



Ralf Lämmel

# Software Languages

Syntax, Semantics,  
and Metaprogramming

 Springer

Ralf Lämmel

# Software Languages

Syntax, Semantics, and Metaprogramming

 Springer

Ralf Lämmel  
Computer Science Department  
Universität Koblenz-Landau  
Koblenz, Germany

ISBN 978-3-319-90798-7      ISBN 978-3-319-90800-7 (eBook)  
<https://doi.org/10.1007/978-3-319-90800-7>

Library of Congress Control Number: 2018942228

© Springer International Publishing AG, part of Springer Nature 2018

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Springer imprint is published by the registered company Springer International Publishing AG part of Springer Nature.

The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

# Copyright and Attribution

## *Copyright for Code*

The code in this book is part of the open-source YAS project:

<http://www.softlang.org/yas>.

YAS is licensed under the [MIT license](#).

Copyright 2016–2018 Ralf Lämmel

Permission is hereby granted, free of charge, to any person obtaining a copy of the YAS code including software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

## *Artwork Credits*

*Cover artwork:*

Wojciech Kwasnik, *The Tower of Software Languages*, 2017.

With the assistance of Archina Void and Daniel Dünker.

Licensed under [CC BY-SA 4.0](#).

Artwork *DMT, acrylic, 2006* by [Matt Sheehy](#) is quoted with the artist’s permission.

*Credits for per-chapter artwork:*

Wojciech Kwasnik.

See the individual chapters for details.

# Contents

- List of Recipes ..... xxvii
  
- Acronyms ..... xxix
  - Fabricated languages ..... xxix
  - Other acronyms ..... xxx
  
- 1 The Notion of a Software Language** ..... 1
  - 1.1 Examples of Software Languages ..... 2
    - 1.1.1 Real-World Software Languages ..... 2
    - 1.1.2 Fabricated Software Languages ..... 3
      - 1.1.2.1 BNL: A Language of Binary Numbers ..... 5
      - 1.1.2.2 BTL: An Expression Language ..... 5
      - 1.1.2.3 BL: A Language for Buddy Relationships ..... 6
      - 1.1.2.4 BFPL: A Functional Programming Language ..... 6
      - 1.1.2.5 BIPL: An Imperative Programming Language ..... 7
      - 1.1.2.6 FSML: A Language for Finite State Machines ..... 7
      - 1.1.2.7 BGL: A Language for Context-Free Grammars ... 9
  - 1.2 Classification of Software Languages ..... 9
    - 1.2.1 Classification by Paradigm ..... 10
    - 1.2.2 Classification by Type System ..... 11
    - 1.2.3 Classification by Purpose ..... 12
    - 1.2.4 Classification by Generality or Specificity ..... 13
    - 1.2.5 Classification by Representation ..... 14
    - 1.2.6 Classification by Notation ..... 15
    - 1.2.7 Classification by Degree of Declarativeness ..... 15
  - 1.3 The Lifecycle of Software Languages ..... 17
    - 1.3.1 Language Definition ..... 18
    - 1.3.2 Language Implementation ..... 20
      - 1.3.2.1 Compilation versus Interpretation ..... 20
      - 1.3.2.2 Architecture of a Compiler ..... 21
      - 1.3.2.3 Classification of Language Processors ..... 22

1.3.2.4	Metaprogramming Systems	25
1.3.2.5	Language Workbenches	26
1.3.3	Language Evolution	27
1.4	Software Languages in Software Engineering	28
1.4.1	Software Re-Engineering	28
1.4.2	Software Reverse Engineering	30
1.4.3	Software Analysis	31
1.4.4	Technological Spaces	33
1.4.5	Model-Driven Engineering	36
	Summary and outline	38
	References	38
<b>2</b>	<b>A Story of a Domain-Specific Language</b>	<b>51</b>
2.1	Language Concepts	52
2.2	Internal DSL	54
2.2.1	Baseline Object Model	54
2.2.2	Fluent API	57
2.2.3	Interpretation	61
2.2.4	Well-Formedness	63
2.3	External DSL	66
2.3.1	Syntax Definition	67
2.3.2	Syntax Checking	68
2.3.3	Parsing	71
2.4	DSL Services	74
2.4.1	Interchange Format	74
2.4.2	Code Generation	76
2.4.3	Visualization	82
	Summary and outline	84
	References	85
<b>3</b>	<b>Foundations of Tree- and Graph-Based Abstract Syntax</b>	<b>87</b>
3.1	Tree-Based Abstract Syntax	88
3.1.1	Trees versus Terms	88
3.1.2	A Basic Signature Notation	89
3.1.3	Abstract Syntax Trees	90
3.1.4	An Extended Signature Notation	91
3.1.5	Illustrative Examples of Signatures	92
3.1.5.1	Syntax of Simple Expressions	92
3.1.5.2	Syntax of Simple Imperative Programs	92
3.1.5.3	Syntax of Simple Functional Programs	93
3.1.5.4	Syntax of Finite State Machines	94
3.1.6	Languages as Sets of Terms	95
3.1.7	Conformance to a Signature	96
3.2	Graph-Based Abstract Syntax	96
3.2.1	Trees versus Graphs	97

- 3.2.2 Languages as Sets of Graphs . . . . . 98
- 3.2.3 A Metamodeling Notation . . . . . 99
- 3.2.4 Conformance to a Metamodel . . . . . 100
- 3.2.5 Illustrative Examples of Metamodels . . . . . 101
  - 3.2.5.1 Syntax of Finite State Machines . . . . . 101
  - 3.2.5.2 Syntax of Simple Functional Programs . . . . . 102
- 3.3 Context Conditions . . . . . 102
- 3.4 The Metametalevel . . . . . 103
  - 3.4.1 The Signature of Signatures . . . . . 104
  - 3.4.2 The Signature of Metamodels . . . . . 105
  - 3.4.3 The Metamodel of Metamodels . . . . . 106
- Summary and outline . . . . . 107
- References . . . . . 108
  
- 4 Representation of Object Programs in Metaprograms . . . . . 109**
  - 4.1 Representation Options . . . . . 110
    - 4.1.1 Untyped Representation . . . . . 110
    - 4.1.2 Universal Representation . . . . . 111
    - 4.1.3 Typeful Representation . . . . . 112
      - 4.1.3.1 Algebraic Data Type-Based Representation . . . . . 112
      - 4.1.3.2 Object-Based Representation . . . . . 114
      - 4.1.3.3 Reference Relationships . . . . . 115
      - 4.1.3.4 Smart Constructors . . . . . 118
    - 4.1.4 Interchange Formats . . . . . 120
      - 4.1.4.1 JSON Representation . . . . . 120
      - 4.1.4.2 XML Representation . . . . . 121
  - 4.2 Conformance Checking . . . . . 122
    - 4.2.1 Language-Specific Conformance Checking . . . . . 122
    - 4.2.2 Generic Conformance Checking . . . . . 123
    - 4.2.3 Schema-Based Conformance Checking . . . . . 125
  - 4.3 Serialization . . . . . 128
  - 4.4 AST-to-ASG Mapping . . . . . 129
  - Summary and outline . . . . . 133
  - References . . . . . 133
  
- 5 A Suite of Metaprogramming Scenarios . . . . . 135**
  - 5.1 Interpretation . . . . . 136
    - 5.1.1 Basics of Interpretation . . . . . 136
    - 5.1.2 Interpretation with Stores . . . . . 138
    - 5.1.3 Interpretation with Environments . . . . . 141
    - 5.1.4 Stepwise Interpretation . . . . . 143
  - 5.2 Compilation . . . . . 146
    - 5.2.1 Architecture of a Compiler . . . . . 147
    - 5.2.2 Translation to Assembly Code . . . . . 148
    - 5.2.3 Translation to Machine Code . . . . . 151

5.3	Analysis	154
5.3.1	Type Checking	154
5.3.2	Well-Formedness Checking	156
5.3.3	Fact Extraction	159
5.4	Transformation	161
5.4.1	Optimization	162
5.4.2	Refactoring	164
5.5	Composition	169
	Summary and outline	173
	References	173
<b>6</b>	<b>Foundations of Textual Concrete Syntax</b>	<b>177</b>
6.1	Textual Concrete Syntax	178
6.1.1	A Basic Grammar Notation	178
6.1.2	Derivation of Strings	179
6.1.3	An Extended Grammar Notation	180
6.1.4	Illustrative Examples of Grammars	181
6.1.4.1	Syntax of Simple Expressions	181
6.1.4.2	Syntax of Simple Imperative Programs	181
6.1.4.3	Syntax of Simple Functional Programs	183
6.1.4.4	Syntax of Finite State Machines	183
6.2	Concrete versus Abstract Syntax	184
6.3	Languages as Sets of Strings	186
6.3.1	Context-Free Grammars	186
6.3.2	The Language Generated by a Grammar	187
6.3.3	Well-Formed Grammars	187
6.3.4	The Notion of Acceptance	188
6.4	Languages as Sets of Trees	188
6.4.1	Concrete Syntax Trees	189
6.4.2	The Notion of Parsing	190
6.4.3	Ambiguous Grammars	190
6.5	Lexical Syntax	192
6.6	The Metametalevel	194
6.6.1	The Signature of Grammars	194
6.6.2	The Signature of Concrete Syntax Trees	196
6.6.3	The Grammar of Grammars	197
6.6.4	The Grammar of Signatures	198
6.6.5	The Grammar of Metamodels	199
	Summary and outline	200
	References	200
<b>7</b>	<b>Implementation of Textual Concrete Syntax</b>	<b>201</b>
7.1	Representations and Mappings	202
7.2	Parsing	204
7.2.1	Basic Parsing Algorithms	204



7.2.1.1	Top-Down Acceptance	204
7.2.1.2	Bottom-Up Acceptance	209
7.2.1.3	Top-Down Parsing	212
7.2.1.4	Bottom-Up Parsing	213
7.2.2	Recursive Descent Parsing	213
7.2.3	Parser Generation	217
7.2.4	Parser Combinators	218
7.3	Abstraction	220
7.3.1	Recursive Descent Parsing	221
7.3.2	Semantic Actions	222
7.3.3	Parser Combinators	224
7.3.4	Text-to-Model	225
7.4	Formatting	226
7.4.1	Pretty Printing Combinators	226
7.4.2	Template Processing	228
7.5	Concrete Object Syntax	231
7.5.1	Quotation	232
7.5.2	Antiquotation	234
	Summary and outline	237
	References	238
<b>8</b>	<b>A Primer on Operational Semantics</b>	<b>241</b>
8.1	Big-step Operational Semantics	242
8.1.1	Metavariables	242
8.1.2	Judgments	242
8.1.3	Inference Rules	243
8.1.4	Derivation Trees	246
8.1.5	Big-Step Style Interpreters	247
8.1.5.1	Aspects of Implementation	247
8.1.5.2	Explicit Model of Failure	250
8.1.5.3	Rule-by-Rule Mapping	251
8.1.6	More Examples of Big-Step Style	253
8.1.6.1	Semantics of Simple Imperative Programs	253
8.1.6.2	Semantics of Simple Functional Programs	257
8.2	Small-Step Operational Semantics	258
8.2.1	Big- versus Small-Step Judgments	259
8.2.2	Normal Form	260
8.2.3	Derivation Sequences	261
8.2.4	Small-Step Style Interpreters	263
8.2.5	More Examples of Small-Step Style	264
8.2.5.1	Semantics of Simple Imperative Programs	264
8.2.5.2	Semantics of Simple Functional Programs	267
8.2.5.3	Semantics of Finite State Machines	269
	Summary and outline	270
	References	270

<b>9</b>	<b>A Primer on Type Systems</b>	271
9.1	Types	272
9.2	Typing Judgments	272
9.3	Typing Rules	273
9.4	Typing Derivations	274
9.5	Type Safety	274
9.6	Type Checking	277
9.7	More Examples of Type Systems	278
9.7.1	Well-Typedness of Simple Imperative Programs	278
9.7.2	Well-Typedness of Simple Functional Programs	284
9.7.3	Well-Formedness of Finite State Machines	287
	Summary and outline	287
	References	288
<b>10</b>	<b>An Excursion into the Lambda Calculus</b>	289
10.1	The Untyped Lambda Calculus	290
10.1.1	Syntax	290
10.1.2	Semantics	291
10.1.3	Substitution	292
10.1.4	Predefined Values and Operations	294
10.1.5	Fixed-Point Computation	295
10.1.6	Interpretation	296
10.1.7	Turing Completeness	298
10.2	The Simply Typed Lambda Calculus	299
10.2.1	Syntax	299
10.2.2	Semantics	300
10.2.3	Type System	300
10.2.4	Type Checking	301
10.2.5	Type Erasure	302
10.3	System $F$	303
10.3.1	Syntax	304
10.3.2	Semantics	305
10.3.3	Type System	306
10.3.4	Type Erasure	307
10.4	Type-System Extensions	309
10.4.1	Records and Variants	309
10.4.2	Structural Type Equivalence	312
10.4.3	Structural Subtyping	312
10.4.4	Nominal Typing	315
	Summary and outline	318
	References	318

<b>11</b>	<b>An Ode to Compositionality</b>	319
11.1	Compositionality	320
11.2	Direct Style	320
11.2.1	Semantic Domains	321
11.2.2	Semantic Functions	321
11.2.3	Semantic Combinators	322
11.2.4	Fixed-Point Semantics	323
11.2.5	Direct-Style Interpreters	325
11.3	Continuation Style	328
11.3.1	Continuations	328
11.3.2	Continuation-Style Interpreters	329
11.3.3	Semantics of Gotos	330
	Summary and outline	333
	References	334
<b>12</b>	<b>A Suite of Metaprogramming Techniques</b>	335
12.1	Term Rewriting	336
12.1.1	Rewrite Rules	336
12.1.2	Encoding Rewrite Rules	338
12.1.3	Normalization	340
12.1.4	Strategic Programming	341
12.1.5	Rewriting-Related concerns	345
12.1.5.1	Other Traversal Idioms	345
12.1.5.2	Concrete Object Syntax	346
12.1.5.3	Graph Rewriting and Model Transformation	346
12.1.5.4	Origin Tracking	346
12.1.5.5	Layout Preservation	346
12.2	Attribute Grammars	347
12.2.1	The Basic Attribute Grammar Formalism	347
12.2.2	Attribute Evaluation	350
12.2.3	Attribute Grammars as Functional Programs	354
12.2.4	Attribute Grammars with Conditions	356
12.2.5	Semantic Actions with Attributes	358
12.3	Multi-Stage Programming	363
12.3.1	Inlining as an Optimization Scenario	364
12.3.2	Quasi-Quotation and Splicing	364
12.3.3	More Typeful Staging	366
12.4	Partial Evaluation	368
12.4.1	The Notion of a Residual Program	368
12.4.2	Interpretation with Inlining	370
12.4.3	Interpreter with Memoization	375
12.5	Abstract Interpretation	380
12.5.1	Sign Detection as an Optimization Scenario	380
12.5.2	Semantic Algebras	381
12.5.3	Concrete Domains	382

12.5.4	Abstract Domains . . . . .	383
12.5.5	Examples of Abstract Interpreters . . . . .	386
12.5.5.1	A Type-Checking Interpreter . . . . .	386
12.5.5.2	A Sign-Detection Interpreter . . . . .	388
	Summary and outline . . . . .	393
	References . . . . .	394
<b>Postface</b>	. . . . .	399
	The importance of Software Language Engineering . . . . .	399
	Software Languages: Key Concepts . . . . .	400
	Omissions in This Book . . . . .	401
	Complementary Textbooks . . . . .	403
	Software Languages in Academia . . . . .	405
	Feedback Appreciated . . . . .	408
	References . . . . .	408
<b>Index</b>	. . . . .	415

# List of Recipes

- 2.1 Recipe (Development of a fluent API) ..... 61
- 2.2 Recipe (Development of an interpreter) ..... 63
- 2.3 Recipe (Development of a constraint checker)..... 66
- 2.4 Recipe (Authoring a grammar) ..... 67
- 2.5 Recipe (Development of a syntax checker) ..... 69
- 2.6 Recipe (Development of a parser) ..... 73
- 2.7 Recipe (Development of a code generator) ..... 82
  
- 3.1 Recipe (Authoring an abstract syntax definition)..... 107
  
- 4.1 Recipe (Implementation of a conformance checker)..... 133
  
- 5.1 Recipe (Development of an interpreter (continued)) ..... 138
- 5.2 Recipe (Development of a software transformation)..... 167
  
- 8.1 Recipe (Implementation of inference rules) ..... 247
  
- 11.1 Recipe (Compositional interpretation) ..... 327
  
- 12.1 Recipe (Design of a strategic program) ..... 345
- 12.2 Recipe (Design of an attribute grammar) ..... 362
- 12.3 Recipe (Design of a multi-stage program) ..... 368
- 12.4 Recipe (Design of a partial evaluator) ..... 379
- 12.5 Recipe (Design of an abstract interpreter) ..... 393

# Acronyms

## Fabricated Languages

In this book, several software languages have been “fabricated” to capture core design aspects of diverse real-world software languages. See Section [1.1.2](#) for a detailed discussion. Here is a summary:

BAL	Basic Assembly Language
BFPL	Basic Functional Programming Language
BGL	Basic Grammar Language
BIPL	Basic Imperative Programming Language
BL	Buddy Language
BML	Basic Machine Language
BNL	Binary Number Language
BSL	Basic Signature Language
BTL	Basic TAPL Language
EFPL	Extended Functional Programming Language
EGL	Extended Grammar Language
EIPL	Extended Imperative Programming Language
EL	Expression Language
ESL	Extended Signature Language
FSML	Finite State Machine Language
MML	MetaModeling Language
TLL	Typed Lambda Language
ULL	Untyped Lambda Language

## Other Acronyms

ADT	abstract data type
AG	attribute grammar
AOP	aspect-oriented programming
ASG	abstract syntax graph
AST	abstract syntax tree
BNF	Backus Naur form
ccpo	chain complete partial order
CFG	context-free grammar
COP	context-oriented programming
CPS	continuation-passing style
CST	concrete syntax tree
DSL	domain-specific language
DSML	domain-specific modeling language
EBNF	extended Backus Naur form
FSM	finite state machine
IDE	integrated development environment
IR	intermediate representation
JIT	just in time
LMS	lightweight modular staging
MDE	model-driven engineering
OO	object oriented/orientation
OOP	object-oriented programming
PEG	parsing expression grammar
RDF	resource description framework
SLR	software language repository
UML	unified modeling language

# Chapter 1

## The Notion of a Software Language



JEAN-MARIE FAVRE.<sup>1</sup>

**Abstract** In this chapter, we characterize the notion of “software language” in a broad sense. We begin by setting out diverse examples of programming, modeling, and specification languages to cover a wide range of use cases of software languages in software engineering. Then, we classify software languages along multiple dimensions and describe the lifecycle of software languages, with phases such as language definition and implementation. Finally, we identify areas in software engineering that involve software languages in different ways, for example, software reverse engineering and software re-engineering.

---

<sup>1</sup> When the “Software Languages” community was formed around 2005–2007, Jean-Marie Favre was perhaps the key pillar and visionary and community engineer. His views and interests are captured very well in publications like these: [105, 104, 106, 100, 103].

**Artwork Credits for Chapter Opening:** This work by Wojciech Kwasnik is licensed under CC BY-SA 4.0. This artwork quotes the artwork *DMT*, acrylic, 2006 by Matt Sheehy with the artist's permission. This work also quotes [https://commons.wikimedia.org/wiki/File:Vincent\\_van\\_Gogh\\_-\\_Zeegezicht\\_bij\\_Les\\_Saintes-Maries-de-la-Mer\\_-\\_Google\\_Art\\_Project.jpg](https://commons.wikimedia.org/wiki/File:Vincent_van_Gogh_-_Zeegezicht_bij_Les_Saintes-Maries-de-la-Mer_-_Google_Art_Project.jpg), subject to the attribution “Vincent van Gogh: Seascape near Les Saintes-Maries-de-la-Mer (1888) [Public domain], via Wikimedia Commons.” This work artistically morphes an image, <https://www.flickr.com/photos/eelcovisser/4772847104>, showing the person honored, subject to the attribution “Permission granted by Eelco Visser for use in this book.”



## 1.1 Examples of Software Languages

In this book, we discuss diverse software languages; we may use them for illustrative purposes, and we may even define or implement them or some subsets thereof. For clarity, we would like to enumerate all these languages here in one place so that the reader will get an impression of the “language-related profile” of this book.

### 1.1.1 Real-World Software Languages

By “real-world language”, we mean a language that exists independently of this book and is more or less well known. We begin with *programming languages* that will be used for illustrative code in this book. We order these languages loosely in terms of their significance in this book.

- *Haskell*<sup>2</sup>: The functional programming language Haskell
- *Java*<sup>3</sup>: The Java programming language
- *Python*<sup>4</sup>: The dynamic programming language Python

We will use some additional software languages in this book; these languages serve the purpose of specification, modeling, or data exchange rather than programming; we order these languages alphabetically.

- *ANTLR*<sup>5</sup>: The grammar notation of the ANTLR technology
- *JSON*<sup>6</sup>: The JavaScript Object Notation
- *JSON Schema*<sup>7</sup>: The JSON Schema language
- *XML*<sup>8</sup>: Extensible Markup Language
- *XSD*<sup>9</sup>: XML Schema Definition

Furthermore, we will refer to diverse software languages in different contexts, for example, for the purpose of language classification in Section 1.2; we order these languages alphabetically.

- *Alloy*<sup>10</sup>: The Alloy specification language
- *CIL*<sup>11</sup>: Bytecode of .NET’s CLR

---

<sup>2</sup> Haskell language: <https://www.haskell.org/>

<sup>3</sup> Java language: [https://en.wikipedia.org/wiki/Java\\_\(programming\\_language\)](https://en.wikipedia.org/wiki/Java_(programming_language))

<sup>4</sup> Python language: <https://www.python.org/>

<sup>5</sup> ANTLR language: <http://www.antlr.org/>

<sup>6</sup> JSON language: <https://en.wikipedia.org/wiki/JSON>

<sup>7</sup> JSON Schema language: <http://json-schema.org/>

<sup>8</sup> XML language: <https://en.wikipedia.org/wiki/XML>

<sup>9</sup> XSD language: [https://en.wikipedia.org/wiki/XML\\_Schema\\_\(W3C\)](https://en.wikipedia.org/wiki/XML_Schema_(W3C))

<sup>10</sup> Alloy language: <http://alloy.mit.edu/alloy/>

<sup>11</sup> CIL language: [https://en.wikipedia.org/wiki/Common\\_Intermediate\\_Language](https://en.wikipedia.org/wiki/Common_Intermediate_Language)

- *Common Log Format*<sup>12</sup>: The NCSA Common log format
- *DocBook*<sup>13</sup>: The DocBook semantic markup language for documentation
- *FOAF*<sup>14</sup>: The friend of a friend ontology
- *INI file*<sup>15</sup>: The INI file format
- *Java bytecode*<sup>16</sup>: Bytecode of the JVM
- *make*<sup>17</sup>: The make tool and its language
- *OWL*<sup>18</sup>: Web Ontology Language
- *Prolog*<sup>19</sup>: The logic programming language Prolog
- *QTFF*<sup>20</sup>: QuickTime File Format
- *RDF*<sup>21</sup>: Resource Description Framework
- *RDFS*<sup>22</sup>: RDF Schema
- *Scala*<sup>23</sup>: The functional OO programming language Scala
- *Smalltalk*<sup>24</sup>: The OO reflective programming language Smalltalk
- *SPARQL*<sup>25</sup>: SPARQL Protocol and RDF Query Language
- *UML*<sup>26</sup>: Unified Modeling Language
- *XPath*<sup>27</sup>: The XML path language for querying
- *XSLT*<sup>28</sup>: Extensible Stylesheet Language Transformations

### 1.1.2 Fabricated Software Languages

In this book, we “fabricated” a few software languages: these are small, idealized languages that have been specifically designed and implemented for the purposes of the book, although in fact these languages are actual or de facto subsets of real-world software languages. The language names are typically acronyms with expansions hinting at the nature of the languages. Language definitions of language-based

---

<sup>12</sup> Common Log Format language: [https://en.wikipedia.org/wiki/Common\\_Log\\_Format](https://en.wikipedia.org/wiki/Common_Log_Format)

<sup>13</sup> DocBook language: <https://en.wikipedia.org/wiki/DocBook>

<sup>14</sup> FOAF language: <http://semanticweb.org/wiki/FOAF.html>

<sup>15</sup> INI file language: [https://en.wikipedia.org/wiki/INI\\_file](https://en.wikipedia.org/wiki/INI_file)

<sup>16</sup> Java bytecode language: [https://en.wikipedia.org/wiki/Java\\_bytecode](https://en.wikipedia.org/wiki/Java_bytecode)

<sup>17</sup> make language: [https://en.wikipedia.org/wiki/Make\\_\(software\)](https://en.wikipedia.org/wiki/Make_(software))

<sup>18</sup> OWL language: [https://en.wikipedia.org/wiki/Web\\_Ontology\\_Language](https://en.wikipedia.org/wiki/Web_Ontology_Language)

<sup>19</sup> Prolog language: <https://en.wikipedia.org/wiki/Prolog>

<sup>20</sup> QTFF language: [https://en.wikipedia.org/wiki/QuickTime\\_File\\_Format](https://en.wikipedia.org/wiki/QuickTime_File_Format)

<sup>21</sup> RDF language: <https://www.w3.org/RDF/>

<sup>22</sup> RDFS language: <https://www.w3.org/TR/rdf-schema/>

<sup>23</sup> Scala language: [https://en.wikipedia.org/wiki/Scala\\_\(programming\\_language\)](https://en.wikipedia.org/wiki/Scala_(programming_language))

<sup>24</sup> Smalltalk language: <https://en.wikipedia.org/wiki/Smalltalk>

<sup>25</sup> SPARQL language: <https://en.wikipedia.org/wiki/SPARQL>

<sup>26</sup> UML language: [https://en.wikipedia.org/wiki/Unified\\_Modeling\\_Language](https://en.wikipedia.org/wiki/Unified_Modeling_Language)

<sup>27</sup> XPath language: <https://en.wikipedia.org/wiki/XPath>

<sup>28</sup> XSLT language: <https://en.wikipedia.org/wiki/XSLT>

software components are available for these languages from the book's repository.<sup>29</sup> The footnotes in the following list link to the repository locations for the languages.

- [\*BAL\*](#): *Basic Assembly Language*
- [\*BFPL\*](#): *Basic Functional Programming Language*
- [\*BGL\*](#): *Basic Grammar Language*
- [\*BIPL\*](#): *Basic Imperative Programming Language*
- [\*BML\*](#): *Binary Machine Language*
- [\*BNL\*](#): *Binary Number Language*
- [\*BSL\*](#): *Basic Signature Language*
- [\*BTL\*](#): *Basic TAPL Language*
- [\*BL\*](#): *Buddy Language*
- [\*EFPL\*](#): *Extended Functional Programming Language*
- [\*EGL\*](#): *Extended Grammar Language*
- [\*EIPL\*](#): *Extended Imperative Programming Language*
- [\*EL\*](#): *Expression Language*
- [\*ESL\*](#): *Extended Signature Language*
- [\*FSML\*](#): *Finite State Machine Language*
- [\*MML\*](#): *Meta Modeling Language*
- [\*TLL\*](#): *Typed Lambda Language*
- [\*Text\*](#): The “language” of text (such as Unicode 8.0 strings)
- [\*ULL\*](#): *Untyped Lambda Language*

In the rest of this section, we quickly introduce some of these languages, thereby providing a first indication of the diversity of language aspects covered by the book.

**Binary Number Language (BNL)** A trivial language of binary numbers with an intended semantics that maps binary to decimal values.

**Basic TAPL Language (BTL)** A trivial expression language in reference to the TAPL textbook (Types and programming languages [210]).

**Buddy Language (BL)** A trivial language for modeling persons in terms of their names and buddy relationships.

**Basic Functional Programming Language (BFPL)** A really simple functional programming language which is an actual syntactic subset of the established programming language *Haskell*.

**Basic Imperative Programming Language (BIPL)** A really simple imperative programming language which is a de-facto subset of the established programming language C.

**Finite State Machine Language (FSML)** A really simple language for behavioral modeling which is variation on statecharts of the established modeling language UML.

**Basic Grammar Language (BGL)** A specification language for concrete syntax, which can also be executed for the purpose of parsing; it is a variation on the established Backus-Naur form (BNF).

---

<sup>29</sup> <http://github.com/softlang/yas>

### 1.1.2.1 BNL: A Language of Binary Numbers

We introduce *BNL* (*Binary Number Language*). This is a trivial language whose elements are essentially the binary numbers. Here are some binary numbers and their associated “interpretations” as decimal numbers:

- **0**: 0 as a decimal number;
- **1**: 1 as a decimal number;
- **10**: 2 as a decimal number;
- **11**: 3 as a decimal number;
- **100**: 4 as a decimal number;
- **101**: 5 as a decimal number;
- **101.01**: 5.25 as a decimal number.

Thus, the language contains integer and rational numbers – only positive ones, as it happens. BNL is a trivial language that is nevertheless sufficient to discuss the most basic aspects of software languages such as *syntax* and *semantics*. A syntax definition of BNL should define valid sequences of digits, possibly containing a period. A semantics definition of BNL could map binary to decimal numbers. We will discuss BNL’s abstract syntax in Chapter 3 and the concrete syntax in Chapter 6.

### 1.1.2.2 BTL: An Expression Language

We introduce *BTL* (*Basic TAPL Language*). This is a trivial language whose elements are essentially expressions over natural numbers and Boolean values. Here is a simple expression:

```
pred if iszero zero then succ succ zero else zero
```

The meaning of such expressions should be defined by expression evaluation. For instance, the expression form `iszero  $e$`  corresponds to a test of whether  $e$  evaluates to the natural number zero; evaluation of the form is thus assumed to return a Boolean value. The expression shown above evaluates to zero because `iszero zero` should compute to true, making the if-expression select the then-branch `succ succ zero`, the predecessor of which is `succ zero`.

An interpreter of BTL expressions should recursively evaluate BTL expression forms. BTL is a trivial language that is nevertheless sufficient to discuss basic aspects of interpretation (Chapter 5), semantics (Chapter 8), and type systems (Chapter 9).

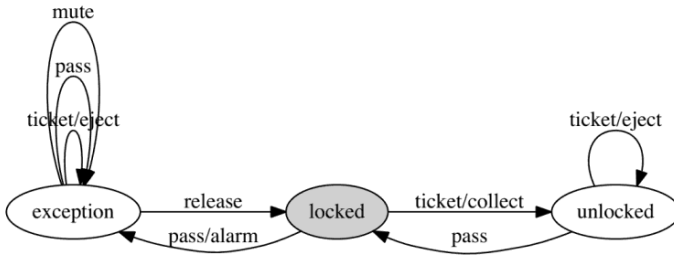


Fig. 1.2 A finite state machine for a turnstile.

The FSM also identifies possible transitions between states triggered by “events,” possibly causing “actions”; see the edges in the visual notation.

These are these states in the turnstile FSM:

- *locked*: The turnstile is locked. No passenger is allowed to pass.
- *unlocked*: The turnstile is unlocked. A passenger may pass.
- *exception*: A problem has occurred and metro personnel need to intervene.

There are input symbols which correspond to the events that a user or the environment may trigger. There are output symbols which correspond to the actions that the state machine should perform upon a transition. These are some of the events and actions of the turnstile FSM:

- Event *ticket*: A passenger enters a ticket into the card reader.
- Event *pass*: A passenger passes through the turnstile, as noticed by a sensor.
- Action *collect*: The ticket is collected by the card reader.
- Action *alarm*: An alarm is turned on, thereby requesting metro personnel.

The meanings of the various transitions should be clear. Consider, for example, the transition from the source state “locked” to the target state “unlocked”, which is annotated by “ticket/collect” to mean that the transition is triggered by entering a ticket and the transition causes ticket collection to happen.

FSML is a domain-specific modeling language (DSML). FSML supports *state-based modeling* of systems. The specification can be executed to simulate possible behaviors of a turnstile. The specification could also be used to generate a code skeleton for controlling an actual turnstile, as part of an actual metro system. FSML is a trivial language that can be used to discuss basic aspects of domain-specific language definition and implementation. For what it matters, languages for state-based behavior are widely established in software and systems engineering. For instance, the established modeling language UML consists, in fact, of several modeling languages; UML’s state machine diagrams are more general than FSML. We will discuss FSML in detail in Chapter 2.

### 1.1.2.7 BGL: A Language for Context-Free Grammars

We introduce *BGL* (*Basic Grammar Language*). This language can be used to define the *concrete textual syntax* of other software languages. Thus, BGL gets us to the metalevel. Here is an illustration of BGL – a definition of the syntax of BNL – the language of binary numbers, as introduced earlier:

```
[number] number : bits rest ; // A binary number
[single] bits : bit ; // A single bit
[many] bits : bit bits ; // More than one bit
[zero] bit : '0' ; // The zero bit
[one] bit : '1' ; // The nonzero bit
[integer] rest : ; // An integer number
[rational] rest : '.' bits ; // A rational number
```

Each line is a grammar production (a rule) with the syntactic category (or the so-called nonterminal) to the left of “:” and its definition to the right of “:”. For instance, the first production defines that a binary number consists of a bit sequence bits for the integer part followed by rest for the optional rational part. The right-hand phrases compose so-called terminals (“0”, “1”, and “.”) and nonterminals (bit, bits, rest, and number) by juxtaposition. The rules are labeled, thereby giving a name to each construct.

BGL is a domain-specific modeling language in that it supports modeling (or specifying or defining) concrete textual syntax. One may “execute” BGL in different ways. Most obviously, one may execute a BGL grammar for the purpose of accepting or *parsing* input according to the syntax defined. BGL, like many other notations for syntax definition, is grounded in the fundamental formalism of *context-free grammars* (CFGs). BGL is a variation on BNF [21]. There exist many real-world notations for syntax definition [277]; they are usually more complex than BGL and may be tied to specific technology, for example, for parsing. We will develop BGL in detail in Chapter 6.

## 1.2 Classification of Software Languages

There are hundreds or even thousands of established software languages, depending on how we count them. It may be useful to group languages in an ontological manner. In particular, a *classification* of software languages (i.e., a language taxonomy) is a useful (if not necessary) pillar of a definition of “software language”.

Wikipedia, which actually uses the term “computer language” at the root of the classification, identifies the following top-level classifiers:<sup>30</sup>

- data-modeling languages;
- markup languages;
- programming languages;
- specification languages;
- stylesheet languages;
- transformation languages.

Any such branch can be classified further in terms of constructs and concepts. For instance, in the case of programming languages, there exist textbooks on programming languages, programming paradigms, and programming language theory such as [199, 232], which identify constructs and concepts. There is also scholarly work on the classification of programming languages [20, 90] and the identification of language concepts and corresponding paradigms [258].

Several classes of software languages (other than programming languages) have been identified, for example, *model transformation languages* [75], *business rule modeling languages* [239], *visual languages* [46, 49, 190], and *architecture description languages* [192]. There is more recent work aimed at the classification of software languages (or computer languages) more broadly [13, 237, 3, 171]. The *101companies* project<sup>31</sup> [102, 101, 173, 166] is also aimed at a taxonomy of software languages, but the results are of limited use, at the time of writing.

In the remainder of this section, we classify software languages along different dimensions. A key insight here is that a single classification tree is insufficient. Multiple inheritance may be needed, or orthogonal dimensions may need to be considered separately.

### 1.2.1 Classification by Paradigm

When focusing on programming languages as a category of software languages, classification may be based on the *programming paradigm*. A paradigm is characterized by a central notion of programming or computing. Here is an incomplete, classical list of paradigms:

**Imperative programming** Assignable (updatable) variables and updatable (in-place) data structures and sequential execution of statements of operations on variables. Typically, procedural abstractions capture statements that describe control flow with basic statements for updates. We exercise imperative programming with the fabricated language BIPL (Section 1.1.2.5) in this book.

---

<sup>30</sup> We show Wikipedia categories based on a particular data-cleaning effort [171]. This is just a snapshot, as Wikipedia is obviously evolving continuously.

<sup>31</sup> <http://101companies.org>

**Functional programming** The application of functions models computation with compound expressions to be reduced to values. Functions are first-class citizens in that functions may receive and return functions; these are higher-order functions. We exercise functional programming with the fabricated language BFPL (Section 1.1.2.4) in this book – however, though higher-order functions are not supported.

**Object-oriented (OO) programming** An object is a capsule of state and behavior. Objects can communicate with each other by sending messages, the same message being implementable differently by different kinds of objects. Objects also engage in structural relationships, i.e., they can participate in whole–part and reference relationships. Objects may be constructed by instantiation of a given template (e.g., a class). *Java* and *C#* are well-known OO programming languages.

**Logic programming** A program is represented as a collection of logic formulae. Program execution corresponds to some kind of proof derivation. For instance, *Prolog* is a well-known logic programming language; computation is based on depth-first, left-to-right proof search through the application of definite clauses.

There exist yet other programming or computing notions that may characterize a paradigm, for example, message passing and concurrency. Many programming languages are, in fact, *multi-paradigm* languages in that they support several paradigms. For instance, *JavaScript* is typically said to be both a functional and an imperative OO programming language and a scripting language. Programming languages may be able to support programming according to a paradigm on the basis of some encoding scheme without being considered a member of that paradigm. For instance, in *Java* prior to version 8, it was possible to encode functional programs in *Java*, while proper support was added only in version 8.

Van Roy offers a rich discussion of programming paradigms [258]. Programming concepts are the basic primitive elements used to construct programming paradigms. Often, two paradigms that seem quite different (for example, functional programming and object-oriented programming) differ by just one concept. The following are the concepts discussed by Van Roy: record, procedure, closure, continuation, thread, single assignment, (different forms of) cell (state), name (unforgeable constant), unification, search, solver, log, nondeterministic choice, (different forms of) synchronization, port (channel), clocked computation. Van Roy identifies 27 paradigms, which are characterized as sets of programming concepts. These paradigms can be clearly related in terms of the concepts that have to be added to go from one paradigm to another.

## 1.2.2 Classification by Type System

Furthermore, languages may also be classified in terms of their typing discipline or type system [210] (or the lack thereof). Here are some important options for programming languages in particular:



- Static typing** The types of variables and other abstractions (e.g., the argument and result types of methods or functions) are statically known, i.e., without executing the program – this is at compile time for compiled languages. For instance, Haskell and Java are statically typed languages.
- Dynamic typing** The types of variables and other abstractions are determined at runtime. A variable’s type is the type of the value that is stored in that variable. A method or function’s type is the one that is implied by a particular method invocation or function application. For instance, Python is a dynamically typed language.
- Duck typing** The suitability of a variable (e.g., an object variable in object-oriented programming) is determined at runtime on the basis of checking for the presence of certain methods or properties. Python uses duck typing.
- Structural typing** The equivalence or subtyping relationship between types in a static typing setting is determined on the basis of type structure, such as the components of record types. Scala supports some form of structural typing.
- Nominal typing** The equivalence or subtyping relationship between types in a static typing setting is determined on the basis of explicit type names and declared relationships between them. Java’s reference types (classes and interfaces including “extends” and “implements” relationships) commit to nominal typing.

### 1.2.3 Classification by Purpose

Languages may be classified on the basis of the *purpose* of the language (its usage) or its elements. Admittedly, the term “purpose” may be somewhat vague, but the illustrative classifiers in Table 1.1 may convey our intuition. We offer two views: the purpose of the language versus that of its elements; these two views are very similar.

**Table 1.1** Classification by the purpose of language elements

Purpose (language)	Purpose (element)	Classifier	Example
Programming	Program	Programming language	Java
Querying	Query	Query language	XPath
Transformation	Transformation	Transformation language	XSLT
Modeling	Model	Modeling language	UML
Specification	Specification	Specification language	Alloy
Data representation	Data	Data format	QTFF (QuickTime file format)
Documentation	Documentation	Documentation language	DocBook
Configuration	Configuration	Configuration language	INI file
Logging	Log	Log format	Common Log Format
...	...	...	...

tion of Section 1.1.2.7. Several of the illustrative languages of Section 1.1 were introduced as string languages.

**Tree language** (See also “markup language” below.) Language elements are represented, viewed, and edited as trees, for example, as XML trees or JSON dictionaries. A tree language is defined in terms of a suitable grammar or data modeling notation, for example, XSD in the case of XML. As it happens, we did not present any tree languages in Section 1.1.2.7, but we will discuss tree-based abstract syntax definitions later for some of the string languages that we have already seen. Tree languages play an important role in language implementation.

**Graph language** Language elements are represented, viewed, and edited as graphs, i.e., more or less constrained collections of nodes and edges. Appropriate grammar and data modeling notations exist for this case as well. The language BL for buddy relationships (Section 1.1.2.3) was introduced as a graph language and we hinted at a visual concrete syntax. A graph language may be coupled with a string or tree language in the sense of alternative representations of the same “conceptual” language. For instance, BL may be represented in a string-, tree-, or graph-based manner.

## 1.2.6 Classification by Notation

One may also distinguish languages in terms of notation; this classification is very similar to the classification by representation:

**Textual (text) language** This is essentially a synonym for “string language”.

**Markup language** Markup, as in XML, is used as the main principle for expressing language elements. The use of markup is one popular notation for tree languages. With an appropriate semantics of identities, markup can also be used as a notation for graphs. Not every tree language relies on markup for the notation. For instance, JSON provides another, more dictionary-oriented notation for tree languages.

**Visual (graphical) language** A visual notation is used. The languages BL for buddy relationships (Section 1.1.2.3) and FSML for state-based modeling (Section 1.1.2.6) were introduced in terms of a visual notation.

## 1.2.7 Classification by Degree of Declarativeness

An (executable) language may be said to be (more or less) *declarative*. It turns out to be hard to identify a consensual definition of declarativeness, but this style of classification is nevertheless common. For instance, one may say that programs (or models) of a declarative language describe more the “what” than the “how”. That is, a declarative program’s semantics is not strongly tied to execution order.

Let us review the languages of Section 1.1:

**Binary Number Language (BNL)** A trivial language.

**Buddy Language (BL)** A trivial language.

**Basic Functional Programming Language (BFPL)** This language is “pure”, i.e., free of side effects. Regardless of the evaluation order of subexpressions, complete evaluation of a main expression should lead to the same result – modulo some constraints to preserve termination. For instance, argument expressions of a function application could be evaluated in different orders without affecting the result. Thus, BFPL is a declarative programming language.

**Basic Imperative Programming Language (BIPL)** This language features imperative variables such that the execution order of statements affects the result of computation. Thus, BIPL is not a declarative programming language.

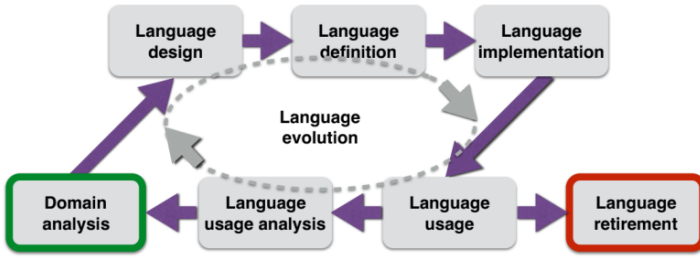
**Finite State Machine Language (FSML)** This language models finite states and event- and action-labeled transitions between states. Its actual semantics or execution order is driven by an event sequence. FSML would usually be regarded as a declarative (modeling) language.

**Basic Grammar Language (BGL)** This grammar notation defines sets of strings in a rule-based manner. Thus, BGL’s most fundamental semantics is “declarative” in the sense that it is purely mathematical, without reference to any operational details. Eventually, we may attach a more or less constrained operational interpretation to BGL so that we can use it for efficient, deterministic parsing. Until that point, BGL would be usually regarded as a declarative (specification) language.

We may also consider subcategories of declarative languages such that it is emphasized how declarativeness is achieved. Two examples may suffice:

**Rule-based language** Programs are composed from rules, where a rule typically combines a condition and an action part. The condition part states when the rule is applicable; the action part states the implications of rule application. Some logic programming languages, for example, Prolog, can be very well considered to be rule-based languages. Some schemes for using functional programming, for example, in interpretation or program transformation, also adopt the rule-based approach. Event-driven approaches may also use rules with an additional “event” component, for example, the “event condition action” (ECA) paradigm, as used in active databases [88].

**Constraint-based language** Programs involve constraints as means of selecting or combining computations. These constraints are aggregated during program execution and constraint resolution is leveraged to establish whether and how given constraints can be solved. For instance, there exist various constraint-logic programming languages which enrich basic logic programming with constraints on sets of algebras for numbers [117].



**Fig. 1.3** The lifecycle of a software language. The nodes denote phases of the lifecycle. The edges denote transitions between these phases. The lifecycle starts with *domain analysis*. The lifecycle ends (theoretically) with *language retirement*. We may enter cycles owing to language evolution.

---

**Exercise 1.2** (Classification of *make*) [Basic level]  
 Study the language used by the well-known *make* utility and argue that the language is declarative and identify what subcategory of declarativeness applies.

---

### 1.3 The Lifecycle of Software Languages

The notion of a software lifecycle can be usefully adopted for languages. That is, a language goes through a lifecycle, possibly iterating or skipping some phases; see Fig. 1.3. These phases are described in some detail as follows:

**Domain analysis** A domain analysis is required to discover the domain that is to be addressed by a new language. A domain analysis answers these questions: What are the central concepts in the domain? For instance, the central concepts are states and transitions between states for FSML, differential equations and their solution in a language for weather forecasts, or layouts and their rendering in a language for HTML/XML stylesheets. These concepts form the foundation of language design and for everything onwards. Arguably, no domain analysis is performed for general-purpose programming languages.

**Language design** The domain concepts are mapped, at the design level of abstraction, to language constructs and language concepts. The emerging language should be classified in terms of paradigm, degree of declarativeness, and other characteristics. A language may be presented as a composition of very specific language constructs as well as reusable language constructs, for example, for basic expressions or modules. The first samples are written so that a syntax emerges, and the transition to the phase of language definition has then begun.

**Language definition** The language design is refined into a language definition. Most notably, the syntax and semantics of the language are defined. Assuming

an executable language definition, a first language implementation (a proof of concept) is made available for experiments, so that the transition to the phases of language implementation and usage has begun.

**Language implementation** The language is properly implemented. Initially, a usable and efficient compiler or interpreter needs to be provided. Eventually, additional *language processors* and tool support may be provided, for example, documentation tools, formatters, style checkers, and refactorings. Furthermore, support for an integrated development environment (IDE) may be implemented.

**Language usage** The language is used in actual software development. That is, language artifacts are “routinely” authored and a body of software artifacts acquire dependencies on the language. This is not explicitly modeled in Fig. 1.3, but the assumption is, of course, that the language implementation is continuously improved and new language processors are made available.

**Language evolution** Language definitions may be revised to incorporate new language features or respond to experience with language usage. Obviously, changes to language definitions imply work on language implementations. Language changes may even break backward compatibility, in which cases these changes will necessitate migration of existing code in those languages.

**Language usage analysis** Language evolution and the systematic improvement of domain analysis as well as language design, definition, and implementation, may benefit from language usage analysis [155, 100, 172], as an empirical element of the lifecycle. By going through the lifecycle in cycles, the language may evolve in different ways. For instance, the language may be extended so that a new version becomes available, which again needs to be implemented and put to use.

**Language retirement** In practice, languages, once adopted, are rarely retired completely, because the costs and risks of retirement are severe impediments. Retirement may still happen in the narrow scope of projects or organizations. In theory, a language may become obsolete, i.e., there are no software artifacts left that depend on that language. Otherwise, *language migration* may be considered. That is, software artifacts that depend on a language are migrated (i.e., transformed manually or automatically) to another language.

Many aspects of these phases, with some particular emphasis on the lifecycle of DSLs are discussed in [133, 272, 273, 197, 214, 56, 98, 94, 78, 265, 229]. In the present book, the focus is on language definition and implementation; we are concerned only superficially with domain analysis, language design, evolution, and retirement.

### 1.3.1 Language Definition

Let us have a deeper look at the lifecycle phase of language definition. A language is defined to facilitate implementation and use of the language. There are these aspects of language definition:

**Syntax** The definition of the syntax consists of rules that describe the valid language elements which may be drawn from different “universes”: the set of all strings (say, text), the set of all trees (of some form, e.g., XML-like trees), or the set of all graphs (of some form). Different kinds of formalisms may be used to specify the rules defining the syntax. We may distinguish *concrete* and *abstract syntax* – the former is tailored towards users who need to read and write language elements, and the latter is tailored towards language implementation. Abstract syntax is discussed in Chapters 3 and 4. Concrete syntax is discussed in Chapters 6 and 7.

**Semantics** The definition of semantics provides a mapping from the syntactic categories of a language (such as statements and expressions) to suitable domains of meanings. The actual mapping can be defined in different ways. For instance, the mapping can be defined as a set of syntax-driven inference rules which model the stepwise execution or reduction of a program; this is known as small-step operational semantics (Chapter 8). The mapping can also be applied by a translation, for example, by a model-to-model transformation in model-driven engineering (MDE).

**Pragmatics** The definition of the pragmatics explains the purpose of language concepts and provides recommendations for their usage. Language pragmatics is often defined only informally through text and samples. For instance, the pragmatics definition for a C-like language with arrays may state that arrays should be used for efficient (constant-time) access to indices in ordered collections of values of the same type. Also, arrays should be favored over (random-access) files or databases for as long as in-memory representation of the entire data structure is reasonable. In modeling languages for finite state machine (e.g., FSML), events proxy for sensors and actions proxy for actors in an embedded system.

**Types** Some languages also feature a *type system* as a part of the language definition. A type system provides a set of rules for assigning or verifying *types*, i.e., properties of language phrases, for example, different expression types such as “int” or “string” in a program with expressions. We speak of *type checking* if the type system is used to check explicitly declared types. We speak of *type inference* if the type system is used additionally to infer missing type declarations. A type system needs to be able to bind names in the sense that any use of an abstraction such as a variable, a method, or a function is linked to the corresponding declaration. Such *name binding* may be defined as part of the type system or they may be defined somewhat separately. We discuss types in detail in Chapter 9. Even when a language does not have an interesting type system, i.e., different types and rules about their use in abstractions, the language may still feature other constraints regarding, for example, the correct use of names. Thus, we may also speak of *well-formedness* more generally, as opposed to *well-typedness* more specifically. For instance, in FSML, the events handled by a given source state must be distinct for the sake of determinism.

When definitions of syntax, types, and semantics are considered formal artifacts such that these artifacts are treated in a formal (mathematical) manner, then we operate within the context of *programming language theory*. A formal approach

pass compilation”) or may be integrated into one phase (“single-pass compilation”). The components are explained more in detail as follows:

**Parser** A parser verifies the conformance of given input (i.e., text) to the syntax rules of a language and represents the input in terms of the structure defined by the rules. A parser performs parsing. Compilers and interpreters begin by parsing. Many other language processors, as discussed below, also involve parsing.

**Semantic analysis** A syntax tree only represents the structure of the source code. For any sort of nontrivial treatment such as code generation, the syntax tree needs to be enriched with attributes and links related to typing and name binding. Names with their bindings and other attributes may be aggregated in a data structure which is referred to as a symbol table or environment.

**Code generator** The enriched syntax tree is translated, more or less directly, into machine code, i.e., code of some actual or virtual machine. In particular, code generation involves resource and storage decisions such as register allocation, i.e., assigning program variables to processor registers of the target machine. In this book, few technicalities of code generation are discussed; this topic is covered perfectly by the literature on compiler construction.

Ideally, the components are described by specifications such as grammars, type systems, name-binding rules, and rewrite systems, as indicated in Fig. 1.4. In practice, the components are often implemented in a more ad hoc fashion.

This is a simplified data flow, because actual compilers may involve additional phases. That is, parsing may consist of several phases in itself: preprocessing; lexical analysis (scanning, lexing, or tokenization); syntax analysis including parse-tree construction and syntax desugaring. Also, there may be extra steps preceding code generation: translation to a (simpler) *intermediate representation* (IR) and IR-level optimization. Further, code generation may also involve *optimization* at the level of the target language and a separation between translation to assembly code, mapping to machine code, and some elements of linking. Finally, code generation may actually rely on translation such that the given input language is translated into a well-defined subset of an existing (programming) language so that an available compiler can be used afterwards.

---

**Exercise 1.3** (An exercise on language implementation) [Basic level]  
*Research the current version of the JDK (Java Development Kit) and identify and characterize at least two language implementations that are part of it.*

---

### 1.3.2.3 Classification of Language Processors

Languages are implemented in many ways other than just regular compilers and interpreters. We use the term “language processor” to refer to any sort of functionality for automated processing of software artifacts in a language-aware manner,

i.e., with more or less awareness of the syntax, types, and semantics of the artifacts. Examples of language processors include documentation generators, refactoring tools, bug checkers, and metrics calculation tools. Language processors often consist of several components and perform processing in phases, as we discussed above for compilers. Rather than classifying language processors directly, let us classify language-based software components. We do not make any claim of completeness for this classification. Several of the classifiers below will reappear in the discussion of the role of software languages across different software engineering areas (Section 1.4):

**Parser or text-to-model transformation** The term “parser” has already been introduced in the context of compilation and interpretation. The term “text-to-model transformation” is specifically used in the MDE community when one wants to emphasize that the result of parsing is not a parse *tree*, but rather a model in the sense of metamodeling, thus potentially involving, for example, references after completing name binding.

**Unparser, formatter, pretty printer, or model-to-text transformation** An artifact is formatted as text, possibly also subject to formatting conventions for the use of spaces and line breaks. Formatting may start from source code (i.e., text), concrete syntax trees (i.e., parse trees), or abstract syntax trees. Formatting is typically provided as a service in an IDE.

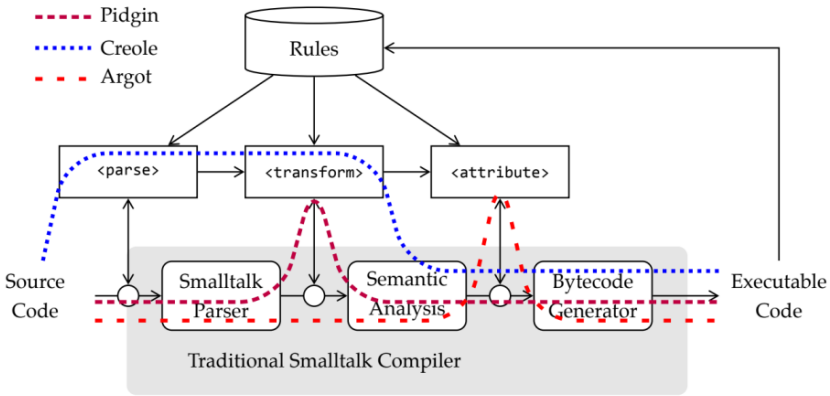
**Preprocessor** As part of parsing, code may be subject to macro expansion and conditional compilation. Such preprocessing may serve the purpose of, for example, configuration management in the sense of software variability and desugaring in the sense of language extension by macros. Interestingly, preprocessing gives rise to a language of its own for the preprocessing syntax such that the preprocessor can be seen as an interpreter of that language; the result type of this sort of interpretation is, of course, text [99]. One may also assume that a base language is extended by preprocessing constructs so that preprocessing can be modeled as a translation from the extended to the base language. In fact, some macro system work in that manner. In practice, preprocessing is often used in an undisciplined (i.e., not completely syntax-aware) manner [29, 18, 184].

**Software transformation or model-to-model transformation** A software transformation is a mapping between software languages. The term “model-to-model transformation” is used in the model transformation and MDE community. We may classify transformations in terms of whether the source and target languages are the same and whether the source and target reside at the same level of abstraction [195]. Thus:

**Exogenous transformation** The source and target languages are different, as in the case of code generation (translation) or language migration.

**Endogenous transformation** The source and target languages are the same, as in the case of program refactoring or compiler optimization [116, 196]. We can further distinguish in-place and out-place transformations [195, 35] in terms of whether the source model is “reused” to produce the target model. (Exogenous transformations are necessarily out-place transformations.)





**Fig. 1.5** The code compilation pipeline of Helvetia, showing multiple interception paths; there are hooks to intercept parsing `<parse>`, AST transformation `<transform>`, and semantic analysis `<attribute>`. Source: [215]. Additional capabilities of Helvetia support editing (coloring), debugging, etc. © 2010 Springer.

and PLT Redex [107]), *compiler frameworks* (e.g., LLVM [180]), and *modeling frameworks* (e.g., AM3 [24]).

Metaprogramming and software language engineering efforts may be “advertised” through *software language repositories* (SLRs) [165], i.e., repositories with components for language processing (interpreters, translators, analyzers, transformers, pretty printers, etc.). Further examples of SLRs include the repositories for Krishnamurthi’s textbook on programming languages [160], Batory’s Prolog-based work on teaching MDE [28], Zaytsev et al.’s software language processing suite (SLPS) [278], and Basciani et al.’s extensible web-based modeling platform MDE-Forge [26].

**1.3.2.5 Language Workbenches**

Metaprogrammers may also be supported in an interactive and integrated fashion. Accordingly, the notion of *language workbenches* [96, 97, 144, 143, 267, 266, 269, 263] encompasses enhanced metaprogramming systems that are, in fact, IDEs for language implementation. A language workbench assumes specialized language definitions that cater for IDE services such as syntax-directed, structural, or projectional editing, coloring, synthesis of warnings and errors, package exploration, quick fixes, and refactorings.

Figure 1.5 illustrates the compilation pipeline of the metaprogramming system Helvetia [214, 215]. In fact, Helvetia is an extensible development environment for embedding DSLs into a host language (Smalltalk) and its tools such as the editor and debugger. Thus, Helvetia is a language workbench.

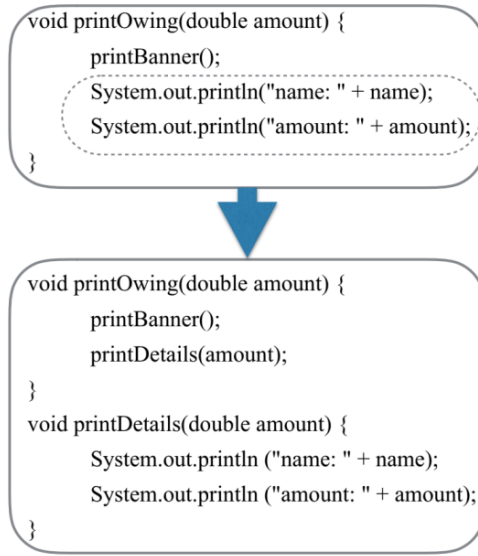


Fig. 1.6 Illustration of the “extract method” refactoring.

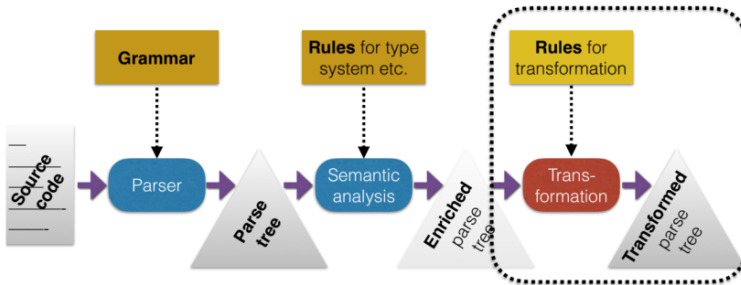
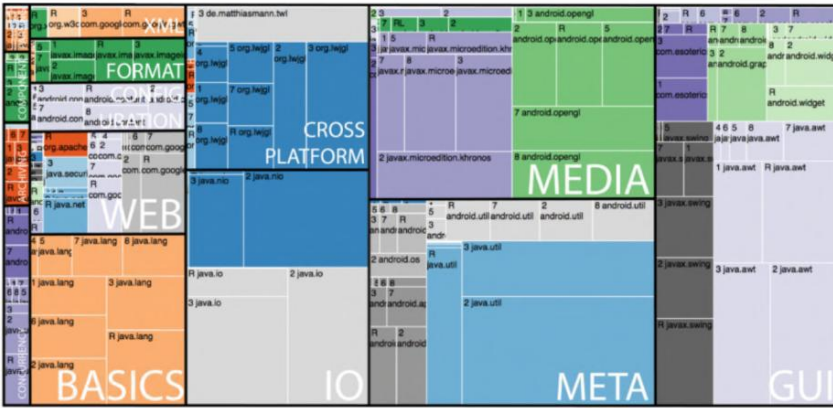


Fig. 1.7 Overall data flow for a re-engineering transformation. We have marked the phase which replaces code generation in the standard data flow for compilation.

228]. Even in the simple example at hand, some constraints have to be met; for example, the extracted statements must not return.

Figure 1.7 shows the overall data flow for a re-engineering transformation as needed, for example, for refactoring or restructuring. This data flow should be compared with the data flow for compilation; see Fig. 1.4. The two data flows share the phases of parsing and semantic analysis. The actual transformation is described (ideally) by declarative rules of a transformation language. Not every re-engineering use case requires a full-blown semantic analysis, which is why we have grayed out slightly the corresponding phase in Fig. 1.7. In fact, not even a proper syntax-aware transformation is needed in all cases, but instead a lexical approach may be applicable [152].



**Fig. 1.8** An API-usage map for an open-source Java project. The complete rectangle (in terms of its size) models the references to all APIs made by all developers. The nested rectangles partition references by domain (e.g., GUI rather than Swing or AWT). The rectangles nested further partition references by API; one color is used per API. Within each such rectangle, the contributions of distinct developers (1, . . . , 8 for the top-eight committers and “R” for the rest) are shown. Source: [4].

### 1.4.2 Software Reverse Engineering

We quote: “reverse engineering is the process of analyzing a subject system to identify the system’s components and their interrelationships and create representations of the system in another form or at a higher level of abstraction” [59]. For instance, we may extract a call graph from a system, thereby identifying call sites (such as packages, files, classes, methods, or functions) and actual calls (such as method or function calls). Reverse engineering may also be concerned with *architecture recovery* [126, 128, 158, 33], for example, the identification of components in a legacy system. Overall, reverse engineering is usually meant to help with program comprehension and to prepare for software re-engineering or to otherwise facilitate software development.

Figure 1.8 shows the visual result of a concrete reverse engineering effort aimed at understanding API usage in Java projects [4]. The tree map groups API references (i.e., source code-level references to API methods) so that we can assess the contributions of different APIs and of individual developers for each API to the project.

Figure 1.9 shows the overall data flow for a reverse engineering component that is based on the paradigm of fact extraction [109, 201, 185, 27]. Just as in the cases of compilation or transformation for re-engineering, we begin with parsing and (possibly customized) semantic analysis. The data flow differs in terms of last phase for fact extraction. The extracted facts can be thought of as sets of tuples, for example, pairs of caller/callee sites to be visualized eventually as a call graph.

Reverse engineering often starts from some sort of fact extraction. Reverse engineering may also involve data analysis based, for example, on relational alge-

Figure 1.10 gives an example of how metrics and simple visualization can be combined to analyze a software process – in this case, a process for the improvement of a grammar [7]. The changes of the values of the metrics can be explained as consequences of the specific grammar revisions applied at the corresponding commit points.

### 1.4.4 Technological Spaces

We quote: “A technological space is a working context with a set of associated concepts, body of knowledge, tools, required skills, and possibilities. It is often associated to a given user community with shared know-how, educational support, common literature and even workshop and conference regular meetings” [161].

For instance, there are the following technological spaces, which we characterize in a keyword style by pointing out associated languages, technologies, and concepts:

**Grammarware** string, grammar, parsing, CST, AST, term, rewriting, ...  
**XMLware** XML, XML infoset, DOM, DTD, XML Schema, XPath, XQuery, XSLT, ...  
**JSONware** JSON, JSON Schema, ...  
**Modelware** UML, MOF, EMF, class diagram, modeling, metamodeling, model transformation, MDE, ...  
**SQLware** table, SQL, relational model, relational algebra, ...  
**RDFware** resource, triple, Linked Data, WWW, RDF, RDFS, OWL, SPARQL, ...  
**Objectware** objects, object graphs, object models, state, behavior, ...  
**Jaware** Java, Java bytecode, JVM, Eclipse, JUnit, ...

We refer to [40] for a rather detailed discussion of one technological space – modelware (MDE). We refer to [89] for a discussion of multiple technological spaces with focus on Modelware and RDFware centric and cursory coverage of grammarware, Jaware, and XMLware and the interconnections between these spaces.

Technological spaces are deeply concerned with software languages:

**Data models** The data in a space conforms to some data model, which can be viewed as a “semantic domain” in the sense of semantics in the context of language definition. For instance, the data model of XML is defined by a certain set of trees, according to the XML infoset [274]; the data model JSON is a dictionary format that is a simple subset of Javascript objects; and the data model of SQLware is the relational model [67].

**Schema languages** Domain- or application-specific data can be defined by appropriate schema-like languages. Schemas are to tree- or graph-based data what (context-free) grammars are to string languages [149]. For instance, the schema language of JSON is JSON Schema [208]; the schema language of grammarware is EBNF [137] in many notational variations [277]; and the schema languages of

tion may be more or less aligned with an assumed ontology of language concepts. Typically, an interpreter-based approach is used for illustration. Examples include Sebesta's "Concepts of Programming Languages" [128], Sethi's "Programming Languages: Concepts and Constructs" [129] and Scott's "Programming Language Pragmatics" [127]. These books also cover, to some extent, programming language theory and compiler construction.

The present book is not concerned with a systematic discussion of programming paradigms and programming language concepts. Nevertheless, the book exercises (in fact, "defines") languages of different paradigms and discusses various language concepts in a cursory manner. This book goes beyond textbooks on programming paradigms by covering metaprogramming broadly, which is not a central concern in textbooks on paradigms.

**Compiler construction** This is the classical subject in computer science that, arguably, comes closest to the subject of software languages. Examples of textbooks on compiler construction and overall programming language implementation include Aho, Lam, Sethi, and Ullman's seminal "Compilers: Principles, Techniques, and Tools" [1], Louden's "Compiler Construction: Principles and Practice" [87], and Appel's product line of textbooks such as Appel and Palsberg's "Modern Compiler Implementation in Java" [3].

The present book briefly discusses compilation (translation), but it otherwise covers compiler construction at best superficially. For instance, lower-level code optimization and code generation are not covered. This book covers language implementation more broadly than textbooks on compiler construction, with regard to both the kinds of software languages and the kinds of language-based software components. Most notably, this book covers metaprogramming scenarios other than compilation, and metaprogramming techniques other than those used in a typical compiler.

**Hybrids** There are a number of books that touch upon several of the aforementioned topics in a significant manner. There is Krishnamurthi's "Programming Languages: Application and Interpretation" [74], which combines programming language theory and programming paradigms in a powerful manner. There is Ranta's "Implementing Programming Languages: An Introduction to Compilers and Interpreters" [116] with coverage of programming paradigms and compiler construction. There is also Stuart's "Understanding Computation: From Simple Machines to Impossible Programs" [137], which is exceptionally broad in scope: it covers various fundamental topics in computer science, including parsing and interpretation; it explains all notions covered to the working Ruby programmer in a pragmatic manner.

The present book aims at a deeper discussion of the implementation and lifecycle of software languages in the broader context of software engineering, with the central topic being metaprogramming in the sense of source-code analysis and manipulation.

**Domain-specific languages** There are some more or less recent textbooks on DSLs. Fowler's "Domain-Specific Languages" [41] discusses relatively basic or mainstream OO techniques and corresponding patterns for language implemen-