

Analyzing and Inferring the Structure of Code Changes

Miryung Kim

Computer Science & Engineering
University of Washington
Final Exam
September 26th, 2008

***Software evolution plays an ever-increasing
role in software development***

Motivating Scenarios

- “This program worked a month ago but is not working now. What changed since then? Which change led to a bug?”
- “Did Bob implement the intended changes correctly?”
- “There’s a merge conflict. What did Alice change?”

Diff Output

Changed Code		
File Name	Status	Lines
DummyRegistry	New	20 lines
AbsRegistry	New	133 lines
JRMPRegistry	Modified	123 lines
JeremieRegistry	Modified	52 lines
JacORBCosNaming	Modified	133 lines
IIOPCosNaming	Modified	50 lines
CmiRegistry	Modified	39 lines
NameService	Modified	197 lines
NameServiceManager	Modified	15 lines
Total Change: 9 files, 723 lines		

```
- public class CmiRegistry implements
NameService {
+ public class CmiRegistry extends
AbsRegistry implements NameService {
-     private int port = ...
-     private String host = null
-     public void setPort (int p) {
-         if (TraceCarol. isDebug()) { ...
-         }
-     }
-     public int getPort() {
-         return port;
-     }
-     public void setHost(String host)
{ ....
```

Check-In Comment

“Common methods go in an abstract class. Easier to extend/maintain/fix”

Changed Code		
File Name	Status	Lines
DummyRegistry	New	20 lines
AbsRegistry	New	133 lines
JRMPRegistry	Modified	123 lines
JeremieRegistry	Modified	52 lines
JacORBCosNaming	Modified	133 lines
IIOPCosNaming	Modified	50 lines
CmiRegistry	Modified	39 lines
NameService	Modified	197 lines
NameServiceManager	Modified	15 lines
Total Change: 9 files, 723 lines		

Why did all these files change together?
Is anything missing in this change?

Limitations

Diff

- Low-level

Natural Language Description (Check-In Comment)

- Often incomplete
- Difficult to trace back to code changes

Research Question

How do we **automatically extract** the differences between two versions into a **concise** and **meaningful** program change representation?

Research Question

How do we **automatically extract** the differences between two versions into a **concise** and **meaningful** program change representation?

- Help programmers reason about code changes **at a high level**
- Enable researchers to study software evolution better

Example Output

All `draw` methods take an additional `int` input argument.

All `setHost` methods in `Service`'s subclasses deleted calls to SQL library **except `NameService` class.**

...

Concise
Easy to note inconsistent changes

Systematic Changes

- Refactoring [Opdyke 92, Griswold 92, Fowler 99...]

“Move related classes from one package to another package”

Systematic Changes

- Refactoring [Opdyke 92, Griswold 92, Fowler 99...]
- API update [Chow&Notkin 96, Henkel&Diwan 05, Dig&Johnson 05...]

“Update an API and all call sites of the API”

Systematic Changes

- Refactoring [Opdyke 92, Griswold 92, Fowler 99...]
- API update [Chow&Notkin 96, Henkel&Diwan 05, Dig&Johnson 05...]
- Crosscutting concerns [Kiczales et. al. 97, Tarr et. al. 99, Griswold 01...]

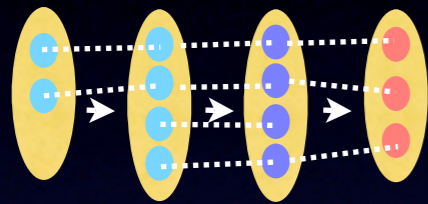
“Adding logging feature throughout code”

Systematic Changes

- Refactoring [Opdyke 92, Griswold 92, Fowler 99...]
- API update [Chow&Notkin 96, Henkel&Diwan 05, Dig&Johnson 05...]
- Crosscutting concerns [Kiczales et. al. 97, Tarr et. al. 99, Griswold 01...]
- Consistent updates on code clones [Miller&Myers 02, Toomim et. al. 04, Kim et. al. 05]

“Apply similar changes to syntactically similar code fragments”

Thesis Overview

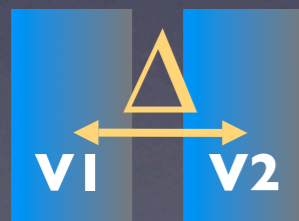


Analyses of Software Evolution

- Evolution of Code Clones



High-level changes are often systematic at a code level



Automatic Inference of High-Level Change Descriptions

- Rule-based Change Representations
- Rule Learning Algorithms

Outline

- Empirical Analyses of Code Clone Evolution [ISESE 04, ESEC/FSE 05]
- Automatic Inference of High-Level Change Descriptions
 - Changes to API Names and Signatures [ICSE 07]
 - Changes to Code Elements and Structural Dependencies
- Future Directions

Code Clones

Code clones are syntactically similar code fragments

```
public void updateFrom (Class c) {  
    String cType = Util.makeType(c.Name  
    ());  
    if (seenClasses.contain(cType)) {  
        return;  
    }  
    seenClasses.add(cType);  
    if (hierarchy!=null) {  
        ....  
    }  
    ...  
}
```

```
public void updateFrom (ClassReader c) {  
    String cType = CTD.convertType  
    (c.Name());  
    if (seenClasses.contain(cType)) {  
        return;  
    }  
    seenClasses.add(cType);  
    if (hierarchy!=null) {  
        ....  
    }  
    ...  
}
```

Found by a clone detector, CCFinder [Kamiya et al. 2002]

Conventional Wisdom about Code Clones

“Code clones must be aggressively refactored because they indicate poor software quality.”

[Fowler 00, Beck 00, Nickell & Smith 03 ...]

```
public void updateFrom (Class c) {  
    String cType = Util.makeType(c.Name  
    ());  
    if (seenClasses.contain(cType)) {  
        return;  
    }  
    seenClasses.add(cType);  
    if (hierarchy!=null) {  
        ....  
    }  
    ...  
}
```

```
public void updateFrom (ClassReader c) {  
    String cType = CTD.convertType  
    (c.Name());  
    if (seenClasses.contain(cType)) {  
        return;  
    }  
    seenClasses.add(cType);  
    if (hierarchy!=null) {  
        ....  
    }  
    ...  
}
```

Found by a clone detector, CCFinder [Kamiya et al. 2002]

A Study of Copy and Paste Programming Practices at IBM

[Kim et al. ISESE 2004]

- To understand programmers' copy and paste coding behavior, I **built an Eclipse plug-in that records edits and replays the captured edits**
- Programmers often **create and manage** code clones with clear intent

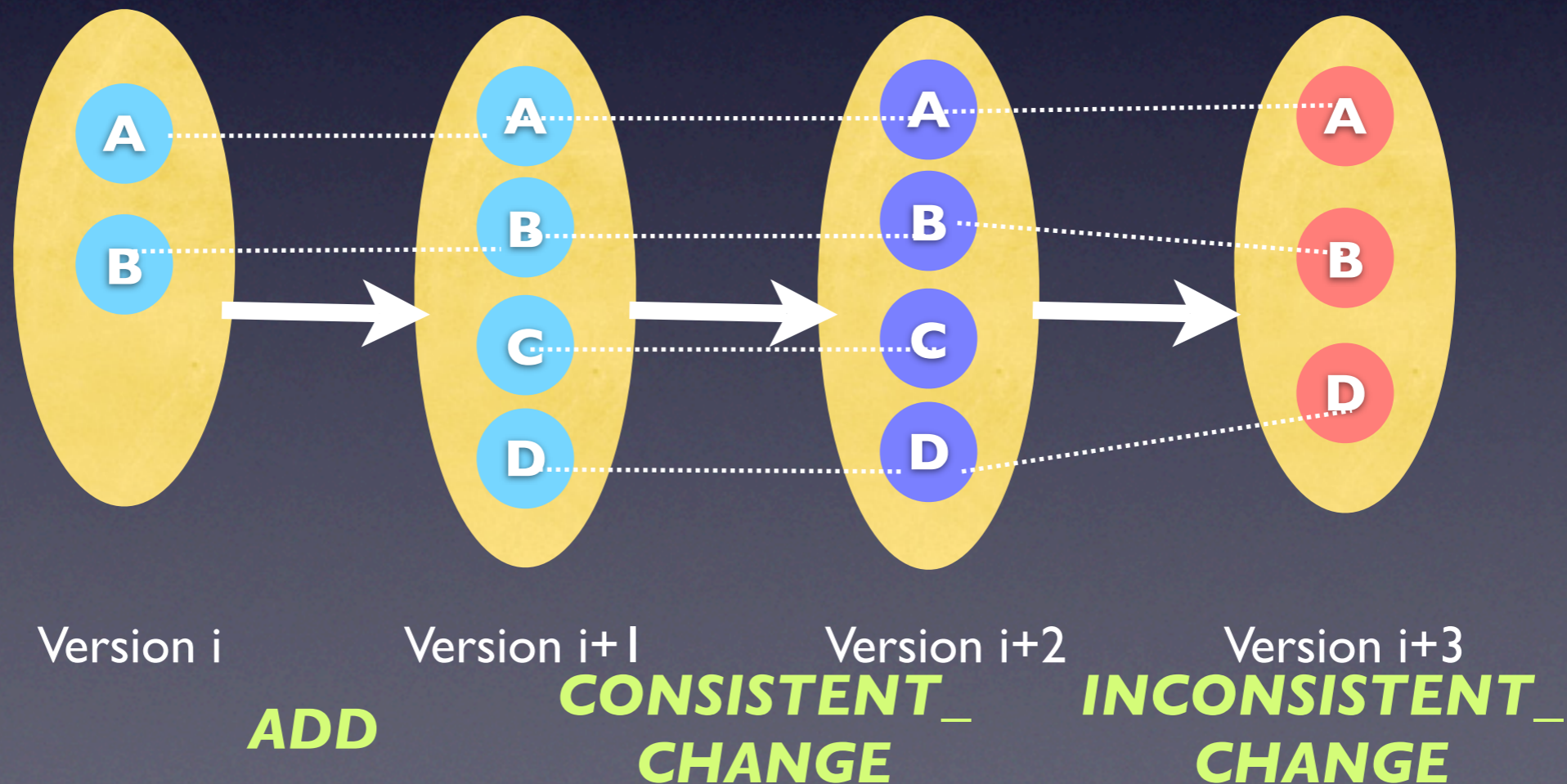
An Empirical Study of Code Clone Genealogies

[Kim et al. ESEC/FSE 2005]

- I developed an approach that ***automatically reconstructs*** the history of code clones ***from a source code repository***
- I studied clone evolution in two Java open source projects, *carol* and *dnsjava*

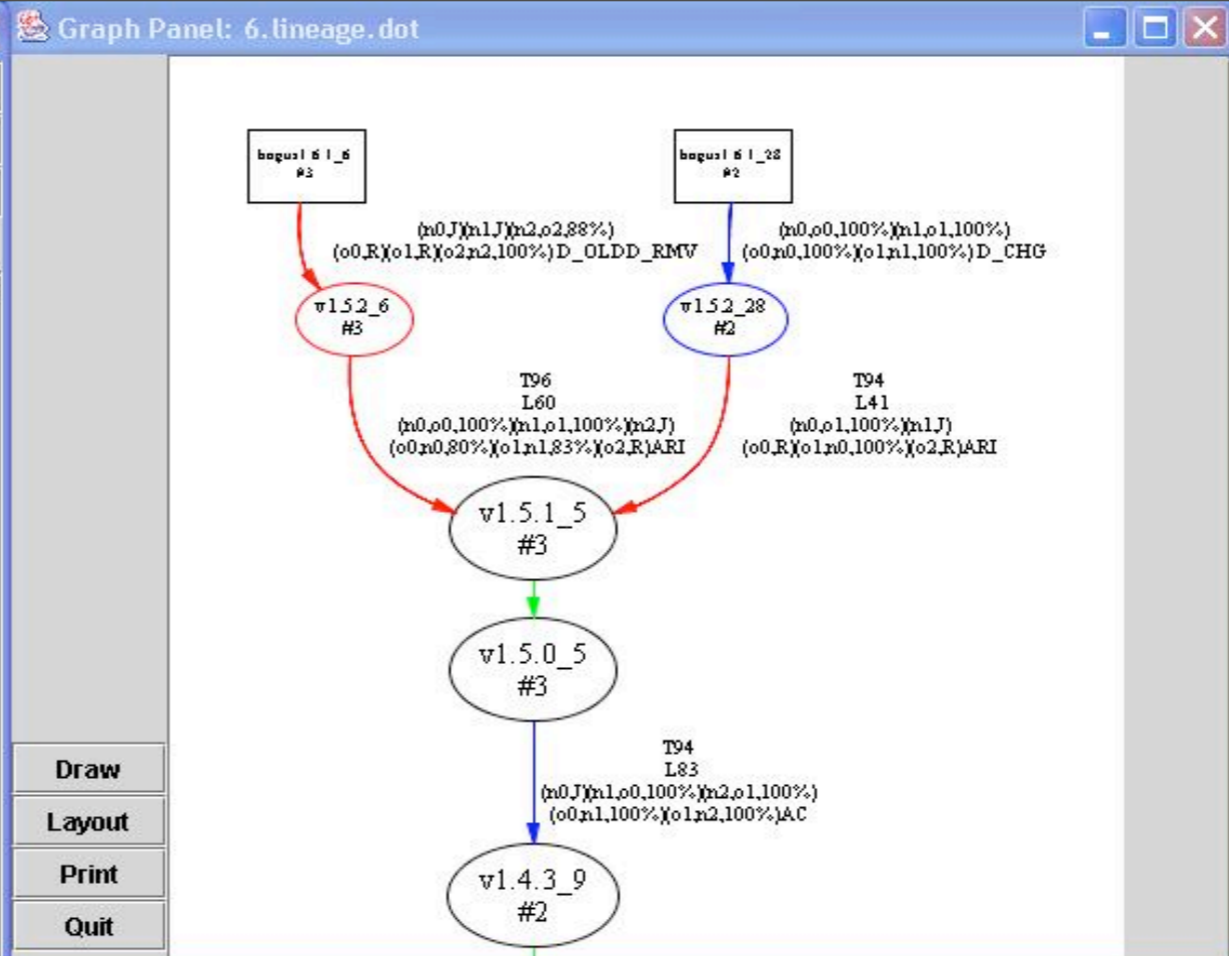
Clone Genealogy

Clone genealogy is a representation that captures clone change patterns over a sequence of program versions



Code	Graph
Postscript	Report
ReleaseStat	LineageStat

8:1.2.0~1.2.4 L:4 S4 Good Fact...
9:1.2.0~1.2.4 L:4 S4 Good Notft...
1.3.0
33:0.9~1.3.0 L:21 C1 S20 Good Fact...
52:1.2.4~1.3.0 L:1 S1 Good Fact...
1.3.1
2:1.1~1.3.1 L:13 C2 S11 Good Notft...
45:0.1~1.3.1 L:33 A1 S32 Good Notft...
1.3.2
13:0.9.2~1.3.2 L:21 C1 S20 Good Notft...
42:0.1~1.3.2 L:34 A3 R2 C3 I2 S77 Good Fact...
5:1.3.0~1.3.2 L:2 S2 Good Notft...
1.3.3
20:1.3.3~1.3.3 L:0 Bad Notft...
3:0.9.1~1.3.3 L:23 S23 Good Notft...
33:1.3.3~1.3.3 L:0 Good Notft...
48:0.1~1.3.3 L:35 C4 S31 Good Fact...
6:1.3.3~1.3.3 L:0 Good Notft...
1.4.0
1.4.1
1.4.2
1.4.3
10:1.4.0~1.4.3 L:3 S3 Good Notft...
1.5.0
1.5.1
13:1.4.0~1.5.1 L:5 S5 Good Notft...
1.5.2
40:1.4.0~1.5.2 L:6 C1 S5 Good Notft...
48:1.3.0~1.5.2 L:10 C1 S9 Good Notft...
57:1.5.0~1.5.2 L:2 S2 Good Notft...
6:1.4.0~1.5.2 L:6 A3 R2 C1 I2 S4 Good Notft...
1.6.1
1.6.2



Group View

Close Compare Write Note Toggle Refactor Toggle Good Trace Forward Trace Backward

1.5.2-6

- 1.4.0:9 Good Notft...
- 1.4.1:9 Good Notft...
- 1.4.2:9 Good Notft...
- 1.4.3:9 Good Notft...
- 1.5.0:5 Good Notft...
- 1.5.1:5 Good Notft...
- 1.5.2:28 Good Notft...
- 1.5.2:6 Good Notft...

```

1.5.2-CERTRecord
1.5.2-DSRecord
1.5.2-KEYBase
1.5.0-CERTRecord
1.5.0-DSRecord
1.5.0-KEYBase

checkObj( arg , arg),
this.key = key;
}
Record
rrFromWire(Name name, int type, int dclass, long ttl, int length,
DataByteInputStream in)
throws IOException
{
KEYRecord rec = new KEYRecord(name, dclass, ttl);
if (in == null)
return rec;
rec.flags = in.readShort();
rec.proto = in.readByte();
rec.alg = in.readByte();
if (length > 4) {
rec.key = new byte[length - 4];
in.read(rec.key);
}
return rec;
}
Record
rdataFromString(Name name, int dclass, long ttl, Tokenizer st, Name origin)
throws IOException

```

6:1.4.0~1.5.2 L:6 A3 R2 C1 I2 S4 Good Notft Control Logic

Contradicting Evidence to Conventional Wisdom

- *Many clones are short-lived, diverging clones*
 - 48-72% of clone genealogies lasted less than 8 check-ins out of over 160 check-ins
 - 26-34% of these clones disappeared due to divergent changes
- *Refactoring cannot remove many long-lived clones*
 - 65-73% of long-lived, consistently changing clones are not easy to refactor using standard refactoring techniques [Folwer 00]

Summary of Studies on Code Clones

By focusing on the **evolutionary aspects of clones**, I found

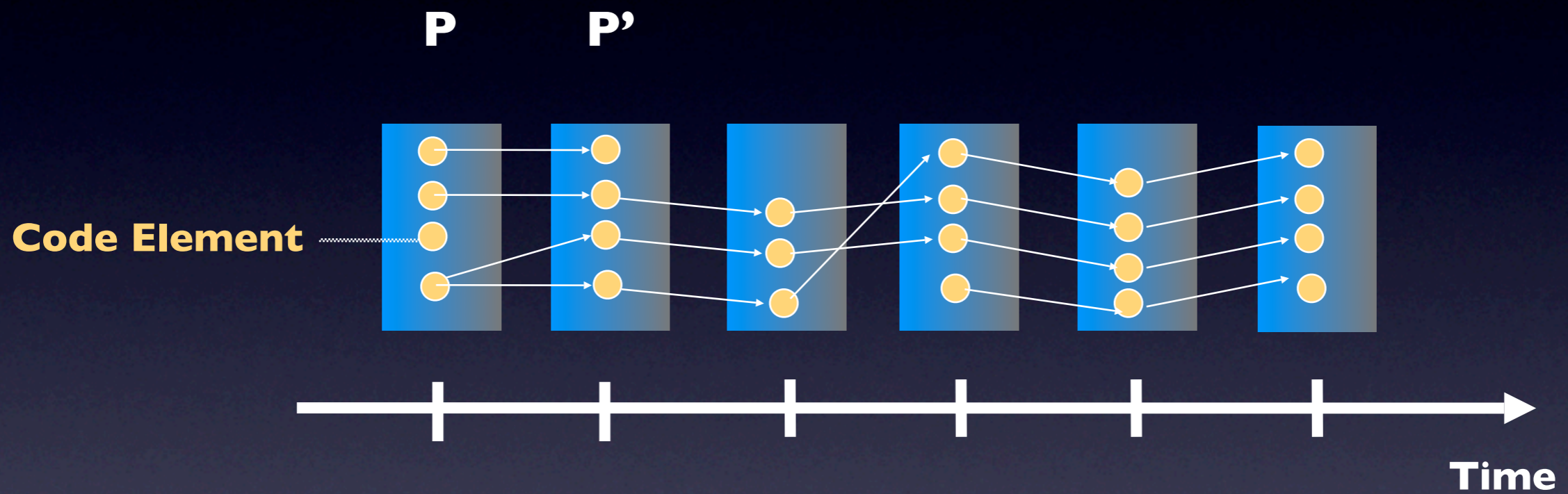
- Clones are inevitable parts of software evolution
- Refactoring may not be applicable to or beneficial for many code clones

My studies shifted research efforts from automatic clone detection to code clone management support (e.g., [Duala-Ekoko & Robillard 07, Krinke 07, Aversano et al. 07, Lozano et al. 07, etc.])

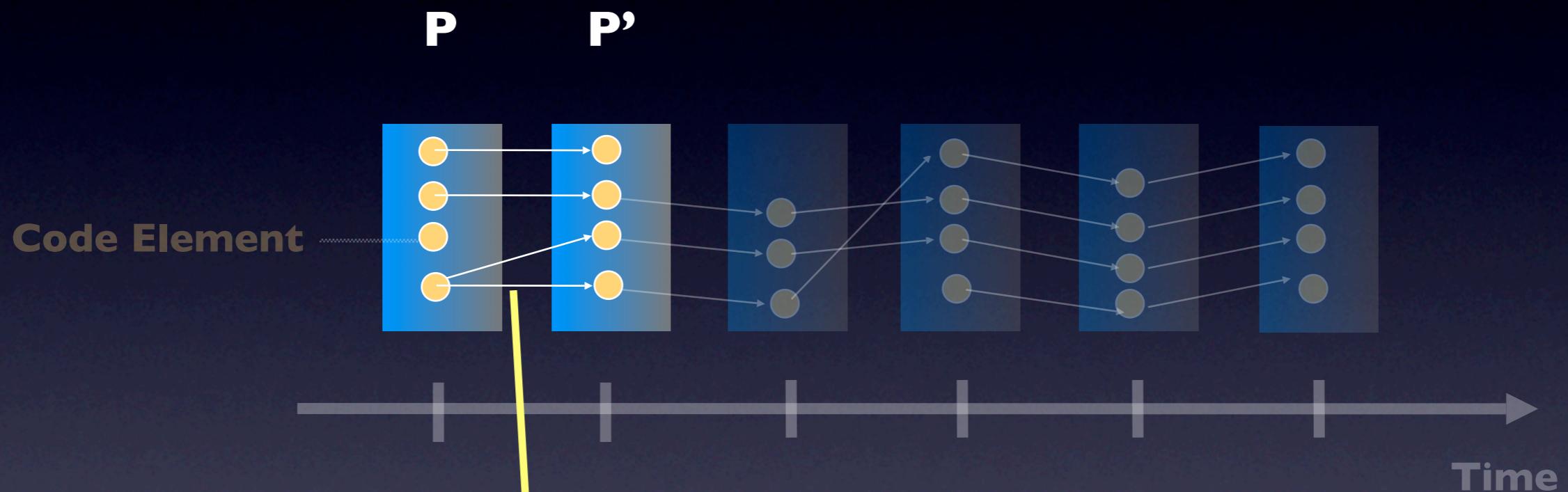
Outline

- Empirical Analyses of Code Clone Evolution
- Automatic Inference of High-Level Change Descriptions
 - Changes to API Names and Signatures
 - Changes to Code Elements and Structural Dependencies
- Future Directions

Motivation: Code Evolution Analyses



Research Question



“How do we automatically match corresponding code elements between two program versions?”

Existing Approaches

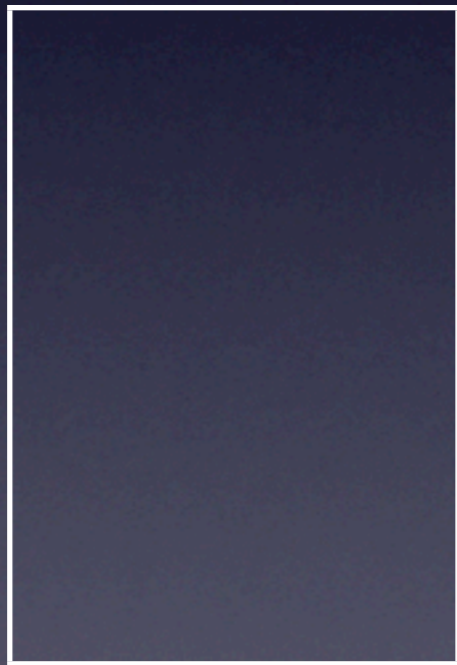
[Kim et al. MSR 2006]

diff, *Syntactic Diff (CDiff)*, *Semantic Diff*, *JDiff*,
BMAT, *origin analysis*, *refactoring*
reconstruction tools, *clone detectors*, etc.

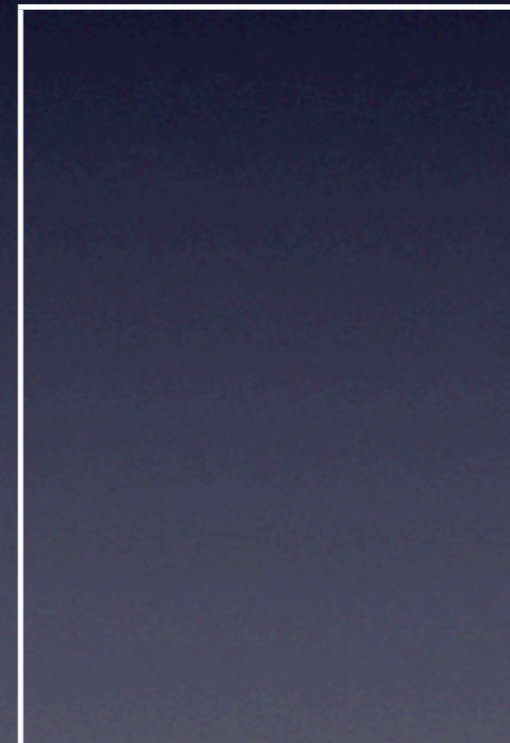
Individually compare code elements
at particular granularities
using similarity measures

Limitations of Existing Approaches

P



P'



Limitations of Existing Approaches

P

<code>Bar.Bar()</code>
<code>Bar.mC(int)</code>
<code>Foo.mA()</code>
<code>Foo.mB()</code>
<code>Foo.mC()</code>
<code>Boo.mA(long)</code>
<code>Boo.mB(long)</code>

P'

<code>Bar.Bar()</code>
<code>Bar.mC(int)</code>
<code>Foo.mA(float)</code>
<code>Foo.mB(float)</code>
<code>Foo.mC()</code>
<code>Bar.mA(long)</code>
<code>Boo.mA(int)</code>
<code>Boo.mB(int)</code>

Limitations of Existing Approaches

P

Bar.Bar ()
Bar.mC (int)
Foo.mA ()
Foo.mB ()
Foo.mC ()
Boo.mA (long)
Boo.mB (long)

P'

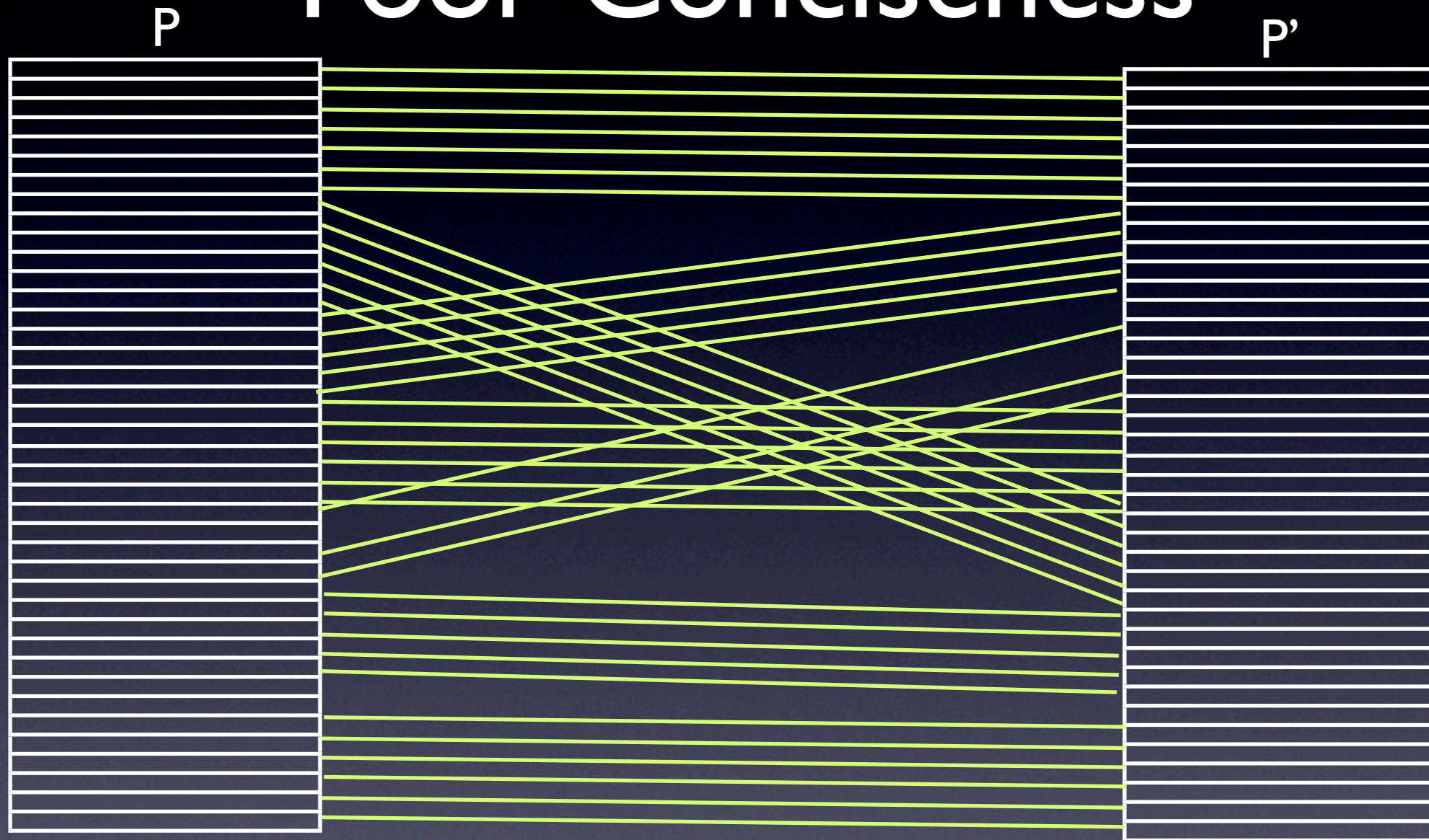
Bar.Bar ()
Bar.mC (int)
Foo.mA (float)
Foo.mB (float)
Foo.mC ()
Bar.mA (long)
Boo.mA (int)
Boo.mB (int)

Limitations of Existing Approaches



Limitation I.

Poor Conciseness



Output is an unstructured, usually lengthy list of matches

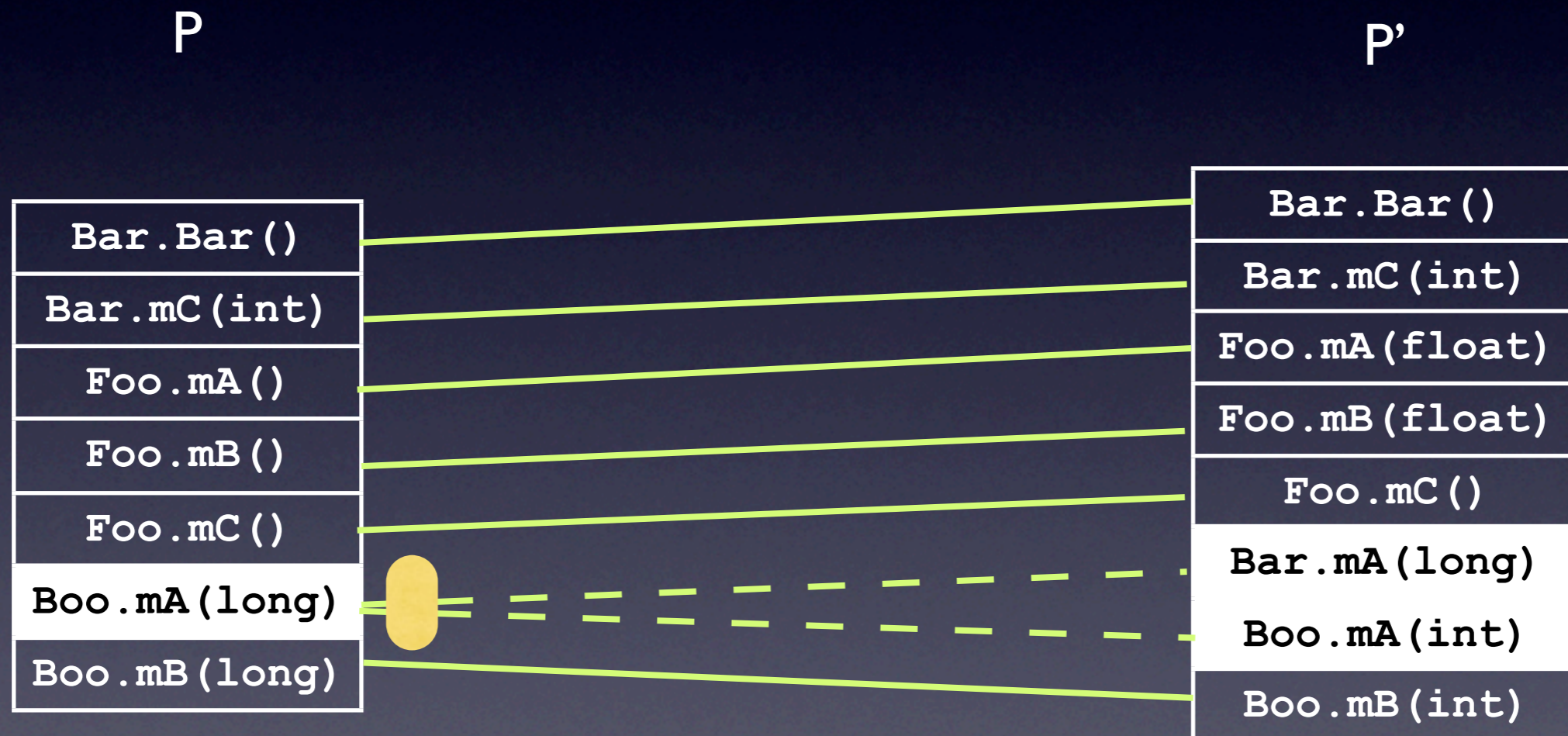
Limitation 2.

Hard to Identify Exception



Difficult to spot inconsistent changes

Limitation 3. Low Recall



Difficult to disambiguate among many potential matches

What is the Core Question?

*Given two program versions (P, P'),
with respect to a particular vocabulary of changes,
find changes from P to P'*

Example Change

P

<code>Factory.createChart()</code>
<code>Factory.createBarChart()</code>
<code>...</code>
<code>Factory.createPieChart()</code>
<code>Factory.createLineChart()</code>

P'

<code>Factory.createChart(int)</code>
<code>Factory.createBarChart(int)</code>
<code>...</code>
<code>Factory.createPieChart()</code>
<code>Factory.createLineChart(int)</code>

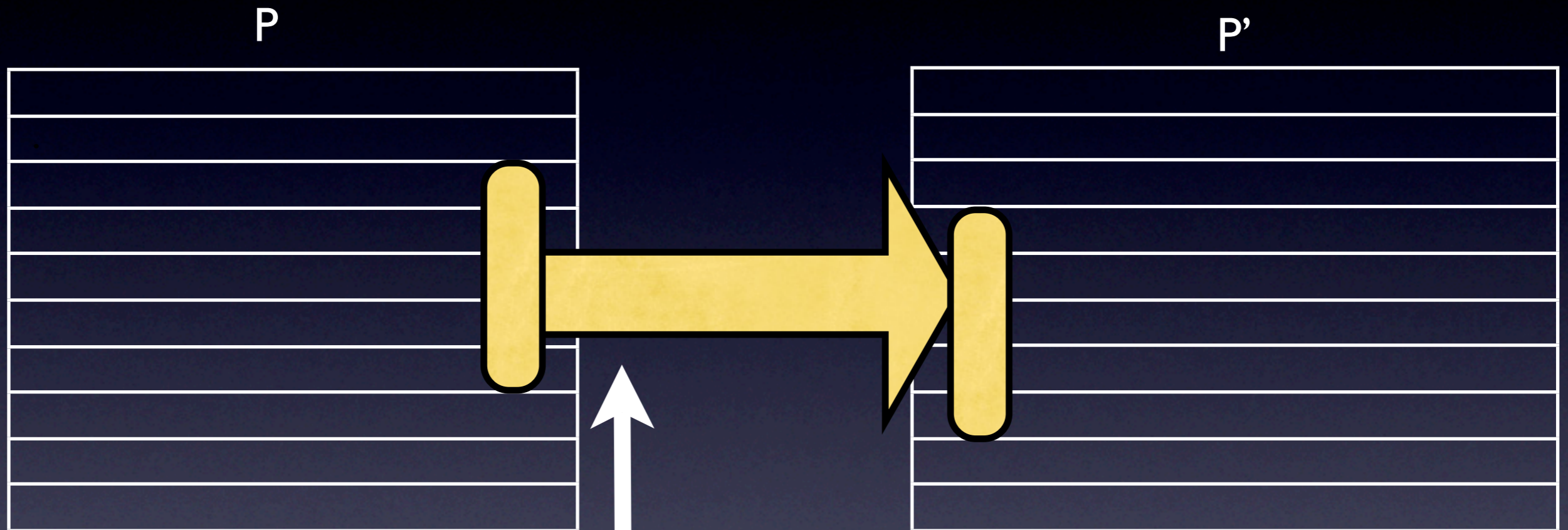
“Add `int` input argument to all chart creation APIs”

Our Rule-based Matching Approach

[Kim et al. ICSE 2007]

- Our **change-rules** can concisely describe **a set of related API-level changes**.
- Our tool **automatically infers** a set of change rules between two versions of a program.

Change-Rule Syntax



```
FOR ALL x:method-header IN  
    scope  
    transformation(x)
```

Scope

- We use a regular expression to denote a set of methods

e.g. `chart.Factory.create*Chart(*)`

API-Level Transformations

- Replace the name of package, class, and method
- Replace the return type
- Modify the input signature, etc.

Example Change-Rule

P

<code>Factory.createChart()</code>
<code>Factory.createBarChart()</code>
<code>...</code>
<code>Factory.createPieChart()</code>
<code>Factory.createLineChart()</code>

P'

<code>Factory.createChart(int)</code>
<code>Factory.createBarChart(int)</code>
<code>...</code>
<code>Factory.createPieChart()</code>
<code>Factory.createLineChart(int)</code>

```
FOR ALL x:method-header IN
  Factory.create*Chart(*)
    argAppend(x, [int])
```

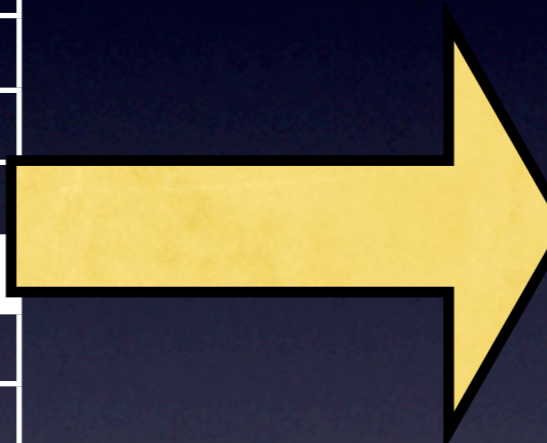
Example Change-Rule

P

Factory.createChart()
Factory.createBarChart()
...
Factory.createPieChart()
Factory.createLineChart()

P'

Factory.createChart(int)
Factory.createBarChart(int)
...
Factory.createPieChart()
Factory.createLineChart(int)



```
FOR ALL x:method-header IN
  Factory.create*Chart(*)
    argAppend(x, [int])
except {Factory.createPieChart()}
```

Algorithm Overview

Input: two versions of a program

Output: a set of change-rules

1. Generate seed matches
2. Generate candidate rules by generalizing seed matches
3. Evaluate and select candidate rules

Step 1: Generate Seed Matches



- Seed matches provide **hints** about likely changes.
- We generate seeds based on textual similarity between two method headers.
- Seed matches need not be all correct matches.

Step 2: Generate Candidate Rules

For each seed $[x, y]$

- Compare x and y and reverse engineer a set of transformations, T .
- Based on x , guess a set of scopes, S .
- Generate candidate rules for each pair in $S \times \text{PowerSet}(T)$.

```
Given a seed match,  
[Foo.getBar(int), Boo.getBar(long)]
```

```
Transformations = {  
  replaceArg(x, int, long)  
  replaceClass(x, Foo, Boo)}
```

```
Scopes = {*.*(*), Foo.*(*) , ... ,  
  *.get*(*) , *.*Bar(*) , ... ,  
  Foo.get*(int) , ... }
```

```
Candidate Rules = {  
  FOR ALL x IN *.*(*)  
    replaceArg(x, int, long),  
  FOR ALL x IN Foo.*(*)  
    replaceClass(x, Foo, Boo), ... ,  
  FOR ALL x IN *.*(*)  
    replaceArg(x, int, long) AND  
    replaceClass(x, Foo, Boo)  
  ... }
```

Step 3: Evaluate and Select Rules

- Greedily select a small subset of candidate rules that explain a large number of matches.
- In each iteration
 - evaluate all candidate rules
 - select a **valid** rule with the most number of matches
 - exclude the matched methods from the set of remaining unmatched methods
- Repeat until no rule can find any additional matches.

Optimizations

- We **create** and **evaluate** rules **on demand**
 1. Candidate rules have subsumption structure
e.g., $*.*.*(*Axis) \subset *.*.*(*)$
 2. The nature of greedy algorithm
- Running time: a few seconds (usual check-ins), average 7 minutes (releases)

Comparative Evaluation

- 3 other tools [Xing and Stroulia 05]
[Weißgerber and Diehl 06] [S. Kim, Pan, and Whitehead 05]
- Evaluation data set (E)
- Precision
 $(|M \cap E| / |M|)$
- Recall
 $(|M \cap E| / |E|)$
- Conciseness

Comparison: Recall & Precision

	programs	Other's Recall	Our Recall	Other's Prec.	Our Prec.
[Xing & Stroulia 05]	jfreechart 18 releases	92%	98%	99%	97%
[Weissgerber & Diehl 06]	jEdit 2715 check-ins	72%	96%	93%	98%
	Tomcat 5096 check-ins	82%	89%	89%	93%
[Kim, Pan & Whitehead 05]	jEdit 1189 check-ins	70%	96%	98%	96%
	ArgoUML 4683 check-ins	82%	95%	98%	94%



Comparison: Recall & Precision

	programs	Other's Recall	Our Recall	Other's Prec.	Our Prec.
[Xing & Stroulia 05]	jfreechart 18 releases	92%	98%	99%	97%
[Weissgerber & Diehl 06]	<p>Precision: 93-98%</p> <p>Recall: 89-98%</p> <p>6-26% higher recall with roughly the same precision</p>				98%
					93%
[Kim, Pan & Whitehead 05]					ArgoUML 4683 check-ins

Comparison: Conciseness

	programs	Other's Results	Our Results	Our Improvement
[Xing & Stroulia 05]	jfreechart 18 releases	4004 refactorings	939 rules	77% decrease in size
[Weissgerber & Diehl 06]	jEdit 2715 check-ins	1218 refactorings	906 rules	26% decrease in size
	Tomcat 5096 check-ins	2700 refactorings	1033 rules	62% decrease in size
[Kim, Pan & Whitehead 05]	jEdit 1189 check-ins	1430 matches	1119 rules	22% decrease in size
	ArgoUML 4683 check-ins	3819 matches	2127 rules	44% decrease in size



Comparison: Conciseness

	programs	Other's Results	Our Results	Our Improvement
[Xing & Stroulia 05]	jfreechart 18 releases	4004 refactorings	939 rules	77% decrease in size
[Weissgerber & Diehl 06]				% decrease in size
				% decrease in size
[Kim, Pan & Whitehead 05]	ArgoUML 4683 check-ins	3819 matches	2127 rules	44% decrease in size

22-77% reduction in the size of matching results

Summary of Code Matching

- Our change-rules **concisely** capture **API-level changes** and identify anomalies to systematic changes
- By inferring such rules, we find method-header level matches with **high recall and precision**

Outline

- Empirical Analyses of Code Clone Evolution
- Automatic Inference of High-Level Change Descriptions
 - Changes to API name and signature
 - Changes to Code Elements and Structural Dependencies (Logical Structural Diff)
- Future Directions

Research Question

“What is a concise change representation beyond API-level refactorings?”

```
public class CmiRegistry implements NameService {  
  
    public void setPort (int p) {  
        ...  
-       SQL.exec(query)  
+       SafeSQL.exec(query)  
  
    }  
  
}
```

```
public class JacORB implements NameService  
{  
    public void setPort (int p) {  
-       if (TraceCarol.isDebug()) {  
        ...  
-       SQL.exec(query)  
+       SafeSQL.exec(query)  
    }  
    ...  
}
```

```
public class LmiRegistry extends AbsRegistry implements NameService {  
-   private int port = ...  
-   private String host = null  
    public void setPort (int p) {  
        ...  
-   SQL.exec(query)  
+   SafeSQL.exec(query)  
    }  
    public int getPort() {  
        return port;  
    }  
    public void setHost(String host)
```

Logical Structural Diff

Abstraction Level	Code elements and structural dependencies (package, type, method, field, overriding, subtyping, method call, field access, and containment)
Scope	Conjunctive logic literal
Transformation	Structural differences Account for changes in method-bodies as well as at a field level
Example Rule	<pre>past_method(m, t) ^ past_subtype("Factory", t) ^ past_calls(m, "render()") => added_calls(m, "Util.log()")</pre>

Logical Structural Diff Algorithm

Output: **logic rules** and **facts** that describe changes to **code elements and structural dependencies**

1. Extract a set of facts from a program using JQuery [Jensen & DeVolder 03]
2. Compute fact-level differences
3. Learn Datalog rules using an inductive logic programming algorithm

Logical Structural Diff Output

- “All `setHost` methods in `Service`'s subclasses in the old version deleted calls to `SQL.exec` except the `setHost` method in the `NameSvc` class.

```
past_subtype("Service", t) ^
past_method(m, "setHost", t)
⇒ deleted calls(m, "SQL.exec")
except t="NameSvc"
```

Quantitative Assessment of *LSDiff*

- 75% of fact-level differences are explained by rules.
- vs. fact-level delta: 9.3 times more concise
- vs. fact-level delta: 9.7 additional contextual facts
- vs. *Diff*: on average 7 rules and 27 facts for 997 lines of changes across 16 files

Focus Group Study

- Pre-screener survey
- Participants: five professional software engineers
 - industry experience ranging from 6 to over 30 years
 - use *diff* and *diff*-based version control system daily
 - review code changes daily except one who did weekly
- One hour structured discussion
 - I worked as a moderator. We also had a note-taker transcribe the discussion. Discussion was audio-taped and transcribed.

Focus Group Hands-On Trial

Carol Revision 430.

SVN check-in message: Common methods go in an abstract class. Easier to extend/maintain/fix

Author: benoif @ Thu Mar 10 12:21:46 2005 UTC

723 lines of changes across 9 files (2 new files and 7 modified files).

Overview

Generated based on LSDiff output.

Inferred Rules

1	(50/50)	<u>By this change, six classes inherit many methods from AbsRegistry class.</u>
2	(32/32)	<u>By this change, six classes implement NameService interface.</u>
3	(6/8)	<u>All methods that are included in JacORB CosNaming class and NameService interface are deleted except start and stop methods.</u>
4	(5/6)	<u>All host fields in the classes that implement NameService interface got deleted except LmiRegistry class.</u>
5	(5/6)	<u>All port fields in the classes that implement NameService interface got deleted except LmiRegistry class.</u>
6	(5/6)	<u>All getHost methods in the classes that implement NameService interface got deleted except LmiRegistry class.</u>

<http://www.cs.washington.edu/homes/miryung/LSDiff/carol429-430.htm>

Focus Group Hands-On Trial

```
46: public class IIOPCosNaming extends AbsRegistry implements NameService {
47:
48:     /**
49:      * Default port number ( 12350 for default)
50:      */
```

All DEFAULT PORT NUMBER fields are added fields except JacORBCosNaming class.

```
51:     private static final int DEFAULT_PORT DEFAULT_PORT_NUMBER = 12350;
```

```
52:
53:     /**
54:      * Sleep time to wait
55:      */
```

```
56:     private static final int SLEEP_TIME = 2000;
```

```
57:
58:     /**
```

```
59:     * port number
```

```
60:     */
```

All port fields in the classes that implement NameService interface got deleted except LmiRegistry

```
61:     private int port = DEFAULT_PORT;
```

```
62:
```

```
63:     /**
```

```
64:     * Hostname to use
```

```
65:     */
```

All host fields in the classes that implement NameService interface got deleted except LmiRegistry

```
66:     private String host = null;
```

Show related changes

<http://www.cs.washington.edu/homes/miryung/LSDiff/carol429-430.htm>

Focus-Group Participants' Comments

“You can’t infer the intent of a programmer,
but this is pretty close.”

“This ‘except’ thing is great!”

“This is cool. I’d use it if we had one.”

“This is a definitely winner tool.”

Focus-Group Participants' Comments

“This looks great for big architectural changes, but I wonder what it would give you if you had lots of random changes.”

“This wouldn't be used if you were just working with one file.”

“This will look for relationships that do not exist.”

Summary of Logical Structural Diff

- We extended our rule-based approach to infer systematic changes **within method bodies**
- LSDiff produces 9.3 times more concise results by identifying 75% of structural differences as systematic changes
- LSDiff complements *diff*
 - by grouping systematic structural differences
 - by detecting potential missed updates.

Outline

- Empirical Analyses of Code Clone Evolution
- Automatic Inference of High-Level Change Descriptions
 - Changes to API name and signature
 - Changes to Structural Dependencies
- Future Directions

Next Steps

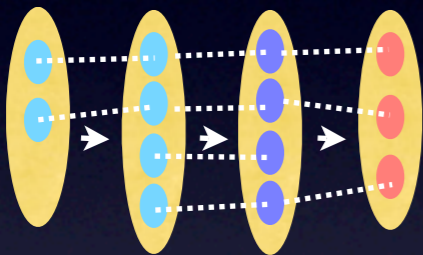
- Present change-rules as English sentences
- Develop higher-order representations
- Use change-rules to improve regression testing
- Search program changes of interests in a source code repository by evaluating programmer-provided rules

My long-term vision is to help programmers
by making software change a first class entity

- Changes in models, requirements, and run-time behavior
- Use change history to help programmers make decisions
 - “When and how should I refactor my program?”

Contributions

Analyses of Software Evolution



- Disproving conventional wisdom about clones
- Insights into systematicness of high-level changes

Automatic Inference of High-Level Change Descriptions



- Rule-based change representations
- Rule learning algorithms

Acknowledgments

David Notkin, Dan Grossman, Rob DeLine, Ken Bube, Gail Murphy, Andreas Zeller, Vibha Sazawal, Tessa Lau, Larry Bergman, Marius Nita, Jonathan Beall, Sunghun Kim, Peter Weissgerber, Zhencheng Xing, Jennifer Bevan, friends@WASP, friends@UW CSE, friends@NCBC, Rob, Cathy, Jackie, Sammy and Alan

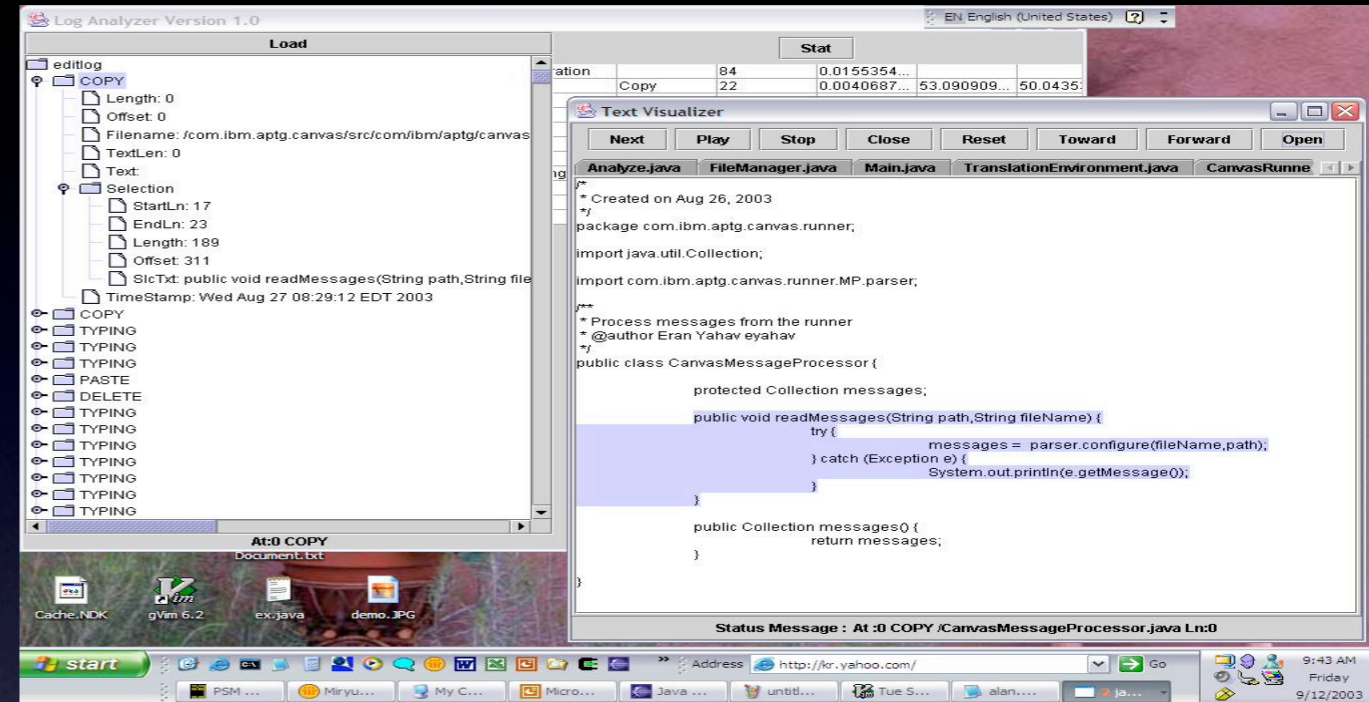
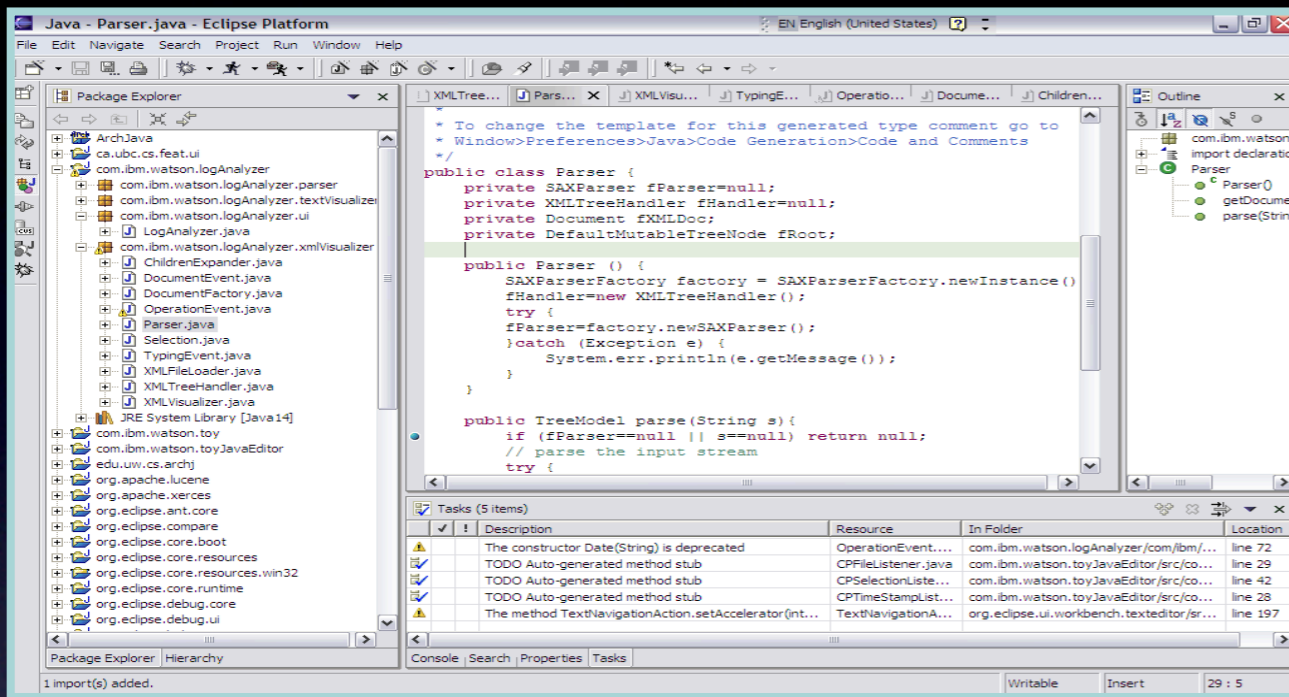
Questions?

Details of Study at IBM

	direct observation	observation using logger and replayer
Subjects	researchers and summer students at IBM T.J. Watson	
No. of subjects	4	5
Hours	about 10 hours	about 50 hours (about 300+ hours of edit logs)
Interviews	questions asked during observation	twice after analysis (30 mins - 1 hour / each)
Programming Languages	Java, C++, and Jython	Java

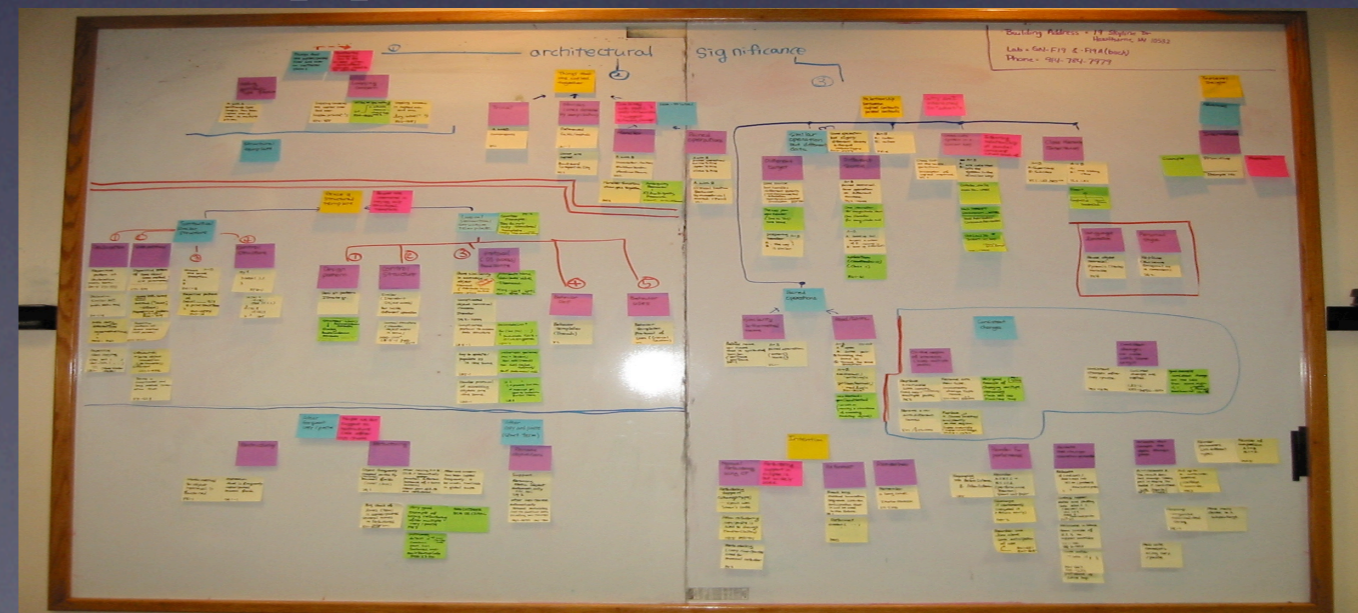
I. Capture Edit Operations from IDE

2. Replay and Reconstruct Editing Context



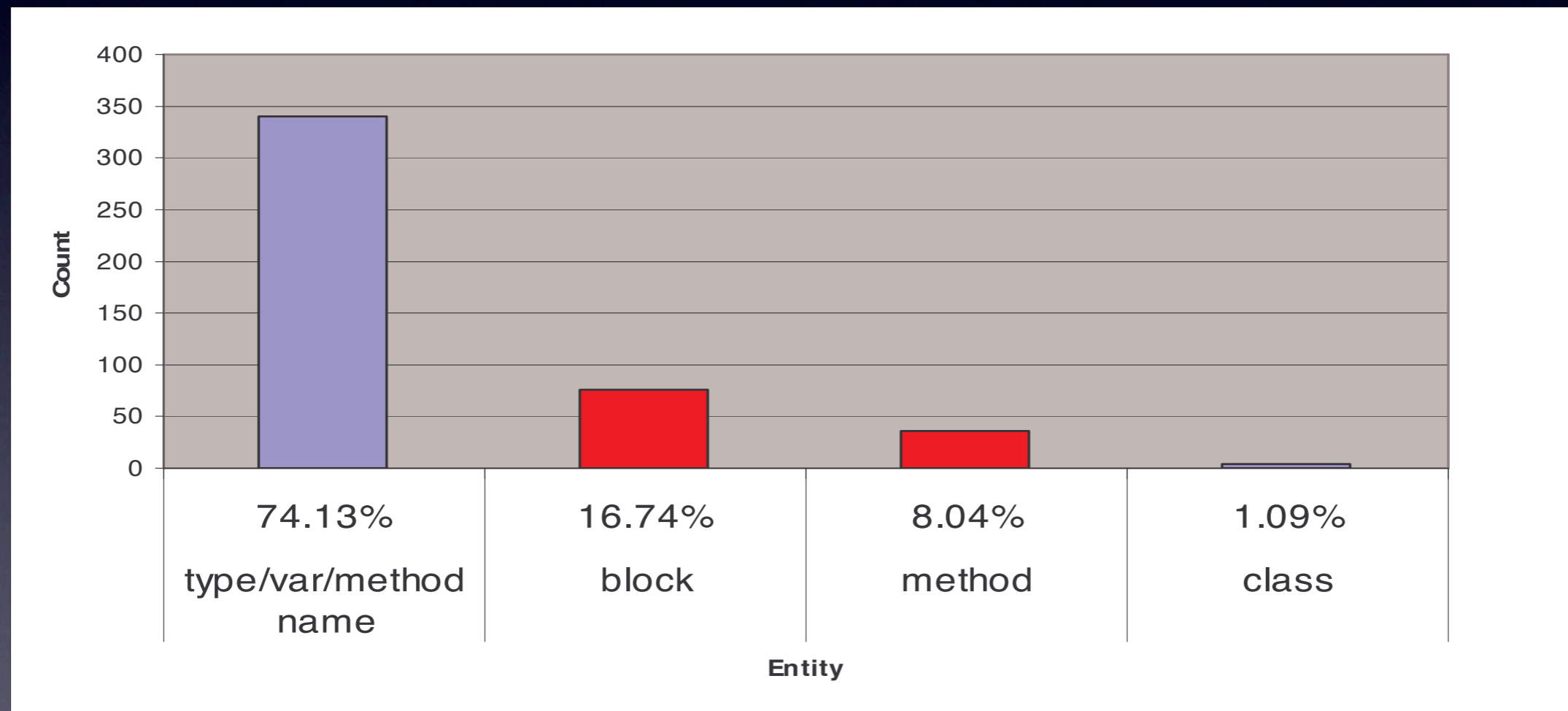
3. Semi-Structured Interviews

4. Create a Taxonomy of Copy&Paste Patterns



Copy and Paste Study Details

- 460 copy and paste instances

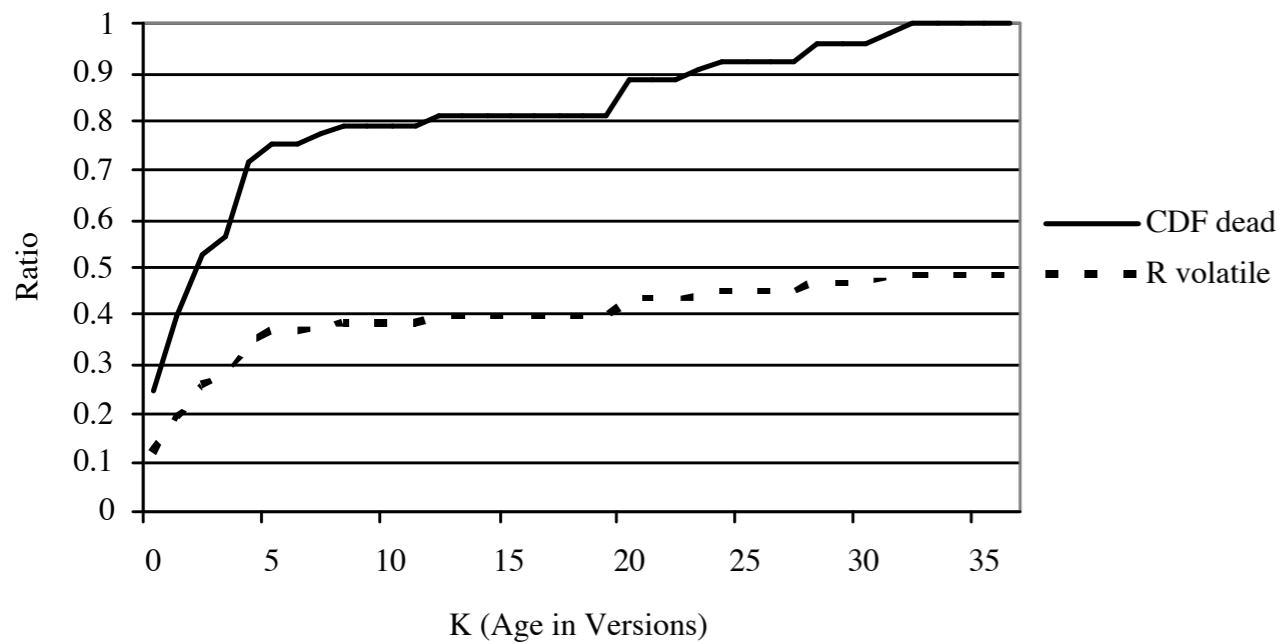


Two Java Open Source Projects

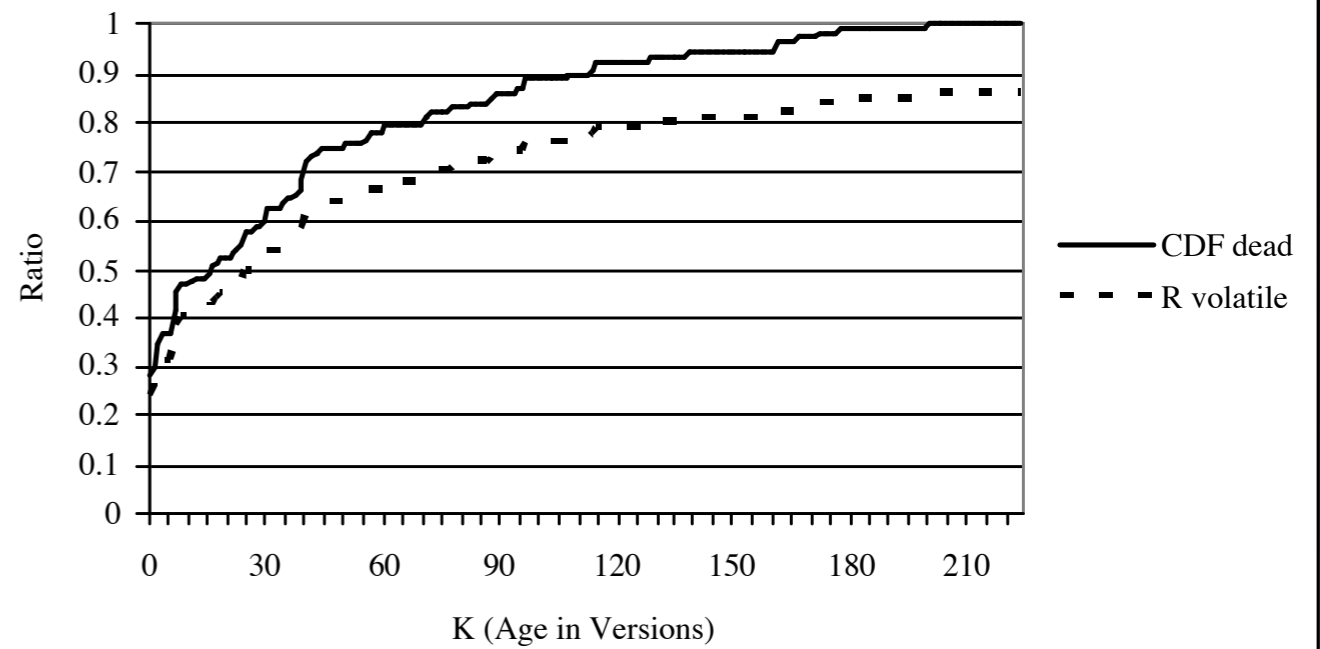
Program	<i>Carol</i>	<i>Dnsjava</i>
LOC	7878 ~ 23731	5756 ~ 21188
Duration	2 years 2 months	5 years 8 months
Check-Ins	164	905

Short-Lived Clones

carol, sim_th=0.3



dnsjava, sim_th=0.3



Why Do Clones Disappear?

Reasons	<i>carol</i>	<i>dnsjava</i>
Divergent changes	26%	34%
Refactoring or removal	67%	45%
Cut off by the threshold	7%	21%

Contrary to conventional wisdom, immediate refactoring may be unnecessary or counterproductive in some cases.

Locally Unfactorable Clones

```
public void exportObject(Remote obj) throws
RemoteException{
    if (TraceCarol.isDebugEnabled()) {
        TraceCarol.debugRmiCarol(
            "MultiPRODelegate.exportObject("...
    }
    try {
        if (init) {
            for (Enumeration e = activePtcls.element(); ..
                ((ObjDlgt)e.nextElement()).exportObject
                (obj);
            }
        }
        catch (Exception e) {
            String msg ="exportObject(Remote obj) fail";
            TraceCarol.error(msg,e);
            throw new RemoteException(msg);
        }
    }
```

```
public void unexportObject(Remote obj) throws
NoSuchObjectException{
    if (TraceCarol.isDebugEnabled()) {
        TraceCarol.debugRmiCarol(
            "MultiPRODelegate.unexportObject("...
    }
    try {
        if (init) {
            for (Enumeration e = activePtcls.element(); ..
                ((ObjDlgt)e.nextElement()).unexportObject
                (obj);
            }
        }
        catch (Exception e) {
            String msg ="unexportObject(Remote obj)
            fail";
            TraceCarol.error(msg,e);
            throw new NoSuchObjectException(msg);
        }
    }
```

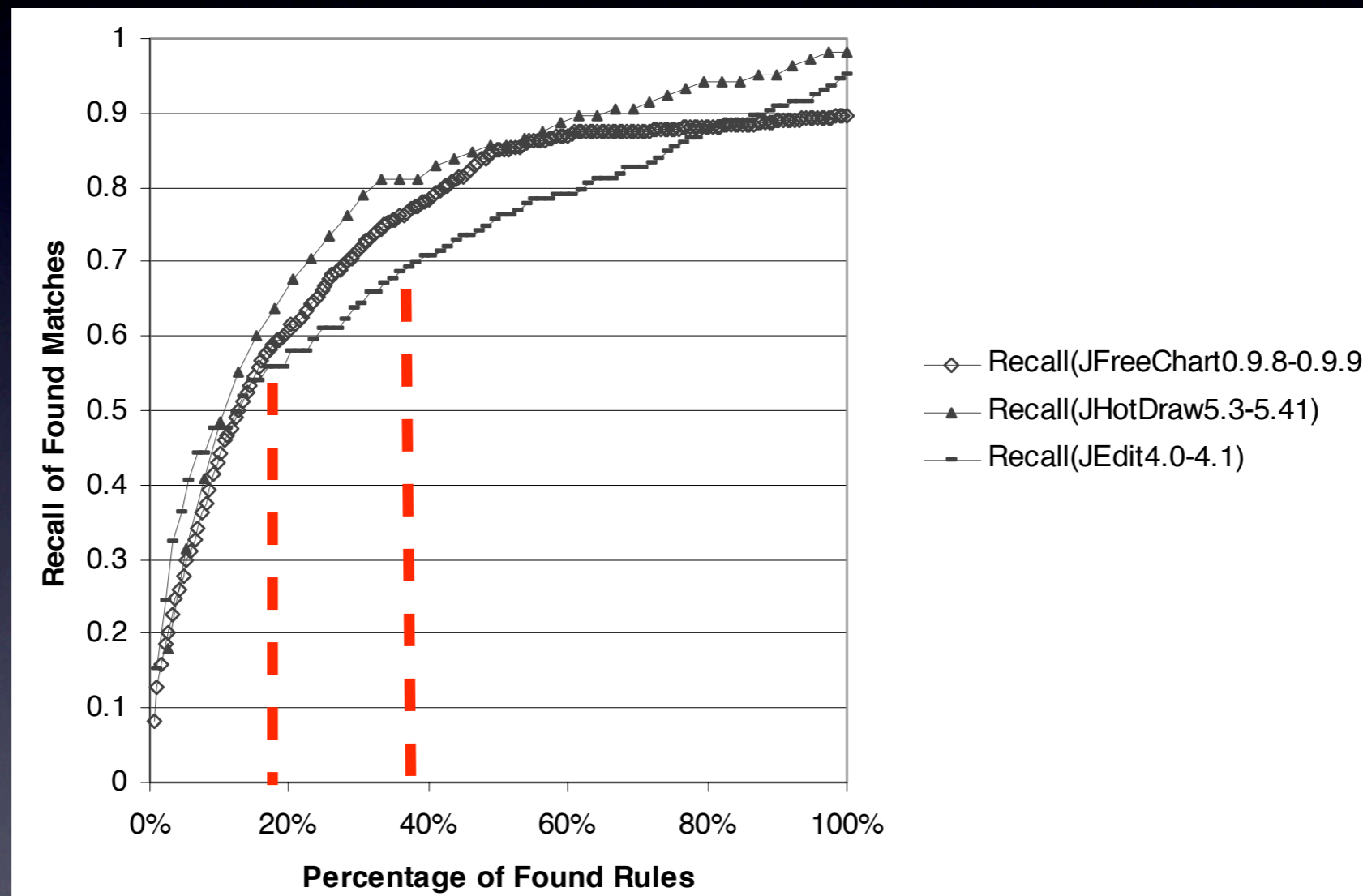
Complexity

- Candidate rules (CR)
 - $D * 2^k * T$, where k is the max number of tokens in a seed, T is the number of transformations, and D is the number of Java methods in the old version
- Greedy Algorithm Complexity = $D * CR^2$

Rule-based Matching Results for Three Release Archives

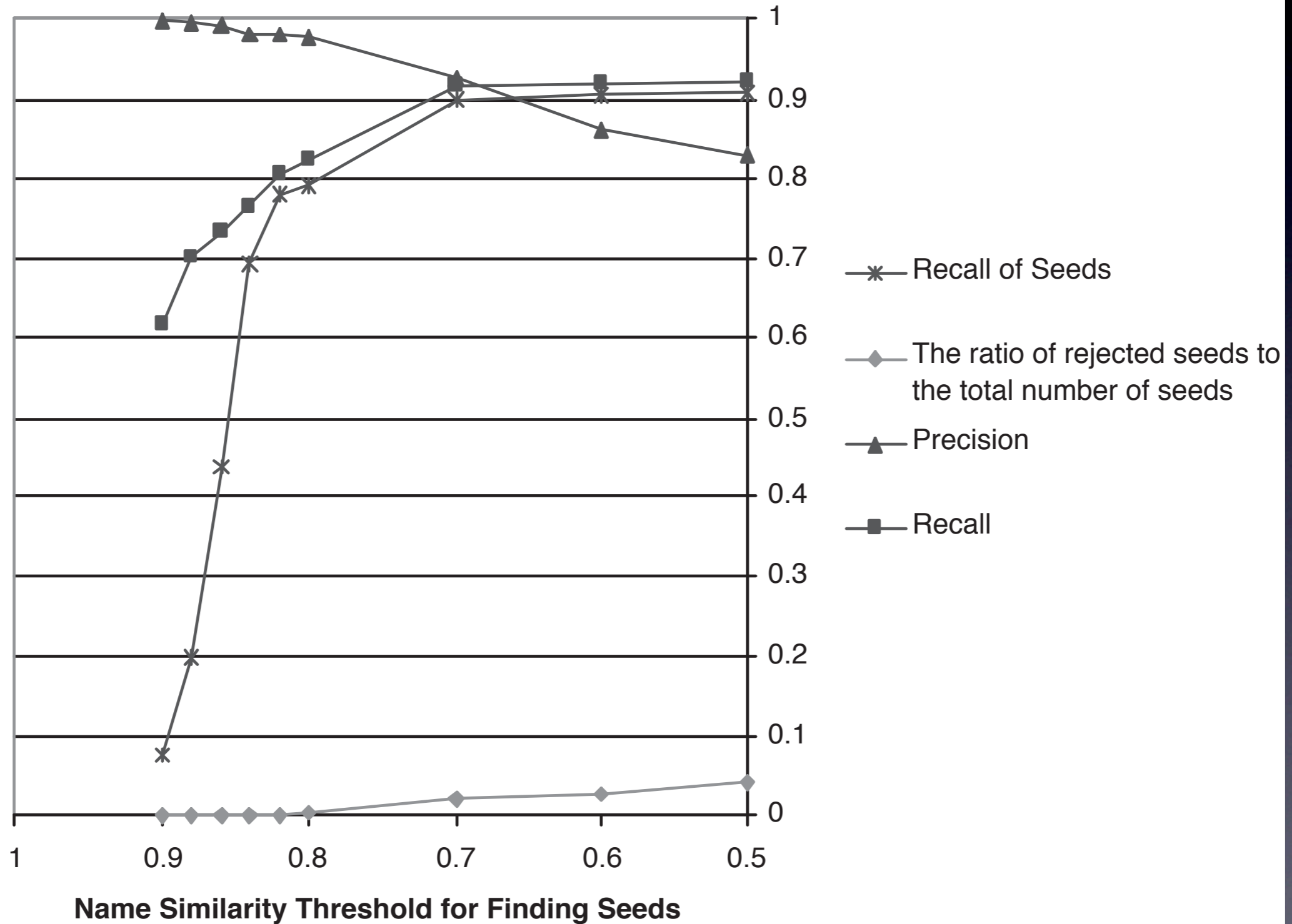
	JFreeChart	jHotDraw	jEdit
	(17 release pairs)	(4 release pairs)	(4 release pairs)
Precision Median (Min ~ Max)	94% (78~100%)	99% (82~100%)	93% (87~95%)
Recall Median (Min ~ Max)	93% (70~100%)	99% (92~100%)	98% (95~100%)
M/R ratio Median (Min ~ Max)	3.50 (1.20~135.23)	2.54 (1.00~244.26)	1.73 (1.23~2.39)

Rule-based Matching Results for Three Release Archives



Top 20% of the rules find over 55% of the matches.
Top 40% of the rules find over 70% of the matches.

Seed Threshold



JFreeChart (www.jfree.org/jfreechart)

The actual release numbers are prefixed with 0.9.

	O	N	O ∩ N	Rule	Match	Prec.	Recall	M/R	Time
4→5	2925	3549	1486	178	1198	0.92	0.92	6.73	21.01
5→6	3549	3580	3540	5	6	1.00	1.00	1.20	<0.01
6→7	3580	4078	3058	23	465	1.00	0.99	20.22	1.04
7→8	4078	4141	0	30	4057	1.00	0.99	135.23	43.06
8→9	4141	4478	3347	187	659	0.91	0.90	3.52	22.84
9→10	4478	4495	4133	88	207	0.99	0.93	2.35	0.96
10→11	4495	4744	4481	5	14	0.79	0.79	2.80	<0.01
11→12	4744	5191	4559	61	113	0.78	0.79	1.85	0.40
12→13	5191	5355	5044	10	145	1.00	0.99	14.50	0.11
13→14	5355	5688	5164	41	134	0.94	0.86	3.27	0.43
14→15	5688	5828	5662	9	21	0.90	0.70	2.33	0.01
15→16	5828	5890	5667	17	77	0.97	0.86	4.53	0.32
16→17	5890	6675	5503	102	285	0.91	0.86	2.79	1.30
17→18	6675	6878	6590	10	61	0.90	1.00	6.10	0.08
18→19	6878	7140	6530	98	324	0.93	0.95	3.31	1.67
19→20	7140	7222	7124	4	14	1.00	1.00	3.50	<0.01
20→21	7222	6596	4454	71	1853	0.99	0.98	26.10	62.99
MED						0.94	0.93	3.50	0.43
MIN						0.78	0.70	1.20	0.00
MAX						1.00	1.00	135.23	62.99

Comparison with Textual Delta

	<i>diff output</i>		<i>LSDiff output</i>	
Median	Changed LOC	Touched Files	Rules	Facts
Carol 10 check-ins	626	11	5	16
Dnsjava 29 releases	354	9	3	23
Lsdiff 10 versions	227	6	1	8
Total (Median)	344	9	2	17

Comparison with Δ FB

	FB _o	FB _n	Δ FB	Rule	Fact	Cvrg.	Csc.	Ad'l.
Carol								
Min	3080	3452	15	1	3	59%	2.3	0.0
Max	10746	10610	1812	36	71	98%	27.5	19.0
Med	9615	9635	97	5	16	87%	5.8	4.0
Avg	8913	8959	426	10	20	85%	9.9	5.5
dnsjava								
Min	3109	3159	4	0	2	0%	1.0	0.0
Max	7200	7204	1500	36	201	98%	36.1	91.0
Med	4817	5096	168	3	24	88%	4.8	0.0
Avg	5144	5287	340	8	37	73%	8.4	14.9
LSD tool								
Min	8315	8500	2	0	2	0%	1.0	0.0
Max	9042	9042	396	6	54	97%	28.9	12.0
Med	8732	8756	142	1	11	91%	9.8	0.0
Avg	8712	8783	172	2	17	68%	11.2	2.3
Med	6650	6712	132	2	17	89%	7.3	0.0
Avg	6632	6732	302	7	27	75%	9.3	9.7

Impact of Varying Parameters

		Rule	Fact	Cvrg.	Csc.	Ad'l.	Time(Min)
m	1	39.6	0	100%	7.4	10.1	2.0
	2	14.6	13.1	92%	10.6	7.4	11.2
	3	9.9	20.4	85%	9.9	5.5	9.1
	4	7.7	25.7	82%	9.1	5.4	8.7
	5	5.7	30	80%	8.5	3.5	7.8
a	0.5	11.1	15.6	89%	10.6	2.1	6.8
	0.625	9.7	17.2	88%	11.0	4.0	7.3
	0.75	9.9	20.4	85%	9.9	5.5	9.0
	0.875	10.8	24.2	78%	8.6	9.1	12.7
	1	13.3	26.2	78%	7.9	12.5	16.5
k	1	7.5	33.8	78%	7.2	0.4	0.7
	2	9.9	20.4	85%	9.9	5.5	9.1

LSDiff Rule Styles

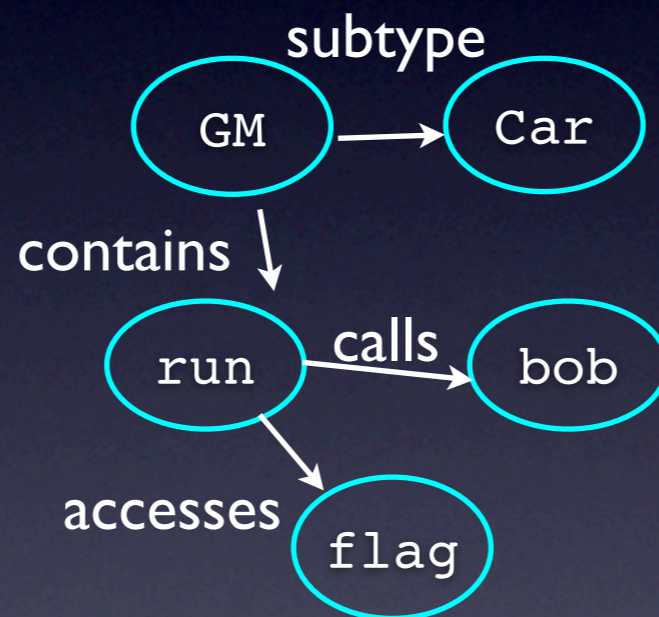
Horn Clause: $A(x) \wedge B(x, y) \wedge C(y) \Rightarrow D(x, y)$

Rule Styles	High-level Change Patterns	Example
<code>past_* => deleted_*</code>	dependency removal, feature deletion, etc.	<code>past_calls(m, "DB.exec") => deleted_calls(m, "DB.exec")</code>
<code>past_* => added_*</code>	consistent updates to clones, etc.	<code>past_accesses("Log.on", m) => added_calls(m, "Log.trace")</code>
<code>current_* => added_*</code>	dependency addition, feature addition, etc.	<code>current_method(m, "getHost", t) ^ current_subtype("Svc", t) => added_calls(m, "Log.trace")</code>
<code>deleted_* => added_*</code> <code>added_* => deleted_*</code>	related code change, API replacement, etc.	<code>deleted_method(m, "getHost", t) => added_inheritedfield("getHost", "Svc", t)</code>

Step I. Extract Facts

```
class GM extends Car
  void run(int c) {
    if (Util.flag)...
    bob();}
}
```

Program



```
type("GM")
subtype("Car", "GM")
method("GM.run", "run", "GM")
accesses("Util.flag", "GM.run")
calls("GM.run", "GM.bob")
```

Fact-base

A fact-base program representation approach has been used by many tools such as JQuery [Jensen&DeVolder 03], CodeQuest [Hajiev et. al. 06], Grok [Holt et. al.], etc.

Step 2. Compute Fact-Level Differences

Old Program (FBo)

past_*

```
subtype  
("Svc", "X")  
subtype  
("Svc", "Y")  
subtype  
("Svc", "Z")  
subtype  
("Svc", "NameSvc")
```

New Program (FBn)

current_*

```
subtype  
("Svc", "X")  
...  
method  
("exec", "X")  
method  
("exec", "Y")  
...
```

Differences (ΔFB)

added_* / deleted_*

```
added_method  
("exec", "X")  
added_method  
("exec", "Y")  
added_method  
("exec", "Z")
```

set difference

Step 3. Learn Rules

- Our rule learner uses a **bounded depth search** algorithm that finds **Datalog rules** in a domain specific form.
- We have input parameters that determine the **validity** of a rule.
 - a : accuracy
 - m : min support

Example.

```
past_calls (x, "foo")  
=> deleted_calls(x, "foo")  
(8/10) a = 0.80, m=8, k=1.
```

Step 3. Learn Rules

```
R := {} // a set of ungrounded rules.
D := reduced  $\Delta$ FB using default winnowing rules
L := {} // a set of valid learned rules.
for each antecedent size,  $i = 0 \dots k$  :
    R := extend all rules in R by adding all
    possible literals.
    for each ungrounded rule,  $r$ :
        for each possible grounded rule  $g$  of  $r$ :
            if (g is valid) L := L U g.
R := select the best  $\beta$  rules in R.
D := D - { facts covered by L}
```

Change Rule Selection Problem

Setting: Suppose that there exist a domain $M = \{m_1, m_2, \dots, m_p\}$, a codomain $N = \{n_1, n_2, \dots, n_q\}$, and a set S of sets $= \{s_1, s_2, s_3, \dots, s_n\}$, where each set s_i is a set of tuples that form a functional binary relation from a domain M' to a codomain N where $M' \subset M$.

Definition: A set s_i has a **conflict** (a, b) with a set s_j if $(a, b) \in s_i$ and $(a, b') \in s_j$ where $b \neq b'$.

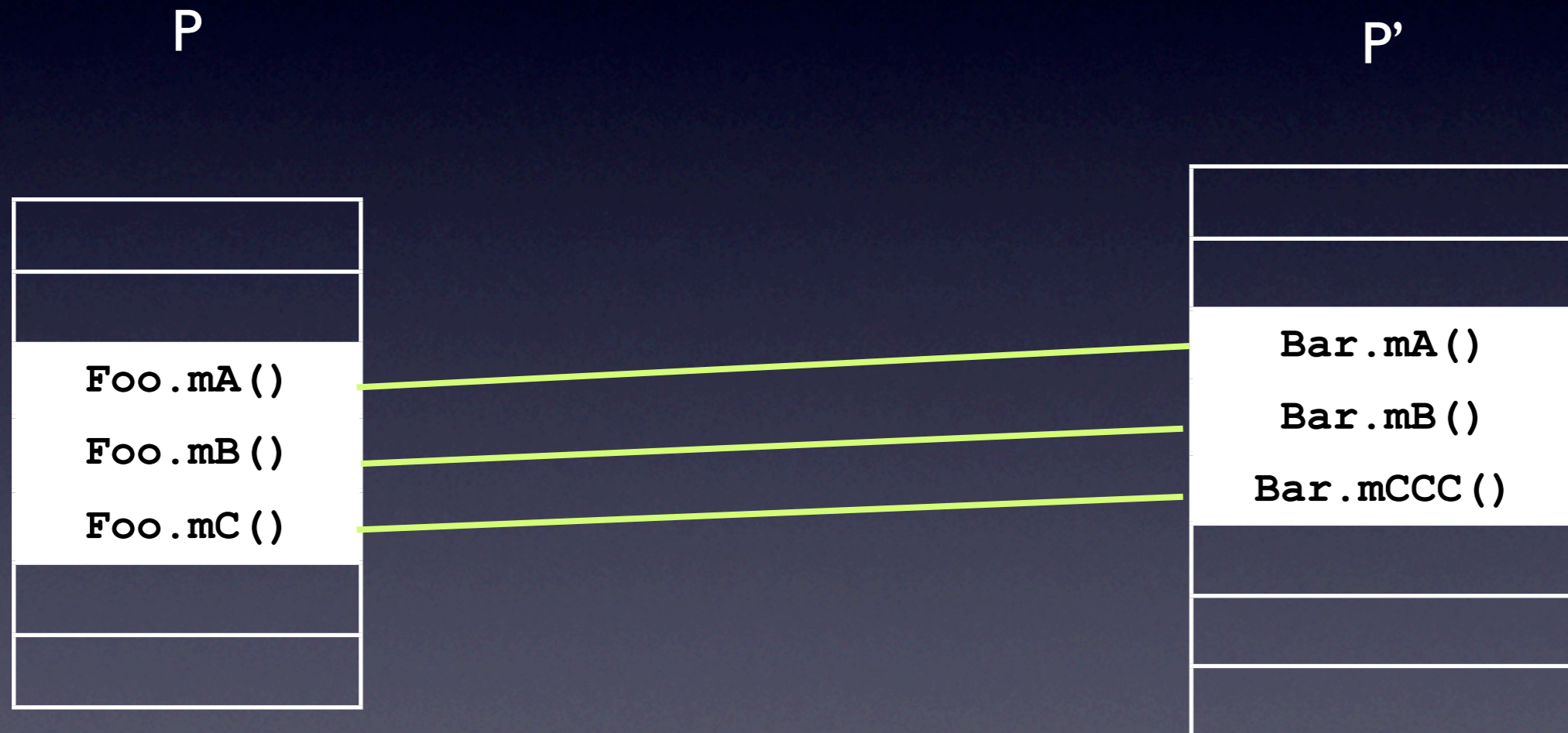
Problem: Given S and e ($0 < e < 1$), find a solution c , where $c \subset S$ and each set s_i in c has at most $e \times |s_i|$ conflicts with other sets in c .

A solution o is an optimal solution if o covers the maximum number of elements in M among all solutions and if o has the minimum cardinality among the solutions that cover the same maximum number of elements in M .

Question: Suppose that the greedy algorithm outputs g , while the optimal solution is o . What is the size of M 's elements covered by g in relation to M 's elements covered by o ? What is $|g|$ in relation to $|o|$?

Limitation of Commutable Change Rules

Rename class **Foo** to **Bar**.
Rename **Bar.mC** to **Bar.mCCC**.



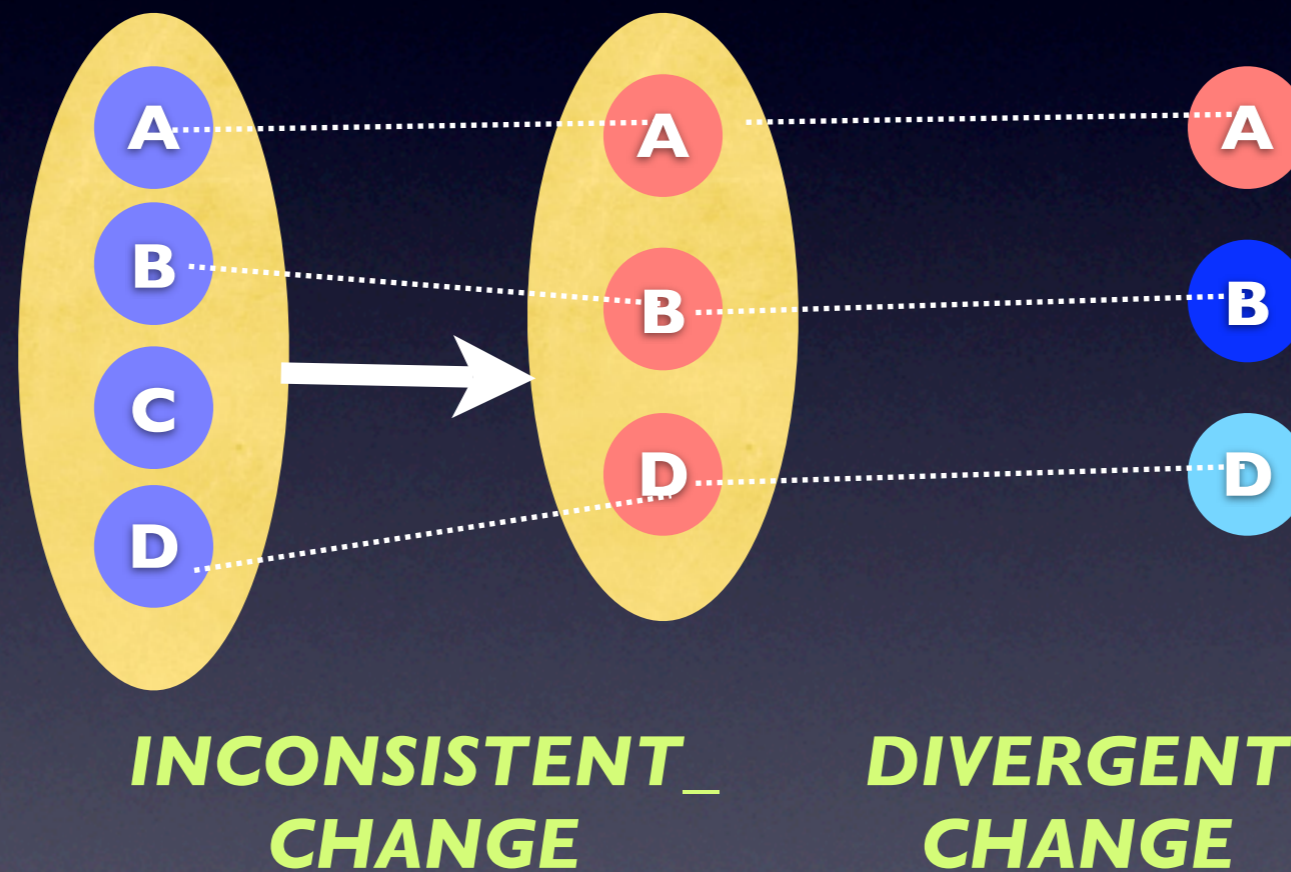
```
for all x in Foo.*(*), classReplace(x, Foo, Bar): 2/3 ✗
```

```
for all x in Foo.mc(), classReplace(x, Foo, Bar) ^ procReplace(x, mC, mCCC) : 1/1 ✓
```

Pros and Cons of Focus-Group Study Method

- + A relative short time period to coordinate, conduct, analyze a study.
- + Probing and clarification are easy
- *Recruiting is difficult*
- *Responses are not independent as it is done in a group setting*

Divergent Change vs. Inconsistent Change Pattern



Divergent change means that a clone group does not exist in the next group. Inconsistent Change is the same as Divergent Change when a clone group has two members.