

# Software Evolution



Miryung Kim, Na Meng, and Tianyi Zhang

**Abstract** Software evolution plays an ever-increasing role in software development. Programmers rarely build software from scratch but often spend more time in modifying existing software to provide new features to customers and fix defects in existing software. Evolving software systems are often a time-consuming and error-prone process. This chapter overviews key concepts and principles in the area of software evolution and presents the fundamentals of state-of-the-art methods, tools, and techniques for evolving software. The chapter first classifies the types of software changes into four types: *perfective* changes to expand the existing requirements of a system, *corrective* changes for resolving defects, *adaptive* changes to accommodate any modifications to the environments, and finally *preventive* changes to improve the maintainability of software. For each type of changes, the chapter overviews software evolution techniques from the perspective of three kinds of activities: (1) applying changes, (2) inspecting changes, and (3) validating changes. The chapter concludes with the discussion of open problems and research challenges for the future.

## 1 Introduction

Software evolution plays an ever-increasing role in software development. Programmers rarely build software from scratch but often spend more time in modifying existing software to provide new features to customers and fix defects in existing

---

All authors have contributed equally to this chapter.

M. Kim · T. Zhang  
University of California, Los Angeles, CA, USA  
e-mail: [miryung@cs.ucla.edu](mailto:miryung@cs.ucla.edu); [tianyi.zhang@cs.ucla.edu](mailto:tianyi.zhang@cs.ucla.edu)

N. Meng  
Virginia Tech, Blacksburg, VA, USA  
e-mail: [nm8247@cs.vt.edu](mailto:nm8247@cs.vt.edu)

software. Evolving software systems are often a time-consuming and error-prone process. In fact, it is reported that 90% of the cost of a typical software system is incurred during the maintenance phase [114] and a primary focus in software engineering involves issues relating to upgrading, migrating, and evolving existing software systems.

The term *software evolution* dates back to 1976 when Belady and Lehman first coined this term. Software evolution refers to the *dynamic behavior* of software systems, as they are maintained and enhanced over their lifetimes [13]. Software evolution is particularly important as systems in organizations become longer-lived. A key notion behind this seminal work by Belady and Lehman is the concept of software system *entropy*. The term entropy, with a formal definition in physics relating to the amount of energy in a closed thermodynamic system, is used to broadly represent a measure of the cost required to change a system or correct its natural disorder. As such, this term has had significant appeal to software engineering researchers, since it suggests a set of reasons for software maintenance. Their original work in the 1970s involved studying 20 user-oriented releases of the IBM OS/360 operating systems software, and it was the first empirical research to focus on the dynamic behavior of a relatively large and mature system (12 years old) at the time. Starting with the available data, they attempted to deduce the nature of consecutive releases of OS/360 and to postulate five *laws* of software evolution: (1) continuing change, (2) increasing complexity, (3) fundamental law of program evolution, (4) conservation of organizational stability, and (5) conservation of familiarity.

Later, many researchers have systematically studied software evolution by measuring concrete metrics about software over time. Notably, Eick et al. [41] quantified the symptoms of *code decay*—*software is harder to change than it should be* by measuring the extent to which each risk factor matters using a rich data set of 5ESS telephone switching system. For example, they measured the number of files changed in each modification request to monitor code decay progress over time. This empirical study has influenced a variety of research projects on mining software repositories.

Now that we accept the fact that software systems go through a *continuing life cycle of evolution* after the initial phase of requirement engineering, design, analysis, testing, and validation, we describe an important aspect of software evolution—*software changes*—in this chapter. To that end, we first introduce the categorization of software changes into four types in Sect. 2. We then discuss the techniques of evolving software from the perspectives of three kinds of activities: (1) change application, (2) change inspection, and (3) change validation. In the following three sections, we provide an organized tour of seminal papers focusing on the abovementioned topics.

In Sect. 3, we first discuss empirical studies to summarize the characteristics of each change type and then overview tool support for applying software changes. For example, for the type of *corrective changes*, we present several studies on the nature and extent of bug fixes. We then discuss automated techniques for fixing bugs such as automated repair. Similarly, for the type of *preventative changes*, we present

empirical studies on refactoring practices and then discuss automated techniques for applying refactorings. Regardless of change types, various approaches could reduce the manual effort of updating software through automation, including source-to-source program transformation, programming by demonstration (PbD), simultaneous editing, and systematic editing.

In Sect. 4, we overview research topics for inspecting software changes. Software engineers other than the change author often perform peer reviews by inspecting program changes and provide feedback if they discover any suspicious software modifications. Therefore, we summarize modern code review processes and discuss techniques for comprehending code changes. This section also overviews a variety of program differencing techniques, refactoring reconstruction techniques, and code change search techniques that developers can use to investigate code changes.

In Sect. 5, we overview research techniques for validating software changes. After software modification is made, developers and testers may create new tests or reuse existing tests, run the modified software against the tests, and check whether the software executes as expected. Therefore, the activity of checking the correctness of software changes involves failure-inducing change isolation, regression testing, and change impact analysis.

## 2 Concepts and Principles

Swanson initially identified three categories of software changes: corrective, adaptive, and perfective [176]. These categories were updated later, and ISO/IEC 14764 instead presents four types of changes: corrective, adaptive, perfective, and preventive [70].

### 2.1 Corrective Change

Corrective change refers software modifications initiated by software defects. A defect can result from design errors, logic errors, and coding errors [172].

- Design errors: software design does not fully align with the requirement specification. The faulty design leads to a software system that either incompletely or incorrectly implements the requested computational functionality.
- Logic errors: a program behaves abnormally by terminating unexpectedly or producing wrong outputs. The abnormal behaviors are mainly due to flaws in software functionality implementations.
- Coding errors: although a program can function well, it takes excessively high runtime or memory overhead before responding to user requests. Such failures may be caused by loose coding or the absence of *reasonable checks* on computations performed.

## 2.2 Adaptive Change

Adaptive change is a change introduced to accommodate any modifications in the environment of a software product. The term **environment** here refers to the totality of all conditions that influence the software product, including business rules, government policies, and software and hardware operating systems. For example, when a library or platform developer may evolve its APIs, the corresponding adaptation may be required for client applications to handle such environment change. As another example, when porting a mobile application from Android to iOS, mobile developers need to apply adaptive changes to translate the code from Java to Swift, so that the software is still compilable and executable on the new platform.

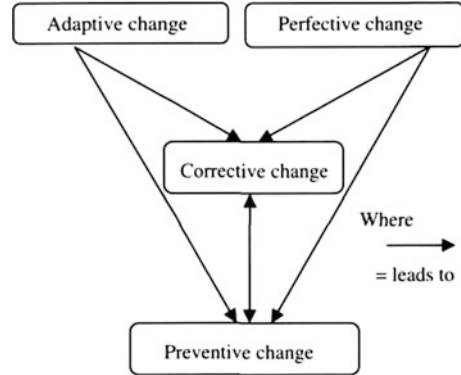
## 2.3 Perfective Change

Perfective change is the change undertaken to expand the existing requirements of a system [55]. When a software product becomes useful, users always expect to use it in new scenarios beyond the scope for which it was initially developed. Such requirement expansion causes changes to either enhance existing system functionality or to add new features. For instance, an image processing system is originally developed to process JPEG files and later goes through a series of perfective changes to handle other formats, such as PNG and SVG. The nature and characteristics of new feature additions are not necessarily easy to define and in fact understudied for that reason. In Sect. 3.3, we discuss a rather well-understood type of perfective changes, called *crosscutting concerns*, and then present tool and language support for adding crosscutting concerns. Crosscutting concerns refer to the *secondary design decisions* such as logging, performance, error handling, and synchronization. Adding these secondary concerns often involves nonlocalized changes throughout the system, due to the *tyranny* of dominant design decisions already implemented in the system. Concerns that are added later may end up being scattered across many modules and thus tangled with one another.

## 2.4 Preventive Change

Preventive change is the change applied to prevent malfunctions or to improve the maintainability of software. According to Lehman's laws of software evolution [108], the long-term effect of corrective, adaptive, and perfective changes is deteriorating the software structure, while increasing entropy. Preventive changes are usually applied to address the problems. For instance, after developers fix some bugs and implement new features in an existing software product, the complexity of

**Fig. 1** Potential relation between software changes [55]. Note: Reprinted from “Software Maintenance: Concepts and Practice” by Penny Grubb and Armstrong A. Takang, Copyright 2003 by World Scientific Publishing Co. Pte. Ltd. Reprinted with permission



source code can increase to an unmanageable level. Through code *refactoring*—a series of behavior-preserving changes—developers can reduce code complexity and increase the readability, reusability, and maintainability of software.

Figure 1 presents the potential relationships between different types of changes [55]. Specifically, both adaptive changes and perfective changes may lead to the other two types of changes, because developers may introduce bugs or worsen code structures when adapting software to new environments or implementing new features.

### 3 An Organized Tour of Seminal Papers: Applying Changes

We discuss the characteristics of *corrective*, *adaptive*, *perfective*, and *preventive* changes using empirical studies and the process and techniques for updating software, respectively, in Sects. 3.1–3.4. Next, regardless of change types, automation could reduce the manual effort of updating software. Therefore, we discuss the topic of automated program transformation and interactive editing techniques for reducing repetitive edits in Sect. 3.5 (Fig. 2).

#### 3.1 Corrective Change

Corrective changes such as bug fixes are frequently applied by developers to eliminate defects in software. There are mainly two lines of research conducted: (1) empirical studies to characterize bugs and corresponding fixes and (2) automatic approaches to detect and fix such bugs. There is no clear boundary between the two lines of research, because some prior projects first make observations about particular kinds of bug fixes empirically and then subsequently leverage their observed characteristics to find more bugs and fix them. Below, we discuss a few

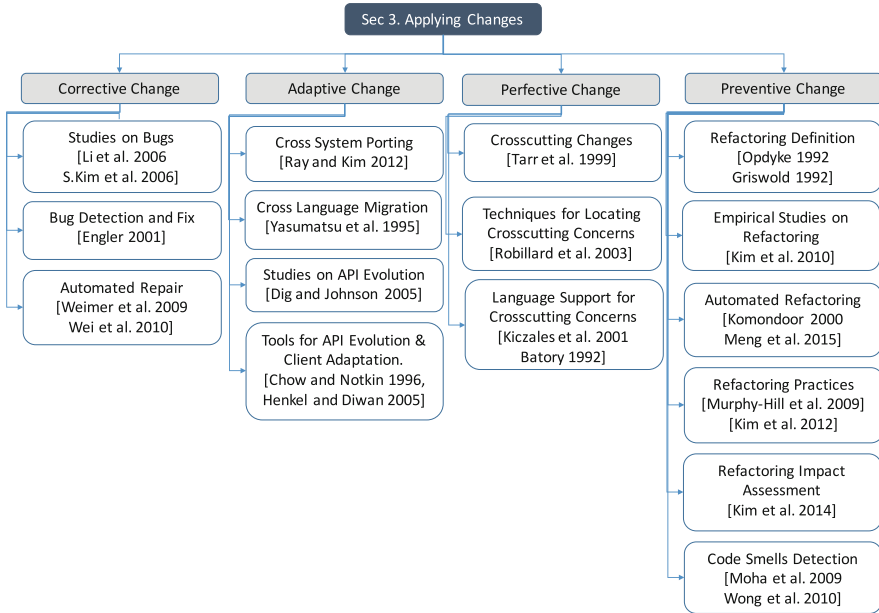


Fig. 2 Applying changes categorized by change type and related research topics

representative examples of empirical studies with such flavor of characterizing and fixing bugs.

### 3.1.1 Empirical Studies of Bug Fixes

In this section, we discuss two representative studies on bug fixes. These studies are not the earliest, seminal works in this domain. Rather, the flavor and style of their studies are representative. Li et al. conducted a large-scale characterization of bugs by digging through bug reports in the wild and by quantifying the extent of each bug type [111]. Kim et al.'s *memory of bug fixes* [90] uses fine-grained bug fix histories to measure the extent of recurring, similar bug fixes and to assess the potential benefit of automating similar fixes based on change history.

Li et al. conducted an empirical study of bugs from two popular open-source projects: Mozilla and Apache HTTP Server [111]. By manually examining 264 bug reports from the Mozilla Bugzilla database, and 209 bug reports from the Apache Bugzilla database, they investigated the root cause, impact, and software components of each software error that exhibited abnormal runtime behaviors. They observed three major root causes: *memory*, *concurrency*, and *semantics*. The memory bugs accounted for 16.3% in Mozilla and 12.2% in Apache. Among memory bugs, NULL pointer dereference was observed as a major cause, accounting for 37.2% in Mozilla and 41.7% in Apache. More importantly, semantic bugs

were observed to be dominant, accounting for 81.1% in Mozilla and 86.7% in Apache. One possible reason is that most semantic bugs are specific to applications. A developer could easily introduce semantic bugs while coding, due to a lack of thorough understanding of software and its requirements. It is challenging to automatically detect or fix such semantic bugs, because diagnosing and resolving them may require a lot of domain-specific knowledge and such knowledge is inherently not generalizable across different systems and applications.

To understand the characteristics and frequency of project-specific bug fixes, Kim et al. conducted an empirical study on the bug fix history of five open-source projects: ArgoUML, Columba, Eclipse, jEdit, and Scarab [90]. With keywords like “Fixed” or “Bugs,” they retrieved code commits in software version history that are relevant to bug fixes, chopped each commit into contiguous code change blocks (i.e., hunks), and then clustered similar code changes. They observed that 19.3–40.3% bugs appeared repeatedly in version history, while 7.9–15.5% of bug-and-fix pairs appeared more than once. The results demonstrated that project-specific bug fix patterns occur frequently enough, and for each bug-and-fix pair, it is possible to both detect similar bugs and provide fix suggestions. Their study also showed history-based bug detection could be complementary to static analysis-based bug detection—the bugs that can be detected by past bug fix histories do not overlap with the bugs that can be detected by a static bug finding tool, PMD [146].

### 3.1.2 Rule-Based Bug Detection and Fixing Approaches

Rule-based bug detection approaches detect and fix bugs based on the assumption that bugs are *deviant program behaviors* that violate implicit programming rules. Then one may ask, where are those implicit rules coming from? Such rules can be written by the developers of bug-finding tools or can be refined based on empirical observation in the wild. For example, Engler et al. define a meta-language for users to easily specify temporal system rules such as “release locks after acquiring them” [44]. They also extend a compiler to interpret the rules and dynamically generate additional checks in the compiler. If any code snippet violates the specified rule(s), the approach reports the snippet as a software bug. Table 1 presents some exemplar system rule templates and instances. With this approach, developers can flexibly define their own rules to avoid some project-specific bugs, without worrying about how to implement checkers to enforce the rules. Engler et al.’s later work enables tool developers to tailor rule templates to a specific system and to check for contradictions and violations [45].

**Table 1** Sample system rule templates and examples from [44]

Rule template	Example
“Never/always do X”	“Do not use floating point in the kernel”
“Do X rather than Y”	“Use memory mapped I/O rather than copying”
“Always do X before/after Y”	“Check user pointers before using them in the kernel”

Another example of rule-based bug detection is CP-Miner, an automatic approach to find copy-paste bugs in large-scale software [110]. CP-Miner is motivated by Chou et al.'s finding that, under the Linux `drivers/i2o` directory, 34 out of 35 errors were caused by copy-paste [24], and based on the insight that when developers copy and paste, they may forget to consistently rename identifiers. CP-Miner first identifies copy-paste code in a scalable way and then detects bugs by checking for a specific rule, e.g., consistent renaming of identifiers.

### 3.1.3 Automated Repair

Automatic program repair generates candidate patches and checks correctness using compilation, testing, and/or specification.

One set of techniques uses *search-based repair* [59] or predefined repair templates to generate many candidate repairs for a bug and then validates them using indicative workloads or test suites. For example, GenProg generates candidate patches by replicating, mutating, or deleting code *randomly* from the existing program [107, 198]. GenProg uses genetic programming (GP) to search for a program variant that retains required functionality but is not vulnerable to the defect in question. GP is a stochastic search method inspired by biological evolution that discovers computer programs tailored to a particular task. GP uses computational analogs of biological mutation and crossover to generate new program variations, in other words program variants. A user-defined fitness function evaluates each variant. GenProg uses the input test cases to evaluate the fitness, and individuals with high fitness are selected for continued evolution. This GP process is successful, when it produces a variant that passes all tests encoding the required behavior and does not fail those encoding the bug.

Another class of strategies in automatic software repair relies on *specifications* or *contracts* to guide sound patch generation. This provides confidence that the output is correct. For example, AutoFix-E generates simple bug fixes from manually prescribed contracts [195]. The key insights behind this approach are to rely on contracts present in the software to ensure that the proposed fixes are semantically sound. AutoFix-E takes an Eiffel class and generates test cases with some automated testing engine first. From the test runs, it extracts object states using Boolean queries. By comparing the states of passing and failing runs, it then generates a fault profile—an indication of what went wrong in terms of an abstract object state. From the state transitions in passing runs, it generates a finite-state behavioral model, capturing the normal behavior in terms of control. Both control and state guide the generation of fix candidates, and only those fixes passing the regression test suite remain.

Some approaches are specialized for particular types of bugs only. For example, FixMeUp inserts missing security checks using inter-procedural analysis, but these additions are very specific and stylized for access-control-related security bugs [173]. As another example, PAR [95] encodes ten common bug fix patterns from Eclipse JDT's version history to improve GenProg. However, the patterns are created manually.



## 3.2 Adaptive Change

Adaptive changes are applied to software, when its environment changes. In this section, we focus on three scenarios of adaptive changes: cross-system software porting, cross-language software migration, and software library upgrade (i.e., API evolution).

Consider an example of cross-system porting. When a software system is installed on a computer, the installation can depend on the configurations of the hardware, the software, and the device drivers for particular devices. To make the software to run on a different processor or an operating system, and to make it compatible with different drivers, we may need adaptive changes to adjust the software to the new environment. Consider another example of cross-language migration where you have software in Java that must be translated to C. Developers need to rewrite software and must also update language-specific libraries. Finally consider the example of API evolution. When the APIs of a library and a platform evolve, corresponding adaptations are often required for client applications to handle such API update. In extreme cases, e.g., when porting a Java desktop application to the iOS platform, developers need to rewrite everything from scratch, because both the programming language (i.e., Swift) and software libraries are different.

### 3.2.1 Cross-System Porting

Software forking—creating a variant product by copying and modifying an existing product—is often considered an ad hoc, low-cost alternative to principled product line development. To maintain such forked products, developers often need to port an existing feature or bug-fix from one product variant to another.

*Empirical Studies on Cross-System Porting* OpenBSD, NetBSD, and FreeBSD have evolved from the same origin but have been maintained independently from one another. Many have studied the BSD family to investigate the extent and nature of cross-system porting. The studies found that (1) the information flow among the forked BSD family is decreasing according to change commit messages [47]; (2) 40% of lines of code were shared among the BSD family [205]; (3) in some modules such as device driver modules, there is a significant amount of adopted code [27]; and (4) contributors who port changes from other projects are highly active contributors according to textual analysis of change commit logs and mailing list communication logs [21].

More recently, Ray et al. comprehensively characterized the temporal, spatial, and developer dimensions of cross-system porting in the BSD family [152]. Their work computed the amount of edits that are ported from other projects as opposed to the amount of code duplication across projects, because not all code clones across different projects undergo similar changes during evolution, and similar changes are not confined to code clones. To identify ported edits, they first built a tool

named as Repertoire that takes *diff* patches as input and compares the content and edit operations of the program patches. Repertoire was applied to total 18 years of NetBSD, OpenBSD, and FreeBSD version history. Their study found that maintaining forked projects involves significant effort of porting patches from other projects—10–15% of patch content was ported from another project's patches. Cross-system porting is periodic, and its rate does not necessarily decrease over time. A significant portion of active developers participate in porting changes from peer projects. Ported changes are less defect-prone than non-porting changes. A significant portion (26–59%) of active committers port changes, but some do more porting work than others. While most ported changes migrate to peer projects in a relatively short amount of time, some changes take a very long time to propagate to other projects. Ported changes are localized within less than 20% of the modified files per release on average in all three BSD projects, indicating that porting is concentrated on a few subsystems.

### 3.2.2 Cross-Language Migration

When maintaining a legacy system that was written in an old programming language (e.g., Fortran) decades ago, programmers may migrate the system to a mainstream general-purpose language, such as Java, to facilitate the maintenance of existing codebase and to leverage new programming language features.

*Cross-Language Program Translation* To translate code implementation from one language to another, researchers have built tools by hard coding the translation rules and implementing any missing functionality between languages. Yasumatsu et al. map compiled methods and contexts in Smalltalk to machine code and stack frames, respectively, and implement runtime replacement classes in correspondence with the Smalltalk execution model and runtime system [208]. Mossienko [127] and Sneed [170] automate COBOL-to-Java code migration by defining and implementing rules to generate Java classes, methods, and packages from COBOL programs. *mppSMT* automatically infers and applies Java-to-C# migration rules using a phrase-based statistical machine translation approach [136]. It encodes both Java and C# source files into sequences of syntactic symbols, called *syntaxemes*, and then relies on the *syntaxemes* to align code and to train sequence-to-sequence translation.

*Mining Cross-Language API Rules* When migrating software to a different target language, API conversion poses a challenge for developers, because the diverse usage of API libraries induces an endless process of specifying API translation rules or identifying API mappings across different languages. Zhong et al. [215] and Nguyen et al. [135, 137] automatically mine API usage mappings between Java and C#. Zhong et al. align code based on similar names and then construct the API transformation graphs for each pair of aligned statements [215]. StaMiner [135] mines API usage sequence mappings by conducting program dependency analysis [128] and representing API usage as a graph-based model [133].

### 3.2.3 Library Upgrade and API Evolution

Instead of building software from scratch, developers often use existing frameworks or third-party libraries to reuse well-implemented and tested functionality. Ideally, the APIs of libraries must remain stable such that library upgrades do not incur corresponding changes in client applications. In reality, however, APIs change their input and output signatures, change semantics, or are even deprecated, forcing client application developers to make corresponding adaptive changes in their applications.

*Empirical Studies of API Evolution* Dig and Johnson manually investigated API changes using the change logs and release notes to study the types of library-side updates that break compatibility with existing client code and discovered that 80% of such changes are refactorings [36]. Xing and Stroulia used UMLDiff to study API evolution and found that about 70% of structural changes are refactorings [203]. Yokomori et al. investigated the impact of library evolution on client code applications using component ranking measurements [210]. Padioleau et al. found that API changes in the Linux kernel led to subsequent changes on dependent drivers, and such collateral evolution could introduce bugs into previously mature code [143]. McDonelle et al. examined the relationship between API stability and the degree of adoption measured in propagation and lagging time in the Android Ecosystem [117]. Hou and Yao studied the Java API documentation and found that a stable architecture played an important role in supporting the smooth evolution of the AWT/Swing API [68]. In a large-scale study of the Smalltalk development communities, Robbes et al. found that only 14% of deprecated methods produce nontrivial API change effects in at least one client-side project; however, these effects vary greatly in magnitude. On average, a single API deprecation resulted in 5 broken projects, while the largest caused 79 projects and 132 packages to break [158].

*Tool Support for API Evolution and Client Adaptation* Several existing approaches semiautomate or automate client adaptations to cope with evolving libraries. Chow and Notkin [25] propose a method for changing client applications in response to library changes—a library maintainer annotates changed functions with rules that are used to generate tools that update client applications. Henkel and Diwan’s CatchUp records and stores refactorings in an XML file that can be replayed to update client code [62]. However, its update support is limited to three refactorings: renaming operations (e.g., types, methods, fields), moving operations (e.g., classes to different packages, static members), or change operations (e.g., types, signatures). The key idea of CatchUp, *record and replay*, assumes that the adaptation changes in client code are exact or similar to the changes in the library side. Thus, it works well for replaying rename or move refactorings or supporting API usage adaptations via inheritance. However, CatchUp cannot suggest programmers how to manipulate the context of API usages in client code such as the surrounding control structure or the ordering between method calls. Furthermore, CatchUp requires that library and client application developers use the same development environment to

record API-level refactorings, limiting its adoption in practice. Xing and Stroulia's Diff-CatchU automatically recognizes API changes of the reused framework and suggests plausible replacements to the obsolete APIs based on the working examples of the framework codebase [204]. Dig et al.'s MolhadoRef uses recorded API-level refactorings to resolve merge conflicts that stem from refactorings; this technique can be used for adapting client applications in case of simple rename and move refactorings occurred in a library [37].

SemDiff [32] mines API usage changes from other client applications or the library itself. It defines an adaptation pattern as a frequent *replacement* of a method invocation. That is, if a method call to  $A.m$  is changed to  $B.n$  in several adaptations,  $B.n$  is likely to be a correct replacement for the calls to  $A.m$ . As SemDiff models API usages in terms of method calls, it cannot support complex adaptations involving multiple objects and method calls that require the knowledge of the surrounding context of those calls. LibSync helps client applications migrate library API usages by learning migration patterns [134] with respect to a partial AST with containment and data dependences. Though it suggests what code locations to examine and shows example API updates, it is unable to transform code automatically. Cossette and Walker found that, while most broken code may be mended using one or more of these techniques, each is ineffective when used in isolation [29].

### 3.3 *Perfective Change*

Perfective change is the change undertaken to expand the existing requirements of a system. Not much research is done to characterize feature enhancement or addition. One possible reason is that the implementation logic is always domain and project-specific and that it is challenging for any automatic tool to predict what new feature to add and how that new feature must be implemented. Therefore, the nature and characteristics of feature additions are understudied.

In this section, we discuss a rather well-understood type of perfective changes, called *crosscutting concerns* and techniques for implementing and managing crosscutting concerns. As programs evolve over time, they may suffer from the *tyranny of dominant decomposition* [180]. They can be modularized in only one way at a time. Concerns that are added later may end up being scattered across many modules and tangled with one another. Logging, performance, error handling, and synchronization are canonical examples of such secondary design decisions that lead to nonlocalized changes.

Aspect-oriented programming languages provide language constructs to allow concerns to be updated in a modular fashion [86]. Other approaches instead leave the crosscutting concerns in a program, while providing mechanisms to document and manage related but dispersed code fragments. For example, Griswold's information transparency technique uses naming conventions, formatting styles, and ordering of code in a file to provide indications about crosscutting concern code that should change together [53].

### 3.3.1 Techniques for Locating Crosscutting Concerns

Several tools allow programmers to automatically or semiautomatically locate crosscutting concerns. Robillard et al. allow programmers to manually document crosscutting concerns using structural dependencies in code [160]. Similarly, the Concern Manipulation Environment allows programmers to locate and document different types of concerns [60]. van Engelen et al. use clone detectors to locate crosscutting concerns [192]. Shepherd et al. locate concerns using natural language program analysis [166]. Breu et al. mine aspects from version history by grouping method calls that are added together [18]. Dagenais et al. automatically infer and represent structural patterns among the participants of the same concern as rules in order to trace the concerns over program versions [33].

### 3.3.2 Language Support for Crosscutting Concerns

*Aspect-oriented programming* (AOP) is a programming paradigm that aims to increase modularity by allowing the separation of crosscutting concerns [181]. Suppose developers want to add a new feature such as logging to log all executed functions. The logging logic is straightforward: printing the function's name at each function's entry. However, manually inserting the same implementation to each function body is tedious and error-prone. With AOP, developers only need to first define the logging logic as **an advice** and then specify the place where to insert the advice (i.e., **pointcut**), such as the entry point of each function. An aspect weaver will read the aspect-oriented code and generate appropriate object-oriented code with the aspects integrated. In this way, AOP facilitates developers to efficiently introduce new program behaviors without cluttering the core implementation in the existing codebase. Many Java bytecode manipulation frameworks implement the AOP paradigm, like ASM [6], Javassist [75], and AspectJ [181], so that developers can easily modify program runtime behaviors without touching source code. The benefit of AOP during software evolution is that crosscutting concerns can be contained as a separate module such as an `aspect` with its `pointcut` and `advice` description and thus reduces the developer effort in locating and updating all code fragments relevant to a particular secondary design decision such as logging, synchronization, database transaction, etc.

*Feature-oriented programming* (FOP) is another paradigm for program generation in software product lines and for incremental development of programs [12]. FOP is closely related to AOP. Both deal with modules that encapsulate crosscuts of classes, and both express program extensions. In FOP, every software is considered as a composition of multiple features or layers. Each feature implements a certain program functionality, while features may interact with each other to collaboratively provide a larger functionality. A software product line (SPL) is a family of programs where each program is defined by a unique composition of features. Formally, FOP considers programs as *values* and program extensions as *functions* [103]. The benefit of FOP is similar to AOP in that secondary design decisions can be encapsulated as

a separate feature and can be composed later with other features using program synthesis, making it easier to add a new feature at a later time during software evolution. Further discussion of program generation techniques for software product lines is described in chapter “Software Reuse and Product Line Engineering.”

### 3.4 Preventive Change

As a software system is enhanced, modified, and adapted to new requirements, the code becomes more complex and drifts away from its original design, thereby lowering the quality of the software. *Refactoring* [159, 52, 140, 122] copes with increasing software complexity by transforming a program from one representation to another while preserving the program’s external behavior (functionality and semantics). Mens et al. present a survey of refactoring research and describe a refactoring process, consisting of the following activities [122]:

1. Identifying where to apply what refactoring(s).
2. Checking that the refactoring to apply preserves program behaviors.
3. Refactoring the code.
4. Assessing the effect of applied refactoring on software quality (e.g., complexity and readability).
5. Maintaining the consistency between refactored code and other related software artifacts, like documentation, tests, and issue tracking records.

Section 3.4.1 describes the definition of refactoring and example transformations. Section 3.4.2 describes empirical studies on refactoring. Section 3.4.3 describes tool support for automated refactoring. Section 3.4.4 describes several studies of modern refactoring practices and the limitations of current refactoring support. Section 3.4.5 describes techniques for assessing the impact of refactoring. Section 3.4.6 describes techniques for identifying opportunities for refactoring.

#### 3.4.1 Definition of Refactoring Operations

Griswold’s dissertation [52] discusses one of the first refactoring operations that automate repetitive, error-prone, nonlocal transformations. Griswold supports a number of restructuring operations: replacing an expression with a variable that has its value, swapping the formal parameters in a procedure’s interface and the respective arguments in its calls, etc. It is important to note that many of these refactoring operations are systematic in the sense that they involve repetitive nonlocal transformations.

Opdyke’s dissertation [140] distinguishes the notion of low-level refactorings from high-level refactorings. High-level refactorings (i.e., composite refactorings) reflect more complex behavior-preserving transformations while low-level refactorings are primitive operations such as creating, deleting, or changing a program entity

or moving a member variable. Opdyke describes three kinds of complex refactorings in detail: (1) creating an abstract superclass, (2) subclassing and simplifying conditionals, and (3) capturing aggregations and components. All three refactorings are systematic in the sense that they contain multiple similar transformations at a code level. For example, creating an abstract superclass involves moving multiple variables and functions common to more than one sibling classes to their common superclass. Subclassing and simplifying conditionals consist of creating several classes, each of which is in charge of evaluating a different conditional. Capturing aggregations and components usually involves moving multiple members from a component to an aggregate object.

While refactoring is defined as behavior-preserving code transformations in the academic literature [122], the de facto definition of refactoring in practice seems to be very different from such rigorous definition. Fowler catalogs 72 types of structural changes in object-oriented programs, but these transformations do not necessarily guarantee behavior preservation [159]. In fact, Fowler recommends developers to write test code first, since these refactorings may change a program's behavior. Murphy-Hill et al. analyzed refactoring logs and found that developers often interleave refactorings with other behavior-modifying transformations [130], indicating that pure refactoring revisions are rare. Johnson's refactoring definition is aligned with these findings—*refactoring improves behavior in some aspects but does not necessarily preserve behavior in all aspects* [79].

### 3.4.2 Empirical Studies of Refactoring

There are contradicting beliefs on refactoring benefits. On one hand, some believe that refactoring improves software quality and maintainability and a lack of refactoring incurs *technical debt* to be repaid in the future in terms of increased maintenance cost [19]. On the other hand, some believe that refactoring does not provide immediate benefits unlike bug fixes and new features during software evolution.

Supporting the view that refactoring provides benefits during software evolution, researchers found empirical evidence that bug fix time decreases after refactoring and defect density decreases after refactoring. More specifically, Carriere et al. found that the productivity measure manifested by the average time taken to resolve tickets decreases after re-architecting the system [22]. Ratzinger et al. developed defect prediction models based on software evolution attributes and found that refactoring-related features and defects have an inverse correlation [151]—if the number of refactorings increases in the preceding time period, the number of defects decreases.

Supporting the opposite view that refactoring may even incur additional bugs, researchers found that code churns are correlated with defect density and that refactorings are correlated with bugs. More specifically, Purushothaman and Perry found that nearly 10% of changes involved only a single line of code, which has less than a 4% chance to result in error, while a change of 500 lines or more has nearly

a 50% chance of causing at least one defect [148]. This result may indicate that large commits, which tend to include refactorings, have a higher chance of inducing bugs. Weißgerber and Diehl found that refactorings often occur together with other types of changes and that refactorings are followed by an increasing number of bugs [196]. Kim et al. investigated the spatial and temporal relationship between API refactorings and bug fixes using a K-revision sliding window and by reasoning about the method-level location of refactorings and bug fixes. They found that the number of bug fixes increases after API refactorings [93].

One reason why refactoring could be potentially error-prone is that refactoring often requires coordinated edits across different parts of a system, which could be difficult for programmers to locate all relevant locations and apply coordinated edits consistently. Several researchers found such evidence from open-source project histories—Kim et al. found the exceptions to systematic change patterns, which often arise from the failure to complete coordinated refactorings [91, 87], cause bugs. Görg and Weißgerber detect errors caused by incomplete refactorings by relating API-level refactorings to the corresponding class hierarchy [51]. Nagappan and Ball found that code churn—the number of added, deleted, and modified lines of code—is correlated with defect density [131]; since refactoring often introduces a large amount of structural changes to the system, some question the benefit of refactoring.

### 3.4.3 Automated Refactoring

The Eclipse IDE provides automatic support for a variety of refactorings, including *rename*, *move*, and *extract method*. With such support, developers do not need to worry about how to check for preconditions or postconditions before manually applying a certain refactoring. Instead, they can simply select the refactoring command from a menu (e.g., *extract method*) and provide necessary information to accomplish the refactoring (e.g., *the name of a new method*). The Eclipse refactoring engine takes care of the precondition check, program transformation, and postcondition check.

During refactoring automation, Opdyke suggests to ensure behavior preservation by specifying *refactoring preconditions* [140]. For instance, when conducting a *create\_method\_function* refactoring, before inserting a member function  $F$  to a class  $C$ , developers should specify and check for five preconditions: (1) the function is not already defined locally. (2) The signature matches that of any inherited function with the same name. (3) The signature of corresponding functions in subclasses matches it. (4) If there is an inherited function with the same name, either the inherited function is not referenced on instances of  $C$  and its subclasses, or the new function is semantically equivalent to the function it replaces. (5)  $F$  will compile as a member of  $C$ . If any precondition is not satisfied, the refactoring should not be applied to the program. These five conditions in Opdyke's dissertation are represented using first-order logic.

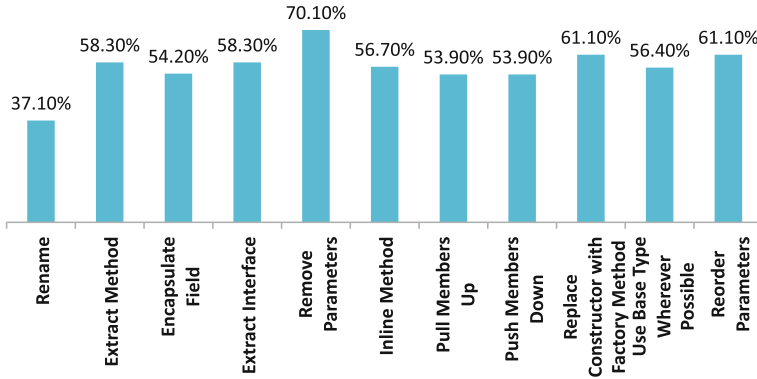


Clone removal refactorings factorize the common parts of similar code by parameterizing their differences using a *strategy* design pattern or a *form template method* refactoring [8, 178, 82, 67, 101]. These tools insert customized calls in each original location to use newly created methods. Juillerat et al. automate *introduce exit label* and *introduce return object* refactorings [82]. However, for variable and expression variations, they define extra methods to mask the differences [8]. Hotta et al. use program dependence analysis to handle gapped clones—trivial differences inside code clones that are safe to factor out such that they can apply the *form template method* refactoring to the code [67]. Krishnan et al. use PDGs of two programs to identify a maximum common subgraph so that the differences between the two programs are minimized and fewer parameters are introduced [101]. RASE is an advanced clone removal refactoring technique that (1) extracts common code; (2) creates new types and methods as needed; (3) parameterizes differences in types, methods, variables, and expressions; and (4) inserts return objects and exit labels based on control and data flow by combining multiple kinds of clone removal transformations [120]. Such clone removal refactoring could lead to an increase in the total size of code because it creates numerous simple methods.

Komondoor et al. extract methods based on the user-selected or tool-selected statements in one method [98, 99]. The *extract method* refactoring in the Eclipse IDE requires contiguous statements, whereas their approach handles noncontiguous statements. Program dependence analysis identifies the relation between selected and unselected statements and determines whether the noncontiguous code can be moved together to form extractable contiguous code. Komondoor et al. apply *introduce exit label* refactoring to handle exiting jumps in selected statements [99]. Tsantalís et al. extend the techniques by requiring developers to specify a variable of interest at a specific point only [188]. They use a block-based slicing technique to suggest a program slice to isolate the computation of the given variable. These automated procedure extraction approaches are focused on extracting code from a single method only. Therefore, they do not handle extracting common code from multiple methods and resolving the differences between them.

### 3.4.4 Real-World Refactoring Practices

Several studies investigated refactoring practices in industry and also examined the current challenges and risks associated with refactoring. Kim et al. conducted a survey with professional developers at Microsoft [94, 96]. They sent a survey invitation to 1290 engineers whose commit messages include a keyword “refactoring” in the version histories of five MS products. Three hundred and twenty-eight of them responded to the survey. More than half of the participants said they carry out refactorings in the context of bug fixes or feature additions, and these changes are generally not semantics-preserving. When they asked about their own definition of refactoring, 46% of participants did not mention preservation of semantics, behavior, or functionality at all. 53% reported that refactorings that they perform do not match the types and capability of transformations supported by existing refactoring engines.



**Fig. 3** The percentage of survey participants who know individual refactoring types but do those refactorings manually [96]

In the same study, when developers are asked “what percentage of your refactoring is done manually as opposed to using automated refactoring tools?”, developers answered they do 86% of refactoring manually on average. Figure 3 shows the percentages of developers who usually apply individual refactoring types manually despite the awareness of automated refactoring tool support. Vakilian et al. [191] and Murphy et al. [129] also find that programmers do not use automated refactoring despite their awareness of the availability of automated refactorings. Murphy-Hill manually inspected source code produced by 12 developers and found that developers only used refactoring tools for 10% of refactorings for which tools were available [130]. For the question, “based on your experience, what are the risks involved in refactorings?”, developers reported regression bugs, code churn, merge conflicts, time taken from other tasks, the difficulty of doing code reviews after refactoring, and the risk of overengineering. 77% think that refactoring comes with a risk of introducing subtle bugs and functionality regression [94].

In a separate study of refactoring tool use, Murphy-Hill et al. gave developers specific examples of when they did not use refactoring tools but could have [130] and asked why. One reason was that developers started a refactoring manually but only partway through realized that the change was a refactoring that the IDE offered—by then, it was too late. Another complaint was that refactoring tools disrupted their workflow, forcing them to use a tool when they wanted to focus on code.

### 3.4.5 Quantitative Assessment of Refactoring Impact

While several prior research efforts have conceptually advanced the benefit of refactoring through metaphors, few empirical studies assessed refactoring impact quantitatively. Sullivan et al. first linked software modularity with option the-

ories [175]. A module provides an option to substitute it with a better one without symmetric obligations, and investing in refactoring activities can be seen as purchasing *options* for future adaptability, which will produce benefits when changes happen and the module can be replaced easily. Baldwin and Clark argued that the modularization of a system can generate tremendous value in an industry, given that this strategy creates valuable options for module improvement [10]. Ward Cunningham drew the comparison between debt and a lack of refactoring: a quick and dirty implementation leaves *technical debt* that incur *penalties* in terms of increased maintenance costs [31]. While these projects advanced conceptual understanding of refactoring impact, they did not quantify the benefits of refactoring.

Kim et al. studied how refactoring impacts inter-module dependencies and defects using the quantitative analysis of Windows 7 version history [96]. Their study finds the top 5% of preferentially refactored modules experience higher reduction in the number of inter-module dependencies and several complexity measures but increase size more than the bottom 95%. Based on the hypothesis that measuring the impact of refactoring requires multidimensional assessment, they investigated the impact of refactoring on various metrics: churn, complexity, organization and people, cohesiveness of ownership, test coverage, and defects.

MacCormack et al. defined modularity metrics and used these metrics to study evolution of Mozilla and Linux. They found that the redesign of Mozilla resulted in an architecture that was significantly more modular than that of its predecessor. Their study monitored design structure changes in terms of modularity metrics without identifying the modules where refactoring changes are made [113]. Kataoka et al. proposed a refactoring evaluation method that compares software before and after refactoring in terms of coupling metrics [84]. Kolb et al. performed a case study on the design and implementation of existing software and found that refactoring improves software with respect to maintainability and reusability [97]. Moser et al. conducted a case study in an industrial, agile environment and found that refactoring enhances quality- and reusability-related metrics [126]. Tahvildari et al. suggested using a catalogue of object-oriented metrics to estimate refactoring impact, including complexity metrics, coupling metrics, and cohesion metrics [177].

### 3.4.6 Code Smells Detection

Fowler describes the concept of *bad smell* as a heuristic for identifying redesign and refactoring opportunities [159]. Example of bad smells include code clone and feature envy. Several techniques automatically identify bad smells that indicate needs of refactorings [186, 187, 190].

Garcia et al. propose several architecture-level bad smells [49]. Moha et al. present the Decor tool and domain specific language (DSL) to automate the construction of design defect detection algorithms [125].

Tsantalis and Chatzigeorgiou's technique identifies *extract method* refactoring opportunities using static slicing [186]. Detection of some specific bad smells such as code duplication has also been extensively researched. Higo et al. propose

the Aries tool to identify possible refactoring candidates based on the number of assigned variables, the number of referred variables, and dispersion in the class hierarchy [64]. A refactoring can be suggested if the metrics for the clones satisfy certain predefined values. Koni-N’Sapu provides refactoring suggestions based on the location of clones with respect to a class hierarchy [100]. Balazinska et al. suggest clone refactoring opportunities based on the differences between the cloned methods and the context of attributes, methods, and classes containing clones [9]. Kataoka et al. use Daikon to infer program invariants at runtime and suggest candidate refactorings using inferred invariants [83]. If Daikon observes that one parameter of a method is always constant, it then suggests a *remove parameter* refactoring. *Breakaway* automatically identifies detailed structural correspondences between two abstract syntax trees to help programmers generalize two pieces of similar code [30].

Gueheneuc et al. detect inter-class design defects [56], and Marinescu identifies design flaws using software metrics [116]. Izurieta and Bieman detect accumulation of non-design-pattern-related code [71]. Guo et al. define domain-specific code smells [57] and investigate the consequence of technical debt [58]. Tsantalis et al. rank clones that have been repetitively or simultaneously changed in the past to suggest refactorings [189]. Wang et al. extract features from code to reflect program context, code smell, and evolution history and then use a machine learning technique to rank clones for refactorings [194].

Among the above tools, we briefly present a few concrete examples of four design smells from Decor [125]. In XERCES, method `handleIncludeElement(XMLAttributes)` of the `org.apache.xerces.xinclude.XIncludeHandler` class is a typical example of *Spaghetti Code*—classes without structure that declare long methods without parameters. A good example of *Blob* (a large controller class that depends on data stored in surrounding data classes) is class `com.aelitis.azureus.core.dht.control.impl.DHTControlImpl` in AZUREUS. This class declares 54 fields and 80 methods for 2965 lines of code. Functional decomposition may occur if developers with little knowledge of object orientation implement an object-oriented system. An interesting example of *Functional Decomposition* is class `org.argouml.uml.cognitive.critics.Init` in ARGOXML, in particular because the name of the class includes a suspicious term, *init*, that suggests a functional programming. The *Swiss Army Knife* code smell is a complex class that offers a high number of services (i.e., interfaces). Class `org.apache.xerces.impl.dtd.DTDGrammar` is a striking example of Swiss Army Knife in XERCES, implementing 4 different sets of services with 71 fields and 93 methods for 1146 lines of code.

Clio detects modularity violations based on the assumptions that multiple types of bad smells are instances of modularity violations that can be uniformly detected by reasoning about modularity hierarchy in conjunction with change locations [200]. They define *modularity violations* as recurring discrepancies between which modules should change together and which modules actually change together according to version histories. For example, when code clones change frequently

together, Clio will detect this problem because the co-change pattern deviates from the designed modular structure. Second, by taking version histories as input, Clio detects violations that happened most recently and frequently, instead of bad smells detected in a single version without regard to the program's evolution context. Ratzinger et al. also detect bad smells by examining change couplings, but their approach leaves it to developers to identify design violations from visualization of change coupling [150].

### 3.5 Automatic Change Application

Regardless of change types, various approaches are proposed to automatically suggest program changes or reduce the manual effort of updating software. In this section, we discuss automated change application techniques including source-to-source program transformation, programming by demonstration (PbD), simultaneous editing, and systematic editing (Fig. 4).

#### 3.5.1 Source Transformation and Languages and Tools

Source transformation tools allow programmers to author their change intent in a formal syntax and automatically update a program using the change script. Most source transformation tools automate repetitive and error-prone program updates. The most ubiquitous and the least sophisticated approach to program transformation is text substitution. More sophisticated systems use program structure information. For example, A\* [102] and TAWK [54] expose syntax trees and primitive data structures. Stratego/XT is based on algebraic data types and term pattern matching [193]. These tools are difficult to use as they require programmers to understand low-level program representations. TXL attempts to hide these low-level details by using an extended syntax of the underlying programming language [26]. Boshernitsan et al.'s iXJ enables programmers to perform systematic code transformations easily by providing a visual language and a tool for describing and prototyping source transformations. Their user study shows that iXJ's visual language is aligned

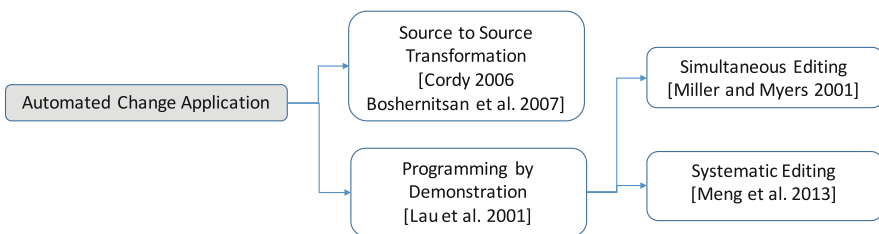


Fig. 4 Automated change application and related research topics

with programmers' mental model of code-changing tasks [16]. Coccinelle [144] allows programmers to safely apply crosscutting updates to Linux device drivers. We describe two seminal approaches with more details.

**Example: TXL** TXL is a programming language and rapid prototyping system specifically designed to support structural source transformation. TXL's source transformation paradigm consists of parsing the input text into a structure tree, transforming the tree to create a new structure tree, and unparsing the new tree to a new output text. Source text structures to be transformed are described using an unrestricted ambiguous context-free grammar in extended Backus-Nauer (BNF) form. Source transformations are described by example, using a set of context-sensitive structural transformation rules from which an application strategy is automatically inferred.

Each transformation rule specifies a *target type* to be transformed, a *pattern* (an example of the particular instance of the type that we are interested in replacing), and a *replacement* (an example of the result we want when we find such an instance). In particular, the pattern is an actual source text example expressed in terms of tokens (terminal symbols) and variables (nonterminal types). When the pattern is matched, variable names are bound to the corresponding instances of their types in the match. Transformation rules can be composed like function compositions.

TXL programs normally consist of three parts, a context-free "base" grammar for the language to be manipulated, a set of context-free grammatical "overrides" (extensions or changes) to the base grammar, and a rooted set of source transformation rules to implement transformation of the extensions to the base language, as shown in Fig. 5. This TXL program overrides the grammar of statements to allow a new statement form. The transformation rule `main` transforms the new form of a statement `V+=E` to an old statement `V:= V+(E)`. In other words, if there are two statements `foo+=bar` and `baz+=boo`, they will be transformed to `foo:=foo+(bar)` and `baz:=baz+(boo)` at the source code level.

**Fig. 5** A simple exemplar TXL file based on [182]

```
% Trivial coalesced addition dialect of Pascal
% Based on standard Pascal grammar
include "Pascal.Grm"

% Overrides to allow new statement forms
redefine statement
    ...
    | [reference] += [expression]
end redefine

% Transform new forms to old
rule main
    replace [statement]
        V [reference] += E [expression]
    by
        V := V + (E)
end rule
```

*Selection pattern:*

```
* expression instance of java.util.Vector (:obj).removeElement(:method)(*
expressions(:args))
```

*Match calls to the removeElement() method where the obj expression is a subtype of java.util.Vector.*

*Transformation action:*

```
$obj$.remove($obj$.indexOf($args$))
```

*Replace these calls with with calls to the remove() method whose argument is the index of an element to remove.*

**Fig. 6** Example iXj transformation

**Example: iXj** iXj’s pattern language consists of a *selection pattern* and a *transformation action*. iXj’s transformation language allows grouping of code elements using a wild-card symbol \*. Figure 6 shows an example selection pattern and a transformation pattern.

To reduce the burden of learning the iXj pattern language syntax, iXj’s visual editor scaffolds this process through from-example construction and iterative refinement; when a programmer selects an example code fragment to change, iXj automatically generates an initial pattern from the code selection and visualizes all code fragments matched by the initial pattern. The initial pattern is presented in a pattern editor, and a programmer can modify it interactively and see the corresponding matches in the editor. A programmer may edit the transformation action and see the preview of program updates interactively.

### 3.5.2 Programming by Demonstration

Programming by demonstration is also called programming by example (PbE). It is an end-user development technique for teaching a computer or a robot new behaviors by demonstrating the task to transfer directly instead of manually programming the task. Approaches were built to generate programs based on the text-editing actions demonstrated or text change examples provided by users [138, 199, 104, 106]. For instance, TELS records editing actions, such as search and replace, and generalizes them into a program that transforms input to output [199]. It leverages heuristics to match actions against each other to detect any loop in the user-demonstrated program.

SMARTedit is a representative early effort of applying PbD to text editing. It automates repetitive text-editing tasks by learning programs to perform them using techniques drawn from machine learning [106]. SMARTedit represents a text-editing program as a series of functions that alter the state of the text editor (i.e., the contents of the file or the cursor position). Like macro-recording systems, SMARTedit learns the program by observing a user performing her task. However, unlike macro-recorders, SMARTedit examines the context in which the user’s actions are performed and learns programs that work correctly in new contexts.

Below, we describe two seminal PBD approaches applied to software engineering to automate repetitive program changes.

*Simultaneous Editing* Simultaneous editing repetitively applies source code changes that are interactively demonstrated by users [124]. When users apply their edits in one program context, the tool replicates the *exact lexical* edits to other code fragments or transforms code accordingly. Linked Editing requires users to first specify the similar code snippets which they want to modify in the same way [184]. As users interactively edit one of these snippets, Linked Editing simultaneously applies the identical edits to other snippets.

*Systematic Editing* Systematic editing is the process of applying similar, but not necessarily identical, program changes to multiple code locations. High-level changes are often systematic—consisting of related transformations at a code level. In particular, crosscutting concerns, refactoring, and API update mentioned in Sects. 3.3, 3.2, and 3.4 are common kinds of systematic changes, because making these changes during software evolution involves tedious effort of locating individual change locations and applying similar but not identical changes. Several approaches have been proposed to infer the general program transformation from one or more code change examples provided by developers [118, 119, 161] and apply the transformation to other program contexts in need of similar changes. Specifically, LASE requires developers to provide multiple similarly changed code examples in Java (at least two) [119]. By extracting the commonality between demonstrated changes and abstracting the changes in terms of identifier usage and control or data dependency constraints in edit contexts, LASE creates a general program transformation, which can both detect code locations that should be changed similarly and suggest customized code changes for each candidate location. For example, in Fig. 7, LASE can take the change example from  $A_{old}$  to  $A_{new}$  as input and apply to the code on  $B_{old}$  to generate  $B_{new}$ . Such change is similar but customized to the code on the right.

$A_{old}$ to $A_{new}$	$B_{old}$ to $B_{new}$
<pre> public IActionBars getActionBars(){ + IActionBars actionBars =   fContainer.getActionBars(); - if (fContainer == null) { + if (actionBars == null &amp;&amp; ! fContainerProvided){   return   Utilities.findActionBars(fComposite ); } - return fContainer.getActionBars(); + return actionBars; </pre>	<pre> public IServiceLocator getServiceLocator(){ + IServiceLocator serviceLocator =   fContainer.getServiceLocator(); - if (fContainer == null) { + if (serviceLocator == null &amp;&amp; ! fContainerProvided){   return   Utilities.findSite(fComposite); } - return fContainer.getServiceLocator(); + return serviceLocator; </pre>

Fig. 7 An example of noncontiguous, abstract edits that can be applied using LASE [119]



## 4 An Organized Tour of Seminal Papers: Inspecting Changes

Section 4.1 presents the brief history of software inspection and discusses emerging themes from modern code review practices. Sections 4.1.1–4.1.5 discuss various methods that help developers better comprehend software changes, including *change decomposition*, *refactoring reconstruction*, *conflict* and *interference* detection, *related change search*, and *inconsistent change detection*. Section 4.2 describes various program differencing techniques that serve as a basis for analyzing software changes. Section 4.3 describes complementary techniques that record software changes during programming sessions (Fig. 8).

### 4.1 Software Inspection and Modern Code Review Practices

To improve software quality during software evolution, developers often perform *code reviews* to manually examine software changes. Michael Fagan from IBM first introduced “code inspections,” in a seminal paper in 1976 [46]. Code inspections are performed at the end of major software development phases, with the aim of finding overlooked defects before moving to the next phase. Software artifacts are circulated a few days in advance and then reviewed and discussed in a series of meetings. The review meetings include the author of an artifact, other developers to assess the

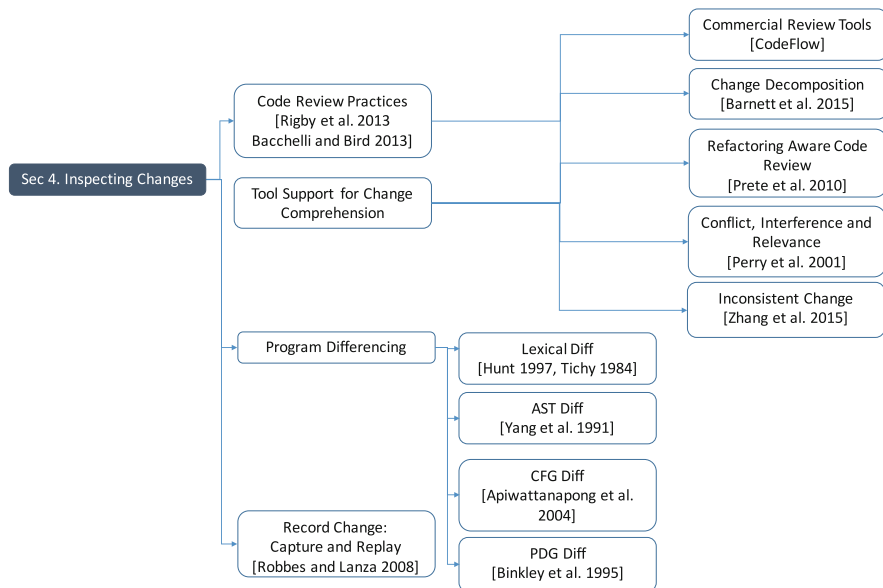


Fig. 8 Change inspection and related research topics

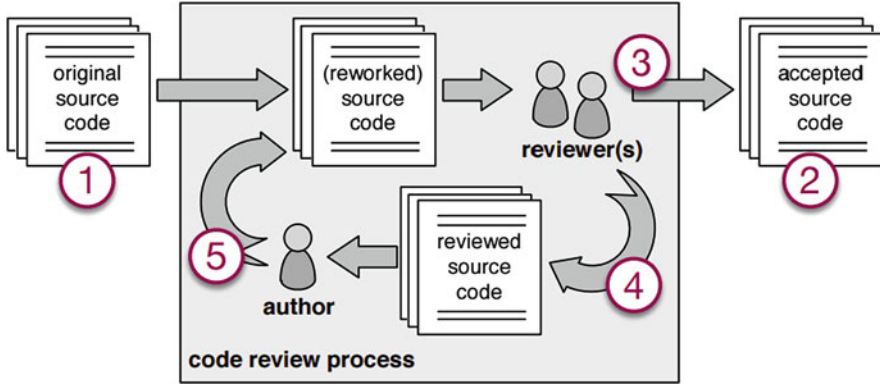


Fig. 9 Modern code review process [14]

artifact, a meeting chair to moderate the discussion, and a secretary to record the discussion. Over the years, code inspections have been proved a valuable method to improve software quality. However, the cumbersome and time-consuming nature of this process hinders its universal adoption in practice [78].

To avoid the inefficiencies in code inspections, most open-source and industrial projects adopt a lightweight, flexible code review process, which we refer to as *modern code reviews*. Figure 9 shows the workflow of modern code reviews. The *author* first submits the *original source code* for review. The *reviewers* then decide whether the submitted code meets the quality acceptance criteria. If not, reviewers can annotate the source code with review comments and send back the *reviewed source code*. The author then revises the code to address reviewers' comments and send it back for further reviews. This process continues till all reviewers accept the revised code.

In contrast to formal code inspections (Fagan style), modern code reviews occur more regularly and informally on program changes. Rigby et al. conducted the first case study about modern code review practices in an open-source software (OSS), Apache HTTP server, using archived code review records in email discussions and version control histories [156]. They described modern code reviews as "early, frequent reviews of small, independent, complete contributions conducted asynchronously by a potentially large, but actually small, group of self-selected experts." As code reviews are practiced in software projects with different settings, cultures, and policies, Rigby and Bird further investigated code review practices using a diverse set of open-source and industrial projects [155]. Despite differences among projects, they found that many characteristics of modern code reviews have converged to similar values, indicating general principles of modern code review practices. We summarize these convergent code review practices as the following.

- *Modern code reviews occur early, quickly, and frequently.* Traditional code inspections happen after finishing a major software component and often last for several weeks. In contrast, modern code reviews happen more frequently and quickly when software changes are committed. For example, the Apache project has review intervals between a few hours to a day. Most reviews are picked up within a few hours among all projects, indicating that reviewers are regularly watching and performing code reviews [155].
- *Modern code reviews often examine small program changes.* During code reviews, the median size of software change varies from 11 to 32 changed lines. The change size is larger in industrial projects, e.g., 44 lines in Android and 78 lines in Chrome, but still much smaller than code inspections, e.g., 263 lines in Lucent. Such small changes facilitate developers to constantly review changes and thus keep up-to-date with the activities of their peers.
- *Modern code reviews are conducted by a small group of self-selected reviewers.* In OSS projects, no reviews are assigned, and developers can select the changes of interest to review. Program changes and review discussions are broadcast to a large group of stakeholders, but only a small number of developers periodically participate in code reviews. In industrial projects, reviews are assigned in a mixed manner—the author adds a group of reviewer candidates and individuals from the group then select changes based on their interest and expertise. On average, two reviewers find an optimal number of defects [155].
- *Modern code reviews are often tool-based.* There is a clear trend toward utilizing review tools to support review tasks and communication. Back in 2008, code reviews in OSS projects were often email-based due to a lack of tool support [156]. In 2013 study, some OSS projects and all industrial projects that they studied used a review tool [155]. More recently, popular OSS hosting services such as GitHub and BitBucket have integrated lightweight review tools to assign reviewers, enter comments, and record discussions. Compared with email-based reviews and traditional software inspections, tool-based reviews provide the benefits of traceability.
- *Although the initial purpose of code review is to find defects, recent studies find that the practices and actual outcomes are less about finding defects than expected.* A study of code reviews at Microsoft found that only a small portion of review comments were related to defects, which were mainly about small, low-level logical issues [7]. Rather, code review provides a spectrum of benefits to software teams, such as knowledge transfer, team awareness, and improved solutions with better practices and readability.

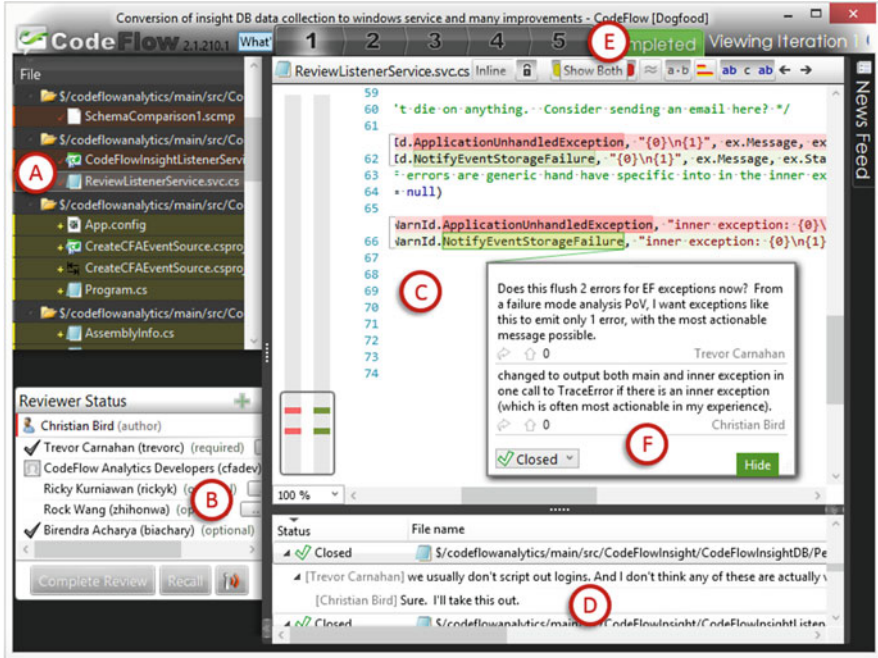


Fig. 10 Example of code review using CodeFlow [17]

### 4.1.1 Commercial Code Review Tools

There is a proliferation of review tools, e.g., Phabricator,<sup>1</sup> Gerrit,<sup>2</sup> CodeFlow,<sup>3</sup> Crucible,<sup>4</sup> and Review Board.<sup>5</sup> We illustrate CodeFlow, a collaborative code review tool at Microsoft. Other review tools share similar functionality as CodeFlow.

To create a review task, a developer uploads changed files with a short description to CodeFlow. Reviewers are then notified via email, and they can examine the software change in CodeFlow. Figure 10 shows the desktop window of CodeFlow. It includes a list of changed files under review (A), the reviewers and their status (B), the highlighted diff in a changed file (C), a summary of all review comments and their status (D), and the iterations of a review (E). If a reviewer would like to provide feedback, she can select a change and enter a comment which is overlaid with the selected change (F). The author and other reviewers can follow up the discussion

<sup>1</sup><http://phabricator.org>.

<sup>2</sup><http://code.google.com/p/gerrit/>.

<sup>3</sup><http://visualstudioextensions.vlasovstudio.com/2012/01/06/codeflow-code-review-tool-for-visual-studio/>.

<sup>4</sup><https://www.atlassian.com/software/crucible>.

<sup>5</sup><https://www.reviewboard.org/>.

by entering comments in the same thread. Typically, after receiving feedback, the author may revise the change accordingly and submit the updated change for additional feedback, which constitutes another review cycle and is termed as an *iteration*. In Fig. 10-E, there are five iterations. CodeFlow assigns a status label to each review comment to keep track of the progress. The initial status is “Active” and can be changed to “Pending,” “Resolved,” “Won’t Fix,” and “Closed” by anyone. Once a reviewer is satisfied with the updated changes, she can indicate this by setting their status to “Signed Off.” After enough reviewers signed off—sign-off policies vary by team—the author can commit the changes to the source repository.

Commercial code review tools facilitate management of code reviews but do not provide deep support for change comprehension. According to Bachhelli et al. [7], understanding program changes and their contexts remains a key challenge in modern code review. Many interviewees acknowledged that it is difficult to understand the rationale behind specific changes. All commercial review tools show the highlighted *textual, line-level diff* of a changed file. However, when the code changes are distributed across multiple files, developers find it difficult to inspect code changes [39]. This obliges reviewers to read changed lines file by file, even when those cross-file changes are done systematically to address the same issue.

#### 4.1.2 Change Decomposition

Prior studies also observe that developers often package program changes of multiple tasks to a single code review [85, 130, 63]. Such large, unrelated changes often lead to difficulty in inspection, since reviewers have to mentally “untangle” them to figure out which subset addresses which issue. Reviewers indicated that they can better understand small, cohesive changes rather than large, tangled ones [156]. For example, a code reviewer commented on Gson revision 1154 saying “I would have preferred to have two different commits: one for adding the new `getFieldNamingPolicy` method and another for allowing overriding of primitives.”<sup>6</sup> Among change decomposition techniques [179, 11], we discuss a representative technique called `CLUSTERCHANGES`.

`CLUSTERCHANGES` is a lightweight static analysis technique for decomposing large changes [11]. The insight is that program changes that address the same issue can be related via implicit dependency such as *def-use* relationship. For example, if a method definition is changed in one location and its call sites are changed in two other locations, these three changes are likely to be related and should be reviewed together. Given a code review task, `CLUSTERCHANGES` first collects the set of definitions for types, fields, methods, and local variables in the corresponding project under review. Then `CLUSTERCHANGES` scans the project for all uses (i.e., references to a definition) of the defined code elements. For instance, any occurrence of a type, field, or method either inside a method or a field initialization is considered

---

<sup>6</sup><https://code.google.com/p/google-gson/source/detail?r=1154>.

to be a use. Based on the extracted def-use information, `CLUSTERCHANGES` identifies three relationships between program changes.

- **Def-use relation.** If the definition of a method or a field is changed, all the uses should also be updated. The change in the definition and the corresponding changes in its references are considered related.
- **Use-use relation.** If two or more uses of a method or a field defined within the change set are changed, these changes are considered related.
- **Enclosing relation.** Program changes in the same method are considered related, under the assumption that (1) program changes to the same method are often related and (2) reviewers often inspect methods atomically rather than reviewing different changed regions in the same method separately.

Given these relations, `CLUSTERCHANGES` creates a partition over the set of program changes by computing a transitive closure of related changes. On the other hand, if a change is not related to any other changes, it will be put into a specific partition of *miscellaneous changes*.

### 4.1.3 Refactoring Aware Code Review

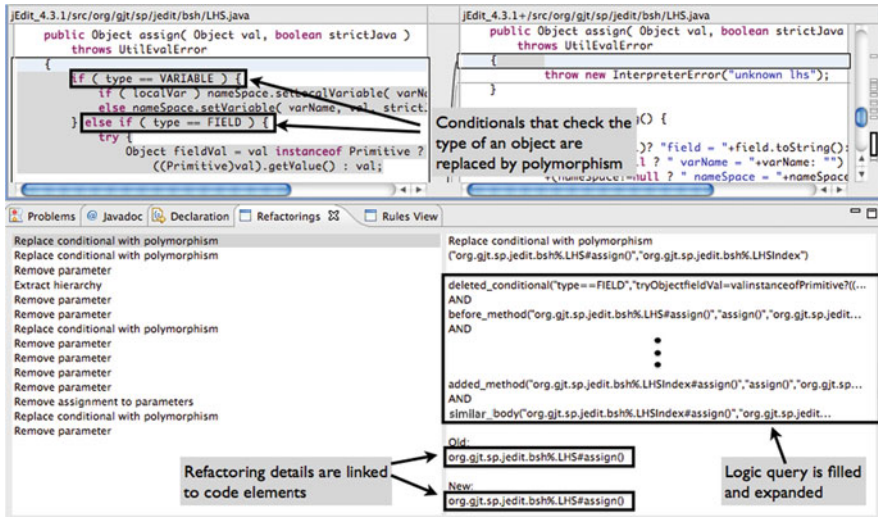
Identifying which refactorings happened between two program versions is an important research problem, because inferred refactorings can help developers understand software modifications made by other developers during peer code reviews. Reconstructed refactorings can be used to update client applications that are broken due to refactorings in library components. Furthermore, they can be used to study the effect of refactorings on software quality empirically when the documentation about past refactorings is unavailable in software project histories.

**Refactoring reconstruction** techniques compare the old and new program versions and identify corresponding entities based on their name similarity and structure similarity [34, 216, 115, 35, 197]. Then based on how basic entities and relations changed from one version to the next, concrete refactoring type and locations are inferred. For example, Xing et al.’s approach [201] UMLDiff extracts class models from two versions of a program, traverses the two models, and identifies corresponding entities based on their name similarity and structure similarity (i.e., similarity in type declaration and uses, field accesses, and method calls). Xing et al. later presented an extended approach to refactoring reconstruction based on change-fact *queries* [202]. They first extract facts regarding design-level entities and relations from each individual source code version. These facts are then pairwise compared to determine how the basic entities and relations have changed from one version to the next. Finally, queries corresponding to well-known refactoring types are applied to the change-fact database to find concrete refactoring instances. Among these refactoring reconstruction techniques, we introduce a representative example of refactoring reconstruction, called RefFinder, in details [147, 92].

*Example: RefFinder* RefFinder is a logic-query-based approach for inferring various types of refactorings in Fowler’s catalog [147]. It first encodes each refactoring type as a structural constraint on the program before and after the refactoring in a template logic rule. It then compares the syntax tree of each version to compute change facts such as `added_subtype`, at the level of code elements (packages, types, methods, and fields), structural dependencies (subtyping, overriding, method calls, and field accesses), and control constructs (while, if statements, and try-catch blocks). It determines a refactoring inference order to find atomic refactorings before composite refactorings (Fig. 11).

For example, consider an *extract superclass* refactoring that extracts common functionality in different classes into a superclass. It finds each *pull-up-method* refactoring and then tests if they combine to an *extract superclass* refactoring. For each refactoring rule, it converts the antecedent of the rule to a logic query and invokes the query on the change-fact database. If the query returns the constant bindings for logic variables, it creates a new logic fact for the found refactoring instance and *writes* it to the fact base. For example, by invoking a query `pull_up_method(?method, ?class, ?superclass) ^ added_type(?superclass)`, it finds a concrete instance of *extract superclass* refactoring. Figure 12 illustrates an example refactoring reconstruction process.

This approach has two advantages over other approaches. First, it analyzes the body of methods including changes to the control structure within method bodies. Thus, it can handle the detection of refactorings such as *replacing conditional code with polymorphism*. Second, it handles composite refactorings, since the approach reasons about which constituent refactorings must be detected first and reason about



**Fig. 11** RefFinder infers a *replace conditionals with polymorphism* refactoring from change facts `deleted_conditional`, `after_subtype`, `before_method`, `added_method` and `similar_body` [92]

pull_up_method template	You have methods with identical results on subclasses; move them to the superclass. $\text{deleted\_method}(m1, n, t1) \wedge \text{after\_subtype}(t2, t1) \wedge \text{added\_method}(m1, n, t2) \Rightarrow$ $\text{pull\_up\_method}(n, t1, t2)$
logic rules	$\text{pull\_up\_method}(m1, t1, t2) \wedge \text{added\_type}(t2) \Rightarrow \text{extract\_superclass}(t1, t2)$
code example	<pre> +public class Customer{ +  chargeFor(start:Date, end:Date) { ... } ...} - public class RegularCustomer{ +public class RegularCustomer extends Customer{ -  chargeFor(start:Date, end:Date){ ... } ...} +public class PreferredCustomer extends Customer{ -  chargeFor(start:Date, end:Date){ ... } // deleted ... } </pre>
found refactorings	<pre> pull_up_method("chargeFor", "RegularCustomer", "Customer") pull_up_method("chargeFor", "PreferredCustomer", "Customer") extract_superclass("RegularCustomer", "Customer") extract_superclass("PreferredCustomer", "Customer") </pre>

**Fig. 12** Reconstruction of *Extract Superclass* refactoring

how those constituent refactorings are knit together to detect higher-level, composite refactorings. It supports 63 out of 72 refactoring types in Fowler’s catalog. As shown in Fig. 11, RefFinder visualizes the reconstructed refactorings as a list. The panel on the right summarizes the key details of the selected refactoring and allows the developer quickly navigate to the associated code fragments.

#### 4.1.4 Change Conflicts, Interference, and Relevance

As development teams become distributed, and the size of the system is often too large to be handled by a few developers, multiple developers often work on the same module at the same time. In addition, the market pressure to develop new features or products makes parallel development no longer an option. A study on a subsystem of Lucent 5ESS telephone found that 12.5% of all changes are made by different developers to the same files within 24 h, showing a high degree of parallel updates [145]. A subsequent study found that even though only 3% of the changes made within 24 h by different developers physically overlapped each other’s changes at a textual level, there was a high degree of semantic interference among parallel changes at a data flow analysis level (about 43% of revisions made within 1 week). They also discovered a significant correlation between files with a high degree of parallel development and the number of defects [165].

Most version control systems are only able to detect most simple types of conflicting changes—changes made on top of other changes [121]. To detect changes that indirectly conflict with each other, some define the notion of *semantic interference* using program slicing on program dependence graphs and integrate non-interfering versions only if there is no overlap between program slices [66]. As another example, some define semantic interference as the overlap between the data-dependence-based impact sets of parallel updates [165].



### 4.1.5 Detecting and Preventing Inconsistent Changes to Clones

Code cloning often requires similar but not identical changes to multiple parts of the system [88], and cloning is an important source of bugs. In 65% of the ported code, at least one identifier is renamed, and in 27% cases, at least one statement is inserted, modified, or deleted [109]. An incorrect adaptation of ported code often leads to porting errors [77]. Interviews with developers confirm that inconsistencies in clones are indeed bugs and report that “nearly every second, unintentional inconsistent changes to clones lead to a fault” [81]. Several techniques find inconsistent changes to similar code fragments by tracking copy-paste code and by comparing the corresponding code and its surrounding contexts [109, 72, 153, 77, 76]. Below, we present a representative technique, called CRITICS.

*Example:* CRITICS CRITICS allows reviewers to interactively detect inconsistent changes through template-based code search and anomaly detection [214]. Given a specified change that a reviewer would like to inspect, CRITICS creates a change template from the selected change, which serves as the pattern for searching similar changes. CRITICS includes *change context* in the template—unchanged, surrounding program statements that are relevant to the selected change. CRITICS models the template as Abstract Syntax Tree (AST) edits and allows reviewers to iteratively customize the template by parameterizing its content and by excluding certain statements. CRITICS then matches the customized template against the rest of the codebase to summarize similar changes and locate potential inconsistent or missing changes. Reviewers can incrementally refine the template and progressively search for similar changes until they are satisfied with the inspection results. This interactive feature allows reviewers with little knowledge of a codebase to flexibly explore the program changes with a desired pattern.

Figure 13 shows a screenshot of CRITICS plugin. CRITICS is integrated with the Compare View in Eclipse, which displays line-level differences per file (see ① in Fig. 13). A user can specify a program change she wants to inspect by selecting the corresponding code region in the Eclipse Compare View. The Diff Template View (see ② in Fig. 13) visualizes the change template of the selected change in a side-by-side view. Reviewers can parameterize concrete identifiers and exclude certain program statements by clicking on the corresponding node in the Diff Template View. Textual Diff Template View (see ⑥ in Fig. 13) shows the change template in a unified format. The Matching Result View summarizes the consistent changes as *similar changes* (see ③ in Fig. 13) and inconsistent ones as *anomalies* (see ④ in Fig. 13).

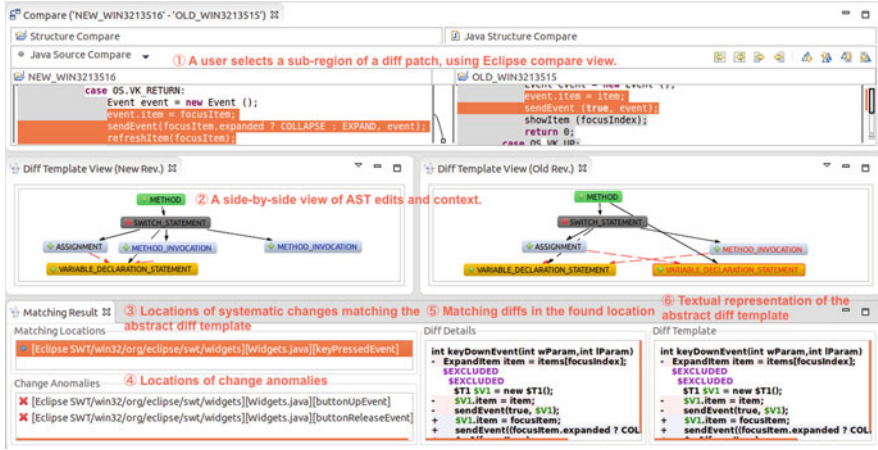


Fig. 13 A screen snapshot of CRITICS’s Eclipse plugin and its features

## 4.2 Program Differencing

Program differencing serves as a basis for analyzing software changes between program versions. The program differencing problem is a dual problem of code matching and is defined as follows:

*Suppose that a program  $P'$  is created by modifying  $P$ . Determine the difference  $\Delta$  between  $P$  and  $P'$ . For a code fragment  $c' \in P'$ , determine whether  $c' \in \Delta$ . If not, find  $c'$ 's corresponding origin  $c$  in  $P$ .*

A code fragment in the new version either contributes to the difference or comes from the old version. If the code fragment has a corresponding origin in the old version, it means that it does not contribute to the difference. Thus, finding the delta between two versions is the same problem as finding corresponding code fragments between two versions.

Suppose that a programmer inserts if-else statements in the beginning of the method  $m\_A$  and reorders several statements in the method  $m\_B$  without changing semantics (see Fig. 14). An intuitively correct matching technique should produce [(p0–c0), (p1–c2), (p2–c3), (p4–c4), (p4–c6), (p5–c7), (p6–c9), (p7–c8), (p8–c10), (p9–c11)] and identify that  $c1$  and  $c5$  are added.

Matching code across program versions poses several challenges. First, previous studies indicate that programmers often disagree about the origin of code elements; low inter-rater agreement suggests that there may be no ground truth in code matching [89]. Second, renaming, merging, and splitting of code elements that are discussed in the context of refactoring reconstruction in Sect. 4.1.3 make the matching problem nontrivial. Suppose that a file `PElmtMatch` changed its name to `PMatching`; a procedure `matchBlck` is split into two procedures `matchDBlck` and `matchCBlck`; and a procedure `matchAST` changed its name to `matchAbstractSyntaxTree`. The intuitively correct matching technique should

**Fig. 14** Example code change

	Past	Current
p0	mA(){	c0 mA(){
p1	if(pred_a){	c1  if(pred_a0){
p2	foo()	c2    if(pred_a){
p3	}	c3      foo()
p4	}	c4      }
p5	mB(b){	c5      }
p6	a:=1	c6      }
p7	b:=b+1	c7  mB(b){
p8	fun(a,b)	c8    b:=b+1
p9	}	c9    a:=1
		c10  fun(a,b)
		c11  }

produce [(PElmtMatch, PMatching), (matchBlck, matchDBlck), (matchBlck, matchCBlck), and (matchAST, matchAbstractSyntaxTree)], while simple name-based matching will consider PMatching, matchDBlck, matchCBlck, and matchAbstract SyntaxTree added and consider PElmtMatch, matchBlck, and matchAST deleted.

Existing code-matching techniques usually employ syntactic and textual similarity measures to match code. They can be characterized by the choices of (1) an underlying program representation, (2) matching granularity, (3) matching multiplicity, and (4) matching heuristics. Below, we categorize program differencing techniques with respect to internal program representations, and we discuss seminal papers for each representation.

### 4.2.1 String and Lexical Matching

When a program is represented as a string, the best match between two strings is computed by finding the longest common subsequence (LCS) [5]. The LCS problem is built on the assumption that (1) available operations are addition and deletion and (2) matched pairs cannot cross one another. Thus, the longest common subsequence does not necessarily include all possible matches when available edit operations include copy, paste, and move. Tichy’s *bdiff* [183] extended the LCS problem by relaxing the two assumptions above: permitting crossing block moves and not requiring one-to-one correspondence.

The line-level LCS implementation, *diff* [69], is fast, reliable, and readily available. Thus, it has served as a basis for popular version control systems such

as CVS. Many evolution analyses are based on *diff* because they use version control system data as input. For example, identification of fix-inducing code snippets is based on line tracking (*file name::function name::line number*) backward from the moment that a bug is fixed [169].

The longest common subsequence algorithm is a dynamic programming algorithm with  $O(mn)$  in time and space, when  $m$  is the line size of the past program and the  $n$  is the line size of the current program. The goal of LCS-based *diff* is to report the minimum number of line changes necessary to convert one file to another. It consists of two phases: (1) computing the length of LCS and (2) reading out the longest common subsequence using a backtrace algorithm. Applying LCS to the example in Fig. 14 will produce the line matching of [(p0–c0), (p1–c1), (p2–c3), (p3–c5), (p4–c6), (p5–c7), (p6–c9), (p8–c10), (p9–c11)]. Due to the assumption of no crossing matches, LCS does not find (p7–c8). In addition, because the matching is done at the line level and LCS does not consider the syntactic structure of code, it produces a line-level match such as (p3–c5) that do not observe the matching block parentheses rule.

#### 4.2.2 Syntax Tree Matching

For software version merging, Yang [206] developed an AST differencing algorithm. Given a pair of functions ( $f_T, f_R$ ), the algorithm creates two abstract syntax trees  $T$  and  $R$  and attempts to match the two tree roots. Once the two roots match, the algorithm aligns  $T$ 's subtrees  $t_1, t_2, \dots, t_i$  and  $R$ 's subtrees  $r_1, r_2, \dots, r_j$  using the LCS algorithm and maps subtrees recursively. This type of tree matching respects the parent-child relationship as well as the order between sibling nodes but is very sensitive to changes in nested blocks and control structures because tree roots must be matched for every level. Because the algorithm respects parent-child relationships when matching code, all matches observe the syntactic boundary of code and the matching block parentheses rule. Similar to LCS, because Yang's algorithm aligns subtrees at the current level by LCS, it cannot find crossing matches caused by code reordering. Furthermore, the algorithm is very sensitive to tree level changes or insertion of new control structures in the middle, because Yang's algorithm performs top-down AST matching.

As another example, Change Distiller [48] uses an improved version of Chawathe et al.'s hierarchically structured data comparison algorithm [23]. Change Distiller takes two abstract syntax trees as input and computes basic tree edit operations such as *insert*, *delete*, *move*, or *update* of tree nodes. It uses *bi-gram string similarity* to match source code statements such as method invocations and uses *subtree similarity* to match source code structures such as if statements. After identifying tree edit operations, Change Distiller maps each tree edit to an atomic AST-level change type.

### 4.2.3 Control Flow Graph Matching

Laski and Szermer [105] first developed an algorithm that computes one-to-one correspondences between CFG nodes in two programs. This algorithm reduces a CFG to a series of single-entry, single-exit subgraphs called hammocks and matches a sequence of hammock nodes using a depth first search (DFS). Once a pair of corresponding hammock nodes is found, the hammock nodes are recursively expanded in order to find correspondences within the matched hammocks.

*Jdiff* [3] extends Laski and Szermer's (LS) algorithm to compare Java programs based on an enhanced control flow graph (ECFG). *Jdiff* is similar to the LS algorithm in the sense that hammocks are recursively expanded and compared but is different in three ways: First, while the LS algorithm compares hammock nodes by the name of a start node in the hammock, *Jdiff* checks whether the ratio of unchanged-matched pairs in the hammock is greater than a chosen threshold in order to allow for flexible matches. Second, while the LS algorithm uses DFS to match hammock nodes, *Jdiff* only uses DFS up to a certain look-ahead depth to improve its performance. Third, while the LS algorithm requires hammock node matches at the same nested level, *Jdiff* can match hammock nodes at a different nested level; thus, *Jdiff* is more robust to addition of while loops or if statements at the beginning of a code segment. *Jdiff* has been used for regression test selection [141] and dynamic change impact analysis [4]. Figure 15 shows the code example and corresponding extended control flow graph representations in Java. Because their representation and matching algorithm is designed to account for dynamic dispatching and exception handling, it can detect changes in the method body of `m3(A a)`, even though it did not have any textual edits: (1) `a.m1()` calls the method definition `B.m()` for the receiver object of type `B` and (2) when the exception type `E3` is thrown, it is caught by the catch block `E1` instead of the catch block `E2`.

CFG-like representations are commonly used in regression test selection research. Rothermel and Harrold [162] traverse two CFGs in parallel and identify a node with unmatched edges, which indicates changes in code. In other words, their algorithm's parallel traversal as soon as it detects changes in a graph structure; thus, this algorithm does not produce deep structural matches between CFGs. However, traversing graphs in parallel is still sufficient for the regression testing problem because it conservatively identifies affected test cases. In practice, regression test selection algorithms [61, 141] require that syntactically changed classes and interfaces are given as input to the CFG matching algorithm.

### 4.2.4 Program Dependence Graph Matching

There are several program differencing algorithms based on a program dependence graph [65, 15, 73].

Horwitz [65] presents a semantic differencing algorithm that operates on a program representation graph (PRG) which combines features of program dependence graphs and static single assignment forms. In her definition, semantic equivalence between two programs  $P1$  and  $P2$  means that, for all states  $\sigma$  such that  $P1$  and  $P2$



**Fig. 15** JDiff change example and CFG representations [4]

halt, the sequence of values produced at  $c_1$  is identical to the sequence of values produced at  $c_2$  where  $c_1$  and  $c_2$  are corresponding locations. Horwitz uses Yang's algorithm [207] to partition the vertices into a group of semantically equivalent vertices based on three properties, (1) the equivalence of their operators, (2) the equivalence of their inputs, and (3) the equivalence of the predicates controlling their evaluation. The partitioning algorithm starts with an initial partition based on the

operators used in the vertices. Then by following flow dependence edges, it refines the initial partition if the successors of the same group are not in the same group. Similarly, it further refines the partition by following control dependence edges. If two vertices in the same partition are textually different, they are considered to have only a *textual change*. If two vertices are in different partitions, they have a *semantic change*. After the partitioning phase, the algorithm finds correspondences between  $P1$ 's vertices and  $P2$ 's vertices that minimize the number of semantically or textually changed components of  $P2$ . In general, PDG-based algorithms are not applicable to popular modern program languages because they can run only on a limited subset of C-like languages without global variables, pointers, arrays, or procedures.

#### 4.2.5 Related Topics: Model Differencing and Clone Detection

A clone detector is simply an implementation of an arbitrary equivalence function. The equivalence function defined by each clone detector depends on a program representation and a comparison algorithm. Most clone detectors are heavily dependent on (1) hash functions to improve performance, (2) parametrization to allow flexible matches, and (3) thresholds to remove spurious matches. A clone detector can be considered as a many-to-many matcher based solely on content similarity heuristics.

In addition to these, several differencing algorithms compare model elements [201, 139, 174, 38]. For example, UMLdiff [201] matches methods and classes between two program versions based on their name. However, these techniques assume that no code elements share the same name in a program and thus use name similarity to produce one-to-one code element matches. Some have developed a general, meta-model-based, configurable program differencing framework [164, 40]. For example, SiDiff [164, 185] allows tool developers to configure various matching algorithms such as identity-based matching, structure-based matching, and signature-based matching by defining how different types of elements need to be compared and by defining the weights for computing an overall similarity measure.

### 4.3 Recording Changes: Edit Capture and Replay

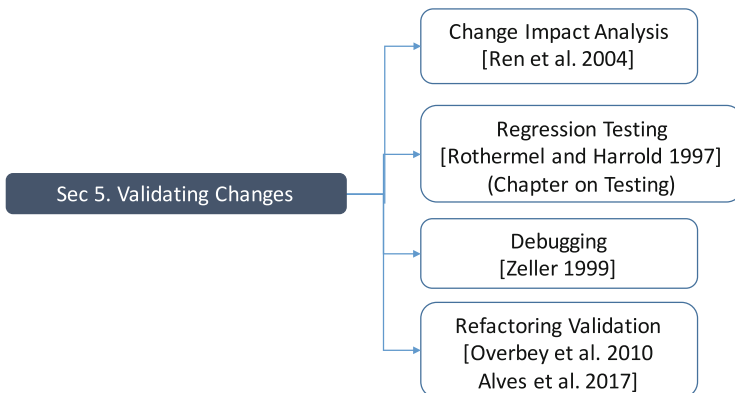
Recorded change operations can be used to help programmers reason about software changes. Several editors or integrated development environment (IDE) extensions capture and replay keystrokes, editing operations, and high-level update commands to use the recorded change information for intelligent version merging, studies of programmers' activities, and automatic updates of client applications. When recorded change operations are used for helping programmers reason about software changes, this approach's limitation depends on the granularity of recorded changes. If an editor records only keystrokes and basic edit operations such as

cut and paste, it is a programmer’s responsibility to raise the abstraction level by grouping keystrokes. If an IDE records only high-level change commands such as refactorings, programmers cannot retrieve a complete change history. In general, capturing change operations to help programmers reason about software change is *impractical* as this approach constrains programmers to use a particular IDE. Below, we discuss a few examples of recording change operations from IDEs:

Spyware is a representative example in this line of work [157]. It is a smalltalk IDE extension to capture AST-level change operations (creation, addition, removal, and property change of an AST node) as well as refactorings. It captures refactorings during development sessions in an IDE rather than trying to infer refactorings from two program versions. Spyware is used to study when and how programmers perform refactorings, but such edit-capture-replay could be used for performing refactoring-aware version merging [37] or updating client applications due to API evolution [62].

## 5 An Organized Tour of Seminal Papers: Change Validation

After making software changes, developers must validate the correctness of updated software. Validation and verification is a vast area of research. In this section, we focus on techniques that aim to identify faults introduced due to software changes. As chapter “Software Testing” discusses the history and seminal work on regression testing in details, we refer the interested readers to that chapter instead. Section 5.1 discusses change impact analysis, which aims to determine the impact of source code edits on programs under test. Section 5.2 discusses how to localize program changes responsible for test failures. Section 5.3 discusses the techniques that are specifically designed to validate refactoring edits under the assumption that software’s external behavior should not change after refactoring (Fig. 16).



**Fig. 16** Change validation and related research topics



### 5.1 Change Impact Analysis

Change impact analysis consists of a collection of techniques for determining the effects of source code modifications and can improve programmer productivity by (a) allowing programmers to experiment with different edits, observe the code fragments that they affect, and use this information to determine which edit to select and/or how to augment test suites; (b) reducing the amount of time and effort needed in running regression tests, by determining that some tests are guaranteed not to be affected by a given set of changes; and (c) reducing the amount of time and effort spent in debugging, by determining a safe approximation of the changes responsible for a given test’s failure.

In this section, we discuss the seminal change impact analysis work, called Chianti, that serves both the purposes of affected test identification and isolation of failure-inducing deltas. It uses a two-phase approach in Fig. 17 [154].

In the first phase, to identify which test cases a developer must rerun on the new version to ensure that all potential regression faults are identified, Chianti takes the old and new program versions  $P_o$  and  $P_n$  and an existing test suite  $T$  as inputs and identifies a set of atomic program changes at the level of methods, fields, and subtyping relationships. It then computes the profile of the test suite  $T$  on  $P_o$  in terms of dynamic call graphs and selects  $T' \subset T$  that guarantees the same regression fault revealing capability between  $T$  and  $T'$ .

In the second phase, Chianti then first runs the selected test cases  $T'$  from the first phase on the new program version  $P_n$  and computes the profile of  $T'$  on  $P_n$  in terms of dynamic call graphs. It then uses both the atomic change set information together with dynamic call graphs to identify which subset of the delta between  $P_o$  and  $P_n$  led to the behavior differences for each failed test on  $P_n$ .

To represent atomic changes, Chianti compares the syntax tree of the old and new program versions and decomposes the edits into atomic changes at a method and field level. Changes are then categorized as added classes (AC), deleted classes (DC), added methods (AM), deleted methods (DM), changed methods (CM), added

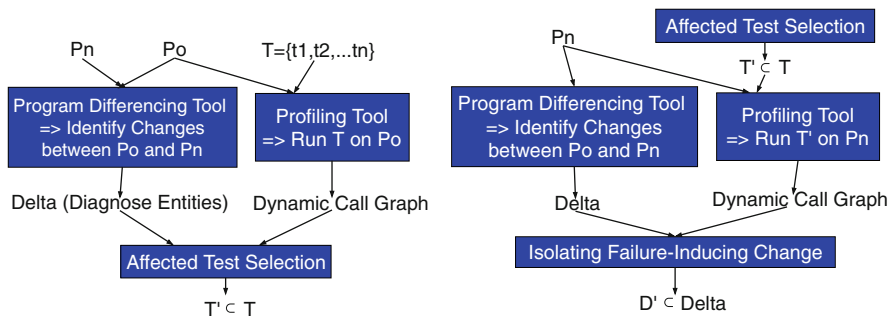


Fig. 17 Chianti change impact analysis: identifying affected tests (left) and identifying affecting change (right) [154]

fields (AF), deleted fields (DF), and lookup (i.e., dynamic dispatch) changes (LC). The LC atomic change category models changes to the dynamic dispatch behavior of instance methods. In particular, an LC change  $LC(Y, X.m())$  models the fact that a call to method  $X.m()$  on an object of type  $Y$  results in the selection of a different method call target.

For example, Fig. 18 shows a software change example and corresponding lists of atomic changes inferred from AST-level comparison. An arrow from an atomic change  $A1$  to an atomic change  $A2$  indicates that  $A2$  is dependent on  $A1$ . For example, the addition of the call  $B.bar()$  in method  $B.foo()$  is the method body change  $CM(B.foo())$  represented as ⑧. This change ⑧ requires the declaration of method  $B.bar()$  to exist first, i.e.,  $AM(B.bar())$  represented as ⑥. This dependence is represented as an arrow from ⑥ to ⑧.

Phase I reports **affected tests**—a subset of regression tests relevant to edits. It identifies a test if its dynamic call graph on the old version contains a node that corresponds to a changed method (CM) or deleted method (DM) or if the call graph contains an edge that corresponds to a lookup change (LC). Figure 18 also shows the dynamic call graph of each test for the old version (left) and the new version (right). Using the call graphs on the left, it is easy to see that (a)  $test1$  is not affected; (b)  $test2$  is affected because its call graph contains a node for  $B.foo()$ , which corresponds to ⑧; and (c)  $test3$  is affected because its call graph contains an edge corresponding to a dispatch to method  $A.foo()$  on an object of type  $C$ , which corresponds to ④.

Phase II then reports **affecting changes**—a subset of changes relevant to the execution of affected tests in the new version. For example, we can compute the affecting changes for  $test2$  as follows. The call graph for  $test2$  in the edited version of the program contains methods  $B.foo()$  and  $B.bar()$ . These nodes correspond to ⑧ and ⑨, respectively. Atomic change ⑧ requires ⑥ and ⑨ requires ⑥ and ⑦. Therefore, the atomic changes affecting  $test2$  are ⑥, ⑦, ⑧, and ⑨. Informally, this means that we can automatically determine that  $test2$  is affected by the addition of field  $B.y$ , the addition of method  $B.bar()$ , and the change to method  $B.foo()$ , but not on any of the other source code changes.

## 5.2 Debugging Changes

The problem of simplifying and isolating failure-inducing input is a long-standing problem in software engineering. *Delta Debugging (DD)* addresses this problem by repetitively running a program with different sub-configurations (subsets) of the input to systematically isolate failure-inducing inputs [211, 212]. DD splits the original input into two halves using a binary search-like strategy and reruns them. DD requires a test oracle function  $test(c)$  that takes an input configuration  $c$  and checks whether running a program with  $c$  leads to a failure. If one of the two halves fails, DD recursively applies the same procedure for only that failure-inducing input configuration. On the other hand, if both halves pass, DD tries

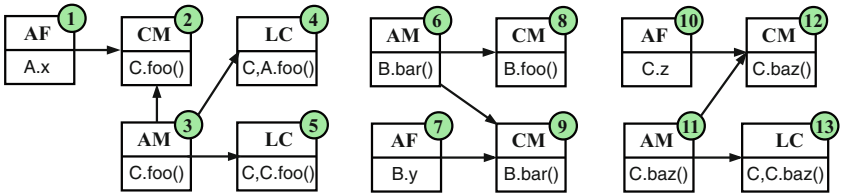
```

class A {
    public A(){ }
    public void foo(){ }
    public int x;
}
class B extends A {
    public B(){ }
    public void foo(){ B.bar(); }
    public static void bar(){ y = 17; }
    public static int y;
}
class C extends A {
    public C(){ }
    public void foo(){ x = 18; }
    public void baz(){ z = 19; }
    public int z;
}
    
```

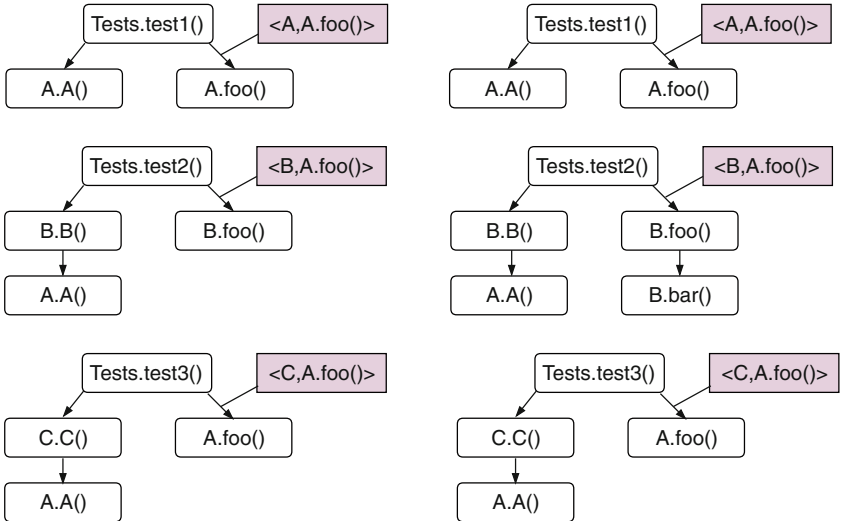
```

class Tests {
    public static void test1(){
        A a = new A();
        a.foo();
    }
    public static void test2(){
        A a = new B();
        a.foo();
    }
    public static void test3(){
        A a = new C();
        a.foo();
    }
}
    
```

(a)



(b)



(c)

**Fig. 18** Chianti change impact analysis. (a) Example program with three tests. Added code fragments are shown in boxes. (b) Atomic changes for the example program, with their inter-dependencies. (c) Call graphs for the tests before and after the changes were applied

different sub-configurations by mixing fine-grained sub-configurations with larger sub-configurations (computed as the complement from the current configuration).

Under the assumption that failure is *monotone*, where  $C$  is a super set of all configurations, if a larger configuration  $c$  is successful, then any of its smaller sub-configurations  $c'$  does not fail, that is,  $\forall c \in C (test(c) = \checkmark \rightarrow \forall c' \subset c (test(c') \neq \times))$ , DD returns a minimal failure-inducing configuration.

This idea of Delta Debugging was applied to isolate failure-inducing changes. It considers all line-level changes between the old and new program version as the candidate set without considering compilation dependences among those changes. In Zeller's seminal paper, "yesterday, my program worked, but today, it does not, why?" Zeller demonstrates the application of DD to isolate program edits responsible for regression failures [211]. DDD 3.1.2, released in December, 1998, exhibited a nasty behavioral change: When invoked with the name of a non-existing file, DDD 3.1.2 dumped core, while its predecessor DDD 3.1.1 simply gave an error message. The DDD configuration management archive lists 116 logical changes between the 3.1.1 and 3.1.2 releases. These changes were split into 344 textual changes to the DDD source. After only 12 test runs and 58 min, the failure-inducing change was found:

```
diff -r1.30 -r1.30.4.1 ddd/gdbinit.C
295,296c296
<
< --- >
string classpath =
getenv("CLASSPATH") != 0 ? getenv("CLASSPATH") : ".";
string classpath = source view->class path();
```

When called with an argument that is not a file name, DDD 3.1.1 checks whether it is a Java class; so DDD consults its environment for the class lookup path. As an "improvement," DDD 3.1.2 uses a dedicated method for this purpose. Unfortunately, the source view pointer used is initialized only later, resulting in a core dump.

*Spectra-Based Fault Localization* Spectrum-based fault localization techniques such as Tarantula [80] statistically compute suspiciousness scores for statements based on execution traces of both passed and failed test cases and rank potential faulty statements based on the derived suspiciousness scores. Researchers have also introduced more suspiciousness computation measures to the realm of fault localization for localizing faulty statements [132, 112] and also developed various automated tool sets which embodies different spectrum-based fault localization techniques [74]. However, such spectrum-based fault localization techniques are not scalable to large evolving software systems, as they compute spectra on all statements in each program version and do not leverage information about program edits between the old and new versions.

To address this problem, FaultTracer [213] combines Chianti-style change impact analysis and Tarantula-style fault localization. To present a ranked list of

potential failure-inducing edits, FaultTracer applies a set of spectrum-based ranking techniques to the affecting changes determined by Chianti-style change impact analysis. It uses a new enhanced call graph representation to measure test spectrum information directly for field-level edits and to improve upon the existing Chianti algorithm. The experimental results show that FaultTracer outperforms Chianti in selecting affected tests (slightly better) as well as in determining affecting changes (with an improvement of approximately 20%). By ranking the affecting changes using spectrum-based profile, it places a real regression fault within a few atomic changes, significantly reducing developers' effort in inspecting potential failure-inducing changes.

### 5.3 Refactoring Validation

Unlike other types of changes, refactoring validation is a special category of change validation. By definition, refactoring must guarantee behavior preservation, and thus the old version's behavior could be compared against the new version's behavior for behavior preservation. Regression testing is the most used strategy for checking refactoring correctness. However, a recent study finds that test suites are often inadequate [149] and developers may hesitate to initiate or perform refactoring tasks due to inadequate test coverage [94]. Soares et al. [171] design and implement SafeRefactor that uses randomly generated test suites for detecting refactoring anomalies.

Formal verification is an alternative for avoiding refactoring anomalies [122]. Some propose rules for guaranteeing semantic preservation [28], use graph rewriting for specifying refactorings [123], or present a collection of refactoring specifications, which guarantee the correctness by construction [142]. However, these approaches focus on improving the correctness of automated refactoring through formal specifications only. Assuming that developers may apply refactoring manually rather, Schaeffer et al. validate refactoring edits by comparing data and control dependences between two program versions [163].

RefDistiller is a static analysis approach [2, 1] to support the inspection of manual refactorings. It combines two techniques. First, it applies predefined templates to identify potential missed edits during manual refactoring. Second, it leverages an automated refactoring engine to identify extra edits that might be incorrect, helping to determine the root cause of detected refactoring anomalies. GhostFactor [50] checks the correctness of manual refactoring, similar to RefDistiller. Another approach by Ge and Murphy-Hill [42] helps reviewers by identifying applied refactorings and letting developers examine them in isolation by separating pure refactorings.

## 6 Future Directions and Open Problems

Software maintenance is challenging and time-consuming. Albeit various research and existing tool support, the global cost of debugging software has risen up to \$312 billion annually [20]. The cost of software maintenance is rising dramatically and has been estimated as more than 90% of the total cost for software [43]. Software evolution research still has a long future ahead, because there are still challenges and problems that cost developers a lot of time and manual effort. In this section, we highlight some key issues in change comprehension and suggestion.

### 6.1 Change Comprehension

Understanding software changes made by other people is a difficult task, because it requires not only the domain knowledge of the software under maintenance but also the comprehension of change intent and the interpretation of mappings between the program semantics of applied changes and those intent. Existing change comprehension tools discussed in Sect. 4.1 and program differencing tools discussed in Sect. 4.2 mainly present the textual or syntactical differences between the before and after versions of software changes. Current large-scale empirical studies on code changes discussed in Sects. 3.1–3.4 also mainly focus on textual or syntactical notion of software changes. However, there is no tool support to automatically summarize the semantics of applied changes or further infer developers' intent behind the changes.

The new advanced change comprehension tools must assist software professionals in two aspects. First, by summarizing software changes with a natural language description, these tools must produce more meaningful commit messages when developers check in their program changes to software version control systems (e.g., SVN, Git) to facilitate other people (e.g., colleagues and researchers) to mine, comprehend, and analyze applied changes more precisely [63]. Second, the generated change summary must provide a second opinion to developers of the changes and enable them to easily check whether the summarized change description matches their actual intent. If there is a mismatch, developers should carefully examine the applied changes and decide whether the changes reflect or realize their original intent.

To design and implement such advanced change comprehension tools, researchers must address several challenges.

1. How should we correlate changes applied in source code, configuration files, and databases to present all relevant changes and their relationships as a whole? For instance, how can we explain why a configuration file is changed together with a function's code body? How are the changes in a database schema correspond to source code changes?

2. How should we map concrete code changes or abstract change patterns to natural language descriptions? For instance, when complicated code changes are applied to improve a program's performance, how can we detect or reveal that intent? How should we differentiate between different types of changes when inferring change intent or producing natural language descriptions accordingly?
3. When developers apply multiple kinds of changes together, such as refactoring some code to facilitate feature addition, can we identify the boundary between the different types of changes? How can we summarize the changes in a meaningful way so that both types of changes are identified and the connection between them is characterized clearly?

To solve these challenges, we may need to invent new program analysis techniques to correlate changes, new change interpretation approaches to characterize different types of changes, and new text mining and natural language processing techniques to map changes to natural language descriptions.

## 6.2 Change Suggestion

Compared with understanding software changes, applying changes is even more challenging and can cause serious problems if changes are wrongly applied. Empirical studies showed that 15–70% of the bug fixes applied during software maintenance were incorrect in their first release [167, 209], which indicates a desperate need for more sophisticated change suggestion tools. Below we discuss some of the limitations of existing automatic tool support and also suggest potential future directions.

**Corrective Change Suggestion** Although various bug fix and program repair tools discussed in Sect. 3.1 detect different kinds of bugs or even suggest bug fixes, the suggested fixes are usually relatively simple. They may focus on single-line bug fixes, multiple if-condition updates, missing APIs to invoke, or similar code changes that are likely to be applied to similar code snippets. However, no existing approach can suggest a whole missing `if`-statement or `while`-loop, neither can they suggest bug fixes that require declaring a new method and inserting the invocation to the new method in appropriate code locations.

**Adaptive Change Suggestion** Existing adaptive change support tools discussed in Sect. 3.2 allow developers to migrate programs between specific previously known platforms (e.g., desktop and cloud). However, it is not easy to extend these tools when a new platform becomes available and people need to migrate programs from existing platforms to the new one. Although cross-platform software development tools can significantly reduce the necessity of platform-to-platform migration tools, these tools are limited to the platforms for which they are originally built. When a new platform becomes available, these tools will undergo significant modifications to support the new platform. In the future, we need extensible program migration

frameworks, which will automatically infer program migration transformations from the concrete migration changes manually applied by developers and then apply the inferred transformations to automate other migration tasks for different target platforms. With such frameworks, developers will not need to manually apply repetitive migration changes.

**Perfective Change Suggestion** There are some programming paradigms developed (e.g., AOP and FOP discussed in Sect. 3.3), which facilitate developers to apply perfective changes to enhance or extend any existing software. However, there is no tool support to automatically suggest what perfective changes to apply and where to apply those changes. The main challenge of creating such tools is that unlike other types of changes, perfective changes usually aim to introduce new features instead of modifying existing features. Without any hint provided by developers, it is almost impossible for any tool to predict what new features to add to the software. However, when developers know what new features they want to add but do not know how to implement those features, some advanced tools can be helpful by automatically searching for relevant open-source projects, identifying relevant code implementation for the queried features, or even providing customized change suggestion to implement the features and to integrate the features into existing software.

**Preventive Change Suggestion** Although various refactoring tools discussed in Sect. 3.4 can automatically refactor code, all the supported refactorings are limited to predefined behavior-preserving program transformations. It is not easy to extend existing refactoring tools to automate new refactorings, especially when the program transformation involves modifications of multiple software entities (i.e., classes, methods, and fields). Some future tools should be designed and implemented to facilitate the extensions of refactoring capabilities. There are also some refactoring tools that suggest refactoring opportunities based on code smells. For instance, if there are many code clones in a codebase, existing tools can suggest a clone removal refactoring to reduce duplicated code. In reality, nevertheless, most of the time developers apply refactorings only when they want to apply bug fixes or add new features, which means that refactorings are more likely to be motivated by other kinds of changes instead of code smells and change history [168]. In the future, with the better change comprehension tools mentioned above, we may be able to identify the trends of developers' change intent in the past and observe how refactorings were applied in combination with other types of changes. Furthermore, with the observed trends, new tools must be built to predict developers' change intent in future and then suggest refactorings accordingly to prepare for the upcoming changes.

### 6.3 Change Validation

In terms of change validation discussed in Sect. 5, there is disproportionately more work being done in the area of validating refactoring (i.e., *preventative changes*),



compared to other types of changes such as *adaptive* and *perfective* changes. Similarly, in the absence of adequate existing tests which helped to discover defects in the first place, it is not easy to validate whether *corrective changes* are applied correctly to fix the defects.

The reason why is that, with the exception of refactoring that has a canonical, straightforward definition of *behavior preserving modifications*, when it comes to other types of software changes, it is difficult to define the updated semantics of software systems. For example, when a developer adds a new feature, how can we know the desired semantics of the updated software?

This problem naturally brings up the needs of having the correct specifications of updated software and having easier means to write such specifications in the context of software changes. Therefore, new tools must be built to guide developers in writing software specifications for the changed parts of the systems. In particular, we see a new opportunity for tool support suggests the template for updated specifications by recognizing the type and pattern of program changes to guide developers in writing updated specifications—Are there common specification patterns for each common type of software changes? Can we then suggest which specifications to write based on common types of program modifications such as API evolution? Such tool support must not require developers to write specifications from scratch but rather guide developers on which specific parts of software require new, updated specifications, which parts of software may need additional tests, and how to leverage those written specifications effectively to guide the remaining areas for writing better specifications. We envision that, with such tool support for reducing the effort of writing specifications for updated software, researchers can build change validation techniques that actively leverage those specifications. Such effort will contribute to expansion of change-type-specific debugging and testing technologies.

## Appendix

The following text box shows selected, recommended readings for understanding the area of software evolution.

## Key References

- Apiwattanapong, T., Orso, A., Harrold, M.J.: A differencing algorithm for object-oriented programs. In: ASE '04: Proceedings of the 19th IEEE International Conference on Automated Software Engineering, pp. 2–13. IEEE Computer Society, Washington (2004)
- Bacchelli, A., Bird, C.: Expectations, outcomes, and challenges of modern code review. In: Proceedings of the 2013 International Conference on Software Engineering, pp. 712–721. IEEE Press, Piscataway (2013)
- Cordy, J.R.: The txl source transformation language. *Sci. Comput. Program.* **61**(3), 190–210 (2006)
- Engler, D.R., Chen, D.Y., Chou, A.: Bugs as inconsistent behavior: a general approach to inferring errors in systems code. In: Symposium on Operating Systems Principles, pp. 57–72 (2001)
- Henkel, J., Diwan, A.: Catchup!: capturing and replaying refactorings to support API evolution. In: ICSE '05: Proceedings of the 27th International Conference on Software Engineering, pp. 274–283. ACM, New York (2005)
- Kim, M., Sazawal, V., Notkin, D., Murphy, G.: An empirical study of code clone genealogies. In: Proceedings of the 10th European Software Engineering Conference Held Jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering, ESEC/FSE-13, pp. 187–196. ACM, New York (2005)
- Kim, M., Zimmermann, T., Nagappan, N.: A field study of refactoring challenges and benefits. In: Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering, FSE '12, pp. 50:1–50:11. ACM, New York (2012)
- Prete, K., Rachatasumrit, N., Sudan, N., Kim, M.: Template-based reconstruction of complex refactorings. In: 2010 IEEE International Conference on Software Maintenance (ICSM), pp. 1–10. IEEE Press, Piscataway (2010)
- Ren, X., Shah, F., Tip, F., Ryder, B.G., Chesley, O.: Chianti: a tool for change impact analysis of java programs. In: OOPSLA '04: Proceedings of the 19th annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, pp. 432–448. ACM, New York (2004)
- Tarr, P., Oshser, H., Harrison, W., Sutton, JSM.: N degrees of separation: multi-dimensional separation of concerns. In: ICSE '99: Proceedings of the 21st International Conference on Software Engineering, pp. 107–119. IEEE Computer Society Press, Los Alamitos (1999)
- Weimer, W., Nguyen, T., Le Goues, C., Forrest, S.: Automatically finding patches using genetic programming. In: Proceedings of the 31st International Conference on Software Engineering, ICSE '09, pp. 364–374. IEEE Computer Society, Washington (2009)
- Zeller, A.: Yesterday, my program worked. today, it does not. Why? In: ESEC/FSE-7: Proceedings of the 7th European Software Engineering Conference Held Jointly with the 7th ACM SIGSOFT International Symposium on Foundations of Software Engineering, pp. 253–267. Springer, London (1999)

## References

1. Alves, E.L.G., Song, M., Kim, M.: Refdistiller: a refactoring aware code review tool for inspecting manual refactoring edits. In: Proceedings of the 22Nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2014, pp. 751–754. ACM, New York (2014)
2. Alves, E.L.G., Song, M., Massoni, T., Machado, P.D.L., Kim, M.: Refactoring inspection support for manual refactoring edits. *IEEE Trans. Softw. Eng.* **PP**(99), 1–1 (2017)
3. Apiwattanapong, T., Orso, A., Harrold, M.J.: A differencing algorithm for object-oriented programs. In: ASE '04: Proceedings of the 19th IEEE International Conference on Automated Software Engineering, pp. 2–13. IEEE Computer Society, Washington (2004)
4. Apiwattanapong, T., Orso, A., Harrold, M.J.: Efficient and precise dynamic impact analysis using execute-after sequences. In: ICSE '05: Proceedings of the 27th International Conference on Software Engineering, pp. 432–441. ACM, New York (2005)
5. Apostolico, A., Galil, Z. (eds.): *Pattern Matching Algorithms*. Oxford University Press, Oxford (1997). Program differencing LCS
6. ASM. <http://asm.ow2.org>
7. Bacchelli, A., Bird, C.: Expectations, outcomes, and challenges of modern code review. In: Proceedings of the 2013 International Conference on Software Engineering, pp. 712–721. IEEE Press, Piscataway (2013)
8. Balazinska, M., Merlo, E., Dagenais, M., Lague, B., Kontogiannis, K.: Partial redesign of java software systems based on clone analysis. In: WCRE '99: Proceedings of the Sixth Working Conference on Reverse Engineering, p. 326. IEEE Computer Society, Washington (1999)
9. Balazinska, M., Merlo, E., Dagenais, M., Lague, B., Kontogiannis, K.: Advanced clone-analysis to support object-oriented system refactoring. In: Proceedings Seventh Working Conference on Reverse Engineering, pp. 98–107 (2000)
10. Baldwin, C.Y., Clark, K.B.: *Design Rules: The Power of Modularity*. MIT Press, Cambridge (1999)
11. Barnett, M., Bird, C., Brunet, J., Lahiri, S.K.: Helping developers help themselves: automatic decomposition of code review changesets. In: Proceedings of the 37th International Conference on Software Engineering–Volume 1, pp. 134–144. IEEE Press, Piscataway (2015)
12. Batory, D., O'Malley, S.: The design and implementation of hierarchical software systems with reusable components. *ACM Trans. Softw. Eng. Methodol.* **1**(4), 355–398 (1992)
13. Belady, L.A., Lehman, M.M.: A model of large program development. *IBM Syst. J.* **15**(3), 225–252 (1976)
14. Beller, M., Bacchelli, A., Zaidman, A., Juergens, E.: Modern code reviews in open-source projects: which problems do they fix? In: Proceedings of the 11th Working Conference on Mining Software Repositories, pp. 202–211. ACM, New York (2014)
15. Binkley, D., Horwitz, S., Reps, T.: Program integration for languages with procedure calls. *ACM Trans. Softw. Eng. Methodol.* **4**(1), 3–35 (1995)
16. Boshernitsan, M., Graham, S.L., Hearst, M.A.: Aligning development tools with the way programmers think about code changes. In: CHI '07: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pp. 567–576. ACM, New York (2007)
17. Bosu, A., Greiler, M., Bird, C.: Characteristics of useful code reviews: an empirical study at microsoft. In: 2015 IEEE/ACM 12th Working Conference on Mining Software Repositories (MSR), pp. 146–156. IEEE, Piscataway (2015)
18. Breu, S., Zimmermann, T.: Mining aspects from version history. In: International Conference on Automated Software Engineering, pp. 221–230 (2006)
19. Brown, N., Cai, Y., Guo, Y., Kazman, R., Kim, M., Kruchten, P., Lim, E., MacCormack, A., Nord, R., Ozkaya, I., Sangwan, R., Seaman, C., Sullivan, K., Zazworka, N.: Managing technical debt in software-reliant systems. In: Proceedings of the FSE/SDP Workshop on Future of Software Engineering Research, FoSER '10, pp. 47–52. ACM, New York (2010)

20. Cambridge University Study States Software Bugs Cost Economy \$312 Billion Per Year. [http://markets.financialcontent.com/stocks/news/read/23147130/Cambridge\\_University\\_Study\\_States\\_Software\\_Bugs\\_Cost\\_Economy\\_\\$312\\_Billion\\_Per\\_Year](http://markets.financialcontent.com/stocks/news/read/23147130/Cambridge_University_Study_States_Software_Bugs_Cost_Economy_$312_Billion_Per_Year)
21. Canfora, G., Cerulo, L., Cimitile, M., Di Penta, M.: Social interactions around cross-system bug fixings: the case of freebsd and opensd. In: Proceeding of the 8th Working Conference on Mining Software Repositories, MSR '11, pp. 143–152. ACM, New York (2011)
22. Carriere, J., Kazman, R., Ozkaya, I.: A cost-benefit framework for making architectural decisions in a business context. In: ICSE '10: Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering, pp. 149–157. ACM, New York (2010)
23. Chawathe, S.S., Rajaraman, A., Garcia-Molina, H., Widom, J.: Change detection in hierarchically structured information. In: SIGMOD '96: Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data, pp. 493–504. ACM, New York (1996)
24. Chou, A., Yang, J., Chelf, B., Hallem, S., Engler, D.: An empirical study of operating systems errors. In: Proceedings of the Eighteenth ACM Symposium on Operating Systems Principles, SOSP '01, pp. 73–88. ACM, New York (2001)
25. Chow, K., Notkin, D.: Semi-automatic update of applications in response to library changes. In: ICSM '96: Proceedings of the 1996 International Conference on Software Maintenance, p. 359. IEEE Computer Society, Washington (1996)
26. Cordy, J.R.: The txl source transformation language. *Sci. Comput. Program.* **61**(3), 190–210 (2006)
27. Cordy, J.R.: Exploring large-scale system similarity using incremental clone detection and live scatterplots. In: 2011 IEEE 19th International Conference on Program Comprehension (2011), pp. 151–160
28. Cornélio, M., Cavalcanti, A., Sampaio, A.: Sound refactorings. *Sci. Comput. Program.* **75**(3), 106–133 (2010)
29. Cossette, B.E., Walker, R.J.: Seeking the ground truth: a retroactive study on the evolution and migration of software libraries. In: FSE '12 Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering. ACM, New York (2012)
30. Cottrell, R., Chang, J.J.C., Walker, R.J., Denzinger, J.: Determining detailed structural correspondence for generalization tasks. In: Proceedings of the the 6th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering, ESEC-FSE '07, pp. 165–174. ACM, New York (2007)
31. Cunningham, W.: The WyCash portfolio management system. In: OOPSLA '92: Addendum to the Proceedings on Object-Oriented Programming Systems, Languages, and Applications (Addendum), pp. 29–30. ACM, New York (1992)
32. Dagenais, B., Robillard, M.P.: Recommending adaptive changes for framework evolution. In: Proceedings of the 30th International Conference on Software Engineering, ICSE '08, pp. 481–490. ACM, New York (2008)
33. Dagenais, B., Breu, S., Warr, F.W., Robillard, M.P.: Inferring structural patterns for concern traceability in evolving software. In: ASE '07: Proceedings of the Twenty-Second IEEE/ACM International Conference on Automated Software Engineering, pp. 254–263. ACM, New York (2007)
34. Demeyer, S., Ducasse, S., Nierstrasz, O.: Finding refactorings via change metrics. In: OOPSLA '00: Proceedings of the 15th ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, pp. 166–177. ACM, New York (2000)
35. Dig, D., Johnson, R.: Automated detection of refactorings in evolving components. In: ECOOP '06: Proceedings of European Conference on Object-Oriented Programming, pp. 404–428. Springer, Berlin (2006)
36. Dig, D., Johnson, R.: How do APIs evolve? A story of refactoring. *J. Softw. Maint. Evol. Res. Pract.* **18**(2), 83–107 (2006)

37. Dig, D., Manzoor, K., Johnson, R., Nguyen, T.N.: Refactoring-aware configuration management for object-oriented programs. In: 29th International Conference on Software Engineering, 2007, ICSE 2007, pp. 427–436 (2007)
38. Duley, A., Spandikow, C., Kim, M.: Vdiff: a program differencing algorithm for verilog hardware description language. *Autom. Softw. Eng.* **19**, 459–490 (2012)
39. Dunsmore, A., Roper, M., Wood, M.: Object-oriented inspection in the face of delocalisation. In: ICSE '00: Proceedings of the 22nd International Conference on Software Engineering, pp. 467–476. ACM, New York (2000). Code inspection, code review, object-oriented, delocalized
40. Eclipse EMF Compare Project description: <http://www.eclipse.org/emft/projects/compare>
41. Eick, S.G., Graves, T.L., Karr, A.F., Marron, J.S., Mockus, A.: Does code decay? Assessing the evidence from change management data. *IEEE Trans. Softw. Eng.* **27**(1), 1–12 (2001)
42. EmersonMurphy-Hill, X.S.: Towards refactoring-aware code review. In: CHASE' 14: 7th International Workshop on Cooperative and Human Aspects of Software Engineering, Co-located with 2014 ACM and IEEE 36th International Conference on Software Engineering (2014)
43. Engelbertink, F.P., Vogt, H.H.: How to save on software maintenance costs. Omnext white paper (2010)
44. Engler, D., Chelf, B., Chou, A., Hallem, S.: Checking system rules using system-specific, programmer-written compiler extensions. In: Proceedings of the 4th Conference on Symposium on Operating System Design & Implementation - Volume 4, OSDI'00. USENIX Association, Berkeley (2000)
45. Engler, D.R., Chen, D.Y., Chou, A.: Bugs as inconsistent behavior: A general approach to inferring errors in systems code. In: Symposium on Operating Systems Principles, pp. 57–72 (2001)
46. Fagan, M.E.: Design and code inspections to reduce errors in program development. *IBM Syst. J.* **38**(2–3), 258–287 (1999). Code inspection, checklist
47. Fischer, M., Oberleitner, J., Ratzinger, J., Gall, H.: Mining evolution data of a product family. In: MSR '05: Proceedings of the 2005 International Workshop on Mining Software Repositories, pp. 1–5. ACM, New York (2005)
48. Fluri, B., Würsch, M., Pinzger, M., Gall, H.C.: Change distilling—tree differencing for fine-grained source code change extraction. *IEEE Trans. Softw. Eng.* **33**(11), 18 (2007)
49. Garcia, J., Popescu, D., Edwards, G., Medvidovic, N.: Identifying architectural bad smells. In: CSMR '09: Proceedings of the 2009 European Conference on Software Maintenance and Reengineering, pp. 255–258. IEEE Computer Society, Washington (2009)
50. Ge, X., Murphy-Hill, E.: Manual refactoring changes with automated refactoring validation. In: 36th International Conference on Software Engineering (ICSE 2014). IEEE, Piscataway (2014)
51. Görg, C., Weißgerber, P.: Error detection by refactoring reconstruction. In: MSR '05: Proceedings of the 2005 International Workshop on Mining Software Repositories, pp. 1–5. ACM Press, New York (2005)
52. Griswold, W.G.: Program restructuring as an aid to software maintenance. PhD thesis, Seattle (1992). UMI Order No. GAX92-03258
53. Griswold, W.: Coping with crosscutting software changes using information transparency. In: Reflection 2001: The Third International Conference on Metalevel Architectures and Separation of Crosscutting Concerns, pp. 250–265. Springer, Berlin (2001)
54. Griswold, W.G., Atkinson, D.C., McCurdy, C.: Fast, flexible syntactic pattern matching and processing. In: WPC '96: Proceedings of the 4th International Workshop on Program Comprehension, p. 144. IEEE Computer Society, Washington (1996)
55. Grubb, P., Takang, A.A.: *Software Maintenance: Concepts and Practice*. World Scientific (2003)
56. Guéhéneuc, Y.-G., Albin-Amiot, H.: Using design patterns and constraints to automate the detection and correction of inter-class design defects. In: Proceedings of the 39th International Conference and Exhibition on Technology of Object-Oriented Languages and Systems (TOOLS39), TOOLS '01, p. 296. IEEE Computer Society, Washington (2001)

57. Guo, Y., Seaman, C., Zazworka, N., Shull, F.: Domain-specific tailoring of code smells: an empirical study. In: Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering - Volume 2, ICSE '10, pp. 167–170. ACM, New York (2010)
58. Guo, Y., Seaman, C., Gomes, R., Cavalcanti, A., Tonin, G., Da Silva, F.Q.B., Santos, A.L.M., Siebra, C.: Tracking technical debt - an exploratory case study. In: 27th IEEE International Conference on Software Maintenance (ICSM), pp. 528–531 (2011)
59. Harman, M.: The current state and future of search based software engineering. In: International Conference on Software Engineering, pp. 342–357 (2007)
60. Harrison, W., Ossher, H., Sutton, S., Tarr, P.: Concern modeling in the concern manipulation environment. In: Proceedings of the 2005 Workshop on Modeling and Analysis of Concerns in Software, pp. 1–5. ACM Press, New York (2005)
61. Harrold, M.J., Jones, J.A., Li, T., Liang, D., Orso, A., Pennings, M., Sinha, S., Spoon, S.A., Gujarathi, A.: Regression test selection for java software. In: OOPSLA '01: Proceedings of the 16th ACM SIGPLAN Conference on Object Oriented Programming, Systems, Languages, and Applications, pp. 312–326. ACM, New York (2001)
62. Henkel, J., Diwan, A.: Catchup!: capturing and replaying refactorings to support API evolution. In: ICSE '05: Proceedings of the 27th International Conference on Software Engineering, pp. 274–283. ACM, New York (2005)
63. Herzig, K., Zeller, A.: The impact of tangled code changes. In: 2013 10th IEEE Working Conference on Mining Software Repositories (MSR), pp. 121–130. IEEE, Piscataway (2013)
64. Higo, Y., Kamiya, T., Kusumoto, S., Inoue, K.: Refactoring support based on code clone analysis. In: PROFES '04: Proceedings of 5th International Conference on Product Focused Software Process Improvement, Kausai Science City, April 5–8, 2004, pp. 220–233 (2004)
65. Horwitz, S.: Identifying the semantic and textual differences between two versions of a program. In: PLDI '90: Proceedings of the ACM SIGPLAN 1990 Conference on Programming Language Design and Implementation, pp. 234–245. ACM, New York (1990)
66. Horwitz, S., Prins, J., Reps, T.: Integrating noninterfering versions of programs. ACM Trans. Program. Lang. Syst. **11**(3), 345–387 (1989)
67. Hotta, K., Higo, Y., Kusumoto, S.: Identifying, tailoring, and suggesting form template method refactoring opportunities with program dependence graph. In: 2012 16th European Conference on Software Maintenance and Reengineering (CSMR), pp. 53–62. IEEE, Piscataway (2012)
68. Hou, D., Yao, X.: Exploring the intent behind API evolution: a case study. In: Proceedings of the 2011 18th Working Conference on Reverse Engineering, WCRE '11, pp. 131–140. IEEE Computer Society, Washington (2011)
69. Hunt, J.W., Szymanski, T.G.: A fast algorithm for computing longest common subsequences. Commun. ACM **20**(5), 350–353 (1977)
70. ISO/IEC 14764:2006: Software engineering software life cycle processes maintenance. Technical report, ISO/IEC (2006)
71. Izurieta, C., Bieman, J.M.: How software designs decay: a pilot study of pattern evolution. In: First International Symposium on ESEM, pp. 449–451 (2007)
72. Jablonski, P., Hou, D.: CREN: a tool for tracking copy-and-paste code clones and renaming identifiers consistently in the IDE. In: Proceedings of the 2007 OOPSLA Workshop on Eclipse Technology eXchange, eclipse '07, pp. 16–20. ACM, New York (2007)
73. Jackson, D., Ladd, D.A.: Semantic diff: a tool for summarizing the effects of modifications. In: ICSM '94: Proceedings of the International Conference on Software Maintenance, pp. 243–252. IEEE Computer Society, Washington (1994)
74. Janssen, T., Abreu, R., Gemund, A.: Zoltar: a toolset for automatic fault localization. In: Proc. of ASE, pp. 662–664. IEEE Computer Society, Washington (2009)
75. Javassist. <http://jboss-javassist.github.io/javassist/>
76. Jiang, L., Misherghi, G., Su, Z., Glondu, S.: Deckard: scalable and accurate tree-based detection of code clones. In: ICSE '07: Proceedings of the 29th International Conference on Software Engineering, pp. 96–105. IEEE Computer Society, Washington (2007)

77. Jiang, L., Su, Z., Chiu, E.: Context-based detection of clone-related bugs. In: ESEC-FSE '07: Proceedings of the the 6th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, pp. 55–64. ACM, New York (2007)
78. Johnson, P.M.: Reengineering inspection. *Commun. ACM* **41**(2), 49–52 (1998)
79. Johnson, R.: Beyond behavior preservation. Microsoft Faculty Summit 2011, Invited Talk, July 2011
80. Jones, J.A., Harrold, M.J., Stasko, J.: Visualization of test information to assist fault localization. In: Proceedings of the 24th International Conference on Software Engineering, ICSE '02, pp. 467–477. ACM, New York (2002)
81. Juergens, E., Deissenboeck, F., Hummel, B., Wagner, S.: Do code clones matter? In: Proceedings of the 31st International Conference on Software Engineering, ICSE '09, pp. 485–495. IEEE Computer Society, Washington (2009)
82. Juillerat, N., Hirsbrunner, B.: Toward an implementation of the “form template method” refactoring. In: SCAM 2007. Seventh IEEE International Working Conference on Source Code Analysis and Manipulation, pp. 81–90. IEEE, Piscataway (2007)
83. Kataoka, Y., Notkin, D., Ernst, M.D., Griswold, W.G.: Automated support for program refactoring using invariants. In: Proceedings of the IEEE International Conference on Software Maintenance (ICSM'01), ICSM '01, pp. 736. IEEE Computer Society, Washington (2001)
84. Kataoka, Y., Imai, T., Andou, H., Fukaya, T.: A quantitative evaluation of maintainability enhancement by refactoring. In: Proceedings of the International Conference on Software Maintenance (ICSM 2002), pp. 576–585. IEEE Computer Society, Washington (2002)
85. Kawrykow, D., Robillard, M.P.: Non-essential changes in version histories. In: Proceedings of the 33rd International Conference on Software Engineering, ICSE '11, pp. 351–360. ACM, New York (2011)
86. Kiczales, G., Hilsdale, E., Hugunin, J., Kersten, M., Palm, J., Griswold, W.G.: An overview of AspectJ. In: Proceedings of the 15th European Conference on Object-Oriented Programming, ECOOP '01, pp. 327–353. Springer, London (2001)
87. Kim, M., Notkin, D.: Discovering and representing systematic code changes. In: Proceedings of the 31st International Conference on Software Engineering, ICSE '09, pp. 309–319. IEEE Computer Society, Washington (2009)
88. Kim, M., Sazawal, V., Notkin, D., Murphy, G.: An empirical study of code clone genealogies. In: Proceedings of the 10th European Software Engineering Conference Held Jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering, ESEC/FSE-13, pp. 187–196. ACM, New York (2005)
89. Kim, S., Pan, K., James Whitehead, J.E.: When functions change their names: automatic detection of origin relationships. In: WCRE '05: Proceedings of the 12th Working Conference on Reverse Engineering, pp. 143–152. IEEE Computer Society, Washington (2005)
90. Kim, S., Pan, K., Whitehead, E.E.J. Jr.: Memories of bug fixes. In: Proceedings of the 14th ACM SIGSOFT International Symposium on Foundations of Software Engineering, SIGSOFT '06/FSE-14, pp. 35–45. ACM, New York (2006)
91. Kim, M., Notkin, D., Grossman, D.: Automatic inference of structural changes for matching across program versions. In: ICSE '07: Proceedings of the 29th International Conference on Software Engineering, pp. 333–343. IEEE Computer Society, Washington (2007)
92. Kim, M., Gee, M., Loh, A., Rachatasumrit, N.: Ref-finder: a refactoring reconstruction tool based on logic query templates. In: FSE '10: Proceedings of the Eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering, pp. 371–372. ACM, New York (2010)
93. Kim, M., Cai, D., Kim, S.: An empirical investigation into the role of refactorings during software evolution. In: ICSE' 11: Proceedings of the 2011 ACM and IEEE 33rd International Conference on Software Engineering (2011)

94. Kim, M., Zimmermann, T., Nagappan, N.: A field study of refactoring challenges and benefits. In: Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering, FSE '12, pp. 50:1–50:11. ACM, New York (2012)
95. Kim, D., Nam, J., Song, J., Kim, S.: Automatic patch generation learned from human-written patches. In: IEEE/ACM International Conference on Software Engineering (2013)
96. Kim, M., Zimmermann, T., Nagappan, N.: An empirical study of refactoring challenges and benefits at microsoft. *IEEE Trans. Softw. Eng.* **40**(7), 633–649 (2014)
97. Kolb, R., Muthig, D., Patzke, T., Yamauchi, K.: Refactoring a legacy component for reuse in a software product line: a case study: practice articles. *J. Softw. Maint. Evol.* **18**, 109–132 (2006)
98. Komondoor, R., Horwitz, S.: Semantics-preserving procedure extraction. In: POPL '00: Proceedings of the 27th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, pp. 155–169. ACM Press, New York (2000)
99. Komondoor, R., Horwitz, S.: Effective, automatic procedure extraction. In: IWPC '03: Proceedings of the 11th IEEE International Workshop on Program Comprehension, p. 33. IEEE Computer Society, Washington (2003)
100. Koni-N'Sapu, G.G.: A scenario based approach for refactoring duplicated code in object-oriented systems. Master's thesis, University of Bern, June 2001
101. Krishnan, G.P., Tsantalis, N.: Refactoring clones: an optimization problem. In: Proceedings of the ICSM, pp. 360–363 (2013)
102. Ladd, D.A., Ramming, J.C.: A\*: a language for implementing language processors. *IEEE Trans. Softw. Eng.* **21**(11), 894–901 (1995)
103. Lammel, R., Saraiva, J., Visser, J. (eds.): Generative and Transformational Techniques in Software Engineering IV, International Summer School, GTTSE 2011, Braga, July 3–9, 2011. Revised Papers. Lecture Notes in Computer Science, vol. 7680. Springer, Berlin (2013)
104. Landauer, J., Hirakawa, M.: Visual AWK: a model for text processing by demonstration. In: Proceedings of the 11th International IEEE Symposium on Visual Languages, VL '95, p. 267. IEEE Computer Society, Washington (1995)
105. Laski, J., Szermer, W.: Identification of program modifications and its applications in software maintenance. In: ICSM 1992: Proceedings of International Conference on Software Maintenance (1992)
106. Lau, T., Wolfman, S.A., Domingos, P., Weld, D.S.: Learning Repetitive Text-Editing Procedures with SMARTedit, pp. 209–226. Morgan Kaufmann, San Francisco (2001)
107. Le Goues, C., Dewey-Vogt, M., Forrest, S., Weimer, W.: A systematic study of automated program repair: fixing 55 out of 105 bugs for \$8 each. In: International Conference on Software Engineering, pp. 3–13 (2012)
108. Lehman, M.M.: On understanding laws, evolution, and conservation in the large-program life cycle. *J. Syst. Softw.* **1**, 213–221 (1984)
109. Li, Z., Lu, S., Myagmar, S., Zhou, Y.: CP-miner: a tool for finding copy-paste and related bugs in operating system code. In: Proceedings of the 6th Conference on Symposium on Operating Systems Design & Implementation - Volume 6, OSDI'04, pp. 20–20. USENIX Association, Berkeley (2004)
110. Li, Z., Lu, S., Myagmar, S., Zhou, Y.: CP-miner: finding copy-paste and related bugs in large-scale software code. *IEEE Trans. Softw. Eng.* **32**(3), 176–192 (2006)
111. Li, Z., Tan, L., Wang, X., Lu, S., Zhou, Y., Zhai, C.: Have things changed now?: An empirical study of bug characteristics in modern open source software. In: Proceedings of the 1st Workshop on Architectural and System Support for Improving Software Dependability, ASID '06, pp. 25–33. ACM, New York (2006)
112. Lo, D., Jiang, L., Budi, A., et al.: Comprehensive evaluation of association measures for fault localization. In: Proceedings of ICSM, pp. 1–10. IEEE, Piscataway (2010)
113. MacCormack, A., Rusnak, J., Baldwin, C.Y.: Exploring the structure of complex software designs: an empirical study of open source and proprietary code. *Manag. Sci.* **52**(7), 1015–1030 (2006)



114. Madhavji, N.H., Ramil, F.J.C., Perry, D.E.: *Software Evolution and Feedback: Theory and Practice*. Wiley, Hoboken (2006)
115. Malpohl, G., Hunt, J.J., Tichy, W.F.: Renaming detection. *Autom. Softw. Eng.* **10**(2), 183–202 (2000)
116. Marinescu, R.: Detection strategies: metrics-based rules for detecting design flaws. In: *Proceedings of the 20th IEEE International Conference on Software Maintenance*, pp. 350–359. IEEE Computer Society, Washington (2004)
117. McDonnell, T., Ray, B., Kim, M.: An empirical study of API stability and adoption in the android ecosystem. In: *2013 29th IEEE International Conference on Software Maintenance (ICSM)*, pp. 70–79 (2013)
118. Meng, N., Kim, M., McKinley, K.S.: Systematic editing: generating program transformations from an example. In: *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '11*, pp. 329–342. ACM, New York (2011)
119. Meng, N., Kim, M., McKinley, K.S.: Lase: locating and applying systematic edits by learning from examples. In: *Proceedings of the 2013 International Conference on Software Engineering, ICSE '13*, pp. 502–511. IEEE Press, Piscataway (2013)
120. Meng, N., Hua, L., Kim, M., McKinley, K.S.: Does automated refactoring obviate systematic editing? In: *Proceedings of the 37th International Conference on Software Engineering - Volume 1, ICSE '15*, pp. 392–402. IEEE Press, Piscataway (2015)
121. Mens, T.: A state-of-the-art survey on software merging. *IEEE Trans. Softw. Eng.* **28**(5), 449–462 (2002)
122. Mens, T., Tourwé, T.: A survey of software refactoring. *IEEE Trans. Softw. Eng.* **30**(2), 126–139 (2004)
123. Mens, T., Van Eetvelde, N., Demeyer, S., Janssens, D.: Formalizing refactorings with graph transformations. *J. Softw. Maint. Evol. Res. Pract.* **17**(4), 247–276 (2005)
124. Miller, R.C., Myers, B.A.: Interactive simultaneous editing of multiple text regions. In: *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, pp. 161–174. USENIX Association, Berkeley (2001)
125. Moha, N., Guéhéneuc, Y.-G., Meur, A.-F.L., Duchien, L.: A domain analysis to specify design defects and generate detection algorithms. In: *Fiadeiro, J.L., Inverardi, P. (eds.) International Conference on FASE, vol. 4961. Lecture Notes in Computer Science*, pp. 276–291. Springer, Berlin (2008)
126. Moser, R., Sillitti, A., Abrahamsson, P., Succi, G.: Does refactoring improve reusability? In: *Proceedings of ICSR*, pp. 287–297 (2006)
127. Mossienko, M.: Automated Cobol to Java recycling. In: *Proceedings Seventh European Conference on Software Maintenance and Reengineering* (2003)
128. Muchnick, S.S.: *Advanced Compiler Design and Implementation*. Morgan Kaufmann, San Francisco (1997)
129. Murphy, G.C., Kersten, M., Findlater, L.: How are Java Software Developers Using the Eclipse IDE? vol. 23, pp. 76–83. IEEE Computer Society Press, Los Alamitos (2006)
130. Murphy-Hill, E., Parnin, C., Black, A.P.: How we refactor, and how we know it. *IEEE Trans. Softw. Eng.* **38**(1), 5–18 (2012)
131. Nagappan, N., Ball, T.: Use of relative code churn measures to predict system defect density. In: *ICSE '05: Proceedings of the 27th International Conference on Software Engineering*, pp. 284–292. ACM, New York (2005)
132. Naish, L., Lee, H., Ramamohanarao, K.: A model for spectra-based software diagnosis. *ACM TOSEM* **20**(3), 11 (2011)
133. Nguyen, T.T., Nguyen, H.A., Pham, N.H., Al-Kofahi, J.M., Nguyen, T.N.: Graph-based mining of multiple object usage patterns. In: *ESEC/FSE '09: Proceedings of the 7th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pp. 383–392. ACM, New York (2009)

134. Nguyen, H.A., Nguyen, T.T., Wilson, G. Jr., Nguyen, A.T., Kim, M., Nguyen, T.N.: A graph-based approach to API usage adaptation. In: Proceedings of the ACM International Conference on Object Oriented Programming Systems Languages and Applications, OOP-SLA '10, pp. 302–321. ACM, New York (2010)
135. Nguyen, A.T., Nguyen, H.A., Nguyen, T.T., Nguyen, T.N.: Statistical learning approach for mining API usage mappings for code migration. In: Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering, pp. 457–468. ACM, New York (2014)
136. Nguyen, A.T., Nguyen, T.T., Nguyen, T.N.: Divide-and-conquer approach for multi-phase statistical migration for source code (t). In: 2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE) (2015)
137. Nguyen, T.D., Nguyen, A.T., Phan, H.D., Nguyen, T.N.: Exploring API embedding for API usages and applications. In: Proceedings of the 39th International Conference on Software Engineering, ICSE '17, pp. 438–449. IEEE Press, Piscataway (2017)
138. Nix, R.: Editing by example. In: Proceedings of the 11th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages, POPL '84, pp. 186–195. ACM, New York (1984)
139. Ohst, D., Welle, M., Kelter, U.: Difference tools for analysis and design documents. In: International Conference on ICSM '03, p. 13. IEEE Computer Society, Washington (2003)
140. Opdyke, W.F.: Refactoring object-oriented frameworks. PhD thesis, Champaign (1992). UMI Order No. GAX93-05645
141. Orso, A., Shi, N., Harrold, M.J.: Scaling regression testing to large software systems. In: SIGSOFT '04/FSE-12: Proceedings of the 12th ACM SIGSOFT Twelfth International Symposium on Foundations of Software Engineering, pp. 241–251. ACM, New York (2004)
142. Overbey, J.L., Foltzler, M.J., Kasza, A.J., Johnson, R.E.: A collection of refactoring specifications for fortran 95. In: ACM SIGPLAN Fortran Forum, vol. 29, pp. 11–25. ACM, New York (2010)
143. Padioleau, Y., Lawall, J.L., Muller, G.: Understanding collateral evolution in linux device drivers. In: Proceedings of the 1st ACM SIGOPS/EuroSys European Conference on Computer Systems 2006, EuroSys '06, pp. 59–71. ACM, New York (2006)
144. Padioleau, Y., Lawall, J., Hansen, R.R., Muller, G.: Documenting and automating collateral evolutions in linux device drivers. In: Proceedings of the 3rd ACM SIGOPS/EuroSys European Conference on Computer Systems 2008, Eurosys '08, pp. 247–260. ACM, New York (2008)
145. Perry, D.E., Siy, H.P., Votta, L.G.: Parallel changes in large-scale software development: an observational case study. *ACM Trans. Softw. Eng. Methodol.* **10**(3), 308–337 (2001)
146. Pmd: <http://pmd.sourceforge.net/>
147. Prete, K., Rachatasumrit, N., Sudan, N., Kim, M.: Template-based reconstruction of complex refactorings. In: 2010 IEEE International Conference on Software Maintenance (ICSM), pp. 1–10. IEEE Press, Piscataway (2010)
148. Purushothaman, R., Perry, D.E.: Toward understanding the rhetoric of small source code changes. *IEEE Trans. Softw. Eng.* **31**(6), 511–526 (2005)
149. Rachatasumrit, N., Kim, M.: An empirical investigation into the impact of refactoring on regression testing. In: ICSM '12: the 28th IEEE International Conference on Software Maintenance, p. 10. IEEE Society, Washington (2012)
150. Ratzinger, J., Fischer, M., Gall, H.: Improving evolvability through refactoring. In: MSR '05 Proceedings of the 2005 International Workshop on Mining Software Repositories, pp. 1–5 (2005)
151. Ratzinger, J., Sigmund, T., Gall, H.C.: On the relation of refactorings and software defect prediction. In: MSR '08: Proceedings of the 2008 International Working Conference on Mining Software Repositories, pp. 35–38. ACM, New York (2008)
152. Ray, B., Kim, M.: A case study of cross-system porting in forked projects. In: Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering, FSE '12, pp. 53:1–53:11. ACM, New York (2012)

153. Ray, B., Kim, M., Person, S., Rungta, N.: Detecting and characterizing semantic inconsistencies in ported code. In: 2013 IEEE/ACM 28th International Conference on Automated Software Engineering (ASE), pp. 367–377 (2013)
154. Ren, X., Shah, F., Tip, F., Ryder, B.G., Chesley, O.: Chianti: a tool for change impact analysis of java programs. In: OOPSLA '04: Proceedings of the 19th annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, pp. 432–448. ACM, New York (2004)
155. Rigby, P.C., Bird, C.: Convergent contemporary software peer review practices. In: Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering, pp. 202–212. ACM, New York (2013)
156. Rigby, P.C., German, D.M., Storey, M.-A.: Open source software peer review practices: a case study of the apache server. In: ICSE '08: Proceedings of the 30th International Conference on Software Engineering, pp. 541–550. ACM, New York (2008)
157. Robbes, R., Lanza, M.: Spyware: a change-aware development toolset. In: ICSE '08: Proceedings of the 30th International Conference on Software Engineering, pp. 847–850. ACM, New York (2008)
158. Robbes, R., Lungu, M., Röthlisberger, D.: How do developers react to API deprecation?: The case of a smalltalk ecosystem. In: Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering, FSE '12, pp. 56:1–56:11. ACM, New York (2012)
159. Roberts, D., Opdyke, W., Beck, K., Fowler, M., Brant, J.: Refactoring: Improving the Design of Existing Code. Addison-Wesley Longman Publishing Co., Inc., Boston (1999)
160. Robillard, M.P., Murphy, G.C.: Feat: a tool for locating, describing, and analyzing concerns in source code. In: ICSE '03: Proceedings of the 25th International Conference on Software Engineering, pp. 822–823. IEEE Computer Society, Washington (2003)
161. Rolim, R., Soares, G., D'Antoni, L., Polozov, O., Gulwani, S., Gheyi, R., Suzuki, R., Hartmann, B.: Learning syntactic program transformations from examples. In: Proceedings of the 39th International Conference on Software Engineering, ICSE '17, pp. 404–415. IEEE Press, Piscataway (2017)
162. Rothermel, G., Harrold, M.J.: A safe, efficient regression test selection technique. *ACM Trans. Softw. Eng. Methodol.* **6**(2), 173–210 (1997)
163. Schaefer, M., de Moor, O.: Specifying and implementing refactorings. In: Proceedings of the ACM International Conference on Object Oriented Programming Systems Languages and Applications, OOPSLA '10, pp. 286–301. ACM, New York (2010)
164. Schmidt, M., Gloetznert, T.: Constructing difference tools for models using the sidiff framework. In: ICSE Companion '08: Companion of the 30th International Conference on Software Engineering, pp. 947–948. ACM, New York (2008)
165. Shao, D., Khurshid, S., Perry, D.: Evaluation of semantic interference detection in parallel changes: an exploratory experiment. In: ICSM 2007. IEEE International Conference on Software Maintenance, pp. 74–83 (2007)
166. Shepherd, D., Fry, Z.P., Hill, E., Pollock, L., Vijay-Shanker, K.: Using natural language program analysis to locate and understand action-oriented concerns. In: AOSD '07: Proceedings of the 6th International Conference on Aspect-Oriented Software Development, pp. 212–224. ACM, New York (2007)
167. Sidiroglou, S., Ioannidis, S., Keromytis, A.D.: Band-aid patching. In: Proceedings of the 3rd Workshop on on Hot Topics in System Dependability, HotDep'07. USENIX Association, Berkeley (2007)
168. Silva, D., Tsantalis, N., Valente, M.T.: Why we refactor? confessions of Github contributors. In: Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2016, pp. 858–870. ACM, New York (2016)
169. Śliwerski, J., Zimmermann, T., Zeller, A.: When do changes induce fixes? In: Proceedings of the 2005 International Workshop on Mining Software Repositories, MSR '05, pp. 1–5. ACM, New York (2005)

170. Sneed, H.M.: Migrating from COBOL to Java. In: Proceedings of the 2010 IEEE International Conference on Software Maintenance (2010)
171. Soares, G.: Making program refactoring safer. In: Proceedings of the 32Nd ACM/IEEE International Conference on Software Engineering - Volume 2, ICSE '10, pp. 521–522 (2010)
172. Software Maintenance and Computers (IEEE Computer Society Press Tutorial). IEEE Computer Society, Los Alamitos (1990)
173. Son, S., McKinley, K.S., Shmatikov, V.: Fix me up: repairing access-control bugs in web applications. In: NDSS Symposium (2013)
174. Soto, M., Münch, J.: Process Model Difference Analysis for Supporting Process Evolution. Lecture Notes in Computer Science, vol. 4257, pp. 123–134. Springer, Berlin (2006)
175. Sullivan, K., Chalasani, P., Sazawal, V.: Software design as an investment activity: a real options perspective. Technical report (1998)
176. Swanson, E.B.: The dimensions of maintenance. In: Proceedings of the 2Nd International Conference on Software Engineering, ICSE '76, pp. 492–497. IEEE Computer Society Press, Los Alamitos (1976)
177. Tahvildari, L., Kontogiannis, K.: A metric-based approach to enhance design quality through meta-pattern transformations. In: Proceedings of the Seventh European Conference on Software Maintenance and Reengineering, CSMR '03, p. 183. IEEE Computer Society, Washington (2003)
178. Tairas, R., Gray, J.: Increasing clone maintenance support by unifying clone detection and refactoring activities. *Inf. Softw. Technol.* **54**(12), 1297–1307 (2012)
179. Tao, Y., Kim, S.: Partitioning composite code changes to facilitate code review. In: 2015 IEEE/ACM 12th Working Conference on Mining Software Repositories (MSR), pp. 180–190. IEEE, Piscataway (2015)
180. Tarr, P., Osher, H., Harrison, W., Sutton, J.S.M.: N degrees of separation: multi-dimensional separation of concerns. In: ICSE '99: Proceedings of the 21st International Conference on Software Engineering, pp. 107–119. IEEE Computer Society Press, Los Alamitos (1999)
181. The AspectJ Project. <https://eclipse.org/aspectj/>
182. The Guided Tour of TXL. <https://www.txl.ca/tour/tour1.html>
183. Tichy, W.F.: The string-to-string correction problem with block moves. *ACM Trans. Comput. Syst.* **2**(4), 309–321 (1984)
184. Toomim, M., Begel, A., Graham, S.L.: Managing duplicated code with linked editing. In: VLHCC '04: Proceedings of the 2004 IEEE Symposium on Visual Languages - Human Centric Computing, pp. 173–180. IEEE Computer Society, Washington (2004)
185. Treude, C., Berlik, S., Wenzel, S., Kelter, U.: Difference computation of large models. In: Proceedings of the the 6th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, ESEC-FSE '07, pp. 295–304. ACM, New York (2007)
186. Tsantalis, N., Chatzigeorgiou, A.: Identification of extract method refactoring opportunities. In: CSMR '09: Proceedings of the 2009 European Conference on Software Maintenance and Reengineering, pp. 119–128. IEEE Computer Society, Washington (2009)
187. Tsantalis, N., Chatzigeorgiou, A.: Identification of move method refactoring opportunities. *IEEE Trans. Softw. Eng.* **35**(3), 347–367 (2009)
188. Tsantalis, N., Chatzigeorgiou, A.: Identification of extract method refactoring opportunities for the decomposition of methods. *J. Syst. Softw.* **84**(10), 1757–1782 (2011)
189. Tsantalis, N., Chatzigeorgiou, A.: Ranking refactoring suggestions based on historical volatility. In: 2011 15th European Conference on Software Maintenance and Reengineering, pp. 25–34 (2011)
190. Tsantalis, N., Chaikalas, T., Chatzigeorgiou, A.: Jdeodorant: identification and removal of type-checking bad smells. In: CSMR '08: Proceedings of the 2008 12th European Conference on Software Maintenance and Reengineering, pp. 329–331. IEEE Computer Society, Washington (2008)
191. Vakilian, M., Chen, N., Negara, S., Rajkumar, B.A., Bailey, B.P., Johnson, R.E.: Use, disuse, and misuse of automated refactorings. In: 2012 34th International Conference on Software Engineering (ICSE), pp. 233–243 (2012)

192. van Engelen, R.: On the use of clone detection for identifying crosscutting concern code. *IEEE Trans. Softw. Eng.* **31**(10), 804–818 (2005). Student Member-Magiel Bruntink and Member-Arie van Deursen and Member-Tom Tourwe
193. Visser, E.: Program transformation with Stratego/XT: rules, strategies, tools, and systems in StrategoXT-0.9. *Domain-Specific Program Generation* **3016**, 216–238 (2004)
194. Wang, W., Godfrey, M.W.: Recommending clones for refactoring using design, context, and history. In: 2014 IEEE International Conference on Software Maintenance and Evolution, pp. 331–340 (2014)
195. Wei, Y., Pei, Y., Furia, C.A., Silva, L.S., Buchholz, S., Meyer, B., Zeller, A.: Automated fixing of programs with contracts. In: Proceedings of the 19th International Symposium on Software Testing and Analysis, ISSTA '10, pp. 61–72. ACM, New York (2010)
196. Weißgerber, P., Diehl, S.: Are refactorings less error-prone than other changes? In: MSR '06: Proceedings of the 2006 International Workshop on Mining Software Repositories, pp. 112–118. ACM, New York (2006)
197. Weißgerber, P., Diehl, S.: Identifying refactorings from source-code changes. In: ASE '06: Proceedings of the 21st IEEE/ACM International Conference on Automated Software Engineering, pp. 231–240. IEEE Computer Society, Washington (2006)
198. Weimer, W., Nguyen, T., Le Goues, C., Forrest, S.: Automatically finding patches using genetic programming. In: Proceedings of the 31st International Conference on Software Engineering, ICSE '09, pp. 364–374. IEEE Computer Society, Washington (2009)
199. Wikipedia. Comparison of BSD operating systems — Wikipedia, the free encyclopedia (2012)
200. Wong, S., Cai, Y., Kim, M., Dalton, M.: Detecting software modularity violations. In: ICSE' 11: Proceedings of the 2011 ACM and IEEE 33rd International Conference on Software Engineering (2011)
201. Xing, Z., Stroulia, E.: UMLDiff: an algorithm for object-oriented design differencing. In: ASE '05: Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering, pp. 54–65. ACM, New York (2005)
202. Xing, Z., Stroulia, E.: Refactoring detection based on UMLDiff change-facts queries. In: WCRE '06: Proceedings of the 13th Working Conference on Reverse Engineering, pp. 263–274. IEEE Computer Society, Washington (2006)
203. Xing, Z., Stroulia, E.: Refactoring practice: how it is and how it should be supported - an eclipse case study. In: ICSM '06: Proceedings of the 22nd IEEE International Conference on Software Maintenance, pp. 458–468. IEEE Computer Society, Washington (2006)
204. Xing, Z., Stroulia, E.: API-evolution support with diff-catchup. *IEEE Trans. Softw. Eng.* **33**(12), 818–836 (2007)
205. Yamamoto, T., Matsushita, M., Kamiya, T., Inoue, K.: Measuring similarity of large software systems based on source code correspondence. In: Proceedings of 2005 Product Focused Software Process Improvement, pp. 530–544 (2005)
206. Yang, W.: Identifying syntactic differences between two programs. *Softw. Pract. Experience* **21**(7), 739–755 (1991)
207. Yang, W., Horwitz, S., Reps, T.: Detecting program components with equivalent behaviors. Technical Report CS-TR-1989-840, University of Wisconsin, Madison (1989)
208. Yasumatsu, K., Doi, N.: SPICE: a system for translating Smalltalk programs into a C environment. *IEEE Trans. Softw. Eng.* **21**(11), 902–912 (1995)
209. Yin, Z., Yuan, D., Zhou, Y., Pasupathy, S., Bairavasundaram, L.: How do fixes become bugs? In: Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering, ESEC/FSE '11, pp. 26–36. ACM, New York (2011)
210. Yokomori, R., Siy, H.P., Noro, M., Inoue, K.: Assessing the impact of framework changes using component ranking. In: Proceedings of ICSM, pp. 189–198. IEEE, Piscataway (2009)
211. Zeller, A.: Yesterday, my program worked. today, it does not. Why? In: ESEC/FSE-7: Proceedings of the 7th European Software Engineering Conference Held Jointly with the 7th ACM SIGSOFT International Symposium on Foundations of Software Engineering, pp. 253–267. Springer, London (1999)

212. Zeller, A.: Automated debugging: are we close? *IEEE Comput.* **34**(11), 26–31 (2001)
213. Zhang, L., Kim, M., Khurshid, S.: Localizing failure-inducing program edits based on spectrum information. In: *Proceedings of ICSM*, pp. 23–32. IEEE, Piscataway (2011)
214. Zhang, T., Song, M., Pinedo, J., Kim, M.: Interactive code review for systematic changes. In: *Proceedings of the 37th International Conference on Software Engineering-Volume 1*, pp. 111–122. IEEE Press, Piscataway (2015)
215. Zhong, H., Thummalapenta, S., Xie, T., Zhang, L., Wang, Q.: Mining API mapping for language migration. In: *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 1*, pp. 195–204. ACM, New York (2010)
216. Zou, L., Godfrey, M.W.: Using origin analysis to detect merging and splitting of source code entities. *IEEE Trans. Softw. Eng.* **31**(2), 166–181 (2005)