

Trace Link Recovery using Static Program Analysis

Bachelorarbeit

zur Erlangung des Grades eines Bachelor of Science
im Studiengang Informatik

vorgelegt von

Maximilian Meffert

Erstgutachter: Prof. Dr. Ralf Lämmel
Institut für Informatik

Zweitgutachter: Msc. Johannes Härtel
Institut für Informatik

Koblenz, im Dezember 2017

Erklärung

Hiermit bestätige ich, dass die vorliegende Arbeit von mir selbständig verfasst wurde und ich keine anderen als die angegebenen Hilfsmittel – insbesondere keine im Quellenverzeichnis nicht benannten Internet-Quellen – benutzt habe und die Arbeit von mir vorher nicht in einem anderen Prüfungsverfahren eingereicht wurde. Die eingereichte schriftliche Fassung entspricht der auf dem elektronischen Speichermedium (CD-Rom).

Ja Nein

Mit der Einstellung der Arbeit in die Bibliothek bin ich ein- ☐ ☐
verstanden.

Der Veröffentlichung dieser Arbeit im Internet stimme ich ☐ ☐
zu.

.....
(Ort, Datum) (Maximilian Meffert)

Zusammenfassung

TBD.

Abstract

TBD.

Acknowledgements

TBD.

Contents

1	Introduction	1
1.1	Motivational Example	1
1.2	Objectives	3
1.3	Approach	4
1.4	Contributions & Non-Contributions	5
1.5	Road-Map	5
2	Background	6
2.1	Axioms of Linguistic Architectures	6
2.1.1	Parthood	7
2.1.2	Fragments	9
2.1.3	Correspondence	10
2.1.4	Conformance	11
2.2	Traceability	12
2.2.1	Traceability Objectives	12
2.2.2	Traceability Terminology	13
2.3	Technological Background	14
2.3.1	MagaL & MegaL/Xtext	14
2.3.2	JAXB (Java Architecture for XML Binding)	15
2.3.3	Hibernate	17
2.3.4	ANTLR (Another Tool For Language Recognition)	17
3	Related Work	19
3.1	Paper Summary	19
3.2	Remarks	21

4	Design	23
4.1	Requirements	23
4.2	Recovery Design	24
4.2.1	Recovery Process	24
4.2.2	Recovery API	24
4.2.3	Recovery System	31
4.3	Megal/Xtext Integration Design	32
5	Implementation	34
5.1	Recovering Fragments	34
5.1.1	Concrete Fragment Models	35
5.1.2	Megamodeling Fragments	38
5.2	Recovering Parthood Links	40
5.3	Recovering Correspondence Links	40
5.4	Recovering Conformance Links	40
6	Mini Case-Study	41
6.1	Example Corpus	41
6.1.1	The 101HRMS Model	41
6.1.2	Linguistic Domains of the Example Corpus	42
7	Conclusion	44
7.1	Future Work	44
	Glossary	45

List of Theorems

1	Axiom (partOf)	7
2	Axiom (Fragment)	9
3	Axiom (correspondsTo)	10
4	Axiom (conformsTo)	12
1	Definition (Trace)	13
2	Definition (Trace Artifact, Source Artifact, Target Artifact)	13
3	Definition (Trace Artifact Type)	13
4	Definition (Trace Link)	13
5	Definition (Trace Link Type)	13
6	Definition (Traceability)	13
7	Definition (Traceability Creation)	14
8	Definition (Trace Recovery)	14

List of Figures

1.1	O/R/X-Mapping Manifestations	2
1.2	JAXB XML/XSD Mapping	2
1.3	Recovery Approach	4
2.1	Simple vs. Proper Parthood	8
2.2	Steps in the JAXB Binding Process [7]	16
4.1	The Recovery Process	24
4.2	The Recovery API	25
4.3	The Abstract Fragment Model	27
4.4	The Syntax Analysis API	28
4.5	IParserFactory Usage Example	28
4.6	AntlrParser Parse Tree Creation	28
4.7	AntlrParser Fragment Creation	29
4.8	Mereological Fragment Analysis API	29
4.9	Comparative Fragment Analysis API	30
4.10	The Recovery System	31
4.11	Integration of the Recovery System int MegaL/Xtext	33
4.12	MegaL Plug-In Instantiation	33
5.1	Idealized Recovery of Java Fragments	36
5.2	Java Fragment AST Model	37
5.3	Construction of Java Fragment ASTs	38
5.4	A Megamodel for Java Fragments	39
6.1	The 101 Human Resource Management System Model	42
6.2	Example Corpus Domains: Java O/R/X	43

Chapter 1

Introduction

ToDo: something about program comprehension

1.1 Motivational Example

Common tasks in software development are implementation of serialization and persistence of domain models. For instance, consider a simple web service which serves data via HTTP (Hypertext Transfer Protocol) as XML (Extensible Markup Language) and stores it in a relational database. We call this an O/R/X-Mapping (Object-Relational- and XML-Mapping) scenario. Given such a system, the same conceptual data, i.e. the domain model, is transported through application tiers in different forms, that is, the same data is represented by various manifestations at a time. Each manifestation involves another software language and technology. Figure 1.1 opposes model- to instance-level syntaxes of different manifestations for a conceptual model for employees in an O/R/X-Mapping scenario. One employee only has a *name* and an *age* attribute. Persistence is presented in tuple notation, data transfer objects of an application tier are represented with a fictional high level programming language object notation and serialization is exemplified with XML and XSD (XML Schema Definition) syntax.

Although figure 1.1 shows a fictionalized example, one can observe structural similarities among the different model and instance representation syntaxes. Such similarities also occur in real-world software engineering through use of conventions which may be predefined in technologies for O/R/X-Mapping. Figure 1.2 displays an annotated Java class and XML/XSD output generated by JAXB

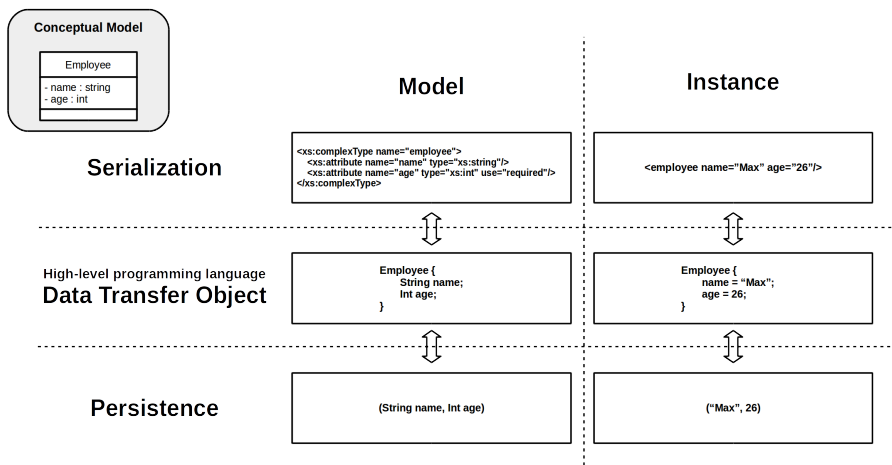


Figure 1.1 O/R/X-Mapping Manifestations

<pre>@XmlRootElement(name="employee") @XmlAccessorType(XmlAccessType.FIELD) public class Employee { @XmlAttribute private int id; @XmlAttribute private String name; @XmlAttribute private int age; @XmlAttribute private double salary; private Department department; private Department managedDepartment; }</pre>	<pre><xs:complexType name="employee"> <xs:attribute name="id" type="xs:int" use="required" /> <xs:attribute name="name" type="xs:string" /> <xs:attribute name="age" type="xs:int" use="required" /> <xs:attribute name="salary" type="xs:double" use="required" /> <xs:sequence> <xs:element ref="department" minOccurs="0" /> <xs:element name="managedDepartment" type="department" minOccurs="0" /> </xs:sequence> </xs:complexType></pre>
<pre><employee name="Max" age="26" salary="55000.0" /></pre>	

Figure 1.2 JAXB XML/XSD Mapping

(Java Architecture for XML Binding). Except for the root element all names for attributes and elements are taken from the Java class. Moreover, the XSD complex type and the Java class share a similar nested structure. Links through similarity can be found within the model-level representations (Java-XSD) and between model- and instance-level manifestations (Java-XML and XSD-XML), for instance:

- `public class Employee {...}` is linked with `<xs:complexType name="employee">...</xs:complexType>` and `<employee .../>`, the latter two are also linked with each other
- `private String name;` is linked with `<xs:attribute name="name" type="xs:string"/>` and `name="Max"`, the latter two are also linked with each other

1.2 Objectives

The aim of this thesis is to provide automated recovery of such similarities as semantic links. Recovered links are inserted into *megamodels* [4] [11] for *linguistic architectures* [10] [24] [19]. Megamodels are models providing a high level of abstraction with other models as modeling elements, e.g. a megamodel may describe the dependencies between metamodels, models and instances. Linguistic architectures intend to describe software systems from a language centric point of view. They model knowledge about software systems in terms of languages, artifacts, technologies, etc. Such entities are interrelated with relationships derived from common software engineering and theoretical computer science vocabulary providing special semantics, e.g. *defines*, *isA*, *instanceOf*, *represents*, *implements*, *realizationOf*, *elementOf*, *subsetOf*.

Linguistic architectures are related to ER Models (Entity-Relationship Models) and ontologies. Recovering semantic links as described above is related to the concept *traceability* and an application of *traceability recovery* [16] (see §2.2). In context of traceability, semantic links may be called *trace links*, however, both establish a relation between two entities denoting a certain meaning.

Trace links among software artifacts denoting that one artifact encodes the same information as the other are called *correspondence* links and denoted *correspondsTo* (see §2.1.3). Trace links between two artifacts denoting that one defines the other, in the sense of a metamodel defining a model, are called *conformance* links and denoted *conformsTo* (see §2.1.4). The main objective is to recover such trace links among well-formed, possibly partial, source code artifacts. We call well-formed partial source code artifacts *fragments* (see §2.1.2).

1.3 Approach

Our approach for recovering trace links utilizes *static program analysis*. We use it in the broad sense that we construct specialized ASTs (Abstract Syntax Trees) for further analysis. Usually static program analysis is used with the purpose of formal verification or measurement of software engineering related properties of source code, e.g. cohesion and coupling metrics. We use static program analysis to semantically link code fragments and uncover knowledge of the analyzed artifacts. In that respect, our approach is related to computing software metrics.

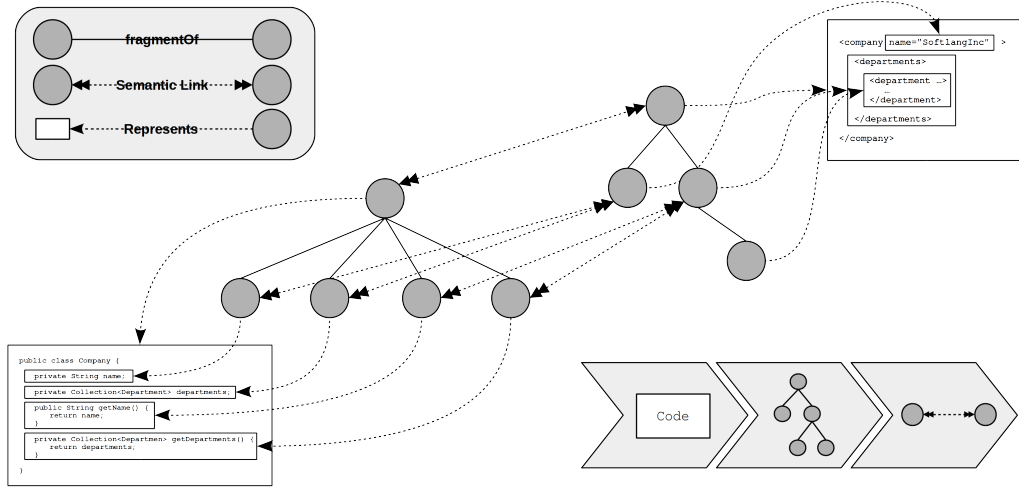


Figure 1.3 Recovery Approach

Figure 1.3 shows a schematic illustration of our recovery approach. It can be summarized with the following three steps:

1. we create two Parse Trees for both inputs respectively
2. Parse Trees are transformed to specialized ASTs; such ASTs are not intended to be used for compilation, instead each node represents a part of code which may be linked during further analysis; a child node represents a piece of code which is properly embedded the code represented by its parent
3. we apply a comparative analysis to both ASTs; while traversing through both trees we check whether each pair of nodes is a trace link

1.4 Contributions & Non-Contributions

This section lists critical contributions and non-contributions of this thesis:

Contribution 1 We contribute a Java based API (Application Programming Interface) for recovering trace links among artifacts and their constituent fragments, providing utilities for static program analysis as described in §1.3 (see §4).

Contribution 2 We contribute an implementation for correspondence and conformance link recovery among JAXB and Hibernate artifacts utilizing the created API and axioms for linguistic architectures (see §5 and §2.1 respectively).

Contribution 3 We contribute a small case-study applying the implemented recovery system to a minimal O/R/X-Mapping program (see §6).

Non-Contribution 1 We do not contribute to the axiomatization of linguistic architectures as described in [24], [11], [23] and [19].

1.5 Road-Map

§2 provides the necessary theoretical and practical background for this thesis. §3 provides a summary on related work and positions this thesis among these works. §4 provides a high level design description for the system developed for this thesis. §5 provides a more detailed description of the implementation for crucial components of the developed system. §6 provides a mini case-study evaluating the developed system. §7 concludes the thesis.

Chapter 2

Background

This chapter provides the necessary theoretical and technological background topics of the thesis.

2.1 Axioms of Linguistic Architectures

This section summarizes axioms of linguistic architectures. Axioms are outlined in [23] and refined in [19]. §2.1.1 introduces the axiomatization for parthood §2.1.2

Axioms are presented in First Order Logic and formalization is taken from [19]. The universe to draw elements from is represented by the Entity-predicate:

$$\forall x. \text{Entity}(x).$$

We assume it holds for everything of interest, hence such things are called *entities*.

Linguistic architectures intend to describe software from language centric point of view, thus we provide specializations of entities for languages:

$$\text{Set}(x) \Rightarrow \text{Entity}(x).$$

$$\text{Language}(x) \Rightarrow \text{Set}(x).$$

There are entities representing sets in a mathematical sense, and there are sets representing (formal-) languages in the sense of theoretical computer science.

On the other hand, we provide specializations for entities involved in software engineering terminology:

$$\text{Artifact}(x) \Rightarrow \text{Entity}(x).$$

$$\text{File}(x) \Rightarrow \text{Artifact}(x).$$

$$\text{Folder}(x) \Rightarrow \text{Artifact}(x).$$

There are entities representing all kinds of digital artifacts, e.g. files and folders. Files represent persistent data resources, locatable either through file systems or web services. Folders represent locatable collections of files. The intended use and semantic of these predicates is not meant to differ from intuitive, every day use.

2.1.1 Parthood

Parthood is the essential relationship when reasoning about correspondence and conformance among artifacts within linguistic architectures [23] [19]. It describes the relation between entities and their constituent parts. The study of parthood and its derivatives is mereology [31] [30]. In the context of linguistic architectures we assume most entities to be composed of several conceptual or physical parts. In short, such entities are the sum of its parts. That is, programs may be compiled from many files, systems consist of several disjoint but dependent components, a Java class is made up of methods and fields, etc. Furthermore such entities are considered to be *mereologically invariant*, i.e. if one part changes, the whole changes as well. Axiom 1 captures parthood at its most basic level.

Axiom 1 (partOf)

$$\text{partOf}(p, w) \Rightarrow \text{Entity}(p) \wedge \text{Entity}(w).$$

$$\text{partOf}(p, w) \Leftarrow p \text{ is a constituent part of } w.$$

Parthood is usually considered to be reflexive, antisymmetric and transitive [31] [30], thus facilitating a partial order:

$\text{partOf}(p, p).$	Reflexivity
$\text{partOf}(p, w) \wedge \text{partOf}(w, p) \Rightarrow p = w.$	Antisymmetry
$\text{partOf}(p, w) \wedge \text{partOf}(w, u) \Rightarrow \text{partOf}(p, u).$	Transitivity

The irreflexive parthood relationship is called proper:

$$\begin{aligned} \text{properPartOf}(p, w) &\Rightarrow \text{Entity}(p) \wedge \text{Entity}(w). \\ \text{properPartOf}(x, y) &\Leftarrow \text{partOf}(x, y) \wedge \neg \text{partOf}(y, x). \end{aligned}$$

Proper parthood is the strict order induced by simple parthood. Figure 2.1 shows a schematic illustration opposing simple to proper parthood. This Venn-style

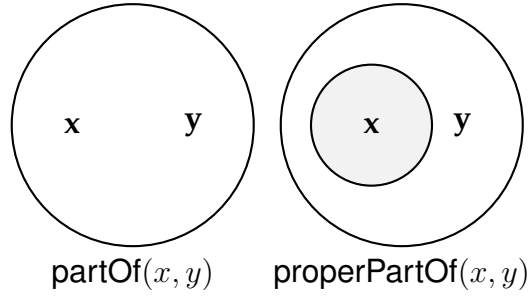


Figure 2.1 Simple vs. Proper Parthood

diagram depicts two scenarios

partOf(x,y) x may be a part of y or vice-versa; x and y cannot be distinguished

properPartOf(x,y) x is certainly a part of y , however y is not a part of x

In general proper parthood is an asymmetric relationship:

$$\text{properPartOf}(x, y) \Rightarrow \neg \text{properPartOf}(y, x). \quad \text{Asymmetry}$$

Mereology also allows for the notion of atomicity [31] [30], i.e. atomic parts which cannot be decomposed in further parts:

$$\text{atomicPart}(x) \Leftarrow \nexists p. \text{properPartOf}(p, x).$$

Atomicity is later used to distinguish cases for correspondence and conformance (see §2.1.3 and §2.1.4).

2.1.2 Fragments

Fragments are a specialization of entities intended to capture the endogenous, mereological decomposition of artifacts [19]. Axiom 2 defines fragments as artifacts, which are neither files nor folders, and are properly embedded in at least one other artifact. For instance, consider the syntactical decomposition of Java classes, i.e. the class fragment contains method and field fragments.

Axiom 2 (Fragment)

$$\begin{aligned} \text{Fragment}(f) &\Rightarrow \text{Artifact}(a) \wedge \neg(\text{File}(f) \vee \text{Folder}(f)). \\ \text{Fragment}(f) &\Rightarrow \exists a. \text{Artifact}(a) \wedge \text{properPartOf}(f, a). \end{aligned}$$

We emphasize the proper parthood here, since not all authors necessarily consider `partOf` to be reflexive [23] [19]. However, we previously distinguished between reflexive and irreflexive parthood and if we were to use reflexive parthood for the definition of the `Fragment`-predicate it would be a tautology:

$$\text{Fragment}(f) \Rightarrow \exists a. \text{Artifact}(a) \wedge \text{partOf}(f, a). \quad \text{Tautology}$$

Since `Fragment(f)` implies `Artifact(f)` and `partOf` is reflexive there is always an artifact for a fragment the latter is a part of, that is the fragment itself.

From the specification of fragments we can first specialize the proper parthood predicate:

$$\begin{aligned} \text{fragmentOf}(f, x) &\Rightarrow \text{Fragment}(f) \wedge \text{Artifact}(x). \\ \text{fragmentOf}(f, x) &\Leftarrow \text{Fragment}(f) \wedge \text{Artifact}(x) \wedge \text{properPartOf}(f, x). \end{aligned}$$

This is just a restriction of domain and range facilitating a special semantic. Secondly, we can describe the process of fragmentation:

$$\text{partOf}(p, w) \wedge \text{Fragment}(w) \Rightarrow \text{Fragment}(p).$$

The fragment nature propagates top-down alongside the order of parthood, i.e. if an entity is a fragment all parts are also fragments [19].

2.1.3 Correspondence

Correspondence is the relationship between two artifacts denoting both represent the same data in the sense that they are mereologically similar [19]. It usually occurs as exogenous relationship, i.e. the involved artifacts do not belong to the same language, for instance XML and JSON (JavaScript Object Notation) may contain the same information.

In order to axiomatize correspondence properly we need to introduce two other relationships first:

represents The relationship denoting that an artifact is a representation or manifestation of an entity.

$$\text{represents}(a, e) \Rightarrow \text{Artifact}(a) \wedge \text{Entity}(e).$$

$$\text{represents}(a, e) \Leftarrow a \text{ is a representation of } e.$$

sameAs The relationship denoting two artifacts represent the same entity.

$$\text{sameAs}(x, y) \Rightarrow \text{Artifact}(x) \wedge \text{Artifact}(y).$$

$$\text{sameAs}(x, y) \Leftarrow \exists e. \text{represents}(x, e) \wedge \text{represents}(y, e).$$

The axiomatization of correspondence from [19] is slightly altered with emphasis of irreflexive proper parthood, Otherwise it would not work with axiom 1.

Axiom 3 (correspondsTo)

$$\text{correspondsTo}(x, y) \Rightarrow \text{Artifact}(x) \wedge \text{Artifact}(y).$$

$$\text{correspondsTo}(x, y) \Leftarrow (\forall px. \text{properPartOf}(px, x) \Rightarrow \exists py. \text{properPartOf}(py, y) \wedge \text{correspondsTo}(px, py))$$

$$\wedge (\forall py. \text{properPartOf}(py, y) \Rightarrow \exists px. \text{properPartOf}(px, x) \wedge \text{correspondsTo}(py, px))$$

$$\vee (\neg \text{properPartOf}(p, x) \vee \text{properPartOf}(p, y)) \wedge \text{sameAs}(x, y).$$

Axiom 3 defines correspondence recursively:

Case I If both artifacts is atomic in the sense it cannot be decomposed in further parts, we check whether both artifacts represent the same data.

Case II Else, if at least one artifact is not atomic, then for each proper part of one artifact has to be a corresponding part in the other and vice versa.

This axiomatization demands a strict one-to-one correspondence in its recursive clause which occurs in reflexive cases, i.e. $\text{correspondsTo}(x, x)$, but may be too strict for real-world applications [23].

2.1.4 Conformance

Conformance is the relationship between two artifacts denoting one satisfies the syntax specification given by the other [19]. This satisfaction is also meant to mereologically decomposable, i.e. there are parts of one artifact specified by a part of the other.

In order to axiomatize conformance properly we need to introduce two other relationships first:

defines The relationship simply denoting an artifact contains the definition of an entity.

$$\text{defines}(a, e) \Rightarrow \text{Artifact}(a) \wedge \text{Entity}(e).$$

$$\text{defines}(a, e) \Leftarrow a \text{ is a definition for } e.$$

elementOf The relationship denoting set-theoretic membership.¹

$$\text{elementOf}(e, s) \Rightarrow \text{Entity}(e) \wedge \text{Set}(s).$$

$$\text{elementOf}(e, s) \Leftarrow e \in s.$$

The axiomatization of conformance from [19] is also slightly altered with emphasis of irreflexive proper parthood, Otherwise it would not work with axiom 1.

¹We do not use the `elementOf` axiom for artifacts and languages from [19], because it introduces the possibility for infinite mutual recursion:

$$\text{elementOf}(a, l) \Rightarrow \text{Artifact}(a) \wedge \text{Language}(l).$$

$$\text{elementOf}(a, l) \Leftarrow \exists s. \text{defines}(s, l) \wedge \text{conformsTo}(a, s)$$

Axiom 4 (conformsTo)

$$\begin{aligned}
\text{conformsTo}(a, d) &\Rightarrow \text{Artifact}(a) \wedge \text{Artifact}(d). \\
\text{conformsTo}(a, d) &\Leftarrow (\forall pa. \text{properPartOf}(pa, a) \wedge \exists pd. \text{properPartOf}(pd, d) \wedge \text{conformsTo}(pa, pd)) \\
&\quad \vee \exists l. \text{defines}(d, l) \wedge \text{elementOf}(a, l).
\end{aligned}$$

Axiom 4 defines conformance recursively:

Case I If both artifacts is atomic in the sense it cannot be decomposed in further parts, we check whether one artifact defines a language the other artifact is an element of.

Case II Else, if at least one artifact is not atomic, then for each proper part of one artifact has to be a proper part in the other artifact it conforms to.

2.2 Traceability

This section provides brief background on the topic of traceability, i.e. its objectives and terminology (see §2.2.1 and §2.2.2 respectively).

2.2.1 Traceability Objectives

Traceability, as a desired property of software systems and their development processes, originates from requirement engineering, i.e. how do we validate or verify that all requirements are met; or more precisely: on what data can we base such validation or verification? [32]

An answer to these questions is motivated by observations on development processes of software systems. Because such a process usually has an iterative and incremental nature, we can reasonably assume engineering activities to leave traces which reflect the performed action.

Modern day traceability objectives are not limited to validation of requirement compliance. It is also regarded as a desirable property for supporting engineering activities, e.g. change impact analysis or dependency analysis [16] [14]. This feeds to the overall topic of program comprehension. We can informally define traceability as "[...] the ability to [...] interrelate uniquely identifiable entities in a way that matters." [12] [32].

2.2.2 Traceability Terminology

This section recapitulates the necessary excerpt of definitions on traceability terms given by [16]. An excessive glossary of traceability terminology can be found in [15].

Definition 1 (Trace)

- **(Noun)** *A specified triplet of elements comprising: a source artifact, a target artifact and a trace link associating the two artifacts. [16]*
- **(Verb)** *The act of following a trace link from a source artifact to a target artifact or vice-versa. [16]*

Definition 2 (Trace Artifact, Source Artifact, Target Artifact) *A traceable unit of data (e.g., a single requirement, a cluster of requirements, a UML class, a UML class operation, a Java class or even a person). A trace artifact is one of the trace elements and is qualified as either a source artifact or as a target artifact when it participates in a trace:*

- **Source Artifact** *The artifact from which a trace originates.*
- **Target Artifact** *The artifact at the destination of a trace.*

[16]

Definition 3 (Trace Artifact Type) *A label that characterizes those trace artifacts that have the same or a similar structure (syntax) and/or purpose (semantics). For example, requirements, design and test cases may be distinct artifact types. [16]*

Definition 4 (Trace Link) *A specified association between a pair of artifacts, one comprising the source artifact and one comprising the target artifact. Trace links are effectively bidirectional. [16]*

Definition 5 (Trace Link Type) *A label that characterizes those trace links that have the same or similar structure (syntax) and/or purpose (semantics). For example, “implements”, “tests”, “refines” and “replaces” may be distinct trace link types. [16]*

Definition 6 (Traceability) *The potential for traces to be established and used. Traceability (i.e., trace “ability”) is thereby an attribute of an artifact or of a collection of artifacts. [16]*

Definition 7 (Traceability Creation) *The general activity of associating two (or more) artifacts, by providing trace links between them. This can be done manually, automatically or semi-automatically, and additional annotations can be provided as desired to characterize attributes of the traces. [16]*

Definition 8 (Trace Recovery) *A particular approach to trace creation that implies the creation of trace links after² the artifacts that they associate have been generated and manipulated. These trace links may be created automatically or semi-automatically using tools. [16]*

2.3 Technological Background

This section provides background on key technologies used by this thesis. §2.3.1 introduces the MegaL modeling language and MegaL/Xtext as corresponding technology. §2.3.2 and §2.3.3 provides an introduction to JAXB and Hibernate respectively. §2.3.4 summarizes the relevant aspects of ANTLR (Another Tool For Language Recognition).

2.3.1 MagaL & MegaL/Xtext

MegaL [24] [10] is a modeling language for linguistic architectures developed by Software Languages Team³ at the University of Koblenz-Landau. MegaL/Xtext [18] [17] is the integration of MegaL into the eclipse IDE (Integrated Development Environment) based on Xtext⁴ developed by the same team.

MegaL provides a textual notation for typed ER Models, i.e. each entity and each relation is strictly typed. This is exemplified by the following listing showing a simple megamodel for citation of scientific papers:

```
Artifact < Entity
File < Artifact
Paper < File
```

²The opposite to trace recovery is **Trace Capture**: "A particular approach to trace creation that implies the creation of trace links concurrently with the creation of the artifacts that they associate. These trace links may be created automatically or semi-automatically using tools." [16]. This is similar to the distinction between static program analysis and dynamic program analysis, thus leaving the scope of this thesis.

³<http://www.softlang.org/> (Retrieved 12th December, 2017)

⁴<http://www.eclipse.org/Xtext/> (Retrieved 12th December, 2017)

```
cites < Paper * Paper

paper1 : Paper
paper2 : Paper
paper1 cites paper2
```

The types, starting with upper-case letters, state that there are files, which are papers. There is a relation `cites`, denoted with a type signature, stating one paper may cite the other. Eventually, there are two concrete paper entity instances `paper1` and `paper2`, where the former cites the latter instance.

MegaL/Xtext provides means for applying additional functionality or semantics to megamodels, e.g. for reasoning or validation on the induced entity graph. For this, MegaL/Xtext implements a plug-in system, which is exemplified in the following listings:

```
MyType < Entity

@Plugin 'MyRelationPlugin'
myRelation < MyType * MyType

MyRelationPlugin : Plugin
MyRelationPlugin = 'classpath:org.softlang.megal.MyRelationPlugin'
```

```
public abstract class GuidedReasonerPlugin extends InjectedReasonerPlugin {
    ...
    protected void guidedDerive(EntityType entityType) throws Throwable {}
    protected void guidedDerive(RelationshipType relationshipType) throws Throwable {}
    protected void guidedDerive(Entity entity) throws Throwable {}
    protected void guidedDerive(Relationship relationship) throws Throwable {}
}
```

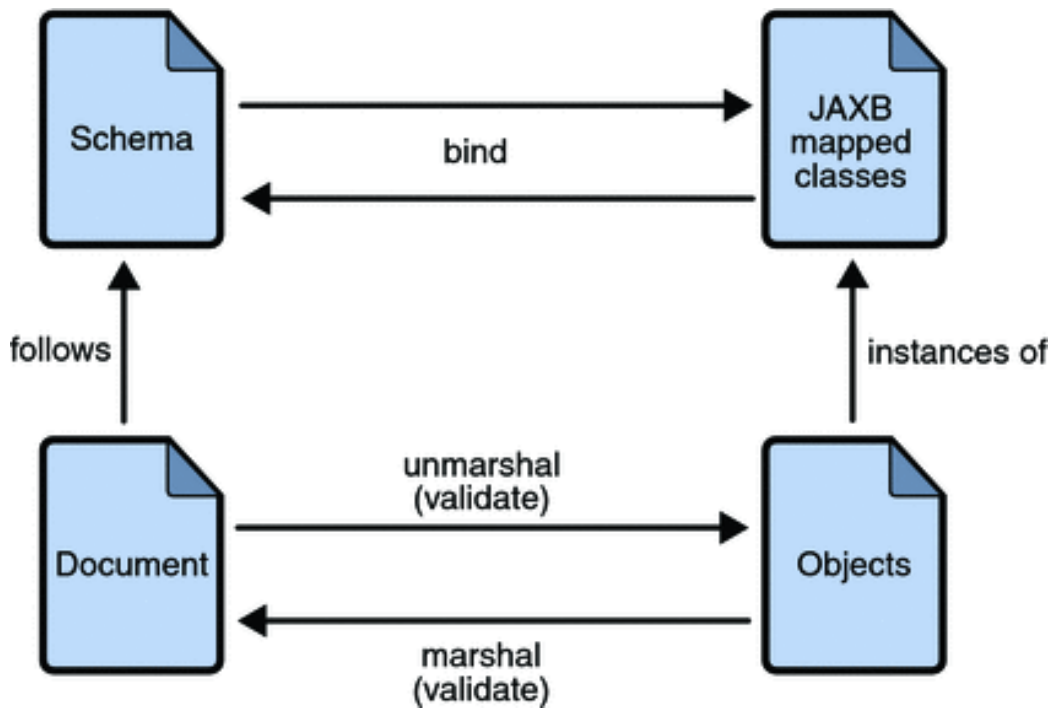
```
public class MyRelationPlugin extends GuidedReasonerPlugin {
    protected void guidedDerive(Relationship relationship) throws Throwable {
        ...
    }
}
```

Given a Java class artifact `MyRelationPlugin.java` implementing special features provided by the MegaL/Xtext API, e.g. the `GuidedReasonerPlugin` class, it can be bound to an entity instance by using its classpath. This plug-in entity can then be used with the `@Plugin` annotation and applied to MegaL elements, i.e. entity types, relationship types, entity instances and relationship instances. `MyRelationPlugin` will be executed during evaluation of the megamodel.

2.3.2 JAXB (Java Architecture for XML Binding)

JAXB [6] [8] is a Java technology developed by the Oracle Corporation for XML serialization and de-serialization of classes and instances. It implements O/X-

Mapping (Object-XML-Mapping), i.e. Java classes can be transformed to XML schema files (XSD files) and instances of classes can be transformed to XML documents. Transformations can also be applied vice-versa. Figure 2.2 shows the O/X-Mapping facilitated by JAXB. The mapping of classes to schemas is called



<https://docs.oracle.com/javase/tutorial/figures/jaxb/jaxb-dataBindingProcess.gif>
(retrieved 11th December, 2017)

Figure 2.2 Steps in the JAXB Binding Process [7]

binding, whereas the mapping of instances to documents is called *marshaling* or *unmarshaling* respectively.

Binding and marshaling is defined through Java annotations. For instance, class and field declarations can be annotated to specify marshaling into specific syntactic XML elements, i.e. element tags or attributes. This is exemplified by the following listing:

```
@XmlElement(name="employee")
@XmlAccessorType(XmlAccessType.FIELD)
public class Employee {
    @XmlAttribute private int id;
    ...
}
```


Annotations can define binding in more detail than shown above, e.g. it is possible to explicitly provide names and types of XML elements. Marshaling itself is triggered using reflection and the `JAXBContext` class:

```
JAXBContext context = JAXBContext.newInstance(Employee.class);
context.createMarshaller().marshal(jaxbElement, writer);
```

More excessive examples and tutorials on JAXB can be found in [6] and in [8].

2.3.3 Hibernate

Hibernate [20] is a Java technology developed by Red Hat Inc. facilitating O/R-Mapping (Object-Relational-Mapping). Using Hibernate, one can define mappings from Java classes and their instances to tables of relational databases and their rows, and vice versa. This allows the underlying framework to generate SQL (Structured Query Language) code for communicating with a relational database system in order to provide means for CRUD operations.

Mappings can be defined using JPA (Java Persistence API) annotations or, traditionally, using XML files, i.e. for each Java class file `MyEntity.java` representing a data entity exists a mapping file `MyEntity.hbm.xml` describing the relational mapping. This is exemplified by the following two listings:

```
public class Employee {
    private int id;
    ...
}
```

```
<class name="Employee" table="EMPLOYEES">
  <id name="id" column="EMPLOYEE_ID">
    ...
  </id>
  ...
</class>
```

Here, the `Employee` class is mapped to the `EMPLOYEE` table its `id` field is mapped to the `EMPLOYEE_ID` column as primary key for the table. More excessive examples and tutorials on Hibernate can be found in [20].

2.3.4 ANTLR (Another Tool For Language Recognition)

ANTLR [29] is a technology providing means for processing structured text using parsing. It is available for Java, C#, JavaScript, Python and many other platforms [27], i.e. ANTLR can generate *LL(*)* parsers for context-free languages in the

target platform's language. Parsers are generated from an EBNF-like grammar notation, e.g. the following "hello world"-grammar taken from [28]:

```
grammar Hello;
hello : 'hello' ID ; // match keyword hello followed by an identifier
ID : [a-z]+ ; // match lower-case identifiers
WS : [ \t\r\n]+ -> skip ; // skip spaces, tabs, newlines
```

ANTLR also provides means for integrating generated parsers into one's project, i.e. means for working with/traversing created Parse Trees. Tree traversal in order to transform Parse Trees to ASTs can be either done traditionally using the Visitor Pattern [13] or using the combination of ANTLR Parse Tree walkers and walker listeners, a variation of the Observer Pattern [13].

Listener interfaces are generated during parser generation. All listeners implement the `ParseTreeListener` interface provided by the ANTLR runtime library:

```
public interface ParseTreeListener {
    void visitTerminal(TerminalNode node);
    void visitErrorNode(ErrorNode node);
    void enterEveryRule(ParserRuleContext ctx);
    void exitEveryRule(ParserRuleContext ctx);
}
```

The following listing shows generated listener interface for the grammar above:

```
public interface HelloListener extends ParseTreeListener {
    void enterHello(HelloContext ctx);
    void exitHello(HelloContext ctx);
}
```

While traversing a Parse Tree, listeners are notified upon entering or exiting tree nodes, called "context" by ANTLR. Traversal is triggered using `ParseTreeWalker` instances, which implement Depth-First Search (DFS). Usage of `ParseTreeWalker` is exemplified in the following listing:

```
...
HelloListener helloListener = ...
ParseTreeWalker parseTreeWalker = new ParseTreeWalker();
parseTreeWalker.walk(helloListener, parseTree);
...
```

An excessive introduction to ANTLR can be found in [29], a simple tutorial can be found in [28].

Chapter 3

Related Work

This chapter provides an non-exhaustive overview over related work regarding trace link recovery approaches. §3.1 summarizes selected papers. §3.2 provides some concluding remarks regarding the summarized approaches and the approach used by this thesis.

3.1 Paper Summary

This section summarizes selected papers regarding their approaches on trace link recovery.

Information Retrieval Methods for Automated Traceability Recovery [9] proposes a semi-automatic approach for recovering trace links:

1. corpus generation using document parsing to extract artifacts of a desired size (paragraphs, classes, methods, etc.)
2. corpus indexing using information retrieval methods in order to create a homogeneous corpus removing stop words and applying other normalization techniques
3. recovery of "candidate" links computing a probabilistic similarity
4. analysis/assessment of candidate links conducted by a human removing false positive links and confirming correct links

It is argued that using probabilistic information retrieval methods reduces the overhead created by development of language specific analysis solutions, e.g. AST generation.

Using Rules for Traceability Creation [33] proposes a fully automatic approach for recovering trace links by normalizing to XML documents and applying XQuery encoded traceability rules on the corpus. This approach also takes ambiguity of natural languages into account by applying a grammatical classification of sentences during preprocessing and considering synonyms during recovery.

Tracing object-oriented code into functional requirements [3] proposes an information retrieval approach for recovering trace links between source code and natural language documentation artifacts. In a provided case study, Java classes are linked to requirement documents based on a probabilistic similarity computed from the set of identifiers found in a class artifact.

Recovering Code to Documentation Links in OO Systems [2] proposes an information retrieval approach for recovering trace links between source code and natural language documentation artifacts based on a probabilistic similarity computed from the set of identifiers found in source code artifacts. This is a predecessor paper of [3].

Mining software repositories for traceability links [21] proposes an approach for recovering trace links among source code and informal documentation artifacts by monitoring their co-change frequency, i.e. artifacts that are changed to concurrently in the history of a version control system. It is argued that a high co-change frequency implies trace links among artifacts.

Combining Textual and Structural Analysis of Software Artifacts for Traceability Link Recovery [26] proposes a combined approach of information retrieval methods for recovering exogenous links between source code and documentation artifacts and classical analysis methods for recovering endogenous links among source code artifacts. Endogenous links among documentation artifacts are proposed to be recovered indirectly, i.e. if source artifacts $S1$ and $S2$ are linked and both are linked to documentation artifacts $D1$ and $D2$ respectively,

then this may be seen as evidence that $D1$ are also linked $D2$. Indirect linking is proposed because the authors argue that information retrieval methods alone are not sufficient for trace link recovery purposes. Documentation artifacts may be linked despite sharing little to no semantic similarities.

Recovering traceability links in software artifact management systems using information retrieval methods [25] provides an excessive case-study evaluating information retrieval methods, namely Latent Semantic Indexing, for trace link recovery. It concludes that information retrieval is not sufficient to recover all links among source code artifacts,

TraceME: Traceability Management in Eclipse [5] introduces TraceME, an eclipse IDE base tool for trace link recovery and management. It utilizes information retrieval techniques, namely the Vector Space Model, for trace link recovery inside eclipse projects.

OpenTrace: An Open Source Workbench for Automatic Software Traceability Link Recovery [1] introduces OpenTrace, a workbench for reproducible experiments with focus on automated trace link recovery. It provides means to apply information retrieval techniques for recovering trace links to large corpora of natural language artifacts.

Traceclipse: An eclipse plug-in for traceability link recovery and management [22] introduces Traceclipse, another eclipse IDE based tool, providing means for trace link recovery using information retrieval techniques, namely the Vector Space Model. The implemented recovery process is semi-automatic in the sense that a user has to inspect proposed links.

3.2 Remarks

The majority of the papers summarized in §3.1 utilize information retrieval techniques for trace link recovery. The approach of this thesis is to use static program analysis techniques for trace link recovery. This notable difference is due to the differing source and target artifacts intended to be recovered. The summarized approaches intend to recover links between source code and informal, natural

language artifacts used for documentation. Because of the latter, applying information retrieval methods seems to be a reasonable strategy. Only [26] argues information retrieval to be insufficient. However, the proposed recovery method (coupling metrics) for inter source code trace links seem to assume a homogeneous corpus. This thesis aims to recover exogenous trace links in a heterogeneous corpus.

Chapter 4

Design

This chapter summarizes the design of the recovery system developed for this thesis. §4.1 summarizes functional requirements of the developed system. §4.2 recapitulates design of the recovery system with respect to the specified requirements. §4.3 summarizes integration in the recovery system with the MegaL/Xtext environment .

4.1 Requirements

This section summarizes functional requirements for the recovery system developed as part of this thesis. All presented requirements are must haves, so no priority differentiation is applied.

Requirement 1 *Fragmen Recovery*. The recovery system has to recover fragments, i.e. syntactically well-formed parts of artifacts (see §2.1.2).

Requirement 2 *Parthood Recovery*. The recovery system has to recover parthood links of fragments (see §2.1.1 and §2.1.2)).

Requirement 3 *Correspondence Recovery*. The recovery system has to recover correspondence links, i.e. links capturing the relation between fragments of artifacts denoting a predefined similarity holds (see §2.1.3).

Requirement 4 *Conformance Recovery*. The recovery system has to recover conformance links, i.e. links capturing the relation between artifacts or fragments denoting that one defines the other (see §2.1.4).

Requirement 5 *MegaL/Xtext Integration.* The recovery system has to run within MegaL/Xtext, i.e. the implementations of requirements 1, 2, 3 and 4 must be compatible and integrated with the MegaL/Xtext plug-in-system.

4.2 Recovery Design

This section summarizes the design of the recovery system developed for this thesis. §4.2.1 will describe the all over process of the recovery system. §4.2.2 will describe the design core API and its components developed for the recovery system. §4.2.3 will describe the design of the actual system for recovering links among O/R/X-Mapping artifacts.

4.2.1 Recovery Process

The recovery process is a straight forward analysis of two artifacts. Figure 4.1 shows a flowchart depicting this. Given two artifacts as input, the recovery pro-

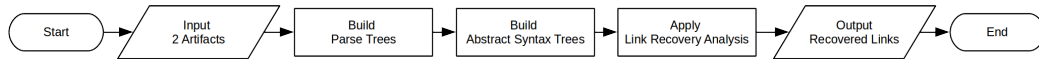


Figure 4.1 The Recovery Process

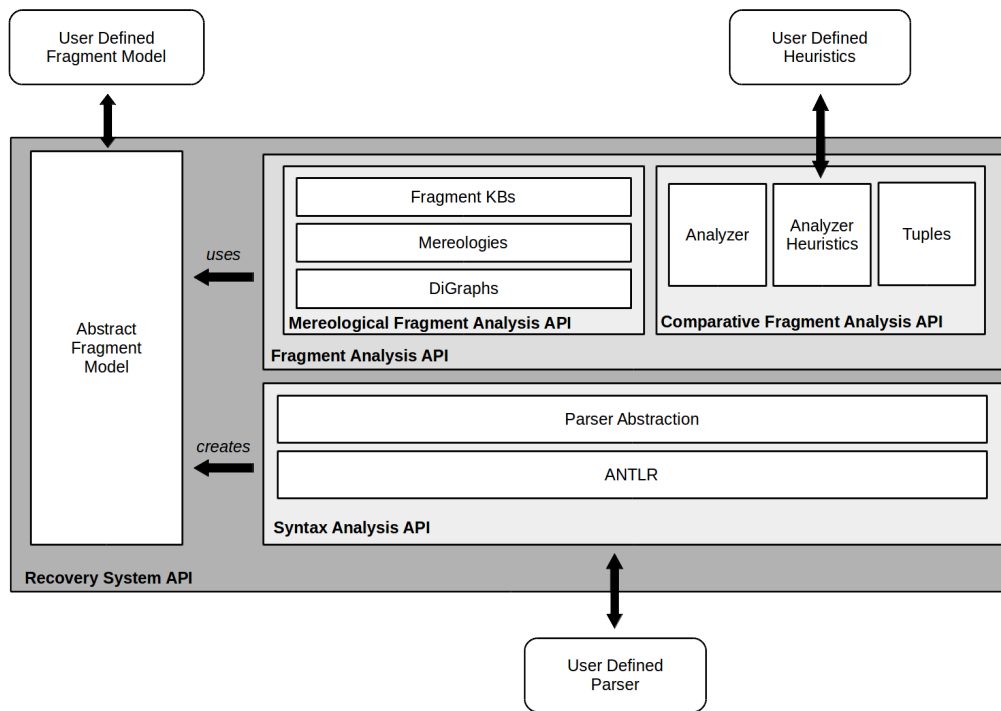
cess works as follows:

1. From each artifact a Concrete Syntax Tree (CST) or Parse Tree is constructed,
2. each Parse Tree is further refined into an AST,
3. both ASTs are compared with each other, i.e. both trees are traversed in a DFS fashion and each pair of nodes is checked whether it can be recovered as link.

Eventually, the set of recovered links serves as output of the process.

4.2.2 Recovery API

The Recovery API is the core of the developed recovery system. It provides generalized data structures and methods for syntactic analysis of artifacts and their



This block diagram depicts the functional outline of the Recovery System API (note, that this does not necessarily correspond to the package outline of its actual implementation).

Figure 4.2 The Recovery API

fragments. Figure 4.2 depicts the API as block diagram. The components of the API are:

Abstract Fragment Model The Abstract Fragment Model serves as base model for all ASTs. It can thought of as Data Transfer Object (DTO) for the analysis components. As user of the API, i have to derive a specific AST from this model. A detailed description follows in §4.2.2.1.

Syntax Analysis API The Syntax Analysis API is the abstraction layer for parsing and AST construction. It is currently backed by ANTLR, but its internal design is loosely coupled, so other parser libraries can be used. As user of the API, i have to implement a parser constructing an AST deriving the Abstract Fragment Model. However, if one uses ANTLR, only AST construction from a CST is required. A detailed description follows in §4.2.2.2.

Fragment Analysis API The Fragment Analysis API consists of two components:

Mereological Fragment Analysis API The Mereological Fragment Analysis API provides components for deriving parthood links between syntactically well-formed fragments. As user of the API, i only have to apply its components to constructed ASTs. A detailed description follows in §4.2.2.3.

Comparative Fragment Analysis API The Comparative Fragment Analysis API provides components for deriving links between different artifacts and their fragments. As user of the API, i have to implement one or more specialized heuristics for deciding which links can be recovered. A detailed description follows in §4.2.2.4.

4.2.2.1 Abstract Fragment Model

The Abstract Fragment Model describes a tree in which each node represents an syntactically well-formed fragment. Figure 4.3 shows an UML (Unified Modeling Language) class diagram of the model. The resulting tree data structure is doubly-linked, i.e. an `IFragment` node aggregates references to its children and to its parent, given it is not the root node. Each fragment `IFragment` contains the text it represents. In order to distinguish fragments which represent the same text, each node also carries the text's position in the artifact through `IFragmentPosition` instances. Positions inside the text are determined by the start and ending line number as well as the corresponding first and last character inside the line. Most of `IFragment`'s relevant code for trees is pre-implemented in the `BaseFragment` abstract class.

4.2.2.2 Syntax Analysis API

The Syntax Analysis API provides abstraction for AST construction. The API itself is really small, it only consists of the `IParser` and `IParserFactory` interfaces. However, its implementation for ANTLR is designed to reduce ANTLR-specific boilerplate code. Figure 4.4 shows the UML class diagram for the ANTLR backed implementation. One can see, that this implementation makes heavy use of the Abstract Factory Pattern [13]. This allows for a quick and easy definition of new parsers as shown in Figure 4.5.

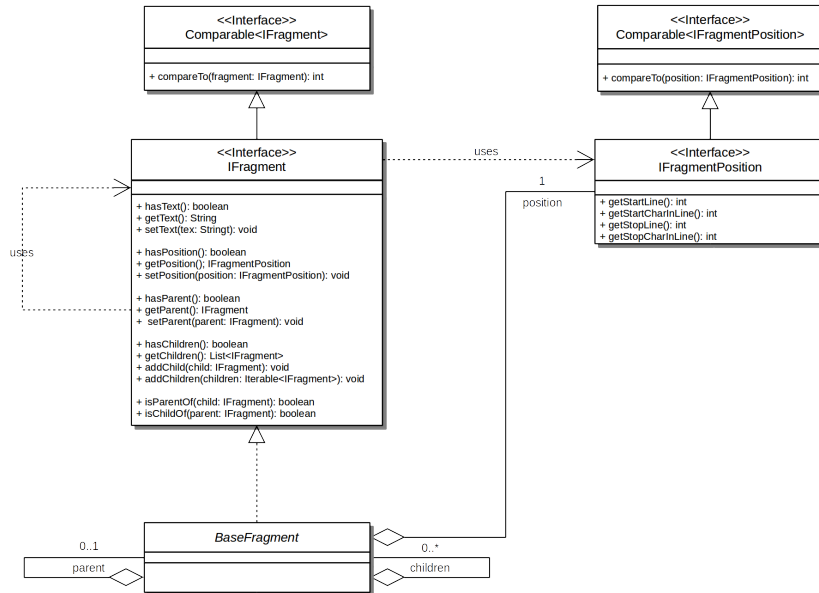


Figure 4.3 The Abstract Fragment Model

The other advantage of the Abstract Factory Pattern is that the instantiation of all relevant classes necessary for the creation of ANTLR Parse Trees takes place in the predefined `AntlrParser` class. This is exemplified by Figure 4.6.

Creation of ASTs or fragment trees is done using the `ParseTreeWalker` and `ParseTreeListener` infrastructure provided by ANTLR as described in §2.3.4. This is an variation of the Observer Pattern. Listeners in conjunction with walkers are an alternative to the Visitor Pattern provided by ANTLR. An instance of `ParseTreeWalker` traverses a Parse Tree using DFS. During traversal, a designated method of `ParseTreeListener` is executed. For AST creation, an API user has to implement the `IFragmentBuildingListener` interface, which is an extension of the `ParseTreeListener` interface. The creation of a fragment tree by `AntlrParser` is exemplified in Figure 4.7.

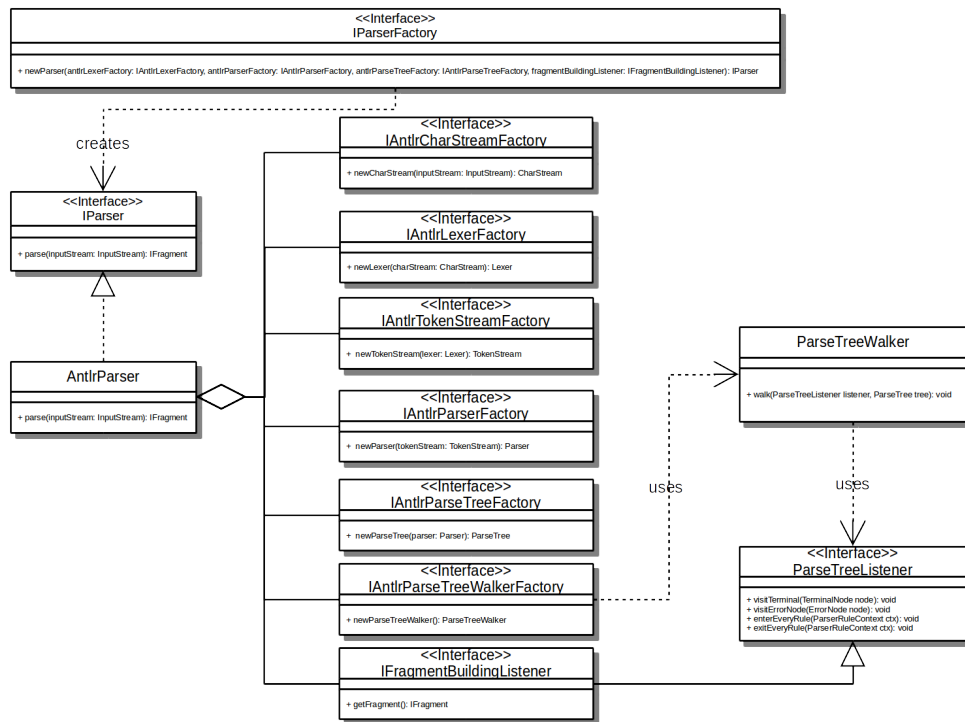


Figure 4.4 The Syntax Analysis API

```

1 ...
2 IParser java8Parser = parserFactory.newParser(
3     java8Lexer::new,
4     java8Parser::new,
5     java8Parser::compilationUnit,
6     java8FragmentBuildingListener::new);
7 ...

```

This example demonstrates the creation of a new parser using an `IParserFactory` instance with Java8 features. Lexer and Parser classes are generated by ANTLR. A suitable listener has to be implemented.

Figure 4.5 `IParserFactory` Usage Example

```

1 ...
2 CharStream charStream = antlrCharStreamFactory.newCharStream(inputStream);
3 Lexer lexer = antlrLexerFactory.newLexer(charStream);
4 TokenStream tokenStream = antlrTokenStreamFactory.newTokenStream(lexer);
5 Parser parser = antlrParserFactory.newParser(tokenStream);
6 ParseTree parseTree = antlrParseTreeFactory.newParseTree(parser);
7 ...

```

This example demonstrates the creation of an ANTLR `ParseTree` instance as implemented by the `AntlrParser` class.

Figure 4.6 `AntlrParser` Parse Tree Creation

```

1 ...
2 ParseTreeWalker parseTreeWalker = antlrParseTreeWalkerFactory.newParseTreeWalker();
3 parseTreeWalker.walk(fragmentBuildingListener, parseTree);
4 IFragment fragment = fragmentBuildingListener.getFragment();
5 ...

```

This example demonstrates the creation of an `IFragment` AST instance as implemented by the `AntlrParser` class.

Figure 4.7 `AntlrParser` Fragment Creation

4.2.2.3 Mereological Fragment Analysis API

The Mereological Fragment Analysis API provides components for deriving parthood links between syntactically well-formed fragments. Direct parthood links are recovered from parent/child relationships within an `IFragment` AST. However, because parthood is transitive, we also need to recover these links. This is done using a digraph data structure upon which we can compute its reflexive transitive closure.

Figure 4.8 shows an UML class diagram depicting the relevant classes and interfaces of the Mereological Fragment Analysis API. At its heart, the API uti-

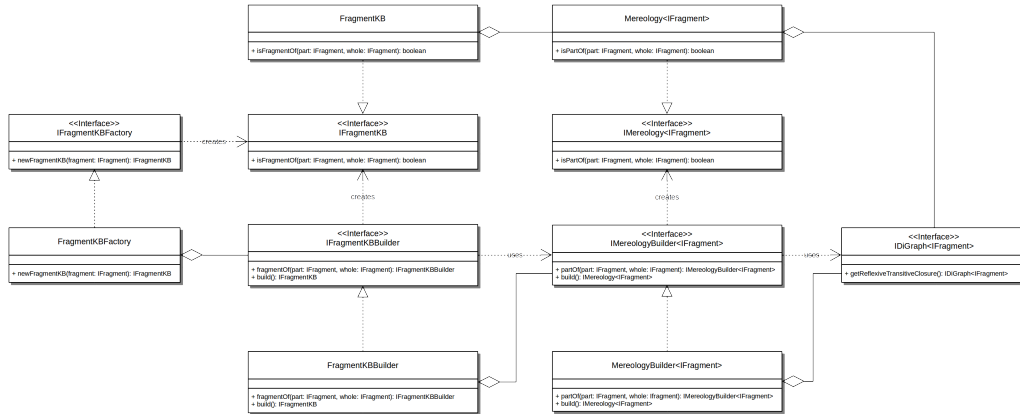


Figure 4.8 Mereological Fragment Analysis API

lizes the generic `IDiGraph` interface whose implementations provide means for interrelating comparable objects, i.e. `IFragment` instances.

This digraph is wrapped by generic `IMereology` implementations, a data structure providing query methods on the topic of mereology, i.e. parthood rela-

tions. Mereologies are constructed using the Builder Pattern [13]. An `IMereologyBuilder` instance adds nodes and edges into its digraph with semantically named methods allowing a descriptive programming style.

`IMereology`'s are then further wrapped by `IFragmentKB` (a Knowledge Base over `IFragment`) implementations in order to avoid dealing with generics throughout the system. This is also done utilizing the Builder Pattern.

The recovery of parthood links is encapsulated through `IFragmentKBFactory` using the Abstract Factory Pattern, which computes `IFragmentKB` instances from an `IFragment` AST as input.

4.2.2.4 Comparative Fragment Analysis API

The Comparative Fragment Analysis API provides components for deriving links from two `IFragment` ASTs generated from different artifacts. Figure 4.9 shows

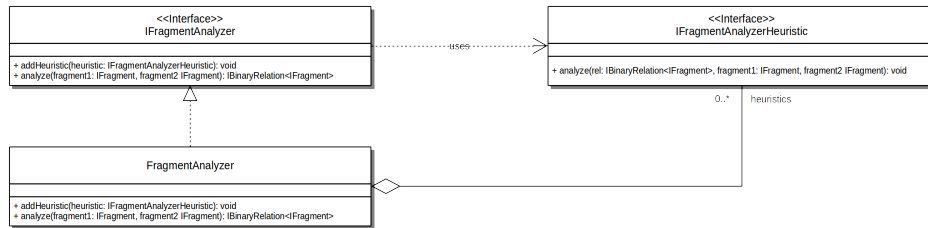


Figure 4.9 Comparative Fragment Analysis API

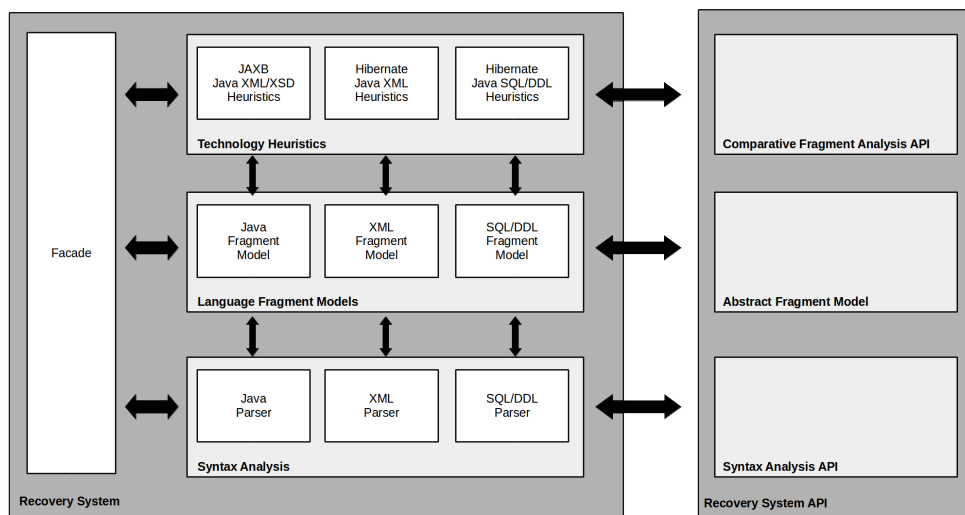
an UML class diagram depicting the relevant interfaces and classes of the API. It utilizes the Strategy Pattern [13] through the `IFragmentAnalyzerHeuristic` interface since concrete behavior of the recovery analysis is for the user to implement. In this context, strategies are called heuristics, because recovery of links is not necessarily based on optimal, i.e. absolutely correct, solutions. Instead strategies may also implement practical solutions with reasonably sufficient results.

`IFragmentAnalyzer` objects allow one to apply multiple heuristics, since there may be more than one practical approach to achieve a goal. Also this allows strategies to keep a relatively small implementation footprint.

Recovered links are captured and stored in an `IBinaryRelation` instance which works like a set of pairs.

4.2.3 Recovery System

The Recovery System is the actual implementation of parthood, correspondence and conformance link recovery for Java based O/R/X-Mapping artifacts. Figure 4.10 shows a block diagram outlining the functional components of the system.



This block diagram depicts the functional outline of the Recovery System and its dependencies to the Recovery System API (note, that this does not necessarily correspond to the package outline of its actual implementation).

Figure 4.10 The Recovery System

The Recovery System utilizes the Syntactic Analysis API described in §4.2.2.2 for Java, XML and SQL/DDDL. For this, the Abstract Fragment Model described in §4.2.2.1 is implemented. Note, fragment models do not implement AST suitable for compilation of the targeted language. Syntactic features unnecessary for the intended analysis, like local variable declarations, arithmetic expressions or invocations, are omitted. Fragment models rather focuses on structural features of the languages at hand.

The Comparative Fragment Analysis API described in §4.2.2.4 is implemented with focus on artifacts of technologies, namely JAXB and Hibernate. It implements heuristics for link recovery between:

- Java and XML for JAXB artifacts, i.e. Java models are serialized as XML
- Java and XSD for JAXB artifacts, i.e. Java models are serialized as XSD
- XML and XSD for JAXB artifacts, i.e. the two previous scenarios occurred
- Java and XML for Hibernate mapping artifacts, i.e. Hibernate uses XML meta-data for O/R-Mapping.
- Java and SQL/DDL for Hibernate generated SQL artifacts, i.e. Hibernate uses Java annotations for O/R-Mapping

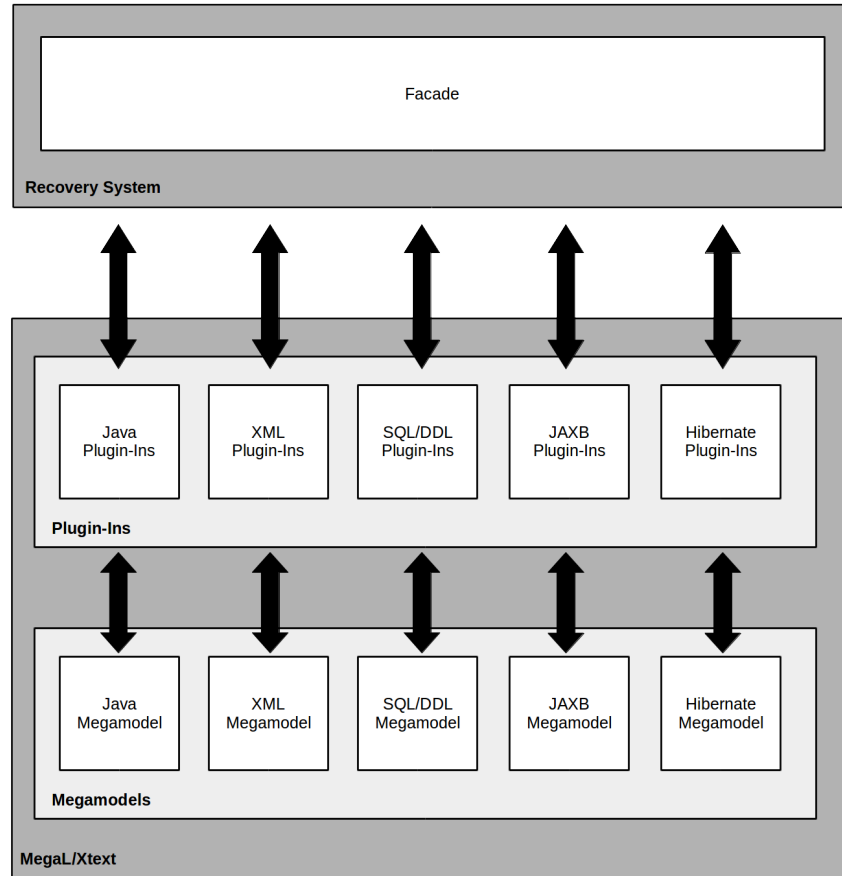
On top of the actual implementation, the Recovery System also utilizes the Facade Pattern [13]. This is to provide a single access point for all implemented analysis features.

4.3 Megal/Xtext Integration Design

This section summarizes the design of the recovery system's integration in the MegaL/Xtext environment. Figure 4.11 shows a block diagram depicting the outline this integration. Note, the Recovery System is implemented as separate project, which is then included as reference (through a `.jar` file) in a MegaL/Xtext project.

MegaL/Xtext provides a plug-in system, which allows to bind special Java classes to `Plugin` entities. Code of such classes is then executed during the evaluation of a megamodel. **ToDo: needs reference** Figure 4.12 exemplifies the declaration of Plug-in withing MegaL.

The Recovery System or rather its facade (see §4.2.3) is used to implement MegaL plug-ins. These Plug-ins are then in turn bound within megamodels, which are divided by language and technology, i.e. there are separate MegaL modules for Java, XML, SQL, Hibernate and JAXB.



This block diagram depicts the functional outline of the Recovery System's integration into the MegaL/Xtext environment (note, that this does not necessarily correspond to the outline of its actual implementation).

Figure 4.11 Integration of the Recovery System into MegaL/Xtext

```

1  ...
2  JavaFragmentRecoveryPlugin : Plugin
3  JavaFragmentRecoveryPlugin realizationOf Java
4  JavaFragmentRecoveryPlugin partOf FileFragmentRecoveryReasonerPlugin
5  JavaFragmentRecoveryPlugin = 'classpath:org.softlang.megal.plugins.impl.java.JavaFragmentRecoveryPlugin'
6  ...

```

This snippet demonstrates the instantiation of `JavaFragmentRecoveryPlugin` as part of `FileFragmentRecoveryReasonerPlugin` in MegaL.

Figure 4.12 MegaL Plug-In Instantiation

Chapter 5

Implementation

This chapter summarizes the implementation of crucial parts of recovery system implemented for this thesis. §5.1 covers the implementation of fragment recovery. §5.2 covers the implementation of parthood link recovery. §5.3 and §5.4 covers the implementation of correspondence and conformance link recovery

5.1 Recovering Fragments

This section summarizes the implementation of fragment recovery. Fragments are syntactically well-formed pieces of code. For instance, consider the following Java class:

```
public class Company {  
    ...  
    private String name;  
    ...  
    public String getName() {  
        return name;  
    }  
    ...  
}
```

It consists, among others, of the following fragments:

- a field declaration fragment for `name`:

```
private String name;
```

- a method declaration fragment for an accessor-method of `name`:

```
public String getName() {  
    return name;  
}
```

- the return statement of the accessor-method:

```
return name;
```

- the class declaration itself can also be considered a fragment.

The task of recovering fragments is to add entities for each of such code pieces into a megamodel. Figure 5.1 exemplifies the recovery of Java¹ fragments in an idealized form.

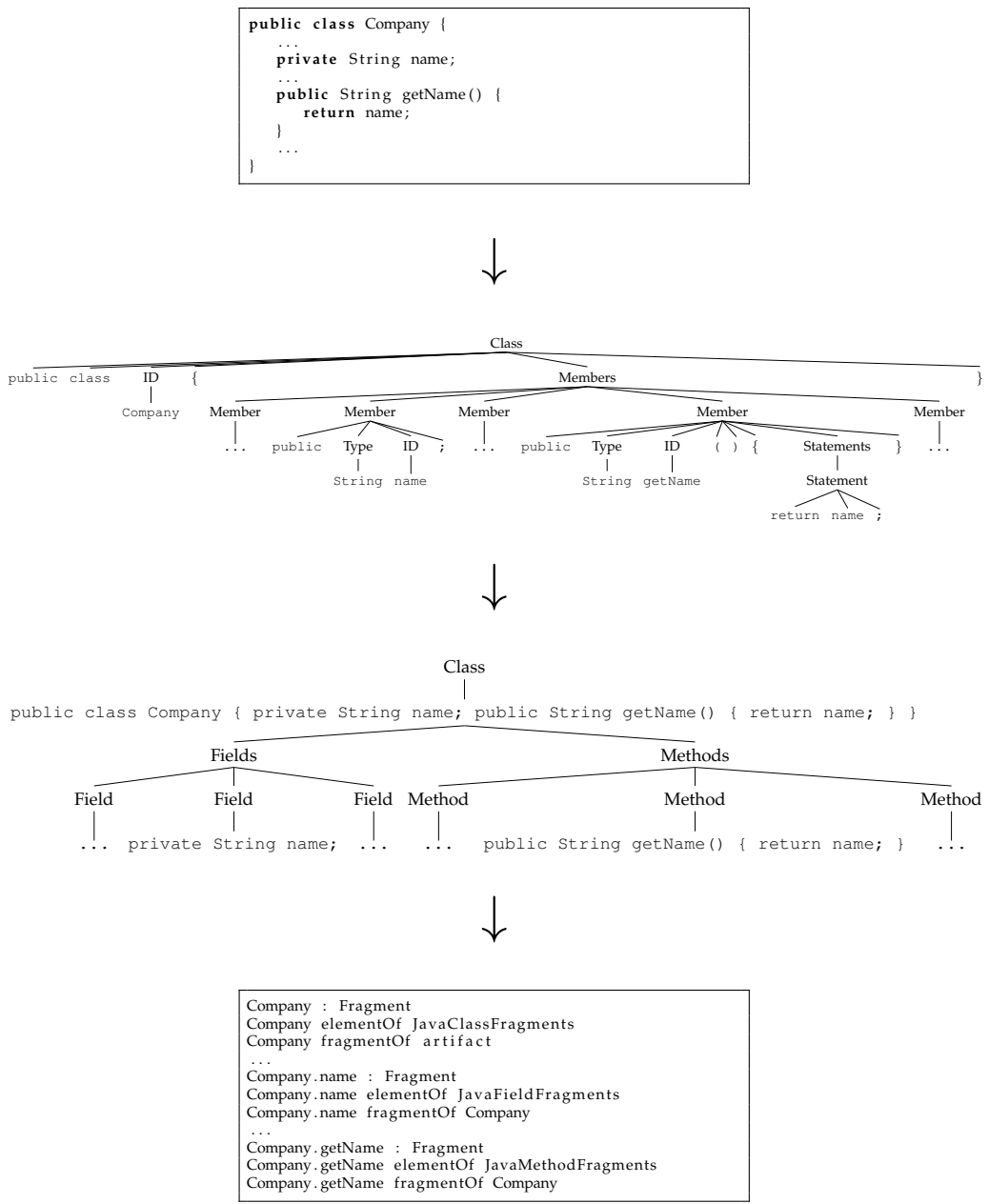
1. a Parse Tree is generated from an artifact; this is done via ANTLR as described in §4.2.2.2 and does not require further explanation, since it is only an application of the most common ANTLR use case
2. then the Parse Tree is transformed into an concrete fragment AST model; this is also done relying on ANTLR tree traversal utilities and described in detail in §5.1.1
3. eventually, nodes of the generated fragment AST are added to a megamodel as entities; the details are described in §5.1.2.

5.1.1 Concrete Fragment Models

As mentioned in §4.2.2.1, Concrete Fragment Models are derivations of the Abstract Fragment Model of the Recovery System API. Figure 5.2 shows an UML class diagram of the implemented Fragment AST for Java. All fragment classes derive from `IFragment` through the base class `JavaFragment`. From here, several specializations are introduced:

- `IdentifiedJavaFragment` for constructs with identifiers
- `ModifiedJavaFragment` for constructs with modifiers like `private`, `public`, `final`, `abstract`, `static`, etc. or annotation meta-data
- `TypedJavaFragment` for constructs with distinct (return-) type like fields, methods or variables

¹Recovery for XML and SQL/DDl is implemented in a similar fashion. If there is a noteworthy difference for other languages it will be explored, otherwise we keep using Java as example domain for the remaining sections of §5.1.



This picture shows an idealized recovery of Java fragments:

code artifact → Parse Tree → fragment AST → megamodel

Both Parse Tree and AST are depicted in a simplified, schematic form.

Figure 5.1 Idealized Recovery of Java Fragments

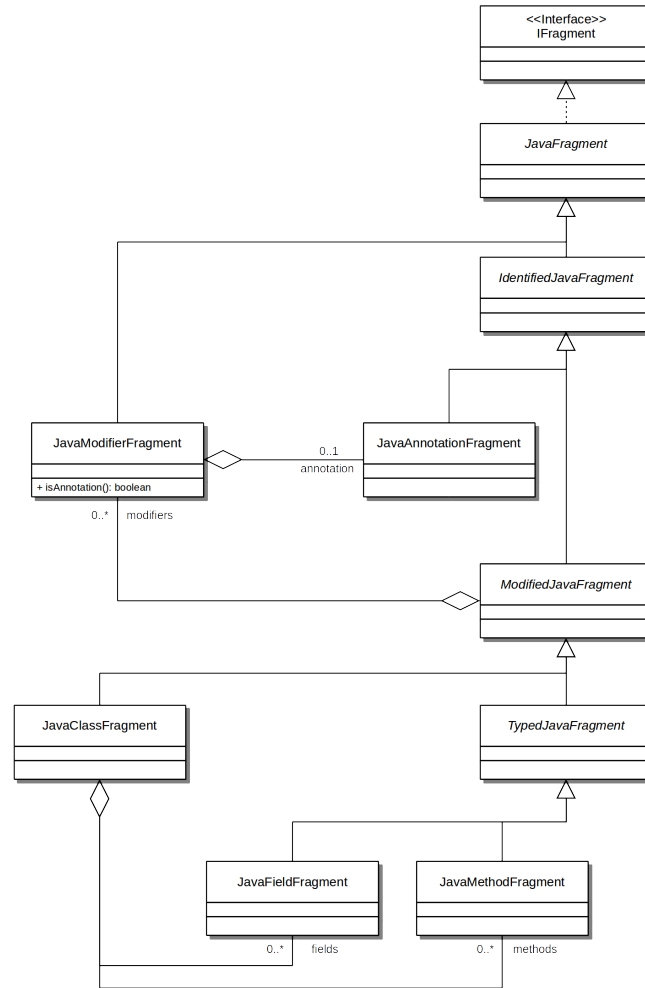


Figure 5.2 Java Fragment AST Model

The fragment model for Java only focuses on structural syntactical features of the language. Everything below method signature level is omitted, because fragments of this level are not important for the scope of this thesis. On the other hand, annotations are captured, since they provide significant information for further recovery analysis. Hibernate and JAXB rely on annotations for O/R/X-Mapping.

The AST is constructed using ANTLR `ParseTreeListener` approach as described in §4.2.2.2. Figure 5.3 shows an excerpt from the listener implementation used to construct the the AST for Java fragments. Here, fragment instances are

```

1  ...
2  @Override
3  public void exitMethodDeclaration (Java8Parser.MethodDeclarationContext ctx) {
4      javaMethodFragments.push(java8FragmentFactory.newJavaMethodFragment(ctx, javaMethodModifierFragments));
5  }
6
7  @Override
8  public void exitNormalClassDeclaration (Java8Parser.NormalClassDeclarationContext ctx) {
9      javaClassFragments.push(java8FragmentFactory.newJavaClassFragment(ctx, javaClassModifierFragments,
10         javaFieldFragments, javaMethodFragments, declaredPackage));
11 }
12 ...

```

This snippet shows an excerpt of the listener used to construct Java fragment ASTs.

Figure 5.3 Construction of Java Fragment ASTs

created when the listener leaves a certain node of the Parse Tree. Creation is delegated to a factory. The listener pushes fragments onto stacks, making the overall implementation work like a pushdown automaton or stack machine.

5.1.2 Megamodeling Fragments

In order to add fragments into a megamodel a suitable linguistic model has to be declared first. Consider the following instance level megamodel of recovered fragments:

```

File < Artifact
Fragment < Artifact
...
elementOf < Set * Set
fragmentOf < Fragment * Artifact
...
aJavaFile : File
aJavaFile = 'Company.java'
...
aJavaFile.Company : Fragment
aJavaFile.Company elementOf JavaClassFragments
aJavaFile.Company fragmentOf aJavaFile
...
aJavaFile.Company.name : Fragment
aJavaFile.Company.name elementOf JavaFieldFragments
aJavaFile.Company.name fragmentOf aJavaFile
aJavaFile.Company.name fragmentOf Company
...
aJavaFile.Company.getName : Fragment
aJavaFile.Company.getName elementOf JavaMethodFragments
aJavaFile.Company.getName fragmentOf aJavaFile
aJavaFile.Company.getName fragmentOf Company

```

This megamodel shows fragments of a Java file `Company.java`. It contains entities denoting a class declaration, a field declaration and a method declaration

fragment. All are declared instances of the entity type `Fragment`. Entity names use dot-notation for implying parthood, i.e.:

$$a.b \Rightarrow b\text{partOf}a$$

Further details of naming-scheme employed for fragment entities are explored in §5.1.2.2.

We use `elementOf` for further specialization and/or generalization of the entitie’s kind². However, the right-hand side of these relationships used to specialize fragments are not simply sets, they are actually modeled as languages as we will see in §5.1.2.1.

5.1.2.1 Fragment Languages

As mentioned in §??, a fragment alone cannot be element of the language the artifact it originated from belongs to, e.g. a Java method alone cannot be accepted by the Java grammar, nor can XML attribute alone be accepted by the XML grammar. Therefore, in order to cleanly add fragments into a megamodel, we need to model fragment languages. Figure 5.4 shows the megamodel for Java fragments. It introduces the special relationship type `fragmentLanguageOf`, denoting the

```

1  ...
2  Language < Set
3  ...
4  fragmentLanguageOf < Language * Language
5  ...
6  Java : Language
7
8  JavaFragments : Language
9  JavaFragments fragmentLanguageOf Java
10
11 JavaClassFragments : Language
12 JavaClassFragments subsetOf JavaFragments
13
14 JavaFieldFragments : Language
15 JavaFieldFragments subsetOf JavaFragments
16
17 JavaMethodFragments : Language
18 JavaMethodFragments subsetOf JavaFragments
19 ...

```

Figure 5.4 A Megamodel for Java Fragments

left-hand side is the language containing all possible fragments the right-hand

²We use the term *kind* to avoid confusion with entity types, although we use it to refer to fragment types.

side's language elements can be deconstructed in. This also implies a super-set relation between the two languages as `fragmentOf` is reflexive like `partOf`, so:

$$A \text{ fragmentLanguageOf } B \Rightarrow B \subseteq A$$

5.1.2.2 Fragment Entity Identifiers

5.2 Recovering Parthood Links

This section summarizes the implementation of parthood link recovery.

5.3 Recovering Correspondence Links

This section the implementation of correspondence and conformance link recovery

5.4 Recovering Conformance Links

This section the implementation of correspondence and conformance link recovery

Chapter 6

Mini Case-Study

This chapter summarizes the methodology used to develop the recovery system for this thesis. In short, the recovery system is developed example-driven using the model of a fictional Human Resource Management System (HRMS) used by the 101wiki¹ for its contributions.

6.1 Example Corpus

The example corpus used to develop the recovery system for this thesis consists of artifacts implementing a fictional HRMS within an O/R/X-Mapping scenario using Java technologies. The model is implemented using plain Java. It is then mapped to plain XML/XSD with JAXB and to SQL/DDDL statements using Hibernate mapping files and/or annotations.

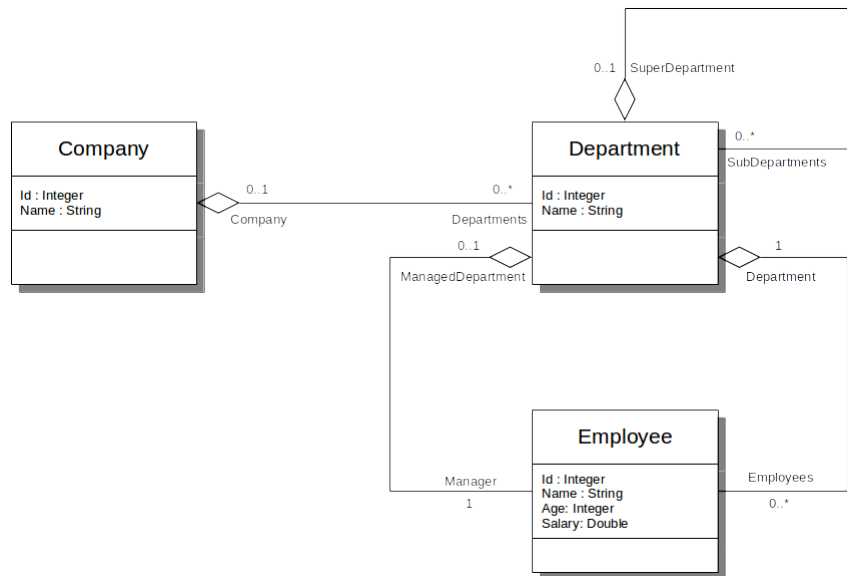
6.1.1 The 101HRMS Model

101wiki Human Resource Management System (101HRMS)² provides a simple model of a company with many departments and employees. Figure 6.1 shows an UML class diagram of a variant of this model.

The 101HRMS model consists of companies attributed with a name. Each company accumulates departments. Each department is also attributed with a

¹<https://101wiki.softlang.org/> (retrieved 12th November, 2017)

²<https://101wiki.softlang.org/101:@system> (retrieved 12th November, 2017)



This UML class diagram depicts the model of the 101HRMS. It consists of simple companies with nested departments and employees mapped to the latter.

Figure 6.1 The 101 Human Resource Management System Model

name, aggregates employees and has one employee acting as manager. Departments can further be refined into sub-departments. Each employee is attributed with a name, an age and a salary. Each entity is also attributed with an ID.

6.1.2 Linguistic Domains of the Example Corpus

The example corpus used to develop the recovery system contains artifacts implementing the 101HRMS model generated or used by Java technologies for O/R/X-Mapping, i.e. a Java model is mapped to plain XML/XSD with JAXB, to a Hibernate mapping file and to SQL/Data Definition Language (DDL) statements. Figure 6.2 shows a schematic illustration of the linguistic domains involved:

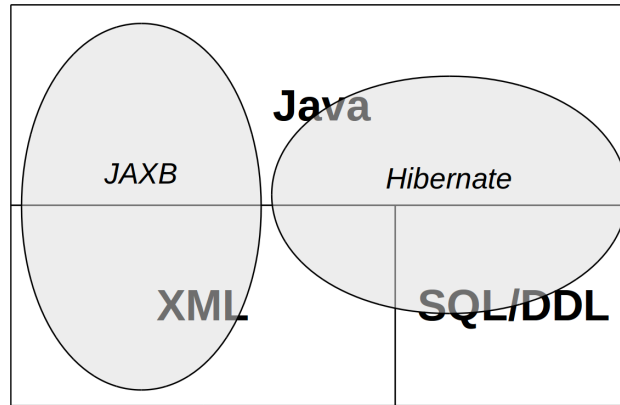
Java The language and technology used to implement the 101HRMS model.

XML The language used to serialize the 101HRMS model.

SQL/DDDL The language used to persist the 101HRMS model.

JAXB The technology used to implement O/X-Mapping of the 101HRMS model.

Hibernate The technology used to implement O/R-Mapping of the 101HRMS model.



This schematic illustration depicts the interrelation among linguistic domains of example corpus used. It depicts languages and technologies for O/R/X-Mapping with Java.

Figure 6.2 Example Corpus Domains: Java O/R/X

The languages (Java, XML & SQL) in Figure 6.2 are displayed as disjoint square sets. Technologies (JAXB & Hibernate) are displayed as oval sets intersecting languages. This is due to their linguistic nature, e.g. JAXB produces specific Java- and XML-Code which does not necessarily intersect with code produced by other technologies. Hibernate intersects all three languages. It uses XML files or Java-Annotations for describing O/R-Mapping of a data-model and generates SQL artifacts according to that mapping. In this sense, technologies create technology-specific subsets of a languages.

Chapter 7

Conclusion

7.1 Future Work

Glossary

101HRMS 101wiki¹ Human Resource Management System². The model used by the 101wiki for its contributions. 41–43, see also HRMS

Abstract Factory Pattern A creational GoF (Gang of Four) pattern used in software design to decouple instantiation from usage of objects. Hides the concrete nature of created instances. 26, 27, 30

ANTLR Another Tool For Language Recognition. 14, 17, 18, 25–28, 35, 37

API Application Programming Interface. 5, 15, 24–27, 29–32, 35, 46

artifact . 3–5, 10–12, 15, 19–23, 36

AST Abstract Syntax Tree: A tree data structure representing the abstract syntax of a parsed text. This tree omits syntactic features like parentheses for grouping or semicolons for sequencing. 3, 4, 18, 20, 24–27, 29–31, 35–38

Builder Pattern A creational GoF pattern used in software design to prevent constructor parameters from piling up. 30

conformance The relation between . 3, 5, 7, 9, 11, 12, 31, 34, 40

correspondence The relation between . 3, 5, 7, 9, 10, 23, 31, 34, 40

CST Concrete Syntax Tree: A tree data structure representing the concrete syntax of a parsed text.. 24, 25

¹<https://101wiki.softlang.org/> (retrieved 12th November, 2017)

²<https://101wiki.softlang.org/101:@system> (retrieved 12th November, 2017)

- DDL** Data Definition Language. Language or subset of a language used to describe structure and content of data. 42
- DFS** The algorithmic concept of traversing a tree or graph data structure 'top-down' until reaching the end of a path before backtracking and traversing another path. 18, 24, 27
- DTO** Data Transfer Object. Objects with no relevant (business) logic of their own. Their sole purpose is to carry data between layers of a software system. 25
- ER Model** Entity-Relationship Model. 3, 14
- Facade Pattern** A structural GoF pattern used in software design to simplify the usage of complex systems or APIs. It provides single access point for such system. Such access points are called facades. 32
- fragment** A syntactically well-formed piece of a possibly larger text. iv, 3–5, 9, 10, 23, 34–39
- GoF** Gang of Four. A group of authors (Erich Gamma, Richard Helm, Ralph Johnson and John Vlissides) publishing on the subject of object-oriented software design. The term may also refer to design patterns described in their book *Design Patterns: Elements of Reusable Object-Oriented Software* [13]. 45–48
- Hibernate** The Hibernate ORM Framework. 5, 14, 17, 32, 37, 41–43
- HRMS** Human Resource Management System. 41
- HTTP** Hypertext Transfer Protocol. 1
- IDE** Integrated Development Environment. 14, 21
- Java** The Java Programming Language and Platform. iv, 2, 5, 7, 9, 15–17, 20, 28, 31, 32, 34–39, 41–43
- JAXB** Java Architecture for XML Binding. 2, 5, 14–17, 32, 37, 41–43
- JPA** Java Persistence API. 17

JSON JavaScript Object Notation. 10

language . 3, 10, 11

linguistic architecture . 3, 14

manifestation . 2, 10, see representation

MegaL The megamodeling language developed by the Softlang Team at the University of Koblenz-Landau for descriptively and prescriptively modeling linguistic architectures of software systems. 14, 15, 32, 33, 47

MegaL/Xtext The Xtext implementation and eclipse IDE integration of MegaL. 14, 15, 23, 24, 32, 33

megamodel . iv, 3, 14, 15, 32, 35, 36, 38, 39

mereology . 7, 8, 29, 30

metamodel . 3

O/R-Mapping Object-Relational-Mapping. 17, 32, 43

O/R/X-Mapping Object-Relational- and XML-Mapping. 1, 2, 5, 24, 31, 37, 41–43

O/X-Mapping Object-XML-Mapping. 15, 16, 42

Observer Pattern A behavioral GoF pattern used in software design to propagate state changes from one object to many dependent objects. 18, 27

ontology . 3

ORM . 46, see O/R-Mapping

Parse Tree . iv, 4, 18, 24, 27, 28, 35, 36, 38, see CST

parthood The relation between an entity and its constituent parts. 6–11, 23, 29–31, 34, 39, 40

representation . 2, 10

SQL Structured Query Language. 17, 32, 42, 43

SQL/DDL The DDL subset of SQL. 31, 32, 35, 41

static program analysis . 3–5, 14, 21

Strategy Pattern A behavioral GoF pattern used in software design to separate behavior from structure. It allows to encapsulate and reuse behavior as part of the configuration of larger constructs. 30

technology . 3, 14

trace . 12

trace link . 3–5

traceability . 3, 12, 13

traceability recovery . 3

UML Unified Modeling Language. 26, 29, 30, 35, 41, 42

Visitor Pattern A behavioral GoF pattern used in software design to separate behavior from structure. Visitors facilitate the extension of behavior without modifying structure. The Visitor Pattern can be used to traverse object graphs. 18, 27

XML Extensible Markup Language. 1, 2, 10, 15–17, 20, 31, 32, 35, 39, 41–43

XSD XML Schema Definition. 2, 16, 32, 41, 42

Bibliography

- [1] E. Angius and R. Witte. “OpenTrace: An Open Source Workbench for Automatic Software Traceability Link Recovery”. In: *2012 19th Working Conference on Reverse Engineering*. Oct. 2012, pp. 507–508. DOI: [10.1109/WCRE.2012.63](https://doi.org/10.1109/WCRE.2012.63).
- [2] G. Antoniol et al. “Recovering code to documentation links in OO systems”. In: *Sixth Working Conference on Reverse Engineering (Cat. No.PR00303)*. Oct. 1999, pp. 136–144. DOI: [10.1109/WCRE.1999.806954](https://doi.org/10.1109/WCRE.1999.806954).
- [3] G. Antoniol et al. “Tracing object-oriented code into functional requirements”. In: *Proceedings IWPC 2000. 8th International Workshop on Program Comprehension*. 2000, pp. 79–86. DOI: [10.1109/WPC.2000.852482](https://doi.org/10.1109/WPC.2000.852482).
- [4] Anya Helene Bagge and Vadim Zaytsev. “Languages, Models and Megamodels”. In: *Post-proceedings of the Seventh Seminar on Advanced Techniques and Tools for Software Evolution, SATToSE 2014, L’Aquila, Italy, 9-11 July 2014*. 2014, pp. 132–143. URL: <http://ceur-ws.org/Vol-1354/paper-12.pdf>.
- [5] G. Bavota et al. “TraceME: Traceability Management in Eclipse”. In: *2012 28th IEEE International Conference on Software Maintenance (ICSM)*. Sept. 2012, pp. 642–645. DOI: [10.1109/ICSM.2012.6405343](https://doi.org/10.1109/ICSM.2012.6405343).
- [6] Oracle Corporation. *Java Architecture for XML Binding*. <http://www.oracle.com/technetwork/articles/javase/index-140168.html>. Retrieved 11th December, 2017.

- [7] Oracle Corporation. *Steps in the JAXB Binding Process*. Retrieved 11th December, 2017. URL: <https://docs.oracle.com/javase/tutorial/figures/jaxb/jaxb-dataBindingProcess.gif>.
- [8] Oracle Corporation. *The Java Tutorials: Java Architecture for XML Binding*. <https://docs.oracle.com/javase/tutorial/jaxb/TOC.html>. Retrieved 11th December, 2017.
- [9] Andrea De Lucia et al. "Information Retrieval Methods for Automated Traceability Recovery". In: *Software and Systems Traceability*. Ed. by Jane Cleland-Huang, Orlena Gotel, and Andrea Zisman. London: Springer London, 2012, pp. 71–98. ISBN: 978-1-4471-2239-5. DOI: [10.1007/978-1-4471-2239-5_4](https://doi.org/10.1007/978-1-4471-2239-5_4). URL: https://doi.org/10.1007/978-1-4471-2239-5_4.
- [10] Jean-Marie Favre, Ralf Lämmel, and Andrei Varanovich. "Modeling the Linguistic Architecture of Software Products". In: *Model Driven Engineering Languages and Systems - 15th International Conference, MODELS 2012, Innsbruck, Austria, September 30-October 5, 2012. Proceedings*. 2012, pp. 151–167. DOI: [10.1007/978-3-642-33666-9_11](https://doi.org/10.1007/978-3-642-33666-9_11). URL: http://dx.doi.org/10.1007/978-3-642-33666-9_11.
- [11] Jean-Marie Favre and Tam Nguyen. "Towards a Megamodel to Model Software Evolution Through Transformations". In: *Electr. Notes Theor. Comput. Sci.* 127.3 (2005), pp. 59–74. DOI: [10.1016/j.entcs.2004.08.034](https://doi.org/10.1016/j.entcs.2004.08.034). URL: <http://dx.doi.org/10.1016/j.entcs.2004.08.034>.
- [12] Richard Freeman Paige et al. "Building Model-Driven Engineering Traceability Classifications". In: (Jan. 2010).
- [13] Erich Gamma et al. *Design Patterns: Elements of Reusable Object-oriented Software*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 1995. ISBN: 0-201-63361-2.
- [14] O. C. Z. Gotel and C. W. Finkelstein. *An analysis of the requirements traceability problem*. Apr. 1994. DOI: [10.1109/ICRE.1994.292398](https://doi.org/10.1109/ICRE.1994.292398).

- [15] Orlena Gotel et al. "Appendix A: Glossary of Traceability Terms (v1.0)". In: *Software and Systems Traceability*. 2012, pp. 413–424.
- [16] Orlena Gotel et al. "Traceability Fundamentals". In: *Software and Systems Traceability*. 2012, pp. 3–22. DOI: [10 . 1007 / 978 - 1 - 4471 - 2239 - 5 _ 1](https://doi.org/10.1007/978-1-4471-2239-5_1). URL: https://doi.org/10.1007/978-1-4471-2239-5_1.
- [17] Johannes Härtel et al. "Interconnected Linguistic Architecture". In: *The Art, Science, and Engineering of Programming Journal* 1 (1 2017). 27 pages.
- [18] Lukas Härtel. "Linguistic architecture on the workbench". Bachelor Thesis. University of Koblenz-Landau, Oct. 2015.
- [19] Marcel Heinz, Ralf Lämmel, and Andrei Varanovich. *Axioms of linguistic architecture*. 9 pages. 2017.
- [20] Red Hat Inc. *Hibernate ORM*. <http://hibernate.org/orm/1>. Retrieved 11th December, 2017.
- [21] H. Kagdi, J. I. Maletic, and B. Sharif. "Mining software repositories for traceability links". In: *15th IEEE International Conference on Program Comprehension (ICPC '07)*. June 2007, pp. 145–154. DOI: [10 . 1109/ICPC.2007.28](https://doi.org/10.1109/ICPC.2007.28).
- [22] Samuel Klock et al. "Traceclipse: An eclipse plug-in for traceability link recovery and management". In: (Jan. 2011).
- [23] Ralf Lämmel. "Coupled software transformations revisited". In: *Proceedings of the 2016 ACM SIGPLAN International Conference on Software Language Engineering, Amsterdam, The Netherlands, October 31 - November 1, 2016*. 2016, pp. 239–252. URL: <http://dl.acm.org/citation.cfm?id=2997366>.
- [24] Ralf Lämmel and Andrei Varanovich. "Interpretation of Linguistic Architecture". In: *Modelling Foundations and Applications - 10th European Conference, ECMFA 2014, Held as Part of STAF 2014, York, UK, July 21-25, 2014. Proceedings*. 2014, pp. 67–82. DOI: [10 . 1007 / 978 -](https://doi.org/10.1007/978-)

- 3-319-09195-2_5. URL: http://dx.doi.org/10.1007/978-3-319-09195-2_5.
- [25] Andrea De Lucia et al. "Recovering Traceability Links in Software Artifact Management Systems Using Information Retrieval Methods". In: *ACM Trans. Softw. Eng. Methodol.* 16.4 (Sept. 2007). ISSN: 1049-331X. DOI: [10.1145/1276933.1276934](https://doi.org/10.1145/1276933.1276934). URL: <http://doi.acm.org/10.1145/1276933.1276934>.
- [26] C. McMillan, D. Poshyvanyk, and M. Revelle. "Combining textual and structural analysis of software artifacts for traceability link recovery". In: *2009 ICSE Workshop on Traceability in Emerging Forms of Software Engineering*. May 2009, pp. 41–48. DOI: [10.1109/TEFSE.2009.5069582](https://doi.org/10.1109/TEFSE.2009.5069582).
- [27] Terence Parr. ANTLR. <http://www.antlr.org/>. Retrieved 12th December, 2017.
- [28] Terence Parr. *Getting Started with ANTLR v4*. <https://github.com/antlr/antlr4/blob/master/doc/getting-started.md>. Retrieved 12th December, 2017.
- [29] Terence Parr. *The Definitive ANTLR 4 Reference*. 2nd. Pragmatic Bookshelf, 2013. ISBN: 1934356999, 9781934356999.
- [30] Achille C. Varzi. "Mereology". In: *The Stanford encyclopedia of philosophy*. Ed. by Edward N. Zalta. Winter 2016 ed. retrieved 27. July 2017. 2016. URL: <https://plato.stanford.edu/archives/win2016/entries/mereology/>.
- [31] Achille C. Varzi. "Parts, Wholes, and Part-Whole Relations: The Prospects of Mereotopology". In: *Data Knowl. Eng.* 20.3 (1996), pp. 259–286. DOI: [10.1016/S0169-023X\(96\)00017-1](https://doi.org/10.1016/S0169-023X(96)00017-1). URL: [http://dx.doi.org/10.1016/S0169-023X\(96\)00017-1](http://dx.doi.org/10.1016/S0169-023X(96)00017-1).
- [32] Stefan Winkler and Jens Pilgrim. "A Survey of Traceability in Requirements Engineering and Model-driven Development". In: *Softw. Syst. Model.* 9.4 (Sept. 2010), pp. 529–565. ISSN: 1619-1366. DOI: [10.](https://doi.org/10.1007/s00378-010-0529-5)

1007/s10270-009-0145-0. URL: <http://dx.doi.org/10.1007/s10270-009-0145-0>.

- [33] Andrea Zisman. "Using Rules for Traceability Creation". In: *Software and Systems Traceability*. Ed. by Jane Cleland-Huang, Orlena Gotel, and Andrea Zisman. London: Springer London, 2012, pp. 147–170. ISBN: 978-1-4471-2239-5. DOI: [10.1007/978-1-4471-2239-5_7](https://doi.org/10.1007/978-1-4471-2239-5_7). URL: https://doi.org/10.1007/978-1-4471-2239-5_7.