



UNIVERSITÄT PADERBORN

Die Universität der Informationsgesellschaft

Fakultät für Elektrotechnik, Informatik und Mathematik

Institut für Informatik

Arbeitsgruppe Softwaretechnik

Bachelorarbeit

Gerichtet an die Arbeitsgruppe Softwaretechnik

zur Erreichung des Grades

Bachelor of Science

Unparsing von Datenstrukturen zur Analyse von C-Präprozessor-Variabilität

von

EUGEN SHULIMOV

Betreut durch:

Paul Maximilian Bittner

Prüfer:

Prof. Dr. Thomas Thüm

Prof. Dr. Stefan Sauer

Paderborn, 6. Januar 2025

Erklärung

Ich versichere, dass ich die Arbeit ohne fremde Hilfe und ohne Benutzung anderer als der angegebenen Quellen angefertigt habe und dass die Arbeit in gleicher oder ähnlicher Form noch keiner anderen Prüfungsbehörde vorgelegen hat und von dieser als Teil einer Prüfungsleistung angenommen worden ist. Alle Ausführungen, die wörtlich oder sinngemäß übernommen worden sind, sind als solche gekennzeichnet.

Ort, Datum

Unterschrift

Zusammenfassung. Es gibt verschiedene Möglichkeiten Variabilität in einem Projekt umzusetzen. Einer dieser Möglichkeiten ist die Nutzung von C-Präprozessor-Annotationen. Dies gestattet uns Variabilität umzusetzen. Es gibt eine Reihe an Analysen und Forschungsarbeiten, Entwickler bei der Umsetzung der Variabilität und deren Analyse zu unterstützen. Dazu werden Tools wie DiffDetective verwendet. Zwar hat DiffDetective einen Parser, aber keinen Unparser. Variation-Trees und Variation-Diffs sind zwei zentralen Datenstrukturen in DiffDetective, um Präprozessorvariabilität und Änderungen daran darzustellen. In dieser Arbeit präsentieren wir einen Unparser für Variation-Trees und Variation-Diffs. Wir haben diesen Algorithmus in DiffDetective implementiert und an einem großen Datensatz validiert. Die von uns gewählten Datensätze sind Vim, sylpheed, gcc und berkeley-db-libdb. Für die Validierung wurde von uns mehrere Korrektheitskriterien für unseren Unparser ausgearbeitet. Damit man feststellen kann, ob ein Unparser syntaktisch oder semantisch korrekt arbeitet.

Inhaltsverzeichnis

1	Einleitung	1
2	Hintergrundwissen	5
2.1	C-Präprozessor	5
2.2	Realisierung von Variabilität mit dem C-Präprozessor	6
3	Unparse-Algorithmus	9
3.1	Variation-Tree und Variation-Diff	9
3.2	Parser	12
3.3	Verlorengelassene Informationen und deren Wiederherstellung	16
3.4	Unparsing	21
3.5	Komplexitätsanalyse der Laufzeit	26
3.6	Zusammenfassung	28
4	Implementierung	31
5	Korrektheit	35
5.1	Korrektheitskriterium	35
5.2	Auswertung	42
5.2.1	Aufbau des Experiments	43
5.2.2	Ergebnisse	43
5.2.3	Diskussion	45
5.3	Zusammenfassung	47
6	Verwandte Arbeiten	49
6.1	Parsen von C-Präprozessor-Annotationen für Variabilitätsanalysen	49
6.2	Dekompilierung von C	50
6.3	Variabilitätsanalysen	50
7	Fazit und Zukünftige Arbeiten	53
	Bibliography	55

Einleitung

Bei der Entwicklung von konfigurierbaren Softwaresystemen, wie zum Beispiel Clone-and-Own, oder Softwareproduktlinien, gibt es im Laufe des Lebenszyklus immer mehr Features. Es ist von Vorteil, eine Möglichkeit zu haben, die Features im Code zu unterscheiden und automatisch zu finden. Einer dieser Vorteile ist, dass einer Produktlinie die Varianten automatisch ableiten kann. Einige Möglichkeiten dazu sind Präprozessor-Annotationen, oder Build-Systeme [2]. Bei der Clone-and-Own-Entwicklung wird für jede Variante der Software eine neue Kopie der gesamten Software angelegt und parallel entwickelt [3]. Dort müssen Features gefunden werden, um diese zu aktualisieren [3, 14, 15, 16, 17, 33].

Der C-Präprozessor ist eine Möglichkeit, Variabilität zu erzeugen [2]. Der C-Präprozessor ist ein Tool, das den Quellcode vor dem Kompilieren manipuliert [2]. Dieses Tool bietet Möglichkeiten zur Dateieinbindung, zu lexikalischen Makros, und zur bedingten Kompilierung [2]. Wie ein mit dem C-Präprozessor annotierter Code aussehen kann, ist in der Abbildung 1.1 Stelle ⑤ zu sehen (Abb. 1.1 St.⑤). Um die Variabilität mithilfe des C-Präprozessors zu erzeugen, brauchen wir dessen Möglichkeit zur bedingten Kompilierung [2]. Dabei können beliebige Formeln in Aussagenlogik über Features im Quellcode mit den C-Präprozessor-Anweisungen `#if`, `#ifdef` und `#ifndef` abgebildet werden [3] (Abb. 1.1 St.⑥). Die Anweisung `#if` ist wie die `if`-Anweisung aus gewohnten Programmiersprachen, `#ifdef` und `#ifndef` sind ähnlich zu `#if` aber reagieren auf Definition oder nicht Definition eines oder mehreren Makros. Als Eingabe erhält der C-Präprozessor einen mit C-Präprozessor-Annotationen annotierten Code. Dieser Eingabe wird gemäß der C-Präprozessor-Annotationen bearbeitet und als Ausgabe erhält man einen Code, welche für die Kompilierung bereit ist.

Die Entwickler sind bei der Entwicklung von konfigurierbarer Software daran interessiert, zu verstehen, wie sich ihre Änderungen auf die Variabilität auswirken und wie die Variabilität von konfigurierbarer Software aussieht [3]. Sonst wenn man das Verständnis über die Auswirkungen der Änderung nicht hat, kann das zu Fehlern und Problemen bei der Entwicklung führen [3, 13, 21, 22, 25, 29]. Dies stellt einen Aspekt einer größeren Aufgabe dar, der Aufrechterhaltung und Weiterentwicklung von Informationen über Variabilität bei Quellcodeänderungen [3]. Für die Entwickler stellt diese Aufgabe eine große Herausforderung dar [3, 23, 24, 26].

Zur Unterstützung der Variabilitätsanalyse kann man Tools verwenden [28, 32], wie zum Beispiel DiffDetective [5]. Der Zweck von DiffDetective ist es, Änderungen im Quellcode und Änderungen der Variabilität darstellbar und den Zusammenhang zwischen ihnen analysierbar

zu machen. DiffDetective stellt einen variabilitätsbezogenen Differencer [3, 5] zur Verfügung, der sich nur auf Aspekte im Code/Text bezieht, welche die Variabilität berücksichtigen. Diese Bibliothek ermöglicht auch die Analyse der Versionshistorie von Softwareproduktlinien [3] und bietet daher einen flexiblen Rahmen für großangelegte empirische Analysen von Git-Versionsverläufen statisch konfigurierbarer Software [4, 5].

Zentral für DiffDetective sind zwei formal verifizierte Datenstrukturen für Variabilität und Änderungen an dieser [3]. Das sind Variation-Trees (Abb. 1.1 St.ⓧ) für variabilitätsbezogenen Code (Abb. 1.1 St.Ⓥ) und Variation-Diffs (Abb. 1.1 St.Ⓨ) für variabilitätsbezogene Diffs (Abb. 1.1 St.Ⓦ). Diese Datenstrukturen sind generisch. Das bedeutet, dass die Datenstrukturen möglichst von der Umsetzung der Variabilität im Code abstrahieren. Also kann eine Umsetzungsmöglichkeit leicht durch eine andere ersetzt werden, zum Beispiel können C-Präprozessor-Annotationen durch Java-Präprozessor-Annotationen ohne oder geringer Änderungen an den Datenstrukturen selbst, ersetzt werden. Ein Variation-Tree ist ein Baum, welcher die Verzweigungen/Variationen eines annotierten Codes darstellt [3, 4, 5]. Ein Variation-Diff ist ein Graph, welcher die Unterschiede zwischen zwei Variation-Trees zeigt [3, 4, 5]. In beiden Fällen werden die Bedingungsknoten, welche Informationen zu Variabilität erhalten, und die Code-Knoten unterschieden. Beim Variation-Diff sind zudem die eingefügten Knoten, die gelöschten Knoten und, die unveränderten Knoten zu unterscheiden.

Das Parsen führt die Eingabe von der konkreten Syntax in die abstrakte Syntax um. In unserem Fall parst DiffDetective C-Präprozessor-Annotationen, dieses kann aber auch auf andere Präprozessor-Annotationen erweitert werden. Beim Parsen wird nur der C-Präprozessor-annotierte Code in seine abstrakte Syntax überführt. Der C- bzw. C++-Code wird als Text behandelt und wird nicht geparst. Das Parsen in DiffDetective funktioniert für Variation-Trees und für Variation-Diffs über einen einzigen gemeinsamen Algorithmus. Der Algorithmus ist an sich für das Parsen von textbasierten Diffs in Variation-Diffs ausgelegt (Abb. 1.1 St.Ⓟ). An den Stellen ① und ⑨ wird anders vorgegangen, da wir als Eingabe C-Präprozessor-Code (Abb. 1.1 St.Ⓥ) haben und als Ausgabe einen Variation-Tree (Abb. 1.1 St.ⓧ). Der gegebene Algorithmus ist für das direkte Parsen von C-Präprozessor-Code nicht ausgelegt. Deshalb wurden dort Umwege verwendet, um diesen Algorithmus anwendbar zu machen und die benötigte Ausgabe zu erzielen. Ein Text kann in ein textbasiertes Diff umgewandelt werden, in dem jede Zeile als unverändert angesehen wird, durch die Bildung eines Diffs mit sich selbst. Dadurch ist es möglich aus C-Präprozessor-Code (Abb. 1.1 St.Ⓥ) ein textbasiertes Diff (Abb. 1.1 St.Ⓦ) zu erzeugen, also wurden die Stelle ⑪ oder ⑫ verwendet. Da jetzt ein textbasiertes Diff vorhanden ist, kann der Algorithmus darauf angewandt werden (Abb. 1.1 St.Ⓟ). Um aus dem erhaltenen Variation-Diff (Abb. 1.1 St.Ⓨ) ein Variation-Tree zu bekommen, muss man die Stelle ④ oder ⑧ aus der Abbildung 1.1 anwenden. So sieht man, dass die Stelle ① durch die Stellen ⑪,⑤,④ und die Stelle ⑨ durch die Stellen ⑫,⑤,⑧ ersetzt werden kann.

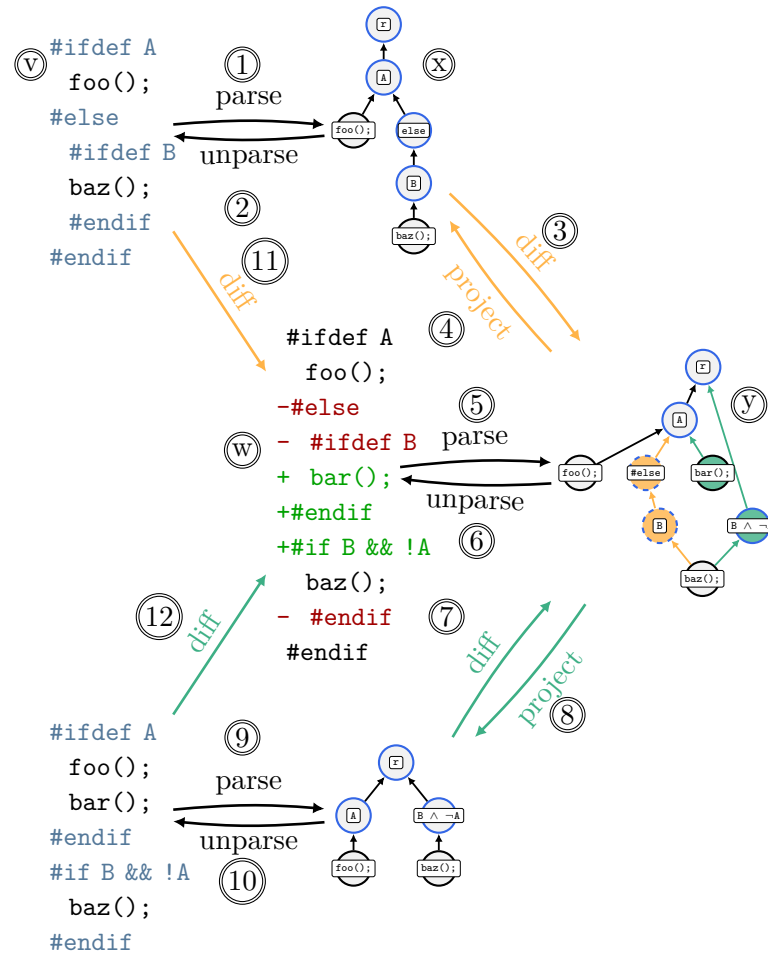


Abbildung 1.1: Überblick über Variability-Aware Differencing [5]

Eine Möglichkeit zur Analyse von Softwareproduktlinien sind Mutation-Tests. Bei Mutation-Tests werden Mutation-Operatoren verwendet, welche aber nur auf der abstrakten Ebene, also auf Variation-Trees, angewandt werden können [1]. Um weiter in der Analyse vorzugehen, muss man von der abstrakten Ebene zu der konkreten Ebene übergehen und hier wird der Unparser gebraucht. Obwohl es einen Parser (Abb. 1.1 St. (1), (5), (9)) für Variation-Trees und Variation-Diffs gibt, gibt es keinen Unparser (Abb. 1.1 St. (2), (6), (10)) für Variation-Trees und Variation-Diffs. Unser Ziel ist es, das zu ändern. Dazu müssen wir einen Unparser entwickeln, welcher auf direktem oder indirektem Wege, Variation-Trees (Abb. 1.1 St. (X)) in C-Präprozessor-Annotierten Code (Abb. 1.1 St. (V)) und Variation-Diffs (Abb. 1.1 St. (Y)) in textbasierte Diffs (Abb. 1.1 St. (W)) überführt. Dabei ist Unparsen eine Überführung aus der abstrakten Syntax in die konkrete Syntax, also ist das Unparsen die Invertierung des Parsens.

Für das Unparsen stellt das Fehlen einiger Informationen, die im annotierten Code vorhanden sein müssen, aber in Variation-Trees bzw. Variation-Diffs nicht vorhanden sind, das größte Problem dar. Diese Informationen sind entweder durch das Parsen verloren gegangen oder waren von Anfang an nicht vorhanden, zum Beispiel wenn Variation-Trees bzw. Variation-Diffs synthetisch erzeugt wurden. Diese Informationen sind die exakte Formel, die ein Mapping-Knoten $\tau(v) = \text{mapping}$ besitzt [3] und die Position von `#endif` und deren Einrückung. Aus diesem Grund müssen wir entweder Annahmen treffen, oder DiffDetective so erweitern, dass er diese

Information explizit speichert. Eine Annahme könnte sein, dass das `#endif` genauso eingerückt ist, wie die Bedingung, zu der es gehört.

In dieser Arbeit werden wir ein Vorgehen zum Unparsen von Variation-Trees und Variation-Diffs in das ursprüngliche Textformat ausarbeiten, dieses Vorgehen in DiffDetective implementieren und anhand von uns ausgearbeiteten Korrektheitskriterien in einer Auswertung die Korrektheit unseren Unparsers bestimmen.

Um die Korrektheit des von uns entwickelten und implementierten Unparser zu bestimmen, kann man wie folgt vorgehen. Zuerst hat man eine Ausgangsdatei. Diese Datei wird geparst und wir erhalten das Ergebnis des Parsens. Als Nächstes wird dieses Ergebnis genommen und ungeparst. Damit bekommen wir das Ergebnis des Unpares. Zum Schluss wird das Ergebnis des Unparesen mit der Ausgangsdatei verglichen und bewertet ob der Unparser korrekt gearbeitet hat. Zum Vergleichen des Ergebnisses des Unparsers und der Ausgangsdatei haben wir drei Korrektheitskriterien definiert. Diese Korrektheitskriterien sind syntaktische Korrektheit, syntaktische Korrektheit ohne Whitespace und semantische Korrektheit. Wir haben mehrere Korrektheitskriterien definiert, da abhängig davon welche, wie viel Information beim Parsen verloren gehen, sind unterschiedlich korrekte Rekonstruktionen möglich. Bei Parsen und Unparsen gibt es eine gegenseitige Abwägung, je mehr Informationen ausgelassen werden, desto leichter ist es einen Parser zu bauen aber dann ist es schwieriger zu unparsen. In die andere Richtung gilt es auch, je mehr sich Informationen gemerkt werden soll, dann ist das Parsen schwieriger aber dann das Unparsen leichter.

Hintergrundwissen

In diesem Kapitel stellen wir Hintergrundwissen zur Verfügung, dieses Wissen ist von Bedeutung für das Verständnis dieser Arbeit. Es handelt sich um C-Präprozessor und einer seiner Einsatzmöglichkeiten. In dem Abschnitt 2.1 dieses Kapitels wird der C-Präprozessor vorgestellt. Seine Möglichkeiten und Anweisungen zusammen mit einem Beispiel. Wie der C-Präprozessor zur Umsetzung der Variabilität im Code genutzt werden kann und welche Bestandteile von dem C-Präprozessor dazu nötig sind, wird im Abschnitt 2.2 des Kapitels erläutert.

2.1 C-Präprozessor

Der C-Präprozessor ist ein Tool, das Quellcode vor dem Kompilieren manipuliert [2]. Dieses Tool bietet Möglichkeiten zur bedingte Kompilierung, zur Dateieinbindung und zur Erstellung lexikalischer Makros [2]. Eine C-Präprozessor-Direktive beginnt mit `#` und geht bis zum ersten Whitespace-Zeichen weiter, optional kann nach der Direktive ein Argument im Rest der Zeile stehen. Der C-Präprozessor hat solche Anweisungen wie, `#include` zum Einbinden von Dateien, um zum Beispiel Header-Dateien wiederzuverwenden. Wie das Aussehen kann, ist in der Abbildung 2.1, Zeile 1 zu sehen. Mit den Anweisungen `#if` (Zeile 6), `#else` (Zeile 10), `#elif` (Zeile 8), `#ifdef` (Zeile 18), `#ifndef` (Zeile 3), und `#endif` (Zeile 5) wird die bedingte Kompilierung erzeugt. Dabei funktionieren `#if`, `#else`, `#elif`, und `#endif` vergleichbar mit dem, was man aus Programmiersprachen und Pseudocode gewohnt ist. `#ifdef` ist ähnlich zu `#if`, wird aber nur dann wahr, wenn das darauf folgende Makro definiert ist. `#ifndef` ist die Negation von `#ifdef`. Makros werden durch die Anweisung `#define` (Zeile 4) erstellt. Der Präprozessor ersetzt dann während seiner Arbeit, den Makronamen durch seine Definition. Während dieser Arbeit kann ein Makro definiert, undefiniert und undefiniert, mit `#undef` (Zeile 2), werden. Der C-Präprozessor hat noch weitere Anweisungen, auf die wir nicht weiter eingehen. Auf weitere Anweisungen wird nicht eingegangen, da diese bezüglich unserer Arbeit nicht relevant sind. Der C-Präprozessor kann in anderen Programmiersprachen verwendet werden. Beispiel für solchen Sprachen sind C++, Assemblersprachen, Fortran und Java. Der Grund dafür ist, dass der C-Präprozessor unabhängig von der zugrundeliegenden Programmiersprache ist.

Der C-Präprozessor-Annotierter Code in der Abbildung 2.1 arbeitet wie folgt. Zuerst in Zeile 1 wird der `studio`-Header eingebunden um mit dem Input und Output zu arbeiten. Danach in Zeile 2 wird der Makro `N` undefiniert. In Zeile 2 wird geprüft ob der Makro `N` existiert, wenn

nicht, wird der mit dem Wert 10 definiert. Als Nächstes in Zeilen 6 bis 12 wird der Makro A abhängig von dem Wert von N definiert. Bei N größer 10 enthält A die Zeichenkette **O.O** bei N gleich 10 wird A auf **;** gesetzt in restlichen Fällen ist A **:(**. In Zeile 14 beginnt die main-Funktion. In Zeile 15 gibt es die Funktion printf, welche Hello world! auf der Konsole ausgibt. Danach in Zeile 17 wird geprüft, ob N definiert wurde, wenn ja, dann bekommt die main-Funktion eine for-Schleife. Diese for-Schleife wird N Mal durchlaufen und gibt bei jedem Schleifendurchlauf auf der Konsole den Makro A aus. In der Zeile 25 gibt es den standardmäßigen return 0.

```

1 | #include <stdio.h>
2 | #undef N
3 | #ifndef N
4 | #define N 10
5 | #endif
6 | #if N > 10
7 | #define A "O.O"
8 | #elif N == 10
9 | #define A ";"
10 | #else
11 | #define A ":(
12 | #endif
13 |
14 |     int main()
15 |     {
16 |         printf("Hello world!");
17 | #ifdef N
18 |         int i;
19 |         for (i = 0; i < N; i++)
20 |         {
21 |             printf(A);
22 |         }
23 | #endif
24 |
25 |         return 0;
26 |     }

```

Abbildung 2.1: Beispiel für C Code mit Präprozessor Anweisungen

2.2 Realisierung von Variabilität mit dem C-Präprozessor

Der C-Präprozessor ist eine Möglichkeit, Variabilität zu erzeugen [2]. Um Variabilität mithilfe des C-Präprozessors zu erzeugen, brauchen wir dessen Möglichkeit zur bedingten Kompilierung [2]. Dies wird mit den C-Präprozessor-Anweisungen **#if** (Abb.2.1 Z6), **#else** (Abb.2.1 Z10), **#elif** (Abb.2.1 Z8), **#ifdef** (Abb.2.1 Z18), **#ifndef** (Abb.2.1 Z2), und **#endif** (Abb.2.1 Z4) bewerkstelligt. Dabei werden Codefragmente von diesen Anweisungen eingeschlossen. Danach, abhängig davon welche Makros definiert sind und auch wie sie definiert sind, werden bestimmte Codefragmente entweder behalten oder entfernt. Die Abbildung 2.2 zeigt, ein von uns erstelltes Beispiel, wie ein mit C-Präprozessor-Annotierter Code aussehen kann. Das Beispiel zeigt, dass die Anweisungen von den C-Präprozessor verschachtelt werden können. Bei der Realisierung von Variabilität wird oft mit Features gearbeitet. Ein Feature ist dabei ein Merkmal oder ein für den Endbenutzer sichtbares Verhalten eines Softwaresystems. Ob eine C-Code Zeile im endgültigen Programm auftaucht oder nicht wird durch die dazugehörige Bedingung bestimmt. Diese Bedingungen werden durch C-Präprozessor-Annotation dargestellt. Es ist möglich, mit den C-Präprozessor Anweisungen eine große Menge an Bedingungen abzubilden [3]. Zur leicht-

```

1 | #ifdef FEATURE_A && FEATURE_B
2 |     foo();
3 |     bar();
4 |     int i = 18
5 | #ifdef FEATURE_D
6 | #define SIZE 200
7 |     foom();
8 | #else
9 | #define SIZE 175
10 |     i = 17;
11 | #endif
12 |     too(i);
13 | #endif
14 | #ifdef FEATURE_C
15 |     baz();
16 | #ifndef FEATURE_B
17 |     sho();
18 | #define SIZE 100
19 | #endif
20 |     bazzz();
21 | #else
22 |     boom();
23 |     broo();
24 | #endif
25 | #if SIZE > 180
26 |     long j;
27 | #elif SIZE < 111
28 |     short j;
29 | #else
30 |     int j;
31 | #endif

```

Abbildung 2.2: Beispiel für Umsetzung der Variabilität mit C-Präprozessor

```

1 |     foo();
2 |     bar();
3 |     int i = 18
4 |     i = 17
5 |     too(i);
6 |     boom();
7 |     broo();
8 |     int j;

```

Abbildung 2.3: Ausgabe des C-Präprozessors wenn Feature A=1, B=1, C=0, D=0

```

1 |     baz();
2 |     sho();
3 |     bazzz();
4 |     short j;

```

Abbildung 2.4: Ausgabe des C-Präprozessors wenn Feature A=0, B=0, C=1, D=0

terer Pflege werden diese Bedingungen oft in eine Folge von mehreren Feature-Annotationen aufgeteilt. Beispiel dafür ist Zeile 7 zu sehen, wo seine Feature-Annotationen `FEATURE_A && FEATURE_B` und `FEATURE_D` sind, aber seine Bedingung `FEATURE_A && FEATURE_B && FEATURE_D`. In diesem Code-Beispiel ist auch die Abhängigkeit einiger Features von anderen zu erkennen, wie in Zeile 5, wo die Auswahl des Features D nur dann möglich ist, wenn auch die Features A und B ausgewählt sind. Die Zeile 17 in Abbildung 2.2 hat als Feature-Annotationen `FEATURE_C` und `!FEATURE_B` die Bedingung dabei aber ist `FEATURE_C && !FEATURE_B`. Nicht nur die Definition von Features, sondern auch die Nicht-Definition kann, einen Einfluss auf das Ergebnis haben, wie in der Zeile 16 zu sehen ist. In dem Beispiel wird, wie bei der Implementierung von Softwareproduktlinien, ein Name pro Feature reserviert. Wenn das Feature dann ausgewählt wird, wird dann ein Makro mit Feature-Namen definiert, mit der Anweisung `#define FEATURE_NAME`. Die Abbildungen 2.3 und 2.4 stellen 2 mögliche Ergebnisse der C-Präprozessor-Ausführung dar. Diese Ergebnisse werden als Varianten bezeichnet. Dabei werden für die Abbildung 2.3 die Features A und B definiert und für die Abbildung 2.4 nur das Feature C. Eine Konfiguration ist eine Auswahl der Features welche ausgewählt und nicht ausgewählt werden. Es ist zu sehen, dass der generierter Code nur in dem Bezeichner j gleich ist und sonst nicht. Das veranschaulicht, wie unterschiedlich das Programm sein kann.

2.2 REALISIERUNG VON VARIABILITÄT MIT DEM C-PRÄPROZESSOR

Unparse-Algorithmus

In diesem Kapitel stellen wir unseren Algorithmus zum Unparsen von Variation-Trees und unsere Methode für das Unparsen von Variation-Diffs vor. Wir beschreiben den theoretischen Hintergrund, welche Bedingungen erfüllt werden müssen, damit der Algorithmus korrekt funktioniert, und die Arbeitsweise des Algorithmus.

Wir beschäftigen uns mit der Definition von Variation-Tree und Variation-Diff, nehmen neue Definitionen für Variation-Tree und Variation-Diff aus anderer Arbeit und erweitern diese. Es werden auf Bedingungen aufgestellt, ohne die der Algorithmus nicht korrekt funktionieren kann. Der von uns entwickelte Algorithmus basiert auf dem Prinzip der Tiefensuche und kann nur für das Unparsen von Variation-Trees verwendet werden. Für das Unparsen von Variation-Diffs reduzieren wir das Problem auf das Unparsen von Variation-Trees indem wir statt Variation-Diff zu unparsen, den Variation-Diff projizieren, dann zwei Variation-Trees unparsen und schließlich ein Diff über das Ergebnis bilden.

Wir fangen an mit Definitionen von Variation-Tree und Variation-Diff in Kapitel 3.1. Danach im Kapitel 3.2 sehen wir den Parser, welcher Variation-Trees und Variation-Diffs erstellt. Im Kapitel 3.3 werden die, während des Parsens, verlorengehende Informationen bestimmt und Möglichkeiten vorgestellt diese zurückzubekommen. Aufbauend auf den vorherigen wird im Kapitel 3.4 unser Algorithmus zum Unparsen von Variation-Trees vorgestellt. Schlussendlich wird im Kapitel 3.5 eine Laufzeitanalyse des Algorithmus durchgeführt.

3.1 Variation-Tree und Variation-Diff

Um verstehen zu können, wie wir Variation-Trees und -Diffs unparsen können, müssen wir uns mit den Einschränkungen und Möglichkeiten des Parsers vertraut machen.

Um Variation-Trees und Variation-Diffs kennenzulernen, betrachten wir zunächst die ursprüngliche Definition der Datenstrukturen [3].

Definition 3.1. *Ein VARIATION-TREE (V, E, r, τ, l) ist ein Baum mit Knoten V , Kanten $E \subseteq V \times V$ und Wurzelknoten $r \in V$. Jede Kante $(x, y) \in E$ verbindet einen Kindknoten x mit seinem Elternknoten y , bezeichnet mit $p(x) = y$. Der Knotentyp $\tau(v) \in \{\text{ARTIFACT}, \text{MAPPING}, \text{ELSE}\}$ identifiziert einen Knoten $v \in V$ entweder als Vertreter eines Implementierungsartefakts, einer Feature-Annotation oder eines else-Zweigs. Das Label $l(v)$ ist eine aussagenlogische Formel,*

wenn $\tau(v) = \text{MAPPING}$, ein Verweis auf ein Implementierungsartefakt, wenn $\tau(v) = \text{ARTIFACT}$, oder leer, wenn $\tau(v) = \text{ELSE}$ ist. Der Wurzelknoten r hat den Typ $\tau(r) = \text{MAPPING}$ und das Label $l(r) = \text{TRUE}$. Ein Knoten e von Typ $\tau(e) = \text{ELSE}$ kann nur unterhalb eines Nichtwurzelknotens v mit dem Typ $\tau(v) = \text{MAPPING}$ platziert werden, dabei hat ein Knoten w von Typ $\tau(w) = \text{MAPPING}$ höchstens einen Knoten u vom Typ $\tau(u) = \text{ELSE}$.

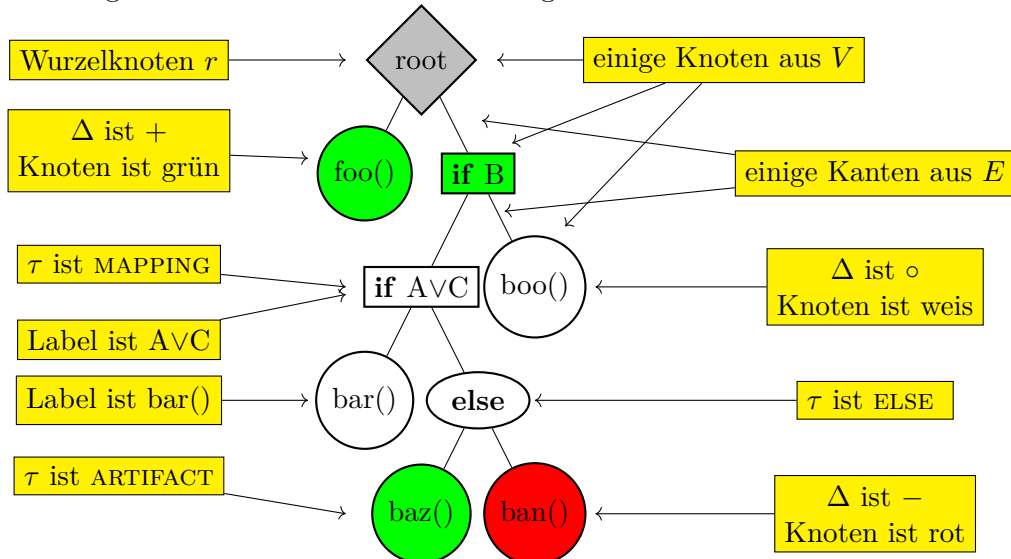
Definition 3.2. Ein VARIATION-DIFF ist ein gerichteter, zusammenhängender, azyklischer Graph $D = (V, E, r, \tau, l, \Delta)$, welcher einen Wurzelknoten hat, mit Knoten V , Kanten $E \subseteq V \times V$, Wurzelknoten $r \in V$, Knotentyp τ , Knotenlabel l und einer Funktion $\Delta : V \cup E \rightarrow \{+, -, \circ\}$, die definiert, ob ein Knoten oder eine Kante hinzugefügt $+$ wurde, entfernt $-$ wurde oder unverändert \circ geblieben ist, so das $\text{PROJECT}(D, t)$ ein Variation-Tree für alle Zeiten $t \in \{a, b\}$ ist.

Definition 3.3. Die Projektion $\text{PROJECT}(D, t)$ für ein Variation-Diff ist ein Variation-Tree, der durch das Entfernen von Δ und den Knoten und Kanten, welche zu der Zeit t nicht vorhanden sind. $\text{PROJECT}((V, E, r, \tau, l, \Delta), t) := (\{v \in V \mid \text{EXISTS}(t, \Delta(v))\}, \{e \in E \mid \text{EXISTS}(t, \Delta(e))\}, r, \tau, l)$

Ob ein Knoten oder eine Kante zu einer gegebenen Zeit existiert oder nicht, wird durch EXISTS bestimmt. EXISTS ist hier genauso definiert wie in [3] und [20].

Definition 3.4. Ob ein Knoten oder eine Kante zu einer gegebenen Zeit existiert oder nicht, stellt EXISTS für $x \in V \cup E$ fest mit $\text{EXISTS}(t, x) := (t = \text{b} \wedge \Delta(x) \neq \text{g}) \vee (t = \text{a} \wedge \Delta(x) \neq -)$.

In der Abbildung unten ist ein Beispiel für ein Variation-Diff gegeben. Es hilft auch bei dem Verständnis von Variation-Trees da diese, ähnlich zu Variation-Diffs sind. Die Abbildung zeigt wie die unterschiedlichen Komponenten des Variation-Diffs visuell dargestellt werden können. Diese Darstellung wurde an der Darstellung eines Variation-Diffs aus DiffDetective [5] angelehnt. In der Abbildung ist ein Variation-Diff und gelbe Kästen mit Pfeilen, welche die Bestandteile des Variation-Diffs beschreiben. Die Form eines Knotens stellt seinen Knotentyp τ dar, runde Knoten haben ARTIFACT als τ , rechteckige Knoten haben MAPPING als τ und elliptische Knoten haben ELSE als τ . Die Farbe eines Knotens stellt sein Δ dar, wiese Knoten haben Δ gleich 0, grüne Knoten haben Δ gleich + und rote Knoten haben Δ gleich $-$. Der Text in den Knoten stellt



den Label dar.

Das ist aber nicht die einzige Möglichkeit Variation-Tree und Variation-Diff zu definieren. In [20] wurden die Variation-Tree und Variation-Diff etwas anders definiert. Obwohl dort auch Variation-Tree und Variation-Diff definiert werden, werden wir in dieser Arbeit die neuen Definitionen als geordneter Variation-Tree und als geordneter Variation-Diff bezeichnen. Diese Definitionen haben eine Eigenschaft, welche wir brauchen um das Unparsen zu bewerkstelligen.

Das ist die Ordnung der Kinderknoten, ohne die wir nicht eindeutig wissen, wie der Inhalt der Knoten einzuordnen ist. Genauer gehen wir darauf in Unterkapitel 3.3 ein. Die Definitionen von geordneten Variation-Trees und geordneten Variation-Diff sehen wie folgt aus:

Definition 3.5. Ein *GEORDNETER VARIATION-TREE* (V, E, r, τ, l, O) ist ein geordneter Baum, bei dem V , E , r , τ und l genauso definiert sind wie in der Definition 3.1. Dazu hat ein geordneter Variation-Tree eine injektive Funktion $O : V \rightarrow \mathbb{N}$, die eine Ordnung für die Kinder eines jeden Knotens jedes Knotens definiert.

Definition 3.6. Ein *GEORDNETER VARIATION-DIFF* $D = (V, E, r, \tau, l, \Delta, O_{\text{before}}, O_{\text{after}})$ ist ein gerichteter, zusammenhängender, azyklischer Graph, bei dem V , E , r , τ , l , und Δ genauso definiert sind wie in der Definition 3.2. Die Reihenfolge der Kinderknoten vor der Änderung O_{before} und nach der Änderung O_{after} sind eine injektive Funktion $O_{\text{before}}, O_{\text{after}} : V \rightarrow \mathbb{N}$. Die Projektionen $\text{project}_O(D, t)$ müssen für alle Zeiten $t \in \{\text{after}, \text{before}\}$ ein Variation-Tree mit demselben Wurzelknoten sein.

Aus Gründen der Eindeutigkeit haben wir auch die Projektion umbenannt, da sich die Definition von $\text{project}_O(D, t)$ von der Definition der Projektion $\text{PROJECT}(D, t)$ aus der Definition 3.3 unterscheidet, wegen der zusätzlichen Informationen, die gespeichert werden.

Definition 3.7. Die Projektion $\text{project}_O(D, t)$ eines geordneten Variation-Diffs $D = (V, E, r, \tau, l, \Delta, O_{\text{before}}, O_{\text{after}})$ zum Zeitpunkt $t \in \{\text{after}, \text{before}\}$ ist definiert als:
 $\text{project}_O(D, t) := (V', E', r, \tau, l, O_t)$, wobei $V' = \{v \in V \mid \text{EXISTS}(t, \Delta(v))\}$,
 $E' = \{e \in E \mid \text{EXISTS}(t, \Delta(e))\}$ und die Existenz von $d \in \{+, -, \circ\}$, zu der Zeit $t \in \{\text{after}, \text{before}\}$ ist in der Definition 3.4 gegeben.

Die Definitionen von normalen Variation-Tree und Variation-Diff sind sehr ähnlich zu den Definitionen von geordneten Variation-Tree und Variation-Diff. Es ist so, da geordnete Variation-Tree bzw. Variation-Diff eine Erweiterung von normale Variation-Tree bzw. Variation-Diff sind. Aus diesem Grund können wir geordnete Variation-Tree bzw. Variation-Diff in normale Variation-Tree bzw. Variation-Diff umwandeln.

Definition 3.8. $\text{reduce}_{\text{OVT}}$ wandelt einen geordneten Variation-Tree (Definition 3.5) in ein normales Variation-Tree (Definition 3.1)

$$\text{reduce}_{\text{OVT}}((V, E, r, \tau, l, O)) := (V, E, r, \tau, l).$$

Definition 3.9. $\text{reduce}_{\text{OVD}}$ wandelt einen geordneten Variation-Diff (Definition 3.6) in ein normales Variation-Diff (Definition 3.2)

$$\text{reduce}_{\text{OVD}}((V, E, r, \tau, l, \Delta, O_{\text{before}}, O_{\text{after}})) := (V, E, r, \tau, l, \Delta).$$

Es bleibt uns nur noch zu zeigen, dass die Reihenfolge der Anwendung nicht von Bedeutung ist. Damit wir die Abbildung 3.1 bekommen, welche zeigt wie geordnete Variation-Trees, geordnete Variation-Diffs, normale Variation-Trees und normale Variation-Diffs transformiert werden können.

Lemma 3.10. Für ein geordneten Variation-Diff D und eine Zeit $t \in \{\text{before}, \text{after}\}$ gilt $\text{reduce}_{\text{OVT}}(\text{project}_O(D, t)) = \text{project}(\text{reduce}_{\text{OVD}}(D), t)$.

$$\begin{aligned} & \text{Beweis. } \text{reduce}_{\text{OVT}}(\text{project}_O(D, t)) \\ &= \text{reduce}_{\text{OVT}}(\text{project}_O((V, E, r, \tau, l, \Delta, O_{\text{before}}, O_{\text{after}}), t)) \\ &= \text{reduce}_{\text{OVT}}((\{v \in V \mid \text{EXISTS}(t, \Delta(v))\}, \{e \in E \mid \text{EXISTS}(t, \Delta(e))\}, r, \tau, l, O_t)) \\ &= (\{v \in V \mid \text{EXISTS}(t, \Delta(v))\}, \{e \in E \mid \text{EXISTS}(t, \Delta(e))\}, r, \tau, l) \\ &= \text{project}((V, E, r, \tau, l, \Delta), t) \\ &= \text{project}(\text{reduce}_{\text{OVD}}(V, E, r, \tau, l, \Delta, O_{\text{before}}, O_{\text{after}}), t) \\ &= \text{project}(\text{reduce}_{\text{OVD}}(D), t) \end{aligned}$$

□

Jetzt haben wir uns mit dem beschäftigt wie Variation-Tree und Variation-Diff zu verstehen sind. Dabei haben wir zwei Definitionen von Variation-Tree bzw. Variation-Diff kennengelernt, die sich sehr ähnlich sind aber auch einen Unterschied haben. Dieser Unterschied ist die Ordnung der Kinderknoten, welche die geordneten Variation-Trees und Variation-Diffs haben und die normalen Variation-Trees und Variation-Diffs nicht. Diese Ordnung ist für das Unparsen notwendig. Dazu haben wir gezeigt das es möglich ist, geordnete Variation-Tree bzw. Variation-Diff in Variation-Tree bzw. Variation-Diff zu überführen. Die notwendige Reduktion und Projektion kann, wie in Abb. 3.1 dargestellt, in beliebiger Reihenfolge angewandt werden. Im späteren Verlauf können wir bei Bedarf diese Unterschiede für unsere Zwecke verwenden. Nachdem wir jetzt wissen was Variation-Trees und Variation-Diffs sind, ist es an der Zeit nachzuvollziehen wie diese gebildet werden.

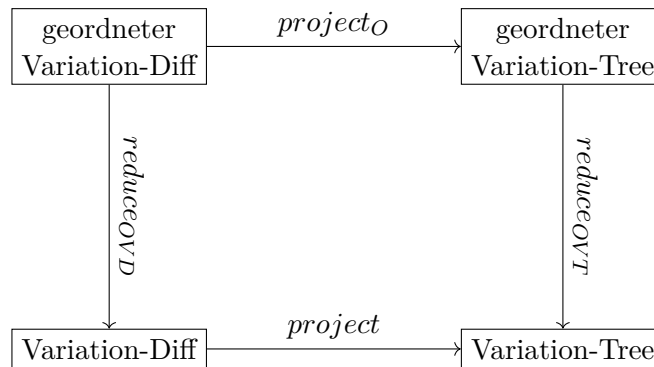


Abbildung 3.1: Transformationen von geordneten Variation-Diff , geordneten Variation-Tree, Variation-Diff und Variation-Tree

3.2 Parser

Jetzt beschäftigen wir uns mit den Parsen, also wie Variation-Trees bzw. Variation-Diffs aus mit C-Präprozessor-Direktiven annotiertem Code bzw. einem textbasierten Diff von solchem Code erstellt werden. Das Verständnis des Parsens ist, für das Verständnis des Unparsens von Bedeutung, da das Unparsen das Parsen invertiert. Dazu schauen wir uns den Parser-Algorithmus von Viegener [31] an, welcher das Parsen von Variation-Diffs aus textbasierten Diffs eingeführt hat.

Der unten stehender Algorithmus überführt einen textbasierten Diff in einen Variation-Diff. Der Algorithmus wurde, außerdem das wir den ins Deutsche übersetzt haben, so wie er ist von Viegener [31] übernommen. Der Dabei werden in dem Algorithmus einige Funktionen verwendet, die nicht so in der Definition vorkamen. Eine dieser Funktionen ist der Code-Typ, welcher die Rolle einer Zeile im Diff repräsentiert. Es kann die Werte *if*, *elif*, *else*, *code*, oder *endif* haben. Bei dem Wert *if* ist gegeben, dass die Zeile eine der Präprozessor Anweisungen *#if*, *#ifdef*, oder *#ifndef* hat. Bei dem Wert *else* ist in der Zeile die Präprozessor Anweisungen *#else*, bei Wert *elif* ist die Präprozessor Anweisungen *#elif* und bei Wert *endif* ist die Präprozessor Anweisungen *#endif* gegeben. Wenn der Wert von Code-Typ *code* ist, dann enthält die Zeile keine Präprozessor Anweisungen, sondern normalen Code. Da Wert *elif* als Erweiterung betrachtet werden kann, wird auf sie nicht weiter in unserer Arbeit eingegangen. Der Code-Typ einer Zeile wird bei der Erstellung eines neuen Knotens in Knotentyp τ überführt. Der Code-Typ *if* wird in τ gleich *MAPPING*, der Code-Typ *code* wird in τ gleich *ARTIFACT* und der Code-Typ *else* wird in τ gleich *ELSE* überführt. Der Code-Typ *endif* hat keine Überführung in Knotentyp τ , diese Code-Type wird nur intern von dem Algorithmus verwendet. Eine andere Funktion ist der Diff-Typ,

welcher für eine Zeile angibt, ob diese Zeile hinzugefügt wurde, entfernt wurde oder unverändert geblieben ist. Der Diff-Typ $+$ sagt, dass die Zeile hinzugefügt wurde, $-$ sagt, dass die Zeile entfernt wurde und \circ sagt, dass die Zeile unverändert geblieben ist. Der Diff-Typ wird, auch wie der Code-Typ, bei der Erstellung eines neuen Knotens in Δ überführt. Dabei wird $+$ in $+$, $-$ in $-$ und \circ in \circ überführt. Der Variation-Diff hat noch einen Knoten welcher keine Widerspiegelung in dem textbasierten Diff enthält, das ist der Wurzelknoten. Der Wurzelknoten repräsentiert den ganzen textbasierten Diff. Er hat als einziger Knoten im Variation-Diff keinen Elternknoten.

Algorithmus 1: Erstellung eines Variation-Diffs aus einem Patch

Data: ein textbasierter Diff
Result: ein Variation-Diff

```

1  erstelle den Wurzelknoten
2  initialisiere ein Stack/Keller before mit dem Wurzelknoten
3  initialisiere ein Stack/Keller after mit dem Wurzelknoten
4
5  foreach Zeile in dem Patch/Diff do
6       $\delta \leftarrow$  identifiziere Diff-Typ der Zeile
7       $\gamma \leftarrow$  identifiziere Code-Typ der Zeile
8       $\sigma \leftarrow$  identifiziere relevante Stacks mithilfe von  $\delta$ 
9
10     if  $\gamma = \text{endif}$  then
11         Entferne, solange Knoten von allen Stacks in  $\sigma$ , bis if-Knoten entfernt wurde
12     else
13         erstelle einen neuen Knoten mit  $\delta$ ,  $\gamma$  und gerichtete Kanten von Elternknoten
            aus  $\sigma$  zu dem neu erstellten Knoten
14         if  $\gamma \neq \text{code}$  then
15             füge den neuen Knoten  $\sigma$  hinzu
16         end
17     end
18 end
  
```

Der Algorithmus arbeitet wie folgt. Ganz am Anfang wird der Wurzelknoten in Zeile 1 erstellt. Danach werden zwei Stacks erstellt und jeweils mit dem Wurzelknoten initialisiert, was in Zeilen 2 und 3 des Algorithmus 1 zu sehen ist. Die Stacks speichern dabei die Elternknoten. Ein Stack speichern die Elternknoten im davor Zustand und der anderer im danach Zustand. Beide Stacks werden mit Wurzelknoten initialisiert, welcher den ganzen Diff repräsentiert und deshalb als einziger Knoten in Variation-Diff keinen Elternknoten hat. In Zeile 5 ist eine Schleife zu sehen, welche über alle Zeilen des textbasierten Diffs geht. Dabei wird für jede Zeile zuerst der Diff-Typ δ in Zeile 6 und dann der Code-Typ γ in Zeile 7 ermittelt. In Zeile 8 werden die relevanten Stacks σ anhand von Diff-Typ δ bestimmt, und zwar wie folgt:

$$\sigma = \begin{cases} \text{Stack } \textit{after} & , \quad \delta = \text{add} \\ \text{Stack } \textit{before} & , \quad \delta = \text{remove} \\ \text{Stacks } \textit{before} \text{ und } \textit{after} & , \quad \delta = \text{none} \end{cases}$$

Diese Informationen werden zum einen dazu gebraucht für Algorithmus interne Berechnungen und zum anderen zur Erstellung von Knoten gebraucht. Danach in Zeile 10 kommen wir zu einer if-Abfrage. Wenn der Code-Typ der bearbeiteten Zeile *endif* entspricht, dann wird aus den relevanten Stacks in σ solange Knoten entfernt bis man einen Knoten mit dem Code-Type γ *if* entfernt hat. Falls beide Stacks relevant sind, muss der if-Knoten in beiden Stacks gefunden werden. Dieses Vorgehen ist notwendig, da eine *endif*-Anweisung auf das Ende des dazugehörigen

if-Blocks oder if-else-Blocks folgen muss. Das führt mit sich, dass die dazugehörigen if-Knoten und else-Knoten keine Elternknoten mehr sein können und aus den Stacks entfernt werden müssen. Wenn der Code-Typ nicht *endif* entspricht, kommen wir in den else-Teil ab Zeile 12 des Algorithmus 1. Dort wird zuerst ein neuer Knoten erstellt, dazu unter anderem wird Diff-Typ δ und Code-Typ γ verwendet. Es werden auch Kanten von Elternknoten aus den Stacks von σ zu diesen neuen Knoten erstellt. Als Nächstes wird in Zeile 14 überprüft, ob der erstellte Knoten nicht von Code-Typ *code* ist. Diese Abfrage ist nötig da nur solche Knoten ein Elternknoten sein können. Den Code-Typ *endif* kann dieser Knoten nicht haben, wegen der if-Abfrage in Zeile 10, welche nicht zulässt, dass ein Knoten mit diesem Typ zu dieser Stelle gelangen kann. Wenn der Knoten nicht von Code-Typ *code* ist, dann wird dieser Knoten den relevanten Stacks aus σ hinzugefügt, sonst wenn der Knoten, den Code-Type *code* hat, wird nichts gemacht.

Der vorgestellte Algorithmus ist für das Parsen von textbasierten Diffs, welche aus mit C-Präprozessor-Annotierten Code entstanden sind, zu Variation-Diffs ausgelegt aber es ist auch möglich den Algorithmus zum Parsen von C-Präprozessor-Annotiertem Code zu einem Variation-Tree zu verwenden. Wir reduzieren das Problem ein Variation-Tree zu parsen auf das Problem ein Variation-Diff zu parsen. Hierzu müssen wir zwei Sachen beachten. Zuerst wäre da die Anpassung der Eingabe, da wir C-Präprozessor-Annotierten Code haben aber der Algorithmus einen textbasierten Diff erwartet. Die zweite Sache wäre die Anpassung der Ausgabe, die Ausgabe des Algorithmus ist ein Variation-Diff, wir brauchen aber einen Variation-Tree. Um die Eingabe gerecht für den Algorithmus zu machen, müssen wir unseren C-Präprozessor-Annotierten Code in ein textbasiertes Diff verwandeln. Dazu bilden wir ein Diff mit unserem C-Präprozessor-Annotierten Code als Davor-Zustand und Danach-Zustand. Danach bekommen ein textbasiertes Diff in dem jede Zeile als unverändert markiert ist. Dabei hat jede Zeile dieses Diffs den Diff-Typ *none*. Da jetzt ein Diff gegeben ist, können wir auf den Diff den Algorithmus anwenden. Die Ausgabe ist dann ein Variation-Diff, welcher in ein Variation-Tree umgewandelt werden muss. Um dies anzustellen, bilden wir eine Projektion des Variation-Diffs auf den Davor- bzw. Danach-Zustand und bekommen einen Variation-Tree. Es ist irrelevant welcher von den beiden Zuständen genommen wird, da der Davor-Zustand gleich dem Danach-Zustand sein soll. Mit den gezeigten Zwischenschritten lässt sich dieser Algorithmus auch für das Parsen von C-Präprozessor-Annotierten Code zu Variation-Trees verwenden.

Wir wollen die Arbeitsweise des Algorithmus veranschaulichen. Dazu wenden wir den Algorithmus, auf das untenstehende, beispielhafte Stück C-Code mit C-Präprozessor-Annotationen an und veranschaulichen das Vorgehen in der Abbildung 3.3. Diese Abbildung zeigt Kasten und den Zeitpunkt der dort gezeigten Information. In einem Kasten ist der Variation-Diff und relevanten Stacks für angegebene Zeile aus dem C-Code unten enthalten. Der Zeitpunkt ist nach Bearbeitung der angegebenen Zeile. Hier ist ein C-Code gegeben. Wir brauchen aber ein textbasierten Diff. Der gezeigte C-Code wird wie gerade beschrieben in ein textbasiertes Diff überführt, somit ist die nötige Eingabe gegeben. Am Anfang des Algorithmus werden die Stacks erstellt und mit dem Wurzelknoten initialisiert. Wir betrachten jetzt die Schleife, die über alle Zeilen des obigen Diffs geht. Wir kommen zur Zeile 1 des C-Codes, dort befindet sich eine normale Codezeile, welche nicht annotiert ist. Es ergibt sich, dass diese Zeile den Code-Typ *code* und den Diff-Typ *none* hat. Alle anderen Zeilen haben auch den Diff-Typ *none*, aus dem Grund wie dieser Diff gebildet wurde und deshalb lassen wir die Erwähnung des Diff-Typs für jede Zeile sein. Dasselbe gilt auch für die relevanten Stacks in σ , da alle Zeilen den Diff-Typ *none* haben, gilt für alle Zeilen auch die gleichen relevanten Stacks und das sind beide. Da diese Zeile nicht Code-Typ *endif* hat, wird ein Knoten mit Code-Type, Diff-Typ, Elternknoten aus den Stacks und dem Inhalt der Zeile erstellt und dem Variation-Diff hinzugefügt, wie das aussieht, ist in

```

1  f();
2  #if(A)
3      #if(B||C)
4          g();
5      #else
6          z();
7  #endif
8  x();
9  #endif

```

Abbildung 3.2: Beispiel für mit C-Präprozessor-Annotierten Code

Abbildung 3.3 in dem Kasten „Nach Z.1“ zu sehen. Der erstellte Knoten hat als Elternknoten den Wurzelknoten, wie es in den Stacks zu sehen ist. In Abbildung 3.1 im Kasten „Nach Z.2“ ist der Variation-Diff und die relevanten Stacks nach der Bearbeitung der Zeile 2 zu sehen. Es wurde ein neuer Knoten erstellt, welcher eine *if*-Anweisung enthält und in beide Stacks wurden dieser Knoten hinzugefügt. Die Schleife wurde fast gleich wie im vorherigen Fall durchgelaufen, außer an der letzten *if*-Abfrage. Diese Abfrage war bei der Zeile 1 *false* dieses Mal, da wir keinen Code-Typ *code* haben, wird diese Abfrage ausgeführt und der neu erstellte Knoten den Stacks hinzugefügt. Der nächste Kasten rechts zeigt den Variation-Diff nach Zeile 3. Der Algorithmus ist genauso wie im vorherigen Fall vorgegangen. Weiter voran wird dem Variation-Diff im nächsten Schritt ein Code-Knoten hinzugefügt, da für die Erstellung dieses Knotens der Code selbst irrelevant ist, wurde hier genauso vorgegangen wie bei der Erstellung eines Code-Knotens in Zeile 1. In der Zeile 5 ist **#else** als Anweisung gegeben. Diese Zeile hat den Code-Typ *else* und somit auch kein *endif*. Es wird in den *else*-Zweig der ersten Abfrage gegangen. Dort wird ein neuer Knoten mit dem Inhalt dieser Zeile erstellt. Der Knoten wird den Stacks hinzugefügt, da der Knoten *else* als Code-Typ hat und nicht *code*, was die innere Abfrage erfüllt (Abbildung 3.3 „Nach Z.5“). In der nächsten Zeile ist wieder eine Codezeile vorhanden und aus der ein Code-Knoten erstellt wird. Wie es danach aussieht, ist in Abbildung 3.3 „Nach Z.6“ zu sehen. Danach in der Zeile 7 treffen wir das erste Mal auf die Anweisung **#endif**, welche den Code-Typ *endif* hat. Damit gelangen wir in den *if*-Teil der ersten Abfrage, welcher sich in der Zeile 11 des Algorithmus 1 befindet. Der Algorithmus entfernt nun von beiden Stacks jeweils so lange Knoten, bis ein *if*-Knoten entnommen wurde. Dabei werden aus den Stacks die Knoten mit *else* und *if*(B||C) entnommen und übrig bleiben der *if*(A) Knoten und der Wurzelknoten. Dieser Schritt verändert den Variation-Diff nicht. Im nächsten Schritt treffen wir wieder auf eine Codezeile und erstellen einen Knoten und fügen den dem Variation-Diff hinzu, welcher in Abbildung 3.3 „Nach Z.8“ zu sehen ist. In der Zeile 9 ist wieder ein **#endif** und wir müssen wieder Knoten aus den Stacks entnehmen. Dieses Mal wird der Knoten *if*(A) entnommen und es bleibt nur der Wurzelknoten übrig. Damit wäre die Arbeit des Algorithmus zu Ende und wir haben als Rückgabewert einen Variation-Diff erhalten. Wir brauchen einen Variation-Tree statt des Variation-Diffs welches wir bekommen haben. Dazu müssen wir eine Projektion auf den Zustand davor oder danach bilden. Schlussendlich erhalten wir einen Variation-Tree, welcher genauso aufgebaut ist, wie der Variation-Diff aus der Abbildung 3.3 Kasten „Nach Z.8“.

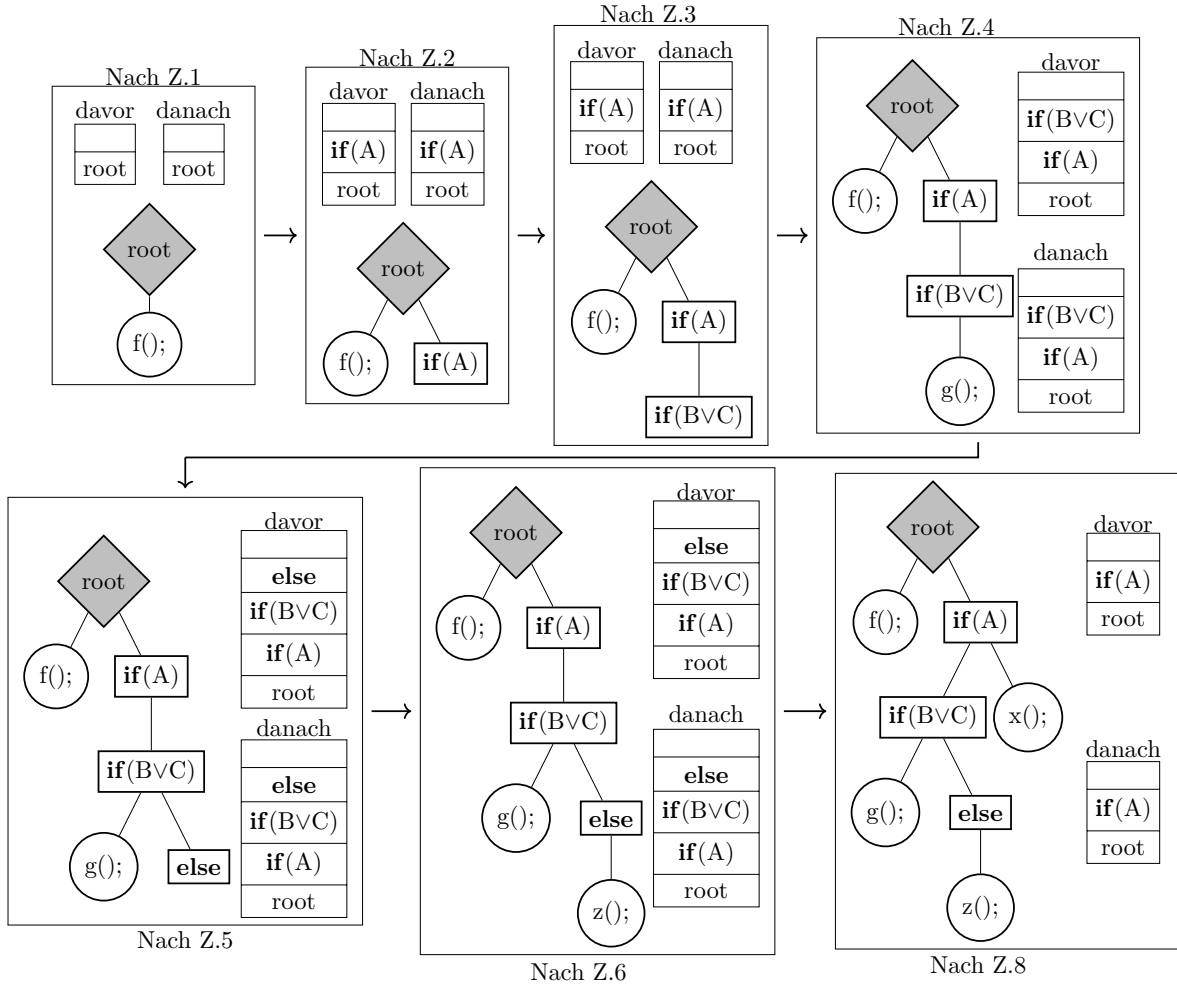


Abbildung 3.3: Beispiel für den Parsen Algorithmus von Viegener

Dieser Algorithmus kann sowohl gewöhnliche Variation-Diffs nach Definition 3.2 als auch geordneten Variation-Diffs nach Definition 3.5 erstellen. Das ist möglich, da die Reihenfolge der Kindknoten in Zeile 12 festgelegt wird.

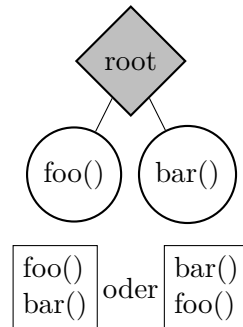
Jetzt wissen wir wie der Parser von Viegener funktioniert und der Parser auch für mehr als nur eine Definition von Variation-Diff zu gebrauchen ist. Dies können wir uns zunutze machen, wenn wir eine eigene Definition von Variation-Diff und Variation-Tree ausarbeiten werden. Es ist nun an der Zeit zu sehen, welche Informationen durch das Parsen und der Definitionen nach verloren gehen und wie wir diese Informationen zurückbekommen können.

3.3 Verlorengelung

Nachdem wir uns mit dem Beschäftigt haben, was Variation-Trees und Variation-Diffs sind und wie dieser aus mit C-Präprozessor-Annotierten Code und textbasierten Diffs erzeugt werden, setzen wir uns damit auseinander, welche Informationen während des Parsens verloren gehen. Wir beschäftigen uns mit der Frage, welche Informationen sind im mit C-Präprozessor-Annotierten Code bzw. textbasierten Diff erhalten, aber nach dem Parsen nicht in Variation-Tree bzw. Variation-Diff zu finden sind und wie wir diese Informationen zurückbekommen können,

um das Unparse zu bewerkstelligen. Die verlorengehende Informationen sind die Ordnung der Zeilen, die Position von Zeilen mit einem `#endif` innerhalb von textbasierten Diff bzw. mit C-Präprozessor-Annotierten Code und der Inhalt aller Zeilen inklusive der Zeilen mit `#endif` und `#else`.

Der Definition von Variation-Tree bzw. Variation-Diff nach haben die Kinder eines Knotens keine Reihenfolge. Deshalb können wir während des Unparsen nicht ermitteln, welcher Knoten zuerst kommt und welcher danach. Zum Beispiel ein Variation-Tree kann mehrdeutig verstanden wäre, was für uns nicht zu gebrauchen ist.



Um das Problem zu umgehen verwenden wir statt der normalen Definition von Variation-Tree und Variation-Diff die Definition von geordneten Variation-Tree und geordneten Variation-Diff. Diese Definition hat eine Ordnung bei den Kinderknoten. Deshalb kann sie uns eine eindeutige Reihenfolge geben, so dass keine Mehrdeutigkeiten in Bezug auf das wie die Knoten eingeordnet sind vorkommt. Das Umsetzen dieser Definition von dem Parser stellt für uns keine Schwierigkeiten dar. Die Definitionen von normalen Variation-Tree bzw. Variation-Diff unterscheiden sich von geordneten Variation-Tree bzw. geordneten Variation-Diff nur um die Ordnung O . Dies hat zur Folge, dass alles andere so wie gewohnt umgesetzt werden kann. Die Ordnung selbst muss dann gesetzt werden, wenn ein neuer Knoten entsteht und er seine Eltern bekommt. Dies findet in der Zeile 12 des Algorithmus 1 statt. Dort stehen keine genaueren Angaben zu Erstellung des Knotens, also kann diese Zeile auch um die Setzung der Ordnung erweitert werden. Damit haben wir unseren ersten Informationsverlust beseitigt.

Als nächstes beschäftigen wir uns mit dem Verlust der Position von `#endif`. Es gibt keinen Knoten innerhalb von Variation-Trees und Variation-Diffs, welcher `#endif` repräsentiert. Deshalb wissen wir nicht, wo sich die Zeilen mit `#endif` befinden müssen wenn wir die Variation-Trees und Variation-Diffs unparse. Zu Beschaffung dieser Information muss die Position von `#endif` aus der gegebenen Information rekonstruiert werden. In dieser Beschreibung wird einfachheitshalber so gehandelt, als gäbe es einen Knoten für `#endif`. Um zu bestimmen, ob das letzte Kind eines Knotens ein `#endif` sein muss, müssen wir prüfen ob der Knotentyp $\tau = \text{MAPPING}$ ist. Wenn das zutrifft, wird ein Dummyknoten als neuer letzter Kindknoten zu diesem Knoten hinzugefügt. Entsprechend wissen wir, dass die Zeile mit dem `#endif` nach dem Inhalt allen anderen Kinderknoten und dessen Unterbäumen kommt. Damit haben wir auch eine Möglichkeit diesen Informationsverlust zu beseitigen.

Als letztes müssen wir uns damit beschäftigen wie wir den Inhalt aller Zeilen aus Variation-Tree oder Variation-Diff bekommen können. Bei Zeilen mit `#endif` und `#else` wäre es möglich auf exakte Gleichheit zum ursprünglichen Text zu verzichten und die entsprechenden Knoten mit vorher festgelegten Textbausteinen zu ersetzen. Der Grund dafür ist, dass die Einrückung von `#endif` und `#else` im Falle von C-Präprozessor nicht dazu führt dass Code fehlerhaft wird.

Es wäre möglich eine Heuristik, wie z.B. das `#endif` und `#else` so weit eingerückt sind, wie der dazugehörige `if` oder das `#else` und `#endif` immer am Zeilenanfang sind, zu verwenden. Bei der Verwendung einer Heuristik, können wir nicht die exakte Gleichheit garantieren, erfordert aber benötigen auch keinen zusätzlichen Aufwand oder Speicher. Jedoch könnten Kommentare in diesen Zeilen nicht wiederhergestellt werden. Bei Knoten mit τ gleich ARTIFACT, könnte man so vorgehen, das man in Label die Zeile abspeichert. Der Definition von Label nach werden dort entweder ein Implementierungsartefakt oder eine aussagenlogische Formel gespeichert. Ein Implementierungsartefakt kann dabei mehr als nur die Codezeile sein. Ein Implementierungsartefakt ist dabei eine identifizierbare Einheit mit beliebiger Granularität innerhalb eines Softwareprojekts [3]. Unsere Arbeit ist darauf ausgelegt, dass wir einen Unparser bereitstellen, welcher aus Variation-Trees bzw. Variation-Diffs mit C-Präprozessor-Annotierten Code bzw. textbasiertes Diff erstellt. Wir könnten fordern, dass das Label als Implementierungsartefakt die Codezeile abspeichert. Für Knoten mit τ gleich ELSE oder MAPPING würde das aber nicht funktionieren. Für Knoten mit τ gleich ELSE wird der Definition nach überhaupt kein Label gespeichert. Für Knoten mit τ gleich MAPPING wird das Label eine aussagenlogische Formel speichern, aber die Bedingungen in der C-Präprozessor-Annotation sind nicht immer eine aussagenlogische Formel und muss deshalb durch Boolean-Abstraction [3] in solche umgewandelt werden. Das kann z.B. so aussehen: `#if A(x) > 3` ist gegeben und nach boolean abstraction sieht es dann so aus `#if A__LB__x__RB__GT__3`. Dabei ist eine zurück Umwandlung nicht garantiert, da wir nicht sicher sein können das z.B `A__LB__x__RB__GT__3` nicht als ein Variablenname gewählt wurde und damit keine Umwandlung benötigt. Es ergibt sich, dass das Label nur für τ gleich ARTIFACT gut die Zeile speichern kann, bei den anderen Zeilen gibt es Schwierigkeiten. Um diese Informationen auch im Variation-Tree und Variation-Diff zu haben, müssen wir diese Informationen explizit speichern. Die Definitionen von Variation-Tree, Variation-Diff, geordneten Variation-Tree und geordneten Variation-Diff sehen aber dafür nichts vor. Deshalb müssen wir die Definition von Variation-Diff bzw. geordneten Variation-Tree und Variation-Diff bzw. geordneten Variation-Diff erweitern. Obwohl wir für τ gleich ARTIFACT die Zeile in dem Label speichern können, werden wir wegen der Einheitlichkeit alle Zeile in der Erweiterung speichern und die nicht unterscheiden. Wir erweitern die Definitionen wie folgt:

Definition 3.11. *Ein SPEICHERNDER VARIATION-TREE (V, E, r, τ, l, M) ist für V, E, r, τ und l so definiert wie in der Definition 3.1 und $M : V \rightarrow (String, String)$ speichert in der ersten Stelle des Tupels die Zeile, welche der Knoten darstellt und an der zweiten Stelle wird für Knoten mit τ gleich MAPPING die dazugehörige Zeile mit `#endif` gespeichert, für die restlichen Knoten wird dort Null gespeichert.*

Definition 3.12. *Ein SPEICHERNDER VARIATION-DIFF $(V, E, r, \tau, l, \Delta, M)$ ist für V, E, r, τ, l und Δ so definiert wie in der Definition 3.2 und $M : V \rightarrow (String, String)$ speichert in der ersten Stelle des Tupels die Zeile, welche der Knoten darstellt und an der zweiten Stelle wird für Knoten mit τ gleich MAPPING die dazugehörige Zeile mit `#endif` gespeichert, für die restlichen Knoten wird dort Null gespeichert.*

Definition 3.13. *Ein GEORDNETER, SPEICHERNDER VARIATION-TREE (V, E, r, τ, l, O, M) ist für V, E, r, τ, l und O so definiert wie in der Definition 3.5 und $M : V \rightarrow (String, String)$ speichert in der ersten Stelle des Tupels die Zeile, welche der Knoten darstellt und an der zweiten Stelle wird für Knoten mit τ gleich MAPPING die dazugehörige Zeile mit `#endif` gespeichert, für die restlichen Knoten wird dort Null gespeichert.*

Definition 3.14. *Ein GEORDNETER, SPEICHERNDER VARIATION-DIFF $(V, E, r, \tau, l, O_{before}, O_{after}, M)$ ist für $V, E, r, \tau, l, O_{before}$ und O_{after} so definiert wie in der Definition 3.6 und $M : V \rightarrow (String, String)$ speichert in der ersten Stelle des Tupels die Zeile, welche der Knoten darstellt*

und an der zweiten Stelle wird für Knoten mit τ gleich `MAPPING` die dazugehörige Zeile mit `#endif` gespeichert, für die restlichen Knoten wird dort `Null` gespeichert.

Dazu nachdem wir die neuen Variation-Trees und Variation-Diffs definiert haben, müssen wir noch entsprechende Projektionen definieren.

Definition 3.15. Die Projektion $project_M(D, t)$ für ein speicherndes Variation-Diff ist das Entfernen von Δ , M und der Knoten und Kanten, welche zu der Zeit t nicht vorhanden sind.
 $project_M((V, E, r, \tau, l, \Delta, M), t) := (\{v \in V \mid EXISTS(t, \Delta(v))\}, \{e \in E \mid EXISTS(t, \Delta(e))\}, r, \tau, l, M)$

Definition 3.16. Die Projektion $project_{OM}(D, t)$ für ein geordnetes, speicherndes Variation-Diff ist das Entfernen von Δ , M und der Knoten und Kanten, welche zu der Zeit t nicht vorhanden sind. Dazu wird nur der Zeit entsprechende Ordnung O_t behalten.

$project_{OM}((V, E, r, \tau, l, \Delta, O_{before}, O_{after}, M), t)$
 $:= (\{v \in V \mid EXISTS(t, \Delta(v))\}, \{e \in E \mid EXISTS(t, \Delta(e))\}, r, \tau, l, O_t, M)$

Da wir die Definitionen erweitert haben, ist es möglich die speichernden Variation-Trees bzw. Variation-Diffs in normale Variation-Trees bzw. Variation-Diffs und geordnete, speichernden Variation-Trees bzw. Variation-Diffs in geordnete Variation-Trees bzw. Variation-Diffs umzuwandeln.

Definition 3.17. $reduce_{MVT}$ wandelt einen speichernden Variation-Tree (Definition 3.11) in einen normale Variation-Tree (Definition 3.1)

$reduce_{MVT}((V, E, r, \tau, l, M)) := (V, E, r, \tau, l).$

Definition 3.18. $reduce_{MVD}$ wandelt einen speichernden Variation-Diff (Definition 3.12) in einen normalen Variation-Diff (Definition 3.2)

$reduce_{MVD}((V, E, r, \tau, l, \Delta, M)) := (V, E, r, \tau, l, \Delta)$

Definition 3.19. $reduce_{MOVT}$ wandelt einen geordneten, speichernden Variation-Tree (Definition 3.13) in einen geordneten Variation-Tree (Definition 3.5)

$reduce_{MOVT}((V, E, r, \tau, l, O, M)) := (V, E, r, \tau, l, O)$

Definition 3.20. $reduce_{MOVD}$ wandelt einen geordneten, speichernden Variation-Diff (Definition 3.14) in einen geordneten Variation-Diff (Definition 3.6)

$reduce_{MOVD}((V, E, r, \tau, l, \Delta, O_{before}, O_{after}, M)) := (V, E, r, \tau, l, \Delta, O_{before}, O_{after})$

Es bleibt uns nur noch zu zeigen das die Reihenfolge der Anwendung von $project$ und $reduce$ irrelevant ist. Dann zusammen mit den gezeigten aus Kapitel 3.1 ist die Abbildung 3.4 gegeben, welche eine Erweiterung der Abbildung 3.1 ist. Die Abbildung 3.4 veranschaulicht die Transformationen von geordneten speichernden Variation-Trees, geordneten speichernden Variation-Diffs, speichernden Variation-Trees, speichernden Variation-Diffs, geordneten Variation-Trees, geordneten Variation-Diffs, normalen Variation-Trees und normalen Variation-Diffs.

Lemma 3.21. Für einen speichernden Variation-Diff D und eine Zeit $t \in \{before, after\}$ gilt $reduce_{MVT}(project_M(D, t)) = project(reduce_{MVD}(D), t)$.

Beweis. $reduce_{MVT}(project_M(D, t))$
 $= reduce_{MVT}(project_M((V, E, r, \tau, l, \Delta, M), t))$
 $= reduce_{MVT}((\{v \in V \mid EXISTS(t, \Delta(v))\}, \{e \in E \mid EXISTS(t, \Delta(e))\}, r, \tau, l, M))$
 $= (\{v \in V \mid EXISTS(t, \Delta(v))\}, \{e \in E \mid EXISTS(t, \Delta(e))\}, r, \tau, l)$
 $= project((V, E, r, \tau, l, \Delta), t)$
 $= project(reduce_{MVD}(V, E, r, \tau, l, \Delta, M), t)$
 $= project(reduce_{MVD}(D), t)$

□

Lemma 3.22. Für einen geordneten, speichernden Variation-Diff D und eine Zeit $t \in \{before, after\}$ gilt $reduce_{MOVT}(project_{OM}(D, t)) = project_O(reduce_{MOVD}(D), t)$.

Beweis. $reduce_{MOVT}(project_{OM}(D, t))$
 $= reduce_{MOVT}(project_{OM}((V, E, r, \tau, l, \Delta, O_{before}, O_{after}, M), t))$
 $= reduce_{MOVT}(\{\{v \in V | EXISTS(t, \Delta(v))\}, \{e \in E | EXISTS(t, \Delta(e))\}, r, \tau, l, O_t\})$
 $= (\{v \in V | EXISTS(t, \Delta(v))\}, \{e \in E | EXISTS(t, \Delta(e))\}, r, \tau, l, O_t)$
 $= project_O((V, E, r, \tau, l, \Delta, O_{before}, O_{after}), t)$
 $= project_O(reduce_{MOVD}(V, E, r, \tau, l, \Delta, O_{before}, O_{after}, M), t)$
 $= project_O(reduce_{MOVD}(D), t)$

□

Nach dem Ganzen ergibt sich, ein folgendes Schaubild:

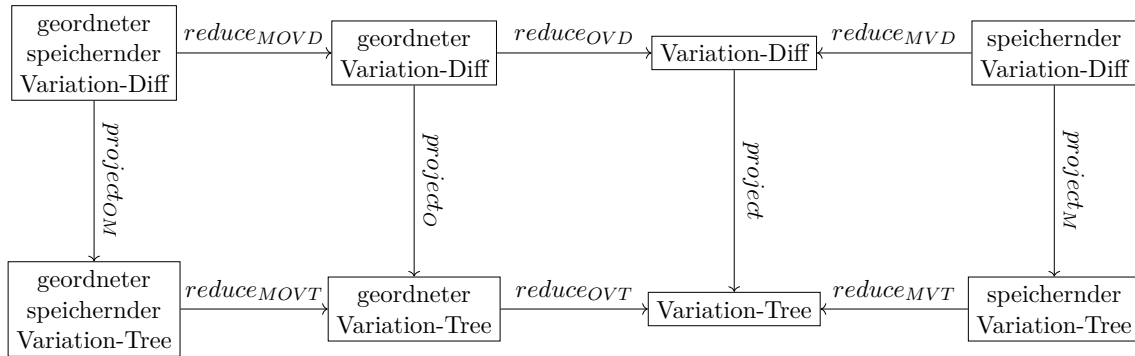


Abbildung 3.4: Transformationen von verschiedenen Variation-Trees und Variation-Diffs

Nachdem wir uns mit den neuen Definitionen beschäftigt haben, müssen wir noch den Parsen anpassen. Damit der Parser auch diese Definitionen umsetzen kann. Dazu muss man die Zeile 11 des Algorithmus etwas erweitern, damit der Algorithmus sich den if-Knoten merkt. Dazu muss man noch eine neue Zeile zum Algorithmus zufügen und das ist die Zeile 12 des Algorithmus 2. Dort wird ein Tupel mit dazugehöriger Zeile des Knotens und der Zeile, mit dem `#endif`, erstellt. Dann wird diese Tupel in M für den gemerkten Knoten gespeichert. Alle anderen Zeilen werden für ihre Knoten jeweils in Zeile 15 des Algorithmus 2 gespeichert. Dort wird ein Tupel mit der Zeile und Null erstellt und für den neu erstellten Knoten gespeichert. Damit sind wir in der Lage die neuen Definitionen auch aus mit C-Präprozessor-Annotierten Code und textbasierten Diffs zu erzeugen.

Algorithmus 2: Erstellung eines Variation-Diffs, welcher sich `#else` und `#endif` merkt, aus einem Patch

Data: ein textbasierter Diff
Result: ein Variation-Diff

```

1  erstelle den Wurzelknoten
2  initialisiere ein Stack/Keller before mit dem Wurzelknoten
3  initialisiere ein Stack/Keller after mit dem Wurzelknoten
4
5  foreach Zeile in dem Patch/Diff do
6       $\delta \leftarrow$  identifiziere Diff-Typ der Zeile
7       $\gamma \leftarrow$  identifiziere Code-Typ der Zeile
8       $\sigma \leftarrow$  identifiziere relevante Stacks mithilfe von  $\delta$ 
9
10     if  $\gamma = \text{endif}$  then
11         Entferne, solange Knoten von allen Stacks in  $\sigma$ , bis ein if-Knoten v entfernt
            wurde
12          $M(v) \leftarrow (M(v)[0], \text{Zeile})$       /* Die erste Stelle von M wird mit dem dort
            enthaltenen Wert wiederbesetzt und die zweite Stelle speichert
            die Zeile mit dem endif */
13     else
14         erstelle einen neuen Knoten v mit  $\delta$ ,  $\gamma$  und gerichtete Kanten von Elternknoten
            aus  $\sigma$  zu dem neu erstellten Knoten
15          $M(v) \leftarrow (\text{Zeile}, \text{Null})$       /* Die erste Stelle von M speichert die
            bearbeitete Zeile und die zweite Stelle wird mit null besetzt */
16         if  $\gamma \neq \text{code}$  then
17             füge den neuen Knoten  $\sigma$  hinzu
18         end
19     end
20 end
```

Mit der Erweiterung der Definitionen können wir den Inhalt der Zeile mit `#endif` oder `#else` bekommen und damit ist dieser Informationsverlust auch beseitigt.

Nachdem wir herausgefunden haben, welche für das Unparsen relevante Information verloren geht und wie man diese Information zurück bekommen kann, sind wir in der Lage das alles zu nutzen, um einen Algorithmus zum Unparsen zu entwickeln.

3.4 Unparsing

Wir haben festgestellt welche Informationen während des Parsens verloren gehen und wie diese zurückzubekommen sind. Nach diesem Schritt sind wir in der Lage das Unparsen zu bewerkstelligen. Darüber geht es in diesem Unterkapitel. Wir stellen unseren Algorithmus für das Unparsen von speichernden, geordneten Variation-Trees und ein Vorgehen zum Unparsen von speichernden, geordneten Variation-Diffs.

Unser Algorithmus ist der Algorithmus 3 und ist nur für das Unparsen von speichernden, geordneten Variation-Tree zu einem mit C-Präprozessor-Annotierten Code bestimmt. Der Algorithmus basiert auf der Tiefensuche. Ein Variation-Tree ist ein Baum, dessen Knoten die Zeilen des mit C-Präprozessor-Annotierten Codes enthalten. Dabei wenn ein Knoten Kinderknoten

hat, bedeutet es das dieser Knoten τ gleich MAPPING oder τ gleich ELSE hat. Dazu wissen wir dadurch, dass die Anweisungen im Kinderknoten von der Anweisung des Elternknotens eingeschlossen werden. Wegen so einer Anordnung ist die Verwendung von Tiefensuche von Vorteil. Die Tiefensuche besucht alle Knoten effizient und dazu geht die Tiefensuche zunächst ein Pfad vollständig in die Tiefe, bevor abzweigende Pfade beschriftet werden. Die Auswahl des nächsten Knotens welcher besucht wird muss etwas geändert werden, damit die Pfade der Reihenfolge nach beschriftet werden. Dazu werden die Kinderknoten eines Knotens den Stack in umgedrehter Reihenfolge hinzugefügt. Es ergibt sich, dass der letzte Kinderknoten auch als letzter drankommt und der erster als erster. Das zusammen ergibt, das unser Algorithmus die Knoten so besucht wie die in den Knoten enthaltene Zeilen angeordnet werden müssen.

Algorithmus 3: Unparsing eines Variation-Tree zu CPP-Annotiertem Code

Data: ein Variation-Tree (V, E, r, τ, l)
Result: mit C-Präprozessor-Annotierten Code

```

1 initialisiere einen leeren Stack/Keller stack
2 initialisiere einen String ergebnis
3 kinder  $\leftarrow \{v \in V \mid (r, v) \in E\}$ 
4 initialisiere ein Array array der Länge  $|kinder|$ 
5 for  $\forall v \in kinder$  do
6   | array[ $O(v)$ ]  $\leftarrow v$ 
7 end
8 for  $i = |kinder| \rightarrow 1$  do
9   | stack.push(array[i])
10 end
11 while stack nicht leer ist do
12   | knoten  $\leftarrow stack.pop()$ 
13   | if  $\tau(knoten) = mapping$  then
14     | erstelle einen Dummyknoten welcher  $M(knoten)[1]$  beinhaltet
15     | stack.push(Dummyknoten)
16   | end
17   | erweitere ergebnis mit dem String  $M(knoten)[0]$ 
18   | kinder  $\leftarrow \{v \in V \mid (knoten, v) \in E\}$ 
19   | initialisiere ein Array array der Länge  $|kinder|$ 
20   | for  $\forall v \in kinder$  do
21     | array[ $O(v)$ ]  $\leftarrow v$ 
22   | end
23   | for  $i = |kinder| \rightarrow 1$  do
24     | stack.push(array[i])
25   | end
26 end
27 return ergebnis

```

Jetzt schauen wir uns den Algorithmus an, nachdem wir sein Konzept besprochen haben. In Zeile 1 des Algorithmus 3 wird ein Stack initialisiert, so wie in der Tiefensuche. Danach in Zeile 2 wird ein String initialisiert, in dem am Ende der gesamte annotierte Code gespeichert wird. Von Zeile 3 bis Zeile 10 werden die Kinderknoten des Wurzelknotens auf den Stack gelegt. Das muss extra gemacht werden, da der Wurzelknoten keine Zeile enthält als Label, sondern true. Damit würde er das Ergebnis verfälschen. Genauer betrachtet wird folgendes gemacht, in der Zeile 3 werden die Kinderknoten des Wurzelknotens bestimmt und als Menge abgespeichert. Dann wird ein Array erstellt, dessen Länge gleich der Anzahl der Kinderknoten des Wurzelknotens

ist. Als Nächstes in Zeile 5 wird über alle Kinderknoten iteriert und die Kinderknoten gemäß ihrer Ordnung im Array abgespeichert. Zum Schluss werden die Kinderknoten in umgekehrter Reihenfolge den Stack hinzugefügt. In Zeile 11 beginnt eine while-Schleife, welche so lange läuft bis der Stack leer wird. In der Schleife wird als Erstes der oberste Knoten von Stack genommen und sich gemerkt, das ist in Zeile 12. Als Nächstes in Zeile 13 wird, geprüft ob τ des Knotens gleich MAPPING ist. Wenn das der Fall ist, dann wird zuerst ein Dummyknoten, welcher die Information zu `#endif` aus `M[1]` enthält, erstellt, dann wird dieser Dummyknoten dem Stack hinzugefügt. Der Dummyknoten hat τ gleich ARTIFACT und im `M[0]` diesen Dummyknotens wird der Inhalt aus `M[1]` gespeichert. Da dieser Knoten nur intern im Algorithmus verwendet wird, kann er solche Werte annehmen, die die Definition nicht vorsieht. Wenn das nicht der Fall ist, wird nichts gemacht und die Abfrage übersprungen. In Zeile 17 wird der Inhalt von `M[0]` String „ergebnis“ hinzugefügt. Ab Zeile 18 bis Zeile 25 werden die Kinderknoten des gerade bearbeiteten Knotens dem Stack in umgedrehter Reihenfolge hinzugefügt. Dort wird genauso vorgegangen wie in den Zeilen 3 bis 10. Damit ist der Inhalt der Schleife abgearbeitet und wir kommen zur Zeile 27, welche den String „ergebnis“ zurückgibt in dem der mit C-Präprozessor-Annotierter Code enthalten ist.

Nachdem ihr mehr über unser Algorithmus erfahren habt, wollen wir seine Arbeitsweise in der Abbildung 3.5 verdeutlichen. In diesem Beispiel werden wir das erhaltene Variation-Tree aus der Abbildung 3.3 unparse. Obwohl dort der Parser ein Variation-Tree erstellt, aber wir ein speichernden, geordneten Variation-Tree brauchen. Ein speichernder, geordneter Variation-Tree ist von dem Aussehen identisch, dem Variation-Tree aus Abbildung 3.3, wenn derselbe mit C-Präprozessor-Annotierter Code mithilfe von Algorithmus 2 geparkt wird. Die Abbildung 3.5 zeigt in kleineren Bildern, entweder den Zustand nach Bearbeitung des Wurzelknotens oder eines Schleifendurchlaufs, und die Reihenfolgen, in der die Bilder betrachtet werden sollen, beginnend mit dem Bild ganz links oben. Ein Bild enthält das speichernde, geordnete Variation-Tree, den Stack, die Ausgabe und welcher Knoten des speichernden, geordneten Variation-Trees bearbeitet wurde. Am Anfang im Bild ganz links oben der Abbildung 3.5 sind wir außerhalb der Schleife und es wird der Wurzelknoten abgearbeitet. Die Kinderknoten des Wurzelknotens werden dem Stack hinzugefügt. Die Ausgabe verändert sich nicht, da der Wurzelknoten keine Zeile enthält. In nächsten Bild sind wir in der Schleife. Der oberste Knoten wird von Stack genommen, der rote Pfeil zeigt auf den. Das ist ein Knoten mit τ gleich ARTIFACT, also wird die if-Abfrage in Zeile 13 des Algorithmus 3 mit falsch beantwortet und übersprungen. Danach in Zeile 17 wird `M[0]` des Knotens der Ausgabe hinzugefügt, was auch bei der Ausgabe im Bild zu sehen ist. Der Inhalt des Knotens ist gleich der ersten Zeile der Ausgabe. Der Knoten ist ein Blattknoten und hat keine Kinderknoten, welche den Stack hinzugefügt werden können. Im Bild danach, wird wieder der oberste Knoten aus dem Stack genommen, deshalb ist der Knoten auf den der rote Pfeil zeigt im Stack des vorherigen Bildes vorhanden und nicht in diesem. Für den betrachteten Knoten gilt τ ist gleich MAPPING. Aus diesem Grund gehen wir in die if-Abfrage aus der Zeile 13 des Algorithmus 3. Dort wird ein Dummyknoten mit `#endif` erstellt und dem Stack hinzugefügt was wir auch im Bild sehen. Außerhalb der Abfrage wird `M[0]` des Knotens wider der Ausgabe hinzugefügt. Die Kinderknoten werden dem Stack hinzugefügt. Im Bild ganz rechts oben wird wie immer der oberste Knoten aus dem Stack genommen welcher mit dem roten Pfeil markiert ist. Dieser Knoten hat τ gleich MAPPING. Es wird gleich wie mit den vorherigen Knoten vorgegangen. Es wird ein Dummyknoten erstellt und dem Stack hinzugefügt. Außerhalb der Abfrage wird wie immer der `M[0]` des Knotens der Ausgabe hinzugefügt. Die Kinderknoten werden dem Stack hinzugefügt. Jetzt sind wir im Bild ganz links in der Mitte. Dieser Knoten hat τ gleich ARTIFACT, deshalb wird hier genauso wie im oberen, zweiten Bild von links vorgegangen. Im nächsten Bild wird ein Knoten mit τ gleich ELSE bearbeitet. Dieser Knoten wurde aus dem Stack genommen. Der Knoten entspricht nicht τ gleich MAPPING, deshalb wird die if-Abfrage in

Zeile 13 des Algorithmus 3 übersprungen. Danach wird der Inhalt von $M[0]$ für diesen Knoten der Ausgabe hinzugefügt. Als Nächstes wird der einzige Kinderknoten dem Stack hinzugefügt. Im nächsten Bild wird wieder ein Knoten mit τ gleich ARTIFACT bearbeitet. Diese Knoten wird, genauso behandelt wie die anderen Knoten mit τ gleich ARTIFACT. Wir sind jetzt bei dem Bild in der Mitte ganz rechts. Dort wird zum ersten Mal ein Dummyknoten mit `#endif` bearbeitet. Der Knoten wird von Stack genommen. Es ist ein Dummyknoten und ist deshalb nicht in den gezeigten speichernde, geordneten Variation-Tree zu finden und es kann kein roter Pfeil auf den zeigen. Es wird geprüft ob der Dummyknoten τ gleich MAPPING hat, das ist nicht der Fall, da dieser Knoten τ gleich ARTIFACT hat. Wir überspringen die Abfrage aus Zeile 13 und gehen zu der Zeile 17 des Algorithmus 3. Dort wird der Inhalt von $M[0]$ für diesen Knoten der Ausgabe hinzugefügt. Dieser Knoten hat keine Kinderknoten, deshalb wird dem Stack auch nichts hinzugefügt. In dem Linken unteren Bild der Abbildung 3.5 wird der mit dem roten Pfeil markierter Knoten bearbeitet. Dieser Knoten hat auch τ gleich ARTIFACT und ist kein Dummyknoten. Deshalb wird dort auch wie in anderen Fällen vorgegangen und wir sehen nichts Neues. Im letzten Bild wird wieder ein Dummyknoten bearbeitet. Das Vorgehen ist in beiden Fällen dasselbe. Damit wären wir mit dem Unparsen fertig. Wir haben das Variation-Tree aus der Abbildung ungeparst und es ist zu sehen, dass die Ausgabe im letzten Bild der Abbildung 3.5 gleich dem mit C-Präprozessor-Annotierten Code ist aus der Abbildung 3.3, die zum Parsen Variation-Tree genutzt wurde.

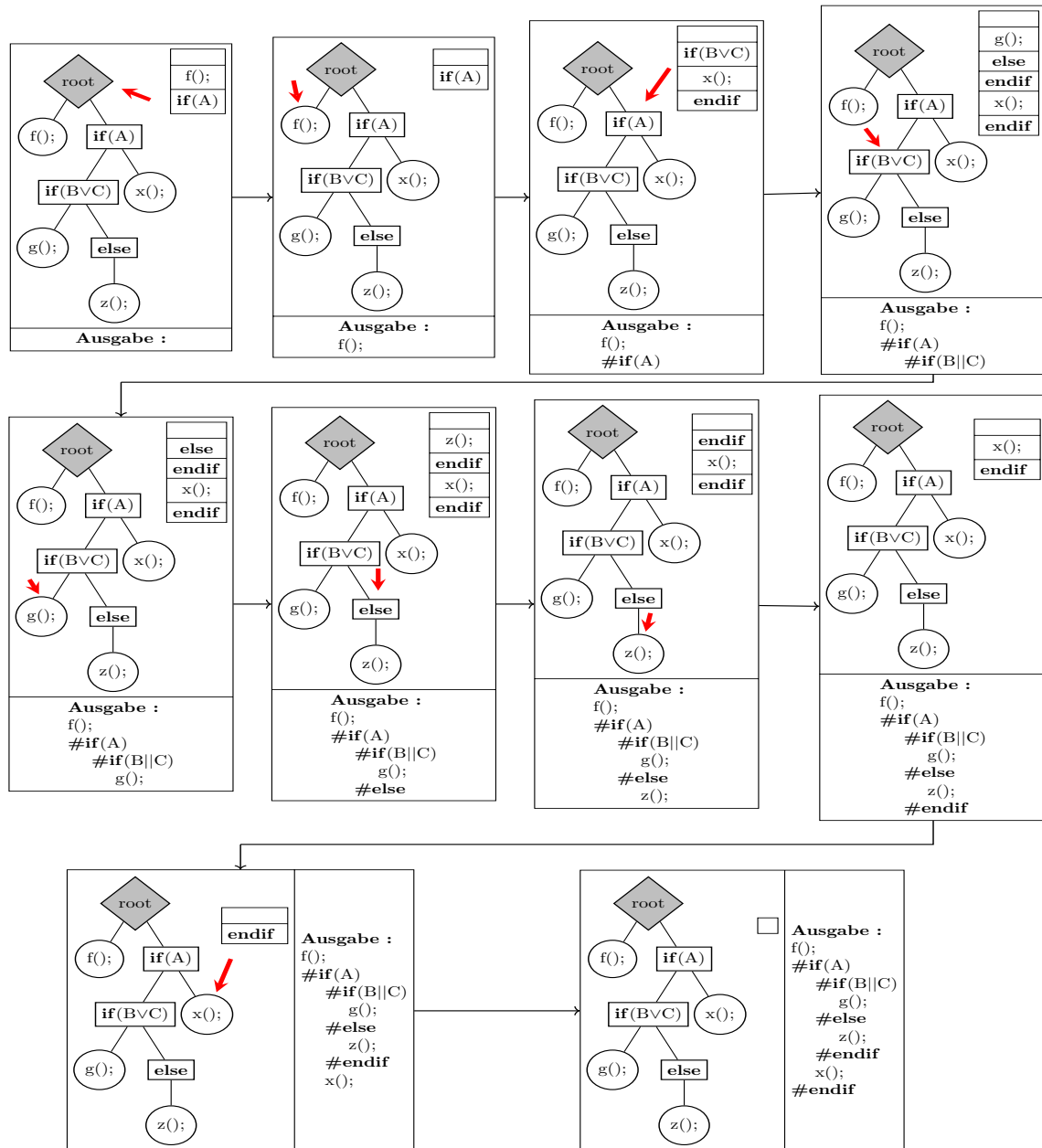


Abbildung 3.5: Beispiel für das Unparsen mithilfe unseren Algorithmus

Wir haben eine Möglichkeit gefunden speichernde, geordnete Variation-Trees zu unparsen. Es bleibt uns noch eine Möglichkeit zum Unparsen von Variation-Diffs zu finden. Eine Möglichkeit zum Unparsen von Variation-Diffs bezieht sich auf speichernde, geordnete Variation-Diffs oder wenn eine Heuristik verwendet wird dann auf geordnete Variation-Diffs. Wir beschreiben das Vorgehen für speichernde, geordnete Variation-Diffs für geordnete Variation-Diffs geht es analog. Zuerst werden zwei Projektionen gebildet, jeweils auf den Zustand-Davor und den Zustand-Danach. Dadurch werden zwei speichernde, geordnete Variation-Trees erstellt. Auf diese Variation-Trees wird unser Algorithmus angewandt. Als Ergebnis erhalten wir zwei mit C-Präprozessor-Annotierte Codes. Zum Schluss wird ein Algorithmus, welcher aus Codes mit C-Präprozessor-Annotation ein textbasiertes Diff erstellen kann, auf das Ergebnis unseren Algorithmus angewandt. Damit haben wir ein Vorgehen zum Unparsen von speichernden, geordneten Variation-Diffs.

Jetzt wissen wir wie man speichernde, geordnete bzw. geordnete Variation -Trees und speichernde, geordnete bzw. geordnete Variation-Trees unparsen kann. Als Nächstes wollen wir eine Komplexitätsanalyse der Laufzeit für unseren Algorithmus durchführen.

3.5 Komplexitätsanalyse der Laufzeit

In diesen Teil des Kapitels wollen wir eine Komplexitätsanalyse der Laufzeit für unseren Algorithmus durchführen. Damit wir eine obere Schranke für seine Effizienz in Abhängigkeit von der Knotenmenge setzen können.

Bei unserer Komplexitätsanalyse der Laufzeit gehen wir von folgenden Laufzeiten für einige Anweisungen aus. Die Zeitkomplexität von Operationen mit dem Stack ist konstant. Das Zugreifen oder Setzen von Inhalten des speichernden, geordneten Variation-Trees wird als Konstanten Zeitaufwand betrachtet, da man τ, l, O, M auch als Variablen eines Knotens sehen kann und der Zugriff oder Setzung von diesen einen konstanten Aufwand hat. Die Konstante Zeitkomplexität hat auch die Verbindung von zwei Strings. Das Erstellen eines Dummyknotens wird auch als konstanten Zeitaufwand betrachtet, da ein Dummyknoten nur als Platzhalter im Stack für die `#endif` Zeile dient und nur $M[0]$ und τ sich merkt. Das Setzen dieser Information ist wie schon angegeben in konstanter Zeit möglich. Wir gehen davon aus, da wir keine komplexen Datenstrukturen haben, dass die Initialisierung in konstanter Zeit möglich ist. In Zeilen 3 und 18 des Algorithmus 4 werden die Kinderknoten bestimmt. Der dazugehörige Zeitaufwand hängt davon ab, wie die Speicherung, der Kinderknoten realisiert ist. Wir gehen davon aus dass die Speicherung von Kinderknoten in Form von Adjazenzliste erfolgt. Diese Annahme ergibt, dass das Erhalten der Kinderknoten eines Knotens eine konstante Zeitkomplexität hat. In Algorithmus 4 ist die Zeitkomplexität, in der O-Notation, für die einzelnen Anweisungen unseren Algorithmus zu sehen.

Algorithmus 4: Komplexitätsanalyse der Laufzeit für Algorithmus 3

Data: ein Variation-Tree (V, E, r, τ, l)
Result: mit C-Präprozessor-Annotierten Code

```

1 initialisiere einen leeren Stack/Keller stack  $O(1)$ 
2 initialisiere einen String ergebnis  $O(1)$ 
3 kinder  $\leftarrow \{v \in V \mid (r, v) \in E\}$   $O(1)$ 
4 initialisiere ein Array array der Länge  $|kinder|$   $O(1)$ 
5 for  $\forall v \in kinder$  do
6   | array[ $O(v)$ ]  $\leftarrow v$   $O(1)$ 
7 end
8 for  $i = |kinder| \rightarrow 1$  do
9   | stack.push(array[i])  $O(1)$ 
10 end
11 while stack nicht leer ist do
12   | knoten  $\leftarrow stack.pop()$   $O(1)$ 
13   | if  $\tau(knoten) = mapping$  then
14     | erstelle einen Dummyknoten welcher  $M(knoten)[1]$  beinhaltet  $O(1)$ 
15     | stack.push(Dummyknoten)  $O(1)$ 
16   | end
17   | erweitere ergebnis mit dem String  $M(knoten)[0]$   $O(1)$ 
18   | kinder  $\leftarrow \{v \in V \mid (knoten, v) \in E\}$   $O(1)$ 
19   | initialisiere ein Array array der Länge  $|kinder|$   $O(1)$ 
20   | for  $\forall v \in kinder$  do
21     | array[ $O(v)$ ]  $\leftarrow v$   $O(1)$ 
22   | end
23   | for  $i = |kinder| \rightarrow 1$  do
24     | stack.push(array[i])  $O(1)$ 
25   | end
26 end
27 return ergebnis  $O(1)$ 

```

Es ist zu sehen, dass alle einzelnen Anweisungen in unserem Algorithmus eine Laufzeit von $O(1)$ haben. Es bleibt uns nur noch zu bestimmen wie viel man die Schleifen durchgelaufen werden. Es ist zu sehen, dass die Schleife in Zeilen 5 bis 7 und die Schleife in Zeilen 8 bis 10 dieselbe Laufzeit haben. n ist die Anzahl aller Knoten. Die Schleife in Zeilen 5 bis 7, wird so viel mal durchgelaufen wie der Wurzelknoten Kinderknoten hat, das bezeichnen wir mit w . Dabei gilt $w < n$ oder noch genauer $w \leq n-1$, da alle Knoten an dem Wurzelknoten hängen können und dabei können es alles Blattknoten sein. Deshalb können wir die Laufzeit der Schleife wie folgt angeben $O(1) * w = O(1 * w) = O(w) = O(n-1) = O(n)$. Dasselbe gilt auch für die Schleife in Zeilen 8 bis 10. Damit ist die Laufzeit der Schleifen zusammen $O(n) + O(n) = O(n + n) = O(2n) = O(n)$. In Zeilen 20 bis 22 haben wir eine Schleife, die über die Kinderknoten des Knotens v geht. Die Schleife wird k_v mal durchlaufen und k_v gibt die Anzahl der Kinderknoten des Knotens v . Unser Algorithmus geht über alle Knoten des Baumes also werden von jedem Knoten die Kinderknoten ermittelt, dessen Anzahl k_v für Knoten v angibt. Bei einem Baum hat ein Knoten nur einen einzigen Elternknoten und wird deshalb auch nur ein mal durchlaufen, weil es auch nur ein einziges Mal den Stack hinzugefügt wird. Es ergibt sich das $w + \sum_{v \in V} k_v = n-1$ gilt. Es gilt, da ein Knoten v k_v Kinderknoten hat, dabei ist w die Anzahl der Kinderknoten des Wurzelknotens. Der Baum ist zusammenhängend also haben alle Knoten außer den Wurzelknoten einen Elternknoten und sind damit auch die Kinderknoten eines anderen Knotens. Deshalb

muss die Aufsummierung alle Kinderknoten die Anzahl alle Knoten außer den Wurzelknoten ergeben, was auch $n-1$ ist. Damit ergibt es sich das die Schleife in Zeilen 20 bis 22 genau $n-1$ -mal insgesamt durchlaufen wird unabhängig davon wie viel mal die äußere Schleife in Zeile 11 durchgelaufen wird. Dasselbe gilt auch für die Schleifen in Zeilen 23 bis 25. Es ergibt sich das die Laufzeit einer Schleife $O(1) * n-1-w = O(1 * (n-1-w)) = O(n-1-w) = O(n)$. Beide Schleifen zusammen haben eine Laufzeit von $O(n) + O(n) = O(n + n) = O(2n) = O(n)$. Diese Laufzeit gilt nicht für einen Durchlauf der äußeren Schleife, sondern für die Arbeitszeit des gesamten Algorithmus. Es bleibt uns nur noch zu schauen wie viel mal die Schleife von Zeile 11 bis Zeile 26 durchlaufen wird. Diese Schleife wird so lange durchlaufen bis der Stack leer wird, dabei wird bei jedem Schleifendurchlauf ein Element aus dem Stack genommen. Wir müssen herausfinden wie viel Elemente insgesamt den Stack hinzugefügt werden. Wir wissen das jeder Kinderknoten des Baumes den Stack hinzugefügt wird, das sind schon $n-1$ Elemente. Es gibt aber noch eine Stelle im Algorithmus, welche Knoten den Stack zufügt und das ist die Zeile 15 des Algorithmus 4. Damit man herausfinden kann wie viel Knoten durch diese Stelle hinzugefügt werden, müssen wir herausfinden wie oft maximal kann die Bedingung in Zeile 13 wahr sein. Das gilt so, da es für jeden Knoten mit τ gleich MAPPING ein Dummyknoten erstellt wird und dem Stack hinzugefügt wird. Wir müssen herausfinden wie viel der $n-1$ Knoten τ gleich MAPPING haben können. Es kann vorkommen das alle $n-1$ Kinderknoten τ gleich MAPPING haben, welches ergibt, dass auch $n-1$ Dummyknoten erstellt und den Stack hinzugefügt werden. In diesem Fall wurden den Stack insgesamt $2(n-1)$ Elemente hinzugefügt und die äußere Schleife wird auch $2(n-1)$ Mal durchlaufen. Die inneren Schleifen werden wie schon oben erwähnt unabhängig von der äußeren Schleife durchlaufen. Deshalb hat auch deren Zeitkomplexität keinen Einfluss auf die Zeitkomplexität der äußeren Schleife. Die if-Anweisung aus Zeilen 13 bis 16 hat eine Laufzeit von $\max(O(1)+O(1),0) + O(1) = \max(O(1),0) + O(1) = O(1) + O(1) = O(1)$. Alle anderen Anweisungen haben auch eine Laufzeit von $O(1)$. Es folgt $(O(1)+O(1)+O(1)+O(1)+O(1))*2(n-1) = O(1)*2(n-1) = O(1*2(n-1)) = O(2(n-1)) = O(n-1) = O(n)$. Eine Laufzeit von $O(n)$ haben alle Schleifen und sind dabei voneinander unabhängig. Eine Laufzeit von $O(1)$ haben die Zeilen 1 bis 4. Diese Zeilen wurden nicht in Schleifen oder anderswo berücksichtigt. Die gesamte Laufzeit unseren Algorithmus ist folgende $O(1) + O(1) + O(1) + O(1) + O(n) + O(n) + O(n) = O(1+1+1+1+n+n+n) = O(3n+4) = O(n)$. Damit konnten wir zeigen, dass die asymptotische Laufzeit unseren Algorithmus bei $O(n)$ liegt.

Die Komplexitätsanalyse der Laufzeit für das Unparsen von speichernden, geordneten Variation-Diffs so wie es gezeigt wurde, muss über drei Algorithmen geschehen. Wir müssen eine Komplexitätsanalyse der Laufzeit für die Projektion durchführen, eine für unseren Algorithmus, was wir auch in oberen Abschnitt gemacht haben, und eine für den Diff-Algorithmus. Die Projektion ist eine Tiefensuche über den Variation-Diff was auch zu einer asymptotischen Laufzeit von $O(n)$ führt. Der Teil, welcher uns Schwierigkeiten bereitet ist, die Komplexitätsanalyse der Laufzeit von dem Diff-Algorithmus. Es können unterschiedliche Algorithmen verwendet werden, welche zu unterschiedlichen Laufzeiten führen. Aus diesem Grund haben wir keine Komplexitätsanalyse der Laufzeit für das Unparsen von speichernden, geordneten Variation-Diffs durchgeführt.

3.6 Zusammenfassung

In diesem Kapitel haben wir uns zuerst mit den Definitionen von Variation-Tree und Variation-Diff beschäftigt. Als Nächstes haben wir uns die Arbeitsweise des Parsers angeschaut. Danach auf der Basis der vorherigen Information haben herausgefunden, welche Informationen während des Parsens verloren gehen und von uns wiederhergestellt werden müssen. Dazu haben

wir die Definitionen von Variation-Tree und Variation-Diff für unsere Zwecke erweitert. Dabei haben wir auch den Parser-Algorithmus so erweitert, dass er von uns definierte speichernde Variation-Trees bzw. geordnete, speichernde Variation-Trees und speichernde Variation-Diffs bzw. geordnete, speichernde Variation-Diffs erstellen kann. Als Nächstes wurde von uns unser Unparse Algorithmus für das Unparsen von Variation-Tree und ein Vorgehen zum Unparsen von Variation-Diffs vorgestellt. Zum Schluss haben wir eine Komplexitätsanalyse der Laufzeit für unsren Algorithmus durchgeführt. Die Komplexitätsanalyse der Laufzeit für unsren Algorithmus hat ergeben, dass der Algorithmus eine asymptotische Laufzeit von $O(n)$ hat.

Implementierung

In diesem Kapitel stellen wir eine mögliche Implementierung unseres Algorithmus vor. Unser Algorithmus wird in das Tool `DiffDetective` eingebaut und erweitert damit seine Möglichkeiten, da `DiffDetective` keinen Unparser für Variation-Trees oder Variation-Diffs hat. `DiffDetective` ist in Java programmiert, deshalb nutzen wir die gleiche Programmiersprache für unseren Unparser. Wir werden unseren Algorithmus in je einer Methode für `VariationTree<T>` und `VariationDiff<T>` implementieren.

Die konkrete Implementierung von Variation-Trees und Variation-Diffs in `DiffDetective` weicht von deren theoretischen Ausarbeitung ab. Für uns wichtig ist, dass der Knotentyp τ nicht eine Funktion ist welche auf Knoten angewandt wird, sondern ein Attribut eines Knotens ist. Dazu werden mehr Fälle als bei dem Knotentyp τ unterscheiden. Jeder Knoten speichert die Kinderknoten in einer Liste und gibt damit auch ihre Reihenfolge an. Variation-Trees und Variation-Diffs in `DiffDetective` sind generische Datenstrukturen. Der generische Teil ist für das Label zuständig. Die für das Label gewählte Klasse muss das Interface `Label` implementieren.

Jedes Label, welches von `VariationTree` und `VariationDiff` werdet wird, muss das Interface `Label` implementieren. Das Wissen darüber können wir uns zunutze machen. Eine Klasse, welche das Interface `Label` implementiert, muss die Zeilen als Liste von Strings darstellen können. In der Umsetzung unseres Algorithmus nutzen wir diese Funktion des Interfaces. Das Label speichert die Zeilen. Deshalb müssen wir das `M`, aus den Definitionen von geordneten, speichernden Variation-Trees und Variation-Diffs, nicht in vollen Umfang realisieren, sondern nur das Speichern der Zeile mit `endif`.

An einigen Stellen muss der Code von `DiffDetective` geändert werden, damit alle für das Unparsen relevante Informationen gespeichert werden. Wir müssen die Klassen `DiffNode` und `VariationTreeNode` um das Attribut `endif` erweitern, welches die Zeile mit dem `endif` speichert. Der Parser muss so geändert werden, dass er die Zeile mit dem `endif` im richtigen Attribut speichert. Die Projektion muss dieses Attribut ebenfalls behandeln, da bei der Projektion die Knoten ihre Klasse wechseln.

In Abbildung 4.1 sehen wir, wie wir unseren Algorithmus umgesetzt haben. Die Funktion `variationTreeUnparser` und `variationDiffUnparser` sind generisch. Der generische Typ `T` muss das Interface `Label` implementieren. Der Input der Funktion `variationTreeUnparser` ist ein `VariationTree` mit Typ `T` und ein Objekt der Klasse `Function` auch mit dem Typ `T` über-

```

1 public static <T extends Label> String variationTreeUnparser(VariationTree<T> tree, Function<
2     List<String>,T> linesToLabel){
3     if(!tree.root().getChildren().isEmpty()) {
4         StringBuilder result = new StringBuilder();
5         Stack<VariationTreeNode<T>> stack = new Stack<>();
6         for (int i = tree.root().getChildren().size() - 1; i >= 0; i--) {
7             stack.push(tree.root().getChildren().get(i));
8         }
9         while (!stack.empty()) {
10             VariationTreeNode<T> node = stack.pop();
11             if (node.isIf()) {
12                 stack.push(new VariationTreeNode<>(NodeType.ARTIFACT, null,null,
13                     linesToLabel.apply(node.getEndIf())));
14             }
15             for (String line : node.getLabel().getLines()) {
16                 result.append(line);
17                 result.append("\n");
18             }
19             for (int i = node.getChildren().size() - 1; i >= 0; i--) {
20                 stack.push(node.getChildren().get(i));
21             }
22             return result.substring(0, result.length() - 1);
23         }else{
24             return "";
25     }
26 }

```

Abbildung 4.1: Unparser für `VariationTree<T>`

geben. Wir haben unsere Funktion generisch gemacht, damit sie für alle Variation-Trees funktioniert. Von dem `VariationTree` werden die Kindknoten des Wurzelknotens bekommen und damit auch alle anderen Knoten bekommen und auch für in einer für uns nutzbaren Reihenfolge. Der zweite Parameter, der Objekt der Klasse Funktion, wird für die Erstellung des Dummyknotens verwendet. Die Funktion `linesToLabel` erhält die Zeilen mit `#endif` als `List<String>` und gibt eine `List<String>` mit den `#endif` Zeilen zurück und hat dabei die generische Klasse `T`. Der Output unserer Funktion `variationTreeUnparser` ist ein String. Dieser String enthält mit C-Präprozessor-Annotierten Code, welcher aus dem gegebenen `VariationTree` erstellt wurde. Am Anfang wird geprüft, ob der Wurzelknoten Kinderknoten hat. Wenn nicht, wird ein leerer String ausgegeben. Sonst wird der Algorithmus ausgeführt. Die Umsetzung ähnelt dem Algorithmus, abgesehen von den Java- und Implementierung-Spezifischen Ausdrucksweise. In der Funktion werden zuerst `StringBuilder` zur Speicherung des Ergebnisses und der Stack initialisiert. Danach werden die Kinderknoten des Wurzelknotens in umgekehrter Reihenfolge auf den Stack gelegt. Als Nächstes beginnt die Schleife, welche so lange durchlaufen wird, bis der Stack leer ist. In der Schleife wird ein Knoten vom Stack genommen. Als Nächstes wird geprüft, ob der Knoten ein If-Knoten ist, in unserem Algorithmus wurde geschaut, ob der τ des Knotens gleich mapping ist. Wenn ja, wird ein Dummyknoten erstellt und dem Stack hinzugefügt. Wenn nein, wird nichts gemacht. Danach so wie in dem Algorithmus 3 werden die Zeilen dieses Knotens dem Ergebnis hinzugefügt. Am Ende des Schleifendurchlaufs werden die Kinderknoten des Knotens in umgekehrter Reihenfolge dem Stack hinzugefügt. Nach der Schliefe wird das Ergebnis als String zurückgegeben.

In der Abbildung 4.2 ist zu sehen, wie wir das Unparsen von Variation-Diffs umgesetzt


```

1 public static <T extends Label> String variationDiffUnparser(VariationDiff<T> diff, Function<List
  <String>,T> linesToLabel) throws IOException {
2     String tree1 = variationTreeUnparser(diff.project(Time.BEFORE), linesToLabel);
3     String tree2 = variationTreeUnparser(diff.project(Time.AFTER), linesToLabel);
4     return JGitDiff.textDiff(tree1, tree2, SupportedAlgorithm.MYERS);
5 }

```

Abbildung 4.2: Unparser für `VariationDiff<T>`

haben. Dabei sind wir wie in dem Kapitel 3 gesagt worden vorgegangen. Wir haben das Problem von Unparsen eines Variation-Diffs auf das Unparsen von Variatio-Trees reduziert. So machen wir die Implementierung des Unparsers für Variation-Diff: Wir projizieren das übergebene `VariationDiff` auf den Zustand-Davor und unparsen dann das erhaltene `VariationTree` mit der vorher gezeigten Funktion. Das Gleiche machen wir auch für den Zustand-Danach. Zuletzt bilden wir, mit einer von `DiffDetectiv` zur Verfügung gestellten Funktion, einen textbasierten Diff und geben ihn als String zurück.

Wir haben eine mögliche Implementierung für unseren Algorithmus in `DiffDetective` vorgestellt. So haben wir `DiffDetective` verbessert und den Benutzern mehr Funktionen zur Verfügung gestellt. Wir müssen noch festlegen, wie wir die Korrektheit unseres Algorithmus überprüfen und eine Auswertung für unsere Funktionen machen.

Korrektheit

In diesem Kapitel stellen wir die Korrektheitskriterien vor, mit deren Hilfe die Korrektheit eines Unparsers bewertet werden kann. In einer empirischen Studie testen wir dann, ob und welche Korrektheitskriterien von unserer Implementierung erfüllt werden.

Es werden jeweils Korrektheitskriterien für textbasierte Diffs und C-Präprozessor-Annotierten Code definiert. Diese sind syntaktische Gleichheit, syntaktische Gleichheit ohne Whitespace und semantische Gleichheit. Die semantische Gleichheit wurde nur für textbasierte Diffs definiert und nicht für C-Präprozessor-Annotierten Code. Es gibt Zusammenhänge zwischen den Korrektheitskriterien: Syntaktische Gleichheit ohne Whitespace wird von der syntaktischen Gleichheit impliziert und aus syntaktischer Gleichheit ohne Whitespace folgt semantische Gleichheit. Bei der Auswertung werden Aufbau und Ergebnisse beschrieben. Basierend darauf werden die Forschungsfragen und Ergebnisse diskutiert. Dabei wird während der Auswertung unsere Implementierung auf die Korrektheitskriterien anhand der Commit-Historien von drei Open-Source Repositories überprüft.

In Kapitel 5.1 werden die Korrektheitskriterien definiert, näher beschrieben und die Zusammenhänge zwischen diesen dargestellt. Die Auswertung, einschließlich des Aufbau des Experiments, der Ergebnisse und der Diskussion, ist in Kapitel 5.2 zu finden.

5.1 Korrektheitskriterium

Nachdem wir unseren Algorithmus in DiffDetective implementiert haben, müssen wir noch feststellen, wie korrekt unsere Implementierung ist. Dazu müssen wir zuerst definieren, wie wir die Korrektheit unserer Implementierung verstehen. Deshalb haben wir Korrektheitskriterien für textbasierte Diffs und C-Präprozessor-Annotierten Code aufgestellt, um die es in diesem Abschnitt geht.

Die Bewertung der Korrektheit unseres Unparser erfolgt wie folgt: Es wird ein C-Präprozessor-Annotierter Code oder textbasierter Diff genommen. Dieser C-Präprozessor-Annotierter Code oder textbasierter Diff wird geparkt und wir erhalten einen Variation-Tree oder ein Variation-Diff. Der Variation-Tree oder Variation-Diff wird im nächsten Schritt ungeparkt und wir erhalten wieder einen C-Präprozessor-Annotierten Code oder einen textbasierten Diff. Zum Schluss wird der Ausgangs-Code mit dem ungeparkten Code beziehungsweise das Ausgangs-Diff mit dem

ungeparsten Diff anhand der von uns definierten Korrektheitskriterien verglichen, um zu entscheiden, ob korrekt ungeparst wurde oder nicht.

Die von uns definierten Korrektheitskriterien sind syntaktische Gleichheit, syntaktische Gleichheit ohne Whitespace und semantische Gleichheit. Es gibt mehrere Korrektheitskriterien, da es zwei Anwendungen gibt, die zusammenarbeiten müssen. Der Parser kann die benötigten Informationen unterschiedlich gut speichern und der Unparser kann an gewissen Stellen mit Heuristiken arbeiten. Dies kann zusammen die Korrektheit des Ergebnisses beeinträchtigen. Die syntaktische Gleichheit haben wir gewählt, da dies der Idealzustand ist, der erreicht werden kann. Die syntaktische Gleichheit ohne Whitespace haben wir gewählt, da man nicht für jeden Anwendungszweck unbedingt einen perfekten Unparser braucht. Dazu können während des Parsens Informationen über Whitespace-Zeichen verloren gehen, welche keinen Einfluss auf das Verständnis des Inhalts haben, aber trotzdem eine perfekte Rekonstruktion verhindern. Die semantische Gleichheit gilt nur für textbasierte Diffs und wurde von uns ausgewählt, da es zurzeit keinen Algorithmus gibt, der aus einem Variation-Diff einen textbasierten Diff erstellt. Es gibt nur das von uns entwickelte Vorgehen, aber bei diesem Vorgehen wird aus zwei C-Präprozessor-Annotierten Codes ein textbasierter Diff erstellt. Der so erhaltene textbasierte Diff kann abhängig davon, welcher Algorithmus zu diesem Zweck verwendet wurde, unterschiedlich aussehen, obwohl die gleichen C-Präprozessor-Annotierten Codes verwendet wurden. Damit wir auch so was überprüfen können, haben wir die semantische Gleichheit eingeführt. In der Tabelle 5.1 sind die formalen Definitionen dieser Korrektheitskriterien zu sehen. *id*, welche in der Tabelle 5.1 verwendet wird, stellt die Identitätsfunktion dar. Die Signaturen anderer in der Tabelle 5.1 verwendeter Funktionen sind folgende: Sei C die Menge aller C-Präprozessor-Annotierten Codes, D die Menge aller textbasierten Diffs, VT die Menge aller Variation-Trees, VD die Menge aller Variation-Diffs und $Z = \{a, b\}$ die Menge mit den Zeiten Davor und Danach. Dann sind die Funktionen wie folgt definiert: $parse_t : C \rightarrow VT$, $unparse_t : VT \rightarrow C$, $deleteWhitespace_t : C \rightarrow C$, $parse_d : D \rightarrow VD$, $unparse_d : VD \rightarrow D$, $textProject : D \times Z \rightarrow C$. Die Formeln in der Tabelle 5.1 zeigen funktional die Definitionen der Korrektheitskriterien. Eine solche Darstellung eignet sich gut für die Definition, da wir in unserem Fall den Parser und Unparser als Funktionen darstellen können und deren Zusammenarbeit als Verkettung. Dazu lassen sich alle anderen Manipulationen als Funktionen darstellen, wie zum Beispiel das Entfernen von Whitespace und die Projektion von textbasierten Diffs. Wir können auch den Zugriff auf einzelne Elemente darstellen, so wie bei der semantischen Gleichheit.

	C-Präprozessor- Annotierter Code	textbasierter Diff
Syntaktische Gleichheit	$unparse_t \circ parse_t = id$	$unparse_d \circ parse_d = id$
Syntaktische Gleichheit ohne Whitespace	$deleteWhitespace_t \circ unparse_t \circ parse_t = deleteWhitespace_t \circ id$	$deleteWhitespace_d \circ unparse_d \circ parse_d = deleteWhitespace_d \circ id$
Semantische Gleichheit	Out of Scope unentscheidbar für C, exponentielles Wachst- tum für C-Präprozessor	Für $\forall t \in \{a, b\}$ und $\forall d \in D$ $x := unparse_d \circ parse_d(d)$ $textProject(x, t) = textProject(id(d), t)$ $\vee deleteWhitespace_t \circ textProject(x, t) = deleteWhitespace_t \circ textProject(id(d), t)$

Tabelle 5.1: Korrektheitskriterien

In diesem Abschnitt sprechen wir über die syntaktische Gleichheit, die in der zweiten Zeile aus der Tabelle 5.1 zu finden ist. Syntaktische Gleichheit bedeutet, dass der zu vergleichende Text in jedem Zeichen identisch ist. Der Vergleich auf syntaktische Gleichheit sieht für C-Präprozessor-Annotierten Code und textbasierte Diffs gleich aus, was in der Abbildung 5.1 zu sehen ist. Hierfür muss der Ausgangs-Code bzw. der textbasierte Diff mit dem Ergebnis nach dem Parsen und Unparsen in jedem Zeichen übereinstimmen. Wie in der Abbildung 5.1 gezeigt, wird ein C-Präprozessor-Annotierter Code bzw. der textbasierte Diff genommen, dann darauf Parser und Unparser angewendet. Das Ergebnis und der C-Präprozessor-Annotierter Code bzw. der textbasierte Diff werden dann jeweils in einen String umgewandelt und diese dann auf Gleichheit geprüft. So wird die syntaktische Gleichheit des C-Präprozessor-Annotierten Code bzw. den textbasierten Diff und des Ergebnisses von Parser und Unparser überprüft.

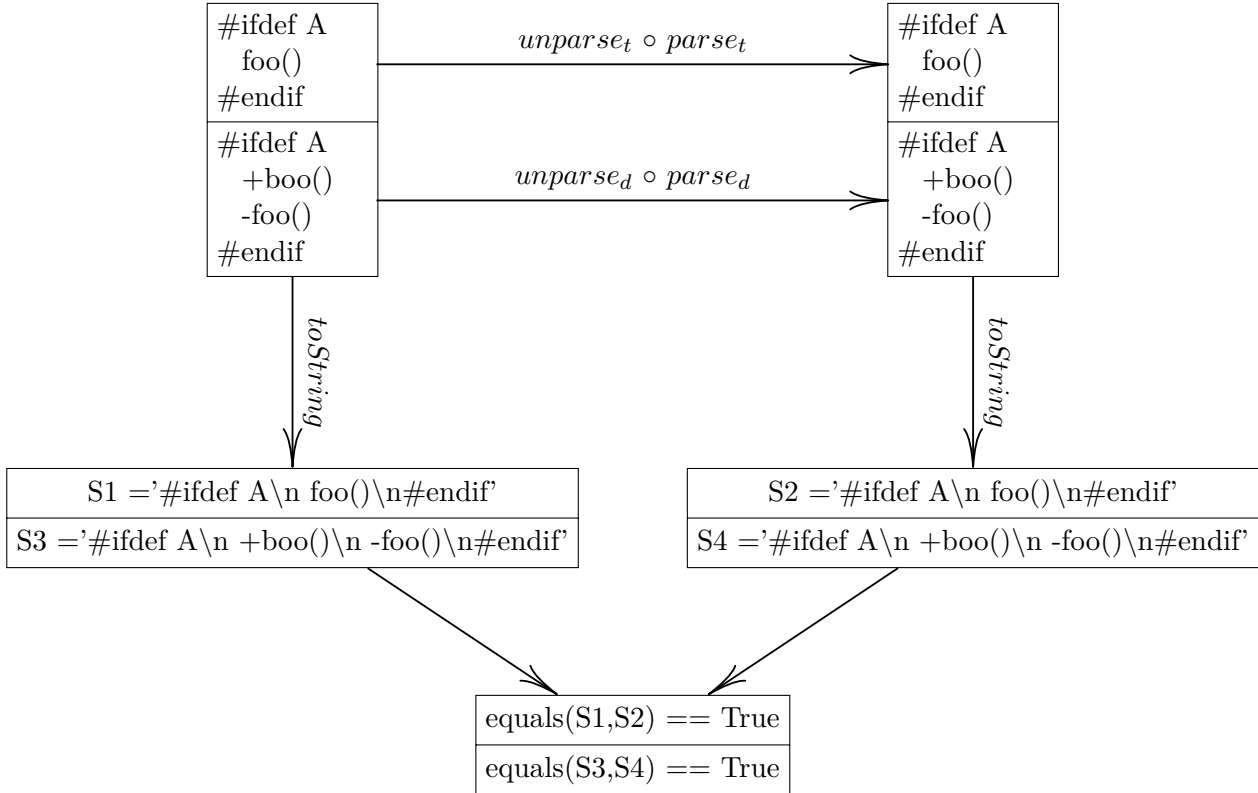


Abbildung 5.1: Beispiel für Syntaktische Gleichheit

Der syntaktischen Korrektheit ohne Whitespace aus der dritten Zeile der Tabelle 5.1 widmen wir uns in diesem Abschnitt. Analog zur syntaktischen Gleichheit ist syntaktische Gleichheit ohne Whitespace für den C-Präprozessor-Annotierten Code und textbasierte Diffs gleich zu verstehen, wie in Abbildung 5.2 zu sehen ist. Bei dieser Art von Korrektheit muss, wie im vorherigen Fall, der Ausgangs-Code mit C-Präprozessor-Annotationen bzw. der textbasierte Diff mit dem Ergebnis nach dem Parsen und Unparsen in jedem Zeichen übereinstimmen, bis auf bestimmten Whitespace. Damit die zu vergleichenden Texte gleich werden, werden einige Schritte unternommen. Diese Schritte sind: das Entfernen aller Whitespace-Zeichen von Anfang einer Zeile bis zu dem ersten Nicht-Whitespace-Zeichen. Danach das Entfernen aller Whitespace-Zeichen nach dem letzten Nicht-Whitespace-Zeichen einer Zeile bis zu dem Zeilenumbruch, der Zeilenumbruch wird nicht entfernt, und das Entfernen aller leeren Zeilen. Im Falle, dass diese Manipulationen mit einem textbasierten Diff durchgeführt werden, wird das Leerzeichen vor dem ersten Nicht-Whitespace-Zeichen nicht entfernt, wenn dieses erste Nicht-Whitespace-Zeichen kein `+` oder `-` ist. Diese Ausnahme ist notwendig, damit die Struktur eines textbasierten Diffs beibehalten wird. Für ein besseres Verständnis ist unsere Implementierung in Java in der Abbildung 5.3 zu sehen. Es werden nicht, wie die Bezeichnung vermuten lässt, alle Whitespace-Zeichen entfernt. Der Grund dafür ist folgender: Wenn alle Whitespace-Zeichen entfernt werden, ist es nicht möglich, festzustellen, ob es Probleme mit den Whitespace-Zeichen gibt, welche für das Verständnis relevant sind. Also zum Beispiel zwischen `#if` und der Bedingung kein Leerzeichen gibt oder die Zeilenumbrüche fehlen, wodurch der Code oder der textbasierte Diff zu einer einzelnen Zeile wird. Solche Fälle werden dann auch als korrekt angesehen, wenn wir alle Whitespace-Zeichen entfernen würden, was wir aber nicht wollen. Die Whitespace-Zeichen, welche innerhalb des Abschnitts vom ersten Nicht-Whitespace-Zeichen bis zu dem letzten Nicht-Whitespace-Zeichen, werden nicht entfernt oder manipuliert. Dazu wird auch der Zeilenumbruch am Ende der nicht leeren Zeilen beibehalten. Die Abbildung 5.2 veranschaulicht dies. Dort sind der Ausgangs-Code

bzw. der textbasierte Diff gegeben. Rechts von dem ist das Ergebnis von Parse und Unparse gegeben. Danach werden diese alle in Strings umgewandelt. Als Nächstes werden die Whitespace-Zeichen, wie oben beschrieben, aus den Strings entfernt und anschließend die auf Äquivalenz geprüft. So wird der C-Präprozessor-Annotierter Code bzw. der textbasierte Diff und das Ergebnis von Parse und Unparse auf syntaktische Gleichheit ohne Whitespace überprüft.

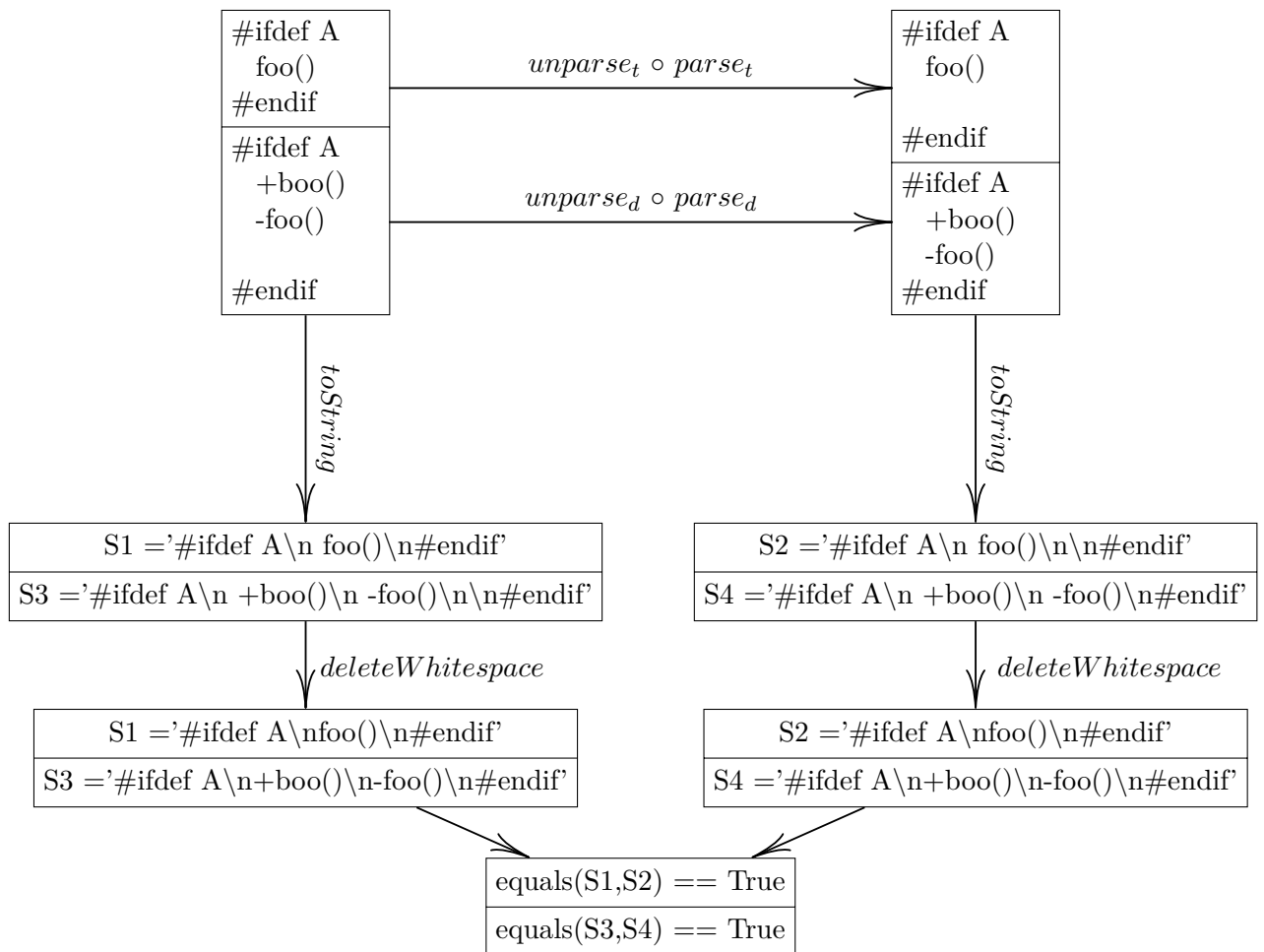


Abbildung 5.2: Beispiel für Syntaktische Gleichheit ohne Whitespace

```

1|    public static String deleteWhitespace(String string, boolean diff){
2|        if(string.isEmpty()){
3|            return "";
4|        }
5|        else {
6|            StringBuilder result = new StringBuilder();
7|            try {
8|                BufferedReader in = new BufferedReader(new StringReader(
9|                    string));
10|                String line = "";
11|                while ((line = in.readLine()) != null) {
12|                    if (!line.replaceAll("\\s+", "").isEmpty()) {
13|                        String temp = line.trim();
14|                        if(diff && !(temp.charAt(0) == '+' || temp.
15|                            charAt(0) == '-')){
16|                            temp = " " + temp;
17|                        }
18|                        result.append(temp);
19|                        result.append("\n");
20|                    }
21|                }
22|            }catch (Exception e){
23|                e.printStackTrace();
24|            }
25|            return result.toString();
26|        }
27|    }

```

Abbildung 5.3: Implementierung von deleteWhitespace

Die semantische Gleichheit von C-Präprozessor-Annotierten Code werden wir nicht betrachten, da wir dafür entscheiden müssen, ob zwei Programme äquivalent sind. Das geht über den Rand unserer Möglichkeiten, da diese Fragestellung unentscheidbar ist und als das Äquivalenzproblem bekannt [8] ist. Mit den C-Präprozessor-Annotationen geht es auch über den Rand unserer Möglichkeiten, da C-Präprozessor-Annotationen hier für Erzeugung der Variabilität verwendet werden. Dabei hat so eine Softwareproduktlinie n Features und im Worst-Case müssen 2^n Varianten der Software betrachtet werden[2], was eine exponentielle Laufzeit bedeutet und somit über den Rand unserer Möglichkeiten geht.

Um die semantische Gleichheit für textbasierte Diffs geht es in diesem Abschnitt. Wie die semantische Gleichheit für textbasierte Diffs zu verstehen ist, ist nicht eindeutig festgelegt. Unsere Interpretation der semantischen Gleichheit für textbasierte Diffs ist an der Gleichheit für Variation-Diffs [4] orientiert. Wir verstehen die semantische Gleichheit wie folgt: Zwei textbasierte Diffs sind semantisch gleich, wenn ihre Projektionen auf den Davor-Zustand bzw. Danach-Zustand syntaktisch gleich oder syntaktisch gleich ohne Whitespace sind. In der Abbildung 5.4 ist dies dargestellt. Dabei ist die Projektion für textbasierte Diffs wie folgt zu verstehen: Ein textbasierter Diff hat Zeilen von drei Typen: unverändert gebliebene Zeilen, gelöschte Zeilen und eingefügte Zeilen. Bei der Projektion werden einige dieser Typen der Zeilen entfernt, einige beibehalten, und so entsteht eine Projektion von textbasierten Diffs auf einen mit C-Präprozessor-Annotierten Code. Dabei wird für die Projektion auf den Zustand-Davor die unveränderten und gelöschten Zeilen beibehalten, während die eingefügten entfernt werden. Für die Projektion auf den Zustand-Danach, die unveränderten und eingefügten Zeilen beibehalten, während die gelöschten Zeilen entfernt werden. Dies verläuft analog zu der Projektion von Variation-Diff zu Variation-Tree.

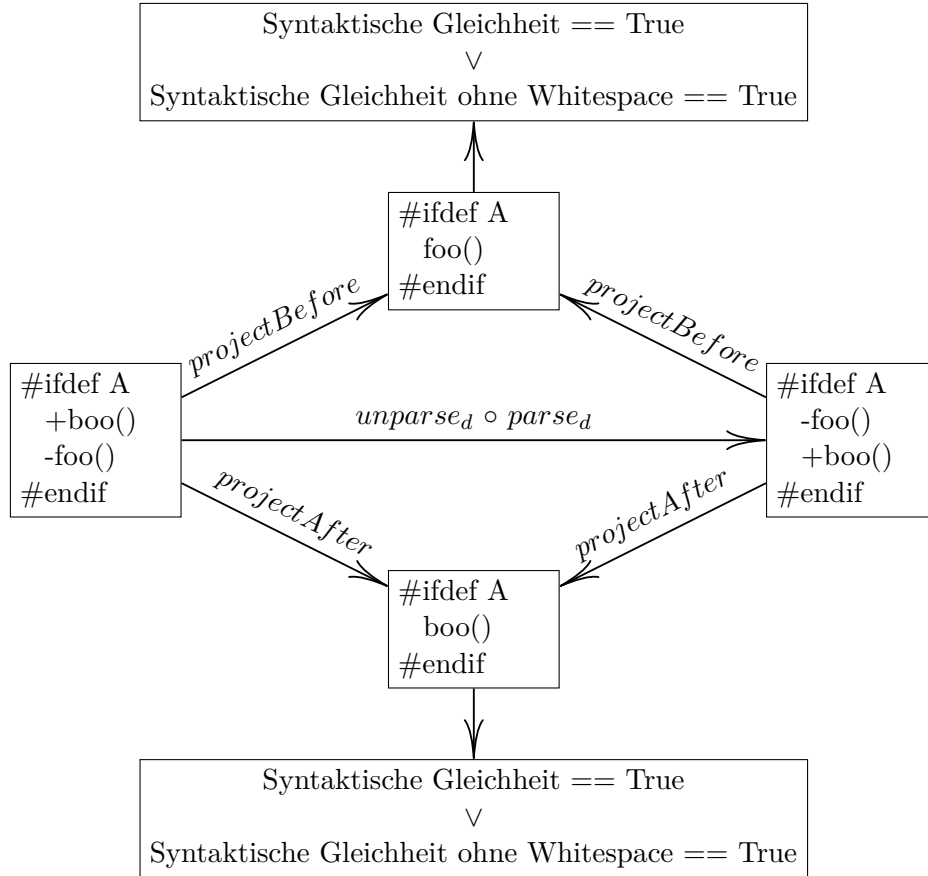


Abbildung 5.4: Beispiel für Semantische Gleichheit

Im Weiteren wollen wir veranschaulichen, wie die Korrektheitskriterien zusammenhängen. Bevor wir aber dazu kommen, müssen wir Hilfsaussagen erläutern. Bei syntaktischer Gleichheit und syntaktischer Gleichheit ohne Whitespace für C-Präprozessor annotiertem Code und textbasierten Diffs werden im Grunde genommen nur Texte verglichen. Aus diesem Grund wird für die Beschreibung der Zusammenhänge zwischen syntaktischer Gleichheit und syntaktischen Gleichheit ohne Whitespace nur auf Text eingegangen und nicht zwischen C-Präprozessor-Annotierten Code und textbasierten Diffs unterschieden.

Der Zusammenhang zwischen syntaktischer Gleichheit und syntaktischer Gleichheit ohne Whitespace ist folgender: Die syntaktische Gleichheit impliziert syntaktische Gleichheit ohne Whitespace. Formal mithilfe von Definition aus der Tabelle 5.1 ist es in dem Beweis unten gezeigt. Diese Aussage stimmt, da syntaktische Gleichheit bedeutet, dass die zwei zu vergleichende Texte in jedem Zeichen und der Position dieser Zeichen identisch sind. Wenn aus den syntaktisch gleichen Texten die Whitespace-Zeichen, wie wir beschrieben haben, entfernt werden, sind diese Texte trotzdem syntaktisch gleich. Der Grund dafür ist, dass bei zwei syntaktisch gleichen Texten alle Whitespace Zeichen an denselben Stellen sind. Das führt mit sich, dass nach deren Entfernung alle anderen Zeichen, welche in beiden Texten gleich sind, auch in beiden Texten identisch verschoben werden und damit auch syntaktische Gleichheit ohne Whitespace besitzen. Damit gilt die Aussage.

Beweis. Da \circ assoziativ ist gilt:

$$deleteWhitespace_t \circ unparsed_t \circ parsed_t = deleteWhitespace_t \circ (unparsed_t \circ parsed_t)$$

$$deleteWhitespace_d \circ unparsed_d \circ parsed_d = deleteWhitespace_d \circ (unparsed_d \circ parsed_d)$$

Wenn syntaktische Gleichheit gilt, ist gegeben:

$$\text{unparse}_t \circ \text{parse}_t = \text{id}$$

$$\text{unparse}_d \circ \text{parse}_d = \text{id}$$

Zu zeigen syntaktische Gleichheit ohne Whitespace ist gegeben:

$$\text{deleteWhitespace}_t \circ \text{unparse}_t \circ \text{parse}_t$$

$$= \text{deleteWhitespace}_t \circ (\text{unparse}_t \circ \text{parse}_t)$$

$$= \text{deleteWhitespace}_t \circ (\text{id})$$

$$= \text{deleteWhitespace}_t \circ \text{id}$$

$$\text{deleteWhitespace}_d \circ \text{unparse}_d \circ \text{parse}_d$$

$$= \text{deleteWhitespace}_d \circ (\text{unparse}_d \circ \text{parse}_d)$$

$$= \text{deleteWhitespace}_d \circ (\text{id})$$

$$= \text{deleteWhitespace}_d \circ \text{id}$$

□

Der Zusammenhang zwischen syntaktischer Gleichheit ohne Whitespace und semantischer Gleichheit ist folgender und gilt in unserem Fall nur für textbasierte Diffs. Der Grund dafür ist, dass wir die semantische Gleichheit für C-Präprozessor-Annotierten Code nicht definiert haben und deshalb auch keine Aussagen darüber treffen können. Die syntaktische Gleichheit ohne Whitespace impliziert semantische Gleichheit. Diese Aussage stimmt, da syntaktische Gleichheit ohne Whitespace bedeutet, dass zwei zu vergleichende textbasierte Diffs, nachdem die Whitespace-Zeichen, wie wir es vorschreiben, entfernt wurden, in jedem Zeichen identisch sind. Das Entfernen der Whitespace-Zeichen, wie wir es vorschreiben, ändert nicht, dass ein textbasierter Diff ein textbasierter Diff ist. Das bedeutet, dass wir dieses textbasierte Diff, bei dem die Whitespace-Zeichen entfernt wurden, projizieren können. Die textbasierten Diffs sind nach dem Entfernen der Whitespace-Zeichen in jedem Zeichen identisch und da diese projiziert werden können, sind auch die Projektionen identisch, was dazu führt, dass die semantische Gleichheit erfüllt wird und damit die Aussage gezeigt ist.

5.2 Auswertung

Im folgenden Abschnitt wollen wir eine Auswertung unserer Implementierung anhand der vorher eingeführten Korrektheitskriterien durchführen. Bei dieser Auswertung wollen wir die drei folgenden Forschungsfragen beantworten.

RQ1 : Erfüllt unser Unparser mindestens ein Korrektheitskriterium für jeden untersuchten Fall?

RQ2 : Wie oft wird welches Korrektheitskriterium von unserem Parser für Variation-Trees eingehalten?

RQ3 : Wie oft wird welches Korrektheitskriterium von unserem Parser für Variation-Diffs eingehalten?

Durch Beantwortung dieser Fragen, wollen wir zu dem Entschluss kommen, ob unser Unparser korrekt ist.

Zuerst im Abschnitt 5.2.1 wird der Aufbau des Experiments vorgestellt. Danach sind die Ergebnisse der Auswertung in Abschnitt 5.2.2 zu sehen. Zum Schluss im Abschnitt 5.2.3 werden die Ergebnisse der Auswertung interpretiert und diskutiert.

5.2.1 Aufbau des Experiments

Bei der Auswertung wird DiffDetective verwendet. Der Aufbau unserer Auswertung ist in der Abbildung 5.5 skizziert und ist folgend aufgebaut: Es werden automatisch die Patches aus dem Git Repositorie extrahiert. Wir arbeiten nur mit Patches, welche geparkt werden können. Die Patches bei denen es nicht gilt, ignorieren wir, da diese für unsere Fragestellung unbedeutend sind. Wir wollen die Korrektheit des Unparsers prüfen, also die Fälle wo wir keinen Variation-Tree oder Variation-Diff bekommen können, müssen wir nicht betrachten. Von Repositories bekommen wir Patches, welche textbasierte Diffs sind, wir haben aber noch keine C-Präprozessor-Annotierte Codes. Damit wir diese bekommen können, müssen wir den textbasierten Diff, welchen wir haben auf den Davor-Zustand und Danach-Zustand projizieren. Damit erhalten wir dann zwei C-Präprozessor-Annotierte Codes. Dadurch wird bei der Auswertung doppelt so viele C-Präprozessor-Annotierte Codes untersucht als textbasierte Diffs. Als Nächstes wird sowohl der textbasierte Diff als auch C-Präprozessor-Annotierte Code geparkt. Der Parser in DiffDetective hat zwei unabhängige Optionen, das sind Collapse-Multiple-Code-Lines und Ignore-Empty-Lines. Bei Collapse-Multiple-Code-Lines werden mehrere hintereinander laufende Zeilen von Code, die keine Annotationen enthalten, in einem Knoten zusammengefasst. Bei Ignore-Empty-Lines werden die leeren Zeilen nicht vom Parser in das Ergebnis aufgenommen und gehen somit verloren. Es gibt vier Möglichkeiten, wie die Optionen für den Parser gesetzt werden können. Für alle diese Möglichkeiten wollen wir betrachten, wie die Korrektheitskriterien eingehalten werden. Deshalb wird das Parsen mit jeder Option-Kombination durchgeführt. Dadurch erhalten wir für ein textbasierten Diff vier Variation-Diffs. Aus einem Patch bekommen wir aber zwei C-Präprozessor-Annotierte Codes, diese werden jeweils unter Verwendung aller Parser-Option-Kombination geparkt, sodass wir schließlich acht Variation-Trees haben. Danach werden diese vier Variation-Diffs und acht Variation-Trees ungeparkt und wir erhalten danach vier textbasierte Diffs und auch C-Präprozessor-Annotierte Codes. Diese textbasierten Diffs und C-Präprozessor-Annotierte Codes werden dann auf die Gleichheit mit den textbasierten Diff und C-Präprozessor-Annotierte Codes, welche für das Parsen verwendet wurden, anhand der von uns definierten Korrektheitskriterien verglichen. Die Ergebnisse nach der Untersuchung aller für uns geltenden Patches werden zusammengetragen.

Als Datenquelle werden von uns die Git-Historien folgender drei Software-Produkt-Linie-Projekte benutzt. Die verwendeten Open-Source-Projekte sind **Vim**, **sylpheed** und **berkeley-db-libdb**. **Vim** ist ein konfigurierbarer Texteditor, welcher auch in den meisten Unix-Systemen und in Apple OS X enthalten ist. **sylpheed** ist ein leichtgewichtiger E-Mail-Client, welcher auf vielen Systemen wie Windows, Linux, BSD, Mac OS X und anderen Unix-ähnlichen Systemen läuft. **berkeley-db-libdb** ist eine eingebettete Key-Value-Datenbankbibliothek. Aus diesen Projekten haben wir nur Dateien den Endungen `.c` und `.cpp` untersucht. Insgesamt haben wir 49401 Patches aus den erwähnten Open-Source-Projekte untersucht.

5.2.2 Ergebnisse

In diesem Abschnitt wollen wir über die Ergebnisse der Auswertung berichten. Während der Auswertung wurden Exceptions ausgelöst. Diese Exceptions wurden dadurch ausgelöst, dass das, was dem Parser übergeben wurde, nicht geparkt werden konnte. Wenn etwas nicht geparkt werden kann, ist es auch nicht für unsere Auswertung relevant und wird von uns ignoriert. Wir können nicht garantieren, dass die Exception nur dadurch ausgelöst wurden, aber keiner der Exceptions hat die Auswertung abgebrochen.

In der Tabelle 5.2 sind die Ergebnisse der Auswertung für C-Präprozessor-Annotierten Code

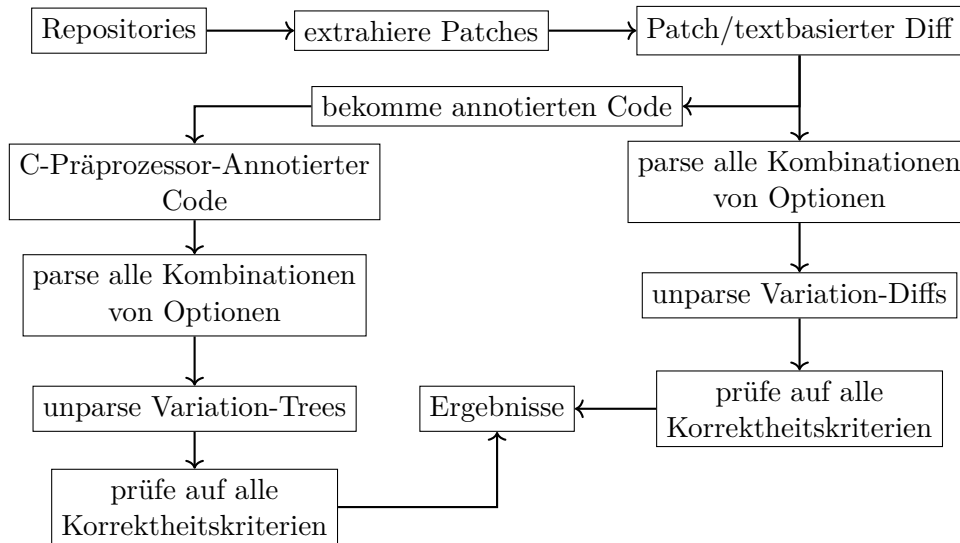


Abbildung 5.5: Aufbau der Auswertung

zu sehen. Die Tabelle zeigt ganz oben die Anzahl der untersuchten C-Präprozessor-Annotierten Codes, das sind 98802. Danach sind auf der linken Seite der Tabelle die Optionen für den Parser zu sehen und ob diese ausgewählt wurden oder nicht. Auf der rechten Seite der Tabelle ist zu sehen, wie oft ein bestimmtes Korrektheitskriterium für die Kombination von Optionen erfüllt wurden. Die Spalte mit dem Fehler zeigt an, wie oft für eine Kombination von Parser-Optionen keiner der von uns definierten Korrektheitskriterien erfüllt wurde.

Anzahl ungeparster C-Präprozessor-Annotierter Codes				98802
Parser-Optionen		Ergebnisse der Auswertung für		
Collapse-Multiple-Code-Lines	Ignore-Empty-Lines	syntaktische Gleichheit	syntaktische Gleichheit ohne Whitespace	Fehler
False	False	0	98802	0
True	False	0	98802	0
False	True	0	98624	178
True	True	0	98624	178

Tabelle 5.2: Ergebnisse der Auswertung für das Unparsen von C-Präprozessor-Annotierten Code

Die Ergebnisse der Auswertung für textbasierte Diffs sind in der Tabelle 5.3 zu sehen. Die Tabelle ist analog zu der Tabelle 5.2 aufgebaut. Oben in der Tabelle ist die Anzahl der untersuchten, textbasierten Diffs zu sehen. Die Anzahl beträgt 49401. Die Anzahl der untersuchten C-Präprozessor-Annotierten Codes aus der Tabelle 5.2 ist doppelt so groß wie die Anzahl der untersuchten textbasierten Diffs. Das liegt daran, dass in unserem Fall jeder C-Präprozessor-Annotierter Code aus textbasierten Diff bekommen wurde. Ein textbasierter Diff besteht aus zwei C-Präprozessor-Annotierten Codes. Wie in der Tabelle 5.2 sind auch in der Tabelle 5.3 die Parser-Optionen und die ausgewählten Kombinationen angegeben. Rechts in der Tabelle sind die Korrektheitskriterien mit dem Fehlerfall zu sehen. In der Tabelle 5.3 ist ein zu untersuchendes Korrektheitskriterium mehr gegeben als in der Tabelle 5.2.

Anzahl ungeparster textbasierter Diffs					49401
Parser-Optionen		Ergebnisse der Auswertung für			
Collapse-Multiple-Code-Lines	Ignore-Empty-Lines	syntaktische Gleichheit	syntaktische Gleichheit ohne Whitespace	semantische Gleichheit	Fehler
False	False	0	46514	49390	11
True	False	0	46514	49390	11
False	True	0	36778	49300	101
True	True	0	36778	49300	101

Tabelle 5.3: Ergebnisse der Auswertung für das Unparsen von textbasierten Diffs

Es wurden 11 textbasierte Diffs als nicht unparsbar vermerkt. So ein Fall tritt auf, wenn für einen textbasierten Diff keiner der Korrektheitskriterien für keine Kombination der Parser-Optionen erfüllt wurde. Wir haben alle diese nicht unparsbaren textbasierten Diffs näher betrachtet. Es hat sich herausgestellt, dass in all diesen Fällen es Probleme mit `#endif`-Annotation gab. Dabei war in jedem dieser Fehlerfälle zu jeder `#if`-Annotation eine `#endif`-Annotation vorhanden.

5.2.3 Diskussion

Nachdem wir die Ergebnisse der Auswertung vorgestellt haben, sind wir in der Lage, die Forschungsfragen zu beantworten. Unsere Forschungsfragen sind folgende: RQ1: Erfüllt unser Unparser mindestens ein Korrektheitskriterium für jeden untersuchten Fall? RQ2: Wie oft wird welches Korrektheitskriterium von unserem Parser für Variation-Trees eingehalten? RQ3: Wie oft wird welches Korrektheitskriterium von unserem Parser für Variation-Diffs eingehalten?

Damit wir diese Fragen RQ1 und RQ2 für das Unparsen von C-Präprozessor-Annotierten Code beantworten können, müssen wir die Tabelle 5.2 betrachten. Für die Beantwortung der Frage RQ1 sind jedoch nur die Zeilen mit den Kombinationen von Parser-Optionen relevant, bei denen Collapse-Multiple-Code-Lines auf False gesetzt ist und Ignore-Empty-Lines auf False sowie Collapse-Multiple-Code-Lines auf True gesetzt ist und Ignore-Empty-Lines auf False. Es sind nur diese beiden Zeilen von Bedeutung, da unsere Unparser darauf ausgerichtet ist, nur mit der Parser-Option Collapse-Multiple-Code-Lines korrekt zu funktionieren. In diesem Fall gehen keine zusätzlichen Informationen verloren, sondern es wird lediglich etwas anders dargestellt, was unseren Unparser nicht stört. Bei der Parser-Option Ignore-Empty-Lines kommt es zu einem zusätzlichen Informationsverlust und um den wiederherzustellen, müssen wir diese Informationen speichern, was dem Sinn dieser Parser-Option widerspricht. Aus diesem Grund ist unser Unparser nicht dafür ausgelegt, beim Setzen der Option Ignore-Empty-Lines zu funktionieren, aber wir wollten trotzdem sehen, wie korrekt er auch bei der Auswahl einer nicht vorgesehenen Option funktioniert. Wenn wir jetzt die ersten beiden Parser-Kombinationen in der Tabelle 5.2 betrachten, stellen wir fest, dass die Werte der beiden Zeilen für das jeweilige Korrektheitskriterium gleich sind. Dabei sehen wir, dass die syntaktische Gleichheit kein einziges Mal erfüllt wurde. Es gibt jedoch auch keinen einzigen Fall, dass der Unparser nicht funktioniert hat. Jedes Mal in der Auswertung konnte unser Unparser die syntaktische Gleichheit ohne Whitespace garantieren. Dies lässt uns zu dem Entschluss kommen, dass unser Unparser für das Unparsen von C-Präprozessor-Annotierten Code korrekt funktioniert. Diese Entscheidung basiert auf der Tatsache, dass bei der Entfernung von Whitespace nicht alle Whitespace-Zeichen entfernt werden, sondern nur bestimmte. Nach dieser Manipulation erhalten wir weiterhin einen C-Präprozessor-Annotierten Code. Wir vermuten, dass der Grund, warum die syntaktische Korrektheit kein einziges Mal erfüllt wurde, darin liegt, dass der Parser, auch wenn die Option Ignore-Empty-Lines

nicht ausgewählt wurde, die letzten Zeilen, wenn sie leer sind, nicht speichert. Dieses Vorgehen des Parsers haben wir bemerkt, als wir unseren Unparser getestet haben. Unser Unparser kann das nicht syntaktisch korrekt wiederherstellen, da diese Informationen nicht in den Variation-Trees bzw. Variation-Diffs vorhanden sind. Aus der Tabelle 5.2 können wir auch entnehmen, wie korrekt unser Unparser bei der ausgewählten Parser-Option Ignore-Empty-Lines arbeitet. Dabei sehen wir, dass auch hier kein einziges Mal die syntaktische Gleichheit erreicht wurde, aber im Gegensatz zu vorherigen Fall gibt es hier Fälle, bei denen unser Unparser nicht korrekt unparsen konnte. Das sind 178 Fälle, welches 0,2% aller untersuchten Fälle darstellt. Für uns ist das Grund genug zu sagen, dass auch bei dieser Option unser Unparser in der überwiegenden Mehrheit der Fälle funktioniert. Bei der Kombination der Parser-Optionen, wo beide Optionen ausgewählt sind, ist zu erkennen, dass die Werte zu der vorherigen Kombination gleich sind. Es sind wieder 178 Fälle gegeben, welche nicht ungeparst werden können. Es lässt sich vermuten, dass es sich um dieselben Fälle, wie im vorherigen Fall handelt. Zu diesem Entschluss kann man kommen, da, wenn nur die Parser-Option Collapse-Multiple-Code-Lines ausgewählt wurde, es keine unparsbaren Fälle gab. Wenn die Parser-Option Ignore-Empty-Lines ausgewählt wurde, hatten wir 178 unparsbare Fälle und bei der Auswahl beider Parser-Optionen gab es wieder 178 unparsbare Fälle. Das lässt vermuten, dass diese unparsbaren Fälle durch die Parser-Option Ignore-Empty-Lines verursacht wurden. Damit lässt sich sagen, dass unser Unparser korrekt für die vorgesehenen Kombinationen von Parser-Optionen funktioniert, aber dass der auch in der überwiegenden Mehrheit der Fälle für nicht vorgesehene Kombinationen von Parser-Optionen arbeitet.

Nachdem wir die Forschungsfrage RQ1 und RQ2 für das Unparsen von C-Präprozessor-Annotierten Code beantwortet haben, müssen wir nun noch die Forschungsfragen RQ1 und RQ3 für das Parsen von textbasierten Diffs beantworten. Dazu gehen wir analog zu dem vorherigen Mal vor und schauen uns die Tabelle 5.3 an. In dieser Tabelle sind für die Beantwortung der Frage RQ1 nur die ersten beiden Kombinationen von Parser-Optionen relevant, aus demselben Grund wie zuvor. Es hat damit zu tun, dass der Unparser für Variation-Diffs intern den Unparser für Variation-Trees verwendet. Es ist zu sehen, dass kein einziges Mal die syntaktische Gleichheit erfüllt wurde. Der Grund dafür könnte derselbe wie für das Unparsen von C-Präprozessor-Annotierten Code sein, da der Parser von C-Präprozessor-Annotierten Code intern einen Parser für textbasierten Diffs verwendet. In beiden Fällen ignoriert der Parser die letzte leere Zeile, auch wenn die Parser-Option Ignore-Empty-Lines nicht ausgewählt wurde. Dabei haben wir jedoch Fälle, die nicht ungeparst werden konnten. Wir haben uns diese Fälle angeschaut und herausgefunden, dass die Probleme bei `#endif` auftreten. Dieses Problem tritt innerhalb eines textbasierten Diffs anhand dieser Fehlerfälle nur maximal zweimal auf. Die betrachteten Probleme beinhalten falsche Einrückung, falsche Positionierung, falsche oder fehlende Kommentare und abweichende Darstellung von `#endif`. Alle Arten von Problemen wurden den textbasierten Diffs entnommen, welche nicht mit den vorgesehenen Parser-Optionen ungeparst werden konnten und aus diesem Grund von uns gespeichert wurden. Es lässt sich vermuten, dass Teile dieser Probleme damit zusammenhängen, dass wir ein textbasiertes Diff nicht selbst unparsen, sondern diese Aufgabe auf das Unparsen von C-Präprozessor-Annotiertem Code reduzieren und mit diesem einen textbasierten Diff bilden. Abhängig davon, welcher Diffing-Algorithmus verwendet wurde, können unterschiedliche textbasierte Diffs herauskommen. Ein weiterer Grund für einige der Probleme könnte sein, dass wir bei der Erweiterung der Parser-Implementierung nicht alles beachtet haben und in bestimmten Fällen, welche wir anhand der Fehlerfälle nicht erkennen konnten, die Information von `#endif` nicht korrekt gespeichert wird. Aus dem Grund, dass die Fehlerfälle weniger als 0,03% aller untersuchten textbasierten Diffs darstellen, können wir sagen, dass in der überwiegenden Mehrheit der Fälle unser Unparser funktioniert. In den

ersten beiden Kombinationen von Parser-Optionen sind noch syntaktische Gleichheit ohne Whitespace und semantische Gleichheit zu sehen; für beide Kombinationen sind die Werte gleich. Dabei sehen wir, dass die Anzahl der Fälle für syntaktische Gleichheit ohne Whitespace nicht gleich der Anzahl der Fälle für semantischen Gleichheit ist. Der Grund dafür könnte der von uns für das Unparsen von textbasierten Diffs verwendete Diffing-Algorithmus sein. Unterschiedliche Diffing-Algorithmen können aus derselben Eingabe unterschiedliche textbasierte Diffs konstruieren. Dies kommt dadurch zustande, dass textbasierte Diffs Zeilen enthalten können, welche vertauscht werden können, ohne dass sich der Sinn des textbasierten Diffs ändert. Dadurch bekommen wir semantische Gleichheit, aber keine syntaktische Gleichheit ohne Whitespace bei der Überprüfung. Für die Kombination von Parser-Optionen, wo Ignore-Empty-Lines verwendet wird, lässt sich sehen, dass unser Unparser es auch hier schafft, größtenteils korrekt unparsen und dass nur ungefähr 0,2% aller betrachteten textbasierten Diffs nicht geschafft wurden korrekt unparsen. Aus der gegebenen Information lässt sich von uns ein Schluss ziehen, dass unser Unparser trotz der Fehler sich in der überwiegenden Mehrheit der Fälle für das Unparsen von textbasierten Diffs geeignet ist.

5.3 Zusammenfassung

In diesem Kapitel haben wir unsere Definition von Korrektheitskriterien für den Unparser vorgestellt. Die Korrektheitskriterien sind syntaktische Gleichheit, syntaktische Gleichheit ohne Whitespace und semantische Gleichheit. Es wurde auch gezeigt, wie diese Korrektheitskriterien zusammenhängen. Dazu haben wir eine empirische Auswertung durchgeführt, um anhand der Korrektheitsriterien zu prüfen, ob unser Unparser korrekt funktioniert. Das Ergebnis der Auswertung ist, dass für den Fall, dass keine Parser-Optionen oder nur Collapse-Multiple-Code-Lines verwendet wird, der Unparser in allen betrachteten Fällen mindestens ein Korrektheitskriterium für das Unparsen von C-Präprozessor-Annotierten Code erfüllt. Wenn beim Unparsen von Variation-Trees, bei deren Erstellung die Parser-Option Ignore-Empty-Lines verwendet wurde, dann beträgt die Korrektheit nur 99,8% aller untersuchten Fälle. Das Ergebnis für das Unparsen von Variation-Diffs sieht folgendermaßen aus: Wenn keine oder nur die Parser-Option Collapse-Multiple-Code-Lines für die Erstellung der Variation-Diffs verwendet wurde, beträgt die Korrektheit mehr als 99,97% aller untersuchten textbasierten Diffs. Wenn die Parser-Option Ignore-Empty-Lines verwendet wurde, dann beträgt die Korrektheit 99,8% aller untersuchten textbasierten Diffs.

Verwandte Arbeiten

6.1 Parsen von C-Präprozessor-Annotationen für Variabilitätsanalysen

DiffDetective ist ein Tool, welches hilft, Änderungen in Quellcode und Änderungen der Variabilität darzustellen und den Zusammenhang zwischen ihnen analysierbar zu machen [3, 5]. Um diesen Zweck zu erfüllen hat DiffDetective verschiedene Features. Mit dieser Arbeit haben wir die Anzahl der Features von DiffDetective erweitert und damit auch seine Einsatzmöglichkeiten vergrößert. DiffDetective überführt C-Präprozessor-Annotierten Code oder textbasierten Diff aus C-Präprozessor-Annotierten Code in einen Variation-Tree oder Variation-Diff [5]. Unparzen ist die Umkehrung von Parsen. Zusammen ermöglichen die Features in beide Richtungen umzuwandeln, von C-Präprozessor-Annotierten Code bzw. textbasierten Diff zu Variation-Tree bzw. Variation-Diff und umgekehrt von Variation-Tree bzw. Variation-Diff zu C-Präprozessor-Annotierten Code bzw. textbasierten Diff. Diese Funktionen müssen aber zur korrekten Arbeit aufeinander abgestimmt werden, wie es unsere Arbeit zeigt.

Es gibt aber auch andere Parser für C-Präprozessor-Annotierten Code. Ein solcher Parser arbeitet statt mit dem C-Präprozessor-Annotierten Code, mit dem Token-Stream und generiert einen abstrakten Syntaxbaum [9, 12]. Der Token-Stream, mit welchen der Parser arbeitet, wird aus dem C-Präprozessor-Annotierten Code generiert, von einem Lexer [9, 12]. Dieser Lexer löst Makros auf und fügt Dateien ein. Das führt mit sich, dass der Parser nicht nur mit C-Präprozessor-Annotationen arbeitet, sondern mit dem ganzen C-Präprozessor-Annotierten Code [9, 12]. Aus diesem Grund ist es nicht möglich unseren Unparser hier zu verwenden, da unser Parser nur C-Präprozessor-Annotationen als Kontrollstrukturen wahrnimmt und den C-Code nur als reinen Text betrachtet, was für diesen Parser nicht der Fall ist.

Es ist auch möglich den C-Präprozessor-Annotierten Code zu den Formeln zu parsen [27]. Dort werden die Bedingungen aller Stellen mit bedingten Kompilierung zusammen verknüpft. Dieses Vorgehen ist in Front-End und Back-End aufgeteilt [27]. In Front-End wird der C-Präprozessor-Annotierte Code analysiert und die C-Präprozessor-Annotationen extrahiert [27]. Danach werden im Back-End auf dieser Grundlage die Formeln erstellt [27]. Die Formeln aus C-Präprozessor-Annotation werden bearbeitet und zusammen verknüpft. Das Ergebnis, dieses Vorgangs sind Formeln, was uns eindeutig mitteilt, dass unserer Unparser hier nicht einsetzbar ist. Obwohl es auch in DiffDetective unser Unparser mit Formeln arbeiten muss, sind die

Formeln für den Unparser so umgesetzt, wie sie in C-Präprozessor-Annotationen sind. Unser Unparser kann deshalb sie auch wie ein Text behandeln und muss nicht die extra bearbeiten.

An diesen Möglichkeiten für das Parsen sehen wir, wie unterschiedlich die Ergebnisse des Parsens sein können, welches mit sich bringt, dass das Einsetzen unseres Unparsers nicht ohne weiteres für andere Parser möglich ist.

6.2 Dekompilierung von C

Dekompilierung ist die Umkehrung der Kompilierung, sowie Unparsen die Umkehrung des Parsens ist. Unparser und die Dekompilierung generieren beide aus ihrer Eingabe einen Text bzw. Code. Die Dekompilierung ist in Front-End, Mid-End und Back-End aufgeteilt [6, 7, 19]. Das Front-End hat als Eingabe den Binärcode und verarbeitet diesen zu einer Zwischendarstellung des Programms [7, 19]. Im Mid-End wird aus der Zwischendarstellung des Programms, die Information zum Kontrollfluss entnommen und eine Kontrollflussgraph erstellt. Zum Schluss im Back-End wird aus der Zwischendarstellung des Programms und dem Kontrollflussgraphen der Code erstellt [7, 19]. Im Back-End wird, ebenso wie bei uns, ein Graph erstellt, und die Zwischendarstellung des Programms wird in Code umgewandelt [7, 19]. Trotzdem, dass in beiden Fällen Graphen verwendet werden, können wir nicht den Ansatz der Dekompilierung verwenden, da bei der Dekompilierung mit dem Binärcode, den Zwischenergebnissen und konkretem C-Code gearbeitet wird [6, 7, 19]. In unseren Fall wird jedoch nur mit den C-Präprozessor-Annotationen interagiert und nicht mit dem C-Code oder dem Binärcode. In unserem Fall ist es sogar nicht der C-Code, da wir den C-Code nur als Text betrachten und mit dem Code nicht explizit arbeiten. In dieser Hinsicht ist unser Unparser generell, solange die C-Präprozessor-Annotationen verwendet werden, kann der restliche Inhalt ein beliebiger Text sein. Bei dem Decompiler ist es anders, da ein Dekompiler speziell für eine Programmiersprache entwickelt wird. Der C-Dekompiler wird für den C-Code entwickelt, der Java-Dekompilr für den Java-Code und so weiter. Der Dekompiler interagiert nicht mit den C-Präprozessor-Annotationen. Unseren Unparser und den Dekompiler kann man als parallele Prozesse betrachten, da diese an verschiedenen Stellen ansetzen.

6.3 Variabilitätsanalysen

Bei Analysen oder Verarbeitungsschritten kann es vorkommen, dass der Code auf abstrakter Ebene verändert wird. Danach muss man Unparsen, um einen Code wieder zubekommen. Ein Beispiel für so etwas sind die Mutation-Tests [1, 18], welche eine der vielversprechendsten Techniken zur Bewertung der Effektivität von Testfällen sind. Die Mutation-Tests lassen sich auch zur Analyse der Testfälle, welche die Variabilität sichern sollten, verwenden [1, 18]. Allgemein bei den Mutation-Tests wird der Code verändert. Diese Veränderungen sind nicht willkürlich, sondern haben gewisse Eigenschaften [1, 18]. Nach diesen Veränderungen erhält man neue, leicht modifizierte Kopien des Ausgangscodes. Auf diese Kopien werden dann die Testfälle für den Ausgangscode verwendet und geschaut, ob die modifizierte Version des Codes einen Fehler erzeugt oder nicht [1, 18]. Dadurch lässt sich Feststellen, wie gut die Testfälle ausgewählt wurden. Dieses Vorgehen lässt sich auch im Kontext der Variabilität verwenden. Dabei werden die Veränderung nur an den Stellen erzeugt, welche für die Variabilität zuständig sind [1, 18]. Im Falle, dass für die Erzeugung der Variabilität der C-Präprozessor verwendet wird, sind das die C-Präprozessor-Annotationen. In Bezug auf unsere Arbeit ist nur dieser Fall für uns von Bedeutung. Dabei wird am Anfang der Code zu einem Variation-Tree geparst. Danach werden die Veränderungen angewandt. An dieser Stelle wird unser Unparser angewendet und überführt das Variation-Tree in ausführbaren Code, auf den die Tests angewandt werden können und die Mutation-Tests werden

durchgeführt.

Ein weiteres Beispiel ist die Analyse, welche Codeblöcke finden soll, die in keiner Variante auftauchen [30]. Bei dieser Analyse werden zwei Datenbanken erstellt. Die erste Datenbank enthält das Variabilität-Model, welches zeigt in welchen Kombinationen die Features ausgewählt werden dürfen [30]. Die zweite Datenbank enthält den Variabilität, welche tatsächlich in Code umgesetzt wurde [30]. Damit dieser erhalten wird, muss man den Code Parsen. Danach ist es möglich für beide Datenbanken jeweils die Feature-Abhängigkeiten in ein binäres Entscheidungsdiagramm umzuwandeln. Anhand dessen man die Erfüllbarkeit überprüfen kann [30]. Es liegt nahe, danach die niemals erfüllbaren Codeblöcke zu entfernen und die restlichen unzu parsen, um einen neuen Code zu bekommen.

Bei Variabilitätsanalysen ist es auch möglich zu prüfen, ob alle Varianten syntaktisch korrekt sind. Dabei muss zuerst der Präprozessor-Annotierter Code zu einem abstrakten Syntaxbaum geparkt werden [11]. Danach wird dieser abstrakter Syntaxbaum von dem Tool CIDE eingelesen. CIDE ist ein Tool, welches erlaubt verschiedene Features verschiedenfarbig einzufärben [11]. Im Zusammenhang mit der Analyse der syntaktischen Korrektheit aller Varianten ist von Bedeutung das CIDE nur das Einfärben der Codefragmente, welche optional für diejenige Programmiersprache sind [11]. Also, sodass bei Entfernen der eingefärbten Codefragmente die Syntax trotzdem gültig bleibt.

Noch eine weitere Variabilitätsanalyse ist die Prüfung der Typkorrektheit für alle Varianten. Vor der Analyse muss der Parser den Präprozessor-Annotierten Code in einen bestimmten Kalkül überführen [10]. Danach wird anhand der Vorschriften des Kalküls geprüft, um Stellen zu finden, wo der Typ nicht eingehalten wurde und dieses beheben [10]. Schließlich kann man den Kalkül zu einem Präprozessor-Annotierten Code unparsen und den weiter bearbeiten.

Fazit und Zukünftige Arbeiten

Wir entwickelten einen Algorithmus zum Unparsen von Variation-Trees und ein Vorgehen zum Unparsen von Variation-Diffs. Variation-Trees und Variation-Diffs sind Datenstrukturen zur Abbildung präprozessor-orientierte Variabilität und der Änderungen daran. Um das zu bewerkstelligen haben wir uns mit Definitionen von Variation-Trees und Variation-Diffs beschäftigt, sowie die gegebenen Definitionen für unsere Zwecke erweitern. Damit wir den Definitionen gerecht werden, haben wir uns mit dem Parser beschäftigt, um den zu erweitern. Den ausgearbeiteten Algorithmus zum Unparsen von Variation-Trees und das Vorgehen zum Unparsen von Variation-Diffs haben wir vorgestellt und in DiffDetective implementiert. Um die Korrektheit unseres Unparsers nachzuweisen, haben wir Korrektheitskriterien entwickelt und unseren Unparser in einer empirischen Auswertung anhand dieser Kriterien getestet. Die ausgearbeiteten Korrektheitskriterien sind syntaktische Gleichheit, syntaktische Gleichheit ohne Whitespace und semantische Gleichheit.

Anhand der durchgeführten Auswertung auf den Datensätzen `Vim`, `sylpheed` und `berkeley-db-libdb` haben wir erkenntlich gemacht für welche Parser-Optionen, wie korrekt unser Unparser funktioniert. Die Auswertung der Korrektheit des Unparsers wurde auf den von uns definierten Korrektheitskriterien durchgeführt. Für die Parser-Optionen, für die unser Unparser für das Unparsen von Variation-Trees entwickelt wurde, funktioniert der korrekt. Dies haben wir an 98802 Testfällen gesehen. Der Unparser für das Unparsen von Variation-Diffs funktioniert leider nicht so gut, 11 der 49401 Testfälle konnten nicht ungeparst werden, welches weniger als 0,03% aller Testfälle darstellt.

Dank dem von uns ausgearbeiteten Unparsers ist es möglich geworden Verfahren durchzuführen, welche zuerst einen Variation-Tree bzw. Variation-Diff erstellen, den bearbeiten und dann dies als C-Präprozessor-Annotierten Code bzw. textbasierten Diff brauchen für weitere Bearbeitungen. Zu solchen Verfahren gehört Mutation-Tests für Variabilität. Damit sind die Möglichkeiten für Analysen der Variabilität, für welche die Entwickler DiffDetective verwenden können, größer geworden.

Wir sehen zwei mögliche Punkte für weitere Forschung. Der erste wäre die Ausarbeitung und Implementierung eines Algorithmus zum Unparsen von Variation-Diffs. Unser Algorithmus ist nur für das Unparsen von Variation-Trees ausgelegt und nicht Variation-Diffs. Damit auch Variation-Diffs ungeparst werden können, reduzieren wir das Problem auf das Unparsen von Variation-Trees. Zwar funktioniert unser Vorgehen, aber wie die Auswertung zeigt, funktioniert

unser Vorgehen nicht immer. Ein eigener Algorithmus zum Unparsen von Variation-Diffs könnte das vielleicht beseitigen. Dazu könnten wir damit die Abhängigkeit von dem Differ-Algorithmus loswerden und die Korrektheit steigern. Als ein wichtiges Problem zum Entwickeln eines Algorithmus zum Unparsen von Variation-Diffs finden wir es, dem Variation-Diff, welcher ein azyklischer Graph ist, zu entnehmen in welcher Reihenfolge die Knoten des Graphen besucht werden sollen. Der zweite Punkt für mögliche zukünftige Arbeiten wäre, dass man sowohl Parser als auch den Unparser auf die Arbeit mit anderen Präprozessoren erweitert. Zurzeit können sowohl Parser als auch der Unparser nur mit dem C-Präprozessor arbeiten. Ein möglicher Kandidat wäre der Java-Präprozessor, bei denen die Anweisungen an den C-Präprozessor angelehnt sind.

Literaturverzeichnis

- [1] M. Al-Hajjaji, F. Benduhn, T. Thüm, T. Leich, and G. Saake. Mutation Operators for Preprocessor-Based Variability. In *Proc. Int'l Workshop on Variability Modelling of Software-Intensive Systems (VaMoS)*, pages 81–88. ACM, 2016.
- [2] S. Apel, D. Batory, C. Kästner, and G. Saake. *Feature-Oriented Software Product Lines*. Springer, 2013.
- [3] P. M. Bittner, C. Tinnes, A. Schultheiß, S. Viegner, T. Kehrer, and T. Thüm. Classifying Edits to Variability in Source Code. In *Proc. Europ. Software Engineering Conf./Foundations of Software Engineering (ESEC/FSE)*, pages 196–208. ACM, 2022.
- [4] P. M. Bittner, A. Schultheiß, S. Greiner, B. Moosherr, S. Krieter, C. Tinnes, T. Kehrer, and T. Thüm. Views on Edits to Variational Software. In *Proc. Int'l Systems and Software Product Line Conf. (SPLC)*, pages 141–152. ACM, 2023.
- [5] P. M. Bittner, A. Schultheiß, B. Moosherr, T. Kehrer, and T. Thüm. Variability-Aware Differencing with DiffDetective. In *Companion Proc. Int'l Conference on the Foundations of Software Engineering (FSE Companion)*, pages 632–636. ACM, 2024.
- [6] Y. Cao, R. Zhang, R. Liang, and K. Chen. Evaluating the effectiveness of decompilers. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ISSTA 2024, page 491–502, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400706127. doi: 10.1145/3650212.3652144.
- [7] C. G. Cifuentes and K. J. Gough. A methodology for decompilation. 1993. URL <https://api.semanticscholar.org/CorpusID:11292649>.
- [8] M. J. Fischer. *Efficiency of Equivalence Algorithms*, pages 153–167. Springer US, Boston, MA, 1972. ISBN 978-1-4684-2001-2. doi: 10.1007/978-1-4684-2001-2_14.
- [9] P. Gazzillo and R. Grimm. SuperC: Parsing All of C by Taming the Preprocessor. In *Proc. ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI)*, pages 323–334. ACM, 2012.
- [10] C. Kästner and S. Apel. Type-Checking Software Product Lines—A Formal Approach. In *Proc. Int'l Conf. on Automated Software Engineering (ASE)*, pages 258–267. IEEE, 2008.
- [11] C. Kästner, S. Apel, and M. Kuhlemann. Granularity in Software Product Lines. In *Proc. Int'l Conf. on Software Engineering (ICSE)*, pages 311–320. ACM, 2008.

- [12] C. Kästner, P. G. Giarrusso, T. Rendel, S. Erdweg, K. Ostermann, and T. Berger. Variability-Aware Parsing in the Presence of Lexical Macros and Conditional Compilation. In *Proc. Conf. on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 805–824. ACM, 2011.
- [13] C. Kästner, S. Apel, T. Thüm, and G. Saake. Type Checking Annotation-Based Product Lines. *Trans. on Software Engineering and Methodology (TOSEM)*, 21(3):14:1–14:39, 2012.
- [14] T. Kehrer, T. Thüm, A. Schultheiß, and P. M. Bittner. Bridging the Gap Between Clone-and-Own and Software Product Lines. In *Proc. Int’l Conf. on Software Engineering (ICSE)*, pages 21–25. IEEE, 2021.
- [15] J. Krüger and T. Berger. An Empirical Analysis of the Costs of Clone- and Platform-Oriented Software Reuse. In *Proc. Europ. Software Engineering Conf./Foundations of Software Engineering (ESEC/FSE)*, pages 432–444. ACM, 2020.
- [16] J. Krüger and T. Berger. Activities and Costs of Re-Engineering Cloned Variants Into an Integrated Platform. In *Proc. Int’l Working Conf. on Variability Modelling of Software-Intensive Systems (VaMoS)*. ACM, 2020.
- [17] E. Kuitert, J. Krüger, S. Krieter, T. Leich, and G. Saake. Getting Rid of Clone-and-Own: Moving to a Software Product Line for Temperature Monitoring. In *Proc. Int’l Systems and Software Product Line Conf. (SPLC)*, pages 179–189. ACM, 2018.
- [18] H. Lackner and M. Schmidt. Towards the Assessment of Software Product Line Tests: A Mutation System for Variable Systems. In *Proc. Workshop on Software Product Line Analysis Tools (SPLat)*, pages 62–69. ACM, 2014.
- [19] Z. Liu and S. Wang. How far we have come: testing decompilation correctness of c decompilers. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis*, ISSTA 2020, page 475–487, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450380089. doi: 10.1145/3395363.3397370.
- [20] B. Moosherr. Constructing Variation Diffs Using Tree Diffing Algorithms. Bachelor’s thesis, University of Ulm, 2023.
- [21] L. Neves, P. Borba, V. Alves, L. Turnes, L. Teixeira, D. Sena, and U. Kulesza. Safe Evolution Templates for Software Product Lines. *J. Systems and Software (JSS)*, 106:42–58, 2015.
- [22] M. Nieke, G. Sampaio, T. Thüm, C. Seidl, L. Teixeira, and I. Schaefer. Guiding the Evolution of Product-Line Configurations. *Software and Systems Modeling (SoSyM)*, 21:225–247, 2022.
- [23] L. Passos, K. Czarnecki, S. Apel, A. Wąsowski, C. Kästner, and J. Guo. Feature-Oriented Software Evolution. In *Proc. Int’l Workshop on Variability Modelling of Software-Intensive Systems (VaMoS)*, pages 1–8. ACM, 2013.
- [24] L. Passos, L. Teixeira, N. Dintzner, S. Apel, A. Wąsowski, K. Czarnecki, P. Borba, and J. Guo. Coevolution of Variability Models and Related Software Artifacts. *Empirical Software Engineering (EMSE)*, 21(4), 2016.
- [25] G. Sampaio, P. Borba, and L. Teixeira. Partially Safe Evolution of Software Product Lines. *J. Systems and Software (JSS)*, 155:17–42, 2019.

- [26] C. Seidl, F. Heidenreich, and U. Aßmann. Co-Evolution of Models and Feature Mapping in Software Product Lines. In *Proc. Int’l Systems and Software Product Line Conf. (SPLC)*, pages 76–85. ACM, 2012.
- [27] J. Sincero, R. Tartler, D. Lohmann, and W. Schröder-Preikschat. Efficient Extraction and Analysis of Preprocessor-Based Variability. In *Proc. Int’l Conf. on Generative Programming and Component Engineering (GPCE)*, pages 33–42. ACM, 2010.
- [28] J. Sprey, C. Sundermann, S. Krieter, M. Nieke, J. Mauro, T. Thüm, and I. Schaefer. SMT-Based Variability Analyses in FeatureIDE. In *Proc. Int’l Working Conf. on Variability Modelling of Software-Intensive Systems (VaMoS)*. ACM, 2020.
- [29] C. Sundermann, M. Nieke, P. M. Bittner, T. Heß, T. Thüm, and I. Schaefer. Applications of #SAT Solvers on Feature Models. In *Proc. Int’l Working Conf. on Variability Modelling of Software-Intensive Systems (VaMoS)*. ACM, 2021.
- [30] R. Tartler, J. Sincero, W. Schröder-Preikschat, and D. Lohmann. Dead or Alive: Finding Zombie Features in the Linux Kernel. In *Proc. Int’l Workshop on Feature-Oriented Software Development (FOSD)*, pages 81–86. ACM, 2009.
- [31] S. Viegner. Empirical Evaluation of Feature Trace Recording on the Edit History of Marlin. Bachelor’s thesis, University of Ulm, 2021.
- [32] B. Zhang and M. Becker. Variability code analysis using the vital tool. In *Proceedings of the 6th International Workshop on Feature-Oriented Software Development, FOSD ’14*, page 17–22, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450329804. doi: 10.1145/2660190.2662113.
- [33] S. Zhou, Ș. Stănciulescu, O. Leßenich, Y. Xiong, A. Wąsowski, and C. Kästner. Identifying Features in Forks. In *Proc. Int’l Conf. on Software Engineering (ICSE)*, pages 105–116. ACM, 2018.