

These notes were originally prepared for our CS course at University Erlangen-Nuremberg (FAU) in Summer 2020. They are directed at 3rd semester CS undergraduates and master students but should be intelligible even for earlier students and could be interesting also for PhD students and for students from adjacent majors. The course is recommended both as a first course in the specialization area Artificial Intelligence as well as a one-off overview on on knowledge representation.

The course was developed in Summer 2020 from scratch and materials were built along the way. It integrated current directions and recent results in research on knowledge representation pulling together materials in an entirely new and original way.

Contents

0	Me	ta-Information	7							
	0.1	Remarks	7							
	0.2	Overview	8							
1	Fundamental Concepts									
	1.1	Abbreviations	9							
	1.2	Motivation	9							
		1.2.1 Knowledge	9							
		1.2.2 Representation and Processing	9							
	1.3	Components of Knowledge	10							
		1.3.1 Syntax and Semantics, Data and Knowledge	10							
		1.3.2 Semantics as Syntax Transformation	11							
		1.3.3 Heterogeneity of Semantics and Knowledge	11							
	1.4	The Tetrapod Model of Knowledge	12							
		1.4.1 Five Aspects of Knowledge	12							
		1.4.2 Relations between the Aspects	12							
2	Ont	ologies	15							
	2.1	General Principles	15							
	2.2	A Basic Ontology Language	16							
	2.3	Representing Ontologies as Triples	19							
	2.4	Writing Ontologies	21							
		2.4.1 The OWL Language	21							
		2.4.2 The Protege Tool	21							
		2.4.3 Exercise 1	21							
	2.5	Ontology Morphisms	22							
3	Syn	tax	25							
	3.1	The 4 Layers of a Formal Language	25							
	3.2	Context-Free Syntax	25							
		3.2.1 Context-Free Grammars	26							
		3.2.2 Inductive Data Types	27							
		3.2.3 Merged Definition	28							
		3.2.4 Contexts	29							
	3.3	Context-Sensitive Syntax	30							
		3.3.1 Principles	30							
		3.3.2 Vocabularies	31							
		3.3.3 Contexts and Variable Binding	31							
		3.3.4 Vocabularies vs. Contexts	32							

4 CONTENTS

	3.4	Implen	entation	33
		3.4.1	Functional Programming Languages	33
		3.4.2	Object-Oriented Programming Languages	34
		3.4.3	Combining Paradigms	35
		3.4.4	Issues Regarding Equivalence	36
		3.4.5		36
		3.4.6		39
4	Ty	pe Syste		41
	4.1	Intrinsi	e vs. Extrinsic Typing	41
		4.1.1	Overview	41
		4.1.2	Combined Definition	42
	4.2	Variabl	e Contexts and Substitutions	43
	4.3	Vocabu	lary Morphisms	44
	4.4	Abstra	et and Concrete Data Types	44
		4.4.1	Abstract vs. Concrete Types	44
		4.4.2	Concrete Data Types	45
		4.4.3	Abstract Data Types	45
	4.5	Abstra	t Data Types in Ontologies	47
		4.5.1	Motivation	47
		4.5.2	Database Schemas as Typed Ontologies	48
				49
5	Dat	ta		51
	5.1	Overvie	w	51
e	Son	nontics		K 2
6		nantics		53
6	Sen 6.1	Kinds o	f Semantics	53
6		Kinds 6.1.1	f Semantics	53 53
6	6.1	Kinds of 6.1.1 6.1.2	f Semantics	53 53 53
6	6.1	Kinds of 6.1.1 6.1.2 Compo	f Semantics	53 53 53 54
6	6.1	Kinds of 6.1.1 6.1.2 Compo	f Semantics	53 53 53 54 57
6	6.1	Kinds 6 6.1.1 6.1.2 Compo Deduct 6.3.1	f Semantics Absolute Semantics Relative Semantics sitionality ave Semantics A Basic Semantic Language: SFOL	53 53 54 57
6	6.1 6.2 6.3	Kinds of 6.1.1 6.1.2 Compo Deduct 6.3.1 6.3.2	f Semantics Absolute Semantics Relative Semantics sitionality ve Semantics A Basic Semantic Language: SFOL Semantics	53 53 54 57 57
6	6.1	Kinds of 6.1.1 6.1.2 Compose Deduct 6.3.1 6.3.2 Concress	f Semantics Absolute Semantics Relative Semantics sitionality ve Semantics A Basic Semantic Language: SFOL Semantics	53 53 54 57 57 57 60
6	6.1 6.2 6.3	Kinds of 6.1.1 6.1.2 Compo Deduct 6.3.1 6.3.2 Concre 6.4.1	f Semantics Absolute Semantics Relative Semantics sitionality ave Semantics A Basic Semantic Language: SFOL Semantics dized Semantics An SQL-Inspired Basic Database Language	53 53 54 57 57 57 60
6	6.1 6.2 6.3	Kinds of 6.1.1 6.1.2 Compo Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2	f Semantics Absolute Semantics Relative Semantics sitionality ve Semantics A Basic Semantic Language: SFOL Semantics Aized Semantics An SQL-Inspired Basic Database Language Semantics	53 53 54 57 57 57 60 60 61
6	6.1 6.2 6.3	Kinds of 6.1.1 6.1.2 Compo Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compu	f Semantics Absolute Semantics Relative Semantics sitionality ave Semantics A Basic Semantic Language: SFOL Semantics dized Semantics An SQL-Inspired Basic Database Language Semantics stational Semantics	53 53 54 57 57 60 61 64
6	6.1 6.2 6.3	Kinds of 6.1.1 6.1.2 Compose Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compose 6.5.1	f Semantics Absolute Semantics Relative Semantics sitionality ave Semantics A Basic Semantic Language: SFOL Semantics dized Semantics dized Semantics An SQL-Inspired Basic Database Language Semantics tational Semantics A Scala-Inspired Basic Programming Language	53 53 54 57 57 60 61 64 65
6	6.16.26.36.46.5	Kinds of 6.1.1 6.1.2 Composed Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compus 6.5.1 6.5.2	f Semantics Absolute Semantics Relative Semantics sitionality we Semantics A Basic Semantic Language: SFOL Semantics dized Semantics An SQL-Inspired Basic Database Language Semantics tational Semantics A Scala-Inspired Basic Programming Language Semantics	53 53 54 57 57 60 61 64 65 65
6	6.1 6.2 6.3	Kinds of 6.1.1 6.1.2 Composed Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compus 6.5.1 6.5.2	f Semantics Absolute Semantics Relative Semantics sitionality we Semantics A Basic Semantic Language: SFOL Semantics Aized Semantics An SQL-Inspired Basic Database Language Semantics tational Semantics A Scala-Inspired Basic Programming Language Semantics	53 53 54 57 57 60 61 64 65
	6.1 6.2 6.3 6.4 6.5	Kinds of 6.1.1 6.1.2 Compose Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compu 6.5.1 6.5.2 Narration	f Semantics Absolute Semantics Relative Semantics Sitionality Eve Semantics A Basic Semantic Language: SFOL Semantics Eized Semantics An SQL-Inspired Basic Database Language Semantics Etational Semantics A Scala-Inspired Basic Programming Language Semantics Example Semantics A Scala-Inspired Basic Programming Language Semantics Example Example Semantics Example Exampl	53 53 54 57 57 57 60 61 64 65 65
	6.1 6.2 6.3 6.4 6.5	Kinds of 6.1.1 6.1.2 Compose Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compu 6.5.1 6.5.2 Narration	f Semantics Absolute Semantics Relative Semantics Sitionality A Basic Semantic Language: SFOL Semantics A Basic Semantics An SQL-Inspired Basic Database Language Semantics Atational Semantics A Scala-Inspired Basic Programming Language Semantics A Scala-Inspired Basic Programming Language Semantics Ver Semantics Ver Semantics Ver Semantics Ver Semantics Ver Semantics Ver Semantics	53 53 54 57 57 57 60 61 64 65 65 68
6	6.1 6.2 6.3 6.4 6.5 6.6 Que	Kinds of 6.1.1 6.1.2 Compo Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compu 6.5.1 6.5.2 Narration of the control of	f Semantics Absolute Semantics Relative Semantics Sitionality Eve Semantics A Basic Semantic Language: SFOL Semantics Lized Semantics An SQL-Inspired Basic Database Language Semantics Lational Semantics A Scala-Inspired Basic Programming Language Semantics Eve	53 53 54 57 57 60 61 64 65 68 71 71
	6.1 6.2 6.3 6.4 6.5 6.6 Que 7.1	Kinds of 6.1.1 6.1.2 Composed Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Comput 6.5.1 6.5.2 Narration of the control	f Semantics Absolute Semantics Relative Semantics Relative Semantics Sitionality Eve Semantics A Basic Semantic Language: SFOL Semantics Aized Semantics An SQL-Inspired Basic Database Language Semantics Atational Semantics A Scala-Inspired Basic Programming Language Semantics Eve Semantics Eve Semantics Eve Semantics Eve Semantics Eve Querying	53 53 54 57 57 60 61 64 65 68 71 72
	6.1 6.2 6.3 6.4 6.5 6.6 Que 7.1	Kinds of 6.1.1 6.1.2 Compose Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compus 6.5.1 6.5.2 Narration of the control	f Semantics Absolute Semantics Relative Semantics Relative Semantics Relative Semantics Relative Semantics Sitionality A Basic Semantic Language: SFOL Semantics A Basic Semantic Language: SFOL Semantics An SQL-Inspired Basic Database Language Semantics Catalonal Semantics A Scala-Inspired Basic Programming Language Semantics A Semantics Semantics Semantics The A Semantic The A Se	53 53 54 57 57 60 61 64 65 68 71 72 72
	6.1 6.2 6.3 6.4 6.5 6.6 Que 7.1 7.2	Kinds of 6.1.1 6.1.2 Composed Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Compus 6.5.1 6.5.2 Narration of the control	f Semantics Absolute Semantics Relative Semantics Sitionality Absolute Semantics Sitionality Absolute Semantics A Basic Semantic Language: SFOL Semantics A Basic Semantics An SQL-Inspired Basic Database Language Semantics Atational Semantics A Scala-Inspired Basic Programming Language Semantics A Scala-Inspired Basic Programming Language Semantics A Se	53 53 54 57 57 60 61 64 65 68 71 72 72 72
	6.1 6.2 6.3 6.4 6.5 6.6 Que 7.1	Kinds of 6.1.1 6.1.2 Compose Deduct 6.3.1 6.3.2 Concre 6.4.1 6.4.2 Comput 6.5.1 6.5.2 Narration of the control	f Semantics Absolute Semantics Relative Semantics Relative Semantics sitionality we Semantics A Basic Semantic Language: SFOL Semantics cized Semantics An SQL-Inspired Basic Database Language Semantics cational Semantics A Scala-Inspired Basic Programming Language Semantics we Semantics ia a Semantics we Querying Method Challenges ized Querying	53 53 54 57 57 60 61 64 65 68 71 72 72

CONTENTS 5

3	For	Formal Systems				
	8.1	Forma	d Systems	75		
		8.1.1	Syntax	7!		
		8.1.2	Translation	7!		
		8.1.3	Interpretation	76		
	8.2	Seman	ntics	76		
		8.2.1	General Case	76		
		8.2.2	Relative Semantics by Translation	7		
		8.2.3	Relative Semantics by Interpretation	77		

6 CONTENTS

Chapter 0

Meta-Information

0.1 Remarks

Important stuff that you should read carefully!

State of these notes I constantly work on my lecture notes. Therefore, keep in mind that:

- I am developing these notes in parallel with the lecture they can grow or change throughout the semester.
- These notes are neither a subset nor a superset of the material discussed in the lecture. On the one handle, they may contain more details than mentioned in the lectures. On the other hand, important material such as background, diagrams, and examples may be part of the lecture but not mentioned in these notes.
- Unless mentioned otherwise, all material in these notes is exam-relevant (in addition to all material discussed in the lectures).

Collaboration on these notes I am writing these notes using LaTeX and storing them in a git repository on GitHub at https://github.com/florian-rabe/Teaching. As an experiment in teaching, I am inviting all of you to collaborate on these lecture notes with me. This would require familiarity with LaTeX as well as Git and GitHub—that is not part of this lecture, but it is an essential skill for you. Ask in the lecture if you have difficulty figuring it out on your own.

By forking and by submitting pull requests for this repository, you can suggest changes to these notes. For example, you are encouraged to:

- Fix typos and other errors.
- Add examples and diagrams that I develop on the board during lectures.
- Add solutions for the homeworks if I did not provide any (of course, I will only integrate solutions after the
 deadline).
- Add additional examples, exercises, or explanations that you came up or found in other sources. If you use material from other sources (e.g., by copying an diagram from some website), make sure that you have the license to use it and that you acknowledge sources appropriately!

I will review and approve or reject the changes. If you make substantial contributions, I will list you as a contributor (i.e., something you can put in your CV).

Any improvement you make will not only help your fellow students, it will also increase your own understanding of the material. Make sure your git commits carry a user name that I can connect to you.)

Other Advice I maintain a list of useful advice for students at https://github.com/florian-rabe/Teaching/blob/master/general/advice_for_students.pdf. It is mostly targeted at older students who work in individual projects with me (e.g., students who work on their BSc thesis). But much of it is useful for you already now or will become useful soon. So have a look.

0.2 Overview

Structure The subsequent *parts* of this course follow the Tetrapod model with one part per aspect. Each of these will describe the concepts, languages, and tools of the respective aspect as well as their relation to other aspects.

The aspects of the Tetrapod are typically handled in individual courses, which describe highly specialized languages and tools in depth. On the contrary, the overall goal of this course will be seeing all of them as different approaches to semantics and knowledge representation. The course will focus on universal principles and their commonalities and differences as well as their advantages and disadvantages.

The subsequent *chapters* of this first part will be dedicated to aspect-independent material. These will not necessarily be taught in the order in which they appear in these notes. Instead, some of them will be discussed in connection to how they are relevant in individual aspects.

Exercises and Running Example Typical practical projects, e.g., the ones that a strong CS graduate might be put in charge of, involve heterogeneous data and knowledge that must be managed using a variety of optimized aspect-specific languages and tools. Interoperability between these is often a major source of inefficiency and bugs. The exercises accompanying the course will mimic this situation: they will be designed around a single large project that requires choosing and integrating methods, languages, and tools from all aspects.

Concretely, this project will be the development of a univis-like system for a university. It will involve heterogeneous data such as course and program descriptions, legal texts, websites, grade tables, and transcript generation code.

Over the course of the semester students will implement a completely functional system applying the lessons of the course. This is very unusual and often impossible for other courses: as any university course must teach many different things from a wide area, it is rarely possible to find a project that requires many and only lessons from a single course. Here KRP is special because its material pervades all aspects of system development.

Chapter 1

Fundamental Concepts

1.1 Abbreviations

knowledge representation and processing KRP the general area of this course knowledge representation language KRL a languages used in KRP

knowledge representation tool KRT a tool implementing a KPL and processing algorithms for it

1.2 Motivation

1.2.1 Knowledge

Human knowledge pervades all sciences including computer science, mathematics, natural sciences and engineering. That is not surprising: "science" is derived from the Latin word "scire" meaning "to know". Similarly, philosophy, from which all sciences derive, is named after the Greek words "philo" meaning loving and "sophia" meaning wisdom, and the for common ending "-logy" is derived from Greek "logos" meaning word (i.e., a representation of knowledge).

In regards to knowledge, computer science is special in two ways: Firstly, many branches of computer science need to understand KRP as a prerequisite for teaching computers to do knowledge-based tasks. In some sense, KRP is the foundation and ultimate goal of all artificial intelligence. Secondly, modern information technology enables all sciences to apply computer-based KRP in order to vastly expand on the domain-specific tasks that can be automated. Currently all sciences are becoming more and more computerized, but most non-CS scientists (and many computer scientists for that matter) lack a systematic education and understanding of IT-KRP. That often leads to bad solutions when domain experts cannot see which KRP solutions are applicable or how to apply them.

1.2.2 Representation and Processing

It is no coincidence that this course uses the phrase "Representation and Processing". In fact, this is an instance of a universal duality. Consider the following table of analogous concept pairs, which could be extended with many more examples:

Representation	Processing
Static	Dynamic
Situation	Change
Be	Become
Data Structures	Algorithms
Set	Function
State	Transition
Space	Time

¹Indeed, a major problem with the currently very successful machine learning-based AI technology is that it remains unclear when and how it does KRP. That can be dangerous because it leads to AI systems recommending decisions without being able to explain why that decision should be trusted.

Again and again, we distinguish a static concept that describes/represents what is a situation/state is and a dynamic concept that describes how it changes. If that change is a computer doing something with or acting on that representation, we speak of "processing".

It is particular illuminating to contrast KRP to the standard CS course on Data Structures and Algorithms (DA).² Generally speaking, DA teaches the methods, and KRP teaches how to apply them. Data structures are a critical prerequisite for representing knowledge. But data structures alone do not capture what the data means (i.e., the knowledge) or if a particular representation makes any sense. Similarly, algorithms are the critical prerequisite for processing knowledge. But while algorithms can be systematically analyzed for efficiency, it is much harder to analyze if an algorithm processes knowledge correctly. The latter requires understanding what the input and output data means.

Capturing knowledge in computers is much harder than developing data structures and algorithms. It is ultimately the same challenge as figuring out if a computer system is working correctly — a problem that is well-known to be undecidable in general and very difficult in each individual case.

1.3 Components of Knowledge

1.3.1 Syntax and Semantics, Data and Knowledge

Four concepts are of particular relevance to understanding knowledge. They form a 2×2 -quadruple of concepts:

All four concepts are primitive, i.e., they cannot be defined in simpler terms. All sciences have few carefully-chosen primitives on which everything builds. This is done most systematically in mathematics (where primitives include set or function). While mathematical primitives as well as some primitives in physics or CS are specified formally, the above four concepts can only be described informally, ultimately appealing to pre-existing human understanding. Moreover, this description is not standardized — different courses may use very different descriptions even they ultimately try to capture the same elusive ideas.

Data (in the narrow sense of computer science) is any object that can be stored in a computer, typically combined with the ability to input/output, transfer, and change the object. This includes bits, strings, numbers, files, etc.

Data by itself is useless because we would have no idea what to do with it. For example, the object O = ((49.5739143, 11.0264941), "2020 - 04 - 21T16 : 15 : 00CEST") is useless data without additional information about its syntax and semantics. Similarly, a file is useless data unless we know which file format it uses.

Syntax is a system of rules that describes which data is **well-formed**. For *O* above the syntax could be "a pair of (a pair of two IEEE double precision floating point numbers) and a string encoding of an time stamp". For a file, the syntax is often indicated by the file name extension, e.g., the syntax of an html file is given in Section 12 of the current HTML standard³.

Syntax alone is useless unless we know what the semantics, i.e., what the data means and thus how to correctly interpret and process the data. For example, the syntax of O allows to check that O is well-formed, i.e., indeed contains two numbers and a timestamp string. That allows rejecting ill-formed data such as ((49.5739143, 11.0264941), "foo"). The HTML syntax allows us to check that a file conforms to the standard.

Semantics is a system of rules that determines the meaning of well-formed data. For example, ISO 8601 specifies that timestamp string refer to a particular date and time in a particular time zone. Further semantics for O might be implicit in the algorithms that produce and consume it: such as "the first component of the pair contains two numbers between 0 and 180 resp. 0 and 360 indicating latitude resp. longitude of a location on earth". Semantics might be multi-staged, and further semantics about O might be that O indicates the location and time of the first lecture of this course. Similarly, Section 14 of the HTML standard specifies the semantics of well-formed HTML files by describing how they are to be rendered in a web browser.

Knowledge is the combining of some data with its syntax and semantics. That allows applying the semantics to obtain the meaning of the data (if syntactically well-formed and signaling an error otherwise). In computer systems,

²The course is typically called "Algorithms and Data Structures", but that is arguably awkward because algorithms can only exist if there are data structure to work with. Compare my notes on that course in this repository, where I emphasize data structures much more than is commonly done in that course.

 $^{^3}$ https://html.spec.whatwg.org/multipage/

- data is represented using primitive data (ultimately the bits provided by the hardware) and encodings of more complex data (bytes, arrays, strings, etc.) in terms of simpler ones,
- syntax is theoretically specified using grammars and practically implemented in programming languages using data structures,
- semantics is represented using algorithms that process syntactically well-formed data,
- knowledge is elusive and often emerges from executing the semantics, e.g., rendering of an HTML file.

1.3.2 Semantics as Syntax Transformation

In order to capture knowledge better in computer systems, we often use two syntax levels: one to represent the data itself and another to represent the knowledge. These can be seen as input and output data. In that case, semantics is a function that translates from the data syntax to the knowledge syntax, and knowledge is the pair of the data and the result of applying the semantics. The following table gives some examples.

Data syntax	Semantics function	Knowledge syntax
SPARQL query	evaluation	result set
SQL query	evaluation	result table
program	compiler	binary code
program expression	interpreter	result value
logical formula	interpretation in a model	mathematical object
HTML document	rendering	graphical representation

Thus, the role of syntax vs. semantics may depend on the context: just like one function's output can be another function's input, one interpretation's knowledge can be another one's syntax. For example, we can first compile a program into binary and then execute it to returns its value.

Such hierarchies of evaluation levels are very common in computer systems. In fact, most state-of-the-art compilers are subdivided into multiple phases each further interpreting the output of the previous one. Thus, if knowledge is represented in computers, it is invariably data itself but relative to a different syntax.

1.3.3 Heterogeneity of Semantics and Knowledge

While it is easy to design languages to represent data in general, it is very difficult to designing KRLs that capture the human-level quality of knowledge. Over the last few decades, the KRP area in computer science has diversified into different subareas that approach this research problem in fundamentally different ways. In fact, KRP in the very general sense of this course is usually not even studied by itself — instead the subareas are so different, specialized, and large that they all sustain their respective university courses and research conferences.

This is related to the fact that the data naturally comes in fundamentally different forms such as graphs, arrays, tables in the sense of relational databases, programs in a programming language, logical formulas, or natural language texts. We speak of **heterogeneous** data. These different forms of data are supported by highly specialized KPTs: graph databases, array databases, relational databases, package databases for programming languages, theorem databases for logics (e.g., the Isabelle Archive of Formal Proofs), databases of research papers (such as the arXiv), and so on.

All of these are very successful for their respective kind of data. And all of them include specifications of semantics and KP algorithms that implement this semantics. But it can vary massively how the semantics is specified and implemented. This has caused major practical problems for tool interoperability: many projects require data in multiple formats and algorithms from multiple tools. But the respective tools are often islands that assume that all data is represented in the tool's language and users do not use outside tools. Therefore, the import/export capabilities of the tools are often limited.

Moreover, transporting data across systems is usually ignorant of the semantics: while each tool takes relatively good care to implement the semantics correctly, there is much less certainty that the semantics is preserved when exchanging data across tools. For a trivial example, consider a tool that measures length in inches vs. a tool that uses centimeters, both using floating point numbers for the data: if they exchange the data, i.e., just the numbers, they may miscommunicate the semantics.⁴

This problem is not easy to fix though. The heterogeneity of data and semantics is so extreme that it is, in some cases, an open theoretical problem how knowledge can be shared at all across tools. The basic idea — exchange the

 $^{^4}$ Problems like this have been involved in major disasters such as the Mars Climate Orbiter.

data in a way that preserves semantics — can be difficult to implement if both tools use entirely different paradigms to specify semantics.

1.4 The Tetrapod Model of Knowledge

The Tetrapod model of knowledge is an ongoing research project by the instructors of this course. A first publication was made in [CFKR21]. The structure of this course will draw heavily on the Tetrapod model to get an overview of the different approaches to KPR and their interoperability problems.

1.4.1 Five Aspects of Knowledge

The Tetrapod model distinguishes five basic **aspects** of knowledge and KPR as described below. For each aspect, there is a variety dedicated KRLs supported by highly optimized KPTs as indicated in the following table:

Aspect	KRLs (examples)	KPTs (examples)
ontologization	ontology languages (OWL), description logics (ALC)	reasoners, SPARQL engines (Virtuoso)
concretization	relational databases (SQL, JSON)	databases (MySQL, MongoDb)
computation	programming languages (C)	interpreters, compilers (gcc)
deduction	logics (HOL)	theorem provers (Isabelle)
narration	document languages (HTML, LaTeX)	editors, viewers

Ontologization focuses on developing and curating a coherent and comprehensive ontology of concepts. This focuses on identifying the central concepts in a domain and their relations. For example, a medical ontology would define concepts for every symptom, disease, and medication and then define relations for which symptoms and medications are related to which disease.

Ontologies typically abstract from the knowledge: they standardize identifiers for the concepts and spell out some properties and relations but do not try to capture all details of the knowledge. Well-designed ontologies can capture exactly the knowledge that different KPTs must share and can thus serve as interoperability layers between them. While organization can use ontology languages such as OWL or RDF, the inherent complexity of formal objects in computer science and mathematics usually requires going beyond general purpose ontology languages (similar to how the programming languages underlying computer algebra systems usually go beyond general purpose programming languages).

Concretization uses languages based on numbers, strings, lists, and records to obtain concrete representations of datasets in order to store and query their properties efficiently. Because concrete objects are so simple and widely used, it is possible and common to build concrete datasets on top of general purpose data representation languages and tools such as JSON or SQL.

Computation uses specification and programming languages to represent algorithmic knowledge.

Deduction uses logics and theorem provers to obtain verifiable correctness.

Narration uses natural language to obtain texts that are easy to understand for humans. Because narrative languages are not well-standardized (apart from general purpose languages such as free text or IATEX), it is common to develop narrative libraries on top of ad-hoc languages that impose some formal structure on top of informal text, such as a fixed tree structure whose leafs are free text or a particular set of IATEX macros that must be used. Narrative libraries can be classified based on whether entries are derived from publications (e.g., one abstract per paper in zbMATH) or mathematical concepts (e.g., one page per concept in nLab).

1.4.2 Relations between the Aspects

The aspects can be visualized as the corners of tetrahedron with ontologization in the center and edges and faces representing solutions that mix two or three aspects as seen in Figure 1.1.

Most approaches try to incorporate all or multiple aspects. But all languages and tools tend to be heavily biased towards and optimized for a single one of the four corner aspects. This is not due to ignorance but because each aspect provides characteristic advantages that are extremely hard to capture at once. In fact, every combination of aspects shares characteristic advantages and disadvantages as sketched in Figure 1.2. For example, deductive and narrative definitions of a function involved well-definedness arguments, and a function defined by a concrete table



Figure 1.1: Tetrapod model of knowledge

				C	haracteris	tic	
Aspect	object	s	ad	vantage	joint ad	vantage	application
					of the o	ther as-	
					pects		
deduction	formal	proofs	co	rrectness	ease of u	ise	verification
computation	progra	$_{ m ms}$	eff	iciency	well-defi	nedness	execution
concretization	concre	te objects	taı	ngibility	abstract	ion	storage/retrieval
narration	texts		fle	xibility	formal	seman-	human understanding
					tics		
		Aspect pa	ir	characte	ristic adv	antage	
		ded./com	р.	rich met	a-theory		
	narr./conc		c. simple languages				
		ded./narr	ded./narr. theorem		s and pro-	ofs	
		comp./com	ıc.	normaliz	zation		
		ded./conc		decidable well-definedness			
		comp./na	rr.	Turing of	completen	ess	

Figure 1.2: Shared properties and advantages of aspects

is trivially well-defined, but a computational definition of a function may throw exceptions when running; but only the latter can store and compute functions efficiently. Consequently, dedicated and mostly disjoint communities have evolved that have produced large aspect-specific datasets.

Chapter 2

Ontologies

2.1 General Principles

Motivation An ontology is an abstract representation of the main concepts in some domain. Here *domain* refers to any area of the real world such as mathematics, biology, diseases and medications, human relationships, etc. Many examples can be found at https://bioportal.bioontology.org/, including the Gene ontology one of the biggest.

Contrary to the other four aspects, ontological knowledge representations do not aim at capturing the entire semantics of the domain objects. Instead, they focus on defining unique identifiers for those objects and describing some of their properties and relations to each other.

We use the word **ontologization** to refer to the process of organizing the knowledge of a domain in ontologies.

Ontologies are most valuable when they are *standardized* (either sanctioned through a formal body or a quasi-standard because everyone uses it). A standard ontology allows everybody in the domain to use the identifiers defined by the ontology in a way that avoids misunderstandings. Thus, in the simplest form, an ontology can be seen as a dictionary defining the technical terms of a domain. For example, the Gene ontology defines identifier GO:000001 to have the formal name "mitochondrion inheritance" and the informal definition "The distribution of mitochondria, including the mitochondrial genome, into daughter cells after mitosis or meiosis, mediated by interactions between mitochondria and the cytoskeleton.".

Ontology Languages An ontology is written in ontology language. Common ontology languages are

- description logics such as ALC,
- the W3C ontology language OWL, which is the standard ontology languages of the semantic web,
- the entity-relationship model, which focuses on modeling rather than formal syntax,
- modeling languages like UML, which is the main ontology language used in software engineering.

Ontology languages are not committed to a particular domain — in the Tetrapod model, they correspond to programming languages and logics, which are similarly uncommitted. Instead, an ontology language is a formal language that standardizes the syntax of how ontologies can be written as well as their semantics.

Ontologies The details of the syntax vary between ontology languages. But as a general rule, every ontology declares

- individual concrete objects that exist in the real world, e.g., "Florian Rabe" or "WuV"
- concept abstract groups of individuals, e.g., "instructor" or "course"
- relation binary relations between two individuals, e.g., "teach"
- **properties** binary relations between an individuals and a concrete value (a number, a date, etc.), e.g., "creditValue"
- concept assertions the statement that a particular individual is an instance of a particular concept
- relation assertions the statement that a particular relation holds about two individuals
- **property assertions** the statement that a particular individual has a particular value for a particular property
- terminological axioms statements about relations between concepts, typically in the form subconcept

of statements like "instructor" ⊑ "person"

All assertions can be understood and spoken as subject-predicate-object triples as follows:

Assertion		Triple	
	Subject	Predicate	Object
concept assertion	"Florian Rabe"	is-a	"instructor"
relation assertion	"Florian Rabe"	"teach"	"WuV"
property assertion	"WuV"	"credit Value"	7.5

This uses a special relation is-a between individuals and concepts. Some languages group is-a with the other binary relations between individuals for simplicity although it is technically a little different.

The possible values of properties must be fixed by the ontology language. Typically, it includes at least standard types such as integers, floating point numbers, and strings. But arbitrary extensions are possible such as dates, RGB-colors, lists, etc. In advanced languages, it is possible that the ontology even introduces its own basic types and values.

Ontologies are often divided into two parts:

- The abstract part contains everything that holds in general independent of which individuals: concepts, relations, properties, and terminological axioms. It describes the general rules how the worlds works without committing to a particular set of inhabitants of the world. This part is commonly called the **TBox** (T for terminological).
- The **concrete** part contains everything that depends on the choice of individuals: individuals and assertion axioms. It populates the world with inhabitants. This part is commonly called the **ABox** (A for assertional).

A separate division into two parts is the following:

- The **signature** part contains everything that introduces a **named entity**: individuals, concepts, relations, and properties.
- The **theory** part contains everything that describes which statements about the named entities are true: assertions and axioms.

Synonyms Because these principles pervade all formal languages, many competing synonyms are used in different domains. Common synonyms are:

Here	OWL	Description logics	ER model	UML	semantics via logics
individual	instance	individual	entity	object, instance	constant
concept	class	concept	entity-type	class	unary predicate
relation	object property	role	role	association	binary predicate
property	data property	(not common)	attribute	field of base type	binary predicate

In particular, the individual-concept relation occurs everywhere and is known under many names:

domain	individual	concept
type theory, logic	constant, term	type
set theory	element	set
database	row	table
$philosophy^1$	object	property
grammar	proper noun	common noun

2.2 A Basic Ontology Language

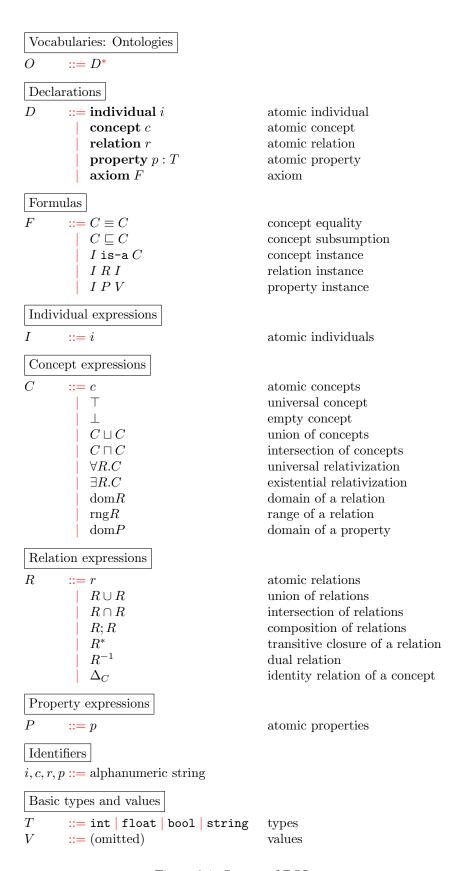


Figure 2.1: Syntax of BOL

We could study practical ontology languages like ALC or OWL now. But those feature a lot of other details that can block the view onto the essential parts. Therefore, we first define a basic ontology language ourselves in order to have full control over the details.

Definition 2.1 (Syntax of BOL). A BOL-ontology is given by the grammar in Fig. 2.1. It is well-formed if

- no identifier is declared twice,
- every property assertion assigns a value of the type required by the property declaration,
- every reference to an atomic individual/concept/relation/property is declared as such.

The above grammar exhibits some general structure that we find throughout formal KR languages. In particular, an ontology consists of **named declarations** of four different kinds of entities as well as some assertions and axioms about them. Each entity declaration clarifies which kind it is (in our case by starting with a keyword) and introduces a new entity identifier. For each kind, there are complex expressions. These are anonymous and built inductively; their base cases are references to the corresponding identifiers. Sometimes (in our case: individuals and properties), the references are the only expressions of the kind. Sometimes (in our case: concepts and relations), there can be many productions for complex expressions. The complex expressions are used to build axioms; in our case, these are the three kinds of assertions and other formulas.

Remark 2.2 (Formulas vs. Assertions). In Fig. 2.1, the three productions in gray are duplicated: they occur both as assertions and as formulas.

We could remove the three productions for assertions and treat them as special cases of axioms. But We keep the duplication here because assertions are often treated differently from the other axioms. They are grouped with the individuals in the ABox whereas the other axioms are seen as part of the TBox. Moreover, when used as assertions, they may have to be interpreted differently than when used as formulas as we will see in Ch. 6.

Alternatively, we could remove the three productions in gray. But then we would lose the ability to talk about formulas that are not true. That will become relevant in Ch. 7.

Axioms The role of the axioms is two-fold:

- They can be used to perform **consequence closure**: a formula may expresses a closure operation that defines assertions that are automatically added to the ontology. That can be difficult as some kind of exhaustive reasoning is needed. For example, if there is a subconcept axiom "instructor" ⊑ "person" and a concept assertion "FlorianRabe" is-a "instructor", we have to add the implied concept assertion "Florian Rabe" is-a "person".
- They can be used to perform **consistency conditions** that must not violated by the ontology. For example, ontologies may contain contradictory assertions or violations of uniqueness constraints such as a person should only have one father or fathers should be male. The axioms exclude such cases. That should succeed if the assertions are already consequence-closed.

But spelling out how that works is part of the semantics, not the syntax.

Example 2.3. We give a simple ontology that could be used to represent knowledge in the context of a university:

```
individual FlorianRabe
individual WuV
concept person
concept male
concept instructor
concept course
relation teach
property creditValue: float
FlorianRabe is—a instructor □ male
WuV is—a course
FlorianRabe teach WuV
WuV creditValue 7.5
male ⊑ person
instructor ⊑ person
dom teach ⊑ instructor
```

```
\begin{array}{ll} {\rm rng} \ {\rm teach} \ \sqsubseteq \ {\rm course} \\ {\rm dom} \ {\rm creditValue} \ \equiv \ {\rm course} \\ {\rm course} \ \sqsubseteq \ \exists \ {\rm teach}^{-1} \ {\rm instructor} \end{array}
```

The axioms are meant to state that males and instructors are persons, teaching is done by instructors to courses, exactly the courses have credits, and (the last axiom) every course is taught by at least one instructor. Whether they actually do mean that, depends on the semantics.

The consequence closure (as defined by the semantics) should add the assertion FlorianRabe is-a person. Alternatively, if we use the axioms for consistency checking, we should add that assertion from the beginning. Otherwise, the axioms would not be true.

If we use axioms for the consequence closure, we can even omit the two concept assertions — they should be inferred using the domain and range axioms for the relation.

The assertion FlorianRabe is—a instructor \sqcap male could also be split into two assertions FlorianRabe is—a instructor and FlorianRabe is—a male. That will be important as some semantics might have difficulties handling all cases. So it can be helpful to use a variant that does not need \sqcap operator.

2.3 Representing Ontologies as Triples

It is common to represent an entire ontology as a set of subject-predicate-object triples. That makes handling ontologies very simple and efficient. This is the preferred representation of the semantic web.

However, while, e.g., relation assertions are naturally triples, not all declarations are, and some tricks may be necessary.

Inferring the Entity Declarations The entity declarations are not naturally triples. But we can usually infer them from the assertions as follows: any identifier that occurs in a position where an entity of a certain kind is expected is assumed to be declared as an entity for that kind.

For example, the individuals are what occurs as the subject of a concept, relation, or property assertion or as the object of a relation assertion. It is conceivable that there are individuals that occur in none of these. But that is unusual because they would be disconnected from everything in the ontology.

If we give TBox and ABox together, this inference approach usually works well. But if we only give a TBox, this would often not allow inferring all entities. The only place where they could occur in the TBox is in the axioms, and it is quite possible to have concept, relation, and property declarations that are not used in the axioms. In fact, it is not unusual not to have any axioms.

Special Predicates To turn declarations into triples, we can use reflection, i.e., the process of talking about our language constructs as if they were data.

Reflection requires introducing some built-in entities that represent the features of the language. In the semantic web area, this is performed using the following entities:

- "rdfs:Resource": a built-in concept of which all individuals are an instance and thus of which every concept is a subconcept
- "rdf:type": a special predicate that relates an entity to its type:
 - an individual to its concept (corresponding to is-a above)
 - other entities to their special type (see below)
- "rdfs:Class": a special class to be used as the type of classes
- "rdf:Property": a special class to be used as the type of properties
- "rdfs:subClassOf": a special relation that relates a subconcept to a superconcept
- "rdfs:domain": a special relation that relates a relation to the concepts of its subjects
- "rdfs:range": a special relation that relates a relation/property to the concept/type of its objects

Here "rdf" and "rdfs" refer to the RDF (Resource Description Framework) and RDFS (RDF Schema) namespaces, which correspond to W3C standards defining those special entities.

Thus, we can represent many and in particular the most important entity declarations as triples:

Assertion	Triple			
	Subject	Predicate	Object	
individual	individual	"rdf:type"	"rdfs:Resource"	
concept	concept	"rdf:type"	"rdf:Class"	
relation	relation	"rdf:type"	"rdf:Property"	
property	property	"rdf:type"	"rdf:Property"	
concept assertion	individual	"rdf:type"	concept	
relation assertion	individual	relation	individual	
property assertion	individual	property	value	
for special forms of	axioms			
$c \sqsubseteq d$	c	"rdfs:subClassOf"	d	
$\operatorname{dom} r \equiv c$	$\mid r \mid$	"rdfs:domain"	c	
$\operatorname{rng} r \equiv c$	r	"rdfs:range"	c	

This is subject to the restriction that only atomic concepts and relations can be handled. For example, only concept assertions can be handled that make an individual an instance of an *atomic* concept. This is particularly severe for axioms, where complex expressions occur most commonly in practice. Here, the special relations allow capturing the most common axioms as triples.

Problems Reflection is subtle and can easily lead to inconsistencies. We can see this in how the approach of RDF(S) special entities breaks the semantics via FOL.

For example, it treats classes both as concepts (when they occur as the object of a concept assertion) and as individuals (when they occur as subject or object of a "rdfs:subClassOf" relation assertion). Similarly, "rdfs:Class" is used both as an individual and as a class. In fact, the standard prescribes that "rdfs:Class" is an instance of itself.

In practice, this is handled pragmatically by using ontologies that make sense. A formal way to disentangle this is to assume that there are two variants of "rdfs:Class", one as an individual and one as a class. The translation must then translate "rdfs:Class" differently depending on how it is used.

It would be better if RDFS were described in a way that is consistent under the implicitly intended FOL semantics. But the more pragmatic approach has the advantage of being more flexible. For example, being able to treat every class, relation, or property also as an individual makes it easy to annotate metadata to them. Metadata is a set of properties such as "rdfs:seeAlso" or "owl:versionInfo", whose subjects can be any entity.

Subject-Centered Representations When giving a set of triples, there are usually a lot of triples with the same subject. For example, we could use a simple concrete syntax with one triple per line and whitespace separating subject, predicate, and object:

```
"FlorianRabe" is—a "instructor"
"FlorianRabe" is—a "male"
"FlorianRabe" "teach" "WuV"
"FlorianRabe" "teach" "KRMT"
"FlorianRabe" "age" 40
"FlorianRabe" "office" "11.137"
```

It is more human-friendly to group these triples in such a way that the subject only has to be listed once. For example, we could use a concrete syntax like this, where the subject occurs first and then predicate-object pairs occur on indented lines:

```
"Florian Rabe"
is—a "instructor"
is—a "male"
"teach" "WuV"
"teach" "KRMT"
"age" 40
"office" "11.137"
```

If the same predicate occurs with multiple values, we can group those as well. For example, we could give the objects for the same predicates as a list following the predicate:

```
"Florian Rabe"
is—a "instructor" "male"
"teach" "WuV" "KRMT"
"age" 40
"office" "11.137"
```

Concrete syntaxes based on the triple representation of ontologies will usually adopt some kind of structure like this. The details may vary.

2.4 Writing Ontologies

2.4.1 The OWL Language

Abstract Syntax and Semantics Due to their central in knowledge representation, a number of languages for ontology writing exist. Most importantly, the syntax and semantics of OWL, including several sublanguages, are standardized by the W3C.

OWL includes a number of built-in special entities. Most importantly, "owl:Thing" corresponds to "rdfs:Resource" as the concept of all individuals.

Concrete Syntax Several concrete syntaxes have been defined and are commonly used for OWL. The OWL2 primer² systematically describes examples in five different concrete syntaxes.

APIs for OWL implement the abstract syntax along with good support for reading/writing ontologies in any of the concrete syntaxes.

2.4.2 The Protege Tool

A widely used tool for writing ontologies in OWL is Protege³.

To get started with Protege without getting confused, we need to continue understand how its key terminology maps to other contexts.

Here	Protege	Edited in WebProtege via
individual	individual	listed in "Individuals" tab
concept	class	listed in "Classes" tab
relation	object property	listed in "Properties" tab
property	data property	listed in "Properties" tab
concept assertion	Type	detail area of the individual in "Individuals" tab
relation assertion	Relationship	detail area of the subject in "Individuals" tab
property assertion	Relationship	detail area of the subject in "Individuals" tab

Protege's interface treats some parts of the ontology specially:

- The "Classes" tab organizes concepts using a tree view based on the subconcept relationship. Superclasses of a class can also be edited directed by listing parents.
- The "Properties" tab organizes properties using a tree view based on the subproperty (i.e., implication, subset) relationship.
- Axioms describing the domain and range of a property can be given directly in its details view.

Note that classes can be in relationships with other classes as well even though that was not considered in the course so far.

2.4.3 Exercise 1

The topic of Exercise 1 is to use Protege to write an OWL ontology for a university.

²https://www.w3.org/TR/2012/REC-owl2-primer-20121211/

 $^{^3}$ https://protege.stanford.edu/

Protege is a graphical editor for the abstract syntax of OWL. Familiarize yourself with the various concrete syntaxes of OWL by writing an ontology that uses every feature once, downloading it in all available concrete syntaxes, and comparing those.

The minimal goal of the exercise session is to get a Hello World example going, at which point the task transitions into homework. There will be no homework submission, but you will use your ontology throughout the course.

You should make sure you understand and setup the process in a way that supports you when you revisit and change your ontology many times throughout the semester.

Other than that, the task is deliberately unconstrained to mimic the typical situation at the beginning of a big project, where it is unclear what the ultimate requirements will be.

2.5 Ontology Morphisms

Ontologies can be used to describe the state of the world. Ontology *morphisms* allow relating different worlds (or different versions of a world) to each other.

Definition 2.4 (Ontology Morphism). Given ontologies V, W, an ontology morphism $m : V \to W$ is an expression of the following grammar

```
 \begin{aligned} & \text{Morphisms} \\ & m ::= A^* \\ & \\ & Definitions \\ & D ::= i := I & \text{individual definition} \\ & | c := C & \text{concept definition} \\ & | r := R & \text{relation definition} \\ & | p := P & \text{property definition} \end{aligned}
```

such that for all axioms F in V, we have that $\overline{m}(F)$ is a theorem of W.

m is well-formed if it contains exactly one definition of every V-identifier with a W-expression.

Definition 2.5 (Homomorphic Extension). Given a morphism $m: V \to W$, we define its homomorphic extension \overline{m} as the mapping from V-expressions to W-expressions that replaces

- every individual reference i with I where i := I is in m
- every concept reference c with C where c := C is in m
- every relation reference r with R where r := R is in m
- every property reference p with P where p := P is in m

Example 2.6. Consider the following ontologies that could be used to model sex and gender.

concept man concept woman axiom man \square woman $\equiv \top$ axiom man \square woman $\equiv \bot$ TransFriendly concept cismale concept sexfemale concept trans

CisNormative

concept cis axiom ...

Then we can give multiple morphisms representing different ways to migrate from the cis-normative to the transfriendly terminology: SexMatters

```
\operatorname{man} \mapsto \operatorname{sexmale}
\operatorname{woman} \mapsto \operatorname{sexfemale}
```

Gender Matters

```
\begin{array}{l} \operatorname{man} \; \mapsto (\operatorname{sexmale} \sqcap \operatorname{cis}) \sqcup (\operatorname{sexfemale} \sqcap \operatorname{trans}) \\ \operatorname{woman} \; \mapsto (\operatorname{sexmale} \sqcap \operatorname{cis}) \sqcup (\operatorname{sexfemale} \sqcap \operatorname{trans}) \end{array}
```

In both cases, we need to prove the homomorphic translations of the axioms.

Chapter 3

Syntax

3.1 The 4 Layers of a Formal Language

Languages are designed, implemented, and evaluated in layers. Fig. 3.1 shows the most common design.

- First, we define a grammar and implement it using a set of inductive types. This is called the AST (abstract syntax tree) of the language. A parser translates string into the corresponding syntax trees, i.e., objects of those inductive types. Objects conforming to the grammar are called **well-formed**.
- Most processing of the language is implemented by inductive functions that traverse the AST. In particular, a context-sensitive inference system is used for type-checking. The type-checking functions take well-formed objects and return booleans, to indicate which are **well-typed**. This defines a context-sensitive sublanguage. That is the language we are really interested in.
- Now the semantics can be implemented: it traverses the AST with the precondition that the input in well-typed. If that succeeds, the objects are called **well-defined**. Ideally, our language satisfies the invariant that for all well-typed objects, the semantics exists. But in some languages well-definedness proofs may be necessary.
- Finally humans evaluate the language by using it. They may report design flaws that feed back into the AST design.

At every layer we detect different issues. Fig. 3.3 and 3.2 give some examples.

Layer	Specified by	Implemented by	Possible error	
Context-Free Syntax	grammar	parser	not derivable from grammar	
Context-Sensitive Syntax	inference system	type checker	symbols not used as declared	KRP
Semantics	inference system, in	terpretation, or translation	undefined semantics	
Pragmatics	human preferences	human judgment	not useful	not KRP

Figure 3.1: The 4 Layers of a Language

Layer	Expression	Issue	Explanation
Context-Free Syntax	1/	syntax error	argument missing
Context-Sensitive Syntax	1/"2"	typing error	argument has wrong type
Semantics	1/0	run-time error	undefined semantics
Pragmatics	1/1	code review comment	unnecessarily complex expression

Figure 3.2: Typical Errors by Layer in a Programming Language

3.2 Context-Free Syntax

Abstractly, context-free syntax is specified using grammars. Concretely, it is implemented using inductive types. In the sequel, we will start with the standard definitions and then make a series of variation to each of these definitions until they become equivalent. The intended equivalence is as follows:

26 CHAPTER 3. SYNTAX

Layer	Expression	Issue	Explanation
Context-Free Syntax	$\forall x$	not well-formed	body missing
Context-Sensitive Syntax	$\forall x.P(y)$	not well-typed	y not declared
Semantics	the $x \in \mathbb{N}$ such that $x < 0$	not well-defined	no such x exists
Pragmatics	$\exists x. x \neq x$	inconsistent	no model exists

Figure 3.3: Typical Errors by Layer in Logic

CFG	IDT
non-terminal	type
production	constructor
non-terminal on left of production	return type of constructor
non-terminals on right of production	arguments types of constructor
terminals on right of production	notation of constructor
words derived from non-terminal N	expressions of type N

Remark 3.1 (Chomsky Hierarchy). The Chomsky hierarchy of grammars has been quite influential and is still taught a lot in introductory courses. But it has little relevance for modern language design.

- CH-0, regular grammars: These are equivalent to regular expressions and finite automata. The latter two are usually better ways to work with these languages than a grammar.
- CH-1, context-free grammars: There are what we consider here.
- CH-2, context-sensitive grammars: Context-sensitive languages are very important. But it is more practical to use context-free grammars and then define a context-sensitive subset later in a different way.
- CH-3, unrestricted grammars: These are Turing-complete and are only relevant theoretically.

3.2.1 Context-Free Grammars

We start with the usual definition:

Definition 3.2 (Context-Free Grammar). Given a set Σ of characters (containing the terminal symbols), a **context-free grammar** consists of

- ullet a set N of names called **non-terminal symbols**
- a set of **productions** each consisting of
 - an element of N, called the **left-hand side**
 - a word over $\Sigma \cup N$, called the **right-hand side**

Example 3.3. Let $\Sigma = \{0, 1, +, \cdot, =, \leq\}$. We give a grammar for arithmetic expressions and formulas about them:

$$E := 0 \\ | 1 \\ | E + E \\ | E \cdot E$$

$$F := E \stackrel{.}{=} E \\ | E \leq E$$

Here we use the BNF style of writing grammars, where the productions are grouped by their left-hand side and written with ::= and ||. We have $N = \{E, F\}$.

First, we give a name to each production of a CFG:

Definition 3.4 (Context-Free Grammar with Named Productions). Given a set Σ of characters (containing the terminal symbols), a **context-free grammar** consists of

- ullet a set N of names called non-terminal symbols
- ullet a set of productions each consisting of
 - a name

- an element of N, called the **left-hand side**
- a word over $\Sigma \cup N$, called the **right-hand side**

Example 3.5. The grammar from above with names written to the right of each production

$$\begin{array}{lll} E ::= 0 & \text{zero} \\ & \mid 1 & \text{one} \\ & \mid E + E & \text{sum} \\ & \mid E \cdot E & \text{product} \\ F ::= E \doteq E & \text{equality} \\ & \mid E \leq E & \text{lessOrEqual} \end{array}$$

This is not common BNF anymore.

Then we add base types to the productions:

Definition 3.6 (Context-Free Grammar with Named Productions and Base Types). Given a set Σ of characters (containing the terminal symbols) and a set T of names (containing the base types allowed in productions), a **context-free grammar** consists of

- \bullet a set N of names called non-terminal symbols
- a set of productions each consisting of
 - a name
 - an element of N, called the **left-hand side**
 - a word over $\Sigma \cup T \cup N$, called the **right-hand side**

The intuition behind base types is that we commonly like to delegate some primitive parts of the grammar to be defined elsewhere. A typical example are literals such as numbers $0, 1, 2, \ldots$: We could give regular expression syntax for digit-strings. Instead, it is nicer to just assume we have a set of base types that we can use to insert an infinite set of literals into the grammar.

Example 3.7. Let Nat be the type of natural numbers and let $T = \{Nat\}$. Then we can improve the grammar from above as follows:

```
\begin{array}{lll} E ::= Nat & \text{literal} \\ \mid E + E & \text{sum} \\ \mid E * E & \text{product} \\ F ::= E \doteq E & \text{equality} \\ \mid E \leq E & \text{lessOrEqual} \end{array}
```

3.2.2 Inductive Data Types

We start with the usual definition:

Definition 3.8 (Inductive Data Type). Given a set of names T (containing the types known in the current context), an *inductive data type* consists of

- a name t, called the **type**,
- a set of **constructors** each consisting of
 - a name n
 - a list of elements of $T \cup \{t\}$, called the **argument** types

Example 3.9. Let Nat be the type of natural numbers and $T = \{Nat\}$. We give an inductive type for arithmetic expressions:

$$E = \mathtt{literal} \, \mathbf{of} \, Nat \quad | \quad \mathtt{sum} \, \mathbf{of} \, E * E \quad | \quad \mathtt{product} \, \mathbf{of} \, E * E$$

Here we use ML-style notation for inductive data types, which separates constructors by | and writes them as name of argument-type-product.

28 CHAPTER 3. SYNTAX

First we generalize to mutually inductive types:

Definition 3.10 (Mutually Inductive Data Types). Given a set T of names (containing the types known in the current context), a family of **mutually inductive data type** consists of

- a set N of names, called the **types**,
- a set of *constructors* each consisting of
 - a name
 - an element of N, called the **return type**
 - a list of elements of $N \cup T$, called the **argument** types

Example 3.11. We extend the type definition from above by adding a second type for formulas. Thus, $N = \{E, F\}$.

```
\begin{array}{c|cccc} E = \mathtt{literal}\,\mathbf{of}\,Nat & | & \mathtt{sum}\,\mathbf{of}\,E*E & | & \mathtt{product}\,\mathbf{of}\,E*E \\ F = \mathtt{equality}\,\mathbf{of}\,E*E & | & \mathtt{lessOrEqual}\,\mathbf{of}\,E*E \end{array}
```

Then we add notations to the constructors:

Definition 3.12 (Mutually Inductive Data Types with Notations). Given a set Σ of characters (containing the terminal symbols) and a set T of names (containing the types known in the current context), a family of **mutually inductive data type with notations** consists of

- a set N of names, called the **types**,
- a set of *constructors* each consisting of
 - a name
 - an element of N, called the **return type**
 - a list of elements of $T \cup N$, called the **argument** types
 - a word over the alphabet $\Sigma \cup T \cup N$ containing the argument types in order and only elements from Σ otherwise, called the **notation** of the constructor

The intuition behind notations is that it can get cumbersome to write all constructor applications as Name(arguments). It is more convenient to attach a notation to them such as

Example 3.13. We extend the type definitions from above by adding notations to each constructor. We use the set $\Sigma = \{+, \cdot, =, \leq\}$ as terminals in the notations.

```
E = \texttt{literal} \ \mathbf{of} \ Nat \ \# \ Nat \quad | \quad \mathsf{sum} \ \mathbf{of} \ E * E \ \# E + E \quad | \quad \mathsf{product} \ \mathbf{of} \ E * E \ \# E \cdot E F = \texttt{equality} \ \mathbf{of} \ E * E \ \# E \ \dot{=} \ E \quad | \quad \texttt{lessOrEqual} \ \mathbf{of} \ E * E \ \# E \leq E
```

Here we write the constructors as name of argument-type-product # notation. It is easy to see that this has introduced redundancy: we can infer the argument types from the notation. So we can just drop the argument types:

```
\begin{array}{c|cccc} E = \mathtt{literal} \ \# \ Nat & | & \mathtt{sum} \ \# \ E + E & | & \mathtt{product} \ \# \ E \cdot E \\ F = \mathtt{equality} \ \# \ E \doteq E & | & \mathtt{lessOrEqual} \ \# \ E \leq E \end{array}
```

3.2.3 Merged Definition

With the variation from above we have arrived at the following equivalence:

Theorem 3.14. Given a set Σ of characters and a set T of names, the following notions are equivalent:

- a family of mutually inductive data types in the context of types T with notations using characters from Σ ,
- a context-free grammar with named productions, terminal symbols from Σ , and base types T.

Proof. The key idea is that

• the types and constructors of the former correspond to the non-terminals and productions of the latter

- for each constructor-production pair
 - the right-hand side of the latter corresponds to the notation of the former,
 - the argument types of the former correspond to the non-terminals occurring on the right-hand side of the latter.

In implementations in programming languages, we often drop the notations. Instead, those are handled, if needed, by special parsing and serialization functions.

However, in an implementation, it is often helpful to additionally give names to each argument of a production/constructor. That yields the following definition:

Definition 3.15 (Context-Free Syntax). Given a set Σ of characters and a set T of names, a context-free syntax consists of

- a set N of names, called the **non-terminals/types**,
- a set of **productions/constructors** each consisting of
 - a name
 - an element of N, called the **left-hand side/return type**
 - a sequence of objects, called the **right-hand side/arguments** which are one of the following
 - * an element of Σ
 - * a pair written (n:t) of a name n, called the **argument name**, and an element $t \in T \cup N$ called the **argument type**.

Example 3.16. Using ad hoc language to write the constructors, our example from above as a context-free syntax could look as follows:

```
E = \texttt{literal} \# (value : Nat) \quad | \quad \texttt{sum} \# (left : E) + (right : E) \quad | \quad \texttt{product} \# (left : E) \cdot (right : E) \\ F = \texttt{equality} \# (left : E) \doteq (right : E) \quad | \quad \texttt{lessOrEqual} \# (left : E) \leq (right : E)
```

3.2.4 Contexts

We assume a context-free language l.

Definition 3.17 (Context). A context is a list of

- grammar terminology: productions N := x
- type terminology: declarations x:N

where each x is a unique name and each N is non-terminal symbol.

The x are called **variables**.

Remark 3.18. Sometimes the grammar itself has specific productions for contexts and variables. In that case, we speak of meta-variable contexts and meta-variables to distinguish them from those of the language.

Definition 3.19 (Expressions in Context). Given a context Γ , a word E derived from non-terminal N that may additionally use the productions of the context is called an *expression of type* N *in context* Γ .

We write this as $\Gamma \vdash_l E : N$.

Definition 3.20 (Substitution). Given two contexts $\vdash_l \Gamma$ and $\vdash_l \Delta$, a **substitution** γ from $\Gamma = x_1 : N_1, \ldots, x_n : N_n$ to Δ is a list $x_1 := e_1, \ldots, x_n := e_n$ where every e_i is an expression of type N_i in context Δ (i.e., $\Delta \vdash_l e_i : N_i$). We write this as $\Delta \vdash_l \gamma : \Gamma$ or as $\vdash_l \gamma : \Gamma \to \Delta$.

Definition 3.21 (Substitution Application). Given an expression $\Gamma \vdash_l E : N$ and a substitution $\vdash_l \gamma : \Gamma \to \Delta$ where $\Gamma = x_1 : N_1, \ldots, x_n : N_n$ and $\gamma = x_1 := e_1, \ldots, x_n := e_n$, we write $E[\gamma]$ for the result of replacing every x_i in E with e_i .

30 CHAPTER 3. SYNTAX

```
Theorem 3.22. If \Gamma \vdash_l E : N \text{ and } \vdash_l \gamma : \Gamma \to \Delta, \text{ then } \Delta \vdash_l E[\gamma] : N.
```

We often want to substitute only a single variable x:N even though E may be defined in a larger context Γ . This is often written E[x:=e]. That is just an abbreviation for $E[\gamma]$, where γ contains x:=e as well as y:=y for every other variable y of Γ .

Notation 3.23. There are many different notations for substitutions that are in common use. Examples include E(x), E[x := e], E[x/e], or [e/x]E. Usually authors of a paper or textbook briefly mention their preferred notation briefly at the beginning.

3.3 Context-Sensitive Syntax

3.3.1 Principles

It is common to define a language as the set of words that can be produced from the syntax, i.e., from a distinguished non-terminal (the start symbol) of the context-free grammar. It is common to define a context-sensitive language as a special case: the set of words that can be produced from a context-sensitive grammar.

This is, however, not helpful in practice. While the above remains the official definition of what the context-sensitive languages are, all practical definitions are entirely different. In fact, context-sensitive grammars are virtually never used to define a specific language. Instead, more restrictive definitions are used that capture more properties of practical languages.

A typical context-sensitive syntax includes the following:

- a context-free syntax,
- a distinguished non-terminal symbol \mathcal{V} , whose words are called **vocabularies**,
- a set of distinguished non-terminal symbols \mathcal{E} , whose words are called \mathcal{E} -expressions,
- a unary predicate $wfv(\Theta)$ on vocabularies Θ ,
- for every vocabulary Θ and every \mathcal{E} , a unary predicate wff $\Theta(E)$ on \mathcal{E} -expressions E.

In case of wfv(Θ), we call Θ well-formed. In case of wff Θ (E), we call E a well-formed \mathcal{E} -expression over Θ .

In practice, it is most useful to think of a context-sensitive language as a family of context-sensitive languages: for every vocabulary, one language containing the expressions.

Remark 3.24 (Terminology). "Well-formed \mathcal{E} -expression over Θ " can be a mouthful. Therefore, it is common to simply say that E is an \mathcal{E} -expression, or that E is a Θ -expression, and expect readers to fill in the details.

It is also common to give the non-terminal \mathcal{E} names, such as "term", "type", "formula", or (confusingly) "expression". Then we simply say "term" instead of "term-expression" and so on.

Example 3.25. We extend the context-free language from Ex. 3.3 to a context-sensitive one as follows:

```
Vocabularies
Voc ::= Decl^*
                                                list of declarations
Declarations
Decl ::= id : Type^* \to Type
                                                typed function symbols
       id: Tupe^* \rightarrow Form
                                                typed predicate symbols
Type ::= Nat \mid String
                                                base types
Expressions
Expr := 0 \mid 1 \mid Expr + Expr \mid Expr * Expr
                                                as before
        id(Expr^*)
                                                application of a function symbol
Formulas
Form := Expr = Expr \mid Expr < Expr
                                                as before
         id(Expr^*)
                                                application of a predicate symbol
```

Voc is the special non-terminal V. The special non-terminals for expressions are Type, Expr, and Form.

The well-formedness predicates check that every identifier is used according to its declaration.

3.3.2 Vocabularies

The vocabularies are typically lists of typically named declarations. They introduce the names that can be used to form expressions. The expression kinds almost always include formulas.

Often declarations contain additional expressions, most importantly types or definitions. In general, all expressions may occur in declarations, but many language systems do not use all of them.

Very different names are used for the vocabularies in different communities. The following table gives an overview:

Aspect	vocabulary Θ	expression kinds \mathcal{E}
Ontologization	ontology	individual, concept, relation, property, formula
Concretization	database schema	cell, row, table, formula
Computation	program	term, type, object, class,
Logic	theory	term, type, formula,
Narration	dictionary	phrases, sentences, texts

The standard library is a fixed vocabulary that is always present. Shifting operators from the grammar into the standard library is a common technique to simplify the language.

Example 3.26. We can rewrite Ex. 3.25 into a much shorter grammar, in which all name-like productions are declared in the standard library:

```
\begin{array}{lll} Voc & ::= Decl^* & \text{list of declarations} \\ Decl & ::= id : Type^* \to Type & \text{typed function symbols} \\ & \mid id : Type^* \to FORM & \text{typed predicate symbols} \\ & \mid id : TYPE & \text{type symbols} \\ Type & ::= id & \text{reference to a type symbol} \\ Expr & ::= id(Expr^*) & \text{application of a function symbol} \\ Form & ::= id(Expr^*) & \text{application of a predicate symbol} \\ \end{array}
```

The standard library is the vocabulary containing the following declarations:

- type symbols: Nat: TYPE, String: TYPE
- function symbols: $0: Nat, 1: Nat, sum: Nat Nat \rightarrow Nat, product: Nat Nat \rightarrow Nat,$
- predicate symbols: $equals: Nat\ Nat \rightarrow FORM, \ lesseq: Nat\ Nat \rightarrow FORM$

3.3.3 Contexts and Variable Binding

In addition to the vocabulary, there is often a second kind of declarations that is part of the context: variables. Variables are usually introduced by *binders*. These are special expression productions that declare a variable and take another expressions, which is the scope of the variable. Somewhat confusingly, the list of variables that are currently in scope is usually (also) called the *context* even though the the overall context consists of the vocabulary and the variable context.

Typical binders are

- Quantifiers like \forall and \exists .
- The let binder.
- Description/choice operators, which return an object with certain properties.
- The block-sequence operator from programming languages often written as a ;-separated sequence of statements wrapped in {}. Here each statement may be a variable declaration, in which case the remainder of the current block is its scope.

If a type system is used, variable declarations V are often typed and written as x : A. Otherwise, they consist only of the variable name x.

32 CHAPTER 3. SYNTAX

Example 3.27. We can extend Ex. 3.26 with expression productions that bind variables:

```
Voc
          ::= Decl^*
                                                list of declarations
Decl
         ::= id : Type^* \to Type
                                                typed function symbols
            id: Type^* \to FORM
                                                typed predicate symbols
            id: TYPE
                                                type symbols
VarDecl := id : Type
                                                typed variables
Context ::= VarDecl^*
                                                contexts
Tupe
         := id
                                                reference to a type symbol
Expr
          := id(Expr^*)
                                                application of a function symbol
             id
                                                reference to a variable
             the VarDecl such that Form
                                                description
             \mathbf{let} \, VarDecl = Expr \, \mathbf{in} \, Expr
                                                let binding
         ::= id(Expr^*)
                                                application of a predicate symbol
Form
             \forall \mathit{VarDecl.Form}
                                                universal quantification
                                                existential quantification
```

For example, in $\forall x : a.F$, the variable x may occur in the formula F. In $let x : int = E_1 in E_2$, the variable x may occur in E_2 but not in E_1 . This must be enforced by the type checker, which must move the variable declaration into the context when traversing into the respective subexpression.

3.3.4 Vocabularies vs. Contexts

The context is often split into two components:

- the vocabulary containing the symbol declarations,
- the context containing the variable declarations.

They are summarized in Fig. 3.4.

Terminology 3.28. The word context is overloaded. It is used both as a generic term for

- any data structure that is carried through a syntax traversal and
- for the special case of such a data structure that collects variable declarations.

Both vocabularies and contexts store declarations that are used later on. Thus the overall context is typically a pair of a vocabulary and a variable context. During context-sensitive traversal, any declaration that is found is remembered in the overall context. Identifier references are only legal if the referenced identifier is in the Context with the appropriate arity, type, etc.

The vocabulary stores all toplevel declarations of the current development. These are often called *symbols* or *constants*. It is often split over multiple files and packages using imports. Users often package up a vocabulary in order to share it with other users. Consequently, it is usually a *set* of declarations and declarations are identified via unique names. To make declaration names from different packages unique, namespaces are used to form qualified names — the namespace is a name for the entire vocabulary, and declarations are referred to by the pair of vocabulary and declaration name. Because the declarations in a vocabulary can be referenced by anybody using that vocabulary, the names cannot be changed easily. Because the vocabulary can become big, it is often implemented as a hash table mapping qualified names to declarations. When checking a vocabulary, each declaration is checked relative to the vocabulary of preceding declarations. (In languages that allow mutual recursion, it may even be that each declaration is checked relative to the entire vocabulary.)

The variable context stores all local *variable* declarations encountered while traversing a single declaration or expression. The scope of each declaration is limited to the remainder of the declaration or subexpression. Consequently, variables can usually be renamed easily without changing the semantics. When checking a declaration, the context is initially empty and grows as variable declarations are traversed. In particular, each expression is checked relative to a vocabulary and a context.

The vocabulary is almost always more expressive than the context: For every a variable declaration, there is a corresponding symbol declaration but not the other way round. Fig. 3.5 gives some examples.

3.4. IMPLEMENTATION 33

	vocabulary	context	
Structure	list of declarations		
Purpose	declare identifiers to be used later on		
Traversal	accumulated during traversal		
Content	previous declarations	variables traversed in current expression	
Needed for	processing declarations and expressions	processing expressions	
Distribution	split over files, packages	local to current declaration	
Modularity	imports, package manager	only local declarations	
Visibility of declarations	entire project	local scope	
Order	usually inessential	essential	
Renaming	change management problem	no change of semantics (α -renaming)	
Duplicate names	disambiguated through namespaces	shadowing	
Size	often large	typically small	
Typical implementation	hash table	list	

Figure 3.4: Commonalities and Differences between the two kinds of Context

	vocabulary	context
	concepts, etc.	none
SFOL	type, function, predicate	only nullary functions
SQL	table, column, row	none
	class, type, function, value	type (restricted), value

Figure 3.5: Kind of declarations in some Languages

3.4 Implementation

Data structures for context-free syntax can be implemented systematically in all programming languages. But, depending on the style of the language, they make drastically different. Below we give the two most important paradigms as examples.

In each case, context-sensitive syntax is implemented as a set of traversal functions over these data structures that return the well-formedness property, i.e., a boolean.

3.4.1 Functional Programming Languages

In a function programming language, inductive data types are a primitive feature. However, notations and named arguments are not available. So helper functions must be used.

The basic recipe is as follows:

- The types and constructors (without the notations and named arguments) are implemented as family of mutually inductive data types.
- For each argument of each constructor, a partial projective function is defined.
- A set of mutually recursive string rendering functions are defined, one for each constructor, that implement the notations.

Example 3.29. We define our example syntax in ML.

First the inductive types (assuming a type Nat already exists in the context):

```
\begin{array}{lll} \mathbf{data}\,E = \mathtt{literal}\,\mathbf{of}\,Nat & | & \mathtt{sum}\,\mathbf{of}\,E*E & | & \mathtt{product}\,\mathbf{of}\,E*E \\ \mathbf{and}\,F = \mathtt{equality}\,\mathbf{of}\,E*E & | & \mathtt{less0rEqual}\,\mathbf{of}\,E*E \end{array}
```

Now the projection functions:

```
\begin{array}{ll} \mathbf{fun} \ \mathtt{literal\_value}(\mathtt{literal}(v)) = SOME \ v \\ | \ \ \mathtt{literal\_value}(\_) = NONE \\ \mathbf{fun} \ \mathtt{sum\_left}(\mathtt{sum}(x,\_)) = SOME \ x \\ | \ \ \mathtt{sum\_left}(\_) = NONE \\ \mathbf{fun} \ \mathtt{sum\_right}(\mathtt{sum}(\_,x)) = SOME \ x \\ | \ \ \mathtt{sum\_right}(\_) = NONE \end{array}
```

34 CHAPTER 3. SYNTAX

and so on for each constructor argument.

Finally, the string rendering functions (assuming a function natToString already exists in the context):

```
\begin{aligned} & \textbf{fun E\_toString}(\texttt{literal}(v)) = natToString \ v \\ & | & \texttt{E\_toString}(\texttt{sum}(x,y)) = \texttt{E\_toString}(x) + " + " + \texttt{E\_toString}(y) \\ & | & \texttt{E\_toString}(\texttt{product}(x,y)) = \texttt{E\_toString}(x) + " \cdot " + \texttt{E\_toString}(y) \\ & \textbf{and F\_toString}(\texttt{equality}(x,y)) = \texttt{E\_toString}(x) + " \doteq " + \texttt{E\_toString}(y) \\ & | & \texttt{F\_toString}(\texttt{lessOrEqual}(x,y)) = \texttt{E\_toString}(x) + " \leq " + \texttt{E\_toString}(y) \end{aligned}
```

Because ML has inductive data types as primitives, pattern-matching on our syntax comes for free. We will get back to that when defining the semantics.

3.4.2 Object-Oriented Programming Languages

In a object-oriented programming language, inductive data types are not available. Therefore, they must be mimicked using classes. On the positive side, this supports arguments names, and notations are a bit easier.

The basic recipe is as follows:

- Each type is implemented as an abstract class.
- Each constructor of type t is implemented as a concrete class that extends the abstract class t.
- The arguments names and type of each constructor c are exactly the argument names and types of the class c. The constructor arguments are stored as fields in the class.
- The abstract classes require a toString method, which is implemented in every concrete class according to its notation.

```
Example 3.30. We define our example syntax in a generic OO-language somewhat similar to Scala.<sup>1</sup> In particular, we assume that the sy
```

```
abstract class E {
 def toString: String
class literal extends E {
  field value: Nat
 constructor (value: Nat) {
    this.value = value
  def toString = value.toString
class sum extends E {
  field left: Nat
  field right: Nat
  constructor (left: E, right: E) {
    this.left = left
    this.right = right
  def toString = left.toString + "+" + right.toString
class product extends E {
  field left: Nat
  field right: Nat
  constructor (left: E, right: E) {
    this.left = left
    this.right = right
  def toString = left.toString + "." + right.toString
```

3.4. IMPLEMENTATION 35

```
abstract class F {
  def toString: String
class equality extends E {
  field left: Nat
  field right: Nat
  constructor (left: E, right: E) {
    this.left = left
    this.right = right
  def toString = left.toString + "=" + right.toString
class product extends E {
  field left: Nat
  field right: Nat
  constructor (left: E, right: E) {
    this.left = left
    this.right = right
  def toString = left.toString + "\leq" + right.toString
```

Because OO-languages do not have inductive data types as primitives, pattern-matching on our syntax requires awkward switch statements. We will get back to that when defining the semantics.

3.4.3 Combining Paradigms

The Scala language combines ideas from functional and OO-programming. That makes its representation of context-free syntax particularly elegant.

In Scala, the constructor arguments are listed right after the class name. These are automatically fields of the class, and a default constructor always exists that defines those fields. That gets rid of a lot of boilerplate.

If we want to make those fields public (and we do because those are the projection functions, we add the keyword val in front of them. But even that is too much boilerplate. So Scala defines a convenience modifier: if we put case in front of the classes corresponding to constructors of our syntax, Scala puts in the val automatically. It also generates a default implementation of toString, which we have to override if we want to implement notations, too. Finally, Scala also generates pattern-matching functions so that we can pattern-match in the same way as in ML.

Then our example becomes (as usual, assuming a class Nat already exists):

```
abstract class E {
  def toString: String
}
case class literal(value: Nat) extends E {
  override def toString = value.toString
}
case class sum(left: Nat, right: Nat) extends E {
  override def toString = left.toString + "+" + right.toString
}
case class product(left: Nat, right: Nat) extends E {
  override def toString = left.toString + "." + right.toString
}
abstract class F {
  def toString: String
}
case class equality(left: Nat, right: Nat) extends E {
```

36 CHAPTER 3. SYNTAX

```
override def toString = left.toString + "=" + right.toString
}
case class lessOrEqual(left: Nat, right: Nat) extends E {
  override def toString = left.toString + "\left" + right.toString
}
```

3.4.4 Issues Regarding Equivalence

There are a few subtle issues.

Equality Forming the same value twice should yield the same result, e.g., literal(1) = literal(1). In ML-style inductive types, that is automatic. In OO-style class-based implementations, we have newliteral(1) = newliteral(1) only if we override the equality check. Scala allows adding the keyword case to a class representing a constructor/production. Among other things, it allows dropping the new and overrides the equality method automatically.

Null Value Object-oriented language automatically provide the null value for every class. Therefore, if we implement a context-free syntax via classes, we always get one spurious value that should not exist. Programmers have to use programming discipline to avoid using it. Occasionally, static analysis tools can non-null annotations to check this discipline.

Sealing When we seal a context-free syntax, it is impossible to add new constructors to it. That guarantees implementors of inductive functions that they know all possible cases and can therefore pattern-match on them. An inductive function on an unsealed datatype can fail when new constructors are added later.

ML-style inductive data types are always sealed. OO-style class-based implementations are never sealed. In Scala, the keyword **sealed** can be added to an abstract class — it forbids any constructors other than the ones present in the same source file.

3.4.5 Implementing Context-Sensitive Checking

We extend the implementation of our example language in Scala from Sect. 3.4.3 to cover the context-sensitive language from Ex. 3.25.

```
abstract class NonTerminal {
    def print (): String
// At the toplevel, every formal language has a vocabulary, which is a list of declarations
// (we can skip the abstract class if the non-terminal has exactly one production)
case class Vocabulary (decls: List [Declaration]) extends NonTerminal {
  // print all declarations and concatenate with separator ", "
  def print() = decls.map(_.print()).mkString(", ")
}
// there can be many different kinds of declarations
abstract class Declaration extends NonTerminal {
  // most declarations have a name, but not all e.g., axioms
  def nameO: Option [String]
case class FunctionSymbolDeclaration(name: String, inputs: List[Type], output: Type) extends
  def nameO = Some(name)
  // print as "name: inputs -> output"
  def print() = name + ": " + inputs.map(_.print()).mkString(" ") + " -> " + output.print()
}
```

3.4. IMPLEMENTATION 37

```
abstract class Type extends NonTerminal
abstract class Expr extends NonTerminal
abstract class Form extends NonTerminal
case class Nat() extends Type {
  def print() = "Nat"
case class String() extends Type {
  def print() = "String"
}
// for every declaration, there is a production to refer to the declared identifier, i.e., to
case class FunctionSymbolReference(name: String, arguments: List[Expr]) extends Expr {
 // print as "name(arguments)"
  def print() = name + "(" + arguments.map(_.print()).mkString(",") + ")"
case class Sum(left: Expr, right: Expr) extends Expr {
  def print() = left.print() + "+" + right.print()
case class Product(left: Expr, right: Expr) extends Expr {
  def print() = left.print() + "*" + right.print()
case class Zero() extends Expr {
  def print() = "zero"
case class One() extends Expr {
  def print() = "one"
case class Equals(left: Expr, right: Expr) extends Form {
  def print() = left.print() + "=" + right.print()
case class LessEq(left: Expr, right: Expr) extends Form {
  def print() = left.print() + "<=" + right.print()
}
```

Then we implement a context-sensitive type-checker for it. It implements the wff-predicates from Def. 8.2.

A type-checker traverses the AST. Any such traversal consists of

- one function for every non-terminal (mutually recursive)
- one case for every constructor
- one recursive call for every constructor argument

We first do the simple case where we only have a single base type. In this case, the base types in the inputs and outputs of the symbols are always the same and only the number of arguments must be checked. In that case, the type-checker is a traversal that

- takes a Vocabulary as an extra argument that's the context that makes everything context-sensitive
- returns Boolean

```
object TypeChecker {
   // check every declaration relative to the vocabulary before it
   // this is the same for virtually every formal language
   def check_Voc(voc: Vocabulary): Boolean = {
```

CHAPTER 3. SYNTAX

38

```
var seenSofar: List [Declaration] = List()
  voc.decls.forall \{d \Rightarrow
     // check the current declaration relative to the ones seen before
     val r = check_Decl(Vocabulary(seenSofar), d)
     // append it to the list of seen declarations
     seenSofar = seenSofar ::: List(d)
  }
}
def check_Decl(voc: Vocabulary, d: Declaration): Boolean = {
 d match {
    case FunctionSymbolDeclaration(n, ins, out) =>
      if (\text{voc.decls.exists}(e \Rightarrow e.\text{nameO} = \text{Some}(n)))
        false // in practice, throw "name already defined" error
         ins.forall(i => check_Type(voc, i)) && check_Type(voc, out)
 }
// compared the the lecture, I've removed function types for simplification
def check_Type(voc: Vocabulary, tp: Type) = {
 tp match {
    case Nat() => true
    case String() => true
}
// In our simplified variant, where we only check arity, expressions are well-typed if
// all subexpressions are and applications use the right number of arguments.
// We will expand on that later.
def check_Expr(voc: Vocabulary, n: Expr): Boolean = {
   n match {
     case Zero() => true
     case One() => true
     case Sum(1,r) => check_Expr(voc, 1) && check_Expr(voc, r)
     case Product(1,r) => check_Expr(voc, 1) && check_Expr(voc, r)
     case FunctionSymbolReference(s, args) =>
       // for the reference to the identifiers in the vocabulary, we have to lookup the ide
       voc.decls.find(d \Rightarrow d.nameO = Some(s)) match {
         case None => false // in practice, throw "symbol not declared" error
         case Some(d) \Rightarrow d match \{
           case FunctionSymbolDeclaration(_, ins,out) =>
             if (ins.length != args.length) {
               false // in practice, throw "wrong number of argumetns" error
                // check every argument recursively
               args.forall(a => check_Expr(voc, a))
        }
      }
  }
}
```

3.4. IMPLEMENTATION 39

```
def check_Form(voc: Vocabulary, f: Form): Boolean = {
    f match {
      case Equals(1,r) => check_Expr(voc, 1) && check_Expr(voc, r)
      case LessEq(1,r) => check_Expr(voc, 1) && check_Expr(voc, r)
  }
}
The following code shows a simple test of the type checker:
object Test {
    def main(args: Array[String]) = {
        // because we do not have a parser, we need to build some objects manually for testin
        // an example vocabulary
        // fib: Nat -> Nat
        val fibDecl = FunctionSymbolDeclaration("fib", List(Nat()), Nat())
        val voc = Vocabulary (List (fib Decl))
        System.out.println(voc.print())
        // check the vocabulary
        System.out.println(TypeChecker.check_Voc(voc))
        // example expressions relative to that vocabulary
        val x = Sum(Zero(), One())
        val y = FunctionSymbolReference("fib", List(x))
        val z = LessEq(x,y)
        System.out.println(z.print())
        // check the expressions
        System.out.println(TypeChecker.check_Form(voc, z))
    }
}
```

3.4.6 Implementing Context-Sensitive Checking with Variables

In languages with variable binding, we have to extend the above implementation: Every checking function for expressions, which used to take a vocabulary, now takes a vocabulary and a variable context.

Chapter 4

Type Systems

4.1 Intrinsic vs. Extrinsic Typing

4.1.1 Overview

We write x : A to say that x has type A. There are two fundamentally different methods for introducing the types A, which are summarized in the following table:

	intrinsic	extrinsic
goes back to λ -calculus by	Church	Curry
general idea	objects carry their type with them	types are designated by the environment
typing is a	function from objects to types	relation between objects and types
objects have	unique type	any number of types
types often interpreted as	disjoint sets	unary predicates on a universal set
type inference for x	uniquely infer A from x	try to find minimal A such that $x:A$
type checking	compare inferred and expected type	prove $x:A$
subtyping $A <: B$	mimicked by casting from A to B	defined by $x : A$ implies $x : B$ for all x
typing decidable	yes unless too expressive	no unless expressivity restricted
typing errors are detected	usually statically (compile-time)	dynamically (run-time)
type of name introduced as	part of declaration	additional axiom
example	individual "WuV":"course"	individual "Wuv", "WuV" is-a "course"
advantages	easy	flexible
	unique type inference	allows subtyping
examples	SFOL, SQL	OWL, Scala, English
	most logics, functional PLs	ontology, OO, natural languages
	many type theories	set theories

Example 4.1 (Extrinsically Typed Ontology Language). In BOL, the objects are the individuals, the types are the concepts, and is-a is the typing relation between them. The typing is extrinsic:

- Individuals and their concept assertions are introduced in separate declarations.
- An individual may be an instance of any number of concepts.
- There is no primary concept that could be returned as the inferred type of an individual.
- Concepts are subject to subtyping $C \sqsubseteq C'$.
- Whether an individual is an instance of a concept, must be checked by reasoning about the is-a relation.

Therefore, all semantics must interpret individuals as elements of a universal collection, and types as unary predicates on that. Specifically, we have

semantics in	universal collection	unary predicate	typing relation i is-a c
FOL	type ι	predicate $c \subseteq \iota$	c(i) true
SQL	table Individuals	table containing ids	id of i in table c
Scala	String	hash set of strings	c.contains (i)
English	proper nouns	common nouns	" i is a c " is true

We can also think of relations as objects. However, BOL cannot express relation types at all, and there is no intrinsic typing. Instead, the domain and range of a relation r are given extrinsically via axioms about dom r and rng r. Like for individuals that allows flexibility as the same relations may have multiple types.

Example 4.2 (Intrinsically Typed Ontology Language). We could define TOL, a typed ontology language that arises as a variant of BOL. The main differences would be

- Individuals are declared with a concept that serves as their type: **individual** i:C.
- Concept assertions are dropped. They are now part of the individual declarations.
- Relations are declared with two concepts for their domain D and range R: relation $r <: D \times R$.
- Properties are declared with a concept for their domain C: **property** $p <: C \times T$.

TOL would make many ontologies more concise. For example, we could simply write

```
concept instructor
concept course
individual FlorianRabe : instructor
teach <: instrctor × course</pre>
```

However, we would lose flexibility. If we want to add the concept "male", it would be difficult to make FlorianRabe have both types. We might be able to remedy that by allowing intersections and declaring individual FlorianRabe: instructor \sqcap male. But even then, we would have to commit to the type of each individual right away — we cannot add different concept assertions for the same individual in different places, a common occurrence in building large ontologies.

Allowing \sqcap would also introduce subtyping. If we are careful in the design of TOL, that may still result in an elegant scalable language. In particular, typing may remain decidable (depending on what other operations we allow). But if we go too far, it may end up so complex that it would have been easier to go with extrinsic typing. That is why we use intrinsic typing only in two related places in BOL:

- The base types and values use an intrinsic type system (whose details we omitted).
- The range of properties is given intrinsically by a base type.

Remark 4.3 (Subtyping). Languages with subtyping usually have to use extrinsic type systems. Typical sources of subtyping are

- explicit subtyping as in $\mathbb{N} <: \mathbb{Z}$
- comprehension/refinement as in $\{x : \mathbb{N} | x \neq 0\}$
- operations like union and intersection on types
- inheritance between classes, in which case subclass = subtype
- anonymous record types as in $\{x: \mathbb{N}, y: \mathbb{Z}\} <: \{x: \mathbb{N}\}$

4.1.2 Combined Definition

Neither intrinsic nor extrinsic typing is strictly better than the other. The choice of type system is a very difficult trade-off when designing a language.

Many practical languages even combine both methods. In that case, an intrinsic system is used for the most important high-level types and an extrinsic system is used to refine (some of) the high-level types:

Definition 4.4 (Type System). A **type system** for a context-free syntax consists of provides for some expression symbols \mathcal{E}

- a non-terminal \mathcal{T} \mathcal{T} -expressions are called the **types** of \mathcal{E} -expressions
- a binary predicate $\vdash_{\Theta} E : T$ on \mathcal{E} -expression E and \mathcal{T} -expressions T if satisfied, we say that E has type T.

We can now recover the intuitions from above as special cases:

- A purely intrinsic type system is one in which the types correspond to the expression symbols themselves. Here each \mathcal{E} -expression has type \mathcal{E} , and the grammar already separates expressions by their (intrinsic) types.
- A purely extrinsic type system has two non-terminals \mathcal{E} and \mathcal{T} . Here all expressions are merged into the same non-terminal, i.e., the grammar does not distinguish expressions by their types. All typing is done extrinsically by the typing predicate.

Example 4.5. We can think of BOL as a combined type system. The intrinsic types are the non-terminals I, C, R, P, and F, which separate the expressions into the five kinds of individuals, concepts, relations, properties, and formulas.

An extrinsic typing relation exists for I: the types for $\mathcal{E} = I$ are those derived from $\mathcal{T} = C$, i.e., the individuals are extrinsically typed by the concepts. The typing relation $\vdash_{\Theta} E : T$ is given by the formula E is $\neg a$ is true relative to Θ .

Example 4.6. In set theory, only a few intrinsic types are used for the high-level grouping of objects. These include at least set and prop. Objects of these intrinsic types are called sets and propositions. Some set theories also use an intrinsic type class. Moreover, types like $set \to prop$ can be allowed as the types of unary predicates on sets. Extrinsic typing is used only for $\mathcal{E} = set$: we have $\mathcal{T} = set$ and the typing relation is the elementhood relation \in between sets.

Example 4.7. In some sense every language has a type system. What is commonly called an untyped language, is technically an intrinsically typed language with only one type.

However, in practice, almost all of this languages have more than one type or some kind weak of additional extrinsic typing. For example, untyped first-order logic really has two intrinsic types: terms and propositions. Python has one intrinsic type of expressions, but these are extrinsically subdivided into lists, dictionaries, integers, etc.

4.2 Variable Contexts and Substitutions

In the presence of a type system, we often use typing in context.

Definition 4.8 (Context). A variable declaration is of the form x : A for an identifier x and a type A. A variable context is a list of variable declarations.

An expression in context Γ may additionally use the variable from Γ . We write $\Gamma \vdash_{\Theta} E : A$ to say that expression E has type A in context Γ over vocabulary Θ .

Example 4.9. Consider the language from Ex. 3.27 with the standard library from Ex. 3.26. Then $x : Nat, y : Nat \vdash equals(x,y) : FORM$ is a formula in context.

Example 4.10. We can also use the non-terminal expression symbols as types. For example, BOL has no variable binding and therefore no contexts in the grammar. But we still have at least contexts for the intrinsic type system given by the grammar. For example, we can use a context c:C,d:C for two concept variables and give the formula $c:C,d:C\vdash c\sqsubseteq d:F$.

In general, for any context-free grammar, we can use the non-terminals as intrinsic types and form contexts. Then x:N means that x is an arbitrary word derived from the non-terminal N. Such contexts are very common when talking about a language. We speak of meta-contexts. For example, in BOL, we might say "Given a concept c and a relation r, \ldots " — that corresponds to using the context c:C,r:R.

Definition 4.11 (Substitution). Given two contexts Γ and Δ , a substitution provides for every variable x:A exactly one variable assignment x/e where e is an expression from Δ .

For a Γ -expression E, we write $E[\gamma]$ for the Δ -expression resulting from substituting every x in E with the corresponding assignment from γ .

A substitution is well-typed, written $\gamma:\Gamma\to\Delta$, if every pair of and x:A in Γ and x/e in γ satisfies $\Delta\vdash e:A[\gamma]$.

Example 4.12 (Identity). The identity substitution id_{Γ} from Γ to itself assigns x/x for every x in Γ .

The mapping $E \mapsto E[id_{\Gamma}]$ is the identity mapping on Γ -expressions.

Example 4.13 (Composition). Given substitutions $\gamma_1: \Gamma_1 \to \Gamma_2$ and $\gamma_2: \Gamma_2 \to \Gamma_3$, the composition $\gamma_1; \gamma_2: \Gamma_1 \to \Gamma_3$ assigns $x/x[\gamma_1][\gamma_2]$ for every x.

The mapping $E \mapsto E[\gamma_1][\gamma_2]$ is the composition of the mappings $E \mapsto E[\gamma_1]$ and $E \mapsto E[\gamma_2]$.

Virtually every type system satisfies the following property. Technically, this must be proved separately for every type system because it is not guaranteed to hold. But in practice, a type system that does not satisfy it is simply broken:

Theorem 4.14 (Type Preservation). Consider a reasonably defined type system. Given a substitution $\gamma: \Gamma \to \Delta$ and an expression $\Gamma \vdash E: A$, we have $\Delta \vdash E[\gamma]: A[\gamma]$.

4.3 Vocabulary Morphisms

For every formal system, we can distinguish three kinds of symbols with associated translations.

symbol	origin	owner	examples	translations
logical symbols = lan-	built into language	language creator	\wedge in SFOL, \sqcup in	language translations
guage symbols			BOL	
non-logical symbols =	declared in vocabulary	vocabulary author	person, WuV	vocabulary morphisms
vocabulary symbols		= language user	from Ex. 2.3	
variables	declared in context	containing declara-	x, y, \dots	substitutions
		tion		

Vocabulary morphisms behave like substitutions except that they map vocabulary symbols to expressions as opposed to variables.

4.4 Abstract and Concrete Data Types

The general thrust of type systems is to shift more and more information into an increasingly complex type system. This is part of a trade-off: the more the type system can do,

- the more requirements can be expressed and violations thereof detected statically,
- the more complex the type system and its documentation and implementation become.

Concrete and abstract data types have proved to be a particularly interesting trade-off on this expressivity-simplicity spectrum and are — in one way or another — part of many type systems.

4.4.1 Abstract vs. Concrete Types

The words *abstract* and *concrete* do not have standard definitions for types. I like the intuitions described below. A type is called **concrete** if its values are

- given by their internal form,
- defined along with the type, typically built from already-known pieces.

In other words, a concrete type is given by an *extensional* description of its objects.

Example 4.15. Simple products are concrete types.

They are introduced by (among other rules)

- $A \times B$ is a type if A and B are
- values of type $A \times B$ are of the form (a, b) for a : A and b : B.

Example 4.16. Enumerations are concrete types that explicitly list a finite set of distinct values.

They can also be seen as the special case of inductive data types in which all constructors take 0 arguments.

Example 4.17. Inductive data types as seen in Def. 3.8 are concrete types. Their values are formed by applying constructors to other values.

I like calling them concrete data types.

A type is called **abstract** if its values are

- given by their externally visible properties,
- defined in any environment that understands the type definition.

In other words, an abstract type is given by an *intensional* description of its objects.

This is the case for abstract data types.

4.4.2 Concrete Data Types

Concrete data types are the same as inductive data types and were already defined rigorously in Sect. 3.2.2.

It is in the nature of concrete types that they do not depend on the situation in which they are used. The definition of the type once and for all defines the possible values. Thus, they do not differ (or differ at most subtly) between languages. Often the only real difference is whether a language has concrete data types or not. Languages that do not natively provide them may differ substantially in how they simulate them though.

4.4.3 Abstract Data Types

Contrary to concrete types, the values of an abstract data types depend on what features the language provides to construct objects that satisfy the abstract requirements. Therefore, the treatment differs widely across languages. The following table gives an overview:

aspect	language	abstract data type
ontologization	UML	class
concretization	SQL	table schema
computation	Scala	class, interface
deduction	various	theory, specification, module, locale
narration	various	emergent feature

Example 4.18 (Classes). A UML class is an abstract data type. Its values are the instances of implementing classes.

A UML class only defines what methods should be available. How they are implemented by specific values of the type is left to the programming languages.

Thus, different programming languages could have different values for the same abstract data type. They certainly look different, e.g., in Java and Scala implementations of the same UML class. But the languages might also be fundamentally different in expressivity, e.g., a Turing-complete programming language might have strictly more values for the same abstract data type than a non-Turing-complete one.

Moreover, which instances actually exist changes during the run time of the program. If we take this into account, the values of the abstract data type are not even fixed within a programming language.

Example 4.19 (Schemas). An SQL table schema is an abstract data type. Its values are the rows.

The schema only defines what types the columns of a table have. Different database systems might theoretically provide different ways to build rows for the table.

However, this does not happen in practice because SQL table columns are typed by base types, which have the same values across database systems. This would be different if we allowed table columns to have function types.

Example 4.20 (Theories). A logical theory (e.g., Monoid) is an abstract data type. Its values are the models of the theory (e.g., for Monoid: $(\mathbb{N}, +, 0)$ or $(\mathbb{N}, *, 1)$).

The theory only defines what operations a model must provide (for Monoid: binary operation and neutral element) and which axioms it must satisfy (for Monoid: associativity, neutrality). How we build the models is left open.

We usually build models in mathematical language and naively assume that fixes the models once and for all. But that is too naive: depending on which mathematical foundation we use (e.g., set theory with or without axioms of choice), we can build different models. Moreover, we can also build models in type theories (which underly many deduction systems such as Coq or Isabelle). We can even build them in programming languages, e.g., by implementing theories as classes (typically moving the axioms into comments).

The choice of language substantially changes what the values of the abstract data type are.

In the sequel we introduce a fairly general definition that subsumes many practical languages.

Definition 4.21 (Abstract Data Type). Consider an arbitrary type system.

An abstract data type (ADT) is

• a **flat** type of the form

$$\{c_1: T_1[=t_i], \ldots, c_n: T_n[=t_i]\}$$

where the c_i are distinct names, the T_i are types, and the t_i are optional and wherever given must have type T_i , or

• a **mixin** type of the form $A_1 * A_2$ for ADTs A_i .

We say that a type system has **internal ADT**s if all ADTs are types (and thus may in particular occur as the T_i in a record type).

The intuition of a mixin A * B is that we merge the fields of A and B. However, this union dependent: if B is flat, its fields may refer to fields introduced in A.

The most important special case of an ADT are classes:

Definition 4.22 (Class). A class definition defines an ADT abbreviation of the form

$$a = a_1 * \dots * a_m * \{c_1 : T_1, \dots, c_n : T_n\}$$

where the a_i are names of previously defined ADTs.

We call the a_i the **superclasses** or **parent classes** and say that a inherits from the a_i . We call the c_i the **fields** or **members** of a.

In an OO-language, a class definition is more commonly written somehow like

```
abstract class a extends a_1 with ... with a_m { c_1 \colon T_1 \vdots c_n \colon T_n }
```

The details can vary, and special care must be taken in programming languages where initialization may have side effects.

Flat ADTs are the standard case, and all mixin ADTs can be simplified into flat ones. This can be seen as a semantics in the sense that the language of flat and mixin ADT is translated to the language of flat ADTs.

Definition 4.23 (Mixin Semantics). The flattening A^{\flat} of an ADT A is defined as follows:

- if A is flat: $A^{\flat} = A$
- if A is of the form $A_1 * ... * A_n$: A^{\flat} arises by concatenating the fields of all A_i^{\flat} where duplicate field names are handled as follows:
 - if the same field (same name, types equal, definitions equal or both absent) occurs more than once,
 only the first occurrence is kept,
 - if the fields $c: T_1[=t_i]$ and $c: T_2[=t_2]$ occur for inequal types T_i , A is ill-formed,
 - if the fields $c: T = t_1$ and $c: T = t_2$ occur for inequal objects t_i , A is ill-formed,
 - if the fields c: T = t and c: T occur, only the defined one is kept (*).

Remark 4.24 (Dependency Between Fields). Our definition sweeps a very important but subtle detail under the rug: in a flat ADT with a field c: T = t, may T and/or t refer to fields declared later? We sketch a few possible answers.

In the simplest case, we forbid such forward references. Then ADTs are very well-behaved. But we have a problem with the case (*) in Def. 4.23: if c:T occurs before c:T=t, we cannot simply drop the former because intermediate fields may refer to c. A straightforward solution would be to declare the ADT to be ill-formed. But unfortunately, this case is very important in practice — it occurs whenever c:T is declared in an abstract class and c:T=t in a concrete class implementing it.

A more common solution is to allow the fields to be mutually recursive. Consider a flat ADT with fields $\Gamma, c : T[=t], \Delta$ where Γ and Δ are lists of fields. Let Γ' and Δ' arise by dropping all definitions. Then we require that

- T must be a well-formed type in context Γ' . Thus, the types may only refer to previous fields.
- t must have type T in context Γ' , $c:T,\Delta'$. Thus, the definitions may be mutually recursive.

This makes the case (*) work. But it comes at the price of recursion, which allows writing non-terminating fields (a feature in a programming language, but potentially undesirable in other settings).

Even so, the mutual-recursion solution is problematic in the presence of dependent types. Here, dropping definitions is not always allowed: T might be well-formed in context Γ , but Γ' might not even be a well-formed context at all. Because OO-languages are usually not dependently-typed, this is not an issue in most settings.

4.5 Abstract Data Types in Ontologies

4.5.1 Motivation

Recall the subject-centered representation of individuals described in Sect. 2.3. Here we introduce an individual together with all assertions of which it is the subject as in

```
individual "FlorianRabe"
is—a "instructor" "male"
"teach" "WuV" "KRMT"
"age" 40
"office" "11.137"
```

It is often desirable to use types to force the presence of such assertions. We might wish require that every instructor teaches a list of things, and has an office. Moreover, we can use types to specify the objects of the respective assertions: we can specify that only courses are taught and that the office is a string. Rather than the relations with subjects "FlorianRabe" just happening to be around as well, the type system would now force their existence and the type of the object. Forgetting to give such an assertion or giving it with the wrong object could be detected statically (i.e., without applying the semantics) and flagged as a typing error.

This leads to the idea of **subject-centered types**. This could looks as follows:

```
concept instructor
  teach course*
  age: int
  office: string

individual "FlorianRabe": "instructor"
  is—a "male"
  "teach" "WuV" "KRMT"
  "age" 40
  "office" "11.137"
```

Now the type "instructor" forces the presence of a list of taught courses (The * is meant to indicate a list.), an integer for the age, and a string for the office.

We can now see that, in fact, every person should have an age, and not just every instructor. Because every instructor is meant to be a person, we could try to capture this as well to avoid redundancy. Moreover, every male is meant to be a person, too.

That leads to the idea of **modular types**. This could look as follows:

```
concept person
   age: int

concept male <: person

concept instructor <: person
   teach course*
   office: string

individual "FlorianRabe": "instructor" □ "male"
   "teach" "WuV" "KRMI"
   "age" 40
   "office" "11.137"
```

Incidentally, that eliminates the need to independently declare relations and properties. Instead, we can treat their occurrences inside the concept definitions as their declarations.

That has the added benefit that two relations/properties of the same name declared in different concepts can be distinguished and can have different types.

4.5.2 Database Schemas as Typed Ontologies

It is very typical that optimized systems are developed for one of the corner aspects. These systems then subsume an ontology language. From the Tetrapod perspective, it is instructive to tease apart the ontology.

In particular, what is typically called an ontology follows the semantic web community. Tetrapodally, they are working on the combination of ontology and concrete data. The former defines the rules of the world (TBox), the latter populates the world (ABox).

We can see the same pattern in the relational database community. Tetrapodally, the SQL schema defines the ontology, and the database holds the concrete data. Fig. 4.1 shows the analogy in detail.

	semantic web	relational databases
ontology aspect	TBox of ontology	SQL schema
conceptual model	knowledge graph	set of tables
concrete data aspect	ABox of ontology	SQL database
concrete data storage	set of triples	set of rows of the tables
concrete data formats	RDF	CSV
concrete data tool	triple store	database implementation
typing	soft/Curry	hard/Church
query language	SPARQL	SQL SELECT query
openness of world	tends to be open	tends to be closed

Figure 4.1: The Two Kinds of Concrete Data Systems

We can then understand relational databases as an ontology+concrete data approach that focuses on abstract data types. Each table is one such type, and each column is a field of the type. The database holds the concrete data, i.e., the rows of the tables.

For example, the SQL schema might declare an ADT by

```
TABLE person(id: ID, name: string, age: int)
TABLE male(id: Ref person)
TABLE instructor(id: ID, parent: Ref person, teach: course*, office: string)
```

Here we write Ref T for the type ID of identifiers if we want to emphasize that it should hold an identifier that is present in the table T.

Concepts like male can be represented as single-column tables holding the identifiers of individuals that are in the concept. Inheritance can be represented by using two separate tables for parent and child type with the latter holding a special column parent holding an identifier in the parent table. For example, concrete data could be inserted as

```
INSERT VALUES (id = X, name = "FlorianRabe", age = 40) INTO person
INSERT VALUES X INTO male
INSERT VALUES (id = Y, parent = X, teach = ["WuV", "KRMT"], office = "11.137")
INTO instructor
```

Usually, relations where the same subject can have multiple objects are not stored as a list-valued column but as a separate two-column table. For example, we would have

```
TABLE instructor(id: ID, parent: Ref person, office: string)
TABLE course(id: ID, name: string)
TABLE teach(subject: Ref person, object: Ref Course)
INSERT VALUES (id = W, name = "WuV") INTO course
INSERT VALUES (id = K, name = "KRMT") INTO course
INSERT VALUES (subect = X, object = W) INTO teach
INSERT VALUES (subect = X, object = K) INTO teach
```

4.5.3 Exercise

The topic of Exercise is to build a relational database schema for a univis-like system.

Chapter 5

Data

5.1 Overview

This chapter section presented ongoing research on developing infrastructure for the semantic representation and interchange of data across systems. It is presented on the slides and the lecture videos.

52 CHAPTER 5. DATA

Chapter 6

Semantics

This chapter introduces different kinds of semantics. The emphasis will be on relative semantics by compositional translations.

We use BOL as a running example and give four different semantics of BOL — using the four other aspects:

Section	Aspect	kind of semantic language	semantic language
6.3	deduction	logic	SFOL
6.4	concretization	database language	SQL
6.5	computation	programming language	Scala
6.6	narration	natural language	English

6.1 Kinds of Semantics

To define a language system, we need to define the well-formedness predicates. Moreover, we need to define the semantics of the well-formed vocabularies and expressions. One way to do that is by translating the syntax into another language and then use existing definitions of well-formedness there.

But we also need to be able to get off the ground, i.e., to define a semantics from scratch when we do not have another language available. This is usually done by giving an inference system for the well-formedness predicates and other judgments such as truth.

We can also do both: if we have a semantics via an inference system and another one via translation, we can show that the latter respects the former. That leads to the concepts of soundness (everything well-formed is translated to something well-formed) and its dual completeness.

6.1.1 Absolute Semantics

Absolute semantics is very hard to define in general. A typical approach is to single out one property of the syntax and then use an inference system to define when that property holds.

For example, for BOL we could use the property whether a formula can be proved from the axiom. The inference system then is the proof calculus that define the provable formulas.

6.1.2 Relative Semantics

Relative semantics uses an additional component as a reference point for semantics.

By Translation Semantics by translation uses a second language as the target of a translation. Both the vocabularies and the expressions are translated, namely to vocabularies and expressions of the target language. Usually, the target language already has a semantics, and the semantics of the original language is obtained by composing the two.

The correspondence for the syntax between context-free grammars and inductive data types can be extended to the semantics. Now we have a correspondence between case-based function definitions and inductive functions.

Aspect	Example
Deductive	$FOL \rightarrow set theory$
Computational	$C++ \rightarrow assembly$
Concrete	$SPARQL \rightarrow SQL$
Narrative	$English \to German$

Figure 6.1: Examples of Semantics by Translation

Definition 6.1. A semantics by translation consists of the following parts:

- \bullet syntax: a language system l
- semantic language: a language system L (from a different or the same aspect as l)
- ullet semantic prefix: a vocabulary P in L that is prefixed to the translation of all vocabularies of l
- vocabulary translation: a function that translates every l-vocabulary T to an L-vocabulary P, $[\![T]\!]$

Critically, the semantic language (which is itself a language system and can thus have a semantics itself) must be a language whose semantics we already know. Therefore, it is often important to give multiple equivalent semantics — choosing a different semantics for different audiences, who might be familiar with different languages.

The role of the semantic prefix P is to define once and for all the L-material that we need in general to interpret l-vocabularies (in the case of BOL: ontologies). It occurs at the beginning of all interpretations of ontologies. In particular, it is equal to the interpretation of empty vocabulary.

By Interpretation Semantics by interpretation uses a situation relative to which expressions are evaluated. The role of the situation is to supply the meaning of all identifiers declared in the vocabulary. Given a situation, inductive functions map all expressions to their semantics. Usually, the result of interpretation is an extremely simple object, whose semantics is self-evident.

Aspect	Situation	Provides
Deductive	model	interpretations of the function/predicate/ symbols
Computational	environment	e.g., standard input/output
Concrete	database	set of objects for each type
Narrative	semantic dictionary	meanings of the words

Figure 6.2: Examples of Situational Semantics by Interpretation

Situational semantics can be seen as a special case of semantics by translation as follows:

- The target language is an implicit assumed background language such as mathematics or the computer hardware. Once this is made explicit, situations can be described as vocabularies of the background language.
- The semantic prefix describes what a situation is in general, i.e., the possible situations for the empty vocabulary.
- The semantics of a vocabulary extends the semantic prefix with the possible ways to interpret the identifiers.
- The semantics of an expression E is the function mapping the situation S to the situational semantics of E under S. This can be expressed as an expression of the target language.

6.2 Compositionality

Introduction There are some general principles shared by all translations:

- Every *l*-declaration is translated to an *L*-declaration for the same name, and ontologies are translated declaration-wise.
- For every non-terminal N of l, there is one inductive function $[\![-]\!]_N$ mapping complex l-expressions derived from N to L-expressions.

55

- The base cases of references to declared l-identifiers are translated to themselves, i.e., to the identifiers of the same name declared in L.
- The other cases are compositional: every case for a complex *l*-expression recurses only into the semantics of the direct subexpressions.

The notion of compositionality captures these properties. An interpretation function is compositional if the interpretation of any kind of expression $E(e_1, \ldots, e_n)$ with subexpressions e_i only depends on E and the interpretation of the e_i , i.e.,

$$[E(e_1,\ldots,e_n)] = [E]([e_1],\ldots,[e_n])$$

for some semantic operation $[\![E]\!]$. Compositionality is also called the substitution property or the homomorphism property. See also Def. 6.11.

More rigorously, we define a compositional translation as follows:

Definition 6.2 (Compositional Semantics). Consider a semantics for syntax grammar l and interpretation function [-].

[-] is compositional if it is defined as follows:

- a family of functions $[\![-]\!]_N$, one for every non-terminal N of l
- for every expressions E derived from N, we put $[\![E]\!] = [\![E]\!]_N$
- each $[-]_N$ is defined by induction on the productions for N
- for each production $N := *(N_1, ..., N_r)$ and all expressions e_i derived from N_i

$$[\![*(e_1,\ldots,e_r)]\!]_N = [\![*]\!]([\![e_1]\!]_{N_1},\ldots,[\![e_r]\!]_{N_r})$$

for some L-expression $[\![*]\!]$

Without loss of generality, we can assume that every production is of the form $N := *(N_1, ..., N_r)$ where the N_i are all the non-terminals on the right-hand side and * is a stand-in for all the terminal symbols.

Compositional Translations of Contexts We can extend every compositional translation to contexts, substitutions, and expressions in contexts:

Definition 6.3. Given a translation [-] as above, for a non-terminal N, we define [N] as the non-terminal from which the translations of N-expressions are derived.

Then we define:

$$[\![x_1:N_1,\ldots,x_n:N_n]\!] := x_1:[\![N_1]\!],\ldots,x_n:[\![N_n]\!]$$
$$[\![x_1:=w_1,\ldots,x_n:=w_n]\!] := x_1:=[\![w_1]\!],\ldots,x_n:=[\![w_n]\!]$$
$$[\![x]\!] := x$$

The requirement of compositionality is critical for two reasons:

- A non-compositional translation could translate l-expressions derived from the same non-terminal N to L-expressions derived from different non-terminals. Then we would not be able to define $[\![N]\!]$.
- The definition [x] := x adds a case to the case distinction in the compositional translation function. Without compositionality, this would not make sense.

Theorem 6.4 (Type Preservation). For a compositional translation as above, we have

$$\Gamma \vdash_{l} w : N$$
 implies $\llbracket \Gamma \rrbracket \vdash_{L} \llbracket w \rrbracket : \llbracket N \rrbracket$

Substitution Theorem The main value of compositionality is the following:

Theorem 6.5 (Substitution Theorem). Consider a compositional semantics.

For every context $\Gamma = x_1 : N_1, \dots, x_r : N_r$, every syntax expression $\Gamma \vdash_l E : N$, and every substitution $\vdash_l \gamma : \Gamma$

$$\llbracket E[\gamma] \rrbracket = \llbracket E \rrbracket [\llbracket \gamma \rrbracket]$$

Formulated without substitutions, this means that for every syntax expression $E(e_1, \ldots, e_r)$ derived from N, where the e_i are subexpression derived from non-terminal N_i , we have

$$[E(e_1,\ldots,e_n)]_N = [E]([e_1]_{N_1},\ldots,[e_n]_{N_r})$$

Simply put, a semantics is compositional iff it is defined by mutually inductive translation functions with only compositional cases. The latter is very easy to check by inspecting the shape of the finitely many cases of the definition. The former is a powerful property because it applies to any of the infinitely many expressions of the syntax.

Non-Compositional Semantics It is highly desirable but not always possible to give a compositional translation. Sometimes a feature of the syntactic language cannot be directly interpreted in the semantic language. In that case, it may still be possible to give a non-compositional translation.

Example 6.6 (Non-Compositional Translation via Sub-Induction). A simple example of non-compositionality is the translation of natural numbers based on zero, one, and addition (i.e., $N::=0 \mid 1 \mid N+N$) into natural numbers based on zero and successor (i.e., $N::=0 \mid succ(N)$): It is straightforward to translate zero and one compositionally:

$$[0] = 0$$
 $[1] = \operatorname{succ}(0)$

Now we would like to translate

$$[\![m+n]\!] = [\![+]\!] ([\![m]\!], [\![n]\!]),$$

but there is no way to define [+] in terms of zero and successor. Instead, we need subcases:

$$[\![m+n]\!] = \begin{cases} [\![m]\!] & \text{if } n = 0\\ \text{succ}([\![m]\!]) & \text{if } n = 1\\ [\![(m+n_1)+n_2]\!] & \text{if } n = n_1+n_2 \end{cases}$$

This corresponds to the usual definition of addition, i.e., [+], by induction.

Other common examples of non-compositional translations are

- several important logical theorems such as
 - cut elimination, which is the translation from sequent calculus with cut to sequent calculus without cut,
 - the deduction theorem, which is the translation from natural deduction to Hilbert calculus,
- almost anything done by an optimizing compiler, e.g., loop unrolling or function inlining,
- query optimization done by a database, e.g., turning a WHERE of a JOIN into a JOIN of WHEREs,
- almost all translations between natural languages, e.g., when words are ambiguous and a different translation must be chosen for the same word based on the context (The introduction of richer intermediate structures like ASTs and functions as values into the translation can recover some compositionality here).

Typical sources of non-compositionality in formal language translations are:

- A case in the translation function requires subcases which inspect the e_i and treat them differently.
- A case in the translation function requires subcases which translate an expression differently based on the context in which it occurs.
- The translation function requires nested inductions, i.e., a case in the translation function (which is already inductive) requires a sub-induction on one of the sub-expressions.
- The semantic prefix is not fixed but depends on the translated object, i.e, the top-level case of the translation scans through the entire argument X to collect all occurrences of a particular feature and then custom-builds the semantic prefix of [X].

See also Ex. 6.12.

Such non-compositional translations are undesirable for multiple reasons:

- The implementation is more complicated and error-prone.
- Reasoning about the translation is more difficult.
- The custom semantic prefix can be large.

But most importantly, non-compositional translations are less robust. Firstly, if we add a production to the syntax, a compositional translation is easy to extend: just add a case to the translation. But a non-compositional translation

may additionally require a new subcase wherever subcases/subinductions are used. Moreover, if a custom semantic prefix is used, its definition may have to be amended, at least it must be rechecked.

Secondly, in practice there are two sources of complex expressions: the ones already mentioned in the language, and the ones used later for other reasons. For example, some complex expressions occur already statically in the definition of a vocabulary V. But others might be appear dynamically later, e.g., when talking about V, proving properties of V, or running queries on V. Thus, the definition of V and the use of complex expressions are decoupled: V is defined statically once and for all, and complex expressions relative to V can be created and used dynamically. But if a custom semantic prefix is used, only the static occurrences inside V can be considered for building the prefix. Thus, it is not possible to translate expressions dynamically unless the semantic prefix is extended all the time while V is used.

6.3 Deductive Semantics

We fix one language that we have already understood and define an interpretation function that maps all complex expression of BOL to the semantic language. For simple ontology languages like BOL, ALC, OWL, etc., it is common to use first-order logic (FOL) as the deductive semantic language. More specifically, we use SFOL, the typed variant of FOL:

6.3.1 A Basic Semantic Language: SFOL

We give typed first-order logic (SFOL) as a language system.

Definition 6.7. Fig. 6.3 gives the context-free grammar. The vocabulary symbol is Thy. The expression symbols are Y, T, and F.

6.3.2 Semantics

Definition 6.8 (Deductive Semantics of BOL). The semantic prefix is the SFOL-theory containing

- a type ι (for individuals),
- additional types and constants corresponding to base types and values of BOL.

Every BOL-ontology O is interpreted as the SFOL-theory P, [0], where [0] is defined in Fig. 6.4.

As foreshadowed above, we can observe some general principles: Every BOL-declaration is translated to an SFOL-declaration for the same name, and ontologies are translated declaration-wise. For every kind of complex BOL-expression, there is one inductive function mapping BOL-expressions to SFOL-expressions. The base cases of references to declared BOL-identifiers are translated to themselves, i.e., to the identifiers of the same name declared in the SFOL-theory. The other cases are compositional: every case for a complex BOL-expression recurses only into the semantics of the direct subexpressions.

```
Vocabularies: theories
Thy := D
Declarations
     := \mathbf{type} \; \mathtt{ID}
                                        type declaration
          \mathbf{fun}\;\mathtt{ID}:Y^{\pmb{*}}\to Y
                                        function symbol declaration
          \operatorname{\mathbf{pred}} \operatorname{\mathtt{ID}} \subseteq Y^*
                                        predicate symbol declaration
          \mathbf{axiom}\; F
                                        axiom
 type expressions
     ::= ID
                                        atomic type
term expressions
T
      ::= ID(T^*)
                                        function symbol applied to arguments
       ID
                                        term variables
 formula expressions
     ::= ID(T^*)
                                        predicate symbol applied to arguments
          T \doteq_Y T
                                        equality of terms at a type
           \top
                                        \operatorname{truth}
           \perp
                                        falsity
           F \wedge F
                                        conjunction
           F \vee F
                                        disjunction
           F \Rightarrow F
                                        implication
           F \Leftrightarrow F
                                        equivalence
           \neg F
                                        negation
          \forall \mathtt{ID}: Y.F
                                        universal quantification at a type
          \exists \mathtt{ID}: Y.F
                                        existential quantification at a type
Identifiers
     ::= alphanumeric string
```

Figure 6.3: Syntax of SFOL

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	BOL Syntax X	Semantics $[\![X]\!]$ in SFOL
BOL declaration individual i concept i unary predicate symbol $i : \iota$ unary predicate symbol $i \subseteq \iota$ binary predicate symbol $i \subseteq \iota$ binary predicate symbol $i \subseteq \iota \times \iota$ axiom $[C]([I])$ $I_1 R I_2$ axiom $[R]([I_1], [I_2])$ $I P V$ axiom $[R]([I_1], [I_2])$ $I I I I I I I I I I I I I I I I I I I$	ontology	SFOL theory
BOL declaration individual i concept i unary predicate symbol $i : \iota$ unary predicate symbol $i \subseteq \iota$ binary predicate symbol $i \subseteq \iota$ binary predicate symbol $i \subseteq \iota \times \iota$ axiom $[C]([I])$ $I_1 R I_2$ axiom $[R]([I_1], [I_2])$ $I P V$ axiom $[R]([I_1], [I_2])$ $I I I I I I I I I I I I I I I I I I I$	D_1,\ldots,D_n	$\llbracket D_1 rbracket, \ldots, \llbracket D_n rbracket$
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	BOL declaration	
$ \begin{array}{llll} \textbf{relation} i & \text{binary predicate symbol } i \subseteq \iota \times \iota \\ \textbf{property} i : T & \text{binary predicate symbol } i \subseteq \iota \times T \\ \textbf{axiom} & & & & & & & & & & & & & & & & & & &$	$\mathbf{individual}i$	nullary function symbol $i:\iota$
$\begin{array}{llll} & \text{binary predicate symbol } i \subseteq \iota \times T \\ I \text{ is-a } C & \text{axiom } \llbracket C \rrbracket (\llbracket I \rrbracket) \\ I_1 R I_2 & \text{axiom } \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{axiom } \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{axiom } \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{axiom } \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{axiom } \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{axiom } \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{Formula} & \text{Formula without free variables} \\ C_1 \sqsubseteq C_2 & \forall x : \iota \cdot \llbracket C_1 \rrbracket (x) \Rightarrow \llbracket C_2 \rrbracket (x) \\ \forall x : \iota \cdot \llbracket C_1 \rrbracket (x) \Rightarrow \llbracket C_2 \rrbracket (x) \\ \text{I is-a } C & \llbracket C \rrbracket (\llbracket I \rrbracket) \\ \parallel R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{I PV} & \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{I PV} & \llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ \text{I Individual} & \text{Terms of type } \iota \\ i & i(x) \\ \hline & & & & & & & & & & \\ C_1 \sqcup C_2 & \llbracket C_1 \rrbracket (x) \vee \llbracket C_2 \rrbracket (x) \\ \forall R.C & \forall y : \iota \cdot \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (y) \\ \exists R.C & \exists y : \iota \cdot \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (y) \\ \exists R.C & \exists y : \iota \cdot \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (y) \\ \text{dom } R & \exists y : \iota \cdot \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (y) \\ \text{dom } R & \exists y : \iota \cdot \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (y) \\ \exists y : \iota \cdot \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (x) \Rightarrow \llbracket $	$\mathbf{concept}i$	unary predicate symbol $i \subseteq \iota$
$\begin{array}{llllllllllllllllllllllllllllllllllll$	${\bf relation}i$	binary predicate symbol $i \subseteq \iota \times \iota$
$I_1 R I_2$ $I P V$ $I P V$ $Axiom [P]([I], [V])$ F $Axiom [F]$ Formula $C_1 \equiv C_2$ $C_1 \sqsubseteq C_2$ $C_1 \sqsubseteq C_2$ $I \text{ is-a } C$ $I_1 R I_2$ $I_2 P V$ $I_3 P V$ $I_4 R I_2$ $I_5 P V$ $I_5 P V$ $I_7 P V$ $I_8 P P P P P P P P P P P$ $I_8 P P P P P P P P P P P P P P P P P P P$	$\mathbf{property}i:T$	binary predicate symbol $i \subseteq \iota \times T$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	I is-a ${\cal C}$	axiom $[\![C]\!]([\![I]\!])$
Formula Formula without free variables $C_1 \equiv C_2$ $\forall x : \iota. \llbracket C_1 \rrbracket (x) \Leftrightarrow \llbracket C_2 \rrbracket (x)$ $\forall x : \iota. \llbracket C_1 \rrbracket (x) \Rightarrow \llbracket C_2 \rrbracket (x)$ I is-a C $\llbracket C \rrbracket (\llbracket I \rrbracket) \rrbracket$ $\llbracket R \rrbracket (\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket)$ $\llbracket P \rrbracket (\llbracket I \rrbracket, \llbracket V \rrbracket)$ Individual Terms of type ι i i i i i i i i i i i i True $ \bot \qquad \qquad$	$I_1 R I_2$	$Axiom [R]([I_1], [I_2])$
Formula $C_1 \equiv C_2 \qquad \forall x : \iota. \llbracket C_1 \rrbracket (x) \Leftrightarrow \llbracket C_2 \rrbracket (x) \\ \forall x : \iota. \llbracket C_1 \rrbracket (x) \Rightarrow \llbracket C_2 \rrbracket (x) \\ \forall x : \iota. \llbracket C_1 \rrbracket (x) \Rightarrow \llbracket C_2 \rrbracket (x) \\ \exists I = C \\ I = C \\$	I P V	$\operatorname{axiom} \ \llbracket P \rrbracket (\llbracket I \rrbracket, \llbracket V \rrbracket)$
$C_1 \equiv C_2 \qquad \forall x : \iota. \llbracket C_1 \rrbracket(x) \Leftrightarrow \llbracket C_2 \rrbracket(x) \\ \forall x : \iota. \llbracket C_1 \rrbracket(x) \Rightarrow \llbracket C_2 \rrbracket(x) \\ \exists s \text{-a } C \qquad \llbracket C \rrbracket(\llbracket I \rrbracket) \\ I_1 R I_2 \qquad \llbracket R \rrbracket(\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket) \\ I \text{-p } V \qquad \llbracket P \rrbracket(\llbracket I \rrbracket, \llbracket V \rrbracket) \\ \end{bmatrix}$ Individual Terms of type ι $i \qquad i$ $Concept \qquad \text{Formula with free variable } x : \iota$ $i \qquad i(x)$ $\top \qquad true$ $\bot \qquad false$ $C_1 \sqcup C_2 \qquad \llbracket C_1 \rrbracket(x) \vee \llbracket C_2 \rrbracket(x) \\ C_1 \sqcap C_2 \qquad \llbracket C_1 \rrbracket(x) \wedge \llbracket C_2 \rrbracket(x) \\ \forall R.C \qquad \forall y : \iota. \llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y) \\ \exists R.C \qquad \exists y : \iota. \llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) \\ \text{dom } R \qquad \exists y : \iota. \llbracket R \rrbracket(y,x) \\ \text{rng } R \qquad \exists y : \iota. \llbracket R \rrbracket(y,x) \\ \text{dom } P \qquad \exists y : I. \llbracket R \rrbracket(x,y) \qquad (T \text{ is type of } P) \\ \text{Relation} \qquad i \qquad i(x,y)$ $R_1 \cup R_2 \qquad \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 \qquad \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 \qquad \exists m : \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^* \qquad (\text{tricky, omitted})$ $\Delta_C \qquad x \doteq y \wedge \llbracket C \rrbracket(x) \\ \text{Formula with free variables } x : \iota, y : T$	F	axiom $[\![F]\!]$
$C_1 \sqsubseteq C_2$ $I \text{ is-a } C$ $I_1 R I_2$ $I P V$ $[P]([I], [V])$ $Individual$ $I Concept$	Formula	Formula without free variables
$\begin{array}{llll} I \text{ is-a } C & & & & & & & & & & & & & & & & & &$	$C_1 \equiv C_2$	$\forall x : \iota. \llbracket C_1 \rrbracket(x) \Leftrightarrow \llbracket C_2 \rrbracket(x)$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$C_1 \sqsubseteq C_2$	$\forall x : \iota. \llbracket C_1 \rrbracket(x) \Rightarrow \llbracket C_2 \rrbracket(x)$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	I is-a ${\cal C}$	$\llbracket C rbracket(\llbracket I rbracket)$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$I_1 R I_2$	$\llbracket R rbracket (\llbracket I_1 rbracket, \llbracket I_2 rbracket)$
Individual i i i Concept i Formula with free variable $x: \iota$ i $i(x)$ \vdash $true$ \vdash $false$ $C_1 \sqcup C_2$ $\llbracket C_1 \rrbracket (x) \vee \llbracket C_2 \rrbracket (x)$ $C_1 \sqcap C_2$ $\llbracket C_1 \rrbracket (x) \wedge \llbracket C_2 \rrbracket (x)$ $\forall R.C$ $\forall y: \iota. \llbracket R \rrbracket (x,y) \Rightarrow \llbracket C \rrbracket (y)$ $\exists R.C$ $\exists y: \iota. \llbracket R \rrbracket (x,y) \wedge \llbracket C \rrbracket (y)$ $dom R$ $\exists y: \iota. \llbracket R \rrbracket (x,y)$ $rng R$ $\exists y: \iota. \llbracket R \rrbracket (y,x)$ $dom P$ $\exists y: T. \llbracket P \rrbracket (x,y)$ $(T \text{ is type of } P)$ Relation i Formula with free variables $x: \iota, y: \iota$ i $i(x,y)$ $R_1 \cup R_2$ $\llbracket R_1 \rrbracket (x,y) \wedge \llbracket R_2 \rrbracket (x,y)$ $R_1; R_2$ $\exists m: \iota. \llbracket R_1 \rrbracket (x,m) \wedge \llbracket R_2 \rrbracket (m,y)$ R^-1 $\llbracket R \rrbracket (y,x)$ R^* (tricky, omitted) Δ_C $x \doteq y \wedge \llbracket C \rrbracket (x)$ Property of type T Formula with free variables $x: \iota, y: T$	I P V	
$\begin{array}{ c c c }\hline \text{Concept} & \text{Formula with free variable } x:\iota\\ i & i(x)\\ \hline & true\\ \hline \bot & false\\ \hline C_1 \sqcup C_2 & \llbracket C_1 \rrbracket(x) \vee \llbracket C_2 \rrbracket(x)\\ \hline C_1 \sqcap C_2 & \llbracket C_1 \rrbracket(x) \wedge \llbracket C_2 \rrbracket(x)\\ \forall R.C & \forall y:\iota.\llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y)\\ \hline \exists R.C & \exists y:\iota.\llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y)\\ \hline \text{dom } R & \exists y:\iota.\llbracket R \rrbracket(y,x)\\ \hline \text{rng } R & \exists y:\iota.\llbracket R \rrbracket(y,x)\\ \hline \text{dom } P & \exists y:T.\llbracket P \rrbracket(x,y) & (T \text{ is type of } P)\\ \hline \hline \text{Relation} & \text{Formula with free variables } x:\iota,y:\iota\\ i & i(x,y)\\ \hline R_1 \cup R_2 & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y)\\ \hline R_1;R_2 & \exists m:\iota.\llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(m,y)\\ \hline R^{-1} & \llbracket R \rrbracket(y,x)\\ \hline R^* & \text{(tricky, omitted)}\\ \hline \Delta_C & x \doteq y \wedge \llbracket C \rrbracket(x)\\ \hline \hline \\ \hline \text{Property of type } T & \text{Formula with free variables } x:\iota,y:T\\ \hline \end{array}$	Individual	Terms of type ι
$i(x)$ $T \qquad true$ $\bot \qquad false$ $C_1 \sqcup C_2 \qquad \llbracket C_1 \rrbracket(x) \vee \llbracket C_2 \rrbracket(x)$ $\forall R.C \qquad \forall y: \iota. \llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y)$ $\exists R.C \qquad \exists y: \iota. \llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y)$ $dom R \qquad \exists y: \iota. \llbracket R \rrbracket(x,y)$ $rng R \qquad \exists y: \iota. \llbracket R \rrbracket(x,y) \qquad (T \text{ is type of } P)$ $Relation \qquad Formula with free variables x: \iota, y: \iota i \qquad i(x,y) R_1 \cup R_2 \qquad \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) R_1; R_2 \qquad \exists m: \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) R^{-1} \qquad \llbracket R \rrbracket(y,x) R^* \qquad (\text{tricky, omitted}) \Delta_C \qquad x \doteq y \wedge \llbracket C \rrbracket(x) Property of type T \qquad Formula with free variables x: \iota, y: T$	i	i
$ \begin{array}{lll} \top & & true \\ \bot & & false \\ C_1 \sqcup C_2 & & \llbracket C_1 \rrbracket(x) \vee \llbracket C_2 \rrbracket(x) \\ C_1 \sqcap C_2 & & \llbracket C_1 \rrbracket(x) \wedge \llbracket C_2 \rrbracket(x) \\ \forall R.C & \forall y: \iota.\llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y) \\ \exists R.C & \exists y: \iota.\llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) \\ \operatorname{dom} R & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ \operatorname{rng} R & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ \operatorname{dom} P & \exists y: T.\llbracket P \rrbracket(x,y) & (T \text{ is type of } P) \\ \hline \text{Relation} & \text{Formula with free variables } x: \iota, y: \iota \\ i & i(x,y) \\ R_1 \cup R_2 & & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 & & \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 & \exists m: \iota.\llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^{-1} & & \llbracket R \rrbracket(y,x) \\ R^* & & \text{(tricky, omitted)} \\ \Delta_C & x \doteq y \wedge \llbracket C \rrbracket(x) \\ \hline \text{Property of type } T & \text{Formula with free variables } x: \iota, y: T \\ \hline \end{array} $	Concept	Formula with free variable $x:\iota$
$\begin{array}{lll} \bot & false \\ C_1 \sqcup C_2 & \llbracket C_1 \rrbracket(x) \vee \llbracket C_2 \rrbracket(x) \\ C_1 \sqcap C_2 & \llbracket C_1 \rrbracket(x) \wedge \llbracket C_2 \rrbracket(x) \\ \forall R.C & \forall y: \iota. \llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y) \\ \exists R.C & \exists y: \iota. \llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) \\ \text{dom } R & \exists y: \iota. \llbracket R \rrbracket(y,x) \\ \text{rng } R & \exists y: \iota. \llbracket R \rrbracket(y,x) \\ \text{dom } P & \exists y: T. \llbracket P \rrbracket(x,y) & (T \text{ is type of } P) \\ \hline \text{Relation} & \text{Formula with free variables } x: \iota, y: \iota \\ i & i(x,y) \\ R_1 \cup R_2 & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 & \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 & \exists m: \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^{-1} & \llbracket R \rrbracket(y,x) \\ R^* & (\text{tricky, omitted}) \\ \Delta_C & x \doteq y \wedge \llbracket C \rrbracket(x) \\ \hline \text{Property of type } T & \text{Formula with free variables } x: \iota, y: T \\ \hline \end{array}$		i(x)
$C_1 \sqcup C_2$ $C_1 \sqcap C_2$ $\forall R.C$ $\exists R.C$ $dom R$ $rng R$ $dom P$ $Relation$ i $i(x,y)$ $R_1 \cup R_2$ $R_1; R_2$ $R_1; R_2$ $R_1; R_2$ $R_1; R_2$ R_2 $R_1; R_2$ $R_1 \sqcup R_2$ $R_1 \sqcup R_2$ $R_1; R_2$ $R_1 \sqcup R_2$ $R_1; R_2$ $R_1 \sqcup R_2$ $R_1 \sqcup R_2$ $R_1 \sqcup R_2$ $R_1; R_2$ $R_1 \sqcup R_2$ $R_1; R_2$ $R_1 \sqcup R_2$ $R_1 \sqcup R_2$ $R_1 \sqcup R_2$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1 \sqcup R_2 \sqcup [R_1] \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1 \sqcup R_2 \sqcup [R_1] \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup R_1 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_1; R_2 \sqcup (x,y) \wedge [R_2] \sqcup (x,y)$ $R_1; R_2 \sqcup (x,y) \wedge $	Τ	true
$ \begin{array}{lll} & & & & & & & & & & & & \\ \forall R.C & & & \forall y: \iota.\llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y) \\ & \exists R.C & & \exists y: \iota.\llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) \\ & & & & \exists y: \iota.\llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) \\ & & & & \exists y: \iota.\llbracket R \rrbracket(x,y) \\ & & & & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ & & & & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ & & & & & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ & & & & & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ & & & & & \exists y: \iota.\llbracket R \rrbracket(y,x) \\ & & & & & \exists y: \iota.\llbracket R \rrbracket(x,y) & (T \text{ is type of } P) \\ \hline & & & & & & \exists y: \iota.\llbracket R \rrbracket(x,y) & (T \text{ is type of } P) \\ \hline & & & & & & & \vdots \\ & & & & & & \vdots \\ & & & &$	\perp	false
$ \forall R.C \qquad \forall y: \iota.\llbracket R \rrbracket(x,y) \Rightarrow \llbracket C \rrbracket(y) $ $ \exists R.C \qquad \exists y: \iota.\llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) $ $ dom R \qquad \exists y: \iota.\llbracket R \rrbracket(x,y) $ $ rng R \qquad \exists y: \iota.\llbracket R \rrbracket(y,x) $ $ dom P \qquad \exists y: T.\llbracket P \rrbracket(x,y) \qquad (T \text{ is type of } P) $ $ Relation \qquad Formula \text{ with free variables } x: \iota, y: \iota $ $ i \qquad \qquad i \qquad i \qquad i \qquad i \qquad i \qquad i \qquad $	$C_1 \sqcup C_2$	$[\![C_1]\!](x) \vee [\![C_2]\!](x)$
$ \exists R.C \\ \operatorname{dom} R \\ \exists y : \iota.\llbracket R \rrbracket(x,y) \wedge \llbracket C \rrbracket(y) \\ \operatorname{rng} R \\ \operatorname{dom} P \\ \exists y : \iota.\llbracket R \rrbracket(y,x) \\ \operatorname{dom} P \\ \exists y : I.\llbracket P \rrbracket(x,y) \\ \exists y : I.\llbracket P \rrbracket(x,y) \\ T \text{ is type of } P) $ Relation Formula with free variables $x : \iota, y : \iota$ if $i(x,y)$ if $i(x,$	$C_1 \sqcap C_2$	
$\begin{array}{lll} \operatorname{dom} R & \exists y : \iota. \llbracket R \rrbracket(x,y) \\ \operatorname{rng} R & \exists y : \iota. \llbracket R \rrbracket(y,x) \\ \operatorname{dom} P & \exists y : T. \llbracket P \rrbracket(x,y) & (T \text{ is type of } P) \\ \hline \text{Relation} & \text{Formula with free variables } x : \iota, y : \iota \\ i & i(x,y) \\ R_1 \cup R_2 & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 & \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 & \exists m : \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^{-1} & \llbracket R \rrbracket(y,x) \\ R^* & (\text{tricky, omitted}) \\ \Delta_C & x \doteq y \wedge \llbracket C \rrbracket(x) \\ \hline \text{Property of type } T & \text{Formula with free variables } x : \iota, y : T \\ \hline \end{array}$	$\forall R.C$	$\forall y : \iota.[\![R]\!](x,y) \Rightarrow [\![C]\!](y)$
$\begin{array}{lll} \operatorname{rng} R & \exists y : \iota. \llbracket R \rrbracket(y,x) \\ \operatorname{dom} P & \exists y : T. \llbracket P \rrbracket(x,y) & (T \text{ is type of } P) \\ \hline \text{Relation} & \operatorname{Formula with free variables } x : \iota, y : \iota \\ i & i(x,y) \\ R_1 \cup R_2 & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 & \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 & \exists m : \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^{-1} & \llbracket R \rrbracket(y,x) \\ R^* & (\operatorname{tricky, omitted}) \\ \Delta_C & x \doteq y \wedge \llbracket C \rrbracket(x) \\ \hline \text{Property of type } T & \operatorname{Formula with free variables } x : \iota, y : T \\ \hline \end{array}$	$\exists R.C$	$\exists y : \iota. \llbracket R \rrbracket(x, y) \land \llbracket C \rrbracket(y)$
$\begin{array}{c c} \operatorname{dom} P & \exists y: T. \llbracket P \rrbracket (x,y) & (T \text{ is type of } P) \\ \hline \text{Relation} & \text{Formula with free variables } x: \iota, y: \iota \\ i & i(x,y) \\ \hline R_1 \cup R_2 & \llbracket R_1 \rrbracket (x,y) \vee \llbracket R_2 \rrbracket (x,y) \\ \hline R_1 \cap R_2 & \llbracket R_1 \rrbracket (x,y) \wedge \llbracket R_2 \rrbracket (x,y) \\ \hline R_1; R_2 & \exists m: \iota. \llbracket R_1 \rrbracket (x,m) \wedge \llbracket R_2 \rrbracket (m,y) \\ \hline R^{-1} & \llbracket R \rrbracket (y,x) \\ \hline R^* & (\text{tricky, omitted}) \\ \hline \Delta_C & x \doteq y \wedge \llbracket C \rrbracket (x) \\ \hline \\ \hline \text{Property of type } T & \text{Formula with free variables } x: \iota, y: T \\ \hline \end{array}$	$\operatorname{dom} R$	$\exists y: \iota. \llbracket R \rrbracket (x,y)$
Relation i Formula with free variables $x:\iota,y:\iota$ i $i(x,y)$ $R_1\cup R_2$ $[\![R_1]\!](x,y)\vee [\![R_2]\!](x,y)$ $R_1\cap R_2$ $[\![R_1]\!](x,y)\wedge [\![R_2]\!](x,y)$ $R_1;R_2$ $\exists m:\iota.[\![R_1]\!](x,m)\wedge [\![R_2]\!](m,y)$ R^{-1} $[\![R]\!](y,x)$ R^* (tricky, omitted) Δ_C $x\doteq y\wedge [\![C]\!](x)$ Property of type T Formula with free variables $x:\iota,y:T$	$\operatorname{rng} R$	$\exists y: \iota. \llbracket R \rrbracket (y, x)$
$ \begin{array}{lll} i & & i(x,y) \\ R_1 \cup R_2 & & & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 & & & \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 & & \exists m : \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^{-1} & & & \llbracket R \rrbracket(y,x) \\ R^* & & & (\text{tricky, omitted}) \\ \Delta_C & & x \doteq y \wedge \llbracket C \rrbracket(x) \\ \hline \text{Property of type T} & \text{Formula with free variables $x : } \iota,y : T \\ \end{array} $		
$\begin{array}{lll} R_1 \cup R_2 & & & \llbracket R_1 \rrbracket(x,y) \vee \llbracket R_2 \rrbracket(x,y) \\ R_1 \cap R_2 & & & \llbracket R_1 \rrbracket(x,y) \wedge \llbracket R_2 \rrbracket(x,y) \\ R_1; R_2 & & \exists m : \iota. \llbracket R_1 \rrbracket(x,m) \wedge \llbracket R_2 \rrbracket(m,y) \\ R^{-1} & & & \llbracket R \rrbracket(y,x) \\ R^* & & & \text{(tricky, omitted)} \\ \Delta_C & & x \doteq y \wedge \llbracket C \rrbracket(x) \\ \hline \text{Property of type T} & \text{Formula with free variables $x : \iota,y : T} \end{array}$	Relation	Formula with free variables $x : \iota, y : \iota$
$R_{1} \cap R_{2}$ $R_{1}; R_{2}$ R^{-1} R^{*} $Droperty of type T$ $R_{1} \cap R_{2}$ $R_{2} \cap R_{2} \cap R_{2}$	i	i(x,y)
$R_1; R_2 \\ R^{-1} \\ R^* \\ \Delta_C \\ \hline Property of type T $	$R_1 \cup R_2$	
$ \begin{array}{ll} R^{-1} & & \llbracket R \rrbracket (y,x) \\ R^* & & \text{(tricky, omitted)} \\ \Delta_C & & x \doteq y \land \llbracket C \rrbracket (x) \\ \hline \text{Property of type T} & \text{Formula with free variables $x:\iota,y:T$} \\ \end{array} $	$R_1 \cap R_2$	
$\begin{array}{ccc} R^* & \text{(tricky, omitted)} \\ \underline{\Delta_C} & x \doteq y \land \llbracket C \rrbracket(x) \\ \hline \text{Property of type } T & \text{Formula with free variables } x : \iota, y : T \end{array}$		
$\begin{array}{c c} \Delta_C & x \doteq y \land \llbracket C \rrbracket(x) \\ \hline \text{Property of type } T & \text{Formula with free variables } x : \iota, y : T \end{array}$		
Property of type T Formula with free variables $x: \iota, y: T$	R^*	
$i \hspace{1cm} \mid i(x,y)$	Property of type T	
	i	i(x,y)

Figure 6.4: Interpretation Function for BOL into SFOL

The consequence closure of SFOL, using the usual semantics of SFOL, induces the desired consequence closure for BOL.

Definition 6.9 (Consequence Closure). We say that a BOL-statement F is a consequence of an ontology O if $\llbracket F \rrbracket$ is an SFOL-theorem of P, $\llbracket O \rrbracket$.

Example 6.10. We interpret the example ontology from Ex. 2.3. Excluding the semantic prefix, it results in

Example 6.11 (Compositionality). The interpretation of BOL is compositional.

For example, consider the case of composition of relations:

$$[R_1; R_2] = \exists m : \iota.[R_1](x, m) \land [R_2](m, y)$$

Here we have an expression $E(e_1,\ldots,e_n)$ with n=2 and E is the ;-operator mapping $(e_1,e_2)\mapsto e_1;e_2$, i.e., R_1 and R_2 are the direct subexpressions of $R_1;R_2$. The semantics is a relatively complicated FOL-formula, but it only depends on $[\![R_1]\!]$ and $[\![R_2]\!]$ — everything else is fixed. We have $[\![;]\!]=(p_1,p_2)\mapsto \exists m:\iota.p_1(x,m)\wedge p_2(m,y)$, i.e., the interpretation of the ;-operator is the function that maps two predicates p_1,p_2 to the formula $\exists m:\iota.p_1(x,m)\wedge p_2(m,y)$. Then we have

$$[R_1; R_2] = [:]([R_1], [R_2]).$$

Example 6.12 (Non-Compositional Translation via Custom Semantic Prefix). In Fig. 6.4, we omitted the case for the translative closure. That was because it is not possible to translate it compositionally into FOL. We can only do it non-compositionally with a custom semantic prefix:

We define the FOL-interpretation of an ontology O by $\llbracket O \rrbracket = P_O, \llbracket O \rrbracket$, where P_O is a custom semantic prefix. P_O is different for every ontology O and is defined as follows:

- 1. We scan through O and collect all occurrences of R^* for any (not necessarily atomic) relation R.
- 2. P_O contains the following declarations for each R:
 - A binary predicate symbol $C_R \subseteq i \times i$. Note that R may be a complex expression; so we have to generate a fresh name C_R here.
 - The axiom $\forall x : \iota, y : \iota R(x, y) \Rightarrow C_R(x, y)$, i.e., C_R extends R.
 - The axiom $\forall x : \iota, y : \iota, z : \iota . C_R(x, y) \land C_R(y, z) \Rightarrow C_R(x, z)$, i.e., C_R is transitive.
- 3. We add the case $[R^*] = C_R(x, y)$ to the interpretation function.

Intuitively, every occurrence of the *-operator is removed from the language and replaced with a fresh name that is axiomatized to have the needed properties. All of these axioms are added to the semantic prefix.

6.4 Concretized Semantics

We give an alternative semantics using a semantic language for concrete data. Specifically we use the database language SQL.

6.4.1 An SQL-Inspired Basic Database Language

We give an SQL-like database language as a formal system.

Definition 6.13. Fig. 6.5 gives the context-free grammar. The vocabulary symbol is S. The expression symbols are T, R, V, and F.

6.4.2 Semantics

Even though this is a very different knowledge aspect, the general principles of the semantics are the same: Every BOL-declaration is translated to an SQL declaration, and ontologies are translated declaration-wise. For every kind of complex expression, there is one inductive function mapping BOL-expressions to SQL-expressions.

In SQL, we can nicely see the difference between declarations and expressions: the former are translated to side effect-ful statements, the latter to side effect-free queries.

Definition 6.14 (Concretized Semantic of BOL). The semantic prefix consists of the following SQL-statements

- \bullet a type ID of identifiers (if not already supported anyway by the underlying database)
- declarations of all base types and values of BOL (if not already supported anyway by the underlying database)
- TABLE individuals (id: ID, name: string), where the id field is unique and automatically generated when inserting values

Every BOL-ontology O is interpreted as the sequence P, [0] of SQL statements, where [0] is defined in Fig. 6.6.

```
Vocabularies: Schemas
  ::= D^*
Declarations
D ::= \mathbf{TABLE} \text{ ID } TT
                                                  table
     INSERT RINTO ID
                                                  row in a table
table types
TT ::= CT^*
                                                  list of column types
CT ::= ID : Y
                                                  column type
table expressions (i.e., table-valued queries)
T ::= ID
                                                  atomic tables
       JOIN T^*
                                                 join of tables
        UNION T^*
                                                  union of tables
       INTER T^*
                                                  intersection of tables
        SELECT DISTINCT CL FROM T
                                                  selection of columns
        T WHERE F
                                                  selection of rows
        T\,\mathbf{AS}\,\mathtt{ID}
                                                  prefix for column names
       T AGGREGATE CA^*
                                                  aggregation
CA := \mathbf{COUNT}(\mathtt{ID})
                                                  count values in a column
       \mathbf{MAX}(\mathtt{ID})
                                                  maximum value in a column
                                                  minimum, sum, etc.
CL := * | (ID | ID AS ID)^*
                                                 list of column names
row expressions (i.e., single row-valued queries)
R ::= VALUES(CD^*)
                                                  explicit row
      T
                                                  table, but may only contain one row
CD ::= ID = V
                                                  column definition
cell expressions (i.e., single column, single row-valued queries)
   ::=R
                                                  row, but may only contain one column
       {\tt ID}(V^*)
                                                  built-in function symbol applied to values
       (base\,values)
                                                  as for BOL
base types
Y ::= (base types)
                                                  as for BOL
formulas
   ::=V
                                                  boolean value
       V = V
                                                  equality of values
        R IN T
                                                  containment of rows in tables
        R NOT IN T
                                                  opposite of containment
                                                  boolean operators
Identifiers
ID ::= alphanumeric string
```

Figure 6.5: Syntax of SQL

BOL Syntax X	Semantics $[X]$ in SQL
ontology	SQL schema (list of statements)
D_1,\ldots,D_n	$\llbracket D_1 rbracket, \ldots, \llbracket D_n rbracket$
BOL declaration $(C, R, P \text{ atomic})$	SQL statement
individual i	INSERT VALUES (name="i") INTO individuals
$\mathbf{concept}i$	TABLE i (id: ID)
${\bf relation}i$	TABLE i (subject: ID, object: ID)
$\mathbf{property}i:Y$	TABLE i (subject: ID, value: Y)
I is-a C	INSERT VALUES (id= $\llbracket I \rrbracket$) INTO C
$I_1 \mathrel{R} I_2$	INSERT VALUES (subject= $\llbracket I_1 \rrbracket$, object= $\llbracket I_2 \rrbracket$) INTO R
I P V	INSERT VALUES (subject= $[I]$, value= V) INTO P
F	consistency check, consequence closure (omitted)
Formula	Query that returns empty result iff formula is true
$C_1 \equiv C_2$	$(\llbracket C_1 \rrbracket \setminus \llbracket C_2 \rrbracket)$ UNION $(\llbracket C_2 \rrbracket \setminus \llbracket C_1 \rrbracket)$
$C_1 \sqsubseteq C_2$	$ig \ \llbracket C_1 rbracket \setminus \llbracket C_2 rbracket$
I is-a C	$\llbracket I rbracket{}{ m IN} \ \llbracket C rbracket{}{ m IN}$
$I_1 \mathrel{R} I_2$	$(\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket)$ IN R
I P V	$(\llbracket I \rrbracket, V)$ IN P
Individual	an identifier from the table individuals
i	SELECT id FROM individuals WHERE name="i"
Concept	SQL table expression for $(id:ID)$
i	SELECT * FROM i
Т	individuals
\perp	individuals WHERE false
$C_1 \sqcup C_2$	$\llbracket C_1 rbracket$ UNION $\llbracket C_2 rbracket$
$C_1 \sqcap C_2$	$\llbracket C_1 rbracket$ INTERSECT $\llbracket C_2 rbracket$
$\forall R.C$	individuals \ (SELECT subject FROM $[\![R]\!]$ WHERE object NOT IN $[\![C]\!]$)
$\exists R.C$	SELECT DISTINCT subject FROM $[\![R]\!]$, $[\![C]\!]$ WHERE object=id
$\operatorname{dom} R$	SELECT DISTINCT subject FROM $[\![R]\!]$
$\operatorname{rng} R$	SELECT DISTINCT object FROM $[\![R]\!]$
$\operatorname{dom} P$	SELECT DISTINCT subject FROM $\llbracket P \rrbracket$
Relation	SQL table expression for (subject: ID, object: ID)
i	SELECT * FROM i
$R_1 \cup R_2$	$[R_1]$ UNION $[R_2]$
$R_1 \cap R_2$	$\llbracket R_1 rbracket$ INTERSECT $\llbracket R_2 rbracket$
$R_1; R_2$	SELECT DISTINCT l.subject, r.object FROM $[\![R_1]\!]$ AS l, $[\![R_2]\!]$ AS r
1	WHERE l.object = $r.subject$
R^{-1}	SELECT object AS subject, subject AS object FROM $[\![R]\!]$
R^*	(tricky, omitted)
Δ_C	SELECT id AS subject, id AS object FROM $[\![C]\!]$
Property of type Y	SQL table expression for (subject: ID, value: Y)
i	SELECT * FROM i

Using the abbreviations: $S \setminus T = \text{SELECT}$ id FROM S WHERE id NOT IN T S, T = JOIN S, T

Figure 6.6: Interpretation Function for BOL into SQL

Remark 6.15 (Limitations). Our interpretation of BOL in SQL is restricted to assertions using only atomic expressions. For example, in the case for I is-a C, we assume that I and C are names. Thus, we have already created an individual for I and a table for C, and we can thus insert the former into the latter. The general case would be more complicated but is much less important in practice. But other expressions very quickly become more difficult.

The interpretation of formulas into SQL is less obvious because SQL is not a logic and therefore does not define a consequence closure. Thus, we can only use axioms for consistency checks in SQL. But that requires first carrying out an explicit consequence closure that adds all implied assertions to the database.

Example 6.16. We interpret the example ontology from Ex. 2.3. Excluding the semantic prefix, the entity declarations and assertions result in the following

```
INSERT VALUES (name="FlorianRabe") INTO individuals
INSERT VALUES (name="WuV") INTO individuals
TABLE person (id: ID)
TABLE male (id: ID)
TABLE instructor (id: ID)
TABLE course (id: ID)
TABLE teach (subject: ID, object: ID)
TABLE creditValue (subject: ID, value: float)
INSERT VALUES (id=2) INTO course
INSERT VALUES (subject=1, object=2) INTO teach
INSERT VALUES (subject=1, value=7.5) INTO creditValue
```

Here we assume that inserting into the table individuals has automatically assigned the ids 1 and 2 to our two individuals.

The concept assertion about FlorianRabe using \sqcap cannot be handled by this semantics. Therefore, we skip that assertion. The two missing assertions

```
INSERT VALUES (id=1) INTO instructor INSERT VALUES (id=1) INTO male
```

must then be provided by performing the consequence closure.

Moreover, the axioms result in the following consistency checks, i.e., queries that should be empty:

```
SELECT * FROM male \ SELECT * FROM person
SELECT * FROM instructor \ SELECT * FROM person
SELECT * FROM (SELECT DISTINCT subject FROM teach) \ SELECT * FROM instructor
SELECT * FROM (SELECT DISTINCT object FROM teach) \ SELECT * FROM course
(SELECT * FROM (SELECT DISTINCT subject FROM creditValue) \ SELECT * FROM course)
    UNION (SELECT * FROM course \ SELECT DISTINCT subject FROM creditValue)
SELECT * FROM course \
    (SELECT DISTINCT subject
    FROM (SELECT object AS subject, subject AS object FROM teach), instructor
    WHERE object=id)
```

Some of these checks will only succeed after performing the consequence closure. In particular, the table person misses the entry 1 for the individual FlorianRabe because the assertion FlorianRabe is-a person is only present as a consequence

6.5 Computational Semantics

We give an alternative semantics using computation, i.e., by using a programming language as the semantic language. Specifically, we use the programming language Scala.

6.5.1 A Scala-Inspired Basic Programming Language

We give a simple programming language as a language system.

Definition 6.17. Fig. 6.7 gives the context-free grammar. The vocabulary symbol is P. The expression symbols are Y, V, and F.

```
Vocabularies: Programs
P ::= D^*
Declarations
D ::= class ID[X^*] extends ID^*\{d^*\}
                                           class definition
    object ID extends ID*\{d^*\}
                                           object definition
d ::= \mathbf{val} \ \mathtt{ID} : Y[=T]
                                           immutable field in a class/object, possibly abstract
    \mathbf{var} \; \mathtt{ID} : Y = T
                                           mutable field in a class, with initial value
 Type expressions
Y ::= ID[X^*]
                                           atomic type (class) applied to type arguments
       ID
                                           built-in type (booleans, int, etc.)
       X
                                           type variable
       Y => Y
                                           function types
 Term expressions
T ::= ID
                                           atomic value (class, value, variable)
       ID
                                           built-in value (boolean operators, etc.)
       T.ID
                                           field access in an object
       T.ID = T
                                           assignment to a mutable field in an object
       ID = T
                                           assignment to a local variable
       T:Y
                                           instance check
       new ID \{d^*\}
                                           new instance of class
       T \doteq_{Y} T
                                           equality of terms at a type
       (\mathtt{ID}:Y) => T
                                           function
       T(T^*)
                                           function applied to values
       \{T^*\}
                                           sequencing (;-operator)
                                           local declaration
       d
       \mathbf{if}(T) T \mathbf{else} T
                                           if-then-else
       while (T)T
                                           while-loop
 Formula expressions
F ::= T
                                           terms of boolean type
Identifiers
ID ::= alphanumeric string
```

Figure 6.7: Syntax of BPL

6.5.2 Semantics

Again, the general principles are the same: Every BOL-declaration is translated to a Scala-declaration, and ontologies are translated declaration-wise to Scala-programs. For every kind of complex expression, there is one inductive function mapping BOL-expressions to Scala-objects.

Definition 6.18 (Computational Semantic of BOL). The semantic prefix consists of the following Scala statements

- classes for all BOL-base types and values for them (if not already present in Scala)
- classes for individuals and hash sets of objects:

```
import scala.collection.mutable.HashSet
val individuals = new HashSet[String]
```

Every BOL-ontology O is interpreted as the Scala program P, $[\![O]\!]$, where $[\![O]\!]$ is defined in Fig. 6.8.

Remark 6.19 (Scala Syntax). In Scala, val x = e evaluates e and stores the result in x. $\{d_1; \ldots; d_n\}$ is evaluated by executing all d_i in order and returning the result of d_n .

(A,B) is the product type $A \times B$ with pairing operator (x,y) and projection functions $\bot 1$ and $\bot 2$. $x \Rightarrow F(x)$ is $\lambda x.F(x)$.

The class HashSet is part of the standard library and offers function += and -= to add/remove elements, contains to test elementhood, and forall, foreach to quantify/iterate over elements.

Types of variables are inferred if omitted.

Remark 6.20 (Limitations). Our interpretation of BOL in Scala has similar problems as the one in SQL. We restrict entities in assertions to be atomic. And we assume that all assertions implied by the consequence closure have already been obtained and added to the ontology.

Example 6.21. We interpret the example ontology from Ex. 2.3. Excluding the semantic prefix, the entity declarations and assertions result in the following

```
individuals += "FlorianRabe"
individuals += "WuV"

val person = new HashSet[String]
val male = new HashSet[String]
val person = new HashSet[String]
val course = new HashSet[String]
val teach = new new HashSet[(String, String)]
val creditValue = new HashSet[(String, float)]
course += WuV
teach += ("FlorianRabe", WuV)
creditValue += (WuV, 7.5)
Like for SQL, the two statements
instructor += "FlorianRabe"
male += "FlorianRabe"
```

must be obtained by consequence closure because we cannot handle the \sqcap assertion. Note that we could easily compute the hash set <code>instructor.diff(male)</code> and add to it. But that would not add anything to the two constituent sets.

If we thing of the axioms as consistency checks, we can translate them to assertions, i.e., Boolean expressions that must be true. We only give some examples:

BOL Syntax X	Semantics $[X]$ in Scala
ontology	Scala program
D_1,\ldots,D_n	object $Program \{P, \llbracket D_1 \rrbracket, \dots, \llbracket D_n \rrbracket \}$
BOL declaration $(C, R, P \text{ atomic})$	Scala declaration
individual i	$val\ i = "i"$; individuals $+= i$
$\mathbf{concept}i$	val i = new HashSet[String]
${f relation}i$	val i = new HashSet[(String,String)]
$\mathbf{property}\ i:T$	val i = new HashSet[(String,T)]
I is-a C	$\llbracket C \rrbracket + = \llbracket I \rrbracket$
$I_1 \mathrel{R} I_2$	$ig[ar{R}ig] += ar{(}ig[ar{I}_1ig], ig[ar{I}_2ig])$
I P V	$\llbracket P \rrbracket += (\llbracket I \rrbracket, \llbracket V \rrbracket)$
F	assertions, consequence closure (omitted)
Formula	Program that evaluates the formula to a Boolean
$C_1 \equiv C_2$	$\{ \text{val } c1 = [C_1]; \text{ val } c2 = [C_2]; $
	$c_1.\text{forall}(x \Rightarrow c_2.\text{contains}(x)) \&\& c_2.\text{forall}(x \Rightarrow c_1.\text{contains}(x))$
$C_1 \sqsubseteq C_2$	$\{ \text{val } c1 = \llbracket C_1 \rrbracket; \text{val } c2 = \llbracket C_2 \rrbracket; c_1.\text{forall}(\mathbf{x} \Rightarrow c_2.\text{contains}(\mathbf{x})) \}$
I is-a C	$\llbracket C \rrbracket$.contains($\llbracket I \rrbracket$)
$I_1 \mathrel{R} I_2$	$\llbracket R \rrbracket$.contains $((\llbracket I_1 \rrbracket, \llbracket I_2 \rrbracket))$
I P V	$\llbracket P \rrbracket$.contains $((\llbracket I \rrbracket, \llbracket V \rrbracket))$
Individual	String object
i	$\mid i \mid$
Concept	HashSet[String] object
i	$\mid i \mid$
Т	individuals
Τ	HashSet[String].empty()
$C_1 \sqcup C_2$	$\llbracket C_1 rbracket$.union($\llbracket C_2 rbracket$)
$C_1 \sqcap C_2$	$\llbracket C_1 rbracket{ m .inter}(\llbracket C_2 rbracket{ m)}$
$\forall R.C$	$\{\text{val c} = [\![C]\!]; \text{ val r} = [\![R]\!]; \text{ val e} = \text{individuals.clone};$
	$r.foreach(x \Rightarrow if (!c.contains(x2)) e = x1); e$
$\exists R.C$	$\{ \text{val c} = [\![C]\!]; \text{ val r} = [\![R]\!]; \text{ val e} = \text{new HashSet[String]}; $
	$r.foreach(x \Rightarrow if (c.contains(x2)) e += x1); e$
$\operatorname{dom} R$	$\{ val \ c = new \ HashSet[String]; \ \llbracket R \rrbracket.foreach(x \Rightarrow c += x._1); \ c \}$
$\operatorname{rng} R$	$\{ val c = new HashSet[String]; [R].foreach(x \Rightarrow c += x2); c \}$
$\operatorname{dom} P$	$\{ val \ c = new \ HashSet[String]; \ \llbracket P \rrbracket.foreach(x \Rightarrow c += x._1); \ c \}$
Relation	HashSet[(String,String)] object
i	$\mid i \mid$
$R_1 \cup R_2$	$\llbracket R_1 rbracket{ m .union}(\llbracket R_2 rbracket)$
$R_1 \cap R_2$	$\llbracket R_1 rbracket{ } \operatorname{inter}(\llbracket R_2 rbracket{ } $
$R_1; R_2$	$\{\text{val r1} = [R_1]; \text{val r2} = [R_2]; \text{val e} = \text{new HashSet}[(\text{String}, \text{String})];$
_ 1	$r1.foreach(x \Rightarrow r2.foreach(y \Rightarrow if (x2 == y1) e += (x1,y2))); e$
R^{-1}	$\{ val \ r = new \ HashSet[(String,String)]; \ \llbracket R \rrbracket. for each(x \Rightarrow r \ += \ (x2,x1)); \ r \}$
R^*	$\{\text{val } \mathbf{r} = \llbracket R \rrbracket; \text{ var } \mathbf{e} = \llbracket \Delta_{\top} \rrbracket; \text{ while } (\text{!equal}(\text{comp}(\mathbf{e},\mathbf{r}), \mathbf{e})) \ \{\mathbf{e} = \text{comp}(\mathbf{e},\mathbf{r})\}; \mathbf{e}\}$
Δ_C	$\{\text{val } \mathbf{r} = \text{new HashSet}[(\text{String}, \text{String})]; \ \llbracket C \rrbracket. \text{foreach}(\mathbf{x} \Rightarrow \mathbf{r} += (\mathbf{x}, \mathbf{x})); \ \mathbf{r}\}$
Property of type T	HashSet[(String,T)] object
i	$\mid i \mid$

Figure 6.8: Interpretation Function for BOL into Scala (assuming functions equal and comp on relations)

```
r
}
val e = new HashSet[String];
r.foreach(x ⇒ if (c.contains(x._2)) e += x._1);
e
};
c1.forall(x ⇒ c2.contains(x))
}
```

6.6 Narrative Semantics

We give an alternative semantics using narration, i.e., by using a natural language as the semantic language. Specifically, we use the natural language English.

Again, the general principles are the same: Every BOL-declaration is translated to an English sentence, and ontologies are translated declaration-wise to English texts. For every kind of complex expression, there is one inductive function mapping BOL-expressions to English phrases.

Definition 6.22 (Narrative Semantic of BOL). The semantic prefix consists of English statements explaining

- the base types of BOL (if they are not universally known),
- that we rely on a lexicon to correctly form plurals (indicated by -s) and verb forms (indicated by -s, -ing, -ed).

Every BOL-ontology O is interpreted as the English text P, $[\![O]\!]$, where $[\![O]\!]$ is defined in Fig. 6.9.

Natural language defines a consequence closure by appealing to consequence in natural language. That is well-defined as long as we express ourselves precisely enough.

Definition 6.23 (Consequence Closure). We say that a BOL-statement F is a consequence of an ontology O if $\llbracket F \rrbracket$ is a consequence of P, $\llbracket O \rrbracket$.

Example 6.24. We interpret the example ontology from Ex. 2.3. Excluding the semantic prefix and the lexicon lookup, it results in the following text:

FlorianRabe is a proper noun.

WuV is a proper noun.

person is a common noun.

male is a common noun.

instructor is a common noun.

course is a common noun.

teach is a transitive verb.

creditValue is a common noun for a property that can take float-values.

FlorianRabe is a instructor and is a male.

WuV is a course.

FlorianRabe teachs WuV.

WuV has creditValue 7.5.

everything that is a male also is a person.

everything that is a instructor also is a person.

everything that teachs is a instructor.

everything that has something that teachs them is a instructor.

has some creditValue is the same as is a course.

everything that is a course also has something that teachs them and who is a instructor.

This English is very clunky of course. Multiple tweaks would be needed to get the grammar right:

- It is "teaches" and "taught" instead of "teachs" and "teached",
- It is "an instructor" instead of "a instructor",

BOL Syntax X	Semantics $[\![X]\!]$ in English
ontology	English text
D_1,\ldots,D_n	$\llbracket D_1 rbracket, \ldots, \llbracket D_n rbracket$
BOL declaration	dictionary entry or true sentence
$\mathbf{individual}i$	i is a proper noun.
$\mathbf{concept}i$	i is a common noun.
${\bf relation}i$	i is a transitive verb.
$\mathbf{property}i:T$	i is a common noun for a property that can take T -values.
I is-a ${\cal C}$	$\llbracket I rbracket{} \llbracket C rbracket{}$.
$I_1 R I_2$	$\llbracket I_1 rbracket \llbracket R rbracket \llbracket I_2 rbracket.$
I P V	$\llbracket I \rrbracket$ has $\llbracket P \rrbracket$ $\llbracket V \rrbracket$.
F	$\llbracket F rbracket$.
Formula	sentence
$C_1 \equiv C_2$	$\llbracket C_1 \rrbracket$ is the same as $\llbracket C_2 \rrbracket$.
$C_1 \sqsubseteq C_2$	everything that $[C_1]$ also $[C_2]$.
I is-a ${\cal C}$	$\llbracket I rbracket{} \llbracket C rbracket{}$.
$I_1 R I_2$	$\llbracket I_1 rbracket \llbracket R rbracket \llbracket I_2 rbracket.$
I P V	$\llbracket I \rrbracket \text{ has } \llbracket P \rrbracket \ \llbracket V \rrbracket.$
Individual	noun phrase (to be used as a subject or object)
i	i
Concept	intransitive verb phrase (to be plugged after a subject)
i	is a i
Τ	is anyone
<u></u>	is no one
$C_1 \sqcup C_2$	$\llbracket C_1 rbracket$ or $\llbracket C_2 rbracket$
$C_1 \sqcap C_2$	$\llbracket C_1 rbracket$ and $\llbracket C_2 rbracket$
$\forall R.C$	$[\![R]\!]$ something only if it $[\![C]\!]$
$\exists R.C$	$[\![R]\!]$ something that $[\![C]\!]$
$\operatorname{dom} R$	$[\![R]\!]$ something
$\operatorname{rng} R$	has something that $[\![R]\!]$ them
$\operatorname{dom} P$	has some $[\![P]\!]$
Relation	transitive verb phrase (to be plugged between subject and object)
i	i + third person singular s
$R_1 \cup R_2$	$\llbracket R_1 rbracket$ or $\llbracket R_2 rbracket$
$R_1 \cap R_2$	$\llbracket R_1 rbracket$ and $\llbracket R_2 rbracket$
$R_1; R_2$	$\llbracket R_1 \rrbracket$ something that $\llbracket R_2 \rrbracket$
R^{-1}	has something who $[\![R]\!]$ them and who is
R^*	$[\![R]\!]$ something that $[\![R]\!]$ something and so on that $[\![R]\!]$
Δ_C	$\llbracket C \rrbracket$ and is the same as
Property of type T	property phrase
i	i

Figure 6.9: Interpretation Function for BOL into English (intransitive VP version)

- Sentences start with upper case letters.
- Proper nouns often have different names in the ontology than in reality, e.g., it should be "Florian Rabe" and "credit value" instead of "Florian Rabe" and "credit Value".

Moreover, the language could be polished in many places. For example, "is a instructor and is a male" could become "is a instructor and a male" with a relatively easy special case treatment, or it could become "is a male instructor" with a more complex semantics that interprets some concepts via nouns and some via adjectives.

Remark 6.25 (Variants of English). We are relatively open as to what kind of English we want to use as the semantic language. The simplest choice would be to use plain English as you could find in a novel or newspaper article. But for many applications (e.g., formal ontologies in the STEM fields), we would rather use STEM English, i.e., English interspersed with formulas, diagrams, and epistemic cues like "Definition", "Theorem", "Proof", and even \square . For this kind of English, \LaTeX is a good target format. We can even use special \LaTeX dialects like sTeX \char [Koh08] where we can capture more of the semantic properties.

Remark 6.26 (Better Language Generation). While the target languages in the other translations are formal languages engineered for regularity and simplicity (in terms of language primitives), natural languages have evolved in practical human communication. As a consequence, the translation in Def. 6.22 results in English that is clumsy at best and non-grammatical in general. We can think of the result as **BOL-pidgin** English.

Let us have a look at some of the problems that appear in both translations:

- We need a lexicon to obtain inflection information and the translation tries to remedy that by appending "s" in various places. This works in some cases but not in others.
- there are many linguistic devices that serve an important role in natural language, but which we are not targeting. An example is plural objects for aggregation. Say we have Pis aC, Mis aC, this would translate to "P is a $[\![C]\!]$, M is a $[\![C]\!]$ " in BOL-pidgin, whereas in natural English we would aggregate this to "P and M are $[\![C]\!]$ s".

A way out is to utilize special systems for dealing with the surface structure of natural language. An example of this is the Grammatical Framework (GF, [Ran11]): it allows specifying a rich formal language of abstract syntax trees for natural language (ASTs) together with language-specific linearizations, which amount to recursive functions that translate ASTs to language-specific strings. GF comes with a large resource library that provides a comprehensive, language-independent AST specification and linearizations for over 35 languages. We will not pursue this here, but there is a special course "Logic-based Natural Language Semantics" at FAU in the Winter Semesters that covers these and related topics. One of the major issues that need to be addressed there and here is the notion of compositionality, which is central to all processing and semantics. We will address it next, and come back to it time and again later.

Chapter 7

Querying via a Semantics

7.1 Overview

Let us assume we have a semantics for our syntax. We again write l for the syntax, L for the semantic, and $\llbracket - \rrbracket$ for the translation function.

We can now use the semantics to answer questions asked in the syntax. Here we use the syntax to phrase a question and the semantics to determine the answer.

We call this *querying*. Contrary to standard practice, we will use that word in a very broad sense that covers all aspects. It is more common to use the word only for concretized querying, where SQL has been developed, which has shaped many intuitions about querying.

Usually, querying requires the syntax to designate some non-terminals as *propositional*. A non-terminal is propositional if the semantics can make its words true. Without a notion of propositions, it is impossible to define what questions and answers even are.

Definition 7.1 (Propositions). A context-free **syntax with propositions** is a context-free syntax with some designated non-terminal symbols.

A semantics with theorems is one that additionally defines some propositions to be theorems. We write $\vdash F$ if F is a theorem.

That definition does not mean that any kind of logic is needed for querying. Many languages use highly restricted notions of propositions that would not generally be considered as logic. For example, languages might use equalities between objects or even equalities between certain objects as the only propositions. The following table gives an overview:

aspect	typical propositions	proposition operators	
ontology language	assertions, concept equality/subsumption		
programming language	equality for some types	boolean operators	
database language	equality for base types	boolean operators	
logic	equality for all types	boolean operators, quantifiers	
natural language	sentences	and, or, some, every,	

Often the development of querying for a language leads to the discovery of omissions in the syntax: certain objects that are helpful to ask questions were omitted from the syntax because they were not needed to describe the data. Then sometimes the syntax is extended with non-terminals or productions that seem like dead code: they are not needed for or not allowed in the official data. The following table gives some examples:

aspect	typical extensions
ontology language	conjunction of assertions
programming language	quantifiers
database language	membership in a table
logic	(already tries to capture all possible propositions)
natural language	(already captures all possible propositions)

Example 7.2 (Propositions in BOL). The obvious choice of propositions for BOL are the formulas.

In Rem. 2.1, we mentioned that the BOL syntax from Fig. 2.2 had some redundant parts that were grayed out. Assertions are needed for writing ontologies only in such that they behave like axioms, i.e., they are automatically true. But for querying BOL, we also need them to behave like formulas so that we can use them as questions, i.e., we must allow them to be true or false.

Moreover, it is common to also allow conjunctions. Therefore, the BOL propositions are the conjunctions of formulas.

In the sequel, we will use each of the four kinds of semantics to

Section	builds on Section	aspect	query	result
7.2	6.3	deduction	proposition	yes/no
7.3	6.4	concretization	proposition with free variables	true ground instances
7.4	6.5	computation	term	value
7.5	6.6	narration	question	answer

Remark 7.3 (Meta-Level Questions). Finally, any semantics admits a meta-level where additional questions can be asked. Examples are asking for the consistency of a theory or the equivalence of two theories/programs/queries. At the next-higher meta-level, we can ask about the completeness of a semantics or the equivalence of two semantics (of which completeness is a special case). These meta-questions can usually not be expressed in the syntax, and we do not consider them a part of querying here. But it is worth mentioning that the need to use yet another language (a meta-language) to ask these questions can be annoying, and some advancements in language design are about trying to integrate them into the syntax. For example, reflection is the process of representing a language in itself so that the language can talk about itself. That way meta-questions become regular questions.

7.2 Deductive Querying

7.2.1 Method

We assume that l is a syntax with propositions and that L is a logic (and thus in particular has propositions) whose semantics has theorems.

It is not guaranteed [-] translates l-propositions to L-propositions. If not, we assume there is some operation True in L that we can use to lift the translations of l-propositions to L-propositions.

A deductive query consists of a proposition. The answer to the query is yes or no.

Thus, the deductive semantics must determine which propositions are theorems, i.e., whether $\vdash_L \text{True}(\llbracket F \rrbracket)$. This is usually done in one of two ways.

Proof theory uses a calculus for L to derive true propositions. Thus, we can say that an l-proposition F is true if the calculus derives $\mathsf{True}(\llbracket F \rrbracket)$. Accordingly, if L has a negation operator \neg , we can say that F is false if the calculus derives $\mathsf{True}(\neg \llbracket F \rrbracket)$.

Model theory uses a second deductive semantics, namely a translation from L to an even stronger deductive language M, usually some form of set theory. Then, we can say that l-propositions are true if the composition of the two translations yields a true M-proposition.

Either way, it is determined whether to answer an *l*-proposition with yes or no.

7.2.2 Challenges

Consistency The L-calculus might derive both F and $\neg F$. In that case L is inconsistent and usually every formula. We usually assume L to be consistent even though we do not always prove that.

Decidability Deductive semantics is usually undecidable, i.e., there is no algorithm that takes in F and always returns yes or no in finite time.

Therefore, deductive querying is very difficult in general. One has to run heuristics (theorem provers) to see if a proof of F or $\neg F$ can be found.

A common compromise is to allow only a restricted set of propositions as queries for which decision procedures exist. However, it can be tricky to find good restrictions, especially if the syntax allows for function symbols and equality.

For example, SFOL is undecidable. But many fragments of SFOL are decidable, such as propositional logic and various fragments in between.

When giving a deductive semantics into SFOL, it is therefore important to check whether the image of [-] falls inside a decidable fragment. This is typically the case for ontology languages.

Completeness Deductive semantics is usually incomplete, i.e., there are unanswered questions. More precisely, the L-calculus typically derives F for some propositions, $\neg F$ for some, but neither for some others. The third kind of proposition cannot be answered by the semantics.

Remark 7.4. The work "complete" is used for two different things in logic.

Firstly, it can be a relation between two semantics, typically proof theory and a model theory. That is the dominant meaning of the word as in, e.g., the completeness theorem for SFOL and Gödel's incompleteness theorem.

Secondly, it can mean that a logic proves or disproves every proposition, i.e., there is no F such that neither F nor $\neg F$ are derivable. That is the sense we use above. This kind of completeness rarely holds, usually only in very restricted circumstances.

Decidability and completeness are essentially the same problem. Specifically, if completeness holds, we already obtain a decision procedure for the logic: to decide the truth of F, enumerate all proofs until a proof of F or $\neg F$ is found. Vice versa, if we have a sound decision procedure, running it on F will prove either F or $\neg F$.

Efficiency Independent of whether the semantics is complete/decidable, theorem proving is typically very expensive.

Therefore, in addition to identifying decidable fragments of a logic, it is desirable to identify *efficiently* decidable fragments. Typically, a semantics meant for efficient practical querying aims for polynomially decidable fragments. This is the case for very simple ontology languages. But it can quickly become exponential if the language of propositions becomes more expressive.

7.3 Concretized Querying

This was discussed on the slides.

7.4 Computational Querying

This was discussed on the slides.

7.5 Narrative Querying

This was discussed on the slides.

Chapter 8

Formal Systems

8.1 Formal Systems

8.1.1 Syntax

Definition 8.1. A formal system consists of

- a set Voc of vocabularies.
- for any vocabulary V, a set Exp_V of expressions

A formal system with morphisms

- additionally provides for any two vocabularies V, W, a set $\operatorname{Voc}(V, W)$ of vocabulary morphisms from V to W
- for any vocabulary morphism $m \in \text{Voc}(V, W)$, a mapping $\text{Exp}_m : \text{Exp}_V \to \text{Exp}_W$

A formal system with typing

- additionally provides a relation $\vdash_V E : E'$ between expression $E, E' \in \text{Exp}_V$.
- such that for $m \in Voc(V, W)$, we have that if $\vdash_V E : E'$, then $\vdash_W Exp_m(E) : Exp_m(E')$.

A formal system with propositions has typing and

- additionally provides an expressions prop, in which case
 - expressions $\vdash_V F$: prop are called *propositions* (or formulas)
 - if $\vdash_V F_{\text{prop}}$, then expressions $\vdash_V P : F$ are called *proofs* of F
 - if some proof P of F exists, we also write $\vdash_V F$
- and now type preservation guarantees for $m \in Voc(V, W)$ we have that if $\vdash_V F$, then $\vdash_W Exp_m(F)$.

A formal system with equality has propositions and additionally provides a proposition $E \doteq_T E'$ for some expressions T and all $E, E' \in \operatorname{Exp}_V(T)$.

The vocabularies are usually lists of named declarations. In that case, the vocabulary morphism from V to W are usually lists of assignments c := E where c is the name of a declaration in V and E is a W-expression. In that case, the mapping Exp_m usually arises by replacing every V-identifier with the W-expression provided by m.

8.1.2 Translation

Definition 8.2. A **translation** from formal system l to formal system L consists of several mappings, all written $\llbracket - \rrbracket$:

- a vocabulary translation $Voc_l \to Voc_L$,
- for any l-vocabulary V, an expression translation $\operatorname{Exp}_V \to \operatorname{Exp}_{\llbracket V \rrbracket}$

If l and L have morphisms, a translation **maps morphisms** if it additionally provides

- for any two l-vocabularies V, W, a morphism translation $Voc(V, W) \to Voc(\llbracket V \rrbracket, \llbracket W \rrbracket)$
- such that for any vocabulary morphism $m \in \text{Voc}(V, W)$ and expression $E \in \text{Exp}_V$, we have $[\![\text{Exp}_m^l(E)]\!] = \text{Exp}_{\lceil m \rceil}^L([\![E]\!])$

If l and L have typing, the translation **preserves typing** if $\vdash_V^l E : E'$ implies $\vdash_{\llbracket V \rrbracket}^L \llbracket E \rrbracket : \llbracket E' \rrbracket$.

If l and L have typing, the translation **preserves propositions** if $[prop^l] = prop^L$.

If l and L have equality, the translation **preserves equality** if $[\![E \doteq_T E']\!] = [\![E]\!] \doteq_{[\![T]\!]} [\![E']\!]$.

The vocabulary translation usually consists of a semantic prefix and a declaration-wise translation: A V-vocabulary $D_1, \ldots D_n$ is translated to P, D'_1, \ldots, D'_n where the D'_i are the translations of the the D_i .

If a translation preserves typing and propositions, it automatically preserves truth, i.e., if $\vdash^L_{\llbracket V \rrbracket}$ implies $\vdash^L_{\llbracket V \rrbracket}$ $\llbracket F \rrbracket$. But the requirement to preserve propositions is often too strong. In practice, it is often better to use the following generalization:

Definition 8.3. If l and L have propositions, a **proposition-lifting** is an operation $\mathtt{True}: \llbracket \mathtt{prop}^l \rrbracket \to \mathtt{prop}^L$. A translation with proposition-lifting **preserves truth** if $\vdash^l_V F$ implies $\vdash^L_{\llbracket V \rrbracket} \mathtt{True} \llbracket F \rrbracket$.

8.1.3 Interpretation

Definition 8.4. An interpretation I of a formal system l consists of the following parts:

- for every l-vocabulary V a set Sit_V , whose elements are called situations,
- for every l-vocabulary V and every situation $S \in \operatorname{Sit}_V$, a function I_S mapping every $E \in \operatorname{Exp}_V$ to its interpretation $[\![E]\!]_S$.

If there is only one interpretation, the value $I_S(E)$ is often written as $[\![E]\!]^S$.

If l has morphisms, the interpretation is called **institutional** if it additionally provides

- for any l-morphism $m: V \to W$, a function $\operatorname{Sit}_m \to \operatorname{Sit}_W \to \operatorname{Sit}_V$
- such that for any $E \in \text{Exp}_V$ and $S \in \text{Sit}_W$, we have that $[\![\text{Exp}_m(E)]\!]^W = [\![E]\!]^{\text{Sit}_m(S)}$.

If l has typing, the interpretation is called **standard** if $\vdash_V^l E : E'$ implies $\llbracket E \rrbracket^S \in \llbracket E' \rrbracket^S$, i.e., types are interpreted as sets and typing as set membership.

If l has propositions, the interpretation is called **classical** if $[prop^l]^S = \{0,1\}$, i.e., propositions are interpreted as truth values.

If l has equality, the interpretation is called **canonical** if $[E \doteq_T E']^S = 1$ iff $[E]^S = [E']^S$

8.2 Semantics

Semantics can be defined absolutely or relatively. In the latter case, both translations and interpretations can be used to assign relative semantics to a formal system. Translation are relative to another formal system, interpretations to a situation.

8.2.1 General Case

Definition 8.5. Given a formal system with propositions, a **deductive semantics** consists of sets $\operatorname{Thm}_V \subseteq \operatorname{Exp}_V(\operatorname{prop})$, called the **theorems**. If $F \in \operatorname{Thm}_V$, we also write $\models_V F$.

If the formal system has typing, the truth judgment defines a deductive semantics by $\operatorname{Thm}_V = \{F \in \operatorname{Exp}_V(\operatorname{prop}) \mid \vdash_V F\}$. In that case, we speak of **soundness** if $\vdash_V F$ implies $\models_V F$ and of **completeness** in the reverse case.

Definition 8.6. Given a formal system with equality, a **computational semantics** consists of a binary relation $\text{Eval}_V \subseteq \text{Exp}_V \times \text{Exp}_V$. We also write $\text{Eval}_V(E) = E'$ as $\vdash_V E \leadsto E'$.

If Eval_V is a function and $\vdash_V E \leadsto E'$ implies $\vdash E \doteq_T E'$, we call it a **normal form**.

If additionally $\vdash_V E \doteq E'$ iff $\text{Eval}_V(E) = \text{Eval}_V(E')$, we call it a **canonical form**.

8.2. SEMANTICS 77

Definition 8.7. Given a formal system with typing and propositions, a **concrete semantics** consists of a function Inst that maps propositions F in context Γ to a set $Inst_V(\Gamma, F)$.

If there is a deductive sematics, the concrete semantics and the deductive are compatible if every $\gamma \in \text{Inst}_V(\Gamma, F)$ satisfies $\vdash_V \gamma : \Gamma$ and $\vdash_V F[\gamma]$.

We speak of an **absolute** semantics if the above are provided by a system of inference rules for the typing and truth judgments.

8.2.2 Relative Semantics by Translation

Consider a formal system translation from l to L with proposition lifting.

Definition 8.8. Given a deductive semantics for L, we define a deductive semantics for l by

$$\models^l_V F$$
 iff $\models^L_{\llbracket V \rrbracket} \mathtt{True}\llbracket F \rrbracket$.

Definition 8.9. Given a computational semantics for L, we define a computational semantics for l by

$$\vdash^l_V E \leadsto E' \qquad \text{iff} \qquad \vdash^L_{\llbracket V \rrbracket} \llbracket E \rrbracket \leadsto \llbracket E' \rrbracket.$$

This is a well-defined evaluation function only if Eval_L is and $\operatorname{Eval}_{\llbracket V \rrbracket}^L(\llbracket E \rrbracket)$ is always in the image of $\llbracket - \rrbracket$.

Definition 8.10. Given a concrete semantics for L, we define a concrete semantics for l by

$$\gamma \in \operatorname{Inst}(\Gamma, F) \qquad \text{iff} \qquad [\![\gamma]\!] \in \operatorname{Inst}([\![\Gamma]\!], \mathsf{True}[\![F]\!]).$$

8.2.3 Relative Semantics by Interpretation

Consider a classical interpretation I of formal system l. We can cast I as a special case of a formal system translation with proposition-lifting under the following practically reasonable assumptions:

- There is (at least implicitly) a formal system L for the background mathematical language.
- The situations in Sit_V are always tuples containing one component for every V-symbol.

In that case we define a translation $\overline{I}: l \to L$ as follows:

- Every vocabulary V is mapped to the L-vocabulary containing for every V-symbol s a declaration of the same name representing the corresponding element in the tuple.
- Every V-expression $E \in \text{Exp}_V(T)$ is mapped to the function $\overline{I}(E)$ that maps $\text{Sit}_V \ni S \mapsto I_S(E)$.
- The proposition lifting is given by

True
$$\overline{I}(F) = 1$$
 iff $I_S(F) = 1$ for all $S \in \operatorname{Sit}_V$.

We can then use that translation to define semantics as above.

In particular, The proposition lifting becomes a lot more obvious if we write [] for I and look at the induced deductive semantics: $\models F$ iff $[\![F]\!]^S = 1$ for all $S \in \operatorname{Sit}_V$.

Bibliography

- [CFKR21] J. Carette, W. Farmer, M. Kohlhase, and F. Rabe. Big Math and the One-Brain Barrier: The Tetrapod Model of Mathematical Knowledge. *The Mathematical Intelligencer*, 43(1):78–87, 2021.
- [Koh08] M. Kohlhase. Using LaTeX as a Semantic Markup Format. Mathematics in Computer Science, 2(2):279–304, 2008.
- [Ran11] A. Ranta. Grammatical Framework: Programming with Multilingual Grammars. CSLI Publications, 2011.