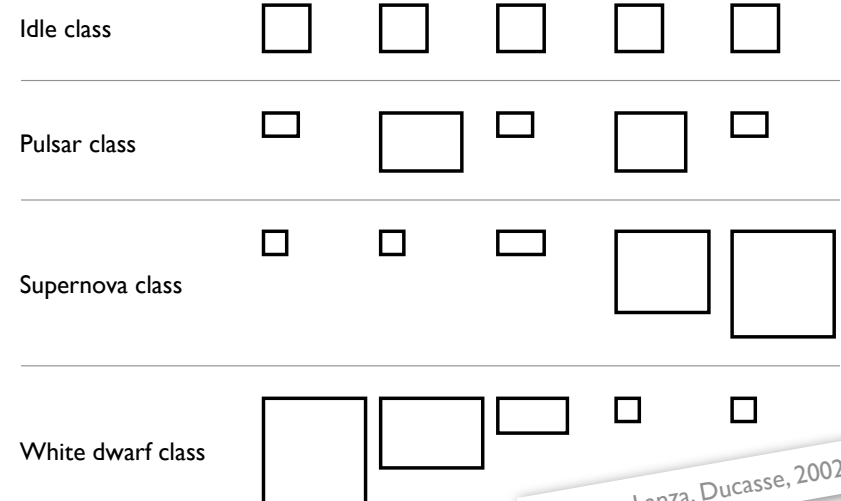


System Complexity shows class hierarchies.

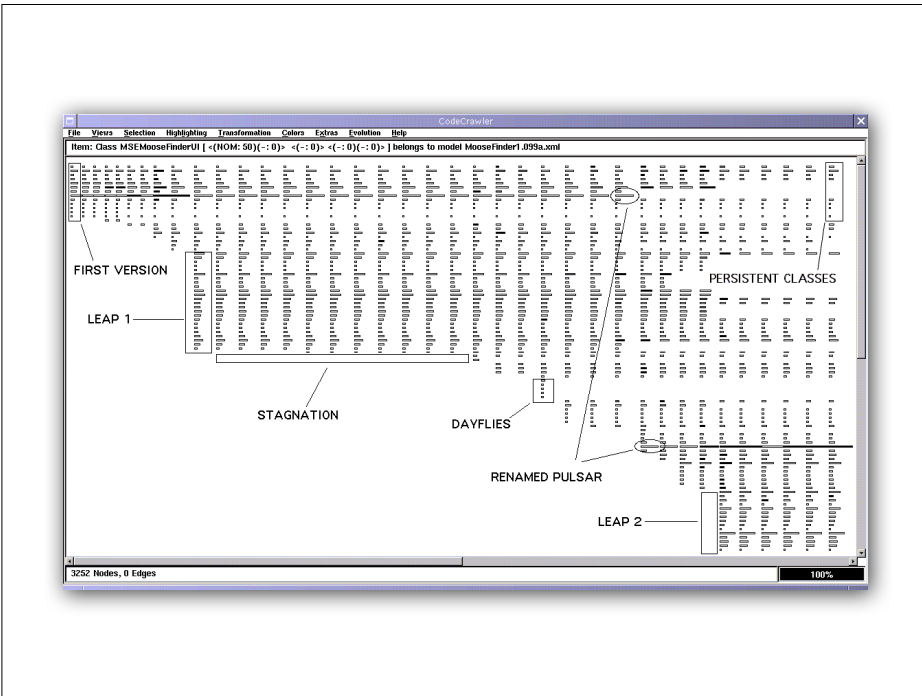
Lanza, Ducasse, 2003



Evolution Matrix shows changes in classes

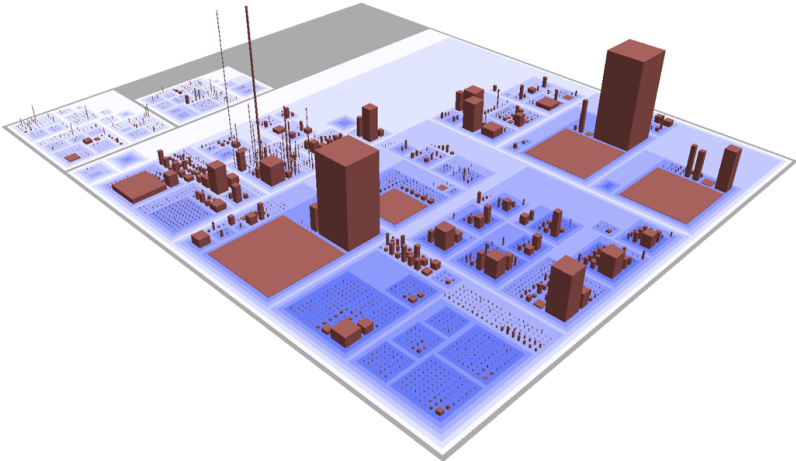


Lanza, Ducasse, 2002



Code City shows where your code lives.

Wettel, Lanza, 2007



classes are buildings grouped in quarters of packages

1 Software Understanding

Trends

- Understanding Systems with High Dynamicity
- Understanding Cross-Language Systems
- Mining Software Repositories

39

1 Software Understanding

Trend 1.1: Understanding Systems with High Dynamicity

Reflection and loading classes at run-time

BAD

- affects points-to analysis
- Solution: DA must complement SA

GOOD

- access to members of classes ;
- support for analysis JVMTI instead of instrumenting code

40

1 Software Understanding

Trend I.2: Understanding Cross-Language Systems

Especially in the case of Web Applications (and .NET)

How to analyze in an **integrated manner** different types of code that coexist

HTML,

Object-oriented,

SQL,

Scripting code

41

1 Software Understanding

Trend I.3: Mining Software Repositories

Analyze versioning systems (study the evolution of software)

Changes correlated with other faultiness [Zeller]

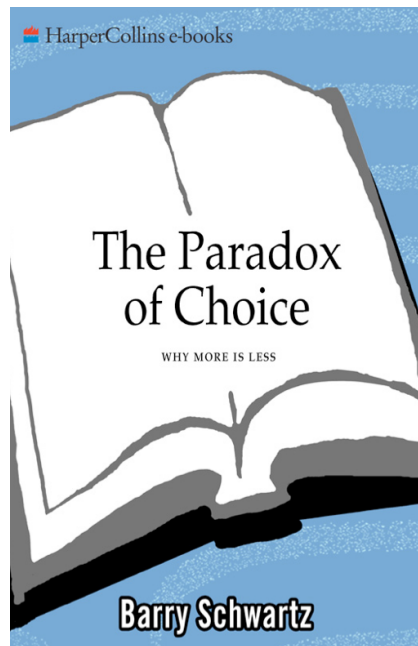
Co-changes correlation with code duplication [Geiger etal]

Co-changes to refine detection of design flaws [Ratiu etal]

Analyze developers behavior

Moved from 1D (SA)
to 2D (SA + DA)
to 3D (SA + DA + Time)

42



43

2 Design Recovery

Achievements

Recovery of UML Models (class and object diagrams)

Identifying Design Patterns (motifs) in code (Guéhéneuc)

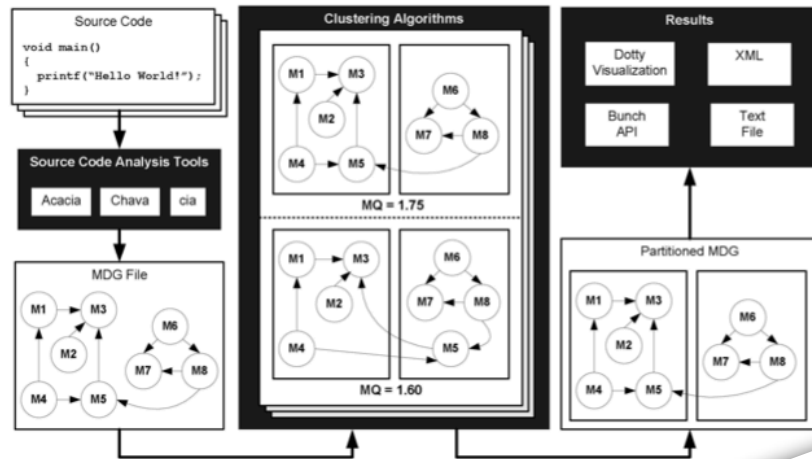
Clustering-Based Architecture Recovery

Feature/Concept Location

44

2 Design Recovery

Achievement 2.2: Clustering-Based Architecture Recovery



Mitchell, Mancoridis, TSE 2006

Formal Concept Analysis (FCA) in a nutshell.

Incidence Table

		properties					
		P1	P2	P3	P4	P5	P6
elements	A	X				X	X
	B	X			X		
	C			X		X	X
	D	X			X		X
	E					X	X

Legend: X a binary relation between the element on the row and the property on the column

from T.Girba "Modeling History to Understand Software Evolution". PhD 2003

Formal Concept Analysis (FCA) in a nutshell.

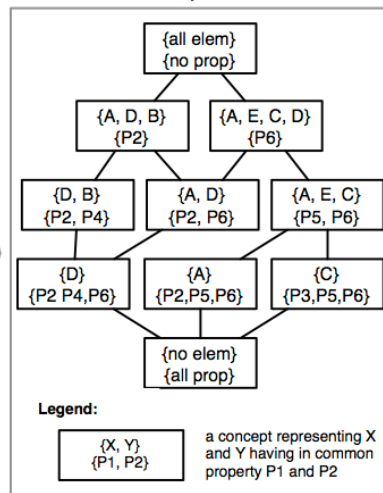
Incidence Table

		properties					
		P1	P2	P3	P4	P5	P6
elements	A	X				X	X
	B	X			X		
	C			X		X	X
	D	X			X		X
	E					X	X

Legend: X a binary relation between the element on the row and the property on the column

FCA

Concept Lattice



from T.Girba "Modeling History to Understand Software Evolution". PhD 2003

2 Design Recovery

Achievement 2.2: Clustering-Based Architecture Recovery

Element = Function

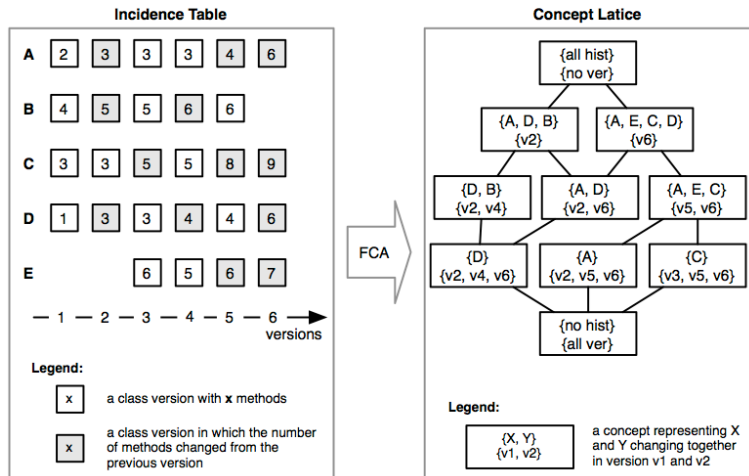
Property = Type used by Function

	returns stack	returns queue	has stack arg.	has queue arg.	uses stack fields	uses queue fields	not queue fields
initStack	✓				✓		✓
initQ		✓				✓	
isEmptyS			✓		✓		✓
isEmptyQ				✓	✓	✓	
push			✓		✓		✓
enq				✓	✓	✓	
pop			✓		✓		✓
deq				✓		✓	

Siff, Reps, TSE 1999

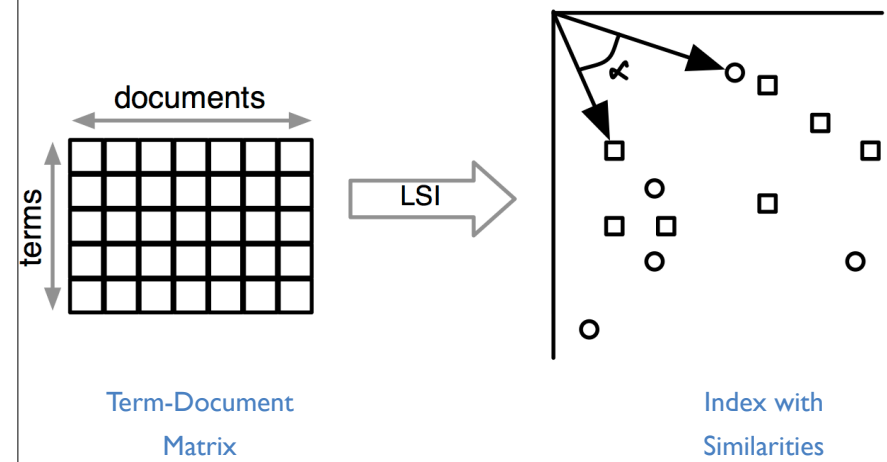
2 Design Recovery

Achievement 2.2: Clustering-Based Architecture Recovery



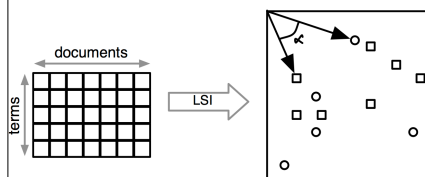
Girba, PhD, 2003

Latent Semantic Indexing (LSI) in a nutshell.



from A. Kuhn et al "Enriching Reverse Engineering with Semantic Clustering", 2005

Latent Semantic Indexing (LSI) in a nutshell.



Benefits

Search Engines

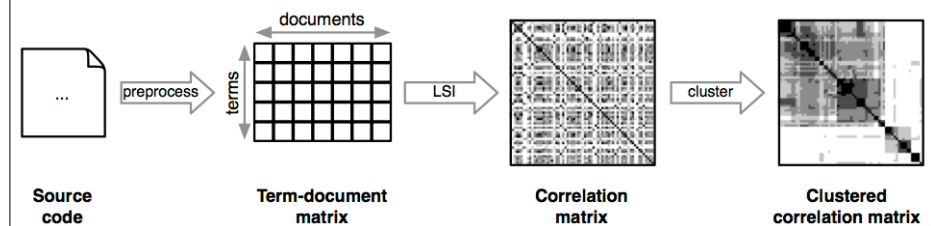
Detect authorship in literature

Automate assignments of papers to reviewers

from A. Kuhn et al "Enriching Reverse Engineering with Semantic Clustering", 2005

2 Design Recovery

Achievement 2.2: Clustering-Based Architecture Recovery



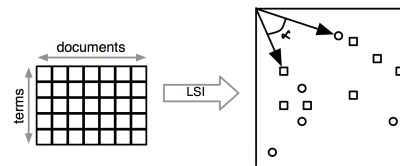
document = classes, methods

terms = identifier names, comments

Kuhn, et al,
WCRE 2005

2 Design Recovery

Achievement 2.3: Feature/Concept Identification and Location



identifying parts of the source code implementing domain concepts

based on combination of

dynamic analysis (Scenario-Based Probabilistic Ranking, good recall) and
static analysis (LSI, good precision)

Poshyvanyk,
ICPC 2006

53

2 Design Recovery

Trends

Extraction of Object Diagrams and Constrains (pre/post conditions)

Migration to Web 2.0 applications

Interactive RE environments

54

2 Design Recovery

Trend 2.3: Interactive RE environments

RE: not only source code info; but also **developers rationales**

2 Problems: (i) incompleteness ; (ii) semi-automatic

Solution: **interactively improve artifact presentation**

- give feedback to RE system (**efficiently!**)
- RE systems should learn (machine learning, GA)
- best in IDEs

55

Issues and Challenges

A

 Continuous RE

B

 RE for SOA and
Autonomic Computing

A Continuous RE

Exploit RE in the Fwd. Eng. process (especially in IDEs)

Benefits

1. **Clearer picture** (understanding) of the developed system
2. **Permanent consistency checks** (design vs. code vs. tests)
3. **Better QA hints**
 - continuous monitoring
 - learning tools for better traceability links

57

B RE for SOA and Autonomic Computing

SOA - separation of software ownership (product) from software use (service)

Autonomic Computing - self-adaptation and self-evolution

58

B RE for SOA and Autonomic Computing

SOA

Each service offers a limited view

Danger: system as an **orchestration** of various services

discrepancy between terminology = **harder to understand**
no access to source-code

59

B RE for SOA and Autonomic Computing

Systems with Autonomic Capabilities

Autonomic Discovery

find and bind new services when default is not ok

Self-Healing

change composition, reconfiguring/repairing system, interrupting exec

Danger: **high dynamism**

pieces composing the execution known only at run-time

60

What about RE Tools as
Composition of Services? :)