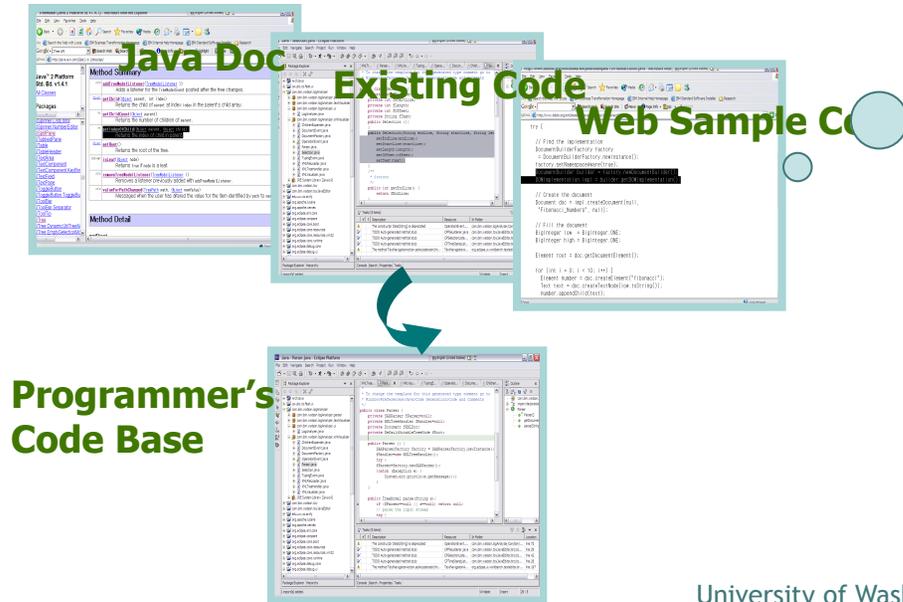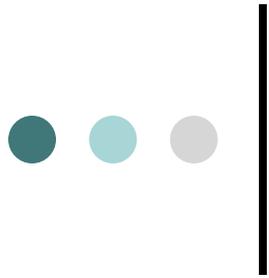# An Ethnographic Study of Copy and Paste Programming Practices in OOPL

Miryung Kim[1], Lawrence Bergman[2], Tessa Lau[2], and David Notkin[1]

Department of Computer Science and Engineering
University of Washington[1],
IBM T.J. Watson Research Center[2]

# Conventional Wisdom

Common but
Bad
Programming Practice

Java Doc

Existing Code

Web Sample Code

Programmer's
Code Base

University of Washington
IBM T.J. Watson Research Center
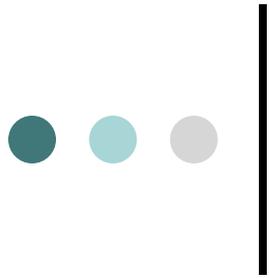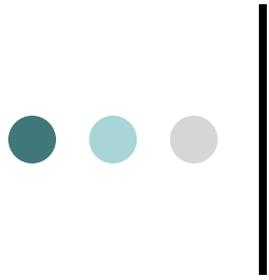
# Contribution

- We address implications of copy and paste (C&P) programming practices.
  - Not only about saving typing.
  - C&P capture design decisions.
  - Programmers actively employ C&P history.
  - With tool support, programmers' intent of C&P can be expressed in a safer and more efficient manner.

# Research Questions

- What are C&P usage patterns?
- Why do people copy and paste code?
- What kind of tool support is needed for C&P usage patterns?

# Outline

- Ethnographic Study: Observation and Analysis
- Taxonomy
- Insights and Tool Ideas

University of Washington
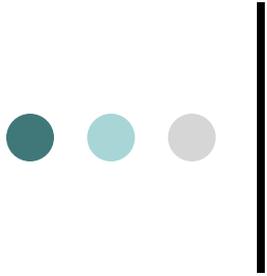IBM T.J. Watson Research Center

# Observation

- preliminary approach
  - direct observation
  - questions asked during observation
  - easy to identify intentions
  - unnatural coding behavior

- final approach
  - logging editing operations with an instrumented text editor
  - replaying off-line
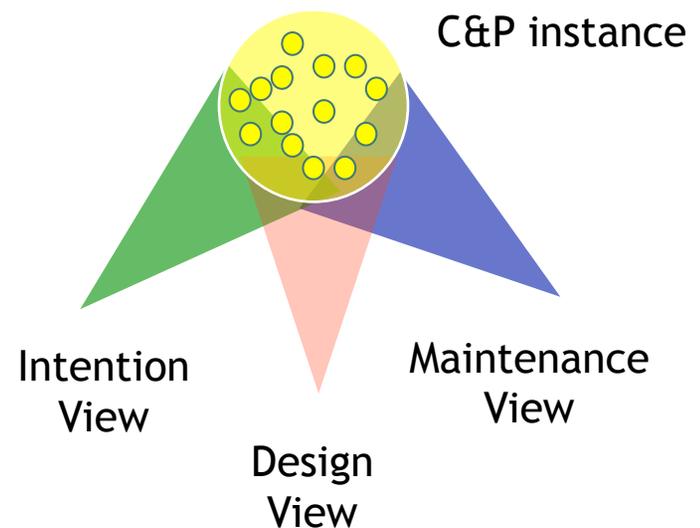  - interviews
  - non-intrusive observation

# Study Setting

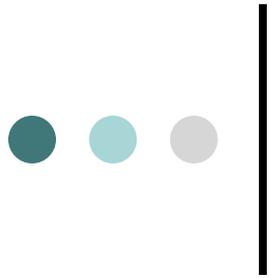| | Direct Observation | Observation using a logger and a replayer |
|---|---|---|
| Subjects | researchers and summer students at IBM T.J. Watson | |
| No. of Subjects | 4 | 5 |
| Hours | about 10 hrs | about 50 hrs |
| Interviews | questions asked during observation | twice after analysis (30 mins – 1 hour/ each) |
| Programming Languages | Java, C++, and Jython | Java |

# Analysis

o contextual inquiry
[Beyer98]

- affinity process:
developing
hypotheses from data
points

o data analysis from
multiple perspectives



C&P instance

Intention
View

Design
View

Maintenance
View

# Outline

- Ethnographic Study: Observation and Analysis
- Taxonomy
- Insights and Tool Ideas

# Programmers' Intentions

- relocate/ regroup/ reorganize
- reorder
- refactoring
- reuse as a structural template
  - syntactic template
  - semantic template

# Example – Syntactic Template

```
static {
        protectedClasses.add("java.lang.Object");
        protectedClasses.add("java.lang.ref.Reference
    $ReferenceHandler");
        protectedClasses.add("java.lang.ref.Reference");
        protectedClasses.add("java.lang.ref.Reference$1");
        protectedClasses.add("java.lang.ref.Reference$Lock");
        protectedMethods.add("java.lang.Thread<init>");
        protectedMethods.add("java.lang.Object<init>");

    protectedMethods.add("java.lang.Thread.getThreadGroup");
}
```

# Semantic Template

Intention

- design patterns

- control structures
  - if – then – else
  - loop construct

- usage of a module
  - data structure access protocols

# Example – Semantic Template: Usage of a Module

**Intention**

```
DOMNodeList *children = doc->getChildNodes();
int numChildren = children->getLength();

for (int i=0; i<numChildren; ++i)
   {
       DOMNode *child = (children->item(i));
       if (child->getNodeType() == DOMNode.ELEMENT_NODE)
       {
           DOMElement *element = (DOMElement*)child;
```

Code Snippets:
traverse over *Elements*
in a *Document*

# Design View

What are <span style="color:green">underlying design decisions</span> that <span style="color:green">induce</span> programmers to C&P in particular patterns?

- Why is text copied and pasted over and over in scattered places?

- Why are blocks of text copied together?

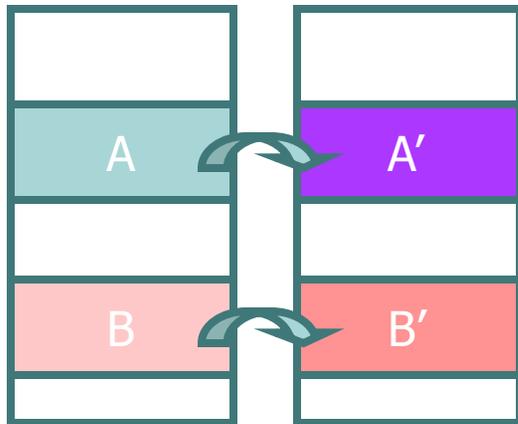- What is the relationship between copied text and pasted text?

# Why is text copied and pasted repeatedly?

**Design**

o lack of modularity

- crosscutting concerns

- example – logging concern

```
if (logAllOperations) {
        try {
        PrintWriter w = getOutput();
        w.write("$$$$$");
        ..
        } catch (IOException e) {
    }
}
```
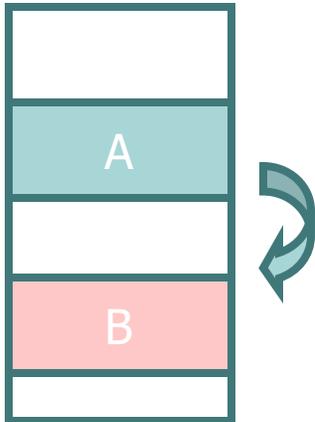
# Why are blocks of text copied together?



- comments
- references fields and constants
- caller method and callee method
- paired operations
  - openFile, closeFile, and writeToFile
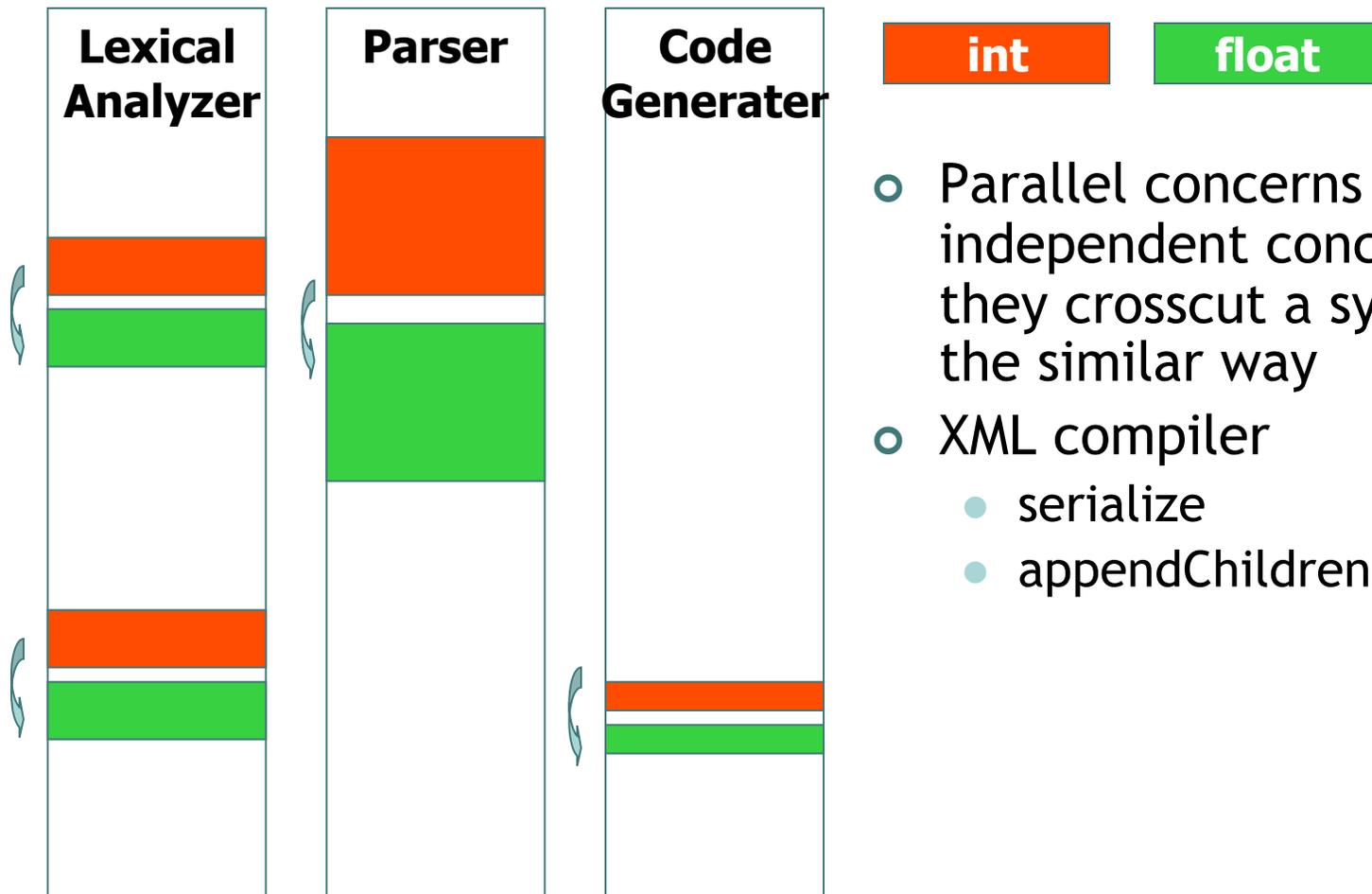  - enterCriticalSection, leaveCriticalSection

**Design**

# What is the relationship between copied and pasted text?

- type dependencies
- similar operations but different data structure
- parallel crosscutting concerns
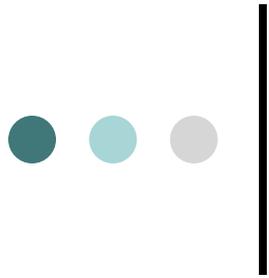  [Griswold01]

# Example - Parallel Crosscutting Concern

**Design**

| Lexical Analyzer | Parser | Code Generater |

**int** **float**

- Parallel concerns are independent concerns but they crosscut a system in the similar way
- XML compiler
  - serialize
  - appendChildren

University of Washington
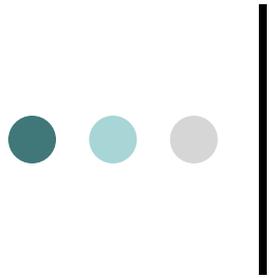IBM T.J. Watson Research Center

# Maintenance Tasks

**Maintenance**

- short term
  - Programmers modify a pasted block to prevent naming conflicts.
  - Programmers remove code fragments irrelevant to the pasted context.
- long term
  - Programmers restructure code after frequent copy and paste of a large text.
  - Programmers tend to apply consistent changes to the code from the same origin.

# Scope and Limitations

- programming languages
  - OOPL vs. functional PL
- development environment
  - Eclipse vs. other editors
- organization characteristics
  - team size, software lifecycle, etc
- duration of study
  - long term vs. short term

# Outline

- Ethnographic Study: Observation and Analysis
- Taxonomy
- Insights and Tool Ideas

# Insights

Limitations of particular programming languages produce unavoidable duplicates in a code base.

University of Washington
IBM T.J. Watson Research Center

# Insights

C&P dependencies are worth observing and maintaining

Tool requirements:

- visualize copied and pasted content
- explicitly maintain and represent C&P dependencies
- allow developers to communicate the intention behind C&P by annotation

University of Washington
IBM T.J. Watson Research Center

# Insights

> Programmers copy an entire code snippet because it contains the structural template that they intend to reuse.

Tool requirements:

- learn a relevant structural template
- assist to modify the portion that is not part of the structural template

# Insights

> Programmers use their memory of C&P history to determine when to restructure code.

Tool requirements:

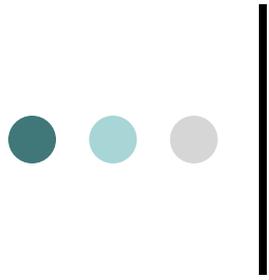- monitor evolution patterns, frequency, and size of code duplicates
- suggest refactoring

# Insights

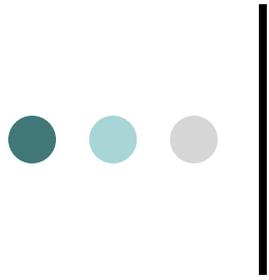> Code snippets originating from the same source are likely to be changed in similar ways.

Tool requirements:

- monitor evolution of structural template within code duplicates
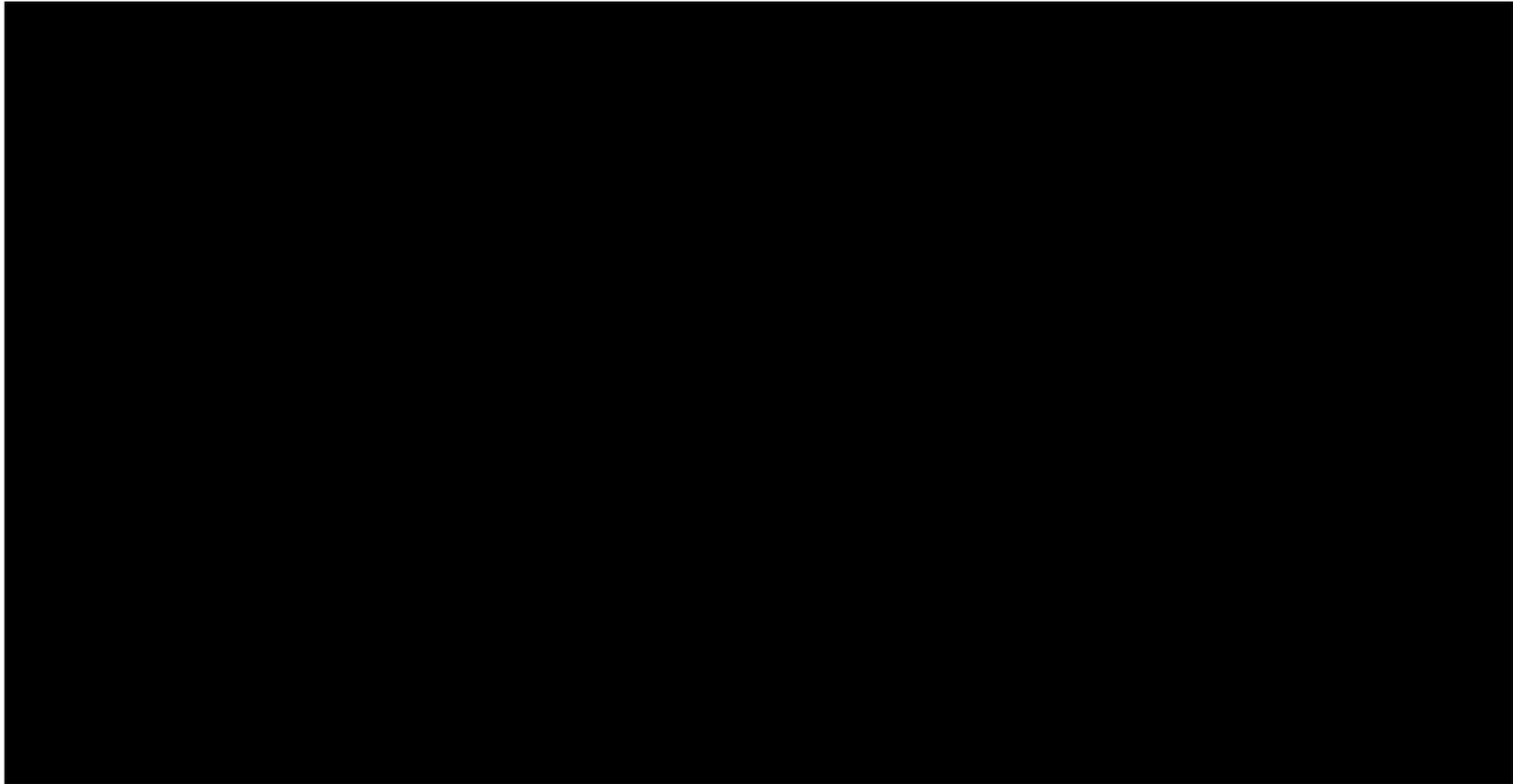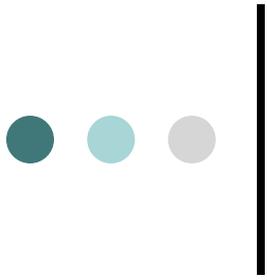- warn programmers when they attempts to change inconsistently

# Related Work

○ study of code reuse [Lange89, Rosson93]

○ information transparency [Griswold01]

○ clone detection [Balazinska02, Baker92, Baxter98, Ducasse99, Kamiya02, Komondoor01, Krinke01]

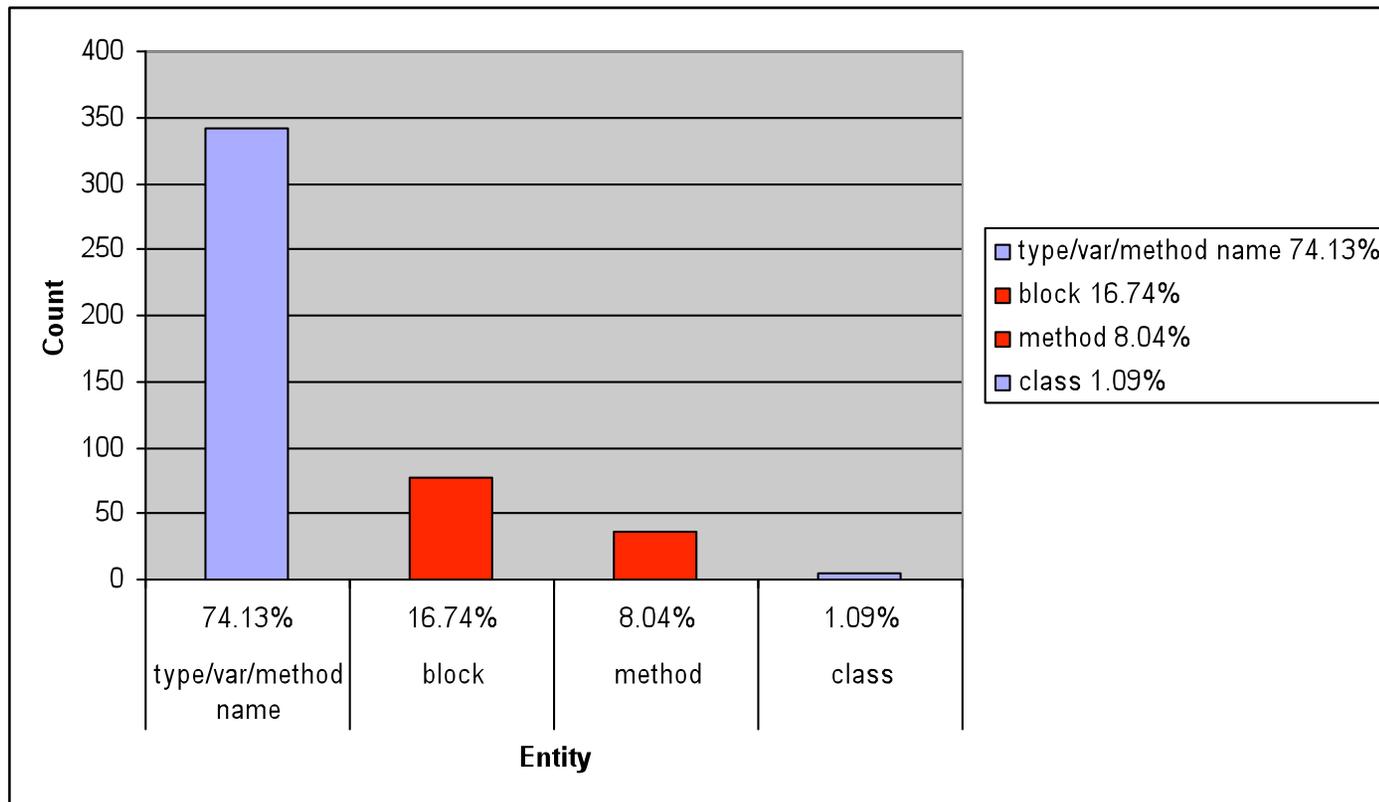○ clone evolution patterns [Lague96, Antoniol02, Rysselberghe04, Godfrey04]

# Conclusion

○ development of the instrumented editor and the replayer

○ study that systematically investigated C&P usage patterns and associated implications

○ proposal of SE tools based on our insights

University of Washington
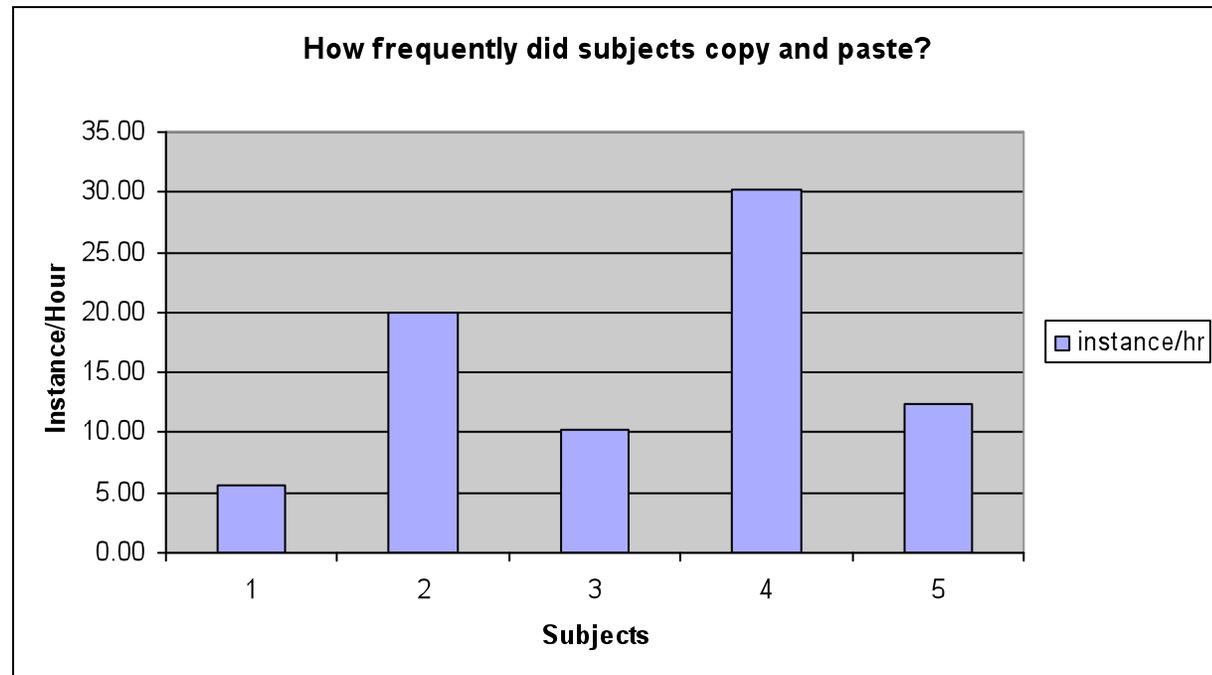IBM T.J. Watson Research Center

# What kind of code snippets do programmers copy and paste?

# How frequently did subjects copy and paste?

- average:
about 16 inst/ hr
- median:
about 12 inst/ hr

**How frequently did subjects copy and paste?**

# How long is the code snippet involved in copy operations?