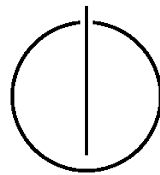


FAKULTÄT FÜR INFORMATIK  
DER TECHNISCHEN UNIVERSITÄT MÜNCHEN

Seminararbeit zum Proseminar Software-Qualität

## **Gapped Clone Detection**

Marius Daniel Schulz



## **Abstract**

In diesem Paper wird die Problematik inkonsistenter Klone in Software-Quellcode beleuchtet. Es wird diskutiert, welche Inkonsistenzen es gibt und welche Auswirkungen diese zur Folge haben. Weiterhin werden Verfahren zur Erkennung und Entfernung von Klonen dargelegt. Schlussendlich wird anhand einer Analyse von fünf in Java implementierten Open-Source-Systemen aufgezeigt, wie inkonsistente Klone in Real-World-Projekten aussehen können.

# Inhaltsverzeichnis

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Einführung</b>                              | <b>4</b>  |
| 1.1      | Begriffsklärung . . . . .                      | 4         |
| 1.1.1    | Clones . . . . .                               | 4         |
| 1.1.2    | Gapped Clones . . . . .                        | 4         |
| 1.2      | Warum kopieren wir Code? . . . . .             | 4         |
| 1.2.1    | Limitierungen der Programmiersprache . . . . . | 4         |
| 1.2.2    | Cross-Cutting Concerns . . . . .               | 4         |
| 1.2.3    | Anfangs unklare Abstraktionen . . . . .        | 5         |
| 1.2.4    | Weitere Gründe . . . . .                       | 5         |
| <b>2</b> | <b>Motivation &amp; Relevanz</b>               | <b>5</b>  |
| 2.1      | Entstehung von Inkonsistenzen . . . . .        | 5         |
| 2.1.1    | Inkonsistenz durch Vorlagen . . . . .          | 6         |
| 2.1.2    | Inkonsistenz aus Vorsicht . . . . .            | 6         |
| 2.1.3    | Inkonsistenz aus Versehen . . . . .            | 6         |
| 2.2      | Wartungsaufwand und Fehlverhalten . . . . .    | 6         |
| 2.2.1    | Empirische Studien . . . . .                   | 7         |
| <b>3</b> | <b>Erkennung von Klonen</b>                    | <b>7</b>  |
| 3.1      | Klon-Typen . . . . .                           | 7         |
| 3.2      | Precision & Recall . . . . .                   | 8         |
| 3.3      | Techniken zur Klonerkennung . . . . .          | 8         |
| 3.3.1    | Textbasiert . . . . .                          | 8         |
| 3.3.2    | Token-basiert . . . . .                        | 8         |
| 3.3.3    | Syntaxbaum-basiert . . . . .                   | 9         |
| 3.3.4    | Abhängigkeitsgraph-basiert . . . . .           | 9         |
| 3.3.5    | Weitere Techniken . . . . .                    | 9         |
| <b>4</b> | <b>Beseitigung von Klonen</b>                  | <b>10</b> |
| 4.1      | Auslagerung von geteiltem Code . . . . .       | 10        |
| 4.2      | Automatische Code-Generierung . . . . .        | 10        |
| 4.3      | Kosten-Nutzen-Verhältnis . . . . .             | 10        |
| <b>5</b> | <b>Analysierte Projekte</b>                    | <b>11</b> |
| <b>6</b> | <b>Anhang: Beispielhafte Klone</b>             | <b>13</b> |

# 1 Einführung

## 1.1 Begriffsklärung

### 1.1.1 Clones

Dem Klonbegriff dieses Papers liegt die Definition von Bettenburg et al. zugrunde, wie er in [3] verwendet wird:

*A code clone is a part of the source code that is identical, or at least highly similar, to another part (clone) in terms of structure and semantics.*

Da die Frage nach semantischer Äquivalenz i.a. nicht entscheidbar ist, interessieren uns an dieser Stelle lediglich diejenigen Klone, die sich in ihrer Programmdarstellung ähneln (engl. *representational clones*).

### 1.1.2 Gapped Clones

Als *gapped clones* bezeichnet man dabei diejenigen Klone, die sich bis auf geringe Abweichungen ähnlich sind; der Begriff „Klon“ darf daher an dieser Stelle nicht als absolut identische Kopie verstanden werden. Im Abschnitt 3.1 wird im Detail auf mögliche Unterschiede zwischen Code-Fragmenten eingegangen, die hier unter dem recht unspezifischen Begriff „geringe Abweichung“ zusammengefasst sind.

## 1.2 Warum kopieren wir Code?

Es gibt diverse Gründe, die einen Entwickler dazu bewegen, existierenden Code zu kopieren und somit Copy & Paste-Programmierung zu betreiben.

### 1.2.1 Limitierungen der Programmiersprache

Bietet die vorliegende Programmiersprache keine Unterstützung für bestimmte Sprachkonstrukte, mit denen sich Code an anderer Stelle wiederverwenden lässt, ist das Kopieren von Code teilweise die einzige Option des Entwicklers.

Ein konkretes Beispiel für eine Sprache, der solch ein Feature fehlte, ist Java. Vor der im September 2004 veröffentlichten Version 5.0 gab es in Java keine Generics. Sollte damals eine bestimmte Funktionalität auf mehreren Datentypen definiert werden, wie z.B. eine Berechnung auf den numerischen Typen `int` und `long`, so musste sie für jeden Datentyp einzeln implementiert werden.

### 1.2.2 Cross-Cutting Concerns

Als Cross-Cutting Concerns bezeichnet man Belange einer Software, die sich, gemäß ihrem Namen, horizontal durch die Anwendung ziehen. Darunter fallen z.B. die folgenden Aspekte:

- Logging,
- Caching,
- Validierung,
- Autorisierung,
- Lokalisierung,
- Memory Management, uvm.

Weil sich diese Funktionalität teilweise nur schwer an einer Stelle zusammenfassen lässt, wird häufig Copy & Paste eingesetzt, um diese Aspekte an unterschiedlichen Stellen zu implementieren, was verstärkt zu Codeduplizierung führt. Lösungsansätze dafür bietet die Aspektorientierte Programmierung [7].

### 1.2.3 Anfangs unklare Abstraktionen

Bei der Entwicklung eines Softwaremoduls sind häufig nicht alle sinnvollen Abstraktionen von Anfang an offensichtlich. Das Entwicklerteam kann daher die Phase der Abstraktion nach hinten schieben, bis sich Muster in der Verwendung der Komponente abzeichnen. So lässt sich vorzeitiges *Overengineering* vermeiden, das die Komplexität der Software unnötig erhöht hätte. Der Preis dafür ist die (temporäre) Redundanz, die der duplizierte Code mit sich bringt, mit all ihren Nachteilen, die im folgenden behandelt werden.

### 1.2.4 Weitere Gründe

Die bisher genannten Gründe waren allesamt technischer Natur, jedoch gibt es darüber hinaus eine Reihe an weiteren Faktoren, die zur Duplizierung von Code beitragen.

Steht das Entwicklungsteam eines Softwareprojektes unter akutem Zeitdruck, kann es als Zwischenlösung lukrativ erscheinen, Code zu duplizieren, anstatt auf saubere Wiederverwendungsmechanismen zu setzen. Findet später kein Refactoring statt, häuft sich *technical debt* in Form von Code-Klonen an. Aus eigener Projekterfahrung kann ich bestätigen, dass solche Provisorien zu ausgeprägter Langlebigkeit tendieren.

Koschke [10] führt an, dass fragwürdige Produktivitätsmetriken von Entwicklern, wie z.B. die Anzahl geschriebener Code-Zeilen, explizit zum Duplizieren von Code verleiten. Schlussendlich ist es ebenfalls denkbar, dass Code Clones aufgrund von Faulheit, Gleichgültigkeit oder mangelndem Bewusstsein der Entwickler stattfindet.

## 2 Motivation & Relevanz

### 2.1 Entstehung von Inkonsistenzen

Dupliziert ein Entwickler ein Stück Code, so ist das kopierte Fragment nun an einer weiteren Stelle redundant vorhanden. Damit unterscheidet es sich aber noch nicht vom Ausgangscode, ist also noch eine konsistente Kopie. Eine Reihe an Ursachen kann jedoch für die

Entstehung von Inkonsistenzen sorgen, wie z.B. Code-Vorlagen, Vorsicht des Entwicklers und Versehen.

### 2.1.1 Inkonsistenz durch Vorlagen

Die Wiederverwendung von existierendem Code per Copy & Paste ist laut Kim et al. [9] die häufigste Ursache fürs Klonen. Dies trifft insbesondere auf Schnittstellen von Funktionsbibliotheken zu, bei deren Verwendung sehr ähnliche Aufrufsmuster entstehen. Nimmt dieser Boilerplate-Code größere Ausmaße an, deutet dies vermutlich auf ein mangelhaftes Design der Schnittstelle hin [11].

Häufig wird jedoch der kopierte Code nicht 1:1 übernommen, sondern abgeändert und an die neuen Bedürfnisse angepasst; es entsteht ein inkonsistenter Klon. Das kopierte Code-Fragment dient demnach als Vorlage, weshalb man von *templating* spricht.

### 2.1.2 Inkonsistenz aus Vorsicht

Auch die Vorsicht eines Entwicklers kann Ursache für inkonsistente Klone sein. Muss ein Entwickler beispielsweise eine Änderung an einem ihm fremden Legacy-System vornehmen, sind ihm oftmals die Auswirkungen seiner Änderung unbekannt. Häufig besitzen diese Systeme keine Tests, was die Unsicherheit noch befördert. In solchen Fällen kann sich der Entwickler dazu entscheiden, seine Anpassungen in einer Kopie vorzunehmen, damit der existierende Code unangetastet bleibt und Breaking Changes vermieden werden.

### 2.1.3 Inkonsistenz aus Versehen

Sowohl bei der Wiederverwendung durch Vorlagen als auch bei der Code-Duplizierung und -Anpassung aus Vorsicht hat der Entwickler inkonsistente Klone bewusst in Kauf genommen. Jedoch kann es leicht passieren, ungewollte Inkonsistenzen herbeizuführen, indem bei einer Änderung an einem Klon nicht alle Instanzen seiner sogenannten *Klongruppe* angepasst werden. Der Entwickler hat so mangels Kenntnis der anderen Ausprägungen dieses Klons unbewusst inkonsistenten Code erzeugt.

## 2.2 Wartungsaufwand und Fehlverhalten

Klone im Allgemeinen führen zu erhöhtem Wartungsaufwand [10], da jede Ausprägung eines Klons gefunden und abgeändert werden muss, um Inkonsistenzen zu vermeiden. Da in der Regel jedoch nicht dokumentiert ist, wo Code geklont worden ist [11], müssen Klone zur Wartungszeit gefunden werden. Schon bei mittelmäßig großen Projekten ist das manuell nicht mehr sinnvoll machbar, sodass automatisierte Unterstützung durch Tooling benötigt wird [11].

Wird die Anpassung bei einem oder mehreren Klonen der Klongruppe vergessen oder übersehen, kann die entstandene Inkonsistenz leicht zu unerwartetem Verhalten des Pro-

gramms führen: Ein Teil der Logik des Systems ist nun an verschiedenen Stellen unterschiedlich implementiert – eine Steilvorlage für Bugs und damit zentraler Bestandteil des Problems.

### 2.2.1 Empirische Studien

Dass das bei weitem kein seltener Fall ist, zeigen Juergens et al. eindrucksvoll in [8]. In ihrer Studie haben die Autoren fünf Software-Systeme untersucht, die in C#, Java und COBOL implementiert wurden. Im Abschnitt „Results“ schlüsseln sie auf, wie viele Klongruppen die Systeme enthielten, welche davon inkonsistent waren und welche wiederum davon unbeabsichtigt waren.

Die beeindruckendste Ziffer schließlich ist der Prozentsatz der unbeabsichtigt inkonsistenten Klone (orig. *unintentionally inconsistent clones*), die zu einem Fehlverhalten (einem *Fault*) des Programms führten: Laut Studie beläuft sich dieser dort als  $\frac{|F|}{|UIC|}$  bezeichnete Wert auf 50%. Intuitiv formuliert heißt das, dass duplizierter Code, der unabsichtlich von seinen anderen Ausprägungen abweicht, in jedem zweiten Fall Ursache für ungewolltes Fehlverhalten ist. Nach ihren Untersuchungen schlussfolgern die Autoren, dass der Anteil an Fehlern in inkonsistenten Klonen höher ist als im Schnitt und bekräftigen somit die Relevanz der Klonfindung.

Andere Autoren auf dem Gebiet bestätigen dieses Ergebnis. Monden et al. [12] berichten von einer höheren Fehlerdichte in Software-Modulen mit großen Klonen. Chou et al. [4] stellten fest, dass Code mit einem Fehler überdurchschnittlich häufig weitere Fehler aufweist, wenn er Klone enthält.

Zusammenfassend lässt sich festhalten, dass Klone – allen voran unbeabsichtigt inkonsistente Klone – einen negativen Einfluss auf die Qualität einer Software haben. Sie machen die Wartung einer Anwendung teurer und erhöhen die Wahrscheinlichkeit, dass sich dabei Fehler einschleichen. Im Sinne der Qualitätssicherung ist es daher im Interesse der Entwickler, Klone im Quellcode automatisiert zu finden.

## 3 Erkennung von Klonen

### 3.1 Klon-Typen

Bevor wir uns verschiedenen Techniken zur Klonerkennung widmen, müssen wir uns bewusst machen, was es für Arten von Klonen gibt und worin sie sich unterscheiden. In der Literatur [10] werden Klone in Software-Quellcode wie folgt kategorisiert:

**Typ 1:** Exakte Kopie bis auf Whitespace, Formatierung und Kommentare.

**Typ 2:** Syntaktisch identische Kopie, bei der Literalwerte abgeändert oder Bezeichner von Variablen, Typen und Funktionen umbenannt wurden.

**Typ 3:** Kopie mit eingefügten, veränderten oder gelöschten Statements.

Als *gapped clones* bezeichnet man Klone vom Typ 3. Das „gapped“ im Namen steht dabei für die Abweichung der Kopien voneinander, also die symbolische „Lücke“.

## 3.2 Precision & Recall

Bei der Betrachtung verschiedener Techniken zur Klonerkennung ist von Interesse, welche Genauigkeit und Trefferquote sie aufweisen. In der englischsprachigen Literatur werden diese Aspekte als *precision* und *recall* bezeichnet. Während *precision* den Anteil relevanter Ergebnisse an allen Suchergebnissen angibt, misst *recall* die Vollständigkeit der gefundenen Resultate. Beim Nachschlagen in der Literatur begegnet man häufig diesen beiden Begriffen, weshalb sie hier der Vollständigkeit halber erwähnt seien.

## 3.3 Techniken zur Klonerkennung

### 3.3.1 Textbasiert

Beim Textvergleich werden zwei Klonkandidaten zeilenweise miteinander verglichen. Damit minimale Änderungen an der Darstellung eines Code-Fragments nicht sofort zum Scheitern des String-Vergleichs führen, werden beim Vergleich häufig Kommentare und Whitespace ignoriert.

Da String-Vergleiche generell vergleichsweise teuer sind, kann zur Steigerung der Performance eine Hash-Funktion eingesetzt werden. Diese berechnet für jede Zeile einen Hash und unterteilt somit die Menge aller Zeilen in verschiedene Partitionen mit dem gleichen Hash-Wert. Der anschließende String-Vergleich wird lediglich auf Zeilen innerhalb der gleichen Partition angewendet [10].

### 3.3.2 Token-basiert

Beim Token-basierten Ansatz sind nicht mehr die einzelnen Zeichen zweier Code-Fragmente die kleinste Einheit, die zum Vergleich herangezogen wird, sondern deren Tokens [1]. Genau wie beim textbasierten Ansatz findet hierbei i.d.R. eine Normalisierung statt, sodass nicht die konkreten Werte eines Tokens betrachtet werden, sondern nur dessen Identität. Damit wird die Token-basierte Erkennung robust gegen Modifizierungen von Whitespace oder Kommentaren, Umbenennung von Identifiern und Änderung von Literalwerten.

Zum effizienten Finden gleicher Token-Sequenzen eignet sich als Datenstruktur ein *suffix tree*, der sich in linearer Laufzeit konstruieren lässt [11]. Klonkandidaten mit bestimmter Mindestlänge lassen sich von der Wurzel ausgehend leicht ablesen, da die Pfade zu den Blättern mit den gemeinsamen Präfixen beschriftet sind; diese Präfixe repräsentieren den geklonten Teil der Code-Fragmente und machen so den Klonkandidaten aus.



### 3.3.3 Syntaxbaum-basiert

Wie der Name vermuten lässt, wird beim Syntaxbaum-basierten Ansatz der abstrakte Syntaxbaum (engl. AST für *abstract syntax tree*) eines Programms erstellt. Die Robustheit der Token-basierten Erkennung gegen (strukturierende) Änderungen ist auch hier gegeben. Der eigentliche Vergleich findet hier auf Teilbäumen des Syntaxbaums statt, die im Sinne der Performancesteigerung gehasht werden können.

Anders als beim Text- oder Token-basierten Ansatz ist für den Aufbau des Syntaxbaums ein tieferes Programmverständnis erforderlich: Es wird ein Parser benötigt, der die Grammatik der verwendeten Programmiersprache versteht. Damit der Parser einen Syntaxbaum erstellen kann, muss der Code in kompilierten Sprachen syntaktisch und semantisch korrekt sein; fehlerhafter Code, der nicht kompiliert, kann also nicht mithilfe eines Syntaxbaums auf Klone untersucht werden.

### 3.3.4 Abhängigkeitsgraph-basiert

Das letzte Verfahren, das in diesem Paper betrachtet werden soll, baut auf dem sogenannten Abhängigkeitsgraph (*program dependence graph*, oder kurz PDG) eines Programms auf. Es macht sich zunutze, dass Zeilen bzw. Statements eines Code-Fragmentes i.d.R. nicht beliebig vertauscht werden können, wenn sie gleichzeitig semantisch korrekt und sinnerhaltend bleiben sollen.

Imperative Programme bestehen aus einer Reihe von Statements, von denen manche Abhängigkeiten untereinander aufweisen. So kann in C-ähnlichen Sprachen eine Variable erst nach ihrer Deklaration verwendet werden, aber nicht vorher. Auch bestehen logische Abhängigkeiten zwischen Statements: Wird ein numerischer Parameter einer Methode erst mit einem Faktor multipliziert und anschließend zu einer Konstanten dazuaddiert, liefern die Operationen in umgekehrter Ausführungsreihenfolge i.a. nicht das gleiche Ergebnis. Die zwei entsprechenden Statements können demnach nicht vertauscht werden, ohne die Programmlogik zu ändern.

Ein *program dependence graph* ermittelt nun die Abhängigkeiten zwischen den Statements eines Code-Fragmentes und repräsentiert diese als Graph. Der Vergleich zweier Code-Einheiten wird dann als Suche nach isomorphen Teilgraphen durchgeführt. Dieses Problem ist jedoch NP-schwer, weshalb in der Praxis approximative Verfahren eingesetzt werden [10].

### 3.3.5 Weitere Techniken

Die Liste der genannten Techniken zur Klonerkennung ist natürlich nicht vollständig; so wurden z.B. Metriken-basierte Ansätze nicht besprochen. Es existiert eine Vielzahl an Algorithmen und Verfahren für dieses Gebiet, aber auch für verwandte Gebiete wie die Plagiatserkennung. Diese würden jedoch den Rahmen dieser kurzen Arbeit sprengen und sind daher in der Literatur nachzulesen.

## 4 Beseitigung von Klonen

Eine Strategie zur Entfernung neuer Klone kann es sein, diese gar nicht erst entstehen zu lassen. Balazinska et al. [2] führen dafür Entwurfsmuster als Ansatz an, der die Anzahl an Klonen durch eine bessere Architektur verringert. Dies ist jedoch nur schwer automatisiert machbar und muss daher i.d.R. manuell erledigt werden [10].

### 4.1 Auslagerung von geteiltem Code

Sind die Klone bereits vorhanden, so können diese häufig durch Refactoring-Verfahren wie *Extract Method* oder *Extract Class* beseitigt werden, sofern die verschiedenen Ausprägungen nicht zu sehr voneinander abweichen [6]. Zum Behandeln von Typ-3-Klonen nennt Koschke beispielhaft Makros und Präprozessor-Anweisungen als Lösungsansatz [10].

Aus meiner eigenen Erfahrung möchte ich an dieser Stelle darauf hinweisen, dass Direktiven für einen Präprozessor – insbesondere Verfahren wie die bedingte Kompilierung – mit Vorsicht zu genießen sind. Auf der einen Seite sind sie schwer zu testen, da der Code in allen möglichen Varianten kompiliert werden muss. Auf der anderen Seite steigern sie die Komplexität eines Moduls, sodass dessen Verständlichkeit leidet.

### 4.2 Automatische Code-Generierung

Abschließend ist es möglich, sich codegenerierende Verfahren zunutze zu machen und den redundanten Code automatisiert erzeugen zu lassen. Da die Generierung nicht manuell stattfindet, ist sie nicht anfällig für Flüchtigkeits- oder Tippfehler. Java setzt diese Methode beispielsweise ein, um den Quellcode für verschiedene Buffer-Klassen wie `ByteBuffer` und `CharBuffer` automatisiert zu erstellen [11]. Dafür kommt Xvcl zum Einsatz, die *XML-Based Variant Configuration Language* [13]. Diese manipuliert konfigurierbare Programmkomponenten gemäß einer XML-Beschreibung und macht so die manuelle Anpassung unnötig.

### 4.3 Kosten-Nutzen-Verhältnis

Damit die Ressourcen, die für die Entfernung von Klonen aufgewendet werden, gerechtfertigt werden können, müssen Kosten und Nutzen der Entfernung in einem günstigen Verhältnis stehen. Es muss vorab situationsabhängig entschieden werden, ob im betroffenen Modul zum gegebenen Zeitpunkt eine Entfernung überhaupt gewünscht ist. Wenn bereits geplant ist, dass sich in naher Zukunft beide Kopien unabhängig voneinander in verschiedene Richtungen weiterentwickeln, macht eine Zusammenführung wenig Sinn. Im Falle solch einer Weiche spricht man daher treffend von *forking* [10].

## 5 Analysierte Projekte

Mithilfe von ConQAT [5], dem *Continuous Quality Assessment Toolkit*, wurden für diese Arbeit die folgenden fünf quelloffenen Java-Projekte auf inkonsistente Klone untersucht:

- Art of Illusion (AoI)
- ArgoUML
- FreeCol
- FreeMind
- JUnit

Für die Analyse wurde der ConQAT-Block `JavaGappedCloneAnalysis` gewählt. Zu Beginn wurden für die Parameter *gap ratio*, *clone minlength* und *max errors* die Werte 0.25, 10 und 5 verwendet. Bei der Suche nach Klonen hat sich schnell herausgestellt, dass die Qualität der gefundenen Resultate maßgeblich von den gewählten Parametern abhängt. Während eine vorgeschriebene Mindestlänge von 10 bei AoI gut funktioniert und längere Klone freigelegt hat, ist dieser Wert deutlich zu hoch für JUnit gewesen: Die meisten Methoden in JUnit sind kurz gehalten und weisen daher nur deutlich kürzere Klone auf. Für ein Projekt, das sich der Qualitätssicherung durch Tests verschrieben hat, ist diese Erkenntnis nicht weiter verwunderlich.

Der mit weit über 100 Zeilen längste gefundene Klon war in AoI enthalten. Bei einer experimentellen Erhöhung des Wertes *max errors* auf 25 fand ihn ConQAT mit 24 Gaps. Im Anhang sind dessen Beginn sowie weitere beispielhafte Klone aus den analysierten Projekten zu finden.

## Literatur

- [1] Brenda S Baker. On finding duplication and near-duplication in large software systems. In *Reverse Engineering, 1995., Proceedings of 2nd Working Conference on*, pages 86–95. IEEE, 1995.
- [2] Magdalena Balazinska, Ettore Merlo, Michel Dagenais, Bruno Lague, and Kostas Kontogiannis. Advanced clone-analysis to support object-oriented system refactoring. In *Reverse Engineering, 2000. Proceedings. Seventh Working Conference on*, pages 98–107. IEEE, 2000.
- [3] Nicolas Bettenburg, Weiyi Shang, Walid M Ibrahim, Bram Adams, Ying Zou, and Ahmed E Hassan. An empirical study on inconsistent changes to code clones at the release level. *Science of Computer Programming*, 77(6):760–776, 2012.
- [4] Andy Chou, Junfeng Yang, Benjamin Chelf, Seth Hallem, and Dawson Engler. *An Empirical Study of Operating Systems Errors*, volume 35. ACM, 2001.

- [5] Florian Deissenboeck, Elmar Juergens, Benjamin Hummel, Stefan Wagner, Benedikt Mas y Parareda, and Markus Pizka. Tool support for continuous quality control. *Software, IEEE*, 25(5):60–67, 2008.
- [6] Richard Fanta and Václav Rajlich. Removing clones from the code. *Journal of Software Maintenance*, 11(4):223–243, 1999.
- [7] Robert E. Filman, Siobhan Clarke, and Tzilla Elrad. *Aspect-Oriented Software Development*. Addison-Wesley, 2005.
- [8] Elmar Juergens, Florian Deissenboeck, Benjamin Hummel, and Stefan Wagner. Do code clones matter? In *Proceedings of the 31st International Conference on Software Engineering*, pages 485–495. IEEE Computer Society, 2009.
- [9] Miryung Kim, Lawrence Bergman, Tessa Lau, and David Notkin. An ethnographic study of copy and paste programming practices in oopl. In *Empirical Software Engineering, 2004. ISESE'04. Proceedings. 2004 International Symposium on*, pages 83–92. IEEE, 2004.
- [10] Rainer Koschke. *Survey of Research on Software Clones*. 2007.
- [11] Rainer Koschke. Similarity in software artifacts and its relation to code generation. April 2013. Code Generation 2013.
- [12] Akito Monden, Daikai Nakae, Toshihiro Kamiya, Shin-ichi Sato, and K-i Matsumoto. Software quality analysis by code clones in industrial legacy software. In *Software Metrics, 2002. Proceedings. Eighth IEEE Symposium on*, pages 87–94. IEEE, 2002.
- [13] National University of Singapore. XML-based variant configuration language. <http://xvcl.comp.nus.edu.sg/cms/>. Abgerufen: 12. Juni 2014.

## 6 Anhang: Beispielhafte Klone

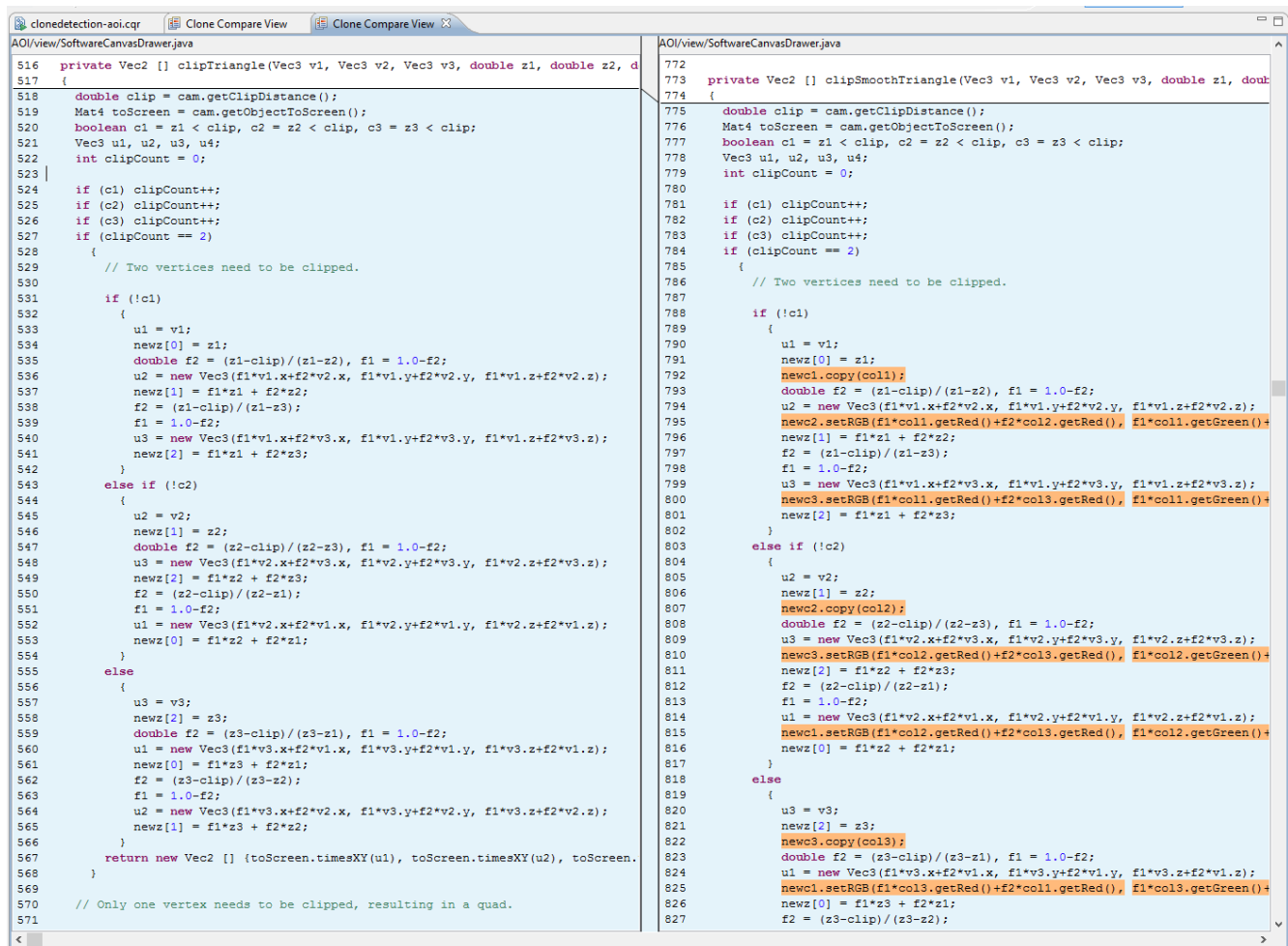


Abbildung 1: Beginn eines langen Klons in AoI

```

FreeMind/freemind/modes/mindmapmode/MindMapController.java
2026     }
2027 }
2028
2029 public boolean extendSelection(MouseEvent e) {
2030     NodeView newlySelectedNodeView = ((MainView) e.getComponent())
2031         .getNodeView();
2032     // MindMapNode newlySelectedNode = newlySelectedNodeView.getModel();
2033     boolean extend = e.isControlDown();
2034     // Fixes Cannot select multiple single nodes
2035     // https://sourceforge.net/tracker/?func=detail&atid=107118&aid=16758296
2036     if (Tools.isMacOsX()) {
2037         extend |= e.isMetaDown();
2038     }
2039     boolean range = e.isShiftDown();
2040     boolean branch = e.isAltGraphDown() || e.isAltDown(); /*
2041                                     * windows alt,
2042                                     * linux altgraph
2043                                     * ....
2044                                     */
2045     boolean retValue = false;
2046
2047     if (extend || range || branch
2048         || !getView().isSelected(newlySelectedNodeView)) {
2049         if (!range) {
2050             if (extend)
2051                 getView().toggleSelected(newlySelectedNodeView);
2052             else
2053                 select(newlySelectedNodeView);
2054             retValue = true;
2055         } else {
2056             retValue = getView().selectContinuous(newlySelectedNodeView);
2057             // /* fc, 25.1.2004: replace getView by controller methods.*/
2058             // if (newlySelectedNodeView != getView().getSelected() &&
2059             //     newlySelectedNodeView.isSiblingOf(getView().getSelected())) {
2060             //     getView().selectContinuous(newlySelectedNodeView);
2061             //     retValue = true;
2062             // } else {
2063             //     /* if shift was down, but no range can be selected, then the
2064             //     // new node is simply selected: */
2065             //     if (!getView().isSelected(newlySelectedNodeView)) {
2066             //         getView().toggleSelected(newlySelectedNodeView);
2067             //         retValue = true;
2068             //     }
2069             // }
2070             if (branch) {
2071                 getView().selectBranch(newlySelectedNodeView, extend);
2072                 retValue = true;
2073             }
2074         }
2075     }
2076     if (retValue) {
2077         e.consume();
2078     }
2079     // Display link in status line
2080     String link = newlySelectedNodeView.getModel().getLink();
2081     link = (link != null ? link : " ");
2082     getController().getFrame().out(link);
2083 }

```

```

FreeMind/freemind/modes/viewmodes/ViewControllerAdapter.java
74     NodeView newlySelectedNodeView = ((MainView) e.getComponent())
75         .getNodeView();
76     // MindMapNode newlySelectedNode = newlySelectedNodeView.getModel();
77     boolean extend = e.isControlDown();
78     boolean range = e.isShiftDown();
79     boolean branch = e.isAltGraphDown() || e.isAltDown(); /*
80                                     * windows alt,
81                                     * linux altgraph
82                                     * ....
83                                     */
84     boolean retValue = false;
85
86     if (extend || range || branch
87         || !getView().isSelected(newlySelectedNodeView)) {
88         if (!range) {
89             if (extend)
90                 getView().toggleSelected(newlySelectedNodeView);
91             else
92                 select(newlySelectedNodeView);
93             retValue = true;
94         } else {
95             retValue = getView().selectContinuous(newlySelectedNodeView);
96             // /* fc, 25.1.2004: replace getView by controller methods.*/
97             // if (newlySelectedNodeView != getView().getSelected() &&
98             //     newlySelectedNodeView.isSiblingOf(getView().getSelected())) {
99             //     getView().selectContinuous(newlySelectedNodeView);
100             //     retValue = true;
101             // } else {
102             //     /* if shift was down, but no range can be selected, then the
103             //     // new node is simply selected: */
104             //     if (!getView().isSelected(newlySelectedNodeView)) {
105             //         getView().toggleSelected(newlySelectedNodeView);
106             //         retValue = true;
107             //     }
108             // }
109             if (branch) {
110                 getView().selectBranch(newlySelectedNodeView, extend);
111                 retValue = true;
112             }
113         }
114     }
115     if (retValue) {
116         e.consume();
117     }
118     // Display link in status line
119     String link = newlySelectedNodeView.getModel().getLink();
120     link = (link != null ? link : " ");
121     getController().getFrame().out(link);
122 }
123     return retValue;
124 }
125
126 public void setFolded(MindMapNode node, boolean folded) {
127     _setFolded(node, folded);
128 }
129
130 public void startupController() {
131     _smar.startupController();

```

Abbildung 2: Möglicherweise unabsichtlich inkonsistenter Bugfix in FreeMind (#1)

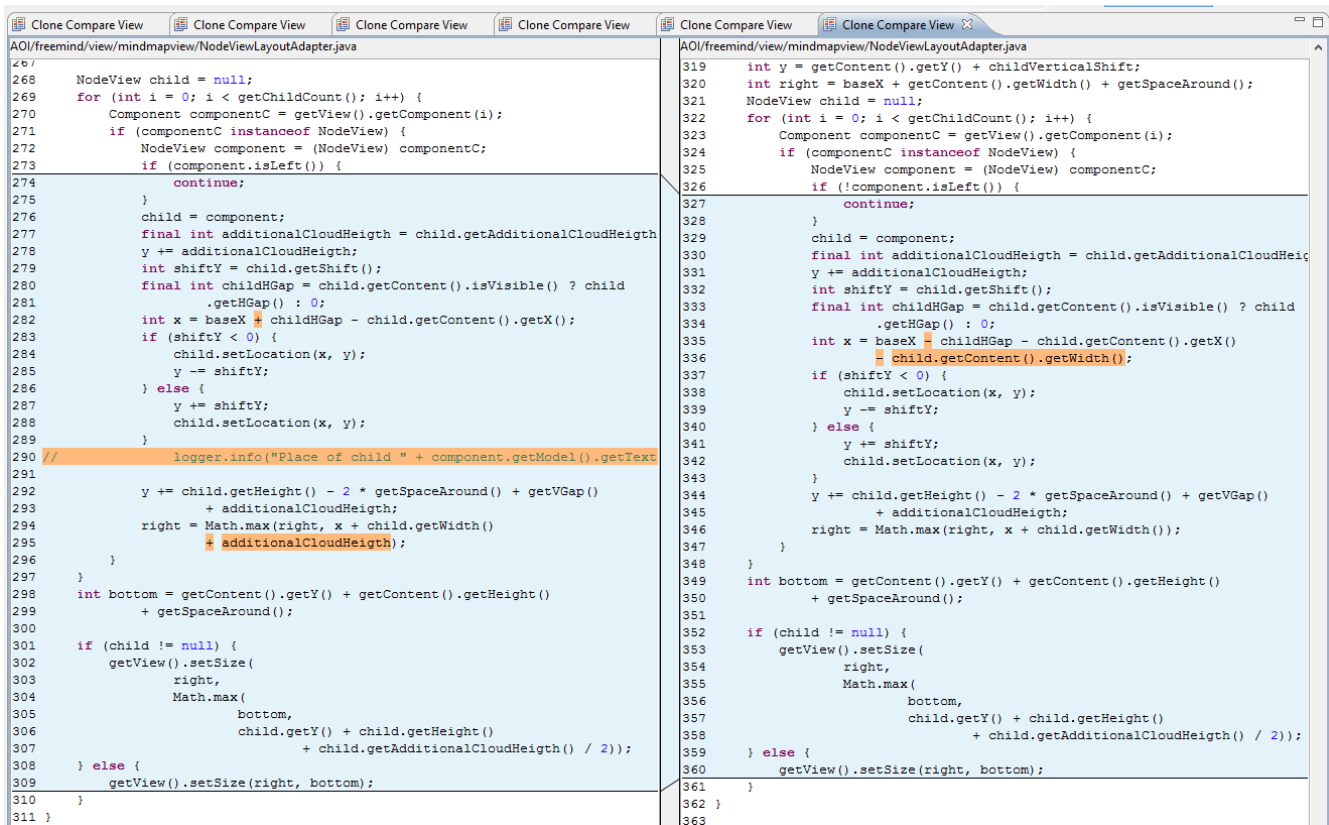


Abbildung 3: Möglicherweise unabsichtlich inkonsistenter Bugfix in FreeMind (#2)

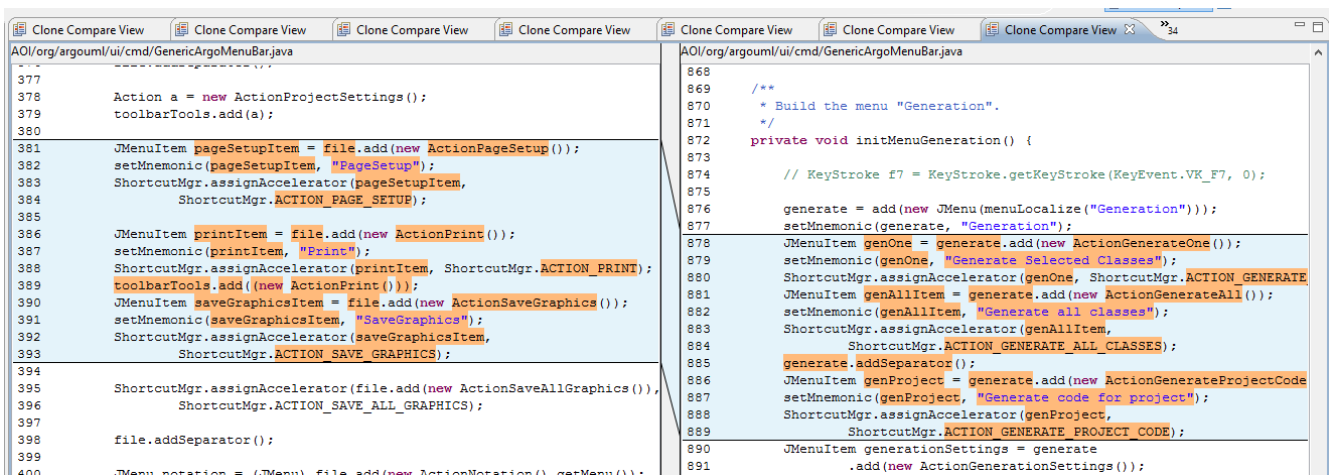


Abbildung 4: Ein false positive in ArgoUML

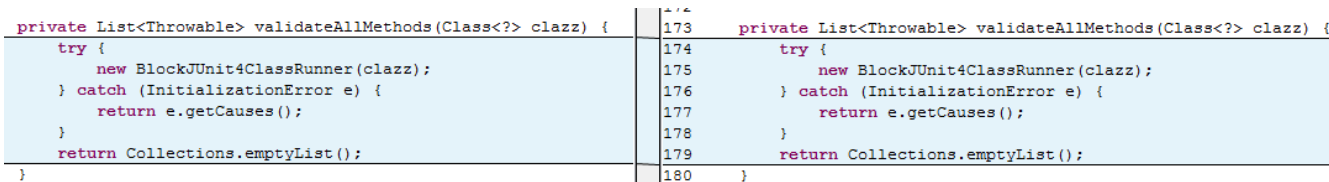


Abbildung 5: Ein uninteressanter Klon ohne Gaps in JUnit